

# Action Dubber: Timing Audible Actions via Inflectional Flow

Wenlong Wan<sup>1,2</sup> Weiying Zheng<sup>3</sup> Tianyi Xiang<sup>4,2,1</sup> Guiqing Li<sup>1</sup> Shengfeng He<sup>2</sup>

<sup>1</sup>South China University of Technology

<sup>2</sup>Singapore Management University

<sup>3</sup>University of Hong Kong

<sup>4</sup>City University of Hong Kong

Video without audio



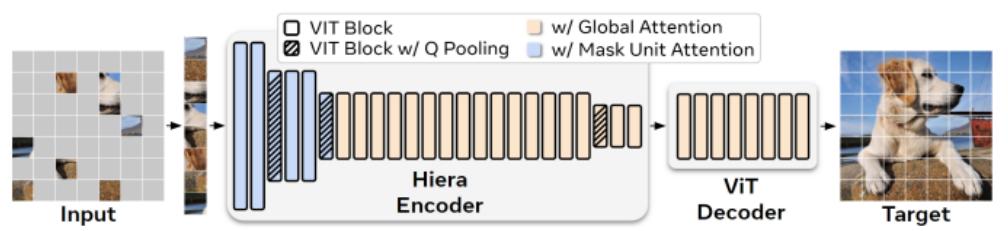
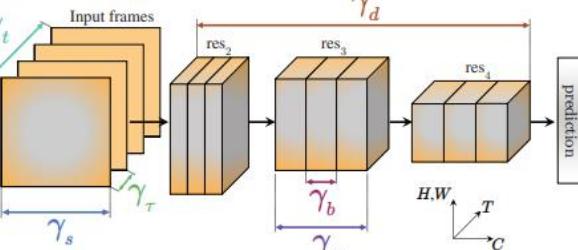
*Problem Definition*

Video without audio



*Problem Definition*

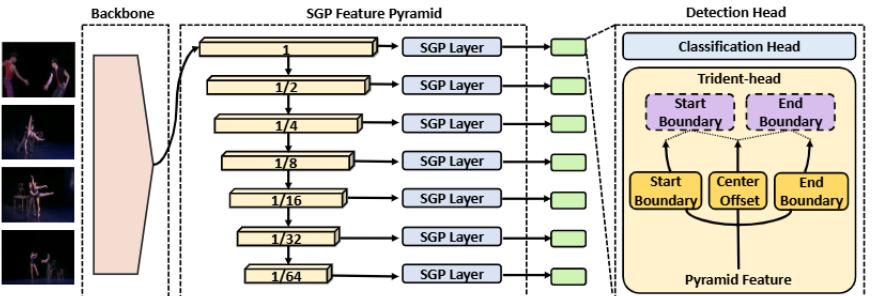
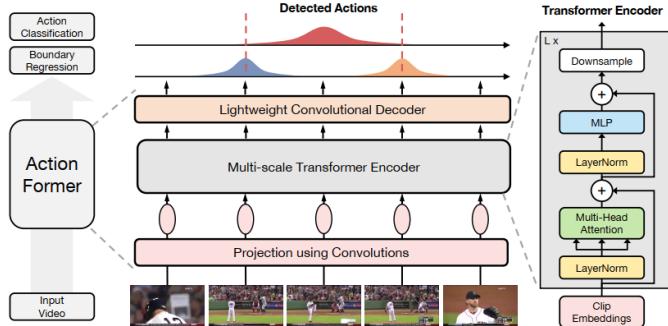
## Video Action Recognition



## Silent Video with Audible Actions



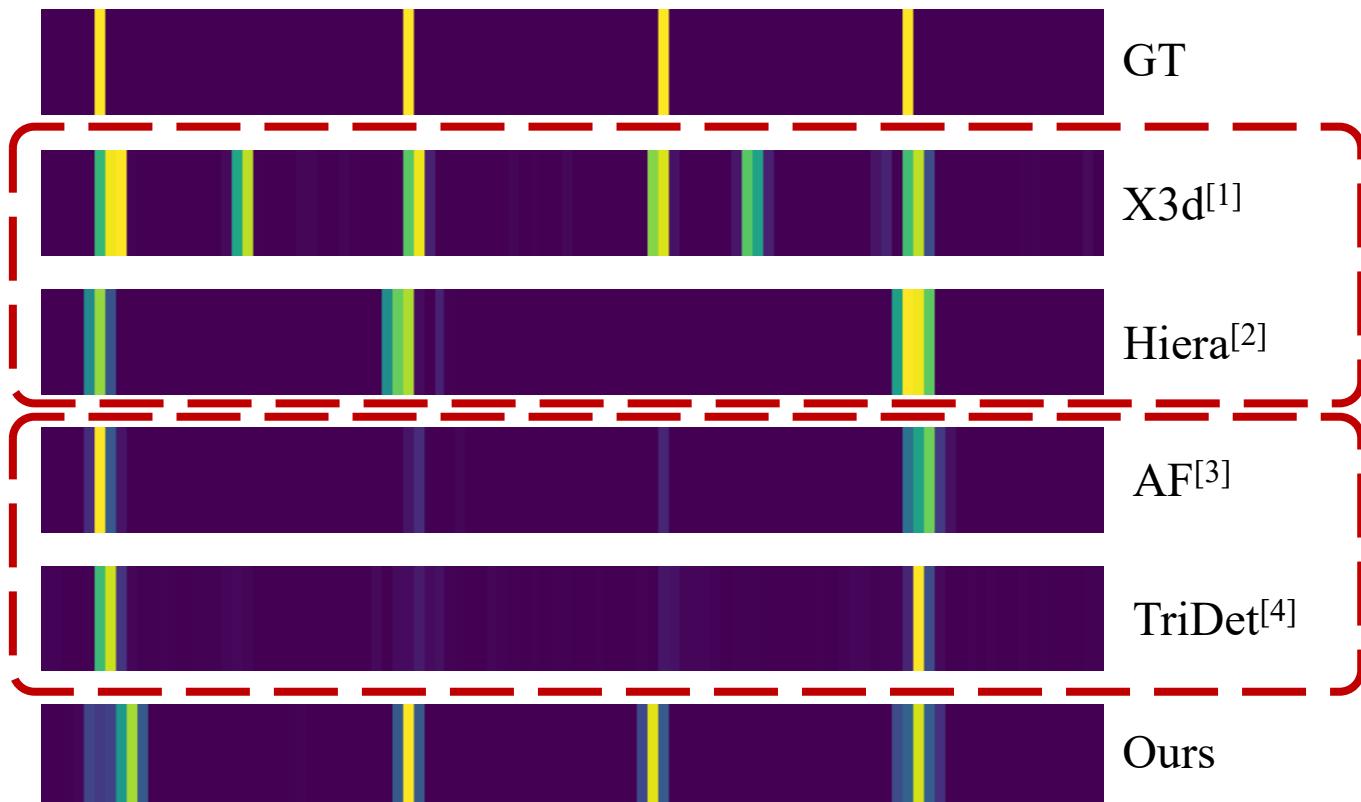
## Temporal Action Localization



ActionFormer<sup>[3]</sup> (ECCV 22')

TriDet<sup>[4]</sup> (CVPR 23')

## Temporal Audio Position



## Motivation

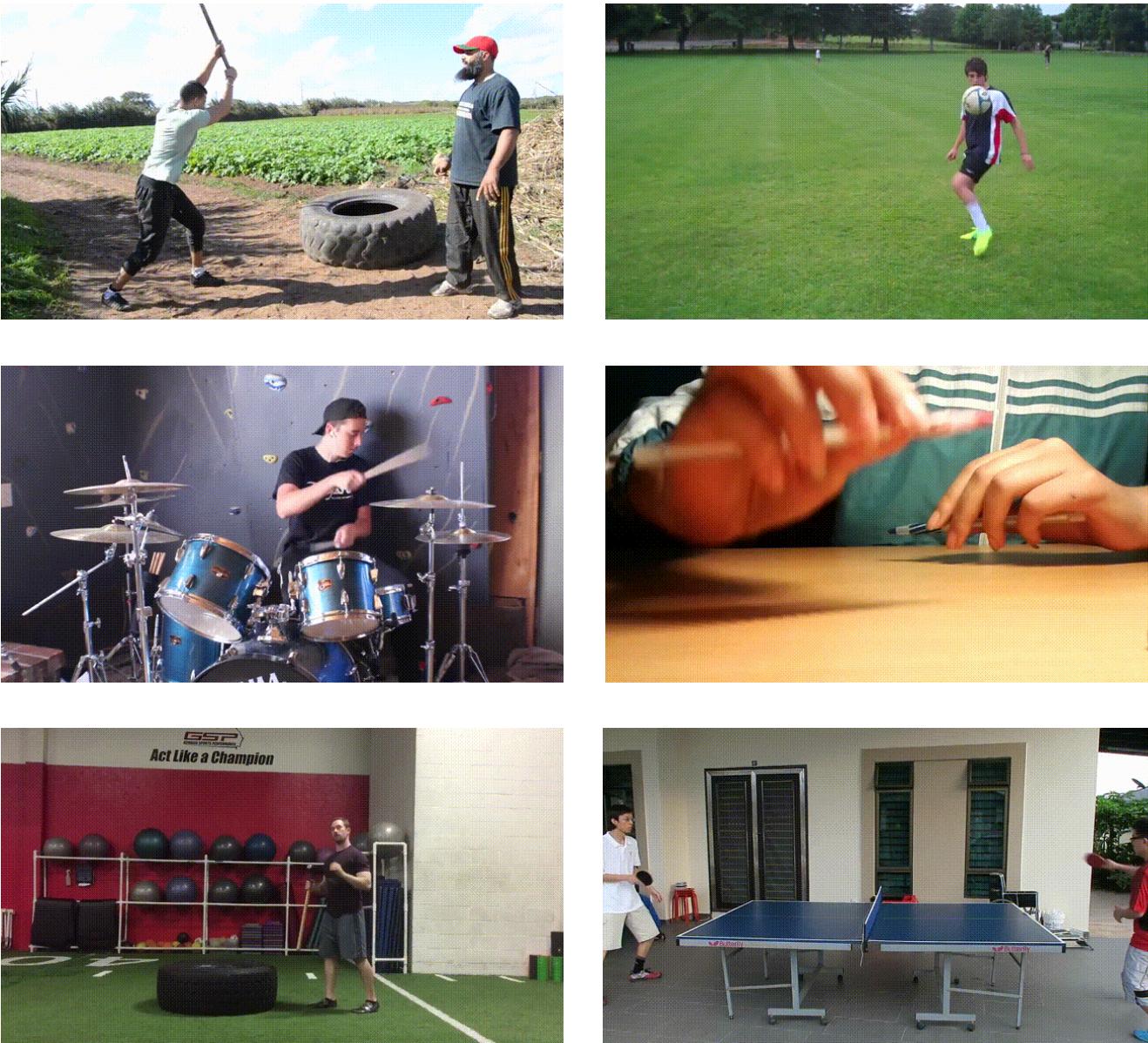
[1] Feichtenhofer, C. (2020). X3d: Expanding architectures for efficient video recognition. In *CVPR*. (pp. 203-213).

[2] Ryali, C., Hu, Y. T., Bolya, D., Wei, C., Fan, H., Huang, P. Y., ... & Feichtenhofer, C. (2023). Hiera: A Hierarchical Vision Transformer without the Bells-and-Whistles. arXiv preprint arXiv:2306.00989.

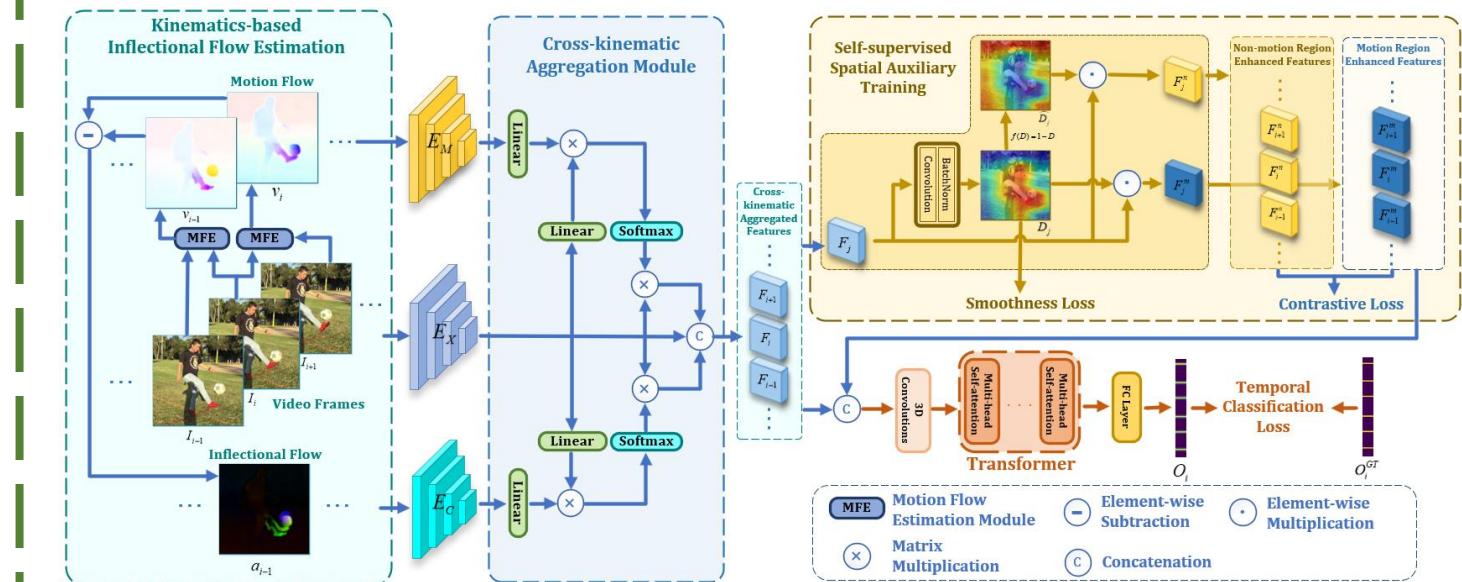
[3] Zhang, C. L., Wu, J., & Li, Y. (2022, October). Actionformer: Localizing moments of actions with transformers. In *ECCV*. (pp. 492-510).

[4] Shi, D., Zhong, Y., Cao, Q., Ma, L., Li, J., & Tao, D. (2023). TriDet: Temporal action detection with relative boundary modeling. In *CVPR*. (pp. 18857-18866).

## *Audible623 Dataset*



## *TA<sup>2</sup>Net*



# Contributions

## Basic Assumption

Audible actions are induced by sudden changes in the forces acting upon an object (  ).



Forces acted on the ball:  
Gravity only

Forces acted on the ball:  
Gravity + Sudden normal force  
(Key frame of action)

Forces acted on the ball:  
Gravity only

\*Air resistance force is too small so it can be ignored.



Discriminative maps

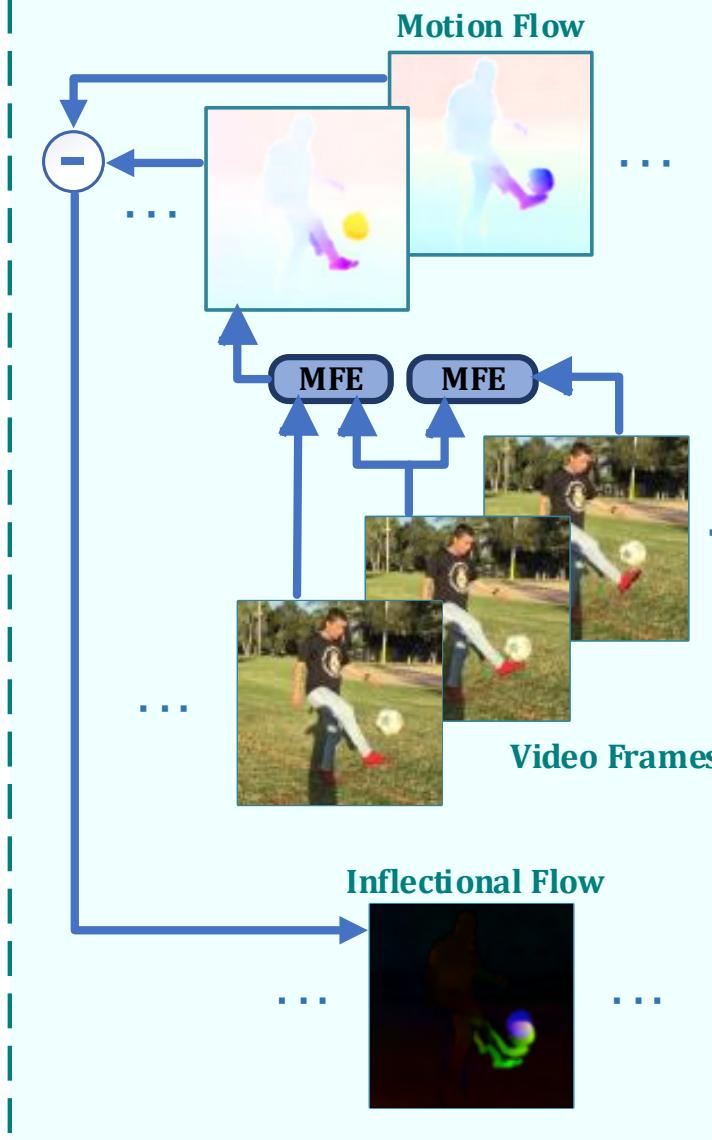
Spatial Modelling (Motion Localization)



Self-supervised  
Spatial Auxiliary  
Training Strategy

*Motion Association*

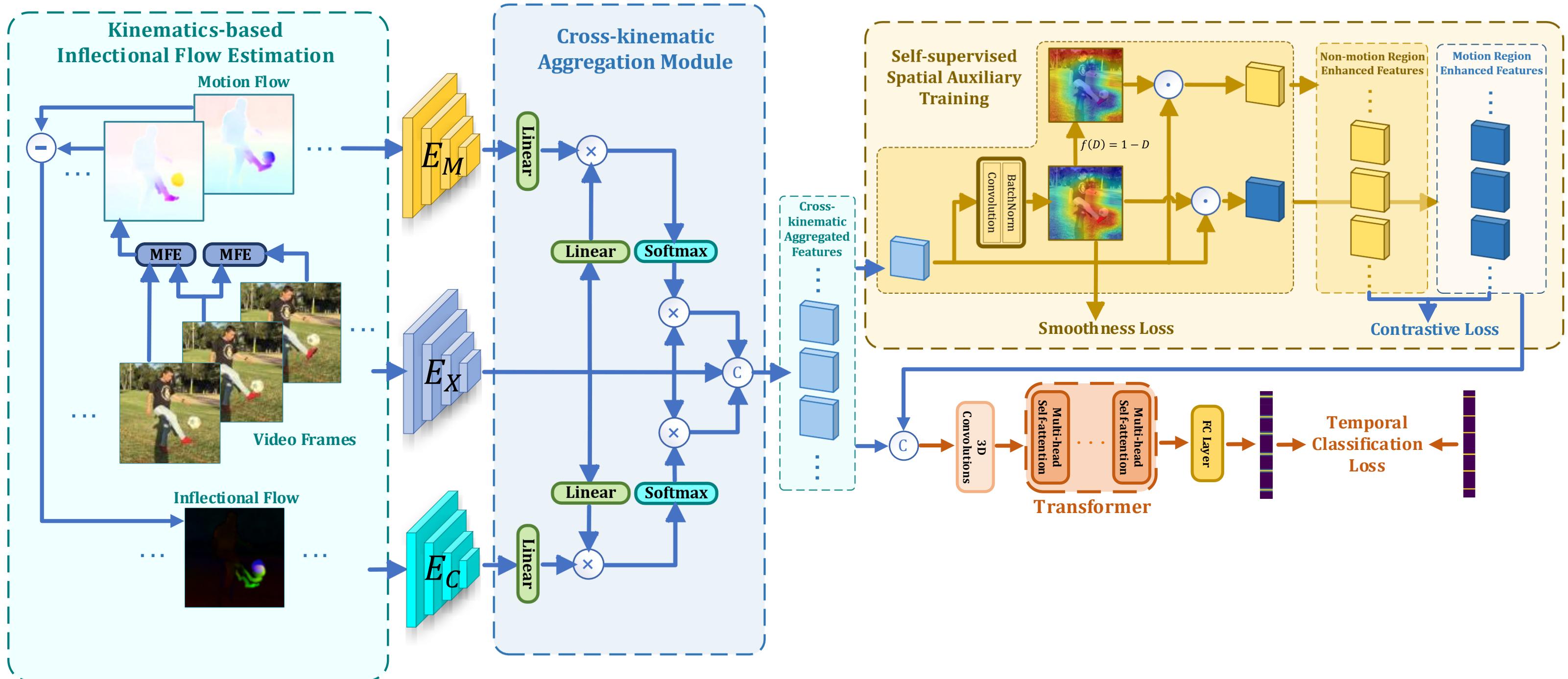
## Kinematics-based Inflectional Flow Estimation

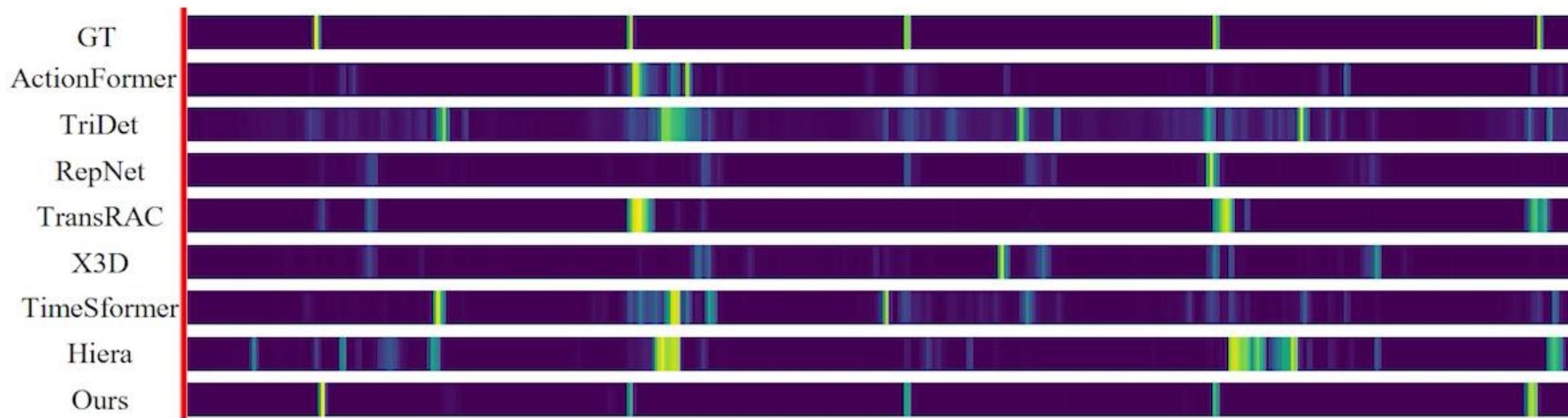


# Overview

*Inflectional Flow Estimation*

# Overview



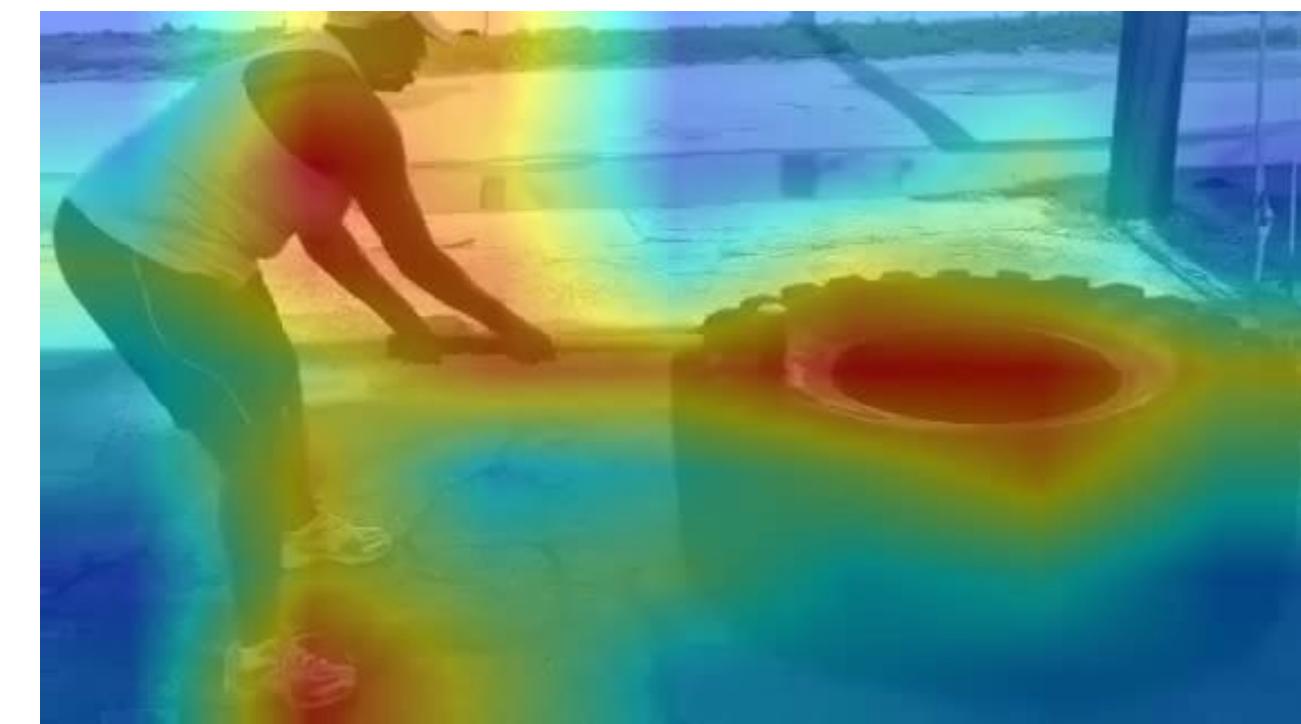
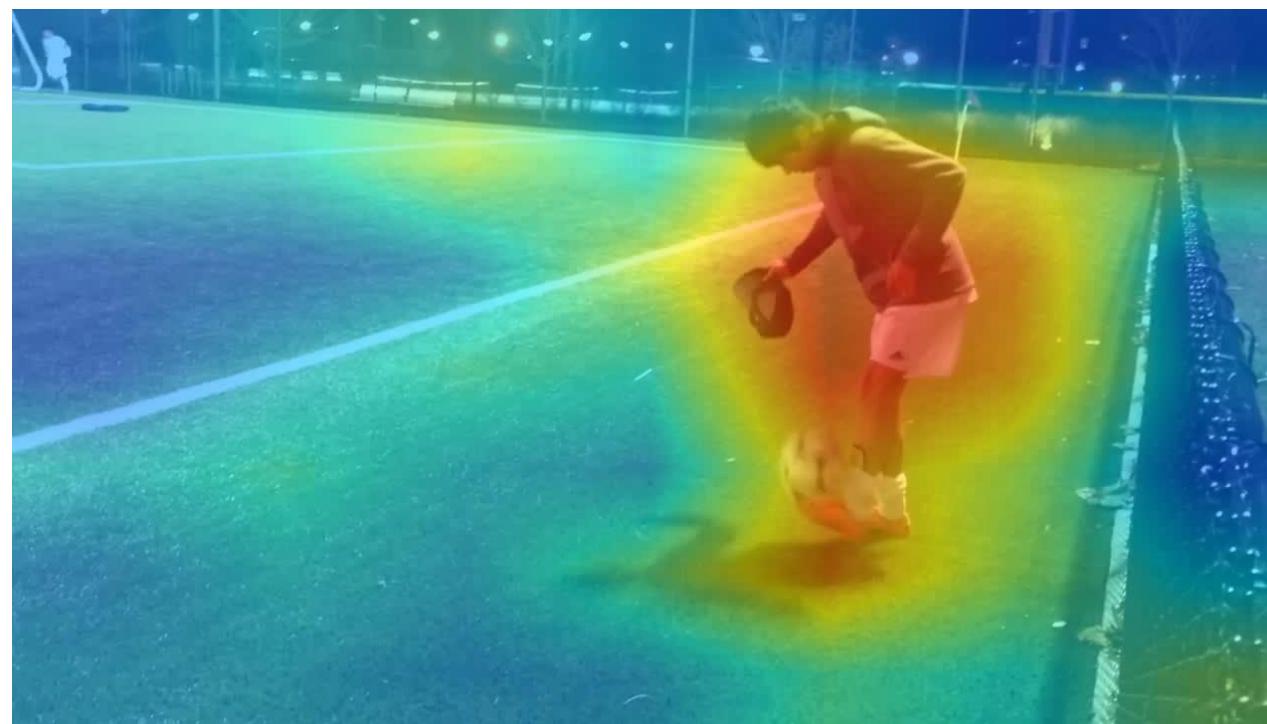


*Temporal Probabilities Predictions*

Method	UCFRep		CountixAV	
	MAE↓	OBO↑	MAE↓	OBO↑
RepNet	0.915	0.074	0.749	0.231
TransRAC	<u>0.594</u>	<b>0.222</b>	<u>0.686</u>	<u>0.255</u>
X3d	1.245	0.037	0.876	0.192
TimeSformer	0.832	0.0	1.551	0.185
Hiera	0.791	0.152	3.648	0.231
<b>Ours</b>	<b>0.588</b>	<u>0.185</u>	<b>0.549</b>	<b>0.346</b>

Table 3: Counting performance on UCFRep and CountixAV datasets. Top-2 results are marked in **bold** and underlined. Our method has not been trained on any counting dataset.

*Counting Performance Evaluation*

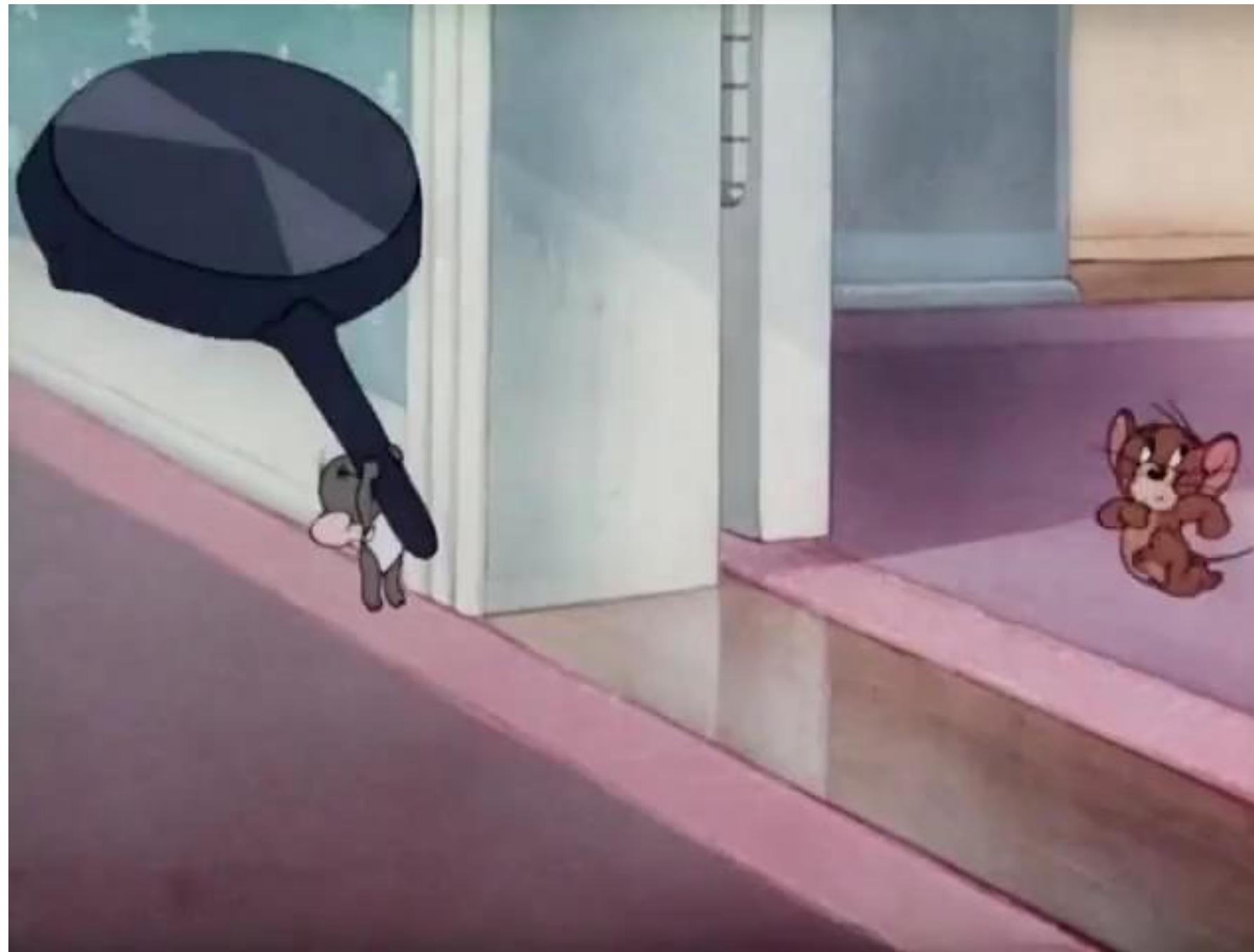


***Self-supervised Spatial Localization Visualization***

# Video Re-dubbing Result



*Video Re-dubbing Result*



*Video Re-dubbing Result*

Thanks for watching!