

Fine-Tuning Pre-Trained Language Models with Gaze Supervision

Anonymous ACL submission

Abstract

Human gaze data provide cognitive information that reflect human language comprehension and has been effectively integrated into a variety of natural language processing (NLP) tasks, demonstrating improved performance over corresponding plain text-based models. In this work, we propose to integrate a gaze module into pre-trained language models (PLMs) at the fine-tuning stage to improve their capabilities to learn representations that are grounded in human language processing. This is done by extending the conventional purely text-based fine-tuning objective with an auxiliary loss to exploit cognitive signals. The gaze module is only included during training, retaining compatibility with existing PLM-based pipelines. We evaluate the proposed approach using two distinct PLMs on the GLUE benchmark and observe that the proposed model improves performance compared to both standard fine-tuning and traditional text augmentation baselines. All code is available on [anonymous_git](#).

1 Introduction

As humans read text, the unconscious cognitive processes that unfold in their minds while comprehending the stimulus text are reflected in their gaze signals (Just and Carpenter, 1980). These gaze signals hold the potential to enhance NLP tasks. Research has focused on using aggregated word-level gaze features to enrich text features (Barrett et al., 2016; Mishra et al., 2016; Hollenstein and Zhang, 2019) or to regularize neural attention mechanisms, making their inductive bias more human-like (Barrett et al., 2018; Sood et al., 2020, 2021).

Moreover, there has been growing interest in adopting non-aggregated scanpaths (i.e., sequence of consecutive fixations) to augment LMs. These scanpaths capture the complete sequential ordering of a reader’s gaze behavior and approximate their attention. Mishra et al. (2017) and Khurana

et al. (2023) employed neural networks to independently encode scanpaths and text, followed by the fusion of the features extracted from both modalities. Yang and Hollenstein (2023) proposed rearranging the contextualized token embeddings produced by PLMs based on the order in which the reader fixates on the words, followed by performing sequence modeling on the reordered sequence. To tackle the issue of gaze data scarcity, Deng et al. (2023a) explored the possibility of augmenting LMs using synthetic scanpaths, generated by a scanpath generation model. Remarkably, synthetic scanpaths demonstrated advantages across various NLP tasks, particularly in settings with limited labeled examples for the downstream task.

In this work, we start from a different perspective and explore utilizing gaze data to improve on the learned representations of PLMs during the fine-tuning stage, without incurring additional computational effort when using the model at application time. To this end, we extend the standard PLM fine-tuning objective with an auxiliary loss by integrating a scanpath module, which serves a dual purpose. First, the auxiliary loss can effectively incorporate human-like gaze signals generated using a scanpath generation model and thus provide informative gradients to guide the LM towards more representative local minima. Second, reordering the token-embedding sequence based on the fixation sequence can diversify textual information, potentially improving generalization performance (Xie et al., 2020). This stands in contrast to heuristic text augmentation strategies, like random word insertion, replacement, swapping, and deletion (Wei and Zou, 2019; Xie et al., 2020). Scanpaths inherently contain cognitive information that better aligns with and complements textual information.

Notably, our proposed gaze module is only active during training (fine-tuning), ensuring alignment with the standard usage of LMs after this stage. This offers two key benefits. First, it fa-

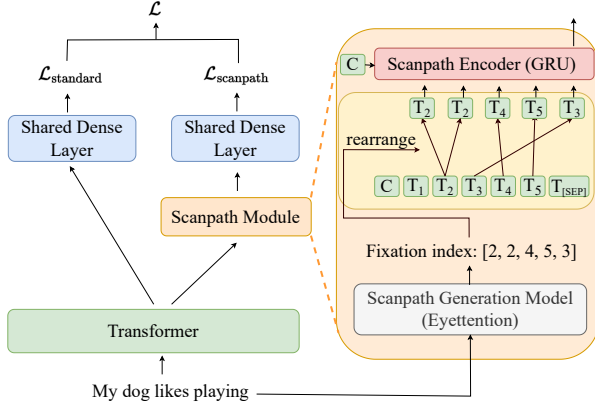


Figure 1: Overall architecture during training. The standard objective is augmented with an auxiliary loss from a scanpath-integrated branch, where token embeddings are rearranged based on the simulated fixation sequence.

cilitates seamless integration with existing LM-based pipelines. Second, at deployment time, it eliminates the need to either collect real-time gaze recordings, which is costly and impractical for most use-cases, or generate synthetic gaze data, which is often computationally challenging for devices with limited computational resources.

On the General Language Understanding Evaluation (GLUE) benchmark, our proposed gaze-augmented fine-tuning outperforms both standard text-only fine-tuning and traditional text augmentation baselines, without incurring additional computational effort at application time.

2 Method

In this section, we start out with a brief description of the conventional fine-tuning procedure for Transformer-based encoders on downstream tasks. Subsequently, we introduce our method, and explain how it incorporates synthetic scanpaths to enhance representation learning of Transformer-based encoders into this fine-tuning procedure. The overall model architecture is illustrated in Figure 1.

Preliminaries Our learning objective is to solve standard multi-class classification or regression problems. We assume access to a Transformer-based PLM like BERT (Devlin et al., 2019) or RoBERTa (Liu et al., 2019). In the conventional fine-tuning approach for downstream tasks, the PLM is adapted to a specific task by fine-tuning all the parameters end-to-end using task-specific inputs and outputs. The final hidden state of the “[CLS]” token typically serves as the aggregated

sentence representation, which is then fed into a newly initialized (series of) dense layer(s) with output neurons corresponding to the number of labels in the task. We minimize the standard cross-entropy loss for classification and mean-squared-error loss for regression, denoted $\mathcal{L}_{\text{standard}}$ in Figure 1.

Scanpath Integration We extend the standard fine-tuning framework by integrating a scanpath module. The design of the scanpath module follows the prior work of Deng et al. (2023a) and Yang and Hollenstein (2023). Specifically, the Transformer encoder produces contextualized token embeddings for a given sentence, with each embedding associated with its position index in the sequence. Simultaneously, a synthetic scanpath (fixation-index sequence) is generated based on the same sentence using the scanpath-generation model Eyettention (Deng et al., 2023b), which has demonstrated effectiveness in simulating human-like scanpaths during reading (see Appendix A for detailed information about the Eyettention model). The scanpath module then rearranges the token-embedding sequence based on the simulated fixation sequence. Subsequently, we use a scanpath encoder, implemented as a layer of Gated Recurrent Units (GRU), to process the reordered sequence. The output from the last step of the scanpath encoder is then forwarded to the subsequent dense layer. For the branch that takes the scanpath into account, we introduce an additional loss term, referred to as $\mathcal{L}_{\text{scanpath}}$ in Figure 1, which represents the cross-entropy loss for classification and the mean-squared-error loss for regression.

Training Objective We combine the standard purely text-based loss and the scanpath-integrated loss with a trade-off factor λ . The final training objective is defined as:

$$\mathcal{L} = \mathcal{L}_{\text{standard}} + \lambda \mathcal{L}_{\text{scanpath}}.$$

The joint optimization of the two branches facilitates the flow of cognitive information from the scanpath module to the Transformer through back-propagation, thereby improving its capability to process and comprehend text. Consequently, during testing, we can remove the scanpath module and generate predictions solely from the Transformer and the final dense layer. This ensures alignment with standard LM usage after the fine-tuning stage, notably preserving its intrinsic efficiency and compatibility.

K	Model	MNLI	QQP	QNLI	SST-2	CoLA	STS-B	MRPC	RTE	Avg
200	BERT	42.10 _{0.46}	62.16 _{1.30}	73.58 _{0.56}	77.68 _{1.71}	18.52 _{4.24}	80.48 _{0.32}	82.12 _{0.43}	54.95 _{0.67}	61.45
	+EDA	47.74 _{1.10}	64.89 _{0.56}	76.23 _{0.34}	80.48 _{1.26}	14.05 _{2.84} †	79.56 _{0.62} †	82.68 _{0.40}	55.74 _{0.30}	62.67
	+SP	42.63 _{0.82}	64.47 _{0.84}	73.83 _{0.44}	81.19 _{0.98}	23.33 _{3.42}	82.01 _{0.28}	82.71 _{0.48}	56.10 _{0.67}	63.28
500	BERT	52.35 _{1.23}	67.33 _{0.29}	77.78 _{0.46}	84.17 _{0.28}	30.29 _{1.86}	83.90 _{0.24}	83.15 _{0.26}	60.43 _{1.07}	67.43
	+EDA	56.37 _{0.88}	68.03 _{0.33}	78.48 _{0.32}	85.37 _{0.17}	28.89 _{1.58} †	83.28 _{0.24} †	84.00 _{0.28}	60.43 _{0.49}	68.11
	+SP	55.40 _{0.61}	67.86 _{0.42}	78.19 _{0.24}	84.22 _{0.52}	35.87 _{1.50}	85.26 _{0.29}	84.52 _{0.46}	61.44 _{0.43}	69.10
1000	BERT	60.51 _{0.66}	69.40 _{0.54}	79.53 _{0.16}	85.25 _{0.51}	39.92 _{0.86}	86.22 _{0.11}	85.42 _{0.23}	63.10 _{1.16}	71.17
	+EDA	61.58 _{0.50}	69.91 _{0.35}	80.49 _{0.16}	86.10 _{0.34}	31.04 _{1.89} †	85.50 _{0.22} †	86.37 _{0.44}	64.26 _{1.16}	70.66†
	+SP	61.75 _{0.32}	70.58 _{0.30}	80.24 _{0.33}	86.70 _{0.09}	42.45 _{0.59}	86.73 _{0.14}	86.77 _{0.69}	63.18 _{1.08}	72.3
200	RoBERTa	40.06 _{0.68}	68.59 _{0.54}	77.21 _{0.60}	88.56 _{0.39}	30.29 _{2.55}	82.84 _{0.43}	83.37 _{0.16}	55.81 _{1.15}	65.84
	+EDA	53.64 _{0.44}	68.84 _{0.71}	77.52 _{0.57}	87.94 _{0.64} †	23.30 _{4.16} †	83.86 _{0.10}	84.05 _{0.49}	58.41 _{1.20}	67.20
	+SP	44.90 _{0.63}	69.05 _{0.69}	78.14 _{0.68}	87.11 _{0.86} †	29.07 _{3.18} †	82.42 _{0.24} †	83.86 _{0.62}	63.03 _{2.58}	67.20
500	RoBERTa	65.20 _{0.46}	73.42 _{0.48}	81.54 _{0.22}	89.61 _{0.35}	39.59 _{0.95}	86.68 _{0.30}	86.09 _{0.36}	62.24 _{1.92}	73.05
	+EDA	64.97 _{0.56} †	71.57 _{0.45} †	81.20 _{0.23} †	89.27 _{0.35} †	36.05 _{2.28} †	86.46 _{0.26} †	87.49 _{0.67}	59.49 _{1.55} †	72.06†
	+SP	64.89 _{0.42} †	73.79 _{0.30}	81.78 _{0.16}	89.75 _{0.30}	39.07 _{1.96} †	86.29 _{0.07} †	87.00 _{0.54}	68.01 _{1.07}	73.82
1000	RoBERTa	70.91 _{0.61}	75.63 _{0.29}	83.43 _{0.12}	90.69 _{0.24}	44.78 _{0.65}	88.06 _{0.19}	88.85 _{0.19}	64.91 _{1.26}	75.91
	+EDA	70.84 _{0.34} †	74.59 _{0.52} †	82.64 _{0.47} †	90.23 _{0.38} †	41.44 _{1.18} †	87.79 _{0.15} †	89.60 _{0.41}	63.25 _{2.00} †	75.05†
	+SP	70.69 _{0.37} †	75.40 _{0.16} †	83.59 _{0.42}	89.91 _{0.35} †	44.43 _{1.88} †	88.12 _{0.17}	89.42 _{0.53}	72.71 _{0.73}	76.78

Table 1: Results on the GLUE benchmark with $K = \{200, 500, 1000\}$ training instances. QQP/MRPC: F1, STS-B: Spearman correlation, CoLA: Matthews correlation, and accuracy for the remaining tasks. We perform 5 runs and report the means along with standard errors. The dagger “†” indicates performance inferior to standard fine-tuning.

3 Experiments

3.1 Evaluation Setup

Data Sets We conduct experiments on the GLUE benchmark (Wang et al., 2018), including sentiment analysis (SST-2), linguistic acceptability (CoLA), similarity and paraphrase tasks (MRPC, STS-B, QQP), and natural language inference tasks (MNLI, QNLI, RTE).

Model and Data Setup We use BERT_{base} and RoBERTa_{base} as the base models in the experiments. We primarily focus on a low-resource setting where only limited labeled examples for the downstream task are available. In such cases, effective fine-tuning strategies are crucial to enable high-capacity LMs to learn more informative representations (Zhang et al., 2021). For each task, we sample a small subset of training instances with sizes $K = \{200, 500, 1000\}$. We take an additional 1,000 instances from the original training set as the development set and use the original development set for testing. Additionally, we consider a high-resource setting where we use the entire training set and report the results on the GLUE development sets. Appendix B gives further details about training and hyper-parameter tuning.

Baselines We compare our proposed method with the standard text-only fine-tuning ($\mathcal{L}_{\text{standard}}$ as the training objective). Moreover, we compare to the EDA method (Wei and Zou, 2019), which randomly inserts, replaces, swaps, and deletes words in the text to augment the training data.

3.2 Results

Low-Resource Performance Table 1 shows that, overall, our scanpath-augmented fine-tuning (+SP) consistently outperforms the standard fine-tuning and EDA baselines, regardless of the number of training instances. We observe performance gains of 2-3% for BERT and 1-2% for RoBERTa over standard fine-tuning. At the per-task level, our method outperforms standard fine-tuning across all tasks in all setups for BERT, and on 5, 5, and 4 out of 8 tasks when trained with 200, 500, and 1,000 instances, respectively, for RoBERTa. The improvements are larger with fewer training instances, indicating the efficacy of our method in low-resource scenarios. Notably, for tasks like CoLA and STS-B where the EDA method yields largely inferior results compared to standard fine-tuning (Model=BERT), our method shows superior performance. This suggests that the scanpath, which inherently contains cognitive information, aligns with and complements textual information effectively.

High-Resource Performance In Table 2, we present the results of different methods when using all training instances. Our scanpath-augmented fine-tuning (+SP) achieves the highest overall performance. While the gains are not as significant as in the low-resource setting for most tasks, notable improvements persist for tasks like CoLA and RTE. In contrast, the EDA method fails to outperform standard fine-tuning overall, aligning with findings from previous research (Longpre et al., 2020).

Model	MNLI	QQP	QNLI	SST-2	CoLA	STS-B	MRPC	RTE	Avg
BERT	83.87	88.02	91.01	92.43	59.90	89.47	90.51	66.79	82.75
+EDA	83.82†	87.53†	90.79†	92.55	56.88†	88.67†	90.94	71.12	82.79
+SP	84.17	88.27	91.38	93.23	64.27	89.61	91.60	71.48	84.25
RoBERTa	87.77	89.03	92.88	94.84	61.48	90.58	93.15	77.98	85.96
+EDA	87.71†	88.58†	92.48†	95.41	58.88†	90.35†	92.93†	76.17†	85.31†
+SP	87.95	89.10	92.97	94.95	63.20	90.55†	92.93†	80.14	86.47

Table 2: Results on the GLUE Development sets with all training samples are used. The dagger “†” indicates performance that is inferior to standard fine-tuning of the PLM.

Model	SST-2	CoLA	MRPC	RTE	Avg.
BERT	92.43	59.90	90.51	66.79	77.41
+SP (-AfterLayer-12)	93.23	64.27	91.60	71.48	80.15
+SP-AfterLayer-11	92.89	63.38	91.19	71.84	79.83
+Pos Emb	93.00	62.91	91.09	70.40	79.35
+SP-AfterLayer-8	93.12	62.44	91.36	70.04	79.24
+Pos Emb	93.12	63.04	91.00	69.68	79.21
+SP-AfterLayer-5	93.12	61.34	90.88	70.40	78.94
+Pos Emb	92.89	61.62	91.03	71.48	79.26
+SP-Emb	93.23	61.11	90.82	68.23	78.35

Table 3: Comparison of the *Scanpath Module* at various model locations: after the n -th Transformer layer (*SP-AfterLayer- n*), and after the embedding layer (*SP-Emb*). We add extra positional embeddings to the token embeddings in the reordered sequence (*+Pos Emb*).

3.3 Ablation Studies

Location of the Scanpath Module We explore the impact of integrating the scanpath module at different feature-representation levels on the model’s performance. Specifically, we experiment with placing the scanpath module after the 11th, 8th, 5th, and embedding layer of the Transformer. In these cases, it is straightforward to use the subsequent Transformer layers to process the scanpath-guided reordered sequence; we therefore remove the scanpath encoder from the module. Moreover, we add extra positional embeddings to the token embeddings after the rearrangement, providing information about the positions of tokens in the sequence.

Table 3 shows that integrating the scanpath module into the model, regardless of its placement, yields improved performance compared to standard text-only fine-tuning. However, placing it at a lower position within the Transformer results in smaller gains. This may be attributed to the top Transformer layers capturing richer semantic information (Jawahar et al., 2019). Placing the scanpath module at the top facilitates better access to this information, potentially aiding in leveraging cognitive information. Furthermore, adding extra positional information to the reordered sequence marginally impacts performance.

Scanpath vs Random Order The core principle of the scanpath module is to utilize the order of fixations to integrate estimated cognitive information into the model. To study whether the observed gains truly arise from the order of fixations, we compare our method which rearranges the token-embedding sequence based on the scanpath to two baselines: (1) shuffling the scanpath ordering, and (2) randomly shuffling the token-embedding sequence. Table 4 shows that shuffling the scanpath results in consistent performance drops across all tasks, indicating the importance of the order of fixations. Furthermore, excluding the scanpath and randomly shuffling BERT token embeddings leads to a large decrease in performance gain, underscoring the importance of both fixated words and their order in enhancing model performance.

Model	SST-2	CoLA	MRPC	RTE	Avg.
BERT	92.43	59.90	90.51	66.79	77.41
+SP	93.23	64.27	91.60	71.48	80.15
+Shuffle SP	93.00	63.81	91.34	71.12	79.82
+Random Shuffle	92.78	60.66	91.42	68.95	78.45

Table 4: Comparison of strategies for reordering token embeddings: scanpath-guided (*SP*), shuffled scanpath-guided (*Shuffle SP*), and (*Random Shuffle*).

4 Conclusion

Our work contributes to the broad effort of enriching NLP models by grounding them in various domains of experience. Our specific focus lies in leveraging scanpath data, demonstrating its vital role in enhancing textual representation learning. By extending the standard PLM fine-tuning objective with a scanpath-integrated loss, we ground the LM in human language processing. Finally, our experiments show that the proposed method surpasses standard fine-tuning and EDA baselines on the GLUE benchmark, pointing to the potentially promising future direction of enriching textual representations with gaze data, especially for low-resource tasks and languages.

Limitations

One limitation of our work is that the scanpath-generation model—Eyettention—was pre-trained on a single eye-tracking corpus with a relatively small sample (see Appendix A). Participants read sentences covering only a single domain and a narrow range of difficulties. This limitation may restrict the knowledge acquired by Eyettention concerning human language processing, thus potentially leading to limited benefits when integrating simulated gaze data into LMs. In our experiments, we observe that our proposed fine-tuning scheme provides fewer benefits to RoBERTa than BERT, even in the low-resources setting. The key distinguishing factor between these models is the scale of unsupervised pre-training. We hypothesize that pre-training provides similar benefits targeted by the simulated cognitive signals generated from Eyettention. Including larger eye-tracking samples for pre-training, possibly covering diverse domains of text reading, could potentially enhance the model’s performance.

Furthermore, it is worth exploring the performance of using other state-of-the-art scanpath generators. Different architectures have been developed recently in the field (Bolliger et al., 2023; Khurana et al., 2023). Exploring the strengths and weaknesses of different scanpath generators when integrated into LMs could provide valuable insight into the development of improved scanpath generators for benefiting NLP tasks.

Ethics Statement

It is essential to acknowledge potential privacy risks in the collection, sharing, and processing of human gaze data. Due to the highly individual nature of eye movements, there exists a possibility of extracting sensitive information such as a participant’s identity (Jäger et al., 2020; Makowski et al., 2021), gender (Sammaknejad et al., 2017) and ethnicity (Blignaut and Wium, 2014) from gaze data, posing a risk of privacy leakage. The use of synthetic gaze data can help alleviate the necessity for large-scale experiments involving human subjects, although some amount of human gaze data remains necessary to train generative models.

References

Maria Barrett, Joachim Bingel, Nora Hollenstein, Marek Rei, and Anders Søgaard. 2018. [Sequence classi-](#)

[fication with human attention](#). In *Proceedings of the 22nd Conference on Computational Natural Language Learning*, pages 302–312, Brussels, Belgium. Association for Computational Linguistics.

Maria Barrett, Joachim Bingel, Frank Keller, and Anders Søgaard. 2016. [Weakly supervised part-of-speech tagging using eye-tracking data](#). In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 579–584, Berlin, Germany. Association for Computational Linguistics.

Yevgeni Berzak, Chie Nakamura, Amelia Smith, Emily Weng, Boris Katz, Suzanne Flynn, and Roger Levy. 2022. CELER: A 365-participant corpus of eye movements in L1 and L2 English reading. *Open Mind*, pages 1–10.

Pieter Blignaut and Daniël Wium. 2014. Eye-tracking data quality as affected by ethnicity and experimental design. *Behavior Research Methods*, 46:67–80.

Lena Bolliger, David Reich, Patrick Haller, Deborah Jakobi, Paul Prasse, and Lena Jäger. 2023. [ScanDL: A diffusion model for generating synthetic scanpaths on texts](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 15513–15538, Singapore. Association for Computational Linguistics.

Kyunghyun Cho, Bart van Merriënboer, Dzmitry Bahdanau, and Yoshua Bengio. 2014. [On the properties of neural machine translation: Encoder–decoder approaches](#). In *Proceedings of SSST-8, Eighth Workshop on Syntax, Semantics and Structure in Statistical Translation*, pages 103–111, Doha, Qatar. Association for Computational Linguistics.

Shuwen Deng, Paul Prasse, David Reich, Tobias Scheffer, and Lena Jäger. 2023a. [Pre-trained language models augmented with synthetic scanpaths for natural language understanding](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 6500–6507, Singapore. Association for Computational Linguistics.

Shuwen Deng, David R Reich, Paul Prasse, Patrick Haller, Tobias Scheffer, and Lena A Jäger. 2023b. Eyettention: An attention-based dual-sequence model for predicting human scanpaths during reading. *Proceedings of the ACM on Human-Computer Interaction*, 7(ETRA):1–24.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

496 Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu,
497 Teven Le Scao, Sylvain Gugger, Mariama Drame,
498 Quentin Lhoest, and Alexander Rush. 2020. [Trans-](#)
499 [formers: State-of-the-art natural language processing](#).
500 In *Proceedings of the 2020 Conference on Empirical*
501 *Methods in Natural Language Processing: System*
502 *Demonstrations*, pages 38–45, Online. Association
503 for Computational Linguistics.

504 Qizhe Xie, Zihang Dai, Eduard Hovy, Thang Luong,
505 and Quoc Le. 2020. Unsupervised data augmenta-
506 tion for consistency training. In *Advances in neural*
507 *information processing systems*, volume 33, pages
508 6256–6268.

509 Duo Yang and Nora Hollenstein. 2023. PLM-AS: Pre-
510 trained language models augmented with scanpaths
511 for sentiment classification. In *Proceedings of the*
512 *Northern Lights Deep Learning Workshop*, Tromsø,
513 Norway.

514 Tianyi Zhang, Felix Wu, Arzoo Katiyar, Kilian Q Wein-
515 berger, and Yoav Artzi. 2021. Revisiting few-sample
516 BERT fine-tuning. In *Proceedings of the 9th Inter-*
517 *national Conference on Learning Representations*
518 *(ICLR)*, Online.

A Model Details

Scanpath Generation Model For the utilization of the scanpath generation model Eyettention, we follow the work of (Deng et al., 2023a). The training process for the Eyettention model is conducted in two phases. First, we pre-train the Eyettention model on the L1 subset of the CELER corpus (Berzak et al., 2022), which comprises eye-tracking recordings collected from native speakers of English during natural reading sentences. Second, the Eyettention model is fine-tuned on downstream NLP tasks. More specifically, in our proposed scanpath-augmented fine-tuning scheme, we fine-tune the Transformer encoder and the Eyettention model, as well as train the scanpath encoder and the final dense layer from scratch. We tailor the parameters of Eyettention for specific downstream tasks, aiming to provide targeted inductive biases. For further details on the Eyettention model, please refer to (Deng et al., 2023b,a).

In our experiments, we evaluate our proposed approach using two distinct PLMs, BERT and RoBERTa, each equipped with its unique tokenizer. The Eyettention model includes a PLM in the text encoder for embedding the stimulus sentence. The generated fixation sequence (token index sequence) is based on the specific tokenizer associated with the PLM used. To facilitate a direct application of the arrangement operation based on the token-embedding sequence and fixation sequence without additional complex conversion, we maintain consistency by using the same PLMs in the Eyettention text encoder when evaluating specific PLMs as our base models. By replacing BERT with RoBERTa in the Eyettention text encoder, we observe a similar validation loss in scanpath prediction on the CELER corpus.

Scanpath Encoder The scanpath encoder is composed of a unidirection GRU layer (Cho et al., 2014) with a hidden size of 768 and a dropout rate of 0.1. We initialize the hidden state of the GRU layer using the [CLS] token outputs from the final layer of the PLMs.

B Training Details

We train all models using the PyTorch (Paszke et al., 2019) library on an NVIDIA A100-SXM4-40GB GPU using the NVIDIA CUDA platform. We use the pre-trained checkpoints from the HuggingFace repository (Wolf et al., 2020) for the language

model BERT_{base} and RoBERTa_{base}. The models are optimized using the AdamW optimizer (Loshchilov and Hutter, 2019). We set the maximum sequence length to 128 and the training batch size to 32.

In the high-resource setting, we train the models for 20 epochs and update the best checkpoint by measuring validation accuracy every 500 steps. For datasets with fewer than 500 steps per epoch, we update and validate at the end of each epoch. We tune the learning rates for BERT from {5e-5, 4e-5, 3e-5, 2e-5} and for RoBERTa from {3e-5, 2e-5, 1e-5} for each task, following the recommendations in the original paper (Devlin et al., 2019; Liu et al., 2019).

In the low-resource setting, we train the models for 10 epochs and save checkpoints every epoch. We use the same learning rate that was found optimal in the high-resource setting for each task. We perform 5 runs with different data seeds ({111,222,333,444,555}) for shuffling, while the seed s=42 is consistently utilized for model training across all models.

In both high-resource and low-resource settings, for our proposed scanpath-augmented fine-tuning method, we conduct a hyperparameter search on the development set to determine the optimal trade-off factor λ for each task, exploring values from {1, 0.7, 0.5, 0.3, 0.1, 0.01, 0.001}. For the EDA baseline, we tune the number of generated augmented sentences added to the original training set, exploring values from {1, 2, 4, 8, 16} based on the recommendations in the original paper (Wei and Zou, 2019).