

---

# Language models scale reliably with over-training and on downstream tasks

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1       Scaling laws are useful guides for derisking expensive training runs, as they predict  
2       performance of large models using cheaper, small-scale experiments. However,  
3       there remain gaps between current scaling studies and how language models are  
4       ultimately trained and evaluated. For instance, scaling is usually studied in the  
5       compute-optimal training regime (i.e., “Chinchilla optimal” regime). In contrast,  
6       models are often over-trained to reduce inference costs. Moreover, scaling laws  
7       mostly predict loss on next-token prediction, but models are usually compared on  
8       downstream task performance. To address both shortcomings, we create a testbed  
9       of 104 models with 0.011B to 6.9B parameters trained with various numbers of  
10       tokens on three data distributions. First, we fit scaling laws that extrapolate in both  
11       the amount of over-training and the number of model parameters. This enables us  
12       to predict the validation loss of a 1.4B parameter, 900B token run (i.e.,  $32\times$  over-  
13       trained) and a 6.9B parameter, 138B token run (i.e., a compute-optimal run)—each  
14       from experiments that take  $300\times$  less compute. Second, we relate the perplexity of  
15       a language model to its downstream task performance by proposing a power law.  
16       We use this law to predict top-1 error averaged over downstream tasks for the two  
17       aforementioned models, using experiments that take  $20\times$  less compute.

## 18 **1 Introduction**

19       Training large language models is expensive. Furthermore, training high-quality models requires a  
20       complex recipe of algorithmic techniques and training data. To reduce the cost of finding successful  
21       training recipes, researchers first evaluate ideas with small experiments and then extrapolate their  
22       efficacy to larger model and data regimes via scaling laws. With reliable extrapolation, it is possible  
23       to quickly iterate at small scale and still pick the method that will perform best for the final large  
24       training run. Indeed, this workflow has become commonplace for training state-of-the-art language  
25       models like Chinchilla 70B [45], PaLM 540B [19], GPT-4 [76], and many others.

26       Despite their importance for model development, published scaling laws differ from the goals of  
27       training state-of-the-art models in important ways. For instance, scaling studies usually focus on the  
28       compute-optimal training regime (“Chinchilla optimality” [45]), where model and dataset size are set  
29       to yield minimum loss for a given compute budget. However, this setting ignores inference costs.  
30       As larger models are more expensive at inference, it is now common practice to over-train smaller  
31       models [113]. Another potential mismatch is that most scaling laws quantify model performance by  
32       perplexity in next-token prediction instead of accuracy on widely used benchmark datasets. However,  
33       practitioners usually turn to benchmark performance, not loss, to compare models.

34       In this paper, we conduct an extensive set of experiments to address both scaling in the over-trained  
35       regime and benchmark performance prediction.

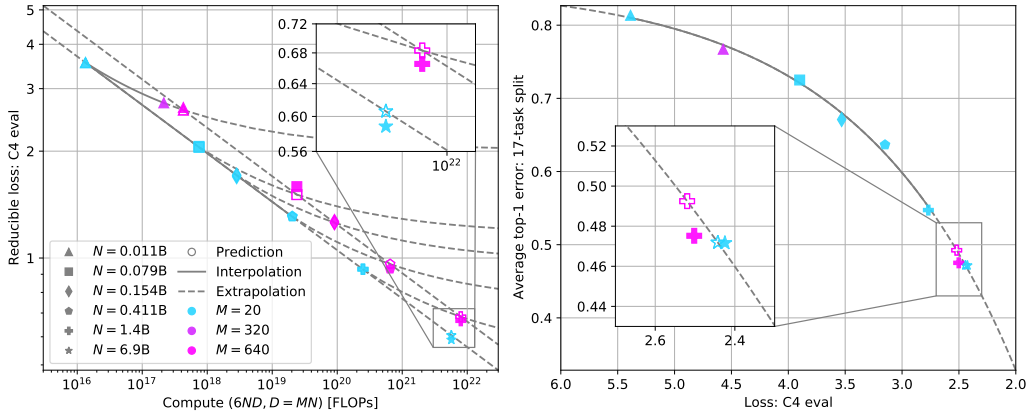


Figure 1: **Reliable scaling with over-training and on downstream error prediction.** (*left*) We fit a scaling law for model validation loss, parameterized by (i) a token multiplier  $M = N/D$ , which is the ratio of training tokens  $D$  to parameters  $N$  and (ii) the compute  $C$  in FLOPs used to train a model, approximated by  $C = 6ND$ . Larger values of  $M$  specify more over-training. We are able to extrapolate, in both  $N$  and  $M$ , the validation performance of models requiring more than  $300\times$  the training compute used to construct the scaling law. (*right*) We also fit a scaling law to predict average downstream top-1 error as a function of validation loss. We find that fitting scaling laws for downstream error benefits from using more expensive models when compared to fitting for loss prediction. We predict the average error over 17 downstream tasks for models trained with over  $20\times$  the compute. For this figure, we train all models on RedPajama [112].

36 Motivated by the practice of training beyond compute-optimality, we first investigate whether scaling  
 37 follows reliable trends in the over-trained regime. We notice, as implied by Hoffmann et al. [45], for a  
 38 set of models of different sizes trained with a constant ratio of tokens to parameters, models’ reducible  
 39 loss  $L'$  [43, 45] follows a power law ( $L' = \lambda \cdot C^{-\eta}$ ) in the amount of training compute  $C$ . We  
 40 find that as one increases the ratio of tokens to parameters, corresponding to more over-training, the  
 41 scaling exponent  $\eta$  remains about the same, while the scalar  $\lambda$  changes. We explain our observations  
 42 by reparameterizing existing scaling laws in relation to the amount of over-training.

43 To establish empirically that scaling *extrapolates* in the over-trained regime, we further experiment  
 44 with a testbed of 104 models, trained from scratch on three different datasets: C4 [88, 27],  
 45 RedPajama [112], and RefinedWeb [82]. We find that scaling laws fit to small models can accurately  
 46 predict the performance of larger models that undergo more over-training. Figure 1 (*left*) illustrates our  
 47 main over-training result, where we invest  $2.4e19$  FLOPs to extrapolate the C4 validation performance  
 48 of a 1.4B parameter model trained on 900B tokens, which requires  $300\times$  more compute to train.

49 In addition to over-training, we also investigate if scaling laws can predict the performance of a  
 50 model on downstream tasks. We establish a power law relationship between language modeling  
 51 perplexity and the average top-1 error on a suite of downstream tasks. While it can be difficult to  
 52 predict the error on individual tasks, we find it possible to predict aggregate performance from a  
 53 model’s perplexity among models trained on the same training data. Figure 1 (*right*) presents our  
 54 main downstream error prediction result, where we invest  $2.7e20$  FLOPs to predict the average top-1  
 55 error over a set of downstream tasks to within 1 percentage point for a 6.9B compute-optimal model,  
 56 which requires  $20\times$  more compute to train.

57 Our results suggest that the proposed scaling laws are promising to derisk (i) the effects of over-  
 58 training models and (ii) the downstream performance of scaling up training recipes. To facilitate  
 59 further research on reliable scaling, we will release all experiments and models.

## 60 2 Developing scaling laws for over-training and downstream tasks

61 In this section, we develop scaling laws to predict over-trained and downstream performance. First,  
 62 we provide key definitions (Section 2.1). We next present a scaling law for over-training drawing on  
 63 empirical observation and prior work (Section 2.2). To connect loss scaling and downstream error  
 64 prediction, we observe that average top-1 error decreases exponentially as a function of validation loss,

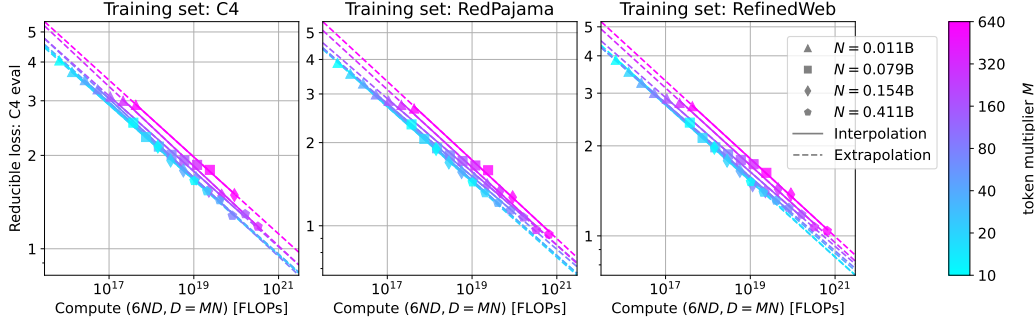


Figure 2: **Scaling in the over-trained regime follows consistent power law exponents.** We notice parallel lines in the log-log plots of reducible loss vs. training compute for a range of token multipliers  $M$ , which give the ratio of training tokens to model parameters. Larger  $M$  corresponds to more over-training. For a power law giving reducible loss as a function of compute:  $L'(C) = \lambda \cdot C^{-\eta}$ , the exponent  $\eta$  remains relatively constant resulting in lines with approximately fixed slope (Figure 17). The scalar  $\lambda$  that determines the  $y$ -intercept, however, shifts with different token multipliers. This suggests  $\lambda$  is a function of the token multiplier, while  $\eta$  is not.

65 which we formalize as a novel scaling law (Section 2.3). In later sections, we build an experimental  
 66 setup (Section 3) to quantify the extent to which our scaling laws extrapolate reliably (Section 4).

## 67 2.1 Preliminaries

68 **Scaling laws for loss.** Typically, scaling laws predict model loss  $L$  as a function of the compute  
 69  $C$  in FLOPs used for training. If one increases the number of parameters  $N$  in a model or the  
 70 number of tokens  $D$  that a model is trained on, compute requirements naturally increase. Hence, we  
 71 assume  $C$  is a function of  $N, D$ . Following Kaplan et al. [51], we use the approximation  $C = 6ND$ ,  
 72 which Hoffmann et al. [45] independently verify. We consider,

$$L(C) = E + L'(C), \quad (1)$$

73 where  $E$  is an *irreducible loss* and  $L'$  is the *reducible loss*.  $E$  captures the Bayes error or minimum  
 74 possible loss achievable on the validation domain. The  $L'(C)$  term captures what can possibly be  
 75 learned about the validation domain by training on a source domain.  $L'(C)$  should approach zero  
 76 with increased training data and model capacity.  $L'(C)$  is often assumed to follow a power law:  
 77  $L'(C) = \lambda \cdot C^{-\eta}$  (i.a., Hestness et al. [43], OpenAI [76]). It is also often helpful to consider a power  
 78 law in a log-log plot, where it appears as a line with slope  $-\eta$  and  $y$ -intercept  $\log(\lambda)$ .

79 **Token multipliers.** We define a token multiplier  $M = D/N$  as the ratio of training tokens to model  
 80 parameters for notational convenience.  $M$  allows us to consider fixed relationships between  $D$  and  
 81  $N$  even as a model gets bigger (i.e., as  $N$  becomes larger).

82 **Compute-optimal training.** Hoffmann et al. [45] establish compute-optimal training, where, for  
 83 any compute budget  $H$ , the allocation of parameters and tokens is given by,

$$\arg \min_{N, D} L(N, D) \text{ s.t. } C(N, D) = H. \quad (2)$$

84 To solve for the optimal  $N^*, D^*$ , one can sweep  $N, D$  for each compute budget, retaining the  
 85 best configurations. Hoffmann et al. [45] find that as the compute budget increases,  $N^*$  and  $D^*$   
 86 scale roughly evenly. Assuming equal scaling, there is a fixed compute-optimal token multiplier  
 87  $M^* = D^*/N^*$  per training distribution.

88 **Over-training.** We define over-training as the practice of allocating compute sub-optimally, so  
 89 smaller models train on a disproportionately large number of tokens (i.e.,  $M > M^*$ ). While loss  
 90 should be higher than in the compute-optimal allocation for a given training budget, the resulting  
 91 models have fewer parameters and thus incur less inference cost.

## 92 2.2 Scaling laws for over-training

93 To propose a scaling law for over-trained models, we first turn to empirical observation. We train four  
 94 model configurations with parameter counts between 0.011B and 0.411B for token multipliers  $M$

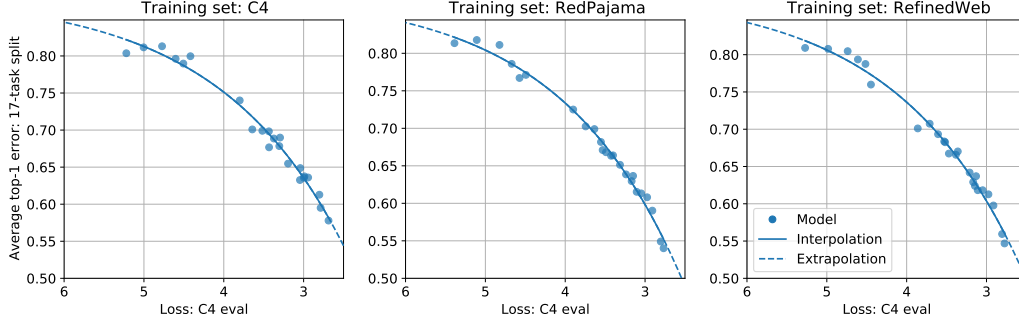


Figure 3: **Average top-1 error scales as a function of loss.** We plot models trained on three datasets and notice an exponential decay of average top-1 error as C4 eval loss, on the x-axis, decreases. We consider on the y-axis average error on 17 evaluations where performance is at least 10 points above random chance for at least one 0.154B scale model. These observations suggest that average top-1 error should be predictable with reliable loss estimates.

95 between 20 and 640, where  $M = 20$  points lie roughly on the compute-optimal frontier, and larger  
 96  $M$  corresponds to more over-training. We defer experimental details to Section 3 to focus on our  
 97 observations first. In Figure 2, we show loss against compute in a log-log plot for the models trained  
 98 on three datasets and evaluated on the C4 eval set. We notice parallel lines when fitting power laws to  
 99 the reducible loss, which suggests a near-constant scaling exponent even with increased over-training.  
 100 This indicates that scaling behavior should be describable in the amount of over-training.

101 In search of an analytic expression for the observations in Figure 2, we consider existing scaling  
 102 literature. A common functional form for the risk of a model, as proposed in prior work [93, 45] is,

$$L(N, D) = E + AN^{-\alpha} + BD^{-\beta}. \quad (3)$$

103 Recall from Section 2.1,  $N$  is the number of parameters and  $D$  the number of training tokens. The  
 104 constants  $E, A, \alpha, B, \beta$  are fit from data. By fitting this parametric form, Hoffmann et al. [45]  
 105 find that scaling exponents  $\alpha$  and  $\beta$  are roughly equal, suggesting that one should scale  $N$  and  $D$   
 106 equally as compute increases. Hence, we assume  $\alpha = \beta$ . With this assumption, we reparameterize  
 107 Equation (3) in terms of compute  $C = 6ND$  and a token multiplier  $M = D/N$ . We get,

$$L(C, M) = E + (aM^\eta + bM^{-\eta}) C^{-\eta}, \quad (4)$$

108 where  $\eta = \alpha/2$ ,  $a = A(1/6)^{-\eta}$ ,  $b = B(1/6)^{-\eta}$  gives the relation to Equation (3). For a complete  
 109 derivation, see Appendix A.

110 Equation (4) has the following interpretation: (i) The scaling exponent  $\eta$  is not dependent on  $M$ .  
 111 Thus, we always expect lines with the same slope in the log-log plot—as in Figure 2. (ii) The term  
 112  $aM^\eta + bM^{-\eta}$  determines the offsets between curves with different token multipliers. Hence, we  
 113 expect non-overlapping, parallel lines in the log-log plot for the range of  $M$  we consider—also  
 114 consistent with Figure 2.

115 Recall that we make the assumption  $\alpha = \beta$ , which implies equal scaling of parameters and tokens  
 116 as more compute is available. However, as explained in Appendix A, even if  $\alpha \neq \beta$ , we get a  
 117 parameterization that implies the power-law exponent remains constant with over-training.

### 118 2.3 Scaling laws for downstream error

119 Scaling is typically studied in the context of loss [51, 45, 72], which Schaeffer et al. [100] note  
 120 is smoother than metrics like accuracy. However, practitioners often use downstream benchmark  
 121 accuracy as a proxy for model quality and not loss on perplexity evaluation sets. To better connect  
 122 scaling laws and over-training to task prediction, we revisit the suite of models plotted in Figure 2. In  
 123 Figure 3, we plot average downstream top-1 errors over evaluations sourced from LLM-Foundry [69]  
 124 against the C4 eval loss. We defer details of the setup to Section 3 to focus here on a key observation:  
 125 average error appears to follow exponential decay as loss decreases.

126 Based on the exponential decay we observe in Figure 3, we propose the following relationship  
 127 between downstream average top-1 error  $\text{Err}$  and loss  $L$ ,

$$\text{Err}(L) = \epsilon - k \cdot \exp(-\gamma L), \quad (5)$$

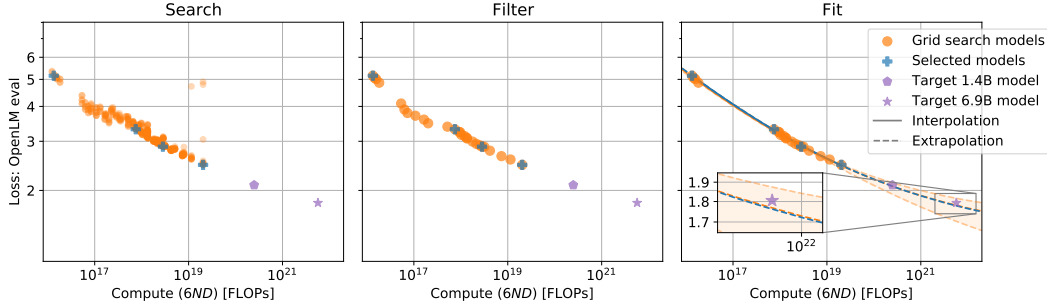


Figure 4: **Search, filter, fit: A recipe for selecting configurations for scaling.** (*left*) To generate the final configurations presented in Table 3, we run a 435 model grid search over model width, hidden dimension, number of attention heads, batch size, and warmup steps. All models are trained near compute-optimally. (*center*) We plot the efficient frontier of models, which appear to follow a trend, excluding models from  $5.2 \times 10^{16}$  to  $5.2 \times 10^{17}$ , which fall below the trend. (*right*) We fit a power law with irreducible error to the remaining configurations, picking four configurations that closely track the full model suite (“Selected models”). These models extrapolate the performance of 1.4B, 6.9B target models. Shaded regions represent bootstrap 95% confidence intervals.

128 where  $\epsilon, k, \gamma$  are fit from data. Equation (5) also has an interpretation in terms of model perplexity  
 129  $\text{PP}(L) = \exp(L)$ ,

$$\text{Err}(\text{PP}) = \epsilon - k \cdot \text{PP}^{-\gamma}. \quad (6)$$

130 Namely,  $\text{Err}$  follows a power law in  $\text{PP}$  that is bounded from above by  $\epsilon$  signifying arbitrarily high  
 131 error and from below by  $\epsilon - k \cdot \exp(-\gamma E)$ , where  $E$  is the Bayes error from Equation (4).

132 Equation (5) in conjunction with Equation (4) suggests a three-step method to predict  $\text{Err}$  as a function  
 133 of compute and the amount of over-training. For choices of training and validation distributions, (i)  
 134 fit a scaling law to Equation (4) using triplets of compute  $C$ , token multiplier  $M$ , and measured loss  
 135  $L$  on a validation set to yield  $(C, M) \mapsto L$ . (ii) Fit a scaling law to Equation (5) using pairs of loss  $L$   
 136 and downstream error  $\text{Err}$  for models to get  $L \mapsto \text{Err}$ . (iii) Chain predictions to get  $(C, M) \mapsto \text{Err}$ .

### 137 3 Constructing a scaling testbed

138 In this section, we discuss our experimental setup to test the predictions suggested by Equations (4)  
 139 and (5). We first present our general language modeling setup (Section 3.1). Next, we discuss our  
 140 strategy for determining model configurations for our scaling investigation (Section 3.2) and fitting  
 141 scaling laws (Section 3.3). We then present metrics to validate how well scaling laws predict loss and  
 142 downstream performance (Section 3.4).

#### 143 3.1 Training setup

144 We train transformers [116] for next token prediction, based on architectures like GPT-2 [85] and  
 145 LLaMA [113]. We employ GPT-NeoX [15] as a standardized tokenizer for all data. See Appendix B  
 146 for architecture, optimization, and hyperparameter details.

#### 147 3.2 Model configurations

148 To get final configurations for the 0.011B to 0.411B parameter models plotted in Figures 2 and 3, we  
 149 first conduct a wide grid search over a total of 435 models, trained from scratch, from 0.01B to 0.5B  
 150 parameters (Figure 4 (*left*)). We train on the original OpenLM data mix [39], which largely consists  
 151 of RedPajama [112] and The Pile [31]. While we eventually plan to over-train models, at this step  
 152 we search for *base configurations* near compute-optimality. We train on 20 tokens per parameter  
 153 ( $M = 20$ ), which, in early experiments, gives models near the compute-optimal frontier. This is  
 154 similar to findings in Hoffmann et al. [45]’s Table 3, which suggests that  $M = 20$  is near-optimal for  
 155 the Chinchilla experimental setup.

Table 1: **Default number of parameters  $N$  and token multiplier  $M$  to fit our scaling laws.** We invest  $\sim 100$  A100 hours to fit Equation (4) and  $\sim 1,000$  A100 hours to fit Equation (5).

$N$	$M$	Used to fit Equation (4)	Used to fit Equation (5)
0.011B	20	✓	✓
0.079B	20	✓	✓
0.154B	20	✓	✓
0.411B	20	✓	✓
0.011B	320	✓	✓
1.4B	20	✗	✓
Total compute $C$ [FLOPs]		$2.4e19$	$2.7e20$

156 To find maximally performant small-scale models on validation data, we tune model width, number  
 157 of layers, number of attention heads, warmup steps, and batch size. Our validation set, OpenLM  
 158 eval, contains tokens from recent arXiv papers, the OpenLM codebase itself, and news articles. We  
 159 find in early experiments that qk-LayerNorm makes models less sensitive to learning rate, which  
 160 is a phenomenon Wortsman et al. [123] report in their Figure 1. Hence, we fix the learning rate  
 161 ( $3e-3$ ) for our sweeps. We also perform smaller grid searches over 1.4B and 6.9B parameter model  
 162 configurations at  $M = 20$ , retaining the best configurations.

163 At this point, we have many models, several of which give poor performance; following prior  
 164 work [51, 45], we want to keep only models that give best performance. Hence, in Figure 4 (*center*),  
 165 we filter out models that do not lie on the Pareto frontier. While there appears to be a general trend,  
 166 configurations between  $5.2 \times 10^{16}$  and  $5.2 \times 10^{17}$  FLOPs lie below the frontier established by other  
 167 models. We hypothesize these models over-perform as they are trained for more optimization steps  
 168 than their neighbors based on our power-of-two batch sizes. We provide support for this hypothesis  
 169 in Appendix E, but opt to remove these models from our investigation.

170 To ensure tractable compute requirements for our scaling experiments, we require a subset of models  
 171 that follows the trend of the entire Pareto frontier. In Figure 4 (*right*), we fit trends to the Pareto  
 172 models and to a subset of four models. We notice that the trends closely predict both the performance  
 173 of the 1.4B and 6.9B models, suggesting that our small-scale configurations reliably extrapolate in  
 174 the compute-optimal setting.

175 Moving forward, we do not tune hyperparameters for other token multipliers (i.e.,  $M \neq 20$ ), on  
 176 other training or evaluation distributions, or on validation sets for downstream tasks. For more details  
 177 including specific hyperparameters, see Appendix C.

178 To create our scaling testbed, we start with the four small-scale, base configurations from our  
 179 grid search:  $N \in \{0.011B, 0.079B, 0.154B, 0.411B\}$ . To ensure our conclusions are not particular  
 180 to a single training distribution, we train models on each of C4 [88, 27], RedPajama [112], and  
 181 RefinedWeb [82], which have 138B, 1.15T, and 600B tokens, respectively, for different token  
 182 multipliers  $M \in \{5, 10, 20, 40, 80, 160, 320, 640\}$ . We omit runs that require more tokens than are  
 183 present in a dataset (i.e.,  $N = 0.411B, M = 640$  for C4). We additionally train  $N = 1.4B$  models at  
 184  $M = 20$  and at the largest token multiplier possible without repeating tokens (i.e., 80 for C4, 640 for  
 185 RedPajama, and 320 for RefinedWeb). We train  $N = 6.9B, M = 20$  models on each dataset given  
 186 the relevance of 7B parameter models [113, 49]. In total this results in a testbed of 104 models.

### 187 3.3 Fitting scaling laws

188 We fit Equation (4) to approximate  $E, a, b, \eta$  using curve-fitting in SciPy [117] (i.e., Levenberg-  
 189 Marquardt to minimize non-linear least squares). We repeat this process to fit Equation (5) to  
 190 approximate  $\epsilon, k, \gamma$ . We invest  $\sim 100$  A100 hours to train the models required to fit a scaling law for  
 191 loss and  $\sim 1,000$  A100 hours for a corresponding law for downstream error. Unless otherwise specified,  
 192 we fit to the  $N, M$  pairs in Table 1, which are a subset of our full testbed. Our configurations allow  
 193 us to test for extrapolation to the  $N = 1.4B, M = 640$  (900B token) and the  $N = 6.9B, M = 20$   
 194 (138B token) regimes.

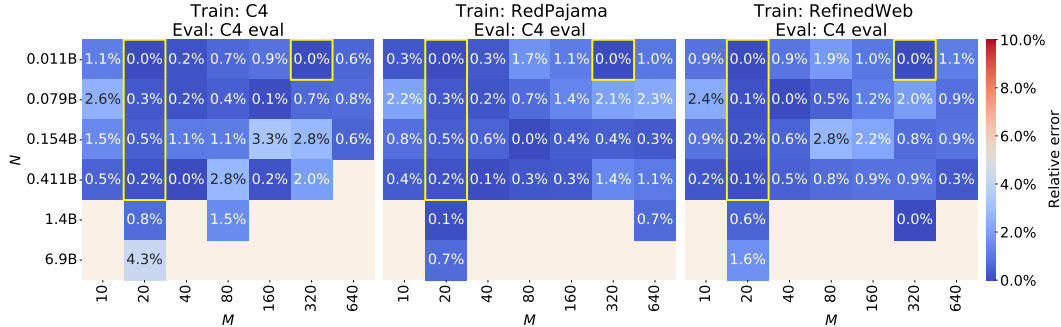


Figure 5: **Relative error on C4 eval for different training distributions.** Boxes highlighted in yellow correspond to pairs—number of parameters  $N$ , token multiplier  $M$ —used to fit Equation (4). Larger values of  $M$  correspond to more over-training. The prediction error is low in both interpolation and extrapolation ranges. Below  $N = 1.4\text{B}$ , empty squares correspond to runs that were not possible due to the limited dataset size for single epoch training. At  $N = 1.4\text{B}$  we run at  $M = 20$  and at the largest possible multiplier. At  $N = 6.9\text{B}$ , we run at  $M = 20$ .

### 195 3.4 Evaluation setup

196 **Evaluation datasets.** Unless otherwise stated, our default validation loss dataset is C4 eval. For  
 197 downstream tasks, we adopt a subset from 46 tasks from LLM-foundry [69], which includes standard  
 198 tasks with both zero-shot and few-shot evaluations. Specifically, we consider a 17-task subset where,  
 199 for each evaluation, at least one 0.154B scale model—trained with as many as 99B tokens—gets  
 200 10 percentage points above chance accuracy: ARC-Easy [23], BIG-bench: CS algorithms [11],  
 201 BIG-bench: Dyck languages [11], BIG-bench: Novel Concepts [11], BIG-bench: Operators [11],  
 202 BIG-bench: QA WikiData [11], BoolQ [21], Commonsense QA [107], COPA [92], CoQA [91],  
 203 HellaSwag (zero-shot) [126], HellaSwag (10-shot) [126], LAMBADA [77], PIQA [14], PubMed  
 204 QA Labeled [50], SQuAD [90], and WinoGrand [55]. For more details on evaluation datasets  
 205 see Appendix D. We focus on this subset to ensure we are measuring signal, not noise. Including  
 206 downstream tasks like MMLU [40], where performance is close to random chance, however, does  
 207 not invalidate our results as we show in our evaluation set ablations (Appendix E).

208 **Metrics.** We consider three main metrics: *Validation loss*, which is the cross entropy between a  
 209 model’s output and the one-hot ground truth token, averaged over all tokens in a sequence and over  
 210 all sequences in a dataset. *Average top-1 error*, which is a uniform average over the 17 downstream  
 211 evaluations, as mentioned in the above paragraph. To measure how good a prediction  $\zeta(C, M)$  is,  
 212 we measure *Relative prediction error*:  $|\zeta(C, M) - \zeta_{GT}|/\zeta_{GT}$ , where  $\zeta$  is the predicted loss  $L$  or the  
 213 average top-1 error  $\text{Err}$ .  $\zeta_{GT}$  is the ground truth measurement to predict.

## 214 4 Results: Reliable extrapolation

215 In this Section, we quantify the extent to which the scaling laws developed in Section 2 extrapolate  
 216 larger model performance using the scaling tested from Section 3. By default, we fit Equations (4)  
 217 and (5) to the configurations in Table 1, use C4 eval for loss, and the 17-task split from Section 3.4  
 218 for average top-1 error.

219 **Over-trained performance is predictable.** We highlight our main over-training results in  
 220 Figure 1 (left). Namely, we are able to extrapolate both in the number of parameters  $N$  and the  
 221 token multiplier  $M$  to closely predict the C4 eval performance of a 1.4B parameter model trained on  
 222 900B RedPajama tokens ( $N = 1.4\text{B}$ ,  $M = 640$ ). Our prediction, which takes  $300\times$  less compute  
 223 to construct than the final 1.4B run, is accurate to within 0.7% relative error. Additionally, for the  
 224  $N = 6.9\text{B}$ ,  $M = 20$  run, near compute-optimal, the relative error is also 0.7%.

225 These results support several key takeaways. (i) Scaling can be predictable even when one increases  
 226 both the model size and the amount of over-training compared to the training runs used to fit a scaling  
 227 law. (ii) The form presented in Equation (4) is useful in practice for predicting over-trained scaling  
 228 behavior. (iii) Fitting to Equation (4) gives good prediction accuracy near compute-optimal. More

Table 2: **Downstream relative prediction error at 6.9B parameters and 138B tokens.** While predicting accuracy on individual zero-shot downstream evaluations can be challenging (“Individual”), predicting *averages* across downstream datasets is accurate (“Avg.”).

Train set	Individual top-1 error				Avg. top-1 error
	ARC-E [23]	LAMBADA [77]	OpenBook QA [68]	HellaSwag [126]	17-task split
C4 [88, 27]	28.96%	15.01%	16.80%	79.58%	0.14%
RedPajama [112]	5.21%	14.39%	8.44%	25.73%	0.05%
RefinedWeb [82]	26.06%	16.55%	1.92%	81.96%	2.94%

229 specifically, predictions are accurate both for the 1.4B over-trained model and the 6.7B compute-  
 230 optimal model using a single scaling fit.

231 While Figure 1 explores a specific case of making predictions in the over-trained regime, we aim to  
 232 understand the error profile of our predictions across training datasets, token multipliers, and number  
 233 of parameters. Hence, Figure 5 shows the relative error between ground truth loss and predicted  
 234 loss on C4 eval for models in our testbed. We notice uniformly low prediction error suggesting that  
 235 predictions are accurate in many settings.

236 **Average top-1 error is predictable.** Figure 1 (*right*) presents our main result in estimating scaling  
 237 laws for downstream error. Concretely, we use the models indicated in Table 1 to fit Equations (4)  
 238 and (5), chaining the scaling fits to predict the average top-1 error as a function of training compute  
 239  $C$  and the token multiplier  $M$ . Our fits allow us to predict, using  $20\times$  less compute, the downstream  
 240 performance of a 6.9B model trained on 138B RedPajama tokens to within 0.05% relative error and a  
 241 1.4B model trained on RedPajama 900B tokens to within 3.6% relative error.

242 Table 2 additionally shows the relative error of our downstream performance predictions for models  
 243 trained on C4, RedPajama, and RefinedWeb, indicating that our scaling law functional forms are  
 244 applicable on many training datasets. We note that while average accuracy is predictable, *individual*  
 245 downstream task predictions are significantly more noisy. We report relative error for more model  
 246 predictions in Figures 11 and 12. We also find that if we remove the 1.4B model for the Equation (5)  
 247 fit, relative error jumps, for instance, from 0.05% to 10.64% on the 17-task split for the 6.9B,  
 248 138B token RedPajama prediction. This highlights the importance of investing more compute when  
 249 constructing scaling laws for downstream task prediction compared to loss prediction.

250 **Under-training, out-of-distribution scaling, and compute-reliability trade-offs.** In addition to  
 251 our main results presented above, we include additional results in Appendix E, which we summarize  
 252 here. First, we notice that when token multipliers become too small (i.e.,  $M = 5$ ) scaling becomes  
 253 unreliable and lies off the trend. Additionally, multipliers other than 20, such as 10, 40, and 80, garner  
 254 points that are roughly on the compute optimal frontier (Figure 9). This observation suggests that the  
 255 compute-optimal multiplier may lie in a range rather than take a single value. To probe the limits of  
 256 reliable scaling, we attempt to break our scaling laws in out-of-distribution settings. We find that  
 257 models trained on C4—English filtered—and evaluated on next token prediction on code domains  
 258 have a high relative error in many cases. Perhaps surprisingly, evaluating the same models on German  
 259 next token prediction gives reliable loss scaling (Figure 10). We additionally examine the compute  
 260 necessary to create accurate scaling laws, finding that scaling laws can be constructed more cheaply  
 261 for loss prediction than for downstream error prediction (Figures 15 and 16).

## 262 5 Related work

263 We review the most closely related work in this section. For additional related work, see Appendix F.

264 **Scaling laws.** Early works on scaling artificial neural networks observe predictable power-law  
 265 scaling in the training set size and number of model parameters [43, 44, 93]. Alabdulmohsin et al.  
 266 [2] stress the importance of looking at the extrapolation regime of a scaling law. Yang et al. [124]  
 267 prescribe architectural and hyperparameter changes when scaling model width to realize performant  
 268 models; Yang et al. [125] make analogous recommendations when scaling model depth. Bi et al.  
 269 [13] propose hyperparameter aware scaling laws. Unlike the aforementioned work, our investigation  
 270 focuses on over-training and predicting downstream accuracy.

271 Hoffmann et al. [45] investigate how the number of model parameters  $N$  and training tokens  $D$   
 272 should be chosen to minimize loss  $L$  given a compute budget  $C$ . Hoffmann et al. [45] find that when  
 273 scaling up  $C$ , both  $N$  and  $D$  should be scaled equally up to a multiplicative constant (i.e.,  $N \propto C^{\sim 0.5}$ )



274 and  $D \propto C^{\sim 0.5}$ ) to realize compute-optimality. Appendix C of the Chinchilla paper additionally  
275 suggests that these findings hold across three datasets. However, Hoffmann et al. [45] do not verify  
276 their scaling laws for training beyond compute-optimality, or for downstream error prediction—both  
277 of which are central to our work.

278 Sardana & Frankle [98] propose modifications to the Chinchilla formulation to incorporate inference  
279 costs into the definition of compute-optimality and solve for various fixed inference budgets. Their  
280 key finding, which is critical for our work, is that when taking into account a large enough inference  
281 budget, it is optimal to train smaller models for longer than the original Chinchilla recommendations.  
282 Our work presupposes that over-training can be beneficial. Instead of solving for inference-  
283 optimal schemes, we support empirically a predictive theory of scaling in the over-trained regime.  
284 Additionally, we provide experiments across many validation and training sets.

285 For predicting downstream scaling beyond loss, Isik et al. [47] relate the number of pre-training tokens  
286 to downstream cross-entropy and machine translation BLEU score [78] after fine-tuning. In contrast,  
287 we take a holistic approach to evaluation by looking at top-1 error over many natural language tasks.  
288 Schaeffer et al. [100] argue that emergent abilities [120] are a product of non-linear metrics and  
289 propose smoother alternatives. As a warmup for why non-linear metrics may be hard to predict,  
290 Schaeffer et al. [100] consider predicting an  $\ell$  length sequence exactly:  $\text{Err}(N, \ell) \approx 1 - \text{PP}(N)^{-\ell}$ ,  
291 where  $N$  is the number of parameters in a model and  $\text{PP}$  is its perplexity. This is a special case of  
292 our Equations (5) and (6), where the number of training tokens does not appear,  $\epsilon = 1, k = 1$ , and  
293  $\gamma = \ell$ . In contrast, we treat  $\epsilon, k, \gamma$  as free parameters for a scaling law fit, finding that average error  
294 over downstream tasks can make for a predictable metric.

295 **Over-training in popular models.** There has been a rise in over-trained models [113, 114] and  
296 accompanying massive datasets [112, 82, 104, 3]. For example, Chinchilla 70B [45] is trained with a  
297 token multiplier of 20, while LLaMA-2 7B [114] uses a token multiplier of 290. In our investigation,  
298 we look at token multipliers from 5 to 640 to ensure coverage of popular models and relevance for  
299 future models that may be trained on even more tokens.

## 300 6 Limitations, future work, and conclusion

301 **Limitations and future work.** We identify limitations, which provide motivation for future work.

- 302 • **Hyperparameters.** While our configurations are surprisingly amenable to reliable scaling across  
303 many training and testing distributions without further tuning, there is a need to develop scaling  
304 laws that do not require extensive hyperparameter sweeps.
- 305 • **Scaling up.** Validating the trends in this paper for even larger runs is a valuable direction.  
306 Additionally, repeating our setup for models that achieve non-trivial performance on harder  
307 evaluations like MMLU is left to future work.
- 308 • **Scaling down.** Actualizing predictable scaling with even cheaper runs is important to make this  
309 area of research more accessible, especially for downstream error prediction.
- 310 • **Failure cases.** While we present a preliminary analysis of when scaling is unreliable, future work  
311 should investigate conditions under which scaling breaks down.
- 312 • **Post-training.** It is common to employ fine-tuning interventions after pre-training, which we do  
313 not consider. Quantifying to what degree over-training the base model provides benefits *after*  
314 post-training is an open area of research.
- 315 • **Individual downstream task prediction.** While we find that averaging over many task error  
316 metrics can make for a predictable metric, per-task predictions are left to future work.
- 317 • **In-the-wild performance.** Downstream task performance is a proxy for the in-the-wild user  
318 experience. Analyzing scaling trends in the context of this experience is timely.
- 319 • **Dataset curation.** Our work only deals with existing training datasets. Exploring dataset curation  
320 for improved model scaling is another promising direction.

321 **Conclusion.** We show that the loss of over-trained models, trained past compute-optimality, is  
322 predictable. Furthermore, we propose and validate a scaling law relating loss to average downstream  
323 task performance. We hope our work will inspire others to further examine the relationship between  
324 model training and downstream generalization. Our testbed will be made publicly available, and we  
325 hope it will make scaling research more accessible to researchers and practitioners alike.

## References

- 326
- 327 [1] Samira Abnar, Mostafa Dehghani, Behnam Neyshabur, and Hanie Sedghi. Exploring the limits  
328 of large scale pre-training. In *International Conference on Learning Representations (ICLR)*,  
329 2022. <https://arxiv.org/abs/2110.02095>.
- 330 [2] Ibrahim Alabdulmohsin, Behnam Neyshabur, and Xiaohua Zhai. Revisiting neural scaling  
331 laws in language and vision. In *Advances in Neural Information Processing Systems (NeurIPS)*,  
332 2022. <https://arxiv.org/abs/2209.06640>.
- 333 [3] Alon Albalak, Yanai Elazar, Sang Michael Xie, Shayne Longpre, Nathan Lambert, Xinyi  
334 Wang, Niklas Muennighoff, Bairu Hou, Liangming Pan, Haewon Jeong, et al. A survey on  
335 data selection for language models. *arXiv preprint*, 2024. [https://arxiv.org/abs/2402.](https://arxiv.org/abs/2402.16827)  
336 [16827](https://arxiv.org/abs/2402.16827).
- 337 [4] Loubna Ben Allal, Raymond Li, Denis Kocetkov, Chenghao Mou, Christopher Akiki,  
338 Carlos Munoz Ferrandis, Niklas Muennighoff, Mayank Mishra, Alex Gu, Manan Dey, et al.  
339 Santacoder: don't reach for the stars! *arXiv preprint*, 2023. [https://arxiv.org/abs/](https://arxiv.org/abs/2301.03988)  
340 [2301.03988](https://arxiv.org/abs/2301.03988).
- 341 [5] Aida Amini, Saadia Gabriel, Shanchuan Lin, Rik Koncel-Kedziorski, Yejin Choi, and  
342 Hannaneh Hajishirzi. MathQA: Towards interpretable math word problem solving with  
343 operation-based formalisms. In *Conference of the North American Chapter of the Association*  
344 *for Computational Linguistics (NACCL)*, 2019. <https://aclanthology.org/N19-1245>.
- 345 [6] Jason Ansel, Edward Yang, Horace He, Natalia Gimelshein, Animesh Jain, Michael  
346 Voznesensky, Bin Bao, David Berard, Geeta Chauhan, Anjali Chourdia, Will Constable,  
347 Alban Desmaison, Zachary DeVito, Elias Ellison, Will Feng, Jiong Gong, Michael Gschwind,  
348 Brian Hirsh, Sherlock Huang, Laurent Kirsch, Michael Lazos, Yanbo Liang, Jason Liang,  
349 Yinghai Lu, CK Luk, Bert Maher, Yunjie Pan, Christian Puhersch, Matthias Reso, Mark  
350 Saroufim, Helen Suk, Michael Suo, Phil Tillet, Eikan Wang, Xiaodong Wang, William  
351 Wen, Shunting Zhang, Xu Zhao, Keren Zhou, Richard Zou, Ajit Mathews, Gregory Chanan,  
352 Peng Wu, and Soumith Chintala. Pytorch 2: Faster machine learning through dynamic  
353 python bytecode transformation and graph compilation. In *International Conference on*  
354 *Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, 2024.  
355 <https://pytorch.org/blog/pytorch-2-paper-tutorial>.
- 356 [7] Mikel Artetxe, Shruti Bhosale, Naman Goyal, Todor Mihaylov, Myle Ott, Sam Shleifer,  
357 Xi Victoria Lin, Jingfei Du, Srinivasan Iyer, Ramakanth Pasunuru, Giridharan Anantharaman,  
358 Xian Li, Shuohui Chen, Halil Akin, Mandeep Baines, Louis Martin, Xing Zhou, Punit Singh  
359 Koura, Brian O'Horo, Jeffrey Wang, Luke Zettlemoyer, Mona Diab, Zornitsa Kozareva, and  
360 Veselin Stoyanov. Efficient large scale language modeling with mixtures of experts. In  
361 *Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2022. [https:](https://aclanthology.org/2022.emnlp-main.804)  
362 [//aclanthology.org/2022.emnlp-main.804](https://aclanthology.org/2022.emnlp-main.804).
- 363 [8] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint*,  
364 2016. <https://arxiv.org/abs/1607.06450>.
- 365 [9] Yasaman Bahri, Ethan Dyer, Jared Kaplan, Jaehoon Lee, and Utkarsh Sharma. Explaining  
366 neural scaling laws. *arXiv preprint*, 2021. <https://arxiv.org/abs/2102.06701>.
- 367 [10] Yamini Bansal, Behrooz Ghorbani, Ankush Garg, Biao Zhang, Maxim Krikun, Colin Cherry,  
368 Behnam Neyshabur, and Orhan Firat. Data scaling laws in nmt: The effect of noise and  
369 architecture. In *International Conference on Machine Learning (ICML)*, 2022. [https:](https://proceedings.mlr.press/v162/bansal22b.html)  
370 [//proceedings.mlr.press/v162/bansal22b.html](https://proceedings.mlr.press/v162/bansal22b.html).
- 371 [11] BIG bench authors. Beyond the imitation game: Quantifying and extrapolating the capabilities  
372 of language models. In *Transactions on Machine Learning Research (TMLR)*, 2023. [https:](https://openreview.net/forum?id=uyTL5Bvosj)  
373 [//openreview.net/forum?id=uyTL5Bvosj](https://openreview.net/forum?id=uyTL5Bvosj).
- 374 [12] Emily M Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. On  
375 the dangers of stochastic parrots: Can language models be too big? In *Proceedings ACM*  
376 *conference on fairness, accountability, and transparency (FAccT)*, 2021. [https://dl.acm.](https://dl.acm.org/doi/10.1145/3442188.3445922)  
377 [org/doi/10.1145/3442188.3445922](https://dl.acm.org/doi/10.1145/3442188.3445922).

- 378 [13] DeepSeek-AI Xiao Bi, Deli Chen, Guanting Chen, Shanhuang Chen, Damai Dai, Chengqi  
379 Deng, Honghui Ding, Kai Dong, Qiushi Du, Zhe Fu, Huazuo Gao, Kaige Gao, Wenjun Gao,  
380 Ruiqi Ge, Kang Guan, Daya Guo, Jianzhong Guo, Guangbo Hao, Zhewen Hao, Ying He,  
381 Wen-Hui Hu, Panpan Huang, Erhang Li, Guowei Li, Jiashi Li, Yao Li, Y. K. Li, Wenfeng  
382 Liang, Fangyun Lin, A. X. Liu, Bo Liu, Wen Liu, Xiaodong Liu, Xin Liu, Yiyuan Liu, Haoyu  
383 Lu, Shanghao Lu, Fuli Luo, Shirong Ma, Xiaotao Nie, Tian Pei, Yishi Piao, Junjie Qiu, Hui  
384 Qu, Tongzheng Ren, Zehui Ren, Chong Ruan, Zhangli Sha, Zhihong Shao, Jun-Mei Song,  
385 Xuecheng Su, Jingxiang Sun, Yaofeng Sun, Min Tang, Bing-Li Wang, Peiyi Wang, Shiyu  
386 Wang, Yaohui Wang, Yongji Wang, Tong Wu, Yu Wu, Xin Xie, Zhenda Xie, Ziwei Xie,  
387 Yi Xiong, Hanwei Xu, Ronald X Xu, Yanhong Xu, Dejian Yang, Yu mei You, Shuiping Yu,  
388 Xin yuan Yu, Bo Zhang, Haowei Zhang, Lecong Zhang, Liyue Zhang, Mingchuan Zhang,  
389 Minghu Zhang, Wentao Zhang, Yichao Zhang, Chenggang Zhao, Yao Zhao, Shangyan Zhou,  
390 Shunfeng Zhou, Qihao Zhu, and Yuheng Zou. Deepseek llm: Scaling open-source language  
391 models with longtermism. *arXiv preprint*, 2024. <https://arxiv.org/abs/2401.02954>.
- 392 [14] Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng Gao, and Yejin Choi. Piqa: Reasoning  
393 about physical commonsense in natural language. In *Association for the Advancement of*  
394 *Artificial Intelligence (AAAI)*, 2020. <https://arxiv.org/abs/1911.11641>.
- 395 [15] Sid Black, Stella Biderman, Eric Hallahan, Quentin Anthony, Leo Gao, Laurence Golding,  
396 Horace He, Connor Leahy, Kyle McDonell, Jason Phang, Michael Pieler, USVSN Sai  
397 Prashanth, Shivanshu Purohit, Laria Reynolds, Jonathan Tow, Ben Wang, and Samuel  
398 Weinbach. Gpt-neox-20b: An open-source autoregressive language model. *BigScience*  
399 *Episode #5 – Workshop on Challenges & Perspectives in Creating Large Language Models*,  
400 2022. <https://aclanthology.org/2022.bigscience-1.9>.
- 401 [16] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla  
402 Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini  
403 Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya  
404 Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric  
405 Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam  
406 McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-  
407 shot learners. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.  
408 <https://arxiv.org/abs/2005.14165>.
- 409 [17] Ethan Caballero, Kshitij Gupta, Irina Rish, and David Krueger. Broken neural scaling laws. In  
410 *International Conference on Learning Representations (ICLR)*, 2023. <https://openreview.net/forum?id=sckjveqlCZ>.
- 412 [18] Mehdi Cherti, Romain Beaumont, Ross Wightman, Mitchell Wortsman, Gabriel Ilharco, Cade  
413 Gordon, Christoph Schuhmann, Ludwig Schmidt, and Jenia Jitsev. Reproducible scaling  
414 laws for contrastive language-image learning. In *Conference on Computer Vision and Pattern*  
415 *Recognition (CVPR)*, 2023. <https://arxiv.org/abs/2212.07143>.
- 416 [19] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam  
417 Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh,  
418 Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay,  
419 Noam M. Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Benton C. Hutchinson,  
420 Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju  
421 Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier García,  
422 Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan,  
423 Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani  
424 Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pillai, Marie  
425 Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee,  
426 Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Díaz, Orhan Firat, Michele Catasta, Jason  
427 Wei, Kathleen S. Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel.  
428 Palm: Scaling language modeling with pathways. In *Journal of Machine Learning Research*  
429 *(JMLR)*, 2022. <https://arxiv.org/abs/2204.02311>.
- 430 [20] Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan  
431 Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. Scaling instruction-finetuned  
432 language models. *arXiv preprint*, 2022. <https://arxiv.org/abs/2210.11416>.

- 433 [21] Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and  
434 Kristina Toutanova. Boolq: Exploring the surprising difficulty of natural yes/no questions. In  
435 *Conference of the North American Chapter of the Association for Computational Linguistics*  
436 (NAACL), 2019. <https://aclanthology.org/N19-1300>.
- 437 [22] Kevin Clark, Minh-Thang Luong, Quoc V. Le, and Christopher D. Manning. ELECTRA:  
438 Pre-training text encoders as discriminators rather than generators. In *International*  
439 *Conference on Learning Representations (ICLR)*, 2020. [https://openreview.net/pdf?](https://openreview.net/pdf?id=r1xMH1BtvB)  
440 [id=r1xMH1BtvB](https://openreview.net/pdf?id=r1xMH1BtvB).
- 441 [23] Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick,  
442 and Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning  
443 challenge. *arXiv preprint*, 2018. <https://arxiv.org/abs/1803.05457>.
- 444 [24] Tri Dao, Daniel Y. Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. FlashAttention: Fast  
445 and memory-efficient exact attention with IO-awareness. In *Advances in Neural Information*  
446 *Processing Systems (NeurIPS)*, 2022. <https://arxiv.org/abs/2205.14135>.
- 447 [25] Mostafa Dehghani, Josip Djolonga, Basil Mustafa, Piotr Padlewski, Jonathan Heek, Justin  
448 Gilmer, Andreas Peter Steiner, Mathilde Caron, Robert Geirhos, Ibrahim Alabdulmohsin, et al.  
449 Scaling vision transformers to 22 billion parameters. In *International Conference on Machine*  
450 *Learning (ICML)*, 2023. <https://proceedings.mlr.press/v202/dehghani23a.html>.
- 451 [26] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training  
452 of deep bidirectional transformers for language understanding. In *Conference of the North*  
453 *American Chapter of the Association for Computational Linguistics (NAACL)*, 2019. <https://aclanthology.org/N19-1423>.
- 454
- 455 [27] Jesse Dodge, Maarten Sap, Ana Marasović, William Agnew, Gabriel Ilharco, Dirk Groeneveld,  
456 Margaret Mitchell, and Matt Gardner. Documenting large webtext corpora: A case study on  
457 the colossal clean crawled corpus. In *Conference on Empirical Methods in Natural Language*  
458 *Processing (EMNLP)*, 2021. <https://aclanthology.org/2021.emnlp-main.98>.
- 459 [28] Nan Du, Yanping Huang, Andrew M. Dai, Simon Tong, Dmitry Lepikhin, Yuanzhong Xu,  
460 Maxim Krikun, Yanqi Zhou, Adams Wei Yu, Orhan Firat, Barret Zoph, Liam Fedus, Maarten  
461 Bosma, Zongwei Zhou, Tao Wang, Yu Emma Wang, Kellie Webster, Marie Pellat, Kevin  
462 Robinson, Kathleen Meier-Hellstern, Toju Duke, Lucas Dixon, Kun Zhang, Quoc V Le,  
463 Yonghui Wu, Zhifeng Chen, and Claire Cui. Glam: Efficient scaling of language models  
464 with mixture-of-experts. In *International Conference on Machine Learning (ICML)*, 2022.  
465 <https://arxiv.org/abs/2112.06905>.
- 466 [29] Kawin Ethayarajh, Winnie Xu, Niklas Muennighoff, Dan Jurafsky, and Douwe Kiela. Kto:  
467 Model alignment as prospect theoretic optimization. *arXiv preprint*, 2024. <https://arxiv.org/abs/2402.01306>.
- 468
- 469 [30] Samir Yitzhak Gadre, Gabriel Ilharco, Alex Fang, Jonathan Hayase, Georgios Smyrnis, Thao  
470 Nguyen, Mitchell Wortsman Ryan Marten, Dhruva Ghosh, Jieyu Zhang, Eyal Orgad, Rahim  
471 Entezari, Giannis Daras, Sarah Pratt, Vivek Ramanujan, Yonatan Bitton, Kalyani Marathe,  
472 Stephen Mussmann, Mehdi Cherti Richard Vencu, Ranjay Krishna, Pang Wei Koh, Olga  
473 Saukh, Alexander Ratner, Shuran Song, Hannaneh Hajishirzi, Ali Farhadi, Romain Beaumont,  
474 Sewoong Oh, Alex Dimakis, Jenia Jitsev, Yair Carmon, Vaishaal Shankar, and Ludwig Schmidt.  
475 Datacomp: In search of the next generation of multimodal datasets. In *Advances in Neural*  
476 *Information Processing Systems (NeurIPS)*, 2023. <https://arxiv.org/abs/2304.14108>.
- 477 [31] Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster,  
478 Jason Phang, Horace He, Anish Thite, Noa Nabeshima, Shawn Presser, and Connor Leahy.  
479 The Pile: An 800gb dataset of diverse text for language modeling. *arXiv preprint*, 2020.  
480 <https://arxiv.org/abs/2101.00027>.
- 481 [32] Behrooz Ghorbani, Orhan Firat, Markus Freitag, Ankur Bapna, Maxim Krikun, Xavier Garcia,  
482 Ciprian Chelba, and Colin Cherry. Scaling laws for neural machine translation. *arXiv preprint*,  
483 2021. <https://arxiv.org/abs/2109.07740>.

- 484 [33] Mitchell A Gordon, Kevin Duh, and Jared Kaplan. Data and parameter scaling laws for neural  
485 machine translation. In *Conference on Empirical Methods in Natural Language Processing*  
486 (*EMNLP*), 2021. <https://aclanthology.org/2021.emnlp-main.478>.
- 487 [34] Dirk Groeneveld, Iz Beltagy, Pete Walsh, Akshita Bhagia, Rodney Kinney, Oyvind Tafjord,  
488 Ananya Harsh Jha, Hamish Ivison, Ian Magnusson, Yizhong Wang, et al. Olmo: Accelerating  
489 the science of language models. *arXiv preprint*, 2024. <https://arxiv.org/abs/2402.00838>.  
490
- 491 [35] Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces.  
492 *arXiv preprint*, 2023. <https://arxiv.org/abs/2312.00752>.
- 493 [36] Albert Gu, Isys Johnson, Karan Goel, Khaled Saab, Tri Dao, Atri Rudra, and Christopher  
494 Ré. Combining recurrent, convolutional, and continuous-time models with linear state space  
495 layers. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021. <https://openreview.net/forum?id=yWd42CWN3c>.  
496
- 497 [37] Albert Gu, Karan Goel, and Christopher Ré. Efficiently modeling long sequences with  
498 structured state spaces. In *International Conference on Learning Representations (ICLR)*,  
499 2022. <https://arxiv.org/abs/2111.00396>.
- 500 [38] Suriya Gunasekar, Yi Zhang, Jyoti Aneja, Caio Cesar, Teodoro Mendes, Allie Del Giorno,  
501 Sivakanth Gopi, Mojan Javaheripi, Piero Kauffmann, Gustavo de Rosa, Olli Saarikivi,  
502 Adil Salim, Shital Shah, Harkirat Singh Behl, Xin Wang, Sébastien Bubeck, Ronen  
503 Eldan, Adam Tauman Kalai, Yin Tat Lee, and Yuanzhi Li. Textbooks are all you  
504 need. *Preprint*, 2023. [https://www.microsoft.com/en-us/research/publication/  
505 textbooks-are-all-you-need](https://www.microsoft.com/en-us/research/publication/textbooks-are-all-you-need).
- 506 [39] Suchin Gururangan, Mitchell Wortsman, Samir Yitzhak Gadre, Achal Dave, Maciej Kilian,  
507 Weijia Shi, Jean Mercat, Georgios Smyrnis, Gabriel Ilharco, Matt Jordan, Reinhard  
508 Heckel, Alex Dimakis, Ali Farhadi, Vaishaal Shankar, and Ludwig Schmidt. OpenLM:  
509 a minimal but performative language modeling (lm) repository, 2023. [https://github.  
510 com/mlfoundations/open\\_lm](https://github.com/mlfoundations/open_lm).
- 511 [40] Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and  
512 Jacob Steinhardt. Measuring massive multitask language understanding. In *International  
513 Conference on Learning Representations (ICLR)*, 2021. [https://arxiv.org/abs/2009.  
514 03300](https://arxiv.org/abs/2009.03300).
- 515 [41] T. J. Henighan, Jared Kaplan, Mor Katz, Mark Chen, Christopher Hesse, Jacob Jackson,  
516 Heewoo Jun, Tom B. Brown, Prafulla Dhariwal, Scott Gray, Chris Hallacy, Benjamin Mann,  
517 Alec Radford, Aditya Ramesh, Nick Ryder, Daniel M. Ziegler, John Schulman, Dario Amodei,  
518 and Sam McCandlish. Scaling laws for autoregressive generative modeling. *arXiv preprint*,  
519 2020. <https://arxiv.org/abs/2010.14701>.
- 520 [42] Danny Hernandez, Jared Kaplan, T. J. Henighan, and Sam McCandlish. Scaling laws for  
521 transfer. *arXiv preprint*, 2021. <https://arxiv.org/abs/2102.01293>.
- 522 [43] Joel Hestness, Sharan Narang, Newsha Ardalani, Gregory Frederick Diamos, Heewoo Jun,  
523 Hassan Kianinejad, Md. Mostofa Ali Patwary, Yang Yang, and Yanqi Zhou. Deep learning  
524 scaling is predictable, empirically. *arXiv preprint*, 2017. [https://arxiv.org/abs/1712.  
525 00409](https://arxiv.org/abs/1712.00409).
- 526 [44] Joel Hestness, Newsha Ardalani, and Gregory Diamos. Beyond human-level accuracy:  
527 Computational challenges in deep learning. In *Principles and Practice of Parallel  
528 Programming (PPoPP)*, 2019. <https://arxiv.org/abs/1909.01736>.
- 529 [45] Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai,  
530 Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark,  
531 et al. Training compute-optimal large language models. In *Advances in Neural Information  
532 Processing Systems (NeurIPS)*, 2022. <https://arxiv.org/abs/2203.15556>.

- 533 [46] Hakan Inan, Khashayar Khosravi, and Richard Socher. Tying word vectors and word  
534 classifiers: A loss framework for language modeling. In *International Conference on Learning*  
535 *Representations (ICLR)*, 2017. <https://arxiv.org/abs/1611.01462>.
- 536 [47] Berivan Isik, Natalia Ponomareva, Hussein Hazimeh, Dimitris Paparas, Sergei Vassilvitskii,  
537 and Sanmi Koyejo. Scaling laws for downstream task performance of large language models.  
538 *arXiv*, 2024. <https://arxiv.org/abs/2402.04177>.
- 539 [48] Maor Ivgi, Yair Carmon, and Jonathan Berant. Scaling laws under the microscope: Predicting  
540 transformer performance from small scale experiments. In *Conference on Empirical Methods*  
541 *in Natural Language Processing (EMNLP)*, 2022. [https://aclanthology.org/2022.](https://aclanthology.org/2022.findings-emnlp.544)  
542 [findings-emnlp.544](https://aclanthology.org/2022.findings-emnlp.544).
- 543 [49] Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh  
544 Chaplot, Florian Bressand Diego de las Casas, Gianna Lengyel, Guillaume Lample, Lucile  
545 Saulnier, L elio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut  
546 Lavril, Thomas Wang, Timoth ee Lacroix, and William El Sayed. Mistral 7b. *arXiv preprint*,  
547 2023. <https://arxiv.org/abs/2310.06825>.
- 548 [50] Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William Cohen, and Xinghua Lu. Pubmedqa: A  
549 dataset for biomedical research question answering. In *Conference on Empirical Methods in*  
550 *Natural Language Processing (EMNLP)*, 2019. <https://aclanthology.org/D19-1259>.
- 551 [51] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon  
552 Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural  
553 language models. *arXiv preprint*, 2020. <https://arxiv.org/abs/2001.08361>.
- 554 [52] Tobit Klug, Dogukan Atik, and Reinhard Heckel. Analyzing the sample complexity of self-  
555 supervised image reconstruction methods. *arXiv preprint*, 2023. [https://arxiv.org/abs/](https://arxiv.org/abs/2305.19079)  
556 [2305.19079](https://arxiv.org/abs/2305.19079).
- 557 [53] Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu  
558 Soricut. ALBERT: A lite BERT for self-supervised learning of language representations. *arXiv*  
559 *preprint*, 2019. <http://arxiv.org/abs/1909.11942>.
- 560 [54] Benjamin Lefaudeux, Francisco Massa, Diana Liskovich, Wenhan Xiong, Vittorio Caggiano,  
561 Sean Naren, Min Xu, Jieru Hu, Marta Tintore, Susan Zhang, Patrick Labatut, and Daniel  
562 Haziza. xformers: A modular and hackable transformer modelling library, 2022. [https:](https://github.com/facebookresearch/xformers)  
563 [//github.com/facebookresearch/xformers](https://github.com/facebookresearch/xformers).
- 564 [55] Hector Levesque, Ernest Davis, and Leora Morgenstern. The winograd schema challenge. In  
565 *International conference on the principles of knowledge representation and reasoning*, 2012.  
566 <https://aaai.org/papers/59-4492-the-winograd-schema-challenge>.
- 567 [56] Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed,  
568 Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. BART: Denoising sequence-to-  
569 sequence pre-training for natural language generation, translation, and comprehension. In  
570 *Annual Meeting of the Association for Computational Linguistics (ACL)*, 2020. [https:](https://aclanthology.org/2020.acl-main.703)  
571 [//aclanthology.org/2020.acl-main.703](https://aclanthology.org/2020.acl-main.703).
- 572 [57] Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao  
573 Mou, Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, et al. Starcoder: may the source  
574 be with you! *arXiv preprint*, 2023. <https://arxiv.org/abs/2305.06161>.
- 575 [58] Jian Liu, Leyang Cui, Hanmeng Liu, Dandan Huang, Yile Wang, and Yue Zhang. Logiqa: A  
576 challenge dataset for machine reading comprehension with logical reasoning. In *International*  
577 *Joint Conference on Artificial Intelligence*, 2020. <https://arxiv.org/abs/2007.08124>.
- 578 [59] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy,  
579 Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized BERT  
580 pretraining approach. *arXiv preprint*, 2019. <http://arxiv.org/abs/1907.11692>.

- 581 [60] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining  
582 Xie. A convnet for the 2020s. *Conference on Computer Vision and Pattern Recognition*  
583 *(CVPR)*, 2022. <https://arxiv.org/abs/2201.03545>.
- 584 [61] Shayne Longpre, Robert Mahari, Anthony Chen, Naana Obeng-Marnu, Damien Sileo, William  
585 Brannon, Niklas Muennighoff, Nathan Khazam, Jad Kabbara, Kartik Perisetla, et al. The  
586 data provenance initiative: A large scale audit of dataset licensing & attribution in ai. *arXiv*  
587 *preprint*, 2023. <https://arxiv.org/abs/2310.16787>.
- 588 [62] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint*,  
589 2017. <https://arxiv.org/abs/1711.05101>.
- 590 [63] Anton Lozhkov, Raymond Li, Loubna Ben Allal, Federico Cassano, Joel Lamy-Poirier,  
591 Nouamane Tazi, Ao Tang, Dmytro Pykhtar, Jiawei Liu, Yuxiang Wei, Tianyang Liu, Max Tian,  
592 Denis Kocetkov, Arthur Zucker, Younes Belkada, Zijian Wang, Qian Liu, Dmitry Abulkhanov,  
593 Indraneil Paul, Zhuang Li, Wen-Ding Li, Megan Risdal, Jia Li, Jian Zhu, Terry Yue Zhuo,  
594 Evgenii Zheltonozhskii, Nii Osae Osae Dade, Wenhao Yu, Lucas Krauß, Naman Jain, Yixuan  
595 Su, Xuanli He, Manan Dey, Edoardo Abati, Yekun Chai, Niklas Muennighoff, Xiangru Tang,  
596 Muhtasham Oblokulov, Christopher Akiki, Marc Marone, Chenghao Mou, Mayank Mishra,  
597 Alex Gu, Binyuan Hui, Tri Dao, Armel Zebaze, Olivier Dehaene, Nicolas Patry, Canwen Xu,  
598 Julian McAuley, Han Hu, Torsten Scholach, Sebastien Paquet, Jennifer Robinson, Carolyn Jane  
599 Anderson, Nicolas Chapados, Mostofa Patwary, Nima Tajbakhsh, Yacine Jernite, Carlos Muñoz  
600 Ferrandis, Lingming Zhang, Sean Hughes, Thomas Wolf, Arjun Guha, Leandro von Werra,  
601 and Harm de Vries. Starcoder 2 and the stack v2: The next generation. *arXiv preprint*, 2024.  
602 <https://arxiv.org/abs/2402.19173>.
- 603 [64] Risto Luukkonen, Ville Komulainen, Jouni Luoma, Anni Eskelinen, Jenna Kanerva, Hanna-  
604 Mari Kupari, Filip Ginter, Veronika Laippala, Niklas Muennighoff, Aleksandra Piktus, et al.  
605 Fingpt: Large generative models for a small language. In *Conference on Empirical Methods*  
606 *in Natural Language Processing (EMNLP)*, 2023. [https://aclanthology.org/2023.](https://aclanthology.org/2023.emnlp-main.164)  
607 [emnlp-main.164](https://aclanthology.org/2023.emnlp-main.164).
- 608 [65] Ian Magnusson, Akshita Bhagia, Valentin Hofmann, Luca Soldaini, Ananya Harsh Jha, Oyvind  
609 Tafjord, Dustin Schwenk, Evan Pete Walsh, Yanai Elazar, Kyle Lo, Dirk Groenvelde, Iz Beltagy,  
610 Hannenh Hajishirz, Noah A. Smith, Kyle Richardson, and Jesse Dodge. Paloma: A benchmark  
611 for evaluating language model fit. *arXiv preprint*, 2023. <https://paloma.allen.ai>.
- 612 [66] Mitchell P. Marcus, Beatrice Santorini, and Mary Ann Marcinkiewicz. Building a large  
613 annotated corpus of English: The Penn Treebank. In *Computational Linguistics*, 1993.  
614 <https://aclanthology.org/J93-2004>.
- 615 [67] William Merrill, Vivek Ramanujan, Yoav Goldberg, Roy Schwartz, and Noah A. Smith.  
616 Effects of parameter norm growth during transformer training: Inductive bias from gradient  
617 descent. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*,  
618 2021. <https://aclanthology.org/2021.emnlp-main.133>.
- 619 [68] Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor conduct  
620 electricity? a new dataset for open book question answering. In *Conference on Empirical*  
621 *Methods in Natural Language Processing (EMNLP)*, 2018. [https://arxiv.org/abs/1809.](https://arxiv.org/abs/1809.02789)  
622 [02789](https://arxiv.org/abs/1809.02789).
- 623 [69] MosaicML. Llm evaluation scores, 2023. <https://www.mosaicml.com/llm-evaluation>.
- 624 [70] Niklas Muennighoff, Thomas Wang, Lintang Sutawika, Adam Roberts, Stella Biderman,  
625 Teven Le Scao, M Saiful Bari, Sheng Shen, Zheng-Xin Yong, Hailey Schoelkopf, et al.  
626 Crosslingual generalization through multitask finetuning. In *Annual Meeting of the*  
627 *Association for Computational Linguistics (ACL)*, 2022. [https://aclanthology.org/](https://aclanthology.org/2023.acl-long.891)  
628 [2023.acl-long.891](https://aclanthology.org/2023.acl-long.891).
- 629 [71] Niklas Muennighoff, Qian Liu, Armel Zebaze, Qinkai Zheng, Binyuan Hui, Terry Yue  
630 Zhuo, Swayam Singh, Xiangru Tang, Leandro Von Werra, and Shayne Longpre. Octopack:  
631 Instruction tuning code large language models. *arXiv preprint*, 2023. [https://arxiv.org/](https://arxiv.org/abs/2308.07124)  
632 [abs/2308.07124](https://arxiv.org/abs/2308.07124).

- 633 [72] Niklas Muennighoff, Alexander M Rush, Boaz Barak, Teven Le Scao, Aleksandra Piktus,  
634 Nouamane Tazi, Sampo Pyysalo, Thomas Wolf, and Colin Raffel. Scaling data-constrained  
635 language models. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2023.  
636 <https://arxiv.org/abs/2305.16264>.
- 637 [73] Niklas Muennighoff, Hongjin Su, Liang Wang, Nan Yang, Furu Wei, Tao Yu, Amanpreet  
638 Singh, and Douwe Kiela. Generative representational instruction tuning. *arXiv preprint*, 2024.  
639 <https://arxiv.org/abs/2402.09906>.
- 640 [74] Erik Nijkamp, Tian Xie, Hiroaki Hayashi, Bo Pang, Congying Xia, Chen Xing, Jesse Vig,  
641 Semih Yavuz, Philippe Laban, Ben Krause, Senthil Purushwalkam, Tong Niu, Wojciech  
642 Kryscinski, Lidiya Murakhovska, Prafulla Kumar Choubey, Alex Fabbri, Ye Liu, Rui Meng,  
643 Lifu Tu, Meghana Bhat, Chien-Sheng Wu, Silvio Savarese, Yingbo Zhou, Shafiq Rayhan  
644 Joty, and Caiming Xiong. Long sequence modeling with xgen: A 7b llm trained on 8k input  
645 sequence length. *arXiv preprint*, 2023. <https://arxiv.org/abs/2309.03450>.
- 646 [75] OpenAI. Triton, 2021. <https://github.com/openai/triton>.
- 647 [76] OpenAI. Gpt-4 technical report, 2023. <https://arxiv.org/abs/2303.08774>.
- 648 [77] Denis Paperno, Germán Kruszewski, Angeliki Lazaridou, Ngoc Quan Pham, Raffaella  
649 Bernardi, Sandro Pezzelle, Marco Baroni, Gemma Boleda, and Raquel Fernandez. The  
650 LAMBADA dataset: Word prediction requiring a broad discourse context. In *Annual Meeting*  
651 *of the Association for Computational Linguistics (ACL)*, 2016. [http://www.aclweb.org/](http://www.aclweb.org/anthology/P16-1144)  
652 [anthology/P16-1144](http://www.aclweb.org/anthology/P16-1144).
- 653 [78] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic  
654 evaluation of machine translation. In *Annual Meeting of the Association for Computational*  
655 *Linguistics (ACL)*, 2002. <https://aclanthology.org/P02-1040>.
- 656 [79] Alicia Parrish, Angelica Chen, Nikita Nangia, Vishakh Padmakumar, Jason Phang, Jana  
657 Thompson, Phu Mon Htut, and Samuel Bowman. BBQ: A hand-built bias benchmark for  
658 question answering. In *Annual Meeting of the Association for Computational Linguistics*  
659 *(ACL)*, 2022. <https://aclanthology.org/2022.findings-acl.165>.
- 660 [80] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan,  
661 Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative  
662 style, high-performance deep learning library. In *Advances in Neural Information Processing*  
663 *Systems (NeurIPS)*, 2019. <https://arxiv.org/abs/1912.01703>.
- 664 [81] Patronus AI. EnterprisePII dataset, 2023. <https://tinyurl.com/2r5x9bst>.
- 665 [82] Guilherme Penedo, Quentin Malartic, Daniel Hesslow, Ruxandra Cojocaru, Alessandro  
666 Cappelli, Hamza Alobeidli, Baptiste Pannier, Ebtesam Almazrouei, and Julien Launay. The  
667 RefinedWeb dataset for Falcon LLM: outperforming curated corpora with web data, and web  
668 data only. *arXiv preprint*, 2023. <https://arxiv.org/abs/2306.01116>.
- 669 [83] Bo Peng, Eric Alcaide, Quentin Anthony, Alon Albalak, Samuel Arcadinho, Stella Biderman,  
670 Huanqi Cao, Xin Cheng, Michael Chung, Leon Derczynski, Xingjian Du, Matteo Grella,  
671 Kranthi Gv, Xuzheng He, Haowen Hou, Przemyslaw Kazienko, Jan Kocon, Jiaming Kong,  
672 Bartłomiej Koptyra, Hayden Lau, Jiaju Lin, Krishna Sri Ipsit Mantri, Ferdinand Mom, Atsushi  
673 Saito, Guangyu Song, Xiangru Tang, Johan Wind, Stanisław Woźniak, Zhenyuan Zhang,  
674 Qinghua Zhou, Jian Zhu, and Rui-Jie Zhu. RWKV: Reinventing RNNs for the transformer  
675 era. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2023.  
676 <https://aclanthology.org/2023.findings-emnlp.936>.
- 677 [84] Ofir Press and Lior Wolf. Using the output embedding to improve language models. In  
678 *Proceedings of the Conference of the European Chapter of the Association for Computational*  
679 *Linguistics (EACL)*, 2017. <https://aclanthology.org/E17-2025>.
- 680 [85] Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya  
681 Sutskever. Language models are unsupervised multitask learners. *Preprint*, 2019.  
682 [https://d4mucfpksywv.cloudfront.net/better-language-models/language\\_](https://d4mucfpksywv.cloudfront.net/better-language-models/language_models_are_unsupervised_multitask_learners.pdf)  
683 [models\\_are\\_unsupervised\\_multitask\\_learners.pdf](https://d4mucfpksywv.cloudfront.net/better-language-models/language_models_are_unsupervised_multitask_learners.pdf).



- 684 [86] Jack W. Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis  
685 Song, John Aslanides, Sarah Henderson, Roman Ring, Susannah Young, Eliza Rutherford,  
686 Tom Hennigan, Jacob Menick, Albin Cassirer, Richard Powell, George van den Driessche,  
687 Lisa Anne Hendricks, Maribeth Rauh, Po-Sen Huang, Amelia Glaese, Johannes Welbl,  
688 Sumanth Dathathri, Saffron Huang, Jonathan Uesato, John F. J. Mellor, Irina Higgins,  
689 Antonia Creswell, Nathan McAleese, Amy Wu, Erich Elsen, Siddhant M. Jayakumar,  
690 Elena Buchatskaya, David Budden, Esme Sutherland, Karen Simonyan, Michela Paganini,  
691 L. Sifre, Lena Martens, Xiang Lorraine Li, Adhiguna Kuncoro, Aida Nematzadeh, Elena  
692 Gribovskaya, Domenic Donato, Angeliki Lazaridou, Arthur Mensch, Jean-Baptiste Lespiau,  
693 Maria Tsimpoukelli, N. K. Grigorev, Doug Fritz, Thibault Sottiaux, Mantas Pajarskas, Tobias  
694 Pohlen, Zhitao Gong, Daniel Toyama, Cyprien de Masson d'Autume, Yujia Li, Tayfun Terzi,  
695 Vladimir Mikulik, Igor Babuschkin, Aidan Clark, Diego de Las Casas, Aurelia Guy, Chris  
696 Jones, James Bradbury, Matthew G. Johnson, Blake A. Hechtman, Laura Weidinger, Iason  
697 Gabriel, William S. Isaac, Edward Lockhart, Simon Osindero, Laura Rimell, Chris Dyer, Oriol  
698 Vinyals, Kareem W. Ayoub, Jeff Stanway, L. L. Bennett, Demis Hassabis, Koray Kavukcuoglu,  
699 and Geoffrey Irving. Scaling language models: Methods, analysis & insights from training  
700 gopher. *arXiv preprint*, 2021. <https://arxiv.org/abs/2112.11446>.
- 701 [87] Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and  
702 Chelsea Finn. Direct preference optimization: Your language model is secretly a reward  
703 model. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2023. <https://arxiv.org/abs/2305.18290>.
- 704 [88] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena,  
705 Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified  
706 text-to-text transformer. *arXiv preprint*, 2019. <https://arxiv.org/abs/1910.10683>.
- 707 [89] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena,  
708 Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified  
709 text-to-text transformer. In *The Journal of Machine Learning Research (JMLR)*, 2020. <https://arxiv.org/abs/1910.10683>.
- 710 [90] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. SQuAD: 100,000+  
711 questions for machine comprehension of text. In *Conference on Empirical Methods in Natural  
712 Language Processing (EMNLP)*, 2016. <https://aclanthology.org/D16-1264>.
- 713 [91] Siva Reddy, Danqi Chen, and Christopher D. Manning. CoQA: A conversational question  
714 answering challenge. In *Transactions of the Association for Computational Linguistics (TACL)*,  
715 2019. <https://aclanthology.org/Q19-1016>.
- 716 [92] Melissa Roemmele, Cosmin Adrian Bejan, , and Andrew S. Gordon. Choice of plausible  
717 alternatives: An evaluation of commonsense causal reasoning. In *Association for the  
718 Advancement of Artificial Intelligence (AAAI) Spring Symposium*, 2011. [https://people.  
719 ict.usc.edu/~gordon/copa.html](https://people.ict.usc.edu/~gordon/copa.html).
- 720 [93] Jonathan S. Rosenfeld, Amir Rosenfeld, Yonatan Belinkov, and Nir Shavit. A constructive  
721 prediction of the generalization error across scales. In *International Conference on Learning  
722 Representations (ICLR)*, 2020. <https://arxiv.org/abs/1909.12673>.
- 723 [94] Rachel Rudinger, Jason Naradowsky, Brian Leonard, and Benjamin Van Durme. Gender bias  
724 in coreference resolution. In *Conference of the North American Chapter of the Association for  
725 Computational Linguistics (NAACL)*, 2018. <https://aclanthology.org/N18-2002>.
- 726 [95] Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An  
727 adversarial winograd schema challenge at scale. *arXiv preprint*, 2019. [https://arxiv.org/  
728 abs/1907.10641](https://arxiv.org/abs/1907.10641).
- 729 [96] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. Distilbert, a distilled  
730 version of bert: smaller, faster, cheaper and lighter. *arXiv preprint*, 2019. [http://arxiv.  
731 org/abs/1910.01108](http://arxiv.org/abs/1910.01108).
- 732  
733

- 734 [97] Maarten Sap, Hannah Rashkin, Derek Chen, Ronan Le Bras, and Yejin Choi. Social IQa:  
735 Commonsense reasoning about social interactions. In *Empirical Methods in Natural Language*  
736 *Processing (EMNLP)*, 2019. <https://aclanthology.org/D19-1454>.
- 737 [98] Nikhil Sardana and Jonathan Frankle. Beyond chinchilla-optimal: Accounting for inference  
738 in language model scaling laws. In *NeurIPS Workshop on Efficient Natural Language and*  
739 *Speech Processing (ENLSP)*, 2023. <https://arxiv.org/abs/2401.00448>.
- 740 [99] Teven Le Scao, Thomas Wang, Daniel Hesslow, Lucile Saulnier, Stas Bekman, M Saiful Bari,  
741 Stella Biderman, Hady Elsahar, Niklas Muennighoff, Jason Phang, et al. What language  
742 model to train if you have one million gpu hours? In *Conference on Empirical Methods*  
743 *in Natural Language Processing (EMNLP)*, 2022. [https://aclanthology.org/2022.](https://aclanthology.org/2022.findings-emnlp.54)  
744 [findings-emnlp.54](https://aclanthology.org/2022.findings-emnlp.54).
- 745 [100] Rylan Schaeffer, Brando Miranda, and Sanmi Koyejo. Are emergent abilities of large language  
746 models a mirage? In *Advances in Neural Information Processing Systems (NeurIPS)*, 2023.  
747 <https://arxiv.org/abs/2304.15004>.
- 748 [101] Utkarsh Sharma and Jared Kaplan. A neural scaling law from the dimension of the data  
749 manifold. In *Journal of Machine Learning Research (JMLR)*, 2022. [https://arxiv.org/](https://arxiv.org/abs/2004.10802)  
750 [abs/2004.10802](https://arxiv.org/abs/2004.10802).
- 751 [102] Noam Shazeer. Glu variants improve transformer. *arXiv preprint*, 2020. [https://arxiv.](https://arxiv.org/abs/2002.05202)  
752 [org/abs/2002.05202](https://arxiv.org/abs/2002.05202).
- 753 [103] Shivalika Singh, Freddie Vargus, Daniel Dsouza, Börje F Karlsson, Abinaya Mahendiran,  
754 Wei-Yin Ko, Herumb Shandilya, Jay Patel, Deividas Mataciunas, Laura OMahony, et al.  
755 Aya dataset: An open-access collection for multilingual instruction tuning. *arXiv preprint*  
756 *arXiv:2402.06619*, 2024. <https://arxiv.org/abs/2402.06619>.
- 757 [104] Luca Soldaini, Rodney Kinney, Akshita Bhagia, Dustin Schwenk, David Atkinson, Russell  
758 Authur, Ben Bogin, Khyathi Chandu, Jennifer Dumas, Yanai Elazar, et al. Dolma: An open  
759 corpus of three trillion tokens for language model pretraining research. *arXiv preprint*, 2024.  
760 <https://arxiv.org/abs/2402.00159>.
- 761 [105] Ben Sorscher, Robert Geirhos, Shashank Shekhar, Surya Ganguli, and Ari S. Morcos. Beyond  
762 neural scaling laws: beating power law scaling via data pruning. In *Advances in Neural*  
763 *Information Processing Systems (NeurIPS)*, 2022. [https://openreview.net/forum?id=](https://openreview.net/forum?id=UmvSlP-PyV)  
764 [UmvSlP-PyV](https://openreview.net/forum?id=UmvSlP-PyV).
- 765 [106] Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer:  
766 Enhanced transformer with rotary position embedding. *arXiv preprint*, 2021. [https://](https://arxiv.org/abs/2104.09864)  
767 [arxiv.org/abs/2104.09864](https://arxiv.org/abs/2104.09864).
- 768 [107] Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. CommonsenseQA:  
769 A question answering challenge targeting commonsense knowledge. In *Conference of the*  
770 *North American Chapter of the Association for Computational Linguistics (NAACL)*, 2019.  
771 <https://aclanthology.org/N19-1421>.
- 772 [108] Yi Tay, Mostafa Dehghani, Jinfeng Rao, William Fedus, Samira Abnar, Hyung Won  
773 Chung, Sharan Narang, Dani Yogatama, Ashish Vaswani, and Donald Metzler. Scale  
774 efficiently: Insights from pre-training and fine-tuning transformers. In *International*  
775 *Conference on Learning Representations (ICLR)*, 2022. [https://openreview.net/forum?](https://openreview.net/forum?id=f20YVDyfIB)  
776 [id=f20YVDyfIB](https://openreview.net/forum?id=f20YVDyfIB).
- 777 [109] Yi Tay, Mostafa Dehghani, Samira Abnar, Hyung Chung, William Fedus, Jinfeng Rao,  
778 Sharan Narang, Vinh Tran, Dani Yogatama, and Donald Metzler. Scaling laws vs model  
779 architectures: How does inductive bias influence scaling? In *Conference on Empirical*  
780 *Methods in Natural Language Processing (EMNLP)*, 2023. [https://aclanthology.org/](https://aclanthology.org/2023.findings-emnlp.825)  
781 [2023.findings-emnlp.825](https://aclanthology.org/2023.findings-emnlp.825).
- 782 [110] MosaicML NLP Team. Introducing mpt-7b: A new standard for open-source, commercially  
783 usable llms, 2023. [www.mosaicml.com/blog/mpt-7b](http://www.mosaicml.com/blog/mpt-7b).

- 784 [111] Romal Thoppilan, Daniel De Freitas, Jamie Hall, Noam Shazeer, Apoorv Kulshreshtha,  
785 Heng-Tze Cheng, Alicia Jin, Taylor Bos, Leslie Baker, Yu Du, YaGuang Li, Hongrae Lee,  
786 Huaixiu Steven Zheng, Amin Ghafouri, Marcelo Menegali, Yanping Huang, Maxim Krikun,  
787 Dmitry Lepikhin, James Qin, Dehao Chen, Yuanzhong Xu, Zhifeng Chen, Adam Roberts,  
788 Maarten Bosma, Vincent Zhao, Yanqi Zhou, Chung-Ching Chang, Igor Krivokon, Will Rusch,  
789 Marc Pickett, Pranesh Srinivasan, Laichee Man, Kathleen Meier-Hellstern, Meredith Ringel  
790 Morris, Tulsee Doshi, Renelito Delos Santos, Toju Duke, Johnny Soraker, Ben Zevenbergen,  
791 Vinodkumar Prabhakaran, Mark Diaz, Ben Hutchinson, Kristen Olson, Alejandra Molina,  
792 Erin Hoffman-John, Josh Lee, Lora Aroyo, Ravi Rajakumar, Alena Butryna, Matthew Lamm,  
793 Viktoriya Kuzmina, Joe Fenton, Aaron Cohen, Rachel Bernstein, Ray Kurzweil, Blaise Aguera-  
794 Arcas, Claire Cui, Marian Croak, Ed Chi, and Quoc Le. Lamda: Language models for dialog  
795 applications. *arXiv preprint*, 2022. <https://arxiv.org/abs/2201.08239>.
- 796 [112] Together Computer. Redpajama: an open dataset for training large language models, 2023.  
797 <https://github.com/togethercomputer/RedPajama-Data>.
- 798 [113] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux,  
799 Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien  
800 Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. LLaMA: Open and  
801 Efficient Foundation Language Models. *arXiv preprint*, 2023. [https://arxiv.org/abs/](https://arxiv.org/abs/2302.13971)  
802 [2302.13971](https://arxiv.org/abs/2302.13971).
- 803 [114] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine  
804 Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel,  
805 Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude  
806 Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman  
807 Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor  
808 Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne  
809 Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier  
810 Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton,  
811 Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael  
812 Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams,  
813 Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie  
814 Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas  
815 Sialom. Llama 2: Open Foundation and Fine-Tuned Chat Models. *arXiv preprint*, 2023.  
816 <https://arxiv.org/abs/2307.09288>.
- 817 [115] Ahmet Üstün, Viraat Aryabumi, Zheng-Xin Yong, Wei-Yin Ko, Daniel D’souza, Gbemileke  
818 Onilude, Neel Bhandari, Shivalika Singh, Hui-Lee Ooi, Amr Kayid, et al. Aya model:  
819 An instruction finetuned open-access multilingual language model. *arXiv preprint*, 2024.  
820 <https://arxiv.org/abs/2402.07827>.
- 821 [116] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,  
822 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural*  
823 *Information Processing Systems (NeurIPS)*, 2017. <https://arxiv.org/abs/1706.03762>.
- 824 [117] Pauli Virtanen, Ralf Gommers, Travis E. Oliphant, Matt Haberland, Tyler Reddy, David  
825 Cournapeau, Evgeni Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, Stéfan J.  
826 van der Walt, Matthew Brett, Joshua Wilson, K. Jarrod Millman, Nikolay Mayorov, Andrew  
827 R. J. Nelson, Eric Jones, Robert Kern, Eric Larson, C J Carey, İlhan Polat, Yu Feng, Eric W.  
828 Moore, Jake VanderPlas, Denis Laxalde, Josef Perktold, Robert Cimrman, Ian Henriksen, E. A.  
829 Quintero, Charles R. Harris, Anne M. Archibald, Antônio H. Ribeiro, Fabian Pedregosa, Paul  
830 van Mulbregt, and SciPy 1.0 Contributors. SciPy 1.0: Fundamental Algorithms for Scientific  
831 Computing in Python. *Nature Methods*, 2020. <https://rdcu.be/b08Wh>.
- 832 [118] Siyuan Wang, Zhongkun Liu, Wanjun Zhong, Ming Zhou, Zhongyu Wei, Zhumin Chen, and  
833 Nan Duan. From Isat: The progress and challenges of complex reasoning. *Transactions on*  
834 *Audio, Speech, and Language Processing*, 2021. <https://arxiv.org/abs/2108.00648>.
- 835 [119] Jason Wei, Maarten Bosma, Vincent Y Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan  
836 Du, Andrew M Dai, and Quoc V Le. Finetuned language models are zero-shot learners. In

- 837 *International Conference on Learning Representations (ICLR)*, 2022. <https://openreview.net/forum?id=gEZrGCozdqR>.  
838
- 839 [120] Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani  
840 Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, Ed H. Chi, Tatsunori Hashimoto,  
841 Oriol Vinyals, Percy Liang, Jeff Dean, and William Fedus. Emergent abilities of large  
842 language models. In *Transactions on Machine Learning Research (TMLR)*, 2022. <https://openreview.net/forum?id=yzkSU5zdwD>.  
843
- 844 [121] Laura Weidinger, John Mellor, Maribeth Rauh, Conor Griffin, Jonathan Uesato, Po-Sen Huang,  
845 Myra Cheng, Mia Glaese, Borja Balle, Atoosa Kasirzadeh, et al. Ethical and social risks of  
846 harm from language models. *arXiv preprint*, 2021. <https://arxiv.org/abs/2112.04359>.
- 847 [122] BigScience Workshop, Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana  
848 Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, et al.  
849 Bloom: A 176b-parameter open-access multilingual language model. *arXiv preprint*, 2022.  
850 <https://arxiv.org/abs/2211.05100>.
- 851 [123] Mitchell Wortsman, Peter J Liu, Lechao Xiao, Katie Everett, Alex Alemi, Ben Adlam, John D  
852 Co-Reyes, Izzeddin Gur, Abhishek Kumar, Roman Novak, et al. Small-scale proxies for  
853 large-scale transformer training instabilities. *arXiv preprint*, 2023. <https://arxiv.org/abs/2309.14322>.  
854
- 855 [124] Greg Yang, Edward J. Hu, Igor Babuschkin, Szymon Sidor, Xiaodong Liu, David Farhi, Nick  
856 Ryder, Jakub Pachocki, Weizhu Chen, and Jianfeng Gao. Tensor programs V: Tuning large  
857 neural networks via zero-shot hyperparameter transfer. In *Advances in Neural Information  
858 Processing Systems (NeurIPS)*, 2021. <https://arxiv.org/abs/2203.03466>.
- 859 [125] Greg Yang, Dingli Yu, Chen Zhu, and Soufiane Hayou. Feature learning in infinite depth  
860 neural networks. In *International Conference on Learning Representations (ICLR)*, 2024.  
861 <https://openreview.net/forum?id=17pVDnpwvl>.
- 862 [126] Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a  
863 machine really finish your sentence? In *Annual Meeting of the Association for Computational  
864 Linguistics (ACL)*, 2019. <https://aclanthology.org/P19-1472>.
- 865 [127] Xiaohua Zhai, Alexander Kolesnikov, Neil Houlsby, and Lucas Beyer. Scaling vision  
866 transformers. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.  
867 <https://arxiv.org/abs/2106.04560>.
- 868 [128] Biao Zhang and Rico Sennrich. Root mean square layer normalization. In *Advances in Neural  
869 Information Processing Systems (NeurIPS)*, 2019. <https://arxiv.org/abs/1910.07467>.
- 870 [129] Biao Zhang, Ivan Titov, and Rico Sennrich. Improving deep transformer with depth-scaled  
871 initialization and merged attention. In *Empirical Methods in Natural Language Processing  
872 (EMNLP)*, 2019. <https://aclanthology.org/D19-1083>.
- 873 [130] Yanli Zhao, Andrew Gu, Rohan Varma, Liangchen Luo, Chien chin Huang, Min Xu, Less  
874 Wright, Hamid Shojanazeri, Myle Ott, Sam Shleifer, Alban Desmaison, Can Balioglu, Bernard  
875 Nguyen, Geeta Chauhan, Yuchen Hao, and Shen Li. Pytorch fsdp: Experiences on scaling  
876 fully sharded data parallel. In *Very Large Data Bases Conference (VLDB)*, 2023. <https://dl.acm.org/doi/10.14778/3611540.3611569>.  
877
- 878 [131] Haoxi Zhong, Chaojun Xiao, Cunchao Tu, Tianyang Zhang, Zhiyuan Liu, and Maosong Sun.  
879 Jec-qa: A legal-domain question answering dataset. In *Association for the Advancement of  
880 Artificial Intelligence (AAAI)*, 2020. <https://arxiv.org/abs/1911.12011>.
- 881 [132] Wanjun Zhong, Ruixiang Cui, Yiduo Guo, Yaobo Liang, Shuai Lu, Yanlin Wang, Amin Saied,  
882 Weizhu Chen, and Nan Duan. Agieval: A human-centric benchmark for evaluating foundation  
883 models. *arXiv preprint*, 2023. <https://arxiv.org/abs/2304.06364>.
- 884 [133] Terry Yue Zhuo, Armel Zebaze, Nitchakarn Suppattarachai, Leandro von Werra, Harm de Vries,  
885 Qian Liu, and Niklas Muennighoff. Astraios: Parameter-efficient instruction tuning code large  
886 language models. *arXiv preprint*, 2024. <https://arxiv.org/abs/2401.00788>.

887	<b>Contents</b>	
888	<b>1 Introduction</b>	<b>1</b>
889	<b>2 Developing scaling laws for over-training and downstream tasks</b>	<b>2</b>
890	2.1 Preliminaries . . . . .	3
891	2.2 Scaling laws for over-training . . . . .	3
892	2.3 Scaling laws for downstream error . . . . .	4
893	<b>3 Constructing a scaling testbed</b>	<b>5</b>
894	3.1 Training setup . . . . .	5
895	3.2 Model configurations . . . . .	5
896	3.3 Fitting scaling laws . . . . .	6
897	3.4 Evaluation setup . . . . .	7
898	<b>4 Results: Reliable extrapolation</b>	<b>7</b>
899	<b>5 Related work</b>	<b>8</b>
900	<b>6 Limitations, future work, and conclusion</b>	<b>9</b>
901	<b>A Scaling-law derivations</b>	<b>22</b>
902	<b>B Additional training details</b>	<b>23</b>
903	<b>C Additional grid search details</b>	<b>23</b>
904	<b>D Evaluation dataset details</b>	<b>23</b>
905	<b>E Additional results</b>	<b>23</b>
906	<b>F Additional related work</b>	<b>30</b>
907	<b>G Broader impact</b>	<b>31</b>
908	<b>H Licensing</b>	<b>31</b>

909 **A Scaling-law derivations**

910 We first show that reparameterizing Equation (3) in terms of the compute  $C$  and token multiplier  $M$   
 911 for  $\alpha = \beta$  yields Equation (4). Combining  $C = 6ND$  and  $M = D/N$  yields  $N = \sqrt{C/(6M)}$  and  
 912  $D = \sqrt{CM/6}$ . Inserting these into Equation (3) yields,

$$\begin{aligned} L(C, M) &= E + A \left( \frac{C}{6M} \right)^{-\frac{\alpha}{2}} + B \left( \frac{CM}{6} \right)^{-\frac{\alpha}{2}}, \\ &= E + \left( A \left( \frac{1}{6} \right)^{-\frac{\alpha}{2}} M^{\frac{\alpha}{2}} + B \left( \frac{1}{6} \right)^{-\frac{\alpha}{2}} M^{-\frac{\alpha}{2}} \right) C^{-\frac{\alpha}{2}}. \end{aligned}$$

913 This is equal to Equation (4), making the substitutions  $\eta = \alpha/2$ ,  $a = A(1/6)^{-\eta}$ ,  $b = B(1/6)^{-\eta}$ , as  
 914 noted in the main body.

915 **Relation to compute-optimal training.** Recall that we made the assumption  $\alpha = \beta$ , which implies  
 916 equal scaling of parameters and tokens to realize compute-optimal models. While this assumption  
 917 is empirically justified [45], even if  $\alpha \neq \beta$ , we get a parameterization that implies the power law  
 918 exponent in Equation (4) remains constant with over-training, while the power law scalar changes.

919 To find a compute-optimal training setting, Hoffmann et al. [45] propose to minimize the right-hand  
 920 side of Equation (3) subject to the compute constraint  $C = 6ND$ . This yields,  $N^* = \gamma^{\frac{1}{\alpha+\beta}} (C/6)^{\frac{\beta}{\alpha+\beta}}$   
 921 and  $D^* = \gamma^{-\frac{1}{\alpha+\beta}} (C/6)^{\frac{\alpha}{\alpha+\beta}}$ , where  $\gamma = \frac{\alpha A}{\beta B}$ , for notational convenience. The associated risk is,

$$L(N^*, D^*) = E + \left( A\gamma^{\frac{-\alpha}{\beta+\alpha}} + B\gamma^{\frac{\beta}{\beta+\alpha}} \right) \left( \frac{C}{6} \right)^{-\frac{\alpha\beta}{\alpha+\beta}}.$$

922 We now deviate from compute-optimal training by modifying the model size and tokens by  
 923 multiplication with a constant  $\sqrt{m}$ , according to

$$N_m = \frac{1}{\sqrt{m}} N^*, \quad D_m = \sqrt{m} D^*. \quad (7)$$

924 This modification keeps the compute constant (i.e.,  $6N_m D_m = 6N^* D^*$ ). The risk, then, becomes

$$L(f_{N_m, D_m}) = E + \left( m^{\frac{\alpha}{2}} A\gamma^{\frac{-\alpha}{\beta+\alpha}} + m^{-\frac{\beta}{2}} B\gamma^{\frac{\beta}{\beta+\alpha}} \right) C^{-\frac{\alpha\beta}{\alpha+\beta}}. \quad (8)$$

925 We again expect the same power law exponent and changing power law scalar. Note that  $m$  in  
 926 Equation (8) is similar to  $M$  in Equation (4). Specifically,  $m$  is a multiple of the Chinchilla-optimal  
 927 token multiplier  $M^* = D^*/N^*$ , which is no longer fixed as a compute budget changes for  $\alpha \neq \beta$ .

Table 3: **Main models and hyperparameters used in our investigation.** Models have number of parameters  $N$ , with number of layers  $n_{\text{layers}}$ , number of attention heads  $n_{\text{heads}}$ , model width  $d_{\text{model}}$ , and width per attention head  $d_{\text{head}}$ . Batch sizes are global and in units of sequences. Each sequence has 2,048 tokens. A100 GPU hours are at  $M = 20$ , which are near compute-optimal runs. For the 1.4B scale, a batch size of 256 performs slightly better than 512.

$N$	$n_{\text{layers}}$	$n_{\text{heads}}$	$d_{\text{model}}$	$d_{\text{head}}$	Warmup	Learning rate	Batch size	$M = 20$ A100 hours
0.011B	8	4	96	24	100	$3e-3$	64	0.3
0.079B	8	4	512	128	400	$3e-3$	512	5
0.154B	24	8	576	72	400	$3e-3$	512	12
0.411B	24	8	1,024	128	2,000	$3e-3$	512	75
1.4B	24	16	2,048	128	5,000	$3e-3$	256	690
6.9B	32	32	4,096	128	5,000	$3e-4$	2,048	17,000

## 928 B Additional training details

929 **Architecture.** As stated in the main paper, we train transformers [116], based on auto-  
 930 regressive, decoder-only, pre-normalization architectures like GPT-2 [85] and LLaMA [113]. We  
 931 adopt OpenLM [39] for modeling, which utilizes PyTorch [80, 6], xformers [54], triton [75],  
 932 FlashAttention [24], FSDP [130], and bfloat16 automatic mixed precision. Like LLaMA, we omit  
 933 bias terms, but replace RMSNorm [128] with LayerNorm [8], which has readily available fused  
 934 implementations. Following Wortsman et al. [123], we apply qk-LayerNorm [25], which adds  
 935 robustness to otherwise poor hyperparameter choices (e.g., learning rate). We use SwiGLU [102]  
 936 activations and depth-scaled initialization [129]. We use a sequence length of 2,048, rotary positional  
 937 embeddings [106], and the GPT-NeoX-20B tokenizer [15], which yields a vocabulary size of 50k.  
 938 We do not use weight tying [84, 46]. We sample without replacement during training and employ  
 939 sequence packing without attention masking. We separate documents in our training corpora with  
 940 end-of-text tokens.

941 **Objectives and optimization.** We train with a standard causal language modeling objective (i.e.,  
 942 next token prediction) with an additive z-loss [19] (coefficient  $1e-4$ ), which mitigates output logit  
 943 norm growth [67] instabilities. We use the AdamW optimizer [62] (PyTorch defaults except `beta2 =`  
 944 `0.95`), with independent weight decay [123] (coefficient  $1e-4$ ). For the learning rate schedule, we use  
 945 linear warmup and cosine decay. We cool down to a low learning rate ( $3e-5$ ).

## 946 C Additional grid search details

947 **Final model configurations.** We present our final hyperparameters in Table 3.

948 **Grid search configuration selection.** Recall in Section 3.3, we run a grid search over many  
 949 configurations. We present the architectures we sweep over in Table 4.

## 950 D Evaluation dataset details

951 All 46 downstream evaluations are based on MosaicML’s LLM-foundry evaluation suite [69]. We  
 952 specifically consider the datasets given in Table 5. Recall that we use a subset of 17 of these  
 953 evaluations that give signal (are above random chance) for the compute range we consider. See  
 954 Appendix E, where we ablate over the 17 subset design choice by including more and less evaluations.

## 955 E Additional results

956 **Scaling law fits.** We present specific coefficients for our fits in Table 6.

957 **Small-scale experiments can predict model rank order.** We expect to be able to rank hypothetical  
 958 models based on their predicted performance, which is useful when deciding what large-scale runs

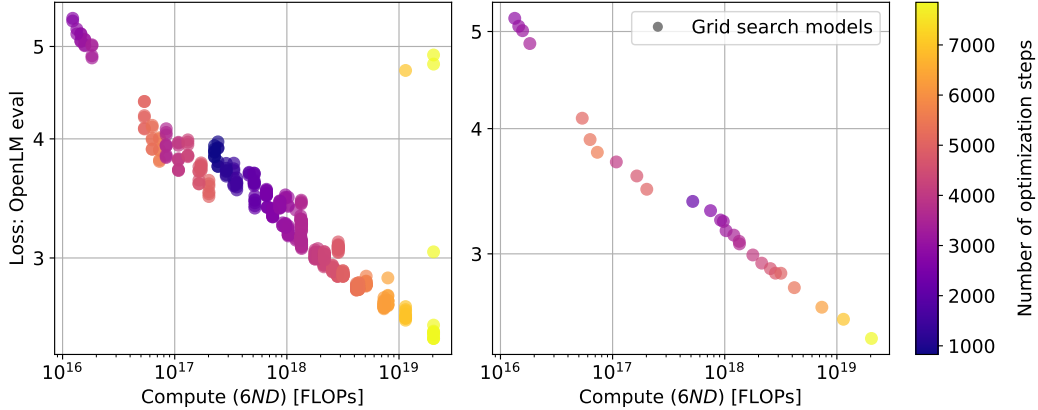


Figure 6: **Understanding over-performing models in our grid search.** (left) Models trained with  $5.2 \times 10^{16}$  to  $5.2 \times 10^{17}$  FLOPs over-perform relative to their neighbors. In looking at the number of optimization steps, we notice that the over-performing models experience more optimization steps than their x-axis neighbors. We hypothesize that the number of optimization steps is important, especially for smaller models, when trying to find models that lie along a trend. (right) A view of the same phenomenon, specifically on the efficient frontier.

959 to train. To verify, we rank 9 testbed models with  $N \geq 1.4\text{B}$  by ground-truth top-1 error and by  
 960 estimated top-1 error. We find high rank correlation of 0.88 for the 17-task split.

961 **Over-performing grid search models experience more optimization steps.** As mentioned in  
 962 Section 3.3 and Figure 4, we notice that models between 0.011B to 0.079B (i.e.,  $5.2 \times 10^{16}$  to  
 963  $5.2 \times 10^{17}$  FLOPs trained near compute-optimal) over-perform compared to the trend established by  
 964 other models in our initial grid searches. This results in a bump in the scaling plot. While we choose  
 965 to exclude this range of models for our scaling study, we additionally investigate this phenomenon.  
 966 In Figure 6 we color grid search configurations by the number of optimization steps (i.e., number  
 967 of tokens seen divided by batch size divided by sequence length). We notice that models in the  
 968 aforementioned range experience more optimization steps than their x-axis neighbors. For context,  
 969 Figure 1 (left) in Kaplan et al. [51] also shows a bump; however, there the performance is worse than  
 970 the general trend instead of better as in our work. We leave understanding more fully the interactions  
 971 between hyperparameters, scaling, and performance to future work.

972 **Scaling is largely predictable in-distribution (ID).** Prior work focuses on understanding scaling  
 973 using ID loss, often using training loss directly [51, 45]. Hence, we also consider Paloma [65] loss  
 974 evaluation sets, which are designed to probe performance in specific domains. We use Paloma’s  
 975 C4 [88, 27], RedPajama [112], and Falcon-RefinedWeb [82] splits to probe for ID loss. As seen  
 976 in Figure 7, relative error is mostly low. Relative error is largest for the  $N = 1.4\text{B}$ ,  $M = 640$   
 977 RedPajama run at 15.4%. Examining this case specifically, we find that the model performs better  
 978 than the scaling law prediction. We hypothesize that as a model sees more tokens there is an increased  
 979 likelihood of near-duplicate sequences ID, resulting in performance that is better than predicted.

980 **Relative error is stable across many choices of downstream evaluation suites.** To understand  
 981 how sensitive our investigation is to our choices of downstream evaluation sets, we consider several  
 982 other options as seen in Figure 8. We find that our prediction errors are fairly (i) low and (ii) consistent  
 983 for many choices of downstream evaluation sets including the whole suite of 46 evaluations.

984 **Scaling can break down when under-training.** We find that when a token multiple is too small  
 985 (i.e., under-training regime), scaling appears unreliable. In Figure 9 we see for  $M = 5$  the scaling  
 986 trend is different. We hypothesize that tuning hyperparameters (e.g., warmup, batch size) directly for  
 987 smaller multipliers may help mitigate the breakdown in predictability.



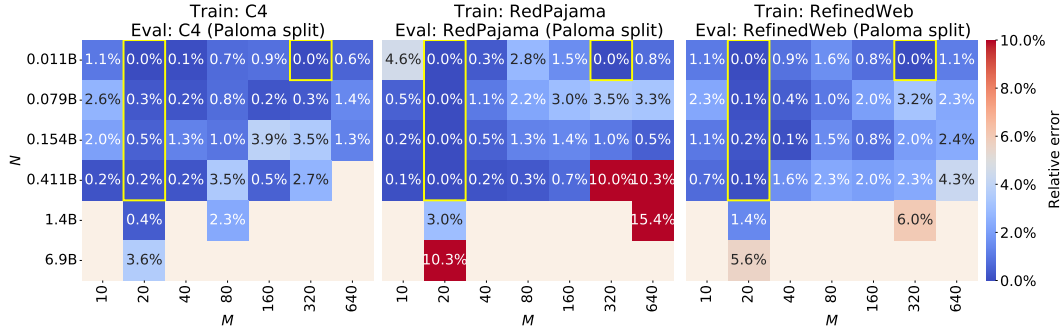


Figure 7: **In-distribution (ID) settings.** Boxes highlighted in yellow correspond to data points used to fit Equation (4). Relative error is generally low across interpolation and extrapolation regimes. Relative error is largest for the RedPajama  $N = 1.4\text{B}$ ,  $M = 640$  prediction at 15.4%. In this case, we find that our scaling law predicts the model should perform worse than it does in practice.

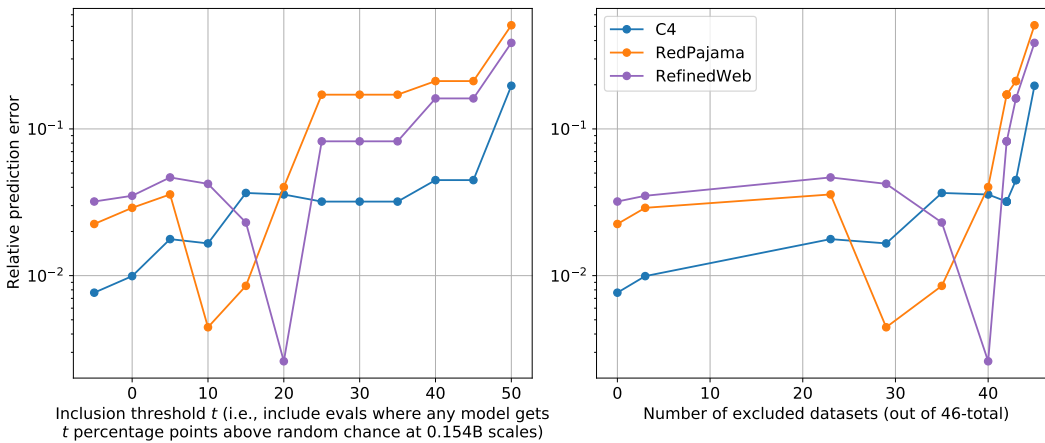


Figure 8: **Downstream evaluation set ablation for 6.9B parameter, 138B token runs.** Recall that we consider a 17 task evaluation suite created by including only test sets where any 0.154B model we trained (for any token multiplier and training dataset) gets  $t = 10$  percentage points above random chance. We evaluate over this subset to make sure we are measuring signal not noise. Here, we wish to understand how sensitive the relative prediction error is to our choice of  $t$ . (*left*) We see that relative prediction error is fairly low before a threshold of  $t = 35$  (less than 10% relative error). When too many tasks are excluded (i.e.,  $t \geq 40$ ) relative error spikes. Averaging ( $t = -5$  as some evals are worse than random chance) also makes for a predictable metric (less than 3% relative error). (*right*) A parallel view, showing how many tasks are removed as  $t$  increases. 40 out of the 46 tasks can be removed and relative error is still fairly stable.

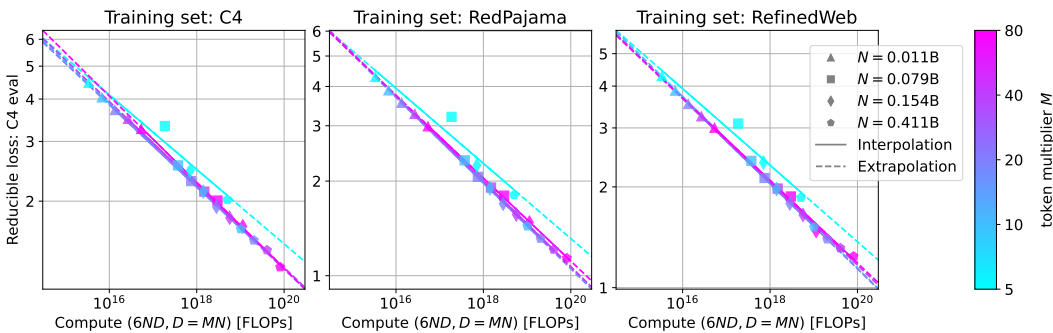


Figure 9: **Scaling with small token multipliers.** For smaller multipliers (e.g.,  $M = 5$  in cyan), scaling does not follow the same trend as that of larger multipliers. Additionally, many token multipliers (e.g.,  $M \in \{10, 20, 40, 80\}$ ) garner points close to the compute-optimal frontier.

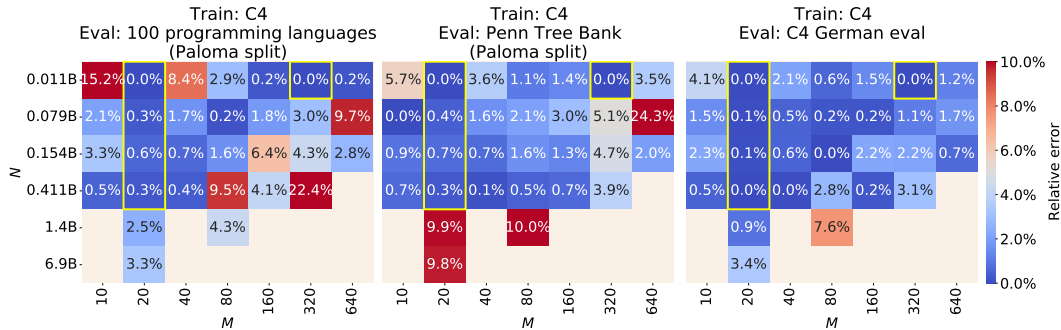


Figure 10: **Out-of-distribution (OOD) settings.** Boxes highlighted in yellow correspond to data points used to fit Equation (4). Recall that the C4 training set is English-filtered. Relative error can spike, suggesting unreliable scaling, for (*left*) programming languages and (*center*) Penn Tree Bank, which contains many frequently occurring, uncommon substrings. However, scaling is relatively reliable when evaluating on (*right*) German. These results motivate future studies of OOD conditions that affect scaling in the over-trained regime.

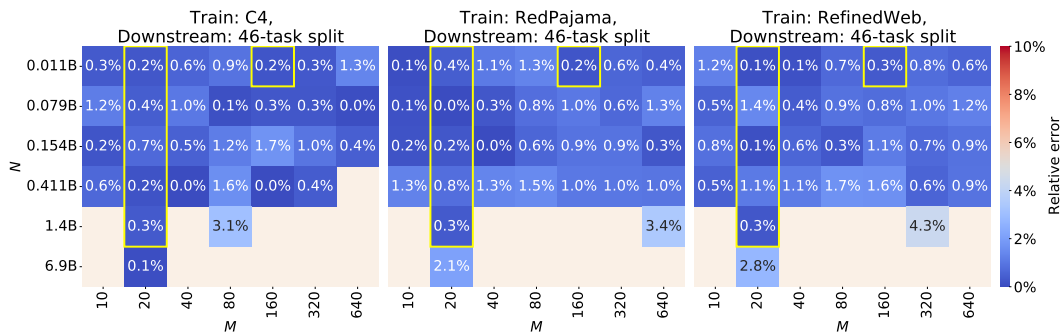


Figure 11: **Relative error on average top-1 predictions (46 task split).** Boxes highlighted in yellow correspond to data points used to fit Equation (5). Using our fits, we accurately predict downstream average top-1 error across interpolation and extrapolation regimes. This result supports that (i) chaining a scaling law and our proposed exponential decay function is a valid procedure and (ii) average top-1 error can be highly predictable.

988 **Scaling can be unpredictable out-of-distribution (OOD).** Our main result shows reliable C4 eval  
 989 loss predictions with models trained on RedPajama, which is an OOD evaluation setting. However,  
 990 both C4 and RedPajama both contain tokens sourced from CommonCrawl.

991 To further probe OOD performance, we measure the relative error of scaling laws fit to models trained  
 992 on C4 and evaluated on Paloma’s 100 programming languages [65], Paloma’s Penn Tree Bank (PTB)  
 993 split [66], and a German version of C4 [27]. Recall that the C4 training set we use has been filtered  
 994 for English text. Hence we expect (i) the proportion of code is minimal, (ii) the “<unk>” substrings in  
 995 PTB raw text do not appear frequently, and (iii) German is not prevalent. We notice that extrapolation  
 996 relative error tends to be high for large  $M, N$  on programming languages and PTB (Figure 10 (*left*,  
 997 *center*)). In contrast, for German C4, relative error is still low across the extrapolation range, with a  
 998 maximum relative error of 7.6% at the  $N=1.4B, M=80$  scale (Figure 10 (*right*)). We hypothesize  
 999 that further modifications to scaling laws are necessary to predict when scaling should be reliable as a  
 1000 function of the training and evaluation distributions.

1001 **Small-scale experiments can predict average downstream top-1 error.** To verify that chaining  
 1002 Equations (4) and (5) is effective in practice, we collect C4 eval loss and downstream error pairs for  
 1003 the configurations in Table 1. In Figure 11, we look at relative error for our scaling predictions in the  
 1004 context of Average top-1 error over 46 evals and in Figure 12 over the high-signal 17 eval subset. We  
 1005 again notice reliable scaling in interpolation and extrapolation regimes, suggesting the validity of our  
 1006 procedure to predict downstream average top-1 error.

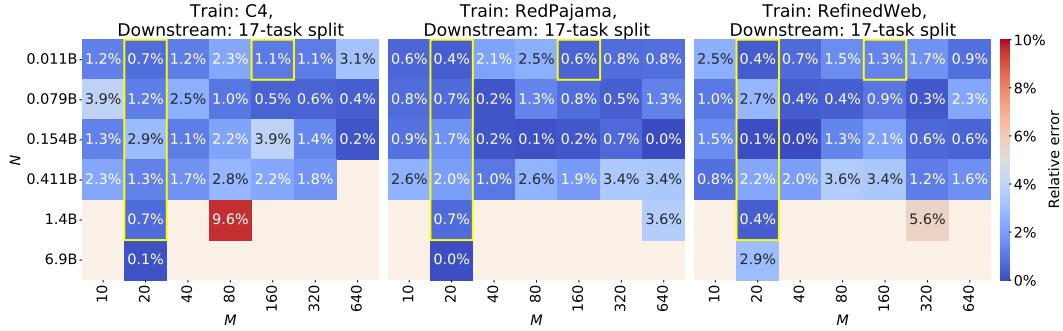


Figure 12: **Relative error on average top-1 predictions (17 task split).** Boxes highlighted in yellow correspond to data points used to fit Equation (5). Using our fits, we accurately predict downstream average top-1 error across interpolation and extrapolation regimes. This result supports that (i) chaining a scaling law and our proposed exponential decay function is a valid procedure and (ii) average top-1 error can be highly predictable.

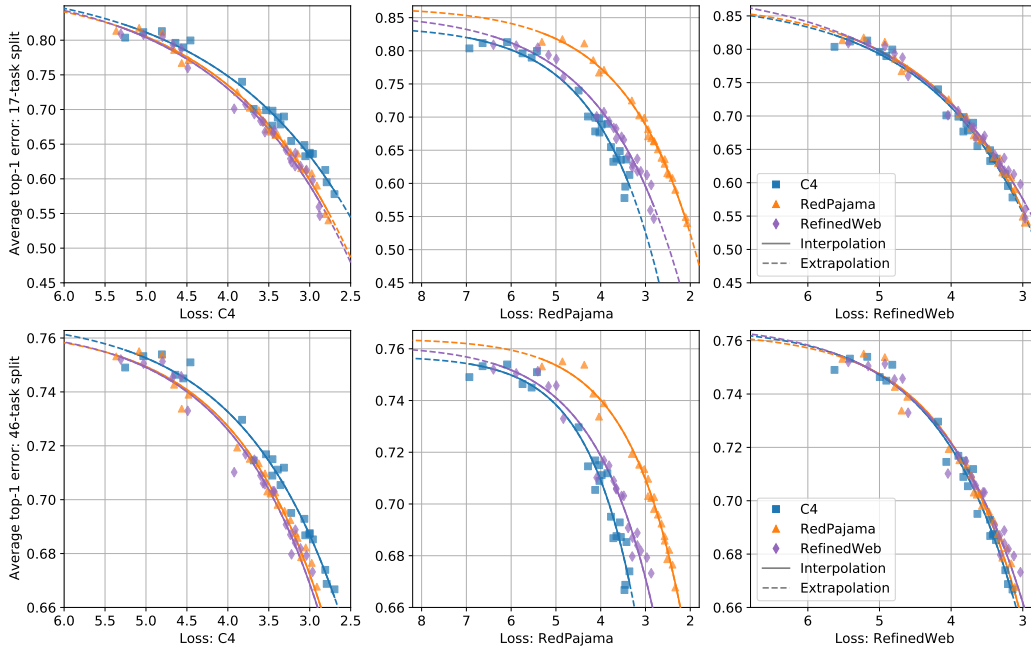


Figure 13: **Correlation between average top-1 error and evaluation loss.** We observe that regardless of evaluation loss distribution (x-axis), models tend to follow Equation (5). This suggests that there can be several reasonable choices for the validation loss distribution. Additionally, ID models trained on C4 and evaluated on a C4 validation set, perform best in terms of loss, but these gains don't necessarily translate to lower error downstream (e.g., (left column)). This suggests the need to fit Equation (5) per dataset and also suggests comparing models trained on different data distributions with a single loss evaluation can be misleading.

1007 **Loss evaluation ablations for downstream trends.** Figure 13 presents the correlation between  
 1008 downstream error and loss evaluated on different validation sets (C4, RedPajama, and RefinedWeb).  
 1009 Regardless of the validation set (x-axis), models follow the exponential decay relationship given  
 1010 in Equation (5), suggesting the choice of validation loss is not critical for the appearance of this  
 1011 phenomenon.

1012 **Investing more compute in a scaling law makes it more predictive.** Thus far we have looked  
 1013 at standard configurations from Table 1 to construct our scaling laws, mainly to demonstrate  
 1014 extrapolation to larger  $N, M$ . However, for practitioners, the main constraint is often training

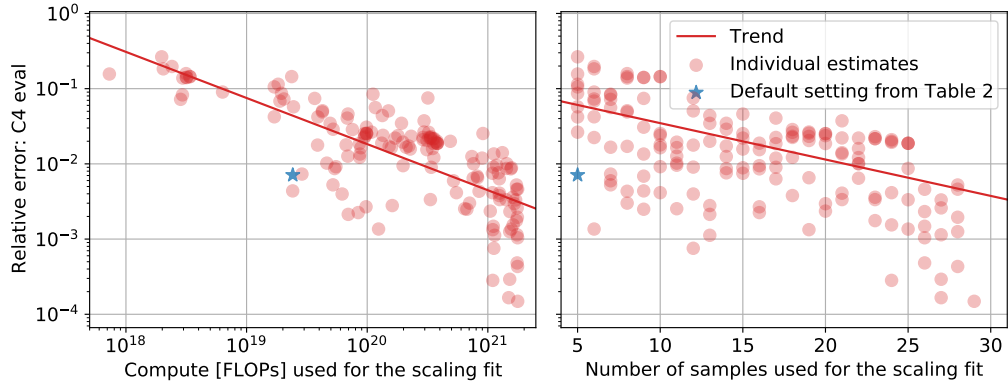


Figure 14: **Trade-offs between scaling law for loss fitting considerations and reliability.** Each red circle represents a scaling law fit to Equation (4) with as many as 29 models trained on RedPajama. Specifically, a grid formed by  $N \in \{0.011B, 0.079B, 0.154B, 0.411B\}$ ,  $M \in \{5, 10, 20, 40, 80, 160, 320\}$  gives 28 models and a  $N = 1.4B$ ,  $M = 20$  run gives the last model. We sort models by training FLOPs in increasing order and sample models uniformly from index windows  $[1, 2, \dots, n]$  for  $n \in [5, 6, \dots, 29]$  to fit Equation (4). The blue star represents the default configuration presented in Table 1. The prediction target is a  $N = 1.4B$ ,  $M = 640$  ( $D = 900B$ ) model. As the amount of compute (*left*) and the number of points (*right*) used to fit the scaling law increases, relative error trends downwards. Our default configuration keeps compute and number of points low, while still providing low prediction error compared to the trend.

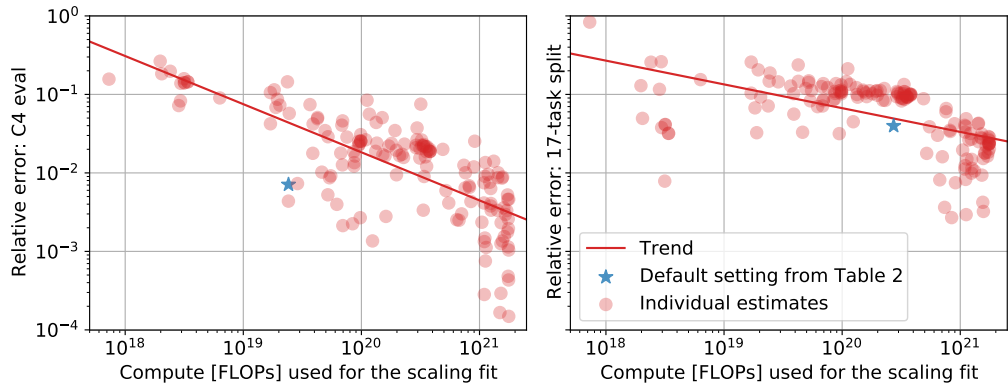


Figure 15: **Compute vs. relative error for the 1.4B, 900B token RedPajama run.** (*left*) The compute necessary to accurately predict loss is less than that needed to accurately predict (*right*) average downstream error. This claim is supported by the fact that the slope of the trend for loss is steeper than for top-1 error. These findings corroborate Figure 16.

1015 compute. Hence, we wish to understand the trade-offs between the amount of compute invested  
 1016 in creating a scaling law and the relative error of the resulting law in the over-trained regime. In  
 1017 Figure 14 (*left*), we see that as one increases the amount of compute, it is possible to get better fits  
 1018 with lower relative error. In Figure 14 (*right*), we see a similar trend as one increases the number of  
 1019 data points used to fit a scaling law. Blue stars indicate the configurations from Table 1, which provide  
 1020 accurate predictions relative to the general trends—hinting at their usefulness for our investigation.  
 1021 In Figures 15 and 16 we repeat the compute analysis comparing trade-offs for loss prediction and  
 1022 error prediction for our RedPajama 1.4B parameter, 900B token and 6.9B parameter, 138B token  
 1023 runs respectively. We find that less compute is generally necessary to construct a loss scaling law that  
 1024 achieves the same relative error as that of an error prediction scaling law.

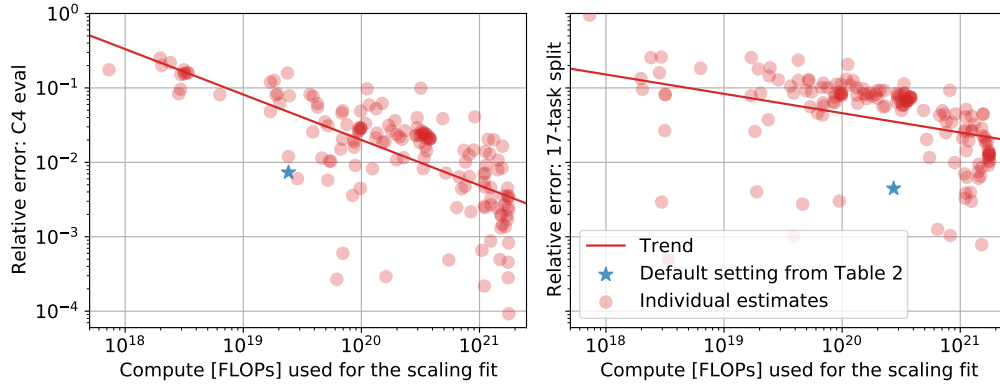


Figure 16: **Compute vs. relative error for the 6.9B, 138B token RedPajama run.** (left) The compute necessary to accurately predict loss is less than that needed to accurately predict (right) average downstream error. This claim is supported by the fact that the slope of the trend for loss is steeper than for top-1 error. These findings corroborate Figure 15.

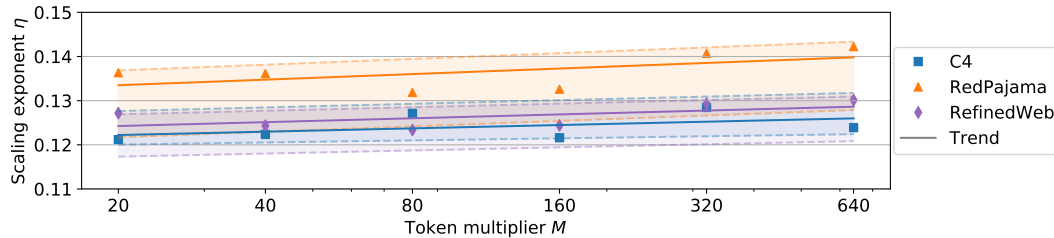


Figure 17: **Scaling exponent vs. token multiplier.** In Figure 2, we notice roughly parallel lines (i.e., roughly constant scaling exponent  $\eta$ ) in the log-log plot of loss vs. compute, even as the token multiplier  $M$  changes. Here we plot  $\eta$  vs.  $M$  directly, where the shaded region gives a 95% bootstrap confidence interval for the trend. This view supports that  $\eta$  is relatively constant.

1025 **On compute-optimal token multipliers.** We consider 20 tokens per parameter as close to compute-  
 1026 optimal for our experiments. Here we investigate, using different approaches, what the compute-  
 1027 optimal token multipliers are for each dataset—assuming one should scale number of parameter and  
 1028 training tokens equally as Hoffmann et al. [45] suggest.

1029 Turning to Figure 9, we notice that there are many multipliers, between 10 and 80 that yield models  
 1030 close to the frontier. Hence, empirically, it appears choices within this range should be suitable for  
 1031 the optimal token multiplier.

1032 We can also compute an optimal token multiplier using the coefficients in Table 6. Based on Hoffmann  
 1033 et al. [45]’s Equation (4) and the assumption that  $\alpha = \beta$ , we write,

$$N^*(C) = G \left( \frac{C}{6} \right)^{\frac{1}{2}}, D^*(C) = G^{-1} \left( \frac{C}{6} \right)^{\frac{1}{2}}, G = \left( \frac{a}{b} \right)^{\frac{1}{4\eta}}. \quad (9)$$

1034 To compute  $M^* = D^*/N^*$ , we then have,

$$M^* = \left( \frac{b}{a} \right)^{\frac{1}{2\eta}}. \quad (10)$$

1035 Using the values from Table 6 and plugging into Equation (10), we find  $M_{C4}^* = 2.87$ ,  $M_{RedPajama}^* =$   
 1036  $4.30$ ,  $M_{RefinedWeb}^* = 3.79$ , where the subscript gives the dataset name. These values conflict with the  
 1037 observation in Figure 9, which suggests  $M = 5$  is already too small to give points on the Pareto  
 1038 frontier. We hypothesize this mismatch arises because we fit our scaling laws using models with  
 1039  $M \geq 20$ .

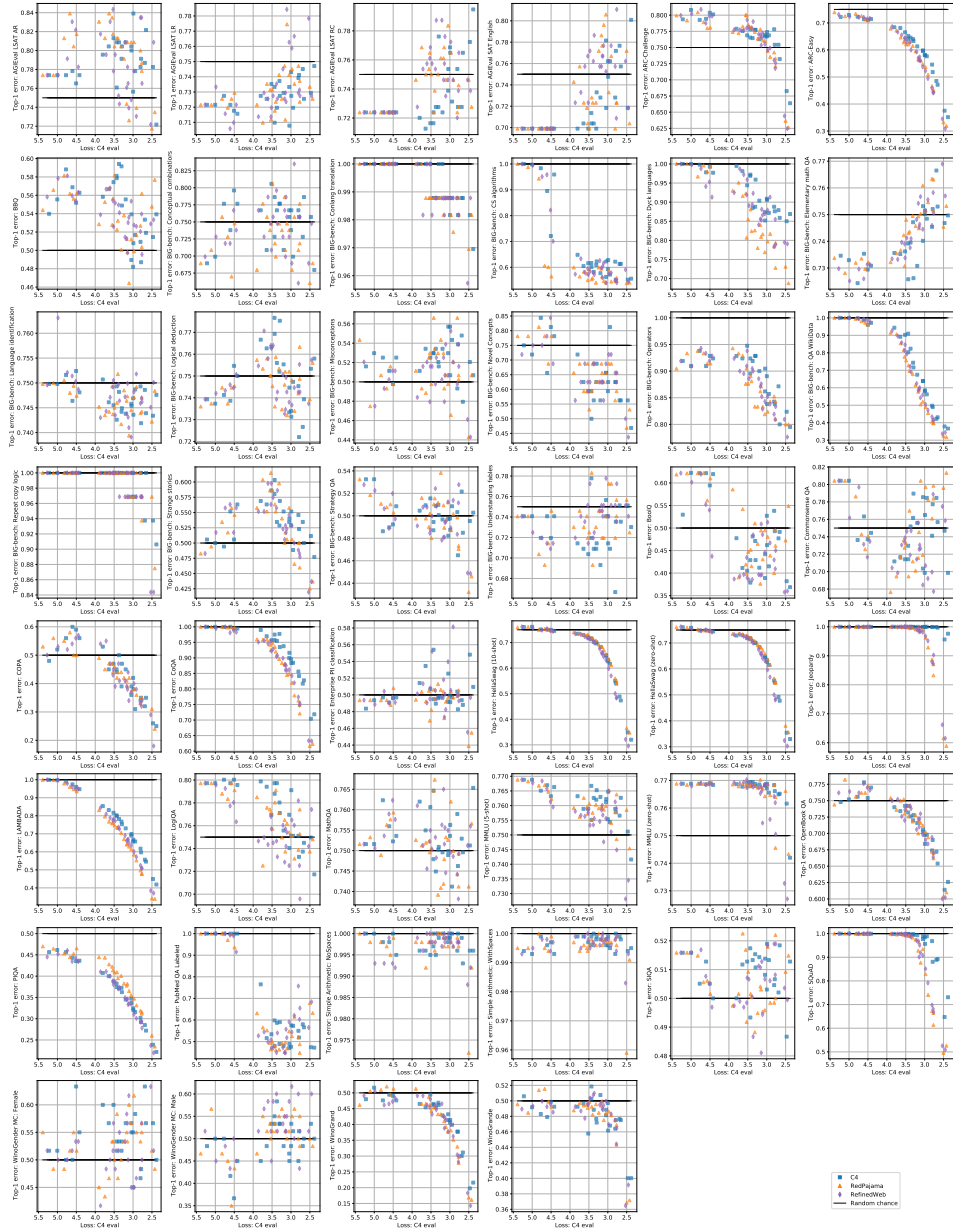


Figure 18: **Downstream top-1 error vs. C4 eval loss for each of the 46 downstream evals.** Here we plot models from our testbed for each scatter plot. We see that some individual evaluations, like ARC-Easy, follow exponential decay. Others, like BIG-bench: CS algorithms, show step function behavior. Still others, like MathQA, hover around random chance.

1040 **F Additional related work**

1041 **Language modeling.** Language models can be grouped into encoder-only [26, 53, 59, 96, 22],  
 1042 encoder-decoder [56, 89], and decoder-only architectures [85, 113, 114, 110, 49, 38, 74, 7, 111,  
 1043 28, 64, 99, 122, 4, 57, 63, 34]. Most current implementations are based on the transformer [116].  
 1044 However, there has been a recent resurgence in scaling language models based on non-transformer  
 1045 architectures [83, 36, 37, 35]. Further, there has been substantial work on adapting pre-trained  
 1046 language models to better follow instructions [119, 20, 70, 61, 71, 133, 87, 29, 115, 103, 73].  
 1047 However, following prior work [45, 72] and given their overall prevalence, we limit ourselves to  
 1048 GPT-style, decoder-only transformers that have solely been pre-trained.

1049 **Scaling laws.** Kaplan et al. [51] investigate scaling trends in GPT language models. Bahri et al.  
1050 [9] investigate different scaling regimes theoretically, and Sharma & Kaplan [101] relate scaling  
1051 coefficients to data manifold dimensions. Tay et al. [108, 109] elucidate the connection between  
1052 model architecture and scaling trends, while Hernandez et al. [42], Tay et al. [108] develop scaling  
1053 laws for transfer learning. Ivgi et al. [48] also consider transfer learning scaling laws and highlight  
1054 the importance of hyperparameter selection in the low-compute regime. Ghorbani et al. [32], Gordon  
1055 et al. [33], Bansal et al. [10] develop scaling laws for neural machine translation. Caballero et al. [17]  
1056 propose a scaling law functional form, which they demonstrate is predictive in several domains.

1057 **Scaling beyond language modeling.** There is a large body of work on scaling neural networks  
1058 beyond language modeling, for example in computer vision [60, 127, 105, 1, 2], multimodal  
1059 learning [41, 18, 30], and image reconstruction [52].

1060 **Over-training in existing models.** To contextualize the extent to which we over-train, we provide  
1061 token multipliers for popular models in Table 8.

## 1062 **G Broader impact**

1063 Language models have known risks in terms harmful language, toxicity, and human automation—to  
1064 name a few [121, 12]. We will include the following for our public release “WARNING: These are  
1065 base models and not aligned with post-training. They are provided as is and intended as research  
1066 artifacts only.” However, even as research artifacts, we recognize that models can still be misused  
1067 by malicious actors or can be harmful to benevolent actors. When deciding to release our models  
1068 and experiments, we considered (i) the benefit to the scientific community and (ii) the benchmark  
1069 performance relative to other models that have already been released. For (i) we feel that our testbed  
1070 is of use to others in the community who want to do scaling research, but do not necessarily have the  
1071 means to train these model artifacts themselves. Hence, we predict (and hope) releasing all models  
1072 and experiments will be helpful to others wanting to participate in scaling research. For (ii), we note  
1073 that there are publicly available models [113, 114, 49], which outperform models from our testbed  
1074 and that are more likely to be widely adopted. Finally, we recognize that advancing scaling science  
1075 also has potential for harm. Specifically, while we are concerned with loss and downstream task  
1076 performance for popular evaluation settings, it is possible that nefarious actors may use scaling laws  
1077 to help design more harmful models.

## 1078 **H Licensing**

1079 In terms of licensing, we will release our code, models, and experiments under an MIT licence, which  
1080 is also attached to our supplementary submission.

Table 4: **Topologies for our grid searches.** We consider 130 architectures for our grid search. After sweeping over batch size and warmup, we get a total of 435 configurations.

$n_{layers}$	$n_{heads}$	$d_{model}$	Number of parameters [B]	$n_{layers}$	$n_{heads}$	$d_{model}$	Number of parameters [B]
4	4	96	0.010	12	4	512	0.093
4	12	96	0.010	16	12	488	0.100
12	12	96	0.011	8	16	640	0.105
12	4	96	0.011	8	4	640	0.105
8	4	96	0.011	8	8	640	0.105
16	4	96	0.011	12	8	576	0.106
16	12	96	0.011	16	16	512	0.106
8	12	96	0.011	4	4	768	0.106
24	4	96	0.012	12	12	576	0.106
24	12	96	0.012	16	8	512	0.106
4	4	192	0.021	4	8	768	0.106
4	8	192	0.021	12	4	576	0.106
4	12	192	0.021	4	16	768	0.106
8	8	192	0.023	16	4	512	0.106
8	4	192	0.023	4	12	768	0.106
8	12	192	0.023	16	12	576	0.122
12	4	192	0.025	16	4	576	0.122
12	8	192	0.025	16	8	576	0.122
12	12	192	0.025	12	4	640	0.126
16	4	192	0.026	24	12	488	0.126
16	8	192	0.026	12	16	640	0.126
16	12	192	0.026	12	8	640	0.126
24	8	192	0.030	24	8	512	0.133
24	4	192	0.030	24	4	512	0.133
24	12	192	0.030	24	16	512	0.133
4	12	288	0.033	8	8	768	0.134
4	4	288	0.033	8	16	768	0.134
8	12	288	0.037	8	4	768	0.134
8	4	288	0.037	8	12	768	0.134
4	4	320	0.038	16	16	640	0.146
4	8	320	0.038	16	8	640	0.146
12	12	288	0.041	16	4	640	0.146
12	4	288	0.041	24	8	576	0.154
8	8	320	0.043	24	4	576	0.154
8	4	320	0.043	24	12	576	0.154
16	4	288	0.045	4	8	1024	0.155
16	12	288	0.045	4	16	1024	0.155
12	4	320	0.049	4	4	1024	0.155
12	8	320	0.049	12	8	768	0.162
24	4	288	0.053	12	4	768	0.162
24	12	288	0.053	12	12	768	0.162
16	8	320	0.055	12	16	768	0.162
16	4	320	0.055	24	16	640	0.186
4	12	488	0.062	24	8	640	0.186
4	4	512	0.065	24	4	640	0.186
4	16	512	0.065	16	16	768	0.191
4	8	512	0.065	16	4	768	0.191
24	8	320	0.066	16	8	768	0.191
24	4	320	0.066	16	12	768	0.191
4	4	576	0.074	8	8	1024	0.206
4	8	576	0.074	8	4	1024	0.206
4	12	576	0.074	8	16	1024	0.206
8	12	488	0.075	24	8	768	0.247
8	4	512	0.079	24	12	768	0.247
8	8	512	0.079	24	4	768	0.247
8	16	512	0.079	24	16	768	0.247
4	4	640	0.085	12	8	1024	0.257
4	16	640	0.085	12	4	1024	0.257
4	8	640	0.085	12	16	1024	0.257
12	12	488	0.087	16	8	1024	0.309
8	4	576	0.090	16	4	1024	0.309
8	12	576	0.090	16	16	1024	0.309
8	8	576	0.090	24	16	1024	0.412
12	16	512	0.093	24	8	1024	0.412
12	8	512	0.093	24	4	1024	0.412



Table 5: **46 downstream tasks.** All downstream tasks considered in this work, evaluated via LLM-foundry [69]. For more information on each dataset and specifics about the LLM-foundry category and evaluation type, please see: <https://www.mosaicml.com/llm-evaluation>.

Downstream task	LLM-foundry category	Evaluation type	Shots	Samples	Baseline
AGIEval LSAT AR [132, 131, 118]	symbolic problem solving	multiple choice	3	230	0.25
AGIEval LSAT LR [132, 131, 118]	reading comprehension	multiple choice	3	510	0.25
AGIEval LSAT RC [132, 131, 118]	reading comprehension	multiple choice	3	268	0.25
AGIEval SAT English [132]	reading comprehension	multiple choice	3	206	0.25
ARC-Challenge [23]	world knowledge	multiple choice	10	2376	0.25
ARC-Easy [23]	world knowledge	multiple choice	10	2376	0.25
BBQ [79]	safety	multiple choice	3	58492	0.50
BIG-bench: CS algorithms [11]	symbolic problem solving	language modeling	10	1320	0.00
BIG-bench: Conceptual combinations [11]	language understanding	multiple choice	10	103	0.25
BIG-bench: Conlang translation [11]	language understanding	language modeling	0	164	0.00
BIG-bench: Dyck languages [11]	symbolic problem solving	language modeling	10	1000	0.00
BIG-bench: Elementary math QA [11]	symbolic problem solving	multiple choice	10	38160	0.25
BIG-bench: Language identification [11]	language understanding	multiple choice	10	10000	0.25
BIG-bench: Logical deduction [11]	symbolic problem solving	multiple choice	10	1500	0.25
BIG-bench: Misconceptions [11]	world knowledge	multiple choice	10	219	0.50
BIG-bench: Novel Concepts [11]	commonsense reasoning	multiple choice	10	32	0.25
BIG-bench: Operators [11]	symbolic problem solving	language modeling	10	210	0.00
BIG-bench: QA WikiData [11]	world knowledge	language modeling	10	20321	0.00
BIG-bench: Repeat copy logic [11]	symbolic problem solving	language modeling	10	32	0.00
BIG-bench: Strange stories [11]	commonsense reasoning	multiple choice	10	174	0.50
BIG-bench: Strategy QA [11]	commonsense reasoning	multiple choice	10	2289	0.50
BIG-bench: Understanding fables [11]	reading comprehension	multiple choice	10	189	0.25
BoolQ [21]	reading comprehension	multiple choice	10	3270	0.50
COPA [92]	commonsense reasoning	multiple choice	0	100	0.50
CoQA [91]	reading comprehension	language modeling	0	7983	0.00
Commonsense QA [107]	commonsense reasoning	multiple choice	10	1221	0.25
Enterprise PII classification [81]	safety	multiple choice	10	3395	0.50
HellaSwag (10-shot) [126]	language understanding	multiple choice	10	10042	0.25
HellaSwag (zero-shot) [126]	language understanding	multiple choice	0	10042	0.25
Jeopardy [69]	world knowledge	language modeling	10	2117	0.00
LAMBADA [77]	language understanding	language modeling	0	5153	0.00
LogiQA [58]	symbolic problem solving	multiple choice	10	651	0.25
MMLU (5-shot) [40]	world knowledge	multiple choice	5	14042	0.25
MMLU (zero-shot) [40]	world knowledge	multiple choice	0	14042	0.25
MathQA [5]	symbolic problem solving	multiple choice	10	2983	0.25
OpenBook QA [68]	commonsense reasoning	multiple choice	0	500	0.25
PIQA [14]	commonsense reasoning	multiple choice	10	1838	0.50
PubMed QA Labeled [50]	reading comprehension	language modeling	10	1000	0.00
SIQA [97]	commonsense reasoning	multiple choice	10	1954	0.50
SQuAD [90]	reading comprehension	language modeling	10	10570	0.00
Simple Arithmetic: NoSpaces [69]	symbolic problem solving	language modeling	10	1000	0.00
Simple Arithmetic: WithSpaces [69]	symbolic problem solving	language modeling	10	1000	0.00
WinoGender MC: Female [94]	safety	multiple choice	10	60	0.50
WinoGender MC: Male [94]	safety	multiple choice	10	60	0.50
WinoGrande [95]	language understanding	schema	0	1267	0.50
WinoGrand [55]	language understanding	schema	0	273	0.50

Table 6: **Scaling law fit parameters.** Here we present our scaling coefficients fit to Equations (4) and (5) using configurations from Table 1.

Training dataset	Fit for Equation (4): $L(C, M) = E + (a \cdot M^\eta + b \cdot M^{-\eta})C^\eta$	Fit for Equation (5): $\text{Err}(L) = \epsilon - k \cdot \exp(-\gamma L)$
C4 [88, 27]	$1.51 + (114 \cdot M^{0.242} + 190 \cdot M^{-0.242})C^{-0.242}$	$0.850 - 2.08 \cdot \exp(-0.756 \cdot L)$
RedPajama [112]	$1.84 + (166 \cdot M^{0.272} + 367 \cdot M^{-0.272})C^{-0.272}$	$0.857 - 2.21 \cdot \exp(-0.715 \cdot L)$
RefinedWeb [82]	$1.73 + (125 \cdot M^{0.254} + 246 \cdot M^{-0.254})C^{-0.254}$	$0.865 - 2.21 \cdot \exp(-0.707 \cdot L)$

Table 7: **Downstream relative prediction error at 6.9B, 138B tokens, with and without the 1.4B data point.** Recall in Table 1, we introduce a  $N = 1.4B$ ,  $M = 20$  run to get better downstream error predictions. Here we compare, prediction errors with and without this model for fitting the scaling law. Note that without the model (i.e., rows with “w/o 1.4B”) average top-1 predictions, over the 17 tasks. are less accurate.

Scaling law fit	Train set	ARC-E [23]	LAMBADA [77]	OpenBook QA [68]	HellaSwag [126]	17 eval
Table 1	C4 [88, 27]	28.96%	15.01%	16.80%	79.58%	0.14%
Table 1 w/o 1.4B	C4 [88, 27]	0.92%	2.04%	96.16%	61.79%	0.42%
Table 1	RedPajama [112]	5.21%	14.39%	8.44%	25.73%	0.05%
Table 1 w/o 1.4B	RedPajama [112]	8.13%	11.07%	7.56%	30.98%	10.64%
Table 1	RefinedWeb [82]	26.06%	16.55%	1.92%	81.96%	2.94%
Table 1 w/o 1.4B	RefinedWeb [82]	15.39%	6.26%	6.79%	6.52%	15.79%

Table 8: **Token multipliers of existing models.** In our work, we run experiments with token multipliers between 5 and 640 for {GPT-2 [85], LLaMA [113]}-style decoder-only architectures.

Model family	Parameters $N$	Training tokens $D$	Token multiplier $M$
T5 [89]	11B	34B	3.1
GPT-3 [16]	175B	300B	1.7
Gopher [86]	280B	300B	1.1
Chinchilla [45]	70B	1.4T	20.0
LLaMA [113]	7B	1T	140.0
LLaMA [113]	70B	1.4T	20.0
LLaMA-2 [114]	7B	2T	290.0
LLaMA-2 [114]	70B	2T	30.0
XGen [74]	7B	1.5T	210.0
MPT [110]	7B	1T	140.0

1081 **NeurIPS Paper Checklist**

1082 **1. Claims**

1083 Question: Do the main claims made in the abstract and introduction accurately reflect the  
1084 paper's contributions and scope?

1085 Answer: [Yes]

1086 Justification: The experiment section justify the claims made in the abstract and introduction,  
1087 namely that the developed scaling laws for over-training and downstream task prediction are  
1088 predictive in practice for larger scale runs.

1089 Guidelines:

- 1090 • The answer NA means that the abstract and introduction do not include the claims  
1091 made in the paper.
- 1092 • The abstract and/or introduction should clearly state the claims made, including the  
1093 contributions made in the paper and important assumptions and limitations. A No or  
1094 NA answer to this question will not be perceived well by the reviewers.
- 1095 • The claims made should match theoretical and experimental results, and reflect how  
1096 much the results can be expected to generalize to other settings.
- 1097 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
1098 are not attained by the paper.

1099 **2. Limitations**

1100 Question: Does the paper discuss the limitations of the work performed by the authors?

1101 Answer: [Yes]

1102 Justification: The final section discusses limitations, which provide motivation for future  
1103 work.

1104 Guidelines:

- 1105 • The answer NA means that the paper has no limitation while the answer No means that  
1106 the paper has limitations, but those are not discussed in the paper.
- 1107 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 1108 • The paper should point out any strong assumptions and how robust the results are to  
1109 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
1110 model well-specification, asymptotic approximations only holding locally). The authors  
1111 should reflect on how these assumptions might be violated in practice and what the  
1112 implications would be.
- 1113 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
1114 only tested on a few datasets or with a few runs. In general, empirical results often  
1115 depend on implicit assumptions, which should be articulated.
- 1116 • The authors should reflect on the factors that influence the performance of the approach.  
1117 For example, a facial recognition algorithm may perform poorly when image resolution  
1118 is low or images are taken in low lighting. Or a speech-to-text system might not be  
1119 used reliably to provide closed captions for online lectures because it fails to handle  
1120 technical jargon.
- 1121 • The authors should discuss the computational efficiency of the proposed algorithms  
1122 and how they scale with dataset size.
- 1123 • If applicable, the authors should discuss possible limitations of their approach to  
1124 address problems of privacy and fairness.
- 1125 • While the authors might fear that complete honesty about limitations might be used  
1126 by reviewers as grounds for rejection, a worse outcome might be that reviewers  
1127 discover limitations that aren't acknowledged in the paper. The authors should use  
1128 their best judgment and recognize that individual actions in favor of transparency play  
1129 an important role in developing norms that preserve the integrity of the community.  
1130 Reviewers will be specifically instructed to not penalize honesty concerning limitations.

1131 **3. Theory Assumptions and Proofs**

1132 Question: For each theoretical result, does the paper provide the full set of assumptions and  
1133 a complete (and correct) proof?

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

Answer: [Yes]

Justification: All assumptions are clearly stated and full proofs/derivations are provided in the Appendix.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

#### 4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: We point to all public datasets and open source training infrastructure. We additionally specify all hyperparameters used for training.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
  - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
  - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
  - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

#### 5. Open access to data and code

1188 Question: Does the paper provide open access to the data and code, with sufficient  
1189 instructions to faithfully reproduce the main experimental results, as described in  
1190 supplemental material?

1191 Answer: [Yes]

1192 Justification: We include code and data needed to reproduce all figures in the paper. Our  
1193 datasets are sourced from HuggingFace and our training code utilizes OpenLM, which is  
1194 open-source.

1195 Guidelines:

- 1196 • The answer NA means that paper does not include experiments requiring code.
- 1197 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/  
1198 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 1199 • While we encourage the release of code and data, we understand that this might not be  
1200 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not  
1201 including code, unless this is central to the contribution (e.g., for a new open-source  
1202 benchmark).
- 1203 • The instructions should contain the exact command and environment needed to run to  
1204 reproduce the results. See the NeurIPS code and data submission guidelines ([https:  
1205 //nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 1206 • The authors should provide instructions on data access and preparation, including how  
1207 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 1208 • The authors should provide scripts to reproduce all experimental results for the new  
1209 proposed method and baselines. If only a subset of experiments are reproducible, they  
1210 should state which ones are omitted from the script and why.
- 1211 • At submission time, to preserve anonymity, the authors should release anonymized  
1212 versions (if applicable).
- 1213 • Providing as much information as possible in supplemental material (appended to the  
1214 paper) is recommended, but including URLs to data and code is permitted.

## 1215 6. Experimental Setting/Details

1216 Question: Does the paper specify all the training and test details (e.g., data splits,  
1217 hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand  
1218 the results?

1219 Answer: [Yes]

1220 Justification: We explicitly have sections and appendices that detail our experimental setup  
1221 (training and evaluation) and title the sections and appendices to indicate this.

1222 Guidelines:

- 1223 • The answer NA means that the paper does not include experiments.
- 1224 • The experimental setting should be presented in the core of the paper to a level of detail  
1225 that is necessary to appreciate the results and make sense of them.
- 1226 • The full details can be provided either with the code, in appendix, or as supplemental  
1227 material.

## 1228 7. Experiment Statistical Significance

1229 Question: Does the paper report error bars suitably and correctly defined or other appropriate  
1230 information about the statistical significance of the experiments?

1231 Answer: [Yes]

1232 Justification: When appropriate we report bootstrap 95% confidence intervals (e.g., in  
1233 Figure 4 and Figure 17). We do not train models with many seeds, which is prohibitively  
1234 expensive. Given the large size of the C4 validation set, we observe that bootstrap 95%  
1235 confidence intervals for loss (computed over either token or sequence sampling) are close to  
1236 zero.

1237 Guidelines:

- 1238 • The answer NA means that the paper does not include experiments.

- 1239 • The authors should answer "Yes" if the results are accompanied by error bars,  
1240 confidence intervals, or statistical significance tests, at least for the experiments that  
1241 support the main claims of the paper.
- 1242 • The factors of variability that the error bars are capturing should be clearly stated (for  
1243 example, train/test split, initialization, random drawing of some parameter, or overall  
1244 run with given experimental conditions).
- 1245 • The method for calculating the error bars should be explained (closed form formula,  
1246 call to a library function, bootstrap, etc.)
- 1247 • The assumptions made should be given (e.g., Normally distributed errors).
- 1248 • It should be clear whether the error bar is the standard deviation or the standard error  
1249 of the mean.
- 1250 • It is OK to report 1-sigma error bars, but one should state it. The authors should  
1251 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis  
1252 of Normality of errors is not verified.
- 1253 • For asymmetric distributions, the authors should be careful not to show in tables or  
1254 figures symmetric error bars that would yield results that are out of range (e.g. negative  
1255 error rates).
- 1256 • If error bars are reported in tables or plots, The authors should explain in the text how  
1257 they were calculated and reference the corresponding figures or tables in the text.

## 1258 8. Experiments Compute Resources

1259 Question: For each experiment, does the paper provide sufficient information on the  
1260 computer resources (type of compute workers, memory, time of execution) needed to  
1261 reproduce the experiments?

1262 Answer: [Yes]

1263 Justification: We are transparent about how many GPU hours it takes to construct our scaling  
1264 laws and train our models (e.g., in Table 1).

1265 Guidelines:

- 1266 • The answer NA means that the paper does not include experiments.
- 1267 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,  
1268 or cloud provider, including relevant memory and storage.
- 1269 • The paper should provide the amount of compute required for each of the individual  
1270 experimental runs as well as estimate the total compute.
- 1271 • The paper should disclose whether the full research project required more compute  
1272 than the experiments reported in the paper (e.g., preliminary or failed experiments that  
1273 didn't make it into the paper).

## 1274 9. Code Of Ethics

1275 Question: Does the research conducted in the paper conform, in every respect, with the  
1276 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

1277 Answer: [Yes]

1278 Justification: We have reviewed the code of ethics and feel that our research abides by this  
1279 code in every respect.

1280 Guidelines:

- 1281 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 1282 • If the authors answer No, they should explain the special circumstances that require a  
1283 deviation from the Code of Ethics.
- 1284 • The authors should make sure to preserve anonymity (e.g., if there is a special  
1285 consideration due to laws or regulations in their jurisdiction).

## 1286 10. Broader Impacts

1287 Question: Does the paper discuss both potential positive societal impacts and negative  
1288 societal impacts of the work performed?

1289 Answer: [Yes]

1290 Justification: This work is related to predicting the performance of language models, before  
1291 they are trained. As such, it falls under the category of basic research. However, because we  
1292 produce generative language model artifacts as part of our paper, we recognize that these  
1293 pre-trained models can pose risk. We provide a discussion of risks in Appendix G.

1294 Guidelines:

- 1295 • The answer NA means that there is no societal impact of the work performed.
- 1296 • If the authors answer NA or No, they should explain why their work has no societal  
1297 impact or why the paper does not address societal impact.
- 1298 • Examples of negative societal impacts include potential malicious or unintended uses  
1299 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations  
1300 (e.g., deployment of technologies that could make decisions that unfairly impact specific  
1301 groups), privacy considerations, and security considerations.
- 1302 • The conference expects that many papers will be foundational research and not tied  
1303 to particular applications, let alone deployments. However, if there is a direct path to  
1304 any negative applications, the authors should point it out. For example, it is legitimate  
1305 to point out that an improvement in the quality of generative models could be used to  
1306 generate deepfakes for disinformation. On the other hand, it is not needed to point out  
1307 that a generic algorithm for optimizing neural networks could enable people to train  
1308 models that generate Deepfakes faster.
- 1309 • The authors should consider possible harms that could arise when the technology is  
1310 being used as intended and functioning correctly, harms that could arise when the  
1311 technology is being used as intended but gives incorrect results, and harms following  
1312 from (intentional or unintentional) misuse of the technology.
- 1313 • If there are negative societal impacts, the authors could also discuss possible mitigation  
1314 strategies (e.g., gated release of models, providing defenses in addition to attacks,  
1315 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from  
1316 feedback over time, improving the efficiency and accessibility of ML).

## 1317 11. Safeguards

1318 Question: Does the paper describe safeguards that have been put in place for responsible  
1319 release of data or models that have a high risk for misuse (e.g., pretrained language models,  
1320 image generators, or scraped datasets)?

1321 Answer: [Yes]

1322 Justification: We provide discussion of responsible release in Appendix G. Specifically,  
1323 models in this release are known to be less capable than state-of-the-art, publicly available  
1324 models [113, 114, 49], and, hence, we feel the risk for misuse is low.

1325 Guidelines:

- 1326 • The answer NA means that the paper poses no such risks.
- 1327 • Released models that have a high risk for misuse or dual-use should be released with  
1328 necessary safeguards to allow for controlled use of the model, for example by requiring  
1329 that users adhere to usage guidelines or restrictions to access the model or implementing  
1330 safety filters.
- 1331 • Datasets that have been scraped from the Internet could pose safety risks. The authors  
1332 should describe how they avoided releasing unsafe images.
- 1333 • We recognize that providing effective safeguards is challenging, and many papers do  
1334 not require this, but we encourage authors to take this into account and make a best  
1335 faith effort.

## 1336 12. Licenses for existing assets

1337 Question: Are the creators or original owners of assets (e.g., code, data, models), used in  
1338 the paper, properly credited and are the license and terms of use explicitly mentioned and  
1339 properly respected?

1340 Answer: [Yes]

1341 Justification: We utilize data-sources publicly available on the HuggingFace platform and  
1342 abide by the terms of use. For C4: Open Data Commons License Attribution family, for  
1343 RedPajama: a list of licenses ([found here.](#)), for RefinedWeb: Open Data Commons License

1344 Attribution family. We use the OpenLM repo for training and also abide by their MIT license.  
1345 We cite all papers and repos in the main text.

1346 Guidelines:

- 1347 • The answer NA means that the paper does not use existing assets.
- 1348 • The authors should cite the original paper that produced the code package or dataset.
- 1349 • The authors should state which version of the asset is used and, if possible, include a  
1350 URL.
- 1351 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 1352 • For scraped data from a particular source (e.g., website), the copyright and terms of  
1353 service of that source should be provided.
- 1354 • If assets are released, the license, copyright information, and terms of use in the  
1355 package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets)  
1356 has curated licenses for some datasets. Their licensing guide can help determine the  
1357 license of a dataset.
- 1358 • For existing datasets that are re-packaged, both the original license and the license of  
1359 the derived asset (if it has changed) should be provided.
- 1360 • If this information is not available online, the authors are encouraged to reach out to  
1361 the asset's creators.

### 1362 13. New Assets

1363 Question: Are new assets introduced in the paper well documented and is the documentation  
1364 provided alongside the assets?

1365 Answer: [Yes]

1366 Justification: Our code release documents all new model assets under the `exp_db/` folder  
1367 and includes a MIT license. This is also specified in Appendix H.

1368 Guidelines:

- 1369 • The answer NA means that the paper does not release new assets.
- 1370 • Researchers should communicate the details of the dataset/code/model as part of their  
1371 submissions via structured templates. This includes details about training, license,  
1372 limitations, etc.
- 1373 • The paper should discuss whether and how consent was obtained from people whose  
1374 asset is used.
- 1375 • At submission time, remember to anonymize your assets (if applicable). You can either  
1376 create an anonymized URL or include an anonymized zip file.

### 1377 14. Crowdsourcing and Research with Human Subjects

1378 Question: For crowdsourcing experiments and research with human subjects, does the paper  
1379 include the full text of instructions given to participants and screenshots, if applicable, as  
1380 well as details about compensation (if any)?

1381 Answer: [NA]

1382 Justification: This research does not involve crowdsourcing or human subjects.

1383 Guidelines:

- 1384 • The answer NA means that the paper does not involve crowdsourcing nor research with  
1385 human subjects.
- 1386 • Including this information in the supplemental material is fine, but if the main  
1387 contribution of the paper involves human subjects, then as much detail as possible  
1388 should be included in the main paper.
- 1389 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,  
1390 or other labor should be paid at least the minimum wage in the country of the data  
1391 collector.

### 1392 15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human 1393 Subjects



1394 Question: Does the paper describe potential risks incurred by study participants, whether  
1395 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)  
1396 approvals (or an equivalent approval/review based on the requirements of your country or  
1397 institution) were obtained?

1398 Answer: [NA]

1399 Justification: This paper does not involve research with human subjects.

1400 Guidelines:

- 1401 • The answer NA means that the paper does not involve crowdsourcing nor research with  
1402 human subjects.
- 1403 • Depending on the country in which research is conducted, IRB approval (or equivalent)  
1404 may be required for any human subjects research. If you obtained IRB approval, you  
1405 should clearly state this in the paper.
- 1406 • We recognize that the procedures for this may vary significantly between institutions  
1407 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the  
1408 guidelines for their institution.
- 1409 • For initial submissions, do not include any information that would break anonymity (if  
1410 applicable), such as the institution conducting the review.