# *Supplementary Material*:
# Few-Shot Temporal Action Localization with Query Adaptive Transformer

Sauradip Nag[1,2]
s.nag@surrey.ac.uk

Xiatian Zhu [1]
xiatian.zhu@surrey.ac.uk

Tao Xiang[1,2]
t.xiang@surrey.ac.uk

[1] Centre for Vision Speech and Signal Processing (CVSSP) University of Surrey, UK

[2] iFlyTek-Surrey Joint Research Centre on Artificial Intelligence

## 1  Meta-Training Algorithm

To facilitate understanding of our method we summarize the training of task-specific snippet classification in Algorithm 1.

## 2  Ablation Study

### 2.1  Qualitative Analysis

For visual analysis, we provide two qualitative examples in Figure 1. To visualize the intra-class variation challenge which our method in particular the proposed query adaptive Transformer aims to address, we show some common examples in Figure 2.

### 2.2  Choice of Video Embedding Layer

We examine which layer of GTAD [■] is a good choice for video snippet embedding. In particular, we test five GTAD layers. The result curve in Figure 3 shows that deeper layers are usually better than shallow ones, suggesting that snippet-level contextual information is useful for action localization. We select the layer-5 as our embedding layer as it has best cost-effectiveness.

---

**Algorithm 1** Pseudo code for Snippet Classification

---

 1: **Input**: Training dataset $D_{base}$, video embedding $\mathbb{VE}$.

 2: **Output**: A query adaptive Transformer with the optimized parameters $\psi^*$.

 3: **Initialize params**: iterations $\rightarrow N_{it}$, episodes $\rightarrow N_{eps}$, shot $\rightarrow K$, epochs $\rightarrow N_{ech}$, learning rate $\rightarrow \eta_\phi, \eta_\psi$.

 4: **Training:**

 5: **for** $i = 0$ to $(N_{ech} - 1)$ **do**

 6:     **for** $j = 0$ to $(N_{eps} - 1)$ **do**

 7:         $\{S, S_{label}, Q, Q_{label}\} \leftarrow Task(\mathcal{D}_{base}, K)$

 8:         $\mathbb{F}_S \leftarrow \mathbb{VE}(S), \mathbb{F}_Q \leftarrow \mathbb{VE}(Q)$

 9:

10:         **Step 1:** ...................................................................

11:         **for** $l = 0$ to $(N_{it} - 1)$ **do**

12:             Obtain logits: $p \leftarrow h_\phi(\mathbb{F}_S)$

13:             Compute loss: $\mathcal{L}_{ce}(p; \phi)$

14:             UPDATE $\phi : \phi_{l+1} = \phi_l - \eta_\phi \nabla_\phi \mathcal{L}_{ce}$

15:         **end for**

16:         $\phi^* \leftarrow \phi_{N_{it}}$

18:         Freeze $\phi^* \rightarrow \phi^*.detach()$

19:

20:         **Step 2** ...................................................................

21:         ADAPT $\phi^*$: $\phi^{**} \leftarrow Trans(\phi^*, X_{se}^q, X_{se}^q)$

22:         Obtain logits: p' $\leftarrow h_{\phi^{**}}(\mathbb{F}_Q)$

23:         Compute loss: $\mathcal{L}_{ce}(p'; \phi^{**})$

24:         META-UPDATE $\psi$: $\psi^* = \psi - \eta_\psi \nabla_\psi \mathcal{L}_{ce}$

25:     **end for**

26:     Save the best $\psi^*$ during meta-training.

27: **end for**

---

# References

[1] Mengmeng Xu, Chen Zhao, David S. Rojas, Ali Thabet, and Bernard Ghanem. G-tad: Sub-graph localization for temporal action detection. In *CVPR*, June 2020.
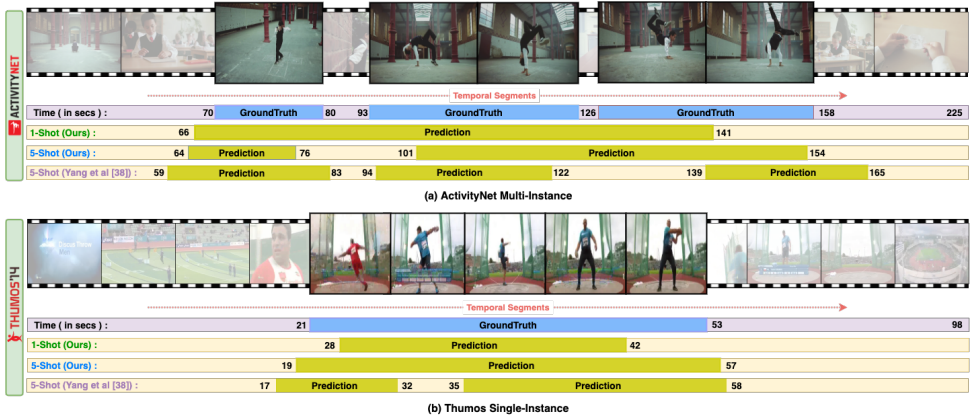
Figure 1: **Qualitative results** of (a) "BreakDancing" class on ActivityNet and (b) "Throw Discuss" class on THUMOS.
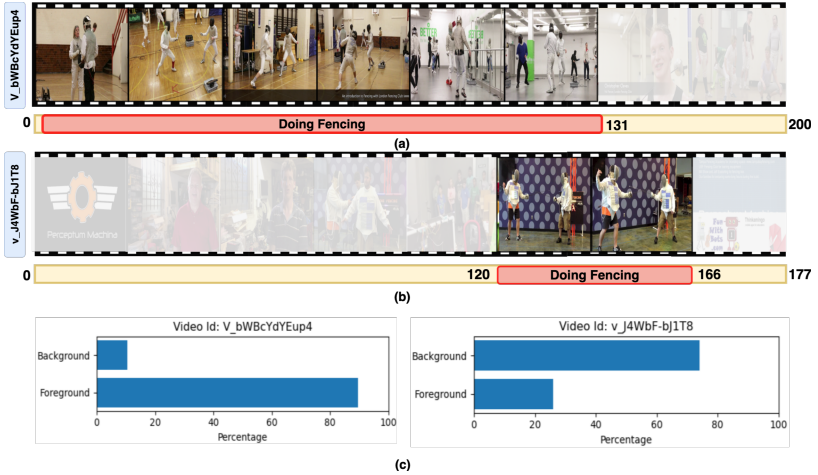


Figure 2: **Intra-class variation** example in the "Doing Fencing" class on ActivityNet-v1.3. As can be seen, the two videos present clear difference in viewpoint, scene setup, background, illumination, as well as instance length (c).
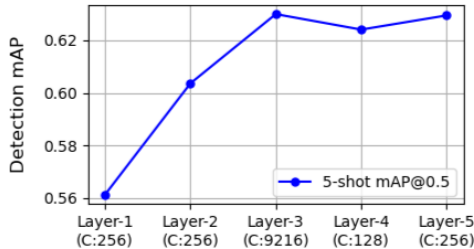


Figure 3: **Ablation of GTAD video embedding layer** in the single-instance setting on ActivityNet-v1.3. The number in round bracket is the embedding dimension.