

Eliciting In-Context Learning in Vision-Language Models for Videos Through Curated Data Distributional Properties

Anonymous ACL submission

Abstract

A major reason behind the recent success of large language models (LLMs) is their *in-context learning* capability, which makes it possible to rapidly adapt them to downstream text-based tasks by prompting them with a small number of relevant demonstrations. While large vision-language models (VLMs) have recently been developed for tasks requiring both text and images, they largely lack in-context learning over visual information, especially in understanding and generating text about videos. In this work, we implement **Emergent In-context Learning on Videos (EILeV)**, a novel training paradigm that induces in-context learning over video and text by capturing key properties of pre-training data found by prior work to be essential for in-context learning in transformers. In our experiments, we show that **EILeV**-trained models outperform other off-the-shelf VLMs in few-shot video narration for novel, rare actions. Furthermore, we demonstrate that these key properties of bursty distributions, skewed marginal distributions, and dynamic meaning each contribute to varying degrees to VLMs' in-context learning capability in narrating procedural videos. Our results, analysis, and **EILeV**-trained models yield numerous insights about the emergence of in-context learning over video and text, creating a foundation for future work to optimize and scale VLMs for open-domain video understanding and reasoning.

1 Introduction

In recent years, the advent of transformer-based (Vaswani et al., 2017) large language models (LLMs) has garnered significant attention in and beyond the AI research community. A central reason for this is their *in-context learning* capability (Brown et al., 2020), which makes it possible to rapidly adapt LLMs to novel tasks by simply prompting them with a few demonstrations. This capability removes the need for the expensive and

arduous task-specific fine-tuning required by earlier language modeling approaches.

While in-context learning has been extensively studied and utilized in purely text-based problems in language understanding, reasoning, and generation, there are myriad potential applications for this rapid post-deployment adaptation in processing *video*. For example, in embodied and task-oriented AI, a major challenge is to recognize novel, rare human actions from video that cannot possibly be completely covered in training data (Perrett et al., 2023; Du et al., 2023; Bao et al., 2023). A vision-language model (VLM) capable of in-context learning over video could address this challenge, as it would only require a few related videos of actions as few-shot, in-context examples to recognize and reason about these novel, rare actions. However, while large VLMs for jointly processing text and images have been developed (Li et al., 2022, 2023c; Dai et al., 2023; Zhu et al., 2023a; Peng et al., 2023; Liu et al., 2023), they are typically not optimized for reasoning over multiple images (i.e., frames), crucial for understanding videos. Meanwhile, a handful of open-source VLMs have recently been developed for video understanding (Zellers et al., 2022; Li et al., 2023b; Zhang et al., 2023; Lin et al., 2023), but they lack in-context learning.

In-context learning in text-only, transformer-based LLMs was initially observed to improve with increased model size, along with the size and diversity of training data (Brown et al., 2020). Later, Chan et al. (2022) identified several distributional properties of the training data as causes for this emergent behavior in transformer-based models: (1) bursty distributions with entities that tend to appear in clusters, (2) skewed marginal distributions with a long tail of infrequent items, and (3) dynamic meaning with label multiplicity. However, as their experiments relied on small transformer-based models trained on synthetic image classification data, it remains unclear whether their findings

084 hold true for VLMs trained on video and text at
085 scale.

086 In this work, we address this question by con-
087 ducting systematic empirical experiments to investi-
088 gate whether these training data distributional prop-
089 erties also elicit in-context learning capabilities in
090 VLMs for video. Specifically, we use various text
091 annotations from Ego4D (Grauman et al., 2022),
092 a popular video dataset, to implement **Emergent**
093 **In-context Learning on Videos (EILeV)**, a novel
094 VLM training method that satisfies all three prop-
095 erties and successfully elicits in-context learning
096 over video and text. In our experiments, we ob-
097 serve that the **EILeV**-trained models outperform
098 other off-the-shelf VLMs in few-shot video nar-
099 ration on rare actions, and that, through careful
100 ablation studies, each property indeed contributes
101 to this in-context learning capability. Furthermore,
102 our analysis yields a host of new insights around
103 the importance of each property in in-context learn-
104 ing for video.

105 The contributions of our work are as follows:
106 (1) we propose **EILeV**, a novel training method
107 that can elicit in-context learning capabilities in
108 VLMs for video and text, (2) we validate through
109 systematic ablation experiments that the same data
110 distributional properties that elicit in-context learn-
111 ing in small transformer-based models also apply
112 to VLMs for videos, and (3) we release a set of
113 **EILeV**-trained VLMs with in-context learning ca-
114 pabilities optimized for egocentric videos.

115 2 Related Work

116 2.1 In-Context Learning

117 Brown et al. (2020) discovered in-context learning
118 in LLMs when creating GPT-3. This was a sig-
119 nificant departure from fine-tuning which involves
120 parameter updates to adapt LLMs to downstream
121 tasks. Instead, in-context learning enables LLMs to
122 be adapted without parameter updates by prompt-
123 ing them with a few examples of a task as part
124 of the input context for text generation. The size
125 of the model and training data were thought to be
126 key to training a model with in-context learning
127 capabilities.

128 More recently, there has been more research on
129 the exact causes of in-context learning. Min et al.
130 (2022) proposed MetaICL, a meta-training frame-
131 work to elicit in-context learning capabilities in
132 text-only language models. MetaICL conditions
133 each example with related in-context examples dur-

134 ing training. Chan et al. (2022) investigated the dis-
135 tributional properties of training data for in-context
136 learning. Their findings showed that there are cer-
137 tain properties that encourage in-context learning in
138 transformer-based models, and massive textual data
139 from the web used to train LLMs naturally have
140 those properties. Furthermore, Reddy (2023) found
141 that in-context learning is driven by the abrupt
142 emergence of an induction head. There have also
143 been works with findings about in-context learning
144 in VLMs. Notably, training large generative VLMs
145 with image-text interleaved data has been shown
146 to be an effective technique to improve model per-
147 formance, especially in tasks involving in-context
148 learning (Alayrac et al., 2022; McKinzie et al.,
149 2024; Wang et al., 2024). Our work combines
150 these insights from prior work around the cause of
151 in-context learning to propose a new VLM training
152 paradigm for video and text, and carefully investi-
153 gates how they contribute to in-context learning.

154 2.2 Vision-Language Models (VLMs)

155 With the recent success of text-only LLMs, there
156 have been various efforts to replicate their success
157 in multimodal settings, especially vision and lan-
158 guage. Two different types of approaches in train-
159 ing generative VLMs have been proposed. The first
160 is to train them from scratch using large text and
161 paired image and text datasets (Hao et al., 2022;
162 Huang et al., 2024; Peng et al., 2023; Lu et al.,
163 2023). This approach allows the most controlla-
164 bility and flexibility as the resulting VLM is not
165 dependent on other pre-trained models that may
166 have undesirable behaviors, but it requires a mas-
167 sive amount of compute and data. In order to ad-
168 dress these challenges, a number of approaches
169 have been proposed to create VLMs by learning a
170 mapping from a frozen pre-trained vision encoder
171 to the input space of a frozen pre-trained LLM
172 (Alayrac et al., 2022; Li et al., 2023b; Zhao et al.,
173 2023; Li et al., 2022, 2023c; Dai et al., 2023; Liu
174 et al., 2023; Zhang et al., 2023; Lin et al., 2023;
175 Yang et al., 2022; Li et al., 2023d; Zhu et al., 2023a;
176 Laurençon et al., 2023; Maaz et al., 2023; Ye et al.,
177 2023; Gong et al., 2023; Zhang et al., 2024).

178 Some of these approaches enable the result-
179 ing VLMs to process videos by representing
180 them as sequences of still frames; however, only
181 Flamingo (Alayrac et al., 2022), Otter (Li et al.,
182 2023b) and Kosmos-2 (Peng et al., 2023) support
183 in-context learning over video and text as a by-
184 product of their large-scale pre-training. In this

| | | | |
|-----|---|---|-----|
| 185 | work, we conduct thorough investigation of how | and lastly discuss how we evaluate the in-context | 232 |
| 186 | key properties of training data achieve in-context | learning capability of VLMs trained on video and | 233 |
| 187 | learning beyond just as a by-product of large-scale | text. | 234 |
| 188 | training. | | |
| 189 | 3 Three Distributional Properties for | 4.1 Problem Definition | 235 |
| 190 | In-Context Learning | We target the task of <i>few-shot video narration</i> using | 236 |
| 191 | Since Brown et al. (2020) discovered in-context | the Ego4D dataset (Grauman et al., 2022). | 237 |
| 192 | learning in text-only LLMs, there has been much | | |
| 193 | research into the cause for in-context learning. In | Few-Shot Video Narration <i>Video narration</i> is a | 238 |
| 194 | particular, Chan et al. (2022) found that three char- | captioning task where given a video, a system must | 239 |
| 195 | acteristics of the training data are important in elic- | generate a text description of the events occurring | 240 |
| 196 | iting in-context learning in transformer-based mod- | in the video. Here, <i>few-shot video narration</i> refers | 241 |
| 197 | els, each of which is abundant in both natural lan- | to the implementation of this task where a VLM | 242 |
| 198 | guage and video data: <i>bursty distributions</i> , <i>skewed</i> | (pre-trained on large-scale video and text data) is | 243 |
| 199 | <i>marginal distributions</i> , and <i>dynamic meaning</i> . | conditioned with one or more example videos and | 244 |
| 200 | | narrations before being prompted to generate a nar- | 245 |
| 201 | Bursty Distributions In-context learning relies | ration for a held-out video clip. If conditioning | 246 |
| 202 | on data where entities appear in clusters, or non- | such a VLM on several example videos and narra- | 247 |
| 203 | uniformly depending on the context. Groups of re- | tions improves the quality of narration, this implies | 248 |
| 204 | lated entities may be mentioned frequently in some | that the VLM is indeed capable of in-context learn- | 249 |
| 205 | contexts, but much more rarely in other contexts. | ing over video and text. | 250 |
| 206 | | | |
| 207 | Skewed Marginal Distributions In-context | Ego4D Ego4D is a popular large-scale dataset of | 251 |
| 208 | learning also relies on data of skewed marginal | egocentric videos that have been densely annotated | 252 |
| 209 | distributions with a long tail of infrequent items | with human-written English narrations, ideal for | 253 |
| 210 | (i.e., a Zipfian distribution). This phenomenon is a | our task. Beyond narrations, the dataset includes | 254 |
| 211 | long-standing challenge in representing language | higher-level class labels for the verbs and nouns | 255 |
| 212 | and images, and has long been observed in text, | associated with each narrated video clip. These an- | 256 |
| 213 | image, and video datasets collected for research. | notations enable systematic ablations for all three | 257 |
| 214 | | distributional properties of training data discovered | 258 |
| 215 | Dynamic Meaning Lastly, in-context learning | by Chan et al. (2022) to facilitate in-context learn- | 259 |
| 216 | relies on dynamic meaning, where a single entity | ing, enabling a systematic study of in-context learn- | 260 |
| 217 | can have multiple possible interpretations, and mul- | ing over video and text in VLMs. These ablations | 261 |
| 218 | multiple entities can map to the same interpretation. | are introduced in Section 4.2. | 262 |
| 219 | In natural language, we observe this property in word | | |
| 220 | senses, homonyms, and synonyms. In the visual | 4.2 Training Paradigm & Ablations | 263 |
| 221 | world, a particular object may be described in multi- | Using Ego4D’s “Forecasting Hands & Objects | 264 |
| 222 | ple valid ways, e.g., synonyms, physical properties, | Master File”, we construct a dataset of interleaved | 265 |
| 223 | and hypernyms. Meanwhile, many distinct objects | text and video that satisfies these properties, and | 266 |
| 224 | may be grouped based on various descriptors. | use it to train and evaluate VLMs. We call this | 267 |
| 225 | | training procedure Emergent In-context Learning | 268 |
| 226 | 4 Problem & Methods | on Videos (EILeV) . EILeV uses the video and | 269 |
| 227 | In this section, we first introduce the target prob- | text data provided by Ego4D to implement all three | 270 |
| 228 | lem and dataset for our evaluations of in-context | distributional properties necessary for in-context | 271 |
| 229 | learning. Next, we introduce EILeV , our training | learning: bursty distributions, skewed marginal dis- | 272 |
| 230 | paradigm which captures all three distributional | tributions, and dynamic meaning. To demonstrate | 273 |
| 231 | properties thought to elicit in-context learning, as | the importance of each distributional property cap- | 274 |
| | well as the ablations we use to validate the im- | tured in EILeV , we use Ego4D’s detailed annota- | 275 |
| | portance of each property in enabling in-context | tions to carefully ablate each property as illustrat- | 276 |
| | learning over video and text. We then introduce | ed in Figure 1. | 277 |
| | the model architecture we apply this paradigm to, | For all experiments, each training data point con- | 278 |
| | | sists of a <i>context</i> with 16 video-narration pairs, | 279 |
| | | and a <i>query</i> with a single video-narration pair. We | 280 |

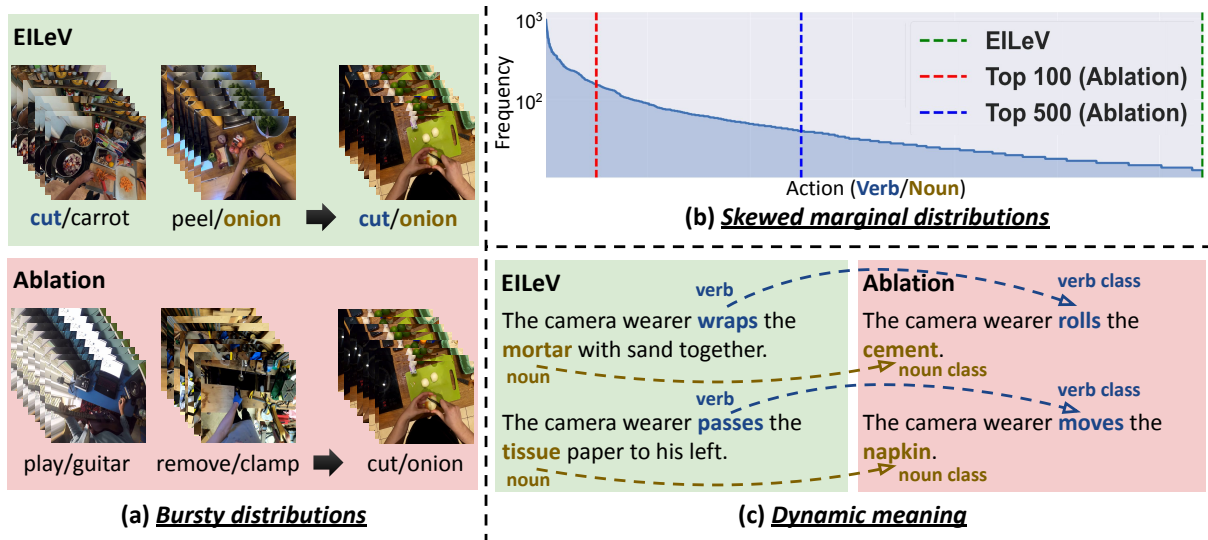


Figure 1: In our proposed training procedure **EILeV**, we ensure that the training data satisfy the following three properties: (a) bursty distributions, (b) skewed marginal distributions, and (c) dynamic meanings. Then, we ablate each property to demonstrate its importance. We ablate property (a) by randomly sampling in-context examples; we ablate property (b) by varying the number of common actions in the training data; we ablate property (c) by canonicalizing verbs and nouns using their corresponding verb and noun classes.

convert the action narrations into question-answer pairs where the narrations are the answers, e.g., e.g., *What is the camera wearer doing? The camera wearer cuts a carrot.* We vary the syntactic form of questions using a set of templates (Appendix C). The training objective is to maximize the likelihood of the sequence of tokens in the ground-truth action narration, conditioned on the context and video clip from the query.

Next, we discuss how each distributional property was incorporated and ablated in **EILeV**.

Bursty Distributions In order to implement bursty distributions in **EILeV**, we take advantage of the annotations in Ego4D, where each video clip is annotated with a verb class and a noun class based on the main action portrayed in the clip. Specifically, we sample video clips and action narrations that share the same verb class as the query for half of the context, and we sample those with the same noun class for the other half. We further ensure that none of the sampled video clips and action narrations match both the verb class and noun class of the query simultaneously. This ensures that the context, while comprising a “burst” of similar concepts, only provides partial information regarding the query. This property can then be ablated by randomly sampling video clips and action narrations without regard to their verb and noun classes. Figure 1 (a) illustrates the two sampling strategies. We can measure the impact of bursty distributions

by training VLMs with each type of context and comparing their in-context learning capabilities.

Skewed Marginal Distributions Like most natural datasets, Ego4D’s verb and noun class labels have a skewed marginal distribution with a long tail of verb-noun pairs, making it ideal for our study. To study how the skewed marginal distributions of training data affect the in-context learning capability of trained models, we first use the verb and noun class annotations from Ego4D to designate the most frequent 80% verb-noun pairs as *common actions* for training, and the remaining 20% as *rare actions* only for evaluation. It is important to note that while none of the rare actions are part of the common action training data, they may still share either verb or noun classes with common actions. For example, if the training data contain common actions (*put, key*) and (*sit, bench*), there may exist a rare action (*put, bench*) in the evaluation data.

To measure how the skewness of marginal distributions in the training data impacts models’ capability to generalize to these novel held-out actions, we then vary the number of common actions in the training data through three experiments. Specifically, we construct a training dataset with only the top 100 common actions (little skewness without a long tail of infrequent actions), one with the top 500 common actions (moderate skewness with a short tail of infrequent actions) and another with all the common actions (highly skewed with a long

tail of infrequent items). We uniformly upsample the datasets with top 100 and top 500 common actions to keep all three training datasets to be the same size. Figure 1 (b) shows how these training datasets with different marginal distributions are constructed. Given these curated training datasets, we can measure the impact of the skewness of the marginal distributions of the training data on trained models’ in-context learning capability.

Dynamic Meaning For dynamic meaning, we rely on the fact that Ego4D’s natural language action narrations contain words of multiple senses, homonyms, and synonyms. To ablate this dynamic meaning property in **EILeV**, we canonicalize verbs and their corresponding objects in the action narrations. Specifically, we prompt an LLM (Llama-2-Chat 7B; Touvron et al., 2023) to replace the verb and its corresponding object of each action narration with their verb and noun class. Figure 1 (c) shows the canonicalization process. We can then measure the impact of dynamic meaning by comparing the in-context learning capability of VLMs trained on data with and without this property.

4.3 Model

To experiment with **EILeV** as discussed above, we adopt a VLM architecture capable of processing sequential data interleaved with both video clips and texts, making it possible to infer patterns and relationships among them and thus support the emergence of in-context learning over them. We initialize our model with BLIP-2 (Li et al., 2023c), a VLM created by learning a transformer-based projection (called a querying transformer or Q-Former) from a frozen pre-trained vision encoder into the input space of a frozen LLM. Since BLIP-2’s original implementation is not able to handle data interleaved with video clips and texts, we follow Hao et al. (2022) to perform simple modifications to enable its frozen language model to serve as a universal interface for video clips and texts.¹ Specifically, we first encode all the video clips by independently encoding sampled frames with BLIP-2’s frozen Vision Transformer (ViT)-based (Dosovitskiy et al., 2021) vision encoder to produce a sequence of vision tokens for each video clip. The sequence of vision tokens is then compressed by

¹While there exist VLMs that already natively support interleaved video and text (Alayrac et al., 2022; Awadalla et al., 2023; Li et al., 2023b), we intentionally chose a VLM that did not to isolate the impact of our **EILeV** training paradigm on VLMs’ in-context learning capability.

BLIP-2’s Q-Former into a fixed-length sequence. The fixed-length sequence is further projected to the word embedding space of the frozen language model of BLIP-2 by a linear layer. It is then interleaved with the text tokens according to the order in which video clips and texts appear in the interleaved data to form the input to the frozen language model. Following the fine-tuning procedure of Li et al. (2023c), we freeze the vision encoder and language model of the BLIP-2 models during training. For all of our experiments, we use BLIP-2 with 2.7 billion parameter OPT (Zhang et al., 2022) as its frozen language model (BLIP-2 OPT-2.7B), and BLIP-2 with XL-size Flan-T5 (Wei et al., 2022) as its frozen language model (BLIP-2 Flan-T5-xl).

4.4 Evaluation

To evaluate our various model ablations, we need a means to measure the quality of action narrations generated by models, and the degree to which in-context learning supports this generation.

4.4.1 Action Narration Generation

One major difficulty in evaluating generative models for the action narration generation task is that there is no single correct way to describe the action in a video clip. In an ideal world, we would rely on human annotators to rate how close a generated action narration is to the ground truth, but the cost to do so would be prohibitive. In order to address this challenge, a number of semantic-similarity-based metrics (Zhang et al., 2019; Reimers and Gurevych, 2019) that correlate closely with human judgment have been proposed, and we take advantage of them in our evaluations. Specifically, we report the performance along semantic similarity-based scores produced by Siamese Sentence-BERT Bi-Encoder (STS-BE; Reimers and Gurevych, 2019). For completeness, we also report ROUGE-L (Lin, 2004), a lexical-based text generation metric.

4.4.2 In-Context Learning Capability

To evaluate the in-context learning capability of trained models for action narration, we vary the number of in-context examples in context-query instances (different numbers of “shots”) and calculate the above text generation metrics for generated action narrations on the test set. If adding more shots improves narration quality under these metrics, this suggests that the VLM is successfully using in-context learning to adapt to the action narration generation task. Within a single experiment

setting, we use the same pre-sampled in-context examples to ensure fair comparison.

5 Experimental Results

In our experiments, we find that the performance of both **EILeV**-trained models strictly increases as more in-context examples (shots) are provided, indicating that **our models successfully acquired in-context learning capabilities during training**. First, in Section 5.1, we establish the in-context learning capability of our models by measuring their performance on rare actions they were not trained on (the key challenge motivating this work), and compare their performance to that of off-the-shelf VLMs. In Sections 5.2, 5.3, and 5.4, we compare their performance to that of models trained on datasets with each key distributional property ablated (as described in Section 4.2) to explore the impact of these training data properties on in-context learning for video and text in VLMs.

5.1 Generalization to Rare Actions

We first compare our **EILeV**-trained models with existing off-the-shelf VLMs in the challenging practical setting that motivated this work: *adaptation to rare actions*. Specifically, we evaluate our models, Kosmos-2 (Peng et al., 2023), and Otter (Li et al., 2023b) on the evaluation set of held-out rare action videos from Ego4D described in Section 4.2.² We choose these two models as they are the only open-source large VLMs that support video input and in-context-learning out-of-the-box at the time of writing. Compared to our **EILeV**-trained models, these models have been trained on far more multi-modal interleaved (MMI) data directly related to in-context learning over video (Table 1), as well as other naturalistic multi-modal and text data from the Internet. They also have far more trainable parameters: Kosmos-2 has 1.6 billion and Otter has 1.3 billion, while our models have 188 million (the same number as BLIP-2). Further, unlike our architectural modification that represents each video with a fixed-length sequence, Kosmos-2 and Otter both treat each video as a sequence of images. For an evaluation representative of the practical usage of VLMs, we do not fine-tune models (which requires prohibitive computing power).

²Our models were not trained on these rare actions, and Kosmos-2 was not trained on Ego4D. While Otter was trained on Ego4D, the video-text training data was not interleaved as proposed for **EILeV**-trained models, and the low frequency of these actions nevertheless poses a significant challenge.

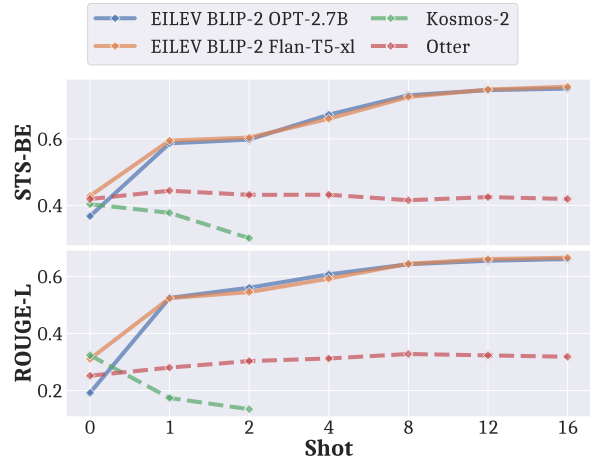


Figure 2: Performance of off-the-shelf VLMs (Kosmos-2 and Otter) on the evaluation set of rare actions for the skewed marginal distributions ablation experiment.

| Model | MMI Dataset Size |
|---|--|
| EILeV BLIP-2 OPT-2.7B & Flan-T5-xl | 115K context-query instances |
| Kosmos-2 | 71M image-text webpages (Huang et al., 2023) |
| Otter | 101.2M image-text webpages (Zhu et al., 2023b) & 2.8M context-query instances (Li et al., 2023a) |

Table 1: Off-the-shelf and **EILeV**-trained VLMs and their multi-modal interleaved (MMI) dataset sizes.

Instead, we rely solely on models’ in-context learning capability to adapt to these rare actions.

Figure 2 shows the results of this evaluation.³ While the zero-shot performance of our **EILeV**-trained models is similar to Kosmos-2 and Otter, **as we provide in-context examples, the performance of our models increases while that of off-the-shelf VLMs does not**. Consequently, our **EILeV-trained VLMs significantly outperform off-the-shelf VLMs**. While Kosmos-2 and Otter have not been fine-tuned on this exact data, they are much larger models trained on an enormous amount of naturalistic data, and their in-context learning capability is a main selling point thought to remove the need for task-specific fine-tuning. Therefore, it is reasonable to expect their performance to improve with more in-context examples

³We can only perform evaluations up to 2-shot with Kosmos-2, as it runs out of its context window beyond 2-shot.

or even outperform our models. This observation underscores that *training smaller VLMs with a focused approach like EILeV can be advantageous for certain use-cases*, such as generating narrations for novel, rare actions, than training large, generalist VLMs on huge naturalistic datasets.

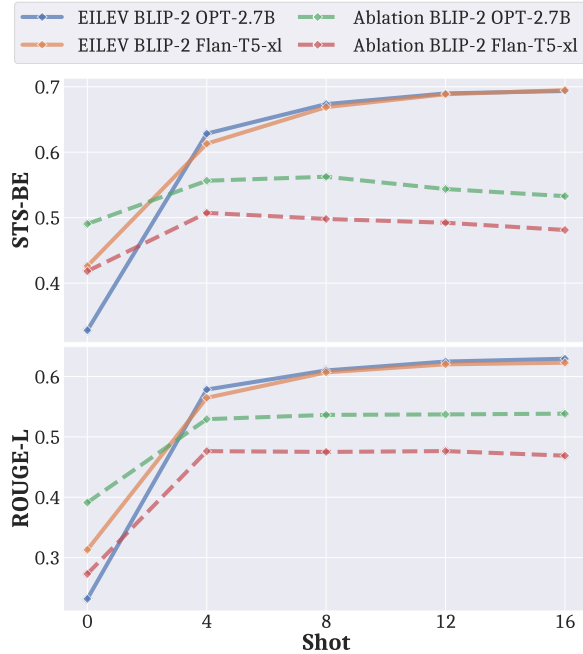


Figure 3: Results for the bursty distributions ablation experiment.

5.2 Bursty Distributions Ablation

Figure 3 shows the results of the bursty distributions ablation experiment. To maintain the same action distributions in both the training and test sets, we use a random train-test split with a ratio of 75/25 for this experiment. Unlike the EILeV-trained models, the performance of the models trained on randomly sampled in-context examples (ablation) initially improves from 0-shot to 4-shot, but tapers or even decreases as more examples are provided. This indicates that they failed to acquire in-context learning capabilities during training, suggesting that **bursty distributions are indeed necessary for in-context learning on video and text**. We hypothesize that the initial improvement in performance from 0-shot to 4-shot is mainly due to the fact that ablation models have learned to mimic lexical characteristics from in-context examples. However, as they have failed to learn to exploit the semantic information from in-context examples due to the lack of bursty distributions in training data, they do not benefit from additional in-context examples.

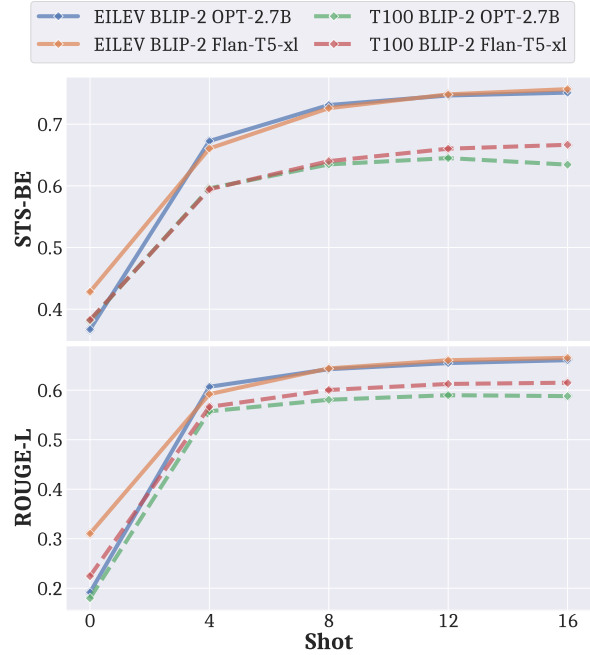


Figure 4: Results for the skewed marginal distributions ablation experiment using a training dataset with top 100 common actions (T100).

5.3 Skewed Marginal Distributions Ablation

Figures 4 and 5 show the results of the skewed marginal distribution ablation experiment. The T100 models trained on data with only the top 100 common actions (little skewness without a long tail of infrequent actions) show a noticeably inferior in-context learning performance to the EILeV-trained models that were trained on the training dataset with all the common actions (highly skewed with a long tail of infrequent items). On the other hand, the T500 models trained on data with the top 500 common actions (moderate skewness with a short tail of infrequent actions) show an in-context learning performance that is only slightly worse than the EILeV-trained models, indicating that **an increased amount of skewness with a long tail of infrequent items makes in-context learning more likely to appear in VLMs**. Further, we observe that the T500 models outperform their respective EILeV-trained models in the 0-shot setting. This is an instance of in-context versus in-weights learning tradeoff (also studied in Chan et al., 2022), a phenomenon where in-context learning capability can reduce pre-trained models' ability to utilize knowledge encoded in their weights during pre-training. Interestingly, we do not observe this pattern with the T100 models, perhaps because the less diverse training data is not representative enough for models to gain sufficient in-weights knowledge.

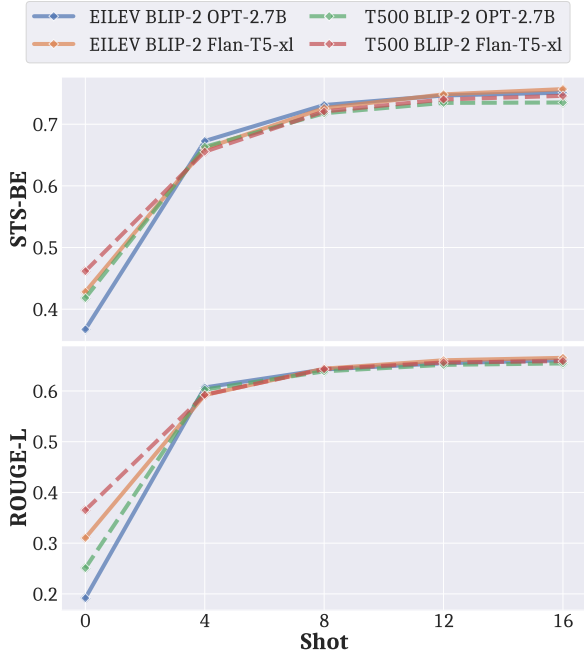


Figure 5: Results for the skewed marginal distributions ablation experiment using a training dataset with top 500 common actions (T500).

5.4 Dynamic Meaning Ablation

Figure 6 shows the results of the dynamic meaning ablation experiment. We use a random train-test split with a ratio of 75/25 for this experiment to maintain the same action distributions in both the training and test sets. The ablation models trained on data with verbs and their corresponding objects canonicalized surprisingly acquire some in-context learning capabilities, but the **EILeV**-trained models mostly outperform them. Since the performance gaps under this ablation are smaller than that of the previous ablations, this suggests that **while dynamic meaning plays a role in the in-context capabilities of a VLM, it contributes less than bursty and skewed marginal distributions do**. Interestingly, however, the performance gap is much more pronounced for STS-BE (semantic similarity metric) than ROUGE-L (lexical metric), suggesting that dynamic meaning contributes more to the model’s ability to extract semantic information from in-context examples than lexical information.

6 Conclusion

In this work, we conducted a first-of-its-kind systematic investigation of in-context learning in vision-language models (VLMs) trained on videos and text. Specifically, we implemented **Emergent In-context Learning on Videos (EILeV)**, a novel training paradigm capturing three key properties of

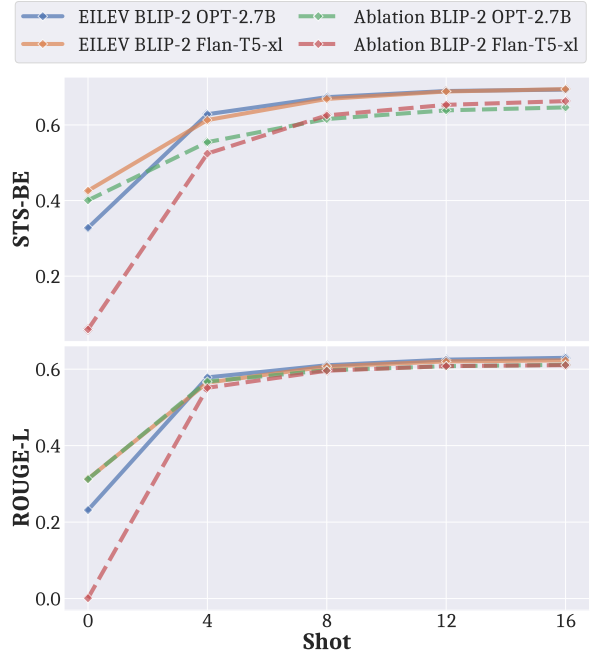


Figure 6: Results for the dynamic meaning ablation experiment.

training data found to induce in-context learning in transformers (Chan et al., 2022): bursty distributions, skewed marginal distributions, and dynamic meaning. In our experiments, we showed that our **EILeV**-trained models exhibit in-context learning capabilities superior to that of off-the-shelf VLMs, as they were significantly more adaptable to novel, rare actions. We demonstrated that all three of these properties are indeed important to optimize the in-context learning capabilities of these models on narrating actions in videos, especially bursty and skewed marginal distributions.

Our work yields new insights about the nature of in-context learning in video and text. For example, we observed that while reducing the skewness of the training data distribution compromised in-context learning capability, it improved in-weights learning in trained models (Chan et al., 2022). We also found that dynamic meaning had a bigger impact on semantic similarity metrics for generated narrations than lexical metrics, suggesting this property is particularly important for acquiring semantic information through in-context learning.

While we focused on action narration in Ego4D (Grauman et al., 2022) as a proof-of-concept, **EILeV** serves as a foundation for the community to build VLMs capable of in-context learning on video and text in broader tasks and domains. We release our **EILeV**-trained models as a resource for future work in egocentric video narration.

7 Limitations

Since our **EILeV**-trained models are optimized and evaluated for action narration generation on ego-centric video using in-context learning, their ability to generalize to diverse, real-world scenarios may be limited. However, this focus was by design and necessity. The primary goal of this work was to verify that the three distributional properties identified by Chan et al. (2022) also elicit in-context learning capabilities in VLMs for videos. To that end, we intentionally chose to use Ego4D, a dataset with sufficient annotations to enable our systematic ablation experiments as a proof of concept. Despite this limitation, **EILeV**-trained models may retain some capability to answer other types of questions due to the use of a frozen language model. Furthermore, **EILeV** is a general training method that can be applied to other tasks given the appropriate data.

Additionally, our models may inherit biases from their frozen language models, making it possible that they could generate harmful content. Before deploying such a system for real-world applications, safety measures like guardrails and training data sanitization are crucial to minimize potential negative impact. On the other hand, since we used the diverse and global data from Ego4D to train our models, this may mitigate possible socio-economic bias found in pre-trained visual representations (Nwatu et al., 2023).

References

Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob L Menick, Sebastian Borgeaud, Andy Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikołaj Białkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karén Simonyan. 2022. [Flamingo: a visual language model for few-shot learning](#). In *Advances in Neural Information Processing Systems*, volume 35, pages 23716–23736. Curran Associates, Inc.

Anas Awadalla, Irena Gao, Josh Gardner, Jack Hessel, Yusuf Hanafy, Wanrong Zhu, Kalyani Marathe, Yonatan Bitton, Samir Gadre, Shiori Sagawa, et al. 2023. [Openflamingo: An open-source framework for training large autoregressive vision-language models](#). *arXiv preprint arXiv:2308.01390*.

Yuwei Bao, Keunwoo Peter Yu, Yichi Zhang, Shane Storks, Itamar Bar-Yossef, Alexander De La Iglesia, Megan Su, Xiao Lin Zheng, and Joyce Chai.

2023. [Can foundation models watch, talk and guide you step by step to make a cake?](#) *arXiv preprint arXiv:2311.00738*.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.

Stephanie Chan, Adam Santoro, Andrew Lampinen, Jane Wang, Aaditya Singh, Pierre Richemond, James McClelland, and Felix Hill. 2022. Data distributional properties drive emergent in-context learning in transformers. *Advances in Neural Information Processing Systems*, 35:18878–18891.

Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. 2023. [Instructblip: Towards general-purpose vision-language models with instruction tuning](#). *Preprint, arXiv:2305.06500*.

Dima Damen, Hazel Doughty, Giovanni Maria Farinella, Antonino Furnari, Jian Ma, Evangelos Kazakos, Davide Moltisanti, Jonathan Munro, Toby Perrett, Will Price, and Michael Wray. 2022. [Rescaling egocentric vision: Collection, pipeline and challenges for epic-kitchens-100](#). *International Journal of Computer Vision (IJCV)*, 130:33–55.

Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. 2021. [An image is worth 16x16 words: Transformers for image recognition at scale](#). In *International Conference on Learning Representations*.

Yuqing Du, Ksenia Konyushkova, Misha Denil, Akhil Raju, Jessica Landon, Felix Hill, Nando de Freitas, and Serkan Cabi. 2023. [Vision-language models as success detectors](#). In *Proceedings of The 2nd Conference on Lifelong Learning Agents*, volume 232 of *Proceedings of Machine Learning Research*, pages 120–136. PMLR.

Tao Gong, Chengqi Lyu, Shilong Zhang, Yudong Wang, Miao Zheng, Qian Zhao, Kuikun Liu, Wenwei Zhang, Ping Luo, and Kai Chen. 2023. [Multimodal-gpt: A vision and language model for dialogue with humans](#). *arXiv preprint arXiv:2305.04790*.

Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary Chavis, Antonino Furnari, Rohit Girdhar, Jackson Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, Miguel Martin, Tushar Nagarajan, Ilija Radosavovic, Santhosh Kumar Ramakrishnan, Fiona Ryan, Jayant Sharma, Michael Wray, Mengmeng Xu, Eric Zhongcong Xu, Chen Zhao, Siddhant Bansal, Dhruv Batra, Vincent Cartillier, Sean Crane, Tien Do, Morrie Doulaty, Akshay Erapalli, Christoph Feichtenhofer, Adriano Fragomeni, Qichen Fu, Christian

| | | |
|-----|--|---|
| 724 | Fuegen, Abrham Gebreselasie, Cristina Gonzalez, James Hillis, Xuhua Huang, Yifei Huang, Wenqi Jia, Weslie Khoo, Jachym Kolar, Satwik Kottur, Anurag Kumar, Federico Landini, Chao Li, Yanghao Li, Zhenqiang Li, Karttikeya Mangalam, Raghava Modhugu, Jonathan Munro, Tullie Murrell, Takumi Nishiyasu, Will Price, Paola Ruiz Puentes, Merey Ramazanovna, Leda Sari, Kiran Somasundaram, Audrey Southerland, Yusuke Sugano, Ruijie Tao, Minh Vo, Yuchen Wang, Xindi Wu, Takuma Yagi, Yunyi Zhu, Pablo Arbelaez, David Crandall, Dima Damen, Giovanni Maria Farinella, Bernard Ghanem, Vamsi Krishna Ithapu, C. V. Jawahar, Hanbyul Joo, Kris Kitani, Haizhou Li, Richard Newcombe, Aude Oliva, Hyun Soo Park, James M. Rehg, Yoichi Sato, Jianbo Shi, Mike Zheng Shou, Antonio Torralba, Lorenzo Torresani, Mingfei Yan, and Jitendra Malik. 2022. Ego4d: Around the World in 3,000 Hours of Ego-centric Video. In <i>IEEE/CVF Computer Vision and Pattern Recognition (CVPR)</i> . | |
| 744 | Yaru Hao, Haoyu Song, Li Dong, Shaohan Huang, Zewen Chi, Wenhui Wang, Shuming Ma, and Furu Wei. 2022. Language models are general-purpose interfaces. <i>arXiv preprint arXiv:2206.06336</i> . | |
| 748 | Shaohan Huang, Li Dong, Wenhui Wang, Yaru Hao, Saksham Singhal, Shuming Ma, Tengchao Lv, Lei Cui, Owais Khan Mohammed, Qiang Liu, et al. 2023. Language is not all you need: Aligning perception with language models. <i>arXiv preprint arXiv:2302.14045</i> . | |
| 754 | Shaohan Huang, Li Dong, Wenhui Wang, Yaru Hao, Saksham Singhal, Shuming Ma, Tengchao Lv, Lei Cui, Owais Khan Mohammed, Barun Patra, et al. 2024. Language is not all you need: Aligning perception with language models. <i>Advances in Neural Information Processing Systems</i> , 36. | |
| 760 | Hugo Laurençon, Lucile Saulnier, Leo Tronchon, Stas Bekman, Amanpreet Singh, Anton Lozhkov, Thomas Wang, Siddharth Karamcheti, Alexander M Rush, Douwe Kiela, et al. 2023. Obelics: An open web-scale filtered dataset of interleaved image-text documents. In <i>Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track</i> . | |
| 768 | Bo Li, Yuanhan Zhang, Liangyu Chen, Jinghao Wang, Fanyi Pu, Jingkang Yang, Chunyuan Li, and Ziwei Liu. 2023a. Mimic-it: Multi-modal in-context instruction tuning. <i>arXiv preprint arXiv:2306.05425</i> . | |
| 772 | Bo Li, Yuanhan Zhang, Liangyu Chen, Jinghao Wang, Jingkang Yang, and Ziwei Liu. 2023b. Otter: A multi-modal model with in-context instruction tuning. <i>arXiv preprint arXiv:2305.03726</i> . | |
| 776 | Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023c. BLIP-2: bootstrapping language-image pre-training with frozen image encoders and large language models. In <i>ICML</i> . | |
| | Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. 2022. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In <i>International Conference on Machine Learning</i> , pages 12888–12900. PMLR. | 780 781 782 783 784 |
| | KunChang Li, Yinan He, Yi Wang, Yizhuo Li, Wenhai Wang, Ping Luo, Yali Wang, Limin Wang, and Yu Qiao. 2023d. Videochat: Chat-centric video understanding. <i>arXiv preprint arXiv:2305.06355</i> . | 785 786 787 788 |
| | Bin Lin, Bin Zhu, Yang Ye, Munan Ning, Peng Jin, and Li Yuan. 2023. Video-llava: Learning united visual representation by alignment before projection. <i>arXiv preprint arXiv:2311.10122</i> . | 789 790 791 792 |
| | Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In <i>Text summarization branches out</i> , pages 74–81. | 793 794 795 |
| | Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. Visual instruction tuning. In <i>NeurIPS</i> . | 796 797 |
| | Ilya Loshchilov and Frank Hutter. 2018. Decoupled weight decay regularization. In <i>International Conference on Learning Representations</i> . | 798 799 800 |
| | Jiasen Lu, Christopher Clark, Sangho Lee, Zichen Zhang, Savya Khosla, Ryan Marten, Derek Hoiem, and Aniruddha Kembhavi. 2023. Unified-io 2: Scaling autoregressive multimodal models with vision, language, audio, and action. <i>arXiv preprint arXiv:2312.17172</i> . | 801 802 803 804 805 806 |
| | Muhammad Maaz, Hanoona Rasheed, Salman Khan, and Fahad Shahbaz Khan. 2023. Video-chatgpt: Towards detailed video understanding via large vision and language models. <i>arXiv preprint arXiv:2306.05424</i> . | 807 808 809 810 811 |
| | Brandon McKinzie, Zhe Gan, Jean-Philippe Fauconnier, Sam Dodge, Bowen Zhang, Philipp Dufter, Dhruvi Shah, Xianzhi Du, Futang Peng, Floris Weers, Anton Belyi, Haotian Zhang, Karanjeet Singh, Doug Kang, Ankur Jain, Hongyu He, Max Schwarzer, Tom Gunter, Xiang Kong, Aonan Zhang, Jianyu Wang, Chong Wang, Nan Du, Tao Lei, Sam Wiseman, Guoli Yin, Mark Lee, Zirui Wang, Ruoming Pang, Peter Grasch, Alexander Toshev, and Yinfei Yang. 2024. Mml: Methods, analysis & insights from multimodal llm pre-training. <i>ArXiv</i> , abs/2403.09611. | 812 813 814 815 816 817 818 819 820 821 822 |
| | Sewon Min, Mike Lewis, Luke Zettlemoyer, and Hananeh Hajishirzi. 2022. Metaicl: Learning to learn in context. In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 2791–2809. | 823 824 825 826 827 828 |
| | Joan Nwatu, Oana Ignat, and Rada Mihalcea. 2023. Bridging the digital divide: Performance variation across socio-economic factors in vision-language models. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 10686–10702, Singapore. Association for Computational Linguistics. | 829 830 831 832 833 834 835 |

| | | |
|-----|--|-----|
| 836 | Zhiliang Peng, Wenhui Wang, Li Dong, Yaru Hao, Shaohan Huang, Shuming Ma, and Furu Wei. 2023. Kosmos-2: Grounding multimodal large language models to the world. <i>arXiv preprint arXiv:2306.14824</i> . | 891 |
| 837 | | 892 |
| 838 | | 893 |
| 839 | | 894 |
| 840 | | 895 |
| 841 | Toby Perrett, Saptarshi Sinha, Tilo Burghardt, Majid Mirhemdi, and Dima Damen. 2023. Use your head: Improving long-tail video recognition. In <i>Computer Vision and Pattern Recognition</i> . | 896 |
| 842 | | 897 |
| 843 | | 898 |
| 844 | | 899 |
| 845 | Gautam Reddy. 2023. The mechanistic basis of data dependence and abrupt learning in an in-context classification task. In <i>The Twelfth International Conference on Learning Representations</i> . | 900 |
| 846 | | 901 |
| 847 | | 902 |
| 848 | | 903 |
| 849 | Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing</i> . Association for Computational Linguistics. | 904 |
| 850 | | 905 |
| 851 | | 906 |
| 852 | | 907 |
| 853 | | 908 |
| 854 | Zhan Tong, Yibing Song, Jue Wang, and Limin Wang. 2022. VideoMAE: Masked autoencoders are data-efficient learners for self-supervised video pre-training. In <i>Advances in Neural Information Processing Systems</i> . | 909 |
| 855 | | 910 |
| 856 | | 911 |
| 857 | | 912 |
| 858 | | 913 |
| 859 | Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. <i>arXiv preprint arXiv:2307.09288</i> . | 914 |
| 860 | | 915 |
| 861 | | 916 |
| 862 | | 917 |
| 863 | | 918 |
| 864 | | 919 |
| 865 | Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. <i>Advances in neural information processing systems</i> , 30. | 920 |
| 866 | | 921 |
| 867 | | 922 |
| 868 | | 923 |
| 869 | | 924 |
| 870 | Alex Jinpeng Wang, Linjie Li, Kevin Qinghong Lin, Jianfeng Wang, Kevin Qinghong Lin, Zhengyuan Yang, Lijuan Wang, and Mike Zheng Shou. 2024. Cosmo: Contrastive streamlined multimodal model with interleaved pre-training . <i>ArXiv</i> , abs/2401.00849. | 925 |
| 871 | | 926 |
| 872 | | 927 |
| 873 | | 928 |
| 874 | | 929 |
| 875 | | 930 |
| 876 | Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V Le. 2022. Finetuned language models are zero-shot learners . In <i>International Conference on Learning Representations</i> . | 931 |
| 877 | | 932 |
| 878 | | 933 |
| 879 | | 934 |
| 880 | | 935 |
| 881 | Antoine Yang, Antoine Miech, Josef Sivic, Ivan Laptev, and Cordelia Schmid. 2022. Zero-shot video question answering via frozen bidirectional language models. In <i>NeurIPS</i> . | 936 |
| 882 | | 937 |
| 883 | | 938 |
| 884 | | 939 |
| 885 | Qinghao Ye, Haiyang Xu, Guohai Xu, Jiabo Ye, Ming Yan, Yiyang Zhou, Junyang Wang, Anwen Hu, Pengcheng Shi, Yaya Shi, et al. 2023. mplug-owl: Modularization empowers large language models with multimodality. <i>arXiv preprint arXiv:2304.14178</i> . | 940 |
| 886 | | 941 |
| 887 | | 942 |
| 888 | | 943 |
| 889 | | 944 |
| 890 | | 945 |
| | Rowan Zellers, Jiasen Lu, Ximing Lu, Youngjae Yu, Yanpeng Zhao, Mohammadreza Salehi, Aditya Kusupati, Jack Hessel, Ali Farhadi, and Yejin Choi. 2022. Merlot reserve: Multimodal neural script knowledge through vision and language and sound. In <i>CVPR</i> . | |
| | Hang Zhang, Xin Li, and Lidong Bing. 2023. Videollama: An instruction-tuned audio-visual language model for video understanding. <i>arXiv preprint arXiv:2306.02858</i> . | |
| | Renrui Zhang, Jiaming Han, Chris Liu, Aojun Zhou, Pan Lu, Hongsheng Li, Peng Gao, and Yu Qiao. 2024. LLaMA-adapter: Efficient fine-tuning of large language models with zero-initialized attention . In <i>The Twelfth International Conference on Learning Representations</i> . | |
| | Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. 2022. Opt: Open pre-trained transformer language models. <i>arXiv preprint arXiv:2205.01068</i> . | |
| | Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. Bertscore: Evaluating text generation with bert. In <i>International Conference on Learning Representations</i> . | |
| | Yue Zhao, Ishan Misra, Philipp Krähenbühl, and Rohit Girdhar. 2023. Learning video representations from large language models. In <i>CVPR</i> . | |
| | Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. 2023a. Minigt-4: Enhancing vision-language understanding with advanced large language models. <i>arXiv preprint arXiv:2304.10592</i> . | |
| | Wanrong Zhu, Jack Hessel, Anas Awadalla, Samir Yitzhak Gadre, Jesse Dodge, Alex Fang, Youngjae Yu, Ludwig Schmidt, William Yang Wang, and Yejin Choi. 2023b. Multimodal C4: An open, billion-scale corpus of images interleaved with text. <i>arXiv preprint arXiv:2304.06939</i> . | |
| | A Additional Experiments | |
| | A.1 Additional Baselines | |
| | We report the performance of three additional baselines on the Ego4D-based dataset used in the main ablation experiments, as well as another dataset constructed from EPIC-KITCHENS-100 (Damen et al., 2022). The first is a naive action classification baseline (“VideoMAE”). Specifically, we fine-tune the “videomae-huge-finetuned-kinetics” variant of VideoMAE (Tong et al., 2022) using the verb and noun class annotations to produce a verb and a noun classifier. The predicted verb and noun classes are then transformed into action narrations using an off-the-shelf LLM (7 billion parameter Llama-2-Chat (Touvron et al., 2023)). Note that | |

this baseline only uses videos as its input, and cannot perform in-context learning. The second are off-the-shelf BLIP-2 models with the architectural modifications from Section 4.3 for interleaved data support (“BLIP-2 OPT-2.7B & Flan-T5-xl”). The third are **EILeV**-trained models with in-context examples ablated, and fine-tune solely on the query (“FT BLIP-2 OPT-2.7B & Flan-T5-xl”).

A.1.1 Results on Ego4D

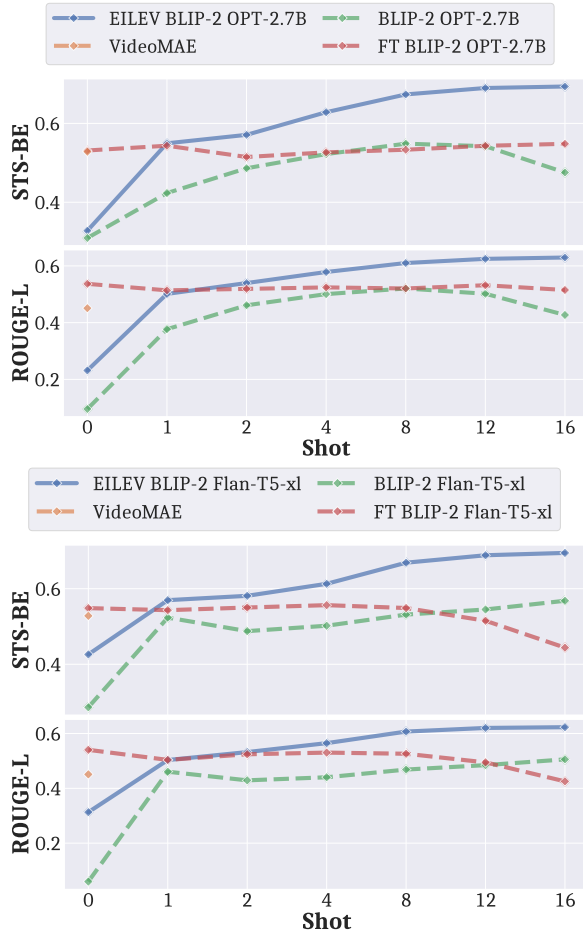


Figure 7: Performance of additional baselines on the Ego4D-based dataset.

Figure 7 reports the performance of the three additional baselines on the Ego4D-based dataset. We use a random train-test split with a ratio of 75/25 for this experiment to maintain the same action distributions in both the training and test sets. The **EILeV**-trained BLIP-2 models demonstrate superior in-context learning capabilities, as their performance improves with an increasing number of shots, ultimately outperforming all baseline models. This is a further indication that **EILeV** has successfully elicited in-context-learning capabilities in them. The VideoMAE and FT BLIP-2 models

exhibit the best performance at 0-shot, suggesting they have the most amount of in-weights knowledge due to their fine-tuning. However, VideoMAE cannot process in-context examples, and its 0-shot performance is quickly outperformed by **EILeV**-trained models with only one in-context example. The performance of FT BLIP-2 models stagnates or even declines as the number of shots increases, highlighting their lack of in-context learning capabilities and the importance of the training data design discussed in Section 4.2. These findings about the performance of different models at 0-shot and subsequent shots align with Chan et al. (2022) observations regarding the “tradeoff between in-context learning and in-weights learning,” where no models could maintain both in their experiments. In our experiment, the **EILeV**-trained BLIP-2 models are optimized for in-context learning, as evidenced by their subpar performance at 0-shot and superior performance with additional shots, whereas the FT BLIP-2 models show the opposite trend. We leave designing training data to find the right balance for future work.

A.1.2 Results on EPIC-KITCHENS-100

Next, we test if **EILeV**-trained BLIP-2 models trained solely on Ego4D can generalize to out-of-distribution actions via in-context learning. Specifically, we evaluate them on the validation split of a different egocentric video dataset, EPIC-KITCHENS-100, without further fine-tuning. Note that there is a significant distributional shift between Ego4D and EPIC-KITCHENS-100 even though they both contain egocentric videos in the kitchen setting as evidenced by the t-SNE plot in Figure 9. All the experimental setups are same as the experiments on the Ego4D-based dataset except evaluation context-query instances are formed by sampling both the context and the query from the validation set of EPIC-KITCHENS-100. Unlike Ego4D, the action narrations from EPIC-KITCHENS-100 are not full sentences, but simple verb-noun phrases. Therefore, we use an LLM (7 billion parameter Llama-2-Chat (Touvron et al., 2023)) to turn the simple verb-noun phrases into full sentences with “the camera wearer” as the subject.

Figure 8 reports the evaluation results. The performance of the **EILeV**-trained BLIP-2 models improves with an increasing number of in-context examples and ultimately outperforms all the baselines. This indicates that these models can generalize to

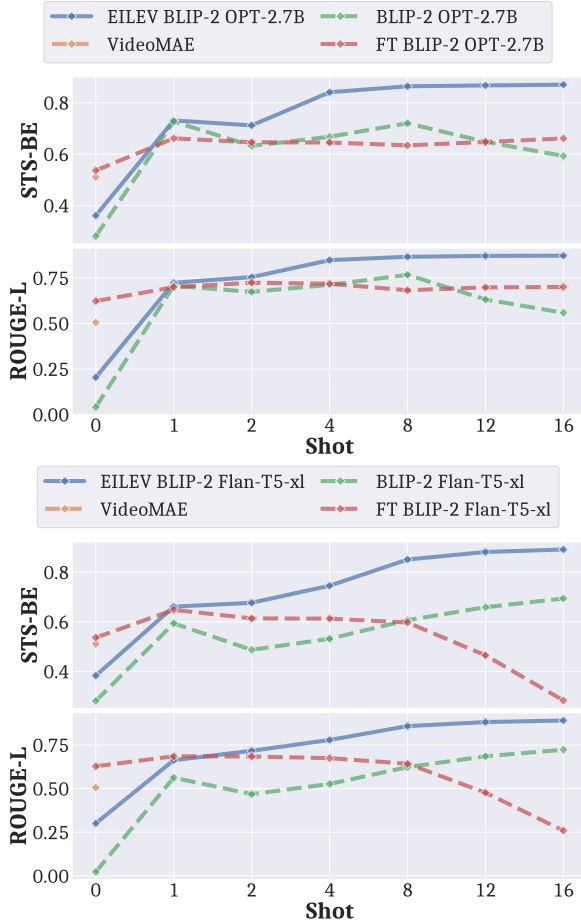


Figure 8: Performance of additional baselines on the EPIC-KITCHENS-100-based dataset

out-of-distribution actions via in-context learning. All the baseline models exhibit similar trends as on the Ego4D-based dataset: they demonstrate the best performance at 0-shot but fail to benefit from the in-context examples.

A.2 In-Context or In-Weights Learning

We now aim to validate that the source of the generalization capabilities demonstrated by the **EILEV**-trained models in Section 5.1 is indeed from in-context learning, not in-weights learning. This is to further reinforce our claim that **EILEV**-trained models can generalize to actions that they have not seen during training, i.e., actions of which they have no direct in-weights knowledge. To that end, we use the frequency of each verb/noun class in the common action training data as the proxy for the knowledge about the verb/noun class encoded into the weights of the model (in-weights learning), and the difference in model performance between 16-shot and 0-shot settings for a particular rare action as the proxy for in-context learning perfor-

mance. If the model relies on in-weights learning for a particular novel, rare action, the difference in performance for that action between 16-shot and 0-shot settings would be correlated to the frequency of the corresponding verb/noun class in the training data. This outcome is not desired, as we want the model to rely on in-context learning for generating accurate narrations of novel, rare actions unseen during training.

Figure 10 shows the scatter plots between the log verb/noun class frequency in the training data and the difference in STS-BE for the corresponding rare action between 16-shot and 0-shot settings for the **EILEV**-trained models. For example, given a rare action (“put”, “bench”), a point on the scatter plot may refer to the log frequency of “put” in the common action training data in the x-axis and the difference in the STS-BE performance of **EILEV** BLIP-2 OPT-2.7B on (“put”, “bench”) between 16-shot and 0-shot. As the scatter plots and their corresponding R^2 values show, there is a minimal linear correlation between the log verb/noun class frequency in the training data and the difference in STS-BE for the corresponding action from in-context learning. This suggests that the **EILEV**-trained models generate accurate narrations for novel, rare actions via in-context learning rather than in-weights learning, as the linear model does not significantly account for the variance in the observed data.

A.3 Context Modeling and In-Context Learning

In this evaluation, we seek to investigate if the **EILEV**-trained models perform correct context modeling by incorporating the relationships between video clips and narrations. To that end, we evaluate the **EILEV**-trained models and the off-the-shelf BLIP-2 baseline models from Section A.1 on shuffled in-context examples where video clips no longer match the action narrations. We then compare their performance from shuffled in-context examples (the treatment group) to the one from unshuffled in-context examples as the control group. If the performance remains unchanged, it implies that the model does not consider the relationships between in-context video clips and action narrations. On the other hand, if the performance decreases, it implies that the model does take the relationships between video clips and action narrations into account, and the mismatch adversely affects its performance. We do not report the results at 0 and 1-shot since shuffling of the in-context video

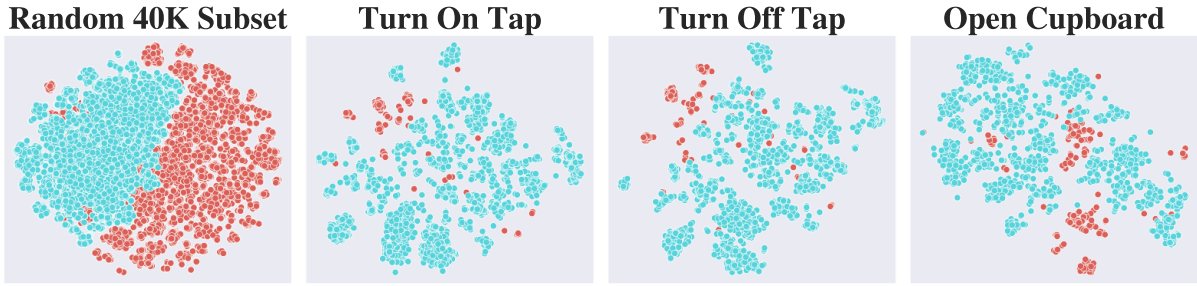


Figure 9: t-SNE plots of the video embeddings from the frozen vision encoder of BLIP-2 OPT-2.7B. Ego4D videos are in red, and EPIC-KITCHENS-100 videos are in blue. Plots for a randomly sampled subset of 40k videos from both and three most common actions from EPIC-KITCHENS-100 are shown. We manually map Ego4D actions to the EPIC-KITCHENS-100 actions.

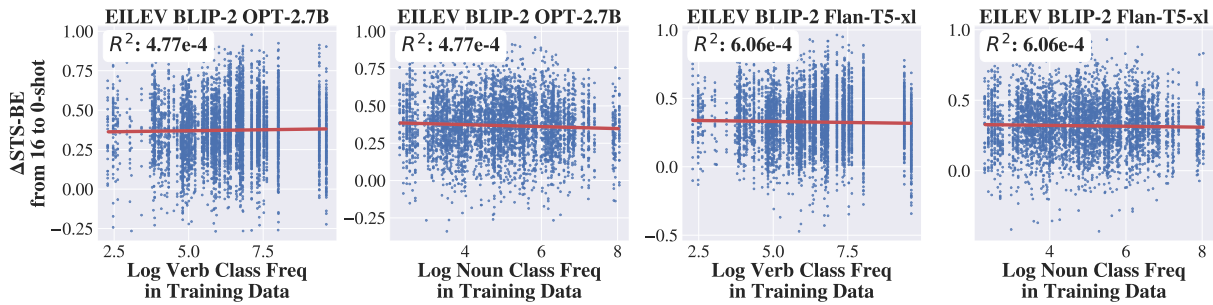


Figure 10: Scatter plots with trend lines and R^2 values between the log verb/noun class frequency in the training data with common actions and the difference in STS-BE (Δ STS-BE) for the corresponding rare action between 16-shot and 0-shot settings for the **EILEV**-trained models.

clips would not have any impact at those settings.

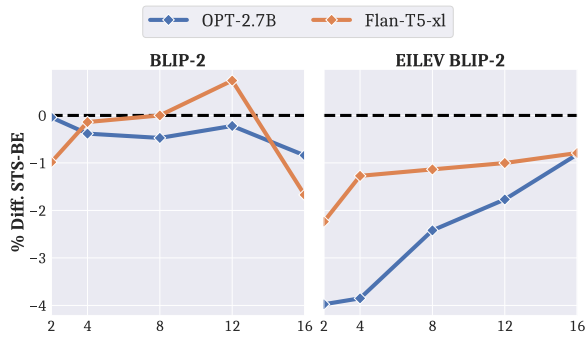


Figure 11: Percentage difference plots between the treatment group with shuffled in-context video clips and the control group. A negative value below the dotted zero line means the STS-BE performance of the treatment group is worse than the control group.

Figure 11 shows the percentage differences in STS-BE from 16-shot to 0-shot between the treatment group and the control group for the **EILEV**-trained models and the off-the-self BLIP-2 models. For the off-the-shelf BLIP-2 models, the percentage differences are small across all shots. This indicates that they rely mostly on the context as a whole rather than the semantic details from the relationships between video clips and action nar-

rations when performing in-context learning. We hypothesize that our proposed architectural modifications (Section 4.3 allow the off-the-shelf BLIP-2 models to tap into the text-only in-context learning capabilities of their frozen language models, which lack the ability to extract semantic details from the relationships between video clips and action narrations. This hypothesis is supported by their subpar in-context learning capabilities from Section A.1, which speaks to the importance of our modifications to the training data. On the other hand, there is a clear drop in performance for the **EILEV**-trained models in terms of the semantic-similarity-based metric STS-BE. This indicates that the **EILEV**-trained models extract detailed semantic information from the correspondence between in-context video clips and action narrations.

B Training Details

In all of our experiments, each video clip is created by taking the four seconds before and after its action narration timestamp, and 8 frames are sampled uniformly from each video clip. The total training batch size is 128 and the optimizer is AdamW (Loshchilov and Hutter, 2018) with the

1121 initial learning rate of 1×10^{-5} , weight decay of
1122 0.05 and a linear scheduler. We train for 5 epochs
1123 on 8 NVIDIA A40 GPUs using distributed data
1124 parallel. We evaluate every 200 steps and select the
1125 model with the lowest loss. The training time is
1126 about a day and a half.

1127 **C Question Templates**

1128 Table 2 shows the question-answer pair templates
1129 we use in our experiments. They are based on the
1130 instruction templates proposed by [Dai et al. \(2023\)](#).

Table 2: List of question-answer pair templates.

| |
|---|
| What is the camera wearer doing? {narration} |
| Question: What is the camera wearer doing? {narration} |
| What is the camera wearer doing? An answer to the question is {narration} |
| Q: What is the camera wearer doing? A: {narration} |
| Given the video, answer the following question. What is the camera wearer doing? {narration} |
| Based on the video, respond to this question: What is the camera wearer doing? Answer: {narration} |
| Use the provided video to answer the question: What is the camera wearer doing? {narration} |
| What is the answer to the following question? "What is the camera wearer doing?" {narration} |
| The question "What is the camera wearer doing?" can be answered using the video. The answer is {narration} |
