



# MomentDiff: Generative Video Moment Retrieval from Random to Real

Pandeng Li,<sup>1</sup> Chen-Wei Xie,<sup>2</sup> Hongtao Xie,<sup>1\*</sup> Liming Zhao,<sup>2</sup> Lei Zhang,<sup>1</sup>  
Yun Zheng,<sup>2</sup> Deli Zhao,<sup>2</sup> Yongdong Zhang<sup>1</sup>

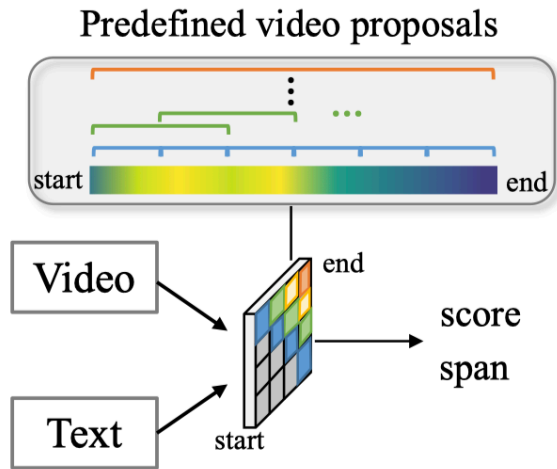
<sup>1</sup> University of Science and Technology of China

<sup>2</sup> Alibaba Group

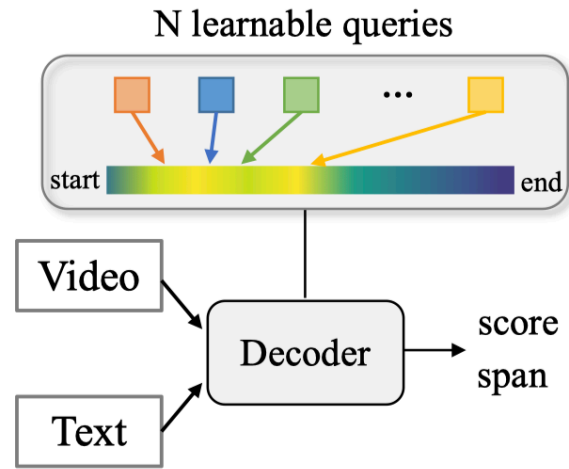
lpd@mail.ustc.edu.cn



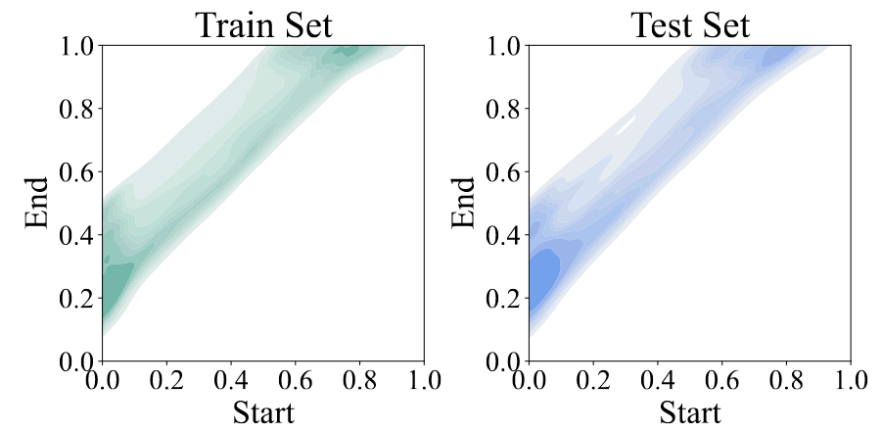
- Temporal location bias



(a) Dense: 2DTAN



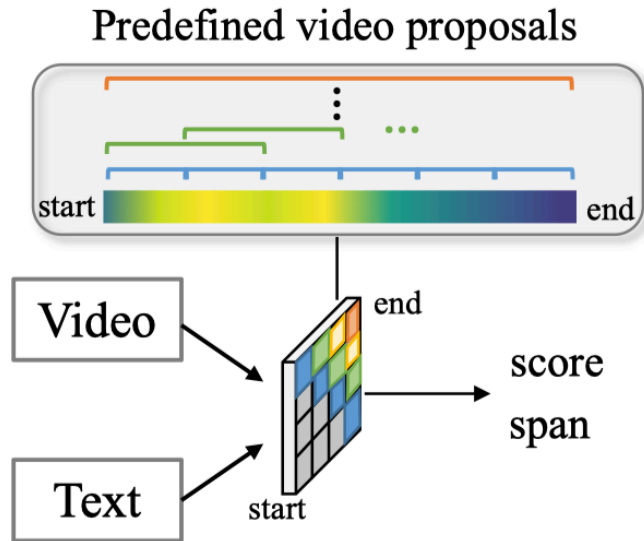
(b) Sparse: MomentDETR



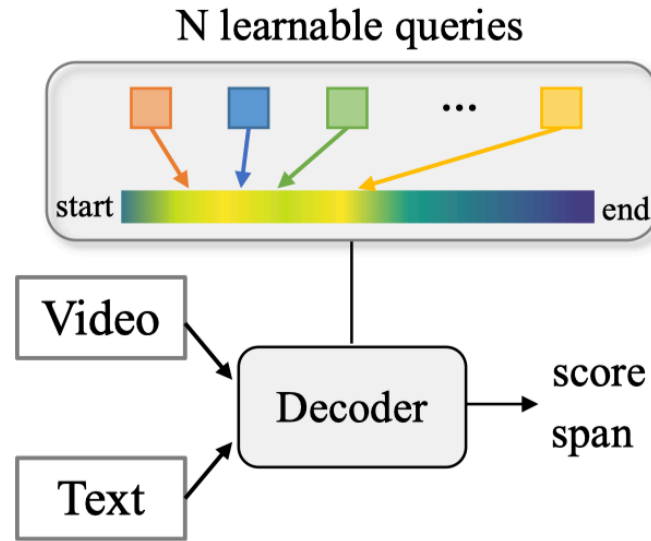
Moment distributions for Charades-STA [1]

[1] Otani, M., Y. Nakashima, E. Rahtu, et al. Uncovering hidden challenges in query-based video moment retrieval. In BMVC. 2020.

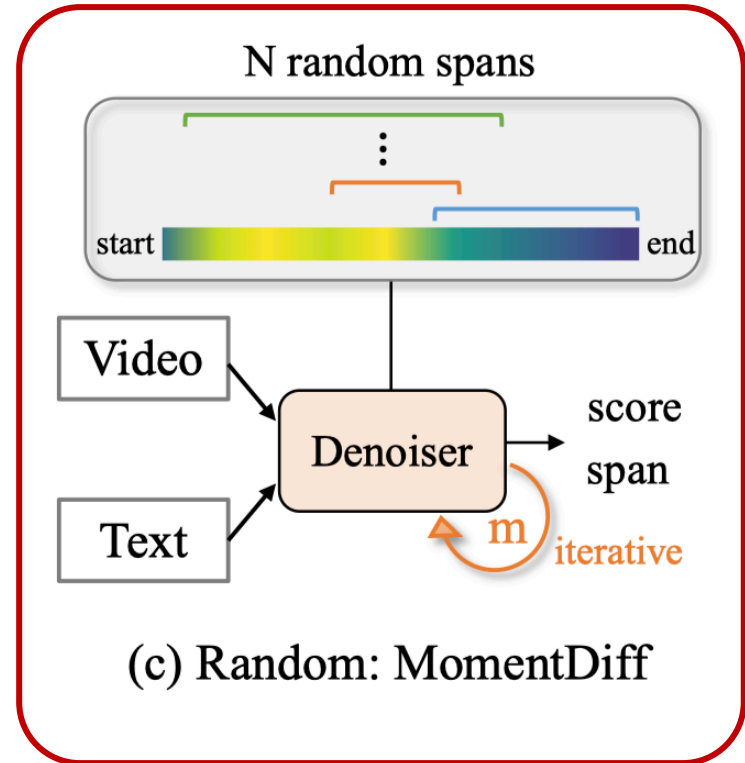
# Our idea



(a) Dense: 2DTAN

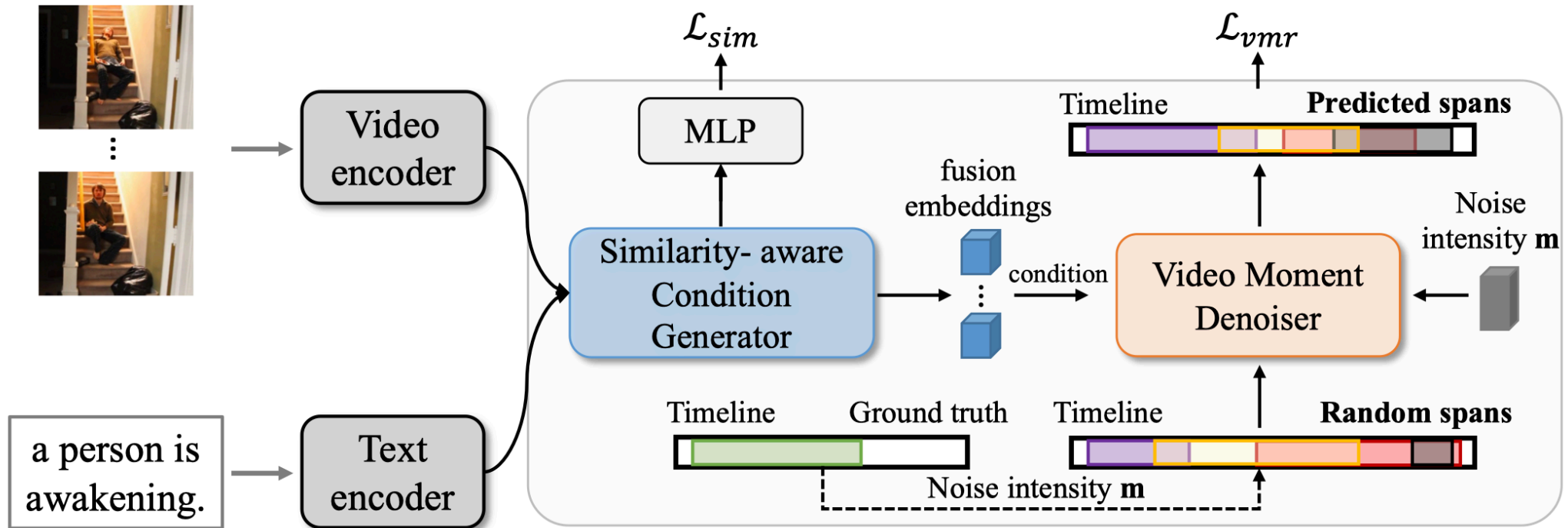


(b) Sparse: MomentDETR

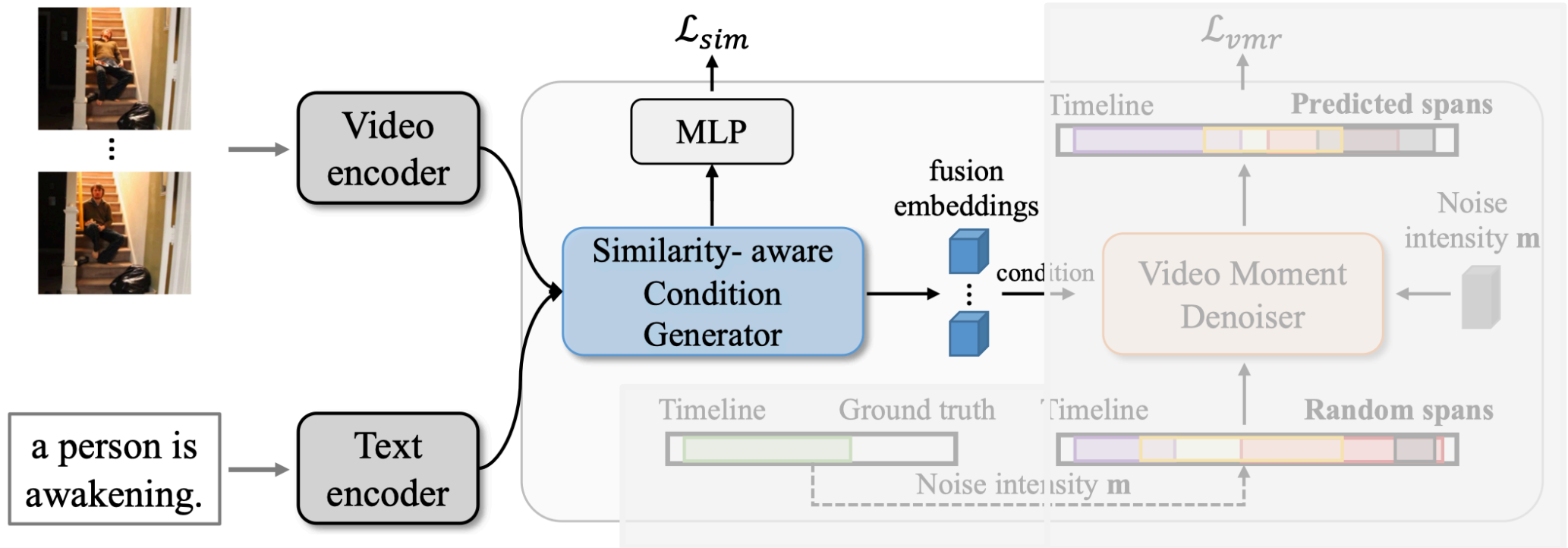


(c) Random: MomentDiff

- **Framework**

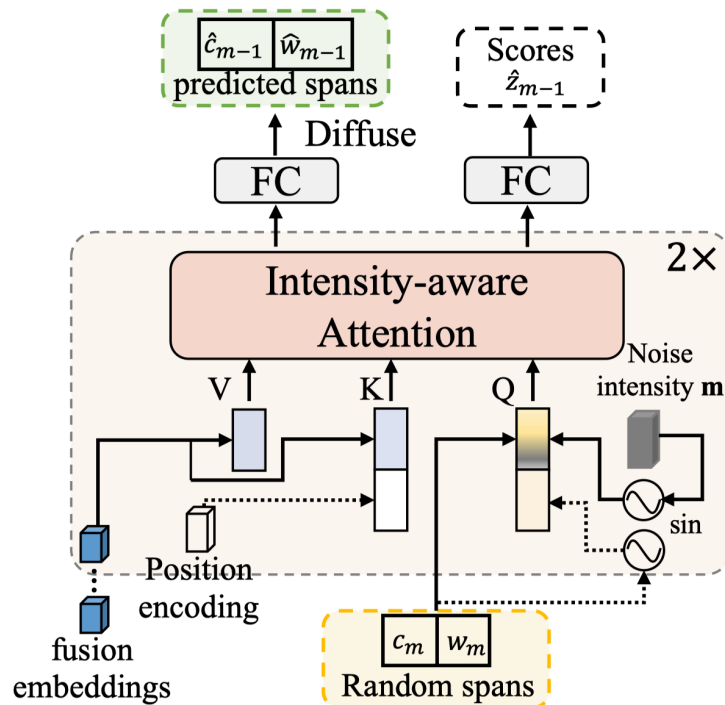


- Similarity-aware Condition Generator (SCG)**



$$\mathcal{L}_{sim} = -\frac{1}{N_v} \sum_{i=1} \mathbf{y}_i * \log(\mathbf{s}_i) + (1 - \mathbf{y}_i) * \log(1 - \mathbf{s}_i) + \frac{1}{N_s} \sum_{j=1} \max(0, \beta + \mathbf{s}_{n_j} - \mathbf{s}_{p_j})$$

- **Video Moment Denoiser (VMD)**



1. Span normalization
2. Span embedding
3. Intensity-aware attention
4. Denoising training

$$\mathcal{L}_{\text{vmr}}(\mathbf{x}_0, f_{\theta}(\mathbf{x}_m, m, \mathbf{F})) = \lambda_{L1} \|\mathbf{x}_0 - \hat{\mathbf{x}}_{m-1}\| + \lambda_{\text{iou}} \mathcal{L}_{\text{iou}}(\mathbf{x}_0, \hat{\mathbf{x}}_{m-1}) + \lambda_{\text{ce}} \mathcal{L}_{\text{ce}}(\hat{\mathbf{z}}_{m-1})$$



## Charades-STA

Table 1: Performance comparisons (%) on the Charades-STA dataset. "\*" denotes that we re-implement the method under the same training scheme. "A" stands for using audio data.

Method	Type	Charades-STA				
		R1@0.5	R1@0.7	MAP@0.5	MAP@0.75	MAP <sub>avg</sub>
MAN [32]	VGG, Glove	41.21	20.54	-	-	-
RaNet* [74]	VGG, Glove	42.91	25.82	53.28	24.41	28.55
2DTAN* [26]	VGG, Glove	41.34	23.91	54.68	24.15	29.26
DORi [77]	VGG, Glove	43.47	26.37	-	-	-
CBLN [43]	VGG, Glove	47.94	28.22	-	-	-
DCM [39]	VGG, Glove	47.80	28.00	-	-	-
MMN* [27]	VGG, Glove	46.93	27.07	58.85	28.16	31.58
MomentDETR* [36]	VGG, Glove	50.54	28.01	57.39	25.62	29.87
MomentDiff	VGG, Glove	<b>51.94</b>	<b>28.25</b>	<b>59.86</b>	<b>29.11</b>	<b>31.66</b>
UMT [37]	VGG+A, Glove	48.44	29.76	58.03	27.46	30.37
MomentDiff	VGG+A, Glove	<b>52.62</b>	<b>29.93</b>	<b>60.69</b>	<b>29.74</b>	<b>31.81</b>
DEBUG [78]	C3D, Glove	37.39	17.69	-	-	-
LPNet [35]	C3D, Glove	40.94	21.13	-	-	-
VSLNet* [29]	C3D, Glove	48.67	30.33	56.88	25.79	30.16
MomentDETR* [36]	C3D, Glove	50.49	29.95	56.27	26.08	29.92
MomentDiff	C3D, Glove	<b>53.79</b>	<b>30.18</b>	<b>59.32</b>	<b>29.85</b>	<b>31.89</b>
MomentDETR* [36]	SF+C, C	53.22	30.87	58.86	26.43	30.43
MomentDiff	SF+C, Glove	55.42	32.17	60.93	32.47	32.59
MomentDiff	SF+C, C	<b>55.57</b>	<b>32.42</b>	<b>61.07</b>	<b>32.51</b>	<b>32.85</b>





## QVHighlights

Table 2: Performance comparisons (%) on QVHighlights with SF+C video features and CLIP text features. "★" notes that we re-implement the method with only segment encode videos. MDE is the abbreviation of MomentDETR [36].

Method	QVHighlights				
	R1@0.5	R1@0.7	MAP@0.5	MAP@0.75	MAP <sub>avg</sub>
MCN [24]	11.41	2.72	24.94	8.22	10.67
CAL [79]	25.49	11.54	23.40	7.65	9.89
XML [80]	41.83	30.35	44.63	31.73	32.14
XML+ [80]	46.69	33.46	47.89	34.67	34.90
MDE★ [36]	53.56	34.09	53.97	28.65	29.39
MomentDiff	<b>57.42</b>	<b>39.66</b>	<b>54.02</b>	<b>35.73</b>	<b>35.95</b>
UMT★† [37]	56.26	40.31	52.77	36.82	35.79
MomentDiff†	<b>58.21</b>	<b>41.48</b>	<b>54.57</b>	<b>37.21</b>	<b>36.84</b>

## TACoS

Table 3: Performance comparisons (%) on TACoS. We adopt C3D features to notes that we re-implement the method with only segment encode videos. MDE is the abbreviation of MomentDETR [36].

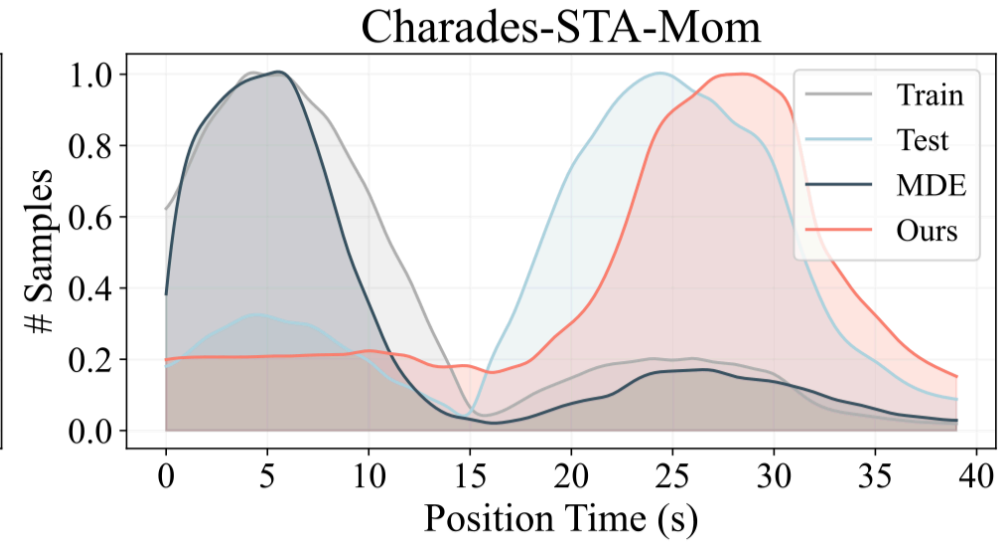
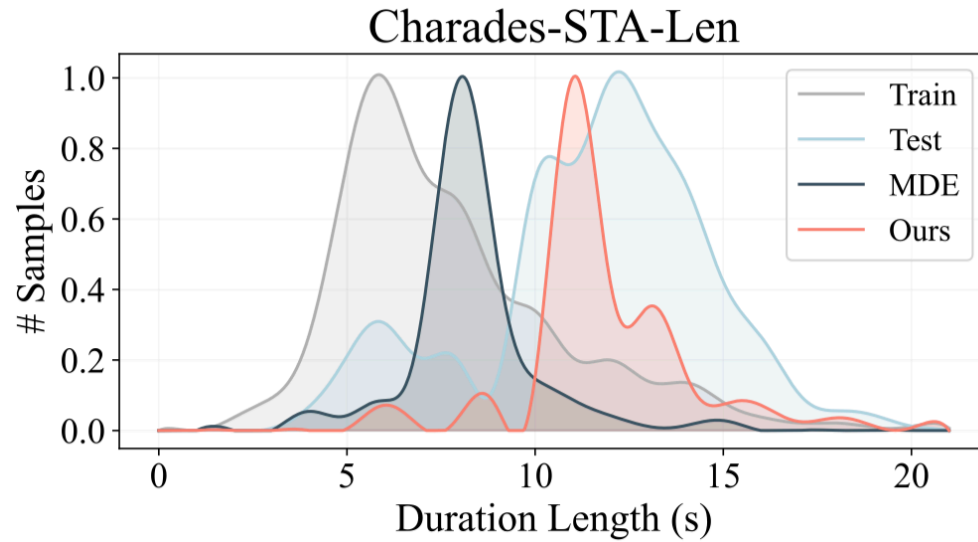
Method	TACoS		
	R1@0.1	R1@0.3	R1@0.5
CTRL [22]	24.32	18.32	13.30
SCDM [31]	-	26.11	21.17
DRN [30]	-	-	23.17
DCL [40]	49.36	38.84	29.07
CBLN [43]	49.16	38.98	27.65
FVMR [41]	53.12	41.48	29.12
RaNet [74]	-	43.34	33.54
MDE★ [36]	41.16	32.21	20.55
MMN★ [27]	51.39	39.24	26.17
MomentDiff	<b>56.81</b>	<b>44.78</b>	<b>33.68</b>



## R1@n and MAP results on public OOD datasets

Method	Charades-CD				ActivityNet-CD			
	R1@0.3	R1@0.5	R1@0.7	MAP	R1@0.3	R1@0.5	R1@0.7	MAP
2DTAN	49.71	28.95	12.78	12.60	40.04	22.07	10.29	12.77
MMN	55.91	34.56	15.84	15.73	44.13	24.69	12.22	15.06
MomentDETR	57.34	41.18	19.31	18.95	39.98	21.30	10.58	12.19
MomentDiff	<b>67.73</b>	<b>47.17</b>	<b>22.98</b>	<b>22.76</b>	<b>45.54</b>	<b>26.96</b>	<b>13.69</b>	<b>16.38</b>

## Statistical distributions on our anti-bias datasets



Temporal moment  $\hat{\mathbf{x}}_0 = (\hat{\mathbf{c}}_0, \hat{\mathbf{w}}_0)$

→ The center time of moment

→ The duration length of moment

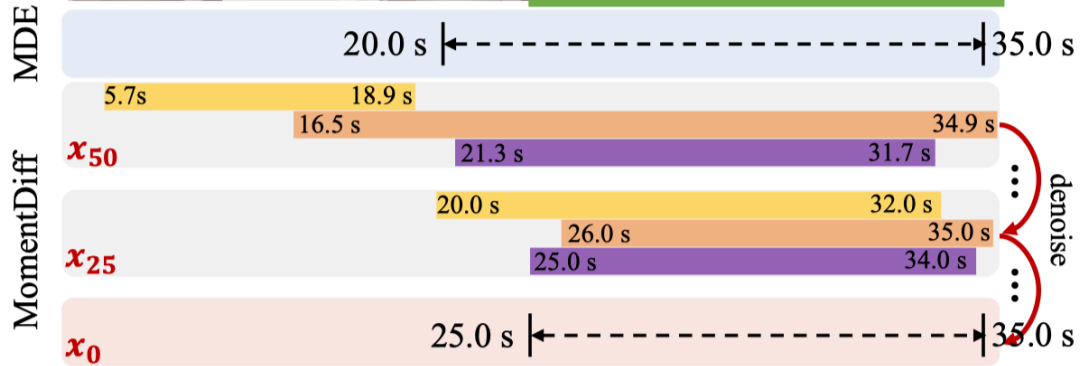


## R1@n and MAP results on anti-bias datasets

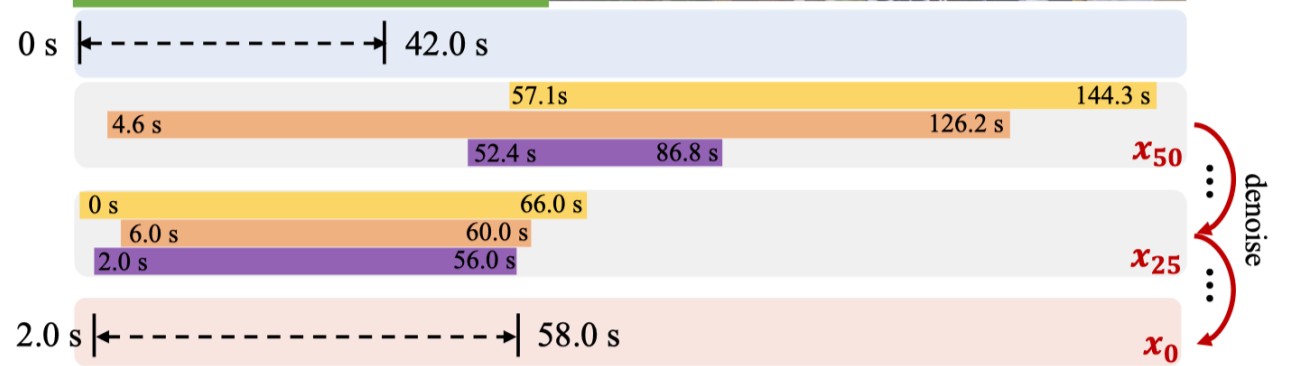
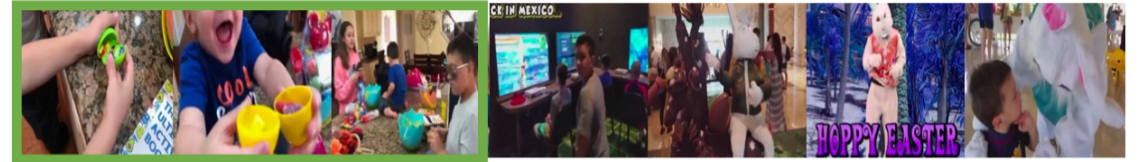
Method	Charades-STA-Len				Charades-STA-Mom			
	R1@0.3	R1@0.5	R1@0.7	MAP	R1@0.3	R1@0.5	R1@0.7	MAP
2DTAN	39.68	28.68	17.72	22.79	27.81	20.44	10.84	17.23
MMN	43.58	34.31	19.94	26.85	33.58	27.20	14.12	19.18
MomentDETR	42.73	34.39	16.12	24.02	29.94	22.16	11.56	18.66
MomentDiff	<b>51.25</b>	<b>38.32</b>	<b>23.38</b>	<b>28.19</b>	<b>48.39</b>	<b>33.59</b>	<b>15.71</b>	<b>21.37</b>

## Visualization of the diffusion process

*Query: The person starts drinking a glass of coffee.*



*Query: Kids checking out their goodies and chocolates during Easter.*





# MomentDiff: Generative Video Moment Retrieval from Random to Real

Thank you!

Github Code



Homepage

