

VideoCon: Robust Video-Language Alignment via Contrast Captions

Hritik Bansal¹ Yonatan Bitton² Idan Szpektor^{2*} Kai-Wei Chang^{1*} Aditya Grover^{1*}
¹UCLA ²Google Research
 {hbansal, kwchang, adityag}@cs.ucla.edu
 {yonatanbitton, szpektor}@google.com

Abstract

Despite being (pre)trained on a massive amount of data, state-of-the-art video-language alignment models are not robust to semantically-plausible contrastive changes in the video captions. Our work addresses this by identifying a broad spectrum of contrast misalignments, such as replacing entities, actions, and flipping event order, which alignment models should be robust against. To this end, we introduce the VideoCon, a video-language alignment dataset constructed by a large language model that generates plausible contrast video captions and explanations for differences between original and contrast video captions. Then, a generative video-language model is fine-tuned with VideoCon to assess video-language entailment and generate explanations. Our VideoCon-based alignment model significantly outperforms current models. It exhibits a 12-point increase in AUC for the video-language alignment task on human-generated contrast captions. Finally, our model sets new state of the art zero-shot performance in temporally-extensive video-language tasks such as text-to-video retrieval (SSv2-Temporal) and video question answering (ATP-Hard). Moreover, our model shows superior performance on novel videos and human-crafted captions and explanations.

1. Introduction

Semantically aligning data points from diverse modalities is a long-standing goal of AI. We focus on video-language alignment, which is challenging due to the complexities involved in understanding of entities, their relationships, and temporal order of the depicted events [17]. Recent models such as VideoCLIP [55], ImageBind [14] learn a shared embedding space. Similarly, generative models such as Flamingo [1], mPLUG-Owl-Video [61] can provide a classification label (e.g., yes/no) when queried about video-language alignment.

Despite large-scale pretraining, prior work [5, 37, 38, 51] highlights that video-language alignment models are not robust to semantically plausible manipulations to an original aligned caption in the form of contrast captions, such as from ‘dog runs away *before* it eats food’ to ‘dog runs away *after* it eats food’. Such pitfalls in robustness questions the trustworthiness of alignment models for large-scale deployment. To mitigate these shortcomings, one possible solution is to scale video-language pairs more for increased diversity during pretraining. However, this is challenging due to the difficulties in sourcing new, high-quality and permissible content, as well as the requirement for substantial storage capacity. Several works [11, 13, 16] have shown that naively training models on web-scale data has diminishing returns on downstream tasks, and emphasize the importance of data quality. Furthermore, the recent studies [28, 62] demonstrate that applying a contrastive objective to the pre-training datasets does not encourage the model to grasp the fine-grained details within image/region-caption data.

To this end, we take a scalable, active strategy to gather high-quality data that is deliberately enriched with the attributes that we want to instill in alignment models. We create a novel dataset, **VideoCon**, to improve the robustness of models. Specifically, the dataset consists of a variety of semantically plausible video-language misalignments in contrast captions. These misalignments include altering *objects (entities), actions, attributes, relations, counts, event orders*, and introducing *hallucinations* (Figure 2). To construct VideoCon, a large language model (PaLM-2 API) takes video-caption pairs as input and generates high-quality contrast captions for a given misalignment type. To make our dataset temporally-challenging, we skipped “easy” video-caption pairs whose alignment could be inferred based on a single frame (image) understanding [9, 26] (§3.1). In addition, the LLM generates natural language explanations (NLEs) [42] to the differences between original and altered captions, which are used for further robust training. We performed human verification on a sample of VideoCon and found that it is of high-quality. Finally, to evaluate the model’s generalization capabilities, we col-

*Equal Advising.

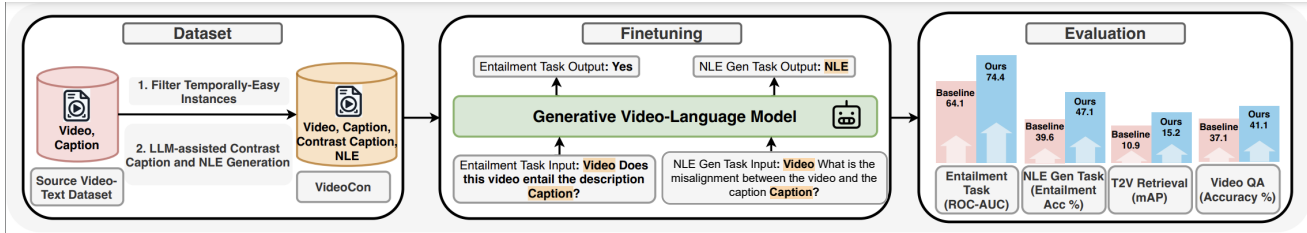


Figure 1. **Overview of our VideoCon approach.** First, aligned video-language pairs are filtered to retain temporally-challenging instances. Then contrast captions and natural language explanations (NLE) are generated by an LLM to create the VideoCon dataset. Second, a video-language alignment model is finetuned with VideoCon on the alignment and NLE tasks. Finally, the finetuned model is evaluated against the baseline model. Our results show that it outperforms the baseline, achieving state-of-the-art results on downstream tasks.

lect human-generated contrast captions and NLEs for the videos sourced from external datasets that did not contribute to VideoCon’s development.

We finetuned a generative video-language model (mPLUG-Owl-Video) on the VideoCon dataset. The trained model surpasses existing video-language alignment models by a large margin on the LLM-generated test set for both video-language alignment and NLE generation tasks. Interestingly, we observed that our finetuned model generalizes to unseen videos and human-generated contrast captions and NLEs, and outperforms the baseline models. For instance, our model’s ROC-AUC exceeds the baseline model by 12 points on the human-generated contrast captions. This indicates that our model has developed a better understanding of the entities, their interactions, action understanding, as well as the temporal order of the events for robust video-language alignment.

We further assessed the effectiveness of robust training via contrast captions on zero-shot downstream video-language tasks such text-to-video retrieval and video question answering on the temporally-challenging and action-intensive SSv2-Temporal [45] and SSv2-Events [5]. Our model achieves state-of-the-art (SOTA) performance, improving on SSv2-Temporal by 4.3 mAP, SSv2-Events by 3.6 mAP points. In addition, our model also achieves SOTA on temporal and causal video question answering in the ATP-Hard dataset, increasing 4% accuracy. This suggests that equipping a model with the knowledge of contrast captions is highly data-efficient and effective in improving its robustness in comparison to scaling the pretraining data. The complete pipeline is illustrated in Figure 1. The dataset and the model will be released upon acceptance.

2. Video Language Alignment

We are interested in assessing the semantic alignment between the video¹ and text data since it powers many prac-

¹Like prior works [33, 55], we use only the video frames (the visual channel) without the soundtrack (the audio channel).

tical applications such as video-text retrieval [57], video generation [7, 47] and video captioning [59]. To this end, [14, 39, 49, 55] designed (image)video-text alignment models that are utilized for evaluating the semantic similarity between the two modalities. However, previous works [5, 37, 38, 51] have questioned their robustness to semantically plausible changes to the video descriptions, termed here *contrast captions*. Our aim is to improve the robustness of video-text alignment models by training on contrast captions with a wide range of misalignments.

Consider a dataset $\mathcal{D} = \{(V_i, T_i, C_i, E_i)\}$ where V_i is a video, T_i is an aligned caption, C_i is a contrast caption which is a perturbation of T_i but misaligns with V_i , and E_i is a natural language explanation for the misalignment between V_i and C_i . We consider two video-language alignment tasks: (a) video-language entailment, (b) natural language explanation.

Video-Language Entailment (VLE) casts video-text alignment as a Visual Entailment (VE) task. VE was originally defined for images as premises and texts as hypothesis [53, 54]. We extend VE definition also for videos as premises, under which a classification model $A_{vle}(V, T)$ predicts whether a video V entails a text T .

Natural Language Explanation (NLE) requires a model, $A_{nle}(V, C)$, to generate an open-ended explanation for the discrepancy between a video V and a non-entailing caption C .

In this paper, we address both VLE and NLE tasks under a multitask setting in which a single video-language generative model generates the binary label for entailment and the open-ended explanation.

3. VideoCon: Contrast Captions Generation for Robust Video-Language Alignment

Our research goal is to measure the impact of a comprehensive dataset on increasing the robustness of video-text align-



Figure 2. **Overview of the VideoCon data generation process from top to bottom.** Specifically, we prompt a large language model (PaLM-2) with the original caption that is grounded in the video, and the intended type of misalignment within the contrast caption. We consider *seven* kinds of misalignments including object, action, attribute, counting, spatial relation, hallucination, and event order flip. We provide a generated contrast caption and the corresponding natural language explanation for each misalignment type.

ment models. To this end, we first collect video-caption pairs where the caption cannot be derived from a single frame of video. We then categorize a wide range of semantically plausible manipulations of video captions. Using an LLM for large-scale computation, contrast captions and related explanations are generated for the defined categories, constructing the VideoCon dataset. Finally, we extend VideoCon to include human-created contrast captions as held-out evaluation on unseen videos. We detail the dataset construction steps below.

3.1. Temporally-Challenging Instance Selection

To construct VideoCon, we start with existing datasets that include natural (real) videos and associated high-quality human-written captions: MSR-VTT [57], VaTeX [48], and TEMPO [17]. MSR-VTT and VaTeX consist of 20 captions and 10 captions per video, respectively, while TEMPO consists of a single caption per video. More dataset details are in Appendix §B.

TEMPO is designed to create temporally-challenging instances, while MSR-VTT and VaTeX contain more general video-caption pairs. For MSR-VTT and VaTeX, we filter out instances, where the caption is highly associated with a single frame in the video based on an image-text alignment model. In such cases, a video-text alignment can leverage shortcuts and align the video to its caption without understanding the temporal or causal relations depicted in the video. We want to filter such instances.

To this end, we employ the End-to-End VNLI model [60] to calculate an alignment score $A_{vle}(V, T)$ between a video $V = \{I_1, I_2, \dots, I_N\}$ and a text T where I_i is a frame from

the video sampled at a rate of 1 frame per second. Formally,

$$A_{vle}(V, T) = \max_i(VNLI(I_i, T)) \quad (1)$$

where $VNLI(I_i, T)$ is the task of visual natural language inference that assesses the extent to which the text T entails the image I_i . There are 20 and 10 captions per video in the MSR-VTT and VaTeX datasets, respectively. We retain 5 captions per video from these datasets with the lowest $A_{vle}(V, T)$, and the remaining captions are filtered out. Post-filtering, the percentage of temporally-challenging instances increased from 36.5% to 81.5% in MSR-VTT, and from 42.6% to 71% in VaTeX.

3.2. Categories of Contrast Captions

We aim for VideoCon to include a wide range of misalignments in its contrast captions. Overall, VideoCon covers *seven* misalignment types, exemplified in Figure 2. We include replacement of *objects* (entities) and *actions* following the analysis in [37, 38], and replacement of *attributes*, *counts*, *relations*, as well as adding unrelated but plausible information to captions as *hallucinations* following [29, 32, 35]’s study of image/text alignment model brittleness. Since most video-text models rely on pretrained image backbones, they are likely to suffer from similar problems. Finally, following [5]’s analysis that video-text models do not understand temporal order of the events, we include *event order flipping* as misalignment type.

3.3. Data Generation using an LLM

To generate contrast captions and corresponding NLE we first assign one of the seven misalignment types (§3.2) to

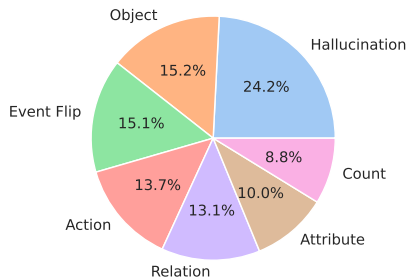


Figure 3. **Distribution of the types of misalignments within the contrast captions of the VideoCon dataset.** We observe that the dataset has good representation for all the kinds of misalignments ranging from 8.8% to 24.2%.

each caption in the input video-text datasets (§3.1) (details in Appendix §C). Then, given a video V and a misalignment type m , we prompt PaLM-2 API² [2] to generate a contrast caption and accompanied explanation (our type-specific prompts are detailed in Appendix §D).

Analyzing the LLM generations, we found that sometimes the output caption C do not contradict the original caption T . For example, a generated contrast caption “a person riding a car” does not contradict the original caption “a person riding a mustang”. To filter such cases, we employ a Natural Language Inference (NLI) model [19] and remove cases in which the contrast caption is assessed as entailed by the original caption $NLI(T, C) > 0.5$. Post-filtering, each tuple (V, T, C, m) is converted to the two instances of video/language entailment task: $A_{en}(V, T) = 1$ and $A_{en}(V, C) = 0$. We present the dataset statistics for the entailment task in Table 1, including train/eval/test splits. In addition, Fig. 3 shows the distribution of misalignment types in the dataset. We observe that VideoCon maintains a high density across the 7 misalignments ranging from 8.8% to 24.2%.

We also found that some generated explanations do not describe the differences between T and C well. For example, the explanation “two friends are not traveling together” does not fully describe the discrepancy between “three friends traveling together” and “two friends are traveling together”. To filter these out, generated examples are removed if $NLI(F(T, C), E) < 0.6$ where $F(T, C)$ is the premise comprising the original and contrast captions. Specifically, premise will be ‘Expected Caption: T Actual Caption: E ’ and hypothesis will be ‘Difference between Expected and Actual Caption: E ’. This filter indicates that the information in the explanation is not entailed by the difference between the two captions. The dataset statistics for the NLE task is presented in Table 1. We refer to the final

²<https://developers.google.com/ai/generativeai/products/palm>

Source	Video-Language Entailment			Natural Language Explanation		
	Train	Val	Test	Train	Val	Test
MSR-VTT	38366	478	16538	15888	206	6788
VaTeX	66480	736	8110	30180	345	3636
TEMPO	10712	7098	2708	4165	2739	1073
Total	115558	8312	27356	50233	3290	11497

Table 1. Statistics for the VLE and NLE tasks in VideoCon.

LLM-generated dataset as VideoCon (LLM).

To assess the quality of VideoCon (LLM), we perform human evaluation on 500 contrast captions and NLEs (details in Appendix E). The human evaluator found 91% of the contrast captions and 89% of the NLEs to be valid, indicating the high-quality of VideoCon (LLM).

3.4. Data Generation using Humans

To study whether a model trained on VideoCon (LLM) generalizes to out-of-distribution videos and its performance on human-generated contrast captions, we randomly selected a set of videos from the validation set of ActivityNet [10]. This dataset consists of captions matched with segments in the video, e.g., “a little boy is climbing on an outside gym” matched to the first 10 seconds of its related video. We extracted video segments with an associated caption. Human workers³ on Amazon MTurk were then shown the video segments and associated captions and were asked to create a semantically plausible contrast caption and a corresponding NLE (more details in Appendix §F). We did not communicate any type of target misalignments to encourage natural diversity of human created contrast captions.

Overall, we collected 570 tuples $(V, T, C_{human}, E_{human})$ where V is the video, T is the original caption, C_{human} is the human-written contrast caption, and E_{human} is the human-written explanations. We denote this dataset by VideoCon (Human). We sample 100 instances from this dataset, and found 93% to be clean. In addition, we observe that many of the human-generated contrast captions perturbing one or more objects (35%) and actions (35%) depicted in the caption. While 8% – 10% of the contrast captions flip the order of the events and attribute of the objects. As this dataset is largely unfiltered, it contains a mix of temporally-easy and challenging instances. We also constructed a more temporally-challenging subset of 290 instances, denoted VideoCon (Human-Hard), by filtering out tuples in which $A_{ve}(V, T) < 0.5$ (Eq. (1)), as in §3.1.

4. Experimental Setup

We next describe our evaluation setting for measuring the impact of VideoCon on video-text alignment modeling.

³A shortlist that passed our qualification test.

4.1. Finetuning with VideoCon

Our goal in constructing VideoCon (LLM) is to improve robustness of video-text alignment models by fine-tuning on this dataset. To this end, we start with the mPLUG-Owl-Video model [61], denoted *Owl-Base*. Its building blocks are CLIP [39] as visual encoder and LLaMA-7B [46] as text encoder/decoder and it was pretrained on VideoChat [27].

Entailment Task:
Given: V (Video), T (Caption), C (Contrast Caption)
Instruction (I): [V] Does this video entail the description [T]?
Response (R): Yes
Instruction (I): [V] Does this video entail the description [C]?
Response (R): No

Figure 4. Entailment task prompt for finetuning.

Natural Language Explanation Generation Task:
Given: V (Video), C (Contrast Caption), E (NLE)
Instruction (I): [V] What is the misalignment between this video and the description [C]?
Response (R): [E]

Figure 5. NLE generation task prompt for finetuning.

To fine-tune *Owl-Base* on VideoCon (LLM), its $\{V, T, C, E\}$ ⁴ tuples were converted into two types of multimodal instruction-response pairs, one for the VLE task (I_{vle}, R) (Fig. 4) and one for the NLE task (I_{nle}, R) (Fig. 5). We then train *Owl-Base* on all instruction pairs from both the tasks with maximum likelihood loss, resulting in a single model *Owl-Con*.

4.2. VideoCon Evaluation Metrics

To evaluate the performance of the *Owl-Con* on video-text alignment we generate *Owl-Con* response to prompt I_{vle} for video V and text $Y \in \{T, C\}$. We then calculate the probability of generating responses $s_y = \text{Owl-Con}(\text{'Yes'} | I_{vle}(V, Y))$ and $s_n = \text{Owl-Con}(\text{'No'} | I_{vle}(V, Y))$, and based on these scores the probability for class 'Yes': $P_{yes}(V, Y) = \frac{s_y}{s_y + s_n}$. Finally, we compute the ROC-AUC score for $P_{yes}(V, Y)$ over the VideoCon (LLM) eval set, with $\{V, T\}$ as label 1 and $\{V, C\}$ as label 0.

To evaluate *Owl-Con* on the NLE task, we prompt it with instruction I_{nle} instantiated on $\{V, C\}$ pairs from the VideoCon (LLM) eval set. We compare the generated explanation \hat{E} to the ground truth E by measuring entailment

⁴V: video, T: original caption, C: contrast caption, E: explanation.

probability $NLI(E, \hat{E})$. In our experiments, we experiment with two NLI automatic metrics: (a) Q^2 score [19], and (b) PaLM-2 API. We performed human evaluation to measure the agreement between the automatic metrics and the human-rating. We found that both metrics achieve high agreement with human assessment (Appendix §H).

4.3. Video-Text Downstream Tasks

We complement the VideoCon intrinsic evaluation over the testset with an extrinsic evaluation over two temporal and action difficult downstream tasks.

We evaluate alignment model performance for *text2video retrieval* over SSv2-Temporal [45] and SSv2-Events [5] datasets. We consider the SSv2-Template captions instead of the label captions since they remove the object-centric bias in model evaluation [26]. We compute input-text/candidate-video alignment score, rank videos and report *mean Average Precision* (mAP). We evaluate alignment model performance for *video question answering* over the ATP-Hard [9] dataset. We cast each question/candidate-answer pair as an imperative statement using PaLM-2 API, measure alignment to the input video and report *Accuracy*. More details on the downstream datasets and the evaluation setup are in Appendix §I.

4.4. Baselines

For the video-text alignment text, we compare *Owl-Con* with the following baselines: (a) End-to-End VNLI as zero-shot *atemporal* model since it does not have access to the temporal order of the video frames, (b) VideoCLIP [55], (c) ImageBind [14], (d) *Owl-Base*, and (e) *Owl-Rand*: *Owl-Base* fine-tuned on VideoCon tuples $\{V, T, \hat{C}, E\}$ where \hat{C} is randomly selected from other captions in the dataset. *Owl-Rand* would indicate if there is merit in the contrast, hard-negative captions in VideoCon. We include additional baselines TACT [5] and VFC [37] for evaluating on the downstream tasks (§5.3).

5. Experiments

We present our intrinsic (VideoCon eval set) and extrinsic (downstream tasks) evaluation results, showing the benefits of VideoCon for robust video-language alignment.

5.1. Performance on VideoCon Entailment Task

We present the ROC-AUC scores of the tested models in Table 2. From the table we see that the baseline models find the VideoCon testset difficult, as reflected by low AUC scores (e.g. *Owl-Base*- 57.2), close to random. Even training on VideoCon train instances, but with "easy" negatives (*Owl-Rand*- 59.7), hardly improves the base models. A significant improvement is achieved with the VNLI-specific model (67), showing that the entailment task is not inherently represented in generic video-language aligned training

Models	VideoCon (LLM) Test	VideoCon (Human)	VideoCon (Human-Hard)
Random	50.0	50.0	50.0
VideoCLIP [55]	53.2	47.3	47.5
ImageBind (Video-Text) [14]	57.1	65.2	63.0
<i>Owl-Base</i> [61]	57.2	66.8	64.1
<i>Owl-Rand</i>	59.7	68.9	65.5
End-to-End VNLI [60]	67.0	72.4	65.0
<i>Owl-Con (Ours)</i>	84.6	78.3	74.4

Table 2. ROC-AUC scores of the tested models for the entailment task on VideoCon test sets.

Models	VideoCon (LLM)		VideoCon (Human)	
	Q^2 entailment	PaLM-2 entailment acc. (%)	Q^2 entailment	PaLM-2 entailment acc.(%)
<i>Owl-Base</i>	0.19	36.8	0.23	39.6
<i>Owl-Con (Ours)</i>	0.50	65.4	0.32	47.1

Table 3. Performance of the tested models on the NLE generation task, measured via entailment metrics.

sets and requires specific training. Yet, the best performance is achieved by training on VideoCon, which addresses the diversity in plausible misalignments and includes “difficult” training examples, reaching 84.6 AUC. This demonstrates the merit of VideoCon for improving video-language alignment robustness. We show qualitative examples for the model predictions in §6.2.

When evaluating on out-of-domain (OOD) data around video types and misalignment distribution, we again see that training with VideoCon offers significant improvement to alignment detection, outperforming all baselines, albeit with smaller relative gains: 17% and 16% improvement compared to *Owl-Base* on (Human) and (Human-Hard) respectively compared to 48% on (LLM) test. In future work, we plan to further diversify the misalignments VideoCon covers to further improve its benefits on OOD cases.

We notice that the performance of the VNLI atemporal model is better than existing video-language alignment models. It might be attributed to its training with contrast captions in [60]. It further highlights that the existing video-language models are not robust in comparison to a atemporal probe on video-language alignment evaluation, corroborating the findings from [9, 26].

5.2. Performance on NLE Generation Task

Table 3 presents the performance of the tested models against the ground-truth on the NLE task, depicting average Q^2 score and PaLM-2 entailment accuracy. The results show that on in-domain VideoCon, *Owl-Con* outperforms *Owl-Base* by an impressive 263% and 178% relative increase on Q^2 score and PaLM-2 accuracy respectively. This indicates the finetuned model can accurately generate NLE that match well with the ground-truth NLE. This indicates that our model can generate accurate NLE for a wide range of misalignments in the video captions, which makes it use-

ful for dense video-language alignment evaluation.

On out-of-domain VideoCon, the improvement is more moderate but still high: 40% and 20% relative increase on Q^2 and PaLM-2 respectively. This is probably due to the more diverse ways humans express explanations compared to LLM prompting. In future work we plan to further address linguistic diversity in explanations for more robust generation and evaluation.

5.3. Performance on Video-Text Downstream Tasks

Models	SSv2-Temporal	SSv2-Events
	mAP	mAP
Random	7.3	3.3
VideoCLIP	9.8	6.4
ImageBind (video-language)	10.5	5.5
<i>Owl-Base</i>	10.9	6.8
TACT [5]	-	7.8
<i>Owl-Rand</i>	12.1	9.9
End-to-End VNLI [60]	14.6	10.4
<i>Owl-Con (Ours)</i>	15.2	11.4

Table 4. Mean Average Precision (mAP) scores for the tested models in the zero-shot text-to-video retrieval tasks.

We next present our results on the two downstream tasks, Text2Video Retrieval and Video Question Answering. Starting with the retrieval task, we report mean Average Precision (mAP) of the tested models on the SSv2-Temporal and SSv2-Events datasets in Table 4. The benefits of training with additional examples tailored for temporal video-language alignment is already evident in the performance of *Owl-Rand*, which improves over the previous SSv2-Events SOTA - TACT with a relative increase of 27%.

However, when training on harder negative contrastive instances, *Owl-Con* achieves a significant improvement,




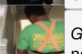

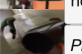






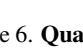

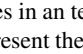
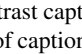
		Success		Failure		
VideoCon LLM		Caption: The adult gets the child's attention after tapping his shoulder	E		Caption: Two people move a large object down a narrow stairway	E
		Contrast Caption: The adult gets the child's attention before tapping	C		Contrast Caption: Two people move a large object up a narrow stairway	E
		GT NLE: First, an adult hand is seen tapping the child's shoulder, then the adult gets the child's attention, not the other way around	E		GT NLE: Men move a large piece of furniture down a narrow staircase, not up a narrow staircase	C
		Predicted NLE: First, the adult taps the child's shoulder, then the adult gets the child's attention, not the other way around	GT		Predicted NLE: Three men are moving a large piece of furniture, not two	GT
VideoCon Human		Caption: We see the group making cookies	E		Caption: A lady puts her hair in a bun on the side	E
		Contrast Caption: We see the group eating cookies	C		Contrast Caption: A lady puts her hair in a braid on the side	E
		GT NLE: The group is making cookies, not eating them	E		GT NLE: The woman puts her hair in a bun, not a braid	C
		Predicted NLE: We see the group cooking cookies, not eating them	GT		Predicted NLE: The lady puts her hair in a braid on top of her head, not on the side	GT

Figure 6. **Qualitative examples for the success (green) and failure (red) modes of our model.** In every example, we present a few video frames in an temporal order from top to bottom, its associated caption, contrast caption, ground-truth NLE from the datasets. Additionally, we present the predicted NLE from our model. The small boxes at the end of caption cells indicate whether our model consider that caption to be grounded in the video. **E** and **C** indicates that the model predicts the caption to entail and contradict to the video, respectively. **E-GT** and **C-GT** indicates the predicted NLE entails and contradicts the ground-truth (GT) NLE, respectively.

outperforming all baselines, with a relative increase over the best baseline End-to-End VNLI model by 7.5% on SSv2-Temporal and 9.6% on SSv2-Events (46% over TACT), setting new SOTA results. This points at the benefits of exposing the model to temporal examples, such as *action* and *event-order*.

Models	Accuracy (%)
CLIP	23.8
VideoCLIP	23.4
ImageBind (video-language)	25.4
TACT [5]	27.6
VFC [37]	31.4
<i>Owl-Base</i>	37.1
<i>Owl-Rand</i>	37.2
End-to-End VNLI [60]	39.0
<i>Owl-Con (Ours)</i>	41.1

Table 5. Accuracy scores for the tested models on the zero-shot video question-answering task on ATP-Hard dataset.

For the Video Question Answering task, we compare the performance of the various models in Table 5. Here too *Owl-Con* achieves SOTA results and outperforms the strongest baseline End-to-End VNLI model with a relative increase of 5.1%. This corroborates the observations in our other experiments, which demonstrate the advantage of the VideoCon datasets, covering various misalignments, especially those pertaining to temporal and causal reasoning over dynamic events. The results also confirm the need for carefully chosen contrastive negative examples, showing that picking negatives at random may mask out the potential benefit of an alignment training set. Finally, the competitive performance of atemporal End-to-End VNLI model on the downstream tasks is surprising and underscores the need for stronger video-language datasets for robust benchmarking.

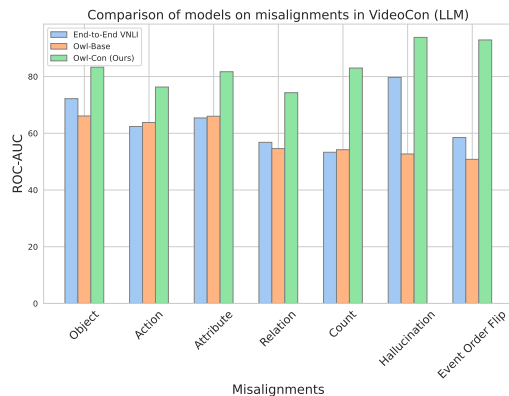


Figure 7. ROC-AUC of End-to-End VNLI, *Owl-Base*, and *Owl-Con* across all types of misalignment in VideoCon (LLM) test set.

6. Analysis

We analyze *Owl-Con*'s performance improvements across the kinds of misalignments in VideoCon. Additionally, we present a few qualitative examples to highlight the success and failure modes of our model.

6.1. Per-misalignment Entailment Results

We compared the ROC-AUC scores of the atemporal End-to-End VNLI, *Owl-Base*, and *Owl-Con* on specific misalignments in the contrast captions from VideoCon (LLM) testset in Figure 7. We observed that *Owl-Con* outperforms the baseline models across all misalignment types. This suggests that our model can reason well about the entities, their relations, and the temporal order of events in the video.

The largest improvement of *Owl-Con* compared to the two baselines is on *event order flip*, indicating that the baselines lack temporal understanding and the VideoCon is efficient in adding this capability to an alignment model. In

addition, on *hallucination* both *Owl-Con* and End-to-End VNLI significantly outperform *Owl-Base*, since both models were explicitly exposed to entailment/non-entailment training data. It is surprising to see that while End-to-End VNLI was trained on significantly more entailment data, much of it human-curated, *Owl-Con* outperforms it with only automatically generated data. This could be due to the better encoding of video in *Owl-Con* compared to the atemporal nature of End-to-End VNLI. Finally, the analysis shows other types of atemporal misalignments that are difficult for End-to-End VNLI to sort out, e.g. *counting* and *relation*, where the training data in VideoCon is useful to improve these capabilities as well. This shows that our approach of detailed analysis of misalignment types of generation of examples for them is effective.

6.2. Qualitative Examples

We highlight a few classification examples of *Owl-Con* in Figure 6. The rows refer to the test source of the instances and the columns refer to the success and failure modes, respectively. In Row1/Column1, we observe that our model provides correct predictions for the entailment between the video and original caption while predicting contradiction for the contrast caption that flips the order of the events i.e., *grabbing attention* and *tapping shoulders*. Interestingly, our model can also provide the accurate NLE when prompted with the video and the contrast caption. This suggests that our model is useful for providing fine-grained details about the video-language alignment. In Row2/Column2, the model confuses ‘buns’ with ‘braids’ in hair and gives a wrong NLE that contradicts the ground-truth. This error, due to its inability to distinguish between objects, might be improved with diverse videos and captions.

7. Related Work

Foundation Models for Video-Language Understanding. Foundation models have emerged for video-language understanding [1, 4, 49, 55, 56] by pre-training on large amount of video-text pairs scraped from the web [6, 36, 58]. Additionally, prior works have either leveraged the pre-trained CLIP model for video-language tasks [12, 33, 34] or adopted a socratic approach [50, 63] to employ LLMs (GPT-3) in reasoning over video captions. We highlight that despite the large-scale training of the video-language foundation models [14, 55, 56], they lack robustness to semantic changes to the captions which severely limits their real-world use for alignment applications. We provide a fix by training models on a novel video-centric VideoCon dataset.

Improving Video-Language Robustness. Prior work [37, 38, 51] highlights that the video-text models cannot comprehend the semantics of the text with focus on manipulating the verb, actions, and entities grounded in the video

description. To improve the temporal understanding, [5] finetunes a pretrained model with temporal order loss. Despite this, their models do not achieve good zero-shot performance on downstream tasks consistently. In our work, we categorize a wide range of plausible misalignments in the contrast captions, and create a temporally-challenging VideoCon dataset.

Video-Language Alignment Evaluation. Many applications such as text-to-video retrieval [15, 48, 57] and text-to-video generation [7, 47] require evaluation of the semantic alignment between the natural language text and raw video. In this work, we indicate that the existing video-text models such as VideoCLIP and ImageBind are not robust to semantic changes in the video captions, which becomes critical for faithful video-text alignment evaluation. In our work, we propose VideoCon and finetune a video-language generative model to perform robust entailment task and provide fine-grained NLE for the observed misalignments between the video and text. In the future, our model can be utilized to enhance alignment through sparse (entailment scores) and dense (fine-grained NLE) feedback [43].

8. Conclusion

We introduced a comprehensive dataset, VideoCon, designed for robust video-text alignment. It features various semantic misalignments and explanations for text-video discrepancies. Through finetuning video-language models on this dataset, we enhanced their performance on complex tasks like text-to-video retrieval and video question answering, achieving state-of-the-art results.

One current limitation and an important future direction is to increase the complexity of the generated contrast captions. Specifically, the model may encounter several misalignments within a single contrast caption. Addressing this issue, the model should be equipped to accurately assign low entailment scores to these contrast captions and consequently generate precise NLEs. An important future direction is to scale VideoCon to larger datasets. Here, we create contrast captions for high-quality captions written by humans for every video, however, the web-scale datasets have low-quality captions that are not well grounded in the video. In this regard, using synthetic data followed by VideoCon-like contrast caption generation can be a plausible approach.

9. Acknowledgement

This material is based on research supported by the ECOLE program under Cooperative Agreement HR00112390060 with the US Defense Advanced Research Projects Agency (DARPA). In addition, the research is partly funded by ONR grant N00014-23-1-2780. Hritik Bansal is supported in part by AFOSR MURI grant FA9550-22-1-0380.

References

- [1] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. *Advances in Neural Information Processing Systems*, 35:23716–23736, 2022. 1, 8
- [2] Rohan Anil, Andrew M Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, et al. Palm 2 technical report. *arXiv preprint arXiv:2305.10403*, 2023. 4
- [3] Lisa Anne Hendricks, Oliver Wang, Eli Shechtman, Josef Sivic, Trevor Darrell, and Bryan Russell. Localizing moments in video with natural language. In *Proceedings of the IEEE international conference on computer vision*, pages 5803–5812, 2017. 2
- [4] Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lučić, and Cordelia Schmid. Vivit: A video vision transformer. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6836–6846, 2021. 8, 1
- [5] Piyush Bagad, Makarand Tapaswi, and Cees GM Snoek. Test of time: Instilling video-language models with a sense of time. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2503–2516, 2023. 1, 2, 3, 5, 6, 7, 8
- [6] Max Bain, Arsha Nagrani, Gül Varol, and Andrew Zisserman. Frozen in time: A joint video and image encoder for end-to-end retrieval. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1728–1738, 2021. 8, 1
- [7] Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler, and Karsten Kreis. Align your latents: High-resolution video synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22563–22575, 2023. 2, 8, 1
- [8] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020. 1
- [9] Shyamal Buch, Cristóbal Eyzaguirre, Adrien Gaidon, Jiajun Wu, Li Fei-Fei, and Juan Carlos Niebles. Revisiting the” video” in video-language understanding. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2917–2927, 2022. 1, 5, 6
- [10] Fabian Caba Heilbron, Victor Escorcia, Bernard Ghanem, and Juan Carlos Niebles. Activitynet: A large-scale video benchmark for human activity understanding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 961–970, 2015. 4, 3
- [11] Alex Fang, Albin Madappally Jose, Amit Jain, Ludwig Schmidt, Alexander Toshev, and Vaishaal Shankar. Data filtering networks. *arXiv preprint arXiv:2309.17425*, 2023. 1
- [12] Han Fang, Pengfei Xiong, Luhui Xu, and Yu Chen. Clip2video: Mastering video-text retrieval via image clip. *arXiv preprint arXiv:2106.11097*, 2021. 8, 1
- [13] Samir Yitzhak Gadre, Gabriel Ilharco, Alex Fang, Jonathan Hayase, Georgios Smyrnis, Thao Nguyen, Ryan Marten, Mitchell Wortsman, Dhruva Ghosh, Jieyu Zhang, et al. Datacomp: In search of the next generation of multimodal datasets. *arXiv preprint arXiv:2304.14108*, 2023. 1
- [14] Rohit Girdhar, Alaaeldin El-Nouby, Zhuang Liu, Mannat Singh, Kalyan Vasudev Alwala, Armand Joulin, and Ishan Misra. Imagebind: One embedding space to bind them all. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15180–15190, 2023. 1, 2, 5, 6, 8
- [15] Raghav Goyal, Samira Ebrahimi Kahou, Vincent Michalski, Joanna Materzynska, Susanne Westphal, Heuna Kim, Valentin Haenel, Ingo Fruend, Peter Yianilos, Moritz Mueller-Freitag, et al. The” something something” video database for learning and evaluating visual common sense. In *Proceedings of the IEEE international conference on computer vision*, pages 5842–5850, 2017. 8, 1, 5
- [16] Suriya Gunasekar, Yi Zhang, Jyoti Aneja, Caio César Teodoro Mendes, Allie Del Giorno, Sivakanth Gopi, Mojan Javaheripi, Piero Kauffmann, Gustavo de Rosa, Olli Saarikivi, et al. Textbooks are all you need. *arXiv preprint arXiv:2306.11644*, 2023. 1
- [17] Lisa Anne Hendricks, Oliver Wang, Eli Shechtman, Josef Sivic, Trevor Darrell, and Bryan Russell. Localizing moments in video with temporal language. *arXiv preprint arXiv:1809.01337*, 2018. 1, 3, 2
- [18] Matthew Honnibal and Ines Montani. spaCy 2: Natural language understanding with Bloom embeddings, convolutional neural networks and incremental parsing. To appear, 2017. 2
- [19] Or Honovich, Leshem Choshen, Roei Aharoni, Ella Neeman, Idan Szpektor, and Omri Abend. Q2: Evaluating factual consistency in knowledge-grounded dialogues via question generation and question answering. *arXiv preprint arXiv:2104.08202*, 2021. 4, 5
- [20] <https://commoncrawl.org/>. 1
- [21] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021. 3
- [22] Yushi Hu, Benlin Liu, Jungo Kasai, Yizhong Wang, Mari Ostendorf, Ranjay Krishna, and Noah A Smith. Tifa: Accurate and interpretable text-to-image faithfulness evaluation with question answering. *arXiv preprint arXiv:2303.11897*, 2023. 1
- [23] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *International conference on machine learning*, pages 4904–4916. PMLR, 2021. 1
- [24] Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan, Fabio Viola, Tim Green, Trevor Back, Paul Natsev, et al. The kinetics human action video dataset. *arXiv preprint arXiv:1705.06950*, 2017. 1

- [25] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. **3**
- [26] Jie Lei, Tamara L Berg, and Mohit Bansal. Revealing single frame bias for video-and-language learning. *arXiv preprint arXiv:2206.03428*, 2022. **1, 5, 6**
- [27] KunChang Li, Yinan He, Yi Wang, Yizhuo Li, Wenhai Wang, Ping Luo, Yali Wang, Limin Wang, and Yu Qiao. Videochat: Chat-centric video understanding. *arXiv preprint arXiv:2305.06355*, 2023. **5**
- [28] Liunian Harold Li, Zi-Yi Dou, Nanyun Peng, and Kai-Wei Chang. Desco: Learning object recognition with rich language descriptions. *arXiv preprint arXiv:2306.14060*, 2023. **1**
- [29] Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. Evaluating object hallucination in large vision-language models. *arXiv preprint arXiv:2305.10355*, 2023. **3**
- [30] Weixin Liang, James Zou, and Zhou Yu. Alice: Active learning with contrastive natural language explanations. *arXiv preprint arXiv:2009.10259*, 2020. **1**
- [31] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36, 2024. **6**
- [32] Jiaying Lu, Jinneng Rao, Kezhen Chen, Xiaoyuan Guo, Yawen Zhang, Baochen Sun, Carl Yang, and Jie Yang. Evaluation and mitigation of agnosia in multimodal large language models. *arXiv preprint arXiv:2309.04041*, 2023. **3**
- [33] Huaishao Luo, Lei Ji, Ming Zhong, Yang Chen, Wen Lei, Nan Duan, and Tianrui Li. Clip4clip: An empirical study of clip for end to end video clip retrieval and captioning. *Neurocomputing*, 508:293–304, 2022. **2, 8, 1**
- [34] Yiwei Ma, Guohai Xu, Xiaoshuai Sun, Ming Yan, Ji Zhang, and Rongrong Ji. X-clip: End-to-end multi-grained contrastive learning for video-text retrieval. In *Proceedings of the 30th ACM International Conference on Multimedia*, pages 638–647, 2022. **8, 1**
- [35] Zixian Ma, Jerry Hong, Mustafa Omer Gul, Mona Gandhi, Irena Gao, and Ranjay Krishna. Crepe: Can vision-language foundation models reason compositionally? In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10910–10921, 2023. **3**
- [36] Antoine Miech, Dimitri Zhukov, Jean-Baptiste Alayrac, Makarand Tapaswi, Ivan Laptev, and Josef Sivic. Howto100m: Learning a text-video embedding by watching hundred million narrated video clips. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 2630–2640, 2019. **8, 1**
- [37] Liliane Momeni, Mathilde Caron, Arsha Nagrani, Andrew Zisserman, and Cordelia Schmid. Verbs in action: Improving verb understanding in video-language models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15579–15591, 2023. **1, 2, 3, 5, 7, 8**
- [38] Jae Sung Park, Sheng Shen, Ali Farhadi, Trevor Darrell, Yejin Choi, and Anna Rohrbach. Exposing the limits of video-text models through contrast sets. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3574–3586, 2022. **1, 2, 3, 8**
- [39] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. **2, 5**
- [40] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International Conference on Machine Learning*, pages 8821–8831. PMLR, 2021. **1**
- [41] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. **1**
- [42] Fawaz Sammani, Tanmoy Mukherjee, and Nikos Deligiannis. Nlx-gpt: A model for natural language explanations in vision and vision-language tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8322–8332, 2022. **1**
- [43] Jérémy Scheurer, Jon Ander Campos, Tomasz Korbak, Jun Shern Chan, Angelica Chen, Kyunghyun Cho, and Ethan Perez. Training language models with language feedback at scale. *arXiv preprint arXiv:2303.16755*, 2023. **8, 1**
- [44] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. *Advances in Neural Information Processing Systems*, 35:25278–25294, 2022. **1**
- [45] Laura Sevilla-Lara, Shengxin Zha, Zhicheng Yan, Vedanuj Goswami, Matt Feiszli, and Lorenzo Torresani. Only time can tell: Discovering temporal data for temporal modeling. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 535–544, 2021. **2, 5**
- [46] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023. **5**
- [47] Jiuniu Wang, Hangjie Yuan, Dayou Chen, Yingya Zhang, Xiang Wang, and Shiwei Zhang. Modelscope text-to-video technical report. *arXiv preprint arXiv:2308.06571*, 2023. **2, 8, 1**
- [48] Xin Wang, Jiawei Wu, Junkun Chen, Lei Li, Yuan-Fang Wang, and William Yang Wang. Vatex: A large-scale, high-quality multilingual dataset for video-and-language research. In *The IEEE International Conference on Computer Vision (ICCV)*, 2019. **3, 8, 1**
- [49] Yi Wang, Kunchang Li, Yizhuo Li, Yinan He, Bingkun Huang, Zhiyu Zhao, Hongjie Zhang, Jilan Xu, Yi Liu, Zun Wang, et al. Internvideo: General video foundation models via generative and discriminative learning. *arXiv preprint arXiv:2212.03191*, 2022. **2, 8, 1**

- [50] Zhenhailong Wang, Manling Li, Ruochen Xu, Luowei Zhou, Jie Lei, Xudong Lin, Shuohang Wang, Ziyi Yang, Chenguang Zhu, Derek Hoiem, et al. Language models with image descriptors are strong few-shot video-language learners. *Advances in Neural Information Processing Systems*, 35: 8483–8497, 2022. [8](#), [1](#)
- [51] Zhenhailong Wang, Ansel Blume, Sha Li, Genglin Liu, Jaemin Cho, Zineng Tang, Mohit Bansal, and Heng Ji. Paxion: Patching action knowledge in video-language foundation models. *arXiv preprint arXiv:2305.10683*, 2023. [1](#), [2](#), [8](#)
- [52] Junbin Xiao, Xindi Shang, Angela Yao, and Tat-Seng Chua. Next-qa: Next phase of question-answering to explaining temporal actions. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9777–9786, 2021. [5](#)
- [53] Ning Xie, Farley Lai, Derek Doran, and Asim Kadav. Visual entailment task for visually-grounded language learning. *arXiv preprint arXiv:1811.10582*, 2018. [2](#)
- [54] Ning Xie, Farley Lai, Derek Doran, and Asim Kadav. Visual entailment: A novel task for fine-grained image understanding. *arXiv preprint arXiv:1901.06706*, 2019. [2](#)
- [55] Hu Xu, Gargi Ghosh, Po-Yao Huang, Dmytro Okhonko, Armen Aghajanyan, Florian Metze, Luke Zettlemoyer, and Christoph Feichtenhofer. Videoclip: Contrastive pre-training for zero-shot video-text understanding. *arXiv preprint arXiv:2109.14084*, 2021. [1](#), [2](#), [5](#), [6](#), [8](#)
- [56] Haiyang Xu, Qinghao Ye, Ming Yan, Yaya Shi, Jiabo Ye, Yuanhong Xu, Chenliang Li, Bin Bi, Qi Qian, Wei Wang, et al. mplug-2: A modularized multi-modal foundation model across text, image and video. *arXiv preprint arXiv:2302.00402*, 2023. [8](#), [1](#)
- [57] Jun Xu, Tao Mei, Ting Yao, and Yong Rui. Msr-vtt: A large video description dataset for bridging video and language. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5288–5296, 2016. [2](#), [3](#), [8](#), [1](#)
- [58] Hongwei Xue, Tiankai Hang, Yanhong Zeng, Yuchong Sun, Bei Liu, Huan Yang, Jianlong Fu, and Baining Guo. Advancing high-resolution video-language representation with large-scale video transcriptions. In *International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. [8](#), [1](#)
- [59] Antoine Yang, Arsha Nagrani, Paul Hongsuck Seo, Antoine Miech, Jordi Pont-Tuset, Ivan Laptev, Josef Sivic, and Cordelia Schmid. Vid2seq: Large-scale pretraining of a visual language model for dense video captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10714–10726, 2023. [2](#)
- [60] Michal Yarom, Yonatan Bitton, Soravit Changpinyo, Roei Aharoni, Jonathan Herzig, Oran Lang, Eran Ofek, and Idan Szpektor. What you see is what you read? improving text-image alignment evaluation. *arXiv preprint arXiv:2305.10400*, 2023. [3](#), [6](#), [7](#), [1](#), [2](#)
- [61] Qinghao Ye, Haiyang Xu, Guohai Xu, Jiabo Ye, Ming Yan, Yiyang Zhou, Junyang Wang, Anwen Hu, Pengcheng Shi, Yaya Shi, et al. mplug-owl: Modularization empowers large language models with multimodality. *arXiv preprint arXiv:2304.14178*, 2023. [1](#), [5](#), [6](#)
- [62] Mert Yuksekgonul, Federico Bianchi, Pratyusha Kalluri, Dan Jurafsky, and James Zou. When and why vision-language models behave like bags-of-words, and what to do about it? In *The Eleventh International Conference on Learning Representations*, 2022. [1](#)
- [63] Andy Zeng, Maria Attarian, Brian Ichter, Krzysztof Choromanski, Adrian Wong, Stefan Welker, Federico Tombari, Aavek Purohit, Michael Ryoo, Vikas Sindhwani, et al. Socratic models: Composing zero-shot multimodal reasoning with language. *arXiv preprint arXiv:2204.00598*, 2022. [8](#), [1](#)
- [64] Luowei Zhou, Chenliang Xu, and Jason J Corso. Towards automatic learning of procedures from web instructional videos. In *AAAI Conference on Artificial Intelligence*, pages 7590–7598, 2018. [7](#)