

# Scaling laws for large language models

# What & Why?

Given a fixed <sup>C</sup>compute budget, how should one trade-off <sup>N</sup>model size and the <sup>D</sup>number of training tokens?

Accounting for both *training* and *inference*, how does one minimise the cost required to produce and serve a high quality model?

An analogous formulation to Moore's Law for transistor density.

# Kaplan & McCandlish et al., 2020



Performance has a power-law relationship with each of  $N$ ,  $D$ , and  $C$  when not bottlenecked by the other two.

$$L(N, D) = \left[ \left( \frac{N_c}{N} \right)^{\frac{\alpha_N}{\alpha_D}} + \left( \frac{D_c}{D} \right) \right]^{\alpha_D}$$

# Kaplan & McCandlish et al., 2020



An estimate of the total non-embedding training compute  $C$  in terms of  $N$ , and  $D$  is given by,

$$C \approx \underbrace{(2N + 2 \times 2N)}_{\text{Compute per token}} \cdot D = 6ND$$

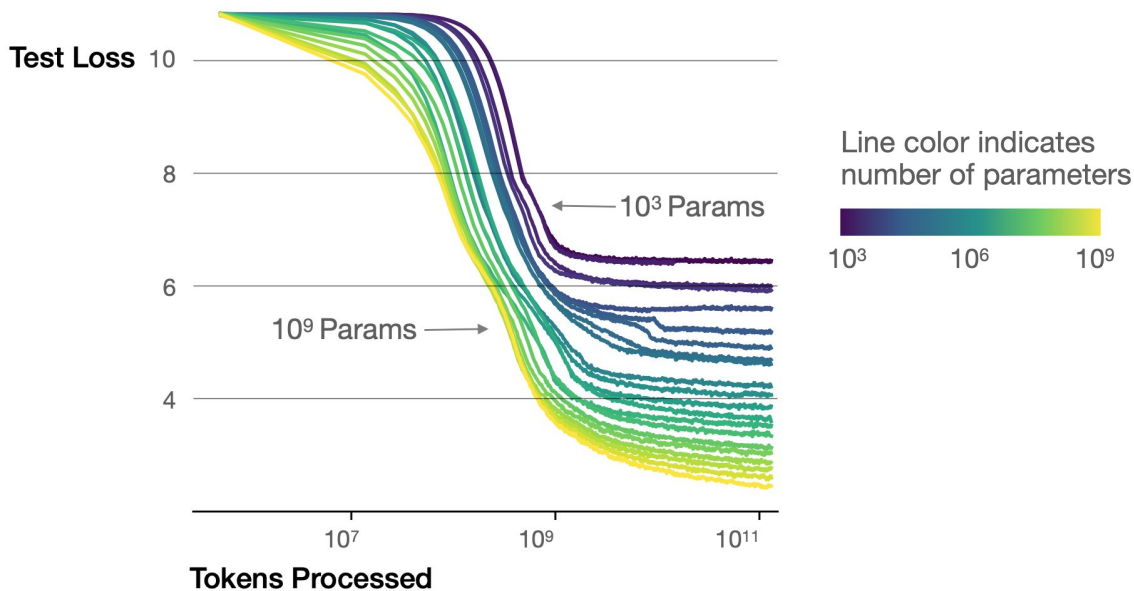
Forward Pass                      Backward Pass <sup>1</sup>

<sup>1</sup> <https://sites.krieger.jhu.edu/jared-kaplan/files/2019/04/ContemporaryMLforPhysicists.pdf>

# Sample Efficiency of Large Models



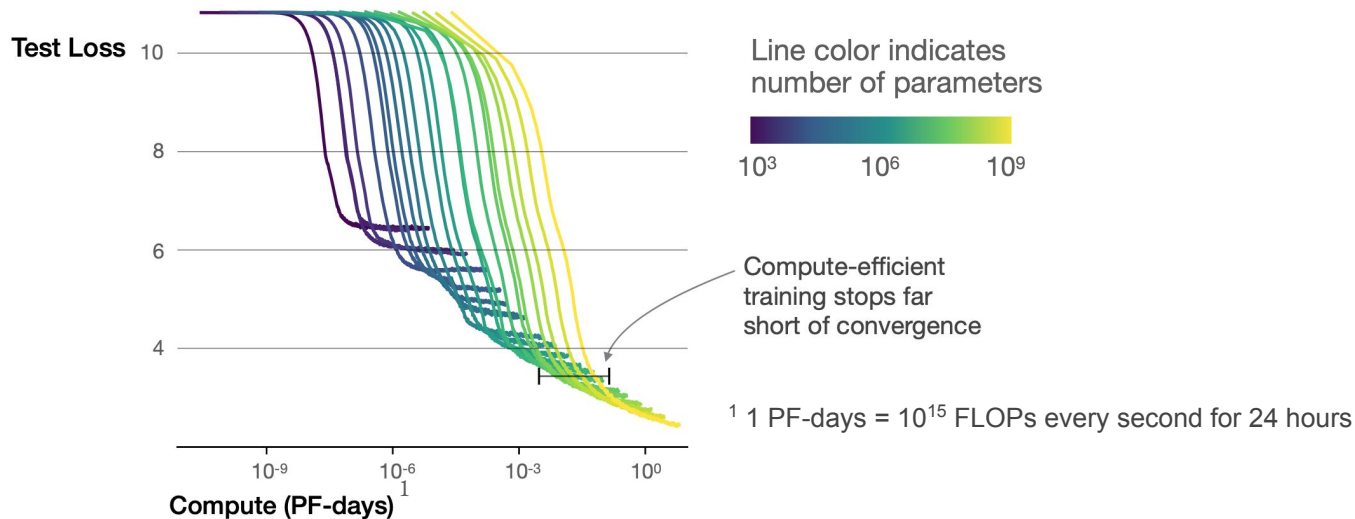
Larger models require fewer samples to reach the same performance.



# Convergence is Inefficient



Optimal performance can be attained by training *very large models* and stopping *significantly short of convergence*.



Kaplan, Jared, et al. "Scaling laws for neural language models." arXiv preprint arXiv:2001.08361 (2020)



# Optimal Allocation of Compute

As compute budget increases, optimal allocation is *larger models*, not more data or longer training.

For example, a 10x increase in compute should be allocated as:

- a 5x increase in model size
  - a 2x increase in data size
- $$\left\{ \begin{array}{l} \sim \text{a } 1.86\text{x increase in batch size} \\ \sim \text{a } 1.07\text{x increase in number of steps} \end{array} \right.$$

# Hoffmann et al., 2022



$$L(N, D) = E + \frac{A}{\textcolor{teal}{N}^\alpha} + \frac{B}{\textcolor{orange}{D}^\beta}$$

Irreducible error      Model error      Data error

The diagram shows the loss function L(N, D) decomposed into three terms. The first term, E, is labeled 'Irreducible error'. The second term, A/N^alpha, is labeled 'Model error', with the variable N highlighted in teal. The third term, B/D^beta, is labeled 'Data error', with the variable D highlighted in orange. Curved lines connect each term to its label below.





# Hoffmann et al., 2022

Optimising the loss function  $L$  under the compute budget  $C = 6ND$ , the compute-optimal model size and dataset size are obtained as,

$$N_{\text{opt}}(C) = G \left( \frac{C}{6} \right)^{\frac{\beta}{\alpha+\beta}} \quad D_{\text{opt}}(C) = G^{-1} \left( \frac{C}{6} \right)^{\frac{\alpha}{\alpha+\beta}}$$
$$\text{where } G = \left( \frac{\alpha A}{\beta B} \right)^{\frac{1}{\alpha+\beta}}$$



# Optimal Allocation of Compute

Unlike Kaplan et al., 2020, Hoffman et al., 2022 recommend that given a 10x increase in compute, the model size and number of training tokens should be scaled in *equal proportions*.

Approach	Coeff. $a$ where $N_{opt} \propto C^a$	Coeff. $b$ where $D_{opt} \propto C^b$
1. Minimum over training curves	0.50 (0.488, 0.502)	0.50 (0.501, 0.512)
2. IsoFLOP profiles	0.49 (0.462, 0.534)	0.51 (0.483, 0.529)
3. Parametric modelling of the loss	0.46 (0.454, 0.455)	0.54 (0.542, 0.543)
<a href="#">Kaplan et al. (2020)</a>	0.73	0.27

# Chinchilla Optimal



"For the compute budget used to train *Gopher*, the optimal model [*Chinchilla*] should be 4 times smaller, while being trained on 4 times more tokens."

MMLU 5-shot accuracy

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

Random	25.0%
Average human rater	34.5%
GPT-3 5-shot	43.9%
<i>Gopher</i> 5-shot	60.0%
<b><i>Chinchilla</i> 5-shot</b>	<b>67.6%</b>
Average human expert performance	89.8%
June 2022 Forecast	57.1%
June 2023 Forecast	63.4%

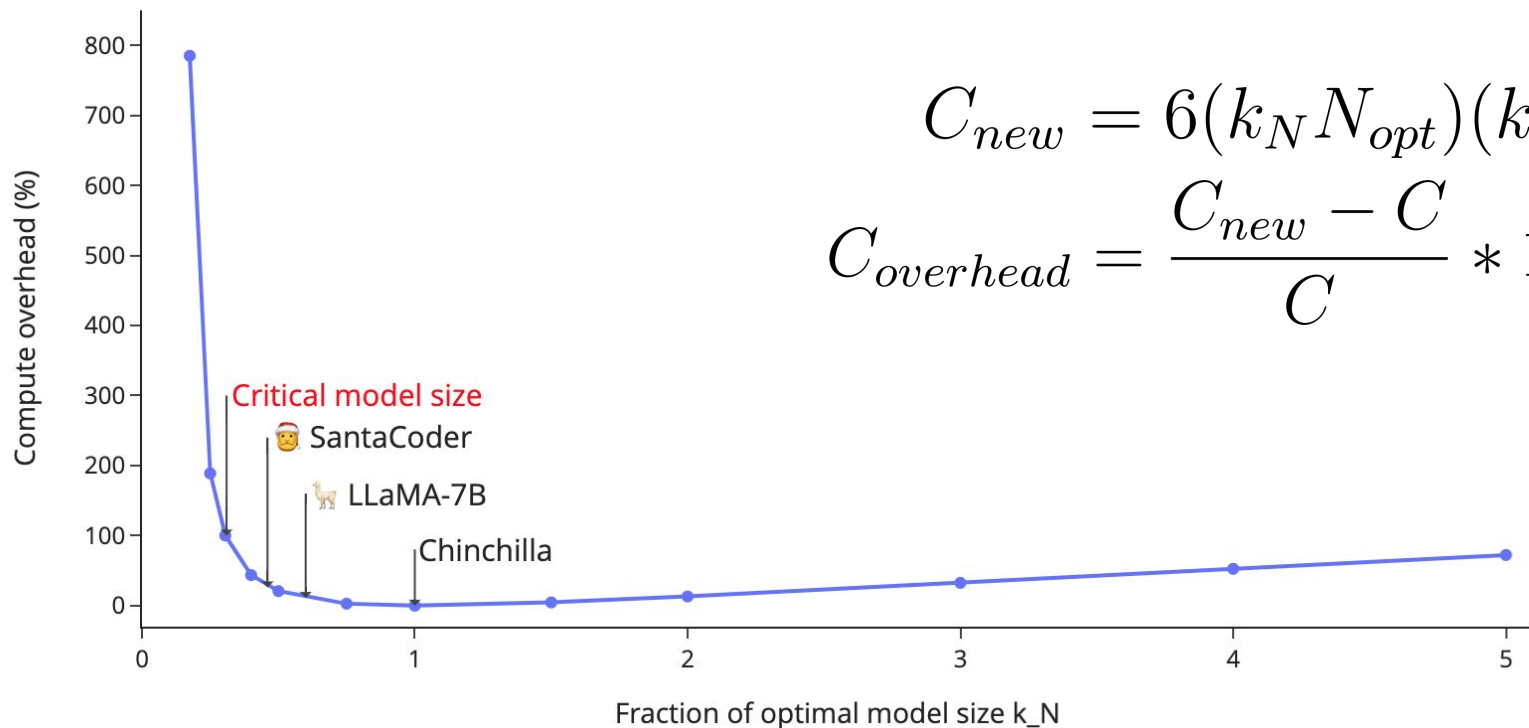
# Compute Overhead

If we reduce the **model size** by  $k_N$ , by how much do we need to increase the **number of tokens** and **compute** to obtain the same loss?

$$k_D = \left( 1 - (k_N^{-\alpha} - 1) \frac{AN_{opt}^{-\alpha}}{BD_{opt}^{-\beta}} \right)^{\frac{1}{-\beta}}$$

$k_D$  turns out to be independent of the **compute budget**!

# Compute Overhead



$$C_{new} = 6(k_N N_{opt})(k_D D_{opt})$$
$$C_{overhead} = \frac{C_{new} - C}{C} * 100.$$

# Training-Inference Compute Trade-off

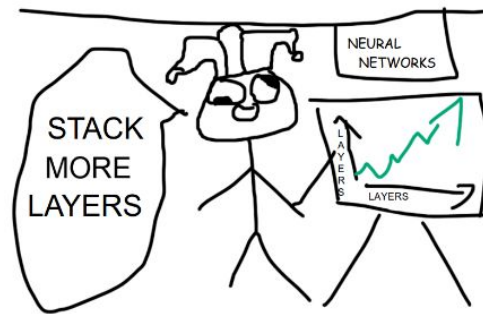
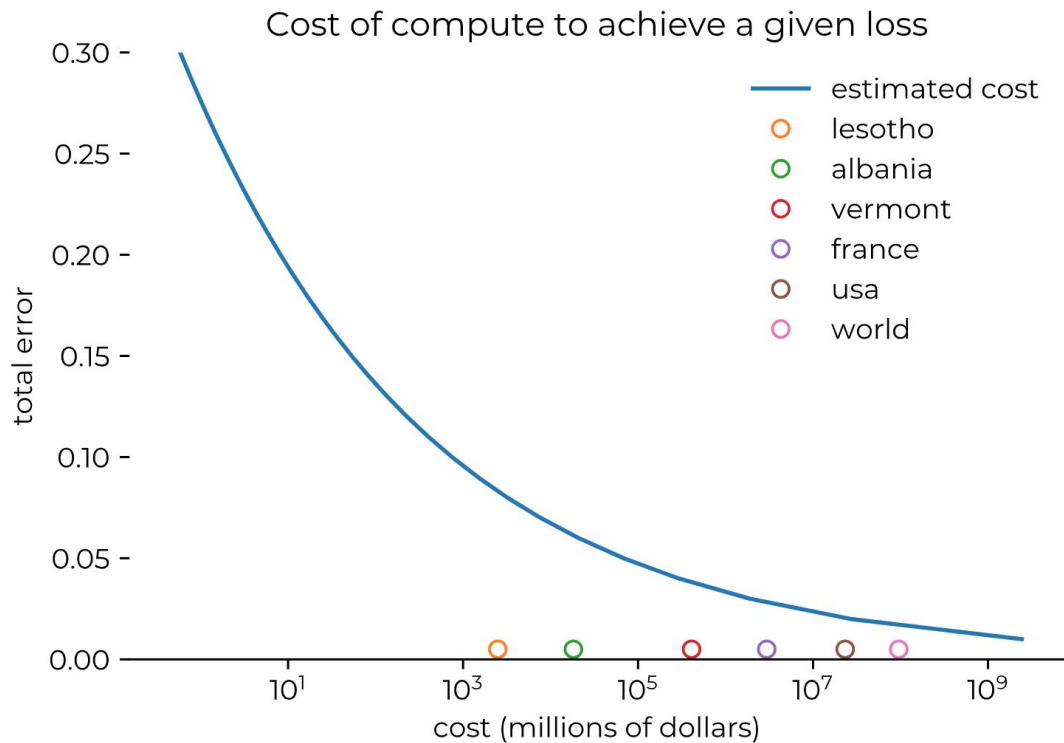
Chinchilla scaling laws only account for the computational cost of *training*.<sup>1</sup>

It may be beneficial to train far beyond Chinchilla optimal to reduce inference costs.



<sup>1</sup> Sardana, Nikhil, and Jonathan Frankle. "Beyond Chinchilla-Optimal: Accounting for Inference in Language Model Scaling Laws." arXiv preprint arXiv:2401.00448 (2023).

# Do we have enough compute?



# Do we have enough tokens?

We are projected to run out of language data between 2030 & 2050.<sup>1</sup>

Model	Stock of data (#words)	Growth rate
Recorded speech	1.46e17	5.2%
	[3.41e16; 4.28e17]	[4.95%; 5.2%]
Internet users	2.01e15	8.14%
	[6.47e14; 6.28e15]	[7.89%; 8.14%]
Popular platforms	4.41e14	8.14%
	[1.21e14; 1.46e15]	[7.89%; 8.14%]
CommonCrawl	9.62e13	16.68%
	[4.45e13; 2.84e14]	[16.41%; 16.68%]
Indexed websites	2.21e14	NA
	[5.16e13; 6.53e15]	
<b>Aggregated model</b>	<b>7.41e14</b>	<b>7.15%</b>
	<b>[6.85e13; 7.13e16]</b>	<b>[6.41%; 17.49%]</b>

<sup>1</sup> Villalobos et al. "Will we run out of data? An analysis of the limits of scaling datasets in Machine Learning." arXiv preprint arXiv:2211.04325 (2022)



# Future

Quality of data

Multi-modal training

Translation of loss into real-world performance

Sparsely activated gated Mixture-of-Expert models