

SLlama: Parameter-Efficient Language Model Architecture for Enhanced Linguistic Competence Under Strict Data Constraints

Anonymous ACL submission

Abstract

Scaling data and model size has driven recent advances in language modeling, but this strategy falters under scenarios with strict data constraints, as in the *BabyLM Challenge*. However, insights from *Chinchilla* paper highlights that smaller models trained on more data outperform larger counterparts trained inadequately, emphasizing the need for compact architectures. Furthermore, while embedding weight tying is a common parameter-saving technique, we find it significantly diminishes linguistic competence in compact models.

In response, we explore alternative architectural strategies that preserve the parameter-efficiency of tied models without sacrificing the representational benefits of untied embeddings. Consequently, we introduce **SLlama** a Llama3 architecture variant which incorporates targeted modifications—*Repeated Reduced Hidden Size and Projection (RRHP)*, *Permuted Weight Attention (PWA)*, *Shared Projection Multi-Layer Perceptron (SPMLP)*, and *Layer Weight Sharing*—to compress transformer components. Without relying on distillation, SLlama achieves a **31.72% improvement** in linguistic knowledge acquisition over the BabyLlama baseline, with a comparable GLUE score and significantly lower parameter count. These results demonstrate that well-designed, compact models can rival larger ones under strict data constraints.

1 Introduction

Large-scale language models (LLMs) have shown remarkable performance across a wide array of natural language understanding tasks. This success is often attributed to the trend of scaling both model size and training data, a strategy epitomized by recent architectures such as GPT-3 and LLaMA. But reliance on massive datasets and billions of parameters poses challenges when data availability is limited—a scenario increasingly relevant in

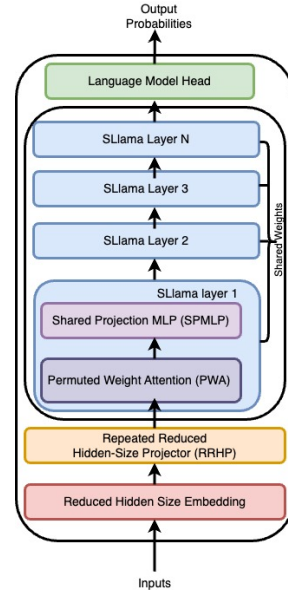


Figure 1: SLlama - Llama Architecture with Reduced Embedding, Repeated Projection, Permuted Weight Attention, Shared Projection MLP and Weight Sharing

controlled research settings like the *BabyLM Challenge*.

However, the *Chinchilla* paper offers a pivotal insight into this problem by demonstrating that, under fixed compute or data budgets, models with fewer parameters but trained on more data tend to outperform larger counterparts trained on less data. In contexts where data resource is barely 10M tokens, it is imperative to design architecturally compact models that can learn efficiently from limited data. This often spurs the adoption of embedding weight tying as a parameter-saving technique. Yet, we find that embedding weight tying impairs the linguistic competence of small models by collapsing distinct representational roles—input encoding and output prediction—into a single shared space. To address this, we investigate architectural strategies that circumvent the need for weight tying while retaining the parameter efficiency of tied models

and the representational flexibility of untied models. Our goal is to develop compact yet competent language models optimized for training on just 10 million tokens—the core constraint of the *BabyLM Challenge*.

Hence, we introduce SLlama, a parameter-efficient variant of the LLaMA3 architecture designed to balance representational capacity with parameter efficiency. SLlama leverages four key architectural innovations to reduce a model’s parameter count and maximize learning from limited data: (1) Repeated Reduced Hidden Size and Projection (RRHP), (2) Permuted Weight Attention (PWA), (3) Shared Projection Multi-Layer Perceptron (SPMLP), and (4) Hidden Layer Weight Sharing.

Crucially, SLlama (2.6M¹) is trained without distillation of teacher models. Despite this, it achieves a 31.72% improvement in linguistic knowledge acquisition over the BabyLlama² baseline (58M), maintains comparable performance on GLUE, and does so with significantly fewer parameters. These results suggest that with thoughtful architectural design, smaller models can not only survive but thrive in data-scarce environments.

1.1 Contributions

Our key contributions are:

1. We demonstrate that embedding weight tying, while widely used for model compression, it negatively impacts the linguistic competence of small models.
2. We propose and evaluate architectural strategies that eliminate the need for weight tying while preserving both compactness of weight tying and representational flexibility of untied weights, achieving a 31.72% improvement over the BabyLlama baseline.
3. Introduction of SLlama— a novel variant of the LLaMA3 architecture tailored for data-constrained settings which combines several transformer compression techniques to optimize performance under a 10M token constraint.

To ensure transparency and reproducibility, we release code, trained models, and evaluation scripts on GitHub and Hugging Face.

¹number of parameters

²A student Llama Model distilled from two teacher models Llama(360M) and GPT-2 (705M)

2 Preliminaries

The BabyLM Challenge. Choshen et al. (2024) hosted a second round of shared task where the volume of training data contains restricted to 10M tokens. The training and evaluation data contain words that children under the age of 5 years would have heard. This was to motivate small-scale pre-training which can be a sandbox for developing novel techniques for improving data efficiency. The resulting models would be evaluated³ on linguistic competence (BLiMP), conceptual understanding (GLUE), and general world knowledge (Ewok).

Among these assessments, BLiMP is of particular interest to us, as we believe that a language model, true to its name, should exhibit meaningful linguistic competence. Moreover, if such competence can be acquired from as few as 10 million tokens, we believe that collecting comparable volumes of data for low-density languages is a feasible goal. This would open the door to training pure language models—those untainted by data from other languages and thus less susceptible to cross-linguistic or cultural bias—for linguistically faithful modeling in low-density settings.

BLiMP Evaluation Unlike HELM (Liang et al., 2022), MMLU (Hendrycks et al., 2020), and FLASK (Cheng et al., 2023), which emphasize high-level task performance and alignment with user intent, BLiMP provides a fine-grained evaluation of core linguistic competence. Although older, BLiMP offers detailed probes into phenomena such as anaphor agreement, argument structure, island effects, irregular forms, and ellipsis—structures fundamental to syntactic and semantic understanding across languages. While recent frameworks reflect the evolving capabilities of large language models, they often obscure fine-grained linguistic diagnostics by focusing on derived abilities like reasoning and discourse. BLiMP, by contrast, foregrounds the grammatical structures that underlie these abilities, offering a clearer lens into a model’s linguistic fluency.

Initial Experiments. Motivated by our interest in acquiring linguistic competence from just 10 million tokens, we adopted a standard approach of sweeping over a range of model configurations—varying hidden sizes from 64 to 1024 and the number of decoder layers from 2 to 10—while

³Since models thrive on experience, the evaluation sets are filtered by the Organizers

tying the embedding layers and language model heads. While we initially expected the largest model to demonstrate the strongest linguistic competence, we were surprised to find that the best-performing model had a hidden size of 1024 and only 4 layers. Across the 24 configurations, we observed only weak correlations between model size and performance. Further details on this are given in the Appendix A.

Following Chinchilla Hoffmann et al. (2022), which recommends doubling training tokens with each doubling of model size (approximate ratio 1:20), a 10M-token budget implies an ideal model size of 5M parameters. In practice, this ratio is often relaxed. Chinchilla itself uses 70B parameters trained on 1.4T tokens. Based on this, we trained two models with a hidden size of 64 and 6 decoder layers: one with 4.4M parameters and untied weights, closely matching the theoretical target, and another with 2.4M parameters and tied weights, reflecting the more conservative, practical design. We show the performance of the two models in Table 1

Model Name	Model Size	BLiMP(%)
Small Tied	2.4M	56.0
Small Untied	4.4M	91.9
Big Tied	120M	64.5
<i>BabyLlama</i>	58M	69.8

Table 1: BLiMP scores for models of different sizes under a 10M token training budget and the baseline model. Note how the untied model extremely out-performs the tied model, earlier larger models and the baseline model. The model in italics is the baseline model.

This early result suggests that weight tying negatively affects the linguistic competence of small language models. While we defer a detailed explanation of this phenomenon to a later section, it is important to acknowledge its impact. Despite this drawback, the parameter savings from weight tying are appealing—achieving comparable performance with a 2.4M-parameter model relative to a 4.4M-parameter model offers clear advantages at scale. To mitigate the adverse effects of embedding weight tying while preserving its parameter efficiency, we introduce several parameter reduction techniques at different transformer blocks: Linear Hidden-Size Reduction and Projection (LHRP), Attention Hidden-Size Reduction and Projection (AHRP), Repeated Reduced Hidden-Size Projection (RRHP), Shared Key Query Atten-

tion (SKQA), Repeat-Reduced-Attention (RRA), Permuted Weight Attention (PWA) and Shared Projection Multi-Layer Perceptron (SPMLP). We adopted existing techniques like Hidden Layer Weight Sharing and intermediate weight reuse. In view of empirical evidences, we streamlined these reduction techniques. The techniques we adopted are collectively named SLlama.

3 Model Reduction Techniques

Recent studies have focused on minimizing the memory footprint of models by reducing parameters within the embedding layer, language model head, and MLP units (Tang et al., 2024; Liu et al., 2024; Zhang et al., 2024b). Although vocabulary size (v) reduction is a common practice (Tang et al., 2024), we chose to maintain the vocabulary size in the Hugging Face Llama3 implementation (Grattafiori et al., 2024). Our investigation of parameter reduction schemes, detailed below, focuses on the embedding layer, Feed Forward Network, and the self-attention blocks of a Transformer model.

3.1 Embedding Parameter Reduction

Inspired by the Mixed Dimension Embeddings (MDE) approach proposed by Pansare et al. (2022) and Ginart et al. (2021), we explored alternatives to embedding weight sharing by reducing the dimensionality of the embedding layer. Specifically, we reduced the hidden size (h) of the embedding layer by a factor of four (h_r). Given that the hidden layers of the decoder are initialized with h , a projection scheme is required to map the reduced embedding dimension to the original hidden size h . We investigated three such projection methods: Linear Hidden-Size Reduction and Projection (LHRP), Attention Hidden-Size Reduction and Projection (AHRP), and Repeated Reduced Hidden-Size and Projection (RRHP). LHRP employs a linear layer as described in Equation 1, effectively reducing the parameters from vh to vh_r . This method technically projects the embedding vector into a larger dimensional space, effectively assuming the relationship between the small and large representations is linear.

AHRP leverages the conventional attention mechanism described in Equation 2. AHRP utilises $vh_r + 2h_r + h^2/r$ parameters instead of vh . Conceptually, AHRP magnifies the cogent dimensions of the smaller representations. Finally, RRHP initial-

izes the embedding layer with the reduced hidden size h_r and repeats the resulting representation r times before feeding it to the decoder layers, effectively repeating the information encoded in the smaller representation r times. This method reduces the parameter count by $3vh_r$.

$$\text{Linear}(x, A) = xA^T + b \quad (1)$$

where:

$$x \in \mathbb{R}^{m \times h_r}$$

$$A \in \mathbb{R}^{h_r \times h}$$

3.2 Self-Attention Parameter Reduction

Optimized attention mechanisms with reduced complexity have shown performance comparable to standard multi-head attention (MHA) (Zhang et al., 2024a; Kitaev et al., 2020). While prior work addresses inference-time KV cache memory, our focus is on reducing the parameter count of self-attention in compact language models. Building on earlier embedding reduction strategies, we propose three lightweight attention variants: Shared Key Query Attention (SKQA), Repeat-Reduced Attention (RRA), and Permuted Weight Attention (PWA).

The design of SKQA stems from the interpretation of the attention mechanism as a similarity selection process, which is particularly relevant in language modeling. The attention weights are computed according to Equation (2), and the attention output is derived using Equation (3). Equation (2) can be viewed as computing a probability distribution of inter-token similarity when K and Q are equivalent. We investigated the feasibility of this similarity-based attention by equating the weights of K and Q; effectively reducing parameter count by h^2 .

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

$$\text{Attn}(Q, K, V) = \text{Attn_weight}(Q, K)V \quad (3)$$

where:

$$Q \in \mathbb{R}^{h_r \times h_r}$$

$$K \in \mathbb{R}^{h_r \times h_r}$$

$$V \in \mathbb{R}^{h_r \times h}$$

RRA, in contrast, was inspired by the Repeated Reduced Hidden-Size and Projection reduction

technique described earlier, that is, $Q, K, V \in \mathbb{R}^{h \times h_r}$ are subsequently repeated. Finally, PWA was motivated by the embedding layer reduction strategy presented by Li et al. (2017); Algorithm 1 illustrates its implementation. PWA effectively reduces memory demand from $4h^2$ to $6h$.

Algorithm 1 Permuted Weight Attention

Require: $h, n, m > 0$

Ensure: $\text{permutation}(n, m) > 3h$

permutes \leftarrow list of permutation(n, m)

$\theta \leftarrow \text{Embedding}(n, h)$

$q_idx \leftarrow \text{permutes}[0:h]$

$k_idx \leftarrow \text{permutes}[h:2h]$

$v_idx \leftarrow \text{permutes}[2h:3h]$

$Q = \text{Linear}(x, \theta[q_idx])$

$K = \text{Linear}(x, \theta[k_idx])$

$V = \text{Linear}(x, \theta[v_idx])$

$\text{attn} = \text{Attn}(Q, K, V)$

3.3 MLP Block Parameter Reduction

The Feed-Forward Network (FFN) in Transformers accounts for a large share of parameters, typically using two linear layers: one expanding the hidden size h to nh (with $n = 3$) and another projecting back to h , totaling $6h^2$ parameters. LLaMA adds a gated projection layer, increasing this to $7h^2$. To reduce this overhead, we propose **Shared Projection MLP (SPMLP)**, which ties the weights of the expansion and projection layers, saving $3h^2$ parameters.

3.4 Inter-Layer Weight Reduction Strategies

To further reduce model size, we explored two common inter-layer weight reduction techniques: layer reuse and weight sharing. Layer reuse (Liu et al., 2024) passes the hidden state through a layer multiple times (in our case, twice). Thus, if layer reuse $r = 2$, the model is initialized with n/r layers where n is the number of layers, effectively reducing model size by $11nh^2/2$ parameters provided no reduction scheme was introduced. On the other hand, Weight sharing (Lan et al., 2020) ties the weights of multiple layers, significantly reducing the number of parameters to $11n_g h^2$ where n_g is the number of groups the layers are divided into. We implemented both techniques, sharing weights across all layers in the model for the weight-sharing approach.

Model Block	Reduction Techniques	Model Size(M)	BLiMP (Sup.) (%)	Ewok (%)	GLUE (%)	Avg. (%)
Embedding Layer	LHRP	2.8814	60.47 (49.22)	57.58	63.26	57.63
	AHRP	2.8820	59.02 (52.80)	56.58	62.41	57.70
	RRHP	2.8803	91.94 (77.61)	57.91	63.57	72.76
Self Attention	PWA	4.3200	91.94 (77.61)	57.52	63.47	72.64
	PWA ^R	2.7800	91.94 (77.61)	57.76	63.02	72.58
	RRA	4.3400	59.20 (52.51)	57.88	63.33	58.23
	RRA ^R	2.8100	62.28 (51.65)	57.87	62.83	58.66
	SKQA	4.4200	91.94 (77.61)	58.25	63.18	72.75
	SKQA ^R	2.8800	91.94 (77.61)	57.71	63.75	72.75
Decoder Layer	Reuse	2.6700	91.94 (77.61)	57.84	63.83	72.81
	Reuse ^S	2.6700	91.94 (77.61)	57.63	62.40	72.41
	Share	2.6300	91.94 (77.61)	57.76	63.14	72.62
	Share ^S	2.6100	91.94 (77.61)	57.22	62.33	72.28

Table 2: The performance of Llama Architecture based models with with reduction techniques at different model blocks. Technique^R utilizes Repeated Reduced Hidden size Projection (RRHP). Technique^S utilizes Shared Projection MLP (SPMLP).

4 Further Experiments

We used a single NVIDIA RTX A6000 to train every model in this study.

Training Data and Hyper-parameter. Our experiments (both those described in Section 2 and this section) utilized the BabyLM challenge dataset Choshen et al. (2024), with a complete data description available in Warstadt et al. (2023a). After initial hyperparameter search, all pretraining employed cosine learning rate decay with minimum and maximum rates of 4×10^{-5} and 4×10^{-4} , respectively. We set the gradient accumulation to 2, batch size to 128, and sequence length to 256. Training runs were conducted for 3,000 iterations.

Baseline Model and Evaluation Tasks The Baby Llama model (Timiryasov and Tastet, 2023), which was among the leading solutions in the original BabyLM challenge and serves as the state-of-the-art baseline for the second BabyLM challenge⁴, was trained using knowledge distillation from two larger teacher models (Llama and GPT2), with the student model reportedly outperforming the teachers.

Evaluations was performed using the pipeline provided by Choshen et al. (2024); Gao et al. (2023), encompassing four tasks: BLiMP, BLiMP supplement (Warstadt et al., 2023c), GLUE (Wang et al., 2019), and Ewok (Ivanova et al., 2024).

⁴<https://github.com/babylm/evaluation-pipeline-2024?tab=readme-ov-file>

These tasks assess linguistic competence (BLiMP), conceptual understanding (GLUE), and general world knowledge (Ewok).

Successive Evaluations of the reduction techniques We evaluated the impact of the reduction techniques in each model block and report the results in Table 2. Linear Hidden Reduction and Projection (LHRP), Attention Hidden Reduction and Projection (AHRP), Repeated Reduced Hidden-Size and Projection (RRHP) are schemes to reduce parameter count at the embedding layer. Shared Key Query Attention (SKQA), Repeat-Reduced-Attention (RRA), and Permuted Weight Attention (PWA) were applied to the self-attention implementation. Shared Projection MLP (SPMLP) was applied to the MLP of each decoder layer. Lastly, intermediate layer reuse and inter-layer weight sharing were applied to the decoder layers.

5 Results

We begin by examining the extent to which individual reduction techniques balance parameter efficiency and model performance. Following this, we turn our attention to the combined application of the techniques which use least parameters. We refer to the combination of those techquies as SLlama. We analyse the results with the BLiMP (Warstadt et al., 2023c) framework.

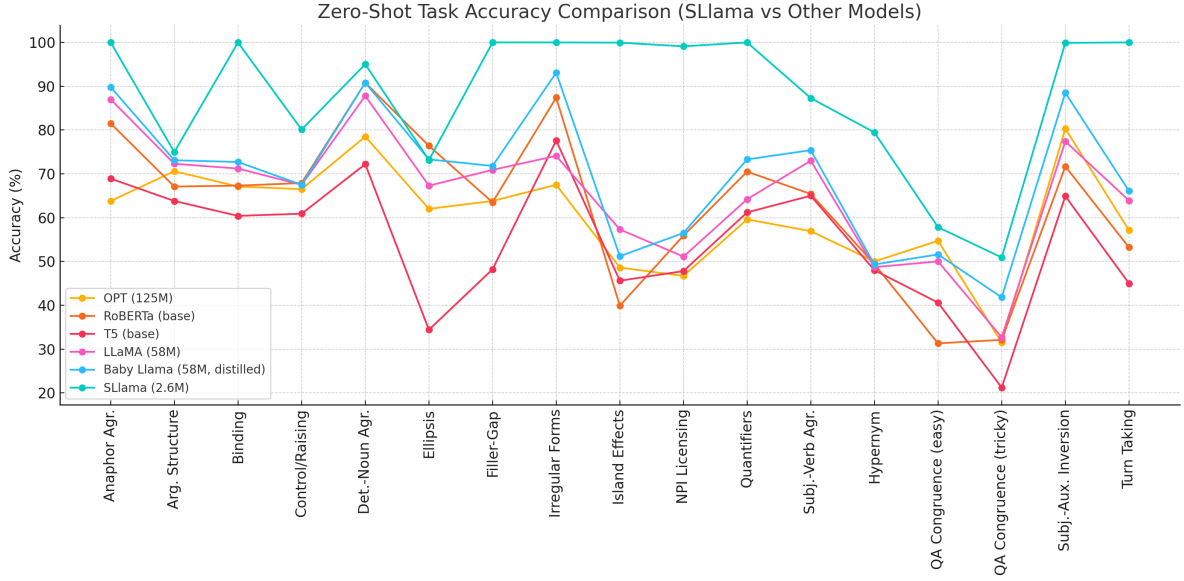


Figure 2: SLlama performance in Zero-shot BLiMP tasks relative to Baby Llama and other larger models.

5.1 Comparison of Reduction Techniques

Of the three reduction techniques applied to the embedding layer, RRHP has the optimal balance of reduction and performance as demonstrated in Table 2. Recall that, for RRHP, we divide the hidden dimension by four then repeat for further processing. This implies that the model learns salient representation of tokens which when repeated, is sufficient to undertake down stream tasks. For further experiments we discarded LHRP and AHRP

At the self attention block, PWA uses the smallest number of parameters while maintaining a competitive overall performance, closely followed by SKQA, as shown in Table 2. Relative to SKQA, PWA reduces parameter count by a larger factor but suffers performance drop. Comparing RRPH to PWA^R and SKQA^R, the performance of the latter only dropped by 0.01 while that of PWA^R dropped by 0.18. We consider this drop as a weakness of PWA. However, it's gain in parameter reduction compensates for it's weakness. We discarded RRA and SKQA from subsequent experiments.

For the MLPs, although we observe a minor decline in overall performance when SPMLP is included in the architecture, the parameter reduction remains compelling, hence, we include SPMLP in the SLlama architecture. Furthermore, Table 2 includes the performance of models that employ intermediate layer weight reuse and layer weight sharing in conjunction with SPMLP. The macro-average scores across all models show minimal variation, thus, the parameter reduction achieved

through weight-sharing presents a compelling advantage. Note that the discrepancies introduced by PWA and SPMLP in the overall performance of RRPH variants emerge from the GLUE scores and not the BLiMP scores. This signifies that our model reduction techniques are optimised for linguistic competence with a potential slight degradation of conceptual competence.

SLlama Architecture The SLlama architecture integrates reduction techniques with least parameter count while preserving competitive⁵ performance. Specifically, SLlama combines Repeated Reduced Hidden Size and Projection (RRHP), Permuted Weight Attention (PWA), Shared Projection Multi-Layer Perceptron (SPMLP), and Layer Weight Sharing to achieve architectural compactness. Compared to a similar configuration of Llama architecture, SLlama achieves a 40% reduction in parameter count without compromising linguistic competence.

5.2 Comparison with Baselines and other Models

We compared SLlama with Baby Llama (58M,distilled), OPT(125M), RoBERTa(base), T5(base), Llama2(58M),GPT-2(705M) in Figure 3 and Figure 2. All models are trained on the same BabyLM challenge dataset. Note the superiority of SLlama architecture over other models in BLiMP

⁵By competitive, we mean the drop in performance is less than 1.0

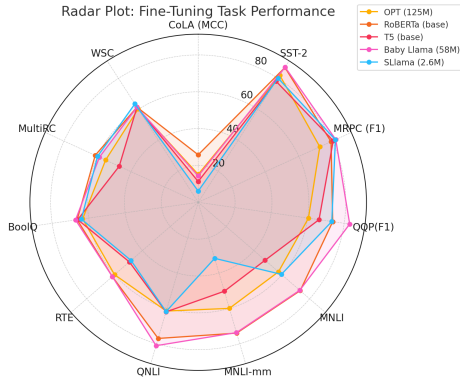


Figure 3: Comparing SLlama with other models on superGLUE

tasks maintaining the prowess of the base Llama architecture despite size reduction. Compared to Baby Llama(58M) (Timiryasov and Tastet, 2023), SLlama(2.6M) has around $20\times$ fewer parameters and improves linguistic competence by 31.72% without any knowledge distillation. It also maintains a comparable GLUE score without hyperparameter tuning.

SLlama’s Strong Generalization Across Core Grammar The BLiMP tasks span syntactic, morphological, semantic, and pragmatic domains. SLlama achieves near-perfect accuracy on core grammatical phenomena such as anaphor agreement, filler-gap dependencies, irregular forms, and quantifier interpretation. It also excels in subject–auxiliary inversion (99.9%) and binding (99.98%). Following the observations of Warstadt et al. (2023c), this performance suggests that SLlama effectively encodes core syntactic dependencies and morphological regularities despite its small size. Such strong generalization indicates that, with targeted architectural reductions, even highly compact models can acquire grammatical competence typically associated with much larger models.

Improvements Over Comparable Models Compared to Baby LLaMA (58M, distilled) and even LLaMA-360M, SLlama frequently outperforms across categories: 1. Filler-gap: SLlama (100%) > Baby LLaMA (71.8%) > LLaMA-360M (70.6%) 2. NPI licensing: SLlama (99.11%) > LLaMA-360M (57.3%) 3. Island effects: SLlama (99.95%) > LLaMA-360M (50.4%) These suggest that scaling down parameters does not necessarily reduce linguistic competence, and may even improve it when guided by effective architectural

design.

6 Discussion

We provide an explanation for the degradation in linguistic performance caused by weight tying and discuss how the employed reduction techniques shed light on language processing dynamics in parameter-efficient architectures.

6.1 Degraded Linguistic Competence with Weight Tying

By weight tying, we refer to the practice of sharing parameters between the input embedding matrix and the language model output head. As demonstrated in Table 1, this technique degrades linguistic competence in small models—a phenomenon warranting further investigation. Notably, the findings of Eldan and Li (2023); Press and Wolf (2017); Mnih and Teh (2012) offer insights that may justify this degradation.

Mnih and Teh (2012) hypothesized that when tying the embedding weights, rows corresponding to semantically similar words should exhibit near-identical representations—such that the input embedding encodes synonyms in a comparable manner, while the output embedding assigns similar score distributions to interchangeable words. Expanding on this, Press and Wolf (2017) empirically demonstrated that tying input and output embeddings produces a joint representation more closely aligned to the output embedding of an untied model.

However, their findings also suggest that untied embeddings evolve into distinct representations. By compressing these distinct roles into a shared space, weight tying limits the model’s ability to retain rich input representations essential to linguistic competence.

Furthermore, Eldan and Li (2023) confirmed that the embedding and shallow layers of a model host most linguistic information. Given that the poor performance of tied Llama are pronounced on linguistic evaluation, we conclude that the drop in performance is due to the observation of Press and Wolf (2017); that is, the embedding aligns more to the output layer and has lost salient linguistic information. Thus, empirically, untying embeddings improves performance on linguistic tasks for small language models.

This raises the question: would linguistic performance improve without reducing the hidden size?

In practice, no—LLaMA models with a 64×6 configuration and those with larger hidden sizes but tied weights perform similarly, as shown in Table 3.

6.2 Implications of the Reduction Techniques

At the embedding layer, LHRP reveals that linguistic information encoded in the embedding layer cannot be linearly projected into a higher-dimensional space without incurring a loss of critical content. Similarly, even the more expressive attention mechanism fails to reliably upscale linguistic representations without degradation. In contrast, the effectiveness of RRHP suggests that repetition, rather than projection, offers a more viable path for preserving and extending learned linguistic representations. Shared Key-Query Attention (SKQA) reframes self-attention as a linguistic operation based on token similarity. It enforces symmetry by sharing the key and query weight matrices. While future work may explore omitting one matrix entirely, such simplifications require careful evaluation. SKQA may also be less effective in asymmetric tasks like machine translation, where source–target distinctions are crucial.

Additionally, while repetition of learned embeddings (as seen in RRHP) has proven effective, our experiments with Reduced Repeated Attention (RRA) demonstrate that modifying the attention-defining neurons—particularly by altering or compressing them—can significantly impair model performance. This highlights a key asymmetry: embedding representations tolerate structural repetition, whereas the attention mechanism is more sensitive to architectural perturbations during language processing.

7 Related Work

As large models like PaLM (Chowdhery et al., 2022) and GPT-3 (Brown et al., 2020) push performance boundaries, their computational demands have prompted interest in data-efficient and compact alternatives. Data efficiency efforts include dataset reduction via k-means clustering (Kaddour, 2023), deduplication (Lee et al., 2022), and high-quality data curation (Mueller and Linzen, 2023; Eldan and Li, 2023; Gunasekar et al., 2023; Huebner et al., 2021), with studies emphasizing the role of data diversity (Lu et al., 2024; Mekala et al., 2024). We build on this by training SLlama under the 10M-token constraint of the BabyLM Challenge (Warstadt et al., 2023b,a; Choshen et al., 2024),

highlighting performance under limited data.

Compression techniques such as ROBE (Desai et al., 2022), MEmCom (Pansare et al., 2022), Mixed Dimension Embeddings (Ginart et al., 2021), and Slim Embeddings (Li et al., 2017) have reduced large embedding table sizes. For transformer models, inter-layer weight sharing and factorized embeddings (Lan et al., 2020) helped reduce BERT’s footprint (Devlin et al., 2019). Concurrently, smaller models like OPT (Zhang et al., 2022), Phi (Gunasekar et al., 2023), and PanGu- π (Tang et al., 2024) show that careful architectural design—often overlooked under fixed-compute assumptions (Kaplan et al., 2020)—can yield competitive performance. SLlama continues this trend, introducing novel reductions that preserve linguistic competence.

Weight sharing, though common (Tang et al., 2024; Lan et al., 2020; Ainslie et al., 2023), has uneven effects. While normalized shared embeddings can mitigate performance loss (Liu et al., 2020), we find that tying input-output embeddings degrades linguistic quality. In contrast, sharing attention weights (e.g., key-query) retains expressivity, suggesting that selective weight sharing is key to balancing efficiency and capability.

8 Conclusion

We introduced SLlama, a parameter-efficient adaptation of the LLaMA architecture designed for data- and scale-constrained settings like the BabyLM Challenge. Combining reduction strategies—Repeated Reduced Hidden Size and Projection (RRHP), Permuted Weight Attention (PWA), Shared Projection MLP (SPMLP), and Layer Weight Sharing—we show that small models can achieve strong linguistic performance without relying on embedding weight tying, which we find degrades linguistic competence.

Our findings suggest that repetition-based projections offer a more robust path for preserving linguistic representations than linear expansion or tied embeddings. Moreover, our analysis of SLlama’s components offers a deeper understanding of how architectural efficiency and linguistic expressivity interact, revealing design principles that extend beyond scaling.

SLlama contributes both a performant architecture and a conceptual framework for future exploration of efficient language models—particularly in low-resource or edge deployment scenarios.

Limitations

While this study demonstrates promising results, several limitations must be considered. Our findings are primarily based on the LLaMA architecture, and while certain trends may generalize, further research is needed to assess the applicability of our techniques across diverse model architectures. Additionally, the BabyLM dataset, while useful for studying small-data training, lacks linguistic diversity, limiting the evaluation of our models to English. Future work should explore performance on more diverse datasets, including low-resource languages, and assess the models' ability to acquire commonsense and factual knowledge.

Moreover, real-world deployment challenges remain, particularly regarding performance on edge devices, where quantization-related degradation has yet to be fully examined. The scalability of our compression techniques to larger models and datasets also requires further investigation. Ultimately, striking an optimal balance between model efficiency and linguistic richness is an ongoing challenge, and future research should focus on refining model reduction strategies to ensure robust language representation while maintaining computational efficiency.

References

- Joshua Ainslie, James Lee-Thorp, Michiel de Jong, Yury Zemlyanskiy, Federico Lebrón, and Sumit Sanghai. 2023. [Gqa: Training generalized multi-query transformer models from multi-head checkpoints](#). *Preprint*, arXiv:2305.13245.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). *Preprint*, arXiv:2005.14165.
- Eric Zelikman Cheng, Eric Zhao, Swaroop Mishra, John Canny, Tianyi Zhang, and James Zou. 2023. [Flask: Fine-grained language model evaluation based on alignment skill sets](#). *arXiv preprint arXiv:2307.10928*.
- Leshem Choshen, Ryan Cotterell, Michael Y. Hu, Tal Linzen, Aaron Mueller, Candace Ross, Alex Warstadt, Ethan Wilcox, Adina Williams, and

Chengxu Zhuang. 2024. [\[call for papers\] the 2nd BabyLM Challenge: Sample-efficient pretraining on a developmentally plausible corpus](#). *Computing Research Repository*, arXiv:2404.06214.

- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. 2022. [Palm: Scaling language modeling with pathways](#). *Preprint*, arXiv:2204.02311.
- Aditya Desai, Li Chou, and Anshumali Shrivastava. 2022. [Random offset block embedding array \(robe\) for criteotb benchmark mlperf dlrm model : 1000× compression and 3.1× faster inference](#). *Preprint*, arXiv:2108.02191.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [Bert: Pre-training of deep bidirectional transformers for language understanding](#). *Preprint*, arXiv:1810.04805.
- Ronen Eldan and Yuanzhi Li. 2023. [Tinystories: How small can language models be and still speak coherent english?](#) *Preprint*, arXiv:2305.07759.
- Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac'h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. 2023. [A framework for few-shot language model evaluation](#).
- A.A. Ginart, Maxim Naumov, Dheevatsa Mudigere, Jiyan Yang, and James Zou. 2021. [Mixed dimension embeddings with application to memory-efficient recommendation systems](#). In *2021 IEEE International Symposium on Information Theory (ISIT)*, page 2786–2791. IEEE Press.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh

732	Goyal, Anthony Hartshorn, Aobo Yang, Archi Mi-	ney Meers, Xavier Martinet, Xiaodong Wang, Xi-	796
733	tra, Archie Sravankumar, Artem Korenev, Arthur	aofang Wang, Xiaoqing Ellen Tan, Xide Xia, Xin-	797
734	Hinsvark, Arun Rao, Aston Zhang, Aurelien Ro-	feng Xie, Xuchao Jia, Xuwei Wang, Yaelle Gold-	798
735	driguez, Austen Gregerson, Ava Spataru, Baptiste	schlag, Yashesh Gaur, Yasmine Babaei, Yi Wen,	799
736	Roziere, Bethany Biron, Binh Tang, Bobbie Chern,	Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao,	800
737	Charlotte Caucheteux, Chaya Nayak, Chloe Bi,	Zacharie Delpierre Coudert, Zheng Yan, Zhengxing	801
738	Chris Marra, Chris McConnell, Christian Keller,	Chen, Zoe Papakipos, Aaditya Singh, Aayushi Sri-	802
739	Christophe Touret, Chunyang Wu, Corinne Wong,	vastava, Abha Jain, Adam Kelsey, Adam Shajnfeld,	803
740	Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Al-	Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand,	804
741	lonsius, Daniel Song, Danielle Pintz, Danny Livshits,	Ajay Menon, Ajay Sharma, Alex Boesenberg, Alexei	805
742	Danny Wyatt, David Esiobu, Dhruv Choudhary,	Baevski, Allie Feinstein, Amanda Kallet, Amit San-	806
743	Dhruv Mahajan, Diego Garcia-Olano, Diego Perino,	gani, Amos Teo, Anam Yunus, Andrei Lupu, An-	807
744	Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy,	dres Alvarado, Andrew Caples, Andrew Gu, Andrew	808
745	Elina Lobanova, Emily Dinan, Eric Michael Smith,	Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchan-	809
746	Filip Radenovic, Francisco Guzmán, Frank Zhang,	dani, Annie Dong, Annie Franco, Anuj Goyal, Aparajita	810
747	Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis An-	Saraf, Arkabandhu Chowdhury, Ashley Gabriel,	811
748	derson, Govind Thattai, Graeme Nail, Gregoire Mi-	Ashwin Bharambe, Assaf Eisenman, Azadeh Yaz-	812
749	alon, Guan Pang, Guillem Cucurell, Hailey Nguyen,	dan, Beau James, Ben Maurer, Benjamin Leonhardi,	813
750	Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan	Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi	814
751	Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Is-	Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Han-	815
752	han Misra, Ivan Evtimov, Jack Zhang, Jade Copet,	cock, Bram Wasti, Brandon Spence, Brani Stojkovic,	816
753	Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park,	Brian Gamido, Britt Montalvo, Carl Parker, Carly	817
754	Jay Mahadeokar, Jeet Shah, Jelmer van der Linde,	Burton, Catalina Mejia, Ce Liu, Changan Wang,	818
755	Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu,	Changkyu Kim, Chao Zhou, Chester Hu, Ching-	819
756	Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang,	Hsiang Chu, Chris Cai, Chris Tindal, Christoph Fe-	820
757	Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park,	ichtenhofer, Cynthia Gao, Damon Civin, Dana Beaty,	821
758	Joseph Rocca, Joshua Johnstun, Joshua Saxe, Jun-	Daniel Kreymmer, Daniel Li, David Adkins, David	822
759	teng Jia, Kalyan Vasuden Alwala, Karthik Prasad,	Xu, Davide Testuggine, Delia David, Devi Parikh,	823
760	Kartikaya Upasani, Kate Plawiak, Ke Li, Kenneth	Diana Liskovich, Didem Foss, Dingkan Wang, Duc	824
761	Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer,	Le, Dustin Holland, Edward Dowling, Eissa Jamil,	825
762	Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Kushal	Elaine Montgomery, Eleonora Presani, Emily Hahn,	826
763	Lakhotia, Lauren Rantala-Yearly, Laurens van der	Emily Wood, Eric-Tuan Le, Erik Brinkman, Este-	827
764	Maaten, Lawrence Chen, Liang Tan, Liz Jenkins,	ban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun,	828
765	Louis Martin, Lovish Madaan, Lubo Malo, Lukas	Felix Kreuk, Feng Tian, Filippas Kokkinos, Firat	829
766	Blecher, Lukas Landzaat, Luke de Oliveira, Madeline	Ozgenel, Francesco Caggioni, Frank Kanayet, Frank	830
767	Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar	Seide, Gabriela Medina Florez, Gabriella Schwarz,	831
768	Paluri, Marcin Kardas, Maria Tsimpoukelli, Mathew	Gada Badeer, Georgia Swee, Gil Halpern, Grant	832
769	Oldham, Mathieu Rita, Maya Pavlova, Melanie Kam-	Herman, Grigory Sizov, Guangyi, Zhang, Guna	833
770	badur, Mike Lewis, Min Si, Mitesh Kumar Singh,	Lakshminarayanan, Hakan Inan, Hamid Shojanaz-	834
771	Mona Hassan, Naman Goyal, Narjes Torabi, Niko-	eri, Han Zou, Hannah Wang, Hanwen Zha, Haroun	835
772	lay Bashlykov, Nikolay Bogoychev, Niladri Chatterji,	Habeeb, Harrison Rudolph, Helen Suk, Henry As-	836
773	Ning Zhang, Olivier Duchenne, Onur Çelebi, Patrick	pegren, Hunter Goldman, Hongyuan Zhan, Ibrahim	837
774	Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vas-	Damlaj, Igor Molybog, Igor Tufanov, Ilias Leontiadis,	838
775	asic, Peter Weng, Prajjwal Bhargava, Pratik Dubal,	Irina-Elena Veliche, Itai Gat, Jake Weissman, James	839
776	Praveen Krishnan, Punit Singh Koura, Puxin Xu,	Geboski, James Kohli, Janice Lam, Japhet Asher,	840
777	Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj	Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jen-	841
778	Ganapathy, Ramon Calderer, Ricardo Silveira Cabral,	nifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy	842
779	Robert Stojnic, Roberta Raileanu, Rohan Maheswari,	Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe	843
780	Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ron-	Cummings, Jon Carvill, Jon Shepard, Jonathan Mc-	844
781	nie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan	Phie, Jonathan Torres, Josh Ginsburg, Junjie Wang,	845
782	Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sa-	Kai Wu, Kam Hou U, Karan Saxena, Kartikay Khand-	846
783	hana Chennabasappa, Sanjay Singh, Sean Bell, Seo-	delwal, Katayoun Zand, Kathy Matosich, Kaushik	847
784	hyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sha-	Veeraraghavan, Kelly Michelena, Keqian Li, Kiran	848
785	ran Narang, Sharath Rapparthu, Sheng Shen, Shengye	Jagadeesh, Kun Huang, Kunal Chawla, Kyle	849
786	Wan, Shruti Bhosale, Shun Zhang, Simon Van-	Huang, Lailin Chen, Lakshya Garg, Lavender A,	850
787	denhende, Soumya Batra, Spencer Whitman, Sten	Leandro Silva, Lee Bell, Lei Zhang, Liangpeng	851
788	Sootla, Stephane Collot, Suchin Gururangan, Syd-	Guo, Licheng Yu, Liron Moshkovich, Luca Wehrst-	852
789	ney Borodinsky, Tamar Herman, Tara Fowler, Tarek	edt, Madian Khabsa, Manav Avalani, Manish Bhatt,	853
790	Sheasha, Thomas Georgiou, Thomas Scialom, Tobias	Martynas Mankus, Matan Hasson, Matthew Lennie,	854
791	Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal	Matthias Reso, Maxim Groshev, Maxim Naumov,	855
792	Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh	Maya Lathi, Meghan Keneally, Miao Liu, Michael L.	856
793	Ramanathan, Viktor Kerkez, Vincent Gouget, Vir-	Seltzer, Michal Valko, Michelle Restrepo, Mihir Pa-	857
794	ginie Do, Vish Vogeti, Vitor Albiero, Vladan Petro-	tel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark,	858
795	vic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whit-	Mike Macey, Mike Wang, Miquel Jubert Hermoso,	859

860	Mo Metanat, Mohammad Rastegari, Munish Bansal,	and Laurent Sifre. 2022. Training compute-optimal	921
861	Nandhini Santhanam, Natascha Parks, Natasha	large language models . <i>Preprint</i> , arXiv:2203.15556.	922
862	White, Navyata Bawa, Nayan Singhal, Nick Egebo,		
863	Nicolas Usunier, Nikhil Mehta, Nikolay Pavlovich	Philip A. Huebner, Elior Sulem, Fisher Cynthia, and	923
864	Laptev, Ning Dong, Norman Cheng, Oleg Chernoguz,	Dan Roth. 2021. BabyBERTa: Learning more gram-	924
865	Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin	mar with small-scale child-directed language . In <i>Pro-</i>	925
866	Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pe-	<i>ceedings of the 25th Conference on Computational</i>	926
867	dro Rittner, Philip Bontrager, Pierre Roux, Piotr	<i>Natural Language Learning</i> , pages 624–646, Online.	927
868	Dollar, Polina Zvyagina, Prashant Ratanchandani,	Association for Computational Linguistics.	928
869	Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel		
870	Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu	Anna A. Ivanova, Aalok Sathe, Benjamin Lipkin, Un-	929
871	Nayani, Rahul Mitra, Rangaprabhu Parthasarathy,	nathi Kumar, Setayesh Radkani, Thomas H. Clark,	930
872	Raymond Li, Rebekkah Hogan, Robin Battey, Rocky	Carina Kauf, Jennifer Hu, R. T. Pramod, Gabriel	931
873	Wang, Russ Howes, Rutu Rinott, Sachin Mehta,	Grand, Vivian Paulun, Maria Ryskina, Ekin Akyürek,	932
874	Sachin Siby, Sai Jayesh Bondu, Samyak Datta, Sara	Ethan Wilcox, Nafisa Rashid, Leshem Choshen,	933
875	Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov,	Roger Levy, Evelina Fedorenko, Joshua Tenenbaum,	934
876	Satadru Pan, Saurabh Mahajan, Saurabh Verma,	and Jacob Andreas. 2024. Elements of world knowl-	935
877	Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lind-	edge (ewok): A cognition-inspired framework for	936
878	say, Shaun Lindsay, Sheng Feng, Shenghao Lin,	evaluating basic world knowledge in language mod-	937
879	Shengxin Cindy Zha, Shishir Patil, Shiva Shankar,	els . <i>Preprint</i> , arXiv:2405.09605.	938
880	Shuqiang Zhang, Shuqiang Zhang, Sinong Wang,		
881	Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala,	Jean Kaddour. 2023. The minipile challenge	939
882	Stephanie Max, Stephen Chen, Steve Kehoe, Steve	for data-efficient language models . <i>Preprint</i> ,	940
883	Satterfield, Sudarshan Govindaprasad, Sumit Gupta,	arXiv:2304.08442.	941
884	Summer Deng, Sungmin Cho, Sunny Virk, Suraj		
885	Subramanian, Sy Choudhury, Sydney Goldman, Tal	Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B.	942
886	Remez, Tamar Glaser, Tamara Best, Thilo Koehler,	Brown, Benjamin Chess, Rewon Child, Scott Gray,	943
887	Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim	Alec Radford, Jeffrey Wu, and Dario Amodei. 2020.	944
888	Matthews, Timothy Chou, Tzook Shaked, Varun	Scaling laws for neural language models . <i>Preprint</i> ,	945
889	Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai	arXiv:2001.08361.	946
890	Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad		
891	Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu,	Nikita Kitaev, Łukasz Kaiser, and Anselm Levskaya.	947
892	Vladimir Ivanov, Wei Li, Wenchen Wang, Wen-	2020. Reformer: The efficient transformer . <i>Preprint</i> ,	948
893	wen Jiang, Wes Bouaziz, Will Constable, Xiaocheng	arXiv:2001.04451.	949
894	Tang, Xiaojian Wu, Xiaolan Wang, Xilun Wu, Xinbo		
895	Gao, Yaniv Kleinman, Yanjun Chen, Ye Hu, Ye Jia,	Zhenzhong Lan, Mingda Chen, Sebastian Goodman,	950
896	Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi,	Kevin Gimpel, Piyush Sharma, and Radu Sori-	951
897	Youngjin Nam, Yu, Wang, Yu Zhao, Yuchen Hao,	cut. 2020. Albert: A lite bert for self-supervised	952
898	Yundi Qian, Yunlu Li, Yuzi He, Zach Rait, Zachary	learning of language representations . <i>Preprint</i> ,	953
899	DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang,	arXiv:1909.11942.	954
900	Zhiwei Zhao, and Zhiyu Ma. 2024. The llama 3 herd		
901	of models . <i>Preprint</i> , arXiv:2407.21783.	Katherine Lee, Daphne Ippolito, Andrew Nystrom,	955
		Chiyuan Zhang, Douglas Eck, Chris Callison-Burch,	956
902	Suriya Gunasekar, Yi Zhang, Jyoti Aneja, Caio	and Nicholas Carlini. 2022. Deduplicating train-	957
903	César Teodoro Mendes, Allie Del Giorno, Sivakanth	ing data makes language models better . <i>Preprint</i> ,	958
904	Gopi, Mojan Javaheripi, Piero Kauffmann, Gustavo	arXiv:2107.06499.	959
905	de Rosa, Olli Saarikivi, Adil Salim, Shital Shah,		
906	Harkirat Singh Behl, Xin Wang, Sébastien Bubeck,	Zhongliang Li, Raymond Kulhanek, Shaojun Wang,	960
907	Ronen Eldan, Adam Tauman Kalai, Yin Tat Lee,	Yunxin Zhao, and Shuang Wu. 2017. Slim embed-	961
908	and Yuanzhi Li. 2023. Textbooks are all you need .	ding layers for recurrent neural language models .	962
909	<i>Preprint</i> , arXiv:2306.11644.	<i>Preprint</i> , arXiv:1711.09873.	963
910	Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou,	Percy Liang, Rishi Bommasani, Tony Lee, et al. 2022.	964
911	Mantas Mazeika, Dawn Song, and Jacob Steinhardt.	Holistic evaluation of language models . In <i>Proce-</i>	965
912	2020. Measuring massive multitask language under-	<i>edings of NeurIPS 2022 Track on Datasets and Bench-</i>	966
913	standing . <i>arXiv preprint arXiv:2009.03300</i> .	<i>marks</i> .	967
914	Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch,	Jinyang Liu, Yujia Zhai, and Zizhong Chen. 2020. Nor-	968
915	Elena Buchatskaya, Trevor Cai, Eliza Rutherford,	malization of input-output shared embeddings in text	969
916	Diego de Las Casas, Lisa Anne Hendricks, Johannes	generation models . <i>Preprint</i> , arXiv:2001.07885.	970
917	Welbl, Aidan Clark, Tom Hennigan, Eric Noland,		
918	Katie Millican, George van den Driessche, Bogdan	Zechun Liu, Changsheng Zhao, Forrest Iandola, Chen	971
919	Damoc, Aurelia Guy, Simon Osindero, Karen Si-	Lai, Yuandong Tian, Igor Fedorov, Yunyang Xiong,	972
920	monyhan, Erich Elsen, Jack W. Rae, Oriol Vinyals,	Ernie Chang, Yangyang Shi, Raghuraman Krish-	973
		namoorthi, Liangzhen Lai, and Vikas Chandra. 2024.	974

975	Mobilellm: Optimizing sub-billion parameter language models for on-device use cases. <i>Preprint</i> , arXiv:2402.14905.	
976		
977		
978	Zhenyan Lu, Xiang Li, Dongqi Cai, Rongjie Yi, Fangming Liu, Xiwen Zhang, Nicholas D. Lane, and Mengwei Xu. 2024. <i>Small language models: Survey, measurements, and insights. Preprint</i> , arXiv:2409.15790.	
979		
980		
981		
982		
983	Dheeraj Mekala, Alex Nguyen, and Jingbo Shang. 2024. <i>Smaller language models are capable of selecting instruction-tuning training data for larger language models. Preprint</i> , arXiv:2402.10430.	
984		
985		
986		
987	Andriy Mnih and Yee Whye Teh. 2012. <i>A fast and simple algorithm for training neural probabilistic language models. Preprint</i> , arXiv:1206.6426.	
988		
989		
990	Aaron Mueller and Tal Linzen. 2023. <i>How to plant trees in language models: Data and architectural effects on the emergence of syntactic inductive biases. Preprint</i> , arXiv:2305.19905.	
991		
992		
993		
994	Niketan Pansare, Jay Katukuri, Aditya Arora, Frank Cipollone, Riyaaz Shaik, Noyan Tokgozoglu, and Chandru Venkataraman. 2022. <i>Learning compressed embeddings for on-device inference. Preprint</i> , arXiv:2203.10135.	
995		
996		
997		
998		
999	Ofir Press and Lior Wolf. 2017. <i>Using the output embedding to improve language models. Preprint</i> , arXiv:1608.05859.	
1000		
1001		
1002	Yehui Tang, Fangcheng Liu, Yunsheng Ni, Yuchuan Tian, Zheyuan Bai, Yi-Qi Hu, Sichao Liu, Shangling Jui, Kai Han, and Yunhe Wang. 2024. <i>Rethinking optimization and architecture for tiny language models. Preprint</i> , arXiv:2402.02791.	
1003		
1004		
1005		
1006		
1007	Inar Timiryasov and Jean-Loup Tastet. 2023. <i>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</i> , pages 279–289, Singapore. Association for Computational Linguistics.	
1008		
1009		
1010		
1011		
1012		
1013		
1014	Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. 2019. <i>Glue: A multi-task benchmark and analysis platform for natural language understanding. Preprint</i> , arXiv:1804.07461.	
1015		
1016		
1017		
1018		
1019	Alex Warstadt, Leshem Choshen, Aaron Mueller, Adina Williams, Ethan Wilcox, and Chengxu Zhuang. 2023a. <i>Call for papers – the babyLM challenge: Sample-efficient pretraining on a developmentally plausible corpus. Preprint</i> , arXiv:2301.11796.	
1020		
1021		
1022		
1023		
1024	Alex Warstadt, Aaron Mueller, Leshem Choshen, Ethan Wilcox, Chengxu Zhuang, Juan Ciro, Rafael Mosquera, Bhargavi Paranjabe, Adina Williams, Tal Linzen, and Ryan Cotterell. 2023b. <i>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</i> , pages 1–34, Singapore. Association for Computational Linguistics.	1030
1025		1031
1026		1032
1027		1033
1028		
1029		
	Alex Warstadt, Alicia Parrish, Haokun Liu, Anhad Mohananey, Wei Peng, Sheng-Fu Wang, and Samuel R. Bowman. 2023c. <i>Blimp: The benchmark of linguistic minimal pairs for english. Preprint</i> , arXiv:1912.00582.	1034
		1035
		1036
		1037
		1038
	Jiale Zhang, Yulun Zhang, Jinjin Gu, Jiahua Dong, Linghe Kong, and Xiaokang Yang. 2024a. <i>Xformer: Hybrid x-shaped transformer for image denoising. Preprint</i> , arXiv:2303.06440.	1039
		1040
		1041
		1042
	Peiyuan Zhang, Guangtao Zeng, Tianduo Wang, and Wei Lu. 2024b. <i>Tinyllama: An open-source small language model. Preprint</i> , arXiv:2401.02385.	1043
		1044
		1045
	Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, Todor Mihaylov, Myle Ott, Sam Shleifer, Kurt Shuster, Daniel Simig, Punit Singh Koura, Anjali Sridhar, Tianlu Wang, and Luke Zettlemoyer. 2022. <i>Opt: Open pre-trained transformer language models. Preprint</i> , arXiv:2205.01068.	1046
		1047
		1048
		1049
		1050
		1051
		1052
		1053
	A Initial Experiments of Model Sweep	1054
	Characterizing the Llama Architecture To isolate the effect of distillation, we conducted experiments to characterize the inherent capabilities of the Llama architecture and to establish the relationship between its key configuration parameters (hidden size, intermediate size, and number of layers) and performance on the aforementioned evaluation tasks. Following the recommendations of Tang et al. (2024), we tie the embedding layer and language model head, a widely used strategy to improve parameter efficiency in small-scale language models. Starting with a hidden size of 64 (to minimize resource consumption), we varied the number of layers from 2 to 12.	1055
		1056
		1057
		1058
		1059
		1060
		1061
		1062
		1063
		1064
		1065
		1066
		1067
		1068
		1069
		1070
		1071
		1072
		1073
		1074
		1075
		1076
		1077
		1078
		1079
		1080
		1081

latter configuration (hidden size 64 and 6 layers). Finally, to ascertain the plausibility of weight tying, we trained a 64 by 6 model with untied weights.

Characterizing the Llama Architecture We present the results of the experiment to characterize the inherent ability of Llama architecture without distillation in Table 1. We observed that the relationship between macroscore and model size is not direct. Further analysis presented in Figure 4, shows the correlation between model size parameters (hidden size and number of layers) and the model’s performance across the different evaluation dimensions (linguistic competence, world knowledge, and conceptual understanding). While statistical significance was generally weak, several trends emerged: 1) a weak but consistent positive correlation between hidden size and BLiMP score (linguistic knowledge); 2) an inconsistent positive relationship between hidden size and GLUE score; 3) a strong and consistent negative correlation between hidden size and world knowledge; 4) an inconsistent positive trend between the number of layers and linguistic competence; 5) a weak positive trend between the number of layers and conceptual understanding; and 6) a noticeable weak negative trend between the number of layers and linguistic competence. While these observations suggest the need to carefully balance horizontal (hidden size) and vertical (number of layers) scaling, particularly while training on limited data, more data is needed to fully concretize these claims. However, the positive impact of increasing layer count for smaller hidden sizes was evident, supporting previous findings of Liu et al. (2024).

The result in 3 influenced our hyper-parameter selection.

B Architectural Comparison

We compare different language model architectures and present them in Table 4. All models are trained on the same dataset but for different epochs.

C SuperGLUE scores

We also include the performance of other architecture reported in other studies in Table 5. While our emphasis is not superGLUE, it is noteworthy to demonstrate that the architecture maintains decent conceptual competence relative to the larger models.

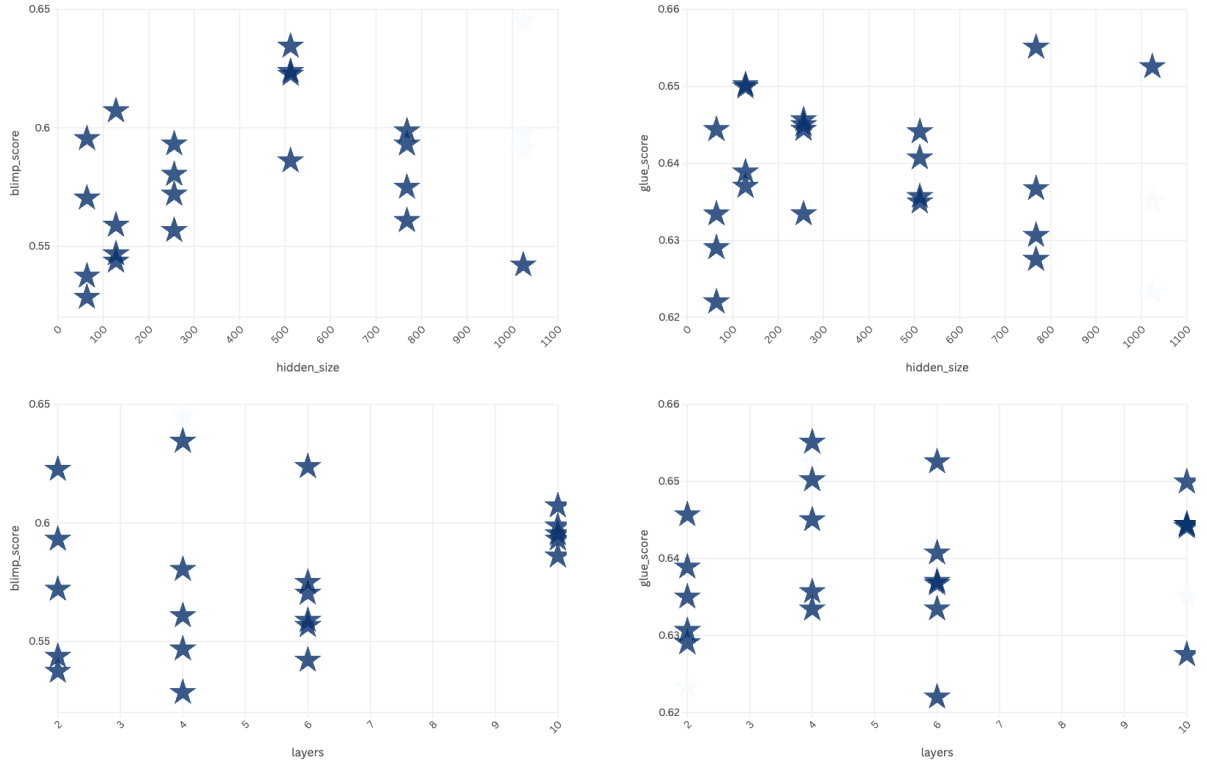


Figure 4: Correlation of hidden size and number of layers to BLiMP and GLUE scores with Spearman correlations of 0.38 and 0.88, respectively.

Model	Hidden Size	Layers	BLiMP	GLUE	EWoK	BLiMP-Sup	Macro Avg.
model1	1024	10	59.75%	63.51%	54.09%	49.45%	56.70%
model2	1024	2	59.10%	62.33%	53.86%	52.31%	56.90%
model3	1024	4	64.53%	65.24%	53.30%	48.47%	57.89%
model4	1024	6	54.21%	65.25%	53.94%	50.87%	56.07%
model5	128	10	60.72%	64.99%	56.14%	48.48%	57.58%
model6	128	2	54.37%	63.89%	55.96%	48.51%	55.68%
model7	128	4	54.68%	65.02%	56.10%	53.65%	57.36%
model8	128	6	55.89%	63.70%	56.96%	48.48%	56.26%
model9	256	10	59.32%	64.44%	55.65%	51.32%	57.68%
model10	256	2	57.20%	64.57%	55.13%	56.06%	58.24%
model11	256	4	58.03%	64.50%	55.72%	50.43%	57.17%
model12	256	6	55.67%	63.34%	55.96%	49.66%	56.16%
model13	512	10	58.60%	64.41%	55.13%	46.85%	56.25%
model14	512	2	62.26%	63.50%	55.36%	53.28%	58.60%
model15	512	4	63.44%	63.57%	55.73%	48.21%	57.74%
model16	512	6	62.37%	64.07%	55.59%	51.58%	58.40%
model17	64	10	59.54%	64.44%	58.01%	49.77%	57.94%
model18	64	2	53.74%	62.90%	57.99%	55.07%	57.42%
model19	64	4	52.85%	63.34%	57.91%	48.00%	55.52%
model20	64	6	57.03%	62.20%	57.02%	48.71%	56.24%
model21	768	10	59.87%	62.75%	54.85%	49.57%	56.76%
model22	768	2	59.31%	63.06%	54.48%	54.54%	57.85%
model23	768	4	56.08%	65.51%	54.07%	53.49%	57.29%
model24	768	6	57.49%	63.67%	53.74%	53.18%	57.02%

Table 3: Evaluation scores across models with varying hidden sizes and number of layers. Best values per metric are in bold.

Task	OPT (125M)	RoBERTa (base)	T5 (base)	LLaMA2 (58M)	LLaMA2 (360M)	GPT-2 (705M)	Baby LLaMA (58M, distilled)	SLlama (2.6M)
Anaphor Agr.	63.80	81.50	68.90	87.00	87.60	89.60	89.80	100.00
Arg. Structure	70.60	67.10	63.80	72.30	73.50	73.50	73.10	74.98
Binding	67.10	67.30	60.40	71.20	72.10	71.50	72.70	99.98
Control/Raising	66.50	67.90	60.90	67.50	67.40	68.40	67.50	80.11
Det.-Noun Agr.	78.50	90.80	72.20	87.80	89.60	87.40	90.80	95.03
Ellipsis	62.00	76.40	34.40	67.30	68.50	69.90	73.30	73.13
Filler-Gap	63.80	63.50	48.20	70.90	70.60	70.20	71.80	100.00
Irregular Forms	67.50	87.40	77.60	74.10	68.90	83.10	93.10	100.00
Island Effects	48.60	39.90	45.60	57.30	50.40	51.60	51.20	99.95
NPI Licensing	46.70	55.90	47.80	51.10	57.30	50.50	56.50	99.11
Quantifiers	59.60	70.50	61.20	64.20	59.00	69.80	73.30	100.00
Subj.-Verb Agr.	56.90	65.40	65.00	73.00	69.70	67.50	75.40	87.29
Hypernym	50.00	49.40	48.00	48.70	49.40	49.20	49.30	79.45
QA Congruence (easy)	54.70	31.30	40.60	50.00	53.10	56.20	51.60	57.81
QA Congruence (tricky)	31.50	32.10	21.20	32.70	41.80	45.50	41.80	50.91
Subj.-Aux. Inversion	80.30	71.70	64.90	77.40	84.30	81.70	88.50	99.90
Turn Taking	57.10	53.20	45.00	63.90	68.60	65.70	66.10	100.00

Table 4: Comparative performance of SLlama and larger models on BLiMP tasks. Results for baseline models (OPT, RoBERTa, T5, LLaMA, GPT-2, and Baby LLaMA) are taken from the original baseline paper. All models, including SLlama, are trained on the same 10M-token dataset.

Task	OPT (125M)	RoBERTa (base)	T5 (base)	Baby Llama (58M)	SLlama (2.6M)
CoLA (MCC)	15.2	25.8	11.3	14.3	6.1
SST-2	81.9	87.0	78.1	87.2	80.1
MRPC (F1)	72.5	79.2	80.5	82.0	81.8
QQP(F1)	60.4	73.7	66.2	83.0	73.2
MNLI	57.6	73.2	48.0	72.9	59.7
MNLI-mm	60.0	74.0	50.3	73.7	31.7
QNLI	61.5	77.0	62.0	81.1	61.8
RTE	60.0	61.6	49.4	61.6	48.2
BoolQ	63.3	66.3	66.0	67.2	64.0
MultiRC	55.2	61.4	47.1	58.9	60.0
WSC	60.2	61.4	61.4	61.4	63.5

Table 5: Performance on super GLUE

. Evaluations are accuracy except specified otherwise. All models are pretrained on the training dataset.

Parameter	Value
eval_interval	250
log_interval	10
eval_iters	150
eval_only	False
wandb_log	True
wandb_project	further_baby_experiments
wandb_run_name	test_llama_c_alone
gradient_accumulation_steps	2
batch_size	128
block_size	block_size
dropout	0.1
learning_rate	4e-4
max_iters	3000
weight_decay	0.0
warmup_iters	200
lr_decay_iters	5000
min_lr	4e-5
train_or_dev	train
out_dir	test_baby

Table 6: Training configuration for SLlama experiments