

HIERARCHICAL REPRESENTATIONS FOR SPATIO-TEMPORAL VISUAL ATTENTION MODELING AND UNDERSTANDING

Presented by
MIGUEL ÁNGEL FERNÁNDEZ TORRES

in partial fulfillment of the requirements for the
Degree of Doctor in Multimedia and Communications

UNIVERSIDAD CARLOS III DE MADRID

Advisors:
DR. IVÁN GONZÁLEZ DÍAZ
DR. FERNANDO DÍAZ DE MARÍA

Leganés, January 2019

Some rights reserved. This work is licensed under the Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License. To view a copy of this license, visit <http://creativecommons.org/licenses/by-nc-sa/4.0/> or send a letter to Creative Commons, PO Box 1866, Mountain View, CA 94042, USA.

We experienced events in an order,
and perceived their relationship as cause and effect.

They experienced all events at once, and perceived a purpose
underlying them all. A minimizing, maximizing purpose.

— Ted Chiang, *Stories of Your Life and Others*

AGRADECIMIENTOS

Ha llegado el momento. Quizás no sea el final que mi imaginación de niño soñó cuando comenzó esta tesis. Quizás no haya sido la aventura que esperábamos.

“Las historias son criaturas salvajes —dijo el monstruo—. Cuando las sueltas, ¿quién sabe los desastres que pueden causar?”

— Patrick Ness, *A Monster Calls*

Pero, en el fondo, sí que ha sido la experiencia que necesitaba. Pasa el tiempo apenas sin darnos cuenta, y la vida nos va llevando por esos lugares que nos quedan por descubrir, que nos hacen falta para crecer. A veces el amor al que hay que renunciar a cambio de una lección aprendida es muy grande, pero vale la pena seguir dando pasos hacia delante cuando finalmente eres consciente del camino que te quedaba por recorrer. Dejando el pasado atrás, y pensando casi por primera vez en el presente más que en el futuro, me siento muy afortunado de haber podido vivir estos diez años de formación en la Universidad Carlos III de Madrid. Han dado para tanto, y todo tan bueno...

Creo que he aprendido a investigar. Y es aquí donde tengo que agradecer enormemente, en primer lugar, a Fernando, el haberme dado la oportunidad de formar parte del Grupo de Procesado Multimedia y descubrir mi pasión por la Visión Artificial. A Iván, porque su talento y dedicación han conseguido sacar lo mejor de mí estos años. Gracias a los dos por la confianza y paciencia que habéis tenido conmigo, por enseñarme a creer en mí.

A los compañeros de laboratorio, por su apoyo en los momentos que más lo necesitaba. A Fernando, por regalarme los mejores consejos y conversaciones, y por preocuparse de cada uno de nosotros. También a Tomás, Amaya, Álex, FerFer, Rubén y Carmen. A Luis, por todas las películas y conciertos compartidos, me siento muy afortunado de haber podido conocerte más allá de estas aulas. A Ascen, porque desde aquella vez que me ayudaste a preparar mi Erasmus en Viena no has dejado de apoyarme en todo lo que necesitaba, siempre es un gusto charlar contigo. Muy especialmente a mis dos compañeros de croquetones. Gracias a mi compañero sobre-saliente Antonio por los momentos más divertidos, por haber sido capaz de volver a sacar el niño que llevo dentro cuando más lo necesitaba, por enseñarme a reírle a la vida. También a Javi. Empezamos juntos esta aventura hace diez años, y me lo has puesto todo tan fácil siempre... gracias por aportar la calma cuando era día

de tormenta. A Borja, Lorena y otros tantos compañeros de departamento junto a los que he crecido estos años. A Harold, por haberme prestado su ayuda para poder terminar esta tesis a tiempo.

He dado mis primeros pasos como profesor, y aquí también tengo que agradecer a muchos de los que he mencionado antes por ser un ejemplo a seguir, y por haber inculcado en mí esta vocación. No hay cosa que me haga más feliz cada semana que bajar a dar clase a un laboratorio y tratar de enseñar todo lo que puedo. Sí, sí la hay: Salir de una clase con esa sensación de haber compartido algo importante con los estudiantes que tengo a mi lado.

También he tenido oportunidad de viajar y conocer mundo. En Viena y West Lafayette he vivido la parte más impresionante de esta aventura. Gracias al Prof. Zygmunt Pizlo por acogerme en Purdue University y por lo mucho que ha aportado desde su perspectiva psicológica a esta tesis. Nunca olvidaré mis viajes a Chicago y Los Ángeles, los paseos por la noche junto a los rascacielos iluminados, todas las personas que fui encontrando por el camino. Tampoco a las que conocí en mi primera experiencia en un congreso en Bucharest. A Song y Souad, con las que he podido volver a encontrarme en Madrid no hace mucho.

Y entre investigación y docencia, también ha habido ratos para disfrutar. Y disfrutando hay veces que se ríe y otras que se llora, muchas de ellas de emoción, otras porque de los amigos también se aprende. Empezando por aquellos que he tenido la oportunidad de conocer en este lugar, o muy cerca, gracias a los compañeros de penúltimas y closing parties: Sergio, Unai, Adrián, Diego. A Cristina y Víctor, sois unos haters y lo sabéis. A mis compICñeros preferidos: Gisela, Raquel y Rafa. A Alba, qué pena que ya no cojamos el tren juntos para perdernos de camino a casa.

Y siguiendo con los que me han acompañado fuera de la universidad. Gracias a mi gran compañero de cruces y cruzadas, y mejor amigo Pedro. ¿Cuál es el próximo destino? A Álvaro, que coge ahora el testigo de esta tesis y comienza la suya. Gracias por estar ahí siempre este último año, por recordarme lo importante que es quererte tal y como eres, eres el amigo que cualquiera querría tener. A Patri, por los muchos pasos que llevamos ya dados juntos, y otras tantas carreras que vamos a compartir pronto, por ser tan buena amiga. A Wilson, al que también le gusta Vetusta, tengo ganas de más conciertos juntos. A Morganne, por esas tardes recorriendo Madrid y hablando de cine.

Para terminar, a mi familia. A mis padres, por haber antepuesto sus intereses a los míos, porque me han dejado siempre hacer lo que más quería. A mi madre, porque hemos hecho esta tesis codo con codo. Gracias por tus abrazos pero, sobre todo, por ser mi ejemplo de perseverancia y compromiso en esta vida. A mi padre, con el que espero tener más tiempo a partir de ahora para compartir todas esas

tecnologías que tanto le gustan. También a mi hermana, de la que admiro tanto su poca pereza, y a la que deseo que su esfuerzo tenga la recompensa que merece muy pronto, y podamos celebrar nuestros éxitos juntos.

Seguramente me estoy dejando a muchas personas fuera de estas líneas, por falta de memoria o espacio, no por ello menos importantes. Gracias a todas ellas por haberse cruzado alguna vez por esta historia, por hacerla única y tan emocionante, por ser fundamentales en esta etapa de crecimiento personal y profesional que hoy cierro con los sentimientos un tanto encontrados. La tristeza por terminar algo tan grande, y las ganas y la ilusión de ver lo que está por llegar. Todavía queda mucho por descubrir. Sigamos caminando.

Miguel Ángel

PUBLISHED AND SUBMITTED CONTENT

Some parts of the following publications are included or extended in this thesis:

1. Fernández-Torres, M. Á., González-Díaz, I., & Díaz-de-María, F. (2016, June). A probabilistic topic approach for context-aware visual attention modeling. In Content-Based Multimedia Indexing (CBMI), 2016 14th International Workshop on (pp. 1-6). IEEE.
<https://doi.org/10.1109/CBMI.2016.7500272>

Parts of this article are extended in Chapters 3 and 4 of the thesis. Whenever material from this source is included in this thesis, it is singled out with an explicit reference.

2. Fernández-Torres, M. Á., González-Díaz, I., & Díaz-de-María, F. (submitted). Probabilistic Topic Model for Context-Driven Visual Attention Understanding. IEEE Transactions on Circuits and Systems for Video Technology.

Parts of this article are included or extended in Chapters 3 and 4 of the thesis. Whenever material from this source is included in this thesis, it is singled out with an explicit reference.

OTHER RESEARCH MERITS

The following publications were part of my Ph.D. research, but the topics they cover are out of the scope of this thesis, so they are not included in it:

1. Fernández-Martínez, F., Hernández-García, A., Fernández-Torres, M. A., González-Díaz, I., García-Faura, Á., & de María, F. D. (2017). Exploiting visual saliency for assessing the impact of car commercials upon viewers. *Multimedia Tools and Applications*, 1-31.
<https://doi.org/10.1007/s11042-017-5339-9>
2. López-Labraca, J., Fernández-Torres, M. Á., González-Díaz, I., Díaz-de-María, F., & Pizarro, Á. (2018). Enriched dermoscopic-structure-based cad system for melanoma diagnosis. *Multimedia Tools and Applications*, 77(10), 12171-12202.
<https://doi.org/10.1007/s11042-017-4879-3>
3. Martínez-Cortés, T., Fernández-Torres, M. Á., Jiménez-Moreno, A., González-Díaz, I., Díaz-de-María, F., Guzmán-De-Villoria, J. A., & Fernández, P. (2014, October). A Bayesian model for brain tumor classification using clinical-based features. In *Image Processing (ICIP), 2014 IEEE International Conference on* (pp. 2779-2783). IEEE.
<https://doi.org/10.1109/ICIP.2014.7025562>

RESUMEN EXTENDIDO

En el siguiente resumen extendido se recogen los aspectos más relevantes de la presente Tesis doctoral. En primer lugar, se presenta la motivación del trabajo realizado. A continuación, se describen los principales objetivos y las contribuciones originales más destacadas. Finalmente, se resumen las conclusiones más relevantes, así como se mencionan posibles líneas futuras de investigación a partir del trabajo llevado a cabo.

MOTIVACIÓN DE LA TESIS

Dentro del marco de la Inteligencia Artificial, la Visión Artificial [1] es una disciplina científica que tiene como objetivo simular automáticamente las funciones del sistema visual humano, tratando de resolver tareas como la localización [2] y el reconocimiento [3] de objetos, la detección de eventos [4] o el seguimiento de objetos [5].

A pesar de la gran variedad de sistemas automáticos que se han desarrollado para resolver estos problemas, algunos de ellos verdaderamente efectivos, la mayor parte requiere procesar grandes cantidades de información visual, lo cual influye negativamente en su eficiencia. A diferencia de estos sistemas, el sistema visual humano es capaz de seleccionar de manera casi inmediata los elementos más importantes para poder interactuar en un contexto dado, al mismo tiempo que se ve atraído por aquellos estímulos más sorprendentes o llamativos. Esto se debe a su función de atención visual, la cual puede entenderse como un proceso de optimización para la percepción y la cognición visual. Si fuéramos capaces de diseñar algoritmos para la interpretación de escenas que llevasen a cabo esta función, podríamos ayudar a usuarios y expertos en escenarios de aplicación complejos, en los que se requiere procesar una gran cantidad de información simultáneamente, tales como la conducción [6], la aeronáutica [7] o la videovigilancia [8]. Esto permitiría disminuir la probabilidad de que se produzcan errores humanos, así como agilizar los procesos de decisión de los expertos.

La atención visual se puede estudiar en dos dominios diferentes: espacial y temporal. Estos dan lugar a definir tres tipos de modelos computacionales: espacial, espacio-temporal y temporal [9]. La mayor parte de los modelos computacionales de atención visual existentes consideran la componente espacial para guiar el procesamiento de la información visual hacia las regiones más llamativas o interesantes de una escena. Además, la información que percibimos en el mundo es dinámica, por lo que es igual de importante modelar cómo cambia

a lo largo del tiempo, lo que permite actualizar la atención espacial en función de las localizaciones seleccionadas previamente, así como seleccionar intervalos temporales de especial interés.

También es habitual distinguir entre dos familias de modelos de atención visual: modelos *Bottom-Up*, basados en las características visuales de la escena; y modelos *Top-Down*, los cuales tienen en cuenta un conocimiento a priori de la escena, o determinadas indicaciones para resolver una tarea [10, 11].

La motivación principal de esta tesis es, por tanto, el estudio y desarrollo de representaciones jerárquicas para el modelado y la interpretación de la atención visual espacio-temporal.

En particular, se proponen dos modelos computacionales de atención visual:

1. En primer lugar, se presenta un modelo generativo probabilístico *top-down* para el modelado y la interpretación de la atención visual en diferentes contextos.
2. En segundo lugar, se implementa una red profunda para el modelado de la atención visual. Esta arquitectura estima, en primer lugar, la atención visual espacio-temporal *top-down*, para finalmente modelar la atención en el dominio temporal. Su diseño está orientado a su aplicación final en un escenario de videovigilancia.

OBJETIVOS Y CONTRIBUCIONES ORIGINALES DE LA TESIS

Tal y como se comenta en el apartado anterior, esta tesis puede dividirse en dos partes.

En la primera parte de la tesis, se introduce nuestra primera aproximación: un modelo generativo probabilístico para el modelado y la interpretación de la atención visual espacio-temporal. El modelo propuesto, al que se ha denominado *visual ATtention Topic Model* o ATOM, es genérico, independiente del escenario de aplicación y está basado en las teorías psicológicas más destacadas sobre la atención visual [10, 11]. Además, considera la relación existente entre factores *bottom-up* y *top-down*.

Partiendo del conocido algoritmo *Latent Dirichlet Allocation* o LDA propuesto por David Blei et al. [12] para el análisis de corpus grandes de información textual, así como teniendo en cuenta dos de sus extensiones supervisadas [13, 14], nuestro modelo define la atención visual espacio-temporal *top-down* como una combinación de subtarefas latentes que, a su vez, se representan mediante combinaciones de características espacio-temporales de bajo, medio y alto nivel.

En particular, esta primera aproximación da lugar a las siguientes contribuciones:

- En primer lugar, se introduce un conjunto amplio de características de bajo nivel para el modelado de la atención visual, tales como el color, la intensidad, la orientación o el movimiento. A continuación, se procede a modelar características de medio y alto nivel, relacionadas con la estimación del movimiento de cámara en la escena y la detección de objetos.
- En segundo lugar, nuestro algoritmo incorpora un nivel intermedio formado por subtareas latentes. Este nivel permite acortar distancias entre la etapa de extracción de características y la de modelado de la atención visual, así como obtener representaciones de la atención más comprensibles y fáciles de interpretar.
- Además, nuestro modelo incorpora una variable categórica binaria que modela la atención visual en cada una de las localizaciones espaciales de una escena. Esta variable nos permite alinear automáticamente las subtareas determinadas por nuestro sistema con la información existente en aquellos lugares de la escena que atraen la atención de los usuarios.

A continuación, se incluye en la tesis un análisis exhaustivo de este algoritmo. Para ello, el modelo ATOM se utiliza para estimar e interpretar la atención visual en diferentes contextos (exteriores, videojuegos, noticias, etc.), definidos dentro de dos amplias bases de datos de vídeo anotadas con fijaciones de los ojos de diferentes sujetos: CRCNS-ORIG [15] y DIEM [16]. Esto permite ilustrar cómo nuestro modelo es capaz de aprender de manera efectiva representaciones jerárquicas de la atención visual adaptadas a diferentes escenarios. Además, los modelos obtenidos se utilizan para estimar la atención visual, comparando su eficiencia con la de otros modelos existentes en el estado del arte.

En la segunda parte de la tesis, se describe nuestra segunda aproximación: una red profunda que permite modelar la atención en el dominio temporal a partir de la atención visual espacio-temporal estimada. Nuestro algoritmo, al que se ha denominado *Spatio-Temporal to Temporal visual ATtention NETwork* o ST-T-ATTEN, modela la atención a lo largo del tiempo considerando una variable basada en las fijaciones proporcionadas por diferentes sujetos desarrollando una misma tarea.

En primer lugar, se introduce la hipótesis fundamental del modelo, la cual establece que la atención en el dominio temporal puede estimarse midiendo la dispersión de la localización de las fijaciones facilitadas por diferentes usuarios. Además, se puede entender la dimensión temporal de la atención como un mecanismo de filtrado, el cual permite identificar intervalos temporales de especial importancia en secuencias de vídeo.

En particular, nuestra segunda aproximación da lugar a las siguientes contribuciones:

- En primer lugar, se presentan tres arquitecturas para la extracción de características de alto nivel que permitan modelar la atención visual. Estas características, basadas en el color, el movimiento y los objetos presentes en la escena, servirán como entrada a nuestro sistema.
- En segundo lugar, se propone un *ground-truth* temporal a nivel de frame basado en las fijaciones de diferentes sujetos. Este *ground-truth* se obtiene atendiendo a la dispersión de las localizaciones fijadas. Además, permite validar la hipótesis fundamental de nuestro sistema, introducida anteriormente. Esta variable servirá para entrenar nuestros modelos para la estimación de la atención en el dominio temporal.
- Finalmente, se puede distinguir dos etapas en nuestro modelo:
 - 1) Una red para la estimación de la atención visual espacio-temporal, basada en una arquitectura de tipo codificador-decodificador convolucional, *Convolutional Encoder Decoder* o CED [17];
 - 2) Una red de tipo *Long Short-Term Memory* o LSTM [18] para el modelado de la dimensión temporal de la atención.

Finalmente, los experimentos de la segunda parte de la tesis tienen como objetivo validar las diferentes configuraciones propuestas para cada una de las etapas del modelo ST-T-ATTEN. Para ello, se hará uso de la base de datos BOSS [19], la cual contiene secuencias de vídeo grabadas en un contexto ferroviario, en las que tienen lugar diferentes eventos anómalos. Nuestro objetivo es determinar la configuración óptima para la arquitectura completa propuesta, así como motivar su uso como mecanismo de filtrado de información visual en un escenario de videovigilancia.

CONCLUSIONES

A lo largo de la tesis se han propuesto dos algoritmos jerárquicos para modelar la atención visual en secuencias de vídeo.

El primer algoritmo, presentado en la primera parte de la tesis y denominado ATOM, es un modelo generativo probabilístico para la estimación e interpretación de la atención visual espacio-temporal. La definición del sistema propuesto es genérica e independiente del escenario de aplicación. Además, el modelo se fundamenta en las teorías psicológicas más importantes acerca de la atención visual [10, 11], las cuales han establecido que la atención visual no se basa directamente en la información proporcionada por los procesos

visuales tempranos, sino en una representación contextual derivada de los mismos.

Utilizando como base el algoritmo LDA [12] y dos de sus extensiones supervisadas [13, 14], ATOM define la atención visual espacio-temporal *top-down* como una combinación de subtareas latentes que, a su vez, se representan mediante combinaciones de características espacio-temporales de bajo, medio y alto nivel. Por tanto, dado un frame en una secuencia de vídeo, el sistema recibe a su entrada un conjunto de mapas de características (color, intensidad, movimiento, etc.). A continuación, define un nivel de subtareas latentes entre la etapa de extracción de características y la de modelado de la atención visual. Finalmente, se utiliza una variable categórica binaria para alinear las subtareas definidas con la información derivada de las fijaciones de diferentes sujetos. Esta variable se genera a partir de un modelo de regresión logística aplicado sobre las subtareas, teniendo en cuenta la proporción en la que las mismas se dan en el frame.

El análisis llevado a cabo en la primera parte de la tesis demuestra la habilidad del modelo ATOM para aprender de manera efectiva, a partir de un conjunto amplio de características visuales, representaciones jerárquicas de la atención visual adaptadas específicamente a diferentes contextos (exteriores, videojuegos, deportes, noticias, etc.). Para ello, se ha hecho uso de dos amplias bases de datos de vídeo, CRCNS-ORIG [15] y DIEM [16]. Por otro lado, los experimentos muestran la facilidad de comprensión de las representaciones de la atención visual obtenidas por nuestro modelo, gracias a su uso de características tradicionales, tales como el color o el movimiento. Además, se observa que la detección de objetos o elementos sencillos como el rostro de las personas o el texto en una pantalla, modelados a continuación a partir de distribuciones espaciales discretas, así como el uso de *Redes Neuronales Convolucionales* o CNNs para obtener mapas de características basados en objetos, incrementa notablemente el rendimiento del sistema, permitiendo mejorar los resultados obtenidos por una gran variedad de métodos en el estado del arte a la hora de estimar la atención visual espacio-temporal.

En la segunda parte de la tesis, se describe nuestra segunda aproximación, denominada ST-T-ATTEN. Con este nuevo modelo se da un paso hacia adelante y se estima la atención en el dominio temporal, a partir de estimaciones de la atención visual espacio-temporal. La hipótesis fundamental de nuestro modelo establece que la atención en el dominio temporal puede estimarse midiendo la dispersión de la localización de las fijaciones proporcionadas por diferentes sujetos. En primer lugar, para demostrar esta hipótesis, se mide la correlación existente entre las secuencias definidas por el movimiento de los ojos de diferentes

sujetos cuando sucede un evento importante o anómalo, dadas las secuencias de vídeo recogidas en la base de datos BOSS [19]. A pesar de que este nivel de atención temporal constituye un indicador muy útil para detectar eventos importantes en escenarios complejos y concurridos, la atención en el dominio temporal ha de ser considerada siempre como un mecanismo de filtrado que permite seleccionar intervalos de tiempo candidatos a contener eventos sospechosos y que, por tanto, reduce el procesamiento posterior que tendría que llevar a cabo un sistema de detección de anomalías. Teniendo en cuenta esta hipótesis, el algoritmo ST-T-ATTEN trata de modelar la atención en el dominio temporal a partir de estimaciones de la atención visual espacio-temporal.

Motivados por el reciente éxito de las CNNs para el aprendizaje de representaciones jerárquicas profundas, así como de las LSTMs para el modelado de series temporales, el algoritmo propuesto se compone de dos etapas. La primera etapa, definida como *Spatio-Temporal visual Attention Network* o ST-ATTEN, consiste en una red de tipo CED que recibe a su entrada tres mapas de características de alto nivel para modelar la atención visual, basados en el color, movimiento y los objetos presentes en la escena. Todos estos mapas se obtienen a partir de CNNs. A continuación, esta arquitectura de tipo codificador-decodificador convolucional permite tanto estimar mapas de atención visual espacio-temporal como obtener representaciones latentes de la atención visual. Además, se proponen dos configuraciones para este módulo del sistema, las cuales se diferencian en las capas inicial y final del codificador y decodificador, respectivamente. En la primera configuración, estas capas son convolucionales, mientras que en la segunda son convolucionales de tipo LSTM.

La segunda etapa de nuestro sistema, denominada *Temporal Attention Network* o T-ATTEN, es una arquitectura de tipo LSTM, la cual permite estimar, para cada frame en una secuencia de vídeo, la atención en el dominio temporal. También se distingue entre dos versiones de T-ATTEN, dependiendo de si éste recibe a su entrada el mapa de atención visual espacio-temporal a la salida del decodificador en ST-ATTEN o, en cambio, la representación latente generada por el codificador.

A continuación, se ha evaluado la arquitectura ST-T-ATTEN propuesta en el escenario de videovigilancia definido por la base de datos BOSS [19], la cual incluye secuencias de vídeo que han sido grabadas en un contexto ferroviario, y que contienen diferentes tipos de eventos anómalos o sospechosos (varios abusos a mujeres, el robo de un teléfono móvil, una pelea entre pasajeros, etc.). El objetivo principal de los experimentos de la segunda parte de la tesis es evaluar las diferentes arquitecturas propuestas para nuestro modelo ST-T-ATTEN. A partir de estos experimentos, se ha determinado que la mejor configuración para nuestra arquitectura consiste, en primer

lugar, en una etapa ST-ATTEN con capas convolucionales, la cual permite fusionar de manera efectiva la información proporcionada por los tres mapas de características a su entrada. Después, el módulo T-ATTEN ofrece similares prestaciones tanto si recibe a su entrada un mapa o una representación latente de la atención visual.

Finalmente, se describen dos aplicaciones potenciales a nivel de usuario de nuestra propuesta. Por un lado, dado un escenario de videovigilancia, la atención estimada en el dominio temporal puede aplicarse para seleccionar en tiempo real las cámaras más importantes en un array de monitores, dirigiendo la atención de los operadores hacia aquellas cámaras que potencialmente muestran anomalías o eventos sospechosos. Por otro lado, esta atención temporal se puede aplicar también en tareas off-line que implican la visualización de una gran cantidad de horas de grabaciones de videovigilancia, reduciendo la cantidad de información que los operadores tienen que procesar. Por tanto, se puede concluir que, introduciendo algunas mejoras al sistema propuesto, éste podría ser capaz de proporcionar a los operadores una experiencia completa de la atención visual, identificando no únicamente las localizaciones más llamativas de la escena, sino también seleccionando intervalos temporales relevantes, de acuerdo con los eventos que han tenido lugar previamente en la escena, así como con los eventos que están sucediendo en otras cámaras al mismo tiempo.

LÍNEAS FUTURAS DE INVESTIGACIÓN

Finalmente, en esta sección se identifican y comentan las líneas futuras de investigación más prometedoras en relación con el trabajo presentado en esta tesis.

Llegados a este punto, no cabe duda acerca de las enormes ventajas que tiene el modelado de la atención visual dentro del campo de la Inteligencia Artificial. Tampoco sobre las infinitas posibilidades que un concepto tan abstracto tiene para el procesamiento y la interpretación de un mundo en el que cada vez se maneja una mayor cantidad de datos. A pesar de la gran variedad de modelos computacionales de atención visual existentes en la literatura, todavía queda mucho camino por recorrer, no sólo para lograr un sistema que modele automáticamente esta función cognitiva, sino también para entender cómo el sistema visual humano lleva a cabo este proceso de optimización.

Teniendo en cuenta los dos paradigmas más populares en la actualidad para el aprendizaje de representaciones, los cuales se basan en el Aprendizaje Profundo, *Deep Learning* o DL, y los Modelos Gráficos Probabilísticos, *Probabilistic Graphical Models* o PGM, nuestras contribuciones han demostrado la importancia tanto de la tarea de *percibir*, desempeñada por representaciones jerárquicas

profundas, como de la habilidad de *deducir*, característica de los modelos PGM, a la hora de modelar e interpretar la atención visual.

En primer lugar, es importante conseguir buenas representaciones del mundo que nos rodea para modelar la atención, y es ahí donde las redes profundas y, en particular, las CNNs, desempeñan un papel fundamental en la percepción automática. Además, dado que la atención visual lleva a cabo no una, sino varias tareas complejas, es fundamental poder interpretar cómo un modelo computacional hace uso de las representaciones jerárquicas proporcionadas por redes profundas. Esto se puede conseguir a partir de métodos probabilísticos que permitan definir relaciones entre las variables observadas. Esta dirección, definida recientemente como *Bayesian Deep Learning* o BDL [20], es la que queremos seguir en trabajos futuros, prestando una especial atención a la aplicación de BDL a los modelos probabilísticos de temas latentes [12-14], los cuales constituyen la base de nuestra primera aproximación para la interpretación de la atención visual: ATOM. Si consiguiéramos definir subtarefas no solamente en el espacio, sino también a lo largo del tiempo, podríamos establecer relaciones entre los conceptos reconocidos en una o varias secuencias de vídeo, tanto en la misma escena como en escenas diferentes.

En segundo lugar, se ha demostrado en la segunda parte de la tesis las importantes ventajas que tiene modelar la atención en el dominio temporal, la cual permite seleccionar intervalos temporales de especial importancia en secuencias de vídeo. Estos intervalos seleccionados ayudan a reducir, además, la carga computacional en posibles aplicaciones a nivel de usuario. Desde esta perspectiva, la atención visual apenas ha sido tratada en el estado del arte hasta la fecha, a pesar de su utilidad para el procesado y el análisis de grandes cantidades de información visual, en aplicaciones como la detección de anomalías.

Una línea de investigación interesante que no se ha tratado en esta tesis es la interpretación de las secuencias definidas por el movimiento de los ojos, lo cual facilitaría la implementación de sistemas de mayor comprensión y utilidad para estimar la variación de la atención visual a lo largo del tiempo. Para ello, creemos que el uso de métodos de aprendizaje por refuerzo o *Reinforcement Learning* puede ser un camino prometedor a seguir [21].

Por último, estamos motivados a continuar estudiando el modelado tanto de la atención visual espacio-temporal como de la atención en el dominio temporal en secuencias de vídeo reproducidas al mismo tiempo, con el objetivo de ayudar a los expertos en escenarios complejos y concurridos. Para ello, en los próximos meses se procederá a anotar bases de datos de vídeo grandes, tales como VIRAT [22] o UCF-Crime [23], con fijaciones de diferentes sujetos, las cuales servirán para un análisis más completo

de la arquitectura ST-T-ATTEN propuesta, así como para introducir posibles mejoras en la misma.

ABSTRACT

This PhD. Thesis concerns the study and development of hierarchical representations for spatio-temporal visual attention modeling and understanding in video sequences. More specifically, we propose two computational models for visual attention. First, we present a generative probabilistic model for context-aware visual attention modeling and understanding. Secondly, we develop a deep network architecture for visual attention modeling, which first estimates top-down spatio-temporal visual attention, and ultimately serves for modeling attention in the temporal domain.

The first part of the thesis introduces our first proposal: a generative probabilistic framework for spatio-temporal visual attention modeling and understanding. The model proposed is generic, independent of the application scenario and founded on the most outstanding psychological studies about attention. Moreover, it considers the existing concurrence between bottom-up and top-down factors.

Drawing in the well-known Latent Dirichlet Allocation method for the analysis of large corpus of data, and some of its supervised extensions, our approach defines task- or context-driven visual attention in video as a mixture of latent sub-tasks, which are in turn represented as combinations of low-, mid- and high-level spatio-temporal features. Latent sub-tasks discovered are automatically aligned to the information drawn from human fixations by means of a categorical variable response, which is generated by a logistic regression model over the sub-task proportions. Therefore, our algorithm incorporates an intermediate level formed by latent sub-tasks, which bridges the gap between features and visual attention, and enables to obtain more comprehensible interpretations of attention guidance.

The experiments related to our first approach demonstrate its ability to successfully learn hierarchical representations of visual attention, specifically adapted to diverse contexts, on the basis of a wide set of features. Besides, results show how our proposal significantly outperforms quite a few competent methods in the literature when estimating visual attention.

The second part of the thesis presents our second proposal: a deep network architecture that takes a step further and goes from spatio-temporal visual attention prediction to attention estimation in the temporal domain. The model is fundamentally supported by the assumption that a measurement of task-driven visual attention in the temporal domain can be drawn from the dispersion of fixation

locations recorded from several observers. Although this temporal level of attention constitutes a useful clue to detect important events in crowded and complex scenarios, attention in the temporal domain should be considered as an early filtering mechanism, which selects candidate time segments to contain suspicious events, and therefore reduces the later processing devoted to the anomaly detection.

Based on this hypothesis, and inspired by the recent success of Convolutional Neural Networks for learning deep hierarchical representations and Long Short-Term Memory Units for time series forecasting, our approach is composed of two stages. On the one hand, the first stage consists of a Convolutional Encoder Decoder network that receives at its input three high-level feature maps for visual attention guidance: RGB-based, motion and objectness. Then, through an encoding-decoding architecture, the network concurrently estimates spatio-temporal visual attention maps and extracts latent representations of visual attention. On the other hand, the second stage involves an architecture based on Long Short-Term Memory Units that estimates, for each frame in a video sequence, a temporal attention response. We propose different configurations for both stages, in order to assess various architectures of our proposal.

Finally, the second approach proposed is evaluated in a video surveillance scenario, which contains video sequences recorded in a railway transport context, with different types of suspicious or anomalous events. In addition, we discuss two potential end-user applications for our proposal. On the one hand, given a surveillance scenario, the estimated temporal attention response could be applied to select in real-time the most outstanding screens from the monitoring array, thus driving operator's attention to scenes that potentially show anomalies or suspicious events. On the other hand, this response could be also applied in off-line tasks which imply reviewing many hours of surveillance recordings, reducing the information to be processed by the operator.

CONTENTS

List of Figures	xxx
List of Tables	xxxv
Acronyms	xxxvi
1 INTRODUCTION	1
1.1 Visual attention	1
1.2 Hierarchical representations for visual attention	2
1.2.1 Feature engineering	3
1.2.2 Machine Learning	4
1.2.3 Representation learning	9
1.3 Goals and context of the thesis	11
1.4 Structure of the thesis and contributions	13
2 A MULTIDISCIPLINARY PERSPECTIVE ON VISUAL ATTENTION	17
2.1 Introduction	17
2.2 Neurophysiological basis of visual attention	18
2.2.1 Human Visual System: eye and brain	18
2.2.2 Visual attention	22
2.3 Psychophysical theories of visual attention	24
2.3.1 The Feature Integration Theory	24
2.3.2 The Guided Search Model	25
2.3.3 What are the attributes that guide attention? . .	27
2.3.4 Eye movements	27
2.4 Computational modeling of visual attention	28
2.4.1 Bottom-up versus top-down approaches	30
2.4.2 Bayesian models	34
2.4.3 Deep Neural Networks	35
2.4.4 Applications	37
3 A GENERATIVE PROBABILISTIC MODEL FOR SPATIO-TEMPORAL VISUAL ATTENTION	39
3.1 Introduction	39
3.2 Related work and main contributions	40
3.3 Feature engineering for visual attention guidance . . .	41
3.3.1 Basic features: color, intensity and orientation .	41
3.3.2 Motion-based features	45
3.3.3 Novelty features	50
3.3.4 Object-based features	51
3.4 Latent Topic Models	55
3.4.1 Bag-of-Words model	56
3.4.2 Latent Dirichlet Allocation	57
3.4.3 Supervised topic models	59
3.4.4 Applications to Computer Vision	61

3.5	Visual Attention Topic Model	62
3.5.1	Model overview	62
3.5.2	Guiding features extraction	66
3.5.3	Inference process	68
3.5.4	Learning sub-tasks for spatio-temporal visual attention estimation	72
4	EXPERIMENTS ON CONTEXT-DRIVEN VISUAL ATTENTION UNDERSTANDING AND PREDICTION	75
4.1	Introduction	75
4.2	Experimental design	75
4.2.1	Databases	76
4.2.2	Experimental setup	77
4.2.3	Evaluation metrics	77
4.2.4	Model initialization	80
4.3	Understanding visual attention as a mixture of sub-tasks	80
4.4	Results on visual attention estimation	85
4.5	Comparison with state-of-the-art methods	87
4.6	Where we are: model strengths and limitations	89
4.7	Conclusions	93
5	DEEP NEURAL NETWORKS FOR MODELING VISUAL ATTENTION IN THE TEMPORAL DOMAIN	95
5.1	Introduction	95
5.2	Related work	96
5.3	Deep Neural Networks	98
5.3.1	Neural Networks	98
5.3.2	Convolutional Neural Networks	108
5.3.3	Encoder-Decoder Networks	113
5.3.4	Recurrent Neural Networks	116
5.4	Feature learning for visual attention guidance	120
5.4.1	RGB-based spatial network	121
5.4.2	Optical flow-based motion network	121
5.4.3	Objectness-based network	122
5.5	Spatio-Temporal to Temporal Visual Attention Network	124
5.5.1	Fundamental hypothesis of the model	124
5.5.2	Fixation-based temporal ground-truth	127
5.5.3	Model overview	130
5.5.4	Spatio-Temporal Visual Attention Network	131
5.5.5	Temporal Attention Network	134
5.5.6	Implementation details	135
6	EXPERIMENTS ON TEMPORAL VISUAL ATTENTION ESTIMATION IN A VIDEO SURVEILLANCE SCENARIO	137
6.1	Introduction	137
6.2	Experimental design	137
6.2.1	Databases	138
6.2.2	Experimental setup	139
6.2.3	Evaluation metrics	139

6.2.4	Training and implementation details	140
6.3	Results on spatio-temporal visual attention estimation with ST-ATTEN	141
6.4	Results on attention estimation in the temporal domain with ST-T-ATTEN	144
6.5	Where we are: towards guiding anomaly detection . .	147
6.6	Conclusions	150
7	CONCLUSIONS AND FUTURE LINES OF RESEARCH	153
7.1	Conclusions	153
7.2	Future lines of research	155
A	DERIVATION OF THE FORMULAS FOR THE ATOM	159
A.1	Expansion of the lower bound	159
A.1.1	Lower bound of the local appearance model . .	160
A.1.2	Lower bound of the visual attention response .	161
A.2	Derivation of the formulas for the variational parameters	162
A.3	Derivation of the formulas for the model parameters .	164
A.4	Derivation of the formulas for the logistic regression model	165
B	EYE-TRACKING DATABASES USED IN THE THESIS	167
B.1	CRCNS-ORIG database	167
B.1.1	Description	167
B.1.2	Video categories	169
B.1.3	Context-aware visual attention understanding .	169
B.2	DIEM database	172
B.2.1	Description	172
B.2.2	Video categories	172
B.2.3	Context-aware visual attention understanding .	175
B.3	BOSS database	178
B.3.1	Description	178
B.3.2	Video sequences	178
	BIBLIOGRAPHY	181

LIST OF FIGURES

Figure 1.1	Graphical representation of Bayesian (a) discriminative and (b) generative supervised models.	6
Figure 1.2	Graphical representation of Bayesian directed generative unsupervised models.	9
Figure 1.3	Visual attention features and representation models covered by the different systems presented in the thesis, classified according to the spatial, spatio-temporal or temporal dimension in video sequences where they are modeled.	12
Figure 1.4	Visual attention features and representation models covered by the different systems presented in the thesis, classified according to the types of processes that they involve.	14
Figure 2.1	Diagram of the human eye.	19
Figure 2.2	Schematic diagram of a biological neuron.	20
Figure 2.3	Diagram of the human brain.	21
Figure 2.4	Diagram representations of (a) the Feature Integration Theory (FIT) and (b) the Guided Search Model (GSM) [11].	25
Figure 2.5	Chronological timeline of the visual attention models in the <i>state-of-the-art</i> reviewed in this thesis.	29
Figure 3.1	Basic, motion-based and novelty feature maps computed for two example frames taken from Videogames (a) and Sports (b) categories from CRCNS-ORIG [15] database.	42
Figure 3.2	Motion parameterization in two example frames taken from a commercials video in the DIEM [16] database.	48
Figure 3.3	Camera motion modeling in an Outdoor frame taken from the CRCNS-ORIG [15] database.	49
Figure 3.4	Harris response for text detection in an example frame taken from a TV news video in the DIEM [16] database.	53
Figure 3.5	Object-based feature maps computed for example frames taken from TVNews (a, b, c, d) and TalkShows (e) categories from CRCNS-ORIG [15] database.	54

Figure 3.6	Bag-of-Words (BoW) model applied to: (a) A corpus of texts from movie reviews, (b) A corpus of images.	56
Figure 3.7	(a) Graphical representation of Latent Dirichlet Allocation (LDA) [12]. (b) Graphical representation of the variational distribution used to approximate the posterior in LDA. . . .	57
Figure 3.8	Graphical representation of Supervised Latent Dirichlet Allocation (sLDA) [13].	60
Figure 3.9	Visual attention modeled in three different scenarios taken from CRCNS-ORIG [15] database as a mixture of several relevant sub-tasks, associated with particular areas of special importance for observers.	63
Figure 3.10	Graphical representation of the proposed visual Attention Topic Model (ATOM) generative model.	65
Figure 3.11	Processing pipelines of the generative probabilistic approach proposed for spatio-temporal visual attention modeling. . .	73
Figure 4.1	TPs and FPs sampled in an example frame taken from a documentary video in DIEM [16] database, according to a probabilistic shuffle map.	78
Figure 4.2	Three most prominent attracting (AT) and inhibiting (IT) sub-tasks inferred by (a) <i>Outdoor</i> and (b) <i>TV News context-aware</i> models learned based on CRCNS-ORIG [15] database.	83
Figure 4.3	Three most prominent attracting (AT) and inhibiting (IT) sub-tasks inferred by (a) <i>Commercials</i> and (b) <i>Sports context-aware</i> models learned based on DIEM [16] database.	84
Figure 4.4	Visual attention maps obtained by ATOM for some example frames from CRCNS-ORIG [15] database.	85
Figure 4.5	Results obtained by the proposed <i>context-generic</i> and <i>context-aware</i> ATOM models in the CRCNS-ORIG [15] database. . .	86
Figure 4.6	Visual attention maps obtained by ATOM for some example frames from DIEM [16] database.	87
Figure 4.7	Results obtained by the proposed <i>context-generic</i> and <i>context-aware</i> ATOM models in the DIEM [16] database.	88

Figure 4.8	Visual attention maps generated by some of the most outstanding methods in the <i>state-of-the-art</i> for some intricate example frames taken from CRCNS-ORIG [15] and DIEM [16] databases.	91
Figure 4.9	Frame sequences taken from CRCNS-ORIG [15] database to analyze some ATOM model drawbacks and define future lines of research.	92
Figure 5.1	a) Typical video surveillance monitoring room. (b) The task of a CCTV operator is to find a possible anomalous event amongst multiple distractors, displayed at the same time in a large array of more than 20 screens.	96
Figure 5.2	Mathematical model of a computational neuron.	99
Figure 5.3	Graphical representations of feed-forward Neural Networks (NNs).	100
Figure 5.4	Graphical representation of the most commonly used activation functions. (a) Sigmoid. (b) Hyperbolic tangent. (c) Rectified Linear Unit (ReLU). (d) Exponential Linear Unit (ELU) represented for different values of α	102
Figure 5.5	Computational graph of a Neural Network (NN) layer, where forward and back-propagation stages in gradient-based learning are represented with green and red arrows, respectively.	106
Figure 5.6	Application of a $k = 3 \times 3$ convolutional kernel over a 6×6 input padded with a 1×1 border of zeros, using a stride of $s = 2$	109
Figure 5.7	Application of a $k = 3 \times 3$ dilated convolutional kernel over a 7×7 input, using a dilation factor of $d = 2$ (1 space between kernel elements).	110
Figure 5.8	Architecture diagrams of the (a) VGG-16 [113] and (b) ResNet-50 [193] networks for image recognition.	112
Figure 5.9	(a) Graphical representation of an Encoder-Decoder Network (EDN). (b) Example diagram of a convolutional encoder-decoder architecture for medical image segmentation.	114
Figure 5.10	Graphical representations of a Recurrent Neural Network (RNN).	117
Figure 5.11	Diagram of a Long Short-Term Memory (LSTM) unit.	119

Figure 5.12	RGB-based, motion and objectness feature maps computed for an example frame taken from BOSS [19] database.	121
Figure 5.13	Visual attention in the temporal domain modeled in a video-surveillance sequence taken from BOSS [19] database.	123
Figure 5.14	Processing pipelines of the Spatio-Temporal to Temporal visual ATtention NETwork (ST-T-ATTEN) proposed.	126
Figure 5.15	Visual attention in the temporal domain modeled in three video-surveillance sequences taken from BOSS [19] database.	129
Figure 5.16	Diagram of the Spatio-Temporal to Temporal visual ATtention NETwork (ST-T-ATTEN) proposed.	131
Figure 5.17	Diagram of the ST-ATTEN for spatio-temporal visual attention estimation, which consists of a CED architecture.	133
Figure 5.18	Architecture diagrams of the ST-T-ATTEN configurations proposed.	134
Figure 6.1	Visual attention maps obtained by ST-ATTEN for some example frames taken from BOSS [19] database.	142
Figure 6.2	Visual attention in the temporal domain \hat{a}_t estimated by ST-T-ATTEN in a video-surveillance sequence taken from BOSS [19] database.	146
Figure 6.3	Visual attention in the temporal domain \hat{a}_t estimated by ST-T-ATTEN in two video-surveillance sequences taken from BOSS [19] database.	148
Figure B.1	CRCNS-ORIG [15] database: <i>Context-Generic</i>	169
Figure B.2	CRCNS-ORIG [15] database: <i>Outdoor</i>	169
Figure B.3	CRCNS-ORIG [15] database: <i>Videogames</i>	170
Figure B.4	CRCNS-ORIG [15] database: <i>Commercials</i>	170
Figure B.5	CRCNS-ORIG [15] database: <i>TV News</i>	170
Figure B.6	CRCNS-ORIG [15] database: <i>Sports</i>	171
Figure B.7	CRCNS-ORIG [15] database: <i>Talk Shows</i>	171
Figure B.8	CRCNS-ORIG [15] database: <i>Others</i>	171
Figure B.9	DIEM [16] database: <i>Context-Generic</i>	176
Figure B.10	DIEM [16] database: <i>TV Shows</i>	176
Figure B.11	DIEM [16] database: <i>Documentaries</i>	176
Figure B.12	DIEM [16] database: <i>Commercials</i>	177
Figure B.13	DIEM [16] database: <i>Talk Shows</i>	177
Figure B.14	DIEM [16] database: <i>Sports</i>	177
Figure B.15	DIEM [16] database: <i>Cooking</i>	178

Figure B.16	DIEM [16] database: <i>TV News</i>	178
-------------	--	-----

LIST OF TABLES

Table 4.1	Categories into which (a) CRCNS-ORIG [15] and (b) DIEM [16] databases are divided. . . .	77
Table 4.2	Comparison with state-of-the-art methods in the CRCNS-ORIG [15] database.	89
Table 4.3	Comparison with state-of-the-art methods in the DIEM [16] database.	90
Table 5.1	Encoder and decoder architectures for the CONV-ST-ATTEN and CONV-LSTM-ST-ATTEN configurations of the Spatio-Temporal visual ATtention NETwork (ST-ATTEN) proposed. . . .	132
Table 6.1	Results obtained on the BOSS [19] database by the proposed ST-ATTEN and other methods for comparison when estimating spatio-temporal visual attention.	143
Table 6.2	Results obtained on the BOSS [19] database by the proposed ST-T-ATTEN when modeling visual attention in the temporal domain. . . .	145
Table 6.3	Results obtained by the proposed ST-T-ATTEN and other comparison methods considered as filtering mechanisms for guiding anomaly detection in the video surveillance scenario defined by the BOSS [19] database.	149
Table B.1	Categories in the CRCNS-ORIG [15] database. Clips included in each category are enumerated, together with their number of frames.	168
Table B.2	Categories in the DIEM [16] database. Clips included in each category are enumerated, together with their number of frames.	172
Table B.3	Videos from the BOSS [19] database for the experiments in Chapter 6. Clips are enumerated together with their number of frames.	179

ACRONYMS

AI	Artificial Intelligence
AIM	Attention based on Information Maximization
AT	Attractive Topic
ATOM	visual Attention TOpic Model
AUC	Area Under ROC Curve
AWS	Adaptive Whitening Saliency
AWS-D	Dynamic Adaptive Whitening Saliency
BDL	Bayesian Deep Learning
BN	Batch Normalization
BoW	Bag-of-Words
BPTT	Back-propagation Through Time
BU	Bottom-Up
CC	Correlation Coefficient
CCTV	Closed-Circuit TeleVision
CED	Convolutional Encoder Decoder
CNN	Convolutional Neural Network
CONV	Convolutional
CPU	Central Processing Unit
CRF	Conditional Random Field
C-A	Context-Aware
C-G	Context-Generic
DBA	Dirichlet-Bernoulli Alignment
DL	Deep Learning
DNN	Deep Neural Network
DoG	Difference of Gaussians
EDN	Encoder-Decoder Network

ELBO	Evidence Lower Bound
ELU	Exponential Linear Unit
EM	Expectation-Maximization
EMA	Exponential Moving Average
FC	Fully-Connected
FIT	Feature Integration Theory
FOA	Focus of Attention
FP	False Positive
GAN	Generative Adversarial Network
GBVS	Graph-Based Visual Saliency
GMM	Gaussian Mixture Model
GPU	Graphics Processing Unit
GRU	Gated Recurrent Unit
GSM	Guided Search Model
GT	Ground-Truth
HMM	Hidden Markov Model
HOG	Histogram of Oriented Gradients
HVS	Human Visual System
ICA	Independent Component Analysis
ICL	Incremental Coding Length
IID	Independent and Identically Distributed
IOR	Inhibition of Return
IFT	Inverse Fourier Transform
ILSVRC	ImageNet Large Scale Visual Recognition Competition
IT	Inhibiting Topic
KL	Kullback-Leibler divergence
LDA	Latent Dirichlet Allocation
LN	Layer Normalization
LSK	Local Steering Kernels

LSTM	Long Short-Term Memory
LTM	Latent Topic Model
MAE	Mean Absolute Error
MBGD	Mini-Batch Gradient Descent
ML	Machine Learning
MLP	Multi-layer Perceptron
MPSE	Mean Pairwise Squared Error
MRF	Markov Random Field
MSE	Mean Squared Error
NLP	Natural Language Processing
NN	Neural Network
NSS	Normalized Scanpath Saliency
NUS	Non-Uniform Sampling
PCA	Principal Component Analysis
PCC	Pearson Correlation Coefficient
PDF	Probability Density Function
PGM	Probabilistic Graphical Models
PLSA	Probabilistic Latent Semantic Analysis
POOL	Pooling
PQFT	Phase spectrum of Quaternion Fourier Transform
RAM	Random Access Memory
ReLU	Rectified Linear Unit
RMSP _{prop}	Root Mean Squared Propagation
RNN	Recurrent Neural Network
ROC	Receiver Operating Characteristic
ROI	Region of Interest
SAM	Saliency Attentive Model
sAUC	Shuffled Area Under ROC Curve
SC	Superior Colliculus

SDSR	Saliency Detection by Self-Resemblance
SGD	Stochastic Gradient Descent
SIFT	Scale-Invariant Feature Transform
sLDA	Supervised Latent Dirichlet Allocation
sNSS	Shuffled Normalized Scanpath Saliency
SM	Saliency Map
SSD	Sum of Squared Differences
ST-ATTEN	Spatio-Temporal visual ATtention NETwork
ST-T-ATTEN	Spatio-Temporal to Temporal visual ATtention NETwork
SUN	Saliency Using Natural Statistics
SVM	Support Vector Machine
T-ATTEN	Temporal ATtention NETwork
TD	Top-Down
TP	True Positive
VAM	Visual Attention Map
V ₁	Primary Visual Cortex
VIDYA	Variable Index Dynamic Average
WMAP	Weighted Maximum Phase Alignment
WTA	Winner-Take-All
1D	one dimensional
2D	two dimensional
3D	three dimensional

INTRODUCTION

This introductory chapter is organized as follows. First, we make a short introduction to visual attention in Section 1.1, enumerating its different types and possible applications. Secondly, in Section 1.2, we discuss the use of hierarchical representations for visual attention, going from feature engineering to feature learning, through a brief description of the types of Machine Learning (ML) covered by the systems presented in the thesis. Then, in Section 1.3, we introduce the main focus of the thesis, which is the study and development of hierarchical representations for spatio-temporal visual attention modeling and understanding. Finally, Section 1.4 summarizes the structure and contributions of this dissertation.

1.1 VISUAL ATTENTION

“The world and its universe are, to anything or anyone with senses, incomprehensibly big data.” (Mark Andrejevic, 2014) [24]

We have been always surrounded by data. However, never before had we lived in such a data-driven world. Nowadays, unstoppable technological advances make possible to capture almost everything, anywhere and anytime, which has resulted in a massive amount of information that is necessary to filter and process.

Within the framework of Artificial Intelligence (AI), Computer Vision [1] emerged in the late 1960s with the objective of automatically simulating the Human Visual System (HVS) functions. Drawing from the visual information captured in digital images and video sequences, this interdisciplinary field seeks to discover good representations of the real-world in order to carry out particular tasks such as object location [2] and recognition [3], event detection [4] or visual tracking [5].

In spite of the wide variety of systems that are continuously released and improved to solve these tasks, some of them truly effective, they still need to process large amounts of visual information for achieving high performances, which dramatically impacts on their efficiency. Human beings, however, inherently select

the most important elements to interact in a context and, besides, are rapidly attracted by striking stimulus. And this is thanks to the visual attention function of the [HVS](#), which can be understood as an optimization process for visual cognition and perception. If we were able to design image-understanding algorithms that accomplish this operation, we could use them to reduce their computational cost. At the same time, we would help users and experts when dealing with applications and complex scenarios which require processing large amounts of information simultaneously, such as driving [\[6\]](#), aviation [\[7\]](#) and video surveillance [\[8\]](#), reducing the probability of human errors and speeding up the decision making processes.

Visual attention can be readily identified in two different domains, spatial and temporal, which allow to define three types of computational models for visual attention: spatial, spatio-temporal and temporal [\[9\]](#). Most of existing models consider a spatial component to guide information processing to conspicuous locations or areas of particular interest in a scene. Moreover, visual information in real world is dynamic, so it is equally important to model how it changes over time, in order to update spatial attention based on previously selected locations, which allows modeling visual attention in a spatio-temporal manner, as well as selecting time segments of special importance.

It is also common to distinguish between two families of visual attention models: *stimulus-driven* Bottom-Up ([BU](#)) models, which are based on visual features of the scene, and *goal-driven* Top-Down ([TD](#)) approaches, which take into account prior knowledge or advanced indications [\[10, 11\]](#). Eye movements play a major role in this second type of models, by providing information about which locations are essential for perception and how long they are fixated [\[11, 25\]](#). Although we live in a spatio-temporal reality, the majority of existing computational models for visual attention are [BU](#) and have been built for still images. What is more, the few available [TD](#) methods have been designed for well-determined scenarios, and are not applicable to other contexts. Finally, there is still room for models that take advantage of the demonstrated concurrence between [BU](#) and [TD](#) factors.

1.2 HIERARCHICAL REPRESENTATIONS FOR VISUAL ATTENTION

At present, most of the computer vision-based applications are addressed via feature-based algorithms, which often imply Machine Learning ([ML](#)) and optimization methods. The performance of these applications is highly dependent on features or representations extracted from the visual information beforehand. Hence, features constitute themselves a major and prevailing area of research, which has rapidly evolved in the last few years, from traditional

handcrafted feature engineering to high-level representation learning [26].

1.2.1 Feature engineering

Traditional *feature engineering* involves transforming the domain knowledge of the data into features or properties common to all objects or items considered in a particular task. This process is difficult, time-consuming and requires expert knowledge [27]. Furthermore, the performance of a ML model significantly depends on the quality and quantity of the features obtained. The better the features are, the simpler and more flexible the model needed will be.

Attending to their semantic meaning, image features can be classified into three groups [28]: low-, mid- and high-level descriptors. Low-level descriptors, such as color histograms, texture and shape features, capture either global or local visual properties, and can be directly extracted from the whole image or local regions, respectively. Mid-level features constitute an intermediate step between low and high level, and rely on a global analysis of low-level descriptors, in order to perform annotation or similarity matching tasks, among others. High-level features represent semantic concepts, interpretable by humans, such as faces, cars or any kind of objects, as well as simple categorizations (“urban vs. countryside”, “indoor vs. outdoor”, etc.).

In the field of visual attention, a great effort has been made from multiple perspectives to determine which features better represent those conspicuous areas of the scene for observers [29]. According to the most widely accepted psychological theories [10, 11], there are three features which mainly attract human attention: intensity or luminance contrast, color and orientation. Then, some other attributes, such as motion, shape or faces, might be useful to develop a system that simulates the HVS. Most computer vision researchers have modeled these properties separately, in order to develop computational mechanisms for predicting visual attention. However, only few works have tried to understand how they are combined to perform this function.

In contrast to feature engineering, *feature or representation learning* [26, 30] encompasses those ML techniques that automatically transform the data at the system input into abstract representations which allow an AI to understand the world around us, improving its performance when solving a particular task. Indeed, representation learning often constitutes a preprocessing stage previous to a prediction problem. In the following subsections, we first briefly describe the types of ML covered by the systems presented in this thesis. Afterwards, we introduce representation learning and discuss the issues that should be addressed by a good representation.

1.2.2 Machine Learning

Two definitions of Machine Learning (ML) are usually highlighted. First, the informal, traditional one stated by Arthur Samuel in 1959: “Machine learning is the field of study that gives computers the ability to learn without being explicitly programmed.” [31]. The more recent, by Tom Mitchell, establishes that “A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P , if its performance at tasks in T , as measured by P , improves with experience E .” [32].

- *Experiences E* are related to the nature of ML methods, which define the way they process the information available in a database or collection or examples. Examples are features measured or extracted from the existing objects or events in the world. For instance, we can measure the value of the pixels of an image, but we can also extract its corresponding edge image.

According to the “no free lunch” theorem of David Wolpert and William Macready [33], there is no ML algorithm universally better than any other, and the goal of ML is thus to determine what is the way of experiencing that provides an AI with the most relevant distributions to understand a particular real-world scenario.

The two main different types of ML are supervised and unsupervised learning. In this thesis, we will attend both to their generative and discriminative paradigms from either a probabilistic or a functional perspective, with the purpose of framing elaborated contributions to spatio-temporal visual attention, based on Latent Topic Models (LTMs) (Chapter 3, section 3.4) and Deep Neural Networks (DNNs) (Chapter 5, section 5.3).

A third less common, but also active, research area of ML is known as *reinforcement learning*, which consists on learning suitable actions to perform in order to maximize a reward function. Reinforcement learning has recently been employed in computer vision for object location [34, 35] and image classification [21]. The authors of the latter reference also tested it to automatically play a simple game.

- *Tasks T* determine how ML algorithms process the examples in a database. Some examples of ML tasks are classification, regression and density estimation. Tasks in computer vision such as image classification, retrieval and segmentation have been tackled both in supervised [36–38] or unsupervised [39, 40] experiences. The paramount tasks we aim to solve in this

thesis are spatio-temporal and temporal visual attention modeling: first, we interpret how spatio-temporal visual attention works in several contexts, and then we apply it in a video surveillance scenario for temporal modeling of attention.

- *Performance measures* P evaluate how a ML algorithm works, and are often tailored to the tasks carried out by the system. This evaluation is performed using a test set of data different from the one used for training the system. Metrics to assess spatio-temporal visual attention and temporal attention estimation are in the scope of Chapters 4 and 6, which cover the experiments undertaken during the thesis.

The interested reader is referred to [41–43] for further insight into the different methods and concepts in ML. In an attempt to include all the existing ML algorithms in a unique taxonomy, Goodfellow et al. propose in [43] an easy recipe: combine a database, a function or probability distribution to approximate, an optimization procedure and a model.

Supervised learning

Predictive or supervised learning is the most common experience of ML. Its goal is to predict the value of a *response variable* vector or target \hat{y} given the value of a vector \hat{x} of input features, by means of a model learned from a *training set* $\{(\mathbf{x}_n, \mathbf{y}_n)\}_{n=1}^N$ of N input-output example pairs sampled from a true data distribution D . This can be obtained via a discriminative or a generative model [44, 45]. Figure 1.1 shows the graphical representation of these two approaches. A graphical structure defines the conditional dependence between variables in a model [46].

While in a classification or recognition task \mathbf{y}_n is a categorical variable from a finite set $\mathbf{y}_n \in \{1, \dots, C\}$, the problem is called regression if \mathbf{y}_n involves one or more continuous variables. Considering that models in Figure 1.1 are parametric, we will illustrate their differences by solving the following supervised classification problem.

Given the training set introduced before, let us denote individually $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ as the set of N input vectors and $\mathbf{Y} = \{\mathbf{y}_1, \dots, \mathbf{y}_N\}$ as their corresponding classes, assumed independently sampled from the same distribution D .

On the one hand, we can address *discriminative approaches* from a deterministic or a probabilistic point of view:

- From a *deterministic* or functional perspective, the objective of supervised applications is to learn the mapping $f : \mathbf{X} \mapsto \mathbf{Y}$ between the input feature space \mathbf{X} and the class space \mathbf{Y} . This

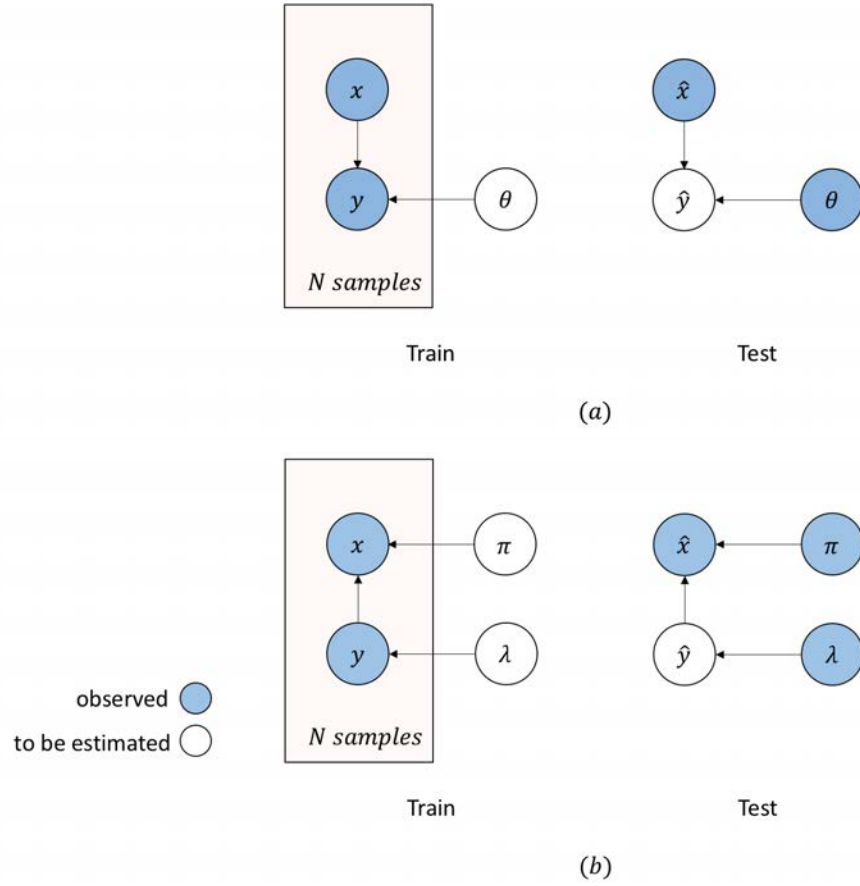


Figure 1.1: Graphical representation of Bayesian (a) discriminative and (b) generative supervised models. Shaded nodes represent N independent input-output example pairs (x, y) , and white nodes indicate parameter vectors to be estimated. Edges show the conditional dependence between variables.

mapping is defined by means of the optimal function f^* , from a set of parametrized functions F , that minimizes the expected value of a loss function L given samples drawn from the true distribution D :

$$f^* = \arg \min_{f \in F} E_{(\mathbf{x}, \mathbf{y}) \sim D} [L(f(\mathbf{x}_n), \mathbf{y}_n)], \quad (1.1)$$

where $E[\cdot]$ stands for the expected value. The loss function L computes the difference between the predicted label $\hat{\mathbf{y}}_n = f(\mathbf{x}_n)$ and the true label \mathbf{y}_n and is chosen according to the task performed, just as evaluation metrics.

Because it is not possible to access to all samples in the true data distribution D , the problem is intractable and can be only optimized considering the available N training samples, assuming that they are Independent and Identically

Distributed (IID), and expecting that they are representative of D , so that it can be expressed as follows:

$$f^* = \arg \min_{f \in F} \sum_{n=1}^N L(f(\mathbf{x}_n), \mathbf{y}_n). \quad (1.2)$$

- Alternatively, from a *probabilistic* point of view, \mathbf{X} and \mathbf{Y} are considered random variables and the objective is achieved by obtaining the conditional distribution $p(\mathbf{y}|\mathbf{x})$ of the target \mathbf{y} given the observations \mathbf{x} . The aim is to determine the class $\hat{\mathbf{y}}$ associated with a new unseen input vector $\hat{\mathbf{x}}$, which implies evaluating the following conditional distribution:

$$p(\hat{\mathbf{y}}|\hat{\mathbf{x}}, \mathbf{X}, \mathbf{Y}) \quad (1.3)$$

This distribution can be represented as $p(\hat{\mathbf{y}}|\hat{\mathbf{x}}, \theta)$, where θ constitutes its corresponding set of parameters. Considering the N independent training samples, the likelihood function is expressed as:

$$p(\mathbf{Y}|\mathbf{X}, \theta) = \prod_{n=1}^N p(\mathbf{y}_n|\mathbf{x}_n, \theta). \quad (1.4)$$

Assuming a prior $p(\theta)$, its product with the likelihood function provides a joint distribution $p(\theta, \mathbf{Y}|\mathbf{X})$ of the parameters θ and the classes \mathbf{Y} given the observations \mathbf{X} . Then, we can obtain the posterior distribution of θ as:

$$p(\theta|\mathbf{X}, \mathbf{Y}) = \frac{p(\theta, \mathbf{Y}|\mathbf{X})}{p(\mathbf{Y}|\mathbf{X})} = \frac{p(\theta)p(\mathbf{Y}|\mathbf{X}, \theta)}{\int p(\theta)p(\mathbf{Y}|\mathbf{X}, \theta)d\theta} \quad (1.5)$$

Marginalizing the predictive distribution with respect to θ weighted by the posterior distribution, we are able to predict $\hat{\mathbf{y}}$ for unseen samples $\hat{\mathbf{x}}$:

$$p(\hat{\mathbf{y}}|\hat{\mathbf{x}}, \mathbf{X}, \mathbf{Y}) = \int p(\hat{\mathbf{y}}|\hat{\mathbf{x}}, \theta)p(\theta|\mathbf{X}, \mathbf{Y})d\theta. \quad (1.6)$$

On the other hand, *generative approaches* learn a probabilistic model of the joint distribution $p(\mathbf{x}, \mathbf{y}|\theta)$ of the feature vector \mathbf{x} and the class label \mathbf{y} , conditioned on a set of parameters $\theta = \{\lambda, \pi\}$. Given a prior probability for the classes $p(\mathbf{y}|\lambda)$ together with a class-conditional density for each class $p(\mathbf{x}|\mathbf{y}, \pi)$, we can express $p(\mathbf{x}, \mathbf{y}|\theta)$ as:

$$p(\mathbf{x}, \mathbf{y}|\theta) = p(\mathbf{y}|\lambda)p(\mathbf{x}|\mathbf{y}, \pi). \quad (1.7)$$

Then, the joint distribution is obtained by drawing from the N independent training samples as follows:

$$p(\mathbf{X}, \mathbf{Y}, \theta) = p(\theta) \prod_{n=1}^N p(\mathbf{x}_n, \mathbf{y}_n | \theta). \quad (1.8)$$

This distribution has to be maximized in order to determine the most probable value of θ . According to Bayes' Rule:

$$p(\mathbf{X}, \mathbf{Y}, \theta) = p(\theta | \mathbf{X}, \mathbf{Y}) p(\mathbf{X}, \mathbf{Y}). \quad (1.9)$$

Hence, maximizing $p(\mathbf{X}, \mathbf{Y}, \theta)$ is equivalent to maximize the posterior distribution $p(\theta | \mathbf{X}, \mathbf{Y})$. The posterior distribution can be then used to evaluate $p(\hat{\mathbf{y}} | \hat{\mathbf{x}}, \mathbf{X}, \mathbf{Y})$ on new samples $\hat{\mathbf{x}}$, in order to make predictions $\hat{\mathbf{y}}$.

The main advantage of generative approaches with respect to discriminative ones is that the joint distribution $p(\mathbf{X}, \mathbf{Y}, \theta)$ models how the data has been generated, which allows to create new synthetic feature vectors $\hat{\mathbf{x}}$ that follow the same distribution than the existing samples. In addition, this implies that these methods can benefit from the mixture of labeled and unlabeled data in semi-supervised frameworks. Well-known examples of generative methods are Gaussian Mixture Models (GMMs) [41] and Hidden Markov Models (HMMs) [47].

Despite this additional capability, traditional learning algorithms showed generative models limited performance to find optimal model parameters, which leads to the true distributions of the data. Hence, discriminative approaches often provide better generalization performances. Support Vector Machines (SVMs) [48] and Neural Networks (NNs) [43] are examples of discriminative methods, as well as the widely used linear and logistic regression models [42], on which trending Deep Neural Networks (DNNs) are based, further explained in section 5.3.

Unsupervised learning

The second main type of ML is called *descriptive* or *unsupervised learning*. Also known as *knowledge discovery*, its objective is to find patterns of interest in the data, given a set of unlabeled inputs $D = \{\mathbf{x}_n\}_{n=1}^N$, by means of hidden variables. It should be noted that, despite the frequent use of this type of variables in unsupervised methods, they can be also arisen by supervised models, such as the Encoder-Decoder Networks (EDNs) introduced in section 5.3.3.

Latent or *hidden* variables \mathbf{z} are representations of the data not directly observed but rather inferred from other variables that can be directly measured. They reduce the dimensionality of the observable data providing a model that explains and makes this information easier to understand. The underlying structures and relations established can be helpful for clustering data into groups, such as in the well-known *K-means* [42] algorithm, and also to reduce the dimensionality of high-dimensional vectors, as in PCA [42].

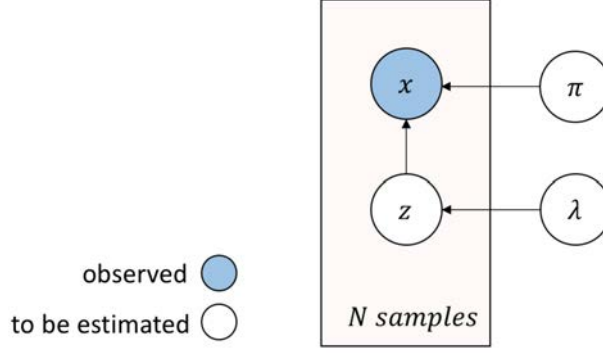


Figure 1.2: Graphical representation of Bayesian directed generative unsupervised models. Shaded nodes represent N independent inputs x , and white nodes indicate hidden variables z and parameter vectors to be inferred. Edges show the conditional dependence between variables.

In this thesis, we will use the well-known Latent Dirichlet Allocation (LDA) [12] directed generative model, which is detailed in section 3.4.2, for visual attention understanding. Directed generative models compute the distribution $p(\mathbf{x}|\theta)$ of the data \mathbf{x} given the parameters $\theta = \{\lambda, \pi\}$ as follows, by means of the prior $p(z|\lambda)$ of the latent variables z and the conditional distribution $p(\mathbf{x}|z, \pi)$ that establishes the relationship between latent and observed variables:

$$p(\mathbf{x}|\theta) = \sum_z p(z|\lambda)p(\mathbf{x}|z, \pi). \quad (1.10)$$

Figure 1.2 shows a basic graphical representation for this type of models. Generative Adversarial Networks (GANs) [49] constitute another recent example of directed generative architecture.

1.2.3 Representation learning

One of the current key challenges of ML is to model and understand complex abstract concepts such as attention or emotion. In the same way that human beings are able to efficiently process information, researchers pursue automatic methodologies capable of separating, given raw input data, useful from irrelevant information, relating it to basic interpretable features (e.g. color, shape), and representing it in a structured or hierarchical way.

According to the outstanding review of Yoshua Bengio et al. about representation learning [26], we should take account of the following aspects in order to achieve a good representation:

- A good representation is one that involves multiple explanatory factors of the observed input, which are useful to solve a particular supervised task.

- A hierarchical organization of explanatory factors is always desirable, which describes the world around us by establishing relationships from less abstract concepts (e.g. movie, film director, actress, etc.), to more abstract ones (e.g. art, entertainment, saliency, etc.).
- Semi-supervised frameworks are helpful to take advantage of the capability of complex unsupervised models, which provide latent representations of the world. They allow to maintain a connection between these hidden representations and our semantic concepts and categories, which contribute to a better understanding of how machines see our reality.
- Making associations between tasks facilitates solving applications for which we do not have enough information annotated or knowledge to interpret the scenarios they imply. Methodologies such as multi-task and transfer learning are in line with this objective [50].
- We should keep in mind the existing correlation between nearby observations, which are often associated with the same semantic or categorical concepts, and change their representations similarly at different spatial and temporal scales.
- Finally, it is important to strive for the simplicity of factor dependencies, which is essential to reach efficient representations.

Nowadays, we can basically identify two paradigms for representation learning: Deep Learning (DL) and Probabilistic Graphical Models (PGM) [20].

Inspired by the hierarchical architecture of the biological neural system, DL methods can be understood as representation methods with multiple layers of representation [30]. Starting with the raw input at the bottom layer, each layer is composed of simple non-linear units that transform its input into a new representation. This representation constitutes the input of a higher, slightly more abstract layer, being the output of the final top layer a lower-dimensional feature at a very high level.

On the other hand, PGM learn a set of latent random variables, and make use of structures that define relationships between these variables, in an attempt to represent distributions over the observed data.

It is worth noting the great contribution of DL representations to perception tasks such as seeing, proved by object recognition and tracking applications [3, 51]; hearing, performed by speech recognition or audio retrieval systems [52, 53]; or reading, carried

out by sentiment analysis and machine translation methods [54, 55]. However, we are still far from a completely understanding of the representations derived from deep architectures. In contrast, PGM have stood out by their ability of thinking and understanding, dealing better with uncertainty than DL, at the expense of performing worse in perception tasks. The integration of both paradigms, which has been denoted as Bayesian Deep Learning (BDL), seems to be the way forward for machine intelligence.

From traditional feature engineering techniques modeling the world around us, to widely adopted feature learning methods at present, we come to the following conclusion: We have been a lot of time trying to teach machines how to define surroundings in our language, by means of handcrafted features based on our experiences. However, machines are not like humans. They have their own language, probably the reason for the success of deep representation learning. We have reached a point where it is quite complex to argue about machine representations or semantics. Their comprehensive capacity sometimes seems to be beyond our scope. Now it is time to understand how machines learn from experiences of this world, shaped like multimedia content, such as audio, images or video sequences, which closely approximate our reality. Will machines perform in tasks like visual attention in a similar way than humans? Will be necessary to let machines choose first their own paths to solve these tasks, and then develop translation mechanisms to interpret them? Without doubt, we are at the beginning of a new promising and exciting era for AI.

1.3 GOALS AND CONTEXT OF THE THESIS

In this section, we discuss the main focus of this thesis, which concerns the study and development of hierarchical representations for spatio-temporal visual attention modeling and understanding.

Specifically, the thesis makes the following two main contributions towards our goals:

1. We introduce a hierarchical generative probabilistic model for context-aware visual attention modeling and understanding.

Our first approach, which we have called visual Attention TOpic Model (ATOM), models visual attention in the spatio-temporal domain by considering the existing concurrence between BU and TD factors.

2. We develop a deep network architecture for visual attention modeling, which is oriented to be applied in a video surveillance scenario.

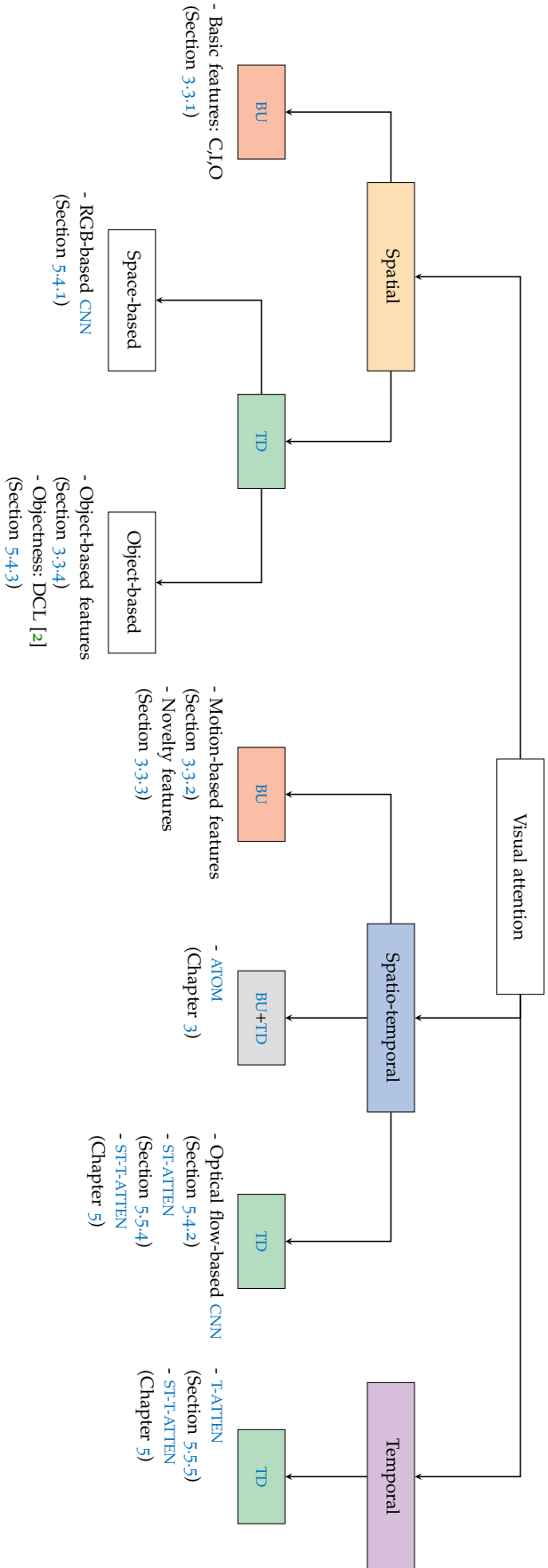


Figure 1.3: Visual attention features and representation models presented in the thesis, classified according to the spatial, spatio-temporal or temporal dimension in video sequences where they are modeled. Sections and chapters where they are described are indicated next to each item.

We have called our second proposal as Spatio-Temporal to Temporal visual ATtention NETwork (**ST-T-ATTEN**). It first estimates **TD** spatio-temporal visual attention, which ultimately serves for modeling visual attention in the temporal domain.

Our particular contributions associated with both systems are mentioned in the next section, which also summarizes the main content of each chapter of the thesis.

In order to contextualize our contributions with respect to the existing types of visual attention models (Section 1.1) and the different methodologies for visual information representation (Section 1.2), we include two diagrams. In both diagrams, next to each item, sections or chapters where features and representation models are explained are indicated.

On the one hand, Figure 1.3 outlines the features and representations for visual attention guidance covered throughout the thesis. They are classified according to the dimension (spatial, spatio-temporal or temporal) that they model in video sequences, and constitute a wide and complete framework to give context to our proposals. Reading from left to right, the diagram goes from spatial through spatio-temporal to temporal representations, which are classified according to the two main families of visual attention models introduced in section 1.1: **BU** and **TD** implementations. In addition, for the case of **TD** spatial methods, the diagram differentiates between space-based features, which rely on the information drawn from eye fixations, and object-based features, related to salient objects in the scene.

On the other hand, features and representation models for visual attention guidance are depicted in Figure 1.4, according to the feature engineering or feature learning processes that they involve. Within the context of feature learning, we distinguish between shallow models, which are those with one or few levels of representation, and deep models with multiple layers of representation, representing the new Computer Vision paradigm. Furthermore, the diagram also reflects the difference between generative and discriminative methods.

1.4 STRUCTURE OF THE THESIS AND CONTRIBUTIONS

In this section, we present the structure of the dissertation, introducing our main scientific contributions in the corresponding chapters.

Chapter 2 makes a review of the most relevant and recent related work in perception and visual attention from a multidisciplinary perspective.

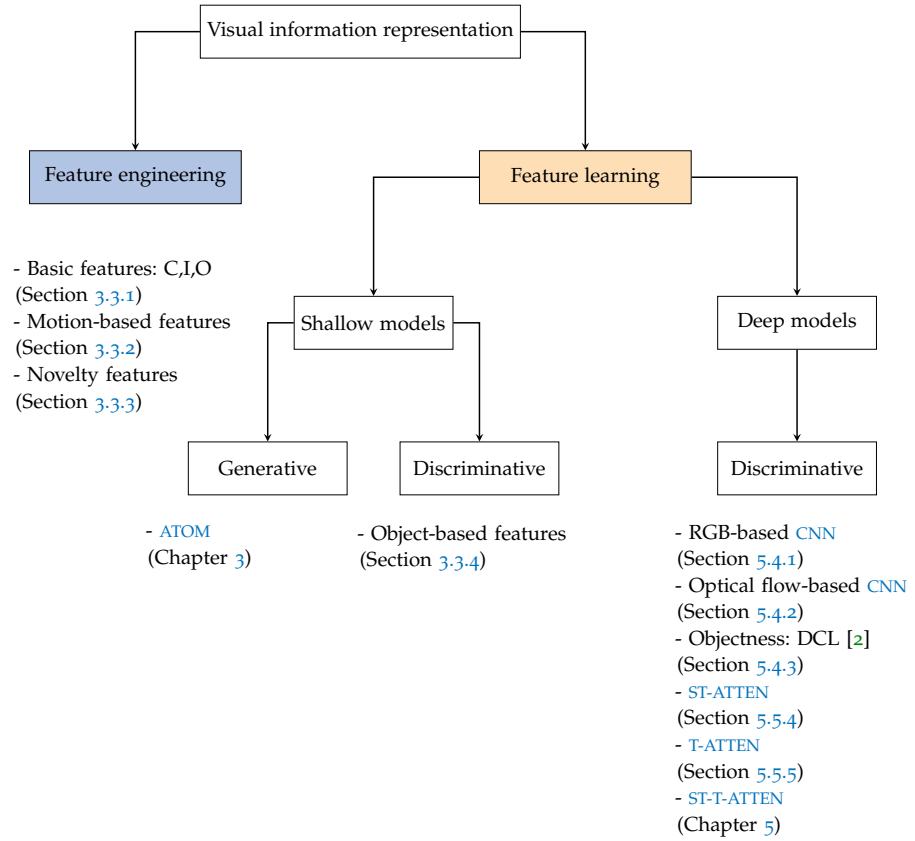


Figure 1.4: Visual attention features and representation models covered by the different systems presented in the thesis, classified according to the types of processes that they involve. Sections and chapters where they are explained are indicated next to each item.

- First, we review visual attention from a neurophysiological perspective, mainly describing the mechanisms of the eye and the brain for visual selection and representation.
- Then, we introduce the most outstanding psychological theories of visual attention, as well as some noticeable studies on the role of eye movements.
- Finally, we summarize some of the existing computational models of visual attention, attending to either Bayesian models or deep learning-based approaches close to the contributions of the thesis.

Then, we have developed two computational systems for visual attention, which constitute the main contributions of this thesis.

Chapter 3 introduces our first proposal: a generative probabilistic framework for spatio-temporal visual attention modeling and understanding.

We first briefly discuss the related work in computational visual attention modeling. The model proposed, which we have called

visual Attention TOpic Model ([ATOM](#)), is generic, independent of the application scenario and founded on the most outstanding psychological studies about attention. Drawing in the well-known Latent Dirichlet Allocation ([LDA](#)) [[12](#)] method for the analysis of large corpus of data and some of its supervised extensions [[13](#), [14](#)], our approach defines task- or context-driven visual attention in video as a mixture of latent sub-tasks, which are in turn represented as combinations of low-, mid- and high-level spatio-temporal features.

In particular, we make the following contributions in this chapter:

- We introduce feature engineering for visual attention guidance, providing a wide set of handcrafted features, which are later used in our experiments. Starting from basic and novelty spatio-temporal low-level features, such as color, intensity, orientation or motion, we move on to describe and model some mid- and high-level features related to camera motion estimation and object detection.
- Then, our algorithm incorporates an intermediate level formed by latent sub-tasks, which bridges the gap between features and visual attention, and enables to obtain more comprehensible interpretations of attention guidance.
- Moreover, we generate a categorical binary response for each spatial location to model visual attention. This allows to automatically align the sub-tasks discovered to a binary response by means of a logistic regression, which fully corresponds to the definition of human fixations.

Chapter [4](#) provides an in-depth analysis of [ATOM](#). For that purpose, our model is used for context-driven visual attention modeling and understanding in two large-scale video databases annotated with eye fixations: CRCNS-ORIG [[15](#)] and DIEM [[16](#)]. We illustrate how our approach successfully learns hierarchical guiding representations adapted to several contexts. Moreover, we analyze the models obtained, as well as perform a comparison with quite a few *state-of-the-art* methods.

Chapter [5](#) describes our second proposal: a deep network architecture that goes from spatio-temporal visual attention prediction to attention estimation in the temporal domain. The system proposed, which we have named Spatio-Temporal to Temporal visual ATtention NETwork ([ST-T-ATTEN](#)), models visual attention over time as a fixation-based response.

First, we review the most relevant and recent works in visual attention estimation applying deep learning-based architectures. Then, we introduce the fundamental hypothesis of the second part of the thesis: attention in the temporal domain can be predicted

using the dispersion of gaze locations recorded from several subjects.

Indeed, visual attention in the temporal domain can be understood as a filtering mechanism, which allows to select time segments of special importance in video sequences. Hence, it could be used to prevent human errors and speed up decision making processes in real applications which require watching large amounts of visual information, such as the task of video surveillance.

We make the following particular contributions in this chapter:

- We describe three feature learning architectures for visual attention guidance, which provide input feature maps to our system: RGB-based spatial, optical flow-based and objectness-based networks.
- We propose a frame-level fixation-based temporal ground-truth, which is computed attending to the dispersion at fixation spatial locations from several subjects. Furthermore, we validate the fundamental hypothesis introduced above. We will use this variable to train our models to estimate attention in the temporal domain.
- Our proposed **ST-T-ATTEN** is built on the combination of two modules: 1) A Spatio-Temporal visual ATtention NETwork (**ST-ATTEN**) for spatio-temporal visual attention estimation, which consists on a Convolutional Encoder Decoder (**CED**) [17] network; 2) A Temporal ATtention NETwork (**T-ATTEN**) for modeling visual attention in the temporal domain, based on Long Short-Term Memory (**LSTM**) [18] units, widely used for time series forecasting.

Chapter 6 describes the experiments conducted to validate the different configurations proposed for the **ST-T-ATTEN** modules. We make use of the BOSS [19] database, which contains videos recorded in a railway transport context with different anomalous events, with the aim of determining the optimal configuration for the whole **ST-T-ATTEN** proposed, as well as motivating its use as an information filtering mechanism in a video surveillance application.

Finally, in Chapter 7, we summarize the conclusions drawn from the main contributions of the thesis, which serve to outline future lines of research.

A MULTIDISCIPLINARY PERSPECTIVE ON VISUAL ATTENTION

2.1 INTRODUCTION

During the 1970s, scientists from several disciplines began to show a great interest in understanding how optical images could be processed to extract useful information about the environment [56]. This resulted in the emergence of vision science.

Vision science [56] is defined as an interdisciplinary branch from cognitive science, which is devoted to the study of visual mental states and processes from different, compatible and complementary perspectives. All of them talk about vision in the common language of computation, by accepting that it may take place not only in living organisms, through eyes and brains, but also when information from cameras is processed in ad-hoc programmed computer systems.

Over the past few decades, psychologists have tried to explain visual perception through a vast amount of theories and models. Moreover, neurophysiologists have made experiments to monitor neuron activity. Furthermore, computational neuroscientists have built neural network architectures to simulate how these neurons represent and react to visual stimuli. Drawing on these findings, computer vision scientists have sought to develop computational models and algorithms which automatically address the cognitive functions involved in attention.

Indeed, a great world full of visible information is opened to us, and the Human Visual System (HVS) has the paramount responsibility of dealing with attentive processes. Due to the limited capacity of the brain to process such a big amount of sensory input, attention involves the inherent search operations that reformulate and optimize generic perception and cognition problems so that they become tractable [57]. Eye movements allow acquiring and tracking visual stimuli, unconsciously highlighting the most conspicuous [58] [59] areas in a particular context, or willingly selecting those that aid to solve a particular task [60].

This thesis presents a framework for visual attention estimation and understanding from a computational view, not only applying existing computer vision techniques, but also contemplating psychological arguments. This chapter makes a review of the most relevant and recent related work in perception and visual attention, bearing in mind all the perspectives differentiated above. The purpose of this *state-of-the-art* is thus to provide a broad overview of the visual attention research by identifying the basis of our work.

CHAPTER OVERVIEW

Starting from the structure and the processes involved in the eyes and the brain, visual attention is reviewed from a neurophysiological perspective in Section 2.2. We discuss the difference between *overt attention*, which implies eye movements and fixations, and *covert attention*, which is more related to the mechanisms of the brain for visual selection and representation. Then, Section 2.3 introduces the most outstanding psychological theories of visual attention, which allude to *early representation* features that guide the attention of observers. Moreover, we also cover some noticeable studies on the the role of eye movements. Finally, we summarize some of the existing computational models of visual attention in Section 2.4, mainly attending to those approaches close to the contributions of the thesis.

2.2 NEUROPHYSIOLOGICAL BASIS OF VISUAL ATTENTION

This section describes the structure of the HVS, attending to the regions of the eye and the nervous system that take part in the process of visual perception. According to the Professor Stephen E. Palmer [56], “*visual perception is an information extraction process that involves the acquisition of knowledge about objects and events in the environment*”. This information comes from the light that is emitted or reflected by objects.

It should be noted that HVS has an extraordinary ability to select only the necessary information in order to interact with a given scenario, being able to infer the rest with sufficient accuracy. Hence, given an image, vision implies a heuristic process to infer the most likely environmental condition that could have produced it.

2.2.1 Human Visual System: eye and brain

Both eyes and brain are essential for visual perception. The complete eye-brain system must perform adequately to obtain trustworthy visual data.

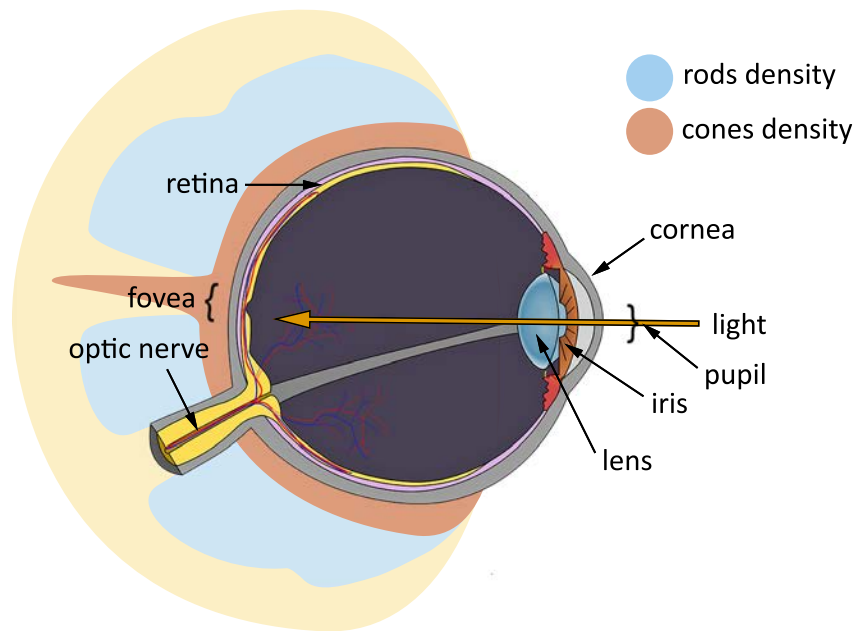


Figure 2.1: Diagram of the human eye. Rods and cones densities are drawn around the retina in blue and red, respectively. Adapted from Wikimedia Commons [61].

The eye

The structure of the human eye is shown in Figure 2.1. Humans have two approximately spherical eyes. They are situated at two holes in the skull called *eye sockets*, which are placed at about the horizontal midline of the head. Eyes are moved by six small and strong *extraocular muscles*, which are responsible for eye movements, allowing to scan different regions of the visual field. They are monitored by several nuclei in the brain stem, via the oculomotor neurons.

Several parts of the eye carry out optical functions. First, eyes collect the light that enters through the *cornea*. The light crosses an opening in the *iris* called *pupil*, behind which the *lens* is located. Finally, incoming light projects an image onto the *retina*, a curved surface at the back of the eye. The retina is composed of more than 100 million light-sensitive cells, known as *photoreceptors*, which transform light into neural activity. There are two types of photoreceptors: rods and cones. As can be seen in Figure 2.1, rods, which are longer and more numerous (about 120 million), are located everywhere in the retina except at its center. They are highly sensitive to light, so they allow us to see at low light levels or *scotopic conditions*. In contrast, most of the cones, which are shorter and fewer (8 million), are clustered in the *fovea*, situated at the center of the retina. They are much less sensitive to light, used under normal lighting or *photopic conditions*, and also in all experiences of color.

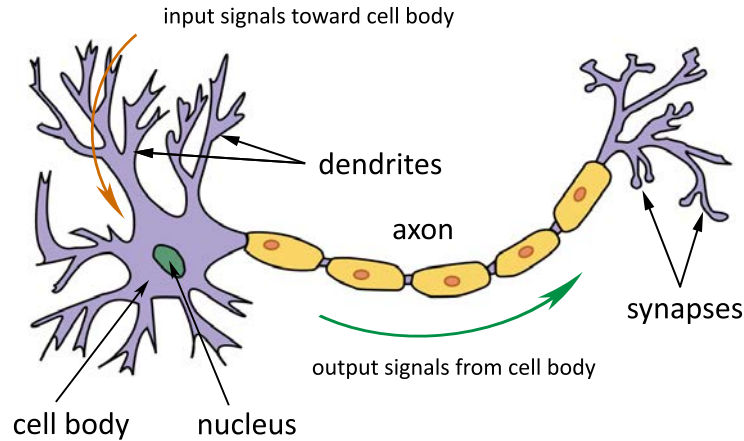


Figure 2.2: Schematic diagram of a biological neuron. Adapted from [56, 62].

Light comprises photons, many small units of energy that produce electrical changes in photoreceptors, and the information travels via the *optic nerves* to the visual centers in the brain.

The brain

From the fovea in each eye, the optic nerves cross over to the opposite side of the brain, leading the information from the left half of the visual field to the right side of the brain and vice versa. The brain processes this information, in order to make it useful for observers.

Neurons constitute the basic computational cells of the brain. The human brain is composed of around 100 billion neurons. As shown in Figure 2.2, a neuron first receives electrical signals coming from other neurons across the *dendrites*. Within the *cell body*, where the *nucleus* is located, these inputs are converted into a series of output spikes that are propagated through its *axon* to other neurons. The *firing rate* of the neuron determines the frequency of the spikes. Finally, *synapses* connect the axon to the dendrites of the following neurons.

There are two pathways on each half of the brain. One of them arrives to the *Superior Colliculus (SC)*, which seems to be involved in the control of eye movements by processing information related to the location of objects in the world; the second and larger pathway goes to the *occipital* or *Primary Visual Cortex (V_1)*.

Nowadays, we have an evidence about the function of the *occipital*, *parietal* and *temporal lobes* of the visual cortex, which are identified in Figure 2.3. The cortical processing begins at V_1 cells, where *spatial receptive fields* respond to visual stimuli, so that a mapping of the information from the retina is produced. According to the scale-space theory of computer vision [63], receptive fields encode simple visual patterns of light, such as oriented edges or color

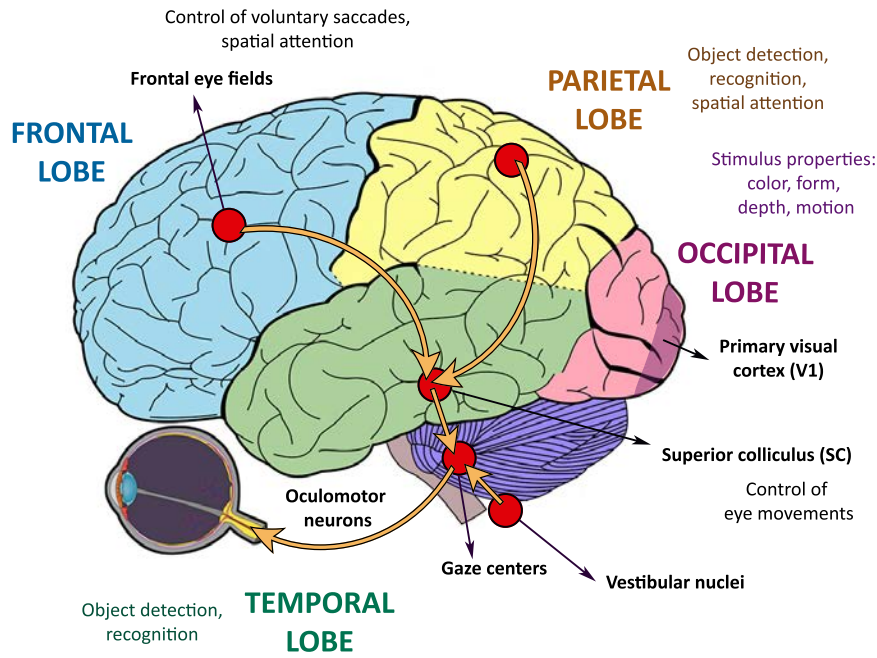


Figure 2.3: Diagram of the human brain. Arrows indicate the connection between the eye and the principal areas in the brain involved in the visual attention process: frontal eye fields and posterior parietal cortex, which guide spatial attention, and the Superior Colliculus (SC), which controls both eye movements and covert shifts of attention. Eye and brain diagrams taken from Wikimedia Commons [61, 67].

blotches, which constitute the first stages of visual processing, by reflecting the symmetry properties of the world that surrounds us. Moreover, a reduced set of operations over these simple patterns allow to obtain a wide variety of complex underlying representations for visual perception. Convolutional Neural Networks (CNNs) [64] emerged in an attempt to reproduce the function of receptive fields, recently demonstrating an astonishing performance in a lot of applications. They will be described in section 5.3.2, as the basis of our contributions for visual attention estimation in the temporal domain in Chapter 5.

The information from V_1 is then projected to other parts of the occipital lobe, and also to areas of the parietal and temporal ones. Some studies [65, 66] suggest that these different regions involve small maps where several properties derived from the retinal stimulation are coded in parallel, such as color, form, depth and motion. Regarding the parietal and temporal cortex, they seem to be responsible for the identification and location of objects, respectively.

2.2.2 Visual attention

Visual perception is inherently selective [56]. We are able either to globally process the information in a scene, or to focus our vision more on some particular objects. We sometimes even attend locally to their specific parts or properties, depending on their importance in the activity we are performing. In addition, we choose quite automatically where to fixate our vision next. For instance, although there may be a lot of appliances in a kitchen, we immediately head towards the fridge if we want to drink something, without looking at the toaster or the dishwasher.

All these strategies for selecting and processing information in the visual field are related to attention. Two different acts of visual attention can be distinguished. First, attention is called *overt* if it is external and observable by others, implying eye movements to fixate from one object to another. Different fixations of a context contain useful information which is shaped like visual images on the retina. Then, part of the fixated information is selected by *covert* attention to be fully processed. Covert selections are, conversely, internal and unobservable by others. They do not imply eye movements, but allow to shift our gaze to peripheral Regions Of Interest (ROIs), chosen from the information processed.

Overt attention: eye movements and fixations

Visual attention can be thus described, on the one hand, as a temporal process that involves a sequence of eye fixations preceded by different types of eye movements. This results in a series of instantaneous spatial locations of the visual axis called *gaze points*.

There are four basic types of *eye movements* [68]. If we look at still images or static objects, we mainly perform saccadic movements to scan over them. *Saccades* are very quick (20-40 ms) voluntarily or involuntarily ballistic jumps between two points of fixation. During a saccade, both eyes drift in the same direction. Moreover, the trajectory of a saccade cannot be changed when the eyes are in motion.

In real-life situations, where either the viewer or the objects are moving, three more types of eye movements can be found:

- *Smooth pursuit movements*, which are slow in comparison with saccades, are used to track the position of moving objects. The ability of the HVS to take clear images from tracked objects depends on how fast they move, being less accurate at speeds higher than 30 degrees per second, when subjects start to use saccades to follow objects.
- *Vergence movements* allow HVS to fixate objects located at different depths. In this type of movement, eyes move in

opposite directions and have an angle of convergence that depends on the distance of the target from the observer. Eyes seeing nearby objects strongly converge.

- *Vestibular movements* are controlled by the vestibular system in the inner ear, and contribute to keep the target fixed on the fovea when the head changes its position and orientation (ego-motion). In this situation, eyes compensate the ego-motion by moving in the opposite direction of the head, normally at its same speed.

When the eyes stop examining the scene, *fixations* take place and the *HVS* takes comprehensive information about what is being looked at. Although it is not often mentioned, fixations are not directly measurable, but composed of minute microsaccades, tremor and drift movements that focus the eyes on the target, generating multiple gaze point samples. They have a particular duration, usually between 50-600 ms, and can reveal meaningful information about attention and understanding. Given a specific context, the time to first fixation in a conspicuous location or target is short, while a long fixation duration may suggest a greater effort to make sense of a stimulus, or an appealing one.

All these movements have disparate neural mechanisms, which are spread in different areas of the brain, as shown in Figure 2.3. First, frontal eye fields in the frontal cortex control the voluntary saccades. On the other hand, both smooth pursuit and vergence movements require visual feedback, so they are monitored by means of information from the motion channels in visual cortex and binocular disparity channels in occipital cortex, respectively. Finally, vestibular movements result from disturbances in the fluid of the semicircular canals of the inner ear. These are monitored by the vestibular system, which connects to the oculomotor neurons to provide them with the correct eye velocity signal. The *SC* is also involved in the control of eye movements, as was mentioned above.

Covert attention: Relation with overt attention

On the other hand, *visual attention* concerns a set of complex covert processes that aid an observer to select and gather the most outstanding information within the visual field, with the aim of successfully solving a cognitive problem in a particular environment.

Covert attention is usually directed at the *ROI* fixated by the eyes. Professor Stephen E. Palmer [56] makes an interesting metaphor in order to explain the relationship between eye movements and covert attention: "*Attention is like an internal eye that can be moved around to sample the visual field much as the eye can be moved around to sample the visual world.*" What is more, there is evidence that eye movements

usually follow attentional movements. Thus, covert attention, viewed as the primary function of visual selection, controls overt saccades, which play an important but supporting role, driving them to the appropriate locations and enhancing the perception of events happening there.

Therefore, it is known that there exist a strong correlation between the areas in the brain that controls eye movements and covert attention: SC controls both eye movements and covert shifts of attention, while frontal eye fields and the posterior parietal cortex are responsible for guiding spatial attention. Arrows in Figure 2.3 indicate the connection between the eye and the principal areas in the brain involved in the visual attention process.

Finally, many studies agree in saying that this frontoparietal network may involve an attentional *priority map* [69], which represents items in the visual world according to their importance in a particular situation. We will discuss more about the utility of this representation of visual attention in the next section, now from a psychophysical perspective.

2.3 PSYCHOPHYSICAL THEORIES OF VISUAL ATTENTION

In this section, we gather the statements of the most outstanding psychological theories of visual attention, which establish relevant features for the perception of objects and support the majority of existing computational attention systems [70].

Two fundamental theories have been the most influential: the Feature Integration Theory (FIT) [10] and the Guided Search Model (GSM) [11]. Based on the foundations of these studies, we can differentiate between two main families of visual attention models: Bottom-Up (BU) or *stimulus-driven* and Top-Down (TD) or *task-oriented*. In addition, it is also worth mentioning the importance of eye movements in scene perception, as explained in the famous classic study of Yarbus [25].

In addition, we refer for the first time to several aspects that have been considered in the design of the probabilistic model for visual attention understanding presented in Chapter 3.

2.3.1 The Feature Integration Theory

Treisman and Gelade introduced the Feature Integration Theory (FIT) [10] in 1980, which states that several features or attributes are identified early, automatically and in parallel across the visual field, while objects are registered separately as a conjunction of these features at a later serial stage. The model is depicted in Figure 2.4(a). A master map of location results from this combination of attributes, which indicates where the objects are, prior to their recognition.

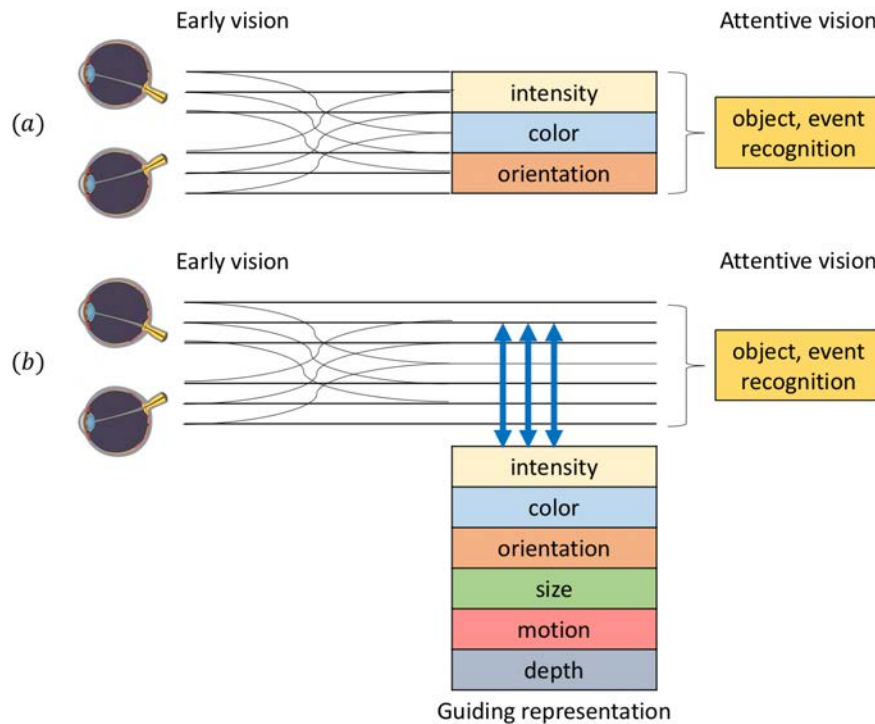


Figure 2.4: Diagram representations of (a) the Feature Integration Theory (FIT) [10], where visual attention is modeled as a combination of intensity, color and orientation pre-attentive features, and (b) the Guided Search Model (GSM) [11], which introduces “guiding representation” as a control mechanism before object recognition, and also mentions new features such as size, motion or depth. Adapted from [71]. Eye diagrams taken from Wikimedia Commons [61].

In light of this theory, Treisman discussed the difference between feature and conjunction search processes, when looking for a target in a scenario full of distractors. In a *feature* or parallel search process, we seek through distractors that differ from the target by a unique feature. In contrast, Treisman asserted that a *conjunction* or serial search assumes a more complex and time-consuming task, since the distractors have one or more features in common with the target. The more attributes aid to discern the target from the distractors, the easier the search for a target is. Nevertheless, if we know the features of the target in advance, conjunction search performance can be improved by inhibiting the features which are exclusive from distractors.

2.3.2 The Guided Search Model

The standard FIT defines parallel and serial search as autonomous processes that cannot share information between them. By contrast, Wolfe’s Guided Search Model (GSM) 2.0 [11] from 1994 claims that

the [FIT](#) attentive serial search has to be guided by useful information in the preattentive parallel processes, which divided their corresponding set of stimulus into distractors and candidate targets.

[GSM](#) thus supports that attention can be guided towards specific targets by modulating gains associated with low-level features. Indeed, visual search is a continuous process and, consequently, the information from the parallel processes to the serial process can be updated over time.

Subsequent works by Wolfe [[71](#), [72](#)] present the idea of “guiding representation” or guidance as a control device located to one side of the main pathway from early vision to object recognition, as shown in [Figure 2.4\(b\)](#). It controls the access to the attentional bottleneck, so the guidance is abstracted from the main pathway despite of not being part of the pathway itself. Thus, the way we see stimulus in the world is different from the representations upon which guidance is founded. Rather than altering the stimulus such as filters would do, this module guides attention as a CCTV operator working at a public building (e.g. a train station or a university) would do. Based on an abstract representation of some notions (e.g. threat, suspicious object), the operator selects some parts of the scenario that receive more attention than others.

Two ways of guidance are possible, Bottom-Up ([BU](#)) and Top-Down ([TD](#)), which correspond to the two main types of computational visual attention systems. [BU](#) attention is fast, involuntary and mainly based on characteristics of the visual scene (*stimulus-driven*) such as color, orientation, motion or depth. By contrast, [TD](#) attention is slow, voluntary and determined by cognitive phenomena like knowledge, expectations or advanced indications (*goal-driven*).

Guidance can be ultimately represented as an activation map. An activation map drawn only by [BU](#) signals is a Saliency Map ([SM](#)) [[73](#)]. Instead, [TD](#) guidance provides a priority or Visual Attention Map ([VAM](#)) that sorts by relevance the items existing in the visual field for selective attention and recognition [[29](#)].

It is also worth noting the concept of scene guidance incorporated in the version 4.0 of [GSM](#) [[74](#)], which results in a non-selective pathway in the visual search process: Observers are able to determine very rapidly the global properties or ‘gist’ of a given scene (e.g. a highway with intense traffic) before they selectively attend to the most conspicuous objects (e.g. a damaged car).

Hence, guidance is not based directly on the information provided by early visual processes but on a coarse and contextual representation derived from them. This interpretation of visual attention supports the main assumption of our model in [Chapter 3](#) and opens the door to the inclusion of an intermediate layer mapping the low level stimuli to an intermediate representation.

2.3.3 *What are the attributes that guide attention?*

A significant number of authors have contributed to determine a wide set of features that might drive attention, as outlined in the excellent review by Wolfe [29, Chapter 2]. In this survey, features are grouped depending on their consensual reliability, which is determined taking into account both the number of studies that support them and the convergence of their demonstrations. Some of these features are more effective than others, not only because of what they measure but also depending on the quality of the visual support where they are computed.

First, Treisman and Gelade's FIT [10] mentioned three basic features: intensity or luminance contrast, color and orientation. These are considered undoubted attributes, so they constitute the basis of almost all the existing models that explain visual attention.

Additionally, Wolfe's Guided Search Model (GSM) 2.0 [11] enumerated other attributes that humans can perceive efficiently and thus could be also considered salient in a scene: motion, scale or size, shape and depth. According to McLeod et al. [75], motion speed and direction could be represented separately. Moreover, it is also unclear if shape should be defined as a whole or, alternatively, as a family of simpler attributes such as curvature, line termination or closure. Then, depth aid to modulate features like size. Near objects stand out from far ones, which seem to be smaller. Wolfe later extended this list [71], raising doubtful or complicated cases such as novelty, faces or other semantic categories (e.g. 'car' , 'dog').

Based on psychologists intuition, we have selected some of these attributes as input to our approach for visual attention understanding, with the purpose of appraising their utility in diverse contexts. Chapter 3 introduces this set of features, as well as the image processing techniques used for their extraction.

2.3.4 *Eye movements*

The role of eye movements in scene perception had already been studied before the introduction of the perception theories referred above. According to the revision of previous research on high-level scene perception made by J. M. Henderson and A. Hollingworth [76], one can figure out what are the procedures that control *where* and *how long* each fixation point tends to remain centered at a particular location to provide a complete understanding of scene.

Yarbus classic study of 1967 [25] showed that, although first few fixations in a scene seem to be controlled by global characteristics, positions of later fixations are not random but landed on regions that are useful or essential for perception. Eyes are either driven by TD factors that direct fixations toward task-driven informative locations

(e.g. cooking, driving) or led to low-level image discontinuities called salient regions (e.g. bright regions, edges). The time the eyes remain in a given region also depends on its visual and semantic properties.

The experience of a complete and integrated visual world is thus based on an abstract representation that covers general information about the scene combined with perceptual information arisen from fixations. By examining eye movements, it could be possible to infer the underlying factors affecting fixations or task at hand, even to interpret observer thoughts [77].

This is the purpose of our model in Chapter 3, which introduces an intermediate level between feature extraction and visual attention computation stages based on the information drawn from fixations. This level consists of latent sub-tasks that can be used to determine why some locations are more conspicuous than others. Thus, rather than directly learn a predictor of human attention over low-level visual features, our method provides a hierarchical interpretation of visual attention, advantageous for further comprehensive analysis.

The experiments conducted by Yarbus have motivated researchers to assess the possibilities of eye tracking for assistance in real applications such as industry control [78], health-care [79] and video surveillance [8]. Experts in these applications have to process a large amount of visual information at the same time, which implies a high cognitive effort that might be reduced by modeling fixations behavior along time. Indeed, it has been observed that there is a strong correlation between fixation patterns of different viewers performing the same task, specially during an important or suspicious event [78, 80].

The latter is the ultimate objective of our system in Chapter 5, which is trained to model a temporal attention response arisen from fixations dispersion across viewers.

2.4 COMPUTATIONAL MODELING OF VISUAL ATTENTION

So far, we have reviewed how vision and visual attention work in humans. Similarly, vision in camera systems depends on the interaction among light, surfaces that reflect light, and a visual system that can detect light. Furthermore, both human eyes and cameras share certain physical similarities, and the same optical function: they gather light reflected from objects in order to obtain a sharply-focused image.

However, while humans acquire knowledge about their surroundings, being able to respond to a given situation, cameras have no perceptual capabilities, so they can not interpret their recordings at all. As cameras are not able to deal with attention processes, computer vision researchers have developed automatic

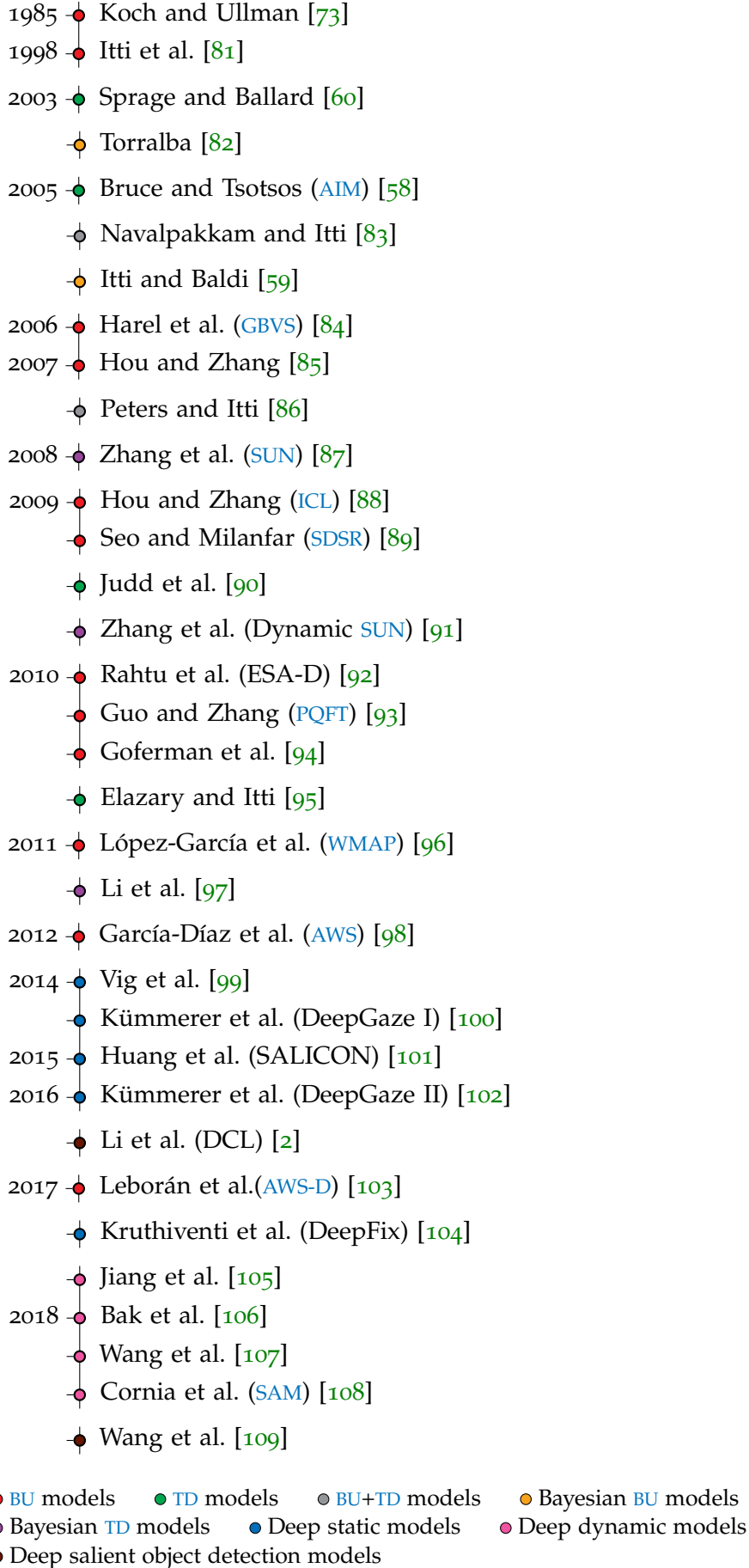


Figure 2.5: Chronological timeline of the visual attention models in the *state-of-the-art* reviewed in this thesis.

systems which efficiently determine the most appealing regions from images or videos, by means of Visual Attention Maps (VAMs).

This section provides a review of some of the existing visual attention algorithms in the *state-of-the-art*, referring both to BU and TD methods. We will specially focus on those that are closely related to our approaches, and also on promising deep learning-based architectures. All the models explained in the following paragraphs are summarized in Figure 2.5 in chronological order. For further information, a exceptional detailed survey on visual attention modeling is presented by Borji et al. [9].

2.4.1 Bottom-up versus top-down approaches

On the basis of Treisman and Gelade's FIT [10], Koch and Ullman [73] presented in 1985 a design to combine early vision features, and defined the concept of SM, which was subsequently mentioned in Wolfe's GSM as a mechanism to model local visual attention driven by the set of visual stimuli in the scene.

The first implementation and verification of a BU model, performed by Itti et al. [81], and incorporating color, intensity and orientation features, would nevertheless come more than ten years later. After that, Harel et al. [84] proposed a saliency algorithm based on graphs, which extracted the same features at different scales. These two representations are the most frequently employed in the literature due to their good performance in a variety of applications [37, 110, 111].

The great majority of visual attention models developed are BU approaches, as the ones explained below. It is also surprising the lack of use and modeling of spatio-temporal and high-level features to address visual attention in real scenarios or videos, although almost half of the methods found in the literature are dynamic ones.

- *Itti et al., 1998* [81]: Related to FIT [10] and in accordance with the biological architecture proposed by Koch and Ullman [73], it was the first existing implementation of a spatial saliency model. It decomposes visual input into a set of topographic feature maps of color, intensity and orientation. The three conspicuity maps are normalized and summed constituting the SM, whose maximum determines the most salient image region.

Additionally, the authors propose a later stage to determine the order in which fixations may occur based on the SM obtained. To this end, the SM is modeled as a two dimensional (2D) layer of leaky integrate-and-fire neurons which feeds into a biologically-plausible Winner-Take-All (WTA) neural network. Each SM neuron excites its corresponding WTA neuron, until the most salient ("the winner") fires. This causes the Focus of

Attention (FOA) to be shifted to the location of the winner neuron. After this, all neurons are reset, and an Inhibition of Return (IOR) mechanism is activated, which either allows the next most salient area to become the following winner or prevents the FOA from shortly reaching an attended region.

- *Harel et al., 2006 [84]*: This spatial BU saliency algorithm based on graphs extracts the same multi-scale features than Itti et. al [81]. Graph-Based Visual Saliency (GBVS) then computes a fully connected graph over all feature map grid locations, by assigning weights between each two nodes that are proportional to the similarity of feature values and their spatial distance. The obtained graphs proceed as Markov chains that define an equivalence relation either between nodes and states or edge weights and transition probabilities. Their associated equilibrium distribution results in activation maps, whose combination gives rise eventually to the SM.
- *Hou and Zhang, 2007 [85]*: Hou and Zhang first proposed a spectral residual saliency model for static images. On the basis that statistical singularities in the spectrum may determine the anomalous regions of a given image, they derive its amplitude and phase. Then, they compute the log-spectrum of the down-sampled image. After that, the spectral residual is obtained by multiplying by a local average filter, and subtracting the result from the original version. Finally, the SM is built in the spatial domain by using the Inverse Fourier Transform (IFT) and squaring the value of each spatial location.
- *Hou and Zhang, 2009 [88]*: Their next approach to saliency, called Incremental Coding Length (ICL), is a principle that tries to maximize the overall entropy gain of a given set of sample visual features, now both in dynamic and static settings. Salient cues are those unexpected feature values that produce an entropy gain in the perception state. Therefore, features with large coding length increments will allow to reach attention selectivity.
- *Seo and Milanfar, 2009 [89]*: Given an image or video, the Saliency Detection by Self-Resemblance (SDSR) spatio-temporal framework computes local regression kernel descriptors. Each pixel or voxel in the SM indicates the statistical likelihood of saliency of a feature matrix given its surrounding feature matrices. According to the authors, the use of Local Steering Kernels (LSK) as features instead of conventional filter responses captures the underlying local structure of the data, even in the presence of significant distortions.

- *Rahtu et al., 2010* [92]: This spatio-temporal method, named ESA-D, proposes a combination of a statistical saliency measure based on contrasts in illumination, color and motion, together with a Conditional Random Field (CRF) model for salient object detection. The motivation of the authors to determine the most appealing areas by minimizing an energy function derived from a CRF is that saliency estimation objective is usually to achieve an object-level segmentation instead of a pixel-level one.
- *Guo and Zhang, 2010* [93]: The authors proposed the Phase spectrum of Quaternion Fourier Transform (PQFT) method, which carries out a quaternion representation of video frames by means of intensity, color and motion features.
- *Goferman et al., 2010* [94]: This context-aware saliency model takes into account four psychological principles of human attention: 1) local low-level features, such as color and contrast; 2) global factors, which stand out features that differ from the norm; 3) visual organization rules related to forms and 4) a high-level face detector.
- *López-García et al., 2011* [96]: The authors proposed the Weighted Maximum Phase Alignment (WMAP) measure as a spatial visual attention estimator with the purpose of significantly accelerating a scene recognition task, preserving its performance. The approach considers both efficient coding, in order to reduce the redundancy of the input data, and the detection of important attributes of the image via local edge phase and energy.
- *García-Díaz et al., 2012* [98]: The Adaptive Whitening Saliency (AWS) model provides a measure of saliency by considering the local energy variability in the *Lab* color space. First, a Gabor bank of filters is applied to the luminance channel *L*, extracting several multioriented multiresolution features. Additionally, a multiscale decomposition of the *a* and *b* color components is computed. Finally, a Principal Component Analysis (PCA) is performed over all these low-level representations in order to decorrelate them, obtaining the local measure of variability that underlies the SM.
- *Leborán et al., 2017* [103]: This approach, called Dynamic Adaptive Whitening Saliency (AWS-D), is an extension of the AWS [98] explained above that computes either static or dynamic saliency maps. The hypothesis of the authors was that saliency has a strong relationship with the variability of the local energy measured over a statistically decorrelated and normalized space. Thus, in order to estimate saliency, the

model looks for the space-time points with maximum variability in the distribution of the local energy across spatio-temporal scales and orientations. In contrast to other spatio-temporal methods, it does not rely on information derived from an explicit background model or the estimation of the optical flow to compute motion features.

In contrast, **TD** architectures are still scarce and very often tailored to well-defined scenarios. In such cases, the evaluation of the whole scheme is performed regardless of the capability of the guidance tool. What is more, most **TD** approaches guide attention towards specific targets by modulating gains associated with low-level stimuli.

- *Sprague and Ballard, 2003* [60]: They presented a reinforcement learning method that combines action selection and visual perception in a sidewalk navigation task.
- *Bruce and Tsotsos, 2005* [58]: The Attention based on Information Maximization (**AIM**) model proposed computes a **VAM** based on the Shannon's self-information measure. At each image region, saliency is the information that the region conveys with respect to its surroundings. The information of the visual feature is inversely proportional to the likelihood of observing it. Consequently, to calculate this feature, the Probability Density Function (**PDF**) has to be estimated. Moreover, in order to reduce the dimensionality of the problem, an Independent Component Analysis (**ICA**) is performed. Thus, given a local image region, the probability of observing the *RGB* values is estimated via the product of the likelihood of the components associated with that region.
- *Judd et al., 2009* [90]: The model is based on a linear Support Vector Machine (**SVM**), taking some image features and human fixations to define salient locations.
- *Elazary and Itti, 2010* [95]: They proposed a more flexible **TD** model that can concurrently select the best features to guide attention and adjust the width of feature detectors.

Finally, although suggested by the prevalent studies about attention [11, 25], just a few works proposed hybrid models incorporating **BU** and **TD** factors.

- *Navalpakkam and Itti, 2005* [83]: The model optimizes the integration of **BU** cues for target detection by maximizing the signal-to-noise ratio of the target versus background.
- *Peters and Itti, 2007* [86]: This model computes a task-dependent map based on the scene gist and the eye

fixations gathered from a video game scenario. BU and TD integration is simply conducted as a multiplication of both components.

2.4.2 Bayesian models

Since our model proposed in Chapter 3 is based on a Bayesian formulation of visual attention, this section briefly introduces some Bayesian or probabilistic approaches found in the *state-of-the-art*.

On the one hand, probabilistic BU algorithms make use of Bayes' rule to combine the features observed with prior constraints:

- *Torralba, 2003* [82]: The author presented a Bayesian approach for visual search tasks. BU saliency stands on a global feature based on the scene gist, which provides a shortcut to detect the presence or absence of objects in an image before exploring it.
- *Itti and Baldi, 2005* [59]: They proposed a probabilistic framework for modeling saliency as "surprise" by computing the KL divergence between the posterior and prior beliefs about image features, either in space or time domains.

On the other hand, Bayesian TD models are characterized by their capacity to learn from data, taking advantage of data statistics to model the underlying attention process and allowing to obtain interpretable relationships between data and visual fixations:

- *Zhang et al., 2008* [87]: This framework understands saliency as the point-wise mutual information between BU local features and TD search target features. The model, known as Saliency Using Natural Statistics (SUN), tries to reproduce the visual experience acquired by an organism. To achieve this, it defines visual saliency as the probability of a searched target at every point in the visual field given the features observed. Using Bayes' rule, and assuming that feature and location are independent, the self-information is taken as definition of BU saliency: the rarer a feature is, the more it will attract our attention. Given this definition, features are calculated as responses of biologically plausible linear filters, such as Difference of Gaussians (DoG) and Gabor filters, as same as in [81], and also as the responses to filters learned from natural images, using ICA.

After this first approach, the model was extended to temporally dynamic scenes in [91], characterizing the video statistics around each pixel using a bank of spatio-temporal filters with separable space-time components.

- *Li et al., 2011* [97]: They provided a multi-task learning approximation for visual attention in video, where different ranking functions for fusing BU and TD maps are learned depending on the scene content.

2.4.3 Deep Neural Networks

Convolutional Neural Networks (CNNs) [64, 112], the current dominant paradigm for many supervised tasks in computer vision, have also been applied to model visual attention achieving promising results, especially in the still image domain. Furthermore, either for SMs refinement based on attention modules or visual attention estimation in video, researchers have recently drawn on recurrent Long Short-Term Memory (LSTM)-based networks [18].

DNNs involve training end-to-end models according to a loss function, using a database of images or videos annotated with GT fixations. They unify feature extraction, fusion and saliency prediction into a single structure. This usually improves the system performance at the expense of making the analysis of these stages more challenging, mainly due to the abstract nature of representations at the deepest layers of these strategies.

Among the first attempts to rely on deep learning for static saliency estimation, the following ones deserve our analysis:

- *Vig et al., 2014* [99]: It constituted one of the first approaches that makes use of CONV layers for saliency prediction, training a softmax classifier on top of them, so that SMs are formulated as generalized Bernoulli distributions.
- *Kümmerer et al., 2014* [100], *2016* [102]: First, in 2014, the authors presented Deep Gaze I, which builds SMs by using the object recognition model of Krizhevsky et al. [3] and a prior distribution to model the central fixation bias [100]. Further on, in 2016, Deep Gaze II [102] applies transfer learning to saliency prediction by fine-tuning a few layers on top of the features from a VGG network [113], also for object recognition.
- *Huang et al., 2015* [101]: The SALICON fine- and coarse-scale model evaluates the use of four commonly-known differentiable saliency metrics as the objective function of a simple CNN architecture, providing image SMs which integrate information at different scales. Furthermore, the authors introduced a large-scale image dataset [114] for training new models, annotated by means of a mouse-tracking procedure, which seemed to correlate well with human fixations.
- *Kruthiventi et al., 2017* [104]: The authors presented DeepFix, a fully CNN built on top of a VGG network [113] for hierarchically

modeling the BU mechanism of visual attention. The network captures semantics at multiple scales and information derived from the global context, and also models center-bias effect in human attention.

- Cornia *et al.*, 2018 [108]: SAM model is able to predict accurate SMs by incorporating a neural attentive mechanism based on convolutional LSTMs [115]. Given a SM, the method refines it by iteratively focusing on the most prominent regions, and also considering the center bias existing in human fixations by learning a set of prior Gaussian maps.

Despite the outstanding performance achieved by these approaches, they still miss some key elements [116], mostly related to mis-detections of people, actions and text, and the relative importance assigned to them when they take place simultaneously.

It should also be pointed out that only a few works have drawn on deep learning to tackle the estimation of visual attention in videos:

- Jiang *et al.*, 2017 [105]: Together with a large-scale eye-tracking database of generic videos, the authors proposed a CNN to learn spatio-temporal features based on object motion, and also a two-layer convolutional LSTM network to smooth the transition between SMs of consecutive frames.
- Bak *et al.*, 2018 [106]: The authors studied the use of dynamic models for saliency prediction in videos, providing several single and two-stream CNNs and evaluating different fusion mechanisms to combine spatial and temporal information. They demonstrated the importance of considering inherent motion information, by training models on estimated optical flow.
- Wang *et al.*, 2018 [107]: Similarly to [105], the authors shared a new large-scale database of videos with fixations, organized by their categories, and then presented a CNN-LSTM architecture with an attention mechanism. The CNN encodes the static saliency information, which allows the LSTM to learn temporal saliency representations for successive video frames.

Finally, CNNs have also been applied to a particular type of saliency, closely related to the object detection task, both in images or videos. These models provide a map of objectness, measuring the probability that each image location belongs to an object. Nevertheless, these networks are often trained on databases with object segmentation masks as GT instead of fixations.

- Li *et al.*, 2016 [2]: The authors presented an end-to-end deep contrast network for salient object detection. The network

consists of a pixel-level FC stream, which generates a SM with pixel-level accuracy, and a segment-wise spatial pooling stream, which improves the modeling of saliency discontinuities along object boundaries. Moreover, on top of these two streams, a CRF model can be applied to improve the spatial coherence and contour localization.

- Wang *et al.*, 2018 [109]: They provided an efficient framework for object detection in videos that captures spatial and temporal saliency information via a short-term analysis consisting of learning from adjacent frame pairs, without the need to compute optical flow.

The reader is also referred to a recently released survey by Ali Borji about deep learning-based models for saliency prediction in images and videos [117], where the author also discusses emerging applications of these architectures and which aspects should be considered in order to improve them.

2.4.4 Applications

Nowadays, computer vision techniques have to deal with millions and millions of available data, just as the HVS does. This is probably the prime reason why the effort in developing computational systems to accomplish the task carried out by visual attention has increased during the last few years.

We can highlight two main purposes for visual attention modeling [117]:

- The first, related to our contributions in Chapters 3 and 4, is to understand *how* visual attention works in humans, trying to describe the behavioral and neural processes involved.
- The second is to predict *where* people look, in order to address traditional image and video applications, such as object [118, 119] and action [120, 121] recognition, video summarization [122], patient diagnosis [123] or image quality assessment [124], in broader and more complex scenarios, while providing efficient solutions and better performances. In line with the second objective, our contributions in Chapters 5 and 6 pursue to facilitate the task of a CCTV operator in a video surveillance scenario, by means of a deep architecture that models attention in the temporal domain.

A GENERATIVE PROBABILISTIC MODEL FOR SPATIO-TEMPORAL VISUAL ATTENTION

3.1 INTRODUCTION

Modern computer vision techniques have to deal with vast amounts of visual data, which requires a computational effort that has often to be accomplished in challenging scenarios. The interest in solving these image and video applications efficiently has led researchers to develop visual attention-based methods to expertly drive the corresponding processing to conspicuous regions that either depend on the context or are based on specific requirements of the task.

In this chapter, we propose a generative hierarchical probabilistic framework for spatio-temporal visual attention understanding and prediction in video. Our model is independent of the application scenario, and founded on the most outstanding psychological studies about attention and eye movements, which support that *guidance* is not based directly on the information provided by early visual processes but on a contextual representation arisen from them.

Drawing from the well-known Latent Dirichlet Allocation (LDA) [12] method for the analysis of large corpus of data, and inspired by some of its supervised extensions [13, 14], our approach defines task- or context-driven visual attention as a mixture of latent sub-tasks, which are in turn modeled as a combination of specific distributions associated with low-, mid- and high-level spatio-temporal features. Learning from fixations gathered from human observers, we incorporate an intermediate level between feature extraction and visual attention estimation that enables to obtain guiding representations.

CHAPTER OVERVIEW

The chapter is organized as follows. First, in Section 3.2, we review the most relevant and recent related work in perception and spatio-temporal visual attention, justify our claims, and present our

main contributions. Then, Section 3.3 presents a broad set of example features which may potentially guide the attention of observers, and will be tested as input for our experiments in Chapter 4. Next, probabilistic Latent Topic Models (LTM) are introduced in Section 3.4, putting special emphasis on LDA and its supervised extensions, on the basis of which our approach is developed. Finally, Section 3.5 describes in detail our generative probabilistic framework for spatio-temporal visual attention understanding and prediction.

3.2 RELATED WORK AND MAIN CONTRIBUTIONS

As it was discussed in Section 2.4.1, most of the visual attention models developed thus far are Bottom-Up (BU) approaches [84, 88, 94, 103], whereas Top-Down (TD) architectures are mostly tailored to scenarios where it is not critical to achieve a good estimation of visual attention to solve a particular task [60, 86]. Only a few works tackle the confluence between BU and TD factors, and there is still a lot of research to be done in real scenarios or videos.

Probabilistic models have an undeniable potential: they are able either to estimate attention or to understand its process [82, 87, 97]. However, their expressive ability is often limited to very simple single-layered fusion schemes built on top of features.

To overcome these shortcomings, we propose a general data-driven hierarchical probabilistic architecture to estimate visual attention in videos, which can be applied to different scenarios and tasks by simply learning from human fixations.

Our LTM-based design was first described in [125]. It introduced an intermediate level formed by latent sub-tasks, which bridges the gap between features and visual attention, and enables to obtain more comprehensive interpretations of attention guidance. These representations provide additional information about how features are combined both in attracting and inhibiting spatial locations. Then, TD visual attention is modeled as a linear regression over a set of learned intermediate sub-tasks rather than over the features themselves. Depending on the context, distinct features could draw visual attention. For instance, motion features are useful to follow players and track objects in outdoor scenes, while color, faces or text are more relevant in TV recordings. However, the fundamental basis of the system is, indeed, generic and independent from the application scenario.

In a recent work [126], we updated this design, making two substantial contributions:

1. We generated a categorical binary response for each spatial location to model visual attention, in contrast to the continuous variable used in our previous approach [125]. The new system now allows to automatically align the sub-tasks discovered to a

binary response by means of a logistic regressor, which fully corresponds to the definition of human fixations.

2. We extended the initial set of basic and novelty spatio-temporal low-level features presented in [125], including and modeling some new mid- and high-level features related to camera motion estimation and object detection, and taking advantage of powerful paradigms such as Convolutional Neural Networks (CNNs). To do the latter, we make use of the features derived from a recently released deep contrast network for salient object detection with pixel-level accuracy [2].

For the sake of simplicity, we will only describe the current extended version of the model in this thesis [126], providing the results on visual attention estimation of the first approach as part of the comparison with the state-of-the-art methods carried out in Section 4.4.

3.3 FEATURE ENGINEERING FOR VISUAL ATTENTION GUIDANCE

According to the most leading psychology theories for computational attention systems [10] [11], different simple features are early and pre-attentively processed in parallel to guide visual search in the human brain (see Section 2.3).

Selective visual attention is built on what it is called the *early representation* [73], a set of conspicuity maps related with some *elementary features* such as color, orientation or motion. These topographical maps do not only surround physical attributes, but also may be explained as relational aspects of these physical characteristics. We may even guide our attention by focusing on mid- and high-level features such as symmetry, faces or text.

In the experiments described in Chapter 4, a wide set of features has been considered. For the sake of completeness, we briefly describe the features in the following subsections.

3.3.1 Basic features: color, intensity and orientation

As stated in FIT [10], the majority of computational models of attention consider three early visual features: *intensity* or *luminance contrast* (I), *color* (C) and *orientation* (O). In this section, we briefly explain and compare how they are represented in the famous model of Itti et al. [81] and its update based on graphs [84], which have been introduced in Section 2.4.1. Due to their easy interpretation, effectiveness and prevalence almost up to the advent of CNNs, we decided to make use of the activation maps from [84] in our experiments on visual attention understanding.

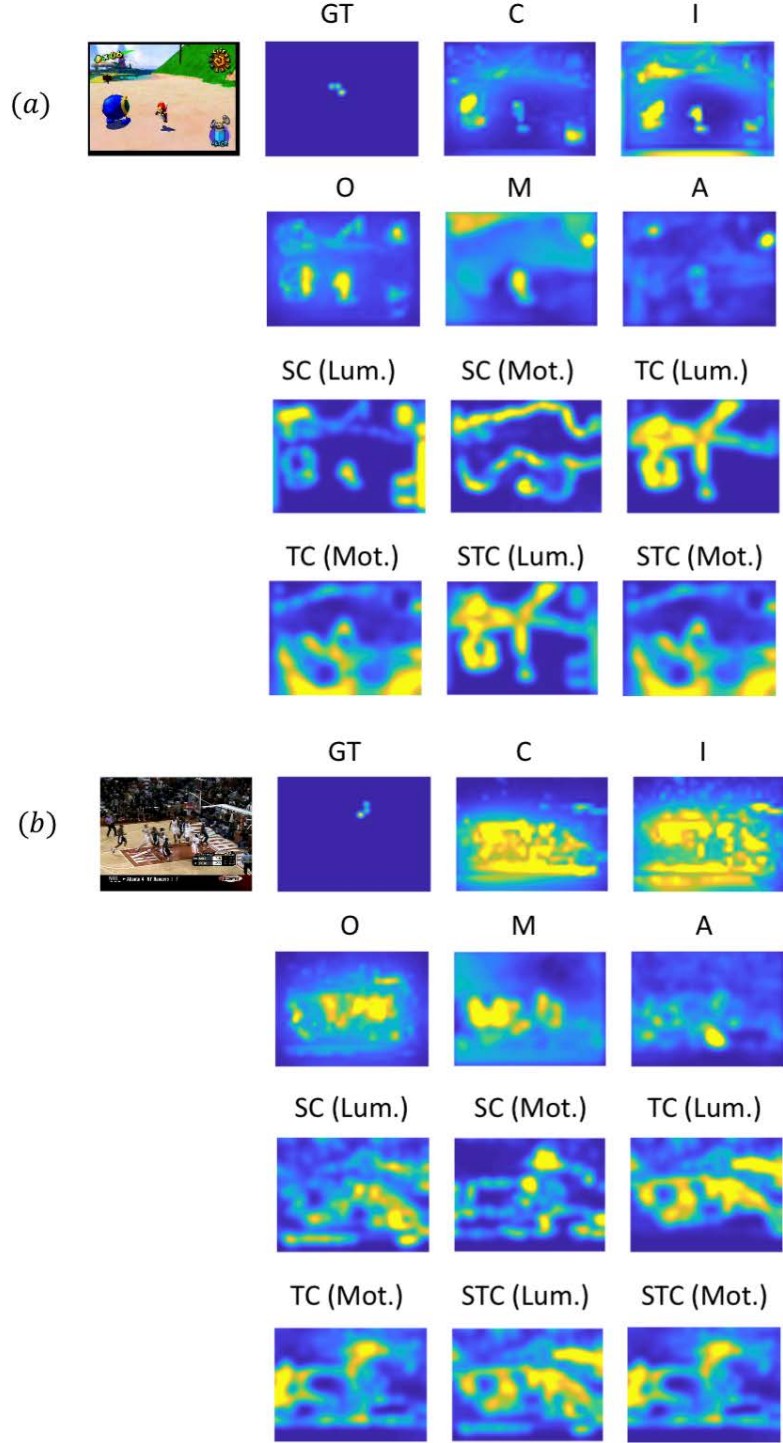


Figure 3.1: Basic, motion-based and novelty feature maps computed for two example frames taken from Videogames (a) and Sports (b) categories from CRCNS-ORIG [15] database. Basic features are color (C), intensity (I) and orientation (O), extracted by using the BU approach for saliency estimation of Harel et al. [84]. Motion-based features are velocity or motion magnitude (M) and acceleration (A). Novelty features are spatial coherence (SC (Lum.), SC (Mot.)), temporal coherence (TC (Lum.), TC (Mot.)) and spatio-temporal coherence (STC (Lum.), STC (Mot.)), computed either over the pixel intensity values \mathcal{I} or the motion phase θ_M .

Feature maps extraction

On the basis of the red (r), green (g) and blue (b) components of a given still image, five channels are obtained, which constitute the basis of the subsequent feature and activation maps. First, an intensity channel I , computed as the linear combination of the three components:

$$I = \frac{r + g + b}{3} \quad (3.1)$$

Then, r , g and b channels are normalized by I to separate hue from intensity, and four additional color channels (red, green, blue and yellow) are calculated:

$$R = r - \frac{g + b}{2} \quad (3.2)$$

$$G = g - \frac{r + b}{2} \quad (3.3)$$

$$B = b - \frac{r + g}{2} \quad (3.4)$$

$$Y = \frac{r + g}{2} - \frac{|r - g|}{2} - b \quad (3.5)$$

A multi-scale process is performed for feature extraction. For that purpose, five Gaussian pyramids ($I(\sigma)$, $R(\sigma)$, $G(\sigma)$, $B(\sigma)$ and $Y(\sigma)$) are generated by consecutively low-pass filtering and sub-sampling each channel. In total, each pyramid is composed of nine spatial scales $\sigma \in [0, 8]$.

It should be noted that Harel et al. [84] method proposes the *DKL* color space [127] as a better alternative to the *RGBY* model. This color space is composed by three axis. The first one represents luminance changes independently from variations in chromaticity, while along the others chromaticity varies without changing the excitation of blue-sensitive or red- and green-sensitive cones, respectively.

Finally, I is convoluted by several oriented Gabor filters $O(\sigma, \theta)$ at different scales σ and with multiple orientations $\theta \in \{0^\circ, 45^\circ, 90^\circ, 135^\circ\}$ in order to extract orientation-based maps.

Once the feature maps have been obtained, an activation map associated with each early visual feature is computed. To this end, the pioneer Itti et al. [81] model proposed a center-surround approach, while Harel et al. [84] presented a graph-based mechanism. Both are explained in the following paragraphs.

Center-surround activation maps formation

Activation maps in [81] arise from the center-surround difference (\ominus) between “center” fine scales c and “surround” coarser scales s . This operation, which tries to simulate the receptive field structure of neurons in the *HVS* [56, 63], allows to detect prominent locations

with respect to their surround. To this effect, it involves an interpolation to the finer scale and a point-by-point subtraction.

First, several intensity contrast maps, which correspond to neurons receptivity to dark spots on bright surrounds and vice versa, are obtained as follows:

$$\mathcal{I}(c, s) = |I(c) \ominus I(s)|, \quad (3.6)$$

where \ominus denotes the across-scale difference between two maps. Secondly, different color maps are concerned with red/green ($\mathcal{RG}(c, s)$) and blue/yellow ($\mathcal{BY}(c, s)$) double opponencies, also perceivable in human visual cortex:

$$\mathcal{RG}(c, s) = |(R(c) - G(c)) \ominus (G(s) - R(s))| \quad (3.7)$$

$$\mathcal{RG}(c, s) = |(B(c) - Y(c)) \ominus (Y(s) - B(s))| \quad (3.8)$$

Local orientation information is also provided by orientation-selective neurons in primary visual cortex. Following the same process, orientation feature maps are computed to encode the local orientation contrast between center and surround scales:

$$\mathcal{O}(c, s, \theta) = |O(c, \theta) \ominus O(s, \theta)| \quad (3.9)$$

In accordance with the hypothesis that similar maps associated with a particular feature compete for saliency, while different features contribute independently to it, three separate activation or conspicuity maps are built for intensity contrast, color and orientation. Feature maps obtained have different dynamic ranges, so a normalization operation $\mathcal{N}(\cdot)$ is applied to them before their combination. This operator not only removes amplitude differences between maps, but also stands out in each map those activation spots whose difference from the average is large. Finally, activation maps are computed by performing an across-scale addition between maps \oplus , reducing each feature map before to the fourth spatial scale considered:

$$\bar{\mathcal{I}} = \bigoplus_c \bigoplus_s \mathcal{N}(\mathcal{I}(c, s)) \quad (3.10)$$

$$\bar{\mathcal{C}} = \bigoplus_c \bigoplus_s [\mathcal{N}(\mathcal{RG}(c, s)) + (\mathcal{BY}(c, s))] \quad (3.11)$$

$$\bar{\mathcal{O}} = \sum_{\theta} \mathcal{N} \left(\bigoplus_c \bigoplus_s \mathcal{N}(\mathcal{O}(c, s, \theta)) \right) \quad (3.12)$$

Graph-based activation maps formation

Alternatively, the Markovian approach presented in [84] tries to imitate the communication between neurons in the visual cortex when processing areas of a scene. Given a feature map M at a particular scale, it establishes a fully-connected directed graph G_A by

connecting every location (i, j) in $M(\sigma)$ with all other locations (p, q) . The dissimilarity $d((i, j)|(p, q))$ of $M(i, j)$ and $M(p, q)$ is defined as:

$$d((i, j)|(p, q)) \triangleq \left| \log \frac{M(i, j)}{M(p, q)} \right| \quad (3.13)$$

Thus, a weight is assigned to the directed edge from location (i, j) to location (p, q) :

$$w_1((i, j)|(p, q)) \triangleq d((i, j)|(p, q)) \cdot F(i - p, j - q) \quad (3.14)$$

where F is related to their spatial distance:

$$F(a, b) \triangleq \exp \left(-\frac{a^2 + b^2}{2\gamma^2} \right), \quad (3.15)$$

being γ a free parameter of the algorithm. A Markov chain is then defined on G_A by normalizing the edge weights of each node to sum to 1, and drawing on the correspondence between locations and states, and edge weights and transition probabilities. The equilibrium distribution of this chain highlights those regions that have high dissimilarity with their surrounding, resulting in the expected conspicuousness map. Finally, the activation map associated with each early visual feature is obtained from the combination of all the normalized conspicuousness maps at different scales, according to a procedure similar to the one applied in [81], which has been explained above. Figure 3.1 includes examples of these feature maps for two frames in CRCNS-ORIG [15] database.

3.3.2 Motion-based features

Motion is undoubtedly another feature that attracts our gaze. It was introduced to model visual attention for the first time in [128]. Given two images, this neurobiological approach considers a motion map in terms of the difference between their corresponding Gabor orientation pyramids, capturing a wide range of object speeds.

The motion-based features used as input for the method presented in this chapter of the thesis draw instead on the optical flow technique proposed in [129] for dense motion estimation. Moreover, a parametric motion model is obtained to estimate camera motion, which also serves to detect moving objects. Both are described below, together with the feature maps used in our attention model: *velocity*, *acceleration* and *camera motion*.

Optical flow estimation

Optical flow [130] computes an independent estimate of motion \mathbf{v}_n at each spatial location $\mathbf{x}_n = (x_n, y_n)$, which can be tackled by minimizing the Sum of Squared Differences (SSD) between the

intensity or brightness of corresponding pixels in two consecutive frames I_{t-1} and I_t in a video:

$$E_{SSD}(\{\mathbf{v}_n\}) = \sum_n [I_t(\mathbf{x}_n + \mathbf{v}_n) - I_{t-1}(\mathbf{x}_n)]^2, \quad (3.16)$$

where $\{\mathbf{v}_n\}$ denotes the whole vector field. In order to efficiently optimize this cost function, an image pyramid is usually built and motion is estimated hierarchically from coarse to fine scales, as first suggested by Lucas and Kanade in [131]. The solution to this function is underconstrained, since we have two variables $\mathbf{v}_n = (u_n, v_n)$ to determine and just one equation per pixel. For each pair of consecutive frames I_{t-1} and I_t , the patch-based typical approach to this problem involves a local summation over overlapping regions \mathbf{x}_n and $(\mathbf{x}_n + \mathbf{v}_n + \Delta\mathbf{v}_n)$, as well as performing gradient descent on Eq. 3.16 using a Taylor series expansion of the displaced image function:

$$\begin{aligned} E_{SSD}(\{\mathbf{v}_n + \Delta\mathbf{v}_n\}) &= \sum_n [I_t(\mathbf{x}_n + \mathbf{v}_n + \Delta\mathbf{v}_n) - I_{t-1}(\mathbf{x}_n)]^2 \\ &\approx \sum_n [I_t(\mathbf{x}_n + \mathbf{v}_n) + \mathbf{J}_t(\mathbf{x}_n + \mathbf{v}_n) \Delta\mathbf{v}_n - I_{t-1}(\mathbf{x}_n)]^2 \\ &= \sum_n [\mathbf{J}_t(\mathbf{x}_n + \mathbf{v}_n) \Delta\mathbf{v}_n + e_n]^2, \end{aligned} \quad (3.17)$$

where

$$\mathbf{J}_t(\mathbf{x}_n + \mathbf{v}_n) = \nabla I_t(\mathbf{x}_n + \mathbf{v}_n) = \left(\frac{dI_t}{dx}, \frac{dI_t}{dy} \right) (\mathbf{x}_n + \mathbf{v}_n) \quad (3.18)$$

is the image gradient or Jacobian at $(\mathbf{x}_n + \mathbf{v}_n)$ and $e_n = I_t(\mathbf{x}_n + \mathbf{v}_n) - I_{t-1}(\mathbf{x}_n)$ is the temporal derivative or brightness change between images.

Horn and Schunck [132] later proposed a regularization-based variational framework to minimize Eq. 3.16 simultaneously over all flow vectors \mathbf{v}_n , which is known as the linearized optical flow constraint:

$$E_{SSD} = \iint [(I_x u + I_y v + I_t)^2 + \alpha (|\nabla u|^2 + |\nabla v|^2)] dx dy, \quad (3.19)$$

denoting $(I_x, I_y) = \mathbf{J}_t(\mathbf{x}_n + \mathbf{v}_n)$ and $I_t = e_n$ spatial and temporal derivatives, respectively. In addition, α is a regularization constant to be determined.

Using as baseline the algorithms in [133, 134], and also including symmetric flow computation, Liu et al. [129] presented a layer-wise optical flow estimation method. Layered motion framework arises from the observation that motion in a scene is often associated with few objects at different depths, so that pixels motion can be estimated more accurately if they are grouped into suitable objects

or layers. According to this approach, the optical flow constraint has three terms. First, a data term, which matches the two consecutive frames I_{t-1} and I_t :

$$E_{data}^{(t)} = \int g * M_t(x, y) |I_t(x + u_t, y + v_t) - I_{t-1}(x, y)| dx dy, \quad (3.20)$$

being g a Gaussian filter, M_t the visible layer mask that indicates which layers' pixels are not occluded, and $\mathbf{v}_t = (u_t, v_t)$ the flow field from I_t to I_{t-1} . Second, a smoothness term, defined as:

$$E_{smooth}^{(t)} = \int (|\nabla u_t|^2 + |\nabla v_t|^2)^\eta dx dy, \quad (3.21)$$

where η constant varies between 0.5 and 1. Finally, symmetric matching is achieved by means of the following term:

$$E_{sym}^{(t)} = \int |u_t(x, y) + u_{t-1}(x + u_t, y + v_t)| + |v_t(x, y) + v_{t-1}(x + u_t, y + v_t)| dx dy. \quad (3.22)$$

Note that $E^{(t-1)}$ terms are defined in a similar way. Thus, the optimization function consists of the sum of these terms:

$$E(\mathbf{v}_t, \mathbf{v}_{t-1}) = \sum_{i=t-1}^t E_{data}^{(i)} + \alpha E_{smooth}^{(i)} + \beta E_{sym}^{(i)}, \quad (3.23)$$

being α and β two additional regularization constants. Flow computation is performed from coarse to fine image pyramid levels, updating the visible layer mask M_t after each scale. Given the flow estimated at the current scale, if $M_{t-1}(\mathbf{x} + \mathbf{v}_t) = 0$ or the symmetry term in Eq. 3.22 at this location is beyond a threshold, $M_t(\mathbf{x}) = 0$ for the next finer scale.

Motion parameterization

From the vector field \mathbf{v}_n computed, the correspondence between the spatial locations \mathbf{x}_n and \mathbf{x}'_n in two consecutive frames can be defined as:

$$\mathbf{x}'_n = \mathbf{x}_n + \mathbf{v}_n \quad (3.24)$$

If N_t is the total number of pixels of each frame, and we express each spatial location $n \in N_t$ in homogeneous coordinates, so that $\bar{\mathbf{x}}_n = (x_n, y_n, 1)^T$, we are able to represent both frames as two $N \times 3$ matrices \mathbf{X} and \mathbf{X}' . Then, a parametric model for camera motion can be obtained as:

$$\mathbf{X}' = \mathbf{X}\mathbf{H} \Rightarrow \mathbf{v} = \mathbf{X}(\mathbf{H} - \mathbf{I}) \Rightarrow \mathbf{v} = \mathbf{X}\mathbf{P} \quad (3.25)$$

being \mathbf{H} and \mathbf{P} the 3×3 matrices that define the geometric transformation between frames and the parametric camera motion

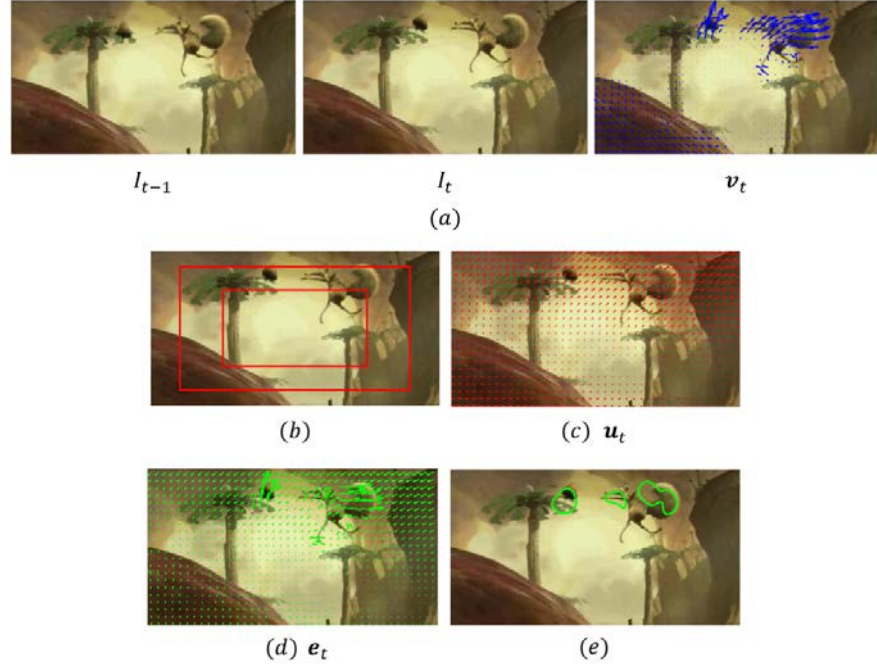


Figure 3.2: Motion parameterization in two example frames taken from a Commercials video in the DIEM [16] database. (a) Frames I_{t-1} and I_t , together with the computed optical flow vector field \mathbf{v}_t . (b) Camera motion modeling. Assuming that moving objects are centered on the image and since optical flow estimation close to the edges is less accurate, we only consider spatial locations within the red inner ring to obtain the parametric camera motion model. (c) Camera motion modeled as a translation vector field \mathbf{u}_t . (d) Residual motion vector field \mathbf{e}_t . (e) Moving objects correspond to the green regions indicated in the image, with residual energy \mathbf{e}_t^2 higher than an empirically determined threshold ζ .

model, respectively. According to an affine motion model [135], \mathbf{P} is defined as:

$$\mathbf{P} = \begin{bmatrix} \mathbf{A} & \mathbf{u} \\ \mathbf{0}^T & 1 \end{bmatrix} \quad (3.26)$$

where \mathbf{A} is a 2×2 non-singular matrix, $\mathbf{u} = (u_x, u_y)$ a translation 2D-vector and $\mathbf{0}$ a null 2D-vector. The transformation has six degrees of freedom, which correspond to the six elements in \mathbf{A} and \mathbf{u} , and is estimated as:

$$\mathbf{P} = \mathbf{X}^+ \mathbf{u} \quad (3.27)$$

being \mathbf{X}^+ the pseudoinverse of the matrix \mathbf{X} of homogeneous coordinates. Assuming that moving objects tend to be centered on the image and since optical flow estimation close to the edges is less accurate, we only consider spatial locations within an inner ring for camera motion modeling, as shown in Figure 3.2(b).

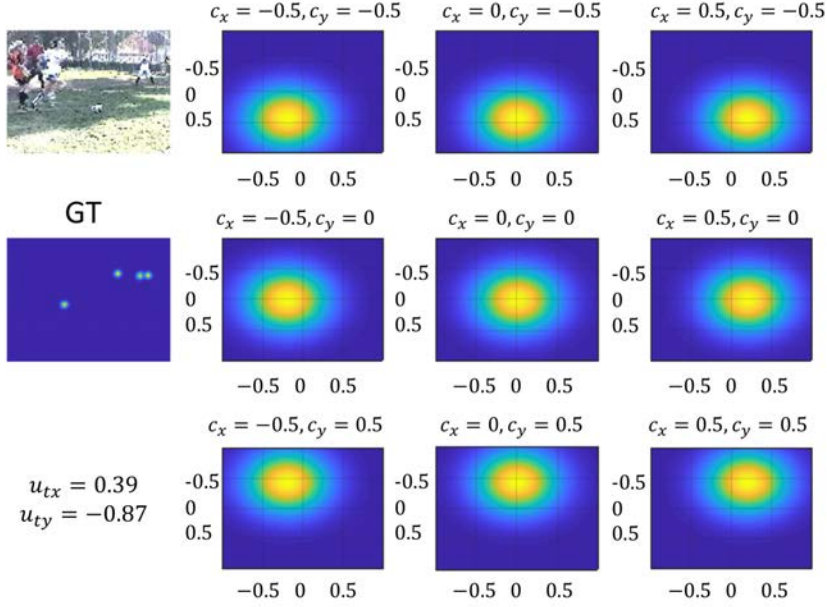


Figure 3.3: Camera motion modeling in an Outdoor frame taken from the CRCNS-ORIG [15] database. Visual attention based on camera motion is modeled by means of a 2D Gaussian distribution $N(\mathbf{c}_{z_n} \odot \mathbf{u}, \Sigma_{z_n})$ over the spatial coordinates, where \mathbf{u} is the translation vector modeling the camera motion. Figure illustrates, given a sample value of \mathbf{u} , the distribution learned for different $\mathbf{c} = (c_x, c_y)$ values, being \mathbf{c} the vector of parameters that establishes a relation between the camera motion and the predicted position of the attention. Variance is set to $\Sigma^2 = \text{diag}(0.25)$ in order to cover a sufficiently wide area in the scene.

Finally, if we want to detect moving objects on the scene, we use \mathbf{P} to compute the energy e_n^2 of the residual motion for each pixel \mathbf{x}_n :

$$e_n^2 = \|\mathbf{v}_n - \mathbf{P}\mathbf{x}_n\|^2 \geq \zeta \quad (3.28)$$

Those pixels with residual energy higher than an empirically determined threshold ζ will correspond to moving objects, as can be seen in the example in Figure 3.2(e).

Feature extraction

Once we have computed optical flow and subtracted camera motion from motion vectors to detect moving objects, we compute two maps based on them. First, velocity or motion magnitude (M), which is calculated using the Euclidean or L2 norm as follows:

$$\mathcal{M}_t = \|\mathbf{v}_t\| \quad (3.29)$$

Then, acceleration (A), which is its absolute derivative:

$$\mathcal{A}_t = \|\mathbf{v}_t - \mathbf{v}_{t-1}\| \quad (3.30)$$

Examples of these two maps are shown in Figure 3.1.

Finally, camera motion may also influence viewers regarding a video. Indeed, as seen in previous studies [136], observers tend to follow the camera motion direction to draw their attention to the new information and objects that emerge in the camera view.

As was defined above, $\mathbf{x}_n = (x_n, y_n)$ is a 2D vector with spatial coordinates x and y associated with the spatial location n . Hence, the visual attention based on camera motion is modeled by means of a 2D Normal distribution over the spatial coordinates:

$$\mathcal{CM}_t \sim N(\mathbf{c} \odot \mathbf{u}_t, \Sigma), \quad (3.31)$$

where $\mathbf{u} = (u_x, u_y)$ is the vector modeling the camera motion as a simple translation whose values are computed from a parametric affine motion model, as described above; \odot stands for the Hadamard product between vectors, and $\mathbf{c} = (c_x, c_y)$ is the vector of parameters that establishes a relation between the camera motion and the predicted position of the attention, and is learned during the training process. Figure 3.3 illustrates, given a sample value of \mathbf{u} , examples of the distribution for several values of the vector \mathbf{c} . Intuitively, higher absolute values of \mathbf{c} point to camera motion as an important feature for visual attention. If \mathbf{c} and \mathbf{u} have the same positive or negative sign, camera motion constitutes an attracting feature, and the 2D Gaussian distribution is shifted in its direction. In contrast, if they have the opposite sign, camera motion inhibits attention, and the distribution is shifted in the opposite direction. The second parameter Σ , which controls the spatial extent of the Gaussian distribution, has been empirically set to $\Sigma = \text{diag}(0.25)$ in order to cover a sufficiently wide area in the scene.

3.3.3 Novelty features

Those regions of the scene that continually change may also attract the attention of observers. In order to highlight them, novelty can be modeled by using the so-called *coherence-based features*, which analyze the distribution of pixel values along space and time in order to detect areas where dispersion is large. To do this, we rely on the work done in [137], extracting *spatial*, *temporal* and *spatio-temporal* coherence maps. In the following definitions, let us consider f_n as the value of a given feature map at the location n with spatial coordinates $\mathbf{x}_n = (x_n, y_n)$, over which a coherence-based feature value is computed.

- Spatial Coherency (SC) identifies regions that belong to a well defined object, or where motion is stable, highlighting most changing regions, which can be more surprising and salient. For each pixel or spatial location n , it is calculated as the

variance with respect to the mean μ_n of its neighbor values in a window of size $N \times N$, with $N = 5$:

$$\mathcal{SC}_n = \frac{1}{N^2} \sum_{m \in R_n} (f_m - \mu_n)^2 \quad (3.32)$$

where R_n stands for the $N \times N$ neighborhood centered in the spatial location n .

- Temporal Coherency (TC) is useful to distinguish between regions with small random motion (e.g. leaves falling from trees) and those with regular motion. Given a pixel n in a current frame t , it is the variance with respect to the mean μ_{tn} of its value across the $T = 7$ preceding frames, as the effect of motion in one frame on the scan path of the eye usually lasts for no more than 5 – 7 frames [138]:

$$\mathcal{TC}_{tn} = \frac{1}{T} \sum_{i \in [t-T+1, t]} (f_{in} - \mu_{tn})^2 \quad (3.33)$$

- Spatio-Temporal Coherency (STC) combines both previous measures and it is calculated for each pixel n as the variance with respect to the mean μ_{tn} of the set of pixel values within a $N \times N$ neighborhood, with $N = 5$, across the $T = 7$ preceding frames:

$$\mathcal{STC}_{tn} = \frac{1}{TN^2} \sum_{i \in [t-T+1, t]} \sum_{m \in R_n} (f_{im} - \mu_{tn})^2 \quad (3.34)$$

In total, 6 maps are computed: three over the pixel intensity values I and three over the motion phase $\theta_{\mathcal{M}} = \arctan \frac{v}{u}$. Examples of these maps are gathered in Figure 3.1.

3.3.4 Object-based features

Despite the questionable conclusions of some psychologists [29, 71] about the inclusion and modeling of faces and other semantic categories as attributes that guide attention, they have been considered in some computational approaches [139, 140]. Moreover, a recent analysis of Bylinskii et al. [116] gathers a list of under-predicted regions when estimating saliency in images, which mainly consists of parts of objects or subjects.

Hence, we have considered detectors for some general-purpose objects that tend to attract visual attention, in order to appraise their utility in some contexts. In particular, cascade object detectors based on the Viola-Jones algorithm [141] are used to detect people's *frontal* (F) and *profile faces* (PF), *upper bodies* (B) and *pedestrians* (P) and a

detector working on the Harris corner response [142] is used to detect *text* (T). Both methods are briefly explained here below. Many detectors for other visual concepts may also be included in our model without effort.

Cascade object detectors

Cascade classifiers [141] are trained to detect objects with invariant aspect ratio at different scales. Therefore, in order to use them to locate an object whose appearance changes significantly when varying its orientation, such as in a face, it would be necessary to train a single detector for each of its views, as it is carried out in our approach by considering either frontal or profile faces. These kind of detectors are outstanding for being extremely fast, at the same time they can achieve high detection rates.

A cascade detector involves several stages. Each of them constitutes an ensemble of weak learners, which are simple classifiers trained using AdaBoost [143]. Given a candidate window x in an image, a weak classifier $h_t(x)$ is a threshold function that depends on a feature value $f_t(x)$. It can be formulated as follows:

$$h_t(x) = \begin{cases} -s_t & \text{if } f_t(x) < \theta_t \\ s_t & \text{otherwise} \end{cases} \quad (3.35)$$

Simple features f_t used, which are founded on the Haar filters introduced in [144], are calculated as the difference between the sum of the pixels within two sub-regions in x . The threshold θ_t and the polarity $s_t \in \pm 1$ are determined in the learning phase of the detector, on the basis of positive and negative samples of the object class for which it is trained. At each stage of the detector, a final strong decision is made as the weighted linear combination of the T decisions made by all weak learners, being the weights α_t inversely proportional to their corresponding training errors:

$$h(x) = \text{sgn} \left(\sum_{t=1}^T \alpha_t h_t(x) \right) \quad (3.36)$$

Under the assumption that the majority of regions of an image covered by a sliding window do not contain an object of interest, early stages of the cascade are designed to rapidly select the most promising sub-windows with a low false negative rate. Indeed, a cascade can be understood as a degenerate decision tree where, if a sub-window is rejected at any stage, no further processing is performed, dramatically decreasing the number of sub-windows to be evaluated. The complexity of the strong classifiers gradually increases until reaching the end of the cascade, with the purpose of achieving a final high detection rate.

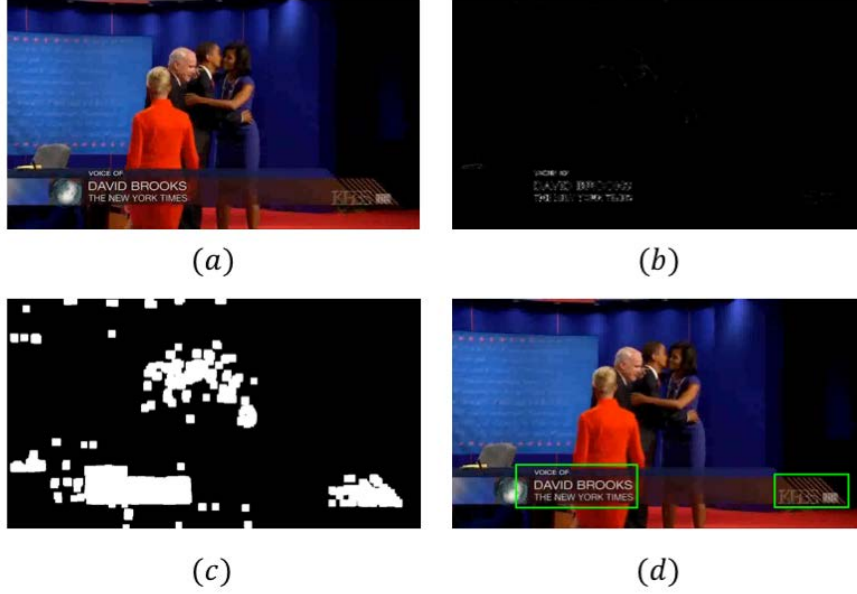


Figure 3.4: Harris response for text detection in an example frame taken from a TV News video in the DIEM [16] database. (a) Example frame. (b) Absolute value of the Harris response computed in order to detect corners, which often correspond to text areas. (c) Binary mask obtained after applying to the response a non-maximum suppression operation consisting of a dilation. (d) Regions in the binary mask are filtered, selecting those which are horizontally or vertically aligned and occupy less than one third of the area of the whole image, which correspond to text bounding boxes.

Harris response for text detection

A simple text detector is proposed based on the commonly used Harris detector [142], which locates corners in an image. Interest points in whose local neighborhood two dominant and different edge directions arise constitute corners, which often correspond to text areas. Corners are invariant to translation, rotation and illumination.

In order to detect corners in a grayscale image I , the following second-moment matrix M is computed, which is derived from its horizontal I_x and vertical I_y directional gradients:

$$M = \begin{bmatrix} I_x I_x & I_x I_y \\ I_x I_y & I_y I_y \end{bmatrix} \quad (3.37)$$

Then, the Harris response HR is calculated as follows:

$$HR = \det(M) - k \operatorname{tr}(M)^2, \quad (3.38)$$

where k is an empirically determined constant. A non-maximum suppression operation consisting of a dilation is performed to find

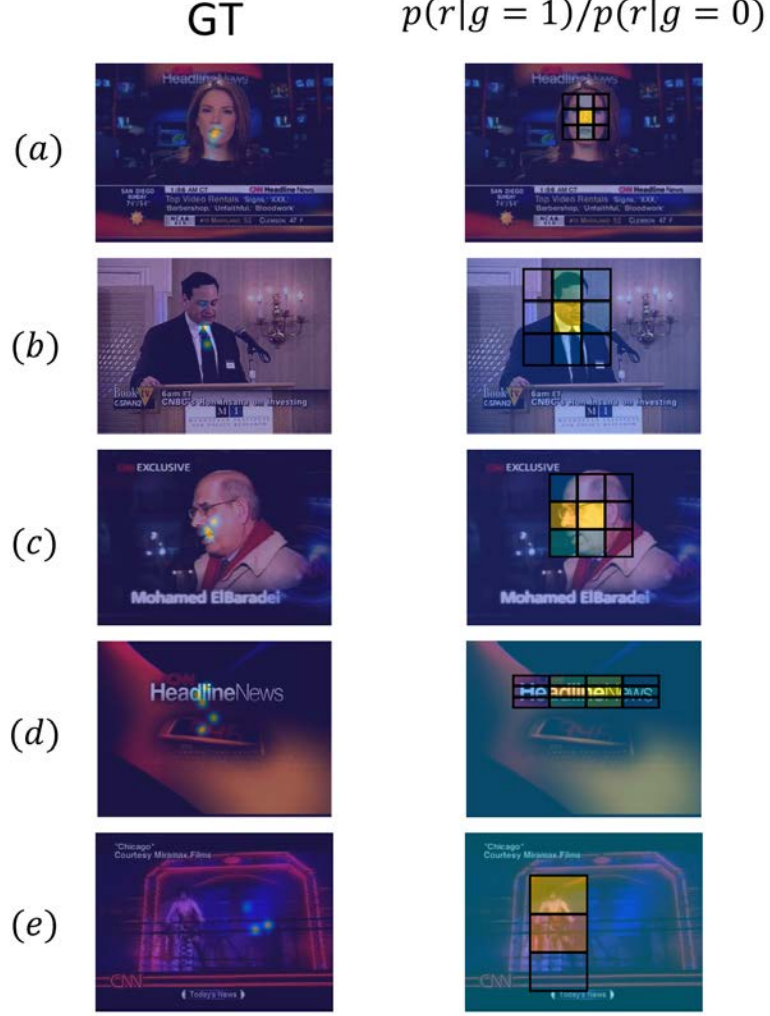


Figure 3.5: Object-based feature maps computed for example frames taken from TVNews (a, b, c, d) and TalkShows (e) categories from CRCNS-ORIG [15] database. (Left) Human fixations do not cover the whole object, but concentrate on particular areas/parts of the objects. (Right) Consequently, and based on the detected bounding box, we have divided the image into a set of subregions $r = 0, 1, \dots, R$. Some of them ($r > 0$) divide the object into several cells (9 for frontal (F) and profile faces (PF), and upper bodies (B); 3 for pedestrians (P) and 12 for text (T)). Moreover, an additional subregion $r = 0$ is considered for the background, covering the rest of the image. Overlay heat maps highlight subregions where probabilities of each object for fixated points ($p(r|g=1)$, being $g \in \{0, 1\}$ the ground truth variable indicating if the spatial location attracts or inhibits the attention) are substantially higher than those for non-fixated points ($p(r|g=0)$). Although the prior probability of objects is fairly lower than the probability of background in the database, it can be seen that objects are quite attractive for observers, due to the significant probability of internal cells given fixated locations.

the local maxima in HR , which results in candidate text regions. Finally, we measure the area, orientation and eccentricity of the candidate regions, in order to obtain one or several bounding boxes. We assume that texts are usually horizontally or vertically aligned, and occupy less than one third of the area of the whole image. The complete process is illustrated in Figure 3.4.

Feature modeling

The output bounding boxes from the detectors described above are used to generate high level spatial feature maps. Visual attention usually points to particular locations within the objects, so this fact has to be considered when modeling these features. Since the size of the detected bounding boxes is often large, if we use a 2D Gaussian centered in the bounding box that contains a particular object, for instance, we are notably emphasizing the center of the object with respect to its surroundings. However, attention may be generally fixed at some elements of the object and not only at its center, such as in the case of faces or pedestrians, where subjects often look at the eyes or upper body part, respectively. Rather than directly considering the detected boxes as the feature maps, we have developed spatially-aware discrete distributions.

As shown in Figure 3.5, given a detected bounding box, we consider a non-uniform grid with $R+1$ cells: R cells subdivide the detected box into r small subregions ($r > 0$), and an additional subregion is considered for the background ($r = 0$). Hence, for a given object l being detected (we keep l as the index of the features, in this case object detections), we model a discrete distribution over the $R+1$ defined cells as $p(r|\beta_l)$, where r is a cell in the grid (which is object dependent), and β_l are the parameters of the discrete distribution for the object category l . The distributions are then factorized for every object category and instance (in case that more than one object of a given category are detected on the same frame). By means of discrete spatial distributions that divide objects into several sub-regions, we are able to learn which parts of the object are more attractive, taking advantage of this knowledge to provide more accurate estimations of visual attention.

3.4 LATENT TOPIC MODELS

As introduced in Section 1.2.2, generative models not only make predictions on unseen data, but also offer an interesting interpretation about how this information was generated. Generative probabilistic LTM_s such as Probabilistic Latent Semantic Analysis (PLSA) [145] or Latent Dirichlet Allocation (LDA) [12] have, besides, an additional advantage: they can be used both in unsupervised and supervised contexts.

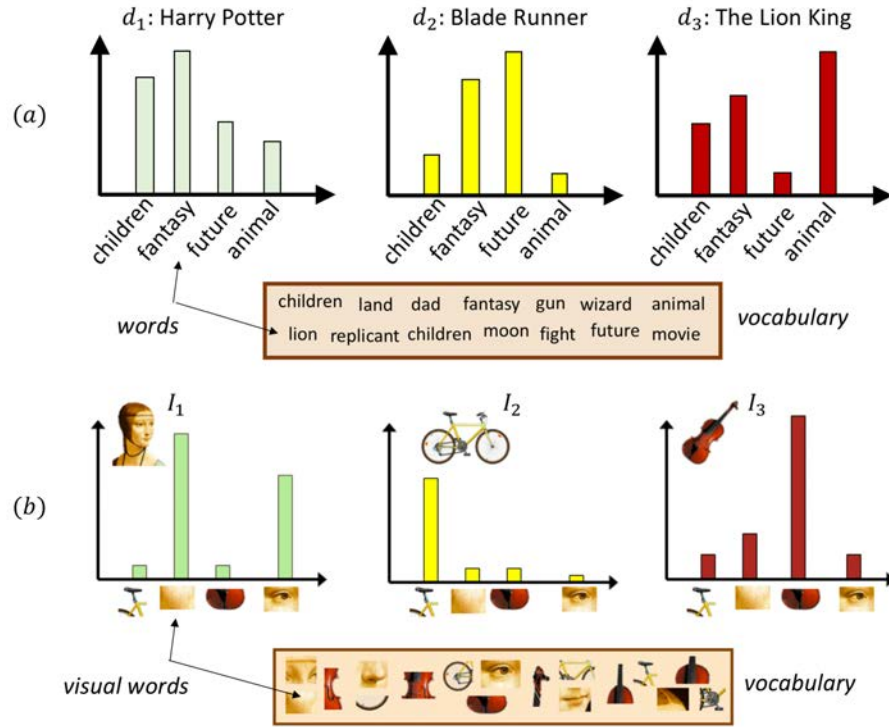


Figure 3.6: Bag-of-Words (BoW) model applied to: (a) A corpus of texts from movie reviews, where each review $\{d_1, d_2, d_3\}$ is represented by means of a histogram of word occurrences; (b) A corpus of images, where each image $\{I_1, I_2, I_3\}$ corresponds to a histogram of visual word occurrences. In both cases, words are taken from a finite vocabulary. Figure adapted from [146].

This section primarily describes the LDA graphical approach by David M. Blei et al. [12], which is the most frequently LTM used, and two of its supervised extensions [13, 14], which are the basis of the contributions to visual attention understanding and estimation presented in this chapter. LTMs are thought to model large collections of discrete data, such as corpus of texts, images or audio tracks. For this reason, we will begin by defining the Bag-of-Words (BoW) notation, employed to organize hierarchically these entities.

3.4.1 Bag-of-Words model

The Bag-of-Words (BoW) model [12, 147] is a hierarchical representation method classically used in Natural Language Processing (NLP) and text categorization, which was later extended to image recognition and retrieval in the computer vision field. Given a large collection of texts, the BoW model defines the following terms, which help to understand the intuition behind the LDA approach:

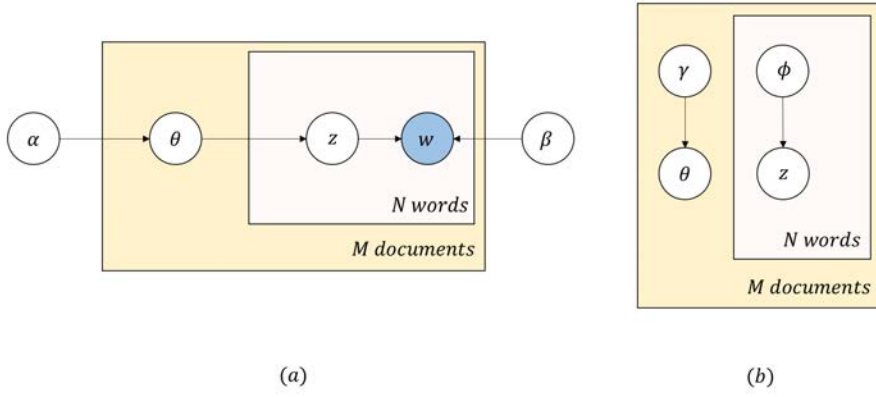


Figure 3.7: (a) Graphical representation of LDA [12]. (b) Graphical representation of the variational distribution used to approximate the posterior in LDA. Shaded nodes represent observed variables, while white nodes denote latent variables to be inferred. Boxes mean independent repetitions, and edges show the dependencies among variables.

- A *word* \mathbf{w} is an item from a finite vocabulary, and constitutes the basic unit of discrete data.
- A *vocabulary* $\{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_V\}$ is a finite collection of V words. Each word \mathbf{w}_v in the vocabulary is represented by a V -vector with a 1 at the position w^v of the word in the vocabulary ($w^v = 1$) and 0 everywhere else ($w^u = 0$ for all $u \neq v$).
- A *document* $\mathbf{d} = (\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_{N_d})$ consists of a sequence of N_d words.
- A *corpus* $\mathcal{D} = \{\mathbf{d}_1, \mathbf{d}_2, \dots, \mathbf{d}_D\}$ is a collection of D documents.

Figure 3.6 shows two examples of application of the BoW model in text and image corpora. Similarly to a document, an *image* \mathbf{I} is decomposed into a set of *keypoints* represented by means of N_I *visual descriptors* $(\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_{N_I})$, which are associated with a finite *vocabulary* of V *visual words*.

3.4.2 Latent Dirichlet Allocation

Latent Dirichlet Allocation (LDA) [12] is a hierarchical Bayesian method initially conceived to model large collections of discrete data, such as text documents. Given a *corpus* of D documents, the model provides an explicit representation of each *document* $\mathbf{d} \in \mathcal{D}$ as a finite mixture over an underlying set of latent topics which, in turn, are modeled by a distribution over words. In fact, a document is defined as a sequence of N_d words denoted by $\mathbf{d} = (\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_{N_d})$, being a *word* an item from a finite vocabulary $\{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_V\}$.

As can be appreciated in the graphical representation of the model shown in Figure 3.7(a), LDA establishes a three-level representation hierarchy. The model first assumes a known and fixed number of topics K in the corpus $\{\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_K\}$. At corpus-level, the K -dimensional Dirichlet variable α sets the global distribution of the topics or topic proportions $\{\alpha_1, \alpha_2, \dots, \alpha_K\}$, being $\alpha_k > 0$, in the whole corpus. In addition, β includes a collection of K V -dimensional discrete or multinoulli variables with the probabilities of each word in the vocabulary. Then, at document-level, the variable θ_d represents the particular topic proportions in each document \mathbf{d} . Finally, at word-level, the variable \mathbf{z}_{dn} stands for the topic associated with each word \mathbf{w}_{dn} in each document \mathbf{d} . \mathbf{z}_{dn} is defined as a K -vector with a 1 at the position of the topic assigned and 0 everywhere else.

Hence, for each document \mathbf{d} in a corpus \mathcal{D} , LDA involves the following generative process. For the sake of simplicity, let us note that we have omitted the document subindex d in those document-dependent variables:

1. Draw the document particular proportions θ of K topics using a corpus-level Dirichlet distribution of parameter α : $\theta|\alpha \sim \text{Dir}(\alpha)$.
2. For each word $\mathbf{w}_n \in N_d$ in the document \mathbf{d} :
 - a) Draw topic assignment $p(\mathbf{z}_n|\theta)$ using a multinomial distribution over the topic proportions θ : $\mathbf{z}_n|\theta \sim \text{Mult}(\theta)$.
 - b) Draw a word \mathbf{w}_n using $p(\mathbf{w}_n|\mathbf{z}_n, \beta)$, which is a multinomial probability conditioned on the topic \mathbf{z}_n .

Given a document \mathbf{d} in the corpus and the corpus-level parameters α and β , the joint distribution of a topic mixture θ , a set of K topics \mathbf{z} and a set of N_d words \mathbf{w} is expressed as follows:

$$p(\theta, \mathbf{z}, \mathbf{w}|\alpha, \beta) = p(\theta|\alpha) \prod_{n=1}^{N_d} p(\mathbf{z}_n|\theta) p(\mathbf{w}_n|\mathbf{z}_n, \beta). \quad (3.39)$$

In order to apply the LDA method, we have to compute the posterior distribution of the latent variables θ, \mathbf{z} given a document \mathbf{d} :

$$p(\theta, \mathbf{z}|\mathbf{w}, \alpha, \beta) = \frac{p(\theta, \mathbf{z}, \mathbf{w}|\alpha, \beta)}{p(\mathbf{w}|\alpha, \beta)}. \quad (3.40)$$

This distribution, nevertheless, is intractable for exact inference due to the coupling between θ and β , so it is necessary to consider an approximate algorithm such as the convexity-based variational inference proposed in [148]. The idea is to make use of Jensen's inequality to arise an adjustable lower bound on the log likelihood, drawing on some variational parameters. These parameters are estimated via an optimization process that tries to find the tightest possible lower bound.

As shown in Figure 3.7(b), by dropping the edges between θ , \mathbf{z} and \mathbf{w} , and also the \mathbf{w} nodes, we achieve the following variational distribution:

$$q(\theta, \mathbf{z} | \gamma, \phi) = q(\theta | \gamma) \prod_{n=1}^{N_d} q(\mathbf{z}_n | \phi_n), \quad (3.41)$$

being the Dirichlet parameter γ and the multinomial parameters (ϕ_1, \dots, ϕ_N) the new free variational parameters.

The optimization problem to find the tightest lower bound to the posterior consists in minimizing the Kullback-Leibler divergence (KL) between the variational distribution and the true posterior $p(\theta, \mathbf{z} | \mathbf{w}, \alpha, \beta)$, and is defined as:

$$(\gamma^*, \phi^*) = \arg \min_{(\gamma, \phi)} D(q(\theta, \mathbf{z} | \gamma, \phi) || p(\theta, \mathbf{z} | \mathbf{w}, \alpha, \beta)) \quad (3.42)$$

The problem can be solved by means of a variational Expectation-Maximization (EM) algorithm, which involves the following iterative procedure:

1. *E-step*: For each document \mathbf{d} , find the optimum values of the variational parameters $\{\gamma_{\mathbf{d}}^*, \phi_{\mathbf{d}}^*\}$.
2. *M-step*: Maximize the obtained lower bound on the log-likelihood with respect to the parameters of the model α and β .

3.4.3 Supervised topic models

So far, we have discussed about the advantages of generative models and the ability of LTM_s to represent texts as a mixture over topics, which can be inferred from a large collection of documents. However, the topics discovered by algorithms like LDA are implicit, so that human expertise is required to arise a comprehensible interpretation of their semantics (e.g. to relate a topic with high probabilities of terms “match”, “players”, “games”, “ball”, with the semantic concept of “sport”).

The objective of supervised extensions of LDA described in the following paragraphs is thus to infer latent topics that not only explain the distribution of words in documents, but also serve to automatically predict a response variable. Both approaches consider response variables y at document-level, as shown in the graphical representation of Figure 3.8.

Supervised Latent Dirichlet Allocation

Supervised Latent Dirichlet Allocation (sLDA) [13] incorporates to LDA a continuous response variable y associated with each

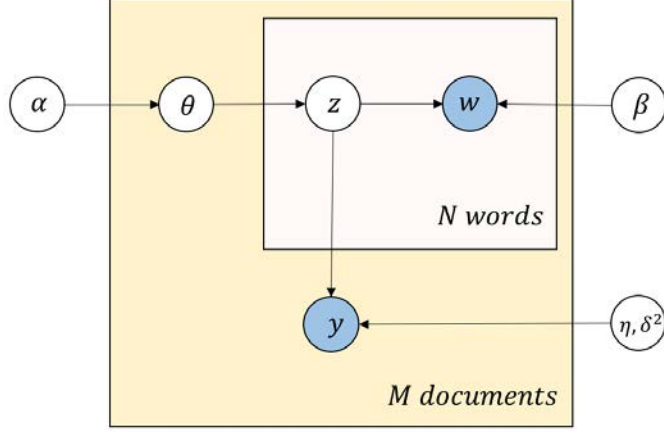


Figure 3.8: Graphical representation of [sLDA](#) [13], which incorporates to [LDA](#) [12] a response variable y , modeled using a linear regression model with parameters $\{\eta, \delta^2\}$. In the case of the [DBA](#) [14] extension of [LDA](#), the representation is similar except for the parameter of the logistic regression model, which is simply η . Shaded nodes represent observed variables, while white nodes denote latent variables to be inferred. Boxes mean independent repetitions, and edges show the dependencies among variables.

document. The documents and the responses are jointly modeled in order to find latent topics that predict the response variable of new unlabeled documents. Examples of applications of [sLDA](#) based on [NLP](#) include the prediction of the numerical rating of a movie review or the number of visits of a website depending on their content.

Figure 3.8(a) shows the graphical representation of [sLDA](#). For each document \mathbf{d} in a corpus \mathcal{D} , the model involves the following generative process. For the sake of simplicity, we have omitted again the document subindex d in document-dependent variables:

1. Draw the document particular proportions θ of K topics using a corpus-level Dirichlet distribution of parameter α : $\theta|\alpha \sim \text{Dir}(\alpha)$.
2. For each word $\mathbf{w}_n \in N_d$ in the document \mathbf{d} :
 - a) Draw topic assignment $p(\mathbf{z}_n|\theta)$ using a multinomial distribution over the topic proportions θ : $\mathbf{z}_n|\theta \sim \text{Mult}(\theta)$.
 - b) Draw a word \mathbf{w}_n using $p(\mathbf{w}_n|\mathbf{z}_n, \beta)$, which is a multinomial probability conditioned on the topic \mathbf{z}_n .
3. Draw a Gaussian response variable $y|\mathbf{z}_{1:N}, \eta, \delta^2 \sim N(\eta^T \bar{\mathbf{z}}, \delta^2)$, based on a linear regression model.

Indeed, in addition to the original [LDA](#), a third step is included at document-level corresponding to the introduced response variable y , which is modeled using a normal linear regression model $N(\eta^T \bar{\mathbf{z}}, \delta^2)$, where $\bar{\mathbf{z}} := (1/N) \sum_{n=1}^N \mathbf{z}_n$ represents the empirical frequencies of

the topics in the document, η is a L -length vector containing the regression coefficients of the response variable, and δ^2 is a dispersion parameter, which provides certain flexibility when modeling the variance of y .

The posterior distribution to solve **sLDA** is $p(\theta, \mathbf{z} | \mathbf{w}, y, \alpha, \beta, \eta, \delta^2)$, and can be again approximated by means of the **KL** optimization problem defined in Eq. 3.42. Now, at the M-step of the algorithm, we also maximize the lower bound on the log-likelihood with respect to the supervision parameters η and δ^2 .

Dirichlet-Bernoulli Alignment

Dirichlet-Bernoulli Alignment (**DBA**) [14] presents an alternative supervised extension to **LDA**, with the purpose of considering multi-class, multi-label and multi-instance classification tasks, where each document from a corpus consists of multiple instances and is related to multiple classes. For instance, a **NLP** application of **DBA** could be, given a social network profile (*document*), the classification of its publications into a set of categories (*topics*).

Hence, each document is modeled as a mixture over a set of predefined classes. Then, each word is generated independently conditioned on the sampled class, and the label of the document is generated conditioned on all the sampled labels used for generating its words.

The graphical representation of **DBA** is similar to the one presented in Figure 3.8 for **sLDA**. However, in contrast to **sLDA**, **DBA** response variable is not continuous but categorical; therefore, a multinomial logistic regression model given by a Bernoulli or softmax distribution $y | \mathbf{z}_{1:N}, \eta \sim \text{Be} \left(\frac{\exp(\eta^T \mathbf{z})}{1 + \exp(\eta^T \mathbf{z})} \right)$ automatically aligns the topics discovered from the data to the predefined classes.

The posterior distribution to solve **DBA** is $p(\theta, \mathbf{z} | \mathbf{w}, y, \alpha, \beta, \eta)$, and can be again approximated by means of the **KL** optimization problem defined in Eq. 3.42. Now, at the M-step of the algorithm, we also maximize the lower bound on the log-likelihood with respect to the supervision parameter η .

3.4.4 Applications to Computer Vision

The description of the **LTM** methods discussed in the previous sections has been focused on the analysis of corpus of texts. However, **LTM** models have been also applied to other types of data [149], such as audio and music [150] or population genetics [151], among others.

What is more, they have been widely-used in computer vision for image retrieval [152], segmentation [39] and captioning [153]. In video processing, they have been applied for action recognition [154]. In most of these applications, image features are quantized into discrete

values so that they take the form of words in documents. However, this is not a straightforward step and requires computing dictionaries of visual words [155, 156].

Notwithstanding this, to our knowledge, the LTM presented in the next section is the first approach proposed for spatio-temporal visual attention on the basis of this type of graphical models.

3.5 VISUAL ATTENTION TOPIC MODEL

In this section, we describe in detail the system proposed for spatio-temporal visual attention understanding and prediction, which we have called visual Attention TOpic Model (ATOM).

3.5.1 Model overview

ATOM generative model is supported by the following assumption [125]:

Task- or context-driven visual attention in video can be modeled as a mixture of several sub-tasks which, in turn, can be represented as combinations of low-, mid- and high-level spatio-temporal features obtained from video frames.

The generative model thus receives as input a set of visual features, which are used to learn several related sub-tasks. These sub-tasks automatically lead the attention of the system to the most appealing areas of a scene. Depending on the scenario, visual attention may be attracted by different events. Our goal is not to detect these events of interest for a particular application, but to efficiently guide the later processing to areas of special importance in the video.

Figure 3.9 illustrates our hypothesis for three different scenarios in CRCNS-ORIG [15] database. First, looking at the contexts given, visual attention may be attracted by different events or elements in the scene: people *running* and *walking* in the case of *Outdoor*; *game character* and *goals* or *items* in *Videogames*; and *players* and *scoreboards* in *Sports*. Note that some contexts may share similar attractions, like *ball*, which is present both on *Outdoor* and *Sports* videos. Our goal is to automatically discover sub-tasks that guide later processing to the areas where those occur. In turn, these sub-tasks can be modeled as combinations of spatio-temporal features. For instance, the use of a motion feature combined with a detected face or pedestrian could be useful to represent the sub-task “Player”. In contrast, the sub-task “Scoreboard” is well-defined by some intensity or color features, together with a detected text.

Probabilistic Latent Topic Models (LTMs), which have been commonly used to extract hidden semantic structures (*latent topics*) from a text corpus, can be helpful to unsupervisedly understand large

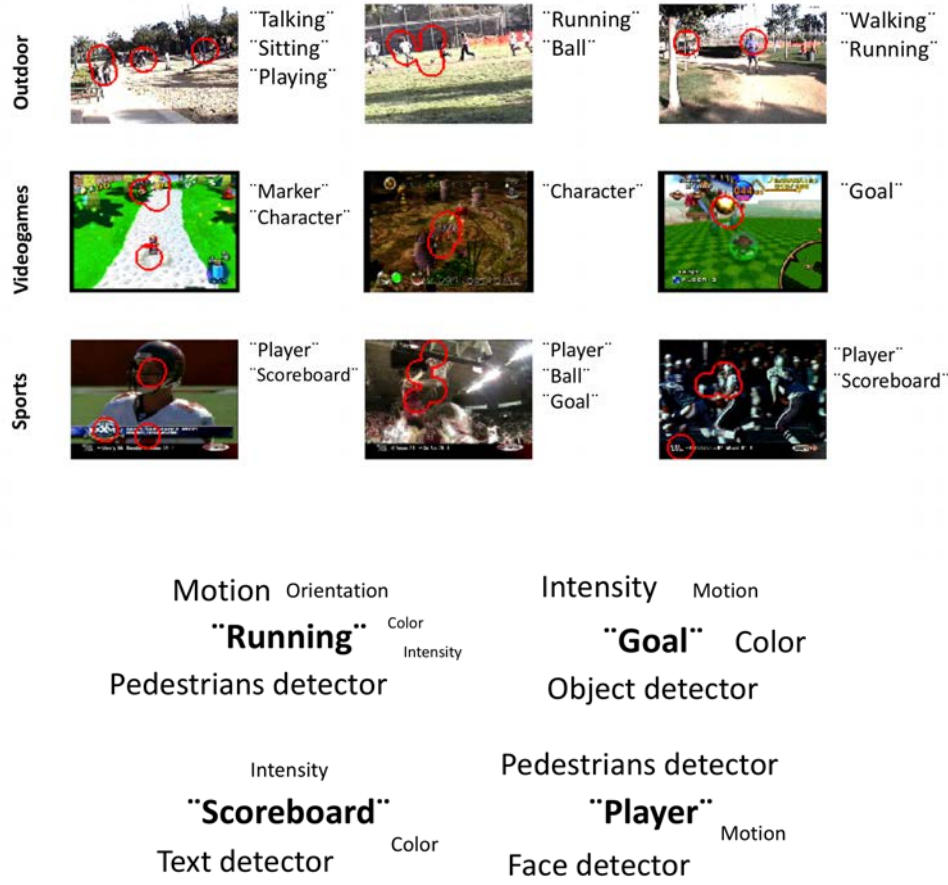


Figure 3.9: Visual attention modeled in three different scenarios taken from CRCNS-ORIG [15] database (*Outdoor*, *Videogames* and *Sports*) as a mixture of several relevant sub-tasks (e.g. "Running", "Goal", "Player", etc.), associated with particular areas of special importance for observers, which are highlighted in the example frames on the left side. Some of them may appear similarly in different contexts, such as "Ball" or "Goal". On the bottom of the figure, word clouds show how some sub-tasks (bold central words) are represented as a combination of features (surrounding words: intensity, motion, detectors, etc.). Feature importance, represented by the font size of each word showing a feature, varies from one sub-task to another. For example, motion information and pedestrian detections are more relevant for "Running"; in contrast, an object detector, along with intensity and color features are more advantageous to represent a "Goal" in a videogame.

amounts of information, such as the human perception features that are quickly and parallelly processed by the brain. Our approaches

involve thus a LTM which relies on the well-known LDA algorithm [12] and its supervised extension DBA [14].

First, by understanding frames as a mixtures over topics, LDA allows to interpret them using unsupervised statistical distributions, which associate each frame to multiple topics with different proportions. In our particular scenario, task-driven visual attention is modeled as a finite mixture over a set of K topics, which represent the sub-tasks contributing to model visual attention, either by attracting or by inhibiting it. Note that both terms, topics and sub-tasks, are used interchangeably along the thesis. In parallel, for a given video frame I_t , a set of L visual descriptors $\mathbf{f} = \{f_1, f_2, \dots, f_L\}$ is computed at each spatial location n , so that the latent topics are in turn modeled as combinations of these features.

The original LDA is completely unsupervised, so that the topics are learned to maximize the likelihood of a corpus, and requires of human knowledge to align topics and semantic concepts. In our case, in contrast, we aim to learn how humans guide their attention to visual stimuli, so that the Ground-Truth (GT) fixations provided by different subjects will drive our training step. Visual attention is thus estimated by means of a logistic regression model over the topic assignments. This logistic regression is in charge of aligning the topics discovered from frames to the information gathered in GT binary fixation maps. Hence, our final model draws on the DBA introduced in [14]. Let us note that the latent nature of the topics remains unchanged in our supervised models, as the human fixations used in the training phase are not supervising the topics but, instead, the binary response variable learned by the logistic regression.

The graphical representation of the model is shown in Figure 3.10. In the same way as LDA, ATOM establishes a three-level representation hierarchy. The model first assumes a known and fixed number of latent topics K in the video corpus $\{\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_K\}$, which represent the sub-tasks that contribute to model visual attention. Let us note that some of these sub-tasks may attract human attention whereas others may inhibit it. At video corpus-level, the K -dimensional Dirichlet variable α sets the global distribution of the sub-tasks or sub-task proportions $\{\alpha_1, \alpha_2, \dots, \alpha_K\}$, being $\alpha_k > 0$, in the whole video corpus. High values of all components α_k of the variable α result in mixtures where all sub-tasks are considered to estimate visual attention in every video frame. In contrast, low values of only some α_k provide more particular mixtures of sub-tasks for each frame, being the attention determined by only few prevailing sub-tasks. Moreover, Γ includes a collection of K L -dimensional variables which define the distribution of each feature l given the topic k . Depending on the nature of features, it is possible to model them using the most suitable

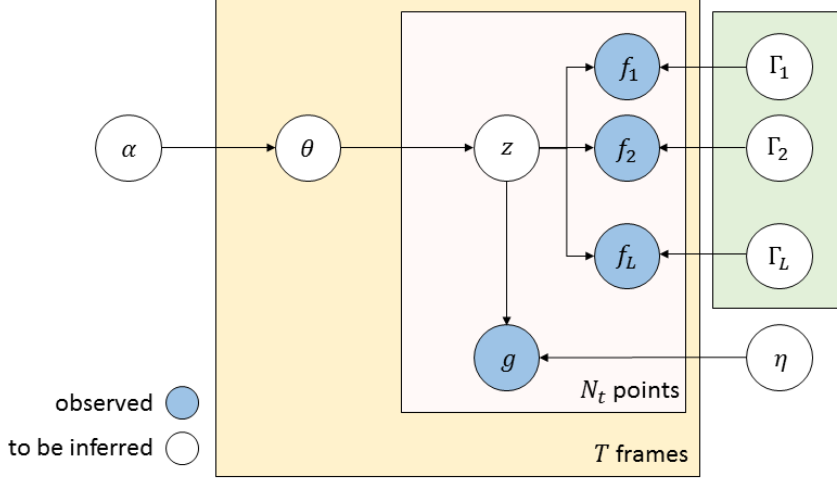


Figure 3.10: Graphical representation of the proposed visual Attention Topic Model (ATOM) generative model. Shaded nodes represent observations from frames, white nodes indicate hidden variables to be inferred, and boxes mean independent repetitions. Edges show the dependencies among variables.

distribution: e.g. normal, exponential, discrete, etc. Then, at frame-level, the variable θ represents the particular sub-task proportions in each frame I_t . Finally, at spatial location-level, the variable \mathbf{z}_n stands for the sub-task associated with each spatial location n in each frame I_t . \mathbf{z}_n is an indexing K -dimensional vector with all zeros in except of a 1 in the position of the selected topic.

The proposed ATOM thus involves the following generative process for each frame I_t in a video corpus $\mathcal{I} = \{I_1, I_2, \dots, I_T\}$. Let us note that, for simplicity, we have removed the sub-index t of the frame in the notation:

1. Draw the frame particular proportions θ of K topics using a corpus-level Dirichlet distribution of parameter α : $\theta | \alpha \sim \text{Dir}(\alpha)$.
2. For each spatial location $n \in N$ in the frame I_t :
 - a) Draw topic assignment using a multinomial distribution over the topic proportions θ : $\mathbf{z}_n | \theta \sim \text{Mult}(\theta)$.
 - b) Represent the local appearance of the spatial location n by drawing L independent visual features f_{ln} using the topic particular distributions $p(f_{ln} | \mathbf{z}_n, \Gamma)$, where Γ includes the parameters of the distributions of the L features, given the selected topic \mathbf{z}_n .
 - c) Draw the binary response variable g_n modeling the visual attention using a logistic regression model given by the following Bernoulli distribution:

$g_n | \mathbf{z}_n, \eta \sim \text{Be} \left(\frac{\exp(g_n \eta^T \mathbf{z}_n)}{1 + \exp(\eta^T \mathbf{z}_n)} \right)$, where η is the parameter vector that models attention based on the selected topic \mathbf{z}_n .

Hence, for each frame I_t , we first generate a particular mixture of these topics θ based on the distribution with the global topic proportions α . Once θ is known, we analyze the different spatial locations of the frame such that, for each n , we first select a sub-task by using the index-variable \mathbf{z}_n . Based on \mathbf{z}_n , we draw the local appearance of the spatial location using the particular feature-topic distribution $f_{nl} | \mathbf{z}_n, \Gamma$, where Γ stands for the parameters of the distributions of the L features considered. Sub-task is thus chosen so that its corresponding distribution parameters are the ones that maximize the likelihood of the visual features observed at this location.

For the sake of simplicity, we assume that $p(\mathbf{z}_n | \theta)$ is independent for all locations n , which makes the solution tractable, both simplifying the definition of the algorithm and, at the same time, improving the system efficiency. In contrast, other approaches such as MRFs [157], applied to image segmentation, are able to capture such spatial constraints. Nonetheless, it should be noted that some of the visual features that we extract for each sampled location (e.g. color, intensity, orientation, CNNs-based) consider beforehand this spatial dependency. Moreover, we assume conditional independence among the L features, so that the joint distribution of features for a particular topic can be factorized into the individual probability distributions $p(f_l | \mathbf{z}, \Gamma)$. Finally, we also generate the attention response g_n by computing the logistic regression model over the selected topics.

3.5.2 Guiding features extraction

Motivated by the general conclusions of psychological theories about attention [10, 11], the general hierarchical probabilistic framework presented may operate over a great number of diverse features. Depending on their nature, they may be modeled using various probability distributions: e.g. *normal*, *exponential*, *discrete*, etc. It should be remarked that our model is not feature-dependent, so that any kind of feature can be incorporated by selecting the appropriate distribution. Furthermore, for each application scenario and based on human fixations, our model will automatically discover which particular features are more and less discriminant to model attention and correspondingly assign appropriate parameters to their distributions. Hence, one could include a broad general set of features as the model will automatically reduce or neglect the

influence of those that do not guide the attention in a particular context.

Hereunder is a list of the features extracted for our experiments in Chapter 4, which correspond to 24 feature maps, including the section where they were explained. Some of the feature maps are handcrafted and allow us to perform a meaningful interpretation of the estimated visual attention. They carry continuous values, and are modeled using a Gaussian probability density function:

Basic features (Section 3.3.1)

1. *Color* (C)
2. *Intensity* (I)
3. *Orientation* (O)

Motion-based features (Section 3.3.2)

4. *Velocity or motion magnitude* (M)
5. *Acceleration* (A)

Novelty features (Section 3.3.3)

6. *Spatial Coherency (Luminance)* (SC (Lum.))
7. *Spatial Coherency (Motion)* (SC (Mot.))
8. *Temporal Coherency (Luminance)* (TC (Lum.))
9. *Temporal Coherency (Motion)* (TC (Mot.))
10. *Spatio-Temporal Coherency (Luminance)* (STC (Lum.))
11. *Spatio-Temporal Coherency (Motion)* (STC (Mot.))

Then, camera motion is modeled as a multivariate Gaussian $\mathcal{CM} \sim N(\mathbf{c} \odot \mathbf{u}, \Sigma)$, as described in Section 3.3.2. Due to the diagonal nature of the covariance matrix Σ , we can decompose it into two independent univariate Gaussians (feature maps 12 and 13).

Next, we have used some object detectors in order to compute high-level spatial feature maps, which are modeled by means of discrete spatial distributions, as explained in Section 3.3.4:

14. *Frontal faces detector* (F)
15. *Upper bodies detector* (B)

16. *Profile faces detector* (PF)
17. *Pedestrians detector* (P)
18. *Text detector* (T)

Finally, we decide to consider 6 feature maps (19-24) derived a CNN, which are modeled using Gaussian distributions. As in other computer vision applications, there is no doubt about the success of CNNs for visual attention modeling. Nevertheless, despite their capability of discovering discriminant high-level visual features, it is still necessary to clarify the relationship between the feature maps derived from CNNs and the psychophysical stimuli that guide attention. This implies the development of complementary modules able to provide this mapping, such as our hierarchical method, which facilitates the integration with such NN schemes. Indeed, our intermediate sub-task level can be placed straightforwardly over the top layers of a deep network.

CNN-based features have been drawn from the Deep Contrast Network for salient object detection recently introduced by Li et al. [2]. For the sake of completeness, and due to its use in our system for modeling visual attention in the temporal domain in Chapter 5, we also take a brief look to this architecture in Section 5.4.3. The reason is twofold: first, they allow modeling more general objects than those identified by previously mentioned detectors; and second, they demonstrate the ability of our model to find efficient and diverse combinations of features that help to understand how visual attention works in a given scenario. We employ the models trained by the authors on a different image dataset, and use the feature maps of the penultimate layer to obtain features modeling general objectness.

3.5.3 Inference process

This section explains the inference process of our probabilistic model. As in the original LDA [12] and its extension [14], exact inference is not possible due to the coupling between the variables θ and \mathbf{z} , which prevents from inferring the posterior distribution of the parameters given the data. Therefore, we propose to use a simplified variational distribution q (that is tractable) and mean-field variational inference, so that the KL between the variational distribution q and the posterior distribution is computed. The proposed variational distribution is as follows:

$$q(\theta, \mathbf{z} | \gamma, \phi_{1:N}) = q(\theta | \gamma) \prod_{n=1}^N q(\mathbf{z}_n | \phi_n) \quad (3.43)$$

that incorporates two new variational parameters: ϕ , which is the parameter of a multinomial distribution $q(\mathbf{z}_n|\phi_n)$, and γ , the parameter of a Dirichlet distribution $q(\theta|\gamma)$. This optimization is equivalent to maximize the Evidence Lower Bound (ELBO) over the log-likelihood of all the frames in the corpus. In particular, using Jensen's inequality, the ELBO of the log-likelihood of a frame can be expressed as:

$$\begin{aligned} \log p(f_{1:N,1:L}, g_{1:N}|\alpha, \Gamma_{1:K,1:L}, \eta) &\geq E_q[\log p(\theta|\alpha)] \\ &+ \sum_{n=1}^N E_q[\log p(\mathbf{z}_n|\theta)] + \sum_{n=1}^N E_q[\log p(f_{n,1:L}|\mathbf{z}_n, \Gamma_{1:K,1:L})] \\ &+ \sum_{n=1}^N E_q[\log p(g_n|\mathbf{z}_n, \eta)] + H(q) \end{aligned} \quad (3.44)$$

where $E_q[\cdot]$ and $H(\cdot)$ are, respectively, the expectation over the variational distribution q and the entropy of a distribution.

The first two terms of Eq. (3.44) and the entropy of the variational distribution are identical to the corresponding terms in ELBO for unsupervised LDA and are described in [12]. The third term is the expected log probability of the features given the related topic model parameters. As was mentioned in Section 3.5.1, we assume conditional independence among features. In the following paragraphs, we particularize this expression for the considered distributions.

- If the feature map f_{nl} is modeled with a univariate *Gaussian distribution* $\Gamma_{1:K,l} \sim \{\mu_{1:K,l}, \sigma_{1:K,l}^2\}$, such as for basic and novelty spatio-temporal features or CNN-based features, the equation for this term is:

$$\begin{aligned} E_q[\log p(f_{nl}|\mathbf{z}_n, \Gamma_{1:K,l})] &= - \sum_{k=1}^K \phi_{nk} \log(\sigma_{kl} \sqrt{2\pi}) \\ &- \sum_{k=1}^K \phi_{nk} \frac{(f_{nl} - \mu_{kl})^2}{2\sigma_{kl}^2} \end{aligned} \quad (3.45)$$

where ϕ_{nk} is the probability that the location n has been drawn by the topic k .

- In the case of camera motion features, the distribution is a multivariate Gaussian $p(\mathbf{x}_n|\mathbf{z}_n, \mu_k, \Sigma_k)$ with $\mu_k = \mathbf{c}_k \odot \mathbf{u}$, being \mathbf{c}_k a parameter to be estimated and $\mathbf{u} = (u, v)$ the camera motion vector. However, due to the diagonal nature of the covariance matrix Σ_k we can decompose it into two independent univariate Gaussians and apply the previous expression.

- In contrast, if the feature is modeled as a *discrete probability distribution* over cells r in a grid, as happens for objects-based features, the expression is:

$$E_q[\log p(r_n | \mathbf{z}_n, \beta_{l z_n})] = \sum_{k=1}^K \phi_{nk} \log(\beta_{kl r_n}) \quad (3.46)$$

where r_n stands for the region in the non-uniform grid defined for the object l that contains the location n , and $\beta_{kl r_n}$ is the value of the discrete distribution in region r_n that contains the point n for the object l and the topic k .

The fourth term includes the visual attention response variable g_n and is drawn as a logistic regression model over the topic assignment \mathbf{z}_n with parameter η :

$$\begin{aligned} E_q[\log p(g_n | \mathbf{z}_n, \eta)] &= E_q \left[\left(g_n - \frac{1}{2} \right) \eta^T \mathbf{z}_n \right] \\ &- E_q \left[\log \left(\exp \left(\frac{\eta^T \mathbf{z}_n}{2} \right) + \exp \left(\frac{-\eta^T \mathbf{z}_n}{2} \right) \right) \right] \end{aligned} \quad (3.47)$$

By taking second derivatives, it can be noticed that the second term above is a convex function in the variable $\eta^{T^2} \mathbf{z}_n^2 = (\eta^T \odot \eta^T)(\mathbf{z}_n \odot \mathbf{z}_n)$, so we can bound it by using the lower bound for logistic function [158], which is the first order Taylor expansion in the variable $\eta^{T^2} \mathbf{z}_n^2$:

$$\begin{aligned} &\log \left(\exp \left(\frac{\eta^T \mathbf{z}_n}{2} \right) + \exp \left(\frac{-\eta^T \mathbf{z}_n}{2} \right) \right) \\ &\geq -\frac{\xi_n}{2} - \log(1 + \exp(-\xi_n)) \\ &\quad - \frac{1}{4\xi_n} \tanh \left(\frac{\xi_n}{2} \right) E_q \left[\eta^{T^2} \mathbf{z}_n^2 - \xi_n^2 \right] \\ &\approx -\frac{\xi_n}{2} - \log(1 + \exp(-\xi_n)) \\ &\quad - \frac{1}{4\xi_n} \tanh \left(\frac{\xi_n}{2} \right) (\eta^{T^2} \phi_n - \xi_n^2) \end{aligned} \quad (3.48)$$

where ϕ_n is the vector of topic proportions ϕ_{nk} in the location n and ξ_n is an additional variational parameter associated with each point n .

It should be noted that, during variational inference, we work on expected values. This means that the indexing variable \mathbf{z}_n is replaced by the variational ϕ_n , which now contains the expected values of the topic assignments given a location n . Therefore, since ϕ_n is a vector with real values (the topic proportions for that sampled location), in practice each location n is in turn modeled as the mixture of sub-tasks that best explains its visual appearance.

Computing the derivatives of the KL with respect to the parameters and setting them equal to zero allows us to obtain the

update equations for the variational procedure. In particular, in the *variational E-step* we must update the variational parameters:

$$\phi_{nk} \propto \frac{\prod_{l=1}^{L_D} \beta_{klr_n}}{\prod_{l=1}^{L_C} \sigma_{kl}} \exp \left[\Psi(\gamma_k) - \Psi \left(\sum_{j=1}^k \gamma_j \right) + \left(g_n - \frac{1}{2} \right) \eta_k - \frac{1}{4\tilde{\xi}_k} \tanh \left(\frac{\tilde{\xi}_k}{2} \right) \eta_k^2 - \sum_{l=1}^{L_C} \frac{(f_{nl} - \mu_{kl})^2}{2\sigma_{kl}^2} \right] \quad (3.49)$$

$$\gamma_k = \alpha_k + \sum_{n=1}^N \phi_{nk} \quad (3.50)$$

$$\tilde{\xi}_{nk} = \eta_k \phi_{nk} \quad (3.51)$$

being L_C and L_D the number of continuous (Gaussian) and discrete features respectively, and $L = L_C + L_D$ the total number of features. Note that we have used the expression $E_q[\log(p(\theta_k|\gamma))] = \Psi(\gamma_k) - \Psi \left(\sum_{j=1}^k \gamma_j \right)$, where $\Psi(\cdot)$ is the digamma function.

In the M-step, we maximize the corpus-level [ELBO](#) with respect to the model parameters $\Gamma_{1:K,1:L}, \eta$, in order to compute their optimal values.

First, parameters μ_{kl} and σ_{kl}^2 are computed for each Gaussian feature l and topic k .

$$\mu_{kl} = \frac{1}{\Delta_{kl}} \sum_{t=1}^T \sum_{n=1}^{N_t} \phi_{tnk} f_{tnl} \quad (3.52)$$

$$\sigma_{kl}^2 = \frac{1}{\Delta_{kl}} \sum_{t=1}^T \sum_{n=1}^{N_t} \phi_{tnk} (f_{tnl} - \mu_{kl})^2 \quad (3.53)$$

where $\Delta_{kl} = \sum_{t=1}^T \sum_{n=1}^{N_t} \phi_{tnk}$ is the normalization factor.

In the case of camera motion, as mentioned above, the parameter is the vector $\mathbf{c}_k = (c_{kx}, c_{ky})$ that multiplies the camera motion vector $\mathbf{u}_t = (u_t, v_t)$ to determine the mean of the Gaussian distribution:

$$\mathbf{c}_k = \frac{\sum_{t=1}^T \sum_{n=1}^{N_t} \phi_{tnk} \mathbf{u}_t \mathbf{x}_{tn}}{\sum_{t=1}^T \sum_{n=1}^{N_t} \phi_{tnk} \mathbf{u}_t^2} \quad (3.54)$$

where $\mathbf{x}_{tn} = (x_{tn}, y_{tn})$ stands for the spatial coordinates vector of the location n in frame t .

Finally, for the case of object-based discrete features, the probabilities β_{klr} of the regions r defined from the outputs of the the object-detector l , and for every topic k are:

$$\beta_{klr} \propto \sum_{t=1}^T \sum_{n=1}^{N_t} \phi_{tnk} 1[r_{nl} = r] \quad (3.55)$$

where $1[r_{nl} = r]$ means that we have a 1 just in case that the point n belongs to the region r (otherwise we have a zero). It is worth noting

that we have added the subindex t when necessary to indicate the frame number in the corpus.

Furthermore, during the training step, we use the [GT](#) response value g_{tn} of all points in the corpus to learn the parameter of the logistic regression model:

$$\eta_k = \frac{2 \sum_{t=1}^T \sum_{n=1}^{N_t} \phi_{tnk} (g_{tn} - \frac{1}{2})}{\sum_{t=1}^T \sum_{n=1}^{N_t} \frac{\phi_{tnk}}{\xi_{nk}} \tanh(\frac{\xi_{nk}}{2})} \quad (3.56)$$

A more comprehensive development of the [ELBO](#) and the previous formulas for parameters estimation is provided in [Appendix A](#).

3.5.4 Learning sub-tasks for spatio-temporal visual attention estimation

As in other supervised approaches, we can distinguish two main stages in our framework, as shown in [Figure 3.11](#). First, in the learning phase, optimal values for the parameters that maximize the [ELBO](#) of the log-likelihood are learned. As we need to learn from annotated data, we first describe how we sample this data from the annotated video datasets. Since we are on a highly unbalanced scenario, in which the areas that attract visual attention are strongly less prominent than those that inhibit it, we need to prevent the later dominating the learning process, which might lead to a poor performance. For that end, we have used the Non-Uniform Sampling ([NUS](#)) strategy proposed in [\[118\]](#), which allows to generate training datasets that balance the number of attracting and non-attracting points. While the first are selected based on the [GT](#) masks computed from human fixations for a given video frame, non-attracting points are sampled from those spatial locations which have not been fixated by viewers in any frame of the same video. In addition, the sampling process also provides the ground truth binary response g_n for each sampled spatial location ($g_n = 1$ for attracting points, and zero otherwise).

Once models are trained, in the test phase, attention is predicted at uniformly spaced locations n in frames. For that end, we remove all terms relating to the supervision (variable g) and estimate the visual attention maps using the expected value of the logistic regression over the topic or sub-task assignments:

$$E[g_n | f_{n,1:L}, \alpha, \Gamma_{1:K}, \eta] \approx \frac{\exp(\eta^T \phi_n)}{1 + \exp(\eta^T \phi_n)} \quad (3.57)$$

In addition, knowing that given a particular frame visual attention is usually focused on small areas of the size occupied by fixations, a histogram equalization procedure is carried out to highlight the most significant regions detected, which helps to improve the system performance.

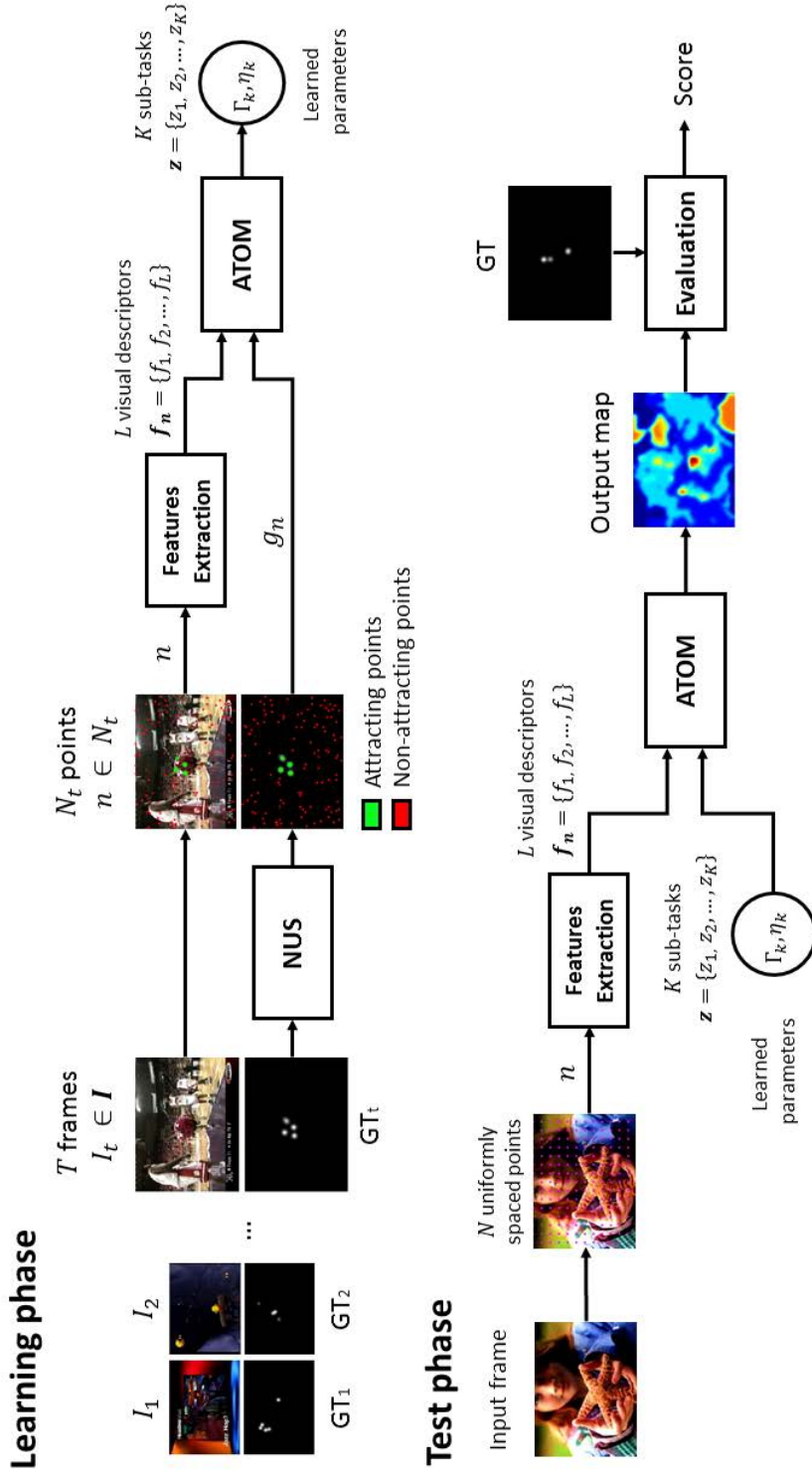


Figure 3.11: Processing pipelines of the approach proposed. First, in the learning phase, we learn the optimal values for the parameters associated to the K sub-tasks that model visual attention. A NUS [18] strategy allows to generate training datasets that balance the number of attracting and non-attracting points. Then, in the test phase, attention is predicted for each frame at N uniformly spaced locations.

EXPERIMENTS ON CONTEXT-DRIVEN VISUAL ATTENTION UNDERSTANDING AND PREDICTION

4.1 INTRODUCTION

In this chapter, we provide an in-depth analysis of our proposal for visual attention modeling described in Chapter 3. We give a meaningful insight about the information reflected in each of the sub-tasks that decompose the visual attention. To this end, we illustrate how our approach successfully learns hierarchical guiding representations adapted to several contexts. Furthermore, we perform a comparison with quite a few methods reported in the literature of visual attention in video.

Experiments show how our proposal successfully learns particularly adapted hierarchical explanations of visual attention in diverse video genres, outperforming several leading models in the literature.

CHAPTER OVERVIEW

First, the experimental design is described in Section 4.2, introducing the databases and the evaluation metrics used to provide the results, and also the initialization of the model. Then, the visual Attention TOpic Model ([ATOM](#)) is used for context-driven visual attention understanding in Section 4.3. Experimental results, together with an analysis of the obtained models and a comparison with *state-of-the-art* methods, are provided in Sections 4.4 and 4.5. Finally, Section 4.6 discusses the model strengths and limitations, and Section 4.7 summarizes our conclusions and motivates and outlines future work.

4.2 EXPERIMENTAL DESIGN

The purpose of our experiments is to demonstrate the ability of the proposed [ATOM](#) to learn meaningful sub-tasks that can be used to understand what guides visual attention in different contexts,

drawing conclusions on whether observers are either driven by similar generic sub-tasks or, in contrast, by certain specific tasks related to each particular scenario. For this reason, we have selected the well-known freely-accessible CRCNS-ORIG [15] and DIEM [16] as benchmark datasets.

4.2.1 Databases

In this section we briefly describe the databases used for the experiments. Further information about the division of the database videos into categories can be found in Appendix B.

CRCNS-ORIG database

CRCNS-ORIG [15] dataset contains eye movement recordings from eight distinct subjects freely watching 50 different video clips (over 46,000 video frames, 25 minutes total, 640×480). Eye traces have been obtained using a 240 Hz ISCAN RK-464 eye-tracker. As set out in Table 4.1(a), clips include complex video stimuli that can be divided into seven categories: *Outdoor*, *Videogames*, *Commercials*, *TV News*, *Sports*, *Talk Shows* and *Others*. Eye fixations of at least 4 subjects are provided for each clip.

The dataset was delivered some years ago with this same intention pursued with our analysis, and has been employed to evaluate a lot of *state-of-the-art* saliency models. However, to our knowledge, none of them had attempted so far to offer a data interpretation such as the one resulted from our approach.

DIEM database

DIEM [16] dataset contains eye movement recordings from over 250 participants freely watching 84 high-definition natural videos (over 240,000 video frames, 134 minutes total, variable dimensions). Eye traces have been obtained using a 1,000 Hz SR Research Eyelink 2000 desktop mounted eye tracker. As is summarized in Table 4.1(b), clips have been classified into seven categories: *TV Shows*, *Documentaries*, *Commercials*, *Talk Shows*, *Sports*, *Cooking* and *TV News*. Eye fixations from approximately 50 subjects are provided for each clip.

In contrast to CRCNS-ORIG [15], DIEM [16] constitutes a greater source of video annotated with *GT* fixations, which serves not only to provide a more truthful result of the method proposed but also to train deeper *CNN*s such as the one presented in the next chapter for motion-based feature maps extraction.

Table 4.1: Categories into which (a) CRCNS-ORIG [15] and (b) DIEM [16] databases are divided.

(a) CRCNS-ORIG [15]			(b) DIEM [16]		
Context	# clips	Frames	Context	# clips	Frames
Outdoor	17	8,357	TV Shows	12	34,271
Videogames	9	15,809	Documentaries	18	56,382
Commercials	4	2,618	Commercials	15	40,558
TV News	7	8,071	Talk Shows	5	8,657
Sports	5	4,851	Sports	20	54,293
Talk Shows	4	4,244	Cooking	7	23,684
Others	4	2,539	TV News	7	22,607
TOTAL	50	46,489	TOTAL	84	240,452

4.2.2 Experimental setup

In order to both assess the performance and gain insight into the latent information provided by the proposed probabilistic method for visual attention estimation, we will compare two different approaches for each database: a) a Context-Generic (C-G) model trained using frames belonging to videos in all the categories; and b) 7 Context-Aware (C-A) models trained on those videos belonging to each category or genre.

The performance over every video in the datasets is evaluated by conducting a 4-fold cross validation procedure, in the case of CRCNS-ORIG [15], and a 5-fold cross validation, in the case of DIEM [16], so that at each iteration some videos are picked for evaluation. For the purpose of avoiding over-fitting, all frames of a video are always grouped together in the same set (train or test).

4.2.3 Evaluation metrics

In parallel to the proposal of computational models for saliency and visual attention, a great effort has been made to evaluate their performance. As summarized in the excellent comprehensive study by Bylinskii et al. [159], this has resulted in a wide variety of metrics based on different assumptions: how the Ground-Truth (GT) fixation map is represented, whether center bias is considered or not or the type of normalization applied to VAM, among others.

In this section, we define those metrics that we have used to assess the performance of the spatio-temporal visual attention methods proposed in this thesis, as well as those taken from the *state-of-the-art*.

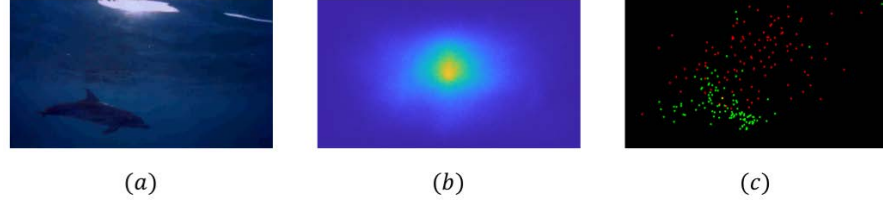


Figure 4.1: **TPs** and **FPs** sampled in an example frame taken from a documentary video in DIEM [16] database. **TPs** correspond to fixations locations, while **FPs** are sampled according to a probabilistic shuffle map with fixations in frames from all other videos in the dataset. (a) Example frame. (b) Shuffle map. (c) **TPs** and **FPs** sampled, indicated in green and red, respectively. Image has been dilated for a better visualization.

The following metrics have been selected according to the suggestions in [159] for models conceived for fixation prediction, which is the aim of the experiments in this chapter, and video surveillance scenarios, as those presented in Chapter 6:

- *Shuffled Area Under ROC Curve (sAUC)*: The Area Under ROC Curve (AUC) [160] is the most used metric in the literature for the evaluation of visual attention models. Given an image or a video frame and its corresponding **GT** fixation map, a **VAM** normalized between 0 and 1 can be understood as the soft output of a binary classifier of fixations; hence, the area under a Receiver Operating Characteristic (ROC) curve, which measures the trade-off between True Positives (**TPs**) and False Positives (**FPs**) at different threshold values, provides a performance score. Depending on how **TPs** and **FPs** are calculated, we can distinguish different AUC implementations [90, 161]. For our experiments, we have chosen a probabilistic Shuffled Area Under ROC Curve (sAUC) metric [162]. This score counteracts the effect of the commonly-observed central fixation bias in scene viewing [163], which is advantageous for those models that consider a center prior. The sAUC chosen to provide our results probabilistically samples **FPs** from fixated locations in other images or videos, instead of uniformly at random.

Despite the effectiveness of the AUCs scores, they are invariant to monotonic transformations of the **VAM**. What is more, attention maps that place different amounts of density at fixated locations receive similar scores as long as they keep fixed the order of the locations. Therefore, it is recommended to supplement them with other metrics.

- *Shuffled Normalized Scanpath Saliency (sNSS)*: The Normalized Scanpath Saliency (NSS) metric, firstly introduced in [164], is

given by the averaged normalized visual attention at fixated locations. Given an attention map VAM , it is computed as follows:

$$NSS = \frac{1}{N_{TP}} \sum_{i=1}^{N_{TP}} \frac{VAM_i - \mu_{VAM}}{\sigma_{VAM}} \quad (4.1)$$

where N_{TP} is the total number of **GT** fixations (**TPs**), and μ_{VAM} and σ_{VAM} represent the mean and standard deviation of the **VAM** values, respectively. **NSS** is sensitive to **FPs** and monotonic transformations of the map, in contrast to classical **AUCs** scores, so it constitutes an interesting complement to these. For the evaluation of the methods considered in the thesis, we make use of the **sNSS** version of this score, proposed by Leborán et al. [103]. Unlike in the original **NSS**, N_{boot} sets of **FPs** sampled from fixated locations in different images or videos are obtained, in order to compute a mean $\mu_{VAMi} = \frac{1}{M} \sum_{m=1}^M VAM_m$ and a standard deviation $\sigma_{VAMi} = \sqrt{\frac{1}{M-1} \sum_{m=1}^M (VAM_m - \mu_{VAMi})^2}$, where $VAM_m \in \{TP \cup FP\}$ and $M = \text{card}(TP \cup FP)$. Then, **sNSS** compensates the center bias effect, by computing the average of N_{boot} scores for each frame:

$$sNSS = \frac{1}{N_{boot} N_{TP}} \sum_{i=1}^{N_{boot}} \sum_{j=1}^{N_{TP}} \frac{VAM_j - \mu_{VAMi}}{\sigma_{VAMi}} \quad (4.2)$$

Positive **sNSS** indicates correspondence between maps above chance, while a high number of **FPs** drives the overall **sNSS** down.

In order to evaluate the performance of visual attention models in a particular video, a probabilistic map that consists of fixations in frames from all other videos in the dataset is used as shuffle map for both scores. Figure 4.1 shows how **TPs** (green locations) and **FPs** (red locations) are sampled in an example frame taken from DIEM [16] database. As can be appreciated in the shuffle map, viewers have a tendency to look at the center of the image, as discussed above. Hence, more **FPs** are sampled close to the center of the frame, which prevents **sAUC** and **sNSS** metrics from being affected by the center bias. Moreover, 95% confidence bounds are provided for both metrics used.

Finally, for comparison purposes, we have considered the three baseline models introduced by Judd et al. in [165]:

- **CHANCE**: The model generates a **VAM** for each frame by randomly selecting some pixels as salient, which leads to a poor performance.

- **CENTER:** The model consists in a stretched symmetric [2D](#) Gaussian distribution centered on the frame, in such a way that closer locations to the center are more salient. This model serves as a good indicator to determine if the evaluation metrics used are affected by center bias or not.
- **H50:** For each frame, the model generates a [VAM](#) that contains the fixations of the 50% of subjects available. It constitutes a good realistic upper bound, which puts into perspective the efficiency of the assessed approaches.

4.2.4 Model initialization

Due to the stochastic nature of our approach, a correct initialization of the parameters is important to both fasten the convergence and reach an optimal model. As the goal is to learn sub-tasks that either attract or inhibit attention, we initialize basic, novelty and [CNNs](#)-based feature distributions as follows: we initialize some topics that inhibit and other that attract visual attention, with $\mu_{kl} = 0$ and $\mu_{kl} = 1$, respectively (remember that our features are maps in the range $[0, 1]$). Then, in order to provide initial variances for the topics, we compute two separate sets of variances with respect to $\mu_{kl} = \{0, 1\}$, from non-attracting and attracting locations respectively. Then, we run a separate *k-means* over the variance values and obtain the corresponding K centroids, one per topic. For camera motion features, the parameters \mathbf{c}_k are randomly initialized with values close to 0 whereas, as we have already mentioned, Σ_k is empirically set to $\Sigma_k = \text{diag}(0.25)$. Finally, discrete distribution features for object detection are initialized uniformly for every region in the non-uniform grid.

Last but not least, the main parameter of the proposed model is the number K of sub-tasks or topics that contribute to model visual attention. For simplicity, we have used the same number of attracting and inhibiting topics in our initialization. As indicated in the next sub-section, $K = 60$ is the number of topics used for the rest of the experiments. Initial global topic proportions α have been empirically set to $\alpha_k = 0.01$.

4.3 UNDERSTANDING VISUAL ATTENTION AS A MIXTURE OF SUB-TASKS

The most outstanding outcome of our probabilistic approach is determined by the topics inferred, which effectively help to interpret how visual attention works. Firstly, by means of the proportions in which those are blended, we can establish which sub-tasks are more prevailing for guidance. We have statistically estimated the

importance of each topic by examining the value η_k of the logistic regression model and the topic proportions ϕ_{nk} obtained for each spatial location n evaluated on the test set, as both variables are linearly related to the model response which generates the visual attention map. In particular, the relevance score of each sub-task k is computed as:

$$\mathcal{S}_k = \eta_k \sum_{n=1}^N \phi_{nk} \quad (4.3)$$

Scores are later normalized between $[-1, 1]$ to simplify the analysis.

Secondly, regarding the distribution parameters learned for features considered as input, we can further study the meaning of the sub-tasks, providing useful information about the most conspicuous regions in a given scenario. For the sake of interpretability, it should be noted that we have not considered CNNs-based features in this analysis, since they constitute very high-level representations at different scales whose content is more difficult to understand. Besides, since we know in advance for all features considered in the experiments that low feature values (close to 0) correspond to non-salient locations in frames, while regions with high feature values (close to 1) are very salient, Gaussians' means are not learned and remain fixed in $\mu_{kl} = 0$ and $\mu_{kl} = 1$ during the whole inference process. Then, we consider topics centered in $\mu_{kl} = 0$ and $\mu_{kl} = 1$ as those topics inhibiting (IT) or attracting attention (AT), respectively. Furthermore, the camera motion distribution has been also removed from the analysis as it has been observed that there is not a strong influence of this feature in any of the categories, since parameters \mathbf{c}_k learned for the most prevailing topics have all similar values. Under this simplified scenario, we can evaluate the relevance of *basic and novelty features*, using their learned standard deviation values σ_{kl} :

$$\mathcal{S}_{kl}^C = \frac{\sigma_l^F}{\sigma_{kl}} \quad (4.4)$$

with values in the range $[0, +\infty)$. Given a sub-task k and under our simplified scenario with fixed means ($\mu_{kl} = \{0, 1\}$), a feature l will be representative if its standard deviation σ_{kl} is lower compared to the deviation σ_l^F measured on areas that correspond with the topic type F (fixated areas if the topic is attracting attention, and viceversa).

Moreover, scores for *object-based features* are calculated by computing the cumulative probability of the cells that lie inside the detected bounding box ($r > 0$, excluding the background cell):

$$\mathcal{S}_{kl}^D = \sum_{r=1}^R \beta_{klr} \quad (4.5)$$

with values between $[0, 1]$.

Scores obtained by the three most noteworthy attraction and inhibition sub-tasks for some video genres in CRCNS-ORIG [15] and DIEM [16] databases are shown in Figures 4.2 and 4.3. Significant sub-tasks deduced for the rest of categories included in these databases are gathered in Appendix B. For each category, *ITs* ($\mu_{kl} = 0$) are represented in red on the left side of the bar graphs, while *ATs* ($\mu_{kl} = 1$) appear in blue on the right side. Then, the relevance score S_k of each sub-task k is indicated on top of its graph. Moreover, sub-tasks are represented as combinations of some of the features described in Section 3.5.2: basic and novelty features, such as *color* (C), *intensity contrast* (I), *orientation* (O), *velocity* (M), *acceleration* (A), *luminance spatial coherence* (SC (Lum.)), *motion spatial coherence* (SC (Mot.)), *luminance temporal coherence* (TC (Lum.)), *motion temporal coherence* (TC (Mot.)), *luminance spatio-temporal coherence* (STC (Lum.)), *motion spatio-temporal coherence* (STC (Mot.)); and object-based features, such as *frontal* (F) and *profile faces* (PF), *upper bodies* (B), *pedestrians* (P) and *text* (T). Each bar is associated to a feature score S_{kl}^C , for basic and novelty features, or S_{kl}^D , for object-based features. High values of scores in *ITs* correspond to inhibiting features, which reduce the attentional response. In contrast, high values of scores in *ATs* highlight those features that are more attracting for each category.

Although the number of topics experimentally determined is quite high ($K = 60$), we have observed that only few of them are responsible of guiding attention most of the time, whereas the rest are intended to refine the estimation, specially in the less prevalent sequences.

As can be seen, different sub-tasks are determined to model visual attention in each scenario, existing an appreciable contrast between well-separated categories such as *Outdoor* or *TV News*, which involve distinctive actions (see Figure 4.2). While *context-generic* models are adjusted to the most prominent events in the databases, which consist of faces noticeable by their color and intensity, and motion objects, *context-aware* models have the ability of attaining more particular and explainable activities. Motion and acceleration features are relevant in *Outdoor* (Figure 4.2a), and *Videogames* (Figure B.3) sub-tasks, which could be related to people or characters walking or running. In contrast, faces and texts are more attractive and predominant in categories like *Commercials* (Figures 4.3a, B.4), *TV News* (Figures 4.2b, B.16) and *Talk Shows* (Figures B.7, B.13). Both motion and faces are eye-catching in *Sports* (Figures 4.3b, B.6) videos, which often consist of real-time edited outdoor scenes to be released on TV. Finally, low values of spatial and temporal coherency features are mostly frequent in *IT*, which implies reducing the attentional response in usual and stable locations over space and time.

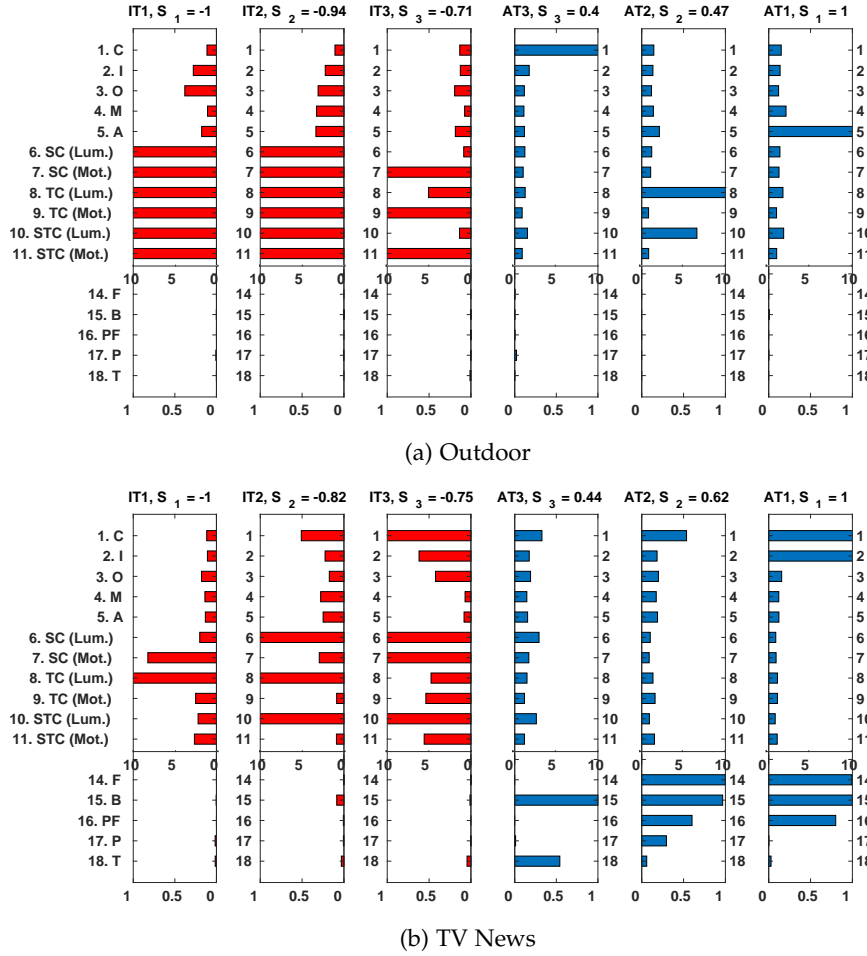


Figure 4.2: Three most prominent attracting (AT) and inhibiting (IT) sub-tasks inferred by (a) *Outdoor* and (b) *TV News* context-aware models learned based on CRCNS-ORIG [15] database. For each category, ITs ($\mu_{kl} = 0$) are shown in red on the left side of the bar graph, while ATs ($\mu_{kl} = 1$) appear in blue on the right side. Then, the relevance score S_k of each sub-task k is indicated on top of its graph. Moreover, sub-tasks are represented as combinations of some of the features described in Section 3.5.2: basic and novelty features, such as *color* (C), *intensity contrast* (I), *orientation* (O), *velocity* (M), *acceleration* (A), *luminance spatial coherence* (SC (Lum.)), *motion spatial coherence* (SC (Mot.)), *luminance temporal coherence* (TC (Lum.)), *motion temporal coherence* (TC (Mot.)), *luminance spatio-temporal coherence* (STC (Lum.)), *motion spatio-temporal coherence* (STC (Mot.)); and object-based features, such as *frontal* (F) and *profile faces* (PF), *upper bodies* (B), *pedestrians* (P) and *text* (T). Each bar is associated to a feature score S_{kl}^C , for basic and novelty features, or S_{kl}^D , for object-based features. High values of scores in ITs correspond to inhibiting features, which reduce the attentional response. In contrast, high values of scores in ATs highlight those features that are more attracting for each category.

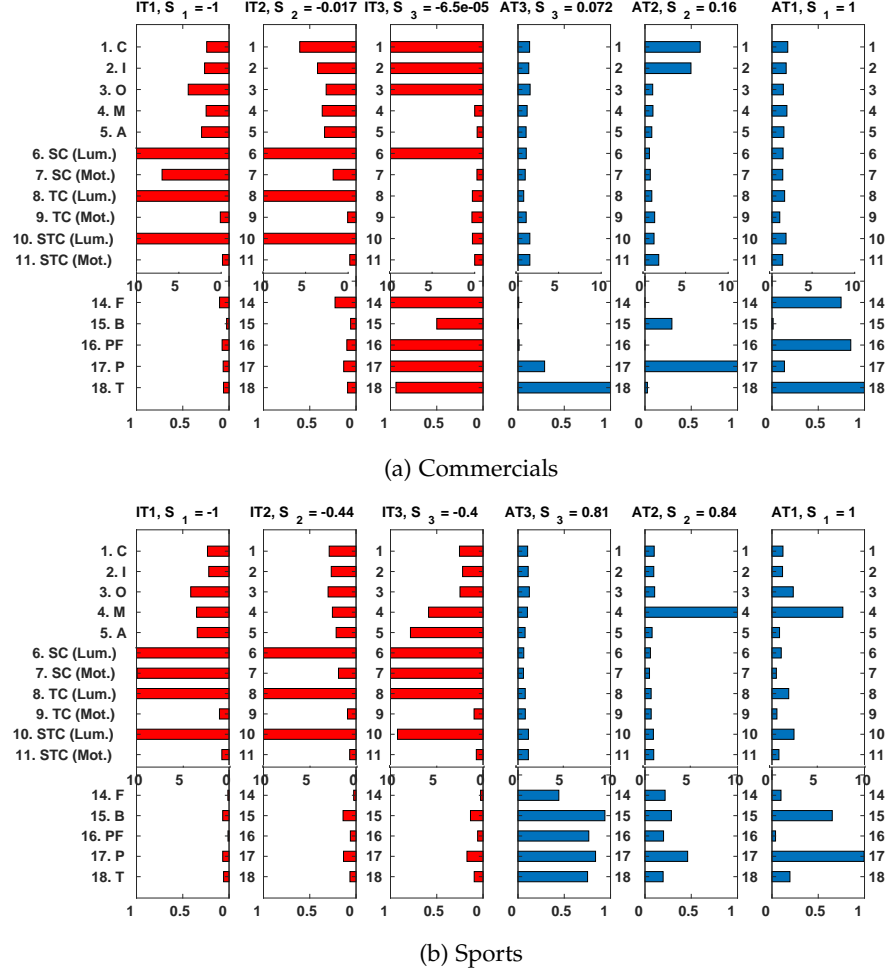


Figure 4.3: Three most prominent attracting (AT) and inhibiting (IT) sub-tasks inferred by (a) *Commercials* and (b) *Sports* context-aware models learned based on DIEM [16] database. For each category, ITs ($\mu_{kl} = 0$) are shown in red on the left side of the bar graph, while ATs ($\mu_{kl} = 1$) appear in blue on the right side. Then, the relevance score S_k of each sub-task k is indicated on top of its graph. Moreover, sub-tasks are represented in the graphs as combinations of some of the features described in Section 3.5.2: basic and novelty features, such as *color* (C), *intensity contrast* (I), *orientation* (O), *velocity* (M), *acceleration* (A), *luminance spatial coherence* (SC (Lum.)), *motion spatial coherence* (SC (Mot.)), *luminance temporal coherence* (TC (Lum.)), *motion temporal coherence* (TC (Mot.)), *luminance spatio-temporal coherence* (STC (Lum.)), *motion spatio-temporal coherence* (STC (Mot.)); and object-based features, such as *frontal* (F) and *profile faces* (PF), *upper bodies* (B), *pedestrians* (P) and *text* (T). Each bar is associated to a feature score S_{kl}^C , for basic and novelty features, or S_{kl}^D , for object-based features. High values of scores in ITs correspond to inhibiting features, which reduce the attentional response. In contrast, high values of scores in ATs highlight those features that are more attracting for each category.

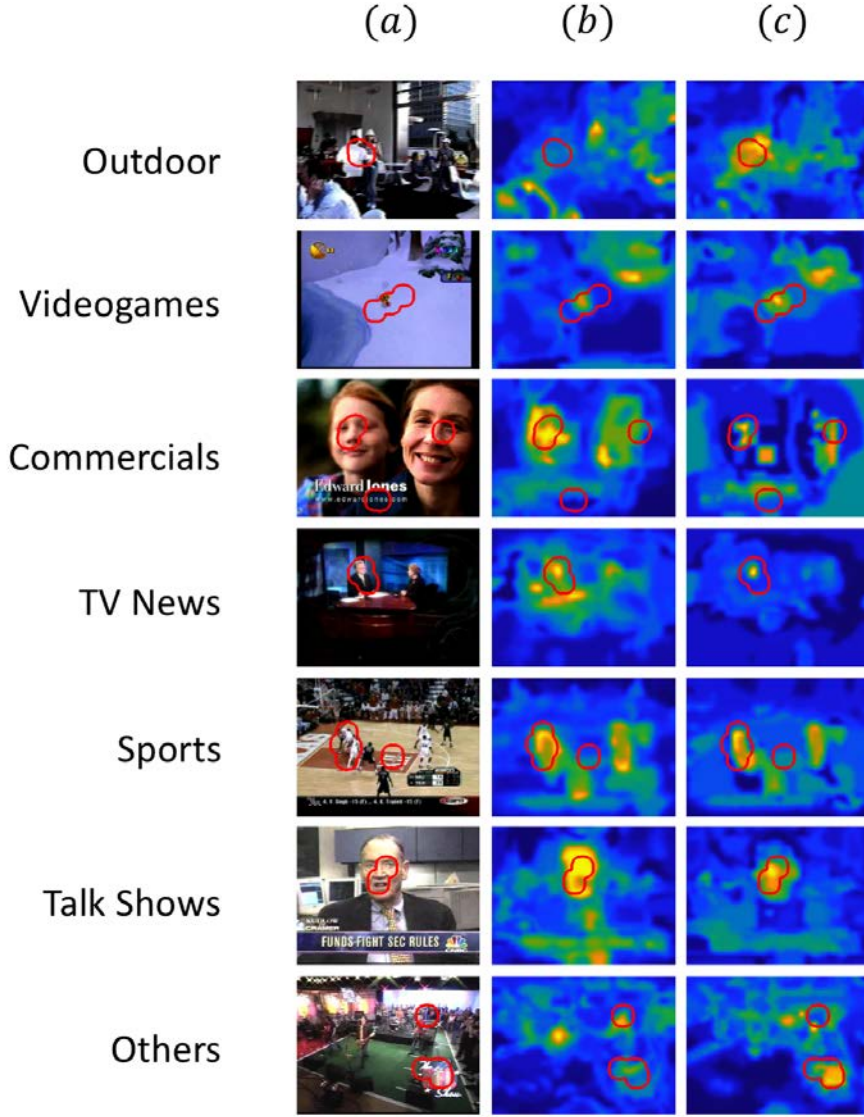


Figure 4.4: Visual attention maps obtained by [ATOM](#) for some example frames from CRCNS-ORIG [15] database. Red boundaries highlight high-density regions of human fixations in the [GT](#) map. (a) Original frames. (b) Context-Generic. (c) Context-Aware.

4.4 RESULTS ON VISUAL ATTENTION ESTIMATION

In this second set of experiments, [CNN](#)s-based features are included and the [ATOM](#) model learns unconstrained Normal distributions without fixating the means.

Results obtained for the two versions of our method in each category and for each database are provided in Figures 4.5 and 4.7, respectively. As can be seen, the *context-aware* models match or outperform the *generic* approach in all genres. Without considering *Others* category in CRCNS-ORIG [15] database, which is more diverse and contains a synthetic saccade test video, best scores are

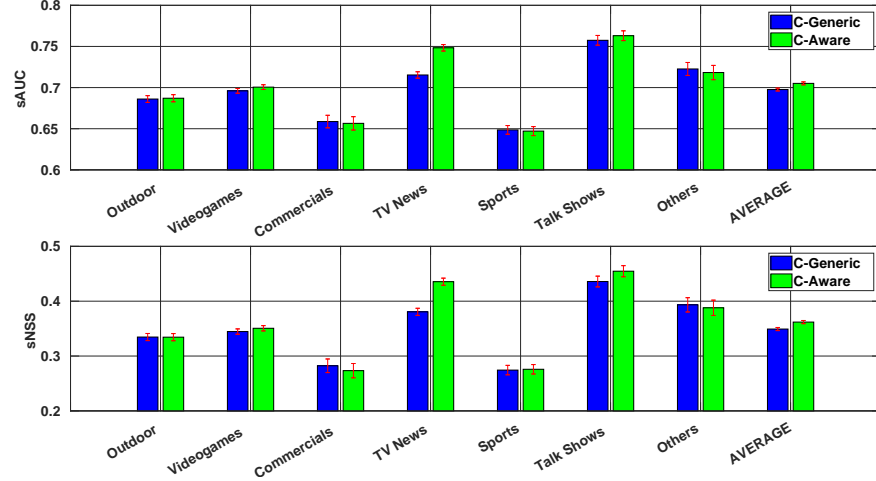


Figure 4.5: Results obtained by the proposed *context-generic* and *context-aware* ATOM models in the CRCNS-ORIG [15] database, which consist of $K = 60$ topics.

obtained for *TV News*, *TVShows* and *Talk Shows* genres, due to the high impact of object detectors (faces, pedestrians) in this genres, as shown in some of the examples provided in Figures 4.4 and 4.6. Scores achieved for *Outdoor* and *Videogames* videos are also remarkable, due to the strong influence assigned to motion-related features. This reinforces the idea that, depending on the context, certain particular sub-tasks aid to guide visual attention. This can be also noticed if we look at the results obtained in categories such as *Others* or *Commercials*, whose associated videos cover a wide variety of contents and thus are not closely related, so consequently it has been hard to find out meaningful topics. In fact, the results for the *context-generic* model in these cases are higher. From our point of view, this undesired effect might come from the fact that C-G has been trained on a wider set of videos than C-A approaches, and therefore has got better generalization. Therefore, it can be concluded that it is necessary to establish well-defined application scenarios where to determine these feature-based representations. In order to provide a fair comparison, we draw on the same number of topics for each of the categories in the dataset chosen, although it has been observed that the performance also depends on the complexity of the scenarios. If we compare the average performance of *context-aware* models with respect to the result obtained by the *context-generic* approach, there is an improvement of 4.6% in terms of sNSS and 1.1% in terms of sAUC, which is closer to the upper threshold given by H50 score.

Thus, we can state that specific *context-aware* representations of visual attention learned over particularized training sets (the training videos belonging to each category) work better than *generic* models learned over larger general datasets (including all video

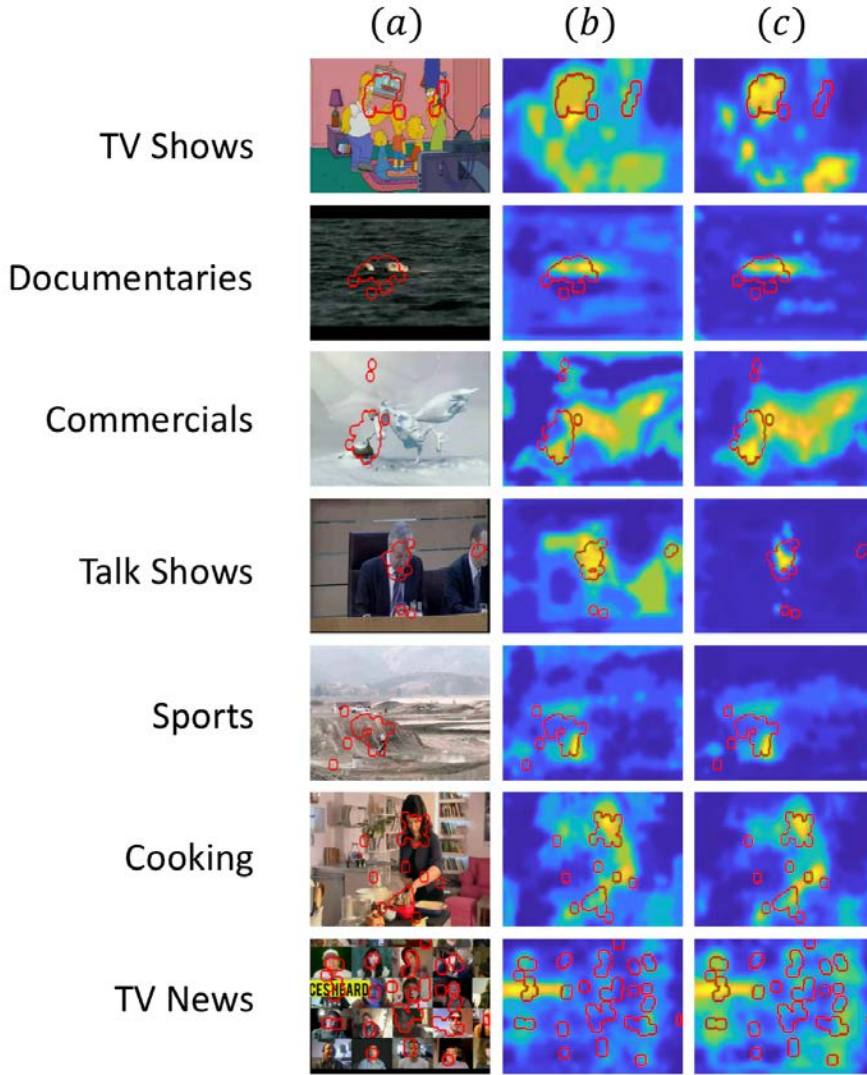


Figure 4.6: Visual attention maps obtained by [ATOM](#) for some example frames from DIEM [16] database. Red boundaries highlight high-density regions of human fixations in the [GT](#) map. (a) Original frames. (b) Context-Generic. (c) Context-Aware.

categories). Based on these results, from now on we will use the *context-aware* version of our algorithm to provide a comparison with other approaches in the state-of-the-art.

4.5 COMPARISON WITH STATE-OF-THE-ART METHODS

With the aim of assessing the performance of our approach in comparison with other methods available in the state-of-the-art, we have selected 17 static and dynamic visual attention models, which are representative of the existing diversity for visual attention prediction: we have included both [BU](#) and [TD](#) or learnable models, a model that uses [CNNs](#) to predict, etc., as well as the three reference

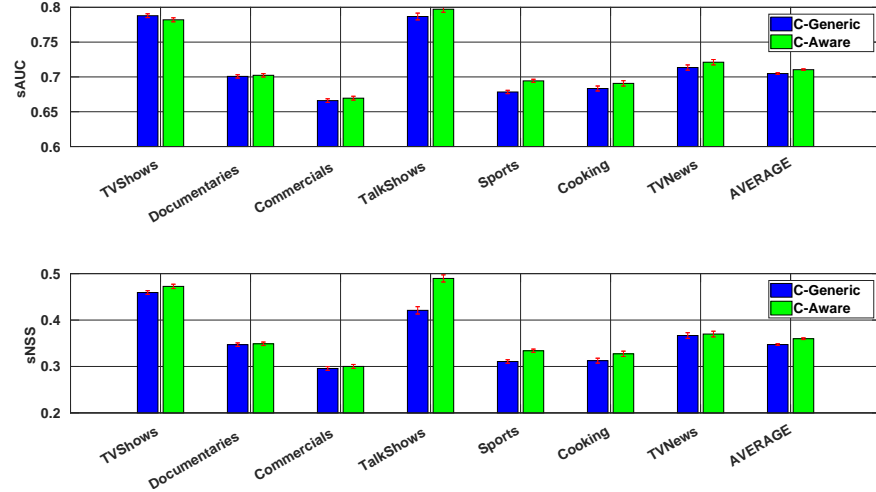


Figure 4.7: Results obtained by the proposed *context-generic* and *context-aware* ATOM models in the DIEM [16] database, which consist of $K = 60$ topics.

models introduced in section 4.2.3 (H50, CHANCE, CENTER). Parameters used are the ones set as default by authors. As can be verified from CENTER baseline, both metrics included in the analysis are not affected by center bias effect.

Tables 4.2 and 4.3 contain all the results obtained for the assessed methods in CRCNS-ORIG [15] and DIEM [16] databases, respectively, together with those reached by the system proposed in Chapter 3 (ATOM). We also include on the list the first approach we presented in [125], which make use of a linear regressor to estimate visual attention instead of the logistic regressor currently employed, as well as other features. Features and number of topics ($K = 40$) taken for this previous configuration are those reported in [125].

The improvement achieved by our model with respect to very recent approaches such as AWS-D [103], DCL [2], WMAP [96] or ICL-D [88] is statistically significant. Moreover, it is also visually noticeable in some intricate cases, as those shown in Figure 4.8, with scenes showing crowds, multiple similar concepts that hamper visual guidance or quick actions.

Finally, we evaluate the computational time on a system with an Intel Core i7-6700K CPU at 4.00GHz and with 32GB of RAM. Regarding our approach, we should distinguish between the learning and the test phase. Both phases involve a feature extraction stage that takes 5.81s per frame, which has not been optimized, and could be highly parallelized by GPUs. Time spent in the learning phase depends on the number of topics of the model trained and the amount of input frames. For instance, training a model with $K = 60$ topics and ~ 3000 frames would take $\sim 45min$. This time can be reduced if the number of topics is decreased to $K = 40$ ($\sim 32min$) or

Table 4.2: Comparison with state-of-the-art methods in the CRCNS-ORIG [15] database.

Model	Learning	$sAUC$	$sNSS$
		$mean (C.I.)^{Rank}$	$mean (C.I.)^{Rank}$
ATOM	YES	0.705 (0.703, 0.707)¹	0.362 (0.359, 0.365)¹
AWS-D [103]	NO	0.700 (0.698, 0.702) ²	0.322 (0.319, 0.325) ³
DCL [2]	YES	0.684 (0.682, 0.686) ³	0.323 (0.320, 0.326) ²
AWS [98]	NO	0.675 (0.674, 0.677) ⁴	0.281 (0.278, 0.285) ⁴
WMAF [96]	NO	0.670 (0.669, 0.672) ⁵	0.236 (0.232, 0.239) ¹²
Hou and Zhang [85]	NO	0.669 (0.667, 0.671) ⁶	0.260 (0.257, 0.263) ⁷
DCL+ [2]	YES	0.666 (0.665, 0.668) ⁷	0.255 (0.251, 0.258) ⁸
ICL-D [88]	NO	0.666 (0.665, 0.668) ⁸	0.217 (0.214, 0.220) ¹⁴
PQFT [93]	NO	0.662 (0.660, 0.663) ⁹	0.243 (0.240, 0.246) ¹¹
Goferman [94]	NO	0.661 (0.659, 0.662) ¹⁰	0.263 (0.260, 0.266) ⁶
SUN [87]	YES	0.654 (0.652, 0.655) ¹¹	0.251 (0.248, 0.254) ⁹
AIM [58]	YES	0.653 (0.652, 0.655) ¹²	0.270 (0.268, 0.273) ⁵
Torralba [82]	NO	0.648 (0.646, 0.650) ¹³	0.251 (0.248, 0.254) ¹⁰
Itti (ST) [81] [84]	NO	0.634 (0.632, 0.636) ¹⁴	0.217 (0.214, 0.220) ¹⁵
Fernández-Torres [125]	YES	0.628 (0.626, 0.630) ¹⁵	0.218 (0.215, 0.221) ¹³
SDSR [89]	NO	0.627 (0.625, 0.628) ¹⁶	0.129 (0.126, 0.132) ¹⁷
GBVS (ST) [84]	NO	0.621 (0.619, 0.623) ¹⁷	0.182 (0.179, 0.186) ¹⁶
ESA-D [92]	NO	0.541 (0.539, 0.543) ¹⁸	0.075 (0.072, 0.078) ¹⁸
H50	NO	0.800 (0.799, 0.802)	0.679 (0.677, 0.681)
CHANCE	NO	0.500 (0.500, 0.500)	-0.000 (-0.000, 0.000)
CENTER	NO	0.509 (0.507, 0.511)	0.057 (0.054, 0.060)

$K = 20$ ($\sim 18min$), which would slightly decrease the performance. Then, in the test phase, the average time per frame is only 0.157s, which is competitive compared to those obtained by the two next best methods, AWS-D [103] (0.075s) and DCL [2] (0.2s).

4.6 WHERE WE ARE: MODEL STRENGTHS AND LIMITATIONS

Despite the improvement reached by the proposed model over the *state-of-the-art* and the compelling information it provides, we are still far from reaching human capacity of almost immediately selecting the most essential elements and areas to reach a full understanding in a given scenario, or to solve a particular task, according to the H50 score reflected in Tables 4.2 and 4.3. Nonetheless, we advocate that the inclusion of an intermediate level between features and visual attention in terms of sub-tasks is a powerful way towards comprehensible guiding representations.

We have demonstrated that some of the traditional basic features used (e.g. color, orientation, motion) are still useful in many cases to predict visual attention in videos. Furthermore, thanks to the object

Table 4.3: Comparison with state-of-the-art methods in the DIEM [16] database.

Model	Learning	sAUC	sNSS
		<i>mean (C.I.)^{Rank}</i>	<i>mean (C.I.)^{Rank}</i>
ATOM	YES	0.710 (0.709, 0.712)¹	0.360 (0.358, 0.362)¹
AWS-D [103]	NO	0.701 (0.700, 0.701) ²	0.319 (0.317, 0.320) ³
DCL [2]	YES	0.695 (0.695, 0.696) ³	0.341 (0.340, 0.342) ²
DCL+ [2]	YES	0.683 (0.682, 0.683) ⁴	0.318 (0.317, 0.319) ⁴
WMAP [96]	NO	0.666 (0.666, 0.667) ⁵	0.233 (0.232, 0.234) ¹²
Hou and Zhang [85]	NO	0.663 (0.662, 0.664) ⁶	0.247 (0.246, 0.248) ⁹
PQFT [93]	NO	0.662 (0.661, 0.663) ⁷	0.235 (0.233, 0.236) ¹¹
Goferman [94]	NO	0.659 (0.658, 0.660) ⁸	0.257 (0.256, 0.258) ⁶
GBVS (ST) [84]	NO	0.653 (0.652, 0.653) ⁹	0.256 (0.255, 0.257) ⁷
AWS [98]	NO	0.652 (0.651, 0.653) ¹⁰	0.271 (0.270, 0.272) ⁵
Itti (ST) [81] [84]	NO	0.638 (0.637, 0.639) ¹¹	0.253 (0.252, 0.254) ⁸
Torralba [82]	NO	0.636 (0.635, 0.636) ¹²	0.232 (0.231, 0.233) ¹³
SUN [87]	YES	0.631 (0.630, 0.631) ¹³	0.220 (0.219, 0.221) ¹⁴
Fernández-Torres [125]	YES	0.630 (0.629, 0.631) ¹⁴	0.219 (0.217, 0.222) ¹⁵
ICL-D [88]	NO	0.629 (0.628, 0.629) ¹⁵	0.154 (0.153, 0.155) ¹⁶
AIM [58]	YES	0.618 (0.617, 0.618) ¹⁶	0.238 (0.237, 0.239) ¹⁰
ESA-D [92]	NO	0.563 (0.562, 0.564) ¹⁷	0.150 (0.149, 0.151) ¹⁷
SDSR [89]	NO	0.531 (0.530, 0.532) ¹⁸	0.022 (0.021, 0.023) ¹⁸
H50	NO	0.827 (0.827, 0.827)	0.662 (0.662, 0.663)
CHANCE	NO	0.500 (0.500, 0.500)	−0.000 (−0.000, 0.000)
CENTER	NO	0.503 (0.502, 0.504)	0.054 (0.052, 0.055)

detectors introduced and the corresponding spatial discrete distributions, we are able to model simple but attractive concepts such as faces or text, putting emphasis on their most noticeable elements. The high performance achieved by these detectors in some categories leads us to reckon the integration of large-scale hierarchical networks for object recognition in future revisions of our model, such as the ones evaluated in the ImageNet Challenge [166]. In addition, there is also a need of a deeper understanding of the scene, establishing relations between recognized concepts both in the same frame or in different frames. This would enable the system to enhance guidance in situations where many conspicuous regions exist and selecting the most significant in the task to be solved, or even an intermediate one (e.g. Figure 4.9(a)); when objects are occluded during few frames (e.g. Figure 4.9(b)); or to determine the sequence of objects or subjects to follow in order to interpret a scene (e.g. Figure 4.9(c)), among others. In other words, we pursue the identification and modeling of sub-tasks, not only over space but also along time.

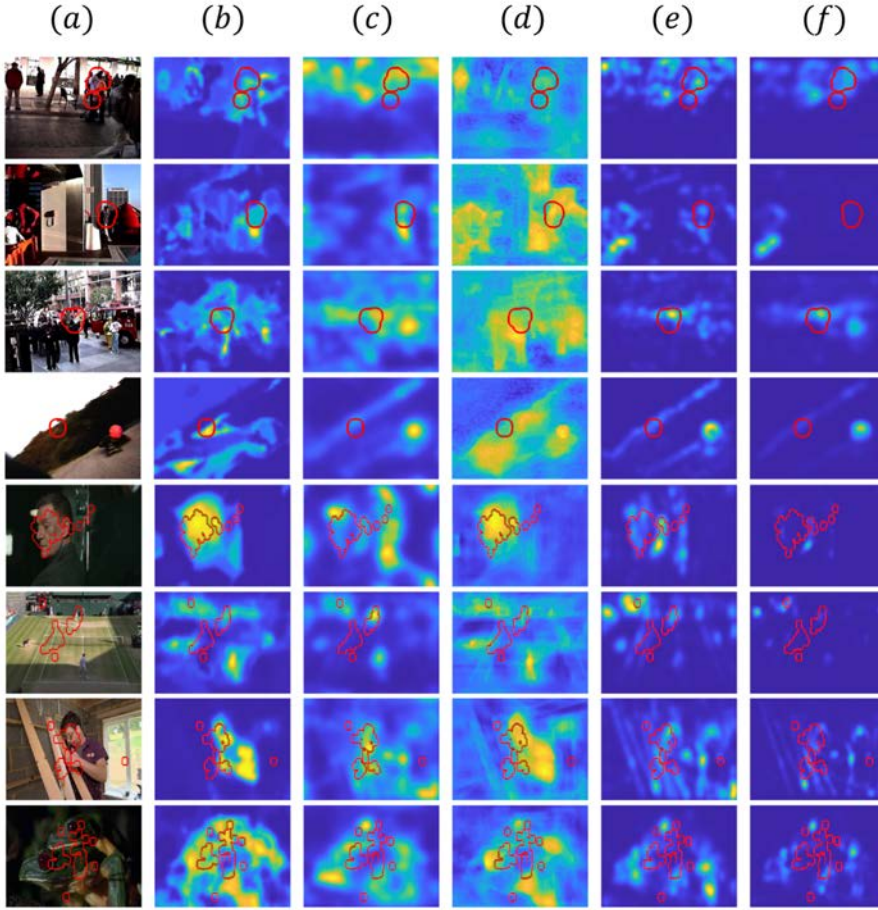


Figure 4.8: Visual attention maps generated by some of the most outstanding methods in the *state-of-the-art* for some intricate example frames taken from CRCNS-ORIG [15] and DIEM [16] databases. Red boundaries highlight high-density regions of human fixations in the GT map. (a) Original frames. (b) ATOM. (c) AWS-D [103]. (d) DCL [2]. (e) WMAP [96]. (f) ICL-D [88].

Finally, the importance of GT eye fixations has to be discussed, both in learning and evaluation stages. As it can be appreciated in some of the examples gathered throughout the chapter, not all fixations contain useful information to train a visual attention system, not only because the occlusions mentioned above, but also due to errors during the eye-tracker acquisition or to the observers' center bias present in many frames. Fixations often fall on edges, not covering completely some objects of interest, such as gaming characters or players, which are essential to infer sub-tasks. Additionally, we might take into account covert attention, which is independent of eye movements and stresses the existence of attention independent of gaze change. Hence, techniques to filter and, if necessary, to extend regions considered as GT should be regarded in upcoming experiments. What is more, existing evaluation metrics do not seem to be appropriate in situations such

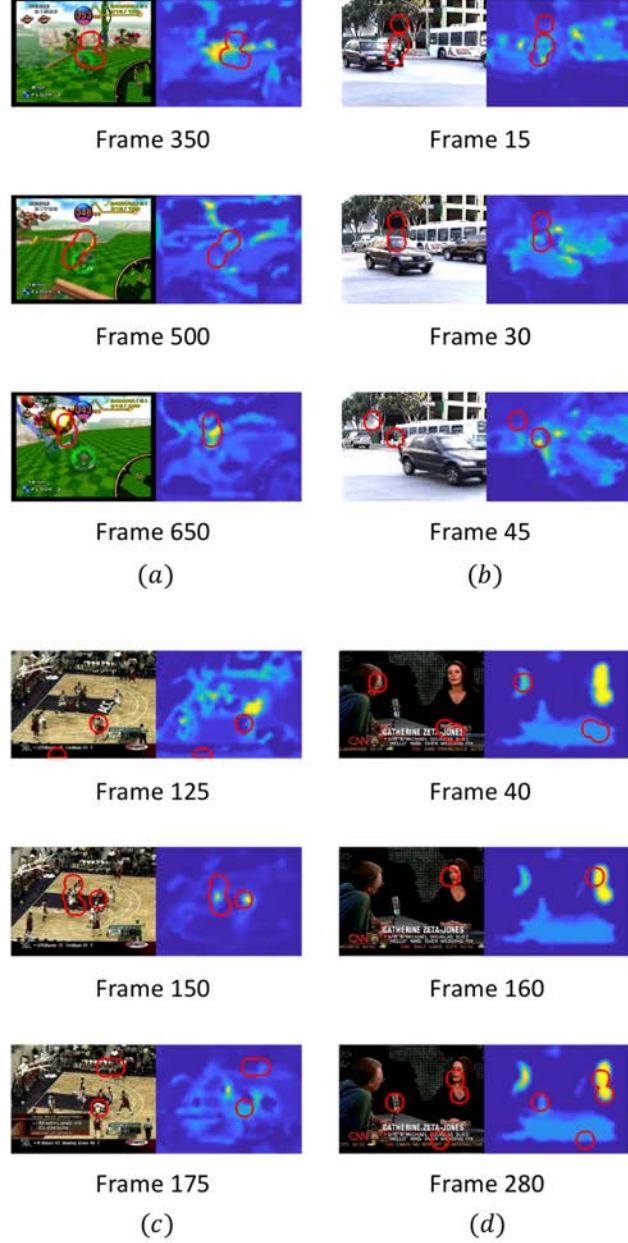


Figure 4.9: Frame sequences taken from CRCNS-ORIG [15] database to analyze some [ATOM](#) model drawbacks and define future lines of research. Red boundaries highlight high-density regions of human fixations in the [GT](#) map, both in original frames and computed visual attention maps. (a) Videogames scenario where many remarkable regions exist, making observers constantly shift their gaze. (b) Outdoor scenario where multiple salient concepts (e.g. car, policeman) overlap each other. (c) Basketball match, in which the sequence of players to follow is decisive to model visual attention. (d) TV talk show, where several quasi-static concepts appear together during a long time lapse and estimated visual attention is either distributed among all or focused in one of them.

as the one shown in Figure 4.9(d), where many remarkable quasi-static concepts appear together during a long time lapse and estimated visual attention is either distributed among all or focused in one of them. If observers' fixations are widely dispersed and attention is switched between various locations, what should be the GT taken for each frame in this case? Should all concepts be considered as attracting during the whole video fragment? We will seek to address these issues in future application scenarios.

4.7 CONCLUSIONS

In Chapters 3 and 4, we have presented a hierarchical probabilistic framework to estimate and understand TD visual attention in videos. Relying on the idea of 'guiding representation' supported by some of the most prevailing psychological theories about visual attention, our ATOM model decomposes it into mixtures of several latent topics or sub-tasks, which are in turn modeled as combinations of low-, mid- and high-level spatio-temporal features obtained from video frames. For that purpose, an intermediate level between feature extraction and visual attention computation phases is introduced, aligning the latent discovered sub-tasks from frames to the information drawn from human fixations. The attention response is thus generated by computing a logistic regression model over topic proportions. It is also worth mentioning that the definition of the method is generic and independent of the input features, which enables an easy adaptation to any application scenario.

The ability of ATOM to successfully learn specifically adapted hierarchical representations of visual attention in diverse contexts has been demonstrated on the basis of a wide set of features. Either classical and easily interpretable feature maps, which have been effective to extract conclusions about the existing scenarios in the well-known CRCNS-ORIG [15] and DIEM [16] databases, or those generated by recently adopted CNNs structures, which allow to capture more complex concepts, have aided to significantly outperform other competent methods in the literature. Moreover, the detection of simple elements such as faces or text, and their modeling through spatial discrete distributions, has led to improve visual attention estimation in certain challenging situations.

Experimental results show the advantage of obtaining comprehensible guiding representations to model visual attention. However, it is still necessary to deepen in some of the stages of the framework, carefully selecting the most meaningful information from fixated regions in the scene, and integrating more robust recognition and understanding techniques that enable to identify more accurate sub-tasks over space and time. To that end, future efforts will be directed towards task-driven approaches, developing

video databases with human fixations to test the usefulness of the system in end-user applications.

DEEP NEURAL NETWORKS FOR MODELING VISUAL ATTENTION IN THE TEMPORAL DOMAIN

5.1 INTRODUCTION

Observers' eye movements constitute a useful source to understand how visual attention works and, consequently, what information should be selected for further processing. Given a complex and crowded scenario, if all observers fix their attention at the same location at the same time, it is very likely that something noticeable is happening.

Photographs in Figure 5.1(a) illustrate a typical video monitoring room. The task of CCTV operators in this scenario is to find a potentially anomalous event (e.g. robberies, road accidents, etc.) amongst multiple distractors (e.g. crowds, similar vehicles, etc.) displayed at the same time in a large array of 20 to more than 500 screens [167], as the one shown in Figure 5.1(b). At what screen should they look each time? Operators often have also to review many hours of surveillance recordings. Moreover, anomalous events seldom happen, which makes these tasks even more difficult to solve. How do they tackle these tasks? Is it possible to develop a system to aid experts to perform them more efficiently? We want to meet these challenges, taking advantage of eye fixations.

In this chapter, inspired by the recent success of Convolutional Neural Networks (CNNs) for learning deep hierarchical image representations and Long Short-Term Memory (LSTM) units for time series forecasting, we propose a network architecture that goes from spatio-temporal visual attention prediction to temporal attention estimation. Visual attention in the temporal domain can be understood as a filtering mechanism, which allows to select time segments of special importance in video sequences.

Supported by the fact that eye fixation sequences of different viewers correlate well when an important or anomalous event happens, our system models visual attention over time as a fixation-based response. Hence, it could be used to prevent human errors and speed up decision making processes in real applications



Figure 5.1: (a) Typical video surveillance monitoring room. Image taken from [168]. (b) The task of a CCTV operator in a video monitoring room is to find a potentially anomalous event amongst multiple distractors, displayed at the same time in a large array of more than 20 screens. Anomalous events seldom happen, which makes this task even more complex to solve. We want to meet this challenge, taking advantage of eye fixations. Image taken from [169].

which require watching large amounts of visual information, at the same time and during long time periods, such as the task of video surveillance.

CHAPTER OVERVIEW

The chapter is organized as follows. First, in Section 5.2, we discuss about the importance of eye tracking to understand the behavior of experts performing real-world tasks. We review also the most relevant and recent work in visual attention estimation applying deep architectures, and present our main contributions. Then, an introduction to deep learning is carried out in the Section 5.3, describing primarily those DNNs and techniques used for our research. Next, three feature learning architectures for attention guidance are described in Section 2.4.3, which will provide input feature maps to our system for modeling attention in the temporal domain. Afterwards, we fully describe in detail the proposed system in Section 5.5. For that purpose, we first introduce our assumptions about task-driven visual attention, making an overview of the complete architecture in Section 5.5.1. Later, we define the process followed to generate a fixation-based temporal GT in Section 5.5.2. Finally, we describe the complete architecture design in Section 5.5.3, and elaborate the two stages of our system in Sections 5.5.4 and 5.5.5.

5.2 RELATED WORK

Following the early experiments of Yarbus in 1967 [25], who concluded that visual attention is ultimately task- or goal-driven,

many other investigations with still stimuli proved this statement [76, 77]; they claim that it might be attainable to infer the attentional processes carried out by the HVS from eye movement sequences. This has motivated researchers during the last two decades to evaluate the possibilities of eye tracking in real applications such as driving safety [6], aviation [7], production and industry control [78], health-care [79] and video surveillance [8].

Visual information is constantly being updated in videos and, consequently, not all the assumptions that hold in still images can be extrapolated to videos recorded in real situations. Indeed, a recent study with almost 150 participants monitoring footage in a Closed-Circuit TeleVision (CCTV) crime scenario [170] has concluded, after conducting experiments with task-oriented and non-task-oriented observers, that the complexity of dynamic environments may decrease the influence of in-advanced instructions, in contrast to what happens with static images.

Despite the above issues, high-valuable information is still appreciated when examining fixations behavior: their typical quasi-random pattern changes just before and during a significant or suspicious event. Moreover, there is a strong correlation between fixations of different viewers when these events happen in a similar scenario [78, 80], both in their location and duration. Taking into consideration this fact, we propose to take a step further by developing a system able to learn this behavior from sequences of fixations, which will be used to model the temporal dimension of visual attention.

The models introduced in this chapter goes from spatio-temporal 2D visual attention maps, used so far for fixation prediction at every frame in a video, and transforms them into temporal 1D visual attention curves. These signals highlight relevant frames in a video, which often correspond to surprising or unusual events, frequently labeled as anomalies. Therefore, our ultimate objective now is not to understand how visual attention works, as in Chapters 3 and 4, but to estimate visual attention in the temporal domain. In contrast to our previous approach, which drew on probabilistic LTM (see Section 3.4), some preliminary experiments suggested the convenience of using CNNs for better modeling spatio-temporal visual attention, as well as LSTM-based architectures for modeling attention in the temporal domain.

Spatio-temporal visual attention in real scenarios is still in its infancy, as we have remarked along this thesis. The exceptional, but still lacking in analysis, performance of CNNs has brought some new approaches to this application. They mainly attend to three attributes: a) spatial RGB-based features; b) motion, modeled either by using optical flow information at the input of a CNN [106] or by

means of recurrent LSTM units [107]; and c) objects, located in maps generated by multi-scale CONV architectures [105].

Visual attention in the temporal domain has been even less tackled in the literature up to date. First works elementarily modeled temporal attention as the mean of the saliency values predicted at each spatial location [171, 172]. Similarly, Ejaz et al. [173] modeled the temporal saliency of a frame as an average of temporal gradients, in order to select key frames in a video summarization application. More recently, Koutras et al. [174] defined a simple Otsu’s threshold operator to transform SMs into saliency curves. Finally, Han et al. proposed in [175] a more sophisticated probabilistic supervised method for temporal visual attention, which aims to estimate movie trailers attractiveness. This model uses the same hypothesis than our work and computes a temporal GT for each video frame based on the fixation dispersion across several observers.

5.3 DEEP NEURAL NETWORKS

Deep learning [176] covers a wide set of computational models with multiple processing layers. Starting from large amounts of raw information and by means of a general-purpose learning procedure, they discover multi-level abstract representations of the data which are able to effectively solve a variety of tasks.

Indeed, during the past five years, Deep Neural Networks (DNNs) [43] have revolutionized multiple applications using ML algorithms, such as speech recognition [52], object detection [3] and natural language understanding [177], dramatically improving the existing state-of-the-art performances. At present, the deep learning field is constantly growing, so new architectures and algorithms are being tested with the objective of either understanding the behavior of DNNs or designing more robust and less time- and computational-consuming systems.

This section makes an introduction to DNNs, mainly providing the description of the modules that comprise CNNs and LSTMs architectures. CNNs were already used for feature extraction in Chapter 3, and now they will constitute the first stages of the LSTM-based system for modeling visual attention in the temporal domain, which is the goal of this chapter. Furthermore, several common strategies for training and optimizing NNs are explained in the different subsections.

5.3.1 Neural Networks

The basic computational cell of the brain is the neuron, as already mentioned in Section 2.2.1. Neurons receive, process and transmit information to carry out different functions, such as visual

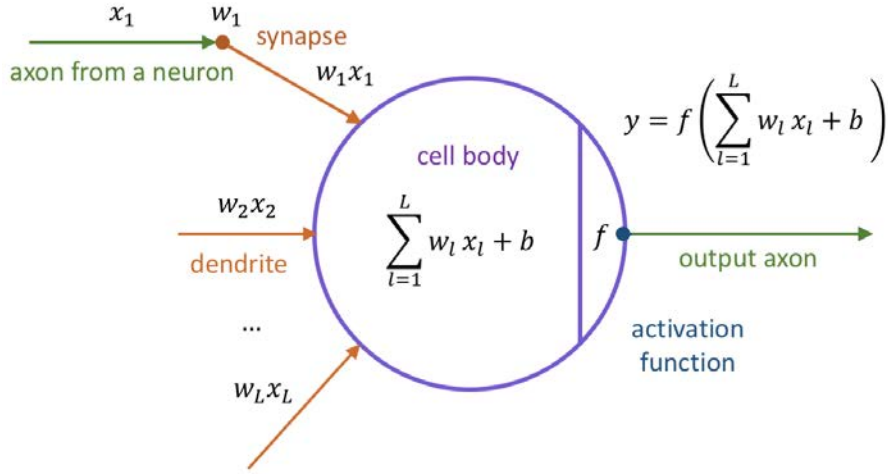


Figure 5.2: Mathematical model of a computational neuron with L inputs $\mathbf{x} = \{x_1, x_2, \dots, x_L\}$ and one output y , composed by a set of weights $\mathbf{w} = \{w_1, w_2, \dots, w_L\}$ and a bias term b . Adapted from [62].

perception, learning and memory [56]. Inspired by this biological fact, Neural Networks (NNs) are directed graphs of computational units, also named neurons.

Following a similar process than the one involved by neurons in the brain, whose schematic diagram was shown in Figure 2.2, a mathematical neuron, represented in Figure 5.2, computes a scalar output from a set of L input signals $\mathbf{x} = \{x_1, x_2, \dots, x_L\}$ [62].

A single neuron or unit is defined as a linear classifier:

$$y = f(\mathbf{w}^T \mathbf{x} + b) = f\left(\sum_{l=1}^L w_l x_l + b\right) \quad (5.1)$$

Indeed, it is composed by a learnable set of weights $\mathbf{w} = \{w_1, w_2, \dots, w_L\}$ and a bias term b that model synaptic strengths and control the excitatory (positive weight) or inhibitory (negative weight) influence of the neuron on subsequent neurons. Moreover, the firing rate of the biological neuron is represented by means of a non-linearity or activation function f . If the activation function corresponds to the sigmoid function ($y = \text{sigm}(\mathbf{w}^T \mathbf{x} + b)$), the operation of a neuron corresponds to a logistic regression.

The perceptron

The perceptron [178] was first introduced by Rosenblatt in 1958 and constitutes the simplest NN model, based on a unique neuron or unit, which allows to solve a binary classification problem. For the case of a multi-class classifier with C classes, a unit per class is needed. The perceptron is now defined in matrix form as follows:

$$\mathbf{y} = f(\mathbf{W}^T \mathbf{x} + \mathbf{b}). \quad (5.2)$$

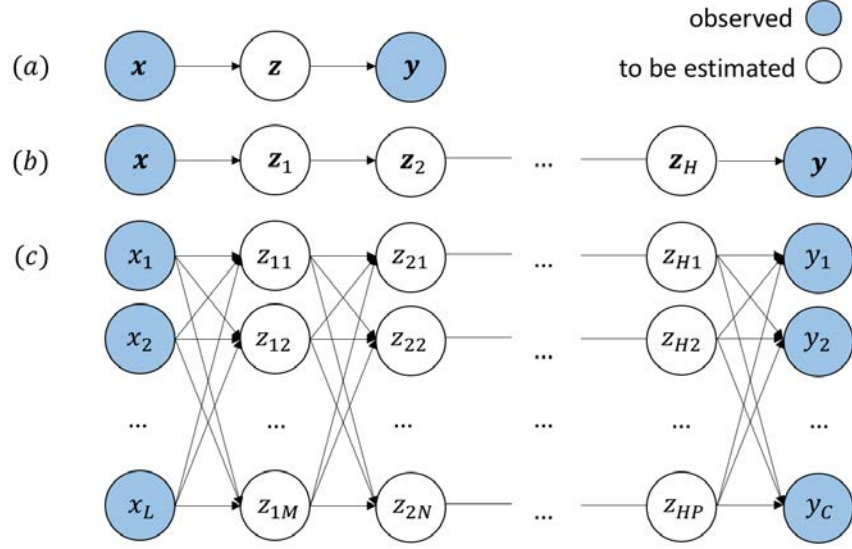


Figure 5.3: Graphical representations of feed-forward NNs with L input features and C output values. (a) Feed-forward network with a hidden layer, represented as a compact graph. Each node corresponds to a vector that contains the layer's activations. (b) Feed-forward network with H hidden layers, represented as a compact graph (c) Feed-forward network with H hidden layers, each of these contains a different number of units $\{M, N, \dots, P\}$, represented as an explicit graph.

The different variables with respect to Eq. (5.1) are \mathbf{y} , which is the vector with the C output values, and the matrix $\mathbf{W} \in \mathbb{R}^{L \times C}$ and the vector $\mathbf{b} \in \mathbb{R}^C$, which contain the weights and biases for the C units, respectively.

During the learning phase of the perceptron, its corresponding weights are iteratively updated based on samples from a training set, with the purpose of minimizing a chosen loss function. For each input sample $(\mathbf{x}_n, \mathbf{y}_n)$ and its predicted output $\hat{\mathbf{y}}_n$, the weight matrix \mathbf{W}_k at step k is updated as follows:

$$\mathbf{W}_k \leftarrow \mathbf{W}_{k-1} + \epsilon(\hat{\mathbf{y}}_n - \mathbf{y}_n)\mathbf{x}_n^T, \quad (5.3)$$

where ϵ denotes the learning rate, a fixed hyper-parameter to be determined that controls how quickly weights are updated.

The process is repeated until the error converges or a different criteria is met. Because all the training samples are processed independently by the algorithm at each step, this method is known as Stochastic Gradient Descent (SGD).

Feed-forward neural networks

Feed-forward neural networks or Multi-layer Perceptrons (MLPs) are directed acyclic graphs of stacked groups of computational units,

organized in layers, which have the ability of learning more complex functions than the ones achieved by perceptrons. A **MLP** consists of an *input layer* \mathbf{x} , which represent the L input features; one or several *hidden layers* \mathbf{z} , with the same or different number of units; and a final output layer \mathbf{y} , with as many units as values to be predicted.

Traditional **NNs** are comprised of Fully-Connected (**FC**) layers. While neurons between adjacent layers share fully pairwise connections, those within the same layer act in parallel and are not connected.

As can be seen in Figure 5.3, **MLPs** can be also represented by graphs with nodes and edges, in the same way as the previously described Bayesian **LTM**s (Chapter 3, Section 3.4). However, it should be noted that they constitute graphical representations of functions instead of distributions.

More formally, the output of a **MLP** with H hidden layers can be expressed as follows:

$$\mathbf{y} = f(\mathbf{W}^T \mathbf{z}_H + \mathbf{b}), \quad (5.4)$$

where \mathbf{y} is its output layer, represented as a vector with the C output values; f , \mathbf{W} and \mathbf{b} are its corresponding activation function, matrix of weights and vector of biases, respectively; and \mathbf{z}_H the vector with the output values of the units in the preceding hidden layer, which is recursively defined as:

$$\mathbf{z}_H = f_H(\mathbf{W}_H^T \mathbf{z}_{H-1} + \mathbf{b}_H). \quad (5.5)$$

Deep Neural Networks (**DNNs**) refer to **MLPs** with a great number of hidden layers and units, in the same way that “deep learning” is the field within **ML** that deals with this type of models. Convolutional Neural Networks (**CNNs**), described in Section 5.3.2, are **DNNs** whose hidden units have local receptive fields, particularly useful to solve computer vision tasks. When feed-forward networks are extended to have feedback connections, they are called Recurrent Neural Networks (**RNNs**), which will be introduced in Section 5.3.4. The following paragraphs cover the most important aspects regarding the definition and training of these architectures.

Architecture design

The complete architecture of a feed-forward network can be summarized by its depth, which is determined by its number of layers; the width of each layer, which is the number of units these layers have; and how these units are connected to each other. Most of the existing **MLPs** models follow a chain-based structure, where each layer is a function of its preceding layer.

As stated by the “no free lunch” theorem [33] introduced in Chapter 1, there is no **ML** algorithm better than any other so, for the

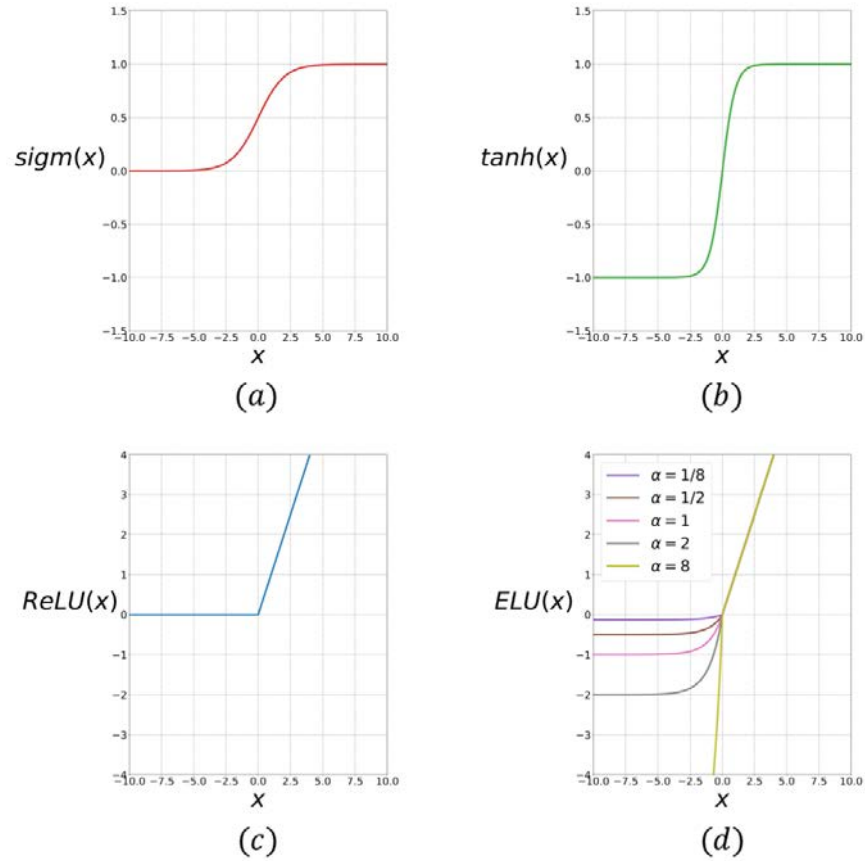


Figure 5.4: Graphical representation of the most commonly used activation functions. (a) Sigmoid. (b) Hyperbolic tangent. (c) Rectified Linear Unit (ReLU). (d) Exponential Linear Unit (ELU) represented for different values of α .

same reason, it is truly complex to define the optimal network structure for a specific application. According to the universal approximation theorem [179], a large DNN with enough capacity or hidden units is able to approximate any continuous function. Despite this, the optimization algorithm used to find the parameters corresponding to that function plays a crucial role in the learning phase and should be carefully validated. What is more, overfitting may occur in a model with higher capacity than needed, which would make it learn the noise in the data instead of the expected underlying relationships that lead to a good generalization. Regularization techniques can be helpful to prevent overfitting, as we will see in the next subsection.

On the other hand, the activation or transfer functions for both the hidden and the output units should be also taken into consideration. The most commonly used functions, which are represented in Figure 5.4, are described here below:

- *Sigmoid*: Defined as $\text{sigm}(x) = 1/(1 + e^{-x})$, the sigmoid non-linearity takes a real-valued scalar and compresses it to the range $[0, 1]$. Although it has been frequently used in the past, it presents two main drawbacks: first, sigmoids saturate at 0 or 1 when x is very negative or positive, respectively, which can hinder the common gradient-based learning; besides, sigmoid outputs are non-zero centered, which may introduce troublesome zig-zagging dynamics in the gradient updates for the weights [62].
- *Hyperbolic tangent*: It compresses a real-valued scalar to the range $[-1, 1]$, and can be seen as a zero-centered scaled version of the sigmoid function; consequently, it is often preferred in practice. In fact, $\tanh(x) = 2\text{sigm}(2x) - 1$.
- *Rectified Linear Unit (ReLU)*: The ReLU function, defined as $f(x) = \max(0, x)$, thresholds the activation at zero, notably accelerating the convergence of optimization methods with respect to the previously described sigmoidal functions, thanks to its simple, linear and non-saturating form. However, ReLU units are sensitive to large gradients. Although this issue can be often avoided using an appropriate learning rate, generalizations from ReLU activation such as *leaky ReLUs* and *maxout* have been tested in an attempt to fix it [43], as well as Exponential Linear Units (ELUs) [180]. As shown in Figure 5.4(d), ELUs have saturated negative values controlled by a hyperparameter α , which push the mean of the activations close to zero and solve the vanishing gradient problem. If the value of α is set too low, ELU and ReLU activations become similar.

ReLU non-linearities are almost always used in most hidden NN layers. Nevertheless, it should be noted that recurrent networks, probabilistic models and autoencoders draw on sigmoidal activation functions, despite their saturation drawbacks, because of their utility when having exploding gradients [43].

Regarding the activation function chosen for the output layer of a NN, it is completely dependent of the task to perform. While sigmoid allows to obtain class scores in a binary classification, *softmax* units, which represent a discrete probability distribution with C possible values, are more suitable for multi-class classifiers. In contrast, linear activations are used in regression problems to predict real values.

Gradient-based learning

The Goodfellow et al.'s recipe for ML [43] introduced in section 1.2.2 is entirely in line with the required specifications for gradient-based

training a DNN: select an optimizer, choose a cost or loss function and obtain a model.

1. *Select an optimizer:* First, the objective of the *optimization algorithm* is to find the weights \mathbf{w} associated with the DNN units that minimize the error between the expected GT values and the values predicted by the non-convex approximate function. For this purpose, iterative, gradient-based optimizers are the preferred option.

In contrast to traditional pure optimization methods, where the cost function is not related to the measure used to evaluate the performance of the system, gradient-based optimizers have in many contexts the advantage of directly minimizing a loss function $J(f(\mathbf{x}; \mathbf{w}), \mathbf{y})$ suited to the task to solve.

The basic algorithm for gradient-based optimization is Mini-Batch Gradient Descent (MBGD) [181]. At each iteration k , this stochastic method estimates the gradient of the approximate function as the average gradient on a small set or batch of M IID samples randomly chosen:

$$\hat{\mathbf{g}} = \nabla_{\mathbf{w}} \left(\frac{1}{M} \sum_{m=1}^M J(f(\mathbf{x}_m; \mathbf{w}), \mathbf{y}_m) \right). \quad (5.6)$$

Then, weights are updated as follows:

$$\mathbf{w}_k \leftarrow \mathbf{w}_{k-1} - \epsilon_k \hat{\mathbf{g}}. \quad (5.7)$$

Similarly to Eq. (5.3), a critical hyper-parameter to determine is the learning rate ϵ_k , which is now represented with the sub-index k because it is common to decay its value along iterations. This allows using a higher rate at the beginning, which prevents becoming stuck at a high cost, and progressively decreasing it in order to avoid significant oscillations in the learning curve.

Another approach to increase MBGD convergence speed is the momentum algorithm [182], which aids to keep the direction and speed at which the parameters are updated in subsequent iterations. Formally, at each iteration k , it computes an Exponential Moving Average (EMA) of the negative past gradients \mathbf{v} to update the weights:

$$\mathbf{v}_k = \alpha \mathbf{v}_{k-1} - \epsilon_k \nabla_{\mathbf{w}} \left(\frac{1}{M} \sum_{m=1}^M J(f(\mathbf{x}_m; \mathbf{w}), \mathbf{y}_m) \right), \quad (5.8)$$

$$\mathbf{w}_k \leftarrow \mathbf{w}_{k-1} + \mathbf{v}_k. \quad (5.9)$$

On the basis of [MBGD](#), two of the most outstanding techniques that incorporate both adaptive learning rates and momentum are Root Mean Squared Propagation ([RMSProp](#)) [[183](#)] and Adam [[184](#)].

Furthermore, it should be noted the importance of properly setting the initial network weights. Parameters from different units need to be initialized to different values, not to converge to the same configuration; this is commonly known as “breaking the symmetry effect”. As a convention, weights are often initialized to small values randomly chosen from a Gaussian or a uniform distribution, and biases are set to zero. While weights should be large enough to propagate information successfully, it is also desirable the use of small values to give a similar prior preference to all units.

Among other optimization strategies it is also remarkable the use of Batch Normalization ([BN](#)) layers [[185](#)], which constitute a method for adaptive reparametrization applicable to any input or hidden layer in a network. The distribution of [DNN](#) layer’s inputs changes during training. This effect, known as *internal covariate shift*, makes this stage substantially difficult, but can be avoided by normalizing mini-batch samples to have mean zero and standard deviation 1. Layer Normalization ([LN](#)) [[186](#)] is an alternative to [BN](#), which consists in normalizing layer activations. It may be useful when dealing with [RNNs](#) and small mini-batches.

2. *Choose a loss function*: Second, the error to minimize is defined by means of a differentiable *loss function*, which is often tailored to the task at hand for a better model fit. The most widely used loss functions correspond to the classical classification and regression tasks. The function associated with a classification task with multiple categories is the cross-entropy. For a problem with C classes and a set of N training samples, it is computed as:

$$H(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_{n=1}^N \sum_{c=1}^C y_n \log(\hat{y}_{nc}), \quad (5.10)$$

being y_n a binary label (1 if corresponds to the samples’s true class) and \hat{y}_{nc} the predicted probability for the class c . In a regression context, either Mean Squared Error ([MSE](#)) or Mean Absolute Error ([MAE](#)) are used, which are defined as follows:

$$MSE(y, \hat{y}) = \frac{1}{N} \sum_{n=1}^N (\hat{y}_n - y_n)^2 \quad (5.11)$$

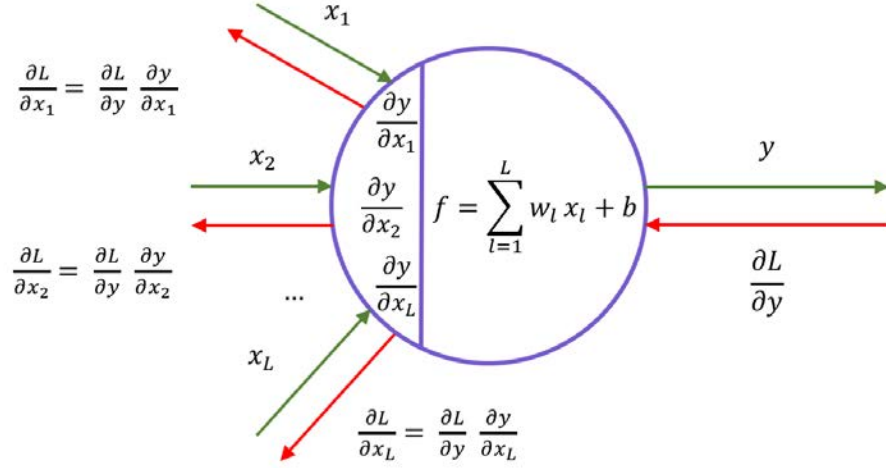


Figure 5.5: Computational graph of a NN layer with L inputs $\{x_1, x_2, \dots, x_L\}$ and one output y , where forward and back-propagation stages in gradient-based learning are represented with green and red arrows, respectively. Operations involved in both stages are indicated next to the arrows. At each iteration of gradient-based learning, the predicted output value is obtained in the forward pass, while during the backward pass the gradient of the approximate function is computed, in order to update the layer weights. Adapted from [62].

$$MAE(y, \hat{y}) = \frac{1}{N} \sum_{n=1}^N |\hat{y}_n - y_n| \quad (5.12)$$

When a DNN is applied to saliency or visual attention estimation, it is common to use a differentiable saliency metric, such as KL divergence, Normalized Scanpath Saliency (NSS) or Correlation Coefficient (CC) [101].

With the purpose of reducing the generalization error of a model, while keeping its training error low, it is prevalent the use of regularization techniques such as *weight decay* or those introduced in Chapter 7 of Goodfellow et al.'s book [43].

Early stopping is probably the most used regularization technique in deep learning. The method considers a validation set composed by some samples unseen during training and monitors, at each iteration, the error in this set as a reference to stop learning when the model starts overfitting. In this way, we can pre-specify a number of iterations after which the training phase finishes if the validation error has not decreased.

Another well-known and useful strategy for regularization is *Dropout* [187], which attempts to reduce overfitting in a network during training by randomly removing a certain percentage of hidden units at each iteration.

3. *Obtain a model:* The *back-propagation algorithm* [188] allows to compute the gradient needed to iteratively update the network weights and obtain a model. During the learning phase of a DNN, we can distinguish between two stages: forward propagation and back-propagation. In contrast, only forward propagation is performed in test.

Whereas in *forward propagation* a mini-batch of M input samples $(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_M)$ flows through a DNN, providing a mini-batch of outputs $(\hat{\mathbf{y}}_1, \hat{\mathbf{y}}_2, \dots, \hat{\mathbf{y}}_M)$, *back-propagation* takes the average loss between the expected GT values $(\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_M)$ and the predicted ones, which flows backwards through the network in order to obtain the gradients.

Figure 5.5 shows a simple computational graph for a NN layer with L inputs \mathbf{x} and one output y , the same represented in Figure 5.2.

- If we move from left to right, the output value y is computed in the forward pass, indicated by green arrows.
- In the backward pass, indicated by red arrows, from right to left, gradients are computed by means of the chain rule of calculus. First, we compute the gradient of the loss J with respect to the output value y : $\frac{\partial J}{\partial y}$. Then, the gradient value for each weight w_l is obtained as follows:

$$\frac{\partial J}{\partial w_l} = \frac{\partial J}{\partial y} \frac{\partial y}{\partial w_l}. \quad (5.13)$$

If the hidden unit participates in a DNN, J is the loss from the subsequent layer, which has to be back-propagated to every previous layer. In this case, we get the gradient value for each input x_l as:

$$\frac{\partial J}{\partial x_l} = \frac{\partial J}{\partial y} \frac{\partial y}{\partial x_l}. \quad (5.14)$$

Back-propagation can also be expressed in vector notation. Let us consider now a vector \mathbf{y} with C output values:

$$\nabla_{\mathbf{w}} J = \left(\frac{\partial \mathbf{y}}{\partial \mathbf{w}} \right)^T \nabla_{\mathbf{y}} J \quad (5.15)$$

$$\nabla_{\mathbf{x}} J = \left(\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \right)^T \nabla_{\mathbf{y}} J, \quad (5.16)$$

where $\frac{\partial \mathbf{y}}{\partial \mathbf{w}}$, $\frac{\partial \mathbf{y}}{\partial \mathbf{x}}$ are $C \times L$ Jacobian matrices of gradients. The same Jacobian-gradient product is performed for each layer in the DNN. It should be noted, nevertheless, that the method is usually applied to tensors of arbitrary

dimensionality rather than vectors, but it is conceptually the same. Given input, output and weight tensors X, Y, W , respectively, we can just treat them as vectors whose indexes have multiple coordinates (e.g. three coordinates for a 3D tensor). Thus, the chain rule is applied as follows:

$$\nabla_W J = \sum_j (\nabla_W Y_j) \frac{\partial J}{\partial Y_j} \quad (5.17)$$

$$\nabla_X J = \sum_j (\nabla_X Y_j) \frac{\partial J}{\partial Y_j}, \quad (5.18)$$

where j is an index variable to represent the complete tuple of coordinates in Y .

5.3.2 Convolutional Neural Networks

Convolutional Neural Networks (CNNs) arised inspired by the neurophysiological work of Hubel and Wiesel in 1962 [189], which discovered that neurons in the early visual system are sensitive to simple patterns of light, such as oriented edges or color blotches. As mentioned in the previous section, CNNs [64, 112] are DNNs composed of hidden units that have local receptive fields, similar to neurons in the primary visual cortex, and designed to take images or data with a grid-like topology as input.

Unlike traditional NNs, a CNN is composed of layers of neurons arranged in a 3D volume. Each layer transforms an input 3D volume into an output 3D volume. The characteristic operation of a CNN is the discrete convolution. Furthermore, CNNs may also involve pooling or down-sampling non-parametric operations.

The use of CNNs has several properties and advantages [43]:

1. *Locality of sparse interactions:* CNNs have sparse weights, which means that they are able to detect small, meaningful features by making use of convolution kernels smaller than the input, not larger than hundreds of pixels. This reduces the memory usage, while fewer operations are required to compute outputs, which improves the efficiency of the model.
2. *Parameter sharing:* In a convolutional layer, the number of parameters is dramatically reduced, because they are shared across the image; in other words, rather than learning a set of parameters for each spatial location, the same convolutional kernel is applied on all of them.
3. *Invariance to translation:* CNNs are translational invariant, so they are able to identify patterns independently of their location in

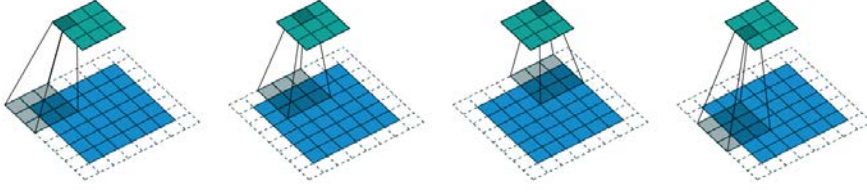


Figure 5.6: Application of a $k = 3 \times 3$ convolutional kernel over a 6×6 input padded with a 1×1 border of zeros, using a stride of $s = 2$. Figure taken from [191].

the input image. Conversely, they are not invariant to rotation or scale changes, and require of other techniques to handle these transformations, such as data augmentation [112, 190].

Architecture design

Common CNNs are built on a sequence of A blocks of B Convolutional (CONV) layers with ReLU activations, sometimes followed by a Pooling (POOL) layer, ending up in a stack of C FC layers, being the last FC layer the one that holds the output predicted values [62]:

$$\begin{aligned} \text{INPUT} &\rightarrow A \times [B \times [\text{CONV} \rightarrow \text{ReLU}] \rightarrow \text{POOL?}] \\ &\rightarrow C \times [\text{FC} \rightarrow \text{ReLU}] \rightarrow \text{FC} \end{aligned}$$

Related to the typical layers of CNNs, CONV layers first involve a set of filters with learnable weights. Let width, height and depth denote the arbitrary dimensions of an input 3D volume. During the forward pass, each filter is spatially slid along the width and the height of the channels or spatial maps stacked on the depth dimension, in order to compute the dot product at each spatial location. The connections between the input and each filter are thus local in space, but full along the depth. The *filter* or *kernel size* k is the receptive field of the neuron, while the size of the output 3D volume depends on three hyper-parameters:

- The number of filters u used corresponds to the *depth* of the resulting output volume.
- The *stride* s with which the filters are slid or down-sampling factor determines the spatial size of the output. Bigger strides result in smaller output volumes.
- *Zero-padding* the input volume is another way of varying the output spatial size, increasing it spatially with a border of one or several pixels.

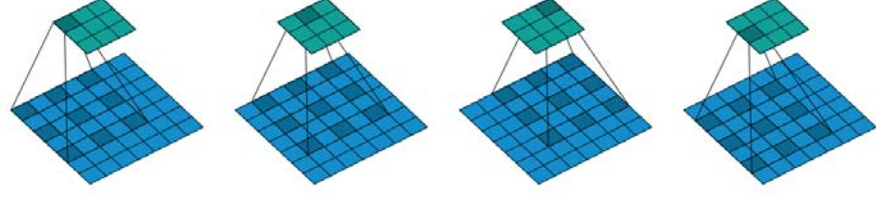


Figure 5.7: Application of a $k = 3 \times 3$ dilated convolutional kernel over a 7×7 input, using a dilation factor of $d = 2$ (1 space between kernel elements). Figure taken from [191].

Given a filter $K_{i,l}$ of size $k = m \times n$, which gives the connection between channel i in the output volume Y and channel l in the input volume X , and is applied centered at the input location j, k , being j and k the row and column positions; assuming that both input and output have the same spatial dimensions, the value that results from a CONV layer at the output location j, k can be expressed more formally in the language of tensors:

$$Y_{i,j,k} = \sum_{l,m,n} X_{l,(j-1) \times s_W + m, (k-1) \times s_H + n} K_{i,l,m,n}, \quad (5.19)$$

where s_W, s_H denote the strides of the convolution along width and height, respectively. Figure 5.6 shows an example of application of a convolutional kernel.

Due to their regular use in CNNs for the extraction of feature maps, we shall introduce a particular type of convolutions, known as *dilated* or “*atrous convolutions*” [192]. They differ from original convolutions in the insertion of spaces between the kernel elements depending on a dilation rate. As can be seen in the example of application in Figure 5.7, a dilation rate d corresponds to $d - 1$ spaces inserted between elements. Dilated convolutions allow to increase the receptive field of the output without increasing the size of the kernel.

Other variants of the basic convolution function, such as transposed convolutional layers or locally connected layers, are introduced in Chapter 9 from [43], and also in the excellent guide to convolution arithmetic by Dumoulin et al. [191].

Secondly, POOL layers are useful to make representations almost invariant to small translations of the input. Moreover, by pooling over outputs from different convolutions, the network learns to become invariant to some transformations. POOL layers perform a function that does not have parameters. The most commonly used POOL layer is *max pooling* (MAX POOL), which involves a maximum operation, but other functions such as the average (AVG POOL) or the Euclidean norm are also considered. POOL layers often entail down-sampling operations ($s > 1$) along width and height, which reduce the dimensions of the feature maps.

Case studies

One of the first successful applications of convolutional networks was *LeNet* [194]. Developed by Yann LeCun in 1990, it is composed of five layers and was used to read digits from zip codes. Throughout the brief history of the application of CNNs in computer vision, three additional case studies should be highlighted, due to their successful achievements in image recognition: *AlexNet* [3], *VGGNet* [113] and *ResNet* [193]. Deeper (8 layers), but with a similar structure to *LeNet*, *AlexNet* is one of the first works that promoted the use of CNNs in computer vision, winner of the ImageNet Large Scale Visual Recognition Competition (ILSVRC) [166] in 2012. The other two aforementioned configurations are used for visual attention feature extraction in this thesis, so they are described hereunder for the sake of completeness.

- *VGGNets* [113] were the state-of-the-art in 2014. By increasing the depth of prior architectures up to 19 layers and using CONV kernels with small receptive fields ($k = 1 \times 1$, $k = 3 \times 3$) and ReLU non-linearities, they outperformed *AlexNet* and other models proposed that year, reaffirming again that networks depth is a critical component for a better performance. The diagram in Figure 5.8(a) shows one of the most used VGG configurations, *VGG-16*, which includes 16 weight layers, 13 $k = 3 \times 3$ CONV layers with stride $s = 1$ and 3 FC layers. As can be seen, some of the CONV layers are followed by $k = 2 \times 2$ MAX POOL layers, with stride $s = 2$.
- *ResNets* [193], winner of ILSVRC [166] in 2015, arose with the aim of reducing the degradation problem caused by the saturation of performance when training deeper NNs. They are built on a series of blocks composed of few layers, among which residual mappings are performed. These mappings consist of adding the input of each block to its output, by making use of skip connections, as can be noticed in the commonly-used 50-layer architecture in Figure 5.8(b). The configuration contains a $k = 7 \times 7$ CONV layer, followed by a MAX POOL layer, both with $s = 2$, and then a stack of fully CONV residual blocks with increasing number of units. The first CONV layer of each block reduces the dimensionality of its input by using a stride of $s = 2$.

At the end of either VGG or *ResNet* networks, the final FC softmax layer contains 1000 channels, associated with the image classes considered in ILSVRC [166]. It should be noted that last FC layers are removed in tasks such as image segmentation or visual attention estimation, which ultimately provide a pixel-wise score map instead of an image-wise score. Moreover, because these applications require

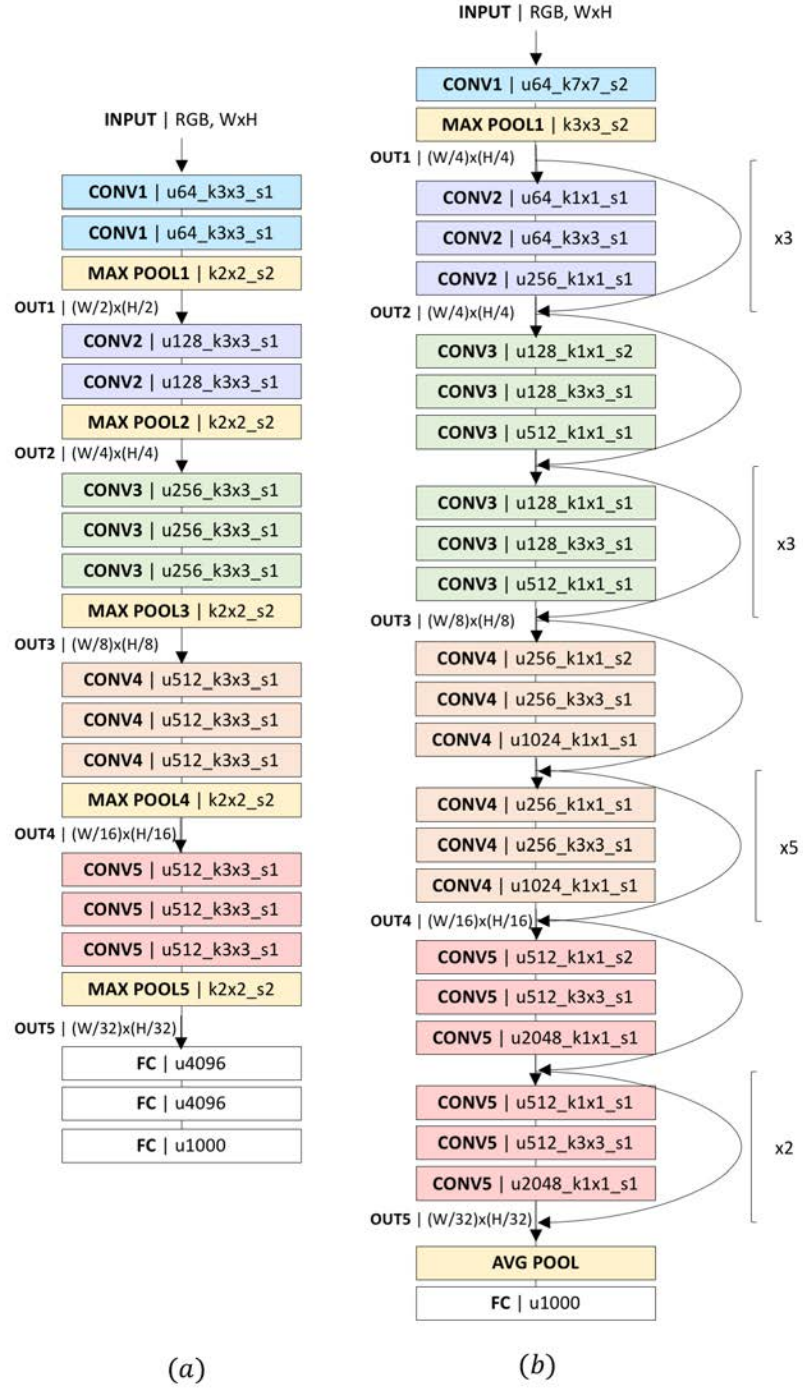


Figure 5.8: Architecture diagrams of the (a) VGG-16 [113] and (b) ResNet-50 [193] networks for image recognition. Layers are defined by their number of units u , kernel size k and stride s with which filters are slid. Given an INPUT image of dimension $W \times H$, output (OUT) sizes are indicated at the end of each block. Last FC softmax layers in both models have 1000 units, which correspond to the image classes defined in ILSVRC [166].

more accurate spatial predictions, some **CONV** layers are often substituted by dilated convolutions, in order to preserve a larger spatial resolution across layers, as we will see in the architectures used for salient feature extraction in Section 5.4.

Transfer learning in CNNs

Transfer learning [50, 195] is a **ML** method which consists in reusing the knowledge acquired while solving a particular task with the aim of addressing a second different but related task.

In order to train a **CNN** to operate in a particular scenario, it is necessary to annotate a large image database, which might become a highly arduous task. That is the reason why an entire **CNN** is seldom trained from scratch [62]. Instead, it is quite common to pretrain the network on a big dataset, such as the well-known ImageNet [196]. This pretrained model is used then as a fixed feature extractor, or its (or some of its) layers are *fine-tuned* on the new smaller database. Smaller learning rates are often used for fine-tuned **CNNs**, assuming that already existing weights are sufficiently good, so it is better not to significantly change them.

5.3.3 Encoder-Decoder Networks

The impressive ability of **DNNs** to capture good hierarchical representations of data has also served to update the way of computing invariant features [17]. These features were traditionally achieved by means of unsupervised methods for dimensionality reduction or clustering (e.g. **PCA**, K-means [42]), as well as through hand-crafted histograms, such as Scale-Invariant Feature Transform (**SIFT**) [197] or Histogram of Oriented Gradients (**HOG**) [198] descriptors.

Encoder-Decoder Networks (**EDNs**) are a special case of feed-forward **NNs**, explicitly designed to learn efficient feature representations. As shown in the graphical representation of Figure 5.9(a), these architectures are composed of two consecutive networks:

1. An *encoder* network, which takes an input \mathbf{x} and, after one or several hidden layers, represents it as a feature code or latent representation $\mathbf{z}_{E_H} = f_{E_H-1}(\mathbf{z}_{E_H-1})$, being E_H its number of hidden layers.
2. A *decoder* network with D_H hidden layers, which reconstructs the feature code \mathbf{z}_{E_H} , in order to produce an output \mathbf{y} tailored to the task to solve. For that purpose, the reconstruction error with respect to a **GT** is measured (e.g. segmentation mask in an image segmentation system, see Figure 5.9(b)).

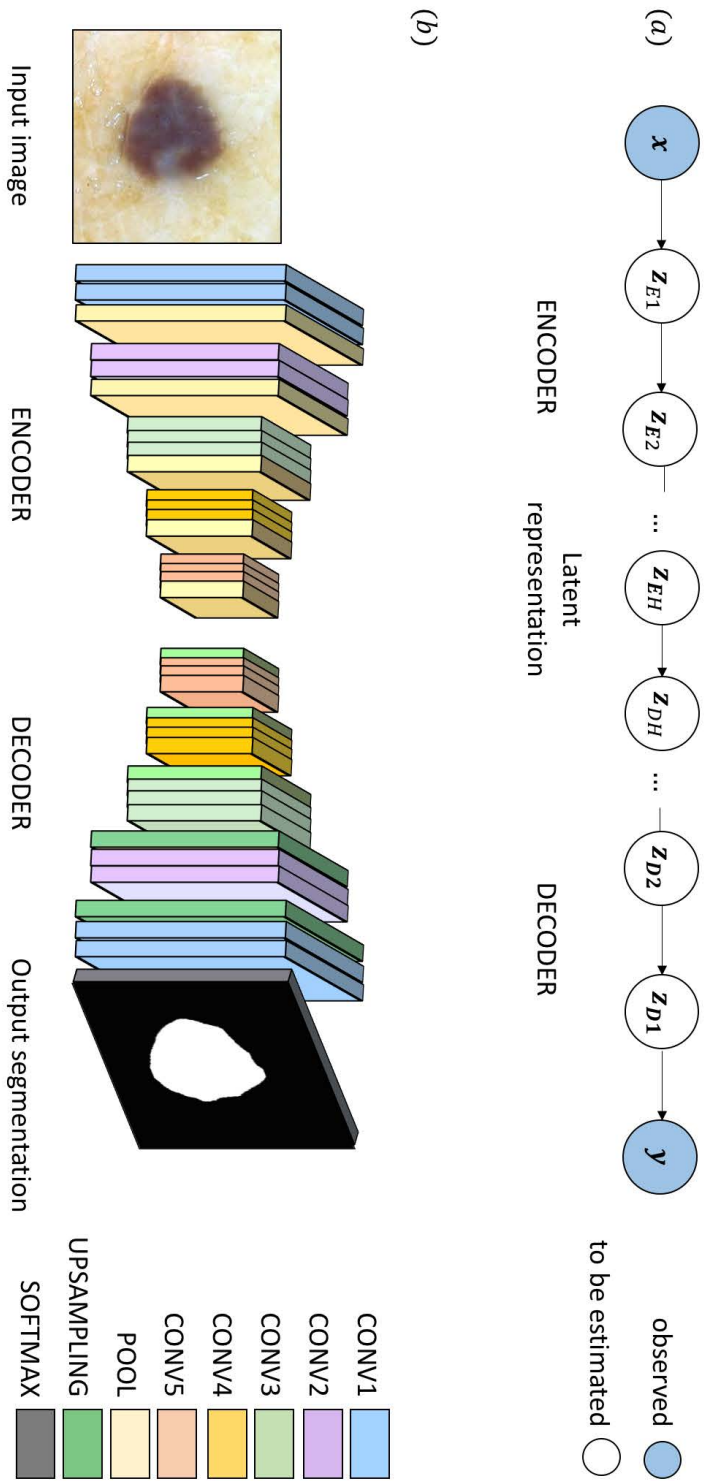


Figure 5.9: (a) Graphical representation of an EDN composed by an encoder with E_H hidden layers, which takes an input x and generates a latent representation z_{E_H} , and a decoder with D_H hidden layers, which produces an output y from the representation, tailored to the task to solve. Each node corresponds to a vector that contains the layer's activations. The EDN has $E_H + D_H$ layers in total. (b) Example diagram of a convolutional encoder-decoder architecture for skin lesion segmentation. Both encoder and decoder networks consist of the 13 convolutional layers in the VGG-16 [113] network represented in Figure 5.8(b).

Following this notation, EDNs have $E_H + D_H$ layers in total. They can be trained by using the same algorithms described for feed-forward networks (see Section 5.3.1). In the particular case when inputs and outputs are equal, those networks are called *autoencoders* [43]; and its encoder and decoder often have the same number of hidden layers ($E_H = D_H$).

We can distinguish between two types of EDNs. The EDN is *undercomplete* when the hidden code dimension is lower than the input dimension, and attempts to extract the most useful or salient properties of the input data. In contrast, the latent representation has a dimension higher than the input in *overcomplete* EDNs.

If the decoder is linear and MSE is considered as loss function, an undercomplete autoencoder learns a similar subspace to the one obtained by PCA. However, a EDN with nonlinear encoder and decoder functions is able to learn more powerful representations than PCA, which is restricted to linear transformations of the data. Besides, if the capacity of the network becomes too high, it might end up copying the training data instead of gathering the most prominent underlying information in the data distribution. In that case, regularization techniques can be helpful to avoid overfitting in overcomplete EDN, providing more sparse and shift-invariant representations [43].

EDNs may involve either FC DNNs, CNNs, RNNs or a mixture of all three. They have been recently applied to solve tasks such as LSTM-based sequence to sequence translation [55], image segmentation [199, 200], or image captioning [201], combining a CNN encoder with a LSTM decoder. Image restoration has also been addressed by using denoising autoencoders with symmetric skip connections [202]. The following is a brief explanation of convolutional EDNs, which we will use in our system described in Section 5.5 for spatio-temporal visual attention estimation.

Convolutional Encoder-Decoder Networks

An image comprises a set of features located in different regions. Given an input image, deep Convolutional Encoder Decoders (CEDs) compute an invariant feature vector that encodes *what* features are in the image (their presence or absence), via a set of transformation parameters, which entail *where* these features are found within the image (their location) [17].

Figure 5.9(b) shows an example diagram of a CED network for object segmentation. On the one hand, the encoder network consists of a series of CONV and POOL sub-sampling layers, just like CNNs. If we look at the encoder of the example, we will notice that it is composed by the 13 CONV layers in the VGG-16 [113] represented in Figure 5.8(b). As explained above in Section 5.3.2, POOL layers help

to manage translation invariance to small spatial shifts in the image and, at the same time, reduce its dimensionality. Therefore, they provide compacted representations which can be useful to efficiently carry out classification or regression problems. On the other hand, the decoder network is constituted by several **CONV** and upsampling layers, which reconstruct the feature vector and converge into a final **CONV** output layer. A pixel-wise binary classification task is performed at the output in order to achieve the desired skin lesion segmentation mask. The decoder of the example has the same number of layers than the encoder, establishing a correspondence between **POOL** and upsampling layers.

5.3.4 Recurrent Neural Networks

Recurrent Neural Networks (**RNNs**) [188, 203] are **NNs** with feedback connections specialized for the processing of sequential data. Widely-used for language modeling in speech recognition [204] since 2010, when they outperformed the standard n-gram models [205], **RNNs** have more recently been used in computer vision applications such as image captioning [206] or video action recognition [207].

RNNs operate on input sequences composed of vectors \mathbf{x}_t , which are denoted by the time step index t , and have a repetitive structure: drawing on cycles, they are able to capture the influence of a present variable $\mathbf{x}^{(t)}$ at a time t on its future value $\mathbf{x}^{(t+1)}$, by learning a set of parameters \mathbf{w} that are shared across the network states. States correspond to the hidden units $\mathbf{z}^{(t)}$ of the network, whose relationship is typically expressed by the following equation:

$$\mathbf{z}^{(t)} = f(\mathbf{z}^{(t-1)}, \mathbf{x}^{(t)}; \mathbf{w}). \quad (5.20)$$

RNNs have two primary advantages [43]. First, they learn a model based on transitions between states, which always has the same input size. Moreover, the model can be independent of the sequence length by making use of the same transition function f with the same parameters \mathbf{w} at every time step, instead of learning a separate model for each one. Given a long sequence, **RNNs** usually work on mini-batches of shorter length τ , similar to how feed-forward networks deal with samples.

It should be noted that we still have not mentioned the output layers that use the information from states to make predictions. Depending on the architecture design, **RNNs** can generate either an output at each time step or a single output for a whole sequence. Besides, they can have recurrent connections between units or, in contrast, from the output at a time step to the hidden units at the subsequent time step.

Figure 5.10 shows both the compact and the time-unfolded graph of a representative **RNN** example. The network maps an input

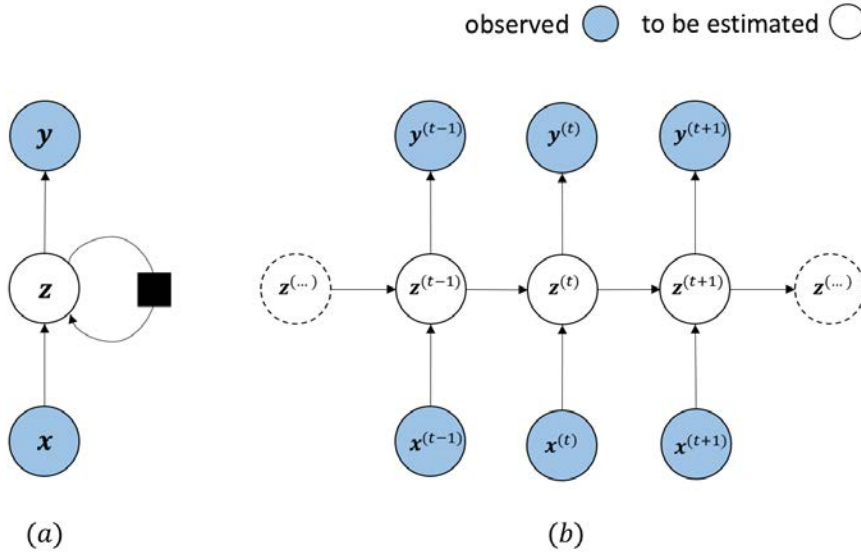


Figure 5.10: Graphical representations of a [RNN](#) that maps an input sequence of \mathbf{x} values to an output sequence \mathbf{y} . (a) [RNN](#) represented as a compact graph. Each node corresponds to a vector that contains the layer's activations. The black square means a delay of a single time step, from the state $\mathbf{z}^{(t)}$ to $\mathbf{z}^{(t+1)}$. (b) [RNN](#) represented as a time-unfolded graph. Each node corresponds to a particular time instance.

sequence \mathbf{x} of values or vectors to an output sequence of \mathbf{y} . Given a initial state $\mathbf{z}^{(0)}$ and a mini-batch of length τ , the following update equations are applied from $t = 1$ to $t = \tau$:

$$\mathbf{z}^{(t)} = g(\mathbf{W}\mathbf{z}^{(t-1)} + \mathbf{U}\mathbf{x}^{(t)} + \mathbf{b}) \quad (5.21)$$

$$\mathbf{y}^{(t)} = h(\mathbf{V}\mathbf{z}^{(t)} + \mathbf{c}) \quad (5.22)$$

where \mathbf{U} , \mathbf{V} and \mathbf{W} are weight matrices that model input-to-hidden, hidden-to-output and hidden-to-hidden connections, respectively; \mathbf{b} and \mathbf{c} are bias vectors; g is the activation functions for the hidden units; and h the activation function at the output.

Gradient-based learning in [RNNs](#)

The unfolded graph in Figure 5.10(b) illustrates how information $\mathbf{x}^{(t)}$ flow in the network during the forward pass from left to right, computing the output values $\mathbf{y}^{(t)}$. The total loss J for the sequence is given by the sum of the loss $J^{(t)}$ at each time step. Then, in the backward pass, gradients are computed for each time step.

The algorithm to compute gradients in [RNNs](#) is called Back-propagation Through Time (BPTT) [43, Section 10.2.2], and applies back-propagation to the unfolded graph. For each node, the gradient is computed recursively, based on the gradients of the following nodes in the graph. In vector notation, given $\frac{\partial J}{\partial J^{(t)}} = 1$, the

gradient $(\nabla_{\mathbf{y}^{(t)}} J)_i$ for each value $y_i^{(t)}$ on the output sequence at time step t can be written as:

$$(\nabla_{\mathbf{y}^{(t)}} J)_i = \frac{\partial J}{\partial J^{(t)}} \frac{\partial J^{(t)}}{\partial y_i^{(t)}} = \frac{\partial J^{(t)}}{\partial y_i^{(t)}}, \quad (5.23)$$

Then, the gradient on the hidden state $\nabla_{\mathbf{z}^{(t)}} J$ is as follows:

$$\nabla_{\mathbf{z}^{(t)}} J = \left(\frac{\partial \mathbf{z}^{(t+1)}}{\partial \mathbf{z}^{(t)}} \right)^T (\nabla_{\mathbf{z}^{(t+1)}} J) + \left(\frac{\partial \mathbf{y}^{(t)}}{\partial \mathbf{z}^{(t)}} \right)^T (\nabla_{\mathbf{y}^{(t+1)}} J). \quad (5.24)$$

Once these gradients are computed for the computational graph associated with a mini-batch of length τ taken from an entire longer sequence, the gradients with respect to the weights matrices \mathbf{U} , \mathbf{V} , and \mathbf{W} , and bias vectors \mathbf{b} and \mathbf{c} can be obtained, so that they can be subsequently updated.

Long Short-Term Memory Units

The main issue when learning long-term dependencies with a [RNN](#) is that gradients propagated over many states are very small or large in magnitude, either vanishing or exploding, dramatically hampering the optimization process. In order to reduce these effects, the scale of the initial weights has to be chosen carefully.

Although the problem of learning long-term dependencies continues being one of the main challenges in deep learning, several techniques have been proposed in order to alleviate it (e.g. *gradient clipping* [208] before the weights update rule produced very large gradient magnitudes). Here we discuss a special type of sequence models, known as *gated RNNs*, whose objective is to define paths through time with non-vanishing and non-exploding derivatives. Among gated [RNNs](#), it is worth mentioning Long Short-Term Memory ([LSTM](#)) units [18], used in our system for modeling attention in the temporal dimension (see Section 5.5) and described below, and Gated Recurrent Units ([GRUs](#)) [209].

Figure 5.11 shows the block diagram of a common [LSTM](#) unit. A [LSTM](#) model is composed by cells characterized by a state $\mathbf{C}^{(t)}$, represented in the diagram by the top horizontal line. Cells involve an internal recurrence or self-loop which complements the outer recurrence of traditional [RNNs](#). This internal recurrence serves to control the flow of information by adding or removing information to/from the cell state, and is defined by means of a system of gating units:

- First, the *forget gate* unit $\mathbf{f}^{(t)}$ decides which information to throw away from the cell state, based on the input $\mathbf{x}^{(t)}$ and the previous hidden units $\mathbf{z}^{(t-1)}$, and using a sigmoid activation function:

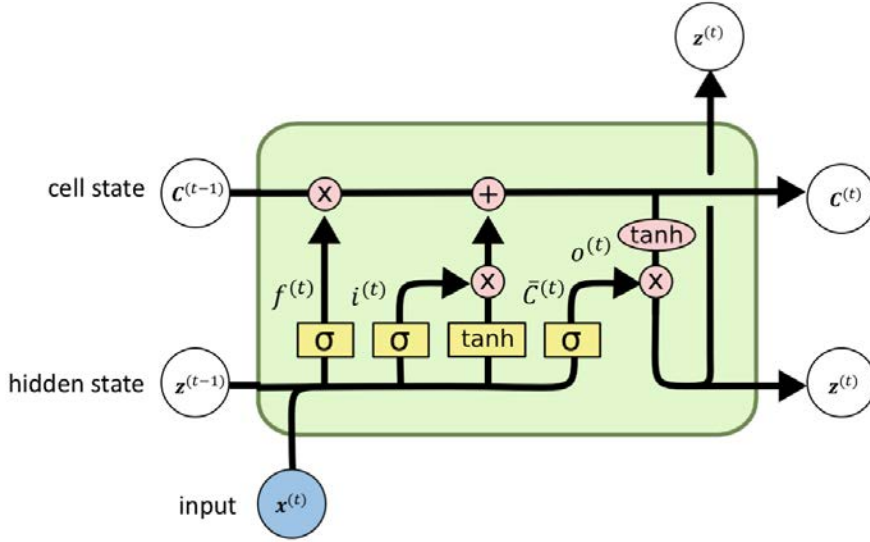


Figure 5.11: Diagram of a **LSTM** unit. Pink circles indicate point-wise operations, while yellow boxes represent, from left to right, forget $f^{(t)}$, input $i^{(t)}$, cell state candidate values $\bar{c}^{(t)}$ and output $o^{(t)}$ gates, which control the flow of information by adding or removing information to/from the cell state $C^{(t)}$, based on the input $\mathbf{x}^{(t)}$, and determine the hidden state $\mathbf{z}^{(t)}$. Adapted from [210].

$$\mathbf{f}^{(t)} = \text{sigm} \left(\mathbf{U}_f \mathbf{x}^{(t)} + \mathbf{W}_f \mathbf{z}^{(t-1)} + \mathbf{b}_f \right) \quad (5.25)$$

where \mathbf{U}_f , \mathbf{W}_f and \mathbf{b}_f are input weights, recurrent weights and biases for the forget gate.

- Then, in order to decide what new information is going to be stored, there are two gates: the *input gate*, which determines the values to be updated:

$$\mathbf{i}^{(t)} = \text{sigm} \left(\mathbf{U}_i \mathbf{x}^{(t)} + \mathbf{W}_i \mathbf{z}^{(t-1)} + \mathbf{b}_i \right) \quad (5.26)$$

and a *tanh* activation, which creates a vector of new candidate values to be added to the cell state:

$$\bar{\mathbf{c}}^{(t)} = \text{tanh} \left(\mathbf{U}_c \mathbf{x}^{(t)} + \mathbf{W}_c \mathbf{z}^{(t-1)} + \mathbf{b}_c \right) \quad (5.27)$$

where \mathbf{U}_i , \mathbf{U}_c are input weights; \mathbf{W}_i , \mathbf{W}_c recurrent weights; and \mathbf{b}_i , \mathbf{b}_c are biases. Based on $\mathbf{f}^{(t)}$, $\mathbf{i}^{(t)}$ and $\bar{\mathbf{c}}^{(t)}$, the *cell state* $\mathbf{C}^{(t)}$ is updated as follows:

$$\mathbf{C}^{(t)} = \mathbf{f}^{(t)} \mathbf{C}^{(t-1)} + \mathbf{i}^{(t)} \bar{\mathbf{c}}^{(t)} \quad (5.28)$$

- Finally, the output value, at the *output gate* $\mathbf{o}^{(t)}$, and the current hidden state $\mathbf{z}^{(t)}$ are computed:

$$\mathbf{o}^{(t)} = \text{sigm}(\mathbf{U}_o \mathbf{x}^{(t)} + \mathbf{W}_o \mathbf{z}^{(t-1)} + \mathbf{b}_o) \quad (5.29)$$

where \mathbf{U}_o , \mathbf{W}_o and \mathbf{b}_o are input weights, recurrent weights and biases again; and

$$\mathbf{z}^{(t)} = \mathbf{o}^{(t)} \tanh(\mathbf{C}^{(t)}). \quad (5.30)$$

Firstly introduced for precipitation nowcasting in [115], convolutional **LSTMs** constitute an extension of **LSTMs** to operate on images and 2D feature maps, and are also used as part of the spatio-temporal visual attention **EDNs** presented in Section 5.5.4. **CONV-LSTMs** simply incorporate convolutional operators either to the input-to-hidden or the hidden-to-hidden transitions, being expressed as follows:

$$\mathbf{F}^{(t)} = \text{sigm}(\mathbf{U}_f * \mathbf{X}^{(t)} + \mathbf{W}_f * \mathbf{Z}^{(t-1)} + \mathbf{b}_f) \quad (5.31)$$

$$\mathbf{I}^{(t)} = \text{sigm}(\mathbf{U}_i * \mathbf{X}^{(t)} + \mathbf{W}_i * \mathbf{Z}^{(t-1)} + \mathbf{b}_i) \quad (5.32)$$

$$\overline{\mathbf{C}}^{(t)} = \tanh(\mathbf{U}_C * \mathbf{X}^{(t)} + \mathbf{W}_C * \mathbf{Z}^{(t-1)} + \mathbf{b}_C) \quad (5.33)$$

$$\mathbf{C}^{(t)} = \mathbf{F}^{(t)} \circ \mathbf{C}^{(t-1)} + \mathbf{I}^{(t)} \circ \overline{\mathbf{C}}^{(t)} \quad (5.34)$$

$$\mathbf{O}^{(t)} = \text{sigm}(\mathbf{U}_o * \mathbf{X}^{(t)} + \mathbf{W}_o * \mathbf{Z}^{(t-1)} + \mathbf{b}_o) \quad (5.35)$$

$$\mathbf{Z}^{(t)} = \mathbf{O}^{(t)} \circ \tanh(\mathbf{C}^{(t)}), \quad (5.36)$$

where $*$ and \circ denote the convolution operator and the Hadamard product, respectively.

5.4 FEATURE LEARNING FOR VISUAL ATTENTION GUIDANCE

In this section, we describe three feature extraction **CNNs** for visual attention guidance. Unlike in our first system, presented in Chapter 3, our main goal now is not to understand how visual attention works in diverse contexts, but to model attention in the temporal domain using spatio-temporal **VAMs** as an input.

For that purpose, we will first make use of three fundamental visual attention feature maps to estimate spatio-temporal visual attention: RGB-based spatial, optical flow-based motion, and objectness-based maps.

While the networks used to obtain spatial and motion feature maps directly learn from **GT** fixation maps, the model for objectness minimizes a loss function depending on **GT** annotated object masks. These feature maps have been obtained by adapting the well-known VGG-16 [113] and ResNet-50 [193] networks, successfully applied first to image recognition, and more recently to saliency estimation, mainly in still images [108].

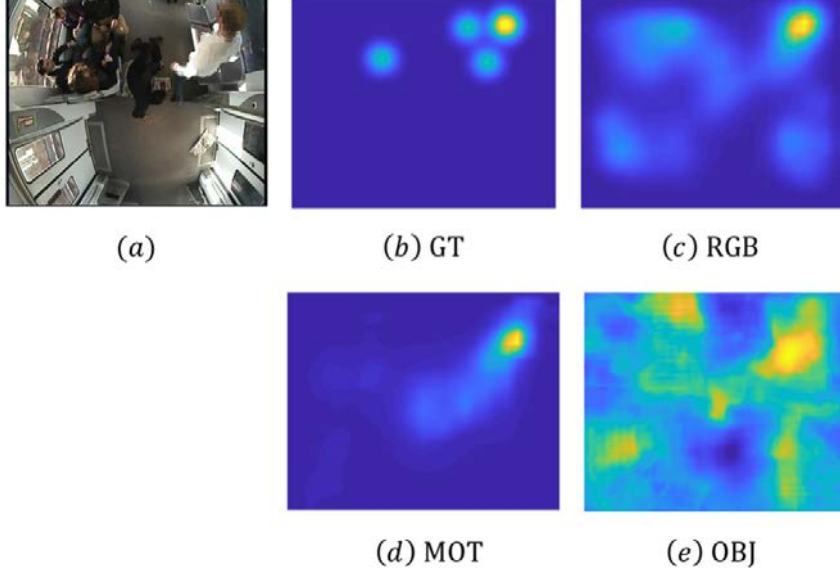


Figure 5.12: RGB-based, motion and objectness feature maps computed for an example frame taken from BOSS [19] database. (a) Original frame. (b) GT fixation map. (c) RGB-based feature map. (d) Motion feature map. (e) Objectness feature map.

5.4.1 RGB-based spatial network

A spatial feature map f_{RGB} is first computed by using a modified version of the ResNet-50 [193] introduced in Section 5.3.2, and a 224×224 RGB image as input. Similarly to the dilated residual convolutional block presented by Marcella Cornia et al. in [108] as part of the SAM model, we first remove the FC layers at the end of the network and then introduce dilated convolutions in either CONV4 or CONV5 blocks with dilation rate $d = 2$ and $d = 4$, respectively. Then, the output tensor of the CONV5 block, with 2048 units, is fed into an additional CONV block composed of two $k = 3 \times 3$ and $k = 1 \times 1$ layers, both with 128 units. Finally, we place a final layer with a unique unit and linear activation, which corresponds to the output spatial map, with dimension 26×26 . Figure 5.12(c) shows an example of a RGB-based feature map associated with a frame taken from BOSS [19] database.

5.4.2 Optical flow-based motion network

Motion feature maps can be achieved by means of CNNs that take as input a pair of subsequent frames or, in contrast, the previously estimated optical flow from these, similarly to our traditional feature extraction process in Section 3.3.2 or to the networks proposed by Cagdas Bak et al. in [106]. The network used for our experiments

has been inspired by the latter work but, unlike this approach and for the sake of continuity, it builds over the modified ResNet-50 [193] architecture used for RGB-based maps to extract a motion feature map f_{MOT} . The model now receives as input an 224×224 image with three channels, corresponding to the horizontal and vertical optical flow components and its associated motion magnitude map. All these channels are first re-scaled to the range $[0, 255]$. Figure 5.12(d) includes an example of a motion feature map for a frame taken from BOSS [19] database.

5.4.3 Objectness-based network

Computational models for salient object detection have also demonstrated the importance of objects regardless their semantic category [92, 109]. Any object in motion is often noticeable and, if more than one object appears on a frame, attention will choose one depending on the conspicuity of their associated low-level properties.

As we anticipated in Section 3.5.2, we make again use of the Deep Contrast Network for general object detection introduced by Li et al. in [2]. Unlike in our previous system, we now only consider the final fused objectness feature map, which we denote as f_{OBJ} and is the output of the DCL model, prior to the CRF-based refinement applied then. In order to achieve this map, the model makes use of two streams:

- The first stream of the architecture is based on the VGG-16 [113] network presented in Section 5.3.2. In the same way as most deep networks for saliency estimation, it includes two additional CONV layers after the final POOL₅ layer, which in fact replace the original FC ones. In addition, it substitutes some of the traditional CONV layers by dilated ones. More precisely, the three CONV₅ layers and the two extra CONV layers incorporated after the final pooling layer have dilation rates $d = 2$ and $d = 4$, respectively. Moreover, with the purpose of implementing a multi-scale version of VGG-16 [113], the authors connect three more CONV blocks to each of the first four MAX POOL layers.
- The second stream is a superpixel or segment-wise pooling network, which models both visual contrast between regions and discontinuities in object boundaries.

The final objectness-based feature map is the linear combination of the four outputs from the multi-scale CONV blocks, the final output map of the first stream and the segment-wise map from the second stream. An example of an objectness-based feature map corresponding to a frame taken from BOSS [19] database is shown in Figure 5.12(e).

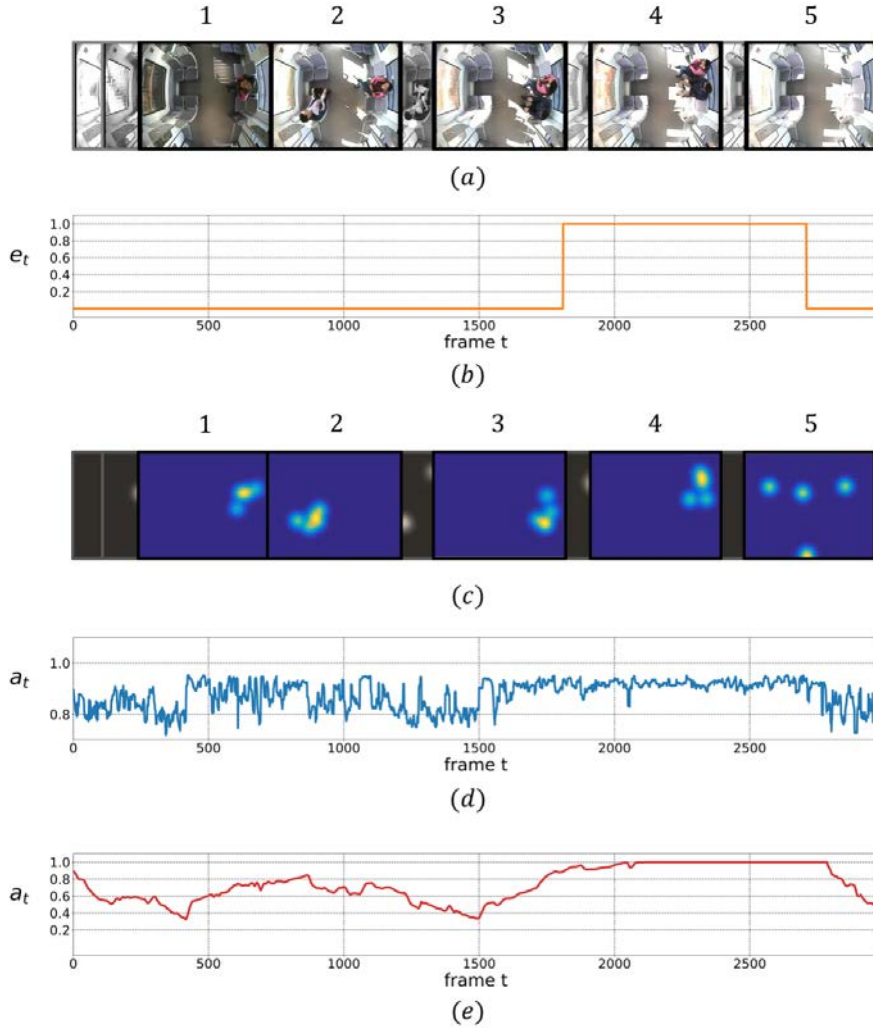


Figure 5.13: Visual attention in the temporal domain modeled in a video-surveillance sequence taken from BOSS [19] database. The sequence shows a woman harassment scene on a train. (a) Sequence scenes. 1. Woman sits on the train. 2. Man gets on the train. 3. Man approaches the woman. 4. Man bothering the woman. 5. Woman and man leave the train. (b) Anomaly detection signal e_t , which is set to 1 when an anomaly happens. (c) Fixation maps. Fixations are gathered from 5 users watching the video. (d) Raw temporal attention response a_t . (e) Filtered temporal attention response a_t . According to the fixations-based signals provided, attention achieves its maximum value just before and at the moment of the harassment (3,4). However, it should be noted that other events also highly attract the attention of observers, such as the moment when the man appears on the scene (2). Therefore, temporal attention response a_t should be considered an early filtering mechanism along time, that allows selecting time segments of special importance, which often match with anomalies.

5.5 SPATIO-TEMPORAL TO TEMPORAL VISUAL ATTENTION NETWORK

In this section, we describe in detail our system for visual attention estimation in the temporal domain, which we have called Spatio-Temporal to Temporal visual ATtention NETwork ([ST-T-ATTEN](#)).

5.5.1 Fundamental hypothesis of the model

The [ST-T-ATTEN](#) proposed for modeling the temporal dimension of visual attention is motivated by the following two assumptions:

1. A measurement of task-driven visual attention in the temporal domain can be drawn studying the fixation dispersion across viewers performing a task in a particular context.
2. Visual attention in the temporal domain can be modeled from an accurate estimation of spatio-temporal visual attention.

First, we have introduced in Section [5.2](#) that there is a significant correlation between eye movement sequences of different observers performing the same task when they perceive a surprising or anomalous event [[78](#), [80](#)]. Therefore:

A measurement of task-driven visual attention in the temporal domain can be drawn studying the fixation dispersion across viewers performing a task in a particular context.

Given a crowded and complex scenario, eye movements constitute a useful source to understand how visual attention works and what information should be selected for further processing. In addition, the temporal level of attention of observers might constitute a useful clue to detect suspicious events or anomalous situations to analyze: If all observers fix their attention at the same location at the same time, it is very likely that something noticeable is happening.

Let us consider a video surveillance scenario such as the one presented in the example in Figure [5.13](#), taken from BOSS [[19](#)] database. The database, further described in Appendix [B](#), contains video sequences with anomalous events, which have been recorded in suburban trains. The sequence in (a) shows a man harassing a woman, while (b) includes an anomaly signal e_t , set to 1 when an anomalous situation happens (see scene 4). If we look at the fixation map in (c), we can notice that all observers are fixing their attention at the location of the anomaly. However, it should be noted that other events are also attracting the attention of viewers, such as the moment when the man gets on the train (see scene 2). This may basically happen due to two reasons: there are no more events on

the scene, or there have not been significant changes in the video for long time until those new events. That is why, similarly to spatio-temporal visual attention, attention in the temporal domain measured by fixations should be always considered as an early filtering mechanism, which allows to select those time segments of special importance in a video, candidates to contain anomalies. These often correlate with anomalous situations, especially in complex videos with multiple simultaneous events. Let us emphasize that our objective is not to detect anomalies, but to develop an information filtering mechanism in order to select relevant time segments where a subsequent anomaly detection system may be more efficiently applied.

Hence, just by recording sequences of fixations from different subjects watching videos and computing a dispersion measurement of fixations across viewers, we could provide a [GT](#) measurement of visual attention in the temporal domain.

Second, we have studied along this thesis that spatio-temporal visual attention maps aim to predict viewers fixations in videos or dynamic scenarios. Therefore:

Visual attention in the temporal domain can be modeled from an accurate estimation of spatio-temporal visual attention.

Spatio-temporal [VAMs](#) can be understood, for each frame in a video, as a [2D](#) probability density function which might provide a temporal attention response similar to the one measured from fixations dispersion across viewers. This motivates us to develop a system to model attention in the temporal domain by taking advantage of spatio-temporal visual attention predictions.

In this work, in line with the current dominant approaches in computer vision, we will make use of [CNNs](#) for spatio-temporal visual attention estimation. In addition, we will assess the performance of [LSTMs](#) for attention estimation in the temporal domain, which have shown impressive results for time series forecasting [\[115\]](#), sequence to sequence learning [\[55\]](#) or image captioning [\[206\]](#), among other applications.

Figure [5.14](#) shows the processing pipelines of our supervised approach. First, during the learning phase, the system receives a set of V training videos and extracts several feature maps for visual attention guidance (RGB, motion and objectness), which become high-level features representing the corresponding frames v_t . Moreover, frames in these videos are annotated with eye fixations from several subjects, which can be represented either as spatio-temporal fixations maps or their corresponding temporal responses. All these inputs are used to learn the optimal values for the parameters w of the architecture proposed, where two stages can be differentiated:

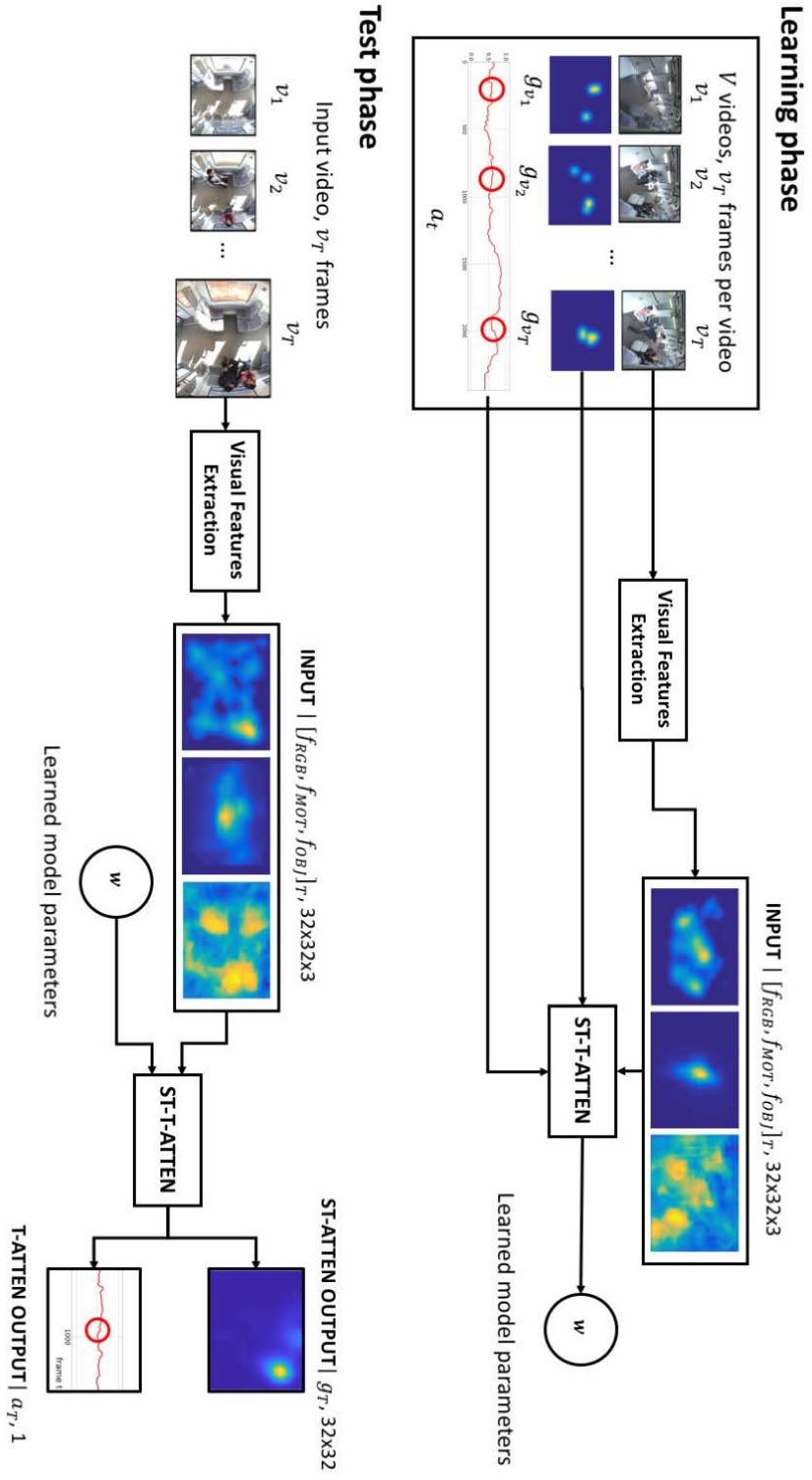


Figure 5.14: Processing pipelines of the ST-T-ATTEN proposed. First, during the learning phase, the system receives a set of training videos and learns the optimal values for its associated parameters. Then, in the test phase, the architecture is qualified to estimate both spatio-temporal visual attention maps and temporal attention responses for new unseen videos.

1. A Spatio-Temporal visual ATtention NETwork (**ST-ATTEN**) consisting on a **CED**, which has the important mission of providing accurate spatio-temporal visual attention maps.
2. A **LSTM**-based Temporal ATtention NETwork (**T-ATTEN**), which will ultimately serve to model attention in the temporal domain.

Then, in the test phase, the **ST-T-ATTEN** is qualified to estimate both spatio-temporal visual attention maps and temporal attention responses for new unseen videos.

5.5.2 Fixation-based temporal ground-truth

Given that the current objective is not to predict a spatial attention response, but a temporal one, we need to generate a frame-level temporal **GT**. As stated in recent behavioral studies [78, 80], and as we introduced in Section 5.5.1, there is a noticeable consistency between observers' eye movements in a scene. Indeed, when an anomalous or suspicious event is happening, gaze locations from different subjects are highly correlated, especially if they are experts or users trained to perform a particular task.

Therefore, on the basis of this fact, we propose a temporal **GT** a_t for each frame v_t , which is computed attending to the dispersion at fixation spatial locations from several subjects. In order to illustrate the **GT** computation process, an example sequence taken from BOSS [19] video surveillance database is shown in Figure 5.13. So far, in Chapter 3, we have defined a binary spatial attention response g_{tn} for each spatial location n in a frame with N_t pixels. This response takes the value of one if the location has been fixated by an observer, and zero otherwise. Hence, given a **GT** soft spatial map g_t that comes from convolving each gaze location with a Gaussian filter, we compute the mean $\mu_{g_t} = (\mu_{g_{tx}}, \mu_{g_{ty}})$ and the standard deviation $\sigma_{g_t} = (\sigma_{g_{tx}}, \sigma_{g_{ty}})$ of the fixation locations:

$$\mu_{g_t} = \frac{\sum_{n=1}^{N_t} g_{tn} \mathbf{x}_{tn}}{\sum_{n=1}^{N_t} g_{tn}} \quad (5.37)$$

$$\sigma_{g_t} = \sqrt{\frac{\sum_{n=1}^{N_t} g_{tn} (\mathbf{x}_{tn} - \mu_{g_{tn}})^2}{\frac{M_t-1}{M_t} \sum_{n=1}^{N_t} g_{tn}}}, \quad (5.38)$$

where $\mathbf{x}_{tn} = (x_{tn}, y_{tn})$ represents the spatial coordinates vector of each location n , and M_t stands for the number of non-zero response locations. Then, the raw temporal attention response a_t can be computed as one minus the weighted mean of the standard deviations along frame width X and height Y :

$$a_t = 1 - \frac{Y\sigma_{g_{tx}} + X\sigma_{g_{ty}}}{X + Y}. \quad (5.39)$$

X and Y are normalizers, which balance the contribution of the standard deviations considered. The response a_t thus takes values between 0, which corresponds to uninteresting frames (maximum fixations dispersion), and 1, which stands out attractive ones (maximum correlation between fixations).

The signal a_t , as shown in Figure 5.13(d), is very noisy. The rationale behind is that observers tend to continuously scan the scene, specially when there are no changes. This noise can be reduced in real-time by applying a first-order infinite impulse response filter. In this work, we make use of an adaptive technique known as Variable Index Dynamic Average (VIDYA) [211], based on an Exponential Moving Average (EMA), in order to filter the temporal attention response.

In addition, the dispersion on which this GT temporal measure is founded may depend on the camera angle and perspective with respect to the objects in the video sequence. For this reason, it is convenient to normalize a_t by subtracting the regular mean, which centers the response around 0, and by dividing it by three times the standard deviation, which covers the ~ 99.7 of the response values, for different camera views. Finally, filtered a_t is clipped to the range $[-1, 1]$ and re-scaled to be in the same interval $[0, 1]$ than the initial raw response. As shown in Figure 5.13(e), this final response a_t is notably softer.

Hypothesis validation

In this section we aim to validate the fundamental hypothesis of this second part of the thesis: attention in the temporal domain can be predicted using the dispersion of gaze locations recorded from several subjects.

Let us demonstrate the existing correlation between anomalies or suspicious events and the GT temporal visual attention response a_t proposed. For that purpose, we make use of GT binary signals e_t , which indicate when anomalies occur in a video sequence, and consider a binary classification problem, where the objective is to classify video frames v_t as anomalous or not, according to their associated filtered a_t . The performance of this proposed fixation-based response is thus assessed by computing the AUC metric [160] for the whole BOSS [19] video surveillance database. The closer is the value of this measure to 1, the higher is the existing correlation between a_t response and the anomaly signal e_t .

After conducting the experiment, we achieve an $AUC = 0.876$, which clearly verifies the correlation between anomalies and a_t . This correlation can be also appreciated in the example in Figure 5.13, where attention achieves its maximum value just before and at the moment of a woman harassment (see scenes 3 and 4). A second

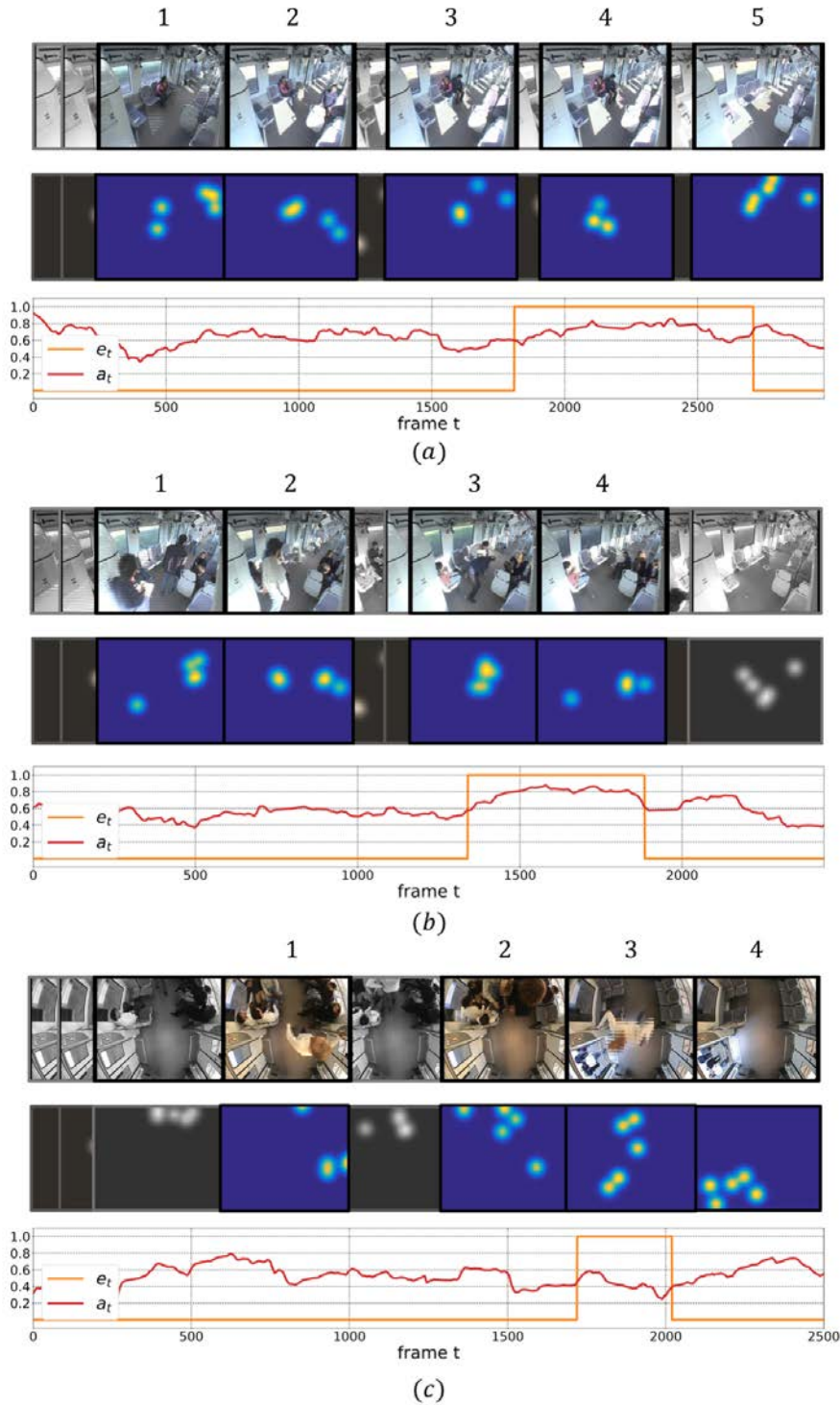


Figure 5.15: Visual attention in the temporal domain a_t modeled in three video-surveillance sequences taken from BOSS [19] database. Anomaly detection signal e_t is also represented, which is set to 1 when an anomaly happens. (a) Woman harassment scene shown in Figure 5.13, taken from a different camera view. (b) Passengers fighting for a newspaper. 1. A first man with a newspaper comes into the wagon. 2. A second man arrives. 3. Second man hits the first and destroys the newspaper. 4. First man lies on the ground. (c) Panic scene. 1. Passengers are warned about an accident. 2-3. Passengers running out of the wagon. 4. Nobody left on the train.

example is provided in Figure 5.15(b), where two passengers fight for a newspaper. There are more people on the wagon when this situation happens (scenes 3 and 4). At that moment, fixation-based attention response becomes high, so it successfully captures the anomaly.

Furthermore, we would also like to discuss why, although highly correlated, visual attention and anomaly are not equal variables, and therefore $AUC < 1$. From a theoretical perspective, there are cases in which consistent visual attention is achieved in the absence of anomalies, such as simple scenarios with few people, as the one in Figure 5.15(a), which is the same woman harassment scene shown in Figure 5.13, taken from a different camera view. Another error case is shown in Figure 5.15(c). The video sequence includes a panic scene, where passengers are warned about an accident (see scene 1) and run out of the wagon (see scenes 2-3). Scene involves all passengers and covers the whole image, so fixations dispersion is high at the moment of the anomaly (scene 3).

Hence, as stated in the previous section, temporal attention can be seen as a filtering mechanism, which often correlates with anomalous situations. This correlation is particularly high in complex scenarios with multiple simultaneous events, which are those that require a greater cognitive effort to be understood. Considering that the proposed temporal attention response a_t constitutes a filtering mechanism to be applied prior to an anomaly detection system, it is critical to obtain a low probability of non-detection ($P_{ND} \simeq 0$) with this signal, while we accept higher false-alarm probabilities ($P_{FA} > 0$).

5.5.3 Model overview

In this section, we overview the Spatio-Temporal to Temporal visual ATtention NETwork (**ST-T-ATTEN**) proposed. The complete architecture of the system is represented in Figure 5.16. Our approach is built on the combination of two modules, which are described in the following sections: 1) A Spatio-Temporal visual ATtention NETwork (**ST-ATTEN**) for spatio-temporal visual attention estimation; 2) A Temporal ATtention NETwork (**T-ATTEN**) for modeling visual attention in the temporal domain.

First, given a video v , the **ST-ATTEN** module consists in a Convolutional Encoder Decoder (**CED**). At each timestep t , it receives an input frame v_t , which is represented by means of high-level feature maps for visual attention guidance, such as the ones described in Section 5.4.3: $[f_{RGB}, f_{MOT}, f_{OBJ}]_{v_t}$. The network is then able to compute, for each frame v_t , either a latent representation \mathbf{z}_t of visual attention or a spatio-temporal visual attention map \hat{g}_t .

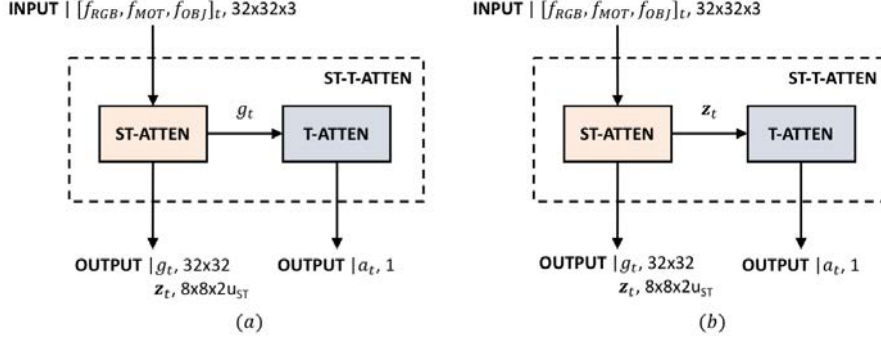


Figure 5.16: Diagram of the **ST-T-ATTEN** proposed. The approach is built on the combination of two modules: 1) A **ST-ATTEN** for spatio-temporal visual attention estimation; 2) A **T-ATTEN** for modeling visual attention in the temporal domain. First, **ST-ATTEN** receives at each timestep t a frame represented by three feature maps $[f_{RGB}, f_{MOT}, f_{OBJ}]_t$. Then, **T-ATTEN** receives as input (a) the spatio-temporal **VAM** g_t or (b) the latent representation z_t obtained at the output of **ST-ATTEN** and estimates a temporal attention response \hat{a}_t .

Secondly, we propose two versions of the **T-ATTEN** module, which are compared. In particular, the **T-ATTEN** receives as input one of the two outputs provided by the **ST-ATTEN**: either the spatio-temporal **VAM** \hat{g}_t (see Figure 5.16(a)) or the latent representation z_t (see Figure 5.16(b)). Then, it estimates, for each frame v_t , a temporal attention response \hat{a}_t .

5.5.4 Spatio-Temporal Visual Attention Network

The first stage of our **ST-T-ATTEN** receives as input, for each frame v_t in a video v , a set of 32×32 feature maps for visual attention guidance, such as the ones introduced in Section 5.4.3 ($[f_{RGB}, f_{MOT}, f_{OBJ}]_{v_t}$), and estimates spatio-temporal visual attention by means of a **CED** architecture, which we have called Spatio-Temporal visual ATtention NETwork (**ST-ATTEN**).

Figure 5.17(a) illustrates the proposed **ST-ATTEN**. As explained in section 5.3.3, **CEDs** can be decomposed into two networks: an encoder and a decoder network. In our system, the first network encodes input feature maps into a latent representation $\mathbf{z} = f_E([f_{RGB}, f_{MOT}, f_{OBJ}]_{v_t}, \theta_E)$, while the second symmetric network transforms this representation into a spatio-temporal visual attention map $\hat{g}_t = f_D(\mathbf{z}_t, \theta_D)$. Due to the low dimensionality of the input features, we decided to make use of dilated convolutions in an attempt to keep the number of parameters limited. Therefore, the encoder network consists of two $k = 3 \times 3$ dilated **CONV** layers with $d = 2$, and u_{ST} and $2u_{ST}$ filters, respectively. These layers correspond to two dilated **CONV** layers in the decoder network with the same

Table 5.1: Encoder and decoder architectures for the CONV-ST-ATTEN and CONV LSTM-ST-ATTEN configurations of the ST-ATTEN proposed.

(a) CONV-ST-ATTEN								
Encoder f_E			Decoder f_D					
Input $[f_{RGB}, f_{MOT}, f_{OBJ}], 32 \times 32 \times 3$			Output $g = f_D(\mathbf{z}, \theta_D), 32 \times 32, \text{KL loss}$					
CONV E1	$(k = 3 \times 3, u_{ST}, s = 1, d = 2)$	MAX POOL 1	$(k = 2 \times 2, s = 2),$	CONV D1	$(k = 3 \times 3, u_{ST}, s = 1, d = 2)$	UPSAMPLING 1	ELU	
CONV E2	$(k = 3 \times 3, 2u_{ST}, s = 1, d = 2)$	MAX POOL 2	$(k = 2 \times 2, s = 2),$	ELU	CONV D2	$(k = 3 \times 3, 2u_{ST}, s = 1, d = 2)$	UPSAMPLING 2	ELU
Latent representation $\mathbf{z} = f_E([f_{RGB}, f_{MOT}, f_{OBJ}], \theta_E), 8 \times 8 \times 2u_{ST}$								
(b) CONV-LSTM-ST-ATTEN								
Encoder f_E			Decoder f_D					
Input $[f_{RGB}, f_{MOT}, f_{OBJ}], 32 \times 32 \times 3$			Output $g = f_D(\mathbf{z}, \theta_D), 32 \times 32, \text{KL loss}$					
CONV-LSTM E1	$(k = 3 \times 3, u_{ST}, d = 2)$	MAX POOL 1	$(k = 2 \times 2, s = 2)$	ELU	CONV-LSTM D1	$(k = 3 \times 3, u_{ST}, d = 2)$	UPSAMPLING 1	ELU
CONV E1	$(k = 3 \times 3, 2u_{ST}, d = 2)$	MAX POOL 2	$(k = 2 \times 2, s = 2)$	ELU	CONV D1	$(k = 3 \times 3, 2u_{ST}, d = 2)$	UPSAMPLING 2	ELU
Latent representation $\mathbf{z} = f_E([f_{RGB}, f_{MOT}, f_{OBJ}], \theta_E), 8 \times 8 \times 2u_{ST}$								

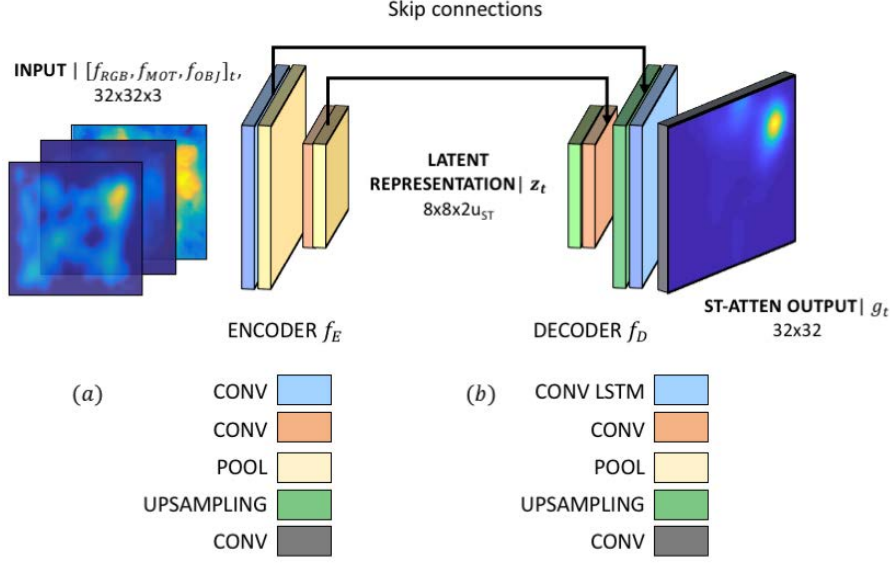


Figure 5.17: Diagram of the **ST-ATTEN** for spatio-temporal visual attention estimation, which consists of a **CED** architecture. Both encoder and decoder have two dilated **CONV** layers with skip connections. Two different configurations are proposed for this stage of the system: (a) **CONV-ST-ATTEN** (b) **CONV-LSTM-ST-ATTEN**.

number of units. After each dilated **CONV** layer in the encoder network, MAX **POOL** with a 2×2 window and $s = 2$ is performed. This operation sub-samples the **CONV** layer output by a factor of 2, which allows to generate representations more robust to spatial translations. In the decoder network, these layers are replaced by upsampling operations. **ELU** activations [180] are introduced after each dilated **CONV** layer. **ELUs** are similar to typically used **ReLU**s, but provide more robustness to noise activations with mean close to zero. Moreover, we introduce skip connections between corresponding dilated **CONV** layers in the encoder and decoder, in an attempt to preserve the spatial resolution of the down-sampled input feature maps. After the decoder network, a final 3×3 **CONV** layer with linear activation generates the output visual attention map \hat{g}_t , which has the same dimensions than the input features (32×32).

We propose two configurations of **ST-ATTEN**: **CONV-ST-ATTEN** (see Figure 5.17(a)), which has been described in the previous paragraph, and **CONV-LSTM-ST-ATTEN** (see Figure 5.17(b)). The architectures of both configurations are further detailed in Table 5.1. The main difference between these two networks relies on the outer layers of the encoder and the decoder, which are **CONV** in the first approach and **CONV-LSTM** in the second one. Although our system already receives dynamic spatio-temporal information from the input optical flow-based feature map f_{MOT} , we include **CONV-LSTMs** in an attempt to model viewers dynamic behavior during the training phase,

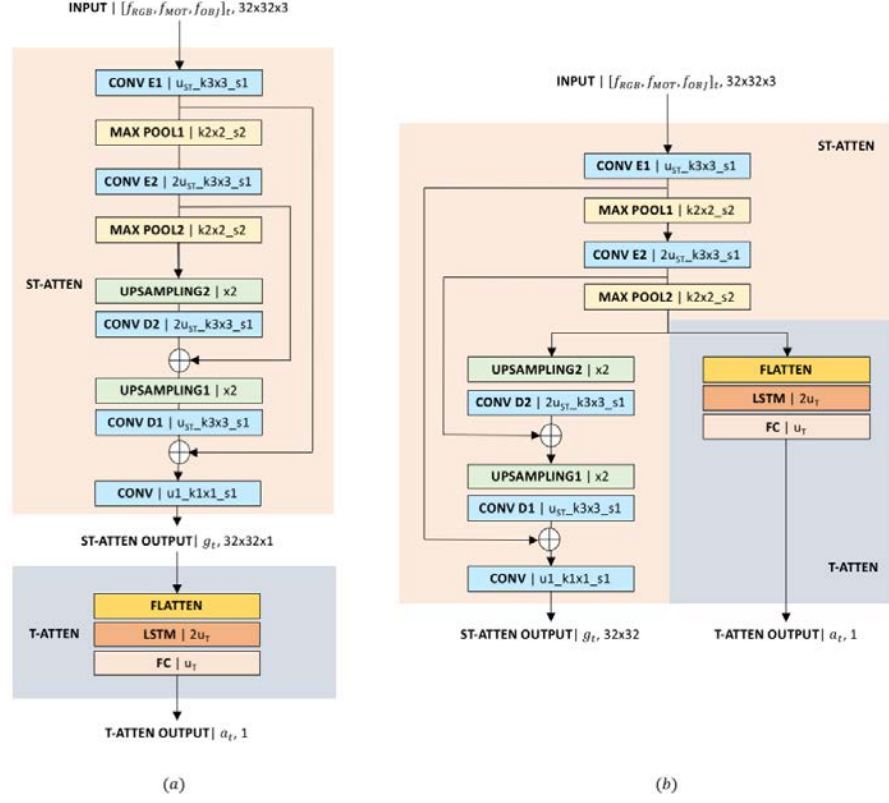


Figure 5.18: Architecture diagrams of the **ST-T-ATTEN** configurations proposed. The network receives as input, at each timestep t , a set of three feature maps $[f_{RGB}, f_{MOT}, f_{OBJ}]_t$, and generates two outputs: a spatio-temporal **VAM** \hat{g}_t and a temporal attention response \hat{a}_t . (a) In the first design, the input to the **T-ATTEN** is the latent representation \mathbf{z}_t computed by the **ST-ATTEN** encoder. (b) In the second configuration, the input to the **T-ATTEN** is the **VAM** \hat{g}_t estimated by the **ST-ATTEN** decoder. Layers are defined by their number of units u , kernel size k and stride s with which the filters are slid.

taking advantage of the spatio-temporal information provided by fixation sequences. This information might be helpful to improve the accuracy of the predicted **VAMs** in complex scenarios where there are more than one object in motion, so that it is necessary to consider previous conspicuous locations for a better attention guidance.

5.5.5 Temporal Attention Network

The second stage of the system, denoted as Temporal ATtention NETwork (**T-ATTEN**), connects with the **ST-ATTEN** explained above and estimates attention in the temporal domain. We expect to take advantage of the effective sequential representations provided by a **LSTM**-based architecture in this stage. Hence, we propose an architecture composed by a layer with u_T **LSTM** units, which allows

the system to learn long-term dependencies between temporal attention responses a_t associated to frames v_t in the same video v , avoiding vanishing or exploding gradients, as explained in Section 5.3.4. This layer is followed by a FC layer, with $u_T/2$ units. Finally, a simple FC layer produces the temporal attention variable \hat{a}_t .

The combination of ST-ATTEN for spatio-temporal visual attention estimation and T-ATTEN for modeling attention in the temporal domain gives rise to our Spatio-Temporal to Temporal visual ATtention NETwork (ST-T-ATTEN). We propose two different configurations for the architecture, as shown in Figure 5.18. Both approaches differ in the connection between ST-ATTEN and T-ATTEN modules. In the first design, the T-ATTEN module is fed at each timestep t with the estimated 32×32 \hat{g}_t at the output of the ST-ATTEN decoder. In contrast, in the second configuration, the $8 \times 8 \times 2u_{ST}$ latent representation \mathbf{z}_t extracted by the encoder network is used as input to the temporal network. For each frame v_t of a video v , both architectures generate at its output a linear temporal attention response \hat{a}_t .

5.5.6 Implementation details

Data preprocessing

First, we normalize feature maps at the input of the ST-T-ATTEN by subtracting the feature mean and dividing by three times the feature standard deviation (which covers the $\sim 99.7\%$ of the data samples). Mean and standard deviation are computed over the training set. Then, we also clip feature maps values to the range $[-1, 1]$. Furthermore, we normalize VAMs at the output of the ST-ATTEN module to sum to 1.

Multitask loss function

In order to train each of the stages of the ST-T-ATTEN, we have considered two different loss functions.

On the one hand, in order to train the ST-ATTEN stage, we make use of the Kullback-Leibler divergence (KL), which is a distribution-based metric, frequently used as a loss function to train CNNs for visual attention prediction due to its proven efficiency [101]. Given a frame v_t with N_t spatial locations, its corresponding fixation map g_t and a predicted visual attention map \hat{g}_t , it is defined as:

$$KL(g_t, \hat{g}_t) = \sum_{n=1}^{N_t} g_{tn} \log \left(\epsilon + \frac{g_{tn}}{\epsilon + \hat{g}_{tn}} \right). \quad (5.40)$$

For each spatial location n , g_{tn} constitutes its associated fixation-based GT, resulting from convolving g_t by a Gaussian filter

with standard deviation equal to one degree of visual angle, in order to obtain a continuous distribution. In addition, \hat{g}_{tn} represents the visual attention predicted for that location. A lower score value indicates a better approximation of the VAM \hat{g}_t to the fixation map g_t .

On the other hand, the T-ATTEN stage is trained by using MSE as loss function (see section 5.3.1, Eq. 5.11). Given a video frame v_t , the MSE between its associated GT temporal attention response a_t and the attention response estimated by our T-ATTEN \hat{a}_t can be written as follows:

$$MSE(a_t, \hat{a}_t) = (\hat{a}_t - a_t)^2 \quad (5.41)$$

Finally, for each frame v_t , the multitask loss function for the overall system is expressed as follows:

$$L_{ST-T-ATTEN}(g_t, \hat{g}_t, a_t, \hat{a}_t) = KL(g_t, \hat{g}_t) + \alpha MSE(a_t, \hat{a}_t), \quad (5.42)$$

where α is a scalar that balances the contribution of the two loss functions, and has been empirically determined, as described in the next chapter.

EXPERIMENTS ON TEMPORAL VISUAL ATTENTION ESTIMATION IN A VIDEO SURVEILLANCE SCENARIO

6.1 INTRODUCTION

The experiments presented in this chapter have as main objective to assess various configurations for training end-to-end the Spatio-Temporal to Temporal visual ATtention NETwork ([ST-T-ATTEN](#)) architecture proposed in Chapter 5. As explained in that chapter, our system ultimately models attention in the temporal domain by aligning spatio-temporal visual attention maps estimated from video frames to frame-level fixation-based temporal attention responses.

To this end, we evaluate our system in the video surveillance scenario defined by the BOSS [19] database, which contains video sequences recorded in a railway transport context, showing different types of suspicious or anomalous events.

CHAPTER OVERVIEW

First, the experimental design is introduced in Section 6.2. Secondly, we determine the optimal architecture for the first stage of the system proposed, [ST-ATTEN](#), in Section 6.3. Then, in Section 6.4, the [ST-ATTEN](#) optimal configuration is used to train end-to-end the complete [ST-T-ATTEN](#), and to provide results on attention estimation in the temporal domain, which allows to discuss the model strenghts and limitations. In Section 6.5, we motivate the use of our system for guiding anomaly detection in video surveillance applications. Finally, Section 6.6 summarizes our conclusions and motivates future work.

6.2 EXPERIMENTAL DESIGN

This section explains the experimental design for the analysis of the [ST-T-ATTEN](#). First, the databases used to train and evaluate its different stages are introduced. Then, we describe the experimental

setup and the evaluation metrics considered to assess the performance of the proposed architectures. Finally, we provide the implementation details related both to the [ST-T-ATTEN](#) and the [CNNs](#) previously presented in Section 5.4, in charge of extracting high-level visual feature maps that become the input to our model.

6.2.1 Databases

SALICON and DIEM

SALICON [114] and DIEM [16] databases have been considered in our experiments to train the modified ResNet-50 [193] models for visual feature extraction introduced in Section 5.4.

On the one hand, in order to obtain RGB-based spatial feature maps, we have made use of SALICON image database to fine-tune a ResNet-50 model pre-trained on ImageNet database [3]. SALICON database contains a set of 10,000 context-generic training images annotated by 60 “free-viewing” observers with a mouse tracking process. Although the consistency between participants’ fixations is lower when compared to the information provided by an eye tracker device, mouse movements constitute a helpful approach to eye tracking in still images, which allow to efficiently annotate very large databases.

On the other hand, DIEM video database, already introduced and analyzed in Chapter 4, has been used to train from scratch the optical flow-based network for motion feature maps extraction.

BOSS

Within the framework of the BOSS project [19], a database with 15 video sequences recorded in RENFE suburban trains in Madrid was released with the aim of developing an efficient transmission system for video-surveillance in a railway transport context. Videos contain events such as a cell phone theft, a fight between passengers, a disease in public and several women harassment. Moreover, two additional sequences with no incidents are included. For each event, three camera views are provided.

In order to evaluate the different architectures for attention estimation in the temporal domain that were proposed in the previous chapter, we have selected the three camera views of 10 sequences from this database, and annotated them with eye fixations. In total, 30 videos (over 84,000 video frames, 56 minutes total, 720×576) have been used. For each video, eye traces from 5 observers have been recorded by using a 250 Hz SMI RED250mobile Eye Tracker system [212]. The complete list of videos annotated for our experiments can be found in Appendix B.

6.2.2 Experimental setup

The experiments presented in this chapter have the objective to assess which of the proposed **ST-T-ATTEN** architectures models better the temporal dimension of attention. For that purpose, we will ultimately evaluate our system when estimating visual attention in the temporal domain.

We conduct our experiments by splitting the 30 videos selected from BOSS [19] database into three folds, each one containing 10 different sequences. In order to avoid over-fitting, the three camera views from the same sequence are grouped together in the same fold. In the following paragraphs, we will describe the process we have followed to train and optimize the different modules involved in our **ST-T-ATTEN**.

First, we have determined the optimal configuration for the first stage of our system: **ST-ATTEN**. To achieve this, we have evaluated the two configurations proposed for this module in Section 5.5.4. We selected the one that provided the best performance in terms of the evaluation metrics mentioned in the next section. Following a 3-fold cross-validation procedure, we have estimated spatio-temporal visual attention maps for each video in a fold, using the remaining two folds for training the network.

Then, we have assessed if our **ST-T-ATTEN** is able to obtain accurate estimations of temporal attention, with the ultimate goal of discussing its utility as a filtering mechanism for subsequent anomaly detection systems in video surveillance scenarios. To do this, we have trained the complete **ST-T-ATTEN** architecture end-to-end. To do so, we initialized the **ST-ATTEN** module using the weights learned in the previous step. As mentioned before, latent representations and **VAMs** extracted by **ST-ATTEN** constitute the input for the next stage of the system: **T-ATTEN**. The evaluation of the complete network has been done by following the same procedure described above: we estimate either spatio-temporal visual attention or attention in the temporal domain in each fold, using the remaining two folds for training the complete architecture.

6.2.3 Evaluation metrics

Spatio-temporal visual attention estimation

In order to evaluate the spatio-temporal **VAMs** predicted by the different **ST-ATTEN** architectures in the first stage, we make use of the **sAUC** and **sNSS** metrics described in Section 4.2.3. Moreover, we include the **KL** metric defined in Eq. 5.40, which is also the loss function used to train the **ST-ATTEN**. Furthermore, for comparison purposes, we consider the three baseline models described in Section 4.2.3: CHANCE, CENTER and H50.

Visual attention estimation in the temporal domain

In addition, since our temporal attention response a_t is a real number in the range $[0, 1]$, we now introduce a rank correlation coefficient for the assessment of the temporal attention responses estimated by the **T-ATTEN** proposed in Chapter 5. The evaluation metric chosen is the well-known Pearson Correlation Coefficient (PCC) [213], which here measures the linear relationship between the fixation-based temporal **GT** $a_{1:T}$ and the estimated temporal attention response $\hat{a}_{1:T}$ of a given set of T video frames. It can be written as follows:

$$PCC = \frac{\sum_{t=1}^T (a_t - \mu_{a_{1:T}})(\hat{a}_t - \mu_{\hat{a}_{1:T}})}{\sqrt{\sum_{t=1}^T (a_t - \mu_{a_{1:T}})^2} \sqrt{\sum_{t=1}^T (\hat{a}_t - \mu_{\hat{a}_{1:T}})^2}}, \quad (6.1)$$

where $\mu_{a_{1:T}}$ and $\mu_{\hat{a}_{1:T}}$ represent the mean values of the **GT** a_t and the estimated attention \hat{a}_t for the considered set of frames, respectively. The coefficient lies in the range $[-1, 1]$, meaning these extreme values an exact negative or positive correlation, while $PCC = 0$ implies no correlation at all.

6.2.4 Training and implementation details

Here we describe the implementation details of the complete system. To begin with, it should be mentioned that we have made use of the Keras framework [214] with Tensorflow backend [215] to build all the networks deployed. Besides, we train our models using a 12GB NVIDIA GeForce GTX TITAN Xp **GPU** on a system with an Intel Core i7-6700K (4.00GHz) **CPU** and 32GB of **RAM**.

Feature extraction networks for visual attention guidance

First, in order to train the ResNet-50 [193] models to compute RGB-based and optical flow-based feature maps, we have considered SALICON [114] image and DIEM [16] video databases, respectively. Besides, we have chosen the **KL** loss function introduced in Eq. 5.40 as loss function. In order to minimize **KL**, we use **SGD**, setting the learning rate to 10^{-4} and using a mini-batch of 10 samples.

On the other hand, we have directly used the model and weights provided by the authors of the selected objectness-based network [2], which has been trained on MSRA-B [216] salient object database.

Spatio-Temporal to Temporal Visual Attention Network

Secondly, **ST-T-ATTEN** is trained in BOSS [19] video surveillance database to model attention in the temporal domain.

Regarding the model weights initialization, we have drawn on Glorot uniform initializer [217], also known as Xavier uniform

initializer, which generates random samples from a uniform distribution.

In addition, some preliminary experiments have shown the utility of Dropout [187] regularization in the ST-ATTEN module. Therefore, we have decided to introduce dropout layers in this stage of the system, which randomly drop, at each iteration, half of the filters ($p = 0.5$) of each CONV layer, both at the encoder and the decoder.

With respect to the loss functions associated to each stage of the system, we have considered the multitask loss function introduced in Section 5.5.6. This loss is a linear combination of the two atomic losses: a) the KL loss presented in Eq. 5.40, used to train the ST-ATTEN stage of the system; and b) the MSE loss defined in Eq. 5.41, used to train the T-ATTEN module of the system. As defined in Eq. 5.42, parameter α balances the contribution of each term (KL-loss, MSE-loss) in the multitask loss. This parameter has been empirically set to $\alpha = 100$. In order to minimize the ST-T-ATTEN loss function $L_{ST-T-ATTEN}$, we consider SGD with a learning rate of 10^{-4} . The network is trained over 10K iterations, using a mini-batch of 256 samples.

6.3 RESULTS ON SPATIO-TEMPORAL VISUAL ATTENTION ESTIMATION WITH ST-ATTEN

In this section, we aim to assess the two architectures for the ST-T-ATTEN module proposed in Section 5.5.4. We evaluate both qualitatively and quantitatively the two configurations: CONV-ST-ATTEN and CONV-LSTM-ST-ATTEN. To that end, we have first validated the parameter u_{ST} associated with both networks. This parameter determines the dimension of the latent representation extracted for visual attention ($8 \times 8 \times 2u_{ST}$). After conducting several experiments with both undercomplete and overcomplete EDNs, we set this parameter to $u_{ST} = 64$, which corresponds to an overcomplete CED, in which the representation \mathbf{z} has a dimension greater than the network input (the three feature maps).

Then, we perform a quantitative evaluation of the two configurations proposed in terms of the sAUC, sNSS and KL metrics. Table 6.1 summarizes the results obtained by the proposed ST-ATTEN on the BOSS [19] database. For the sake of comparison, we include the results achieved by the three considered feature maps: RGB-based (RGB), motion (MOT) and objectness (OBJ). Moreover, we provide the results offered by two simple fusion approaches: a map computed by averaging the three features (AVERAGE) and a map obtained by learning a linear combination over them (LIN. COMB). Finally, the three reference models introduced in Section 4.2.3 (H50, CHANCE, CENTER) are also included for comparison. As can be verified from CENTER baseline, both sAUC and sNSS metrics are not affected by center bias.

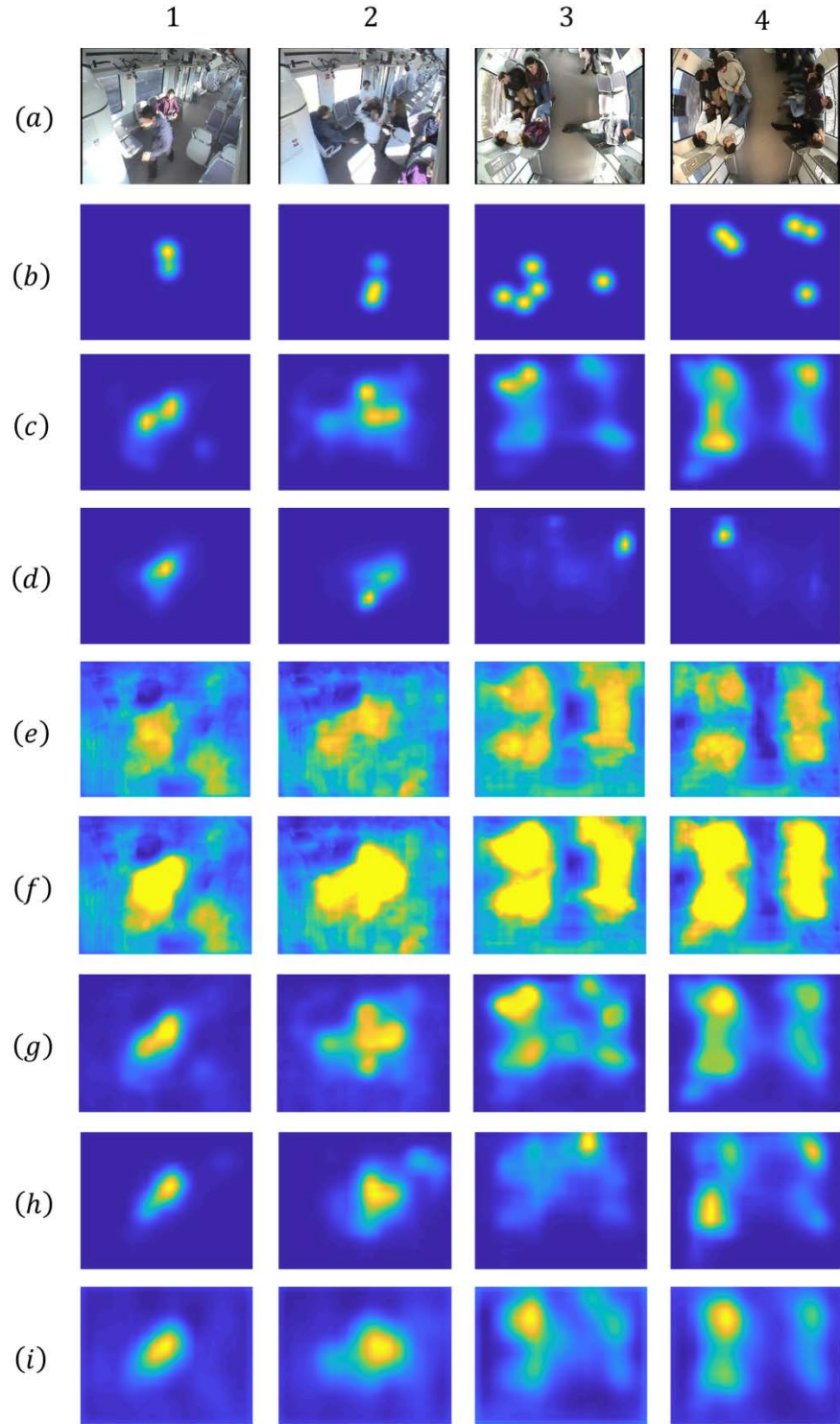


Figure 6.1: Visual attention maps obtained by [ST-ATTEN](#) for some example frames taken from BOSS [\[19\]](#) database. (a) Original frames. (b) [GT](#) fixations map. (c) RGB-based feature map. (d) Motion feature map. (e) Objectness feature map. (f) [VAM](#) obtained from averaging the three feature maps. (g) [VAM](#) obtained by learning a linear combination over the three feature maps. (h) [CONV-ST-ATTEN](#) [VAM](#). (i) [CONV-LSTM-ST-ATTEN](#) [VAM](#).

Table 6.1: Results obtained on the BOSS [19] database by the proposed ST-ATTEN and other methods for comparison when estimating spatio-temporal visual attention.

Model	<i>sAUC</i>	<i>sNSS</i>	<i>KL</i>
	<i>mean (C.I.)</i>	<i>mean (C.I.)</i>	<i>mean (C.I.)</i>
RGB	0.714 (0.710, 0.718)	0.374 (0.367, 0.381)	1.938 (1.922, 1.954)
MOT	0.588 (0.583, 0.592)	0.127 (0.119, 0.134)	2.755 (2.730, 2.781)
OBJ	0.667 (0.663, 0.671)	0.303 (0.296, 0.310)	2.257 (2.249, 2.264)
AVERAGE	0.702 (0.698, 0.710)	0.362 (0.354, 0.369)	2.039 (2.031, 2.048)
LIN. COMB.	0.702 (0.698, 0.706)	0.363 (0.356, 0.370)	1.894 (1.883, 1.905)
CONV	0.750 (0.746, 0.755)	0.424 (0.416, 0.431)	1.563 (1.547, 1.578)
CONV-LSTM	0.748 (0.744, 0.751)	0.410 (0.403, 0.416)	1.610 (1.601, 1.627)
H50	0.826 (0.826, 0.826)	0.692 (0.692, 0.693)	2.137 (2.135, 2.139)
CHANCE	0.500 (0.500, 0.500)	−0.000 (−0.001, 0.000)	4.423 (4.335, 4.338)
CENTER	0.500 (0.499, 0.502)	0.014 (0.012, 0.017)	4.337 (4.335, 4.338)

According to the results in the table, our CED configurations for spatio-temporal visual attention estimation successfully learn non-linear fusion schemes for the three feature maps considered and notably outperform the two baseline fusion approaches (AVERAGE, LIN. COMB.). Furthermore, it can be concluded that the first configuration proposed, which only makes use of CONV layers, is the one that offers a slightly better performance. However, whereas the performance achieved by ST-ATTEN is close to H50 score in terms of sAUC, we are still far from reaching highly accurate estimations of spatio-temporal visual attention according to sNSS metric. This is probably due to some of the aspects discussed by Bylinskii et al. in [116], mainly related to the modeling of high-level concepts, such as objects of action or gaze, which are not explicitly modeled by our approach. It is also worth pointing out that H50 model performs worse than our approaches in terms of distribution-based KL metric. From our point of view, this is not a surprise, as predicting visual attention from fixations of 50% of subjects available is hard when fixation dispersion across observers is high, due to the absence of a clear and conspicuous event on the scene. This happens in the majority of frames of the database, which do not contain anomalous events.

Finally, we have also assessed qualitatively the VAMs obtained by the different methods. Figure 6.1 shows the output maps of the two configurations proposed for some example frames taken from BOSS [19] database. GT eye fixation density maps are also displayed (b). For the sake of comparison, we include the three feature maps used as input to ST-ATTEN (c,d,e). Moreover, we provide the output VAMs

obtained by the two baseline fusion approaches: AVERAGE (f) and LIN. COMB. (g). As can be seen, our architectures (h,i) provide the most accurate estimations of visual attention in the shown cases, providing better attention representations than baseline fusion approaches and individual features. In addition, they are able to learn more sophisticated representations of visual attention in complex situations, as the one shown in the fourth example provided, in contrast to maps obtained from a linear combination, which gives a higher weight to the feature that best model visual attention amongst the three maps considered (RGB-based, according to the results in Table 6.1). However, it can be appreciated that our model fails in estimating attention in crowded scenes with several people, as in the example gathered on the third column.

If we compare the maps estimated by the two proposed CED configurations, it can be noticed that those obtained by a CONV-LSTM architecture (i) are in general smoother than the ones computed by a fully CONV one (h), probably due to the influence of information from previous frames stored in the LSTM state cells. This information may help to reduce noise when feature maps associated to the current frame do not constitute good estimations of the visual attention. For instance, in the third case, estimated motion is concentrated on a very small location and, therefore, it is less attracting than passengers on the train. While CONV-ST-ATTEN map is strongly affected by this motion, CONV-LSTM-ST-ATTEN uses information stored from previous frames to keep the attention on passengers.

However, given the fact that quantitative results achieved by both configurations are quite similar, and that the use of CONV-LSTMs is more computationally demanding, from now on we will make use of the CONV-ST-ATTEN configuration, with $u_{ST} = 64$ filters, as input to the T-ATTEN stage of the system.

6.4 RESULTS ON ATTENTION ESTIMATION IN THE TEMPORAL DOMAIN WITH ST-T-ATTEN

Once we selected the optimal configuration for the ST-ATTEN module, we have trained end-to-end the whole ST-T-ATTEN architecture proposed, as described in Section 6.2. Results obtained on the BOSS [19] database are summarized here below.

First, it is worth mentioning that training the whole architecture barely changes the performance achieved by the ST-ATTEN module; consequently, for the sake of simplicity, we have omitted its analysis in this section, focusing our discussion in the ultimate goal of our system, which is to estimate attention in the temporal domain, by means of the T-ATTEN module. The parameter u_T of this module has

Table 6.2: Results obtained on the BOSS [19] database by the proposed ST-T-ATTEN when modeling visual attention in the temporal domain.

Input \ Architecture	FC	LSTM + FC
	PCC	PCC
VAM	0.166	0.321
LR	0.106	0.323
GT	0.278	0.467

been empirically set to $u_T = 256$, which determines the number of units of its corresponding LSTM (u_T) and FC ($u_T/2$) layers.

As described in Section 5.5.5, we have evaluated two configurations of the T-ATTEN architecture: 1) one that takes the VAM estimated by the ST-ATTEN decoder as input (VAM), and 2) one that works with the latent representation (LR) extracted by the encoder network. Besides, for comparison purposes, we have also considered a baseline T-ATTEN architecture composed of a unique FC layer with u_T units. Finally, we have also trained two reference T-ATTENs (GT), that work with GT VAMs g_t computed using GT fixations recorded from subjects. These two approaches provide a theoretical upper bound of the performance of the system, simulating scenarios in which VAMs are optimal.

Table 6.2 presents the results obtained on the BOSS [19] database by our proposed ST-T-ATTEN for visual attention estimation in the temporal domain, as well as by the models considered for comparison. Results are presented in terms of the PCC (see Section 6.2.3). As can be appreciated, LSTM-based T-ATTEN architectures notably outperform FC ones, which confirms the benefits of using LSTM units to model short and long-term temporal relationships between video frames. Regarding the performance of our approaches, both of them provide a similar PCC score, in spite of the advantages that latent representations have shown in several multi-task learning paradigms [218]. Therefore, we can conclude that, according to the experiments taken on the BOSS [19] database, the optimal configuration for the ST-T-ATTEN proposed consists in a CONV-ST-ATTEN, and a T-ATTEN with LSTM and FC layers. The latter can be fed either with the VAMs computed at the output of the ST-ATTEN decoder or the latent representations generated by the encoder network.

Although we are still far from achieving an exact correlation with the GT temporal attention response a_t that we aim to estimate, it should be noted that PCCs obtained by our models are not far from those provided by the theoretical bounds. As can be seen in the example sequence in Figure 6.3, the temporal attention responses \hat{a}_t

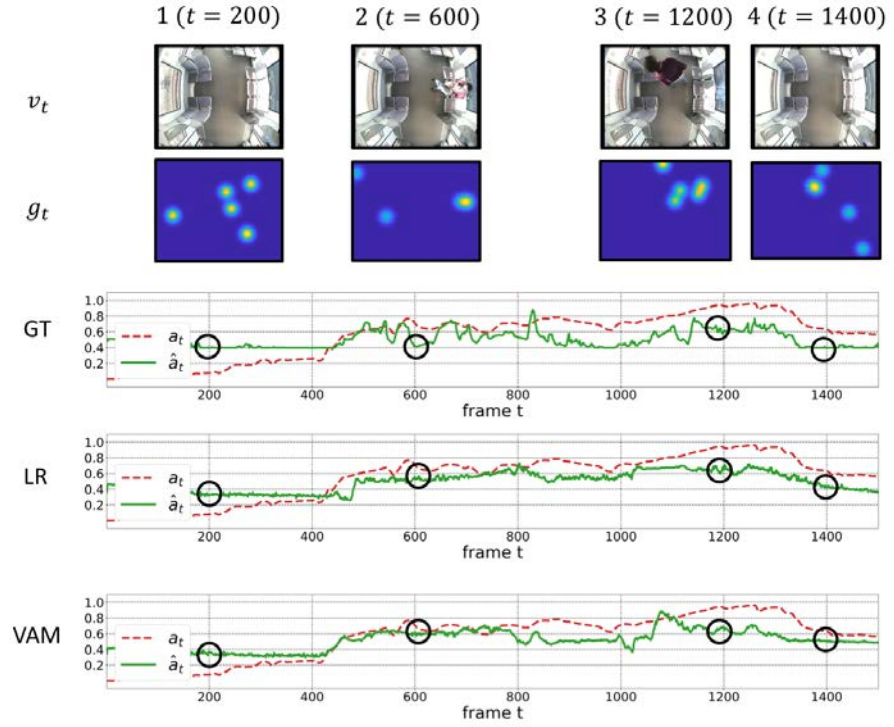


Figure 6.2: Visual attention in the temporal domain \hat{a}_t estimated by [ST-T-ATTEN](#) in a video-surveillance sequence taken from BOSS [19] database. The sequence shows the theft of a mobile phone. 1. Empty train wagon. 2. Woman sits on the train. 3. Somebody stole woman's mobile phone. 4. Woman has gone to report the incident. On the first two rows, some video frames over time v_t are shown, together with their associated [GT](#) fixations map g_t . Then, the temporal attention response \hat{a}_t estimated by each of the three [LSTM](#)-based [T-ATTEN](#) evaluated ([GT](#), [LR](#) and [VAM](#)) is displayed. For the sake of comparison, the [GT](#) temporal attention response a_t is also plotted.

estimated by our [LSTM](#)-based approach have a lower dynamic range compared to the expected [GT](#) fixation-based responses a_t . This may be due to several reasons:

- Features for visual attention guidance considered might not be accurate enough for modeling spatio-temporal visual attention. We should improve them, or even incorporate additional ones, in order to handle crowded and complex scenes, such as the ones shown on the third and fourth columns of Figure 6.1.
- Besides, it has been observed that when the scene is static or there are not obvious conspicuous locations on it, the [ST-ATTEN](#) is not always able to model the pseudo-random nature of eye fixation sequences, which approximately corresponds to [VAMs](#) that distribute visual attention equally in all spatial locations.

Instead, it estimates the same **VAM** for static frames with the same content, which may result in a flat temporal response.

- We have trained **ST-T-ATTEN** with a database that contains few and similar anomalous events, which might not be sufficient to demonstrate our second assumption in Section 5.5.1. To address this issue, we aim to make use of large-scale video surveillance databases such as VIRAT [22] or UCF-Crime [23] for a more complete analysis of the system.

Moreover, given that the **T-ATTEN** has not been able to accurately estimate visual attention in the temporal domain even when the **VAMs** are optimal, we come to the following conclusion, which serves to lead further work: our system, which makes use of traditional **CNNs** layers and the widely-used **MSE** loss for regression, does not compute adequately the function that maps **GT** fixations maps g_t with temporal attention responses a_t . Therefore, future efforts will be made towards designing **CNNs** layers and a loss function tailored to the problem we want to solve with **T-ATTEN** (estimation of attention in the temporal domain), in order to improve our system to efficiently guide anomaly detection in video surveillance scenarios. An alternative to **MSE** that would perhaps be worth testing is the Mean Pairwise Squared Error (**MPSE**), which is a pairwise ranking loss that measures the differences between all possible pairs of corresponding temporal attention responses estimated (\hat{a}_i, \hat{a}_j) and **GT** fixation-based responses (a_i, a_j) in each mini-batch.

6.5 WHERE WE ARE: TOWARDS GUIDING ANOMALY DETECTION

As it was introduced in Chapter 5, visual attention in the temporal domain can be understood as an information filtering mechanism which allows to select candidate time segments to contain anomalous events. This constitutes a very interesting application of visual attention, which would substantially decrease the cognitive effort made by **CCTV** operators in video monitoring.

As a first approximation to the use of the **ST-T-ATTEN** proposed for guiding and reducing the computational cost of an anomaly detection task, we can consider the same binary classification problem posed in Section 5.5.2 for the context defined by the BOSS [19] database, but now measuring the existing correlation between anomalies e_t and the estimated temporal attention response \hat{a}_t by the best **ST-T-ATTEN** configuration determined in the previous section. Moreover, for the sake of comparison, we provide the results obtained by two simple baseline methods: AGGREGATION SUM and AGGREGATION MAX. They also estimate attention in the temporal domain from **VAMs** extracted by **ST-ATTEN**, but aggregating the spatial dimension using a sum (AGGREGATION SUM) or a

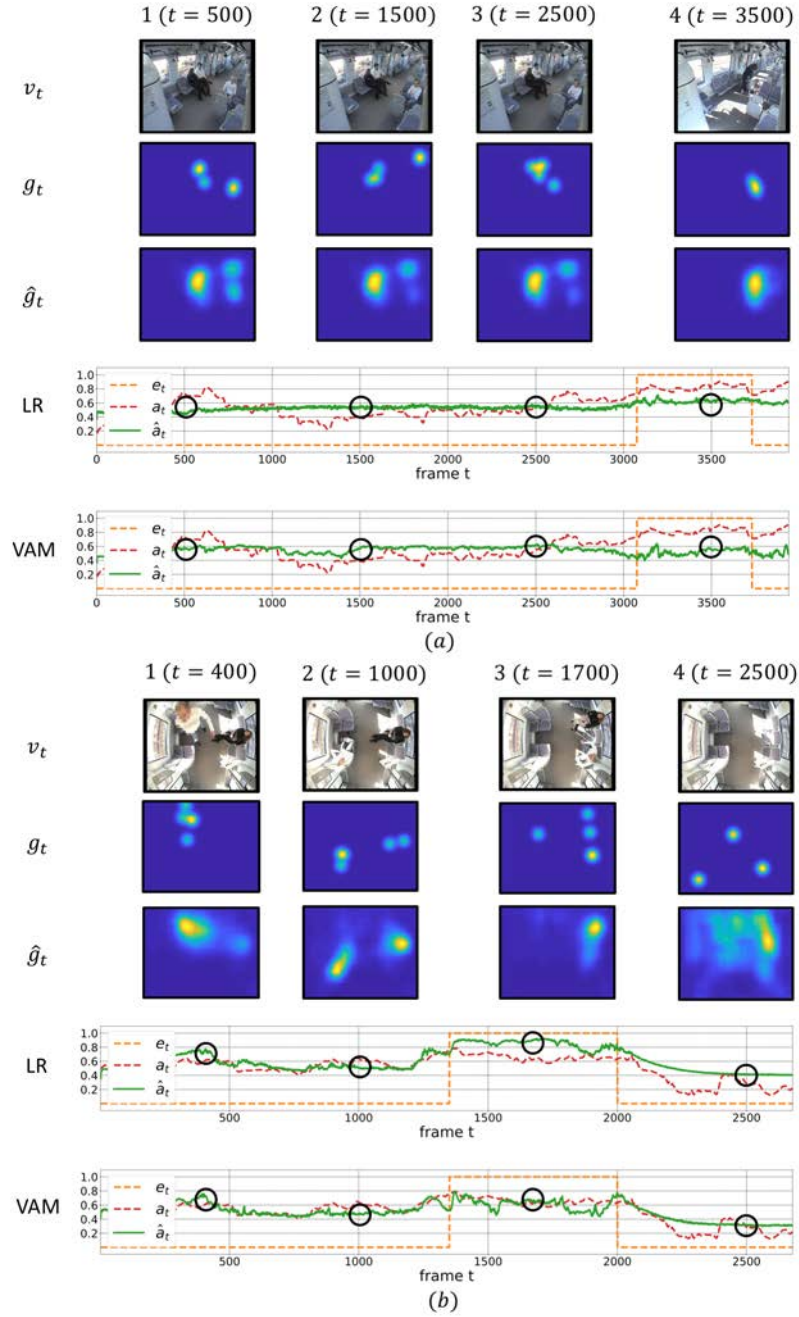


Figure 6.3: Visual attention in the temporal domain \hat{a}_t estimated by ST-T-ATTEN in two video-surveillance sequences taken from BOSS [19] database. On the first two rows, some video frames over time v_t are shown, together with their associated GT fixations map g_t . Then, the spatio-temporal VAMs \hat{g}_t estimated by the ST-ATTEN module of the system are displayed. Finally, the temporal attention response \hat{a}_t estimated by the two LSTM-based T-ATTEN proposed (LR and VAM) is represented. For the sake of comparison, the GT temporal attention response a_t and the anomaly detection signal e_t are also plotted. (a) The sequence shows two men fighting. 1,2. Three passengers are sitting on the train. 3. Two passengers start arguing. 4. Two men fight, and a woman tries to stop them. (b) The sequence shows a woman harassment scene. 1. Man sits on the train. 2. Man starts talking to the woman. 3. Man approaches the woman. 4. Woman and man left the train.

Table 6.3: Results obtained by the proposed [ST-T-ATTEN](#) and other comparison methods considered as filtering mechanisms for guiding anomaly detection in the video surveillance scenario defined by the BOSS [\[19\]](#) database.

Model	AUC
\hat{a}_t	0.703
AGGREGATION MAX	0.613
AGGREGATION SUM	0.543
a_t	0.876

maximum (AGGREGATION MAX) operator. Table [6.3](#) summarizes the results obtained by [ST-T-ATTEN](#) and the methods considered for comparison. After conducting the experiments, we achieve an $AUC = 0.703$, which notably outperforms the two baseline approaches. Besides, despite the AUC obtained by [ST-T-ATTEN](#) is lower to the one achieved by the [GT](#) temporal attention response a_t ($AUC = 0.876$), it constitutes a very promising result, given the complexity of the task to address. Furthermore, it is worth recalling that there is still room for improvement in [ST-T-ATTEN](#), as discussed in the previous section.

In order to illustrate two potential applications of the [ST-T-ATTEN](#) architecture in a video surveillance scenario, we provide two more sequences taken from BOSS [\[19\]](#) database in Figure [6.3](#). Either spatio-temporal [VAMs](#) computed by [ST-ATTEN](#) module or temporal attention responses estimated by [T-ATTEN](#) are shown. Moreover, in order to bring visual attention closer to the context of anomaly detection, we plot again the anomaly detection signals e_t , as we did in Figures [5.13](#) and [5.15](#) to validate the fundamental hypothesis of our approach.

On the one hand, let us imagine a situation where the two video sequences in Figure [6.3](#) are being shown in real-time in two screens belonging to a large array of camera views. In such a situation, the attention estimated in the temporal domain could be applied to select or highlight the most outstanding screen from the monitoring array at every time, thus driving operator's attention to scenes that potentially show anomalies or suspicious events.

Furthermore, the estimated temporal attention response could be also applied in off-line tasks which imply reviewing many hours of surveillance recordings, e.g. a surveillance operator inspecting hours of video footage searching for a particular event or anomaly. In this case, our system might reduce the amount of information to be processed by the operator. As a first demonstration of this application, let us note that identifying a unique frame for each

anomaly would be sufficient for providing a **CCTV** operator with temporal indicators to review footage faster and more efficiently. Considering this fact and given the temporal attention responses estimated by **ST-T-ATTEN** on the BOSS [19] database, it would be necessary only to retrieve approximately the 16% of the frames of the complete database in order to locate over time the 95% of the existing anomalous events.

Therefore, we can conclude that, with some adjustments, the **ST-T-ATTEN** proposed might be able to estimate, given one or multiple camera views, spatio-temporal visual attention maps and temporal attention responses at the same time. Our system would thus provide **CCTV** operators a complete experience of visual attention by highlighting the most conspicuous locations in a given scene and, besides, the most relevant time segments, according not only to previous events in the scene, but also to events happening in different camera views at the same time.

6.6 CONCLUSIONS

In Chapters 5 and 6, we have presented a deep network architecture for visual attention modeling, which goes from **CNN**-based spatio-temporal visual attention prediction to **LSTM**-based attention estimation in the temporal domain. Our model is fundamentally supported by the assumption that a measurement of task-driven visual attention in the temporal domain can be drawn from the dispersion of gaze locations recorded from several subjects. Indeed, the temporal level of attention of observers constitutes an important clue to detect suspicious events or anomalous situations in crowded and complex scenarios. However, it should be borne in mind that, similarly to spatio-temporal visual attention, attention in the temporal domain has to be considered as an early filtering mechanism, which allows to select time segments candidate to contain events of special importance, and therefore reduce the complexity of subsequent anomaly detection systems or to drive the attention of human operators to particular cameras in complex multi-camera **CCTV** systems.

Experimental results have determined the optimal configuration for the **ST-T-ATTEN** proposed. First, it is composed by a **CONV ST-ATTEN** stage, which successfully fuses the information provided by three different visual feature maps at its input: RGB-based, motion and objectness. Then, the **ST-ATTEN** module connects with a **LSTM**-based **T-ATTEN** architecture, which models attention in the temporal domain. After evaluating the system on the BOSS [19] database, either the **T-ATTEN** fed with **VAMs** obtained at the output of the **ST-ATTEN** decoder or the one that receives as input the latent representations extracted by the encoder network have resulted in

similar performances in terms of the PCC score. However, there is still room for the analysis and the improvement of our approach. To that end, future work will address the annotation of a large-scale video surveillance dataset with eye fixations to draw some better conclusions about the behavior of the system, with the ultimate aim of demonstrating its usefulness for guiding anomaly detection in a video surveillance application.

CONCLUSIONS AND FUTURE LINES OF RESEARCH

7.1 CONCLUSIONS

In this thesis we have proposed two hierarchical frameworks for visual attention modeling in video sequences. Visual attention can be modeled in two different domains, spatial and temporal, which leads to three types of computational models: spatial, spatio-temporal and temporal. First, spatial models highlight locations of particular interest in a frame by frame basis. Second, modeling attention in the temporal domain allows either to update spatial attention based on previously selected locations (spatio-temporal) or to select time segments of special importance in a video (temporal).

In Chapter 3, we have presented our first approach, which is called visual Attention TOpic Model ([ATOM](#)) [126]. Our proposal involves a hierarchical generative probabilistic model for spatio-temporal visual attention prediction and understanding. The definition of the system proposed is generic and independent of the application scenario. Moreover, it is founded on the most outstanding psychological studies about attention [10, 11], which hold that attention guidance is not based directly on the information provided by early visual processes but on a contextual representation arisen from them.

Relying on the well-known Latent Dirichlet Allocation ([LDA](#)) [12] and its supervised extensions [13, 14], [ATOM](#) defines task- or context-driven visual attention in video as a mixture of several sub-tasks which, in turn, can be represented as a combination of low-, mid- and high-level spatio-temporal features obtained from video frames. Therefore, given a video frame, the algorithm receives a set of visual feature maps (color, intensity, motion, object-based, etc.) as input. Then, an intermediate level of latent sub-tasks between feature extraction and visual attention modeling is introduced. Finally, latent sub-tasks are aligned to the information drawn from human fixations by means of a categorical variable

response, which is generated by a logistic regression model over the sub-task proportions.

In Chapter 4, we have demonstrated the ability of [ATOM](#) to successfully learn hierarchical representations of visual attention specifically adapted to diverse contexts (outdoors, video games, sports, TV news, etc.), on the basis of a wide set of features. For that purpose, we have made use of the well-known large-scale CRCNS-ORIG [15] and DIEM [16] databases. Experiments have shown the advantage of our comprehensive guiding representations based on handcrafted features to understand how visual attention works in different scenarios. In addition, modeling simple eye-catching elements, such as faces or text, through spatial discrete distributions, as well as considering object-based representations learned by recently adopted [CNNs](#), our proposal significantly outperforms quite a few competent methods in the literature when estimating visual attention.

In Chapter 5, we have introduced our second proposal, which is named Spatio-Temporal to Temporal visual ATtention NETwork ([ST-T-ATTEN](#)). This second approach takes a step further and goes from spatio-temporal visual attention estimation to attention estimation in the temporal domain. The model is fundamentally supported by the assumption that a measurement of task-driven visual attention in the temporal domain can be drawn from the dispersion of fixation locations recorded from several observers. First, to demonstrate this hypothesis, we have measured the existing correlation between eye fixation sequences of different viewers when an important or anomalous event happens on the BOSS [19] database. Although this temporal level of attention constitutes a useful clue to detect important events in crowded and complex scenarios, attention in the temporal domain should always be considered as an early filtering mechanism, which selects candidate time segments to contain suspicious events, and therefore reduces the later processing devoted to the anomaly detection. Based on this hypothesis, we have developed [ST-T-ATTEN](#), which attempts to model attention in the temporal domain from estimations of spatio-temporal visual attention.

Inspired by the recent success of Convolutional Neural Networks ([CNNs](#)) for learning deep hierarchical representations and [LSTM](#) units for time series forecasting, the proposed [ST-T-ATTEN](#) is composed of two stages. The first stage, which is denoted as Spatio-Temporal visual ATtention NETwork ([ST-ATTEN](#)), consists of a Convolutional Encoder Decoder ([CED](#)) network that receives at its input three high-level feature maps for visual attention guidance (RGB-based, motion and objectness), all of them computed by deep [CNNs](#). Then, through an encoding-decoding architecture, the network concurrently estimates spatio-temporal [VAMs](#) and extracts

latent representations of visual attention. We have proposed two configurations for this module of the system. They differ in the outer layers of the encoder and the decoder, which are **CONV** in the first approach and **CONV-LSTM** in the second one.

The second stage of **ST-T-ATTEN**, which is called Temporal ATtention NETwork (**T-ATTEN**), involves a **LSTM**-based architecture that estimates, for each frame in a video sequence, a temporal attention response. We have also distinguished between two versions of **T-ATTEN**, depending on the input variable: either the **VAM** at the output of the **ST-ATTEN** or the latent representations generated by the encoder.

In Chapter 6, the proposed **ST-T-ATTEN** architecture has been evaluated in a video surveillance scenario defined by the BOSS [19] database, which contains video sequences recorded in a railway transport context, with different types of suspicious or anomalous events (several women harassment, a cell phone theft, a passengers fight, etc.). The main purpose of our experiments has been to assess various architectures of our proposal. Experiments have concluded that the best performing architecture is composed by a **CONV ST-ATTEN** stage, which successfully fuses the information provided by the three input feature maps. Then, either the **T-ATTEN** fed with the **VAMs** obtained at the output of the **ST-ATTEN** decoder or the latent representations extracted by its associated encoder have resulted in similar performances in terms of the **PCC** score. However, there is still room for the improvement of our system, as discussed in the experimental section.

Finally, we have also discussed two potential end-user applications for our proposal. On the one hand, given a video surveillance scenario, the temporal attention response could be applied to select in real-time the most outstanding screens from the monitoring array, thus driving operator's attention to scenes that potentially show anomalies or suspicious events. On the other hand, the estimated response could be also applied in off-line tasks which imply reviewing many hours of surveillance recordings, reducing the information to be processed by the operator. With some adjustments, our system might be able to provide **CCTV** operators a complete experience of visual attention, not only highlighting the most conspicuous locations in a scene, but also selecting the most relevant time segments, according to both previous events in the scene and events happening in different camera views at the same time.

7.2 FUTURE LINES OF RESEARCH

Lastly, we conclude this thesis by identifying and discussing potential future lines of research related to our contributions.

At this point, there is no doubt about the great benefits of visual attention modeling in the framework of Artificial Intelligence (AI), nor about the infinite possibilities that such an abstract concept opens for the processing and understanding of this big data world. Despite the wide variety of computational models of visual attention existing in the literature, much remains to be done, not only to meet a system that automatically addresses this cognitive function, but also to understand how HVS carries out this optimization process.

Turning to the two popular representation learning paradigms introduced at the beginning of this thesis, Deep Learning (DL) and Probabilistic Graphical Models (PGM), our contributions have shown the importance of either the task of *seeing*, performed by DL representations, or the ability of *thinking*, characteristic of PGM, for visual attention modeling and understanding.

First, it is important to achieve good representations of the world that surrounds us for attention guidance, and it is here where DNNs architectures and, in particular, CNNs, play an essential role in machine perception. In addition, given that visual attention involves not only one, but several complex tasks, it is paramount to understand how computational visual attention deals with the hierarchical representations provided by DNNs, through probabilistic methods that explain relationships between the observed variables. This direction, recently set by Bayesian Deep Learning (BDL) [20], is the one that we plan to follow in our future research, paying special attention to BDL for topic models, which constitutes a revision to probabilistic Latent Topic Models (LTMs) [12–14], on the basis of which our hierarchical ATOM framework for visual attention understanding has been built. Discovering sub-tasks, not only over space but also along time, will allow establishing relationships between recognized concepts in one or multiple video sequences, both in the same scene or in different ones.

Secondly, in the latter part of the thesis, we have demonstrated the major advantages of modeling visual attention in the temporal domain, selecting video segments of special importance, which subsequently help to reduce the computational burden of subsequent end-user applications. Visual attention has been barely tackled from this perspective in the literature up to date, in spite of its usefulness for the processing and analysis of vast amounts of visual information in applications such as anomaly detection.

One interesting research line we have not covered in this thesis is the interpretation of eye movement sequences, establishing relationships between the content of fixated locations. This would allow to develop more comprehensible and valuable systems for estimating the variation of visual attention over time. Reinforcement learning methods seem a promising way of addressing this challenge [21].

Finally, we are highly motivated to model spatio-temporal visual attention, as well as attention in the temporal domain, given multiple video sequences played at the same time, with the aim of assisting experts in crowded and complex scenarios. For that purpose, we will soon proceed to annotate large-scale video surveillance databases, such as VIRAT [22] or UCF-Crime [23], with human fixations, which will serve for a further analysis and the improvement of the deep [ST-T-ATTEN](#) architecture proposed.

A

DERIVATION OF THE FORMULAS FOR THE ATOM

In this appendix, we provide the derivation of the formulas for the visual Attention TOpic Model (ATOM) presented in Chapter 3.

A.1 EXPANSION OF THE LOWER BOUND

As introduced in section 3.5.3, the optimization of the probabilistic model for visual attention understanding and prediction proposed in the thesis implies maximizing the Evidence Lower Bound (ELBO) over the log-likelihood of all the frames in a corpus of videos. In particular, using Jensen’s inequality, the ELBO of the log-likelihood of a frame with N spatial locations can be expressed as:

$$\begin{aligned} \log p(f_{1:N,1:L}, g_{1:N} | \alpha, \Gamma_{1:K,1:L}, \eta) &\geq E_q[\log p(\theta | \alpha)] \\ &+ \sum_{n=1}^N E_q[\log p(\mathbf{z}_n | \theta)] + \sum_{n=1}^N E_q[\log p(f_{n,1:L} | \mathbf{z}_n, \Gamma_{1:K,1:L})] \\ &+ \sum_{n=1}^N E_q[\log p(g_n | \mathbf{z}_n, \eta)] + H(q) \end{aligned} \quad (\text{A.1})$$

where L is the number of visual descriptors computed as input for the models, K is the number of sub-tasks or topics inferred, $E_q[\cdot]$ is the expectation over the variational distribution q and $H(\cdot)$ is the entropy of the variational distribution.

The first two terms in the ELBO and the entropy of the variational distribution are identical to the corresponding terms in the ELBO for unsupervised LDA [12]:

$$\begin{aligned}
E_q[\log p(\theta|\alpha)] &= \log \Gamma \left(\sum_{k=1}^K \alpha_k \right) - \sum_{k=1}^K \log \Gamma(\alpha_k) \\
&\quad + \sum_{k=1}^K (\alpha_k - 1) E_q[\log \theta_k]
\end{aligned} \tag{A.2}$$

$$E_q[\log p(\mathbf{z}_n|\theta)] = \sum_{k=1}^K \phi_{nk} E_q[\log \theta_k] \tag{A.3}$$

$$\begin{aligned}
H(q) &= - \sum_{n=1}^N \sum_{k=1}^K \phi_{nk} \log \phi_{nk} - \log \Gamma \left(\sum_{k=1}^K \gamma_k \right) \\
&\quad + \sum_{k=1}^K \log \Gamma(\gamma_k) - \sum_{k=1}^K (\gamma_k - 1) E_q[\log \theta_k],
\end{aligned} \tag{A.4}$$

where the expectation of the log of the multinomial random variable θ_k is:

$$E_q[\log \theta_k] = \Psi(\gamma_k) - \Psi \left(\sum_{j=1}^K \gamma_j \right), \tag{A.5}$$

being $\Psi(\cdot)$ the digamma function.

The third and fourth terms are derived in the following subsections.

A.1.1 Lower bound of the local appearance model

The third term is the expected log probability of the data given the related topic model parameters. We assume conditional independence among features. In the following paragraphs, we derive this expression for the different distributions considered.

If the feature map f_{nl} is modeled with a univariate *Gaussian distribution* $\Gamma_{1:K,l} \sim \{\mu_{1:K,l}, \sigma_{1:K,l}^2\}$, such as for basic and novelty spatio-temporal features or CNN-based features, the equation for this term is:

$$\begin{aligned}
E_q[\log p(f_{nl}|\mathbf{z}_n, \Gamma_{1:K,l})] &= - \sum_{k=1}^K \phi_{nk} \log(\sigma_{kl} \sqrt{2\pi}) \\
&\quad - \sum_{k=1}^K \phi_{nk} \frac{(f_{nl} - \mu_{kl})^2}{2\sigma_{kl}^2}
\end{aligned} \tag{A.6}$$

where ϕ_{nk} is the probability that the location n has been drawn by the topic k .

In the case of camera motion features, the distribution is a multivariate Gaussian $p(\mathbf{x}_n|\mathbf{z}_n, \mu_k, \Sigma_k)$ with $\mu_k = \mathbf{c}_k \odot \mathbf{u}$, being \mathbf{c}_k a parameter to be estimated and $\mathbf{u} = (u, v)$ the camera motion vector. However, due to the diagonal nature of the covariance matrix Σ_k we

can decompose it into two independent univariate Gaussian distributions and apply the previous expression:

$$E_q[\log p(\mathbf{x}_n | \mathbf{z}_n, \mathbf{c}_k)] = - \sum_{k=1}^K \phi_{nk} \log(\sqrt{2\pi\Sigma_k}) - \sum_{k=1}^K \phi_{nk} \frac{(\mathbf{x}_n - \mathbf{c}_k \mathbf{u})^T (\mathbf{x}_n - \mathbf{c}_k \mathbf{u})}{2\Sigma_k}. \quad (\text{A.7})$$

where $\mathbf{x}_n = (x_n, y_n)$ is a vector with the spatial coordinates of the location n . As explained in section 3.3.2, Σ_k , which controls the spatial extent of the Gaussian distribution, has been empirically set to $\Sigma_k = \text{diag}(0.25)$ in order to cover a sufficiently wide area in the scene.

In contrast, if the feature is modeled as a *discrete probability distribution* over cells r in a grid, as happens for objects-based features, the expression is:

$$E_q[\log p(r_n | \mathbf{z}_n, \beta_{l z_n})] = \sum_{k=1}^K \phi_{nk} \log(\beta_{kl r_n}), \quad (\text{A.8})$$

where r_n stands for the region in the non-uniform grid defined for the object l that contains the location n , and $\beta_{kl r_n}$ is the value of the of the discrete distribution in that region for the object l and the topic k .

A.1.2 Lower bound of the visual attention response

The fourth term includes the visual attention response variable g_n , which is generated from a Bernoulli distribution, i.e.,

$$p(g_n | \pi_n) = (\pi_n)^{g_n} (1 - \pi_n)^{(1-g_n)}, \quad (\text{A.9})$$

where π is a logistic regression model based on a weighted empirical average of the Dirichlet realization $\eta^T \mathbf{z}_n$, being η the parameter vector that models attention based on the selected topic \mathbf{z}_n :

$$p(g_n | \mathbf{z}_n, \eta) = \frac{\exp(g_n \eta^T \mathbf{z}_n)}{1 + \exp(\eta^T \mathbf{z}_n)}. \quad (\text{A.10})$$

Thus, the Bernoulli distribution is as follows:

$$p(g_n | \mathbf{z}_n, \eta) \sim Be \left(\frac{\exp(g_n \eta^T \mathbf{z}_n)}{1 + \exp(\eta^T \mathbf{z}_n)} \right). \quad (\text{A.11})$$

According to [158], the logistic function in Eq. A.10 can be symmetrized as follows:

$$p(g_n | \mathbf{z}_n, \eta) = \frac{\exp((g_n - \frac{1}{2}) \eta^T \mathbf{z}_n)}{\exp(\frac{\eta^T \mathbf{z}_n}{2}) + \exp(\frac{-\eta^T \mathbf{z}_n}{2})}. \quad (\text{A.12})$$

Then, the expected log probability of the response variable given the topic assignments is expressed as:

$$E_q[\log p(g_n | \mathbf{z}_n, \eta)] = E_q \left[\left(g_n - \frac{1}{2} \right) \eta^T \mathbf{z}_n \right] - E_q \left[\log \left(\exp \left(\frac{\eta^T \mathbf{z}_n}{2} \right) + \exp \left(\frac{-\eta^T \mathbf{z}_n}{2} \right) \right) \right] \quad (\text{A.13})$$

By taking second derivatives, it can be noticed that the second term above, which can be denoted as $E_q[f(\eta^T \mathbf{z}_n)]$, is a convex function in the variable $\eta^{T^2} \mathbf{z}_n^2 = (\eta^T \odot \eta^T)(\mathbf{z}_n \odot \mathbf{z}_n)$, so we can bound it by using the lower bound for the logistic function

$$f(\eta^T \mathbf{z}_n) \geq f(\xi_n) + \frac{\partial f(\xi_n)}{\partial \xi_n^2} (\eta^{T^2} \mathbf{z}_n^2 - \xi_n^2), \quad (\text{A.14})$$

which is the first order Taylor expansion in the variable $\eta^{T^2} \mathbf{z}_n^2$:

$$\begin{aligned} & \log \left(\exp \left(\frac{\eta^T \mathbf{z}_n}{2} \right) + \exp \left(\frac{-\eta^T \mathbf{z}_n}{2} \right) \right) \\ & \geq -\frac{\xi_n}{2} - \log(1 + \exp(-\xi_n)) \\ & \quad - \frac{1}{4\xi_n} \tanh \left(\frac{\xi_n}{2} \right) E_q \left[\eta^{T^2} \mathbf{z}_n^2 - \xi_n^2 \right] \\ & \approx -\frac{\xi_n}{2} - \log(1 + \exp(-\xi_n)) \\ & \quad - \frac{1}{4\xi_n} \tanh \left(\frac{\xi_n}{2} \right) (\eta^{T^2} \phi_n - \xi_n^2), \end{aligned} \quad (\text{A.15})$$

where ϕ_n is the vector of topic proportions ϕ_{nk} in the location n and ξ_n is an additional variational parameter associated to each point n .

It should be noted that, during variational inference, we work on expected values. This means that the indexing variable \mathbf{z}_n is replaced by the variational ϕ_n , which now contains the expected values of the topic assignments given a location n . Therefore, since ϕ_n is a vector with real values (the topic proportions for that sampled location), in practice each location n is in turn modeled as the mixture of sub-tasks that best explains its visual appearance.

A.2 DERIVATION OF THE FORMULAS FOR THE VARIATIONAL PARAMETERS

This section includes the complete derivation of the update equation of the variational multinomial ϕ , which is computed in the E-step of the inference process.

First, we begin with the lower bound that depends on ϕ , incorporating a Lagrange parameter λ to ensure that $\sum_{k=1}^K \phi_{nk} = 1$:

$$\begin{aligned}
ELBO_\phi &= \sum_{n=1}^N \sum_{k=1}^K \phi_{nk} E_q[\log \theta_k] + \sum_{n=1}^N E_q[\log p(f_{n,1:L} | \mathbf{z}_n, \Gamma_{1:K,1:L})] \\
&+ \sum_{n=1}^N E_q[\log p(g_n | \mathbf{z}_n, \eta)] + H(q) \\
&= \sum_{n=1}^N \sum_{k=1}^K \phi_{nk} \left(\Psi(\gamma_k) - \Psi \left(\sum_{j=1}^K \gamma_j \right) \right) \\
&- \sum_{n=1}^N \sum_{k=1}^K \sum_{l=1}^{L_C} \phi_{nk} \left(\log(\sigma_{kl} \sqrt{2\pi}) + \frac{(f_{nl} - \mu_{kl})^2}{2\sigma_{kl}^2} \right) \\
&+ \sum_{n=1}^N \sum_{k=1}^K \sum_{l=1}^{L_D} \phi_{nk} \log(\beta_{klr_n}) \\
&+ \sum_{n=1}^N \sum_{k=1}^K \left(\left(g_n - \frac{1}{2} \right) \eta_k \phi_{nk} - \frac{1}{4\tilde{\xi}_{nk}} \tanh \left(\frac{\xi_{nk}}{2} \right) (\eta_k^2 \phi_{nk} - \tilde{\xi}_{nk}^2) \right) \\
&- \sum_{n=1}^N \sum_{k=1}^K \phi_{nk} \log \phi_{nk} - \log \Gamma \left(\sum_{k=1}^K \gamma_k \right) \\
&+ \sum_{k=1}^K \log \Gamma(\gamma_k) - \sum_{k=1}^K (\gamma_k - 1) E_q[\log \theta_k] \\
&+ \sum_{n=1}^N \lambda_n \left(\sum_{k=1}^K \phi_{nk} - 1 \right), \tag{A.16}
\end{aligned}$$

being L_C and L_D the number of continuous (Gaussian) and discrete features, respectively, and $L = L_C + L_D$ the total number of features.

If we take the derivative of the [ELBO](#) with respect to ϕ_{nk} :

$$\begin{aligned}
\frac{\partial ELBO_{\phi_{nk}}}{\partial \phi_{nk}} &= \Psi(\gamma_k) - \Psi \left(\sum_{j=1}^K \gamma_j \right) \\
&- \sum_{l=1}^{L_C} \left(\log(\sigma_{kl} \sqrt{2\pi}) + \frac{(f_{nl} - \mu_{kl})^2}{2\sigma_{kl}^2} \right) + \sum_{l=1}^{L_D} \log(\beta_{klr_n}) \\
&+ \left(g_n - \frac{1}{2} \right) \eta_k - \frac{1}{4\tilde{\xi}_{nk}} \tanh \left(\frac{\xi_{nk}}{2} \right) \eta_k^2 - \log \phi_{nk} - 1 + \lambda_n \tag{A.17}
\end{aligned}$$

and set it to zero, we obtain the equation for updating the multinomial parameter:

$$\begin{aligned}
\phi_{nk} &\propto \frac{\prod_{l=1}^{L_D} \beta_{klr_n}}{\prod_{l=1}^{L_C} \sigma_{kl}} \exp \left[\Psi(\gamma_k) - \Psi \left(\sum_{j=1}^K \gamma_j \right) + \right. \\
&\quad \left(g_n - \frac{1}{2} \right) \eta_k - \frac{1}{4\tilde{\xi}_k} \tanh \left(\frac{\xi_k}{2} \right) \eta_k^2 - \\
&\quad \left. \sum_{l=1}^{L_C} \frac{(f_{nl} - \mu_{kl})^2}{2\sigma_{kl}^2} \right]. \tag{A.18}
\end{aligned}$$

In addition, it should be noted that the equations corresponding to the variational Dirichlet γ and the Dirichlet parameters α are not

included here, because they are identical to those in the original LDA [12].

A.3 DERIVATION OF THE FORMULAS FOR THE MODEL PARAMETERS

This section includes the complete derivation of the update equations for the model parameters computed in the M-step of the inference process, given a corpus of T video frames, each one with N_t spatial locations.

First, parameters μ_{kl} and σ_{kl}^2 are computed for each Gaussian feature l and topic k .

- The ELBO that depends on μ_{kl} is:

$$ELBO_{\mu_{kl}} = - \sum_{t=1}^T \sum_{n=1}^{N_t} \sum_{k=1}^K \phi_{nk} \frac{(f_{tnl} - \mu_{kl})^2}{2\sigma_{kl}^2}. \quad (\text{A.19})$$

Computing its derivative with respect to μ_{kl} gives:

$$\frac{\partial ELBO_{\mu_{kl}}}{\partial \mu_{kl}} = \sum_{t=1}^T \sum_{n=1}^{N_t} \phi_{tnk} \frac{(f_{tnl} - \mu_{kl})}{\sigma_{kl}^2} \quad (\text{A.20})$$

Setting it to zero, we obtain the update equation:

$$\mu_{kl} = \frac{1}{\Delta_{kl}} \sum_{t=1}^T \sum_{n=1}^{N_t} \phi_{tnk} f_{tnl} \quad (\text{A.21})$$

- The ELBO that depends on σ_{kl}^2 is:

$$\begin{aligned} ELBO_{\sigma_{kl}^2} = & - \sum_{t=1}^T \sum_{n=1}^{N_t} \sum_{k=1}^K \phi_{nk} \log(\sigma_{kl} \sqrt{2\pi}) \\ & - \sum_{t=1}^T \sum_{n=1}^{N_t} \sum_{k=1}^K \phi_{nk} \frac{(f_{tnl} - \mu_{kl})^2}{2\sigma_{kl}^2}. \end{aligned} \quad (\text{A.22})$$

Computing its derivative with respect to σ_{kl}^2 gives:

$$\begin{aligned} \frac{\partial ELBO_{\sigma_{kl}^2}}{\partial \sigma_{kl}^2} = & - \sum_{t=1}^T \sum_{n=1}^{N_t} \frac{\phi_{nk}}{\sigma_{kl}} \\ & + \sum_{t=1}^T \sum_{n=1}^{N_t} \frac{\phi_{nk}}{\sigma_{kl}} \frac{(f_{tnl} - \mu_{kl})^2}{\sigma_{kl}^2} \end{aligned} \quad (\text{A.23})$$

Setting it to zero, we obtain the update equation:

$$\sigma_{kl}^2 = \frac{1}{\Delta_{kl}} \sum_{t=1}^T \sum_{n=1}^{N_t} \phi_{tnk} (f_{tnl} - \mu_{kl})^2 \quad (\text{A.24})$$

For both μ_{kl} and σ_{kl}^2 , $\Delta_{kl} = \sum_{t=1}^T \sum_{n=1}^{N_t} \phi_{tnk}$ is the normalization factor.

In the case of camera motion, as mentioned above, the parameter is the vector $\mathbf{c}_k = (c_{kx}, c_{ky})$ that multiplies the camera motion vector $\mathbf{u}_t = (u_t, v_t)$ to determine the mean of the Gaussian distribution. The [ELBO](#) that depends on this parameter is:

$$ELBO_{\mathbf{c}_k} = - \sum_{t=1}^T \sum_{n=1}^{N_t} \sum_{k=1}^K \phi_{tnk} \frac{(\mathbf{x}_{tn} - \mathbf{c}_k \mathbf{u}_t)^T (\mathbf{x}_{tn} - \mathbf{c}_k \mathbf{u}_t)}{2\Sigma_k}, \quad (\text{A.25})$$

where $\mathbf{x}_{tn} = (x_{tn}, y_{tn})$ stands for the spatial coordinates vector of the location n in frame t . By computing its derivative with respect to \mathbf{c}_k :

$$\frac{\partial ELBO_{\mathbf{c}_k}}{\partial \mathbf{c}_k} = \sum_{t=1}^T \sum_{n=1}^{N_t} \phi_{tnk} \frac{(\mathbf{u}_t \mathbf{x}_{tn} - \mathbf{u}_t^2 \mathbf{c}_k)}{\Sigma_k} \quad (\text{A.26})$$

and setting it to zero, we obtain the following update equation:

$$\mathbf{c}_k = \frac{\sum_{t=1}^T \sum_{n=1}^{N_t} \phi_{tnk} \mathbf{u}_t \mathbf{x}_{tn}}{\sum_{t=1}^T \sum_{n=1}^{N_t} \phi_{tnk} \mathbf{u}_t^2}. \quad (\text{A.27})$$

Finally, for the case of object-based discrete features, the [ELBO](#) that depends on the probabilities β_{klr} of the regions r defined on the object-detector l and for every topic k is:

$$ELBO_{\beta_{klr}} = \sum_{t=1}^T \sum_{n=1}^{N_t} \sum_{k=1}^K \phi_{tnk} \log(\beta_{klr_n}) + \sum_{k=1}^K \lambda_{kl} \left(\sum_{r=1}^R \beta_{klr_n} - 1 \right). \quad (\text{A.28})$$

where we have added the Lagrange multipliers λ_{kl} . By computing its derivative with respect to β_{klr} :

$$\frac{\partial ELBO_{\beta_{klr}}}{\partial \beta_{klr}} = \sum_{t=1}^T \sum_{n=1}^{N_t} \frac{\phi_{tnk}}{\beta_{klr}} + \lambda_{kl}. \quad (\text{A.29})$$

Setting this derivative to zero gives the following update equation:

$$\beta_{klr} \propto \sum_{t=1}^T \sum_{n=1}^{N_t} \phi_{tnk} 1[r_{nl} = r] \quad (\text{A.30})$$

where $1[r_{nl} = r]$ means that we have a 1 just in case the region of the point n for the detector l is r (otherwise we have a zero).

A.4 DERIVATION OF THE FORMULAS FOR THE PARAMETERS OF THE LOGISTIC REGRESSION MODEL

This section includes the complete derivation of the update equations for the parameters of the logistic regression model proposed to estimate visual attention over the underlying topics obtained, either in the E-step or in the M-step of the inference process.

- In the E-step, the variational parameter ξ has to be updated. The [ELBO](#) that depends on ξ_{nk} is:

$$\begin{aligned} ELBO_{\xi_{nk}} = & - \sum_{t=1}^T \sum_{n=1}^{N_t} \sum_{k=1}^K \left(\frac{\xi_{tnk}}{2} + \log(1 + \exp(-\xi_{tnk})) \right) \\ & + \sum_{t=1}^T \sum_{n=1}^{N_t} \sum_{k=1}^K \frac{1}{4\xi_{tnk}} \tanh\left(\frac{\xi_{tnk}}{2}\right) (\eta_k^2 \phi_{tnk} - \xi_{tnk}^2), \end{aligned} \quad (\text{A.31})$$

which corresponds to the lower bound in Eq. [A.15](#). This lower bound is exact if $\xi_{nk}^2 = \eta_k^2 \phi_{nk}$. Consequently, the update equation for this parameter is:

$$\xi_{nk} = \eta_k \phi_{nk}. \quad (\text{A.32})$$

- In the M-step, we update the parameter η , attending to the lower bound that depends on it:

$$\begin{aligned} ELBO_{\eta_k} = & \sum_{t=1}^T \sum_{n=1}^{N_t} \sum_{k=1}^K \left(g_{tn} - \frac{1}{2} \right) \eta_k \phi_{tnk} \\ & - \sum_{t=1}^T \sum_{n=1}^{N_t} \sum_{k=1}^K \frac{1}{4\xi_{tnk}} \tanh\left(\frac{\xi_{tnk}}{2}\right) (\eta_k^2 \phi_{tnk} - \xi_{tnk}^2). \end{aligned} \quad (\text{A.33})$$

Computing the derivative with respect to η_k :

$$\begin{aligned} \frac{\partial ELBO_{\eta_k}}{\partial \eta_k} = & \sum_{t=1}^T \sum_{n=1}^{N_t} \left(g_{tn} - \frac{1}{2} \right) \phi_{tnk} \\ & - \sum_{t=1}^T \sum_{n=1}^{N_t} \frac{1}{2\xi_{tnk}} \tanh\left(\frac{\xi_{tnk}}{2}\right) \eta_k \phi_{tnk}. \end{aligned} \quad (\text{A.34})$$

Setting it to zero, we obtain the update equation:

$$\eta_k = \frac{2 \sum_{t=1}^T \sum_{n=1}^{N_t} \phi_{tnk} (g_{tn} - \frac{1}{2})}{\sum_{t=1}^T \sum_{n=1}^{N_t} \frac{\phi_{tnk}}{\xi_{tnk}} \tanh(\frac{\xi_{tnk}}{2})}. \quad (\text{A.35})$$

B

EYE-TRACKING DATABASES USED IN THE THESIS

This appendix summarizes the databases used in Chapters 4 and 6 from this thesis, with the aim of providing the list of videos included in each of the categories established for our experiments.

It should be noted that SALICON [114] database has also been considered for RGB-based feature maps extraction in Chapter 6, where it has been described. This database is not covered in this appendix because it was not necessary neither to divide its images into categories nor to annotate it with more observers' fixations, so it has been used as in other related works in the state-of-the-art.

Last but not least, the three most significant attraction (AT) and inhibition (IT) sub-tasks determined by the spatio-temporal model for visual attention understanding presented in Chapter 3 are provided for all the video genres in CRCNS-ORIG [15] and DIEM [16] databases. For further comprehension of the diagrams provided, the reader is referred to section 4.3.

B.1 CRCNS-ORIG DATABASE

B.1.1 *Description*

CRCNS-ORIG [15] dataset contains eye movement recordings from eight distinct subjects freely watching 50 different video clips (over 46,000 video frames, 25 minutes total, 640×480). Eye traces have been obtained using a 240 Hz ISCAN RK-464 eye-tracker. Eye fixations of at least 4 subjects are provided for each clip.

Table B.1: Categories in the CRCNS-ORIG [15] database. Clips included in each category are enumerated, together with their number of frames.

Clip name	Frames		
beverly01	490		
beverly03	481		
beverly05	546		
beverly06	521		
beverly07	357		
beverly08	237		
monica03	1,526		
monica04	640	Clip name	Frames
monica05	611	gamecube02	1,819
monica06	164	gamecube04	2,083
standard01	254	gamecube05	213
standard02	515	gamecube06	2,440
standard03	309	gamecube13	898
standard04	612	gamecube16	2,814
standard05	483	gamecube17	2,114
standard06	434	gamecube18	1,999
standard07	177	gamecube23	1,429
TOTAL	8,357	TOTAL	15,809

(a) Outdoor, 17 clips (b) Videogames, 9 clips

		Clip name	Frames		
		tv-news01	918		
		tv-news02	1,058	Clip name	Frames
Clip name	Frames	tv-news03	1,444	tv-sport01	579
tv-adso1	1,077	tv-news04	491	tv-sport02	444
tv-adso2	387	tv-news05	1,341	tv-sport03	1,460
tv-adso3	841	tv-news06	1,643	tv-sport04	982
tv-adso4	313	tv-news09	1,176	tv-sport05	1,386
TOTAL	2,618	TOTAL	8,071	TOTAL	4,851

(c) Commercials, 4 clips (d) TV News, 7 clips (e) Sports, 5 clips

Clip name	Frames	Clip name	Frames
tv-talk01	1,651	saccadetest	516
tv-talk03	783	tv-action01	567
tv-talk04	1,258	tv-announce01	434
tv-talk05	552	tv-musico1	1,022
TOTAL	4,244	TOTAL	2,539

(f) Talk Shows, 4 clips (g) Others, 4 clips

B.1.2 Video categories

For our experiments on context-driven visual attention understanding and prediction presented in Chapter 4, we have divided the dataset into seven categories: *Outdoor*, *Videogames*, *Commercials*, *TV News*, *Sports*, *Talk Shows* and *Others*. Videos included in each category are enumerated in Table B.1.

B.1.3 Context-aware visual attention understanding

Figures B.2-B.8 illustrate the three most prominent attraction (AT) and inhibition (IT) sub-tasks determined by the above-mentioned approach for modeling visual attention in all database contexts. Moreover, for the sake of comparison, significant sub-tasks that define a *context-generic* model trained on frames from the whole database are also provided in Figure B.1.

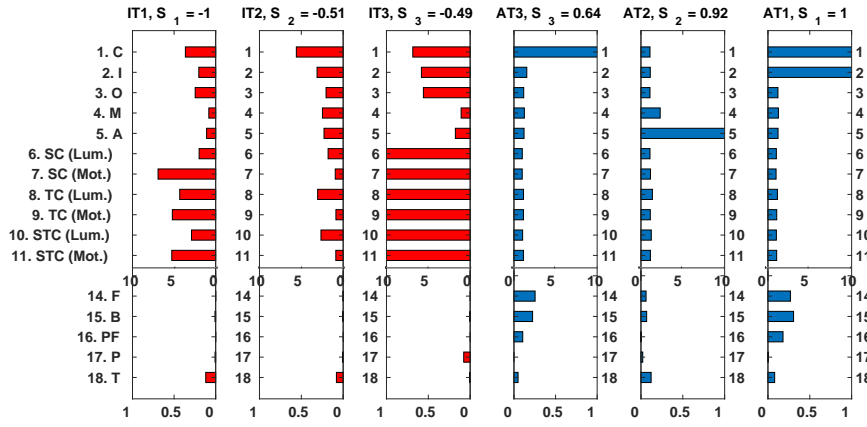


Figure B.1: CRCNS-ORIG [15] database: *Context-Generic*

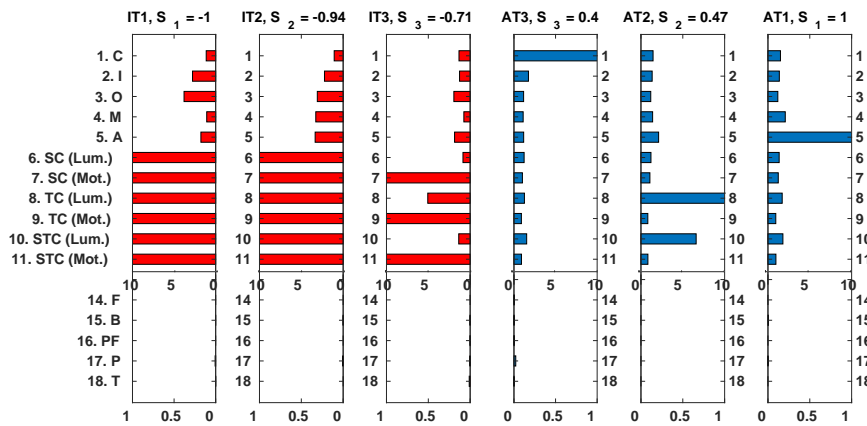
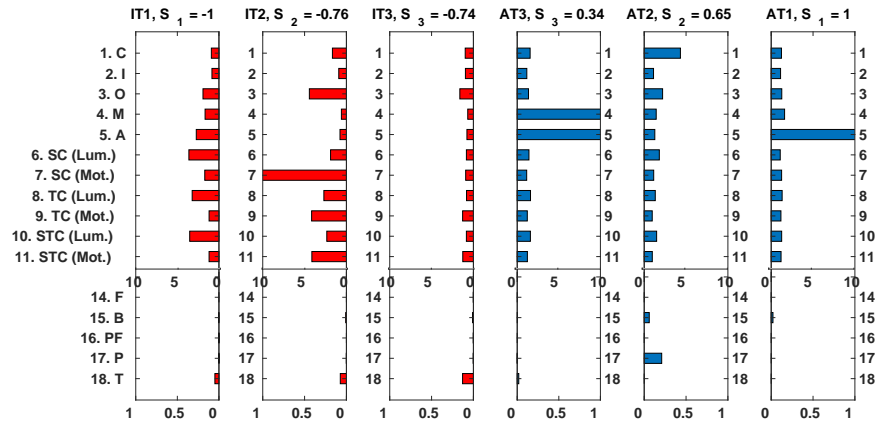
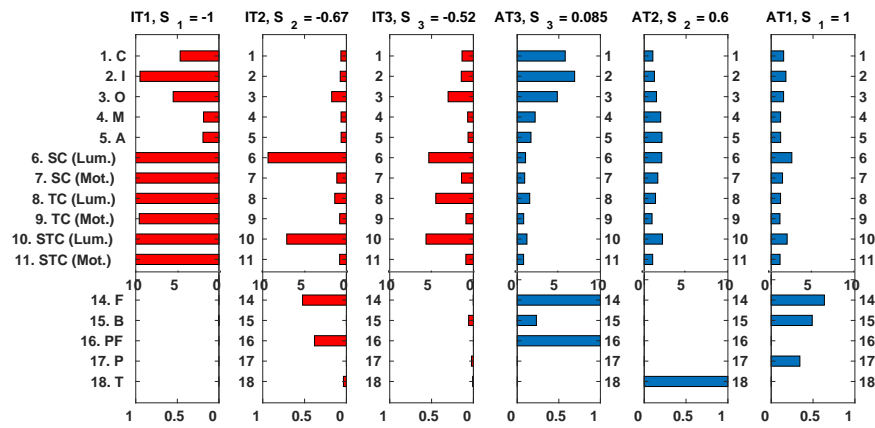
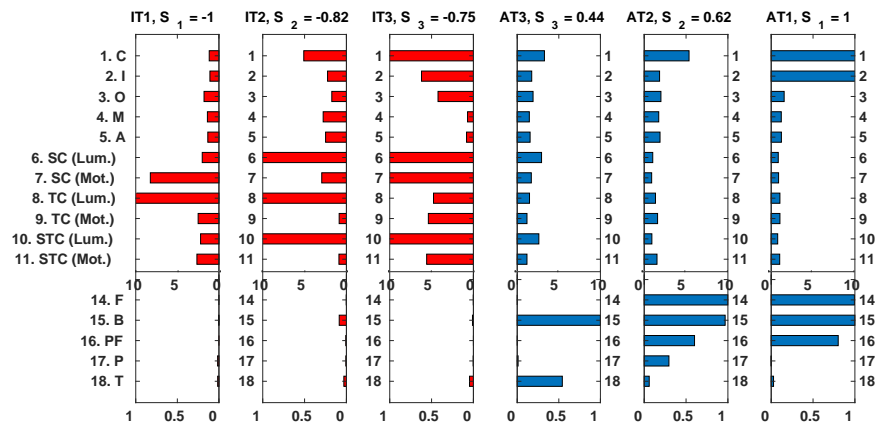
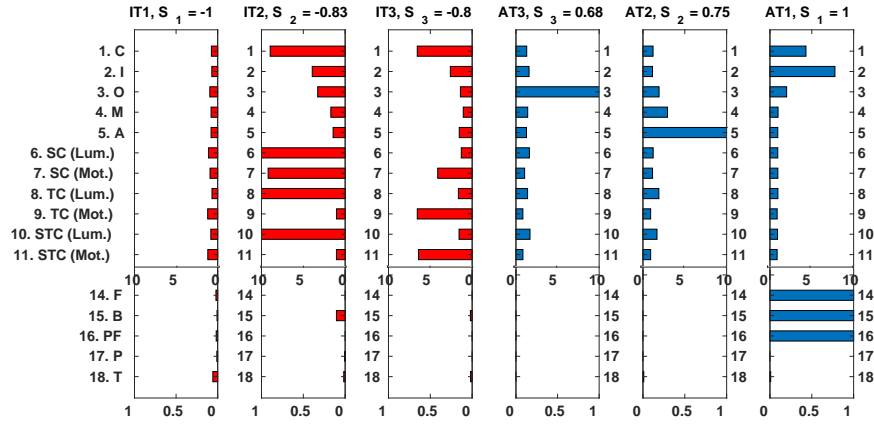
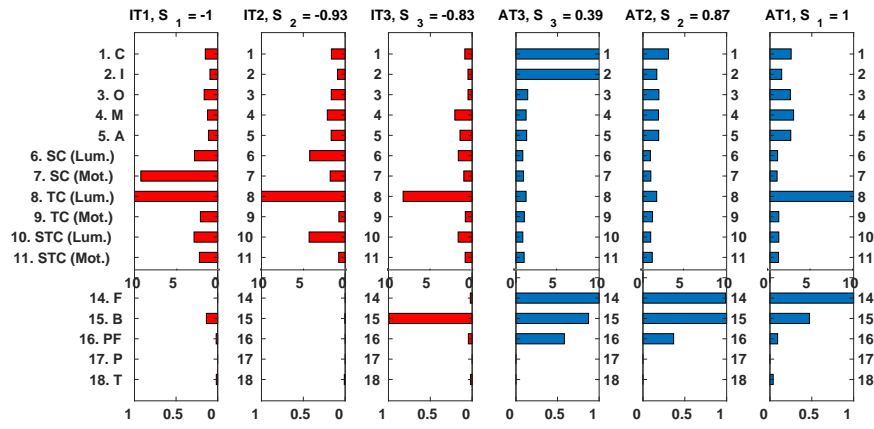
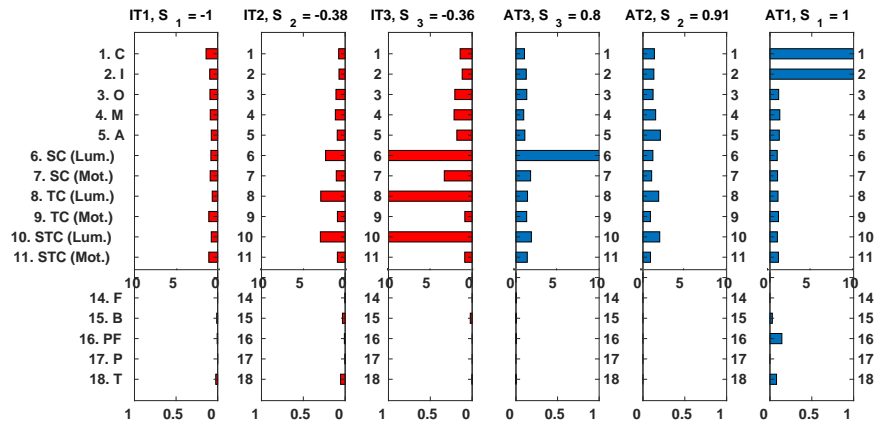


Figure B.2: CRCNS-ORIG [15] database: *Outdoor*

Figure B.3: CRCNS-ORIG [15] database: *Videogames*Figure B.4: CRCNS-ORIG [15] database: *Commercials*Figure B.5: CRCNS-ORIG [15] database: *TV News*

Figure B.6: CRCNS-ORIG [15] database: *Sports*Figure B.7: CRCNS-ORIG [15] database: *Talk Shows*Figure B.8: CRCNS-ORIG [15] database: *Others*

B.2 DIEM DATABASE

B.2.1 *Description*

DIEM [16] dataset contains eye movement recordings from over 250 participants freely watching 84 high-definition natural videos (over 240,000 video frames, 134 minutes total, variable dimensions). Eye traces have been obtained using a 1,000 Hz SR Research Eyelink 2000 desktop mounted eye tracker. Eye fixations from approximately 50 subjects are provided for each clip.

B.2.2 *Video categories*

For the experiments on context-driven visual attention understanding and prediction presented in Chapter 4, database clips have been classified into seven categories: *TV Shows*, *Documentaries*, *Commercials*, *Talk Shows*, *Sports*, *Cooking* and *TV News*, as enumerated in Table B.2.

Table B.2: Categories in the DIEM [16] database. Clips included in each category are enumerated, together with their number of frames.

Clip name	Frames
50_people_brooklyn_1280x720	3,669
50_people_brooklyn_no_voices_1280x720	3,669
50_people_london_1280x720	3,840
50_people_london_no_voices_1280x720	3,840
DIY_SOS_1280x712	1,200
home_movie_Charlie_bit_my_finger_again_960x720	1,661
one_show_1280x712	1,430
stewart_lee_1280x712	2,412
tv_graduates_1280x720	4,045
tv_ketch2_672x544	2,286
tv_the_simpsons_860x528	3,642
tv_uni_challenge_final_1280x712	2,577
TOTAL	34,271

(a) TV Shows, 12 clips

Clip name	Frames
Antarctica_landscape_1246x720	2,135
BBC_life_in_cold_blood_1278x710	3,401
BBC_wildlife_eagle_930x720	3,960
BBC_wildlife_serpent_1280x704	1,038
BBC_wildlife_special_tiget_1276x720	4,320
artic_bears_1066x710	2,786
documentary_adrenaline_rush_1280x720	3,282
documentary_coral_reef_adventure_1280x720	2,969
documentary_discoverers_1280x720	4,560
documentary_dolphins_1280x720	3,181
documentary_mystery_nile_1280x720	2,604
documentary_planet_earth_1280x704	5,082
hummingbirds_closeups_960x720	4,217
hummingbirds_narrator_960x720	1,162
nightlife_in_mozambique_1280x580	1,421
planet_earth_jungles_frogs_1280x704	4,371
planet_earth_jungles_monkeys_1280x704	4,475
university_forum_construction_ionic_1280x720	1,418
TOTAL	56,382

(b) Documentaries, 18 clips

Clip name	Frames
advert_bbc4_bees_1024x576	1,217
advert_bbc4_library_1024x576	1,202
advert_bravia_paint_1280x720	2,167
advert_iphone_1272x720	900
game_trailer_bullet_witch_1280x720	3,720
game_trailer_ghostbusters_1280x720	3,103
game_trailer_lego_indiana_jones_1280x720	3,314
game_trailer_wrath_lich_king_shortened_subtitles_1280x548	5,420
harry_potter_6_trailer_1280x544	2,928
movie_trailer_alice_in_wonderland_1280x682	2,538
movie_trailer_ice_age_3_1280x690	3,283
movie_trailer_quantum_of_solace_1280x688	2,998
music_gummybear_880x720	888
music_red_hot_chili_peppers_shortened_1024x576	5,597
music_trailer_nine_inch_nails_1280x720	1,283
TOTAL	40,558

(c) Commercials, 15 clips

Clip name	Frames
ami_ib4010_closeup_720x576	1,080
ami_ib4010_left_720x576	1,067
ami_is1000a_closeup_720x576	1,262
ami_is1000a_left_720x576	1,270
scottish_parliament_1152x864	3,978
TOTAL	8,657

(d) Talk Shows, 5 clips

Clip name	Frames
basketball_of_sorts_960x720	3,476
one_show_1280x712	900
pingpong_angle_shot_960x720	1,170
pingpong_closeup_rallies_960x720	3,300
pingpong_long_shot_960x720	3,772
pingpong_miscues_1080x720	1,371
pingpong_no_bodies_960x720	4,371
sport_F1_slick_tyres_1280x720	2,259
sport_barcelona_extreme_1280x720	1,721
sport_cricket_ashes_2007_1252x720	2,574
sport_football_best_goals_976x720	2,478
sport_golf_fade_a_driver_1280x720	2,410
sport_poker_1280x640	3,480
sport_scramblers_1280x720	1,525
sport_slam_dunk_1280x720	5,747
sport_surfing_in_thurso_900x720	2,357
sport_wimbledon_baltacha_1280x704	5,818
sport_wimbledon_federer_final_1280x704	2,772
sport_wimbledon_magic_wand_1280x704	1,768
sport_wimbledon_murray_1280x704	2,627
sports_kendo_1280x710	2,768
TOTAL	54,293

(e) Sports, 20 clips

Clip name	Frames
chilli_plasters_1280x712	3,697
growing_sweetcorn_1280x712	2,223
hairy_bikers_cabbage_1280x712	3,121
hydraulics_1280x712	3,611
nigella_chocolate_pears_1280x712	5,393
scottish_starters_1280x712	3,123
spotty_trifle_1280x712	2,516
TOTAL	23,684

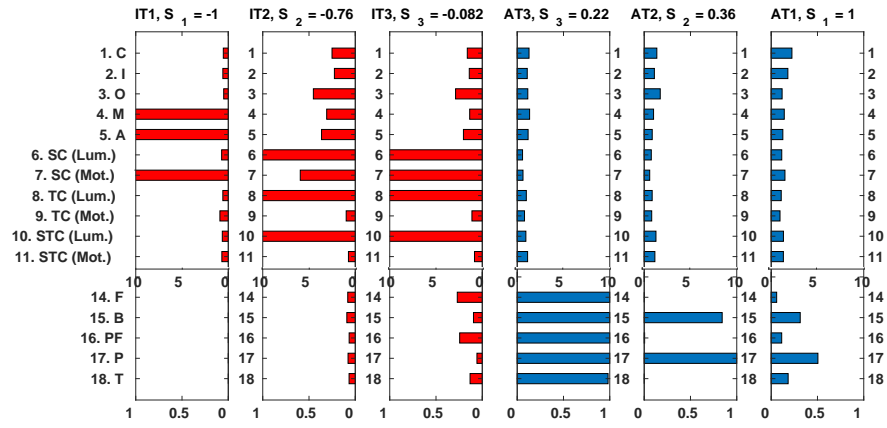
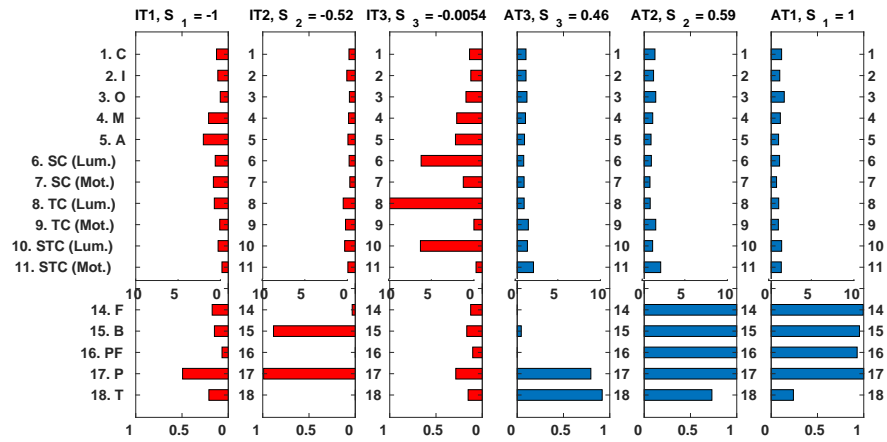
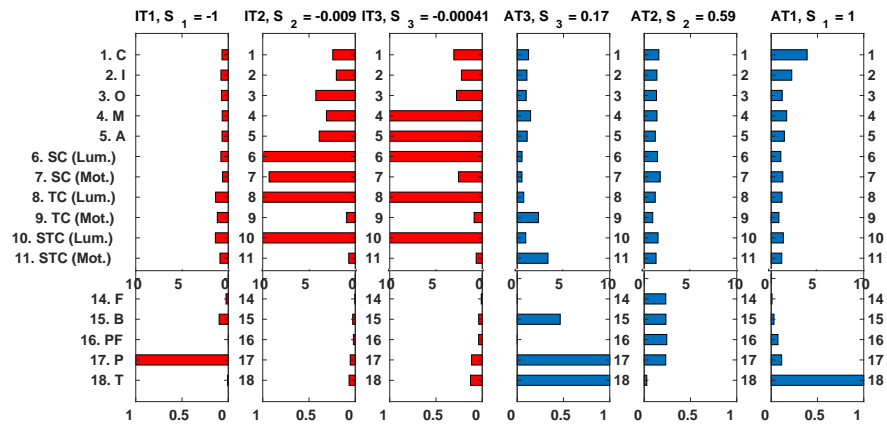
(f) Cooking, 7 clips

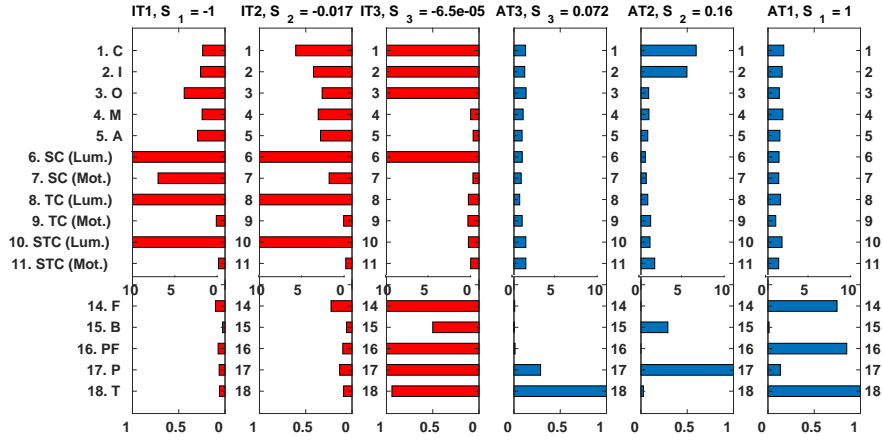
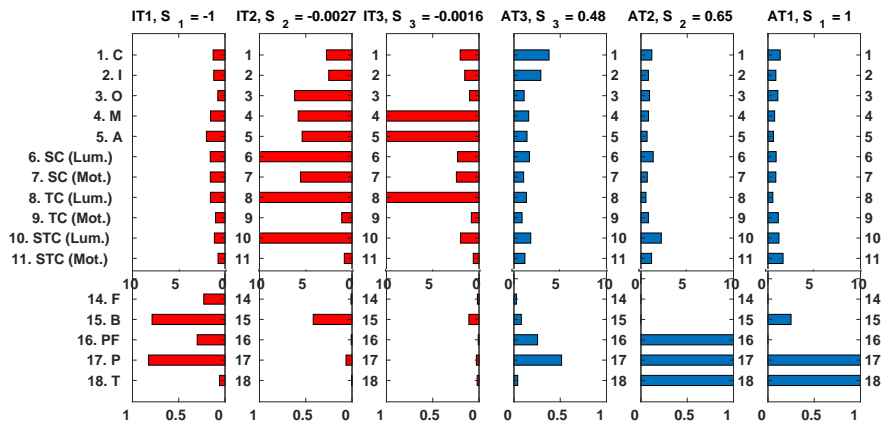
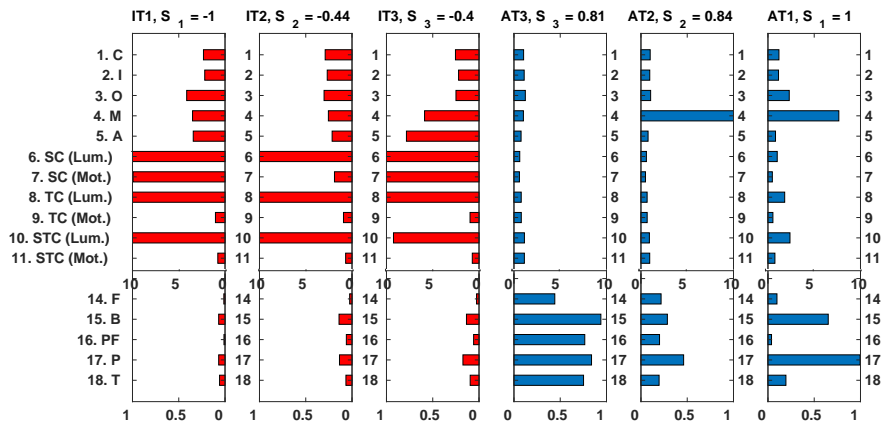
Clip name	Frames
news_newsnight_othello_720x416	2,295
news_sherry_drinking_mice_768x576	1,999
news_tony_blair_resignation_720x540	1,413
news_us_election_debate_1080x600	2,572
news_video_republic_960x720	6,276
news_wimbledon_macenroe_shortened_1024x576	4,980
TOTAL	22,607

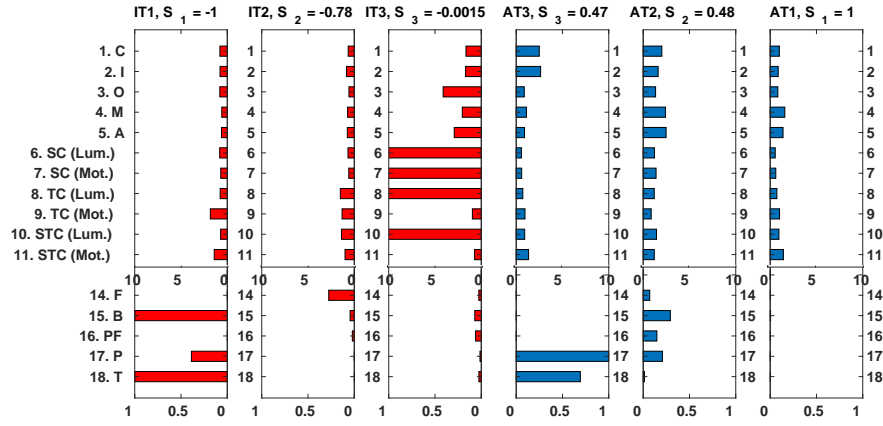
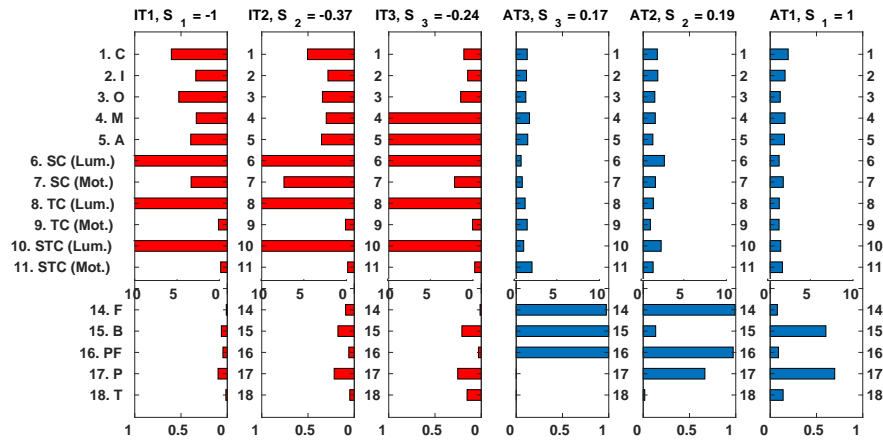
(g) TV News, 7 clips

B.2.3 Context-aware visual attention understanding

Figures B.10-B.16 illustrate the most noteworthy attraction (AT) and inhibition (IT) sub-tasks determined by the above-mentioned approach for modeling visual attention in all database contexts. Moreover, for the sake of comparison, significant sub-tasks that define a *context-generic* model trained on frames from the whole database are also provided in Figure B.9.

Figure B.9: DIEM [16] database: *Context-Generic*Figure B.10: DIEM [16] database: *TV Shows*Figure B.11: DIEM [16] database: *Documentaries*

Figure B.12: DIEM [16] database: *Commercials*Figure B.13: DIEM [16] database: *Talk Shows*Figure B.14: DIEM [16] database: *Sports*

Figure B.15: DIEM [16] database: *Cooking*Figure B.16: DIEM [16] database: *TV News*

B.3 BOSS DATABASE

B.3.1 Description

Within the framework of the BOSS project [19], a database with 15 video sequences recorded in RENFE suburban trains from Madrid was released, with the aim of developing an efficient transmission system for video-surveillance in a railway transport context. Videos contain events such as a cell phone theft, a passengers fight, a disease in public and several women harassment. Moreover, two additional sequences with no incidents are included. For each event, three camera views are provided.

B.3.2 Video sequences

In order to evaluate the architectures for visual attention modeling in the temporal domain proposed in Chapter 5, we have selected the

three camera views of 10 sequences from this database to be annotated with eye fixations. In total, 30 videos (over 84,000 video frames, 56 minutes total, 720×576) have been used, which are enumerated in Table B.3. For each video, eye traces from 5 observers have been recorded by using a 250 Hz SMI RED250mobile Eye Tracker system [212].

Table B.3: Videos from the BOSS [19] database for the experiments in Chapter 6. Clips are enumerated together with their number of frames.

Clip name	Frames
Cell_phone_Spanish.Cam1	1,501
Cell_phone_Spanish.Cam2	1,501
Cell_phone_Spanish.Cam3	1,501
Checkout_French.Cam1	3,941
Checkout_French.Cam2	2,810
Checkout_French.Cam3	2,843
Disease_Public.Cam1	3,082
Disease_Public.Cam2	3,086
Disease_Public.Cam3	3,088
Harass_French.Cam1	2,674
Harass_French.Cam2	2,679
Harass_French.Cam3	2,679
Harass2_French.Cam1	2,976
Harass2_French.Cam2	2,976
Harass2_French.Cam3	2,976
Harass_Spanish.Cam1	2,976
Harass_Spanish.Cam2	2,976
Harass_Spanish.Cam3	2,976
Newspaper_Spanish.Cam1	2,438
Newspaper_Spanish.Cam2	2,438
Newspaper_Spanish.Cam3	2,438
No_Event.Cam1	2,630
No_Event.Cam2	2,635
No_Event.Cam3	2,636
No_Event2.Cam1	4,001
No_Event2.Cam2	4,001
No_Event2.Cam3	4,001
Panic.Cam1	2,501
Panic.Cam2	2,501
Panic.Cam3	2,501
TOTAL	83,962

BIBLIOGRAPHY

- [1] David Forsyth and Jean Ponce. *Computer Vision: A Modern Approach*. 2011.
- [2] G. Li and Y. Yu. "Deep Contrast Learning for Salient Object Detection". In: *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2016, pp. 478–487.
- [3] Alex Krizhevsky, Ilya Sutskever and Geoffrey E Hinton. "ImageNet classification with deep convolutional neural networks". In: *Advances in neural information processing systems*. 2012, pp. 1097–1105.
- [4] F. de-la Calle-Silos, I. González-Díaz and F. Díaz de María. "Mid-level feature set for specific event and anomaly detection in crowded scenes". In: *2013 IEEE International Conference on Image Processing*. 2013, pp. 4001–4005. DOI: [10.1109/ICIP.2013.6738824](https://doi.org/10.1109/ICIP.2013.6738824).
- [5] Arnold WM Smeulders, Dung M Chu, Rita Cucchiara, Simone Calderara, Afshin Dehghan and Mubarak Shah. "Visual tracking: An experimental survey". In: *IEEE transactions on pattern analysis and machine intelligence* 36.7 (2014), pp. 1442–1468.
- [6] Anuj Kumar Pradhan, Kim R Hammel, Rosa DeRamus, Alexander Pollatsek, David A Noyce and Donald L Fisher. "Using eye movements to evaluate effects of driver age on risk perception in a driving simulator". In: *Human factors* 47.4 (2005), pp. 840–852.
- [7] Kavyaganga Kilingaru, Jeffrey W Tweedale, Steve Thatcher and Lakhmi C Jain. "Monitoring pilot "situation awareness"". In: *Journal of Intelligent & Fuzzy Systems* 24.3 (2013), pp. 457–466.
- [8] Christina J Howard, Iain D Gilchrist, Tom Troscianko, Ardhendu Behera and David C Hogg. "Task relevance predicts gaze in videos of real moving scenes". In: *Experimental brain research* 214.1 (2011), p. 131.
- [9] Ali Borji and Laurent Itti. "State-of-the-Art in Visual Attention Modeling". In: *IEEE Trans. Pattern Anal. Mach. Intell.* 35.1 (Jan. 2013), pp. 185–207. ISSN: 0162-8828. DOI: [10.1109/TPAMI.2012.89](https://doi.org/10.1109/TPAMI.2012.89).
- [10] A. M. Treisman and G. Gelade. "A Feature-Integration Theory of Attention." In: *Cognitive Psychology* 12 (1980), pp. 97–136.

- [11] Jeremy M Wolfe. "Guided search 2.0 a revised model of visual search". In: *Psychonomic bulletin & review* 1.2 (1994), pp. 202–238.
- [12] David M. Blei, Andrew Y. Ng and Michael I. Jordan. "Latent Dirichlet Allocation". In: *J. Mach. Learn. Res.* 3 (Mar. 2003), pp. 993–1022. ISSN: 1532-4435.
- [13] Jon D. Mcauliffe and David M. Blei. "Supervised Topic Models". In: *Advances in Neural Information Processing Systems* 20. Ed. by J. C. Platt, D. Koller, Y. Singer and S. T. Roweis. Curran Associates, Inc., 2008, pp. 121–128.
- [14] Shuang-Hong Yang, Hongyuan Zha and Bao-Gang Hu. "Dirichlet-bernoulli alignment: A generative model for multi-class multi-label multi-instance corpora". In: *Advances in neural information processing systems*. 2009, pp. 2143–2150.
- [15] L. Itti and R. Carmi. *Eye-tracking data from human volunteers watching complex video stimuli*. bu;eye;mod. 2009.
- [16] Parag K Mital, Tim J Smith, Robin L Hill and John M Henderson. "Clustering of gaze during dynamic scene viewing is predicted by motion". In: *Cognitive Computation* 3.1 (2011), pp. 5–24.
- [17] Fu Jie Huang, Y-Lan Boureau, Yann LeCun et al. "Unsupervised learning of invariant feature hierarchies with applications to object recognition". In: *Computer Vision and Pattern Recognition, 2007. CVPR'07. IEEE Conference on*. IEEE. 2007, pp. 1–8.
- [18] Sepp Hochreiter and Jürgen Schmidhuber. "Long short-term memory". In: *Neural computation* 9.8 (1997), pp. 1735–1780.
- [19] BOSS European project.
[http : / / www . multitel . be / image / research - development / research - projects / boss . php](http://www.multitel.be/image/research-development/research-projects/boss.php). Accessed: 2016-09-30.
- [20] Hao Wang and Dit-Yan Yeung. "Towards bayesian deep learning: A survey". In: *arXiv preprint arXiv:1604.01662* (2016).
- [21] Volodymyr Mnih, Nicolas Heess, Alex Graves et al. "Recurrent models of visual attention". In: *Advances in neural information processing systems*. 2014, pp. 2204–2212.
- [22] Sangmin Oh, Anthony Hoogs, Amitha Perera, Naresh Cuntoor, Chia-Chih Chen, Jong Taek Lee, Saurajit Mukherjee, JK Aggarwal, Hyungtae Lee, Larry Davis et al. "A large-scale benchmark dataset for event recognition in surveillance video". In: *Computer Vision and Pattern Recognition (CVPR), 2011 IEEE conference on*. IEEE. 2011, pp. 3153–3160.

- [23] Waqas Sultani, Chen Chen and Mubarak Shah. "Real-world Anomaly Detection in Surveillance Videos". In: *arXiv preprint arXiv:1801.04264* (2018).
- [24] Mark Andrejevic. "The big data divide". In: *International Journal of Communication* 8 (2014), p. 17.
- [25] A. L. Yarbus. *Eye Movements and Vision*. Plenum. New York., 1967.
- [26] Yoshua Bengio, Aaron Courville and Pascal Vincent. "Representation learning: A review and new perspectives". In: *IEEE transactions on pattern analysis and machine intelligence* 35.8 (2013), pp. 1798–1828.
- [27] Andrew Ng. *Machine Learning and AI via Brain simulations*.
- [28] Jean Martinet and Ismail Elsayad. "Mid-level image descriptors". In: *Intelligent Multimedia Databases and Information Retrieval: Advancing Applications and Technologies*. IGI Global, 2012, pp. 46–60.
- [29] Jeremy M Wolfe. "Approaches to Visual Search: Feature Integration Theory and Guided Search". In: *The Oxford Handbook of Attention* (2014), p. 11.
- [30] Yoshua Bengio. "Learning deep architectures for AI". In: *Foundations and trends® in Machine Learning* 2.1 (2009), pp. 1–127.
- [31] Arthur L Samuel. "Some studies in machine learning using the game of checkers". In: *IBM Journal of research and development* 3.3 (1959), pp. 210–229.
- [32] Thomas M. Mitchell. *Machine Learning*. 1st ed. New York, NY, USA: McGraw-Hill, Inc., 1997. ISBN: 0070428077, 9780070428072.
- [33] David H Wolpert and William G Macready. "No free lunch theorems for optimization". In: *IEEE transactions on evolutionary computation* 1.1 (1997), pp. 67–82.
- [34] Juan C Caicedo and Svetlana Lazebnik. "Active object localization with deep reinforcement learning". In: *Proceedings of the IEEE International Conference on Computer Vision*. 2015, pp. 2488–2496.
- [35] Miriam Bellver Buena, Xavier Giro-i Nietob, Ferran Marquesb and Jordi Torresa. "Hierarchical object detection with deep reinforcement learning". In: *Deep Learning for Image Processing Applications* 31 (2017), p. 164.

- [36] Javier López-Labracca, Miguel Ángel Fernández-Torres, Iván González-Díaz, Fernando Díaz-de María and Ángel Pizarro. "Enriched dermoscopic-structure-based cad system for melanoma diagnosis". In: *Multimedia Tools and Applications* 77.10 (2018), pp. 12171–12202.
- [37] F. Fernández-Martínez, A. Hernández-García, M. A. Fernández-Torres, I. González-Díaz, Á. García-Faura and F. Díaz de María. "Exploiting visual saliency for assessing the impact of car commercials upon viewers". In: *Multimedia Tools and Applications* 77.15 (2018), pp. 18903–18933. ISSN: 1573-7721. DOI: [10 . 1007 / s11042 - 017 - 5339 - 9](https://doi.org/10.1007/s11042-017-5339-9). URL: <https://doi.org/10.1007/s11042-017-5339-9>.
- [38] Tomás Martínez-Cortés, Miguel Ángel Fernández-Torres, Amaya Jiménez-Moreno, Iván González-Díaz, Fernando Díaz-de María, Juan Adán Guzmán-De-Villoria and Pilar Fernández. "A Bayesian model for brain tumor classification using clinical-based features". In: *Image Processing (ICIP), 2014 IEEE International Conference on*. IEEE. 2014, pp. 2779–2783.
- [39] Iván González-Díaz and Fernando Díaz de María. "A region-centered topic model for object discovery and category-based image segmentation". In: *Pattern Recognition* 46.9 (2013), pp. 2437 –2449. ISSN: 0031-3203. DOI: <https://doi.org/10.1016/j.patcog.2013.01.034>. URL: <http://www.sciencedirect.com/science/article/pii/S0031320313000733>.
- [40] Hui Zhang, Jason E Fritts and Sally A Goldman. "Image segmentation evaluation: A survey of unsupervised methods". In: *computer vision and image understanding* 110.2 (2008), pp. 260–280.
- [41] Christopher M. Bishop. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Berlin, Heidelberg: Springer-Verlag, 2006. ISBN: 0387310738.
- [42] Kevin P. Murphy. *Machine learning : a probabilistic perspective*. 1st ed. MIT Press, Aug. 2013. ISBN: 0262018020. URL: <http://www.amazon.com/exec/obidos/redirect?tag=citeulike07-20&path=ASIN/0262018020>.
- [43] Ian Goodfellow, Yoshua Bengio, Aaron Courville and Yoshua Bengio. *Deep learning*. Vol. 1. MIT press Cambridge, 2016.
- [44] Andrew Y Ng and Michael I Jordan. "On discriminative vs. generative classifiers: A comparison of logistic regression and naive bayes". In: *Advances in neural information processing systems*. 2002, pp. 841–848.

- [45] Christopher M. Bishop and Julia Lasserre. "Generative or Discriminative? Getting the Best of Both Worlds". In: *Bayesian Statistics 8*. Ed. by J. M. Bernardo, M. J. Bayarri, J. O. Berger, A. P. Dawid, D. Heckerman, A. F. M. Smith and M. West. International Society for Bayesian Analysis. Oxford University Press, 2007, pp. 3–24. URL: [http : / / research . microsoft . com / en - us / um / people / cmbishop / downloads / Bishop - Valencia - 07.pdf](http://research.microsoft.com/en-us/um/people/cmbishop/downloads/Bishop-Valencia-07.pdf).
- [46] Daphne Koller and Nir Friedman. *Probabilistic graphical models: principles and techniques*. MIT press, 2009.
- [47] Leonard E Baum and Ted Petrie. "Statistical inference for probabilistic functions of finite state Markov chains". In: *The annals of mathematical statistics* 37.6 (1966), pp. 1554–1563.
- [48] Corinna Cortes and Vladimir Vapnik. "Support-vector networks". In: *Machine learning* 20.3 (1995), pp. 273–297.
- [49] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville and Yoshua Bengio. "Generative adversarial nets". In: *Advances in neural information processing systems*. 2014, pp. 2672–2680.
- [50] Sinno Jialin Pan, Qiang Yang et al. "A survey on transfer learning". In: *IEEE Transactions on knowledge and data engineering* 22.10 (2010), pp. 1345–1359.
- [51] Naiyan Wang and Dit-Yan Yeung. "Learning a deep compact image representation for visual tracking". In: *Advances in neural information processing systems*. 2013, pp. 809–817.
- [52] Geoffrey Hinton, Li Deng, Dong Yu, George E Dahl, Abdel-rahman Mohamed, Navdeep Jaitly, Andrew Senior, Vincent Vanhoucke, Patrick Nguyen, Tara N Sainath et al. "Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups". In: *IEEE Signal processing magazine* 29.6 (2012), pp. 82–97.
- [53] Philippe Hamel and Douglas Eck. "Learning Features from Music Audio with Deep Belief Networks." In: *ISMIR*. Vol. 10. Utrecht, The Netherlands. 2010, pp. 339–344.
- [54] Lei Zhang, Shuai Wang and Bing Liu. "Deep learning for sentiment analysis: A survey". In: *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery* (), e1253.
- [55] Ilya Sutskever, Oriol Vinyals and Quoc V Le. "Sequence to sequence learning with neural networks". In: *Advances in neural information processing systems*. 2014, pp. 3104–3112.

- [56] Stephen E. Palmer. *Vision science : photons to phenomenology*. Stephen E. Palmer.; "A Bradford book."; Bibliography: Includes bibliographical references (p. [737]-769) and indexes. Cambridge, Mass.: MIT Press, 1999, p. 810. ISBN: 0262161834 Thanks for using Barton, the MIT Libraries' catalog [http](http://).
- [57] John K. Tsotsos. *A Computational Perspective on Visual Attention*. 1st. The MIT Press, 2011. ISBN: 0262015412, 9780262015417.
- [58] Neil D. B. Bruce and John K. Tsotsos. "Saliency Based on Information Maximization". In: *NIPS*. 2005.
- [59] L. Itti and P. F. Baldi. "Bayesian Surprise Attracts Human Attention". In: *Advances in Neural Information Processing Systems, Vol. 19 (NIPS*2005)*. Cambridge, MA: MIT Press, 2006, pp. 547-554.
- [60] Nathan Sprague and Dana Ballard. "Eye movements for reward maximization". In: *Advances in neural information processing systems*. 2003, None.
- [61] Wikimedia Commons. *Human Eye Transverse Cut Unlabeled*. File: Human Eye Transverse Cut Unlabeled.jpg. 2016. URL: https://commons.wikimedia.org/wiki/File:Human_Eye_Transverse_Cut_Unlabeled.jpg.
- [62] Justin Johnson Fei-Fei Li and Serena Yeung. *Lecture notes in CS231n: Convolutional Neural Networks for Visual Recognition*. 2018. URL: <http://cs231n.stanford.edu/>.
- [63] Tony Lindeberg. "A computational theory of visual receptive fields". In: *Biological cybernetics* 107.6 (2013), pp. 589-635.
- [64] Yann LeCun, Bernhard Boser, John S Denker, Donnie Henderson, Richard E Howard, Wayne Hubbard and Lawrence D Jackel. "Backpropagation applied to handwritten zip code recognition". In: *Neural computation* 1.4 (1989), pp. 541-551.
- [65] Semir M Zeki. "Functional specialisation in the visual cortex of the rhesus monkey". In: *Nature* 274.5670 (1978), p. 423.
- [66] Margaret Livingstone and David Hubel. "Segregation of form, color, movement, and depth: anatomy, physiology, and perception". In: *Science* 240.4853 (1988), pp. 740-749.
- [67] Wikimedia Commons. *Brain diagram without text*. File: Brain diagram without text.svg. 2008. URL: https://commons.wikimedia.org/wiki/File:Brain_diagram_without_text.svg.
- [68] D. Purves, D. Fitzpatrick, L.C. Katz, A.S. Lamantia, J.O. McNamara, S.M. Williams and G.J. Augustine. *Neuroscience*. Sinauer Associates, 2001. ISBN: 9780878937431. URL: <https://books.google.es/books?id=F4pTPwAACAAJ>.

- [69] James W Bisley. "The neural basis of visual attention". In: *The Journal of physiology* 589.1 (2011), pp. 49–57.
- [70] Simone Frintrop, Erich Rome and Henrik I Christensen. "Computational visual attention systems and their cognitive foundations: A survey". In: *ACM Transactions on Applied Perception (TAP)* 7.1 (2010), p. 6.
- [71] J. M. Wolfe and T. S. Horowitz. "What attributes guide the deployment of visual attention and how do they do it?" In: *Nature Reviews Neuroscience* 5.6 (2004), pp. 495–501.
- [72] Jeremy M Wolfe. "Guided search 4.0". In: *Integrated models of cognitive systems* (2007), pp. 99–119.
- [73] C. Koch and S. Ullman. "Shifts in selective visual attention: towards the underlying neural circuitry." In: *Human neurobiology* 4.4 (1985), pp. 219–227. ISSN: 0721-9075.
- [74] Jeremy M Wolfe, Melissa L-H Võ, Karla K Evans and Michelle R Greene. "Visual search in scenes involves selective and nonselective pathways". In: *Trends in cognitive sciences* 15.2 (2011), pp. 77–84.
- [75] Peter McLeod, Jon Driver and Jennie Crisp. "Visual search for a conjunction of movement and form is parallel". In: *Nature* 332.6160 (1988), p. 154.
- [76] John M Henderson and Andrew Hollingworth. "High-level scene perception". In: *Annual review of psychology* 50.1 (1999), pp. 243–271.
- [77] Ali Borji and Laurent Itti. "Defending Yarbus: Eye movements reveal observers' task". In: *Journal of vision* 14.3 (2014), pp. 29–29.
- [78] Chandresh Sharma, Punitkumar Bhavsar, Babji Srinivasan and Rajagopalan Srinivasan. "Eye gaze movement studies of control room operators: A novel approach to improve process safety". In: *Computers & Chemical Engineering* 85 (2016), pp. 43–57.
- [79] Frouke Hermens, Rhona Flin and Irfan Ahmed. "Eye movements in surgery: A literature review". In: *Journal of Eye movement research* 6.4 (2013).
- [80] Christina Jayne Howard, Tom Troscianko, Iain D Gilchrist, Ardhendu Behera and David C Hogg. "Suspiciousness perception in dynamic scenes: a comparison of CCTV operators and novices". In: *Frontiers in human neuroscience* 7 (2013), p. 441.

- [81] Laurent Itti, Christof Koch and Ernst Niebur. "A Model of Saliency-Based Visual Attention for Rapid Scene Analysis". In: *IEEE Trans. Pattern Anal. Mach. Intell.* 20.11 (Nov. 1998), pp. 1254–1259. ISSN: 0162-8828. DOI: [10.1109/34.730558](https://doi.org/10.1109/34.730558).
- [82] Antonio Torralba. "Modeling global scene factors in attention". In: *J. Opt. Soc. Am. A* 20.7 (2003), pp. 1407–1418. DOI: [10.1364/JOSAA.20.001407](https://doi.org/10.1364/JOSAA.20.001407).
- [83] V. Navalpakkam and L. Itti. "An Integrated Model of Top-Down and Bottom-Up Attention for Optimizing Detection Speed". In: *Computer Vision and Pattern Recognition, 2006 IEEE Computer Society Conference on*. Vol. 2. 2006, pp. 2049–2056. DOI: [10.1109/CVPR.2006.54](https://doi.org/10.1109/CVPR.2006.54).
- [84] Jonathan Harel, Christof Koch and Pietro Perona. "Graph-based visual saliency". In: *Advances in Neural Information Processing Systems* 19. MIT Press, 2007, pp. 545–552.
- [85] X. Hou and L. Zhang. "Saliency Detection: A Spectral Residual Approach". In: *2007 IEEE Conference on Computer Vision and Pattern Recognition*. 2007, pp. 1–8. DOI: [10.1109/CVPR.2007.383267](https://doi.org/10.1109/CVPR.2007.383267).
- [86] Robert J Peters and Laurent Itti. "Beyond bottom-up: Incorporating task-dependent influences into a computational model of spatial attention". In: *2007 IEEE Conference on Computer Vision and Pattern Recognition*. IEEE. 2007, pp. 1–8.
- [87] Lingyun Zhang, Matthew H. Tong, Tim K. Marks, Honghao Shan and Garrison W. Cottrell. "SUN: A Bayesian framework for saliency using natural statistics". In: *Journal of Vision* 8.7 (2008), p. 32. DOI: [10.1167/8.7.32](https://doi.org/10.1167/8.7.32). eprint: [/data/Journals/JOV/933536/jov-8-7-32.pdf](https://data.journals.jov.org/933536/jov-8-7-32.pdf).
- [88] Xiaodi Hou and Liqing Zhang. "Dynamic visual attention: Searching for coding length increments". In: *Advances in neural information processing systems*. 2009, pp. 681–688.
- [89] Hae Jong Seo and Peyman Milanfar. "Static and space-time visual saliency detection by self-resemblance". In: *Journal of Vision* 9.12 (2009), p. 15. DOI: [10.1167/9.12.15](https://doi.org/10.1167/9.12.15). eprint: [/data/journals/jov/932859/jov-9-12-15.pdf](https://data/journals/jov/932859/jov-9-12-15.pdf).
- [90] Tilke Judd, Krista Ehinger, Frédo Durand and Antonio Torralba. "Learning to predict where humans look". In: *2009 IEEE 12th International Conference on Computer Vision*. IEEE. 2009, pp. 2106–2113.

- [91] Lingyun Zhang, Matthew H Tong and Garrison W Cottrell. "SUNDAY: Saliency using natural statistics for dynamic analysis of scenes". In: *Proceedings of the 31st annual cognitive science conference*. AAAI Press Cambridge, MA. 2009, pp. 2944–2949.
- [92] Esa Rahtu, Juho Kannala, Mikko Salo and Janne Heikkilä. "Segmenting Salient Objects from Images and Videos". In: *Proceedings of the 11th European Conference on Computer Vision: Part V. ECCV'10*. Heraklion, Crete, Greece: Springer-Verlag, 2010, pp. 366–379. ISBN: 3-642-15554-5, 978-3-642-15554-3.
- [93] C. Guo and L. Zhang. "A Novel Multiresolution Spatiotemporal Saliency Detection Model and Its Applications in Image and Video Compression". In: *IEEE Transactions on Image Processing* 19.1 (2010), pp. 185–198. ISSN: 1057-7149. DOI: [10.1109/TIP.2009.2030969](https://doi.org/10.1109/TIP.2009.2030969).
- [94] S. Goferman, L. Zelnik-Manor and A. Tal. "Context-Aware Saliency Detection". In: *IEEE Transactions on Pattern Analysis and Machine Intelligence* 34.10 (2012), pp. 1915–1926. ISSN: 0162-8828. DOI: [10.1109/TPAMI.2011.272](https://doi.org/10.1109/TPAMI.2011.272).
- [95] Lior Elazary and Laurent Itti. "A Bayesian model for efficient visual search and recognition". In: *Vision Research* 50.14 (2010). Visual Search and Selective Attention, pp. 1338 –1352. ISSN: 0042-6989. DOI: <http://dx.doi.org/10.1016/j.visres.2010.01.002>.
- [96] Fernando López-García, Xosé Ramón Fdez-Vidal, Xosé Manuel Pardo and Raquel Dosil. "Scene recognition through visual attention and image features: A comparison between sift and surf approaches". In: *Object Recognition. InTech*, 2011.
- [97] J. Li, Y. Tian, T. Huang and W. Gao. "Multi-Task Rank Learning for Visual Saliency Estimation". In: *IEEE Transactions on Circuits and Systems for Video Technology* 21.5 (2011), pp. 623–636. ISSN: 1051-8215. DOI: [10.1109/TCSVT.2011.2129430](https://doi.org/10.1109/TCSVT.2011.2129430).
- [98] Antón Garcia-Diaz, Víctor Leborán, Xosé R. Fdez-Vidal and Xosé M. Pardo. "On the relationship between optical variability, visual saliency, and eye fixations: A computational approach". In: *Journal of Vision* 12.6 (2012), p. 17. DOI: [10.1167/12.6.17](https://doi.org/10.1167/12.6.17). eprint: [/data/journals/jov/933494/i1534-7362-12-6-17.pdf](https://www.journalofvision.org/data/journals/jov/933494/i1534-7362-12-6-17.pdf).
- [99] Eleonora Vig, Michael Dorr and David Cox. "Large-scale optimization of hierarchical features for saliency prediction in natural images". In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2014, pp. 2798–2805.

- [100] Matthias Kümmerer, Lucas Theis and Matthias Bethge. "Deep gaze i: Boosting saliency prediction with feature maps trained on imagenet". In: *arXiv preprint arXiv:1411.1045* (2014).
- [101] Xun Huang, Chengyao Shen, Xavier Boix and Qi Zhao. "Salicon: Reducing the semantic gap in saliency prediction by adapting deep neural networks". In: *Proceedings of the IEEE International Conference on Computer Vision*. 2015, pp. 262–270.
- [102] Matthias Kümmerer, Thomas SA Wallis and Matthias Bethge. "DeepGaze II: Reading fixations from deep features trained on object recognition". In: *arXiv preprint arXiv:1610.01563* (2016).
- [103] V. Leboran, A. Garcia-Diaz, X. Fdez-Vidal and X. Pardo. "Dynamic Whitening Saliency". In: *IEEE Transactions on Pattern Analysis and Machine Intelligence* PP.99 (2016), pp. 1–1. ISSN: 0162-8828. DOI: [10.1109/TPAMI.2016.2567391](https://doi.org/10.1109/TPAMI.2016.2567391).
- [104] Srinivas SS Kruthiventi, Kumar Ayush and R Venkatesh Babu. "Deepfix: A fully convolutional neural network for predicting human eye fixations". In: *IEEE Transactions on Image Processing* 26.9 (2017), pp. 4446–4456.
- [105] Lai Jiang, Mai Xu and Zulin Wang. "Predicting Video Saliency with Object-to-Motion CNN and Two-layer Convolutional LSTM". In: *arXiv preprint arXiv:1709.06316* (2017).
- [106] Cagdas Bak, Aysun Kocak, Erkut Erdem and Aykut Erdem. "Spatio-temporal saliency networks for dynamic saliency prediction". In: *IEEE Transactions on Multimedia* 20.7 (2018).
- [107] Wenguan Wang, Jianbing Shen, Fang Guo, Ming-Ming Cheng and Ali Borji. "Revisiting Video Saliency: A Large-scale Benchmark and a New Model". In: *CoRR* abs/1801.07424 (2018). arXiv: [1801.07424](https://arxiv.org/abs/1801.07424). URL: <http://arxiv.org/abs/1801.07424>.
- [108] M. Cornia, L. Baraldi, G. Serra and R. Cucchiara. "Predicting Human Eye Fixations via an LSTM-Based Saliency Attentive Model". In: *IEEE Transactions on Image Processing* 27.10 (2018), pp. 5142–5154. ISSN: 1057-7149. DOI: [10.1109/TIP.2018.2851672](https://doi.org/10.1109/TIP.2018.2851672).
- [109] Wenguan Wang, Jianbing Shen and Ling Shao. "Video salient object detection via fully convolutional networks". In: *IEEE Transactions on Image Processing* 27.1 (2018), pp. 38–49.
- [110] Tong Yubing, Faouzi Alaya Cheikh, Fahad Fazal Elahi Guraya, Hubert Konik and Alain Trémeau. "A spatiotemporal saliency model for video surveillance". In: *Cognitive Computation* 3.1 (2011), pp. 241–263.

- [111] Xin Wang, Qi Lv, Bin Wang and Liming Zhang. "Airport detection in remote sensing images: a method based on saliency map". In: *Cognitive neurodynamics* 7.2 (2013), pp. 143–154.
- [112] Patrice Y Simard, Dave Steinkraus and John C Platt. "Best practices for convolutional neural networks applied to visual document analysis". In: *null*. IEEE. 2003, p. 958.
- [113] Karen Simonyan and Andrew Zisserman. "Very deep convolutional networks for large-scale image recognition". In: *arXiv preprint arXiv:1409.1556* (2014).
- [114] M. Jiang, S. Huang, J. Duan and Q. Zhao. "SALICON: Saliency in Context". In: *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2015, pp. 1072–1080. DOI: [10.1109/CVPR.2015.7298710](https://doi.org/10.1109/CVPR.2015.7298710).
- [115] SHI Xingjian, Zhourong Chen, Hao Wang, Dit-Yan Yeung, Wai-Kin Wong and Wang-chun Woo. "Convolutional LSTM network: A machine learning approach for precipitation nowcasting". In: *Advances in neural information processing systems*. 2015, pp. 802–810.
- [116] Zoya Bylinskii, Adrià Recasens, Ali Borji, Aude Oliva, Antonio Torralba and Frédo Durand. "Where Should Saliency Models Look Next?" In: *Computer Vision - ECCV 2016 - 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part V*. 2016, pp. 809–824. DOI: [10.1007/978-3-319-46454-1_49](https://doi.org/10.1007/978-3-319-46454-1_49).
- [117] Ali Borji. "Saliency Prediction in the Deep Learning Era: An Empirical Investigation". In: *arXiv preprint arXiv:1810.03716* (2018).
- [118] Iván González-Díaz, Vincent Buso and Jenny Benois-Pineau. "Perceptual modeling in the problem of active object recognition in visual scenes". In: *Pattern Recognition* 56 (2016), pp. 129–141.
- [119] Z. Ren, S. Gao, L. T. Chia and I. W. H. Tsang. "Region-Based Saliency Detection and Its Application in Object Recognition". In: *IEEE Transactions on Circuits and Systems for Video Technology* 24.5 (2014), pp. 769–779. ISSN: 1051-8215. DOI: [10.1109/TCSVT.2013.2280096](https://doi.org/10.1109/TCSVT.2013.2280096).
- [120] T. V. Nguyen, Z. Song and S. Yan. "STAP: Spatial-Temporal Attention-Aware Pooling for Action Recognition". In: *IEEE Transactions on Circuits and Systems for Video Technology* 25.1 (2015), pp. 77–86. ISSN: 1051-8215. DOI: [10.1109/TCSVT.2014.2333151](https://doi.org/10.1109/TCSVT.2014.2333151).

- [121] X. Wang, L. Gao, J. Song and H. Shen. "Beyond Frame-level CNN: Saliency-Aware 3-D CNN With LSTM for Video Action Recognition". In: *IEEE Signal Processing Letters* 24.4 (2017), pp. 510–514. ISSN: 1070-9908. DOI: [10.1109/LSP.2016.2611485](https://doi.org/10.1109/LSP.2016.2611485).
- [122] Yu-Fei Ma, Lie Lu, Hong-Jiang Zhang and Mingjing Li. "A User Attention Model for Video Summarization". In: *Proceedings of the Tenth ACM International Conference on Multimedia*. MULTIMEDIA '02. Juan-les-Pins, France: ACM, 2002, pp. 533–542. ISBN: 1-58113-620-X. DOI: [10.1145/641007.641116](https://doi.org/10.1145/641007.641116).
- [123] Souad Chaabouni, Jenny Benois-pineau, François Tison, Chokri Ben Amar and Akka Zemmari. "Prediction of visual attention with deep CNN on artificially degraded videos for studies of attention of patients with Dementia". In: *Multimedia Tools and Applications* 76.21 (2017), pp. 22527–22546. ISSN: 1573-7721. DOI: [10.1007/s11042-017-4796-5](https://doi.org/10.1007/s11042-017-4796-5). URL: <https://doi.org/10.1007/s11042-017-4796-5>.
- [124] H. Liu and I. Heynderickx. "Visual Attention in Objective Image Quality Assessment: Based on Eye-Tracking Data". In: *IEEE Transactions on Circuits and Systems for Video Technology* 21.7 (2011), pp. 971–982. ISSN: 1051-8215. DOI: [10.1109/TCSVT.2011.2133770](https://doi.org/10.1109/TCSVT.2011.2133770).
- [125] M. A. Fernández-Torres, I. González-Díaz and F. Díaz de María. "A probabilistic topic approach for context-aware visual attention modeling". In: *2016 14th International Workshop on Content-Based Multimedia Indexing (CBMI)*. 2016, pp. 1–6. DOI: [10.1109/CBMI.2016.7500272](https://doi.org/10.1109/CBMI.2016.7500272).
- [126] M. A. Fernández-Torres, I. González-Díaz and F. Díaz de María. "Probabilistic Topic Model for Context-Driven Visual Attention Understanding". In: *IEEE Transactions on Circuits and Systems for Video Technology* (submitted).
- [127] Andrew M Derrington, John Krauskopf and Peter Lennie. "Chromatic mechanisms in lateral geniculate nucleus of macaque." In: *The Journal of physiology* 357.1 (1984), pp. 241–265.
- [128] Laurent Itti, Nitin Dhavale and Frederic Pighin. "Realistic avatar eye and head animation using a neurobiological model of visual attention". In: *Applications and Science of Neural Networks, Fuzzy Systems, and Evolutionary Computation VI*. Vol. 5200. International Society for Optics and Photonics. 2003, pp. 64–79.

- [129] C. Liu. "Beyond Pixels: Exploring New Representations and Applications for Motion Analysis". PhD thesis. Massachusetts Institute of Technology, May 2009.
- [130] Richard Szeliski. *Computer vision: algorithms and applications*. Springer Science & Business Media, 2010.
- [131] Bruce D Lucas, Takeo Kanade et al. "An iterative image registration technique with an application to stereo vision". In: (1981).
- [132] Berthold KP Horn and Brian G Schunck. "Determining optical flow". In: *Artificial intelligence* 17.1-3 (1981), pp. 185–203.
- [133] Andrés Bruhn, Joachim Weickert and Christoph Schnörr. "Lucas/Kanade Meets Horn/Schunck: Combining Local and Global Optic Flow Methods". In: *International Journal of Computer Vision* 61.3 (), pp. 211–231. ISSN: 1573-1405. DOI: [10.1023/B:VISI.0000045324.43199.43](https://doi.org/10.1023/B:VISI.0000045324.43199.43).
- [134] Thomas Brox, Andrés Bruhn, Nils Papenberg and Joachim Weickert. "High Accuracy Optical Flow Estimation Based on a Theory for Warping". In: Springer, 2004, pp. 25–36.
- [135] Richard Hartley and Andrew Zisserman. *Multiple view geometry in computer vision*. Cambridge university press, 2003.
- [136] Golnaz Abdollahian, Zygmunt Pizlo and Edward J Delp. "A study on the effect of camera motion on human visual attention". In: *2008 15th IEEE International Conference on Image Processing*. IEEE. 2008, pp. 693–696.
- [137] D. Mahapatra, S. O. Gilani and M. K. Saini. "Coherency Based Spatio-Temporal Saliency Detection for Video Object Segmentation". In: *IEEE Journal of Selected Topics in Signal Processing* 8.3 (2014), pp. 454–462. ISSN: 1932-4553. DOI: [10.1109/JSTSP.2014.2315874](https://doi.org/10.1109/JSTSP.2014.2315874).
- [138] Richard A Abrams and Shawn E Christ. "Motion onset captures attention". In: *Psychological Science* 14.5 (2003), pp. 427–432.
- [139] Juan Xu, Ming Jiang, Shuo Wang, Mohan S Kankanhalli and Qi Zhao. "Predicting human gaze beyond pixels". In: *Journal of vision* 14.1 (2014), pp. 28–28.
- [140] Farhan Baluch and Laurent Itti. "Mining videos for features that drive attention". In: *Multimedia data mining and analytics*. Springer, 2015, pp. 311–326.

- [141] Paul Viola and Michael Jones. "Rapid object detection using a boosted cascade of simple features". In: *Computer Vision and Pattern Recognition, 2001. CVPR 2001. Proceedings of the 2001 IEEE Computer Society Conference on*. Vol. 1. IEEE. 2001, pp. I-511.
- [142] Chris Harris and Mike Stephens. "A combined corner and edge detector." In: Citeseer. 1988.
- [143] Yoav Freund and Robert E Schapire. "A decision-theoretic generalization of on-line learning and an application to boosting". In: *Journal of computer and system sciences* 55.1 (1997), pp. 119-139.
- [144] Constantine P Papageorgiou, Michael Oren and Tomaso Poggio. "A general framework for object detection". In: *Computer vision, 1998. sixth international conference on*. IEEE. 1998, pp. 555-562.
- [145] Thomas Hofmann. "Unsupervised learning by probabilistic latent semantic analysis". In: *Machine learning* 42.1-2 (2001), pp. 177-196.
- [146] Li Fei-Fei, Rob Fergus and Antonio Torralba. *Recognizing and learning object categories*. ICCV short course. 2009.
- [147] Josef Sivic and Andrew Zisserman. "Efficient visual search of videos cast as text retrieval". In: *IEEE transactions on pattern analysis and machine intelligence* 31.4 (2009), pp. 591-606.
- [148] Michael I Jordan, Zoubin Ghahramani, Tommi S Jaakkola and Lawrence K Saul. "An introduction to variational methods for graphical models". In: *Machine learning* 37.2 (1999), pp. 183-233.
- [149] David M Blei. "Probabilistic topic models". In: *Communications of the ACM* 55.4 (2012), pp. 77-84.
- [150] Samuel Kim, Shrikanth Narayanan and Shiva Sundaram. "Acoustic topic model for audio information retrieval". In: *Applications of Signal Processing to Audio and Acoustics, 2009. WASPAA'09. IEEE Workshop on*. IEEE. 2009, pp. 37-40.
- [151] Jonathan K Pritchard, Matthew Stephens and Peter Donnelly. "Inference of population structure using multilocus genotype data". In: *Genetics* 155.2 (2000), pp. 945-959.
- [152] Li Fei-Fei and Pietro Perona. "A bayesian hierarchical model for learning natural scene categories". In: *Computer Vision and Pattern Recognition, 2005. CVPR 2005. IEEE Computer Society Conference on*. Vol. 2. IEEE. 2005, pp. 524-531.

- [153] David M Blei and Michael I Jordan. "Modeling annotated data". In: *Proceedings of the 26th annual international ACM SIGIR conference on Research and development in informaion retrieval*. ACM. 2003, pp. 127–134.
- [154] Xiaogang Wang, Xiaoxu Ma and WEL Grimson. "Unsupervised Activity Perception in Crowded and Complicated Scenes Using Hierarchical Bayesian Models". In: *IEEE Transactions on Pattern Analysis and Machine Intelligence* 31.3 (2009), pp. 539–555.
- [155] Shiliang Zhang, Qi Tian, Gang Hua, Qingming Huang and Shipeng Li. "Descriptive visual words and visual phrases for image applications". In: *Proceedings of the 17th ACM international conference on Multimedia*. ACM. 2009, pp. 75–84.
- [156] Shiliang Zhang, Qingming Huang, Gang Hua, Shuqiang Jiang, Wen Gao and Qi Tian. "Building contextual visual vocabulary for large-scale image applications". In: *Proceedings of the 18th ACM international conference on Multimedia*. ACM. 2010, pp. 501–510.
- [157] Zoltan Kato and Ting-Chuen Pong. "A Markov random field image segmentation model for color textured images". In: *Image and Vision Computing* 24.10 (2006), pp. 1103 –1114. ISSN: 0262-8856. DOI: <https://doi.org/10.1016/j.imavis.2006.03.005>. URL: <http://www.sciencedirect.com/science/article/pii/S0262885606001223>.
- [158] Tommi S Jaakkola and Michael I Jordan. "Bayesian parameter estimation via variational methods". In: *Statistics and Computing* 10.1 (2000), pp. 25–37.
- [159] Zoya Bylinskii, Tilke Judd, Aude Oliva, Antonio Torralba and Frédo Durand. "What do different evaluation metrics tell us about saliency models?" In: *arXiv preprint arXiv:1604.03605* (2016).
- [160] Tom Fawcett. "An introduction to ROC analysis". In: *Pattern recognition letters* 27.8 (2006), pp. 861–874.
- [161] Ali Borji, Dicky N Sihite and Laurent Itti. "Quantitative analysis of human-model agreement in visual saliency modeling: A comparative study". In: *IEEE Transactions on Image Processing* 22.1 (2013), pp. 55–69.
- [162] Jianming Zhang and Stan Sclaroff. "Exploiting surroundedness for saliency detection: a Boolean map approach". In: *IEEE Trans. Pattern Anlaysis and Machine Intellegence (TPAMI)* (2015).

- [163] Benjamin W Tatler. "The central fixation bias in scene viewing: Selecting an optimal viewing position independently of motor biases and image feature distributions". In: *Journal of vision* 7.14 (2007), pp. 4-4.
- [164] Robert J Peters, Asha Iyer, Laurent Itti and Christof Koch. "Components of bottom-up gaze allocation in natural images". In: *Vision research* 45.18 (2005), pp. 2397-2416.
- [165] Tilke Judd, Frédo Durand and Antonio Torralba. "A Benchmark of Computational Models of Saliency to Predict Human Fixations". In: *MIT Technical Report*. 2012.
- [166] Olga Russakovsky et al. "ImageNet Large Scale Visual Recognition Challenge". In: *International Journal of Computer Vision (IJCV)* 115.3 (2015), pp. 211-252. DOI: [10.1007/s11263-015-0816-y](https://doi.org/10.1007/s11263-015-0816-y).
- [167] Matthew James Stainer, Kenneth C Scott-Brown and Ben Tatler. "Looking for trouble: A description of oculomotor search strategies during live CCTV operation." In: *Frontiers in human neuroscience* 7 (2013), p. 615.
- [168] Wikimedia Commons. *CCTV control room monitor wall*. File: *CCTV control room monitor wall.jpg*. 2017. URL: https://commons.wikimedia.org/wiki/File:CCTV_control_room_monitor_wall.jpg.
- [169] Wikimedia Commons. *Cameratoezichtcentrale politie nederland*. File: *Cameratoezichtcentrale politie nederland.jpg*. 2017. URL: https://commons.wikimedia.org/wiki/File:Cameratoezichtcentrale_politie_nederland.jpg.
- [170] Gemma Graham, James D Sauer, Lucy Akehurst, Jenny Smith and Anne P Hillstrom. "CCTV observation: The effects of event type and instructions on fixation behaviour in an applied change blindness task". In: *Applied Cognitive Psychology* 32.1 (2018), pp. 4-13.
- [171] Yu-Fei Ma, Xian-Sheng Hua, Lie Lu and Hong-Jiang Zhang. "A generic framework of user attention model and its application in video summarization". In: *IEEE transactions on multimedia* 7.5 (2005), pp. 907-919.
- [172] Jiang Peng and Qin Xiao-Lin. "Keyframe-based video summary using visual attention clues". In: *IEEE MultiMedia* 2 (2009), pp. 64-73.
- [173] Naveed Ejaz, Irfan Mehmood and Sung Wook Baik. "Efficient visual attention based framework for extracting key frames from videos". In: *Signal Processing: Image Communication* 28.1 (2013), pp. 34-44.

- [174] Petros Koutras, Georgia Panagiotaropoulou, Antigoni Tsiami and Petros Maragos. "Audio-Visual Temporal Saliency Modeling Validated by fMRI Data". In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*. 2018, pp. 2000–2010.
- [175] Junwei Han, Liye Sun, Xintao Hu, Jungong Han and Ling Shao. "Spatial and temporal visual attention prediction in videos using eye movement data". In: *Neurocomputing* 145 (2014), pp. 140–153.
- [176] Yann LeCun, Yoshua Bengio and Geoffrey Hinton. "Deep learning". In: *nature* 521.7553 (2015), p. 436.
- [177] Ronan Collobert, Jason Weston, Léon Bottou, Michael Karlen, Koray Kavukcuoglu and Pavel Kuksa. "Natural language processing (almost) from scratch". In: *Journal of Machine Learning Research* 12.Aug (2011), pp. 2493–2537.
- [178] Frank Rosenblatt. "The perceptron: a probabilistic model for information storage and organization in the brain." In: *Psychological review* 65.6 (1958), p. 386.
- [179] Kurt Hornik. "Approximation capabilities of multilayer feedforward networks". In: *Neural networks* 4.2 (1991), pp. 251–257.
- [180] Djork-Arné Clevert, Thomas Unterthiner and Sepp Hochreiter. "Fast and Accurate Deep Network Learning by Exponential Linear Units (ELUs)". In: *CoRR abs/1511.07289* (2015). arXiv: 1511.07289. URL: <http://arxiv.org/abs/1511.07289>.
- [181] Léon Bottou. "Online learning and stochastic approximations". In: *On-line learning in neural networks* 17.9 (), p. 142.
- [182] Boris T Polyak. "Some methods of speeding up the convergence of iteration methods". In: *USSR Computational Mathematics and Mathematical Physics* 4.5 (1964), pp. 1–17.
- [183] Tijmen Tieleman and Geoffrey Hinton. *Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude*. Tech. rep. 2012.
- [184] Diederik P Kingma and Jimmy Ba. "Adam: A method for stochastic optimization". In: *arXiv preprint arXiv:1412.6980* (2014).
- [185] Sergey Ioffe and Christian Szegedy. "Batch normalization: Accelerating deep network training by reducing internal covariate shift". In: *arXiv preprint arXiv:1502.03167* (2015).
- [186] Jimmy Lei Ba, Jamie Ryan Kiros and Geoffrey E Hinton. "Layer normalization". In: *arXiv preprint arXiv:1607.06450* (2016).

- [187] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever and Ruslan Salakhutdinov. "Dropout: a simple way to prevent neural networks from overfitting". In: *The Journal of Machine Learning Research* 15.1 (2014), pp. 1929–1958.
- [188] David E Rumelhart, Geoffrey E Hinton and Ronald J Williams. "Learning representations by back-propagating errors". In: *nature* 323.6088 (1986), p. 533.
- [189] David H Hubel and Torsten N Wiesel. "Receptive fields, binocular interaction and functional architecture in the cat's visual cortex". In: *The Journal of physiology* 160.1 (1962), pp. 106–154.
- [190] Alex Hernández-García and Peter König. "Data augmentation instead of explicit regularization". In: *arXiv preprint arXiv:1806.03852* (2018).
- [191] Vincent Dumoulin and Francesco Visin. "A guide to convolution arithmetic for deep learning". In: *arXiv preprint arXiv:1603.07285* (2016).
- [192] Fisher Yu and Vladlen Koltun. "Multi-scale context aggregation by dilated convolutions". In: *arXiv preprint arXiv:1511.07122* (2015).
- [193] Kaiming He, Xiangyu Zhang, Shaoqing Ren and Jian Sun. "Deep residual learning for image recognition". In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016, pp. 770–778.
- [194] Yann Le Cun, Ofer Matan, Bernhard Boser, John S Denker, Don Henderson, Richard E Howard, Wayne Hubbard, LD Jacket and Henry S Baird. "Handwritten zip code recognition with multilayer networks". In: *[1990] Proceedings. 10th International Conference on Pattern Recognition*. Vol. 2. IEEE. 1990, pp. 35–40.
- [195] S. Thrun and L. Pratt. *Learning to Learn*. Springer US, 2012. ISBN: 9781461555292. URL: https://books.google.es/books?id=X_jpBwAAQBAJ.
- [196] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li and Li Fei-Fei. "Imagenet: A large-scale hierarchical image database". In: *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*. Ieee. 2009, pp. 248–255.
- [197] David G Lowe. "Distinctive image features from scale-invariant keypoints". In: *International journal of computer vision* 60.2 (2004), pp. 91–110.
- [198] Navneet Dalal and Bill Triggs. "Histograms of oriented gradients for human detection". In: *Computer Vision and Pattern Recognition, 2005. CVPR 2005. IEEE Computer Society Conference on*. Vol. 1. IEEE. 2005, pp. 886–893.

- [199] Vijay Badrinarayanan, Alex Kendall and Roberto Cipolla. "Segnet: A deep convolutional encoder-decoder architecture for image segmentation". In: *arXiv preprint arXiv:1511.00561* (2015).
- [200] Olaf Ronneberger, Philipp Fischer and Thomas Brox. "U-net: Convolutional networks for biomedical image segmentation". In: *International Conference on Medical image computing and computer-assisted intervention*. Springer. 2015, pp. 234–241.
- [201] Jonghwan Mun, Minsu Cho and Bohyung Han. "Text-Guided Attention Model for Image Captioning." In: *AAAI*. 2017, pp. 4233–4239.
- [202] Xiao-Jiao Mao, Chunhua Shen and Yu-Bin Yang. "Image restoration using convolutional auto-encoders with symmetric skip connections". In: *arXiv preprint arXiv:1606.08921* (2016).
- [203] Alex Graves. "Supervised sequence labelling". In: *Supervised sequence labelling with recurrent neural networks*. Springer, 2012, pp. 5–13.
- [204] Tomáš Mikolov, Martin Karafiát, Lukáš Burget, Jan Černocký and Sanjeev Khudanpur. "Recurrent neural network based language model". In: *Eleventh Annual Conference of the International Speech Communication Association*. 2010.
- [205] Joshua T Goodman. "A bit of progress in language modeling". In: *Computer Speech & Language* 15.4 (2001), pp. 403–434.
- [206] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich Zemel and Yoshua Bengio. "Show, attend and tell: Neural image caption generation with visual attention". In: *International conference on machine learning*. 2015, pp. 2048–2057.
- [207] Shikhar Sharma, Ryan Kiros and Ruslan Salakhudinov. "Action recognition using visual attention". In: *arXiv preprint arXiv:1511.04119* (2015).
- [208] Razvan Pascanu, Tomas Mikolov and Yoshua Bengio. "On the difficulty of training recurrent neural networks". In: *International Conference on Machine Learning*. 2013, pp. 1310–1318.
- [209] Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk and Yoshua Bengio. "Learning phrase representations using RNN encoder-decoder for statistical machine translation". In: *arXiv preprint arXiv:1406.1078* (2014).
- [210] Christopher Olah. *August 27, 2015." Understanding LSTM Networks "*. 2017.

- [211] T.S. Chande and S. Kroll. *The new technical trader: boost your profit by plugging into the latest indicators*. Wiley Finance. Wiley, 1994. ISBN: 9780471597803. URL: <https://books.google.es/books?id=uPMJAQAAMAAJ>.
- [212] SMI Eye Tracking company. <https://www.smivision.com/>. Accessed: 2018-09-20.
- [213] Jacob Benesty, Jingdong Chen, Yiteng Huang and Israel Cohen. "Pearson Correlation Coefficient". In: *Noise Reduction in Speech Processing*. Berlin, Heidelberg: Springer Berlin Heidelberg, 2009, pp. 1–4. ISBN: 978-3-642-00296-0. DOI: 10.1007/978-3-642-00296-0_5. URL: https://doi.org/10.1007/978-3-642-00296-0_5.
- [214] François Chollet et al. *Keras*. <https://keras.io>. 2015.
- [215] Martín Abadi et al. *TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems*. Software available from tensorflow.org. 2015. URL: <https://www.tensorflow.org/>.
- [216] Tie Liu, Zejian Yuan, Jian Sun, Jingdong Wang, Nanning Zheng, Xiaoou Tang and Heung-Yeung Shum. "Learning to detect a salient object". In: *IEEE Transactions on Pattern analysis and machine intelligence* 33.2 (2011), pp. 353–367.
- [217] Xavier Glorot and Yoshua Bengio. "Understanding the difficulty of training deep feedforward neural networks". In: *Proceedings of the thirteenth international conference on artificial intelligence and statistics*. 2010, pp. 249–256.
- [218] Yu Zhang and Qiang Yang. "A survey on multi-task learning". In: *arXiv preprint arXiv:1707.08114* (2017).

COLOPHON

This document was typeset using the typographical look-and-feel classicthesis developed by André Miede and Ivo Pletikosić. The style was inspired by Robert Bringhurst's seminal book on typography "*The Elements of Typographic Style*". classicthesis is available for both L^AT_EX and LyX:

<https://bitbucket.org/amiede/classicthesis/>

We gratefully acknowledge the support of NVIDIA Corporation with the donation of the NVIDIA GeForce GTX TITAN Xp GPU used for part of the research covered in this thesis.

Final Version as of 30th November 2018.