

Multi-Head Self-Attention with Role-Guided Masks

Dongsheng Wang*, Casper Hansen*, Lucas Chaves Lima, Christian Hansen, Maria Maistro, Jakob Grue Simonsen, and Christina Lioma

Department of Computer Science, University of Copenhagen
{wang,c.hansen,lcl,chrh,mm,simonsen,c.lioma}@di.ku.dk

Abstract. The state of the art in learning meaningful semantic representations of words is the Transformer model and its attention mechanisms. Simply put, the attention mechanisms learn to attend to specific parts of the input dispensing recurrence and convolutions. While some of the learned attention heads have been found to play linguistically interpretable roles, they can be redundant or prone to errors. We propose a method to guide the attention heads towards roles identified in prior work as important. We do this by defining role-specific masks to constrain the heads to attend to specific parts of the input, such that different heads are designed to play different roles. Experiments on text classification and machine translation using 7 different datasets show that our method outperforms competitive attention-based, CNN, and RNN baselines.

Keywords: Self-Attention · Transformer · Text Classification

1 Introduction

The Transformer model has had great success in various tasks in Natural Language Processing (NLP). For instance, the state of the art is dominated by models such as BERT [5] and its extensions: RoBERTa [12], ALBERT [9], SpanBERT [8], SemBERT [24], and SciBERT [2], all of which are Transformer-based architectures. Due to this, recent studies have focused on developing approaches to understand how attention heads digest input texts, aiming to increase the interpretability of the model [4,14,20]. The findings of those analyses are aligned: while some attention heads of the Transformer often play linguistically interpretable roles [4,20], others are found to be less important and can be pruned without significantly impacting (indicating redundancy), or even improving (indicating potential errors contained in pruned heads), effectiveness [14,20].

While the above studies show that the effectiveness of the attention heads is, in part, derived from different head roles, only scant prior work analyze the impact of *explicitly* adopting roles for the multiple heads. Such an explicit guidance would force the heads to spread the attention on different parts of the input with the aim of reducing redundancy. This motivates the following research question: *What is the impact of explicitly guiding attention heads?*

To answer this question, we define role-specific masks to guide the attention heads to attend to different parts of the input, such that different heads are designed to play

* Equal contribution

different roles. We first choose important roles based on findings from recent studies on interpretable Transformers roles; then we produce masks with respect to those roles; and finally the masks are incorporated into self-attention heads to guide the attention computation. Experimental results on both text classification and machine translation on 7 different datasets show that our approach outperforms competitive attention-based, CNN, and RNN baselines.

2 Related Work

The Transformer [19] was originally proposed as an encoder-decoder model, but has also been used successfully for transfer learning tasks, especially after being pre-trained on massive amounts of unlabeled texts. At the heart of the transformer lies the notion of multi-head self-attention, where the attention of each head is computed as:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (1)$$

where Q , is the query, K is the key, V is the value, and d_k is the key dimension. The input to each head is a head-specific linear projection, and the Transformer uses multi-heads such that the attention for each head is concatenated for a single output.

Recently, efforts have been made to explore how the Transformer attends over different parts of the input texts [4,7,20]. Clark et al. [4] investigate each attention head’s linguistic roles, and find that particular heads refer to specific aspects of syntax. Voita et al. [20] study the importance of the different heads using layer-wise relevance propagation (LRP) [6], and characterize them based on the role they perform. Furthermore, Voita et al. [20] find that not all heads are equally important and choose to prune the heads using a L_0 regularizer, finding that most of the non-pruned heads have specialized roles.

Scant prior work exists on guiding the attention heads to have a specific purpose. Strubell et al. [18] train the multi-head model with the first head attending to a single syntactic parent token, while the rest being regular attention heads. In contrast, we explore multiple more complex predefined roles grounded in head roles discovered in recent work. Sennrich and Haddow [15] incorporate linguistic features (e.g. sub-word tags, POS tags, etc.) as additional features into an attention encoder and decoder model for the task of machine translation, in order to enrich the model. In contrast, our method also makes use of linguistic features, but instead of enriching the input, we use these linguistic features to define the role-specific masks for guiding the attention heads.

3 Multi-head Attention with Guided Masks

We incorporate role-specific masks for self-attention heads, constraining them to attend to specific parts of the input. By doing this, we aim to reduce the redundancy between the heads, and force the heads to have roles identified in previous work as important. Then, we adopt a weighted gate layer to aggregate the heads.

We first define the multi-head self-attention with role-specific masks in Section 3.1 followed by a description of each role in Section 3.2. We denote our final attention guided Transformer model as Transformer-Guided-Attn.

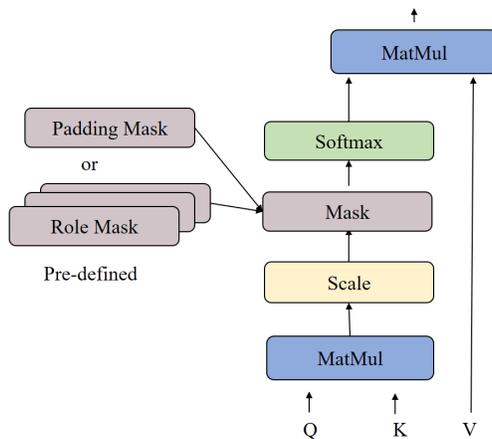


Fig. 1. Scaled-dot product with role mask or padding mask.

3.1 Multi-head attention

We incorporate a role-specific mask into a masked attention head (mh) as:

$$\text{mh}(Q, K, V, M_r) = \text{softmax} \left(\frac{QK^T + M_r}{\sqrt{d_k}} \right) V \quad (2)$$

where M_r is a role-specific mask used to constrain the attention head. For an input of length n , M_r is an n -by- n matrix where each element is either $-\infty$ (ignore) or 0 (include). For multi-head self-attention, we introduce N role-specific masks for the first N heads out of a total of H heads ($N \leq H$). If N is strictly less than H , then the remaining heads are regular attention heads. Based on this, the multi-head attention can be expressed as:

$$\text{MultiHead}(Q, K, V) = \text{Concat}(mh_1, mh_2, \dots, mh_N, h_{N+1}, h_{N+2}, \dots, h_H) W^O \quad (3)$$

where mh_i is the head with a role mask, and h_i is a regular head computed using Eq. (1).

A visualization of using the masks is shown in Figure 1, where we associate the standard padding mask to regular attention heads. The padding masks ensure that inputs shorter than the model allowed length are padded to fit the model.

3.2 Mask roles

We adopt the roles detected as important by Voita et al. [20] and Clark et al. [4]. We categorize them as 1) specialized (rare words and separators), 2) syntactic (dependency syntax and major relations), and 3) window (relative position) roles (see [10,11] for a linguistic basis of this categorisation). We include the separator role as Clark et al. [4] found that over half of BERT’s attention, in layer 6-10, focus on separators. We describe these 5 specific roles below, which are used for creating role-specific masks.

Rare words (RareW) The rare words role refers to the least frequent tokens in a text.

As defined by Voita et al. [20], we compute IDF (inversed document frequency)

scores for all tokens and use the 10 percent least frequent tokens (highest 10 percent values according to IDF) in the sentence as the target attentions.

Separator (Sep_{rat}) The separator role guides the head to point to only separators.

We extend the separator from $\{[SEP], [START], [END]\}$ to common punctuation of {comma, semicolon, dot, question mark, exclamation point}.

Dependency syntax (DepSyn) Dependency syntax role guides the head to attend to tokens with syntactic dependency relations. We assume this role can guide the head to attend to those—not adjacent—but still relevant tokens, complementary to the RelPos role (see below).

Major syntactic relations (MajRel) The major syntactic relations role guides the head to attend to the tokens involved with major syntactic relations. The four major relations defined by Voita et al. [20] are NSUBJ, DOBJ, AMOD, and ADVMOD.

Relative Position (RelPos) The relative position role guides the head to look at adjacent tokens, corresponding to scanning the text with a centered window of size 3.

For each role, we generate the guided mask for each input sentence by first producing an n -by- n matrix with all values as $-\infty$ (corresponding to ignoring all tokens initially). Then, we change the value of position (i, j) into 0.0, referring to the query token i with respect to the guided key token j , depending on the mask role.

4 Experiment

We experimentally compare our Transformer-Guided-Attn model to competitive baselines across 7 datasets in the tasks of text classification and machine translation. We make the source code publicly available on GitHub¹.

4.1 Classification tasks

We consider two different classification tasks: sentiment analysis and topic classification. We compare our methods against six competitive baselines: the original Transformer [19]; multi-scale CNNs [22]; RNNs (BiLSTM) [3]; directional Self-attention (DiSAN) [17] that incorporates temporal order and multi-dimensional attention into the Transformer; phrase-level self-attention (PSAN) [23] which performs self-attention across words inside a phrase; and Transformer-Complex-Order [21] that incorporates sequential order into the Transformer to capture ordered relationships between token positions. For the baselines implemented by us (marked in the Tables), we tune them as described in the original papers. For our Transformer-Guided-Attn, we consider a simple, but effective, way of selecting the combination of role-specific masks: For each layer, we fix 5 attention heads to be guided by the specific roles specified in Section 3.2, and let the remaining be regular heads. We tune the number of layers from $\{2, 4, 6, 8\}$ and number of additional regular heads from $\{1, 3\}$.

Dataset. The statistics of the datasets are shown in Table 1. We use the same splits as done by Wang et al. [21].

¹ <https://github.com/dswang2011/guided-attention-transformer>

Table 1. Classification dataset statistics. CV means 10-fold cross validation.

| Dataset | Train | Test | Task | Vocab. | Class |
|---------|-------|------|------------------|--------|-------|
| CR | 4k | CV | Product review | 6k | 2 |
| TREC | 5.4k | 0.5k | Question | 10k | 6 |
| SUBJ | 10k | CV | Subjectivity | 21k | 2 |
| MPQA | 11k | CV | Opinion polarity | 6k | 2 |
| MR | 11.9k | CV | Movie review | 20k | 2 |
| SST | 67k | 2.2k | Movie review | 18k | 2 |

Table 2. Machine translation results. * marks scores reported from other papers.

| Method | BLEU |
|----------------------------------|-------------|
| Transformer [19] | 34.3 |
| AED + Linguistic [15] * | 28.4 |
| AED + BPE [16] * | 34.2 |
| Tensorized Transformer [13] * | 34.9 |
| Transformer-Complex-Order [21] * | 35.8 |
| Transformer-Guided-Attn (ours) | 38.8 |

Table 3. Classification results (accuracy %). * marks scores reported from other papers.

| Method | CR | TREC | SUBJ | MPAQ | MR | SST |
|--|-------------|-------------|-------------|-------------|-------------|-------------|
| Transformer [19] | 82.0 | 91.8 | 93.2 | 88.6 | 77.7 | 81.8 |
| Multi-scale CNNs [22] | 81.2 | 93.1 | 93.3 | 89.1 | 77.8 | 80.9 |
| BiLSTM [3] | 82.6 | 92.4 | 93.6 | 88.9 | 78.4 | 81.1 |
| DiSAN (Directional Self-Attention) [17]* | 84.1 | 88.3 | 92.2 | 89.5 | 79.7 | 82.9 |
| PSAN (phrase-level Self-Attention) [23]* | 84.2 | 89.1 | 91.9 | 89.9 | 80.0 | 83.8 |
| Transformer-Complex-Order [21]* | 80.6 | 89.6 | 89.5 | 86.3 | 74.6 | 81.3 |
| Transformer-Guided-Attn (ours) | 84.4 | 93.6 | 93.8 | 90.7 | 80.0 | 84.2 |

Results. As shown in Table 3, we observe consistent improvements compared to the best baseline for each dataset, except on MR where we perform as well as PSAN. Compared to the original Transformer model, we obtain accuracy gains of up to 2.96%, depending on the dataset, thus showing a notable performance impact from guiding the attention heads. Compared to DiSAN and PSAN, our proposed Transformer-Guided-Attn obtains consistent improvements over the original Transformer across all datasets, while DiSAN and PSAN both have lower performance for TREC and SUBJ.

4.2 Translation Task

We use the standard WMT 2016 English-German dataset [16] and use four baselines: Attentional encoder-decoder (AED) [15] with linguistic features including morphological, part-of-speech, and syntactic dependency labels as additional embedding space; AED with Byte-pair encoding (BPE) [16] subword segmentation for open-vocabulary translation; the tensorized Transformer [13]; and the Transformer-Complex-order [21]. The first two models are extensions on top of the basic AED [1]. For the models we implement, we follow the same tuning as in the classification experiments. We evaluate the machine translation performance using the Bilingual Evaluation Understudy (BLEU) measure.

Results. Our Transformer-Guided-Attn consistently outperforms the competitive baselines. Specifically, we observe gains of 8.2% compared to the best baseline, Transformer-Complex-Order, and close to 13% compared to the original Transformer. These gains are even larger than the results for the classification experiments, thus highlight-

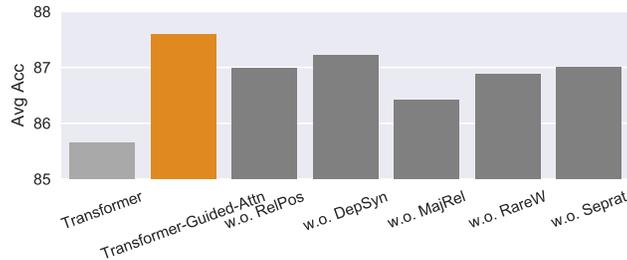


Fig. 2. Ablation study of Transformer-Guided-Attn when dropping each role individually.

ing a significant performance impact from guiding the attention heads for the task of machine translation.

4.3 Ablation study

We now consider the performance impact associated with each role-specific mask. For each classification dataset, we run configurations of our Transformer-Guided-Attn with each role-specific mask excluded once and replaced with a default padding mask used in the Transformer. The average accuracy drop associated with excluding each role-specific mask is shown in Figure 2, which also includes the average accuracy of the Transformer and our Transformer-Guided-Attn using all role-specific masks. We observe that the removal of each role has a negative impact on performance, where the major syntactic relations role (MajRel) has the largest impact. Thus, collectively all roles contribute to the performance of the full Transformer-Guided-Attn model.

5 Conclusion

We presented Transformer-Guided-Attn, a method to explicitly guide the attention heads of the Transformer using role-specific masks. The motivation of this explicit guidance is to force the heads to spread their attention on different parts of the input with the aim of reducing redundancy among the heads. Our experiments demonstrated that incorporating multiple role masks into multi-head attention can consistently improve performance on both classification and machine translation tasks.

As future work, we plan to explore additional roles for masking, as well as evaluating the impact of including it for pre-training language representation models such as BERT [5].

6 Acknowledgments

This work is supported by the European Union’s Horizon 2020 research and innovation programme under the Marie Skłodowska-Curie grant agreement No. 721321 (QUARTZ project) and No. 893667 (METER project).

References

1. Bahdanau, D., Cho, K., Bengio, Y.: Neural machine translation by jointly learning to align and translate. In: Bengio, Y., LeCun, Y. (eds.) 3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings (2015), <http://arxiv.org/abs/1409.0473>
2. Beltagy, I., Cohan, A., Lo, K.: Scibert: Pretrained contextualized embeddings for scientific text. CoRR **abs/1903.10676** (2019), <http://arxiv.org/abs/1903.10676>
3. Bin, Y., Yang, Y., Shen, F., Xu, X., Shen, H.T.: Bidirectional long-short term memory for video description. In: Proceedings of the 24th ACM international conference on Multimedia. pp. 436–440 (2016)
4. Clark, K., Khandelwal, U., Levy, O., Manning, C.D.: What does BERT look at? an analysis of BERT’s attention. arXiv preprint arXiv:1906.04341 (2019)
5. Devlin, J., Chang, M., Lee, K., Toutanova, K.: BERT: pre-training of deep bidirectional transformers for language understanding. In: Burstein, J., Doran, C., Solorio, T. (eds.) Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1. pp. 4171–4186. Association for Computational Linguistics (2019). <https://doi.org/10.18653/v1/n19-1423>, <https://doi.org/10.18653/v1/n19-1423>
6. Ding, Y., Liu, Y., Luan, H., Sun, M.: Visualizing and understanding neural machine translation. In: Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). pp. 1150–1159 (2017)
7. Hoover, B., Strobel, H., Gehrmann, S.: exbert: A visual analysis tool to explore learned representations in transformers models. arXiv preprint arXiv:1910.05276 (2019)
8. Joshi, M., Chen, D., Liu, Y., Weld, D.S., Zettlemoyer, L., Levy, O.: Spanbert: Improving pre-training by representing and predicting spans. *Transactions of the Association for Computational Linguistics* **8**, 64–77 (2020)
9. Lan, Z., Chen, M., Goodman, S., Gimpel, K., Sharma, P., Soricut, R.: Albert: A lite BERT for self-supervised learning of language representations. In: International Conference on Learning Representations (2020)
10. Lioma, C., Blanco, R.: Part of speech based term weighting for information retrieval. In: Boughanem, M., Berrut, C., Mothe, J., Soulé-Dupuy, C. (eds.) Advances in Information Retrieval, 31th European Conference on IR Research, ECIR 2009, Toulouse, France, April 6-9, 2009. Proceedings. Lecture Notes in Computer Science, vol. 5478, pp. 412–423. Springer (2009). https://doi.org/10.1007/978-3-642-00958-7_37, https://doi.org/10.1007/978-3-642-00958-7_37
11. Lioma, C., van Rijsbergen, C.J.K.: Part of speech n-grams and information retrieval. *French Review of Applied Linguistics, Special issue on Information Extraction and Linguistics XIII*(2008/1), 9–22 (2008), https://www.cairn-int.info/article-E_RFLA_131_0009--part-of-speech-n-grams-and-information.htm
12. Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., Stoyanov, V.: Roberta: A robustly optimized bert pretraining approach. arXiv preprint arXiv:1907.11692 (2019)
13. Ma, X., Zhang, P., Zhang, S., Duan, N., Hou, Y., Zhou, M., Song, D.: A tensorized transformer for language modeling. In: Advances in Neural Information Processing Systems. pp. 2229–2239 (2019)
14. Michel, P., Levy, O., Neubig, G.: Are sixteen heads really better than one? In: Advances in Neural Information Processing Systems. pp. 14014–14024 (2019)

15. Sennrich, R., Haddow, B.: Linguistic input features improve neural machine translation. In: Proceedings of the First Conference on Machine Translation, WMT 2016, colocated with ACL 2016, August 11-12, Berlin, Germany. pp. 83–91. The Association for Computer Linguistics (2016). <https://doi.org/10.18653/v1/w16-2209>, <https://doi.org/10.18653/v1/w16-2209>
16. Sennrich, R., Haddow, B., Birch, A.: Edinburgh neural machine translation systems for WMT 16. In: Proceedings of the First Conference on Machine Translation: Volume 2, Shared Task Papers. pp. 371–376. Association for Computational Linguistics, Berlin, Germany (Aug 2016). <https://doi.org/10.18653/v1/W16-2323>, <https://www.aclweb.org/anthology/W16-2323>
17. Shen, T., Zhou, T., Long, G., Jiang, J., Pan, S., Zhang, C.: Disan: Directional self-attention network for rnn/cnn-free language understanding. In: Thirty-Second AAAI Conference on Artificial Intelligence (2018)
18. Strubell, E., Verga, P., Andor, D., Weiss, D., McCallum, A.: Linguistically-informed self-attention for semantic role labeling. In: Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing. pp. 5027–5038. Association for Computational Linguistics, Brussels, Belgium (Oct-Nov 2018). <https://doi.org/10.18653/v1/D18-1548>, <https://www.aclweb.org/anthology/D18-1548>
19. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., Polosukhin, I.: Attention is all you need. In: Advances in neural information processing systems. pp. 5998–6008 (2017)
20. Voita, E., Talbot, D., Moiseev, F., Sennrich, R., Titov, I.: Analyzing multi-head self-attention: Specialized heads do the heavy lifting, the rest can be pruned. In: Korhonen, A., Traum, D.R., Màrquez, L. (eds.) Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers. pp. 5797–5808. Association for Computational Linguistics (2019). <https://doi.org/10.18653/v1/p19-1580>, <https://doi.org/10.18653/v1/p19-1580>
21. Wang, B., Zhao, D., Lioma, C., Li, Q., Zhang, P., Simonsen, J.G.: Encoding word order in complex embeddings. In: 8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020. OpenReview.net (2020), <https://openreview.net/forum?id=Hke-WTVtwr>
22. Wang, D., Simonsen, J.G., Larsen, B., Lioma, C.: The Copenhagen team participation in the factuality task of the competition of automatic identification and verification of claims in political debates of the clef-2018 fact checking lab. CLEF (Working Notes) **2125** (2018)
23. Wu, W., Wang, H., Liu, T., Ma, S.: Phrase-level self-attention networks for universal sentence encoding. In: Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing. pp. 3729–3738 (2018)
24. Zhang, Z., Wu, Y., Zhao, H., Li, Z., Zhang, S., Zhou, X., Zhou, X.: Semantics-aware BERT for language understanding. CoRR **abs/1909.02209** (2019), <http://arxiv.org/abs/1909.02209>