



哈爾濱工業大學(深圳)  
HARBIN INSTITUTE OF TECHNOLOGY, SHENZHEN



鹏城实验室  
PENG CHENG LABORATORY



# MambaVLT: Time-Evolving Multimodal State Space Model for Vision-Language Tracking

CVPR 2025 Highlight



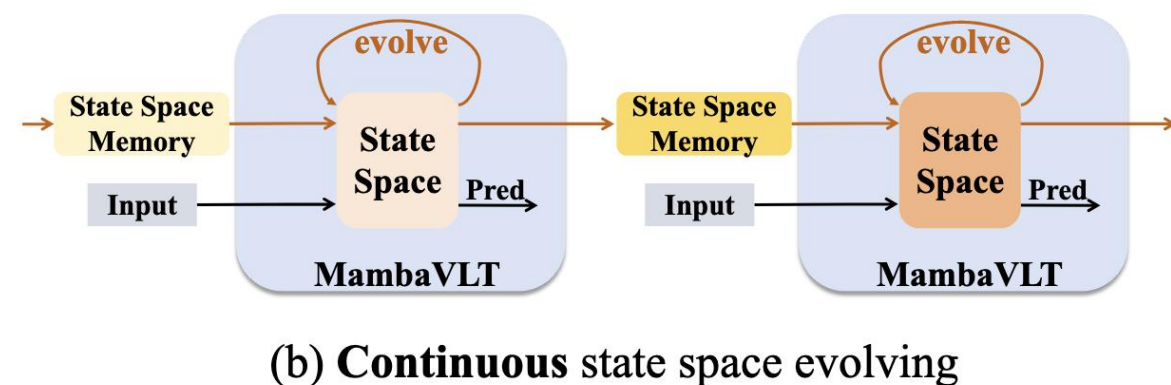
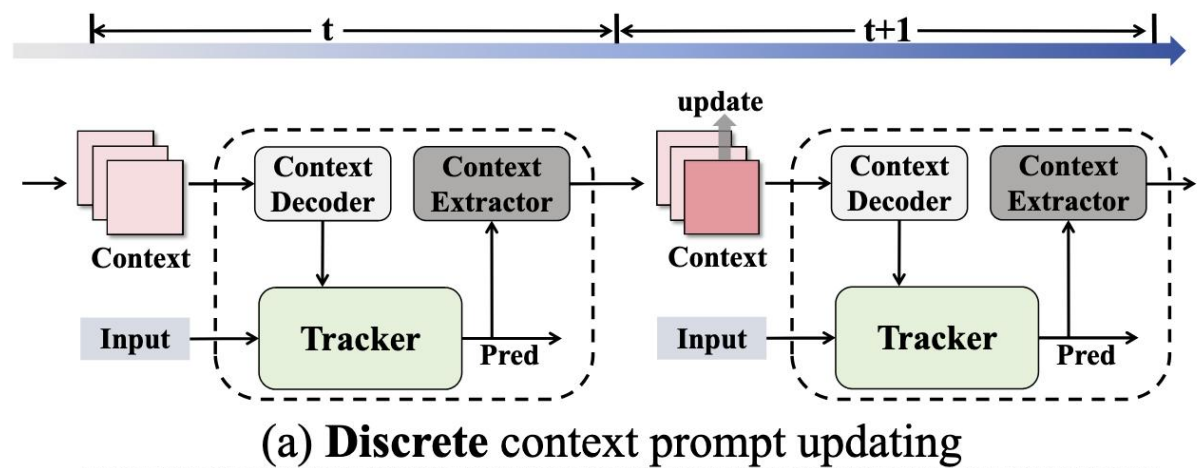
Xinqi Liu<sup>1,†</sup>, Li Zhou<sup>1,†</sup>, Zikun Zhou<sup>2,\*</sup>, Jianqiu Chen<sup>1</sup>, and Zhenyu He<sup>1,\*</sup>

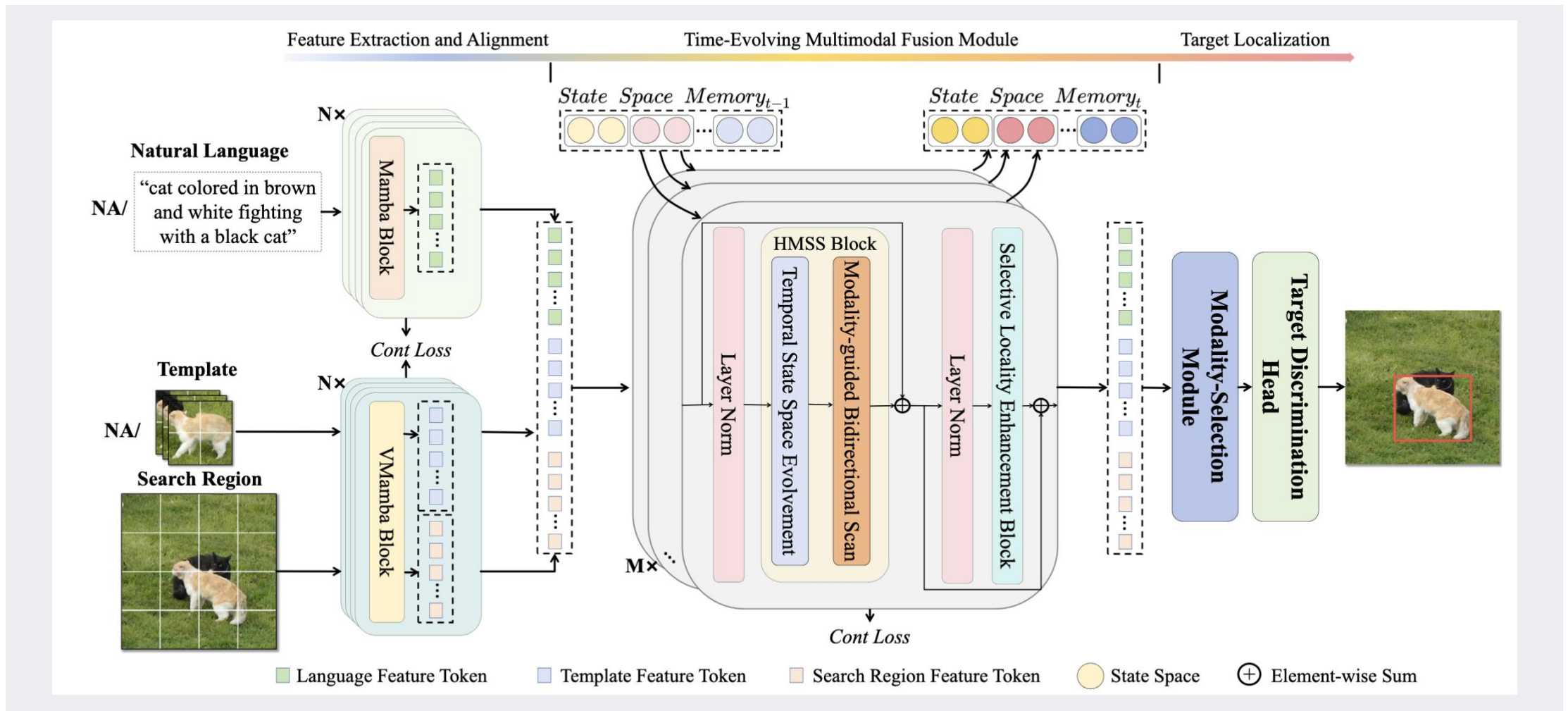
1 Harbin Institute of Technology, Shenzhen 2 Pengcheng Laboratory

†Equal Contribution \*Corresponding Author

# Motivation

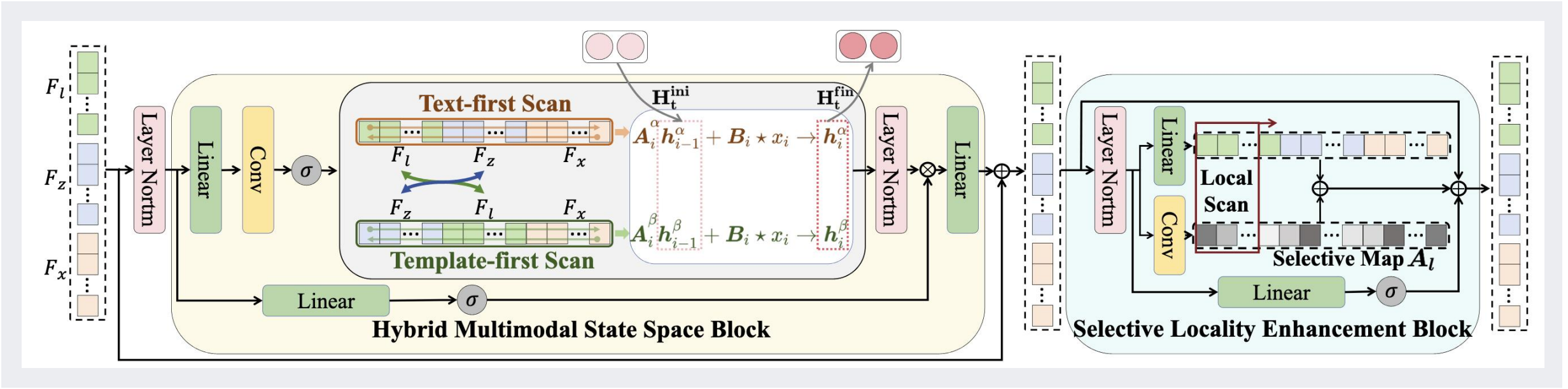
- Traditional vision-language trackers update references in **separate, discrete steps** which hinges on accurate predictions, causing error accumulation and an inability to fully exploit temporal cues. To overcome these limitations, we propose MambaVLT which adopts a **continuous, time-evolving state space mechanism** that can retain and update multimodal reference features in the video with linear complexity.





## Method

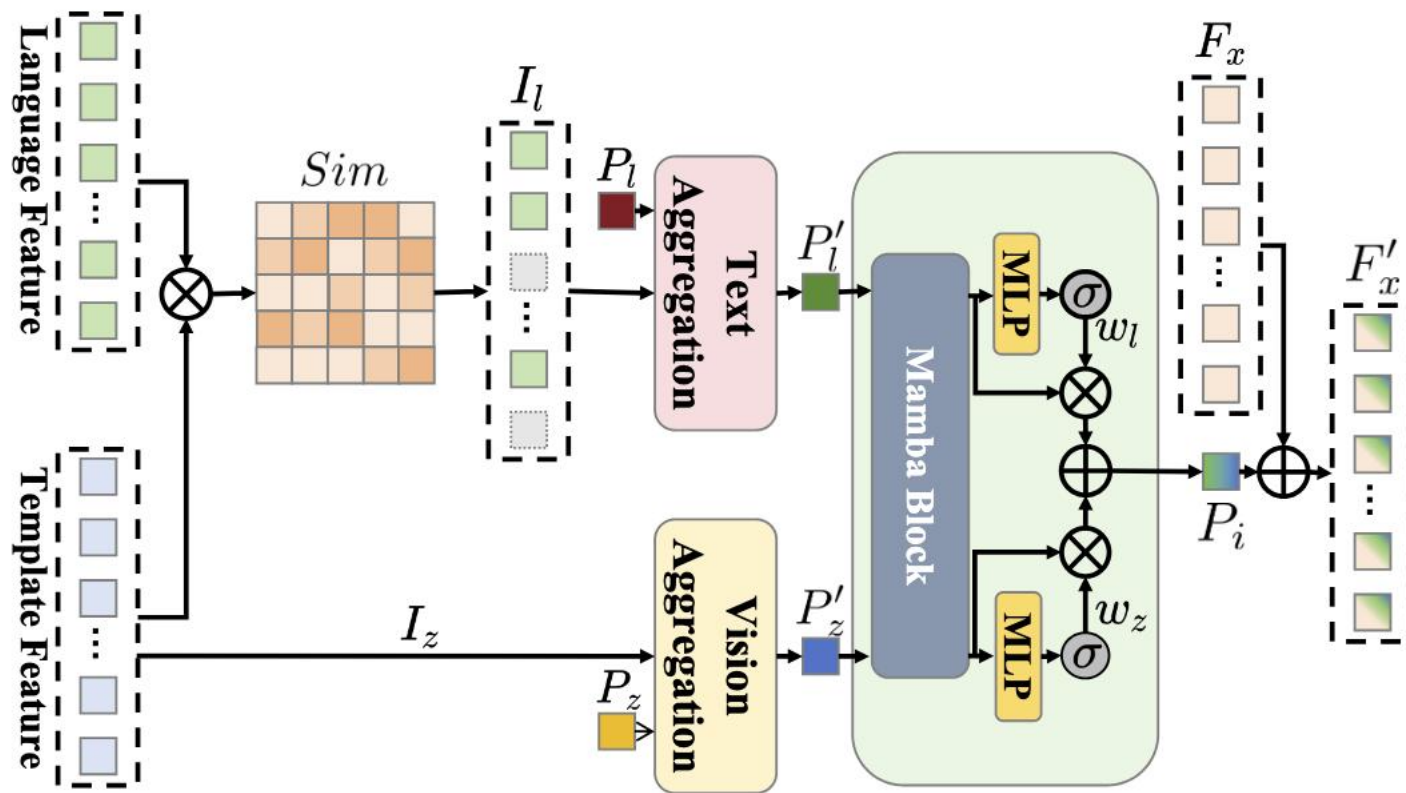
- Overview of the MambaVLT. Given various modality reference settings, features are initially extracted and aligned, then forwarded to the time-evolving multimodal fusion module. Subsequently, these features are input into the localization module to obtain precise localization information.



## Method

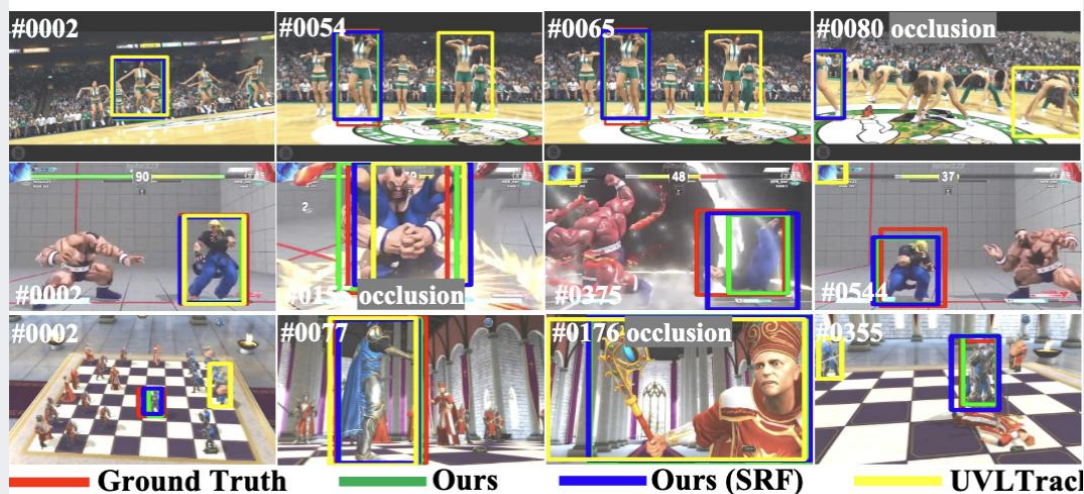
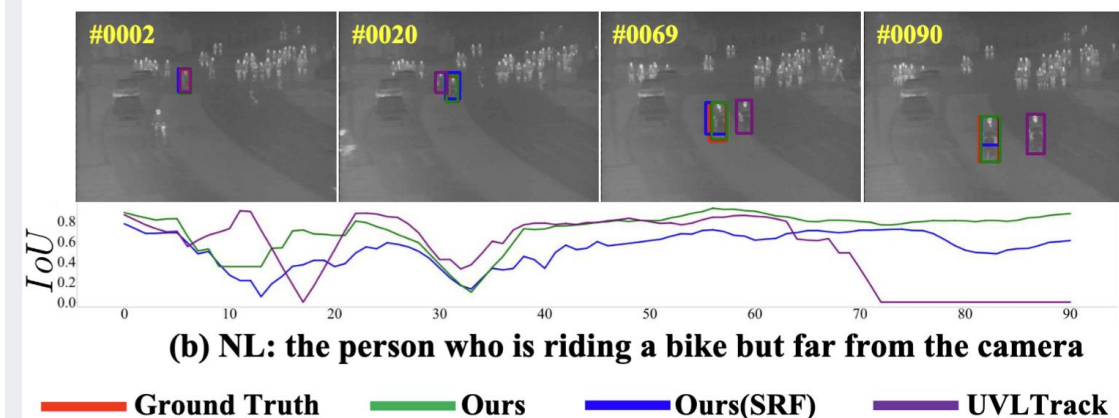
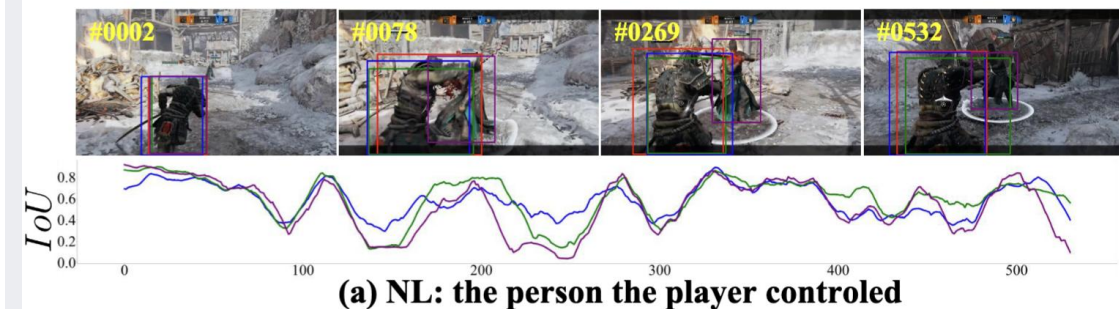
- **Time-Evolving Multimodal Fusion**
  - **Hybrid State-Space Block** captures long-term temporal information by the time-evolving state space, based on which it models multimodal features and updates target reference information by a hybrid scan.
  - **Selective Locality Enhancement Block** performs a sliding window scan with a selective map to enhance multimodal features of the current time stamp through a global receptiveness.

# Methods

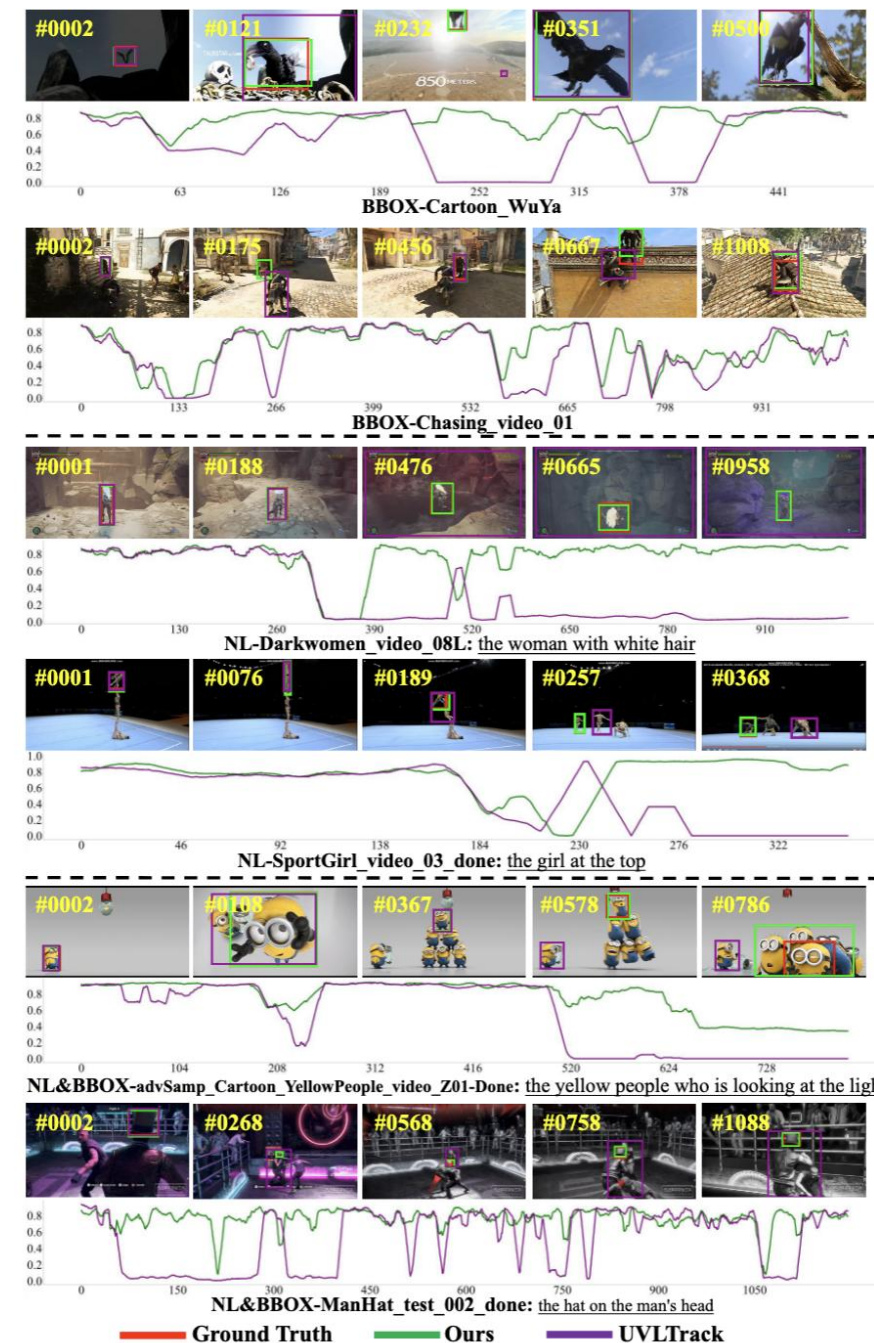


- In different tracking frames, the reliability of the language and template features may vary due to the target motion and appearance changes. Therefore, We further employ a **Modality-Selection Module** to selectively fuse multimodal reference features for search region feature refining.





## Results





哈爾濱工業大學(深圳)  
HARBIN INSTITUTE OF TECHNOLOGY, SHENZHEN



鹏城实验室  
PENG CHENG LABORATORY



*Thanks for Your Watching!*