

Transparent Human Evaluation for Image Captioning

Anonymous ACL submission

Abstract

We establish a rubric-based human evaluation protocol for image captioning models. Our scoring rubrics and their definitions are carefully developed based on machine- and human-generated captions on the MSCOCO dataset. Each caption is evaluated along two main dimensions in a tradeoff (*precision* and *recall*) as well as other aspects that measure the text quality (*fluency*, *conciseness*, and *inclusive language*). Our evaluations demonstrate several critical problems of the current evaluation practice. Human-generated captions show substantially higher quality than machine-generated ones, especially in coverage of salient information (i.e., recall), while most automatic metrics say the opposite. Our rubric-based results reveal that CLIPScore, a recent metric that uses image features, better correlates with human judgments than conventional text-only metrics because it is more sensitive to recall. We hope that this work will promote a more transparent evaluation protocol for image captioning and its automatic metrics.

1 Introduction

Recent progress in large-scale training has pushed the state of the art in vision-language tasks (Li et al., 2020; Zhang et al., 2021, *inter alia*). One of these tasks is image captioning, whose objective is to generate a caption that describes the given image. The performance in image captioning has been primarily measured in automatic metrics (e.g., CIDEr, Vedantam et al., 2015; SPICE, Anderson et al., 2016) on popular benchmarks, such as MSCOCO (Lin et al., 2014) and Flickr8k (Hodosh et al., 2013). Use of these metrics is justified based on their correlations with human judgments collected in previous work (Hodosh et al., 2013; Elliott and Keller, 2014; Kilickaya et al., 2017, *inter alia*).

Continuous use of these previous human judgments, however, raises significant concerns for development of both captioning models and auto-



Machines	P	R	CIDEr
<i>A red fire hydrant spewing water on a street.</i>	5	3	139.2
<i>A red fire hydrant spraying water on a street.</i>	5	3	205.2
Human			
<i>A busted red fire hydrant spewing water all over a street creating a rainbow.</i>	5	5	120.5

Figure 1: These machine captions are *precise* (in the scale of 1–5) but lose points in recall (i.e., coverage of salient information); they both ignore the rainbow in the picture. Automatic metrics, such as CIDEr, do not capture this failure.

matic metrics because of their lack of transparency. In previous work, annotators (crowdworkers, typically) rate image captions directly (Hodosh et al., 2013), pairwise (Vedantam et al., 2015), or along multiple dimensions such as thoroughness (Aditya et al., 2015) and truthfulness (Yatskar et al., 2014). These scoring judgments depend highly on individual annotators’ discretion and understanding of the annotation scheme (Freitag et al., 2021; Clark et al., 2021), making it difficult to decompose, interpret, and validate annotations. This lack of transparency also makes it difficult to interpret evaluation results for downstream applications where some aspects are particularly important (e.g., accessibility for people with visual impairments; Gleason et al., 2019, 2020). Further, these annotations were done only on relatively old models (e.g., MSCOCO leaderboard submissions in 2015; Anderson et al., 2016). Correlations of automatic metrics with human judgments can break down especially when model types change (Callison-Burch et al., 2006),

063 or generation models become increasingly pow- 110
064 erful (Ma et al., 2019; Edunov et al., 2020). We 111
065 thus develop an up-to-date, transparent human eval- 112
066 uation protocol to better understand how current 113
067 models perform and how automatic metrics are 114
068 correlated when applied to current models. 115

069 At the core of our rubrics are two main scores 116
070 in a tradeoff: *precision* and *recall* (Fig. 1). The 117
071 former measures accuracy of the information in a 118
072 caption, and the latter assesses how much of the 119
073 salient information in the image is covered. We 120
074 then penalize a caption if we find a problem in *flu-* 121
075 *ency*, *conciseness*, or *inclusive language*. Two or 122
076 more authors evaluate every instance and collabo- 123
077 rate to resolve disagreements, ensuring high quality 124
078 of the annotations. We assess outputs from four 125
079 strong models as well as human-generated refer- 126
080 ence captions from MSCOCO. We call our scores 127
081 THUMB 1.0 (Transparent **H**uman **B**enchmark), 128
082 and release them publicly.¹ Our key findings from 129
083 the evaluations are: 130

- 084 • Machine-generated captions from recent models 131
085 have been claimed to achieve superhuman perfor- 132
086 mance using popular automatic metrics (human 133
087 performance is ranked at the 250th place in the 134
088 MSCOCO leaderboard),² but they still show sub- 135
089 stantially lower quality than human-generated 136
090 ones. 137
- 091 • Machines fall short of humans, especially in re- 138
092 call (Fig. 1), but most automatic metrics say the 139
093 opposite. 140
- 094 • Human performance is underestimated in the cur- 141
095 rent leaderboard paradigm, and there is still much 142
096 room for improvement on MSCOCO captioning. 143
- 097 • CLIPScore and RefCLIPScore (Hessel et al., 144
098 2021), recently proposed metrics that use im- 145
099 age features, improve correlations particularly in 146
100 recall. While they fail to score human generation 147
101 much higher than machine one, they capture an 148
102 aspect that is less reflected in text-only metrics. 149
- 103 • Currently available strong captioning models gen- 150
104 erate highly fluent captions. Fluency evaluation 151
105 is thus no longer crucial in ranking these models. 152

106 2 Evaluation Protocol 153

107 We establish a transparent evaluation protocol for 154
108 image captioning models. Our rubrics and rules are 155
109 developed through discussions among all annota-

¹Anonymized.

²<https://competitions.codalab.org/competitions/3221#results>.

tors (first four authors of this paper).

111 2.1 Evaluation Setups and Quality Control 112

112 We used images from the test data in the stan- 113
114 dard *Karpathy split* (Karpathy and Fei-Fei, 2015) 115
116 of the MSCOCO dataset (Lin et al., 2014). The 117
118 dataset consists of 113K, 5K, and 5K train/dev/test 119
120 everyday-scene photos sampled from Flickr. We 121
122 randomly sampled 500 test images and prepared 123
124 one human- and four machine-generated captions 125
126 for every image (§2.3). We first performed de- 127
128 velopmental evaluations of 250 captions for 50 129
130 images and created rubrics. We then proceeded 131
132 with the rest of the captions. For every image, cap- 133
134 tions were shuffled, and thus annotators did not 134
135 know which caption corresponded to which model, 135
136 thereby avoiding a potential bias from knowledge 136
137 about the models. We conducted two-stage anno- 137
138 tations: the first annotator scores all captions for 138
139 given images, and the second annotator checks and 139
140 modifies the scores when necessary. After the de- 140
141 velopmental phase, the κ coefficient (Cohen, 1960) 141
142 was 0.86 in precision and 0.82 in recall for the 142
143 rest of the evaluated captions (§2.2.1).³ The first 143
144 four authors of this paper conducted all evaluations; 144
145 none of them are color blind or low vision, two are 145
146 native English speakers, and one is a graduate stu- 146
147 dent in linguistics. We finally ensured that at least 147
148 one native speaker evaluated the fluency of every 148
149 caption (§2.2.2), meaning that if a caption is anno- 149
150 tated by the two non-native speakers, one native 150
151 speaker checks the fluency in an additional round. 151

141 2.2 THUMB 1.0 142

142 We base our evaluations on two main scores 143
143 (**precision** and **recall**) and three types of penalty 144
144 (**fluency**, **conciseness**, and **inclusive language**) 145
145 The overall score is computed by averaging pre- 146
146 cision and recall and deducting penalty points. 146

147 2.2.1 Main Scores 148

148 The two main scores are assessed in the scale of 149
149 1–5. They balance information accuracy and cover- 150
150 age. See §A.4 for score distribution histograms. 150

151 **Precision** Precision (P) measures how precise the 151
152 caption is given the image. For instance, Caption 1- 152
153 B in Table 1 is perfectly precise, while 1-A (*dog vs.* 153
154 *otter, one vs. two frisbees*) and 1-C (*three vs. two* 154

³Furthermore, we found that a third annotator did not change the scores for all 100 captions randomly sampled for meta-evaluations, confirming the sufficiently high quality of our two-stage annotations.






Image	Caption	P	R	Flu.	Total
	1-A: Up-Down <i>A dog playing with a frisbee on the ground.</i>	3	4	0	3.5
	1-B: VinVL-base <i>A otter is laying on the sand next to two frisbees.</i>	5	4	0.1	4.4
	1-C: VinVL-large <i>A small animal laying on a rock with three frisbees.</i>	4	3	0	3.5
	2-A: Up-Down <i>A close up of a plate of broccoli.</i>	5	3	0	4
	2-B: Unified-VLP, VinVL-base, VinVL-large <i>A plate of pasta and broccoli on a table.</i>	4	4	0	4
	2-C: Human <i>A multi colored dish with broccoli and white twisted pasta in it.</i>	5	5	0.1	4.9
	3-A: Unified-VLP <i>A little girl holding a video game controller.</i>	3	4	0	3.5
	3-B: VinVL-large <i>A little girl is blow drying her hair on a couch.</i>	4	5	0	4.5
	3-C: Human <i>A little girl holding a blow dryer next to her head.</i>	5	5	0	5
	4-A: Up-Down <i>A black cat laying in a red suitcase.</i>	3	5	0	4
	4-B: Unified-VLP, VinVL-base, VinVL-large <i>A black cat sitting on top of a red suitcase.</i>	5	5	0	5
	4-C: Human <i>A large black cat laying on top of a pink piece of luggage.</i>	4	5	0	4.5
	5-A: Up-Down, Unified-VLP <i>A man standing in front of a display of donuts.</i>	3	2	0	2.5
	5-B: VinVL-large <i>A woman standing behind a counter at a donut shop.</i>	5	3	0	4
	5-C: Human <i>Woman selling doughnuts with doughnut stock in the background.</i>	5	5	0.3	4.7

Table 1: Example evaluations of machine- and human-generated captions. None of these captions get penalties in conciseness and inclusive language. Evaluated captioning models are described in §2.3

frisbees) are not precise. Precision guards against hallucinations from the language model (*table* in 2-B) that are known to be common failures of image captioning models (Rohrbach et al., 2018). The score of 4 is reserved for relatively minor issues, such as attributes that are almost correct (e.g., *pink* vs. *red* in 4-C, Table 1) or cases where the caption does not contradict with the image but is not guaranteed to be true (e.g., it is unclear whether the girl is sitting on a couch in 3-B). In addition to objects themselves, precision deals with information like properties, attributes, occasions, locations, and relations between objects (e.g., *in a red suitcase* vs. *on a red suitcase* in 4-A).

Recall Recall (R) measures how much of the salient information (e.g., objects, attributes, and relations) from the image is covered by the caption. This includes color (e.g., color of the frisbees in 1-A, 1-B, and 1-C) and guards against generic,

uninformative captions. For instance, an otter is *a small animal*, and thus *small animal* is *precise* (1-C); however, it is much less informative than saying an otter. Similarly, Caption 5-B only says a woman is standing behind a counter at a donut shop, but she is selling donuts, not buying or looking at donuts, which is salient information from the picture. We do not take a point off if missing information is already expected from the caption (e.g., a double-decker bus is typically red). We often find it useful to take a generative approach when evaluating recall: *what image does the caption lead us to imagine?* When the caption entails many potential images that substantially diverge from the given image, the recall score should be low.

2.2.2 Penalties

Fluency Fluency (Flu.) measures the quality of captions as English text regardless of the given im-

age. Initially, we scored fluency in the scale of 1–5, similar to P and R, but we found most captions from modern neural network models were highly fluent. Thus, we instead decided to take points off from the average of P and R if there’s a fluency problem to account for minor issues that are much less problematic than losing one P/R point. The four annotators had extensive discussions and developed rubrics for fluency. Similar to recent work on professional evaluations for machine translation (Freitag et al., 2021), we evaluated under the following principle: if a fluency problem is expected to be easily corrected by a text postprocessing algorithm (e.g., grammatical error correction: Yuan and Briscoe, 2016; Sakaguchi et al., 2017), the penalty should be 0.1. This includes obvious misspellings or grammatical errors (e.g., *A otter* in 1-B) and missing determiners/hyphens (*multi colored* in 2-C). 0.5+ points were subtracted for more severe problems, such as duplication (e.g., *A display case of donuts and doughnuts*), ambiguity (e.g., *A cat is on a table with a cloth on it*), and broken sentences (e.g., *A large concrete sign small buildings behind it*). See Table 6 in §A.1 for more extensive fluency rubrics. Note that the average fluency penalty was 0.01; this confirms that fluency is no longer crucial in ranking models for MSCOCO captioning and contrasts with human evaluations previously done for older captioning models.

Conciseness The scores so far do not take into account conciseness of captions. Specifically, a model could simply increase all scores by describing every detail in a picture. For instance, the following caption is overly repetitive: *a woman lying on her back with knees bent on a beach towel under a multicolored, striped beach umbrella, surrounded by sand, and with clear blue sky above*. We subtract 0.5 points for these captions. Note that most machine captions were short, and this penalty was only applied to two human-generated captions. It might become more crucial for future models with a more powerful object detection module that catches many objects in the picture.

Inclusive Language We found that some instances substantially diverge from inclusive language, raising a concern for downstream applications. In these cases, we added a penalty: 0.5 points were deducted for a subjective comment about appearance (e.g., *very pretty girl*), and 2 points for more severe problems (e.g., *beautiful breasts*).

2.2.3 Rules of THUMB

In our development phase, we established the following additional rules to clarify our annotation scheme.

Avoiding Double Penalties When an error is accounted for in precision, we correct the error before scoring the recall, thereby avoiding penalizing the precision and recall for the same mistake. For example, P=3 is given to Caption 1-A in Table 1 because of its wrong detection (*dog* vs. *otter*; *one* vs. *two frisbees*), but we score the recall assuming that the caption is now *an otter playing with two frisbees on the ground*. This ensures that a generic, useless caption, such as *there is something on something* (P=5, R=1), would be ranked considerably lower than *a dog on the beach with two pink and yellow frisbees* (P=3, R=5). Similarly, the wrong detection in 5-A (*man* vs. *woman*) is handled only in precision. Note that such error correction is not applicable to hallucinations because there is no alignment between a part of the image and a hallucinated object (e.g., *table* in 2-B). This rule departs from the definition of recall in SPICE (Anderson et al., 2016), an automatic metric that measures the F_1 score in scene graphs predicted from reference and generated captions; their alignment is limited to WordNet synonyms (Miller, 1995). This means that classifying an otter as a dog or even a small animal would result in cascading errors both in precision and recall, overrating captions that completely overlook the otter or ones that make a more severe classification error (e.g., miscategorize the otter as a car, compared to a dog).

Object Counts as Attributes All counts are considered as object attributes, and wrong counts are handled in precision. This simplifies the distinction between precision and recall. For instance, both *a frisbee* (1-A) and *three frisbees* (1-C) are precision problems, while saying *some frisbees* would be a recall problem when it is clear that there are exactly two frisbees. Note that this is in line with SPICE, which treats object counts as attributes in a scene graph, rather than duplicating a scene graph for every instance of an object (Anderson et al., 2016).

Black and White Photo MSCOCO contains black and white or gray-scale pictures. Some captions explicitly mention that they are black and white, but we disregard this difference in our evaluations. The crowdsourcing instructions for creating reference captions do not specify such cases (Chen

et al., 2015). Further, we can potentially run post-processing to determine whether it is black and white to modify the caption accordingly, depending on the downstream usage.

Text Processing Image captioning models often differ slightly in text preprocessing. As a result, we found that generated captions were sometimes slightly different in format (e.g., tokenized or detokenized; lowercased or not). For better reproducibility, we follow the spirit of SACREBLEU (Post, 2018), which has become the standard package to compute BLEU scores for machine translation: all evaluations, including automatic metrics, should be done on clean, untokenized text, independently of preprocessing design choices. We apply the following minimal postprocessing to the model outputs and human captions.

- Remove unnecessary spaces at the start or end of every caption.
- Uppercase the first letter.
- Add a period at the end if it doesn't exist, and remove a space before a period if any.

We keep the postprocessing minimal for this work and encourage future model developers to follow the standard practice in machine translation: every model has to output clean, truecased, untokenized text that is ready to be used in downstream modules. This also improves the transparency and reproducibility of automated evaluations (Post, 2018).

2.3 Evaluated Captions

We evaluated the following four strong models from the literature as well as human-generated captions. They share similar pipeline structure: object detection followed by crossmodal caption generation. They vary in model architecture, (pre)training data, model size, and (pre)training objective. Evaluating captions from them will enable us to better understand what has been improved and what is still left to future captioning models.

- **Up-Down** (Anderson et al., 2018) trains Faster R-CNN (Ren et al., 2015) on the Visual Genome dataset (Krishna et al., 2016) for object detection. It then uses an LSTM-based crossmodal generation model.
- **Unified-VLP** (Zhou et al., 2020) uses the same object detection model as Up-Down. The transformer-based generation model is initialized with base-sized BERT (Devlin et al., 2019) and further pretrained with 3M images from Conceptual Captions (Sharma et al., 2018).

- **VinVL-base** and **VinVL-large** (Zhang et al., 2021) train a larger-scale object detection model with the ResNeXt-152 C4 architecture (Xie et al., 2017) on ImageNet (Deng et al., 2009). The transformer generation model is initialized with BERT and pretrained with 5.7M images.

- **Human** randomly selects one from the five human-generated reference captions in MSCOCO. Those captions were created by crowdworkers on Amazon Mechanical Turk (Chen et al., 2015).

Further details are described in §A.3 of Appendix.

3 Results and Analysis

We present results and analysis from our evaluations. Our transparent evaluations facilitate assessments and analysis of both captioning models (§3.1) and automatic metrics (§3.2).

3.1 Comparing Models

Seen in Table 2 (left section) is the model performance that is averaged over the 500 test images and broken down by the rubric categories. Overall, Human substantially outperforms all machines in the P, R, and total scores. In particular, we see a large gap between Human and the machines in recall (e.g., Human 4.35 vs. VinVL-large 3.97). This contrasts with the automatic metric-based ranking of the MSCOCO leaderboard, where Human is ranked at the 250th place.⁴ This result questions claims about human parity or superhuman performance on MSCOCO image captioning. The four machine captioning models are ranked in the expected order, though the small difference between VinVL-large and VinVL-base suggests that simply scaling up models would not lead to a substantial improvement. We see that the three models that are initialized with pretrained BERT (VinVL-large/base, Unified-VLP) are particularly fluent, but the problem is small in the other models as well.

While we compute representative, total scores, our transparent rubrics allow for adjusting weighting of the categories depending on the application of interest. For instance, in the social media domain, recall can be more important than precision to make captions engaging to users (Shuster et al., 2019). To assess the models indepen-

⁴The official leaderboard ranks submissions using CIDEr (Vedantam et al., 2015) with 40 references on the hidden test data. We use the public Karpathy test split instead, but we suspect the same pattern would hold on the hidden data as well, given the large gap between machines and Human.

Model	THUMB 1.0						Automatic Metrics						
	P↑	R↑	Flu.↓	Con.↓	Inc.↓	Total↑	BLEU	ROUGE	BERT-S	SPICE	CIDEr	CLIP-S	RefCLIP-S
Human	4.82	4.35	0.019	0.02	0.00	4.56 ^{+0.03} _{-0.03}	26.2	50.4	0.938	23.7	111.5	0.791	0.834
VinVL-large	4.54	3.97	0.005	0.00	0.00	4.25 ^{+0.04} _{-0.04}	33.3	56.5	0.946	26.4	141.8	0.784	0.834
VinVL-base	4.47	3.95	0.001	0.00	0.00	4.21 ^{+0.04} _{-0.04}	32.3	55.9	0.945	25.6	138.4	0.779	0.830
Unified-VLP	4.35	3.77	0.004	0.00	0.00	4.06 ^{+0.04} _{-0.04}	31.6	55.8	0.945	24.3	128.5	0.771	0.821
Up-Down	4.29	3.50	0.014	0.00	0.00	3.88 ^{+0.05} _{-0.05}	28.4	52.2	0.939	21.0	110.7	0.746	0.803

Table 2: Performance of image captioning models with respect to THUMB 1.0 (left) and automatic metrics (right). All scores are averaged over 500 images randomly sampled from the Karpathy test split. P: precision; R: recall; Flu.: fluency; Con.: conciseness; Inc.: inclusive language. 90% confidence intervals for total scores are calculated by bootstrapping (Koehn, 2004). All reference-based metrics take as input the same four crowdsourced captions that are not used in Human for fair comparisons.

dently of these aggregation decisions, we count the number of times when each model outperforms/underperforms all the others both in P and R (*strictly* best/worst, Table 3). We see patterns consistent with Table 2. For example, Human is most likely to be strictly best and least likely to be strictly worst. This suggests that machine captioning models would still fall short of crowdworkers in a wide range of downstream scenarios.

Model	Human	Vin-large	Vin-base	U-VLP	Up-Down
# Best ↑	327	180	161	112	74
# Worst ↓	65	128	150	190	269

Table 3: # times when each captioning model is *strictly* best/worst in the caption set (i.e., best/worst both in precision and recall).

Metric	w/o Human			w/ Human		
	P	R	Total	P	R	Total
RefCLIP-S	0.34	0.27	0.44	0.31	0.26	0.41 ^{+0.05} _{-0.05}
RefOnlyC	0.42	0.14	0.41	0.37	0.11	0.34 ^{+0.04} _{-0.05}
CLIP-S	0.18	0.27	0.32	0.17	0.28	0.32 ^{+0.05} _{-0.05}
CIDEr	0.27	0.18	0.33	0.21	0.11	0.23 ^{+0.04} _{-0.04}
BERT-S	0.27	0.18	0.33	0.20	0.10	0.21 ^{+0.04} _{-0.04}
SPICE	0.26	0.15	0.30	0.20	0.09	0.21 ^{+0.04} _{-0.04}
ROUGE-L	0.26	0.17	0.31	0.18	0.07	0.18 ^{+0.04} _{-0.04}
BLEU	0.21	0.13	0.25	0.15	0.04	0.13 ^{+0.04} _{-0.04}

Table 4: Instance-level correlations of automatic evaluation scores. RefCLIP-S and CLIP-S use image features unlike the others, and all but CLIP-S require references. All of these reference-based metrics use the same subset of four captions as in Table 2 that exclude Human. All metrics had correlations lower than 0.1 for fluency.

3.2 Comparing Automatic Metrics

While carefully-designed human judgments like ours should be considered more reliable, automatic

metrics allow for faster development cycles. Our transparent evaluations can also be used to analyze how these automatic metrics correlate with different aspects of image captioning. Table 2 (right section) shows automatic scores of the captioning models over 7 popular metrics for image captioning. CLIP-S (Hessel et al., 2021) is a *referenceless* metric that uses image features from CLIP (Radford et al., 2021), a crossmodal retrieval model trained on 400M image-caption pairs from the web. RefCLIP-S augments CLIP-S with similarities between the generated and reference captions. All other metrics, such as SPICE (Anderson et al., 2016) and CIDEr (Vedantam et al., 2015), only use reference captions without image features.

These automatic metrics generally agree with our evaluations in ranking the four machines, but completely disagree in the assessment of Human. Most metrics rank Human near the bottom, showing that they are not reliable in evaluating high-quality, human-generated captions. The two metrics with powerful image and text features (CLIP-S and RefCLIP-S) give high scores to Human compared to the other metrics, but they still fail to score Human substantially higher than VinVL-large. This suggests that automatic metrics should be regularly updated as our models become stronger (and perhaps more similar to humans), and raises a significant concern about the current practice that fixes evaluation metrics over time.

Seen in Table 4 are instance-level Pearson correlation scores between automatic scores and our evaluations.⁵ We also add an ablation study: RefOnlyC removes image features from RefCLIP-S to

⁵Instance-level Pearson correlations with human judgments were often computed in prior work to compare automatic metrics for image captioning (e.g., Hessel et al., 2021). An alternative is system-level correlations, but they would be uninformative with five systems only.

quantify the effect of image features. We consider two types of scenarios: one *with* Human and one *without*. Correlations drop from the latter to the former for all metrics and aspects except CLIP-S, again showing that the metrics are not reliable in assessing human-generated captions. Interestingly, CLIP-S correlates best in recall (0.28 w/ Human) but suffers in precision (0.17 w/ Human). RefOnlyC, in contrast, achieves the best correlations in P at the expense of R. RefCLIP-S balances the two and achieves the best correlation in total scores. This indicates that the CLIP image features particularly help assess coverage of salient information that can be ignored in some reference captions from crowdworkers.⁶ Prior work (Hessel et al., 2021) found that SPICE can still improve correlations when combined with CLIP-S, even though CLIP-S better correlates with human judgments than SPICE. This implies that image-based and reference-only metrics capture different aspects of image captioning. Our analysis indeed agrees with their finding and, further, identifies that recall is one such aspect. For an extensive description of these metrics and their configurations, see §A.2 of Appendix.

3.3 Machine vs. Human Examples

Table 5 provides examples that contrast machine- and human-generated captions. We see that machine-generated captions ignore salient information or make critical errors for these images. These problems often occur in relatively rare cases: a tennis player is showing excitement rather than hitting a ball; a bride and groom are cutting a wedding cake; a boy is wearing a tie without a shirt; a man is putting clothing and a tie on a dummy instead of a person. But these situations are exactly the most important information because of their *atypicality* (Feinglass and Yang, 2021). This illustrates fundamental problems of current image captioning models that are left to future work.

4 Related Work

Human Evaluations for Image Captioning Several prior works conducted human evaluations for

⁶The low recall correlations of reference-only metrics can be partly because the maximum (as opposed to minimum or average) is typically taken over multiple reference captions (e.g., BERTScore, Zhang et al., 2020). Nevertheless, this alone does not explain the recall gap from image-based metrics because RefCLIP-S also takes the maximum score over all references. Future work can explore the relation between precision/recall and different treatments of multiple references.

image captioning with varying models, datasets, and annotation schemes. Much work used crowdworkers from Amazon Mechanical Turk on Flickr-based datasets, including the PASCAL (Rashtchian et al., 2010), Flickr8k/30k (Hodosh et al., 2013; Young et al., 2014), and MSCOCO datasets. Annotators scored the overall quality directly (Kulkarni et al., 2011; Hodosh et al., 2013), pairwise (Vedantam et al., 2015), or along multiple dimensions, such as truthfulness/correctness (Yatskar et al., 2014; Anderson et al., 2016), thoroughness (Aditya et al., 2015), relevance (Yang et al., 2011; Li et al., 2011), and grammaticality/readability (Mitchell et al., 2012; Elliott and Keller, 2013). There are similarities between our rubrics and previous annotations, but our framework defines every dimension in a decomposable way through discussions among all annotators, while focusing on outputs from strong models currently available. Apart from these conventional Flickr-based datasets, some other work evaluated image captions for social media (engagingness, Shuster et al., 2019; accessibility for Twitter users with vision impairments, Gleason et al., 2019, 2020) and news articles (Biten et al., 2019). Our transparent evaluations would enable us to adjust the aggregation method based on the nature of downstream applications. More specializing categories can be added for these applications in later versions (e.g., THUMB 2.0).

Human Evaluations for Other Generation Tasks

Much previous work explored human evaluations for other language generation tasks than image captioning. The WMT shared task (Barrault et al., 2020) conducts human evaluations of state-of-the-art machine translation systems every year; participants or crowdworkers directly rate a translation in a 100-point scale, which is a method developed by Graham et al. (2013, 2014, 2017). GENIE takes a similar approach but hosts human evaluations in leaderboards for machine translation, summarization, and commonsense reasoning (Khashabi et al., 2021). Kryscinski et al. (2019) and Fabbri et al. (2021) assessed many summarization models in a similar annotation scheme to the DUC 2006/2007 evaluations (Dang, 2006). Our transparent evaluation framework is inspired by rubric-based machine translation judgments by professional translators (Freitag et al., 2021), which resulted in different system rankings than the WMT evaluations. As top-performing models and automatic metrics are becoming increasingly similar across various natu-



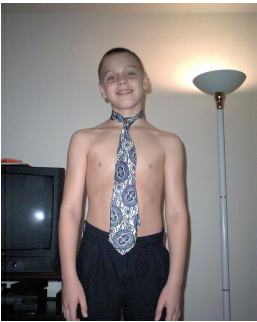

Image	Caption	P	R	Flu.	Total
	6-A: Up-Down <i>A man holding a tennis racquet on a tennis court.</i>	5	3	0	4
	6-B: Unified-VLP, VinVL-base, VinVL-large <i>A man holding a tennis racket on a tennis court.</i>	5	3	0	4
	6-C: Human <i>A tennis player shows controlled excitement while a crowd watches.</i>	5	5	0	5
	7-A: Up-Down <i>A person cutting a cake with a knife.</i>	3	3	0	3
	7-B: Unified-VLP <i>A person cutting a wedding cake with a knife.</i>	3	5	0	4
	7-C: VinVL-base <i>A couple of cakes on a table with a knife.</i>	5	3	0	4
	7-D: VinVL-large <i>A woman cutting a cake with a knife.</i>	3	3	0	3
	7-E: Human <i>Bride and grooms arms cutting the wedding cake with fruit on top.</i>	5	5	0.1	4.9
	8-A: Up-Down <i>A young boy wearing a blue shirt and a blue tie.</i>	3	3	0	3
	8-B: Unified-VLP <i>A young boy wearing a shirt and a tie.</i>	3	3	0	3
	8-C: VinVL-base <i>A young boy wearing a tie standing in front of a lamp.</i>	5	3	0	4
	8-D: VinVL-large <i>A young man wearing a tie and a shirt.</i>	3	3	0	3
	8-E: Human <i>A man wearing only a tie standing next to a lamp.</i>	4	5	0	4.5
	9-A: Up-Down <i>A couple of men standing next to each other.</i>	2	2	0	2
	9-B: Unified-VL <i>Two men standing in a room.</i>	2	2	0	2
	9-C: VinVL-base <i>A couple of men standing in a room.</i>	2	2	0	2
	9-D: VinVL-large <i>Two men standing next to each other in a room.</i>	2	2	0	2
	9-E: Human <i>A man standing next to a dummy wearing clothes.</i>	5	3	0	4

Table 5: Examples that contrast machine- and human-generated captions. All machine-generated captions overlook or misinterpret salient information: the excitement the tennis player expresses, the bride and groom cutting a wedding cake, the boy not wearing a shirt, and the man putting a tie on a dummy. None of these captions are penalized for conciseness or inclusive language. See §A.5 in Appendix for more examples.

527 ral language generation tasks, our findings on im- 536
528 age captioning may be useful for other generation 537
529 tasks as well. 538

530 5 Conclusion 539

531 We developed THUMB 1.0, transparent evalua- 540
532 tions for the MSCOCO image captioning task. We 541
533 refined our rubrics through extensive discussions 542
534 among all annotators, and ensured the high quality 543
535 by two-stage annotations. Our evaluations demon- 544
545

536 strated critical limitations of current image cap- 537
538 tioning models and automatic metrics. While re- 538
539 cent image-based metrics show promising improve- 539
540 ments, they are still unreliable in assessing high- 540
541 quality captions from crowdworkers. We hope that 541
542 our annotation data will help future development 542
543 of better captioning models and automatic metrics, 543
544 and this work will become a basis for transparent 544
545 human evaluations for the image captioning task 545
and beyond.

546
547
548
549
550

551
552
553

554
555
556
557
558

559
560
561
562
563
564
565
566
567

568
569
570
571

572
573
574

575
576
577
578

579
580
581
582

583
584
585

586
587

588
589
590

591
592
593
594

595
596
597
598

References

Somak Aditya, Yezhou Yang, Chitta Baral, Cornelia Fermüller, and Yiannis Aloimonos. 2015. [From images to sentences through scene description graphs using commonsense reasoning and knowledge.](#)

Peter Anderson, Basura Fernando, Mark Johnson, and Stephen Gould. 2016. [SPICE: semantic propositional image caption evaluation.](#) In *Proc. of ECCV*.

Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang. 2018. [Bottom-up and top-down attention for image captioning and visual question answering.](#) In *Proc. of CVPR*.

Loïc Barrault, Magdalena Biesialska, Ondřej Bojar, Marta R. Costa-jussà, Christian Federmann, Yvette Graham, Roman Grundkiewicz, Barry Haddow, Matthias Huck, Eric Joanis, Tom Kocmi, Philipp Koehn, Chi-kiu Lo, Nikola Ljubešić, Christof Monz, Makoto Morishita, Masaaki Nagata, Toshiaki Nakazawa, Santanu Pal, Matt Post, and Marcos Zampieri. 2020. [Findings of the 2020 conference on machine translation \(WMT20\).](#) In *Proc. of WMT*.

Ali Furkan Biten, Lluís Gómez, Marçal Rusiñol, and Dimosthenis Karatzas. 2019. [Good news, everyone! context driven entity-aware captioning for news images.](#) In *Proc. of CVPR*.

Chris Callison-Burch, Miles Osborne, and Philipp Koehn. 2006. [Re-evaluating the role of BLEU in machine translation research.](#) In *Proc. of EACL*.

Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollár, and C. Lawrence Zitnick. 2015. [Microsoft COCO captions: Data collection and evaluation server.](#)

Elizabeth Clark, Tal August, Sofia Serrano, Nikita Haduong, Suchin Gururangan, and Noah A. Smith. 2021. [All that’s ‘human’ is not gold: Evaluating human evaluation of generated text.](#) In *Proc. of ACL*.

Jacob Cohen. 1960. [A coefficient of agreement for nominal scales.](#) *Educational and Psychological Measurement*.

Hoa Trang Dang. 2006. [Overview of DUC 2006.](#) In *Proc. of DUC*.

Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. [ImageNet: A large-scale hierarchical image database.](#) In *Proc. of CVPR*.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding.](#) In *Proc. of NAACL*.

Sergey Edunov, Myle Ott, Marc’Aurelio Ranzato, and Michael Auli. 2020. [On the evaluation of machine translation systems trained with back-translation.](#) In *Proc. of ACL*.

Desmond Elliott and Frank Keller. 2013. [Image description using visual dependency representations.](#) In *Proc. of EMNLP*.

Desmond Elliott and Frank Keller. 2014. [Comparing automatic evaluation measures for image description.](#) In *Proc. of ACL*.

Alexander R Fabbri, Wojciech Kryściński, Bryan McCann, Caiming Xiong, Richard Socher, and Dragomir Radev. 2021. [SummEval: Re-evaluating summarization evaluation.](#) *TACL*.

Joshua Feinglass and Yezhou Yang. 2021. [SMURF: SeMantic and linguistic UndeRstanding fusion for caption evaluation via typicality analysis.](#) In *Proc. of ACL*.

Markus Freitag, George Foster, David Grangier, Viresh Ratnakar, Qijun Tan, and Wolfgang Macherey. 2021. [Experts, errors, and context: A large-scale study of human evaluation for machine translation.](#)

Cole Gleason, Patrick Carrington, Cameron Tyler Cassidy, Meredith Ringel Morris, Kris M. Kitani, and Jeffrey P. Bigham. 2019. ["It’s almost like they’re trying to hide it": How user-provided image descriptions have failed to make Twitter accessible.](#) In *Proc. of WWW*.

Cole Gleason, Amy Pavel, Emma McCamey, Christina Low, Patrick Carrington, Kris M. Kitani, and Jeffrey P. Bigham. 2020. [Twitter a11y: A browser extension to make twitter images accessible.](#) In *Proc. of CHI*.

Yvette Graham, Timothy Baldwin, Alistair Moffat, and Justin Zobel. 2013. [Continuous measurement scales in human evaluation of machine translation.](#) In *Proc. of LAW*.

Yvette Graham, Timothy Baldwin, Alistair Moffat, and Justin Zobel. 2014. [Is machine translation getting better over time?](#) In *Proc. of EACL*.

Yvette Graham, Timothy Baldwin, Alistair Moffat, and Justin Zobel. 2017. [Can machine translation systems be evaluated by the crowd alone.](#) *Natural Language Engineering*.

Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. 2021. [CLIPScore: A reference-free evaluation metric for image captioning.](#) In *Proc. of EMNLP*.

Sepp Hochreiter and Jürgen Schmidhuber. 1997. [Long short-term memory.](#) *Neural Computation*.

Micah Hodosh, Peter Young, and Julia Hockenmaier. 2013. [Framing image description as a ranking task: Data, models and evaluation metrics.](#) *JAIR*.

Andrej Karpathy and Li Fei-Fei. 2015. [Deep visual-semantic alignments for generating image descriptions.](#) In *Proc. of CVPR*.

651	Daniel Khashabi, Gabriel Stanovsky, Jonathan Bragg,	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-	704
652	Nicholas Lourie, Jungo Kasai, Yejin Choi, Noah A.	Jing Zhu. 2002. BLEU: a method for automatic eval-	705
653	Smith, and Daniel S. Weld. 2021. GENIE: A leader-	uation of machine translation . In <i>Proc. of ACL</i> .	706
654	board for human-in-the-loop evaluation of text gener-		
655	ation .	Matt Post. 2018. A call for clarity in reporting BLEU	707
		scores . In <i>Proc. of WMT</i> .	708
656	Mert Kilickaya, Aykut Erdem, Nazli Ikizler-Cinbis, and	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya	709
657	Erkut Erdem. 2017. Re-evaluating automatic metrics	Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sas-	710
658	for image captioning . In <i>Proc. of EACL</i> .	try, Amanda Askell, Pamela Mishkin, Jack Clark,	711
		Gretchen Krueger, and Ilya Sutskever. 2021. Learn-	712
659	Philipp Koehn. 2004. Statistical significance tests for	ing transferable visual models from natural language	713
660	machine translation evaluation . In <i>Proc. of EMNLP</i> .	supervision .	714
661	Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin John-	Cyrus Rashtchian, Peter Young, Micah Hodosh, and	715
662	son, Kenji Hata, Joshua Kravitz, Stephanie Chen,	Julia Hockenmaier. 2010. Collecting image annota-	716
663	Yannis Kalanditis, Li-Jia Li, David A Shamma,	tions using Amazon’s Mechanical Turk . In <i>Proc. of</i>	717
664	Michael Bernstein, and Li Fei-Fei. 2016. Visual	<i>NAACL Workshop on Creating Speech and Language</i>	718
665	Genome: Connecting language and vision using	<i>Data with Amazon’s Mechanical Turk</i> .	719
666	crowdsourced dense image annotations . <i>IJCV</i> .		
667	Wojciech Kryscinski, Nitish Shirish Keskar, Bryan Mc-	Shaoqing Ren, Kaiming He, Ross Girshick, and Jian	720
668	Cann, Caiming Xiong, and Richard Socher. 2019.	Sun. 2015. Faster R-CNN: Towards real-time object	721
669	Neural text summarization: A critical evaluation . In	detection with region proposal networks . In <i>Proc. of</i>	722
670	<i>Proc. of EMNLP</i> .	<i>NeurIPS</i> .	723
671	Girish Kulkarni, Visruth Premraj, Sagnik Dhar, Sim-	Anna Rohrbach, Lisa Anne Hendricks, Kaylee Burns,	724
672	ing Li, Yejin Choi, Alexander C Berg, and Tamara L	Trevor Darrell, and Kate Saenko. 2018. Object hallu-	725
673	Berg. 2011. Baby talk: Understanding and generat-	cination in image captioning . In <i>Proc. of EMNLP</i> .	726
674	ing simple image descriptions . In <i>Proc. of CVPR</i> .		
675	Siming Li, Girish Kulkarni, Tamara L Berg, Alexan-	Keisuke Sakaguchi, Matt Post, and Benjamin	727
676	der C Berg, and Yejin Choi. 2011. Composing sim-	Van Durme. 2017. Grammatical error correction with	728
677	ple image descriptions using web-scale n-grams . In	neural reinforcement learning . In <i>Proc. of IJCNLP</i> .	729
678	<i>Proc. of CoNLL</i> .		
679	Xiujun Li, Xi Yin, Chunyuan Li, Xiaowei Hu,	Sebastian Schuster, Ranjay Krishna, Angel Chang,	730
680	Pengchuan Zhang, Lei Zhang, Lijuan Wang,	Li Fei-Fei, and Christopher D. Manning. 2015. Gener-	731
681	Houdong Hu, Li Dong, Furu Wei, Yejin Choi,	ating semantically precise scene graphs from text-	732
682	and Jianfeng Gao. 2020. Oscar: Object-semantics	ual descriptions for improved image retrieval . In	733
683	aligned pre-training for vision-language tasks . In	<i>Proc. of VL</i> .	734
684	<i>Proc. of ECCV</i> .		
685	Chin-Yew Lin. 2004. ROUGE: A package for automatic	Piyush Sharma, Nan Ding, Sebastian Goodman, and	735
686	evaluation of summaries . In <i>Proc. of Text Summa-</i>	Radu Soricut. 2018. Conceptual Captions: A cleaned,	736
687	<i>rization Branches Out</i> .	hypernymed, image alt-text dataset for automatic im-	737
		age captioning . In <i>Proc. of ACL</i> .	738
688	Tsung-Yi Lin, Michael Maire, Serge J. Belongie,	Kurt Shuster, Samuel Humeau, Hexiang Hu, Antoine	739
689	Lubomir D. Bourdev, Ross B. Girshick, James Hays,	Bordes, and Jason Weston. 2019. Engaging image	740
690	Pietro Perona, Deva Ramanan, Piotr Dollár, and	captioning via personality . In <i>Proc. of CVPR</i> .	741
691	C. Lawrence Zitnick. 2014. Microsoft COCO: com-		
692	mon objects in context . In <i>Proc. of ECCV</i> .	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob	742
		Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz	743
693	Qingsong Ma, Johnny Wei, Ondřej Bojar, and Yvette	Kaiser, and Illia Polosukhin. 2017. Attention is all	744
694	Graham. 2019. Results of the WMT19 metrics	you need . In <i>Proc. of NeurIPS</i> .	745
695	shared task: Segment-level and strong MT systems		
696	pose big challenges . In <i>Proc. of WMT</i> .	Ramakrishna Vedantam, C. Lawrence Zitnick, and Devi	746
		Parikh. 2015. CIDEr: Consensus-based image de-	747
697	George A. Miller. 1995. WordNet: A lexical database	scription evaluation . In <i>Proc. of CVPR</i> .	748
698	for English .		
699	Margaret Mitchell, Jesse Dodge, Amit Goyal, Kota Ya-	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien	749
700	maguchi, Karl Stratos, Xufeng Han, Alyssa Mensch,	Chaumond, Clement Delangue, Anthony Moi, Pier-	750
701	Alex Berg, Tamara Berg, and Hal Daumé III. 2012.	ric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz,	751
702	Midge: Generating image descriptions from com-	Joe Davison, Sam Shleifer, Patrick von Platen, Clara	752
703	puter vision detections . In <i>Proc. of EACL</i> .	Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven	753
		Le Scao, Sylvain Gugger, Mariama Drame, Quentin	754
		Lhoest, and Alexander Rush. 2020. HuggingFace’s	755
		transformers: State-of-the-art natural language pro-	756
		cessing . In <i>Proc. of EMNLP: System Demonstra-</i>	757
		<i>tions</i> .	758

759	Saining Xie, Ross B. Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. 2017. Aggregated residual transformations for deep neural networks . In <i>Proc. of CVPR</i> .	
760		
761		
762		
763	Yezhou Yang, Ching Teo, Hal Daumé III, and Yiannis Aloimonos. 2011. Corpus-guided sentence generation of natural images . In <i>Proc. of EMNLP</i> .	
764		
765		
766	Mark Yatskar, Michel Galley, Lucy Vanderwende, and Luke Zettlemoyer. 2014. See no evil, say no evil: Description generation from densely labeled images . In <i>Proc. of *SEM</i> .	
767		
768		
769		
770	Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. 2014. From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions . <i>TACL</i> .	
771		
772		
773		
774	Zheng Yuan and Ted Briscoe. 2016. Grammatical error correction using neural machine translation . In <i>Proc. of NAACL</i> .	
775		
776		
777	Pengchuan Zhang, Xiujun Li, Xiaowei Hu, Jianwei Yang, Lei Zhang, Lijuan Wang, Yejin Choi, and Jianfeng Gao. 2021. VinVL: Making visual representations matter in vision-language models . In <i>Proc. of CVPR 2021</i> .	
778		
779		
780		
781		
782	Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. BERTScore: Evaluating text generation with BERT . In <i>Proc. of ICLR</i> .	
783		
784		
785		
786	Luowei Zhou, Hamid Palangi, Lei Zhang, Houdong Hu, Jason J. Corso, and Jianfeng Gao. 2020. Unified vision-language pre-training for image captioning and VQA . In <i>Proc. of AAAI</i> .	
787		
788		
789		

A Appendix 790

A.1 Fluency Rubrics 791

Table 6 presents our fluency rubrics. They were developed by the first four authors (two of whom were native English speakers, and one was a graduate student in linguistics). Generally, if a fluency problem is expected to be easily corrected by a text postprocessing algorithm, the penalty is 0.1. More severe errors (e.g., broken sentence and ambiguity) are penalized more. 792
793
794
795
796
797
798
799

A.2 Automatic Metrics 800

Here we discuss details and configurations of the automatic metrics used in §3.2. CLIPScore and RefCLPScore use image features from CLIP (Radford et al., 2021), a crossmodal retrieval model trained on 400M image-caption pairs from the web. All the other five metrics only use reference captions. 801
802
803
804
805
806

BLEU BLEU (Papineni et al., 2002) is a precision-oriented metric and measures n-gram overlap between the generated and reference captions. We use the SACREBLEU implementation of BLEU-4 and get sentence-level scores (Post, 2018).⁷ 807
808
809
810
811
812

ROUGE ROUGE (Lin, 2004) measures the number of overlapping n-grams between the generated and reference captions. We use the HuggingFace implementation of ROUGE-L (Wolf et al., 2020). 813
814
815
816

CIDEr CIDEr (Vedantam et al., 2015) measures the cosine similarity between the n-gram counts of the generated and reference captions with TF-IDF weighting. We use the implementation from the pycocoevalcap package.⁸ 817
818
819
820
821

SPICE SPICE (Anderson et al., 2016) predicts scene graphs from the generated and reference captions using the Stanford scene graph parser (Schuster et al., 2015). It then measures the F_1 score between scene graphs from the generated and reference captions. WordNet Synsets are used to cluster synonyms (Miller, 1995). We again use the implementation from the pycocoevalcap package. 822
823
824
825
826
827
828
829

BERTScore BERTScore (Zhang et al., 2020) aligns tokens between the generated and reference captions using contextual word representations from BERT (Devlin et al., 2019). We use 830
831
832
833

⁷<https://github.com/mjpost/sacreBLEU/blob/v1.2.12/sacrebleu.py#L999>.

⁸<https://github.com/salaniz/pycocoevalcap>.

Fluency Error Type	Penalty	Example
Obvious spelling error, one vs. two words	0.1	<i>cel phone, surf board</i>
Grammatical error that can be easily fixed	0.1	<i>a otter</i>
Casing issue	0.1	<i>tv, christmas</i>
Hyphenation	0.1	<i>horse drawn carriage</i>
Interpretable but unnatural wording	0.1	<i>double decked bus</i>
Non-trivial punctuation	0.2	<i>A bird standing in the wooded area with leaves all around.</i>
Misleading spelling error	0.5	<i>A good stands in the grass next to the water. (good→goose)</i>
Duplication	0.5	<i>A display case of donuts and doughnuts.</i>
Ambiguity	0.5	<i>A cat is on a table with a cloth on it.</i>
Awkward construction	0.1–0.5	<i>There is a freshly made pizza out of the oven.</i>
Broken sentence	0.5+	<i>A large concrete sign small buildings behind it.</i>

Table 6: Fluency penalty rubrics.

the HuggingFace implementation and compute the F_1 score. As in Zhang et al. (2020), we take the maximum score over all reference captions.

CLIPScore CLIPScore (Hessel et al., 2021) is the only *referenceless* metric out of the 7 metrics. It measures the cosine similarity between the generated caption and given image using the representations from CLIP. It is shown to correlate better with human judgments from prior work, compared to previous reference-based metrics (Hessel et al., 2021). We use the official implementation by the authors.⁹

RefCLIPScore RefCLIPScore augments CLIPScore with the maximum similarity between the generated and reference captions. We again use the official implementation.

A.3 Evaluated Captions

We evaluated the following four strong models from the literature as well as human-generated captions. They share similar pipeline structure but vary in model architecture, (pre)training data, model size, and (pre)training objective. Evaluating captions from them will enable us to better understand what has been improved and what is still left to future captioning models.

Up-Down The bottom-up and top-down attention model (Up-Down, Anderson et al., 2018) performs pipelined image captioning: *object detection* that finds objects and their corresponding image regions and *crossmodal generation* that predicts a caption based on the features from object detection. The bottom-up attention finds salient image regions during object detection, and the top-down

one attends to these regions during crossmodal generation. Up-Down uses Faster R-CNN (Ren et al., 2015) and LSTMs (Hochreiter and Schmidhuber, 1997) for object detection and crossmodal generation respectively. Faster R-CNN is trained with the Visual Genome dataset (Krishna et al., 2016), and the crossmodal generation model is trained on the MSCOCO dataset. We generate captions for the test data with a model optimized with crossentropy.¹⁰

Unified-VLP Unified-VLP (Zhou et al., 2020) also runs a pipeline of object detection and crossmodal generation. Faster R-CNN and the transformer architecture (Vaswani et al., 2017) are used for object detection and crossmodal generation respectively. Similar to Up-Down, the Faster R-CNN object detection model is trained with the Visual Genome dataset. The transformer generation model, on the other hand, is initialized with base-sized BERT (Devlin et al., 2019) and pretrained on the Conceptual Captions dataset (3M images, Sharma et al., 2018) with the masked and left-to-right language modeling objectives for the captions. The crossmodal generation model is then finetuned on the MSCOCO dataset. We apply beam search of size 5 to the model with CIDEr optimization.

VinVL-base, VinVL-large VinVL with Oscar (Li et al., 2020; Zhang et al., 2021) performs a similar pipeline of object detection, followed by crossmodal generation. The crossmodal model is initialized with BERT (Devlin et al., 2019) as in Unified-VLP but uses detected object tags to encourage alignments between image features and word representations. The object detection model

⁹<https://github.com/jmhessel/pycocoevalcap>.

¹⁰https://vision-explorer.allenai.org/image_captioning.

with the ResNeXt-152 C4 architecture (Xie et al., 2017) is pretrained with ImageNet (Deng et al., 2009) and trained on 2.5M images from various datasets. The transformer-based crossmodal generator is initialized with BERT, pretrained with 5.7M images, and finetuned for MSCOCO captioning. We use VinVL-base and VinVL-large that are both finetuned with CIDEr optimization¹¹ and generate captions with beam search of size 5.

Human In addition to machine-generated captions from the four models, we assessed the quality of human-generated reference captions from MSCOCO. This will allow us to understand the performance gap between machines and humans, as well as the quality of crowdsourced captions. Human-generated captions were created using Amazon Mechanical Turk (Chen et al., 2015). Crowdworkers were only given the following instructions (Chen et al., 2015):

- Describe all the important parts of the scene.
- Do not start the sentences with “There is.”
- Do not describe unimportant details.
- Do not describe things that might have happened in the future or past.
- Do not describe what a person might say.
- Do not give people proper names.
- The sentences should contain at least 8 words.

Every image has five human-generated captions, and we randomly selected one for each to evaluate. We found, however, a non-negligible number of noisy captions in the MSCOCO dataset from annotation spammers. We often find subjective adjectives (e.g., *very nice/clean/cute*) or words that diverge from *inclusive language* in reference captions, probably because crowdworkers increased the number of words in captions effortlessly (see the last instruction item that says captions have to have 8+ words). To better estimate the performance of a human that invests reasonable effort into the captioning task, we resampled a caption for 13% of the test images, which would have been given a total score lower than 4.0.

A.4 Score Distributions

Seen in Fig. 2 are distributions of precision and recall scores for human and machine-generated captions. We see that the precision distribution looks similar between Human and machines, but

not recall. This provides further support for our claim that current machines fall short of humans particularly in recall.

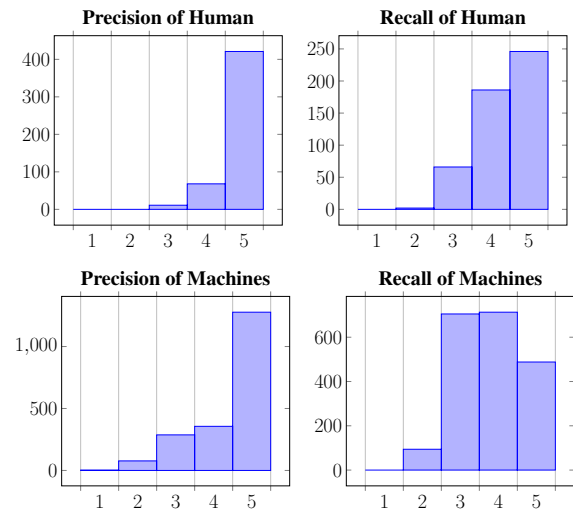


Figure 2: Precision/recall histograms for human- and machine-generated captions.

A.5 Additional Machine vs. Human Examples

Table 7 provides an additional example that contrasts machine- and human-generated captions. All machines generate generic captions and ignore the most important information that a traditional Thanksgiving dinner is being served on the table.

¹¹https://github.com/microsoft/Oscar/blob/master/VinVL_MODEL_ZOO.md#Image-Captioning-on-COCO.


Image	Caption	P	R	Total
	10-A: Up-Down <i>A table that has some food on it.</i>	5	2	3.5
	10-B: Unified-VLP <i>A table with plates of food on a table.</i>	5	2	3.5
	10-C: VinVL-base <i>A red table topped with plates of food and bowls of food.</i>	4	2	3
	10-D: VinVL-large <i>A table with a turkey and other food on it.</i>	5	3	4
	10-E: Human <i>A table set for a traditional Thanksgiving dinner.</i>	5	5	5

Table 7: Additional example that contrasts machine- and human-generated captions. Similar to Table 5, machine-generated captions ignore the most salient information: Thanksgiving dinner. None of these captions are penalized for fluency, conciseness, or inclusive language.