







Temporally Grounding Instructional Diagrams in Unconstrained Videos

Jiahao Zhang¹ Frederic Z. Zhang² Cristian Rodriguez² Yizhak Ben-Shabat^{1,3} Anoop Cherian⁴ Stephen Gould¹

¹The Australian National University ²The Australian Institute for Machine Learning ³Technion Israel Institute of Technology ⁴Mitsubishi Electric Research Labs ¹{first.last}@anu.edu.au ²{first.last}@adelaide.edu.au ³sitzikbs@gmail.com ⁴cherian@merl.com







Problem Statement

Given an untrimmed video that has been evenly divided into N clips for feature extraction, we denote these clips by $\{V_i\}_{i=1}^N$, and the set of M diagrams corresponding to the video by $\{I_i\}_{i=1}^M$, we aim to develop a model capable of accurately predicting the **timespan** of each diagram $\mathbf{t}=(t_{S},t_{e})$, where $t_{\mathcal{S}}$ and $t_{\mathcal{E}}$ are the normalized start and end time of a segment.

2. Motivation & Contributions

Motivation

Dataset

Need to model multiple unbiased segments per video.

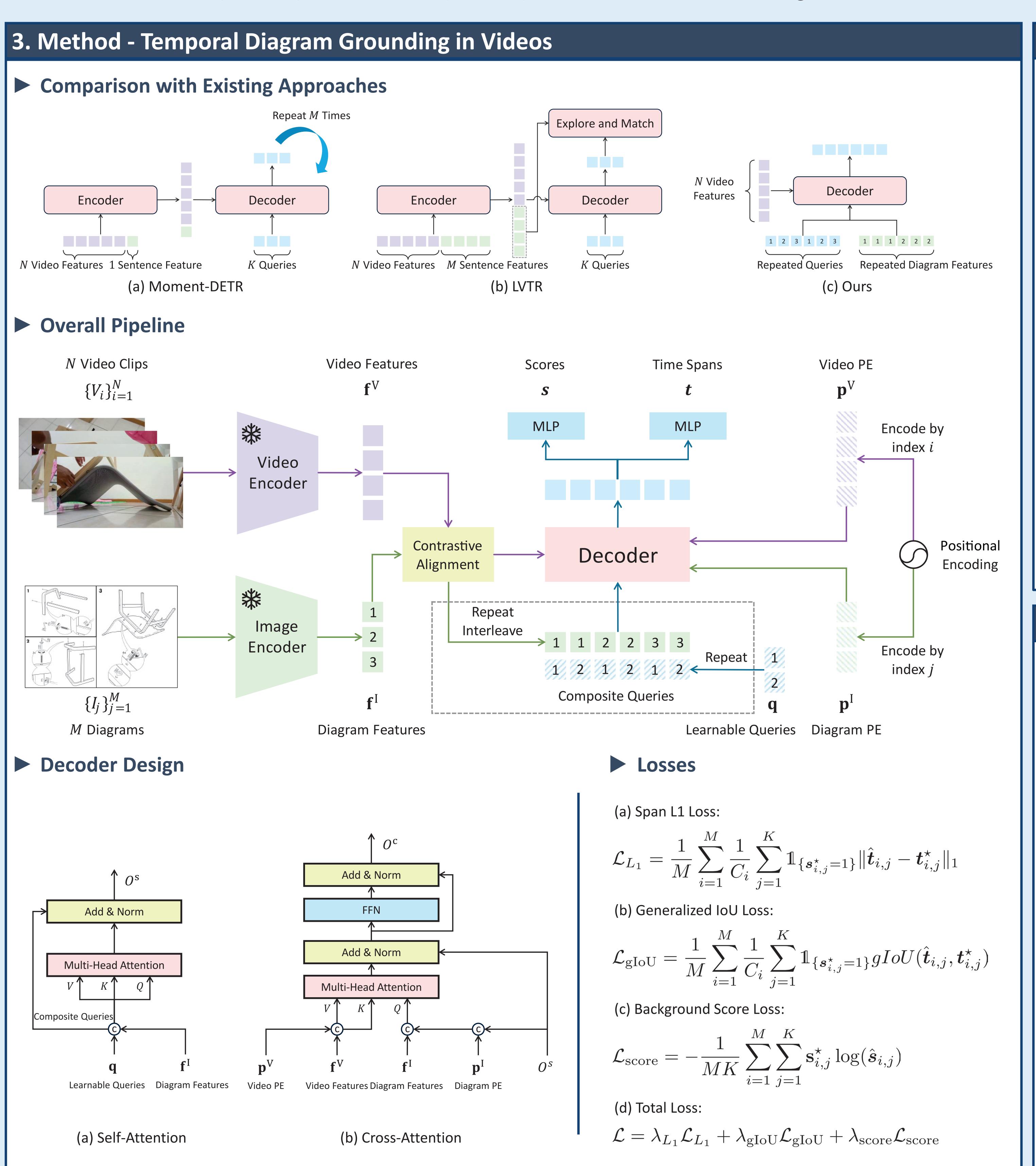
Year

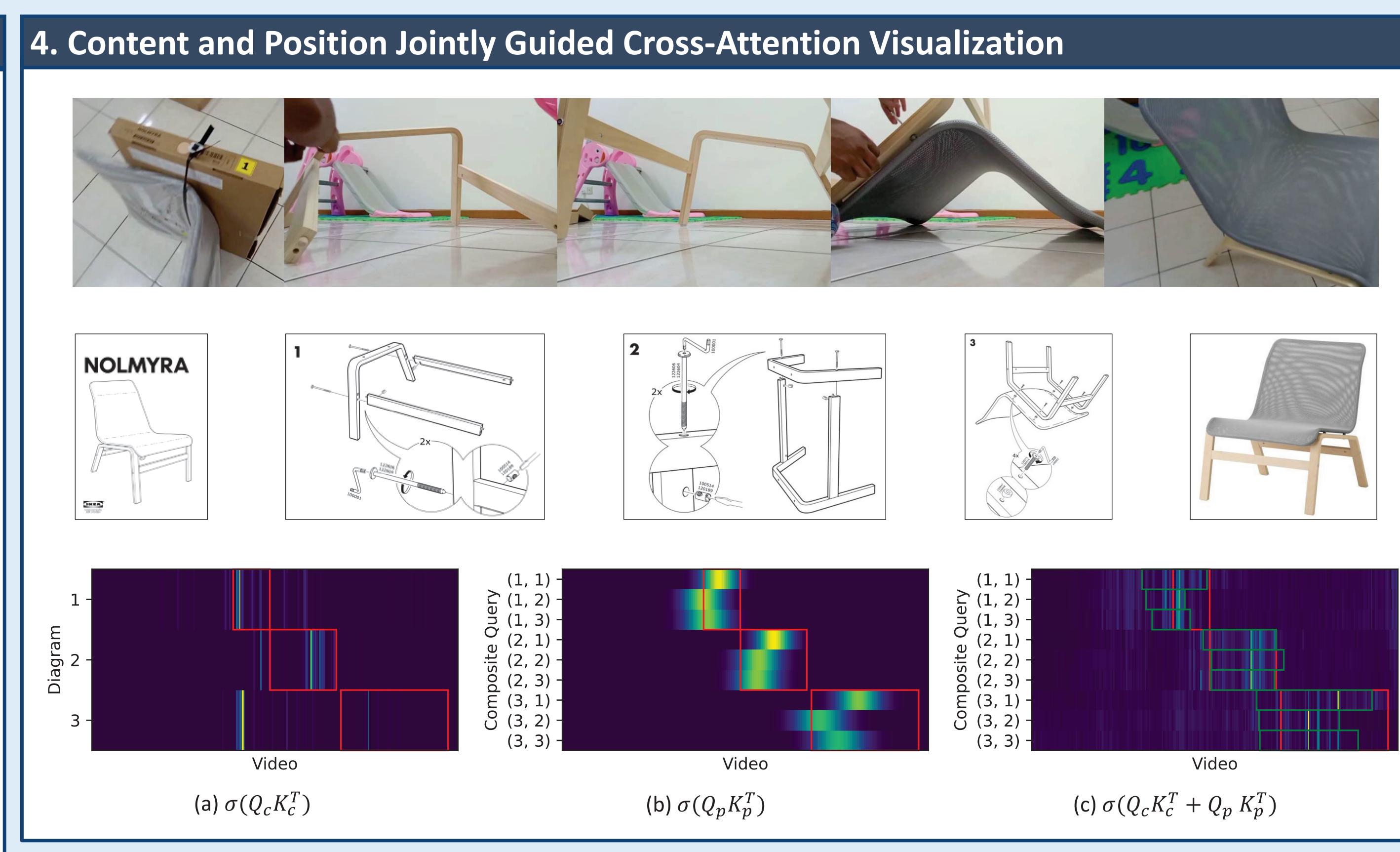
DiDeMo	2017	Max 30s	3.87	
Charades-STA	2017	Avg. 30.60s	2.42 4.82	
ActivityNet Captions	2017	Avg. 117.60s		
YouCook II	2017	Avg. 5.27 (Max 10) min	7.7	
kea Assembly in the Wild	2023	Avg. 11 (Min 1 - Max 79) min	15.57	
1.0 0.8 - 0.6 -	0.6 -	1.0 0.8 - 0.6 - UH 0.4 -	0.8 - 0.6 - 0.6 -	
0.2	0.0	0.2 -	0.4 - 0.2 - 0.4 0.6 0.8	

Video Duration

of Segments per Video

- Contributions
- 1. Sequential Grounding with Composite Queries: Developed a detection transformer that uses composite queries combining both content and positional priors to handle varying query lengths.
- 2. Enhanced Attention Mechanisms: Designed specialized self-attention masks and optimized cross-attention value choices for better grounding.
- State-of-the-Art Performance: Achieved significant improvements on IAW and YouCook2, demonstrating effectiveness for both diagram and sentence sequence grounding.





5. Final Result

Ikea Assembly in the Wild dataset.

Method	N/I1 -	R@1, IoU=			T - T T	
Medilou	Mode	0.3	0.5	0.7	mIoU	
Random	_	1.809	0.254	0.057	4.801	
LVTR	All	11.26	4.591	1.112	7.515	
2D-TAN conv	One	31.24	18.94	8.030	20.51	
2D-TAN pool	One	32.94	20.02	8.170	21.21	
Moment DETR	One	34.00	18.34	7.290	16.60	
Ours w/ Moment DETR	All	37.79	22.74	9.140	23.86	
EaTR	One	38.48	22.77	9.540	24.75	
Ours w/ EaTR	All	42.02	26.45	11.54	27.27	



Method	Text	R@1, IoU=			mIoU
IVICUIIOU		0.3	0.5	0.7	
DORi LocFormer	- BERT	$43.36 \\ 46.76$	$\frac{30.47}{31.33}$	$18.24 \\ 15.81$	$30.46 \\ 30.92$
ExCL TMLGA DORi Ours w/ EaTR	BERT **	26.63 34.77 42.27 52.95	16.15 23.05 29.90 36.28	8.51 12.49 18.38 18.50	18.87 24.42 29.92 35.32

