

# CONTEXT-AWARE TEMPERATURE FOR LANGUAGE MODELING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Current practices to apply temperature scaling assume either a fixed, or a manually-crafted dynamically changing schedule. However, our studies indicate that the individual optimal trajectory for each class can change with the context. To this end, we propose context-aware temperature, a generalized approach to provide an individual optimal temperature trajectory over the context for each vocabulary, while allowing the temperature to be learned along with the remaining model parameters during training. Experiment results confirm that the proposed method significantly improves state-of-the-art language models, achieving a perplexity of 19.90 on Penn Treebank, 33.88 on WikiText-2, and 4.7 on WikiText-103.<sup>1</sup>

## 1 INTRODUCTION

Temperature scaling has been widely used in various domains, such as natural language processing, model calibration, and knowledge distillation, to adjust the smoothness of a distribution, and achieve great performances or control the attribute of generated outputs (Norouzi et al., 2016; Hu et al., 2017; Caccia et al., 2018). Generally, a temperature scalar (or vector)  $\tau$  is applied as a denominator to the logits, and then the divided logits pass through a Softmax layer to yield a probability distribution. If the temperature value  $\tau \rightarrow \infty$ , the distribution becomes more uniform, thus increasing the uncertainty. Contrarily, when  $\tau \rightarrow 0$ , the distribution collapses to a point mass. Although temperature scaling has been justified to achieve great success, existing implementations are limited. They assume either a constant value, or a manually-defined schedule for the temperature (Jang et al., 2017; Ma et al., 2017; Xie et al., 2019). Most importantly, none of them studies the effects on different word tokens when the temperature changes.

To this end, we propose *context-aware temperature*, which is capable of generating an unique temperature value for each token. Through taking the change of contexts into consideration, context-aware temperature serves as a generalized method to provide an optimal temperature for each word. Figure 1 illustrates the temperature trajectories of the tokens during the course of training. The word tokens along the Y-axis are from an input sentence, and the tokens along the X-axis are a selected subset of the vocabulary. At each row  $i$ , we show the contextual temperature values generated for the tokens along the X-axis, taking the context of the  $i$ -th input token in consideration. As shown in the figure, the temperatures of tokens gradually “heat up” or “cool down” as we sequentially process the input tokens. For instance, at row 5 (for the input token “old”), we see that “years” and “months” have obviously heated-up temperatures, while the temperature of “old” is slightly cooling down. These temperatures at row 5 are generated by our proposed model, considering the context of preceding tokens “pierre <unk> N years”, with the goal of predicting the token “old”.

Furthermore, the diagonal pattern indicates that our context-aware model assigns high temperature values to the token that just occurred, effectively suppressing the same token to appear repeatedly in a sentence. Such context-aware behavior of our model intuitively matches human intuition that words do not usually occur consecutively in a sentence. With these observations, we argue that existing methods limit themselves to some fixed schedules, and thus have great difficulty to generalize. The proposed context-aware temperature provides a more generalized temperature mechanism.

Our experimental results on language modeling datasets, including Penn Treebank, WikiText-2 and WikiText-103, demonstrate significant improvements. We also conduct comprehensive analyses

<sup>1</sup>We will release the code upon paper acceptance.



Here  $\tau$  is a vector with a same dimension  $|V|$  as  $\mathbf{z}$ , and is normally set to 1. Below we divide related works into three categories: (a) constant temperature, where each element in  $\tau$  has the same value, (b) dynamic temperature over training iterations, where  $\tau$  is constant in one single iteration, but has different values in every iteration, and (c) dynamic temperature over word position, where  $\tau$  is dynamic in one iteration, and besides,  $\tau$  could change in every iteration.

**Constant Temperature.** Earlier works can be traced back to model distillation (Hinton et al., 2014), where  $\tau$  is constant and chosen empirically. Constant temperature is also used during training (Norouzi et al., 2016; Ma et al., 2017) to control the degree of augmentation. Other works that incorporate  $\tau$  at inference time include model calibration (Guo et al., 2017), and controlling the trade-off between quality and diversity in text generation tasks (Caccia et al., 2018).

**Dynamic Temperature Over Training Iterations.** Most works adopt dynamic temperature through a manually-crafted schedule. Notably, Hu et al. (2017) uses an approximation based on Softmax with a decreasing temperature to enable gradient propagation. Similar techniques are adopted in gumbel-softmax (Jang et al., 2017). In addition, Zhang et al. (2018) shows that with a heating up temperature scaling, embedding vectors are more compact.

**Dynamic Temperature Over Attention.** Another work that is closely related to our work, is the adaptive temperature over an attention model (Lin et al., 2018). The model learns to output a dynamic temperature to control the softness of the attention. The temperature is learned based on the information of decoding at the current step and the attention in previous steps.

### 3 METHODS

#### 3.1 CONTEXT-AWARE TEMPERATURE

Context-aware temperature is a mechanism that chooses the optimal temperature for each vocabulary token, by considering the ‘‘context’’ of a token  $x_t$ , which is the history  $x_{1:t-1}$ .

In this work, the temperature vector  $\tau \in \mathbb{R}^{|V|}$  is generated based on the hidden representations learned from a non-linear mapping function  $f$  (as described in Equation 1). The function  $f$  can be any sequential models. In this work, we explore two parameterizations of the function  $f$ , based on two different models, AWD-LSTM (Merity et al., 2018) and Transformer-XL (Dai et al., 2019). As an example, Figure 2 illustrates the architecture of our model when implemented on top of the MoS model (Yang et al., 2018). In this figure, the above-mentioned function  $f$  is the output of the AWD-LSTM layer, and the components colored in blue constitute the proposed mechanism of context-aware temperature, which we will describe next.

The output from  $f$  is then multiplied with  $W_\tau$  to produce the logits of the temperatures, denoted as  $\mathbf{z}_\tau$ , which subsequently are scaled using a Softmax layer and two parameters  $(\alpha, \beta)$ :

$$\tau = \frac{\sigma(\mathbf{z}_\tau) + \alpha}{\beta}, \text{ where } \mathbf{z}_\tau = f(x_{1:t-1}; \theta^{emb}, \theta^{model})^T \cdot W_\tau. \quad (4)$$

Note that while the values of  $(\alpha, \beta)$  could be assigned manually, they could be learnt by the model itself. We choose to implement them as learnable parameters in this work.

#### 3.2 CONTEXT-AWARE TEMPERATURE MOS

The first model that we explored is the *Context-aware Temperature Mixture of Softmaxes (CT-MoS)*, which extends the MoS model for language modeling with our proposed context-aware temperature mechanism. Equation 5 shows the modified mechanism with the context-aware temperatures.

$$P_{\text{CT-MoS}}(x_t = k | x_{1:t-1}; \Theta) = \sum_m^M \pi_m \cdot \sigma(\mathbf{z}_m \oslash \tau) \quad (5)$$

where  $\oslash$  represents element-wise division.  $\Theta = \cup_{m=1}^M W_m \cup \theta^{emb} \cup \theta^{model} \cup W_\tau$  is the set of all of the parameters of CT-MoS. And,  $\pi_m$  are the weights for the logits, which are defined by the

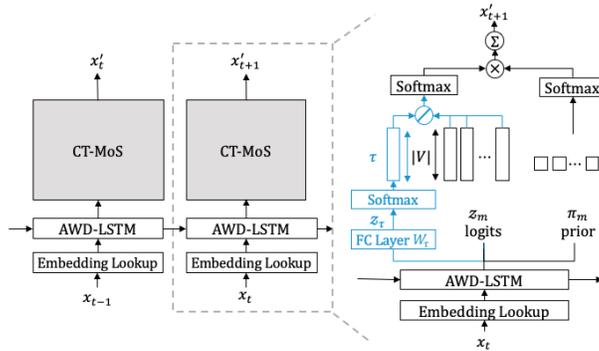


Figure 2: The architecture of the proposed CT-MoS model. Black components are those the same as the MoS model, while the blue ones are the newly added ones in the proposed approach.

MoS model. Please note that using a single matrix  $W_\tau$  may increase the number of parameters significantly, and thus in practice we factorize it into two matrices  $W_{\tau_1}$  and  $W_{\tau_2}$ . Figure 2 highlights the difference between the proposed CT-MoS model and the MoS model.

Compared to prior works, the proposed context-aware temperature is not only dynamic over training iterations, but also dynamic over the word positions in a sentence (see Figure 1). That is, for the same token at different positions with different context, the proposed method would learn a different temperature vector dependent on the token’s history context.

### 3.3 CONTEXT-AWARE TEMPERATURE TRANSFORMER-XL

Besides experimenting with the MoS model, we apply the context-aware temperature mechanism on the Transformer-XL (Dai et al., 2019) model for language modeling. Compared to MoS, Transformer-XL is a purely self-attentive model. We apply the context-aware temperature on the logits  $\mathbf{z}_{\text{TXL}}$  of the Transformer-XL model, and denote the resulting model as “CT-Transformer-XL”.

$$P_{\text{CT-Transformer-XL}}(x_t = k | x_{1:t-1}; \Theta) = \sigma(\mathbf{z}_{\text{TXL}} \odot \boldsymbol{\tau}) \tag{6}$$

where  $\Theta = \theta_{emb} \cup \theta_{model} \cup W_\tau$ .

### 3.4 EFFECTS OF CONTEXT-AWARE TEMPERATURE

In this section, we discuss how context-aware temperature effects the logits  $\mathbf{z}$  and the temperature itself  $\boldsymbol{\tau}$ , through illustrating their corresponding gradients. We consider a language modeling task with a small vocabulary of only two words, i.e.,  $|V| = 2$ . In this setting, the dimensionality of logits  $\mathbf{z}$  is 2 and so is that of the temperature vector  $\boldsymbol{\tau}$ . The range of  $\boldsymbol{\tau}$  is set to  $[0, 1]$ , that is,  $(\alpha, \beta) = (0, 1)$ .

**Gradient of the logit** At a given input, assume that the ground-truth token is  $i = 0$ , thus the cross-entropy loss is  $L = -\ln p_0$ , where  $p_0$  is the model’s output probability for word token  $i = 0$ . In this case, the gradients of logits  $\mathbf{z}$  are illustrated in Figure 3 (the derivation of the gradients is in Appendix A.1). In Figure 3(a), we consider the gradient of  $z_1$  of token  $i = 1$ . When no temperature mechanism is applied, it can be seen that the magnitude of the gradient is bounded within  $[0, 1]$  (as shown by the red mesh), with the largest magnitude 1 (most aggressive update) happens when the probability  $p_1$  is closer to 1 (note that the ideal value for  $p_1$  should be close to 0, as the other token  $i = 0$  is assumed to be the ground truth).

On the other hand, when the context-aware temperature is applied, the gradient can now be set to a much more substantial magnitude by setting a smaller temperature value  $\tau_1$  (as shown by the blue mesh). The additional flexibility of magnifying the gradient by temperature enables the model to be more efficient (or aggressive) in adjusting the parameters to reduce the training loss.

**Gradient of the temperature** Next, we analyze how the model with context-aware temperature updates the value of the temperature. We recall that Equation 4 defines the formulation of the

temperature. We leave the detailed derivation of the gradient of the temperature in Appendix A.1. Figure 4 visualizes the gradients  $\frac{\partial L}{\partial z_{\tau_0}}$  and  $\frac{\partial L}{\partial z_{\tau_1}}$  when  $z_0$  is either positive or negative. We will use the case in Figure 4(a) (where  $z_0 \in \mathbb{R}^-$ ) as an example to explain the effect of the gradient on the temperature. In this case, if  $p_0$  is close to 0, then we expect that the model should be more aggressive to update the parameters (since we assume that the ground-truth class is  $i = 0$ , meaning that  $p_0$  should ideally be close to 1). This aggressive response is indeed visible in the figure, showing values of larger magnitude when  $p_0 \rightarrow 0$ .

The amount of update on the temperature  $\tau_0$  also depends on its current value. In the same Figure 4(a), let’s consider a fixed  $p_0$ , say,  $p_0 = 0.1$ . In this case, as described above, the model wants to increase the value of  $p_0$  to closer to 1. To do that, the model will attempt to increase the term  $e^{z_0/\tau_0}$ . Since  $\tau_0$  is a positive value in  $[0, 1]$  and  $z_0 < 0$  in this case, to maximize  $z_0/\tau_0$ , the optimal is to have the temperature value  $\tau_0 \rightarrow 1$ . When the model is updating the value of  $\tau_0$ , if its current value is already close to 1, then the gradient will be small, since it is already close to the optimal value. On the other hand, if the temperature  $\tau_0$  is still far from 1, then the gradient will update it’s value more aggressively. This behavior is exactly visualized in the figure.

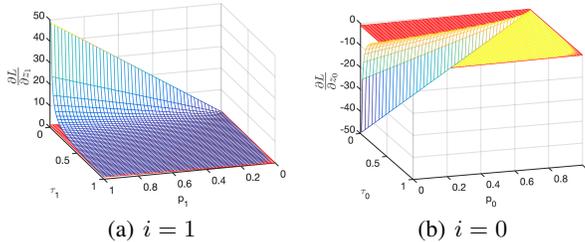


Figure 3: Gradients of loss with respect to (a) logit  $z_0$  and (b) logit  $z_1$ . In each figure, the x-axis is the probability  $p_i$  of class  $i$ , y-axis is the temperature value  $\tau_i$  of class  $i$ , and z-axis is the gradient  $\frac{\partial L}{\partial z_i}$ . The colorful mesh represents the gradients when the context-aware temperature is applied, while the red mesh represents the case without temperature.

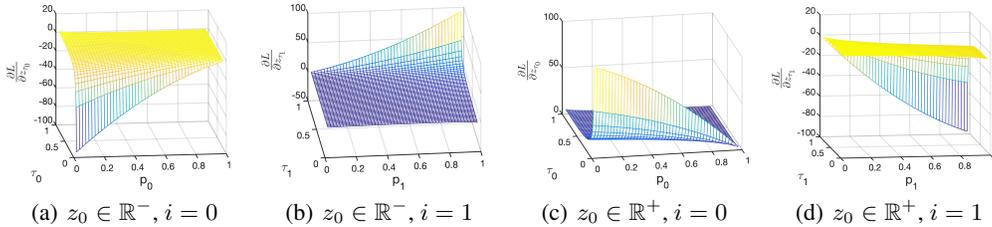


Figure 4: Gradients of loss with respect to  $z_{\tau_i}$ . The x-axis is the probability  $p_i$  of class  $i$ , y-axis is the temperature value  $\tau_i$  of class  $i$ , and z-axis is the gradient  $\frac{\partial L}{\partial z_{\tau_i}}$ .

## 4 EXPERIMENTS

### 4.1 EXPERIMENTAL SETUPS

We evaluate context-aware temperature on three datasets for language modeling: Penn Treebank, WikiText-2 and WikiText-103 (Marcus et al., 1993; Mikolov et al., 2011; Merity et al., 2017).

**CT-MoS.** We conduct experiments on PTB and WT2 using one and four 1080 Ti GPUs, respectively. The environment is PyTorch (Paszke et al., 2017). We follow the training configurations as reported in the MoS paper and their github<sup>2</sup>. For both data sets, we use the same number of param-

<sup>2</sup><https://github.com/zihangdai/mos>

eters as MoS. That is, we use three layers of LSTM with embedding sizes of 960-960-620 on PTB; and for WT2, three layers of LSTM with embedding sizes of 1150-1150-650. The only difference is on the choice of the optimizer, where we use the Adam optimizer with a learning rate of  $1e^{-4}$ .

**CT Transformer-XL.** The experiments on WikiText-103 is conducted with four 1080 Ti GPUs. The baseline model is Transformer-XL, and the environment we use is PyTorch. We follow the training configurations as reported in the official Transformer-XL github<sup>3</sup>. There are two training configurations: standard and large. We adopt the standard one, which has 16 layers, 10 attention heads and the hidden size of 410. Adaptive input representations (Baevski & Auli, 2018) and adaptive softmax (Grave et al., 2016) are used.

## 4.2 MAIN RESULTS

**Penn Treebank.** We compare the proposed method with MoS and AWD-LSTM, including the performance with and without dynamic evaluation (Krause et al., 2017). Since the original MoS model has approximately 22M parameters, to make a fair comparison, we augment its number of parameters to have 24M parameters. We increase the size of each layer proportionally, giving the word embedding size  $d = 300$ , and making the sizes of the three LSTM layers be 1030, 1030 and 670. We denote this augmented model as MoS+.

In Table 1, experimental results show that our CT-MoS model outperforms AWD-LSTM, MoS, and MoS+ on both validation and test sets. Our model achieves 22.92 perplexity on the validation set and 19.90 on the test set without dynamic evaluation. When using dynamic evaluation, the proposed CT-MoS model also achieves significant better perplexities of 10.63 and 10.04.

Table 1: Perplexity comparison on the Penn Treebank dataset. † indicates using dynamic evaluation.

Model	Paras	Validation	Test
AWD-LSTM	24M	60.7	58.8
AWD-LSTM †	24M	51.6	51.1
MoS	22M	58.08	55.97
MoS †	22M	48.33	47.69
MoS+	24M	57.96	55.78
MoS+ †	24M	48.89	48.28
CT-MoS	24M	<b>22.92</b>	<b>19.90</b>
CT-MoS †	24M	<b>10.63</b>	<b>10.04</b>

**WikiText-2.** Table 2 presents our experimental results for WT2. Similar to the results in Table 1, CT-MoS also outperforms the state-of-the-art models in this case. Compared to the MoS model, CT-MoS achieves great improvements, with perplexity 34.81 and 33.88 on the validation and test sets, respectively. With dynamic evaluation, CT-MoS obtains even better perplexity of 14.51 and 13.93. To conduct a fair comparison, we enlarge the number of parameters of MoS to have 45M parameters. The size of each layer is increased proportionally: the embedding size is  $d = 360$  and sizes of the three LSTM layers are 1320-1320-760. We denote the enlarged model as WT2+.

**WikiText-103.** Table 3 shows the experimental results on the WikiText-103 data set. We report the performances of both standard and large Transformer-XL models. The better results achieved by CT-Transformer-XL not only indicates the effectiveness of context-aware temperature, but also demonstrates its applicability on modeling the long-term dependency in an article-based corpus.

## 4.3 ABLATION STUDIES

To further illustrate the effects of the context-aware temperature, we compare it with conventional temperature scaling method, that is, using a constant temperature. We experiment with different

<sup>3</sup><https://github.com/kimiyoungh/transformer-xl>

Table 2: Perplexity comparison on the WikiText-2 dataset. † indicates using dynamic evaluation.

Model	Paras	Validation	Test
AWD-LSTM	33M	69.1	66.0
AWD-LSTM †	33M	46.4	44.3
MoS	35M	66.01	63.33
MoS †	35M	42.41	40.68
MoS+	45M	65.33	62.66
MoS+ †	45M	42.73	40.74
CT-MoS	45M	<b>34.81</b>	<b>33.88</b>
CT-MoS †	45M	<b>14.51</b>	<b>13.93</b>

Table 3: Perplexity comparison on the WikiText-103 dataset.

Model	Paras	Validation	Test
Transformer XL Standard	151M	24.0	23.1
Transformer-XL Large	257M	18.3	18.2
CT Transformer-XL	261M	<b>4.1</b>	<b>4.7</b>

constant values  $\{0.5, 1, 2, 4\}$  on the MoS model. Results in Table 4 show that employing context-aware temperature provides better performance than the conventional method. This indicates that the proposed method indeed has merits on adjusting the model parameters, by considering the context and providing a dynamic and optimized temperature value for each token.

Table 4: Effects of context-awareness. Experiments are conducted on the PTB dataset.

Model	Validation	Test
MoS ( $\tau = 0.5$ )	60.73	58.38
MoS ( $\tau = 1.0$ )	58.08	55.97
MoS ( $\tau = 2.0$ )	57.25	55.11
MoS ( $\tau = 4.0$ )	57.39	55.21
CT-MoS	<b>22.92</b>	<b>19.90</b>

Furthermore, we conduct two ablation studies regarding (a) the model size, and (b) the value range of the temperature vector. Details are left in Appendix A.2 and A.3, respectively.

#### 4.4 ANALYSIS

**Language Uncertainties** Another interesting observation of the proposed context-aware temperature is the correlation between a word’s relative position in a sentence and its learned temperature. Figure 5 shows the means of temperature vectors at different positions in a sentence. We analyze sentences of three different lengths. For instance, the purple line corresponds to our analysis on sentences that are longer than 21 words and shorter than 35 words. For each of these sentences, we consider positions in three disjoint segments: the first 7, middle 7 and last 7 words. The three segments are concatenated to form a “normalized” sentence. This pre-processing ensures the positions of a token only have relative effects to the analysis, and are not effected by the sentence length.

In Figure 5, we observe that the temperature value is low at beginning positions of a sentence. As the position gets further away, the averaged temperature value first has a sharp increase and then decreases at positions near the end. Our intuition is as the following: at the beginning positions the model has little confidence, since there is limited information in the history. The model recognizes this fact and learns a uniform temperature vector over the tokens, to be less assertive and assigns relatively uniform probability over the tokens. As the history builds up at later positions, the model

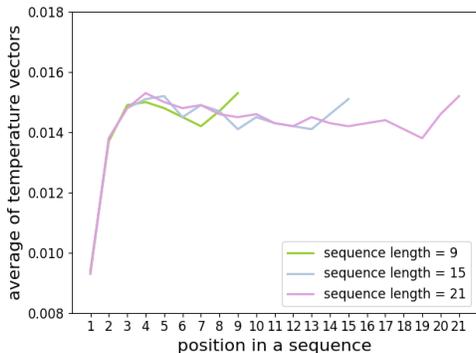


Figure 5: Means of the temperature vectors over positions in a sentence. The average temperature is low at the beginning of a sentence and gradually increases towards latter positions.

Table 5: Analysis of model performance on a sample from the PTB dataset.  $\tau$  denotes the temperature of a certain token. More examples can be found in Appendix A.4.

<b>Reference</b>	\$ N million of general obligation <b>veterans</b> ' <b>tax</b> notes series N via competitive bid		
<b>CT-MoS</b>	\$ N million of general obligation <b>veterans</b> ' <b>tax</b> notes series N via a bid		
<b>MoS</b>	\$ N million of general obligation <b>bonds</b> <b>bonds</b> <b>bonds</b> bonds series N via a bid		
<b>CT-MoS top-3</b>	<b>veterans 0.874</b>	tax 0.082	<unk> 0.034
<b>MoS top-3</b>	<b>bonds 0.495</b>	revenue 0.062	<b>veterans 0.031</b>
<b>Temperature</b>	<b>veterans</b> $\tau = 0.0011$		<b>bonds</b> $\tau = 0.0099$
<b>CT-MoS top-3</b>	<b>tax 0.943</b>	N 0.056	<unk> 0.001
<b>MoS top-3</b>	<b>bonds 0.135</b>	revenue 0.101	<b>tax 0.064</b>
<b>Temperature</b>	<b>tax</b> $\tau = 0.0063$		<b>bonds</b> $\tau = 0.0096$

becomes more confident about the next token and outputs a more spiky probability distribution on plausible tokens. The formation of the spiky distribution is done by having high temperature values for *implausible* tokens. These higher temperature values increases the average temperature value as shown in Figure 5. The uncertainties drop sharply as we reach the end of the sentence.

**Case Study: Effectiveness.** In Table 5, we present a sentence from the PTB dataset to illustrate the differences between MoS and CT-MoS, against the Reference. We highlight two differences in red and blue colors. At the location highlighted in red, we see that CT-MoS successfully predicts the answer “veterans” with a high probability of 0.874. On the other hand, MoS predicts “bonds”, which deviates from the ground truth. The word “veterans” has a temperature of 0.0011, which is much smaller than that of the word “bonds” (0.0099). This contributes to CT-MoS chooses “veterans” over “bonds”. Similar observations can be found at the blue location. We also observe that MoS tends to output a same word repeatedly, such as the four “bonds” tokens predicted for this example. The context-aware temperature mechanism does not have this issue. As illustrated in Figure 1, the context-aware temperature inclines to suppress the token appearing at previous time step by raising the temperature value, which effectively discouraging the same token to be predicted consecutively.

## 5 CONCLUSION AND FUTURE WORK

We propose a fully automated temperature mechanism, which learns an optimal temperature for each individual token based on the history of the context. The proposed *context-aware temperature* obtains strong results on various widely-adopted language modeling datasets. Our work opens up the research directions along the line of fully automated temperature mechanism in various NLP tasks, such as summarization, machine translation, and dialogue generation.

## REFERENCES

- Alexei Baevski and Michael Auli. Adaptive input representations for neural language modeling. *arXiv preprint arXiv:1809.10853*, 2018.
- Massimo Caccia, Lucas Caccia, William Fedus, Hugo Larochelle, Joelle Pineau, and Laurent Charlin. Language gans falling short. In *NIPS Workshop*, 2018.
- Zihang Dai, Zhilin Yang, Yiming Yang, Jaime Carbonell, Quoc V. Le, and Ruslan Salakhutdinov. Transformer-xl: Attentive language models beyond a fixed-length context. In *ACL*, 2019.
- Yann N Dauphin, Angela Fan, Michael Auli, and David Grangier. Language modeling with gated convolutional networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pp. 933–941. JMLR, 2017.
- Edouard Grave, Armand Joulin, Moustapha Cissé, David Grangier, and Hervé Jégou. Efficient softmax approximation for gpus. *arXiv preprint arXiv:1609.04309*, 2016.
- Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q. Weinberger. On calibration of modern neural networks. In *ICML*, 2017.
- Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. In *NIPS*, 2014.
- Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8): 1735–1780, 1997.
- Zhiting Hu, Zichao Yang, Xiaodan Liang, Ruslan Salakhutdinov, and Eric P. Xing. Toward controlled generation of text. In *ICML*, 2017.
- Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax. In *ICLR*, 2017.
- Ben Krause, Emmanuel Kahembwe, Iain Murray, and Steve Renals. Dynamic evaluation of neural sequence models. *arXiv preprint arXiv:1709.07432*, 2017.
- Junyang Lin, Xu Sun, Xuancheng Ren, Muyu Li, and Qi Su. Learning when to concentrate or divert attention: Self-adaptive attention temperature for neural machine translation. In *EMNLP*, 2018.
- Xuezhe Ma, Pengcheng Yin, Jingzhou Liu, Graham Neubig, and Eduard Hovy. Softmax q-distribution estimation for structured prediction: A theoretical interpretation for raml. *arXiv preprint arXiv:1705.07136*, 2017.
- Mitchell P Marcus, Mary Ann Marcinkiewicz, , and Beatrice Santorini. Building a large annotated corpus of english: The penn treebank. In *Computational linguistics*, 1993.
- Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture models. In *ICLR*, 2017.
- Stephen Merity, Nitish Shirish Keskar, and Richard Socher. Regularizing and optimizing lstm language models. In *ICLR*, 2018.
- Tomáš Mikolov, Anoop Deoras, Stefan Kombrink, Lukáš Burget, and Jan Černocký. Empirical evaluation and combination of advanced language modeling techniques. In *INTERSPEECH*, 2011.
- Mohammad Norouzi, Samy Bengio, Zhifeng Chen, Navdeep Jaitly, Mike Schuster, Yonghui Wu, and Dale Schuurmans. Reward augmented maximum likelihood for neural structured prediction. In *NIPS*, 2016.
- Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in PyTorch. In *NIPS Autodiff Workshop*, 2017.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. *OpenAI Blog*, 2019.

Mohammad Shoeybi, Mostofa Patwary, Raul Puri, Patrick LeGresley, Jared Casper, and Bryan Catanzaro. Megatron-lm: Training multi-billion parameter language models using model parallelism. *arXiv preprint arXiv:1909.08053*, 2019.

Sirui Xie, Hehui Zheng, Chunxiao Liu, and Liang Lin. Snas: Stochastic neural architecture search. In *ICLR*, 2019.

Zhilin Yang, Zihang Dai, Ruslan Salakhutdinov, and William W. Cohen. Breaking the softmax bottleneck: A high-rank RNN language model. In *ICLR*, 2018.

Xu Zhang, Felix Xinnan Yu, Svebor Karaman, Wei Zhang, and Shih-Fu Chang. Heated-up softmax embedding. *arXiv preprint arXiv:1809.04157*, 2018.

## A APPENDIX

### A.1 PARTIAL DERIVATIVES OF LOSS TO LOGITS OF TEMPERATURE

Take the case of two classes as example, assume that the ground-truth class is  $i = 0$ . In this case, the loss  $L$  is  $-\ln p_0$ , where  $p_0$  is the probability of class 0 output by the model. The probabilities of the two classes are  $\mathbf{p} = \sigma(\mathbf{z} \odot \boldsymbol{\tau}) = [p_0, p_1]$ . Let  $\mathbf{u} = \mathbf{z} \odot \boldsymbol{\tau} = [u_0, u_1]$ , and  $u_i = z_i/\tau_i$ . Then, we have

$$p_i = \frac{e^{u_i}}{e^{u_0} + e^{u_1}}. \quad (7)$$

$\boldsymbol{\tau}$  is defined in Equation 4, and is essentially  $\sigma(\mathbf{z}_\tau)$ . Therefore,

$$\tau_i = \frac{e^{z_{\tau_i}}}{e^{z_{\tau_0}} + e^{z_{\tau_1}}}. \quad (8)$$

The gradients of the loss with respect to logits  $z_0$  and  $z_1$  are

$$\frac{\partial L}{\partial z_i} = \frac{\partial(-\ln p_0)}{\partial z_i} = \begin{cases} (p_i - 1)\frac{1}{\tau_i} & i = 0 \\ p_i\frac{1}{\tau_i} & i \neq 0 \end{cases} \quad (9)$$

The gradient of  $z_{\tau_0}$  is calculated as below

$$\begin{aligned} \frac{\partial L}{\partial z_{\tau_0}} &= \frac{\partial L}{\partial p_0} \frac{\partial p_0}{\partial u_0} \frac{\partial u_0}{\partial \tau_0} \frac{\partial \tau_0}{\partial z_{\tau_0}} + \frac{\partial L}{\partial p_0} \frac{\partial p_0}{\partial u_1} \frac{\partial u_1}{\partial \tau_1} \frac{\partial \tau_1}{\partial z_{\tau_0}} \\ &= -p_1 \frac{-z_0}{\tau_0^2} \frac{e^{z_{\tau_0} + z_{\tau_1}}}{(e^{z_{\tau_0}} + e^{z_{\tau_1}})^2} + p_1 \frac{-z_1}{\tau_1^2} \frac{-e^{z_{\tau_0} + z_{\tau_1}}}{(e^{z_{\tau_0}} + e^{z_{\tau_1}})^2} \\ &= \frac{1}{\tau_0} p_1 z_0 \tau_1 + \frac{1}{\tau_1} p_1 z_1 \tau_0 \end{aligned} \quad (10)$$

Similarly, the gradient of  $z_{\tau_1}$  is calculated as

$$\begin{aligned} \frac{\partial L}{\partial z_{\tau_1}} &= \frac{\partial L}{\partial p_0} \frac{\partial p_0}{\partial u_0} \frac{\partial u_0}{\partial \tau_0} \frac{\partial \tau_0}{\partial z_{\tau_1}} + \frac{\partial L}{\partial p_0} \frac{\partial p_0}{\partial u_1} \frac{\partial u_1}{\partial \tau_1} \frac{\partial \tau_1}{\partial z_{\tau_1}} \\ &= -p_1 \frac{-z_0}{\tau_0^2} \frac{-e^{z_{\tau_0} + z_{\tau_1}}}{(e^{z_{\tau_0}} + e^{z_{\tau_1}})^2} + p_1 \frac{-z_1}{\tau_1^2} \frac{e^{z_{\tau_0} + z_{\tau_1}}}{(e^{z_{\tau_0}} + e^{z_{\tau_1}})^2} \\ &= -\frac{1}{\tau_0} p_1 z_0 \tau_1 - \frac{1}{\tau_1} p_1 z_1 \tau_0 \end{aligned} \quad (11)$$

### A.2 MODEL SIZE

We have demonstrated the results of MoS+ in both Table 1 and 2. Here we notice that MoS+ has a similar perplexity compared to MoS, indicating that directly increasing model parameters cannot improve the performance. Similar observation and results are also reported by Yang et al. (2018). This shows that the improvements brought by CT-MoS are more than the mere growth of parameters.

As for Transformer-XL, two results achieved by different sizes of the model are reported. We compare with the large one (257M parameters), which has a comparable parameter size to our model (261M parameters). Results in Table 6 demonstrate the significant improvements brought by the proposed method. Again, this shows that the mere growth of parameters do not account for the improvements. In addition, we report the performance of several state-of-the-art models, including Megatron-LM (Shoeybi et al., 2019) and the widely-known GPT-2 (Radford et al., 2019) models. Please note that these models use extra training data during the training procedure. However, the proposed model outperforms these baseline models without using any extra data.

Table 6: Perplexity comparison on the WikiText-103 dataset.

Model	Paras	Validation	Test	extra training data
Transformer XL Standard	151M	24.0	23.1	$\times$
Transformer-XL Large	257M	18.3	18.2	$\times$
CT Transformer-XL	261M	<b>4.1</b>	<b>4.7</b>	$\times$
Megatron-LM	8300M	-	10.81	$\checkmark$
GPT-2 Full	1542M	-	17.48	$\checkmark$
GPT-2 Medium	355M	-	26.37	$\checkmark$

### A.3 TEMPERATURE NORMALIZATION RANGE

The value range of the temperature vector is bounded by a Softmax layer and two parameters  $(\alpha, \beta)$ . Throughout all the experiments, we let  $(\alpha, \beta)$  be learnable parameters and find their values to be  $(0, 1)$ . To further examine the learned value range, we conducted experiments with manually-defined temperature value ranges. The experiments are conducted on the Penn Treebank dataset.

Table 7 shows that the learned value ranges bring better performance than others. Even if we use a manually defined range that is same as the learned range  $(0, 1)$ , using the learned  $(\alpha, \beta)$  slightly outperforms. Furthermore, the proposed method achieves great improvements when the temperature range is either wide or limited. For instance, using  $(\alpha, \beta) = (1, 2)$  gives the range of  $(2, 4)$  and the test perplexity of 41.09; using  $(\alpha, \beta) = (0.2, 5)$  gives the range of  $(1, 6)$  and the test perplexity of 22.12. Both results are obviously better than the test perplexity (47.69) of baseline MoS model.

Table 7: Performance of different value ranges of the temperature vector.

$(\alpha, \beta)$	Validation	Test
$(0, 1)$	23.04	20.08
$(0.01, 2)$	23.12	20.09
$(0.2, 5)$	24.47	22.12
$(1, 3)$	43.97	41.09
$(1, 2)$	44.71	41.42
learned	<b>22.92</b>	<b>19.90</b>

### A.4 CASE STUDIES: EFFECTIVENESS

We present two other cases on the Penn Treebank dataset to visualize the effect of context-aware temperature, in Table 8 and 9. We highlight the five differences in red, brown, purple, blue and teal colors, respectively. Let’s first see Table 8. At the colored locations, we see that the proposed model successfully predicts the correct words, while MoS fails to do so. When taking a look at the temperature value for the specified word, we see that CT-MoS gives the ground-truth tokens smaller temperature values, whereas giving the non ground-truth tokens (such as “the”, “debt”, “its” and “in”) high temperature values. This contributes to the results that our model chooses the correct tokens over the wrong ones. Similar observations can be found in Table 9.

Table 8: Analysis of model performance on a sample from the PTB dataset.

<b>Reference</b>	<b>in august</b> resorts international inc. which sold more than \$ N million of <b>junk</b> bonds suspended <b>interest payments</b>		
<b>CT-MoS</b>	<b>in august</b> the international inc. filed filed \$ than \$ N million of <b>junk</b> bonds under <b>interest payments</b>		
<b>MoS</b>	<b>the the</b> the had inc. said has N than N N billion of <b>debt</b> bonds to <b>its in</b>		
<b>CT-MoS top-3</b>	<b>in 0.413</b>	<b>the 0.163</b>	more 0.094
<b>MoS top-3</b>	<b>the 0.322</b>	<b>in 0.076</b>	a 0.044
<b>Temperature</b>		<b>in</b> $\tau = 0.0094$	<b>the</b> $\tau = 0.0122$
<b>CT-MoS top-3</b>	<b>august 0.418</b>	september 0.246	<b>the 0.153</b>
<b>MoS top-3</b>	<b>the 0.299</b>	addition 0.087	a 0.083
<b>Temperature</b>		<b>august</b> $\tau = 0.0043$	<b>the</b> $\tau = 0.0111$
<b>CT-MoS top-3</b>	<b>junk 0.677</b>	u.s. 0.107	the 0.066
<b>MoS top-3</b>	<b>debt 0.150</b>	assets 0.147	high-yield 0.099
<b>Temperature</b>		<b>junk</b> $\tau = 0.0056$	<b>debt</b> $\tau = 0.0096$
<b>CT-MoS top-3</b>	<b>interest 0.323</b>	<b>its 0.212</b>	from 0.099
<b>MoS top-3</b>	<b>its 0.273</b>	the 0.166	a 0.096
<b>Temperature</b>		<b>interest</b> $\tau = 0.0046$	<b>its</b> $\tau = 0.0068$
<b>CT-MoS top-3</b>	<b>payments 0.800</b>	<b>in 0.053</b>	rates 0.042
<b>MoS top-3</b>	<b>in 0.270</b>	<b>payments 0.240</b>	rates 0.072
<b>Temperature</b>		<b>payments</b> $\tau = 0.0038$	<b>in</b> $\tau = 0.0095$

One thing worth noticing is that, as mentioned in Section 4.4, the MoS model tends to predict a same word repeatedly. In Table 8, the MoS model consecutively predicts “the the the”. This is again in accordance with our observations that the proposed CT-MoS model tends to give the previous token a high temperature value, so as to suppress the corresponding probability.

#### A.5 CASE STUDY: TEMPERATURE AND TOKEN POSITION.

Another aspect to examine how context-aware temperature works is to look at the change of the temperature from a specific token across different positions in a sentence. In Table 10, we present a sample from the PTB dataset, and highlight the occurrences of the word “mortgage” in red. As the position changes, the proposed method chooses a different temperature value, adjusting its confidence of the model’s belief. In this case, the temperature at the third occurrence of “mortgage” is 0.0023, and gradually decreases at subsequent occurrences. Such a decrease indicates that the model gains more confidence in making the prediction, most likely due to richer information from the longer history context.

Table 9: Analysis of model performance on a sample from the PTB dataset.

<b>Reference</b>	imperial <b>corp.</b> based in <b>san</b> diego is the <b>parent</b> of <b>imperial</b> savings & loan		
<b>CT-MoS</b>	imperial <b>corp.</b> based in <b>san</b> diego said attempting <b>parent</b> of <b>imperial</b> corp. & loan		
<b>MoS</b>	the <b>said</b> <b>said</b> in <b>&lt;unk&gt;</b> francisco said a <b>largest</b> of <b>the</b> bank & loan		
<b>CT-MoS top-3</b>	<b>corp.</b> 0.443	is 0.272	savings 0.104
<b>MoS top-3</b>	<b>said</b> 0.189	's 0.071	savings 0.067
<b>Temperature</b>	<b>corp.</b> $\tau = 0.0060$		<b>said</b> $\tau = 0.0073$
<b>CT-MoS top-3</b>	<b>based</b> 0.370	is 0.219	<b>said</b> 0.207
<b>MoS top-3</b>	<b>said</b> 0.333	a 0.071	's 0.070
<b>Temperature</b>	<b>based</b> $\tau = 0.0046$		<b>said</b> $\tau = 0.0089$
<b>CT-MoS top-3</b>	<b>san</b> 0.718	<b>&lt;unk&gt;</b> 0.171	imperial 0.029
<b>MoS top-3</b>	<b>&lt;unk&gt;</b> 0.270	new 0.097	los 0.070
<b>Temperature</b>	<b>san</b> $\tau = 0.0065$		<b>&lt;unk&gt;</b> $\tau = 0.103$
<b>CT-MoS top-3</b>	<b>parent</b> 0.517	first 0.364	<b>&lt;unk&gt;</b> 0.077
<b>MoS top-3</b>	<b>largest</b> 0.161	<b>&lt;unk&gt;</b> 0.097	parent 0.072
<b>Temperature</b>	<b>parent</b> $\tau = 0.0040$		<b>largest</b> $\tau = 0.0107$
<b>CT-MoS top-3</b>	<b>imperial</b> 0.750	<b>the</b> 0.124	<b>&lt;unk&gt;</b> 0.093
<b>MoS top-3</b>	<b>the</b> 0.150	<b>&lt;unk&gt;</b> 0.132	american 0.075
<b>Temperature</b>	<b>imperial</b> $\tau = 0.0012$		<b>the</b> $\tau = 0.0113$

Table 10: Analysis of the temperature for one specific word but at different positions.

<b>CT-MoS</b>	N N standard conventional fixed-rate <b>mortgages(1)</b> N N N N rate capped one-year adjustable rate <b>mortgages(2)</b> <b>&lt;eos&gt;</b> federal national <b>mortgage(3)</b> association fannie mae <b>&lt;eos&gt;</b> posted yields on N year <b>mortgage(4)</b> commitments for delivery within N days priced at par <b>&lt;eos&gt;</b> N N standard conventional fixed-rate <b>mortgages(5)</b> N N N rate capped one-year adjustable rate <b>mortgages(6)</b> <b>&lt;eos&gt;</b>					
<b>Temperature <math>\tau</math></b>	(1) 0.0020	(2) 0.0009	(3) 0.0023	(4) 0.0019	(5) 0.0010	(6) 0.0010