

Optimization of image description metrics using policy gradient methods

Siqi Liu^{*1}, Zhenhai Zhu², Ning Ye², Sergio Guadarrama², and Kevin Murphy²

siqi.liu@cs.ox.ac.uk

{zhenhai, nye, sguada, kpmurphy}@google.com

¹Department of Computer Science, University of Oxford

²Google

Abstract

In this paper, we propose a novel training procedure for image captioning models based on policy gradient methods. This allows us to directly optimize for the metrics of interest, rather than just maximizing likelihood of human generated captions. We show that by optimizing for standard metrics such as BLEU, CIDEr, METEOR and ROUGE, we can develop a system that improve on the metrics and ranks first on the MSCOCO image captioning leader board, even though our CNN-RNN model is much simpler than state of the art models. We further show that by also optimizing for the recently introduced SPICE metric, which measures semantic quality of captions, we can produce a system that significantly outperforms other methods as measured by human evaluation. Finally, we show how we can leverage extra sources of information, such as pre-trained image tagging models, to further improve quality.

1 Introduction

Image description (captioning) is the task of describing the visual content of an image using one or more sentences. This has many applications, including text-based image retrieval, accessibility for blind users, and human-robot interaction. There are many ways to perform this task (see [5] for a recent review). In this paper, we focus on improving the current state of the art method, which is based on encoder-decoder neural networks. In such a model, the image content is encoded using a convolutional neural network (CNN), and then the text is generated, one word at a time, using a recurrent neural network (RNN).

One major flaw with current approaches is that they

use (penalized) maximum likelihood estimation (MLE) for training. That is, the model parameters θ are trained to maximize

$$\begin{aligned} L(\theta) &= \frac{1}{N} \sum_{n=1}^N \log p(\mathbf{y}^n | \mathbf{x}^n, \theta) \\ &= \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} \log p(y_t^n | y_{1:t-1}^n, \mathbf{x}^n, \theta) \end{aligned}$$

where \mathbf{x}^n is the n 'th image, $\mathbf{y}^n = (y_1^n, \dots, y_{T_n}^n)$ is the ground truth caption of the n 'th image and N is the total number of labeled examples. Depending on the dataset, there could be multiple ground truth captions for a given image although it does not affect our derivation of the algorithm. We ignore this distinction in this paper for notational clarity.

The main trouble with this objective is that each prediction is conditioned on the previously observed words from the ground truth. However, at test time, the model will be fed its own predictions. This discrepancy has been called “exposure bias” [14], and can lead to poor performance at inference time. In particular, errors made at an intermediate step can quickly accumulate over the following steps, so the model will likely diverge from desired trajectories. One approach to mitigate this, known as “scheduled sampling”, was proposed in [4]; however, this method has been shown to be statistically inconsistent [9], although it improves model performance in practice. Another problem with MLE is that it only evaluates how much probability is assigned to the true sequence; the quality of all other generated sequences is ignored.

A better approach is to optimize for the same metrics that we use at test time, such as BLEU [13], CIDEr [19], METEOR [3], and ROUGE [10]; we shall call these metrics “BCMR” for short. Although these metrics are not differentiable, we can use policy gradient (PG) methods

^{*}The major part of this work was done while Siqi Liu was an intern at Google.

[16] to optimize them, by treating the score of a candidate sentence as analogous to a reward signal in a reinforcement learning setting. In such a framework, the RNN decoder acts like a stochastic policy, where choosing an action corresponds to generating a word.

The idea of using PG to optimize BLEU score for image captioning has been previously explored in [14]. These authors used a special case of PG known as REINFORCE [22], which they combined with MLE to create a learning method called “MIXER”. However, in their PG method, they implicitly assumed each intermediate action (word) in a partial sequence has the same reward as the sequence-level reward, which is not true in general.

In this paper, we extend this prior work in several ways. First, we use an improved implementation of PG, which estimates the expected future reward of each intermediate action using Monte Carlo rollouts, as in [28]. We also use a learned parametric baseline estimator to reduce the variance of the estimate. The proposed training algorithm is able to surpass the previous state of the art on the MSCOCO captioning leaderboard [11], despite the simplicity of the captioning model itself (Section 4).

Second, we extend our technique to directly optimize the recently introduced SPICE metric [1], which has been shown to correlate much more closely with human judgement of semantic quality than previous BCMR metrics. We show that incorporating this metric results in captions that are deemed significantly better when judged by human raters.

Finally, we show that our proposed method is orthogonal to other techniques proposed in recent literature. As a proof of concept, we show how to integrate an image tagging system [6], which has been pre-trained on a large quantity of noisily labeled web images, as an additional input to the decoder (beyond just the CNN features). We show that this method enables more effective use of such additional inputs, compared to MLE baselines.

2 Related work

There is an extensive body of work on image captioning, which is too large to review here. See [5] for a recent summary. The most relevant piece of prior work is [21], who proposed one of the first encoder-decoder networks for image captioning, called “Show and Tell” (ST), also known as “Neural Image Captioner” (NIC). As our contribution is orthogonal to the underlying model architecture, we use the same ST model for most of this paper.

Numerous extensions to the basic encoder-decoder framework have been proposed, many of which try to enrich the encoder’s representation, so that the image is not represented just by the global features from a CNN. One common approach is to use an attention mechanism, which

lets the decoder focus on specific parts of the input image. This approach is used in [24], who proposed the “Show, Attend and Tell” model. See also [25] for a more recent extension, that applies an attentional RNN on top of the encoder before passing to the decoder. In this paper, we do not use attention mechanisms; nevertheless, we are able to outperform such models, by optimizing better learning objectives.

Another approach is to add high level image visual attributes to the model, enriching the image encoding beyond the CNN. This usually takes the form of an image tagger [23, 26], but can also use specialized modules, such as object detectors [20], or face detection and landmark recognition [18].

There is some recent work on optimizing sequence level objective functions, instead of maximum likelihood training, but mostly in the context of machine translation (see e.g. [12, 28, 2, 15]). Recent work in visual explanation such as [8] employed REINFORCE to optimize sequence level reward such that generated captions are class discriminative, which is orthogonal to our goal. As far as we know, the only paper to explore objectives beyond maximum likelihood for image captioning is the MIXER method of [14]. However, we significantly outperform this method, as we show in Section 4, by using an improved and more general PG algorithm.

3 Methods

In this section, we explain our approach in more detail. First we discuss the policy gradient algorithm, which can be used to optimize any kind of reward function. Next we discuss which reward function to use. We then discuss the model itself, which is a standard CNN-RNN. Finally we discuss an extension to the basic model, which uses a pre-trained image tagger for improved performance.

3.1 Training using policy gradient

At time step t , we pick a discrete action, which corresponds to choosing a word $g_t \in \mathcal{V}$, using a stochastic policy or generator $\pi_\theta(g_t|s_t, \mathbf{x})$, where $s_t = g_{1:t-1}$ is the sequence of words chosen so far, \mathbf{x} is the image and θ are the parameters of the model. Note that in our case the state transition function is deterministic: we simply append the word chosen at time t to get $s_{t+1} = s_t; g_t$, where $u; v$ is the concatenation of the strings u and v .

When we reach the end of the sequence (i.e., once the generator emits the end-of-sentence marker), we get a reward of $R(g_{1:T}|\mathbf{x}^n, \mathbf{y}^n)$, which is the score for producing caption $g_{1:T}$ given image \mathbf{x}^n and ground truth caption (or set of captions) \mathbf{y}^n . This reward can be any function, such as BCMR or SPICE.

Although rewards are only computed once we have reached the end of the sentence, it is useful to have an indication of the quality of a partial sequence. So we define the value of a state (partial sequence) as its expected future reward:

$$V_\theta(g_{1:t}|\mathbf{x}^n, \mathbf{y}^n) = E_{g_{t+1:T}}[R(g_{1:t}; g_{t+1:T}|\mathbf{x}^n, \mathbf{y}^n)] \quad (1)$$

where the expectation is w.r.t. $g_{t+1:T} \sim \pi_\theta(\cdot|g_{1:t}, \mathbf{x}^n)$.

Our goal is optimize the average reward starting from the initial (empty) state s_0 , averaged over the examples in the training set. Hence we define the following objective, which we wish to maximize:

$$J(\theta) = \frac{1}{N} \sum_{n=1}^N V_\theta(s_0|\mathbf{x}^n, \mathbf{y}^n) \quad (2)$$

We now discuss how to optimize Eqn. (2). For simplicity, we will consider a single example n , so we will drop the \mathbf{x}^n and \mathbf{y}^n notation. To compute the gradient of $J(\theta)$, we can use the policy gradient theorem from [16]. In the special case of deterministic transition functions, this theorem simplifies as shown below (see [2] for a proof):

$$\nabla_\theta V_\theta(s_0) = E_{g_{1:T}} \left[\sum_{t=1}^T \sum_{g_t \in \mathcal{V}} \nabla_\theta \pi_\theta(g_t|g_{1:t-1}) Q_\theta(g_{1:t-1}, g_t) \right] \quad (3)$$

where we define the Q function for a state-action pair as follows:

$$Q_\theta(g_{1:t-1}, g_t) = E_{g_{t+1:T}}[R(g_{1:t-1}; g_t; g_{t+1:T})] \quad (4)$$

We can approximate the gradient of the value function with M sample paths, $g_{1:T}^m \sim \pi_\theta$, generated from our policy. This gives

$$\begin{aligned} \nabla_\theta V_\theta(s_0) &\approx \frac{1}{M} \sum_{m=1}^M \sum_{t=1}^T E_{g_t} [\nabla_\theta \log \pi_\theta(g_t|g_{1:t-1}^m) \\ &\quad \times Q_\theta(g_{1:t-1}^m, g_t)] \end{aligned} \quad (5)$$

where the expectation is w.r.t. $g_t \sim \pi_\theta(g_t|g_{1:t-1}^m)$, and where we have exploited the fact that

$$\nabla_\theta \pi_\theta(a|s) = \pi_\theta(a|s) \frac{\nabla_\theta \pi_\theta(a|s)}{\pi_\theta(a|s)} = \pi_\theta(a|s) \nabla_\theta \log \pi_\theta(a|s)$$

If we use $M = 1$, we can additionally replace the E_{g_t} with the value in the sample path, g_t^m , as in REINFORCE. In our experiment, we used $M = 1$ and we subsequently drop the superscript m in the rest of this paper for notational clarity.

The only remaining question is how to estimate the function $Q(s_t, g_t)$. For this, we will follow [28] and use Monte Carlo rollouts. In particular, we first sample K

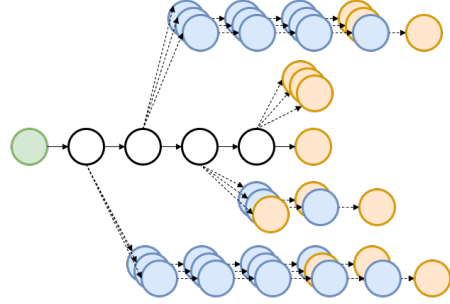


Figure 1: The value of each action is estimated as the average rewards received by its $K = 3$ rollout sequences. Solid arrows indicate the sequence of actions being evaluated. The tokens in green and yellow are respectively BOS (beginning of sequence) and EOS (end of sequence) tokens. Sequences in blue are rollout sequences sampled from partial sequences. Note that rollout sequences do not always have the same length, as they are separately sampled from a stochastic policy.

continuations of the sequence $s_t; g_t$ to get $g_{t+1:T}^k$. Then we compute the average

$$Q(g_{1:t-1}, g_t) \approx \frac{1}{K} \sum_{k=1}^K R(g_{1:t-1}; g_t; g_{t+1:T}^k) \quad (6)$$

If we are in a terminal state, we define $Q(g_{1:T}, \text{EOS}) = R(g_{1:T})$. This process is illustrated in Figure 1 where we show the case of $K = 3$. In English, we estimate how good a particular word choice g_t is by averaging over all complete sequences sampled according to the current policy, conditioned on the partial sequence $g_{1:t-1}$ sampled from the current policy so far.

The above gradient estimator is an unbiased but high variance estimator. One way to reduce its variance is to estimate the expected baseline reward $E_{g_t}[Q(g_{1:t-1}, g_t)]$ using a parametric function; we will denote this baseline as $B_\phi(g_{1:t-1})$. We then subtract this baseline from $Q_\theta(g_{1:t-1}, g_t)$ to get the following estimate for the gradient (using $M = 1$ sample paths):

$$\begin{aligned} \nabla_\theta V_\theta(s_0) &\approx \sum_{t=1}^T \sum_{g_t} [\pi_\theta(g_t|s_t) \nabla_\theta \log \pi_\theta(g_t|s_t) \\ &\quad \times (Q_\theta(s_t, g_t) - B_\phi(s_t))] \end{aligned} \quad (7)$$

where $s_t = g_{1:t-1}$. Subtracting the baseline does not affect the validity of the estimated gradient, but reduces its variance. Here, we simply refer to prior work ([29], [22]) for a full derivation of this property. We train the parameters ϕ of the baseline estimator to minimize the following loss:

$$L_\phi = \sum_t E_{s_t} E_{g_t} (Q_\theta(s_t, g_t) - B_\phi(s_t))^2 \quad (8)$$

In our experiments, the baseline estimator is an MLP which takes as input the hidden state of the RNN at step t . To avoid creating a feedback loop, we do not back-propagate gradients through the hidden state from this loss.

In language generation settings, a major challenge facing PG algorithms is the large action space. This is the case in our task, where the action space corresponds to the entire vocabulary of 8,855 symbols. To help “warm start” the training, we pre-train the RNN decoder (stochastic policy) using a cross-entropy loss, before switching to PG training. This prevents the agent from performing random walks through exponentially many possible paths at the beginning of the training.

The overall algorithm is summarized in Algorithm 1. Note that the Monte Carlo rollouts only require a forward pass through the RNN, which is much more efficient than the forward-backward pass needed for the CNN. Additionally the rollouts can be also be done in parallel for multiple sentences. Consequently, PG training is only about twice as slow as MLE training.

Algorithm 1: PG training algorithm

```

1 Input:  $\mathcal{D} = \{(\mathbf{x}^n, \mathbf{y}^n) : n = 1 : N\}$ ;
2 Train  $\pi_\theta(g_{1:T}|x)$  using MLE on  $\mathcal{D}$ ;
3 Train  $B_\phi$  using MC estimates of  $Q_\theta$  on a small subset
  of  $\mathcal{D}$ ;
4 for each epoch do
5   for example  $(\mathbf{x}^n, \mathbf{y}^n)$  do
6     Generate sequence  $g_{1:T} \sim \pi_\theta(\cdot|\mathbf{x}^n)$ ;
7     for  $t = 1 : T$  do
8       Compute  $Q(g_{1:t-1}, g_t)$  for  $g_t$  with  $K$ 
        Monte Carlo rollouts, using (6);
9       Compute estimated baseline  $B_\phi(g_{1:t-1})$ ;
10    Compute  $\mathcal{G}_\theta = \nabla_\theta V_\theta(s_0)$  using (7);
11    Compute  $\mathcal{G}_\phi = \nabla_\phi L_\phi$ ;
12    SGD update of  $\theta$  using  $\mathcal{G}_\theta$ ;
13    SGD update of  $\phi$  using  $\mathcal{G}_\phi$ ;
```

3.2 Reward Functions for the Policy Gradient

We can use our PG method to optimize many different reward functions. Common choices include BLEU, CIDEr, METEOR and ROUGE. Code for all of these metrics is available as part of the MSCOCO evaluation toolkit.¹ We decided to use a weighted combination of all of these. Since these metrics are not on the same scale, we chose in our experiments the set of weights such that all metrics have approximately the same magnitude. In

¹ <https://github.com/tylin/coco-caption>.

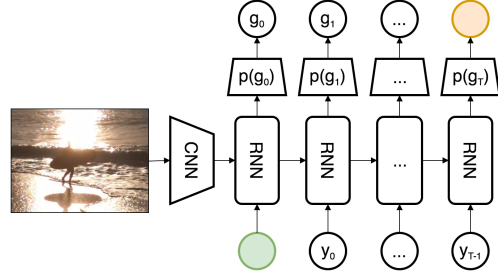


Figure 2: Model architecture of Show and Tell image captioning system [21]. The tokens in green and yellow are respectively BOS (beginning of sequence) and EOS (end of sequence) tokens. At testing time, output from previous time step g_{t-1} is used as input in lieu of y_{t-1} .

particular, we chose 0.5, 0.5, 1.0, 1.0, 1.0, 5.0, 2.0 for BLEU-1, BLEU-2, BLEU-3, BLEU-4, CIDEr, METEOR and ROUGE respectively. Optimizing this weighted combination of BCMR gives state-of-the-art results on the MSCOCO test set, as we discuss in Section 4.2.

One problem with the BCMR metrics is that they are not well correlated with human judgment individually [1]. We therefore also tried optimizing the recently introduced SPICE metric [1], which better reflects human estimates of quality. In particular, SPICE is the only metric that ranks humans above algorithms in terms of captioning quality on the MSCOCO benchmark. We use the open source release of the SPICE code² to evaluate the metric.

Interestingly, we have found that just optimizing SPICE tended to result in captions which are very detailed (as desired), but which often had many repeated phrases, as we show in Section 4. This is because SPICE measures semantic similarity (in terms of a scene graph) between sets of sentences, but does not pay attention to syntactical factors (modulo the requirement that the generated sentence be parseable). We therefore combined SPICE with the CIDEr metric (considered as the best of the standard automatic metrics for MSCOCO), a combination we call SPIDER for short. Based on preliminary experiments, we decided to use an equal weighting for both.

3.3 Encoder-decoder architecture

We use a CNN-RNN architecture similar to the one proposed in the original Show-Tell paper [21]. A high-level diagram is shown in Figure 2. Each symbol in the vocabulary is embedded as a 512 dimensional dense word embedding vector, whose values are initialized randomly.

The encoder CNN is implemented as an Inception-V3 [17] network pretrained on ImageNet³. The

² <https://github.com/peteanderson80/SPICE>.

³ We used the open-source implementation available at:

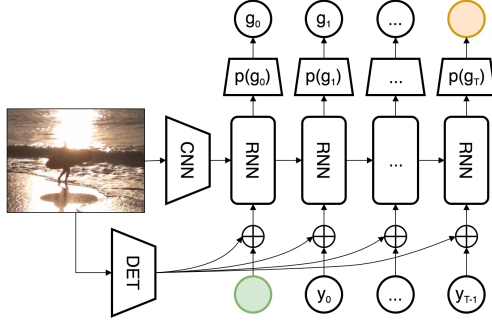


Figure 3: Model architecture of attribute augmented model. The tokens in green and yellow are respectively BOS (beginning of sequence) and EOS (end of sequence) tokens. At testing time, output from previous time step g_{t-1} is used as input in lieu of y_{t-1} .

RNN decoder is a one-layer LSTM with a state size of 512 units, initialized randomly. Each image is encoded by Inception-V3 as a dense feature vector of dimension 2,048 which is then projected to 512 dimension with a linear layer and used as the initial state of RNN decoder.

At training time, we always feed in the ground truth symbol to the RNN decoder; at inference time, the sampled output is fed to the RNN as the next input symbol. We use a greedy decoding procedure, in which we pick the most probable word at each time step. Note that this is suboptimal, since

$$\prod_{t=1}^T \max_{g_t} p(g_t | g_{1:t-1}) \leq \max_{g_{1:T}} \prod_{t=1}^T p(g_t | g_{1:t-1}) \quad (9)$$

We can better approximate the globally optimal sequence by using beam search. However, we find that models trained with PG do not benefit from this in practice, since they tend to concentrate their probability mass on a very small number of actions at each step.

3.4 Adding visual tags

In this section, we describe a simple extension to the basic model from Section 3.3. In particular, we leverage an in-house image tagging system, which is trained on a large, but noisily labeled, dataset called JFT (see [6] for details). This system returns a vector of probability scores $\mathbf{p} = \{p_0, p_1, \dots, p_N\}$ for a set of $N = 3,713$ classes, where $p_i \in [0, 1]$ indicates the probability of presence of a visual attribute of class i in the image. Typically 60–70 tags have a non-zero confidence score for each image. See Figure 1 for some example outputs from this system. A high-level diagram of this architecture is shown in Figure 3.

https://github.com/tensorflow/models/blob/master/slim/nets/inception_v3.py

Given the visual tag score vector, we compute the input to the decoder RNN at time t by combining the scores from the tagger with the previous word in the sequence as follows:

$$\mathbf{x}_t = \mathbf{T}_w \mathbf{w}_t + \mathbf{T}_p \mathbf{p} \quad (10)$$

where \mathbf{w}_t is the embedding of the t 'th word, $\mathbf{T}_w \in \mathcal{R}^{d_h \times d_w}$ maps the word embedding to the RNN state space, \mathbf{p} is the vector of tag probabilities, and $\mathbf{T}_p \in \mathcal{R}^{d_h \times d_p}$ maps from the $d_p = N$ tags to the RNN state space. We shall call our model Show-Tell-Tag.

Our attribute augmented model is similar to LSTM-A₅, recently introduced in [26]. The main difference lies in the choice of CNN (we used Inception-V3 as opposed to ResNet) as well as the vocabulary of visual tags (we train on an internal dataset with 3k noisy tags, whereas they use multi-instance learning to train a classifier to predict the top 1k words in MSCOCO). (See also [23] for some related work on using visual tags for image captioning.)

4 Results

In this section, we report results obtained by different methods on the MSCOCO dataset. This has 82,081 training images, and 40,137 validation images, each with at least 5 ground truth captions. Following standard practice for methods that evaluate on the MSCOCO test server, we hold out a small subset of 1,665 validation images for hyper-parameter tuning, and use the remaining combined training and validation set for training.

We preprocess the text data by lower casing, and replacing words which occur less than 4 times in the 82k training set with UNK; this results in a vocabulary size of 8,855 (identical to the one used in [21]). At training time, we keep all captions to their maximum lengths. At testing time, the generated sequences are truncated to 30 symbols in all experiments.

We report results of the following 7 systems:

- MLE: the Show-Tell model, trained with maximum likelihood estimation. This is our baseline approach, and is close to state of the art. Note that our implementation gives similar results to [21], which uses the same model, but was trained with scheduled sampling.
- MLE-TAG: the Show-Tell-Tag model trained with MLE.
- PG-BCMR: the Show-Tell model trained using PG to optimize the BCMR metric.
- PG-BCMR-TAG: the Show-Tell-Tag model trained using PG to optimize the BCMR metric.

- PG-SPICE: the Show-Tell model trained using PG to optimize the SPICE metric.
- PG-SPIDER: the Show-Tell model trained using PG to optimize the SPIDER metric (with equal weight on SPICE and on CIDEr).
- PG-SPIDER-TAG: the Show-Tell-Tag model trained using PG to optimize the SPIDER metric (with equal weight on SPICE and on CIDEr).

4.1 Qualitative Analysis

Table 1 shows some example captions generated by the 7 different systems. We see that PG-SPICE tends to generate ungrammatical sentences, with a lot of repeated phrases. This is because SPICE measures how well the scene graph induced by a sentence matches the ground truth scene graph, but is relatively insensitive to syntactic quality. However, we also see that by combining SPICE with CIDEr, we get much better results. Henceforth we shall only consider the SPIDER metric, and will ignore pure SPICE.

4.2 Results using automatic metrics on MSCOCO

In this section, we quantitatively evaluate the methods using the MSCOCO online evaluation server⁴.

Table 2 shows the performance of various state-of-the-art models in the literature for image captioning. In particular, we show the best performing models according to the official MSCOCO C-5 leaderboard at the time of writing (Nov 2016). We also report the results of 6 experiments we ran, which include all our models except for PG-SPICE.

We start by discussing the results of our baseline ST models, without using the image tagger. We see that all the PG methods outperform MLE training. We also see that the PG methods significantly outperform all previous methods, even the ones which use more sophisticated models, such as those based on attention (Montreal/Toronto, ATT, Review Net), those that use more complex decoders (Berkeley LRCN), and those that use high-level visual attributes (MSM@MSRA, ATT).

Our PG method also outperforms MIXER [14], which is the only prior work (to the best of our knowledge) which uses PG for image captioning. Note, however, that the BLEU-4 metric quoted for MIXER is not directly comparable to our numbers, as the authors only performed evaluation on 5,000 images from the MSCOCO validation set, not on the official test server. (Their code for image captioning is not available, so we have not been able to do an apples-to-apples comparison.)

⁴mscoco.org/dataset/#captions-leaderboard

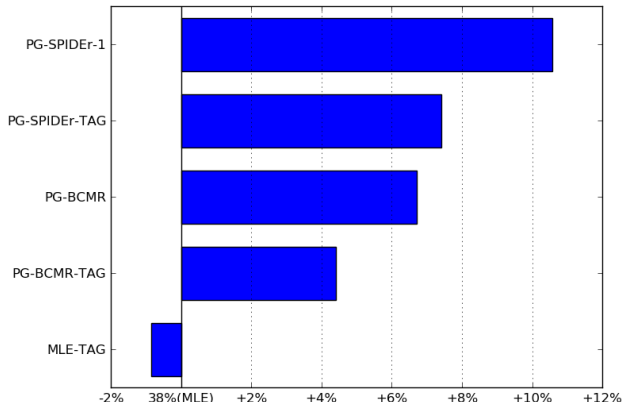


Figure 4: Results of human evaluation on 492 images randomly sampled from the MSCOCO test set. We report the difference in percentage of “not bad” captions for each method compared to baseline 38% of MLE model.

Finally, we see that adding visual tags consistently improves the performance of PG-SPIDER (especially on the CIDEr metric), but not when using MLE or PG-BCMR training. A detailed investigation into why this is the case is left to future work.

4.3 Results using human metrics on MSCOCO

Since we know that BCMR metrics are not well correlated with human measures of quality, we decided to evaluate the methods using human raters. In particular, we use a crowd sourcing platform, using raters who have prior experience with evaluating image captioning and other computer vision models. We showed each image-caption pair to 3 different raters, and asked them to evaluate it on a 4 point scale, depending on whether the caption is “bad”, “okay”, “good” or “excellent”.⁵ We then take the majority vote to get the final rating. If no majority is found, the rating is considered unknown, and this image is excluded from the analysis.

To simplify the presentation, we just focus on measuring the fraction of captions that are classified as “not bad”, which we interpret as the union of “okay”, “good” and “excellent”. (The reason for this is that we want to be sure

⁵ The definitions of these terms, which we gave to raters, is as follows. Excellent: “The caption correctly, specifically and completely describes the foreground/main objects/events/theme of the image.” Good: “The caption correctly and specifically describes most of the foreground/main objects/events/theme of the image, but has minor mistakes in some minor aspects.” Okay: “The caption correctly describes some of the foreground/main objects/events/theme of the image, but is not specific to the image and has minor mistakes in some minor aspects.” Bad: “The caption misses the foreground/main objects/events/theme of the image or the caption contains obviously hallucinated objects/activities/relationships.”






Images with Visual Tags	Ground Truth Captions	Generated Captions
 <ul style="list-style-type: none"> Vehicle: 0.98 Emergency: 0.97 Fire department: 0.96 Car: 0.92 	<ol style="list-style-type: none"> 1. a red and yellow fire truck and some buildings 2. An overhead view shows a fire engine in the street. 3. A red and yellow fire truck with ladders on top 4. A firetruck is parked in the street in between stop lights. 5. A fire truck (ladder truck) drives down a street in the city. 	<ul style="list-style-type: none"> • MLE: a red and white bus is driving down the street • MLE-TAG: a fire truck is driving down the street . • PG-BCMR: a red bus driving down a city street . • PG-BCMR-TAG: a fire truck driving down a street . • PG-SPICE: a red double decker bus on a city street on a street with a bus on the street with a bus on the street in front of a bus on • PG-SPIDER: a red fire truck is on a city street. • PG-SPIDER-TAG: a fire truck is on the street .
 <ul style="list-style-type: none"> Umbrella: 0.79 Field: 0.73 Beauty: 0.66 Plant: 0.61 	<ol style="list-style-type: none"> 1. an umbrella in a field of flowers 2. a big blue umbrella out in front of a field 3. A blue umbrella is placed near a field of wheat. 4. a blue umbrella juts from the big corn field 5. A blue umbrella sits next to a wheat field. 	<ul style="list-style-type: none"> • MLE: a blue umbrella sitting on top of a lush green field . • MLE-TAG: a blue umbrella sitting in the grass next to a tree . • PG-BCMR: a blue umbrella sitting in the grass . • PG-BCMR-TAG: a person holding an umbrella in the grass . • PG-SPICE: a blue umbrella on a black and white umbrella on a grassy field with a person holding a kite in a field with a kite in the sky with a • PG-SPIDER: a blue umbrella sitting on top of a field . • PG-SPIDER-TAG: a blue umbrella sitting in the grass .
 <ul style="list-style-type: none"> Person: 0.98 Mammal: 0.95 Bicycle: 0.93 Vehicle: 0.92 	<ol style="list-style-type: none"> 1. A dog rides in a cart pulled by a man on a bike. 2. A man on a bike pulling a cart with a dog inside. 3. A dog in a basket being carried by its owner. 4. A man riding a bike with a wooden trainer attached and a dog riding in it. 5. A man on a bike pulling a dog in a cart. 	<ul style="list-style-type: none"> • MLE: a man and a dog are sitting in a cart • MLE-TAG: a large brown dog on a leash is being pulled by a man . • PG-BCMR: a man is standing in front of a dog . • PG-BCMR-TAG: a man sitting on a truck with a dog . • PG-SPICE: a group of people on a large group of people on a street with a dog on a cart on a street with a large group of people on a • PG-SPIDER: a man standing next to a dog in a truck . • PG-SPIDER-TAG: a man riding a dog on a back of a truck .
 <ul style="list-style-type: none"> Person: 0.98 Room: 0.84 Event: 0.79 Adult: 0.75 	<ol style="list-style-type: none"> 1. A group of people converse in an office setting. 2. A group of people playing a game with remote controllers. 3. Four young people have crowded into a small office. 4. A group of people standing next to each other in a room. 5. a group of people standing next to each other with some of them holding video game controllers 	<ul style="list-style-type: none"> • MLE: a group of people standing around a living room . • MLE-TAG: a group of people in a room with a video game controller . • PG-BCMR: a group of people standing in a room . • PG-BCMR-TAG: a group of people sitting in a room . • PG-SPICE: a group of people in a room with a man in a chair holding a nintendo wii remote in a living room with a man in a chair holding a • PG-SPIDER: a group of people playing a video game in a living room . • PG-SPIDER-TAG: a group of people playing a video game in a room .
 <ul style="list-style-type: none"> Person: 0.98 White: 0.92 Black: 0.87 Photograph: 0.87 	<ol style="list-style-type: none"> 1. A man looking through a book on top of a table. 2. A man sitting on a bed looking at a book 3. a man is flipping through a book on a bed 4. A man sitting on a bed flipping through pages of a book. 5. A man in a black jacket is flipping through a large book. 	<ul style="list-style-type: none"> • MLE: a man sitting in front of a laptop computer . • MLE-TAG: a person sitting down with a book in front of him . • PG-BCMR: a man sitting in front of a book . • PG-BCMR-TAG: a man sitting at a table with a computer . • PG-SPICE: a man sitting in front of a book and a laptop on a table with a laptop computer on top of a table with a laptop computer on top of • PG-SPIDER: a man sitting at a table with a book . • PG-SPIDER-TAG: a man sitting in front of a book .

Table 1: Example captions from different models on MSCOCO hold-out validation images. For clarity, only top-4 visual tags are shown per example.

Submissions	CIDEr-D	Meteor	ROUGE-L	BLEU-1	BLEU-2	BLEU-3	BLEU-4
MSM@MSRA [26]	0.984	0.256	0.542	0.739	0.575	0.436	0.330
Review Net [25]	0.965	0.256	0.533	0.720	0.550	0.414	0.313
ATT [27]	0.943	0.250	0.535	0.731	0.565	0.424	0.316
Google [21]	0.943	0.254	0.530	0.713	0.542	0.407	0.309
Berkeley LRCN [7]	0.921	0.247	0.528	0.718	0.548	0.409	0.306
MIXER* [14]	-	-	-	-	-	-	0.292
Montreal/Toronto [24]	0.865	0.241	0.516	0.705	0.528	0.383	0.277
Ours: MLE	0.947	0.251	0.531	0.724	0.552	0.405	0.294
Ours: PG-BCMR	1.013	0.257	0.55	0.754	0.591	0.445	0.332
Ours: PG-SPIDEr	1.000	0.251	0.544	0.743	0.578	0.433	0.322
Ours: MLE-TAG	0.949	0.253	0.531	0.722	0.551	0.405	0.296
Ours: PG-BCMR-TAG	1.002	0.254	0.550	0.753	0.593	0.447	0.333
Ours: PG-SPIDEr-TAG	1.042	0.255	0.551	0.751	0.591	0.445	0.331

Table 2: Automatic evaluation on the official MSCOCO C-5 test split. (*) Note that the quoted BLEU-4 score for MIXER is computed on a subset of 5000 MSCOCO validation set, and hence is not directly comparable to other evaluations results in this table.

p-value($X > Y$)	PG-SPIDEr-TAG	PG-BCMR	PG-BCMR-TAG	MLE	MLE-TAG
PG-SPIDEr	0.090	0.014	0.001	<0.001	<0.001
PG-SPIDEr-TAG	-	0.423	0.106	0.002	0.003
PG-BCMR	-	-	0.157	0.003	0.002
PG-BCMR-TAG	-	-	-	0.040	0.017
MLE	-	-	-	-	0.401

Table 3: p-values derived from a pairwise sign test applied to human ratings on 492 images from MSCOCO test set. Statistically significant comparisons (at the 0.05 level) are shown in bold. X and Y correspond to rows and columns respectively.

that any image captioning system we use does not make bad or embarrassing mistakes.)

As a sanity check, we first evaluated 505 ground truth captions from the validation set. Humans said that 87% of these captions were “not bad”. Some of the 13% of captions that were labeled “bad” do indeed contain errors⁶, due to the fact that MSCOCO captions were generated by AMT workers who are not perfect. On the other hand, some captions seem reasonable to us, but did not meet the strict quality criteria our raters were looking for. In any case, 87% is an approximate upper bound on performance we can hope to achieve on the MSCOCO test set.

We then randomly sampled 492 images from the test set (for which we do not have access to the ground truth captions), and generated captions from all of them using our 6 systems, and sent them for human evaluation. Figure 4 shows the fraction of captions that are “not bad” compared to the MLE baseline of 38%. We draw the following conclusions:

- All methods are far below the human ceiling of 87%, and no system is good enough for use in the real world.
- All PG methods outperform MLE training by a very significant margin (see Table 3 for pairwise p-value analysis). This is because the PG methods optimize metrics that are much more closely related to caption quality than the likelihood score.
- PG-SPIDER outperforms PG-BCMR by a 4% margin, despite the fact that PG-BCMR outperforms PG-SPIDER on all the automatic metrics. This is because SPIDER captures both fluency and semantic properties of the caption, both of which human raters are told to pay attention to, whereas BCMR is a more syntactic measure.
- Adding the image tagger seems to hurt performance, even though PG-SPIDER-TAG is superior to PG-SPIDER according to automatic metrics. We are not entirely sure of the reason for this, and plan to investigate it in future work.

5 Conclusion and future work

In this paper, we have shown how using policy gradient methods to optimize the BCMR and SPICE metrics results in much better performance than standard maximum likelihood training. However, we have also seen that none of these existing metrics is perfect; in fact, to get our best results (as judged by humans) we had to combine the SPICE and CIDEr metrics. In future work, we would like

⁶ For example, some captions contain repeated words, e.g., “this is an image image of a modern kitchen”. Others contain typos, e.g., “a blue and white truck with pants in it’s flat bed”. And some do not make semantic sense, e.g., “A crowd of people parked near a double decker bus”.

to explore the possibility of using sequence GANs [28] to automatically learn to distinguish good from bad captions, without having to specify the quality metric explicitly.

We would also like to explore alternatives to the simple PG method we used in this paper. For example, we might investigate an actor-critic method, similar to [2], or the ML-PG hybrid described in [12].

Finally, from a computer vision point of view, we would like to explore the use of more sophisticated representations of the scene, going beyond the generic encoding provided by the CNN. Our use of visual tags was a first step in that direction, but this did not seem to provide much benefit. We would like to investigate this further, by examining performance on other datasets, and using methods that are more sophisticated than image tagging, such as object detectors (cf. [20]).

References

- [1] P. Anderson, B. Fernando, M. Johnson, and S. Gould. Spice: Semantic propositional image caption evaluation. In *ECCV*, 2016. 2, 4
- [2] D. Bahdanau, P. Brakel, K. Xu, A. Goyal, R. Lowe, J. Pineau, A. Courville, and Y. Bengio. An Actor-Critic algorithm for sequence prediction. *Arxiv*, 24 July 2016. 2, 3, 9
- [3] S. Banerjee and A. Lavie. METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In *Proc. ACL workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*, volume 29, pages 65–72, 2005. 1
- [4] S. Bengio, O. Vinyals, N. Jaitly, and N. Shazeer. Scheduled sampling for sequence prediction with recurrent neural networks. In *NIPS*, 2015. 1
- [5] R. Bernardi, R. Cakici, D. Elliott, A. Erdem, E. Erdem, N. Ikizler-Cinbis, F. Keller, A. Muscat, and B. Plank. Automatic description generation from images: A survey of models, datasets, and evaluation measures. *J. of AI Research*, 55:409–442, 2016. 1, 2
- [6] F. Chollet. Information-theoretical label embeddings for large-scale image classification. *Arxiv*, 19 July 2016. 2, 5
- [7] J. Donahue, L. Anne Hendricks, S. Guadarrama, M. Rohrbach, S. Venugopalan, K. Saenko, and T. Darrell. Long-term recurrent convolutional networks for visual recognition and description. In *CVPR*, 2015. 8
- [8] L. A. Hendricks, Z. Akata, M. Rohrbach, J. Donahue, B. Schiele, and T. Darrell. Generating visual explanations. *ECCV*, 2016. 2
- [9] F. Huszár. How (not) to train your generative model: Scheduled sampling, likelihood, adversary? *Arxiv*, 16 Nov. 2015. 1
- [10] C.-Y. Lin and E. Hovy. Automatic evaluation of summaries using n-gram co-occurrence statistics. In *NAACL*, pages 71–78, 2003. 1
- [11] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. Lawrence Zitnick. Microsoft

- COCO: Common objects in context. In *ECCV*, 1 May 2014. 2
- [12] M. Norouzi, S. Bengio, Z. Chen, N. Jaitly, M. Schuster, Y. Wu, and D. Schuurmans. Reward augmented maximum likelihood for neural structured prediction. In *NIPS*, 2016. 2, 9
- [13] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu. BLEU: A method for automatic evaluation of machine translation. In *Proc. ACL*, pages 311–318, 2002. 1
- [14] M. Ranzato, S. Chopra, M. Auli, and W. Zaremba. Sequence level training with recurrent neural networks. *Arxiv*, 2015. 1, 2, 6, 8
- [15] S. Shen, Y. Cheng, Z. He, W. He, H. Wu, M. Sun, and Y. Liu. Minimum risk training for neural machine translation. In *Proc. ACL*, 2016. 2
- [16] R. S. Sutton, D. Mc Allester, S. Singh, and Y. Mansour. Policy gradient methods for reinforcement learning with function approximation. In *NIPS*, 1999. 1, 3
- [17] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna. Rethinking the inception architecture for computer vision. *arXiv preprint arXiv:1512.00567*, 2015. 4
- [18] K. Tran, X. He, L. Zhang, J. Sun, C. Carapcea, C. Thrasher, C. Buehler, and C. Sienkiewicz. Rich image captioning in the wild. In *CVPR*, 2016. 2
- [19] R. Vedantam, C. Lawrence Zitnick, and D. Parikh. Cider: Consensus-based image description evaluation. In *CVPR*, pages 4566–4575, 2015. 1
- [20] S. Venugopalan, L. A. Hendricks, M. Rohrbach, R. Mooney, T. Darrell, and K. Saenko. Captioning images with diverse objects. *Arxiv*, 24 June 2016. 2, 9
- [21] O. Vinyals, A. Toshev, S. Bengio, and D. Erhan. Show and tell: A neural image caption generator. In *CVPR*, 2015. 2, 4, 5, 8
- [22] R. J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine Learning J.*, 8(3-4):229–256, 1 May 1992. 2, 3
- [23] Q. Wu, C. Shen, A. van den Hengel, L. Liu, and A. Dick. What value high level concepts in vision to language problems? In *CVPR*, 2016. 2, 5
- [24] K. Xu, J. Ba, R. Kiros, K. Cho, A. Courville, R. Salakhutdinov, R. Zemel, and Y. Bengio. Show, attend and tell: Neural image caption generation with visual attention. In *ICML*, 2015. 2, 8
- [25] Z. Yang, Y. Yuan, Y. Wu, R. Salakhutdinov, and W. W. Cohen. Review networks for caption generation. In *NIPS*, 2016. 2, 8
- [26] T. Yao, Y. Pan, Y. Li, Z. Qiu, and T. Mei. Boosting image captioning with attributes. In *OpenReview*, 2016. 2, 5, 8
- [27] Q. You, H. Jin, Z. Wang, C. Fang, and J. Luo. Image captioning with semantic attention. *arXiv preprint arXiv:1603.03925*, 2016. 8
- [28] L. Yu, W. Zhang, J. Wang, and Y. Yu. SeqGAN: Sequence generative adversarial nets with policy gradient. *Arxiv*, 18 Sept. 2016. 2, 3, 9
- [29] W. Zaremba and I. Sutskever. Reinforcement learning neural turing machines-revised. *arXiv preprint arXiv:1505.00521*, 2015. 3