

Softmax Bottleneck Makes Language Models Unable to Represent Multi-mode Word Distributions

Anonymous ACL submission

Abstract

Neural language models (LMs) such as GPT-2 estimate the probability distribution over the next word by a softmax over the vocabulary. The softmax layer produces the distribution based on the dot products of a single hidden state and the embeddings of words in the vocabulary. However, we discover that this single hidden state cannot produce all probability distributions regardless of the LM size or training data size because the single hidden state embedding cannot be close to the embeddings of all the possible next words simultaneously when there are other interfering word embeddings between them. In this work, we demonstrate the importance of this limitation both theoretically and practically. Our work not only deepens our understanding of *softmax bottleneck* and *mixture of softmax* (MoS) but also inspires us to propose *multi-facet softmax* (MFS) to address the limitations of MoS. Extensive empirical analyses confirm our findings and show that against MoS, the proposed MFS achieves two-fold improvements in the perplexity of GPT-2 and BERT.

“The greater the ambiguity, the greater the pleasure.” — Milan Kundera

1 Introduction

Recently, researchers have found that transformer-based language models (LMs), such as GPT-2, can learn to generate better as their sizes grow (Radford et al., 2019; Brown et al., 2020; Kaplan et al., 2020). One natural question arises: Do modern language modeling architectures still have restrictions in their ability to represent the appropriate distribution over next words or masked words?

In this paper, we discover that, when predicting the probabilities of becoming the next word given an ambiguous context, GPT-2 is often incapable of assigning the highest probabilities to the appropriate non-synonym candidates. For example, given

the input prompt “After debating whether to bow to the *woman* or the *king* first, the jester decided on the [MASK]”, we would expect the distribution over the [MASK] fillers to put high probabilities on “*woman*” or “*king*” or their synonyms. However, GPT-2 might incorrectly output the *king* and “*queen*” as the top two candidates as in Figure 1.

In the final softmax layer of GPT-2, the log probabilities of the *woman* and *king* are computed based on the dot product between a single hidden state embedding and the global word embeddings of the *woman* and *king*, respectively. To have the highest but similar dot products for the two options, the transformer encoder in GPT-2 wants to output the hidden state that is close to the average of the *woman* embedding and the *king* embedding. However, the words *queen*, *king*, *woman*, and *man* tend to form a parallelogram in the embedding space (Mikolov et al., 2013; Ethayarajh et al., 2019; Wang et al., 2019), which means the *man* and *queen* also have a similar average. Therefore, GPT-2 are forced to also output the *man* or *queen* when it wants to output the *woman* or *king*.

The problem not only happens to GPT-2 or the words whose embeddings form a parallelogram shape. Even though the hidden state embedding of LMs are contextualized, the embedding of each word in the softmax layer is global and static during the inference time. Globally dissimilar words could all become the suitable next word in a context while other interfering words might be between them, which makes the ideal next word embedding distribution to have multiple modes and cannot be modeled by the single embedding representation.

In this work, we propose theorems showing that given any LM using the output softmax layer, when there are more than N word embeddings in a $N - 1$ dimensional subspace/hyperplane (e.g., 4 embeddings in a two dimensional plane), we can always find a set of possible next words (e.g., *woman* and *king*) such that there are some other interfering

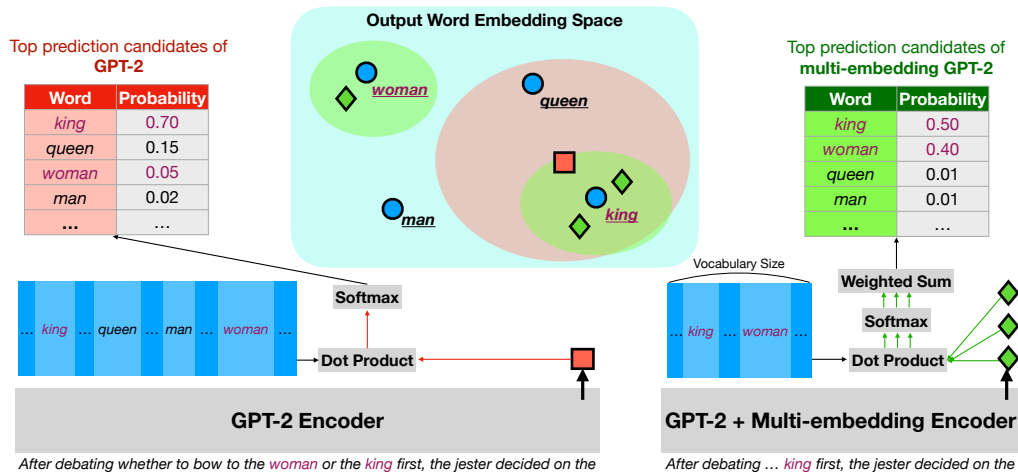


Figure 1: Comparison between the softmax layers using a single embedding and multiple embeddings when the next word should be either *woman* or *king*. In **GPT-2** and **multi-embedding GPT-2**, the hidden states of the context are visualized by the **single facet** \blacksquare and **multiple facets** \blacklozenge , respectively. The word embeddings are visualized using \bullet . GPT-2 cannot output *woman* and *king* as the top two words because *queen* and *man* are close to the middle of *woman* and *king*. The improvement in this type of ambiguous context will be quantified in section 5.

words between them (e.g., *man* or *queen*).

Recently, *mixture of softmax* (MoS) (Yang et al., 2018) regains attention as one of the few effective architecture modifications for transformer LM (Narang et al., 2021; Anonymous, 2021). In the meanwhile, Parthiban et al. (2021) show that the *softmax bottleneck* (Yang et al., 2018) theory is not sufficient to explain the improvement of MoS. Our theorems not only provide geometrical intuitions of why and when the multiple embedding representation such as MoS would do better but also suggest that the *softmax bottleneck* might not be completely solved even if we adopt a very large hidden state size. For example, no matter how large the hidden state size is, as long as $\text{queen} - \text{king} = \text{woman} - \text{man}$ in the embedding space, the LMs cannot output a pair of words in the longer diagonal of the parallelogram as the top two output words.

After better understanding why *mixture of softmax* (MoS) works well, we propose two enhancements over MoS. The first enhancement considers the hidden states of multiple positions and multiple transformer layers when determining the probability in each softmax; the second enhancement uses different contextualized embeddings to compute the probabilities of different subsets of words.

The resulting method, *multi-facet softmax* (MFS), significantly outperforms the MoS and the GPT-2 with the softmax layer on the perplexity for predicting the next word, especially in ambiguous context and non-English text in OpenWeb-

Text (Radford et al., 2019). Finally, we also show that MFS could improve the performance of GPT-2 on ProtoQA (Borotko et al., 2020), a commonsense question answering dataset where each question has multiple acceptable answers.

We summarize our theoretical, methodological, and empirical contributions as follows.

- **Theory:** We show the softmax layer using a single embedding is sometimes not be able to output an appropriate rank of probabilities on a set of words with linearly dependent embeddings.
- **Method:** Addressing two weaknesses in MoS (Yang et al., 2018), we propose *multi-facet softmax* (MFS), a new alternative to the output softmax layer. MFS can replace the softmax in pre-trained LMs to better handle ambiguous contexts without re-training the LMs from scratch.
- **Analysis:** Our comprehensive empirical analyses discover and explain several phenomena, such as a) why using multiple embeddings is usually better than the single embedding with the non-linearity, b) why the improvement is larger in ambiguous contexts, less common languages, or GPT-2 compared to BERT, and c) why increasing hidden state size boosts the capability of distinguishing similar words.

2 Theoretical Limitations of the Single Embedding in the Softmax Layer

In this section, we first review the softmax layer of GPT-2 formally and explain why $\text{queen} - \text{king} =$

woman - man still tends to hold in contextualized LMs. Next, we present our theoretical analyses, which generalize the *woman* and *king* example by showing that the candidate words in a low dimensional subspace would induce the impossibility of ranking some candidates on top of other candidates.

2.1 Background

The LMs typically use a softmax layer to predict $P_S(x|c_t)$, the probability of the next word x given the context at the t th position c_t :

$$P_S(x|c_t) = \frac{\exp(\mathbf{h}_{c_t}^T \mathbf{w}_x)}{\sum_{x'} \exp(\mathbf{h}_{c_t}^T \mathbf{w}_{x'})}, \quad (1)$$

where \mathbf{h}_{c_t} is the t th hidden state in the context c , and \mathbf{w}_x is the output word embedding for the word x (i.e., the linear weights that project the hidden state to the logit of the word x). Yang et al. (2018) point out that the log probability distribution over all the words in the vocabulary V is $\log(P_S(x|c_t))|_{x \in V} = \mathbf{h}_{c_t}^T \mathbf{w}_x - \log(\sum_{x'} \exp(\mathbf{h}_{c_t}^T \mathbf{w}_{x'}))|_{x \in V}$. The distribution is a linear projection from the hidden state \mathbf{h}_{c_t} with dimension D , so the degree of freedom in the distribution is only D (i.e., there cannot be more than D linearly independent log distributions). We call this restriction *softmax bottleneck* theory.

During training, the ideal output word embedding \mathbf{w}_x should be close to the hidden states of the contexts \mathbf{h}_{c_t} that co-occur with the word x while far away from the other hidden states. This objective is similar to the objective function of Word2Vec (Mikolov et al., 2013) except that the context embeddings are contextualized (Kong et al., 2020; Li et al., 2020).

If a context c_t has a higher chance to co-occur with *queen* compared to *king*, the context also has a higher chance to co-occur with *woman* compared to *man* to a similar degree. This is the main reason that makes queen - king = woman - man in the Word2Vec space (Ethayarajh et al., 2019). Therefore, the same linear relations tend to hold in the output word embedding space of GPT-2 as well (Wang et al., 2019).

2.2 Structural Weakness Theorems from Linear Dependency

In addition to words satisfying the analogy relations, the following theorems imply that any linear dependency among the words causes the difficulties of LM in ranking the words in an arbitrary order.

For example, woman + king = queen + man makes a LM unable to output *woman* and *king* as the top two words in Figure 1.

Theorem 1. *If the nonzero output embeddings of N words are linearly dependent and on one side of a plane through the origin, the output softmax layer cannot rank the N words with an arbitrary order according to their probabilities.*

Here, we provide an intuitive justification: if N embeddings are in a subspace whose dimension is smaller than $N - 1$, the N embeddings are going to be linearly dependent and some set of words cannot have the top dot products due to the limited degree of freedom in the subspace. In Appendix D, we formally prove the theorem by identifying the sets of words that cannot be ranked top by the single embedding representation.

In practice, linear dependency holds approximately instead of exactly. For example, woman = queen + man - king + ϵ . In this practical condition, the following theorem states that the logits of the bottom words (i.e., *man* and *queen*) cannot be much smaller than the logits of the top words (i.e., *woman* and *king*).

Theorem 2. *Let the output word embeddings in the set $W = \{\mathbf{w}_i \neq \mathbf{0} | i = 1 \dots N\}$ satisfy $\mathbf{w}_1 = a_2 \mathbf{w}_2 + \dots + a_N \mathbf{w}_N + \epsilon$, where the constant a_2, \dots, a_N are neither all zero nor all negative and $\|\epsilon\| < \epsilon$. Then, there must be a non-trivial partition $P = \{G, S\}$ of W such that there is no hidden state $\|\mathbf{h}\| \leq r$ and a threshold $\tau \geq r\epsilon$ that make $\min_{\mathbf{w}_g \in G} \mathbf{h}^T \mathbf{w}_g \geq (1 + \delta)\tau$ and $\max_{\mathbf{w}_s \in S} \mathbf{h}^T \mathbf{w}_s < \tau$, where $\delta = \frac{2}{1 + \sum_{i=2 \dots N} |a_i|}$.*

Its proof can be found in Appendix D and Appendix B.1 estimates ϵ in GPT-2.

Even though, theoretically-speaking, outputting *woman* and *king* as the top two words is possible due to the appearance of ϵ , GPT-2 may not successfully learn to output the optimal \mathbf{h} and the optimal hidden state for these four words could lead to the wrong probabilities of the other words. Consequently, GPT-2 sometimes still ranks *queen* or *man* higher than *woman* or *king* in practice.

3 Multi-facet Softmax

Using multiple embeddings is a natural solution of modeling a multi-mode distribution. For instance, we can use three embeddings to capture the high probability on the *woman* and *king* but low probability on the *man* and *queen* in Figure 1.

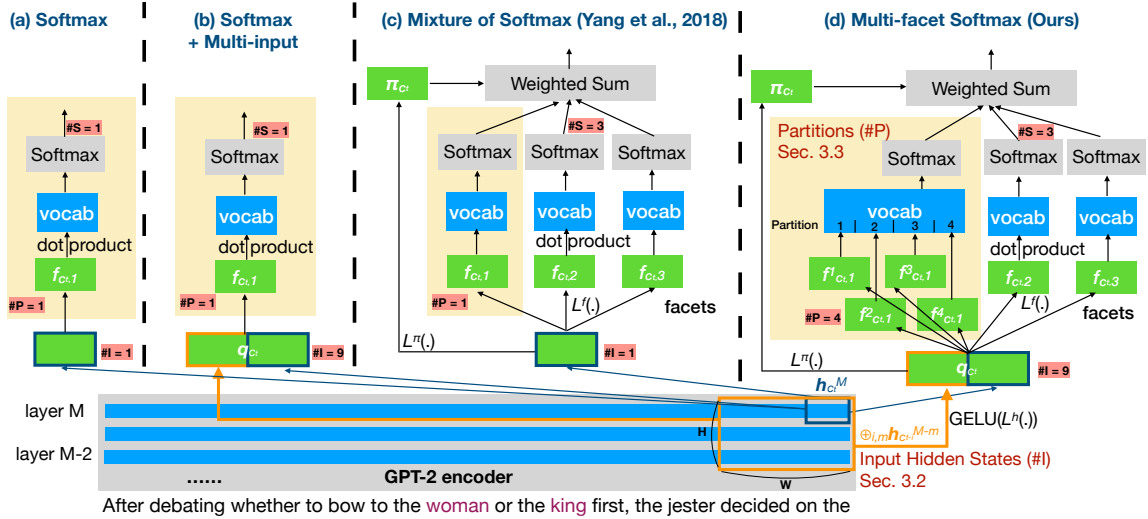


Figure 2: Comparison between different architectures. The $\#S$, $\#I$, and $\#P$ are the number of softmaxes, input hidden states, and partitions, respectively. The green boxes refer to embeddings/vectors. The vocab means the embeddings of all words in the vocabulary. \oplus refers to concatenation. L^h , L^f , and L^π are linear projection layers.

Inspired by our geometric analysis on the limitation of the single embedding, we improve the state-of-the-art multiple embedding solution, *mixture of softmax* (MoS) (Yang et al., 2018) by two enhancements: multiple input hidden states and multiple partitions on the vocabulary.

3.1 Mixture of Softmax

Yang et al. (2018) propose *mixture of softmax* (MoS) to allow a LSTM-based (Hochreiter and Schmidhuber, 1997) LM to produce more linearly independent log probability distributions of the output words given different contexts. As in Figure 2 (c), the MoS first uses multiple linear layers L_k^f to project a hidden state \mathbf{h}_{c_t} into multiple facet embeddings $\mathbf{f}_{c_t,k} = L_k^f(\mathbf{h}_{c_t})$.¹ The multiple facets $\mathbf{f}_{c_t,k}$ and softmaxes would lead to multiple probability distributions, and output probability is the weighted average of the distributions:

$$P_{MoS}(x|c_t) = \sum_{k=1}^K \pi_{c_t,k} \frac{\exp(\mathbf{f}_{c_t,k}^T \mathbf{w}_x)}{\sum_{x'} \exp(\mathbf{f}_{c_t,k}^T \mathbf{w}_{x'})}. \quad (2)$$

The prior weights $\pi_{c_t,k} = \frac{\exp(L_k^\pi(\mathbf{h}_{c_t}))}{\sum_{k'} \exp(L_{k'}^\pi(\mathbf{h}_{c_t}))}$, where L_k^π is another linear projection for dynamically generating the weights and the projection goes through a softmax to ensure $\sum_{k=1}^K \pi_{c_t,k} = 1$.

¹We remove the tanh layer in the original MoS to improve its performance on GPT-2. See Appendix F.1 for details.

3.2 Multiple Input Hidden States

To model the multi-mode distribution, the facets (i.e., the embeddings for different softmaxes) should be able to move freely. For example, in Figure 1, we have three facets but only have two modes, so the two embeddings are very close to the word *king*. However, when we want to output three dissimilar top words such as the *king*, *woman*, and *knight*, one of the facets should be moved to be near to the embedding of the *knight*.

Therefore, we want our solution to satisfy two properties: a) the linear transformation matrix in L_k^f should have a full rank to avoid limiting the degree of freedom in each facet, and b) the relative location of the facets should be context-dependent. MoS cannot satisfy both properties. If the first one is satisfied, the input hidden state is uniquely determined by a facet (e.g., $\mathbf{h}_{c_t} = (L_1^f)^{-1}(\mathbf{f}_{c_t,1})$). Then, there exist a global transformation between two facets (e.g., $\mathbf{f}_{c_t,2} = L_2^f \left((L_1^f)^{-1}(\mathbf{f}_{c_t,1}) \right)$), which violates the second property. In other words, since the facet embeddings are the projection of a single hidden state, the total degree of freedom in all facet embeddings cannot exceed the dimension of the hidden state.

Our solution to this issue is using more input hidden states to construct the facets. As the orange box in Figure 2, we first concatenate a $W \times H$ block of input hidden states into $\oplus_{i=0 \dots W-1, m=1 \dots H} \mathbf{h}_{c_t-i}^{M-m}$, where $M - m$ is the Transformer layer index and $t - i$ is the index of the i th to the last word in the

context. The $W \times H$ is fixed as 3×3 in this paper. We make its dimension the same as the original hidden state $\mathbf{h}_{c_t}^M$ using a linear layer L^h plus a GELU activation function (Hendrycks and Gimpel, 2016). Then, we concatenate it with the original hidden state to form a new input hidden state

$$\mathbf{q}_{c_t} = \mathbf{h}_{c_t}^M \oplus \text{GELU} \left(L^h \left(\bigoplus_{i,m} \mathbf{h}_{c_{t-i}}^{M-m} \right) \right). \quad (3)$$

The new input hidden state is passed through the linear transformation L_k^f to compute the facets $\mathbf{f}_{c_t,k} = L_k^f(\mathbf{q}_{c_t})$ and our prior weights $\pi_{c_t,k} = \frac{\exp(L_k^\pi(\mathbf{q}_{c_t}))}{\sum_{k'} \exp(L_{k'}^\pi(\mathbf{q}_{c_t}))}$. Since the dimension of \mathbf{q}_{c_t} is larger than the dimension of $\mathbf{f}_{c_t,k}$, the inverse function $(L_k^f)^{-1}$ no longer exists.

3.3 Multiple Partitions

The next word distribution could have many modes. However, using many softmaxes significantly increases our computational burden because we need to compute the dot product between each facet and all the word embeddings in our vocabulary.

Inspired by our analysis, we propose to split all the words in the vocabulary into multiple partitions and use different facets for different partitions. For example, if we can put any word from $\{\textit{queen}, \textit{man}, \textit{woman}, \textit{king}\}$ into one partition and the rest of the words into another partition, we no longer have *queen* - *king* = *woman* - *man* in either of the partitions. In this method, each word only belongs to one partition, so we only need to compute one dot product for each word. Thus, the extra computational cost only comes from the extra linear projections for preparing the facets.

In many contexts c_t , the distribution of the next word has only a single mode and the global similarity between words may be useful. Using the multiple partitions alone might lose the similarity information between words in different partitions. Therefore, we propose to only replace the first softmax layer in MoS with the multiple partition method to learn the global similarity of words in different partitions using the other softmaxes. The architecture is illustrated in Figure 2 (d). Formally, we compute the probability using

$$P_{MP}(x|c_t) = \pi_{c_t,1} \frac{\exp((\mathbf{f}_{c_t,1}^{j_x})^T \mathbf{w}_x)}{\sum_{x'} \exp((\mathbf{f}_{c_t,1}^{j_{x'}})^T \mathbf{w}_{x'})} + \sum_{k=2}^K \pi_{c_t,k} \frac{\exp(\mathbf{f}_{c_t,k}^T \mathbf{w}_x)}{\sum_{x'} \exp(\mathbf{f}_{c_t,k}^T \mathbf{w}_{x'})}, \quad (4)$$

where j_x is the partition index that the word x belongs to and $\mathbf{f}_{c_t,1}^{j_x}$ is the facet for the j_x th partition. *Multi-facet softmax* (MFS) is equipped with multiple input hidden states and multiple partitions.

4 Language Modeling Experiments

We evaluate different LM architectures by comparing their capability of predicting the next word in Wikipedia 2021 and a subset of OpenWebText (Radford et al., 2019). In addition to perplexity, we also compare their mean reciprocal ranks (MRR) in Appendix C.1. The size of the training, validation, and testing set are 96%, 2%, and 2% of the whole corpus. After loading the pre-trained GPT-2 models, we train the *GPT-2 Small* for 1 epoch and *GPT-2 Medium* for 0.4 epochs. We also test our methods on BERT in Appendix B.2.

4.1 Baselines

We set different numbers of softmaxes, input hidden states, and partitions in our MFS framework to construct our baselines. The configuration of different baselines could be seen in Table 1.

Softmax (GPT-2): Using a single softmax, input hidden state, and partition as in Figure 2 (a) and Equation 1. The baseline is the same as the original GPT-2 except that we add one more linear layer that converts the hidden state $\mathbf{h}_{c_t}^M$ to the facet embedding $\mathbf{f}_{c_t,1}$ as in other methods.

SigSoftmax (Kanai et al., 2018): The same as *Softmax* except when predicting the next word, Kanai et al. (2018) add some non-linearity into the softmax layer by multiplying the exponent and sigmoid of the logits.

Softmax + Multi-input: Letting *Softmax* access multiple input hidden states as in Figure 2 (b) and Equation 3. The method is similar to Tenney et al. (2019); Fan et al. (2020), and Tay et al. (2021).

MoS (Yang et al., 2018): *MoS (3)* is the *mixture of softmax* with 3 facets/softmaxes, whose probability comes from Equation 2. We also run the MoS with 4 softmaxes in *GPT-2 Small* and call the model *MoS (4)*.

DOC (Takase et al., 2018): Similar to our enhancement using multiple input hidden states, *direct output connection* (DOC) makes each of their facets coming from a different input hidden state.

Other configurations include **Softmax + Multi-partition**, which adds four partitions into the softmax, **MFS – Multi-partition**, which uses only one partition in *MFS* and could also be viewed as *MoS*

Table 1: Perplexity comparison between MFS (Ours) and baselines. #S, #I, #P are the number of softmaxes (i.e., K), input hidden states, and partitions, respectively. The top four baselines use a single softmax. OWT and Wiki are the test set perplexity of OpenWebText and Wikipedia 2021, respectively. The standard errors of all models are smaller than 0.02 perplexity. We also compare the number of parameters and the inference time on one batch.

Models ↓	Configuration			GPT-2 Small				GPT-2 Medium			
	#S	#I	#P	Size	Time	OWT	Wiki	Size	Time	OWT	Wiki
Softmax (GPT-2)	1	1	1	125.0M	84ms	18.72	24.06	355.9M	212ms	15.89	20.34
SigSoftmax (Kanai et al., 2018)	1	1	1	125.0M	91ms	18.63	24.06	355.9M	221ms	16.07	20.65
Softmax + Multi-input	1	9	1	130.9M	87ms	18.50	23.89	366.4M	219ms	15.76	20.29
Softmax + Multi-partition	1	1	4	126.8M	88ms	18.77	24.08	359.0M	218ms	15.89	20.30
MoS (Yang et al., 2018) (4)	4	1	1	126.2M	152ms	18.61	23.77	359.0M	299ms	15.75	20.08
MoS (Yang et al., 2018) (3)	3	1	1	126.2M	130ms	18.63	23.81	358.0M	270ms	15.79	20.11
DOC (Takase et al., 2018)	3	3	1	126.2M	130ms	18.69	24.02	358.0M	270ms	15.88	20.34
MFS – Multi-partition	3	9	1	133.4M	133ms	18.37	23.56	370.6M	276ms	15.65	20.06
MFS – Multi-input	3	1	4	128.0M	134ms	18.60	23.72	361.1M	275ms	15.71	20.08
MFS (Ours)	3	9	4	136.8M	138ms	18.29	23.45	376.9M	283ms	15.64	20.02

Table 2: Perplexity of the *GPT-2 Small* in OpenWebText. The percentages of the perplexity reduction compared to Softmax are presented in the parentheses.

Ratio in Corpus →	Non-English 14%	English 86%
Softmax	13.50 (0.0%)	19.23 (0.0%)
MoS (Yang et al., 2018) (3)	13.19 (2.3%)	19.16 (0.4%)
MFS – Multi-partition	12.98 (3.8%)	18.91 (1.7%)
MFS (Ours)	12.83 (5.0%)	18.83 (2.1%)

+ *Multi-input*, and **MFS – Multi-input**, which uses only one input hidden state to generate all facets.

4.2 Results

Table 1 shows that applying **MFS** to *GPT-2 Small* achieves more than 15% of the perplexity improvement between *GPT-2 Small* and *GPT-2 Medium*, while only increases 5% of their size differences. Except **Softmax + Multi-partition**, adding multiple input hidden states or partitions in different configurations significantly boost the performances. In Appendix B.3, we further show that the improvement of **MFS** over **Softmax** could even become 3-5 times larger in top 5-10% the most ambiguous contexts compared to the rest of the contexts, which suggest that some improvements indeed come from successfully modeling multi-mode distribution.

MFS usually doubles the perplexity improvements between **MoS (3)** and **Softmax** but the running time of **MFS** remains similar to **MoS (3)** because **MFS** only needs a few more linear layers, which is more efficient than adding one more softmax as in **MoS (4)**. **DOC** is worse than **MoS (3)**. This may be due to a starvation problem: the facet from the last hidden state $h_{c_t}^M$ has the prior probability close to 1 and receives most of the gradients. Finally, compared with **Softmax**, the mixed results

in **SigSoftmax** suggest that adding non-linearity into the softmax layer without modeling the multi-mode distribution might not always improve the models (Parthiban et al., 2021).

In Table 3, we present three contexts from the validation set of different datasets and compare the top three predictions of **MFS** and **Softmax** on *GPT-2 Small*. In OpenWebText and Wikipedia 2021, we can see that **Softmax** misses the correct answer in its top three predictions.

OpenWebText is mostly composed of English text, but some non-English text in the corpus allows us to compare the capability of different models in a multi-lingual setting. Table 2 shows that multiple embeddings improve the perplexity of the non-English text more than the perplexity of the English text. We hypothesize that the distribution of the next non-English word is more likely to be multi-mode because GPT-2 learns the global token embeddings mostly in the English contexts, which could make the embeddings of similar tokens in non-English contexts far away.

5 Evaluation on Ambiguous Templates

We synthesize a dataset using templates (Ribeiro et al., 2020) to verify whether the softmax layer in the original GPT-2 really has difficulty in learning to output the bimodal distribution in Figure 1 and whether the multiple embedding methods could overcome the problem. First, we collect the four words with semantic analogy relations in Google analogy dataset (Mikolov et al., 2013). Next, we insert two out of the four words into our manually written templates to form the contexts such as the ones in the last column of Table 3. The templates we used could be found in Appendix F.3. The

Table 3: Prediction visualization using a context in each dataset. We show the top three words with the highest prediction probabilities of each method. In the last three rows, we visualize the outputs of the softmax grey boxes in Figure 2 (d), which model different modes of the next word distribution. The prediction target is boldfaced in the context and the predictions. ## indicates there is no space before the word.

Corpus →	OpenWebText	Wikipedia 2021	Analogy in Templates (section 5)
Input Context	... The Elastic Endpoint Security and Elastic SIEM solutions mentioned in this post are now referred to as Elastic	... law and chance working together cannot generate CSI, either. Moreover, he claims that CSI	I went to Paris and Germany before, and I love one of the places more, which is Germany
Softmax (GPT-2)	the 0.087, E 0.043, End 0.039	the 0.174, this 0.054, if 0.038	Paris 0.893, France 0.045, Germany 0.033
MFS (Ours)	Elastic 0.220, the 0.089, EC 0.033	CSI 0.186, the 0.140, there 0.033	Paris 0.544, Germany 0.389, France 0.064
MFS Softmax 1	end 0.051, the 0.043, security 0.023	the 0.191, law 0.127, if 0.053	Paris 0.979, France 0.013, Germany 0.007
MFS Softmax 2	Elastic 0.652, EC 0.080, ES 0.046	the 0.191, there 0.049, this 0.047	Paris 1.000 Berlin 0.000 ##Paris 0.000
MFS Softmax 3	the 0.193, E 0.040, a 0.014	CSI 0.677, law 0.029, laws 0.019	Germany 0.852, France 0.139, China 0.004

Table 4: Perplexity comparison of different *GPT-2 Small* models on the words with different types of analogy relations. The validation set (valid) includes all four types of relations.

Analogy Relation Types → Models ↓	Diagonal (e.g., <i>king</i> or <i>woman</i>)					Edge (e.g., <i>king</i> or <i>queen</i>)				
	valid	capital-common	capital-world	city-in-state	family	valid	capital-common	capital-world	city-in-state	family
Softmax (GPT-2)	2.30	3.30	2.00	2.25	2.95	2.11	2.42	1.91	2.26	2.38
MoS (Yang et al., 2018) (3)	1.75	2.18	1.60	1.85	2.82	1.87	2.26	1.70	2.04	2.27
MFS – Multi-partition	1.72	2.13	1.59	1.82	2.52	1.84	2.23	1.72	1.96	2.16
MFS (Ours)	1.74	2.15	1.59	1.82	2.63	1.92	2.28	1.78	2.00	2.24

two words can be either the diagonal words (e.g., *king* and *woman*) or the edge word (e.g., *king* and *queen*) in the parallelogram. Finally, we create a dataset with 122k training contexts, 250k validation contexts, and 122k testing contexts, where the word pairs in the testing set are unseen in the training set to see whether the model could learn to output the bimodal distribution in a general way.²

We load the models pre-trained on OpenWebText and continue fine-tuning the models on the last word of each sentence for 10 epochs. We report the testing performances of the best model selected by the validation loss. Since the sets of the word pairs in the training and testing set are disjoint, updating the output word embedding would make GPT-2 solve the task by memorizing/overfitting the training set quickly and lead to much worse testing performances. Thus, we freeze the output word embedding during the training.

Table 4 indicates that when the possible next words are the diagonal words, the **Softmax** model performs much worse compared to other multiple embedding alternatives. In the edge word dataset, the multiple embedding solutions are still better but have a much smaller gap. **MFS – Multi-partition** slightly improves **MoS**. We hypothesize the reason

²The setting is realistic because any related words could become the next word in some ambiguous contexts and all the words are related in a certain way (Sigman and Cecchi, 2002). We cannot expect the training corpora to contain the ambiguous contexts with so many possible next words.

is that multiple input hidden states could help the facets to be moved more freely. Finally, multiple partitions seem to cause slight overfitting in this bimodal distribution prediction task.

We visualize the predictions in the last column of Table 3. We can see two of the softmaxes are close to *Pairs* and the remaining one is close to *German*, while **Softmax** overestimates the probability of *Paris* and ranks *France* higher than the *German*. The result verifies that the correct probability distribution of the words in some ambiguous context is hard to learn using **Softmax**.

6 Answering Ambiguous Questions

ProtoQA (Boratko et al., 2020) is a question answering dataset built for evaluating the common-sense reasoning ability of language models. Each question in ProtoQA is ambiguous and leads to a distribution of possible answers. For instance, the answer of “Name something that people usually do before they leave for work?” is “Shower 0.43, Breakfast 0.30, ...”. The paper discovers that by reformulating the question answering task as a context (e.g., “One thing people usually do before they leave for work is ...”), GPT-2 could generate the possible answers by sampling the next words from its word prediction distribution.

The dataset gives us a chance to directly compare the quality of the distributions generated by different LMs in Table 5. After pretraining *GPT-2 Medium* on the OpenWebText, we fine-tune them

Table 5: ProtoQA performances. All the numbers except perplexity are the percentages of the predictions that match the ground truth exactly on the crowdsourced development set. Max answers top k implies only evaluating the top k answers. Max incorrect top k indicates only evaluating the top answers that contain k errors. The best average performances are highlighted and the standard errors are reported as the confidence interval.

Models ↓	Perplexity on Scraped Development Set	Max Answers				Max Incorrect		
		Top 1	Top 3	Top 5	Top 10	Top 1	Top 3	Top 5
Softmax (GPT-2)	1.5432 ± 0.0003	34.1 ± 0.8	35.2 ± 0.5	37.8 ± 0.4	45.0 ± 0.5	18.3 ± 0.4	30.7 ± 0.5	38.5 ± 0.6
MoS (Yang et al., 2018) (3)	1.5407 ± 0.0004	33.9 ± 0.8	36.0 ± 0.6	37.7 ± 0.6	44.9 ± 0.4	18.3 ± 0.4	31.7 ± 0.6	38.2 ± 0.6
MFS – Multi-partition	1.5411 ± 0.0003	34.3 ± 0.7	36.7 ± 0.7	38.1 ± 0.5	45.2 ± 0.4	19.4 ± 0.4	32.0 ± 0.5	38.6 ± 0.3
MFS (Ours)	1.5402 ± 0.0005	34.1 ± 0.6	36.7 ± 0.5	38.6 ± 0.4	45.4 ± 0.5	19.7 ± 0.4	32.1 ± 0.4	39.7 ± 0.4

using the training data in ProtoQA for 2 epochs. We repeat the fine-tuning 5 times and compare their average perplexity in our validation set. Next, we generate 150 sentences starting from each context and compare the generated answers with the ground truth distribution. For each fine-tuned model, we repeat the generation evaluation 3 times and report the average accuracy of the resulting 15 trials.

We can see that the multiple softmaxes, input hidden states, and partitions usually improve the quality of prediction distribution, and the proposed **MFS**, which combines all modifications, achieves the best performances.

7 Related Work

Yang et al. (2018) propose the concept of *softmax bottleneck*, which points out that the dot product in the softmax layer restricts the representation power of outputting arbitrary conditional probabilities. It also proposes *MoS* to break the softmax bottleneck in an RNN-based LM. Kanai et al. (2018) and Ganea et al. (2019) add nonlinearities into the softmax layer to break the bottleneck more efficiently, but the approaches gain less improvement compared to *MoS*.

A limitation of the aforementioned previous work is that they do not tell us which kinds of sentences would be affected by the bottleneck more and whether the order of the top few next words would be affected, which are the main research questions of our work. Contrary to the previous belief that a large hidden state dimension would eliminate the softmax bottleneck, our theorems suggest that some words in a low dimensional subspace could still make the single embedding in the softmax layer become a bottleneck of arbitrarily ranking the output words. Furthermore, our geometric analyses provide an intuitive explanation about why breaking the bottleneck using multiple embeddings leads to better performances compared to only adding the non-linearity.

Demeter et al. (2020) also analyze the structural weakness of the softmax layer from a geometric perspective. They discover that the words with high prior frequencies could stop the LMs from assigning the high probabilities to rare words. The weakness is different from the softmax bottleneck investigated in this paper. Our work shows that the softmax layer could still prevent the LMs from outputting some top words even if all the possible next words have the same prior frequency.

An alternative to model the multi-mode distribution is to use multiple embeddings to represent each output word (Miao et al., 2019). Compared to *MoS* or our approach that use multiple embeddings to represent each hidden state of the context, their method requires many extra parameters to store different senses of each output word. Another type of related model (Shazeer et al., 2017; Fedus et al., 2021) dynamically routes the signals to different experts (i.e., feed-forward networks) and aggregates their outputs. The methodology is similar to *MoS* and our approach, but they add the mixture-of-experts layer inside each layer of the Transformer encoder while the proposed *MFS* is an alternative to the output softmax layer.

8 Conclusion

When the ideal distribution in the output word embedding space has multiple modes, GPT-2 cannot learn to correctly rank the words in all the modes as the top next words. This shows that the single embedding in the softmax layer, which is used nearly universally by current LMs, constitutes a performance upper bound of predicting the next/masked word. To address the systematic failure caused by these structural weaknesses, we propose *multi-facet softmax* (*MFS*). In our experiments, we confirm that the *MFS* significantly outperforms the standard softmax layer and alleviates the *softmax bottleneck* in the transformer-based LMs such as GPT-2 better than *mixture of softmax* (*MoS*).

9 Ethical and Broader Impact

This work studies a general limitation of LMs and proposes solutions. The proposed theory can help us to understand that some types of hallucinations, mistakes, or biases of LMs could come from *softmax bottleneck* and their incapability of modeling the correct distribution. For example, there are 60% of male characters and 40% of female characters in our training corpus. The language generation model might be forced to assign more than 60% probability to male characters as being much more likely to output *king* than *woman* in Figure 1.

Recently, Narang et al. (2021); Anonymous (2021) show that MoS is one of the few architecture modifications of transformer-based LM that can provide consistent improvements in downstream applications. Our work provides a fundamental reason why the multiple embedding representation is better, which could inspire more future studies that propose a better multiple-embedding architecture to improve LMs (e.g., multi-lingual BERT) or downstream applications. As examples, we list several possible future directions in Appendix G.

Finally, a better LM could lead to both positive and negative societal impacts, but they are not the focus of this paper. Generally speaking, this paper deepens our understanding of the weaknesses of modern LMs and we believe the knowledge can help us to design a better LM that increases the positive impacts and reduces the negative impacts in the future.

References

Anonymous. 2021. [Scaling laws vs model architectures: How does inductive bias influence scaling? an extensive empirical study on language tasks](#). In *ACL ARR Blind Submission*.

Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. 2016. Layer normalization. *arXiv preprint arXiv:1607.06450*.

Michael Boratko, Xiang Li, Tim O’Gorman, Rajarshi Das, Dan Le, and Andrew McCallum. 2020. [ProtoQA: A question answering dataset for prototypical common-sense reasoning](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1122–1136, Online. Association for Computational Linguistics.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child,

Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.

Xingyu Cai, Jiayi Huang, Yuchen Bian, and Kenneth Church. 2021. [Isotropy in the contextual embedding space: Clusters and manifolds](#). In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net.

Hyung Won Chung, Thibault Févry, Henry Tsai, Melvin Johnson, and Sebastian Ruder. 2021. [Rethinking embedding coupling in pre-trained language models](#). In *International Conference on Learning Representations, ICLR 2021*.

David Demeter, Gregory Kimmel, and Doug Downey. 2020. [Stolen probability: A structural weakness of neural language models](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2191–2197, Online. Association for Computational Linguistics.

Kawin Ethayarajh, David Duvenaud, and Graeme Hirst. 2019. [Towards understanding linear word analogies](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3253–3262, Florence, Italy. Association for Computational Linguistics.

Angela Fan, Thibaut Lavril, Edouard Grave, Armand Joulin, and Sainbayar Sukhbaatar. 2020. [Addressing some limitations of transformers with feedback memory](#). *arXiv preprint arXiv:2002.09402*.

William Fedus, Barret Zoph, and Noam Shazeer. 2021. [Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity](#). *arXiv preprint arXiv:2101.03961*.

Octavian Ganea, Sylvain Gelly, Gary Bécigneul, and Aliaksei Severyn. 2019. [Breaking the softmax bottleneck via learnable monotonic pointwise nonlinearities](#). In *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pages 2073–2082. PMLR.

Jun Gao, Di He, Xu Tan, Tao Qin, Liwei Wang, and Tie-Yan Liu. 2019. [Representation degeneration problem in training natural language generation models](#). In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net.

689	Dan Hendrycks and Kevin Gimpel. 2016. Gaussian error linear units (gelus). <i>arXiv preprint arXiv:1606.08415</i> .	<i>Lake Tahoe, Nevada, United States</i> , pages 3111–3119.	745 746
692	Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. <i>Neural computation</i> , 9(8):1735–1780.	George A Miller. 1995. Wordnet: a lexical database for english. <i>Communications of the ACM</i> , 38(11):39–41.	747 748 749
695	Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2020. The curious case of neural text degeneration . In <i>8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020</i> . OpenReview.net.	Sharan Narang, Hyung Won Chung, Yi Tay, William Fedus, Thibault Fevry, Michael Matena, Karishma Malkan, Noah Fiedel, Noam Shazeer, Zhenzhong Lan, et al. 2021. Do transformer modifications transfer across implementations and applications? <i>arXiv preprint arXiv:2102.11972</i> .	750 751 752 753 754 755
700	Sekitoshi Kanai, Yasuhiro Fujiwara, Yuki Yamanaka, and Shuichi Adachi. 2018. Sigsoftmax: Reanalysis of the softmax bottleneck . In <i>Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada</i> , pages 284–294.	Dwarak Govind Parthiban, Yongyi Mao, and Diana Inkpen. 2021. On the softmax bottleneck of recurrent language models . In <i>Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2021, Thirty-Third Conference on Innovative Applications of Artificial Intelligence, IAAI 2021, The Eleventh Symposium on Educational Advances in Artificial Intelligence, EAAI 2021, Virtual Event, February 2-9, 2021</i> , pages 13640–13647. AAAI Press.	756 757 758 759 760 761 762 763 764
707	Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. Scaling laws for neural language models. <i>arXiv preprint arXiv:2001.08361</i> .	Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners.	765 766 767
712	Lingpeng Kong, Cyprien de Masson d’Autume, Lei Yu, Wang Ling, Zihang Dai, and Dani Yogatama. 2020. A mutual information maximization perspective of language representation learning . In <i>8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020</i> . OpenReview.net.	Sara Rajaei and Mohammad Taher Pilehvar. 2021. A cluster-based approach for improving isotropy in contextual embedding space . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 2: Short Papers), Virtual Event, August 1-6, 2021</i> , pages 575–584. Association for Computational Linguistics.	768 769 770 771 772 773 774 775 776
719	Bohan Li, Hao Zhou, Junxian He, Mingxuan Wang, Yiming Yang, and Lei Li. 2020. On the sentence embeddings from pre-trained language models . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 9119–9130, Online. Association for Computational Linguistics.	Marco Tulio Ribeiro, Tongshuang Wu, Carlos Guestrin, and Sameer Singh. 2020. Beyond accuracy: Behavioral testing of NLP models with CheckList . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 4902–4912, Online. Association for Computational Linguistics.	777 778 779 780 781 782 783
726	Ilya Loshchilov and Frank Hutter. 2019. Decoupled weight decay regularization . In <i>7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019</i> . OpenReview.net.	Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc V. Le, Geoffrey E. Hinton, and Jeff Dean. 2017. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer . In <i>5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings</i> . OpenReview.net.	784 785 786 787 788 789 790 791
731	Ning Miao, Hao Zhou, Chengqi Zhao, Wenxian Shi, and Lei Li. 2019. Kernelized bayesian softmax for text generation . In <i>Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada</i> , pages 12487–12497.	Vered Shwartz, Enrico Santus, and Dominik Schlechtweg. 2017. Hypernyms under siege: Linguistically-motivated artillery for hypernymy detection . In <i>Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers</i> , pages 65–75, Valencia, Spain. Association for Computational Linguistics.	792 793 794 795 796 797 798 799
738	Tomás Mikolov, Ilya Sutskever, Kai Chen, Gregory S. Corrado, and Jeffrey Dean. 2013. Distributed representations of words and phrases and their compositionality . In <i>Advances in Neural Information Processing Systems 26: 27th Annual Conference on Neural Information Processing Systems 2013. Proceedings of a meeting held December 5-8, 2013,</i>		

800 Mariano Sigman and Guillermo A Cecchi. 2002. *Global organization of the wordnet lexicon. Proceedings of the National Academy of Sciences*, 99(3):1742–1747. 857

801 858

802 859

803 860

804 Sho Takase, Jun Suzuki, and Masaaki Nagata. 2018. *Direct output connection for a high-rank language model. In Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4599–4609, Brussels, Belgium. Association for Computational Linguistics. 861

805 862

806 863

807 864

808 865

809 866

810 Yi Tay, Mostafa Dehghani, Vamsi Aribandi, Jai Gupta, Philip Pham, Zhen Qin, Dara Bahri, Da-Cheng Juan, and Donald Metzler. 2021. Omninet: Omnidirectional representations from transformers. *arXiv preprint arXiv:2103.01075*. 867

811 868

812 869

813 870

814 871

815 Ian Tenney, Patrick Xia, Berlin Chen, Alex Wang, Adam Poliak, R. Thomas McCoy, Najoung Kim, Benjamin Van Durme, Samuel R. Bowman, Dipanjan Das, and Ellie Pavlick. 2019. *What do you learn from context? probing for sentence structure in contextualized word representations. In 7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net. 872

816 873

817 874

818 875

819 876

820 877

821 878

822 879

823 880

824 Yile Wang, Leyang Cui, and Yue Zhang. 2019. How can bert help lexical semantics tasks? *arXiv preprint arXiv:1911.02929*. 881

825 882

826 883

827 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. *Transformers: State-of-the-art natural language processing. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics. 884

828 885

829 886

830 887

831 888

832 889

833 890

834 891

835 892

836 893

837 894

838 895

839 Zhilin Yang, Zihang Dai, Ruslan Salakhutdinov, and William W. Cohen. 2018. *Breaking the softmax bottleneck: A high-rank RNN language model. In 6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net. 896

840 897

841 898

842 899

843 900

844 901

845 902

846 Zhilin Yang, Thang Luong, Ruslan Salakhutdinov, and Quoc V. Le. 2019. *Mixtape: Breaking the softmax bottleneck efficiently. In Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 15922–15930. 903

847 904

848 905

849 906

850 907

851 908

852 909

853 Avishai Zagoury, Einat Minkov, Idan Szpektor, and William W. Cohen. 2021. *What’s the best place for an AI conference, vancouver or _____: Why completing comparative questions is difficult. In* 910

854 911

855 912

856 913

A Appendix Overview

To demonstrate the wide applicability of our approaches, we conduct more experiments such as applying MFS to BERT in Appendix B. We also show more results and conduct more analyses in Appendix C to further support our conclusions. Next, we provide technical details including the proof of our theorems in Appendix D, the method details in Appendix E, and the experiment details in Appendix F. Finally, in Appendix G, we list several directions that could be further studied in the future.

B More Experiments

We conduct the following five extra experiments to measure the linear dependency among word embeddings in LMs, extend our multi-facet approaches to BERT, confirm the source of the improvement comes from modeling multi-mode distribution, and extend our synthetic experiments to include the output candidate words that have various types of relations and to include the template that favors the single embedding representation.

B.1 Linear Dependency among Words

As we demonstrate in our theorems, the linear dependency in the word embedding imposes a performance upper bound on LMs. In this experiment, we would like to explore whether the word embeddings are still linearly dependent in a larger LM. Besides, Theorem 2 shows that when N words are linearly dependent after moving one of the embeddings with ϵ distance, the LM with the output softmax layer cannot output a large logit margin between two subsets of the N words. We also want to know how small ϵ typically are in the pretrained word embedding and compare the ϵ from different subsets or different LMs.

We randomly select N words in *GPT-2 Small* and *GPT-2 XL* and use the minimal eigenvalue of the matrix formed by their N word embeddings to estimate the ϵ value.³ The top 2 curves in Figure 3 depict the average of minimal eigenvalues from 1,000 sampled N word sets. As we expect, the eigenvalues decrease as N increases (i.e., easier to become linear dependent in a bigger subset of words). As the hidden state size grows from 768

³We normalize all the word embeddings by the average of their magnitudes to fairly compare the distances in *GPT-2 Small* and *GPT-2 XL*.

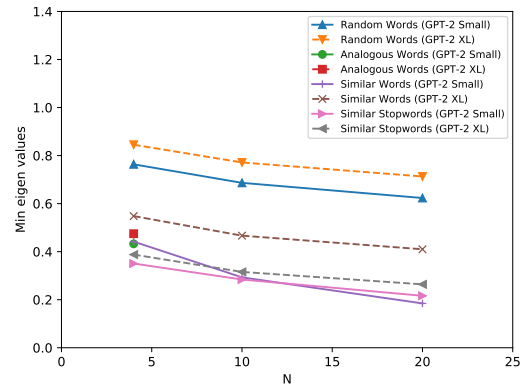


Figure 3: Minimal eigenvalues to indicate the linear dependency among different groups of N word embeddings in *GPT-2 Small* and *GPT-2 XL*.

in *GPT-2 Small* to 1,600 in *GPT-2 XL*, the minimal eigenvalues increase.

As in section 5, we select the 4 words with analogy relation from Google analogy dataset and plot the minimal eigenvalue (averaged over all 4 word sets with an analogy relation) in Figure 3. We can see that the values become much lower and the value of *GPT-2 XL* is only slightly higher than the value of *GPT-2 Small*, which shows that the analogous words still tend to have linearly dependent embeddings in a large LM.

Finally, we select the N similar words by finding the nearest $N - 1$ words of each query word. We exclude the query word pieces whose first character is not a space, and let the query word pieces be either all the rest of the word pieces or only stop words. Figure 3 shows that the minimal eigenvalues are close to the values of analogous words.

Intuitively speaking, the similarly low minimal eigenvalues mean that globally similar words tend to have similar probabilities in every context. Our Theorem 2 formally describes a structural weakness of the output softmax layer in terms of distinguishing the similar words. The low minimal eigenvalues and our theory support the recent empirical finds that LM models tend to be confused by the similar words (Zagoury et al., 2021) and further suggest that the problem is more obvious especially when the size of LM is small, the number of possible next word N is large, or the next word candidates include stop words. This provides a potential explanation why the candidates often include stop words when multiple embeddings outperform the single embedding in Table 3 and Ta-

Table 6: Perplexity of models building on BERT in Wikipedia 2021.

BERT base after training on 100k batches		
Softmax (S11P1)	SigSoftmax (S11P1)	
5.8699	5.8749	
Softmax + Multi-input (S11P1)	Softmax + Multi-partition (S11P4)	
5.8520	5.8656	
MoS (Yang et al., 2018) (4) (S41P1)	MoS (Yang et al., 2018) (3) (S31P1)	DOC (Takase et al., 2018) (S31P1)
5.8523	5.8535	5.8547
MFS – Multi-partition (S31P1)	MFS – Multi-input (S31P4)	MFS (S31P4)
5.8231	5.8536	5.8231
BERT large after training on 30k batches		
Softmax (S11P1)	SigSoftmax (S11P1)	
4.8355	4.8354	
Softmax + Multi-input (S11P1)	Softmax + Multi-partition (S11P4)	
4.8305	4.8363	
MoS (Yang et al., 2018) (4) (S41P1)	MoS (Yang et al., 2018) (3) (S31P1)	DOC (Takase et al., 2018) (S31P1)
4.8268	4.8291	4.8231
MFS – Multi-partition (S31P1)	MFS – Multi-input (S31P4)	MFS (S31P4)
4.8111	4.8287	4.8109

Table 7: Prediction visualization using a context in each dataset. Each row visualizes a model as in Table 3. The models are built on *GPT-2 Medium* in OpenWebText and Wikipedia and on *GPT-2 Small* in the synthesized dataset. *MFS Avg* shows the words that are closest to the average facet embedding in *MFS*. See the details in subsection B.3. We underline the words that appear in the top predictions of both *MFS* and *MFS Avg*.

Corpus →	OpenWebText	Wikipedia 2021	Similar Nouns in Templates
Input Context	... "Part of the Clinton inevitability strategy was to lock down the usual suspects in left-liberal policy," said Dan Nexon, a Georgetown professor who served as one of those informal Sanders advisers. Nex	... The projective line over the dual numbers was described by Josef Grünwald in 1906. This ring includes a nonzero nilpotent "n" satisfying. The plane of dual numbers has a project	There are the militia and the enemy in front of a woman, and she decides to pursue the militia
Softmax (GPT-2)	He 0.014, But 0.011, The 0.007	finite 0.062, hom 0.059, project 0.034	enemy 0.860, militia 0.111, Militia 0.005
MFS (Ours)	Nex 0.013, He 0.012, But 0.011	project 0.096, hom 0.049, dual 0.046	enemy 0.535, militia 0.433, enemies 0.029
MFS Avg	" <u>He</u> , <u>But</u> , The, In, And, (. It	<u>hom</u> , <u>dual</u> , finite, non, ", complex, unit	militia , enemy, Militia, enemies , militias
MFS Softmax 1	But 0.005, He 0.004, The 0.002	project 0.201, dual 0.075, finite 0.030	enemy 0.772, militia 0.189, Militia 0.017
MFS Softmax 2	Nex 0.260, " 0.028, He 0.023	hom 0.093, unit 0.040, non 0.037	militia 0.938, Militia 0.062, militias 0.000
MFS Softmax 3	He 0.025, But 0.022, The 0.014	finite 0.065, map 0.041, plane 0.030	enemy 1.000, enemies 0.000, foe 0.003

943 ble 7. Notice that although *GPT-2 XL* has a better
944 ability to distinguish similar words, it would have
945 difficulty in arbitrarily ranking 20 similar words as
946 having the difficulty in ranking 4 analogous words.
947 Similarly, we expect that the GPT-3 (Brown et al.,
948 2020) would still suffer from the *softmax bottleneck*
949 as long as the N is large enough.

950 B.2 Language Modeling using BERT

951 To demonstrate that our proposed method could
952 improve the LMs other than GPT-2, we apply *multi-*
953 *facet softmax*, **MFS**, to BERT. We test the model on
954 Wikipedia 2021 and the validation size is 0.25% of
955 the whole corpus. After loading pretrained model,
956 we train *bert_base_cased* for 100k batches and
957 *bert_large_cased* for 30k batches.

958 The results are presented in Table 6. First, **MoS**
959 outperforms **Softmax** on BERT. The results sup-
960 port the finding of Narang et al. (2021) that the
961 *softmax bottleneck* not only exists in the next word
962 prediction tasks but also in the masked word predic-

963 tion tasks. Similar to GPT-2, **MFS** at least doubles
964 the improvement of **MoS**. The most improvement
965 over **MoS** comes from using multiple input hidden
966 states while adding multiple partitions yield a small
967 or no improvement. Finally, the improvement be-
968 tween **MFS** and **Softmax** is around 4.5%, which is
969 much smaller than 15% in GPT-2.

970 The smaller improvement supports the conclu-
971 sion of our geometric analyses that the multi-mode
972 ambiguity intensifies *softmax bottleneck*. We only
973 observe the one-directional context before the next
974 target word in GPT-2, but we can observe the bi-
975 directional context surrounding the masked target
976 word in BERT. Thus, compared to next word predic-
977 tion, the multi-mode ambiguity of the masked word
978 prediction occurs less frequently when the masking
979 probability is small (e.g., 15% in BERT). Since the
980 masked word distribution only has a single mode
981 most of the time but we sometimes still want the
982 distribution to have multiple modes, multiple input
983 hidden states can improve the performance by help-

Table 8: The loss improvement comparison between the *Improvement Models* and *Reference Models*. The models are named using their number of softmaxes, input hidden states, and partitions. Thus, S3I9P4 is *MFS*, S3I9P1 is *MFS – Multi-partition*, S1I9P1 is *Softmax + Multi-input*, S3I1P1 is *MoS (3)*, and S1I1P1 is *Softmax*. *Multi-mode Percentage* is the percentage of the contexts where the *Improvement Models* output multi-mode distribution. *Multi-mode Loss Improvement* refers to the average improvement when *Improvement Models* outputs multi-mode distribution and *Other Loss Improvement* refers to the improvement of the contexts where the facets of *Improvement Models* are close to each other. *Improvement Ratio* divides *Multi-mode Loss Improvement* by *Other Loss Improvement*.

Corpus →	OpenWebText					Wikipedia 2021				
Improvement Model	S3I9P4	S3I9P4	S3I9P4	S3I9P1	S3I1P1	S3I9P4	S3I9P4	S3I9P4	S3I9P1	S3I1P1
Reference Model	S3I9P1	S3I1P1	S1I1P1	S1I9P1	S1I1P1	S3I9P1	S3I1P1	S1I1P1	S1I9P1	S1I1P1
Multi-mode Percentage (%)	10.03	10.03	10.03	4.81	3.24	5.85	5.85	5.85	2.66	3.05
Multi-mode Loss Improvement	0.0248	0.0474	0.0649	0.0203	0.0110	0.0282	0.0644	0.1000	0.0472	0.0295
Other Loss Improvement	0.0035	0.0158	0.0211	0.0086	0.0064	0.0033	0.0128	0.0219	0.0136	0.0100
Improvement Ratio	7.01	3.00	3.08	2.34	1.71	8.63	5.04	4.57	3.47	2.94

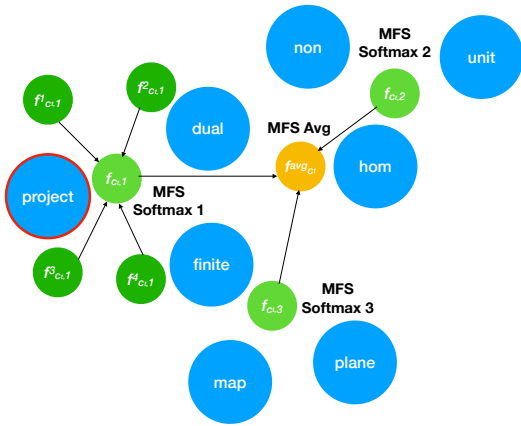


Figure 4: Illustration of the MFS predictions given the Wikipedia context in the second column of Table 7. The green circles mean the facet embeddings from MFS. The orange circle is the average of the facet embeddings (**MFS Avg**). The blue circles are the word embeddings that are close to the facet embeddings and **MFS Avg**. The word *project* is highlighted because it is the next word in our ground truth.

ing the facets to move more freely. On the other hand, multiple partitions are less useful because the distribution rarely has more than three modes.

B.3 Analysis of Improvement on Multi-mode Distribution

To confirm that the perplexity improvements actually come from modeling the multi-mode distribution, we define a metric to measure how multi-mode a distribution is, and then we can compare the perplexity improvement from multi-mode distributions and the improvement from the distributions that are close to a single-mode distribution.

For the method with multiple embeddings, we first compute the weighted average of all the facets

$\mathbf{f}_{c_t}^{avg} = \sum_{k=1}^K \pi_{c_t,k} \mathbf{f}_{c_t,k}$, where we lower the influence of k th facet embedding $\mathbf{f}_{c_t,k}$ with lower prior weight $\pi_{c_t,k}$ and $\mathbf{f}_{c_t,1} = \frac{1}{J} \sum_{j=1}^J \mathbf{f}_{c_t,1}^j$ if J partitions are used. Figure 4 illustrates $\mathbf{f}_{c_t}^{avg}$ and $\mathbf{f}_{c_t,k}$ using the example in the second column of Table 7.

We visualize the new average facet using the words that are closest to the $\mathbf{f}_{c_t}^{avg}$ in the **MFS Avg** row of Table 7. We can see that the prediction of **MFS Avg** is different from **MFS** but similar to **Softmax**. This means there are indeed some other words between the actual next word and the other possibilities, which makes the prediction of **MFS** multi-mode.

Next, to quantify the difference between **MFS** and **MFS Avg**, we define *multi-mode ratio* as $\frac{\sum_{b=1}^T P_M(y_b|c_t)}{\sum_{b=1}^T P_M(x_b|c_t)}$, where P_M could be either P_{MoS} from equation 2 or P_{MP} from equation 4. $\{y_1, \dots, y_T\}$ is the set of words with embeddings closest to $\mathbf{f}_{c_t}^{avg}$ and $\{x_1, \dots, x_T\}$ is the set of words with highest $P_M(x_b|c_t)$. Using the Wikipedia context in Table 7 as an example, the word *project* is retrieved by **MFS** but not by **MFS Avg**, so its *multi-mode ratio* for $T = 2$ is $\frac{P_{MFS}(hom|c_t) + P_{MFS}(dual|c_t)}{P_{MFS}(project|c_t) + P_{MFS}(hom|c_t)} = \frac{0.049 + 0.046}{0.096 + 0.049} \approx 0.66$. Figure 4 illustrates the relation between the **MFS Softmax k** and **MFS Avg**.

When the ratio is closer to 1, the context is less ambiguous and the prediction is closer to a single-mode distribution. We set $T = 20$ and call the prediction with *multi-mode ratio* smaller than 0.9 multi-mode distribution and in Table 8,⁴ we compare the loss (i.e., log of the perplexity) improve-

⁴We also tried $T=5$ or 10 and the trends are similar. If we set the threshold smaller than 0.9, the improvement ratios (e.g., **MFS** over **MoS**) would increase but the multi-mode percentages would decrease.

Table 9: Perplexity comparison of different models on the similar words or dissimilar words. The models are based on *GPT-2 Small* and trained in OpenWebText.

Models ↓	Dissimilar Words			Similar Words		
	Testing	Validation	Training	Testing	Validation	Training
Softmax	1.97	1.98	1.95	2.16	2.16	2.17
MoS (3)	1.81	1.80	1.69	2.05	2.05	1.87
MFS – Multi-partition	1.78	1.79	1.70	2.04	2.06	1.88
MFS	1.79	1.79	1.69	2.02	2.05	1.89

Table 10: MRR (mean reciprocal rank) of different models in OpenWebText. Larger is better.

GPT-2 Small after 1 epoch			
Softmax (S11P1)	SigSoftmax (S11P1)		
0.5494	0.5489		
Softmax + Multi-input (S11P1)	Softmax + Multi-partition (S11P4)		
0.5508	0.5492		
MoS (Yang et al., 2018) (4) (S41P1)	MoS (Yang et al., 2018) (3) (S31P1)	DOC (Takase et al., 2018) (S31P1)	
0.5501	0.5499	0.5494	
MFS – Multi-partition (S31P1)	MFS – Multi-input (S31P4)	MFS (S31P4)	
0.5515	0.5502	0.5519	
GPT-2 Medium after 0.4 epoch			
Softmax (S11P1)	SigSoftmax (S11P1)		
0.5665	0.5650		
Softmax + Multi-input (S11P1)	Softmax + Multi-partition (S11P4)		
0.5677	0.5665		
MoS (Yang et al., 2018) (4) (S41P1)	MoS (Yang et al., 2018) (3) (S31P1)	DOC (Takase et al., 2018) (S31P1)	
0.5674	0.5672	0.5665	
MFS – Multi-partition (S31P1)	MFS – Multi-input (S31P4)	MFS (S31P4)	
0.5685	0.5677	0.5685	

ments in the multi-mode distributions and the improvements in the nearly single-mode distributions.

Table 8 shows that all the multiple embedding approaches have larger loss improvements when outputting multi-mode distributions. The table shows the results based on *GPT-2 Small* and the same analysis using *GPT-2 Medium* also show the same trend. As we use multiple input hidden states and partitions, the differences would be enlarged. Especially when we compare **MFS** and **MFS – Multi-partition**, the loss improvements of highly ambiguous context is 7 or 8 times larger than the other loss improvements, which means a large portion of the overall improvement lies on a small percentage of ambiguous contexts. For the multi-mode distribution in Wikipedia, the loss improvement between **MFS** and **Softmax** could reach 0.10, which is close to the improvement between *GPT-2 Small* and *Medium* (0.16). Thus, we expect that if the corpus has more ambiguous contexts, **MFS** could achieve larger overall loss improvement.

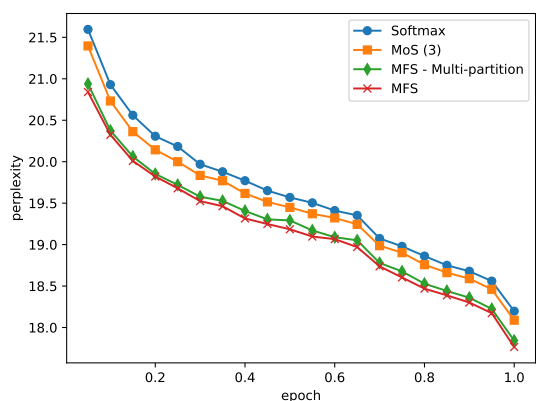
B.4 Template-based Analysis on Similar or Dissimilar Nouns

To know whether the single embedding also has trouble modeling the distribution over nouns without the analogy relation, we let the different models

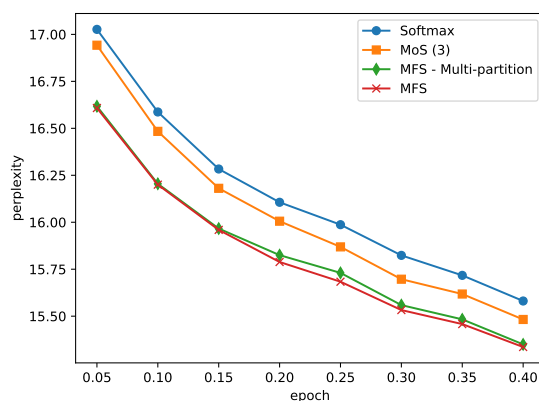
learn to assign similarly high probabilities to two related nouns in our templates. One example in our synthesized dataset is “*I love the **banana** and the **lemon**, and my favorite is the [MASK]*”. The nouns come from a hypernymy detection benchmark (Shwartz et al., 2017) containing 25,498 noun pairs. The relations between nouns in the benchmark include synonym, antonym, attribute, meronym, hypernym, coordination, event, or random. We further split the noun pairs into two datasets based on their cosine similarity in the output word embedding space of our **Softmax** baseline. The pairs with the cosine similarity higher than the medium of all cosine similarities are put into the similar word set and the other pairs are put into the dissimilar word set.

The results are presented in Table 9. In terms of the training, validation, and testing perplexity, multi-embedding approaches consistently outperform the single-embedding baselines, though the margins are smaller than those from the analogous words. Moreover, the improvement gap is larger when the nouns are dissimilar. We hypothesize that as the word embeddings of nouns become further away from each other, the next word distribution is more likely to be multi-mode and thus could be

1083	better captured by multiple embeddings.		
1084	B.5 Adversarial Template Analysis		
1085	To test whether the proposed methods still can ef-		
1086	fectively utilize the information from the global		
1087	word embeddings, we design an adversarial tem-		
1088	plate to create the contexts that can only be com-		
1089	pleted by averaging the global word embeddings.		
1090	For example, “ <i>Miami is not in Wisconsin but is in</i>		
1091	<i>[MASK]=Florida</i> ”.		
1092	In this task, the validation perplexity of Softmax ,		
1093	MoS , MFS – Multi-partition , and MFS are 2.50,		
1094	2.59, 2.54, and 2.88, respectively. Since multiple		
1095	embeddings are not required, it is not surprising		
1096	that Softmax performs the best. Nevertheless, the		
1097	differences are smaller than the differences in Ta-		
1098	ble 4 . We believe that the similar losses are due		
1099	to the fact that multiple embeddings are a general-		
1100	ization of the single embedding, so GPT-2 could		
1101	learn to generate the same embedding for all facets		
1102	to mimic the behavior of single embedding if re-		
1103	quired.		
1104	The significantly worse performance of MFS		
1105	here is caused by the multiple partition technique.		
1106	This result supports our motivation of combining		
1107	multiple partitions with multiple softmaxes and		
1108	shows that multiple partitions handle ambiguous		
1109	contexts better (as shown in Table 8) by sacrificing		
1110	some global word embedding structures. Never-		
1111	theless, a corpus usually has more ambiguous con-		
1112	texts than the adversarial context tested here, so		
1113	using multiple embeddings and multiple partitions		
1114	performs better in Wikipedia and OpenWebText		
1115	overall.		
1116	C More Results		
1117	We provide more numbers and analyses of the am-		
1118	biguous template experiments and language mod-		
1119	eling experiments.		
1120	C.1 Ranking Metric in Language Modeling		
1121	Experiments		
1122	We would like to verify that our perplexity improve-		
1123	ments come from not only the slight probability		
1124	differences of each candidate but also the better		
1125	ranks of the candidates. Thus, in Table 10 , we eval-		
1126	uate different models using mean reciprocal rank		
1127	(MRR). Similar to the perplexity, the MRR im-		
1128	provement from Softmax to MFS is around 15%		
1129	of the MRR improvement from <i>GPT-2 Small</i> to		
1130	<i>GPT-2 Medium</i> , which is similar to the percent-		
	age of perplexity improvement. This suggests that		1131
	MFS could lead to not only a better probability pre-		1132
	diction but also a better candidate rank prediction.		1133
	C.2 Perplexity Curves in Language Modeling		1134
	Experiments		1135
	In Table 1 , we only show the testing perplexity at		1136
	the end of our training. In Figure 6 , we plot the val-		1137
	idation perplexity decay curves during the training		1138
	on OpenWebText. We can see that the performance		1139
	ranking of each model is stable during the training,		1140
	while the improvement of each enhancement may		1141
	vary. For example, in <i>GPT-2 Medium</i> , the improve-		1142
	ment of MFS over MFS – Multi-partition is more		1143
	obvious in epoch 0.25 compared to epoch 0.4.		1144
	C.3 Perplexity Curves in Template Analysis		1145
	In Table 4 , we only show the lowest validation		1146
	perplexity after each of the ten epochs. In Figure 5 ,		1147
	we plot the training and validation perplexity decay		1148
	curves.		1149
	The curves tell us that the multi-embedding		1150
	models perform better in both training and valida-		1151
	tion perplexity. As we train the single-embedding		1152
	models longer, the validation perplexity increases		1153
	quickly, which implies that using a single embed-		1154
	ding to model multi-mode distribution could cause		1155
	severe overfitting when we predict the next word		1156
	given an ambiguous context.		1157
	C.4 Stability in Language Modeling		1158
	Experiments		1159
	In our case, training our model requires a huge		1160
	amount of GPU resources for us, so it is not very		1161
	feasible to train multiple times using multiple ran-		1162
	dom seeds. We indeed try to use different random		1163
	seeds for a few models and we confirm that the val-		1164
	idation loss difference is at least ten times smaller		1165
	than the improvement of different models.		1166
	To verify that our testing dataset is large enough		1167
	to provide stable perplexity, we randomly split the		1168
	testing dataset into 10 subsets and compute the		1169
	standard error of the average testing perplexity of		1170
	the 10 subsets. We find that the standard error is		1171
	less than 0.02 perplexity in all models and datasets		1172
	in Table 1 . The standard error is much smaller than		1173
	most of the improvements, which means our testing		1174
	dataset is large enough to make the reported per-		1175
	plexity stable. The consistent improvements during		1176
	the whole training process in Figure 5 further sup-		1177
	port the stability of our experiments.		1178

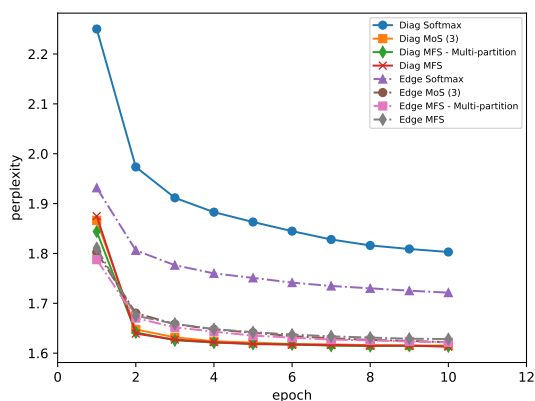


(a) Curves on *GPT-2 Small*

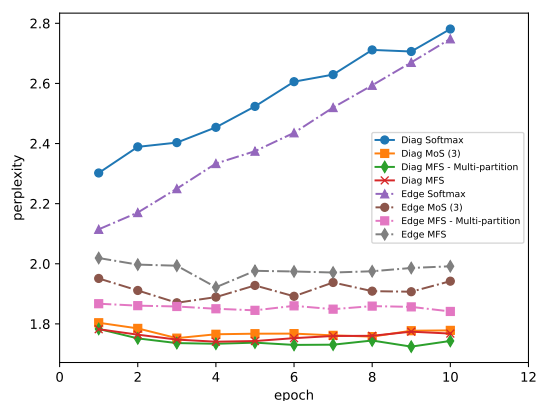


(b) Curves on *GPT-2 Medium*

Figure 5: The perplexity curves for the language modeling tasks using the validation set of OpenWebText.



(a) Perplexity in the training data



(b) Perplexity in the validation data

Figure 6: The perplexity curves from different models for the ambiguous template analysis

Table 11: ProtoQA performances on the crowdsourced development sets. The matching between prediction and ground truth is done by WordNet. All the numbers are percentages. Max answers top k implies only evaluating the top k answers from different LMs. Max incorrect top k indicates only evaluating the top answers that contain k errors. The highest average performances are highlighted and the standard errors are reported as the confidence interval.

Models ↓	Max Answers				Max Incorrect		
	Top 1	Top 3	Top 5	Top 10	Top 1	Top 3	Top 5
Softmax (GPT-2)	36.5 ± 0.7	39.7 ± 0.5	43.5 ± 0.4	52.2 ± 0.6	20.9 ± 0.4	37.7 ± 0.6	46.7 ± 0.6
MoS (Yang et al., 2018) (3)	36.6 ± 0.8	40.2 ± 0.6	43.2 ± 0.6	52.1 ± 0.4	21.3 ± 0.6	38.4 ± 0.5	45.9 ± 0.6
MFS – Multi-partition	37.7 ± 0.7	42.0 ± 0.6	44.6 ± 0.5	52.6 ± 0.3	22.9 ± 0.4	39.5 ± 0.5	47.4 ± 0.4
MFS	36.9 ± 0.7	41.6 ± 0.7	44.4 ± 0.6	52.3 ± 0.6	23.1 ± 0.5	39.7 ± 0.6	46.9 ± 0.6

C.5 ProtoQA Results using WordNet

In Table 5, we report the metrics using exact matching. In Table 11, we report the metrics that match the prediction with the ground truth using WordNet (Miller, 1995) and find the scores show a similar trend.

C.6 Perplexity Improvement versus Model Size

Kaplan et al. (2020) empirically demonstrate that increasing the model size would decrease the loss and their relation follows a scaling law. That is, we can plot the log of model size (i.e., parameter number) versus its loss as in Figure 7, and if a new

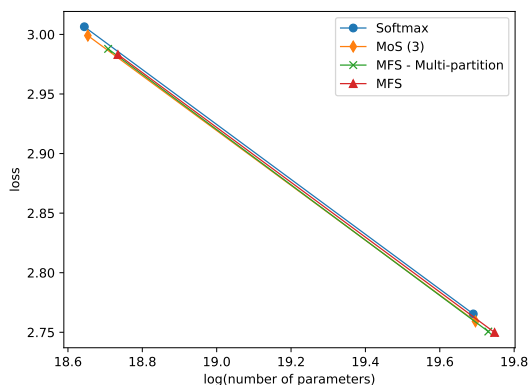


Figure 7: The log of model size versus the log of perplexity in the text set of OpenWebText. The group of points on the left comes from the models based on *GPT-2 Small*. The group of points on the right comes from the models based on *GPT-2 Medium*. The models are trained for 0.4 epoch.

LM model could result in lines that are closer to the origin than the baselines, the new model is better in terms of the loss than only increasing the model size of the baselines.

From Figure 7, we can see that the approaches using multiple embedding are better than the **Softmax** baseline using single embedding. Although the lines formed by **MFS – Multi-partition** and **MFS** are not always closer to the origin than **MoS**, our perplexity improvement from adding multiple input hidden states or multiple partitions cannot be solely explained by their extra parameters for several reasons:

- Compared to **MoS**, the line formed by **MFS – Multi-partition** becomes slightly closer to origin when the model size is close to *GPT-2 Medium*.
- The improvement of **MFS – Multi-partitions (S3I9P1)** is larger than the improvement of **Softmax + Multi-input (S1I9P1)** plus the improvement of **MoS (S3I1P1)** in BERT and GPT-2. For example, in BERT base, the perplexity improvement of **Softmax + Multi-input**, **MoS (3)**, and **MFS – Multi-partitions** are 0.018, 0.016, and 0.047, respectively.
- Our multi-mode analyses in subsection B.3 indicate that our enhancements, especially using multiple partitions, capture the multi-mode distribution better. We expect that the overall perplexity improvement would be larger if the corpus contains more ambiguous contexts. We

also conduct a preliminary experiment to confirm the claim. We add more ambiguous contexts into Wikipedia 2016 by mapping all the uppercased words into the `[UNK]` token. That is, we add another mode corresponding to the `[UNK]` token in many context positions. Then, we train and test the uncased BERT in this synthesized dataset. We found that the improvement of **MFS – Multi-partition** in this case can do significantly better than simply increasing the model size.

- Our enhancements only require some extra linear layers, which are usually more efficient than increasing the model size (e.g., by adding another transformer layer).
- Unlike increasing the model size, keep increasing the number of input hidden states or the number of partitions would lead to a smaller improvement. This suggests that **MFS** cannot keep storing more and more knowledge into its extra linear layers as in the architecture using a larger hidden state size or a deeper transformer encoder.

C.7 More Visualization

In Table 3, we compare the prediction of **MFS** and **Softmax** on *GPT-2 Small*. In the first two columns of Table 7, we present the examples from the models built on *GPT-2 Medium* in OpenWebText and Wikipedia 2021. We can see a similar pattern. The embedding of the correct answer is different from the embeddings of other possibilities, so **Softmax** assigns lower probabilities to the correct answer, while **MFS** does much better. This suggests that a larger model such as *GPT-2 Medium* suffers from the *softmax bottleneck* in a similar way.

In the last column of Table 7, we visualize an example in another synthetic experiment described in subsection B.4. We can see that although there may not be any words between the appropriate candidates, the prediction of **Softmax** may still be biased toward one option much more than the other, while the prediction of **MFS** is much closer to the equally likely bimodal distribution we created in the training data.

D Proof of Theorems

To prove Theorem 1, we first introduce a lemma. Assuming in the word embedding of GPT-2, $\underline{woman} + \underline{king} = \underline{queen} + \underline{man}$, we want to show that the GPT-2 cannot output *woman* and *king* as

the top two words in this lemma. This means we cannot find a hidden state \mathbf{h} and a threshold τ such that $\mathbf{h}^T \mathbf{w}_{\text{woman}} \geq \tau$ and $\mathbf{h}^T \mathbf{w}_{\text{king}} \geq \tau$ but $\mathbf{h}^T \mathbf{w}_{\text{queen}} < \tau$ and $\mathbf{h}^T \mathbf{w}_{\text{man}} < \tau$. This example could be generalized into the following Lemma and Theorems. We can generalize the example as follows:

Lemma 1. *Let the output word embeddings in the set $W = \{\mathbf{w}_{l_j} \neq \mathbf{0} | j = 1 \dots L\} \cup \{\mathbf{w}_{r_j} \neq \mathbf{0} | j = 1 \dots R\}$ satisfy $-a_{l_1} \mathbf{w}_{l_1} - \dots - a_{l_L} \mathbf{w}_{l_L} = a_{r_1} \mathbf{w}_{r_1} + \dots + a_{r_R} \mathbf{w}_{r_R}$, where their coefficient $-a_{l_1}, \dots, -a_{l_L}, a_{r_1}, \dots, a_{r_R}$ are all positive constants and $-a_{l_1} - \dots - a_{l_L} \geq a_{r_1} + \dots + a_{r_R}$. Then, there is no hidden state \mathbf{h} and a threshold τ that make $\min_{\mathbf{w}_g \in G} \mathbf{h}^T \mathbf{w}_g \geq \tau$ and $\max_{\mathbf{w}_s \in S} \mathbf{h}^T \mathbf{w}_s < \tau$, where $G = \{\mathbf{w}_{l_j} | j = 1 \dots L\}$ and $S = \{\mathbf{w}_{r_j} | j = 1 \dots R\}$.*

Proof. To prove by contradiction, we assume there is a \mathbf{h} such that $\forall \mathbf{w}_{l_j} \in G, \mathbf{h}^T \mathbf{w}_{l_j} \geq \tau$ and $\forall \mathbf{w}_{r_j} \in S, \mathbf{h}^T \mathbf{w}_{r_j} < \tau$. Thus, we can get $-a_{l_1} \mathbf{h}^T \mathbf{w}_{l_1} - \dots - a_{l_L} \mathbf{h}^T \mathbf{w}_{l_L} \geq -a_{l_1} \tau - \dots - a_{l_L} \tau \geq (a_{r_1} + \dots + a_{r_R}) \tau > a_{r_1} \mathbf{h}^T \mathbf{w}_{r_1} + \dots + a_{r_R} \mathbf{h}^T \mathbf{w}_{r_R}$, which contradicts to $-a_{l_1} \mathbf{w}_{l_1} - \dots - a_{l_L} \mathbf{w}_{l_L} = a_{r_1} \mathbf{w}_{r_1} + \dots + a_{r_R} \mathbf{w}_{r_R}$. \square

We can rephrase the condition and the conclusion to have our [Theorem 1](#).

Theorem 1. *If the nonzero output embeddings of N words are linearly dependent and on one side of a plane through the origin, the output softmax layer cannot rank the N words with an arbitrary order according to their probabilities.*

Proof. Let the set $W = \{\mathbf{w}_i \neq \mathbf{0} | i = 1 \dots N\}$ contain the embeddings of the N words. Based on the premise, we can write $\mathbf{0} = a_1 \mathbf{w}_1 + \dots + a_N \mathbf{w}_N$ and $\min_{\mathbf{w}_i \in W} \mathbf{h}_0^T \mathbf{w}_i > 0$, where \mathbf{h}_0 is a normal vector of the plane. At least one of the a_i is negative. Otherwise, we will get the contradiction $0 = \mathbf{h}_0^T \mathbf{0} = a_1 \mathbf{h}_0^T \mathbf{w}_1 + \dots + a_N \mathbf{h}_0^T \mathbf{w}_N \geq (a_1 + \dots + a_N) \min_{\mathbf{w}_i \in W} \mathbf{h}_0^T \mathbf{w}_i > 0$. Similarly, at least one of a_i is positive. We can move all the terms in $\mathbf{0} = a_1 \mathbf{w}_1 + \dots + a_N \mathbf{w}_N$ with negative a_i to the left as $-a_{l_1} \mathbf{w}_{l_1} - \dots - a_{l_L} \mathbf{w}_{l_L} = a_{r_1} \mathbf{w}_{r_1} + \dots + a_{r_R} \mathbf{w}_{r_R}$. If $-a_{l_1} - \dots - a_{l_L} \geq a_{r_1} + \dots + a_{r_R}$, we choose $G = \{\mathbf{w}_{l_j} | j = 1 \dots L\}$. Otherwise, we choose $G = \{\mathbf{w}_{r_j} | j = 1 \dots R\}$

Based on Lemma 1, we know that there is a partition $P = \{G, S\}$ such that we cannot have the logits of the words in G be always larger than the logits of the other words in S , so all the words in

G cannot have the probabilities larger than every probability of the words in S . \square

Next, we would like to generalize our [Theorem 1](#) by using a more practical condition where the word embeddings are almost linearly dependent. Notice that the theorem needs to assume the magnitude of the hidden state is limited. Otherwise, the margin could be arbitrarily magnified. In practice, the magnitude is not arbitrarily large in GPT-2 and BERT because a too large magnitude of hidden state could magnify the gradients too much to have a stable training process.

Theorem 2. *Let the output word embeddings in the set $W = \{\mathbf{w}_i \neq \mathbf{0} | i = 1 \dots N\}$ satisfy $\mathbf{w}_1 = a_2 \mathbf{w}_2 + \dots + a_N \mathbf{w}_N + \boldsymbol{\varepsilon}$, where the constant a_2, \dots, a_N are neither all zero nor all negative and $\|\boldsymbol{\varepsilon}\| < \epsilon$. Then, there must be a non-trivial partition $P = \{G, S\}$ of W such that there is no hidden state $\|\mathbf{h}\| \leq r$ and a threshold $\tau \geq r\epsilon$ that make $\min_{\mathbf{w}_g \in G} \mathbf{h}^T \mathbf{w}_g \geq (1 + \delta)\tau$ and $\max_{\mathbf{w}_s \in S} \mathbf{h}^T \mathbf{w}_s < \tau$, where $\delta = \frac{2}{1 + \sum_{i=2 \dots N} |a_i|}$.*

Proof. We can first move all the terms with negative a_i to the left as $\mathbf{w}_1 - a_{l_1} \mathbf{w}_{l_1} - \dots - a_{l_L} \mathbf{w}_{l_L} = a_{r_1} \mathbf{w}_{r_1} + \dots + a_{r_R} \mathbf{w}_{r_R} + \boldsymbol{\varepsilon}$. We perform proof by contradiction, so we assume the logits of the words in G can always be larger than $(1 + \delta)\tau$ and the logits of the words in S can always be smaller than τ .

Case 1: $1 - a_{l_1} - \dots - a_{l_L} \geq a_{r_1} + \dots + a_{r_R}$, so $1 - a_{l_1} - \dots - a_{l_L} \geq \frac{1 + \sum_{i=2 \dots N} |a_i|}{2}$. We choose $G = \{\mathbf{w}_1, \mathbf{w}_{l_1}, \dots, \mathbf{w}_{l_L}\}$ and $S = \{\mathbf{w}_{r_1}, \dots, \mathbf{w}_{r_R}\}$. Thus, we can get $\mathbf{h}^T \boldsymbol{\varepsilon} \leq \|\mathbf{h}\| \|\boldsymbol{\varepsilon}\| \leq r\epsilon \leq \tau$ and

$$\mathbf{h}^T \mathbf{w}_1 - a_{l_1} \mathbf{h}^T \mathbf{w}_{l_1} - \dots - a_{l_L} \mathbf{h}^T \mathbf{w}_{l_L} \quad (5)$$

$$\geq (1 - a_{l_1} - \dots - a_{l_L})(1 + \delta)\tau \quad (6)$$

$$= (1 - a_{l_1} - \dots - a_{l_L}) \left(1 + \frac{2}{1 + \sum_{i=2 \dots N} |a_i|}\right) \tau \quad (7)$$

$$\geq (1 - a_{l_1} - \dots - a_{l_L}) \left(1 + \frac{1}{1 - a_{l_1} - \dots - a_{l_L}}\right) \tau \quad (8)$$

$$= (1 - a_{l_1} - \dots - a_{l_L} + 1)\tau \quad (9)$$

$$\geq (a_{r_1} + \dots + a_{r_R} + 1)\tau \quad (10)$$

$$> a_{r_1} \mathbf{h}^T \mathbf{w}_{r_1} + \dots + a_{r_R} \mathbf{h}^T \mathbf{w}_{r_R} + \mathbf{h}^T \boldsymbol{\varepsilon}, \quad (11)$$

which contradict with $\mathbf{w}_1 - a_{l_1} \mathbf{w}_{l_1} - \dots - a_{l_L} \mathbf{w}_{l_L} = a_{r_1} \mathbf{w}_{r_1} + \dots + a_{r_R} \mathbf{w}_{r_R} + \boldsymbol{\varepsilon}$.

Case 2: $1 - a_{l_1} - \dots - a_{l_L} < a_{r_1} + \dots + a_{r_R}$. We choose $G = \{\mathbf{w}_{r_1}, \dots, \mathbf{w}_{r_R}\}$ and $S =$

$\{\mathbf{w}_1, \mathbf{w}_{l_1}, \dots, \mathbf{w}_{l_L}\}$. Therefore,

$$a_{r_1} \mathbf{h}^T \mathbf{w}_{r_1} + \dots + a_{r_R} \mathbf{h}^T \mathbf{w}_{r_R} \quad (12)$$

$$\geq (a_{r_1} + \dots + a_{r_R}) \left(1 + \frac{2}{1 + \sum_{i=2 \dots N} |a_i|}\right) \tau \quad (13)$$

$$> (a_{r_1} + \dots + a_{r_R}) \left(1 + \frac{1}{a_{r_1} + \dots + a_{r_R}}\right) \tau \quad (14)$$

$$= (a_{r_1} + \dots + a_{r_R} + 1) \tau \quad (15)$$

$$> (1 - a_{l_1} - \dots - a_{l_L} + 1) \tau \quad (16)$$

$$> \mathbf{h}^T \mathbf{w}_1 - a_{l_1} \mathbf{h}^T \mathbf{w}_{l_1} - \dots - a_{l_L} \mathbf{h}^T \mathbf{w}_{l_L} - \mathbf{h}^T \boldsymbol{\varepsilon}. \quad (17)$$

□

E Method Details

When replacing the softmax layer in the pretrained LMs, we found that the initialization of the extra linear layers should make the initial prediction of LMs close to the prediction using a softmax layer, which is the architecture used in the pretraining. The initialization is especially important for BERT. To achieve the goal, we initialize the weights of linear layer such that different facets are almost identical at the beginning and let the LMs gradually learn to output diverse facets during the training. Specifically, we can write the linear layer on the new hidden state $L_k^f(\mathbf{q}_{c_t})$ as

$$\begin{aligned} \mathbf{f}_{c_t, k} &= L_k^f(\mathbf{q}_{c_t}) \\ &= \mathbf{L}_k^I \mathbf{h}_{c_t}^M + \mathbf{L}_k^B \text{GELU} \left(L^h(\oplus_{i, m} \mathbf{h}_{c_t-i}^{M-m}) \right) + \mathbf{b}. \end{aligned} \quad (18)$$

We initialize \mathbf{L}_k^I as an identity matrix, $\mathbf{b} \leftarrow \mathbf{0}$, and $\mathbf{L}_k^B \leftarrow \mathcal{U}(-\epsilon, \epsilon)$, where \mathcal{U} is the uniform distribution and $\epsilon = 0.00005$ if $k \neq K$. Otherwise, $\epsilon = 0$. Consequently, all the facets $\mathbf{f}_{c_t, k}$ are initially close to the last hidden state of the original GPT-2 $\mathbf{h}_{c_t}^M$. Our baselines (e.g., **Softmax**, **MoS**, and **DOC**) also adopt the same way to initialize their weights.

When partitioning the vocabulary, we simply let the j th facet handle the word with index $J \times n + j$ (e.g., the first partition includes the words with indexes 0, 4, 8, ... when the number of partitions $J = 4$). The partition way is easy to be implemented in PyTorch and it won't significantly increase computational time because PyTorch supports the dilated variable access without needing to copy the whole output word embedding matrix.

We implement our models based on huggingface⁵ (Wolf et al., 2020) and we will release our code to provide more details in our methods.

E.1 Architecture Differences in BERT

The architecture of **MFS** for BERT is mostly the same as the one for GPT-2 and the differences are described in this subsection.

In GPT-2 the block of input hidden state is right-aligned with the last word to prevent seeing the ground truth. On the other hand, the block in BERT is centered at the masked word.

The softmax layer of BERT is slightly different from that of GPT-2. For example, BERT adds a bias term after the dot product between the hidden state and the output word embedding. We keep the bias term in our experiments. Besides, the pretrained BERT has a language modeling head including a linear layer, a GELU (Gaussian Error Linear Unit) layer (Hendrycks and Gimpel, 2016), and a layernorm layer (Ba et al., 2016), so instead of adding an extra linear layer as in GPT-2, we just use different language modeling heads to create different facets in BERT. All the heads are initialized using the weights in the pretrained BERT except that the linear layer is initialized as in Equation 18 when the multiple input hidden states are used and the corresponding linear weights $\mathbf{L}_k^B \leftarrow \mathcal{U}(-\epsilon, \epsilon)$, where $\epsilon = 0.05$ if $k \neq K$. Otherwise, $\epsilon = 0$.

F Experimental Details

In this section, we describe some details of our experimental setup. We will release our codes to reproduce our results once the paper is accepted.

F.1 Baselines

The **MoS** (Yang et al., 2018) and **DOC** (Takase et al., 2018) are originally designed for RNN-based LM. To improve their methods on pretrained Transformer-based LM and make their results more comparable to **MFS**, we change some of their implementation details.

MoS originally has a tanh layer before the softmax layers. However, we found that adding tanh hurts the performances of all methods we tested, especially the **Softmax** and **MoS** baselines. For example, after adding tanh and training GPT-2 Small for 0.4 epoch on Wikipedia, the validation perplexity degradation of **Softmax** is from 25.70 to 26.15, the degradation of **MoS** is from 25.42 to 25.83, and

⁵<https://huggingface.co/>

the degeneration of **MFS** is from 25.06 to 25.12. We suspect this is because GPT-2 is pretrained without the tanh layer and the tanh removes the magnitude of facets, which could be viewed as the inverse of the temperature in the softmax layer. Therefore, we remove the tanh layer in all of our experiments. From the theoretical perspective, adding tanh does not invalidate our motivation because tanh is invertible. It would just make the global transformation nonlinear. Therefore, our motivation of making the facets move freely by inputting multiple hidden states still holds even after adding tanh.

In **DOC**, we use the hidden states of the last three transformer layers to compute the three facets and we set $\lambda_\beta = 0$. Each facet is only determined by one layer of hidden state, so the first two facets cannot access the last hidden state. We found that the model quickly learns to only use the last facet because only the last hidden state is trained to perform the LM task in the pretrained models. This prevents the first two facets from getting any gradients and causes a starvation problem.

We tried an aggressive dropout trick to solve the starvation problem in **DOC**. If one of the softmaxes does not assign the highest probability to any of the correct next words in a batch, we consider that the corresponding facet starves, so we drop the other facets with some probability to ensure this starved facet receives some gradients and gradually gets back on track. However, our preliminary experiment suggests that the dropout trick cannot improve the perplexity of **DOC**. The dropout probability is either too low to solve the starvation problem or too high to preserve the knowledge learnt from pre-training. Thus, we do not adopt this trick in our final experiment.

F.2 Language Modeling

We download Wikipedia using http://medialab.di.unipi.it/wiki/Wikipedia_Extractor and OpenWebText using <https://github.com/jcpeterson/openwebtext>. For Wikipedia, we preprocess the text using <https://github.com/attardi/wikiextractor>. For OpenWebText, we download the pre-filtered URLs in 2017 and 2018 and scrape the text on April 2021. When splitting the corpus into training, validation, and testing sets, we do not shuffle the data. Instead, we use the text near the end of the corpus as the validation and test set to reduce the information

leakage. To ensure every model is trained on the same data and accelerate the training in our machines, we split the training data into 20 consecutive partitions and load only one partition at a time during training. For BERT, we perform the sentence segmentation using SpaCy⁶ and input one sentence into BERT at a time.

We set our hyperparameters (e.g., $W \times H = 3 \times 3$ when using multiple input hidden states) based on the validation performance in Wikipedia 2016, the resulting model size, and the memory constraint in GPUs. We use AdamW (Loshchilov and Hutter, 2019) optimizer and set the learning rate as $1e-5$ and do not use the warm-up because the training starts from the pretrained models. The sequence length (i.e., bptt) is set as 200 for GPT-2 and 256 for BERT. The batch size are set as 4 for *GPT-2 Small*, 16 for *GPT-2 Large*, 120 for *BERT base*, and 128 for *BERT large*.

The analyses in Table 2 and Table 8 use the first 4000 sequences in the validation dataset and all the methods are based on *GPT-2 Small*. We use PYCLD2⁷ to distinguish between English and non-English text.

We use NVIDIA GeForce RTX 2080 for training *GPT-2 Small* and *BERT base*, GeForce RTX 8000 for training *GPT-2 Medium*, Tesla M40 for training *BERT large*. Since we start from the pretrained LM, we can finish training each LM within 2 weeks using 1 GPU for *GPT-2 Small*, *BERT base*, and *GPT-2 Medium*, and using 4 GPUs for training *BERT large*.

When testing the inference time in Table 1, we average the time of running NVIDIA TITAN X on 10,000 batches, where each batch contains 4 sequences with length 200.

When visualizing the prediction in Table 3, we exclude the non-ASCII symbol prediction from the top word list of all models.

F.3 Ambiguous Templates Analysis

Among the semantic relations in Google analogy dataset, we choose three different relations between locations: *capital-common-countries*, *capital-world*, *city-in-state*, and one relation between people: *family*. We exclude the *currency* category because their instance often does not form a parallelogram in the word embedding space (Ethayarajh et al., 2019). The templates we use are listed

⁶<https://spacy.io/>

⁷<https://github.com/aboSamoor/pycltd2>

Table 12: The templates used in the analysis. The first four templates are for the analogy relations from the *capital-common-countries*, *capital-world*, and *city-in-state* categories. The next four templates are for the analogy relations from the *family* category. The final four templates are for similar or dissimilar nouns.

Dataset ↓	Templates
Anology (Person or Person)	Between the \$ARG1 and the \$ARG2, I decided to first talk to the [MASK] The \$ARG1 and the \$ARG2 are my favorites, and I especially love the [MASK] The \$ARG1 and the \$ARG2 happily live together. One day, bad luck happens to the [MASK] The \$ARG1 and the \$ARG2 stay at my house, and I need to take care of the [MASK]
Anology (Location or Location)	I went to \$ARG1 and \$ARG2 before, and I love one of the places more, which is [MASK] \$ARG1 and \$ARG2 are my favorites, and I especially love [MASK] My uncle used to live in \$ARG1 and \$ARG2 but now, he is selling his house in [MASK] The traveler plans to visit \$ARG1 and \$ARG2, and the traveler first arrives in [MASK]
Similarity (Noun or Noun)	I love the \$ARG1 and the \$ARG2, and my favorite is the [MASK] Yesterday, a man encountered the \$ARG1 and the \$ARG2. Today, he again saw the [MASK] There are the \$ARG1 and the \$ARG2 in front of a woman, and she decides to pursue the [MASK] If you can choose the \$ARG1 or the \$ARG2, would you choose the [MASK]

in Table 12. For the *family* category, our templates assume the words are not pronouns, so we exclude the set of four words that include *he* or *she*.

For each of the four words in an analogy instance (e.g., *queen* : *king* = *woman* : *man*), we would create 32 training or testing sequences⁸ based on the diagonal words such as *king* or *woman*. Similarly, we would create 64 sequences in the edge datasets. Some words contain multiple word pieces and we average the losses of all word pieces during training and testing.

We split the synthesized sequences based on their word pair overlapping. First, we randomly sample half of the word pairs (e.g., *king* and *queen*) in each category as our training pairs. If both of the word pairs in an analogy instance are training pairs, the instance is put into our training set. If only one of the word pairs is a training pair, the instance would belong to our validation set. The rest of the instances form our testing set. During the fine-tuning, we evaluate a model using the validation set after each epoch and select the model based on its best validation perplexity.

F.4 ProtoQA Evaluation

In our experiments, we use the scraped development set as our validation set and the crowdsourced development set as our test set. We do not test our methods on the test set of ProtoQA because the result of every submission would show up in their leaderboard and we do not want to overwhelm the leaderboard with our 15 trials.

⁸2 (diagonal words) × 4 (templates) × 2 (word orders in the template) × 2 (possible next words)

Due to our limited GPU resources, we compare the methods built on *GPT-2 medium* rather than *GPT-2 large*. To maximize the perplexity of the *GPT-2 medium* model using Softmax on the scraped development set, we fine-tune our models using learning rate 3e-5 and warmup step 500.

The original paper (Boratko et al., 2020) does not consider the frequency of the answer during the fine-tuning (i.e., the most possible answer and the least possible answer of each question appear in the training data with the same chance). In terms of the performance of the scraped development set, we find that weighting each answer based on the square root of its frequency is better than weighting each answer uniformly as in the original paper or weighting each answer based on its frequency, so we use the square root weighting to finetuning all our models.

During testing time, each model generates the answers using Nucleus Sampling (Holtzman et al., 2020) with $p = 0.9$ and temperature = 1. Then, we collect all the words before the first period as an answer and drop the generated sentences without a period.

G Future Work

Capturing the next word distribution well given an ambiguous context could be important in some downstream applications. A next step could be investigating whether multiple facets lead to a better language generation model for the applications. For example, we would like to know whether breaking the *softmax bottleneck* could reduce the hallucination of LMs (e.g., outputting *queen* when the reasonable next words should be *king* or *woman*)

1614 and increase the coherence of the generated text. 1665
1615 We also want to more systematically investigate 1666
1616 whether modeling multi-mode distribution could 1667
1617 help LMs to reduce the undesired bias and to better 1668
1618 distinguish similar words (Zagoury et al., 2021) as 1669
1619 in subsection B.4. 1670

1620 Narang et al. (2021); Anonymous (2021) find 1671
1621 that *MoS* can significantly improve the BERT-like 1672
1622 LMs on natural language understanding (NLU) 1673
1623 tasks when the LMs are trained from scratch. Al- 1674
1624 though we find that the perplexity improvement of 1675
1625 multi-embedding BERT is not as large as multi- 1676
1626 embedding GPT-2, pretraining using multiple em- 1677
1627 beddings does not decrease the inference speed of 1678
1628 the BERT encoder on NLU tasks. This motivates 1679
1629 the future studies of seeing if *MFS* also provides a 1680
1630 larger improvement than *MoS* in NLU tasks. 1681

1631 Table 2 suggests that multiple embeddings im- 1682
1632 prove more in a non-English context. We wonder 1683
1633 whether multiple embeddings are more beneficial 1684
1634 to the LMs that are trained on a non-English do- 1685
1635 minating corpus. Chung et al. (2021) discover that 1686
1636 using a larger output embedding dimension im- 1687
1637 proves the multilingual BERT. An interesting re- 1688
1638 search question is whether the improvement comes 1689
1639 from alleviating the *softmax bottleneck* and whether 1690
1640 *MFS* could also lead to similar improvements in 1691
1641 multilingual benchmarks. 1692

1642 The hidden state size of GPT-3 175B (Brown 1693
1643 et al., 2020) is huge (12,288). An interesting ques- 1694
1644 tion is whether some sets of output word embed- 1695
1645 dings in GPT-3 are still in a low-dimensional sub- 1696
1646 space and whether the *softmax bottleneck* is still 1697
1647 a prominent problem on the road of pursuing gen- 1698
1648 eral intelligence when such a large hidden state 1699
1649 dimension is used. We also would like to know if 1700
1650 models using multiple facets could reach a similar 1701
1651 performance by a smaller hidden state size. 1702

1652 Recently, Gao et al. (2019); Rajaei and Pilehvar 1703
1653 (2021); Cai et al. (2021) point out the structure in 1704
1654 the contextual embedding space prevents it from 1705
1655 having an isotropic property. Our study and Deme- 1706
1656 ter et al. (2020) show that the structure in the word 1707
1657 embedding space only models the global similar- 1708
1658 ity between words and prevents the LM from out- 1709
1659 putting arbitrary context-dependent word distribu- 1710
1660 tions. We would like to know if we can discover a 1711
1661 new LM architecture with a better contextual/word 1712
1662 embedding space that could better model context- 1713
1663 dependent word similarities and balance it with the 1714
1664 global word similarities. 1715

1665 *Mixtape* (Yang et al., 2019) is another efficient 1666
1667 solution to the *softmax bottleneck*, whose hidden 1668
1669 state for each word is the weighted average of the 1670
1671 facets where the weights are dynamically predicted. 1672
1673 If only using one softmax (i.e., $K = 1$), our mul- 1674
1675 tiple partition method could be viewed as a spe- 1676
1677 cial case of *Mixtape* that uses a global and bina- 1678
1679 rized weight to prevent complications of predicting 1680
1681 weights of each word. Our results indicate that mul- 1682
1683 tiple partitions need to be combined with multiple 1684
1685 softmax layers in order to gain consistent perfor- 1686
1687 mance improvement. A potential future direction 1688
1689 is to compare *MFS* and *Mixtape* on the transformer- 1690
1691 based LMs or combine the ideas from *MFS* and 1692
1693 *Mixtape* to gain further improvements. 1694

1695 The results in Kong et al. (2020) suggest that 1696
1697 predicting n-gram could be better than predicting 1698
1699 individual words in BERT in some applications. 1700
1701 The total number of possible n-gram is several or- 1702
1703 ders of magnitude higher than the number of indi- 1704
1705 vidual tokens in the vocabulary. In addition, the lin- 1706
1707 ear dependency among n-gram might be common. 1708
1709 For example, the embedding of the *brown color* 1709
1710 + *a dog* may be similar to the embedding of the 1710
1711 *brown dog*. The problem would be more serious as 1711
1712 the length of the prediction sequence (n) increases, 1712
1713 so predicting the next sentence using a single em- 1713
1714 bedding might suffer from the *softmax bottleneck* 1714
1715 even more. Therefore, our solutions to *softmax* 1715
1716 *bottleneck* may lead to a better phrase represen- 1716
1717 tation or sentence representation in this type of 1717
1718 self-supervised pretraining. 1718

1719 Finally, language modeling is only an example of 1719
1720 extreme classification. The nearly ubiquitous usage 1720
1721 of single embedding representation in the classi- 1721
1722 fication, self-supervised models (e.g., contrastive 1722
1723 learning models), or recommendation problems 1723
1724 provides many research opportunities. We believe 1724
1725 that our theoretical results could guide researchers 1725
1726 to identify the potential applications where the *soft-* 1726
1727 *max bottleneck* is serious and multi-embedding rep- 1727
1728 resentation is accordingly helpful. 1728