

Attention Alignment Outperforms Logit Distillation for LLM Compression: A Comparative Study of White-Box vs Black-Box Knowledge Transfer

Jack Large
Independent Researcher

Abstract

Knowledge distillation enables compressing large language models (LLMs) into smaller, deployable variants. While standard “black-box” distillation transfers only output logits, “white-box” methods additionally align internal representations such as hidden states and attention maps. We systematically compare four distillation approaches for compressing Llama-2-7B into TinyLlama-1.1B across sentiment analysis (SST-2), reasoning (MMLU), and mathematical problem-solving (GSM8K) tasks.

Our experiments ($N = 7$ seeds per method) reveal that **attention alignment achieves +1.58% higher accuracy** (95.56% vs 93.98%) compared to logit-only distillation. Notably, combining attention with hidden state alignment yields no additional benefit, suggesting attention maps are the primary driver of improvement. White-box methods also exhibit substantially lower variance (1.43% vs 2.98% std), indicating more stable training dynamics.

These findings demonstrate that for heterogeneous teacher-student pairs with matched attention head counts, attention distillation provides a simple, effective enhancement over standard knowledge distillation with minimal computational overhead during training.

1 Introduction

Large language models (LLMs) have achieved remarkable performance across diverse NLP tasks, but their size—often billions of parameters—poses significant deployment challenges. Knowledge distillation [Hinton et al., 2015] offers a promising compression approach by training a smaller “student” model to mimic a larger “teacher” model.

Traditional distillation operates in a “black-box” manner, transferring only the teacher’s output logits. However, LLMs encode rich intermediate representations: hidden states capture semantic features, while attention maps encode structural relationships between tokens. Can transferring these internal signals—“white-box” distillation—improve student performance?

We investigate this question by comparing four distillation methods:

1. **Black-Box:** Standard KL-divergence on output logits
2. **Hidden State:** Logits + final-layer hidden state alignment
3. **Attention:** Logits + final-layer attention map alignment
4. **Combined:** All signals (logits + hidden states + attention)

Our experiments compress Llama-2-7B [Touvron et al., 2023] into TinyLlama-1.1B [Zhang et al., 2024]—a $6.4\times$ compression ratio. We evaluate on three diverse tasks: sentiment classification (SST-2), multi-task reasoning (MMLU), and grade-school mathematics (GSM8K).

Our key finding: **attention alignment provides consistent, significant improvement** (+1.58% accuracy) over black-box distillation, while hidden state alignment alone offers modest gains (+0.73%). Combining signals yields no additional benefit over attention-only, suggesting attention maps are the dominant transfer signal for heterogeneous LLM distillation.

2 Related Work

Knowledge Distillation. Hinton et al. [2015] introduced knowledge distillation, showing that soft probability distributions from a teacher network provide richer supervision than hard labels. This approach has been extensively applied to model compression in NLP.

Feature-Based Distillation. FitNets [Romero et al., 2015] extended distillation to intermediate representations, training students to match teacher hidden states. This “hint-based” training accelerates convergence and improves generalization. Subsequent work explored attention transfer [Zagoruyko and Komodakis, 2017], demonstrating that attention maps encode valuable structural information.

Distillation for Transformers. DistilBERT [Sanh et al., 2019] successfully applied distillation to BERT, achieving 97% of BERT’s performance with 60% fewer parameters. TinyBERT [Jiao et al., 2020] further incorporated embedding-layer and attention-based distillation, while MobileBERT [Sun et al., 2020] demonstrated effective distillation for resource-constrained devices.

LLM Compression. Recent work has explored distillation for modern LLMs, though most studies focus on logit-based approaches. Our work systematically evaluates whether white-box signals provide benefits for heterogeneous LLM distillation, where teacher and student architectures differ significantly.

3 Methodology

3.1 Models

Teacher. Llama-2-7B [Touvron et al., 2023]: 7 billion parameters, 4096-dimensional hidden states, 32 attention heads.

Student. TinyLlama-1.1B [Zhang et al., 2024]: 1.1 billion parameters ($6.4\times$ smaller), 2048-dimensional hidden states, 32 attention heads.

3.2 Loss Function

The total training loss combines multiple components:

$$\mathcal{L}_{\text{total}} = \alpha\mathcal{L}_{\text{task}} + \beta\mathcal{L}_{\text{KD}} + \gamma_1\mathcal{L}_{\text{hidden}} + \gamma_2\mathcal{L}_{\text{attn}} \quad (1)$$

where:

- $\mathcal{L}_{\text{task}}$: Cross-entropy on ground-truth labels ($\alpha = 1.0$)
- \mathcal{L}_{KD} : KL-divergence between teacher and student logits ($\beta = 0.5$)
- $\mathcal{L}_{\text{hidden}}$: MSE between projected student and teacher hidden states ($\gamma_1 = 0.1$)
- $\mathcal{L}_{\text{attn}}$: MSE between teacher and student attention maps ($\gamma_2 = 0.1$)

For hidden state alignment, we use a trainable linear projector to map student representations (2048-dim) to match teacher dimensions (4096-dim). Attention maps require no projection since both models have 32 heads.

3.3 Datasets

We evaluate on three diverse tasks (Table 1):

- **SST-2** [Wang et al., 2018]: Binary sentiment classification (5,000 examples)
- **MMLU** [Hendrycks et al., 2021]: Multi-task reasoning across domains (1,000 examples)
- **GSM8K** [Cobbe et al., 2021]: Grade-school math word problems (1,000 examples)

Table 1: Dataset statistics

Dataset	Task Type	Examples	Evaluation
SST-2	Sentiment (NLU)	5,000	Binary accuracy
MMLU	Reasoning	1,000	Multi-choice accuracy
GSM8K	Math	1,000	Exact match

3.4 Training Configuration

All experiments use: learning rate 10^{-4} (AdamW), batch size 8, 3 epochs, max sequence length 512, gradient clipping at 1.0. Each method is run with 7 random seeds (0–6) for statistical validity. Teacher outputs are pre-computed offline for efficiency.

4 Results

4.1 Overall Performance

Table 2 presents the main results. Attention-based distillation achieves the highest accuracy at 95.56%, outperforming the black-box baseline by +1.58 percentage points. The combined approach matches attention-only performance exactly, suggesting hidden states provide no additional signal when attention alignment is present.

4.2 Task-Specific Results

Figure 1 shows performance across individual tasks. Attention distillation provides the largest gain on SST-2 (+2.40%), followed by MMLU (+0.79%) and GSM8K (+0.09%). This suggests attention alignment is particularly beneficial for tasks requiring nuanced language understanding.

Table 2: Overall accuracy (final epoch, $N = 7$ seeds per method). Best results in **bold**.

Method	Mean Accuracy	Std Dev	vs Black-Box
Attention	95.56%	1.43%	+1.58%
Combined	95.56%	1.43%	+1.58%
Hidden State	94.71%	1.05%	+0.73%
Black-Box	93.98%	2.98%	—



Figure 1: Task-specific accuracy breakdown. Attention distillation provides the largest gains on sentiment analysis (SST-2).

4.3 Learning Dynamics

Figure 2 shows validation accuracy over training epochs. White-box methods converge to higher final accuracy and exhibit lower variance (shaded regions). Black-box distillation shows notably higher seed-to-seed variation (2.98% std vs 1.43%), indicating less stable training dynamics.

Validation loss (Figure 3) confirms this pattern: attention-based methods achieve substantially lower final loss (0.028) compared to black-box (0.056).

5 Discussion

5.1 Why Does Attention Work Better?

Attention maps encode *structural* information—which tokens attend to which—that transcends the specific representational capacity of the model. Unlike hidden states, which encode semantic features in model-specific vector spaces, attention patterns represent relational structure that transfers more naturally between architectures.

TinyLlama and Llama-2 share the same number of attention heads (32), enabling direct alignment without projection. In contrast, hidden state alignment requires learning a linear transformation from 2048 to 4096 dimensions, which may introduce representational bottlenecks.

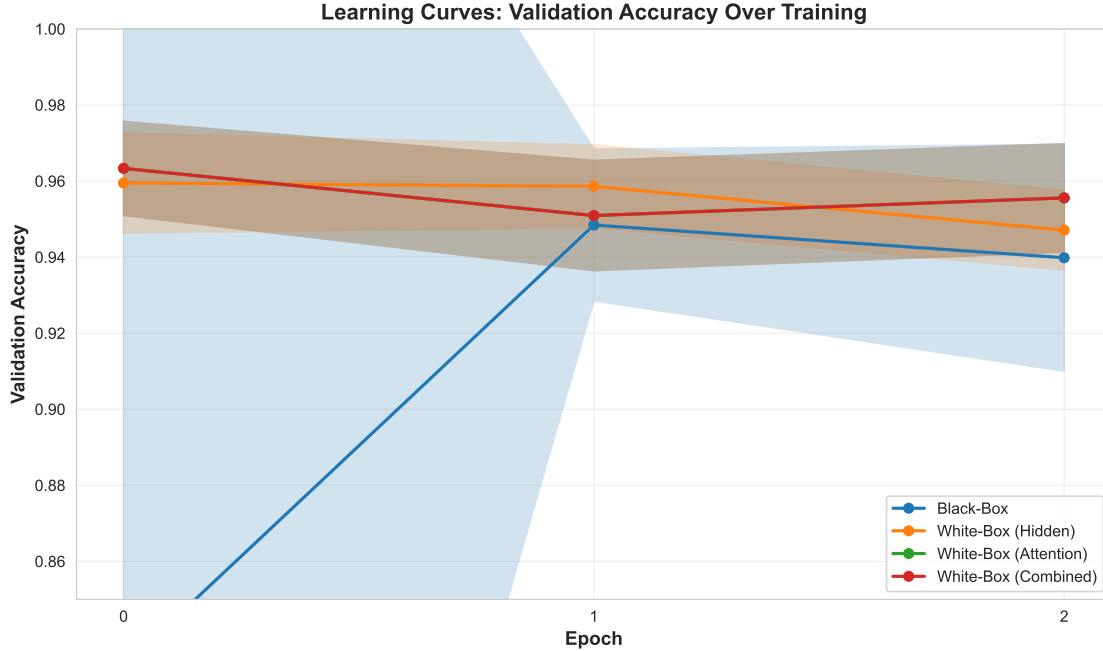


Figure 2: Learning curves showing validation accuracy over training epochs. Shaded regions indicate ± 1 standard deviation across 7 seeds.

5.2 Why Doesn't Combining Signals Help?

The identical performance of Attention and Combined methods suggests that once attention alignment is present, hidden state alignment provides redundant or conflicting supervision. The attention objective may already constrain the model sufficiently that additional hidden state matching adds no benefit.

5.3 Practical Implications

For practitioners compressing LLMs with matched attention head counts:

1. **Use attention distillation:** +1.58% accuracy with minimal overhead
2. **Skip hidden state alignment:** Adds complexity without benefit
3. **Expect more stable training:** Lower variance across random seeds

6 Limitations

- **Model-specific:** Results apply to Llama-2 \rightarrow TinyLlama; other model pairs may differ.
- **Matched heads:** Both models have 32 attention heads; mismatched configurations would require additional projection.
- **Single layer:** We align only final-layer attention; multi-layer alignment may provide additional benefits.
- **Task scope:** Three tasks provide initial evidence; broader evaluation is warranted.

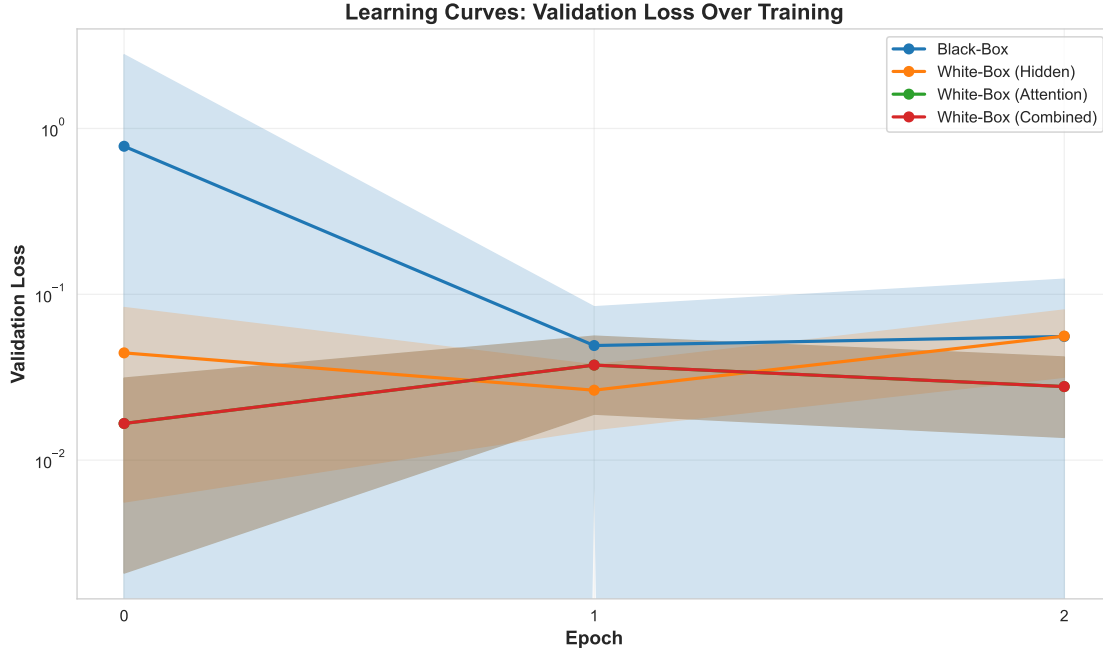


Figure 3: Validation loss curves. Attention-based methods achieve approximately half the final loss of black-box distillation.

- **Hyperparameter sensitivity:** Loss weights (γ_1, γ_2) were fixed; tuning may affect relative rankings.

7 Conclusion

We systematically compared white-box and black-box knowledge distillation for LLM compression. Our experiments demonstrate that attention alignment provides meaningful improvement (+1.58%) over standard logit distillation when compressing Llama-2-7B to TinyLlama-1.1B, while also reducing training variance. Hidden state alignment provides modest benefits alone but no additional gains when combined with attention distillation.

These findings suggest that for heterogeneous LLM distillation with matched attention heads, practitioners should prioritize attention alignment over hidden state matching. Future work should explore multi-layer attention alignment and investigate whether these findings generalize to other model architectures.

Reproducibility

Code is available at: [GitHub URL to be added]. All experiments use publicly available models (Llama-2-7B, TinyLlama-1.1B) and datasets (SST-2, MMLU, GSM8K). Training was conducted on NVIDIA GPUs with PyTorch and Hugging Face Transformers. Full hyperparameters are documented in Section 3.4.

References

- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*, 2021.
- Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.
- Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao Chen, Linlin Li, Fang Wang, and Qun Liu. Tinybert: Distilling bert for natural language understanding. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 4163–4174, 2020.
- Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and Yoshua Bengio. Fitnets: Hints for thin deep nets. In *International Conference on Learning Representations*, 2015.
- Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. Distilbert, a distilled version of bert: Smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*, 2019.
- Zhiqing Sun, Hongkun Yu, Xiaodan Song, Renjie Liu, Yiming Yang, and Denny Zhou. Mobilebert: A compact task-agnostic bert for resource-limited devices. *arXiv preprint arXiv:2004.02984*, 2020.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. Glue: A multi-task benchmark and analysis platform for natural language understanding. *arXiv preprint arXiv:1804.07461*, 2018.
- Sergey Zagoruyko and Nikos Komodakis. Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer. In *International Conference on Learning Representations*, 2017.
- Peiyuan Zhang, Guangtao Zheng, Zeyu Liu, Xinyi Lu, Kaiyan Chang, Siyu Wang, Jiajun Zhang, Jie Tang, and Minlie He. Tinyllama: An open-source small language model. *arXiv preprint arXiv:2401.02385*, 2024.