

This paper is a submitted and accepted version for publication in:
IEEE Signal Processing Letters.

Murtaza, F., Yousaf, M.H., Velastin, S. A., Qian, Y. (En prensa). End-to-End Temporal Action Detection using Bag of Discriminant Snippets (BoDS). Supplementary Material. *IEEE Signal Processing Letters*.

DOI:

©2018 IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works.

End-to-End Temporal Action Detection using Bag of Discriminant Snippets (BoDS) Supplementary Material

I. ANALYSIS AND VISUALIZATION OF KEY-SNIPPETS

In this section we provide the analysis and visualization of some key-snippets (Fig. 1) and their corresponding weights which show the discriminating power of these key-snippets. During encoding of the candidate proposals i.e. Bag of Discriminant Snippets (BoDS), the weights of each key-snippet will be treated as the contribution of this particular key-snippet in the final classification as shown in Fig. 2. Although, three out of seven snippets are assigned to the key-snippet of ‘Class 1’ its contribution is less as it has weight = 0.57 which is less than the weights of key-snippets of ‘Class 2’.

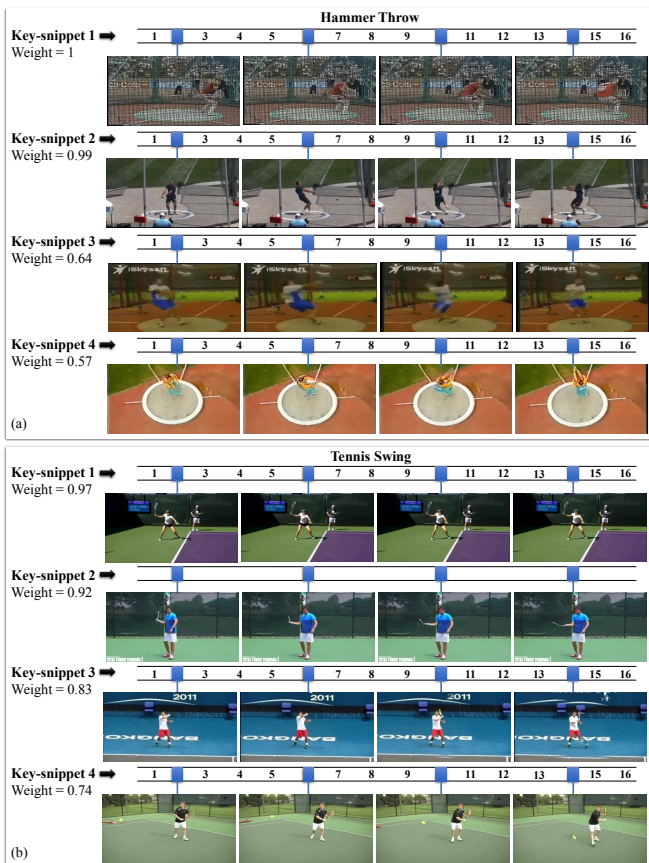


Fig. 1: Visualization of key-snippets for (a) Hammer Throw and (b) Tennis Swing action taken from validation set of Thumos14 dataset. Each key-snippet is composed of 16 frames, but we show only four frames per-key snippet for better visualization.

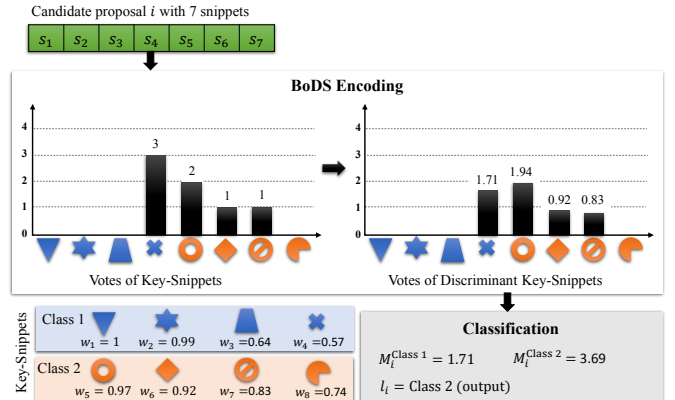


Fig. 2: Illustration of classification of a proposal i into one of the action class based on BoDS encoding. First, the snippets of the proposal i are assigned to the nearest key-snippet and votes of each key-snippet is calculated. Second, the votes of the discriminant key-snippets are calculated by multiplying with their corresponding weights. Finally, the class label of the given proposal i is obtained based upon majority voting scheme.

II. SUBSETS OF ACTIVITYNET

Table I reports results for the five top-level subsets of the ActivityNet-v1.3 dataset. We note that the subset ‘sports and exercises’ achieves the highest mAP due to well defined temporal ordering. In contrast, other subsets achieve lowest mAP as actions in these subsets have unstructured nature and temporal ordering. The names of 21 ‘sports’ actions used in the ‘Comparison’ section are given in Table II.

TABLE I: mAP(%) for five subsets of ActivityNet-v1.3. The numbers of actions in each subset are shown in brackets.

Subset	0.3	0.4	0.5	0.6	0.7
Household (45)	23.6	18.2	15.2	11.7	10.1
Personal care (19)	14.4	12.0	8.0	6.7	5.7
Eating and drinking (11)	17.4	13.2	12.0	6.6	5.0
Socializing and leisure (37)	35.6	30.0	26.0	22.6	18.9
Sports and exercises (88)	49.6	43.5	37.4	32.6	28.1

TABLE II: Names of action classes in ‘sports’ subset of ActivityNet-v1.3.

Archery, Bowling, Bungee jumping, Clean and jerk, Cricket, Curling, Discus throw, Dodgeball, Doing motocross, Hammer throw, High jump, Javelin throw, Long jump, Paintball, Playing kickball, Pole vault, Shot put, Skateboarding, Starting a campfire, Triple jump, Volleyball
