Evaluating distillation methods for data-efficient syntax learning

Takateru Yamakoshi¹, Thomas L. Griffiths², R. Thomas McCoy*³, Robert D. Hawkins*¹

Stanford University, ²Princeton University, ³Yale University

{takateru,rdhawkins}@stanford.edu, tomg@princeton.edu, tom.mccoy@yale.edu

Abstract

Data-efficient training requires strong inductive biases. To the extent that transformer attention matrices encode syntactic relationships, we would predict that knowledge distillation (KD) targeting attention should selectively accelerate syntax acquisition relative to conventional logit-based KD. To test this hypothesis, we train GPT-2 student models on datasets ranging from 10K to 5M sentences using both distillation methods, evaluating them on both syntactic benchmarks and perplexity. Surprisingly, while logit-based KD dramatically improves data-efficiency, attention-based KD provides minimal benefit even for syntactic tasks. This suggests that output distributions provide sufficient supervisory signal for syntax acquisition, indicating that syntactic knowledge may be distributed throughout the network rather than localized in attention patterns.¹

syntactic structure attention-based distillation attention-based distillation

Figure 1: We compare knowledge distillation through logits and attention matrices as a form of inductive bias. We hypothesize that syntactic structures reflected in attention could provide an inductive bias that selectively enhances the student's performance in syntactic phenomena.

1 Introduction

Modern language models successfully capture many aspects of human linguistic abilities, from the fundamentals of grammar (Warstadt et al., 2020; Linzen and Baroni, 2021; Hu et al., 2024) to more sophisticated uses of world knowledge (Ivanova et al., 2024; Yamakoshi et al., 2023). However, they achieve these capabilities only after training on vastly more data than human children receive during language acquisition (Frank, 2023), motivating research into inductive biases (Warstadt et al., 2023), predispositions that guide learning toward particular solutions with less data. These biases include architectural modifications (Sartran et al., 2022), curriculum learning strategies (Diehl Martinez et al., 2023), training objectives (Ahuja et al., 2025), and specialized weight initialization techniques (Bencomo et al., 2025).

In this paper, we use knowledge distillation to study which aspects of a model's learned representations are most critical for scaffolding particular linguistic capabilities. We focus specifically on learning syntax – an ability classically theorized to require strong (innate) biases (Chomsky, 1965; McCoy et al., 2020). Previous research has shown that syntactic information is encoded in the attention mechanism of transformer models (Clark et al., 2019; Ravishankar et al., 2021), and that constraining these attention matrices can serve as an effective inductive bias for syntax (Nguyen et al., 2020; Qian et al., 2021; Yoshida and Oseki, 2022; Sartran et al., 2022).

If syntax is primarily localized in attention matrices, then KD targeting these representations should be sufficient to transfer syntactic abilities, potentially even outperforming conventional logit-based distillation. Moreover, if attention serves as the

¹Code is available at https://github.com/taka-yamakoshi/attention_structures.

primary locus of syntactic information, we would predict that attention-based KD should show selective advantages for syntactic tasks relative to general language modeling performance, providing a targeted inductive bias for syntax acquisition under data-limited conditions. In other words, we use distillation as an analytical tool to probe whether syntactic knowledge is indeed localized in specific model components or distributed more broadly throughout the network (Figure 1).

To investigate these two possibilities, we conducted controlled experiments using a pretrained GPT-2 model (Radford et al., 2019) as teacher, training architecturally-identical students on datasets ranging from 10K to 5M sentences. Our contributions are twofold. First, we demonstrate that conventional logit-based distillation drastically reduces the amount of data required for learning syntax, reaching teacher-level performance with only 500K sentences of training data. Second, surprisingly, attention-based KD provides limited benefit for syntactic tasks despite prior evidence that these matrices encode syntactic structure. Our work illustrates how knowledge distillation can serve as a powerful analytical tool for understanding which aspects of a model's representations are effective for achieving data-efficiency with respect to specific linguistic capabilities.

2 Related Work

2.1 Knowledge distillation

Knowledge distillation consists of three main approaches (Gou et al., 2021): response-based KD, which aligns the output distributions of teacher and student models; feature-based KD, which matches internal representations to transfer detailed computational patterns; and relation-based KD, which preserves relational structures across multiple samples. In this work, we employ both response-based KD through logits and feature-based KD through attention to investigate their relative effectiveness for transferring syntactic knowledge.

While KD was initially developed for model compression, its applications have been expanded in several directions. For example, Furlanello et al. (2018) demonstrated that distilling knowledge to a student of identical architecture can actually improve performance. Others have used KD to facilitate transfer between architecturally different models (Kuncoro et al., 2019, 2020; Abnar et al., 2020), showing that inductive biases from special-

ized architectures can be distilled into more general ones. Finally, recent work has explored KD for data-efficient training, using ensembles of teacher models to improve student performance on limited data (Timiryasov and Tastet, 2023; Samuel, 2023; Yam and Paek, 2024). Our approach maintains architectural consistency between teacher and student, and uses a single pre-trained model as the teacher, in order to isolate the effects of different distillation mechanisms on syntactic competencies.

2.2 How transformers represent syntax

Understanding how transformers capture syntactic structures has been a central question in interpretability research. Numerous studies have identified attention matrices as repositories of syntactic information, with certain attention heads specializing in tracking specific syntactic relations (Clark et al. 2019; Vig and Belinkov 2019; Htut et al. 2019; Ravishankar et al. 2021; Lee et al. 2024; cf. Hassid et al. 2022). Others have shown that incorporating explicit syntactic guidance into attention patterns can improve performance on syntactic tasks (Strubell et al., 2018; Sachan et al., 2021; Bugliarello and Okazaki, 2020; Wang et al., 2019b; Bai et al., 2021; Chen et al., 2024). Recent work has also investigated the data requirements for acquiring syntactic knowledge, with some studies finding that pre-training on small, developmentally plausible corpora can lead to syntax acquisition with the right inductive biases (Warstadt et al., 2023; Huebner et al., 2021). However, the precise mechanisms through which transformers acquire syntactic knowledge, and the relative contributions of different elements of the architecture, remain open questions.

3 Approach

We ask whether distillation through attention provides a stronger inductive bias for syntax acquisition compared to conventional distillation through logits. To investigate this question, we conducted controlled experiments using an identical GPT-2 small architecture (Radford et al., 2019) for both the teacher and student models. We focused on GPT-2 small as it is the smallest available pretrained language model with adequate performance, enabling feasible experiments. The teacher model was a fully pre-trained GPT-2, while the student models were trained from scratch on different subsets of the BabyLM dataset (Warstadt et al., 2023),

ranging from 10K to 5M sentences. By varying the dataset size, we assessed how different distillation methods affect data efficiency. All results reported are averages across three random seeds. Complete training details are provided in Appendix A.

3.1 Distillation via logits

We first established the baseline effectiveness of conventional KD through output distributions. Following Kim and Rush (2016), we implemented word-level KD where the student model learns to match the teacher's output probability distributions. Let $P_t(w|w_{< i})$ and $P_s(w|w_{< i})$ be the conditional probability of the word w at the i-th token calculated by the teacher and the student model respectively. The auxiliary loss for distillation $\mathcal{L}_{\text{logits}}$ for each sentence with length N was defined as

$$\mathcal{L}_{\text{logits}} = \frac{1}{N} \sum_{i=1}^{N} \sum_{w \in V} P_t(w|w_{< i}) \log P_s(w|w_{< i}),$$

where V is the vocabulary. This formulation is equivalent to calculating the forward KL divergence between teacher and student distributions at each token position and taking the average. This auxiliary loss was then added to the standard crossentropy loss \mathcal{L}_{CE} with a coefficient α controlling the strength of distillation:

$$\mathcal{L} = \mathcal{L}_{CE} + \alpha \mathcal{L}_{logits}$$
.

Based on preliminary experiments testing different values of α (Figure S1), we found that $\alpha=10$ led to optimal performance and fixed it at this value for all logit-based distillation experiments.

3.2 Distillation via attention

To test our hypothesis that attention matrices might provide a stronger inductive bias for syntax acquisition, we implemented feature-based KD targeting the attention mechanisms directly. Previous studies have shown that KD through attention matrices is an effective method to perform model compression both in computer vision (Zagoruyko and Komodakis, 2017; Wang et al., 2022; Li et al., 2024) and natural language processing (Aguilar et al., 2020). Following these studies, we calculated the auxiliary loss $\mathcal{L}_{\rm attn}$ as the mean squared error between the attention matrices of the teacher and the student. Let $A_t(l,h)$ and $A_s(l,h)$ be the attention matrices of the head h at layer l calculated by the

teacher and the student model, respectively.

$$\mathcal{L}_{\text{attn}} = \frac{1}{L} \frac{1}{H} \sum_{l=1}^{L} \sum_{h=1}^{H} \text{MSE}(A_t(l,h) - A_s(l,h)),$$

where L and H are the number of layers and heads. As with logit-based distillation, this auxiliary loss was added to the cross-entropy loss with a coefficient $\alpha=1$, based on preliminary experiments (Figure S1).

3.3 Evaluation

To test our hypothesis about the relative effectiveness of different distillation approaches for syntax acquisition, we evaluated models on both syntactic benchmarks and a conventional language modeling metric. If attention matrices encode critical syntactic information not fully captured in output distributions, then attention-based distillation should show selective advantages on syntactic tasks, especially when training data is limited. For syntactic evaluation, we used three datasets based on minimal pairs:

- Linzen (Linzen et al., 2016; Gulordava et al., 2018) tests subject-verb agreement across various syntactic constructions.
- **BLiMP** (Warstadt et al., 2020) tests 67 distinct tasks across 12 syntactic phenomena.
- **Zorro** (Huebner et al., 2021) tests simple syntactic tasks that align with the developmental nature of our training data.

For each item in these benchmarks, we computed the log probability of both sentences and counted the model as correct if it assigns a higher probability to the grammatically acceptable variant. To ensure we capture overall language modeling capability (beyond syntax), we also measured perplexity on the BabyLM test split. This dual evaluation allows us to distinguish between general improvements in language modeling and selective enhancements in syntactic competence, helping to determine whether different distillation methods provide domain-specific inductive biases or general learning benefits.

4 Results

Before testing the effects of KD on syntactic performance, we check how well each KD approach achieves its objective. As shown in Figure S2, logitbased KD successfully enables the student model

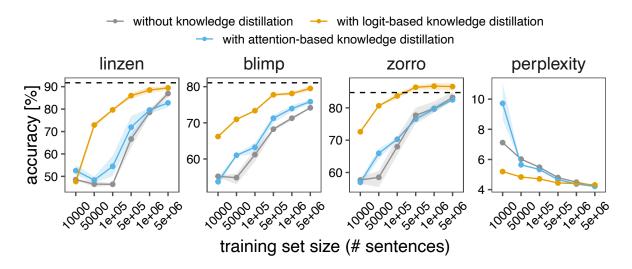


Figure 2: Performance of the students trained on datasets with different sizes. Linzen, BLiMP, and Zorro are targeted syntactic evaluations, while perplexity quantifies general language modeling performance. Bands show the bootstrapped 95% CI across three random seeds. Dashed lines indicate the performance of the teacher.

to achieve a much lower KL divergence from the teacher model, and attention-based KD enables the student model to achieve a much more similar attention pattern to the teacher model. Next, we turn to our main question: how does each KD method affect the linguistic abilities of the student models?

4.1 Logit-based KD improves data efficiency

Figure 2 shows the performance of students trained with and without logit-based KD, across varying dataset sizes. Baseline student models require 5M+sentences to approach teacher performance on syntactic benchmarks without any distillation. Logit-based KD resulted in substantial improvements on both syntactic benchmarks and perplexity. With just 500K sentences (approx. 5M tokens), the students approached the performance of the teacher. The impact of logit-based KD was particularly pronounced with smaller datasets, where inductive biases are most crucial. For models trained on just 50K-100K sentences, KD provided a >20% boost in performance on the Linzen benchmark, elevating models from chance-level performance (50%).

Interestingly, some students outperformed the teacher on the Zorro benchmark. This may reflect the domain alignment between the student's training data and the benchmark, which uses the vocabulary from the BabyLM dataset, whereas the teacher's training data was a general Internet-based corpus. This result suggests that distillation can combine the teacher's knowledge and the domain-specific property of the student's training data.

4.2 Attention-based KD has limited effect

Contrary to our hypothesis that attention matrices provide a stronger inductive bias for syntax acquisition, Figure 2 shows that attention-based KD offered limited benefits compared to logit-based KD, even though it leads to better alignment in attention (Figure S2). This pattern held consistently across all dataset sizes tested, suggesting that the syntactic information encoded in attention matrices may not provide substantial advantages beyond what is already captured in output distributions. Moreover, attention-based KD had higher perplexity than the baseline without KD for the smallest dataset, suggesting attention loss may create conflicting gradients with the language modeling objective.

To determine whether attention-based KD benefits particular aspects of syntax, we performed fine-grained evaluations across grammatical phenomena. Figure S3 breaks down performance by tasks, and Figure S4 by phenomena, for the BLiMP benchmark. Despite considerable variation in the teacher's performance across these tasks and phenomena, the relative performance pattern of different distillation approaches remained remarkably consistent. Interestingly, however, attention-based KD was competitive with logit-based KD in some phenomena, particularly in ellipsis. Similar patterns, including the competitive performance of attention-based KD in ellipsis, were also observed for the Zorro benchmark (Figure S5).

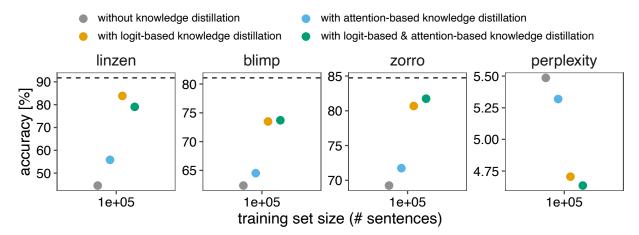


Figure 3: Combining both distillation types does not outperform logit-based distillation alone. Each point represents a single run. All runs were done on the 100K-sentence dataset. Dashed lines indicate the teacher's performance.

4.3 Logits may encode sufficient syntactic information

One possible explanation for these findings is that logit-based distillation might indirectly align attention patterns, making explicit attention distillation redundant (Hewitt and Manning 2019; Murty et al. 2022; Wu et al. 2024; Simon et al. 2025). A preliminary analysis supports this hypothesis: when both KD methods are combined in the same objective, performance remains similar to purely logit-based KD (Figure 3), suggesting no unique contribution from attention. If output distributions provide sufficient signal to scaffold data-efficient syntax learning, it is possible that syntax might be distributed throughout the network rather than being localized primarily in attention patterns.

5 Discussion

Our results reveal a striking contrast in the ability to improve data-efficiency among different KD methods. While KD via logits enabled student models to achieve teacher-level syntactic performance with just 500K sentences, KD via attention matrices – despite their capacity to encode syntactic structures – offered only marginal benefits.

Our results contrast with previous studies showing the benefits of attention distillation (Aguilar et al., 2020; Wang et al., 2020). We identify two potential sources of this divergence: the evaluation metric, and the training objective. First, we specifically focused on syntactic task performance, whereas prior work used more general benchmarks such as SQuAD and GLUE. This raises the possibility that attention patterns may actually be more important for semantics than purely syntac-

tic processing. Second, previous studies focused on model compression (reduction in parameter counts), which is different from our objective (reduction in training set size, using identical model architectures for the teacher and student). For the compression objective, explicit attention alignment may help compensate for reduced capacity, but may be unnecessary (or hurt performance) when the size of the student model is equal to the teacher model.

One key advantage of KD is that it requires minimal assumptions about the specific form of inductive biases. In fact, our results demonstrate that strong syntactic performance can be achieved without relying on explicit grammatical rules. On the other hand, KD-based approaches present certain challenges. KD can be computationally intensive, requiring forward passes through the teacher model for the entire training dataset (see Appendix B for the rough estimates), and the inductive biases transferred via KD are less interpretable than those from explicit grammar-based approaches (Sartran et al., 2022).

Taken together, our findings highlight how feature-based KD can serve as a powerful analytical tool to investigate which features are necessary or sufficient for specific capabilities. Effective distillation through a particular feature suggests that it contains enough information to serve as an inductive bias for the target capability. Our results suggest that the information contained in attention matrices was not a strong enough inductive bias for syntax acquisition, but future work must systematically compare different feature-based KD methods to better understand how different linguistic competencies are encoded within transformer representations.

Limitations

Our evaluation focused specifically on syntactic benchmarks, motivated by previous work showing that attention matrices encode syntactic information and that syntactically-guided attention constraints serve as effective inductive biases. While this targeted approach allowed us to directly address questions about syntax acquisition, it limits the generalizability of our findings to other linguistic competencies. Different aspects of linguistic knowledge may be encoded preferentially in different components of transformer architectures, and distillation methods might show varying effectiveness across other linguistic domains, from semantics and pragmatics to discourse representation. Future work should evaluate attention-based KD on a broader range of benchmarks spanning diverse capabilities, such as SuperGLUE (Wang et al., 2019a) for language understanding and EWOK (Ivanova et al., 2024) for world knowledge. A more comprehensive evaluation would allow researchers to determine whether the relative efficacy of different distillation methods varies across linguistic domains. It's possible that attention-based KD might provide stronger benefits for capabilities other than syntax, such as long-range semantic dependencies or pragmatic reasoning.

Moreover, there may be cross-linguistic differences in the way linguistic knowledge is encoded in the transformer architecture. For example, languages with richer morphology may benefit more from the structural information in the attention, while languages with a more flexible word order may benefit less. Future work should include non-English languages to test these hypotheses.

We acknowledge that there can be multiple ways to implement attention-based KD, and that our analysis does not rule out the possibility of other, more effective methods. We experimented with an alternative method by performing kernel density estimation using a precomputed batch of attention matrices, which did not outperform the reported method. Other methods include specifically targeting a subset of attention heads or dynamically mapping attention matrices from multiple heads instead of a one-to-one match (as in Zhao et al. (2024)). Detailed exploration of these other methods will be left to future work.

Additionally, our experiments used a single pretrained model (GPT-2) as the teacher. Exploring different teacher architectures and model sizes would help determine the generalizability of our findings across different model families and capabilities. Finally, our exploration of feature-based distillation was limited to attention matrices; future work could investigate other internal representations such as hidden states, feed-forward network activations, or combinations of these features.

Ethics Statement

All datasets (BabyLM, Linzen, BLiMP, and Zorro) and the model (GPT-2) used in this paper were employed according to their intended usage. BabyLM consists of the following publicly available datasets (Warstadt et al., 2023):

- CHILDES² (MacWhinney, 2000)
- British National Corpus³ (Consortium, 2007)
- Children's Book Test (Hill et al., 2016)
- Children's Stories Text Corpus (Bensaid et al., 2021)
- Project Gutenberg (Gerlach and Font-Clos, 2020)
- OpenSubtitles (Lison and Tiedemann, 2016)
- QED (Abdelali et al., 2014)
- Wikipedia
- Simple English Wikipedia
- Switchboard Corpus (Godfrey et al., 1992)

While we used knowledge distillation to distill the inductive biases required for data-efficient syntax learning, KD can also transfer the biases embedded in the teacher. When training student models using KD, we need to consider the biases of the teacher as well as those in the training dataset.

Acknowledgments

We thank Yohei Oseki, Ryo Yoshida, and the Language Computational Cognitive Science Laboratory at UTokyo for thoughtful discussions and comments. This work was supported by a Seed Grant from the Princeton-UTokyo Strategic Partnership.

²CC BY-NC-SA 3.0 License

³BNC License

References

- Ahmed Abdelali, Francisco Guzman, Hassan Sajjad, and Stephan Vogel. 2014. The AMARA Corpus: Building Parallel Language Resources for the Educational Domain. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*, pages 1856–1862.
- Samira Abnar, Mostafa Dehghani, and Willem Zuidema. 2020. Transferring Inductive Biases through Knowledge Distillation. *arXiv preprint arXiv:2006.00555*.
- Gustavo Aguilar, Yuan Ling, Yu Zhang, Benjamin Yao, Xing Fan, and Chenlei Guo. 2020. Knowledge Distillation from Internal Representations. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 7350–7357.
- Kabir Ahuja, Vidhisha Balachandran, Madhur Panwar, Tianxing He, Noah A. Smith, Navin Goyal, and Yulia Tsvetkov. 2025. Learning Syntax Without Planting Trees: Understanding Hierarchical Generalization in Transformers. *Transactions of the Association for Computational Linguistics*, 13:121–141.
- Jason Ansel, Edward Yang, Horace He, Natalia Gimelshein, Animesh Jain, Michael Voznesensky, Bin Bao, Peter Bell, David Berard, Evgeni Burovski, and 1 others. 2024. PyTorch 2: Faster machine learning through dynamic python bytecode transformation and graph compilation. In *Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems*, Volume 2, pages 929–947.
- Jiangang Bai, Yujing Wang, Yiren Chen, Yaming Yang, Jing Bai, Jing Yu, and Yunhai Tong. 2021. Syntax-BERT: Improving pre-trained transformers with syntax trees. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 3011–3020, Online. Association for Computational Linguistics.
- Gianluca Bencomo, Max Gupta, Ioana Marinescu, R. Thomas McCoy, and Thomas L. Griffiths. 2025. Teasing Apart Architecture and Initial Weights as Sources of Inductive Bias in Neural Networks. In *Proceedings of the 47th Annual Meeting of the Cognitive Science Society*, pages 4558–4565.
- Eden Bensaid, Mauro Martino, Benjamin Hoover, and Hendrik Strobelt. 2021. FairyTailor: A Multimodal Generative Framework for Storytelling. *arXiv* preprint arXiv:2108.04324.
- Emanuele Bugliarello and Naoaki Okazaki. 2020. Enhancing Machine Translation with Dependency-Aware Self-Attention. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1618–1627.
- Angelica Chen, Ravid Shwartz-Ziv, Kyunghyun Cho, Matthew L Leavitt, and Naomi Saphra. 2024. Sudden Drops in the Loss: Syntax Acquisition, Phase Transitions, and Simplicity Bias in MLMs. In *International Conference on Learning Representations*.

- Noam Chomsky. 1965. *Aspects of the Theory of Syntax*. MIT Press.
- Kevin Clark, Urvashi Khandelwal, Omer Levy, and Christopher D. Manning. 2019. What Does BERT Look at? An Analysis of BERT's Attention. In *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 276–286, Florence, Italy. Association for Computational Linguistics.
- BNC Consortium. 2007. The British National Corpus, XML Edition.
- Richard Diehl Martinez, Zébulon Goriely, Hope Mc-Govern, Christopher Davis, Andrew Caines, Paula Buttery, and Lisa Beinborn. 2023. CLIMB Curriculum Learning for Infant-inspired Model Building. In *Proceedings of the BabyLM Challenge at the 27th Conference on Computational Natural Language Learning*, pages 112–127, Singapore. Association for Computational Linguistics.
- Michael C. Frank. 2023. Bridging the data gap between children and large language models. *Trends in Cognitive Sciences*, 27(11):990–992.
- Tommaso Furlanello, Zachary Lipton, Michael Tschannen, Laurent Itti, and Anima Anandkumar. 2018. Born Again Neural Networks. In *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 1607–1616. PMLR.
- Martin Gerlach and Francesc Font-Clos. 2020. A Standardized Project Gutenberg Corpus for Statistical Analysis of Natural Language and Quantitative Linguistics. *Entropy*, 22(1):126.
- John J Godfrey, Edward C Holliman, and Jane Mc-Daniel. 1992. SWITCHBOARD: Telephone speech corpus for research and development. In Acoustics, speech, and signal processing, ieee international conference on, volume 1, pages 517–520. IEEE Computer Society.
- Jianping Gou, Baosheng Yu, Stephen J Maybank, and Dacheng Tao. 2021. Knowledge Distillation: A Survey. *International Journal of Computer Vision*, 129(6):1789–1819.
- Kristina Gulordava, Piotr Bojanowski, Edouard Grave, Tal Linzen, and Marco Baroni. 2018. Colorless Green Recurrent Networks Dream Hierarchically. In Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers), pages 1195–1205, New Orleans, Louisiana. Association for Computational Linguistics.
- Michael Hassid, Hao Peng, Daniel Rotem, Jungo Kasai, Ivan Montero, Noah A. Smith, and Roy Schwartz. 2022. How Much Does Attention Actually Attend? Questioning the Importance of Attention in

- Pretrained Transformers. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 1403–1416, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- John Hewitt and Christopher D. Manning. 2019. A Structural Probe for Finding Syntax in Word Representations. 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 4129–4138, Minneapolis, Minnesota. Association for Computational Linguistics.
- Felix Hill, Antoine Bordes, Sumit Chopra, and Jason Weston. 2016. The Goldilocks principle: Reading children's books with explicit memory representations. In *International Conference on Learning Representations*.
- Phu Mon Htut, Jason Phang, Shikha Bordia, and Samuel R Bowman. 2019. Do attention heads in BERT track syntactic dependencies? *arXiv* preprint *arXiv*:1911.12246.
- Jennifer Hu, Kyle Mahowald, Gary Lupyan, Anna Ivanova, and Roger Levy. 2024. Language models align with human judgments on key grammatical constructions. *Proceedings of the National Academy of Sciences*, 121(36):e2400917121.
- Philip A. Huebner, Elior Sulem, Fisher Cynthia, and Dan Roth. 2021. BabyBERTa: Learning More Grammar With Small-Scale Child-Directed Language. In *Proceedings of the 25th Conference on Computational Natural Language Learning*, pages 624–646, Online. Association for Computational Linguistics.
- Anna A Ivanova, Aalok Sathe, Benjamin Lipkin, Unnathi Kumar, Setayesh Radkani, Thomas H Clark, Carina Kauf, Jennifer Hu, RT Pramod, Gabriel Grand, and 1 others. 2024. Elements of World Knowledge (EWOK): A cognition-inspired framework for evaluating basic world knowledge in language models. arXiv preprint arXiv:2405.09605.
- Yoon Kim and Alexander M. Rush. 2016. Sequence-Level Knowledge Distillation. In *Proceedings of* the 2016 Conference on Empirical Methods in Natural Language Processing, pages 1317–1327, Austin, Texas. Association for Computational Linguistics.
- Adhiguna Kuncoro, Chris Dyer, Laura Rimell, Stephen Clark, and Phil Blunsom. 2019. Scalable Syntax-Aware Language Models Using Knowledge Distillation. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3472–3484, Florence, Italy. Association for Computational Linguistics.
- Adhiguna Kuncoro, Lingpeng Kong, Daniel Fried, Dani Yogatama, Laura Rimell, Chris Dyer, and Phil Blunsom. 2020. Syntactic Structure Distillation Pretraining for Bidirectional Encoders. *Transactions of the Association for Computational Linguistics*, 8:776–794.

- Isabelle Lee, Joshua Lum, Ziyi Liu, and Dani Yogatama. 2024. Causal Interventions on Causal Paths: Mapping GPT-2's Reasoning From Syntax to Semantics. In *Causality and Large Models @NeurIPS 2024*.
- Alexander C. Li, Yuandong Tian, Beidi Chen, Deepak Pathak, and Xinlei Chen. 2024. On the Surprising Effectiveness of Attention Transfer for Vision Transformers. In *Advances in Neural Information Processing Systems*, volume 37, pages 113963–113990.
- Tal Linzen and Marco Baroni. 2021. Syntactic Structure from Deep Learning. *Annual Review of Linguistics*, 7(Volume 7, 2021):195–212.
- Tal Linzen, Emmanuel Dupoux, and Yoav Goldberg. 2016. Assessing the Ability of LSTMs to Learn Syntax-Sensitive Dependencies. Transactions of the Association for Computational Linguistics, 4:521– 535.
- Pierre Lison and Jörg Tiedemann. 2016. OpenSubtitles2016: Extracting Large Parallel Corpora from Movie and TV Subtitles. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16)*, pages 923–929, Portorož, Slovenia. European Language Resources Association (ELRA).
- Ilya Loshchilov and Frank Hutter. 2019. Decoupled Weight Decay Regularization. In *International Conference on Learning Representations*.
- Brian MacWhinney. 2000. *The CHILDES project: The database*, volume 2. Psychology Press.
- R. Thomas McCoy, Robert Frank, and Tal Linzen. 2020. Does Syntax Need to Grow on Trees? Sources of Hierarchical Inductive Bias in Sequence-to-Sequence Networks. *Transactions of the Association for Computational Linguistics*, 8:125–140.
- Shikhar Murty, Pratyusha Sharma, Jacob Andreas, and Christopher D Manning. 2022. Characterizing intrinsic compositionality in transformers with tree projections. *arXiv preprint arXiv:2211.01288*.
- Xuan-Phi Nguyen, Shafiq Joty, Steven Hoi, and Richard Socher. 2020. Tree-Structured Attention with Hierarchical Accumulation. In *International Conference on Learning Representations*.
- Peng Qian, Tahira Naseem, Roger Levy, and Ramón Fernandez Astudillo. 2021. Structural Guidance for Transformer Language Models. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 3735–3745, Online. Association for Computational Linguistics.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, and 1 others. 2019. Language Models are Unsupervised Multitask Learners. *OpenAI blog*, 1(8):9.

- Vinit Ravishankar, Artur Kulmizev, Mostafa Abdou, Anders Søgaard, and Joakim Nivre. 2021. Attention Can Reflect Syntactic Structure (If You Let It). In Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume, pages 3031–3045, Online. Association for Computational Linguistics.
- Devendra Sachan, Yuhao Zhang, Peng Qi, and William L. Hamilton. 2021. Do Syntax Trees Help Pre-trained Transformers Extract Information? In Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume, pages 2647–2661, Online. Association for Computational Linguistics.
- David Samuel. 2023. Mean BERTs make erratic language teachers: the effectiveness of latent bootstrapping in low-resource settings. In *Proceedings of the BabyLM Challenge at the 27th Conference on Computational Natural Language Learning*, pages 221–237, Singapore. Association for Computational Linguistics.
- Laurent Sartran, Samuel Barrett, Adhiguna Kuncoro, Miloš Stanojević, Phil Blunsom, and Chris Dyer. 2022. Transformer Grammars: Augmenting Transformer Language Models with Syntactic Inductive Biases at Scale. *Transactions of the Association for Computational Linguistics*, 10:1423–1439.
- Pablo J. Diego Simon, Emmanuel Chemla, Jean-Remi King, and Yair Lakretz. 2025. Probing Syntax in Large Language Models: Successes and Remaining Challenges. In Second Conference on Language Modeling.
- Emma Strubell, Patrick Verga, Daniel Andor, David Weiss, and Andrew McCallum. 2018. Linguistically-Informed Self-Attention for Semantic Role Labeling. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 5027–5038, Brussels, Belgium. Association for Computational Linguistics.
- Inar Timiryasov and Jean-Loup Tastet. 2023. Baby Llama: knowledge distillation from an ensemble of teachers trained on a small dataset with no performance penalty. In *Proceedings of the BabyLM Challenge at the 27th Conference on Computational Natural Language Learning*, pages 279–289, Singapore. Association for Computational Linguistics.
- Jesse Vig and Yonatan Belinkov. 2019. Analyzing the Structure of Attention in a Transformer Language Model. In *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 63–76, Florence, Italy. Association for Computational Linguistics.
- Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2019a. SuperGLUE: A Stickier Benchmark for General-Purpose Language Understanding Systems. In *Advances in Neural Information Processing Systems*, volume 32.

- Kai Wang, Fei Yang, and Joost van de Weijer. 2022. Attention Distillation: self-supervised vision transformer students need more guidance. *arXiv* preprint *arXiv*:2210.00944.
- Wenhui Wang, Furu Wei, Li Dong, Hangbo Bao, Nan Yang, and Ming Zhou. 2020. MiniLM: Deep Self-Attention Distillation for Task-Agnostic Compression of Pre-Trained Transformers. In *Advances in Neural Information Processing Systems*, volume 33, pages 5776–5788.
- Yaushian Wang, Hung-Yi Lee, and Yun-Nung Chen. 2019b. Tree Transformer: Integrating Tree Structures into Self-Attention. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 1061–1070, Hong Kong, China. Association for Computational Linguistics.
- Alex Warstadt, Aaron Mueller, Leshem Choshen, Ethan Wilcox, Chengxu Zhuang, Juan Ciro, Rafael Mosquera, Bhargavi Paranjabe, Adina Williams, Tal Linzen, and Ryan Cotterell. 2023. Findings of the BabyLM Challenge: Sample-Efficient Pretraining on Developmentally Plausible Corpora. In *Proceedings of the BabyLM Challenge at the 27th Conference on Computational Natural Language Learning*, pages 1–34, Singapore. Association for Computational Linguistics.
- Alex Warstadt, Alicia Parrish, Haokun Liu, Anhad Mohananey, Wei Peng, Sheng-Fu Wang, and Samuel R. Bowman. 2020. BLiMP: The Benchmark of Linguistic Minimal Pairs for English. *Transactions of the Association for Computational Linguistics*, 8:377–392.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, and 3 others. 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.
- Cindy Wu, Ekdeep Singh Lubana, Bruno Kacper Mlodozeniec, Robert Kirk, and David Krueger. 2024. What Mechanisms Does Knowledge Distillation Distill? In *Proceedings of UniReps: the First Workshop on Unifying Representations in Neural Models*, volume 243 of *Proceedings of Machine Learning Research*, pages 60–75. PMLR.
- Hong Meng Yam and Nathan Paek. 2024. Teaching Tiny Minds: Exploring Methods to Enhance Knowledge Distillation for Small Language Models. In *The 2nd BabyLM Challenge at the 28th Conference on Computational Natural Language Learning*, pages 302–307, Miami, FL, USA. Association for Computational Linguistics.

Takateru Yamakoshi, James McClelland, Adele Goldberg, and Robert Hawkins. 2023. Causal interventions expose implicit situation models for commonsense language understanding. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 13265–13293, Toronto, Canada. Association for Computational Linguistics.

Ryo Yoshida and Yohei Oseki. 2022. Composition, Attention, or Both? In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 5822–5834, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Sergey Zagoruyko and Nikos Komodakis. 2017. Paying More Attention to Attention: Improving the Performance of Convolutional Neural Networks via Attention Transfer. In *International Conference on Learning Representations*.

Tianyang Zhao, Kunwar Yashraj Singh, Srikar Appalaraju, Peng Tang, Vijay Mahadevan, R. Manmatha, and Ying Nian Wu. 2024. No Head Left Behind – Multi-Head Alignment Distillation for Transformers. 38(7):7514–7524.

A Training details

Table S1 shows hyperparameters used in our experiments. The BabyLM preprocessing pipeline⁴ was used to clean the dataset. Since the dataset has one sentence per line, we used the number of sentences as the measure of dataset size rather than the number of words or tokens. All train runs had the same number of training steps (156,250 steps) except for those for the largest dataset size (5,000,000 sentences), which had 312,500 steps to make sure the model sees the entire dataset more than once. We used AdamW optimizer (Loshchilov and Hutter, 2019) with a linear warm-up for 1% of the total number of training steps.

We used Hugging Face transformers (version 4.45.2; Apache License 2.0) (Wolf et al., 2020) and PyTorch (version 2.4.1; BSD-style license ⁵) (Ansel et al., 2024) to train and evaluate models. We used the original OpenAI checkpoint provided by Hugging Face transformers⁶ as the teacher. Experiments took approximately 750 GPU hours with NVIDIA RTX A6000 GPUs.

B Computational costs

Since we used the same model size for both the teacher and the student, our KD experiment required twice the amount of compute compared

n_layers	12
n_heads	12
hidden_size	768
intermediate_size	3072
max # tokens	128
batch size	32
learning rate	0.0002

Table S1: Hyperparameters

with a simple language model training. Regarding the difference between logit-based and attentionbased KDs, there was minimal empirical difference in training time (both took 15 hours using RTX A6000). Theoretically, logit-based KD adds 32 MFLOPs to the standard forward pass of GPT-2, whereas attention-based KD adds 7.1 MFLOPs. The difference stems from (1) different numbers of entries in each representation, and (2) different numbers of operations involved in calculating the loss, both of which are higher for logit-based KD. First, logits have (sequence length)*(vocabulary size)=128*50257=6.4M entries and attention matrices have (sequence length)*(sequence length)*(# layers)*(# heads)=128*128*12*12=2.4M entries. Second, logit-based KD takes five operations: softmax for the teacher (exponential+sum+division), multiplication, and sum, whereas attention-based KD takes three operations: subtraction, square, and sum, resulting in the above estimates.

⁴https://github.com/babylm/babylm_data_ preprocessing

⁵https://github.com/pytorch/pytorch/blob/main/ LICENSE

⁶https://huggingface.co/openai-community/gpt2

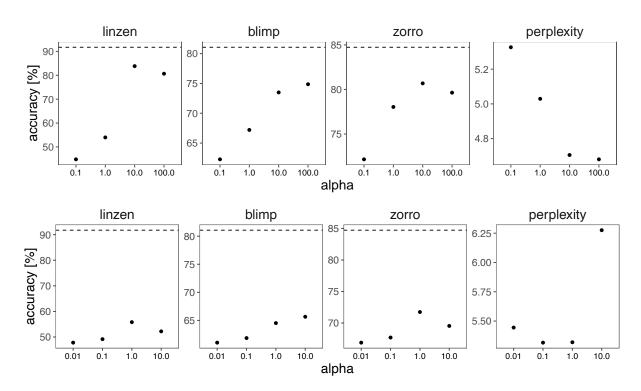


Figure S1: Choosing the best regularization coefficient for logit-based KD (top row) and attention-based KD (bottom row). We used the training set size of 100,000 sentences. Each point shows the result of a single run. Dashed lines show the teacher's performance.

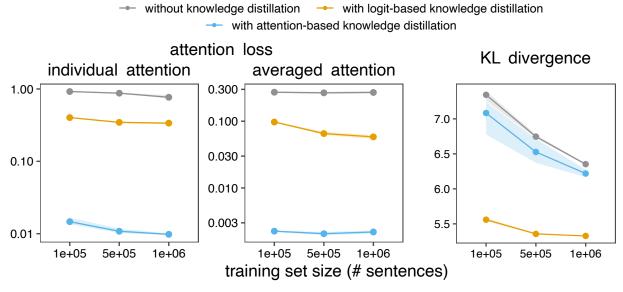


Figure S2: Auxiliary losses evaluated on the BLiMP dataset. We randomly selected 3 items from each task (3*67=201 in total). Unlike attention-based knowledge distillation, logit-based knowledge distillation does not align the internal computations, which leaves the possibility that similar attention patterns are implemented in both the teacher and the student by different attention heads. To account for this, we calculated the loss using the attention matrices averaged across layers and heads (middle), in addition to the loss used in training (left) as described in Section 3.2. Y-axis of the left two panels are on the log scale.

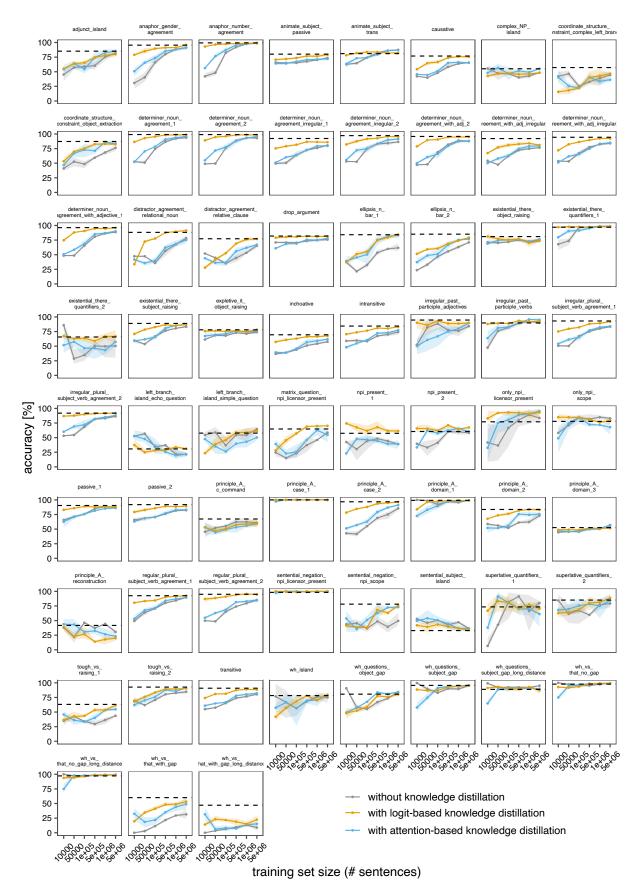


Figure S3: Performance on BLiMP split into tasks. Ribbons show the bootstrapped 95% CI across three random seeds. Dashed lines show the teacher's performance.

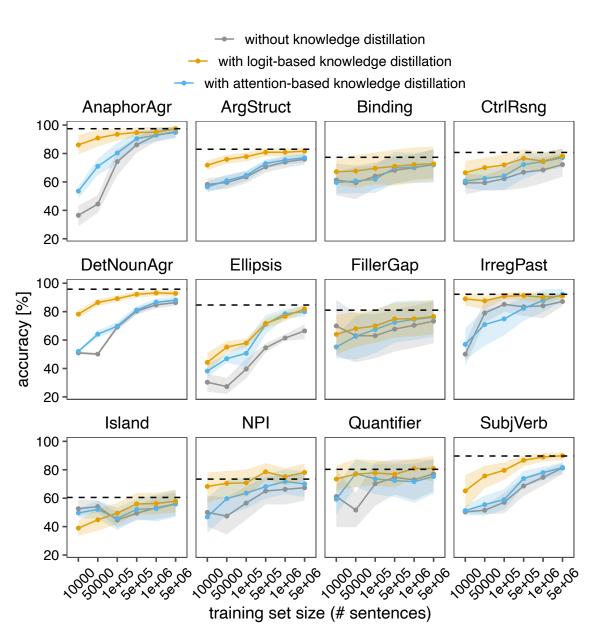


Figure S4: Performance on BLiMP split into phenomena. Ribbons show the bootstrapped 95% CI across three random seeds. Dashed lines show the teacher's performance.

without knowledge distillation with logit-based knowledge distillation with attention-based knowledge distillation

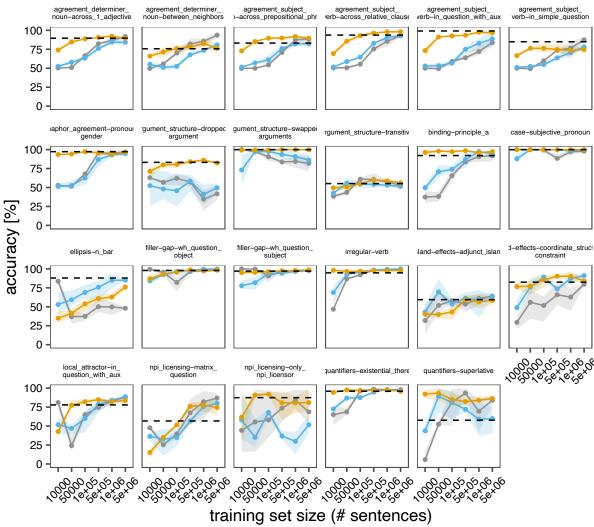


Figure S5: Performance on Zorro split into tasks. Ribbons show the bootstrapped 95% CI across three random seeds. Dashed lines show the teacher's performance.