

This manuscript is currently submitted to Computer Vision and Image Understanding Journal

arXiv:2204.05976v2 [cs.CV] 13 Apr 2022

Video Captioning: a comparative review of where we are and which could be the route

Daniela Moctezuma¹, Tania Ramírez-delReal^{1,2}, Guillermo Ruiz^{1,2}, and Othón González-Chávez¹

¹Centro de Investigación en Ciencias de Información Geoespacial AC, Circuito Tecnopolio II , Aguascalientes, 20313, Mexico

²Consejo Nacional de Ciencia y Tecnología (CONACyT), Av. Insurgentes Sur 1582, Ciudad de Mexico, 03940, Mexico

April 14, 2022

Abstract

Video captioning is the process of describing the content of a sequence of images capturing its semantic relationships and meanings. Dealing with this task with a single image is arduous, not to mention how difficult it is for a video (or images sequence). The amount and relevance of the applications of video captioning are vast, mainly to deal with a significant amount of video recordings in video surveillance, or assisting people visually impaired, to mention a few. To analyze where the efforts of our community to solve the video captioning task are, as well as what route could be better to follow, this manuscript presents an extensive review of more than 105 papers for the period of 2016 to 2021. As a result, the most-used datasets and metrics are identified. Also, the main approaches used and the best ones. We compute a set of rankings based on several performance metrics to obtain, according to its performance, the best method with the best result on the video captioning task. Finally, some insights are concluded about which could be the next steps or opportunity areas to improve dealing with this complex task.

1 Introduction

The incredible explosion of data from the Internet (images, text, videos, etc.) challenges us to deal with it efficiently and appropriately. One of the essential tasks of the research community is generating data understanding technology, for instance, to comprehend a video sequence from surveillance systems. The number of applications that could be beneficial with this technology is enormous, for instance, content retrieval systems, smart video surveillance, and computer-human interface systems, among others.

Nowadays, machine learning methods have well-defined areas, *e.g.*, natural language processing, and artificial vision. One of the most famous and complex tasks is image captioning, which describes the image content using natural language sentences like humans would do. We experience, as humans, that describing visual content is very simple, but for machines, it is a difficult task because there are many aspects to consider, such as the objects in the image, the relations between them, the semantic meaning, and context data, along with others.

There are tremendous advances in image captioning tasks, mainly tackled by deep learning approaches. In this sense, the video captioning task emerged with little attention but showed exciting results with the advent of deep learning algorithms. Besides the complexity of describing a single image's content, describing a video, or images sequences is more challenging because the

time variable becomes crucial to determine the relationship between objects, detecting the actions, and so on.

Nowadays, several approaches are dealing with this task, for instance those exploiting temporal cues [Shi et al.(2019)Shi, Cai, Joty and Gu], [Mun et al.(2019)Mun, Yang, Ren, Xu and Han], [Guo et al.(2019)Guo, Zhang and Gao], motion [Li et al.(2018b)Li, Yao, Pan, Chao and Mei, Wang et al.(2018d)Wang, Wang, Huang, Wang and Tan], [Long et al.(2018)Long, Gan and De Melo], action recognition [Ramanishka et al.(2016)Ramanishka, Das, Park, Venugopalan, Hendricks, Rohrbach and Saenko], [Shetty and Laaksonen(2016)], [Hu et al.(2019)Hu, Chen, Zha and Wu], people trajectories [Qi et al.(2019)Qi, Wang, Li and Luo], [Zhang and Peng(2019b)], events detection [Krishna et al.(2017)Krishna, Hata, Ren, Fei-Fei and Carlos Niebles], [Mun et al.(2019)Mun, Yang, Ren, Xu and Han], [Zhang et al.(2019c)Zhang, Xu, Ouyang and Tan], optical flow [Chen and Jiang(2019)], [Zhou et al.(2018b)Zhou, Zhou, Corso, Socher and Xiong], audio [Xu et al.(2017)Xu, Yao, Zhang and Mei], [Chen et al.(2017)Chen, Chen, Jin and Hauptmann], [Lee(2019)], speech recognition [Iashin and Rahtu(2020)], to name a few. As it will show in Section 4.1 the most fundamental approaches are based on several deep learning architectures, and most of the works are designed as an encoder-decoder model for visual and text data. Another essential aspect presented in many of the analyzed approaches is the well-known attention mechanism that was designed to improve the performance of the encoder-decoder model, specifically on the machine translation task.

A difference between image and video captioning is that in the video, there is a dependency between images to understand the meaning of its content. The whole sequence (or at least a subset of it) must be processed to generate the sentences describing it, on contrary in image captioning, there is no such dependence. Regardless of that, both tasks are currently challenging in machine learning and the natural language processing research community. A possible improvement in video captioning methods' performance could be enhancing the caption generation. In [Shi et al.(2020)Shi, Cai, Gu and Joty], the authors mention that there is a considerable advance in image's coding, but lower performance in the quality of the caption generation process.

In this review, we attempt to show the current state of the video captioning task; the primary datasets used to measure the performance of the proposed approaches, the performance metrics employed, the best-published results, and a deep discussion and analysis of them. To do this, we searched, analyzed, organized, and compared several recent papers and their reported results more comprehensively. This review was done analyzing papers from 2016 to 2021, and most of them were published in conferences or journals. Nevertheless, a few (very few) were considered even though they were published on the arXiv platform¹ to avoid dismissing nothing of the recent work. Through this analysis, some conclusions have also been settled, highlighted some possible improvements, and finally, a complete overall comparison between the analyzed works to find the best solution for the studied period, is presented.

The manuscript is organized as follows, Section 2 explains in formal terms, what is the video captioning task and which are its parts or components. Section 3 incorporates all the elements involved in evaluating the proposed methods or approaches, that means the datasets and the performance metrics. The analysis and discussion are done in Section 4, some aspects to be improved are described in Section 5 as well as some possible applications. Finally, Section 6 gives some conclusions derived from our literature revision and comparison. In Appendix 6 we present tables with all our raw reviewed data.

2 Problem definition

The main goal of video captioning is to enable computers to understand what is happening on a video and build a solid relationship between that content and its corresponding natural language description [Yan et al.(2019)Yan, Tu, Wang, Zhang, Hao, Zhang and Dai]. Video captioning could be seen as the automated task of generating a collection of natural language sentences that describe or explain the video's content [Islam et al.(2021)Islam, Dash, Seum, Raj, Hossain and Shah].

¹<https://arxiv.org/>

As any machine learning problem, the video captioning task can be formulated as follows. Given a pair of (xs, ys) from dataset \mathcal{D} , where xs is a set or sequence of images $xs = xs_1, xs_2, \dots, xs_n$, and ys is a sequence of words, $ys = (w_1, w_2, \dots, w_m)$ that describes the content of xs , find a model that maximizes $p(w_1, w_2, \dots, w_m | xs_1, xs_2, \dots, xs_n)$, that is, find the highest probability of descriptive power ys according with the respective sequence of images xs for all pairs (xs, ys) in the dataset \mathcal{D} .

Figure 1 shows a general overview of the typical solutions for the video captioning task. Most of the works analyzed in this review, have an encoder-decoder framework. The encoder extracts the video features; then, the decoder translates these attributes into the text to generate the descriptions. Different techniques are used to obtain the features. One of them is convolutional neural networks (CNN) in two and three dimensions (C3D); also, an attention mechanism is applied for temporal and spatial characteristics. Furthermore, other features are considered, such as audio, optical flow, maps, object detection, etc. The decoder does the translation applying recurrent neural networks (RNN), specifically, the variant long short-term memory (LSTM). Transformers are another architecture used, as well as the gated recurrent units (GRUs). The video captioning

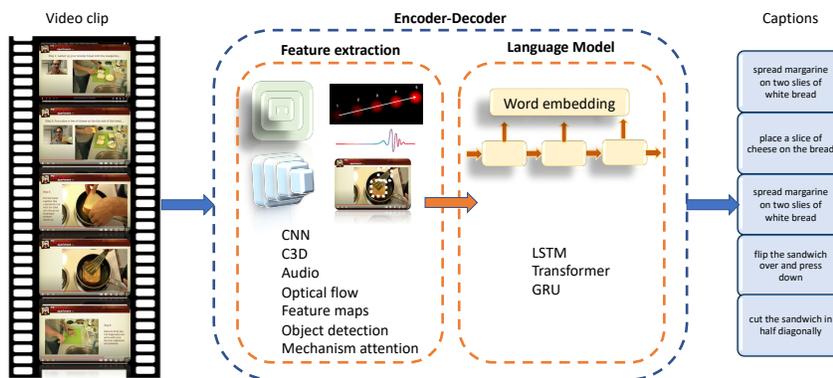


Figure 1: General overview of a typical solutions for video captioning task

models continuously improve to perform better in the training process until the generated captions become similar to the annotated ones. The attention mechanism allows supplementary data to prepare the model to predict the sentences.

A recent review published by [Islam et al.(2021)Islam, Dash, Seum, Raj, Hossain and Shah] organizes the video captioning methods into three classes, traditional methods, machine learning, and deep learning methods. Furthermore, the authors describe parts of the solution of the video captioning task, such as object and action recognition, objects trajectories, etc.; in our review, we will only focus on the overall video captioning task. Another review presented by [Kumar and Mathew(2020)] described many video captioning problems, particularly in extracting multi-modal features, vanishing gradients, and the training time and high resources with deep learning techniques. Also, Kumar et al. explained the possible solutions found in the literature; the papers reported employing more information related to visual, movement, audio, and semantic data contained in the video. Finally, they presented the techniques to extract the features for each data type, where deep learning is widely used, mainly VGG, Resnet, and Inception architectures. Furthermore, other methods help to save temporal information and obtain a feature map for an image sequence; likewise, they concluded that data such as the audio and the known category in videos improve the task of video captioning. A very descriptive work on the several components of proposed approaches is presented by Jain et al. in [Jain et al.(2022)Jain, Al-Turjman, Chaudhary, Nayar, Gupta and Kumar].

In this paper, we tried to detect the best current solution for the video captioning problem, and more than 105 papers were analyzed in a period from 2016 to 2021. Several aspects were observed,

the number of datasets used, the main metrics employed, along with the reported results. We determined which strategy could be considered the best-published solution in this topic in the period considered with several ranking approaches.

3 Performance evaluation

As with any computational task in artificial intelligence, the performance evaluation of the different approaches is crucial to know all the solutions' strengths and weaknesses. Some aspects, such as datasets, metrics, and general conditions, are necessary for an adequate assessment. The following section briefly describes all the datasets, from most to least used in the papers reviewed.

3.1 Datasets

The datasets for video captioning are varied, and the majority of them are publicly available; they mainly belong to cooking or movie clips. This subsection is highlighted through the description and current dataset status regarding its availability and organization of their annotation files.

MSVD (MS-Video YoutubeClips). MSVD [Chen and Dolan(2011)] or YoutubeClips is one of the first datasets, proposed in the year 2011 and generated for the video captioning task. The dataset acquisition was made with Amazon's Mechanical Turk assistance. It contains 1,970 clips and 70,028 sentences, each with 8.7 words on average; the total video duration is 5.3 hours, and it has 13,010 different terms. Although this dataset was able to be downloaded in the past, at this time its link is disabled on the official Microsoft site.

TRECvideo Data (TRECVID). TRECVID [Khan et al.(2012)Khan, Nawab and Gotoh] is conformed by videos about news, meeting, crowd, grouping, traffic, music, sports, animals, and humans interacting with objects. The dataset supplies 140 videos, and there are 20 segments for each. The description for clips has four to six sentences; other information for annotations is keywords and title. The dataset is available on the official site², and the clips and annotations are updated annually for different tasks.

TACoS-MultiLevel [Rohrbach et al.(2014)Rohrbach, Rohrbach, Qiu, Friedrich, Pinkal and Schiele]. It is based on TACoS (Saarbrücken Corpus of Textually Annotated Cooking Scenes) [Regner et al.(2013)Regner, Rohrbach, Wetzl, Thater, Schiele and Pinkal]. The dataset consists of 185 videos and 52,478 sentences in total. The videos' content is about cooking different dishes, and their descriptions are conformed by 3 to 5 sentences. The download is through the official site³, and it is essential to comment that the videos are associated with MPII Cooking 2 Dataset [Rohrbach et al.(2015b)Rohrbach, Rohrbach, Regner, Amin, Andriluka, Pinkal and Schiele].

Youtube2text. Youtube2Text [Guadarrama et al.(2013)Guadarrama, Krishnamoorthy, Malkarinenkar, Venugopalan, Mooney, Darrell and Saenko] is an emerging subset of MSVD; the main difference is the specific use of the English language. Also, Guadarrama et al. proposed a unique split in adjacent videos for training and testing, 1300 and 670, respectively. Contrary to the MSVD dataset, the derived English corpus can be found on its website⁴.

MPII-Movie Description corpus (MPII). MPII [Rohrbach et al.(2015a)Rohrbach, Rohrbach, Tandon and Schiele] is a collection of videos belonging to movies. The primary purpose of this dataset is to provide audio descriptions for movies. It is composed of 94 videos, with 68337 clips and 68375 sentences, and the duration is 73.6 hours. By filling out a request form, videos and annotations are available under demand.

Montreal Video Annotation Dataset (M-VAD). M-VAD [Torabi et al.(2015)Torabi, Pal, Larochelle and Courville] incorporates 84.6 hours of 92 DVDs with 48,986 clips and 55,904 describing sentences. Its data was collected by Descriptive Video Service (DVS) encoded in DVDs. By filling out a submission format, video clips and annotations are available under request.

²trecvid.nist.gov

³www.mpi-inf.mpg.de/departments/computer-vision-and-machine-learning/research/vision-and-language/tacos-multi-level-corpus

⁴www.cs.utexas.edu/users/ml/clamp/videoDescription

MSRVideo to Text (MSR-VTT). MSR-VTT [Xu et al.(2016)Xu, Mei, Yao and Rui] is one of the most extensive datasets for the video captioning task. MSR-VTT supplies 10K web video clips and 200K sentences in total; therefore, 20 different sentences describe each clip; and the length of each one is 10 to 30 seconds; the total duration is about 41.2 hours. The video content is mainly related to gaming, sports, and movies. The official page of this dataset does not have a link that allows its download.

Large Scale Movie Description Challenge (LSMDC). LSMDC [Rohrbach et al.(2017)Rohrbach, Torabi, Rohrbach, Tandon, Pal, Larochelle, Courville and Schiele] is a subset of videos from M-VAD and MPII; also, it includes more movies. The dataset is split to balance the movie genres; the purpose is diversity in the vocabulary. LSMDC 2016 contains 101,046 and 7,408 training and validation clips, respectively. The 2021 version is currently available ⁵, which must be requested by sending a form to access it.

Charades [Sigurdsson et al.(2016)Sigurdsson, Varol, Wang, Farhadi, Laptev and Gupta]. To build this dataset, the Hollywood in Homes process, and Amazon Mechanical Turk tools were used. It can also be used for action classification and localization tasks. The dataset contains 9,848 videos with 30.1 seconds on average length and 27,847 textual descriptions. It is available through its official site⁶.

ActivityNet Captions [Krishna et al.(2017)Krishna, Hata, Ren, Fei-Fei and Carlos Niebles]. It incorporates 20k videos, and the duration is 849 hours with 100k full descriptions, is based on the ActivityNet [Caba Heilbron et al.(2015)Caba Heilbron, Escorcia, Ghanem and Carlos Niebles] dataset for human activity understanding. Each clip has a duration of 10 minutes, and it contains 3.65 phrases on average. The descriptions belong to a specific time. This dataset is available for download on its official site ⁷.

Visual Attentive Script (VAS). VAS [Yu et al.(2017a)Yu, Choi, Kim, Yoo, Lee and Kim] contains 144 video clips of 15 seconds, and each video has three sentences. This dataset is not available to download.

YouCookII [Zhou et al.(2018a)Zhou, Xu and Corso] is a dataset with two annotations per video segment, the first description explains a recipe, and the second is a verification. The dataset contains 2,000 clips describing 89 recipes, 13,829 sentences, and the duration is 176 hours in total; each video has 3-16 segments. The dataset and all its information are available on its official site ⁸.

Fine-grained Sports Narrative (FSN). FSN [Yu et al.(2018)Yu, Cheng, Ni, Wang, Zhang and Yang] contains 2,000 videos from Youtube with an average of 3.16 sentences for each, particularly basketball playing in the NBA. The annotations consist of timestamps and a descriptive paragraph. The dataset is not available for download.

VATEX [Wang et al.(2019b)Wang, Wu, Chen, Li, Wang and Wang] contains 41,250 videos and 825,000 sentences. The videos consist of human activity, and it is in two languages, English and Chinese. The data collection is from Kinetics-600, and it is combined with YouTube videos associated with human activities, the duration of each clip is about ten seconds. Amazon Mechanical Turkers made the captions according to the videos, 10 for each language. The dataset is available through its official site ⁹.

Sports Video Captioning Dataset-Volleyball (SVCDV). SVCDV [Qi et al.(2019)Qi, Wang, Li and Luo] is another dataset related to sports, specifically volleyball. This dataset contains 55 videos from Youtube composed of 4,830 clips, with an average of 9.2 sentences for each, which means 44,436 phrases in total. The dataset is not available for download.

Object-oriented captions. Object-oriented captions [Zhu et al.(2020)Zhu, Hwang, Ma, Chen and Guo] dataset is based on ActivityNet Captions, explicitly a subset related to the action where a game is played. The authors proposed making a new annotation to have direct relations expressed in the sentences between the objects. The dataset contains 75 videos and 534 sentences; furthermore, the descriptions are more extensive because it includes more words, verbs, and ad-

⁵<https://sites.google.com/site/describingmovies/>

⁶allenai.org/data/charades

⁷cs.stanford.edu/people/ranjaykrishna/densevid/

⁸youcook2.eecs.umich.edu

⁹[eric-xw.github.io/vatex-website](https://github.com/eric-xw/vatex-website)

jectives than MSR-VTT, MSVD, ActivityNet Captions, and FSN datasets. The dataset is not available for download.

LIRIS human activities dataset. LIRIS [Wolf et al.(2014)Wolf, Lombardi, Mille, Celiktutan, Jiu, Dogan, Eren, Baccouche, Dellandréa, Bichot et al.] is a dataset specialized in human activity, particularly surveillance and office conditions. The original dataset contains 828 actions without captions, and it is organized into two subsets named D1 and D2. In recent work [Inácio et al.(2021)Inácio, Gutoski, Lazzaretti and Lopes], subset D2 is refreshed with descriptions of the scene action, and it supplies 167 videos with 367 annotations related to spatial, temporal, and description content. This dataset has an official site¹⁰, and it is possible to download the clips; the caption sentences are not available in any repository.

UET Video Surveillance (UETVS). UETVS [Dilawari et al.(2021)Dilawari, Khan, Al-Otaibi, Rehman, Rahman and Nam] is a collected database for smart surveillance purposes, it contains 1200 videos with 3-6 sentences for each video, and the descriptions were acquired from professional English writers. The dataset is not available for download.

AGRIINTRUSION [Dilawari et al.(2021)Dilawari, Khan, Al-Otaibi, Rehman, Rahman and Nam]. This dataset contains 100 videos from YouTube, and the context is the agricultural environment. The dataset is not available for download.

All the aforementioned datasets were used to assess the proposed methods observed in our literature revision. Nevertheless, there are more and less used datasets, the most used detected from 2016 to 2021 years, are MSVD and MSR-VTT, maybe due to the diversity and variety in videos and captions. It is essential to highlight that the published articles use the terms MSVD, Youtube2Text, and YoutubeClips to refer to the same dataset; however, Youtube2Text is a subset that only contains the English language; in addition, it is important to point out the fact that the authors proposed a specific split for the evaluation of the models. Another widely used dataset is ActivityNet Captions, mainly due to annotations in a specific time. Charades has also begun to take off to evaluate proposals and will probably have more use since it is of recent creation.

Figure 2 shows a representative sample for two datasets; Figure 2 (a) is for the MSR-VTT dataset; the clip contains 20 different sentences describing a segment of a video; Figure 2 (b) belongs to a sample of ActivityNet Captions, which is a dataset used in dense-captioning, the clips can have distinct phrases to describe actions in overlapped segments.

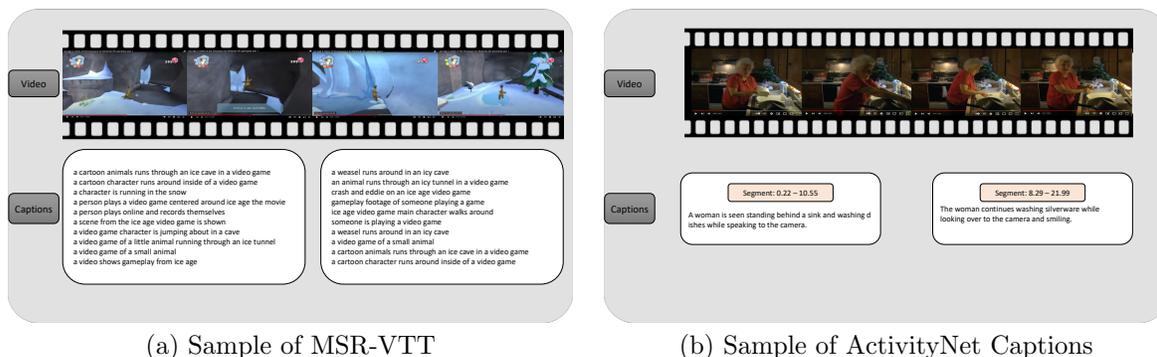


Figure 2: Two representative samples from benchmarks for video captioning datasets with ground-truth captions

In summarizing, Table 1 shows essential information about the datasets, especially the year of the creation, the videos and clips they contain, and finally, the total number of sentences. It is necessary to mention that the UETVS and AGRINTRUSION datasets are not in the table because the total number of sentences is not available, and they are not published in a repository or website. The bolded text datasets correspond to those most used, considering a threshold of

¹⁰projet.liris.cnrs.fr/voir/activities-dataset

more than seven citations. The complete analysis presented in this work is focused on these seven most-used datasets.

Table 1: Datasets for video captioning. * Available online,** Available through request.

Dataset	Year	Videos	Clips	Sentences
MSVD [Cohen and Dolan(2011)]	2011	1,970	1,970	70,028
TRECVID [*] [Khan et al.(2012)Khan, Navech and Gotoh]	2012	140	140	1,820
Youtube2Text [*] [Gandarrama et al.(2013)Gandarrama, Krishnamoorthi, Malakarandhar, Venugopalan, Mooney, Darrell and Saezko]	2013	1,970	1,970	70,028
TACoS-MultiLevel [*] [Rohrbach et al.(2014)Rohrbach, Rohrbach, Qiu, Friedrich, Pankel and Schiele]	2014	185	16,145	52,478
MPL-MD [*] [Rohrbach et al.(2015a)Rohrbach, Rohrbach, Tandon and Schiele]	2015	94	68,337	68,375
M-VAD ^{**} [Rohrbach et al.(2017)Rohrbach, Torabi, Rohrbach, Tandon, Pal, Larochelle, Courville and Schiele]	2015	92	48,986	55,904
Charades [*] [Sigurdsson et al.(2016)Sigurdsson, Vardi, Wang, Farhadi, Laptev and Gupta]	2016	9,848	9,848	27,847
MRE-VTT [Xu et al.(2016)Xu, Jie, Yao and Rui]	2016	7,180	10,000	200,000
LSMDC ^{**} [Rohrbach et al.(2017)Rohrbach, Torabi, Rohrbach, Tandon, Pal, Larochelle, Courville and Schiele]	2016	128K	128K	128K
ActivityNet Captions [*] [Krishna et al.(2017)Krishna, Hata, Ren, Fei-Fei and Carlos Niebles]	2017	20,000	100,000	100,000
VAS [Yu et al.(2017a)Yu, Chae, Kim, Yoo, Lee and Kim]	2017	144	144	402
YouTube2Caption [*] [Zhou et al.(2015a)Zhou, Xu and Corso]	2018	2,000	2,000	13,829
FSN [Yu et al.(2018)Yu, Cheng, Ni, Wang, Zhang and Yang]	2018	2,000	2,000	6,520
VATEX [*] [Wang et al.(2019b)Wang, Wu, Chen, Li, Wang and Wang]	2019	41,250	41,250	325,000
SVCDV [Qi et al.(2019)Qi, Wang, Li and Luo]	2019	55	4,830	44,436
Object-oriented captions [Zhu et al.(2020)Zhu, Hwang, Ma, Chen and Guo]	2020	75	75	531
LIRIS_V2 with descriptions [Jaisin et al.(2021)Jaisin, Gotohko, Lanzaretti and Lopes]	2021	167	167	367

3.1.1 Performance Metrics

As in the image captioning task, the metrics commonly used to evaluate the performance of video captioning are BLEU, METEOR, ROUGE-L, and CIDEr. These metrics are used to evaluate the performance between the system’s results against a human interpretation. The main metrics are described in the following sections.

BLEU - Bilingual Evaluation Understudy. BLEU [Papineni et al.(2002)Papineni, Roukos, Ward and Zhu] is one of the most used metrics to evaluate video captioning task. BLEU analyzes co-occurrences of n-grams between the candidate and reference sentences. These matches are positionally independent, and hence the more matches, the better the result. The n-gram is a sequence of n elements; in this case, the elements are the words (or tokens if previous tokenization is performed). Hence, we can calculate BLEU-1 splitting in unigrams, BLEU-2 for bigrams, and so on.

The clipped n-gram counts for all the candidate sentences and divides by the number of candidate n-grams in the test corpus to calculate the precision score p_n . To calculate several results of several configurations of n-grams, BLEU uses the geometric mean, and to avoid shorter candidates, BLEU applies a penalty factor. The final calculation is performed as:

$$BLEU = BP * exp \left(\sum_{n=1}^N w_n \log p_n \right), \quad (1)$$

where p_n is the precision score, BP is the brevity penalty, N is the longest sequence of n-grams to consider, and w_n are the weights for each different n-gram length.¹¹ BLEU values range from 0 to 1, being 0 worst case and 1 best possible value. For more details about BLEU calculation, please see [Papineni et al.(2002)Papineni, Roukos, Ward and Zhu] reference.

METEOR - Metric for Evaluation of Translation with Explicit Ordering. METEOR [Banerjee and Lavie(2005)] arose to manage the drawbacks in BLEU. The score results of matching words to make a relation translation; this can measure the satisfactorily ordered in the translation.

The score of METEOR considers an F_{mean} and a $Penalty$. F_{mean} is based on precision (P) and recall (R) for unigram matches, and it is considered a harmonic mean. The Equation 2 shows this mathematic relation.

$$F_{mean} = \frac{10PR}{R + 9P} \quad (2)$$

Specifically, P is the ratio of the unigrams of reference and the unigrams obtained for the model, and R is the ratio of the mapped unigrams of the reference to the total.

The final consideration is the $Penalty$; it consists in mapping the fewest feasible quantity of chunks in the sentence to unigrams; if the chunks increase, the penalty does the same. The equation 3 expresses this relation.

¹¹The original paper [Papineni et al.(2002)Papineni, Roukos, Ward and Zhu] uses $N = 4$ and uniform weights, $w_n = 1/N$.

$$Penalty = 0.5 * \left(\frac{no. \text{ chunks}}{no. \text{ unigrams_matched}} \right)^3 \quad (3)$$

Finally, the METEOR score is computed as the equation 4 shows.

$$METEOR = F_{mean} * (1 - Penalty) \quad (4)$$

ROUGE - Recall-Oriented Understudy for Gisting Evaluation. ROUGE [Lin(2004)] is a metric evaluation that measures the number of overlapping in n-gram, term sequences, and word pairs for the captions generated. This metric has four versions: ROUGE-N, ROUGE-L, ROUGE-W, and ROUGE-S. ROUGE-L is used in the video captioning task, and it estimates the longest common subsequence (LCS). The first step is to obtain the Recall (R_{lcs}) and Precision (P_{lcs}) associated with LCS.

R_{lcs} is the ratio of the maximum length of reference (s_r) and predicted (s_p) sentence $LCS(s_r, s_p)$, and the total words of the reference sentence (see equation 5). P_{lcs} is very similar to the recall, but the length consideration is related to the predicted sentence (see equation 6).

$$R_{lcs} = \frac{LCS(s_r, s_p)}{len(s_r)}, \quad (5)$$

$$P_{lcs} = \frac{LCS(s_r, s_p)}{len(s_p)}. \quad (6)$$

Finally, the ROUGE-L score is calculated by the equation 7:

$$ROUGE-L = \frac{(1 + \beta^2)R_{lcs}P_{lcs}}{R_{lcs} + \beta^2P_{lcs}}, \quad (7)$$

where β is a constant number to consider the precision.

CIDeR - Consensus-based Image Description Evaluation. CIDeR measures the similarity of the generated captions versus a set of ground truth sentences created by humans [Vedantam et al.(2015)Vedantam, Lawrence Zitnick and Parikh]. CIDeR arose as an answer to the necessity for a metric that could better correlate with human judgment, employing a *consensus* between the several proposed captions, or references, in a dataset. In this case, one can have two generated sentences called B and C, and a reference sentence called A. CIDeR tries to solve the question of which of the two sentences, B or C, is more similar to A?

To evaluate the generated caption (or candidate sentence c_i) for an image I_i versus a set of image descriptions or references $S_i = s_{i1}, \dots, s_{im}$, with CIDeR, firstly, all words are converted into its root form. Each sentence is represented using the n-grams approach, originally in [Vedantam et al.(2015)Vedantam, Lawrence Zitnick and Parikh]¹². The n-gram w_k is a set of one or more sequential words. The TF-IDF approach is used to generate CIDeR metric. See [Robertson(2004)] for TF-IDF details.

The CIDeR_n for n-grams of size n is calculated with the average cosine similarity between the candidate caption and the references as follows:

$$CIDeR_n(c_i, S_i) = \frac{1}{m} \sum_j \frac{g^n(c_i) \cdot g^n(s_{ij})}{\|g^n(c_i)\| \|g^n(s_{ij})\|}, \quad (8)$$

where $g^n(c_i)$ is a vector formed by $g_k(c_i)$ corresponding to all the n-grams of size n , and $\|g^n(c_i)\|$ is the magnitude of the vector $g^n(c_i)$, is the same for $g^n(s_{ij})$. Finally, CIDeR combines all the results from n values for n-grams as follows:

$$CIDeR(c_i, S_i) = \sum_{n=1}^N w_n CIDeR_n(c_i, S_i), \quad (9)$$

¹²Only values of 1 to 4 for n are considered.

3.1.2 Performance Metrics Examples

To exemplify, Table 2 shows a list of the different values that can be calculated for a single predicted caption, given five different human-annotated captions references. This sample sentence was extracted from the MS-COCO dataset [Lin et al.(2014)Lin, Maire, Belongie, Hays, Perona, Ramanan, Dollár and Zitnick] for image annotation, but the same procedure can be extrapolated to video captioning.

Table 2: Different values of metrics for a single proposed caption considering five references.

Predicted caption	References
A box of different types of doughnuts on a table.	Four donuts in a box with a variety of frostings.
	A close up of many different kinds of doughnuts.
	A group of donuts sitting in a box.
	An image of a box of donuts in a box.
	A group of different types of doughnuts in a box

Metric	Value
BLEU-1 [Papineni et al.(2002)Papineni, Roukos, Ward and Zhu]	0.800000
BLEU-2 [Papineni et al.(2002)Papineni, Roukos, Ward and Zhu]	0.730297
BLEU-3 [Papineni et al.(2002)Papineni, Roukos, Ward and Zhu]	0.643660
BLEU-4 [Papineni et al.(2002)Papineni, Roukos, Ward and Zhu]	0.525382
METEOR [Banerjee and Lavie(2005)]	0.405500
ROUGE-L [Lin(2004)]	0.600000
CIDEr [Vedantam et al.(2015)Vedantam, Lawrence Zitnick and Parikh]	2.203552

As can be seen, the same predicted caption could have very different values regarding the metric employed.

4 Results and discussion

The presented results are based on the review and analysis of 105 published papers in journals and conferences related to video captioning; additional reviews papers were also consulted. The range of years of publication of the reviewed papers is from 2016 to 2021. Our review was done by selecting specific points from all papers, these points are title, year, method’ description, datasets used, results reached, and where the paper was published. We did the analysis detailed in the following sections from all this information.

From the 105 papers, we counted 20 different datasets reported; some datasets were subsets from others. Please see section 3.1 for more details about the descriptions of each dataset reviewed.

Figure 3 shows a histogram of all the datasets found in our literature review. Here, it can be seen that the most cited dataset is MSVD, with 66 times used in different papers. The second is MSR-VTT, used in 61 different works. It can be observed that these two datasets are, by far, the most utilized. Nevertheless, there is a block of five datasets with several citations higher than or equal to 8, and there are also many datasets with few cites; exactly 13 datasets were referenced in less than four works.

As mentioned, the analysis was done from 2016 to 2021, and to observe how many manuscripts were published over the years, Figure 4 shows the number of publications per year being 2019, the year with more papers published related to the video captioning task. In 2020 and 2021, only 12 and 7 manuscripts, were published, respectively. This could be due to the COVID-19 pandemic since more than 23 papers were published the previous year.

To follow a simple analysis scheme, first, the frequency of usage of all the performance metrics was calculated. As a result, Figure 5 displays the percentage of usage of each reported metric. Here, it can be seen that the most used are METEOR, followed by BLEU-4 and CIDEr-D, with second and third place, respectively. BLEU-1, BLEU-2, and BLEU-3 have a similar frequency, and in the label “Other” we joined those metrics with less than two reported results; these are the Average-Recall (AR), SPICE, FCE, RE, and Self-BLEU.

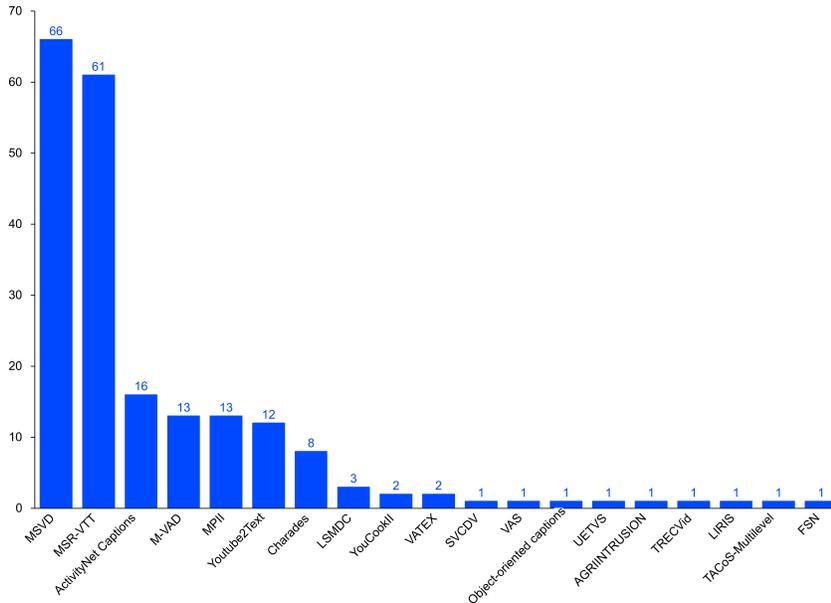


Figure 3: Number of citations for each dataset

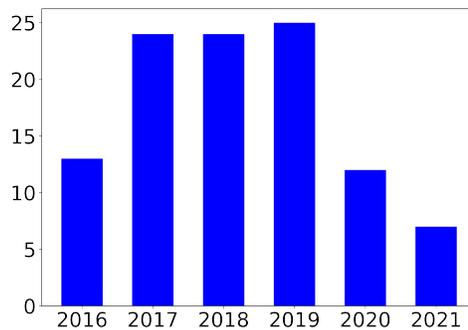


Figure 4: Number of papers related to video captioning task published by year, since 2016 to 2021

To analyze in an organized way, we sorted the results in each dataset based on the METEOR metric, which is the most used (see Figure 5), at least in the reviewed literature. After the results were sorted, we took the top five in each dataset. Table 3 shows the best five results reported in the seven most-used datasets. Here, is observed that the best ranked methods are not necessarily the most recent ones. The best result on the MSVD dataset was reached by [Perez-Martin et al.(2021)Perez-Martin, Bustos and Pérez] on the three compared metrics. The same for the MSR-VTT dataset, where Gao et al. [Gao et al.(2019)Gao, Li, Song and Shen] achieved the best result in all metrics. Nevertheless, this is not the case for the ActivityNet Captions dataset, where [Xiong et al.(2018)Xiong, Dai and Lin] obtained the best result in the METEOR and BLEU-4 metrics and the worst result on CIDEr-D. Only the METOR metric was used to report results for the M-VAD dataset because there were missing values in both BLEU-4 and CIDEr-D. The best result on the MPII dataset was achieved by [Pan et al.(2017)Pan, Yao, Li and Mei], for the Youtube2Text dataset the best result on the METEOR and CIDEr-D was reached by [Chen et al.(2018a)Chen, Li, Zhang and Huang], but this was not the case for BLEU-4. Finally, on the Charades dataset, [Zhao et al.(2019)Zhao, Li and Lu] obtained the best score if we only consider the METEOR metric, but for BLEU-4 and CIDEr-D, that is not the case.

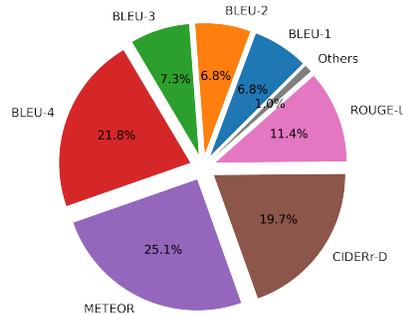


Figure 5: Metrics usage over 105 papers revised

From these reported results, we can observe that there is no simple way to choose the best method because we have multiple metrics. Another argument for this, is that most of the reviewed manuscripts usually employ just one or two datasets in their evaluation, making it difficult to declare an overall winner.

Table 3: The best five results according to METEOR on the seven most used datasets

MSVD	Year	Work	METEOR	BLEU-4	CIDEr-D
	2021	[Perez-Martin et al.(2021)Perez-Martin, Bustos and Pérez]	41.90	64.40	111.50
	2018	[Pu et al.(2018)Pu, Min, Gan and Carin]	38.03	54.27	78.31
	2020	[Pan et al.(2020)Pan, Cai, Huang, Lee, Gaidon, Adeli and Niebles]	36.90	52.20	93.00
	2021	[Chen and Jiang(2021)]	36.90	55.80	74.50
	2020	[Zhang et al.(2020)Zhang, Shi, Yuan, Li, Wang, Hu and Zhu]	36.40	54.30	95.20
MSR-VTT	Year	Work	METEOR	BLEU-4	CIDEr-D
	2019	[Gao et al.(2019)Gao, Li, Song and Shen]	33.50	54.30	72.80
	2021	[Perez-Martin et al.(2021)Perez-Martin, Bustos and Pérez]	30.40	46.40	51.90
	2018	[Pu et al.(2018)Pu, Min, Gan and Carin]	29.98	45.01	51.41
	2019	[Hou et al.(2019)Hou, Wu, Zhao, Luo and Jin]	29.70	42.30	49.10
	2019	[Chen et al.(2019)Chen, Jin, Chen and Hauptmann]	29.61	44.91	51.80
ActivityNet Captions	Year	Work	METEOR	BLEU-4	CIDEr-D
	2018	[Xiong et al.(2018)Xiong, Dai and Lin]	14.75	8.45	14.15
	2019	[Mun et al.(2019)Mun, Yang, Ren, Xu and Han]	13.07	1.28	43.48
	2020	[Iashin and Rahr(2020)]	11.72	2.86	NA
	2019	[Hou et al.(2019)Hou, Wu, Zhao, Luo and Jin]	11.30	1.90	44.20
	2019	[Zhang et al.(2019)Zhang, Xu, Ouyang and Tan]	10.71	1.64	31.41
M-VAD	Year	Work	METEOR	BLEU-4	CIDEr-D
	2017	[Pasumuru and Baussal(2017)]	7.40	NA	NA
	2017	[Pan et al.(2017)Pan, Yao, Li and Mei]	7.40	NA	NA
	2017	[Baraldi et al.(2017)Baraldi, Grassia and Cocchiara]	7.30	NA	NA
	2018	[Xu et al.(2018)Xu, Han, Hong and Tian]	7.20	NA	NA
	2018	[Pu et al.(2018)Pu, Min, Gan and Carin]	7.12	2.08	9.14
MPH	Year	Work	METEOR	BLEU-4	CIDEr-D
	2017	[Pan et al.(2017)Pan, Yao, Li and Mei]	28.8	40.8	47.1
	2018	[Xu et al.(2018)Xu, Liu, Wong, Zhang, Nie, Su and Kankanhalli]	7.9	1.9	NA
	2019	[Zhao et al.(2019)Zhao, Li and Lu]	7.8	NA	NA
	2019	[Xu et al.(2019)Xu, Liu, Nie and Su]	7.7	0.8	NA
	2016	[Pan et al.(2016)Pan, Mei, Yao, Li and Rui]	7.3	NA	NA
Y2T	Year	Work	METEOR	BLEU-4	CIDEr-D
	2018	[Chen et al.(2018)Chen, Li, Zhang and Huang]	35.23	53.21	86.76
	2018	[Zoflaghari et al.(2018)Zoflaghari, Singh and Brox]	35.00	53.50	85.80
	2020	[Wei et al.(2020)Wei, Mi, Hu and Chen]	34.40	46.80	85.70
	2017	[Hori et al.(2017)Hori, Hori, Lee, Zhang, Hashim, Hershey, Marks and Sumi]	34.30	56.80	72.40
	2017	[Chen et al.(2017)Chen, Chen, Jin and Hauptmann]	34.21	47.56	79.57
Charades	Year	Work	METEOR	BLEU-4	CIDEr-D
	2019	[Zhao et al.(2019)Zhao, Li and Lu]	19.7	12.9	18.8
	2017	[Li et al.(2017)Li, Zhao, Lu et al.]	19.1	12.7	18.3
	2018	[Zhao et al.(2018)Zhao, Li, Lu et al.]	19.0	13.3	18.0
	2018	[Wang et al.(2018)Wang, Chen, Wu, Wang and Wang]	18.7	18.8	23.6
	2019	[Hu et al.(2019)Hu, Chen, Zha and Wu]	18.4	14.5	23.7

For these reasons, we present a set of box plots to observe the behavior of each of these three most used metrics over all the datasets. Figure 6 shows the behavior of BLEU-4, METEOR, and CIDEr-D per dataset, over MSVD, MSR-VTT, ActivityNet Captions, M-VAD, MPH, Youtube2Text, and Charades datasets, respectively. All the reported results over this revision of 105 works are grouped by dataset and metric. For instance, on the MSVD dataset, we can observe the median, min, and max values on each metric. Using the BLEU-4 metric, the median value is around 50, but most of the works that reported results on BLEU-4 are below this median; the same happens for METEOR, and the CIDEr-D metric. These aspects could be observed in the rest of the box plots from each dataset, and several insights could be noted. For instance, the METEOR metric has less variation on Youtube2Text despite all three metrics having the same amount of reported results. The same happens with the Charades, and ActivityNet Captions datasets. Another outlier behavior is from the CIDEr-D metric on the M-VAD dataset, and this is because (as we can see in Table 3) only two results were reported using the CIDEr-D metric.

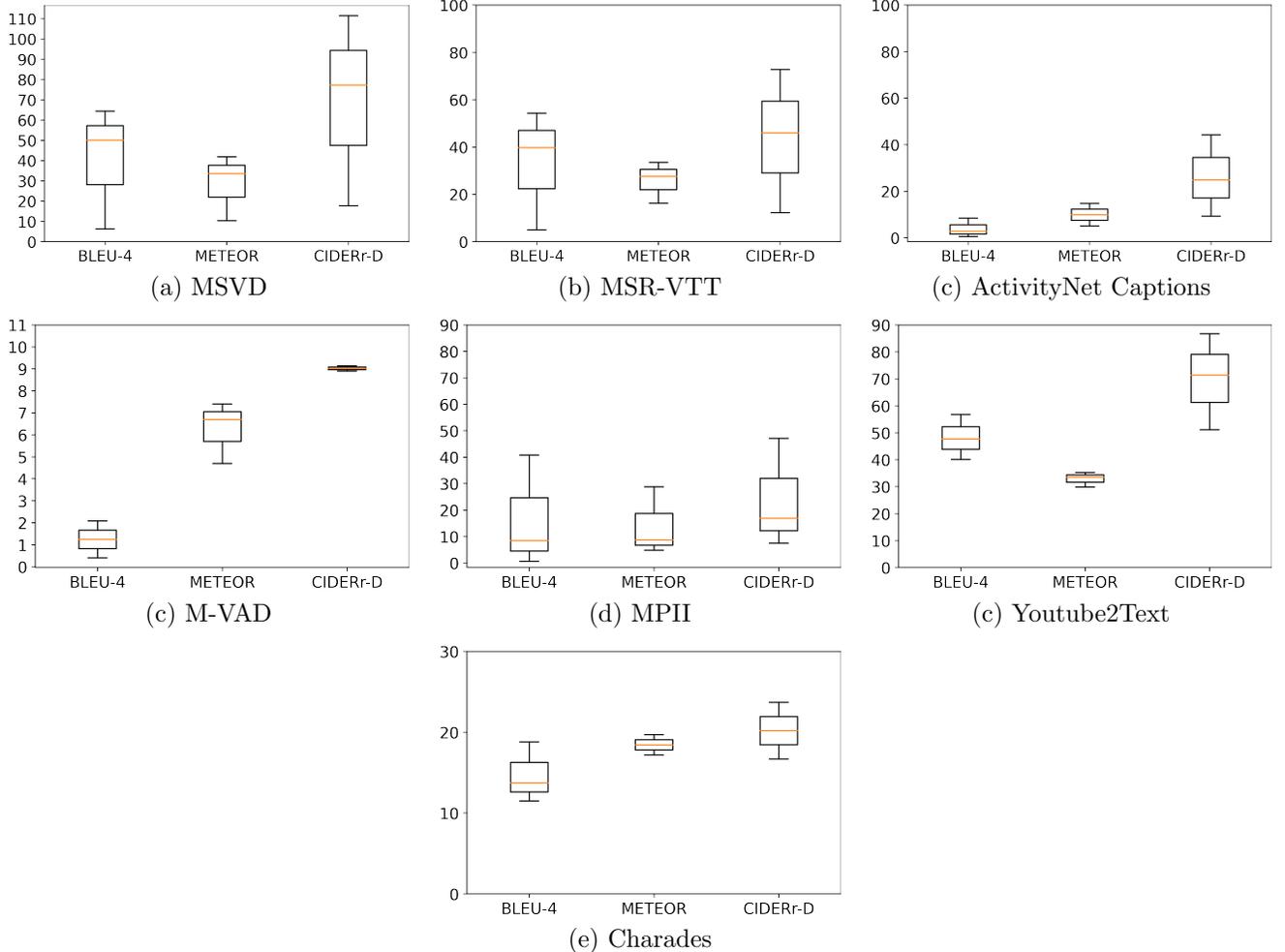


Figure 6: Results in most used datasets with three most used metrics

Trying to present the evolution of all the reported results over the years, in Figures 7, 8, and 9, we show the best result over the years from each dataset using METEOR, BLEU-4, and CIDEr-D metrics. In Figure 7, it can be seen that the highest results always were achieved on the MSVD dataset (red bar in the figure). It is also one of the datasets used in all the studied years, along with MSR-VTT. The lowest results on the METEOR metric, were achieved by MVAD and MPII datasets related papers. For the BLEU-4 metric, the MSVD dataset got the highest results again over the period studied. Contrary to METEOR, the lowest scores were obtained with ActivityNet Captions and M-VAD datasets related works. With the CIDEr-D metric, once again, the MSVD dataset achieved the highest reported results and the lowest on MPII and MVAD datasets. As noted, models trained using MSVD reached the best result in 2021, but this is not the case with the other datasets that achieved the best reported results in different years, from 2017, 2018, and 2019.

In order to obtain a global ranking, we computed individual rankings on the three most used metrics. We gathered the results on METEOR, BLEU-4, and CIDEr-D from all the reviewed works in the seven most-used datasets. This ranking procedure considers the top 5 results according to the metric. Analyzing the set of papers, we can compute the result of each of them in the seven most-used datasets. That means, if a paper achieves the first position, 5 points will be assigned to it; the second position grants 4 points, and go on until the fifth position, with only one point; in this way, we get the paper with the highest number of points. Finally, the sum of

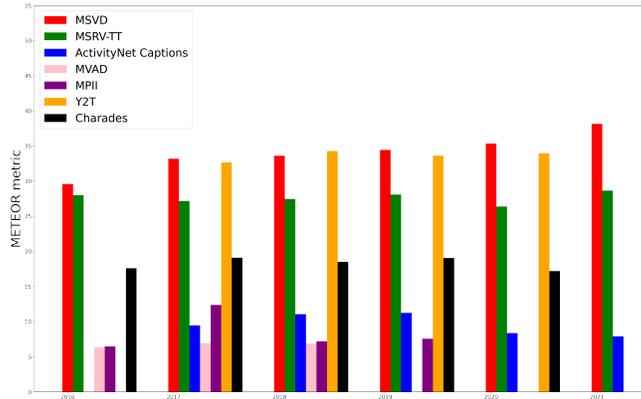


Figure 7: All datasets through years 2016-2021 using METEOR metric

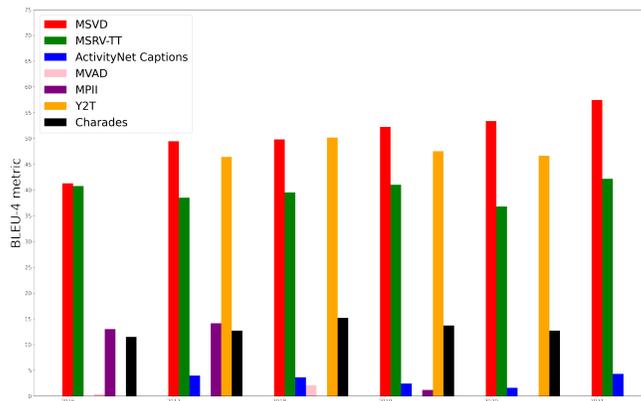


Figure 8: All datasets through years 2016-2021 using BLEU-4 metric

the points is divided by 7 (total number of datasets used) to normalize them.

These ranking results can be seen in Figure 10(a), (b), and (c), respectively. The best work over all the datasets with the METEOR metric was the proposed by [Zhao et al.(2019)Zhao, Li and Lu] (marked in red); the higher the bar value the better the ranking is, so those works with lower bar height have the worst performance. Similar behavior from BLEU-4, in which two works reported best results, the proposed in [Xiao and Shi(2019)], and in [Chen et al.(2019)Chen, Jin, Chen and Hauptmann]. And, using CIDERr-D, the best method ranked is the proposed by [Chen et al.(2019)Chen, Jin, Chen and Hauptmann]

Since not all the works utilized all the datasets to test their solutions, it is not straightforward to rank results over all the datasets; hence, a new ranking was calculated, joining the prior three ones through the average of its previous values, and then sorted them to highest to lowest ones. That is, in the case of one paper having first place in one dataset but also second place in another dataset, it will have a better ranking than others only citing one dataset. The new ranking considers the result in each dataset and how many datasets the work used. The final ranking is shown in Figure 10 (d) from these computations. With the above, a unique winner could be selected; in this case, the best result considering all the metrics and datasets is achieved by the method proposed by Xiong et al. (see [Xiong et al.(2018)Xiong, Dai and Lin]). It is important to note that it is easy to observe the best results with the ranking.

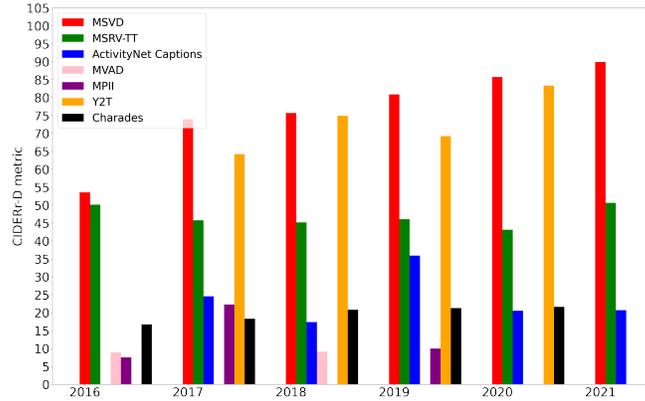
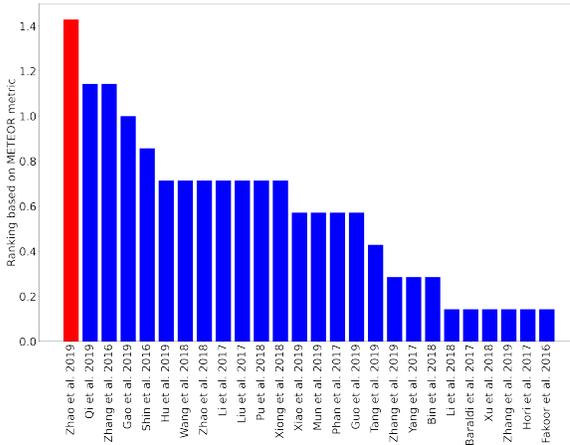
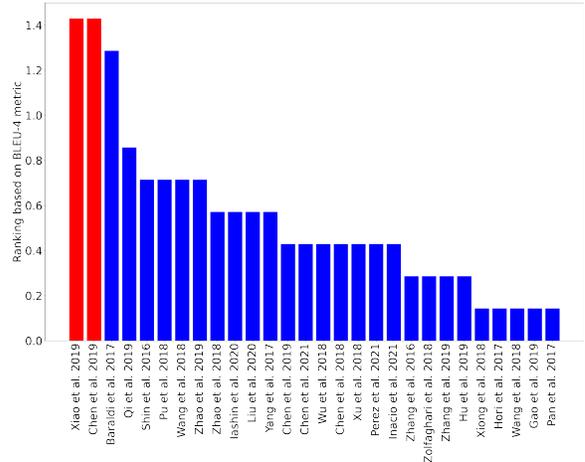


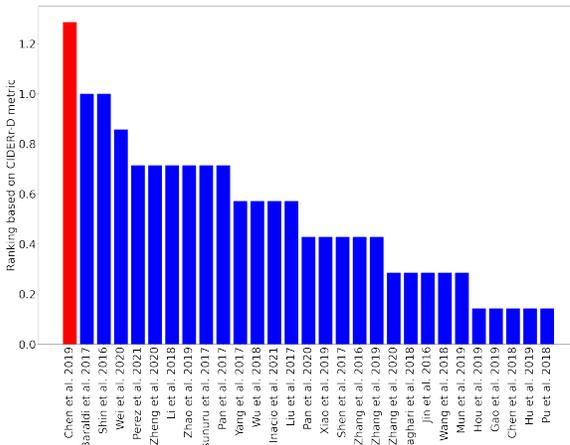
Figure 9: All datasets through years 2016-2021 using CIDEr-D metric



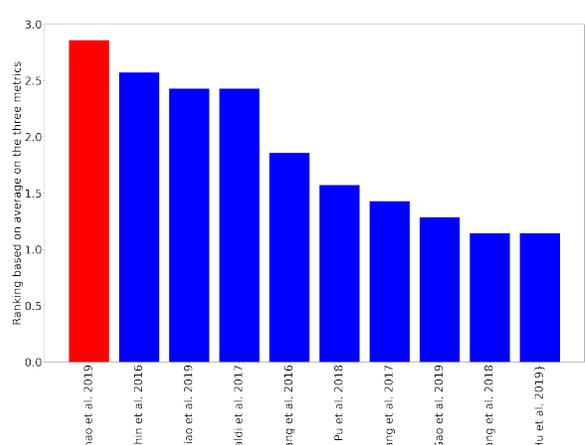
(a) Ranking based on METEOR



(b) Ranking based on BLEU-4



(c) Ranking based on CIDEr-D



(d) Ranking based on average of the 3 metrics

Figure 10: Results based on METEOR metrics in most used datasets

4.1 About approaches

Once the best method in the literature review is determined, it is also essential to know what approaches are prevalent.

The revision extracted the reported results for each dataset, the methods employed, and the follow-up of the references cited or compared. In this way, we state that most of the reported results related to video captioning are considered in this study. The first analysis identifies the technique used in the encoder-decoder framework due to most of the papers utilizing similar methods adding or proposing an innovation to achieve better performance. The convolutional neural networks (CNN) are present in most of the works; the second most employed network is the C3D, and in the last years, the R-CNN architecture has been increasingly used due to the capability to detect objects. Furthermore, the process of audio has been presented through the MFCC technique [Memon et al.(2009)Memon, Lech and He, Sahidullah and Saha(2009)].

Other methods have been used, and it is crucial to highlight that the total numbers presented here do not correspond to the 105 papers; the main reason is the fusion of features. Figure 11 shows a graphic distribution of approaches utilized to extract features in the decoder for video captioning. It is indispensable to mention that a few papers fusion different techniques for this purpose.

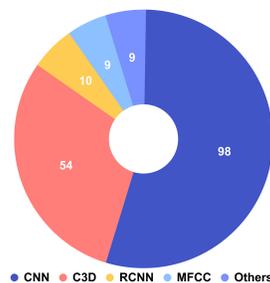


Figure 11: Encoders used over 105 papers revised

The approaches for a decoder in video captioning are particularly emphasized in employing LSTM architectures, and this technique is a specific arrangement of a recurrent neural network (RNN).

However, GRUs and transformers have arisen in recent years; other methods have emerged but do not have much presence yet. Figure 12 displays a visual distribution of practices employed to covert the features of the decoder in captions.

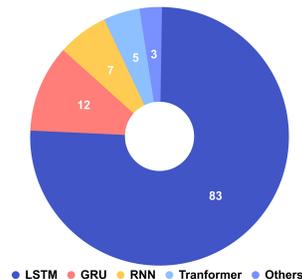


Figure 12: Decoders used over 105 papers revised

To visualize the most frequent methods employed in the reviewed works, Figure 13 shows in a word cloud the most frequent words which correspond to decoder, LSTM, encoder, CNN, visual attention, RNN, and semantic, among others.

but at the end, a unique caption must describe the entire video content. To deal with that, the authors used first a gender annotator from Stanford CoreNLP [Manning et al.(2014)Manning, Surdeanu, Bauer, Finkel, Bethard and McClosky], which assigns likely gender to tokens, then the co-reference resolution was employed with POS (Part of Speech) tagger and lemmatization strategies. Another issue to deal with is the appropriate transition in captions with connective words (e.g., “then” word); to solve that, they ran a CKY parser for PCFG trained on Penn Treebank [Marcinkiewicz(1994)]. Later, the instances, not considering the connective words, are computed as 300-dimensional vectors using the Sentence2Vec [Le and Mikolov(2014)] approach. Ultimately, to do the caption generation, they perform matching of vectorization using the L2 distance. The most outstanding feature of the proposal of Shin et al. is to exploit the content of the video by segments of different sizes, to look for action, not classifying them, but localizing them. The above involves a lot of specific and cumbersome (at least we think) tasks to training specifically. One may think, if there is no action in the video, what will be the result? In that case, they only employ the middle frame to generate the video caption. As well as the solution proposed in the first place (see [Zhao et al.(2019)Zhao, Li and Lu]), Shin et al. tested their approach with three different datasets.

Finally, the third-place approach by Xiao et al. [Xiao and Shi(2019)] employs a reinforced adaptive attention model (RAAM). In this case, they used CNN and Bi-LSTM networks to generate the video features. They try to avoid attention to each word. Furthermore, Xiao et al. proposed a mixed loss training using the policy gradient not only to minimize the word-level cross-entropy, but also the sentence level. Their main contributions lie in the three unidirectional LSTM layers for the decoding task. As well as their combination of visual and contextual information regarding the next-word generation. In this case, Xiao et al. employed the two most used datasets for training, MSVD, and MPII. In general, the improvements of the architecture proposed by Xiao et al. are concrete; we think they take care of more detailed aspects than the first and second-best approaches analyzed in this review.

5 Areas to be improved

There are some opportunity areas that could be attended. Some datasets that have been tested achieving low results could be the source of the lack of generalization of the proposed methods to deal with this task. Some video content features are sometimes missed or non-used, especially those related to speech recognition and movement cues in general.

The lack of exploitation of prior upstanding works related to video action recognition, video motions, trajectories analysis, object, and classification is poor. We saw that best works take advantage of many motion features related to time or temporal cues. Another essential aspect that could be improved is the evaluation metrics used to report performance (and therefore rank the models), because the most used are those not related to semantic aspect; we are aware of new metrics such as BLEURT [Sellam et al.(2020)Sellam, Das and Parikh], BERTScore [Zhang et al.(2019a)Zhang, Kishore, Wu, Weinberger and Artzi], BERTr [Mathur et al.(2019)Mathur, Baldwin and Cohn], CLIPScore [Hessel et al.(2021)Hessel, Holtzman, Forbes, Bras and Choi] which use semantic or contextual embeddings to evaluate the semantic relationship of words. These approaches take advantage of considering the whole sentence for their embeddings (*e.g.*, BERTScore’s bidirectionality) being able to distinguish between polysemous words and homonyms, due to the context of the description (or references in the training phase).

As it had been exposed, the first generation, such as the BLEU metric, relies on measuring the superficial similarity between the sentences, like lexical or grammatical matching instead of semantics. In this case, none of the 105 revised papers used semantically related metrics. All the reviewed publications evaluate their approaches based on the classical first-generation (n-grams based) metrics. Nowadays, although BLEU or CIDEr-D metrics are still commonly used on tasks such as Image Captioning, the new generation of metrics is starting to have more relevance, but in Video Captioning it is not the case yet.

Through all the papers reviewed, there are not many application problems. Most of the works

consulted only apply their methods to the publicly available datasets leaving their work on a basic-science level. We find a huge opportunity to take these models into real-life scenarios to assess their performance in such tasks. We think only reporting performance in terms of known standard metrics is the first step, but further efforts must be done to apply such models to actual practical implementations.

Among the few works that go a step beyond, the video captioning task is mainly applied to the description or synopsis of films to help people with visual impairment; or for the description of recipes. Its use has also been extended to narrate actions in sports, but there are few related works mentioned. Moreover, other areas have been barely explored, for example, smart video surveillance. This application would be helpful in the security sector since a natural language report could be generated on outstanding events or unusual actions.

We acknowledge that state-of-the-art is still far from the point where there would be a *one-size-fits-all* solution for different video captioning tasks. For instance, the description of videos with lots of information in a short time-lapse implies one type of challenge, since there are more sequences to process and detect changes and generate enough captions from it, for example, in sports narration. On the other hand, in surveillance videos that have a longer duration, the scenes to describe can be comparatively shorter; but the challenge is to detect the relevant events to be captioned and do so in real-time, to take proper immediate action.

We can also mention automatic subtitling of videos, scenes, or camera footage for visually impaired people, where the need for real-time performance meets the requirements of extracting meaningful information while discarding not relevant actions. If we take this a step forward, aiding people to navigate cumbersome environments demands the models' safety and robustness, maximizing true-positive rates for risk detection and minimizing the false-positive rates for user comfort.

We do not find any glance of tackling the needs mentioned above at this moment and on the reviewed papers. Hence, we believe huge improvements should be sought in this direction to understand each approach's range of applications truly. We leave this review as an academic effort to gain insight into tracing an overview of the video captioning task, but we would like to contribute to the future by taking these models and research to real-life environments to understand their performance better.

6 Conclusions

In this manuscript was presented a comparative review with a revision from the years 2016 to 2021. In this literature revision, more than 105 papers were analyzed and 105 works were used for a comparison between them. Also, most used datasets and metrics are evaluated and described. Through this period, we compared all the works based on different metrics, namely, METEOR, BLEU-4, and CIDEr-D. They were also compared taking into account the reported dataset used to evaluate the proposed methods. To obtain a winner over this analyzed period, we generate an easy ranking scheme where those works achieving the best results on each dataset got high scores, and with more datasets used in their experiments, the score they get is higher. We calculate, based on this ranking procedure, the best three methods, which are the ones reported by Zhao et al. [Zhao et al.(2019)Zhao, Li and Lu], Shin et al. [Shin et al.(2016)Shin, Ohnishi and Harada], and Xiao et al. [Xiao and Shi(2019)], first, second, and third place, respectively. Finally, after the complete analysis, some insights and improvement opportunities are mentioned.

CRedit authorship contribution statement

Daniela Moctezuma: Conceptualization, Methodology, Data Curation, Software, Investigation, Writing – original draft, Visualization, Funding acquisition. **Tania Ramírez-delReal:** Conceptualization, Methodology, Data Curation, Investigation, Writing – original draft, Visualization. **Guillermo Ruiz:** Writing – review & editing, Visualization, Validation. **Othón González:**

Visualization, Writing – review & editing, Visualization, Validation.

Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Acknowledgments

This work has been done through CONACYT (National Council of Science and Technology from Mexico) support with Ciencia Básica grant with project ID A1-S-34811.

Appendix

Tables show all the papers reviewed organized by dataset.

References

- [Aafaq et al.(2019)Aafaq, Akhtar, Liu, Gilani and Mian] Aafaq, N., Akhtar, N., Liu, W., Gilani, S.Z., Mian, A., 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, pp. 12487–12496.
- [Ballas et al.(2015)Ballas, Yao, Pal and Courville] Ballas, N., Yao, L., Pal, C., Courville, A., 2015. Delving deeper into convolutional networks for learning video representations. arXiv preprint arXiv:1511.06432 .
- [Banerjee and Lavie(2005)] Banerjee, S., Lavie, A., 2005. Meteor: An automatic metric for mt evaluation with improved correlation with human judgments, in: Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization, pp. 65–72.
- [Baraldi et al.(2017)Baraldi, Grana and Cucchiara] Baraldi, L., Grana, C., Cucchiara, R., 2017. Hierarchical boundary-aware neural encoder for video captioning, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 1657–1666.
- [Bin et al.(2018)Bin, Yang, Shen, Xie, Shen and Li] Bin, Y., Yang, Y., Shen, F., Xie, N., Shen, H.T., Li, X., 2018. Describing video with attention-based bidirectional lstm. IEEE transactions on cybernetics 49, 2631–2641.
- [Caba Heilbron et al.(2015)Caba Heilbron, Escorcia, Ghanem and Carlos Niebles] Caba Heilbron, F., Escorcia, V., Ghanem, B., Carlos Niebles, J., 2015. Activitynet: A large-scale video benchmark for human activity understanding, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 961–970.
- [Chen and Dolan(2011)] Chen, D., Dolan, W.B., 2011. Collecting highly parallel data for paraphrase evaluation, in: Proceedings of the 49th annual meeting of the association for computational linguistics: human language technologies, pp. 190–200.
- [Chen et al.(2018a)Chen, Li, Zhang and Huang] Chen, M., Li, Y., Zhang, Z., Huang, S., 2018a. Tvt: Two-view transformer network for video captioning, in: Asian Conference on Machine Learning, PMLR. pp. 847–862.

- [Chen et al.(2017)Chen, Chen, Jin and Hauptmann] Chen, S., Chen, J., Jin, Q., Hauptmann, A., 2017. Video captioning with guidance of multimodal latent topics, in: Proceedings of the 25th ACM international conference on Multimedia, pp. 1838–1846.
- [Chen and Jiang(2019)] Chen, S., Jiang, Y.G., 2019. Motion guided spatial attention for video captioning, in: Proceedings of the AAAI conference on artificial intelligence, pp. 8191–8198.
- [Chen and Jiang(2021)] Chen, S., Jiang, Y.G., 2021. Motion guided region message passing for video captioning, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 1543–1552.
- [Chen et al.(2019)Chen, Jin, Chen and Hauptmann] Chen, S., Jin, Q., Chen, J., Hauptmann, A.G., 2019. Generating video descriptions with latent topic guidance. *IEEE Transactions on Multimedia* 21, 2407–2418.
- [Chen et al.(2018b)Chen, Wang, Zhang and Huang] Chen, Y., Wang, S., Zhang, W., Huang, Q., 2018b. Less is more: Picking informative frames for video captioning, in: Proceedings of the European conference on computer vision (ECCV), pp. 358–373.
- [Daskalakis et al.(2018)Daskalakis, Tzelepi and Tefas] Daskalakis, E., Tzelepi, M., Tefas, A., 2018. Learning deep spatiotemporal features for video captioning. *Pattern Recognition Letters* 116, 143–149.
- [Dilawari et al.(2021)Dilawari, Khan, Al-Otaibi, Rehman, Rahman and Nam] Dilawari, A., Khan, M.U.G., Al-Otaibi, Y.D., Rehman, Z.u., Rahman, A.u., Nam, Y., 2021. Natural language description of videos for smart surveillance. *Applied Sciences* 11, 3730.
- [Dong et al.(2016)Dong, Li, Lan, Huo and Snoek] Dong, J., Li, X., Lan, W., Huo, Y., Snoek, C.G., 2016. Early embedding and late reranking for video captioning, in: Proceedings of the 24th ACM international conference on Multimedia, pp. 1082–1086.
- [Dong et al.(2017)Dong, Su, Zhu and Zhang] Dong, Y., Su, H., Zhu, J., Zhang, B., 2017. Improving interpretability of deep neural networks with semantic information, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 4306–4314.
- [Fakoor et al.(2016)Fakoor, Mohamed, Mitchell, Kang and Kohli] Fakoor, R., Mohamed, A.r., Mitchell, M., Kang, S.B., Kohli, P., 2016. Memory-augmented attention modelling for videos. *arXiv preprint arXiv:1611.02261* .
- [Gan et al.(2017)Gan, Gan, He, Pu, Tran, Gao, Carin and Deng] Gan, Z., Gan, C., He, X., Pu, Y., Tran, K., Gao, J., Carin, L., Deng, L., 2017. Semantic compositional networks for visual captioning, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 5630–5639.
- [Gao et al.(2017)Gao, Guo, Zhang, Xu and Shen] Gao, L., Guo, Z., Zhang, H., Xu, X., Shen, H.T., 2017. Video captioning with attention-based lstm and semantic consistency. *IEEE Transactions on Multimedia* 19, 2045–2055.
- [Gao et al.(2019)Gao, Li, Song and Shen] Gao, L., Li, X., Song, J., Shen, H.T., 2019. Hierarchical lstms with adaptive attention for visual captioning. *IEEE transactions on pattern analysis and machine intelligence* 42, 1112–1131.
- [Guadarrama et al.(2013)Guadarrama, Krishnamoorthy, Malkarnenkar, Venugopalan, Mooney, Darrell and Saenko] Guadarrama, S., Krishnamoorthy, N., Malkarnenkar, G., Venugopalan, S., Mooney, R., Darrell, T., Saenko, K., 2013. Youtube2text: Recognizing and describing arbitrary activities using semantic hierarchies and zero-shot recognition, in: Proceedings of the IEEE international conference on computer vision, pp. 2712–2719.

- [Guo et al.(2019)Guo, Zhang and Gao] Guo, Y., Zhang, J., Gao, L., 2019. Exploiting long-term temporal dynamics for video captioning. *World Wide Web* 22, 735–749.
- [Hessel et al.(2021)Hessel, Holtzman, Forbes, Bras and Choi] Hessel, J., Holtzman, A., Forbes, M., Bras, R.L., Choi, Y., 2021. Clipscore: A reference-free evaluation metric for image captioning. *arXiv preprint arXiv:2104.08718* .
- [Hori et al.(2017)Hori, Hori, Lee, Zhang, Harsham, Hershey, Marks and Sumi] Hori, C., Hori, T., Lee, T.Y., Zhang, Z., Harsham, B., Hershey, J.R., Marks, T.K., Sumi, K., 2017. Attention-based multimodal fusion for video description, in: *Proceedings of the IEEE international conference on computer vision*, pp. 4193–4202.
- [Hou et al.(2019)Hou, Wu, Zhao, Luo and Jia] Hou, J., Wu, X., Zhao, W., Luo, J., Jia, Y., 2019. Joint syntax representation learning and visual cue translation for video captioning, in: *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 8918–8927.
- [Hu et al.(2019)Hu, Chen, Zha and Wu] Hu, Y., Chen, Z., Zha, Z.J., Wu, F., 2019. Hierarchical global-local temporal modeling for video captioning, in: *Proceedings of the 27th ACM International Conference on Multimedia*, pp. 774–783.
- [Iashin and Rahtu(2020)] Iashin, V., Rahtu, E., 2020. Multi-modal dense video captioning, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pp. 958–959.
- [Inácio et al.(2021)Inácio, Gutoski, Lazzaretti and Lopes] Inácio, A.D.S., Gutoski, M., Lazzaretti, A.E., Lopes, H.S., 2021. Osvidcap: A framework for the simultaneous recognition and description of concurrent actions in videos in an open-set scenario. *IEEE Access* 9, 137029–137041.
- [Islam et al.(2021)Islam, Dash, Seum, Raj, Hossain and Shah] Islam, S., Dash, A., Seum, A., Raj, A.H., Hossain, T., Shah, F.M., 2021. Exploring video captioning techniques: A comprehensive survey on deep learning methods. *SN Computer Science* 2, 1–28.
- [Jain et al.(2017)Jain, Agarwalla, Agrawal and Mitra] Jain, A.K., Agarwalla, A., Agrawal, K.K., Mitra, P., 2017. Recurrent memory addressing for describing videos., in: *CVPR Workshops*, pp. 1–8.
- [Jain et al.(2022)Jain, Al-Turjman, Chaudhary, Nayar, Gupta and Kumar] Jain, V., Al-Turjman, F., Chaudhary, G., Nayar, D., Gupta, V., Kumar, A., 2022. Video captioning: a review of theory, techniques and practices. *Multimedia Tools and Applications* , 1–35.
- [Ji and Wang(2021)] Ji, W., Wang, R., 2021. A multi-instance multi-label dual learning approach for video captioning. *ACM Transactions on Multimedia Computing Communications and Applications* 17, 1–18.
- [Jin et al.(2016)Jin, Chen, Chen, Xiong and Hauptmann] Jin, Q., Chen, J., Chen, S., Xiong, Y., Hauptmann, A., 2016. Describing videos using multi-modal fusion, in: *Proceedings of the 24th ACM international conference on Multimedia*, pp. 1087–1091.
- [Khan et al.(2012)Khan, Nawab and Gotoh] Khan, M.U.G., Nawab, R.M.A., Gotoh, Y., 2012. Natural language descriptions of visual scenes corpus generation and analysis, in: *Proceedings of the Joint Workshop on Exploiting Synergies Between Information Retrieval and Machine Translation (ESIRMT) and Hybrid Approaches to Machine Translation (HyTra)*, pp. 38–47.
- [Krishna et al.(2017)Krishna, Hata, Ren, Fei-Fei and Carlos Niebles] Krishna, R., Hata, K., Ren, F., Fei-Fei, L., Carlos Niebles, J., 2017. Dense-captioning events in videos, in: *Proceedings of the IEEE international conference on computer vision*, pp. 706–715.

- [Kumar and Mathew(2020)] Kumar, A., Mathew, R., 2020. A review of methods for video captioning. Available at SSRN .
- [Le and Mikolov(2014)] Le, Q., Mikolov, T., 2014. Distributed representations of sentences and documents, in: International conference on machine learning, PMLR. pp. 1188–1196.
- [Lee(2019)] Lee, J.Y., 2019. Deep multimodal embedding for video captioning. *Multimedia Tools and Applications* 78, 31793–31805.
- [Lei and Huang(2021)] Lei, Z., Huang, Y., 2021. Video captioning based on channel soft attention and semantic reconstructor. *Future Internet* 13, 55.
- [Li and Gong(2019)] Li, L., Gong, B., 2019. End-to-end video captioning with multitask reinforcement learning, in: 2019 IEEE Winter Conference on Applications of Computer Vision (WACV), IEEE. pp. 339–348.
- [Li et al.(2018a)Li, Guo and Fang] Li, W., Guo, D., Fang, X., 2018a. Multimodal architecture for video captioning with memory networks and an attention mechanism. *Pattern Recognition Letters* 105, 23–29.
- [Li et al.(2017)Li, Zhao, Lu et al.] Li, X., Zhao, B., Lu, X., et al., 2017. Mam-rnn: Multi-level attention model based rnn for video captioning., in: IJCAI, pp. 2208–2214.
- [Li et al.(2019)Li, Zhou, Chen and Gao] Li, X., Zhou, Z., Chen, L., Gao, L., 2019. Residual attention-based lstm for video captioning. *World Wide Web* 22, 621–636.
- [Li et al.(2018b)Li, Yao, Pan, Chao and Mei] Li, Y., Yao, T., Pan, Y., Chao, H., Mei, T., 2018b. Jointly localizing and describing events for dense video captioning, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 7492–7500.
- [Lin(2004)] Lin, C.Y., 2004. Rouge: A package for automatic evaluation of summaries, in: Text summarization branches out, pp. 74–81.
- [Lin et al.(2014)Lin, Maire, Belongie, Hays, Perona, Ramanan, Dollár and Zitnick] Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L., 2014. Microsoft coco: Common objects in context, in: European conference on computer vision, Springer. pp. 740–755.
- [Liu et al.(2017a)Liu, Xu, Wong, Li, Su and Kankanhalli] Liu, A.A., Xu, N., Wong, Y., Li, J., Su, Y.T., Kankanhalli, M., 2017a. Hierarchical & multimodal video captioning: Discovering and transferring multimodal knowledge for vision to language. *Computer Vision and Image Understanding* 163, 113–125.
- [Liu et al.(2020)Liu, Ren and Yuan] Liu, S., Ren, Z., Yuan, J., 2020. Sibnet: Sibling convolutional encoder for video captioning. *IEEE transactions on pattern analysis and machine intelligence* .
- [Liu et al.(2017b)Liu, Li and Shi] Liu, Y., Li, X., Shi, Z., 2017b. Video captioning with listwise supervision, in: Proceedings of the AAAI Conference on Artificial Intelligence, pp. 4197–42203.
- [Long et al.(2018)Long, Gan and De Melo] Long, X., Gan, C., De Melo, G., 2018. Video captioning with multi-faceted attention. *Transactions of the Association for Computational Linguistics* 6, 173–184.
- [Manning et al.(2014)Manning, Surdeanu, Bauer, Finkel, Bethard and McClosky] Manning, C.D., Surdeanu, M., Bauer, J., Finkel, J.R., Bethard, S., McClosky, D., 2014. The stanford corenlp natural language processing toolkit, in: Proceedings of 52nd annual meeting of the association for computational linguistics: system demonstrations, pp. 55–60.

- [Marcinkiewicz(1994)] Marcinkiewicz, M.A., 1994. Building a large annotated corpus of english: The penn treebank. Using Large Corpora , 273.
- [Mathur et al.(2019)Mathur, Baldwin and Cohn] Mathur, N., Baldwin, T., Cohn, T., 2019. Putting evaluation in context: Contextual embeddings improve machine translation evaluation, in: Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, pp. 2799–2808.
- [Memon et al.(2009)Memon, Lech and He] Memon, S., Lech, M., He, L., 2009. Using information theoretic vector quantization for inverted mfcc based speaker verification, in: 2009 2nd International Conference on Computer, Control and Communication, IEEE. pp. 1–5.
- [Mun et al.(2019)Mun, Yang, Ren, Xu and Han] Mun, J., Yang, L., Ren, Z., Xu, N., Han, B., 2019. Streamlined dense video captioning, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 6588–6597.
- [Nian et al.(2017)Nian, Li, Wang, Wu, Ni and Xu] Nian, F., Li, T., Wang, Y., Wu, X., Ni, B., Xu, C., 2017. Learning explicit video attributes from mid-level representation for video captioning. Computer Vision and Image Understanding 163, 126–138.
- [Pan et al.(2020)Pan, Cai, Huang, Lee, Gaidon, Adeli and Niebles] Pan, B., Cai, H., Huang, D.A., Lee, K.H., Gaidon, A., Adeli, E., Niebles, J.C., 2020. Spatio-temporal graph for video captioning with knowledge distillation, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 10870–10879.
- [Pan et al.(2016a)Pan, Xu, Yang, Wu and Zhuang] Pan, P., Xu, Z., Yang, Y., Wu, F., Zhuang, Y., 2016a. Hierarchical recurrent neural encoder for video representation with application to captioning, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 1029–1038.
- [Pan et al.(2016b)Pan, Mei, Yao, Li and Rui] Pan, Y., Mei, T., Yao, T., Li, H., Rui, Y., 2016b. Jointly modeling embedding and translation to bridge video and language, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 4594–4602.
- [Pan et al.(2017)Pan, Yao, Li and Mei] Pan, Y., Yao, T., Li, H., Mei, T., 2017. Video captioning with transferred semantic attributes, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 6504–6512.
- [Papineni et al.(2002)Papineni, Roukos, Ward and Zhu] Papineni, K., Roukos, S., Ward, T., Zhu, W.J., 2002. Bleu: a method for automatic evaluation of machine translation, in: Proceedings of the 40th annual meeting of the Association for Computational Linguistics, pp. 311–318.
- [Pasunuru and Bansal(2017)] Pasunuru, R., Bansal, M., 2017. Multi-task video captioning with video and entailment generation. arXiv preprint arXiv:1704.07489 .
- [Pei et al.(2019)Pei, Zhang, Wang, Ke, Shen and Tai] Pei, W., Zhang, J., Wang, X., Ke, L., Shen, X., Tai, Y.W., 2019. Memory-attended recurrent network for video captioning, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 8347–8356.
- [Perez-Martin et al.(2021)Perez-Martin, Bustos and Pérez] Perez-Martin, J., Bustos, B., Pérez, J., 2021. Improving video captioning with temporal composition of a visual-syntactic embedding, in: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision, pp. 3039–3049.
- [Phan et al.(2017)Phan, Henter, Miyao and Satoh] Phan, S., Henter, G.E., Miyao, Y., Satoh, S., 2017. Consensus-based sequence training for video captioning. arXiv preprint arXiv:1712.09532 .

- [Pu et al.(2018)Pu, Min, Gan and Carin] Pu, Y., Min, M.R., Gan, Z., Carin, L., 2018. Adaptive feature abstraction for translating video to text, in: Thirty-Second AAAI Conference on Artificial Intelligence, pp. 1–13.
- [Qi et al.(2019)Qi, Wang, Li and Luo] Qi, M., Wang, Y., Li, A., Luo, J., 2019. Sports video captioning via attentive motion representation and group relationship modeling. *IEEE Transactions on Circuits and Systems for Video Technology* 30, 2617–2633.
- [Ramanishka et al.(2016)Ramanishka, Das, Park, Venugopalan, Hendricks, Rohrbach and Saenko] Ramanishka, V., Das, A., Park, D.H., Venugopalan, S., Hendricks, L.A., Rohrbach, M., Saenko, K., 2016. Multimodal video description, in: Proceedings of the 24th ACM international conference on Multimedia, pp. 1092–1096.
- [Regneri et al.(2013)Regneri, Rohrbach, Wetzel, Thater, Schiele and Pinkal] Regneri, M., Rohrbach, M., Wetzel, D., Thater, S., Schiele, B., Pinkal, M., 2013. Grounding action descriptions in videos. *Transactions of the Association for Computational Linguistics* 1, 25–36.
- [Robertson(2004)] Robertson, S., 2004. Understanding inverse document frequency: on theoretical arguments for idf. *Journal of documentation* .
- [Rohrbach et al.(2014)Rohrbach, Rohrbach, Qiu, Friedrich, Pinkal and Schiele] Rohrbach, A., Rohrbach, M., Qiu, W., Friedrich, A., Pinkal, M., Schiele, B., 2014. Coherent multi-sentence video description with variable level of detail, in: German conference on pattern recognition, Springer. pp. 184–195.
- [Rohrbach et al.(2015a)Rohrbach, Rohrbach, Tandon and Schiele] Rohrbach, A., Rohrbach, M., Tandon, N., Schiele, B., 2015a. A dataset for movie description, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 3202–3212.
- [Rohrbach et al.(2017)Rohrbach, Torabi, Rohrbach, Tandon, Pal, Larochelle, Courville and Schiele] Rohrbach, A., Torabi, A., Rohrbach, M., Tandon, N., Pal, C., Larochelle, H., Courville, A., Schiele, B., 2017. Movie description. *International Journal of Computer Vision* 123, 94–120.
- [Rohrbach et al.(2015b)Rohrbach, Rohrbach, Regneri, Amin, Andriluka, Pinkal and Schiele] Rohrbach, M., Rohrbach, A., Regneri, M., Amin, S., Andriluka, M., Pinkal, M., Schiele, B., 2015b. Recognizing fine-grained and composite activities using hand-centric features and script data. *International Journal of Computer Vision* , 1–28URL: <http://dx.doi.org/10.1007/s11263-015-0851-8>, doi:10.1007/s11263-015-0851-8.
- [Sahidullah and Saha(2009)] Sahidullah, M., Saha, G., 2009. On the use of distributed dct in speaker identification, in: 2009 Annual IEEE India Conference, IEEE. pp. 1–4.
- [Sellam et al.(2020)Sellam, Das and Parikh] Sellam, T., Das, D., Parikh, A., 2020. BLEURT: Learning robust metrics for text generation, in: Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, Association for Computational Linguistics, Online. pp. 7881–7892. URL: <https://aclanthology.org/2020.acl-main.704>, doi:10.18653/v1/2020.acl-main.704.
- [Shen et al.(2017)Shen, Li, Su, Li, Chen, Jiang and Xue] Shen, Z., Li, J., Su, Z., Li, M., Chen, Y., Jiang, Y.G., Xue, X., 2017. Weakly supervised dense video captioning, in: Proceedings of the IEEE conference on Computer Vision and Pattern Recognition, pp. 1916–1924.
- [Shetty and Laaksonen(2016)] Shetty, R., Laaksonen, J., 2016. Frame-and segment-level features and candidate pool evaluation for video caption generation, in: Proceedings of the 24th ACM international conference on Multimedia, pp. 1073–1076.

- [Shi et al.(2020)Shi, Cai, Gu and Joty] Shi, X., Cai, J., Gu, J., Joty, S., 2020. Video captioning with boundary-aware hierarchical language decoding and joint video prediction. *Neurocomputing* 417, 347–356.
- [Shi et al.(2019)Shi, Cai, Joty and Gu] Shi, X., Cai, J., Joty, S., Gu, J., 2019. Watch it twice: Video captioning with a refocused video encoder, in: *Proceedings of the 27th ACM International Conference on Multimedia*, pp. 818–826.
- [Shin et al.(2016)Shin, Ohnishi and Harada] Shin, A., Ohnishi, K., Harada, T., 2016. Beyond caption to narrative: Video captioning with multiple sentences, in: *2016 IEEE International Conference on Image Processing (ICIP)*, IEEE. pp. 3364–3368.
- [Sigurdsson et al.(2016)Sigurdsson, Varol, Wang, Farhadi, Laptev and Gupta] Sigurdsson, G.A., Varol, G., Wang, X., Farhadi, A., Laptev, I., Gupta, A., 2016. Hollywood in homes: Crowdsourcing data collection for activity understanding, in: *European Conference on Computer Vision*, Springer. pp. 510–526.
- [Song et al.(2018)Song, Guo, Gao, Li, Hanjalic and Shen] Song, J., Guo, Y., Gao, L., Li, X., Hanjalic, A., Shen, H.T., 2018. From deterministic to generative: Multimodal stochastic rnns for video captioning. *IEEE transactions on neural networks and learning systems* 30, 3047–3058.
- [Song et al.(2017)Song, Guo, Gao, Liu, Zhang and Shen] Song, J., Guo, Z., Gao, L., Liu, W., Zhang, D., Shen, H.T., 2017. Hierarchical lstm with adjusted temporal attention for video captioning. *arXiv preprint arXiv:1706.01231* .
- [Suin and Rajagopalan(2020)] Suin, M., Rajagopalan, A., 2020. An efficient framework for dense video captioning, in: *Proceedings of the AAAI Conference on Artificial Intelligence*, pp. 12039–12046.
- [Tang et al.(2019)Tang, Wang and Li] Tang, P., Wang, H., Li, Q., 2019. Rich visual and language representation with complementary semantics for video captioning. *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)* 15, 1–23.
- [Torabi et al.(2015)Torabi, Pal, Larochelle and Courville] Torabi, A., Pal, C., Larochelle, H., Courville, A., 2015. Using descriptive video services to create a large data source for video annotation research. *arXiv preprint arXiv:1503.01070* .
- [Tu et al.(2017)Tu, Zhang, Liu and Yan] Tu, Y., Zhang, X., Liu, B., Yan, C., 2017. Video description with spatial-temporal attention, in: *Proceedings of the 25th ACM international conference on Multimedia*, pp. 1014–1022.
- [Vedantam et al.(2015)Vedantam, Lawrence Zitnick and Parikh] Vedantam, R., Lawrence Zitnick, C., Parikh, D., 2015. Cider: Consensus-based image description evaluation, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 4566–4575.
- [Venugopalan et al.(2016)Venugopalan, Hendricks, Mooney and Saenko] Venugopalan, S., Hendricks, L.A., Mooney, R., Saenko, K., 2016. Improving lstm-based video description with linguistic knowledge mined from text. *arXiv preprint arXiv:1604.01729* .
- [Wang et al.(2019a)Wang, Ma, Zhang, Jiang, Wang and Liu] Wang, B., Ma, L., Zhang, W., Jiang, W., Wang, J., Liu, W., 2019a. Controllable video captioning with pos sequence guidance based on gated fusion network, in: *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 2641–2650.
- [Wang et al.(2018a)Wang, Ma, Zhang and Liu] Wang, B., Ma, L., Zhang, W., Liu, W., 2018a. Reconstruction network for video captioning, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 7622–7631.

- [Wang and Schmid(2013)] Wang, H., Schmid, C., 2013. Action recognition with improved trajectories, in: Proceedings of the IEEE international conference on computer vision, pp. 3551–3558.
- [Wang et al.(2018b)Wang, Jiang, Ma, Liu and Xu] Wang, J., Jiang, W., Ma, L., Liu, W., Xu, Y., 2018b. Bidirectional attentive fusion with context gating for dense video captioning, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 7190–7198.
- [Wang et al.(2018c)Wang, Wang, Huang, Wang and Tan] Wang, J., Wang, W., Huang, Y., Wang, L., Tan, T., 2018c. Hierarchical memory modelling for video captioning, in: Proceedings of the 26th ACM international conference on Multimedia, pp. 63–71.
- [Wang et al.(2018d)Wang, Wang, Huang, Wang and Tan] Wang, J., Wang, W., Huang, Y., Wang, L., Tan, T., 2018d. M3: Multimodal memory modelling for video captioning, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 7512–7520.
- [Wang et al.(2018e)Wang, Chen, Wu, Wang and Wang] Wang, X., Chen, W., Wu, J., Wang, Y.F., Wang, W.Y., 2018e. Video captioning via hierarchical reinforcement learning, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 4213–4222.
- [Wang et al.(2019b)Wang, Wu, Chen, Li, Wang and Wang] Wang, X., Wu, J., Chen, J., Li, L., Wang, Y.F., Wang, W.Y., 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, pp. 4581–4591.
- [Wei et al.(2020)Wei, Mi, Hu and Chen] Wei, R., Mi, L., Hu, Y., Chen, Z., 2020. Exploiting the local temporal information for video captioning. *Journal of Visual Communication and Image Representation* 67, 102751.
- [Wolf et al.(2014)Wolf, Lombardi, Mille, Celiktutan, Jiu, Dogan, Eren, Baccouche, Dellandréa, Bichot et al.] Wolf, C., Lombardi, E., Mille, J., Celiktutan, O., Jiu, M., Dogan, E., Eren, G., Baccouche, M., Dellandréa, E., Bichot, C.E., et al., 2014. Evaluation of video activity localizations integrating quality and quantity measurements. *Computer Vision and Image Understanding* 127, 14–30.
- [Wu and Han(2018)] Wu, A., Han, Y., 2018. Multi-modal circulant fusion for video-to-language and backward., in: *IJCAI*, p. 8.
- [Wu et al.(2019)Wu, Han, Yang, Hu and Wu] Wu, A., Han, Y., Yang, Y., Hu, Q., Wu, F., 2019. Convolutional reconstruction-to-sequence for video captioning. *IEEE Transactions on Circuits and Systems for Video Technology* 30, 4299–4308.
- [Wu et al.(2018)Wu, Li, Cao, Ji and Lin] Wu, X., Li, G., Cao, Q., Ji, Q., Lin, L., 2018. Interpretable video captioning via trajectory structured localization, in: Proceedings of the IEEE conference on Computer Vision and Pattern Recognition, pp. 6829–6837.
- [Xiao and Shi(2019)] Xiao, H., Shi, J., 2019. Video captioning with adaptive attention and mixed loss optimization. *IEEE Access* 7, 135757–135769.
- [Xiao and Shi(2020)] Xiao, H., Shi, J., 2020. Video captioning with text-based dynamic attention and step-by-step learning. *Pattern Recognition Letters* 133, 305–312.
- [Xiong et al.(2018)Xiong, Dai and Lin] Xiong, Y., Dai, B., Lin, D., 2018. Move forward and tell: A progressive generator of video descriptions, in: Proceedings of the European Conference on Computer Vision (ECCV), pp. 468–483.

- [Xu et al.(2016)Xu, Mei, Yao and Rui] Xu, J., Mei, T., Yao, T., Rui, Y., 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, pp. 5288–5296.
- [Xu et al.(2017)Xu, Yao, Zhang and Mei] Xu, J., Yao, T., Zhang, Y., Mei, T., 2017. Learning multimodal attention lstm networks for video captioning, in: Proceedings of the 25th ACM international conference on Multimedia, pp. 537–545.
- [Xu et al.(2019)Xu, Liu, Nie and Su] Xu, N., Liu, A.A., Nie, W., Su, Y., 2019. Multi-guiding long short-term memory for video captioning. *Multimedia Systems* 25, 663–672.
- [Xu et al.(2018a)Xu, Liu, Wong, Zhang, Nie, Su and Kankanhalli] Xu, N., Liu, A.A., Wong, Y., Zhang, Y., Nie, W., Su, Y., Kankanhalli, M., 2018a. Dual-stream recurrent neural network for video captioning. *IEEE Transactions on Circuits and Systems for Video Technology* 29, 2482–2493.
- [Xu et al.(2018b)Xu, Han, Hong and Tian] Xu, Y., Han, Y., Hong, R., Tian, Q., 2018b. Sequential video vlad: Training the aggregation locally and temporally. *IEEE Transactions on Image Processing* 27, 4933–4944.
- [Yan et al.(2019)Yan, Tu, Wang, Zhang, Hao, Zhang and Dai] Yan, C., Tu, Y., Wang, X., Zhang, Y., Hao, X., Zhang, Y., Dai, Q., 2019. Stat: Spatial-temporal attention mechanism for video captioning. *IEEE transactions on multimedia* 22, 229–241.
- [Yang et al.(2018)Yang, Zhou, Ai, Bin, Hanjalic, Shen and Ji] Yang, Y., Zhou, J., Ai, J., Bin, Y., Hanjalic, A., Shen, H.T., Ji, Y., 2018. Video captioning by adversarial lstm. *IEEE Transactions on Image Processing* 27, 5600–5611.
- [Yang et al.(2017)Yang, Han and Wang] Yang, Z., Han, Y., Wang, Z., 2017. Catching the temporal regions-of-interest for video captioning, in: Proceedings of the 25th ACM international conference on Multimedia, pp. 146–153.
- [Yu et al.(2018)Yu, Cheng, Ni, Wang, Zhang and Yang] Yu, H., Cheng, S., Ni, B., Wang, M., Zhang, J., Yang, X., 2018. Fine-grained video captioning for sports narrative, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 6006–6015.
- [Yu et al.(2016)Yu, Wang, Huang, Yang and Xu] Yu, H., Wang, J., Huang, Z., Yang, Y., Xu, W., 2016. Video paragraph captioning using hierarchical recurrent neural networks, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 4584–4593.
- [Yu et al.(2017a)Yu, Choi, Kim, Yoo, Lee and Kim] Yu, Y., Choi, J., Kim, Y., Yoo, K., Lee, S.H., Kim, G., 2017a. Supervising neural attention models for video captioning by human gaze data, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 490–498.
- [Yu et al.(2017b)Yu, Ko, Choi and Kim] Yu, Y., Ko, H., Choi, J., Kim, G., 2017b. End-to-end concept word detection for video captioning, retrieval, and question answering, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 3165–3173.
- [Yu and Han(2021)] Yu, Z., Han, N., 2021. Accelerated masked transformer for dense video captioning. *Neurocomputing* 445, 72–80.
- [Yuan et al.(2020)Yuan, Ma, Wang and Zhu] Yuan, Y., Ma, L., Wang, J., Zhu, W., 2020. Controllable video captioning with an exemplar sentence, in: Proceedings of the 28th ACM International Conference on Multimedia, pp. 1085–1093.

- [Zanfir et al.(2016)Zanfir, Marinoiu and Sminchisescu] Zanfir, M., Marinoiu, E., Sminchisescu, C., 2016. Spatio-temporal attention models for grounded video captioning, in: asian conference on computer vision, Springer. pp. 104–119.
- [Zhang and Tian(2016)] Zhang, C., Tian, Y., 2016. Automatic video description generation via lstm with joint two-stream encoding, in: 2016 23rd International Conference on Pattern Recognition (ICPR), IEEE. pp. 2924–2929.
- [Zhang and Peng(2019a)] Zhang, J., Peng, Y., 2019a. Hierarchical vision-language alignment for video captioning, in: International Conference on Multimedia Modeling, Springer. pp. 42–54.
- [Zhang and Peng(2019b)] Zhang, J., Peng, Y., 2019b. Object-aware aggregation with bidirectional temporal graph for video captioning, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 8327–8336.
- [Zhang et al.(2019a)Zhang, Kishore, Wu, Weinberger and Artzi] Zhang, T., Kishore, V., Wu, F., Weinberger, K.Q., Artzi, Y., 2019a. Bertscore: Evaluating text generation with bert. arXiv preprint arXiv:1904.09675 .
- [Zhang et al.(2019b)Zhang, Wang, Ma and Liu] Zhang, W., Wang, B., Ma, L., Liu, W., 2019b. Reconstruct and represent video contents for captioning via reinforcement learning. IEEE transactions on pattern analysis and machine intelligence 42, 3088–3101.
- [Zhang et al.(2017)Zhang, Gao, Zhang, Zhang, Li and Tian] Zhang, X., Gao, K., Zhang, Y., Zhang, D., Li, J., Tian, Q., 2017. Task-driven dynamic fusion: Reducing ambiguity in video description, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 3713–3721.
- [Zhang et al.(2020)Zhang, Shi, Yuan, Li, Wang, Hu and Zha] Zhang, Z., Shi, Y., Yuan, C., Li, B., Wang, P., Hu, W., Zha, Z.J., 2020. Object relational graph with teacher-recommended learning for video captioning, in: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp. 13278–13288.
- [Zhang et al.(2019c)Zhang, Xu, Ouyang and Tan] Zhang, Z., Xu, D., Ouyang, W., Tan, C., 2019c. Show, tell and summarize: Dense video captioning using visual cue aided sentence summarization. IEEE Transactions on Circuits and Systems for Video Technology 30, 3130–3139.
- [Zhao et al.(2019)Zhao, Li and Lu] Zhao, B., Li, X., Lu, X., 2019. Cam-rnn: Co-attention model based rnn for video captioning. IEEE Transactions on Image Processing 28, 5552–5565.
- [Zhao et al.(2018)Zhao, Li, Lu et al.] Zhao, B., Li, X., Lu, X., et al., 2018. Video captioning with tube features., in: IJCAI, pp. 1177–1183.
- [Zheng et al.(2020)Zheng, Wang and Tao] Zheng, Q., Wang, C., Tao, D., 2020. Syntax-aware action targeting for video captioning, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 13096–13105.
- [Zhou et al.(2018a)Zhou, Xu and Corso] Zhou, L., Xu, C., Corso, J.J., 2018a. Towards automatic learning of procedures from web instructional videos, in: Thirty-Second AAAI Conference on Artificial Intelligence, pp. 7590–7598.
- [Zhou et al.(2018b)Zhou, Zhou, Corso, Socher and Xiong] Zhou, L., Zhou, Y., Corso, J.J., Socher, R., Xiong, C., 2018b. End-to-end dense video captioning with masked transformer, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 8739–8748.
- [Zhu et al.(2020)Zhu, Hwang, Ma, Chen and Guo] Zhu, F., Hwang, J.N., Ma, Z., Chen, G., Guo, J., 2020. Understanding objects in video: Object-oriented video captioning via structured trajectory and adversarial learning. IEEE Access 8, 169146–169159.

- [Zhu et al.(2017)Zhu, Xu and Yang] Zhu, L., Xu, Z., Yang, Y., 2017. Bidirectional multirate reconstruction for temporal modeling in videos, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 2653–2662.
- [Zolfaghari et al.(2018)Zolfaghari, Singh and Brox] Zolfaghari, M., Singh, K., Brox, T., 2018. Eco: Efficient convolutional network for online video understanding, in: Proceedings of the European conference on computer vision (ECCV), pp. 695–712.

Work	BLEU-1	BLEU-2	BLEU-3	BLEU-4	METEOR	CIDERr-D	ROUGE-L	Average-Recall (AR)	SPICE	FCE	RE	Self-BLEU
[Shin et al.(2016)Shin, Ohnishi and Harada]	-	-	-	6.30	10.40	17.70	-	-	-	-	-	-
[Pan et al.(2016a)Pan, Xu, Yang, Wu and Zhuang]	81.10	68.60	57.80	46.70	33.90	-	-	-	-	-	-	-
[Yu et al.(2016)Yu, Wang, Huang, Yang and Xu]	81.50	70.40	60.50	49.90	32.60	65.80	-	-	-	-	-	-
[Zanfir et al.(2016)Zanfir, Marinou and Sminchisescu]	81.50	70.80	61.50	50.60	32.40	-	-	-	-	-	-	-
[Fakoor et al.(2016)Fakoor, Mohamed, Mitchell, Kang and Kohli]	79.40	67.10	56.80	46.10	31.80	62.70	-	-	-	-	-	-
[Pan et al.(2016b)Pan, Mei, Yao, Li and Rui]	78.80	66.00	55.40	45.30	31.00	-	-	-	-	-	-	-
[Venugopalan et al.(2016)Venugopalan, Hendricks, Mooney and Saenko]	-	-	-	42.10	31.40	-	-	-	-	-	-	-
[Ballas et al.(2015)Ballas, Yao, Pal and Courville]	-	-	-	43.26	31.60	68.01	-	-	-	-	-	-
[Zhang and Tian(2016)]	-	-	-	-	31.10	-	-	-	-	-	-	-
[Yang et al.(2017)Yang, Han and Wang]	-	-	-	51.10	33.60	74.80	-	-	-	-	-	-
[Baraldi et al.(2017)Baraldi, Grana and Cucchiara]	-	-	-	42.50	32.40	63.50	-	-	-	-	-	-
[Xu et al.(2017)Xu, Yao, Zhang and Mei]	82.30	71.10	61.80	52.30	33.60	70.40	-	-	-	-	-	-
[Pasunuru and Bansal(2017)]	-	-	-	54.50	36.00	92.40	72.80	-	-	-	-	-
[Zhang et al.(2017)Zhang, Gao, Zhang, Zhang, Li and Tian]	-	-	-	45.80	33.30	73.00	69.70	-	-	-	-	-
[Gao et al.(2017)Gao, Guo, Zhang, Xu and Shen]	81.80	70.80	61.10	50.80	33.30	74.80	-	-	-	-	-	-
[Liu et al.(2017a)Liu, Li and Shi]	80.20	69.00	60.10	51.10	32.60	-	-	-	-	-	-	-
[Pan et al.(2017)Pan, Yao, Li and Mei]	-	-	-	54.50	36.00	92.40	72.80	-	-	-	-	-
[Tu et al.(2017)Tu, Zhang, Liu and Yan]	82.60	71.40	61.60	51.10	32.70	67.50	-	-	-	-	-	-
[Gan et al.(2017)Gan, Gan, He, Pu, Tran, Gao, Carin and Deng]	-	-	-	51.10	33.50	77.70	-	-	-	-	-	-
[Jain et al.(2017)Jain, Agarwalla, Agrawal and Mitra]	-	-	-	45.70	31.90	57.30	-	-	-	-	-	-
[Dong et al.(2017)Dong, Su, Zhu and Zhang]	-	-	-	44.60	29.70	-	-	-	-	-	-	-
[Song et al.(2017)Song, Guo, Gao, Liu, Zhang and Shen]	82.90	72.20	63.00	53.00	33.60	73.80	-	-	-	-	-	-
[Liu et al.(2017a)Liu, Xu, Wong, Li, Su and Kankanhalli]	78.00	65.20	54.90	44.30	32.10	68.40	68.90	-	-	-	-	-
[Zhu et al.(2017)Zhu, Xu and Yang]	-	-	-	49.45	33.39	75.45	-	-	-	-	-	-
[Song et al.(2018)Song, Guo, Gao, Li, Hanjalic and Shen]	82.90	72.60	63.50	53.30	33.80	74.80	-	-	-	-	-	-
[Wang et al.(2018c)Wang, Wang, Huang, Wang and Tan]	-	-	-	52.90	33.80	74.50	-	-	-	-	-	-
[Zhang and Peng(2019a)]	83.10	73.00	64.30	55.10	35.30	83.30	-	-	-	-	-	-
[Daskalakis et al.(2018)Daskalakis, Tzelepi and Tefas]	78.11	66.43	55.93	45.02	33.80	63.28	69.62	-	-	-	-	-
[Chen et al.(2018b)Chen, Wang, Zhang and Huang]	-	-	-	52.30	33.30	76.50	69.60	-	-	-	-	-
[Wang et al.(2018d)Wang, Wang, Huang, Wang and Tan]	82.45	72.43	62.78	52.82	33.31	-	-	-	-	-	-	-
[Wu et al.(2018)Wu, Li, Cao, Ji and Lin]	-	-	-	51.70	34.00	74.90	-	-	-	-	-	-
[Long et al.(2018)Long, Gan and De Melo]	83.00	71.90	63.00	52.00	33.50	72.10	-	-	-	-	-	-
[Yang et al.(2018)Yang, Zhou, Ai, Bin, Hanjalic, Shen and Ji]	-	-	-	42.90	30.40	-	-	-	-	-	-	-
[Wang et al.(2018a)Wang, Ma, Zhang and Liu]	-	-	-	52.30	34.10	80.30	69.80	-	-	-	-	-
[Li et al.(2018a)Li, Guo and Fang]	-	-	-	48.00	31.60	68.80	-	-	-	-	-	-
[Wu and Han(2018)]	-	-	-	46.46	33.72	75.46	-	-	-	-	-	-
[Bin et al.(2018)Bin, Yang, Shen, Xie, Shen and Li]	79.00	60.50	48.40	37.30	30.30	-	-	-	-	-	-	-
[Xu et al.(2018b)Xu, Han, Hong and Tian]	-	-	-	51.00	35.15	86.04	-	-	-	-	-	-
[Pu et al.(2018)Pu, Min, Gan and Carin]	-	-	-	54.27	38.03	78.31	-	-	-	-	-	-
[Yan et al.(2019)Yan, Tu, Wang, Zhang, Hao, Zhang and Dai]	-	-	-	52.00	33.33	73.80	-	-	-	-	-	-
[Wu et al.(2019)Wu, Han, Yang, Hu and Wu]	-	-	-	54.10	35.15	82.75	-	-	-	-	-	-
[Xu et al.(2018a)Xu, Liu, Wong, Zhang, Nie, Su and Kankanhalli]	-	-	-	53.00	34.70	79.40	65.90	-	-	-	-	-
[Shi et al.(2019)Shi, Cai, Joty and Gu]	-	-	-	51.70	34.30	86.70	71.90	-	-	-	-	-
[Xiao and Shi(2019)]	84.20	74.10	65.00	55.40	35.60	85.50	71.10	-	-	-	-	-
[Aafaq et al.(2019)Aafaq, Akhtar, Liu, Gilani and Mian]	-	-	-	47.90	35.90	78.10	71.50	-	-	-	-	-
[Tang et al.(2019)Tang, Wang and Li]	82.80	71.70	62.40	52.40	35.70	84.30	72.20	-	-	-	-	-
[Zhang and Peng(2019b)]	-	-	-	56.90	36.20	90.60	-	-	-	-	-	-
[Xu et al.(2019)Xu, Liu, Nie and Su]	82.10	71.60	61.40	53.00	32.90	75.10	69.80	-	-	-	-	-
[Chen and Jiang(2019)]	-	-	-	53.40	35.00	86.70	-	-	-	-	-	-
[Pei et al.(2019)Pei, Zhang, Wang, Ke, Shen and Tai]	-	-	-	48.60	35.10	92.20	71.90	-	-	-	-	-
[Hou et al.(2019)Hou, Wu, Zhao, Luo and Jia]	-	-	-	52.80	36.10	87.80	71.80	-	-	-	-	-
[Hu et al.(2019)Hu, Chen, Zha and Wu]	86.80	75.00	65.10	54.70	35.20	91.30	72.50	-	-	-	-	-
[Gao et al.(2019)Gao, Li, Song and Shen]	83.30	73.60	64.60	54.30	33.50	72.80	-	-	-	-	-	-
[Guo et al.(2019)Guo, Zhang and Gao]	83.80	73.80	64.50	54.50	34.50	79.30	-	-	-	-	-	-
[Li and Gong(2019)]	-	-	-	50.30	34.10	87.50	70.80	-	-	-	-	-
[Li et al.(2019)Li, Zhou, Chen and Gao]	82.80	72.30	63.10	53.40	34.30	72.90	-	-	-	-	-	-
[Wang et al.(2019a)Wang, Ma, Zhang, Jiang, Wang and Liu]	-	-	-	42.00	28.20	48.70	61.60	-	-	-	-	-
[Xiao and Shi(2020)]	-	-	-	54.10	36.10	86.10	72.40	-	-	-	-	-
[Shi et al.(2020)Shi, Cai, Gu and Joty]	-	-	-	51.70	33.00	71.00	-	-	-	-	-	-
[Pan et al.(2020)Pan, Cai, Huang, Lee, Gaidon, Adeli and Niebles]	-	-	-	52.20	36.90	93.00	73.90	-	-	-	-	-
[Liu et al.(2020)Liu, Ren and Yuan]	-	-	-	55.70	35.50	88.80	72.60	-	-	-	-	-
[Zhang et al.(2019b)Zhang, Wang, Ma and Liu]	-	-	-	52.30	34.10	80.30	69.80	-	-	-	-	-
[Zhang et al.(2020)Zhang, Shi, Yuan, Li, Wang, Hu and Zha]	-	-	-	54.30	36.40	95.20	73.90	-	-	-	-	-
[Lei and Huang(2021)]	-	-	-	52.20	35.60	83.70	72.70	-	-	-	-	-
[Chen and Jiang(2021)]	-	-	-	55.80	36.90	74.50	98.50	-	-	-	-	-
[Perez-Martin et al.(2021)Perez-Martin, Bustos and Pérez]	-	-	-	64.40	41.90	111.50	79.50	-	-	-	-	-

Table 4: All results on MSVD dataset

Work	BLEU-1	BLEU-2	BLEU-3	BLEU-4	METEOR	CIDERr-D	ROUGE-L	Average-Recall (AR)	SPICE	FCE	RE	Self-BLEU
[Dong et al.(2016)Dong, Li, Lan, Huo and Snoek]	-	-	-	38.70	26.90	45.90	58.70	-	-	-	-	-
[Ramanishka et al.(2016)Ramanishka, Das, Park, Venugopalan, Hendricks, Rohrbach and Saenko]	-	-	-	40.70	28.60	46.50	61.00	-	-	-	-	-
[Shetty and Laaksonen(2016)]	-	-	-	41.11	27.70	46.40	59.60	-	-	-	-	-
[Jin et al.(2016)Jin, Chen, Chen, Xiong and Hauptmann]	-	-	-	42.60	28.80	61.70	46.70	-	-	-	-	-
[Nian et al.(2017)Nian, Li, Wang, Wu, Ni and Xu]	-	-	-	41.70	28.90	51.40	62.10	-	-	-	-	-
[Xu et al.(2017)Xu, Yao, Zhang and Mei]	-	-	-	36.50	26.50	41.00	59.80	-	-	-	-	-
[Pasumuru and Bansal(2017)]	-	-	-	40.80	28.80	47.10	60.20	-	-	-	-	-
[Zhang et al.(2017)Zhang, Gao, Zhang, Zhang, Li and Tian]	-	-	-	35.50	28.20	42.70	59.10	-	-	-	-	-
[Gao et al.(2017)Gao, Guo, Zhang, Xu and Shen]	-	-	-	38.00	26.10	43.20	-	-	-	-	-	-
[Chen et al.(2017)Chen, Chen, Jin and Hauptmann]	-	-	-	39.20	27.50	48.70	60.30	-	-	-	-	-
[Shen et al.(2017)Shen, Li, Su, Li, Chen, Jiang and Xue]	-	-	-	33.70	25.90	56.90	32.60	-	-	-	-	-
[Tu et al.(2017)Tu, Zhang, Liu and Yan]	-	-	-	37.40	26.60	41.50	-	-	-	-	-	-
[Song et al.(2017)Song, Guo, Gao, Liu, Zhang and Shen]	-	-	-	38.30	26.30	-	-	-	-	-	-	-
[Liu et al.(2017a)Liu, Xu, Wong, Li, Su and Kankanhalli]	-	-	-	37.10	26.70	41.00	59.00	-	-	-	-	-
[Phan et al.(2017)Phan, Henter, Miyao and Satoh]	82.8	69.5	56.20	44.10	29.10	49.70	62.40	-	-	-	-	-
[Hori et al.(2017)Hori, Hori, Lee, Zhang, Harsham, Hershey, Marks and Sumi]	-	-	-	39.70	25.50	40.00	-	-	-	-	-	-
[Song et al.(2018)Song, Guo, Gao, Li, Hanjalic and Shen]	-	-	-	39.80	26.10	40.90	59.30	-	-	-	-	-
[Wang et al.(2018c)Wang, Wang, Huang, Wang and Tan]	-	-	-	39.90	28.30	40.90	-	-	-	-	-	-
[Chen et al.(2018b)Chen, Wang, Zhang and Huang]	-	-	-	41.30	27.70	44.10	59.80	-	-	-	-	-
[Wang et al.(2018d)Wang, Wang, Huang, Wang and Tan]	73.6	59.3	48.26	38.13	26.58	-	-	-	-	-	-	-
[Long et al.(2018)Long, Gan and De Melo]	-	-	-	41.30	28.70	48.00	61.70	-	-	-	-	-
[Wang et al.(2018c)Wang, Chen, Wu, Wang and Wang]	-	-	-	41.30	28.70	48.00	61.70	-	-	-	-	-
[Yang et al.(2018)Yang, Zhou, Ai, Bin, Hanjalic, Shen and Ji]	-	-	-	36.00	26.10	-	-	-	-	-	-	-
[Chen et al.(2018a)Chen, Li, Zhang and Huang]	-	-	-	42.46	28.24	48.53	61.07	-	-	-	-	-
[Wang et al.(2018a)Wang, Ma, Zhang and Liu]	-	-	-	39.10	26.60	42.70	59.30	-	-	-	-	-
[Li et al.(2018a)Li, Guo and Fang]	76.1	62.1	49.10	37.50	26.40	-	-	-	-	-	-	-
[Wu and Han(2018)]	-	-	-	38.10	27.20	42.10	-	-	-	-	-	-
[Bin et al.(2018)Bin, Yang, Shen, Xie, Shen and Li]	78.9	60.4	46.10	33.90	26.20	-	-	-	-	-	-	-
[Pu et al.(2018)Pu, Min, Gan and Carin]	-	-	-	45.01	29.98	51.41	-	-	-	-	-	-
[Yan et al.(2019)Yan, Tu, Wang, Zhang, Hao, Zhang and Dai]	-	-	-	39.30	27.10	43.80	-	-	-	-	-	-
[Wu et al.(2019)Wu, Han, Yang, Hu and Wu]	-	-	-	38.10	27.20	42.10	-	-	-	-	-	-
[Xu et al.(2018a)Xu, Liu, Wong, Zhang, Nie, Su and Kankanhalli]	-	-	-	42.30	29.40	46.10	62.30	-	-	-	-	-
[Shi et al.(2019)Shi, Cai, Joty and Gu]	-	-	-	43.20	28.00	48.30	62.00	-	-	-	-	-
[Qi et al.(2019)Qi, Wang, Li and Luo]	-	-	-	36.70	25.90	33.90	-	-	-	-	-	-
[Aafaq et al.(2019)Aafaq, Akhtar, Liu, Gilani and Mian]	-	-	-	38.30	28.40	48.10	60.70	-	-	-	-	-
[Tang et al.(2019)Tang, Wang and Li]	81.1	67.2	53.70	41.40	29.00	48.90	61.30	-	-	-	-	-
[Zhang and Peng(2019b)]	-	-	-	41.40	28.20	46.90	-	-	-	-	-	-
[Xu et al.(2019)Xu, Liu, Nie and Su]	-	-	-	40.80	27.50	45.40	60.70	-	-	-	-	-
[Chen and Jiang(2019)]	-	-	-	45.40	28.60	50.10	-	-	-	-	-	-
[Pei et al.(2019)Pei, Zhang, Wang, Ke, Shen and Tai]	-	-	-	40.40	28.10	47.10	60.70	-	-	-	-	-
[Hou et al.(2019)Hou, Wu, Zhao, Luo and Jia]	-	-	-	42.30	29.70	49.10	62.80	-	-	-	-	-
[Chen et al.(2019)Chen, Jin, Chen and Hauptmann]	-	-	-	44.91	29.61	51.80	62.81	-	6.85	-	-	-
[Gao et al.(2019)Gao, Li, Song and Shen]	83.3	73.6	64.60	54.30	33.50	72.80	-	-	-	-	-	-
[Guo et al.(2019)Guo, Zhang and Gao]	77.6	64.0	51.30	39.90	27.10	43.80	-	-	-	-	-	-
[Li and Gong(2019)]	-	-	-	40.40	27.90	48.30	61.00	-	-	-	-	-
[Lee(2019)]	-	-	-	36.60	23.80	27.10	52.40	-	-	-	-	-
[Zhao et al.(2019)Zhao, Li and Liu]	-	-	-	36.20	27.90	38.80	58.80	-	-	-	-	-
[Li et al.(2019)Li, Zhou, Chen and Gao]	77.1	62.1	48.70	37.00	26.90	40.70	-	-	-	-	-	-
[Wang et al.(2019a)Wang, Ma, Zhang, Jiang, Wang and Liu]	-	-	-	41.70	27.90	48.40	61.00	-	-	-	-	-
[Xiao and Shi(2020)]	-	-	-	44.60	28.70	48.60	62.20	-	-	-	-	-
[Shi et al.(2020)Shi, Cai, Gu and Joty]	-	-	-	41.80	26.70	45.30	-	-	-	-	-	-
[Pan et al.(2020)Pan, Cai, Huang, Lee, Gaidon, Adeli and Niebles]	-	-	-	40.50	28.30	47.10	60.90	-	-	-	-	-
[Liu et al.(2020)Liu, Ren and Yuan]	-	-	-	41.20	27.80	48.60	60.80	-	-	-	-	-
[Zhang et al.(2019b)Zhang, Wang, Ma and Liu]	-	-	-	39.20	27.50	48.70	60.30	-	-	-	-	-
[Zhang et al.(2020)Zhang, Shi, Yuan, Li, Wang, Hu and Zha]	-	-	-	43.60	28.80	50.90	62.10	-	-	-	-	-
[Wei et al.(2020)Wei, Mi, Hu and Chen]	-	-	-	38.50	26.90	43.70	-	-	-	-	-	-
[Yuan et al.(2020)Yuan, Ma, Wang and Zhu]	-	-	-	5.01	16.25	12.22	30.41	-	-	-	-	-
[Lei and Huang(2021)]	-	-	-	41.30	28.20	48.60	61.90	-	-	-	-	-
[Chen and Jiang(2021)]	-	-	-	41.70	28.90	51.40	62.10	-	-	-	-	-
[Perez-Martin et al.(2021)Perez-Martin, Bustos and Pérez]	-	-	-	46.40	30.40	51.90	64.70	-	-	-	-	-
[Ji and Wang(2021)]	-	-	-	39.30	27.10	-	59.50	-	-	-	-	-

Table 5: All results on MSR-VTT dataset

Work	BLEU-1	BLEU-2	BLEU-3	BLEU-4	METEOR	CIDERr-D	ROUGE-L	Average-Recall (AR)	SPICE	FCE	RE	Self-BLEU
[Inácio et al.(2021)Inácio, Gutoski, Lazzaretti and Lopes]	-	-	-	4.320	9.98	30.50	21.33	-	-	-	-	-
[Yu and Han(2021)]	-	-	-	-	5.82	10.87	-	-	-	-	-	-
[Iashin and Rahtu(2020)]	-	-	5.83	2.860	11.72	-	-	-	-	-	-	-
[Yuan et al.(2020)Yuan, Ma, Wang and Zhu]	-	-	-	0.530	5.11	9.31	11.11	-	-	-	-	-
[Suin and Rajagopalan(2020)]	-	-	2.87	1.350	6.21	13.82	-	53.40	-	-	-	-
[Wang et al.(2018b)Wang, Jiang, Ma, Liu and Xu]	18.99	8.84	4.41	2.300	9.60	12.68	19.10	-	-	-	-	-
[Mun et al.(2019)Mun, Yang, Ren, Xu and Han]	28.02	12.05	4.41	1.280	13.07	43.48	-	-	-	-	-	-
[Zhou et al.(2018b)Zhou, Zhou, Corso, Socher and Xiong]	-	-	4.76	2.230	9.56	-	-	52.95	-	-	-	-
[Qi et al.(2019)Qi, Wang, Li and Luo]	26.60	13.90	8.20	4.900	9.90	24.60	-	-	-	-	-	-
[Zhang et al.(2019c)Zhang, Xu, Ouyang and Tan]	22.76	10.12	4.26	1.640	10.71	31.41	22.85	-	-	-	-	-
[Hou et al.(2019)Hou, Wu, Zhao, Luo and Jia]	-	-	-	1.900	11.30	44.20	22.40	-	-	-	-	-
[Li et al.(2018b)Li, Yao, Pan, Chao and Mei]	19.57	9.90	4.55	1.622	10.33	25.24	-	-	-	-	-	-
[Xiong et al.(2018)Xiong, Dai and Lin]	39.11	22.26	13.52	8.450	14.75	14.15	14.75	-	-	-	17.59	45.8
[Krishna et al.(2017)Krishna, Hata, Ren, Fei-Fei and Carlos Niebles]	26.45	13.48	7.12	3.980	9.46	24.56	-	-	-	-	-	-
[Zhang et al.(2019b)Zhang, Wang, Ma and Liu]	-	-	-	1.740	10.47	38.43	23.49	-	-	-	-	-

Table 6: All results on ActivityNet Captions dataset

Work	BLEU-1	BLEU-2	BLEU-3	BLEU-4	METEOR	CIDERr-D	ROUGE-L	Average-Recall (AR)	SPICE	FCE	RE	Self-BLEU
[Shin et al.(2016)Shin, Ohnishi and Harada]	-	-	-	0.40	4.70	8.90	-	-	-	-	-	-
[Pan et al.(2016a)Pan, Xu, Yang, Wu and Zhuang]	-	-	-	-	6.80	-	-	-	-	-	-	-
[Pan et al.(2016b)Pan, Mei, Yao, Li and Rui]	-	-	-	-	6.70	-	-	-	-	-	-	-
[Venugopalan et al.(2016)Venugopalan, Hendricks, Mooney and Saenko]	-	-	-	-	6.70	-	-	-	-	-	-	-
[Zhang and Tian(2016)]	-	-	-	-	6.70	-	-	-	-	-	-	-
[Yang et al.(2017)Yang, Han and Wang]	-	-	-	-	6.90	-	-	-	-	-	-	-
[Baraldi et al.(2017)Baraldi, Grana and Cucchiara]	-	-	-	-	7.30	-	-	-	-	-	-	-
[Nian et al.(2017)Nian, Li, Wang, Wu, Ni and Xu]	-	-	-	-	5.70	-	-	-	-	-	-	-
[Pasumuru and Bansal(2017)]	-	-	-	-	7.40	-	-	-	-	-	-	-
[Liu et al.(2017b)Liu, Li and Shi]	-	-	-	-	6.90	-	-	-	-	-	-	-
[Pan et al.(2017)Pan, Yao, Li and Mei]	-	-	-	-	7.40	-	-	-	-	-	-	-
[Yang et al.(2018)Yang, Zhou, Ai, Bin, Hanjalic, Shen and Ji]	-	-	-	-	6.30	-	-	-	-	-	-	-
[Xu et al.(2018b)Xu, Han, Hong and Tian]	-	-	-	-	7.20	-	-	-	-	-	-	-
[Pu et al.(2018)Pu, Min, Gan and Carin]	-	-	-	2.08	7.12	9.14	-	-	-	-	-	-

Table 7: All results on MVAD dataset

Work	BLEU-1	BLEU-2	BLEU-3	BLEU-4	METEOR	CIDERr-D	ROUGE-L	Average-Recall (AR)	SPICE	FCE	RE	Self-BLEU
[Shin et al.(2016)Shin, Ohnishi and Harada]	-	-	-	13.0	4.8	7.5	-	-	-	-	-	-
[Baraldi et al.(2017)Baraldi, Grana and Cucchiara]	-	-	-	1.0	7.0	10.8	16.7	-	-	-	-	-
[Nian et al.(2017)Nian, Li, Wang, Wu, Ni and Xu]	-	-	-	-	6.6	-	-	-	-	-	-	-
[Pan et al.(2017)Pan, Yao, Li and Mei]	-	-	-	40.8	28.8	47.1	60.2	-	-	-	-	-
[Xu et al.(2018a)Xu, Liu, Wong, Zhang, Nie, Su and Kankanhalli]	-	-	-	1.9	7.9	-	-	-	-	-	-	-
[Xiao and Shi(2019)]	-	-	-	0.9	6.8	10.0	-	-	-	-	-	-
[Xu et al.(2019)Xu, Liu, Nie and Su]	-	-	-	0.8	7.7	-	-	-	-	-	-	-
[Zhao et al.(2019)Zhao, Li and Lu]	-	-	-	-	7.8	-	-	-	-	-	-	-
[Yang et al.(2018)Yang, Zhou, Ai, Bin, Hanjalic, Shen and Ji]	-	-	-	-	7.2	-	-	-	-	-	-	-
[Liu et al.(2017a)Liu, Xu, Wong, Li, Su and Kankanhalli]	16.9	5.4	1.6	0.6	7.1	8.9	17.0	-	-	-	-	-
[Pan et al.(2016b)Pan, Mei, Yao, Li and Rui]	-	-	-	-	7.3	-	-	-	-	-	-	-
[Venugopalan et al.(2016)Venugopalan, Hendricks, Mooney and Saenko]	-	-	-	-	6.8	-	-	-	-	-	-	-
[Zhang and Tian(2016)]	-	-	-	-	7.0	-	-	-	-	-	-	-

Table 8: All results on MPII dataset

Work	BLEU-1	BLEU-2	BLEU-3	BLEU-4	METEOR	CIDERr-D	ROUGE-L	Average-Recall (AR)	SPICE	FCE	RE	Self-BLEU
[Nian et al.(2017)Nian, Li, Wang, Wu, Ni and Xu]	-	-	-	40.10	29.90	51.10	-	-	-	-	-	-
[Li et al.(2017)Li, Zhao, Lu et al.]	80.1	66.1	54.7	41.30	32.20	53.90	68.80	-	-	-	-	-
[Chen et al.(2017)Chen, Chen, Jin and Hauptmann]	-	-	-	47.56	34.21	79.57	70.45	-	-	-	-	-
[Hori et al.(2017)Hori, Hori, Lee, Zhang, Harsham, Hershey, Marks and Sumi]	-	-	-	56.80	34.30	72.40	-	-	-	-	-	-
[Zhao et al.(2018)Zhao, Li, Lu et al.]	77.6	67.1	55.4	43.80	32.60	52.20	69.30	-	-	-	-	-
[Chen et al.(2018a)Chen, Li, Zhang and Huang]	-	-	-	53.21	35.23	86.76	-	-	-	-	-	-
[Zolfaghari et al.(2018)Zolfaghari, Singh and Brox]	-	-	62.6	53.50	35.00	85.80	-	-	-	-	-	-
[Qi et al.(2019)Qi, Wang, Li and Luo]	81.2	69.7	61.3	50.90	33.50	70.30	-	-	-	-	-	-
[Chen et al.(2019)Chen, Jin, Chen and Hauptmann]	-	-	-	49.26	33.91	83.20	70.95	-	5.44	-	-	-
[Zhao et al.(2019)Zhao, Li and Lu]	80.3	67.6	56.0	42.40	33.40	54.30	69.40	-	-	-	-	-
[Zheng et al.(2020)Zheng, Wang and Tao]	-	-	-	46.50	33.50	81.00	69.40	-	-	-	-	-
[Wei et al.(2020)Wei, Mi, Hu and Chen]	-	-	-	46.80	34.40	85.70	-	-	-	-	-	-

Table 9: All results on Youtube2Text dataset

Work	BLEU-1	BLEU-2	BLEU-3	BLEU-4	METEOR	CIDERr-D	ROUGE-L	Average-Recall (AR)	SPICE	FCE	RE	Self-BLEU
[Fakoor et al.(2016)Fakoor, Mohamed, Mitchell, Kang and Kohli]	50.0	31.1	18.8	11.5	17.6	16.7	-	-	-	-	-	-
[Li et al.(2017)Li, Zhao, Lu et al.]	50.6	31.7	21.3	12.7	19.1	18.3	-	-	-	-	-	-
[Wu et al.(2018)Wu, Li, Cao, Ji and Lin]	-	-	-	13.5	17.8	20.8	-	-	-	-	-	-
[Zhao et al.(2018)Zhao, Li, Lu et al.]	50.7	31.3	19.7	13.3	19.0	18.0	-	-	-	-	-	-
[Wang et al.(2018e)Wang, Chen, Wu, Wang and Wang]	64.4	44.3	29.4	18.8	18.7	23.6	31.2	-	-	-	-	-
[Hu et al.(2019)Hu, Chen, Zha and Wu]	-	-	-	14.5	18.4	23.7	-	-	-	-	-	-
[Zhao et al.(2019)Zhao, Li and Lu]	51.3	32.1	21.7	12.9	19.7	18.8	-	-	-	-	-	-
[Wei et al.(2020)Wei, Mi, Hu and Chen]	-	-	-	12.7	17.2	21.6	-	-	-	-	-	-

Table 10: All results on Charades dataset

Work	BLEU-1	BLEU-2	BLEU-3	BLEU-4	METEOR	CIDERr-D	ROUGE-L	Average-Recall (AR)	SPICE	FCE	RE	Self-BLEU
[Yu et al.(2016)Yu, Wang, Huang, Yang and Xu]	60.8	49.6	38.5	30.5	28.7	1.62	-	-	-	-	-	-

Table 11: All results on TACoS-MultiLevel dataset

Work	BLEU-1	BLEU-2	BLEU-3	BLEU-4	METEOR	CIDERr-D	ROUGE-L	Average-Recall (AR)	SPICE	FCE	RE	Self-BLEU
[Yu et al.(2017b)Yu, Ko, Choi and Kim]	13.5	4.4	1.7	0.8	7.1	10.0	15.9	-	-	-	-	-
[Yu et al.(2017a)Yu, Choi, Kim, Yoo, Lee and Kim]	16.8	5.5	2.1	-	7.2	9.3	15.6	-	-	-	-	-
[Gao et al.(2019)Gao, Li, Song and Shen]	-	-	-	0.7	5.6	10.4	14.6	-	-	-	-	-

Table 12: All results on LSMDC dataset

Work	BLEU-1	BLEU-2	BLEU-3	BLEU-4	METEOR	CIDERr-D	ROUGE-L	Average-Recall (AR)	SPICE	FCE	RE	Self-BLEU
[Yu et al.(2017a)Yu, Choi, Kim, Yoo, Lee and Kim]	30.6	12.5	4.9	-	8.4	8.4	22.9	-	-	-	-	-

Table 13: All results on VAS dataset

Work	BLEU-1	BLEU-2	BLEU-3	BLEU-4	METEOR	CIDeRr-D	ROUGE-L	Average-Recall (AR)	SPICE	FCE	RE	Self-BLEU
[Zhou et al.(2018b)Zhou, Zhou, Corso, Socher and Xiong]	-	-	-	1.42	11.20	-	-	-	-	-	-	-
[Yu and Han(2021)]	-	-	-	-	2.43	4.88	-	-	-	-	-	-

Table 14: All results on YouCook2 dataset

Work	BLEU-1	BLEU-2	BLEU-3	BLEU-4	METEOR	CIDeRr-D	ROUGE-L	Average-Recall (AR)	SPICE	FCE	RE	Self-BLEU
yu2018fine	55.88	46.12	39.21	34.45	27.57	2.61	53.5	-	39.1	19.44	-	-

Table 15: All results on FSN dataset

Work	BLEU-1	BLEU-2	BLEU-3	BLEU-4	METEOR	CIDeRr-D	ROUGE-L	Average-Recall (AR)	SPICE	FCE	RE	Self-BLEU
[Zhang et al.(2020)Zhang, Shi, Yuan, Li, Wang, Hu and Zha]	-	-	-	32.1	22.2	49.7	48.9	-	-	-	-	-
[Chen and Jiang(2021)]	-	-	-	34.2	23.5	57.6	50.3	-	-	-	-	-

Table 16: All results on VATEX dataset

Work	BLEU-1	BLEU-2	BLEU-3	BLEU-4	METEOR	CIDeRr-D	ROUGE-L	Average-Recall (AR)	SPICE	FCE	RE	Self-BLEU
[Qi et al.(2019)Qi, Wang, Li and Luo]	-	-	-	31.76	26.07	2.91	51.62	-	-	-	-	-

Table 17: All results on SVCDV dataset

Work	BLEU-1	BLEU-2	BLEU-3	BLEU-4	METEOR	CIDeRr-D	ROUGE-L	Average-Recall (AR)	SPICE	FCE	RE	Self-BLEU
[Zhu et al.(2020)Zhu, Hwang, Ma, Chen and Guo]	49.2	35.5	27.2	21.2	21	54.3	47.5	-	-	-	-	-

Table 18: All results on Object-oriented captions dataset

Work	BLEU-1	BLEU-2	BLEU-3	BLEU-4	METEOR	CIDeRr-D	ROUGE-L	Average-Recall (AR)	SPICE	FCE	RE	Self-BLEU
[Inácio et al.(2021)Inácio, Gutoski, Lazzaretti and Lopes]	-	-	-	69.54	49.34	354.04	84.05	-	-	-	-	-

Table 19: All results on LIRIS dataset

Work	BLEU-1	BLEU-2	BLEU-3	BLEU-4	METEOR	CIDeRr-D	ROUGE-L	Average-Recall (AR)	SPICE	FCE	RE	Self-BLEU
[Dilawari et al.(2021)Dilawari, Khan, Al-Otaibi, Rehman, Rahman and Nam]	-	-	-	-	33.9	-	0.72	-	-	-	-	-

Table 20: All results on TRECVID dataset