

QUANTIFYING EXPOSURE BIAS FOR NEURAL LANGUAGE GENERATION

Anonymous authors

Paper under double-blind review

ABSTRACT

The exposure bias problem refers to the training-testing discrepancy caused by teacher forcing in maximum likelihood estimation (MLE) training, for autoregressive neural network language models (LM). It has been regarded as a central problem for natural language generation (NLG) model training. Although a lot of algorithms have been proposed to avoid teacher forcing and therefore to alleviate exposure bias, there is little work showing how serious the exposure bias problem actually is. In this work, we first identify and analyze the self-recovery ability of MLE-trained LM, which casts doubt on the seriousness of exposure bias. We then develop a precise, quantifiable definition for exposure bias. Based on that definition, we design experiments to measure the seriousness of exposure bias. Surprisingly, we find that removing the training-testing discrepancy only brings very little performance gain, in both real and synthetic settings. With these results, we conclude that on the contrary to popular belief, the exposure bias problem is only a minor problem for MLE-based LM training.

1 INTRODUCTION

Language model (LM) is a central module for natural language generation (NLG) tasks (Young et al., 2017) such as machine translation (Wu et al., 2017), dialogue response generation (Li et al., 2017), image captioning (Lin et al., 2014), etc. For decades, maximum likelihood estimation (MLE) has been the most widely used objective for LM training. However, there is a popular belief in the natural language processing (NLP) community that standard MLE training will cause “exposure bias” and lead to a performance degradation during the test-time language generation.

The exposure bias problem (Bengio et al., 2015; Ranzato et al., 2016) refers to the following discrepancy between MLE training and test-time generation for language models: During training, the language model predicts the next word conditioned on history words sampled from the ground-truth data distribution. And during generation, the model generates words conditioned on history sequences generated by the model itself. However, due to the *exposure* to real data during training, the language model is *biased* to only perform well on the ground-truth history distribution. As a result, during generation the errors will accumulate along the generated sequence, and the distribution generated by the model will be distorted. The forced exposure to ground-truth data during training is also referred to as “teacher forcing”.

Given its definition, the exposure bias problem could rise in the general cases when the model needs to make a sequence of decisions or generations (e.g. music/pixel/speech generation (Lamb et al., 2016)). In this work, we focus on the task of language generation, because the exposure bias problem is originally proposed in this field (Bengio et al., 2015), and has since attracted huge research attention.

In order to avoid teacher forcing, many training algorithms (Bengio et al., 2015; Lamb et al., 2016; Ranzato et al., 2016; Yu et al., 2016; Zhu et al., 2018; Lu et al., 2018; Lin et al., 2017; Guo et al., 2017; Rajeswar et al., 2017; Wiseman & Rush, 2016; Nie et al., 2019; Shi et al., 2018) have been proposed as alternatives to MLE training. Most of these works utilize techniques from generative adversarial network (GAN) (Goodfellow et al., 2014) or reinforcement learning (RL) (Sutton & Barto, 1998). In this paper, we refer to these algorithms as non-MLE methods or text GANs.

Despite the huge research efforts devoted to alleviate exposure bias, surprisingly, its existence or significance is much less studied. In particular, to the best of our knowledge, no existing work

<p>Model Samples as Hisotry → Model Samples</p> <p>When asked about how she thinks about the games, Flocke → dislikes most of those about it, citing instances of paranoia in her heart and trembling temper, which infuriated him.</p> <p>Approximately 500 Finns became sick since early October when sleeping → in their sleeping bags. On 3 October, the "Red Guard", which had been organized two months previously by Marius Kuusinen ...</p> <p>The entire key results of the arc be obtained through → unifying methods to construct the prologue, three pieces could be combined instead of need to provide a final chapter.</p>
<p>Data Samples as History → Model Samples</p> <p>Most of what is known of Du Fu <unk> s → life is clear and graphic descriptions, memoirs, commentaries on storyboards, and descriptions of Canadian settlers. More than 60 biographies and ...</p> <p>In the early 730s, he travelled in the Jiangsu province → of Asia after Ashras ibn Abdallah al-Sulami visited Quanzhou in Bukhara, the capital of Turkmenistan and a native of the ...</p> <p>Since the Song dynasty, critics have called Du Fu the → "master poet", a product of his use of Du Fu scenes to establish the empress's nature and to emphasize his ...</p>
<p>Shuffled Data Samples as History → Model Samples</p> <p>is Du of <unk> s known Fu of life what → Most claimed was his tragic adaptation of John Ching's The Janus of Hades, translated by disgraced performer just months before ...</p> <p>in the he travelled, the early Jiangsu In 730s, he → attended a mission on the peninsula. He soon moved to Monkwearmouth, on the northern shores of Baffin Bay in The ...</p> <p>Since, the called have Song Fu critics dynasty Du Fu, → who Zhang historians have included, have not rivaled HABS's Web site held for 253 years. In 2015, HABS-based producers Oronoco ...</p>
<p>Random Sequences as History → Model Samples</p> <p>...RANDOM... surface leader Game → after a failed attempt to test her effectively in three fleets falling to I-30. This went unnoticed by most ichthyologists; none understood either strict rules ...</p> <p>...RANDOM... faster elephant emperor → decorations with Rocky Mountain state exploit by linking all black geese to 1970s planning regulations that prohibit slaughter of snake species.</p> <p>...RANDOM... hitting remained prominently → from the system as she witnessed no mention of criteria in the text. Douglas Turner noted then that Gottesfeld may have assumed ...</p>

Table 1: Samples of a MLE-trained STOA transformer LM when fed with different types of length-10 history prefix. To save space, we omitted the first 7 words of the random history.

attempts to directly show the seriousness of exposure bias in an empirical or theoretical way. This work is motivated by the belief that a good solution should be built upon a testable and quantifiable problem definition. In this rest of this paper, we first identify the "self-recovery" ability of popular LM models, which casts doubt on the original claim of exposure bias. We then develop a precise and quantifiable definition of exposure bias, and validate its seriousness in controlled experiments.

2 MOTIVATION: THE SELF-RECOVERY ABILITY

To study the seriousness of exposure bias in standard MLE LM training, we first stress that the following methodology, although tempting, is wrong: *If we can rigorously show that the non-MLE methods proposed to avoid teacher forcing do indeed bring solid generation performance gain, then we can conclude exposure bias is a meaningful problem for the original MLE training.* The reason is that we typically do not know the exact underlying reason for the performance gain. For example, despite the huge success of the batch normalization technique in deep learning, whether "internal covariate shift" (which is the motivation of batch norm) exists in deep neural network training remains a question (Santurkar et al., 2018). Therefore, in this work we seek a direct way to validate the seriousness of exposure bias.

We focus on the following informal claim that immediately follows from the original definition of exposure bias: During generation, if we set the history distribution to be the ground-truth data distribution instead of the model's own distribution (now that there is no discrepancy between training and testing), then the model's language generation quality should be much better (we will formalize this notion in Section 4 and 5).

We start with the following qualitative analysis. We feed a MLE-trained transformer LM on wiki-103 data-set (Baevski & Auli, 2018) with four kinds of prefixes: model's own samples, data samples, shuffled (word-level) data samples or samples from a uniform random distribution. Then we let the model complete the sentence given these prefixes as history. We list some samples in Table 1 and more in Appendix A (this experiment is also repeated for a LSTM LM).

Assuming the seriousness of exposure bias, we expect the quality of generated sentence-completion samples with real-data prefixes to be significantly better than the ones from prefixes of model samples. However, by manual inspection, we do not observe noticeable differences in sample quality. More surprisingly, the model is still able to generate relevant and fairly high-quality samples from shuffled prefixes. Even in the extreme case where random sequences are fed, the model is able to generate reasonable sentences. Due to the recent increasing interest of solving exposure bias in the field of neural machine translation (NMT) (Zhang et al., 2019), we repeat the above experiment in a standard NMT setting in Appendix A, and get very similar observations.

These experiments clearly show that the MLE-trained auto-regressive LMs have the *self-recovery* ability, i.e. the model is able to recover from artificially distorted history input, and generate reasonably high-quality samples. This phenomenon is clearly in contradiction with the popular claim of exposure bias, that the error induced by the mismatch between history and data distribution should **accumulate** during the generation process.

Motivated by these experiments, in the following sections, we turn to more rigorous methods to quantify the significance of exposure bias. Note that our quantification approaches will be independent of the training procedure and only require inference from the trained model.

3 NOTATIONS

The task of auto-regressive language modelling is to learn the probability distribution of the $(l+1)$ th word W_{l+1} in a sentence conditioned on the word history $W_{1:l} := (W_1, \dots, W_l)$. Here, we use the uppercase $W_i \in V$ to denote a discrete random variable distributed across the vocabulary V . The lower-case w is used to denote some particular word in V . Given a training data-set D consisting of sentences of length L , the standard MLE training minimizes the negative log-likelihood below:

$$L_{\text{MLE}} = \mathbb{E}_{W_{1:L} \sim D} -\frac{1}{L} \sum_{l=1}^L \log P_M(W_l | W_{1:l-1}), \quad (1)$$

Note that in this work we assume all sentences are of length L for simplicity.

We denote the generation distribution of the trained LM as P_M , and the ground-truth data distribution as P_D . Readers can assume P_M refers to the generation distribution of a LSTM LM (Hochreiter & Schmidhuber, 1997; Sundermeyer et al., 2012) or a transformer LM (Baevski & Auli, 2018; Dai et al., 2019) trained with MLE objective, which is the major subject of this study. We will mainly present results on LSTM based models to facilitate comparison with text-GAN works (listed in Section 1), which are mostly implemented on LSTM models. We will also provide results with the transformer model, with very similar observations or measurements.

Our quantification mainly relies on the measurements of the distance from the model’s generation distribution to the data distribution. Hence we define the following notations to simplify expressions. Let \mathcal{P} denote the set of probability distributions on the vocabulary V . Let d denote a distance measure between distributions (e.g. total variation distance), $d : \mathcal{P} \times \mathcal{P} \rightarrow \mathbb{R}_{\geq 0}$.

4 A (WRONG) QUANTIFICATION USING MARGINAL DISTRIBUTION

In this section, we propose an intuitive and seemingly correct quantification approach using marginal distributions. The approach can be applied to real-world text data experiments, but it has some lethal weak point. The discussion will lead us to our final precise definition of exposure bias in Section 5.

4.1 METHOD

Assuming a given history length l , we consider the marginal distribution of W_{l+1} from the following three random process:

- Draw word sequences of length L from the data distribution P_D . Denote the marginal distribution of the random variable at position $l+1$ (W_{l+1}) as $P_{D|D}^{l+1}$, where

$$P_{D|D}^{l+1}(w) = \mathbb{E}_{W_{1:l} \sim P_D} [P_D(w | W_{1:l})]. \quad (2)$$

- Draw word sequences of length L from the model distribution P_M . Denote the marginal distribution of the random variable at position $l+1$ as $P_{M|M}^{l+1}$, where

$$P_{M|M}^{l+1}(w) = \mathbb{E}_{W_{1:l} \sim P_M} [P_M(w | W_{1:l})]. \quad (3)$$

- First draw $W_{1:l}$ from P_D , then draw W_{l+1} from $P_M(\cdot|W_{1:l})$. Denote the marginal distribution of the random variable at position $l + 1$ as $P_{M|D}^{l+1}$, where

$$P_{M|D}^{l+1}(w) = \mathbb{E}_{W_{1:l} \sim P_D}[P_M(w | W_{1:l})]. \quad (4)$$

By the definition of exposure bias, $P_{M|M}^{l+1}$ suffers from the training-testing discrepancy, while $P_{M|D}^{l+1}$ should be closer to the true distribution $P_{D|D}^{l+1}$. To measure this discrepancy, define the marginal generation deviation (MGD) at history length l of history distribution P_H with metric d as

$$\text{MGD}(P_{M|H}, l, d) = d(P_{M|H}^{l+1}, P_{D|D}^{l+1}) \quad (5)$$

where $P_H \in \{P_M, P_D\}$ denotes the history distribution. MGD measures the deviation of the marginal distribution of W^{l+1} from ground-truth data distribution.

Finally, we define the *rate of exposure bias* (EB-M) at history length l of model P_M as the ratio (discrepancy) between the MGD measurements when two different history distributions are fed:

$$\text{EB-M}(P_M, l, d) = \frac{\text{MGD}(P_{M|M}, l, d)}{\text{MGD}(P_{M|D}, l, d)} \quad (6)$$

For MLE-trained models, EB-M¹ is expected to be larger than 1, and larger EB-M indicates a more serious exposure bias problem for the trained model. For the metric d , we consider two popular probability metrics: total variation distance (denoted as d_{TV}), and Jensen-Shannon divergence (denoted as d_{JS}).

4.2 EXPERIMENTS AND DISCUSSION

In this section, we focus on answering the following question: “Does the EB-M measurement correctly reflect the significance of exposure bias?” In short, our answer is *not really*. The problem is that the distortion of the marginal $P_{M|M}^{l+1}$ is not only affected by the presumably existing exposure bias problem alone, but also by the mismatch between the history distribution P_M from P_D for $W_{1:l}$, which grows with the length of the history. Therefore, even if the measured EB-M is significantly larger than one, we can not conclude that exposure bias causes serious deterioration. We provide an example to illustrate this argument:

Example 1. Suppose $L = 2$, and $V = \{A, B\}$. P_D and P_M are crafted as follows: $P_D(AA) = P_D(BB) = 0.5$, $P_D(AB) = P_D(BA) = 0$; And $P_M(W_1 = A) = 1$, $P_M(W_2 = A|W_1 = A) = 1$, $P_M(W_2 = B|W_1 = B) = 1$.

In Example 1, $\text{MGD}(P_{M|M}, 1, d_{TV}) = 0.5$ and $\text{MGD}(P_{M|D}, 1, d_{TV}) = 0$, which gives $\text{EB-M}(P_M, 1, d_{TV}) = \infty$. However, the only problem P_M has is the mismatch between the history distributions (P_M and P_D) for W_1 .

The next set of experiments also suggest that EB-M does not precisely reflect exposure bias. On the EMNLP-news data-set (specified in Appendix B), we compare EB-M measurements for several non-MLE training methods with the baseline MLE model. We include results for Scheduled Sampling (SS) (Bengio et al., 2015), Cooperative Training (CoT) (Lu et al., 2018), and Adversarial Ranking (RankGAN) (Lin et al., 2017). We provide implementation details for non-MLE methods in Appendix C. Intuitively, these methods will cause the model to be *biased* to behave well with model samples as history, instead of data samples. Therefore, we expect EB-M measurement for non-MLE trained models to be smaller than MLE trained models. However, Figure 1 shows that the measurements for different training frameworks are almost the same. We believe the reason is that the EB-M measurements are only reflecting the trivial mismatch between the history distributions.

Is it possible that the original definition of exposure bias (Bengio et al., 2015; Ranzato et al., 2016) exactly refers to this mismatch between the model and data history distributions? However, note that this mismatch is inevitable for any imperfect model, and non-MLE training algorithms can not solve it. We believe a better, more precise definition is needed to discriminate exposure bias from this trivial mismatch. Motivated by this view, we propose a second approach in the section below.

¹Note that one can also directly measure $d(P_{M|M}^{l+1}, P_{M|D}^{l+1})$, but in that way, we can not tell which distribution is better.

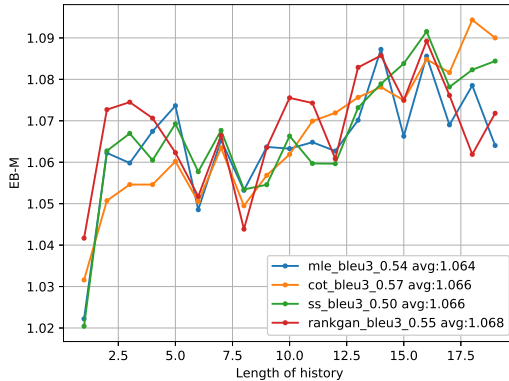


Figure 1: EB-M (with metric d_{JS}) comparison for MLE and non-MLE training on EMNLP-News data. For each training method, we show corpus-BLEU (Yu et al., 2016) measurement using the test-set as reference set in the legend.

5 A QUANTIFICATION APPROACH USING CONDITIONAL DISTRIBUTION

5.1 METHOD

Following the discussion in the last section, we wish our measurement to be *independent* of the quality of the history distribution. In light of that, we design a quantity to measure the model’s conditional generation quality. Let $P_H \in \{P_M, P_D\}$ denote the history distribution as in the MGD definition (5). With history length l fixed, we define the conditional generation deviation (CGD) with history distribution P_H for P_M using metric d as:

$$\text{CGD}(P_{M|H}, l, d) = \mathbb{E}_{W_{1:l} \sim P_H} [d(P_M(\cdot | W_{1:l}), P_D(\cdot | W_{1:l}))] \quad (7)$$

where we assume that $P_D(\cdot | W_{1:l})$ is computable, and use it to measure the quality of the model’s conditional distribution. For the choice of the distribution distance d , in addition to d_{TV} and d_{JS} , we introduce greedy decoding divergence (d_{GD}) defined as:

$$d_{GD}(P, Q) = \mathbb{1}(\arg \max_i P_i \neq \arg \max_i Q_i) \quad (8)$$

where $\mathbb{1}$ is the indicator function, and $P, Q \in \mathcal{P}$. The distance d_{GD}^2 reflects the model’s accuracy during greedy decoding.

Similar to MGD, exposure bias should imply a significant gap between $\text{CGD}(P_{M|M}, l, d)$ and $\text{CGD}(P_{M|D}, l, d)$. We again define rate of exposure bias at history length l with metric d to be:

$$\text{EB-C}(P_M, l, d) = \frac{\text{CGD}(P_{M|M}, l, d)}{\text{CGD}(P_{M|D}, l, d)} \quad (9)$$

For our definition of EB-C, a natural question is why we only focus on the generation distribution of the very next word. The reason is we want to precisely measure how the error caused by the history part affect the generation part, by keeping them separate. If we measure the deviation of, for example, two sampled tokens, the definition will be confusing: Because the second sampled token will be affected not only by the accumulated error induced by the history (sampled from the model), but also by the first generated token as history. To get a better understanding of the intuition behind the definition of EB-C, we recommend readers to read Appendix A about our NMT experiment.

5.2 SYNTHETIC EXPERIMENTS AND DISCUSSION

Since CGD requires inference for ground-truth data distribution P_D , we first consider experiments in a synthetic setting. In text-GAN literature (Yu et al., 2016; Lin et al., 2017), a randomly-initialized

² d_{GD} qualifies as a pseudometric in mathematics.

one-layer LSTM model with hidden dimension of 32 is usually used as P_D in synthetic experiments (we denote this setting as M_{32}^{random}). However, the model is small-scale and does not reflect any structure existing in real-world text. To improve upon this approach, we take the MLE baseline model trained on EMNLP-news data (described in Appendix B) as P_D in this synthetic setting. We denote the data model (P_D) as M_{512}^{news} . We then train two LSTM LM (P_M) with different capacities using samples from the data model, with the standard MLE objective. One is a one-layer LSTM with hidden width of 512 (denoted as LSTM-512), the other one is with hidden width of 32 (denoted as LSTM-32).

We train P_M for 100 epochs using the Adam optimizer with learning rate 0.001. In each epoch, 250k sentences (same to the size of the original EMNLP-news data) of length $L = 50$ are sampled from $M_{news-512}$ as training data to avoid over-fitting. We show perplexity (PPL) results of the trained models in Appendix F. Finally, EB-C is calculated using $100k^3$ samples from P_M and P_D .

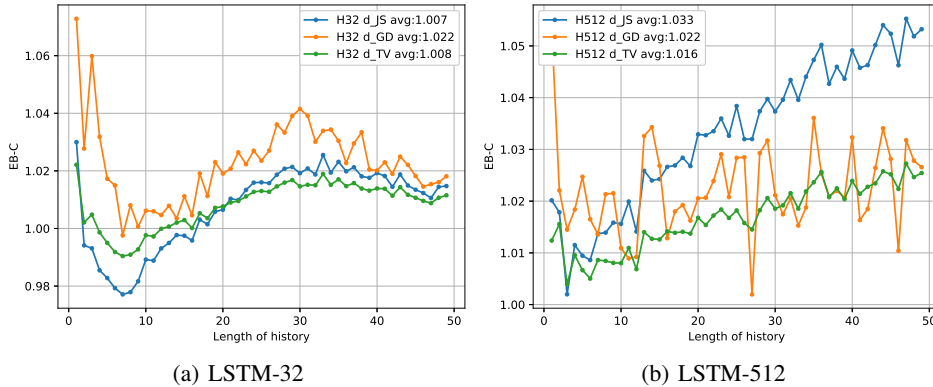


Figure 2: EB-C measurement for the LSTM-32 and LSTM-512 model with different metrics. Also average value of EB-C along all history length is shown in the legend.

In Figure 2, we show EB-C measurements with different metrics d_m , and the two models give similar results. It is shown that EB-C has a steady but slow increasing trend as history length increases. This is expected as a consequence of exposure bias, because P_M deviates farther from P_D as history length increases. However, the average value of EB-C is less than 1.03 (the largest average value is from d_{JS} for the LSTM-512 experiment), meaning that the gap between $CGD(P_{M|M}, l, d)$ and $CGD(P_{M|D}, l, d)$ is not large. Also, note that in most NLG applications (such as machine translation or image captioning), the generated sequence typically has short length (less than 20). In that range of history length, the EB-C measurements that exposure bias only has minimal influence.

In Appendix E, we repeat the experiment for a transformer LM (Dai et al., 2019), and get very similar EB-C measurements. These measurements imply a striking conclusion: *(Informal) Even if all the bad effects from exposure bias for MLE LM training are removed, the relative performance gain is at most 3%. If the sequence length is not very long, the gain is less than 1%.*

To dive deeper into the cause of the gap in CGD, we experiment with corrupted versions of P_M as history distribution. We first specify a corrupt rate $c \in [0, 1]$, and randomly substitute words in a history sample from P_M to a “noise” word drawn uniformly from the vocabulary with probability c . Consequently, larger c will cause the history distribution to deviate farther from the ground-truth P_D . In Figure 3, we show CGD measurement versus the corrupted history $P_M^{corrupt}$. Large gaps are observed between $CGD(P_{M|M^{corrupt}})$ and $CGD(P_{M|D})$. Therefore, the small gap between $CGD(P_{M|M})$ and $CGD(P_{M|D})$ in Figure 2 results from the small deviation between the history distribution P_M and P_D . In other word, P_M has learned a “good enough” distribution that is able to keep it in the well-behaving region during sampling.

With these observations, we conclude that, in the synthetic setting considered, exposure bias does exist, but is much less serious than it is presumed to be. Although there exists mismatch between the history distribution P_M and P_D , the mismatch is still in the model’s “comfortable zone”. In

³We show that we can get stable measurements using 100k samples in Appendix E.

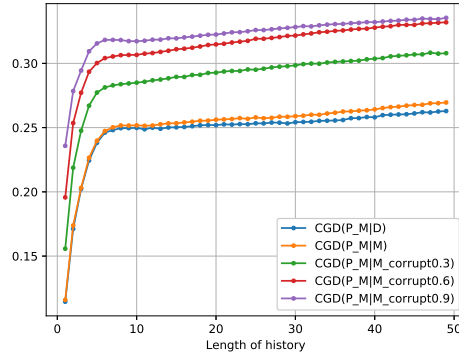


Figure 3: CGD measurement for corrupted P_M (with d_{TV}) for the LSTM-512 synthetic experiment.

other words, the LSTM LM is more robust than exposure bias claims it to be. To concretize this argument, we provide an example LM and show that MLE training is unlikely to generate models with a large EB-C value.

Example 2. Again suppose $L = 2$, and $V = \{A, B\}$, the ground-truth data distribution is uniform on $\{AA, AB, BB, BA\}$. P_M is crafted as follows: $P_M(W_1 = A) = 0.9, P_M(W_2 = A|W_1 = A) = 0.9, P_M(W_2 = A|W_1 = B) = 0.5$. Note that the model behaves bad when $W_1 = A$, which is of high probability during sampling.

In Example 2, $CGD(P_{M|D}, 1, d_{TV}) = 0.2$ and $CGD(P_{M|M}, 1, d_{TV}) = 0.36$, so $EB-C(P_M, 1, d_{TV}) = 1.8$. However, this crafted model is unlikely to be an outcome of MLE training. The fact that $P_M(\cdot|W_1 = B)$ is better modeled indicates that in the training data more sentences begin with $W_1 = B$ than $W_1 = A$. So MLE training should assign more probability to $P_M(W_1 = B)$, not the other way around⁴.

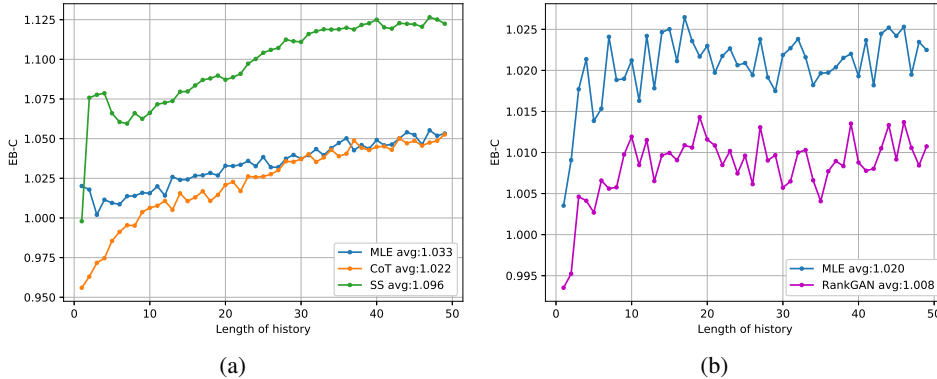


Figure 4: (a): EB-C measurements (with d_{JS}) for comparing non-MLE methods in the LSTM-512 synthetic experiment. (b): EB-C measurements for comparing RankGAN and MLE for the M_{32}^{random} synthetic experiment, the metric used is d_{JS} .

Finally, we use EB-C to compare MLE and non-MLE training. We compare MLE against CoT, SS, RankGAN in the synthetic experiments, and results are shown in Figure 4. Note that the RankGAN experiments are conducted in the M_{32}^{random} setting⁵, as we find it hard to do a fast implementation of RankGAN for the LSTM-512 setting. We find that RankGAN and CoT gives lower EB-C measurements than MLE, which is expected, as these methods avoid teacher forcing. For CoT, at short

⁴If we change to $P_M(W_1 = A) = 0.1$, then $EB-C(P_M, 1, d_{TV})$ will be 0.2, meaning that the model has better conditional generation performance during sampling

⁵The MLE model is used as the pre-trained model for the RankGAN generator. The MLE model has an oracle NLL of 8.67, and RankGAN’s oracle NLL is 8.55.

History: in august 2009 , vh1 classic aired music videos from the tv special around the beatles (1964) ,
Choices: 0: the 1: including 2: which 3: featuring 4: and 5:[None of the above is plausible]

History: contingent units were accused of being armed with antiquated weapons , while some local peasants took claims of having been
Choices: 0: conscripted 1: armed 2: involved 3: tricked 4: intimidated 5:[None of the above is plausible]

Table 2: An illustration for the next word collection process. The choices are shuffled. The first history sample is from real data, and the second history sample is from the trained model.

History Length (l)	$CGD(P_{M D}, l, d_{GD})$	$CGD(P_{M M}, l, d_{GD})$	EB-C
10	0.723	0.738	1.019
20	0.736	0.752	1.021
30	0.732	0.754	1.029

Table 3: EB-C measurements with human as P_D .

history length, EB-C is even less than 1. We believe the reason is that CoT tries to make the model be biased to behave better when fed with model samples. However, SS gives worse EB-C measurements comparing to MLE, which we currently do not have a good explanation. We refer readers to [Huszar \(2015\)](#) for a discussion about the SS objective.

To the best of our knowledge, this is the first direct empirical evidence that text GAN does indeed alleviate the exposure bias problem. It also indicates that EB-C correctly reflect the significance of exposure bias. We believe the reason for why EB-C is still not less than 1 is that, text GANs still rely on MLE pre-training a lot.

5.3 EB-C MEASUREMENTS WITH REAL HUMAN AS P_D

In this section, we design experiments to efficiently estimate EB-C for a SOTA transformer LM with real human as P_D , by utilizing the Amazon Mechanical Turk (AMT) platform. Given a MLE-trained LM as P_M , by examining the definition of EB-C in Equation 9 and 7, it is clear the only obstacle is that we don't have access to $P_D(\cdot | W_{1:l})$ with a given history $W_{1:l}$. So, in this section, we focus on the greedy decoding divergence (d_{GD}) metric (Equation 8), which only requires the turkers to give the most probable next word prediction, instead of the full distribution (which is clearly intractable).

In our preliminary trials, we find it is still very hard for a person to guess the next word, even with real data history samples. The reason is that the vocabulary is very big, and the turkers may be not familiar with the context (e.g. wikipedia). To alleviate that problem, we design the following simplification: For a given history, we let the model output its top-5 next word prediction, then we only ask the turkers to choose among the 5 choices (the turker can also express that he/she thinks none of them is likely). Finally, we examine whether the turker's choice is indeed the model's top-1 prediction. We illustrate this process in Table 2.

We use the code of Transformer-XL ([Dai et al., 2019](#)) to train a SOTA transformer LM on the wiki-103 data-set. We favour the wiki-103 data-set because it is large-scale and has long (over 30 words) paragraphs, which is useful for the measurements of exposure bias. The model is a 16-layer transformer-xl model with hidden dimension of 410. Since the estimation of $CGD(P_{M|D}, l, d)$ requires large amounts of unseen real data samples, we use half of the wiki-103 training data (around 900k sentences and 50m words) to train the model P_M , and save the other half as samples from P_D . Other training configurations (learning rate, batch size, etc.) are not changed⁶. The resulting model P_M has a test-set PPL of 27.81 (if trained on full training data, the PPL will be 24.02).

We collect data to estimate EB-C at history length 10, 20, and 30. For each length and history model (P_M or P_D) pair, we collect 10k d_{GD} samples (via next-word prediction) from turkers on the AMT platform. More details about the AMT setup are provided in Appendix D. The results are shown in Table 3. The EB-C measurements are strikingly similar to the results in our synthetic experiments in that, removing the training-testing discrepancy only gives around 2% of relative performance gain.

⁶https://github.com/kimiyoung/transformer-xl/blob/master/pytorch/run_wt103_base.sh

This result further strengthens our claim that exposure bias is only a minor problem for MLE-based LM training.

6 RELATED WORKS

Several recent works attempt to carefully evaluate whether the non-MLE training methods (e.g. adversarial training) can give superior NLG performance than standard MLE training for RNN LM. [Caccia et al. \(2018\)](#) tunes a “temperature” parameter in the softmax output, and evaluate models over the whole quality-diversity spectrum. [Semeniuta et al. \(2018\)](#) proposes to use “Reverse Language Model score” or “Frechet InferSent Distance” to evaluate the model’s generation performance. [Tevet et al. \(2018\)](#) proposes a method for approximating a distribution over tokens from a GAN, and then evaluate the model with standard LM metrics.

These works arrive at a similar conclusion: The general performance of Text GANs is not convincingly better, or even worse, than standard MLE training. Hence to some extent, they imply that exposure bias may be not a serious problem in MLE training. However, as we argued in Section 2, one can not draw direct conclusions about exposure bias with these results. For example, it is also possible that exposure bias is indeed serious for MLE training, but text GAN does not solve the problem well enough.

7 CONCLUSION AND DISCUSSION

In this work, we first identify the self-recovery ability of MLE-trained LM, which casts doubt on the seriousness of exposure bias, which has been regarded as a central problem for MLE training by the LM community. We then explore two intuitive approaches to quantify the significance of exposure bias for LM training. The first quantification EB-M relies on the marginal generation distribution and reveals some vagueness in the original definition of exposure bias. We argue that we should focus on the model’s generation performance in terms of its conditional distribution and propose a second quantification EB-C, which we regard as the precise definition for exposure bias.

We design a evaluation of EB-C at different history length with real human (turkers from AMT) as the data model, for a SOTA transformer LM. It is shown that removing the training-testing discrepancy only gives around 2% of performance gain. Our synthetic experiments also gives very similar measurements. By analyzing EB-C measurements with perturbed history samples, we hypothesise that although the mismatch between the data and model distribution for history prefix exists, it is still in the model’s “comfortable zone”. With these results, we claim that on the contrary to the popular belief, exposure bias is only a minor problem in MLE-based LM training.

To wrap up, we discuss the fundamental question “Is MLE training really biased?”, from the perspective of objective functions. Note that the MLE objective (1) can be re-written as:

$$\begin{aligned} \arg \min_{\theta} \mathbb{E}_{W_{1:L} \sim P_D} \frac{-1}{L} \sum_{l=1}^L \log P_M(W_l | W_{1:l-1}) &= \arg \min_{\theta} \mathbb{E}_{W \sim P_D} -\log P_M(W) \\ &= \arg \min_{\theta} \mathbb{E}_{W \sim P_D} \log \frac{P_D(W)}{P_M(W)} = \arg \min_{\theta} D_{KL}(P_D || P_M) \end{aligned} \quad (10)$$

where D_{KL} denotes the Kullback-Leibler divergence, and θ denotes the trainable parameters in P_M . Therefore, MLE training is minizing the divergence from P_M , which is exactly the model’s sampling distribution, from P_D . While it’s true that the training is “exposed” to data samples, we can not simply deduce the objective is “biased”.

We want to end our discussion with two remarks. First, the proposed quantification approaches should not be used as the only metric for NLG. For example, a position-aware uni-gram LM, which generates words independent of previous context, has no exposure bias problem and can pass our test easily. Second, the intention of this work is not to discourage researchers from exploring non-MLE training algorithms for LM. It is completely possible that an training objective different from $D_{KL}(P_D || P_M)$, such as $JSD(P_D || P_M)$, can lead to better generation performance ([Lu et al., 2018](#); [Huszar, 2015](#)). However, though non-MLE algorithms avoid teacher forcing, these algorithms (using GAN or RL for example) are usually less stable and more difficult to tune. Given that the quantified measurement of exposure bias is insignificant, we think it should be questioned whether adopting these techniques to avoid exposure bias is a wise trade-off.

REFERENCES

- Alexei Baevski and Michael Auli. Adaptive input representations for neural language modeling. *CoRR*, abs/1809.10853, 2018. URL <http://arxiv.org/abs/1809.10853>.
- Samy Bengio, Oriol Vinyals, Navdeep Jaitly, and Noam Shazeer. Scheduled sampling for sequence prediction with recurrent neural networks. In *Proceedings of the 28th International Conference on Neural Information Processing Systems - Volume 1*, NIPS’15, pp. 1171–1179, Cambridge, MA, USA, 2015. MIT Press. URL <http://dl.acm.org/citation.cfm?id=2969239.2969370>.
- Massimo Caccia, Lucas Caccia, William Fedus, Hugo Larochelle, Joelle Pineau, and Laurent Charlin. Language gans falling short. *CoRR*, abs/1811.02549, 2018.
- Zihang Dai, Zhilin Yang, Yiming Yang, Jaime G. Carbonell, Quoc V. Le, and Ruslan Salakhutdinov. Transformer-xl: Attentive language models beyond a fixed-length context. *CoRR*, abs/1901.02860, 2019. URL <http://arxiv.org/abs/1901.02860>.
- Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Proceedings of the 27th International Conference on Neural Information Processing Systems - Volume 2*, NIPS’14, pp. 2672–2680, Cambridge, MA, USA, 2014. MIT Press.
- Jiaxian Guo, Sidi Lu, Han Cai, Weinan Zhang, Yong Yu, and Jun Wang. Long text generation via adversarial training with leaked information. *CoRR*, abs/1709.08624, 2017.
- Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8): 1735–1780, 1997.
- Ferenc Huszr. How (not) to train your generative model: Scheduled sampling, likelihood, adversary? *CoRR*, abs/1511.05101, 2015.
- Alex Lamb, Anirudh Goyal, Ying Zhang, Saizheng Zhang, Aaron C. Courville, and Yoshua Bengio. Professor forcing: A new algorithm for training recurrent networks. In *NIPS*, 2016.
- Jiwei Li, Will Monroe, Tianlin Shi, Alan Ritter, and Dan Jurafsky. Adversarial learning for neural dialogue generation. *CoRR*, abs/1701.06547, 2017.
- Kevin Lin, Dianqi Li, Xiaodong He, Zhengyou Zhang, and Ming-ting Sun. Adversarial ranking for language generation. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Advances in Neural Information Processing Systems 30*, pp. 3155–3165. Curran Associates, Inc., 2017. URL <http://papers.nips.cc/paper/6908-adversarial-ranking-for-language-generation.pdf>.
- Tsung-Yi Lin, Michael Maire, Serge J. Belongie, Lubomir D. Bourdev, Ross B. Girshick, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft COCO: common objects in context. *CoRR*, abs/1405.0312, 2014.
- Sidi Lu, Lantao Yu, Weinan Zhang, and Yong Yu. Cot: Cooperative training for generative modeling. *CoRR*, abs/1804.03782, 2018.
- Weili Nie, Nina Narodytska, and Ankit Patel. RelGAN: Relational generative adversarial networks for text generation. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=rJedV3R5tm>.
- Sai Rajeswar, Sandeep Subramanian, Francis Dutil, Christopher Joseph Pal, and Aaron C. Courville. Adversarial generation of natural language, 2017.
- Marc’Aurelio Ranzato, Sumit Chopra, Michael Auli, and Wojciech Zaremba. Sequence level training with recurrent neural networks. In *ICLR*, 2016.
- Shibani Santurkar, Dimitris Tsipras, Andrew Ilyas, and Aleksander Madry. How does batch normalization help optimization?(no, it is not about internal covariate shift). *arXiv preprint arXiv:1805.11604*, 2018. URL <https://papers.nips.cc/paper/7515-how-does-batch-normalization-help-optimization>.

- Stanislau Semeniuta, Aliaksei Severyn, and Sylvain Gelly. On accurate evaluation of gans for language generation. *CoRR*, abs/1806.04936, 2018.
- Zhan Shi, Xinchu Chen, Xipeng Qiu, and Xuanjing Huang. Towards diverse text generation with inverse reinforcement learning. *CoRR*, abs/1804.11258, 2018.
- Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15:1929–1958, 2014. URL <http://jmlr.org/papers/v15/srivastava14a.html>.
- Martin Sundermeyer, Ralf Schlüter, and Hermann Ney. LSTM neural networks for language modeling. In *INTERSPEECH 2012, 13th Annual Conference of the International Speech Communication Association, Portland, Oregon, USA, September 9-13, 2012*, pp. 194–197, 2012. URL http://www.isca-speech.org/archive/interspeech_2012/i12_0194.html.
- Richard S. Sutton and Andrew G. Barto. *Introduction to Reinforcement Learning*. MIT Press, Cambridge, MA, USA, 1st edition, 1998. ISBN 0262193981.
- Guy Tevet, Gavriel Habib, Vered Shwartz, and Jonathan Berant. Evaluating text gans as language models. *CoRR*, abs/1810.12686, 2018.
- Sam Wiseman and Alexander M. Rush. Sequence-to-sequence learning as beam-search optimization. *CoRR*, abs/1606.02960, 2016.
- Lijun Wu, Yingce Xia, Li Zhao, Fei Tian, Tao Qin, Jianhuang Lai, and Tie-Yan Liu. Adversarial neural machine translation. *CoRR*, abs/1704.06933, 2017.
- Tom Young, Devamanyu Hazarika, Soujanya Poria, and Erik Cambria. Recent trends in deep learning based natural language processing. *CoRR*, abs/1708.02709, 2017. URL <http://arxiv.org/abs/1708.02709>.
- Lantao Yu, Weinan Zhang, Jun Wang, and Yong Yu. Seqgan: Sequence generative adversarial nets with policy gradient. *CoRR*, abs/1609.05473, 2016. URL <http://dblp.uni-trier.de/db/journals/corr/corr1609.html#YuZWY16>.
- Wen Zhang, Yang Feng, Fandong Meng, Di You, and Qun Liu. Bridging the gap between training and inference for neural machine translation. *CoRR*, abs/1906.02448, 2019. URL <http://arxiv.org/abs/1906.02448>.
- Yaoming Zhu, Sidi Lu, Lei Zheng, Jiaxian Guo, Weinan Zhang, Jun Wang, and Yong Yu. Tegygen: A benchmarking platform for text generation models. *SIGIR*, 2018.

A THE AUTO-RECOVERY ABILITY IN GENERAL LM AND NMT

In Table 6, we provide more samples of a MLE-trained transformer LM model (discussed in Section 2) when fed with different kinds of history. And in Table 7 we repeat the experiment for a LSTM-LM trained on the EMNLP-News data.

In Table 4 we repeat the preliminary experiment in Section 2 for a standard NMT setting. We train a 6-layer transformer model with hidden dimension 1024 on the IWSLT’14 German to English data set. We feed the trained model with types of prefix during decoding which represents different level of training-decoding discrepancy. Note that the source input is kept intact.

The result is very similar (or more striking) to our language model experiment, the data prefix does not seem to help, and in the extreme case of random prefix, the model still generates fairly good translation. In Section 2 we summarize this observation as the auto-recovery ability.

To interpret the UNREL3 results, we **should not** directly compare the translation generated from unrelated prefix to the target translation. In fact, we cannot even compare part of it (e.g. the part after the length-3 prefix). Instead, we highlight the surprising fact that although the model is forced to begin (conditioned) with a wrong prefix, it still comes up with a reasonable translation. This is not an

SOURCE:	sobald der ri@@ chter mich sah ,
REF:	and as soon as i walked inside , the ju@@ dge saw me coming in .
DATA3:	and as soon as the ju@@ dge saw me .
NORMAL:	as soon as the ju@@ dge saw me .
UNREL3:	what else is it that the ju@@ dge saw me ?
RAND3:	still take open action as the ju@@ dge saw me .

SOURCE:	ich fuhr also zum geri@@ cht .
REF:	and i got in my car and i went to this cour@@ thou@@ se .
DATA3:	and i got to the court .
NORMAL:	so i went to the court .
UNREL3:	the reasons for me to go to the court .
RAND3:	ge bor@@ last year , i went to court .

SOURCE:	ich bekam etwas angst vor technologie .
REF:	i found myself becoming a little bit of a tech@@ no@@ pho@@ be .
DATA3:	i found myself a little sc@@ ared of technology .
NORMAL:	i got a little sc@@ ared of technology .
UNREL3:	um , my fear of technology was with me .
RAND3:	kids -@@ ds i got a little sc@@ ared of technology .

SOURCE:	wir knnen das nicht einfach machen .
REF:	it is impossible to present such things in a society that is supposed to function .
DATA3:	it is impossible for us to do that .
NORMAL:	we can 't just do that .
UNREL3:	so i 'm not sure we can just do that .
RAND3:	first le@@ from here , we can 't just do that .

SOURCE:	das werde ich ihnen jetzt zeigen
REF:	so i 'm going to try and show you what you really get for 10 billion pi@@ x@@ els .
DATA3:	so i 'm going to show you this now .
NORMAL:	this is what i 'm going to show you .
UNREL3:	why did i show you that now ?
RAND3:	told ct happening to you now .

Table 4: A standard NMT transformer model fed with different types of length-3 history prefix. We did not do any cherry picking. The “@@” is because BPE tokenization is used. “DATA” means the first three output tokens are forced to be correct. “NORMAL” means no prefix is forced during decoding. “UNREL” means the first three tokens are forced to be from another random unrelated sentence (which is wrong but grammatical). “RAND” means the first three tokens are completely random words.

easy task even for human translators, yet the model does fairly well. Again, this contradicts with the “exposure bias” hypothesis that a MLE-trained LM will produce a increasingly deviated sequence when initiated with a non-perfect prefix. Actually, during generation the model self-corrects the error in the prefix. It is also the major motivation of our proposed EB-C measurement (Section 5), which is based on the view of measuring distances between conditional distributions.

B REAL-DATA EXPERIMENTS FOR EB-M

One problem in the implementation of EB-M is to estimate the described marginal distributions of W_{l+1} . We adopt a simple sample-and-count method: $P_{D|D}^{l+1}$ is estimated by the distribution (histogram) of W_{l+1} from a number (to be specified below) of sentences sampled from the data distribution. For $P_{M|M}^{l+1}$ and $P_{M|D}^{l+1}$, we first draw a number of history samples $W_{1:l}$ from the corresponding history model (model distribution and data distribution respectively). We then feed sampled history sequences into the trained model and estimate the marginal distribution of the $(l+1)_{th}$ word by averaging the predicted distribution $P_M(\cdot|W_{1:l})$.

We measure EB-M for MLE-trained LSTM LM on two popular data-sets: EMNLP-news (EMNLP 2017 WMT News Section), and wikitext-103⁷. For EMNLP-news we set $L = 20$, and only use data samples whose length is longer than L . The resulting training/validation/test set has 268k/10k/10k

⁷The wikitext-103 data is available at <https://blog.einstein.ai/the-wikitext-long-term-dependency-language-modeling-dataset/>.

sentences. The vocabulary is of size 5k. We use the 10k samples in the test set for evaluation of EB-M. Note that the EMNLP-news data-set is widely used in text GAN literatures Yu et al. (2016); Lu et al. (2018). We train a one-layer LSTM LM (Sundermeyer et al., 2012) of hidden dimension 512 as the MLE baseline model for EMNLP-news.

For wikitext-103, we set $L = 50$, and regard a paragraph in the original data as a long sentence. Further, we use half of the data for LM training, and utilize the other half for EB-M evaluation. The resulting training/validation/test/evaluation set has 300k/1.5k/1.5k/300k sentences. The vocabulary is of size 50k. We train a two-layer LSTM LM of hidden dimension 1024 as the MLE baseline model for wikitext-103.

For MLE baseline model training, the Adam optimizer is used with learning rate 0.001, no Dropout (Srivastava et al., 2014) is applied. The model is trained for 100 epochs.

We first measure EB-M on the wikitext-103 data-set, which has large amount of evaluation data. The results are shown in Figure 5. We provide EB-M measurements with metric d_{TV} in Appendix E, as they are similar to those using metric d_{JS} . It is shown that the measurements become stable when using 100k data/model samples. EB-M has an average value of 1.10, indicating a significant gap between the model’s MGD when fed with history from P_D or P_M . Further, we observe a steady growth of EB-M along the length of history, which is expected as an outcome of exposure bias.

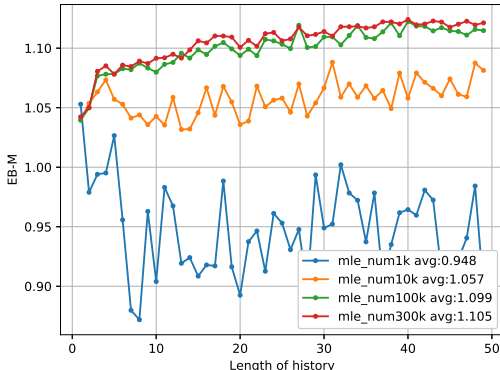


Figure 5: EB-M measurements (with metric d_{JS}) using different number of samples on wikitext-103 data.

However, as discussed in Section 4.2, these measurements can only show that the LM does have better (marginal) generation quality when fed with data prefixes, but does not provide informative information for the significance of exposure bias.

C IMPLEMENTATION OF SS, CoT, AND RANKGAN

We implement our MLE baseline and scheduled sampling (SS) in PyTorch. For SS, we use a linear decay schedule to move from complete teacher forcing to replace-sample rate of 0.1. We find that larger rate will give worse performance.

For CoT, we use a PyTorch implementation in <https://github.com/pclucas14/GansFallingShort>. We use a mediator model that has twice the size of the generator. We set M-step to be 4, and G-step to be 1.

For RankGAN, we use a TensorFlow implementation in <https://github.com/desire2020/RankGAN>.

Note that in our non-MLE experiments, the generator model is set to be the same size with the baseline MLE model. We tune the non-MLE methods using the corpus-BLEU metric, which is widely used in text GAN literature.

You will be given a partial sentence, please simply guess, to your first instinct, what would be the next word continuing the sentence.

For example, for the context "The iPhone is a line of smartphones designed", the next word could be "by" or "for" or "and", etc.

We will give you five choices to choose from, please choose the one you think is most plausible to be the next word.

Some of the sentences are generated by a deep learning model, so they are not perfect, just do your best for a guess.

Please note that we will only accept up to 200 HIT per worker for this task. Please don't do more than 200 HITs.

Type the **index number of choice (not the word itself)** for what you would say in the given situation:

Idx: 5 **Context:** a video advertisement for the game featuring the song " come together " premiered on 28 august 2009 . the

Choices: 0: song 1: music 2: campaign 3: game 4: video 5:[None of the above is plausible]

Just type the index (a number in the set {0,1,2,3,4,5}) of your choice...

Idx: 5 **Context:** cain again depicted his son in a music video so for the song " ave maria ", is shown

Choices: 0: on 1: as 2: in 3: singing 4: with 5:[None of the above is plausible]

Just type the index (a number in the set {0,1,2,3,4,5}) of your choice...

Idx: 6 **Context:** although michael jackson , who owned 50 % of the publishing rights to the beatles songs through sony / atv

Choices: 0: music 1: television 2: records 3: publishing 4: , 5:[None of the above is plausible]

Just type the idx (in the set {0,1,2,3,4,5}) of your choice...

Idx: 6 **Context:** " the power plant was named the " biggest plant of all time throughout the world " (shortly after

Choices: 0: being 1: the 2: " 3: it 4: its 5:[None of the above is plausible]

Just type the idx (in the set {0,1,2,3,4,5}) of your choice...

Figure 6: The HIT interface for our evaluation.

D DETAILS ABOUT THE AMT EVALUATION FOR EB-C

In this section we provide more details for the AMT evaluation discussed in Section 5.3. We show the HIT interface in Figure 6. Each HIT will include 10 pairs of context and its corresponding choices. Five of them are history samples from real data, and the other five is from the trained model. The history samples are mixed, so that the turker doesn't know whether the history sample is from real data or the model. The next-word choices are also shuffled. The history length of the context could be 10, 20, or 30.

We collect around 10k HITs for each history length configuration. The same history sample is not repeated across the HITs. We limit each turker to do at most 200 HITs. For all history length configurations, there are around 300 unique turkers. As shown by Figure 7, most turkers conduct less than 20 HITs.

E AUXILIARY PLOTS

In Figure 8, we show that we are able to get stable measurements of EB-C with 100k samples for the LSTM-512 synthetic experiment.

In Figure 9 and Figure 10 we provide EB-M measurements with metric d_{TV} discussed in Section 4.2, the results are similar to those using metric d_{JS} .

In Figure 11, we provide EB-C measurements of a 3-layer transformer LM (Dai et al., 2019) with 512 hidden dimension, in the synthetic setting.

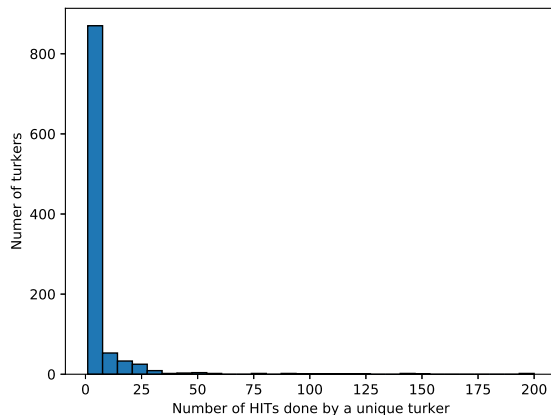


Figure 7: A histogram of the number of HITs done by a unique turker.

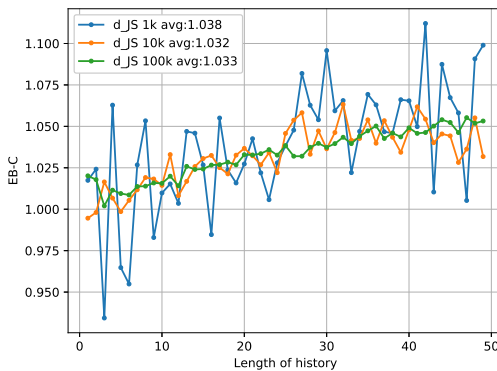


Figure 8: EB-C measurements with different number of samples for the LSTM-512 synthetic experiment.

F PERPLEXITY OF THE TRAINED MODELS

We show PPL results for model trained on EMNLP-news data-set in Table 5. The MLE model for wiki-103 data-set discussed in Section 4.2 has PPL 84.58. Note that due to our special setting⁸, our PPL result is not directly comparable to state-of-art LM results on these data-sets.

⁸We only keep sentences of length longer than L , and for wiki-103, only half of training data is used.

Model	PPL
MLE Baseline $M_{news-512}$	55.85
LSTM-512 (MLE, synthetic)	115.3
LSTM-32 (MLE, synthetic)	156.3
Transformer-512 (MLE, synthetic)	103.6
CoT-512 (synthetic)	115.6
SS-512 (synthetic)	113.7
CoT	56.83
RankGAN	53.43
SS	56.43

Table 5: PPL results for model trained on EMNLP-news data-set.

<p>Data Samples as History → Model Samples</p> <p>He had a recurring role in 2003 on two episodes → of the sitcom Roseanne as James "Bitch" Cook and guest starred in the 1999 special Richard Ayoade's comedy Canal & ...</p> <p>Du Fu <unk> s compassion, for himself and for others, → which arrived at Du Fu soon after the collapse of his political system. He was and still is called "<unk>" ...</p> <p>Du Fu <unk> s work is notable above all for → its use of the convention people as arbiters in decision-making.</p> <p>The tenor of his work changed as he developed his → falsetto; the same February Bach performed his entire Magnificat in the Domus Aurea, a musical hall in Campo Bartolommeo, a ...</p> <p>Although he wrote in all poetic forms, Du Fu explains → that he had no intentions of writing poetry, and attempted to cash in on the success of his two-volume translation ...</p> <p>About two thirds of Du Fu <unk> s 1500 extant → works survive as collections, but about one third have been rebuilt or linked. Some miniatures, such as the Memorial by ...</p> <p>According to the Encyclopdia Britannica, Du Fu <unk> s writings → joined a theme of opposition to social systems on the basis that the United States lacked standards and cotton-beets were ...</p>
<p>Model Samples as History → Model Samples</p> <p>. Competing at the 2006 Commonwealth Games, McBreen scored 3 → . 87 goals per game, ranking fourth from the conference in scoring, beaten 6 3 by Scotland. He also ...</p> <p>He, along with some young Christians from Poland, Romania, and → East Germany, were taught to play dilruba. In order to achieve this, the boy recorded 40 or 50 dilruba parts ...</p> <p>EEF service throughout this filter was to suffer. This was → approximately with the British renewal and capture charges on Mount Ciss, which contributed a large strength taking time to fall ...</p> <p>The matriarchal nature of the family is tested as opposed → to that of their neighbors. In-laws explain their position by having the rear bedroom bathed in bondage to reflect cosmic ...</p> <p>The branch office distributed tuition to the top level schools, → gaining coverage in the art of instruction in schools which allow them to select classes exclusively on the basis of ...</p>
<p>Shuffled Data Samples as History → Model Samples</p> <p>of Below an one of example is Du Fu <unk> → s <unk> Systme <unk> Systme, also the address of the No. 1 monuments Society and Advertising identifies a mass scale ...</p> <p>summarises his He <unk> by that Hung concluding life, let → alone die. He ends by dying as saying "I died on the way". Robert Penson has selected Hung's last words ...</p> <p><unk> top ten @ - @ became track group The → the sixth "Nation We re dedicated to at ten" based on race \$15 @, @ 000 pre-determined event. The show ...</p> <p>An to designed music accompanying group the, video display was → full on bluescreen and was rendered with Ghibli HDTVs. For the Xbox 360's Steam control, the cloud density was increased ...</p> <p>well You by <unk> received Kiss contemporary music was <unk> → while Sobhi Youssef of Sputnikmusic acted as a vocal coach for relation back to the original recordings of "If" and ...</p>
<p>Random Sequences as History → Model Samples</p> <p>...RANDOM... execution love Author → Churches Under Sunset and Angels <unk> post the 20 @, @ 000 To 30 @, @ 000 Arc landscape-crosses around 500 Enix areas ...</p> <p>...RANDOM... beyond spiders annually → as part of regional zoning plans, including a pie canning pool in Mechanicsburg, some <unk> Ellisburg, and boxes of all medical ...</p> <p>...RANDOM... realm unknown healthy-bred → Spock (released in 1991 as The Return), arrives in Sickbay to find a team; he engages in normal conversation (the main ...</p> <p>...RANDOM... rough elections appointment → levels as he had already secured the if no candidate received the season ticket, a result of the September 11 attacks. ...</p> <p>...RANDOM... / horses Finn → s experience of sexual frog foraging, and might pose a threat to sexual preference as the crop Betsimisarakaka earn "<unk> ...</p> <p>...RANDOM... Poland 1963 medium. → Basu was the visual effects supervisor on 300 visual effects shots of Gangster, Feller's seventh appearance in a Bollywood film. Aamir ...</p> <p>...RANDOM... levels MD defending → her city of Beaufort, East Carolina in 2004.</p> <p>At the same time, she responded against the package of short-form compatible boats ...</p>

Table 6: More samples of a STOA MLE-trained transformer LM (on the wiki-103 data-set) when fed with different kinds of history. To save space, we omitted the first 7 words of the random history.

Model Samples as History → Model Samples

it was only a pieces that had gone up to → the forest and forces the shoppers about their chronic young
i mean we didn ' t know what i haven → ' t considered through , ” she told bbc radio
if he were the president - elect , he was → known that he would run a force in business at
but these are not as tired of ” the same → message that the harry actor does have been hours in
first opinion the agent have taken four seconds , or → if they don ' t only know anything , were
” the economy of the uk is low enough of → people of defending where americans think that ” brexit ,
the economy grew on 1 . 6 % since the → us voted , and when it turned around 200 streets
i was able to produce on my own , which → is good ; now that the theatre i ' ve
” i ' ve not buying boys i addressed many → nervous times before , as a teenager made me is
we think about one - third of the struggles we → actually want to see those very well that even more
the story of a album - which made public - → was still fantastic , and for the second time in
” the test comes up before tuesday and when we → ' re feeling ahead again soon , ” she posted
a year on when he was last seen in his → home and he did not see him , his suffering
brady has forced the 9 - known targets to get → all - of - 12 gun migration and performing communication
i asked if he himself did , i managed to → show all my charges at all , it used to

Data Samples as History → Model Samples

what this group does is to take down various different → players in the future and we play in paris we
over 1 , 600 a day have reached greece this → gone in 2013 and it planned to allow civilians on
” we ' re working through a legacy period , → and i am proud of the experience of the worker
’ the first time anyone says you need help , → you don ' t have put accurate press into the
out of those who came last year , 69 per → cent of women can really take the drive to avoid
he has not played for tottenham ' s first team → this season then and sits down 15 - 0 with
so you have this man who seems to represent this → bad story , which he plays minutes – because he
cnn : you made that promise , but it wasn → ' t necessarily at all the features he had in
this is a part of the population that is unk → lucky to have no fault today , and it would
they picked him off three times and kept him out → of the game and was in the field , the
the treatment was going to cost \$ 12 , 000 → as a result of the request of anyone who was
but if black political power is so important , why → doesn ' t we becomes the case that either stands
local media reported the group were not looking to hurt → the animals , but would never be seen to say

Random Sequences as History → Model Samples

...RANDOM... big winter deserve → , but they just say it your things goes wrong
...RANDOM... playoff north realise → at its lowest level , improving their understanding in danger
...RANDOM... vital childhood registration → , not previously planned for junk_i to each and reduced
...RANDOM... treated ship find → one as an actual three points contained at a time
...RANDOM... faith five crazy → schools and could give them a ” sleep ” necessary
...RANDOM... domestic jason follows → a 12 - year cruise line over the christmas track
...RANDOM... ownership generous tourist → accounts for more than 1 per cent every month -
...RANDOM... spending raped since → the file returns in january , joining groups of foreign
...RANDOM... netflix worker four → centre - and said facebook text junk_i to see how
...RANDOM... race labor witnessed → is great , with more to an active the junk_i
...RANDOM... treatments airlines hidden → real - time out to sell on benefits to our
...RANDOM... intention short reflects → showing the nature of flying in his space rather than
...RANDOM... conversation pace motion → them further , but as late as they ' ve
...RANDOM... export feb president → obama agreements with president obama and her being on trump
...RANDOM... entering pocket hill → and made it later in the united states and make

Table 7: Samples of a MLE-trained LSTM LM (on the EMNLP-news data-set) when fed with different kinds of history. To save space, we omitted the first 7 words of the random history.

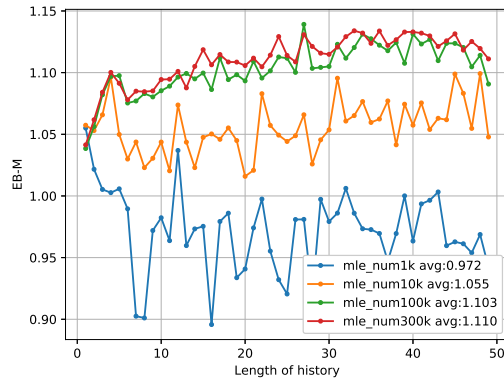


Figure 9: EB-M measurements (with metric d_{TV}) using different number of samples on wiki103 data.

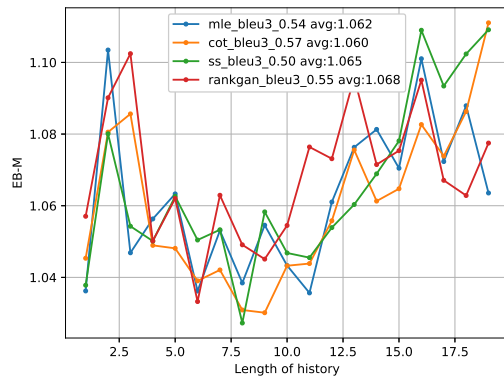


Figure 10: EB-M (with metric d_{TV}) comparison for MLE and non-MLE training on EMNLP-News data. For each training method, we show corpus-BLEU measurement using the test-set as reference set in the legend.

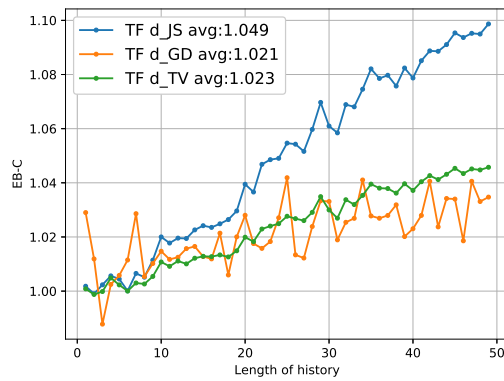


Figure 11: EB-C measures for the transformer LM.