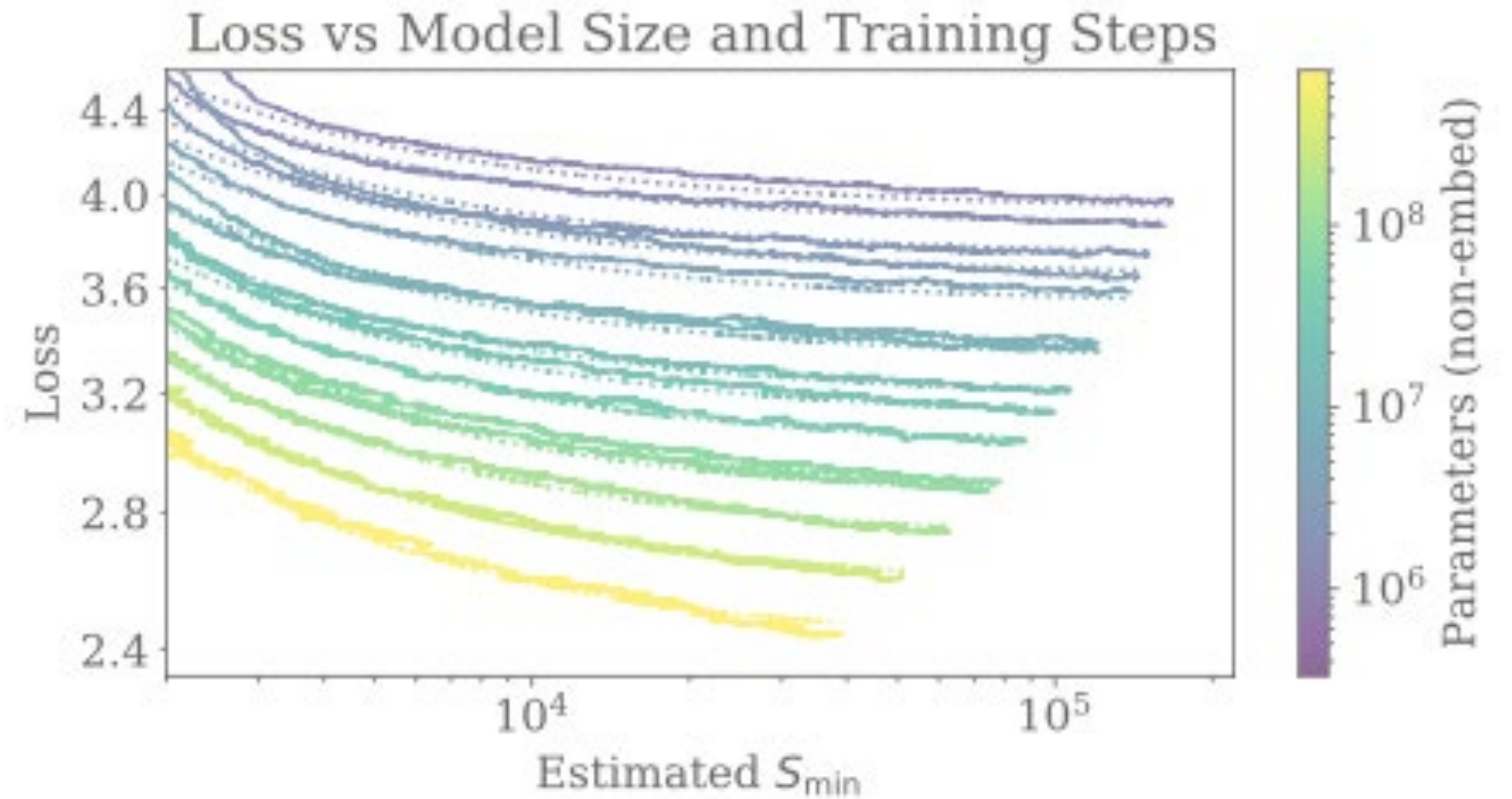
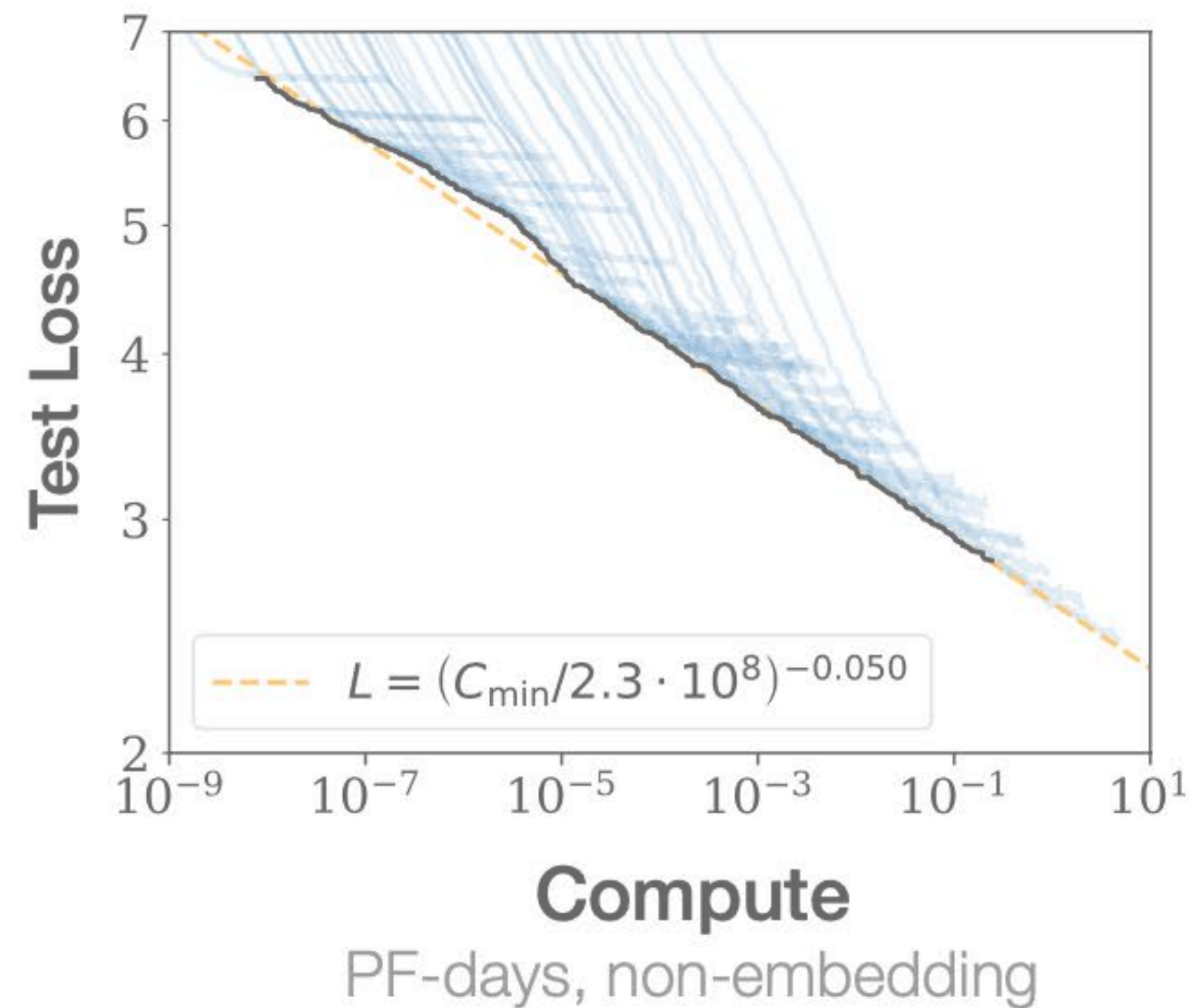


Lecture 20: Scaling Laws

Speaker: Phillip Isola



AlexNet to AlphaGo Zero: A 300,000x Increase in Compute (Log Scale)

Petaflop/s-days



The Bitter Lesson

© Rich Sutton. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <https://ocw.mit.edu/help/faq-fair-use/>

Rich Sutton

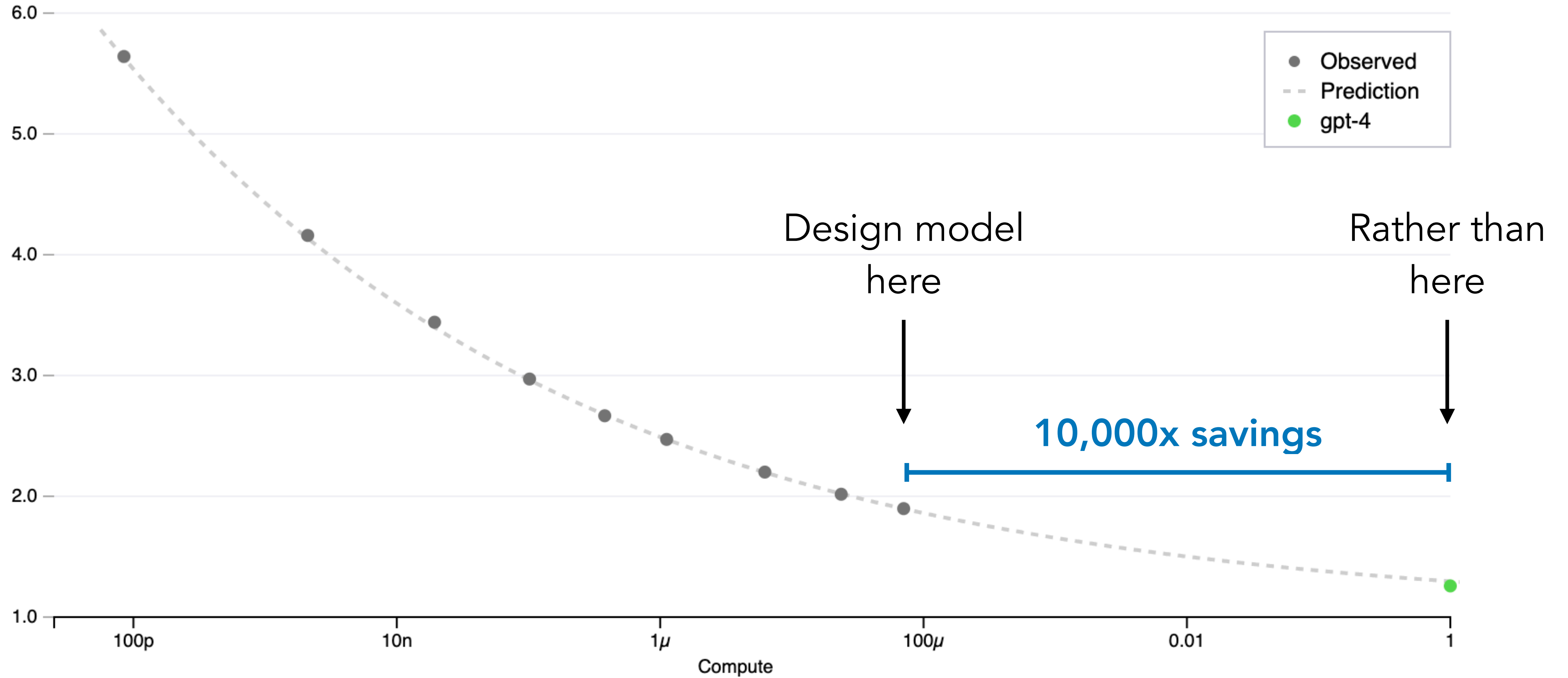
March 13, 2019

The biggest lesson that can be read from 70 years of AI research is that general methods that leverage computation are ultimately the most effective, and by a large margin. The ultimate reason for this is Moore's law, or rather its generalization of continued exponentially falling cost per unit of computation. Most AI research has been conducted as if the computation available to the agent were constant (in which case leveraging human knowledge would be one of the only ways to improve performance) but, over a slightly longer time than a typical research project, massively more computation inevitably becomes available. Seeking an improvement that makes a difference in the shorter term, researchers seek to leverage their human knowledge of the domain, but the only thing that matters in the long run is the leveraging of computation. These two need not run counter to each other, but in practice they tend to. Time spent on one is time not spent on the other. There are psychological commitments to investment in one approach or the other. And the human-knowledge approach tends to complicate methods in ways that make them less suited to taking advantage of general methods leveraging computation. There were many examples of AI researchers' belated learning of this bitter lesson, and it is instructive to review some of the most prominent.

If you had a given budget of compute, what model would you train on how much data?

OpenAI codebase next word prediction

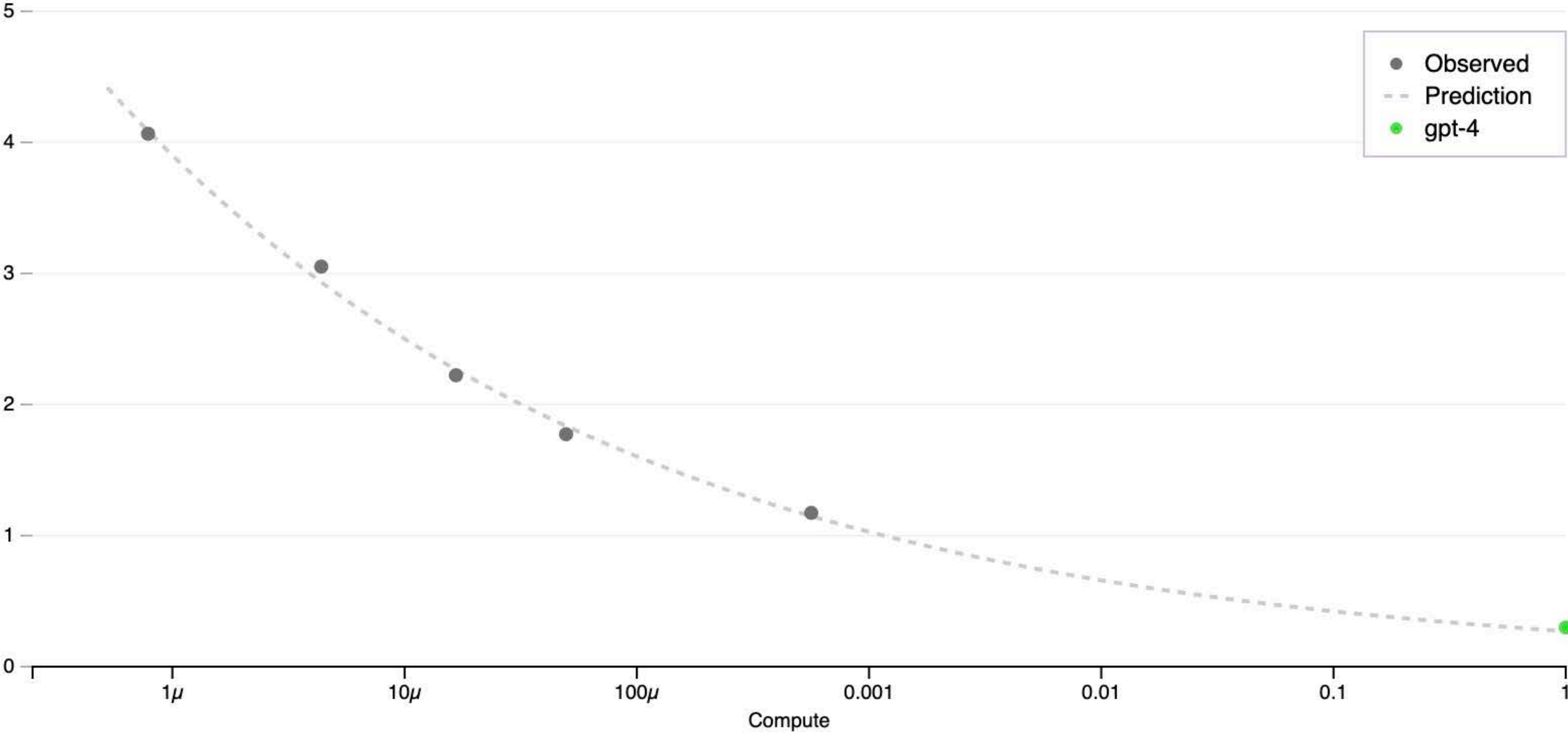
Bits per word



<https://openai.com/research/gpt-4>

Capability prediction on 23 coding problems

– Mean Log Pass Rate



<https://openai.com/research/gpt-4>

Scaling laws

minimize *test-loss(model size, data, batch size, steps,...)* **such that** *compute* **within budget**

for a given architecture, how much data / what model size would be needed for a target performance?

can we extrapolate from smaller to larger models / experiments??

$$N_{opt}(C), D_{opt}(C) = \underset{N,D \text{ s.t. } \text{FLOPs}(N,D)=C}{\operatorname{argmin}} L(N, D)$$

How does the test loss scale as a function of data, #parameters etc?

- autoregressive Transformer
- data: WebText2, 96GB of text, $2.29 \cdot 10^{10}$ tokens
1024 token context
- vary:
 - model size: 768 - 1.5 billion non-embedding parameters
 - data: 22M - 23B tokens
 - shape (depth, width, attention heads, ...)
 - batch size

- compute = 6 x #parameters x batchsize x #steps
- computation: in PF-days (petaflop/s-days)
1 PF-day = $8.64 \cdot 10^{19}$ FLOPs

Scaling Laws for Neural Language Models

Jared Kaplan *

Johns Hopkins University, OpenAI

jaredk@jhu.edu

Sam McCandlish*

OpenAI

sam@openai.com

Tom Henighan

OpenAI

henighan@openai.com

Tom B. Brown

OpenAI

tom@openai.com

Benjamin Chess

OpenAI

bchess@openai.com

Rewon Child

OpenAI

rewon@openai.com

Scott Gray

OpenAI

scott@openai.com

Alec Radford

OpenAI

alec@openai.com

Jeffrey Wu

OpenAI

jeffwu@openai.com


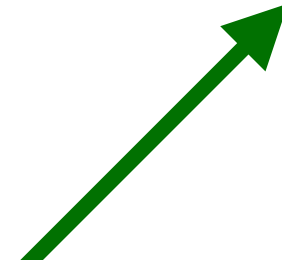
Dario Amodei

OpenAI

damodei@openai.com

Finding #1: power law relationships

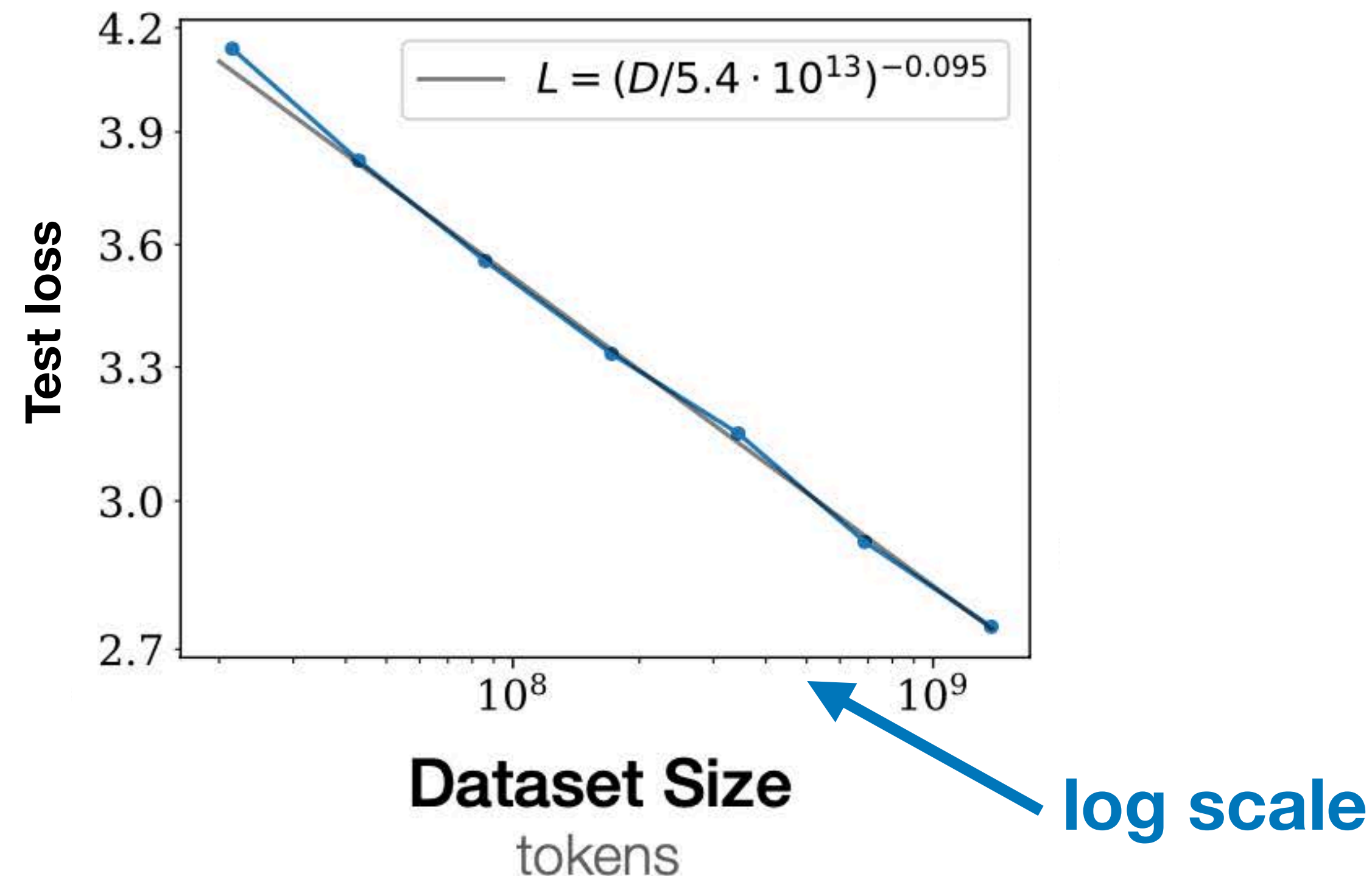
$$L(X) = (1/X)^\alpha = X^{-\alpha}$$

test loss  **resource** 
(data/parameters/
compute

- scale $X \rightarrow 2X$ \Rightarrow Loss $\rightarrow 2^{-\alpha}$ Loss

How does the test loss scale as a function of data?

- #parameters, compute: “infinite”

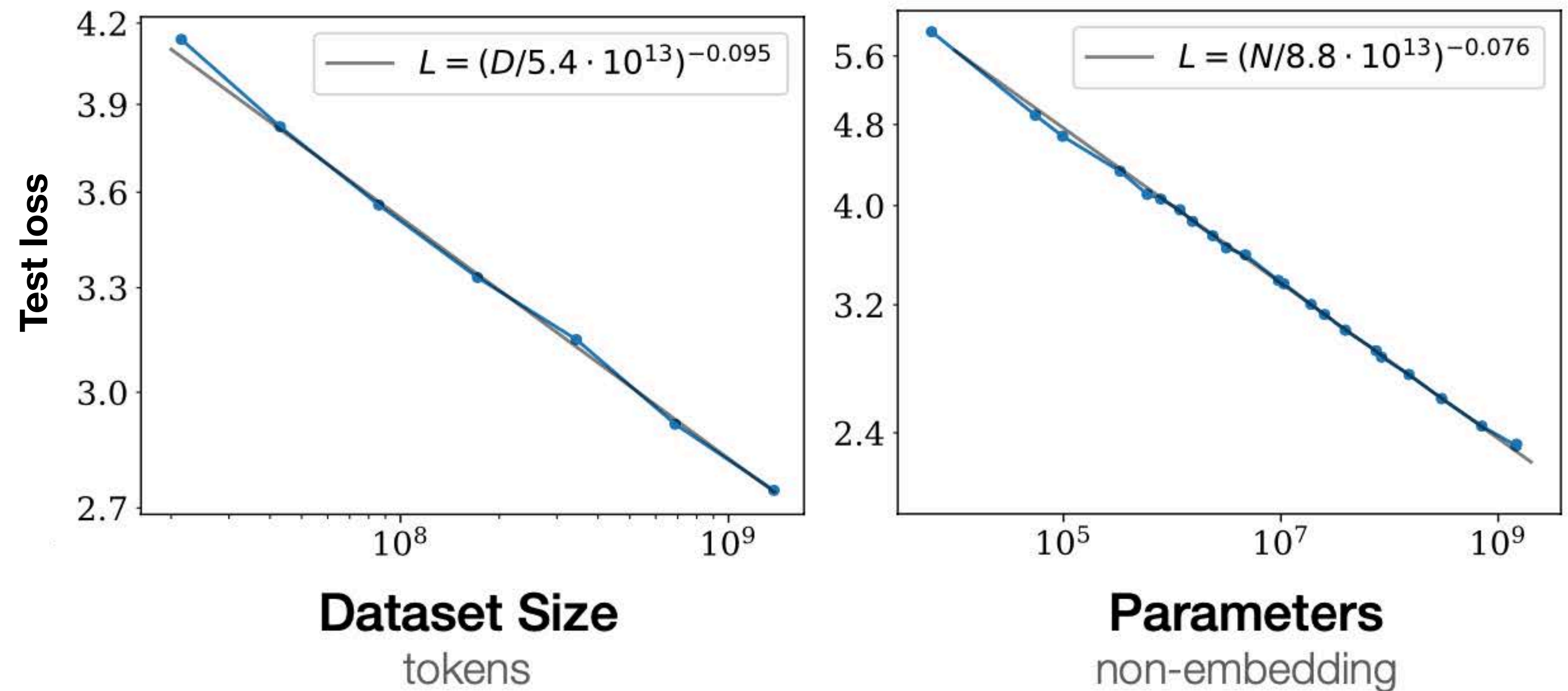


- power law relationship:

$$L(D) = (D_c/D)^{\alpha_D} ; \quad \alpha_D \sim 0.095, \quad D_c \sim 5.4 \times 10^{13} \text{ (tokens)}$$

How does the test loss scale as a function of #parameters?

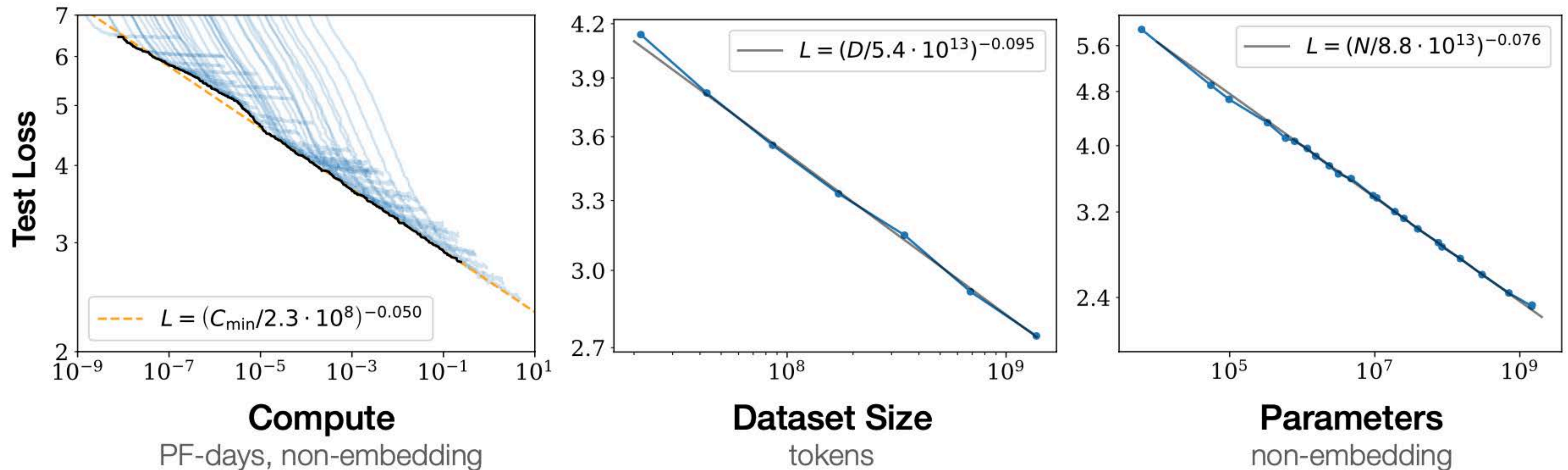
- #tokens, compute: “infinite”



$$L(N) = (N_c/N)^{\alpha_N} ; \quad \alpha_N \sim 0.076, \quad N_c \sim 8.8 \times 10^{13} \text{ (non-embedding parameters)}$$

How does the test loss scale with available compute?

- power law for all of them



$$\underline{L(C_{\min}) = (C_c^{\min}/C_{\min})^{\alpha_C^{\min}}}; \quad \alpha_C^{\min} \sim 0.050, \quad C_c^{\min} \sim 3.1 \times 10^8 \text{ (PF-days)}$$

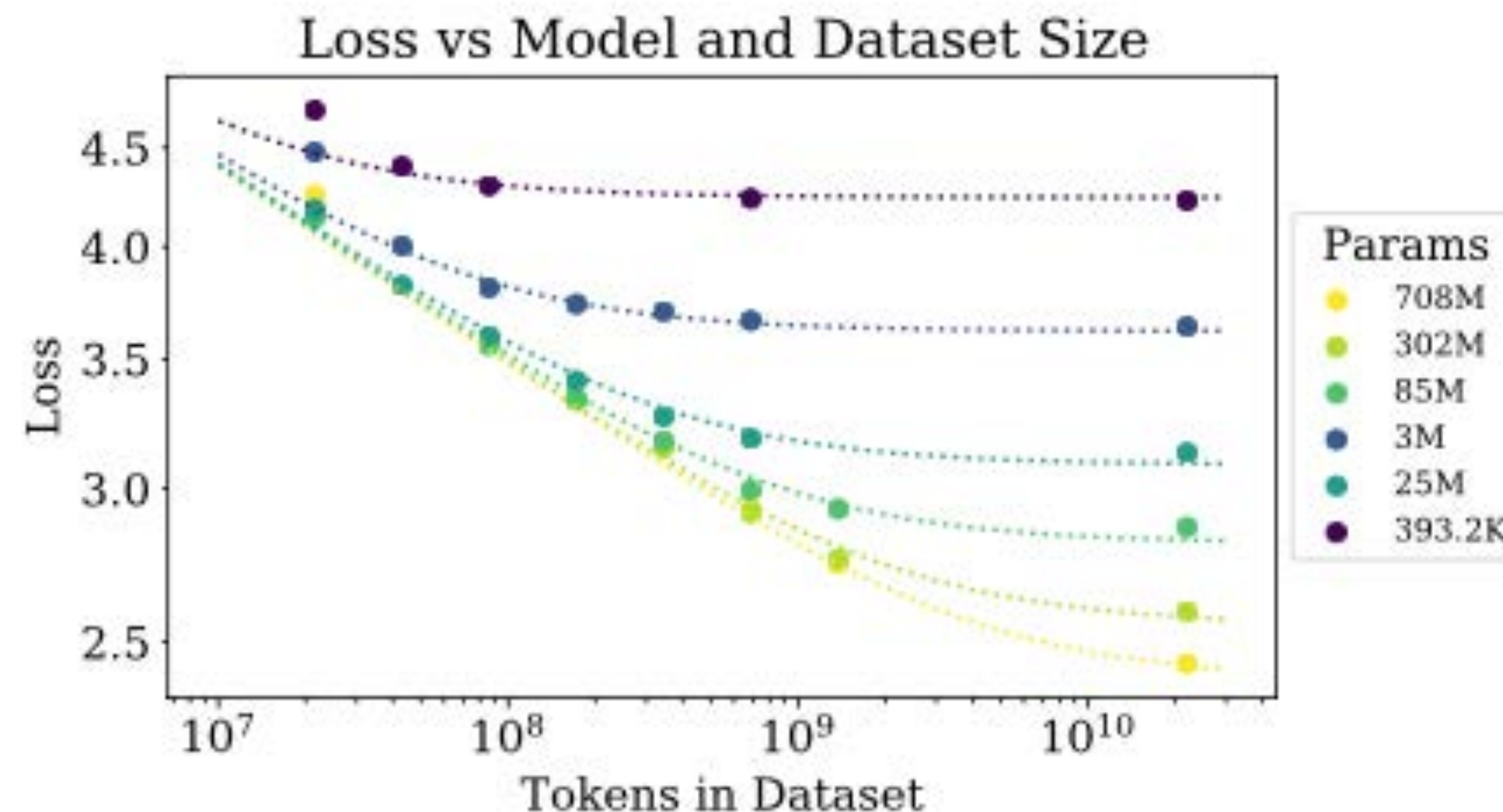
Finding #2: simultaneous dependence on multiple factors

- “Performance improves predictably as long as we scale up N and D in tandem, but enters a regime of diminishing returns if either N or D is held fixed while the other increases.”

- Loss as function of #parameters N & data D :

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

constants



Parameter	α_N	α_D	N_c	D_c
Value	0.076	0.103	6.4×10^{13}	1.8×10^{13}

Finding #2: simultaneous dependence on multiple factors

- “Performance improves predictably as long as we scale up N and D in tandem, but enters a regime of diminishing returns if either N or D is held fixed while the other increases.”

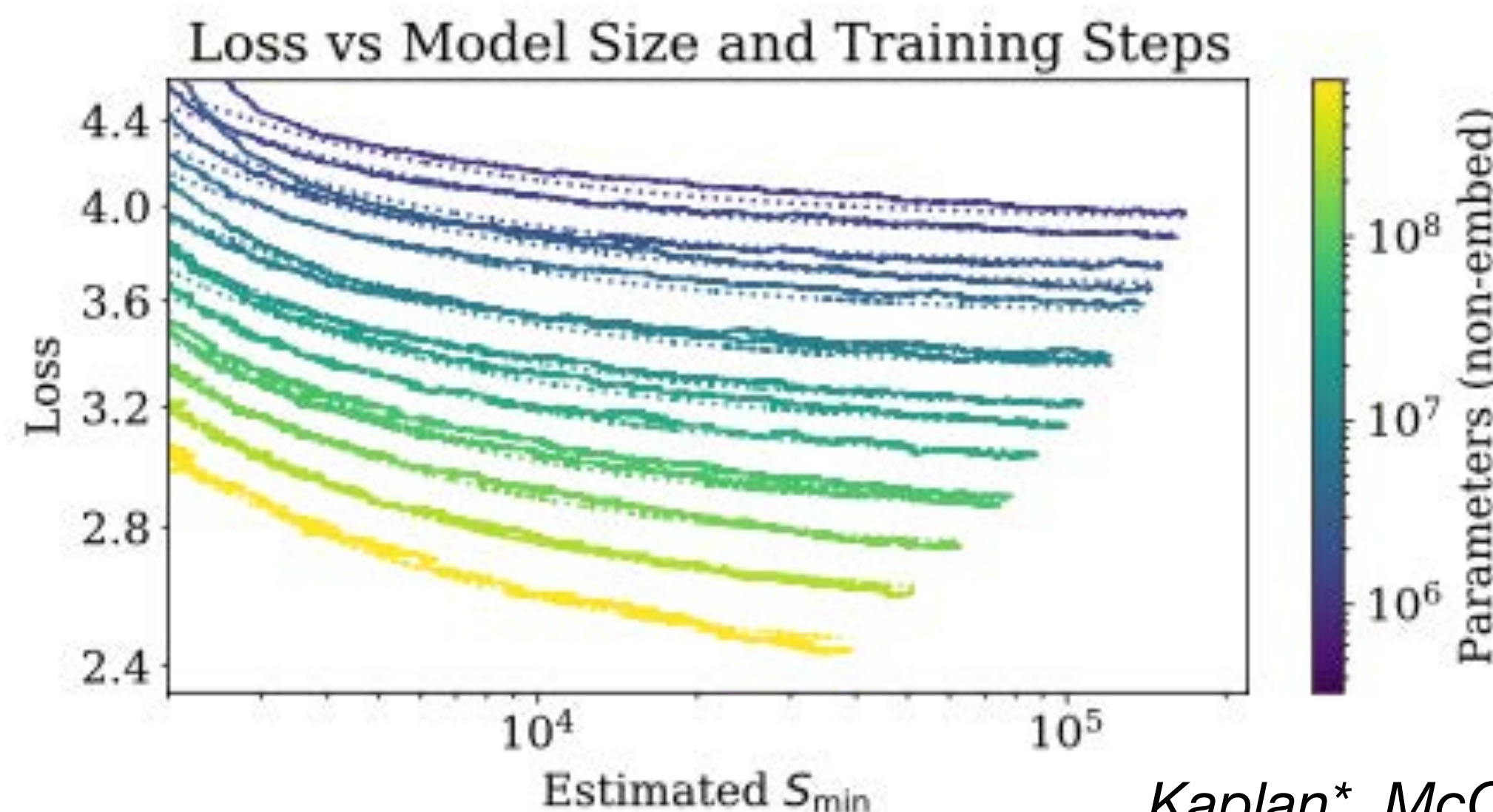
- Loss as function of #parameters N & data D :

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

- Loss as function of #parameters N & #steps S :

$$L(N, S_{\min}) = \left(\frac{N_c}{N} \right)^{\alpha_N} + \left(\frac{S_c}{S_{\min}} \right)^{\alpha_S}$$

constants



Finding #2: simultaneous dependence on multiple factors

- “Performance improves predictably as long as we scale up N and D in tandem, but enters a regime of diminishing returns if either N or D is held fixed while the other increases.”

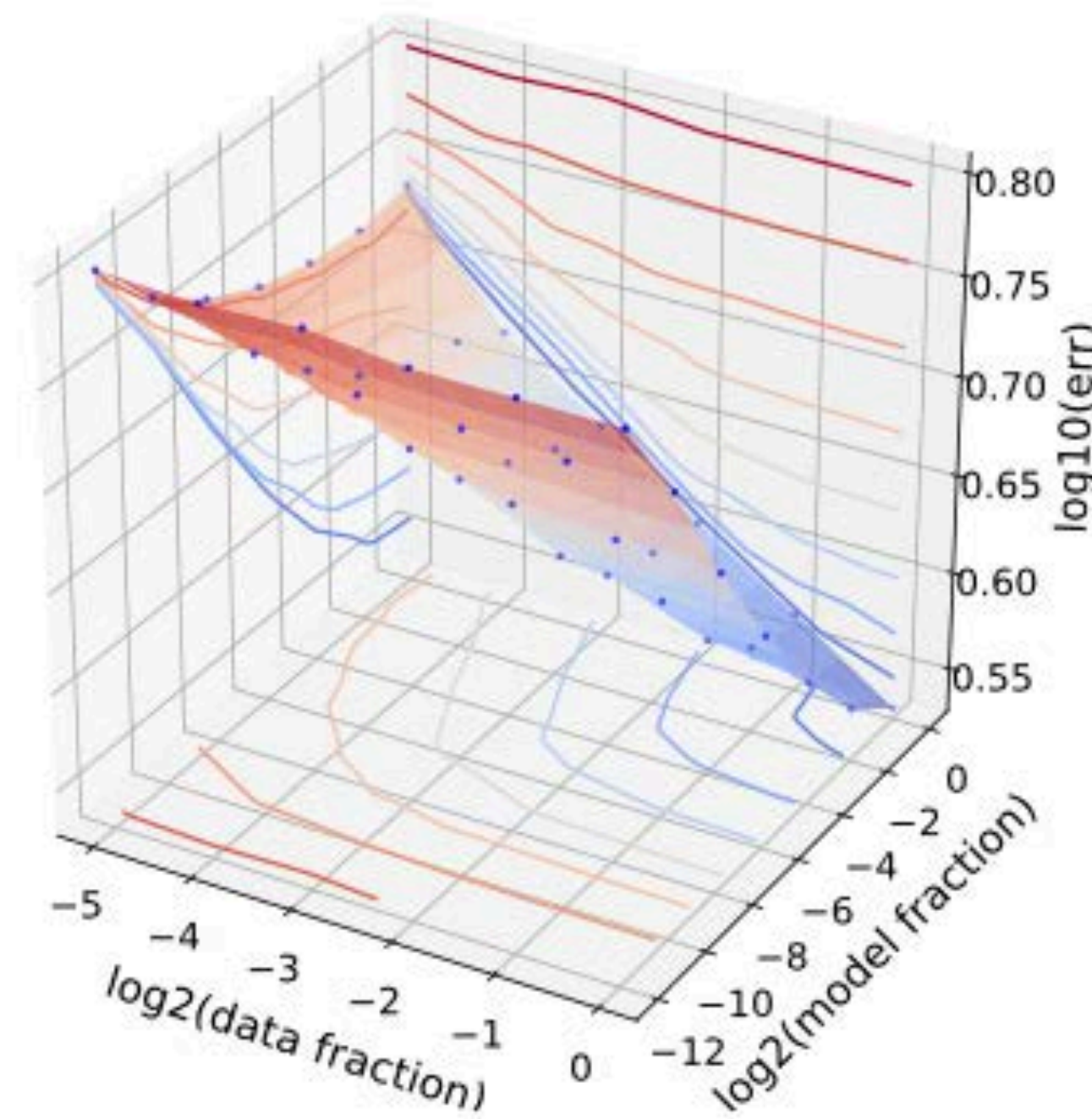
- Loss as function of #parameters N & data D :
$$L(N, D) = \left[\left(\frac{N_c}{N} \right)^{\frac{\alpha_N}{\alpha_D}} + \frac{D_c}{D} \right]^{\alpha_D}$$
- Loss as function of #parameters N & #steps S :

$$L(N, S_{\min}) = \left(\frac{N_c}{N} \right)^{\alpha_N} + \left(\frac{S_c}{S_{\min}} \right)^{\alpha_S}$$

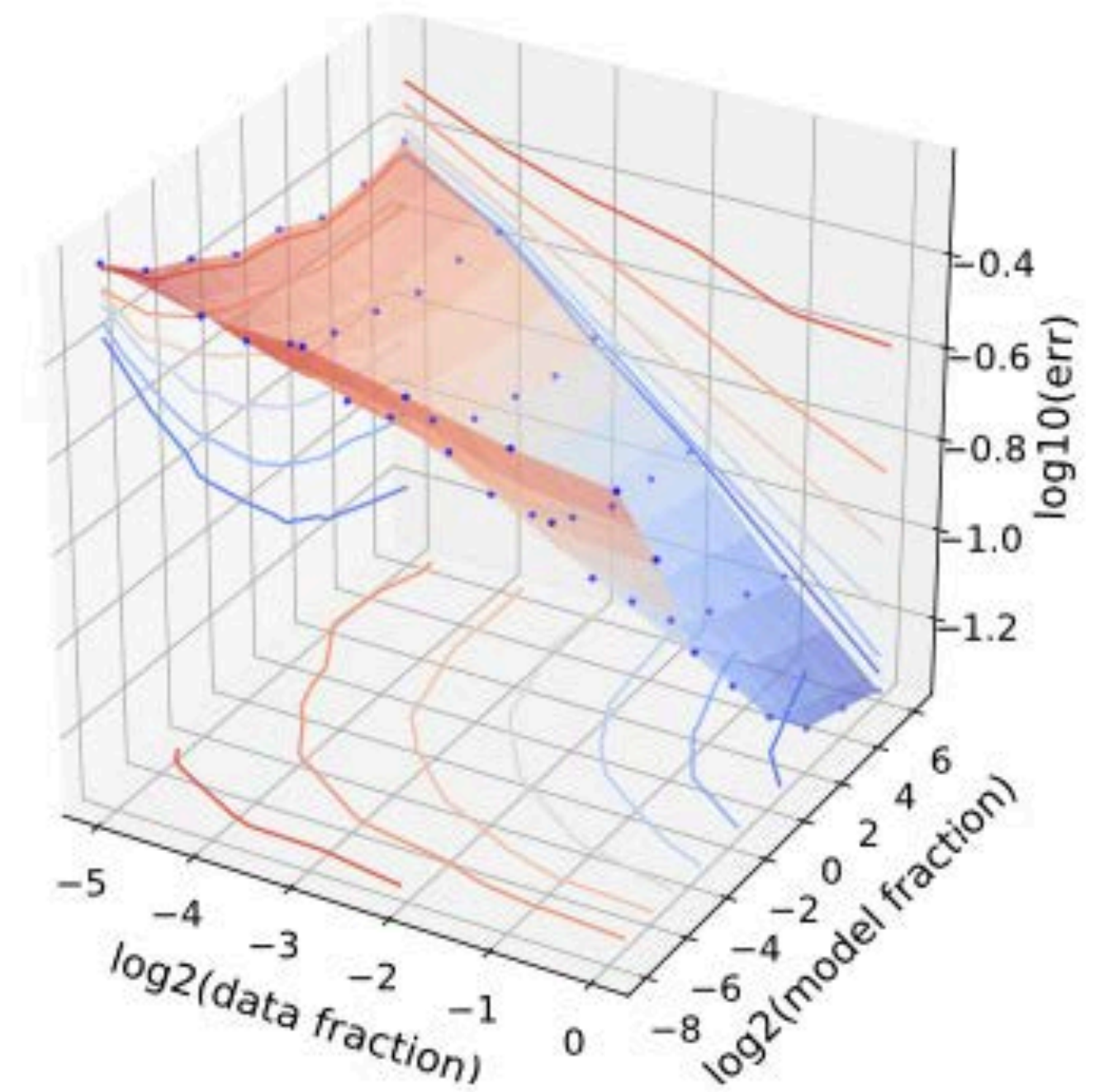
- from this, they conclude to scale, with compute C :
 $\text{\#params} \propto C^{0.73}, \text{ data} = \text{batch-size} \cdot \text{\#steps} \propto C^{0.27}$

Take these exact equations/numbers with a grain of salt!

Similar findings in other works



(a) Wiki103 error (cross entropy) landscape.



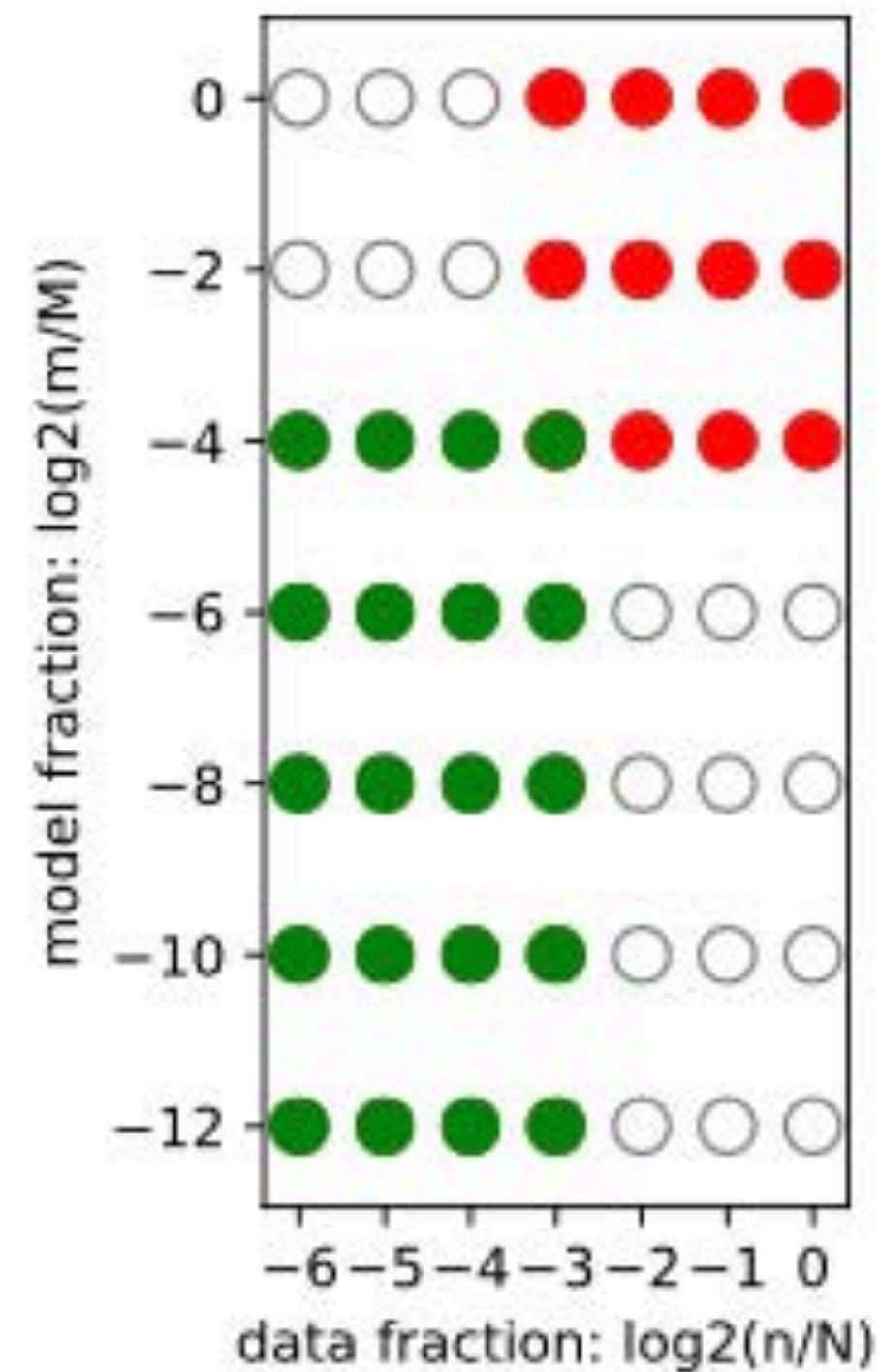
(b) CIFAR10 error (top1) landscape.

- *Rosenfeld et al:*
for vision, language, fit function

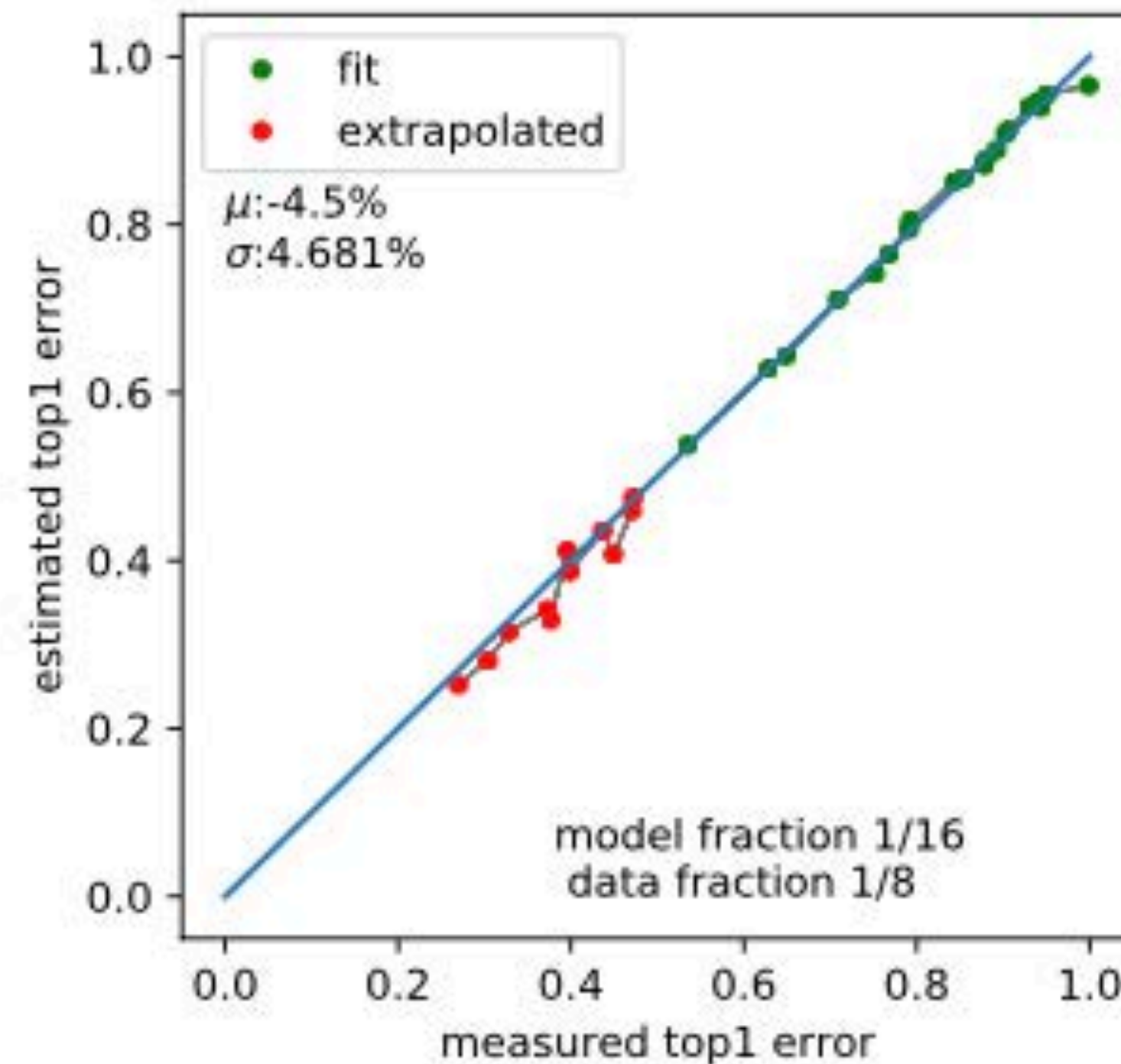
$$L(N, D) \approx \frac{a(N)}{D^{\alpha_1(N)}} + \frac{b(N)}{N^{-\alpha_2(D)}} + c_\infty$$

Extrapolation is possible!

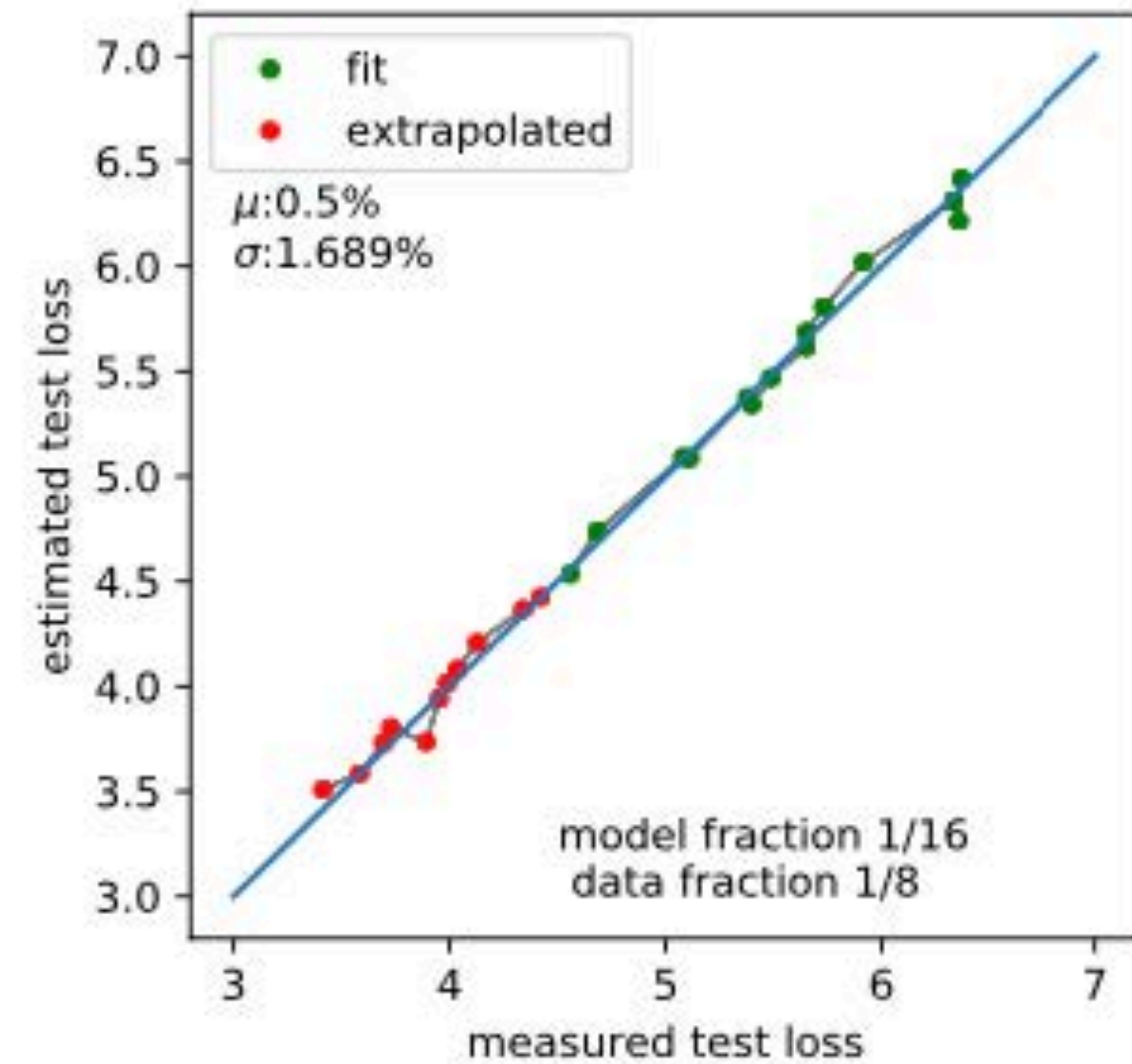
- from *Rosenfeld et al*



(a) Illustration.



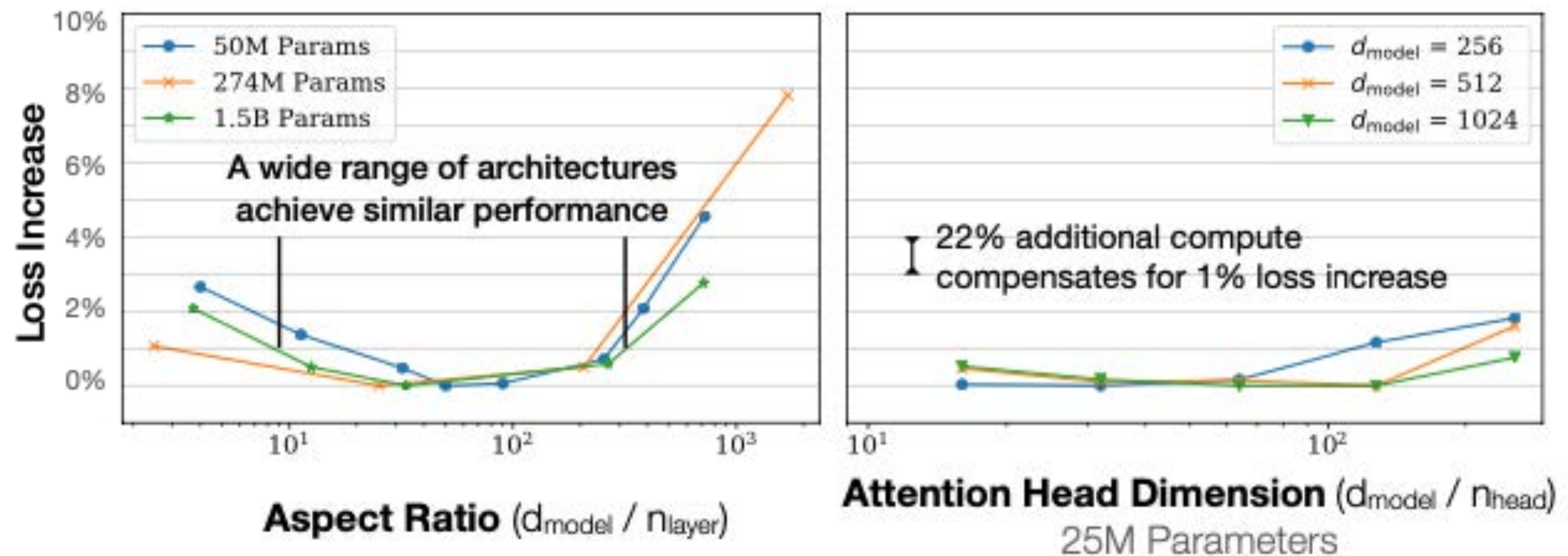
(b) Extrapolation on ImageNet



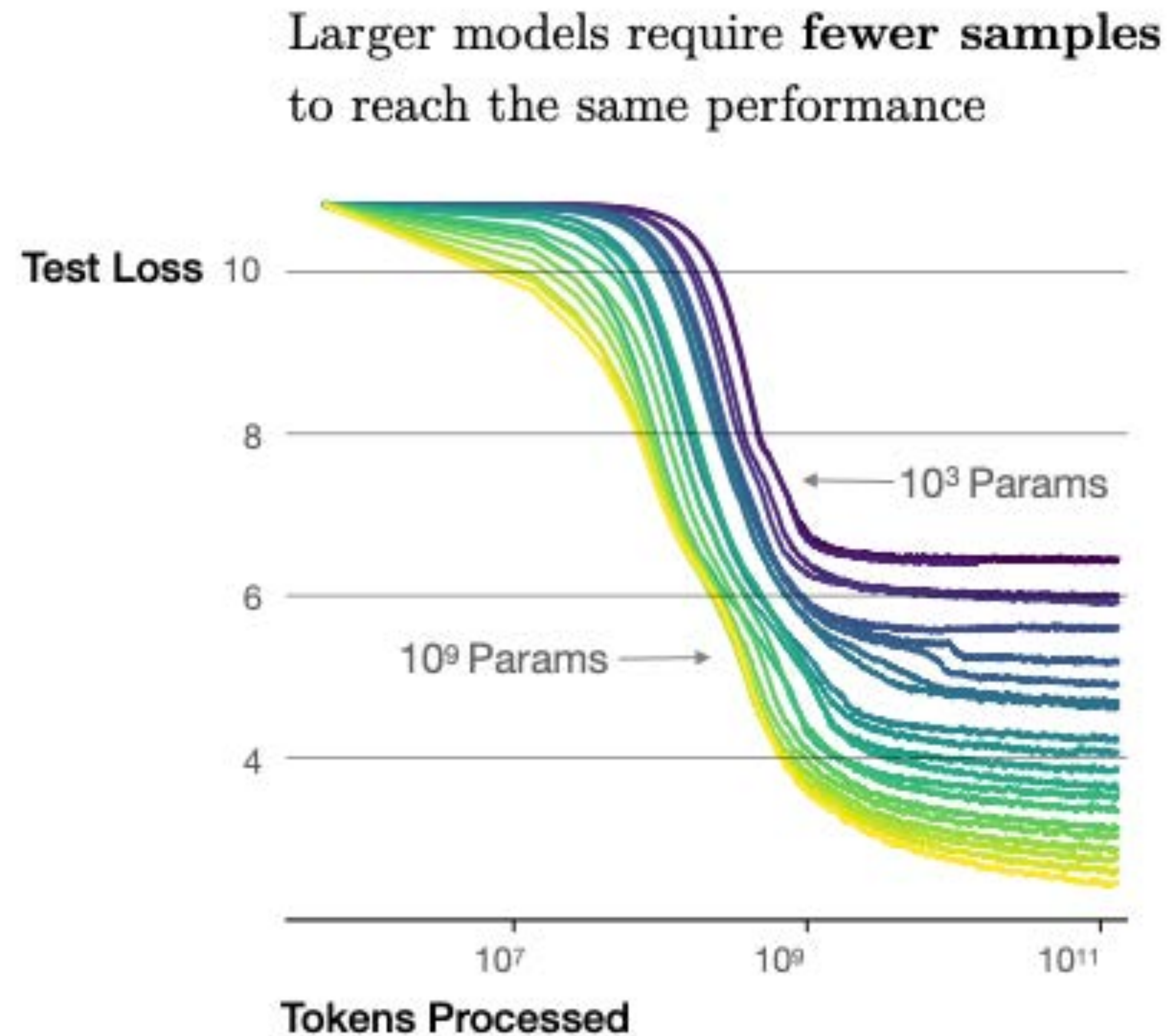
(c) Extrapolation on WikiText-103.

Finding #3: dependence on model shape

- “Performance depends strongly on scale, weakly on model shape”

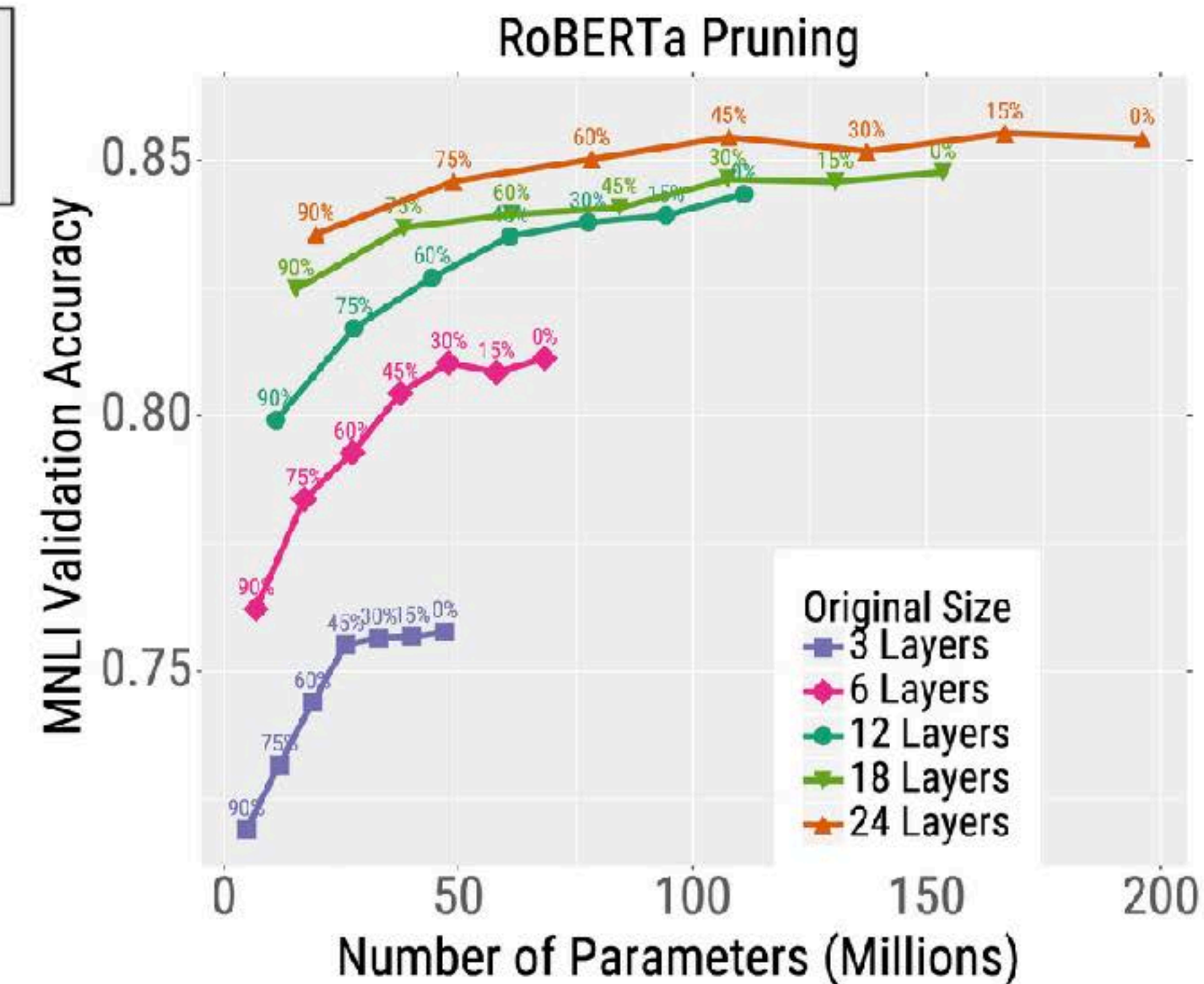
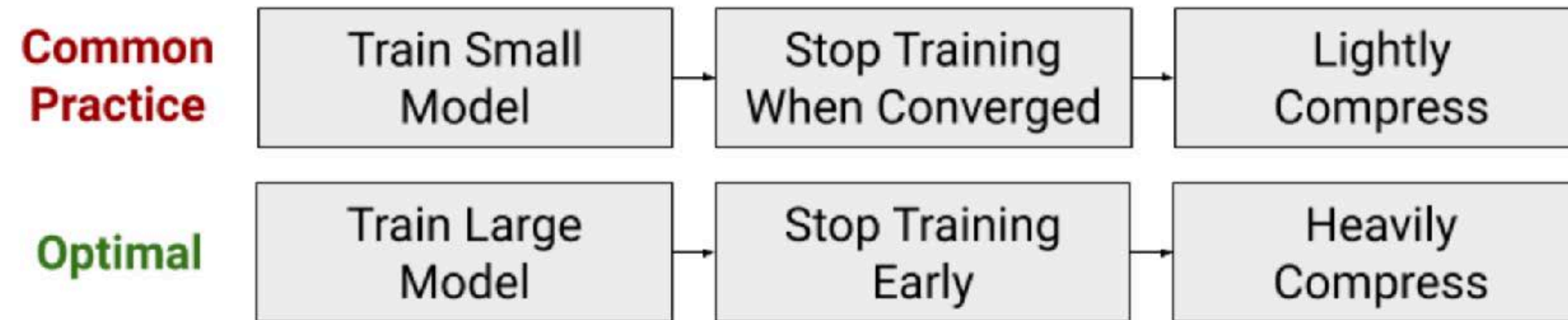


Finding #4: larger models require fewer samples to reach the same performance



Similar findings...

- *Li et al, 2020*: larger models train faster, not much loss via compression



Take-aways from *Kaplan et al.*

- test loss scales as **power law** as a function of data, parameter count, compute
- test loss is more sensitive to parameter count than data: with more compute, **increase model by a larger factor than dataset size**
- don't train to completion, rather make your model larger
- in general, scaling laws can help to “optimally” allocate computation

“Caveats” for *Kaplan et al*

- “... the scaling with D at very large model size still remains mysterious. Without a theory or systematic understanding of the corrections to our scaling laws, it’s difficult to determine in what circumstances they can be trusted”
- didn’t thoroughly investigate small data regime, poor fits of $L(N,D)$ in that regime. Didn’t investigate regularization or data augmentation.
- didn’t tune all hyperparameters; possibly better learning rates for short training runs

Implications

- subsequently: larger models, scale data only moderately

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



spend compute on model size ...



Pocket gopher courtesy of LeonardoWeiss on Wikipedia Commons. Used under CC BY.

Dinosaur © Universal Pictures. Chinchilla © source unknown. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <https://ocw.mit.edu/help/faq-fair-use/>

Finding the right predictions is not easy...

... or on data?



Training Compute-Optimal Large Language Models

Jordan Hoffmann*, Sebastian Borgeaud*, Arthur Mensch*, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals and Laurent Sifre*

*Equal contributions

Two years later...

under a given compute budget. We find that current large language models are significantly under-trained, a consequence of the recent focus on scaling language models whilst keeping the amount of training data constant. **By training over 400 language models ranging from 70 million to over 16 billion parameters on 5 to 500 billion tokens, we find that for compute-optimal training, the model size and the number of training tokens should be scaled equally: for every doubling of model size the number of training tokens should also be doubled.** We test this hypothesis by training a predicted compute-optimal model, *Chinchilla*, that uses the same compute budget as *Gopher* but with 70B parameters and 4× more data. *Chinchilla* uniformly and significantly outperforms *Gopher* (280B), GPT-3 (175B), Jurassic-1 (178B), and Megatron-Turing NLG (530B) on a large range of downstream evaluation tasks. This also means that *Chinchilla* uses substantially less compute for fine-tuning and inference, greatly

- still assume a power law, but find somewhat different equations / scaling
- still do not train to minimum loss, but train for many more tokens (& steps), increase #parameters and data equally: “smaller” models, train longer

Differences

- different learning rate schedules (not explored in Kaplan et al), adapted to dataset size
- larger models (> 500M parameters)

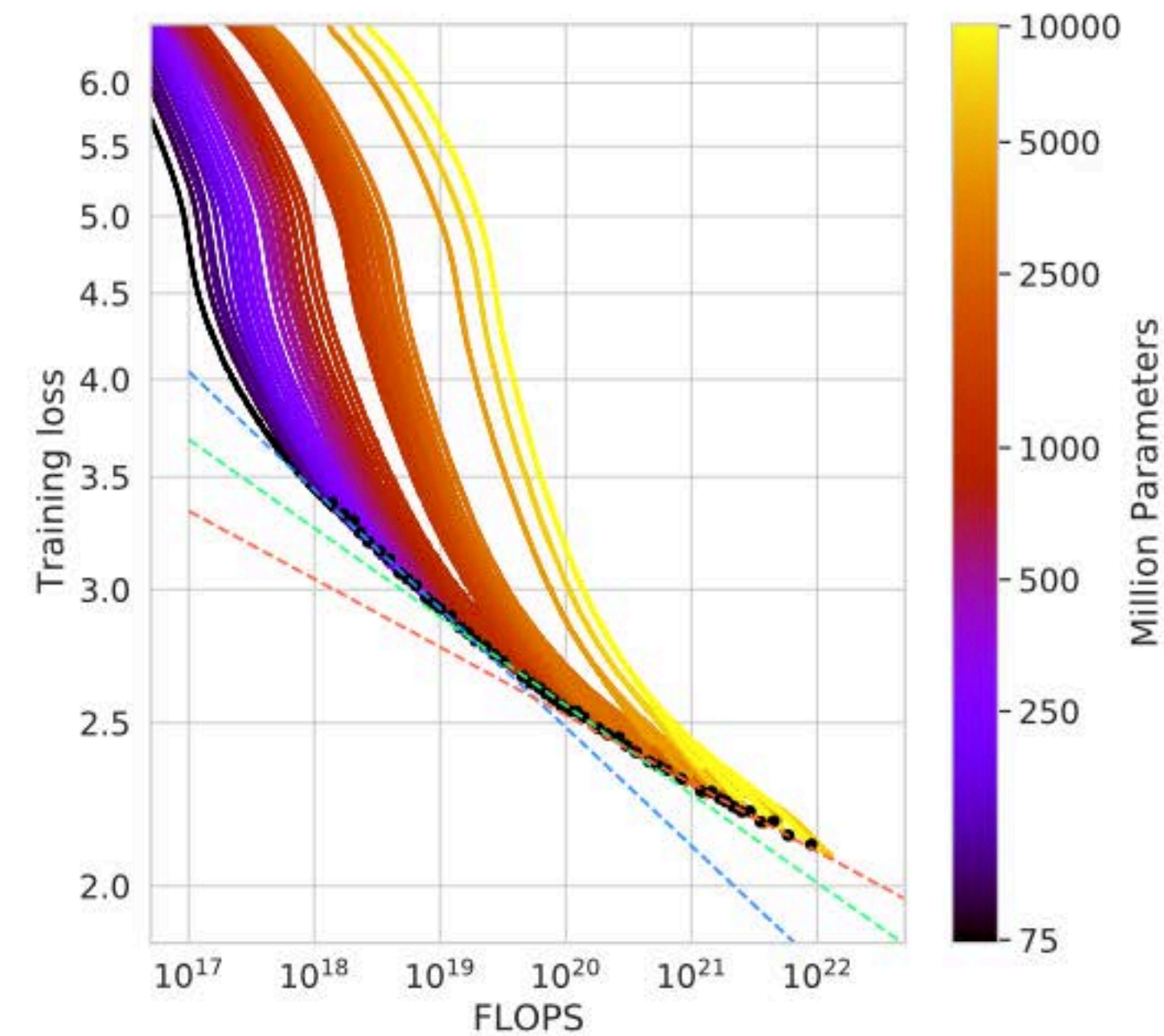


Figure A5 | **Training curve envelopes.** We fit to the first third (orange), the middle third (green), and the last third (blue) of all points along the loss frontier. We plot only a subset of the points.

Approaches to estimate loss as a function of data and parameter count

- minimize $L(N,D)$ with constraint $FLOPs(N,D) = budget$
- Approach 1: for each fixed model size, record entire training curve (=different #training tokens, as they use 1 epoch only)
- Approach 2: fixed #FLOPs, use final training loss for different N,D
- Approach 3: fit parametric loss function

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

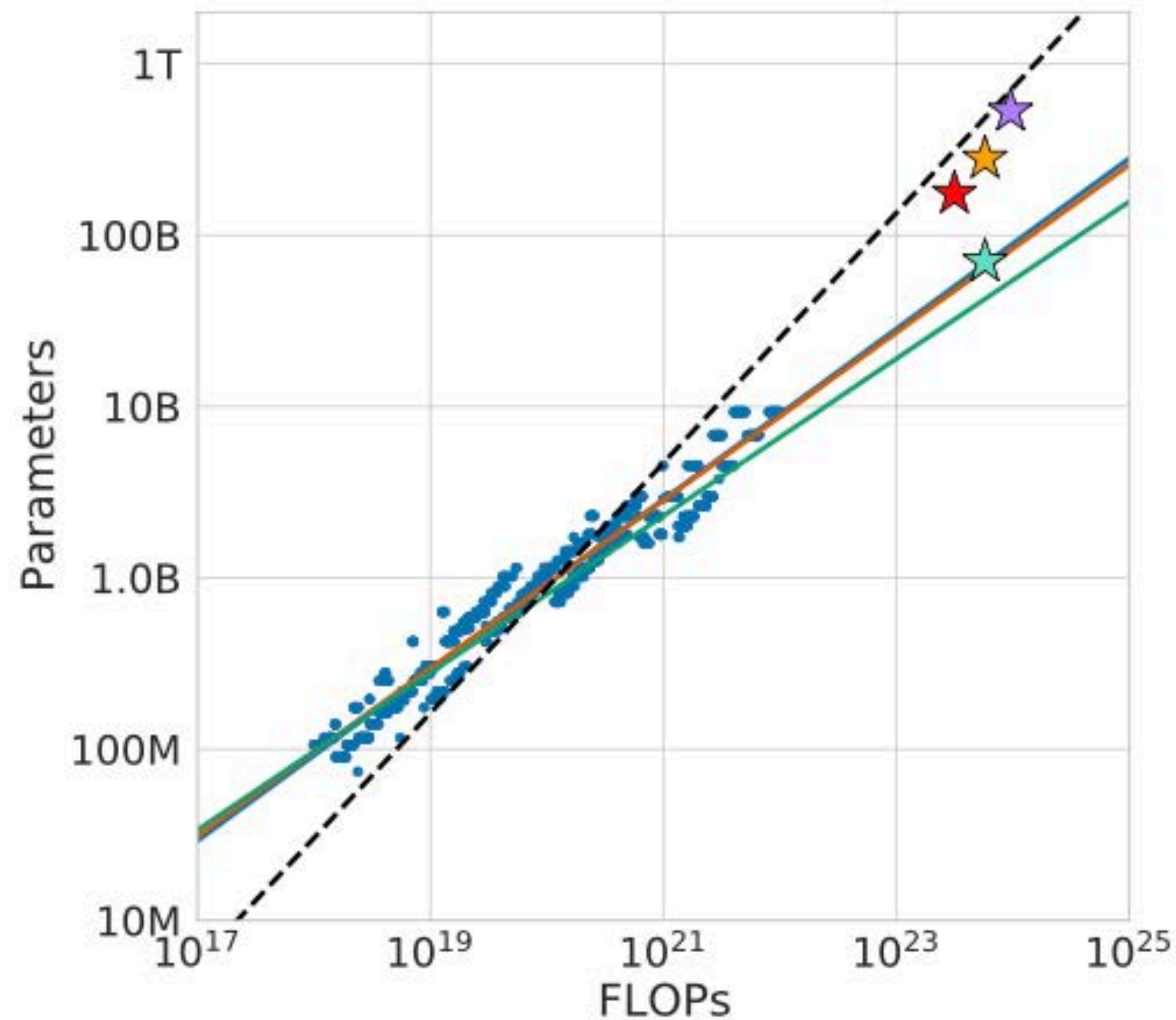
- for all 3: optimal scaling of #parameters N / data D with compute C is roughly

$$N \propto C^{0.49}, D \propto C^{0.51} \quad \text{vs Kaplan et al: } \#params \propto C^{0.73}, \text{ data} \propto C^{0.27}$$

=> existing large models would need much more training data.

Better allocation: **smaller model, more data/steps**

Predictions: extrapolation



For a given #FLOPs, how many parameters should we allocate?
(compute = 6 x data x params)

- Approach 1
- Approach 2
- Approach 3
- Kaplan et al (2020)
- Chinchilla (70B)
- Gopher (280B)
- GPT-3 (175B)
- Megatron-Turing NLG (530B)

Discussion / limitations

- use (smoothed) training loss as a proxy for test loss
- expensive experiments: only two large-scale training runs (Gopher, Chinchilla), no additional tests at intermediate scales
- assume a power law for $L(N)$ — but some concavity at high compute budgets, so it may not hold universally
- only one epoch of training (each data point seen once)

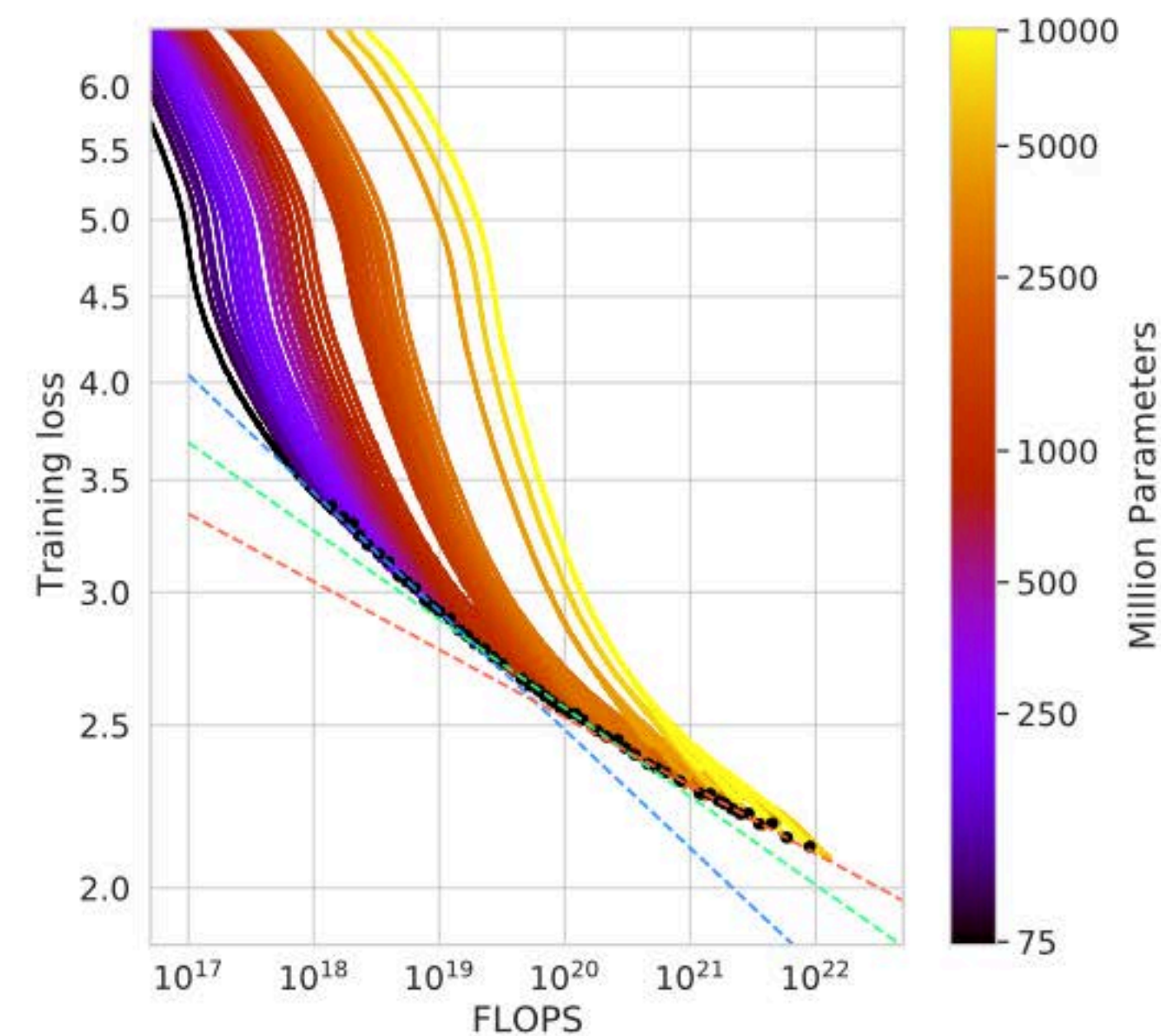


Figure A5 | **Training curve envelopes.** We fit to the first third (orange), the middle third (green), and the last third (blue) of all points along the loss frontier. We plot only a subset of the points.

Many more on scaling laws...

- many earlier works, too (*Tan & Le 2019, Cho et al 2015, Miceli Barone et al 2017, Johnson et al 2018, Hestness et al 2017*)
- vision models tend to have better scaling in model size than transformers
 $L(N) \propto 1/N^\alpha$ for vision models (CNNs), with larger $\alpha \approx 0.5$ (*Rosenfeld et al 2019, Sharma et al*)
- power law for $L(D)$ across vision, language, speech tasks (*Hestness et al 2017*)
- power laws in theoretical bounds: model (*Yarotsky, 2018*), data size (*Liang et al, 2019*)
-

Can we explain scaling laws? One hypothesis

A hypothesis: data manifold

Scaling Laws from the Data Manifold Dimension

Utkarsh Sharma

*Department of Physics and Astronomy
Johns Hopkins University
Baltimore, MD 21218, USA*

USHARMA7@JHU.EDU

Jared Kaplan

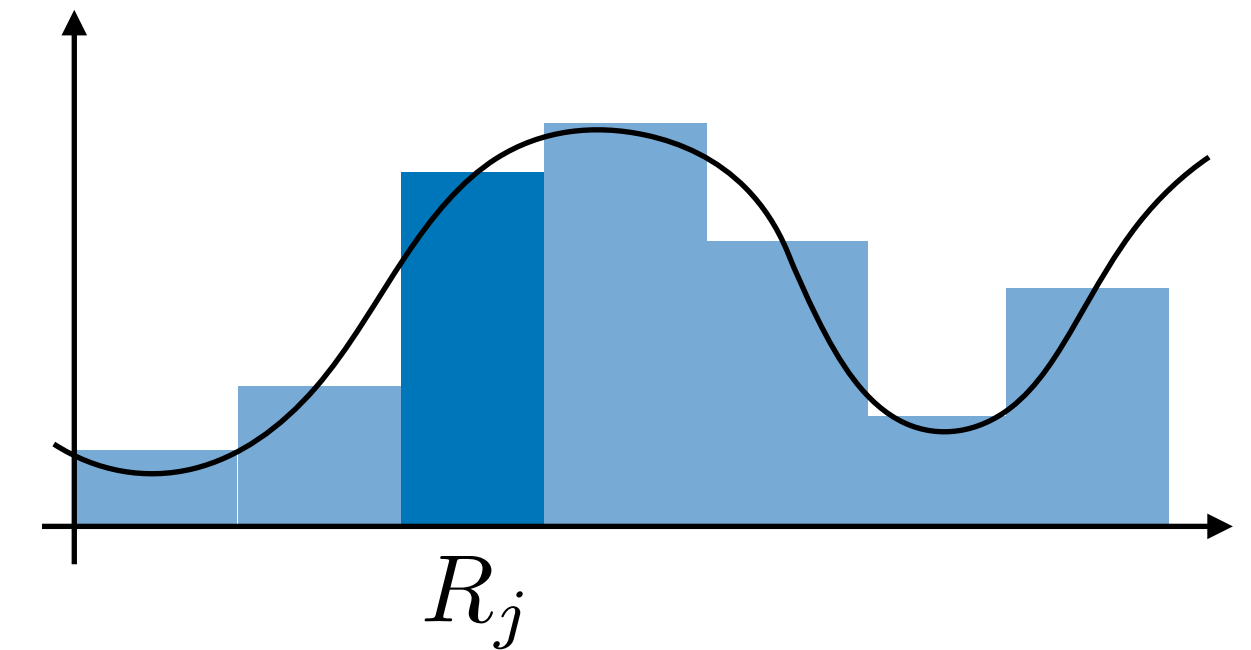
*Department of Physics and Astronomy
Johns Hopkins University
Baltimore, MD 21218, USA*

JAREDK@JHU.EDU

- Idea: neural networks are doing regression on an embedded data manifold

A toy model

- Toy model: approximate Lipschitz function $f(x)$ on $[0,1]^d$ by a piecewise constant function $c(x)$.
- Use hypercubes with side length s :
 $N = s^{-d}$ regions (values).



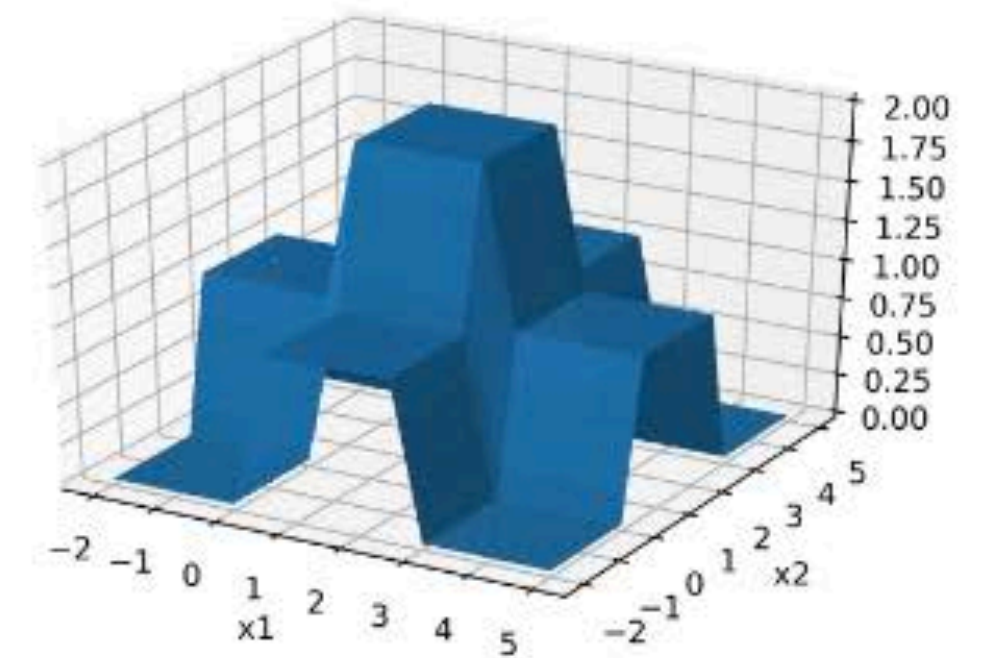
- MSE scales as $L(s) = \int_0^1 |f(x) - c(x)|^2 dx \lesssim \lambda^2(s^2)$

Lipschitz constant

- So, up to constant factors, $L(N) \lesssim \frac{1}{N^{2/d}}$

- For piecewise linear approximations (ReLU), MSE and cross entropy scale as s^4 , so

$$L(N) \lesssim N^{-4/d}$$



Relation to neural networks: hypothesis

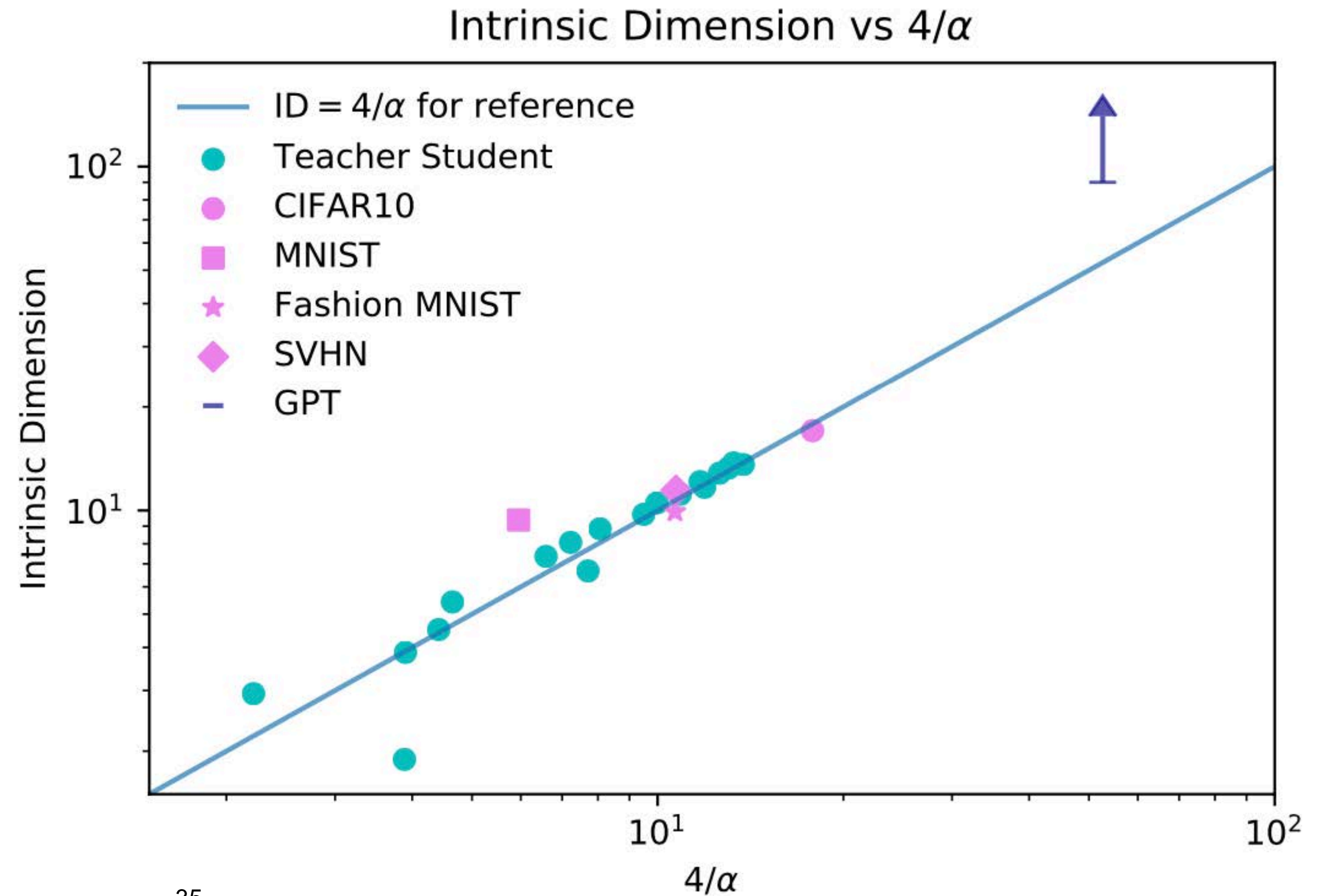
- For piecewise linear approximations (ReLU), MSE and cross entropy scale as s^4 , so $L(N) \lesssim N^{-4/d}$
- Common belief: NN map data into a low-dimensional “manifold” that depends on data, loss/task
- If we take d to be intrinsic dimension of data manifold, then we get a scaling law

$$L(N) \propto \frac{1}{N^\alpha} \text{ with } \alpha \approx 4/d$$

- Suggests that scaling exponent is strongly related to data and task: different models will scale similarly on the same data

Empirical evidence

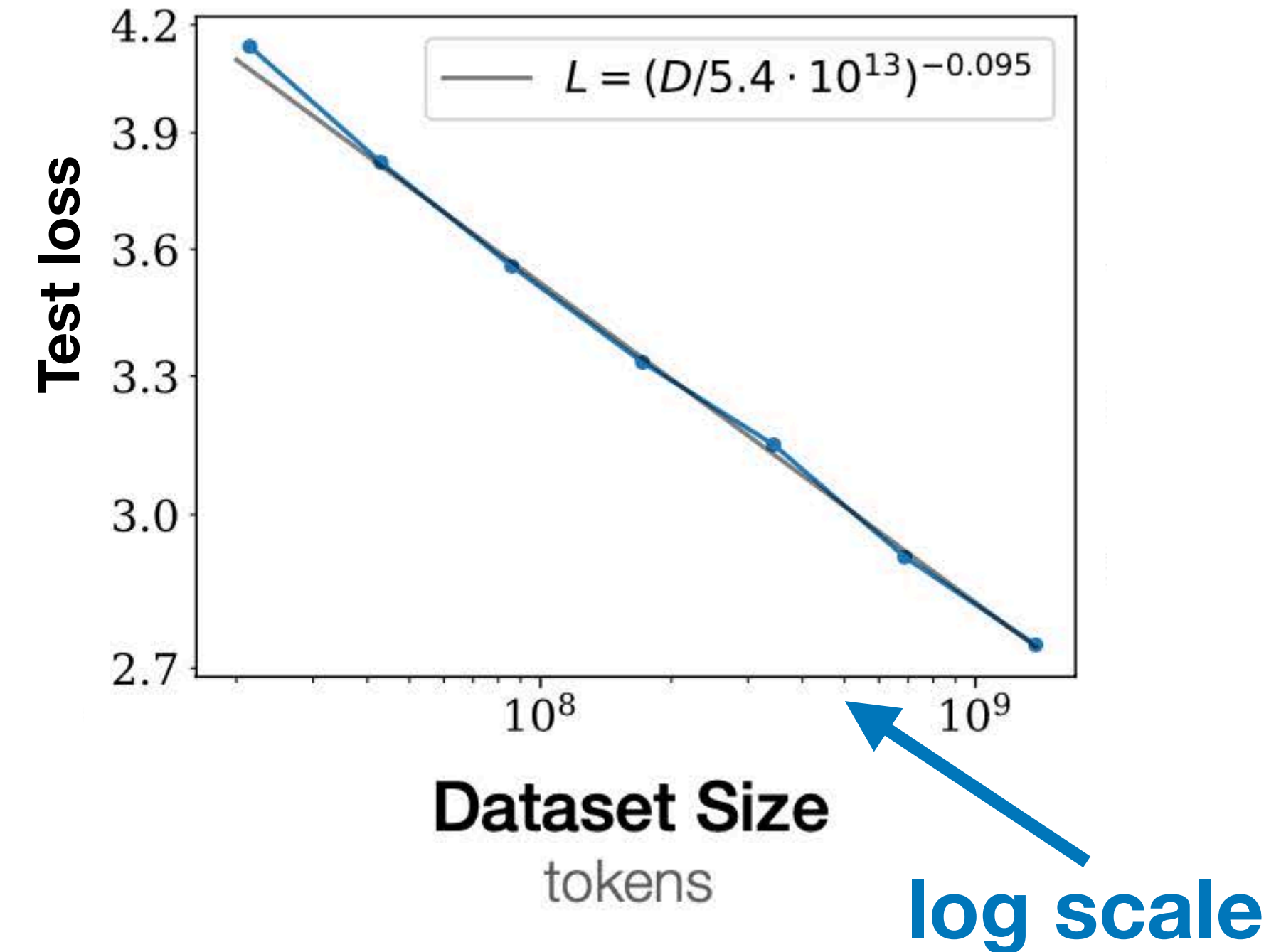
- Hypothesis: $L(N) \propto \frac{1}{N^\alpha}$ with $\alpha \approx 4/d$
- Experiment: measure intrinsic dimension of embedding and relate to scaling exponent α
- Good fit for vision models, only upper bound for transformer: $\alpha \gtrsim 4/d$ from upper bound on $L(N)$



Breaking power law scaling

Breaking power law scaling

- Power laws with small exponents: a small improvement in loss may need an order of magnitude more data/parameters
- Diminishing returns: suggests that many training examples are highly redundant
- Pruning intelligently can lead to better scaling (e.g. Sorscher et al, 2022)

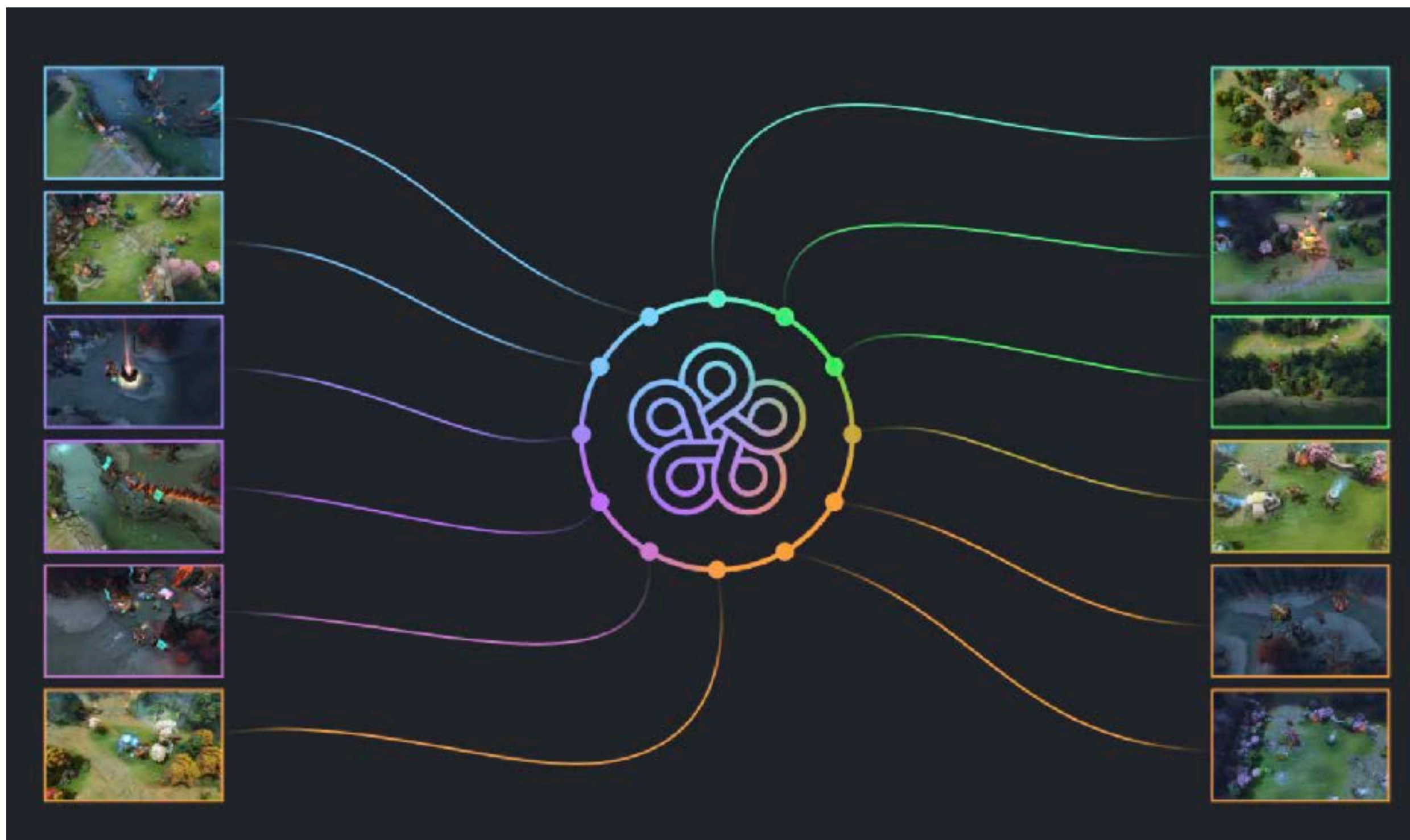


Critical batch size

What should the batch size be?

- Bigger batches are more costly but provide better estimates of the full-batch gradient
- $\text{FLOPS}(16 \text{ SGD steps with batch size } 16) = \text{FLOPS}(1 \text{ SGD step with batch size } 256)$. Which yields better performance?
- FLOPs aren't everything: one compute node can only handle some max batch size at once (in parallel). If you use bigger batches then you have to split up across nodes or do sequential computation. Splitting incurs communication costs. Sequential computation incurs time cost.

Left © OpenAI. Right © OpenAI, McCandlish, et al. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <https://ocw.mit.edu/help/faq-fair-use/>



An Empirical Model of Large-Batch Training

Sam McCandlish*
OpenAI
sam@openai.com

Jared Kaplan
Johns Hopkins University, OpenAI
jaredk@jhu.edu

Dario Amodei
OpenAI
damodei@openai.com

and the OpenAI Dota Team[†]

Abstract

In an increasing number of domains it has been demonstrated that deep learning models can be trained using relatively large batch sizes without sacrificing data efficiency. However the limits of this massive data parallelism seem to differ from domain to domain, ranging from batches of tens of thousands in ImageNet to batches of millions in RL agents that play the game Dota 2. To our knowledge there is limited conceptual understanding of why these limits to batch size differ or how we might choose the correct batch size in a new domain. In this paper, we demonstrate that a simple and easy-to-measure statistic called the *gradient noise scale* predicts the largest useful batch size across many domains and applications, including a number of supervised learning datasets (MNIST, SVHN, CIFAR-10, ImageNet, Billion Word), reinforcement learning domains (Atari and Dota), and even generative model training (autoencoders on SVHN). We find that the noise scale increases as the loss decreases over a training run and depends on the model size primarily through improved model performance. Our empirically-motivated theory also describes the tradeoff between compute-efficiency and time-efficiency, and provides a rough model of the benefits of adaptive batch-size training.

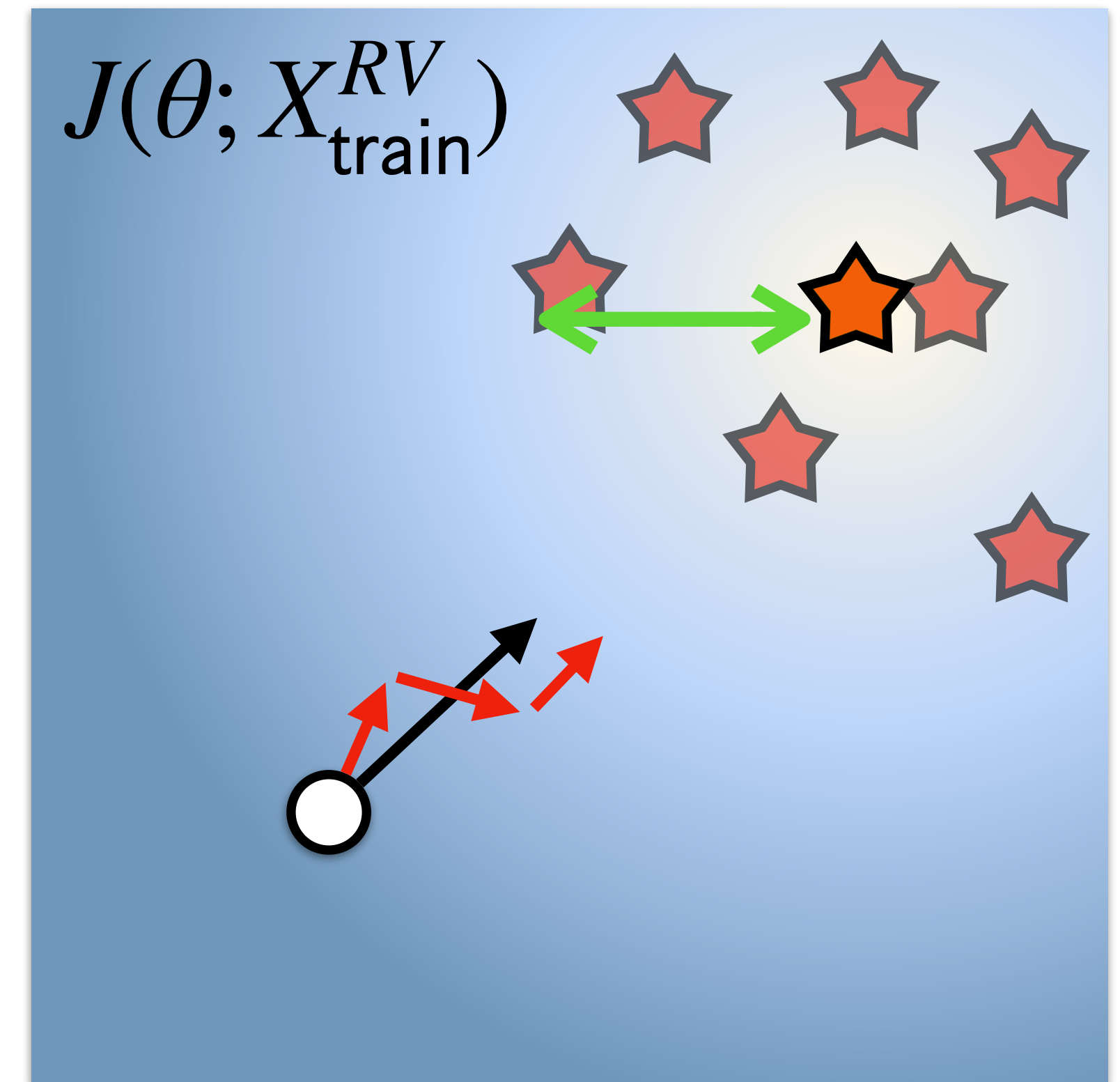
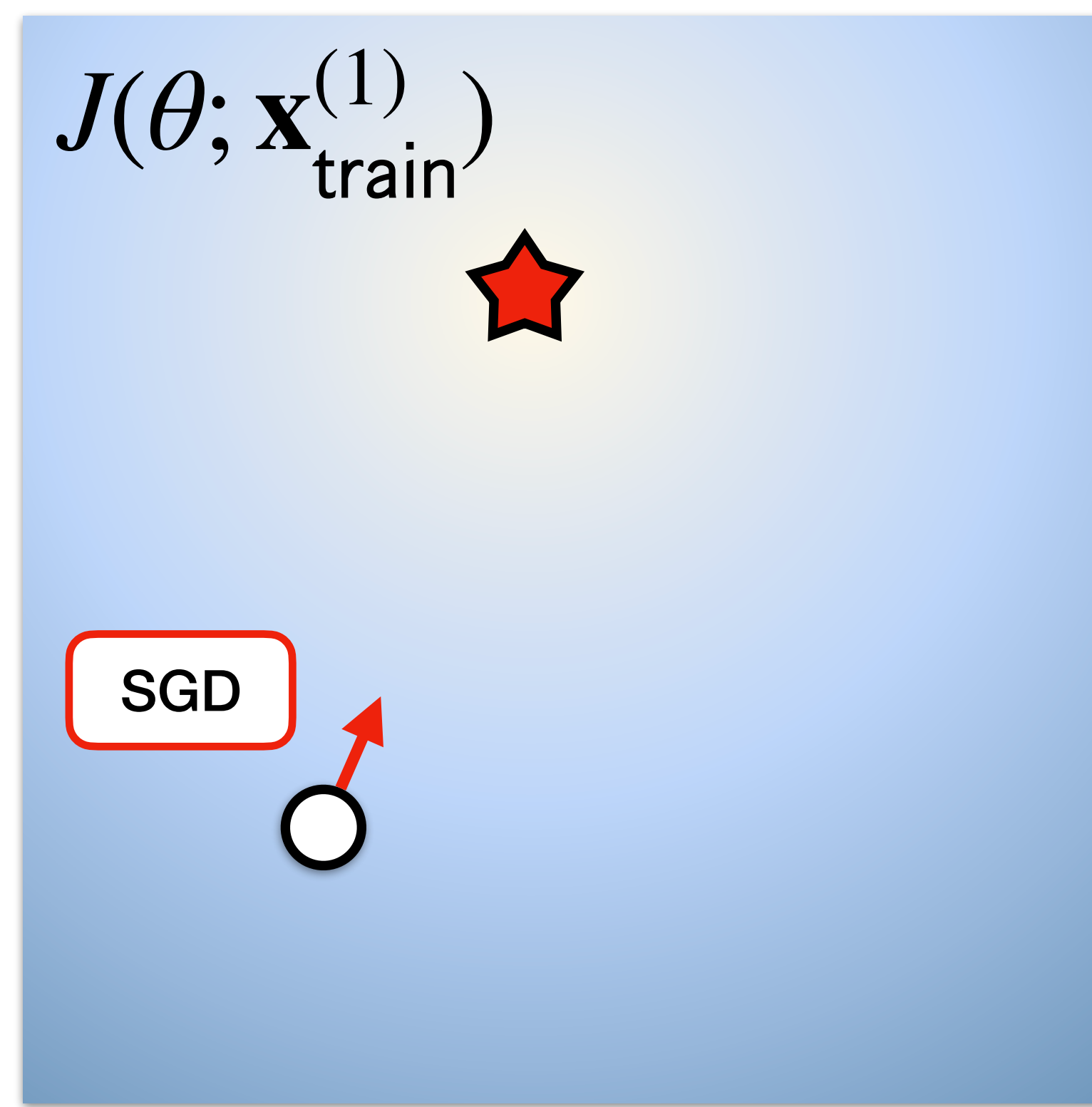
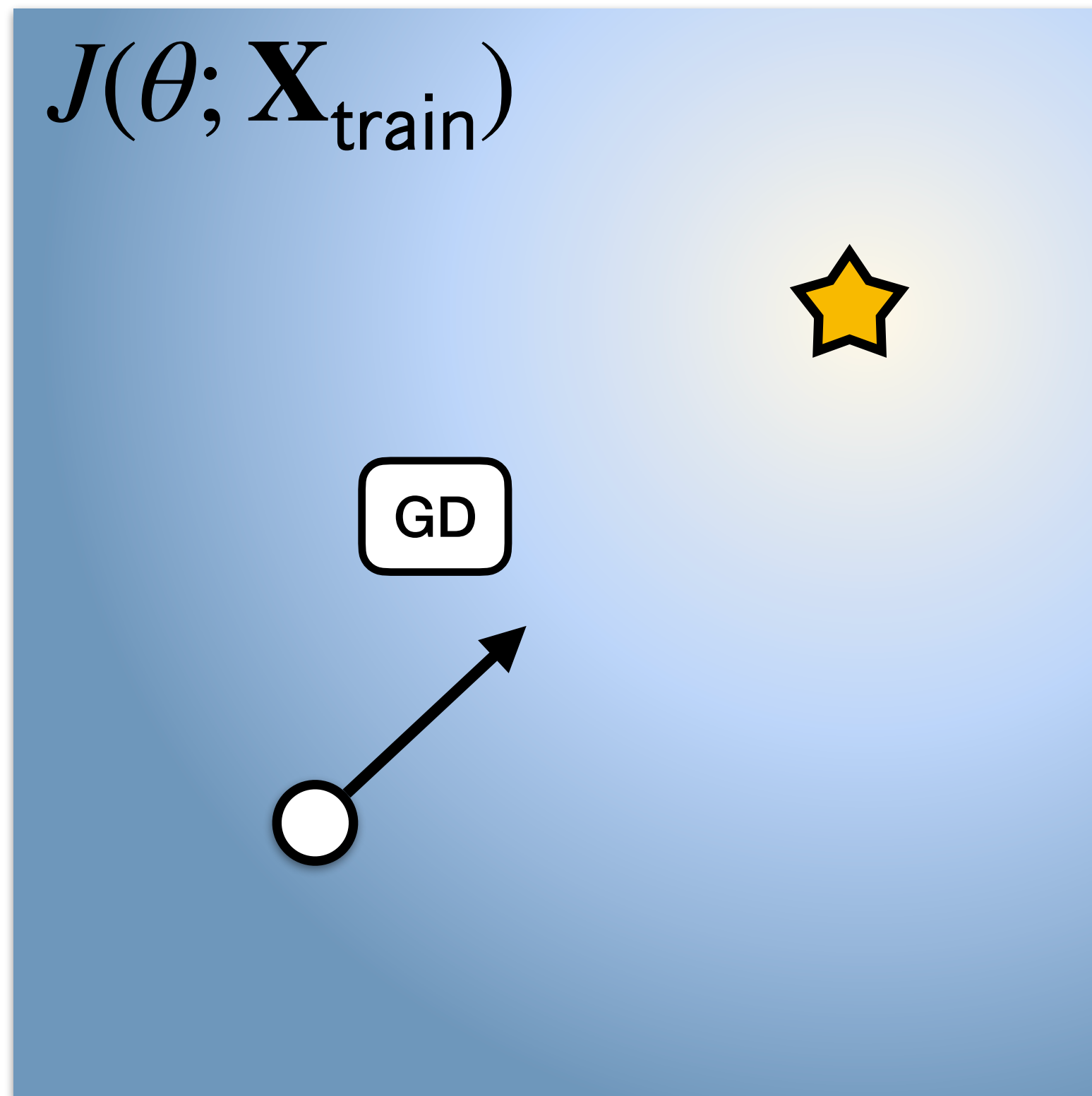
<https://openai.com/research/openai-five-defeats-dota-2-world-champions>

<https://cdn.openai.com/dota-2.pdf>

Batch size: ~2 million

<https://arxiv.org/abs/1812.06162>

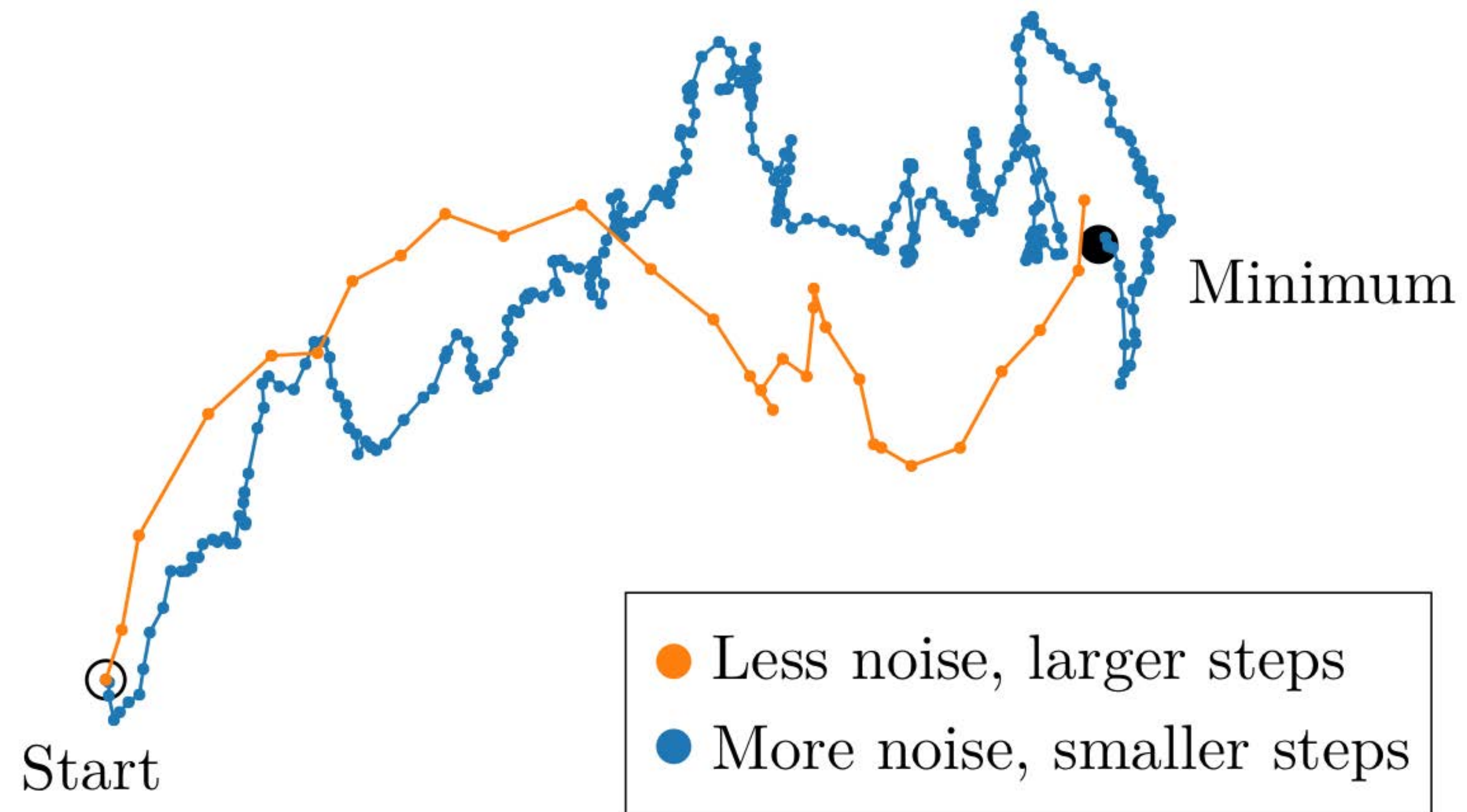
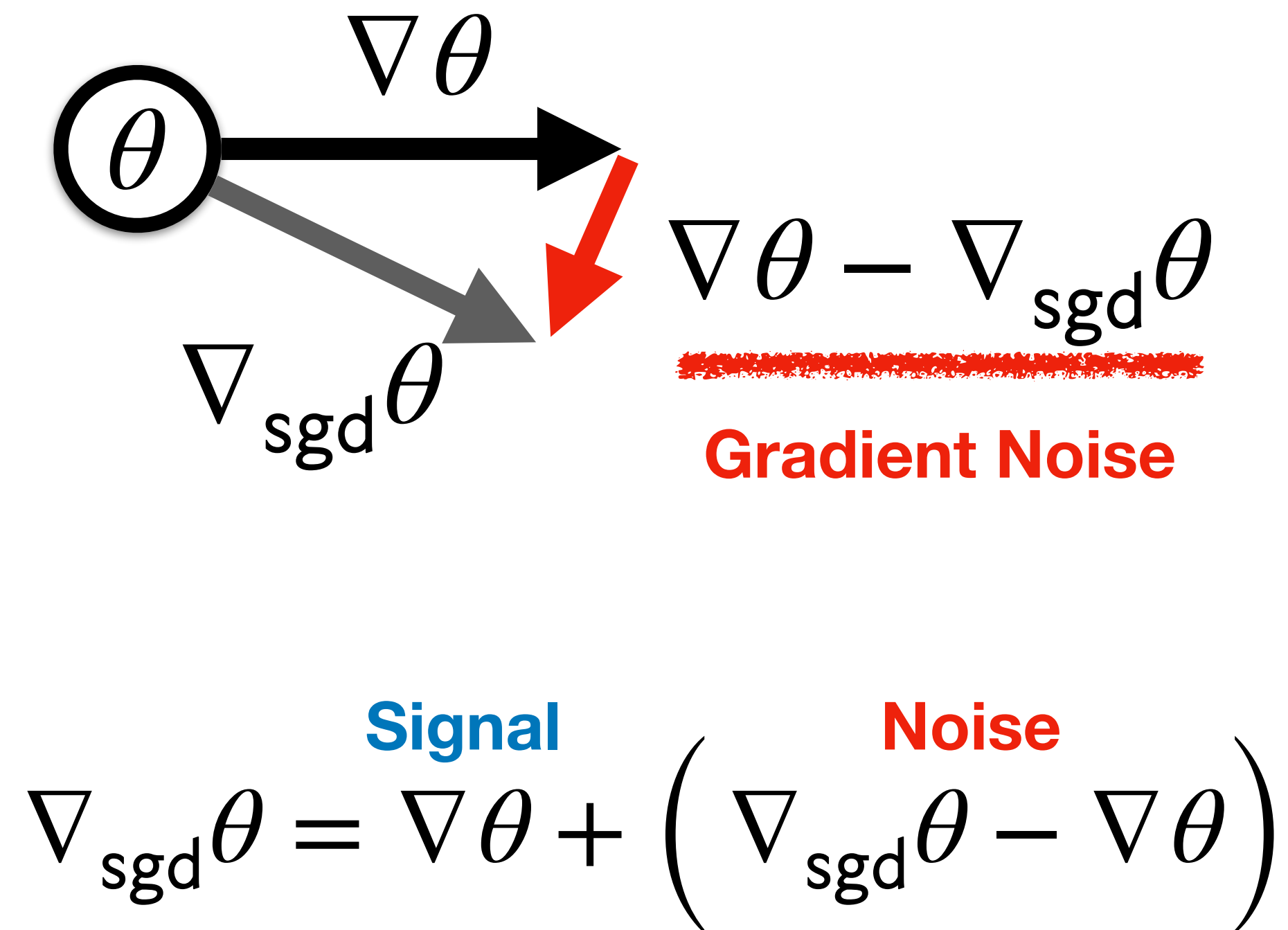
Variability of solution as a function
given your mini-batch size



[slide credit: adapted from Minyoung Huh]

McCandlish₄₁Kaplan, Amodei, "An Empirical Model of Large-Batch Training", 2018

Gradient noise



How much progress can one step of SGD make?

Optimal learning rate for SGD?

$$\mathbb{E} [\mathcal{L}(\theta + \eta G_{\text{sgd}})] = \mathcal{L}(\theta) - \eta |G|^2 + \frac{\eta^2}{2} \left(G^T H G + \frac{\text{trace}(H \Sigma)}{B} \right)$$

covariance of G_{sgd} , with $B=1$
↙

Take the derivative and solve for the learning rate

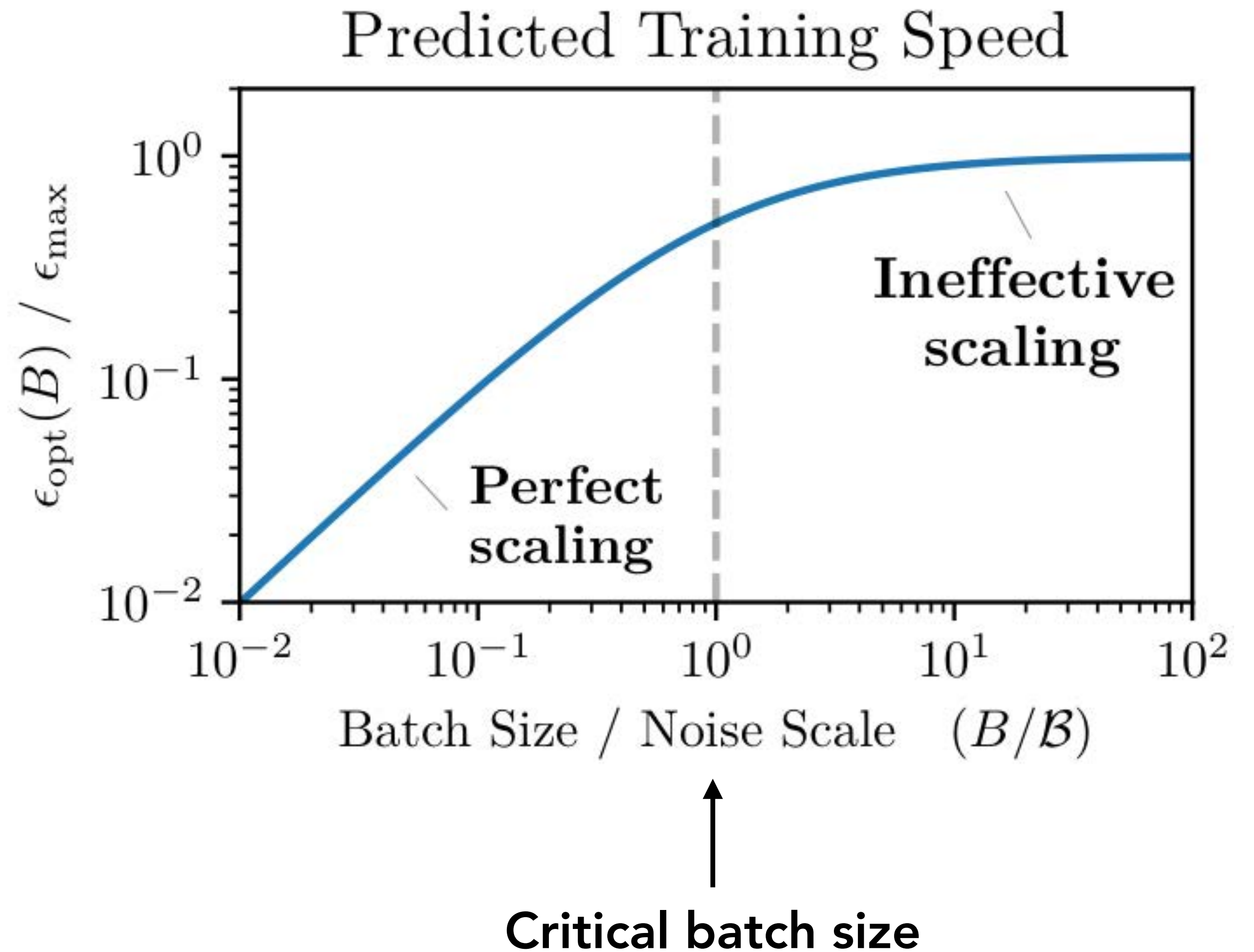
$$\eta_{\text{opt}} = \frac{1}{1 + \mathcal{B}_{\text{noise}}/B} \eta_{\text{max}} \quad \text{where} \quad \mathcal{B}_{\text{noise}} = \frac{\text{trace}(H \Sigma)}{G^T H G}$$

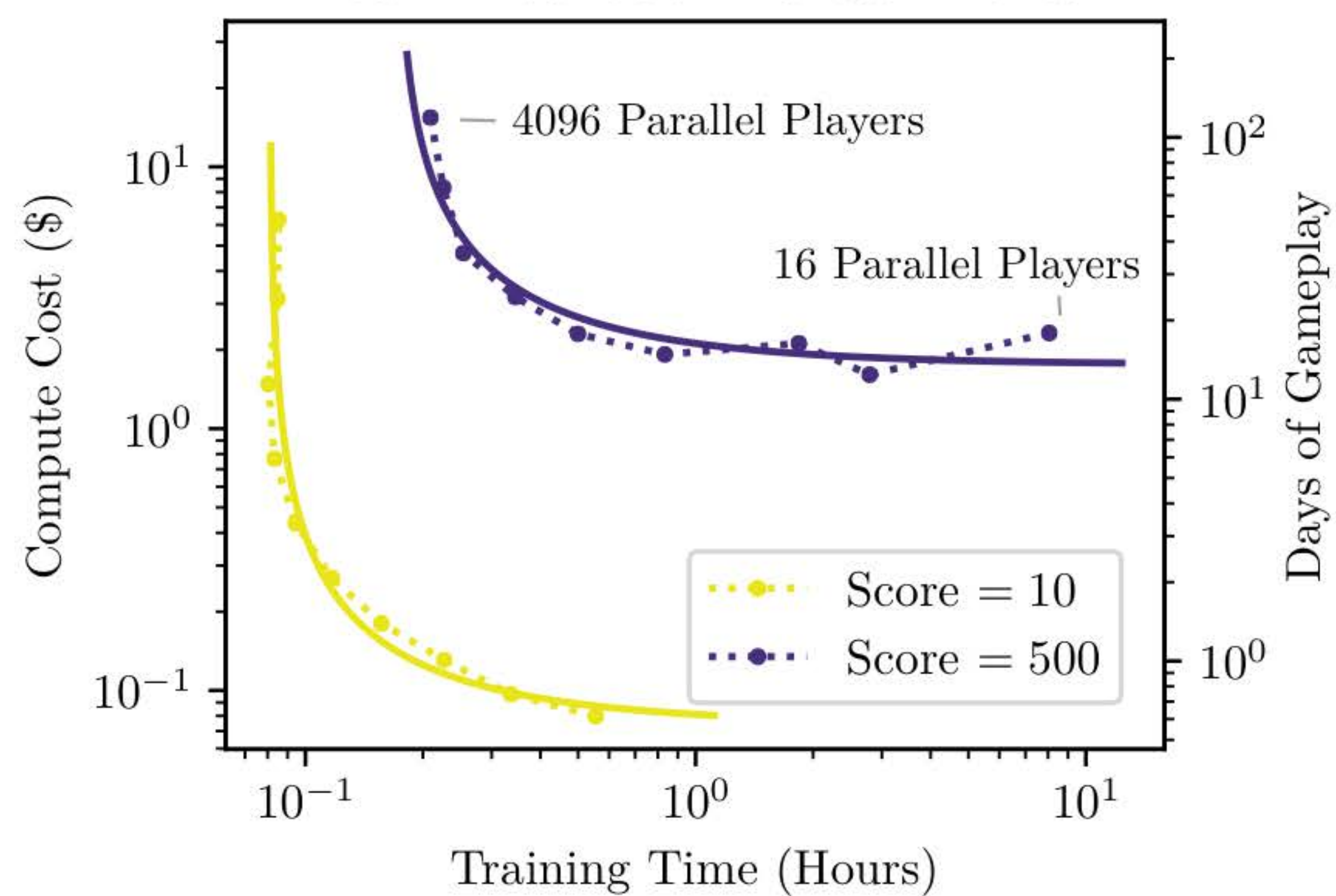
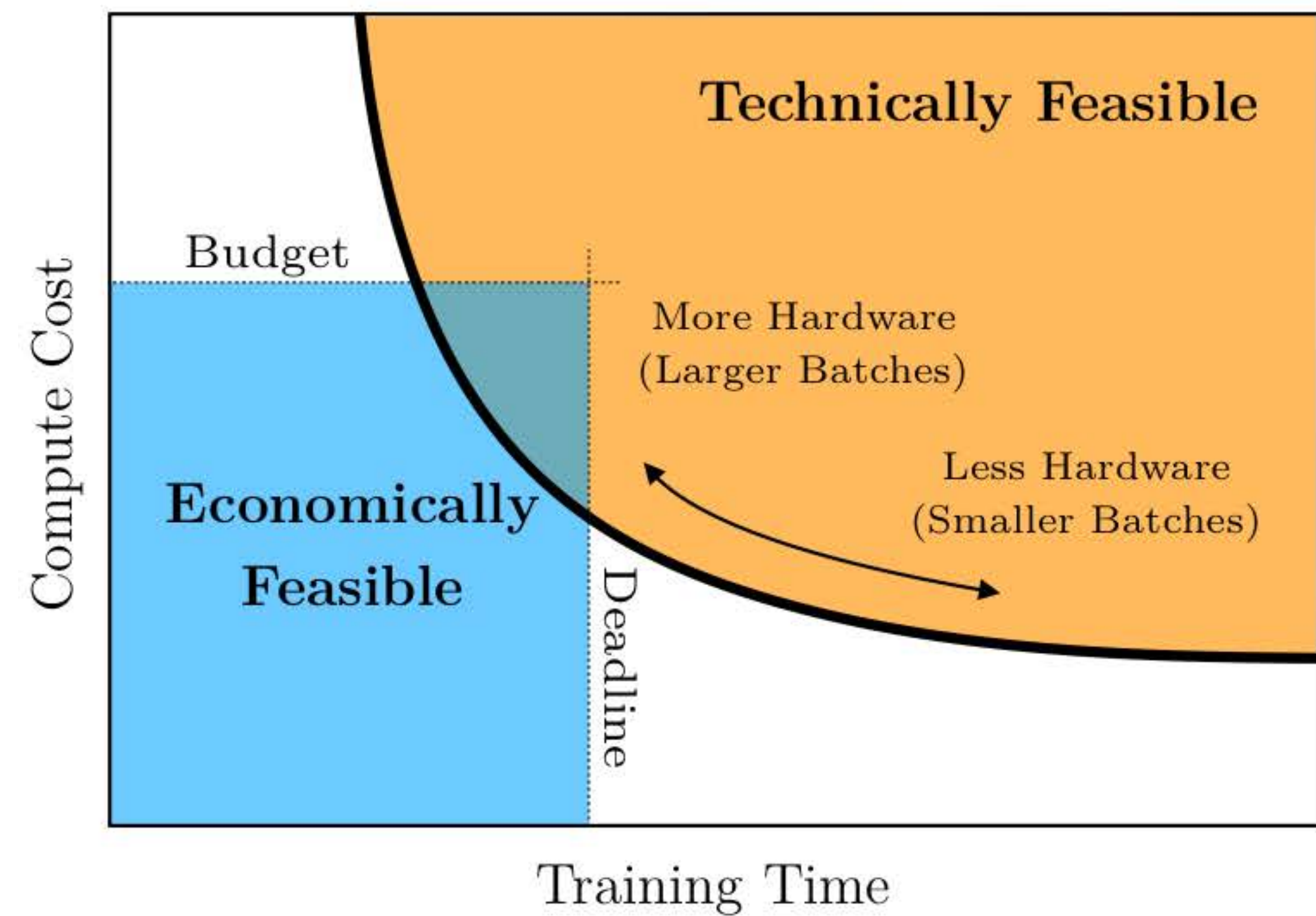
Plug the equation learning rate back into the equation

$$\Delta \mathcal{L}_{\text{opt}} = \frac{1}{1 + \mathcal{B}_{\text{noise}}/B} \Delta \mathcal{L}_{\text{max}} \quad \text{where} \quad \Delta \mathcal{L}_{\text{max}} = \frac{1}{2} \frac{|G|^4}{G^T H G}$$

- One step of SGD can make progress proportional to:

$$\frac{1}{1 + \mathcal{B}_{\text{noise}}/B}$$
- For small B , second term of denominator dominates, so increasing B linearly increases progress
- For large B , first term of denominator dominates, so increasing B has little effect

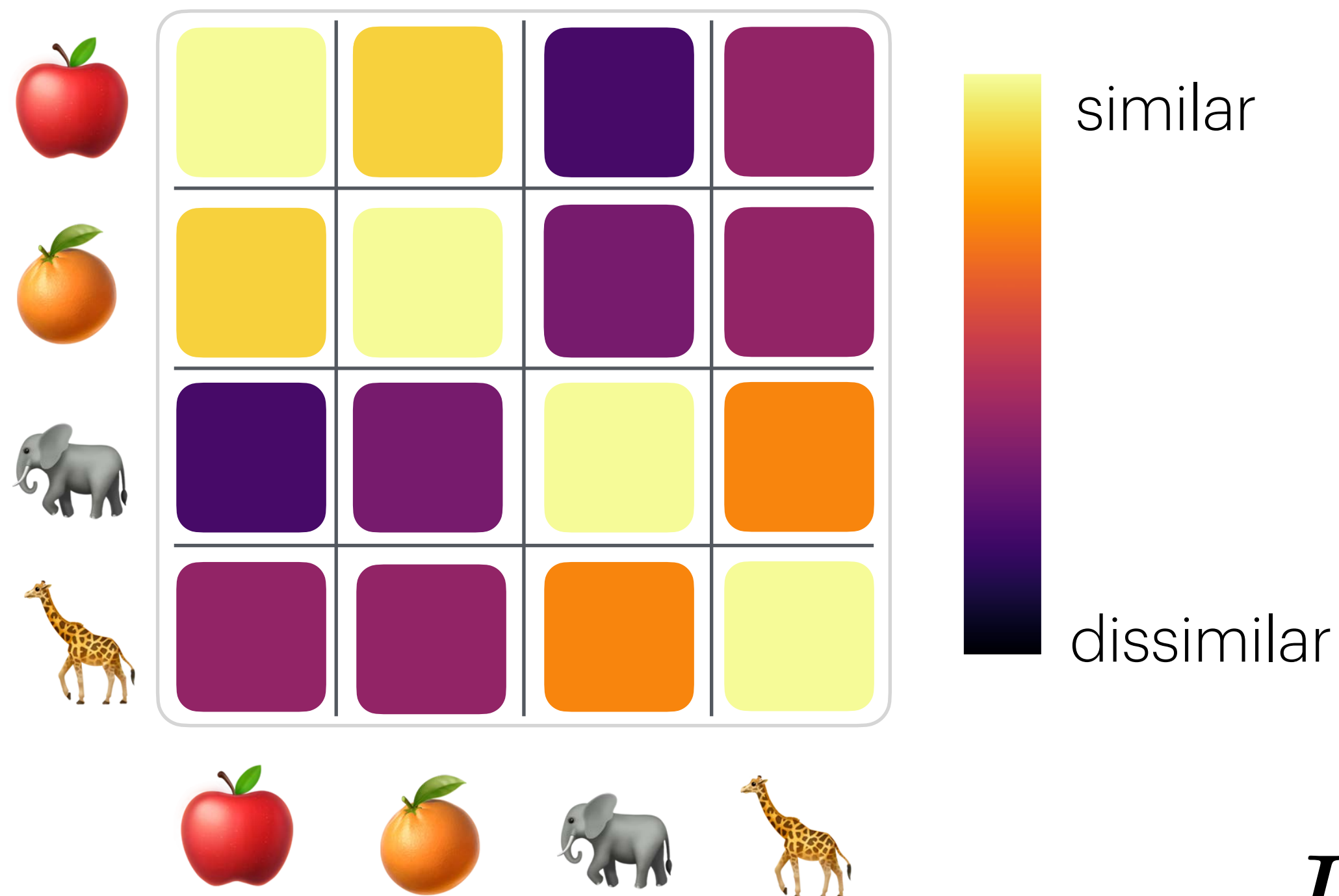




Representational convergence

Characterizing representations using kernels

K_{vision}



Restrict our attention to **vector embeddings**

$$f : \mathcal{X} \rightarrow \mathbb{R}^n$$

Characterize a representation
in terms of its **kernel**

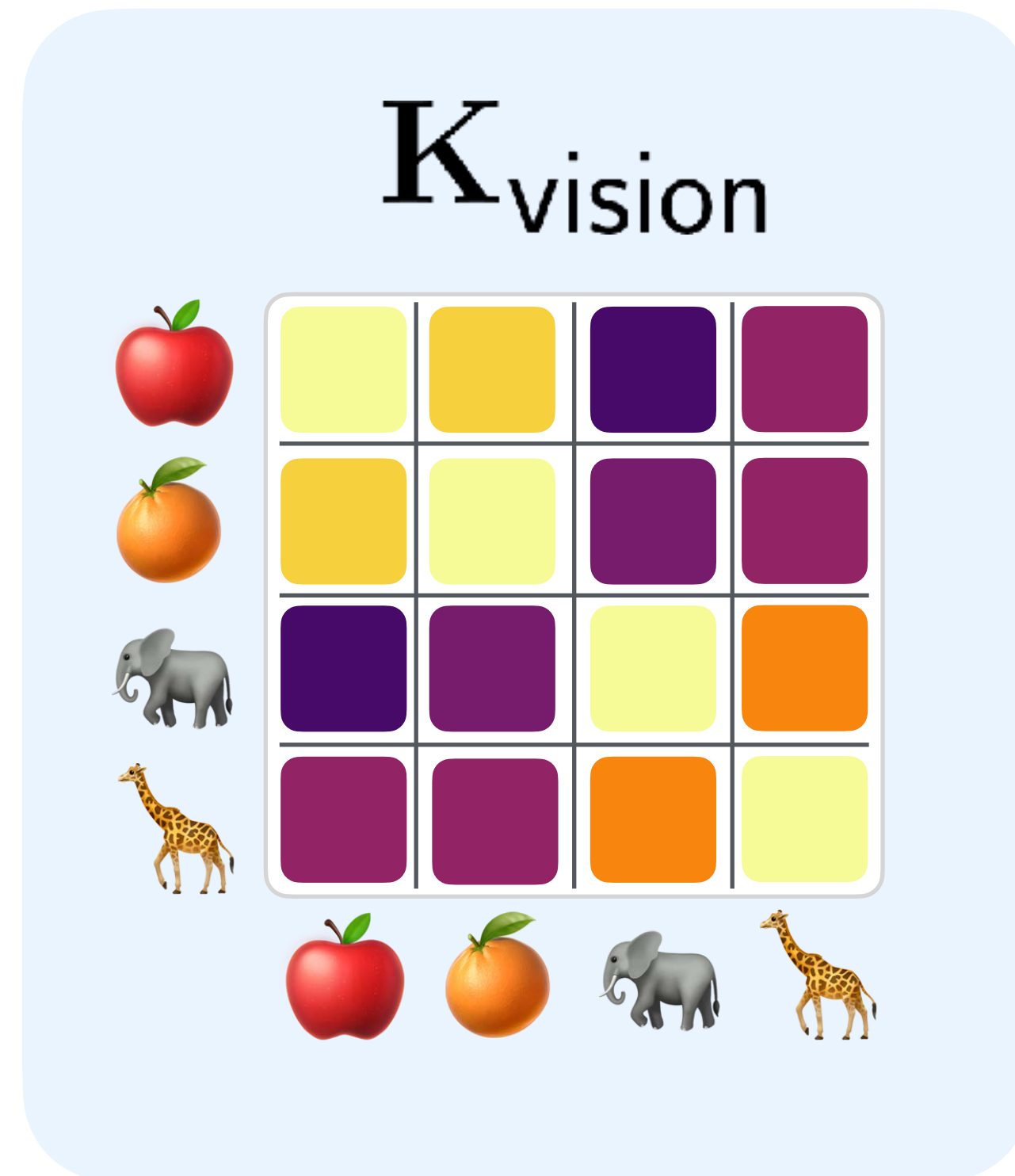
$$K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$$

$$K(x_i, x_j) = \langle f(\text{apple}), f(\text{orange}) \rangle$$

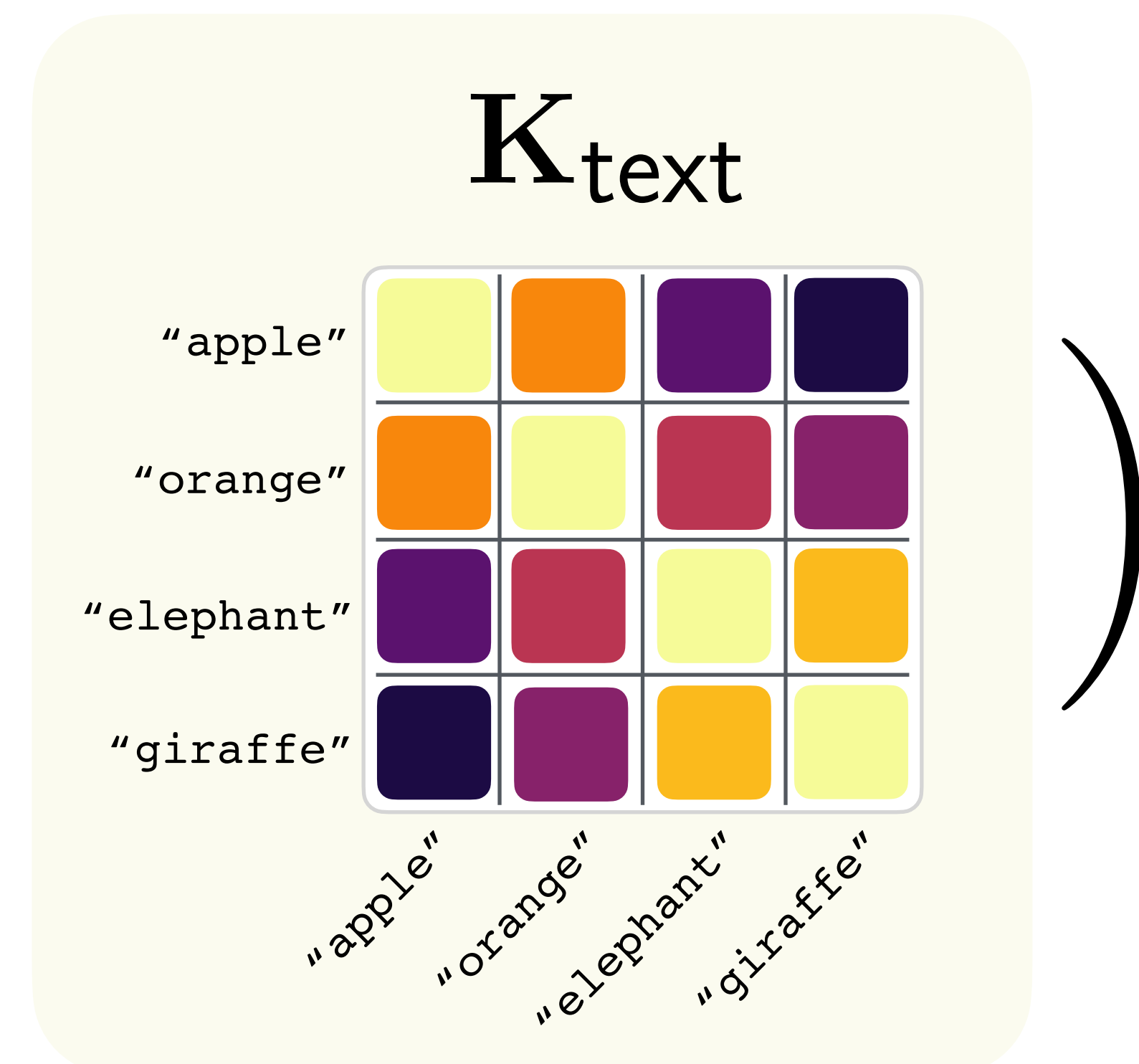
Icons © source unknown. All rights reserved. This content is excluded from our Creative Commons license. For more information, see

<https://ocw.mit.edu/help/faq-fair-use/>

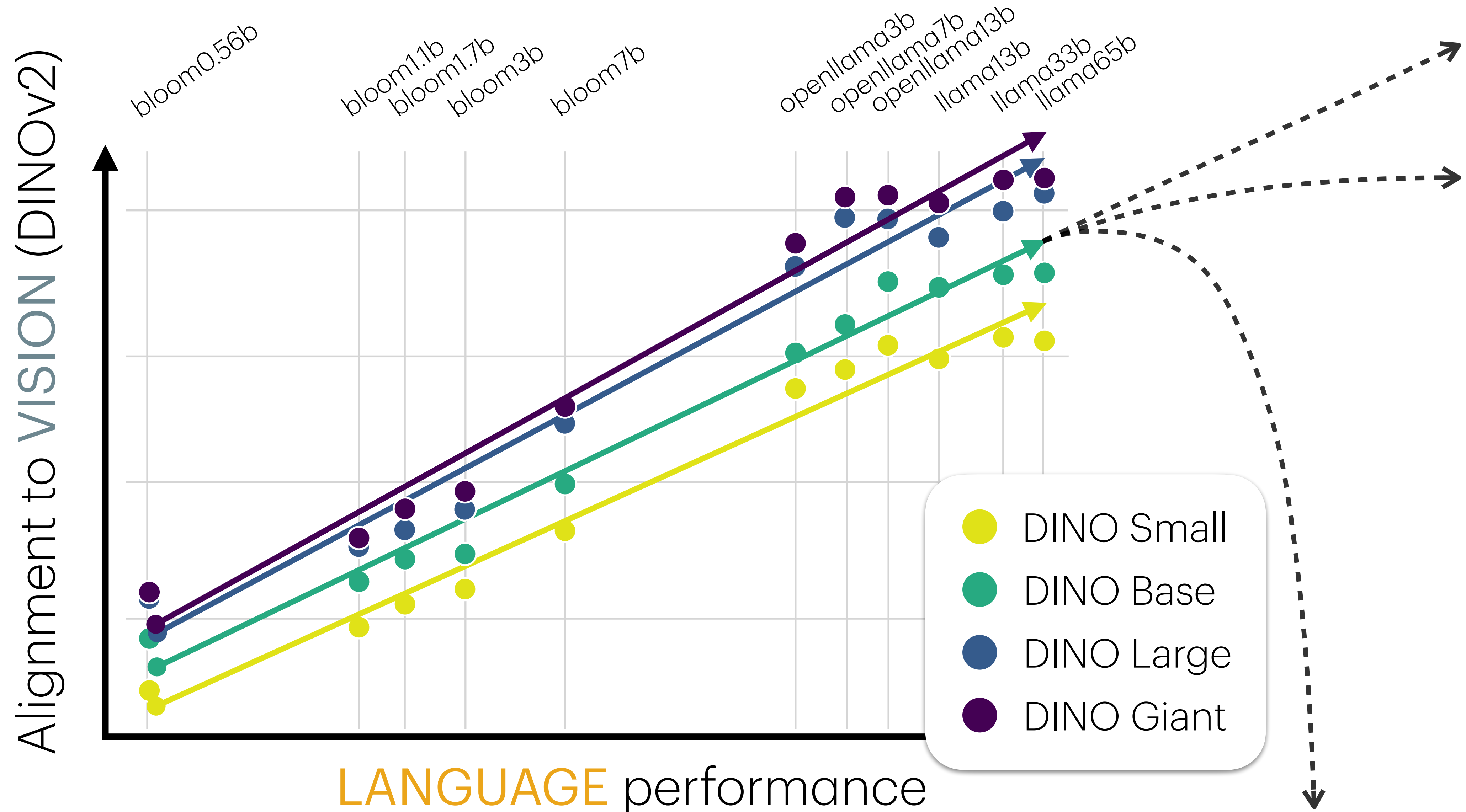
sim (

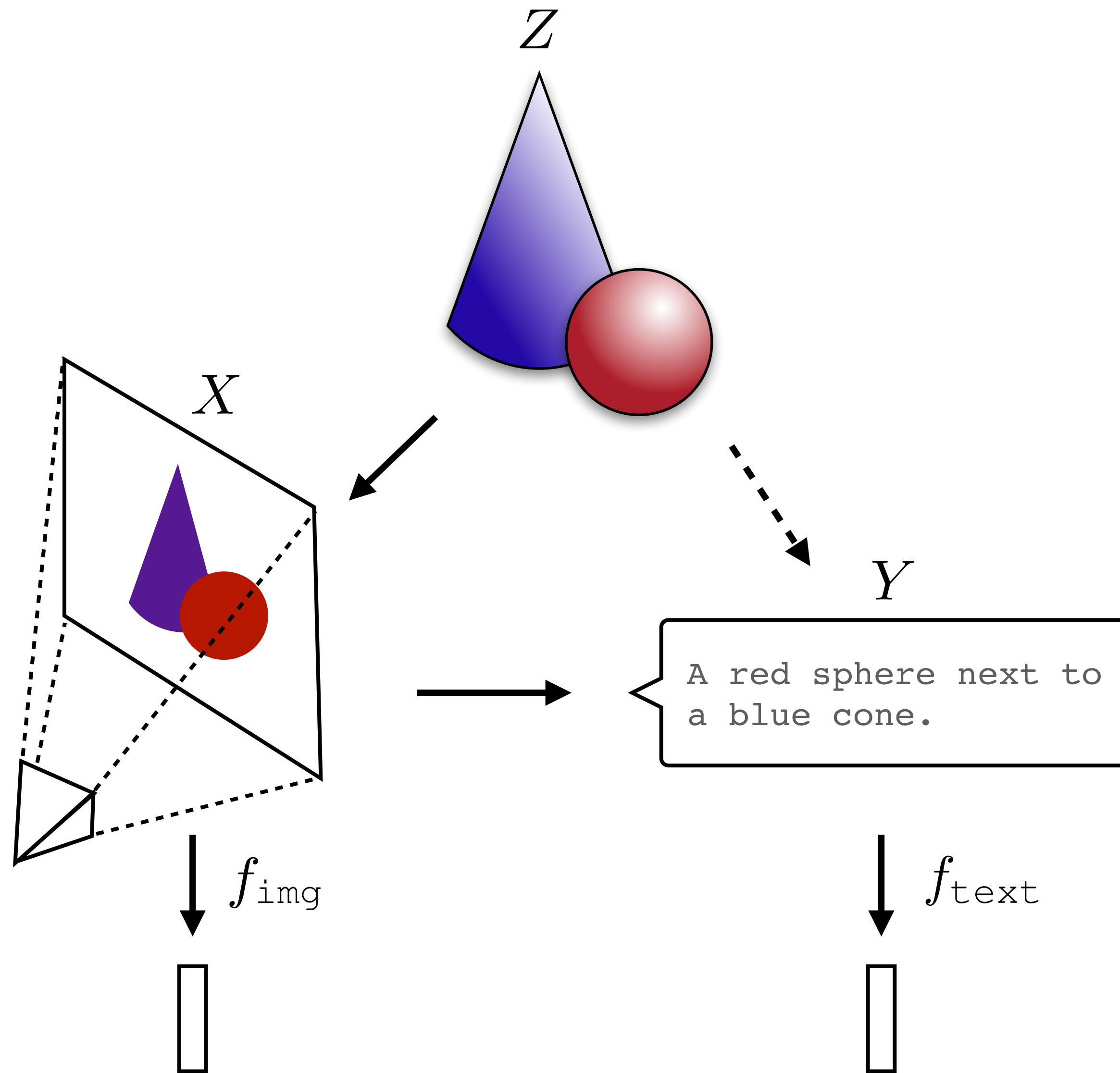


,



Scaling laws for cross-modal alignment





Courtesy of Huh, et al. Used under CC BY.

Huh*, Cheung*, Wang*, Isola*, "The Platonic Representation Hypothesis," ICML 2024

Summary

- Neural scaling laws: predict test loss as a function of resources and model hyperparameters
- Allows to “optimally” allocate compute resources
- Power law scaling in model size, data size across variety of models and tasks
- Actual parameters hard to measure, large extrapolation nontrivial

MIT OpenCourseWare

<https://ocw.mit.edu>

6.7960 Deep Learning

Fall 2024

For information about citing these materials or our Terms of Use, visit: <https://ocw.mit.edu/terms>