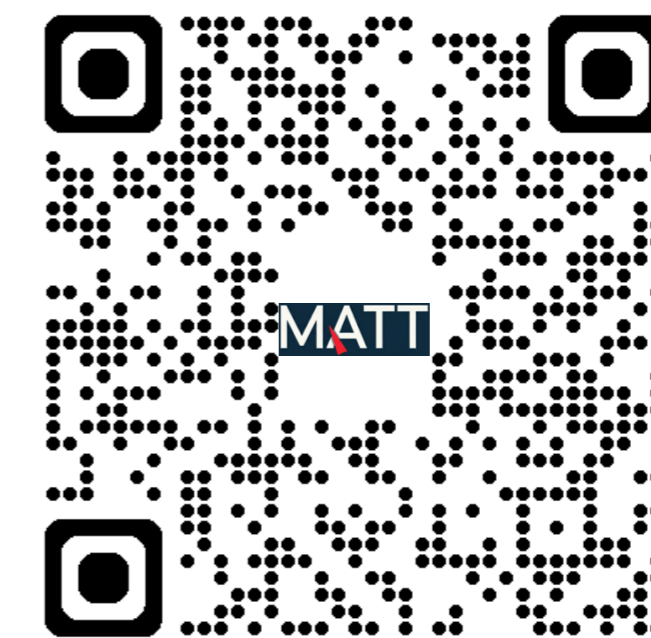


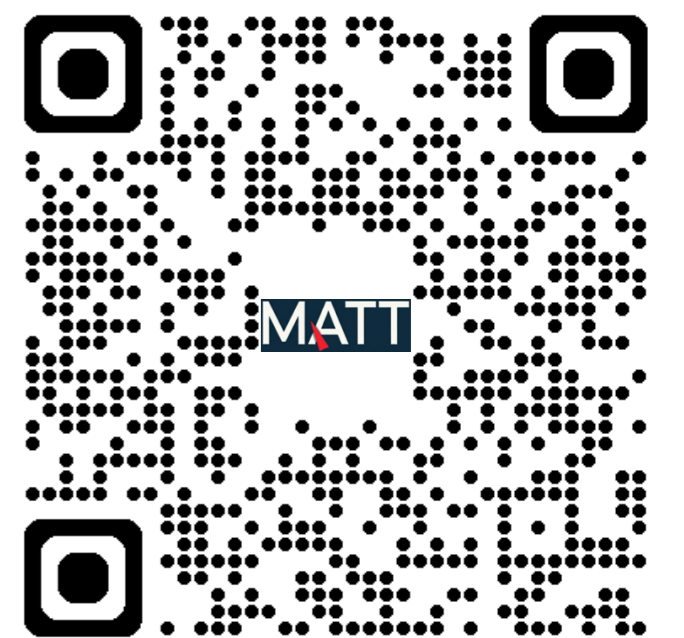
# Mistake Attribution: Fine-Grained Mistake Understanding in Egocentric Videos

Yayuan Li<sup>1</sup>, Aadit Jain<sup>1</sup>, Filippos Bellos<sup>1</sup>, Jason Corso<sup>1,2</sup>

<sup>1</sup>University of Michigan, <sup>2</sup>Voxel51



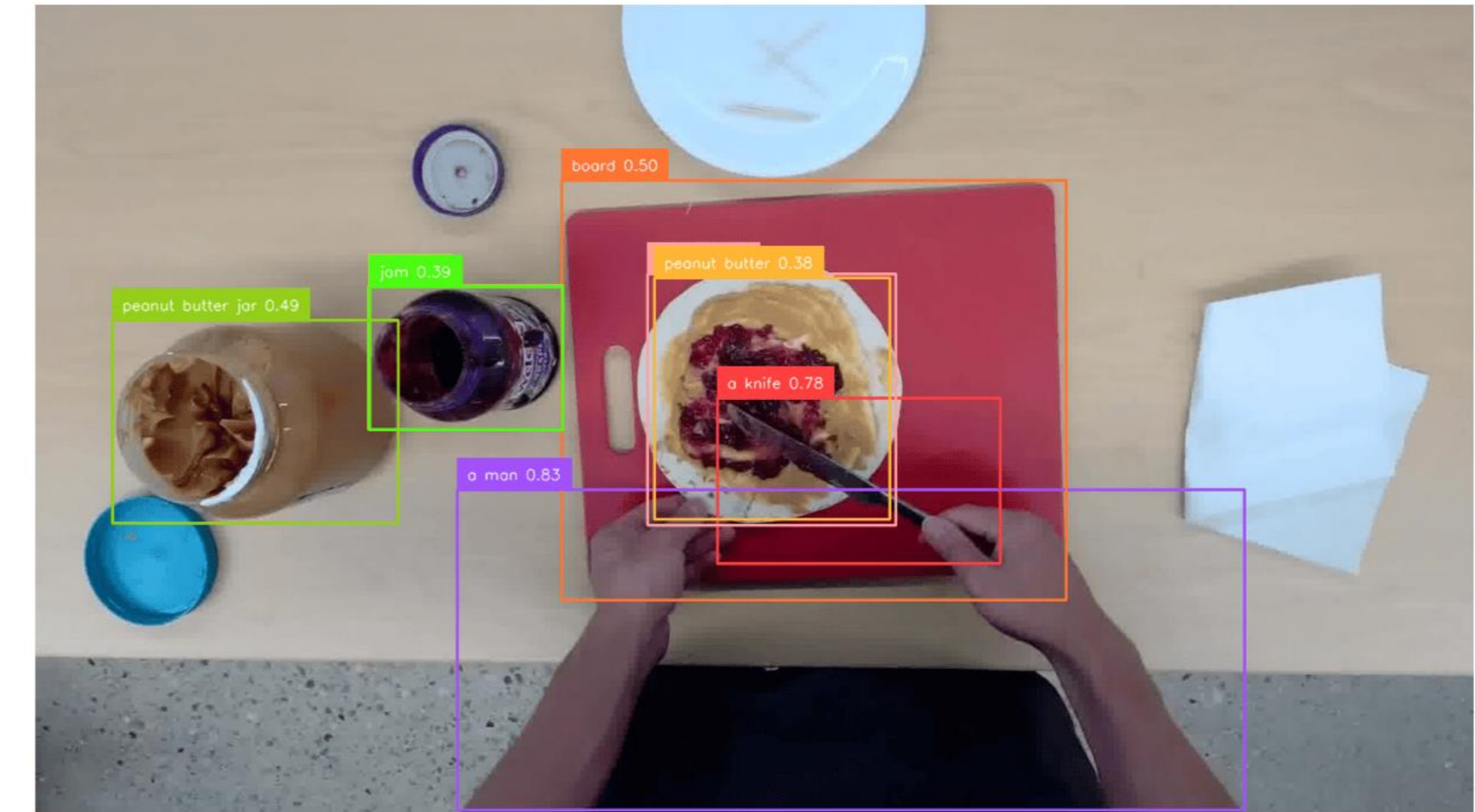
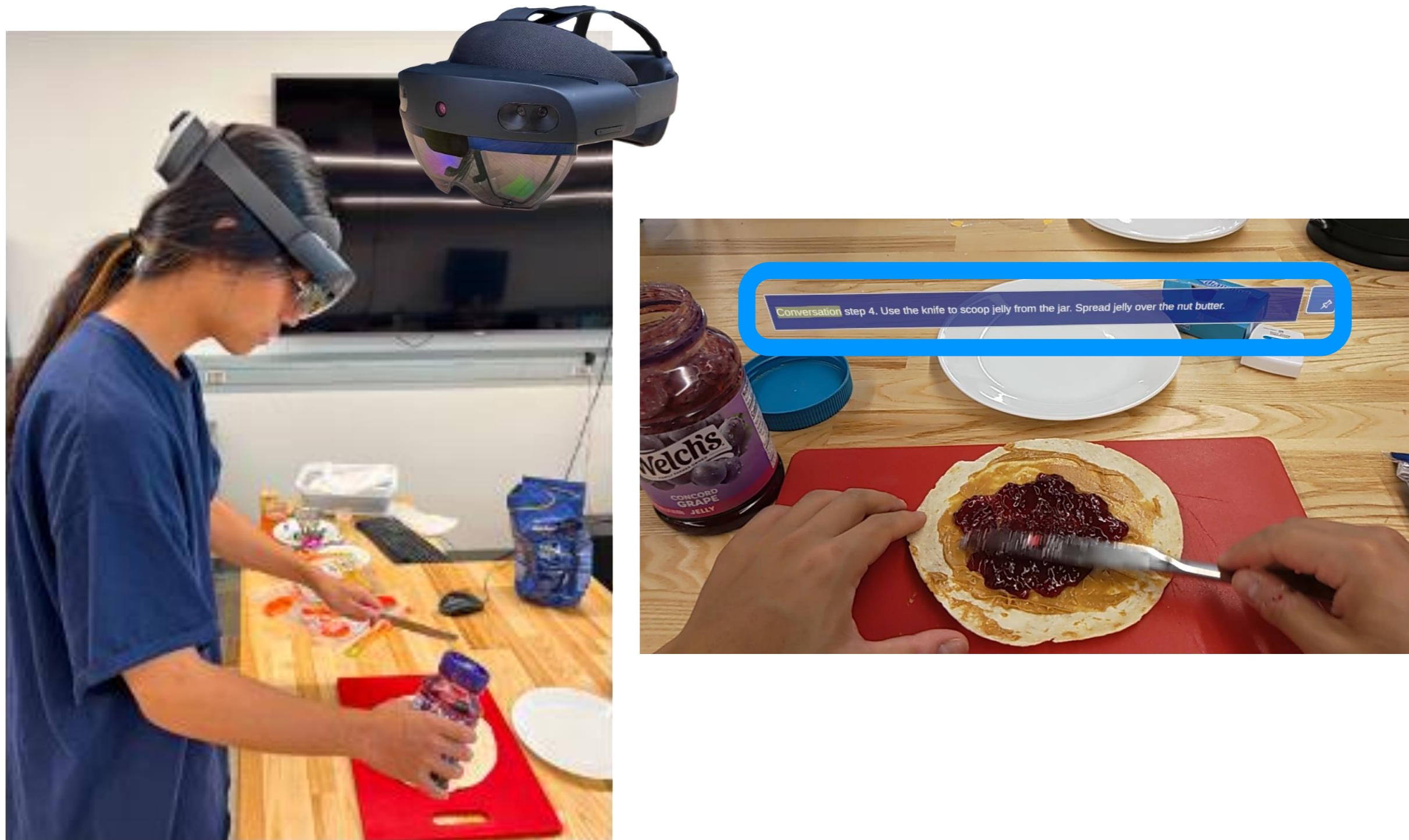
Page



Contact

# Mistake Understanding — What? & Why?

Envisioning Instructional AI<sup>[1,2,3]</sup>:  
Assist human to complete physical tasks



Understanding human activity is fundamental

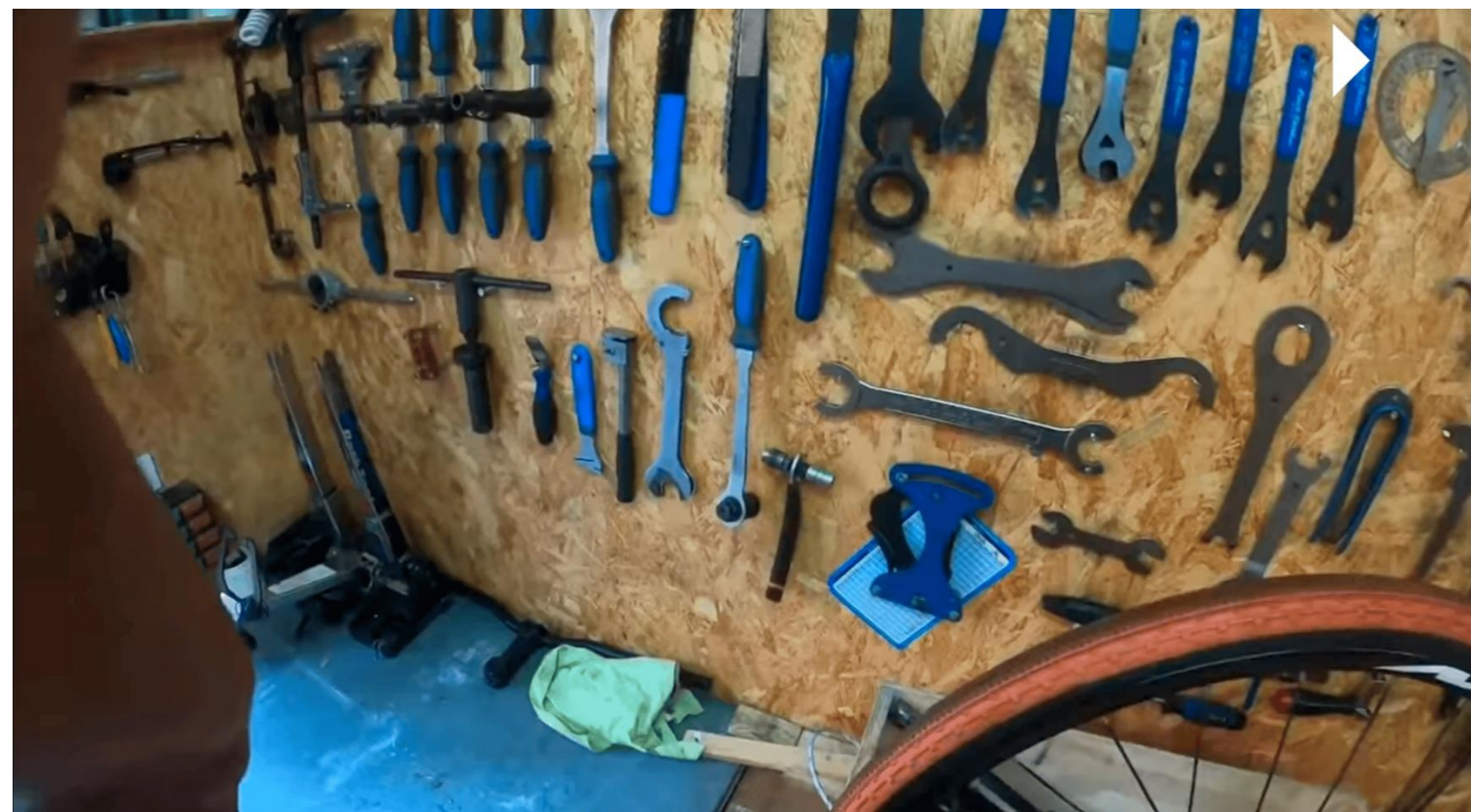
Human mistake is an under-explored activity in development of Instructional AI.

1. Bellos, Filippas\*, Yayuan Li\*, Cary Shu, Ruey Day, Jeffrey Siskind, and Jason Corso. "Towards Effective Human-in-the-Loop Assistive AI Agents." In Proceedings of the IEEE/CVF International Conference on Computer Vision (Workshop), pp. 2513-2522. 2025.
2. Castelo, Sonia, Joao Rulff, Erin McGowan, Bea Steers, Guande Wu, Shaoyu Chen, Iran Roman et al. "Argus: Visualization of ai-assisted task guidance in ar." IEEE transactions on visualization and computer graphics 30, no. 1 (2023): 1313-1323.
3. Bao, Yuwei, Keunwoo Yu, Yichi Zhang, Shane Storks, Xiao Zheng, and Joyce Chai. "Can Foundation Models Watch, Talk and Guide You Step by Step to Make a Cake?." the Association for Computational Linguistics: EMNLP 2023, pp. 12325-12341. 2023.

# Existing Mistake Understanding

Lack of diagnostic details

Attempt Video



Instruction Text: “Grab a chain whip”

Mistake Detection:

Yes

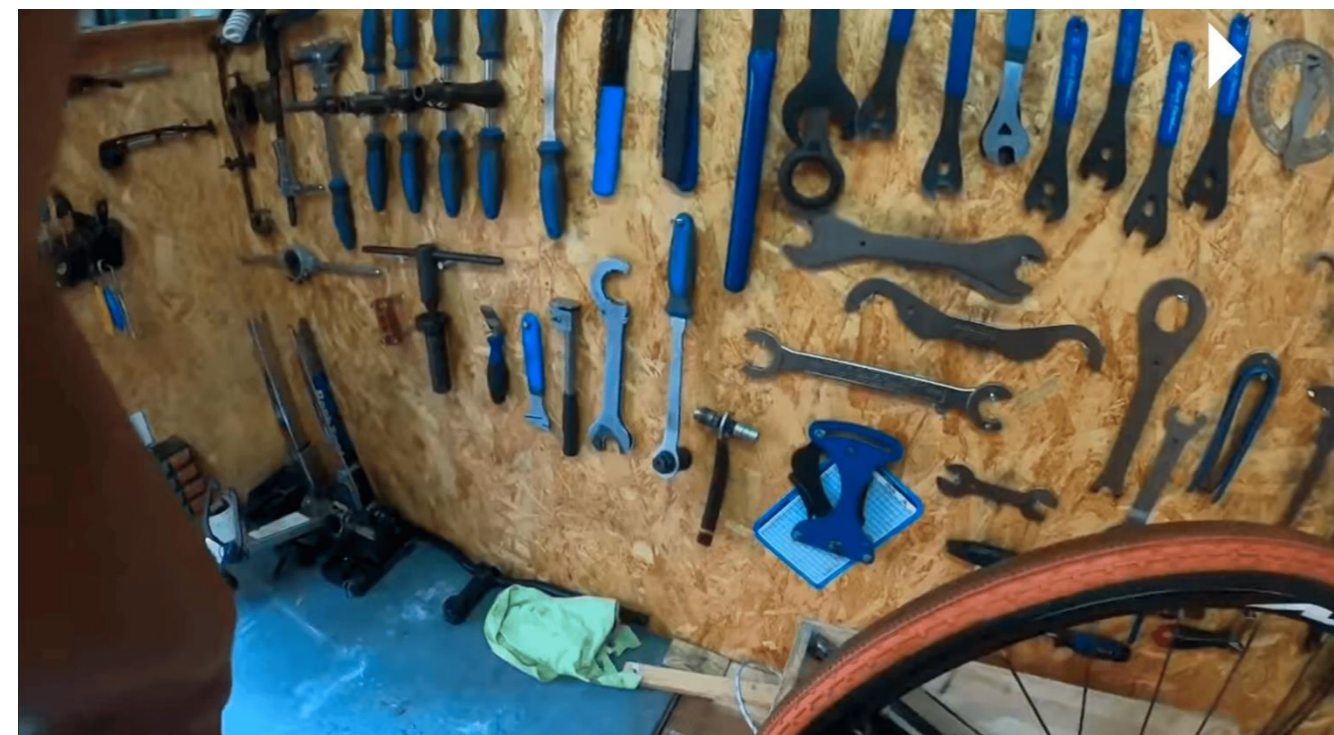
Mistake Recognition:

“Modification Mistake”

# Mistake Attribution — Problem Formulation

## Input

Grab a chain whip  
**Instruction Text**



**Attempt Video**

## Mistake ATtribution

### Semantic Attribution

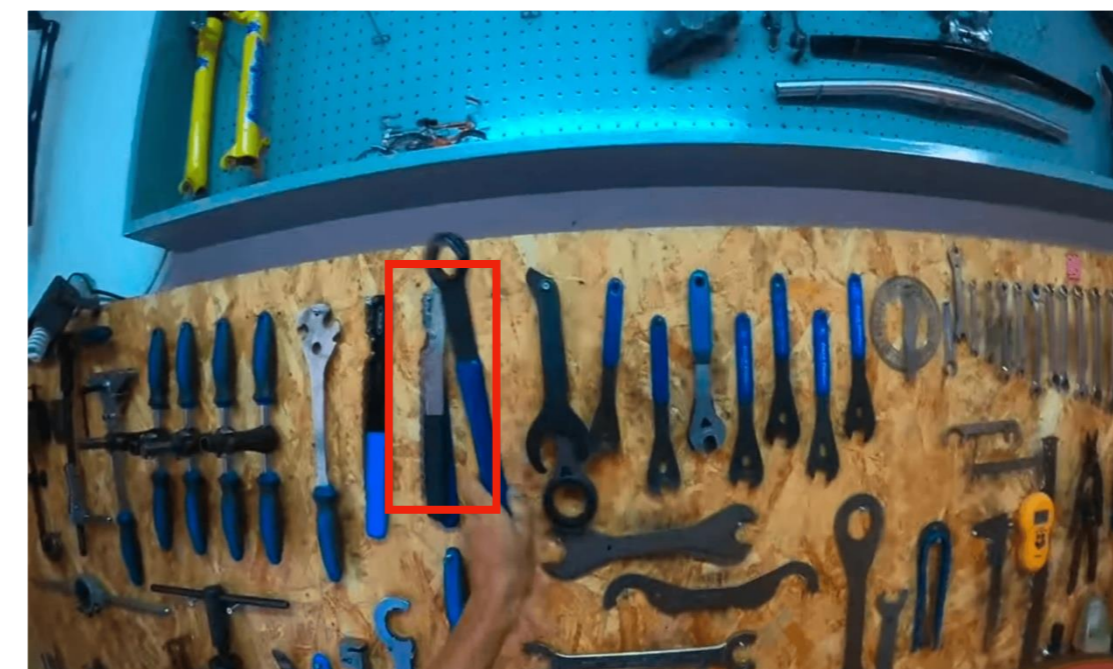
Grab a chain whip  
|  
Correct Predicate  
|  
Mistaken Object

### Temporal Attribution



Point-of-No-Return (PNR):  $t=3.2s$

### Spatial Attribution



# Challenges

- Benchmark:  
Hard to manually collect large scale mistake dataset and attribution annotation
  
- Method:  
Lack of a unified models to produce semantic (what), temporal (when), spatial (where) attribution

# Challenges

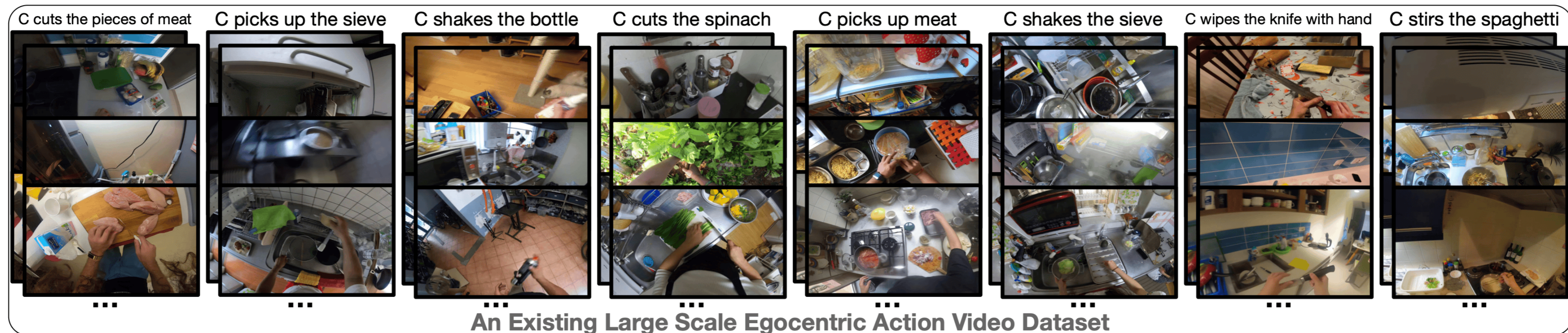
- Benchmark:  
Hard to manually collect large scale mistake dataset and attribution annotation

Dataset	# Samples		Semantic Var.		Visual Var.		Annotation			
	Total	By activity	Activ.	Dm.	Env.	Part.	Det.	Sem.	Temp.	Spa.
EgoPER [23]	599	9.7	62	1	2	11	✓	▲	✗	▲
Assembly101 [4, 40]	707	2.0	358	1	1	53	✓	▲	✗	✗
HoloAssist [46]	7,562	0.91	8,285	2	-	222	✓	▲	▲	✗
CaptionCook4D [37]	1,964	5.6	352	1	10	8	✓	▲	✗	✗
Epic-Tent [18]	626	52.2	12	1	-	24	✓	▲	✗	▲
EPIC-KITCHENS	89,975	5.0	18,003	1	45	37	✗	✗	✗	▲
Ego4D	155,367	2.1	75,423	14+	100+	931	✗	✗	▲	▲

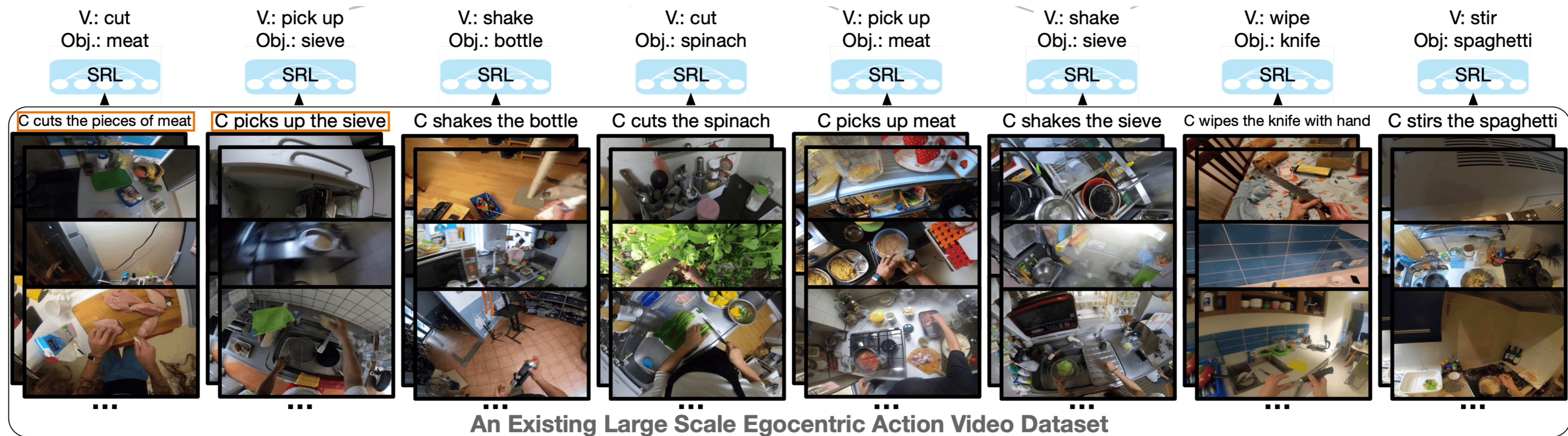
# MisEngine

Automatically compose mistakes from existing action datasets

The essence of mistakes is the divergence of video from text

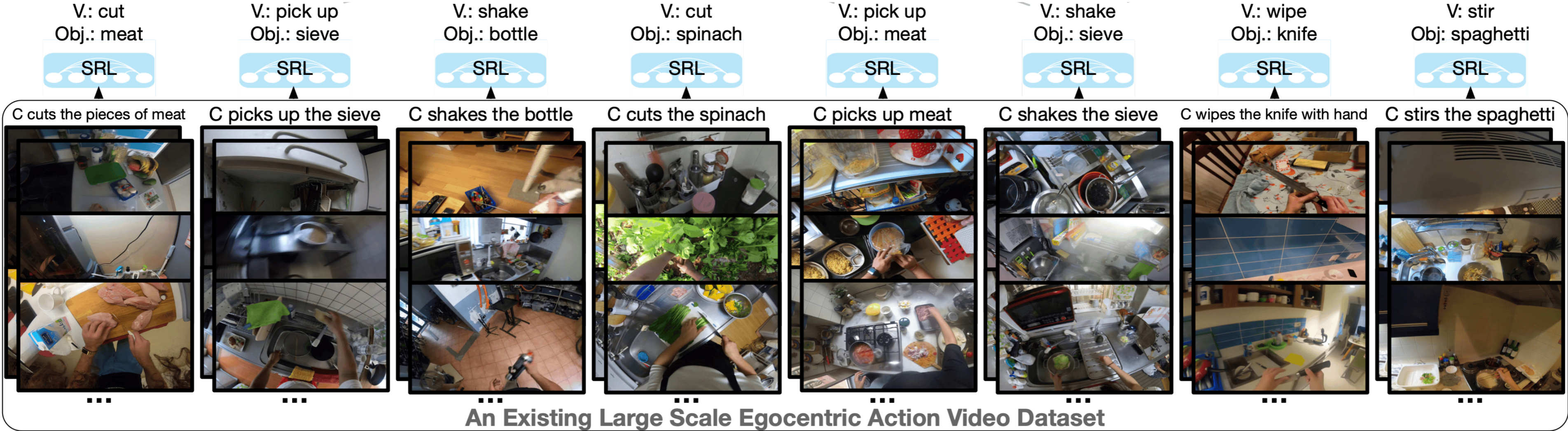
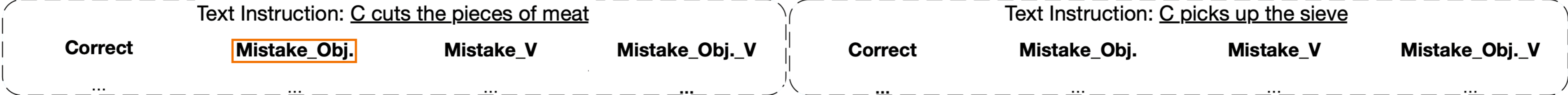


# MisEngine

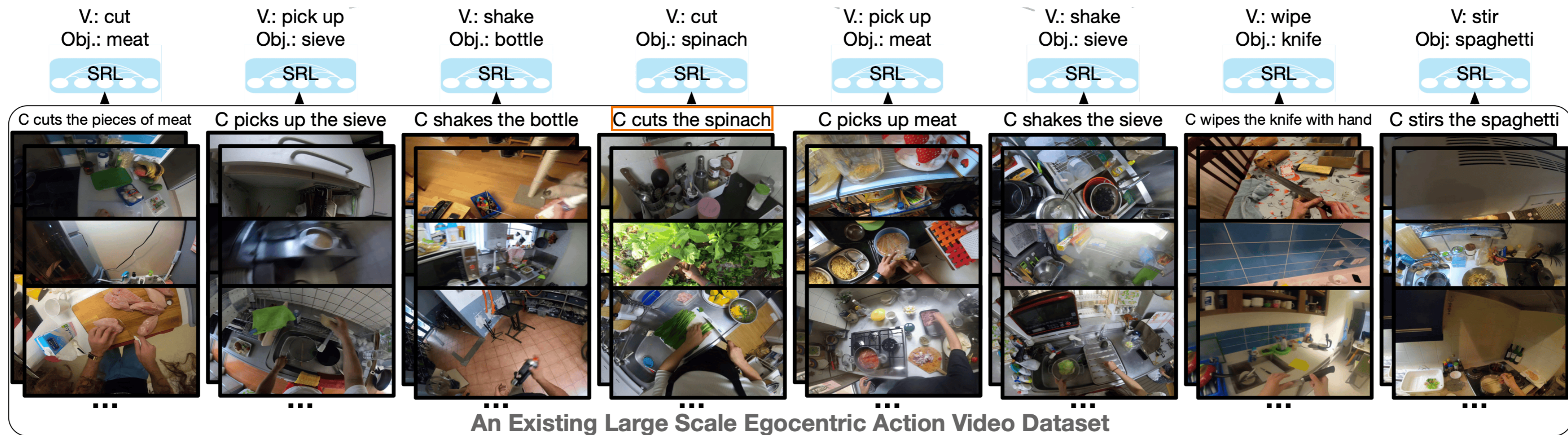
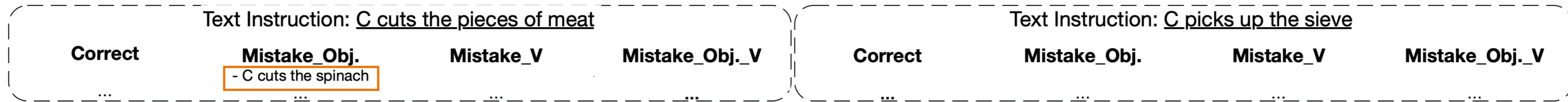




# MisEngine

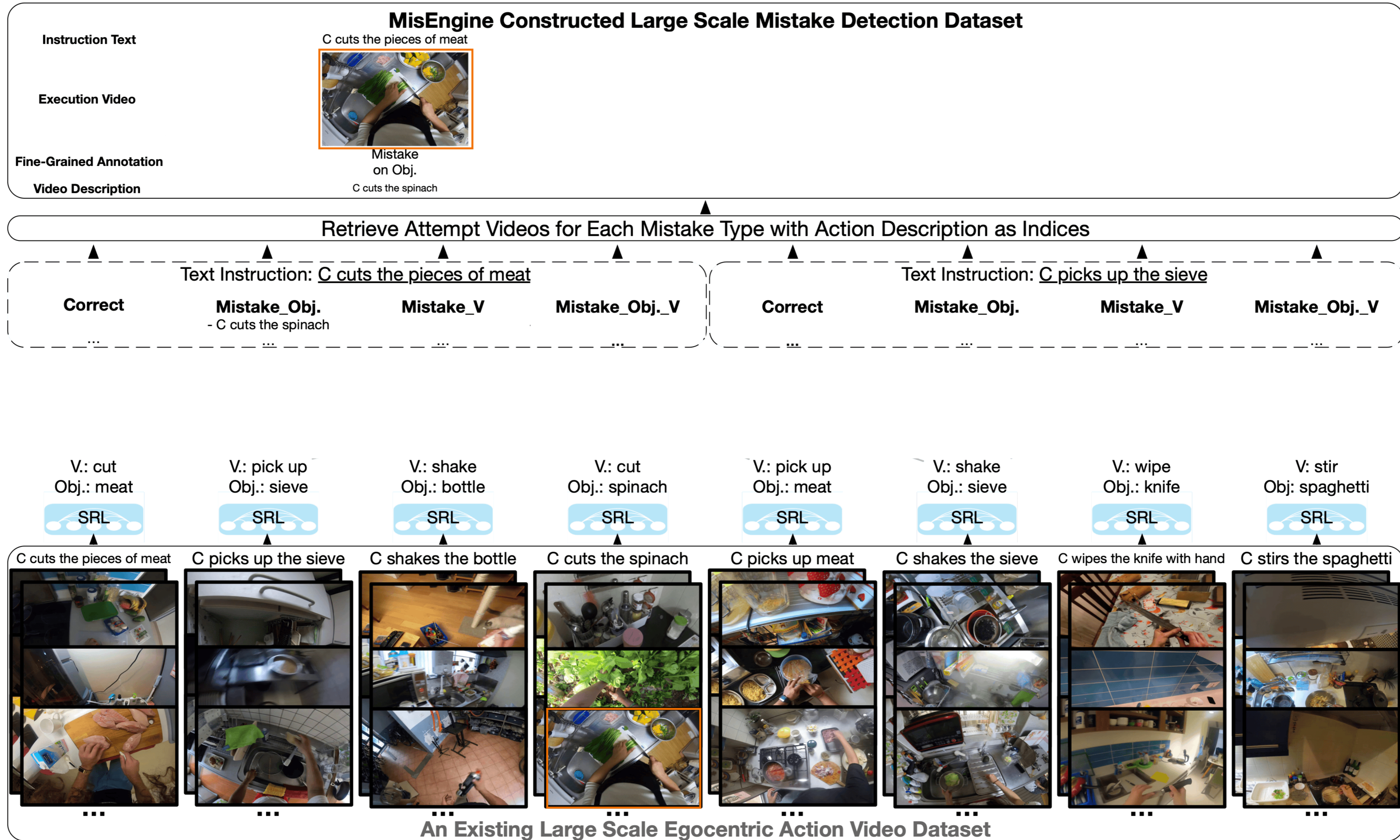


# MisEngine

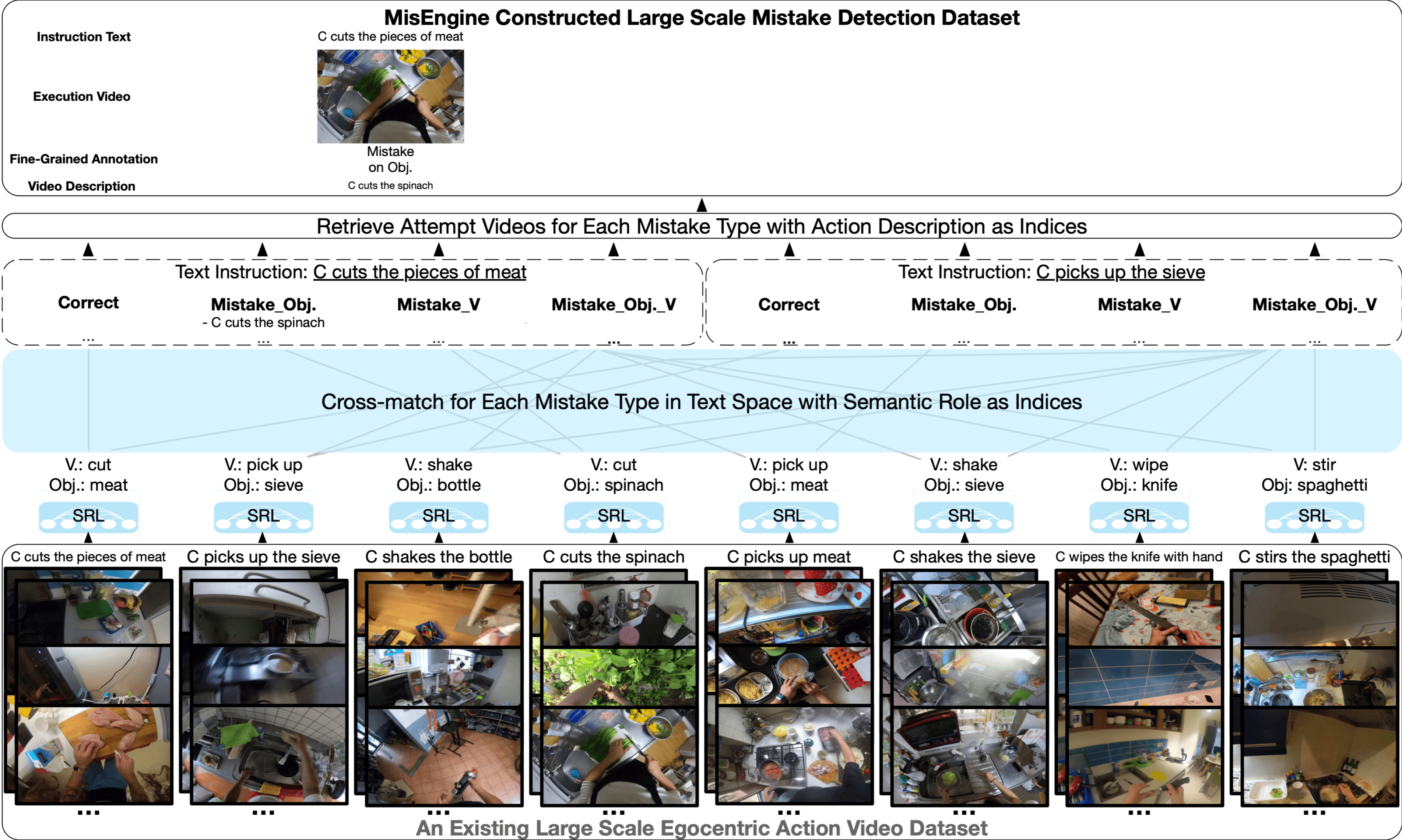




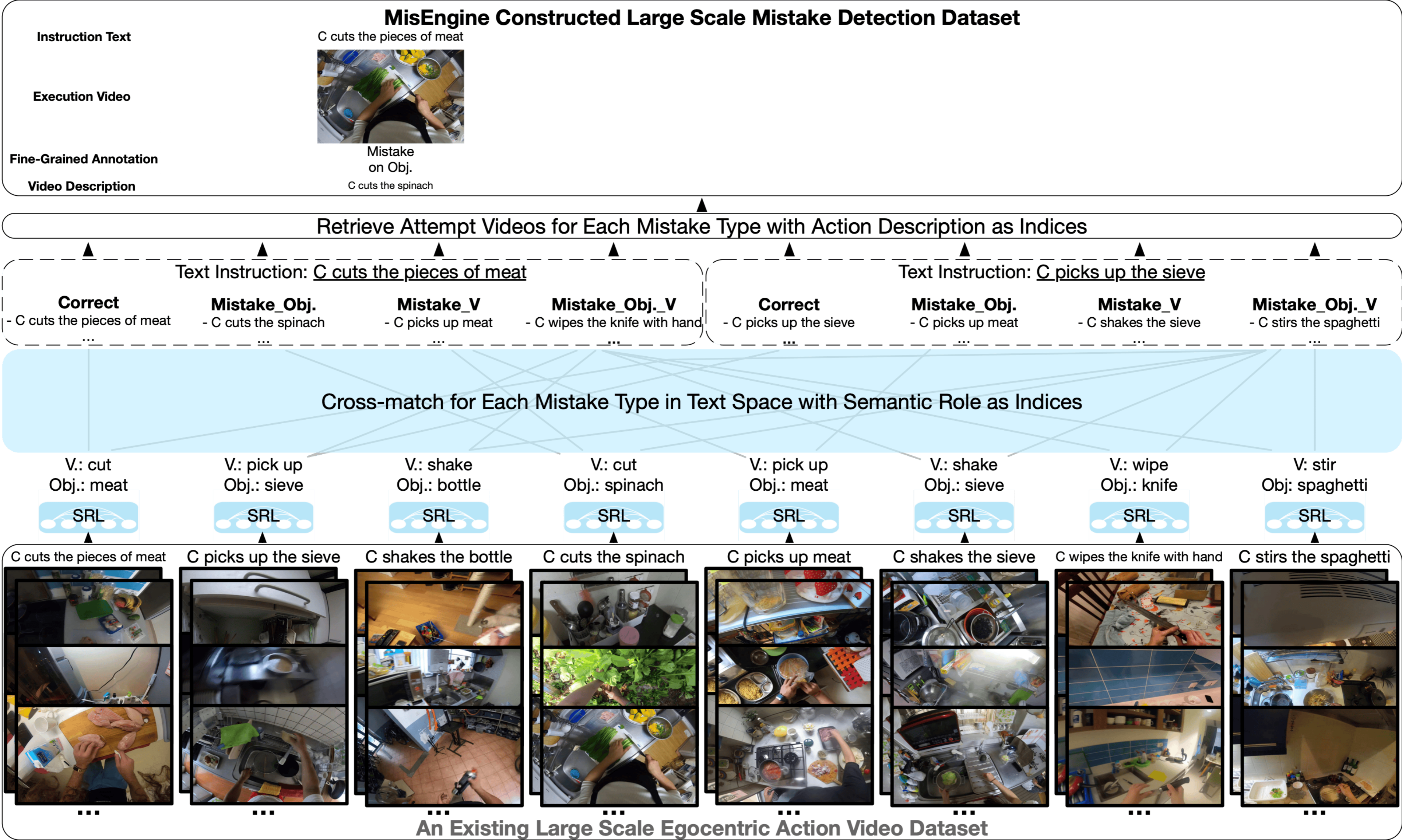
# MisEngine



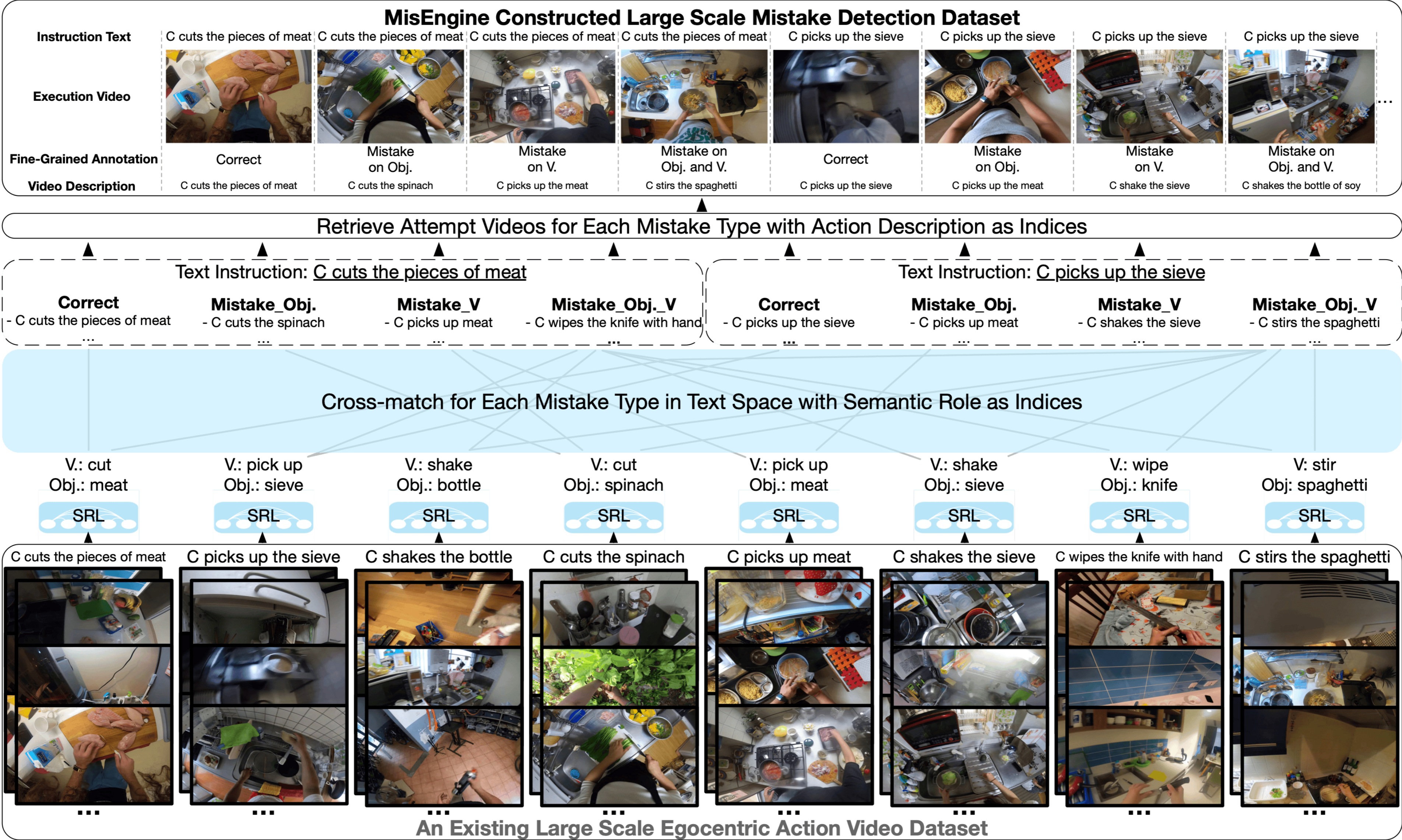
# MisEngine

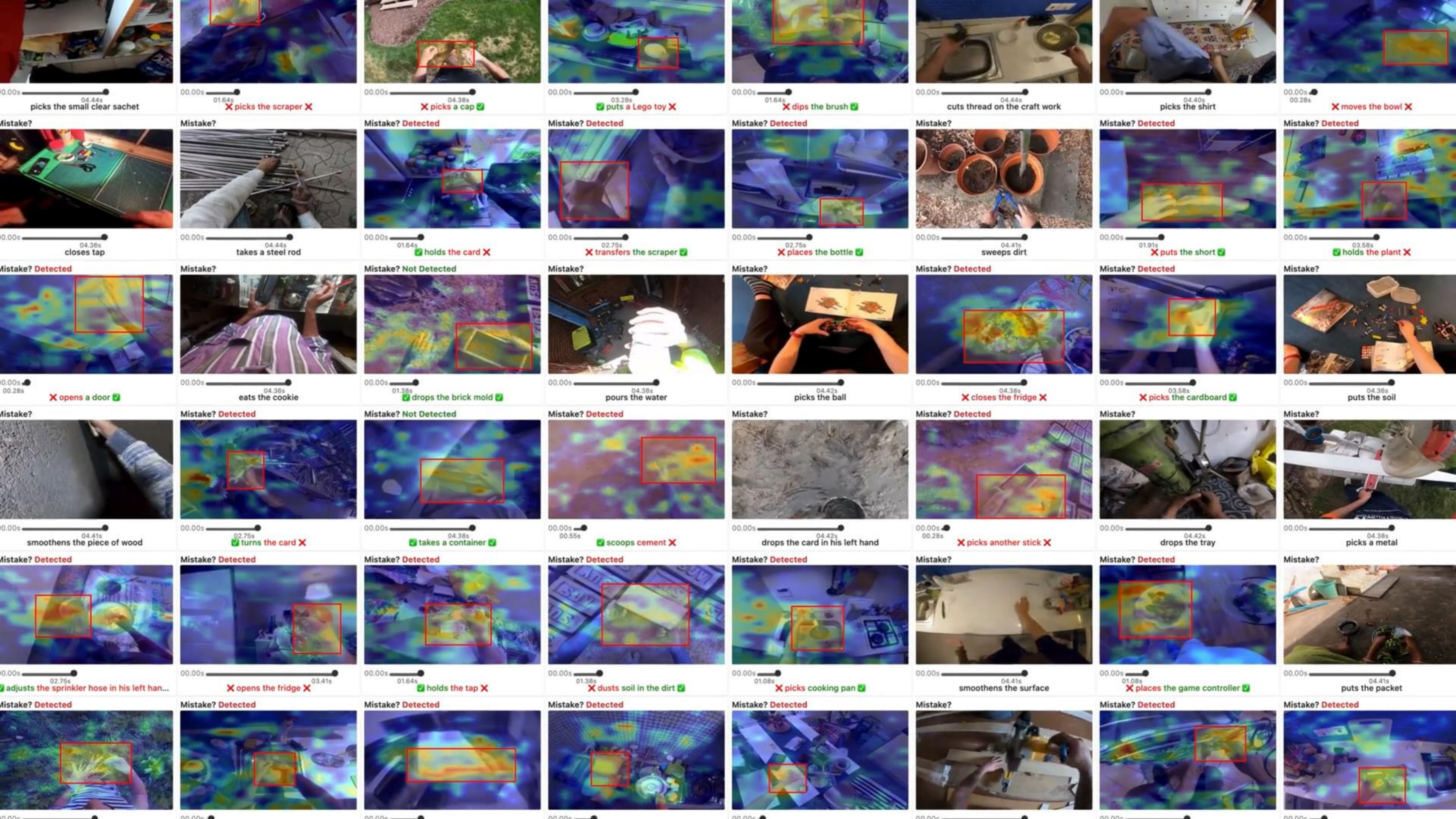


# MisEngine



# MisEngine





# MisEngine

Inherit scale and diversity from existing action datasets

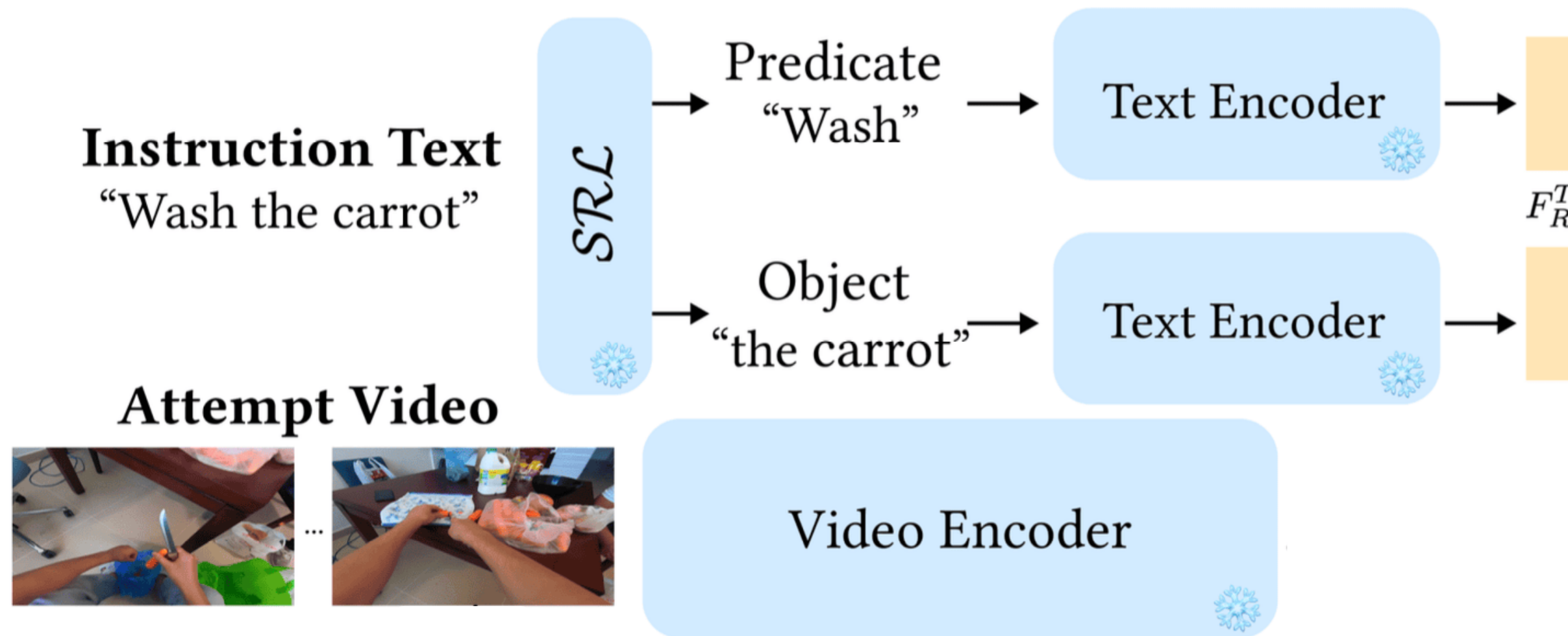
The composition implies mistake detection & semantic attribution annotation

Inherit temporal and spatial attribution from original datasets

Dataset	# Samples		Semantic Var.		Visual Var.		Annotation			
	Total	By activity	Activ.	Dm.	Env.	Part.	Det.	Sem.	Temp.	Spa.
EgoPER [23]	599	9.7	62	1	2	11	✓	▲	✗	▲
Assembly101 [4, 40]	707	2.0	358	1	1	53	✓	▲	✗	✗
HoloAssist [46]	7,562	0.91	8,285	2	-	222	✓	▲	▲	✗
CaptionCook4D [37]	1,964	5.6	352	1	10	8	✓	▲	✗	✗
Epic-Tent [18]	626	52.2	12	1	-	24	✓	▲	✗	▲
EPIC-KITCHENS	89,975	5.0	18,003	1	45	37	✗	✗	✗	▲
Ego4D	155,367	2.1	75,423	14+	100+	931	✗	✗	▲	▲
EPIC-KITCHENS-M	221,094	18.0	12,283	1	45	37	✓	✓	✗	▲
Ego4D-M	257,584	16.0	16,099	14+	100+	248	✓	✓	✓	✓

# MisFormer

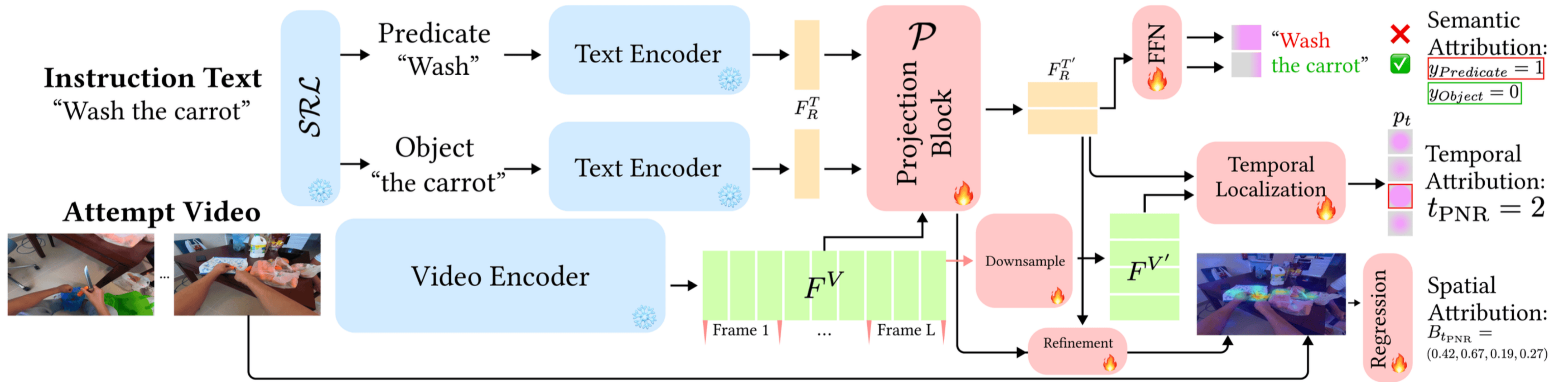
Pretrained VLM for feature extraction



# MisFormer

Pretrained VLM for feature extraction

Tailored attribution heads



# Quantitative Results

MisFormer, an unified model, outperforms baselines for individual task

Dataset	Method	Average		Predicate		Object	
		Acc. $\uparrow$	F1 $\uparrow$	Acc. $\uparrow$	F1 $\uparrow$	Acc. $\uparrow$	F1 $\uparrow$
EK	LLaVA-Vid [49]	71.71	70.34	64.10	64.55	79.32	76.13
	Vid-Chat [29]	64.17	63.27	54.45	52.11	73.89	74.43
	ChatGPT [34]	77.66	77.23	66.23	65.11	89.09	89.35
	MisFormer (Ours) <sup>†</sup>	<b>84.91</b>	<b>84.78</b>	<b>76.35</b>	<b>75.31</b>	<b>93.47</b>	<b>94.25</b>
	MisFormer (Ours)	<b>84.13</b>	<b>83.89</b>	<b>76.83</b>	<b>76.43</b>	<b>91.43</b>	<b>91.34</b>
	Ego	LLaVA-Vid [49]	42.83	40.12	48.10	47.95	37.56
	Vid-Chat [29]	36.47	34.36	36.73	36.84	36.21	31.88
	ChatGPT [34]	52.45	50.95	53.42	50.20	51.48	51.70
	MisFormer (Ours) <sup>†</sup>	<b>59.37</b>	<b>55.41</b>	<b>55.44</b>	<b>53.03</b>	<b>63.33</b>	<b>57.79</b>
	MisFormer (Ours)	<b>62.03</b>	<b>56.24</b>	<b>58.40</b>	<b>55.22</b>	<b>65.66</b>	<b>57.25</b>

Semantic Attribution

Method	MAE (frames) $\downarrow$	MAE (seconds) $\downarrow$	
EgoMotion-COMPASS [24]	48.96	1.632	
EgoT2 [48]	24.48	0.816	
MisFormer (Ours)	<b>19.14</b>	<b>0.638</b>	

**Temporal Attribution**

Method	mIoU (%) $\uparrow$	CD (%) $\downarrow$	BSE (%) $\downarrow$
MediaPipe-U [28]	49.88	13.47	20.98
SSDA [25]	<b>64.54</b>	<b>7.21</b>	<b>12.34</b>
MisFormer (Ours)	<u>59.21</u>	<u>10.36</u>	<u>16.27</u>

Spatial Attribution

# Results: Semantic Attribution

Prune chili peppers



Instruction Text

Attempt Video

Prune chili peppers

Mistaken Predicate  Mistaken Object

ChatGPT-4o<sup>1</sup>

Prune chili peppers

Correct Predicate  Mistaken Object

LLaVa-Vid<sup>2</sup>

Prune chili peppers

Mistaken Predicate  Correct Object

GroundTruth

Prune chili peppers

Mistaken Predicate  Correct Object

**MisFormer  
(Ours)<sup>3</sup>**

1. OpenAI. (2025). GPT-4o System Card. arXiv:2410.21276.  
2. Bin Lin, Yang Ye, Bin Zhu, Jiayi Cui, Munan Ning, Peng Jin, Li Yuan. Video-LLaVA: Learning United Visual Representation by Alignment Before Projection. EMNLP 2024.  
3. Yayuan Li, Aditya Jain, Foteini Bellos, Jason Corso. Mistake Attribution: Fine-Grained Mistake Understanding in Egocentric Videos. CVPR 2026.

# Results: Temporal Attribution

Prune chili peppers

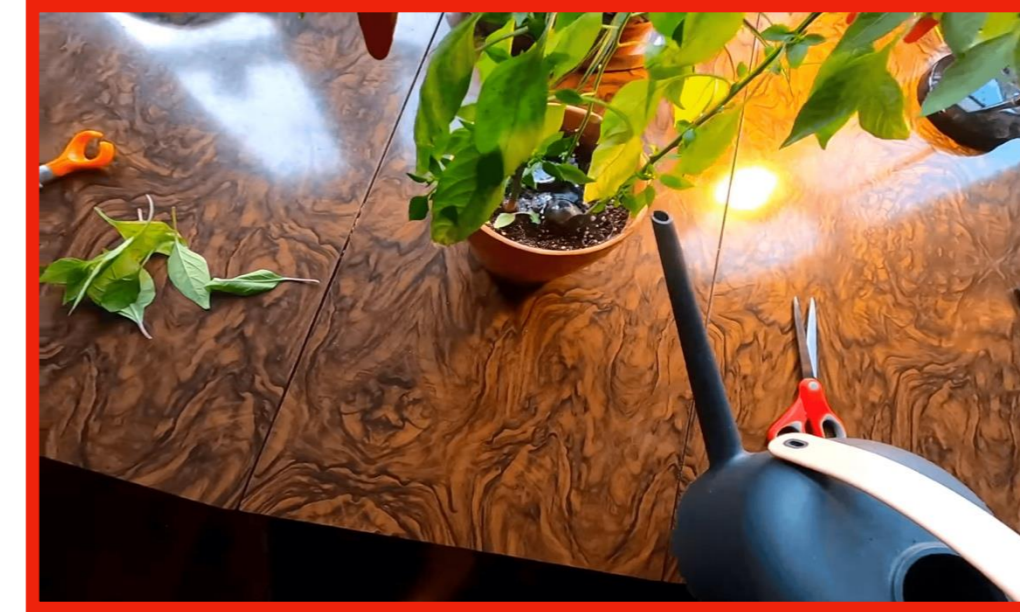


Instruction Text

Attempt Video



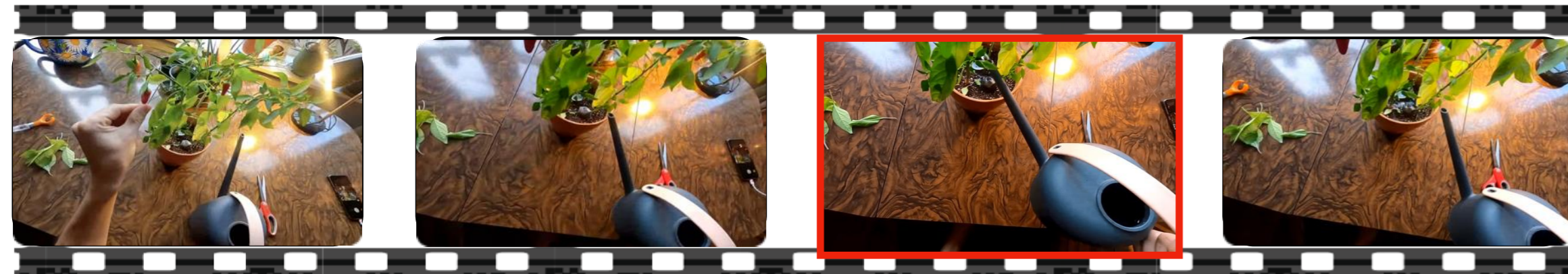
EgoMotion-COMPASS<sup>1</sup>: t=1.2s (Error=17.3%)



EgoT2<sup>2</sup>: t=3.4s (Error=24.5%)



**MisFormer (Ours)<sup>3</sup>: t=2.5s (Error=7.7%)**



Ground Truth Point of No Return: t=2.1s

1. Feichtenhofer et al. Masked Autoencoders for Egocentric Video Understanding. Ego4D Challenge Workshop, CVPR 2022.  
2. Zihui Xue, Yale Song, Kristen Grauman, Lorenzo Torresani. Egocentric Video Task Translation. CVPR 2023.  
3. Yayuan Li, Aditya Jain, Foteini Bellos, Jason Corso. Mistake Attribution: Fine-Grained Mistake Understanding in Egocentric Videos. CVPR 2026.

# Results: Spatial Attribution

Prune chili peppers



Point of No Return:  $t=2.1s$

Instruction Text

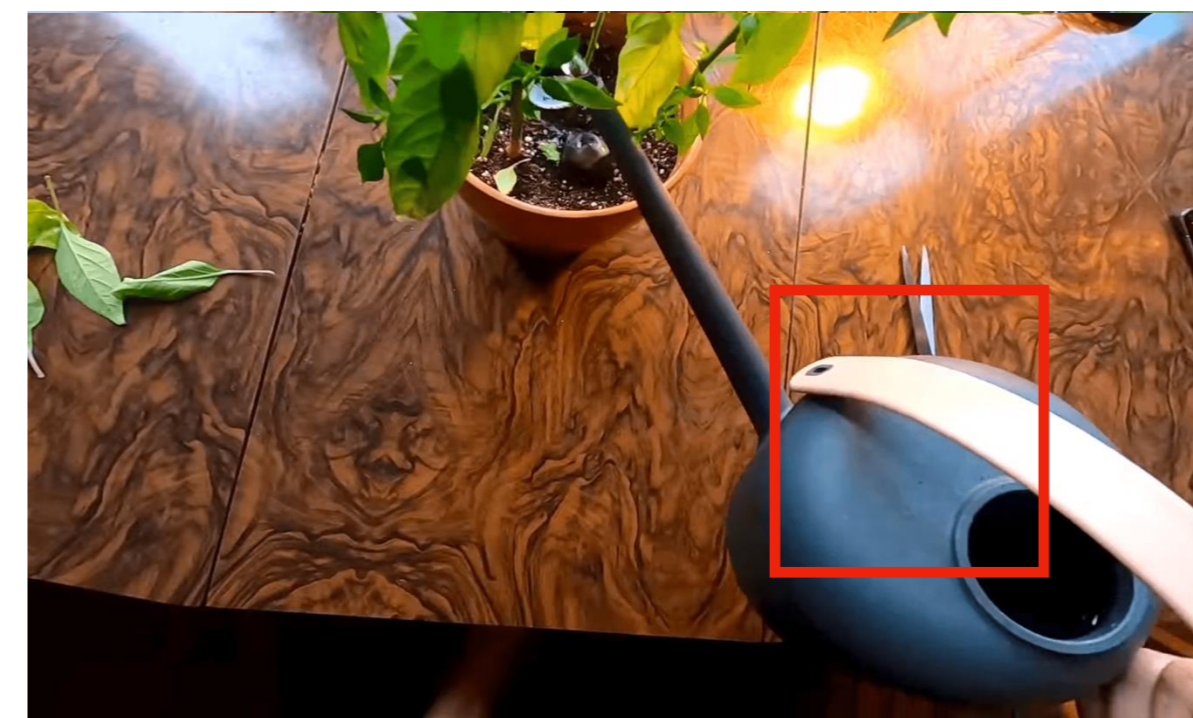
Attempt Video



Ground Truth



MediaPipe-U<sup>1</sup>: Not detected



SSDA<sup>2</sup>: IoU=12.1%



**MisFormer(Ours): IoU=48.5%**

1. Camillo Lugaresi, et al.. MediaPipe: A Framework for Perceiving and Processing Reality. CVPR Workshop, 2019.

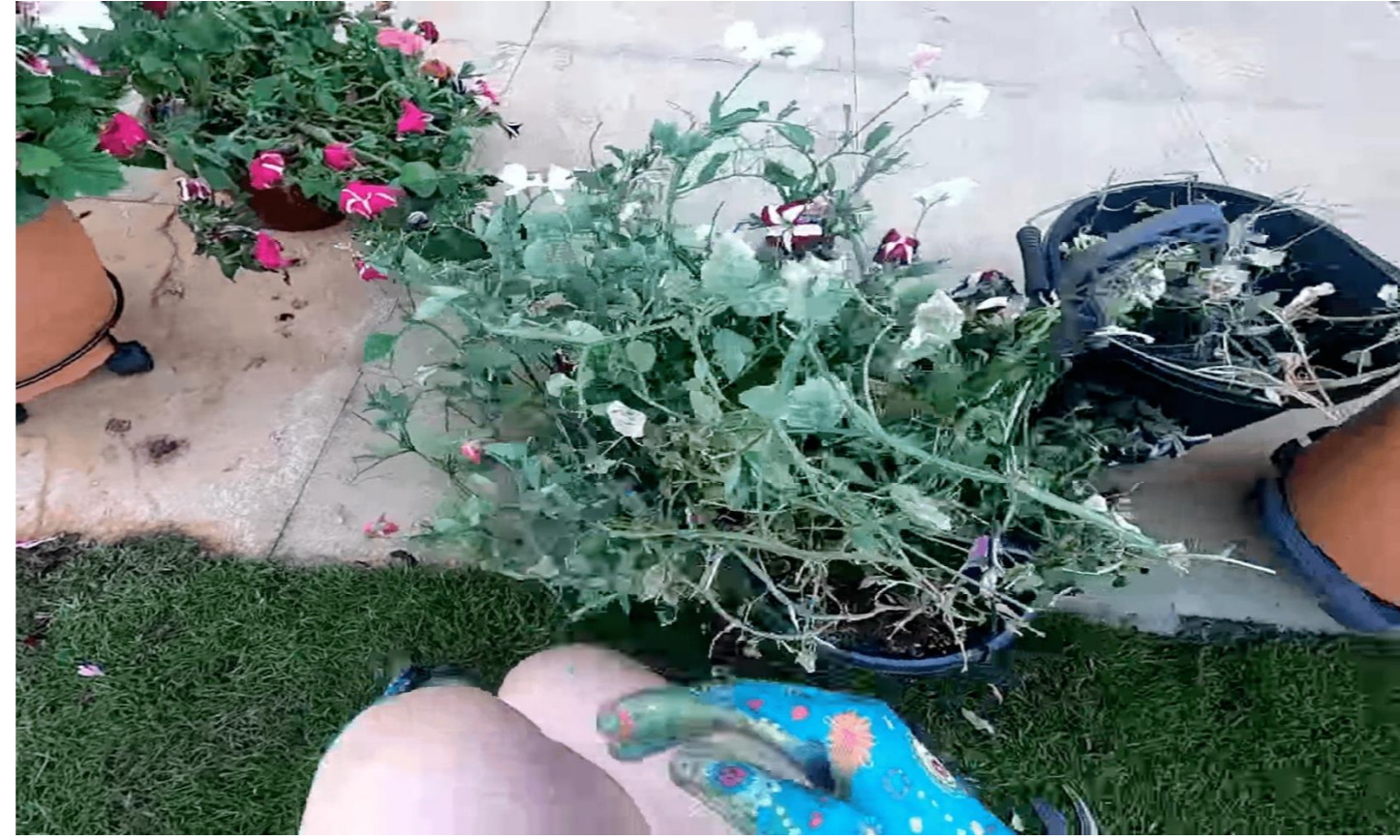
2. Rosario Leonardi, Antonino Furnari, Francesco Ragusa, Giovanni Maria Farinella. Are Synthetic Data Useful for Egocentric Hand-Object Interaction Detection? ECCV 2024.

3. Yayuan Li, Aditya Jain, Foteini Bellos, Jason Corso. Mistake Attribution: Fine-Grained Mistake Understanding in Egocentric Videos. CVPR 2026.

# More Results

Semantic

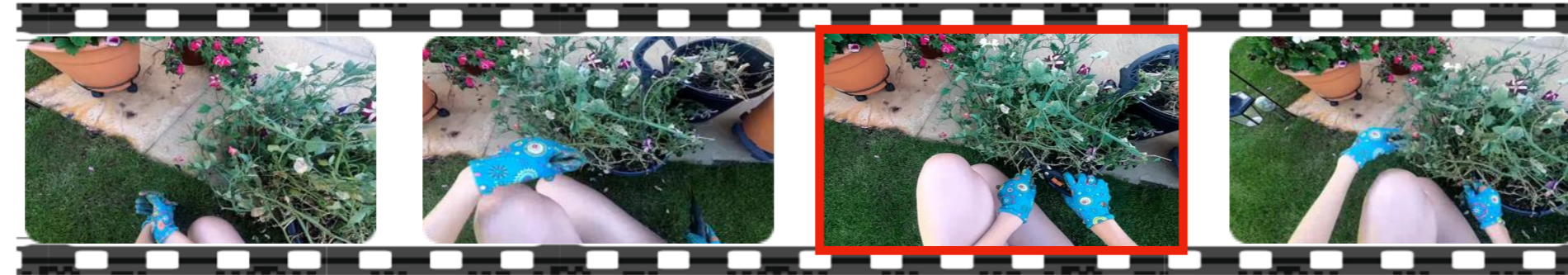
Tie stem



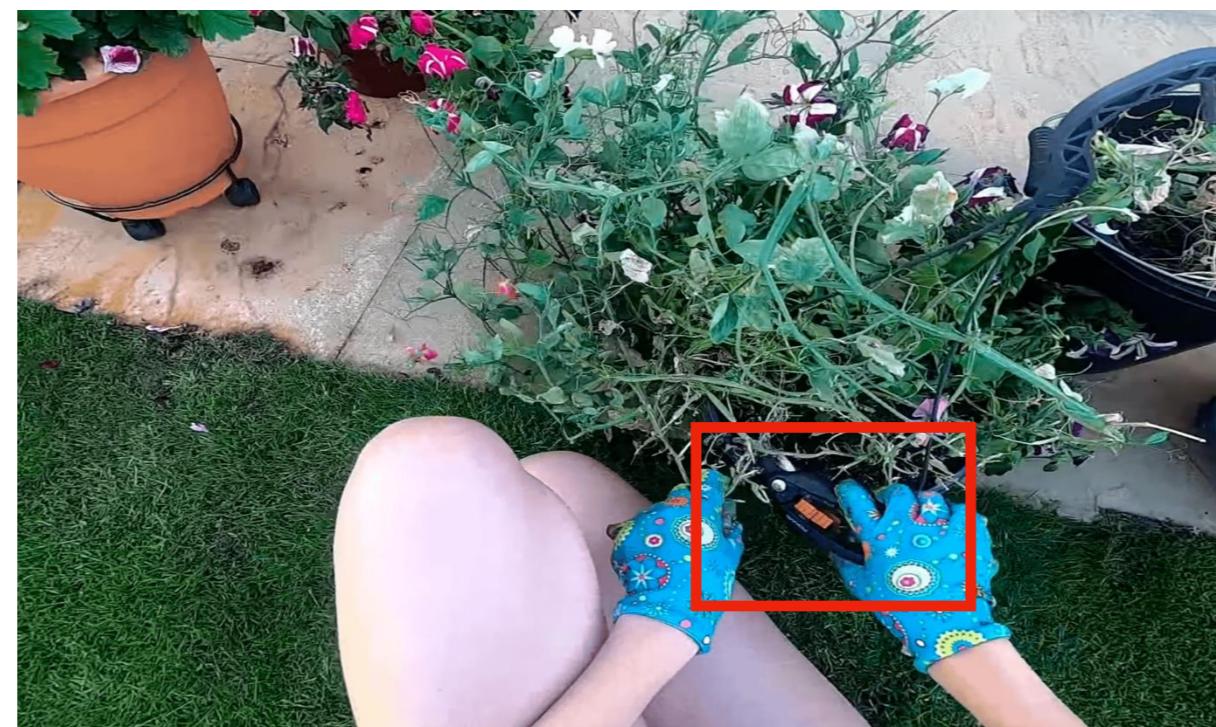
Take plate



Temporal



Spatial



Mistake?



00.00s  
00.04s

picks the small clear sachet

Mistake?



00.00s  
00.01s

picks the scraper

Mistake?



00.00s  
00.00s

dips the brush

Mistake?



00.00s  
00.00s

picks the shirt

Mistake?



00.00s  
00.00s

moves the bowl

Mistake?



00.00s  
00.00s

places the bottle

Mistake?



00.00s  
00.00s

opens a door

Mistake?



00.00s  
00.00s

smoothens the piece of wood

Mistake?



00.00s  
00.00s

takes a container

Mistake?



00.00s  
00.00s

scoops cement

Mistake?



00.00s  
00.00s

picks another stick

Mistake?



00.00s  
00.00s

dusts soil in the dirt

# Follow-up



Paper



Contact

## Mistake Attribution: Fine-Grained Mistake Understanding in Egocentric Videos

Yayuan Li<sup>1</sup>, Aadit Jain<sup>1</sup>, Filippos Bellos<sup>1</sup>, Jason Corso<sup>1,2</sup>

<sup>1</sup>University of Michigan, <sup>2</sup>Voxel51