

Surveying the Landscape of Image Captioning Evaluation: A Comprehensive Taxonomy, Trends, and Metrics Analysis

Uri Berger^{◊†} and Gabriel Stanovsky[◊] and Omri Abend[◊] and Lea Frermann[†]

[◊]School of Computer Science and Engineering, The Hebrew University of Jerusalem, Israel

{uri.berger2, gabriel.stanovsky, omri.abend}@mail.huji.ac.il

[†]School of Computing and Information Systems, University of Melbourne, Australia

lea.frermann@unimelb.edu.au

Abstract

The task of image captioning has recently been gaining popularity, and with it the complex task of evaluating the quality of image captioning models. In this work, we present the first survey and taxonomy of over 70 different image captioning metrics and their usage in hundreds of papers, specifically designed to help users select the most suitable metric for their needs. We find that despite the diversity of proposed metrics, the vast majority of studies rely on only five popular metrics, which we show to be weakly correlated with human ratings. We hypothesize that combining a diverse set of metrics can enhance correlation with human ratings. As an initial step, we demonstrate that a linear regression-based ensemble method, which we call ENSEMB-EVAL, trained on one human ratings dataset, achieves improved correlation across five additional datasets, showing there is a lot of room for improvement by leveraging a diverse set of metrics.¹

1 Introduction

Evaluating the output of image captioning is challenging for several reasons. First, as in other generative tasks, multiple outputs may be valid for the same input, as illustrated in Figure 1, where different valid captions for the same image have no overlap in content words. Second, it involves bridging across modalities, requiring evaluators to compare text and images, unlike most generative tasks that only involve textual information. Attending to the difficulty and importance of image captioning is the sheer volume of metrics proposed for this task, and their usage in hundreds of image captioning models (Figure 2).

In this work, we survey the different approaches proposed for evaluating image captioning, and

provide a first taxonomy comprising over 70 metrics, which were largely developed independently of one another. We examine usage patterns of metrics over the past 14 years and demonstrate that combining multiple metrics enhances correlation with human ratings.

We begin by exploring various definitions of image captioning presented in previous studies and examining their impact on evaluation (Section 2). Then, in Section 3, we systematically examine all relevant papers from 15 major venues in NLP, vision and machine learning between 2010 and 2024. This approach yields a body of work consisting of 71 different automatic evaluation metrics and 5 human evaluation paradigms, used in over 300 papers.

Next, we organize both automatic metrics and human evaluation in a principled taxonomy (Section 4). Our taxonomy is the first to comprehensively cover all automatic metrics used in image captioning research. We design the taxonomy dimensions with the needs of metric users in mind, focusing on dimensions such as the candidate’s property each metric aims to measure. We also propose the first taxonomy for human evaluation metrics, categorized into groups such as comparative evaluation and scale rating evaluation.

We then survey the use of evaluation methods across 314 papers (Section 5). Interestingly, we find that despite the wealth of different metrics proposed for the task, the vast majority of examined papers use only five simple metrics (BLEU, METEOR, ROUGE, CIDEr, SPICE), although there has been a recent increase in the adoption of alternative metrics. Furthermore, we find that the use of human evaluation is declining in recent years and that significance and inter-annotator agreement are rarely reported.

Finally, in a series of evaluations, we show that the five most popular metrics show only a

¹Our code and data are available at github.com/uriberger/caption_evaluation.



Lego figures in the middle of the desert
 A pyramid with statues of ancient kings
 Some people and a dog in front of an old monument

Figure 1: An image with various valid captions that have no overlap in content words, exemplifying some of the challenges in evaluating image captioning (taken from Crossmodal-3600; Thapliyal et al., 2022).

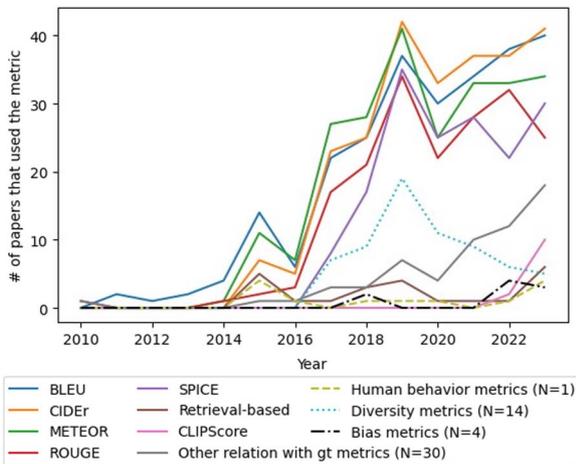


Figure 2: Metrics usage over the years. Metrics used in ≤ 1 papers are omitted. $N =$ indicates number of metrics in this category.

weak correlation with human ratings, whereas lesser-known metrics show a much higher correlation. This underscores the value of a systematic and comprehensive survey in identifying these lesser-known but more effective metrics. We then show that an ensemble of selected metrics we name ENSEMBEVAL, optimized for diversity using a feature selection algorithm, achieves enhanced correlation with human ratings. This paves the way for future research on metrics that consider multiple aspects simultaneously.

To recap, our contributions are threefold. First, we introduce the first comprehensive taxonomy of automatic captioning evaluation metrics categorized by user-oriented aspects, helping future

model developers select the most suitable metric based on the specific property they want their model to excel in. Second, we examine key trends in image captioning evaluation, highlighting a shift from reliance on five simple metrics toward alternative approaches. Finally, we demonstrate that a simple linear regression-based ensemble method improves correlation with human ratings and provide the code publicly for future use.

2 Task Definition

Image captioning can be motivated by diverse goals, and what counts as a ‘good’ caption heavily depends on the system’s purpose and context. Accurate evaluation requires explicitly defining a system’s purpose, a step often overlooked or only implicitly addressed in previous work. We start with a review of common definitions.

Early research framed image captioning as a means to assess vision-and-language models based on image-caption semantics, rather than practical applications (Hodosh et al., 2013). More recently, the task has often been defined simply as describing the visual content of an image or scene in natural language (Stefanini et al., 2022; Nallapaneni and Konakanchi, 2023; González-Chávez et al., 2024). Yet, this definition does not specify what visual content should be described or the intended nature of the description. These aspects depend on both the *purpose* and the *context* of the description, which we now discuss in detail.

Descriptions serve different **purposes**, each with distinct requirements. For example, captions designed for accessibility should *replace* the image, while captions in news articles typically *augment* or *link* the content of the image and article (Kreiss et al., 2022). In practice, most studies focus on image *replacement*, i.e., generating descriptions intended to serve as substitutes for the image, as defined by Kreiss et al. (2022).

The exact purpose of image captioning is often implicit in image captioning papers. Hodosh et al. (2013) discuss the task in detail, and adopt a number of fundamental theoretical distinctions as to the purposes of captioning, taken from library science (Shatford, 1986; Jaimes and Chang, 1999). Specifically, they argue that captions that focus on conceptual descriptions (rather than low-level ‘perceptual’ descriptions) are the most relevant for image understanding, the most commonly

Community	Conferences
NLP	AAACL, ACL, CoNLL, EACL EMNLP, NAACL, TACL, *SEM
CV	CVPR, ICCV, ECCV
ML	Neurips, ICML, ICLR, AAAI

Table 1: Venues included in our review. NLP: Natural Language Processing; CV: Computer Vision; ML: Machine Learning.

tackled purpose in NLP. Within this category, an important sub-distinction is between specific captions, that identify entities by their names (e.g., Great Pyramid of Giza), and generic ones (e.g., that instead identify an object as a pyramid). Image understanding usually (albeit, not universally) targets the latter.

The second key factor is the **context** in which a caption is generated. For instance, the top caption in Figure 1 may suit Lego’s Instagram account, while the middle caption might be more appropriate for a travel guide.

While our survey covers all works on captioning evaluation—irrespective of purpose or context—our experiments on metric performance and ensembling (Section 6) focus on the *replacement* task, specifically.

3 Literature Survey

We perform a systematic review and gather usage statistics of metrics in conference papers on English image captioning. Overall, we examined 314 papers, and identified 71 distinct metrics.

Scope. We focus on a large set of 15 venues spanning three distinct research communities (Table 1). Figure 3 shows the annual number of papers for each community. We begin collecting data starting from 2024 backwards until reaching the first year with no relevant papers, which is 2009.

Search Strategy. To identify image captioning related papers within a venue, we search the venue proceedings for papers with the substrings *caption* or *description* in their titles, assuming relevant titles are likely to include variations of the phrases *image captioning* or *image descriptions*.

Filtering. Next, we manually filter out papers that we deem irrelevant to our review. Specifically, we exclude: 1) papers on other types of

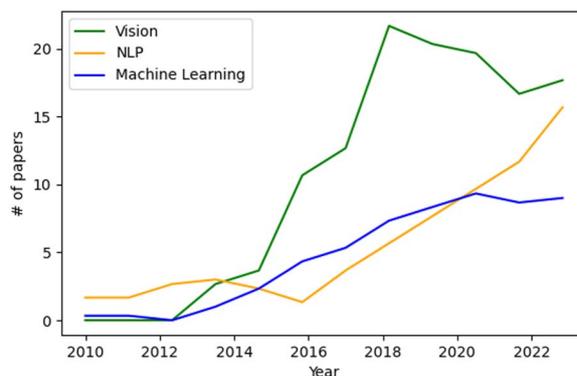


Figure 3: Annual number of image captioning papers per community, smoothed by convolving with a window size of 3 years.

captioning such as video captioning; 2) papers unrelated to captioning that happen to have *caption* or *description* in their titles; 3) papers that focus on languages other than English; and 4) papers that did not conduct any experiments (e.g., review papers).

Collected Data. We collect two types of information. First, we intend to describe the types of metrics used in previous work. For each metric introduced or used in the papers we examined, we document its name, implementation details and the paper in which it was introduced. To the best of our knowledge, we are the first to systematically collect such a comprehensive set of metrics, as previous research has focused on a few well-known metrics (see Section 8).

Second, we seek to analyze the patterns of metrics usage. Therefore, we record which metrics were used in each examined paper.²

4 Taxonomy

Here, we introduce a taxonomy of the metrics from the systematic review outlined in Section 3.

This work is novel in two respects: We are the first to comprehensively include all metrics proposed and used in previous studies, rather than focusing on a limited set of well-known metrics. Second, we are the first to describe a taxonomy for human evaluation methods.

4.1 Automatic Evaluation

Motivation. Given the large number of 71 automatic evaluation methods identified in our survey,

²We only document metrics that were used in the main body of the paper, omitting metrics used in the Appendix.

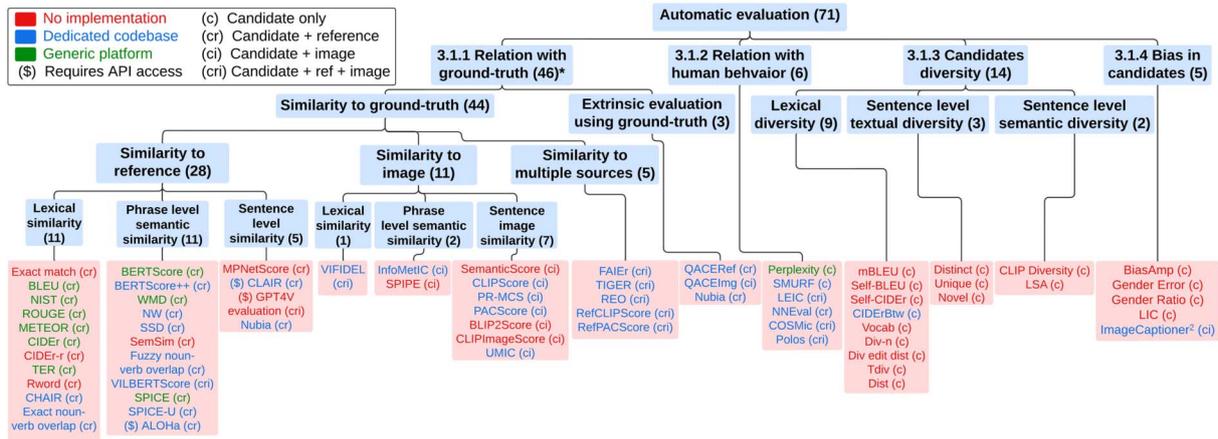


Figure 4: A taxonomy of automatic metrics for image captioning evaluation. For each category, we list the number of metrics in parentheses. Top categories are preceded by the relevant section number. Blue rectangles contain categories, red contain metrics. * – Nubia is assigned to two distinct categories under *Relation with ground-truth*.

it is a challenge to stay updated on the latest metrics. As a result, the common practice is to rely on previously used metrics, typically the five dominant ones (BLEU, CIDEr, METEOR, ROUGE, and SPICE) as shown in Figure 2.

To address this issue, we propose a taxonomy that helps users select the most suitable metric. We categorize metrics based on the specific property they assess in the candidate caption.

Method. We begin by categorizing the metrics based on the evaluated property. For example, BLEU, BERTScore, and CLIPScore are classified under *lexical similarity to reference*, *phrase level semantic similarity to reference*, and *semantic similarity to image*, respectively.³

Next, we incrementally group together pairs of categories based on their similarity to one another. For instance, *lexical similarity to reference* and *phrase level semantic similarity to reference* are both grouped under *similarity to reference*, while *semantic similarity to image* falls under *similarity to image*. We repeat this process iteratively until only one overarching category remains.

Finally, following the process of Nickerson et al. (2013), we extend the taxonomy with three additional dimensions, again motivated by the overarching goal to support users in their choice of metrics:

- **Input type:** The data provided to the metric (candidate only, candidate + reference, candidate + image, candidate + reference + image).

candidate + image, candidate + reference + image).

- **Implementation availability:** Availability of the metric’s implementation (none, dedicated codebase, or integration within a generic platform like HuggingFace).
- **API access requirements:** Whether the metric relies on a closed model API access.

The initial taxonomy was constructed by the first author, followed by validation from the remaining authors. Any disagreements were resolved through majority voting. The final taxonomy is visually presented in Figure 4. In the following sections, we describe the various automatic evaluation methods.

4.1.1 Relation with Ground-Truth

The largest and most widely used class of automatic methods use a ground-truth source (either a reference caption or the image) as a basis for evaluating the candidate. We categorize these methods based on the relation they examine (similarity/extrinsic) and the nature of the ground-truth (reference, image, both) they use.

Similarity to Reference. Early automatic methods, still the most common today, compare the candidate caption to a reference caption.

Lexical Similarity: Some metrics compare the candidate and reference captions at the word level using a straightforward textual comparison, most naïvely using exact match (Kang et al., 2023). More advanced methods include computing n -gram overlap (*BLEU*: Papineni et al., 2002;

³One metric (*Yngve score*, [Liu et al., 2019a]) did not align with any category, another (*Nubia*, [Kane et al., 2020]) fit into two categories.

NIST: Doddington, 2002; *ROUGE*: Lin, 2004; *METEOR*: Banerjee and Lavie, 2005), n -gram TF-IDF (*CIDEr*: Vedantam et al., 2015; *CIDEr-r*: Oliveira dos Santos et al., 2021), and calculating minimal edits to match the reference (*TER*, Snover et al., 2006). Other metrics include the number of reference words in the candidate (*Rword*, Cho et al., 2022) and semantic comparisons including synonyms (*CHAIR*, Rohrbach et al., 2018). Finally, some studies measure the precision and recall of specific word categories, including parts of speech (e.g., *Exact noun/verb overlap* Chan et al., 2023a), objects (Wang et al., 2021c) and named entities (common in newsimage captioning, e.g., Zhang and Wan, 2023).

Phrase-level Semantic Similarity: Others compare the semantics of sentence elements, most commonly using phrase embeddings, either with context (*BERTScore*: Zhang et al., 2020a; *BERTScore++*: Yi et al., 2020) or without (*WMD*: Kusner et al., 2015; *NW*: Cornia et al., 2019; *SSD*: Takmaz et al., 2020; *SemSim*: Nag Chowdhury et al., 2021; *Fuzzy noun/verb overlap*: Chan et al., 2023a). *VILBERTScore* (Lee et al., 2020) enriches phrase embeddings by injecting visual information. The *SPICE* family (*SPICE*: Anderson et al., 2016; *SPICE-U*: Wang et al., 2020c) compares the components of scene graphs of the candidate and reference captions. *ALOHA* (Petryk et al., 2024) prompts Large Language Models (LLMs) to identify object phrases in both candidate and references and computes the similarity of these phrases' embeddings.

Sentence-level Similarity: Some metrics compare semantics at the sentence level. *MPNetScore* (Black et al., 2024) compares sentence level embeddings. Others (*CLAIR*: Chan et al., 2023b; *GPT4V evaluation*: Ge et al., 2024) prompt LLMs to compare the candidate and reference caption. *Nubia* (Kane et al., 2020) explicitly trains a model to evaluate sentence similarity.

Extracted Information Similarity: In specific domains (e.g., medical image captioning), it is common to extract relevant information from both the reference and candidate captions for comparison (e.g., a list of diseases from the captions in clinical image captioning; Nishino et al., 2022).

Similarity to Image. Recent work suggests comparing the candidate caption to the image

instead of a reference caption. These metrics are also known as reference-free metrics.

Lexical Similarity: *VIFIDEL* (Madhyastha et al., 2019) compares object phrases in the candidate with names of objects identified in the image.

Phrase-level Semantic Similarity: *InfoMetIC* (Hu et al., 2023a) compares text and visual token embeddings. *SPIPE* (Xie et al., 2022) is a variation of *SPICE* where the scene graphs of the candidate and image are compared.

Sentence-image Similarity: Several recent works embedded sentences and images in a joint visual-textual space and computed the similarity of candidate-image embeddings (*Semantic Score*: Dognin et al., 2019; *CLIPScore*: Hessel et al., 2021; *PR-MCS*: Kim et al., 2023; *PACScore*: Sarto et al., 2023; *BLIP2Score*: Zeng et al., 2024). *CLIPImageScore* (Ge et al., 2024) use text-to-image methods to generate an image based on the candidate and computes the embedding similarity between this generated image and the original image. *UMIC* (Lee et al., 2021b) explicitly trains models to predict similarity of candidate caption and image.

Retrieval-based Methods: A popular evaluation protocol involves using text-to-image retrieval on the test set (e.g., Kornblith et al., 2023). In this setting, a candidate caption is matched with each image in the test set using an image-text matching model (most commonly CLIP), and the images are ranked by their matching score with the candidate. Then, recall@ n is applied as the evaluation metric, with typical values for n being 1, 5, and 10.

Similarity to Multiple Sources. Some methods use both the image and reference captions as sources for comparison. *FAIEr* (Wang et al., 2021b) compares the candidate scene graph with a fusion of the image and reference scene graphs. *TIGER* (Jiang et al., 2019b) and *REO* (Jiang et al., 2019a) compare the candidate-image joint embedding vector with the reference-image joint embedding vector. *RefCLIPScore* (Hessel et al., 2021) and *RefPACScore* (Sarto et al., 2023) compute the harmonic mean of candidate-image similarity with candidate-reference similarity.

Extrinsic Evaluation Using Ground-truth. Several NLP tasks focus on understanding the

relation between two input sources. Some metrics use models trained for this task with the candidate and a ground-truth source as the inputs.

Question Generation and Question Answering. Lee et al. (2021a) generate questions from the candidate and score it by how well a question-answering model answers using the reference caption (*QACERef*) or the image (*QACEImg*).

Natural Language Inference (NLI). Nubia (Kane et al., 2020) uses NLI models with the reference as the premise and the candidate as the hypothesis.

4.1.2 Relation With Human Behavior

The primary aim of captioning systems developed in recent years is to imitate human behavior, both in the text they generate and in the evaluation process. Metrics in this category strive to explicitly assess these properties.

Candidate Fluency. Some works (e.g., Ou et al., 2023) evaluate the candidate’s *perplexity* by employing an off-the-shelf large language model (commonly GPT-2, Radford et al., 2019) to compute the probability of the candidate’s tokens. *SMURF* (Feinglass and Yang, 2021) computes the activation flow through a Transformer model self-attention layers, assumed to increase when the candidate differs greatly from typical captions.

Is the Candidate Human-like? *LEIC* (Cui et al., 2018) and *NNEval* (Sharif et al., 2018) measure how likely the candidate is to deceive a model that discriminates between human-generated and machine-generated sentences.

Human Rating Prediction. *COSMic* (Inan et al., 2021) and *Polos* (Wada et al., 2024) explicitly train a model to predict human scores.

4.1.3 Candidate Diversity

Several methods measure the diversity of generated captions. All assume a set of generated captions S for which diversity is measured, where S can either consist of candidates for the entire test set, a set of similar images or a single image.

Lexical Diversity: While some works use lexical similarity methods such as BLEU (*mBLEU*: Shetty et al., 2017; *self-BLEU*: Zhu et al., 2018) and CIDEr (*Self-CIDEr*: Wang and Chan, 2019;

CIDErBtw: Wang et al., 2020a) on pairs of captions from S , others specifically designed methods to measure lexical diversity in S (*Vocab*: Shetty et al., 2017; *Div-n*: Shetty et al., 2017; *diversity edit distance*: Dai et al., 2018a; *Tdiv*: Liu et al., 2019a; *Dist*: Liu et al., 2019a).

Sentence-level Textual Diversity: Other methods compute the percentage of distinct candidates in S (*Distinct*, AKA *Unique*, Wang et al., 2017a), or the percentage of sentences not seen in the training set (*Novel*, Wang et al., 2017a).

Sentence-level Semantic Diversity: Recent methods propose to compare the semantics of candidates as an indicator for diversity. *LSA* (Wang and Chan, 2019) computes a matrix K where $K_{i,j}$ is the dot-product between the bag-of-words vectors of captions i, j in S , and calculates the singular vector decomposition (SVD) of K . *CLIP diversity* (Li et al., 2023b) computes the similarity of the CLIP embeddings of captions in S .

4.1.4 Bias in Candidates

Various metrics have been proposed to quantify the extent of bias evident in captioning models. *BiasAmp* (Zhao et al., 2017) measures the amplification of bias by the model compared to the training set by comparing the correlation between predefined attributes and activities (e.g., female-cooking) in model- and human-generated captions. Hendricks et al. (2018) introduce two metrics: *Gender error*, which assumes that images are labeled as male or female and calculates the number of gender misclassified words in the candidates, and *Gender ratio*, which defines male or female sentences based on the inclusion of predefined gender-related words and computes the ratio of male candidates to female candidates. Others measure bias amplification by training models to predict protected attributes given a caption (*LIC*, Hirota et al., 2022) or an image-caption pair (*ImageCaptioner*², Abdelrahman et al., 2024), and computing the difference between accuracies when training on candidates versus references.

4.2 Human Evaluation

While automatic evaluation offers clear benefits in terms of scale and consistency, it still serves as a surrogate for human evaluation. To ensure that the improvements shown by automatic methods are genuine, many image captioning studies also

apply human evaluation methods to a subset of the data.

The human evaluation taxonomy development process followed a similar approach to that in automatic evaluation (Section 4.1). In this case, a single dimension emerged: the evaluation framework.

Scale Rating. Several studies direct human participants to evaluate candidates on a discrete scale based on various attributes, such as relevance (Maeda et al., 2023), fluency (Wu et al., 2023b) and descriptiveness (Yue et al., 2023).

Comparative. A different approach presents pairs of captions to human participants and asks them to decide which one is better without knowing their source. Early studies compared candidates to references (Kuznetsova et al., 2014; Yatskar et al., 2014), while recent research compares candidates with captions generated by baseline models (Tanaka et al., 2024; Ge et al., 2024).

Yes/No Questions. Certain studies involve human participants answering yes/no questions such as whether the candidate exhibits human-like qualities (Yao et al., 2019) or describes all the objects in the image (Chen et al., 2022).

Retrieval. Another popular approach is to ask participants to perform text-to-image retrieval given the candidate caption (Wang and Chan, 2019; Ou et al., 2023).

Answer Questions Using the Candidate. Nie et al. (2020) prompts participants with questions and requires them to answer using the candidate provided, without showing the images.

5 Metrics Usage Analysis

In this section, we analyze the data from Section 3 to identify usage patterns and present the main findings, excluding 2024 as the data was collected before the year ended. Appendix B lists metric usage for all 314 reviewed papers.

5.1 Automatic Evaluation

A Set of 5 Metrics Dominates. Figure 2 displays the number of papers using different metrics each year. Given the large number of metrics identified (71), we only plot usage for the seven most common metrics, clustering all other metrics by the categories from Section 4.

The figure shows that since 2015, five metrics (BLEU, CIDEr, METEOR, ROUGE, and SPICE; henceforth, the five dominant metrics) have been used substantially more frequently than all other metrics. This trend has begun to shift in recent years, with an increase in the use of other metrics. Notably, four of the dominant metrics (all except SPICE) are lexical similarity metrics, criticized for failing to capture semantic similarity (Giménez and Márquez, 2007). Figure 5(A) illustrates a decline in the usage of lexical metrics in recent years relative to overall metric usage, excluding diversity and bias metrics from this count.

The Number of Metrics per Paper Has Settled at 4–5. Figure 5(B) shows a gradual increase in the mean number of metrics used per paper each year until reaching approximately 4.5 in 2019, and remaining in the range of 4–5 since then.

5.2 Human Evaluation

Human Evaluation Usage is Decreasing. Figure 5(C) shows the fraction of papers with human evaluation each year. From 2010 to 2014, all papers conducted human evaluation. However, the publication of large image-caption datasets made this more challenging, and beginning at 2015 (a year after MSCOCO’s release) we see a general decline in the use of human evaluation, likely due to improvements in automatic evaluation metrics.

Comparative is the Most Common Human Framework. Figure 5(D) presents the number of papers using each human evaluation framework per year. The comparative framework is most frequently used, closely followed by scale rating.

Significance and Agreement Are Rarely Reported. Of 94 papers that used human evaluation (29.9% of the papers we documented), only 19 reported significance, which is crucial in human evaluation because, due to the high cost of human labor, only a small subset of the test set is assessed. Moreover, results may be unreliable if annotator’s agreement is low; however, only 6 of these 94 papers report annotator agreement.

6 Experiments

Prior captioning evaluation work assumed an automatic metric’s quality depends on its alignment with human evaluations. To facilitate this, human rating datasets were created, where participants

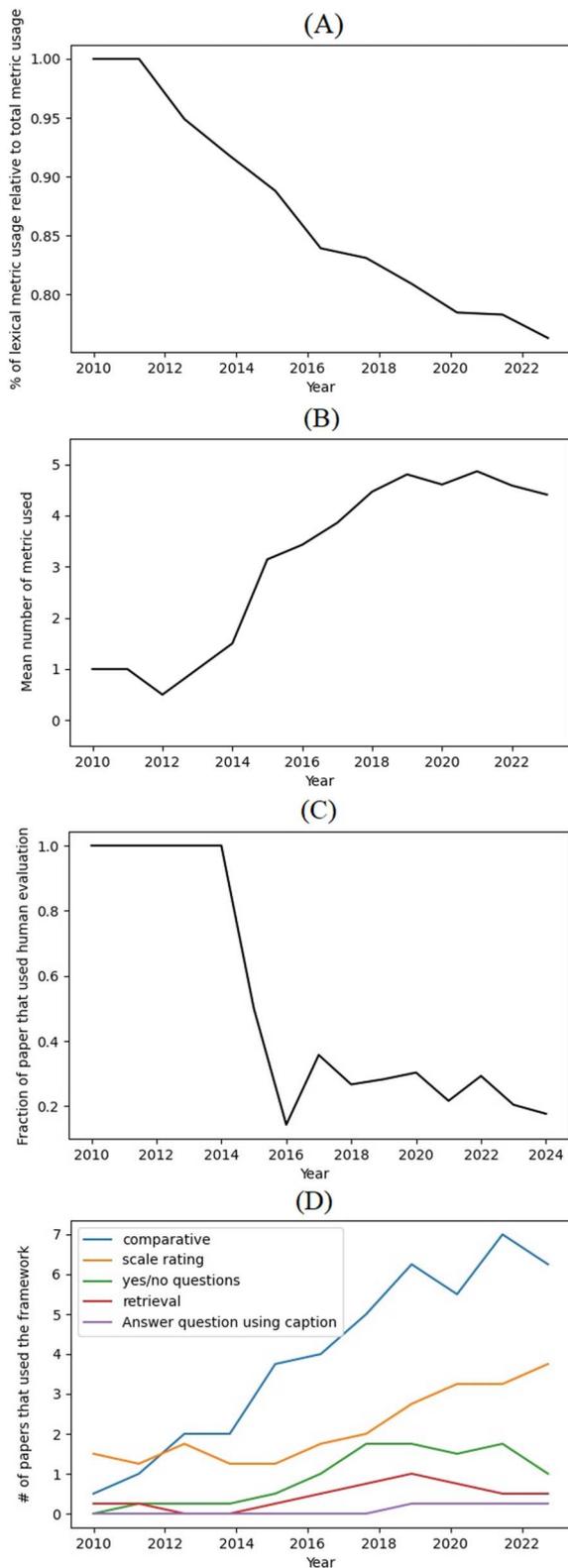


Figure 5: Metrics usage analysis plots. (A) Percentage of lexical metric usage relative to total metric usage each year, smoothed by convolving with a window size of four years. (B) Mean number of metrics used per paper each year. (C) Fraction of papers that performed human evaluation each year. (D) Human evaluation frameworks usage over the years, smoothed by convolving with a window size of four years.

rated candidate captions. A common approach (among others, Anderson et al., 2016; Hessel et al., 2021) is to evaluate metrics by applying them to the candidates in these datasets and measuring the correlation of their scores with the human ratings. We follow this convention but extend prior research by incorporating metrics that have not been previously studied in our experiments.

Throughout our experiments, we focus on the *replacement* subtask (see Section 2). This definition—producing descriptions intended to serve as substitutes for the image—aligns most closely with the human rating datasets, where annotators are instructed to evaluate whether the text accurately reflects the image’s content without adding information not present in the image.⁴

6.1 Experimental Setup

Datasets. We experiment with the following 7 human rating datasets.

Flickr8k-Expert and Flickr8k-CF (Hodosh et al., 2013) include human ratings for captions on the 1000 images of the Flickr8k test set. The captions are human-generated captions selected from the test set, where each image is associated with a set of captions (not necessarily initially generated for that image). Flickr8k-Expert includes 5,822 captions rated by a small and controlled group of native English speakers on a scale from 1 to 4. Flickr8k-CF comprises 47,830 captions, each rated by three or more crowd-sourced workers using a binary label (describes/doesn’t describe the image). For each caption, we take the percentage of positive ratings as the final ground truth score.

Composite (Aditya et al., 2015) includes 11,985 human ratings for captions of 3,925 images from Flickr8k, Flickr30k, and MSCOCO. Each image is paired with two model-generated and one original reference captions. Each caption is rated for correctness and thoroughness on a scale from 1 to 5. Following previous studies (Anderson et al., 2016), we use the correctness rating (also termed “relevance” in some prior research) as the score.

THumB (Kasai et al., 2022) features human ratings of captions for 500 images sourced from MSCOCO. Each image is associated with four model-generated and one original reference captions. Annotators provide ratings for precision and recall (on a scale of 1 to 5), deducting some

⁴See Appendix A.1 for annotators instructions for existing human rating datasets.

penalty points from their average in case of lack of fluency, conciseness and include language to reach the overall score. We take this overall score as the score for this caption.

Polaris (Wada et al., 2024) contains 131,020 human ratings for captions generated for 13,691 images from MSCOCO and nocaps (Agrawal et al., 2019). Ratings were provided on a five-point scale and then transformed to values in the range [0,1] using min-max normalization. The captions were generated by 10 different models. The dataset is partitioned into training, validation, and test sets.

Pascal50S (Vedantam et al., 2015) contains 4,000 pairs of candidate captions for images from the Pascal Sentence dataset (Rashtchian et al., 2010), evenly distributed across four categories: both captions are human-generated for the image (HC), both are human-generated but only one corresponds to the image (HI), one is human-generated while the other is machine-generated (HM), both are machine-generated (MM). Each pair has 48 human judgments indicating which of the two candidates is more similar to a reference caption. Following previous work (Hessel et al., 2021), we take the majority vote for each candidate pair (ties are broken randomly) and randomly select 5 references out of the available 48 for reference-based metrics.

The Reformulations dataset (Berger et al., 2025) includes model-generated captions for 1,405 images from MSCOCO and Flickr30k, along with human-generated reformulations, i.e., corrected versions of the captions (if any errors exist). Evaluation metrics are expected to favor the caption after human correction. The dataset contains 5,208 caption-reformulation pairs; we omit 864 pairs where the reformulation is identical to the caption.

More details on the version we used for each dataset can be found in Appendix A.1.

Metrics. We select a subset of the 71 identified automatic metrics for evaluation. First, we exclude bias metrics, as bias is not captured by the human rating datasets. Next, we filter out all metrics that receive multiple candidates as input (such as diversity metrics), as the human ratings are provided for a single candidate. Finally, we exclude all metrics that rely on a closed model API access (CLAIR, GPT4V evaluation, ALOHa), focusing on publicly accessible evaluation methods.

After filtering, 56 metrics remain. To ensure fair comparison, we conduct all the experiments

Cluster	Metrics
1	BLEU1, BLEU2, BLEU3, BLEU4, CIDEr, Exact NO, Fuzzy NO, METEOR, ROUGE, SPICE
2	BLIP2Score, CLIPImageScore, CLIPScore, MPNetScore, PACScore, RefCLIPScore, RefPACScore, Polos
3	Exact VO, Fuzzy VO

Table 2: Metrics clustered by mutual correlation.

for each metric ourselves, even if previous studies have reported results for this metric. Due to the effort required and the scope of this study, we focus on metrics frequently employed in image captioning research. Specifically, we select metrics used in at least two papers per year on average since their publication.⁵ This narrows our selection to 20 metrics (Table 2).⁶ For details on the implementation of each metric, refer to Appendix A.2.

6.2 Metric Mutual Correlation

To examine relationships among the selected metrics, we analyze their mutual correlations using Spectral Clustering with $N = 3$ clusters,⁷ where the adjacency between two metrics is defined as their mutual correlation on the Polaris train set. The resulting grouping is shown in Table 2. Cluster 1 primarily includes lexical similarity metrics, cluster 2 consists of Transformer-based metrics, and cluster 3 contains verb overlap metrics, which exhibit low correlation with all other metrics.

6.3 ENSEMBLEVAL: An Ensemble of Evaluation Methods

Our survey reveals that while a large number of metrics (71, as detailed in our literature review, Section 3) have been developed to assess diverse properties (outlined in our taxonomy, Section 4), only a small subset from a few categories is widely used (as shown in our usage analysis, Section 5). Furthermore, popular metric groups have faced increasing criticism in recent years. For instance, lexical similarity metrics are often

⁵Not including the publication year and 2024.

⁶We omit two metrics for which we could not obtain a full implementation: PR-MCS and InfoMetC.

⁷The number of clusters was chosen to maximize the Silhouette score.

criticized for failing to capture semantic similarity effectively (Giménez and Márquez, 2007), while the recently popularized reference-free metrics have been challenged for focusing solely on visual grounding errors while neglecting caption implausibility errors (Ahmadi and Agrawal, 2024).

Consequently, we hypothesize that a more diverse evaluation approach, capable of capturing multiple aspects of a candidate caption, could enhance correlation with human ratings. As an initial step toward this goal, we explore a straightforward implementation: an ensemble method combining existing metrics through a linear combination.⁸ We train the linear coefficients on one human rating dataset and evaluate the ensemble on additional datasets. If capturing multiple aspects indeed improves alignment with human judgments, the learned coefficients should generalize across datasets, resulting in improved correlation with human ratings. We refer to our proposed ensemble approach as `ENSEMBEVAL`.

Due to strong correlations among certain metrics, using all metrics in the ensemble might create redundancy. Therefore, we employ a sequential feature selection algorithm to identify a subset of metrics that accurately predict human ratings. At each step, the algorithm adds the best metric to the ensemble based on the cross-validation score of a linear regression estimator of human ratings. The process stops when the estimator’s score does not increase by at least ε . We use the Polaris train set⁹ for selecting the metrics and the validation set to determine the optimal ε value (0.0001).

Next, we standardize the metrics values to the range $[0, 1]$ and find the metrics coefficients by training a linear regression model to predict human ratings, again using the Polaris train dataset. The rescaling ensures that the coefficients are meaningfully comparable.¹⁰ More details on our feature selection and linear regression procedures can be found in Appendix A.3.

Table 3 displays the selected metrics, their cluster assignments (Section 6.2), and their linear regression coefficients. The Transformer-based metrics (cluster 2) have the largest coefficients, making them the most influential in predicting

⁸Future research can explore more sophisticated multi-aspect metrics, which lies beyond the scope of this work.

⁹This is the only human rating dataset that is split into train/val/test.

¹⁰Rescaling did not impact the performance of our model.

Metric	Cluster	Coef
Polos	2	0.55
BLIP2Score	2	0.40
PACScore	2	0.29
Exact NO	1	0.08
BLEU1	1	0.07
Fuzzy VO	3	0.02
ROUGE	1	−0.07
CIDEr	1	−0.15
RefCLIPScore	2	−0.17

Table 3: Metrics selected using feature selection to predict human ratings, their cluster numbers (as shown in Table 2), and linear regression coefficients.

the candidate score. Negative coefficients are observed for some metrics in clusters 1 and 2, likely to mitigate redundancy from multiple metrics within the same cluster. For instance, BLIP2Score, Polos, and PACScore all receive positive coefficients but are highly correlated. To balance this redundancy, a fourth metric from the same cluster (RefCLIPScore) is assigned a negative coefficient.

Considering the frequent use of the five dominant metrics in prior research, we include an ensemble of these metrics (the dominant ensemble) as a baseline. As before, we determine the coefficients for this ensemble using linear regression.

6.4 Correlation with Human Ratings

We now compute the correlation between metric scores and human ratings across all datasets containing human ratings (all but Pascal50S and the Reformulations datasets).¹¹ We follow previous work (Anderson et al., 2016; Hessel et al., 2021; Kasai et al., 2022; Wada et al., 2024) and use Kendall-C correlation for Flickr8k-Expert, Composite, and Polaris; Kendall-B correlation for Flickr8k-CF, and Pearson correlation for THUMB.

Table 4 presents the results. We omit results for five metrics¹² that were proposed ad hoc in an experiment (rather than in a dedicated paper) and showed weak performance.

¹¹For Polaris we use the test set.

¹²CLIPImageScore, Exact NO, Exact VO, Fuzzy NO, Fuzzy VO.

Metric	Flickr8k (Expert)	Flickr8k (CF)	Composite	THumB	Polaris
BLEU4	30.8	16.9	30.6	10.4	46.2
BLEU3	31.5	17.0	30.9	11.8	46.6
BLEU2	32.5	17.9	31.1	15.8	47.1
BLEU1	32.3	17.9	31.4	19.5	45.4
ROUGE	32.3	19.9	32.5	18.7	46.2
METEOR	41.8	22.3	39.0	18.5	51.2
CIDEr	43.9	24.6	37.5	22.4	52.1
Dominant ensemble	43.7	23.3	38.5	23.1	52.3
SPICE	44.9	24.4	40.4	21.0	50.9
CLIPScore	51.4	34.4	53.8	31.9	51.5
PACScore	54.3	36.0	55.7	31.4	52.4
MPNetScore	54.7	35.3	54.7	40.4	53.6
RefCLIPScore	53.0	36.4	55.4	40.7	54.1
RefPACScore	55.9	37.6	57.3	42.4	55.2
BLIP2Score	52.5	36.7	61.5	44.9	53.7
Polos	56.4	37.8	58.0	43.4	57.8
ENSEMBEVAL (ours)	58.5	38.7	61.7	46.6	58.8
Δ from 2nd best	+2.1	+0.9	+0.2	+1.7	+1.0

Table 4: Correlation with human ratings across different datasets. The best performing metric is in **bold**.

An Ensemble Improves over Individual Metrics. ENSEMBEVAL demonstrates the highest correlation with human ratings across all datasets, indicating that integrating various aspects of caption quality yields a more “human-like” score.

Dominant Metrics Are Weaker than Ecent Alternatives. The five dominant metrics, as well as their ensemble, correlate less with human ratings compared to all other examined metrics.

Lesser-known Metrics Prove to Be Valuable. While some metrics were introduced in dedicated papers, others were proposed ad hoc in the experiments section. Previous studies comparing metrics performance have only discussed the former, missing strong correlations with human ratings observed with ad hoc metrics like MPNetScore and BLIP2Score. This underscores the importance of a systematic review of all relevant papers.

6.5 Accuracy on Pairwise Comparison

We also perform a pairwise comparison task on the Pascal50S and Reformulations datasets to assess metrics’ accuracy in assigning a higher score to the human-preferred candidate in each pair. For Pascal50S, we use majority vote to indicate human

preference and report the mean and standard deviation across five random instances of tie-breaking and reference selection. For Reformulations, we consider the reformulated caption as preferred.

Results are presented in Table 5. ENSEMBEVAL attained the highest score on the Reformulations dataset, as well as when comparing human- and machine-generated captions in Pascal 50 (HM). In all other cases MPNetScore and BLIP2Score performed best. As noted earlier, these metrics were absent from previous studies since they were not introduced in dedicated papers.

6.6 Why Does the Ensemble Improve Correlation?

To understand why ENSEMBEVAL improves correlation over individual metrics, we manually examine cases where scores of individual metrics deviate from human ratings while the ensemble score aligns more closely.

These cases often arise from disagreements on which entities to describe. If reference captions omit an entity mentioned in the candidate caption, reference-based metrics assign a low score, while reference-free metrics give a high one. Humans typically rate such captions moderately, as they are accurate but may not align with their descriptive

Metric	Pascal50S					Reformulations
	HC	HI	HM	MM	Mean	
BLEU4	59.8 ± 0.7	92.3 ± 0.7	85.6 ± 0.5	57.3 ± 1.2	73.7 ± 0.5	49.8
BLEU3	60.6 ± 0.9	93.0 ± 0.7	88.1 ± 0.5	57.6 ± 1.0	74.8 ± 0.6	51.8
ROUGE	63.1 ± 0.8	95.2 ± 0.6	91.9 ± 0.3	60.3 ± 1.0	77.6 ± 0.4	53.1
BLEU2	62.9 ± 1.7	94.0 ± 0.4	90.1 ± 0.2	58.6 ± 1.0	76.4 ± 0.7	54.9
BLEU1	62.8 ± 1.3	95.1 ± 0.5	91.5 ± 0.4	59.8 ± 0.7	77.3 ± 0.2	56.0
CIDEr	65.6 ± 1.7	98.0 ± 0.4	90.9 ± 0.3	64.9 ± 1.1	79.8 ± 0.4	54.3
SPICE	61.2 ± 1.9	94.3 ± 0.7	85.5 ± 0.7	49.5 ± 0.7	72.6 ± 0.7	67.6
METEOR	64.4 ± 1.7	97.6 ± 0.4	94.1 ± 0.8	65.3 ± 0.9	80.4 ± 0.6	66.3
Dominant ensemble	65.5 ± 1.7	97.7 ± 0.3	93.0 ± 0.5	68.0 ± 0.7	81.1 ± 0.6	66.0
RefPACScore	67.6 ± 0.7	99.6 ± 0.1	96.0 ± 0.2	75.7 ± 0.6	84.7 ± 0.3	73.3
RefCLIPScore	64.1 ± 1.2	99.6 ± 0.1	95.8 ± 0.4	72.9 ± 0.6	83.1 ± 0.4	74.9
BLIP2Score	60.5 ± 0.2	99.8 ± 0.0	96.3 ± 0.0	71.2 ± 0.5	82.0 ± 0.1	76.1
Polos	69.5 ± 1.2	99.6 ± 0.0	97.0 ± 0.3	78.5 ± 0.6	86.1 ± 0.1	73.0
MPNetScore	71.9 ± 0.6	99.8 ± 0.1	96.3 ± 0.5	79.0 ± 0.6	86.7 ± 0.2	72.7
PACScore	60.4 ± 0.1	99.3 ± 0.0	96.8 ± 0.0	72.9 ± 0.4	82.4 ± 0.1	77.2
CLIPScore	56.1 ± 0.2	99.3 ± 0.0	96.3 ± 0.0	71.2 ± 0.4	80.7 ± 0.1	80.8
ENSEMBEVAL (ours)	68.7 ± 0.6	99.8 ± 0.1	98.3 ± 0.3	77.3 ± 0.3	86.0 ± 0.2	81.4

Table 5: Accuracy in pairwise comparison. In Pascal50S we report mean and standard deviation across five random instances of tie-breaking and reference selection. HC: both captions are human-generated for the image, HI: both are human-generated but only one corresponds to the image, HM: one caption is human-generated while the other is machine-generated, MM: both captions are machine-generated. In each dataset, best scoring metric is in **bold**.

choices. The ensemble, combining both metric types, aligns more closely with human ratings. This issue is especially relevant beyond English, as research shows speakers of different languages tend to describe different entities (Berger and Ponti, 2025). See Figure 6 for an example.

7 Software Package

Our empirical analysis of the redundancy/complementarity of the different metrics, reveals that an ensemble of several methods may be considerably better correlated with human rankings of the captions. To this aim, we release a software package, ENSEMBEVAL, under the permissive MIT license.¹³

The software package implements the ensemble metric presented in Section 6.3. It also straightforwardly supports the use of alternative sets of weights. The software package can easily be extended with additional metrics, either reference-based or reference-free. The package is therefore usable not only for replication of the results presented in this paper, but also as a generic tool for ensembling image captioning evaluation metrics.

¹³github.com/uriberger/caption_evaluation.

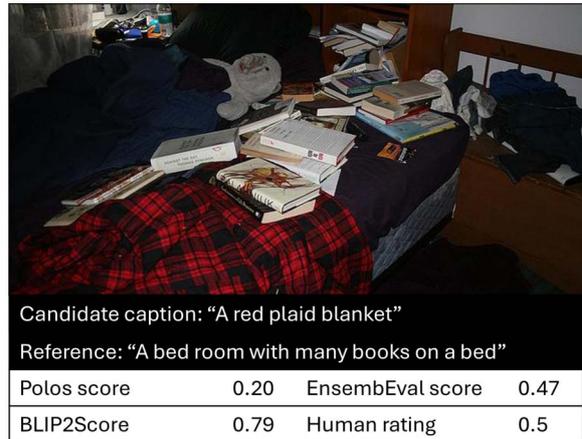


Figure 6: An image, corresponding candidate caption, reference example and evaluation scores. Since no reference captions mention the blanket, Polos (a reference-based metric) assigns a low score. However, as the caption accurately describes the image, BLIP2Score (a reference-free metric) gives a high score. The ensemble score falls in between, aligning best with human ratings. Image and caption taken from the Polaris dataset.

8 Related Work

8.1 Image Captioning Surveys

To the best of our knowledge, no previous work has exclusively surveyed image captioning

evaluation methods. However, most image captioning surveys include a dedicated section for automatic evaluation methods, typically mentioning the five dominant metrics along with one or two additional ones (e.g., Hossain et al., 2019; Liu et al., 2019b; Ghandi et al., 2023; Sharma and Padha, 2023).

A few studies (Stefanini et al., 2022; Xu et al., 2023) provide a more detailed taxonomy of automatic metrics. Our work differs from them in two ways: First, our systematic review identifies strong metrics not covered previously. Second, while these studies group metrics into loosely defined categories (e.g., “standard” metrics), we clearly define our categorization criteria based on the properties each metric aims to measure.

We are also the first to develop a taxonomy for human evaluation frameworks. Before our study, Bernardi et al. (2016) discussed human evaluation but only mentioned the scale rating framework.

Prior to our study, Staniūtė and Šešok (2019) reported metric usage from selected papers; however, they manually selected popular papers rather than identifying them systematically. Similar to our work, two studies (Sharma, 2021; Al-Shamayleh et al., 2024) systematically identify image captioning papers and report metric usage, but cover much fewer papers (79 and 41, respectively) and metrics (3 and 8, respectively), and do not report trends or provide in-depth analysis.

8.2 Image Captioning Metrics Analysis

Several previous studies have analyzed the effectiveness of image captioning metrics. Hodosh et al. (2013) compare BLEU and ROUGE scores to human ratings and find that these metrics fall short in measuring content quality. Elliott and Keller (2014) compute the correlation of lexical similarity metrics with human ratings and find that METEOR had the strongest correlation. Kilickaya et al. (2017) compare the five dominant metrics and WMD, finding that n -gram based metrics exhibit lower performance than SPICE and WMD.

Recently, as reference-free metrics like CLIP-Score gain prominence, there has been growing criticism directed at this line of research. Deutsch et al. (2022) argue that reference-free metrics can be exploited at test time to find outputs that maximize their scores. Ahmadi and Agrawal (2024) find these metrics sensitive to visual grounding errors but not to caption implausibility.

9 Conclusion

We conducted a comprehensive survey of image captioning evaluation methods, yielding two key outcomes: a taxonomy of captioning metrics, including overlooked high-performing ones (e.g., BLIP2Score) and an analysis revealing that most papers rely on five metrics that have only a weak correlation with human ratings. We further showed that an ensemble of several existing metrics improves correlation with human ratings, and proposed a simple weighted ensemble method to this effect. We release a software package that implements this ensemble methods and facilitates future developments in image captioning metric ensembling methods.

Acknowledgments

This work was supported in part by the Israel Science Foundation (grant no. 2424/21), by a grant from the Israeli Planning and Budgeting Committee (PBD), and by the HUJI-UoM joint PhD program.

References

- Eslam Abdelrahman, Pengzhan Sun, Li Erran Li, and Mohamed Elhoseiny. 2024. Image Captioner2: Image captioner for image captioning bias amplification assessment. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 20902–20911. <https://doi.org/10.1609/aaai.v38i19.30080>
- Somak Aditya, Yezhou Yang, Chitta Baral, Cornelia Fermuller, and Yiannis Aloimonos. 2015. From images to sentences through scene description graphs using commonsense reasoning and knowledge. *arXiv preprint arXiv:1511.03292*.
- Harsh Agrawal, Karan Desai, Yufei Wang, Xinlei Chen, Rishabh Jain, Mark Johnson, Dhruv Batra, Devi Parikh, Stefan Lee, and Peter Anderson. 2019. nocaps: Novel object captioning at scale. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*. <https://doi.org/10.1109/ICCV.2019.00904>
- Saba Ahmadi and Aishwarya Agrawal. 2024. An examination of the robustness of reference-free

- image captioning evaluation metrics. In *Findings of the Association for Computational Linguistics: EACL 2024*, pages 196–208, St. Julian's, Malta. Association for Computational Linguistics.
- Hiba Ahsan, Daivat Bhatt, Kaivan Shah, and Nikita Bhalla. 2021. Multi-modal image captioning for the visually impaired. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Student Research Workshop*, pages 53–60, Online. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2021.naacl-srw.8>
- Ahmet Aker and Robert Gaizauskas. 2010. Generating image descriptions using dependency relational patterns. In *Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics*, pages 1250–1258, Uppsala, Sweden. Association for Computational Linguistics.
- Ahmad Sami Al-Shamayleh, Omar Adwan, Mohammad A. Alsharaiah, Abdelrahman H. Hussein, Qasem M. Kharma, and Christopher Ifeanyi Eke. 2024. A comprehensive literature review on image captioning methods and metrics based on deep learning technique. *Multimedia Tools and Applications*, 83(12):34219–34268. <https://doi.org/10.1007/s11042-024-18307-8>
- Malihe Alikhani, Piyush Sharma, Shengjie Li, Radu Soricut, and Matthew Stone. 2020. Cross-modal coherence modeling for caption generation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6525–6535, Online. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2020.acl-main.583>
- Aliki Anagnostopoulou, Mareike Hartmann, and Daniel Sonntag. 2023. Towards adaptable and interactive image captioning with data augmentation and episodic memory. In *Proceedings of The Fourth Workshop on Simple and Efficient Natural Language Processing (SustaiNLP)*, pages 245–256, Toronto, Canada (Hybrid). Association for Computational Linguistics. <https://doi.org/10.18653/v1/2023.sustainlp-1.19>
- Peter Anderson, Basura Fernando, Mark Johnson, and Stephen Gould. 2016. Spice: Semantic propositional image caption evaluation. In *Computer Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part V 14*, pages 382–398. Springer. https://doi.org/10.1007/978-3-319-46454-1_24
- Peter Anderson, Basura Fernando, Mark Johnson, and Stephen Gould. 2017. Guided open vocabulary image captioning with constrained beam search. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 936–945, Copenhagen, Denmark. Association for Computational Linguistics. <https://doi.org/10.18653/v1/D17-1098>
- Peter Anderson, Stephen Gould, and Mark Johnson. 2018a. Partially-supervised image captioning. In *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc.
- Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang. 2018b. Bottom-up and top-down attention for image captioning and visual question answering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. <https://doi.org/10.1109/CVPR.2018.00636>
- Jyoti Aneja, Harsh Agrawal, Dhruv Batra, and Alexander Schwing. 2019. Sequential latent spaces for modeling the intention during diverse image captioning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*. <https://doi.org/10.1109/ICCV.2019.00436>
- Jyoti Aneja, Aditya Deshpande, and Alexander G. Schwing. 2018. Convolutional image captioning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. <https://doi.org/10.1109/CVPR.2018.00583>
- Satanjeev Banerjee and Alon Lavie. 2005. ME-TTEOR: 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, pages 65–72, Ann Arbor, Michigan. Association for Computational Linguistics.
- Manuele Barraco, Sara Sarto, Marcella Cornia, Lorenzo Baraldi, and Rita Cucchiara. 2023. With a little help from your own past: Prototypical memory networks for image captioning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 3021–3031. <https://doi.org/10.1109/ICCV51070.2023.00282>
- Uri Berger, Omri Abend, Lea Frermann, and Gabriel Stanovsky. 2025. Improving image captioning by mimicking human reformulation feedback at inference-time. *arXiv preprint arXiv:2501.04513*.
- Uri Berger and Edoardo Ponti. 2025. Cross-lingual and cross-cultural variation in image descriptions. In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 9453–9465, Albuquerque, New Mexico. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2025.naacl-long.478>
- Raffaella Bernardi, Ruket Cakici, Desmond Elliott, Aykut Erdem, Erkut Erdem, Nazli Ikizler-Cinbis, Frank Keller, Adrian Muscat, and Barbara Plank. 2016. Automatic description generation from images: A survey of models, datasets, and evaluation measures. *Journal of Artificial Intelligence Research*, 55:409–442. <https://doi.org/10.1613/jair.4900>
- Romain Bielawski and Rufin VanRullen. 2023. CLIP-based image captioning via unsupervised cycle-consistency in the latent space. In *Proceedings of the 8th Workshop on Representation Learning for NLP (RepL4NLP 2023)*, pages 266–275, Toronto, Canada. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2023.repl4nlp-1.22>
- Ali Furkan Biten, Lluís Gomez, Marçal Rusinol, and Dimosthenis Karatzas. 2019. Good news, everyone! Context driven entity-aware captioning for news images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. <https://doi.org/10.1109/CVPR.2019.01275>
- Alexander Black, Jing Shi, Yifei Fan, Tu Bui, and John Collomosse. 2024. VIXEN: Visual text comparison network for image difference captioning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 846–854. <https://doi.org/10.1609/aaai.v38i2.27843>
- Emanuele Bugliarello and Desmond Elliott. 2021. The role of syntactic planning in compositional image captioning. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 593–607, Online. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2021.eacl-main.48>
- Michele Cafagna, Kees van Deemter, and Albert Gatt. 2022. Understanding cross-modal interactions in V&L models that generate scene descriptions. In *Proceedings of the Workshop on Unimodal and Multimodal Induction of Linguistic Structures (UM-IoS)*, pages 56–72, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics. <https://doi.org/10.18653/v1/2022.umios-1.6>
- Daigang Cai, Lichen Zhao, Jing Zhang, Lu Sheng, and Dong Xu. 2022. 3DJCG: A unified framework for joint dense captioning and visual grounding on 3d point clouds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 16464–16473. <https://doi.org/10.1109/CVPR52688.2022.01597>
- Tingjia Cao, Ke Han, Xiaomei Wang, Lin Ma, Yanwei Fu, Yu-Gang Jiang, and Xiangyang Xue. 2020. Feature deformation meta-networks in image captioning of novel objects. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(07):10494–10501. <https://doi.org/10.1609/aaai.v34i07.6620>
- David Chan, Austin Myers, Sudheendra Vijayanarasimhan, David Ross, and John Canny. 2023a. IC3: Image captioning by committee consensus. In *Proceedings of the 2023 Conference on Empirical Methods in Natural*

- Language Processing*, pages 8975–9003, Singapore. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2023.emnlp-main.556>
- David Chan, Suzanne Petryk, Joseph Gonzalez, Trevor Darrell, and John Canny. 2023b. CLAIR: Evaluating image captions with large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 13638–13646, Singapore. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2023.emnlp-main.841>
- Arjun Chandrasekaran, Devi Parikh, and Mohit Bansal. 2018. Punny captions: Witty word-play in image descriptions. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 770–775, New Orleans, Louisiana. Association for Computational Linguistics. <https://doi.org/10.18653/v1/N18-2121>
- Soravit Changpinyo, Bo Pang, Piyush Sharma, and Radu Soricut. 2019. Decoupled box proposal and featurization with ultrafine-grained semantic labels improve image captioning and visual question answering. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 1468–1474, Hong Kong, China. Association for Computational Linguistics. <https://doi.org/10.18653/v1/D19-1155>
- Chen Chen, Shuai Mu, Wanpeng Xiao, Zexiong Ye, Liesi Wu, and Qi Ju. 2019a. Improving image captioning with conditional generative adversarial nets. *Proceedings of the AAAI Conference on Artificial Intelligence*, 33(01):8142–8150. <https://doi.org/10.1609/aaai.v33i01.33018142>
- Cheng-Kuan Chen, Zhufeng Pan, Ming-Yu Liu, and Min Sun. 2019b. Unsupervised stylish image description generation via domain layer norm. *Proceedings of the AAAI Conference on Artificial Intelligence*, 33(01):8151–8158. <https://doi.org/10.1609/aaai.v33i01.33018151>
- Fuhai Chen, Rongrong Ji, Jiayi Ji, Xiaoshuai Sun, Baochang Zhang, Xuri Ge, Yongjian Wu, Feiyue Huang, and Yan Wang. 2019c. Variational structured semantic inference for diverse image captioning. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc.
- Fuhai Chen, Rongrong Ji, Xiaoshuai Sun, Yongjian Wu, and Jinsong Su. 2018a. Groupcap: Group-based image captioning with structured relevance and diversity constraints. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. <https://doi.org/10.1109/CVPR.2018.00146>
- Hui Chen, Guiguang Ding, Sicheng Zhao, and Jungong Han. 2018b. Temporal-difference learning with sampling baseline for image captioning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 32(1). <https://doi.org/10.1609/aaai.v32i1.12263>
- Hongge Chen, Huan Zhang, Pin-Yu Chen, Jinfeng Yi, and Cho-Jui Hsieh. 2018c. Attacking visual language grounding with adversarial examples: A case study on neural image captioning. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2587–2597, Melbourne, Australia. Association for Computational Linguistics. <https://doi.org/10.18653/v1/P18-1241>
- Jun Chen, Han Guo, Kai Yi, Boyang Li, and Mohamed Elhoseiny. 2022. VisualGPT: Data-efficient adaptation of pretrained language models for image captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 18030–18040. <https://doi.org/10.1109/CVPR52688.2022.01750>
- Jia Chen and Qin Jin. 2020. Better captioning with sequence-level exploration. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. <https://doi.org/10.1109/CVPR42600.2020.01090>
- Jianfu Chen, Polina Kuznetsova, David Warren, and Yejin Choi. 2015. Déjà image-captions: A corpus of expressive descriptions in repetition. In *Proceedings of the 2015 Conference*

- of the North American Chapter of the Association for Computational Linguistics: *Human Language Technologies*, pages 504–514, Denver, Colorado. Association for Computational Linguistics. <https://doi.org/10.3115/v1/N15-1053>
- Jia Chen, Yike Wu, Shiwan Zhao, and Qin Jin. 2021a. Language resource efficient learning for captioning. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 1887–1895, Punta Cana, Dominican Republic. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2021.findings-emnlp.162>
- Long Chen, Zhihong Jiang, Jun Xiao, and Wei Liu. 2021b. Human-like controllable image captioning with verb-specific semantic roles. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 16846–16856. <https://doi.org/10.1109/CVPR46437.2021.01657>
- Long Chen, Hanwang Zhang, Jun Xiao, Liqiang Nie, Jian Shao, Wei Liu, and Tat-Seng Chua. 2017a. SCA-CNN: Spatial and channel-wise attention in convolutional networks for image captioning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. <https://doi.org/10.1109/CVPR.2017.667>
- Minghai Chen, Guiguang Ding, Sicheng Zhao, Hui Chen, Qiang Liu, and Jungong Han. 2017b. Reference based LSTM for image captioning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 31(1). <https://doi.org/10.1609/aaai.v31i1.11198>
- Shi Chen and Qi Zhao. 2018. Boosted attention: Leveraging human attention for image captioning. In *Proceedings of the European Conference on Computer Vision (ECCV)*. https://doi.org/10.1007/978-3-030-01252-6_5
- Shizhe Chen, Qin Jin, Peng Wang, and Qi Wu. 2020. Say as you wish: Fine-grained control of image caption generation with abstract scene graphs. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. <https://doi.org/10.1109/CVPR42600.2020.00998>
- Sijin Chen, Hongyuan Zhu, Xin Chen, Yinjie Lei, Gang Yu, and Tao Chen. 2023a. End-to-end 3d dense captioning with Vote2Cap-DETR. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11124–11133. <https://doi.org/10.1109/CVPR52729.2023.01070>
- Tianlang Chen, Zhongping Zhang, Quanzeng You, Chen Fang, Zhaowen Wang, Hailin Jin, and Jiebo Luo. 2018d. “factual” or “emotional”: Stylized image captioning with adaptive learning and attention. In *Proceedings of the European Conference on Computer Vision (ECCV)*. https://doi.org/10.1007/978-3-030-01249-6_32
- Xinlei Chen and C. Lawrence Zitnick. 2015. Mind’s eye: A recurrent visual representation for image caption generation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. <https://doi.org/10.1109/CVPR.2015.7298856>
- Xinpeng Chen, Lin Ma, Wenhao Jiang, Jian Yao, and Wei Liu. 2018e. Regularizing rnns for caption generation by reconstructing the past with the present. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. <https://doi.org/10.1109/CVPR.2018.00834>
- Zhenyu Chen, Ali Gholami, Matthias Niessner, and Angel X. Chang. 2021c. Scan2Cap: Context-Aware Dense Captioning in RGB-D Scans. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3193–3203. <https://doi.org/10.1109/CVPR46437.2021.00321>
- Zhenyu Chen, Ronghang Hu, Xinlei Chen, Matthias Nießner, and Angel X. Chang. 2023b. UniT3D: A unified transformer for 3D dense captioning and visual grounding. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 18109–18119. <https://doi.org/10.1109/ICCV51070.2023.01660>
- Jaemin Cho, Seunghyun Yoon, Ajinkya Kale, Franck Deroncourt, Trung Bui, and Mohit Bansal. 2022. Fine-grained image captioning with CLIP reward. In *Findings of the Association for Computational Linguistics:*

- NAACL 2022, pages 517–527, Seattle, United States. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2022.findings-naacl.39>
- Cesc Chunseong Park, Byeongchang Kim, and Gunhee Kim. 2017. Attend to you: Personalized image captioning with context sequence memory networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Reuben Cohn-Gordon, Noah Goodman, and Christopher Potts. 2018. Pragmatically informative image captioning with character-level inference. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 439–443, New Orleans, Louisiana. Association for Computational Linguistics. <https://doi.org/10.18653/v1/N18-2070>
- Marcella Cornia, Lorenzo Baraldi, and Rita Cucchiara. 2019. Show, control and tell: A framework for generating controllable and grounded captions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8307–8316. <https://doi.org/10.1109/CVPR.2019.00850>
- Marcella Cornia, Matteo Stefanini, Lorenzo Baraldi, and Rita Cucchiara. 2020. Meshed-memory transformer for image captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. <https://doi.org/10.1109/CVPR42600.2020.01059>
- Yin Cui, Guandao Yang, Andreas Veit, Xun Huang, and Serge Belongie. 2018. Learning to evaluate image captioning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. <https://doi.org/10.1109/CVPR.2018.00608>
- Bo Dai, Sanja Fidler, and Dahua Lin. 2018a. A neural compositional paradigm for image captioning. In *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc.
- Bo Dai, Sanja Fidler, Raquel Urtasun, and Dahua Lin. 2017. Towards diverse and natural image descriptions via a conditional gan. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*. <https://doi.org/10.1109/ICCV.2017.323>
- Bo Dai and Dahua Lin. 2017. Contrastive learning for image captioning. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.
- Bo Dai, Deming Ye, and Dahua Lin. 2018b. Rethinking the form of latent states in image captioning. In *Proceedings of the European Conference on Computer Vision (ECCV)*. https://doi.org/10.1007/978-3-030-01228-1_18
- Riccardo Del Chiaro, Bartłomiej Twardowski, Andrew Bagdanov, and Joost van de Weijer. 2020. RATT: Recurrent attention to transient tasks for continual image captioning. In *Advances in Neural Information Processing Systems*, volume 33, pages 16736–16748. Curran Associates, Inc.
- Chaorui Deng, Ning Ding, Mingkui Tan, and Qi Wu. 2020. Length-controllable image captioning. In *Computer Vision – ECCV 2020*, pages 712–729, Cham. Springer International Publishing. https://doi.org/10.1007/978-3-030-58601-0_42
- Aditya Deshpande, Jyoti Aneja, Liwei Wang, Alexander G. Schwing, and David Forsyth. 2019. Fast, diverse and accurate image captioning guided by part-of-speech. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. <https://doi.org/10.1109/CVPR.2019.01095>
- Roberto Dessì, Michele Bevilacqua, Eleonora Gualdoni, Nathanaël Carraz Rakotonirina, Francesca Franzon, and Marco Baroni. 2023. Cross-domain image captioning with discriminative finetuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6935–6944. <https://doi.org/10.1109/CVPR52729.2023.00670>
- Daniel Deutsch, Rotem Dror, and Dan Roth. 2022. On the limitations of reference-free evaluations of generated text. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*,

- pages 10960–10977, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2022.emnlp-main.753>
- Jacob Devlin, Hao Cheng, Hao Fang, Saurabh Gupta, Li Deng, Xiaodong He, Geoffrey Zweig, and Margaret Mitchell. 2015. Language models for image captioning: The quirks and what works. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 100–105, Beijing, China. Association for Computational Linguistics. <https://doi.org/10.3115/v1/P15-2017>
- George Doddington. 2002. Automatic evaluation of machine translation quality using n-gram co-occurrence statistics. In *Proceedings of the second international conference on Human Language Technology Research*, pages 138–145. <https://doi.org/10.3115/1289189.1289273>
- Pierre Dognin, Igor Melnyk, Youssef Mroueh, Jerret Ross, and Tom Sercu. 2019. Adversarial semantic alignment for improved image captions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. <https://doi.org/10.1109/CVPR.2019.01071>
- Jeffrey Donahue, Lisa Anne Hendricks, Sergio Guadarrama, Marcus Rohrbach, Subhashini Venugopalan, Kate Saenko, and Trevor Darrell. 2015. Long-term recurrent convolutional networks for visual recognition and description. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. <https://doi.org/10.1109/CVPR.2015.7298878>
- Desmond Elliott and Frank Keller. 2013. Image description using visual dependency representations. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1292–1302, Seattle, Washington, USA. Association for Computational Linguistics. <https://doi.org/10.18653/v1/D13-1128>
- Desmond Elliott and Frank Keller. 2014. Comparing automatic evaluation measures for image description. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 452–457, Baltimore, Maryland. Association for Computational Linguistics. <https://doi.org/10.3115/v1/P14-2074>
- Yi Ma, Li Deng, and Geoffrey Zweig. 2015. Image captioning with a deep neural network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. <https://doi.org/10.1109/CVPR.2015.7298754>
- Yi Ma, Li Deng, and Geoffrey Zweig. 2015. From captions to visual concepts and back. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. <https://doi.org/10.1109/CVPR.2015.7298754>
- Zhihao Fan, Zhongyu Wei, Siyuan Wang, and Xuanjing Huang. 2019. Bridging by word: Image grounded vocabulary construction for visual captioning. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 6514–6524, Florence, Italy. Association for Computational Linguistics. <https://doi.org/10.18653/v1/P19-1652>
- Hao Fang, Saurabh Gupta, Forrest Iandola, Rupesh K. Srivastava, Li Deng, Piotr Dollar, Jianfeng Gao, Xiaodong He, Margaret Mitchell, John C. Platt, C. Lawrence Zitnick, and Geoffrey Zweig. 2015. From captions to visual concepts and back. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. <https://doi.org/10.1109/CVPR.2015.7298754>
- Zhiyuan Fang, Jianfeng Wang, Xiaowei Hu, Lin Liang, Zhe Gan, Lijuan Wang, Yezhou Yang, and Zicheng Liu. 2022. Injecting semantic concepts into end-to-end image captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 18009–18019. <https://doi.org/10.1109/CVPR52688.2022.01748>
- Junjie Fei, Teng Wang, Jinrui Zhang, Zhenyu He, Chengjie Wang, and Feng Zheng. 2023a. Transferable decoding with visual entities for zero-shot image captioning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 3136–3146.
- Zhengcong Fei. 2021a. Memory-augmented image captioning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(2):1317–1324. <https://doi.org/10.1609/aaai.v35i2.16220>

- Zhengcong Fei. 2021b. Partially non-autoregressive image captioning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(2): 1309–1316. <https://doi.org/10.1609/aaai.v35i2.16219>
- Zhengcong Fei. 2022. Attention-aligned transformer for image captioning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 36(1):607–615. <https://doi.org/10.1609/aaai.v36i1.19940>
- Zhengcong Fei, Mingyuan Fan, Li Zhu, Junshi Huang, Xiaoming Wei, and Xiaolin Wei. 2023b. Uncertainty-aware image captioning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 37(1):614–622. <https://doi.org/10.1609/aaai.v37i1.25137>
- Zhengcong Fei, Xu Yan, Shuhui Wang, and Qi Tian. 2022. DeeCap: Dynamic early exiting for efficient image captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 12216–12226. <https://doi.org/10.1109/CVPR52688.2022.01190>
- Joshua Feinglass and Yezhou Yang. 2021. SMURF: SeMantic and linguistic UndeRstanding fusion for caption evaluation via typicality analysis. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 2250–2260, Online. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2021.acl-long.175>
- Steven Y. Feng, Kevin Lu, Zhuofu Tao, Malihe Alikhani, Teruko Mitamura, Eduard Hovy, and Varun Gangal. 2022. Retrieve, caption, generate: Visual grounding for enhancing common-sense in text generation models. *Proceedings of the AAAI Conference on Artificial Intelligence*, 36(10):10618–10626. <https://doi.org/10.1609/aaai.v36i10.21306>
- Yansong Feng and Mirella Lapata. 2010. How many words is a picture worth? Automatic caption generation for news images. In *Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics*, pages 1239–1249, Uppsala, Sweden. Association for Computational Linguistics.
- Yang Feng, Lin Ma, Wei Liu, and Jiebo Luo. 2019. Unsupervised image captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. <https://doi.org/10.1109/CVPR.2019.00425>
- Zhongtian Fu, Kefei Song, Luping Zhou, and Yang Yang. 2024. Noise-aware image captioning with progressively exploring mismatched words. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(11):12091–12099. <https://doi.org/10.1609/aaai.v38i11.29097>
- Chuang Gan, Zhe Gan, Xiaodong He, Jianfeng Gao, and Li Deng. 2017a. Stylenet: Generating attractive visual captions with styles. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Zhe Gan, Chuang Gan, Xiaodong He, Yunchen Pu, Kenneth Tran, Jianfeng Gao, Lawrence Carin, and Li Deng. 2017b. Semantic compositional networks for visual captioning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. <https://doi.org/10.1109/CVPR.2017.127>
- Junlong Gao, Shiqi Wang, Shanshe Wang, Siwei Ma, and Wen Gao. 2019a. Self-critical n-step training for image captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Lianli Gao, Kaixuan Fan, Jingkuan Song, Xianglong Liu, Xing Xu, and Heng Tao Shen. 2019b. Deliberate attention networks for image captioning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 33(01):8320–8327. <https://doi.org/10.1609/aaai.v33i01.33018320>
- Yiqi Gao, Xinglin Hou, Yuanmeng Zhang, Tiezheng Ge, Yuning Jiang, and Peng Wang. 2022. CapOnImage: Context-driven dense-captioning on image. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 3449–3465, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2022.emnlp-main.226>
- Hongwei Ge, Zehang Yan, Kai Zhang, Mingde Zhao, and Liang Sun. 2019. Exploring overall

- contextual information for image captioning in human-like cognitive style. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*. <https://doi.org/10.1109/ICCV.2019.00184>
- Yunhao Ge, Xiaohui Zeng, Jacob Samuel Huffman, Tsung-Yi Lin, Ming-Yu Liu, and Yin Cui. 2024. Visual fact checker: Enabling high-fidelity detailed caption generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14033–14042. <https://doi.org/10.1109/CVPR52733.2024.01331>
- Taraneh Ghandi, Hamidreza Pourreza, and Hamidreza Mahyar. 2023. Deep learning approaches on image captioning: A review. *ACM Computing Surveys*, 56(3):1–39. <https://doi.org/10.1145/3617592>
- Jesús Giménez and Lluís Màrquez. 2007. Linguistic features for automatic evaluation of heterogenous MT systems. In *Proceedings of the Second Workshop on Statistical Machine Translation*, pages 256–264, Prague, Czech Republic. Association for Computational Linguistics. <https://doi.org/10.3115/1626355.1626393>
- Othón González-Chávez, Guillermo Ruiz, Daniela Moctezuma, and Tania Ramirez-delReal. 2024. Are metrics measuring what they should? An evaluation of image captioning task metrics. *Signal Processing: Image Communication*, 120:117071. <https://doi.org/10.1016/j.image.2023.117071>
- Jiuxiang Gu, Jianfei Cai, Gang Wang, and Tsuhan Chen. 2018a. Stack-captioning: Coarse-to-fine learning for image captioning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 32(1). <https://doi.org/10.1609/aaai.v32i1.12266>
- Jiuxiang Gu, Shafiq Joty, Jianfei Cai, and Gang Wang. 2018b. Unpaired image captioning by language pivoting. In *Proceedings of the European Conference on Computer Vision (ECCV)*.
- Jiuxiang Gu, Shafiq Joty, Jianfei Cai, Handong Zhao, Xu Yang, and Gang Wang. 2019. Unpaired image captioning via scene graph alignments. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*. <https://doi.org/10.1109/ICCV.2019.01042>
- Jiuxiang Gu, Gang Wang, Jianfei Cai, and Tsuhan Chen. 2017. An empirical study of language cnn for image captioning. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*. <https://doi.org/10.1109/ICCV.2017.138>
- Longteng Guo, Jing Liu, Peng Yao, Jiangwei Li, and Hanqing Lu. 2019. MSCap: Multi-style image captioning with unpaired stylized text. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. <https://doi.org/10.1109/CVPR.2019.00433>
- Longteng Guo, Jing Liu, Xinxin Zhu, Peng Yao, Shichen Lu, and Hanqing Lu. 2020. Normalized and geometry-aware self-attention network for image captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. <https://doi.org/10.1109/CVPR42600.2020.01034>
- Tszhang Guo, Shiyu Chang, Mo Yu, and Kun Bai. 2018. Improving reinforcement learning based image captioning with natural language prior. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 751–756, Brussels, Belgium. Association for Computational Linguistics. <https://doi.org/10.18653/v1/D18-1083>
- Zixin Guo, Tzu-Jui Wang, and Jorma Laaksonen. 2022. CLIP4IDC: CLIP for image difference captioning. In *Proceedings of the 2nd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 12th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 33–42, Online only. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2022.aacl-short.5>
- Danna Gurari, Yinan Zhao, Meng Zhang, and Nilavra Bhattacharya. 2020. Captioning images taken by people who are blind. In *Computer Vision – ECCV 2020*, pages 417–434, Cham. Springer International Publishing. https://doi.org/10.1007/978-3-030-58520-4_25

- Sen He, Hamed R. Tavakoli, Ali Borji, and Nicolas Pugeault. 2019. Human attention in image captioning: Dataset and analysis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*. <https://doi.org/10.1109/ICCV.2019.00862>
- Lisa Anne Hendricks, Kaylee Burns, Kate Saenko, Trevor Darrell, and Anna Rohrbach. 2018. Women also snowboard: Overcoming bias in captioning models. In *Proceedings of the European Conference on Computer Vision (ECCV)*. https://doi.org/10.1007/978-3-030-01219-9_47
- Lisa Anne Hendricks, Subhashini Venugopalan, Marcus Rohrbach, Raymond Mooney, Kate Saenko, and Trevor Darrell. 2016. Deep compositional captioning: Describing novel object categories without paired training data. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. <https://doi.org/10.1109/CVPR.2016.8>
- Simao Herdade, Armin Kappeler, Kofi Boakye, and Joao Soares. 2019. Image captioning: Transforming objects into words. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc.
- Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. 2021. CLIP-Score: A reference-free evaluation metric for image captioning. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7514–7528, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2021.emnlp-main.595>
- Yusuke Hirota, Yuta Nakashima, and Noa Garcia. 2022. Quantifying societal bias amplification in image captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 13450–13459. <https://doi.org/10.1109/CVPR52688.2022.01309>
- Yusuke Hirota, Yuta Nakashima, and Noa Garcia. 2023. Model-agnostic gender debiased image captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 15191–15200. <https://doi.org/10.1109/CVPR52729.2023.01458>
- Micah Hodosh, Peter Young, and Julia Hockenmaier. 2013. Framing image description as a ranking task: Data, models and evaluation metrics. *Journal of Artificial Intelligence Research*, 47:853–899. <https://doi.org/10.1613/jair.3994>
- Ukyo Honda, Yoshitaka Ushiku, Atsushi Hashimoto, Taro Watanabe, and Yuji Matsumoto. 2021. Removing word-level spurious alignment between images and pseudo-captions in unsupervised image captioning. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 3692–3702, Online. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2021.eacl-main.323>
- MD Zakir Hossain, Ferdous Sohel, Mohd Fairuz Shiratuddin, and Hamid Laga. 2019. A comprehensive survey of deep learning for image captioning. *ACM Computing Surveys (CSUR)*, 51(6):1–36. <https://doi.org/10.1145/3295748>
- Mehrdad Hosseinzadeh and Yang Wang. 2021. Image change captioning by learning from an auxiliary task. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2725–2734. <https://doi.org/10.1109/CVPR46437.2021.00275>
- Jingyi Hou, Xinxiao Wu, Xiaoxun Zhang, Yayun Qi, Yunde Jia, and Jiebo Luo. 2020. Joint commonsense and relation reasoning for image and video captioning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(07):10973–10980. <https://doi.org/10.1609/aaai.v34i07.6731>
- Anwen Hu, Shizhe Chen, Liang Zhang, and Qin Jin. 2023a. InfoMetIC: An informative metric for reference-free image caption evaluation. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3171–3185, Toronto, Canada. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2023.acl-long.178>

- Xiaowei Hu, Zhe Gan, Jianfeng Wang, Zhengyuan Yang, Zicheng Liu, Yumao Lu, and Lijuan Wang. 2022. Scaling up vision-language pre-training for image captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 17980–17989.
- Xiaowei Hu, Xi Yin, Kevin Lin, Lei Zhang, Jianfeng Gao, Lijuan Wang, and Zicheng Liu. 2021. Vivo: Visual vocabulary pre-training for novel object captioning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(2):1575–1583. <https://doi.org/10.1609/aaai.v35i2.16249>
- Yushi Hu, Hang Hua, Zhengyuan Yang, Weijia Shi, Noah A. Smith, and Jiebo Luo. 2023b. PromptCap: Prompt-guided image captioning for VQA with GPT-3. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 2963–2975.
- Lun Huang, Wenmin Wang, Jie Chen, and Xiao-Yong Wei. 2019a. Attention on attention for image captioning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*. <https://doi.org/10.1109/ICCV.2019.00473>
- Lun Huang, Wenmin Wang, Yaxian Xia, and Jie Chen. 2019b. Adaptively aligned image captioning via adaptive attention time. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc.
- Qiuyuan Huang, Pengchuan Zhang, Dapeng Wu, and Lei Zhang. 2018. Turbo learning for captionbot and drawingbot. In *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc.
- Xiaoke Huang, Jianfeng Wang, Yansong Tang, Zheng Zhang, Han Hu, Jiwen Lu, Lijuan Wang, and Zicheng Liu. 2024. Segment and caption anything. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 13405–13417. <https://doi.org/10.1109/CVPR52733.2024.01273>
- EunJeong Hwang and Vered Shwartz. 2023. MemeCap: A dataset for captioning and interpreting memes. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 1433–1445, Singapore. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2023.emnlp-main.89>
- Mert Inan, Piyush Sharma, Baber Khalid, Radu Soricut, Matthew Stone, and Malihe Alikhani. 2021. COSMic: A coherence-aware generation metric for image descriptions. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 3419–3430, Punta Cana, Dominican Republic. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2021.findings-emnlp.291>
- Alejandro Jaimes and Shih-Fu Chang. 1999. Conceptual framework for indexing visual information at multiple levels. In *Internet Imaging*, volume 3964, pages 2–15. SPIE. <https://doi.org/10.1117/12.373443>
- Jiayi Ji, Yunpeng Luo, Xiaoshuai Sun, Fuhai Chen, Gen Luo, Yongjian Wu, Yue Gao, and Rongrong Ji. 2021. Improving image captioning by leveraging intra- and inter-layer global representation in transformer network. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(2):1655–1663. <https://doi.org/10.1609/aaai.v35i2.16258>
- Xu Jia, Efstratios Gavves, Basura Fernando, and Tinne Tuytelaars. 2015. Guiding the long-short term memory model for image caption generation. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*. <https://doi.org/10.1109/ICCV.2015.277>
- Ming Jiang, Junjie Hu, Qiuyuan Huang, Lei Zhang, Jana Diesner, and Jianfeng Gao. 2019a. REO-relevance, extraneous, omission: A fine-grained evaluation for image captioning. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 1475–1480, Hong Kong, China. Association for Computational Linguistics. <https://doi.org/10.18653/v1/D19-1156>
- Ming Jiang, Qiuyuan Huang, Lei Zhang, Xin Wang, Pengchuan Zhang, Zhe Gan, Jana Diesner, and Jianfeng Gao. 2019b. TIGER: Text-to-image grounding for image caption

- evaluation. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2141–2152, Hong Kong, China. Association for Computational Linguistics. <https://doi.org/10.18653/v1/D19-1220>
- Wenhao Jiang, Lin Ma, Xinpeng Chen, Hanwang Zhang, and Wei Liu. 2018a. Learning to guide decoding for image captioning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 32(1). <https://doi.org/10.1609/aaai.v32i1.12283>
- Wenhao Jiang, Lin Ma, Yu-Gang Jiang, Wei Liu, and Tong Zhang. 2018b. Recurrent fusion network for image captioning. In *Proceedings of the European Conference on Computer Vision (ECCV)*. https://doi.org/10.1007/978-3-030-01216-8_31
- Yang Jiao, Shaoxiang Chen, Zequn Jie, Jingjing Chen, Lin Ma, and Yu-Gang Jiang. 2022. MORE: multi-order relation mining for dense captioning in 3D scenes. In *Computer Vision – ECCV 2022*, pages 528–545, Cham. Springer Nature Switzerland. https://doi.org/10.1007/978-3-031-19833-5_31
- Justin Johnson, Andrej Karpathy, and Li Fei-Fei. 2016. DenseCap: Fully convolutional localization networks for dense captioning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. <https://doi.org/10.1109/CVPR.2016.494>
- Hassan Kane, Muhammed Yusuf Kocyigit, Ali Abdalla, Pelkins Ajano, and Mohamed Coulibali. 2020. NUBIA: NeUral based interchangeability assessor for text generation. In *Proceedings of the 1st Workshop on Evaluating NLG Evaluation*, pages 28–37, Online (Dublin, Ireland). Association for Computational Linguistics.
- Wooyoung Kang, Jonghwan Mun, Sungjun Lee, and Byungseok Roh. 2023. Noise-aware learning from web-crawled image-text data for image captioning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 2942–2952. <https://doi.org/10.1109/ICCV51070.2023.00275>
- Andrej Karpathy and Li Fei-Fei. 2015. Deep visual-semantic alignments for generating image descriptions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. <https://doi.org/10.1109/CVPR.2015.7298932>
- Jungo Kasai, Keisuke Sakaguchi, Lavinia Dunagan, Jacob Morrison, Ronan Le Bras, Yejin Choi, and Noah A. Smith. 2022. Transparent human evaluation for image captioning. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3464–3478, Seattle, United States. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2022.naacl-main.254>
- Lei Ke, Wenjie Pei, Ruiyu Li, Xiaoyong Shen, and Yu-Wing Tai. 2019. Reflective decoding network for image captioning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*. <https://doi.org/10.1109/ICCV.2019.00898>
- Mert Kilickaya, Aykut Erdem, Nazli Ikizler-Cinbis, and Erkut Erdem. 2017. Re-evaluating automatic metrics for image captioning. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers*, pages 199–209, Valencia, Spain. Association for Computational Linguistics. <https://doi.org/10.18653/v1/E17-1019>
- Dong-Jin Kim, Jinsoo Choi, Tae-Hyun Oh, and In So Kweon. 2019a. Dense relational captioning: Triple-stream networks for relationship-based captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Dong-Jin Kim, Jinsoo Choi, Tae-Hyun Oh, and In So Kweon. 2019b. Image captioning with very scarce supervised data: Adversarial semi-supervised learning approach. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2012–2023, Hong Kong, China. Association for Computational Linguistics. <https://doi.org/10.18653/v1/D19-1208>

- Hoeseong Kim, Jongseok Kim, Hyungseok Lee, Hyunsung Park, and Gunhee Kim. 2021a. Viewpoint-agnostic change captioning with cycle consistency. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 2095–2104.
- Hyoungun Kim, Abhay Zala, Graham Burri, and Mohit Bansal. 2021b. FIXMYPOSE: Pose correctional captioning and retrieval. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(14):13161–13170. <https://doi.org/10.1609/aaai.v35i14.17555>
- Yongil Kim, Yerin Hwang, Hyeongu Yun, Seunghyun Yoon, Trung Bui, and Kyomin Jung. 2023. PR-MCS: Perturbation robust metric for MultiLingual image captioning. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 12237–12258, Singapore. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2023.findings-emnlp.819>
- Simon Kornblith, Lala Li, Zirui Wang, and Thao Nguyen. 2023. Guiding image captioning models toward more specific captions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 15259–15269. <https://doi.org/10.1109/ICCV51070.2023.01400>
- Elisa Kreiss, Fei Fang, Noah Goodman, and Christopher Potts. 2022. Concadia: Towards image-based text generation with a purpose. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 4667–4684, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2022.emnlp-main.308>
- Chia-Wen Kuo and Zsolt Kira. 2022. Beyond a pre-trained object detector: Cross-modal textual and visual context for image captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 17969–17979.
- Chia-Wen Kuo and Zsolt Kira. 2023. HAAV: Hierarchical aggregation of augmented views for image captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11039–11049. <https://doi.org/10.1109/CVPR52729.2023.01062>
- Matt Kusner, Yu Sun, Nicholas Kolkin, and Kilian Weinberger. 2015. From word embeddings to document distances. In *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pages 957–966, Lille, France. PMLR.
- Polina Kuznetsova, Vicente Ordonez, Alexander Berg, Tamara Berg, and Yejin Choi. 2012. Collective generation of natural image descriptions. In *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 359–368, Jeju Island, Korea. Association for Computational Linguistics.
- Polina Kuznetsova, Vicente Ordonez, Alexander Berg, Tamara Berg, and Yejin Choi. 2013. Generalizing image captions for image-text parallel corpus. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 790–796, Sofia, Bulgaria. Association for Computational Linguistics.
- Polina Kuznetsova, Vicente Ordonez, Tamara L. Berg, and Yejin Choi. 2014. TreeTalk: Composition and compression of trees for image descriptions. *Transactions of the Association for Computational Linguistics*, 2:351–362. https://doi.org/10.1162/tacl_a_00188
- Iro Laina, Christian Rupprecht, and Nassir Navab. 2019. Towards unsupervised image captioning with shared multimodal embeddings. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*. <https://doi.org/10.1109/ICCV.2019.00751>
- Remi Lebret, Pedro Pinheiro, and Ronan Collobert. 2015. Phrase-based image captioning. In *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pages 2085–2094, Lille, France. PMLR.
- Hwanhee Lee, Thomas Scialom, Seunghyun Yoon, Franck Deroncourt, and Kyomin Jung. 2021a. QACE: Asking questions to evaluate an image caption. In *Findings of the*

- Association for Computational Linguistics: EMNLP 2021*, pages 4631–4638, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Hwanhee Lee, Seunghyun Yoon, Franck Dernoncourt, Trung Bui, and Kyomin Jung. 2021b. UMIC: An unreferenced metric for image captioning via contrastive learning. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 220–226, Online. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2021.acl-short.29>
- Hwanhee Lee, Seunghyun Yoon, Franck Dernoncourt, Doo Soon Kim, Trung Bui, and Kyomin Jung. 2020. ViLBERTScore: Evaluating image caption using vision-and-language BERT. In *Proceedings of the First Workshop on Evaluation and Comparison of NLP Systems*, pages 34–39, Online. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2020.eval4nlp-1.4>
- Dongxu Li, Junnan Li, Hung Le, Guangsen Wang, Silvio Savarese, and Steven C. H. Hoi. 2023a. LAVIS: A one-stop library for language-vision intelligence. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, pages 31–41, Toronto, Canada. Association for Computational Linguistics.
- Guang Li, Linchao Zhu, Ping Liu, and Yi Yang. 2019a. Entangled transformer for image captioning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*.
- Jiaxuan Li, Duc Minh Vo, Akihiro Sugimoto, and Hideki Nakayama. 2024a. EVCap: Retrieval-augmented image captioning with external visual-name memory for open-world comprehension. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 13733–13742. <https://doi.org/10.1109/CVPR52733.2024.01303>
- Linghui Li, Sheng Tang, Lixi Deng, Yongdong Zhang, and Qi Tian. 2017. Image caption with global-local attention. *Proceedings of the AAAI Conference on Artificial Intelligence*, 31(1). <https://doi.org/10.1609/aaai.v31i1.11236>
- Nannan Li, Zhenzhong Chen, and Shan Liu. 2019b. Meta learning for image captioning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 33(01):8626–8633. <https://doi.org/10.1609/aaai.v33i01.33018626>
- Runjia Li, Shuyang Sun, Mohamed Elhoseiny, and Philip Torr. 2023b. OxfordTVG-HIC: Can machine make humorous captions from images? In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 20293–20303.
- Siming Li, Girish Kulkarni, Tamara L. Berg, Alexander C. Berg, and Yejin Choi. 2011. Composing simple image descriptions using web-scale n-grams. In *Proceedings of the Fifteenth Conference on Computational Natural Language Learning*, pages 220–228. Portland, Oregon, USA. Association for Computational Linguistics.
- Wei Li, Linchao Zhu, Longyin Wen, and Yi Yang. 2023c. DeCap: Decoding clip latents for zero-shot captioning via text-only training.
- Wenyan Li, Jonas Lotz, Chen Qiu, and Desmond Elliott. 2024b. The role of data curation in image captioning. In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1074–1088, St. Julian’s, Malta. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2024.eacl-long.65>
- Xiangyang Li, Shuqiang Jiang, and Jungong Han. 2019c. Learning object context for dense captioning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 33(01):8650–8657. <https://doi.org/10.1609/aaai.v33i01.33018650>
- Yehao Li, Yingwei Pan, Ting Yao, and Tao Mei. 2022. Comprehending and ordering semantics for image captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 17990–17999. <https://doi.org/10.1109/CVPR52688.2022.01746>

- Yehao Li, Ting Yao, Yingwei Pan, Hongyang Chao, and Tao Mei. 2019d. Pointing novel objects in image captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81. Barcelona, Spain. Association for Computational Linguistics.
- Bing Liu, Dong Wang, Xu Yang, Yong Zhou, Rui Yao, Zhiwen Shao, and Jiaqi Zhao. 2022. Show, deconfound and tell: Image captioning with causal inference. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 18041–18050. <https://doi.org/10.1109/CVPR52688.2022.01751>
- Chenxi Liu, Junhua Mao, Fei Sha, and Alan Yuille. 2017a. Attention correctness in neural image captioning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 31(1). <https://doi.org/10.1609/aaai.v31i1.11197>
- Fenglin Liu, Xuancheng Ren, Yuanxin Liu, Houfeng Wang, and Xu Sun. 2018a. simNet: Stepwise image-topic merging network for generating detailed and comprehensive image captions. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 137–149, Brussels, Belgium. Association for Computational Linguistics.
- Fuxiao Liu, Yinghan Wang, Tianlu Wang, and Vicente Ordonez. 2021. Visual news: Benchmark and challenges in news image captioning. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6761–6771, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2021.emnlp-main.542>
- Junhao Liu, Kai Wang, Chunpu Xu, Zhou Zhao, Ruifeng Xu, Ying Shen, and Min Yang. 2020. Interactive dual generative adversarial networks for image captioning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(07):11588–11595. <https://doi.org/10.1609/aaai.v34i07.6826>
- Lixin Liu, Jiajun Tang, Xiaojun Wan, and Zongming Guo. 2019a. Generating diverse and descriptive image captions using visual paraphrases. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*. <https://doi.org/10.1109/ICCV.2019.00434>
- Siqi Liu, Zhenhai Zhu, Ning Ye, Sergio Guadarrama, and Kevin Murphy. 2017b. Improved image captioning via policy gradient optimization of spider. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*. <https://doi.org/10.1109/ICCV.2017.100>
- Xiaoxiao Liu, Qingyang Xu, and Ning Wang. 2019b. A survey on deep neural network-based image captioning. *The Visual Computer*, 35(3):445–470. <https://doi.org/10.1007/s00371-018-1566-y>
- Xihui Liu, Hongsheng Li, Jing Shao, Dapeng Chen, and Xiaogang Wang. 2018b. Show, tell and discriminate: Image captioning by self-retrieval with partially labeled data. In *Proceedings of the European Conference on Computer Vision (ECCV)*. https://doi.org/10.1007/978-3-030-01267-0_21
- Zhiyue Liu, Jinyuan Liu, and Fanrong Ma. 2024. Improving cross-modal alignment with synthetic pairs for text-only image captioning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(4):3864–3872. <https://doi.org/10.1609/aaai.v38i4.28178>
- Di Lu, Spencer Whitehead, Lifu Huang, Heng Ji, and Shih-Fu Chang. 2018. Entity-aware image caption generation. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4013–4023, Brussels, Belgium. Association for Computational Linguistics. <https://doi.org/10.18653/v1/D18-1435>
- Jiasen Lu, Caiming Xiong, Devi Parikh, and Richard Socher. 2017. Knowing when to look: Adaptive attention via a visual sentinel for image captioning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. <https://doi.org/10.1109/CVPR.2017.345>
- Jianjie Luo, Yehao Li, Yingwei Pan, Ting Yao, Jianlin Feng, Hongyang Chao, and

- Tao Mei. 2023a. Semantic-conditional diffusion networks for image captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 23359–23368.
- Ruotian Luo, Brian Price, Scott Cohen, and Gregory Shakhnarovich. 2018. Discriminability objective for training descriptive captions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. <https://doi.org/10.1109/CVPR.2018.00728>
- Tiang Luo, Chris Rockwell, Honglak Lee, and Justin Johnson. 2023b. Scalable 3D captioning with pretrained models. In *Advances in Neural Information Processing Systems*, volume 36, pages 75307–75337. Curran Associates, Inc.
- Yunpeng Luo, Jiayi Ji, Xiaoshuai Sun, Liujuan Cao, Yongjian Wu, Feiyue Huang, Chia-Wen Lin, and Rongrong Ji. 2021. Dual-level collaborative transformer for image captioning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(3): 2286–2293. <https://doi.org/10.1609/aaai.v35i3.16328>
- Chih-Yao Ma, Yannis Kalantidis, Ghassan AlRegib, Peter Vajda, Marcus Rohrbach, and Zsolt Kira. 2020. Learning to generate grounded visual captions without localization supervision. In *Computer Vision – ECCV 2020*, pages 353–370, Cham. Springer International Publishing. https://doi.org/10.1007/978-3-030-58523-5_21
- Feipeng Ma, Yizhou Zhou, Fengyun Rao, Yueyi Zhang, and Xiaoyan Sun. 2024a. Image captioning with multi-context synthetic data. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(5):4089–4097. <https://doi.org/10.1609/aaai.v38i5.28203>
- Ziping Ma, Furong Xu, Jian Liu, Ming Yang, and Qingpei Guo. 2024b. SyCoCa: Symmetrizing contrastive captioners with attentive masking for multimodal alignment. *arXiv preprint arXiv:2401.02137*.
- Pranava Madhyastha, Josiah Wang, and Lucia Specia. 2019. VIFIDEL: Evaluating the visual fidelity of image descriptions. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 6539–6550, Florence, Italy. Association for Computational Linguistics. <https://doi.org/10.18653/v1/P19-1654>
- Pranava Swaroop Madhyastha, Josiah Wang, and Lucia Specia. 2018. End-to-end image captioning exploits distributional similarity in multimodal space. In *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 381–383, Brussels, Belgium. Association for Computational Linguistics. <https://doi.org/10.18653/v1/W18-5455>
- Koki Maeda, Shuhei Kurita, Taiki Miyanishi, and Naoaki Okazaki. 2023. Query-based image captioning from multi-context 360degree images. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 6940–6954, Singapore. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2023.findings-emnlp.463>
- Shweta Mahajan and Stefan Roth. 2020. Diverse image captioning with context-object split latent spaces. In *Advances in Neural Information Processing Systems*, volume 33, pages 3613–3624, Curran Associates, Inc.
- Junhua Mao, Wei Xu, Yi Yang, Jiang Wang, Zhiheng Huang, and Alan Yuille. 2015. Deep captioning with multimodal recurrent neural networks (M-RNN).
- Rebecca Mason and Eugene Charniak. 2014a. Domain-specific image captioning. In *Proceedings of the Eighteenth Conference on Computational Natural Language Learning*, pages 11–20, Ann Arbor, Michigan. Association for Computational Linguistics. <https://doi.org/10.3115/v1/W14-1602>
- Rebecca Mason and Eugene Charniak. 2014b. Nonparametric method for data-driven image captioning. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 592–598, Baltimore, Maryland. Association for Computational Linguistics. <https://doi.org/10.3115/v1/P14-2097>
- Alexander Mathews, Lexing Xie, and Xuming He. 2016. Senticap: Generating image descriptions with sentiments. *Proceedings of the AAAI Conference on Artificial Intelligence*, 30(1). <https://doi.org/10.1609/aaai.v30i1.10475>

- Alexander Mathews, Lexing Xie, and Xuming He. 2018. SemStyle: Learning to generate stylised image captions using unaligned text. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. <https://doi.org/10.1109/CVPR.2018.00896>
- Effrosyni Mavroudi and René Vidal. 2022. Weakly-supervised generation and grounding of visual descriptions with conditional generative models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 15544–15554. <https://doi.org/10.1109/CVPR52688.2022.01510>
- Luke Melas-Kyriazi, Alexander Rush, and George Han. 2018. Training for diversity in image paragraph captioning. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 757–761, Brussels, Belgium. Association for Computational Linguistics. <https://doi.org/10.18653/v1/D18-1084>
- Zihang Meng, David Yang, Xuefei Cao, Ashish Shah, and Ser-Nam Lim. 2022. Object-centric unsupervised image captioning. In *Computer Vision – ECCV 2022*, pages 219–235, Cham. Springer Nature Switzerland. https://doi.org/10.1007/978-3-031-20059-5_13
- Victor Milewski, Marie-Francine Moens, and Iacer Calixto. 2020. Are scene graphs good enough to improve image captioning? In *Proceedings of the 1st Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 10th International Joint Conference on Natural Language Processing*, pages 504–515, Suzhou, China. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2020.aacl-main.50>
- Suvir Mirchandani, Licheng Yu, Mengjiao Wang, Animesh Sinha, Wenwen Jiang, Tao Xiang, and Ning Zhang. 2022. FaD-VLP: Fashion vision-and-language pre-training towards unified retrieval and captioning. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 10484–10497, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2022.emnlp-main.716>
- Margaret Mitchell, Jesse Dodge, Amit Goyal, Kota Yamaguchi, Karl Stratos, Xufeng Han, Alyssa Mensch, Alex Berg, Tamara Berg, and Hal Daumé III. 2012. Midge: Generating image descriptions from computer vision detections. In *Proceedings of the 13th Conference of the European Chapter of the Association for Computational Linguistics*, pages 747–756, Avignon, France. Association for Computational Linguistics.
- Abdelrahman Mohamed, Fakhraddin Alwajih, El Moatez Billah Nagoudi, Alcides Inciarte, and Muhammad Abdul-Mageed. 2023. Violet: A vision-language model for Arabic image captioning with gemini decoder. In *Proceedings of ArabicNLP 2023*, pages 1–11, Singapore (Hybrid). Association for Computational Linguistics. <https://doi.org/10.18653/v1/2023.arabicnlp-1.1>
- Youssef Mohamed, Faizan Farooq Khan, Kilichbek Haydarov, and Mohamed Elhoseiny. 2022. It is okay to not be okay: Overcoming emotional bias in affective image captioning by contrastive data collection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 21263–21272. <https://doi.org/10.1109/CVPR52688.2022.02058>
- Jonghwan Mun, Minsu Cho, and Bohyung Han. 2017. Text-guided attention model for image captioning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 31(1). <https://doi.org/10.1609/aaai.v31i1.11237>
- Sreyasi Nag Chowdhury, Rajarshi Bhowmik, Hareesh Ravi, Gerard de Melo, Simon Razniewski, and Gerhard Weikum. 2021. Exploiting image–text synergy for contextual image captioning. In *Proceedings of the Third Workshop on Beyond Vision and Language: Integrating Real-world Knowledge (LANTErn)*, pages 30–37, Kyiv, Ukraine. Association for Computational Linguistics.
- Shunta Nagasawa, Yotaro Watanabe, and Hitoshi Iyatomi. 2021. Validity-based sampling and smoothing methods for multiple reference image captioning. In *Proceedings*

- of the *Third Workshop on Multimodal Artificial Intelligence*, pages 36–41, Mexico City, Mexico. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2021.maiworkshop-1.6>
- Sai Suprabhanu Nallapaneni and Subrahmanyam Konakanchi. 2023. A comprehensive analysis of real-world image captioning and scene identification. *Journal of Electrical Electronics Engineering*, 2(3). <https://doi.org/10.33140/JEEE.02.03.14>
- Edwin G. Ng, Bo Pang, Piyush Sharma, and Radu Soricut. 2021. Understanding guided image captioning performance across domains. In *Proceedings of the 25th Conference on Computational Natural Language Learning*, pages 183–193, Online. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2021.conll-1.14>
- Khanh Nguyen, Ali Furkan Biten, Andres Mafla, Lluís Gomez, and Dimosthenis Karatzas. 2023. Show, interpret and tell: Entity-aware contextualised image captioning in wikipedia. *Proceedings of the AAAI Conference on Artificial Intelligence*, 37(2):1940–1948. <https://doi.org/10.1609/aaai.v37i2.25285>
- Kien Nguyen, Subarna Tripathi, Bang Du, Tanaya Guha, and Truong Q. Nguyen. 2021. In defense of scene graphs for image captioning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 1407–1416. <https://doi.org/10.1109/ICCV48922.2021.00144>
- Minh Thu Nguyen, Duy Phung, Minh Hoai, and Thien Huu Nguyen. 2020. Structural and functional decomposition for personality image captioning in a communication game. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 4587–4593, Online. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2020.findings-emnlp.411>
- Van-Quang Nguyen, Masanori Suganuma, and Takayuki Okatani. 2022. GRIT: Faster and better image captioning transformer using dual visual features. In *Computer Vision – ECCV 2022*, pages 167–184, Cham. Springer Nature Switzerland. https://doi.org/10.1007/978-3-031-20059-5_10
- Robert C. Nickerson, Upkar Varshney, and Jan Muntermann. 2013. A method for taxonomy development and its application in information systems. *European Journal of Information Systems*, 22(3):336–359. <https://doi.org/10.1057/ejis.2012.26>
- Allen Nie, Reuben Cohn-Gordon, and Christopher Potts. 2020. Pragmatic issue-sensitive image captioning. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1924–1938, Online. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2020.findings-emnlp.173>
- Mitja Nikolaus, Mostafa Abdou, Matthew Lamm, Rahul Aralikkatte, and Desmond Elliott. 2019. Compositional generalization in image captioning. In *Proceedings of the 23rd Conference on Computational Natural Language Learning (CoNLL)*, pages 87–98, Hong Kong, China. Association for Computational Linguistics. <https://doi.org/10.18653/v1/K19-1009>
- Toru Nishino, Yasuhide Miura, Tomoki Taniguchi, Tomoko Ohkuma, Yuki Suzuki, Shoji Kido, and Noriyuki Tomiyama. 2022. Factual accuracy is not enough: Planning consistent description order for radiology report generation. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 7123–7138, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2022.emnlp-main.480>
- David Nukrai, Ron Mokady, and Amir Globerson. 2022. Text-only training for image captioning using noise-injected CLIP. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 4055–4063, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2022.findings-emnlp.299>
- Gabriel Oliveira dos Santos, Esther Luna Colombini, and Sandra Avila. 2021. CIDER-R: Robust consensus-based image description evaluation. In *Proceedings of the Seventh Workshop on Noisy User-generated Text (W-NUT 2021)*, pages 351–360, Online. Association for

- Computational Linguistics. <https://doi.org/10.18653/v1/2021.wnut-1.39>
- Vicente Ordonez, Girish Kulkarni, and Tamara Berg. 2011. Im2Text: Describing images using 1 million captioned photographs. In *Advances in Neural Information Processing Systems*, volume 24. Curran Associates, Inc.
- Jiefu Ou, Benno Krojer, and Daniel Fried. 2023. Pragmatic inference with a CLIP listener for contrastive captioning. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 1904–1917, Toronto, Canada. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2023.findings-acl.120>
- Yingwei Pan, Ting Yao, Yehao Li, and Tao Mei. 2020. X-linear attention networks for image captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. <https://doi.org/10.1109/CVPR42600.2020.01098>
- George Pantazopoulos, Alessandro Suglia, and Arash Eshghi. 2022. Combine to describe: Evaluating compositional generalization in image captioning. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics: Student Research Workshop*, pages 115–131, Dublin, Ireland. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2022.acl-srw.11>
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: A method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics. <https://doi.org/10.3115/1073083.1073135>
- Dong Huk Park, Trevor Darrell, and Anna Rohrbach. 2019. Robust change captioning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*. <https://doi.org/10.1109/ICCV.2019.00472>
- Hyeryun Park, Kyungmo Kim, Jooyoung Yoon, Seongkeun Park, and Jinwook Choi. 2020. Feature difference makes sense: A medical image captioning model exploiting feature difference and tag information. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: Student Research Workshop*, pages 95–102, Online. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2020.acl-srw.14>
- Marco Pedersoli, Thomas Lucas, Cordelia Schmid, and Jakob Verbeek. 2017. Areas of attention for image captioning. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*. <https://doi.org/10.1109/ICCV.2017.140>
- Suzanne Petryk, David Chan, Anish Kachinthaya, Haodi Zou, John Canny, Joseph Gonzalez, and Trevor Darrell. 2024. ALOHa: A new measure for hallucination in captioning models. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 2: Short Papers)*, pages 342–357, Mexico City, Mexico. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2024.naacl-short.30>
- Yunchen Pu, Zhe Gan, Ricardo Henao, Xin Yuan, Chunyuan Li, Andrew Stevens, and Lawrence Carin. 2016. Variational autoencoder for deep learning of images, labels and captions. In *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc.
- Yayun Qi, Wentian Zhao, and Xinxiao Wu. 2024. Relational distant supervision for image captioning without image-text pairs. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(5):4524–4532. <https://doi.org/10.1609/aaai.v38i5.28251>
- Yu Qin, Jiajun Du, Yonghua Zhang, and Hongtao Lu. 2019. Look back and predict forward in image captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. <https://doi.org/10.1109/CVPR.2019.00856>
- Haoyi Qiu, Zi-Yi Dou, Tianlu Wang, Asli Celikyilmaz, and Nanyun Peng. 2023. Gender biases in automatic evaluation metrics for image captioning. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 8358–8375,

- Singapore. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2023.emnlp-main.520>
- Longtian Qiu, Shan Ning, and Xuming He. 2024. Mining fine-grained image-text alignment for zero-shot captioning via text-only training. *Proceedings of the AAI Conference on Artificial Intelligence*, 38(5):4605–4613. <https://doi.org/10.1609/aaai.v38i5.28260>
- Tingyu Qu, Tinne Tuytelaars, and Marie-Francine Moens. 2024. Visually-aware context modeling for news image captioning. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 2927–2943, Mexico City, Mexico. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2024.naacl-long.162>
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.
- Abisek Rajakumar Kalarani, Pushpak Bhattacharyya, Niyati Chhaya, and Sumit Shekhar. 2023. ‘Let’s not quote out of context’: Unified vision-language pretraining for context assisted image captioning. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 5: Industry Track)*, pages 695–706, Toronto, Canada. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2023.acl-industry.67>
- Rita Ramos, Emanuele Bugliarello, Bruno Martins, and Desmond Elliott. 2024. PAELLA: Parameter-efficient lightweight language-agnostic captioning model. In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 3549–3564, Mexico City, Mexico. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2024.findings-naacl.225>
- Rita Ramos, Desmond Elliott, and Bruno Martins. 2023a. Retrieval-augmented image captioning. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 3666–3681, Dubrovnik, Croatia. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2023.eacl-main.266>
- Rita Ramos, Bruno Martins, and Desmond Elliott. 2023b. LMCap: Few-shot multilingual image captioning by retrieval augmented language model prompting. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 1635–1651, Toronto, Canada. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2023.findings-acl.104>
- Rita Ramos, Bruno Martins, Desmond Elliott, and Yova Kementchedjheva. 2023c. Small-Cap: Lightweight image captioning prompted with retrieval augmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2840–2849. <https://doi.org/10.1109/CVPR52729.2023.00278>
- Cyrus Rashtchian, Peter Young, Micah Hodosh, and Julia Hockenmaier. 2010. Collecting image annotations using Amazon’s Mechanical Turk. In *Proceedings of the NAACL HLT 2010 Workshop on Creating Speech and Language Data with Amazon’s Mechanical Turk*, pages 139–147, Los Angeles. Association for Computational Linguistics.
- Yuchen Ren, Zhendong Mao, Shancheng Fang, Yan Lu, Tong He, Hao Du, Yongdong Zhang, and Wanli Ouyang. 2023. Crossing the gap: Domain generalization for image captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2871–2880. <https://doi.org/10.1109/CVPR52729.2023.00281>
- Zhou Ren, Xiaoyu Wang, Ning Zhang, Xutao Lv, and Li-Jia Li. 2017. Deep reinforcement learning-based image captioning with embedding reward. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. <https://doi.org/10.1109/CVPR.2017.128>
- Steven J. Rennie, Etienne Marcheret, Youssef Mroueh, Jerret Ross, and Vaibhava Goel. 2017. Self-critical sequence training for image captioning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern*

- Recognition (CVPR)*. <https://doi.org/10.1109/CVPR.2017.131>
- Anna Rohrbach, Lisa Anne Hendricks, Kaylee Burns, Trevor Darrell, and Kate Saenko. 2018. Object hallucination in image captioning. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4035–4045, Brussels, Belgium. Association for Computational Linguistics. <https://doi.org/10.18653/v1/D18-1437>
- Anna Rohrbach, Marcus Rohrbach, Siyu Tang, Seong Joon Oh, and Bernt Schiele. 2017. Generating descriptions with grounded and co-referenced people. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. <https://doi.org/10.1109/CVPR.2017.447>
- Dan Ruta, Andrew Gilbert, Pranav Aggarwal, Naveen Marri, Ajinkya Kale, Jo Briggs, Chris Speed, Hailin Jin, Baldo Faieta, Alex Filipkowski, Zhe Lin, and John Collomosse. 2022. StyleBabel: Artistic style tagging and captioning. In *Computer Vision – ECCV 2022*, pages 219–236, Cham. Springer Nature Switzerland. https://doi.org/10.1007/978-3-031-20074-8_13
- Fawaz Sammani and Luke Melas-Kyriazi. 2020. Show, edit and tell: A framework for editing image captions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. <https://doi.org/10.1109/CVPR42600.2020.00486>
- Sara Sarto, Manuele Barraco, Marcella Cornia, Lorenzo Baraldi, and Rita Cucchiara. 2023. Positive-augmented contrastive learning for image and video captioning evaluation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6914–6924. <https://doi.org/10.1109/CVPR52729.2023.00668>
- Naeha Sharif, Lyndon White, Mohammed Bennamoun, and Syed Afaq Ali Shah. 2018. NNEval: Neural network based evaluation metric for image captioning. In *Proceedings of the European Conference on Computer Vision (ECCV)*. https://doi.org/10.1007/978-3-030-01237-3_3
- Himanshu Sharma. 2021. A survey on image captioning datasets and evaluation metrics. In *IOP Conference Series: Materials Science and Engineering*, volume 1116, page 012184. IOP Publishing. <https://doi.org/10.1088/1757-899X/1116/1/012184>
- Himanshu Sharma and Devanand Padha. 2023. A comprehensive survey on image captioning: From handcrafted to deep learning-based techniques, a taxonomy and open research issues. *Artificial Intelligence Review*, 56(11):13619–13661. <https://doi.org/10.1145/3492865>
- Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. 2018. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2556–2565, Melbourne, Australia. Association for Computational Linguistics. <https://doi.org/10.18653/v1/P18-1238>
- Sara Shatford. 1986. Analyzing the subject of a picture: A theoretical approach. *Cataloging & Classification Quarterly*, 6(3):39–62. <https://doi.org/10.1300/J104v06n03.04>
- Tingke Shen, Amlan Kar, and Sanja Fidler. 2019. Learning to caption images through a lifetime by asking questions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*. <https://doi.org/10.1109/ICCV.2019.01049>
- Rakshith Shetty, Marcus Rohrbach, Lisa Anne Hendricks, Mario Fritz, and Bernt Schiele. 2017. Speaking the same language: Matching machine to human captions by adversarial training. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*. <https://doi.org/10.1109/ICCV.2017.445>
- Jiahe Shi, Yali Li, and Shengjin Wang. 2021a. Partial off-policy learning: Balance accuracy and diversity for human-oriented image captioning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 2187–2196. <https://doi.org/10.1109/ICCV48922.2021.00219>
- Xiangxi Shi, Xu Yang, Jiuxiang Gu, Shafiq Joty, and Jianfei Cai. 2020a. Finding it at

- another side: A viewpoint-adapted matching encoder for change captioning. In *Computer Vision – ECCV 2020*, pages 574–590, Cham. Springer International Publishing. https://doi.org/10.1007/978-3-030-58568-6_34
- Zhan Shi, Hui Liu, Martin Renqiang Min, Christopher Malon, Li Erran Li, and Xiaodan Zhu. 2021b. Retrieval, analogy, and composition: A framework for compositional generalization in image captioning. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 1990–2000, Punta Cana, Dominican Republic. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2021.findings-emnlp.171>
- Zhan Shi, Hui Liu, and Xiaodan Zhu. 2021c. Enhancing descriptive image captioning with natural language inference. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 269–277, Online. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2021.acl-short.36>
- Zhan Shi, Xu Zhou, Xipeng Qiu, and Xiaodan Zhu. 2020b. Improving image captioning with better use of caption. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7454–7464, Online. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2020.acl-main.664>
- Kurt Shuster, Samuel Humeau, Hexiang Hu, Antoine Bordes, and Jason Weston. 2019. Engaging image captioning via personality. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. <https://doi.org/10.1109/CVPR.2019.01280>
- Oleksii Sidorov, Ronghang Hu, Marcus Rohrbach, and Amanpreet Singh. 2020. TextCaps: A dataset for image captioning with reading comprehension. In *Computer Vision – ECCV 2020*, pages 742–758, Cham. Springer International Publishing. https://doi.org/10.1007/978-3-030-58536-5_44
- Matthew Snover, Bonnie Dorr, Rich Schwartz, Linnea Micciulla, and John Makhoul. 2006. A study of translation edit rate with targeted human annotation. In *Proceedings of the 7th Conference of the Association for Machine Translation in the Americas: Technical Papers*, pages 223–231, Cambridge, Massachusetts, USA. Association for Machine Translation in the Americas.
- Lingyun Song, Jun Liu, Buyue Qian, and Yihe Chen. 2019. Connecting language to images: A progressive attention-guided network for simultaneous image captioning and language grounding. *Proceedings of the AAAI Conference on Artificial Intelligence*, 33(01):8885–8892. <https://doi.org/10.1609/aaai.v33i01.33018885>
- Zeliang Song, Xiaofei Zhou, Zhendong Mao, and Jianlong Tan. 2021. Image captioning with context-aware auxiliary guidance. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(3):2584–2592. <https://doi.org/10.1609/aaai.v35i3.16361>
- Raimonda Staniūtė and Dmitrij Šešok. 2019. A systematic literature review on image captioning. *Applied Sciences*, 9(10):2024. <https://doi.org/10.3390/app9102024>
- Matteo Stefanini, Marcella Cornia, Lorenzo Baraldi, Silvia Cascianelli, Giuseppe Fiameni, and Rita Cucchiara. 2022. From show to tell: A survey on deep learning-based image captioning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(1):539–559. <https://doi.org/10.1109/TPAMI.2022.3148210>, PubMed: 35130142
- Qing Sun, Stefan Lee, and Dhruv Batra. 2017. Bidirectional beam search: Forward-backward inference in neural sequence models for fill-in-the-blank image captioning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. <https://doi.org/10.1109/CVPR.2017.763>
- Ece Takmaz, Sandro Pezzelle, Lisa Beinborn, and Raquel Fernández. 2020. Generating image descriptions via sequential cross-modal alignment guided by human gaze. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*,

- pages 4664–4677, Online. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2020.emnlp-main.377>
- Kohtaro Tanaka, Kohei Uehara, Lin Gu, Yusuke Mukuta, and Tatsuya Harada. 2024. Content-specific humorous image captioning using incongruity resolution chain-of-thought. In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 2348–2367, Mexico City, Mexico. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2024.findings-naacl.152>
- Hamed R. Tavakoli, Rakshith Shetty, Ali Borji, and Jorma Laaksonen. 2017. Paying attention to descriptions generated by image captioning models. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*. <https://doi.org/10.1109/ICCV.2017.272>
- Ashish V. Thapliyal, Jordi Pont Tuset, Xi Chen, and Radu Soricut. 2022. Crossmodal-3600: A massively multilingual multimodal evaluation dataset. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 715–729, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2022.emnlp-main.45>
- Alasdair Tran, Alexander Mathews, and Lexing Xie. 2020. Transform and tell: Entity-aware news image captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. <https://doi.org/10.1109/CVPR42600.2020.01305>
- Yunbin Tu, Liang Li, Li Su, Zheng-Jun Zha, Chenggang Yan, and Qingming Huang. 2023. Self-supervised cross-view representation reconstruction for change captioning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 2805–2815. <https://doi.org/10.1109/ICCV51070.2023.00263>
- Yunbin Tu, Liang Li, Chenggang Yan, Shengxiang Gao, and Zhengtao Yu. 2021a. R³Net:relation-embedded representation reconstruction network for change captioning. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 9319–9329, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2021.emnlp-main.735>
- Yunbin Tu, Tingting Yao, Liang Li, Jiedong Lou, Shengxiang Gao, Zhengtao Yu, and Chenggang Yan. 2021b. Semantic relation-aware difference representation learning for change captioning. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 63–73, Online. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2021.findings-acl.6>
- Yoshitaka Ushiku, Masataka Yamaguchi, Yusuke Mukuta, and Tatsuya Harada. 2015. Common subspace for model and similarity: Phrase learning for caption generation from images. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*. <https://doi.org/10.1109/ICCV.2015.306>
- Ramakrishna Vedantam, Samy Bengio, Kevin Murphy, Devi Parikh, and Gal Chechik. 2017. Context-aware captions from context-agnostic supervision. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. <https://doi.org/10.1109/CVPR.2017.120>
- Ramakrishna Vedantam, C. Lawrence Zitnick, and Devi Parikh. 2015. Cider: Consensus-based image description evaluation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4566–4575. <https://doi.org/10.1109/CVPR.2015.7299087>
- Subhashini Venugopalan, Lisa Anne Hendricks, Marcus Rohrbach, Raymond Mooney, Trevor Darrell, and Kate Saenko. 2017. Captioning images with diverse objects. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. <https://doi.org/10.1109/CVPR.2017.130>
- Gilad Vered, Gal Oren, Yuval Atzmon, and Gal Chechik. 2019. Joint optimization for cooperative image captioning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*. <https://doi.org/10.1109/ICCV.2019.00899>
- Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. 2015. Show and tell: A neural image caption generator. In *Proceedings of*

- the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. <https://doi.org/10.1109/CVPR.2015.7298935>
- Duc Minh Vo, Hong Chen, Akihiro Sugimoto, and Hideki Nakayama. 2022. Noc-rek: Novel object captioning with retrieved vocabulary from external knowledge. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 18000–18008.
- Duc Minh Vo, Quoc-An Luong, Akihiro Sugimoto, and Hideki Nakayama. 2023. A-cap: Anticipation captioning with commonsense knowledge. In *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10824–10833. <https://doi.org/10.1109/CVPR52729.2023.01042>
- Yuiga Wada, Kanta Kaneda, Daichi Saito, and Komei Sugiura. 2024. Polos: Multimodal metric learning from human feedback for image captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 13559–13568. <https://doi.org/10.1109/CVPR52733.2024.01287>
- Dalin Wang, Daniel Beck, and Trevor Cohn. 2019a. On the role of scene graphs in image captioning. In *Proceedings of the Beyond Vision and LAnguage: inTEgrating Real-world kNowledge (LANTERN)*, pages 29–34, Hong Kong, China. Association for Computational Linguistics. <https://doi.org/10.18653/v1/D19-6405>
- Jing Wang, Jinhui Tang, Mingkun Yang, Xiang Bai, and Jiebo Luo. 2021a. Improving ocr-based image captioning by incorporating geometrical relationship. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1306–1315. <https://doi.org/10.1109/CVPR46437.2021.00136>
- Jiuniu Wang, Wenjia Xu, Qingzhong Wang, and Antoni B. Chan. 2020a. Compare and reweight: Distinctive image captioning using similar images sets. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part I 16*, pages 370–386. Springer. https://doi.org/10.1007/978-3-030-58452-8_22
- Josiah Wang, Pranava Swaroop Madhyastha, and Lucia Specia. 2018. Object counts! Bringing explicit detections back into image captioning. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 2180–2193, New Orleans, Louisiana. Association for Computational Linguistics. <https://doi.org/10.18653/v1/N18-1198>
- Li Wang, Zechen Bai, Yonghua Zhang, and Hongtao Lu. 2020b. Show, recall, and tell: Image captioning with recall mechanism. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(07):12176–12183. <https://doi.org/10.1609/aaai.v34i07.6898>
- Liwei Wang, Alexander Schwing, and Svetlana Lazebnik. 2017a. Diverse and accurate image description using a variational auto-encoder with an additive gaussian encoding space. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.
- Ning Wang, Jiajun Deng, and Mingbo Jia. 2024. Cycle-consistency learning for captioning and grounding. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(6):5535–5543. <https://doi.org/10.1609/aaai.v38i6.28363>
- Ning Wang, Jiangrong Xie, Hang Luo, Qinglin Cheng, Jihao Wu, Mingbo Jia, and Linlin Li. 2023a. Efficient image captioning for edge devices. *Proceedings of the AAAI Conference on Artificial Intelligence*, 37(2):2608–2616. <https://doi.org/10.1609/aaai.v37i2.25359>
- Ning Wang, Jiahao Xie, Jihao Wu, Mingbo Jia, and Linlin Li. 2023b. Controllable image captioning via prompting. *Proceedings of the AAAI Conference on Artificial Intelligence*, 37(2):2617–2625. <https://doi.org/10.1609/aaai.v37i2.25360>
- Qingzhong Wang and Antoni B. Chan. 2019. Describing like humans: On diversity in image captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. <https://doi.org/10.1109/CVPR.2019.00432>
- Sijin Wang, Ziwei Yao, Ruiping Wang, Zhongqin Wu, and Xilin Chen. 2021b. FAIEr:

- Fidelity and adequacy ensured image caption evaluation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 14050–14059. <https://doi.org/10.1109/CVPR46437.2021.01383>
- Ting Wang, Weidong Chen, Yuanhe Tian, Yan Song, and Zhendong Mao. 2023c. Improving image captioning via predicting structured concepts. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 360–370, Singapore. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2023.emnlp-main.25>
- Weixuan Wang, Zhihong Chen, and Haifeng Hu. 2019b. Hierarchical attention network for image captioning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 33(01):8957–8964. <https://doi.org/10.1609/aaai.v33i01.33018957>
- Yiyu Wang, Jungang Xu, and Yingfei Sun. 2022b. End-to-end transformer based model for image captioning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 36(3):2585–2594. <https://doi.org/10.1609/aaai.v36i3.20160>
- Yufei Wang, Zhe Lin, Xiaohui Shen, Scott Cohen, and Garrison W. Cottrell. 2017b. Skeleton key: Image captioning by skeleton-attribute decomposition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. <https://doi.org/10.1109/CVPR.2017.780>
- Yufei Wang, Ian Wood, Stephen Wan, and Mark Johnson. 2021c. ECOL-R: Encouraging copying in novel object captioning with reinforcement learning. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 1222–1234, Online. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2021.eacl-main.104>
- Yuxuan Wang, Difei Gao, Licheng Yu, Weixian Lei, Matt Feiszli, and Mike Zheng Shou. 2022a. Geb+: A benchmark for generic event boundary captioning, grounding and retrieval. In *Computer Vision – ECCV 2022*, pages 709–725, Cham. Springer Nature Switzerland. https://doi.org/10.1007/978-3-031-19833-5_41
- Zeyu Wang, Berthy Feng, Karthik Narasimhan, and Olga Russakovsky. 2020c. Towards unique and informative captioning of images. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part VII 16*, pages 629–644. Springer. https://doi.org/10.1007/978-3-030-58571-6_37
- Zhen Wang, Long Chen, Wenbo Ma, Guangxing Han, Yulei Niu, Jian Shao, and Jun Xiao. 2022c. Explicit image caption editing. In *Computer Vision – ECCV 2022*, pages 113–129, Cham. Springer Nature Switzerland. https://doi.org/10.1007/978-3-031-20059-5_7
- Jianzong Wu, Xiangtai Li, Henghui Ding, Xia Li, Guangliang Cheng, Yunhai Tong, and Chen Change Loy. 2023a. Betrayed by captions: Joint caption grounding and generation for open vocabulary instance segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 21938–21948.
- Mingrui Wu, Xuying Zhang, Xiaoshuai Sun, Yiyi Zhou, Chao Chen, Jiabin Gu, Xing Sun, and Rongrong Ji. 2022. DIFNet: Boosting visual information flow for image captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 18020–18029. <https://doi.org/10.1109/CVPR52688.2022.01749>
- Shengqiong Wu, Hao Fei, Wei Ji, and Tat-Seng Chua. 2023b. Cross2StrA: Unpaired cross-lingual image captioning with cross-lingual cross-modal structure-pivoted alignment. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2593–2608, Toronto, Canada. Association for Computational Linguistics.
- Yujia Xie, Luowei Zhou, Xiyang Dai, Lu Yuan, Nguyen Bach, Ce Liu, and Michael Zeng. 2022. Visual clues: Bridging vision and language foundations for image paragraph captioning. In *Advances in Neural Information Processing Systems*, volume 35, pages 17287–17300. Curran Associates, Inc.

- Guanghai Xu, Shuaicheng Niu, Mingkui Tan, Yucheng Luo, Qing Du, and Qi Wu. 2021. Towards accurate text-based image captioning with content diversity exploration. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 12637–12646. <https://doi.org/10.1109/CVPR46437.2021.01245>
- Jitao Xu, François Buet, Josep Crego, Elise Bertin-Lemée, and François Yvon. 2022. Joint generation of captions and subtitles with dual decoding. In *Proceedings of the 19th International Conference on Spoken Language Translation (IWSLT 2022)*, pages 74–82, Dublin, Ireland (in-person and online). Association for Computational Linguistics. <https://doi.org/10.18653/v1/2022.iwslt-1.7>
- Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich Zemel, and Yoshua Bengio. 2015. Show, attend and tell: Neural image caption generation with visual attention. In *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pages 2048–2057, Lille, France. PMLR.
- Liming Xu, Quan Tang, Jiancheng Lv, Bochuan Zheng, Xianhua Zeng, and Weisheng Li. 2023. Deep image captioning: A review of methods, trends and future challenges. *Neurocomputing*, 546:126287. <https://doi.org/10.1016/j.neucom.2023.126287>
- Semih Yagcioglu, Erkut Erdem, Aykut Erdem, and Ruket Cakici. 2015. A distributed representation based query expansion approach for image captioning. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 106–111, Beijing, China. Association for Computational Linguistics. <https://doi.org/10.3115/v1/P15-2018>
- Kun Yan, Lei Ji, Huaishao Luo, Ming Zhou, Nan Duan, and Shuai Ma. 2021. Control image captioning spatially and temporally. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 2014–2025, Online. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2021.acl-long.157>
- Bang Yang, Fenglin Liu, Xian Wu, Yaowei Wang, Xu Sun, and Yuexian Zou. 2023a. MultiCapCLIP: Auto-encoding prompts for zero-shot multilingual visual captioning. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 11908–11922, Toronto, Canada. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2023.acl-long.664>
- Cheng-Fu Yang, Yao-Hung Hubert Tsai, Wan-Cyuan Fan, Russ R. Salakhutdinov, Louis-Philippe Morency, and Frank Wang. 2022. Paraphrasing is all you need for novel object captioning. In *Advances in Neural Information Processing Systems*, volume 35, pages 6492–6504. Curran Associates, Inc.
- Fan Yang, Shalini Ghosh, Emre Barut, Kechen Qin, Prashan Wanigasekara, Chengwei Su, Weitong Ruan, and Rahul Gupta. 2024. Masking latent gender knowledge for debiasing image captioning. In *Proceedings of the 4th Workshop on Trustworthy Natural Language Processing (TrustNLP 2024)*, pages 227–238, Mexico City, Mexico. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2024.trustnlp-1.19>
- Li-Chuan Yang, Chih-Yuan Yang, and Jane Yung-jen Hsu. 2021a. Object relation attention for image paragraph captioning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(4):3136–3144. <https://doi.org/10.1609/aaai.v35i4.16423>
- Linjie Yang, Kevin Tang, Jianchao Yang, and Li-Jia Li. 2017. Dense captioning with joint inference and visual context. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. <https://doi.org/10.1109/CVPR.2017.214>
- Xu Yang, Chongyang Gao, Hanwang Zhang, and Jianfei Cai. 2021b. Auto-parsing network for image captioning and visual question answering. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*,

- pages 2197–2207. <https://doi.org/10.1109/ICCV48922.2021.00220>
- Xu Yang, Jiawei Peng, Zihua Wang, Haiyang Xu, Qinghao Ye, Chenliang Li, Songfang Huang, Fei Huang, Zhangzikang Li, and Yu Zhang. 2023b. Transforming visual scene graphs to image captions. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 12427–12440, Toronto, Canada. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2023.acl-long.694>
- Xu Yang, Kaihua Tang, Hanwang Zhang, and Jianfei Cai. 2019a. Auto-encoding scene graphs for image captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. <https://doi.org/10.1109/CVPR.2019.01094>
- Xu Yang, Yongliang Wu, Mingzhuo Yang, Haokun Chen, and Xin Geng. 2023c. Exploring diverse in-context configurations for image captioning. In *Advances in Neural Information Processing Systems*, volume 36, pages 40924–40943. Curran Associates, Inc.
- Xu Yang, Hanwang Zhang, and Jianfei Cai. 2019b. Learning to collocate neural modules for image captioning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*. <https://doi.org/10.1109/ICCV.2019.00435>
- Xuewen Yang, Svebor Karaman, Joel Tetreault, and Alejandro Jaimes. 2021c. Journalistic guidelines aware news image captioning. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 5162–5175, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2021.emnlp-main.419>
- Xuewen Yang, Heming Zhang, Di Jin, Yingru Liu, Chi-Hao Wu, Jianchao Tan, Dongliang Xie, Jue Wang, and Xin Wang. 2020. Fashion captioning: Towards generating accurate descriptions with semantic rewards. In *Computer Vision – ECCV 2020*, pages 1–17, Cham. Springer International Publishing. https://doi.org/10.1007/978-3-030-58601-0_1
- Zhengyuan Yang, Yijuan Lu, Jianfeng Wang, Xi Yin, Dinei Florencio, Lijuan Wang, Cha Zhang, Lei Zhang, and Jiebo Luo. 2021d. Tap: Text-aware pre-training for text-VQA and text-caption. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8751–8761. <https://doi.org/10.1109/CVPR46437.2021.00864>
- Zhilin Yang, Ye Yuan, Yuexin Wu, William W. Cohen, and Russ R. Salakhutdinov. 2016. Review networks for caption generation. In *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc.
- Zhuolin Yang, Wei Ping, Zihan Liu, Vijay Korthikanti, Weili Nie, De-An Huang, Linxi Fan, Zhiding Yu, Shiyi Lan, Bo Li, Mohammad Shoeybi, Ming-Yu Liu, Yuke Zhu, Bryan Catanzaro, Chaowei Xiao, and Anima Anandkumar. 2023d. Re-ViLM: Retrieval-augmented visual language model for zero and few-shot image captioning. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 11844–11857, Singapore. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2023.findings-emnlp.793>
- Li Yao, Nicolas Ballas, Kyunghyun Cho, John R. Smith, and Yoshua Bengio. 2016. Empirical performance upper bounds for image and video captioning. In *ICLR*.
- Linli Yao, Weiyang Wang, and Qin Jin. 2022. Image difference captioning with pre-training and contrastive learning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 36(3):3108–3116. <https://doi.org/10.1609/aaai.v36i3.20218>
- Ting Yao, Yingwei Pan, Yehao Li, and Tao Mei. 2017a. Incorporating copying mechanism in image captioning for learning novel objects. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. <https://doi.org/10.1109/CVPR.2017.559>
- Ting Yao, Yingwei Pan, Yehao Li, and Tao Mei. 2018. Exploring visual relationship for image captioning. In *Proceedings of the European Conference on Computer Vision (ECCV)*.

https://doi.org/10.1007/978-3-030-01264-9_42

- Ting Yao, Yingwei Pan, Yehao Li, and Tao Mei. 2019. Hierarchy parsing for image captioning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*. <https://doi.org/10.1109/ICCV.2019.00271>
- Ting Yao, Yingwei Pan, Yehao Li, Zhaofan Qiu, and Tao Mei. 2017b. Boosting image captioning with attributes. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*. <https://doi.org/10.1109/ICCV.2017.524>
- Mark Yatskar, Michel Galley, Lucy Vanderwende, and Luke Zettlemoyer. 2014. See no evil, say no evil: Description generation from densely labeled images. In *Proceedings of the Third Joint Conference on Lexical and Computational Semantics (*SEM 2014)*, pages 110–120, Dublin, Ireland. Association for Computational Linguistics and Dublin City University. <https://doi.org/10.3115/v1/S14-1015>
- Yanzhi Yi, Hangyu Deng, and Jinglu Hu. 2020. Improving image captioning evaluation by considering inter references variance. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 985–994, Online. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2020.acl-main.93>
- Guojun Yin, Lu Sheng, Bin Liu, Nenghai Yu, Xiaogang Wang, and Jing Shao. 2019. Context and attribute grounded dense captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. <https://doi.org/10.1109/CVPR.2019.00640>
- Quanzeng You, Hailin Jin, Zhaowen Wang, Chen Fang, and Jiebo Luo. 2016. Image captioning with semantic attention. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. <https://doi.org/10.1109/CVPR.2016.503>
- Licheng Yu, Eunbyung Park, Alexander C. Berg, and Tamara L. Berg. 2015. Visual Madlibs: Fill in the blank description generation and question answering. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*. <https://doi.org/10.1109/ICCV.2015.283>
- Zhihao Yuan, Xu Yan, Yinghong Liao, Yao Guo, Guanbin Li, Shuguang Cui, and Zhen Li. 2022. X-trans2cap: Cross-modal knowledge transfer using transformer for 3D dense captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8563–8573. <https://doi.org/10.1109/CVPR52688.2022.00837>
- Zihao Yue, Anwen Hu, Liang Zhang, and Qin Jin. 2023. Learning descriptive image captioning via semipermeable maximum likelihood estimation. In *Advances in Neural Information Processing Systems*, volume 36, pages 79124–79141. Curran Associates, Inc.
- Zequan Zeng, Yan Xie, Hao Zhang, Chiyu Chen, Bo Chen, and Zhengjue Wang. 2024. MeaCap: Memory-augmented zero-shot image captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14100–14110. <https://doi.org/10.1109/CVPR52733.2024.01337>
- Zequan Zeng, Hao Zhang, Ruiying Lu, Dongsheng Wang, Bo Chen, and Zhengjue Wang. 2023. Conzic: Controllable zero-shot image captioning by sampling-based polishing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 23465–23476. <https://doi.org/10.1109/CVPR52729.2023.02247>
- Hongkuan Zhang, Saku Sugawara, Akiko Aizawa, Lei Zhou, Ryohei Sasano, and Koichi Takeda. 2022a. Cross-modal similarity-based curriculum learning for image captioning. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 7599–7606, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2022.emnlp-main.516>
- Junzhe Zhang and Xiaojun Wan. 2023. Exploring the impact of vision features in news image captioning. In *Findings of the Association for Computational Linguistics: ACL*

- 2023, pages 12923–12936, Toronto, Canada. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2023.findings-acl.818>
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020a. Bertscore: Evaluating text generation with bert. In *International Conference on Learning Representations*.
- Wei Zhang, Yue Ying, Pan Lu, and Hongyuan Zha. 2020b. Learning long- and short-term user literal-preference with multimodal hierarchical transformer network for personalized image caption. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(05):9571–9578. <https://doi.org/10.1609/aaai.v34i05.6503>
- Wenqiao Zhang, Haochen Shi, Jiannan Guo, Shengyu Zhang, Qingpeng Cai, Juncheng Li, Sihui Luo, and Yueting Zhuang. 2022b. Magic: Multimodal relational graph adversarial inference for diverse and unpaired text-based image captioning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 36(3):3335–3343. <https://doi.org/10.1609/aaai.v36i3.20243>
- Wenqiao Zhang, Haochen Shi, Siliang Tang, Jun Xiao, Qiang Yu, and Yueting Zhuang. 2021a. Consensus graph representation learning for better grounded image captioning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(4):3394–3402. <https://doi.org/10.1609/aaai.v35i4.16452>
- Xuying Zhang, Xiaoshuai Sun, Yunpeng Luo, Jiayi Ji, Yiyi Zhou, Yongjian Wu, Feiyue Huang, and Rongrong Ji. 2021b. RSTNet: Captioning with adaptive attention on visual and non-visual words. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 15465–15474. <https://doi.org/10.1109/CVPR46437.2021.01521>
- Dora Zhao, Angelina Wang, and Olga Russakovsky. 2021. Understanding and evaluating racial biases in image captioning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 14830–14840. <https://doi.org/10.1109/ICCV48922.2021.01456>
- Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. 2017. Men also like shopping: Reducing gender bias amplification using corpus-level constraints. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 2979–2989, Copenhagen, Denmark. Association for Computational Linguistics. <https://doi.org/10.18653/v1/D17-1323>
- Sanqiang Zhao, Piyush Sharma, Tomer Levinboim, and Radu Soricut. 2019. Informative image captioning with external sources of information. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 6485–6494, Florence, Italy. Association for Computational Linguistics. <https://doi.org/10.18653/v1/P19-1650>
- Wentian Zhao, Xinxiao Wu, and Xiaoxun Zhang. 2020. MemCap: Memorizing style knowledge for image captioning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(07):12984–12992. <https://doi.org/10.1609/aaai.v34i07.6998>
- Yu Zhao, Hao Fei, Wei Ji, Jianguo Wei, Meishan Zhang, Min Zhang, and Tat-Seng Chua. 2023. Generating visual spatial description via holistic 3D scene understanding. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 7960–7977, Toronto, Canada. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2023.acl-long.442>
- Yu Zhao, Jianguo Wei, ZhiChao Lin, Yueheng Sun, Meishan Zhang, and Min Zhang. 2022. Visual spatial description: Controlled spatial-oriented image-to-text generation. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 1437–1449, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2022.emnlp-main.93>
- Ervine Zheng and Qi Yu. 2023. Evidential interactive learning for medical image captioning. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of

- Proceedings of Machine Learning Research*, pages 42478–42491. PMLR.
- Yue Zheng, Yali Li, and Shengjin Wang. 2019. Intention oriented image captions with guiding objects. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. <https://doi.org/10.1109/CVPR.2019.00859>
- Wenjie Zhong and Yusuke Miyao. 2021. Leveraging partial dependency trees to control image captions. In *Proceedings of the Second Workshop on Advances in Language and Vision Research*, pages 16–21, Online. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2021.alvr-1.3>
- Yiwu Zhong, Liwei Wang, Jianshu Chen, Dong Yu, and Yin Li. 2020. Comprehensive image captioning via scene graph decomposition. In *Computer Vision – ECCV 2020*, pages 211–229, Cham. Springer International Publishing. https://doi.org/10.1007/978-3-030-58568-6_13
- Zhipeng Zhong, Fei Zhou, and Guoping Qiu. 2023. Aesthetically relevant image captioning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 37(3):3733–3741. <https://doi.org/10.1609/aaai.v37i3.25485>
- Luowei Zhou, Hamid Palangi, Lei Zhang, Houdong Hu, Jason Corso, and Jianfeng Gao. 2020a. Unified vision-language pre-training for image captioning and VQA. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(07):13041–13049. <https://doi.org/10.1609/aaai.v34i07.7005>
- Mingyang Zhou, Grace Luo, Anna Rohrbach, and Zhou Yu. 2022. Focus! relevant and sufficient context selection for news image captioning. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 6078–6088, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2022.findings-emnlp.450>
- Yuanen Zhou, Meng Wang, Daqing Liu, Zhenzhen Hu, and Hanwang Zhang. 2020b. More grounded image captioning by distilling image-text matching model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. <https://doi.org/10.1109/CVPR42600.2020.00483>
- Yucheng Zhou and Guodong Long. 2023. Style-aware contrastive learning for multi-style image captioning. In *Findings of the Association for Computational Linguistics: EACL 2023*, pages 2257–2267, Dubrovnik, Croatia. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2023.findings-eacl.169>
- Yaoming Zhu, Sidi Lu, Lei Zheng, Jiaxian Guo, Weinan Zhang, Jun Wang, and Yong Yu. 2018. Taxygen: A benchmarking platform for text generation models. In *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval, SIGIR '18*, pages 1097–1100, New York, NY, USA. Association for Computing Machinery. <https://doi.org/10.1145/3209978.3210080>

A Experiments: Additional Information

A.1 Datasets

Flickr8k. We use the official version of the Flickr8k dataset.¹⁴ We follow previous work (Anderson et al., 2016; Hessel et al., 2021) when handling references that are also included in the candidate set: In Flickr8k-Expert we remove these sentences from the candidate set (158 examples), and in Flickr8k-CF we remove these sentences from the reference set.

In Flickr8k-expert, annotators are instructed to rank the candidate on a scale of 1 to 4, with the following specifications:

4 = Sentence describes the image: Sentence contains no errors, everything described in the sentence appears in the image.

3 = Sentence almost describes the image: Sentence contains only a single moderate error or a small number of minor errors, major details described in the sentence appear in the image.

2 = Sentence barely describes the image: Sentence contains a major error or many moderate or minor errors, only some minor details described in the sentence appear in the image.

1 = Sentence does not describe the image at all: No details described in the sentence appear in the image, sentence is unrelated to image.

In Flickr8k-CF, annotators are shown images with ten sentences below each images, and are asked to select all sentences that could describe the image, with the following specifications:

Yes- The sentence is a good enough description of the image (minor details may be wrong, and not everything that appears in the image has to be described).

No- The sentence does not describe the image because it contains major errors or is completely unrelated to the image.

Composite. We use the official version of the Composite dataset.¹⁵ While previous work mentioned that there are three candidates per image and a single score per candidate, we find that the dataset contains four candidates per image for MSCOCO and Flickr30k, and two scores per candidate (correctness and thoroughness). We find that results from prior work was most closely replicated when we discarded the fourth candidate when one exists and used the correctness score as the total score. Also, following Hessel et al. (2021) we remove sentences that appear both as references and as candidates from the reference set.

Instructions for the correctness rating (which we use in this study) are as follows:

Descriptions that correctly describe the image content with higher precision have better correctness ratings.

Score 1: The description has no relevance to the image.

Score 2: The description have only weak relevance to the image.

Score 3: The description have some relevance to the image.

Score 4: The description relates closely to the image.

Score 5: The description relates perfectly to the image.

THumB. We use the official version of the THumB dataset.¹⁶ Following the practice in other datasets, we remove sentences that appear both as references and as candidates from the reference set.

In this study, the authors performed the annotations themselves, relying on general descriptions for each aspect of the rating rather than providing specific instructions.

Polaris. We use the HuggingFace version of the Polaris dataset.¹⁷

Annotators were instructed to rate the candidate on a scale of 1 to 5, with the following specifications:

Evaluate the extent to which the sentences capture the features of the image, considering fluency, relevance and descriptiveness.

¹⁴www.kaggle.com/datasets/dibyansudiptiman/flickr-8k.

¹⁵imagesdg.wordpress.com/image-to-scene-description-graph.

¹⁶github.com/jungokasai/THumB.

¹⁷huggingface.co/datasets/yuwd/Polaris.



Figure 7: Image provided as an example in the Polaris dataset instructions.

Use a five-point scale for your evaluation, employing a point addition or deduction system.

For descriptiveness, rate a caption as ‘Excellent’ if it is both detailed and accurate.

For relevance, award points if the caption accurately depicts any present object, even if irrelevant words are included. Deduct points for each irrelevant word.

For fluency, assign lower scores if there are grammatical errors.

For example, consider the caption ‘A man wearing a helmet is skateboarding with a dog behind him.’ with the given image (Figure 7). This caption scores high in relevance as it accurately describes the main object in the image. However, the mention of ‘a dog,’ which is absent from the image, reduces its descriptiveness. Despite its high fluency, this warrants a ‘Fair’ (3 on a scale of 1 to 5) rating overall.

Pascal50S. The link to download the official version of the dataset in the CIDEr blog post [Link] seems to be broken; we thank (REMOVED FOR ANONIMYZATION) for providing us with the original files.

A.2 Metrics Implementations.

The five dominant metrics. For the five most dominant metrics we use the pycocoevalcap implementation. [Link]

Exact/fuzzy noun/verb overlap. We use an implementation [Link] provided to us by the author.

CLIPImageScore. Since no implementation is provided by the authors, we implement the metric ourselves. We use the SD-XL 1.0-base model [Link] to generate the image based on the candidate.

(Ref) CLIP and PAC Score, Polos. We use the official implementations. [CLIP] [PAC] [Polos]

BLIP2Score. Since no implementation is provided by the authors, we implement the metric ourselves. We use the BLIP2 image-text-matching model from the LAVIS package (Li et al., 2023a).

MPNetScore. Since no implementation is provided by the authors, we implement the metric ourselves. We use the all-mpnet-base-v2 model. [Link]

A.3 Feature Selection and Linear Regression

We use the scikit-learn [Link] implementations of feature selection and linear regression. We use forward-selection in feature selection.

A.4 Reproducibility Issues

In our experiments in Section 6.4 we tried our best to replicate the results reported by previous studies. We now describe the cases where we were unable to. We compare our results to Hessel et al. (2021), Kasai et al. (2022), Sarto et al. (2023), Hu et al. (2023a) and Wada et al. (2024).

The five dominant metrics. Across some of the datasets, we found small differences between our values and previous reported values, no larger than 0.2. One exception is the THUMB dataset, where prior work (Kasai et al., 2022; Hu et al., 2023a) report no digits after the decimal point, so the differences are no larger than 1.0, except for BLEU4, where the difference was larger because the authors used the sacrebleu implementation while we used the pycocoevalcap implementation, consistent with the common practice in image captioning.

(Ref) CLIP and PAC Score. Across all datasets our results are identical to those previous reported.

One exception is the Polaris dataset, where our results for these metrics differed from those of Wada et al. (2024). Since we used the HuggingFace version of the dataset and the official implementations for these metrics, we assume this discrepancy results from the authors using a more preliminary version of the dataset.

B Metric Usage Records

We now list all the papers we examined and the metric categories they used. Abbreviations: LexSim: Lexical similarity; RefPhSim: Candidate-reference phrase semantic similarity; RefSenSim: Candidate-reference sentence-level similarity; InfoSim: Extracted information similarity; ImPhSim: Candidate-image phrase-level semantic similarity; SenImSim: Sentence-image similarity; Ret: Retrieval-based methods; MulSim: Candidate-multiple source similarity; Flu: Fluency, Div: Diversity; Yngve: Yngve score; Bias: Bias.

2024. *NLP* Yang et al. (2024): LexSim, RefPhSim, Bias, Ramos et al. (2024): LexSim, Tanaka et al. (2024): Div, Qu et al. (2024): LexSim, Li et al. (2024b): LexSim, SenImSim

Vision Ge et al. (2024): RefSenSim, SenImSim, Huang et al. (2024): LexSim, RefPhSim, Zeng et al. (2024): LexSim, RefPhSim, SenImSim, Li et al. (2024a): LexSim, RefPhSim

Machine Learning Ma et al. (2024b): LexSim, RefPhSim, Black et al. (2024): LexSim, RefSenSim, Liu et al. (2024): LexSim, RefPhSim, Ma et al. (2024a): LexSim, Qi et al. (2024): LexSim, Qiu et al. (2024): LexSim, RefPhSim, Wang et al. (2024): LexSim, RefPhSim, Fu et al. (2024): LexSim, RefPhSim.

2023. *NLP* Wang et al. (2023c): LexSim, RefPhSim, Qiu et al. (2023): LexSim, RefPhSim, SenImSim, Hwang and Shwartz (2023): LexSim, RefPhSim, Chan et al. (2023a): LexSim, RefPhSim, Ret, Maeda et al. (2023): LexSim, RefPhSim, SenImSim, MulSim, Yang et al. (2023d): LexSim, RefPhSim, Mohamed et al. (2023): LexSim, Wu et al. (2023b): LexSim, Yang et al. (2023a): LexSim, RefPhSim, Yang et al. (2023b): LexSim, RefPhSim, Rajakumar Kalarani et al. (2023): LexSim, Ramos et al. (2023b): LexSim, Ou et al. (2023): Flu, Zhang and Wan (2023): LexSim, Bielawski and VanRullen (2023): LexSim, Anagnostopoulou et al. (2023): LexSim, Zhao et al. (2023): LexSim, RefPhSim, Ramos et al. (2023a): LexSim, RefPhSim, Zhou and Long (2023): LexSim, RefPhSim, Flu.

Vision Wu et al. (2023a): LexSim, Li et al. (2023b): Flu, Div, Kornblith et al. (2023): LexSim, SenImSim, Ret, MulSim, Hu et al. (2023b): LexSim, RefPhSim, Tu et al. (2023): LexSim, RefPhSim, Fei et al. (2023a): LexSim, RefPhSim, Fan et al. (2023): LexSim, RefPhSim, Kang et al. (2023): LexSim, RefPhSim, SenImSim, Ret, Barraco et al. (2023): LexSim, RefPhSim, SenImSim, MulSim, Hu et al. (2023b): LexSim, RefPhSim, Chen et al. (2023b): LexSim, Dessì et al. (2023): LexSim, RefPhSim, Vo et al. (2023): LexSim, RefPhSim, SenImSim, Ret, MulSim, Luo et al. (2023a): LexSim, RefPhSim, Zeng et al. (2023): LexSim, RefPhSim, SenImSim, MulSim, Div, Kuo and Kira (2023): LexSim, RefPhSim, Chen et al. (2023a): LexSim, Ramos et al. (2023c): LexSim, RefPhSim, Hirota et al. (2023): LexSim, RefPhSim, SenImSim, Bias, Ren et al. (2023): LexSim, RefPhSim.

Machine Learning Yang et al. (2023c): LexSim, Yue et al. (2023): SenImSim, Ret, Flu, Div, Luo et al. (2023b): SenImSim, Ret, Zheng and Yu (2023): LexSim, Li et al. (2023c): LexSim, RefPhSim, Fei et al. (2023b): LexSim, RefPhSim, Zhong et al. (2023): LexSim, RefPhSim, Nguyen et al. (2023): LexSim, RefPhSim, Wang et al. (2023a): LexSim, RefPhSim, Wang et al. (2023b): LexSim, RefPhSim.

2022. *NLP* Gao et al. (2022): LexSim, Div, Zhang et al. (2022a): LexSim, RefPhSim, Mirchandani et al. (2022): LexSim, Nukrai et al. (2022): LexSim, Zhou et al. (2022): LexSim, Zhao et al. (2022): LexSim, RefPhSim, Nishino et al. (2022): LexSim, RefPhSim, InfoSim, Cafagna et al. (2022): LexSim, RefPhSim, Cho et al. (2022): LexSim, RefPhSim, SenImSim, Ret, MulSim, Pantazopoulos et al. (2022): LexSim, RefPhSim, Xu et al. (2022): LexSim, Guo et al. (2022): LexSim.

Vision Jiao et al. (2022): LexSim, Wang et al. (2022a): LexSim, RefPhSim, Ruta et al. (2022): LexSim, Wang et al. (2022c): LexSim, RefPhSim, Nguyen et al. (2022): LexSim, RefPhSim, Meng et al. (2022): LexSim, RefPhSim, Chen et al. (2022): LexSim, Hirota et al. (2022): LexSim, Bias, Cai et al. (2022): LexSim, Wu et al. (2022): LexSim, RefPhSim, Chen et al. (2022): LexSim, Yuan et al. (2022): LexSim, Li et al. (2022): LexSim, RefPhSim, Fei et al. (2022): LexSim, RefPhSim, Liu et al. (2022): LexSim, RefPhSim, Vo et al. (2022): LexSim, RefPhSim, Mohamed et al. (2022): LexSim, Fang et al. (2022): LexSim, RefPhSim, Hu et al. (2022): LexSim, RefPhSim, Kuo and Kira (2022): LexSim, RefPhSim, Mavroudi and Vidal (2022): LexSim, RefPhSim.

Machine Learning Yang et al. (2022): LexSim, RefPhSim, Yang et al. (2022): LexSim, RefPhSim, SenImSim, Div, Xie et al. (2022): ImPhSim, Fei (2022): LexSim, RefPhSim, Zhang et al. (2022b): LexSim, RefPhSim, Div, Yao et al. (2022): LexSim, Wang et al. (2022b): LexSim, RefPhSim, Feng et al. (2022): LexSim, RefPhSim, Flu.

2021. *NLP* Yang et al. (2021c): LexSim, Liu et al. (2021): LexSim, Tu et al. (2021a): LexSim, RefPhSim, Chen et al. (2021a): LexSim, Shi et al. (2021b): LexSim, RefPhSim, Div, Ng et al. (2021): LexSim, RefPhSim, Div, Yan et al. (2021): LexSim, Shi et al. (2021c): LexSim, RefPhSim, Ret, Tu et al. (2021b): LexSim, RefPhSim, Bugliarello and Elliott (2021): LexSim, RefPhSim, Wang et al. (2021c): LexSim, RefPhSim, Honda et al. (2021): LexSim, RefPhSim, Nag Chowdhury et al. (2021): LexSim, RefPhSim, Ahsan et al. (2021): LexSim, RefPhSim, Zhong and Miyao (2021): LexSim, RefPhSim, Nagasawa et al. (2021): LexSim, RefPhSim.

Vision Shi et al. (2021a): LexSim, Div, Kim et al. (2021a): LexSim, RefPhSim, Zhao et al. (2021): LexSim, RefPhSim, Yang et al. (2021b): LexSim, RefPhSim, Nguyen et al. (2021): LexSim, RefPhSim, Yang et al. (2021d): LexSim, Chen et al. (2021b): LexSim, RefPhSim, Div, Hosseinzadeh and Wang (2021): LexSim, RefPhSim, Xu et al. (2021): LexSim, RefPhSim, Div, Wang et al. (2021a): LexSim, RefPhSim, Chen et al. (2021c): LexSim, Zhang et al. (2021b): LexSim, RefPhSim.

Machine Learning Hu et al. (2021): LexSim, RefPhSim, Ji et al. (2021): LexSim, RefPhSim, Fei (2021b): LexSim, RefPhSim, Fei (2021a): LexSim, RefPhSim, Song et al. (2021): LexSim, RefPhSim, Luo et al. (2021): LexSim, RefPhSim, Zhang et al. (2021a): LexSim, RefPhSim, Yang et al. (2021a): LexSim, Kim et al. (2021b): LexSim.

2020. *NLP* Milewski et al. (2020): LexSim, RefPhSim, Alikhani et al. (2020): LexSim, Shi et al. (2020b): LexSim, RefPhSim, Park et al. (2020): LexSim, Nie et al. (2020): LexSim, Nguyen et al. (2020): LexSim, RefPhSim, Takmaz et al. (2020): LexSim, RefPhSim.

Vision Wang et al. (2020a): LexSim, RefPhSim, Ret, Div, Sidorov et al. (2020): LexSim, RefPhSim, Wang et al. (2020c): LexSim, RefPhSim, Yang et al. (2020): LexSim, RefPhSim, Deng et al. (2020): LexSim, RefPhSim, Zhong et al. (2020): LexSim, RefPhSim, Div, Shi et al. (2020a): LexSim, RefPhSim, Gurari et al. (2020): LexSim, RefPhSim, Ma et al. (2020): LexSim, RefPhSim, Chen et al. (2020): LexSim, RefPhSim, Sammani and Melas-Kyriazi (2020): LexSim, RefPhSim, Guo et al. (2020): LexSim, RefPhSim, Chen and Jin (2020): LexSim, RefPhSim, Div, Cornia et al. (2020): LexSim, RefPhSim, Pan et al. (2020): LexSim, RefPhSim, Zhou et al. (2020b): LexSim, RefPhSim, Tran et al. (2020): LexSim.

Machine Learning Del Chiaro et al. (2020): LexSim, Mahajan and Roth (2020): LexSim, RefPhSim, Div, Zhang et al. (2020b): LexSim, RefPhSim, Zhou et al. (2020a): LexSim, RefPhSim, Zhao et al.

(2020): LexSim, Flu, Wang et al. (2020b): LexSim, RefPhSim, Liu et al. (2020): LexSim, Hou et al. (2020): LexSim, RefPhSim, Cao et al. (2020): LexSim, RefPhSim.

2019. *NLP* Zhao et al. (2019): LexSim, Fan et al. (2019): LexSim, Nikolaus et al. (2019): LexSim, RefPhSim, Changpinyo et al. (2019): LexSim, RefPhSim, Kim et al. (2019b): LexSim, RefPhSim, Wang et al. (2019a): LexSim, RefPhSim.

Vision Shen et al. (2019): LexSim, Div, Liu et al. (2019a): LexSim, RefPhSim, Ret, Div, Yngve, Ge et al. (2019): LexSim, Yao et al. (2019): LexSim, RefPhSim, Yang et al. (2019b): LexSim, RefPhSim, Aneja et al. (2019): LexSim, RefPhSim, Div, Park et al. (2019): LexSim, RefPhSim, Huang et al. (2019a): LexSim, RefPhSim, Laina et al. (2019): LexSim, RefPhSim, Ke et al. (2019): LexSim, RefPhSim, He et al. (2019): LexSim, Vered et al. (2019): LexSim, RefPhSim, Ret, Li et al. (2019a): LexSim, RefPhSim, Agrawal et al. (2019): LexSim, RefPhSim, Gu et al. (2019): LexSim, RefPhSim, Cornia et al. (2019): LexSim, RefPhSim, Zheng et al. (2019): LexSim, RefPhSim, Dognin et al. (2019): LexSim, SenImSim, Feng et al. (2019): LexSim, RefPhSim, Wang and Chan (2019): LexSim, RefPhSim, Div, Guo et al. (2019): LexSim, Flu, Yin et al. (2019): LexSim, Kim et al. (2019a): LexSim, Ret, Gao et al. (2019a): LexSim, RefPhSim, Qin et al. (2019): LexSim, RefPhSim, Yang et al. (2019a): LexSim, RefPhSim, Deshpande et al. (2019): LexSim, RefPhSim, Div, Biten et al. (2019): LexSim, RefPhSim, Li et al. (2019d): LexSim, RefPhSim, Shuster et al. (2019): LexSim, RefPhSim.

Machine Learning Herdade et al. (2019): LexSim, RefPhSim, Chen et al. (2019c): LexSim, RefPhSim, Div, Huang et al. (2019b): LexSim, RefPhSim, Wang et al. (2019b): LexSim, RefPhSim, Song et al. (2019): LexSim, Li et al. (2019b): LexSim, RefPhSim, Ret, Li et al. (2019c): LexSim, Gao et al. (2019b): LexSim, RefPhSim, Chen et al. (2019a): LexSim, RefPhSim, Chen et al. (2019b): LexSim, RefPhSim.

2018. *NLP* Sharma et al. (2018): LexSim, RefPhSim, Chen et al. (2018c): LexSim, Madhyastha et al. (2018): LexSim, RefPhSim, Liu et al. (2018a): LexSim, RefPhSim, Guo et al. (2018): LexSim, Melas-Kyriazi et al. (2018): LexSim, Lu et al. (2018): LexSim, Div, Wang et al. (2018): LexSim, RefPhSim, Cohn-Gordon et al. (2018): Ret, Chandrasekaran et al. (2018): NONE.

Vision Chen and Zhao (2018): LexSim, Liu et al. (2018b): LexSim, RefPhSim, Ret, Div, Gu et al. (2018b): LexSim, Div, Dai et al. (2018b): LexSim, RefPhSim, Hendricks et al. (2018): LexSim, Bias, Jiang et al. (2018b): LexSim, RefPhSim, Chen et al. (2018d): LexSim, Yao et al. (2018): LexSim, RefPhSim, Chen et al. (2018e): LexSim, RefPhSim, Luo et al. (2018): LexSim, RefPhSim, Ret, Mathews et al. (2018): LexSim, RefPhSim, Flu, Chen et al. (2018a): LexSim, Aneja et al. (2018): LexSim, RefPhSim, Anderson et al. (2018b): LexSim, RefPhSim.

Machine Learning Dai et al. (2018a): LexSim, RefPhSim, Div, Huang et al. (2018): LexSim, RefPhSim, Anderson et al. (2018a): LexSim, RefPhSim, Jiang et al. (2018a): LexSim, Div, Chen et al. (2018b): LexSim, Gu et al. (2018a): LexSim, RefPhSim.

2017. *NLP* Anderson et al. (2017): LexSim, RefPhSim.

Vision Shetty et al. (2017): LexSim, RefPhSim, Div, Liu et al. (2017b): LexSim, Gu et al. (2017): LexSim, RefPhSim, Pedersoli et al. (2017): LexSim, Tavakoli et al. (2017): LexSim, Yao et al. (2017b): LexSim, RefPhSim, Dai et al. (2017): LexSim, RefPhSim, Vedantam et al. (2017): LexSim, Gan et al. (2017a): LexSim, Ren et al. (2017): LexSim, Lu et al. (2017): LexSim, RefPhSim, Chunseong Park et al. (2017): LexSim, Yang et al. (2017): NONE, Gan et al. (2017b): LexSim, Chen et al. (2017a): LexSim, Venugopalan et al. (2017): LexSim, Yao et al. (2017a): LexSim, Sun et al. (2017): LexSim, Rennie et al. (2017): LexSim, Wang et al. (2017b): LexSim, RefPhSim, Div, Rohrbach et al. (2017): LexSim.

Machine Learning Dai and Lin (2017): LexSim, Ret, Wang et al. (2017a): LexSim, RefPhSim, Div, Liu et al. (2017a): LexSim, Chen et al. (2017b): LexSim, Li et al. (2017): LexSim, Mun et al. (2017): LexSim.

2016. NLP

Vision Hendricks et al. (2016): LexSim, Johnson et al. (2016): LexSim, Ret, You et al. (2016): LexSim.

Machine Learning Yang et al. (2016): LexSim, Pu et al. (2016): LexSim, Flu, Yao et al. (2016): LexSim, Mathews et al. (2016): LexSim.

2015. NLP Devlin et al. (2015): LexSim, Flu, Yagcioglu et al. (2015): LexSim, Chen et al. (2015): LexSim.

Vision Jia et al. (2015): LexSim, Ushiku et al. (2015): LexSim, Yu et al. (2015): LexSim, Chen and Lawrence Zitnick (2015): LexSim, Ret, Flu, Vinyals et al. (2015): LexSim, Ret, Fang et al. (2015): LexSim, Flu, Donahue et al. (2015): LexSim, Ret, Karpathy and Fei-Fei (2015): LexSim, Ret

Machine Learning Xu et al. (2015): LexSim, Lebrete et al. (2015): LexSim, Mao et al. (2015): LexSim, Ret, Flu.

2014. NLP Mason and Charniak (2014b): LexSim, Mason and Charniak (2014a): LexSim, Yatskar et al. (2014): LexSim, Kuznetsova et al. (2014): LexSim.

2013. NLP Kuznetsova et al. (2013): LexSim, Elliott and Keller (2013): LexSim.

2012. NLP Kuznetsova et al. (2012): LexSim, Mitchell et al. (2012): NONE.

2011. NLP Li et al. (2011): LexSim.

Machine Learning Ordonez et al. (2011): LexSim.

2010. NLP Feng and Lapata (2010): LexSim, Aker and Gaizauskas (2010): LexSim.