



## Motivation

ActivityNet-Captions



19/30s - 63%  
"Walking the dog"

Ego4D

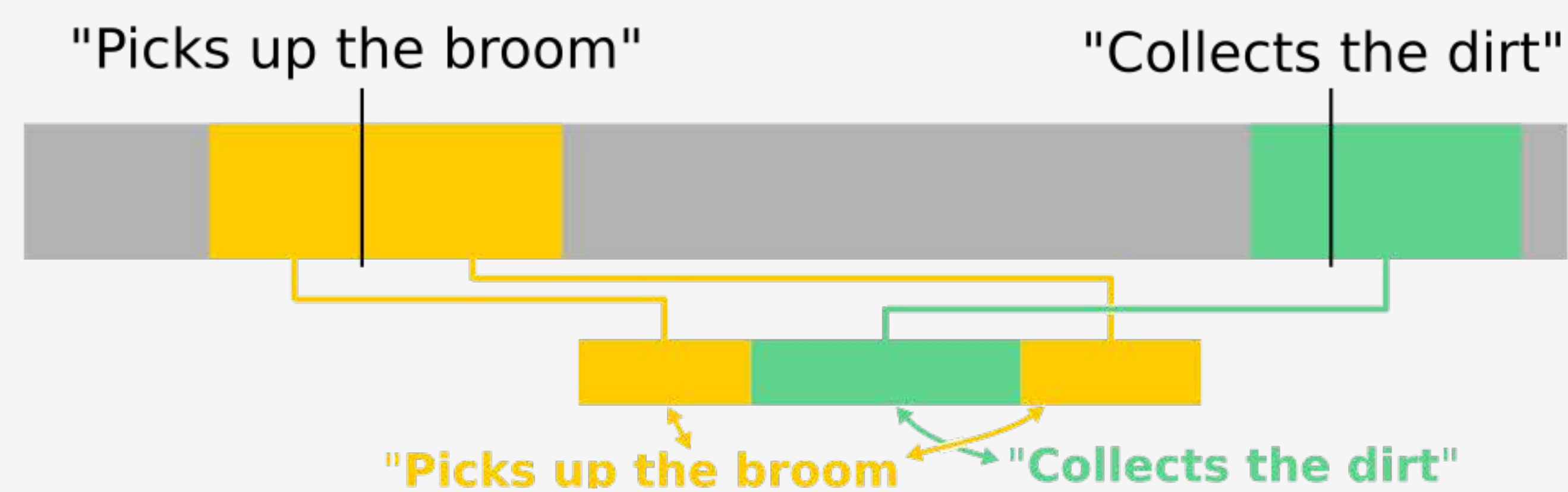


3/3314s - 0.1%  
"Opens the drawer"

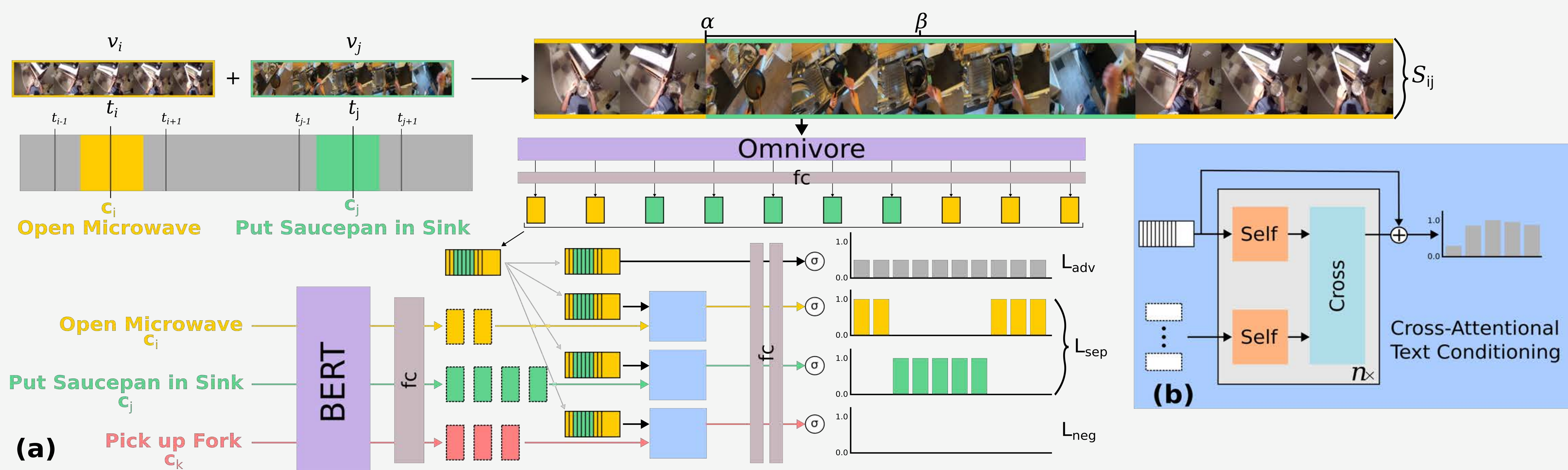
Comparison of Temporal Sentence Grounding Datasets

	Total Vid Duration	Avg. Vid Duration	Avg. Mom. Duration	Annotations / Video	Total Annotations	Avg. Coverage ↓
ANet-Captions [1]	487.6h	2.0min	37.1s	4.9	72k	30.90%
Charades-STA [2]	57.1h	0.5min	8.1s	2.3	16k	27.00%
DiDeMo [3]	88.7h	0.5min	6.5s	3.9	41k	21.70%
TACoS [4]	10.1h	4.8min	27.9s	143.6	18k	9.70%
Ego4D [5]	234.9h	17.3min	2.0s	214.1	223k	<b>0.19%</b>
EPIC-Kitchens [6]	73.4h	8.9min	3.1s	134.1	67k	<u>0.58%</u>

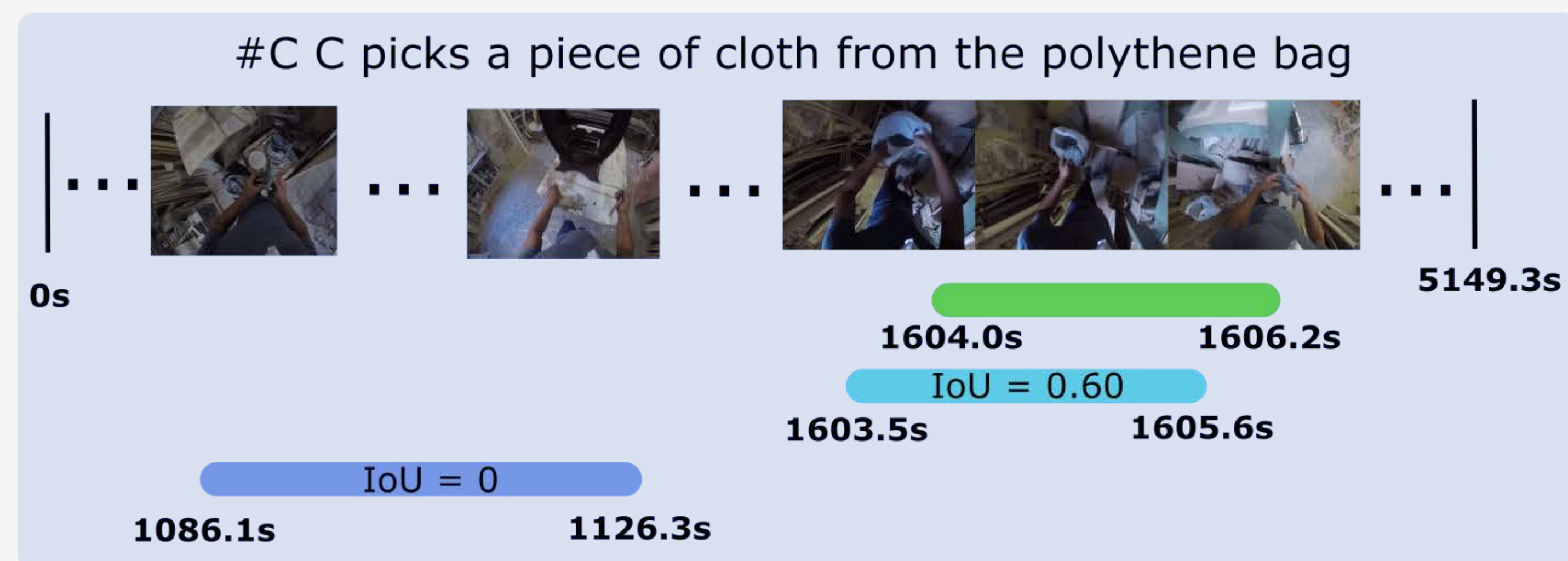
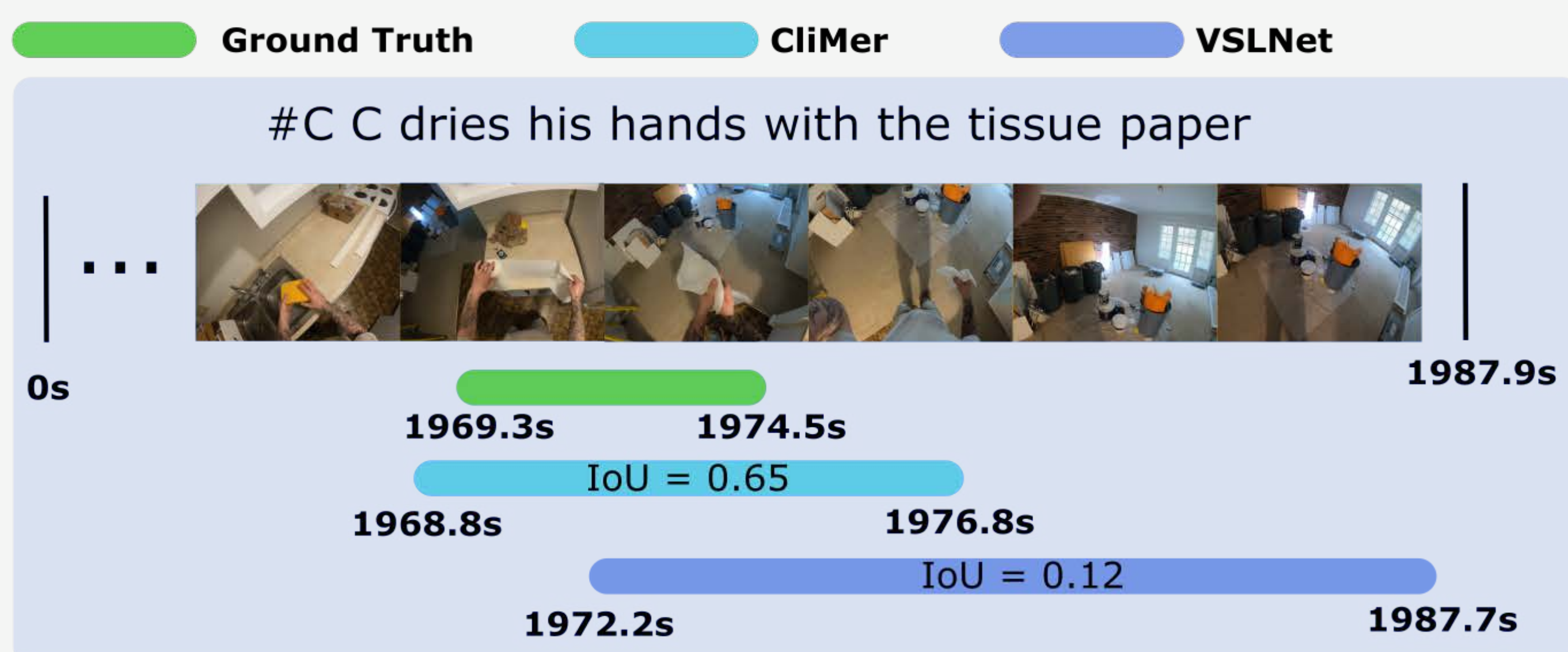
## Clip Merging & Rough Timestamps



## CliMer



## Results - Qualitative



## Results - Quantitative

	Ego4d			
	IoU=0.1	IoU=0.3	IoU=0.5	mR
Random Baseline	0.47	0.09	0.02	0.19
VSLNet [1]	7.32	3.16	1.35	3.94
CliMer	<b>9.68</b>	<b>5.03</b>	<b>2.24</b>	<b>5.65</b>

	EPIC-Kitchens			
	IoU=0.1	IoU=0.3	IoU=0.5	mR
Random Baseline	0.78	0.13	0.03	0.31
VSLNet [1]	19.32	8.76	3.90	10.66
CliMer	<b>22.20</b>	<b>11.57</b>	<b>5.25</b>	<b>13.01</b>

## Ablation - Clip Merging & Hard Negatives

Video	Merged Segment	Hard Negatives	IoU=0.1	IoU=0.3	IoU=0.5	mR
Video 1	✓	-	7.28	3.09	1.08	3.82
Video 2	-	✓	7.23	4.42	2.20	4.62
	✓	✓	<b>9.68</b>	<b>5.03</b>	<b>2.24</b>	<b>5.65</b>

## Conclusion

- Clip merging proves to be an effective method of training using clips segmented from rough timestamps.
- We explore Ego4D and EPIC-Kitchens for Temporal Sentence Grounding showcasing their difficulty.
- Qualitatively CliMer shows a better ability to pick out precise moments for given sentence queries.