

Dual-Level Decoupled Transformer for Video Captioning

Yiqi Gao^{1,2} Xinglin Hou³ Wei Suo^{1,2} Mengyang Sun^{1,2}
Tiezheng Ge³ Yuning Jiang³ Peng Wang^{1,2}

¹School of Computer Science, Northwestern Polytechnical University

²National Engineering Lab for Integrated Aero-Space-Ground-Ocean Big Data Application Technology

³Alibaba Group

{gyqjz,suwei1994,sunmenmian}@mail.nwpu.edu.cn

{xinglin.hxl,tiezheng.gtz,mengzhu.jyn}@alibaba-inc.com

peng.wang@nwpu.edu.cn

ABSTRACT

Video captioning aims to understand the spatio-temporal semantic concept of the video and generate descriptive sentences. The de-facto approach to this task dictates a text generator to learn from *offline-extracted* motion or appearance features from *pre-trained* vision models. However, these methods may suffer from the so-called "*couple*" drawbacks on both *video spatio-temporal representation* and *sentence generation*. For the former, "*couple*" means learning spatio-temporal representation in a single model(3DCNN), resulting the problems named *disconnection in task/pre-train domain* and *hard for end-to-end training*. As for the latter, "*couple*" means treating the generation of visual semantic and syntax-related words equally. To this end, we present \mathcal{D}^2 - a dual-level decoupled transformer pipeline to solve the above drawbacks: (i) for video spatio-temporal representation, we decouple the process of it into "first-spatial-then-temporal" paradigm, releasing the potential of using dedicated model(e.g. image-text pre-training) to connect the pre-training and downstream tasks, and makes the entire model end-to-end trainable. (ii) for sentence generation, we propose *Syntax-Aware Decoder* to dynamically measure the contribution of visual semantic and syntax-related words. Extensive experiments on three widely-used benchmarks (MSVD, MSR-VTT and VATEX) have shown great potential of the proposed \mathcal{D}^2 and surpassed the previous methods by a large margin in the task of video captioning.

1 INTRODUCTION

Video captioning, which aims to understand spatio-temporal relation inside the video and describe it with natural language sentences, is a fundamental research task for multi-modal video-and-language understanding. It becomes an emerging requirement with the rapid emergence of videos in our lives. To generate good captions for videos, it involves not only the understanding of spatio-temporal semantics in videos but also expressing these factors into a natural language. Existing works [1, 46–48] mainly adopt a two-stage framework for video captioning: firstly extracting visual representation of the video using offline feature extractor (3DCNN and object detector), and then decode natural sentences based on these fixed features (Fig 1 (top)).

Despite being reasonable, these methods still suffer from the so-called "*couple*" drawbacks on both video spatio-temporal representation and sentence generation process. For video spatio-temporal representation, "*couple*" means that the learning process of spatio-temporal semantic is restricted in a single model(Fig 1), *i.e.*, 3D CNN,

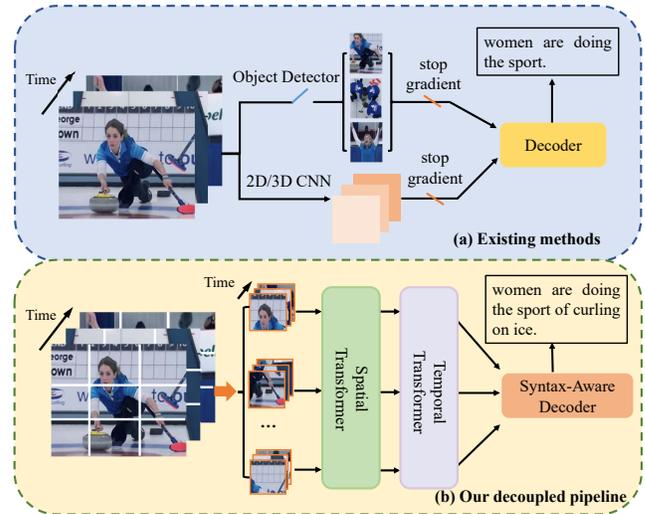


Figure 1: Top: The 3D-CNN or the 2D variant with complex temporal fusion block, along with an detection module to focus on the interested objects, in conventional two-stage video captioning model is generally infeasible for end-to-end finetuning due to the memory constraint on long video. **Bottom:** The proposed \mathcal{D}^2 model overcomes this limit with the unified attention modules to align the language information and the decoupled spatial and temporal features which are represented by the residual coding. Benefiting from the modularized design, we can instantiate each block with dedicated pre-trained model and jointly optimize the entire system. Please refer to Fig 2 for the detailed design.

which brings two main limitations: (i) disconnection in task/pre-train domain: offline feature extractors are often trained on tasks and domains different from the target tasks. For example, 3D CNN is generally trained from pure video data without any textual input on action recognition task, while being applied to video captioning. Large-scale video-text pre-training offers a way to mitigate this issue. However, compared with its image counterpart, the collection of the video-text dataset is much more complex and its noise is also much larger [26]. This makes the video-text pre-trained models are difficult to play a big role in the video captioning tasks. Besides, convolutional kernels are specifically designed to capture short-range spatio-temporal information, they cannot model long-range

dependencies that extend beyond the receptive field. (ii): end-to-end training: due to the memory and computation limitation, it is also infeasible to directly plug these feature extractors into a video captioning framework for end-to-end fine-tuning, causing the disconnection between pre-trained feature and downstream tasks; For sentences generation, "*couple*" means that existing decoding methods pay equal attention to visual semantic ("*woman*" in Fig 1) and syntax-related words ("*are*" in Fig 1) during the whole sentence generation process, making the generation process unreasonable.

To tackle the above drawbacks, we propose \mathcal{D}^2 , a dual-level decoupled pure transformer pipeline for end-to-end video captioning. In terms of video spatio-temporal representation, we *decouple* the learning process into "first-spatial-then-temporal" paradigm. technically, we firstly use a 2D vision transformer to generate a spatial representation for each frame, then, we propose a *Residual-Aware Temporal Block(RATB)* to build the temporal relationship between each frame. This brings us two main advantages over the two limitations mentioned above: (i): compared to the video-text dataset, image-text dataset is easy to access, releasing the potential of applying image-text pre-training to our 2D vision transformer, which is more suitable to multimodal tasks like video captioning. (ii) due to the lightweight nature of the 2D model (as opposed to 3D CNN), we can easily perform end-to-end training, building the connection between pre-trained feature and downstream video captioning tasks. In terms of sentence generation, we *decouple* the generation process of visual semantic and syntax-related words. Concretely, using the syntactic prior provided by the pre-trained language model, our *Syntax-Aware Decoder(SAD)* dynamically measures the contribution of visual features and syntactic prior information for the generation of each word, facilitating a more reasonable and fine-grained video captioning generation.

Experimentally, we conduct a series of ablation studies on different modules for *decouple* spatio-temporal representation learning, as well as modules for *decoupling* visual semantic and syntax-related words generation, gaining insights on the performance of our novel *decouple* pipeline in video captioning task. Our \mathcal{D}^2 , when tested on the three benchmarks(MSRVTT, MSVD and VATEX), outperforms existing methods on all metrics by a large margin. Our main contributions are summarized as follows.

- We propose \mathcal{D}^2 , a novel transformer pipeline, which decouple the process of *video spatio-temporal representation learning* and *sentence generation*. For video spatio-temporal representation learning, our pipeline decouple the previous offline spatio-temporal representation learning into "**first spatial then temporal**" paradigm, addressing the problems of *disconnection in task/pre-train domain* and *end-to-end training*.
- For caption generation, our model decouples the generation of visual semantic and syntax-related words, adaptively measuring its contribution, resulting in a more reasonable and fine-grained video captioning generation process.
- We show that \mathcal{D}^2 surpasses all previous methods for video captioning, achieving a new state-of-the-art on MSVD, MSR-VTT and VATEX benchmarks.

2 RELATED WORK

2.1 Image and Video Representation Learning With Transformer

Inspired by the development of natural language processing, the proposed ViT [12] firstly introduces transformer into image classification and achieves surprising performance. This motivates many researchers to conduct a more in-depth study of the transformer as the backbone network. Different from CNNs, transformers are not limited by the receptive field and can obtain more comprehensive contextual information. Meanwhile, due to the characteristics of the attention mechanism to dynamically generate attention coefficients for different instances, the expressive ability of transformers is also stronger. Considering the above advantages, researchers are increasingly applying transformers to extract feature representations, both in the field of images and videos. For the image domain, multi-scale features [24, 39, 44] are expanded by introducing pyramid structure into the transformer, allowing the transformer-based backbone to better adapt to the downstream vision tasks. Several work [9, 24, 45] also focuses on the balance between attention span and computational overhead by adding local windows or filtering high-value patches.

When considering the video domain, timing information should be added. TimeSformer [4] studies five different variants of space-time attention and suggests a factorized space-time attention for its strong speed-accuracy trade-off. ViViT [2] examines four factorized designs of spatial and temporal attention for the pre-trained ViT model. Similar to our work, they adopt divided attentions on spatial and temporal with two-path transformer models. However, they focus on patch-level attention for video action recognition. We mainly investigate frame-level spatio-temporal semantic representation for video captioning tasks. Video Swin-Transformer [25] designs a video recognition framework by transferring the thought of swin-transformer from space domain to space-time domain. Although this method has achieved certain results, the design of 3D local window still makes the overall framework complex.

The training of backbone networks often needs to be driven by a large amount of data, which usually consumes a lot of computing resources. In this paper, we design an end-to-end framework to avoid the problem of offline characterization, making the pipeline more concise.

2.2 Video Captioning

Over the past years, we have witnessed the great development of video captioning. Before the emergence of neural networks, template-based methods that used the concepts of Verb, Subject and Object (SVO) [3, 20] have become dominant. While in the era of deep learning, a broad collection of methods have been proposed [33, 35] which mainly adopt an encoder-decoder framework. These works tend to extract images' features with CNNs and adopt RNNs to generate descriptions. Venugopalan [35] firstly introduces an LSTM to generate the mean pooled representations over all frames. Wang [37] tries to enhance the quality of generated captions by reproducing the frame features from decoding hidden states. And Buch [5] connects the same output from the two hidden

layers of the opposite direction at a particular time step to further refine the descriptions. However, the original implementation equipped with RNN is difficult to capture long-term dependencies since only the last hidden state takes part in the final generation. The proposed attention mechanism [14, 17, 34, 43] alleviates this deficiency. For each output that the decoder generates, it has access to the entire input sequence at the temporal level. Note that different regions of each video frame also have different contributions to the prediction of the final word. For example, the description object is more important than the background. Therefore, it is also necessary to introduce spatial attention. Li [22] adopts two layers of spatio-temporal dynamic attention for video subtitles. When calculating the spatial attention weight of a specific frame, it will also consider the attention weight of the previous frame. In this way, the spatial attention map is linked across time. Wang [38] tries to learn a model to distinguish the foreground from the background in the video without explicit supervision, and calculates the significance score from the spatial feature map to separate the foreground and background according to the threshold. Recent work typically follows the pipeline where off-the-shelf 2D and 3D is used to extract spatio-temporal feature for video representation. Unlike previous works, we mainly focus on decouple the spatio-temporal representation learning process for better learning.

3 OUR APPROACH

In this section, we introduce the proposed video captioning framework \mathcal{D}^2 , which is illustrated in Figure 2. Compared with existing method, \mathcal{D}^2 decouple the video captioning process both in video representation and caption generation phase. We begin with introducing the *spatial encoder*. And then, the *Residual-Aware Temporal Block*(RATB) and the *Syntax-Aware Decoder*(SAD) are elaborated. The training details will be given at the end of this section.

3.1 Spatial Encoder

Given a video \mathcal{V} , the video captioning task aims to generate a caption $y = \{y_1, \dots, y_T\}$ to describe semantic concepts in \mathcal{V} , where y_t denotes the t -th word in the caption. In the following, we introduce how a frame representation is obtained through a spatial encoder.

The video clips $v_i \in \mathcal{V}$ are represented as a sequence of frames (images) in our paper. Specifically, the video clip v_i is composed of $|v_i|$ sampled frames such that $|v_i| = \{v_i^1, v_i^2, \dots, v_i^{|v_i|}\}$. Unlike previous methods [36, 46–48] which extract clip features using pre-trained CNN [6, 16] or object detector [31]. Our \mathcal{D}^2 model is trained on pixels directly via taking the frames as input in an end-to-end manner.

In order to get the video representation, we first extract the frames from each video clip and then encode them through our spatial encoder to get a sequence of frame features. In this paper, we adopt a 2D ViT-B/16 [12] as our spatial encoder and mainly consider using a sequence of frame representation as video representation. The ViT first extracts non-overlapping image patches, then performs a linear projection to project them into 1D tokens. The transformer architecture is used to model the interaction between each patch of the input frame to get the final representation. We use the output from the [CLS] token as the frame representation. For the input frame sequence of video $v_i = \{v_i^1, v_i^2, \dots, v_i^{|v_i|}\}$, the

generated representation can denote as $Z_i = \{z_i^1, z_i^2, \dots, z_i^{|v_i|}\}$. Benefiting from the modularized design of model, we can instantiate our spatial encoder with dedicated image-text pre-trained model. Specifically, we utilize the recent one-stage image-text pre-trained model ALBEF [21] to initialize our spatial encoder. The impact of different weight initialization strategies is examined in our ablation.

3.2 Residual-Aware Temporal Block

Considering that viewing video representation as a sequence of frame representations may ignore temporal dependency, we propose Residual-Aware Temporal Block(RATB) to build both long-range and Residual level temporal dependency between each frame. Since two successive frames contain content displacement, which reflects the actual actions, we explicitly propose a Residual Attention Mechanism to extend the input frame representation and guide the temporal transformer to encode more motion-related representations. Specifically, we adopt transformed residual of frame representation between adjacent time stamps to describe the motion change, which is formulated as:

$$\mathbf{Z}_d = \delta \left(\left\{ Z_f^1 - Z_f^0, Z_f^2 - Z_f^1, \dots, Z_f^{m-1} - Z_f^{m-2} \right\} + \mathbf{P} \right) \quad (1)$$

where \mathbf{P} is the positional embedding, Z_f^{m-1} and Z_f^m are two adjacent frame representations, δ is 1-layer transformer encoder layer, and \mathbf{Z}_d is the difference representations. We insert difference representations between every frames as below:

$$\mathbf{Z}_n = \left\{ Z_f^0, Z_d^1, Z_f^1, Z_d^2, Z_f^2, \dots, Z_d^{m-1}, Z_f^{m-1} \right\} + \mathbf{P} + \mathbf{T} \quad (2)$$

where \mathbf{Z}_n is the output of our RATB, \mathbf{P} is the positional embedding, \mathbf{T} is the type embedding. With the design of residual-level attention and frame-level temporal attention, atomic actions in short term segments can be contextualized with the rest of the video, thus to be fully disambiguated.

Compared with the 3D-CNN counterpart, which is inherently limited in capturing long-range dependencies by means of aggregation of shorter-range information, our RATB built temporal dependency at both shorter-range and long-range levels.

3.3 Syntax-Aware Decoder

Existing decoding methods in video captioning were suffering from the semantics and syntax coupling drawback. \mathcal{D}^2 solve this problem by dynamic fusing semantic features and syntactic features during decoding. We begin with the syntactic feature extraction and then introduce the semantic and syntactic decoding method.

3.3.1 Syntactic Feature Extraction. Unlike previous works using part-of-speech(POS) [36] tagging tools, we choose a pre-trained language model to provide syntactic prior information. To be specific, in order to get syntactic features during the decoding stage, we build a pretrained 12-layers GPT-based language model to extract syntactic features. The language model was first pre-trained by using a large corpus of documents, and then fine-tuned by using captions sentences only. Some previous works [10] have proven that the pre-trained language model could retain syntactic information inside some of their attention heads. So in that way, we could obtain dense syntactic features containing not only POS tagging

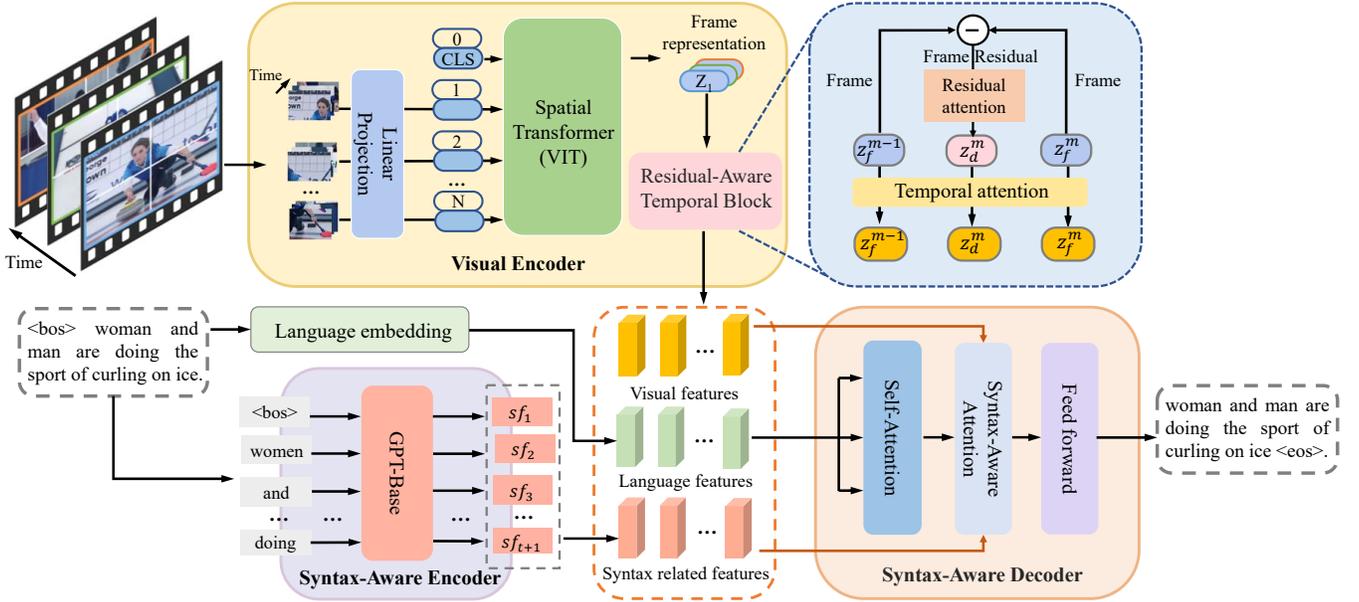


Figure 2: The overall architecture of our proposed \mathcal{D}^2 . \mathcal{D}^2 decouples the learning process in both video and caption side. On the video side, spatial representation is obtained using *Spatial Transformer (ViT)*. With a carefully designed *Residual-Aware Temporal Block*, temporal dependency is built both on *residual* and *long-range* level. On the caption generation side, \mathcal{D}^2 decouple the generation of visual semantic and syntax-related words by a *Syntax-Aware Decoder*, in which semantic and syntactic feature is fused dynamically during the whole generation process. Benefiting from the *neat* and *modularized* design, we can instantiate each module with dedicated pre-trained model and jointly optimize the entire system for task-specific fine-tuning.

information but also the whole syntactic information by using the pretrained language model.

Specifically, given the word sequence $Y = (\langle \text{bos} \rangle, y_1, y_2, \dots, y_N)$, our syntax-aware encoder is asked to predict this offset sequence $\hat{Y} = (\hat{y}_1, \hat{y}_2, \dots, \hat{y}_N, \langle \text{eos} \rangle)$ by one forward process. This entire process can be expressed as follows:

$$sf = \text{GPT}(Y + pos) \quad (3)$$

$$\hat{Y} = \log_ \text{softmax}(FF(sf)) \quad (4)$$

where pos is position embedding, FF is feedforward network, \hat{Y} is the output distribution of predicted words, sf is our syntactic feature. This language model is trained with XE loss. This can be expressed as:

$$sf_t \leftarrow \text{SAE}(Y_{<t}), sf_t \in \mathbb{R}^{d_{\text{model}}} \quad (5)$$

where SAE is our Syntax-Aware Encoder.

3.3.2 Visual Semantic and Syntactic Decoding. After extracting syntactic features, we combine it with the output of our RATB as our candidates during decoding. In our opinion the visual encoder could provide more semantic information, so the decoder could choose semantic or syntactic clues flexibility by using the attention mechanism. Besides, we combine the visual encoder output features, the pretrained language model output features and the current word embedding together as the query vector of the attention module. By this way, we hope the decoder could avoid generating trivial

words. This can be formulated as follows:

$$h_t = \text{Decoder}([Z_n; sf_t], Y_{<t}) \quad (6)$$

$$q_{i,t} = h_t W_i^Q, k_{i,t} = [Z_n; sf_t] W_i^K, v_{i,t} = [Z_n; sf_t] W_i^V \quad (7)$$

$$\text{head}_{i,t} = \text{softmax}(q_{i,t} k_{i,t}^T) v_{i,t} \quad (8)$$

$$\text{head}_i = \text{Concate}(\text{head}_{i,1}, \dots, \text{head}_{i,M}) \quad (9)$$

$$\text{att} = \text{Concate}(\text{head}_1, \dots, \text{head}_h) W^O \quad (10)$$

where Z_n is the output of our RATB, sf_t is the output of syntactic encoder, $q_{i,t}$, $k_{i,t}$, $v_{i,t}$ are the query, key matrix and value matrix for the t time step word in head i respectively, $\text{head}_{i,t}$ is the attention result at t -th timestep, head_i is the attention result in head i , att is the final attention coefficient for sequence generation.

3.4 Training details

The captioning model is typically trained by the cross-entropy loss (XE) given the ground-truth pair (\mathcal{V}, y^*) .

$$L_{XE}(\theta) = - \sum_{t=1}^T \log(p_{\theta}(y_t^* | y_{1:t-1}^*, \mathcal{V})), \quad (11)$$

where θ is the parameters of our model, $y_{1:T}^*$ is the target ground truth sequence.

To address the exposure bias and target mismatch problem in XE, we directly optimize the non-differentiable metric with Self-Critical Sequence Training [32]:

$$L_{RL}(\theta) = -E_{y_{1:T}p_{\theta}} [r(y_{1:T})], \quad (12)$$

where the reward $r(\cdot)$ is the CIDEr-D score.

Besides, following [11], we use the mean of rewards rather than greedy decoding to baseline the reward. The gradient expression for one sample is formulated as:

$$b = \frac{1}{k} \left(\sum_i^k r(y_i) \right), \quad (13)$$

$$\nabla_{\theta} L_{RL}(\theta) \approx -\frac{1}{k} \sum_{i=1}^k \left((r(y_{1:T}^i) - b) \nabla_{\theta} \log p_{\theta}(y_{1:T}^i) \right), \quad (14)$$

where k is the number of the sampled sequences, $y_{1:T}^i$ is the i -th, and b is the mean of the rewards obtained by the sampled sequences.

4 EXPERIMENTS

In this section, we conduct experiments to verify the effectiveness of the proposed methods. Firstly, We introduce three widely-used datasets: MSVD [7], MSR-VTT [42] and recent VATEX [40]. Then implementation details including hyper-parameters and techniques are illustrated. After that, we make comparisons between our methods and the state-of-the-arts. More ablation studies are also discussed in the final part.

4.1 Datasets

4.1.1 MSVD. The MSVD dataset contains 1970 video clips and roughly 80000 English sentences. It was firstly developed in 2010 and used to test and train the translational and paraphrase algorithms. Similar to the prior work [36], we separate the dataset into 1,200 train, 100 validation ,and 670 test videos.

4.1.2 MSR-VTT. The MSR-VTT dataset contains 10,000 video clips from the YouTube website. we follow the standard splits in [42] for fair comparison which separates the dataset into 6,513 training, 497 validation and 2,990 test videos.

4.1.3 VATEX. The VATEX dataset is the most recently released large-scale dataset that contains 41,269 videos. Each video is annotated with 10 English and 10 Chinese descriptions. We utilize English corpora in our experiments. According to the official splits [40], the dataset is divided into 25,991 training, 3,000 validation, and 6,000 public testing.

4.2 Implementation Details.

For the sentences longer than 30 words are truncated (50 for VATEX); the punctuations are removed (for VATEX are retained); all words are converted into lower case.

Our model is trained on pixels directly via taking the frames as input. We divide video into 12 clips. If not otherwise stated, we randomly sample a single frame from each clip for training, and use the middle frame for inference. The impact of different number of clips is examined in ablation. The whole framework is optimized with ADAM [19] optimizer. We set the initial learning rate of the RATB and SAD to 10^{-4} , the spatial encoder to 10^{-5} . We pre-train our 2D visual encoder with the similar spirit as [21]. For the other components, the parameters are randomly initialized with Xavier init. Our training is started with cross entropy optimization. If the

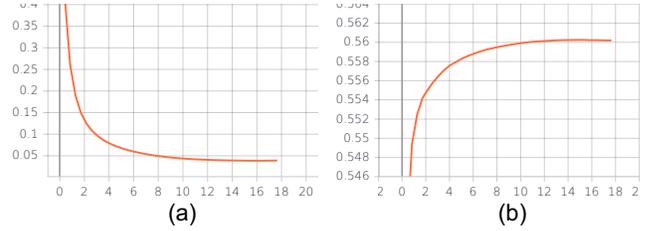


Figure 3: (a): Validation loss during Syntax-Aware Encoder fine-tuning. (b): Accuracy on validation set.

cider value drops for 5 consecutive epochs, it will turn to self-critical sequence training. The training process stops when the cider value drops for 5 consecutive epochs in self-critical sequence training.

We use input image size 256×256 . Our model is implemented in PyTorch [28] and transformers [41]. We set the batch size to 48. The dimension of our spatial encoder is 768. We set the dimension of Residual-Aware Temporal Block and Syntax-Aware Decoder to 512. To avoid over-fitting, we exploit dropout operation after the multi-head self-attention layer and the FFN of each transformer encoder layer. The dropout ratio is set to 0.1 by default.

To guarantee the quality of generated "syntax realted features", we build our *Syntax-Aware Encoder* based on GPT-Base provided by pytorch-transformers [41]. We then fine-tune our *Syntax-Aware Encoder* on caption sentences to obtain syntactic information. Fig 3 shows the validation loss and validation accuracy on VATEX datasets, demonstrating the ability of our *Syntax-Aware Encoder* to encode syntactic information.

4.3 Comparison to State-of-the-Arts

We compare our model with the following state-of-the-art method, which all use offline feature-extractor (either 3DCNN or object detector) for spatio-temporal representation.

- MGSA [?]: LSTM-based model for motion-guided caption generation.
- FCVC [13]: Fully convolutional network for coarse-to-fine caption generation.
- POS-CG [36]: Part-of-speech guided caption generation.
- POS-VCT [18]: This method takes pos for caption generation.
- MARN [29]: This method leverages memory network to capture cross-video contents.
- PMI-CAP [8]: This method learns pairwise modality interactions to better exploit complementary information for each pair of modalities in video.
- SGN [33]: This method uses a semantic grouping network to capture the most discriminating word phrases.
- NACF [?]: This method uses a coarse-to-fine decoding method to better capture visual words in video.
- OA-BTG [46]: This method uses a object-aware aggregation bidirectional temporal graph (OA-BTG) to capture detailed temporal dynamics for salient objects in video.
- GRU-EVE [1]: A LSTM based model which capture spatial-temporal dynamics by Short Fourier Transform.
- STG-KD [27]: This method distills the knowledge of spatial-temporal object interactions from a spatial-temporal graph.

Model	Ref	3DCNN	Detector	MSVD				MSRVTT			
				B@4	M	R	C	B@4	M	R	C
MGSA [?]	AAAI19	✓	✗	53.4	35.0	-	86.7	45.4	28.6	-	50.1
FCVC [13]	AAAI19	✓	✗	53.1	34.8	71.8	79.8	-	-	-	-
POS-CG [36]	ICCV19	✓	✗	52.5	34.1	71.3	88.7	42.0	28.2	61.6	48.7
POS-CG+RL [36]	ICCV19	✓	✗	53.9	34.9	72.1	91.0	41.3	28.7	62.1	53.4
POS-VCT [18]	ICCV19	✓	✗	52.8	36.1	71.8	87.8	42.3	29.7	62.8	49.1
MARN [29]	CVPR19	✓	✗	48.6	35.1	71.9	92.2	40.4	28.1	60.7	47.1
PMI-CAP [8]	ECCV20	✓	✗	54.7	36.4	-	95.2	42.1	28.7	-	49.4
PMI-CAP+Audio [8]	ECCV20	✓	✗	-	-	-	-	43.9	29.5	-	50.6
SGN [33]	AAAI21	✓	✗	52.8	35.5	72.9	94.3	40.8	28.3	60.8	49.5
NACF [?]	AAAI21	✓	✗	55.6	36.2	-	96.3	42.0	28.7	-	51.4
OA-BTG [46]	CVPR19	✓	✓	56.9	36.2	-	90.6	41.4	28.2	-	46.9
GRU-EVE [1]	CVPR19	✓	✓	47.9	35.0	71.5	78.1	38.3	28.4	60.7	48.1
STG-KD [27]	CVPR20	✓	✓	52.2	36.9	73.9	93.0	40.5	28.3	60.9	47.1
SAAT [49]	CVPR20	✓	✓	46.5	33.5	69.4	81.0	40.5	28.2	60.9	49.1
ORG-TRL [48]	CVPR20	✓	✓	54.3	36.4	73.9	95.2	43.6	28.8	62.1	50.9
OPEN-BOOK [47]	CVPR21	✓	✓	-	-	-	-	42.8	29.3	61.7	52.9
O2NA [23]	ACL21	✓	✓	55.4	37.4	74.5	96.4	41.6	28.5	62.4	51.1
\mathcal{D}^2	-	✗	✗	56.9	38.4	75.1	99.2	44.5	30.0	63.3	56.3

Table 1: Performance comparisons on MSVD and MSRVTT benchmarks. B@4,C,M, and R denote BLEU@4, CIDEr-D, METEOR, and ROUGE_L, respectively. - means not available.

- SAAT [49]: This method learns actions by simultaneously referring to the subject and video dynamics.
- ORG-TRL [48]: This method uses a object relation graph to build the object interactions.
- OPEN-BOOK[47]:They use “retrieval and copy” pipeline to help caption generation.
- O2NA [23]: This method uses a Object-Oriented Non-Autoregressive approach (O2NA).

Model	B@4	M	R	C
VATEX (ICCV19) [40]	28.7	21.9	47.2	45.6
ORG-TRL (CVPR20) [48]	32.1	22.2	48.9	49.7
NSA (CVPR20) [15]	31.0	22.7	49.0	57.1
OPENBOOK (CVPR21) [47]	33.9	23.7	50.2	57.5
\mathcal{D}^2	35.1	25.1	51.3	60.9

Table 2: Performance comparisons on VATEX benchmark. B@4, C, M, R denote BLEU@4, CIDEr-D, METEOR, and ROUGE_L, respectively.

Table 1 reports the video caption performances of different models on the MSVD and MSR-VTT datasets. We can see that our \mathcal{D}^2 consistently exhibits better performance than the others. To be specific, on MSVD, our \mathcal{D}^2 surpasses all the other methods for all metrics by a large margin even for the strongest competitor O2NA [23]. The CIDEr-D score of our method reaches 99.2%, which advances O2NA [23] by 2.9%. As for MSRVTT, our \mathcal{D}^2 achieves better performance than all the other methods for all metrics. Compared with the strongest competitor OPENBOOK [47], we advance it by 6.4%(56.3 vs. 52.9). The boost of performance on both MSVD and

MSRVTT demonstrates the advantages of our \mathcal{D}^2 which uses *decoupled* pipeline on both video spatio-temporal representation learning and sentence generation. Importantly, since CIDEr-D weights the n-grams that are relevant to the video content, demonstrating that our model generates more video-relevant captions.

Besides, we also report the results of our model on the public test set of the recent published VATEX dataset as shown in Table 2. Compared with these SOTA methods, our model achieves the best performance on *all metrics*. Specifically, our \mathcal{D}^2 reaches 60.9% in terms of CIDEr-D, which advances the strongest competitor OPEN-BOOK [47] by 5.9%, proving the effectiveness of our *decoupled* pipeline.

4.4 Ablation Analysis

In this section, we conduct several ablation studies on the VATEX and MSRVTT datasets to determine some hyper-parameter and demonstrate the effectiveness of our proposed module. Since our method directly takes video frame as input and does end-to-end training; before doing the experiments about the main module, in section 4.4.1, we first give the experiments about sample rates of input clips. Then, we give the main ablation experiments in section 4.4.2 about our proposed *Residual-Aware Temporal Block(RATB)* and *Syntax-Aware Decoder(SAD)*. To fully exploit the effectiveness of the module in *Residual-Aware Temporal Block(RATB)*, we did some experiments and give results in section 4.4.3. Finally, we give the ablation experiments results about image-text pre-training and end-to-end training in section 4.4.4 and section 4.4.5 respectively.

4.4.1 Ablation on number of input clips. In order to investigate the impact of the different number of input clips, as shown in Table 3, we set the number of clips to 8, 12, 16 for our model respectively.



GT: a train is traveling through the european country side

Baseline: a train is moving on the field

Ours: a train is moving on the road in the forest



GT: a pair of young people play a challenging game of ping pong

Baseline: two people are playing games

Ours: two people are playing ping pong on the table



GT: a black suit man is speaking from a studio

Baseline: a man is talking to someone

Ours: a man in a suit is talking on the news



GT: a woman mixing various ingredients in a bowl together

Baseline: a woman is cooking

Ours: a woman is mixing ingredients into a bowl

Figure 4: Visualization of generations on MSR-VTT with GT and our \mathcal{D}^2 . Compared with baseline, our \mathcal{D}^2 can generate more richer and descriptive captions.

Number of clips	MSRVTT		VATEX	
	BLEU@4	CIDEr-D	BLEU@4	CIDEr-D
8	44.1	55.5	34.2	59.2
12	44.5	56.3	35.1	60.9
16	44.3	56.5	34.3	60.9

Table 3: Ablation on number of input clips.

The performance cannot increase with more clips and we use 12 for our experiments.

Method	Module		MSRVTT		VATEX	
	RATB	SAD	B@4	C	B@4	C
Baseline			43.3	52.1	33.9	55.2
Baseline	✓		43.9	54.2	34.2	57.3
Baseline		✓	44.1	53.9	34.4	57.1
Baseline	✓	✓	44.5	56.3	35.1	60.9

Table 4: Ablation of main design of our \mathcal{D}^2 . Spatial encoder together with vanilla decoder is selected as our baseline.

4.4.2 Ablation of main design of our \mathcal{D}^2 . We study the benefits of RATB and SAD of our method. Specifically, the spatial encoder with image-text pre-training and vanilla transformer-based decoder are combined as our baseline model. As shown in Table 4, our *Residual-Aware Temporal Block*(RATB) boost CIDEr score by 1.9%(2.1%) on VATEX(MSRVTT) dataset, demonstrating the importance of building temporal dependency for video frames. Compared Row 1 and Row 3 in Table 4, we observe that our *Syntax-Aware Decoder*(SAD) boost CIDEr-D score from 55.2(52.1) to 57.1(53.9) on VATEX(MSRVTT) dataset, proving the effectiveness of decouple the sentence generation. Combining the above two modules, the best performance(60.9 on VATEX; 56.3 on MSRVT) is achieved on both VATEX and MSRVT datasets.

Method	Module		MSRVTT		VATEX	
	RA	TA	B@4	C	B@4	C
Baseline			42.8	52.1	33.9	55.2
Baseline	✓		43.3	53.2	34.3	56.1
Baseline		✓	43.1	53.5	34.0	56.4
Baseline	✓	✓	43.9	54.2	34.2	57.3

Table 5: Ablation of Residual-Aware Temporal Block. RA, TA denote Residual-Level Attention, Temporal-Level Attention, respectively. Spatial encoder together with vanilla decoder is selected as our baseline.

Init Method	MSRVTT		VATEX	
	BLEU@4	CIDEr-D	BLEU@4	CIDEr-D
Resnet152(cls)	27.4	22.6	17.8	17.2
ViT-B/16(cls)	40.2	49.2	29.1	51.4
ViT-B/16(image-text)	43.3	52.1	33.9	55.2

Table 6: Ablation study on different type of spatial encoder and different init method.

4.4.3 Ablation of Residual-Aware Temporal Block. Compared with existing methods, our Residual-Aware Temporal Block builds temporal dependency between frames in both short-range and long-range. Specifically, *Residual-Level Attention*(RA) is adopted to encode short-range(atomic action) relation. Together with *Temporal-Level Attention*, atomic actions in short-range segments can be contextualized with the rest of the video. Table 5 presents the effect of two sub-modules. The results show that RA or TA is available as an effective sub-module and combining these two sub-modules boost the CIDEr-D score by 0.9 or 1.2(1.1 or 1.4) on VATEX(MSRVTT).

4.4.4 Ablation of image-text pre-training. Due to the better accessibility of image-text datasets (as opposed to the video-text dataset,

which is noisy and hard to collect), it is trivial for our spatial encoder to enrich visual semantics via pre-training on image-text pairs (e.g., MSCOCO). In this section, we study the effect of different types of spatial encoder and its init method. Results are summarized in Table 6. We observed that there is a huge performance boost (row1 and row2) when replacing resnet152(cls) with ViT(cls), demonstrating that ViT has much more frame representation ability than the resnet. We hypothesize that it is mostly because the information loses via meanpooling operation [30]. Besides, when we conduct image-text pre-training on image-text dataset to our spatial encoder, the performance continues boosting, demonstrating that image-text pre-trained model benefits the video captioning task for its semantic spatial representation.

for the generation of each word, obtaining a more reasonable and fine-grained video captioning generation. We evaluate our methods on MSVD, MSR-VTT and VATEX, all leading to the best performance on all metrics. This sets the new state-of-the-art and our model could be the new backbone model for video captioning.

Method	MSRVTT			VATEX		
	BLEU@4	METEOR	CIDEr	BLEU@4	METEOR	CIDEr
fix	43.6	29.2	54.8	34.3	24.4	58.6
e2e	44.5	30.0	56.3	35.1	25.1	60.9

Table 7: Ablation of the benefit of end-to-end training.

4.4.5 Ablation of End-to-End Training. One of the main difference between our method and previous is that we decouple the process of spatio-temporal learning: we use a 2D transformer to obtain the spatial contextual information then build temporal correlation using our designed module, giving us the benefit of end-to-end training (as opposed to 3DCNN which is hard to fine-tune). We conduct the ablation on both MSRVTT and VATEX datasets with cross-entropy loss. As it can be seen from Tab 7, end-to-end training boosts all metrics on MSRVTT (1.5 CIDEr score) and VATEX (2.3 CIDEr score), proving the effectiveness of our pipeline.

4.5 Qualitative Analysis

We show some examples in Fig. 4. It can be seen, the content of captions generated by our model is richer than the baseline model, and more activity associations are involved. For instance, the example at down-left shows that the baseline model can only understand the general meaning of the video. By contrast, our model can recognize more detailed activities (e.g. playing ping pong). The rest of the examples have similar characteristics.

5 CONCLUSION

In this paper, we propose \mathcal{D}^2 , a dual-level decoupled pure transformer pipeline for end-to-end video captioning. We address the "couple" drawbacks on both video spatio-temporal representation learning and sentence generation. For video spatio-temporal representation, we present "first-spatial-then-temporal" paradigm. technically, we use a 2D vision transformer to generate a spatial representation for each frame and build both short and long range temporal dependency with our proposed *Residual-Aware Temporal Block (RATB)*. For sentence generation, we *decouple* the generation process of visual semantic and syntax-related words. Specifically, utilizing the syntactic prior provided by pre-trained language model, our proposed *Syntax-Aware Decoder (SAD)* dynamically measures the contribution of visual features and syntactic prior information

REFERENCES

- [1] Nayyer Aafaq, Naveed Akhtar, Wei Liu, Syed Zulqarnain Gilani, and Ajmal Mian. 2019. Spatio-temporal dynamics and semantic attribute enriched visual encoding for video captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 12487–12496.
- [2] Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lučić, and Cordelia Schmid. 2021. Vivit: A video vision transformer. *arXiv preprint arXiv:2103.15691* (2021).
- [3] Andrei Barbu, Alexander Bridge, Zachary Burchill, Dan Coroian, Sven Dickinson, Sanja Fidler, Aaron Michaux, Sam Mussman, Siddharth Narayanaswamy, Dhaval Salvi, et al. 2012. Video in sentences out. *arXiv preprint arXiv:1204.2742* (2012).
- [4] Gedas Bertasius, Heng Wang, and Lorenzo Torresani. 2021. Is Space-Time Attention All You Need for Video Understanding? *arXiv preprint arXiv:2102.05095* (2021).
- [5] Shyamal Buch, Victor Escorcía, Chuanqi Shen, Bernard Ghanem, and Juan Carlos Niebles. 2017. Sst: Single-stream temporal action proposals. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*. 2911–2920.
- [6] Joao Carreira and Andrew Zisserman. 2017. Quo vadis, action recognition? a new model and the kinetics dataset. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 6299–6308.
- [7] David Chen and William B Dolan. 2011. Collecting highly parallel data for paraphrase evaluation. In *Proceedings of the 49th annual meeting of the association for computational linguistics: human language technologies*. 190–200.
- [8] Shaoxiang Chen, Wenhao Jiang, Wei Liu, and Yu-Gang Jiang. 2020. Learning modality interaction for temporal sentence localization and event captioning in videos. In *European Conference on Computer Vision*. Springer, 333–351.
- [9] Xiangxiang Chu, Zhi Tian, Yuqing Wang, Bo Zhang, Haibing Ren, Xiaolin Wei, Huaxia Xia, and Chunhua Shen. 2021. Twins: Revisiting the Design of Spatial Attention in Vision Transformers. In *Advances in neural information processing systems*.
- [10] Kevin Clark, Urvashi Khandelwal, Omer Levy, and Christopher D Manning. 2019. What does bert look at? an analysis of bert’s attention. *arXiv preprint arXiv:1906.04341* (2019).
- [11] Marcella Cornia, Matteo Stefanini, Lorenzo Baraldi, and Rita Cucchiara. 2020. Meshed-memory transformer for image captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10578–10587.
- [12] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiuhua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. 2020. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929* (2020).
- [13] Kun Cheng Fang, Lian Zhou, Cheng Jin, Yuejie Zhang, Kangnian Weng, Tao Zhang, and Weiguo Fan. 2019. Fully convolutional video captioning with coarse-to-fine and inherited attention. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 33. 8271–8278.
- [14] Lianli Gao, Zhao Guo, Hanwang Zhang, Xing Xu, and Heng Tao Shen. 2017. Video captioning with attention-based LSTM and semantic consistency. *IEEE Transactions on Multimedia* 19, 9 (2017), 2045–2055.
- [15] Longteng Guo, Jing Liu, Xinxin Zhu, Peng Yao, Shichen Lu, and Hanqing Lu. 2020. Normalized and geometry-aware self-attention network for image captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10327–10336.
- [16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.
- [17] Chiori Hori, Takaaki Hori, Teng-Yok Lee, Ziming Zhang, Bret Harsham, John R Hershey, Tim K Marks, and Kazuhiko Sumi. 2017. Attention-based multimodal fusion for video description. In *Proceedings of the IEEE international conference on computer vision*. 4193–4202.
- [18] Jingyi Hou, Xinxiao Wu, Wentian Zhao, Jiebo Luo, and Yunde Jia. 2019. Joint syntax representation learning and visual cue translation for video captioning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 8918–8927.
- [19] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).
- [20] Atsuhiko Kojima, Takeshi Tamura, and Kunio Fukunaga. 2002. Natural language description of human activities from video images based on concept hierarchy of actions. *International Journal of Computer Vision* 50, 2 (2002), 171–184.
- [21] Junnan Li, Ramprasaath R Selvaraju, Akhilesh Deepak Gotmare, Shafiq Joty, Caiming Xiong, and Steven Hoi. 2021. Align before Fuse: Vision and Language Representation Learning with Momentum Distillation. *arXiv preprint arXiv:2107.07651* (2021).
- [22] Xuelong Li, Bin Zhao, Xiaoqiang Lu, et al. 2017. MAM-RNN: Multi-level Attention Model Based RNN for Video Captioning. In *IJCAI*, Vol. 2017. 2208–2214.
- [23] Fenglin Liu, Xuancheng Ren, Xian Wu, Bang Yang, Shen Ge, and Xu Sun. [n. d.]. O2NA: An Object-Oriented Non-Autoregressive Approach for Controllable Video Captioning. ([n. d.]).
- [24] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. 2021. Swin transformer: Hierarchical vision transformer using shifted windows. *arXiv preprint arXiv:2103.14030* (2021).
- [25] Ze Liu, Jia Ning, Yue Cao, Yixuan Wei, Zheng Zhang, Stephen Lin, and Han Hu. 2021. Video swin transformer. *arXiv preprint arXiv:2106.13230* (2021).
- [26] Antoine Miech, Jean-Baptiste Alayrac, Lucas Smaira, Ivan Laptev, Josef Sivic, and Andrew Zisserman. 2020. End-to-end learning of visual representations from uncurated instructional videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 9879–9889.
- [27] Boxiao Pan, Haoye Cai, De-An Huang, Kuan-Hui Lee, Adrien Gaidon, Ehsan Adeli, and Juan Carlos Niebles. 2020. Spatio-temporal graph for video captioning with knowledge distillation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10870–10879.
- [28] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. 2019. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems* 32 (2019), 8026–8037.
- [29] Wenjie Pei, Jiyuan Zhang, Xiangrong Wang, Lei Ke, Xiaoyong Shen, and Yu-Wing Tai. 2019. Memory-attended recurrent network for video captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 8347–8356.
- [30] Maithra Raghu, Thomas Unterthiner, Simon Kornblith, Chiyuan Zhang, and Alexey Dosovitskiy. 2021. Do Vision Transformers See Like Convolutional Neural Networks? *arXiv preprint arXiv:2108.08810* (2021).
- [31] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. 2016. Faster R-CNN: towards real-time object detection with region proposal networks. *IEEE transactions on pattern analysis and machine intelligence* 39, 6 (2016), 1137–1149.
- [32] Steven J Rennie, Etienne Marcheret, Youssef Mroueh, Jerret Ross, and Vaibhava Goel. 2017. Self-critical sequence training for image captioning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 7008–7024.
- [33] Hobin Ryu, Sunghun Kang, Haeyong Kang, and Chang D Yoo. 2021. Semantic Grouping Network for Video Captioning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 35. 2514–2522.
- [34] Yunbin Tu, Xishan Zhang, Bingtao Liu, and Chenggang Yan. 2017. Video description with spatial-temporal attention. In *Proceedings of the 25th ACM international conference on Multimedia*. 1014–1022.
- [35] Subhashini Venugopalan, Huijuan Xu, Jeff Donahue, Marcus Rohrbach, Raymond Mooney, and Kate Saenko. 2014. Translating videos to natural language using deep recurrent neural networks. *arXiv preprint arXiv:1412.4729* (2014).
- [36] Bairui Wang, Lin Ma, Wei Zhang, Wenhao Jiang, Jingwen Wang, and Wei Liu. 2019. Controllable video captioning with pos sequence guidance based on gated fusion network. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2641–2650.
- [37] Bairui Wang, Lin Ma, Wei Zhang, and Wei Liu. 2018. Reconstruction network for video captioning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 7622–7631.
- [38] Huiyun Wang, Youjiang Xu, and Yahong Han. 2018. Spotting and aggregating salient regions for video captioning. In *Proceedings of the 26th ACM international conference on Multimedia*. 1519–1526.
- [39] Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. 2021. Pyramid Vision Transformer: A Versatile Backbone for Dense Prediction Without Convolutions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*. 568–578.
- [40] Xin Wang, Jiawei Wu, Junkun Chen, Lei Li, Yuan-Fang Wang, and William Yang Wang. 2019. Vatec: A large-scale, high-quality multilingual dataset for video-and-language research. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 4581–4591.
- [41] Thomas Wolf, Julien Chaumond, Lysandre Debut, Victor Sanh, Clement Delangue, Anthony Moi, Pierric Cistac, Morgan Funtowicz, Joe Davison, Sam Shleifer, et al. 2020. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*. 38–45.
- [42] Jun Xu, Tao Mei, Ting Yao, and Yong Rui. 2016. Msr-vtt: A large video description dataset for bridging video and language. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 5288–5296.
- [43] Li Yao, Atousa Torabi, Kyunghyun Cho, Nicolas Ballas, Christopher Pal, Hugo Larochelle, and Aaron Courville. 2015. Describing videos by exploiting temporal structure. In *Proceedings of the IEEE international conference on computer vision*. 4507–4515.
- [44] Li Yuan, Yunpeng Chen, Tao Wang, Weihao Yu, Yujun Shi, Zi-Hang Jiang, Francis E.H. Tay, Jiashi Feng, and Shuicheng Yan. 2021. Tokens-to-Token ViT: Training Vision Transformers From Scratch on ImageNet. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*. 558–567.
- [45] Li Yuan, Qibin Hou, Zihang Jiang, Jiashi Feng, and Shuicheng Yan. 2021. Volo: Vision outlooker for visual recognition. *arXiv preprint arXiv:2106.13112* (2021).
- [46] Junchao Zhang and Yuxin Peng. 2019. Object-aware aggregation with bidirectional temporal graph for video captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 8327–8336.

- [47] Ziqi Zhang, Zhongang Qi, Chunfeng Yuan, Ying Shan, Bing Li, Ying Deng, and Weiming Hu. 2021. Open-book Video Captioning with Retrieve-Copy-Generate Network. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 9837–9846.
- [48] Ziqi Zhang, Yaya Shi, Chunfeng Yuan, Bing Li, Peijin Wang, Weiming Hu, and Zheng-Jun Zha. 2020. Object relational graph with teacher-recommended learning for video captioning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 13278–13288.
- [49] Qi Zheng, Chaoyue Wang, and Dacheng Tao. 2020. Syntax-aware action targeting for video captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 13096–13105.