

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/268689555>

# CIDEr: Consensus-based Image Description Evaluation

Article · November 2014

Source: arXiv

---

CITATIONS

530

---

READS

873

3 authors, including:



[Devi Parikh](#)

Virginia Polytechnic Institute and State University

219 PUBLICATIONS 18,090 CITATIONS

[SEE PROFILE](#)

# CIDEr: Consensus-based Image Description Evaluation

Ramakrishna Vedantam  
Virginia Tech  
vrama91@vt.edu

C. Lawrence Zitnick  
Microsoft Research  
larryz@microsoft.com

Devi Parikh  
Virginia Tech  
parikh@vt.edu

## Abstract

*Automatically describing an image with a sentence is a long-standing challenge in computer vision and natural language processing. Due to recent progress in object detection, attribute classification, action recognition, etc., there is renewed interest in this area. However, evaluating the quality of descriptions has proven to be challenging. We propose a novel paradigm for evaluating image descriptions that uses human consensus. This paradigm consists of three main parts: a new triplet-based method of collecting human annotations to measure consensus, a new automated metric that captures consensus, and two new datasets: PASCAL-50S and ABSTRACT-50S that contain 50 sentences describing each image. Our simple metric captures human judgment of consensus better than existing metrics across sentences generated by various sources. We also evaluate five state-of-the-art image description approaches using this new protocol and provide a benchmark for future comparisons.*

## 1. Introduction

Recent advances in object recognition [12], attribute classification [18], action classification [21, 7] and crowd-sourcing [35] have increased the interest in solving higher level scene understanding problems. One such problem is generating human-like descriptions of an image. In spite of the growing interest in this area, the evaluation of novel sentences generated by automatic approaches remains challenging. Evaluation is critical for measuring progress and spurring improvements in the state of the art. This has already been shown in various problems in computer vision, such as detection [10, 5], segmentation [10, 22], and stereo [34].

Existing evaluation metrics for image description attempt to measure several desirable properties. These include grammaticality, saliency (covering main aspects), correctness/truthfulness, etc. Using human studies, these properties may be measured, e.g. on separate *one to five* [25, 32, 37, 8] or *pairwise* scales [38]. Unfortunately, combining these various results into one measure of sentence quality is

difficult. Alternatively, other works [17, 15] ask subjects to judge the overall quality of a sentence. However what humans like often does not correspond to what is human-like.<sup>1</sup>

Common to all of these approaches is the desire to measure the “human-likeness” of a sentence [19, 9]. That is, does an automatically generated sentence sound like a sentence that was written by a human? Rather than using indirect methods for measuring this, we propose to measure human-likeness directly. We introduce a novel consensus-based evaluation protocol, which measures the similarity of a sentence to the majority, or *consensus* of how most people describe the image (Fig. 1). One realization of this evaluation protocol uses human subjects to judge sentence similarity to the ground truth sentences. The question “Which of two sentences is most similar to another sentence?” is posed to the subjects. The resulting quality score is based on how often a sentence is labeled as being *more* similar to a human-generated sentence. The relative nature of the question helps make the task more objective. We encourage the reader to review how a similar protocol has been used in [36] to capture human perception of image similarity. These annotation protocols for similarity may be understood as instantiations of 2AFC (two alternative forced choice) [3], a popular modality in psychophysics.

Since human studies are expensive, hard to reproduce, and slow to evaluate, automatic evaluation measures are commonly desired. To be useful in practice, automated metrics should agree well with human judgment. Some popular metrics used for image description evaluation are BLEU [28] (precision-based) from the machine translation community and ROUGE [39] (recall-based) from the summarization community. Unfortunately, these metrics have been shown to correlate weakly with human judgment [17, 8, 4, 15]. For the task of judging the overall quality of a description, the METEOR [8] metric has shown better correlation with human subjects. Other metrics rely on the ranking of captions [15] and cannot evaluate novel image descriptions.

We propose a new automatic *consensus* metric of image description quality – CIDEr (Consensus-based Image De-

<sup>1</sup>This is a subtle but important distinction. We show qualitative examples of this in the supplementary material.

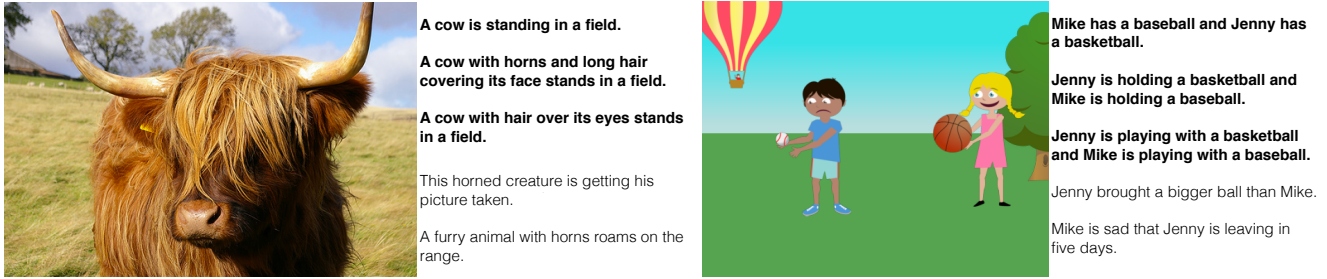


Figure 1: Images from our PASCAL-50S (left) and ABSTRACT-50S (right) datasets with a subset of corresponding (human) sentences. Sentences shown in **bold** are representative of the consensus descriptions for these images. We propose to capture such descriptions with our evaluation protocol.

scription Evaluation). Our metric measures the similarity of a generated sentence against a set of ground truth sentences written by humans. Our metric shows high agreement with our consensus-based measure using human subjects. Using sentence similarity, the notions of grammaticality, saliency, importance and accuracy (precision and recall) are inherently captured by our metric.

Existing datasets popularly used to evaluate image description approaches have a maximum of only five descriptions per image [30, 15, 27]. However, we find that five sentences are not sufficient for measuring how a “majority” of humans would describe an image. Thus, to accurately measure consensus, we collect two new evaluation datasets containing 50 descriptions per image – PASCAL-50S and ABSTRACT-50S. The PASCAL-50S dataset is based on the popular UIUC Pascal Sentence Dataset, which has 5 descriptions per image. This dataset has been used for both training and testing in numerous works [25, 17, 11, 32]. The ABSTRACT-50S dataset is based on the dataset of Zitnick and Parikh [40]. While previous methods have only evaluated using 5 sentences, we explore the use of 1 to ~50 sentences. Interestingly, we find that most metrics improve in performance with more sentences except BLEU.

**Contributions:** In this work, we propose a consensus-based evaluation protocol for image descriptions. We introduce a new annotation modality for human judgment, a new automated metric, and two new datasets. We compare the performance of five state-of-the-art machine generation approaches [25, 17, 11, 32]. Our code and datasets will be made publicly available. Finally, to facilitate the adoption of this protocol, we will set up an evaluation server where image description methods may be evaluated.

## 2. Related Work

**Textual and Visual Elements:** Numerous papers have studied the relationship between language constructs and image content. Berg *et al.* [2] characterize the relative importance of objects (nouns). Zitnick and Parikh [40] study relationships between visual and textual features by creating a synthetic Abstract Scenes Dataset. Other works have

modeled prepositional relationships [13], attributes (adjectives) [18, 29], and visual phrases (*i.e.* visual elements that co-occur) [33]. Recent works have utilized techniques in deep learning to learn joint embeddings of text and image fragments [16].

**Image Description Generation:** Various methods have been explored for generating full descriptions for images. Broadly, the techniques are either retrieval- [11, 27, 15] or generation-based [25, 17, 38, 32]. While some retrieval-based approaches use global retrieval [11], others retrieve text phrases and stitch them together in an approach inspired by extractive summarization [27]. Generative approaches have explored creating sentences by inference over image detections and text-based priors [17] or exploiting word co-occurrences using syntactic trees [25]. Rohrbach *et al.* [32] propose a machine translation approach that goes from an intermediate semantic representation to sentences. Yang *et al.* propose a slot filling approach [37] based on methods used for text summarization. Some other approaches include [14, 19]. Most of the approaches use the UIUC Pascal Sentence Dataset [11, 17, 25, 32, 14] for evaluation. Recent work has also looked at generating descriptions of images with dense human visual annotation [38]. In this work we focus on the problem of evaluating approaches using our novel consensus-based protocol.

**Automated Evaluation:** Automated evaluation metrics have been used in many domains within Artificial Intelligence (AI), such as statistical machine translation and text summarization. Some of the popular metrics in machine translation include those based on precision, such as BLEU [28] and those based on precision as well as recall, such as METEOR [1]. While BLEU (BiLingual Evaluation Understudy) has been the most popular metric, its effectiveness has been repeatedly questioned [17, 8, 4, 15]. A popular metric in the summarization community is ROUGE [39] (Recall Oriented Understudy of Gisting Evaluation). This metric is primarily recall-based and thus has a tendency to reward long sentences with high recall. These metrics have been shown to have weak to moderate correlation with human judgment [8]. Recently, METEOR has been used for image description evaluation with more promising re-

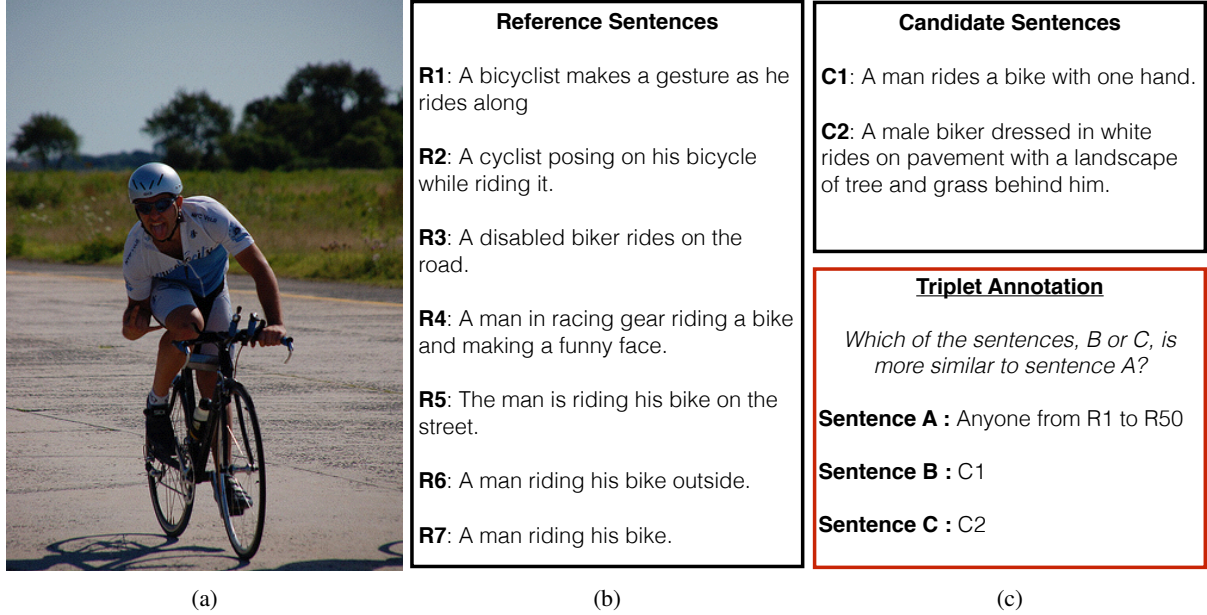


Figure 2: We show an illustration of our triplet annotation modality. Given an image (a), with reference sentences (b) and a pair of candidate sentences (c), we match them with a reference sentence one by one to form triplets. Subjects are shown these 50 triplets (red box) on Amazon Mechanical Turk and asked to pick which sentence (B or C) is more similar to sentence A.

sults [9]. Another metric proposed by Hodosh *et al.* [15] can only evaluate ranking-based approaches, it cannot evaluate novel sentences. We propose a consensus-based metric that rewards a sentence for being similar to the majority of human written descriptions. Our metric can evaluate both retrieved as well as generated (novel) sentences.

**Datasets:** Numerous datasets have been proposed for studying the problem of generating image descriptions. The most popular dataset is the UIUC Pascal Sentence Dataset [30]. This dataset contains 5 human written descriptions for 1,000 images. This dataset has been used by a number of approaches for training and testing. The SBU captioned photo dataset [27] contains one description per image for a million images, mined from the web. These are commonly used for training image description approaches. Approaches are then tested on a query set of 500 images with one sentence each. We demonstrate that more sentences per image are essential for reliable automatic evaluation. The Abstract Scenes dataset [40] contains cartoon-like images with two descriptions. The recently released MSCOCO dataset [20] contains five sentences for a collection of over 100K images. Other datasets of images and associated descriptions include [26, 15]. In this work, we introduce two new datasets. First is the PASCAL-50S dataset where we collected 50 sentences per image for the 1,000 images from UIUC Pascal Sentence dataset. The second is the ABSTRACT-50S dataset where we collected 50 sentences for a subset of 500 images from the Abstract Scenes dataset.

The rest of this paper is organized as follows. We first give details of our triplet human annotation modality

(Sec. 3). Then we provide the details of our consensus-based automated metric, CIDEr (Sec. 4). In Sec. 5 we provide the details of our two new image-sentence datasets, PASCAL-50S and ABSTRACT-50S. Our contributions of triplet annotation, metric and dataset make consensus-based image description evaluation feasible. Our results (Sec. 7) demonstrate that our automated metric and our proposed datasets capture consensus better than existing choices.

All our human studies are performed on the Amazon Mechanical Turk (AMT). Subjects are restricted to the United States, and other qualification criteria are imposed based on worker history<sup>2</sup>.

### 3. Consensus Interface

Given an image and a collection of human generated *reference* sentences describing it, the goal of our consensus-based protocol is to measure the similarity of a *candidate* sentence to a majority of how most people describe the image (*i.e.* the *reference* sentences). In this section, we describe our human study protocol for generating ground truth consensus scores. In Sec. 7, these ground truth scores are used to evaluate several automatic metrics including our proposed CIDEr metric.

An illustration of our human study interface is shown in Fig. 2. Subjects are shown three sentences: A, B and C. They are asked to pick which of two sentences (B or C) is most similar to sentence A. Sentences B and C are two candidate sentences, while sentence A is a reference sen-

<sup>2</sup>Approval rate greater than 95%, Minimum 500 HITs approved



tence. For each choice of B and C, we form triplets using all the reference sentences for an image. We provide no explicit concept of “similarity”. Interestingly, even though we do not say that the sentences are image descriptions, some workers commented that they were imagining the scene to make the choice. The relative nature of the task – “Which of the two sentences, B or C, is more similar to A?” – helps make the assessment more objective. That is, it is easier to judge if one sentence is more similar than another to a sentence, than to provide an absolute rating from 1 to 5 of the similarity between two sentences [3].

We collect three human judgments for each triplet. For every triplet, we take the majority vote of the three judgments. For each pair of candidate sentences (B, C), we assign B the winner if it is chosen as more similar by a majority of triplets, and similarly for C. These pairwise relative rankings are used to evaluate the performance of the automated metrics. That is, when automatic metrics give both sentences B and C a score, we check whether B received a higher score or C. Accuracy is computed as the proportion of candidate pairs on which humans and the automatic metric agree on which of the two sentences is the winner.

#### 4. CIDEr Metric

Our goal is to automatically evaluate for image  $I_i$  how well a candidate sentence  $c_i$  matches the consensus of a set of image descriptions  $S_i = \{s_{i1}, \dots, s_{im}\}$ . All words in the sentences (both candidate and references) are first mapped to their stem or root forms. That is, “fishes”, “fishing” and “fished” all get reduced to “fish.” We represent each sentence using the set of  $n$ -grams present in it. An  $n$ -gram  $\omega_k$  is a set of one or more ordered words. In this paper we use  $n$ -grams containing one to four words.

Intuitively, a measure of consensus would encode how often  $n$ -grams in the candidate sentence are present in the reference sentences. Similarly,  $n$ -grams not present in the reference sentences should not be in the candidate sentence. Finally,  $n$ -grams that commonly occur across all reference sentences should be given lower weight, since they are likely to be less informative. To encode this intuition, we perform a Term Frequency Inverse Document Frequency (TF-IDF) weighting for each  $n$ -gram [31]. The number of times an  $n$ -gram  $\omega_k$  occurs in a reference sentence  $s_{ij}$  is denoted by  $h_k(s_{ij})$  or  $h_k(c_i)$  for the candidate sentence  $c_i$ . We compute the TF-IDF weighting  $g_k(s_{ij})$  for each  $n$ -gram  $\omega_k$  using:

$$g_k(s_{ij}) = \frac{h_k(s_{ij})}{\sum_{\omega_l \in \Omega} h_l(s_{ij})} \log \left( \frac{|I|}{\sum_{I_p \in I} \min(1, \sum_q h_k(s_{pq}))} \right), \quad (1)$$

where  $\Omega$  is the vocabulary of all  $n$ -grams and  $I$  is the set of all images in the dataset. The first term measures the TF

of each  $n$ -gram  $\omega_k$ , and the second term measures the rarity of  $\omega_k$  using its IDF. Intuitively, TF places higher weight on  $n$ -grams that frequently occur in the reference sentences describing an image, while IDF reduces the weight of  $n$ -grams that commonly occur across all descriptions. That is, the IDF provides a measure of word saliency by discounting popular words that are likely to be less visually informative. The IDF is computed using the logarithm of the number of images in the dataset  $|I|$  divided by the number of images for which  $\omega_k$  occurs in any of its reference sentences.

Our CIDEr <sub>$n$</sub>  score for  $n$ -grams of length  $n$  is computed using the average cosine similarity between the candidate sentence and the reference sentences, which accounts for both precision and recall:

$$\text{CIDEr}_n(c_i, S_i) = \frac{1}{m} \sum_j \frac{\mathbf{g}^n(c_i) \cdot \mathbf{g}^n(s_{ij})}{\|\mathbf{g}^n(c_i)\| \|\mathbf{g}^n(s_{ij})\|}, \quad (2)$$

where  $\mathbf{g}^n(c_i)$  is a vector formed by  $g_k(c_i)$  corresponding to all  $n$ -grams of length  $n$  and  $\|\mathbf{g}^n(c_i)\|$  is the magnitude of the vector  $\mathbf{g}^n(c_i)$ . Similarly for  $\mathbf{g}^n(s_{ij})$ .

We use higher order (longer)  $n$ -grams to capture grammatical properties as well as richer semantics. We combine the scores from  $n$ -grams of varying lengths as follows:

$$\text{CIDEr}(c_i, S_i) = \sum_{n=1}^N w_n \text{CIDEr}_n(c_i, S_i), \quad (3)$$

Empirically, we found that uniform weights  $w_n = 1$  work the best. We use  $N = 4$ .

We also experimented with non-uniform  $w_n$  that weighs the longer  $n$ -grams more, soft word-level semantic similarity measures such as WordNet similarity [24] and word2vec similarity [23], as well as various relevance weighting schemes from information retrieval [31]. We found that the simple metric we present here performs the best. Interestingly, we find that word-level semantic similarity helps at fewer number of sentences. These soft-similarities help overcome the lack of data. But with more sentences, the noise introduced by existing similarity measures holds the metric back. Moreover, simplicity is critical to ensure that the metric is interpretable, transparent, and easy to adopt.

#### 5. New Datasets

We propose two new datasets – PASCAL-50S and ABSTRACT-50S – for evaluating methods that generate novel image descriptions. Both the datasets have 50 reference sentences per image for 1,000 and 500 images, respectively. These are intended as “testing” datasets, crafted to enable consensus-based evaluation. For a list of training datasets, we encourage the reader to explore [20, 27]. The PASCAL-50S dataset uses all 1,000 images from the UIUC Pascal Sentence Dataset [30] whereas the ABSTRACT-50S dataset uses 500 random images from the Abstract Scenes

Dataset [40]. The Abstract Scenes Dataset contains scenes made from clipart objects. Our two new datasets are very different both visually and in the type of image descriptions produced.

Our goal is to collect image descriptions that are objective and representative of the image content. Subjects are shown an image and a text box, asking them to “Describe what is going on in the image”. We ask subjects to capture the main aspects of the scene and provide descriptions that others are also likely to provide. This includes writing descriptions rather than “dialogs” or overly descriptive sentences. Workers were told that a good description should help others recognize the image from a collection of similar images. Instructions also mentioned that work with poor grammar would be rejected. Snapshots of our interface can be found in the supplementary material. Overall, we had 465 subjects for ABSTRACT-50S and 683 subjects for PASCAL-50S datasets. We ensure that each sentence for an image is written by a different subject.

The average sentence length for the ABSTRACT-50S dataset is 10.59 words compared to 8.8 words for PASCAL-50S. This is indicative of a denser semantic sampling in the ABSTRACT-50S dataset as compared to the PASCAL-50S dataset.

## 6. Experimental Setup

The goals of our experiments are two-fold:

- Evaluating how well our proposed metric CIDEr captures human judgement of consensus, as compared to existing metrics.
- Comparing existing state-of-the-art automatic image description approaches in terms of how well the descriptions they produce match human consensus of how to describe images.

We first detail how we select candidate sentences to evaluate using CIDEr (and other metrics). We then list the various existing metrics that we compare CIDEr to. We then list the various automatic image description approaches and our experimental set up to compare them.

**Candidate Sentences:** On ABSTRACT-50S, we use 48 of our 50 sentences as reference sentences (sentence A in our triplet annotation). The remaining 2 sentences per image can be used as candidate sentences. We form 400 pairs of candidate sentences (B and C in our triplet annotation). These include two kinds of pairs. The first are 200 human-human correct pairs (HC), where we pick two human sentences describing the same image. The second kind are 200 human-human incorrect pairs (HI), where one of the sentences is a human description for the image and the other is also a human sentence but describing some other image from the dataset picked at random.

For PASCAL-50S, our candidate sentences come from a diverse set of sources: human sentences from the UIUC

Pascal Sentence Dataset as well as machine-generated sentences from five state-of-the-art image description methods. These span both retrieval-based and generation-based methods: Midge [25], Babytalk [17], Story [11], and two versions of Translating Video Content to Natural Language Descriptions [32] (Video and Video+)<sup>3</sup>. We form 4,000 pairs of candidate sentences (again, B and C for our triplet annotation). These include four types of pairs (1,000 each). The first two are human-human correct (HC) and human-human incorrect (HI) similar to ABSTRACT-50S. The third are human-machine (HM) pairs formed by pairing a human sentence describing an image with a machine generated sentence describing the same image. Finally, the fourth are machine-machine (MM) pairs, where we compare two machine generated sentences describing the same image. We pick the machine generated sentences randomly, so that each method participates in roughly equal number of pairs, on a diverse set of images. Ours is the first work to perform a comprehensive evaluation across these different kinds of sentences.

For consistency, we drop two reference sentences for the PASCAL-50S evaluations so that we evaluate on both datasets (ABSTRACT-50S and PASCAL-50S) with a maximum of 48 reference sentences.

**Metrics:** The existing metrics used in the community for evaluation of image description approaches are BLEU [28] and ROUGE [39]<sup>4</sup>. BLEU is precision-based and ROUGE is recall-based. More specifically, image description methods have used versions of BLEU called BLEU<sub>1</sub> and BLEU<sub>4</sub>, and a version of ROUGE called ROUGE<sub>1</sub>. A recent survey paper [9] has used a different version of ROUGE called ROUGE<sub>S</sub>, as well as a machine translation metric called METEOR [1]. We now briefly describe these metrics along with other metrics that exist in the NLP community, but have not been used for image description evaluation so far. More details can be found in the supplementary material. **BLEU** stands for BiLingual Evaluation Understudy [28]. It computes an  $n$ -gram based precision for the candidate sentence with respect to the references. The key idea of BLEU is to compute precision by *clipping*. Clipping computes precision for a word, based on the maximum number of times it occurs in any reference sentence. Thus, a candidate sentence saying “The The The”, would get credit for saying only one “The”, if the word occurs utmost once across individual references. Versions of BLEU called BLEU<sub>1</sub>, BLEU<sub>2</sub>, BLEU<sub>3</sub> and BLEU<sub>4</sub> can be computed by varying the lengths of the  $n$ -grams used. Another version of BLEU, called BLEU<sub>Overall</sub><sup>5</sup> computes a geometric mean of

<sup>3</sup>We thank the authors of these approaches for making their outputs available to us.

<sup>4</sup>Another metric proposed by Hodosh *et al.* [15] can only evaluate ranking based approaches and not approaches that generate novel sentences.

<sup>5</sup>We say BLEU<sub>Overall</sub> to follow convention. The corresponding CIDEr<sub>Overall</sub> score is just referred to as CIDEr.

the  $n$ -gram scores from 1 to 4. It adds a “brevity-penalty” to discourage overly short sentences. It is typically used for evaluating machine translation approaches in the NLP community. **ROUGE** stands for Recall Oriented Understudy of Gisting Evaluation [39]. It computes  $n$ -gram based recall for the candidate sentence with respect to the references. It is a popular metric for summarization evaluation. Similar to BLEU, versions of ROUGE can be computed by varying the  $n$ -gram count. Two other versions of ROUGE are ROUGE<sub>S</sub> and ROUGE<sub>L</sub>. These compute an F-measure with a recall bias using *skip-bigrams* and *longest common subsequence* respectively, between the candidate and each reference sentence. Skip-bigrams are all pairs of ordered words in a sentence, sampled non-consecutively. Given these scores, they return the maximum score for the set of references as the judgment of quality. **METEOR** stands for Metric for Evaluation of Translation with Explicit ORdering [1]. Similar to ROUGE<sub>L</sub> and ROUGE<sub>S</sub>, it also computes the F-measure based on matches, and returns the maximum score over a set of references as its judgment of quality. However, it resolves word-level correspondences in a more sophisticated manner, using exact matches, stemming and semantic similarity. It optimizes over matches minimizing *chunkiness*. Chunkiness implies that matches should be consecutive, wherever possible. It also sets parameters favoring recall over precision in its F-measure computation. We implement all the metrics, except for METEOR, for which we use [6] (version 1.5).

**Machine Approaches:** We comprehensively evaluate which machine generation methods are best at matching consensus sentences. For this experiment, we select a subset of 100 images from the UIUC Pascal Sentence Dataset which have outputs for all the five machine description methods used in our evaluation: Midge [25], Babytalk [17], Story [11], and two versions of Translating Video Content to Natural Language Descriptions [32] (Video and Video+). For each image, we form all  ${}^5C_2$  pairs of machine-machine sentences. This ensures that each machine approach gets compared to all other machine approaches on each image. This gives us 1,000 pairs. We form triplets by “tripling” each pair with 20 random reference sentences. We collect human judgement of consensus using our triplet annotation modality as well as evaluate our proposed automatic consensus metric CIDEr. In both cases, we count the fraction of times a machine description method beats another method in terms of being more similar to the reference sentences. To the best of our knowledge, we are the first work to perform an exhaustive evaluation of the state-of-the-art in image description, across retrieval- and generation-based methods. We will set up an evaluation server that will enable future image description approaches to evaluate their approach on our proposed datasets.

## 7. Results

In this section we evaluate the effectiveness of our consensus-based metric CIDEr on the PASCAL-50S and ABSTRACT-50S datasets. We begin by exploring how many sentences are sufficient for reliably evaluating our consensus metric. Next, we compare our metric against several other commonly used metrics, and quantify machine performance by comparing to how well humans perform on the same task. Using CIDEr we evaluate several existing automatic image description approaches and compare the results to human judgement. Finally, we evaluate human performance at predicting consensus and compare it to (CIDEr).

### 7.1. How many sentences are enough?

We begin by analyzing how the number of reference sentences affects the accuracy of automated metrics. To quantify this, we collect 100 sentences for a subset of 50 randomly sampled images from the UIUC Pascal Sentence Dataset. We then pool human-human correct, human-machine, machine-machine and human-human incorrect sentence pairs (179 in total) and get triplet annotations. This gives us the ground truth consensus score for all pairs. We evaluate BLEU<sub>1</sub>, ROUGE<sub>1</sub> and CIDEr<sub>1</sub> with up to 100 reference sentences used to score the candidate sentences. We find that the accuracy improves significantly for the first 10 sentences (Fig. 3) for all metrics. From 1 to 5 sentences, the agreement for ROUGE<sub>1</sub> improves from 0.61 to 0.73. Both ROUGE<sub>1</sub> and CIDEr<sub>1</sub> continue to improve until reaching 50 sentences, after which the results begin to saturate. Curiously, BLEU<sub>1</sub> shows a decrease in performance with more sentences. BLEU does a max operation over sentence level matches, and thus as more sentences are used, the likelihood of matching a lower quality reference sentence increases. Based on this pilot, we collect 50 sentences per image for our ABSTRACT-50S and PASCAL-50S datasets. For the remaining experiments we explore results using 1 to 50 sentences.

### 7.2. Accuracy of Automated Metrics

We compare the performance of CIDEr, BLEU, ROUGE and METEOR with the human consensus scores in Fig. 4. That is, for each metric we compute the scores for two candidate sentences. The metric is correct if the sentence with higher score is the same as the sentence chosen by our human studies as being more similar to the reference sentences. The candidate sentences are both human and machine generated. For BLEU and ROUGE we show both their popular versions and the version we found to give best performance. We sample METEOR at fewer points due to high run-time. For a more comprehensive evaluation across different versions of each metric, please check the supplementary material.

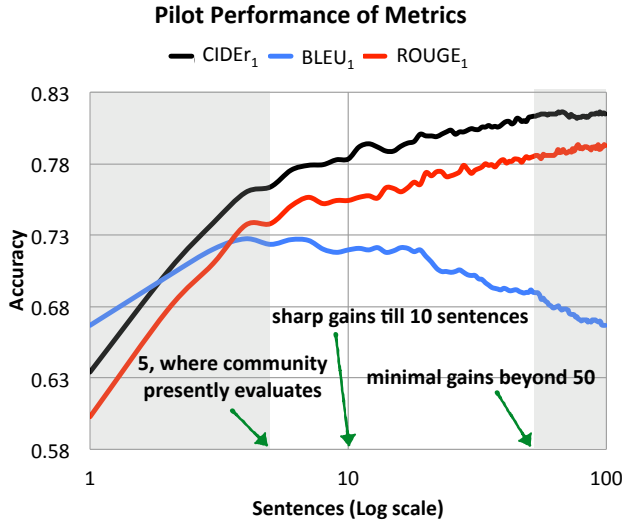
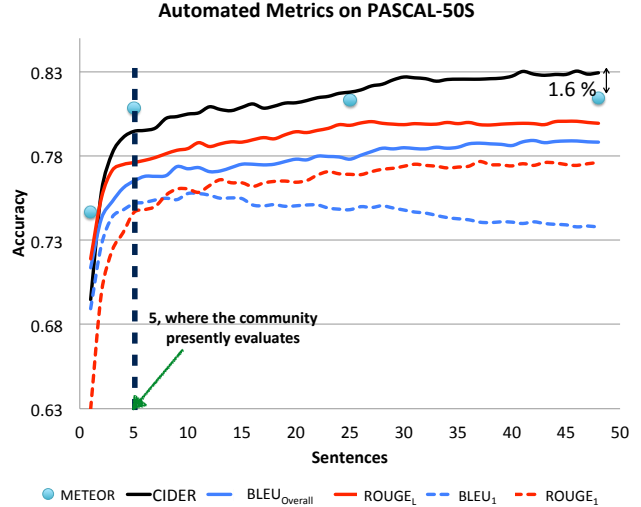


Figure 3: We show accuracy (y-axis) versus *log* number of sentences (x-axis) for our pilot study. We note that the gains saturate after 50 sentences. Our proposed metric (CIDEr<sub>1</sub>) performs better than others (BLEU<sub>1</sub> and ROUGE<sub>1</sub>).

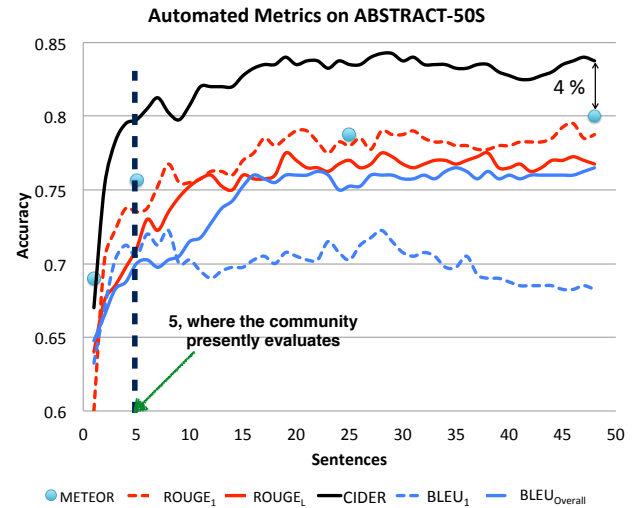
We find that CIDEr is the best performing metric, on both ABSTRACT-50S as well as PASCAL-50S. It is followed by METEOR on each dataset. Even using only 5 sentences, both CIDEr and METEOR perform well in comparison to BLEU and ROUGE. CIDEr beats METEOR at 5 sentences on ABSTRACT-50S, whereas METEOR does better at five sentences on PASCAL-50S. This is because METEOR incorporates soft-similarity, which helps when using fewer sentences. Popular metrics like ROUGE<sub>1</sub> and BLEU<sub>1</sub> are not as good at capturing consensus. METEOR, despite its sophistication does a max across reference scores, which limits its ability to utilize larger numbers of reference sentences. CIDEr provides consistent scores across both the datasets, giving 83% and 83.75% accuracy on PASCAL-50S and ABSTRACT-50S respectively. In contrast METEOR drops by 1.4% in accuracy on the semantically dense ABSTRACT-50S dataset.

Considering previous papers only used 5 reference sentences per image for evaluation, the relative boost in performance is substantial. Using BLEU or ROUGE at 5 sentences, we obtained 74% and 75% accuracy on PASCAL-50S. With CIDEr at 48 sentences, we achieve 83% accuracy. This brings automated evaluation much closer to human performance (89.92%). More details on how human performance is computed are provided in Sec. 7.4.

We next show the best performing versions of the metrics CIDEr, BLEU, ROUGE and METEOR on PASCAL-50S and ABSTRACT-50S, respectively, for different kinds of candidate pairs (Table 1). As discussed in Sec. 5 we have four kinds of pairs: (human–human correct) HC, (human–human incorrect) HI, (human–machine) HM, and (machine–machine) MM. We find that out of six cases, our



(a) PASCAL-50S



(b) ABSTRACT-50S

Figure 4: Accuracy of automated metrics (y-axis) plotted against number of reference sentences (x-axis). Metrics currently used for evaluating image descriptions are shown in *dashed* lines. Other existing metrics and our proposed metric are in *solid* lines. Note that present datasets only allow evaluation at 5 sentences. CIDEr is the best performing metric on both datasets followed by METEOR. METEOR is sampled at fewer points, due to high run-time.

proposed automated metric is best in five. We show significant gains on the challenging MM and HC tasks that involve differentiating between fine-grained differences between sentences (two machine generated sentences and two human generated sentences). This encouraging result provides evidence that the CIDEr metric will continue to perform well as automatic metrics continue to improve. On the easier tasks of judging consensus on HI and HM pairs, all methods perform well.



Metric	PASCAL-50S				ABSTRACT-50S	
	HC	HI	HM	MM	HC	HI
BLEU	62.1	97.5	93.4	62.1	61.0	92.0
ROUGE	63.9	98.5	95.1	63.9	67.5	90.0
METEOR	66.1	99.3	<b>96.3</b>	64.0	66.5	93.5
CIDEr	<b>69.5</b>	<b>99.7</b>	92.0	<b>70.7</b>	<b>71.0</b>	<b>94.5</b>

Table 1: Results on four kinds of pairs for PASCAL-50S and two kinds of pairs for ABSTRACT-50S. The best performing method is shown in **bold**. Note: we use ROUGE<sub>L</sub> for PASCAL-50S and ROUGE<sub>1</sub> for ABSTRACT-50S, and BLEU<sub>Overall</sub> version of BLEU

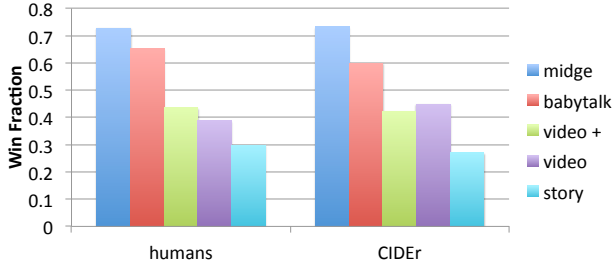


Figure 5: Fraction of times a machine generation approach wins against the other four (y-axis), plotted for human annotations and our automated metric, CIDEr.

### 7.3. Which automatic image description approaches produce consensus descriptions?

We have shown that CIDEr and our new dataset containing 50 sentences per image provides a more accurate metric over previous approaches. We now use it to evaluate existing state-of-the-art automatic image description approaches. Our methodology for conducting this experiment is described in Sec. 6. Our results are shown in Fig. 5. We show the fraction of times an approach is rated better than other approaches on the y-axis. We note that Midge [25] is rated as having the best consensus by both humans and CIDEr, followed by Babytalk [17]. Story [11] is the lowest ranked, by both humans and CIDEr. Humans and CIDEr differ on the ranking of the two video approaches (Video and Video+) [32]. We calculate the pearson’s correlation between the fraction of wins for a method on human annotations and using CIDEr. We find that humans and CIDEr agree with a high correlation (0.98). We show a sampling of the outputs from each of these image description methods in the supplementary material.

### 7.4. Human Performance

In our final set of experiments we measure human performance at predicting which of two candidate sentences better matches the consensus. Human performance puts into context how clearly consensus is defined, and provides a loose bound on how well we can expect automated metrics to perform. We evaluate both human and machine performance at

predicting consensus on all 4,000 pairs from PASCAL-50S dataset and 400 pairs from the ABSTRACT-50S dataset described in Sec. 6. To create the same experimental set up for both humans and machines, we obtain ground truth consensus for each of the pairs using our triplet annotation on 24 references out of 48. For predicting consensus, humans (via triplet annotations) and machines both use the remaining 24 sentences as reference sentences. We find that the best machine performance is 82.30 on PASCAL-50S using CIDEr, in contrast to human performance which is at 89.92. On the ABSTRACT-50S dataset, CIDEr is at 81.75 accuracy, whereas human performance is at 82.50. In terms of rank correlation, CIDEr is at 68.12 (human 79.08) and 70.44 (human 91.66) on ABSTRACT-50S and PASCAL-50S respectively. METEOR is known to have a correlation of  $\sim 52$  [9].

## 8. Discussion

The popular automated metrics in use, BLEU<sub>1</sub> and ROUGE<sub>1</sub>, are both inadequate to capture consensus. This is because BLEU was originally introduced for the machine translation task and ROUGE was introduced for the text summarization task. BLEU focuses on high precision while ROUGE prefers high recall. Sentence generation has aspects of both tasks – “translating” visual content to language and “summarizing” the contents of an image into a succinct description. We want highly salient or important objects and events to be described (good recall) while at the same time keeping the description succinct (precise), as dictated by the pragmatics of language.

Our evaluation, across five state-of-the-art approaches (Sec. 7.3) can be used as a benchmark for future image description approaches. To enable this, we will set up an evaluation server where authors can upload descriptions of images in our dataset and obtain an evaluation score using a variety of metrics. This effort can be grown in collaboration with others from the community. The problem of image description is receiving increased attention, and we believe the time is ripe for systematic benchmarks and evaluation protocols to further aid progress.

## 9. Conclusion

In this work we proposed a consensus-based evaluation protocol for image description evaluation. Our protocol enables an objective comparison of machine generation approaches based on their “human-likeness”, without having to make arbitrary calls on weighing content, grammar, saliency, *etc.* with respect to each other. We introduce an annotation modality for measuring consensus, a metric CIDEr for automatically computing consensus that outperforms existing metrics, and two datasets, PASCAL-50S and ABSTRACT-50S with 50 sentences per image. Combined, these contributions enable systematic evaluation of image description approaches.

## References

- [1] S. Banerjee and A. Lavie. Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. pages 65–72, 2005. 2, 5, 6
- [2] A. C. Berg, T. L. Berg, H. D. III, J. Dodge, A. Goyal, X. Han, A. Mensch, M. Mitchell, A. Sood, K. Stratos, and K. Yamaguchi. Understanding and predicting importance in images. In *CVPR*. IEEE, 2012. 2
- [3] R. Bogacz, E. Brown, J. Moehlis, P. Holmes, and J. D. Cohen. The physics of optimal decision making: a formal analysis of models of performance in two-alternative forced-choice tasks. *Psychol Rev*, 113(4):700–765, Oct. 2006. 1, 4
- [4] C. Callison-burch and M. Osborne. Re-evaluating the role of bleu in machine translation research. In *In EACL*, pages 249–256, 2006. 1, 2
- [5] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A Large-Scale Hierarchical Image Database. In *CVPR09*, 2009. 1
- [6] M. Denkowski and A. Lavie. Meteor universal: Language specific translation evaluation for any target language. In *Proceedings of the EACL 2014 Workshop on Statistical Machine Translation*, 2014. 6
- [7] P. K. Dokania, A. Behl, C. V. Jawahar, and P. M. Kumar. Learning to rank using high-order information. *ECCV*, 2014. 1
- [8] D. Elliott and F. Keller. Image description using visual dependency representations. In *EMNLP*, pages 1292–1302. ACL, 2013. 1, 2
- [9] D. Elliott and F. Keller. 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, June 2014. Association for Computational Linguistics. 1, 2, 5, 8
- [10] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The PASCAL Visual Object Classes Challenge 2010 (VOC2010) Results. <http://www.pascal-network.org/challenges/VOC/voc2010/workshop/index.html>. 1
- [11] A. Farhadi, M. Hejrati, M. A. Sadeghi, P. Young, C. Rashtchian, J. Hockenmaier, and D. Forsyth. Every picture tells a story: Generating sentences from images. In *Proceedings of the 11th European Conference on Computer Vision: Part IV, ECCV’10*, 2010. 2, 5, 6, 8
- [12] P. F. Felzenszwalb, R. B. Girshick, D. McAllester, and D. Ramanan. Object detection with discriminatively trained part based models. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 32(9):1627–1645, 2010. 1
- [13] A. Gupta and L. S. Davis. Beyond nouns: Exploiting prepositions and comparative adjectives for learning visual classifiers. In D. A. Forsyth, P. H. S. Torr, and A. Zisserman, editors, *ECCV (1)*, volume 5302 of *Lecture Notes in Computer Science*, pages 16–29. Springer, 2008. 2
- [14] A. Gupta, Y. Verma, and C. Jawahar. Choosing linguistics over vision to describe images. 2012. 2
- [15] M. Hodosh, P. Young, and J. Hockenmaier. Framing image description as a ranking task: Data, models and evaluation metrics. *J. Artif. Intell. Res. (JAIR)*, 47:853–899, 2013. 1, 2, 3, 5
- [16] A. Karpathy, A. Joulin, and L. Fei-Fei. Deep fragment embeddings for bidirectional image sentence mapping. *CoRR*, 2014. 2
- [17] G. Kulkarni, V. Premraj, S. Dhar, S. Li, Y. Choi, A. C. Berg, and T. L. Berg. Baby talk: Understanding and generating image descriptions. In *Proceedings of the 24th CVPR*, 2011. 1, 2, 5, 6, 8
- [18] C. H. Lampert, H. Nickisch, and S. Harmeling. Learning to detect unseen object classes by betweenclass attribute transfer. In *In CVPR*, 2009. 1, 2
- [19] S. Li, G. Kulkarni, T. L. Berg, A. C. Berg, and Y. Choi. Composing simple image descriptions using web-scale n-grams. In *Proceedings of the Fifteenth Conference on Computational Natural Language Learning*, CoNLL ’11, pages 220–228, Stroudsburg, PA, USA, 2011. Association for Computational Linguistics. 1, 2
- [20] T. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft COCO: Common objects in context. In *ECCV*, 2014. 3, 4
- [21] S. Maji, L. Bourdev, and J. Malik. Action recognition from a distributed representation of pose and appearance. In *IEEE International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2011. 1
- [22] D. Martin, C. Fowlkes, D. Tal, and J. Malik. A database of human segmented natural images and its application to evaluating segmentation algorithms and measuring ecological statistics. In *Proc. 8th Int’l Conf. Computer Vision*, volume 2, pages 416–423, July 2001. 1
- [23] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean. Distributed representations of words and phrases and their compositionality. In *Advances in Neural Information Processing Systems 26*. 2013. 4
- [24] G. A. Miller. Wordnet: A lexical database for english. *Commun. ACM*. 4
- [25] M. Mitchell, X. Han, and J. Hayes. Midge: Generating descriptions of images. In *Proceedings of the Seventh International Natural Language Generation Conference, INLG ’12*, pages 131–133, Stroudsburg, PA, USA, 2012. Association for Computational Linguistics. 1, 2, 5, 6, 8
- [26] H. Mller, P. Clough, T. Deselaers, and B. Caputo. *Image-CLEF: Experimental Evaluation in Visual Information Retrieval*. Springer Publishing Company, Incorporated, 1st edition, 2010. 3
- [27] V. Ordonez, G. Kulkarni, and T. L. Berg. Im2text: Describing images using 1 million captioned photographs. In *Neural Information Processing Systems (NIPS)*, 2011. 2, 3, 4
- [28] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu. Bleu: A method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics, ACL ’02*, pages 311–318, Stroudsburg, PA, USA, 2002. Association for Computational Linguistics. 1, 2, 5
- [29] D. Parikh and K. Grauman. Relative Attributes. In *ICCV*, 2011. 2
- [30] C. Rashtchian, P. Young, M. Hodosh, and J. Hockenmaier. 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*, CSLDAMT '10, Stroudsburg, PA, USA, 2010. Association for Computational Linguistics. [2](#), [3](#), [4](#)
- [31] S. Robertson. Understanding inverse document frequency: On theoretical arguments for idf. *Journal of Documentation*, 60:2004, 2004. [4](#)
- [32] M. Rohrbach, W. Qiu, I. Titov, S. Thater, M. Pinkal, and B. Schiele. Translating video content to natural language descriptions. In *IEEE International Conference on Computer Vision (ICCV)*, December 2013. [1](#), [2](#), [5](#), [6](#), [8](#)
- [33] M. A. Sadeghi and A. Farhadi. Recognition using visual phrases. 2011. [2](#)
- [34] D. Scharstein and R. Szeliski. A taxonomy and evaluation of dense two-frame stereo correspondence algorithms. *Int. J. Comput. Vision*, 2002. [1](#)
- [35] A. Sorokin and D. Forsyth. Utility data annotation with amazon mechanical turk. In *Computer Vision and Pattern Recognition Workshops, 2008. CVPRW '08. IEEE Computer Society Conference on*, pages 1–8, June 2008. [1](#)
- [36] O. Tamuz, C. Liu, S. Belongie, O. Shamir, and A. T. Kalai. Adaptively learning the crowd kernel. In *In ICML11*, 2011. [1](#)
- [37] Y. Yang, C. L. Teo, H. D. III, and Y. Aloimonos. Corpus-guided sentence generation of natural images. In *EMNLP. ACL*, 2011. [1](#), [2](#)
- [38] M. Yatskar, M. Galley, L. Vanderwende, and L. Zettlemoyer. 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)*, page 110120, Dublin, Ireland, August 2014. Association for Computational Linguistics and Dublin City University. [1](#), [2](#)
- [39] C. yew Lin. Rouge: a package for automatic evaluation of summaries. pages 25–26, 2004. [1](#), [2](#), [5](#), [6](#)
- [40] C. L. Zitnick and D. Parikh. Bringing semantics into focus using visual abstraction. In *CVPR*, 2013. [2](#), [3](#), [4](#)