

Consensus-based Sequence Training for Video Captioning

Sang Phan Gustav Eje Henter Yusuke Miyao Shin’ichi Satoh
National Institute of Informatics, Japan
{plsang, gustav, yusuke, satoh}@nii.ac.jp

Abstract

Captioning models are typically trained using the cross-entropy loss. However, their performance is evaluated on other metrics designed to better correlate with human assessments. Recently, it has been shown that reinforcement learning (RL) can directly optimize these metrics in tasks such as captioning. However, this is computationally costly and requires specifying a baseline reward at each step to make training converge. We propose a fast approach to optimize one’s objective of interest through the REINFORCE algorithm. First we show that, by replacing model samples with ground-truth sentences, RL training can be seen as a form of weighted cross-entropy loss, giving a fast, RL-based pre-training algorithm. Second, we propose to use the consensus among ground-truth captions of the same video as the baseline reward. This can be computed very efficiently. We call the complete proposal Consensus-based Sequence Training (CST). Applied to the MSRVT video captioning benchmark, our proposals train significantly faster than comparable methods and establish a new state-of-the-art on the task, improving the CIDEr score from 47.3 to 54.2.

1. Introduction

The goal of video captioning is to automatically generate informative and human-like text descriptions of given videos. However, video content is generally very diverse, and so are human-annotated captions. In Fig. 1 we show examples of different descriptions of the same video in the training set of a popular video-captioning dataset. We see that not all captions are created equal – different persons may pay attention to different segments and aspects of the video, and sometimes annotators may even make mistakes.

Since there is no single “canonical caption” for a given video, several works propose automatic evaluation metrics that aim to correlate with human judgment [20, 6, 18, 32], enabling the quality of different captions to be compared. Fig. 1 displays the CIDEr score [32], which aims to measure the consensus among annotators. However, there has been less effort to develop captioning systems that directly

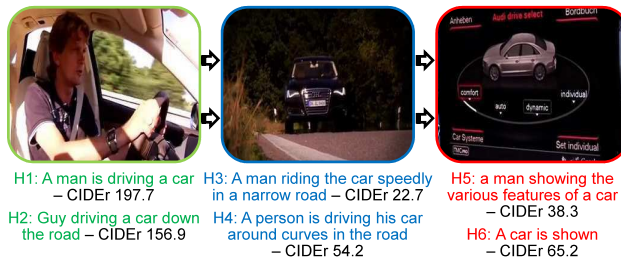


Figure 1: Examples of human-annotated captions and their different CIDEr consensus scores for a single video.

aim to generate human-like captions by performing well on these metrics. Instead, most video-captioning systems are trained by the principle of Maximum Likelihood Estimation (MLE), also known as Cross-Entropy (XE) minimization. This is, unfortunately, not ideal: several works [2, 10, 30] have shown that MLE/XE emphasizes accurately describing outliers, instead of the most typical case. This makes the outcome of XE training sensitive to unusual or aberrant captions, rather than optimizing for stable output around the human consensus of what an appropriate caption would be.

A compelling alternative to MLE/XE is to maximize the objective of interest directly. This can be done by through the reinforcement learning (RL) framework [29] using a method such as REINFORCE [33]. In RL, the score of a candidate sentence is used as a reward signal, and the model attempts to maximize this reward. If the chosen reward metric is the CIDEr score [32], which is designed to approximate human judgment of appropriateness and consensus, the model may be taught to generate more human-like captions. However, REINFORCE training is difficult because it requires estimating the gradient of the expected reward, which is not stable using mini-batch estimation.

To our knowledge, only one published work, Pasunuru & Bansal [21], has applied REINFORCE to video captioning. In the process, they obtained top scores on a standard captioning benchmark. They used a MIXER scheme [24], that gradually mixes RL training into XE training, to stabilize the learning. This is known to be very sensitive to the

annealing schedule used [19].

MLE/XE and RL training each have pros and cons. XE training is much faster than RL, but is unlikely to produce human-like captions. RL training can directly optimize objective functions that reward human-likeness, but the gradient with respect to the expected reward is unstable without costly variance reduction. In this paper, we propose a new training approach that combines the advantages of both training schemes. Our main contributions are:

1. We show that training with the XE objective is a special case of REINFORCE training when applied to sentences obtained from the ground-truth data, rather than samples from the model being trained.
2. Based on the previous insight, we introduce a simple weighted XE pre-training scheme (WXE) that approximates RL training. Our proposal is equally stable and fast to train as conventional XE, yet provides a 2-point improvement in terms of the CIDEr metric. Our pre-training addresses a well-known problem of XE known as objective mismatch [24], but (unlike standard RL) not a related issue known as exposure bias (see Sec. 2), thus allowing the losses incurred by these two important problems to be assessed separately.
3. We further propose to fine-tune our pre-trained models through full REINFORCE, but using the consensus scores of ground-truth captions as the baseline reward. We call this the Self-Consensus Baseline, or SCB. SCB trains twice as fast as estimating the baseline reward through the greedy method [25], and provides greater CIDEr score improvements. Together, pre-training and fine-tuning address both objective mismatch and exposure bias, and form our proposed Consensus-based Sequence Training, or CST.

We conduct thorough experiments to demonstrate the advantage of CST on the large-scale MSRVT [34] benchmark. Our results on the MSRVT video-captioning dataset improve the results in terms of CIDEr from 47.3 (XE) to 54.2 (full CST), surpassing the best previously published score of 51.7 to establish a new state-of-the-art on the task.

2. Related Work

Issues with Standard Captioning Models In line with other deep text-generation models, captioning models are typically trained to maximize the likelihood of the next word given the previous ground-truth input words (i.e., MLE/XE). This approach has two main drawbacks:

First is the *objective mismatch* problem. This arises because the objective function optimized at training time is not the true objective we are interested in, or rather, is not equal to the metric used to quantify success on the task. Conventional maximum likelihood, in particular, is especially poorly suited for practical applications. Mathematically, it is similar to minimizing a weighted squared error

in which accuracy on the least probable captions is given the greatest weight [2]. As highlighted in [10], this is a poor fit for data-generation applications like captioning, since it means that accuracy at the maximum-probability point – the greedy output eventually returned by the system – is given the lowest priority of all during training. It also is sensitive to poor or erroneous captions in the data; it is not *statistically robust* [12]. Randomly sampled (rather than greedy) output is also affected, and models fit using MLE are likely to produce samples that humans perceive as unnatural [30].

To overcome the objective mismatch, one would want to train models to directly optimize the performance metric of interest, which in captioning typically is a discrete NLP metric such as BLEU [20], METEOR [6], ROUGE-L [18] or CIDEr [32]. Unfortunately, these metrics operate on full sentences (in contrast to model output, which is generated iteratively over words) and are not differentiable, which prevents training using stochastic gradient descent.

The second problem is *exposure bias* [24], which is the input distribution mismatch between training and testing time. During training, the model is only exposed to word sequences from the training data (so called teacher forcing), while at testing time, the model instead only has access to its own predictions; therefore, whenever the model encounters a state it has never been exposed to before, it may behave unpredictably for the rest of the output sequence. Methods such as scheduled sampling [3] propose to gradually expose models to input words from the model distribution rather than the ground truth. However, a major limitation of this and related approaches like professor-forcing [17] is that, at each time-step, the output word needs to be same as in the ground truth, implicitly discouraging the model from generating novel captions.

REINFORCE Baseline Estimators Recently, the REINFORCE [33] method has been applied to captioning tasks [24, 19, 9, 25, 1, 21]. This is an RL algorithm that can optimize any metric of interest and trains on sampled sequences, thus overcoming the limitations of maximum likelihood, scheduled sampling, and professor forcing above. Training with REINFORCE is considerably difficult because the gradient of the expected reward is estimated using a single sample, which exhibits high variance. By estimating an expected baseline reward and subtracting this baseline from the estimated gradient, one can reduce the variance considerably without changing the expected gradient [33]. Most existing work use parametric functions to estimate the expected baseline reward [24, 19]; this can be implemented as a separate neural network which takes the hidden state of the captioning model (neural network) as input at each time step. Instead of estimating the reward signal, Self-Critical Sequence Training (SCST) [25] utilizes the (greedy) output of the current model at inference time as the baseline, at the cost of having to compute the score of

each baseline sequence. Taking a cue from the up-down model [1], our proposal trains on multiple captions of the same video simultaneously; however, we propose a Self-Consensus Baseline (SCB), which simply computes the average reward among ground-truth captions as the baseline, rather than generating and scoring a new greedy baseline for every sentence and epoch during training as in [25].

Consensus-based Caption Evaluation As human evaluation is expensive, captioning performance is commonly assessed using automatic measures that approximate human judgments. We will use the leading CIDEr metric [32], designed to quantify similarity to the consensus among captions better than other metrics. Optimizing CIDEr should encourage output that is similar to typical human captions.

3. Captioning Models

Most captioning models are based on encoder-decoder neural networks. The encoder takes pre-extracted multi-modal features from the video and embed these into a fixed-size space, with the final video representation v being the concatenation of all the embedded features. The decoder receives video features from the encoder and generates an output sentence using a Recurrent Neural Network (RNN). Among RNNs, the Long Short-Term Memory (LSTM) [11] architecture is widely used, since it yields state-of-the-art results in many different applications. At a high level, the computation at each time-step of the LSTM can be expressed in terms of the previously generated word w_{t-1} and hidden state vector h_{t-1} as follows:

$$h_t = LSTM(h_{t-1}, w_{t-1}) \quad (1)$$

Let $o_t = W_o h_t$ be the vocabulary-sized output of the LSTM, let θ denote the parameters of the model, and suppose h_t maintains information on the input video and the previous words, the distribution of the next word is then given by a softmax function:

$$p_\theta(w_t|h_t) = \text{softmax}(o_t) \quad (2)$$

At training time, w_t is typically a ground-truth token; however, one can also perform sampling from the $p_\theta(w_t|h_t)$ distribution. At testing time, w_t is often selected by a greedy, iterative approach such as beam search [28], which aims to identify the most likely caption (word sequence).

3.1. Cross Entropy Training (XE)

Let $w = (w_1, w_2, \dots, w_T)$ be a target, ground-truth sequence. Typically the parameters θ of the model are learned by minimizing the cross entropy (XE) loss:

$$L^{XE}(\theta) = -\log p_\theta(w) \quad (3)$$

$$= -\sum_{t=1}^T \log p_\theta(w_t|h_t) \quad (4)$$

3.2. REINFORCE Training (RL)

The captioning process can be formulated as a Reinforcement Learning problem [24], where the agent is the LSTM language model interacting with an environment of words and image/video features. The model p_θ defines a policy network, in which each action corresponds to predicting the next word.

Mathematically, the goal of REINFORCE training is to minimize the negative expected reward of model samples:

$$L^{RL}(\theta) = -\mathbb{E}_{w^s \sim p_\theta}[r(w^s)] = -\sum_{w^s} r(w^s) p_\theta(w^s), \quad (5)$$

where $w^s = (w_1^s, w_2^s, \dots, w_T^s)$, with w_t^s being the word sampled from the model p_θ at time t . To minimize the negative expected reward, one might differentiate (5) as:

$$\nabla_\theta \mathbb{E}_{w^s \sim p_\theta}[r(w^s)] = \mathbb{E}_{w^s \sim p_\theta}[\nabla_\theta r(w^s)] \quad (6)$$

However, this direct approach runs into a problem if the reward function $r(w^s)$ is non-differentiable, i.e., $r(w^s)$ is not a continuous function with respect to θ . Alternatively, one can differentiate the sum in (5) as:

$$\nabla_\theta \mathbb{E}_{w^s \sim p_\theta}[r(w^s)] = \nabla_\theta \sum_{w^s} r(w^s) p_\theta(w^s) \quad (7)$$

$$= \sum_{w^s} r(w^s) \nabla_\theta p_\theta(w^s) \quad (8)$$

$$= \sum_{w^s} p_\theta(w^s) r(w^s) \frac{\nabla_\theta p_\theta(w^s)}{p_\theta(w^s)} \quad (9)$$

$$= \mathbb{E}_{w^s \sim p_\theta}[r(w^s) \nabla_\theta \log p_\theta(w^s)]; \quad (10)$$

this expression does not require differentiating $r(w^s)$ w.r.t. θ and can be estimated through Monte Carlo sampling.

In practice, gradients estimated based on (10) are unstable and it is necessary to perform variance reduction using a baseline estimator b (cf. [24]):

$$\nabla_\theta L^{RL}(\theta) = -\mathbb{E}_{w^s \sim p_\theta}[(r(w^s) - b) \nabla_\theta \log p_\theta(w^s)] \quad (11)$$

$$\approx -(r(w^s) - b) \nabla_\theta \log p_\theta(w^s) \quad (12)$$

It is easy to show that the gradient in (11) is unchanged as long as b does not depend on the sample w^s .

4. Consensus-based Sequence Training

This section details the novel contributions of this paper. First, we present a new theoretical connection between XE and RL training. We then expand on this insight to formulate an entire training scheme that utilizes the consensus among the multiple ground-truth captions for each video to improve training in several ways. As a bonus, the two stages of our scheme disentangle the objective mismatch and the exposure bias problems described in [24].

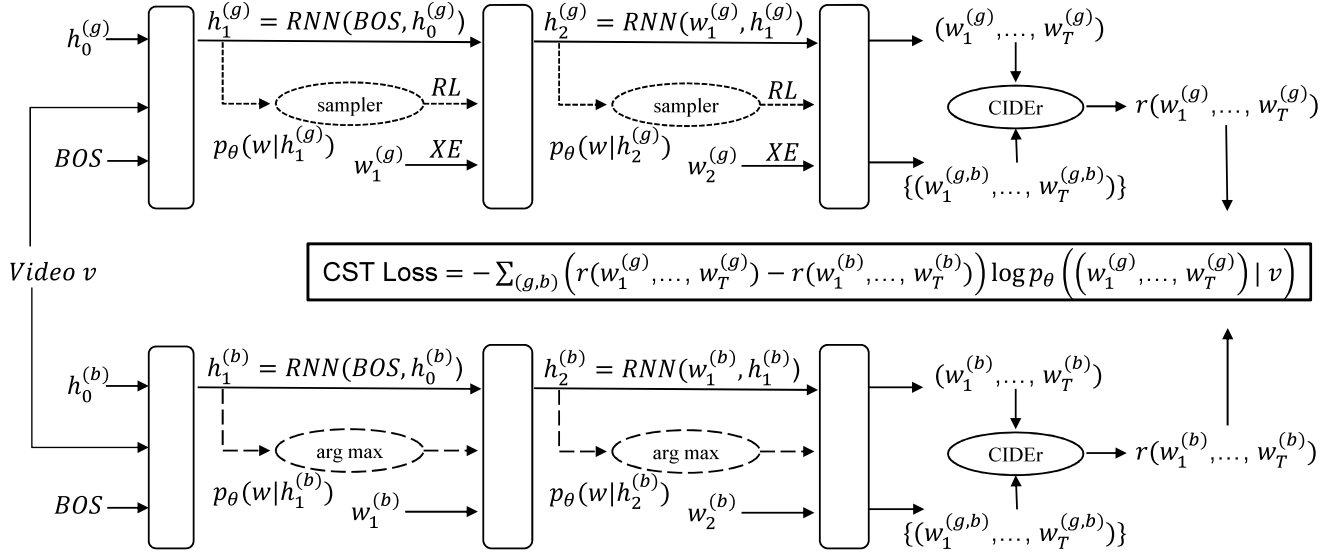


Figure 2: Overview of Consensus-based Sequence Training (CST). CST can be used to train both XE and RL schemes. Note that the baseline reward is estimated using the consensus scores of *training* captions, e.g., (g) and (b) , which is very efficient. Dashed connections indicate the SCST method [25].

4.1. Relation between XE Training and RL Training

Using the chain rule, the derivatives of both the XE and RL loss w.r.t. the parameters θ are:

$$\frac{\partial L(\theta)}{\partial \theta} = \sum_{t=1}^T \frac{\partial L(\theta)}{\partial o_t} \frac{\partial o_t}{\partial \theta} \quad (13)$$

The partial derivative of each loss w.r.t. the softmax input o_t is given by [24]:

$$\frac{\partial L^{XE}(\theta)}{\partial o_t} = p_\theta(w_t | h_t) - 1(w_t) \quad (14)$$

$$\frac{\partial L^{RL}(\theta)}{\partial o_t} = (r(w^s) - b)(p_\theta(w_t^s | h_t) - 1(w_t^s)) \quad (15)$$

The gradient is the difference between the prediction and the one-hot representation, represented by the $1(\cdot)$ indicator function, of the target word.¹ In RL training, this gradient is further weighted by a stabilized reward $r(w^s) - b$. In that case the sampled word w_t^s can be either encouraged or penalized, in contrast to the ground-truth word w_t that's always being encouraged in XE training.

There are interesting similarities between Eq. (15) and Eq. (14). In particular, in the situation $w_t^s = w_t$ we obtain:

$$\frac{\partial L^{RL}(\theta)}{\partial o_t} = (r(w) - b) \frac{\partial L^{XE}(\theta)}{\partial o_t}, \quad (16)$$

¹We provide a full derivation in the supplementary material.

Therefore, when training RL on the *ground-truth captions*, we can compute the terms in the $L^{RL}(\theta)$ loss as follows:

$$L^{RL}(\theta | w) = (r(w) - b)L^{XE}(\theta | w), \quad (17)$$

Eq. (17) suggests that RL training on ground-truth captions (which happens when the captioning system behaves like the human annotators) is a generalization of XE training where each sentence loss is weighted by a reward. More generally, we expect that systems with human-like output can be trained to an approximate optimum on ground-truth captions w_t without having to draw samples w_t^s and evaluate their reward. In cases where sampling or reward calculation is a computational bottleneck (like with CIDEr), this can produce a substantial speed improvement. The trade-off is that no exploration of non-ground-truth sentences is performed, so resulting models may still suffer from exposure bias; however, the weighting still mitigates objective mismatch, meaning that the $w_t^s \approx w_t$ approximation neatly separates the two issues originally recognized in XE, and represents a middle ground between XE and REINFORCE.

Mathematically speaking, how can a weighted cross-entropy be helpful for learning human consensus? Consider a video v described by the caption set $S = \{w^{(1)}, \dots, w^{(N_v)}\}$, N_v being the number of captions. Then:

$$L_v^{RL}(\theta) = -\frac{1}{N_v} \sum_{i=1}^{N_v} (r(w^{(i)}) - b^{(i)}) \log p_\theta(w^{(i)}) \quad (18)$$

Intuitively, this loss increases the probability of captions that have high-rewards. If we choose r to be the CIDEr met-

ric – designed to measure the consensus among captions – this loss assigns lower weight to the model’s log-likelihood performance on captions with low consensus score (low reward), reducing their impact on the training outcome, in favor of consistently generating more typical captions. This differs from conventional, unweighted XE/MLE, where models typically incur severe log-likelihood penalties for failing to explain bad or atypical captions in training data.

4.2. Self-consensus Baseline

Notably, our pre-training (unlike traditional RL) is stable to train without a baseline b , since the use of ground-truth captions significantly constrains the sentence space. Full RL training, however, still requires a baseline estimator for variance reduction. Observing that (1) the theoretically optimal baseline b is the expected reward of the current policy and (2) leading automatic captioning systems achieve broadly similar CIDEr scores as the average of held-out human captions (52.9 vs. 50.2 in [15]), we propose using the rewards of ground-truth sentences as the baseline. Specifically, we use the average reward of all the sentences describing the same video to compute b . Using the notation in 4.1, the baseline reward for each video v is given by:

$$b(v) = \frac{1}{N_v} \sum_{i=1}^{N_v} r(w^{(i)}) \quad (19)$$

The cost to estimate our baseline during training is negligible as the rewards can be pre-computed. Since this baseline leverages the mutual consensus among captions of the same video, we call it the Self-Consensus Baseline, or SCB.

4.3. Full Training Scheme Proposal

We unite our proposals of an improved pre-training stage and consensus-based SCB in RL training under the label Consensus-based Sequence Training, or CST. Figure 2 depicts how the SCB is computed in CST and the main difference with SCST [25], while Algorithm 1 summarizes the complete CST method. At first, we propose to “warm start” the model by pre-training it on the ground-truth data using the weighted cross-entropy loss defined in Eq. (18). This is similar in implementation to standard XE training, except that the log-likelihood of each ground-truth word now is weighted by the CIDEr score of the corresponding caption. We call this method WXE, for Weighted Cross-Entropy. Warm-starting is typical RL practice to prevent models from making random actions in the exponentially large sentence-distribution space at the beginning of training, except that conventional implementations pre-train the model using the unweighted XE loss. As we show in Sec. 6, our new WXE warm-start strategy gives better objective scores than XE training while being equally fast, and provides a better starting point for later, full REINFORCE training.

Algorithm 1: CST

Data: $\mathbb{D} = \{(v^n, w^n)\}$, with $n = 1 \dots N$, is the set of video-caption pairs;

- 1 Pre-compute the CIDEr scores of all captions in \mathbb{D} ;
- 2 **for each epoch do**
- 3 **for each video v do**
- 4 **if WXE (Pre-)Training then**
- 5 Get caption w from the ground truth;
- 6 Get $r(w)$ from the pre-computed scores;
- 7 Set $b = 0$;
- 8 **else if RL Training then**
- 9 Generate sample sequence $w \sim p_\theta(\cdot|v)$;
- 10 Compute $r(w)$ using CIDEr;
- 11 Compute baseline reward b using Eq. (19);
- 12 Compute $L_v^{RL}(\theta)$ using Eq. (18);
- 13 Back-propagate to compute $\frac{\partial L^{RL}(\theta)}{\partial \theta}$;
- 14 SGD-update θ using $\frac{\partial L^{RL}(\theta)}{\partial \theta}$;

After pre-training with WXE, we perform full RL training, in which candidate captions are sampled from the estimated model distribution. This is significantly slower per epoch due to the need to compute the CIDEr scores of model samples, but enables the model to explore more of the possible sentences that can be generated.

By dividing training into two steps, (1) pre-training using weighted cross-entropy on the ground truth; and (2) RL training on samples from the model, we address both issues of conventional XE training. In (1), we handle the objective mismatch by directly optimizing the CIDEr metric, though still on the ground-truth data. The resulting performance illustrates how much we may improve the captioning model by simply changing the training objective from XE towards CIDEr, without affecting training time. In (2), we allow the model to generate samples and provide feedback through the reward signal. This additionally addresses the exposure bias of MLE, since the sequences evaluated now are sampled from the captioning model.

5. Experimental Setup

5.1. Dataset and Features

We evaluated the performance of our CST method and its variants on the Microsoft Video-to-Text dataset (MSRVTT) [34] popular in captioning. This dataset is composed of 10,000 videos in 20 different categories. Each video is annotated by 20 different persons. We used the data splits provided in [34] which contain 6,513, 497, and 2,990 videos for training, validation, and testing, respectively.

The following features from different modalities were extracted from the videos and used as system inputs:

- ResNet [8] is an extremely deep network to extract static image features. We used the ResNet 200-layer model trained on the ImageNet dataset to extract image feature from the fc7 layer (2,048-dim) every 8 frames.
- C3D [31] captures short-term motion over 16 video frames. We extracted features at the fc6 layer of a model trained on the Sports-1M dataset, and applied mean pooling to represent each video chunk.
- MFCC [5] acoustic features were extracted from 25 ms audio segments with a frame step of 10 ms. The 13 first MFCC coefficients and their Δ and Δ^2 features were then encoded with a codebook size of 256 clusters using VLAD [13], yielding a 9,984-dim representation.
- Category embedding. Each video in the MSRVT dataset belongs to one of 20 categories such as music, food, or gaming [34]. We used the word-embedding model proposed in GloVe [22] to encode the category label of each video into a 300-dimensional vector.

5.2. Implementation Details

A vocabulary of 10,533 words was created from all words appearing more than thrice in the data. Three special tokens were added to the vocabulary: $\langle \text{BOS} \rangle$, $\langle \text{EOS} \rangle$, and $\langle \text{UNK} \rangle$, representing the beginning and end of a sentence, and unknown tokens, respectively. We did not apply any text pre-processing except for lowercase conversion. Word vectors of size 512 were used to encode previously-generated words (actions) for input to the LSTMs when generating the next word. We use a single-layer unidirectional LSTM with a 512-dimensional hidden state for captioning. Following [35], we inserted a rate-0.5 dropout layer on the input and output of each LSTM unit during training.

In each iteration, the encoder loaded a mini-batch of 64 videos with pre-extracted multimodal features. Each feature was then projected into an embedding using a linear layer of dimension 512. We trained on all 20 ground-truth/sampled captions of each video at the same time. The Adam [16] optimizer was used for training, with a fixed learning rate of 1×10^{-4} in all experiments. XE/WXE pre-training was ran for maximum 50 epochs; full RL training continued until epoch 200, with the best result reported.

We used beam search [28] with a beam size of 5 to find highly probable output sentences for evaluation. Training and evaluation were implemented with PyTorch, with objective scores calculated using the MS COCO toolkit [4].²

5.3. Training Methods Considered

For the experiments, we trained the same basic captioning LSTM using a number of variants of CST, gradually stepping from conventional XE/MLE training to full CST. The different training configurations are introduced below, with an overview provided in Table 1:

XE: This LSTM was trained on our extracted features and network setup using conventional maximum-likelihood, as a comparison point, and a bottom line against which to judge our improvements.

CST_GT_None: This is the CST pre-training paradigm, in which the reward-weighted cross entropy in Eq. (18) is used for training; for this reason, this training method is also referred to as **WXE**, for short. This weighting alleviates the objective mismatch problem of XE. The difference from full REINFORCE is that the loss is evaluated on ground-truth captions (“GT”) rather than captions sampled from the model being trained. (The word “None” signifies that no baseline b was used.)

CST_MS_SCB: This is a full REINFORCE (RL) training paradigm, where training gradients are computed on model samples (“MS”) as in Eq. (12). Using samples instead of ground-truth captions should address both the objective mismatch and the exposure bias problems simultaneously. For faster convergence, CST_MS_SCB training was always initialized from CST_GT_None/WXE above. Since REINFORCE training is unstable out of the box in video captioning, we used a baseline b based on self-consensus (SCB introduced in Sec. 4.2) – the average held-out CIDEr score of ground-truth captions – to stabilize training. This is significantly simpler than the MIXER approach [24] used in [21], needing no mixing scheme tuning or annealing.

CST_MS_Greedy: This is another full RL training paradigm, identical to CST_MS_SCB except that instead of basing b on the consensus among ground-truth captions, b was set equal to the score (reward) of a highly-likely sentence generated by iteratively (greedily) generating the most probable word under the model, as proposed in [25].

SCST: This is the same as CST_MS_Greedy, except that (like in [25]) full RL training starts from conventional XE instead of WXE. It thus constitutes a reimplementaion of [25], but applied to video rather than image captioning.

6. Experimental Results

6.1. Performance on Different Feature Sets

To begin with, we performed a preliminary experiment comparing conventional XE and CIDEr-weighted WXE when trained on each of the features in Sec. 5.1 separately, as well as on the full, multimodal feature set. The resulting systems were compared in terms of seven different objective metrics (four BLEU scores and three other metrics, including CIDEr), with the numerical results tabulated in supplementary material. Except when training only on the Category features, WXE-models outperformed XE baselines in all aspects, confirming the efficacy of the WXE pre-training scheme over conventional XE. Since the full feature set gave the best performance across the board, we used this input feature set in all subsequent experiments.

²<https://github.com/tylin/coco-caption>

Table 1: Overview of the most important training methods considered in the experiments.

Label	Start from	Sequences used in training	Baseline b	Comment
XE	Scratch	Ground truth	None	Unweighted log-likelihood; bottom line
CST_GT_None	Scratch	Ground truth (“GT”)	None	Weighted log-likelihood; a.k.a. “WXE”
CST_MS_Greedy	WXE	Model samples (“MS”)	Greedy [25]	
CST_MS_SCB	WXE	Model samples	SCB	“Full CST”
CIDEnt-RL [21]	XE	Truth & samples (MIXER)	Lin. reg. [24]	Previous best
SCST [25]	XE	Model samples	Greedy	Our implementation

Table 2: Per-batch training time for different methods.

Model	Wall-clock time (s)
XE	0.677
CST_GT_None	0.594
CST_MS_Greedy	10.377
CST_MS_SCB	5.256

6.2. Training Time

Having selected the feature set, we trained systems using each of the methods in Sec. 5.3. Table 2 reports the wall-clock time per batch for training each of these systems except SCST, which performed similarly to CST_MS_Greedy. The major computational bottleneck is evaluating CIDEr scores. The WXE and XE methods are by far the fastest (and essentially equally fast), since any required CIDEr scores can be pre-computed. The computation of the reward (CIDEr score) of sentences sampled in the inner loop of CST_MS makes these full RL schemes run much slower. The fact that the Greedy baseline of [25] also must compute the CIDEr score of the generated baseline sentence essentially doubles the computational load of CST_MS_Greedy and SCST over our proposed CST_MS_SCB.

6.3. Comparison Against the State of the Art

Table 3 reviews the leading results from previous literature on the MSRVTT video captioning benchmark in terms of four different scoring metrics (BLEU [20], METEOR [6], ROUGE_L [18], and CIDEr [32]), and compares this to the performance of XE and CIDEr-optimizing models trained as described in Sec. 5.3. On three of the four metrics, v2t_navigator [14] – a system based on cross-entropy minimization – achieves the highest score. On the CIDEr metric, our primary interest since it correlates best with human judgments, the situation is different, however: the best score is reported by CIDEnt-RL [21], the only video-captioning publication we are aware of that optimizes CIDEr score directly, through RL. Their score, 51.7, is the highest CIDEr score reported on this benchmark thus far.

Table 3: Comparison between CST variants and leading previous results achieved on the MSRVTT test set.

Method	BLEU_4	METEOR	ROUGE_L	CIDEr
XE and WXE training				
v2t_navigator [14]	42.6	28.8	61.7	46.7
Aalto [27]	39.8	26.9	59.8	45.7
VideoLAB [23]	39.1	27.7	60.6	44.1
ruc-uva [7]	38.7	26.9	58.7	45.9
Dense Caption [26]	41.4	28.3	61.1	48.9
Our XE	43.0	28.7	61.7	47.3
CST_GT_None	44.1	29.1	62.4	49.7
Full RL training				
CIDEnt-RL [21]	40.5	28.4	61.4	51.7
SCST	41.3	28.1	61.9	51.9
CST_MS_Greedy	42.2	28.9	62.3	53.4
CST_MS_SCB	41.9	28.7	62.1	53.6
CST_MS_SCB (*)	41.4	28.8	62.2	54.2

Among the methods implemented for this paper, our conventional XE system performs very similarly (and often slightly better than) v2t_navigator on all metrics. The CST-derived methods improve further compared to this baseline. In particular, our proposed pre-training CST_GT_None (a.k.a. WXE) outperforms all other methods on the metrics that are not CIDEr, while also improving the CIDEr score by more than 2 points. Nonetheless, we are most interested in the CIDEr metric, since it better captures human consensus. We see that fine-tuning CST_GT_None using the two REINFORCE-schemes implemented leads to further improved CIDEr scores. Starting full RL training from WXE rather than XE (CST_GT_Greedy vs. SCST) actually improved performance on all measures, establishing the utility of the WXE pre-training scheme. The highest CIDEr score was attained by our self-consensus baseline. This is appealing, since SCB represents a simple idea and runs about twice as fast as CST_MS_Greedy. Our full CST proposal gives a CIDEr score of 53.6, almost 4 points greater than WXE, and 2 points superior to the best score

Table 4: Performance of CST_GT_None when optimizing different metrics (different reward functions).

Optimization metric	BLEU_1	BLEU_2	BLEU_3	BLEU_4	METEOR	ROUGE_L	CIDEr
XE	82.1	68.6	55.2	43.0	28.7	61.7	47.3
BLEU_4	80.6	67.1	53.9	41.9	28.4	61.1	45.5
METEOR	81.6	68.1	54.8	42.6	28.7	61.4	47.7
ROUGE_L	82.3	69.1	55.8	43.9	29.0	62.1	48.6
CIDEr	82.8	69.5	56.2	44.1	29.1	62.4	49.7

from prior literature. In the last run, denoted by (*), we experimented with computing the SCB from the 20 sequences sampled during training, instead of the 20 GT sequences.³ This reached a new state-of-the-art CIDEr score of 54.2.

6.4. CST for Other Rewards

Even though CST was conceived for optimizing metrics that reward consensus (like CIDEr), we also investigated the efficacy of the CST_GT_None training scheme in optimizing reward functions based on other objective scoring metrics. As shown in Table 4, CIDEr and ROUGE_L are the only two reward functions where CST outperforms XE under all evaluation metrics. Among those two, CIDEr optimization scores the best.

6.5. Qualitative Analysis of Generated Captions

Figure 3 shows examples of generated captions and their CIDEr scores for videos where human reference annotations had a high standard deviation in held-out CIDEr score. These are videos where some human captions are close to a consensus, while others deviate considerably from it. On this type of video, the XE training is likely to be unduly influenced by outlying examples, leading to poor output, as seen in the figure. In contrast, CST, especially with full RL, tends to generate more appropriate captions.

7. Conclusion

This paper proposed a Consensus-based Sequence Training (CST) scheme to generate video captions. CST is a variant of the REINFORCE algorithm, but exploits the multiple existing training-set captions in several new ways. First, based on a new theoretical insight, CST performs an RL-like pre-training, but with captions from the training data replacing model samples. This alleviates the objective mismatch issue in conventional XE training while being equally fast. Second, CST applies REINFORCE for fine-tuning using the consensus (average reward) among training captions as the baseline b . This fine-tuning additionally removes the exposure bias, and trains twice as fast as a greedy baseline

³The gradient estimate in (11) may now be biased, but the computation time is unchanged as we anyhow need to compute CIDEr for the samples.



GT: (H) A man is making a fire with tree limbs (480); (L) Man talking about building a fire (215.9); CIDEr std. = 103.0

XE: A man is playing a video game (2.9)

CST_GT_None: A man is playing a video game (2.9)

CST_MS_SCB: A man is talking about a fire (176.7)

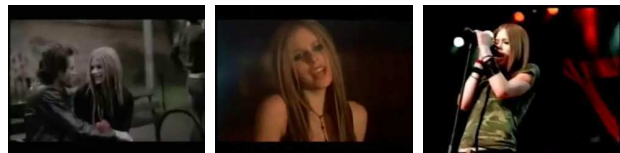


GT: (H) Kids eat and react to fried food (291.5); (L) Children are eating food (4.6); CIDEr std. = 93.2

XE: A man is talking to a woman in a kitchen (0.1)

CST_GT_None: a group of people are eating food (87.3)

CST_MS_SCB: two kids are eating food (103.9)



GT: (H) A female singer is a small figure near the front edge of a dark stage (261.9); (L) In the night party the girl is singh the song she is amzing beauty (4.9); CIDEr std. = 88.1

XE: A group of people are singing and dancing (36.0)

CST_GT_None: A group of people are singing in a music video (47.2)

CST_MS_SCB: A woman is singing a song in a music video (117.3)

Figure 3: Example captions for videos with a broad range of CIDEr scores (from L to H) among human annotators.

estimator. The two stages of CST allow objective mismatch and exposure bias to be assessed separately, and together establish a new state-of-the-art on the task.

References

- [1] P. Anderson, X. He, C. Buehler, D. Teney, M. Johnson, S. Gould, and L. Zhang. Bottom-up and top-down attention for image captioning and VQA. *arXiv preprint*

- arXiv:1707.07998*, 2017. 2, 3
- [2] A. Basu, I. R. Harris, N. L. Hjort, and M. C. Jones. Robust and efficient estimation by minimising a density power divergence. *Biometrika*, 85(3):549–559, 1998. 1, 2
- [3] S. Bengio, O. Vinyals, N. Jaitly, and N. Shazeer. Scheduled sampling for sequence prediction with recurrent neural networks. In *NIPS*, pages 1171–1179, 2015. 2
- [4] X. Chen, H. Fang, T.-Y. Lin, R. Vedantam, S. Gupta, P. Dollár, and C. L. Zitnick. Microsoft COCO captions: Data collection and evaluation server. *arXiv preprint arXiv:1504.00325*, 2015. 6
- [5] S. Davis and P. Mermelstein. Comparison of parametric representations for monosyllabic word recognition in continuously spoken sentences. *IEEE Transactions on Acoustics, Speech, and Signal Processing*, 28(4):357–366, 1980. 6
- [6] M. Denkowski and A. Lavie. Meteor Universal: Language specific translation evaluation for any target language. In *WMT@ACL*, 2014. 1, 2, 7
- [7] J. Dong, X. Li, W. Lan, Y. Huo, and C. G. M. Snoek. Early embedding and late reranking for video captioning. In *ACM MM*, pages 1082–1086, 2016. 7
- [8] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016. 6
- [9] L. A. Hendricks, Z. Akata, M. Rohrbach, J. Donahue, B. Schiele, and T. Darrell. Generating visual explanations. In *ECCV*, pages 3–19, 2016. 2
- [10] G. E. Henter, S. Ronanki, O. Watts, M. Wester, Z. Wu, and S. King. Robust TTS duration modelling using DNNs. In *ICASSP*, pages 5130–5134, 2016. 1, 2
- [11] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural Computation*, 9(8):1735–1780, 1997. 3
- [12] P. J. Huber. *Robust Statistics*. Springer, New York, NY, 2nd edition, 2011. 2
- [13] H. Jegou, F. Perronnin, M. Douze, J. Sánchez, P. Perez, and C. Schmid. Aggregating local image descriptors into compact codes. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 34(9):1704–1716, 2012. 6
- [14] Q. Jin, J. Chen, S. Chen, Y. Xiong, and A. Hauptmann. Describing videos using multi-modal fusion. In *ACM MM*, pages 1087–1091, 2016. 7
- [15] Q. Jin, S. Chen, J. Chen, and A. Hauptmann. Knowing yourself: Improving video caption via in-depth recap. In *ACM MM*, 2017. 5
- [16] D. Kingma and J. Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015. 6
- [17] A. M. Lamb, A. Goyal, Y. Zhang, S. Zhang, A. Courville, and Y. Bengio. Professor forcing: A new algorithm for training recurrent networks. In *NIPS*, pages 4601–4609, 2016. 2
- [18] C.-Y. Lin. ROUGE: A package for automatic evaluation of summaries. In *ACL Workshop*, 2004. 1, 2, 7
- [19] S. Liu, Z. Zhu, N. Ye, S. Guadarrama, and K. Murphy. Improved image captioning via policy gradient optimization of SPIDER. In *ICCV*, 2017. 2
- [20] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu. BLEU: a method for automatic evaluation of machine translation. In *ACL*, pages 311–318, 2002. 1, 2, 7
- [21] R. Pasunuru and M. Bansal. Reinforced video captioning with entailment rewards. In *EMNLP*, pages 990–996, 2017. 1, 2, 6, 7
- [22] J. Pennington, R. Socher, and C. D. Manning. GloVe: Global vectors for word representation. In *EMNLP*, pages 1532–1543, 2014. 6
- [23] V. Ramanishka, A. Das, D. H. Park, S. Venugopalan, L. A. Hendricks, M. Rohrbach, and K. Saenko. Multimodal video description. In *ACM MM*, pages 1092–1096, 2016. 7
- [24] M. Ranzato, S. Chopra, M. Auli, and W. Zaremba. Sequence level training with recurrent neural networks. In *ICLR*, 2016. 1, 2, 3, 4, 6, 7, 10
- [25] S. J. Rennie, E. Marcheret, Y. Mroueh, J. Ross, and V. Goel. Self-critical sequence training for image captioning. In *CVPR*, 2017. 2, 3, 4, 5, 6, 7, 10, 11
- [26] Z. Shen, J. Li, Z. Su, M. Li, Y. Chen, Y.-G. Jiang, and X. Xue. Weakly supervised dense video captioning. In *CVPR*, 2017. 7
- [27] R. Shetty and J. Laaksonen. Frame- and segment-level features and candidate pool evaluation for video caption generation. In *ACM MM*, pages 1073–1076, 2016. 7
- [28] I. Sutskever, O. Vinyals, and Q. V. Le. Sequence to sequence learning with neural networks. In *NIPS*, pages 3104–3112, 2014. 3, 6
- [29] R. S. Sutton and A. G. Barto. *Reinforcement learning: An introduction*. MIT Press, 1998. 1
- [30] L. Theis, A. van den Oord, and M. Bethge. A note on the evaluation of generative models. *arXiv preprint arXiv:1511.01844*, 2015. 1, 2
- [31] D. Tran, L. Bourdev, R. Fergus, L. Torresani, and M. Paluri. Learning spatiotemporal features with 3D convolutional networks. In *ICCV*, pages 4489–4497, 2015. 6
- [32] R. Vedantam, C. Lawrence Zitnick, and D. Parikh. CIDEr: Consensus-based image description evaluation. In *CVPR*, pages 4566–4575, 2015. 1, 2, 3, 7
- [33] R. J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine Learning*, 8(3-4):229–256, 1992. 1, 2
- [34] J. Xu, T. Mei, T. Yao, and Y. Rui. MSR-VTT: A large video description dataset for bridging video and language. In *CVPR*, pages 5288–5296, 2016. 2, 5, 6
- [35] W. Zaremba, I. Sutskever, and O. Vinyals. Recurrent neural network regularization. *arXiv preprint arXiv:1409.2329*, 2014. 6

Supplemental Material

A. Further Training Details of CST

We encountered the same problem reported in [25] when training our CST models. More specifically, the generated captions tends to end unexpectedly such as “a girl is applying makeup to her”, “a woman is talking to a”, and “a woman is talking about a”. We suspect the models get difficult in predicting the correct ending token, while those bigrams such as “to her”, “to a” and “about a” were common among the captions, thus influences the CIDEr rewards. As suggested by [25], we also include $\langle \text{EOS} \rangle$ at the end of each caption in the reference set when computing CIDEr rewards. This simple amendment substantially prevents the model from generating those incomplete sentences.

B. Derivation of Eq. (14)

In [24] the authors did not provide full derivations of Eq. (14) and Eq. (15). For the sake of completeness, in this section and the next section we provide full derivations for both equations.

Because o_t is softmax input (the logit), we can compute the probability distribution over the vocabulary at time step t as follows:

$$p_\theta(w_t|h_t) = \frac{\exp(o_t)}{\sum \exp(o_t)} \quad (20)$$

In Eq (4), to compute the XE loss, we also use $p_\theta(w_t|h_t)$ in the context of the probability of the ground-truth word w_t . Recall that $1(w_t)$ is the one-hot representation of w_t , the XE loss can be rewritten w.r.t. o_t as follows:

$$L^{XE}(\theta) = - \sum_{t=1}^T \log \frac{\exp(o_t^T 1(w_t))}{\sum \exp(o_t)}, \quad (21)$$

Note that $\frac{\partial L^{XE}(\theta)}{\partial o_t}$ depends only on the softmax input o_t . Therefore:

$$\begin{aligned} \frac{\partial L^{XE}(\theta)}{\partial o_t} &= - \frac{\partial}{\partial o_t} \left(\log \frac{\exp(o_t^T 1(w_t))}{\sum \exp(o_t)} \right) \\ &= \frac{\partial \log \sum \exp(o_t)}{\partial o_t} - \frac{\partial o_t^T 1(w_t)}{\partial o_t} \\ &= \frac{1}{\sum \exp(o_t)} \frac{\partial \sum \exp(o_t)}{\partial o_t} - 1(w_t) \\ &= \frac{\exp(o_t)}{\sum \exp(o_t)} - 1(w_t) \\ &= p_\theta(w_t|h_t) - 1(w_t) \end{aligned} \quad (22)$$

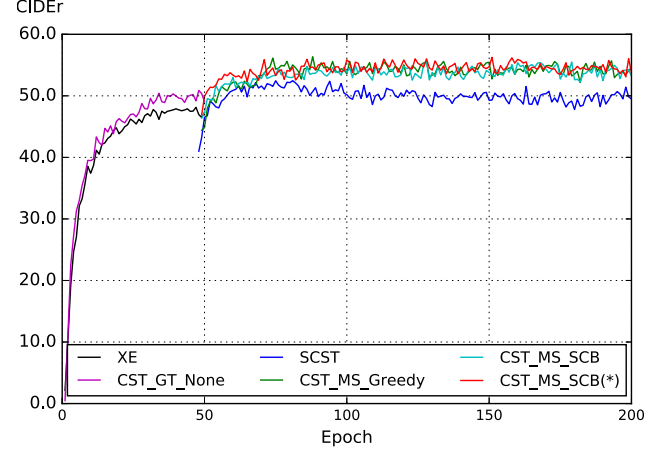


Figure 4: Performance on the validation set during training of different methods.

C. Derivation of Eq. (15)

From Eq (11), we can compute $\frac{\partial L^{RL}(\theta)}{\partial \theta}$ by a single sample sequence w^s :

$$\begin{aligned} \frac{\partial L^{RL}(\theta)}{\partial \theta} &= -(r(w^s) - b) \frac{\partial \log p_\theta(w^s)}{\partial \theta} \\ &= -(r(w^s) - b) \sum_{t=1}^T \frac{\partial \log p_\theta(w_t^s|h_t)}{\partial \theta} \end{aligned} \quad (23)$$

Here we also use $p_\theta(w_t^s|h_t)$ in the context of the probability of the sampled word w_t^s . More precisely, $\frac{\partial L^{RL}(\theta)}{\partial \theta}$ can be computed in term of o_t and the one-hot representation $1(w_t^s)$ as follows.

$$\frac{\partial L^{RL}(\theta)}{\partial \theta} = -(r(w^s) - b) \sum_{t=1}^T \frac{\partial}{\partial o_t} \left(\log \frac{\exp(o_t^T 1(w_t^s))}{\sum \exp(o_t)} \right) \quad (24)$$

Note that the reward function $r(\cdot)$ and the baseline b are treated as constants during back-propagation, and $\frac{\partial L^{RL}(\theta)}{\partial o_t}$ depends only on the logit o_t . Therefore:

$$\frac{\partial L^{RL}(\theta)}{\partial o_t} = -(r(w^s) - b) \frac{\partial}{\partial o_t} \left(\log \frac{\exp(o_t^T 1(w_t^s))}{\sum \exp(o_t)} \right) \quad (25)$$

Following the same steps in (22) to take the gradient of the second term, we come up with Eq. (15).

D. Performance on the Validation Set

We check the CIDEr performance on the validation set for every epoch and report the performance in Fig 4. Our

Table 5: Performance comparison of our WXE (CST_GT_None) model and standard XE model on different features.

Input feature(s)	Method	BLEU_1	BLEU_2	BLEU_3	BLEU_4	METEOR	ROUGE_L	CIDEr
Resnet	XE	75.6	61.0	48.1	37.2	26.7	58.0	39.9
	CST_GT_None	75.8	61.3	48.5	37.5	26.6	58.3	41.5
C3D	XE	77.2	62.4	48.9	37.2	26.7	58.5	39.8
	CST_GT_None	77.5	62.9	49.6	38.1	27.1	59.0	43.4
MFCC	XE	68.9	49.8	36.4	26.1	21.8	51.3	17.9
	CST_GT_None	69.5	51.1	38.1	27.7	22.0	51.8	19.8
Category	XE	70.2	52.1	38.8	28.4	22.7	53.0	25.5
	CST_GT_None	69.1	50.1	37.2	27.2	22.6	52.2	24.2
ResNet + C3D + MFCC + Category	XE	82.1	68.6	55.2	43.0	28.7	61.7	47.3
	CST_GT_None	82.8	69.5	56.2	44.1	29.1	62.4	49.7

WXE model (CST_GT_None) almost always reaches better check points than the standard XE model (in the first 50 epochs). When switching from XE training to RL training, all the RL models get difficult to converge in the first few iterations, but quickly recover and reach higher optimal levels. When start RL training from our WXE model, both the SCB and greedy baselines have similar convergence rates; however, our SCB baselines are two times faster to compute. Performance of SCST [25] model in our implementation, which was started training from the XE model, is clearly inferior to the RL models which was started from the pre-train WXE model.

E. WXE vs. XE Training

In Table 5 we compare the performance of our WXE (CST_GT_None) with the standard XE training on different video features. As we see, WXE significantly improves CIDEr metrics for all feature except Category feature. For example, the performance gain on Resnet, C3D and MFCC features are 1.6, 3.6 and 1.9 CIDEr respectively. We also observe that optimizing for CIDEr metric will also improve other metrics in most of the cases. Performance of Category feature is not reliable to assess the effectiveness of our WXE method because there are only 20 different categories, and many videos have the same category label. Therefore, this feature is not discriminative and should not be used alone to evaluate the captioning model. In summary, our proposed WXE method performs better than the standard XE method. Especially, once the rewards such as CIDEr scores are pre-computed, the WXE loss is as fast to compute as the standard XE loss.