
Dense Backpropagation Improves Routing for Sparsely-Gated Mixture-of-Experts

Anonymous Author(s)
Affiliation
Address
email

Abstract

1 Sparsely-gated Mixture-of-Experts (MoEs) have proven to be more effi-
2 cient than dense Transformers because they can dynamically activate a
3 subset of their overall parameters by *routing* tokens to selected “experts”,
4 allowing practitioners to scale up model parameter counts without sig-
5 nificantly increasing total compute. However, current MoE training ap-
6 proaches only update the router with a sparse gradient and suffer from
7 issues such as load imbalance. We propose a new router that can receive
8 a dense gradient update from a sparse forward pass. Our method adds
9 minimal overhead, but improves on the common Top-K routing in both
10 performance and load balance.

11 1 Introduction

12 Large-scale pretraining hinges on scaling up the number of parameters in a model, because
13 models with more parameters are more sample-efficient and require less training to reach
14 the same performance as smaller models [Kaplan et al., 2020, Hoffmann et al., 2022]. Most
15 academic research has adopted the dense Transformer architecture [Vaswani et al., 2023]
16 because its performance scales well with parameters and data. However, a sparsely ac-
17 tivated Mixture-of-Experts (MoE) Transformer architecture [Shazeer et al., 2017] has been
18 used by many industry deployments [Team et al., 2024, xAI, 2024, Databricks, 2024, Jiang
19 et al., 2024, Snowflake, 2024, DeepSeek-AI et al., 2024] because MoEs have been shown to
20 scale better than dense Transformers [Clark et al., 2022, Du et al., 2022, Lepikhin et al., 2020,
21 Fedus et al., 2022]. This is not unique to dense Transformers; MoEs based on state-space
22 modules rather than dense Transformer blocks outperform networks of only state-space
23 modules [Lieber et al., 2024, Liquid, 2024]. MoEs learn a *routing* function that selectively
24 activates the TopK subset of their parallel MLP modules, or *experts*, most relevant to a
25 given input. This conditionally sparse activation [Jacobs et al., 1991, Jordan and Jacobs,
26 1994] allows the model parameter count to be increased multiplicatively without signifi-
27 cantly increasing the cost of training or inference.

28 However, this same router that unlocks superior scaling also presents a challenge for MoEs,
29 because the router does not receive a gradient update from experts that it does not activate,
30 and may not learn to route a token to its appropriate expert. One critical issue is the load
31 imbalance problem — where a few experts are over-utilized — which leads to inefficient
32 training and resource usage [Zoph et al., 2022, Zhou et al., 2022]. We propose a new router
33 that can receive a dense gradient update from a sparse forward pass to address the instabil-
34 ity issues arising from sparse routing. Our method adds minimal overhead, but improves
35 on the common Top-K routing in both performance and load balance.

36 2 Background & Related Work

37 **MoEs.** The MoE layer replaces the feedforward networks (FFN) of transformers and
 38 consists of two components : **1)** N FFNs (*experts*), $E_0(x), E_1(x), \dots, E_N(x)$ and **2)** a router
 39 that assigns tokens to experts. Each input to the MoE layer is processed by K experts where
 40 $K < N$ and is thus the source of sparsity in MoEs. The K experts are chosen by the router,
 41 which is a learnable component that maps each token to a set of weights over the experts.
 42 The router performs a linear transformation $\mathbb{R}^{d_{\text{token}}} \rightarrow \mathbb{R}^N$ which produces logits; these are
 43 normalized using softmax, resulting in a probability distribution over the experts. With the
 44 router’s linear transformation parameterized by a matrix W , we can represent the expert
 45 weights π in the following way:

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

46 Once we have these expert weights, we apply a routing function to decide which of K
 47 experts to route and process this token through.

48 **Top-K routing.** A standard method to select K out of N experts given the expert weights
 49 is to select the experts corresponding to the K highest weights. Top-K routing [Fedus et al.,
 50 2022] passes the token to the K selected experts and averages the expert outputs using these
 51 weights to produce the final output. Experts not selected by the Top-K routing function do
 52 not process the token, and this introduces sparsity in MoEs. By representing the K chosen
 53 experts as the set A , we can express the output of the MoE layer as:

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

54 Thus, the expert weights have a dual purpose : They are used by the routing function to
 55 decide which of the K experts to process a token through, and also provide the weights for
 56 combining the outputs of the expert.

57 The Top-K routing scheme makes the MoE layer desirable for training large, compute-
 58 efficient neural networks. It allows models to be scaled up, by way of increasing the total
 59 number of experts, while keeping the compute per token constant (as it is a function of K
 60 and not N).

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

$$\frac{\partial y}{\partial W} = \frac{\partial y}{\partial \pi} \frac{\partial \pi}{\partial W} \quad (3)$$

65 The steps in Eq. (1) are easily differentiable as they consist of linear operations and acti-
 66 vations. Thus, the first term in Eq. (3), $\frac{\partial y}{\partial \pi}$, is straightforward to compute. Eq. (2), how-
 67 ever, isn’t differentiable because the Top-K expert selection is a discrete function: given the
 68 continuous router weights $\pi \in \mathbb{R}^N$, the set of selected experts A is one of $\binom{N}{K}$ combina-
 69 tions. One way to get around backpropagation of nondifferentiable operations is to use the
 70 straight-through estimator [Bengio et al., 2013], which treats the operator as the identity
 71 function. In this setting, the Top-K routing function is bypassed and Eq. (2) becomes the
 72 dot product between π and the vector of all $E_i(x)$ with the following gradient:

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

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

77 **Related Works.** Previous work has tried to address the issue of routing in MoEs. Separate
 78 from Top-K is the Sinkhorn routing method [Clark et al., 2022]. Fedus et al. [2022] which

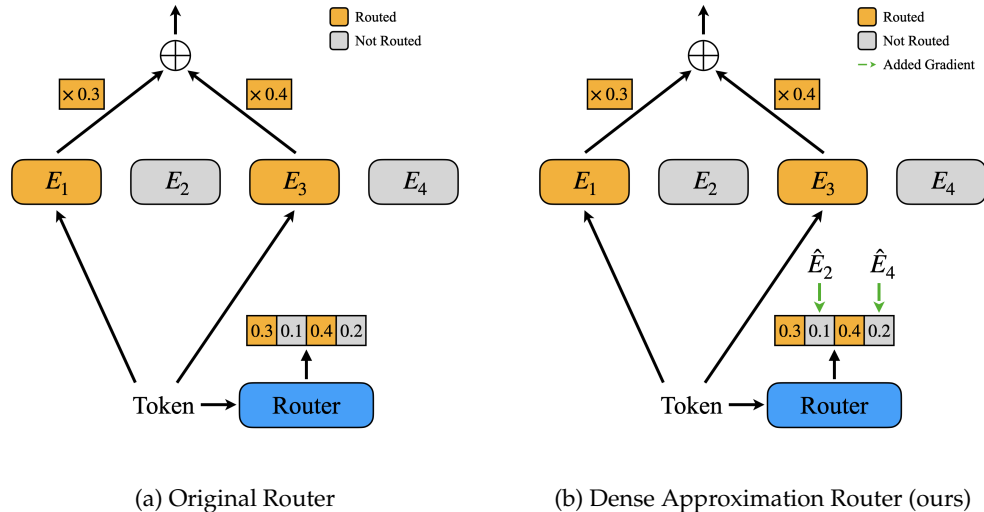


Figure 1: **Overview of Routing with Dense Approximations.** The original mixture of experts router only receives gradients corresponding to experts the token is routed to, because there is no output from other experts. Our approach provides the router with a complete (dense) gradient, by approximating 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.

79 proposes an auxiliary loss that encourages load balancing. Dai et al. [2024] propose mul-
 80 tiple additional auxiliary loss terms. Recently, Wang et al. [2024] propose learning biases
 81 rather than an auxiliary load balancing loss. Even more recently, Phi-3.5-MoE Abdin et al.
 82 [2024] uses SparseMixer [Liu et al., 2024, 2023], another estimator for $\partial y / \partial \pi$ not involv-
 83 ing straight-through. Our approach is to still use straight-through, but *approximate* these
 84 additional expert outputs.

85 3 Designing a New Routing Method

86 In this section we design a new router that can receive a dense gradient update while
 87 being sparsely activated. In a standard MoE, the embedding corresponding to expert i in
 88 the routing layer (i.e. the i th row of the routing weight matrix) receives no gradient update
 89 from a token x if x is not routed to expert i . This is because $E_i(x)$ is never computed, so
 90 it provides no upstream gradient. This corresponds to experts that are not in the top K
 91 being omitted in Eq. (2). We apply an approximation $\hat{E}_i(x)$ as a substitute for the upstream
 92 gradient, so that the router can receive some non-zero signal corresponding to this expert.
 93 Thus, the router can factor in outputs from all experts when learning to route each token.

94 3.1 Approximating Expert Activations

95 To approximate the dense gradient in Eq. (4), we must approximate $E_i(x)$ for every expert
 96 i that a token x was not passed to. Although we have no information about what the
 97 function E_i looks like for x , when training with large token batch sizes it is very likely that
 98 we have outputs of E_i for many other tokens. We develop two general approaches to to
 99 develop an estimator $\hat{E}_i(x)$, using the expert outputs of other relevant tokens. *Expert group*
 100 *approximation:* We first apply a single approximation to a large group of tokens that were
 101 not routed to expert i . This is efficient, but it may not necessarily be a viable approximation
 102 for any specific x . However, we hypothesize that this is a good estimator for the expert
 103 output across the entire batch - this is sufficient as we will only need an approximation for
 104 the batch gradient to update the router. *Attention approximation:* Our secondary approach
 105 produces an expert output approximation specifically for each token (see Appendix A).

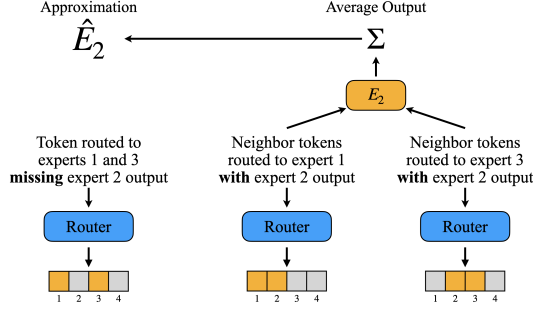


Figure 2: **Architecture of the Expert Group Approximation method.** In this example, we have 4 experts with $K = 2$. Consider all inputs routed to experts 1 and 3, characterized by the routing decision $R = \{1, 3\}$. As described in Figure 1b, we need to approximate these inputs’ activations for all 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.

106 3.2 Notation on Expert Routing

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

116 3.3 Expert Group Approximation

117 We primarily consider the case where we approximate the expert output $E_i(x)$ for many
 118 tokens at a time. For a token x , we want to approximate outputs of experts that x was
 119 not routed to, i.e. $E_i(x)$ where $i \notin R(x)$. We hypothesize that tokens being routed to the
 120 same expert is a strong indicator of similarity between the tokens. This is supported by
 121 our empirical observations in Appendix B.2. We develop an approximation for $E_i(x)$ by
 122 aggregating outputs of E_i for tokens that were routed to both expert i and an expert x was
 123 routed to. Formally, we consider an alternate routing decision $R' = \{i, j, \cdot\}, j \in R(x)$ that
 124 consists of one expert x is routed to, the expert i we wish to approximate, and any other
 125 experts (if $K > 2$). Then, the adjacent token space $X_{R'}$ will consist of tokens that are very
 126 similar to x by virtue of having similar routing decisions (see Fig. 9). Moreover, they will be
 127 routed to expert i , and we hypothesize that their outputs $\sum_{x' \in X_{R'}} E_i(x')$ will approximately
 128 represent $E_i(x)$. We can aggregate such outputs over all possible routing decisions:

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

129 We apply a single aggregate approximation for each routing decision to all tokens with that
 130 routing decision. Note that we only compute N^2 individual sums as that is the number of
 131 possible combinations $\{i, j, \cdot\}$. In Fig. 2 we visualize this method for $K = 2$.

132 **4 Evaluation**

133 **4.1 Evaluation**

134 **Main Result.** Our main result compares the Expert Group Approximation, which performs
 135 a dense update of the router weights by approximating the dense gradient, to baseline Top-
 136 K routing. Details on model training are provided in Appendix C. In Table 1 we find that
 137 our lightweight approximation method improves performance by a similar amount as ac-
 138 tivating an additional expert (that is, going from $K = 2$ to $K = 3$), without the additional
 139 computational overhead during training and inference of actually needing to use the pa-
 140 rameters of a third expert. The choice of $K = 2$ follows Zoph et al. [2022].

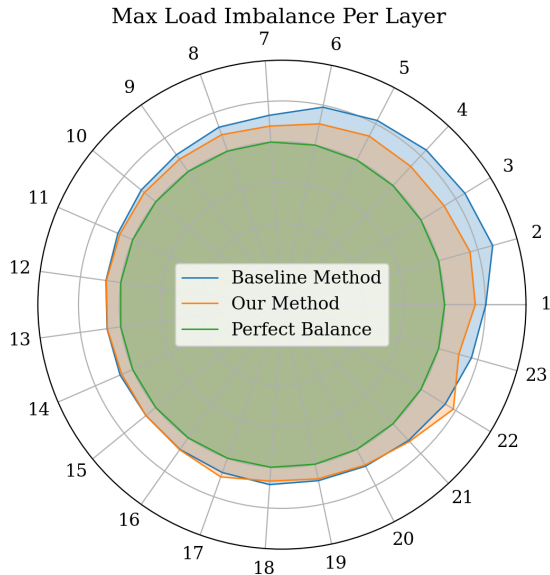


Figure 3: **Load Balancing using Expert Group Approximation.** We define maximum load imbalance at each layer as $\max_i(N \cdot f_i)$ where f_i is the fraction of tokens routed to each of N experts. The green ring indicates perfect balance, where each expert receives $1/N$ fraction of tokens; the outer ring indicates a maximum imbalance of 1.5. We record maximum imbalance after training on 3 billion tokens. By sending a complete gradient signal to the router, our model has better distribution of load than the baseline.

141 **Load Balance.** Our method improves over the baseline in perplexity, but the reason for
 142 this is in how it improves the routing distribution. Without gradient signal for unactivated
 143 experts, top-K routing may not be able to learn a balanced distribution across experts. This
 144 would lead to many more tokens being routed to some experts than others. In Fig. 3 we
 145 validate that the baseline top-K ($K = 2$) routing has an “imbalanced load”, as measured
 146 by the proportion of tokens being routed to different experts (labeled by color) relative to
 147 the baseline (dotted red line) of an even distribution of tokens across experts. Our method

Table 1: Our expert group approximation obtains the best validation perplexity after 20B tokens, achieving the same performance as $K = 3$ without activating an additional expert.

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

148 improves load balance, which may be one cause for improved performance and is of inde-
149 pendent interest on its own because it will lead to greater efficiency during inference.

150 **Ablations.** We conduct further ablations on design choices and efficiency in Ap-
151 pendix C.1.

152 5 Discussion

153 We propose a training method for MoEs to improve load balancing and language mod-
154 eling performance. By approximating the signal of a dense mixture-of-experts layer, the
155 MoE router is able to learn a better distribution of routing inputs to different experts. This
156 approximated dense signal unlocks the possibility for more sparse MoEs at training and
157 inference time. Whereas typical TopK routing would provide too sparse of a signal to learn
158 a stable routing distribution, our method demonstrates significant improvements in load
159 balancing and perplexity in very sparse configurations.

160 **Limitations.** The scope of our evaluation is limited; we only train models for at most 20B
161 tokens, and the largest MoE we train has fewer than 1B active parameters. Furthermore,
162 we only report the validation perplexity on a held-out subset of the training dataset and do
163 not report any benchmark scores. The scope of problems caused by routers includes load
164 balance and inability to handle distribution shifts during finetuning, but we only analyze
165 the impact of our method on load balance and do not know whether it actually makes it
166 easier to finetune MoEs. We plan to address these limitations in a future version of this
167 work.

168 **Future Work.** Our methods are somewhat unique in that they scale with the token batch
169 size per GPU, and improvements in memory efficiency therefore are critical. Developing
170 and integrating kernels to reduce the memory requirements of the MoE itself will allow us
171 to use larger microbatches. Another avenue for future work is developing entirely custom
172 kernels using our methods in order to reduce the computational overhead of approximat-
173 ing the dense router gradient.

174 **References**

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

178 Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai,
179 Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan
180 Clark, Tom Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan
181 Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol
182 Vinyals, and Laurent Sifre. Training compute-optimal large language models, 2022. URL
183 <https://arxiv.org/abs/2203.15556>.

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

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

190 Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Gar-
191 rett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, Soroosh Mariooryad, Yifan
192 Ding, Xinyang Geng, Fred Alcober, Roy Frostig, Mark Omernick, Lexi Walker, Cos-
193 min Paduraru, Christina Sorokin, Andrea Tacchetti, Colin Gaffney, Samira Daruki, Ol-
194 can Sercinoglu, Zach Gleicher, Juliette Love, Paul Voigtlaender, Rohan Jain, Gabriela
195 Surita, Kareem Mohamed, Rory Blevins, Junwhan Ahn, Tao Zhu, Kornraphop Kaw-
196 intiranon, Orhan Firat, Yiming Gu, Yujing Zhang, Matthew Rahtz, Manaal Faruqui,
197 Natalie Clay, Justin Gilmer, JD Co-Reyes, Ivo Penchev, Rui Zhu, Nobuyuki Morioka,
198 Kevin Hui, Krishna Haridasan, Victor Campos, Mahdis Mahdieh, Mandy Guo, Samer
199 Hassan, Kevin Kilgour, Arpi Vezer, Heng-Tze Cheng, Raoul de Liedekerke, Siddharth
200 Goyal, Paul Barham, DJ Strouse, Seb Noury, Jonas Adler, Mukund Sundararajan, Sharad
201 Vikram, Dmitry Lepikhin, Michela Paganini, Xavier Garcia, Fan Yang, Dasha Valter,
202 Maja Trebacz, Kiran Vodrahalli, Chulayuth Asawaroengchai, Roman Ring, Norbert Kalb,
203 Livio Baldini Soares, Siddhartha Brahma, David Steiner, Tianhe Yu, Fabian Mentzer,
204 Antoine He, Lucas Gonzalez, Bibo Xu, Raphael Lopez Kaufman, Laurent El Shafey,
205 Junhyuk Oh, Tom Hennigan, George van den Driessche, Seth Odoom, Mario Lucic,
206 Becca Roelofs, Sid Lall, Amit Marathe, Betty Chan, Santiago Ontanon, Luheng He, Den-
207 nis Teplyashin, Jonathan Lai, Phil Crone, Bogdan Damoc, Lewis Ho, Sebastian Riedel,
208 Karel Lenc, Chih-Kuan Yeh, Aakanksha Chowdhery, Yang Xu, Mehran Kazemi, Ehsan
209 Amid, Anastasia Petrushkina, Kevin Swersky, Ali Khodaei, Gowoon Chen, Chris Larkin,
210 Mario Pinto, Geng Yan, Adria Puigdomenech Badia, Piyush Patil, Steven Hansen, Dave
211 Orr, Sebastien M. R. Arnold, Jordan Grimstad, Andrew Dai, Sholto Douglas, Rishika
212 Sinha, Vikas Yadav, Xi Chen, Elena Gribovskaya, Jacob Austin, Jeffrey Zhao, Kaushal Pa-
213 tel, Paul Komarek, Sophia Austin, Sebastian Borgeaud, Linda Friso, Abhimanyu Goyal,
214 Ben Caine, Kris Cao, Da-Woon Chung, Matthew Lamm, Gabe Barth-Maron, Thais Kago-
215 hara, Kate Olszewska, Mia Chen, Kaushik Shivakumar, Rishabh Agarwal, Harshal God-
216 hia, Ravi Rajwar, Javier Snaider, Xerxes Dotiwalla, Yuan Liu, Aditya Barua, Victor Un-
217 gureanu, Yuan Zhang, Bat-Orgil Batsaikhan, Mateo Wirth, James Qin, Ivo Danihelka,
218 Tulsee Doshi, Martin Chadwick, Jilin Chen, Sanil Jain, Quoc Le, Arjun Kar, Madhu
219 Gurumurthy, Cheng Li, Ruoxin Sang, Fangyu Liu, Lampros Lamprou, Rich Munoz,
220 Nathan Lintz, Harsh Mehta, Heidi Howard, Malcolm Reynolds, Lora Aroyo, Quan
221 Wang, Lorenzo Blanco, Albin Cassirer, Jordan Griffith, Dipanjan Das, Stephan Lee,
222 Jakub Sygnowski, Zach Fisher, James Besley, Richard Powell, Zafarali Ahmed, Dominik
223 Paulus, David Reitter, Zalan Borsos, Rishabh Joshi, Aedan Pope, Steven Hand, Vittorio
224 Selo, Vihan Jain, Nikhil Sethi, Megha Goel, Takaki Makino, Rhys May, Zhen Yang,
225 Johan Schalkwyk, Christina Butterfield, Anja Hauth, Alex Goldin, Will Hawkins, Evan
226 Senter, Sergey Brin, Oliver Woodman, Marvin Ritter, Eric Noland, Minh Giang, Vijay
227 Bolina, Lisa Lee, Tim Blyth, Ian Mackinnon, Machel Reid, Obaid Sarvana, David Sil-
228 ver, Alexander Chen, Lily Wang, Loren Maggiore, Oscar Chang, Nithya Attaluri, Gre-
229 gory Thornton, Chung-Cheng Chiu, Oskar Bunyan, Nir Levine, Timothy Chung, Ev-

230 genii Eltyshev, Xiance Si, Timothy Lillicrap, Demetra Brady, Vaibhav Aggarwal, Boxi
231 Wu, Yuanzhong Xu, Ross McIlroy, Kartikeya Badola, Paramjit Sandhu, Erica Moreira,
232 Wojciech Stokowiec, Ross Hemsley, Dong Li, Alex Tudor, Pranav Shyam, Elahe Rahim-
233 toroghi, Salem Haykal, Pablo Sprechmann, Xiang Zhou, Diana Mincu, Yujia Li, Ravi
234 Addanki, Kalpesh Krishna, Xiao Wu, Alexandre Frechette, Matan Eyal, Allan Dafoe,
235 Dave Lacey, Jay Whang, Thi Avrahami, Ye Zhang, Emanuel Taropa, Hanzhao Lin,
236 Daniel Toyama, Eliza Rutherford, Motoki Sano, HyunJeong Choe, Alex Tomala, Chal-
237 lence Safranek-Shrader, Nora Kassner, Mantas Pajarskas, Matt Harvey, Sean Sechrist,
238 Meire Fortunato, Christina Lyu, Gamaleldin Elsayed, Chenkai Kuang, James Lottes,
239 Eric Chu, Chao Jia, Chih-Wei Chen, Peter Humphreys, Kate Baumli, Connie Tao, Ra-
240 jkumar Samuel, Cicero Nogueira dos Santos, Anders Andreassen, Nemanja Rakićević,
241 Dominik Grewe, Aviral Kumar, Stephanie Winkler, Jonathan Caton, Andrew Brock, Sid
242 Dalmia, Hannah Sheahan, Iain Barr, Yingjie Miao, Paul Natsev, Jacob Devlin, Feryal
243 Behbahani, Flavien Prost, Yanhua Sun, Artiom Myaskovsky, Thanumalayan Sankara-
244 narayana Pillai, Dan Hurt, Angeliki Lazaridou, Xi Xiong, Ce Zheng, Fabio Pardo, Xi-
245 aowei Li, Dan Horgan, Joe Stanton, Moran Ambar, Fei Xia, Alejandro Lince, Mingqiu
246 Wang, Basil Mustafa, Albert Webson, Hyo Lee, Rohan Anil, Martin Wicke, Timothy
247 Dozat, Abhishek Sinha, Enrique Piqueras, Elahe Dabir, Shyam Upadhyay, Anudhyan
248 Boral, Lisa Anne Hendricks, Corey Fry, Josip Djolonga, Yi Su, Jake Walker, Jane La-
249 banowski, Ronny Huang, Vedant Misra, Jeremy Chen, RJ Skerry-Ryan, Avi Singh, Shruti
250 Rijhwani, Dian Yu, Alex Castro-Ros, Beer Changpinyo, Romina Datta, Sumit Bagri,
251 Arnar Mar Hrafnkelsson, Marcello Maggioni, Daniel Zheng, Yury Sulsky, Shaobo Hou,
252 Tom Le Paine, Antoine Yang, Jason Riesa, Dominika Rogozinska, Dror Marcus, Dalia El
253 Badawy, Qiao Zhang, Luyu Wang, Helen Miller, Jeremy Greer, Lars Lowe Sjos, Azade
254 Nova, Heiga Zen, Rahma Chaabouni, Mihaela Rosca, Jiepu Jiang, Charlie Chen, RuiBo
255 Liu, Tara Sainath, Maxim Krikun, Alex Polozov, Jean-Baptiste Lespiau, Josh Newlan,
256 Zeyncep Cankara, Soo Kwak, Yunhan Xu, Phil Chen, Andy Coenen, Clemens Meyer,
257 Katerina Tsihlias, Ada Ma, Juraj Gottweis, Jinwei Xing, Chenjie Gu, Jin Miao, Chris-
258 tian Frank, Zeynep Cankara, Sanjay Ganapathy, Ishita Dasgupta, Steph Hughes-Fitt,
259 Heng Chen, David Reid, Keran Rong, Hongmin Fan, Joost van Amersfoort, Vincent
260 Zhuang, Aaron Cohen, Shixiang Shane Gu, Anhad Mohanane, Anastasija Ilic, Taylor
261 Tobin, John Wieting, Anna Bortsova, Phoebe Thacker, Emma Wang, Emily Caveness,
262 Justin Chiu, Eren Sezener, Alex Kaskasoli, Steven Baker, Katie Millican, Mohamed El-
263 hawaty, Kostas Aisopos, Carl Lebsack, Nathan Byrd, Hanjun Dai, Wenhao Jia, Matthew
264 Wiethoff, Elnaz Davoodi, Albert Weston, Lakshman Yagati, Arun Ahuja, Isabel Gao,
265 Golan Pundak, Susan Zhang, Michael Azzam, Khe Chai Sim, Sergi Caelles, James Keel-
266 ing, Abhanshu Sharma, Andy Swing, YaGuang Li, Chenxi Liu, Carrie Grimes Bostock,
267 Yamini Bansal, Zachary Nado, Ankesh Anand, Josh Lipschultz, Abhijit Karmarkar, Lev
268 Proleev, Abe Ittycheriah, Soheil Hassas Yeganeh, George Polovets, Aleksandra Faust,
269 Jiao Sun, Alban Rustemi, Pen Li, Rakesh Shivanna, Jeremiah Liu, Chris Welty, Feder-
270 erico Lebron, Anirudh Baddepudi, Sebastian Krause, Emilio Parisotto, Radu Soricut,
271 Zheng Xu, Dawn Bloxwich, Melvin Johnson, Behnam Neyshabur, Justin Mao-Jones, Ren-
272 shen Wang, Vinay Ramasesh, Zaheer Abbas, Arthur Guez, Constant Segal, Duc Dung
273 Nguyen, James Svensson, Le Hou, Sarah York, Kieran Milan, Sophie Bridgers, Wik-
274 tor Gworek, Marco Tagliasacchi, James Lee-Thorp, Michael Chang, Alexey Guseynov,
275 Ale Jakse Hartman, Michael Kwong, Ruizhe Zhao, Sheleem Kashem, Elizabeth Cole,
276 Antoine Miech, Richard Tanburn, Mary Phuong, Filip Pavetic, Sebastien Cevey, Ra-
277 mona Comanescu, Richard Ives, Sherry Yang, Cosmo Du, Bo Li, Zizhao Zhang, Mariko
278 Iinuma, Clara Huiyi Hu, Aurko Roy, Shaan Bijwadia, Zhenkai Zhu, Danilo Martins,
279 Rachel Saputro, Anita Gergely, Steven Zheng, Dawei Jia, Ioannis Antonoglou, Adam
280 Sadovsky, Shane Gu, Yingying Bi, Alek Andreev, Sina Samangooei, Mina Khan, Tomas
281 Kocisky, Angelos Filos, Chintu Kumar, Colton Bishop, Adams Yu, Sarah Hodkinson,
282 Sid Mittal, Premal Shah, Alexandre Moufarek, Yong Cheng, Adam Bloniarz, Jaehoon
283 Lee, Pedram Pejman, Paul Michel, Stephen Spencer, Vladimir Feinberg, Xuehan Xiong,
284 Nikolay Savinov, Charlotte Smith, Siamak Shakeri, Dustin Tran, Mary Chesus, Bernd
285 Bohnet, George Tucker, Tamara von Glehn, Carrie Muir, Yiran Mao, Hideto Kazawa,
286 Ambrose Slone, Kedar Soparkar, Disha Shrivastava, James Cobon-Kerr, Michael Shar-
287 man, Jay Pavagadhi, Carlos Araya, Karolis Misiunas, Nimesh Ghelani, Michael Laskin,
288 David Barker, Qiuqia Li, Anton Briukhov, Neil Houlsby, Mia Glaese, Balaji Lakshmi-

289 narayanan, Nathan Schucher, Yunhao Tang, Eli Collins, Hyeontaek Lim, Fangxiaoyu
290 Feng, Adria Recasens, Guangda Lai, Alberto Magni, Nicola De Cao, Aditya Siddhant,
291 Zoe Ashwood, Jordi Orbay, Mostafa Dehghani, Jenny Brennan, Yifan He, Kelvin Xu,
292 Yang Gao, Carl Saroufim, James Molloy, Xinyi Wu, Seb Arnold, Solomon Chang, Ju-
293 lian Schrittwieser, Elena Buchatskaya, Soroush Radpour, Martin Polacek, Skye Giordano,
294 Ankur Bapna, Simon Tokumine, Vincent Hellendoorn, Thibault Sottiaux, Sarah Cogan,
295 Aliaksei Severyn, Mohammad Saleh, Shantanu Thakoor, Laurent Shefey, Siyuan Qiao,
296 Meenu Gaba, Shuo yiin Chang, Craig Swanson, Biao Zhang, Benjamin Lee, Paul Kis-
297 han Rubenstein, Gan Song, Tom Kwiatkowski, Anna Koop, Ajay Kannan, David Kao,
298 Parker Schuh, Axel Stjerngren, Golnaz Ghiasi, Gena Gibson, Luke Vilnis, Ye Yuan, Fe-
299 lipe Tiengo Ferreira, Aishwarya Kamath, Ted Klimentko, Ken Franko, Kefan Xiao, Indro
300 Bhattacharya, Miteyan Patel, Rui Wang, Alex Morris, Robin Strudel, Vivek Sharma, Pe-
301 ter Choy, Sayed Hadi Hashemi, Jessica Landon, Mara Finkelstein, Priya Jhakra, Justin
302 Frye, Megan Barnes, Matthew Mauger, Dennis Daun, Khuslen Baatarsukh, Matthew
303 Tung, Wael Farhan, Henryk Michalewski, Fabio Viola, Felix de Chaumont Quiry, Char-
304 line Le Lan, Tom Hudson, Qingze Wang, Felix Fischer, Ivy Zheng, Elspeth White, Anca
305 Dragan, Jean baptiste Alayrac, Eric Ni, Alexander Pritzel, Adam Iwanicki, Michael Is-
306 ard, Anna Bulanova, Lukas Zilka, Ethan Dyer, Devendra Sachan, Srivatsan Srinivasan,
307 Hannah Muckenhirn, Honglong Cai, Amol Mandhane, Mukarram Tariq, Jack W. Rae,
308 Gary Wang, Kareem Ayoub, Nicholas FitzGerald, Yao Zhao, Woohyun Han, Chris Al-
309 berti, Dan Garrette, Kashyap Krishnakumar, Mai Gimenez, Anselm Levskaya, Daniel
310 Sohn, Josip Matak, Inaki Iturrate, Michael B. Chang, Jackie Xiang, Yuan Cao, Nishant
311 Ranka, Geoff Brown, Adrian Hutter, Vahab Mirrokni, Nanxin Chen, Kaisheng Yao,
312 Zoltan Egyed, Francois Galilee, Tyler Liechty, Praveen Kallakuri, Evan Palmer, Sanjay
313 Ghemawat, Jasmine Liu, David Tao, Chloe Thornton, Tim Green, Mimi Jasarevic, Sharon
314 Lin, Victor Cotruta, Yi-Xuan Tan, Noah Fiedel, Hongkun Yu, Ed Chi, Alexander Neitz,
315 Jens Heitkaemper, Anu Sinha, Denny Zhou, Yi Sun, Charbel Kaed, Brice Hulse, Swaroop
316 Mishra, Maria Georgaki, Sneha Kudugunta, Clement Farabet, Izhak Shafran, Daniel Vlas-
317 sic, Anton Tsitsulin, Rajagopal Ananthanarayanan, Alen Carin, Guolong Su, Pei Sun,
318 Shashank V, Gabriel Carvajal, Josef Broder, Iulia Comsa, Alena Repina, William Wong,
319 Warren Weilun Chen, Peter Hawkins, Egor Filonov, Lucia Loher, Christoph Hirschschall,
320 Weiyi Wang, Jingchen Ye, Andrea Burns, Hardie Cate, Diana Gage Wright, Federico Pic-
321 cinini, Lei Zhang, Chu-Cheng Lin, Ionel Gog, Yana Kulizhskaya, Ashwin Sreevatsa,
322 Shuang Song, Luis C. Cobo, Anand Iyer, Chetan Tekur, Guillermo Garrido, Zhuyun
323 Xiao, Rupert Kemp, Huaixiu Steven Zheng, Hui Li, Ananth Agarwal, Christel Ngani,
324 Kati Goshvadi, Rebeca Santamaria-Fernandez, Wojciech Fica, Xinyun Chen, Chris Gor-
325 golewski, Sean Sun, Roopal Garg, Xinyu Ye, S. M. Ali Eslami, Nan Hua, Jon Simon,
326 Pratik Joshi, Yelin Kim, Ian Tenney, Sahitya Potluri, Lam Nguyen Thiet, Quan Yuan, Flor-
327 rian Luisier, Alexandra Chronopoulou, Salvatore Scellato, Praveen Srinivasan, Minmin
328 Chen, Vinod Koverkathu, Valentin Dalibard, Yaming Xu, Brennan Saeta, Keith Ander-
329 son, Thibault Sellam, Nick Fernando, Fantine Huot, Junehyuk Jung, Mani Varadarajan,
330 Michael Quinn, Amit Raul, Maigo Le, Ruslan Habalov, Jon Clark, Komal Jalan, Kalesha
331 Bullard, Achintya Singhal, Thang Luong, Boyu Wang, Sujeevan Rajayogam, Julian Eisen-
332 schlos, Johnson Jia, Daniel Finchelstein, Alex Yakubovich, Daniel Balle, Michael Fink,
333 Sameer Agarwal, Jing Li, Dj Dvijotham, Shalini Pal, Kai Kang, Jaclyn Konzelmann, Jen-
334 nifer Beattie, Olivier Dousse, Diane Wu, Remi Crocker, Chen Elkind, Siddhartha Reddy
335 Jonnalagadda, Jong Lee, Dan Holtmann-Rice, Krystal Kallarackal, Rosanne Liu, Denis
336 Vnukov, Neera Vats, Luca Invernizzi, Mohsen Jafari, Huanjie Zhou, Lilly Taylor, Jennifer
337 Prendki, Marcus Wu, Tom Eccles, Tianqi Liu, Kavya Kopparapu, Francoise Beaufays,
338 Christof Angermueller, Andreea Marzoca, Shourya Sarcar, Hilal Dib, Jeff Stanway, Frank
339 Perbet, Nejc Trdin, Rachel Sterneck, Andrey Khorlin, Dinghua Li, Xihui Wu, Sonam
340 Goenka, David Madras, Sasha Goldshtein, Willi Gierke, Tong Zhou, Yaxin Liu, Yannic
341 Liang, Anais White, Yunjie Li, Shreya Singh, Sanaz Bahargam, Mark Epstein, Sujoy Basu,
342 Li Lao, Adnan Ozturel, Carl Crous, Alex Zhai, Han Lu, Zora Tung, Neeraj Gaur, Alanna
343 Walton, Lucas Dixon, Ming Zhang, Amir Globerson, Grant Uy, Andrew Bolt, Olivia
344 Wiles, Milad Nasr, Ilia Shumailov, Marco Selvi, Francesco Piccinno, Ricardo Aguilar,
345 Sara McCarthy, Misha Khalman, Mrinal Shukla, Vlado Galic, John Carpenter, Kevin Vil-
346 lela, Haibin Zhang, Harry Richardson, James Martens, Matko Bosnjak, Shreyas Rammo-
347 han Belle, Jeff Seibert, Mahmoud Alnahlawi, Brian McWilliams, Sankalp Singh, Annie

348 Louis, Wen Ding, Dan Popovici, Lenin Simicich, Laura Knight, Pulkit Mehta, Nishesh
349 Gupta, Chongyang Shi, Saaber Fatehi, Jovana Mitrovic, Alex Grills, Joseph Pagadora,
350 Dessie Petrova, Danielle Eisenbud, Zhishuai Zhang, Damion Yates, Bhavishya Mittal,
351 Nilesh Tripuraneni, Yannis Assael, Thomas Brovelli, Prateek Jain, Mihajlo Velimirovic,
352 Canfer Akbulut, Jiaqi Mu, Wolfgang Macherey, Ravin Kumar, Jun Xu, Haroon Qureshi,
353 Gheorghe Comanici, Jeremy Wiesner, Zhitao Gong, Anton Ruddock, Matthias Bauer,
354 Nick Felt, Anirudh GP, Anurag Arnab, Dustin Zelle, Jonas Rothfuss, Bill Rosgen, Ashish
355 Shenoy, Bryan Seybold, Xinjian Li, Jayaram Mudigonda, Goker Erdogan, Jiawei Xia,
356 Jiri Simsa, Andrea Michi, Yi Yao, Christopher Yew, Steven Kan, Isaac Caswell, Carey
357 Radebaugh, Andre Elisseeff, Pedro Valenzuela, Kay McKinney, Kim Paterson, Albert
358 Cui, Eri Latorre-Chimoto, Solomon Kim, William Zeng, Ken Durden, Priya Ponnappalli,
359 Tiberiu Sosea, Christopher A. Choquette-Choo, James Manyika, Brona Robenek, Harsha
360 Vashisht, Sebastien Pereira, Hoi Lam, Marko Velic, Denese Owusu-Afriyie, Katherine
361 Lee, Tolga Bolukbasi, Alicia Parrish, Shawn Lu, Jane Park, Balaji Venkatraman, Alice Tal-
362 bert, Lambert Rosique, Yuchung Cheng, Andrei Sozanschi, Adam Paszke, Praveen Ku-
363 mar, Jessica Austin, Lu Li, Khalid Salama, Wooyeol Kim, Nandita Dukkupati, Anthony
364 Baryshnikov, Christos Kaplanis, XiangHai Sheng, Yuri Chervonyi, Caglar Unlu, Diego
365 de Las Casas, Harry Askham, Kathryn Tunyasuvunakool, Felix Gimeno, Siim Poder,
366 Chester Kwak, Matt Miecnikowski, Vahab Mirrokni, Alek Dimitriev, Aaron Parisi, Dang-
367 gyi Liu, Tomy Tsai, Toby Shevlane, Christina Kouridi, Drew Garmon, Adrian Goedeck-
368 emeyer, Adam R. Brown, Anitha Vijayakumar, Ali Elqursh, Sadegh Jazayeri, Jin Huang,
369 Sara Mc Carthy, Jay Hoover, Lucy Kim, Sandeep Kumar, Wei Chen, Courtney Biles, Gar-
370 rett Bingham, Evan Rosen, Lisa Wang, Qijun Tan, David Engel, Francesco Pongetti, Dario
371 de Cesare, Dongseong Hwang, Lily Yu, Jennifer Pullman, Srini Narayanan, Kyle Levin,
372 Siddharth Gopal, Megan Li, Asaf Aharoni, Trieu Trinh, Jessica Lo, Norman Casagrande,
373 Roopali Vij, Loic Matthéy, Bramandia Ramadhana, Austin Matthews, CJ Carey, Matthew
374 Johnson, Kremena Goranova, Rohin Shah, Shereen Ashraf, Kingshuk Dasgupta, Rasmus
375 Larsen, Yicheng Wang, Manish Reddy Vuyyuru, Chong Jiang, Joana Ijazi, Kazuki Osawa,
376 Celine Smith, Ramya Sree Boppana, Taylan Bilal, Yuma Koizumi, Ying Xu, Yasemin Al-
377 tun, Nir Shabat, Ben Bariach, Alex Korchemniy, Kiam Choo, Olaf Ronneberger, Chimezie
378 Iwuanyanwu, Shubin Zhao, David Soergel, Cho-Jui Hsieh, Irene Cai, Shariq Iqbal, Mar-
379 tin Sundermeyer, Zhe Chen, Elie Bursztein, Chaitanya Malaviya, Fadi Biadsy, Prakash
380 Shroff, Inderjit Dhillon, Tejasi Latkar, Chris Dyer, Hannah Forbes, Massimo Nicosia, Vi-
381 taly Nikolaev, Somer Greene, Marin Georgiev, Pidong Wang, Nina Martin, Hanie Sedghi,
382 John Zhang, Praseem Banzal, Doug Fritz, Vikram Rao, Xuezhi Wang, Jiageng Zhang,
383 Viorica Patraucean, Dayou Du, Igor Mordatch, Ivan Jurin, Lewis Liu, Ayush Dubey,
384 Abhi Mohan, Janek Nowakowski, Vlad-Doru Ion, Nan Wei, Reiko Tojo, Maria Abi Raad,
385 Drew A. Hudson, Vaishakh Keshava, Shubham Agrawal, Kevin Ramirez, Zhichun Wu,
386 Hoang Nguyen, Ji Liu, Madhavi Sewak, Bryce Petrini, DongHyun Choi, Ivan Philips,
387 Ziyue Wang, Ioana Bica, Ankush Garg, Jarek Wilkiewicz, Priyanka Agrawal, Xiaowei
388 Li, Danhao Guo, Emily Xue, Naseer Shaik, Andrew Leach, Sadh MNM Khan, Julia
389 Wiesinger, Sammy Jerome, Abhishek Chakladar, Alek Wenjiao Wang, Tina Ornduff,
390 Folake Abu, Alireza Ghaffarkhah, Marcus Wainwright, Mario Cortes, Frederick Liu,
391 Joshua Maynez, Andreas Terzis, Pouya Samangouei, Riham Mansour, Tomasz Kepa,
392 François-Xavier Aubet, Anton Algymr, Dan Banica, Agoston Weisz, Andras Orban,
393 Alexandre Senges, Ewa Andrejczuk, Mark Geller, Niccolo Dal Santo, Valentin Anklin,
394 Majd Al Merey, Martin Baeuml, Trevor Strohman, Junwen Bai, Slav Petrov, Yonghui
395 Wu, Demis Hassabis, Koray Kavukcuoglu, Jeffrey Dean, and Oriol Vinyals. Gemini 1.5:
396 Unlocking multimodal understanding across millions of tokens of context, 2024. URL
397 <https://arxiv.org/abs/2403.05530>.

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

400 Databricks. Dbrx, 2024. URL [https://www.databricks.com/blog/](https://www.databricks.com/blog/introducing-dbrx-new-state-art-open-llm)
401 [introducing-dbrx-new-state-art-open-llm](https://www.databricks.com/blog/introducing-dbrx-new-state-art-open-llm).

402 Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary,
403 Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian
404 Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lample, L lio Renard Lavaud,

405 Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang,
406 Szymon Antoniak, Teven Le Scao, Théophile Gervet, Thibaut Lavril, Thomas Wang,
407 Timothée Lacroix, and William El Sayed. Mixtral of experts, 2024. URL <https://arxiv.org/abs/2401.04088>.
408

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

411 DeepSeek-AI, Aixin Liu, Bei Feng, Bin Wang, Bingxuan Wang, Bo Liu, Chenggang Zhao,
412 Chengqi Deng, Chong Ruan, Damai Dai, Daya Guo, Dejian Yang, Deli Chen, Dongjie Ji,
413 Erhang Li, Fangyun Lin, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang,
414 Hanwei Xu, Hao Yang, Haowei Zhang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui
415 Li, Hui Qu, J. L. Cai, Jian Liang, Jianzhong Guo, Jiaqi Ni, Jiashi Li, Jin Chen, Jingyang
416 Yuan, Junjie Qiu, Junxiao Song, Kai Dong, Kaige Gao, Kang Guan, Lean Wang, Lecong
417 Zhang, Lei Xu, Leyi Xia, Liang Zhao, Liyue Zhang, Meng Li, Miaojun Wang, Mingchuan
418 Zhang, Minghua Zhang, Minghui Tang, Mingming Li, Ning Tian, Panpan Huang, Peiyi
419 Wang, Peng Zhang, Qihao Zhu, Qinyu Chen, Qiushi Du, R. J. Chen, R. L. Jin, Ruiqi Ge,
420 Ruizhe Pan, Runxin Xu, Ruyi Chen, S. S. Li, Shanghao Lu, Shangyan Zhou, Shanhuang
421 Chen, Shaoqing Wu, Shengfeng Ye, Shirong Ma, Shiyu Wang, Shuang Zhou, Shuiping
422 Yu, Shunfeng Zhou, Size Zheng, T. Wang, Tian Pei, Tian Yuan, Tianyu Sun, W. L. Xiao,
423 Wangding Zeng, Wei An, Wen Liu, Wenfeng Liang, Wenjun Gao, Wentao Zhang, X. Q.
424 Li, Xiangyue Jin, Xianzu Wang, Xiao Bi, Xiaodong Liu, Xiaohan Wang, Xiaojin Shen,
425 Xiaokang Chen, Xiaosha Chen, Xiaotao Nie, Xiaowen Sun, Xiaoxiang Wang, Xin Liu,
426 Xin Xie, Xingkai Yu, Xinnan Song, Xinyi Zhou, Xinyu Yang, Xuan Lu, Xuecheng Su,
427 Y. Wu, Y. K. Li, Y. X. Wei, Y. X. Zhu, Yanhong Xu, Yanping Huang, Yao Li, Yao Zhao,
428 Yaofeng Sun, Yaohui Li, Yaohui Wang, Yi Zheng, Yichao Zhang, Yiliang Xiong, Yilong
429 Zhao, Ying He, Ying Tang, Yishi Piao, Yixin Dong, Yixuan Tan, Yiyuan Liu, Yongji Wang,
430 Yongqiang Guo, Yuchen Zhu, Yuduan Wang, Yuheng Zou, Yukun Zha, Yunxian Ma,
431 Yuting Yan, Yuxiang You, Yuxuan Liu, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhen
432 Huang, Zhen Zhang, Zhenda Xie, Zhewen Hao, Zhihong Shao, Zhiniu Wen, Zhipeng Xu,
433 Zhongyu Zhang, Zhuoshu Li, Zihan Wang, Zihui Gu, Zilin Li, and Ziwei Xie. Deepseek-
434 v2: A strong, economical, and efficient mixture-of-experts language model, 2024. URL
435 <https://arxiv.org/abs/2405.04434>.

436 Aidan Clark, Diego de las Casas, Aurelia Guy, Arthur Mensch, Michela Paganini, Jordan
437 Hoffmann, Bogdan Damoc, Blake Hechtman, Trevor Cai, Sebastian Borgeaud, George
438 van den Driessche, Eliza Rutherford, Tom Hennigan, Matthew Johnson, Katie Millin-
439 can, Albin Cassirer, Chris Jones, Elena Buchatskaya, David Budden, Laurent Sifre, Si-
440 mon Osindero, Oriol Vinyals, Jack Rae, Erich Elsen, Koray Kavukcuoglu, and Karen
441 Simonyan. Unified scaling laws for routed language models, 2022. URL <https://arxiv.org/abs/2202.01169>.
442

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

447 Dmitry Lepikhin, HyoukJoong Lee, Yuanzhong Xu, Dehao Chen, Orhan Firat, Yanping
448 Huang, Maxim Krikun, Noam Shazeer, and Zhifeng Chen. Gshard: Scaling giant models
449 with conditional computation and automatic sharding. *arXiv preprint arXiv:2006.16668*,
450 2020.

451 William Fedus, Barret Zoph, and Noam Shazeer. Switch transformers: Scaling to trillion
452 parameter models with simple and efficient sparsity, 2022. URL [https://arxiv.org/
453 abs/2101.03961](https://arxiv.org/abs/2101.03961).

454 Opher Lieber, Barak Lenz, Hofit Bata, Gal Cohen, Jhonathan Osin, Itay Dalmedigos,
455 Erez Safahi, Shaked Meir, Yonatan Belinkov, Shai Shalev-Shwartz, Omri Abend, Raz
456 Alon, Tomer Asida, Amir Bergman, Roman Glozman, Michael Gokhman, Avshalom
457 Manevich, Nir Ratner, Noam Rozen, Erez Shwartz, Mor Zusman, and Yoav Shoham.
458 Jamba: A hybrid transformer-mamba language model, 2024. URL [https://arxiv.
459 org/abs/2403.19887](https://arxiv.org/abs/2403.19887).

- 460 Liquid. Liquid, 2024. URL <https://www.liquid.ai/liquid-foundation-models>.
- 461 Robert A. Jacobs, Michael I. Jordan, Steven J. Nowlan, and Geoffrey E. Hinton. Adaptive
462 Mixtures of Local Experts. *Neural Computation*, 3(1):79–87, 03 1991. ISSN 0899-7667. doi:
463 10.1162/neco.1991.3.1.79. URL <https://doi.org/10.1162/neco.1991.3.1.79>.
- 464 M. I. Jordan and R. A. Jacobs. Hierarchical mixtures of experts and the em algorithm.
465 In Maria Marinaro and Pietro G. Morasso, editors, *ICANN '94*, pages 479–486, London,
466 1994. Springer London. ISBN 978-1-4471-2097-1.
- 467 Barret Zoph, Irwan Bello, Sameer Kumar, Nan Du, Yanping Huang, Jeff Dean, Noam
468 Shazeer, and William Fedus. St-moe: Designing stable and transferable sparse expert
469 models, 2022. URL <https://arxiv.org/abs/2202.08906>.
- 470 Yanqi Zhou, Tao Lei, Hanxiao Liu, Nan Du, Yanping Huang, Vincent Zhao, Andrew Dai,
471 Zhifeng Chen, Quoc Le, and James Laudon. Mixture-of-experts with expert choice rout-
472 ing, 2022. URL <https://arxiv.org/abs/2202.09368>.
- 473 Yoshua Bengio, Nicholas Léonard, and Aaron Courville. Estimating or propagating gra-
474 dients through stochastic neurons for conditional computation, 2013. URL <https://arxiv.org/abs/1308.3432>.
- 476 Damai Dai, Chengqi Deng, Chenggang Zhao, R. X. Xu, Huazuo Gao, Deli Chen, Jiashi
477 Li, Wangding Zeng, Xingkai Yu, Y. Wu, Zhenda Xie, Y. K. Li, Panpan Huang, Fuli
478 Luo, Chong Ruan, Zhifang Sui, and Wenfeng Liang. Deepseekmoe: Towards ulti-
479 mate expert specialization in mixture-of-experts language models, 2024. URL <https://arxiv.org/abs/2401.06066>.
- 481 Lean Wang, Huazuo Gao, Chenggang Zhao, Xu Sun, and Damai Dai. Auxiliary-loss-free
482 load balancing strategy for mixture-of-experts, 2024. URL <https://arxiv.org/abs/2408.15664>.
- 484 Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed Awadallah, Ammar Ahmad Awan,
485 Nguyen Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, Alon Benhaim,
486 Misha Bilenko, Johan Bjorck, Sébastien Bubeck, Martin Cai, Qin Cai, Vishrav Chaudhary,
487 Dong Chen, Dongdong Chen, Weizhu Chen, Yen-Chun Chen, Yi-Ling Chen, Hao Cheng,
488 Parul Chopra, Xiyang Dai, Matthew Dixon, Ronen Eldan, Victor Fragoso, Jianfeng Gao,
489 Mei Gao, Min Gao, Amit Garg, Allie Del Giorno, Abhishek Goswami, Suriya Gunasekar,
490 Emman Haider, Junheng Hao, Russell J. Hewett, Wenxiang Hu, Jamie Huynh, Dan Iter,
491 Sam Ade Jacobs, Mojan Javaheripi, Xin Jin, Nikos Karampatziakis, Piero Kauffmann,
492 Mahoud Khademi, Dongwoo Kim, Young Jin Kim, Lev Kurilenko, James R. Lee, Yin Tat
493 Lee, Yuanzhi Li, Yunsheng Li, Chen Liang, Lars Liden, Xihui Lin, Zeqi Lin, Ce Liu,
494 Liyuan Liu, Mengchen Liu, Weishung Liu, Xiaodong Liu, Chong Luo, Piyush Madan, Ali
495 Mahmoudzadeh, David Majercak, Matt Mazzola, Caio César Teodoro Mendes, Arindam
496 Mitra, Hardik Modi, Anh Nguyen, Brandon Norick, Barun Patra, Daniel Perez-Becker,
497 Thomas Portet, Reid Pryzant, Heyang Qin, Marko Radmilac, Liliang Ren, Gustavo
498 de Rosa, Corby Rosset, Sambudha Roy, Olatunji Ruwase, Olli Saarikivi, Amin Saied,
499 Adil Salim, Michael Santacroce, Shital Shah, Ning Shang, Hiteshi Sharma, Yelong Shen,
500 Swadheen Shukla, Xia Song, Masahiro Tanaka, Andrea Tupini, Praneetha Vaddamanu,
501 Chunyu Wang, Guanhua Wang, Lijuan Wang, Shuohang Wang, Xin Wang, Yu Wang,
502 Rachel Ward, Wen Wen, Philipp Witte, Haiping Wu, Xiaoxia Wu, Michael Wyatt, Bin
503 Xiao, Can Xu, Jiahang Xu, Weijian Xu, Jilong Xue, Sonali Yadav, Fan Yang, Jianwei Yang,
504 Yifan Yang, Ziyi Yang, Donghan Yu, Lu Yuan, Chenruidong Zhang, Cyril Zhang, Jian-
505 wen Zhang, Li Lyna Zhang, Yi Zhang, Yue Zhang, Yunan Zhang, and Xiren Zhou. Phi-3
506 technical report: A highly capable language model locally on your phone, 2024. URL
507 <https://arxiv.org/abs/2404.14219>.
- 508 Liyuan Liu, Young Jin Kim, Shuohang Wang, Chen Liang, Yelong Shen, Hao Cheng, Xi-
509 aodong Liu, Masahiro Tanaka, Xiaoxia Wu, Wenxiang Hu, Vishrav Chaudhary, Zeqi Lin,
510 Chenruidong Zhang, Jilong Xue, Hany Awadalla, Jianfeng Gao, and Weizhu Chen. Grin:
511 Gradient-informed moe, 2024. URL <https://arxiv.org/abs/2409.12136>.

512 Liyuan Liu, Jianfeng Gao, and Weizhu Chen. Sparse backpropagation for moe training,
513 2023. URL <https://arxiv.org/abs/2310.00811>.

514 Noam Shazeer. Glu variants improve transformer, 2020. URL <https://arxiv.org/abs/2002.05202>.

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

521 Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. Layer normalization, 2016. URL
522 <https://arxiv.org/abs/1607.06450>.

523 Jianlin Su, Yu Lu, Shengfeng Pan, Ahmed Murtadha, Bo Wen, and Yunfeng Liu. Roformer:
524 Enhanced transformer with rotary position embedding, 2023. URL <https://arxiv.org/abs/2104.09864>.

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

530 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle,
531 Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal,
532 Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev,
533 Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson,
534 Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte
535 Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller,
536 Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Niko-
537 laidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, David Esiobu,
538 Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hup-
539 kes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith,
540 Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Ander-
541 son, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen,
542 Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel
543 Kloumann, Ishan Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranne,
544 Jason Park, Jay Mahadeokar, Jeet Shah, Jelder van der Linde, Jennifer Billock, Jenny
545 Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao
546 Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua
547 Saxe, Junteng Jia, Kalyan Vasuden Alwala, Kartikeya Upasani, Kate Plawiak, Ke Li,
548 Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley
549 Chiu, Kunal Bhalla, Lauren Rantala-Yearry, Laurens van der Maaten, Lawrence Chen,
550 Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher,
551 Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh,
552 Manohar Paluri, Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie
553 Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Nar-
554 jes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Olivier Duchenne,
555 Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng,
556 Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing
557 He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Sil-
558 veira Cabral, Robert Stojnic, Roberta Raileanu, Rohit Girdhar, Rohit Patel, Romain
559 Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui
560 Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun So-
561 nia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Rapparthi, Sheng Shen,
562 Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer
563 Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar
564 Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speck-
565 bacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vign-
566 nesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vladan

567 Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet,
 568 Xiaodong Wang, Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao Jia, Xuwei Wang, Yaelle
 569 Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue
 570 Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Pa-
 571 pakipos, Aaditya Singh, Aaron Grattafiori, Abha Jain, Adam Kelsey, Adam Shajnfeld,
 572 Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex
 573 Boesenberg, Alex Vaughan, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit San-
 574 gani, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, An-
 575 drew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Franco, Aparajita
 576 Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisen-
 577 man, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang,
 578 Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden
 579 Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo,
 580 Carl Parker, Carly Burton, Catalina Mejia, Changhan Wang, Changkyu Kim, Chao
 581 Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichten-
 582 hofer, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, Danny Wyatt, David Ad-
 583 kins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem
 584 Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine
 585 Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Erik Brinkman, Esteban Ar-
 586 caute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Firat Ozgenel,
 587 Francesco Caggioni, Francisco Guzmán, Frank Kanayet, Frank Seide, Gabriela Medina
 588 Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Govind Thattai,
 589 Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hamid Sho-
 590 janazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph,
 591 Helen Suk, Henry Aspegren, Hunter Goldman, Ibrahim Damlaj, Igor Molybog, Igor Tu-
 592 fanov, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Japhet
 593 Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy
 594 Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon
 595 Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai
 596 Wu, Kam Hou U, Karan Saxena, Karthik Prasad, Kartikay Khandelwal, Katayoun Zand,
 597 Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kun Huang, Ku-
 598 nal Chawla, Kushal Lakhota, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Le-
 599 andro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca
 600 Wehrstedt, Madian Khabza, Manav Avalani, Manish Bhatt, Maria Tsimpoukelli, Marty-
 601 nas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim
 602 Naumov, Maya Lathi, Meghan Keneally, Michael L. Seltzer, Michal Valko, Michelle Re-
 603 strepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike
 604 Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal,
 605 Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal,
 606 Nick Egebo, Nicolas Usunier, Nikolay Pavlovich Laptev, Ning Dong, Ning Zhang, Nor-
 607 man Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent,
 608 Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr
 609 Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad
 610 Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra,
 611 Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Rohan Maheswari, Russ
 612 Howes, Ruty Rinott, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sar-
 613 gun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Verma, Seiji Yamamoto, Sharadh Ra-
 614 maswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy
 615 Zha, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal,
 616 Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve
 617 Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Sungmin Cho, Sunny Virk, Suraj
 618 Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best,
 619 Thilo Kohler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy
 620 Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mo-
 621 han, Vinay Satish Kumar, Vishal Mangla, Vítor Albiero, Vlad Ionescu, Vlad Poenaru,
 622 Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes
 623 Bouaziz, Will Constable, Xiaocheng Tang, Xiaofang Wang, Xiaojian Wu, Xiaolan Wang,
 624 Xide Xia, Xilun Wu, Xinbo Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin
 625 Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yuchen Hao, Yundi Qian, Yuzi

- 626 He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, and Zhi-
627 wei Zhao. The llama 3 herd of models, 2024. URL [https://arxiv.org/abs/2407.](https://arxiv.org/abs/2407.21783)
628 21783.
- 629 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization, 2019. URL
630 <https://arxiv.org/abs/1711.05101>.
- 631 Adam Ibrahim, Benjamin Thérien, Kshitij Gupta, Mats L. Richter, Quentin Anthony, Tim-
632 othée Lesort, Eugene Belilovsky, and Irina Rish. Simple and scalable strategies to contin-
633 ually pre-train large language models, 2024. URL [https://arxiv.org/abs/2403.](https://arxiv.org/abs/2403.08763)
634 08763.
- 635 Alex Andonian, Quentin Anthony, Stella Biderman, Sid Black, Preetham Gali, Leo Gao, Eric
636 Hallahan, Josh Levy-Kramer, Connor Leahy, Lucas Nestler, Kip Parker, Michael Pieler,
637 Jason Phang, Shivanshu Purohit, Hailey Schoelkopf, Dashiell Stander, Tri Songz, Curt
638 Tigges, Benjamin Thérien, Phil Wang, and Samuel Weinbach. GPT-NeoX: Large Scale
639 Autoregressive Language Modeling in PyTorch, 9 2023. URL [https://www.github.](https://www.github.com/eleutherai/gpt-neox)
640 [com/eleutherai/gpt-neox](https://www.github.com/eleutherai/gpt-neox).
- 641 Trevor Gale, Deepak Narayanan, Cliff Young, and Matei Zaharia. Megablocks: Efficient
642 sparse training with mixture-of-experts, 2022. URL [https://arxiv.org/abs/2211.](https://arxiv.org/abs/2211.15841)
643 15841.

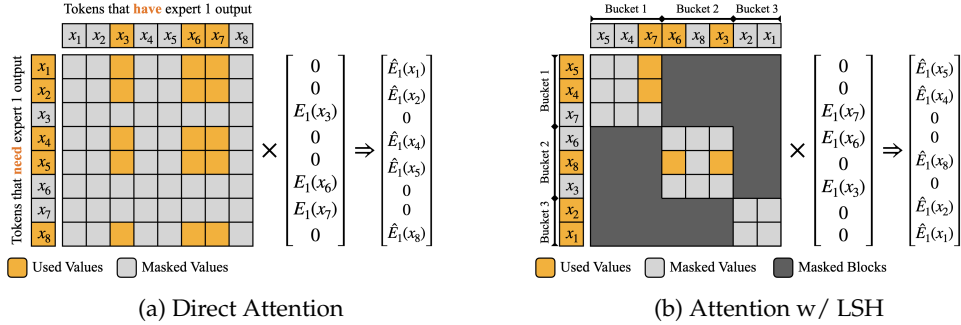


Figure 4: **Attention scores of direct and LSH attention methods.** For each expert, we define an attention head which uses queries corresponding to inputs not routed to the expert, and keys corresponding to inputs routed to the expert. Grey entries denote queries and keys which do not meet this criteria, and whose attention scores are masked out. The attention scores of each head will then be multiplied by the values, corresponding to expert outputs of tokens routed to that 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 entirely skip a majority of the attention computation.

644 A Token-specific Approximations Using Attention

645

646 A.1 Global Attention

647 Our Expert Group Approximation computes an approximation, for each expert, for all to-
 648 kens routed to it from each other expert, and in this manner computes N^2 approximations.
 649 However, we may want to actually compute an approximation for specific tokens. Con-
 650 sider tokens belonging to the set $x \in X_{\{i,\cdot\}}^C$, i.e. tokens *not* routed to expert i . We want to
 651 approximate $E_i(x)$ for such x . At a high level, we want to search for similar tokens to x ,
 652 select their expert outputs $E_i(x_j)$, and aggregate these outputs as a weighted linear com-
 653 bination. A well-known approach to this problem is attention. We want to query with
 654 all tokens *not* routed to expert i , i.e. $X_{\{i,\cdot\}}^C$. The keys will correspond to tokens that *were*
 655 routed to expert i , i.e. $X_{\{i,\cdot\}}$. And the values will be the expert outputs of these relevant in-
 656 puts. Fig. 4a (left) outlines how we compute an approximation using multi-head attention,
 657 where each head corresponds to approximating for a single expert.

658 A.2 Sparse Attention using LSH

659 Computing attention across all tokens on an accelerator is computationally expensive, and
 660 we do not need attention scores for *all* the tokens to compute the approximation, just for
 661 the most similar tokens to x . With a block-sparse attention mask, we can greatly reduce the
 662 attention computation especially when most of the computed scores would be redundant.
 663 In Fig. 4b we outline our attention approximation that uses locality-sensitive hashing (LSH)
 664 to group tokens into buckets, with a high probability that the nearest neighbors to a token
 665 will lie in the same bucket. The attention mask now has an additional condition: the query
 666 index q and key index k must correspond to tokens in the same bucket. We sort the QKV
 667 into groups based on their assigned buckets to encourage a block-diagonal attention mask,
 668 and verify that this sparsity reduces the runtime of our attention approximation. Note that
 669 as exemplified in Fig. 4b, it is possible that some tokens receive no approximation because
 670 there are no keys to query in the bucket. In this case, we set the approximation to 0.

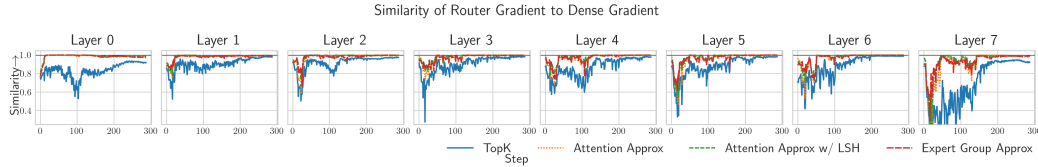


Figure 5: **Accuracy in approximating the dense router gradient for each approach.** This is recorded using a model with 8 experts and $K = 2$. The *dense gradient* of the output with respect to the router weights is artificially computed at each step by passing inputs through a dense mixture of experts layer, where all experts are selected. This is done independently from 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 using a standard TopK router, along with using each of our proposed router modifications. Our approaches are much more accurate and stable in approximating the dense router gradient.

671 B Approximation Statistics

672

673 B.1 Approximation Fidelity

674 We verify that our method is indeed faithfully approximating the dense router gradient
 675 i.e. the gradient to the router if all experts were activated. We track the dense gradient by
 676 routing to all experts and backpropagating only on the MoE output (independent of the
 677 full forward pass). This dense gradient is compared to the actual router gradient for each
 678 of our approaches in Fig. 5. We also observe a major difference between the gradient of the
 679 standard Top-K router and our approach.

680 The differences in our approaches become clear as we scale the model to become more
 681 sparse. We expand to $N = 32$ experts while maintaining $K = 2$ in Fig. 6 and find that it is
 682 more difficult to approximate the true dense router gradient. While all of our approaches
 683 sufficiently approximate the dense gradient with $N = 8$ experts, the performance gap be-
 684 tween them is apparent with $N = 32$. The expert group approximation and LSH attention
 685 methods are significantly better than the direct attention method, and this is also consistent
 686 with our validation results in Table 1. This is likely due to the heuristics we apply to re-
 687 strict our approximation to only the most relevant tokens: the expert group approximation
 688 requires inputs to have an expert in common, and LSH requires inputs to be similar. More-
 689 over, the gap between our methods and Top-K is wider with 32 experts. We believe that
 690 in larger models with even more experts, our method will yield increasingly significant
 691 improvements over Top-K routing.

692 In Fig. 7 we reproduce the gradient similarity plots with SparseMixer [Liu et al., 2023].
 693 Surprisingly, we find that SparseMixer is the worst approximation of the dense gradient
 694 across the board. Initial experiments also validate that SparseMixer does not outperform
 695 any of the other methods.

696 In Fig. 8 we provide an additional analysis of the gradient norm of our approximation
 697 compared to the dense gradient. We include statistics for SparseMixer as well. This logging
 698 is also done with $N = 8$ and $K = 2$. The Top-K gradient has significantly lower norm than
 699 the dense gradient, and the SparseMixer is an order of magnitude lower in many cases.
 700 Our methods closely approximate the dense gradient norm consistently; replicating both
 701 the direction and magnitude suggests that we are sufficiently approximating the dense
 702 gradient entirely.

703 B.2 Empirical Observations on Input Similarity

704 Our methods operate on the assumption that expert outputs for an input can be approx-
 705 imated by taking outputs from other similar inputs. We observe this during training by

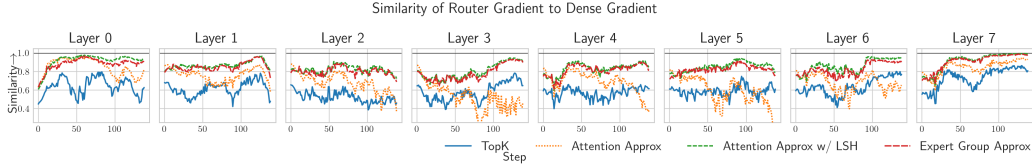


Figure 6: **Dense router gradient approximation accuracy with fine-grained experts.** We implement fine-grained experts as in DeepSeekMoE [DeepSeek-AI et al., 2024] to observe the behavior of our approximation methods across more experts while keeping parameter count fixed. In this example, the model 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 train batch size.

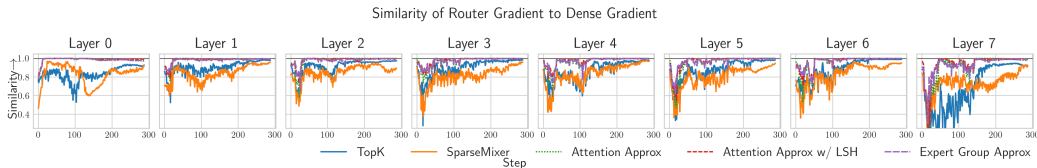


Figure 7: **Comparison of gradient approximation with our methods and SparseMixer.** We provide additional results showing how our gradient approximations compare against SparseMixer, another method to estimate the dense gradient.

706 partitioning each batch of inputs based on the experts that they are routed to. For each expert
 707 we compute cosine similarity among all possible pairs of inputs routed to the expert
 708 and cosine similarity between expert outputs of the corresponding pairs. We specifically
 709 track the expert output similarity when these inputs have a cosine similarity > 0.75 . In
 710 Fig. 10 we demonstrate that when inputs are very similar, they tend to have very similar
 711 expert outputs on average. Thus, we can approximate a missing expert output for a token
 712 by taking a nearby token’s expert output.

713 Moreover, our expert group approximation method specifically assumes that being routed
 714 to the same expert is a proxy indicating similarity. For this method to work, it must be
 715 the case that two inputs that share experts in common are similar on average. Another
 716 desirable property is that these inputs have a similarity above some threshold (> 0.75) with
 717 a very high probability. Then, when approximating an expert output for a token, it suffices
 718 to use the output for another token routed to one of the same experts. We demonstrate the
 719 two above properties empirically in Fig. 9. Inputs with one expert in common are not only
 720 very similar on average, they are also very similar with a high probability. This suggests
 721 that our gradient estimator using the expert group approximation method is both accurate
 722 and consistent.

723 C Experimental Setup

724 **Model Architecture.** We train an MoE with 24 blocks, a hidden dimension of 1024, and 8
 725 experts, for a total of 2B parameters, 780M of which are activated when we use the standard
 726 $K = 2$ top-K routing. We use SwiGLU [Shazeer, 2020] MLPs following Llama [Touvron
 727 et al., 2023], using an expansion factor such that the intermediate size of the MLP is 2816,
 728 16 attention heads with dimension is 64, LayerNorm [Ba et al., 2016] and RoPE [Su et al.,
 729 2023].

730 **Dataset.** We train on FineWeb [Penedo et al., 2024] with the Llama3 tokenizer [Dubey et al.,
 731 2024]. We split it into train, validation, and test splits and report the validation perplexity.

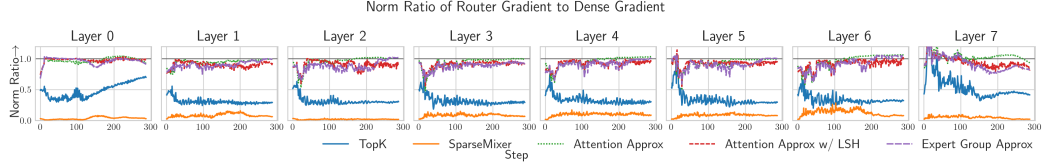
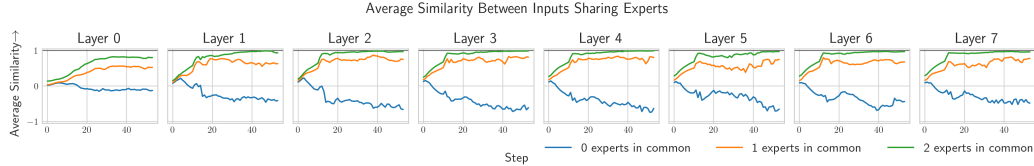
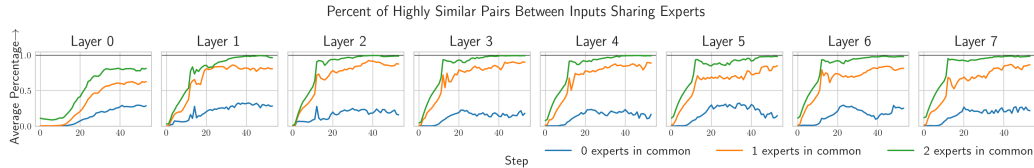


Figure 8: **Comparison of gradient norms relative to the dense gradient.** When computing the dense gradient, we also record its L_2 norm and log the ratio of this to the L_2 norms of the actual router gradients during training. Our methods produce router gradients with approximately the same magnitude. Along with the results showing strong cosine similarity, this suggests that we are almost perfectly approximating the dense gradient.



(a) Input Similarity



(b) Proportion of Highly Similar Inputs

Figure 9: **Similarity of Inputs Routed to Same Experts.** In a model with $N = 8$ experts and $K = 2$, we consider three distinct groups of input pairs: those were not routed to any of the same experts, those that have exactly one expert in common, and those that have both experts in common. Fig. 9a denotes the cosine similarity, on average, between inputs of each group. As expected, we see that inputs that are routed to both the same experts become highly similar, especially as training progresses. We also see inputs with no experts in common diverge in terms of similarity. However, inputs that have just one expert in common are still very similar, regardless of the other expert they are routed to. Moreover, Fig. 9b shows that a high percentage of inputs are highly similar — we define “highly similar” as having a cosine similarity > 0.75 . This suggests that having at least one expert in common is a consistent indicator of similarity across groups of inputs.

732 **Hyperparameters.** We use the AdamW optimizer [Loshchilov and Hutter, 2019]. We use
 733 the modified cosine learning rate schedule from Ibrahim et al. [2024]. We set the minimum
 734 learning rate to 6×10^{-5} , the max learning rate to 6×10^{-4} , and the number of warmup it-
 735 erations to 1000. We use a sequence length of 2048 and a global batch size of 1024, resulting
 736 in a global token batch size of 2^{21} . The total number of iterations is 10,000 so that we train
 737 on $20B$ tokens, roughly following the compute-optimal [Hoffmann et al., 2022] number of
 738 training tokens for a $1B$ dense model. We set the auxiliary loss [Fedus et al., 2022] to 0.01.

739 **Implementation.** We train with the gpt-neox library [Andonian et al., 2023] integrated
 740 with Megablocks [Gale et al., 2022]. The TFLOPS vary depending on the method and the
 741 number of experts chosen; for simplicity, we do not account for the router or the number
 742 of experts activated when reporting the TFLOPS, so that the number of flops we count in a
 743 forward and backward pass is the same as a dense model.

744 We plot validation results throughout training in Fig. 11.

745 C.1 Ablating Design Choices

746

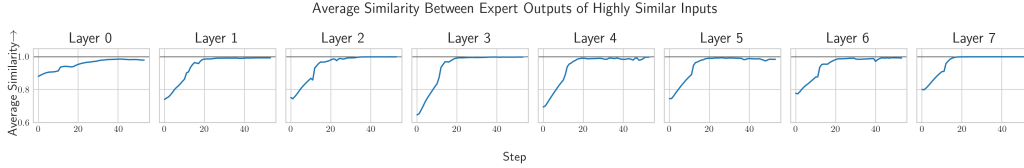


Figure 10: **Similarity of Expert Outputs With Similar Inputs.** For each expert in a model with $N = 8$ experts and $K = 2$, we consider the similarity of expert outputs when the inputs are “highly similar” i.e. with cosine similarity > 0.75 . After a few training steps, the average similarity is very high and approaches the maximum value of 1. This supports our assumption about the expert networks being Lipschitz continuous, as similar inputs indeed produce very similar expert outputs.

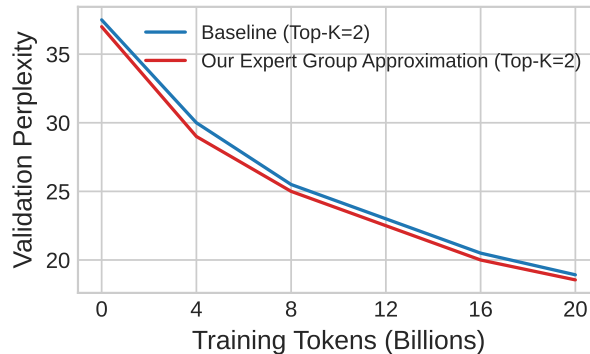


Figure 11: We plot the validation perplexity on FineWeb for the baseline Top-K router and our expert group approximation router. Without incurring significant overhead, we improve over the baseline.

747 We now present additional results on ablating our main design choices.

748 **Expert Group Approximation.** We consider two variations on the expert group approx-
 749 imation method. As a reminder, in this method we construct a mask of shape
 750 $experts, experts$ for each token. The row is the expert that token was routed to, and the
 751 column is the expert we want an approximation for. When we take the product of this
 752 mask and the router scores, we can weight each row by the probability corresponding to
 753 the expert we want an approximation for, or weight each approximation by the probability
 754 for the expert we’re using the approximation for. The former should give us more “accu-
 755 rate” approximations, because it will prioritize tokens that are more likely to be routed to
 756 the expert we want an approximation for. The latter should give us more “viable” approx-
 757 imations, because it weights by closeness to the space we’re using the approximation for.
 758 We compare these methods to the baseline in Table 2. Neither method improves over the
 759 baseline, but we think this may warrant further investigation.

Routing Method	Validation Perplexity
Expert Group Approx.	20.81
“Accurate”	20.97
“Viable”	21.14

Table 2: Ablating design choices in the expert group approximation method. Validation perplexity is reported after 12B tokens.

760 **Comparing Different Approximation Methods.** We use the Expert Group Approxima-
 761 tion method for our main results because it is lightweight, easy to implement, and provides
 762 good performance. However, the other two methods we consider also outperform the top-
 763 K ($K = 2$) baseline. Indeed, as we showed in Fig. 5, the Attention+LSH method seems

764 to obtain a better approximation of the true dense gradient. The primary reason why
 765 we report our main results with Expert Grouping is because the Expert Group Approx-
 766 imation method requires no additional memory overhead. This allows us to use larger
 767 microbatches, and therefore there are more tokens on each GPU that we can use for the
 768 approximation. In Table 3 we find that even with a microbatchsize $4\times$ smaller than that of
 769 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 3: Comparison of activated experts, routing methods, and validation perplexity after training on 20B tokens.

770 **Method Overhead.** We have already outlined the implementation of the Expert Group
 771 Approximation method, which only requires materializing two additional tensors of size
 772 $experts, experts$ and $experts, micro_batch_size \times sequence_length$. In Table 4 we compare the
 773 throughput of our method to the baseline, and find that even with an unoptimized method,
 774 we achieve 97.7% of the throughput of the baseline. We anticipate that we can further close
 775 this gap by directly modifying the gradient in the backward pass, rather than performing
 776 the approximation in the forward pass as we currently do and letting PyTorch’s autograd
 777 compute the gradient.

Routing Method	TFLOPS
Top-K (K=2)	73.4
Expert Group Approx.	71.7

Table 4: Comparing the throughput of the baseline and our method.