# Rethinking Exposure Bias in Adversarial Language Modeling

**Yifan Xu** *, **Kening Zhang** *, **Haoyu Dong, Yuezhou Sun, Wenlong Zhao & Zhuowen Tu**
University of California San Diego
{yix081,kez040, had002, yus174, wez094,ztu}@ucsd.edu

## Abstract

Exposure bias refers to the phenomenon that a language model trained under the teacher forcing schema may perform poorly at the inference stage when its predictions are conditioned on its outputs that diverge from the training corpus. Although several adversarial training methods have been proposed to avoid teacher forcing, lacking a clear evaluation for the exposure bias remains a concern. The contribution of our work is two-fold. (1) We propose to evaluate exposure bias based on the quality of sentence generated in the sentence completion task. (2) We adopt two strategies, *multi-range reinforcing* and *multi-entropy sampling*, to stabilize adversarial training, and show an improvement over the competing models with regards to the sentence completion task and corpus BLEUs.

## 1 Introduction

Likelihood-based language models with deep neural networks have been widely adopted to tackle the language modeling tasks (Graves et al., 2013; Karpathy and Fei-Fei, 2015; Bahdanau et al., 2014). By far, one of the most popular training strategies is *teacher forcing*, which is derived from the general maximum likelihood estimation (MLE) principle (Williams and Zipser, 1989). Under the teacher forcing schema, the language model makes predictions conditioned on the ground-truth inputs. This is susceptible to so-called *exposure bias*: a model may perform poorly at the inference stage, once its prefix diverges from the previously learned data (Bengio et al., 2015). However, there is little work on how to expose and quantify such performance degeneration in text generation.

A common strategy to mitigate the exposure bias problem is to impose additional supervision upon the model's self-generated output via adversarial training. The actor-critic (AC) method (Konda and Tsitsiklis, 2000) and SeqGAN (Yu et al., 2017) introduce an additional critic network to offer rewards on a language model's self-generated sequences. Therefore, the language model can later, at the inference stage, predict robustly with its previous outputs. One issue in adversarial training is that the signal from the critic network is very sparse, which leads to stability issues. The second issue is about the non-stationary sampled data with strongly correlated online updates (Pfau and Vinyals, 2016; Mnih et al., 2016). Due to these problems, existing language GANs (Yu et al., 2017; Lin et al., 2017; Guo et al., 2017) have a risk of compromising generation diversity (Caccia et al., 2018). This paper makes the following contributions:

1. We propose to evaluate the exposure bias for a language model by performing the sentence completion task using the ground truth prefix.

2. We introduce a new approach, *multi-entropy sampling* and *multi-range reinforcing* (MEMR), to overcome the difficulties during adversarial training, which demonstrates a significant improvement over the competing models in the corpus BLEUs metrics, as well as our proposed measures in sentence completion.

## 2 Related Works

A common measure quantifying the exposure bias is still absent. Existing works often show performance gains by introducing adversarial training but questions remain if such gains indeed result in the reduction of the exposure bias (Bahdanau et al., 2016; Yu et al., 2017). Later works add generation diversity into consideration (Shi et al., 2018; Caccia et al., 2018; Alihosseini et al., 2019) or take a perspective from traditional language modeling aspects (Tevet et al., 2018). A closely related work to our evaluation measure is (He et al., 2019). The difference is that He et al. (2019) requires inference for ground truth data distribution with experiments performed using synthetic data.

---

* Equal contribution.

An early work addressing the exposure bias problem is (Bengio et al., 2015) in which a curriculum learning approach called *scheduled sampling* is proposed by gradually replacing the ground-truth tokens with the model's predictions. In recent RL-inspired works, Ranzato et al. (2015) adopt the REINFORCE algorithm (Sutton et al., 2000) to directly optimize the test-time evaluation score. Bahdanau et al. (2016) employ a similar approach by training a critic network to predict the metric score for the actor's generated sequence of tokens. In parallel, a language version of generative adversarial networks (GANs) (Goodfellow et al., 2014), SeqGAN, is introduced in (Yu et al., 2017). SeqGAN consists of a generator pre-trained under MLE and a discriminator pre-trained to discern the generator's distribution from the real data. Follow-up works such as RankGAN (Lin et al., 2017) and LeakGAN (Guo et al., 2017) alter the training objectives or model architectures to enhance the guidance. RankGAN (Lin et al., 2017) replaces the binary reward with a relative ranking score. Leak-GAN (Guo et al., 2017) allows the discriminator to "leak" its internal states to the generator at intermediate steps. Shi et al. (2018) model a reward function using inverse reinforcement learning (IRL).

# 3 Exposure Bias Evaluation

## 3.1 Exposure Bias

Cross-entropy loss adopted in teacher forcing is equivalent to minimizing the *forward* KL divergence $D_{KL}(P||Q_\theta)$ between data distribution $P$ and model distribution $Q_\theta$. However, during the inference stage, the model is often evaluated based on the quality of its generated samples. The evaluation metrics or human experts can be seen as surrogates of the data distribution $P$, so what they measure is the *reverse* KL divergence $D_{KL}(Q_\theta||P)$.

In Bayesian inference, there is a well-known difference between $D_{KL}(P||Q)$ and $D_{KL}(Q||P)$ (MacKay, 2003). Minimizing $D_{KL}(P||Q)$ encourages the model to cover all the modes in the training data, which will result in over-generalization in the extreme case. In contrast, minimizing $D_{KL}(Q||P)$ prefers the model to concentrate on the largest mode while ignoring the others, which tends to cause mode collapse (Huszár, 2015). In our language modeling task, an LSTM strives to cover the entire data distribution at the cost of over-generalization. It is more likely to produce prefixes different from those seen at the training stage, and the fact that this model has never learned to predict based on these prefixes potentially leads to the exposure bias .

## 3.2 Sentence Completion Task

In this section, we form a sentence completion task to evaluate the exposure bias. Given a sentence prefix $X_{1:k}$ of length K drawn from a data distribution $P$, we apply a language model $Q_\theta$ to perform sentence completion until final the $T$ step, starting from such prefix.

- If the prefix $X_{1:k}$ is sampled from a seen distribution $P_{seen}$, then the exposure bias for the sentence completion task should be relatively low, where

$$Q_\theta(X_{k:T}|P_{seen}) = \mathbb{E}_{X_{1:k} \sim P_{seen}} Q_\theta(X_{k:T}|X_{1:k})$$

- If the prefix $X_{1:k}$ comes from an unseen data distribution $P_{unseen}$, then the exposure bias for the task can be critical, where

$$Q_\theta(X_{k:T}|P_{unseen}) = \mathbb{E}_{X_{1:k} \sim P_{unseen}} Q_\theta(X_{k:T}|X_{1:k})$$

Based on the definition for the exposure bias, $Q_\theta(X_{k:T}|P_{unseen})$ should suffer more from the training-testing deviation than $Q_\theta(X_{k:T}|P_{seen})$. Also, such performance degeneration should be more significant when prefix $k$ grows longer in both scenarios. These two hypotheses are confirmed by our result in Figure 1.

As a measurement to assess model's generation quality, forward Corpus BLEU, *BLEU*$_\text{F}$, is evaluated. Because precision is the primary concern, we set softmax temperature $\tau = 0.5$ to sample high-confidence sentences from model's distribution.

Based on the task completion task results in Figure 1, we observe that original SeqGAN (Yu et al., 2017) shows more stable result although many text GAN variants are proposed later, which is unexpected. Therefore, our method MEMR is motivated to improve SeqGAN by introducing denser reward signal from the critic network and further stabilizing the adversarial training.

# 4 Method Description

## 4.1 Actor-Critic Training

Actor-Critic methods (ACs) formulates language modeling as a generalized Markov Decision Process (MDP) problem, where the actor learns to optimize its policy guided by the critic, while the critic learns to optimize its value function based on the actor's output and external reward information. As Pfau and Vinyals (2016) points out, GAN methods can be seen as a special case of AC where the critic aims to distinguish the actor's generation from real data and the actor is optimized in an opposite direction to the critic.

In this work, we use a standard single-layer LSTM as the actor network. The training objective is to maximize the model's expected end rewards with policy gradient (Sutton et al., 2000):

$$\mathcal{L}(\theta) = -\mathbb{E}_{X_{1:T} \sim \pi_\theta} \sum_{t=1}^{T} Q_\phi(x_t, h_t) \log \pi_\theta(x_t|h_t)$$

In practice, we perform a Monte-Carlo (MC) search with roll-out policy following Yu et al. (2017) to sample complete sentences starting from each location in a predicted sequence and compute their end rewards. Empirically, we found out that the maximum, instead of average, of rewards in the MC search better represents each token's actor value and yields better results during training. Therefore, we compute the action value by:

$$Q_\phi(x_t, h_t) = \max_{X_{t:T} \in MC^\theta(X_{1:t}, T)} Q_\phi(X_{1:T})$$

Then, We use a convolutional neural network (CNN) as the critic to predict the expected rewards for current generated prefix:

$$\mathcal{L}(\phi) = -\mathbb{E}_{X_{1:T} \sim \pi_\theta} (r(X_{1:T}) - Q_\phi(X_{1:T}))^2$$

## 4.2 MEMR

During the experiment, we observe a certain level of instability for the learned models. In the previous literature, two major factors behind the training instability are the sparse reward from critic network and the update correlation in the sampling process (Pfau and Vinyals, 2016; Mnih et al., 2016; Volodymyr et al., 2013). We address these problems using the following strategies:

**Multi-Entropy Sampling:** Language GANs can be seen as online RL methods, where the language model is updated from data generated by a single policy. Most sampled sentences in MC search are highly correlated. Similar to Xu et al. (2019), we empirically observe that increasing the range of the entropy of the actor's sample distribution during training is beneficial to the adversarial training performance. Specifically, we alternate the temperature $\tau$ in the softmax to generate samples under different behavior policies. During the critic's training, the ground-truth sequences are assigned a perfect target value of 1. The samples obtained with $\tau < 1$ are supposed to contain lower entropy, thus they receive a higher target value close to 1. Those samples obtained with $\tau > 1$ contain higher entropy, and the target value is closer to 0. This mechanism decorrelates updates during sequential sampling by sampling from multiple diverse entropy distributions synchronously.

**Multi-Range Reinforcing:** Our idea of multi-range supervision takes inspiration from deeply-supervised nets (DSNs) (Lee et al., 2015). By design, lower layers in a CNN have smaller receptive fields, allowing them to make better use of local patterns. Differently from DSNs (Lee et al., 2015) which disregard all intermediate predictions in the end, we average the reward predictions from multiple intermediate layers of the critic network with the final output, which attend to local n-grams rather than the whole complete sentence. This is a solution to the reward sparseness, as the language model can receive averaged reward with more local information.

## 4.3 Effectiveness of Multi-Range Reinforcing and Multi-Entropy Sampling

Table 1 demonstrates the effectiveness of multi-entropy sampling (ME) and multi-range reinforcing (MR). We observe that ME improves BLEU$_{F5}$ (precision) significantly while MEMR further enhances BLEU$_{F5}$ (precision) and BLEU$_{F5}$ (recall). Detailed explanations of these metrics can be found in Section 5.2.

| Architecture | BLEU$_{F5}$ | BLEU$_{B5}$ |
|---|---|---|
| TF | $15.4 \pm 0.17$ | $30.5 \pm 0.08$ |
| AC | $13.8 \pm 0.16$ | $30.3 \pm 0.13$ |
| AC (with ME) | $22.4 \pm 0.25$ | $30.0 \pm 0.09$ |
| AC (with MEMR ) | $24.5 \pm 0.14$ | $31.6 \pm 0.10$ |

Table 1: Effectiveness of the proposed ME and MEMR strategies on EMNLP2017 WMT News Dataset

## 5 Experiment

### 5.1 Datasets

We perform evaluations on two datasets: *EMNLP2017 WMT News* [1] and *Google-small*, a subset of Google One Billion Words [2].

- *EMNLP2017 WMT News* is provided in (Zhu et al., 2018), a benchmarking platform for text GANs. The entire dataset is split into a training set of 195,010 sentences, a validation set of 83,576 sentences, and a test set of 10,000 sentences. The vocabulary size is 5,254 and the average sentence length is 27.

- *Google-small* is sampled and pre-processed from the Google One Billion Words. It contains a training set of 699,967 sentences, a validation set of 200,000 sentences, and a test set of 99,985 sentences. The vocabulary size is 61,458 and the average sentence length is 29.

### 5.2 BLEU metric

We adopt three variations of BLEU metric from Shi et al. (2018). *BLEU$_F$*, or forward BLEU, is a metric for precision, and *BLEU$_B$*, or backward BLEU, is a metric for recall. *BLEU$_{HA}$* computes the harmonic mean of both BLEU. These three metrics take both

---

[1] https://github.com/geek-ai/Texygen
[2] http://www.statmt.org/lm-benchmark/

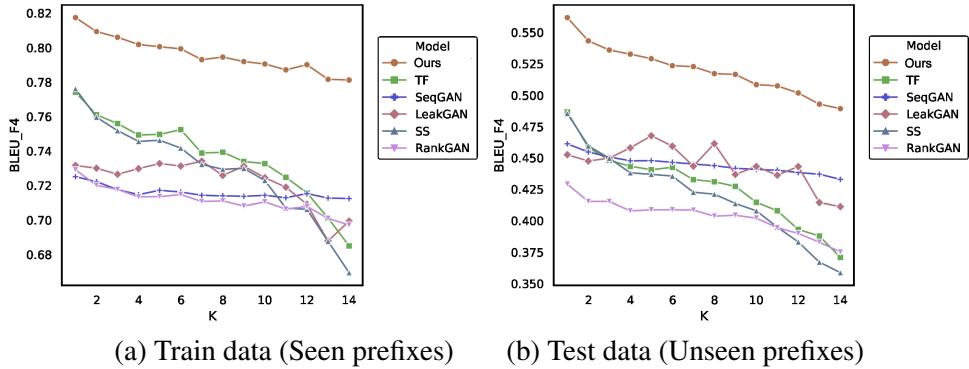(a) Train data (Seen prefixes)　　(b) Test data (Unseen prefixes)

Figure 1: **Sentence Completion Task** results based on prefixes from training and testing datasets on EMNLP2017 WMT News [Higher is better]. In each experiment, the data source for the prefixes is used as the reference to calculate $BLEU_{F4}$.

| Model | EMNLP2017 WMT | | | Google-small | | |
| --- | --- | --- | --- | --- | --- | --- |
| | $BLEU_{F5}$ | $BLEU_{B5}$ | $BLEU_{HA5}$ | $BLEU_{F5}$ | $BLEU_{B5}$ | $BLEU_{HA5}$ |
| TEACHER FORCING (TF) | $15.4 \pm 0.11$ | $30.5 \pm 0.05$ | $20.5 \pm 0.10$ | $9.6 \pm 0.03$ | $12.9 \pm 0.02$ | $11.00 \pm 0.02$ |
| SCHEDULED SAMPLING (SS) (Bengio et al., 2015) | $12.1 \pm 0.14$ | $30.3 \pm 0.06$ | $17.3 \pm 0.14$ | $6.2 \pm 0.04$ | $10.7 \pm 0.02$ | $7.8 \pm 0.04$ |
| SEQGAN (Yu et al., 2017) | $16.6 \pm 0.09$ | $28.7 \pm 0.37$ | $21.0 \pm 0.11$ | $20.7 \pm 0.02$ | $14.4 \pm 0.02$ | $17.0 \pm 0.01$ |
| RANKGAN (Lin et al., 2017) | $17.7 \pm 0.14$ | $30.1 \pm 0.06$ | $22.3 \pm 0.11$ | $21.4 \pm 0.06$ | $12.7 \pm 0.03$ | $15.9 \pm 0.02$ |
| LEAKGAN (Guo et al., 2017) | $19.8 \pm 0.11$ | $31.6 \pm 0.04$ | $24.4 \pm 0.10$ | - | - | - |
| MEMR (ours) | $\textbf{24.5} \pm 0.08$ | $\textbf{31.6} \pm 0.06$ | $\textbf{27.9} \pm 0.07$ | $\textbf{22.0} \pm 0.07$ | $\textbf{15.8} \pm 0.02$ | $\textbf{18.4} \pm 0.03$ |

Table 2: Corpus BLEUs Results on EMNLP2017 WMT News and the Google-small dataset. The 95 % confidence intervals from multiple trials are reported. [†] the Google-small was not tested in (Guo et al., 2017) and we are unable to train LeakGAN on this dataset using the official code due to its training complexity (taking 10+ hours per epoch).

diversity and quality into consideration. A model with severe mode collapse or diverse but incorrect outputs receives low scores.

## 5.3　Implementation Details

We implement a standard single-layer LSTM as the generator (actor) and a eight-layer CNN as the discriminator (critic). The LSTM has embedding dimension 32 and hidden dimension 256. The CNN consists of 8 layers with filter size 3, where the 3rd, 5th, and 8th layers are directly connected to the output layer for multi-range supervision. Other parameters are consistent with Zhu et al. (2018). Adam optimizer is deployed for both critic and actor with learning rate $10^{-4}$ and $5 \cdot 10^{-3}$ respectively. The target values for the critic network are set to [0, 0.2, 0.4, 0.6, 0.8] for samples generated by the LSTM with softmax temperatures [0.5, 0.75, 1.0, 1.25, 1.5].

## 6　Results

Based on the sentence completion results in Figure 1, all models decrease in precision of generated text (reflected via $BLEU_{F4}$) as the fed-in prefix length ($K$) increases, but the effect is stronger on the unseen test data, revealing the existence of exposure bias. Nonetheless, our model trained under ME and MR yields the best sentence quality and a relatively moderate performance decline.

Although TF and SS demonstrate higher $BLEU_{F5}$ performance with shorter prefixes, their

sentence qualities drop drastically on the test dataset with longer prefixes. On the other hand, GANs begin with lower $BLEU_{F4}$ precision scores but demonstrate less performance decay as the prefix grows longer and gradually outperform TF. This robustness against unseen prefixes exhibits that supervision from a learned critic can boost a model's stability in completing unseen sequences. The better generative quality in TF and the stronger robustness against exposure bias in GANs are two different objectives in language modeling, but they can be pursued at the same time. Our model's improvement in both perspectives exhibit one possibility to achieve the goal.

We also report Corpus BLEUs to reflect the quality and diversity of generated text in Table 2 with competing models on EMNLP2017 WMT News and Google-small. Our model, MEMR, outperforms the others in Corpus BLEUs, indicating a high diversity and quality in its sample distribution.

## 7　Conclusion

We propose to use the sentence completion task to reveal exposure bias in text generation. Further, we overcome the hurdles in adversarial training with *multi-range reinforcing* and *multi-entropy sampling* (MEMR), which shows an improvement in the sentence completion task and Corpus BLEUs.

## References

Danial Alihosseini, Ehsan Montahaei, and Mahdieh Soleymani Baghshah. 2019. Jointly measuring diversity and quality in text generation models. In *Proceedings of the Workshop on Methods for Optimizing and Evaluating Neural Language Generation*, pages 90–98.

Dzmitry Bahdanau, Philemon Brakel, Kelvin Xu, Anirudh Goyal, Ryan Lowe, Joelle Pineau, Aaron Courville, and Yoshua Bengio. 2016. An actor-critic algorithm for sequence prediction. *arXiv preprint arXiv:1607.07086*.

Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2014. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*.

Samy Bengio, Oriol Vinyals, Navdeep Jaitly, and Noam Shazeer. 2015. Scheduled sampling for sequence prediction with recurrent neural networks. In *Advances in Neural Information Processing Systems*, pages 1171–1179.

Massimo Caccia, Lucas Caccia, William Fedus, Hugo Larochelle, Joelle Pineau, and Laurent Charlin. 2018. Language gans falling short. *arXiv preprint arXiv:1811.02549*.

Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2014. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680.

Alex Graves, Abdel-rahman Mohamed, and Geoffrey Hinton. 2013. Speech recognition with deep recurrent neural networks. In *Acoustics, speech and signal processing (icassp), 2013 ieee international conference on*, pages 6645–6649. IEEE.

Jiaxian Guo, Sidi Lu, Han Cai, Weinan Zhang, Yong Yu, and Jun Wang. 2017. Long text generation via adversarial training with leaked information. *arXiv preprint arXiv:1709.08624*.

Tianxing He, Jingzhao Zhang, Zhiming Zhou, and James Glass. 2019. Quantifying exposure bias for neural language generation. *arXiv preprint arXiv:1905.10617*.

Ferenc Huszár. 2015. How (not) to train your generative model: Scheduled sampling, likelihood, adversary? *arXiv preprint arXiv:1511.05101*.

Andrej Karpathy and Li Fei-Fei. 2015. Deep visual-semantic alignments for generating image descriptions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3128–3137.

Vijay R Konda and John N Tsitsiklis. 2000. Actor-critic algorithms. In *Advances in neural information processing systems*, pages 1008–1014.

Chen-Yu Lee, Saining Xie, Patrick Gallagher, Zhengyou Zhang, and Zhuowen Tu. 2015. Deeply-supervised nets. In *Artificial Intelligence and Statistics*, pages 562–570.

Kevin Lin, Dianqi Li, Xiaodong He, Zhengyou Zhang, and Ming-Ting Sun. 2017. Adversarial ranking for language generation. In *Advances in Neural Information Processing Systems*, pages 3155–3165.

David JC MacKay. 2003. *Information theory, inference and learning algorithms*. Cambridge university press.

Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. 2016. Asynchronous methods for deep reinforcement learning. In *International conference on machine learning*, pages 1928–1937.

David Pfau and Oriol Vinyals. 2016. Connecting generative adversarial networks and actor-critic methods. *arXiv preprint arXiv:1610.01945*.

Marc'Aurelio Ranzato, Sumit Chopra, Michael Auli, and Wojciech Zaremba. 2015. Sequence level training with recurrent neural networks. *arXiv preprint arXiv:1511.06732*.

Zhan Shi, Xinchi Chen, Xipeng Qiu, and Xuanjing Huang. 2018. Towards diverse text generation with inverse reinforcement learning. *arXiv preprint arXiv:1804.11258*.

Richard S Sutton, David A McAllester, Satinder P Singh, and Yishay Mansour. 2000. Policy gradient methods for reinforcement learning with function approximation. In *Advances in neural information processing systems*, pages 1057–1063.

Guy Tevet, Gavriel Habib, Vered Shwartz, and Jonathan Berant. 2018. Evaluating text gans as language models. *arXiv preprint arXiv:1810.12686*.

Mnih Volodymyr, Koray Kavukcuoglu, David Silver, Alex Graves, and Ioannis Antonoglou. 2013. Playing atari with deep reinforcement learning. In *NIPS Deep Learning Workshop*.

Ronald J Williams and David Zipser. 1989. A learning algorithm for continually running fully recurrent neural networks. *Neural computation*, 1(2):270–280.

Yifan Xu, Lu Dai, Udaikaran Singh, Kening Zhang, and Zhuowen Tu. 2019. Neural program synthesis by self-learning. *arXiv preprint arXiv:1910.05865*.

Lantao Yu, Weinan Zhang, Jun Wang, and Yong Yu. 2017. Seqgan: Sequence generative adversarial nets with policy gradient. In *AAAI*, pages 2852–2858.

Yaoming Zhu, Sidi Lu, Lei Zheng, Jiaxian Guo, Weinan Zhang, Jun Wang, and Yong Yu. 2018. Texygen: A benchmarking platform for text generation models. *arXiv preprint arXiv:1802.01886*.