

DENSE BACKPROPAGATION IMPROVES ROUTING FOR SPARSELY-GATED MIXTURE-OF-EXPERTS

Anonymous authors

Paper under double-blind review

ABSTRACT

Mixture of Experts (MoE) pretraining is more scalable than dense Transformer pretraining, because MoEs learn to route inputs to a sparse set of their feedforward parameters. However, this means that MoEs only receive a sparse backward update, leading to problems such as router load imbalance where some experts receive more tokens than others. We present a lightweight approximation method that gives the MoE a dense gradient while only sparsely activating its parameters. A key insight into the design of our method is that at scale, many tokens not routed to a given expert may nonetheless lie in the span of tokens that were routed to that expert, allowing us to create an approximation for the expert output of that token from existing expert outputs. Our dense backpropagation outperforms standard TopK routing across multiple settings without significantly increasing runtime.

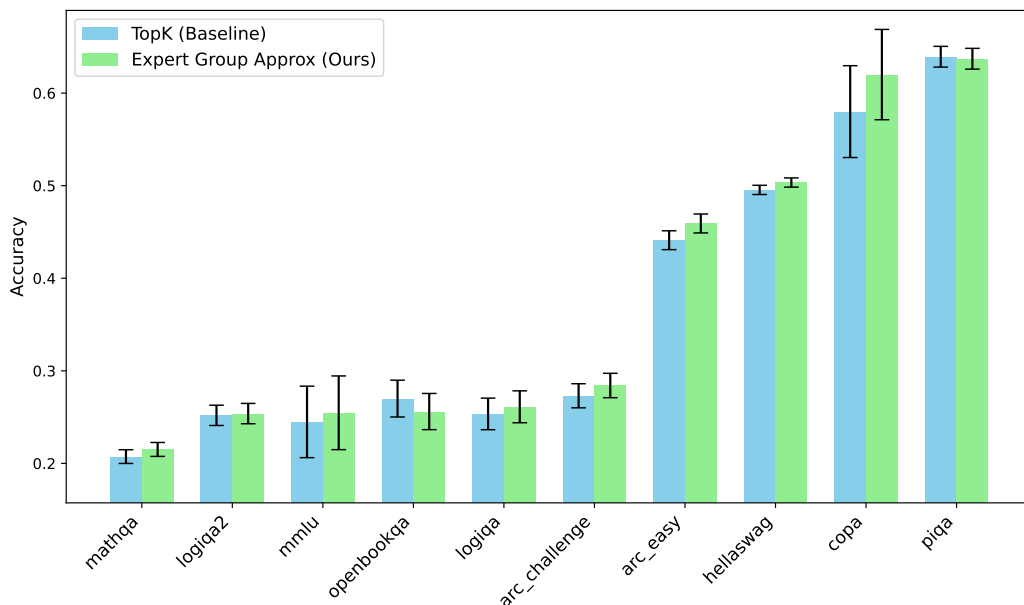


Figure 1: Our dense backpropagation method (green) improves over the baseline (blue) that uses Top-K routing across a range of standard benchmarks. See Section 4 for the experimental setup.

1 INTRODUCTION

Large-scale pretraining hinges on scaling up the number of parameters in a model, because models with more parameters are more sample-efficient and require less training to reach the same performance as smaller models (Kaplan et al., 2020; Hoffmann et al., 2022). The most common model architecture is a dense Transformer architecture (Vaswani et al., 2023) because its performance scales well with parameters and data. However, a sparsely activated Mixture-of-Experts (MoE) Transformer architecture (Shazeer et al., 2017) has been used by many industry deployments (Team et al.,

2024; xAI, 2024; Databricks, 2024; Jiang et al., 2024; Snowflake, 2024; DeepSeek-AI et al., 2024) because MoEs have been shown to scale even better than dense Transformers (Clark et al., 2022; Du et al., 2022; Lepikhin et al., 2020; Fedus et al., 2022). MoEs learn a *routing* function that selectively activates the TopK subset of their modules, or *experts*, most relevant to a given input. This conditionally sparse activation (Jacobs et al., 1991; Jordan & Jacobs, 1994) allows us to multiplicatively increase the model parameter count without significantly increasing the cost of training or inference.

The sparse router enables MoEs to scale, but it also presents a challenge, because the router does not receive a gradient update from experts that it does not activate, and may not learn to route a token to its appropriate expert. This may cause load imbalance— where a few experts are over-utilized — leading to inefficient training and resource usage (Zoph et al., 2022; Zhou et al., 2022).

In this work we propose a new router that can receive a dense gradient update from a sparse forward pass to address the instability issues arising from sparse routing. Our method adds minimal overhead, but improves on the common Top-K routing in both performance and load balance.

2 BACKGROUND & RELATED WORK

MoEs. The MoE layer replaces the feedforward networks (FFN) of transformers and consists of two components : **1)** N FFNs (*experts*), $E_0(x), E_1(x), \dots, E_N(x)$ and **2)** a router that assigns tokens to experts. Each input to the MoE layer is processed by K experts where $K < N$, and this is the source of sparsity in MoEs. The K experts are chosen by the router, which is a learnable component that maps each token to a set of weights over the experts. The router performs a linear transformation $\mathbb{R}^{d_{\text{token}}} \rightarrow \mathbb{R}^N$ which produces logits; these are normalized using softmax, resulting in a probability distribution over the experts. With the router’s linear transformation parameterized by a matrix W , we can represent the expert weights π in the following way:

$$\pi \in \mathbb{R}^N = \text{Softmax}(Wx) \tag{1}$$

Once we have these expert weights, we apply a routing function to decide which of K experts to route and process this token through. We consider TopK because it is the most popular.

Top-K routing. A standard method to select K out of N experts given the expert weights is to select the experts corresponding to the K highest weights. Top-K routing (Fedus et al., 2022) passes the token to the K selected experts and averages the expert outputs using these weights to produce the final output. Experts not selected by the Top-K routing function do not process the token, and this introduces sparsity in MoEs. By representing the K chosen experts as the set A , we can express the output of the MoE layer as an average of expert outputs weighted by the router scores:

$$y = \sum_{i \in A} \pi_i E_i(x) \tag{2}$$

The expert weights serve two roles. They are used by the routing function to decide which of the K experts to process a token through, and also provide the weights for combining the outputs of the expert. The Top-K routing scheme makes the MoE layer desirable for training large, compute-efficient neural networks. It allows models to be scaled up, by way of increasing the total number of experts, while keeping the compute per token constant (as it is a function of K and not N).

The Router Gradient. Consider the gradient of the MoE layer’s output y with respect to the router parameters W . We express y as a function of W by combining Eq. (1) and Eq. (2). With the chain rule, we can backpropagate through this function by considering the gradient at each respective step:

$$\frac{\partial y}{\partial W} = \frac{\partial y}{\partial \pi} \frac{\partial \pi}{\partial W} \tag{3}$$

The first term in Eq. (3), $\frac{\partial y}{\partial \pi}$, is straightforward to compute because the steps in Eq. (1) are easily differentiable, as they consist of linear operations and activations. But Eq. (2) is not differentiable because Top-K expert selection transforms the continuous router weights $\pi \in \mathbb{R}^N$ into a discrete set of selected experts A with $\binom{N}{K}$ possible values. One way to get around backpropagation of nondifferentiable operations is to use the straight-through estimator (Bengio et al., 2013), which treats the operator as the identity function. With straight-through we bypass the Top-K routing function and Eq. (2) becomes the dot product between π and the vector of all $E_i(x)$ with the following gradient:

108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161

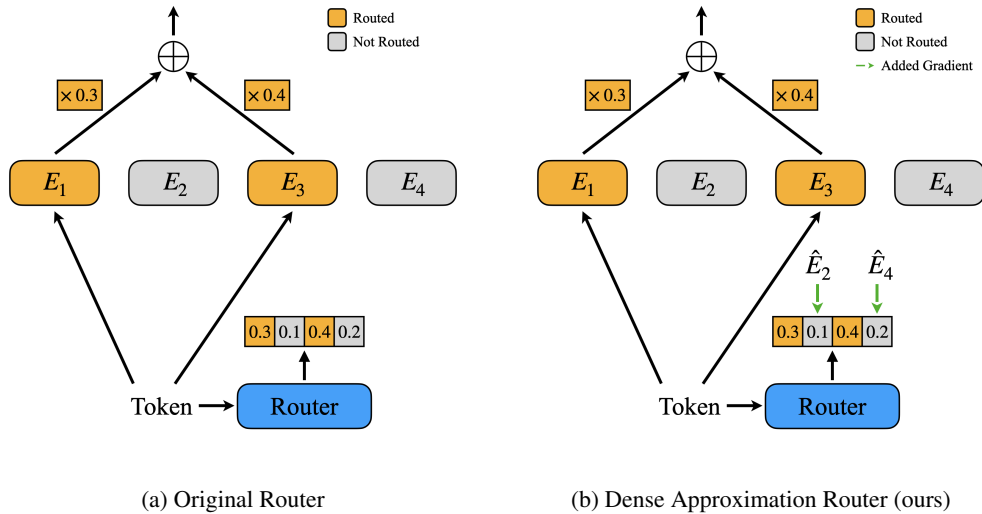


Figure 2: **Overview of Routing with Dense Approximations.** The original mixture of experts router only receives gradients corresponding to experts the token is routed to, because there is no output from other experts. Our approach provides the router with a complete (dense) gradient, by approximating the activations of experts that a token is not routed to. As indicated by the dashed green arrows, the approximated gradients are not actually connected to the token in the computation graph; instead, they are artificially applied in the backward pass.

$$\frac{\partial y}{\partial \pi} = [E_1(x), E_2(x) \dots E_N(x)] \quad (4)$$

This gradient requires the output of *all* of the experts for that token. Passing a token through all the experts will destroy the sparsity of the MoE layer, ruining the scalability of this architecture. In this work, we develop methods for applying the straight-through estimator while maintaining the sparsity of the MoE layer by *approximating* the output of the experts not selected by Top-K routing.

Related Works. Previous work has tried to address the issue of routing in MoEs. Separate from Top-K is the Sinkhorn routing method (Clark et al., 2022). Fedus et al. (2022) propose an auxiliary loss that encourages load balancing. Dai et al. (2024) propose multiple additional auxiliary loss terms. Recently, Wang et al. (2024b) propose learning biases rather than an auxiliary load balancing loss. Even more recently, Phi-3.5-MoE (Abdin et al., 2024) uses SparseMixer (Liu et al., 2024; 2023), another estimator for $\partial y / \partial \pi$ not involving straight-through. We provide a comparison between our approach and SparseMixer later in Section 4.3. Our approach is to still use straight-through, but *approximate* these additional expert outputs. We now present our method.

3 DESIGNING DENSE BACKPROPAGATION

In this section we design a new backpropagation method that can send a dense gradient to the sparsely activated router and expert parameters without doing additional forward passes. In a standard MoE, both the embedding corresponding to expert i in the routing layer (i.e. the i th row of the routing weight matrix) and the weights of expert i receive no gradient update from a token x if x is not routed to expert i . This is because $E_i(x)$ is never computed, so it provides no upstream gradient. This corresponds to experts that are not in the top K being omitted in Eq. (2). We apply an approximation $\hat{E}_i(x)$ as a substitute for the upstream gradient so that the router and expert weights can receive some non-zero signal corresponding to this expert. Then, the router can factor in outputs from all experts when learning to route each token, and the expert can be updated based on all tokens. Ultimately, instead of K/N expert parameters being updated per token, every token will update more parameters in both the router and the experts. Updating the router should improve load balance and the routing distribution.

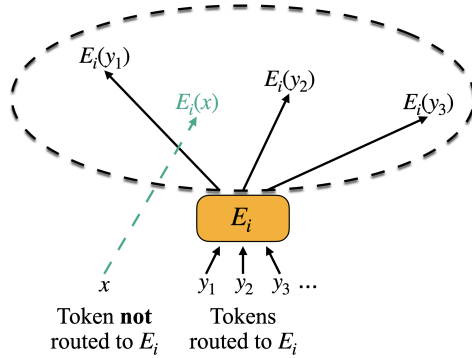


Figure 3: **Motivation for approximating expert outputs.** Consider an expert E_i . In every training batch, some fraction of tokens Y will be routed to E_i . With a large training batch size, we expect more tokens y_j to be routed to E_i . It is likely, then, that the outputs $E_i(y_j)$ will span the output space of this expert. Thus we expect that at the pretraining scale, a linear combination of expert outputs $\sum_j \alpha_j E_i(y_j)$ can approximate $E_i(x)$, for a token x which was not routed to E_i .

3.1 APPROXIMATING EXPERT ACTIVATIONS WITH DENSE BACKPROPAGATION

To approximate the dense gradient in Eq. (4), we must approximate $E_i(x)$ for every expert i that a token x was not passed to. Although we have no information about what the function E_i looks like for x , when training with large token batch sizes it is very likely that we have outputs of E_i for many other tokens. We hypothesize that although $E_i(x)$ is not computed for a given x , the output lies in the span of other tokens’ activations for E_i . In other words, $\hat{E}_i(x) = \sum_j c_j E_i(x_j)$ for some other tokens x_j . Our approximation method relies on finding these relevant x_j and the corresponding weights c_j to develop an approximation for a token. We visualize this general approach in Fig. 3.

Dense Backpropagation via Expert Group Approximation. Our dense backpropagation method that we dub *expert group approximation* produces an estimator $\hat{E}_i(x)$ from the expert outputs of other relevant tokens. We apply a single approximation to a large group of tokens that were not routed to expert i . This is efficient, but it may not be a good approximation for any specific x . Instead, this is an estimator for the expert output across the entire batch - this is sufficient as we will only need an approximation for the batch gradient to update the router during training. In dense backprop, we select relevant x_j and c_j purely based on the router output.

This expert approximation is used solely to estimate the dense gradient in equation 4. We do not use this dense approximation in the forward pass, as that is inconsistent with the sparse forward pass we expect at inference. Additionally, applying this approximation at inference is infeasible because of the much smaller batch size. However, we do use this approximation to update the experts along with the router. As we will show, updating all the experts is a crucial component of our method, and something prior work has not done. In particular, Liu et al. (2023) focuses on the router gradient.

3.2 NOTATION ON EXPERT ROUTING

Let $R(x)$ be the set of indices corresponding to the K experts that a token x is routed to. This can be thought of as the *routing decision* for x , based on the selected experts A in Eq. (2). For example, in a top- k sparse mixture of experts layer with $N = 8$ experts and $K = 2$, x routed to the first and last experts will have $R(x) = \{1, 8\}$. Note that $R(x)$ can have $\binom{N}{K}$ possible discrete outputs. With this notation, we can partition all tokens X based on their routing decisions and denote X_R as the subset of tokens routed to experts indexed by R . In the preceding example, the token routed to the first and last experts would belong to the set $X_{\{1,8\}}$. Some of our methods involve denoting whether a token was routed to a set of experts instead of its exact routing decision. We denote tokens routed to expert i along with any other experts as $X_{\{i,\cdot\}}$. For example, $X_{\{1,8\}} = X_{\{1,\cdot\}} \cap X_{\{8,\cdot\}}$

216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269

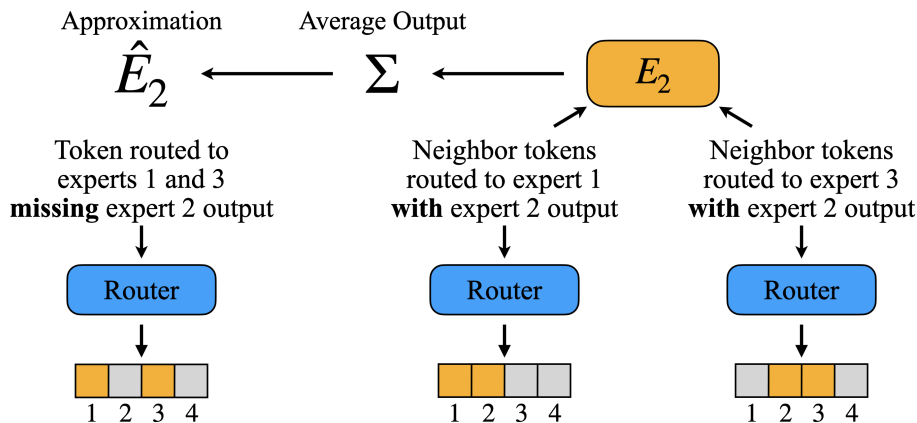


Figure 4: **Architecture of the Expert Group Approximation method.** In this example, we have 4 experts with $K = 2$. Consider all inputs routed to experts 1 and 3, characterized by the routing decision $R = \{1, 3\}$. As described in Figure 2b, we need to approximate these inputs’ activations for all other experts. In approximating expert 2, for example, we collect all inputs x' with a routing decision similar to R specifically including expert 2: $R' = \{1, 2\}$ and $R' = \{2, 3\}$. In general there will be K such adjacent groups. The aggregation of these inputs’ activations for expert 2 is used to approximate expert 2 for all inputs routed to experts 1 and 3.

3.3 EXPERT GROUP APPROXIMATION

Our method is based on approximating the expert output $E_i(x)$ for a group of multiple tokens at once. For a token x , we want to approximate outputs of experts that x was not routed to, i.e. $E_i(x)$ where $i \notin R(x)$. We hypothesize that tokens being routed to the same expert is a strong indicator of similarity between the tokens. This is because tokens routed to the same expert E_i both yield high dot products with the embedding W_i corresponding to that expert in the routing layer (see Eq. (1)). We develop an approximation for $E_i(x)$ by aggregating outputs of E_i for tokens that were routed to both expert i , and other experts that x was routed to. Formally, we consider an alternate routing decision $R' = \{i, j, \cdot\}, j \in R(x)$ that consists of one expert j that x is routed to, the expert i we wish to approximate, and any other experts (if $K > 2$). Then, the adjacent token space $X_{R'}$ will consist of tokens that are very similar to x by virtue of having similar routing decisions. Moreover, they will have an output for expert i , and we hypothesize that their outputs $\sum_{x' \in X_{R'}} E_i(x')$ will approximately represent $E_i(x)$. We can aggregate such outputs over all possible routing decisions:

$$\forall x \in X_R : \hat{E}_i(x) = \frac{1}{K} \sum_{j \in R} \frac{1}{|X_{\{i, j, \cdot\}}|} \sum_{x' \in X_{\{i, j, \cdot\}}} E_i(x') \quad (5)$$

We apply a single aggregate approximation for each routing decision to all tokens with that routing decision. In other words, each token belongs to K “groups” corresponding to the K experts it was routed to. Each expert has N groups, and each group j receives an approximation for expert i based on tokens routed to experts i, j . This gives us N^2 total approximations. When we approximate expert i , a token receives one approximation from each group, scaled by $\frac{1}{K}$ to take the average. In Fig. 4 we visualize this method for $K = 2$. We can pick different constants for the number of activated experts and the number of groups, so our method is applicable to $K = 1$, but we keep these constants at $K = 2$ throughout the paper because it is a common choice across prior work.

As mentioned above, we use this approximation to update both the router and the experts. This is necessary to maintain consistency, as updating the router densely but the experts sparsely can cause them to diverge. We also omit the approximation from the forward pass to only apply it in the backward pass. Consider an expert approximation based on other token outputs $\hat{E}_i(x) = \sum_j c_j E_i(x_j)$ assuming $\sum_j c_j = 1$. We produce such approximations for all experts in $R(x)^C$ that

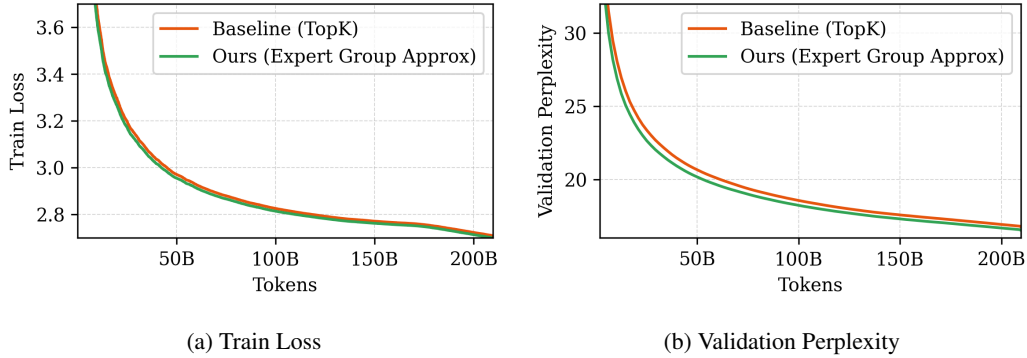


Figure 5: **Training and Validation Results.** We train a $2B$ total parameter MoE with 32 fine-grained experts, using the standard Top-K=2 router ($470M$ active params) and our expert group approximation method, each on 200B tokens of the Fineweb dataset. Without incurring significant overhead, we improve over the baseline.

x was not routed to: $y' = \sum_{i \in R(x)^c} \hat{E}_i(x)$. We update the original MoE output y from Eq. (2):

$$y := y + y' - \text{sg}(y') \quad (6)$$

where sg is the stop-gradient operator. As a result, our forward pass output is unchanged. During backpropagation, the expert approximations $\hat{E}_i(x)$ fill in the missing entries of the dense router gradient in equation 4. Moreover, the gradient to an expert i from the token x_j used to approximate $\hat{E}_i(x)$ is increased by the upstream gradient $\partial \hat{E}_i(x)$ times the weight c_j . In expectation, roughly K/N tokens are routed to an expert i so $(N - K)/N$ tokens will add gradients $\partial \hat{E}_i$ based on their approximations. Overall, this will lead to a dense gradient approximation for each expert since it will receive gradients as if it processed $K/N + (N - K)/N = N/N$ of all tokens instead of K/N . We all-reduce the approximation gradients across data parallel workers. This gives us a large sample size to estimate the dense gradient, regardless of the micro batch size on each worker.

4 EVALUATION

We train a fine-grained (DeepSeek-AI et al., 2024) MoE with 32 experts. The model has a hidden dim of 1024, with 2B total parameters. Each expert is a bottleneck MLP that has intermediate size 704. 470M parameters are activated when doing the standard $K = 2$ top-K routing (Zoph et al., 2022). We ablate the model architecture, hidden dim, number of total experts, and number of active experts through the course of the evaluation in Table 2 and Table 3. We train on 200B tokens from FineWeb (Penedo et al., 2024) with the Llama3 tokenizer (Dubey et al., 2024). We split it into train, validation, and test splits and report the validation perplexity. We defer other implementation details (learning rate schedule, initializations) to Appendix B.

4.1 MAIN RESULTS

Main Result. Our main result compares the Expert Group Approximation, which performs a dense update of the router weights by approximating the dense gradient, to baseline Top-K routing. In Fig. 5 we plot both training loss and validation set perplexity over the course of training. Our Expert Group Approximation method yields improvements over the baseline Top-K routing method throughout training, and this improvement is still present after training on 200B tokens. We further evaluate the performance of our Expert Group Approximation on standard benchmark tasks (that Penedo et al. (2024) notes are “high signal” for models trained on Fineweb) in Table 1 and find that our method provides a modest improvement over Top-K across multiple benchmarks.

Table 1: **Benchmarks.** Our method improves over the baseline across a set of benchmarks.

	mathqa	logiqa2	mmlu	openbookqa	logiqa	arc_challenge	arc_easy	hellaswag	copa	piqa	Average
Baseline	20.7 _{0.01}	25.2 _{0.01}	24.5 _{0.04}	27.0 _{0.02}	25.3 _{0.02}	27.3 _{0.01}	44.1 _{0.01}	49.5 _{0.00}	58.0 _{0.05}	63.9 _{0.01}	36.6 _{0.02}
Ours	21.5 _{0.01}	25.4 _{0.01}	25.5 _{0.04}	25.6 _{0.02}	26.1 _{0.02}	28.4 _{0.01}	45.9 _{0.01}	50.3 _{0.00}	62.0 _{0.05}	63.7 _{0.01}	37.4 _{0.02}
Improvement	+0.8%	+0.2%	+1.0%	-1.4%	+0.8%	+1.1%	+1.8%	+0.8%	+4.0%	-0.2%	+0.9%

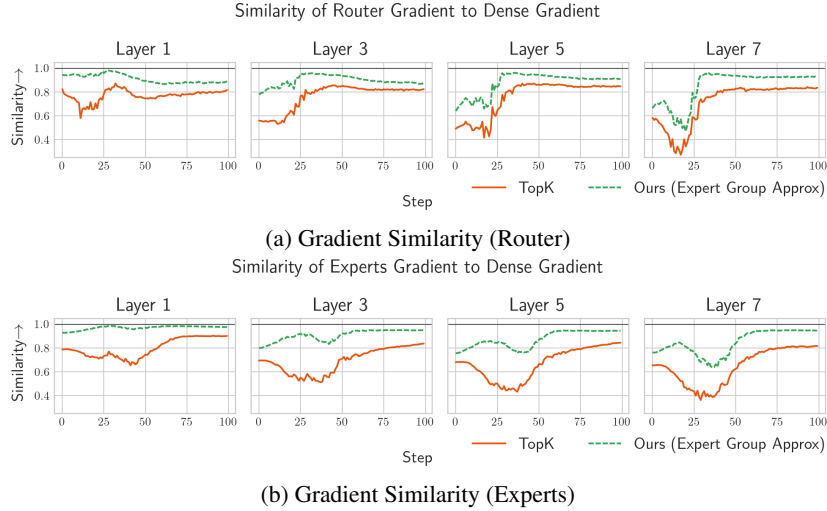


Figure 6: **Approximation fidelity.** We train MoEs with 32 fine-grained experts, recording the similarity with the true dense gradient (for the router and the experts) for the standard Top-K=2 baseline and our method. The gradients returned from our method are significantly closer to the true dense gradient. While prior work such as Liu et al. (2023) only aimed to densely update the router, we also densely update the experts, and we see that we actually do a better job of updating the experts than the router. Because Liu et al. (2023) do not design their method to approximate the true dense gradient of the router, we do not report their gradient similarity here (although it is low, their method achieves routing imbalance and performance in an entirely different mechanism to ours). *Our method faithfully reconstructs the true dense gradient signals.*

4.2 EMPIRICAL ANALYSIS

Gradient Approximation. We verify that our method is indeed faithfully approximating the dense router gradient i.e. the gradient to the router if all experts were activated. We track the dense gradient by routing to all experts and backpropagating only on the MoE output (independent of the full forward pass). In Fig. 6a we compare this dense gradient to the actual router gradient for each of our approaches and the standard Top-K router, and similarly for the expert weights in Fig. 6b. We plot the cosine similarity between the dense gradient and our approximated gradient, where the relatively high similarities our methods achieve indicate we are more accurately approximating the dense gradient. We also observe a major difference between the gradient of the standard Top-K router and our approach in later layers. We record the gradient similarity from the start of training until the trends stabilize; if we were to extend these plots out, we would see the similarities stay relatively stable. This validates our motivation to try and approximate the true dense gradient.

Load Balance. One reason for our method’s superior performance is in how it improves the routing distribution. Without a gradient signal for unactivated experts, top-K routing may not learn a balanced distribution across experts, and route more tokens to some experts than others. On the other hand, our method sends a signal to the router even for experts that a specific token was not routed to. The router now receives more information to make better decisions on how to route inputs. Load imbalancing leads to some experts processing a larger fraction of tokens than others, and thus having to learn a representation for more tokens. although all experts share the same number of parameters. Balancing expert load thus allows for more efficient utilization of MoE parameters.

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

Max Load Imbalance Per Layer

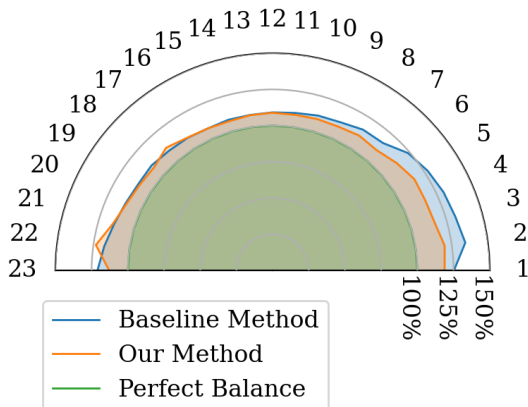


Figure 7: **Load Balancing using Expert Group Approximation.** We define maximum load imbalance at each layer as $\max_i(N \cdot f_i)$ where f_i is the fraction of tokens routed to each of N experts. The green ring indicates perfect balance, where each expert receives $1/N$ fraction of tokens; the outer ring indicates a maximum imbalance of 150%. We record maximum imbalance after training on 3 billion tokens. Our method improves load balance by sending a complete gradient to the router.

In Fig. 7 we validate that the baseline top- K ($K = 2$) routing has an “imbalanced load”, as measured by the proportion of tokens being routed to different experts relative to the optimal balance (green line) of an even distribution of tokens across experts. Our method improves load balance, which may be one cause for improved performance and is of independent interest on its own because it may lead to more efficient inference (although we do not profile inference workloads in this work).

4.3 ABLATIONS

Table 2: **Method scaling.** Our method’s training perplexity improvement after 10B tokens scales with the number of total experts, batch size, and number of active experts. The 8-expert MoE and 32-expert MoE both have $2B$ total parameters and use $K = 2$. However, in the 32-expert MoE, each MLP is a bottleneck layer. The 32-expert MoE therefore has just $470M$ active params, while the 8-expert MoE has $780M$ active params. We use the best number of experts (32) for the batch size scaling experiments, and we use the best batch size (2^{21}) for the expert scaling experiments. We also vary the number of active experts between $K = 2$ and $K = 4$ for the 32-expert MoE and find that our method provides consistent improvements across different values of K .

	Num. Experts		Batch Size			Active Experts	
	8	32	2^{19}	2^{20}	2^{21}	2	4
Baseline	22.03	22.23	27.54	24.51	22.23	22.23	21.54
Ours	21.80	21.91	27.51	24.37	21.91	21.91	21.22
Improvement	1.1%	1.5%	0.1%	0.6%	1.5%	1.5%	1.5%

Scaling. In Table 2 we show that our method’s improvement over the baseline increases as we scale up training. As we increase the number of experts from 8 in a standard MoE architecture, to 32 fine-grained experts, our improvement over the TopK baseline increases from 1.1% to 1.5%. This is because the 32-expert MoE has *more sparsity*, and therefore the gap between the baseline and the true dense gradient (that we are trying to approximate with our method) is larger. Our method constructs an approximation for unactivated experts from the rest of the tokens in the batch, so as the total training batch size increases, the likelihood that a linear combination of other tokens can

reproduce these missing expert outputs increases, as shown in Fig. 3. As we increase the train batch size, our improvement over the baseline increases from 0.1% at 2^{19} to 1.5% at 2^{21} . DeepSeek-AI et al. (2024) start with a batch size of $\approx 2^{23}$ and ramp up to a batch size of $\approx 2^{25}$, so the batch sizes we consider are still significantly smaller than those used to train the SOTA MoEs. *The improvement of our method over the baseline may be marginal at the current scale, but we believe that our method will continue to improve over the baseline as we scale up our models and batch sizes.*

Table 3: **Method variations.** We report perplexity after training an 8-expert $2B$ total-parameter MoE for $10B$ tokens. The “Accurate” variation on our method weights tokens based on how likely they were to be routed to the expert we’re constructing an approximation for. The “Viable” variation weights by the probability of the expert we’re using the approximation for. The “Accurate” variation is the best version of our method, because it gives us more accurate approximations.

Routing Method	Val. PPL ($K = 1$)	Improvement	Val. PPL ($K = 2$)	Improvement
Baseline	23.91	–	23.16	–
Expert Group Approx.	23.62	0.29%	22.62	0.54%
“Accurate”	23.48	0.43%	22.50	0.66%
“Viable”	23.49	0.42%	22.63	0.53%

Variations on the method. Thus far we have presented results where the approximation we use for a group of tokens is just a simple average of the expert outputs for similar tokens. However, when we construct the approximation we can weight tokens based on their router scores. In Table 3 we compare two variations on our main method. As a reminder, in this method we construct a mask of shape $experts, experts$ for each token. The row is the expert that token was routed to, and the column is the expert we want an approximation for. When we take the product of this mask and the router scores, we can weight each row by the probability corresponding to the expert we want an approximation for, or weight each approximation by the probability for the expert we’re using the approximation for. The former should give us more “accurate” approximations, because it will prioritize tokens that are more likely to be routed to the expert we want an approximation for. The latter should give us more “viable” approximations, because it weights by closeness to the space we’re using the approximation for. We find that both variants improve over our original method, with the “accurate” variation proving to be the best across $K = 1$ and $K = 2$. *Stacking the “accurate” variant on top of our method would further improve our main results.*

Scaling with Sparsity. We present most results in this paper with $K = 2$ because Zoph et al. (2022) finds that it improves performance over $K = 1$ with minimal overhead, and thus multiple prior works use $K = 2$. However, we can see in Table 3 that our method is *also better* at $K = 1$. Virtually zero MoE papers train with $K = 1$; the closest is Sun et al. (2024), and they use a shared expert (DeepSeek-AI et al., 2024) alongside the $K = 1$ expert, so in fact 2 experts are active at each layer. We believe this can be a new capability that our method unlocks, which in turn will lead to even greater sparsity and therefore more efficient MoEs.

Comparison to Sparsemixer. Liu et al. (2023; 2024) use Sparsemixer, which estimates the true router gradient without straight-through. Our method is distinct in that we also provide gradient updates for the expert weights (as seen in Fig. 6b). Liu et al. (2024) note that Sparsemixer lags behind TopK (which our method always outperforms) for the first $0.5T$ tokens, likely due to the noise that Sparsemixer adds. In the $10B$ tokens that we were able to run Sparsemixer for, our method significantly outperforms it.

Further Ablations. We conduct further ablations on design choices in Appendix C.1.

4.4 EFFICIENCY

Although our method outperforms the baseline in all settings we consider, to ensure a fair comparison it is crucial that our method does not significantly increase wall-clock time. In Table 4 we report the throughput in samples per second and TFLOPS of TopK, Sparsemixer, and our method while training a fine-grained MoE with 32 experts and $K = 2$ on a single GPU (for reproducibility). We find that the overhead of Sparsemixer and our method decreases as we increase the hidden size. The $2B$ model we train has 1024 hidden size, 2048 hidden size is an $8B$ model, etc. As the hidden size increases, the proportion of time spent in the MLP matmuls increases, and the overhead of our model

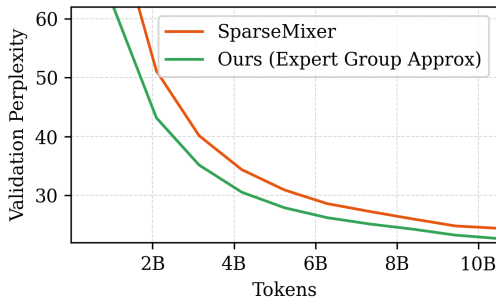


Figure 8: **SparseMixer comparison.** We train a fine-grained MoE with 32 experts with our method and SparseMixer (Liu et al., 2023). Our method, which outperforms the TopK baseline from the start of training, outperforms SparseMixer. Liu et al. (2024) report that SparseMixer requires many ($0.5T$) tokens to catch up to the TopK baseline.

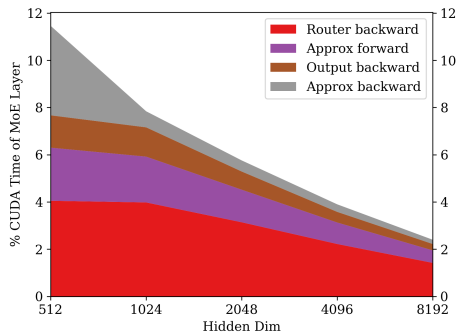


Figure 9: **CUDA overhead scaling.** The overhead of our method, as measured in the percent of the CUDA time it occupies, decreases as we increase the hidden size of the model. Our main experiments are conducted with a hidden size of 1024, resulting in 470M active parameters and 2B hidden parameters. A hidden size of 4096 would be commonly trained model.

Table 4: **Throughput Comparison.** We compare the throughput between different routing methods across hidden dimensions. Throughput is measured in tokens per second (sequence length 2048) on a single GPU. Overhead is calculated relative to TopK routing.

Hidden Dim	Tokens per second			Overhead vs TopK		
	TopK	SparseMixer	Ours	TopK	SparseMixer	Ours
1024	87,039	81,018	75,449	-	-6.92%	-13.32%
2048	12,288	11,469	11,790	-	-6.67%	-4.05%
4096	10,404	10,180	10,240	-	-2.15%	-1.57%

is no longer significant compared to the total MoE layer runtime. The same is true for SparseMixer, although our method is slightly faster than SparseMixer for hidden sizes 2048 and larger.

We further analyze this in Fig. 9, where we record the percent of CUDA time of the MoE layer that is spent in each part of our method. We develop Triton (Tillet et al., 2019) kernels to efficiently implement our method. Our code is currently open-source and will be linked here upon publication. For larger models, most of the overhead comes from the “Router backward” kernel where we compute the gradients to pass to update the router.

The overhead we report in Table 4 does not account for communication overhead. On our cluster, about half of the iteration time for the 1024 hidden size MoE is spent in NCCL operations when training on 8 nodes, and as we increase the number of nodes, the communication costs dominate following the trend of Wang et al. (2024a). Although performance varies dramatically across clusters, in our training setup our method adds near-zero overhead when training MoE with hidden size greater than 2048 (that is, more than 10B parameters).

5 DISCUSSION

We propose a training method for MoEs to improve load balancing and language modeling performance. By approximating the signal of a dense mixture-of-experts layer, the MoE router is able to learn a better distribution of routing inputs to different experts. This approximated dense signal unlocks the possibility for more sparse MoEs at training and inference time. Whereas typical Top-K routing would provide too sparse of a signal to learn a stable routing distribution, our method demonstrates significant improvements in load balancing in very sparse configurations. This unlocks better performance on standard language benchmarks, and has little overhead.

540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593

6 REPRODUCIBILITY STATEMENT

We provide detailed explanations of all our proposed methods in order to encourage reproducibility of our work: see Fig. 4 and Fig. 14. We use open-source LLM training libraries ([Andonian et al., 2023](#); [Gale et al., 2022](#)) for implementing our experiments, and we train our models using a publicly available dataset ([Penedo et al., 2024](#)). Our methods are built on top of these open-source libraries, and we plan to release our code publicly in the future. We also provide the details of our experiments in Section 4, and further details on hyperparameters are included in Appendix B.

594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647

REFERENCES

- Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed Awadallah, Ammar Ahmad Awan, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, Alon Benhaim, Misha Bilenko, Johan Bjorck, Sébastien Bubeck, Martin Cai, Qin Cai, Vishrav Chaudhary, Dong Chen, Dongdong Chen, Weizhu Chen, Yen-Chun Chen, Yi-Ling Chen, Hao Cheng, Parul Chopra, Xiyang Dai, Matthew Dixon, Ronen Eldan, Victor Fragoso, Jianfeng Gao, Mei Gao, Min Gao, Amit Garg, Allie Del Giorno, Abhishek Goswami, Suriya Gunasekar, Emman Haider, Junheng Hao, Russell J. Hewett, Wenxiang Hu, Jamie Huynh, Dan Iter, Sam Ade Jacobs, Mojan Javaheripi, Xin Jin, Nikos Karampatziakis, Piero Kauffmann, Mahoud Khademi, Dongwoo Kim, Young Jin Kim, Lev Kurilenko, James R. Lee, Yin Tat Lee, Yuanzhi Li, Yunsheng Li, Chen Liang, Lars Liden, Xihui Lin, Zeqi Lin, Ce Liu, Liyuan Liu, Mengchen Liu, Weishung Liu, Xiaodong Liu, Chong Luo, Piyush Madan, Ali Mahmoudzadeh, David Majercak, Matt Mazzola, Caio César Teodoro Mendes, Arindam Mitra, Hardik Modi, Anh Nguyen, Brandon Norrick, Barun Patra, Daniel Perez-Becker, Thomas Portet, Reid Pryzant, Heyang Qin, Marko Radmilac, Liliang Ren, Gustavo de Rosa, Corby Rosset, Sambudha Roy, Olatunji Ruwase, Olli Saarikivi, Amin Saied, Adil Salim, Michael Santacrose, Shital Shah, Ning Shang, Hiteshi Sharma, Yelong Shen, Swadheen Shukla, Xia Song, Masahiro Tanaka, Andrea Tupini, Praneetha Vaddamanu, Chunyu Wang, Guanhua Wang, Lijuan Wang, Shuohang Wang, Xin Wang, Yu Wang, Rachel Ward, Wen Wen, Philipp Witte, Haiping Wu, Xiaoxia Wu, Michael Wyatt, Bin Xiao, Can Xu, Jiahang Xu, Weijian Xu, Jilong Xue, Sonali Yadav, Fan Yang, Jianwei Yang, Yifan Yang, Ziyi Yang, Donghan Yu, Lu Yuan, Chenruidong Zhang, Cyril Zhang, Jianwen Zhang, Li Lyna Zhang, Yi Zhang, Yue Zhang, Yunan Zhang, and Xiren Zhou. Phi-3 technical report: A highly capable language model locally on your phone, 2024. URL <https://arxiv.org/abs/2404.14219>.
- Alex Andonian, Quentin Anthony, Stella Biderman, Sid Black, Preetham Gali, Leo Gao, Eric Hallahan, Josh Levy-Kramer, Connor Leahy, Lucas Nestler, Kip Parker, Michael Pieler, Jason Phang, Shivanshu Purohit, Hailey Schoelkopf, Dashiell Stander, Tri Songz, Curt Tigges, Benjamin Thérien, Phil Wang, and Samuel Weinbach. GPT-NeoX: Large Scale Autoregressive Language Modeling in PyTorch, 9 2023. URL <https://www.github.com/eleutherai/gpt-neox>.
- Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. Layer normalization, 2016. URL <https://arxiv.org/abs/1607.06450>.
- Yoshua Bengio, Nicholas Léonard, and Aaron Courville. Estimating or propagating gradients through stochastic neurons for conditional computation, 2013. URL <https://arxiv.org/abs/1308.3432>.
- Aidan Clark, Diego de las Casas, Aurelia Guy, Arthur Mensch, Michela Paganini, Jordan Hoffmann, Bogdan Damoc, Blake Hechtman, Trevor Cai, Sebastian Borgeaud, George van den Driessche, Eliza Rutherford, Tom Hennigan, Matthew Johnson, Katie Millican, Albin Cassirer, Chris Jones, Elena Buchatskaya, David Budden, Laurent Sifre, Simon Osindero, Oriol Vinyals, Jack Rae, Erich Elsen, Koray Kavukcuoglu, and Karen Simonyan. Unified scaling laws for routed language models, 2022. URL <https://arxiv.org/abs/2202.01169>.
- Damai Dai, Chengqi Deng, Chenggang Zhao, R. X. Xu, Huazuo Gao, Deli Chen, Jiashi Li, Wangding Zeng, Xingkai Yu, Y. Wu, Zhenda Xie, Y. K. Li, Panpan Huang, Fuli Luo, Chong Ruan, Zhifang Sui, and Wenfeng Liang. Deepseekmoe: Towards ultimate expert specialization in mixture-of-experts language models, 2024. URL <https://arxiv.org/abs/2401.06066>.
- Databricks. Dbrx, 2024. URL <https://www.databricks.com/blog/introducing-dbrx-new-state-art-open-llm>.
- DeepSeek-AI, Aixin Liu, Bei Feng, Bin Wang, Bingxuan Wang, Bo Liu, Chenggang Zhao, Chengqi Dengr, Chong Ruan, Damai Dai, Daya Guo, Dejian Yang, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Hanwei Xu, Hao Yang, Haowei Zhang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Li, Hui Qu, J. L. Cai, Jian Liang, Jianzhong Guo, Jiaqi Ni, Jiashi Li, Jin Chen, Jingyang Yuan, Junjie Qiu, Junxiao Song, Kai Dong, Kaige Gao, Kang Guan, Lean Wang, Lecong Zhang, Lei Xu, Leyi Xia, Liang Zhao, Liyue

648 Zhang, Meng Li, Miaojun Wang, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Mingming
649 Li, Ning Tian, Panpan Huang, Peiyi Wang, Peng Zhang, Qihao Zhu, Qinyu Chen, Qiushi Du, R. J.
650 Chen, R. L. Jin, Ruiqi Ge, Ruizhe Pan, Runxin Xu, Ruyi Chen, S. S. Li, Shanghao Lu, Shangyan
651 Zhou, Shanhuang Chen, Shaoqing Wu, Shengfeng Ye, Shirong Ma, Shiyu Wang, Shuang Zhou,
652 Shuiping Yu, Shunfeng Zhou, Size Zheng, T. Wang, Tian Pei, Tian Yuan, Tianyu Sun, W. L.
653 Xiao, Wangding Zeng, Wei An, Wen Liu, Wenfeng Liang, Wenjun Gao, Wentao Zhang, X. Q.
654 Li, Xiangyue Jin, Xianzu Wang, Xiao Bi, Xiaodong Liu, Xiaohan Wang, Xiaojin Shen, Xiaokang
655 Chen, Xiaosha Chen, Xiaotao Nie, Xiaowen Sun, Xiaoxiang Wang, Xin Liu, Xin Xie, Xingkai
656 Yu, Xinnan Song, Xinyi Zhou, Xinyu Yang, Xuan Lu, Xuecheng Su, Y. Wu, Y. K. Li, Y. X. Wei,
657 Y. X. Zhu, Yanhong Xu, Yanping Huang, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Li, Yaohui
658 Wang, Yi Zheng, Yichao Zhang, Yiliang Xiong, Yilong Zhao, Ying He, Ying Tang, Yishi Piao,
659 Yixin Dong, Yixuan Tan, Yiyuan Liu, Yongji Wang, Yongqiang Guo, Yuchen Zhu, Yudian Wang,
660 Yuheng Zou, Yukun Zha, Yunxian Ma, Yuting Yan, Yuxiang You, Yuxuan Liu, Z. Z. Ren, Zehui
661 Ren, Zhangli Sha, Zhe Fu, Zhen Huang, Zhen Zhang, Zhenda Xie, Zhewen Hao, Zhihong Shao,
662 Zhiniu Wen, Zhipeng Xu, Zhongyu Zhang, Zhuoshu Li, Zihan Wang, Zihui Gu, Zilin Li, and
663 Ziwei Xie. Deepseek-v2: A strong, economical, and efficient mixture-of-experts language model,
664 2024. URL <https://arxiv.org/abs/2405.04434>.

664 Nan Du, Yanping Huang, Andrew M Dai, Simon Tong, Dmitry Lepikhin, Yuanzhong Xu, Maxim
665 Krikun, Yanqi Zhou, Adams Wei Yu, Orhan Firat, et al. Glam: Efficient scaling of language
666 models with mixture-of-experts. In *International Conference on Machine Learning*, pp. 5547–
667 5569. PMLR, 2022.

668 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha
669 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony
670 Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark,
671 Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere,
672 Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris
673 Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong,
674 Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny
675 Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino,
676 Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael
677 Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson,
678 Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah
679 Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan
680 Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Ma-
681 hadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy
682 Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak,
683 Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Al-
684 wala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini,
685 Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Lauren Rantala-Yeary, Laurens van der
686 Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo,
687 Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Man-
688 nat Singh, Manohar Paluri, Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova,
689 Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal,
690 Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Olivier Duchenne, Onur
691 Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhar-
692 gava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong,
693 Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic,
694 Roberta Raileanu, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sum-
695 baly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa,
696 Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang,
697 Sharath Rapparthi, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende,
698 Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney
699 Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom,
700 Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta,
701 Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vladan Petro-
v, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang,
Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao Jia, Xuwei Wang, Yaelle Goldschlag, Yashesh Gaur,
Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre

702 Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aaron Grattafiori, Abha
703 Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay
704 Menon, Ajay Sharma, Alex Boesenberg, Alex Vaughan, Alexei Baevski, Allie Feinstein, Amanda
705 Kallet, Amit Sangani, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew
706 Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Franco, Aparajita
707 Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh
708 Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De
709 Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Bran-
710 don Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina
711 Mejia, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai,
712 Chris Tindal, Christoph Feichtenhofer, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li,
713 Danny Wyatt, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana
714 Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil,
715 Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Erik Brinkman, Esteban Ar-
716 caute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Firat Ozgenel, Francesco
717 Caggioni, Francisco Guzmán, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella
718 Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Govind Thattai, Grant Herman, Grigory
719 Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hamid Shojanazeri, Han Zou, Hannah Wang,
720 Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Gold-
721 man, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Irina-Elena Veliche, Itai Gat, Jake Weissman,
722 James Geboski, James Kohli, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer
723 Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe
724 Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie
725 Wang, Kai Wu, Kam Hou U, Karan Saxena, Karthik Prasad, Kartikay Khandelwal, Katayoun
726 Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kun Huang, Kunal
727 Chawla, Kushal Lakhota, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva,
728 Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian
729 Khabza, Manav Avalani, Manish Bhatt, Maria Tsimpoukelli, Martynas Mankus, Matan Hasson,
730 Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Ke-
731 neally, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel
732 Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mo-
733 hammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navy-
734 ata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikolay Pavlovich Laptev, Ning Dong,
735 Ning Zhang, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli,
736 Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux,
737 Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao,
738 Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Raymond Li,
739 Rebekkah Hogan, Robin Battey, Rocky Wang, Rohan Maheswari, Russ Howes, Ruty Rinott,
740 Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Sa-
741 tadru Pan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lind-
742 say, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shiva Shankar, Shuqiang Zhang, Shuqiang
743 Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen
744 Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Sungmin Cho,
745 Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser,
746 Tamara Best, Thilo Kohler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Tim-
747 othy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan,
748 Vinay Satish Kumar, Vishal Mangla, Vitor Albiero, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu
749 Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Con-
750 stable, Xiaocheng Tang, Xiaofang Wang, Xiaoqian Wu, Xiaolan Wang, Xide Xia, Xilun Wu,
751 Xinbo Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi,
752 Youngjin Nam, Yu, Wang, Yuchen Hao, Yundi Qian, Yuzi He, Zach Rait, Zachary DeVito, Zef
753 Rosnbrick, Zhaoduo Wen, Zhenyu Yang, and Zhiwei Zhao. The llama 3 herd of models, 2024.
754 URL <https://arxiv.org/abs/2407.21783>.

751 William Fedus, Barret Zoph, and Noam Shazeer. Switch transformers: Scaling to trillion parameter
752 models with simple and efficient sparsity, 2022. URL <https://arxiv.org/abs/2101.03961>.

753 Trevor Gale, Deepak Narayanan, Cliff Young, and Matei Zaharia. Megablocks: Efficient sparse
754 training with mixture-of-experts, 2022. URL <https://arxiv.org/abs/2211.15841>.

756 Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza
757 Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hen-
758 nigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy,
759 Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and Laurent Sifre.
760 Training compute-optimal large language models, 2022. URL [https://arxiv.org/abs/
761 2203.15556](https://arxiv.org/abs/2203.15556).

762 Adam Ibrahim, Benjamin Thérien, Kshitij Gupta, Mats L. Richter, Quentin Anthony, Timothée
763 Lesort, Eugene Belilovsky, and Irina Rish. Simple and scalable strategies to continually pre-train
764 large language models, 2024. URL <https://arxiv.org/abs/2403.08763>.

765 Robert A. Jacobs, Michael I. Jordan, Steven J. Nowlan, and Geoffrey E. Hinton. Adaptive Mixtures
766 of Local Experts. *Neural Computation*, 3(1):79–87, 03 1991. ISSN 0899-7667. doi: 10.1162/
767 neco.1991.3.1.79. URL <https://doi.org/10.1162/neco.1991.3.1.79>.

768 Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris
769 Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, Gi-
770 anna Lengyel, Guillaume Bour, Guillaume Lample, Léo Renard Lavaud, Lucile Saulnier, Marie-
771 Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le
772 Scao, Théophile Gervet, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed.
773 Mixtral of experts, 2024. URL <https://arxiv.org/abs/2401.04088>.

774 M. I. Jordan and R. A. Jacobs. Hierarchical mixtures of experts and the em algorithm. In Maria
775 Marinaro and Pietro G. Morasso (eds.), *ICANN '94*, pp. 479–486, London, 1994. Springer Lon-
776 don. ISBN 978-1-4471-2097-1.

777 Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child,
778 Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language
779 models, 2020. URL <https://arxiv.org/abs/2001.08361>.

780 Dmitry Lepikhin, HyoukJoong Lee, Yuanzhong Xu, Dehao Chen, Orhan Firat, Yanping Huang,
781 Maxim Krikun, Noam Shazeer, and Zhifeng Chen. Gshard: Scaling giant models with conditional
782 computation and automatic sharding. *arXiv preprint arXiv:2006.16668*, 2020.

783 Liyuan Liu, Jianfeng Gao, and Weizhu Chen. Sparse backpropagation for moe training, 2023. URL
784 <https://arxiv.org/abs/2310.00811>.

785 Liyuan Liu, Young Jin Kim, Shuohang Wang, Chen Liang, Yelong Shen, Hao Cheng, Xiaodong
786 Liu, Masahiro Tanaka, Xiaoxia Wu, Wenxiang Hu, Vishrav Chaudhary, Zeqi Lin, Chenruidong
787 Zhang, Jilong Xue, Hany Awadalla, Jianfeng Gao, and Weizhu Chen. Grin: Gradient-informed
788 moe, 2024. URL <https://arxiv.org/abs/2409.12136>.

789 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization, 2019. URL <https://arxiv.org/abs/1711.05101>.

790 Toan Q. Nguyen and Julian Salazar. Transformers without tears: Improving the normalization
791 of self-attention. 2019. doi: 10.5281/ZENODO.3525484. URL [https://zenodo.org/
792 record/3525484](https://zenodo.org/record/3525484).

793 Guilherme Penedo, Hynek Kydlíček, Loubna Ben allal, Anton Lozhkov, Margaret Mitchell, Colin
794 Raffel, Leandro Von Werra, and Thomas Wolf. The fineweb datasets: Decanting the web for the
795 finest text data at scale, 2024. URL <https://arxiv.org/abs/2406.17557>.

796 Noam Shazeer. Glu variants improve transformer, 2020. URL [https://arxiv.org/abs/
797 2002.05202](https://arxiv.org/abs/2002.05202).

798 Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc Le, Geoffrey Hinton,
799 and Jeff Dean. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer,
800 2017. URL <https://arxiv.org/abs/1701.06538>.

801 Snowflake. Arctic, 2024. URL [https://www.snowflake.com/en/blog/
802 arctic-open-efficient-foundation-language-models-snowflake/](https://www.snowflake.com/en/blog/arctic-open-efficient-foundation-language-models-snowflake/).

810 Jianlin Su, Yu Lu, Shengfeng Pan, Ahmed Murtadha, Bo Wen, and Yunfeng Liu. Roformer: En-
811 hanced transformer with rotary position embedding, 2023. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2104.09864)
812 [2104.09864](https://arxiv.org/abs/2104.09864).

813 Xingwu Sun, Yanfeng Chen, Yiqing Huang, Ruobing Xie, Jiaqi Zhu, Kai Zhang, Shuaipeng Li, Zhen
814 Yang, Jonny Han, Xiaobo Shu, Jiahao Bu, Zhongzhi Chen, Xuemeng Huang, Fengzong Lian,
815 Saiyong Yang, Jianfeng Yan, Yuyuan Zeng, Xiaoqin Ren, Chao Yu, Lulu Wu, Yue Mao, Jun Xia,
816 Tao Yang, Suncong Zheng, Kan Wu, Dian Jiao, Jinbao Xue, Xipeng Zhang, Decheng Wu, Kai Liu,
817 Dengpeng Wu, Guanghui Xu, Shaohua Chen, Shuang Chen, Xiao Feng, Yigeng Hong, Junqiang
818 Zheng, Chengcheng Xu, Zongwei Li, Xiong Kuang, Jianglu Hu, Yiqi Chen, Yuchi Deng, Guiyang
819 Li, Ao Liu, Chenchen Zhang, Shihui Hu, Zilong Zhao, Zifan Wu, Yao Ding, Weichao Wang, Han
820 Liu, Roberts Wang, Hao Fei, Peijie Yu, Ze Zhao, Xun Cao, Hai Wang, Fusheng Xiang, Mengyuan
821 Huang, Zhiyuan Xiong, Bin Hu, Xuebin Hou, Lei Jiang, Jianqiang Ma, Jiajia Wu, Yaping Deng,
822 Yi Shen, Qian Wang, Weijie Liu, Jie Liu, Meng Chen, Liang Dong, Weiwen Jia, Hu Chen, Feifei
823 Liu, Rui Yuan, Huilin Xu, Zhenxiang Yan, Tengfei Cao, Zhichao Hu, Xinhua Feng, Dong Gu,
824 Tinghao Yu, Yangyu Tao, Feng Zhang, Jianchen Zhu, Chengzhong Xu, Xirui Li, Chong Zha, Wen
825 Ouyang, Yinben Xia, Xiang Li, Zekun He, Rongpeng Chen, Jiawei Song, Ruibin Chen, Fan Jiang,
826 Chongqing Zhao, Bo Wang, Hao Gong, Rong Gan, Winston Hu, Zhanhui Kang, Yong Yang,
827 Yuhong Liu, Di Wang, and Jie Jiang. Hunyuan-large: An open-source moe model with 52 billion
828 activated parameters by tencent, 2024. URL <https://arxiv.org/abs/2411.02265>.

829 Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer,
830 Damien Vincent, Zhufeng Pan, Shibo Wang, Soroosh Mariooryad, Yifan Ding, Xinyang Geng,
831 Fred Alcober, Roy Frostig, Mark Omernick, Lexi Walker, Cosmin Paduraru, Christina Sorokin,
832 Andrea Tacchetti, Colin Gaffney, Samira Daruki, Olcan Sercinoglu, Zach Gleicher, Juliette Love,
833 Paul Voigtlaender, Rohan Jain, Gabriela Surita, Kareem Mohamed, Rory Blevins, Junwhan Ahn,
834 Tao Zhu, Kornrathop Kawintiranon, Orhan Firat, Yiming Gu, Yujing Zhang, Matthew Rahtz,
835 Manaal Faruqi, Natalie Clay, Justin Gilmer, JD Co-Reyes, Ivo Penchev, Rui Zhu, Nobuyuki
836 Morioka, Kevin Hui, Krishna Haridasan, Victor Campos, Mahdis Mahdieh, Mandy Guo, Samer
837 Hassan, Kevin Kilgour, Arpi Vezer, Heng-Tze Cheng, Raoul de Liedekerke, Siddharth Goyal,
838 Paul Barham, DJ Strouse, Seb Noury, Jonas Adler, Mukund Sundararajan, Sharad Vikram, Dmitry
839 Lepikhin, Michela Paganini, Xavier Garcia, Fan Yang, Dasha Valter, Maja Trebacz, Kiran Vo-
840 drahalli, Chulayuth Asawaroengchai, Roman Ring, Norbert Kalb, Livio Baldini Soares, Sid-
841 dhartha Brahma, David Steiner, Tianhe Yu, Fabian Mentzer, Antoine He, Lucas Gonzalez, Bibo
842 Xu, Raphael Lopez Kaufman, Laurent El Shafey, Junhyuk Oh, Tom Hennigan, George van den
843 Driessche, Seth Odoom, Mario Lucic, Becca Roelofs, Sid Lall, Amit Marathe, Betty Chan, San-
844 tiago Ontanon, Luheng He, Denis Teplyashin, Jonathan Lai, Phil Crone, Bogdan Damoc, Lewis
845 Ho, Sebastian Riedel, Karel Lenc, Chih-Kuan Yeh, Aakanksha Chowdhery, Yang Xu, Mehran
846 Kazemi, Ehsan Amid, Anastasia Petrushkina, Kevin Swersky, Ali Khodaei, Gowoon Chen, Chris
847 Larkin, Mario Pinto, Geng Yan, Adria Puigdomenech Badia, Piyush Patil, Steven Hansen, Dave
848 Orr, Sebastien M. R. Arnold, Jordan Grimstad, Andrew Dai, Sholto Douglas, Rishika Sinha, Vikas
849 Yadav, Xi Chen, Elena Gribovskaya, Jacob Austin, Jeffrey Zhao, Kaushal Patel, Paul Komarek,
850 Sophia Austin, Sebastian Borgeaud, Linda Friso, Abhimanyu Goyal, Ben Caine, Kris Cao, Da-
851 woon Chung, Matthew Lamm, Gabe Barth-Maron, Thais Kagohara, Kate Olszewska, Mia Chen,
852 Kaushik Shivakumar, Rishabh Agarwal, Harshal Godhia, Ravi Rajwar, Javier Snaider, Xerxes
853 Dotiwala, Yuan Liu, Aditya Barua, Victor Ungureanu, Yuan Zhang, Bat-Orgil Batsaikhan, Ma-
854 teo Wirth, James Qin, Ivo Danihelka, Tulsee Doshi, Martin Chadwick, Jilin Chen, Sanil Jain,
855 Quoc Le, Arjun Kar, Madhu Gurusurthy, Cheng Li, Ruoxin Sang, Fangyu Liu, Lampros Lam-
856 prou, Rich Munoz, Nathan Lintz, Harsh Mehta, Heidi Howard, Malcolm Reynolds, Lora Aroyo,
857 Quan Wang, Lorenzo Blanco, Albin Cassirer, Jordan Griffith, Dipanjan Das, Stephan Lee, Jakub
858 Sygnowski, Zach Fisher, James Besley, Richard Powell, Zafarali Ahmed, Dominik Paulus, David
859 Reitter, Zalan Borsos, Rishabh Joshi, Aedan Pope, Steven Hand, Vittorio Selo, Vihan Jain, Nikhil
860 Sethi, Megha Goel, Takaki Makino, Rhys May, Zhen Yang, Johan Schalkwyk, Christina Butter-
861 field, Anja Hauth, Alex Goldin, Will Hawkins, Evan Senter, Sergey Brin, Oliver Woodman, Mar-
862 vin Ritter, Eric Noland, Minh Giang, Vijay Bolina, Lisa Lee, Tim Blyth, Ian Mackinnon, Machel
863 Reid, Obaid Sarvana, David Silver, Alexander Chen, Lily Wang, Loren Maggiore, Oscar Chang,
Nithya Attaluri, Gregory Thornton, Chung-Cheng Chiu, Oskar Bunyan, Nir Levine, Timothy
Chung, Evgenii Eltyshev, Xiance Si, Timothy Lillcrap, Demetra Brady, Vaibhav Aggarwal, Boxi
Wu, Yuanzhong Xu, Ross McIlroy, Kartikeya Badola, Paramjit Sandhu, Erica Moreira, Wojciech
Stokowiec, Ross Hemsley, Dong Li, Alex Tudor, Pranav Shyam, Elahe Rahimtoroghi, Salem

864 Haykal, Pablo Sprechmann, Xiang Zhou, Diana Mincu, Yujia Li, Ravi Addanki, Kalpesh Krishna,
865 Xiao Wu, Alexandre Frechette, Matan Eyal, Allan Dafoe, Dave Lacey, Jay Whang, Thi Avrahami,
866 Ye Zhang, Emanuel Taropa, Hanzhao Lin, Daniel Toyama, Eliza Rutherford, Motoki Sano, Hyun-
867 Jeong Choe, Alex Tomala, Chalence Safranek-Shrader, Nora Kassner, Mantas Pajarskas, Matt
868 Harvey, Sean Sechrist, Meire Fortunato, Christina Lyu, Gamaleldin Elsayed, Chenkai Kuang,
869 James Lottes, Eric Chu, Chao Jia, Chih-Wei Chen, Peter Humphreys, Kate Baumli, Connie Tao,
870 Rajkumar Samuel, Cicero Nogueira dos Santos, Anders Andreassen, Nemanja Rakićević, Do-
871 minik Grewe, Aviral Kumar, Stephanie Winkler, Jonathan Caton, Andrew Brock, Sid Dalmia,
872 Hannah Sheahan, Iain Barr, Yingjie Miao, Paul Natsev, Jacob Devlin, Feryal Behbahani, Flavien
873 Prost, Yanhua Sun, Artiom Myaskovsky, Thanumalayan Sankaranarayanan Pillai, Dan Hurt, An-
874 geliki Lazaridou, Xi Xiong, Ce Zheng, Fabio Pardo, Xiaowei Li, Dan Horgan, Joe Stanton,
875 Moran Ambar, Fei Xia, Alejandro Lince, Mingqiu Wang, Basil Mustafa, Albert Webson, Hyo
876 Lee, Rohan Anil, Martin Wicke, Timothy Dozat, Abhishek Sinha, Enrique Piqueras, Elahe Dabir,
877 Shyam Upadhyay, Anudhyan Boral, Lisa Anne Hendricks, Corey Fry, Josip Djolonga, Yi Su,
878 Jake Walker, Jane Labanowski, Ronny Huang, Vedant Misra, Jeremy Chen, RJ Skerry-Ryan,
879 Avi Singh, Shruti Rijhwani, Dian Yu, Alex Castro-Ros, Beer Changpinyo, Romina Datta, Sumit
880 Bagri, Arnar Mar Hrafnkelsson, Marcello Maggioni, Daniel Zheng, Yury Sulsky, Shaobo Hou,
881 Tom Le Paine, Antoine Yang, Jason Riesa, Dominika Rogozinska, Dror Marcus, Dalia El Badawy,
882 Qiao Zhang, Luyu Wang, Helen Miller, Jeremy Greer, Lars Lowe Sjøs, Azade Nova, Heiga Zen,
883 Rahma Chaabouni, Mihaela Rosca, Jiepu Jiang, Charlie Chen, RuiBo Liu, Tara Sainath, Maxim
884 Krikun, Alex Polozov, Jean-Baptiste Lespiau, Josh Newlan, Zeynep Cankara, Soo Kwak, Yun-
885 han Xu, Phil Chen, Andy Coenen, Clemens Meyer, Katerina Tsihlias, Ada Ma, Juraj Gottweis,
886 Jinwei Xing, Chenjie Gu, Jin Miao, Christian Frank, Zeynep Cankara, Sanjay Ganapathy, Ishita
887 Dasgupta, Steph Hughes-Fitt, Heng Chen, David Reid, Keran Rong, Hongmin Fan, Joost van
888 Amersfoort, Vincent Zhuang, Aaron Cohen, Shixiang Shane Gu, Anhad Mohanane, Anastasija
889 Ilic, Taylor Tobin, John Wieting, Anna Bortsova, Phoebe Thacker, Emma Wang, Emily Caveness,
890 Justin Chiu, Eren Sezener, Alex Kaskasoli, Steven Baker, Katie Millican, Mohamed Elhawaty,
891 Kostas Aisopos, Carl Lebsack, Nathan Byrd, Hanjun Dai, Wenhao Jia, Matthew Wiethoff, El-
892 naz Davoodi, Albert Weston, Lakshman Yagati, Arun Ahuja, Isabel Gao, Golan Pundak, Su-
893 san Zhang, Michael Azzam, Khe Chai Sim, Sergi Caelles, James Keeling, Abhanshu Sharma,
894 Andy Swing, YaGuang Li, Chenxi Liu, Carrie Grimes Bostock, Yamini Bansal, Zachary Nado,
895 Ankesh Anand, Josh Lipschultz, Abhijit Karmarkar, Lev Proleev, Abe Ittycheriah, Soheil Has-
896 sas Yeganeh, George Polovets, Aleksandra Faust, Jiao Sun, Alban Rrustemi, Pen Li, Rakesh
897 Shivanna, Jeremiah Liu, Chris Welty, Federico Lebron, Anirudh Baddepudi, Sebastian Krause,
898 Emilio Parisotto, Radu Soricut, Zheng Xu, Dawn Bloxwich, Melvin Johnson, Behnam Neyshabur,
899 Justin Mao-Jones, Renshen Wang, Vinay Ramasesh, Zaheer Abbas, Arthur Guez, Constant Segal,
900 Duc Dung Nguyen, James Svensson, Le Hou, Sarah York, Kieran Milan, Sophie Bridgers, Wiktor
901 Gworek, Marco Tagliasacchi, James Lee-Thorp, Michael Chang, Alexey Guseynov, Ale Jakse
902 Hartman, Michael Kwong, Ruizhe Zhao, Sheleem Kashem, Elizabeth Cole, Antoine Miech,
903 Richard Tanburn, Mary Phuong, Filip Pavetic, Sebastien Cevey, Ramona Comanescu, Richard
904 Ives, Sherry Yang, Cosmo Du, Bo Li, Zizhao Zhang, Mariko Iinuma, Clara Huiyi Hu, Aurko Roy,
905 Shaan Bijwadia, Zhenkai Zhu, Danilo Martins, Rachel Saputro, Anita Gergely, Steven Zheng,
906 Dawei Jia, Ioannis Antonoglou, Adam Sadovsky, Shane Gu, Yingying Bi, Alek Andreev, Sina
907 Samangoeei, Mina Khan, Tomas Kocisky, Angelos Filos, Chintu Kumar, Colton Bishop, Adams
908 Yu, Sarah Hodgkinson, Sid Mittal, Premal Shah, Alexandre Moufarek, Yong Cheng, Adam Blo-
909 niarz, Jaehoon Lee, Pedram Pejman, Paul Michel, Stephen Spencer, Vladimir Feinberg, Xuehan
910 Xiong, Nikolay Savinov, Charlotte Smith, Siamak Shakeri, Dustin Tran, Mary Chesus, Bernd
911 Bohnet, George Tucker, Tamara von Glehn, Carrie Muir, Yiran Mao, Hideto Kazawa, Ambrose
912 Slone, Kedar Soparkar, Disha Shrivastava, James Cobon-Kerr, Michael Sharman, Jay Pavagadhi,
913 Carlos Araya, Karolis Misiunas, Nimesh Ghelani, Michael Laskin, David Barker, Qijia Li, An-
914 ton Briukhov, Neil Houlsby, Mia Glaese, Balaji Lakshminarayanan, Nathan Schucher, Yunhao
915 Tang, Eli Collins, Hyeontaek Lim, Fangxiaoyu Feng, Adria Recasens, Guangda Lai, Alberto
916 Magni, Nicola De Cao, Aditya Siddhant, Zoe Ashwood, Jordi Orbay, Mostafa Dehghani, Jenny
917 Brennan, Yifan He, Kelvin Xu, Yang Gao, Carl Saroufim, James Molloy, Xinyi Wu, Seb Arnold,
Solomon Chang, Julian Schrittwieser, Elena Buchatskaya, Soroush Radpour, Martin Polacek,
Skye Giordano, Ankur Bapna, Simon Tokumine, Vincent Hellendoorn, Thibault Sottiaux, Sarah
Cogan, Aliaksei Severyn, Mohammad Saleh, Shantanu Thakoor, Laurent Shefey, Siyuan Qiao,
Meenu Gaba, Shuo yiin Chang, Craig Swanson, Biao Zhang, Benjamin Lee, Paul Kishan Ruben-
stein, Gan Song, Tom Kwiatkowski, Anna Koop, Ajay Kannan, David Kao, Parker Schuh, Axel

918 Stjerngren, Golnaz Ghiasi, Gena Gibson, Luke Vilnis, Ye Yuan, Felipe Tiengo Ferreira, Aish-
919 warya Kamath, Ted Klimenko, Ken Franko, Kefan Xiao, Indro Bhattacharya, Miteyan Patel, Rui
920 Wang, Alex Morris, Robin Strudel, Vivek Sharma, Peter Choy, Sayed Hadi Hashemi, Jessica
921 Landon, Mara Finkelstein, Priya Jhakra, Justin Frye, Megan Barnes, Matthew Mauger, Dennis
922 Daun, Khuslen Baatarsukh, Matthew Tung, Wael Farhan, Henryk Michalewski, Fabio Viola, Fe-
923 lix de Chaumont Quitry, Charline Le Lan, Tom Hudson, Qingze Wang, Felix Fischer, Ivy Zheng,
924 Elspeth White, Anca Dragan, Jean baptiste Alayrac, Eric Ni, Alexander Pritzel, Adam Iwan-
925 icki, Michael Isard, Anna Bulanova, Lukas Zilka, Ethan Dyer, Devendra Sachan, Srivatsan Srin-
926 vasan, Hannah Muckenhirn, Honglong Cai, Amol Mandhane, Mukarram Tariq, Jack W. Rae, Gary
927 Wang, Kareem Ayoub, Nicholas FitzGerald, Yao Zhao, Woohyun Han, Chris Alberti, Dan Gar-
928 rette, Kashyap Krishnakumar, Mai Gimenez, Anselm Levskaya, Daniel Sohn, Josip Matak, Inaki
929 Iturrate, Michael B. Chang, Jackie Xiang, Yuan Cao, Nishant Ranka, Geoff Brown, Adrian Hut-
930 ter, Vahab Mirrokni, Nanxin Chen, Kaisheng Yao, Zoltan Egyed, Francois Galilee, Tyler Liechty,
931 Praveen Kallakuri, Evan Palmer, Sanjay Ghemawat, Jasmine Liu, David Tao, Chloe Thornton,
932 Tim Green, Mimi Jasarevic, Sharon Lin, Victor Cotruta, Yi-Xuan Tan, Noah Fiedel, Hongkun
933 Yu, Ed Chi, Alexander Neitz, Jens Heitkaemper, Anu Sinha, Denny Zhou, Yi Sun, Charbel
934 Kaed, Brice Hulse, Swaroop Mishra, Maria Georgaki, Sneha Kudugunta, Clement Farabet, Izhak
935 Shafran, Daniel Vlasic, Anton Tsitsulin, Rajagopal Ananthanarayanan, Alen Carin, Guolong Su,
936 Pei Sun, Shashank V, Gabriel Carvajal, Josef Broder, Iulia Comsa, Alena Repina, William Wong,
937 Warren Weilun Chen, Peter Hawkins, Egor Filonov, Lucia Loher, Christoph Hirsenschall, Weiyi
938 Wang, Jingchen Ye, Andrea Burns, Hardie Cate, Diana Gage Wright, Federico Piccinini, Lei
939 Zhang, Chu-Cheng Lin, Ionel Gog, Yana Kulizhskaya, Ashwin Sreevatsa, Shuang Song, Luis C.
940 Cobo, Anand Iyer, Chetan Tekur, Guillermo Garrido, Zhuyun Xiao, Rupert Kemp, Huaixiu Steven
941 Zheng, Hui Li, Ananth Agarwal, Christel Ngani, Kati Goshvadi, Rebeca Santamaria-Fernandez,
942 Wojciech Fica, Xinyun Chen, Chris Gorgolewski, Sean Sun, Roopal Garg, Xinyu Ye, S. M. Ali
943 Eslami, Nan Hua, Jon Simon, Pratik Joshi, Yelin Kim, Ian Tenney, Sahitya Potluri, Lam Nguyen
944 Thiet, Quan Yuan, Florian Luisier, Alexandra Chronopoulou, Salvatore Scellato, Praveen Srin-
945 vasan, Minmin Chen, Vinod Koverkathu, Valentin Dalibard, Yaming Xu, Brennan Saeta, Keith
946 Anderson, Thibault Sellam, Nick Fernando, Fantine Huot, Junehyuk Jung, Mani Varadarajan,
947 Michael Quinn, Amit Raul, Maigo Le, Ruslan Habalov, Jon Clark, Komal Jalan, Kalesha Bullard,
948 Achintya Singhal, Thang Luong, Boyu Wang, Sujeevan Rajayogam, Julian Eisenschlos, Johnson
949 Jia, Daniel Finchelstein, Alex Yakubovich, Daniel Balle, Michael Fink, Sameer Agarwal, Jing Li,
950 Dj Dvijotham, Shalini Pal, Kai Kang, Jaclyn Konzelmann, Jennifer Beattie, Olivier Dousse, Di-
951 ane Wu, Remi Crocker, Chen Elkind, Siddhartha Reddy Jonnalagadda, Jong Lee, Dan Holtmann-
952 Rice, Krystal Kallarackal, Rosanne Liu, Denis Vnukov, Neera Vats, Luca Invernizzi, Mohsen
953 Jafari, Huanjie Zhou, Lilly Taylor, Jennifer Prendki, Marcus Wu, Tom Eccles, Tianqi Liu, Kavya
954 Kopparapu, Francoise Beaufays, Christof Angermueller, Andreea Marzoca, Shourya Sarcar, Hi-
955 lal Dib, Jeff Stanway, Frank Perbet, Nejc Trdin, Rachel Sterneck, Andrey Khorlin, Dinghua Li,
956 Xihui Wu, Sonam Goenka, David Madras, Sasha Goldshtein, Willi Gierke, Tong Zhou, Yaxin
957 Liu, Yannie Liang, Anais White, Yunjie Li, Shreya Singh, Sanaz Bahargam, Mark Epstein, Sujoy
958 Basu, Li Lao, Adnan Ozturel, Carl Crous, Alex Zhai, Han Lu, Zora Tung, Neeraj Gaur, Alanna
959 Walton, Lucas Dixon, Ming Zhang, Amir Globerson, Grant Uy, Andrew Bolt, Olivia Wiles, Mi-
960 lad Nasr, Iliia Shumailov, Marco Selvi, Francesco Piccinno, Ricardo Aguilar, Sara McCarthy,
961 Misha Khalman, Mrinal Shukla, Vlado Galic, John Carpenter, Kevin Vilella, Haibin Zhang,
962 Harry Richardson, James Martens, Matko Bosnjak, Shreyas Rammohan Belle, Jeff Seibert, Mah-
963 moud Alnahlawi, Brian McWilliams, Sankalp Singh, Annie Louis, Wen Ding, Dan Popovici,
964 Lenin Simicich, Laura Knight, Pulkit Mehta, Nishesh Gupta, Chongyang Shi, Saaber Fatehi, Jo-
965 vana Mitrovic, Alex Grills, Joseph Pagadora, Dessie Petrova, Danielle Eisenbud, Zhishuai Zhang,
966 Damion Yates, Bhavishya Mittal, Nilesh Tripuraneni, Yannis Assael, Thomas Brovelli, Prateek
967 Jain, Mihajlo Velimirovic, Canfer Akbulut, Jiaqi Mu, Wolfgang Macherey, Ravin Kumar, Jun Xu,
968 Haroon Qureshi, Gheorghe Comanici, Jeremy Wiesner, Zhitao Gong, Anton Ruddock, Matthias
969 Bauer, Nick Felt, Anirudh GP, Anurag Arnab, Dustin Zelle, Jonas Rothfuss, Bill Rosgen, Ashish
970 Shenoy, Bryan Seybold, Xinjian Li, Jayaram Mudigonda, Goker Erdogan, Jiawei Xia, Jiri Simsa,
971 Andrea Michi, Yi Yao, Christopher Yew, Steven Kan, Isaac Caswell, Carey Radebaugh, Andre
Elisseeff, Pedro Valenzuela, Kay McKinney, Kim Paterson, Albert Cui, Eri Latorre-Chimoto,
Solomon Kim, William Zeng, Ken Durden, Priya Ponnappalli, Tiberiu Sosea, Christopher A.
Choquette-Choo, James Manyika, Brona Robenek, Harsha Vashisht, Sebastien Pereira, Hoi Lam,
Marko Velic, Denese Owusu-Afriyie, Katherine Lee, Tolga Bolukbasi, Alicia Parrish, Shawn
Lu, Jane Park, Balaji Venkatraman, Alice Talbert, Lambert Rosique, Yuchung Cheng, Andrei

972 Sozanschi, Adam Paszke, Praveen Kumar, Jessica Austin, Lu Li, Khalid Salama, Wooyeol Kim,
973 Nandita Dukkkipati, Anthony Baryshnikov, Christos Kaplanis, XiangHai Sheng, Yuri Chervonyi,
974 Caglar Unlu, Diego de Las Casas, Harry Askham, Kathryn Tunyasuvunakool, Felix Gimeno, Siim
975 Poder, Chester Kwak, Matt Miecnikowski, Vahab Mirrokni, Alek Dimitriev, Aaron Parisi, Dan-
976 gyi Liu, Tomy Tsai, Toby Shevlane, Christina Kouridi, Drew Garmon, Adrian Goedeckemeyer,
977 Adam R. Brown, Anitha Vijayakumar, Ali Elqursh, Sadegh Jazayeri, Jin Huang, Sara Mc Carthy,
978 Jay Hoover, Lucy Kim, Sandeep Kumar, Wei Chen, Courtney Biles, Garrett Bingham, Evan
979 Rosen, Lisa Wang, Qijun Tan, David Engel, Francesco Pongetti, Dario de Cesare, Dongseong
980 Hwang, Lily Yu, Jennifer Pullman, Srini Narayanan, Kyle Levin, Siddharth Gopal, Megan Li,
981 Asaf Aharoni, Trieu Trinh, Jessica Lo, Norman Casagrande, Roopali Vij, Loic Matthey, Braman-
982 dia Ramadhana, Austin Matthews, CJ Carey, Matthew Johnson, Kremena Goranova, Rohin Shah,
983 Shereen Ashraf, Kingshuk Dasgupta, Rasmus Larsen, Yicheng Wang, Manish Reddy Vuyyuru,
984 Chong Jiang, Joana Ijazi, Kazuki Osawa, Celine Smith, Ramya Sree Boppana, Taylan Bilal, Yuma
985 Koizumi, Ying Xu, Yasemin Altun, Nir Shabat, Ben Bariach, Alex Korchemnyy, Kiam Choo, Olaf
986 Ronneberger, Chimezie Iwuanyanwu, Shubin Zhao, David Soergel, Cho-Jui Hsieh, Irene Cai,
987 Shariq Iqbal, Martin Sundermeyer, Zhe Chen, Elie Bursztein, Chaitanya Malaviya, Fadi Biadys,
988 Prakash Shroff, Inderjit Dhillon, Tejasi Latkar, Chris Dyer, Hannah Forbes, Massimo Nicosia,
989 Vitaly Nikolaev, Somer Greene, Marin Georgiev, Pidong Wang, Nina Martin, Hanie Sedghi, John
990 Zhang, Praseem Banzal, Doug Fritz, Vikram Rao, Xuezhong Wang, Jiageng Zhang, Viorica Pa-
991 traucean, Dayou Du, Igor Mordatch, Ivan Jurin, Lewis Liu, Ayush Dubey, Abhi Mohan, Janek
992 Nowakowski, Vlad-Doru Ion, Nan Wei, Reiko Tojo, Maria Abi Raad, Drew A. Hudson, Vaishakh
993 Keshava, Shubham Agrawal, Kevin Ramirez, Zhichun Wu, Hoang Nguyen, Ji Liu, Madhavi Se-
994 wak, Bryce Petrini, DongHyun Choi, Ivan Philips, Ziyue Wang, Ioana Bica, Ankush Garg, Jarek
995 Wilkiewicz, Priyanka Agrawal, Xiaowei Li, Danhao Guo, Emily Xue, Naseer Shaik, Andrew
996 Leach, Sath MNM Khan, Julia Wiesinger, Sammy Jerome, Abhishek Chakladar, Alek Wenjiao
997 Wang, Tina Ornduff, Folake Abu, Alireza Ghaffarkhah, Marcus Wainwright, Mario Cortes, Fred-
998 erick Liu, Joshua Maynez, Andreas Terzis, Pouya Samangouei, Riham Mansour, Tomasz Kepa,
999 François-Xavier Aubet, Anton Algymr, Dan Banica, Agoston Weisz, Andras Orban, Alexandre
1000 Senges, Ewa Andrejczuk, Mark Geller, Niccolo Dal Santo, Valentin Anklin, Majd Al Mery,
1001 Martin Baeuml, Trevor Strohm, Junwen Bai, Slav Petrov, Yonghui Wu, Demis Hassabis, Koray
1002 Kavukcuoglu, Jeffrey Dean, and Oriol Vinyals. Gemini 1.5: Unlocking multimodal understanding
1003 across millions of tokens of context, 2024. URL <https://arxiv.org/abs/2403.05530>.

1004 Philippe Tillet, H. T. Kung, and David Cox. Triton: an intermediate language and compiler for
1005 tiled neural network computations. In *Proceedings of the 3rd ACM SIGPLAN International
1006 Workshop on Machine Learning and Programming Languages*, MAPL 2019, pp. 10–19, New
1007 York, NY, USA, 2019. Association for Computing Machinery. ISBN 9781450367196. doi:
1008 10.1145/3315508.3329973. URL <https://doi.org/10.1145/3315508.3329973>.

1009 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
1010 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Ar-
1011 mand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation
1012 language models, 2023. URL <https://arxiv.org/abs/2302.13971>.

1013 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez,
1014 Lukasz Kaiser, and Illia Polosukhin. Attention is all you need, 2023. URL <https://arxiv.org/abs/1706.03762>.

1015 Guanhua Wang, Chengming Zhang, Zheyu Shen, Ang Li, and Olatunji Ruwase. Domino: Elim-
1016 inating communication in llm training via generic tensor slicing and overlapping, 2024a. URL
1017 <https://arxiv.org/abs/2409.15241>.

1018 Hongyu Wang, Shuming Ma, Li Dong, Shaohan Huang, Dongdong Zhang, and Furu Wei. Deep-
1019 net: Scaling transformers to 1,000 layers, 2022. URL <https://arxiv.org/abs/2203.00555>.

1020
1021 Lean Wang, Huazuo Gao, Chenggang Zhao, Xu Sun, and Damai Dai. Auxiliary-loss-free load
1022 balancing strategy for mixture-of-experts, 2024b. URL <https://arxiv.org/abs/2408.15664>.

1023
1024 xAI. Grok-1, 2024. URL [https://github.com/xai-org/grok-1?tab=](https://github.com/xai-org/grok-1?tab=readme-ov-file)
1025 [readme-ov-file](https://github.com/xai-org/grok-1?tab=readme-ov-file).

1026 Yanqi Zhou, Tao Lei, Hanxiao Liu, Nan Du, Yanping Huang, Vincent Zhao, Andrew Dai, Zhifeng
1027 Chen, Quoc Le, and James Laudon. Mixture-of-experts with expert choice routing, 2022. URL
1028 <https://arxiv.org/abs/2202.09368>.
1029
1030 Barret Zoph, Irwan Bello, Sameer Kumar, Nan Du, Yanping Huang, Jeff Dean, Noam Shazeer, and
1031 William Fedus. St-moe: Designing stable and transferable sparse expert models, 2022. URL
1032 <https://arxiv.org/abs/2202.08906>.
1033
1034
1035
1036
1037
1038
1039
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

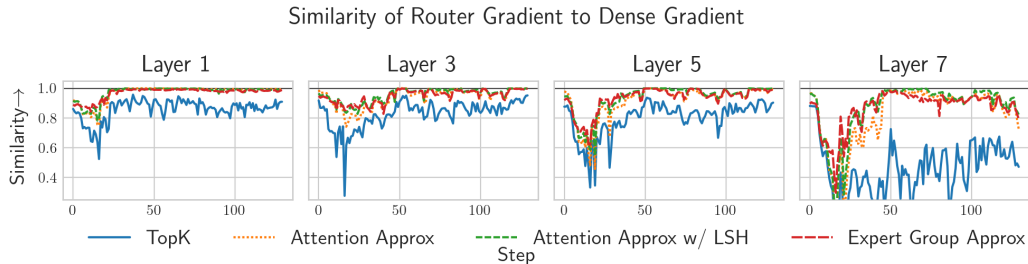


Figure 10: **Accuracy in approximating the dense router gradient for each approach.** We use a model with 8 experts and $K = 2$. We artificially compute the *dense gradient* of the output with respect to the router weights at each step by passing inputs through a dense mixture of experts layer, where all experts are selected. We do this independent of the actual forward pass computation, while using the same set of MoE parameters. The similarity between this dense gradient and the actual gradient propagated to the router indicates how well the router is learning from all experts. We plot this similarity with a standard Top-K router, and with each of our proposed router modifications. Our approaches are much more accurate and stable in approximating the dense router gradient.

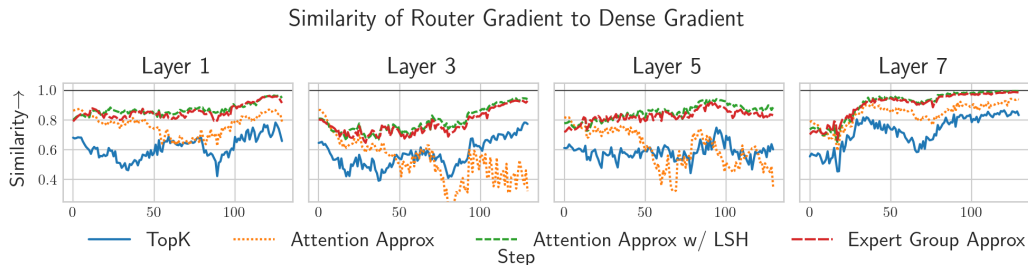


Figure 11: **Dense router gradient approximation accuracy with fine-grained experts.** We implement fine-grained experts as in DeepSeekMoE (DeepSeek-AI et al., 2024) to observe the behavior of our approximation methods across more experts while keeping parameter count fixed. In this example, the model now has 32 experts with $K = 2$. With more experts, it becomes increasingly hard to approximate the dense gradient, and the difference between our methods and the Top-K router is more apparent. Moreover, we can clearly compare the efficacy of each method and see that the attention approximation with LSH is the best. Note the average number of tokens per expert also decreases by a factor of 4 as well, and we would expect even better performance in our approximations by scaling the local batch size per GPU, because we don’t communicate tokens across ranks.

A APPROXIMATION STATISTICS

A.1 FURTHER GRADIENT ANALYSIS

The differences in our approaches become clear as we scale the model to become more sparse. We expand to $N = 32$ experts while maintaining $K = 2$ in Fig. 11 and find that it is more difficult to approximate the true dense router gradient. While all of our approaches sufficiently approximate the dense gradient with $N = 8$ experts, the performance gap between them is apparent with $N = 32$. The expert group approximation and LSH attention methods are significantly better than the direct attention method, and this is also consistent with our validation results in Table 6. This is likely due to the heuristics we apply to restrict our approximation to only the most relevant tokens: the expert group approximation requires tokens to have an expert in common, and LSH requires tokens to be similar. Moreover, the gap between our methods and Top-K is wider with 32 experts. We believe that in larger models with even more experts, our method will yield increasingly significant improvements over Top-K routing. We provide further analysis on our approximations in Appendix A. In Fig. 12 we provide an additional analysis of the gradient norm of our approximation compared to the dense gradient. This logging is also done with $N = 8$ and $K = 2$. The Top-K gradient has

Table 5: Our expert group approximation obtains the best validation perplexity after 20B tokens on a standard MoE architecture, achieving the same performance as $K = 3$ without activating an additional expert, meaning that we achieve the best performance-FLOPS tradeoff.

Activated Experts	Routing Method	Validation Perplexity
$K = 1$	Baseline	19.61
$K = 2$	Baseline	18.92
$K = 3$	Baseline	18.56
$K = 2$	Expert Group Approx. (Ours)	18.55

significantly lower norm than the dense gradient. Our methods closely approximate the dense gradient norm consistently; replicating both the direction and magnitude suggests that we are sufficiently approximating the dense gradient entirely.

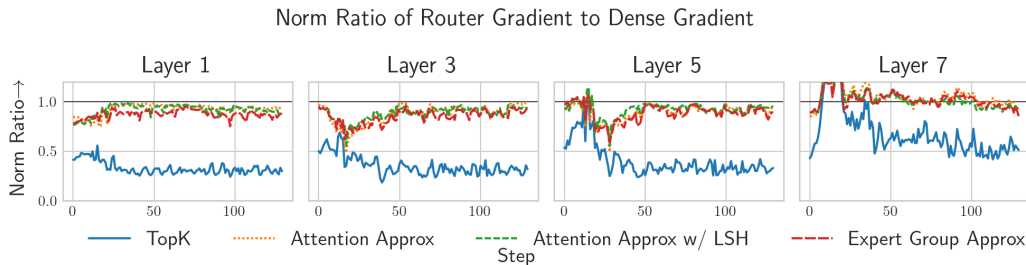


Figure 12: **Comparison of gradient norms relative to the dense gradient.** When computing the dense gradient, we also record its L_2 norm and log the ratio of this to the L_2 norms of the actual router gradients during training. Our methods produce router gradients with approximately the same magnitude. Along with the results showing strong cosine similarity, this suggests that we are almost perfectly approximating the dense gradient.

In Fig. 5 we presented results with a Deepseek-style finegrained MoE. In Table 5 we present results with a conventional architecture. This model has 2B hidden parameters and has 780M active parameters. Our method outperforms the baseline. As we will show, our method improves performance by improving load balance. Load balancing is especially important early on in training – this is where our method shows the largest improvement in Fig. 5. Improving the load balance of experts yields benefits akin to actually activating more experts. In Table 5 we find that our lightweight approximation method improves performance by a similar amount as activating an additional expert (that is, going from $K = 2$ to $K = 3$), without the additional computational overhead during training and inference of actually needing to use the parameters of a third expert. The choice of $K = 2$ in all experiments follows Zoph et al. (2022).

B HYPERPARAMETERS AND IMPLEMENTATION DETAILS

Model Configuration. Both MoEs have 24 blocks and a hidden dimension of 1024. The Deepseek MoE has fine-grained experts. Each expert is a bottleneck MLP of shape (1024, 704). The conventional MoE has an expansion factor such that the intermediate size of the MLP is 2816. We use SwiGLU (Shazeer, 2020) MLPs following Llama (Touvron et al., 2023), 16 attention heads with dimension of 64, LayerNorm (Ba et al., 2016) and RoPE (Su et al., 2023). **Hyperparameters.** We use the initialization from Wang et al. (2022) for residual branch merge layers and the initialization from Nguyen & Salazar (2019) for all other layers. We use the AdamW optimizer (Loshchilov & Hutter, 2019). We use the modified cosine learning rate schedule from Ibrahim et al. (2024). We use a sequence length of 2048 and a global batch size of 1024, resulting in a global token batch size of 2^{21} . We set the auxiliary loss (Fedus et al., 2022) to 0.01.

Implementation. We train with the gpt-neox library (Andonian et al., 2023) integrated with Megablocks (Gale et al., 2022). The TFLOPS vary depending on the method and the number of

1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295

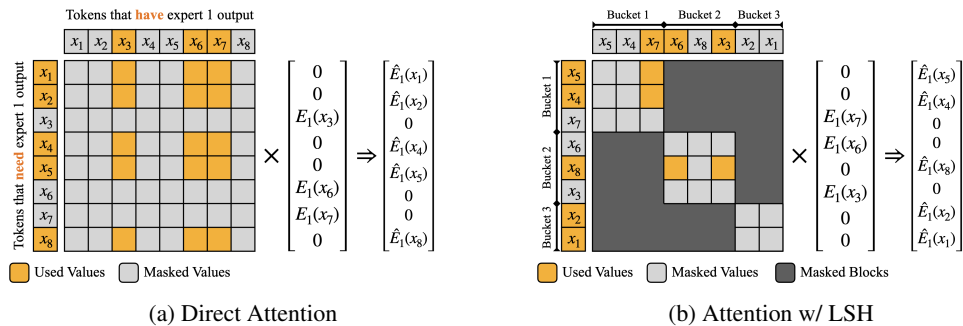


Figure 14: **Attention scores of direct and LSH attention methods.** For each expert, we define an attention head that uses queries corresponding to inputs not routed to the expert, and keys corresponding to inputs routed to the expert. Grey entries denote queries and keys that do not meet this criteria, and whose attention scores are masked out. We multiply the attention scores of each head with the values, which are expert outputs of tokens routed to that attention head’s expert. This implementation is common to both the direct and LSH attention method. In the latter, we further optimize the attention calculation by sorting inputs into buckets based on cosine similarity. This creates a block-sparse attention map, allowing kernels to skip most of the attention computation.

attention implementations to higher dimensions could make this token-wise approximation method more feasible.

C.1 ABLATING DESIGN CHOICES

We ablate our main design choices and show that our main method does not compromise throughput.

Comparing Different Approximation Methods. We use the Expert Group Approximation method for our main results because it is lightweight, easy to implement, and provides good performance. However, the other two methods we consider also outperform the top-K ($K = 2$) baseline. Indeed, as we showed in Fig. 11, the Attention+LSH method seems to obtain a better approximation of the true dense gradient. The primary reason why we report our main results with Expert Grouping is because the Expert Group Approximation method requires no additional memory overhead. This allows us to use larger microbatches, and therefore there are more tokens on each GPU that we can use for the approximation. In Table 6 we find that even with a microbatchsize $4\times$ smaller than that of the Expert Group method, the Attention+LSH method is competitive.

Routing Method	Microbatchsize	Validation Perplexity
Attention	4	18.72
Attention+LSH	4	18.64
Expert Group	16	18.55
Baseline	16	18.92

Table 6: Comparison of routing methods after training on 20B tokens.