

TeaForN: Teacher-Forcing with N-grams

Sebastian Goodman

Google Research
Venice, CA 90291
seabass@google.com

Nan Ding

Google Research
Venice, CA 90291
dingnan@google.com

Radu Soricut

Google Research
Venice, CA 90291
rsoricut@google.com

Abstract

Sequence generation models trained with teacher-forcing suffer from issues related to exposure bias and lack of differentiability across timesteps. Our proposed method, Teacher-Forcing with N-grams (TeaForN), addresses both these problems directly, through the use of a stack of N decoders trained to decode along a secondary time axis that allows model-parameter updates based on N prediction steps. TeaForN can be used with a wide class of decoder architectures and requires minimal modifications from a standard teacher-forcing setup. Empirically, we show that TeaForN boosts generation quality on one Machine Translation benchmark, WMT 2014 English-French, and two News Summarization benchmarks, CNN/Dailymail and Gigaword.

1 Introduction

Many state-of-the-art sequence generation models are trained using a technique called teacher-forcing (Goodfellow et al., 2016). Teacher-forcing is popular because it improves sample efficiency and provides training stability, but models trained with teacher-forcing are known to suffer from issues such as exposure bias (Venkatraman et al., 2015; Bengio et al., 2015; Ding and Soricut, 2017) and a lack of differentiability across timesteps (i.e., training updates made when decoding at time-step t cannot fully propagate to time-step $t - 1$). Previous attempts to address these issues include scheduled sampling (Bengio et al., 2015), parallel N-gram prediction (Yan et al., 2020), and sampling from previous predictions (Zhang et al., 2019).

Our proposed method, Teacher-Forcing with N-grams (TeaForN), imposes few requirements on the decoder architecture and does not require curriculum learning or sampling model outputs. TeaForN fully embraces the teacher-forcing paradigm and

extends it to N-grams, thereby addressing the problem at the level of teacher-forcing itself.

The advent of large-scale pretraining has pushed the state-of-the-art on Natural Language benchmarks to impressive heights, often showing gains across many tasks at once (Devlin et al., 2019; Raffel et al., 2019; Zhang et al., 2019). A negative consequence of this is the tendency towards large, data-hungry models, which have a negative impact on energy-consumption and accessibility (Strubell et al., 2019), as well as higher latency and production costs. As such, it is of increasing importance to develop techniques that counteract these tendencies. While TeaForN does increase training cost moderately, it can be used to drive down latency and inference cost, which dominate the overall cost of a production model.

Many sequence generation models use beam search to improve generation quality (Vaswani et al., 2017; Raffel et al., 2019; Zhang et al., 2019; Yan et al., 2020). In contrast with greedy decoding, beam search keeps the k most-likely candidates at each decoding timestep. While beam search has proven to be a reliable technique for improving output quality, previous work has shown that beam search actually degrades performance for sufficiently large k (Koehn and Knowles, 2017). In addition, the inference cost of a model increases linearly with k , due to the need for multiple decodings. We conduct an analysis of the effect of beam size on models trained both with and without TeaForN. We show that models trained with TeaForN require a smaller beam size to reach similar performance, a property that can achieve significant cost-savings.

Our experiments show that TeaForN can boost performance on both Machine Translation and News Summarization tasks, provided there is sufficient model capacity. With TeaForN, Transformer_{big} (Vaswani et al., 2017) improves by +.5 SacreBLEU (Post, 2018) on the WMT14 En-Fr

benchmark with beam search and +.3 without. When using TeaForN for summarization, PEGASUS_{large} (Zhang et al., 2019) improves by +.3 ROUGE-L on the Gigaword benchmark (Rush et al., 2015) and by +.2 on the CNN/Dailymail benchmark (Hermann et al., 2015). Further, PEGASUS_{large} trained with TeaForN matches the prior ROUGE-L scores on these benchmarks *without beam search*, representing an 8x reduction in decoder inference cost.

2 Related Work

One of the standard approaches to sequence-learning training is Maximum-likelihood Estimation (MLE). Although widely used in large array of applications, MLE estimation for sequence learning suffers from the exposure-bias problem (Venkatraman et al., 2015; Ranzato et al., 2015). Exposure-bias produces brittle models due to training procedures during which the models are only exposed to their training data distribution but not to their own predictions. Possible solutions to the exposure-bias problem in neural-network settings have used “data as demonstrator” (Venkatraman et al., 2015) and “scheduled sampling” (Bengio et al., 2015) approaches. Although improving model performance in practice, such proposals have been shown to be statistically inconsistent (Huszar, 2015), and still need to perform MLE-based warm-start training, rendering such solutions unsatisfactory. Along similar lines, the “professor forcing” (Lamb et al., 2016) method uses adversarial domain adaptation to encourage network dynamics to be the same during training and inference, though it requires sampling sequences during training.

A different approach, based on reinforcement learning methods, achieves sequence learning following a policy-gradient (PG) method (Sutton et al., 1999). It directly attacks the exposure-bias problem by having the training models exposed exclusively to their own predictions while scoring them using reward functions. However, this approach introduces another issue, related to the large discrepancy between the model prediction distribution and the reward function’s values, which is especially acute during the early training stages when the predicted outputs are all equally bad. As a result, this method also requires a warm-start phase in which the model distribution achieves some local maximum with respect to a reward-free objective (e.g., MLE), followed by a model refinement phase in

which reward-based PG updates are used to refine the model (Ranzato et al., 2015; Wu et al., 2016; Liu et al., 2017). Although such combinations achieve better results in practice compared to pure likelihood-based approaches, they are unsatisfactory from a theoretical and modeling perspective, as well as inefficient from a speed-to-convergence perspective. A pure PG formulation that side-steps these issues is (Ding and Soricut, 2017), which allows for both cold-start training as well as more efficient convergence properties.

The PG-based approaches have an inherent complexity that stems from the use of quirky reward functions such as ROUGE (Lin, 2004) or CIDEr (Vedantam et al., 2015), which forfeits the advantage of sample efficiency as they often cannot be efficiently computed using current accelerators like TPUs (You et al., 2019). MLE-based approaches appear to be favored due to efficiency properties, and the search for training methods that produce less brittle models is still on-going.

Another closely related idea is End-to-End Backprop (E2E) (Ranzato et al., 2015), which has a similar goal of naturally approximating sequence level training by propagating smooth model predictions instead of groundtruth inputs. TeaForN differs from E2E in several key ways. First, TeaForN learns jointly from both groundtruth and model predictions as inputs throughout the entire training duration, whereas E2E requires a training schedule to transition from groundtruths to model predictions. Second, TeaForN supports methods other than k-max for computing smooth model predictions, two of which we explore as a part of our work. Third, we introduce the notion of a discount factor, which weights the importance of immediate predictions higher than that of future predictions.

Another such work is (Yan et al., 2020), which proposes a modified Transformer for parallel N-gram prediction. While their work does address the issue of strong local correlations caused by teacher-forcing, it does not address exposure bias, as it always trains on groundtruth inputs.

Also related are models such as the one proposed by (Strubell et al., 2017), which uses a stack of dilated convolutions to iteratively refine model predictions. Though architecturally similar, TeaForN only uses the stack at training time and solves for a fundamentally different problem.

Our TeaForN method maintains the efficiency advantages of MLE-based approaches, while ad-

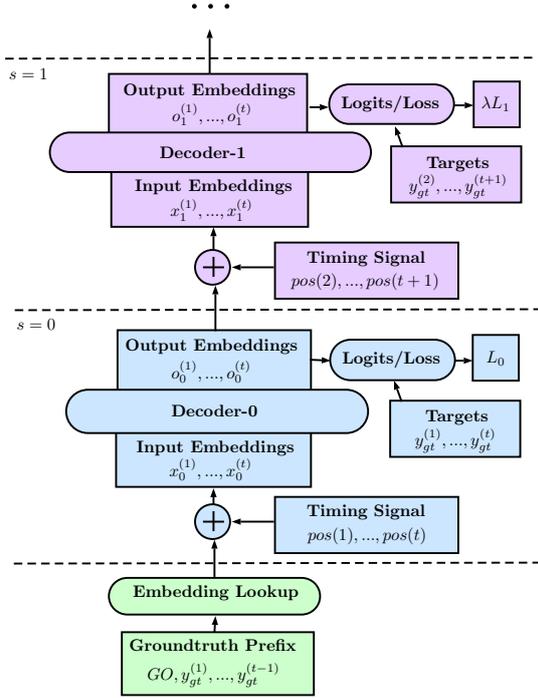


Figure 1: An illustration of TeaForN training, wherein each decoder after the first uses the outputs of the previous decoder as inputs. Decoder weights may be shared across layers in order to address exposure bias.

addressing both exposure bias and the issue of differentiability across timesteps. In addition, it is general enough to be used on a wide class of autoregressive decoders, including RNN (Hochreiter and Schmidhuber, 1997; Chung et al., 2014) and Transformer (Vaswani et al., 2017) decoders, though our experiments focus on the Transformer.

3 TeaForN

Autoregressive sequence decoders are trained to minimize the negative log likelihood of the groundtruth tokens $y_{gt}^{(t)}$. During training, previous *groundtruth tokens* are used as decoder inputs for predicting the next token. If we define the embedding matrix to be E of size $V \times D$, where V is the vocabulary size and D is the embedding size, and the embedding of the groundtruth token as $x^{(t)} = E[y_{gt}^{(t-1)}, :] := e_{gt}^{(t-1)}$ of size D , then the standard teacher-forcing loss is equal to,

$$\begin{aligned} L &= - \sum_{t=1}^T \log(P(y_{gt}^{(t)} | y_{gt}^{(0)}, \dots, y_{gt}^{(t-1)})) \\ &= - \sum_{t=1}^T \log(P(y_{gt}^{(t)} | x^{(1)}, \dots, x^{(t)})) \end{aligned}$$

where we define $y_{gt}^{(0)} = GO$ as the starting token.

The class probability distribution P is typically modeled as softmax-normalized logits, which are a linear projection of decoder output $o^{(t)}$ of size D onto the class embeddings using output projection matrix W of size $V \times D$:

$$P(y^{(t)} | x^{(1)}, \dots, x^{(t)}) = \text{softmax}(W o^{(t)}).$$

To reduce the model parameter size, it is standard to share the parameters of the output projection matrix and the embedding matrix, such that $E = W$.

During inference, groundtruth tokens become unavailable. Therefore, previous tokens from the *model predictions* are used as decoder inputs for decoding the next token. The discrepancy between training time and inference time input distributions causes models to suffer from exposure bias, meaning that they do not learn to correct for past decoding errors (Bengio et al., 2015).

TeaForN addresses exposure bias by learning jointly how to predict from both groundtruth and past model predictions as inputs. TeaForN setups consist of a stack of N decoders, as illustrated in Figure 1. At position t , the first decoder (Decoder-0) takes input from the embedding of the previous groundtruth token $x_0^{(t)} = e_{gt}^{(t-1)}$ and learns to predict the target token $y_{gt}^{(t)}$, same as in teacher-forcing. The next decoder (Decoder-1) takes input from the output $x_1^{(t)} = o_0^{(t)}$ of the first decoder and learns to predict the next target token $y_{gt}^{(t+1)}$.

More formally, let us use subscript $s \in [0, N)$ to denote the offset within the decoder stack. We define the input to decoder s at time t as:

$$x_s^{(t)} = \text{pos}(t + s) + \begin{cases} e_{gt}^{(t-1)} & s = 0 \\ o_{s-1}^{(t)} & s > 0 \end{cases}$$

where $\text{pos}(t + s)$ is a timing signal that is added to the inputs for models such as the Transformer (Vaswani et al., 2017). This term may be omitted for models that do not expect it.

The training loss of Decoder- s at time t is the negative log likelihood of the $(t + s)^{th}$ element in the groundtruth sequence:

$$L_s = - \sum_{t=1}^T \log(P(y_{gt}^{(t+s)} | x_s^{(1)}, \dots, x_s^{(t)}; \theta_s))$$

and the total TeaForN training loss is the sum of decoder losses

$$L = \sum_{s=0}^{N-1} \lambda^{s-1} L_s$$

where $\lambda \in (0, 1]$ is a discount factor needed to weigh the risk of harming next-word accuracy against the benefits of TeaForN. During inference, TeaForN uses only the first decoder (Decoder-0) in the stack; the rest are discarded.

The intuition behind TeaForN is as follows. Under standard teacher-forcing, the decoder output $o^{(t)}$ only learns to predict the groundtruth label $y_{gt}^{(t)}$, while outputs that favor other classes are considered equally bad, and will be penalized by the loss. This is not reasonable because classes carrying similar meanings to the groundtruth label do not change the meaning of the sequence significantly, and may still lead to the correct prediction for the next label. Under TeaForN, the decoder output $o^{(t)}$ is also used as the input of a secondary decoder for decoding the next position. Therefore, all outputs that result in predicting the next groundtruth label $y_{gt}^{(t+1)}$ will have lower loss and therefore be differentiated from other outputs.

In our experiments, we allow the decoder parameters to be either shared ($\theta_0 = \theta_s, \forall s$) or unshared. In a shared-weight configuration, the model learns to predict the next groundtruth label from the class that the same model predicted in the previous position. This is similar to the inference time condition, so we expect shared-weight TeaForN to address exposure-bias better than unshared-weight TeaForN. Shared-weight configurations also have performance advantages such as lower memory consumption and faster training.

Since TeaForN solves for a more difficult problem than teacher-forcing, we expect it to work better for models with higher capacity. We later show evidence of this by comparing results for two model sizes on Machine Translation.

It is straightforward to show that TeaFor1 (N=1) and teacher-forcing are equivalent, as the inputs to the first TeaForN decoder are groundtruth sequence embeddings and $\lambda^0 = 1$. Thus, TeaForN is a natural extension of teacher-forcing to N-grams.

3.1 Embedded Top-k Stacked Decoder Input

Previously, our TeaForN model directly used the decoder output of the $(s - 1)$ -th stack as the input of the decoder of the s -th stack:

$$x_s^{(t)} = o_{s-1}^{(t)}. \quad (1)$$

This is an approximation to the inference-time decoder input, which (for greedy decoding) is

$$x_s^{(t)} = E[\text{argmax}(W o_{s-1}^{(t)}, :)], \quad (2)$$

where $\text{argmax}(x)$ returns the index of the V -dim vector with the maximum value.

Inspired by the End-to-End Backprop (E2E) (Ranzato et al., 2015), we also consider the following alternative decoder input,

$$x_s^{(t)} = E^\top \text{softmax}(\text{top}_k(W o_{s-1}^{(t)})), \quad (3)$$

where top_k is a function which keeps the top-k values of the vector, and masks out the others.

It is easy to verify, when $k = 1$, Eq. (3) reduces to Eq. (2); when $k = V$,

$$x_s^{(t)} = E^\top \text{softmax}(W o_{s-1}^{(t)}).$$

Compared to Eq. (1), Eq. (3) is more computationally expensive, as it involves additional embedding matrix multiplications and/or a top-k sorting. Furthermore, we would like to emphasize a critical difference between the TeaForN and E2E (Ranzato et al., 2015). In TeaForN, the 0-th stack of every position is always clamped to the groundtruth input, while for E2E the groundtruth is completely thrown away after warm-up training. The groundtruth clamping allows the TeaForN to avoid the warm-up training which is necessary for E2E.

4 Experimental Results

Our empirical study of TeaForN is comprised of two sections. First, we present experiments on Machine Translation using the well-known Transformer model (Vaswani et al., 2017). Second, we show results for News Summarization, for which we use PEGASUS (Zhang et al., 2019), a state-of-the-art pretrained text summarization model.

We perform minimal hyperparameter tuning over the course of these experiments. This can be partly credited to the underlying models being well-tuned already, but also to TeaForN, which works out-of-the-box without much hyperparameter tuning. One exception is the tuning of the number of training steps, as we found that the number of steps used by previous settings is sometimes insufficient.

4.1 Machine Translation

In this section, we study the effects of applying TeaForN to a well-known Transformer-based Machine Translation model. We present results for two size variants of the model, Transformer_{base} and Transformer_{big} (Vaswani et al., 2017). The differences are summarized in Table 2.

We use the same WMT14 language-pair benchmarks originally reported in the Transformer paper:

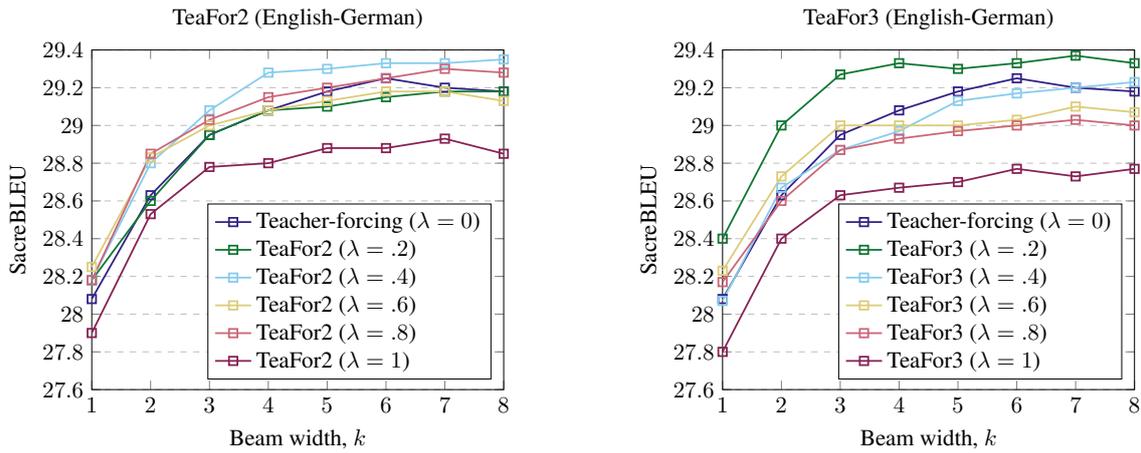


Figure 3: Beam width vs. SacreBLEU on the WMT 2014 English-German benchmark, for discount factors $\lambda \in \{0, .2, .4, .6, .8, 1\}$. Left and right plots show TeaFor2 and TeaFor3, respectively. Reported scores are averages over $n=3$ independent training runs, with error bars omitted for readability. See Appendix Tables 11 and 13 for results with standard error measurements.

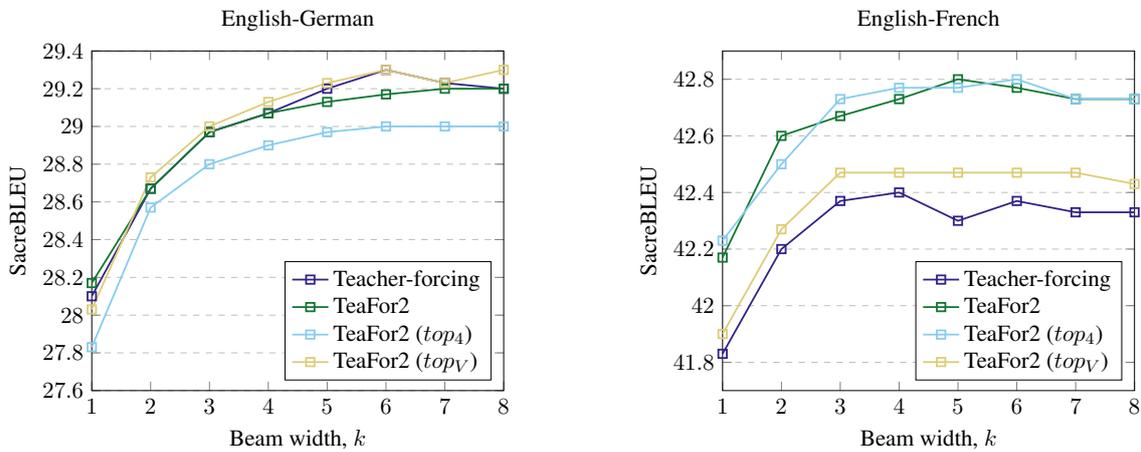


Figure 4: Beam width vs. SacreBLEU on WMT 2014 benchmarks comparing approximation methods Top-K. Reported scores are averages over $n=3$ independent training runs, with error bars omitted for readability. See Appendix Tables 15 and 16 for results with standard error measurements.

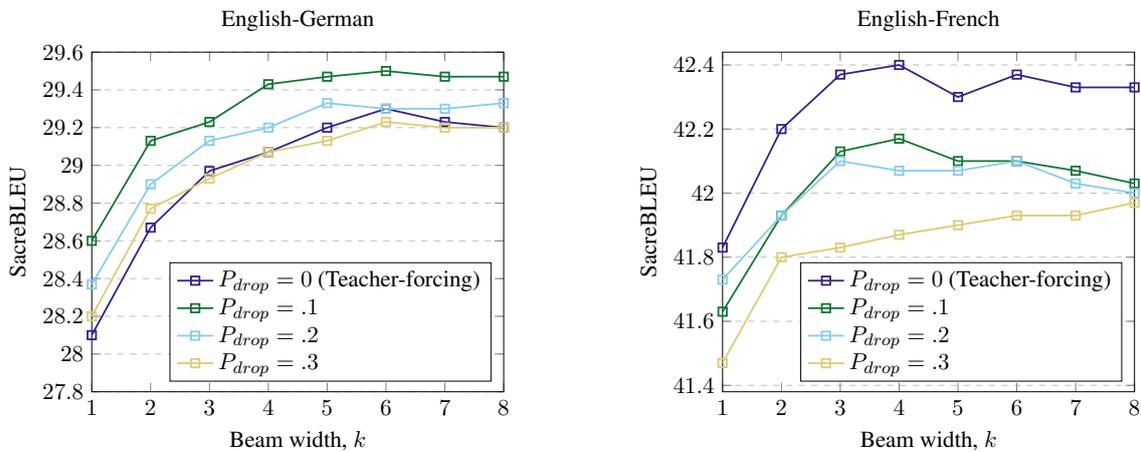


Figure 5: Beam width vs. SacreBLEU on WMT 2014 benchmarks using Word Drop Regularization. Reported scores are averages over $n=3$ independent training runs. See Appendix Tables 17 and 18 for results with standard error measurements.

A Appendix

- RS Sutton, D McAllester, S Singh, and Y Mansour. 1999. Policy gradient methods for reinforcement learning with function approximation. In *NIPS*.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Proceedings of NeurIPS*.
- Ramakrishna Vedantam, C. Lawrence Zitnick, and Devi Parikh. 2015. CIDEr: Consensus-based image description evaluation. In *Proceedings of CVPR*.
- Arun Venkatraman, Martial Hebert, and J. Andrew Bagnell. 2015. Improving multi-step prediction of learned time series models. In *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence*, pages 3024–3030. AAAI Press.
- Y. Wu, M. Schuster, and al. 2016. Google’s neural machine translation system: Bridging the gap between human and machine translation. *CoRR*, abs/1609.08144.
- Yu Yan, Weizhen Qi, Yeyun Gong, Dayiheng Liu, Nan Duan, Jiusheng Chen, Ruofei Zhang, and Ming Zhou. 2020. Prophetnet: Predicting future n-gram for sequence-to-sequence pre-training.
- Yang You, Zhao Zhang, Cho-Jui Hsieh, James Demmel, and Kurt Keutzer. 2019. Fast deep neural network training on distributed systems and cloud tpus. *IEEE Trans. Parallel Distrib. Syst.*, 30(11):2449–2462.
- Jingqing Zhang, Yao Zhao, Mohammad Saleh, and Peter J. Liu. 2019. Pegasus: Pre-training with extracted gap-sentences for abstractive summarization.
- Wen Zhang, Yang Feng, Fandong Meng, Di You, and Qun Liu. 2019. Bridging the Gap between Training and Inference for Neural Machine Translation. *arXiv e-prints*, page arXiv:1906.02448.

