



# Post-Processing Temporal Action Detection



**Sauradip Nag<sup>1,2</sup>**



**Xiatian Zhu<sup>1,3</sup>**



**Yi-Zhe Song<sup>1,2</sup>**



**Tao Xiang<sup>1,2</sup>**

<sup>1</sup>CVSSP, University of Surrey, UK

<sup>2</sup>IFlyTek-Surrey Joint Research Centre for Artificial Intelligence, UK

<sup>3</sup>People Centered AI, University of Surrey, UK

## Task Description : Revisiting Temporal Action Detection



Untrimmed Video

**What** is the Activity ?  
("Playing Ice Hockey")

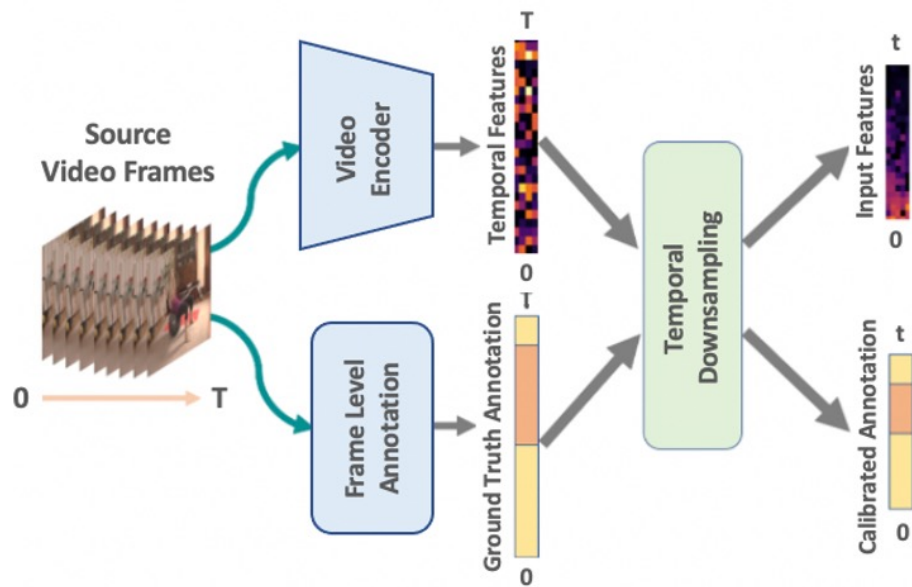
**When** is the Activity Occuring ?  
(13 s – 28 s )

**Sub-Task 1 :**  
Action **Classification**

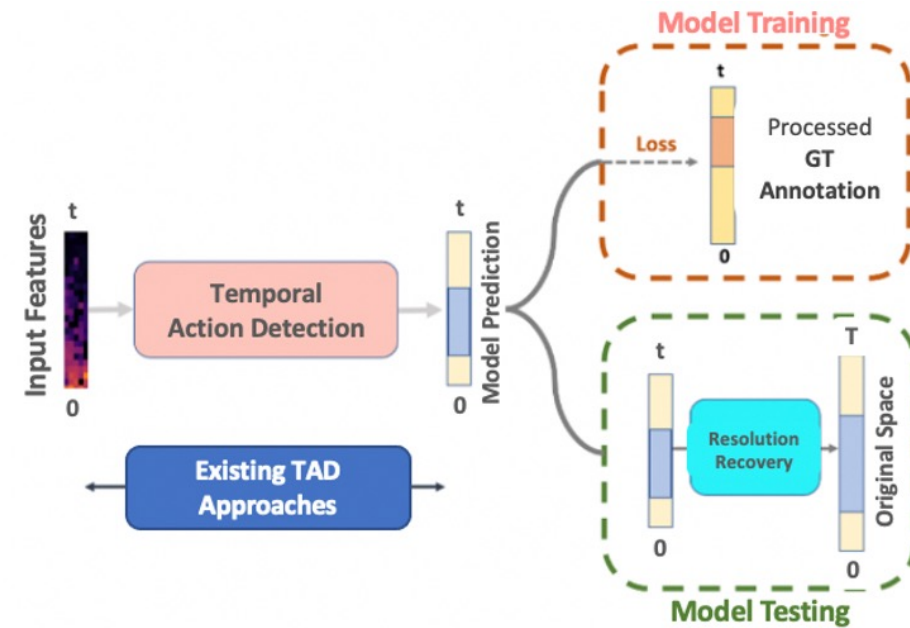
**Sub-Task 2 :**  
Temporal **Regression**



## Current Approaches: Data Pre-Processing and Model Estimation



(a) Data Preprocessing



(b) Model Estimation

Downsample the Temporal Resolution ; Detect Actions in downsampled resolution ; Inference in original resolution



## Quantization Error : How does it Arise ?



**Quantization Error** (in red) is introduced at start/end points of action due to **ceiling/floor/round operations**.

We solve this by **fitting Gaussian Distribution** at start/end points.



## Quantization Error : How to Solve ?



GT Label



Existing Snippet Level Prediction



Our Sub-Snippet Level Post-Processing



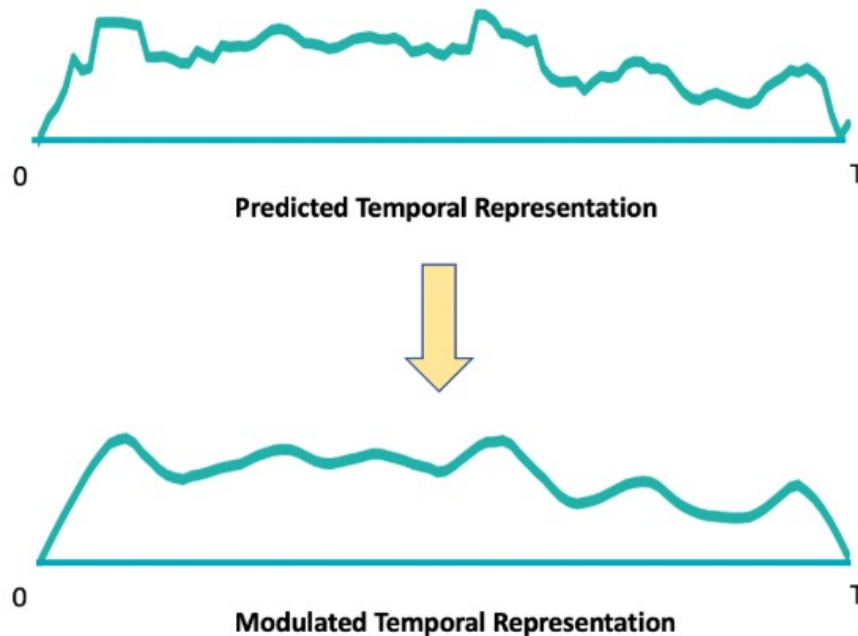
Instead of **Standard Snippet Level Prediction**  
we predict at **Sub-Snippet Level**

Each Snippet may contain action/background frame



# Gaussian Approximated Post-Processing

## Step 1: Temporal Smoothing



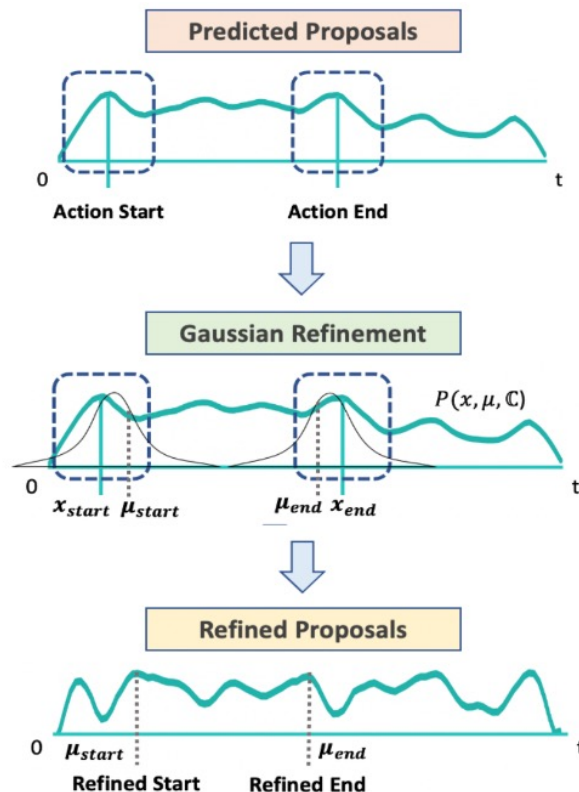
Given a Action Temporal Snippet Representation ,

We smooth the distribution using **Gaussian Kernel** to **avoid multiple peaks**.



# Gaussian Approximated Post-Processing

## Step 2: Gaussian Refinement and Calibration



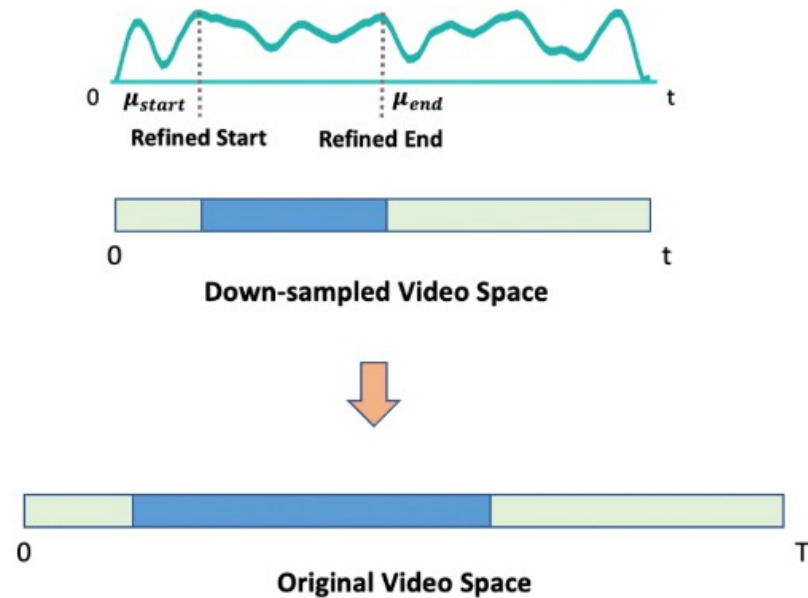
Given a predicted boundary point (start/end) :

- We fit **Gaussian Distribution** at the boundary point
- We find the **optimal Shift** using **Taylor's Expansion**
- The newly shifted start/end points are at **sub-snippet level**



# Gaussian Approximated Post-Processing

## Step 3: Temporal Resolution Recovery



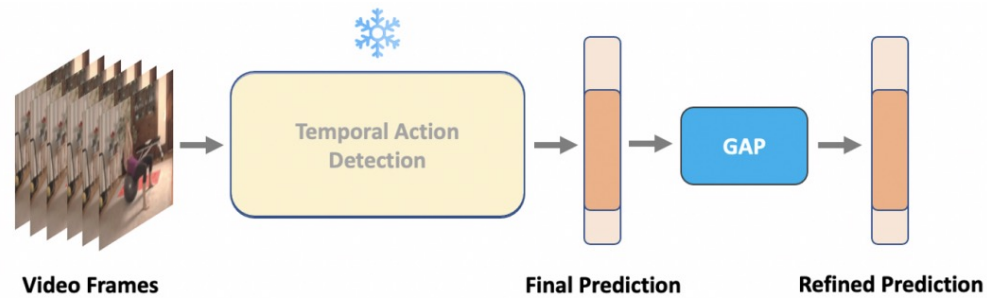
We **recover the temporal resolution** by multiplying the Video Duration

The **refined start/end point** reduces the **quantization error** at **sub-snippet level**

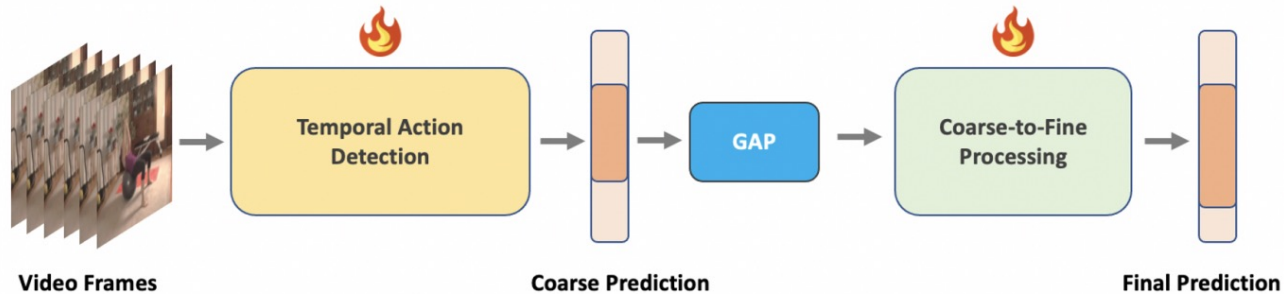




# GAP Integration : Plug-And-Play Module



(a)



(b)

Can be used during both

**(a) Model Retraining**

**(b) Inference Post-Processing**



## GAP Improves on Standard Supervised Action Detection

Category	Method	ActivityNet				THUMOS14			
		mAP				mAP			
		0.5	0.75	0.95	Avg	0.3	0.5	0.7	Avg
Anchor-based	MUSES	50.0	34.9	6.5	34.0	68.9	56.9	31.0	53.4
	<b>MUSES + GAP</b>	<b>50.3</b>	<b>35.5</b>	<b>6.9</b>	<b>34.3</b>	<b>69.3</b>	<b>57.8</b>	<b>31.9</b>	<b>53.8</b>
	PBRNet	53.9	34.9	8.9	35.0	58.5	51.3	29.5	-
	<b>PBRNet + GAP</b>	<b>54.4</b>	<b>35.4</b>	<b>9.2</b>	<b>35.2</b>	<b>59.2</b>	<b>51.9</b>	<b>30.0</b>	-
Anchor-Free	BMN	50.1	34.8	8.3	33.9	56.0	38.8	20.5	38.5
	<b>BMN + GAP</b>	<b>50.5</b>	<b>35.2</b>	<b>8.6</b>	<b>34.3</b>	<b>56.6</b>	<b>39.4</b>	<b>21.0</b>	<b>38.9</b>
	GTAD	50.4	34.6	9.0	34.1	54.5	40.2	23.4	39.3
	<b>GTAD + GAP</b>	<b>50.8</b>	<b>34.9</b>	<b>9.2</b>	<b>34.4</b>	<b>55.0</b>	<b>40.5</b>	<b>23.8</b>	<b>39.6</b>
	DCAN	51.8	35.9	9.4	35.4	68.2	54.1	32.6	-
	<b>DCAN + GAP</b>	<b>52.4</b>	<b>36.4</b>	<b>9.6</b>	<b>35.8</b>	<b>68.6</b>	<b>54.6</b>	<b>33.0</b>	-
	RTDNet	47.2	30.7	8.6	30.8	68.3	51.9	23.7	-
	<b>RTDNet + GAP</b>	<b>47.7</b>	<b>31.1</b>	<b>8.8</b>	<b>31.2</b>	<b>68.8</b>	<b>52.3</b>	<b>24.2</b>	-
	AFSD	52.4	35.3	6.5	34.4	67.3	55.5	31.1	52.0
	<b>AFSD + GAP</b>	<b>53.0</b>	<b>35.9</b>	<b>7.1</b>	<b>34.8</b>	<b>68.0</b>	<b>56.1</b>	<b>31.5</b>	<b>52.5</b>
	ActionFormer	53.5	36.2	8.2	35.6	82.1	71.0	43.9	66.8
<b>ActionFormer + GAP</b>	<b>53.9</b>	<b>36.4</b>	<b>8.5</b>	<b>36.0</b>	<b>82.3</b>	<b>71.4</b>	<b>44.2</b>	<b>66.9</b>	
React	-	-	-	-	69.2	57.1	35.6	55.0	
<b>React + GAP</b>	-	-	-	-	<b>69.5</b>	<b>57.3</b>	<b>35.7</b>	<b>55.2</b>	
Proposal-Free	TAGS	56.3	36.8	9.6	36.5	68.6	57.0	31.8	52.8
	<b>TAGS + GAP</b>	<b>56.7</b>	<b>37.2</b>	<b>9.8</b>	<b>36.7</b>	<b>69.1</b>	<b>57.4</b>	<b>32.0</b>	<b>53.0</b>

SOTA on standard benchmarks  
Like **THUMOS14** and **ActivityNet**



# GAP Also Improves on Any-Shot/ Any-Supervision Action Detection

## GAP on Weakly Supervised Setting

Model	mAP					Avg
	0.3	0.4	0.5	0.6	0.7	
ASL [5]	51.8	-	31.1	-	11.4	32.2
ASL [5] + GAP	<b>53.0</b>	-	<b>31.7</b>	-	<b>11.5</b>	<b>32.4</b>
CoLA [13]	51.5	41.9	32.2	22.0	13.1	40.9
CoLA [13] + GAP	<b>51.8</b>	<b>42.2</b>	<b>32.4</b>	<b>22.2</b>	<b>13.2</b>	<b>41.0</b>
TS-PCA [4]	52.4	43.5	34.6	23.7	12.6	-
TS-PCA [4] + GAP	<b>52.9</b>	<b>44.0</b>	<b>34.9</b>	<b>24.0</b>	<b>12.8</b>	-
CO2-Net [3]	54.5	45.7	38.3	26.4	13.4	-
CO2-Net [3] + GAP	<b>54.9</b>	<b>46.0</b>	<b>38.8</b>	<b>27.1</b>	<b>14.0</b>	-
ASM-Loc [2]	57.1	46.8	36.6	25.2	13.4	45.1
ASM-Loc [2] + GAP	<b>58.1</b>	<b>47.5</b>	<b>37.1</b>	<b>25.6</b>	<b>13.8</b>	<b>45.5</b>
DELU [1]	56.5	47.7	40.5	27.2	15.3	46.4
DELU [1] + GAP	<b>57.0</b>	<b>48.1</b>	<b>40.9</b>	<b>27.6</b>	<b>15.5</b>	<b>46.6</b>

## GAP on Few-Shot Setting

Shot	Models	mAP			
		0.5	0.7	0.9	Avg
1	QAT [10]	44.9	29.2	11.2	25.9
	QAT [10] + GAP	<b>45.8</b>	<b>30.0</b>	<b>11.8</b>	<b>26.5</b>
5	QAT [10]	51.8	32.6	11.9	30.2
	QAT [10] + GAP	<b>52.2</b>	<b>32.9</b>	<b>12.1</b>	<b>30.4</b>

## GAP on Zero-Shot Setting

Models	mAP			
	0.5	0.75	0.95	Avg
Baseline	28.0	16.4	1.2	16.0
Baseline <sup>†</sup> + GAP	<b>28.7</b>	<b>16.8</b>	<b>1.7</b>	<b>16.5</b>
Baseline + GAP	<b>28.2</b>	<b>16.6</b>	<b>1.3</b>	<b>16.2</b>
STALE [9]	32.1	20.7	5.9	20.5
STALE [9] + GAP	<b>32.4</b>	<b>21.1</b>	<b>6.2</b>	<b>20.8</b>



## GAP is effective when Temporal Resolution is Small

Method	Temporal Resolution	mAP			
		0.5	0.75	0.95	Avg
BMN [9]	25	44.7	27.9	7.0	28.1
<b>BMN+GAP</b>	25	<b>45.5</b>	<b>28.4</b>	<b>7.3</b>	<b>28.5</b>
BMN	100	50.1	34.8	8.3	33.9
<b>BMN+GAP</b>	100	<b>50.5</b>	<b>35.2</b>	<b>8.6</b>	<b>34.3</b>
BMN	400	50.9	34.9	8.1	34.0
<b>BMN+GAP</b>	400	<b>51.1</b>	<b>35.0</b>	<b>8.2</b>	<b>34.1</b>

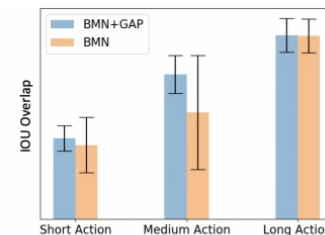
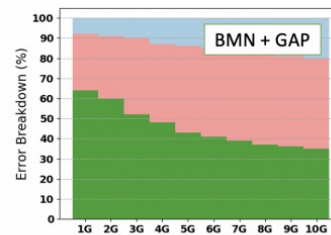
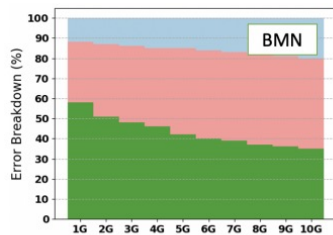
As Temporal Resolution is Low,  
**More Chances of Ambiguity** among Action Boundaries

**GAP is effective** in such scenarios

As Temporal Resolution Increases,  
Duration per snippet increases, **Ambiguity reduces**

**GAP is less effective** in such scenarios

Error-Sensitivity Analysis shows that



## GAP does not bottleneck efficiency

Method	Inference Time	Speed
AFSD [8]	0.29 sec	1931 FPS
AFSD + <b>GAP</b>	0.31 sec	1792 FPS

GAP **does not** learn any extra parameters.

GAP may slightly **increase the training time** but negligible.

GAP has **almost negligible effect** on inference time.





Scan the QR Code for Code

## Thank You for Listening

For any questions, contact : [s.nag@surrey.ac.uk](mailto:s.nag@surrey.ac.uk)

