



Section 09

Xiang Pan

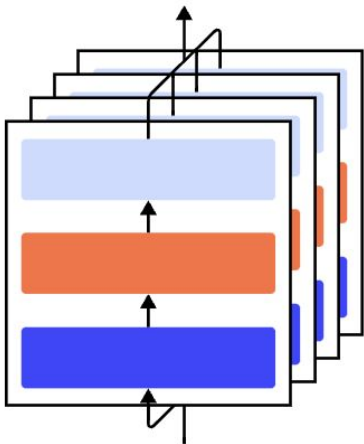
11/06/2023

PART 01

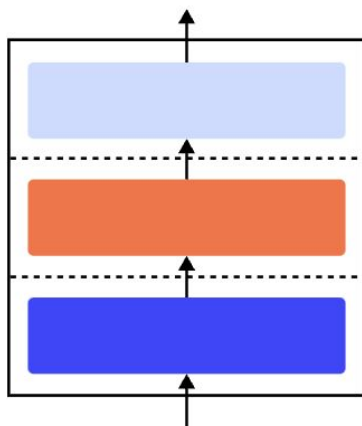
LLM Training

LLM Training

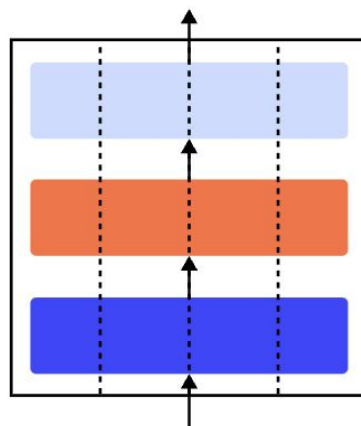
Data Parallelism



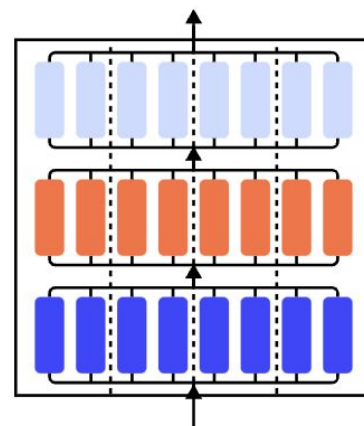
Pipeline Parallelism



Tensor Parallelism



Expert Parallelism



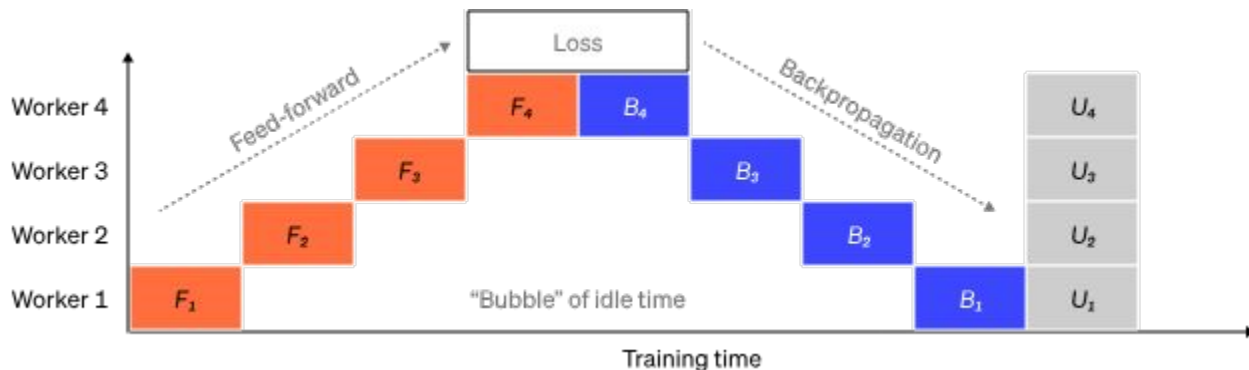
An illustration of various parallelism strategies on a three-layer model. Each color refers to one layer and dashed lines separate different GPUs.

Data Parallelism

- **Model Fits into a single GPU memory**
 - Keep multiple copies of parameters
- **Steps:**
 - independently compute the gradient on each worker;
 - average the gradients across workers;
 - **synchronous overhead!**
 - asynchronous synchronization schemes?
 - independently compute the same new parameters on each worker.

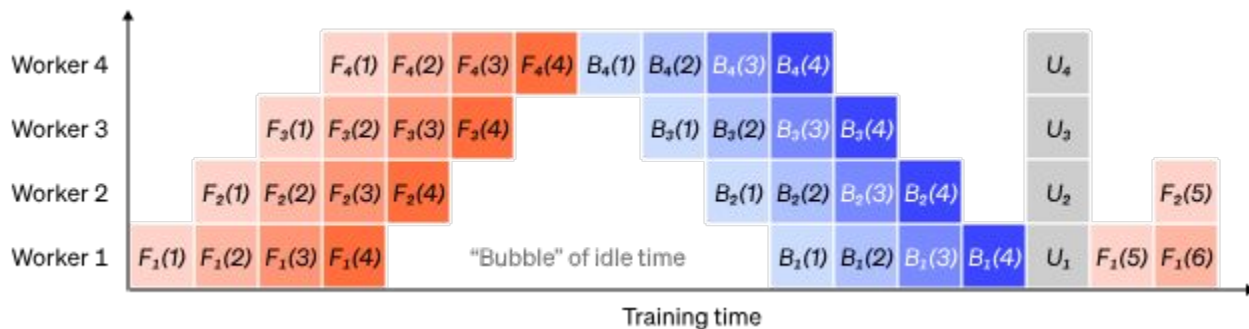
Pipeline Parallelism

- waiting for outputs from the previous machine



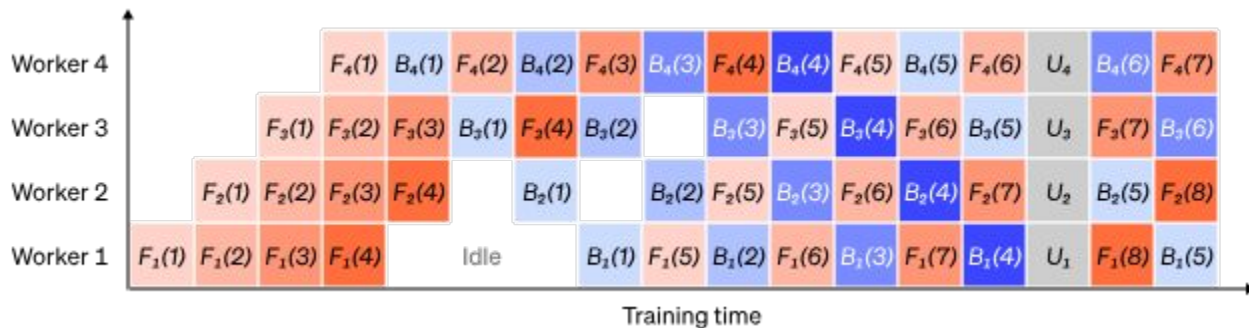
Pipeline Parallelism

□ GPipe



Pipeline Parallelism

□ PipeDream



Tensor Parallelism

- Transformer/CNN -> Matrix Multiplication
- Sequence parallelism:

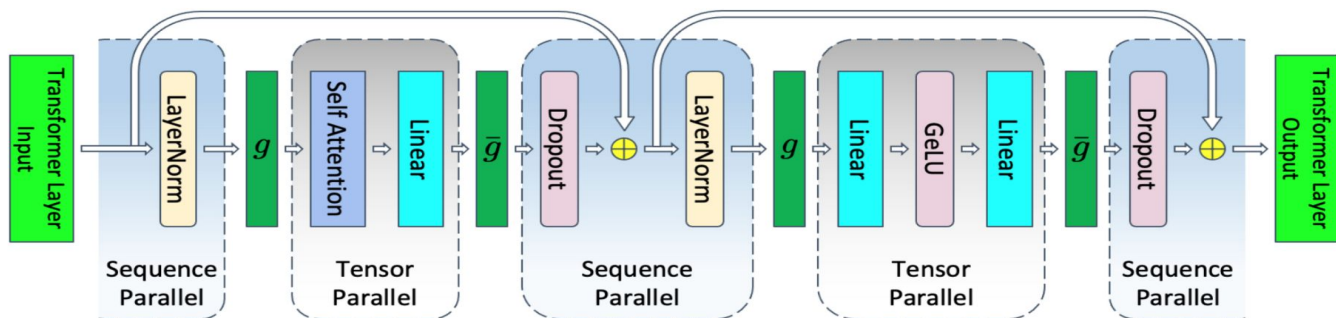
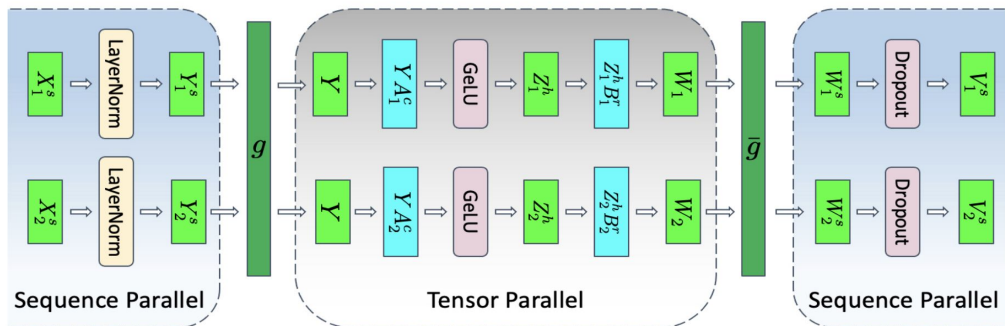


Figure 5: Transformer layer with tensor and sequence parallelism. g and \bar{g} are conjugate. g is all-gather in the forward pass and reduce-scatter in the backward pass. \bar{g} is reduce-scatter in forward pass and all-gather in backward pass.

Tensor Parallelism



$$[Y_1^s, Y_2^s] = \text{LayerNorm}([X_1^s, X_2^s]),$$

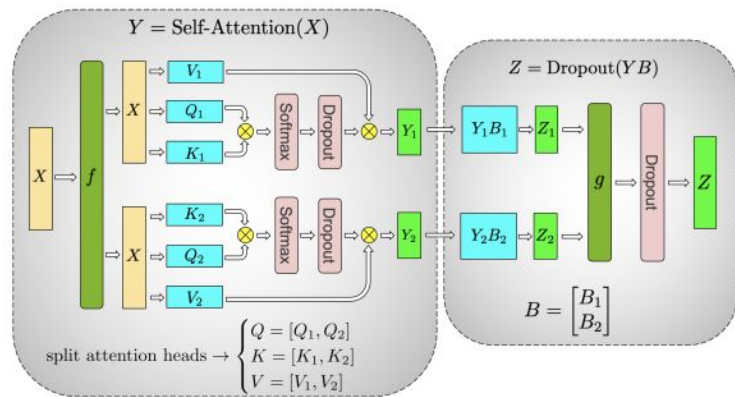
$$Y = g(Y_1^s, Y_2^s),$$

$$[Z_1^h, Z_2^h] = [\text{GeLU}(Y A_1^c), \text{GeLU}(Y A_2^c)],$$

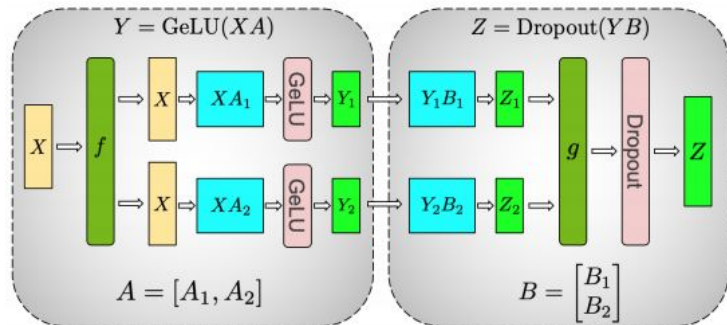
$$W_1 = Z_1^h B_1^r \text{ and } W_2 = Z_2^h B_2^r,$$

$$[W_1^s, W_2^s] = \bar{g}(W_1, W_2),$$

$$[V_1^s, V_2^s] = [\text{Dropout}(W_1^s), \text{Dropout}(W_2^s)].$$

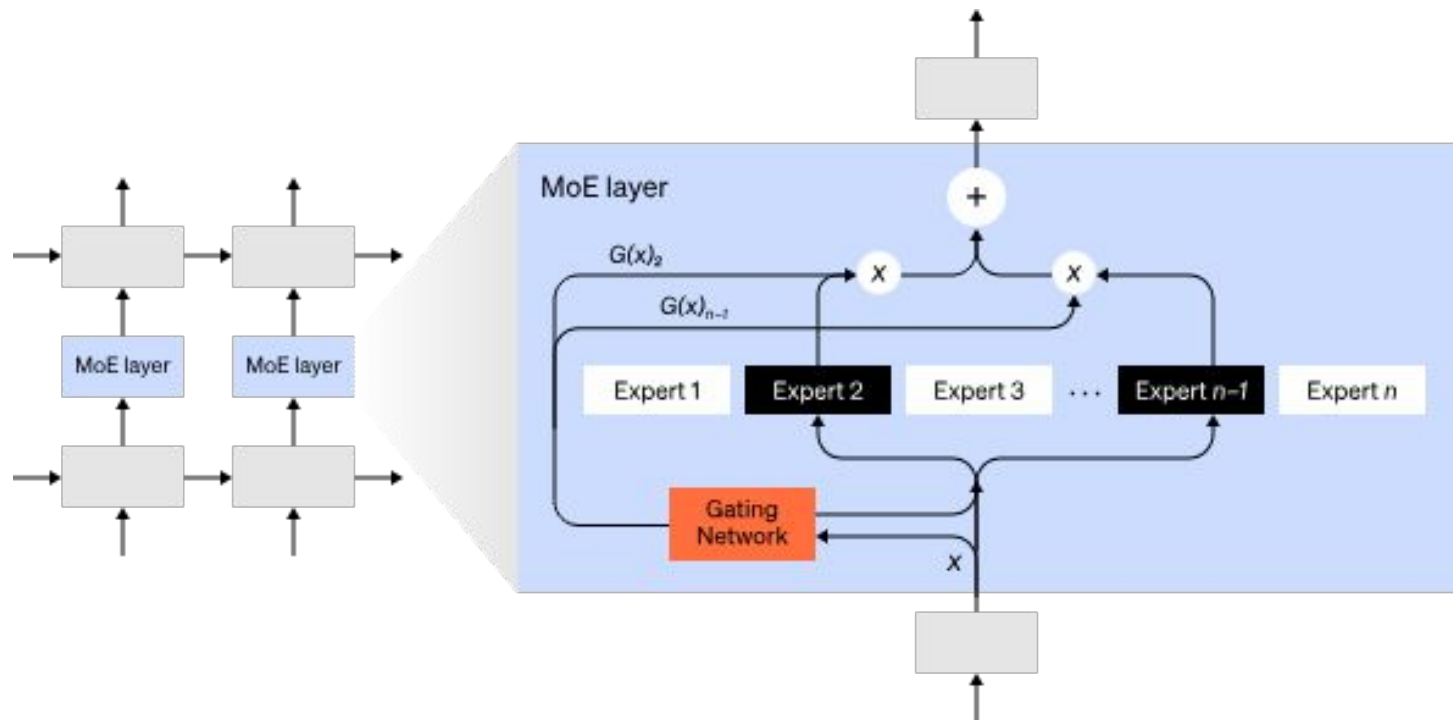


(b) Self-Attention



(a) MLP

MOE

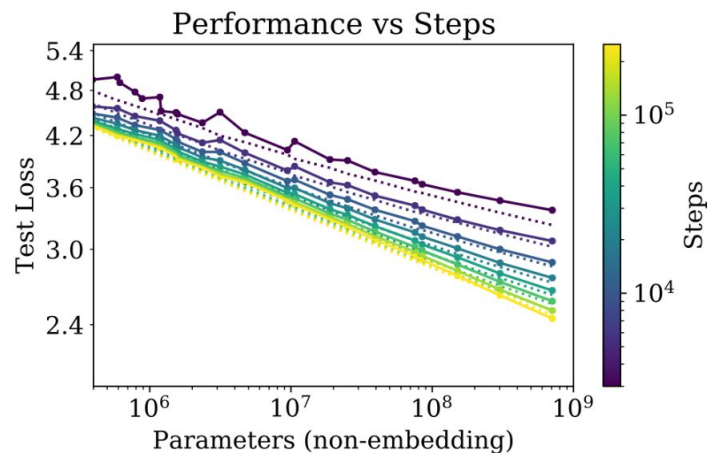
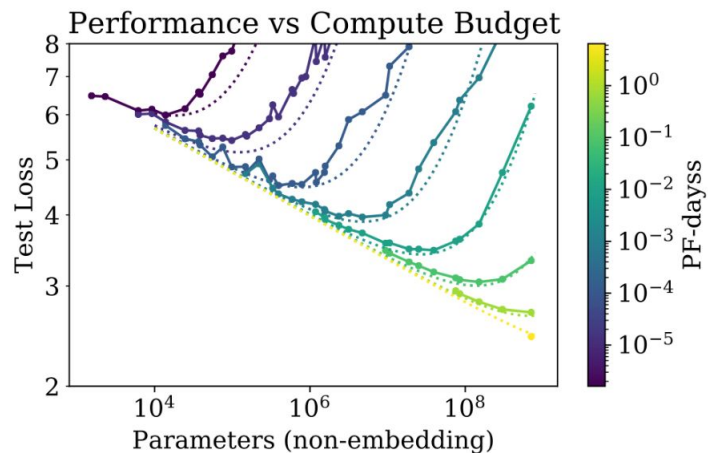


PART 02

Scaling Law for LLM Training

Scaling Law

Why do we care about scaling law?



Optimal Configuration for LLM

Given a fixed FLOPs budget, how should one trade-off model size and the number of training tokens?

FLOPs (Token, Parameter)

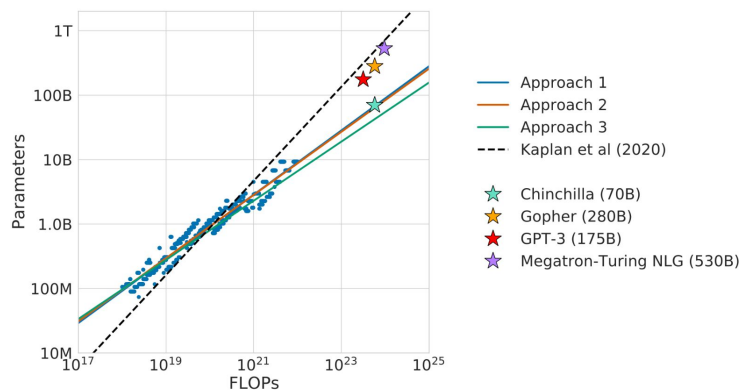


Figure 1 | **Overlaid predictions.** We overlay the predictions from our three different approaches, along with projections from Kaplan et al. (2020). We find that all three methods predict that current large models should be substantially smaller and therefore trained much longer than is currently done. In Figure A3, we show the results with the predicted optimal tokens plotted against the optimal number of parameters for fixed FLOP budgets. Chinchilla outperforms Gopher and the other large models (see Section 4.2).

- Approach 1: Fix model sizes and vary number of training tokens
- Approach 2: IsoFLOP profiles
- Approach 3: Fitting a parametric loss function

Approach 1: Fix model sizes and vary number of training tokens

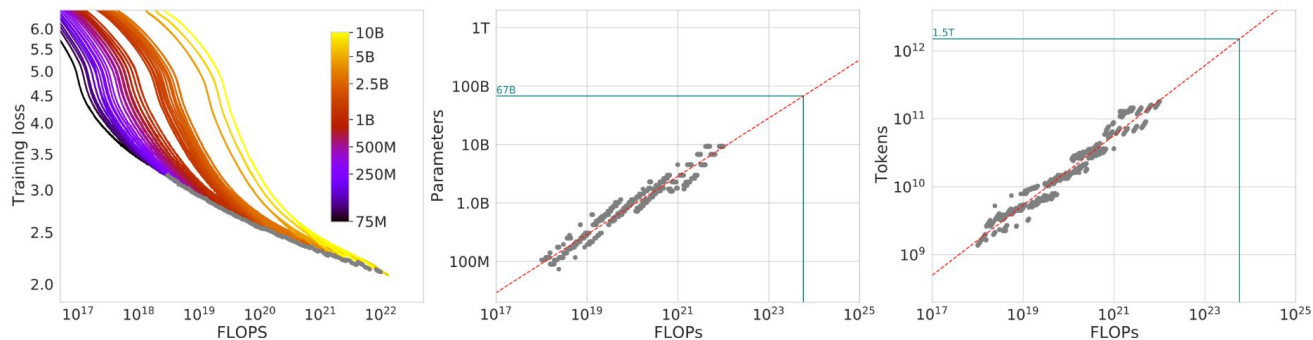


Figure 2 | **Training curve envelope.** On the **left** we show all of our different runs. We launched a range of model sizes going from 70M to 10B, each for four different cosine cycle lengths. From these curves, we extracted the envelope of minimal loss per FLOP, and we used these points to estimate the optimal model size (**center**) for a given compute budget and the optimal number of training tokens (**right**). In green, we show projections of optimal model size and training token count based on the number of FLOPs used to train *Gopher* (5.76×10^{23}).

Approach 2: IsoFLOP profiles

For a given FLOP budget, what is the optimal parameter count?

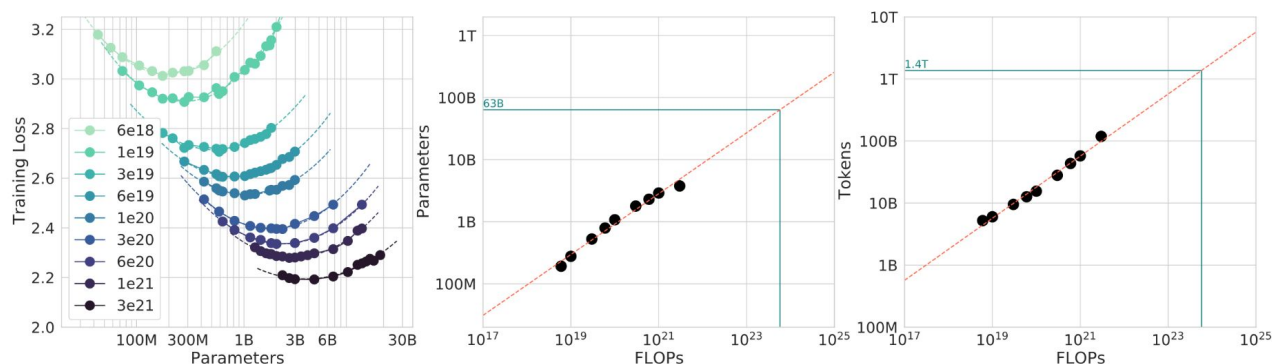
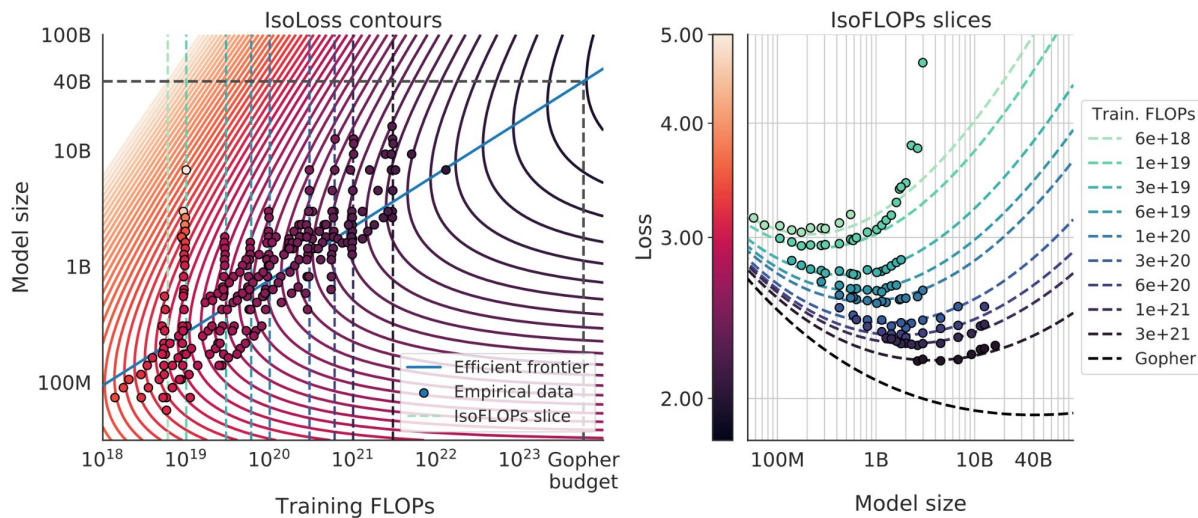


Figure 3 | **IsoFLOP curves.** For various model sizes, we choose the number of training tokens such that the final FLOPs is a constant. The cosine cycle length is set to match the target FLOP count. We find a clear valley in loss, meaning that for a given FLOP budget there is an optimal model to train (**left**). Using the location of these valleys, we project optimal model size and number of tokens for larger models (**center** and **right**). In green, we show the estimated number of parameters and tokens for an *optimal* model trained with the compute budget of *Gopher*.

Approach 3: Fitting a parametric loss function

$$\hat{L}(N, D) \triangleq E + \frac{A}{N^\alpha} + \frac{B}{D^\beta}.$$

N: Model Size
D: Token



FLOPs (Token, Parameter)

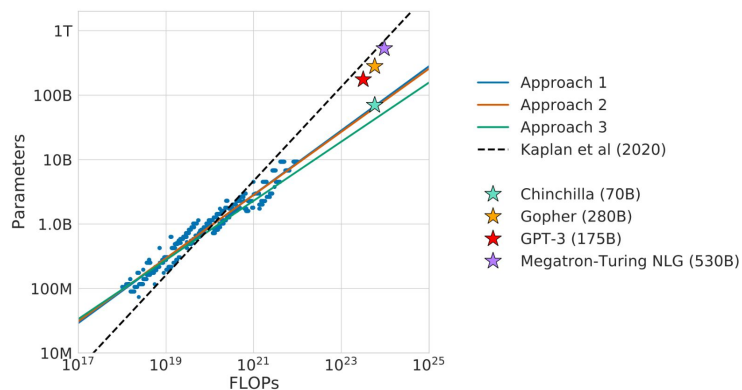


Figure 1 | **Overlaid predictions.** We overlay the predictions from our three different approaches, along with projections from Kaplan et al. (2020). We find that all three methods predict that current large models should be substantially smaller and therefore trained much longer than is currently done. In Figure A3, we show the results with the predicted optimal tokens plotted against the optimal number of parameters for fixed FLOP budgets. Chinchilla outperforms Gopher and the other large models (see Section 4.2).

- Approach 1: Fix model sizes and vary number of training tokens
- Approach 2: IsoFLOP profiles
- Approach 3: Fitting a parametric loss function

FLOPs (Token, Parameter)

Model	Size (# Parameters)	Training Tokens
LaMDA (Thoppilan et al., 2022)	137 Billion	168 Billion
GPT-3 (Brown et al., 2020)	175 Billion	300 Billion
Jurassic (Lieber et al., 2021)	178 Billion	300 Billion
<i>Gopher</i> (Rae et al., 2021)	280 Billion	300 Billion
MT-NLG 530B (Smith et al., 2022)	530 Billion	270 Billion
<i>Chinchilla</i>	70 Billion	1.4 Trillion