# EVOQUER: Enhancing Temporal Grounding with Video-Pivoted Back Query Generation

**Yanjun Gao[1], Lulu Liu[1], Jason Wang[1], Xin Chen[2], Huayan Wang[2], Rui Zhang[1]**

Pennsylvania State University[1], Kwai Inc[2]

{yug125,lzl5409,jjw6188,rmz5227}@psu.edu,
xinchen.hawaii@gmail.com, wanghuayan@kuaishou.com

## Abstract

Temporal grounding aims to predict a time interval of a video clip corresponding to a natural language query input. In this work, we present EVOQUER, a temporal grounding framework incorporating an existing text-to-video grounding model and a video-assisted query generation network. Given a query and an untrimmed video, the temporal grounding model predicts the target interval, and the predicted video clip is fed into a video translation task by generating a simplified version of the input query. EVOQUER forms closed-loop learning by incorporating loss functions from both temporal grounding and query generation serving as feedback. Our experiments on two widely used datasets, Charades-STA and ActivityNet, show that EVOQUER achieves promising improvements by 1.05 and 1.31 at R@0.7. We also discuss how the query generation task could facilitate error analysis by explaining temporal grounding model behavior.

## 1 Introduction

Temporal grounding aims to find the time interval in an untrimmed video that expresses the same meaning as a natural language query. It locates the video content that semantically corresponds to a natural language query, addressing the temporal, semantic alignment between language and vision. It is broadly applicable in many tasks such as visual storytelling (Lukin et al., 2018; Huang et al., 2016), video caption generation (Krishna et al., 2017; Long et al., 2018), and video machine translation (Wang et al., 2019b).

Recent work on temporal grounding has achieved significant progress (Mun et al., 2020; Chen and Jiang, 2019; Chen et al., 2018; Zhang et al., 2019; Gao et al., 2017). They emphasize modeling the semantic mapping of verbs and nouns in the text query to visual clues such as actions and objects that indicate the candidate time intervals. However, most of them only employ a uni-direction
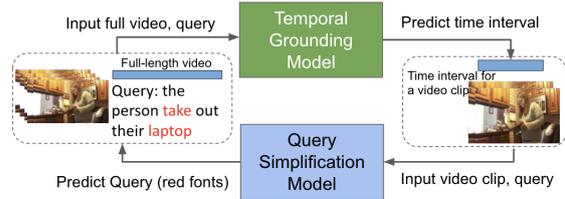


Figure 1: An overview of EVOQUER as a closed-loop system pipeline.

single-task learning flow. To strengthen the learning and facilitate error analysis, we explore the possibility of enhancing the temporal grounding model with related tasks. To this end, we borrow the idea of feedback-error-learning from control theory and computational neuroscience (Kawato, 1990; Gomi and Kawato, 1993). Using a closed-loop system, the control network learns to correct its error from feedback and gains stronger supervision to stimulate learning. We investigate if a temporal grounding model can be improved by incorporating another network to generate feedback in a closed-loop learning fashion.

We propose a novel framework, EVOQUER (*En*hancing Temporal Grounding with *VideO*-Pivoted Back *QUER*y Generation), integrating a text-to-video and a video-to-text flow, as shown in Figure 1. Specifically, we adapt a video-pivoted query simplification task that simplifies the query to shorter phrases with verbs and noun phrases only. Instead of re-generating the full queries, the query simplification task is smaller in the problem size and thus could serve as an auxiliary task to assist the main task. Furthermore, we incorporate visual pivots in query simplification to provide more fine-grained semantic discrepancy associated with words (Chen et al., 2019b; Lee et al., 2019). The pipeline pairs a state-of-the-art temporal grounding model LGI (Mun et al., 2020) with a video machine translation model (Wang et al., 2019b) for query simplification. Given a query and an untrimmed video, the pipeline predicts the
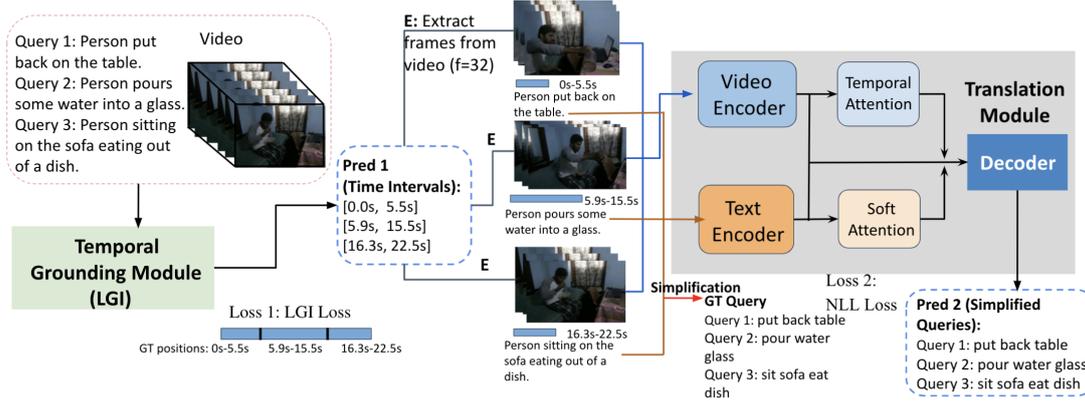
Figure 2: Our EVOQUER framework combines a LGI model for temporal grounding and a translation module that outputs simplified queries.

time interval, feeds the predicted video clips with the original query to the translation model, and generates a simplified query. EVOQUER generates two losses individually from the temporal grounding and query simplification task, and combines the loss to update all network components, giving stronger supervision signals. On two temporal grounding datasets, Charades-STA (Gao et al., 2017) and ActivityNet (Krishna et al., 2017), EVOQUER outperforms the original model by 1.05 and 1.31 at R@0.7, demonstrating its effectiveness. Our analysis of query simplification output indicates that this video-to-text auxiliary task casts light on explaining temporal grounding model behavior.

## 2 Related Work

Existing temporal grounding approaches can be split into three categories: strongly supervised (Anne Hendricks et al., 2017; Gao et al., 2017; Liu et al., 2018; Chen et al., 2018, 2019a; Chen and Jiang, 2019; Ge et al., 2019; Ghosh et al., 2019; Zhang et al., 2019; Yuan et al., 2019; Mun et al., 2020; Rodriguez et al., 2020), weakly supervised (Tan et al., 2021), and reinforcement learning (Wang et al., 2019a; He et al., 2019). The primary example of strongly supervised approaches is LGI (Mun et al., 2020) that achieves state-of-the-art performance on the Charades-STA and ActivityNet datasets. It uses word-level and sentence-level attention to predict time intervals. Within our closed-loop framework, we utilize this LGI algorithm as a black box to achieve the best performance on supervised temporal grounding. Recent progress also uses transformers with language and vision pretraining (Radford et al., 2021; Lei et al., 2021; Luo et al., 2021).

Another similar task is text-to-video moment retrieval that focuses on grounding between query and video (Xu et al., 2019; Lin et al., 2020; Liu et al., 2018), framing the task as retrieving video frames, slightly different from temporal grounding. Furthermore, our framework is also relevant to video captioning which aims to generate a description of text given a video (Das et al., 2013; Yao et al., 2015; Venugopalan et al., 2015a,b; Xu et al., 2015; Zhou et al., 2019, 2018a). Recent developments have utilized end-to-end transformer models for video captioning (Zhou et al., 2018b).

## 3 EVOQUER Framework

We design a closed-loop framework for temporal grounding such that the model receives (1) supervision in predicting time intervals and (2) feedback from the output video features extracted from the prediction. To achieve this, EVOQUER involves two components: a temporal grounding module and a translation module. The temporal grounding module predicts time intervals given an untrimmed video and a query. The translation module takes input from queries and video features trimmed by the predicted intervals, and outputs a simplified query with only verbs and nouns. We use the LGI model (Mun et al., 2020) for temporal grounding, which achieves state-of-the-art performance using supervised learning. For query simplification, we use the video machine translation framework VMT (Wang et al., 2019b) whose performance is competitive in video-assisted bilingual translation.

Our pipeline is presented in Figure 2. The input to the framework is an untrimmed video and a set of queries. Following Mun et al. (2020), we use I3D frame-based features for video representation and an embedding layer inside a text encoder for word representation. Given the video features

and queries, LGI predicts time intervals with the content corresponding to a given query. Next, we extract frames from videos trimmed by the predicted interval to represent the content of the video clip. To maintain the continuity of the content, we extract 32 frames per video clip in a way that the content of the trimmed videos is evenly distributed across all 32 frames. Since the camera used captures 24 frames per second, a 32-frame video roughly spans 1.3 seconds. We feed the extracted video features and input query into a translation module consisting of two biLSTM-based encoders and an LSTM-based decoder with attention. Video hidden states and text hidden states are sent individually to two attention modules, while being concatenated into one vector representation and sent to the decoder as initial hidden states. In the attention network, temporal attention is learned through video features, and soft attention through query hidden states. The attention is fed into the decoder as context representation.

Instead of learning to decode the original query, we want the model to focus on the words that distinguish the video content: verbs and nouns. In the Charades dataset, annotators tend to use various verb tenses when describing the video activities. For example, both *"closes the door"* and *"closing the door"* are used on the same video content. Therefore, we lemmatize the words, label the query with part-of-speech (POS) tags, and extract verbs and nouns as simplified versions of the queries. The decoder learns to predict simplified queries and computes a negative log-likelihood (NLL) loss at the end of the decoding. Finally, we combine the NLL loss from query simplification with the LGI loss from time interval prediction to train the networks jointly in an end-to-end fashion.

## 4 Experiments and Results

**Datasets** We evaluate our framework on Charades-STA (Gao et al., 2017) and ActivityNet (Krishna et al., 2017), two widely used benchmark data sets for temporal grounding. We follow the dataset setting in (Mun et al., 2020), where both datasets are set with train/valid/test as 50%, 25%, and 25% respectively. The dataset statistics are reported in Table 1.

The two datasets vary greatly on most of the statistics. We think ActivityNet is a more challenging dataset, as it requires the decoder in EVOQUER to predict correct words from a much bigger vo-

| Dataset | Charades-STA | ActivityNet |
|---|---|---|
| Num Queries | 27,847 | 71,957 |
| Num Videos | 9,848 | 20,000 |
| Avg Video Len (Sec) | ~30 | ~120 |
| Input Query $|V|$ | 1,140 | 11,125 |
| Simpl. Query $|V|$ | 560 | 5,946 |
| Simpl. #Tks per Query | 2.31 | 4.12 |

Table 1: Statistics contrasting Charades-STA and ActivityNet and the impact of simplification on each dataset.

| Data | Model | R@0.3 | R@0.5 | R@0.7 | mIoU |
|---|---|---|---|---|---|
| char | LGI | 71.54 | **58.08** | 34.68 | 50.28 |
| | EVOQUER | **71.57** | 57.81 | **35.73** | **50.48** |
| anet | LGI | 57.76 | 40.38 | 22.60 | 40.65 |
| | EVOQUER | **59.21** | **42.02** | **23.91** | **41.61** |

Table 2: Results on Charades-STA (char) and ActivityNet (anet) from the LGI model and EVOQUER.

cabulary whose size is almost 10 times bigger than Charades-STA. In our experiments on ActivityNet, we find that the EVOQUER converges at a higher NLL loss than Charades-STA and it fails to produce good quality simplified queries, which we suspect is attributed to the harder decoding task.

**Evaluation metrics** We adopt two conventional temporal grounding metrics: *R@tIoU* measuring recall at different thresholds (0.3, 0.5, and 0.7) for temporal intervals between ground truth and prediction; *mIoU* reporting the average of temporal interval recall from all threshold levels. For query simplification, we evaluate the predicted queries with two metrics. Jaccard similarity measures intersection over union between words in ground truth and in prediction. Since it does not penalize for duplicated words, Jaccard similarity gives us a rough estimation for the quality of translation output. BLEU (Papineni et al., 2002) is a standard evaluation metric for machine translation that measures n-gram word overlap. Most of the simplified queries are two-word length, thus we report BLEU unigram and bigram.

**Temporal grounding results** Table 2 presents results on the Charades-STA and ActivityNet test set from a re-trained LGI model and EVOQUER models.[1] Compared to LGI, EVOQUER shows improvement on R@0.7 and mIoU, especially 1.05 and 1.31 on R@0.7, the hardest threshold for temporal interval overlap.

Table 3 presents statistics of samples where our model show improvements and drops compared to

---

[1]Using the codes from the author's GitHub and the parameters presented in the original paper, we train the LGI model on Charades-STA train set from scratch. We suspect the difference between our replication and results presented in the paper is attributed to initialization.

| | Both >= R@0.3 | | | Both |
| | EVOQUER ↑ | EVOQUER ↓ | Same | <R@0.3 |
|---|---|---|---|---|
| char | 441 | 362 | 1,347 | 777 |
| anet | 4,268 | 3,124 | 8,074 | 10,538 |

Table 3: Counts of samples that are scored by R@tIoU with four categories from comparison between EVOQUER and LGI model. Three of the categories are from samples where both models achieve recall >= threshold 0.3: samples that are improved (EVOQUER ↑), samples with performance drops (EVOQUER ↓), and equal performance with at least R@0.3 (Same). The fourth category is when both perform below R@0.3 (Both <R@0.3).

LGI. We divide the samples into four categories according to their recall: when EVOQUER ranks in a higher threshold than the LGI (e.g. 0.7 vs 0.5), when EVOQUER ranks lower than the LGI, when both have the same recalls that are at least R@0.3, and when both scores rank below R@0.3. EVO-QUER shows 441 samples on Charades-STA and 4,268 samples of ActivityNet that are above 0.3 recall threshold, with 79 and 1,144 samples of absolute improvement. There remains a large number of samples to be improved (Both <R@0.3)). Recall that query simplification task serves as an auxiliary task for the temporal grounding model thus the performance upper bound of EVOQUER could be limited by LGI. We consider the EVOQUER improvements promising after seeing a fair number of samples being improved.

**Translation output analysis** We evaluate translation outputs using BLEU and Jaccard scores between ground truths and predictions. On Charades-STA test set, EVOQUER achieved 51.98 Jaccard Similarity, 53.04 on BLEU unigram, and 42.47 on BLEU bigram. We will show how the query simplification output helps the error analysis of temporal grounding model in case study. On ActivityNet, we fail at training EVOQUER to generate reasonable simplified queries. We suspect two reasons for this failure: 1) the decoder vocabulary on ActivityNet is much larger than Charades-STA, which makes the simplification task harder, as we mentioned previously; 2) the current design of translation module is too simple to handle features from the longer predicted time interval and the input queries. Nonetheless, the results on Charades-STA indicates the potential benefit of our framework, shown in case studies next.

**Case study** We show output examples of predicted intervals and simplified queries to understand the model performance. Figure 3 shows two video clips trimmed by the ground truth interval,



Input query: person *closing* the *door* to the entryway (Simplified query: close door)
Model prediction: Evoquer: close door book door

Input query: a person *holding* a *glass opens* the *refrigerator* (Simplified query: hold glass open refrigerator)
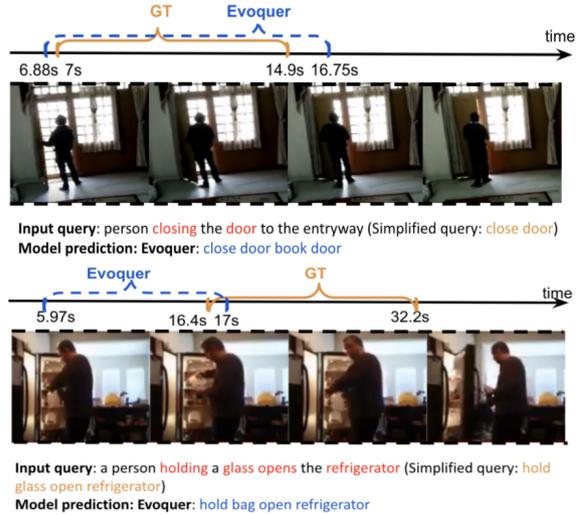Model prediction: Evoquer: hold bag open refrigerator

Figure 3: Two example video clips trimmed by ground truth intervals. In the first example (top), EVOQUER successfully predicts time interval and simplified queries as ground truth. In the second example (bottom), EVOQUER fails to predict time interval.

the queries, and predicted simplification. In the first example, EVOQUER predicts interval overlapping with ground truth and correctly translates the verb and noun *close door*. Judging from the video content, the door was already closed; thus, an *open door* action must occur before the *close door*. In the second case, EVOQUER predicts an interval rarely intersecting with ground truth. We review the video and find that at 5.97s, the person in the video starts the action *open the refrigerator door* and pours milk into a glass. Additionally, at 16.4s, he finishes *pouring* and puts the milk back into the refrigerator (shown as the first picture of Figure 3 bottom). Meanwhile, he is holding the glass and leaving the refrigerator door open. Although EVOQUER fails to intersect with the gold standard, it captures the action *open the door* at 5.97s, showing its capability in understanding the video content. We suspect EVOQUER thinks that the person is holding a bag instead of a gallon of milk since both are white in color and similar in size. Thus, it predicts *hold bag* instead of *hold glass*.

## 5 Conclusion

We propose a novel framework, EVOQUER, for temporal grounding that incorporates a query simplification task. It forms closed-loop learning and provides feedback to the temporal grounding model and enhances the learning. Our experiments demonstrate promising results on predicting time intervals and query simplification. Future work will explore more settings and extend to other datasets.

# References

Lisa Anne Hendricks, Oliver Wang, Eli Shechtman, Josef Sivic, Trevor Darrell, and Bryan Russell. 2017. Localizing moments in video with natural language. In *Proceedings of the IEEE international conference on computer vision*, pages 5803–5812.

Jingyuan Chen, Xinpeng Chen, Lin Ma, Zequn Jie, and Tat-Seng Chua. 2018. Temporally grounding natural sentence in video. In *Proceedings of the 2018 conference on empirical methods in natural language processing*, pages 162–171.

Jingyuan Chen, Lin Ma, Xinpeng Chen, Zequn Jie, and Jiebo Luo. 2019a. Localizing natural language in videos. In *Proceedings of the AAAI Conference on Artificial Intelligence*.

Shaoxiang Chen and Yu-Gang Jiang. 2019. Semantic proposal for activity localization in videos via sentence query. In *Proceedings of the AAAI Conference on Artificial Intelligence*.

Shizhe Chen, Qin Jin, and Jianlong Fu. 2019b. From words to sentences: A progressive learning approach for zero-resource machine translation with visual pivots. *arXiv preprint arXiv:1906.00872*.

Pradipto Das, Chenliang Xu, Richard F Doell, and Jason J Corso. 2013. A thousand frames in just a few words: Lingual description of videos through latent topics and sparse object stitching. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2634–2641.

Jiyang Gao, Chen Sun, Zhenheng Yang, and Ram Nevatia. 2017. Tall: Temporal activity localization via language query. In *Proceedings of the IEEE international conference on computer vision*, pages 5267–5275.

Runzhou Ge, Jiyang Gao, Kan Chen, and Ram Nevatia. 2019. Mac: Mining activity concepts for language-based temporal localization. In *2019 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 245–253. IEEE.

Soham Ghosh, Anuva Agarwal, Zarana Parekh, and Alexander Hauptmann. 2019. Excl: Extractive clip localization using natural language descriptions. *arXiv preprint arXiv:1904.02755*.

Hiroaki Gomi and Mitsuo Kawato. 1993. Neural network control for a closed-loop system using feedback-error-learning. *Neural Networks*, 6(7):933–946.

Dongliang He, Xiang Zhao, Jizhou Huang, Fu Li, Xiao Liu, and Shilei Wen. 2019. Read, watch, and move: Reinforcement learning for temporally grounding natural language descriptions in videos. In *Proceedings of the AAAI Conference on Artificial Intelligence*.

Ting-Hao Huang, Francis Ferraro, Nasrin Mostafazadeh, Ishan Misra, Aishwarya Agrawal, Jacob Devlin, Ross Girshick, Xiaodong He, Pushmeet Kohli, Dhruv Batra, et al. 2016. Visual storytelling. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1233–1239.

Mitsuo Kawato. 1990. Feedback-error-learning neural network for supervised motor learning. In *Advanced neural computers*, pages 365–372. Elsevier.

Ranjay Krishna, Kenji Hata, Frederic Ren, Li Fei-Fei, and Juan Carlos Niebles. 2017. Dense-captioning events in videos. In *International Conference on Computer Vision (ICCV)*.

Jason Lee, Kyunghyun Cho, and Douwe Kiela. 2019. Countering language drift via visual grounding. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4376–4386.

Jie Lei, Linjie Li, Luowei Zhou, Zhe Gan, Tamara L Berg, Mohit Bansal, and Jingjing Liu. 2021. Less is more: Clipbert for video-and-language learning via sparse sampling. *arXiv preprint arXiv:2102.06183*.

Zhijie Lin, Zhou Zhao, Zhu Zhang, Qi Wang, and Huasheng Liu. 2020. Weakly-supervised video moment retrieval via semantic completion network. In *Proceedings of the AAAI Conference on Artificial Intelligence*.

Meng Liu, Xiang Wang, Liqiang Nie, Xiangnan He, Baoquan Chen, and Tat-Seng Chua. 2018. Attentive moment retrieval in videos. In *The 41st international ACM SIGIR conference on research & development in information retrieval*, pages 15–24.

Xiang Long, Chuang Gan, and Gerard De Melo. 2018. Video captioning with multi-faceted attention. *Transactions of the Association for Computational Linguistics*, 6:173–184.

Stephanie Lukin, Reginald Hobbs, and Clare Voss. 2018. A pipeline for creative visual storytelling. In *Proceedings of the First Workshop on Storytelling*, pages 20–32, New Orleans, Louisiana. Association for Computational Linguistics.

Huaishao Luo, Lei Ji, Ming Zhong, Yang Chen, Wen Lei, Nan Duan, and Tianrui Li. 2021. Clip4clip: An empirical study of clip for end to end video clip retrieval. *arXiv preprint arXiv:2104.08860*.

Jonghwan Mun, Minsu Cho, and Bohyung Han. 2020. Local-global video-text interactions for temporal grounding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10810–10819.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318.

Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. *arXiv preprint arXiv:2103.00020*.

Cristian Rodriguez, Edison Marrese-Taylor, Fatemeh Sadat Saleh, Hongdong Li, and Stephen Gould. 2020. Proposal-free temporal moment localization of a natural-language query in video using guided attention. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 2464–2473.

Reuben Tan, Huijuan Xu, Kate Saenko, and Bryan A Plummer. 2021. Logan: Latent graph co-attention network for weakly-supervised video moment retrieval. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 2083–2092.

Subhashini Venugopalan, Marcus Rohrbach, Jeffrey Donahue, Raymond Mooney, Trevor Darrell, and Kate Saenko. 2015a. Sequence to sequence-video to text. In *Proceedings of the IEEE international conference on computer vision*, pages 4534–4542.

Subhashini Venugopalan, Huijuan Xu, Jeff Donahue, Marcus Rohrbach, Raymond Mooney, and Kate Saenko. 2015b. Translating videos to natural language using deep recurrent neural networks. In *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1494–1504, Denver, Colorado. Association for Computational Linguistics.

Weining Wang, Yan Huang, and Liang Wang. 2019a. Language-driven temporal activity localization: A semantic matching reinforcement learning model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 334–343.

Xin Wang, Jiawei Wu, Junkun Chen, Lei Li, Yuan-Fang Wang, and William Yang Wang. 2019b. Vatex: A large-scale, high-quality multilingual dataset for video-and-language research. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4581–4591.

Huijuan Xu, Kun He, Bryan A Plummer, Leonid Sigal, Stan Sclaroff, and Kate Saenko. 2019. Multi-level language and vision integration for text-to-clip retrieval. In *Proceedings of the AAAI Conference on Artificial Intelligence*.

Huijuan Xu, Subhashini Venugopalan, Vasili Ramanishka, Marcus Rohrbach, and Kate Saenko. 2015. A multi-scale multiple instance video description network. *arXiv preprint arXiv:1505.05914*.

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*, pages 4507–4515.

Yitian Yuan, Tao Mei, and Wenwu Zhu. 2019. To find where you talk: Temporal sentence localization in video with attention based location regression. In *Proceedings of the AAAI Conference on Artificial Intelligence*.

Da Zhang, Xiyang Dai, Xin Wang, Yuan-Fang Wang, and Larry S Davis. 2019. Man: Moment alignment network for natural language moment retrieval via iterative graph adjustment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1247–1257.

Luowei Zhou, Yannis Kalantidis, Xinlei Chen, Jason J Corso, and Marcus Rohrbach. 2019. Grounded video description. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6578–6587.

Luowei Zhou, Chenliang Xu, and Jason Corso. 2018a. Towards automatic learning of procedures from web instructional videos. In *Proceedings of the AAAI Conference on Artificial Intelligence*.

Luowei Zhou, Yingbo Zhou, Jason J Corso, Richard Socher, and Caiming Xiong. 2018b. End-to-end dense video captioning with masked transformer. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8739–8748.

# A   Appendix

For ActivityNet, we tried different settings while tuning the hyper parameters. The frames are always 32 in these experiments. In the first experiment, We set the learning rate to be 0.00004 and the batch size to be 64. We update the learning rate every 150 epochs, and run 500 epochs in total. In the second experiment, we change the batch size to be 128, and update the learning rate every 200 epochs. We run 600 epochs in total. In the last two experiments, we try a different learning rate, which is 0.0004. We run the experiment for another 600 epochs with all the other parameters to be the same as in experiment 1 and 2.