## APPENDIX FOR ACTIONS INSPIRE EVERY MOMENT: ONLINE ACTION AUGMENTED DENSE VIDEO CAPTIONING

## Anonymous authors

Paper under double-blind review

## A APPENDIX

006

007

012

013

A.1 ADDITIONAL IMPLEMENTATION DETAILS

We utilize the CLIP model pretrained on the LAION-2B dataset. Specifically, we use its ViT-Large model (303M parameters) and its 12-layer Transformer text model (128M parameters). The CLIP-initialized ViT is kept frozen throughout all stages of training described below.

We further pretrain the CLIP text model on the image captioning task using the same LAION-2B dataset, with batch size 1024 for 0.2 epochs. We use the Adam optimizer with momentum 0.9, an initial learning rate (LR) of 5e-5, 5000 warmup steps, linear LR decay, weight decay 1e-2.

For dense video captioning, the token reduction Transformer and autoregressive Transformer modules are added. Each module consists of 8 layers with a model dimension of 512, and 32 M parameters. In total, the entire model contains approximately 500M parameters.

025 As described in section 3.3, we apply image-based simulated video pretraining on the entire model, 026 including the newly added modules. To simulate video sequences, we sample 3 to 5 images from 027 the LAION-2B dataset, repeating each image multiple times to form a 16-frame sequence. To create 028 smoother transitions, we blend pixels at the boundaries by applying a weighted sum of two images, 029 using a randomly selected blending ratio  $\alpha \in [0.1, 0.9]$ , e.g., blending pixels of images A and B as  $\alpha A + (1 - \alpha)B$ . We apply random augmentations to each frame to avoid overly monotonous 031 sequences. This pretraining follows the same segment-by-segment autoregressive framework for 032 online dense video captioning. We use a batch size of 32 and train the model for 100000 steps. The optimizer is Adam with momentum 0.9, an initial LR of 1e-4, 5000 warmup steps, cosine LR decay, 033 and a weight decay of 1e-5. 034

When finetuning on dense video captioning, the model is trained for 20000 steps with a batch size 16. We again use the Adam optimizer with momentum 0.9, an initial LR of 1e-4, 5000 warmup steps, cosine LR decay and a weight decay of 1e-5. For time tokenization, we use relative time tokens following Vid2Seq (Yang et al., 2023). We quantize a video of duration T frames into B = 32equally spaced time bins.

For inference, we follow the standard protocol to use beam search, with a beam size of 6 and temperature 1, followed by temporal NMS with a threshold of 0.7 to remove overlapping intervals.

042 043 044

A.2 COMPUTATIONAL COST OF OUR MODEL

Our model uses 410 GFLOPs per segment, totaling 6560 GFLOPs for 16 segments. The retrieval operates on the precomputed text embeddings.

047 048

049

A.3 PROMPTING FOR ACTION TEXT CORPUS CONSTRUCTION

We construct a corpus of action phrases to serve as contextual priors during the dense video captioning process. These phrases are designed to capture key actions or events in each video segment and
identify relevant objects. To build this corpus, we draw from two main sources: 1) Captions from
the training splits of dense video captioning datasets: ViTT, YouCook2, and ActivityNet, where we
use only the text captions excluding the video frames. 2) A broader, less domain-specific corpus

from the HowTo100M dataset (Miech et al., 2019), again using only the subtitle text without the 055 video content. 056

Unlike previous methods that use raw video captions, which tend to be lengthy and unfocused (Xu 057 et al., 2024), we summarize these captions into concise action phrases using the publicly available 058 language model, Gemma (Team et al. (2024), huggingface.co/google/gemma-2-27b).

To improve the extraction process, we refine the prompt to ensure concise and consistently formatted 060 action phrases. Specifically, we emphasize singular nouns, avoid numerical terms, and enforce a 061 strict format of <action verb(ing)> <target object (if any)>. For example, our prompt is: Your goal 062 is to summarize the input sentence using as few words as possible. Focus on the words describing 063 actions or events. Use singular nouns, avoid articles and numeric terms. Respond in the format of 064 <action verb (ing)> <target object (if any)>. Input: {raw caption}. Answer: 065

For HowTo100M subtitles, which are often longer, we adjust the prompt to focus on extracting a 066 single main action or event: The input is video subtitle text. Choose the main action or event in the 067 video and summarize it using as few words as possible. Focus on the words describing actions or 068 events. Use singular nouns, avoid articles and numeric terms. Respond in the format of *<action* 069 *verb(ing)* <*target object (if any)*. *Input: {video subtitles}. Answer:* 070

071 These prompts effectively generate concise action phrases. After processing all text in each source, 072 we deduplicate phrases by merging those with the same set of words, regardless of word order. After filtering out some least frequent phrases, we obtain 30,000 action phrases for each corpus. 073

074 075

076

A.4 ADDITIONAL ABLATIONS

077 We conduct additional ablation studies using the same setup described in section 4.2, where we use 16 frames per video, with 1 frame per segment at a resolution of  $256 \times 256$  pixels, and report the 079 results on the ViTT dataset. In this ablation, the mixed training (Table 6) and image-based simulated 080 video pretraining (table 8) are *not* used, unless otherwise noted. 081

082 Effect of text decoder pretraining. In Table A1, we show the effect of pretraining the text decoder 083 for image captioning using the LAION-2B dataset (section 3.1). While image captioning pretraining 084 improves performance, our model still performs reasonably well without it. Notably, most recent 085 dense video captioning methods (Yang et al., 2023; Wang et al., 2024; Ren et al., 2024; Wu et al., 2024; Zhou et al., 2024) employ language pretraining for their text decoders, and we follow this 086 approach to enhance the captioning performance. 087

088 Ablation on the size of action text corpus. Table A2 presents the effect of varying the size of 089 the action text corpus. We randomly sample subsets of the corpus at 1%, 10%, 50%, and 100%. 090 Increasing the corpus size improves performance, with the most notable gains observed up to 50%. 091 This shows that our constructed corpus is effective in covering a broad range of action phrases for 092 retrieval augmentation. 093

Ablation on the number of segments. Table A3 studies the effect of the number of segments. 094 Overall, more frames improves performance. Across different segment configurations, our model 095 performs robustly overall, with 16 segments yielding the best results. 096

097 Ablation on ViTT, YouCook2, ActivityNet datasets. In Table A4, we evaluate the effects of our 098 key method components. The baseline refers to the online model described in section 4.1. We ob-099 serve both our main contributions – online action-augmentation (section 3.2) and image-based simulated video pretraining (section 3.3) – make complimentary improvements to performance across 100 all three benchmarks: ViTT, YouCook2, and ActivityNet. 101

102 103

- A.5 COMPARISON OF APPROACHES IN EXISTING METHODS 104
- 105

Table A5 compares various strategies used in existing methods, focusing on key aspects such as 106

support for online video captioning, reliance on video-text pretraining, and the backbone models 107 employed.

108		text decoder	S	С	М						
109		~	9.9	37.2	9.6						
110		Х	I	8.3	30.4	8.0					
111	Table A1: Ablation on tex	at model pro	e <b>training</b> o	n the i	mage o	caption	ing tas	sk. Ro	esults	on Vi	TT.
112		size of actior	n text corpus	S	С	М					
113	_	19	%	8.7	33.0	8.4					
114		10	%	9.3	35.2	9.1					
115		50 <sup>4</sup> 100	% )%	9.7	36.7	9.5					
116	T 11 40 1		• •		51.2	).0 D					
117	Table A2: 1	Effect of size	e of action	text c	orpus.	Resul	ts on V	/11°1.			
118		# segments	# frames	S	С	М					
119		8	8	9.3	36.2	8.9					
120		8	16	9.6	37.0	9.1					
101		16 16	16	9.9	37.2	9.6					
121		32	32	0.7	36.8	9.8					
122		32	52	9.1	50.8	9.0					
123	Table A3: Number of segm	ents which o	controls the	e numł	per of c	lecodir	ig outp	outs. 1	Result	s on V	/iTT.
124		I	ViTT			YouCoo	k2	1	Activ	vitvNet	r
125	method	S	C M	<b>F</b> 1	S	C N	F1	S	C	M	F1

method	S	С	Μ	F1	S	С	Μ	F1	S	С	M	F1	
baseline online model	7.6	27.7	7.2	34.0	5.3	27.7	6.8	22.4	5.4	31.6	9.8	45.8	
online action-augmentation	9.8	37.4	9.4	35.5	6.9	39.4	8.0	25.5	6.7	34.6	11.2	46.2	
image-based simulated video pretraining	8.9	34.1	8.5	37.8	6.1	35.0	7.2	27.8	6.2	33.7	10.4	47.6	
both combined	10.8	39.1	10.3	39.2	8.0	45.6	9.3	30.7	7.5	36.4	12.1	49.9	

Table A4: Ablation of our method on ViTT, YouCook2, ActivityNet datasets. This ablation uses S=16 segments per video, L=1 frame per segment at a resolution of  $256 \times 256$  pixels.

method	online	video-text pretraining	backbone
E2ESG (Zhu et al., 2022)	N	Ø	C3D
PDVC (Wang et al., 2021)	N	Ø	TSN
OmniViD (Wang et al., 2024)	N	Kinetics	VideoSwin + Bart
TimeChat (Ren et al., 2024)	N	YT-Temporal, ViTT, ActivityNet, etc.	Eva-CLIP-G + Llama-7B
Vid2Seq † (Yang et al., 2023)	N	YT-Temporal-1B	CLIP-L + Bert-B
DoYou (Kim et al., 2024)	N	ø	CLIP-L
DIBS (Wu et al., 2024)	N	Howto100M	CLIP-L
Streaming (Zhou et al., 2024)	Y	YT-Temporal-1B	CLIP-L + Bert-B
AIEM (ours)	Y	Ø	CLIP-L

Table A5: Comparison to the state-of-the-art on dense video captioning.

## References

130

142

143

- Minkuk Kim, Hyeon Bae Kim, Jinyoung Moon, Jinwoo Choi, and Seong Tae Kim. Do you remember? dense video captioning with cross-modal memory retrieval. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 13894–13904, 2024.
- Antoine Miech, Dimitri Zhukov, Jean-Baptiste Alayrac, Makarand Tapaswi, Ivan Laptev, and Josef Sivic. Howto100m: Learning a text-video embedding by watching hundred million narrated video clips. In *ICCV*, 2019.
- Shuhuai Ren, Linli Yao, Shicheng Li, Xu Sun, and Lu Hou. Timechat: A time-sensitive multimodal large language model for long video understanding. In *Proceedings of the IEEE/CVF Conference* on Computer Vision and Pattern Recognition, pp. 14313–14323, 2024.
- Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya
  Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. Gemma: Open
  models based on gemini research and technology. *arXiv preprint arXiv:2403.08295*, 2024.
- Junke Wang, Dongdong Chen, Chong Luo, Bo He, Lu Yuan, Zuxuan Wu, and Yu-Gang Jiang.
   Omnivid: A generative framework for universal video understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 18209–18220, 2024.
- 161 Teng Wang, Ruimao Zhang, Zhichao Lu, Feng Zheng, Ran Cheng, and Ping Luo. End-to-end dense video captioning with parallel decoding. In *ICCV*, 2021.

162 163 164	Hao Wu, Huabin Liu, Yu Qiao, and Xiao Sun. Dibs: Enhancing dense video captioning with un- labeled videos via pseudo boundary enrichment and online refinement. In <i>Proceedings of the</i> <i>IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 18699–18708, 2024.
165	Lilen Vy, Vifei Hyang, Juplin Hay, Cya Chan, Vysiis Zhang, Dyi Fang, and Waidi Via, Dataisyal
166	augmented exponentric video cantioning. In Proceedings of the IEEE/CVE Conference on Com-
167 168	puter Vision and Pattern Recognition, pp. 13525–13536, 2024.
169	Antoine Vang Arsha Nagrani Paul Hongsuck Seo. Antoine Miech, Jordi Pont-Tuset, Ivan Lantev
170	Iosef Sivic and Cordelia Schmid Vid2seq: Large-scale pretraining of a visual language model
171	for dense video captioning. CVPR, 2023.
172	
173	Xingyi Zhou, Anurag Arnab, Shyamal Buch, Shen Yan, Austin Myers, Xuehan Xiong, Arsha Na-
174	grani, and Cordelia Schmid. Streaming dense video captioning. In <i>Proceedings of the IEEE/CVF</i>
175	Conference on Computer Vision and Pattern Recognition, pp. 18243–18252, 2024.
176	Wanrong Zhu, Bo Pang, Ashish V. Thapliyal, William Yang Wang, and Radu Soricut. End-to-end
177	dense video captioning as sequence generation. In COLING, 2022.
178	
179	
180	
181	
182	
183	
184	
185	
186	
187	
188	
189	
190	
191	
192	
193	
194	
195	
196	
197	
198	
199	
200	
201	
202	
203	
204	
205	
200	
208	
209	
210	
211	
212	
213	
214	
215	