

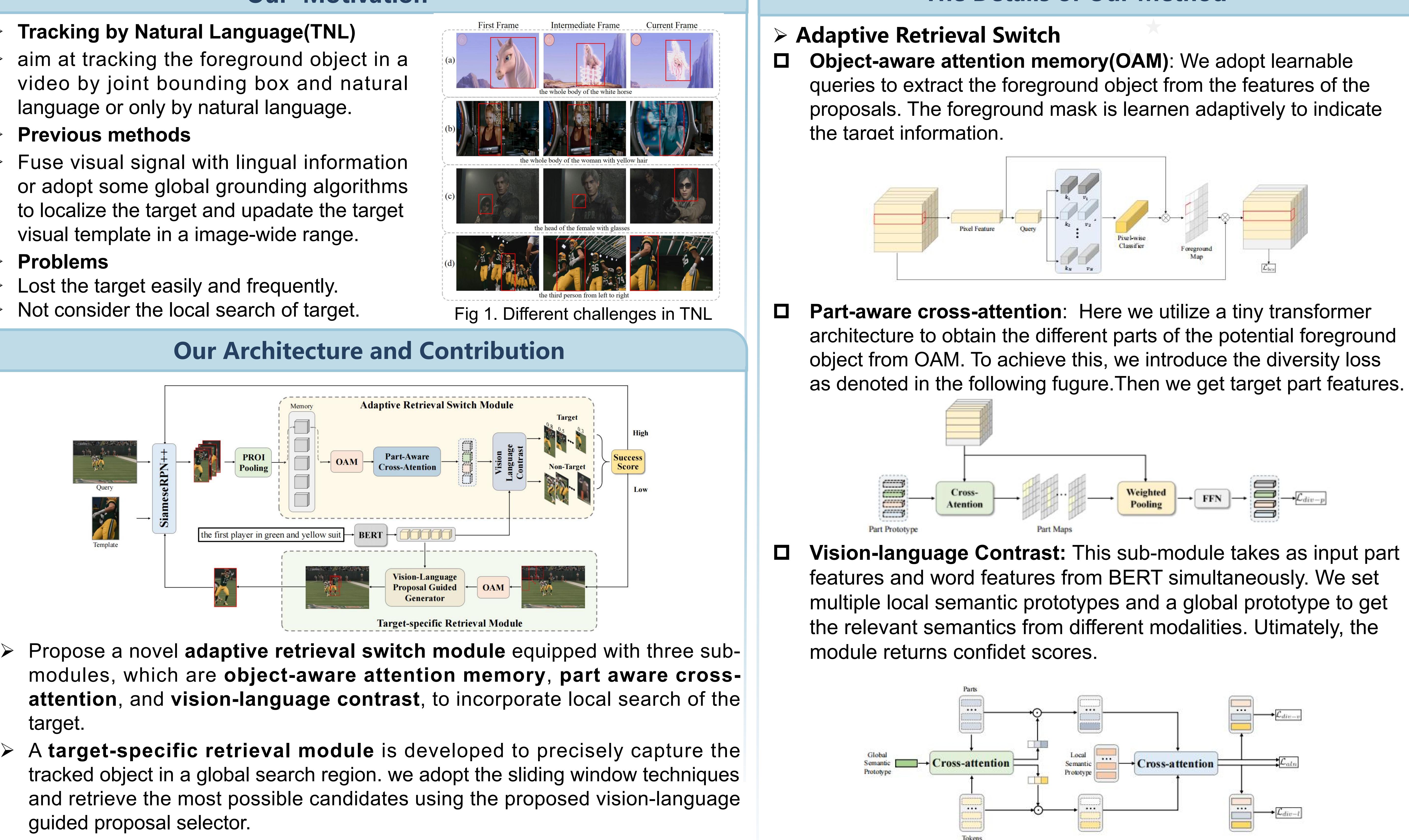
# **Cross-modal Target Retrieval for Tracking by Natural Language**

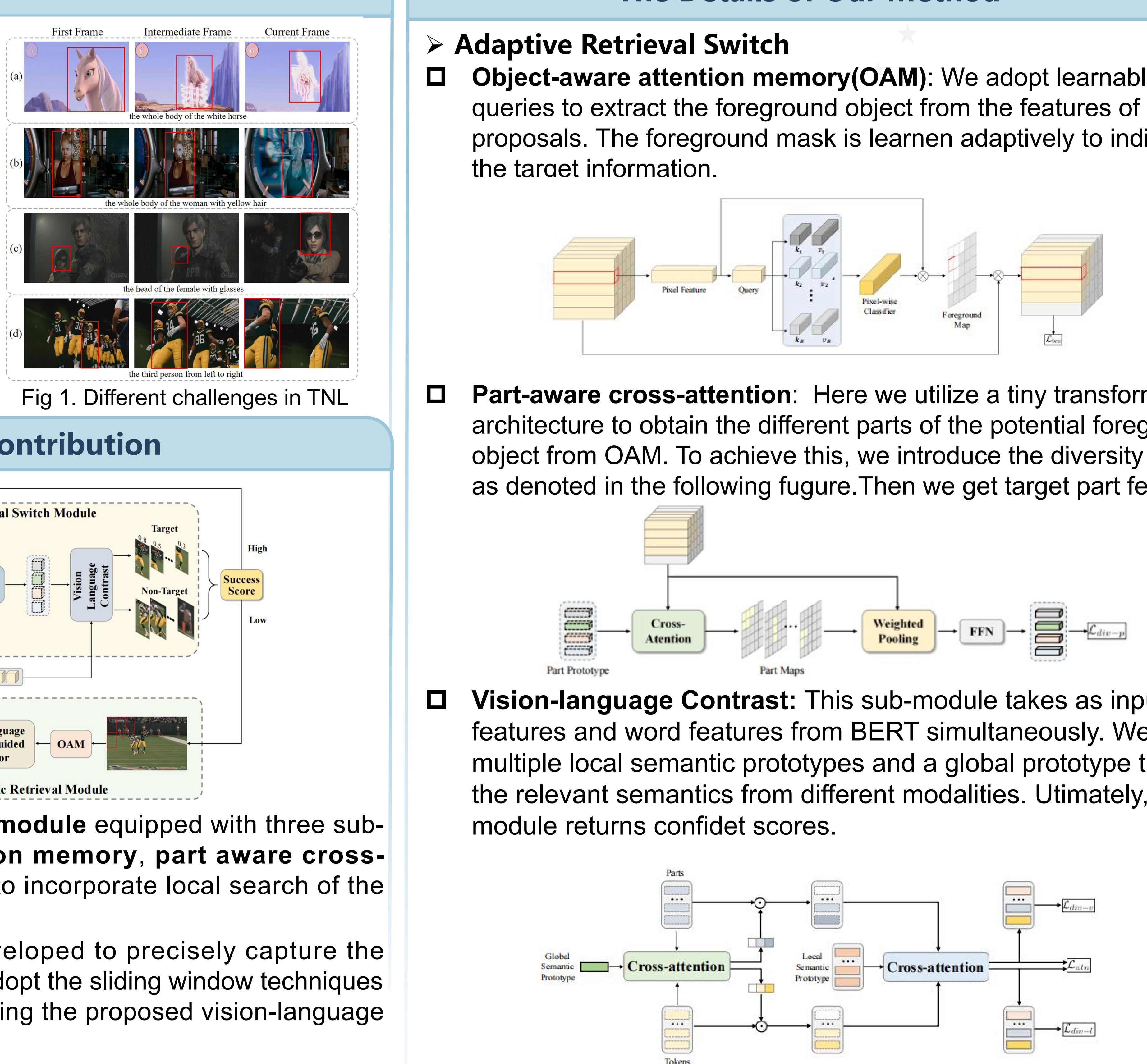
## **Our Motivation**

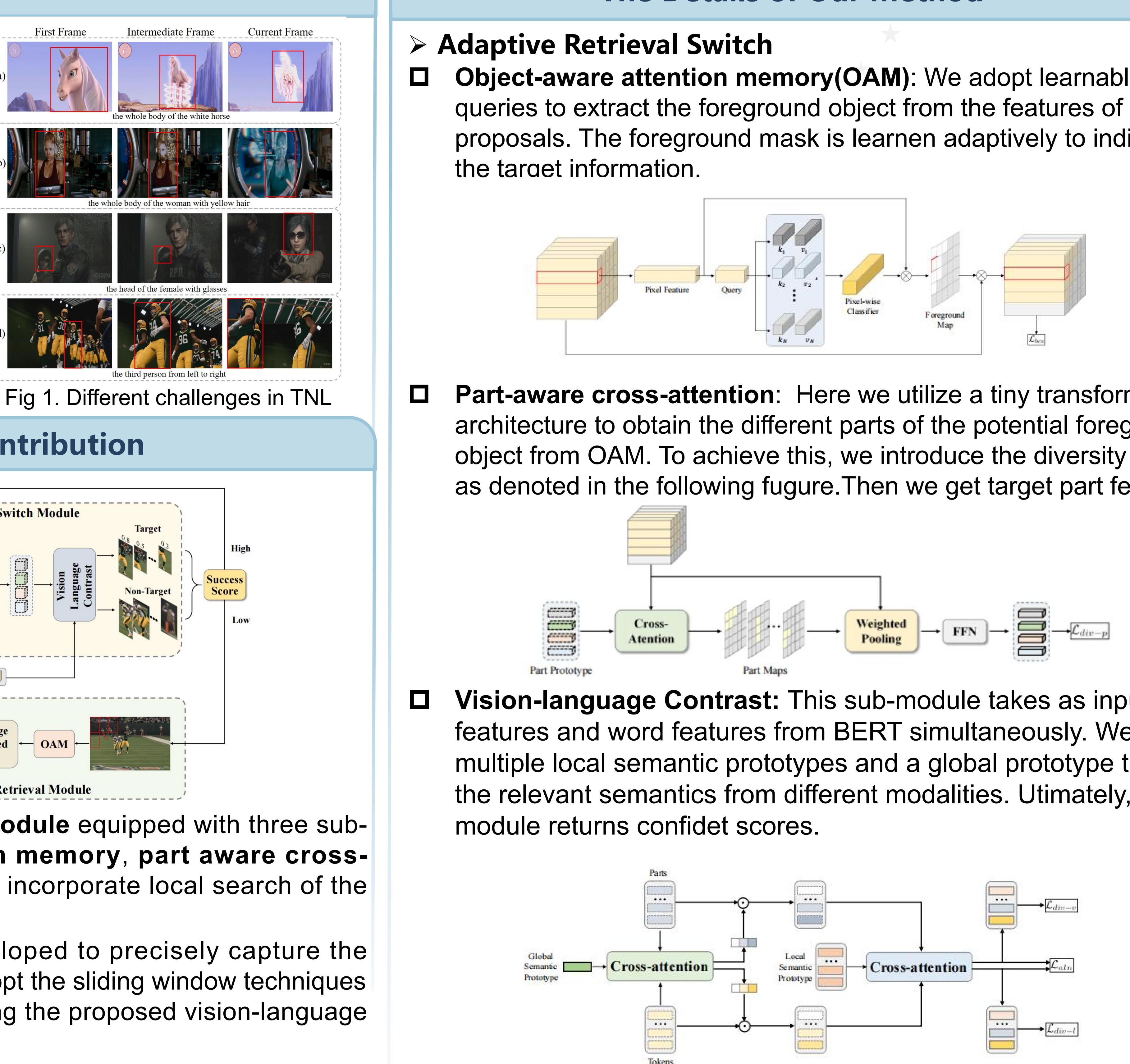
 $\triangleright$  aim at tracking the foreground object in a language or only by natural language.

visual template in a image-wide range.

- Lost the target easily and frequently.
- > Not consider the local search of target.







Jiamin Wu, Yihao Li, Jun Yu, Zhongpen Cai, Yuwen Pan

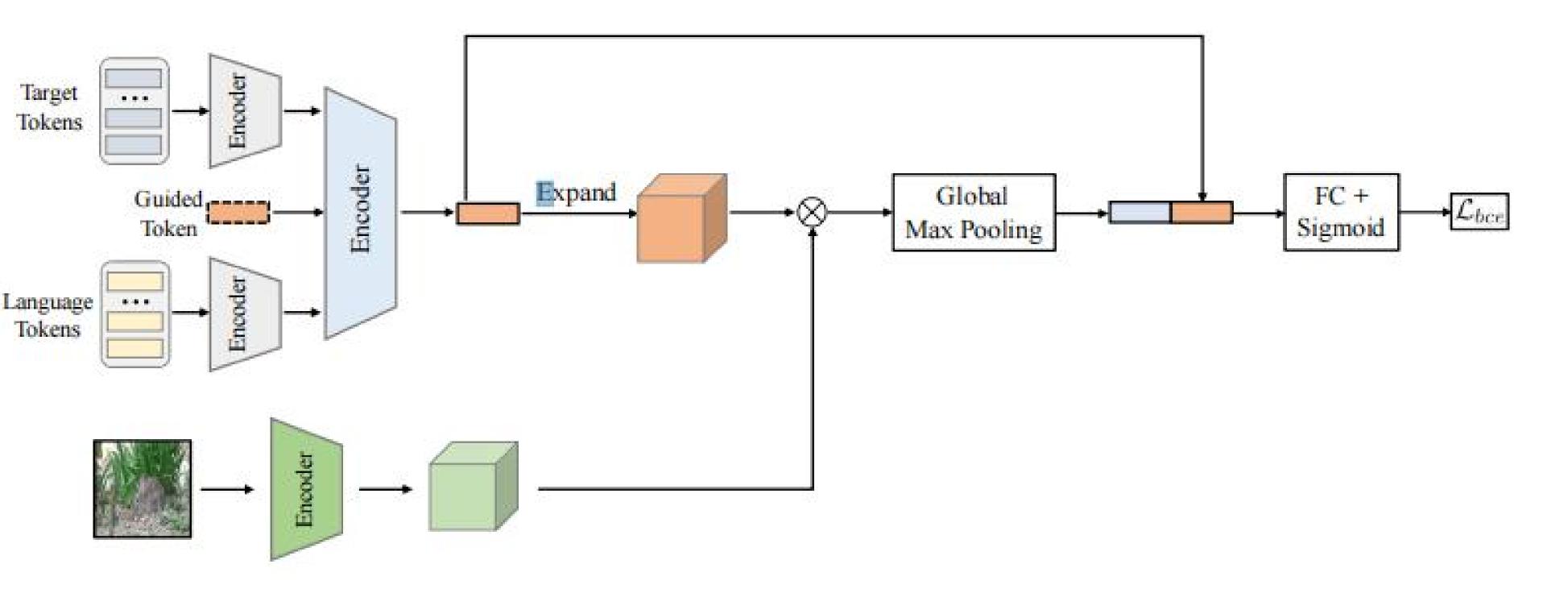
University of Science and Technology of China

## The Details of Our Method



## **Target-specific Retrieval**

We take the output from vision-language contrast module as indicator of switch to global search.We also adopt the window technique to get proposals of different location. Then we test each proposal with the proposed module as below.



## Experiments

### Tab 1. Results in different benchmarks.

Method	<b>OTB-Lang</b>		LaSOT		TNL2K	
	NL	NL+BBox	NL	NL+BBox	NL	NL+BBox
Li et al. [26]	0.29   0.25	0.72   0.55		-	_	-
eng et al. [13]	0.56   0.54	0.73   0.67	-	0.56   0.50	-	0.27   0.34   0.25
eng et al. [12]	$0.78 \mid 0.54$	0.79   0.61	0.28   0.28	0.35   0.35		0.27   0.33   0.25
ang et al. [40]	-	0.89   0.65	-3	0.30   0.27	-	-
GTI [45]	-	0.73   0.58	s	0.47   0.47		-
ang et al. [41]	0.24   0.19	0.88   0.68	0.49   0.51	0.55   0.51	0.06   0.11   0.11	0.42   0.50   0.42
Ours	0.72   0.53	0.91   0.69	$0.51 \mid 0.52$	$0.56 \mid 0.53$	$0.09 \mid 0.15 \mid 0.14$	$0.45 \mid 0.52 \mid 0.44$

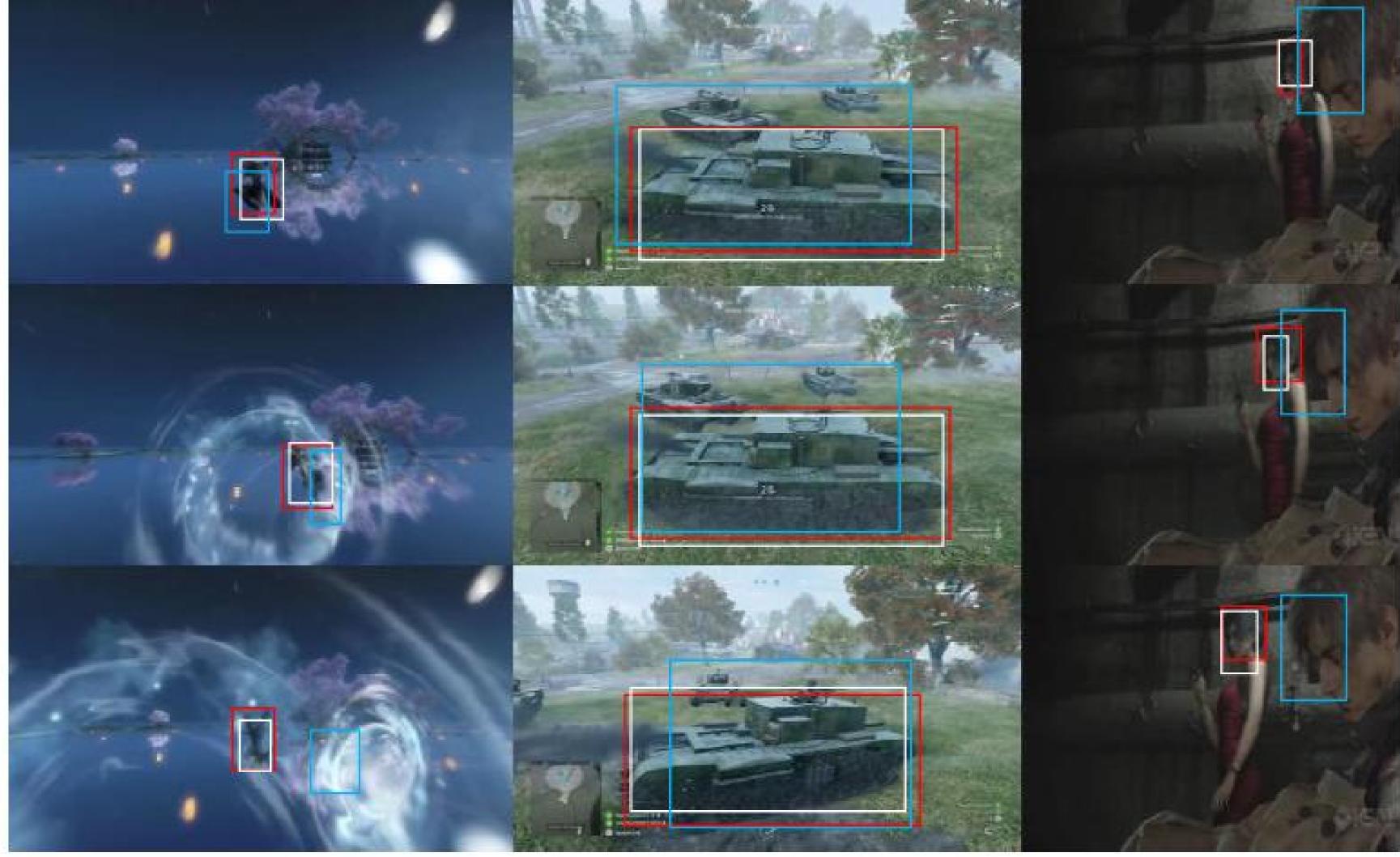


Fig 2. Visualization of our method.