
VideoEval: Vision-Centric Benchmark Suite for Low-Cost Evaluation of Video Foundation Model

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 With the accumulation of high-quality data and advancements in visual pretrain-
2 ing paradigms, recent Video Foundation Models (VFMs) have made significant
3 progress, demonstrating remarkable performance on popular video understanding
4 benchmarks. However, conventional benchmarks (e.g. Kinetics) and evaluation
5 protocols are limited by their relatively poor diversity, high evaluation costs, and
6 saturated performance metrics. In this work, we introduce a comprehensive bench-
7 mark suite to address these issues, namely **VideoEval**. We establish the **Video Task**
8 **Adaption Benchmark (VidTAB)** and the **Video Embedding Benchmark (VidEB)**
9 from two perspectives: evaluating the task adaptability of VFMs under few-shot
10 conditions and assessing their feature embedding’s direct applicability to down-
11 stream tasks. With VideoEval, we conduct a large-scale study of 20 popular
12 open-source vision foundation models. Our study reveals some insightful findings,
13 1) overall, current VFMs exhibit weak generalization across diverse tasks, 2) in-
14 creasing video data, whether labeled or in video-text pairs, does not necessarily
15 improve task performance, 3) the effectiveness of some pre-training paradigms
16 may not be fully validated in previous benchmarks, and 4) combining different pre-
17 training paradigms can help develop models with better generalization capabilities.
18 We believe this study serves as an important complement to the current evaluation
19 methods for VFMs and offers valuable insights for future research directions.

20 1 Introduction

21 The field of deep learning is experiencing a significant paradigm shift due to the emergence of
22 foundation models (FMs). These models, exemplified by BERT [1], GPT [2, 3, 4], CLIP [5] and
23 Stable Diffusion [6], are trained on massive and diverse data at scale and demonstrate remarkable
24 adaptability to a broad spectrum of downstream tasks.

25 In the realm of video understanding, early researchers train backbone networks [7, 8, 9, 10] using
26 visual classification tasks on large-scale labeled datasets like ImageNet [11] and Kinetics [12]. How-
27 ever, the high cost associated with labeled data promotes the development of self-supervised learning
28 methods that capitalize on unlabeled data for visual pre-training [13, 14, 15, 16, 17]. Furthermore,
29 researchers delve into multimodal pre-training utilizing large-scale visual-text pairs [18, 19, 20, 21],
30 thereby enhancing their models’ capabilities and demonstrating impressive zero-shot performance.
31 Overall, fueled by the accumulation of high-quality image and video data and advancements in visual
32 pre-training paradigms, Video Foundation Models (VFMs) witness remarkable progress in recent
33 years. A new generation of VFMs [15, 16, 22, 23, 24, 25, 26] emerges, demonstrating outstanding
34 performance on conventional video understanding benchmarks.

35 The rapid development of VFMs raises the problem: *How to evaluate a video foundation model?*
36 In image realm, Previous works assess the generalization capability of Image Foundation Models

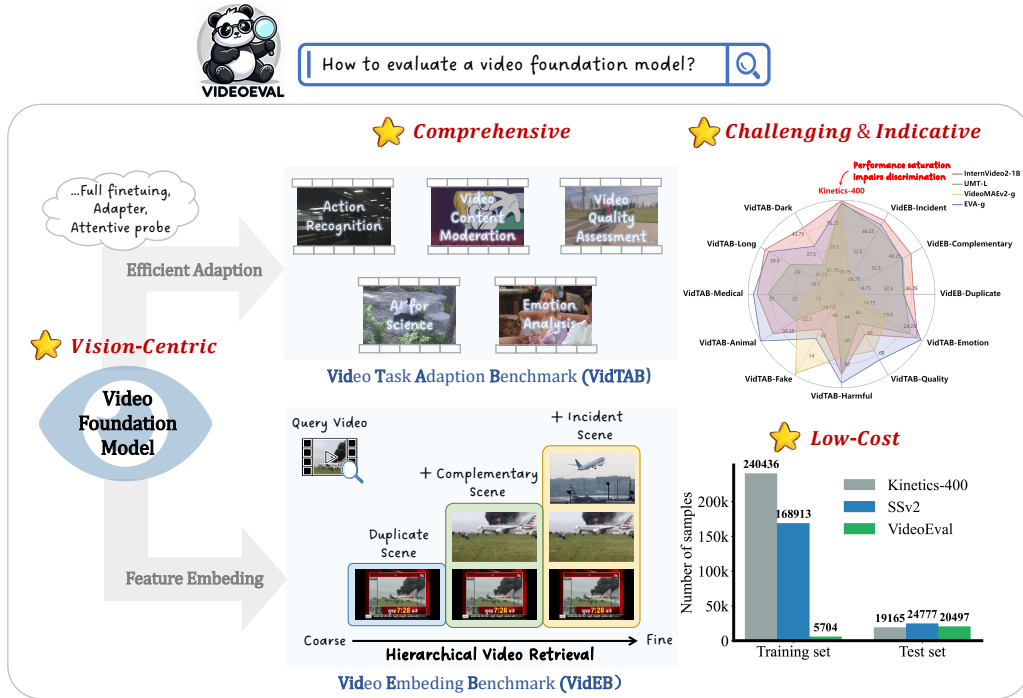


Figure 1: **Overview of VideoEval.** We propose a novel, vision-centric evaluation method for video foundation models that is comprehensive, challenging, indicative, and low-cost.

37 (IFMs) by evaluating their performance on numerous downstream visual tasks, encompassing diverse
 38 scenarios and evaluation protocols [27, 28, 29, 30, 31, 32, 33]. However, previous works primarily
 39 evaluates VFMs through benchmarks focusing on action recognition tasks [16, 23, 34]. Some
 40 studies [25, 26, 24] have also considered combining language models to evaluate performance on
 41 multimodal tasks. There are **several problems with current evaluation methods**: (1) Benchmarks
 42 like Kinetics [12], Something [35] and AVA [36], which focus on action recognition, overlook other
 43 video understanding scenarios (e.g., video quality assessment), limiting their applicability in evaluat-
 44 ing the generalization capabilities of visual foundation backbones across diverse video understanding
 45 applications. (2) The performance of VFMs on conventional benchmark [37] has reached a saturation
 46 point (90% Top-1 accuracy), making it challenging to differentiate between the true capabilities of
 47 different VFMs. (3) The high validation costs associated with conventional evaluation protocols,
 48 which often necessitate end-to-end training on the entire dataset, pose a significant challenge, particu-
 49 larly for large VFMs. (4) Incorporating language models may introduce bias when evaluating VFMs,
 50 as performance differences might stem from the language model rather than the VFMs itself.

51 To tackle these problems, we build a comprehensive benchmark suite for evaluation of VFMs, namely
 52 VideoEval. As shown in Figure 1, our method has the following key features: **Comprehensive**:
 53 First, we created the Video Task Adaptation Benchmark (VidTAB) to evaluate the adaptability of
 54 VFMs to unseen tasks with limited samples. We collected public datasets from various video task
 55 domains, including action recognition in special scenarios, AI for science, video content moderation,
 56 video quality/aesthetic assess, and emotion analysis. From these domains, we constructed eight
 57 adaptation tasks and developed evaluation protocols and adaptation methods suitable for current
 58 VFMs. Additionally, to assess the capability of VFMs’ feature embedding for downstream applica-
 59 tions, we created the Video Embedding Benchmark (VidEB), which includes four tasks that evaluate
 60 embedding at different granularities. **Challenging & Indicative**: Due to the diversity of test data
 61 and the effectiveness of our evaluation protocols, our VideoEval can effectively distinguish between
 62 various VFMs that perform similarly on traditional benchmarks, providing deeper insights into their
 63 true capabilities. **Low-cost**: Thanks to our training-light few-shot evaluation and training-free feature
 64 embedding evaluation protocols, VideoEval requires significantly fewer training samples compared to
 65 previous benchmarks, while maintaining a comparable number of testing samples to ensure accurate
 66 and stable evaluations. **Vision-centric**: Our evaluation focuses solely on the Video FMs themselves,
 67 avoiding the introduction of biases that may arise from incorporating language models.

68 Based on VideoEval, we evaluate 20 open-source vision foundation models, including VFMs,
 69 Image Foundation Models (IFMs), and IFMs with image-to-video methods. **Our main findings as**
 70 **following:** First, current VFMs still struggle to adapt to unseen video tasks with limited training
 71 samples. Second, while more data and larger models generally improve performance, augmenting
 72 video training data can sometimes negatively affect certain tasks. Third, the effectiveness of certain
 73 pre-training paradigms, such as VideoMAEv2 [22], may not have been adequately validated in
 74 previous benchmarks. Finally, combining multiple pre-training paradigms can lead to models with
 75 better generalization capabilities, such as performing multimodal contrastive learning after unimodal
 76 visual self-supervised pre-training [21, 26].

Table 1: **Comparison of VFMs Benchmark.** "Num. training" denotes number of training samples, "Num. test" denotes number of test samples, and "Beyond Action" denotes the tasks in this benchmark extend beyond action understanding. Compared to previous benchmarks, our VideoEval framework achieves more comprehensive and reliable evaluations at a lower cost.

Benchmark	Num. training	Num. test	Beyond Action	Task Diversity	Domain Diversity	VFMs-specific protocol
<i>Single-dataset Benchmarks</i>						
Kinetics-400 [37]	240,436	19,165	✗	✗	✗	✗
Sth-Sth V2 [38]	168,913	24,777	✗	✗	✗	✗
Moment-in-Time [39]	791,246	33,898	✗	✗	✗	✗
UCF101 [40]	9,537	3,783	✗	✗	✗	✗
<i>Multi-dataset Benchmarks</i>						
SEVERE [41]	868,446	144,830	✗	✓	✓	✗
BEAR [42]	240,236	140,436	✗	✓	✓	✗
VideoGLUE [34]	1,896,621	239,011	✗	✓	✓	✓
VideoEval	5,704	20,497	✓	✓	✓	✓

77 2 Related work

78 **Video foundation models** With the continuous growth of image [43, 44, 45] and video data [46,
 79 20, 47, 48, 49] and advancements in pre-training paradigms, research on Video Foundation Models
 80 (VFMs) has progressed rapidly. Current VFMs are primarily built around two pre-training paradigms:
 81 masked video modeling based on unimodal video data [15, 16, 22, 17, 50, 51, 52] and video-text
 82 contrastive learning based on multimodal visual-text pairs [18, 53, 19, 54, 20]. Some works [25,
 83 21, 24] combine these paradigms, enabling VFMs to extend further into multimodal understanding.
 84 Additionally, some studies introduce modalities like audio and speech on top of video and text [47,
 85 48, 26], further expanding the capabilities of VFMs. Recently, InternVideo2 [26] leverages mature
 86 pre-training paradigms and large-scale high-quality data to scale VFMs to 6 billion parameters,
 87 achieving remarkable performance improvements.

88 **Evaluation of VFMs** Previous works primarily utilize action recognition benchmarks focused on
 89 appearance and motion [12, 38, 36] to evaluate VFMs. To enhance evaluation diversity, some studies
 90 explore richer domains and tasks [55, 42, 56], but they remain limited to action recognition tasks. The
 91 InternVideo series [25, 26] and VideoGLUE [34] attempt to provide a more comprehensive evaluation
 92 of VFMs by expanding the number of benchmarks and evaluation protocols. However, these efforts
 93 are still based on existing benchmarks and incurred high validation costs. In contrast, our work
 94 considers the characteristics and application scenarios of VFMs, offering a comprehensive and low-
 95 cost evaluation solution through task definition and evaluation protocols, aimed at rapidly verifying
 96 the generalization capabilities of VFMs—a crucial aspect currently lacking in the community’s
 97 development of these models.

98 3 Building VideoEval

99 We argue that a powerful video foundation model should possess two key capabilities: (1) strong task
 100 adaptation ability, i.e., the ability to *adapt to diverse, unseen tasks with limited training samples*, and
 101 (2) the capacity to *extract feature embedding that retain and distill key information from videos*, di-
 102 rectly supporting various downstream tasks. From these perspectives, we construct VideoEval, which
 103 includes the Video Task Adaptation Benchmark (VidTAB) and the Video Embedding Benchmark
 104 (VidEB). By creating diverse task scenarios and employing efficient evaluation methods, VideoEval
 105 can quickly and comprehensively assess the generalization ability of VFMs in video understanding. In

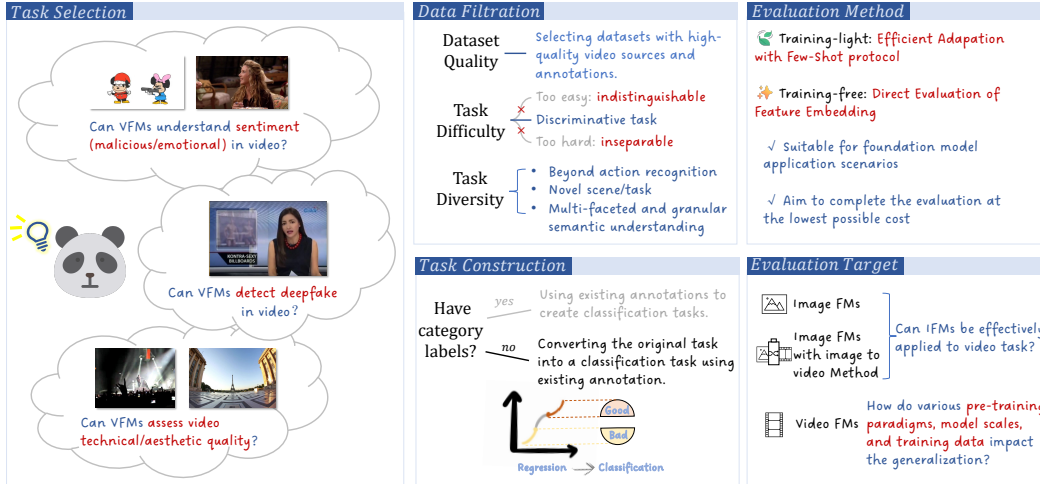


Figure 2: Illustration of building VideoEval.

Table 2: Task details of VideoEval. All videos are collected from the public datasets for building tasks of VidTAB and VidEB.

Domain	Task	Source	Task Description
<i>Video Task Adaptation Benchmark (VidTAB)</i>			
Action Recognition in Special Scenarios	Action Recognition in Dark Scene	ARID [57]	Recognizing 11 distinct human actions in dark scenarios. e.g. Run / Walk / Drink
	Action Recognition in Long Video	BreakFast [58]	Classifying 10 types of long-duration cooking videos. e.g. Milk / Tea / Sandwich
AI for Science	Medical Surgery	SurgicalActions160 [59]	Classifying 16 surgical actions in gynecologic laparoscopy. e.g. Knotting / Suction / Injection
	Animal Behavior	Animal Kingdom [60]	Classifying 12 behaviors of wild animals from diverse environmental footage. e.g. Flying / Chirping / Preening
Video Content Moderation	Fake Face	FaceForensics++ [61]	Determine whether the faces in the video have been tampered with by AI technology (such as DeepFake). e.g. Origin video / Video with fake face
	Harmful Content	mob [62]	Detecting 3 degrees of malicious content within videos. e.g. Obscene / Indecent activity / Violent activity
Video Quality Assessment	Quality Assess	DOVER [63]	Evaluating videos from an aesthetic and technical perspective and categorizing them into low and high quality. e.g. Low quality / High quality
Emotion Analysis	Emotion Analysis	CAER [64]	Classifying 7 different human emotions in video. e.g. Happy / Fear / Anger
<i>Video Embedding Benchmark (VidEB)</i>			
Scene Understanding in Temporal Contexts	Duplicate Scene Retrieval	FIVR5K [65]	Retrieve Duplicate Scene Videos (DSV): Videos captured by the same camera and sharing at least one scene (without considering any application transformations).
	Complementary Scene Retrieval	FIVR5K [65]	Retrieve Complementary Scene Videos (CSV): Retrieve a portion of the same spatiotemporal segment captured from different perspectives.
	Incident Scene Retrieval	FIVR5K [65]	Retrieve Incident Scene Videos (ISV): The same event is close in both space and time, but there are no overlapping videos.
	Copy Detection	DVSC23 [66]	Detecting edited versions of the same source video. Given a query inserted with one or more copied segments, detect the source video from the database.

106 this section, we present our VideoEval in detail. The construction pipeline for VideoEval is illustrated
 107 in Figure 2, and the evaluation tasks we ultimately constructed are presented in Table 2.

108 3.1 Video Task Adaption Benchmark

109 **Collecting diverse dataset from public source.** Previous benchmarks primarily focus on evaluating
 110 video models based on human actions, overlooking many other tasks requiring video understanding.
 111 Therefore, we consider five different application scenarios: **1) Action Recognition in Special Scenarios (Action):** While previous benchmarks have extensively examined action recognition tasks, our focus
 112 here is to assess VFMs’ capabilities in recognizing actions within special scenarios. **2) AI for Science (Science):** Referencing previous work [24], we classify tasks related to medicine and natural sciences
 113 as a category. **3) Video Content Moderation (Safety):** We group tasks related to identifying harmful or
 114 misleading information in video content. **4) Video Quality Assessment (Quality):** We categorize more
 115 subjective tasks into this group. The goal is to assess VFMs’ ability to learn low-level information
 116 and human aesthetic preferences. **5) Emotion Analysis (Emotion):** We group tasks related to human
 117 emotion analysis into this category to evaluate VFMs’ ability to understand human emotions.
 118
 119

120 **Constructing the adaptation task based on the existing annotations.** Classification tasks are
 121 straightforward and well-defined, with strong classification performance often indicating robust
 122 feature learning. Therefore, they are suitable for evaluating video foundation models. We construct
 123 adaptation classification tasks based on the collected data and annotations as follow: **1) Remove Low-**
 124 **Quality Video Datasets:** We manually exclude datasets with videos that have low resolution (below

Table 3: **Task difficulty assessment based on visual language models.** For tasks with fewer categories, such as Fake Face (n=2) and Quality Assess (n=2), random guessing can lead to high accuracy, which may result in a lower apparent proportion of hard samples. Therefore, the zero-shot classification accuracy of the models should also be considered when making task selection.

ratio %	Dark Scene	Long Video	Medical Surgery	Animal Behavior	Fake Face	Harmfull Content	Quality Assess	Emotion Analysis
Easy	18.45	24.57	0.00	19.18	39.06	28.78	53.04	7.21
Spatial	19.00	20.44	4.17	20.86	20.72	24.56	51.24	5.01
Temporal	20.09	22.39	19.79	23.90	4.89	22.76	13.26	27.06
Hard	36.90	26.28	62.50	35.58	9.00	20.17	3.04	47.15

240p), low frame rate (below 15fps), insufficient quantity (fewer than 150 videos per category), or low annotation accuracy (below 90%). **2) Select Discriminative Tasks:** For task difficulty screening, we first evaluate zero-shot classification performance using CLIP-L [5], EVA-g [67], ViCLIP-L [20], and Internvideo2-1B [26]. We then classify samples as follows: *Easy*: Samples that are correctly classified by three or more models. *Spatial*: Samples that are correctly classified by both CLIP and EVA. *Temporal*: Samples that are correctly classified by at least one of ViCLIP or Internvideo2-1B, but not by CLIP and EVA. *Hard*: Samples that are incorrectly classified by all models. We use the zero-shot classification accuracy of the models and the aforementioned proportions as references for task selection. Based on this, we choose tasks with lower zero-shot classification accuracy, higher proportions of Hard and Temporal samples, and lower proportions of Easy samples. The proportions of each type of sample in the tasks we ultimately selected can be found in Table 3. **3) Control the Number of Categories:** For datasets that originally include category labels, such as ARID [57] and Animal Kingdom [60], we select categories with sufficient samples to ensure evaluation accuracy and stability. We also control the final number of categories to avoid making the adaptation task overly difficult. We observed that both zero-shot testing and few-shot experiments based on current VFMs show that when the number of categories is too high, models often perform no better than random guessing. Although this issue may be mitigated as VFMs improve, we currently need to control the number of categories to effectively showcase differences between models. We select the main categories for each task and limit the number of categories to around 10 (based on few-shot experiments). **4) Handling Multi-label and Regression Tasks:** For datasets that are not originally classification tasks, we transform the tasks into classification tasks. For example, for DOVER [63], which is used for video aesthetics and technical quality assessment (a regression task), we assume that videos with quality scores in the top 40% are "high-quality videos" and those with scores in the bottom 40% are "low-quality videos", thus converting the original task into a binary classification task. In total, we construct eight classification tasks to evaluate the adaptation capabilities of video foundation models.

Determining the evaluation protocol. Previous studies [25, 26, 34] typically train video models using entire samples of training set, and most popular benchmarks have large training sample sizes. We argue that this evaluation method overlooks the examination of the adaptation capability of VFMs. As illustrated in Figure 3, under the scenario of using full training samples, the differences between VFMs are difficult to discern. However, under a low-sample protocol, different foundation models exhibit varying degrees of task adaptation capabilities. We observe that for tasks such as Action Recognition in Dark Scenes, which VFMs usually excel at, there are significant differences in adaptation capabilities among different models when training samples are extremely limited (4 shot and 16 shot). As the number of samples gradually increases to 100 shot, these differences diminish. Conversely, for more challenging tasks like Emotion Analysis, the performances of different models are uniformly weak when training samples are extremely limited, showing no discernible differences until a certain number of training samples (100 shot) are reached, at which point different models begin to demonstrate distinct adaptation capabilities. Therefore, to account for the adaptation capabilities of models with different numbers of training samples, we define a task adaptation capability evaluation score (TA-score):

$$TA - score = \frac{Acc^{4s} + Acc^{16s} + Acc^{100s}}{3} \quad (1)$$

Where Acc^{4s} , Acc^{16s} , Acc^{100s} represent the model's top-1 accuracy for 4-shot, 16-shot, and 100-shot classifications, respectively. Unless otherwise specified, we will use TA-score to denote the performance of various tasks in VidTAB.

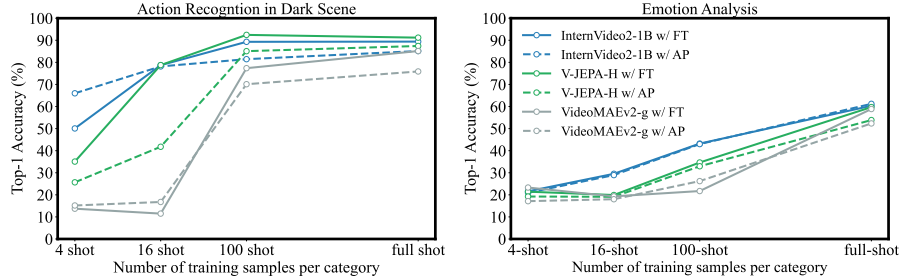


Figure 3: **Performance comparison on different training data scales.** We evaluate the performance variation of multiple video foundation models across tasks from two different domains as the scale of the training data changed. 'FT' and 'AP' denote full finetuning and attentive probe, respectively.

Table 4: **Comparison of adaptation method** All results are obtained using A100-80G with PyTorch-builtin mixed precision, using a batch size of 4 to measure Cuda memory and training time. "Dark" and "Emotion" denote the tasks of Action Recognition in Dark Scenes and Emotion Analysis, respectively. We show the result of V-JEPA-H [23] here,

Adaptation method	Tunable Params (M)	Cuda Memory (G)	Training Time (h)	Dark TA-score	Emotion TA-score
full finetuning	663.7	52.1	1.0	68.8	25.3
adapter	52.6	45.0	1.0	62.4	24.7
attentive probe	19.7	6.4	0.4	54.7	23.8
linear probe	0.0	6.0	0.3	12.9	16.2

169 **Identifying efficient adaptation method for evaluation.** We also need to identify how to adapt
 170 the foundation models to the corresponding task. Previous work [68, 69, 70, 71, 72] has explored
 171 various strategies for efficient adapting the foundation models. Here, we consider several of the most
 172 common and popular methods: **Full Finetuning:** Fine-tuning all the parameters of the pre-trained
 173 model. **Adapter:** Freezing the pre-trained model and inserting learnable low-rank adapter [73]
 174 modules into each block of the pre-trained model for adaptation. **Attentive Probe:** Freezing the
 175 pre-trained model and adding an additional learnable cross-attention block at the end of the model to
 176 achieve attentive pooling, followed by a linear projection for classification. **Linear Probe:** Directly
 177 using the features from the pre-trained model, performing mean pooling, and then using a linear
 178 projection for classification. We evaluate the performance of these adaptation methods based on the
 179 V-JEPA-H model, as shown in Table 4. Full finetuning and adapter exhibited the best adaptation
 180 performance, but incurred high training costs. Linear probe was highly efficient but showed weak
 181 adaptation performance. Attentive probe offered a good trade-off between efficiency and adaptation
 182 performance. Therefore, in subsequent evaluation experiments, we employed attentive probe to adapt
 183 various vision foundation models.

184 3.2 Video Embedding Benchmark

185 The main application domains of video embeddings we considering include: Label-Level: Classifica-
 186 tion and Action Retrieval. Instance-Level: Retrieval, Copy Detection and Ranking. For label-level
 187 tasks, VidTAB has already provided a flexible way to evaluate models. Therefore, VidEB aims to
 188 assess existing models at a finer semantic level, focusing on instance-level tasks. Although ranking
 189 tasks are common in recommendation system scenarios, they are influenced by user information and
 190 interactions, in addition to video data. Based on prior research [74], using frozen embeddings for
 191 video features does not consistently improve recommendation tasks (resulting in minimal or even
 192 negative effects). Thus, we have narrowed the final dataset scope to instance-level retrieval and copy
 193 detection. Apart from the traditional classification tasks, the evaluation of representations typically
 194 involves standard benchmarks such as video action retrieval [75, 76, 77], which primarily rely on
 195 class labels. However, this approach often overlooks the overall scene context and exhibits an overlap
 196 with recognition tasks. In contrast, inspired by previous works [78, 66, 79, 80, 81], we establish more
 197 rigorous criteria for embedding evaluation in Table 2. Specifically, we require the model to determine
 198 the priority and retrieve individual samples based on the overall similarity, rather than solely relying

199 on class labels. This evaluation protocol provides a more comprehensive assessment of the model’s
200 capability to encapsulate subtle visual information.

201 **Evaluation protocol.** To facilitate fine-grained embedding evaluation, we incorporate two tasks
202 for assessment: **(1) Hierarchical Video Retrieval** aims to retrieve videos from a database that
203 closely matches the query video in terms of scene, viewpoint, and temporal context. According
204 to previous work [65], videos related to the query are categorized into three levels based on their
205 similarity to the query: Duplicate Scene Videos (DSVs), Complementary Scene Videos (CSVs), and
206 Incident Scene Videos (ISVs), as shown in Table 2: Consequently, the retrieval tasks are structured
207 into three hierarchical levels: *Duplicate Scene Video Retrieval*: only DSVs are positive instances.
208 *Complementary Scene Video Retrieval*: both DSVs and CSVs are positive instances. *Incident Scene*
209 *Video Retrieval*: DSVs, CSVs, and ISVs are all positive instances. For the evaluation metric, we
210 follow [65] to utilize the mean Average Precision (mAP) to assess the quality of video ranking. **(2)**
211 **Video Copy Detection** aims to detect edited copies of the query video. Instead of the ranking/retrieval
212 task where all video pairs need to be sorted according to video embedding similarity, it is required
213 to identify a set of video pairs that contain edited versions of the given query. Following [66], we
214 consider the micro-AP (μ AP) as our evaluation metric that operates on all queries jointly and takes
215 the confidence scores into account.

216 4 Benchmarking Video Foundation Models

217 4.1 Targets and details of evaluation

218 **Evaluation targets** We evaluate twenty open-source vision foundation models. Including: (1)
219 twelve video foundation models, covering *different pre-training paradigms, model scales, and*
220 *training data scales*, to analyze the impact of these factors on the generalization capability of
221 foundation models. (2) five image foundation models to observe *how much generalization capability*
222 *trained on image data can exhibit in video understanding*. (3) three image-to-video methods based
223 on image foundation models to assess the *effectiveness of current efficient transfer methods*.

224 **Implementation details** All models take 8 frames (16 frames if the model has temporal downsam-
225 pling), with each frame being 224x224 in size as input. For VidTAB, to ensure fair comparison and
226 efficient assessment, we train all models for the same number of epochs and made minor adjustments
227 to the hyperparameters to ensure convergence. For VidEB, all models take 16 frames, with each frame
228 being 224x224 in size as input. In hierarchical video retrieval, the similarity of video-level embedding
229 determines the ranking of retrieval results. In video copy detection, each sample is segmented into 5
230 clips. The detection confidence score for the entire video is derived from the maximum frame-wise
231 similarity computed for each query-reference pair. See the Appendix for more details.

232 4.2 Results on VidTAB

233 **Zero-shot evaluation** To preliminarily assess the characteristics and difficulty of the dataset, we
234 first evaluate the zero-shot performance of the eight tasks we created using two image language
235 models and two video language models. As shown in the top section of Table 3, both image and
236 video models demonstrated some level of performance for action-related tasks, with video models
237 exhibiting relatively higher performance. For tasks involving low-level information understanding,
238 such as Quality Assessment task, image models performed significantly better. In contrast, for other
239 tasks involving scenarios typically unseen in training data, such as medical surgery videos or Safety
240 Review tasks requiring complex semantic reasoning, all models exhibited almost no performance.

241 **Main results** Table 5 presents the evaluation results on VidTAB. We summarize our findings as
242 follows. **On the whole, (1)** Despite exhibiting a degree of generalization capability, *current vision*
243 *FMs still struggle to adapt to unseen video tasks with limited training samples*. VFMs outperform
244 IFMs, particularly in tasks related to action and behavior understanding. However, IFMs exhibit
245 superior performance on more novel tasks, specifically in the domains of safety and quality, especially
246 when combined with image-to-video adaptation techniques. **(2)** The *adaptation performance of*
247 *models generally increases with the growth of data and model size*, as observed by the improvements
248 observed from V-JEPA-L to V-JEPA-H (+1.5) and ViCLIP-L-10M to ViCLIP-L-200M (+1.3).

Table 5: **Evaluating state-of-the-art FMs on the VidTAB.** The best and second-best results of foundation models are noted by **blue** and underline, respectively. 'I', 'V', and 'IV' denote image data, video data, and mixed image-video data, respectively. Data marked in gray indicates that the model uses a model trained on that data as initialization. 'K710ft' indicates that the model was fine-tuned with supervision using the labeled action recognition dataset Kinetics-710 (0.66M). Considering the random error in few-shot experiments, we conducted 3-fold experiments for both 4-shot and 16-shot settings, and used their mean as the final result. We also provide the results of full finetuning in the appendix.

	# Params. (M)	# Pt. Data	Average	Action		Science		Safety		Quality	Emotion
				Dark.	Long.	Medical.	Animal.	Harmful.	Fake.	Quality.	Emotion.
Random	-	-	22.7	9.1	10.0	6.3	8.3	33.3	50.0	50.0	14.3
<i>Zero-shot performance of visual language models</i>											
CLIP-L [5]	428	I-400M	35.7	29.2	34.6	12.5	32.9	42.1	56.3	65.5	12.9
EVA-CLIP-g [67]	1365	I-2B	36.0	32.8	37.2	9.4	28.5	39.6	52.8	69.5	17.9
ViCLIP-L [20]	428	I-400M+V-200M	33.6	26.2	37.5	8.3	29.3	32.1	52.2	53.9	29.0
InternVideo2 _{stage2} [26]	1350	IV-1.1M+IV-25.5M	40.6	37.1	40.2	11.5	45.2	59.1	51.3	56.1	24.3
<i>Image Foundation Model</i>											
CLIP-L [5]	316	I-400M	43.2	31.9	37.8	32.3	37.4	54.2	58.2	66.6	27.6
SigLiP-SO [82]	444	I-4.11B	43.3	27.6	38.4	36.5	35.8	53.3	58.5	67.8	28.5
EVA-g [83]	1035	I-2B	45.8	40.2	47.1	34.4	41.0	51.8	55.2	68.1	29.0
DINOv2-L [84]	317	I-142M	42.7	40.8	45.0	39.6	36.1	38.9	52.2	63.2	25.6
DINOv2-g [84]	1165	I-142M	44.4	37.8	46.4	42.7	36.0	48.5	53.2	64.3	26.3
<i>Image Foundation Model with image-to-video adaptation method</i>											
ST-Adapter-CLIP-L [70]	328	I-400M	46.5	42.4	44.3	31.2	40.1	47.4	64.6	<u>71.5</u>	30.4
AIM-CLIP-L [71]	328	I-400M	48.8	41.5	50.0	38.5	40.2	46.4	69.5	73.7	<u>30.6</u>
Zero2V-CLIP-L [72]	303	I-400M	46.3	40.3	47.0	31.2	40.2	46.1	<u>65.2</u>	69.9	30.5
<i>Video Foundation Model</i>											
ViCLIP-L-10M [20]	316	I-400M+V-10M	41.8	31.2	42.7	30.2	35.3	47.9	53.9	66.2	26.9
ViCLIP-L-200M [20]	316	I-400M+V-200M	43.3	38.2	44.6	30.2	37.9	47.4	54.9	65.9	27.5
VideoMAEv1-L [16]	316	V-0.24M	43.3	45.6	30.8	31.2	37.4	56.5	51.9	68.7	24.0
VideoMAEv1-H [16]	651	V-0.24M	44.7	45.5	31.0	35.4	38.6	55.8	51.8	70.5	29.1
VideoMAEv2-g [22]	1037	V-1.35M	37.8	35.2	18.3	18.8	33.7	<u>59.6</u>	50.9	64.7	21.6
VideoMAEv2-g ^{K710ft} [22]	1037	V-1.35M+K710ft	<u>54.0</u>	76.4	<u>72.6</u>	<u>50.0</u>	42.4	43.8	56.9	63.2	27.0
UMT-L _{stage1} [21]	316	V-0.66M	40.6	34.3	35.4	30.0	34.2	45.6	53.6	64.7	27.0
UMT-L _{stage2} [21]	316	V-0.66M+IV-25M	44.0	34.2	43.9	22.9	39.4	63.9	53.0	67.3	27.4
V-JEPA-L [23]	318	V-2M	43.5	50.4	34.3	39.6	39.7	43.9	51.7	66.7	21.4
V-JEPA-H [23]	653	V-2M	45.1	53.8	37.6	35.4	40.4	47.3	53.0	68.1	25.1
InternVideo2-1B _{stage1} [26]	1037	IV-1.1M	46.1	45.2	50.3	33.3	38.7	52.3	53.5	65.9	29.3
InternVideo2-1B _{stage1} [26]	1037	IV-1.1M+K710ft	56.7	<u>75.6</u>	77.5	53.1	<u>45.4</u>	47.2	55.5	66.2	33.2
InternVideo2-1B _{stage2} [26]	1037	IV-1.1M+IV-25.5M	53.6	66.0	71.1	38.5	50.0	53.6	54.7	64.3	30.3

249 **For the pre-training data, (3)** While augmenting video training data is generally beneficial, it
250 can negatively impact the performance on some tasks. For both VideoMAEv2-g and InternVideo2-
251 1B_{stage1}, fine-tuning on Kinetics-710 data significantly enhances Action-related tasks, but consis-
252 tently degrades certain Safety and Quality tasks. Similar findings are observed with ViCLIP-L,
253 where post-pretraining on a large-scale video dataset improves Action-related tasks but diminishes
254 performance in other domains (Science, Safety, Quality, Emotion). It could be attributed to the
255 limited diversity of the current video training data. **(4)** For models trained on single-modal visual data,
256 incorporating additional weak-supervised post-pretraining with visual-text data leads to significant
257 improvements in adaptation capabilities. This is evident in the performance gains observed from
258 UMT-L_{stage1} to UMT-L_{stage2} (+3.6) and from InternVideo2-1B_{stage1} to InternVideo2-1B_{stage2}
259 (+8.0). Interestingly, this finding contradicts previous conclusions drawn from commonly used action
260 recognition benchmarks, suggesting that these benchmarks may introduce bias. **For the pre-training**
261 **paradigms of model, (5)** The effectiveness of pre-training paradigms in scaling model size might not
262 be adequately validated on popular action recognition benchmarks. While VideoMAEv2 successfully
263 scaled a model to 1B parameters using the dual masking strategy [22], its adaptation performance
264 (37.7 vs 44.4) significantly declined compared to VideoMAEv1-H. Interestingly, VideoMAEv2-g
265 demonstrated remarkable performance after fine-tuning on Kinetics-710 (0.66M), suggesting that the
266 abundant labeled data may have compensated for the shortcomings of its pre-training performance.
267 **(6)** Single-modal self-supervised pre-training paradigms exhibit superior data efficiency compared to
268 multimodal weakly-supervised pre-training paradigms. Notably, V-JEPA and VideoMAEv1, trained
269 solely on relatively small-scale unlabeled video data via self-supervised pre-training, demonstrate
270 comparable or even superior performance to ViCLIP, which is trained on a massive dataset of video-

Table 6: **Evaluation of State-of-the-Art Foundation Models on the VidEB Dataset.** "K400pt" and "K400ft" denote that the model is pre-trained and fine-tuned, respectively, using the labeled action recognition dataset Kinetics-400 (0.31M). MCL: Multi-modal Contrastive Learning, SCL: Self-supervised Contrastive Learning, MVM: Masked Video Modeling, SFT: Supervised Fine-tuning. Other notations are consistent with those in Table 5.

	Pretrain Tasks	# Pretrain Data	Average	Scene			
				Duplicate	Complementary	Incident	Copyright
<i>Image Foundation Model</i>							
CLIP-L [5]	MCL	I-400M	43.0	41.1	46.4	52.0	32.3
EVA-g [83]	MCL	I-2B	37.1	41.4	46.1	51.7	9.3
SigLiP-SO [82]	MCL	I-4.11B	38.6	40.6	45.5	51.5	16.9
DINOv2-L [84]	SCL	I-142M	45.6	49.0	53.5	54.3	25.6
DINOv2-g [84]	SCL	I-142M	48.6	50.5	55.1	56.0	32.8
<i>Video Foundation Model</i>							
VideoMAEv1-L [16]	MVM	K400pt	12.9	14.5	15.1	13.2	8.8
VideoMAEv1-L-K400ft [16]	MVM+SFT	K400pt+ft	27.4	27.6	30.2	30.3	21.6
VideoMAEv2-g [22]	MVM	V-1.35M	11.6	14.8	15.4	13.4	2.8
VideoMAEv2-g-K710ft [22]	MVM+SFT	V-1.35M+K710ft	37.4	33.8	37.1	37.1	41.7
UMT-L _{stage1} [21]	MVM	V-0.66M	41.1	42.2	46.6	49.6	25.7
UMT-L _{stage1} -K710ft [21]	MVM+SFT	V-0.66M+K710ft	29.0	26.4	29.4	30.3	30.0
UMT-L _{stage2} [21]	MVM+MCL	V-0.66M+IV-25M	34.2	33.4	37.3	40.6	25.4
V-JEPA-L [23]	MVM	V-2M	19.7	21.3	23.9	21.7	12.0
V-JEPA-H [23]	MVM	V-2M	20.2	21.5	23.7	21.2	14.3
InternVideo2-1B _{stage1} [26]	MVM	IV-1.1M	50.4	47.3	52.1	54.9	47.3
InternVideo2-1B _{stage1} -K710ft [26]	MVM+SFT	IV-1.1M+K710ft	33.9	30.5	34.2	34.1	36.9
InternVideo2-1B _{stage2} [26]	MVM+MCL	IV-1.1M+IV-25.5M	34.6	32.4	36.8	39.9	29.3

271 text pairs. **In addition, (7) Effective adaptation method for FMs is crucial.** Three image-to-video
 272 methods based on CLIP-L achieved significant performance improvements compared to using an
 273 attentive probe directly. We believe this represents a promising avenue for future research.

274 4.3 Results on VidEB

275 The main results of VidEB are presented in Table 6. We evaluate the embedding performance using
 276 different pre-training paradigms for IFMs and VFMs as frozen feature extractors. Surprisingly, **IFMs**
 277 **performs better than most VFMs**, likely due to the existing gap in spatial modeling capabilities
 278 between VFMs and IFMs. **For the pre-training paradigms of the model, (1) The contrastive learn-**
 279 **ing (CL) based approach consistently excels in embedding evaluation.** Due to CL’s emphasis on the
 280 relationships between samples during training, DINOv2, which focuses solely on vision, outperforms
 281 vision-language contrastive methods like CLIP across multiple tasks. **(2) The effectiveness of masked**
 282 **video modeling is closely tied to the targets it reconstructs or aligns with.** With higher semantic
 283 richness, it shows progressive improvements in embedding quality for VideoMAE-L, V-JEPA-L, and
 284 UMT-L_{stage1}. **(3) Vision-centric pretraining outperforms Multi-modal pretraining in vision-centric**
 285 **scenarios.** Comparing UMT-L_{stage1} and InternVideo2-1B_{stage1} with their multi-modal counterparts
 286 UMT-L_{stage2} and InternVideo2-1B_{stage2}, the introduction of visual-text pair data in multi-stage
 287 training does not enhance performance in vision-centric scenarios. This is also consistent with the
 288 performance differences observed between DINO and CLIP-style pre-training methods. Additionally,
 289 we assess the **impact of fine-tuning on the embedding evaluation of these pre-trained models. (4)**
 290 **Labels bring new semantic information or disrupt existing finer-grained semantic information.** The
 291 performance variations after fine-tuning differ based on the pre-training strategy. For UMT-L_{stage1}
 292 and InternVideo2-1B_{stage1}, fine-tuning leads to a significant drop in performance (-12.1 for UMT
 293 and -16.5 for InternVideo) due to the introduction of more singular label information, which causes
 294 catastrophic forgetting. In contrast, VideoMAE and VideoMAEv2 show substantial performance
 295 gains (+14.5 and +25.8, respectively) because the low-level semantics learned during pre-training are
 296 less abstract and benefit more from the addition of high-level label information.

297 5 Conclusions

298 We present VideoEval, a comprehensive benchmark suite for efficiently evaluating the VFMs. To this
 299 end, we establish VidTAB, which explores suitable evaluation tasks and protocols for VFMs from
 300 the perspective of assessing their adaptability to unknown tasks with limited samples. Additionally,
 301 we create VidEB to evaluate the capability of VFMs’ feature embedding in directly supporting
 302 downstream tasks. Utilizing VideoEval, we conduct a large-scale study involving 20 popular open-
 303 source vision foundation models, providing valuable insights for future research directions.

304 References

- 305 [1] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of
306 deep bidirectional transformers for language understanding. *ArXiv*, abs/1810.04805, 2018.
- 307 [2] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
308 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
309 few-shot learners. In *NeurIPS*, 2020.
- 310 [3] OpenAI. Gpt-4 technical report. *ArXiv*, abs/2303.08774, 2023.
- 311 [4] OpenAI. Gpt-4v(ision) system card. [https://api.semanticscholar.org/CorpusID:
312 263218031](https://api.semanticscholar.org/CorpusID:263218031), 2023.
- 313 [5] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
314 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
315 models from natural language supervision. In *ICML*, 2021.
- 316 [6] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer.
317 High-resolution image synthesis with latent diffusion models, 2021.
- 318 [7] Christoph Feichtenhofer, Haoqi Fan, Jitendra Malik, and Kaiming He. Slowfast networks for
319 video recognition. In *ICCV*, 2019.
- 320 [8] Gedas Bertasius, Heng Wang, and Lorenzo Torresani. Is space-time attention all you need for
321 video understanding? In *ICML*, 2021.
- 322 [9] Ze Liu, Jia Ning, Yue Cao, Yixuan Wei, Zheng Zhang, Stephen Lin, and Han Hu. Video swin
323 transformer. In *CVPR*, 2022.
- 324 [10] Haoqi Fan, Bo Xiong, Kartikeya Mangalam, Yanghao Li, Zhicheng Yan, Jitendra Malik, and
325 Christoph Feichtenhofer. Multiscale vision transformers. In *ICCV*, 2021.
- 326 [11] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale
327 hierarchical image database. In *CVPR*, 2009.
- 328 [12] Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijaya-
329 narasimhan, Fabio Viola, Tim Green, Trevor Back, Paul Natsev, et al. The kinetics human
330 action video dataset. *ArXiv*, abs/1705.06950, 2017.
- 331 [13] Tian Pan, Yibing Song, Tianyu Yang, Wenhao Jiang, and Wei Liu. Videomoco: Contrastive
332 video representation learning with temporally adversarial examples. In *CVPR*, 2021.
- 333 [14] Chen Wei, Haoqi Fan, Saining Xie, Chao-Yuan Wu, Alan Yuille, and Christoph Feichtenhofer.
334 Masked feature prediction for self-supervised visual pre-training. In *CVPR*, 2022.
- 335 [15] Christoph Feichtenhofer, Haoqi Fan, Yanghao Li, and Kaiming He. Masked autoencoders as
336 spatiotemporal learners. *NeurIPS*, 2022.
- 337 [16] Zhan Tong, Yibing Song, Jue Wang, and Limin Wang. VideoMAE: Masked autoencoders are
338 data-efficient learners for self-supervised video pre-training. In *NeurIPS*, 2022.
- 339 [17] Rui Wang, Dongdong Chen, Zuxuan Wu, Yinpeng Chen, Xiyang Dai, Mengchen Liu, Yu-Gang
340 Jiang, Luwei Zhou, and Lu Yuan. Bevt: Bert pretraining of video transformers. *CVPR*, 2022.
- 341 [18] Hu Xu, Gargi Ghosh, Po-Yao Huang, Dmytro Okhonko, Armen Aghajanyan, Florian Metze,
342 Luke Zettlemoyer, and Christoph Feichtenhofer. Videoclip: Contrastive pre-training for zero-
343 shot video-text understanding. In *EMNLP*, 2021.
- 344 [19] Shen Yan, Tao Zhu, Zirui Wang, Yuan Cao, Mi Zhang, Soham Ghosh, Yonghui Wu, and
345 Jiahui Yu. Video-text modeling with zero-shot transfer from contrastive captioners. *ArXiv*,
346 abs/2212.04979, 2022.
- 347 [20] Yi Wang, Yinan He, Yizhuo Li, Kunchang Li, Jiashuo Yu, Xin Jian Ma, Xinyuan Chen, Yaohui
348 Wang, Ping Luo, Ziwei Liu, Yali Wang, Limin Wang, and Y. Qiao. Internvid: A large-scale
349 video-text dataset for multimodal understanding and generation. *ICLR*, 2024.

- 350 [21] Kunchang Li, Yali Wang, Yizhuo Li, Yi Wang, Yinan He, Limin Wang, and Yu Qiao. Unmasked
351 teacher: Towards training-efficient video foundation models. In *ICCV*, 2023.
- 352 [22] Limin Wang, Bingkun Huang, Zhiyu Zhao, Zhan Tong, Yinan He, Yi Wang, Yali Wang, and
353 Yu Qiao. Videomae v2: Scaling video masked autoencoders with dual masking. In *CVPR*, 2023.
- 354 [23] Adrien Bardes, Quentin Garrido, Jean Ponce, Xinlei Chen, Michael Rabbat, Yann LeCun, Mido
355 Assran, and Nicolas Ballas. V-jepa: Latent video prediction for visual representation learning.
356 2023.
- 357 [24] Long Zhao, Nitesh B Gundavarapu, Liangzhe Yuan, Hao Zhou, Shen Yan, Jennifer J Sun, Luke
358 Friedman, Rui Qian, Tobias Weyand, Yue Zhao, et al. Videoprism: A foundational visual
359 encoder for video understanding. *ArXiv*, abs/2402.13217, 2024.
- 360 [25] Yi Wang, Kunchang Li, Yizhuo Li, Yinan He, Bingkun Huang, Zhiyu Zhao, Hongjie Zhang,
361 Jilan Xu, Yi Liu, Zun Wang, Sen Xing, Guo Chen, Junting Pan, Jiashuo Yu, Yali Wang,
362 Limin Wang, and Yu Qiao. Internvideo: General video foundation models via generative and
363 discriminative learning. *ArXiv*, abs/2212.03191, 2022.
- 364 [26] Yi Wang, Kunchang Li, Xinhao Li, Jiashuo Yu, Yinan He, Guo Chen, Baoqi Pei, Rongkun
365 Zheng, Jilan Xu, Zun Wang, et al. Internvideo2: Scaling video foundation models for multimodal
366 video understanding. *ArXiv*, abs/2403.15377, 2024.
- 367 [27] Xiaohua Zhai, Joan Puigcerver, Alexander Kolesnikov, Pierre Ruysen, Carlos Riquelme, Mario
368 Lucic, Josip Djolonga, André Susano Pinto, Maxim Neumann, Alexey Dosovitskiy, Lucas
369 Beyer, Olivier Bachem, Michael Tschannen, Marcin Michalski, Olivier Bousquet, Sylvain Gelly,
370 and Neil Houlsby. The visual task adaptation benchmark. *Arxiv*, abs/1910.04867, 2019.
- 371 [28] Dan Hendrycks, Kevin Zhao, Steven Basart, Jacob Steinhardt, and Dawn Song. Natural
372 adversarial examples. In *CVPR*, 2021.
- 373 [29] Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishaal Shankar. Do imagenet
374 classifiers generalize to imagenet? In *ICML*, 2019.
- 375 [30] Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo,
376 Rahul Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, Dawn Song, Jacob Steinhardt, and Justin
377 Gilmer. The many faces of robustness: A critical analysis of out-of-distribution generalization.
378 In *ICCV*, 2021.
- 379 [31] Haohan Wang, Songwei Ge, Zachary C. Lipton, and Eric P. Xing. Learning robust global
380 representations by penalizing local predictive power. In *NeurIPS*, 2019.
- 381 [32] Badr Youbi Idrissi, Diane Bouchacourt, Randall Balestriero, Ivan Evtimov, Caner Hazir-
382 bas, Nicolas Ballas, Pascal Vincent, Michal Drozdal, David Lopez-Paz, and Mark Ibrahim.
383 Imagenet-x: Understanding model mistakes with factor of variation annotations. In *ICLR*, 2023.
- 384 [33] Micah Goldblum, Hossein Souri, Renkun Ni, Manli Shu, Viraj Prabhu, Gowthami Somepalli,
385 Prithvijit Chattopadhyay, Mark Ibrahim, Adrien Bardes, Judy Hoffman, Rama Chellappa,
386 Andrew Gordon Wilson, and Tom Goldstein. Battle of the backbones: A large-scale comparison
387 of pretrained models across computer vision tasks. In *NeurIPS*, 2023.
- 388 [34] Liangzhe Yuan, Nitesh Bharadwaj Gundavarapu, Long Zhao, Hao Zhou, Yin Cui, Lu Jiang,
389 Xuan Yang, Menglin Jia, Tobias Weyand, Luke Friedman, et al. Videoglue: Video general
390 understanding evaluation of foundation models. *ArXiv*, abs/2307.03166, 2023.
- 391 [35] Raghav Goyal, Samira Ebrahimi Kahou, Vincent Michalski, Joanna Materzynska, Susanne
392 Westphal, Heuna Kim, Valentin Haenel, Ingo Fründ, Peter Yianilos, Moritz Mueller-Freitag,
393 Florian Hoppe, Christian Thureau, Ingo Bax, and Roland Memisevic. The “something something”
394 video database for learning and evaluating visual common sense. In *ICCV*, 2017.
- 395 [36] Chunhui Gu, Chen Sun, Sudheendra Vijayanarasimhan, Caroline Pantofaru, David A. Ross,
396 George Toderici, Yeqing Li, Susanna Ricco, Rahul Sukthankar, Cordelia Schmid, and Jitendra
397 Malik. Ava: A video dataset of spatio-temporally localized atomic visual actions. *CVPR*, 2017.

- 398 [37] Will Kay, João Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijaya-
399 narasimhan, Fabio Viola, Tim Green, Trevor Back, Apostol Natsev, Mustafa Suleyman, and
400 Andrew Zisserman. The kinetics human action video dataset. *ArXiv*, abs/1705.06950, 2017.
- 401 [38] Raghav Goyal, Samira Ebrahimi Kahou, Vincent Michalski, Joanna Materzynska, Susanne
402 Westphal, Heuna Kim, Valentin Haenel, Ingo Fruend, Peter Yianilos, Moritz Mueller-Freitag,
403 et al. The " something something " video database for learning and evaluating visual common
404 sense. In *ICCV*, 2017.
- 405 [39] Mathew Monfort, Bolei Zhou, Sarah Adel Bargal, Alex Andonian, Tom Yan, Kandan Ramakr-
406 ishnan, Lisa M. Brown, Quanfu Fan, Dan Gutfreund, Carl Vondrick, and Aude Oliva. Moments
407 in time dataset: One million videos for event understanding. *TPAMI*, 2020.
- 408 [40] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. UCF101: A dataset of 101 human
409 actions classes from videos in the wild. *Arxiv*, abs/1212.0402, 2012.
- 410 [41] Fida Mohammad Thoker, Hazel Doughty, Piyush Bagad, and Cees GM Snoek. How severe is
411 benchmark-sensitivity in video self-supervised learning? In *ECCV*, 2022.
- 412 [42] Andong Deng, Taojiannan Yang, and Chen Chen. A large-scale study of spatiotemporal
413 representation learning with a new benchmark on action recognition. In *ICCV*, 2023.
- 414 [43] Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A
415 cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *ACL*, 2018.
- 416 [44] Soravit Changpinyo, Piyush Sharma, Nan Ding, and Radu Soricut. Conceptual 12m: Pushing
417 web-scale image-text pre-training to recognize long-tail visual concepts. In *CVPR*, 2021.
- 418 [45] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman,
419 Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick
420 Schramowski, Srivatsa Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmarczyk,
421 and Jenia Jitsev. LAION-5B: an open large-scale dataset for training next generation image-text
422 models. In *NeurIPS*, 2022.
- 423 [46] Max Bain, Arsha Nagrani, Gül Varol, and Andrew Zisserman. Frozen in time: A joint video
424 and image encoder for end-to-end retrieval. In *ICCV*, 2021.
- 425 [47] Sihan Chen, Xingjian He, Longteng Guo, Xinxin Zhu, Weining Wang, Jinhui Tang, and Jing
426 Liu. VALOR: vision-audio-language omni-perception pretraining model and dataset. *Arxiv*,
427 abs/2304.08345, 2023.
- 428 [48] Sihan Chen, Handong Li, Qunbo Wang, Zijia Zhao, Mingzhen Sun, Xinxin Zhu, and Jing Liu.
429 Vast: A vision-audio-subtitle-text omni-modality foundation model and dataset. *NeurIPS*, 2024.
- 430 [49] Tsai-Shien Chen, Aliaksandr Siarohin, Willi Menapace, Ekaterina Deyneka, Hsiang-wei Chao,
431 Byung Eun Jeon, Yuwei Fang, Hsin-Ying Lee, Jian Ren, Ming-Hsuan Yang, and Sergey
432 Tulyakov. Panda-70m: Captioning 70m videos with multiple cross-modality teachers. *Arxiv*,
433 abs/2402.19479, 2024.
- 434 [50] Rui Wang, Dongdong Chen, Zuxuan Wu, Yinpeng Chen, Xiyang Dai, Mengchen Liu, Lu Yuan,
435 and Yu-Gang Jiang. Masked video distillation: Rethinking masked feature modeling for
436 self-supervised video representation learning. In *CVPR*, 2023.
- 437 [51] Rohit Girdhar, Alaaeldin El-Nouby, Mannat Singh, Kalyan Vasudev Alwala, Armand Joulin,
438 and Ishan Misra. Omnimaes: Single model masked pretraining on images and videos. In *CVPR*,
439 2023.
- 440 [52] Chaitanya Ryali, Yuan-Ting Hu, Daniel Bolya, Chen Wei, Haoqi Fan, Po-Yao Huang, Vaibhav
441 Aggarwal, Arkabandhu Chowdhury, Omid Poursaeed, Judy Hoffman, Jitendra Malik, Yanghao
442 Li, and Christoph Feichtenhofer. Hiera: A hierarchical vision transformer without the bells-and-
443 whistles. In *ICML*, 2023.
- 444 [53] Jinpeng Wang, Yixiao Ge, Rui Yan, Yuying Ge, Kevin Qinghong Lin, Satoshi Tsutsui, Xudong
445 Lin, Guanyu Cai, Jianping Wu, Ying Shan, Xiaohu Qie, and Mike Zheng Shou. All in one:
446 Exploring unified video-language pre-training. In *CVPR*, 2023.

- 447 [54] Feng Cheng, Xizi Wang, Jie Lei, David J. Crandall, Mohit Bansal, and Gedas Bertasius. Vindlu:
448 A recipe for effective video-and-language pretraining. *ArXiv*, abs/2212.05051, 2022.
- 449 [55] Fida Mohammad Thoker, Hazel Doughty, Piyush Bagad, and Cees G. M. Snoek. How severe is
450 benchmark-sensitivity in video self-supervised learning? In *ECCV*, 2022.
- 451 [56] Madeline Chantry Schiappa, Naman Biyani, Prudvi Kamtam, Shruti Vyas, Hamid Palangi, Vib-
452 hav Vineet, and Yogesh S Rawat. A large-scale robustness analysis of video action recognition
453 models. In *CVPR*, 2023.
- 454 [57] Yuecong Xu, Jianfei Yang, Haozhi Cao, Kezhi Mao, Jianxiong Yin, and Simon See. Arid: A
455 new dataset for recognizing action in the dark. In *IJCAI*, 2021.
- 456 [58] Hilde Kuehne, Ali Arslan, and Thomas Serre. The language of actions: Recovering the syntax
457 and semantics of goal-directed human activities. In *CVPR*, 2014.
- 458 [59] Klaus Schoeffmann, Heinrich Husslein, Sabrina Kletz, Stefan Petschornig, Bernd Muenzer,
459 and Christian Beecks. Video retrieval in laparoscopic video recordings with dynamic content
460 descriptors. *Multimedia Tools and Applications*, 77:16813–16832, 2018.
- 461 [60] Xun Long Ng, Kian Eng Ong, Qichen Zheng, Yun Ni, Si Yong Yeo, and Jun Liu. Animal
462 kingdom: A large and diverse dataset for animal behavior understanding. In *CVPR*, 2022.
- 463 [61] Andreas Rossler, Davide Cozzolino, Luisa Verdoliva, Christian Riess, Justus Thies, and Matthias
464 Nießner. Faceforensics++: Learning to detect manipulated facial images. In *ICCV*, 2019.
- 465 [62] Syed Hammad Ahmed, Muhammad Junaid Khan, HM Qaisar, and Gita Sukthankar. Malicious
466 or benign? towards effective content moderation for children’s videos. *ArXiv*, abs/2305.15551,
467 2023.
- 468 [63] Haoning Wu, Erli Zhang, Liang Liao, Chaofeng Chen, Jingwen Hou, Annan Wang, Wenxiu
469 Sun, Qiong Yan, and Weisi Lin. Exploring video quality assessment on user generated contents
470 from aesthetic and technical perspectives. In *ICCV*, 2023.
- 471 [64] Jiyoung Lee, Seungryong Kim, Sunok Kim, Jungin Park, and Kwanghoon Sohn. Context-aware
472 emotion recognition networks. In *ICCV*, 2019.
- 473 [65] Giorgos Kordopatis-Zilos, Symeon Papadopoulos, Ioannis Patras, and Ioannis Kompatsiaris.
474 Fivr: Fine-grained incident video retrieval. *IEEE Transactions on Multimedia*, 21, 2019.
- 475 [66] Ed Pizzi, Giorgos Kordopatis-Zilos, Hiral Patel, Gheorghe Postelnicu, Sugosh Nagavara Ravin-
476 dra, Akshay Gupta, Symeon Papadopoulos, Giorgos Toliass, and Matthijs Douze. The 2023
477 video similarity dataset and challenge. *Computer Vision and Image Understanding*, 2024.
- 478 [67] Quan Sun, Yuxin Fang, Ledell Yu Wu, Xinlong Wang, and Yue Cao. Eva-clip: Improved
479 training techniques for clip at scale. *ArXiv*, abs/2303.15389, 2023.
- 480 [68] Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin de Laroussilhe,
481 Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning
482 for NLP. In *ICML*, 2019.
- 483 [69] Bruce XB Yu, Jianlong Chang, Haixin Wang, Lingbo Liu, Shijie Wang, Zhiyu Wang, Junfan
484 Lin, Lingxi Xie, Haojie Li, Zhouchen Lin, et al. Visual tuning. *ACM Computing Surveys*, 2023.
- 485 [70] Junting Pan, Ziyi Lin, Xiatian Zhu, Jing Shao, and Hongsheng Li. Parameter-efficient image-to-
486 video transfer learning. *arXiv*, abs/2206.13559, 2022.
- 487 [71] Taojiannan Yang, Yi Zhu, Yusheng Xie, Aston Zhang, Chen Chen, and Mu Li. Aim: Adapting
488 image models for efficient video action recognition. In *ICLR*, 2023.
- 489 [72] Xinhao Li and Limin Wang. Zeroi2v: Zero-cost adaptation of pre-trained transformers from
490 image to video. *ArXiv*, abs/2310.01324, 2023.

- 491 [73] Jonas Pfeiffer, Andreas Rücklé, Clifton Poth, Aishwarya Kamath, Ivan Vulić, Sebastian Ruder,
492 Kyunghyun Cho, and Iryna Gurevych. Adapterhub: A framework for adapting transformers. In
493 *EMNLP*, 2020.
- 494 [74] Yongxin Ni, Yu Cheng, Xiangyan Liu, Junchen Fu, Youhua Li, Xiangnan He, Yongfeng Zhang,
495 and Fajie Yuan. A content-driven micro-video recommendation dataset at scale. *arXiv preprint*
496 *arXiv:2309.15379*, 2023.
- 497 [75] Tengda Han, Weidi Xie, and Andrew Zisserman. Self-supervised co-training for video represen-
498 tation learning. In *NeurIPS*, 2020.
- 499 [76] Dejing Xu, Jun Xiao, Zhou Zhao, Jian Shao, Di Xie, and Yueting Zhuang. Self-supervised
500 spatiotemporal learning via video clip order prediction. In *CVPR*, 2019.
- 501 [77] Tengda Han, Weidi Xie, and Andrew Zisserman. Memory-augmented dense predictive coding
502 for video representation learning. In *ECCV*, 2020.
- 503 [78] Bryan A Plummer, Liwei Wang, Chris M Cervantes, Juan C Caicedo, Julia Hockenmaier, and
504 Svetlana Lazebnik. Flickr30k entities: Collecting region-to-phrase correspondences for richer
505 image-to-sentence models. In *ICCV*, 2015.
- 506 [79] Xiao Wu, Alexander G Hauptmann, and Chong-Wah Ngo. Practical elimination of near-
507 duplicates from web video search. In *Proceedings of the 15th ACM international conference on*
508 *Multimedia*, pages 218–227, 2007.
- 509 [80] Yu-Gang Jiang, Yudong Jiang, and Jiajun Wang. Vcdb: a large-scale database for partial copy
510 detection in videos. In *ECCV*, 2014.
- 511 [81] Matthijs Douze, Giorgos Tolias, Ed Pizzi, Zoë Papakipos, Lowik Chanussot, Filip Radenovic,
512 Tomas Jenicek, Maxim Maximov, Laura Leal-Taixé, Ismail Elezi, et al. The 2021 image
513 similarity dataset and challenge. *arXiv preprint arXiv:2106.09672*, 2021.
- 514 [82] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for
515 language image pre-training. In *ICCV*, 2023.
- 516 [83] Yuxin Fang, Wen Wang, Binhui Xie, Quan Sun, Ledell Wu, Xinggang Wang, Tiejun Huang,
517 Xinlong Wang, and Yue Cao. Eva: Exploring the limits of masked visual representation learning
518 at scale. In *CVPR*, 2023.
- 519 [84] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov,
520 Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning
521 robust visual features without supervision. *ArXiv*, abs/2304.07193, 2023.

522 **NeurIPS Paper Checklist**

523 **1. Claims**

524 Question: Do the main claims made in the abstract and introduction accurately reflect the
525 paper's contributions and scope?

526 Answer: [Yes]

527 Justification: Our main claims made in the abstract and introduction accurately reflect the
528 paper's contributions and scope.

529 Guidelines:

- 530 • The answer NA means that the abstract and introduction do not include the claims
531 made in the paper.
- 532 • The abstract and/or introduction should clearly state the claims made, including the
533 contributions made in the paper and important assumptions and limitations. A No or
534 NA answer to this question will not be perceived well by the reviewers.
- 535 • The claims made should match theoretical and experimental results, and reflect how
536 much the results can be expected to generalize to other settings.
- 537 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
538 are not attained by the paper.

539 **2. Limitations**

540 Question: Does the paper discuss the limitations of the work performed by the authors?

541 Answer: [Yes]

542 Justification: We place it in Appendix Section D.

543 Guidelines:

- 544 • The answer NA means that the paper has no limitation while the answer No means that
545 the paper has limitations, but those are not discussed in the paper.
- 546 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 547 • The paper should point out any strong assumptions and how robust the results are to
548 violations of these assumptions (e.g., independence assumptions, noiseless settings,
549 model well-specification, asymptotic approximations only holding locally). The authors
550 should reflect on how these assumptions might be violated in practice and what the
551 implications would be.
- 552 • The authors should reflect on the scope of the claims made, e.g., if the approach was
553 only tested on a few datasets or with a few runs. In general, empirical results often
554 depend on implicit assumptions, which should be articulated.
- 555 • The authors should reflect on the factors that influence the performance of the approach.
556 For example, a facial recognition algorithm may perform poorly when image resolution
557 is low or images are taken in low lighting. Or a speech-to-text system might not be
558 used reliably to provide closed captions for online lectures because it fails to handle
559 technical jargon.
- 560 • The authors should discuss the computational efficiency of the proposed algorithms
561 and how they scale with dataset size.
- 562 • If applicable, the authors should discuss possible limitations of their approach to
563 address problems of privacy and fairness.
- 564 • While the authors might fear that complete honesty about limitations might be used by
565 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
566 limitations that aren't acknowledged in the paper. The authors should use their best
567 judgment and recognize that individual actions in favor of transparency play an impor-
568 tant role in developing norms that preserve the integrity of the community. Reviewers
569 will be specifically instructed to not penalize honesty concerning limitations.

570 **3. Theory assumptions and proofs**

571 Question: For each theoretical result, does the paper provide the full set of assumptions and
572 a complete (and correct) proof?

573 Answer: [NA]

574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627

Justification: Our paper don't have theory assumptions and proofs

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: See <https://github.com/MCG-NJU/VideoEval>.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678

Answer: [Yes]

Justification: See <https://github.com/MCG-NJU/VideoEval>.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: See Appendix and <https://github.com/MCG-NJU/VideoEval>.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: For VidTAB, we referred to the setting of the previous few shot work and repeated the few shot experiment three times to reduce randomness.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).

- 679 • It should be clear whether the error bar is the standard deviation or the standard error
680 of the mean.
- 681 • It is OK to report 1-sigma error bars, but one should state it. The authors should
682 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
683 of Normality of errors is not verified.
- 684 • For asymmetric distributions, the authors should be careful not to show in tables or
685 figures symmetric error bars that would yield results that are out of range (e.g. negative
686 error rates).
- 687 • If error bars are reported in tables or plots, The authors should explain in the text how
688 they were calculated and reference the corresponding figures or tables in the text.

689 8. Experiments compute resources

690 Question: For each experiment, does the paper provide sufficient information on the com-
691 puter resources (type of compute workers, memory, time of execution) needed to reproduce
692 the experiments?

693 Answer: [Yes]

694 Justification: See <https://github.com/MCG-NJU/VideoEval> for training details.

695 Guidelines:

- 696 • The answer NA means that the paper does not include experiments.
- 697 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
698 or cloud provider, including relevant memory and storage.
- 699 • The paper should provide the amount of compute required for each of the individual
700 experimental runs as well as estimate the total compute.
- 701 • The paper should disclose whether the full research project required more compute
702 than the experiments reported in the paper (e.g., preliminary or failed experiments that
703 didn't make it into the paper).

704 9. Code of ethics

705 Question: Does the research conducted in the paper conform, in every respect, with the
706 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

707 Answer: [Yes]

708 Justification: Our paper follows the NeurIPS Code of Ethics in every respect.

709 Guidelines:

- 710 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 711 • If the authors answer No, they should explain the special circumstances that require a
712 deviation from the Code of Ethics.
- 713 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
714 eration due to laws or regulations in their jurisdiction).

715 10. Broader impacts

716 Question: Does the paper discuss both potential positive societal impacts and negative
717 societal impacts of the work performed?

718 Answer: [Yes]

719 Justification: See Appendix Section D.

720 Guidelines:

- 721 • The answer NA means that there is no societal impact of the work performed.
- 722 • If the authors answer NA or No, they should explain why their work has no societal
723 impact or why the paper does not address societal impact.
- 724 • Examples of negative societal impacts include potential malicious or unintended uses
725 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
726 (e.g., deployment of technologies that could make decisions that unfairly impact specific
727 groups), privacy considerations, and security considerations.

- 728
- 729
- 730
- 731
- 732
- 733
- 734
- 735
- 736
- 737
- 738
- 739
- 740
- 741
- 742
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
 - The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
 - If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

743 **11. Safeguards**

744 Question: Does the paper describe safeguards that have been put in place for responsible
745 release of data or models that have a high risk for misuse (e.g., pretrained language models,
746 image generators, or scraped datasets)?

747 Answer: [NA]

748 Justification: We don't provided new model.

749 Guidelines:

- 750
- 751
- 752
- 753
- 754
- 755
- 756
- 757
- 758
- 759
- The answer NA means that the paper poses no such risks.
 - Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
 - Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
 - We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

760 **12. Licenses for existing assets**

761 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
762 the paper, properly credited and are the license and terms of use explicitly mentioned and
763 properly respected?

764 Answer: [Yes]

765 Justification: See Appendix.

766 Guidelines:

- 767
- 768
- 769
- 770
- 771
- 772
- 773
- 774
- 775
- 776
- 777
- 778
- 779
- The answer NA means that the paper does not use existing assets.
 - The authors should cite the original paper that produced the code package or dataset.
 - The authors should state which version of the asset is used and, if possible, include a URL.
 - The name of the license (e.g., CC-BY 4.0) should be included for each asset.
 - For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
 - If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
 - For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.

780 • If this information is not available online, the authors are encouraged to reach out to
781 the asset’s creators.

782 13. New assets

783 Question: Are new assets introduced in the paper well documented and is the documentation
784 provided alongside the assets?

785 Answer: [Yes]

786 Justification: See <https://github.com/MCG-NJU/VideoEval>, we provide our annota-
787 tions here.

788 Guidelines:

- 789 • The answer NA means that the paper does not release new assets.
- 790 • Researchers should communicate the details of the dataset/code/model as part of their
791 submissions via structured templates. This includes details about training, license,
792 limitations, etc.
- 793 • The paper should discuss whether and how consent was obtained from people whose
794 asset is used.
- 795 • At submission time, remember to anonymize your assets (if applicable). You can either
796 create an anonymized URL or include an anonymized zip file.

797 14. Crowdsourcing and research with human subjects

798 Question: For crowdsourcing experiments and research with human subjects, does the paper
799 include the full text of instructions given to participants and screenshots, if applicable, as
800 well as details about compensation (if any)?

801 Answer: [NA]

802 Justification: Our paper does not involve crowdsourcing nor research with human subjects

803 Guidelines:

- 804 • The answer NA means that the paper does not involve crowdsourcing nor research with
805 human subjects.
- 806 • Including this information in the supplemental material is fine, but if the main contribu-
807 tion of the paper involves human subjects, then as much detail as possible should be
808 included in the main paper.
- 809 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
810 or other labor should be paid at least the minimum wage in the country of the data
811 collector.

812 15. Institutional review board (IRB) approvals or equivalent for research with human 813 subjects

814 Question: Does the paper describe potential risks incurred by study participants, whether
815 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
816 approvals (or an equivalent approval/review based on the requirements of your country or
817 institution) were obtained?

818 Answer: [NA]

819 Justification: Our paper does not involve crowdsourcing nor research with human subjects.

820 Guidelines:

- 821 • The answer NA means that the paper does not involve crowdsourcing nor research with
822 human subjects.
- 823 • Depending on the country in which research is conducted, IRB approval (or equivalent)
824 may be required for any human subjects research. If you obtained IRB approval, you
825 should clearly state this in the paper.
- 826 • We recognize that the procedures for this may vary significantly between institutions
827 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
828 guidelines for their institution.
- 829 • For initial submissions, do not include any information that would break anonymity (if
830 applicable), such as the institution conducting the review.

831
832
833
834
835
836
837
838
839
840
841
842

16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [NA]

Justification: We only use LLM for simple work like writing.

Guidelines:

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>) for what should or should not be described.