
Understanding Compute-Parameter Trade-offs in Sparse Mixture-of-Expert Language Models

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Scaling language model capacity is crucial for achieving better performance, as it
2 allows these models to capture more complex patterns and representations. Empir-
3 ically, increasing model size and compute improves outcomes; however, the rela-
4 tionship between model parameters and compute per example, and their combined
5 contribution to capacity, is not yet fully understood. We explore this relationship
6 through sparse Mixture-of-Expert models (MoEs), which allow scaling the num-
7 ber of parameters without proportionally increasing the FLOPs per example. We
8 investigate how varying the sparsity level, i.e., the ratio of non-active to total pa-
9 rameters, affects model performance in terms of both pretraining and downstream
10 objectives. We find that under different constraints (e.g. parameter and total train-
11 ing compute), there is an optimal level of sparsity that improves both training
12 efficiency and model performance. These results provide a clearer understanding
13 of the impact of sparsity in scaling laws for MoEs and complement existing works
14 in this area, offering insights for designing more efficient architectures.

15 1 Introduction

16 Empirical scaling laws for language model pretraining [15, 14, 19, 23, 13, 4, 27, 17] have demon-
17 strated that proportionally increasing model capacity, along with data and total compute budget,
18 consistently decreases pretraining loss, improves downstream task performance [8, 3, 1] and un-
19 locks emergent capabilities [24]. A recurring notion in these studies is that model capacity is well
20 quantified by the total number of model parameters. However, the number of parameters is not
21 the only means to increase model capacity—*compute per example* (i.e., a *fixed-sized input*), mea-
22 sured in FLOPs, also plays a significant role. In fact, several mechanisms [22, 7, 25, 12, 6] allow
23 for independent variation of the number of parameters or FLOPs per example within a model. For
24 instance, Mixture-of-Experts (MoE) models[22] introduce “FLOP-free parameters” by leveraging
25 sparsity, where only a subset of expert modules is activated for each input. Under specific con-
26 ditions, the total number of parameters can serve as a reasonable relative estimator of FLOPs per
27 example. Therefore, using the number of parameters as a measure of model capacity in scaling law
28 studies is appropriate. However, in scenarios or for architectures where the number of parameters
29 and FLOPs per example are not inherently linked, it is essential to jointly consider the effects of
30 these variables on scaling model capacity. We thus ask “*Can we draw scaling laws for the opti-
31 mal trade-off between parameter count and FLOPs per example?*” To address this question, we
32 study sparse Mixture-of-Expert Transformers (MoEs) [22, 16, 10, 28, 18] in the context of language
33 modeling.

34 Existing scaling law studies for MoEs, investigate the role of variables like number and granularity
35 of experts, underlying dense model size and inference compute in predicting the performance of the
36 models under different conditions such as training or inference compute optimality [9, 4, 27, 17]. In

37 this paper, we focus on the interaction between FLOPs per example and total parameter count, and
38 their impact on model performance in MoEs, through a large-scale empirical study.

39 We define sparsity as the ratio of inactive experts to the total number of experts, which indirectly
40 controls FLOPs per example in MoEs. We evaluate loss and downstream metrics for different spar-
41 sities, model sizes, and compute budgets terms. Our findings are summarized as follows:

- 42 • **Effect of Sparsity on Scaling Laws for Optimal Model Size:** For any specific sparsity level,
43 performance of the models as a function of their size exhibits parabolic behavior under a fixed
44 training compute budget. i.e., the model reaches its optimal performance at a vertex, that indicates
45 optimal model size. Under these conditions:
 - 46 – The optimal active number of parameters decreases as the sparsity level increases, leading
47 to smaller FLOPs per example and more efficient inference even though the total number of
48 parameters increases (see §2.1).
 - 49 – While the trend of increasing active number of parameters is similar across all training com-
50 pute budgets; the optimal active number of parameters decrease more rapidly with sparsity
51 as the training compute budget increases (see §3).
- 52 • **Optimal Sparsity for Fixed Model Size:** For any given number of parameters and under a fixed
53 training compute budget, model performance as a function of sparsity exhibits a parabolic pattern,
54 reaching its peak at an optimal sparsity level (see §2.2). Specifically, the optimal sparsity level:
 - 55 – Increases with the total number of parameters approaching 1.0 for larger models. i.e., if
56 a model is relatively small for a given training compute budget, sparsifying it more than a
57 threshold will hurt its performance. On the other hand, if a model is relatively large for a
58 given compute budget, further sparsifying it helps as it leads to increase in the number of
59 tokens the model is trained on under the given training budget constraints (see §2.2).
 - 60 – Decreases across all model sizes as the training compute budget decreases (see §D.1 and
61 §D.2).
- 62 • **Effect of Sparsity on Downstream Performance:** Models with similar pretraining perplexity
63 have similar downstream task performance regardless of sparsity. For reading comprehension
64 tasks (e.g., CoQA [21], SQuAD [20]), denser models perform better, potentially due to their
65 higher inference-time compute than a perplexity-matched sparse model. Alternative strategies to
66 increase inference time compute dynamically [25, 12] may address this gap (see §4).

67 Ultimately, this paper highlights the crucial role of the sparsity variable in scaling laws for MoEs,
68 emphasizing that the most efficient model configuration requires joint optimization of model size,
69 total training budget and sparsity level and the optimal balance between FLOPs per example and
70 parameter count in MoEs depends on both the computational constraints and the main objective.

71 2 The Interplay between Model Parameters and Sparsity in MoEs

72 Is there an optimal trade-off between parameter count and FLOPs per example in MoEs under the
73 setting where the training compute budget is fixed?

74 Intuitively, under infinite data setting, scaling model capacity along with the training compute budget
75 leads to performance improvements. Previous scaling law studies suggest that, conditioned on a
76 training compute budget measured in FLOPs denoted by C , the optimal number of parameters,
77 $N^*(C)$, exhibits a power-law relationship with C [14]:

$$N^*(C) = \arg \min_N \mathcal{L}(N; C) \propto C^a \quad (1)$$

78 Our goal is to study how to optimally trade-off FLOPs per example (i.e. a fixed-sized input) and total
79 parameters as shown in Equation 2. Instead of FLOPs vs. parameters, we investigate the relationship
80 between Sparsity S and total number of parameters N , as S is the variable in MoEs that indirectly
81 impacts FLOPs per example.¹ Essentially, for models with the same N , the model with a higher S
82 will have fewer active parameters N_a , resulting in fewer FLOPs per example. For more details on

¹We use the active number of parameters as a proxy for FLOPs per example, as $6N_aD$ provides a good estimate of the total FLOP count for MoEs; see Appendix C for details.

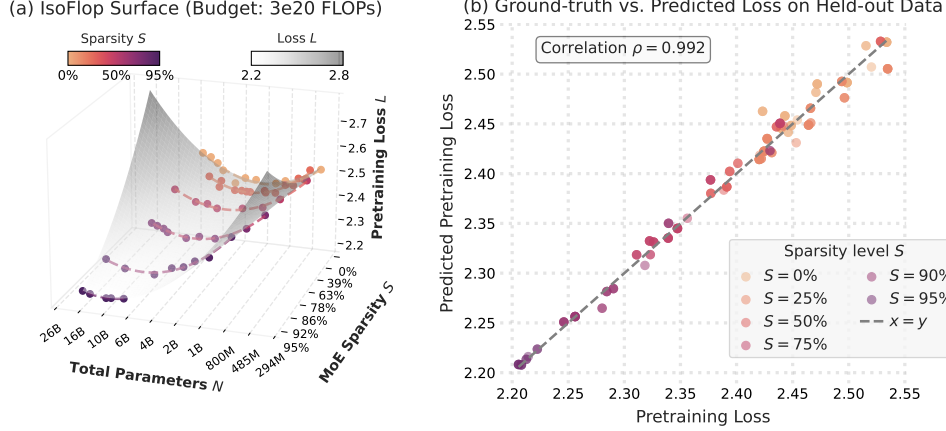


Figure 1: **IsoFLOP surface over observed pretraining loss L , model size N and sparsity S .** We fit a polynomial mapping N , S and their interaction to L using empirical data to obtain plot (a) from which we observe that for fixed compute budget the loss is decreasing with increased model sparsity. The plot on the right shows the goodness-of-fit from which we observe that the predictions and observed loss values are highly correlated with a small prediction error. (see Figure 6 in Appendix D.1 for other total training compute budgets.)

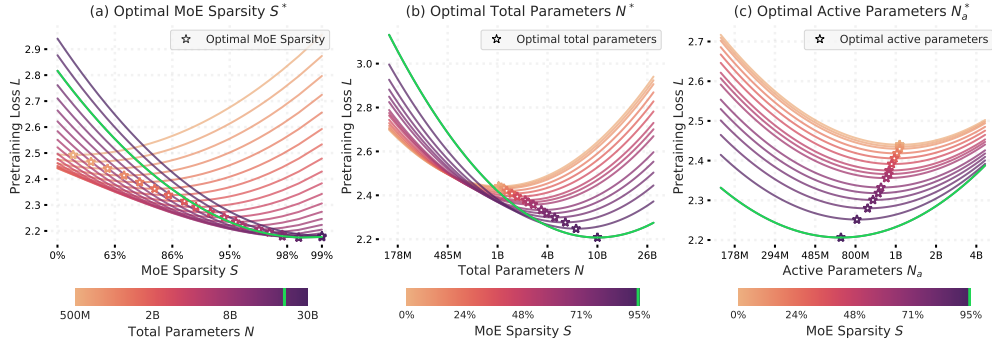


Figure 2: **IsoFLOP slices along Sparsity and Model Size.** We use fitted isoFLOP surfaces (Section 2) to analyze how sparsity S and model size N impact the loss L for a fixed compute budget. We identify optimal points by (a) fixing N and varying S , (b) fixing S and varying N and (c) fixing S and varying active parameters N_a . Observe that (a) the optimal sparsity S increases with increasing model size N and converges to 1 while (b) and (c) show that the optimal model size N and active parameter count N_a increase and decrease respectively with increasing sparsity levels. (see Figure 7 in Appendix D.1 for other total training compute budgets.)

83 the notations and experimental settings see Appendix A and Appendix B.

$$(N^*, S^*) = \arg \min_{N, S} \mathcal{L}(N, S; C) \quad (2)$$

84 To simplify the problem of understanding the joint role of N and S in predicting \mathcal{L} , we break the
85 problem, Equation 2, into two parts:

86 1. "How does the sparsity level impact the scaling laws of the relationship between N and C for
87 training-compute optimal models?" To address this question in §2.1, we fix S and vary N , study-
88 ing how optimal N and N_a change for different values of S :

$$N^* = \arg \min_N \mathcal{L}(N; C, S) \quad (3)$$

89 2. "Is there an optimal balance between total and active number of parameters under fixed training-
90 compute budget?" To address this question in §2.2, we fix N and vary S , studying how optimal
91 S changes across different values of N :

$$S^* = \arg \min_S \mathcal{L}(S; C, N) \quad (4)$$

92 As the first step, considering a fixed training compute budget C , we fit a 3D surface, referred to
93 as the IsoFLOP surface, in Figure 1a, using a polynomial function, following approach II of Hoff-
94 mann et al. [14]. We include sparsity and fit a single IsoFLOP surface across all data points, rather

95 than fitting separate curves for fixed sparsity levels or model sizes. We conducted a grid search to
96 determine the optimal polynomial degree for N , S , and the interaction term $N \times S$, finding that
97 a degree of $(2, 2, 2)$ resulted in the lowest cross-validation error. Both N and S are in log space
98 (see Appendix B for more details). Figure 1b illustrates the goodness of fit, demonstrating a strong
99 correlation and low predictive error.

100 As seen in Figure 1a, the IsoFLOP surface plot is parabolic along model size, suggesting that the
101 findings of Hoffmann et al. [14] extend to MoEs across different sparsity levels, i.e., $\mathcal{L}(N; C, S)$ is
102 parabolic, with its optimal solution located at the turning point. However, along sparsity, pretraining
103 loss decreases monotonically, indicating that, for the same compute budget, sparser models achieve
104 better pretraining loss. To better understand these observations, we examine slices of the IsoFLOP
105 surface along the axes of S and N separately in §2.1 and §2.2, respectively.

106 2.1 Optimal Model Size for Fixed Sparsity Level

107 Here we examine how sparsity influences scaling laws governing the relationship between N and C
108 for training-compute optimal models, i.e. how does N^* , for a given C, S (Equation 3), change as we
109 increase S ? Looking at slices of the IsoFLOP surface along the model size dimension, in Figure 2b
110 and (c), we observe how the IsoFLOP curves shift along loss and model size. Considering the
111 training-compute optimal model, for a fixed compute budget, loss decreases as we increase sparsity.
112 Furthermore, while sparser models have larger N compared to denser models, as seen in Figure 2b,
113 they have a smaller active parameter count N_a ; hence, fewer FLOPs per example. More parameters
114 in total increase the capacity of the sparser models to fit the data, while fewer FLOPs per example
115 allow the model to be trained with more tokens, i.e., higher D , for the same compute budget.

116 2.2 Optimal Sparsity Level for Fixed Model Size

117 Understanding the dynamics between number of parameters and FLOPs per example is essential
118 for training models with smaller inference cost under constraint training budget. This leads us to a
119 fundamental question: Is there an optimal balance between the total and active number of parameters
120 under a fixed training-compute budget? This section investigates this question. Specifically, we ask:
121 Given N and C , How does S^* change as we increase N ?

122 To address this, we look into slices of the IsoFLOP surface along sparsity, Figure 2a. As we can
123 see in this figure, given a fixed compute budget for training and fixed model size there $\mathcal{L}(S; N, C)$
124 exhibits a parabolic profile, reaching its optimum value at the vertex where $S = S^*$. We observe in
125 Figure 2a, generally, for smaller models, models with $N < N_{th}$, increasing the sparsity level, and
126 for larger models, models with $N > N_{th}$, increasing sparsity has a positive impact. More accurately,
127 for a fixed compute budget the optimal sparsity level increases with model size and converges to 1
128 as the model size grows (see Figure 8 in §D.2 in the Appendix for more details).

129 If the model size has more parameters than a threshold N_{th} it is favorable to sparsify as much as
130 possible. Note that the model with the lowest loss is not the largest sparsest model, i.e., there is a
131 compute optimal model size even after MoE sparsity is introduced, and increasing total number of
132 parameters would lead to under-training if training compute budget is fixed.

133 These results highlight the importance of balancing the number of parameters with FLOPs per ex-
134 ample. Intuitively, when the total number of parameters is small, higher sparsity results in fewer
135 active parameters, and thus fewer FLOPs per example. We speculate that this reduction in FLOPs
136 per example may lead to inefficiencies during both training and inference. Conversely, when the
137 total number of parameters is large, a fixed compute budget may not allow sufficient training on
138 enough tokens to make use of the model’s additional capacity.

139 3 Impact of Training Compute Budget on the Interaction between Model 140 Parameters and Sparsity

141 Does increasing compute budget impacts the interaction between the parameters and compute per
142 example and how they contribute to model’s capacity? In other words, does the recipe for optimally
143 increasing model capacity change as we scale up the training budget?

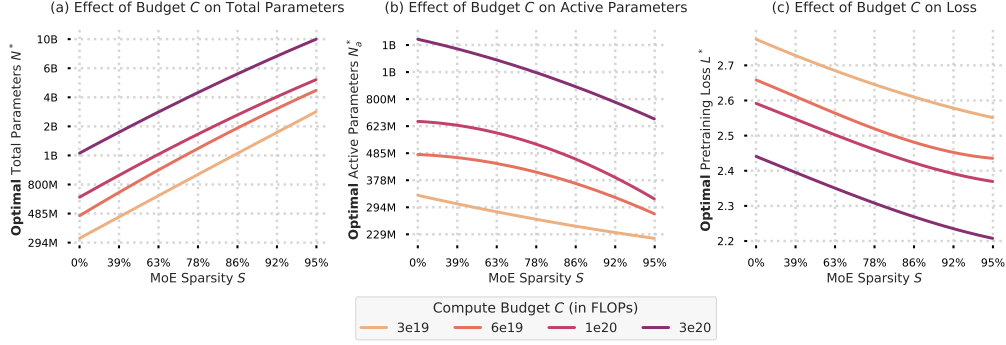


Figure 3: **Effect of compute budget on model size, number of active parameters and loss with sparsity.** Over all budgets considered, we observe that (a) the optimal model size N increases with sparsity, (b) the optimal number of active parameters N_a decreases with sparsity, and (c) the loss L decreases with sparsity.

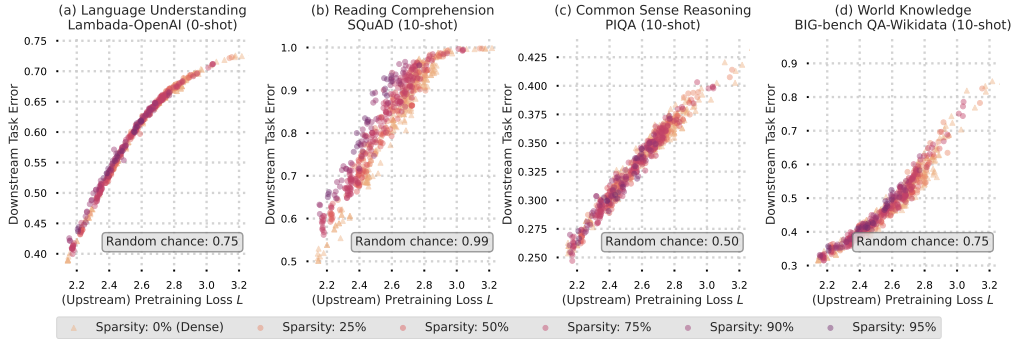


Figure 4: **Effect of sparsity on downstream vs upstream performance.** Downstream error shows a tight relationship with pretraining (“upstream”) loss across downstream tasks across all sparsity levels.

144 To answer this question, in Figure 3 we illustrate the trends for changing the total number of param-
 145 eters, N^* , the number of active parameters, N_a^* , and the loss, L^* , with sparsity level across different
 146 compute budgets.

147 Figure 3c shows that the optimal sparsity level approaches 1 across all compute budgets used in our
 148 experiments. There is no significant difference observed in the slope of the loss vs sparsity curves
 149 across different training compute budgets used in our experiments. This observation suggests that
 150 there is no diminishing effect of sparsity on the pretraining loss as we increase training compute
 151 budget, i.e., if there is no constraint on the model size, sparsity improves the performance of the
 152 model across all training budgets. As shown in §2.2, when model size in terms of total number of
 153 parameters is fixed, optimal sparsity level does not always approach 1.0, and it decreases as we
 154 increase the training compute budget (see Appendix D.2 in Figure 8).

155 Furthermore, as we see in Figures 3a and 3b, across all training compute budgets, we see a consistent
 156 trend of increasing N and decreasing N_a for compute optimal models as sparsity level increases.
 157 However, the rate of decreasing N_a , which can be interpreted as inference cost, increases as we
 158 increase training compute budget. This means, at larger training compute budgets, the benefit of
 159 reducing compute in terms of FLOPs per example amplifies. i.e., the gains of increasing sparsity
 160 level to reduce inference cost becomes more significant at larger training budgets.

161 4 Effect of MoE Sparsity on Downstream Task Performance

162 In this section, we study how sparsity affects the relationship between upstream and downstream
 163 performance of MoEs. In other words, does sparsity impact the relative gains from improvements
 164 in pretraining tasks on downstream tasks?

165 We use downstream tasks from the evaluation suite in `llm-foundry`² for benchmarking our pre-
 166 trained models. The downstream tasks are divided into four pre-defined categories namely: language

²Github repository: <https://github.com/mosaicml/llm-foundry>

167 understanding, world knowledge, reading comprehension, and symbolic reasoning to help us sys-
168 tematically test whether the downstream vs upstream performance trend remains the same or is
169 different as we vary sparsity values.

170 We observe from Figure 4a (language understanding), Figure 4c (common sense reasoning) and
171 Figure 4d (world knowledge) that there is a tight relationship between upstream (pretraining) loss
172 and downstream performance (error) across all these tasks. However, Figure 4b (reading compre-
173 hension) shows an example of a task where models with higher sparsity transfer worse compared to
174 denser models. This decrease in the transfer performance of sparser models on these tasks maybe
175 due to the lower inference-time compute in sparser models over their denser counterparts for similar
176 pretraining loss. Further analysis is needed to verify this intuition. If fewer FLOPs per example
177 is the reason behind worse transfer performance in sparser models, this effect might diminish at
178 larger total training compute budget. Moreover, one can leverage approaches like chain-of-thought
179 reasoning to independently increase FLOPs per example during inference time

180 While our results may indicate that there maybe no additional benefit obtained via sparsity in MoEs,
181 we caution the reader that this suggestion maybe an artifact of the scale of our experiments. In the
182 end, since, as shown in §2, sparser models are more efficient both in terms of training and inference
183 cost (when measured in terms of theoretical FLOPs); we can reach a better pretraining performance
184 with higher sparsity levels at a lower cost, which can translates to better downstream performance.

185 5 Conclusion

186 This paper underscores the role of sparsity in the scaling laws for Mixture-of-Expert Transformers
187 (MoEs), showing that the most efficient model configuration depends on balancing model size, train-
188 ing compute, and sparsity level. The optimal recipe for balancing FLOPs per example and parameter
189 count in MoEs depends on the objective as well other resource constraints. Our findings indicate
190 that sparsity, as a knob that controls FLOPs per example in MoEs, is a powerful mechanism for
191 optimizing model performance under constrained training compute budgets. By balancing the total
192 number of parameters, compute, and sparsity, MoEs can be scaled more effectively. These insights
193 provide valuable guidance for scaling language models, especially for MoEs, where the trade-offs
194 between parameters and FLOPs must be carefully managed.

195 MoEs were originally introduced to allow increasing model capacity without a significant increase
196 in inference cost. Our experiments show that under fixed total training compute budget increasing
197 sparsity in MoEs leads to smaller FLOPs per example, higher number of parameters, and lower pre-
198 training loss simultaneously. In other words, in the context of MoEs, if there are no constraints on
199 the total number of parameters, increasing the capacity of the model through parameter count seem
200 to be the optimal strategy if lower pretraining loss is the main goal. On the other hand, when com-
201 paring how well the pretraining performance transfers to various downstream tasks, denser models
202 seem to be better on certain types of task that potentially rely on deeper processing of the input vs
203 the knowledge stored in the parameters of the model. This potentially signals the importance of the
204 role of FLOPs per example in increasing the capacity of the model during inference. Furthermore,
205 under conditions where memory, i.e., number of total parameters, is a constraint, we find that there
206 is an optimal sparsity value that depends both on the total number of parameters and total training
207 compute budget.

208 It is also noteworthy that, in this paper, we have prioritized training compute-optimal models, in
209 contrast to many published results on large language models (LLMs), which often rely on over-
210 trained models. As a result, the performance of the models we use for the analysis in this paper is
211 not directly comparable to those of other studies, where they overtrain smaller language models, to
212 reduce the cost of inference relative to training.

213 Future work will examine the optimal balance of FLOPs per example and parameter count with more
214 emphasis and in depth analysis on performance of the models on different types of downstream
215 tasks, as well as investigating how the finding on the role of sparsity in MoEs extend to model
216 architectures or approaches with different mechanisms to change FLOPs per example and number
217 of trainable parameters of the models independently. More specifically, an interesting follow-up
218 is to investigate the scaling behaviors of the models which allow negative sparsity values (through
219 parameter sharing).

References

- [1] BIG-bench authors. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. URL <https://openreview.net/forum?id=uyTL5Bvosj>.
- [2] S. Black, S. Biderman, E. Hallahan, Q. Anthony, L. Gao, L. Golding, H. He, C. Leahy, K. McDonnell, J. Phang, et al. Gpt-neox-20b: An open-source autoregressive language model. *arXiv preprint arXiv:2204.06745*, 2022.
- [3] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei. Language models are few-shot learners. In H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc., 2020. URL https://proceedings.neurips.cc/paper_files/paper/2020/file/1457c0d6bfcb4967418bfb8ac142f64a-Paper.pdf.
- [4] A. Clark, D. d. I. Casas, A. Guy, A. Mensch, M. Paganini, J. Hoffmann, B. Damoc, B. Hechtman, T. Cai, S. Borgeaud, G. v. d. Driessche, E. Rutherford, T. Hennigan, M. Johnson, K. Millican, A. Cassirer, C. Jones, E. Buchatskaya, D. Budden, L. Sifre, S. Osindero, O. Vinyals, J. Rae, E. Elsen, K. Kavukcuoglu, and K. Simonyan. Unified scaling laws for routed language models. In *Proceedings of the 39th International Conference on Machine Learning*. PMLR, 2022.
- [5] T. Computer. Redpajama: An open source recipe to reproduce llama training dataset. <https://github.com/togethercomputer/RedPajama-Data>, Apr. 2023. Accessed: YYYY-MM-DD.
- [6] R. Csord’as, K. Irie, J. Schmidhuber, C. Potts, and C. D. Manning. Moeut: Mixture-of-experts universal transformers. *ArXiv*, abs/2405.16039, 2024. URL <https://api.semanticscholar.org/CorpusID:270063139>.
- [7] M. Dehghani, S. Gouws, O. Vinyals, J. Uszkoreit, and L. Kaiser. Universal transformers. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=HyzdRiR9Y7>.
- [8] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In J. Burstein, C. Doran, and T. Solorio, editors, *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1423. URL <https://aclanthology.org/N19-1423>.
- [9] N. Du, Y. Huang, A. M. Dai, S. Tong, D. Lepikhin, Y. Xu, M. Krikun, Y. Zhou, A. W. Yu, O. Firat, B. Zoph, L. Fedus, M. Bosma, Z. Zhou, T. Wang, Y. E. Wang, K. Webster, M. Pellat, K. Robinson, K. S. Meier-Hellstern, T. Duke, L. Dixon, K. Zhang, Q. V. Le, Y. Wu, Z. Chen, and C. Cui. Glam: Efficient scaling of language models with mixture-of-experts. *ArXiv*, abs/2112.06905, 2021. URL <https://api.semanticscholar.org/CorpusID:245124124>.
- [10] W. Fedus, B. Zoph, and N. Shazeer. Switch transformers: scaling to trillion parameter models with simple and efficient sparsity. *J. Mach. Learn. Res.*, 23(1), jan 2022. ISSN 1532-4435.
- [11] T. Gale, D. Narayanan, C. Young, and M. Zaharia. MegaBlocks: Efficient Sparse Training with Mixture-of-Experts. *Proceedings of Machine Learning and Systems*, 5, 2023.
- [12] S. Goyal, Z. Ji, A. S. Rawat, A. K. Menon, S. Kumar, and V. Nagarajan. Think before you speak: Training language models with pause tokens. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=ph04CRkPdC>.

- 271 [13] T. Henighan, J. Kaplan, M. Katz, M. Chen, C. Hesse, J. Jackson, H. Jun, T. B. Brown, P. Dhari-
272 wal, S. Gray, C. Hallacy, B. Mann, A. Radford, A. Ramesh, N. Ryder, D. M. Ziegler, J. Schul-
273 man, D. Amodei, and S. McCandlish. Scaling laws for autoregressive generative modeling.
274 *arXiv preprint arXiv: Arxiv-2010.14701*, 2020.
- 275 [14] J. Hoffmann, S. Borgeaud, A. Mensch, E. Buchatskaya, T. Cai, E. Rutherford, D. de Las Casas,
276 L. A. Hendricks, J. Welbl, A. Clark, T. Hennigan, E. Noland, K. Millican, G. van den
277 Driessche, B. Damoc, A. Guy, S. Osindero, K. Simonyan, E. Elsen, O. Vinyals, J. Rae,
278 and L. Sifre. An empirical analysis of compute-optimal large language model training.
279 In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Ad-
280 vances in Neural Information Processing Systems*, volume 35, pages 30016–30030. Curran
281 Associates, Inc., 2022. URL [https://proceedings.neurips.cc/paper_files/paper/
282 2022/file/c1e2faff6f588870935f114ebe04a3e5-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/c1e2faff6f588870935f114ebe04a3e5-Paper-Conference.pdf).
- 283 [15] J. Kaplan, S. McCandlish, T. Henighan, T. B. Brown, B. Chess, R. Child, S. Gray, A. Radford,
284 J. Wu, and D. Amodei. Scaling laws for neural language models. *CoRR*, abs/2001.08361,
285 2020. URL <https://arxiv.org/pdf/2001.08361.pdf>.
- 286 [16] D. Lepikhin, H. Lee, Y. Xu, D. Chen, O. Firat, Y. Huang, M. Krikun, N. Shazeer, and Z. Chen.
287 {GS}hard: Scaling giant models with conditional computation and automatic sharding. In
288 *International Conference on Learning Representations*, 2021. URL [https://openreview.
289 net/forum?id=qrwe7XHTmYb](https://openreview.net/forum?id=qrwe7XHTmYb).
- 290 [17] J. Ludziejewski, J. Krajewski, K. Adamczewski, M. Pióro, M. Krutul, S. Antoniak, K. Ciebiera,
291 K. Król, T. Odrzygóźdź, P. Sankowski, M. Cygan, and S. Jaszczur. Scaling laws for fine-
292 grained mixture of experts. In *ICLR 2024 Workshop on Mathematical and Empirical Un-
293 derstanding of Foundation Models*, 2024. URL [https://openreview.net/forum?id=
294 Iizr8qwH7J](https://openreview.net/forum?id=Iizr8qwH7J).
- 295 [18] N. Muennighoff, L. Soldaini, D. Groeneveld, K. Lo, J. Morrison, S. Min, W. Shi, P. Walsh,
296 O. Tafjord, N. Lambert, Y. Gu, S. Arora, A. Bhagia, D. Schwenk, D. Wadden, A. Wettig,
297 B. Hui, T. Dettmers, D. Kiela, A. Farhadi, N. A. Smith, P. W. Koh, A. Singh, and H. Hajishirzi.
298 Olmoe: Open mixture-of-experts language models, 2024. URL [https://arxiv.org/abs/
299 2409.02060](https://arxiv.org/abs/2409.02060).
- 300 [19] OpenAI. Gpt-4 technical report. *PREPRINT*, 2023.
- 301 [20] P. Rajpurkar, R. Jia, and P. Liang. Know what you don’t know: Unanswerable questions
302 for SQuAD. In I. Gurevych and Y. Miyao, editors, *Proceedings of the 56th Annual Meeting
303 of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 784–789,
304 Melbourne, Australia, July 2018. Association for Computational Linguistics. doi: 10.18653/
305 v1/P18-2124. URL <https://aclanthology.org/P18-2124>.
- 306 [21] S. Reddy, D. Chen, and C. D. Manning. CoQA: A conversational question answering chal-
307 lenge. *Transactions of the Association for Computational Linguistics*, 7:249–266, 2019. doi:
308 10.1162/tacl_a_00266. URL <https://aclanthology.org/Q19-1016>.
- 309 [22] N. Shazeer, A. Mirhoseini, K. Maziarz, A. Davis, Q. Le, G. Hinton, and J. Dean. Outra-
310 geously large neural networks: The sparsely-gated mixture-of-experts layer. In *International
311 Conference on Learning Representations*, 2017. URL [https://openreview.net/forum?
312 id=B1ckMDqlg](https://openreview.net/forum?id=B1ckMDqlg).
- 313 [23] G. Team, R. Anil, S. Borgeaud, Y. Wu, J.-B. Alayrac, J. Yu, R. Soricut, J. Schalkwyk, A. M.
314 Dai, A. Hauth, et al. Gemini: A family of highly capable multimodal models, 2024. URL
315 <https://arxiv.org/abs/2312.11805>.
- 316 [24] J. Wei, Y. Tay, R. Bommasani, C. Raffel, B. Zoph, S. Borgeaud, D. Yogatama, M. Bosma,
317 D. Zhou, D. Metzler, E. H. Chi, T. Hashimoto, O. Vinyals, P. Liang, J. Dean, and W. Fedus.
318 Emergent abilities of large language models. *Transactions on Machine Learning Research*,
319 2022. ISSN 2835-8856. URL <https://openreview.net/forum?id=yzkSU5zdWd>. Survey
320 Certification.

- 321 [25] J. Wei, X. Wang, D. Schuurmans, M. Bosma, brian ichter, F. Xia, E. H. Chi, Q. V. Le, and
322 D. Zhou. Chain of thought prompting elicits reasoning in large language models. In A. H. Oh,
323 A. Agarwal, D. Belgrave, and K. Cho, editors, *Advances in Neural Information Processing*
324 *Systems*, 2022. URL https://openreview.net/forum?id=_VjQlMeSB_J.
- 325 [26] M. Wortsman, P. J. Liu, L. Xiao, K. Everett, A. Alemi, B. Adlam, J. D. Co-Reyes, I. Gur,
326 A. Kumar, R. Novak, et al. Small-scale proxies for large-scale transformer training instabilities.
327 *arXiv preprint arXiv:2309.14322*, 2023.
- 328 [27] L. Yun, Y. Zhuang, Y. Fu, E. P. Xing, and H. Zhang. Toward inference-optimal mixture-of-
329 expert large language models. *arXiv preprint arXiv:2404.02852*, 2024.
- 330 [28] B. Zoph, I. Bello, S. Kumar, N. Du, Y. Huang, J. Dean, N. Shazeer, and W. Fedus. ST-MoE: de-
331 signing stable and transferable sparse expert models. *arXiv preprint arXiv:2202.08906*, 2022.

332

333
334

Appendices

335	A Preliminaries	11
336	A.1 Mixture-of-Expert (MoE) Transformers	11
337	B Experimental Setup	12
338	C Estimating Mixture-of-Expert (MoE) FLOPs	14
339	D Additional Analysis	16
340	D.1 Interplay between parameters and FLOPs per example	16
341	D.2 Effect of training budget and model size on optimal MoE sparsity	16
342	D.3 Effect of sparsity on downstream task performance	17
343	D.4 Comparing IsoFLOP Surface Analysis with Independent 2d IsoFLOPs	17

344
345
346

347 A Preliminaries

348 In this section, we provide a brief overview of Mixture-of-Expert (MoE) Transformers.

349 A.1 Mixture-of-Expert (MoE) Transformers

350 Mixture-of-Experts Transformers modify the standard transformer architecture by introducing in the
351 MLP layer. In this design, the experts are MLP (Multi-Layer Perceptron) modules that follow the
352 attention mechanism and are selectively activated for each token. A gating mechanism determines
353 which MLP experts are most relevant for each token, ensuring that only a subset of experts (top-k)
354 is active at any given time, while the rest remain inactive. Below, we provide the notations used
355 throughout the paper for various terms related to training MoEs.

356 **Total and Active Parameters:** In MoEs, we distinguish between total and active parameters,
357 denoted by N and N_a , respectively. The total parameter count, N , includes all parameters of the
358 network, encompassing both the experts and the rest of the architecture. The active parameter count,
359 N_a , refers to the parameters associated with the active portion of the experts, along with the rest of
360 the network that is always utilized.

361 **Top-k Expert Selection:** In MoEs, the gating mechanism assigns tokens to a subset of experts
362 using a top-k selection process, where k denotes the number of experts activated for each token. The
363 gate computes a relevance score for each expert, and the top k experts with the highest scores are
364 selected and activated. This selective activation limits the computational overhead by ensuring that
365 only a fraction of the experts are used per token.

366 **Expansion Factor and Granularity:** The expansion factor, typically denoted by E , represents
367 the increase in model capacity due to the inclusion of multiple experts, measured as a multiplicative
368 factor relative to the base dense model. The granularity, G , determines the size of each expert
369 relative to the size of the MLP module in the base dense model. The total number of experts in the
370 model is given by $E \times G$, where E scales the capacity and G controls the level of granularity.

371 **Sparsity (S):** In general, sparsity is defined as the ratio of inactive to total parameters. However,
372 in the context of MoEs, we focus on the sparsity of the MLP modules specifically. Therefore, we
373 define the sparsity level as the ratio of inactive to total experts, given by:

$$S = \frac{\text{number of non-active experts}}{\text{number of total experts}}. \quad (5)$$

374 This definition provides an interpretable measure of sparsity but cannot be directly used to calculate
375 the active parameter count N_a due to the contribution of other parameters in the model that remain
376 unsparisified.

377 B Experimental Setup

378 We train and evaluate auto-regressive sparse Mixture-of-Experts (MoE) language models of varying
379 sizes and configurations on subsets of the RedPajamaV1 dataset [5]. The key variables we explore
380 in our experiments are total model parameters N , training compute budget C , and the MoE sparsity
381 S .

382 **Pre-training data.** Our models are pre-trained on subsets of the RedPajamaV1 dataset³ [5], which
383 attempts to replicate the LLaMA pre-training data recipe and comprises 1.2 trillion tokens from
384 sources such as Common Crawl, C4, GitHub, and Wikipedia. In all our experiments, the effective
385 dataset size is adjusted based on the training compute budget C and the model size N . We tokenize
386 the data using the GPT-NeoX tokenizer [2], which has a vocabulary size of 50,432 tokens.

387 **Model and tokenizer.** We use auto-regressive transformer-based MoE language models in order
388 to study compute-parameter trade-offs by varying MoE sparsity. We use the Megablocks library [11]
389 to train dropless MoEs in which the routing mechanism ensures that all tokens are efficiently routed
390 without being dropped due to routing capacity constraints.

391 **Optimizer and scheduler.** We optimize our models using the scale-free Adam optimizer⁴ with
392 variable learning rate, a weight decay of 1×10^{-5} , and fixed Adam-specific parameters $\beta =$
393 $(0.9, 0.95)$ and $\varepsilon = 1 \times 10^{-8}$. We use a learning rate scheduler consisting of a linear warm-up
394 phase followed by a cosine decay. The warm-up phase increases the learning rate from 0 to the
395 base learning rate over a fraction of the total training steps (selected from $\{0.1, 0.05, 0.02\}$). After
396 warm-up, the learning rate decays following a cosine schedule for the remaining training steps.

397 **Fitting isoFLOP surfaces.** Recall that in Section 2, we fit isoFLOP surfaces to predict pretraining
398 loss L as a polynomial function of model size N and MoE sparsity S for a fixed training budget C .
399 The polynomial function takes the form

$$L(N, S) = a\hat{N}^\alpha + b\hat{S}^\beta + c(\hat{N} \cdot \hat{S})^\gamma + d \quad (6)$$

400 where $\hat{N} = \log N$ and $\hat{S} = -\log(1 - S)$ —we find that applying log transformations improves
401 the fit of the resulting isoFLOP surface. Through a grid search over the polynomial coefficients
402 $\alpha, \beta, \gamma \in \{0, 1, 2, 3, 4\}$, we found that the best fit was obtained for $\alpha = \beta = \gamma = 2$, i.e., a
403 quadratic polynomial. We evaluate the fitted isoFLOP surfaces in Figure 1 by (a) re-running the
404 fitting procedure $k = 100$ times on randomly sub-sampled data and (b) evaluating the Pearson
405 correlation between the true and predicted pretraining loss values on a set of held-out data points.

406 **Hyperparameters.** We fix a subset of hyperparameters for which changing values in preliminary
407 experiments (a) did not significantly improve pre-training loss, (b) the optimal value remained the
408 same across several model configurations, or (c) in order to reduce the search space (i.e., limited
409 compute resources). Specifically, we first opted to use z -router loss [28] and qk -normalization [26]
410 in order to stabilize training for large MoEs. Second, we fixed MoE router jitter noise to 0, as it did
411 not improve performance. We also fixed our batch size to 2048 for all model sizes.

412 We swept over hyperparameters that, when adjusted, (a) significantly improved pre-training loss and
413 (b) the optimal values varied across different model configurations. We increase the MoE sparsity
414 by decreasing the number of active experts and/or increasing the number of total experts. We also
415 varied the MoE granularity [17], MoE load balancing regularizer, Adam learning rate, and linear
416 warm-up steps (fraction) in order to improve pre-training loss. The table below summarizes our
417 hyperparameter sweeps:

³GitHub repository: <https://github.com/togethercomputer/RedPajama-Data>

⁴Scale-free Adam: <https://fabian-sp.github.io/posts/2024/02/decoupling/>

Table 1: Hyperparameter configurations and search spaces

Hyperparameter	Configuration	Search Space
Sparsity Level	Tuned	{0, 25, 50, 75, 90, 95, 98}%
Number of Total Experts	Tuned	Adjusted depending on sparsity
Number of Active Experts	Tuned	Adjusted depending on sparsity
Granularity	Tuned	{1, 2}
Learning Rate	Tuned	[0.003, 0.002, 0.001]
Load Balancing Factor	Tuned	{0.02, 0.05}
Warm-up Steps	Tuned	{2, 5, 10}%
Batch Size	Constant	2048
Jitter Noise	Constant	0
z-Loss	Constant	0
z-Router Loss	Constant	0.001
QK Norm	Constant	Applied

418 C Estimating Mixture-of-Expert (MoE) FLOPs

419 Similar to prior work on scaling laws (e.g., [15, 14, 17]), we use theoretical FLOP estimates as
 420 proxies for training and inference costs of language models. In this section, we (a) outline our
 421 methodology for estimating FLOPs for MoEs and (b) show that the proposed estimator closely
 422 approximates empirical FLOPs of large-scale MoEs.

423 **Setup and notation.** Consider an MoE model with n_{layers} MoE layers, each with an embedding
 424 dimension of d_{model} . We denote the number of total experts and active experts in each MoE layer
 425 by E_{total} and E_{active} respectively. Following Ludziejewski et al. [17], we let G denote the MoE
 426 granularity, which defaults to 1 and controls the size of each expert relative to the size of a feed-
 427 forward layer in an equivalent dense transformer. In our experiments, we use a vocabulary size
 428 $n_{\text{vocab}} = 50,432$ context length n_{ctx} of 2048 and use GLU modules (Gated Linear Units) [22] over
 429 feed-forward modules as the architecture of choice for MoE experts. We also set the (a) hidden
 430 dimension of each GLU expert d_{ffn} to $4 \cdot d_{\text{model}}$ and (b) instantiate MoEs where the number of
 431 attention heads n_{heads} times the dimensionality for each head d_{head} equals d_{model} , i.e., $n_{\text{heads}}d_{\text{head}} =$
 432 d_{model} .

433 **Estimating module-specific FLOPs.** To estimate the FLOPs of a given MoE model, we first
 434 individually estimate the FLOPs per token incurred by a forward *and* backward pass through every
 435 module in MoEs. Then, we aggregate these estimates to obtain the final estimator for the FLOPs per
 436 token incurred by a forward *and* backward pass through the model.

437 Like in prior work [15, 14], we take a two-step approach to estimate module-specific FLOPs. Given
 438 a module, we first estimate the number of parameters in the module and then scale this with an
 439 appropriate constant corresponding to the number of add-multiply operations per parameter through
 440 a forward and backward pass of the given module. We also omit non-leading terms such as non-
 441 linearities, biases, and layer normalization in our estimation. We estimate the FLOPs per token for
 442 attention modules, MoE routers, MoE experts, and the final un-embedding layer as follows:

- 443 1. **Attention module.** We estimate the FLOPs incurred via the QKV (and final) projections,
 444 attention logits, and attention values of all heads in a multi-head attention module as follows.
 - 445 • *QKV (and final) projections.* These projections involve $4 \cdot d_{\text{model}}n_{\text{heads}}d_{\text{heads}} = 4d_{\text{model}}^2$
 446 parameters. Following Kaplan et al. [15], we use the multiplicative constant $C = 6$ to
 447 account for the add-multiply operations per parameter in a forward and backward pass
 448 through linear modules, resulting in a FLOPs-per-token estimate of $4 \cdot C \cdot d_{\text{model}}^2$.
 - 449 • *Attention logits.* The FLOPs required to compute the attention logits for all n_{ctx} tokens
 450 equals $C \cdot n_{\text{ctx}}^2 d_{\text{model}}$ FLOPs, making the FLOP-per-token estimate equal to $C \cdot n_{\text{ctx}} d_{\text{model}}$.
 - 451 • *Attention values.* The computation of attention values requires a per-token weighted sum
 452 over $n_{\text{ctx}} d_{\text{model}}$ -dimensional vectors, making the estimate $C \cdot n_{\text{ctx}} d_{\text{model}}$.
- 453 2. **MoE module.** Given an MoE layer, we estimate the FLOPs incurred by its router and all
 454 experts separately.
 - 455 • *Router.* The MoE routing linearly maps a d_{model} -dimensional token embedding to a E_{total} -
 456 dimensional logit vector, which is subsequently used to map the token to E_{active} active
 457 experts. Following Ludziejewski et al. [17], we use a multiplicative constant $R = 14$ that
 458 accounts for the add-multiply-route operations per router parameter. The resulting FLOP
 459 estimate equals $R \cdot d_{\text{model}} E_{\text{total}}$
 - 460 • *Experts.* Each MoE experts corresponds to a GLU module [22] with $d_{\text{ffn}} = 4 \cdot d_{\text{model}}$. Since
 461 there are E_{active} active experts with granularity G , each involving three linear projections,
 462 this results in a FLOP estimate of $1/G \cdot 3 \cdot E_{\text{active}} \cdot C \cdot d_{\text{model}} d_{\text{ffn}} = 12C/G \cdot E_{\text{active}} \cdot d_{\text{model}}^2$.
- 463 3. **Un-embedding layer.** The un-embedding linear layer maps the final d_{model} -dimensional em-
 464 bedding of a token to n_{vocab} -dimensional logits, making the FLOPs-per-token $C \cdot n_{\text{vocab}} d_{\text{model}}$.

465 **Estimating MoE FLOPs.** We can aggregate the module-level FLOP estimates described above to
 466 estimate the FLOPs per token required for a single forward and backward pass through a given MoE
 467 model as follows:

$$n_{\text{layer}} \left(4C d_{\text{model}}^2 + 2C d_{\text{model}} n_{\text{ctx}} + 12C/G E_{\text{active}} d_{\text{model}}^2 + R d_{\text{model}} E_{\text{total}} \right) + C n_{\text{vocab}} d_{\text{model}}$$

468 When $E_{\text{total}}/d_{\text{model}}$ is small, which is typically the case in practice, the FLOPs induced by MoE routing
 469 can be ignored as they contribute negligibly to the estimator. This allows us to simplify the estimator
 470 to:

$$\text{MoE FLOPs per token} := C \cdot n_{\text{layers}} d_{\text{model}}^2 \left(4 + \frac{2n_{\text{ctx}}}{d_{\text{model}}} + \frac{12E_{\text{active}}}{G} + \frac{n_{\text{vocab}}}{d_{\text{model}} n_{\text{layers}}} \right) \quad (7)$$

471 **Evaluating $6N_a D$ as a FLOPs-per-token estimator in MoE Models** For standard dense trans-
 472 formers, the FLOPs are often estimated as $6ND$ [15, 14]. Given that D is fixed and not adjusted
 473 dynamically, N can serve as a reliable relative estimator of FLOPs per token for dense transformer
 474 models.

475 To adapt the $6ND$ estimator for MoE models, we replace N with N_a (the active number of
 476 parameters)—the number of parameters used in every forward and backward pass. In Figure 5,
 477 we evaluate the accuracy of the $6N_a D$ estimator by plotting the ratio between the MoE FLOPs esti-
 478 mator described in Equation 7 and $6N_a D$ as a function of model size N and a fixed context length
 479 $D = 2048$. The results show that, across all sparsity levels, the ratio remains close to one, and the
 480 gap between the two estimators decreases as model size N increases.

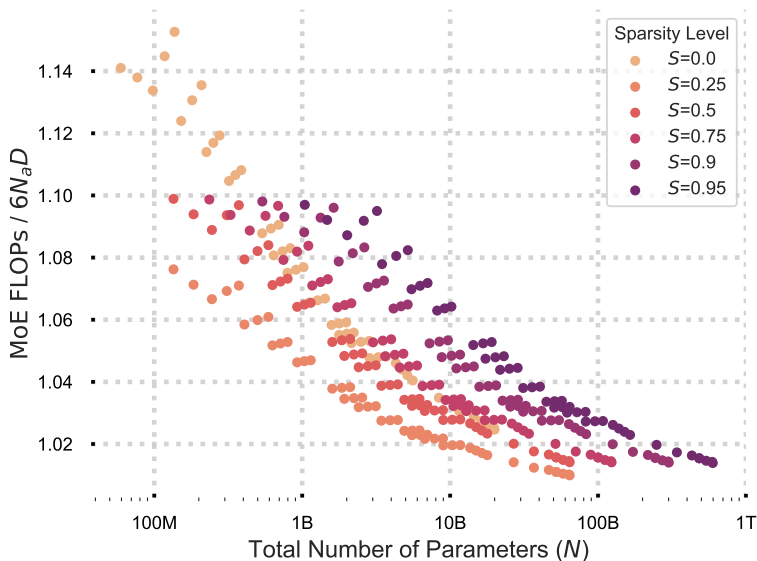


Figure 5: **Accuracy of $6N_a D$ FLOPs Estimator for MoEs.** Ratio of the MoE FLOPs estimator (Equation 7) to the $6N_a D$ estimator as a function of the total number of parameters, for a fixed context length of $D = 2048$, used in our experiments.

481 **D Additional Analysis**

482 **D.1 Interplay between parameters and FLOPs per example**

483 Recall that in Section 2, we showed that isoFLOP curves were predictive of pretraining loss for
484 different parameter counts and sparsity levels. In this section, we show similar results with additional
485 training compute budgets.

- 486 1. In Figure 6, we first show that IsoFLOP surfaces mapping model size N and sparsity level S to
487 pre-training loss L are predictive for all training compute budgets that we consider, ranging from
488 $3e19$ to $1e21$ FLOPs.
- 489 2. In Figure 7, we analyze the fitted IsoFLOP surfaces (one for each training budget) and find that
490 the (a) effect of model size N on optimal MoE sparsity S^* and (b) the effect of MoE sparsity S
491 on the optimal total and active parameters, N^* and N_a^* , is similar for all training budgets.

492 **D.2 Effect of training budget and model size on optimal MoE sparsity**

493 Recall that Section 3, we demonstrated how the relationship between optimal total parameters N^* ,
494 optimal active parameters $N_{a,s}^*$, and optimal pretraining loss L predictably changes as a function of
495 sparsity S and training budget C . In this section, we use the fitted isoFLOP surfaces to analyze how
496 the optimal MoE sparsity S^* changes as a function of total parameters N and training budget C , as
497 shown in Figure 8. Our main findings are:

- 498 • Across all training budgets (ranging from $3e19$ to $3e20$ FLOPs), increasing the total parameters
499 N leads to an increase in the optimal sparsity level S^* .
- 500 • For a fixed model size (i.e., total parameters N), increasing the training budget C generally re-
501 duces the optimal sparsity level S^* .
- 502 • The relationship between model size N and optimal S^* is not linear. For smaller models (up
503 to about $500 \cdot 10^6$ parameters), the optimal sparsity remains at 0 (i.e., dense) for most compute
504 budgets.

505 **D.3 Effect of sparsity on downstream task performance**

506 In Section 4, we analyzed the relationship between upstream pre-training loss and downstream task
507 performance across different MoE sparsity levels. We found that language understanding and world
508 knowledge tasks generally showed a strong correlation between upstream and downstream perfor-
509 mance, while reading comprehension tasks seemed to favor denser models to some extent.

510 In this section, we provide additional plots for a broader range of tasks within each category to
511 further support our findings. We consider the following tasks:

- 512 • **Common Sense Reasoning:** PIQA, CommonSenseQA, OpenBookQA, COPA
- 513 • **Language Understanding:** LAMBADA, HellaSwag, Winograd, Winogrande
- 514 • **Reading Comprehension:** SQuAD, CoQA, BoolQ
- 515 • **World Knowledge:** TruthfulQA, ARC-Easy, ARC-Challenge

516 Figure 9 shows the relationship between upstream pre-training loss and downstream task perfor-
517 mance for these additional tasks. Each row corresponds to a task category and each subplot repre-
518 sents a different task, with points colored according to MoE sparsity S . The x -axis represents the
519 upstream pre-training loss, while the y -axis shows the downstream task performance metric (usually
520 accuracy or error rate). These results supplement our main findings from Section 4:

- 521 • We observe consistent trends across tasks within each category, with language understanding and
522 world knowledge tasks showing strong correlations between upstream and downstream perfor-
523 mance regardless of sparsity.
- 524 • Reading comprehension tasks continue to show a slight advantage for denser models, while com-
525 mon sense reasoning tasks (which can be considered part of the symbolic problem-solving cate-
526 gory) show more varied relationships between upstream and downstream performance.

527 **D.4 Comparing IsoFLOP Surface Analysis with Independent 2d IsoFLOPs**

528 Recall that in Section 2, we used IsoFLOP surfaces that predict pre-training loss across varying
529 parameter counts and sparsity levels to understand how optimal sparsity and optimal model size
530 depend on each other.

531 In this section, we evaluate whether these findings remain consistent when we do not rely on fitted
532 IsoFLOP surfaces. Specifically, similar to Approach II in Hoffmann et al. [14], we directly fit
533 univariate quadratic functions that map model size N to pre-training loss L , independently for each
534 sparsity level and training compute budget. We then assess these univariate fits to determine whether
535 our findings in Section 2 hold.

- 536 • In Figure 10, each row shows how the optimal total and active parameters change as a function of
537 MoE sparsity for fixed training budgets. As in our findings from Section 2 (Figure 2), increasing
538 sparsity increases the optimal total parameters while decreasing the optimal active parameters.
539 Moreover, larger compute budgets still result in higher optimal total and active parameters, re-
540 gardless of the sparsity level.
- 541 • Furthermore, in Figure 11, we observe that across all training compute budgets, increasing sparsity
542 reduces the optimal pre-training loss. This is consistent with the trends identified in Section 3
543 (Figure 3), thereby validating our earlier results.

544

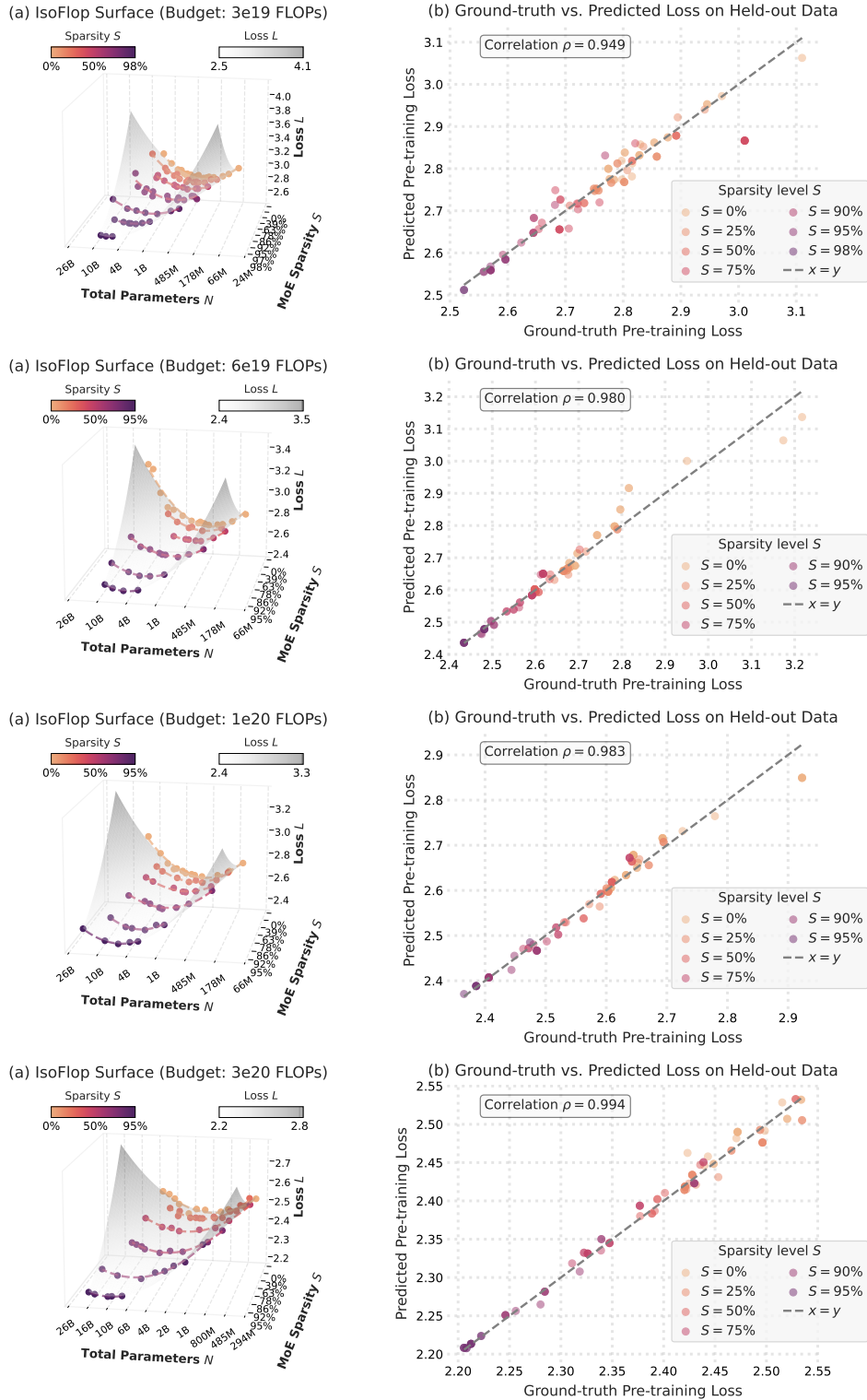


Figure 6: **IsoFLOP surfaces over total parameters N , MoE sparsity S , and pretraining loss L for different compute budgets.** The rows correspond to IsoFLOP surface fitted using models trained with a budget of 3e19, 6e19, 1e20, 3e20, and 1e21. The subplots on the left visualize IsoFLOP surfaces mapping total parameters N and sparsity level S to pretraining loss L . The subplots on the right correlate the ground-truth pretraining loss with the estimated pretraining loss on held-out data. Taken together, these results show that isoFLOP surfaces are accurate proxies for understanding how model size and MoE sparsity jointly impact pretraining loss.

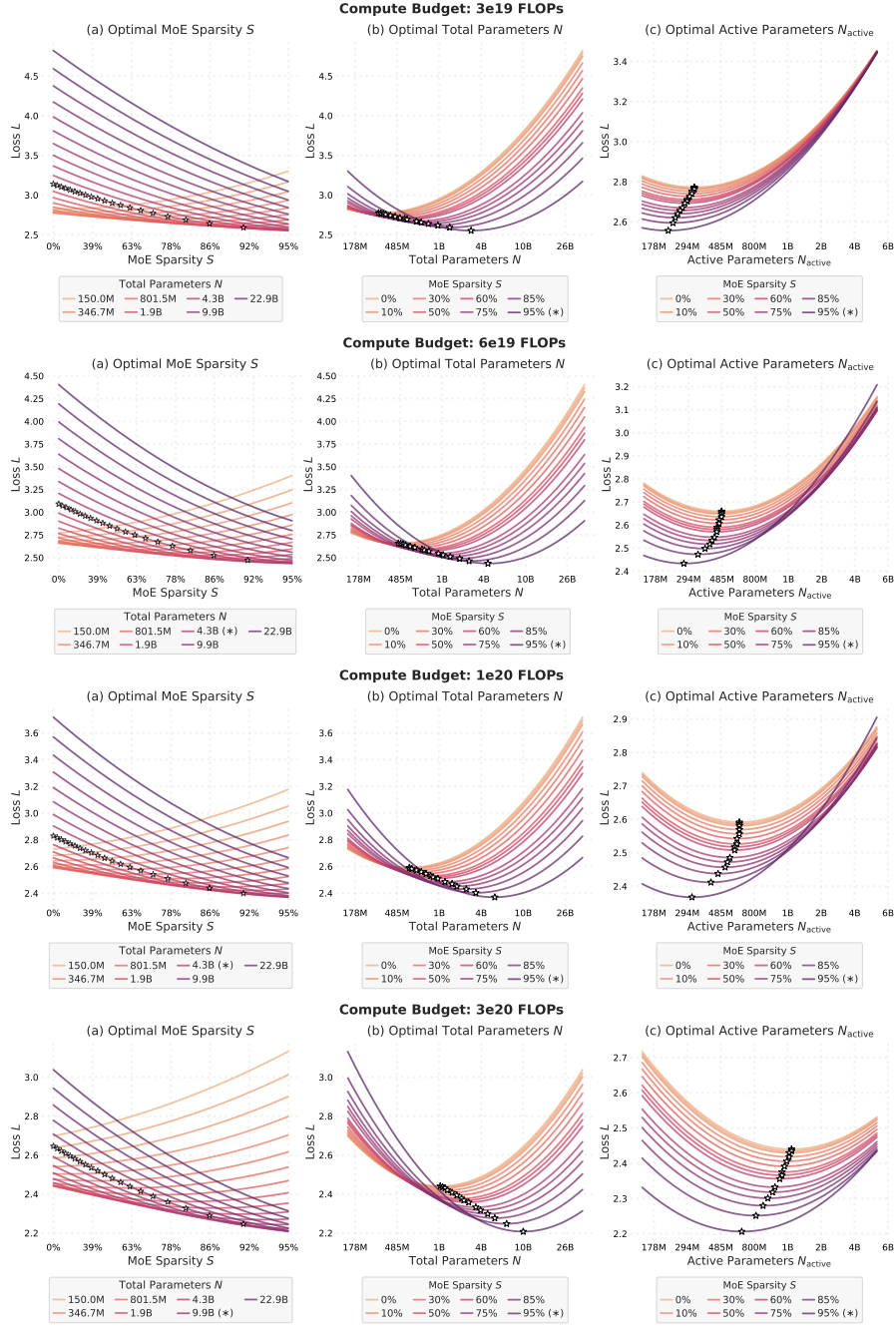


Figure 7: **Optimal MoE configurations predictably change with training compute budget.** Each row corresponds to an analysis of how optimal MoE sparsity S^* , total parameters N^* , and active parameters N_{active}^* change for a given training budget. The subplots on the left show that (a) increasing the training budget increases the model size N (denoted with black dots) with the minimum pretraining loss and (b) for models smaller than a threshold (which increases with training budget), dense models (i.e., 0% sparsity) fare better than sparse MoEs. The subplots in the second and third panel show that (a) increasing MoE sparsity increases the optimal total parameters N^* and decreases the optimal active parameters N_{active}^* . In both cases, for a fixed sparsity level, increasing the budget shifts increases the optimal total and active parameters.

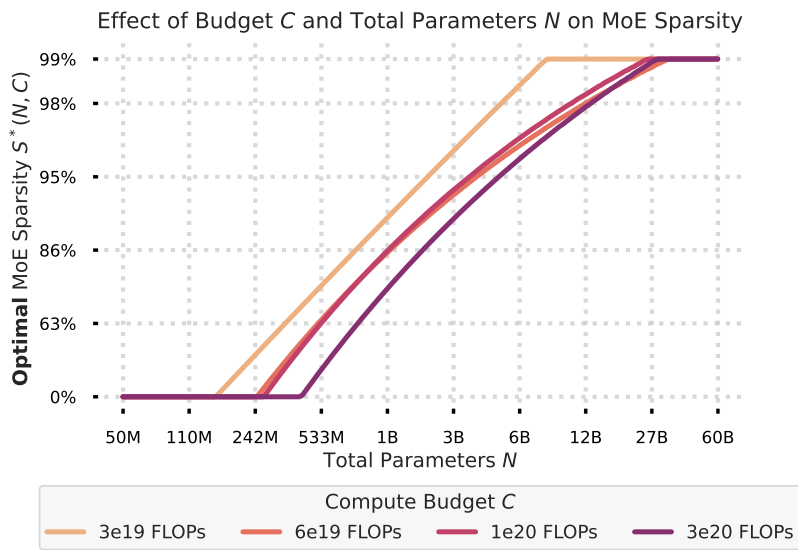


Figure 8: **Effect of training budget C and total parameters N on MoE sparsity.** This figure shows how the optimal MoE sparsity S^* changes with respect to the total number of parameters N and the training budget C . The x -axis represents the total parameters N on a logarithmic scale and the y -axis shows the optimal MoE sparsity S^* .

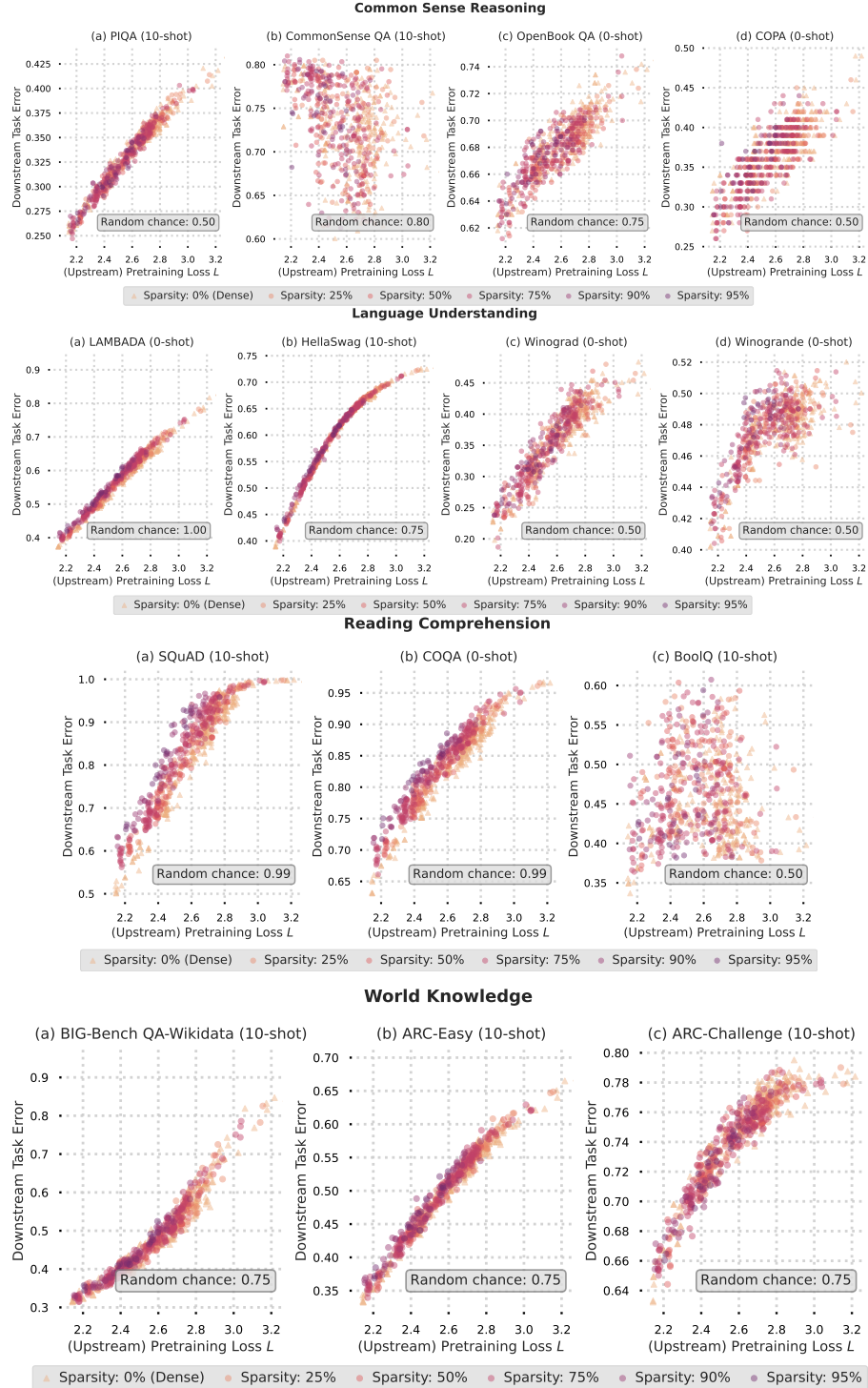


Figure 9: **Downstream task performance vs. upstream pre-training loss.** Each subplot shows the relationship between upstream pre-training loss (x-axis) and downstream task performance (y-axis) for a specific task. Similar to our results in Section 4, we find that the MoE sparsity level does not change the relationship between upstream pre-training loss and downstream task performance.

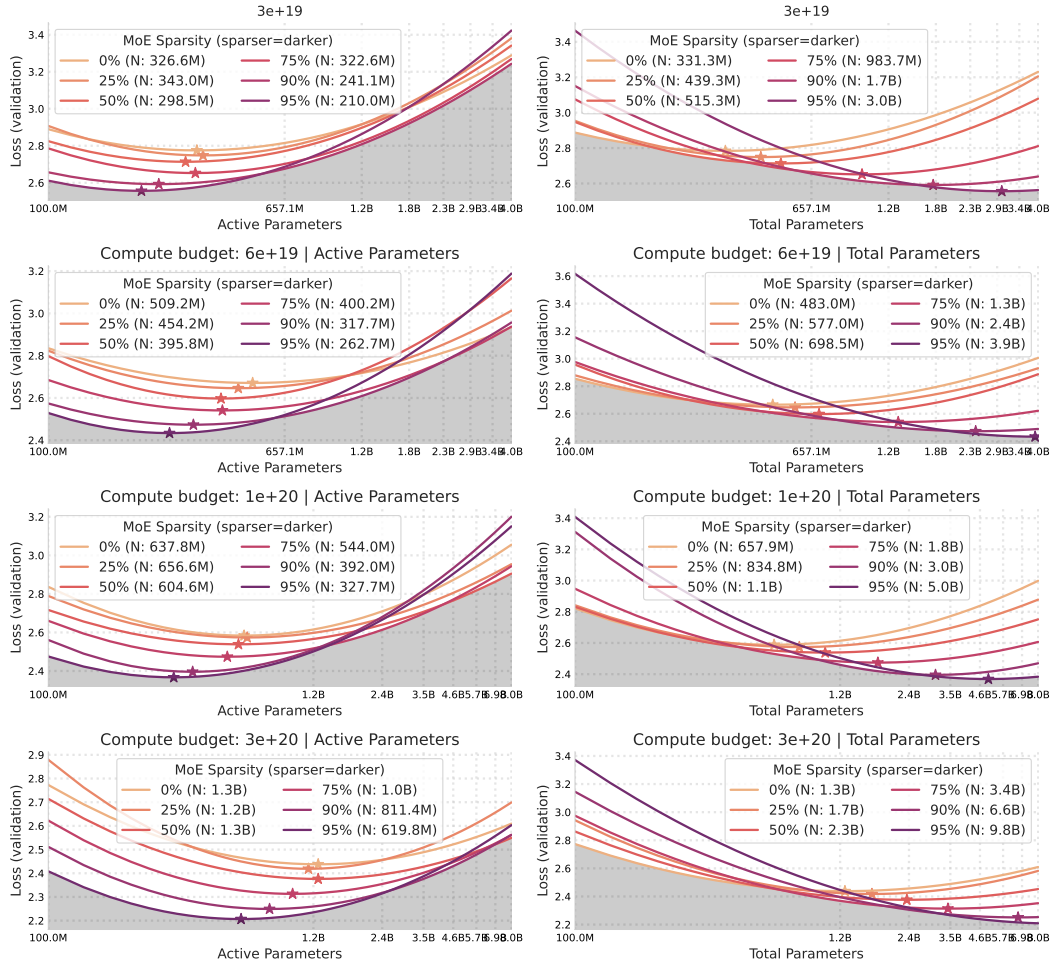


Figure 10: **Effect of MoE sparsity on optimal total and active parameters across different training compute budgets.** Each row shows the change in total and active parameters as a function of sparsity level for fixed training budgets. Increasing sparsity leads to an increase in the optimal total parameters while reducing the optimal active parameters, consistent with our findings in Section 2 (Figure 2). Larger training compute budgets result in higher optimal (total and active) parameters across all sparsity levels.

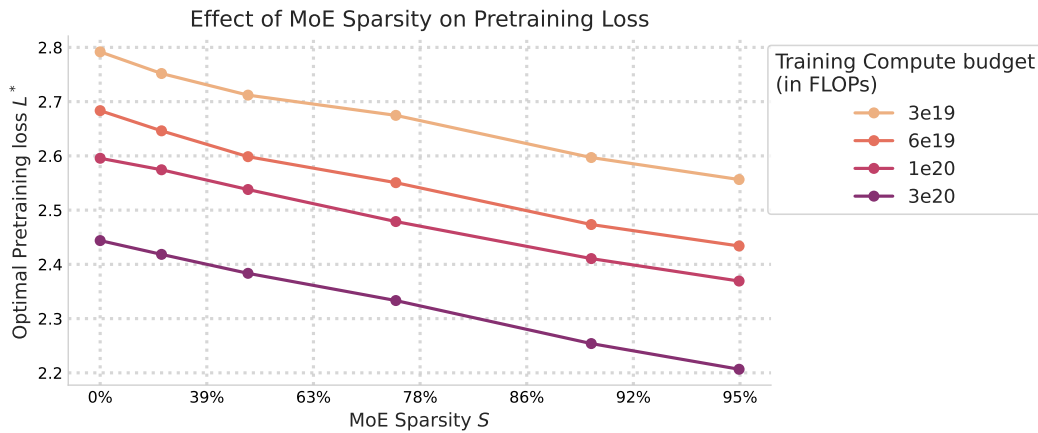


Figure 11: **Effect of MoE sparsity on pretraining loss across different training compute budgets.** As sparsity increases, the validation loss decreases for all compute budgets, with larger budgets (darker lines) achieving lower losses at each sparsity level. This trend is consistent with the findings from Section 3, demonstrating that increasing sparsity reduces the optimal pretraining loss across all compute budgets.