

LoLCATs: ON LOW-RANK LINEARIZING OF LARGE LANGUAGE MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

Recent works show we can linearize large language models (LLMs)—swapping the quadratic attentions of popular Transformer-based LLMs with subquadratic analogs, such as linear attention—avoiding the expensive pretraining costs. However, linearizing LLMs often significantly degrades model quality, still requires training over billions of tokens, and remains limited to smaller 1.3B to 7B LLMs. We show that the subquadratic analogs used in prior work struggle to approximate the original softmax attention layer. We thus propose Low-rank Linear Conversion via Attention Transfer (LoLCATs), a simple two-step method that improves LLM linearizing quality with orders of magnitudes less memory and compute: (1) the “*attention transfer*” training step uses our new linear attention architecture, designed to improve the approximation fidelity, and minimizes the MSE between the original and new layer’s attention outputs, (2) we adjust for any approximation errors by simply using *low-rank* adaptation (LoRA). LoLCATs significantly improves linearizing quality, training efficiency, and scalability. LoLCATs produces state-of-the-art subquadratic LLMs from Llama 3 8B and Mistral 7B v0.1, leading to 20+ points of improvement on 5-shot MMLU, with only 0.2% of past methods’ model parameters and 0.4% of their training tokens. Finally, we apply LoLCATs to create the first linearized 70B and 405B LLMs (50× larger than prior work). When compared with prior methods under the same compute budgets, LoLCATs significantly improves linearizing quality, closing the gap between linearized and original Llama 3.1 70B and 405B LLMs by 78.7% and 77.4% on 5-shot MMLU.

1 INTRODUCTION

“Linearizing” large language models (LLMs)—or converting existing Transformer-based LLMs into attention-free or subquadratic alternatives—has shown promise for scaling up efficient architectures. While many such architectures offer complexity-level efficiency gains, like *linear-time* and *constant-memory* generation, they are often limited to smaller models pretrained on academic budgets (Gu & Dao, 2023; Peng et al., 2023; Yang et al., 2023; Arora et al., 2024; Beck et al., 2024). In a complementary direction, linearizing aims to start with openly available LLMs—*e.g.*, those with 7B+ parameters trained on trillions of tokens (AI, 2024; Jiang et al., 2023)—and (i) swap their softmax attentions with subquadratic analogs, before (ii) further finetuning to recover quality. This holds exciting promise for quickly scaling up subquadratic capabilities in modern LLMs.

However, to better realize this promise and allow anyone to convert LLMs into subquadratic models, we desire methods that are (1) **quality-preserving**, *e.g.*, recovering the zero-shot abilities of modern LLMs; (2) **parameter and token efficient**, to linearize LLMs on widely accessible compute; and (3) **highly scalable**, to support the various 70B+ LLMs available today (Touvron et al., 2023a;b).

Existing methods present opportunities to improve all three criteria. On quality, despite using motivated subquadratic analogs such as RetNet-inspired linear attentions (Sun et al., 2023; Mercat et al., 2024) or state-space model (SSM)-based Mamba layers (Gu & Dao, 2023; Yang et al., 2024; Wang et al., 2024), prior works significantly reduce performance on popular LM Evaluation Harness tasks (LM Eval) (Gao et al., 2023) (up to 23.4-28.2 pts on 5-shot MMLU (Hendrycks et al., 2020)). On parameter and token efficiency, to adjust for architectural differences, prior methods update *all* model parameters in at least one stage of training (Mercat et al., 2024; Wang et al., 2024; Yang et al.,

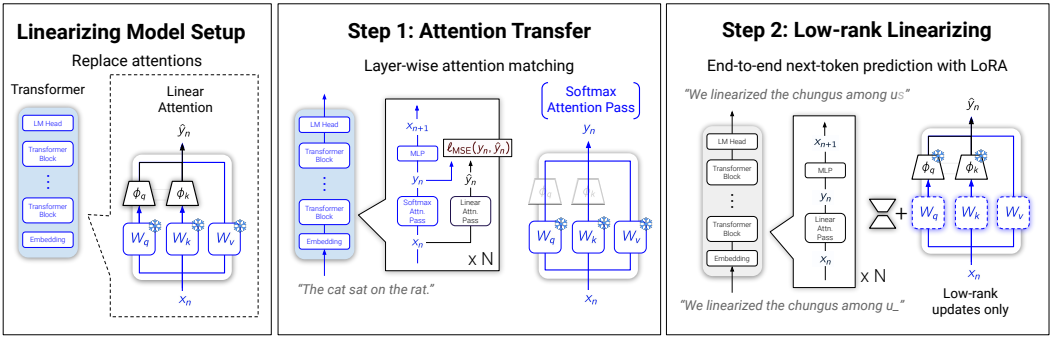


Figure 1: **LOLCATS framework**. We linearize LLMs by (1) training attention analogs to approximate softmax attentions (attention transfer), before swapping attentions and (2) minimally adjusting (with LoRA).

2024), and use 100B tokens to linearize 7B LLMs. On scalability, these training costs make linearizing larger models on academic compute more difficult; existing works only linearize up to 8B LLMs. This makes it unclear how to support linearizing 70B to 405B LLMs (Dubey et al., 2024).

In this work, we thus propose **LOLCATS (Low-rank Linear Conversion with Attention Transfer)**, a simple approach to improve the quality, efficiency, and scalability of linearizing LLMs. As guiding motivation, we ask if we can linearize LLMs by simply reducing architectural differences, *i.e.*,

1. Starting with simple softmax attention analogs such as linear attention (Eq. 2), and *training* their parameterizations explicitly to approximate softmax attention (“**attention transfer**”).
2. Subsequently only training with low-cost finetuning to adjust for any approximation errors, *e.g.*, with low-rank adaptation (LoRA) (Hu et al., 2021) (“**low-rank linearizing**”).

In evaluating this hypothesis, we make several contributions. First, to better understand linearizing feasibility, we empirically study attention transfer and low-rank linearizing with existing linear attentions. While intuitive—by swapping in perfect subquadratic softmax attention approximators, we could get subquadratic LLMs with no additional training—prior works suggest linear attentions struggle to match softmax expressivity (Keles et al., 2023; Qin et al., 2022) or need full-model updates to recover linearizing quality (Kasai et al., 2021; Mercat et al., 2024). In contrast, we find that while *either* attention transfer or LoRA alone is insufficient, we can rapidly recover quality by simply doing *both* (Figure 3, Table 1). At the same time, we do uncover quality issues related to attention-matching architecture and training. With prior linear attentions, the best low-rank linearized LLMs still significantly degrade in quality vs. original Transformers (up to 42.4 pts on 5-shot MMLU). With prior approaches that train all attentions jointly (Zhang et al., 2024), we also find that later layers can result in 200× the MSE of earlier ones (Figure 7). We later find this issue aggravated by larger LLMs; joint training for Llama 3.1 405B’s 126 attention layers fails to linearize LLMs.

Next, to resolve these issues and improve upon our original criteria, we detail LOLCATS’ method components. For **quality**, we generalize prior notions of learnable linear attentions to sliding window + linear attention variants. These remain subquadratic to compute yet consistently yield better attention transfer via lower mean-squared error (MSE) on attention outputs. For **parameter and token efficiency**, we maintain our simple 2-step framework of (1) training subquadratic attentions to match softmax attentions, before (2) adjusting for any errors via only LoRA. For **scalability**, we use finer-grained “block-by-block” training. We split LLMs into blocks of k layers before jointly training attentions only within each block to improve layer-wise attention matching. We pick k to balance the speed of training blocks in parallel with the memory of saving hidden state outputs of prior blocks (as inputs for later ones). We provide a simple cost model to navigate these tradeoffs.

Finally, in experiments, we validate that LOLCATS improves on each of our desired criteria.

- On **quality**, when linearizing popular LLMs such as Mistral-7B and Llama 3 8B, LOLCATS significantly improves past linearizing methods (by 0.2–8.0 points (pts) on zero-shot LM Eval tasks; +17.8 pts on 5-shot MMLU). With Llama 3 8B, LOLCATS for the first time closes the zero-shot LM Eval gap between linearized and Transformer models (73.1 vs 73.7 pts), while supporting 3× higher throughput and 64× larger batch sizes vs. popular FlashAttention-2 (Dao, 2023) implementations (generating 4096 token samples on an 80GB H100). We further validate

Name	Architecture	Quality Preserving	Parameter Efficient	Token Efficient	Validated at Scale
Pretrained	Attention	✓✓	✗✗	✗✗	✓✓
SUPRA	Linear Attention	✗	✗	✓	✓
Mohawk	Mamba (2)	✗	✗	✓	✗
Mamba in Llama	Mamba (2)	✗	✗	✓	✓
LoLCATs	Softmax-Approx. Linear Attention	✓	✓	✓✓	✓✓

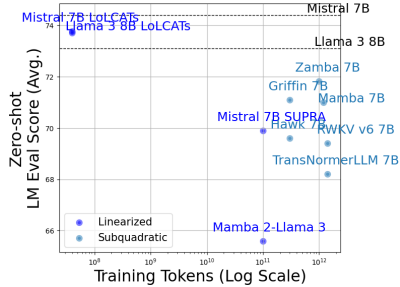


Figure 2: **Linearizing comparison.** LoLCATs significantly improves LLM linearizing quality and training efficiency.

LoLCATs-linearizing as a high-quality training method, outperforming strong 7B subquadratic LLMs (RWKV-v6 (Peng et al., 2024), Mamba (Gu & Dao, 2023), Griffin (De et al., 2024)) and hybrids (StripedHyena (Poli et al., 2023a), Zamba (Glorioso et al., 2024)) trained from scratch by 1.8 to 4.7 pts on average over popular LM Eval tasks.

- On **parameter and token-efficiency**, by only training linear attention feature maps in Stage 1, while only using LoRA on linear attention projections in Stage 2, LoLCATs enables these gains while updating only $<0.2\%$ of past linearizing methods’ model parameters. This also only takes 40M tokens, *i.e.*, 0.003% and 0.04% of prior pretraining and linearizing methods’ token counts.
- On **scalability**, with LoLCATs we scale up linearizing to support the Llama 3.1 70B and 405B parameter LLMs (Dubey et al., 2024). LoLCATs presents the first viable approach to linearizing larger LLMs, creating the first linearized 70B and 405B LLMs with only 9.5 hours on an 8×80 GB H100 node for Llama 3.1 70B and 16 hours across 24 80GB H100s for Llama 3.1 405B. This is in total less than half the compute reported by prior methods to linearize 8B models (5 days on 8×80 GB H100s) (Wang et al., 2024). Furthermore, under these computational constraints, LoLCATs significantly improves quality vs. prior linearizing approaches (Kasai et al., 2021; Mercat et al., 2024). With Llama 3.1 70B and 405B, we close 78.7% and 77.4% of the 5-shot MMLU gap between Transformers and linearized variants respectively.

2 PRELIMINARIES

To motivate LoLCATs, we first go over Transformers, attention, and linear attention. We then briefly discuss related works on linearizing Transformers and Transformer-based LLMs.

Transformers and Attention. Popular LLMs such as Llama 3 8B (AI@Meta, 2024) and Mistral 7B (Jiang et al., 2023) are decoder-only Transformers, with repeated blocks of multi-head *softmax attention* followed by MLPs (Vaswani et al., 2017). For one head, attention models outputs $\mathbf{y} \in \mathbb{R}^{l \times d}$ from inputs $\mathbf{x} \in \mathbb{R}^{l \times d}$ (where l is sequence length, d is head dimension) with query, key, and value weights $\mathbf{W}_q, \mathbf{W}_k, \mathbf{W}_v \in \mathbb{R}^{d \times d}$. In causal language modeling, we compute $\mathbf{q} = \mathbf{x}\mathbf{W}_q$, $\mathbf{k} = \mathbf{x}\mathbf{W}_k$, $\mathbf{v} = \mathbf{x}\mathbf{W}_v$, before getting attention weights \mathbf{a} and outputs \mathbf{y} via

$$a_{n,i} = \frac{\exp(\mathbf{q}_n^\top \mathbf{k}_i / \sqrt{d})}{\sum_{i=1}^n \exp(\mathbf{q}_n^\top \mathbf{k}_i / \sqrt{d})}, \quad \mathbf{y}_n = \sum_{i=1}^n a_{n,i} \mathbf{v}_i \quad (1)$$

Multi-head attention maintains inputs, outputs, and weights for each head, *e.g.*, $\mathbf{x} \in \mathbb{R}^{h \times l \times d}$ or $\mathbf{W}_q \in \mathbb{R}^{h \times d \times d}$ (h being number of heads), and computes Eq. 1 for each head. In both cases, we compute final outputs by concatenating \mathbf{y}_n across heads, before using output weights $\mathbf{W}_o \in \mathbb{R}^{hd \times hd}$ to compute $\mathbf{y}_n \mathbf{W}_o \in \mathbb{R}^{l \times hd}$. While expressive, causal softmax attention requires all $\{\mathbf{k}_i, \mathbf{v}_i\}_{i \leq n}$ to compute \mathbf{y}_n . For long context or large batch settings, this growing *KV cache* can incur prohibitive memory costs even with state-of-the-art implementations such as FlashAttention (Dao, 2023).

Linear Attention. To get around this, Katharopoulos et al. (2020) show a similar attention operation, but with *linear* time and *constant* memory over generation length (linear time and space when processing inputs). To see how, note that softmax attention’s exponential is a kernel function $\mathcal{K}(\mathbf{q}, \mathbf{k})$, which in general can be expressed as the dot product of feature maps $\phi : \mathbb{R}^d \mapsto \mathbb{R}^{d'}$.

Swapping $\exp(\mathbf{q}^\top \mathbf{k} / \sqrt{d})$ with $\phi(\mathbf{q})^\top \phi(\mathbf{k})$ in Eq. 1 gives us *linear attention* weights and outputs:

$$\hat{a}_{n,i} = \frac{\phi(\mathbf{q}_n)^\top \phi(\mathbf{k}_i)}{\sum_{i=1}^n \phi(\mathbf{q}_n)^\top \phi(\mathbf{k}_i)}, \quad \hat{\mathbf{y}}_n = \sum_{i=1}^n \hat{a}_{i,n} \mathbf{v}_i \quad (2)$$

Rearranging terms via matrix product associativity, we get

$$\hat{\mathbf{y}}_n = \sum_{i=1}^n \frac{\phi(\mathbf{q}_n)^\top \phi(\mathbf{k}_i) \mathbf{v}_i}{\sum_{i=1}^n \phi(\mathbf{q}_n)^\top \phi(\mathbf{k}_i)} = \frac{\phi(\mathbf{q}_n)^\top \left(\sum_{i=1}^n \phi(\mathbf{k}_i) \mathbf{v}_i^\top \right)}{\phi(\mathbf{q}_n)^\top \sum_{i=1}^n \phi(\mathbf{k}_i)} \quad (3)$$

This lets us compute both the numerator $\mathbf{s}_n = \sum_{i=1}^n \phi(\mathbf{k}_i) \mathbf{v}_i^\top$ and denominator $\mathbf{z}_n = \sum_{i=1}^n \phi(\mathbf{k}_i)$ as recurrent “KV states”. With $\mathbf{s}_0 = \mathbf{0}, \mathbf{z}_0 = \mathbf{0}$, we recurrently compute linear attention outputs as

$$\hat{\mathbf{y}}_n = \frac{\phi(\mathbf{q}_n)^\top \mathbf{s}_n}{\phi(\mathbf{q}_n)^\top \mathbf{z}_n} \quad \text{for } \mathbf{s}_n = \mathbf{s}_{n-1} + \phi(\mathbf{k}_n) \mathbf{v}_n^\top \quad \text{and} \quad \mathbf{z}_n = \mathbf{z}_{n-1} + \phi(\mathbf{k}_n) \quad (4)$$

Eq. 3 lets us compute attention over an input sequence of length n in $\mathcal{O}(ndd')$ time and space, while Eq. 4 lets us compute n new tokens in $\mathcal{O}(ndd')$ time and $\mathcal{O}(dd')$ memory. Especially during generation, when softmax attention has to compute new tokens sequentially anyway, Eq. 4 enables time and memory savings if $d' < (\text{prompt length} + \text{prior generated tokens})$.

Linearizing Transformers. To combine efficiency with quality, various works propose different ϕ , (e.g., $\phi(x) = 1 + \text{ELU}(x)$ as in Katharopoulos et al. (2020)). However, they typically train linear attention Transformers from scratch. We build upon recent works that *swap* the softmax attentions of *existing* Transformers with linear attention before finetuning the modified models with next-token prediction to recover language modeling quality. These include methods proposed for LLMs (Mercat et al., 2024), and those for smaller Transformers—e.g., 110M BERTs (Devlin et al., 2018)—reasonably adaptable to modern LLMs (Kasai et al., 2021; Mao, 2022; Zhang et al., 2024).

3 METHOD: LINEARIZING LLMs WITH LOLCATs

We now study how to build a high-quality and highly efficient linearizing method. In Section 3.1, we present our motivating framework, which aims to (1) learn good softmax attention approximators with linear attentions and (2) enable low-rank adaptation for recovering linearized quality. In Section 3.2, we find that while this attention transfer works surprisingly well for low-rank linearizing with existing linear attentions, on certain tasks, it still results in sizable quality gaps compared to prior methods. We also find that attention-transfer quality strongly corresponds with the final linearized model’s performance. In Section 3.3, we use our learned findings to overcome prior issues, improving attention transfer to subsequently improve low-rank linearizing quality.

3.1 A FRAMEWORK FOR LOW-COST LINEARIZING

In this section, we present our initial LOLCATs framework for linearizing LLMs in an effective yet efficient manner. Our main hypothesis is that by first learning linear attentions that approximate softmax, we can then swap these attentions in as drop-in subquadratic replacements. We would then only need a minimal amount of subsequent training—e.g., that is supported by low-rank updates—to recover LLM quality in a cost-effective manner effectively. We thus proceed in two steps.

1. **Parameter-Efficient Attention Transfer.** For each softmax attention in an LLM, we aim to learn a closely-approximating linear attention, *i.e.*, one that computes attention outputs $\hat{\mathbf{y}} \approx \mathbf{y}$ for all natural inputs \mathbf{x} . We treat this as a feature map learning problem, learning ϕ to approximate softmax. For each head and layer, let ϕ_q, ϕ_k be query, key feature maps. Per head, we compute:

$$\mathbf{y}_n = \underbrace{\sum_{i=1}^n \frac{\exp(\mathbf{q}_n^\top \mathbf{k}_i / \sqrt{d})}{\sum_{i=1}^n \exp(\mathbf{q}_n^\top \mathbf{k}_i / \sqrt{d})} \mathbf{v}_i}_{\text{Softmax Attention}}, \quad \hat{\mathbf{y}}_n = \underbrace{\sum_{i=1}^n \frac{\phi_q(\mathbf{q}_n)^\top \phi_k(\mathbf{k}_i)}{\sum_{i=1}^n \phi_q(\mathbf{q}_n)^\top \phi_k(\mathbf{k}_i)} \mathbf{v}_i}_{\text{Linear Attention}} \quad (5)$$

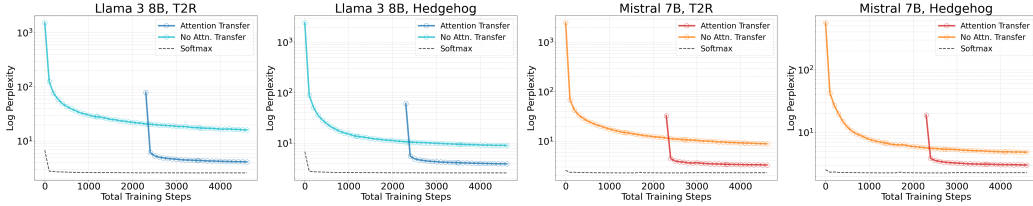


Figure 3: **Attention transfer training efficiency.** Even accounting for initial training steps, low-rank linearizing with attention transfer still consistently achieves lower perplexity faster across feature maps and LLMs.

for all $n \in [l]$ with input $\in \mathbb{R}^{l \times d}$, and train ϕ_q, ϕ_k to minimize sample mean squared error (MSE)

$$\ell_{\text{MSE}} = \frac{1}{MH} \sum_{m=1}^M \sum_{h=1}^H \ell_{\text{MSE}}^{h,m}, \quad \ell_{\text{MSE}}^{h,m} = \frac{1}{d} \sum_{n=1}^d (\mathbf{y}_n - \hat{\mathbf{y}}_n)^2 \quad (6)$$

i.e., jointly for each head h in layer m . Similar to past work (Kasai et al., 2021; Zhang et al., 2024), rather than manually design ϕ , we parameterize each $\phi: \mathbb{R}^d \mapsto \mathbb{R}^{d'}$ as a learnable layer:

$$\phi_q(\mathbf{q}) := f(\mathbf{q}\tilde{\mathbf{W}}_q + \tilde{\mathbf{b}}_q), \quad \phi_k(\mathbf{k}) := f(\mathbf{k}\tilde{\mathbf{W}}_k + \tilde{\mathbf{b}}_k)$$

Here $\tilde{\mathbf{W}} \in \mathbb{R}^{d \times d'}$ and $\tilde{\mathbf{b}} \in \mathbb{R}^{d'}$ are trainable weights and optional biases, $f(\cdot)$ is a nonlinear activation, and d' is an arbitrary feature dimension (set to equal head dimension d in practice).

- Low-rank Adjusting.** After training the linearizing layers, we replace the full-parameter training of prior work with low-rank adaptation (LoRA) (Hu et al., 2021). Like prior work, to adjust for the modifying layers and recover language modeling quality, we now train the modified LLM end-to-end over tokens to minimize a sample next-token prediction loss $\ell_{\text{xent}} = -\sum \log P_{\Theta}(\mathbf{u}_{t+1} | \mathbf{u}_{1:t})$. Here P_{Θ} is the modified LLM, Θ is the set of LLM parameters, and we aim to maximize the probability of true \mathbf{u}_{t+1} given past tokens $\mathbf{u}_{1:t}$ (Fig. 1 right). However, rather than train all LLM parameters, we only train the swapped linear attention $\mathbf{W}_q, \mathbf{W}_k, \mathbf{W}_v, \mathbf{W}_o$ with LoRA. Instead of full-rank updates, $\mathbf{W}' \leftarrow \mathbf{W} + \Delta\mathbf{W}$, LoRA decomposes $\Delta\mathbf{W}$ as the product of two low-rank matrices $\mathbf{B}\mathbf{A}$, $\mathbf{B} \in \mathbb{R}^{d \times r}$, $\mathbf{A} \in \mathbb{R}^{r \times d}$. For parameter efficiency, we aim to pick small $r \ll d$.

Training footprint and efficiency. Both steps remain parameter-efficient. For Step 1, optimizing Eq. 6 is similar to a layer-by-layer cross-architecture distillation. We compute layer-wise (\mathbf{x}, \mathbf{y}) as pretrained attention inputs and outputs, using an LLM forward pass over natural language samples (Fig. 1 middle). However, to keep our training footprint low, we freeze the original pretrained attention layer’s parameters and simply insert new ϕ_q, ϕ_k after $\mathbf{W}_q, \mathbf{W}_k$ in each softmax attention (Fig. 1 left). We compute outputs $\mathbf{y}, \hat{\mathbf{y}}$ with the same attention weights in separate passes (choosing either Eq. 1 or Eq. 3; Fig. 1 middle). For Llama 3 8B or Mistral 7B, training ϕ_q, ϕ_k with $d' = 64$ then only takes 32 layers \times 32 heads \times 2 feature maps \times (128 \times 64) weights \approx 16.8M trainable weights (0.2% of LLM sizes). For Step 2, LoRA with $r = 8$ on all attention projections suffices for state-of-the-art quality. This updates just $<0.09\%$ of 7B parameter counts.

3.2 BASELINE STUDY: ATTENTION TRANSFER AND LOW-RANK LINEARIZING

We now aim to understand if attention transfer and low-rank adjusting are sufficient for linearizing LLMs. It is unclear whether these simple steps can lead to high-quality LLMs, given that prior works default to more involved approaches (Mercat et al., 2024; Yang et al., 2024; Wang et al., 2024). They use linearizing layers featuring GroupNorms (Wu & He, 2018) and decay factors (Sun et al., 2023), or alternate SSM-based architectures (Gu & Dao, 2023; Dao & Gu, 2024). They also all use full-LLM training after swapping in the subquadratic layers. In contrast, we find that simple linear attentions can lead to viable linearizing, with attention transfer + LoRA obtaining competitive quality on 4 / 6 popular LM Eval tasks.

Experimental Setup. We test the LOLCATS framework by linearizing two popular base LLMs, Llama 3 8B (AI, 2024) and Mistral 7B v0.1 (Jiang et al., 2023). For linearizing layers, we study two feature maps used in prior work (Table 2). To support the rotary positional embeddings (RoPE) (Su

Attention	Llama 3 8B				Mistral 7B			
	T2R		Hedgehog		T2R		Hedgehog	
	PPL@0	PPL@2/4	PPL@0	PPL@2/4	PPL@0	PPL@2/4	PPL@0	PPL@2/4
No ✗	1539.39	16.05	2448.01	9.02	2497.13	8.85	561.47	4.87
Yes ✓	79.33	4.11	60.86	3.90	32.78	3.29	18.94	3.04

Table 1: Alpaca validation set perplexity (PPL) of linearized LLMs, comparing attention transfer, no LoRA adjusting (PPL@0) and PPL after training (PPL@2/4; 2 with attention transfer, 4 without, for equal total steps).

Model	Tokens (B)	PiQA	ARC-E	ARC-C	HS	WG	MMLU
Llama 3 8B	-	79.9	80.1	53.3	79.1	73.1	66.6
→ Mamba2	100	76.8	74.1	48.0	70.8	58.6	43.2
→ Hedgehog	0.04	77.4	71.1	40.6	66.5	54.3	24.2
Mistral 7B	-	82.1	80.9	53.8	81.0	74.0	62.4
→ SUPRA	100	80.4	75.9	45.8	77.1	70.3	34.2
→ Hedgehog	0.04	79.3	76.4	45.1	73.1	57.5	28.2

Figure 4: **Linearizing comparison on LM Eval.** Task names in Table 4. Acc. norm: ARC-C, HS. Acc. otherwise. 5-shot MMLU. 0-shot otherwise.

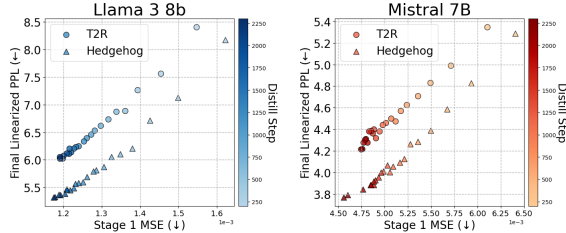


Figure 5: **Attention MSE vs. PPL.** Across feature maps, LLMs; lower MSE coincides with better linearized quality.

et al., 2024) in these LLMs, we apply the feature maps ϕ after RoPE,¹ i.e., computing query features $\phi_q(q) = f(\text{RoPE}(q)\tilde{W}_q + \tilde{b})$. For linearizing data, we wish to see if LOLCATs with a small amount of data can recover general zero-shot and instruction-following LLM abilities. We use the 50K samples of a cleaned Alpaca dataset², due to its ability to improve general instruction-following in 7B LLMs despite its relatively small size (Taori et al., 2023). Following Zhang et al. (2024), we train all feature maps jointly. We include training code and implementation details in App. ??).

To study the effects of attention transfer and low-rank linearizing across LLMs and linear attention architectures, we evaluate their validation set perplexity (Table 1, Fig. 3) and downstream LM Eval zero-shot quality (Table 4). We train both stages with the same data, evaluate with early stopping, and use either two epochs for both attention transfer and LoRA adjusting or four epochs with either alone ($\approx 40M$ total training tokens). For LoRA, we use $r = 8$ as a popular default (Hu et al., 2021), which results in training 0.2% of LLM parameter counts.

Feature Map	$\phi(q)$ (same for k)	Weight Shapes
T2R	$\text{ReLU}(q\tilde{W} + \tilde{b})$	$\tilde{W} : (128, 128), \tilde{b} : (128,)$
Hedgehog	$[\text{SM}_d(q\tilde{W}) \oplus \text{SM}_d(-q\tilde{W})]$	$\tilde{W} : (128, 64)$

Table 2: **Learnable feature maps.** Transformer to RNN (T2R) from Kasai et al. (2021), Hedgehog from Zhang et al. (2024), both \oplus (concat) and SM_d (softmax) apply over feature dimension.

Attention Transfer + LoRA Enables Fast LLM Linearizing. In Table 1 and Fig. 3, we report the validation PPL of linearized LLMs, ablating attention transfer and LoRA adjusting. We find that while attention transfer alone is often insufficient (c.f., PPL@0, Table 1), a single low-rank update rapidly recovers performance by 15–75 PPL (Fig. 3), where training to approximate softmax leads to up to 11.9 lower PPL than no attention transfer. Somewhat surprisingly, this translates to performing competitively with prior linearizing methods that train *all* model parameters (Mercat et al., 2024; Wang et al., 2024) (within 5 accuracy points on 4 / 6 popular LM Eval tasks; Table 4), while *only* training with 0.04% of their token counts and 0.2% of their parameter counts. The results suggest we can linearize 7B LLMs at orders-of-magnitude less training costs than previously shown.

LoL SAD: Limitations of Low-Rank Linearizing. At the same time, we note quality limitations with the present framework. While sometimes close, low-rank linearized LLMs perform worse than full-parameter alternatives and original Transformers on 5 / 6 LM Eval tasks (up to 42.4 points on 5-shot MMLU; Table 4). To understand the issue, we study whether the attention transfer stage can produce high-fidelity linear approximations of softmax attention. We note three observations:

1. Attention transfer quality (via output MSE) strongly ties to low-rank linearized quality (Fig. 5).

¹Unlike prior works that apply ϕ before RoPE (Mercat et al., 2024; Su et al., 2024), our choice preserves the linear attention kernel connection, where we can hope to learn ϕ_q, ϕ_k for $\exp(q^\top k' / \sqrt{d}) \approx \phi_q(q)^\top \phi_k(k')$.

²<https://huggingface.co/datasets/yahma/alpaca-cleaned>

324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377

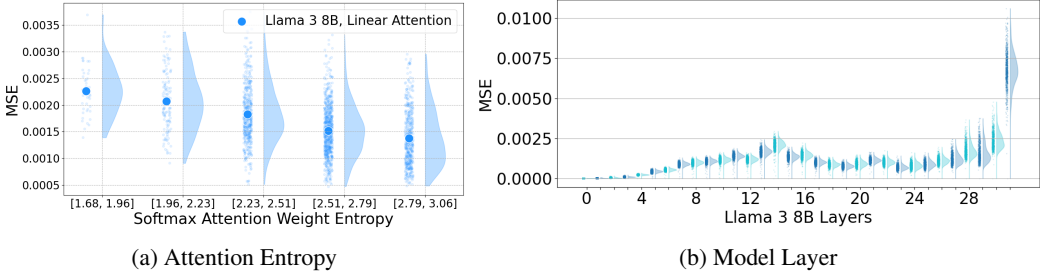


Figure 6: **Sources of Attention Transfer Error** with Llama 3 8B. We find two potential sources of attention transfer difficulty: (a) low softmax attention entropy and (b) attentions in later layers.

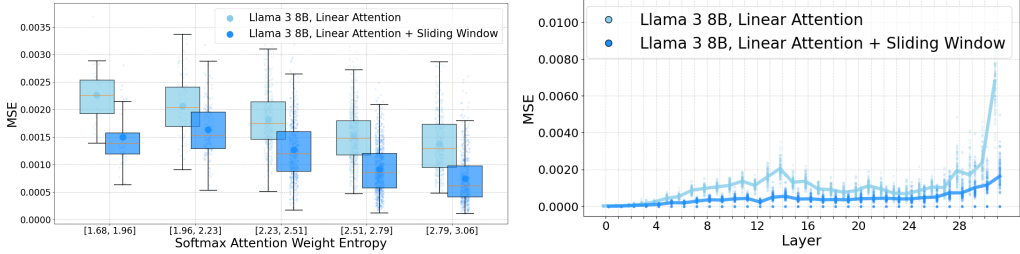


Figure 7: **Improving Attention matching MSE**. Linearizing with linear + sliding window attention better matches LLM softmax attentions (lower MSE) over attention entropy values and LLM layers.

2. Larger attention output MSEs coincide with lower softmax attention weight entropies (Fig. 6a).
3. Larger MSEs also heavily concentrate in attention input samples from later layers (Fig. 6b).

We thus hypothesize we can improve quality by reducing MSE in two ways. First, to approximate samples with lower softmax attention weight entropies—*i.e.*, “spikier” distributions more challenging to capture with linear attentions (Zhang et al., 2024)—we may need better attention-matching architectures. Second, to address the difficulty of learning certain attention layers, we may need more fine-grained layer-wise supervision instead of jointly training all layers.

3.3 LOLCATS: IMPROVED LOW-RANK LINEARIZING

We now introduce two simple improvements in architecture (Section 3.3.1) and linearizing procedure (Section 3.3.2) to improve low-rank linearizing quality.

3.3.1 ARCHITECTURE: GENERALIZING LEARNABLE LINEAR ATTENTIONS

As described, we can apply our framework with any linear attentions with learnable ϕ (e.g., T2R and Hedgehog, Figure 3). However, to improve attention-matching quality, we introduce a hybrid ϕ parameterization combining linear attention and *sliding window* attention. Motivated by prior works that show quality improvements when combining attention layers with linear attentions (Arora et al., 2024; Munkhdalai et al., 2024), we combine *short sliding windows* of softmax attention (Beltagy et al., 2020; Zhu et al., 2021) (size 64 in experiments) followed by linear attention in a single layer. This allows attending to all prior tokens for each layer while keeping the entire LLM subquadratic. For window size w and token indices $[1, \dots, n - w, \dots, n]$, we apply the softmax attention over the w most recent tokens, and compute attention outputs \hat{y}_n as

$$\hat{y}_n = \frac{\sum_{i=n-w+1}^n \gamma \exp(\mathbf{q}_n^\top \mathbf{k}_i / \sqrt{d} - c_n) \mathbf{v}_i + \phi_q(\mathbf{q}_n)^\top (\sum_{j=1}^{n-w} \phi_k(\mathbf{k}_j) \mathbf{v}_j^\top)}{\sum_{i=n-w+1}^n \gamma \exp(\mathbf{q}_n^\top \mathbf{k}_i / \sqrt{d} - c_n) + \phi_q(\mathbf{q}_n)^\top (\sum_{j=1}^{n-w} \phi_k(\mathbf{k}_j)^\top)} \quad (7)$$

γ is a learnable mixing term, and c_n is a stabilizing constant as in log-sum-exp calculations ($c_n = \max_i \{\mathbf{q}_n^\top \mathbf{k}_i / \sqrt{d} : i \in [n - w + 1, \dots, n]\}$). Like before, we can pick any learnable ϕ .

Subquadratic efficiency. The hybrid layer retains linear time and constant memory generation. For n -token prompts, we initially require $\mathcal{O}(w^2d)$ and $\mathcal{O}((n - w)dd')$ time and space for window

Model	Training Tokens (B)	PiQA	ARC-e	ARC-c (norm)	HellaSwag (norm)	Wino-grande	MMLU (5-shot)	Avg.	Avg. (no MMLU)
Mistral 7B	-	82.1	80.9	53.8	81.0	74.0	62.4	72.4	74.4
Mistral 7B SUPRA	100	80.4	75.9	45.8	77.1	70.3	34.2	64.0	69.9
Mistral 7B Hedgehog	0.04	79.3	76.4	45.1	73.1	57.5	28.2	59.9	66.3
Mistral 7B LoLCATs (Ours)	0.04	81.1	81.1	52.9	80.5	73.2	52.2	70.2	73.8
Llama 3 8B	-	79.9	80.1	53.3	79.1	73.1	66.6	72.0	73.1
Mamba2-Llama 3	100	76.8	74.1	48.0	70.8	58.6	43.2	61.9	65.6
Mamba2-Llama 3, 50% Attn.	100	81.5	78.8	58.2	78.4	69.0	56.7	70.4	73.2
Llama 3 8B Hedgehog	0.04	77.4	71.1	40.6	66.5	54.3	24.2	55.7	62.0
Llama 3 8B LoLCATs (Ours)	0.04	80.6	81.8	53.5	79.1	73.4	54.9	70.6	73.7

Table 3: **LoLCATs comparison among linearized 7B+ LLMs.** Among linearized 7B+ LLMs, LoLCATs-linearized Mistral 7B and Llama 3 8B consistently achieve best or 2nd-best performance on LM Eval tasks (only getting 2nd best to Mamba-Transformer hybrids). LoLCATs closes the Transformer quality gap by 73.8% (Mistral 7B) and 86.1% (Llama 3 8B) (average over all tasks; numbers except Hedgehog cited from original works), despite only using 40M tokens to linearize (a $2,500\times$ improvement in tokens-to-model efficiency).

and linear attention respectively, attending over a w -sized KV-cache and computing KV and K-states (Eq. 4). For generation, we only need $\mathcal{O}(w^2d + dd')$ time and space for every token. We evict the KV-cache’s first \mathbf{k} , \mathbf{v} , compute $\phi_k(\mathbf{k})$, and add $\phi_k(\mathbf{k})\mathbf{v}^\top$ and $\phi_k(\mathbf{k})$ to KV and K-states respectively.

3.3.2 TRAINING: LAYER (OR BLOCK)-WISE ATTENTION TRANSFER

We describe the training approach and provide a simplified model to show its cost-quality tradeoffs. To improve layer-wise quality, instead of computing the training loss (Eq. 6) over all $m \in [M]$ for a model with M layers, we compute the loss over k -layer blocks, and train each block independently:

$$\ell_{\text{MSE}}^{\text{block}} = \frac{1}{kH} \sum_{m=i}^{i+k} \sum_{h=1}^H \ell_{\text{MSE}}^{h,m} \quad (\text{for blocks starting at layers } i = 0, k, 2k, \dots) \quad (8)$$

We can choose k to balance cost and quality. There are several approaches for training block-wise, including via joint training with separate optimizer groups per block or by sequentially training separate blocks. The primary costs-tradeoffs between these two approaches are:

- **Compute.** Increasing k increases the compute required per block. While the joint training of Llama 3.1 405B in 16-bit precision uses multiple NVIDIA H100 $8\times 80\text{GB}$ nodes, separate blocks of $k = 9$ or fewer layers fits on a *single* H100 80GB GPU, at sequence length 1024.
- **Memory and training time.** The total amount of memory required is $2 \times T \times d \times \frac{L}{k}$ for total training tokens T , model dimension d , number of layers L and 2-byte (16-bit) precision. At the Llama 3.1 405B scale, saving outputs per-layer ($k = 1$) for 40M tokens would require 165TB of disk space. Sequentially saving the outputs and training the blocks increases total training time.

4 EXPERIMENTS

Through experiments, we study: (1) if LoLCATs linearizes LLMs with higher quality than existing subquadratic alternatives and linearizations, and higher generation efficiency than original Transformers (Section 4.1); (2) how ablations on attention transfer loss, subquadratic architecture, and parameter and token counts impact LLM downstream quality (Section 4.2); (3) how LoLCATs’ quality and efficiency holds up to 70B and 405B LLMs, where we linearize and compare model quality across the complete Llama 3.1 family (Section 4.3).

4.1 MAIN RESULTS: LoLCATs EFFICIENTLY RECOVERS QUALITY IN LINEARIZED LLMs

In our main evaluation, we linearize the popular base Llama 3 8B (AI, 2024) and Mistral 7B (Jiang et al., 2023) LLMs. We first test if LoLCATs can efficiently create high-quality subquadratic LLMs from strong base Transformers, comparing to existing linearized LLMs from prior methods. We also test if LoLCATs can create subquadratic LLMs that outperform modern Transformer alternatives pretrained from scratch. For space, we defer linearizing training details to App. A.

Model	Training Tokens (B)	PiQA	ARC-e	ARC-c (norm)	HellaSwag (norm)	Wino-grande	MMLU (5-shot)	Avg.	Avg. (no MMLU)
Transformer									
Gemma 7B	6000	81.9	81.1	53.2	80.7	73.7	62.9	72.3	74.1
Mistral 7B	8000 ³	82.1	80.9	53.8	81.0	74.0	62.4	72.4	74.4
Llama 3 8B	15000	79.9	80.1	53.3	79.1	73.1	66.6	72.0	73.1
Subquadratic									
Mamba 7B	1200	81.0	77.5	46.7	77.9	71.8	33.3	64.7	71.0
RWKV-6 World v2.1 7B	1420	78.7	76.8	46.3	75.1	70.0	-	69.4	69.4
TransNormerLLM 7B	1400	80.1	75.4	44.4	75.2	66.1	43.1	64.1	68.2
Hawk 7B	300	80.0	74.4	45.9	77.6	69.9	35.0	63.8	69.6
Griffin 7B	300	81.0	75.4	47.9	78.6	72.6	39.3	65.8	71.1
Hybrid									
StripedHyena-Nous-7B	-	78.8	77.2	40.0	76.4	66.4	26.0	60.8	67.8
Zamba 7B	1000	81.4	74.5	46.6	80.2	76.4	57.7	69.5	71.8
Linearized									
Mistral 7B LoLCATs (Ours)	0.04	81.1	81.1	52.9	80.5	73.2	52.2	70.2	73.8
Llama 3 8B LoLCATs (Ours)	0.04	80.6	81.8	53.5	79.1	73.4	54.9	70.6	73.7

Table 4: **LoLCATs comparison to pretrained subquadratic LLMs.** LoLCATs-linearized Mistral 7B and Llama 3 8B further outperform pretrained subquadratic Transformer alternatives by 0.1 to 9.4 points (Avg.), while only training 0.2% of the model parameters on 0.013% to 0.003% of their pretraining token counts.

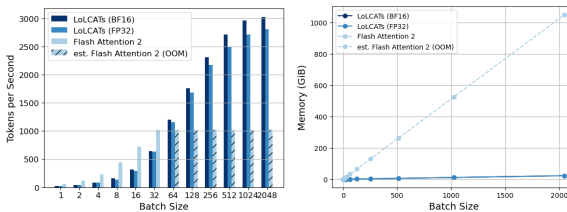
In Table 4, we report results on six popular LM Evaluation Harness (LM Eval) tasks (Gao et al., 2023). Compared to recent linearizing methods, LoLCATs significantly improves quality and training efficiency across tasks and LLMs. On quality, LoLCATs closes 73.8% and 81.1% of the Transformer-linearizing gap for Mistral 7B and Llama 3 8B respectively, notably improving 5-shot MMLU by 63.8% and 50% against next best fully subquadratic models. On efficiency, we achieve these results despite only training <0.2% of model parameters via LoRA versus prior full-parameter training and use 40M tokens versus the prior 100B (0.04% of the latter, a 2500× improvement in “tokens-to-model” efficiency). Among all 7B LLMs, LoLCATs-linearized LLMs further outperform strong subquadratic Transformer alternatives, including RNNs or linear attentions (RWKV-v6 (Peng et al., 2024), Hawk (De et al., 2024), Griffin (De et al., 2024), TransNormer (Qin et al., 2023)), state-space models (SSMs) (Mamba (Gu & Dao, 2023)), and hybrid architectures with some full attention (StripedHyena (Poli et al., 2023b), Zamba (Glorioso et al., 2024)).

4.2 LoLCATs COMPONENT PROPERTIES AND ABLATIONS

We next validate that LoLCATs linearizing enable subquadratic efficiency, and study how each of LoLCATs’ components contribute to these linearizing quality gains.

Subquadratic Generation Throughput and Memory.

We measure the generation throughput and memory of LoLCATs LLMs, validating that linearizing LLMs can significantly improve their generation efficiency. We use the popular Llama 3 8B HuggingFace checkpoint⁴, and compare LoLCATs implemented in HuggingFace Transformers with the supported FlashAttention-2 (FA2) implementation (Dao, 2023). We benchmark on a single 80GB H100 and benchmark two LoLCATs implementations with the Hedgehog feature map and (linear + sliding window) attention in FP32 and BF16. In Fig. 8a and Fig. 8b, we report the effect of scaling batch size on throughput and memory. We measure throughput as (newly generated tokens × batch size / total time), using 128 token prompts and 4096 token generations. As batch size scales, LoLCATs-linearized LLMs achieve significantly higher throughput than FA2. We note this is primarily due to lower memory, where FA2 runs out of memory at batch size 64. Meanwhile, LoLCATs supports up to 3000 tokens / second with batch size 2048 (Fig. 8a), only maintaining a fixed “KV state” as opposed to the growing KV cache in all attention implementations (Fig. 8b).



(a) Batch size vs. Throughput (b) Batch size vs. Memory

⁴<https://huggingface.co/meta-llama/Meta-Llama-3-8B>

486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

Feature Map	LM Eval Metric	Swap & Finetune	+Attention Transfer	+Sliding Window, +Attention Transfer	+ Sliding Window, No Attention Transfer
Hedgehog	Avg. Zero-Shot	44.20	55.32	70.60	68.78
	MMLU (5-shot)	23.80	23.80	52.50	45.80
T2R	Avg. Zero-Shot	38.84	54.83	68.28	39.52
	MMLU (5-shot)	23.20	23.10	40.70	23.80

Table 5: **LOLCATS component ablations**, linearizing Llama 3 8B over 1024-token sequences. Default configuration highlighted. Across feature maps, LOLCATS’ attention transfer and sliding window increasingly improve linearized LLM quality.

	PiQA acc	ARC Easy acc	ARC Challenge (acc norm)	HellaSwag (acc norm)	WinoGrande acc	MMLU (5-shot) acc
Llama 3.1 70B	83.10	87.30	60.60	85.00	79.60	78.80
Linearized, no attn. transfer	75.70	70.10	39.10	77.40	58.60	26.60
LoLCATs (Ours)	82.10	85.00	60.50	84.60	73.70	67.70
Llama 3.1 405B	85.58	87.58	66.21	87.13	79.40	82.98
Linearized, no attn. transfer	84.44	86.62	64.33	86.19	79.87	33.86
LoLCATs (Ours)	85.58	88.80	67.75	87.41	80.35	71.90

Table 6: **Linearizing Llama 3.1 70B and 405B**. Among the first linearized 70B and 405B LLMs (via low-rank linearizing), LOLCATs significantly improves zero- and few-shot quality.

Ablations. We study how adding the attention transfer and linear + sliding window attention in LOLCATs contribute to downstream linearized LLM performance, linearizing Llama 3 8B over 1024-token long samples (Table 5). We start with standard linear attentions (Hedgehog, Zhang et al. (2024); T2R, Kasai et al. (2021)), using the prior linearizing procedure of just swapping attentions and finetuning the model to predict next tokens (Mercat et al., 2024). We then add either (i) attention transfer, (ii) linear + sliding window attentions, or (iii) both, and report the average LM Eval score over the six popular zero-shot tasks in Table 4 and 5-shot MMLU accuracy. Across feature maps, we validate the LOLCATs combination leads to best performance.

4.3 SCALING UP LINEARIZING TO 70B AND 405B LLMs

We finally use LOLCATs to scale up linearizing to Llama 3.1 70B and 405B models. In Table 6, we find that LOLCATs provides the first practical solution for linearizing larger LLMs, achieving significant quality improvements over prior linearizing approaches (Mercat et al., 2024). For Llama 3.1 70B, we achieve a 41.1 point improvement in 5-shot MMLU accuracy. For Llama 3.1 405B, LOLCATs similarly achieves a 38.0 point improvement over prior methods. These results highlight LOLCATs ability to linearize large-scale models with greater efficiency and improved performance, showing for the first time that we can scale up linearizing to 70B+ LLMs.

5 CONCLUSION

We propose LOLCATs, an efficient LLM linearizing method that (1) trains attention analogs—such as linear attentions and linear attention + sliding window hybrids—to approximate an LLM’s self-attentions, before (2) swapping the attentions and only finetuning the replacing attentions with LoRA. We exploit the fidelity between these attention analogs and softmax attention, where we reduce the problem of linearizing LLMs to learning to approximate softmax attention in a subquadratic analog. Furthermore, we demonstrate that via an MSE-based attention output-matching loss, we are able to train such attention analogs to approximate the “ground-truth” softmax attentions in practice. On popular zero-shot LM Evaluation harness benchmarks and 5-shot MMLU, we find this enables producing high-quality, high-inference efficiency LLMs that outperform prior Transformer alternatives while only updating 0.2% of model parameters and requiring 0.003% of the training tokens to achieve similar quality with LLM pretraining. Our findings significantly improve linearizing quality and accessibility, allowing us to create the first linearized 70B and 405B LLMs.

540 ETHICS STATEMENT

541
542 Our work deals with improving the efficiency of open-weight models. While promising for benefi-
543 cial applications, increasing their accessibility also raises concerns about potential misuse. Bad
544 actors could leverage our technique to develop LLMs capable of generating harmful content, spread-
545 ing misinformation, or enabling other malicious activities. We focus primarily on base models, but
546 acknowledge that linearizing could also be used on instruction-tuned LLMs; research on whether
547 linearizing preserves guardrails is still an open question. We acknowledge the risks and believe in
548 the responsible development and deployment of efficient and widely accessible models.

549 REPRODUCIBILITY

550
551 We include experimental details in Appendix A, including sample code for the linearizing architec-
552 ture and training (Appendix C).

554 REFERENCES

- 555
556 Mistral AI. Mixtral of experts. *Mistral AI — Frontier AI in your hands*, May 2024. URL <https://mistral.ai/news/mixtral-of-experts/>.
- 557
558 AI@Meta. Llama 3 model card. 2024. URL [https://github.com/meta-llama/](https://github.com/meta-llama/llama3/blob/main/MODEL_CARD.md)
559 [llama3/blob/main/MODEL_CARD.md](https://github.com/meta-llama/llama3/blob/main/MODEL_CARD.md).
- 560
561 Simran Arora, Sabri Eyuboglu, Michael Zhang, Aman Timalsina, Silas Alberti, Dylan Zinsley,
562 James Zou, Atri Rudra, and Christopher Ré. Simple linear attention language models balance
563 the recall-throughput tradeoff. *arXiv preprint arXiv:2402.18668*, 2024.
- 564
565 Maximilian Beck, Korbinian Pöppel, Markus Spanring, Andreas Auer, Oleksandra Prudnikova,
566 Michael Kopp, Günter Klambauer, Johannes Brandstetter, and Sepp Hochreiter. xlstm: Extended
567 long short-term memory. *arXiv preprint arXiv:2405.04517*, 2024.
- 568
569 Iz Beltagy, Matthew E Peters, and Arman Cohan. Longformer: The long-document transformer.
570 *arXiv preprint arXiv:2004.05150*, 2020.
- 571
572 Together Computer. Redpajama: An open source recipe to reproduce llama training dataset, 2023.
573 URL <https://github.com/togethercomputer/RedPajama-Data>.
- 574
575 Tri Dao. Flashattention-2: Faster attention with better parallelism and work partitioning. *arXiv*
576 *preprint arXiv:2307.08691*, 2023.
- 577
578 Tri Dao and Albert Gu. Transformers are ssms: Generalized models and efficient algorithms through
579 structured state space duality. *arXiv preprint arXiv:2405.21060*, 2024.
- 580
581 Soham De, Samuel L Smith, Anushan Fernando, Aleksandar Botev, George Cristian-Muraru, Al-
582 bert Gu, Ruba Haroun, Leonard Berrada, Yutian Chen, Srivatsan Srinivasan, et al. Griffin: Mix-
583 ing gated linear recurrences with local attention for efficient language models. *arXiv preprint*
584 *arXiv:2402.19427*, 2024.
- 585
586 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep
587 bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- 588
589 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha
590 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models.
591 *arXiv preprint arXiv:2407.21783*, 2024.
- 592
593 Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster,
594 Laurence Golding, Jeffrey Hsu, Alain Le Noac’h, Haonan Li, Kyle McDonell, Niklas Muen-
595 nighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lin-
596 tang Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. A framework
597 for few-shot language model evaluation, 12 2023. URL [https://zenodo.org/records/](https://zenodo.org/records/10256836)
598 [10256836](https://zenodo.org/records/10256836).

- 594 Bogdan Gliwa, Iwona Mochol, Maciej Biesek, and Aleksander Wawer. SAMSum corpus: A human-
595 annotated dialogue dataset for abstractive summarization. In Lu Wang, Jackie Chi Kit Cheung,
596 Giuseppe Carenini, and Fei Liu (eds.), *Proceedings of the 2nd Workshop on New Frontiers in Sum-*
597 *marization*, pp. 70–79, Hong Kong, China, November 2019. Association for Computational Lin-
598 guistics. doi: 10.18653/v1/D19-5409. URL <https://aclanthology.org/D19-5409>.
- 599 Paolo Glorioso, Quentin Anthony, Yury Tokpanov, James Whittington, Jonathan Pilault, Adam
600 Ibrahim, and Beren Millidge. Zamba: A compact 7b ssm hybrid model. *arXiv preprint*
601 *arXiv:2405.16712*, 2024.
- 602 Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv*
603 *preprint arXiv:2312.00752*, 2023.
- 604 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and
605 Jacob Steinhardt. Measuring massive multitask language understanding. *arXiv preprint*
606 *arXiv:2009.03300*, 2020.
- 607 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,
608 and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint*
609 *arXiv:2106.09685*, 2021.
- 610 Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot,
611 Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al.
612 Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.
- 613 Jungo Kasai, Hao Peng, Yizhe Zhang, Dani Yogatama, Gabriel Ilharco, Nikolaos Pappas, Yi Mao,
614 Weizhu Chen, and Noah A. Smith. Finetuning pretrained transformers into RNNs. In *Proceed-*
615 *ings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 10630–
616 10643, Online and Punta Cana, Dominican Republic, November 2021. Association for Computa-
617 tional Linguistics. doi: 10.18653/v1/2021.emnlp-main.830. URL <https://aclanthology.org/2021.emnlp-main.830>.
- 618 Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers are
619 rnns: Fast autoregressive transformers with linear attention. In *International conference on ma-*
620 *chine learning*, pp. 5156–5165. PMLR, 2020.
- 621 Feyza Duman Keles, Pruthvi Mahesakya Wijewardena, and Chinmay Hegde. On the computational
622 complexity of self-attention. In *International Conference on Algorithmic Learning Theory*, pp.
623 597–619. PMLR, 2023.
- 624 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint*
625 *arXiv:1711.05101*, 2017.
- 626 Huanru Henry Mao. Fine-tuning pre-trained transformers into decaying fast weights. In *Proceed-*
627 *ings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 10236–
628 10242, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Lin-
629 guistics. doi: 10.18653/v1/2022.emnlp-main.697. URL <https://aclanthology.org/2022.emnlp-main.697>.
- 630 Jean Mercat, Igor Vasiljevic, Sedrick Keh, Kushal Arora, Achal Dave, Adrien Gaidon, and Thomas
631 Kollar. Linearizing large language models. *arXiv preprint arXiv:2405.06640*, 2024.
- 632 Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mix-
633 ture models. In *International Conference on Learning Representations*, 2017. URL <https://openreview.net/forum?id=Byj72udxe>.
- 634 Tsendsuren Munkhdalai, Manaal Faruqui, and Siddharth Gopal. Leave no context behind: Efficient
635 infinite context transformers with infini-attention. April 2024.
- 636 Guilherme Penedo, Quentin Malartic, Daniel Hesslow, Ruxandra Cojocaru, Alessandro Cappelli,
637 Hamza Alobeidli, Baptiste Pannier, Ebtesam Almazrouei, and Julien Launay. The refinedweb
638 dataset for falcon llm: outperforming curated corpora with web data, and web data only. *arXiv*
639 *preprint arXiv:2306.01116*, 2023.

- 648 Bo Peng, Eric Alcaide, Quentin Anthony, Alon Albalak, Samuel Arcadinho, Huanqi Cao, Xin
649 Cheng, Michael Chung, Matteo Grella, Kranthi Kiran GV, et al. Rwkv: Reinventing rns for
650 the transformer era. *arXiv preprint arXiv:2305.13048*, 2023.
- 651 Bo Peng, Daniel Goldstein, Quentin Anthony, Alon Albalak, Eric Alcaide, Stella Biderman, Eugene
652 Cheah, Teddy Ferdinan, Haowen Hou, Przemysław Kazienko, et al. Eagle and finch: Rwkv with
653 matrix-valued states and dynamic recurrence. *arXiv preprint arXiv:2404.05892*, 2024.
- 654 Michael Poli, Stefano Massaroli, Eric Nguyen, Daniel Y Fu, Tri Dao, Stephen Baccus, Yoshua
655 Bengio, Stefano Ermon, and Christopher Ré. Hyena hierarchy: Towards larger convolutional
656 language models. *arXiv preprint arXiv:2302.10866*, 2023a.
- 657 Michael Poli, Jue Wang, Stefano Massaroli, Jeffrey Quesnelle, Ryan Carlow, Eric Nguyen,
658 and Armin Thomas. StripedHyena: Moving Beyond Transformers with Hybrid Signal
659 Processing Models, 12 2023b. URL [https://github.com/togethercomputer/
660 stripedhyena](https://github.com/togethercomputer/stripedhyena).
- 661 Zhen Qin, Xiaodong Han, Weixuan Sun, Dongxu Li, Lingpeng Kong, Nick Barnes, and Yiran
662 Zhong. The devil in linear transformer. In *Proceedings of the 2022 Conference on Empirical
663 Methods in Natural Language Processing*, pp. 7025–7041, Abu Dhabi, United Arab Emirates,
664 December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.emnlp-main.
665 473. URL <https://aclanthology.org/2022.emnlp-main.473>.
- 666 Zhen Qin, Dong Li, Weigao Sun, Weixuan Sun, Xuyang Shen, Xiaodong Han, Yunshen Wei, Bao-
667 hong Lv, Fei Yuan, Xiao Luo, et al. Scaling transnormer to 175 billion parameters. *arXiv preprint
668 arXiv:2307.14995*, 2023.
- 669 Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever.
670 Language models are unsupervised multitask learners. 2019. URL [https://api.
671 semanticscholar.org/CorpusID:160025533](https://api.semanticscholar.org/CorpusID:160025533).
- 672 Uri Shaham, Elad Segal, Maor Ivgi, Avia Efrat, Ori Yoran, Adi Haviv, Ankit Gupta, Wenhan Xiong,
673 Mor Geva, Jonathan Berant, and Omer Levy. SCROLLS: Standardized CompaRison over long
674 language sequences. In *Proceedings of the 2022 Conference on Empirical Methods in Natural
675 Language Processing*, pp. 12007–12021, Abu Dhabi, United Arab Emirates, December 2022.
676 Association for Computational Linguistics. URL [https://aclanthology.org/2022.
677 emnlp-main.823](https://aclanthology.org/2022.emnlp-main.823).
- 678 Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: En-
679 hanced transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2024.
- 680 Yutao Sun, Li Dong, Shaohan Huang, Shuming Ma, Yuqing Xia, Jilong Xue, Jianyong Wang, and
681 Furu Wei. Retentive network: A successor to transformer for large language models. *arXiv
682 preprint arXiv:2307.08621*, 2023.
- 683 Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy
684 Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model.
685 https://github.com/tatsu-lab/stanford_alpaca, 2023.
- 686 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
687 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and
688 efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023a.
- 689 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-
690 lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open founda-
691 tion and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023b.
- 692 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
693 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural informa-
694 tion processing systems*, 30, 2017.
- 695 Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman.
696 Glue: A multi-task benchmark and analysis platform for natural language understanding. *arXiv
697 preprint arXiv:1804.07461*, 2018.

702 Junxiong Wang, Daniele Paliotta, Avner May, Alexander M Rush, and Tri Dao. The mamba in the
703 llama: Distilling and accelerating hybrid models. *arXiv preprint arXiv:2408.15237*, 2024.

704
705 Yuxin Wu and Kaiming He. Group normalization. In *Proceedings of the European conference on*
706 *computer vision (ECCV)*, pp. 3–19, 2018.

707 Songlin Yang, Bailin Wang, Yikang Shen, Rameswar Panda, and Yoon Kim. Gated linear attention
708 transformers with hardware-efficient training. *arXiv preprint arXiv:2312.06635*, 2023.

709
710 Songlin Yang, Bailin Wang, Yu Zhang, Yikang Shen, and Yoon Kim. Parallelizing linear transform-
711 ers with the delta rule over sequence length. *arXiv preprint arXiv:2406.06484*, 2024.

712 Michael Zhang, Kush Bhatia, Hermann Kumbong, and Christopher Re. The hedgehog & the por-
713 cupine: Expressive linear attentions with softmax mimicry. In *The Twelfth International Confer-*
714 *ence on Learning Representations*, 2024. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=4g0212N2Nx)
715 [4g0212N2Nx](https://openreview.net/forum?id=4g0212N2Nx).

716
717 Chen Zhu, Wei Ping, Chaowei Xiao, Mohammad Shoeybi, Tom Goldstein, Anima Anandkumar,
718 and Bryan Catanzaro. Long-short transformer: Efficient transformers for language and vision.
719 *Advances in neural information processing systems*, 34:17723–17736, 2021.

721 A EXPERIMENTAL DETAILS

722
723 For linearizing layers, we replace softmax attentions with hybrid linear + terraced window analogs
724 (Section 3.3.1), using Hedgehog’s feature map for its prior quality [Zhang et al. \(2024\)](#). For lineariz-
725 ing data, we report results using the Alpaca-linearized models. We also tried a more typical pre-
726 training corpus (1B tokens⁵ of RedPajama [Computer \(2023\)](#)), but found comparable performance
727 when controlling for number of token updates. To linearize, we simply train all feature maps in
728 parallel for two epochs with learning rate 1e-2, before applying LoRA on the attention projection
729 layers for two epochs with learning rate 1e-4. By default, we use LoRA rank $r = 8$, amounting to
730 training <0.09% of all model parameters. For both stages, we train with early stopping, AdamW
731 optimizer [Loshchilov & Hutter \(2017\)](#), and packing into 1024-token sequences with batch size 8.

733 B RELATED WORK

734
735 In this work, we build upon both approaches explicitly proposed to linearize LLMs [Mercat et al.](#)
736 [\(2024\)](#), as well as prior methods focusing on smaller Transformers reasonably adaptable to modern
737 LLMs [Kasai et al. \(2021\)](#); [Mao \(2022\)](#); [Zhang et al. \(2024\)](#). We highlight two approaches most
738 related to LOLCATs and their extant limitations next.

739
740 **Scalable UPtraining for Recurrent Attention (SUPRA).** [Mercat et al. \(2024\)](#) linearize LLMs
741 by swapping softmax attentions with linear attentions similar to Retentive Network (RetNet) lay-
742 ers [Sun et al. \(2023\)](#), before jointly training all model parameters on the RefinedWeb pretraining
743 dataset [Penedo et al. \(2023\)](#). In particular, they suggest that linearizing LLMs with the vanilla linear
744 attention in Eq. 2 is unstable, and swap attentions with

$$745 \hat{y}_n = \text{GroupNorm} \left(\sum_{i=1}^n \gamma^{n-i} \phi(\mathbf{q}_n)^\top \phi(\mathbf{k}_i) \mathbf{v}_i \right) \quad (9)$$

746
747
748 GroupNorm [Wu & He \(2018\)](#) is used as the normalization in place of the $\sum_{i=1}^n \phi(\mathbf{q}_n)^\top \phi(\mathbf{k}_i)$ de-
749 nominator in Eq. 2, γ is a decay factor as in RetNet, and ϕ is a modified *learnable* feature map from
750 Transformer-to-RNN (T2R) [Kasai et al. \(2021\)](#) with rotary embeddings [Su et al. \(2024\)](#). In other
751 words, $\phi(\mathbf{x}) = \text{RoPE}(\text{ReLU}(\mathbf{x}\mathbf{W} + \mathbf{b}))$ with $\mathbf{W} \in \mathbb{R}^{d \times d}$ and $\mathbf{b} \in \mathbb{R}^d$ as trainable weights and
752 biases. With this approach, they recover zero-shot capabilities in linearized Llama 2 7B [Touvron](#)
753 [et al. \(2023b\)](#) and Mistral 7B [Jiang et al. \(2023\)](#) models on popular LM Evaluation Harness [Gao](#)
754 [et al. \(2023\)](#) and SCROLLS [Shaham et al. \(2022\)](#) tasks.

755 ⁵[https://huggingface.co/datasets/togethercomputer/](https://huggingface.co/datasets/togethercomputer/RedPajama-Data-1T-Sample)
[RedPajama-Data-1T-Sample](https://huggingface.co/datasets/togethercomputer/RedPajama-Data-1T-Sample)

Hedgehog. Zhang et al. (2024) show we can train linear attentions to approximate softmax attentions, improving linearized model quality by swapping in the linear attentions as learned drop-in replacements. They use the standard linear attention (Eq. 2), where query, key, value, and output projections (the latter combining outputs in multi-head attention (Vaswani et al., 2017)) are first copied from an existing softmax attention. They then specify learnable feature maps $\phi(\mathbf{x}) = [\text{softmax}(\mathbf{x}\mathbf{W} + \mathbf{b}) \oplus \text{softmax}(-\mathbf{x}\mathbf{W} - \mathbf{b})]$ (where \oplus denotes concatenation, and both \oplus and the softmax are applied over the *feature dimension*) for \mathbf{q} and \mathbf{k} in each head and layer, and train ϕ such that linear attention weights $\hat{\mathbf{a}}$ match a Transformer’s original softmax weights \mathbf{a} . Given some sample data, they update ϕ with a cross-entropy-based distillation to minimize:

$$\mathcal{L}_n = - \sum_{i=1}^n \frac{\exp(\mathbf{q}_n^\top \mathbf{k}_i / \sqrt{d})}{\sum_{i=1}^n \exp(\mathbf{q}_n^\top \mathbf{k}_i / \sqrt{d})} \log \frac{\phi(\mathbf{q}_n)^\top \phi(\mathbf{k}_i)}{\sum_{i=1}^n \phi(\mathbf{q}_n)^\top \phi(\mathbf{k}_i)} \quad (10)$$

as the softmax and linear attention weights are both positive and sum to 1. As they focus on task-specific linearization (e.g., GLUE classification (Wang et al., 2018) or WikiText-103 language modeling (Merity et al., 2017)), for both attention and model training they use task-specific training data. By doing this “attention distillation”, they show significant linearized quality improvements over T2R on both smaller Transformers (e.g., 110M parameter BERTs Devlin et al. (2018) and 125M GPT-2s Radford et al. (2019)), and Llama 2 7B for a specific SAMSUM summarization task Gliwa et al. (2019).

C CODE IMPLEMENTATION

We include sample code for implementing LOLCATs with HuggingFace Transformers API.

```

780 1 def compute_loss(self, model: nn.Module, data: dict[torch.Tensor], **
781     kwargs: any,):
782     """
783     Attention distillation ("attention transfer")
784     - For each layer and head, get attentions and train to
785     minimize some combo of MSE and cross-entropy loss
786     """
787     input_seq_len = data['input_ids'].shape[-1]
788     inputs = {'input_ids': data['input_ids'].to(model.device)}
789
790     # Get softmax attention outputs
791     with torch.no_grad():
792         # Set base_inference to True to use FlashAttention
793         for layer in traverse_layers(model):
794             layer.self_attn.base_inference = True
795         # Get hidden states
796         true_outputs = model(**inputs, output_attentions=True,
797             use_cache=False,)
798         # Save attention layer inputs and outputs in outputs.attentions
799         # attn_inputs = [a[0] for a in true_outputs.get('attentions')]
800         # attn_outputs = [a[1] for a in true_outputs.get('attentions')]
801         true_attn_io = true_outputs.get('attentions') # layer-wise attn
802         inputs and outputs
803         true_outputs = true_outputs.get('logits').cpu()
804         for layer in traverse_layers(model):
805             layer.self_attn.base_inference = False
806
807     # Get trainable subquadratic attention outputs
808     attention_type = getattr(layer.self_attn, 'attention_type', None)
809     past_key_values = get_attention_cache(attention_type)
810
811     total_seq_len = 0
812     position_ids = torch.arange(input_seq_len).view(1, -1)
813
814     loss_mse = 0
815     for layer_idx, layer in enumerate(traverse_layers(model)):
816         attn_input, attn_output = true_attn_io[layer_idx]

```

```

810 36         attn_preds = layer.self_attn(attn_input.to(model.device),
811 37             attention_mask=None,
812 38             position_ids=position_ids,
813 39             past_key_value=past_key_values)[1]
814 40         # MSE on layer outputs
815 41         loss_mse += criterion_mse(attn_preds, attn_output)
816 42         loss_mse = loss_mse / (layer_idx + 1) * self.mse_factor
817 43         loss = loss_mse

```

Listing 1: Attention Distillation Code

```

819
820 1 class LolcatsLlamaAttention(nn.Module):
821 2     """
822 3     Hedgehog attention implementation initialized from a
823 4     `LlamaAttention` or `MistralAttention` object (base_attn)
824 5
825 6     Most of the arguments are directly tied to argparse args
826 7
827 8     Note that we don't currently support padding.
828 9     """
829 10    def __init__(self,
830 11                base_attn: nn.Module, # like LlamaAttention
831 12                feature_map: str,
832 13                feature_map_kwargs: dict,
833 14                layer_idx: Optional[int] = None,
834 15                max_layer_idx: Optional[int] = None,
835 16                feature_map_mlp: Optional[str] = None,
836 17                feature_map_mlp_kwargs: Optional[dict] = None,
837 18                tie_qk_fmap: Optional[bool] = False,
838 19                rotary_config: Optional[dict] = None,
839 20                attention_type: Optional[str] = 'hedgehog_llama',
840 21                mask_value: int = 0,
841 22                eps: float = 1e-12,):
842 23        super().__init__()
843 24
844 25        self.mask_value = mask_value
845 26        self.eps = eps
846 27        self.layer_idx = (layer_idx if layer_idx is not None
847 28                        else base_attn.layer_idx)
848 29        self.max_layer_idx = max_layer_idx
849 30
850 31        self.rotary_config = rotary_config
851 32
852 33        self.tie_qk_fmap = tie_qk_fmap
853 34        self.init_feature_map_(feature_map, feature_map_kwargs,
854 35                               feature_map_mlp, feature_map_mlp_kwargs)
855 36        self.init_weights_(base_attn)
856 37
857 38    def init_feature_map_(self,
858 39                        feature_map: str,
859 40                        feature_map_kwargs: dict,
860 41                        feature_map_mlp: str = None,
861 42                        feature_map_mlp_kwargs: dict = None):
862 43        """
863 44        Initialize feature map
864 45        """
865 46        if feature_map_mlp is not None:
866 47            feature_map_kwargs['num_heads'] = self.num_heads
867 48            feature_map_kwargs['head_dim'] = self.head_dim
868 49            feature_map_kwargs['dtype'] = self.q_proj.weight.dtype
869 50            feature_map_kwargs['device'] = self.q_proj.weight.device
870 51            feature_map_mlp = init_feature_map_mlp(feature_map_mlp,
871 52                                                  feature_map_mlp_kwargs)
872 53        self.feature_map_q = init_feature_map_act(name=feature_map,

```



```

864 54                                     mlp=feature_map_mlp,
865 55                                     **feature_map_kwargs)
866 56     if self.tie_qk_fmap: # tie mlp weights for query and key feature
867 57     maps
868 57         self.feature_map_k = self.feature_map_q
869 58     else:
870 59         self.feature_map_k = copy.deepcopy(self.feature_map_q)
871 60
872 61 def init_weights_(self, base_attn: nn.Module):
873 62     """
874 63     Initialize module layers, weights, positional dependencies, etc.
875 64     """
876 65     self.attention_dropout = 0 # We don't use dropout
877 66     self.hidden_size = base_attn.hidden_size
878 67     self.num_heads = base_attn.num_heads
879 68     self.head_dim = base_attn.head_dim
880 69     self.num_key_value_heads = base_attn.num_key_value_heads
881 70     self.num_key_value_groups = base_attn.num_key_value_groups
882 71
883 72     self.q_shape = [self.num_heads, self.head_dim]
884 73     self.k_shape = [self.num_key_value_heads, self.head_dim]
885 74     self.v_shape = [self.num_key_value_heads, self.head_dim]
886 75
887 76     self.max_position_embeddings = base_attn.max_position_embeddings
888 77     device = base_attn.q_proj.weight.device
889 78     scaling_factor = getattr(base_attn.rotary_emb, 'scaling_factor',
890 79     1.)
891 80     if self.rotary_config is None:
892 81         self.rotary_emb = get_rotary_embeddings(
893 82             rope_scaling_type=None,
894 83             head_dim=self.head_dim,
895 84             max_position_embeddings=base_attn.rotary_emb.
896 85             max_position_embeddings,
897 86             rope_theta=base_attn.rotary_emb.base,
898 87             rope_scaling_factor=scaling_factor,
899 88             device=device,
900 89         )
901 90     else:
902 91         if 'device' not in self.rotary_config:
903 92             self.rotary_config['device'] = device
904 93             self.rotary_emb = get_rotary_embeddings(**self.rotary_config)
905 94
906 95     # Just initialize with original weights
907 96     device = base_attn.q_proj.weight.device
908 97     self.q_proj = base_attn.q_proj
909 98     self.k_proj = base_attn.k_proj
910 99     self.v_proj = base_attn.v_proj
911 100     self.o_proj = base_attn.o_proj
912 101     del base_attn
913 102
914 103 def linear_attention(self, q: torch.Tensor, k: torch.Tensor, v: torch
915 104 .Tensor) -> Tuple[torch.Tensor, Optional[torch.Tensor], Optional[
916 105 Tuple[torch.Tensor]]]:
917 106     """
918 107     Compute linear attention with CUDA kernel implementation from
919 108     fast-transformers
920 109     """
921 110     dtype = q.dtype
922     y = causal_dot_product(q.contiguous().to(dtype=torch.float32),
923                             k.contiguous().to(dtype=torch.float32),
924                             v.contiguous().to(dtype=torch.float32)).to
925                             (dtype=dtype))
926     y = y / (torch.einsum("bhld,bhld->bhl", q, k.cumsum(dim=2)) +
927             self.eps)[..., None]
928     return y, None, None

```

```

918 111
919 112     def forward(self,
920 113             hidden_states: torch.Tensor,
921 114             attention_mask: Optional[torch.Tensor] = None,
922 115             position_ids: Optional[torch.LongTensor] = None,
923 116             past_key_value: Optional[Tuple[int, torch.Tensor, torch.
Tensor]] = None,
924 117             output_attentions: bool = False,
925 118             use_cache: bool = False,
926 119             **kwargs) -> Tuple[torch.Tensor, Optional[torch.Tensor],
Optional[Tuple[torch.Tensor]]]:
927 120         """
928 121         Forward pass modified from transformers.models.mistral.
929 122         modeling_mistral (v4.36)
930 123         """
931 124         b, l, _ = hidden_states.size()
932 125         q = self.q_proj(hidden_states)
933 126         k = self.k_proj(hidden_states)
934 127         v = self.v_proj(hidden_states)
935 128         kv_seq_len = k.shape[-2]
936 129
937 130         q = q.view(b, l, *self.q_shape).transpose(1, 2)
938 131         k = k.view(b, l, *self.k_shape).transpose(1, 2)
939 132         v = v.view(b, l, *self.v_shape).transpose(1, 2)
940 133
941 134         if past_key_value is not None:
942 135             kv_seq_len += past_key_value[0].shape[-2]
943 136
944 137         cos, sin = self.rotary_emb(k, seq_len=kv_seq_len)
945 138         q, k = apply_rotary_pos_emb(q, k, cos, sin, position_ids)
946 139
947 140         k = repeat_kv(k, self.num_key_value_groups)
948 141         v = repeat_kv(v, self.num_key_value_groups)
949 142         q, k = self.feature_map_q(q), self.feature_map_k(k)
950 143
951 144         if attention_mask is not None and q.shape[2] > 1:
952 145             lin_attn_mask = attention_mask[:, None, :, None]
953 146             k = k.masked_fill(~lin_attn_mask, self.mask_value)
954 147
955 148         if past_key_value is not None:
956 149             kv_state = past_key_value.kv_states[self.layer_idx]
957 150             k_state = past_key_value.k_states[self.layer_idx]
958 151
959 152             y_true, _, _ = self.linear_attention(q, k, v)
960 153             past_key_value.update(k, v, self.layer_idx)
961 154         else:
962 155             y_true, _, _ = self.linear_attention(q, k, v)
963 156
964 157         y_true = y_true.transpose(1, 2).contiguous().view(b, l, self.
hidden_size)
965 158         y_true = self.o_proj(y_true)
966 159         attn_weights = None
967 160
968 161         return y_true, attn_weights, past_key_value

```

Listing 2: LoLCATs Attention Implementation

```

966 1 class HedgehogFeatureMap(nn.Module):
967 2     """
968 3     Final 'activation' of feature map. Can probably be combined with
969 4     `HedgehogFeatureMapMLP` below
970 5
971 6     Full feature map is like f(xW + b)
972 7     -> This is the 'f' part
973 8     """

```

```

972 9 def __init__(self,
973 10     head_dim_idx: int = -1,
974 11     eps: float = 1e-12,
975 12     mlp: nn.Module = None,
976 13     halfspace: bool = False,
977 14     ):
978 15     super().__init__()
979 16     self.head_dim_idx = head_dim_idx
980 17     self.eps = eps
981 18     self.mlp = mlp if mlp is not None else nn.Identity()
982 19     self.activation = (self.halfspace_activation if halfspace
983 20                       else self.fullspace_activation)
984 21
985 22 def fullspace_activation(self, x: torch.Tensor):
986 23     return torch.cat([
987 24         torch.softmax(x, dim=self.head_dim_idx),
988 25         torch.softmax(-x, dim=self.head_dim_idx)
989 26     ], dim=self.head_dim_idx).clamp(min=self.eps)
990 27
991 28 def halfspace_activation(self, x: torch.Tensor):
992 29     return torch.softmax(x, dim=self.head_dim_idx).clamp(min=self.eps)
993 30 )
994 31
995 32 def forward(self, x: torch.Tensor):
996 33     """
997 34     Assume x.shape is (batch_size, n_heads, seq_len, head_dim)
998 35     """
999 36     return self.activation(self.mlp(x))
1000 37
1001 38 class HedgehogFeatureMapMLP(nn.Module):
1002 39     """
1003 40     Learnable MLP in feature map.
1004 41
1005 42     Full feature map is like f(xW + b)
1006 43     -> This is the 'W' and (optional) 'b' part
1007 44     """
1008 45     def __init__(self,
1009 46         num_heads: int,
1010 47         head_dim: int, # input dim
1011 48         feature_dim: int, # output dim
1012 49         dtype: torch.dtype,
1013 50         device: torch.device,
1014 51         skip_connection: bool = False,
1015 52         bias: bool = False):
1016 53         super().__init__()
1017 54         self.num_heads = num_heads
1018 55         self.head_dim = head_dim
1019 56         self.feature_dim = feature_dim
1020 57         self.dtype = dtype
1021 58         self.device = device
1022 59         self.skip_connection = skip_connection
1023 60         self.bias = bias
1024 61         self.init_weights_()
1025 62
1026 63     def init_weights_(self):
1027 64         """
1028 65         Initialize W and b
1029 66         """
1030 67         self.weight = nn.Parameter(torch.zeros(
1031 68             (self.num_heads, self.head_dim, self.feature_dim),
1032 69             dtype=self.dtype, device=self.device,
1033 70         ))
1034 71         nn.init.kaiming_uniform_(self.weight)
1035 72

```

```
1026 73     if self.bias:
1027 74         self.bias = nn.Parameter(torch.zeros(
1028 75             (1, self.num_heads, 1, self.feature_dim),
1029 76             dtype=self.dtype, device=self.device,
1030 77         ))
1031 78         nn.init.kaiming_uniform_(self.bias)
1032 79     else:
1033 80         self.bias = 0. # hack
1034 81
1035 82     def forward(self, x: torch.Tensor):
1036 83         """
1037 84         Assume x.shape is (batch_size, num_heads, seq_len, head_dim)
1038 85         """
1039 86         _x = torch.einsum('hdf,bhld->bhlf', self.layer, x) + self.bias
1040 87         return x + _x if self.skip_connection else _x
```

Listing 3: Hedgehog Learnable Feature Map Implementation

1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
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
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133

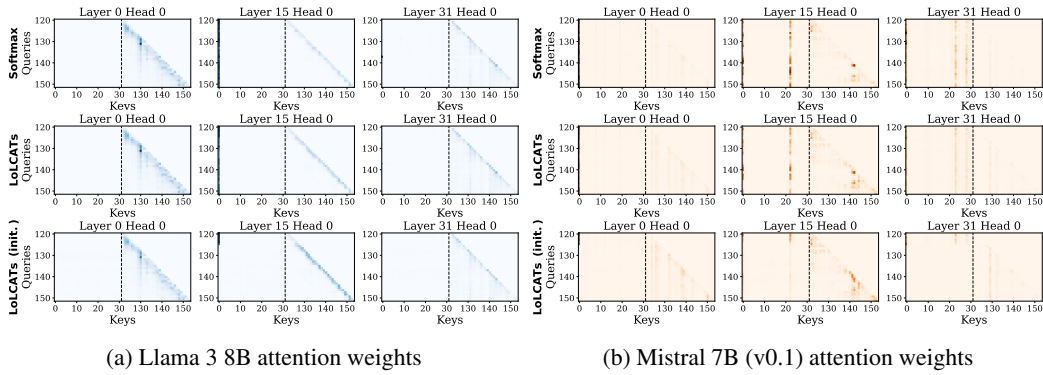


Figure 9: **Attention Transfer.** For both Llama 3 8B and Mistral 7B v0.1 LLMs, LoLCATs attention transfer trains subquadratic attentions that match original attention weights, despite only supervising based on attention layer outputs. They also learn to recover weights outside of the softmax windows, *c.f.* trained versus initialized (init.) attentions between queries at positions 130 - 150 and keys at positions 0 - 32.

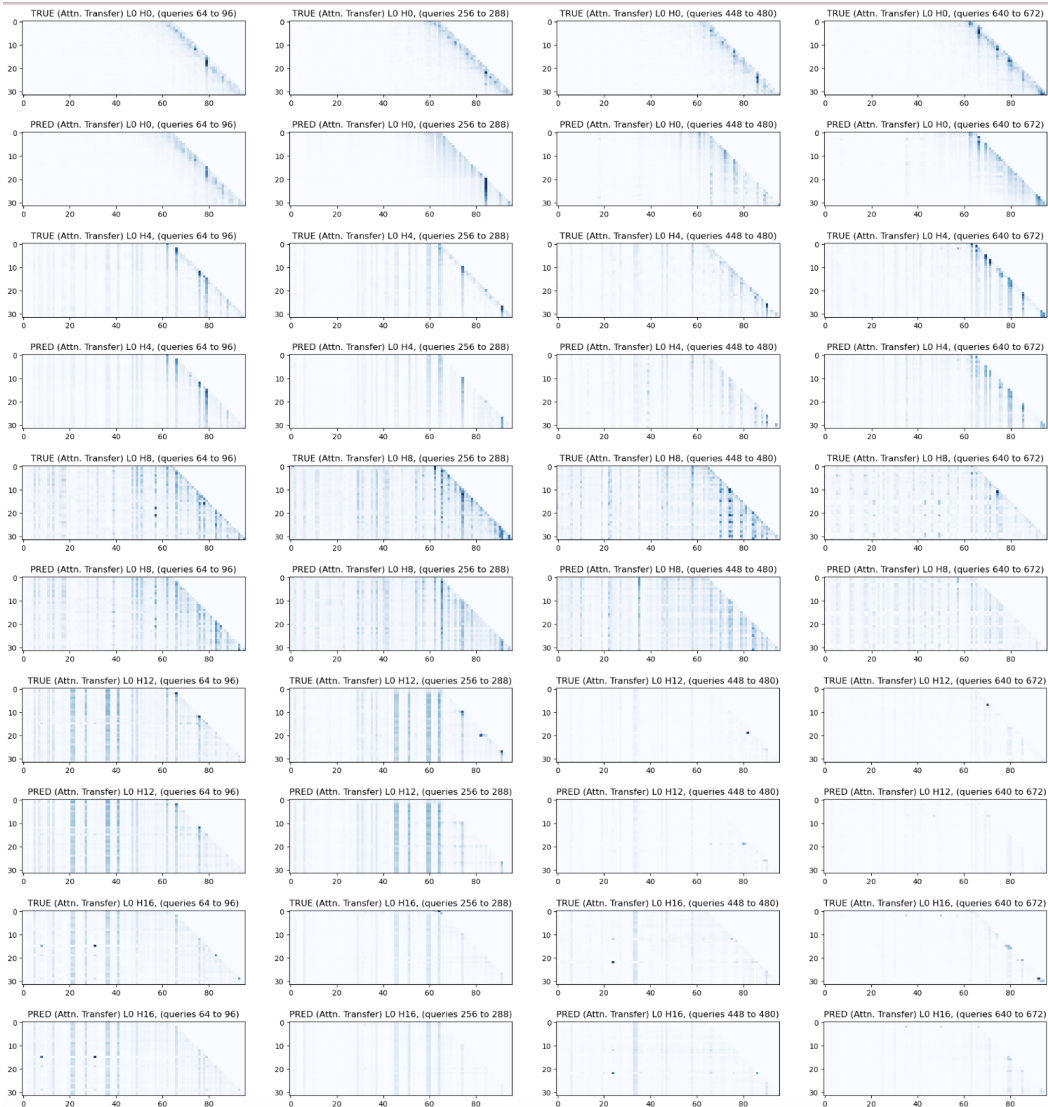


Figure 10: Linear attention (PRED) and softmax attention (TRUE) weights for hedgehog learned feature map, with attention transfer.

1134
 1135
 1136
 1137
 1138
 1139
 1140
 1141
 1142
 1143
 1144
 1145
 1146
 1147
 1148
 1149
 1150
 1151
 1152
 1153
 1154
 1155
 1156
 1157
 1158
 1159
 1160
 1161
 1162
 1163
 1164
 1165
 1166
 1167
 1168
 1169
 1170
 1171
 1172
 1173
 1174
 1175
 1176
 1177
 1178
 1179
 1180
 1181
 1182
 1183
 1184
 1185
 1186
 1187



Figure 11: Linear attention (PRED) and softmax attention (TRUE) weights for hedgehog learned feature map, without attention transfer.