

An Empirical Study of End-to-End Video-Language Transformers with Masked Visual Modeling

Tsu-Jui Fu^{†*}, Linjie Li^{‡*}, Zhe Gan[‡], Kevin Lin[‡], William Yang Wang[†], Lijuan Wang[‡], Zicheng Liu[‡]

[†]UC Santa Barbara [‡]Microsoft

{tsu-juifu, william}@cs.ucsb.edu

{lindsey.li, zhe.gan, keli, lijuanw, zliu}@microsoft.com

Abstract

Masked visual modeling (MVM) has been recently proven effective for visual pre-training. While similar reconstructive objectives on video inputs (e.g., masked frame modeling) have been explored in video-language (VidL) pre-training, the pre-extracted video features in previous studies cannot be refined through MVM during pre-training, and thus leading to unsatisfactory downstream performance. In this work, we systematically examine the potential of MVM in the context of VidL learning. Specifically, we base our study on a fully end-to-end Video-Language Transformer (VIOLET), which mitigates the disconnection between fixed video representations and MVM training. In total, eight different reconstructive targets of MVM are explored, from low-level pixel values and oriented gradients to high-level depth maps, optical flow, discrete visual tokens and latent visual features. We conduct comprehensive experiments and provide insights on the factors leading to effective MVM training. Empirically, we show VIOLET pre-trained with MVM objective achieves notable improvements on 13 VidL benchmarks, ranging from video question answering, video captioning, to text-to-video retrieval.

1. Introduction

Video, containing multiple modalities in nature, has been used as an epitome to test how AI systems perceive. Video-language (VidL) research aims at extending this ability to convey perception via language. Popular VidL tasks were introduced, such as text-to-video retrieval [80, 30, 59], video question answering [26, 79], and video captioning [80, 7]. Recent progresses in VidL learning mostly focus on VidL pre-training [63, 52, 88] with video-text matching [40, 84] and masked language modeling [12]. There have also been attempts on similar masked modeling on vision inputs. For example, masked frame modeling [40]

aims to recover masked frame representations. However, the pre-extracted video features cannot be refined during pre-training, which may limit its effectiveness.

Meanwhile, self-supervised vision pre-training has been proven highly effective by reconstructing the masked image patches through raw pixel values [22, 78], discrete visual tokens [4, 86], or visual-semantic features [75, 76]. However, they all only focus on the visual modality. It is unknown how masked visual modeling (MVM) objectives can help VidL learning, especially given that the paired language inputs can already provide high-level semantics.

Motivated by this, we conduct a comprehensive study of MVM for VidL learning. As Figure 1, we base our study on a fully end-to-end Video-Language Transformer (named VIOLET [17]), and study a broad spectrum of MVM targets, including RGB pixel values (Pixel), histogram of oriented gradients (HOG), depth maps (Depth), optical flow (Flow), discrete visual tokens (VQ), spatial-focused image features (SIF), temporal-aware video features (TVF), and multimodal features (MMF). During pre-training, we mask out some proportions of the video input along both spatial and temporal dimensions, and the model learns to recover the MVM targets for these masked patches. Equipped with another two standard pre-training tasks (*i.e.*, video-text matching (VTM) and masked language modeling (MLM)), we empirically verify the effectiveness of different MVM variants on downstream VidL tasks.

Our study reveals that: (i) spatial-focused image features (SIF) is the most effective MVM target on video-text inputs; and (ii) the effects of different MVM targets on downstream VidL tasks are not shared between video-text and image-text inputs. For example, SIF extracted from the same model brings a large drop on downstream VidL performance when pre-trained with image-text pairs. In addition, we conduct comprehensive analyses of the masking strategy and ratio, combination of different MVM targets, to shed light on effective MVM training for VidL learning.

In summary, our contributions are three-fold. (i) We present an empirical study of masked visual modeling for

* means equal contribution.

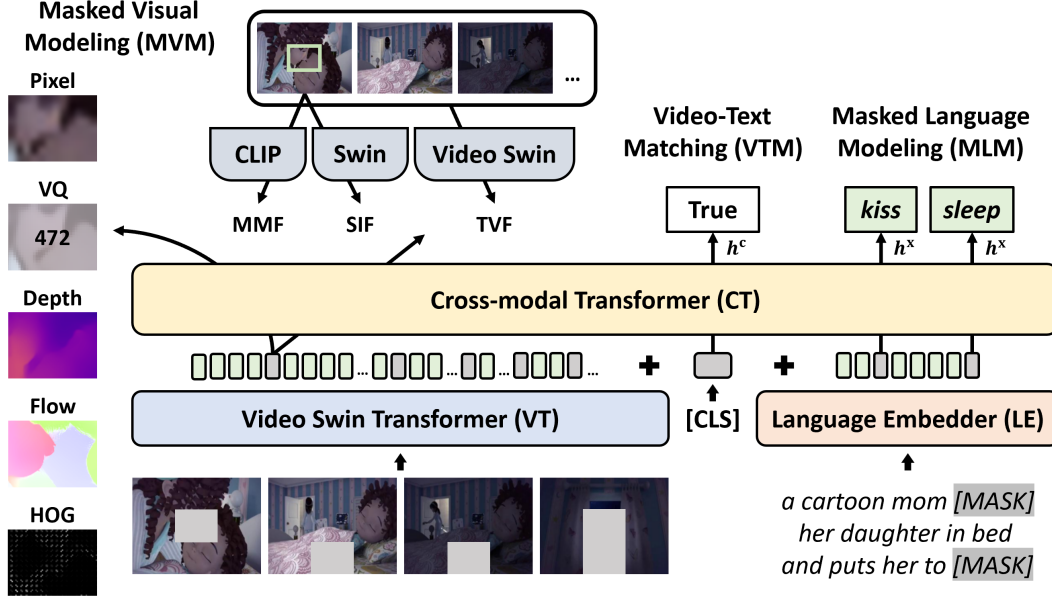


Figure 1: We systematically explore *eight* masked visual modeling (MVM) targets for end-to-end video-language (VidL) pre-training, including RGB pixel values (Pixel), histogram of oriented gradients (HOG), depth maps (Depth), optical flow (Flow), discrete visual tokens (VQ), spatial-focused image features (SIF), temporal-aware video features (TVF), and multi-modal features from CLIP (MMF). Besides MVM, the proposed VIOLET model is pre-trained along with video-text matching (VTM) and masked language modeling (MLM).

video-language pre-training; (ii) We conduct comprehensive analyses with extensive experimental results to shed lights on effective MVM training; and (iii) VIOLET pre-trained with MVM objective achieves strong performance on 13 VidL datasets over 3 popular tasks, covering video question answering, video captioning and text-to-video retrieval. Concretely, compared to models trained on the same 5M pre-training dataset, VIOLET with effective MVM pre-training brings notable mean improvements of +5.4% accuracy on video question answering, +6.6% recall on text-to-video retrieval, and +11.4 CIDEr on video captioning.

2. Related Work

Video-Language Understanding. Joint video-language (VidL) understanding [42, 45, 27, 33, 18, 54] aims at interpreting the physical world via both vision and text perception. Researchers have explored such capability on VidL tasks including text-to-video retrieval [80, 30, 59, 38, 40], video question answering [26, 79, 36, 37], moment retrieval [24, 20, 30, 38], and video captioning [74, 87, 80, 59]. Prior arts before the large-scale pre-training era [19, 85, 34, 15, 33, 37] leverage offline extracted video features [28, 72, 6, 77, 16, 11, 23, 31, 2]. Later on, VidL pre-trained models [63, 88, 40, 52] built on the above pre-extracted features have shown promising results. To enhance the performance, there have been parallel interests in bringing in more modalities from raw video inputs [18, 60, 44] and end-to-end training [51, 35, 84, 3],

aiming to elevate video representations for VidL modeling.

Masked Visual Modeling (MVM). Aligned with the success of transformer-based [70] language pre-training [32, 46], image-text pre-training [9, 64] and video-text pre-training [29, 82, 81] have shown promising results on diverse vision-language (VL) tasks. Popular VL pre-training tasks include visual-text matching (VTM) and masked language modeling (MLM), which are directly adapted from language pre-training [12]. Similar masked modeling on visual inputs [9, 40, 14] has also been introduced to VL pre-training, but are not as useful. Among the literature of vision pre-training itself, MAE [22, 67] and SimMIM [78] reconstruct the pixels of the masked image patches to enhance visual representation. BEiT [4], iBOT [86], VIM-PAC [65], and BEVT [73] adopt a BERT-like pre-training strategy to recover the missing visual tokens. On the other hand, MaskFeat [75] and MVP [76] consider latent features for MVM, including hand-crafted HOG features and image features extracted from pre-trained CLIP models [55]. Unlike previous works exploring MVM on uni-modal data, in this study, we conduct a comprehensive investigation on how different MVM targets can help VidL learning.

3. Method

We first describe the problem formulation in Section 3.1, and then detail the overall framework of the proposed VIOLET model in Section 3.2. Finally, we discuss eight differ-

ent target features considered for masked visual modeling (MVM) in Section 3.3.

3.1. Problem Setting

Given a large-scale video-language (VidL) dataset D , we aim to pre-train a VidL transformer to learn effective video-text representations. The learned representations can be transferred to downstream tasks for performance improvement. Different from existing works that focus on MVM for pure vision problems [4, 22, 86], we study MVM as a VidL pre-training task. Given a video-text pair $(\mathcal{V}, \mathcal{X})$ where \mathcal{V} is a sequence of video frames and \mathcal{X} is a sequence of word tokens. As shown in Figure 1, we randomly mask out some portions of the input frames \mathcal{V} , and learn to predict the target features corresponding to the masked patches. To output a correct prediction, the model will have to resort to other relevant video frames \mathcal{V} and/or text tokens \mathcal{X} . This facilitates cross-modality learning for better VidL understanding.

In addition, we employ the commonly used VidL pre-training objectives, including video-text matching (VTM) and masked language modeling (MLM), where VTM aims to predict whether an input video-text pair is matched or not, while MLM aims to predict the masked word tokens from the surrounding context.¹ Our overall pre-training objective can be written as:

$$\mathcal{L} = \mathcal{L}_{\text{MVM}} + \mathcal{L}_{\text{VTM}} + \mathcal{L}_{\text{MLM}}, \quad (1)$$

where \mathcal{L}_{MVM} , \mathcal{L}_{VTM} , \mathcal{L}_{MLM} are the MVM, VTM and MLM objectives, respectively.

3.2. End-to-End Video-Language Transformer

We conduct our empirical study using an end-to-end Video-Language Transformer (VIOLET [17]). As Figure 1, VIOLET contains 3 components: Video Swin Transformer (VT), Language Embedder (LE), and Cross-modal Transformer (CT). VIOLET takes video \mathcal{V} and sentence \mathcal{X} as inputs. Sparse-sampled frames $\{f_1, f_2, \dots\}$ from \mathcal{V} are first segmented into a set of video patches, and then processed by VT to compute video features $v = \{v_1, v_2, \dots\}$. LE extracts the word embeddings $w = \{w_1, w_2, \dots\}$ for each word token $\{x_1, x_2, \dots\}$ in \mathcal{X} . Then, CT performs cross-modal fusion on top of v and w to produce joint VidL representations $h = [h^v, h^c, h^x]$, where h^v, h^c, h^x denote the hidden representations of video patches, the special [CLS] token, and other word tokens. We pre-train VIOLET in an end-to-end manner with all three objectives in Equation 1 on top of the outputs h from CT.

3.3. Target Features

Masked visual modeling (MVM) is a generic masked feature prediction task, where we mask out some of the vi-

sual input patches, and then predict the target features corresponding to the masked ones. Thus, a core design of MVM is the target features, which enables VIOLET learning a desired aspect of visual modeling. While MVM has been explored in pure vision tasks [4, 22, 75], it remains an open question whether MVM can facilitate the interactions between video and language modalities. In this study, we investigate *what design of MVM is effective in the context of video-language pre-training?*

Following [78, 75], we employ a simple linear layer or 2-layer MLP as the prediction head for MVM, to project the hidden video representations (h^v , of hidden size 768) from CT to the same dimension as the MVM targets. The default MVM loss is the l_1 loss, unless specified otherwise. Next, we introduce the considered target features in details.

RGB Pixel Values (Pixel). We treat the normalized RGB pixel values as one of the candidate target features. During MVM, VIOLET learns to reconstruct the pixel values of the masked patches. The linear MVM head projects h^v into the same dimension as the raw video frame patch ($H \times W \times 3$).

Histogram of Oriented Gradients (HOG). HOG [10] is a pioneer feature descriptor that describes the gradients of orientations of the image. While HOG has been proven effective for visual pre-training [75], it is unknown whether it can benefit VidL pre-training. We extract HOG features in a dense grid level, and use such feature descriptors as the prediction targets of MVM. The HOG feature map is of the same size as the input video frame, but with channel size 1. The linear MVM prediction head projects h^v to the same dimension as HOG for the video frame patch ($H \times W \times 1$).

Depth Maps (Depth). Since depth maps usually contains finer-grained details of the object shapes and general scene layout of the foreground objects, it is worth exploring whether depth maps can be used to improve the scene/object understanding capability of a VidL pre-trained model. To obtain such MVM target, we employ a pre-trained dense prediction transformer (DPT) [57] to perform monocular depth estimation given an input video frame. The linear prediction head used for Depth is the same as the one for HOG, as both targets are of channel size 1.

Optical Flow (Flow). Optical flow is commonly used in motion analysis and video understanding. Here, we analyze whether apparent velocity of objects can benefit VidL pre-training. We employ a pre-trained recurrent all-pairs field transforms (RAFT) [66] to compute optical flow given the consecutive video frames. We directly use the estimated optical flow values as the prediction target, and supervise the MVM training with l_1 loss. To obtain the MVM predictions, we concatenate the hidden video representations computed by CT on consecutive frames, and employ a linear layer to project the concatenated video representations

¹Refer to the Appendix for detailed formulation of VTM and MLM.

Pre-training Tasks	MVM Target	TGIF-Frame	DiDeMo-Retrieval			
		Acc.	R1	R5	R10	AveR
VTM+MLM	None	68.1	28.7	57.0	69.7	51.8
+MVM	RGB Pixel Values	68.3 (+0.2)	29.2 (+0.5)	58.6 (+1.6)	70.1 (+0.4)	52.6 (+0.8)
	Histogram of Oriented Gradients [10]	67.3 (-0.8)	26.6 (-2.1)	54.9 (-2.1)	68.1 (-1.6)	49.8 (-2.0)
	Depth Maps (DPT-L [57])	68.0 (-0.1)	27.3 (-1.4)	55.0 (-2.0)	68.3 (-1.4)	50.2 (-1.6)
	Optical Flow (RAFT-L [66])	67.6 (-0.5)	30.3 (+1.6)	58.0 (+1.0)	70.3 (+0.3)	52.9 (+1.1)
	Spatial-focused Image Features (Swin-B [47])	68.8 (+0.7)	35.4 (+6.7)	62.4 (+5.2)	74.9 (+6.3)	57.6 (+5.8)
	Temporal-aware Video Features (VidSwin-L [48])	68.0 (-0.1)	32.8 (+4.1)	60.5 (+3.5)	73.0 (+3.3)	55.4 (+3.6)
	Discrete Visual Tokens (DALL-E [56])	68.4 (+0.3)	28.1 (-0.6)	56.6 (-0.4)	69.4 (-0.5)	51.3 (-0.5)
	Multimodal Features (CLIP-ViT-B [55])	67.7 (-0.4)	29.8 (+1.1)	57.8 (+0.8)	68.5 (-1.2)	52.1 (+0.3)

Table 1: **Comparing target features for MVM applied to video-text data.** All variants are pre-trained on WebVid [3] for 5 epochs. Masking is performed randomly (RM) with ratio of 15%. The final pre-training setting is highlighted in gray.

(of hidden size 768×2) to the same dimension as the estimated optical flow target for a given patch ($H \times W \times 2$).

Discrete Visual Tokens (VQ). In addition to continuous MVM targets, we also consider the discrete variational auto-encoder (dVAE) [69, 56] to quantize video inputs. dVAE is learned to tokenize images into discrete visual tokens q from a finite dictionary, and then reconstruct the original visual scene based on q , where q should have a one-to-one correspondence with the input image patches spatially. We first adopt dVAE to tokenize the t^{th} video frame f_t into q_t : $q_t = \text{dVAE}(f_t)$, and then a 2-layer MLP is used to project h_v into the finite VQ vocabularies. As VQ token is discrete, we can model MVM with VQ as a classification problem, and adopt the cross-entropy loss to optimize the MVM training, following [4].

Spatial-focused Image Features (SIF). We investigate whether image features can be useful for improving VidL pre-training. We employ a well-known vision transformer (such as Swin Transformer [47]) to extract the grid features given an input image. We then normalize the extracted grid features and consider them as ground-truth MVM targets. Likewise, we adopt a 2-layer MLP to project h_v to the same dimension as the image feature target.

Temporal-aware Video Features (TVF). We also study the impact of video features to VidL pre-training. We employ pre-trained video transformer (such as Video Swin Transformer [48]) to compute temporal-aware features for this analysis. Given a set of video frames, we use the transformer to extract video features in the form of space-time cubes, and then apply l_1 regression between normalized video features and MVM predictions from a 2-layer MLP head of the masked video patches.

Multimodal Features (MMF). We further study if the features learned via multimodal pre-training can benefit VidL pre-training. We utilize the vision branch of the ViT-Base backbone [13] in CLIP [55] to extract such multimodal features, and use the normalized features as the prediction targets in MVM pre-training. Again, we apply l_1 regression

between the MVM predictions made via a 2-layer MLP head and the MMF targets for the masked patches.

4. Study: Target Features for MVM

Settings. We pre-train VIOLET on WebVid-2.5M [3] for 5 epochs, and report accuracy on TGIF-Frame [26] for video question answering and R1/R5/R10/AveR on DiDeMo [25] for text-to-video retrieval.² We initialize our Video Swin Transformer (VT) with VideoSwin-Base [48], pre-trained on Kinetics-600 [28]. Language Embedder (LE) and Cross-modal Transformer (CT) are initialized from pre-trained BERT-Base [12]. During pre-training, we sparsely sample 4 video frames and randomly crop them into 224x224 to split into patches with $H = W = 32$. For all downstream tasks, we adopt the same video frame size and patch size but 5 sparse-sampled frames. We keep the training recipe (e.g., optimizer settings, masking ratio, training schedule, etc.) consistent across all targets, which we find generally good in practice.³ For MVM targets that involve a teacher model, we use official models released by the authors. We compare models pre-trained with 8 different MVM variants to the baseline pre-trained with only VTM and MLM. Our goal is to find the best MVM target features that can provide the largest performance improvement over this baseline. Results are summarized in Table 1. We first categorize the MVM targets into 4 groups, and discuss their performance in details.

One-stage Visual Targets. We include *Pixel* and *HOG*, as they do not require training a deep neural network in advance to extract these features. Compared to the baseline without MVM objective, regressing the explicit RGB colors contributes to a relatively small gain of +0.2% on TGIF-Frame and +0.8% on AveR for DiDeMo Retrieval. In contrast, HOG renders degradation on downstream video-

²We base our ablation experiments on these two representative datasets for fast iteration, our main results are reported on 13 benchmarks in Section 6. Details about downstream adaptation are included in the Appendix.

³Refer to the Appendix for more on training details.

MVM Targets	TGIF-Frame	DiDeMo-Retrieval				
	Acc.	R1	R5	R10	AveR	
Pixel	68.3	29.2	58.6	70.1	52.6	
Flow	67.6	30.3	58.0	70.3	52.9	
SIF	68.8	35.4	62.4	74.9	57.6	
SIF + Pixel	68.8	31.8	60.4	73.0	55.1	
SIF + Flow	68.7	34.4	61.5	72.8	56.3	

Table 2: **Combining MVM targets.** All variants are pre-trained on WebVid [3] for 5 epochs, using RM with 15% as the masking strategy. The final pre-training setting is highlighted in gray.

language (VidL) performance (-0.8% on TGIF-Frame and -2.0% on DiDeMo-Retrieval). We hypothesize that this is due to the missing color information in HOG features, which is critical in VidL understanding.

Supervised Pseudo-label Targets. We include *Depth Maps (Depth)* and *Optical Flow (Flow)*. Intuitively, Depth and Flow can be considered as continuous pseudo “labels”, which are made by models trained to perform depth and optical flow estimation [57, 66]. Depth does not improve over baseline with VTM+MLM. The nature of depth maps are to separate the foreground from the background, thus may guide the model to ignore information from the background, even when they are relevant for solving downstream VidL tasks (-0.1% on TGIF-Frame, -1.6% on DiDeMo Retrieval). Flow only focuses on the moving part between frames, while ignores the spatial details of static components, thus fail on more spatially-focused TGIF-Frame task (-0.5%). We also find that the optical flow estimation model easily fails with sparse sampling strategy, which is widely adopted in VidL pre-training.⁴

Supervised Visual Feature Targets. We include continuous features extracted from the last layers of image classification model [47] (i.e., *Spatial-focused Image Features (SIF)*) and action recognition model [48] (i.e., *Temporal-aware Video Features (TVF)*). We consider regressing supervised features from Swin-B or VidSwin-L as a type of knowledge distillation from unimodal models to our VIO-LET. SIF achieves significant improvement over baseline (+0.7% on TGIF-Frame and +5.8% on AveR for DiDeMo-Retrieval). In contrast, TVF fails to improve TGIF-Frame accuracy (-0.1%), though it brings notable improvement on retrieval performance (+3.6% on AveR). By distilling the knowledge from Swin-B, we enforce the model to focus more on spatial details of each frame, which we hypothesize is the main reason behind the large performance improvement. As previous study [5] pointed out, existing VidL benchmarks largely test on spatial understanding about the key frame of the video, with only a fractional of examples actually testing on temporal reasoning over multiple frames.

⁴Please find visualization examples in the Appendix.

Image Features Model	Train Data	IN-1K	DiDeMo-Retrieval				
		ACC@1	TGIF-Frame Acc.	R1	R5	R10	AveR
ResNet-50 [23]	IN-1K	76.1	67.3	29.1	58.1	69.3	52.2
Swin-T [47]	IN-1K	81.2	68.9	33.8	63.6	74.2	57.2
Swin-B	IN-1K	83.5	68.3	34.9	63.4	73.9	57.4
Swin-B	IN-22K	85.2	68.8	35.4	62.4	74.9	57.6
Swin-L	IN-22K	86.3	68.2	33.2	62.4	72.6	56.1

Table 3: **Comparing different image feature targets for MVM.** All variants are pre-trained on WebVid [3] with VTM+MLM+MVM (SIF) for 5 epochs, using RM with 15% as the masking strategy. The final pre-training setting is highlighted in gray.

Self-supervised Multimodal Feature Targets. We use *Discrete Visual Tokens (VQ)* from DALL-E [56] and continuous *Multimodal Features (MMF)* extracted from CLIP [55]. Both models are pre-trained on large-scale image-text datasets, usually much more expensive than all other targets. Both targets improve the performance by a slight margin on only one task. VQ that can capture patch-level semantics, benefits TGIF-Frame (+0.3%) which mostly focuses on scene understanding. While MMF from CLIP, contrastively pre-trained to measure the high-level similarity between the entire image and text sentence, is helpful for DiDeMo-Retrieval (+0.3% on AveR).

Summary. Among all targets, regressing RGB values (Pixel) and distilling features from Swin-B [47] (SIF) are the only two that produce consistent gains over the baseline on both downstream tasks. MVM with SIF achieves the best performance, with a gain of +0.7% on TGIF-Frame and +5.8% on AveR for DiDeMo-Retrieval over the baseline. Therefore, we use SIF as the default target for MVM in the following sections, unless specified otherwise.

5. Analyses of MVM

Combining MVM Targets. As different MVM targets focus on different aspects of visual modeling, a naive way to enable model with different visual capabilities is to combine them together. Specifically, the model pre-training can be supervised by more than one MVM loss, which are simply added together to be backpropagated. In Table 2, we find there is no merit in combining different MVM targets, leading to worse downstream performance than using SIF alone. When combining the best two targets found in Table 1: Pixel+SIF, it performs better than Pixel only, but does not improve over using SIF alone. We hypothesize that the explicit details of pixel values may conflict with the high-level visual semantics summarized in the grid features from the image classifier. We further try to combine SIF with Flow in the hope of enforcing both temporal and spatial reasoning over video inputs. In addition, Flow is a better candidate than other targets, as it demonstrates some advantages on retrieval performance in Table 1, and it is a different type

Masking Strategy	Time Cost hours	TGIF-Frame	DiDeMo-Retrieval			
		Acc.	R1	R5	R10	AveR
RM	8.0	68.8	35.4	62.4	74.9	57.6
BM	8.0	69.0	35.9	63.3	74.6	57.9
AM	34.5	68.4	31.5	59.9	72.0	54.7
RM+BM	8.0	68.7	36.4	64.2	74.4	58.3
RM+AM	20.5	68.8	33.7	63.2	73.5	56.8
BM+AM	20.5	68.9	35.6	61.9	74.4	57.3
RM+BM+AM	17.0	68.6	34.7	62.0	74.8	57.2

Table 4: Impact of **masking strategy of MVM**. All variants are pre-trained on WebVid [3] with VTM+MLM+MVM (SIF) for 5 epochs. The masking ratio is set as 15% for all masking strategies. The final pre-training setting is highlighted in gray .

of targets from SIF, compared to temporal-aware video features. The results are consistent, with improvements over optical flow only; while the performance drops, compared to SIF alone. Though our results are not encouraging, we believe how to effectively combine different MVM targets is an interesting direction for future study.

MVM Target Extractors vs. Downstream Performance.

In Table 3, we explore different image classification models as the MVM target extractor for SIF, and verify whether stronger image classification model enables better VidL performance. We compare ResNet-50 [23], Swin-Tiny/Base/Large [47], trained on ImageNet-1K (IN1K) or ImageNet-22K (IN-22K) [11]. Our results suggest that downstream VidL performance is not directly proportional to image classification accuracy. However, distilling grid features from Swin architecture is evidently more effective than that from ResNet-50, as Swin models share similar inductive bias as the VideoSwin backbone in VIOLET.

Masking Strategy. We investigate the effect of different masking strategies in Table 4, including random masking (RM), blockwise masking (BM), attended masking (AM), and their combinations. Below, we introduce each masking strategy in details.

- **Random Masking (RM).** Following the conventional practice in MLM, we randomly select a certain percentage p_m of video frame patches from the whole video inputs to be masked. In Table 5, we explore different masking ratios (p_m), and empirically find $p_m = 30\%$ gives the best downstream performance.
- **Blockwise Masking (BM).** To make MVM relying less on similar neighbor patches, we adopt blockwise masking [65, 4] that masks blocks of video patches along spatial-temporal dimension rather than independently masking randomly sampled patches for each frame. Specifically, we randomly sample an (H', W', T') as a masking block, where all $H' \times W'$ visual patches in the following T' consecutive frames will be masked; we re-

p_m	TGIF-Frame	DiDeMo-Retrieval			
	Acc.	R1	R5	R10	AveR
15%	68.8	35.4	62.4	74.9	57.6
30%	68.8	36.2	64.0	74.5	58.2
45%	68.9	35.6	61.9	74.4	57.3
60%	68.1	34.1	63.9	74.6	57.5
75%	68.3	35.4	62.4	74.2	57.3

Table 5: Impact of **masking ratio of MVM**. All variants are pre-trained on WebVid [3] with VTM+MLM+MVM (SIF) for 5 epochs, using RM as the masking strategy. The final pre-training setting is highlighted in gray .

peat this process until $> p_m$ of video patches are masked to perform MVM pre-training.

- **Attended Masking (AM).** Attended masking tries to put more weights on the more important elements based on the attention weights computed by Cross-modal Transformer (CT). A similar idea has been explored in [84] for MLM. Here, we extend AM to both visual and textual modalities. We first keep the video-text inputs intact, feed them into CT to compute the attention weights, to decide which portions in video and text are more important. We then select the top p_m of most-attended patches/tokens to be masked in video-text inputs for MVM and MLM.

To combine different masking strategies, we randomly apply one masking method for each video-text pair in a batch. Results in Table 4 suggest that TGIF-Frame can slightly benefit from BM, and combining BM with RM leads to the best retrieval performance on DiDeMo. As video usually presents analogous visual patterns in spatial-temporal neighbors (*i.e.*, nearby patches within current frame or neighboring frames), when masking patches independently (*i.e.*, RM), these neighbors can make the masked patches easy to recover, and may lead to spurious success in MVM evaluation. By masking a block (*i.e.*, BM) instead of individual patches, the model cannot merely rely on similar neighboring visual cues but requires actual visual reasoning to recover a group of missing patterns. Combining BM with RM leads to more diverse dropout patterns in video inputs, which is in analogy to data augmentation.

In addition, AM and combinations with AM are not effective for both downstream tasks. It is also worth noting that AM greatly increase the training time (4 times more than RM/BM), due to the additional forward pass needed to compute the attention weights. In our implementation, we optimize the three losses altogether in the same forward-backward pass. Hence, the performance drop with AM may be due to the important elements (*e.g.*, visual patches containing the main object or content words) are more likely to be masked together and leaving the less relevant elements

Pre-training Tasks	MVM Target	TGIF-Frame	DiDeMo-Retrieval			
		Acc.	R1	R5	R10	AveR
ITM+MLM	None	69.8	36.4	64.3	74.7	58.4
+MVM	RGB Pixel Values	69.7 (-0.1)	35.8 (-0.6)	64.4 (+0.1)	74.9 (+0.2)	58.4
	Histogram of Oriented Gradients [10]	69.8	34.9 (-1.5)	64.4 (+0.1)	75.1 (+0.4)	58.1 (-0.3)
	Depth Maps (DPT-L [57])	69.6 (-0.2)	32.3 (-4.1)	63.8 (-0.5)	74.2 (-0.5)	56.9 (-1.5)
	Spatial-focused Image Features (Swin-B [47])	69.7 (-0.1)	31.6 (-4.8)	60.5 (-3.8)	72.5 (-2.2)	54.9 (-3.5)
	Discrete Visual Tokens (DALL-E [56])	69.8	34.4 (-2.0)	62.6 (-1.7)	75.1 (+0.4)	57.4 (-1.0)
	Multimodal Features (CLIP-ViT-B [55])	69.8	33.6 (-2.8)	62.9 (-1.4)	75.6 (+0.9)	57.4 (-1.0)

Table 6: **Comparing target features for MVM applied to image-text data.** All variants are pre-trained on CC3M [62] for 5 epochs. Masking is performed randomly (RM) with ratio of 15%.

Pre-training Tasks	MVM Target		TGIF-Frame	DiDeMo-Retrieval			
	WebVid2.5M	CC3M	Acc.	R1	R5	R10	AveR
VTM+MLM	None	None	69.7	36.7	66.5	76.6	59.9
+MVM	Spatial-focused Image Features (Swin-B [47])	None	71.1	38.8	69.6	80.0	62.8
	Spatial-focused Image Features (Swin-B)	Pixel	71.3	39.7	69.3	78.4	62.5

Table 7: **Combining MVM target features for both video-text and image-text data.** All variants are pre-trained on WebVid2.5M [3] + CC3M [62] for 5 epochs. The final pre-training setting is highlighted in gray.

(e.g., scene background or stop words) intact, which will especially make the learning of video-text matching harder.

Applying MVM to Image-Text Data. As image can be considered as a special case of video with temporal size 1, video-language (VidL) pre-training can take advantages of image-text data, which has been proven successful in [35, 3]. The current trend in VidL pre-training is to leverage both video-text data and image-text data. Therefore, we repeat the experiments in Section 4 and examine which MVM targets work the best on downstream VidL tasks, when pre-trained on image-text data only. We remove optical flow and temporal-aware video features from this study, as the inputs are static images. In Table 6, we pre-train VIOLET on CC3M [62] for 5 epochs and report results on TGIF-Frame and DiDeMo-Retrieval. The performance trend with different MVM targets are not consistent with that observed on video-text data. Pixel is able to largely preserve the baseline (VTM+MLM) performance, while other MVM targets lead to different degrees of performance drop, especially on retrieval. Without visual implications from neighbor frames as video, MVM is more challenging to learn on image data. On the other hand, MVM over an image may easily fit in static visual representation, which could hurt video temporal reasoning and not benefit downstream VidL learning.

Combining Video-Text Data with Image-Text Data. We further follow [3, 39] to use both video-text data and image-text data for pre-training, and investigate different ways to combine MVM targets on image and video data in Table 7. Note that we adopt the best training strategy found in the above investigations, that is, using spatial-focused image feature (SIF) as MVM target for video inputs, and using blockwise masking (BM) + random masking (RM) with

masking ratio of 30% as the masking strategy. As the best MVM target (Pixel) on image data does not show improvement over the baseline without MVM objective in Table 6, we explore with/without MVM objective on images in this combined pre-training. For the baseline with VTM+MLM only, we simply remove the MVM objective on both image and video data, while keeping the rest training settings. Under the strict fair comparison, we observe adding MVM objectives contributes to $>+0.4\%$ gains on TGIF-Frame and $>+2.6\%$ increase on AveR for DiDeMo-Retrieval. Comparing with or without MVM objective on images, they achieve comparable performance on both tasks. Therefore, in our final setting, we only apply MVM objective on video data.

6. Main Results

In this section, we present the main results of VIOLET on 13 video-language (VidL) tasks, and compare them with prior arts. Table 8 shows the comparison on **video question answering (QA)** and **video captioning**. We observe that VIOLET is effective in learning transferable knowledge for the downstream tasks. For example, considering pre-training data at a similar scale (*i.e.*, $\leq 5M$), as shown in the top rows of Table 8, VIOLET achieves better results than prior arts, including ALPRO [39], ClipBERT [35], and SwinBERT [43], across all considered video QA and video captioning benchmarks. Specifically, when pre-training with the exact same data (*i.e.*, WebVid2.5M [3] + CC3M [62]), VIOLET surpasses ALPRO by 2.4% accuracy on MSRVTT-QA and 8.4% accuracy on MSVD-QA, respectively.

We also compare with other models pre-trained on significantly larger scale of video-text pairs. As shown in the

Method	# Pretrain videos/images	TGIF [26]			MSRVTT[80]		LSMDC [68]		MSVD [7]	Captioning	
		Act.	Trans.	Frame	MC [83]	QA [79]	MC	FiB	QA [79]	MSRVTT	MSVD
ClipBERT [35]	0.2M	82.8	87.8	60.3	88.2	37.4	-	-	-	-	-
ALPRO [39]	5M	-	-	-	-	42.1	-	-	46.3	-	-
SwinBERT [43]	-	-	-	-	-	-	-	-	-	53.8	120.6
<i>Models pre-trained on more data</i>											
JustAsk [81]	69M	-	-	-	-	41.5	-	-	46.3	-	-
MERLOT [84]	180M	94.0	96.2	69.5	90.9	43.1	81.7	52.9	-	-	-
All-in-one [71]	283M	95.5	94.7	66.3	92.3	46.8	84.4	-	48.3	-	-
MV-GPT [61]	53M	-	-	-	-	41.7	-	-	-	60.0	-
VIOLET	5M	94.8	99.0	72.8	97.6	44.5	84.4	56.9	54.7	58.0	139.2

Table 8: Comparison with SOTA on **video question answering** and **video captioning**. Accuracy and CIDEr scores are reported for QA and captioning. VIOLET is pre-trained on WebVid2.5M [3]+CC3M [62] with VTM+MLM+MVM (SIF on videos) for 10 epochs. We gray out methods that use significantly more pre-training data for fair comparison.

Method	# Pretrain videos/images	MSRVTT [80]			DiDeMo [25]			LSMDC [59]		
		R1	R5	R10	R1	R5	R10	R1	R5	R10
ClipBERT [35]	0.2M	22.0	46.8	59.9	20.4	48.0	60.8	-	-	-
Frozen [3]	5M	31.0	59.5	70.5	31.0	59.8	72.4	15.0	30.8	39.8
ALPRO [39]	5M	33.9	60.7	73.2	35.9	67.5	78.8	-	-	-
BridgeFormer [21]	5M	37.6	64.8	75.1	37.0	62.2	73.9	17.9	35.4	44.5
<i>Models pre-trained on more data</i>										
HERO [40]	136M	16.8	43.4	57.7	-	-	-	-	-	-
All-in-one [71]	138M	37.9	68.1	77.1	32.7	61.4	73.5	-	-	-
Clip4Clip [50]	400M	42.1	71.9	81.4	43.4	70.2	80.6	21.6	41.8	49.8
VIOLET	5M	37.2	64.8	75.8	47.9	76.5	84.1	24.0	43.5	54.1

Table 9: Comparison with SOTA on **text-to-video retrieval** tasks. All results are reported on R1/5/10. We gray out methods that use significantly more pre-training data for fair comparison.

bottom rows of Table 8, although we use less pre-training data than others, VIOLET still achieves comparable or better performance than those large-scale pre-training models.

We observe similar findings on video captioning. On MSRVTT captioning, VIOLET is only 2 points behind MV-GPT [61] pre-trained with 53M video-text pairs, which is 10 times larger than ours (5M). In addition, MV-GPT leverages ASR transcripts to enhance the captioning performance, while our captioning model takes only video frames as inputs and outputs the video caption.⁵ We believe augmenting VIOLET with additional modalities, such as audio or ASR transcripts, can further improve captioning performance, which we leave as future work.

Table 9 presents the comparison on **text-to-video retrieval**. When pre-training with the same datasets (*i.e.*, WebVid2.5M [3] + CC3M [62]), VIOLET shows across-the-board improvements with all metrics considered on DiDeMo and LSMDC. It is worth noting that our method performs comparably to BridgeFormer [21] on MSRVTT-Retrieval. BridgeFormer adopts a noun/verb masking strategy during pre-training, which is specially aligned to the

simple sentences in MSRVTT. However, it cannot show similar effects on DiDeMo and LSMDC due to more complex texts with multiple nouns/verbs. In contrast, the studied MVM can achieve a comprehensive enhancement in VidL learning and lead to notable improvements (+10.9% R1 on DiDeMo and +6.1% R1 on LSMDC).

7. Conclusion

We initiate the first empirical study on adopting masked visual modeling (MVM) for video-language (VidL) learning. We explore diverse MVM objectives upon end-to-end Video-Language Transformer (VIOLET), including low-level pixel space, high-level visual semantics, and extracted latent features. We show that image-text and video-text data may not share the same MVM target. Specifically, spatial-focused image feature (SIF) works the best on video-text inputs; only RGB pixel values can preserve the baseline performance without MVM for static images, paired with texts. Our analyses on different combinations of MVM targets, various SIF target extractors, and varying masking strategies/ratios shed light on effective MVM design. VIOLET pre-trained with MVM objective achieves strong performance on 3 popular VidL tasks, including video ques-

⁵Details about downstream finetuning on captioning benchmarks are included in the Appendix.

Weight Init.	MVM	TGIF-Frame	DiDeMo-Retrieval			
		Acc.	R1	R5	R10	AveR
Random	×	55.9	5.6	19.9	29.8	18.5
	✓	56.5	7.4	22.9	33.8	21.4
VidSwin-B [47])	×	68.1	28.7	57.0	69.7	51.8
	✓	68.8	35.1	63.3	73.1	57.2

Table 10: Impact of **weight initialization of video backbone**. All variants are pre-trained on WebVid [3] for 5 epochs. The MVM target is spatial-focused image features (SIF) from Swin-B [47]). The final pre-training setting is highlighted in gray .

MVM Loss	TGIF-Frame	DiDeMo-Retrieval			
	Acc.	R1	R5	R10	AveR
l_1	68.8	35.4	62.4	74.9	57.6
l_2	68.8	33.0	60.1	71.9	55.0

Table 11: Impact of **MVM loss type**. All variants are pre-trained on WebVid [3] with VTM+MLM+MVM (SIF) for 5 epochs, using RM as the masking strategy with ratio of 15%. The final pre-training setting is highlighted in gray .

tion answering, video captioning and text-to-video retrieval, across 13 VidL benchmarks.

A. Additional Results

Weight Initialization of Video Backbone. We compare the mask visual modeling (MVM) using spatial-focused image features (SIF) by different initialized video backbones in Table 10. At first, although the used video transformer (VT) is randomly initialized, the MVM training still enhances the visual representation and benefits the downstream video-language (VidL) tasks. Furthermore, MVM can also boost better initialized VT from VidSwin-B and lead to a comprehensive increase. Specifically, the improvement gap is more significant than random initialization. With a better initialized VT, we can learn better from MVM and enlarge its effectiveness during pre-training.

Type of MVM Loss. We compare the type of loss function for the MVM training by using least absolute deviations (l_1) or least square errors (l_2) in Table 11. It is well known that the l_1 loss can be resistant to outlier data. We show that MVM through l_1 is also more robust and leads to better performance on both video question answering and text-to-video retrieval than the l_2 loss.

MVM Prediction Head. We investigate the prediction head for MVM in Table 12. As a result, a single linear layer is not enough to model the complicated distilling MVM features. Therefore, we follow VTM and MLM to use 2-layer MLP as the prediction head for MVM.

MVM Head	TGIF-Frame	DiDeMo-Retrieval			
	Acc.	R1	R5	R10	AveR
1 Linear Layer	68.8	31.3	60.1	72.8	54.7
2-layer MLP	68.8	35.4	62.4	74.9	57.6

Table 12: Impact of **MVM prediction head**. All variants are pre-trained on WebVid [3] with VTM+MLM+MVM (SIF) for 5 epochs, using RM as the masking strategy with ratio of 15%. The final pre-training setting is highlighted in gray .

MVM Target	TGIF-Frame	DiDeMo-Retrieval			
	Acc.	R1	R5	R10	AveR
TVF (VidSwin-L [48])	68.0	32.8	60.5	73.0	55.4
TVF (VidSwin-B)	67.5	25.8	55.0	68.0	49.6

Table 13: **Temporal-aware video feature (TVF) target models vs. downstream performance**. All variants are pre-trained on WebVid [3] with VTM+MLM+MVM (TVF) for 5 epochs, using RM as the masking strategy with ratio of 15%.

TVF target extractors vs. downstream performance. We compare distilling video features from VidSwin-B vs. VidSwin-L (the default setting in the main text) in Table 13. Here, for experiments with VidSwin-B, the same VidSwin-B weight is used to initialize the video backbone and to extract the MVM target. Hence, the MVM objective can be easily minimized by simply ignoring the text inputs, which conflict with the other objectives. This variant is in principal similar to masked frame modeling in HERO [40], the key difference lies in whether the video backbone in VIO-LET is refined during pre-training.

Additional Exploration in Combining MVM Targets. We explore the additional combination of distilling MVM targets in Table 14. MVM with SIF has an obvious advantage over TVF only on both video question answering and text-to-video retrieval. While, considering SIF+TVF seems not to bring a robust improvement, especially decreasing text-to-video retrieval. The previous study [5] shows that current VidL benchmarks primarily focus on spatial understanding of the key frame from videos. Furthermore, combining TVF with SIF will result in an excessive training overhead. Accordingly, we choose SIF as our final pre-training setting.

Qualitative Results. Figure 2 shows good and bad examples of optical flow predictions made by RAFT-L [66] with sparsely sampled frames. In 2b, the top example shows zoom in shots, and the bottom one shows moving shots, where all content in the current frame move, which is the main reason behind the failure in optical flow estimation.

We also show the visualizations of zero-shot text-to-video retrieval on MSRVT (Figure 3), DiDeMo (Figure 4),

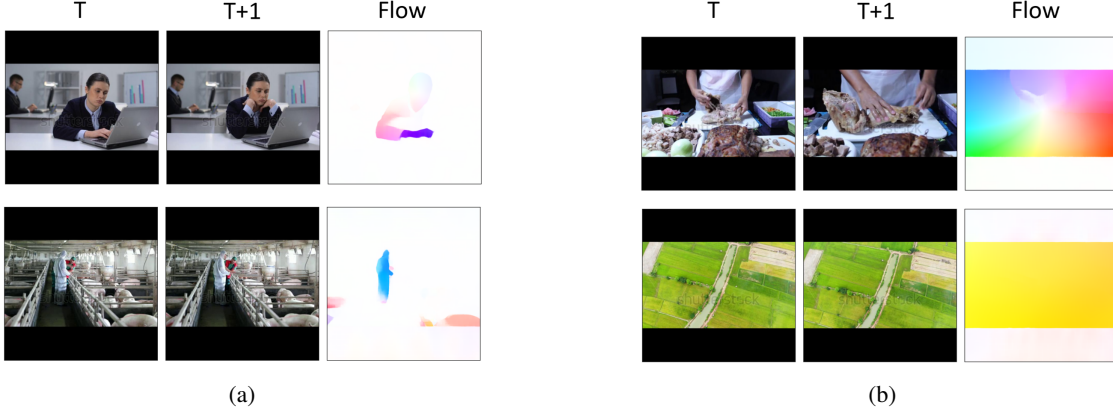


Figure 2: **Visualization of optical flow (Flow) predictions by RAFT-L [66] with sparsely sampled frames.** We show examples of good cases in (a) and bad cases in (b).

MVM Targets	TGIF-Frame	DiDeMo-Retrieval				
	Acc.	R1	R5	R10	AveR	
SIF	68.8	35.4	62.4	74.9	57.6	
TVF	68.0	32.8	60.5	73.0	55.4	
SIF + TVF	69.2	33.8	63.0	74.4	57.1	

Table 14: **Combining target features for MVM.** All variants are pre-trained on WebVid [3] for 5 epochs. The final pre-training setting is highlighted in gray.

and LSMDC (Figure 5) to demonstrate that MVM can help video understanding from different domains, such as gaming, animation, human activity, or movie scene.

B. Additional Pre-training Details

Vidoe-Text Matching (VTM). VTM enhances the cross-modal fusion via modeling the alignments between visual and textual inputs. At each training step, we randomly replace the corresponding text \mathcal{X}_{pos} for a given video \mathcal{V} with the text description \mathcal{X}_{neg} from a different video in the same batch. Both the positive pair $(\mathcal{V}, \mathcal{X}_{\text{pos}})$ and negative pair $(\mathcal{V}, \mathcal{X}_{\text{neg}})$ are modeled by Cross-modal Transformer (CT), and VTM is to tell them apart from the global VidL representation h^c of the $[\text{CLS}]$ token. In particular, h^c will be processed by a fully-connected layer (FC^{VTM}) to learn contrastively through classification:

$$b^{\text{pos}} = \text{FC}^{\text{VTM}}(h^c_{\text{pos}}), b^{\text{neg}} = \text{FC}^{\text{VTM}}(h^c_{\text{neg}}),$$

$$\mathcal{L}_{\text{VTM}} = -\frac{1}{B} \sum_i^B \log \frac{b_i^{\text{pos}}}{b_i^{\text{pos}} + \sum b_i^{\text{neg}}}, \quad (2)$$

where h^c_{pos} or h^c_{neg} is h^c of positive or negative pairs.

Masked Language Modeling (MLM). In MLM, we randomly mask out some word tokens with a probability of

15%.⁶ The goal is to recover these masked word tokens x from the joint VidL features h modeled by CT. Specifically, the corresponding h^x for these masked tokens are fed in a fully-connected layer (FC^{MLM}) and projected to the discrete token space for classification:

$$x'_i = \text{FC}^{\text{MLM}}(h^x_i),$$

$$\mathcal{L}_{\text{MLM}} = -\mathbb{E} \left[\frac{1}{|\mathcal{M}^{\text{MLM}}|} \sum_{i \in \mathcal{M}^{\text{MLM}}} \log P(x_i | x'_i) \right], \quad (3)$$

where \mathcal{M}^{MLM} denotes the index set of masked word tokens.

Implementation Details Our implementation of VIOLET is based on PyTorch [53]. As discussed in the main text and supported by the additional experimental results above, our final pre-training setting is (i) VTM+MLM+MVM (with MVM target as spatial-focused image features from Swin-B [47], applied on video-text inputs only) as the pre-training tasks; (ii) 2-layer MLP as the MVM prediction head and l_1 regression as the MVM loss; and (iii) blockwise masking + random masking with masking ratio of 30% as the masking strategy. We adopt AdamW [49] as the optimizer with a warmup learning rate schedule of 5e-5 peak learning rate, betas of (0.9, 0.98), and weight decay of 1e-3 for all pre-training experiments. We pre-train VIOLET on 32 NVIDIA V100 GPUs with a batch size of 28 per GPU. Pre-training with 10 epochs on WebVid2.5M [3] + CC3M [62] takes about 27 hours to finish. We present the training settings for all finetuning experiments in the next section.

C. Experimental Setup of Downstream Tasks

We evaluate our pre-trained VIOLET on 3 popular video-language tasks, including text-to-video retrieval, video question answering and video captioning, across 13

⁶Following BERT [12], We replace 80% of masked word tokens as the $[\text{MASK}]$ token, 10% as a random token, and 10% as its original token.

downstream datasets. For text-to-video retrieval, we report model performance on MSRVT [80], DiDeMo [25], and LSMDC [59] and use Recall at K (R@K, K=1,5,10) as the evaluation metric. For video question answering, we consider datasets in both multiple-choice and open-ended settings, including TGIF-Action, TGIF-Transition, TGIF-Frame [26], MSRVT-MC [83], MSRVT-QA, MSVD-QA [79], LSMDC-MC and LSMDC-FiB [68]. We evaluate our models using accuracy. For video captioning, we report CIDER scores on MSRVT and MSVD.

We follow the standard training/validation/testing splits of the original datasets. If not otherwise stated, we sparsely sample $T = 5$ video frames and adopt video frame size 224 with patch size $H = W = 32$. Similar to pre-training, we use AdamW [49] to fine-tune VIOLET for each downstream task with a warmup learning rate schedule of $2e-5$ peak learning rate, betas of (0.9, 0.98), and weight decay of $1e-3$. All finetuning experiments are conducted on Microsoft Azure [1] adopting mixed-precision training with DeepSpeed [58].⁷ All video data are pre-processed by evenly extracting 32 frames to avoid expensive decoding on-the-fly. During training, we randomly sample T frames from 32 frames, resize the shorter side of all frames to 224 and random crop (224x224) at the same location for all the frames in a given video. During inference, we evenly sample T frames from 32 frames and center crop (224x224) for all the sampled video frames.

C.1. Text-To-Video Retrieval

For text-to-video retrieval, similar to visual-text matching (VTM) during pre-training, we treat corresponding video-text pairs in the same batch as positives and all other pairwise combinations as negatives. We adopt a fully-connected (FC) layer (FC^{T2V}) over the VidL representation h^c of the [CLS] token to learn through classification:

$$b^{\text{pos}} = FC^{T2V}(h_{\text{pos}}^c), b^{\text{neg}} = FC^{T2V}(h_{\text{neg}}^c),$$

$$\mathcal{L}_{T2V} = -\frac{1}{B} \sum_i \log \frac{b_i^{\text{pos}}}{b_i^{\text{pos}} + \sum b_i^{\text{neg}}}, \quad (4)$$

where h_{pos}^c or h_{neg}^c is h^c of positive or negative pairs. In particular, we use pre-trained FC^{VTM} for zero-shot text-to-video retrieval and to initialize FC^{T2V} for further fine-tuning on each downstream text-to-video retrieval task.

MSRVT [80] contains 10K YouTube videos with 200K human annotations. For fair comparison [3, 35], we train on 9K training+validation splits and evaluate on the 1K-A testing split. We adopt batch size 20 per GPU and train for 10 epochs.

⁷We conduct retrieval finetuning on 8 x 80GB A100 GPUs to enable larger batch size, while all other finetuning experiments are conducted on 8 x 32GB V100 GPUs.

VideoQA	Task	#Option
Multiple-Choice	TGIF-Action [26]	5
	TGIF-Transition [26]	5
	MSRVT-MC [83]	5
	LSMDC-MC [68]	5
Open-Ended	TGIF-Frame [26]	-
	MSRVT-QA [79]	-
	MSVD-QA [8]	-
	LSMDC-FiB [68]	-

Table 15: Summary of **video question answering** tasks. For open-ended Video QA, we do not limit the answer vocabulary to a fixed answer candidate set.

DiDeMo [25] consists of 10K videos annotated with 40K sentences from Flickr. Following [3, 35], we concatenate all sentences from the same video into a paragraph and perform paragraph-to-video retrieval for DiDeMo. We adopt batch size 16 per GPU and train for 10 epochs.

LSMDC [59] contains 118K video clips from 202 movies. Each clip has a caption from movie scripts or descriptive video services. Following [3, 52], we evaluate on 1K testing clips that disjoint from the training+validation splits. We adopt batch size 20 per GPU and train for 5 epochs.

C.2. Video Question Answering

We test our model on video question answering (QA) tasks in both multiple-choice and open-ended settings, as summarized in Table 15. We follow LAVENDER [41] to formulate Video QA as Masked Language Modeling due to its superior performance. For multiple-choice QA tasks, we concatenate question with all answer options and add a [MASK] to form the input text ($Q+A0+A1+A2+A3+A4+[MASK]$). We adopt the same Masked Language Modeling (MLM) layer used during pre-training upon h^x to predict the word token corresponding to the answer index (e.g., 0, 1, 2, 3, 4). Similarly, for open-ended QA tasks, we apply MLM over the input ($Q+[MASK]$). Cross-entropy loss is used to supervise the downstream finetuning over the whole word vocabulary.

TGIF-Action, TGIF-Transition, and TGIF-Frame [26] require spatial-temporal reasoning to answer questions regarding GIF videos in TGIF-QA. Specifically, we aim to test our model along three dimensions: (i) **Action**: to recognize the repeated action; (ii) **Transition**: to identify the transition between the before and after states; (iii) **Frame**: to answer questions about a specific frame from the GIF video. Among them, TGIF-Action and TGIF-Transition are collected under multiple-choice setting, and TGIF-Frame is an open-ended video QA task with free-form answers. We adopt batch size 24 and train for 56/20/10 epochs for Action/Transition/Frame, respectively.

MSRVTT-MC [83] and **MSRVTT-QA** [79] are created based on videos and captions in MSRVTT [80]. MSRVTT-MC is a multiple-choice task with videos as questions, and captions as answers. Each video contains 5 captions, with only one positive match. This setting can be viewed as video-to-text retrieval, hence we simply evaluate the model trained on MSRVTT-Retrieval. MSRVTT-QA contains 243K open-ended questions over 10K videos. We adopt batch size 24 per GPU and training epochs 8.

MSVD-QA [79] consists of 47K open-ended questions over 2K videos, based on video-caption pairs from MSVD [8]. We adopt batch size 24 per GPU and train for 10 epochs.

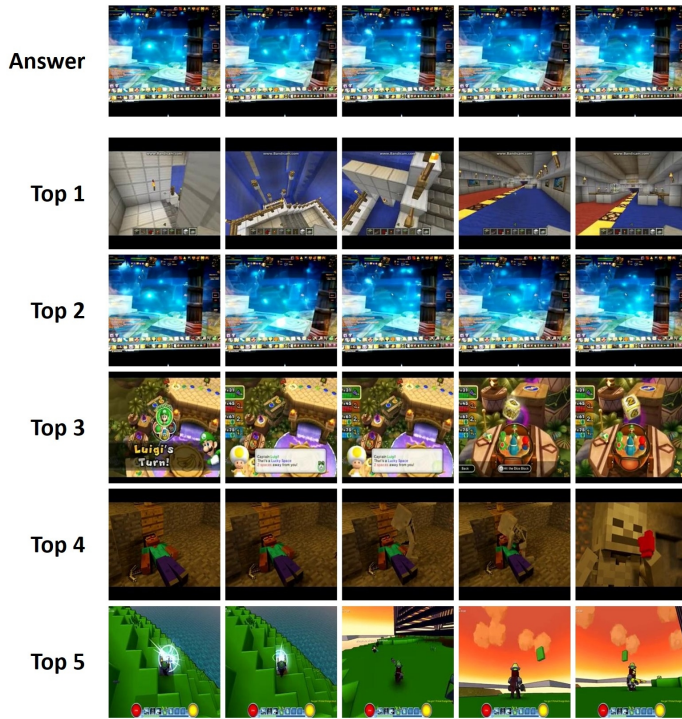
LSMDC-MC and **LSMDC-FiB** [68] are built from LSMDC dataset [59]. Similar to MSRVTT-MC, LSMDC-MC requires the model to select the only positive caption that describes the video from 5 caption candidates, and can be formulated as video-to-text retrieval. LSMDC-FiB replaces a word in the question sentence with the [BLANK] token, and requires the model to recover the missing word. We regard LSMDC-FiB as an open-ended Video QA task. In particular, we replace the [BLANK] token with [MASK] token, and use the MLM prediction head over the representation h_x of the [MASK] token to predict the correct answer. We adopt batch size 24 per GPU and train for 10 epochs.

C.3. Video Captioning

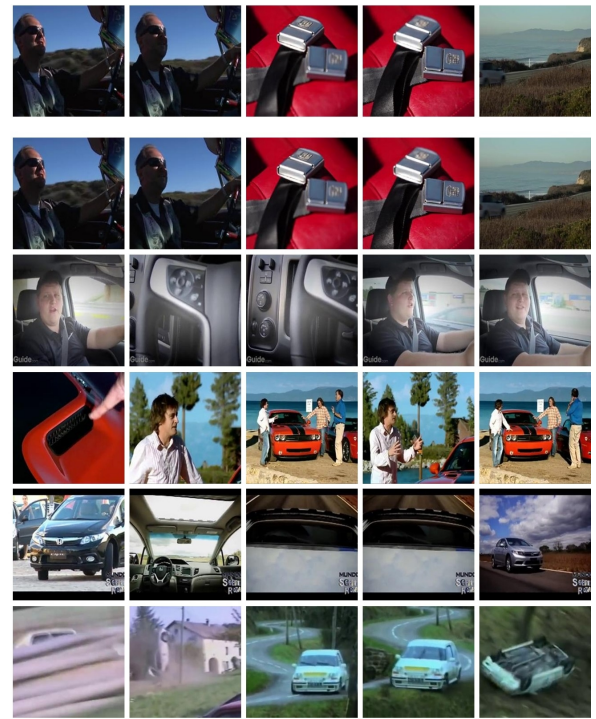
For video captioning, we evaluate on MSRVTT [80] and MSVD [7]. **MSRVTT** consists of 10K videos with 20 captions per video, and **MSVD** contains 2K videos, with 40 captions per video. We follow the standard captioning splits to train/evaluate with VIOLET.

The captioning finetuning is formulated as masked language modeling (MLM) with a causal attention mask so that the current word token only attends to the tokens before it, but not the ones after it, following SwinBERT [43]. During training, we set the probability of random masking caption tokens to be 0.15, the same as what is used in MLM during pre-training. We adopt batch size 24 per GPU and train for 20 epochs. During inference, we generate the captions auto-regressively. At each generation step, a [MASK] token is appended to the previously generated tokens, and the model will predict the current tokens based on the learned embedding at the [MASK] token position. We perform generation until the model outputs a [SEP], which is defined as the sentence ending token or when it reaches the maximum generation step 50.

"There is someone playing a game in a computer."



"A man is driving a car through the countryside."



"Cartoon show for kids."



"A video of a rock group performing one of their songs."

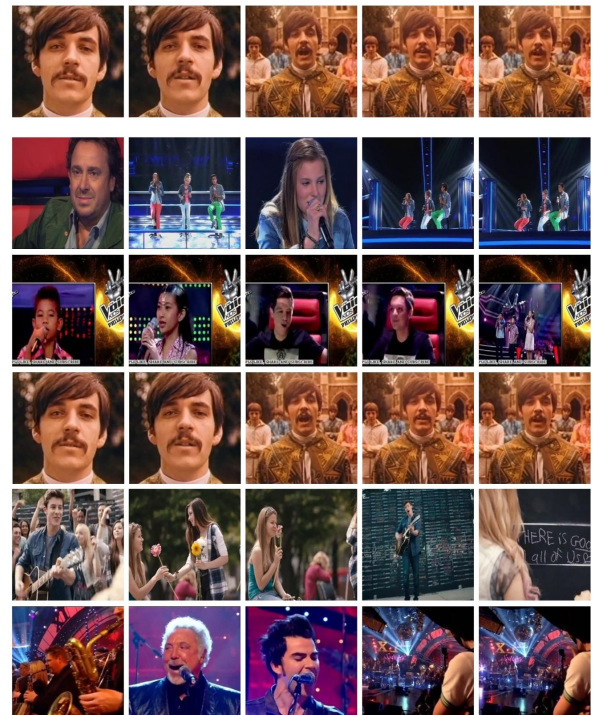


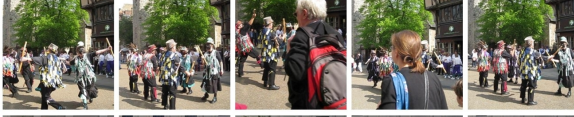
Figure 3: Qualitative examples of zero-shot text-to-video retrieval on MSRVT [80].

"We first see people walking and the crowd. We see two people walk in front of audience. A man ' lady are seen walking through the scene."

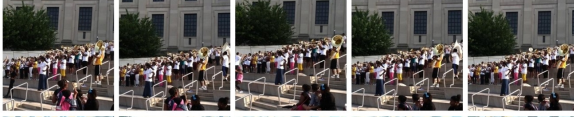
Answer



Top 1



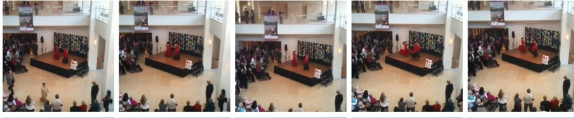
Top 2



Top 3



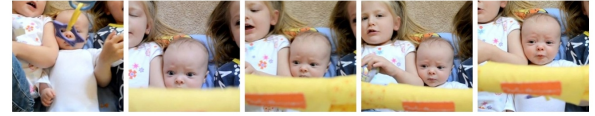
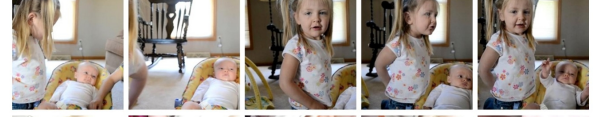
Top 4



Top 5



"The baby's face is the only face in the frame for a brief time. Then the two girls on either side of him appear. A blond hair girl kneels down beside a baby."



"There is a parade and a man goes by with a tuba. Camera zooms in on a statue the crowd is carrying. A parade of people in purple carrying something passes."

Answer



Top 1



Top 2



Top 3



Top 4



Top 5



"Camera zooms in to stage. Blue light flashes up from the stage for first time. Lights change from red to green the second time."



Figure 4: Qualitative examples of **zero-shot text-to-video retrieval** on DiDeMo [25].

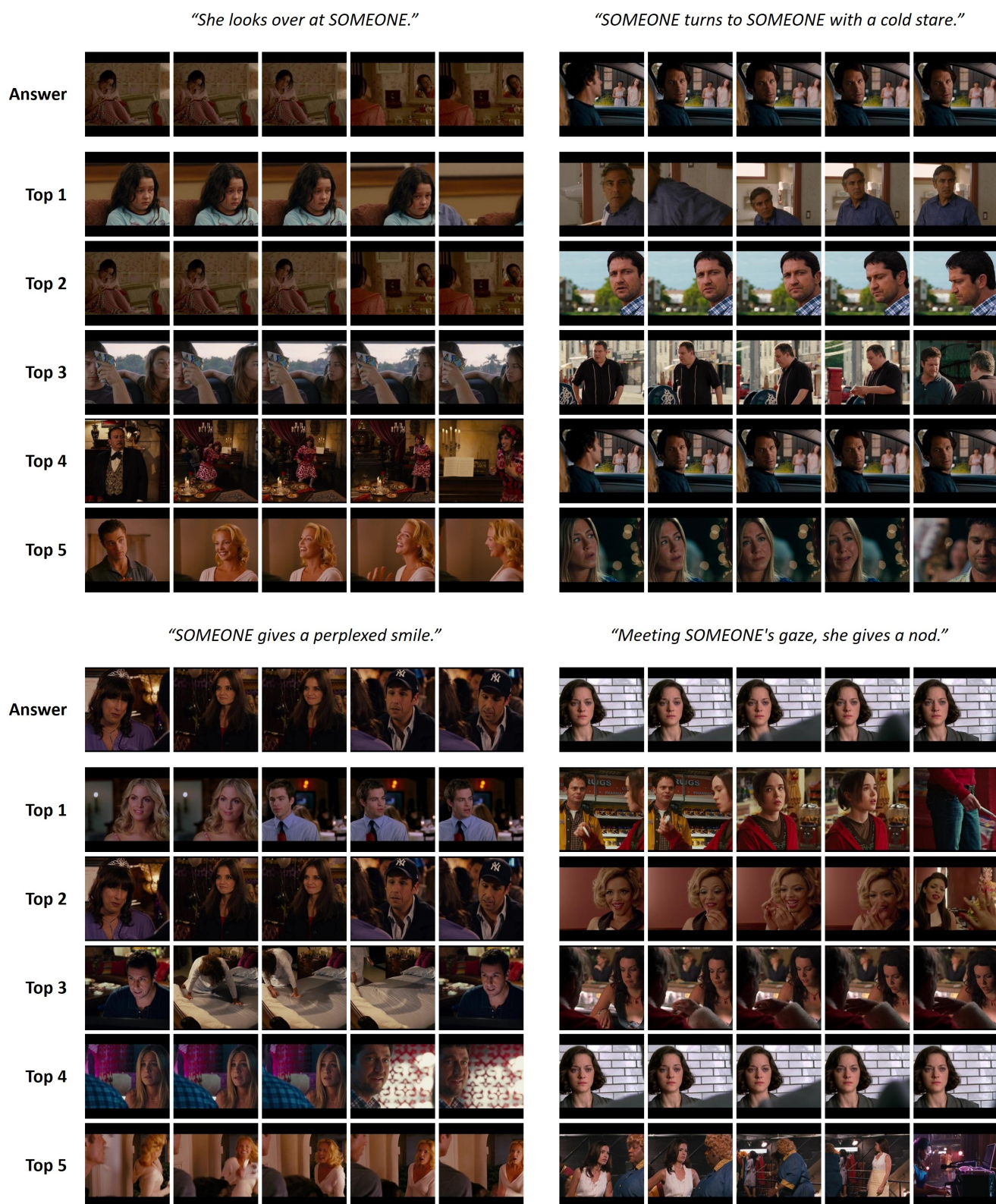


Figure 5: Qualitative examples of **zero-shot text-to-video retrieval** on LSMDC [59].

References

- [1] Microsoft Azure. <https://azure.microsoft.com/>. 11
- [2] Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang. Bottom-Up and Top-Down Attention for Image Captioning and Visual Question Answering. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2
- [3] Max Bain, Arsha Nagrani, Gül Varol, and Andrew Zisserman. Frozen in Time: A Joint Video and Image Encoder for End-to-End Retrieval. In *International Conference on Computer Vision (ICCV)*, 2021. 2, 4, 5, 6, 7, 8, 9, 10, 11
- [4] Hangbo Bao, Li Dong, and Furu Wei. BEiT: BERT Pre-Training of Image Transformers. In *International Conference for Learning Representations (ICLR)*, 2022. 1, 2, 3, 4, 6
- [5] Shyamal Buch, Cristóbal Eyzaguirre, Adrien Gaidon, Jiajun Wu, Li Fei-Fei, and Juan Carlos Niebles. Revisiting the "Video" in Video-Language Understanding. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 5, 9
- [6] Joao Carreira and Andrew Zisserman. Quo Vadis, Action Recognition? A New Model and the Kinetics Dataset. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 2
- [7] David L. Chen and William B. Dolan. Collecting Highly Parallel Data for Paraphrase Evaluation. In *ACL*, 2011. 1, 8, 12
- [8] David L. Chen and William B. Dolan. Collecting Highly Parallel Data for Paraphrase Evaluation. In *Annual Meetings of the Association for Computational Linguistics (ACL)*, 2011. 11, 12
- [9] Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. UNITER: UNiversal Image-TExt Representation Learning. In *European Conference on Computer Vision (ECCV)*, 2020. 2
- [10] Navneet Dalal and Bill Triggs. Histograms of Oriented Gradients for Human Detection. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2005. 3, 4, 7
- [11] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: a Large-Scale Hierarchical Image Database. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2009. 2, 6
- [12] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Conference of the North American Chapter of the Association for Computational Linguistics (NAACL)*, 2019. 1, 2, 4, 10
- [13] 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. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. In *International Conference for Learning Representations (ICLR)*, 2021. 4
- [14] Zi-Yi Dou, Yichong Xu, Zhe Gan, Jianfeng Wang, Shuohang Wang, Lijuan Wang, Chenguang Zhu, Pengchuan Zhang, Lu Yuan, Nanyun Peng, et al. An empirical study of training end-to-end vision-and-language transformers. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2
- [15] Chenyou Fan, Xiaofan Zhang, Shu Zhang, Wensheng Wang, Chi Zhang, and Heng Huang. Heterogeneous Memory Enhanced Multimodal Attention Model for Video Question Answering. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2
- [16] Christoph Feichtenhofer, Haoqi Fan, Jitendra Malik, and Kaiming He. SlowFast Networks for Video Recognition. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2
- [17] Tsu-Jui Fu, Linjie Li, Zhe Gan, Kevin Lin, William Yang Wang, Lijuan Wang, and Zicheng Liu. VIOLET: End-to-End Video-Language Transformers with Masked Visual-token Modeling. In *arXiv:2111.1268*, 2021. 1, 3
- [18] Valentin Gabeur, Chen Sun, Karteek Alahari, and Cordelia Schmid. Multi-modal Transformer for Video Retrieval. In *European Conference on Computer Vision (ECCV)*, 2020. 2
- [19] Jiyang Gao, Runzhou Ge, Kan Chen, and Ram Nevatia. Motion-Appearance Co-Memory Networks for Video Question Answering. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2
- [20] Jiyang Gao, Chen Sun, Zhenheng Yang, and Ram Nevatia. TALL: Temporal Activity Localization via Language Query. In *International Conference on Computer Vision (ICCV)*, 2017. 2
- [21] Yuying Ge, Yixiao Ge, Xihui Liu, Dian Li, Ying Shan, Xiaohu Qie, and Ping Luo. BridgeFormer: Bridging Video-text Retrieval with Multiple Choice Questions. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 8
- [22] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked Autoencoders Are Scalable Vision Learners. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 1, 2, 3
- [23] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep Residual Learning for Image Recognition. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 2, 5, 6
- [24] Lisa Anne Hendricks, Oliver Wang, Eli Shechtman, Josef Sivic, Trevor Darrell, and Bryan Russell. Localizing Moments in Video with Natural Language. In *International Conference on Computer Vision (ICCV)*, 2017. 2
- [25] Lisa Anne Hendricks, Oliver Wang, Eli Shechtman, Josef Sivic, Trevor Darrell, and Bryan Russell. Localizing Moments in Video with Natural Language. In *International Conference on Computer Vision (ICCV)*, 2017. 4, 8, 11, 14
- [26] Yunseok Jang, Yale Song, Youngjae Yu, Youngjin Kim, and Gunhee Kim. TGIF-QA: Toward Spatio-Temporal Reasoning in Visual Question Answering. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 1, 2, 4, 8, 11
- [27] Jianwen Jiang, Ziqiang Chen, Haojie Lin, Xibin Zhao, and Yue Gao. Divide and Conquer: Question-Guided Spatio-

- Temporal Contextual Attention for Video Question Answering. In *AAAI Conference on Artificial Intelligence (AAAI)*, 2020. [2](#)
- [28] Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan, Fabio Viola, Tim Green, Trevor Back, Paul Natsev, Mustafa Suleyman, and Andrew Zisserman. The Kinetics Human Action Video Dataset. In *arXiv:1705.06950*, 2017. [2](#), [4](#)
- [29] Seonhoon Kim, Seohyeong Jeong, Eunbyul Kim, Inho Kang, and Nojun Kwak. Self-supervised Pre-training and Contrastive Representation Learning for Multiple-choice Video QA. In *AAAI Conference on Artificial Intelligence (AAAI)*, 2021. [2](#)
- [30] Ranjay Krishna, Kenji Hata, Frederic Ren, Li Fei-Fei, and Juan Carlos Niebles. Dense-Captioning Events in Videos. In *International Conference on Computer Vision (ICCV)*, 2017. [1](#), [2](#)
- [31] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A. Shamma, Michael S. Bernstein, and Fei-Fei Li. Visual Genome: Connecting Language and Vision Using Crowdsourced Dense Image Annotations. In *International Journal of Computer Vision (IJCV)*, 2017. [2](#)
- [32] Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. ALBERT: A Lite BERT for Self-supervised Learning of Language Representations. In *International Conference for Learning Representations (ICLR)*, 2020. [2](#)
- [33] Thao Minh Le, Vuong Le, Svetha Venkatesh, and Truyen Tran. Hierarchical Conditional Relation Networks for Video Question Answering. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. [2](#)
- [34] Jie Lei, Tamara L Berg, and Mohit Bansal. QVHighlights: Detecting Moments and Highlights in Videos via Natural Language Queries. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2021. [2](#)
- [35] Jie Lei, Linjie Li, Luowei Zhou, Zhe Gan, Tamara L. Berg, Mohit Bansal, and Jingjing Liu. Less is More: ClipBERT for Video-and-Language Learning via Sparse Sampling. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. [2](#), [7](#), [8](#), [11](#)
- [36] Jie Lei, Licheng Yu, Mohit Bansal, and Tamara L. Berg. TVQA: Localized, Compositional Video Question Answering. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2018. [2](#)
- [37] Jie Lei, Licheng Yu, Tamara L. Berg, and Mohit Bansal. TVQA+: Spatio-Temporal Grounding for Video Question Answering. In *Annual Meeting of the Association for Computational Linguistics (ACL)*, 2020. [2](#)
- [38] Jie Lei, Licheng Yu, Tamara L. Berg, and Mohit Bansal. TVR: A Large-Scale Dataset for Video-Subtitle Moment Retrieval. In *European Conference on Computer Vision (ECCV)*, 2020. [2](#)
- [39] Dongxu Li, Junnan Li, Hongdong Li, Juan Carlos Niebles, and Steven C.H. Hoi. Align and Prompt: Video-and-Language Pre-training with Entity Prompts. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. [7](#), [8](#)
- [40] Linjie Li, Yen-Chun Chen, Yu Cheng, Zhe Gan, Licheng Yu, and Jingjing Liu. HERO: Hierarchical Encoder for Video+Language Omni-representation Pre-training. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2020. [1](#), [2](#), [8](#), [9](#)
- [41] Linjie Li, Zhe Gan, Kevin Lin, Chung-Ching Lin, Zicheng Liu, Ce Liu, and Lijuan Wang. Lavender: Unifying video-language understanding as masked language modeling. *arXiv preprint arXiv:2206.07160*, 2022. [11](#)
- [42] Linjie Li, Jie Lei, Zhe Gan, Licheng Yu, Yen-Chun Chen, Rohit Pillai, Yu Cheng, Luowei Zhou, Xin Eric Wang, William Yang Wang, Tamara Lee Berg, Mohit Bansal, Jingjing Liu, Lijuan Wang, and Zicheng Liu. VALUE: A Multi-Task Benchmark for Video-and-Language Understanding Evaluation. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2021. [2](#)
- [43] Kevin Lin, Linjie Li, Chung-Ching Lin, Faisal Ahmed, Zhe Gan, Zicheng Liu, Yumao Lu, and Lijuan Wang. SwinBERT: End-to-End Transformers with Sparse Attention for Video Captioning. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. [7](#), [8](#), [12](#)
- [44] Song Liu, Haoqi Fan, Shengsheng Qian, Yiru Chen, Wenkui Ding, and Zhongyuan Wang. HiT: Hierarchical Transformer with Momentum Contrast for Video-Text Retrieval. In *arXiv:2103.15049*, 2021. [2](#)
- [45] Yang Liu, Samuel Albanie, Arsha Nagrani, and Andrew Zisserman. Use What You Have: Video Retrieval Using Representations From Collaborative Experts. In *British Machine Vision Conference (BMVC)*, 2020. [2](#)
- [46] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. RoBERTa: A Robustly Optimized BERT Pretraining Approach. In *arXiv:1907.11692*, 2019. [2](#)
- [47] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin Transformer: Hierarchical Vision Transformer using Shifted Windows. In *International Conference on Computer Vision (ICCV)*, 2021. [4](#), [5](#), [6](#), [7](#), [9](#), [10](#)
- [48] Ze Liu, Jia Ning, Yue Cao, Yixuan Wei, Zheng Zhang, Stephen Lin, and Han Hu. Video Swin Transformer. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. [4](#), [5](#), [9](#)
- [49] Ilya Loshchilov and Frank Hutter. Decoupled Weight Decay Regularization. In *International Conference for Learning Representations (ICLR)*, 2019. [10](#), [11](#)
- [50] 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. In *arXiv:2104.08860*, 2021. [8](#)
- [51] Antoine Miech, Jean-Baptiste Alayrac, Lucas Smaira, Ivan Laptev, Josef Sivic, and Andrew Zisserman. End-to-End Learning of Visual Representations from Uncurated Instructional Videos. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. [2](#)
- [52] 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 *International Conference on Computer Vision (ICCV)*, 2019. 1, 2, 11
- [53] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. PyTorch: An Imperative Style, High-Performance Deep Learning Library. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2019. 10
- [54] Mandela Patrick, Po-Yao Huang, Yuki Asano, Florian Metze, Alexander Hauptmann, Joao Henriques, and Andrea Vedaldi. Support-set bottlenecks for video-text representation learning. In *International Conference for Learning Representations (ICLR)*, 2021. 2
- [55] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning Transferable Visual Models From Natural Language Supervision. In *International Conference on Machine Learning (ICML)*, 2021. 2, 4, 5, 7
- [56] 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 (ICML)*, 2021. 4, 5, 7
- [57] Rene Ranftl, Alexey Bochkovskiy, and Vladlen Koltun. Vision Transformers for Dense Prediction. In *International Conference on Computer Vision (ICCV)*, 2021. 3, 4, 5, 7
- [58] Jeff Rasley, Samyam Rajbhandari, Olatunji Ruwase, and Yuxiong He. DeepSpeed: System optimizations enable training deep learning models with over 100 billion parameters. In *KDD*, 2020. 11
- [59] Anna Rohrbach, Marcus Rohrbach, Niket Tandon, and Bernt Schiele. A Dataset for Movie Description. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015. 1, 2, 8, 11, 12, 15
- [60] Andrew Rouditchenko, Angie Boggust, David Harwath, Brian Chen, Dhiraj Joshi, Samuel Thomas, Kartik Audhkhasi, Hilde Kuehne, Rameswar Panda, Rogerio Feris, Brian Kingsbury, Michael Picheny, Antonio Torralba, and James Glass. AVLnet: Learning Audio-Visual Language Representations from Instructional Videos. In *INTER-SPEECH*, 2021. 2
- [61] Paul Hongsuck Seo, Arsha Nagrani, Anurag Arnab, and Cordelia Schmid. End-to-end Generative Pretraining for Multimodal Video Captioning. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 8
- [62] Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual Captions: A Cleaned, Hypernymed, Image Alt-text Dataset For Automatic Image Captioning. In *Annual Meeting of the Association for Computational Linguistics (ACL)*, 2018. 7, 8, 10
- [63] Chen Sun, Austin Myers, Carl Vondrick, Kevin Murphy, and Cordelia Schmid. VideoBERT: A Joint Model for Video and Language Representation Learning. In *International Conference on Computer Vision (ICCV)*, 2019. 1, 2
- [64] Hao Tan and Mohit Bansal. LXMERT: Learning Cross-Modality Encoder Representations from Transformers. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2019. 2
- [65] Hao Tan, Jie Lei, Thomas Wolf, and Mohit Bansal. VIM-PAC: Video Pre-Training via Masked Token Prediction and Contrastive Learning. In *arXiv:2106.11250*, 2021. 2, 6
- [66] Zachary Teed and Jia Deng. RAFT: Recurrent All-Pairs Field Transforms for Optical Flow. In *European Conference on Computer Vision (ECCV)*, 2020. 3, 4, 5, 9, 10
- [67] Zhan Tong, Yibing Song, Jue Wang, and Limin Wang. VideoMAE: Masked Autoencoders are Data-Efficient Learners for Self-Supervised Video Pre-Training. In *arXiv:2203.12602*, 2022. 2
- [68] Atousa Torabi, Niket Tandon, and Leonid Sigal. Learning Language-Visual Embedding for Movie Understanding with Natural-Language. In *arXiv:1609.08124*, 2016. 8, 11, 12
- [69] Aaron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu. Neural Discrete Representation Learning. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2017. 4
- [70] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention Is All You Need. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2017. 2
- [71] Alex Jinpeng Wang, Yixiao Ge, Rui Yan, Yuying Ge, Xudong Lin, Guanyu Cai, Jianping Wu, Ying Shan, Xiaohu Qie, and Mike Zheng Shou. All in One: Exploring Unified Video-Language Pre-training. In *arXiv:2203.07303*, 2022. 8
- [72] Limin Wang, Yuanjun Xiong, Zhe Wang, Yu Qiao, Dahua Lin, Xiaoou Tang, and Luc Van Gool. Temporal Segment Networks: Towards Good Practices for Deep Action Recognition. In *European Conference on Computer Vision (ECCV)*, 2016. 2
- [73] Rui Wang, Dongdong Chen, Zuxuan Wu, Yinpeng Chen, Xiyang Dai, Mengchen Liu, Yu-Gang Jiang, Luwei Zhou, and Lu Yuan. BEVT: BERT Pretraining of Video Transformers. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2
- [74] 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 *International Conference on Computer Vision (ICCV)*, 2019. 2
- [75] Chen Wei, Haoqi Fan, Saining Xie, Chao-Yuan Wu, Alan Yuille, and Christoph Feichtenhofer. Masked Feature Prediction for Self-Supervised Visual Pre-Training. In *arXiv:2112.09133*, 2022. 1, 2, 3
- [76] Longhui Wei, Lingxi Xie, Wengang Zhou, Houqiang Li, and Qi Tian. MVP: Multimodality-guided Visual Pre-training. In *arXiv:2203.05175*, 2022. 1, 2
- [77] Saining Xie, Chen Sun, Jonathan Huang, Zhuowen Tu, and Kevin Murphy. Rethinking Spatiotemporal Feature Learning: Speed-Accuracy Trade-offs in Video Classification. In *European Conference on Computer Vision (ECCV)*, 2018. 2

- [78] Zhenda Xie, Zheng Zhang, Yue Cao, Yutong Lin, Jianmin Bao, Zhuliang Yao, Qi Dai, and Han Hu. SimMIM: A Simple Framework for Masked Image Modeling. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 1, 2, 3
- [79] Dejing Xu, Zhou Zhao, Jun Xiao, Fei Wu, Hanwang Zhang, Xiangnan He, and Yueting Zhuang. Video Question Answering via Gradually Refined Attention over Appearance and Motion. In *ACM Multimedia (ACMMM)*, 2017. 1, 2, 8, 11, 12
- [80] Jun Xu, Tao Mei, Ting Yao, and Yong Rui. MSR-VTT: A Large Video Description Dataset for Bridging Video and Language. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 1, 2, 8, 11, 12, 13
- [81] Antoine Yang, Antoine Miech, Josef Sivic, Ivan Laptev, and Cordelia Schmid. Just Ask: Learning to Answer Questions from Millions of Narrated Videos. In *International Conference on Computer Vision (ICCV)*, 2021. 2, 8
- [82] Zekun Yang, Noa Garcia, Chenhui Chu, Mayu Otani, Yuta Nakashima, and Haruo Takemura. BERT Representations for Video Question Answering. In *Winter Conference on Applications of Computer Vision (WACV)*, 2020. 2
- [83] Youngjae Yu, Jongseok Kim, and Gunhee Kim. A Joint Sequence Fusion Model for Video Question Answering and Retrieval. In *European Conference on Computer Vision (ECCV)*, 2018. 8, 11, 12
- [84] Rowan Zellers, Ximing Lu, Jack Hessel, Youngjae Yu, Jae Sung Park, Jize Cao, Ali Farhadi, and Yejin Choi. MERLOT: Multimodal Neural Script Knowledge Models. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2021. 1, 2, 6, 8
- [85] Bowen Zhang, Hexiang Hu, and Fei Sha. Cross-Modal and Hierarchical Modeling of Video and Text. In *European Conference on Computer Vision (ECCV)*, 2018. 2
- [86] Jinghao Zhou, Chen Wei, Huiyu Wang, Wei Shen, Cihang Xie, Alan Yuille, and Tao Kong. iBOT: Image BERT Pre-Training with Online Tokenizer. In *International Conference on Learning Representations (ICLR)*, 2022. 1, 2, 3
- [87] Luowei Zhou, Chenliang Xu, and Jason J. Corso. Towards Automatic Learning of Procedures from Web Instructional Videos. In *AAAI Conference on Artificial Intelligence (AAAI)*, 2018. 2
- [88] Linchao Zhu and Yi Yang. ActBERT: Learning Global-Local Video-Text Representations. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 1, 2