

A Multi-modal Global Instance Tracking Benchmark (MGIT): Better Locating Target in Complex Spatio-temporal and Causal Relationship

Shiyu Hu^{1,2}, Dailing Zhang^{1,2}, Meiqi Wu³, Xiaokun Feng^{1,2}, Xuchen Li³, Xin Zhao^{1,2}, Kaiqi Huang^{1,2,5}

PROCESSING 1 School of Artificial Intelligence, University of Chinese Academy of Sciences; 2 Institute of Automation, Chinese Academy of Sciences; 3 School of Computer Science and Technology, University of Chinese Academy of Sciences; 4 School of Computer Science, Beijing University of Posts and Telecommunications; 5 Center for Excellence in Brain Science and Intelligence Technology, Chinese Academy of Sciences

Long sequences with complex spatial-

Video temporal variation and casual relationship

Coupling human causal reasoning ability into

multi-modal information











Language description: the second arrow from left to right

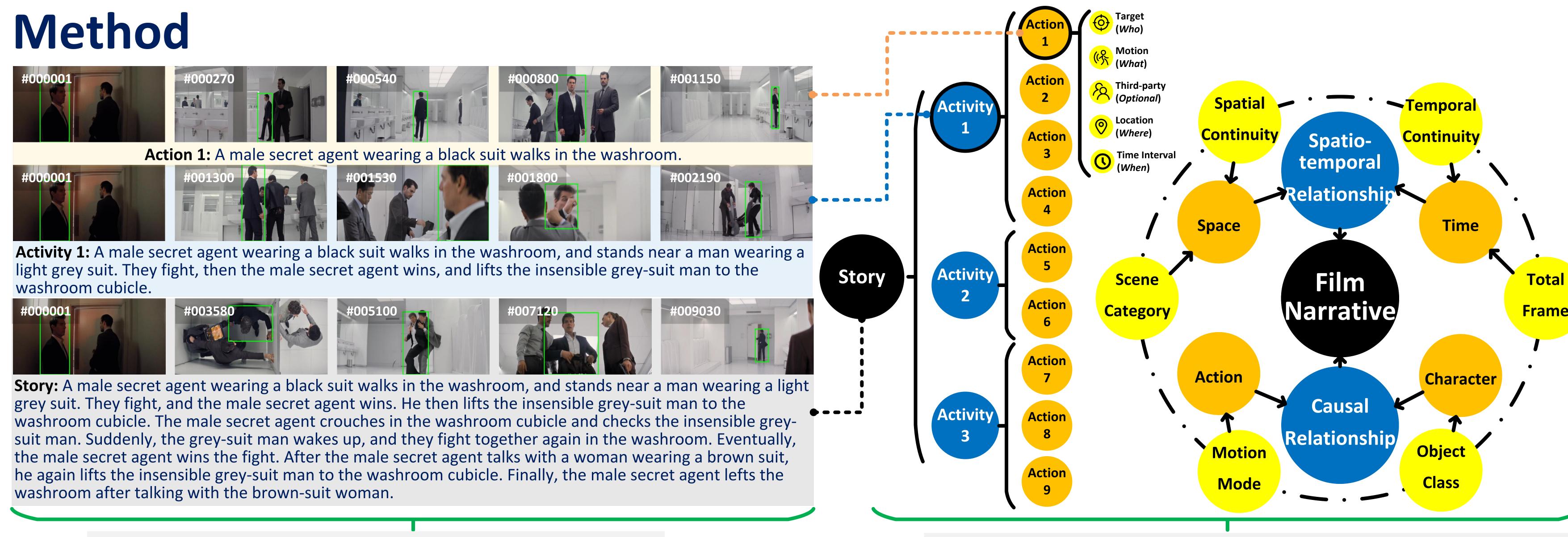




Simple semantic descriptions Lang- (multiple qualified targets, no qualified target)



Existing trackers always perform poorly in complex environments (e.g., longer videos with more complicated narrative content)



Contributions

- We propose a new multimodal benchmark named MGIT. It consists of 150 long videos with a total of 2.03 million frames, and the average length of a single sequence is 5-22 times longer than existing multimodal benchmarks.
- We design a multi-granular annotation strategy for providing scientific semantic information.
- We execute comparative experiments on other Multi-granular annotation strategy based on benchmarks, and conduct Language the hierarchical structure of human cognitive detailed experimental analyses on MGIT.



Tracker	OTB-Lang [1]		TNL2k 3		LaSOT [2]		LaSOText 17		LaSOTSub		LaSOTNLC		MGIT	
	PRE	SR	PRE	SR	PRE	SR	PRE	SR	PRE	SR	PRE	SR	PRE	SR
SNLT 46	0.848	0.666	0.081	0.100	0.475	0.459	0.306	0.262	0.527	0.495	0.513	0.483	0.004	0.036
VLT_SCAR [42]	0.898	0.739	0.556	0.497	0.677	0.630	0.503	0.428	0.670	0.633	0.659	0.633	0.124	0.177
VLT_TT [42]	0.931	0.764	0.583	0.539	0.714	0.670	0.549	0.465	0.707	0.660	0.721	0.662	0.324	0.474
JointNLT [18]	0.856	0.653	0.598	0.552	0.640	0.607	0.457	0.398	0.624	0.583	0.707	0.651	0.433	0.603
d case analysis Results of different trackers on MGIT														1GIT

Tracker

SiamCAR [1]

PrDiMP 12

TransT [39]

OSTrack [15]

GRM [16]

SNLT 46

VLT_TT [42]

Conclusions

KeepTrack [13

MixFormer [14]

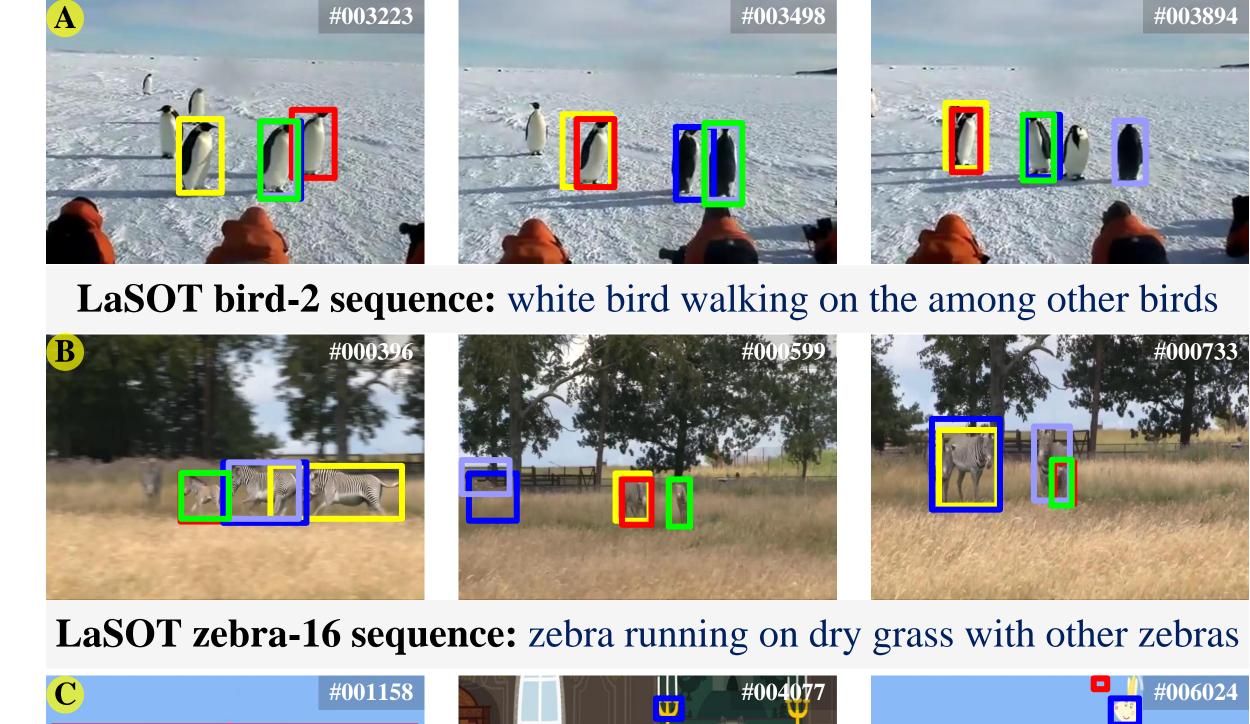
VLT_SCAR [42] SNN

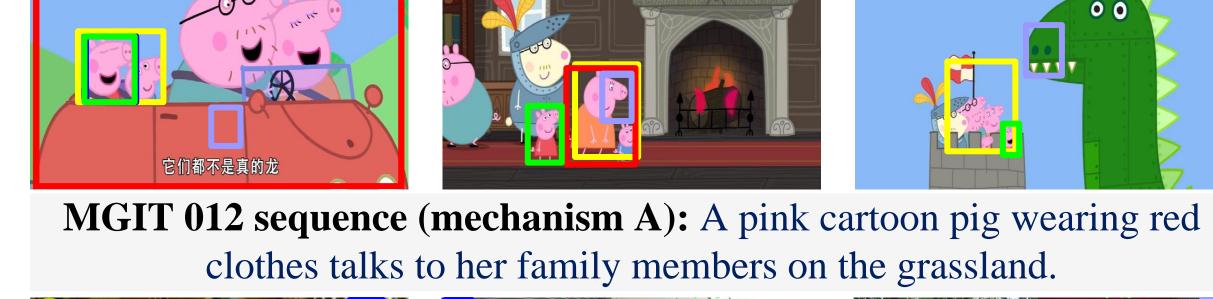
SiamRCNN [10]

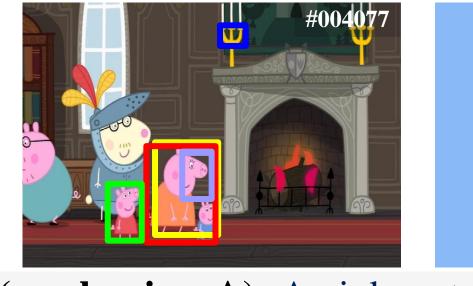
Helping algorithms understand video

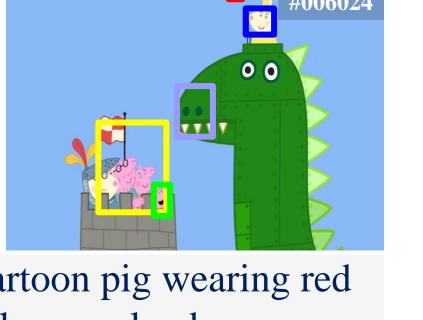
content from a multi-modal perspective

Bad case analysis











MGIT 006 sequence (mechanism A): A skateboard is slid by a man in black on the playground.

VLT_TT (NeurIPS22) SNLT (CVPR21)

 MGIT is a complex environment, the annotation strategy is a feasible solution for coupling human understanding into semantic labels.

NL&BBox

NL&BBox

NL&BBox

NL&BBox

(mechanism B-E)

Action (B)

Activity (C)

Activity (C)

Activity (C)

Activity (C)

Architecture

SNN+CF

SNN+CF

Transformer

Transformer

Fransforme

Transformer

SNN

 Existing trackers should improve the capability for processing long text and aligning multi-modal information.

