

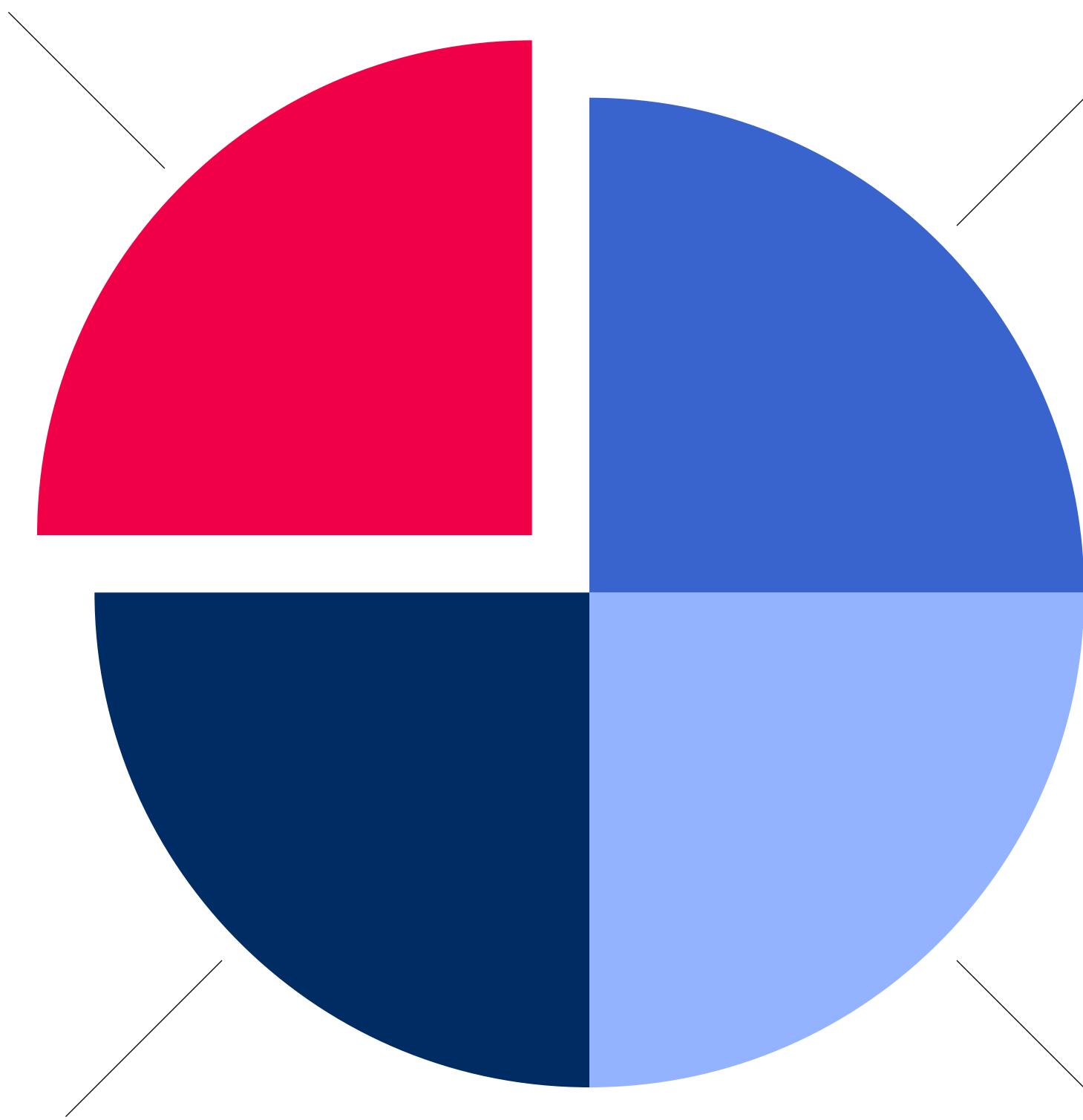
CS-461

Foundation Models and Generative AI

Learning at Scale: Supervised, Self-Supervised, and Beyond

Charlotte Bunne, Fall Semester 2025/26

Learning Principles

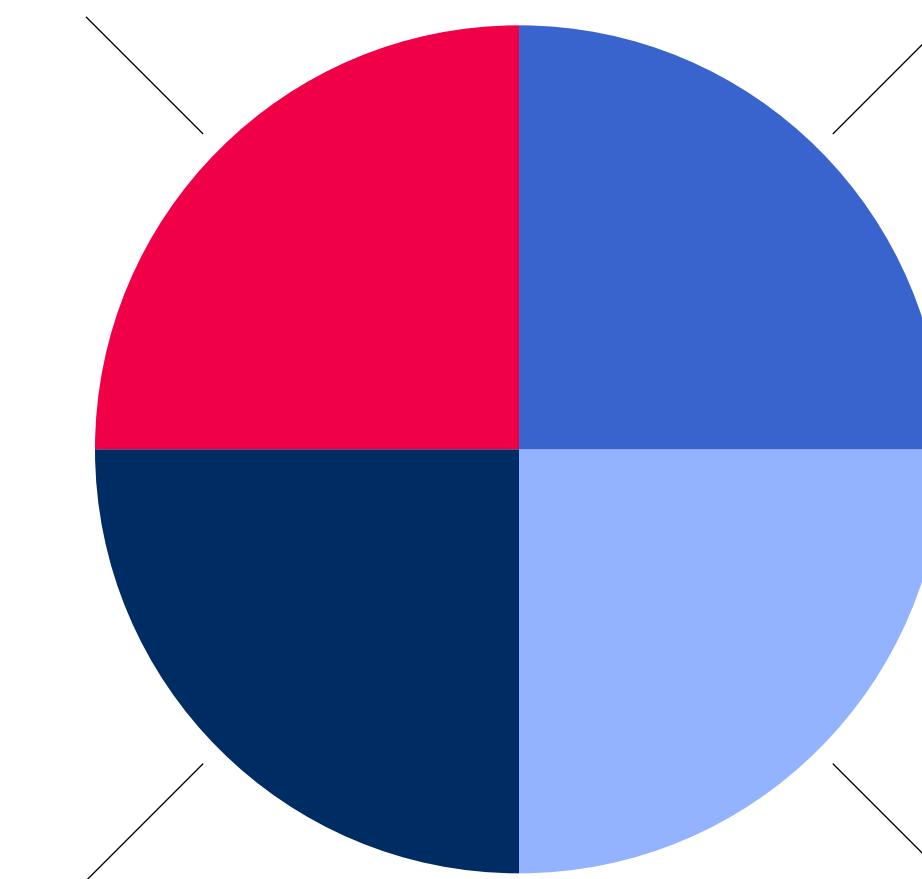


Expected Risk Minimization

All learning minimizes expected risk.

The difference is in what we know!

$$\min_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\ell(\theta; x, y)]$$



- **For supervised:** We have data samples x with labels y . Loss $\ell(f_{\theta}(x), y)$.
- **For unsupervised:** We *only* have data samples x . Loss defined on x alone, i.e., $\ell(f_{\theta}(x))$.
- **For semi-supervised:** We have data samples x and few labeled but many unlabeled samples.
The loss is a mixture of $\ell_{\text{sup}} + \lambda \cdot \ell_{\text{unsup}}$.
- **For self-supervised:** We only have data samples x but will create *synthetic labels* from x .
We will create the task from x itself, i.e., $\ell(f_{\theta}(x_{\text{partial}}), x_{\text{hidden}})$.

Supervised Learning

Supervised

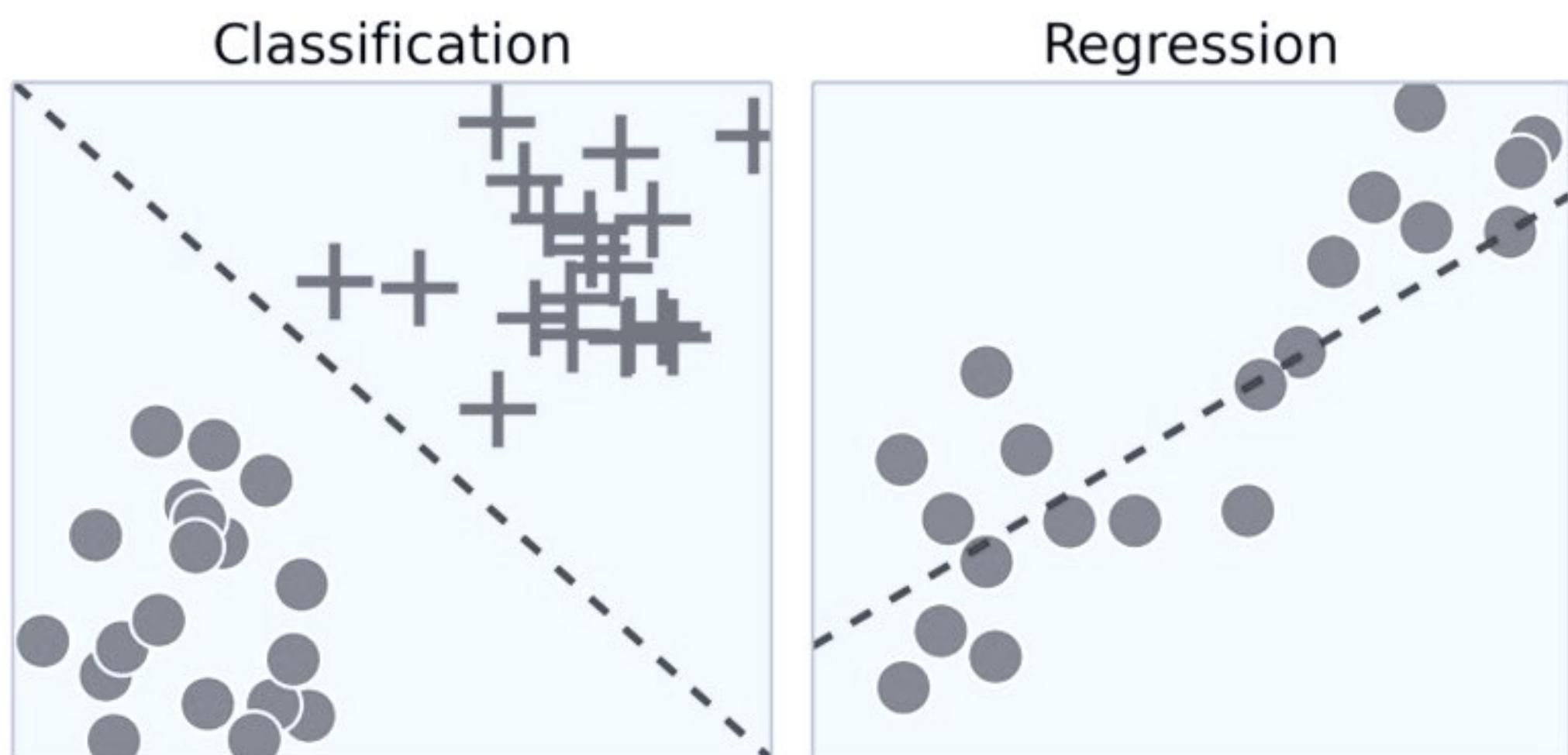
Goal: Learn a mapping $x \mapsto y$ from labelled data.

Data: $((x_1, y_1), \dots, (x_n, y_n))$, i.e., pairs of data points and corresponding label.

Model: find function f with parameters θ such that $f_\theta : x \mapsto f_\theta(x)$.

If label y is discrete: **classification**

y is continuous: **regression**



Supervised Learning

Supervised

Data: $((x_1, y_1), \dots, (x_n, y_n))$, i.e., pairs of data points and corresponding label.

Model: find function f with parameters θ such that $f_\theta : x \mapsto f_\theta(x)$.

If label y is discrete: **classification**

y is continuous: **regression**

Loss: For classification, e.g., **Cross-Entropy Loss or Negative Log-Likelihood (NLL)**

$$\mathcal{L}_{\text{CE}} = -\frac{1}{n} \sum_{i=1}^n \sum_{k=1}^K y_{i,k} \log (\hat{y}_{i,k})$$

K classes

$y_{i,k}$: true label

$\hat{y}_{i,k} = p_\theta(k | x_i)$: predicted probability for class k

For binary classification e.g., **Binary Cross-Entropy (BCE)**

$$\mathcal{L}_{\text{BCE}} = -\frac{1}{n} \sum_{i=1}^n [y_i \log (\hat{y}_i) + (1 - y_i) \log (1 - \hat{y}_i)]$$

Supervised Learning

Supervised

Data: $((x_1, y_1), \dots, (x_n, y_n))$, i.e., pairs of data points and corresponding label.

Model: find function f with parameters θ such that $f_\theta : x \mapsto f_\theta(x)$.

If label y is discrete: **classification**

y is continuous: **regression**

Loss: For regression, e.g., **Mean Squared Error (MSE)**

$$\mathcal{L} = \frac{1}{n} \sum_i (y_i - f_\theta(x_i))^2$$

or Mean Absolute Error (MAE), Huber loss (for outliers), etc.

Unsupervised Learning

Goal: Discover structure in unlabeled data x .

Unsupervised

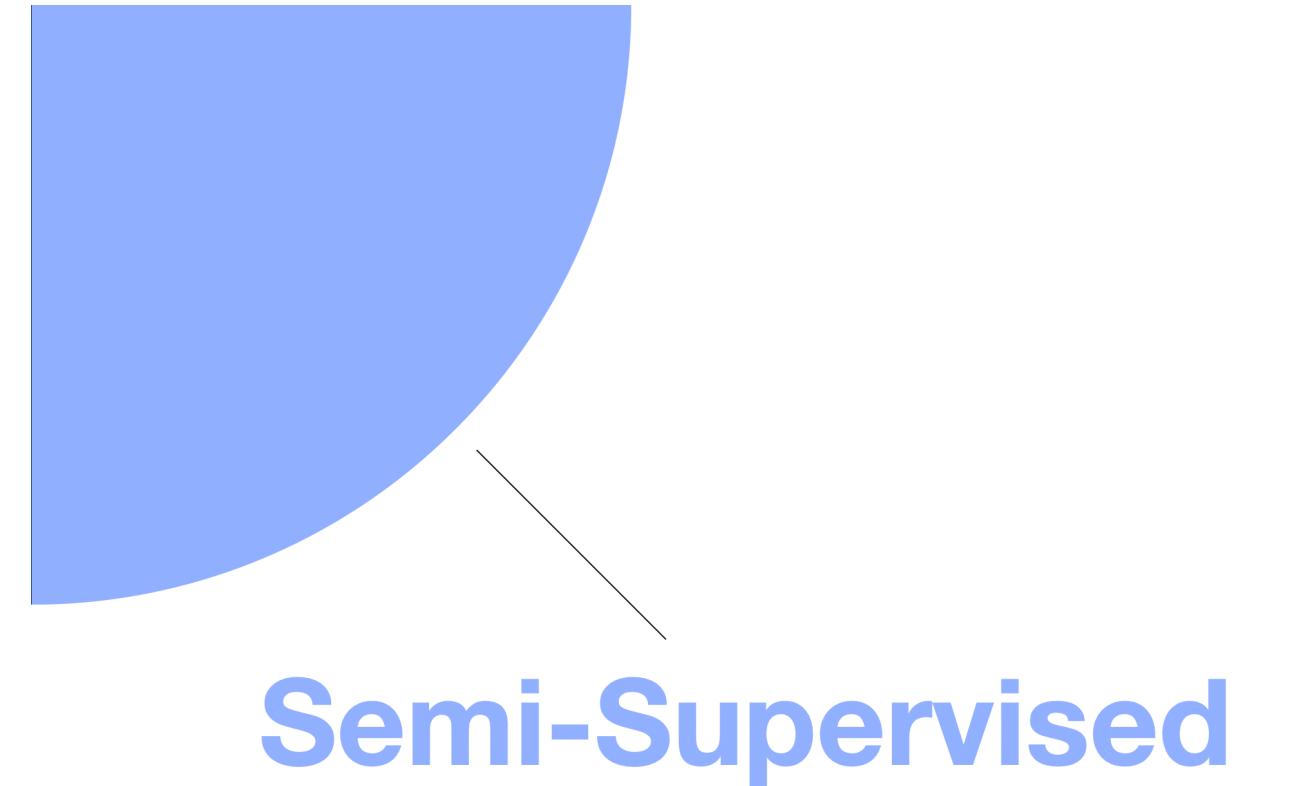
... clustering.

... dimensionality reduction.

... **generative modeling**.

Lecture 3-5: Generative Models

Semi-Supervised Learning



Goal: Combine small labeled and large unlabeled datasets.

$$\mathcal{L} = \mathcal{L}_{\text{supervised}} + \lambda \cdot \mathcal{L}_{\text{unsupervised}}$$

$\mathcal{L}_{\text{supervised}}$: Standard loss on labeled data

$\mathcal{L}_{\text{unsupervised}}$: **Consistency/pseudo-label loss** on unlabeled data

Key Idea: Adding unlabeled data does more than just increase dataset size.

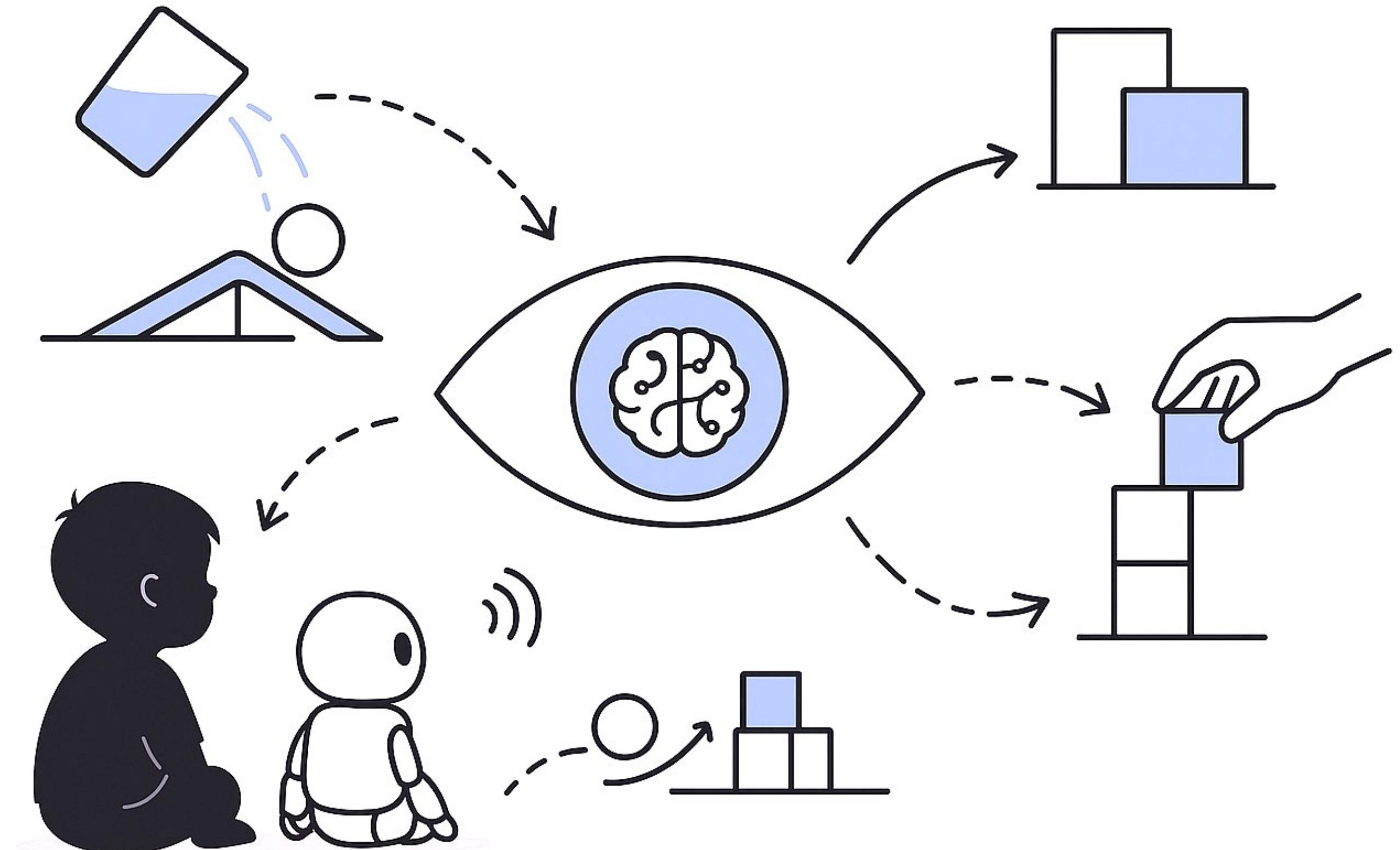
The unsupervised learning component reveals the underlying structure of the data distribution, enabling more refined decision boundaries than labeled data alone could provide.

- Key Assumptions:**
1. **Smoothness:** Similar inputs should have similar outputs
 2. **Cluster:** Data points in same cluster likely share same label
 3. **Manifold:** Data lies on lower-dimensional manifold in high-dimensional space

Self-Supervised Learning

Self-Supervised

“ The dark matter of intelligence. – LeCun and Misra (2021)



As babies, we learn how the world works largely by observation.

Later in life, we observe the world, act on it, observe again, and build hypotheses to explain how our actions change our environment by trial and error.

Generalized knowledge about the world, or common sense, forms the bulk of biological intelligence in both humans and animals.

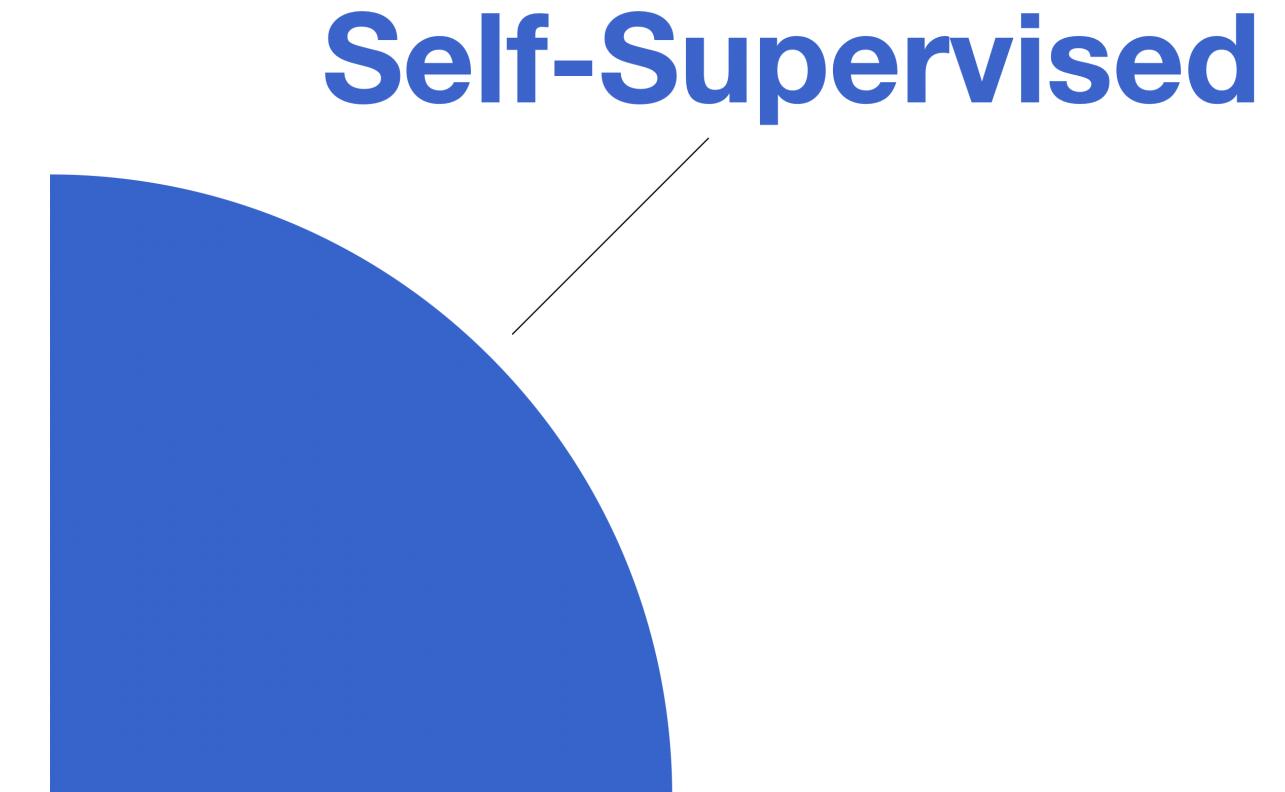
Common sense is the dark matter of artificial intelligence.

Self-Supervised Learning

... is the basic training principle of foundation models.



Thomas G. Dietterich
@tdietterich



I propose that we adopt the term "Large Self-Supervised Models (LSSMs)" as a replacement for "Foundation Models" and "LLMs". "LLMs" don't capture non-linguistic data and "Foundation Models" is too grandiose. Thoughts? [@percyliang](#)

2:58 am · 13 Aug 2022



Percy Liang ✅ @percyliang · 13 Aug 2022

The beauty of language is that you can have multiple terms that highlight different aspects of the same object. You don't have to choose. I use "LLM" to talk about LLMs, "self-supervised" for their construction, and "foundation model" for their function. No term can be replaced.

1

4

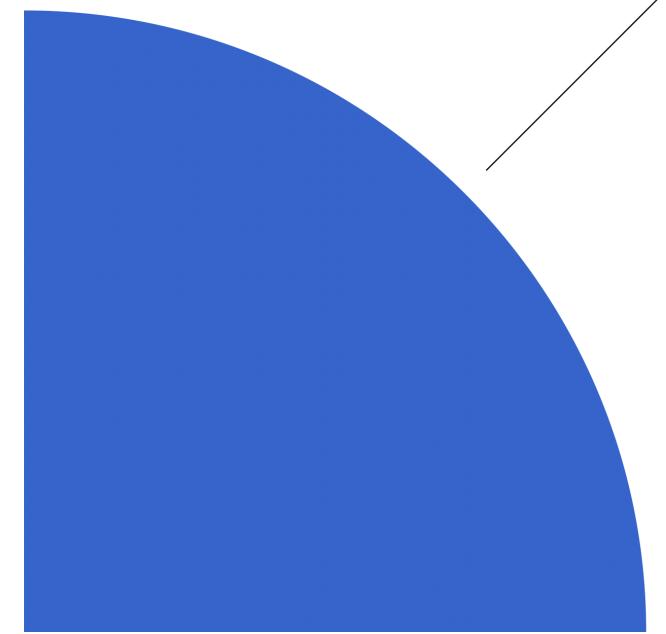
33



Self-Supervised Learning

Self-Supervised

Goal: Create surrogate tasks from raw data, no human labels.



Data: Unlabeled data (x_1, \dots, x_n) but we create supervision from the data itself.

Core Idea: Design **pretext tasks** that force the model to learn useful representations without labels y .

Key Insight: The data contains its own supervisory signal through structure, context, transformations.

How? Use unlabeled data for initial pretraining of a model to learn general features, which are then adapted to specific tasks through fine-tuning with a smaller labeled dataset.

Self-Supervised Learning

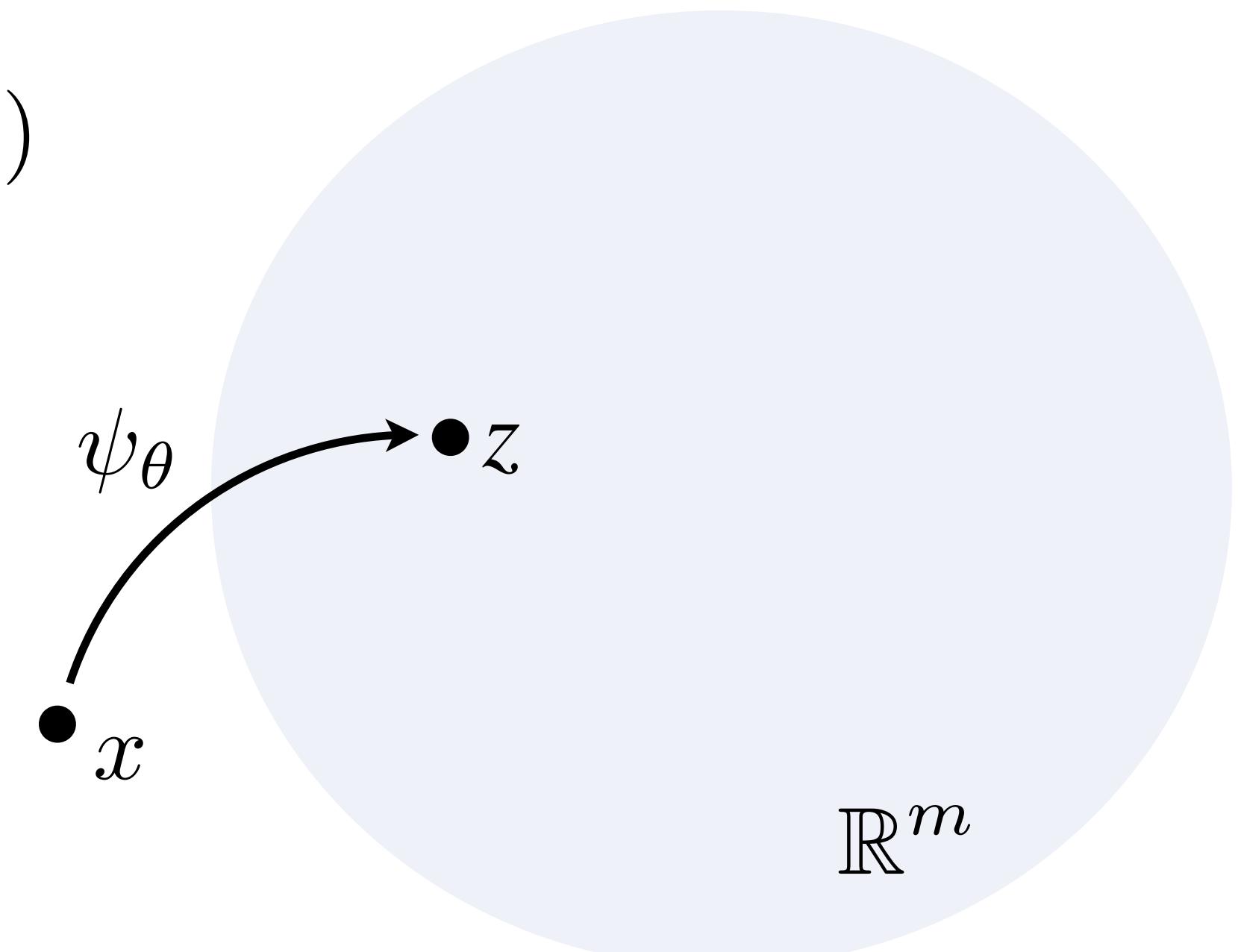
Self-Supervised

1. Pretraining Data (x_1, \dots, x_n)

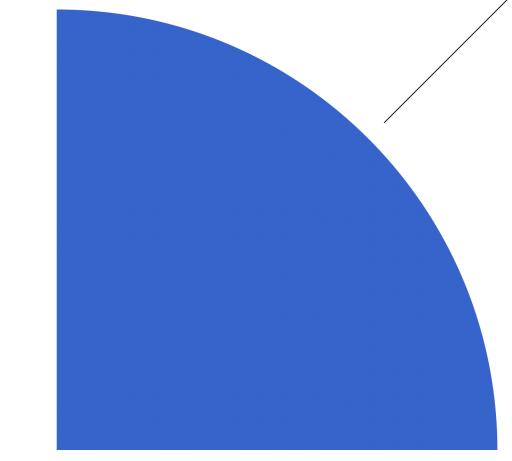
Model $\psi_\theta : x \mapsto \psi_\theta(x) \in \mathbb{R}^m$ ← *representation / embedding*

Optimize $\mathcal{L}_{\text{pre}}(\theta) = \frac{1}{n} \sum_{i=1}^n \ell_{\text{pre}}(\psi_\theta, x_i)$

to obtain $\hat{\theta}$
and pretrained model $\psi_{\hat{\theta}}$



Self-Supervised Learning



2. Evaluation

Data $((x_1^{\text{task}}, y_1^{\text{task}}), \dots, (x_{n_t}^{\text{task}}, y_{n_t}^{\text{task}}))$

n_t : number of downstream task samples

$n_t = 0$ **zero-shot learning** Model performs new tasks without any task-specific examples, using only pretrained representations

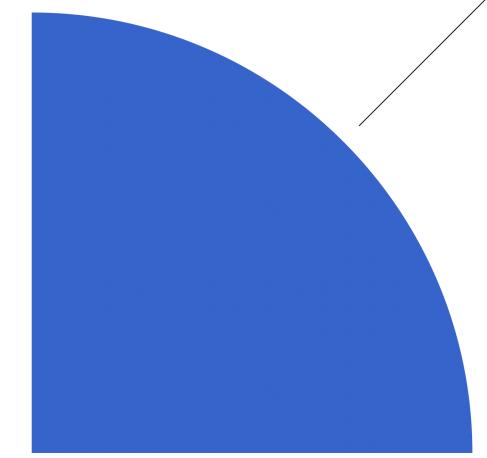
$n_t = \text{few}$ **few-shot learning** Model adapts to new tasks with just a handful of labeled examples (typically 1-10 per class)

→ **Tests true generalization:**

Can the model understand and solve novel tasks beyond its training distribution?

Lecture 10: Emergent Behaviors

Self-Supervised Learning



2. Evaluation Data $((x_1^{\text{task}}, y_1^{\text{task}}), \dots, (x_{n_t}^{\text{task}}, y_{n_t}^{\text{task}}))$

Downstream task via **linear probing / linear head**

prediction model $w^\top \psi_{\hat{\theta}}(x), \quad w \in \mathbb{R}^m$

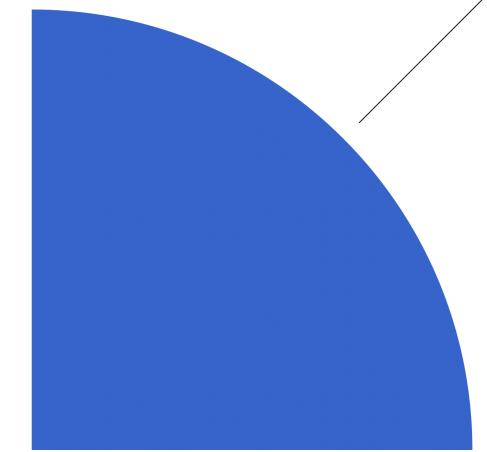
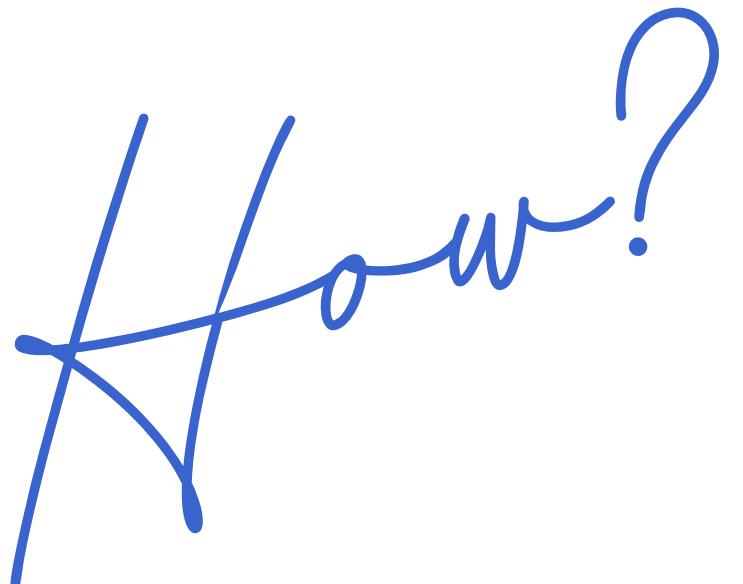
→ train w via $\min_w \frac{1}{n_t} \sum_{i=1}^{n_t} \ell_{\text{task}}(y_i^{\text{task}}, w^\top \psi_{\hat{\theta}}(x_i^{\text{task}}))$

$\hat{\theta} \star$: *foundation model parameters are frozen*

3. Finetuning Optimize both w and θ on downstream task
Initialize $\theta \leftarrow \hat{\theta}$ and w randomly

Lecture 5-9: Architectures

Self-Supervised Learning



1. Contrastive Learning

2. Masked Learning

3. Autoregressive Learning

4. Clustering-Based Methods

5. Distillation-Based Methods

Self-Supervised Learning: Contrastive Learning

Principle: Learn representations by **contrasting** similar (positive) and dissimilar (negative) pairs.

Defining Positive and Negative Pairs Without Labels



Positive Pairs: Different views of the same semantic content

- Created through augmentations that preserve meaning
- **Assumption:** these transformations do not change the underlying concept



Negative Pairs: Views from different instances

- Typically all other samples in the batch
- **Assumption:** random samples are likely semantically different

Self-Supervised Learning: Contrastive Learning

Defining Positive and Negative Pairs Without Labels

... in **vision**

+ Positive Pairs

original



color-jitter



reflection



rotation



- Negative Pairs

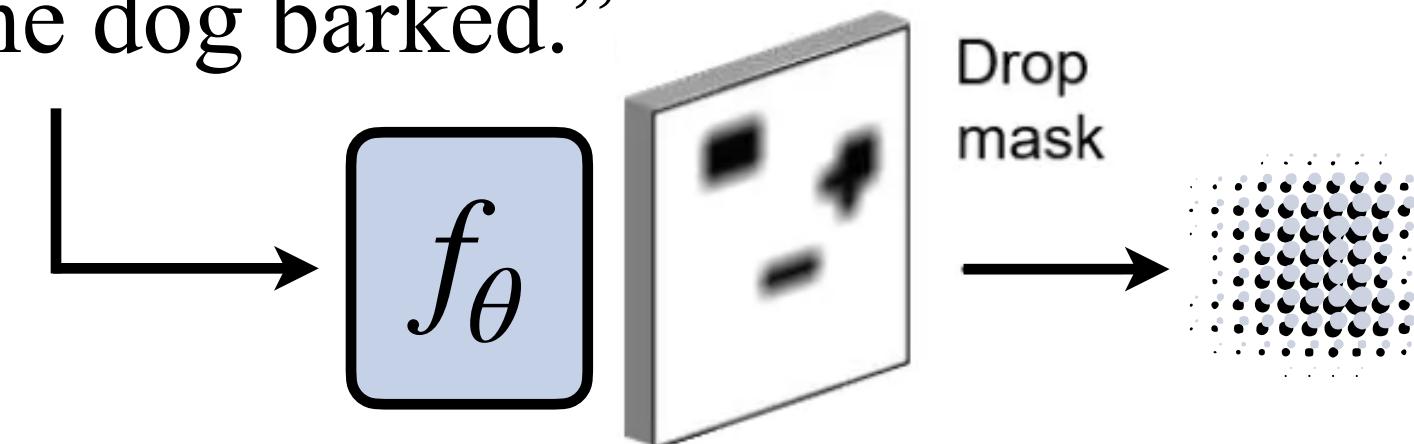
different image



Why? A cropped dog photo is still a dog; different photos are likely different objects.

... in **language**

“The dog barked.”



embedding of the same sentence
with different dropout masks

“The cat meowed.” different sentence

Why? “The dog barked” \approx “The dog woofed”
 \neq random other sentences

Self-Supervised Learning: Contrastive Learning

Self-Supervised

Evolution of Contrastive Learning Objectives

Early Approaches: Single Pairs

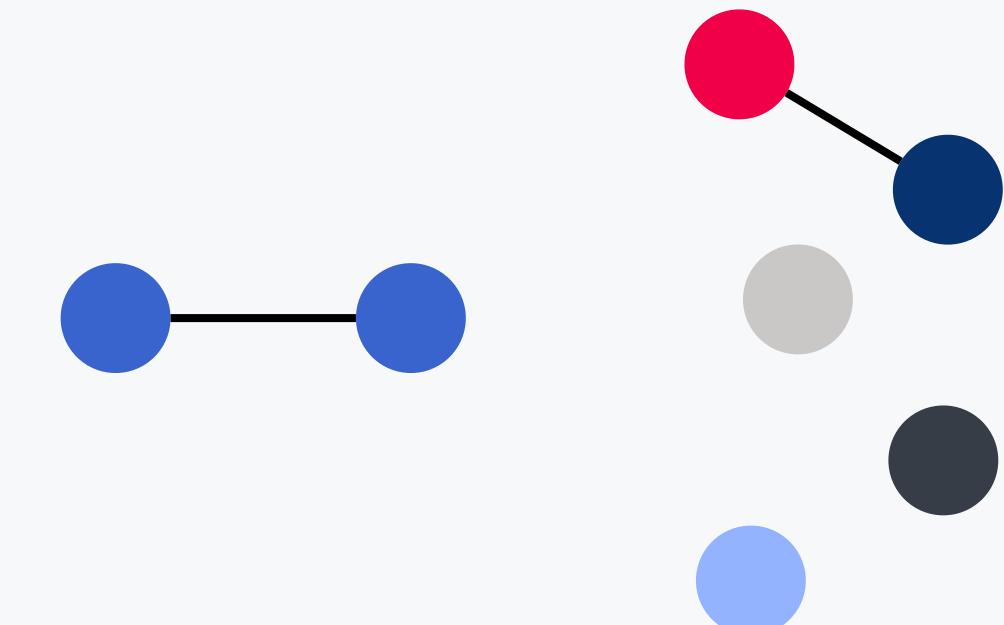
Contrastive Loss

$$\begin{aligned}\mathcal{L}_{\text{contr}}(\mathbf{x}_i, \mathbf{x}_j, \theta) = & \mathbb{1}[y_i = y_j] \|f_\theta(\mathbf{x}_i) - f_\theta(\mathbf{x}_j)\|_2^2 \\ & + \mathbb{1}[y_i \neq y_j] \max(0, \epsilon - \|f_\theta(\mathbf{x}_i) - f_\theta(\mathbf{x}_j)\|_2)^2\end{aligned}$$

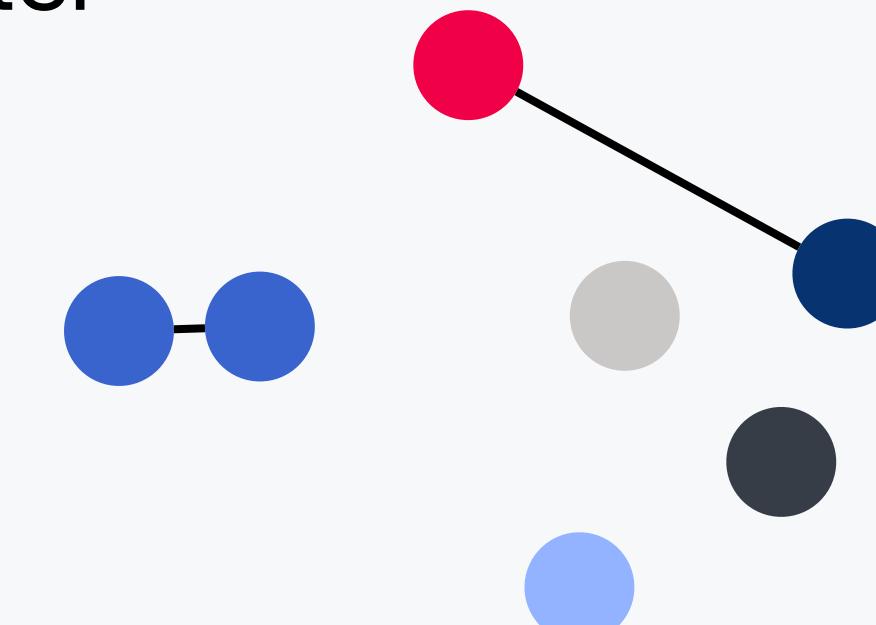
hyperparameter defining the lower bound distance
between samples of different classes

Chopra et al., (2005)

before



after



Self-Supervised Learning: Contrastive Learning

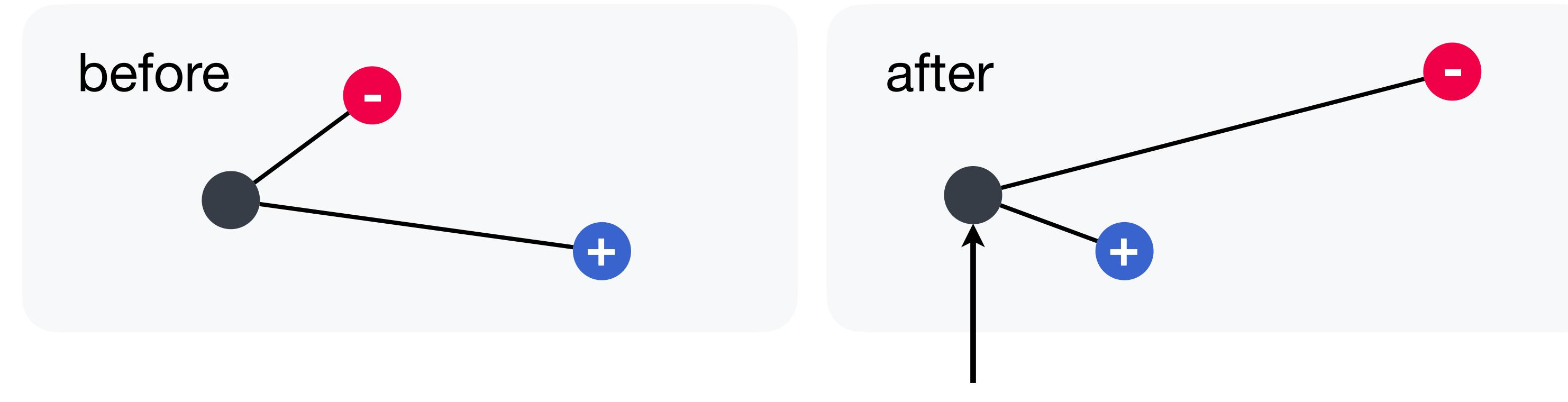
Triplet loss considers an anchor, a positive, and a negative sample simultaneously.
The contrastive loss works with pairs (positive or negative) independently.

→ Triplet loss is less greedy, as it considers relative distances between all three samples.

Triplet Loss

$$\mathcal{L}_{\text{triplet}} (\mathbf{x}, \mathbf{x}^+, \mathbf{x}^-) = \sum_{\mathbf{x} \in \mathcal{X}} \max \left(0, \|f(\mathbf{x}) - f(\mathbf{x}^+)\|_2^2 - \|f(\mathbf{x}) - f(\mathbf{x}^-)\|_2^2 + \epsilon \right)$$

margin parameter
set as minimum
offset between
distances of similar
vs. dissimilar pairs

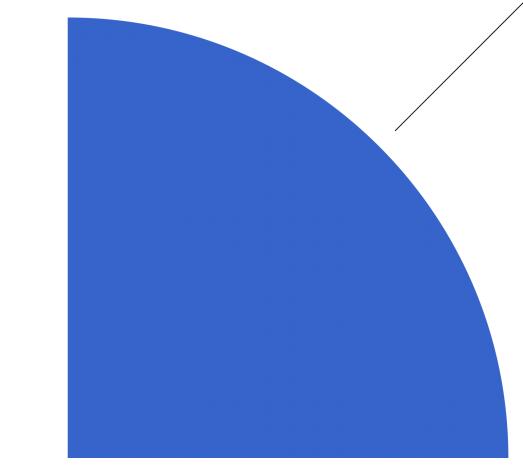


Schroff et al., (2015)

Self-Supervised Learning: Contrastive Learning

Modern Approaches: Multiple Pairs

Generalizes triplet loss to include comparison with multiple negative samples.

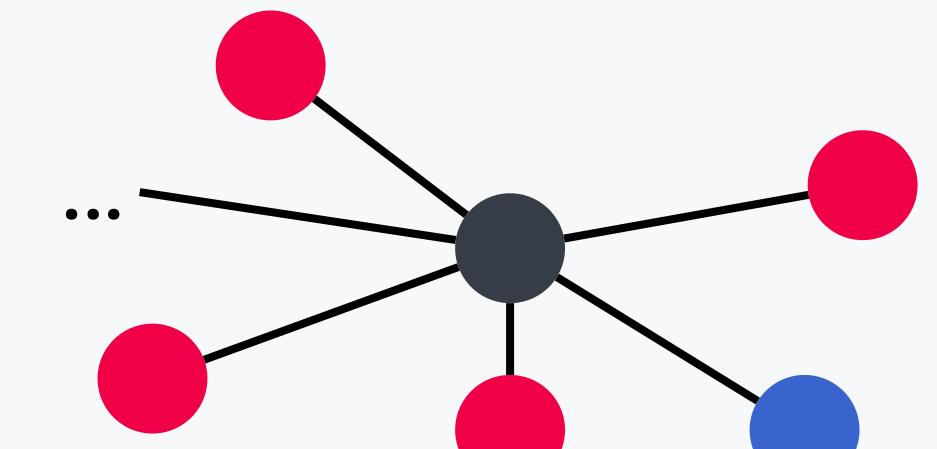


N-Pair Loss

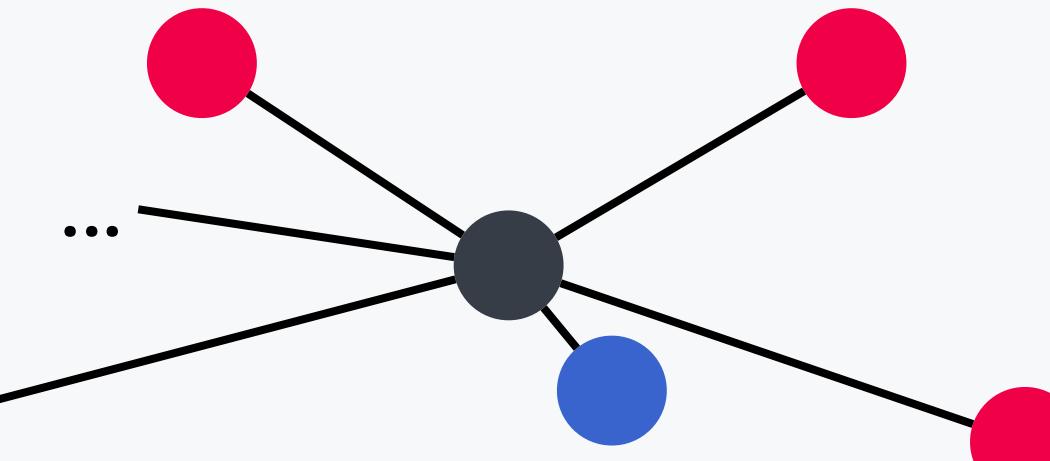
$$\begin{aligned}\mathcal{L}_{N\text{-pair}} \left(\mathbf{x}, \mathbf{x}^+, \left\{ \mathbf{x}_i^- \right\}_{i=1}^{N-1} \right) &= \log \left(1 + \sum_{i=1}^{N-1} \exp \left(f(\mathbf{x})^\top f(\mathbf{x}_i^-) - f(\mathbf{x})^\top f(\mathbf{x}^+) \right) \right) \\ &= -\log \frac{\exp \left(f(\mathbf{x})^\top f(\mathbf{x}^+) \right)}{\exp \left(f(\mathbf{x})^\top f(\mathbf{x}^+) \right) + \sum_{i=1}^{N-1} \exp \left(f(\mathbf{x})^\top f(\mathbf{x}_i^-) \right)}\end{aligned}$$

Sohn (2016)

before

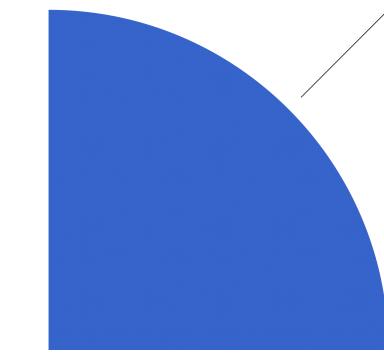


after



Note: If we only sample one negative sample per class, it is equivalent to the softmax loss for multi-class classification.

Self-Supervised Learning: Contrastive Learning



InfoNCE Loss

$$\mathcal{L}_{\text{InfoNCE}} = -\mathbb{E}_{c \sim p(c)} \mathbb{E}_{x \sim p(x|c)} \mathbb{E}_{\{x_i\} \sim p(x)} \left[\log \frac{\exp(f'(x)^\top f(c)/\tau)}{\exp(f'(x)^\top f(c)/\tau) + \sum_{i=1}^{N-1} \exp(f'(x_i)^\top f(c)/\tau)} \right]$$

Derivation: Noise Contrastive Estimation (NCE)

i.e., a method for estimating parameters of a statistical model

Idea: Instead of learning $p_{\text{data}}(x)$, it learns to distinguish data from noise ← **contrastive!**

How? Logistic regression with binary crossentropy loss

log-odds ratio

$$s_\theta(x) = \log \frac{p_\theta(x)}{p_n(x)}$$

NCE Loss

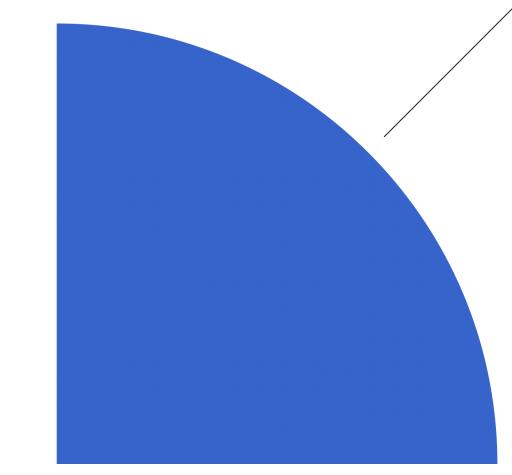
$$\mathcal{L}_{\text{NCE}} = -\mathbb{E}_{x \sim p_{\text{data}}} [\log \sigma(s_\theta(x))] - k \cdot \mathbb{E}_{x \sim p_n} [\log (1 - \sigma(s_\theta(x)))]$$

number of noise samples

per real sample

Self-Supervised Learning: Contrastive Learning

Generalize binary classification in NCE to N-way classification.
InfoNCE asks 'which one of these N samples is the positive one?'



InfoNCE Loss

$$\mathcal{L}_{\text{InfoNCE}} = -\mathbb{E}_{c \sim p(c)} \mathbb{E}_{x \sim p(x|c)} \mathbb{E}_{\{x_i\} \sim p(x)} \left[\log \frac{\exp(f'(x)^\top f(c)/\tau)}{\exp(f'(x)^\top f(c)/\tau) + \sum_{i=1}^{N-1} \exp(f'(x_i)^\top f(c)/\tau)} \right]$$

- encoders f, f'
- context c (anchor) and positive sample x (e.g., its augmentation)
- positive pair: (c, x) where $x \sim p(x | c)$ and negative samples: $\{x_i\}_{i=1}^{N-1}$ where $x_i \sim p(x)$
- score function: $s(x, c) = \exp(f'(x)^\top f(c)/\tau)$

van den Oord et al., (2018)

↑
temperature that controls how "sharp" or
"smooth" the similarity distribution is

Exercise 1 · Task 1

→ connection to
information theory

Self-Supervised Learning: Contrastive Learning

Self-Supervised

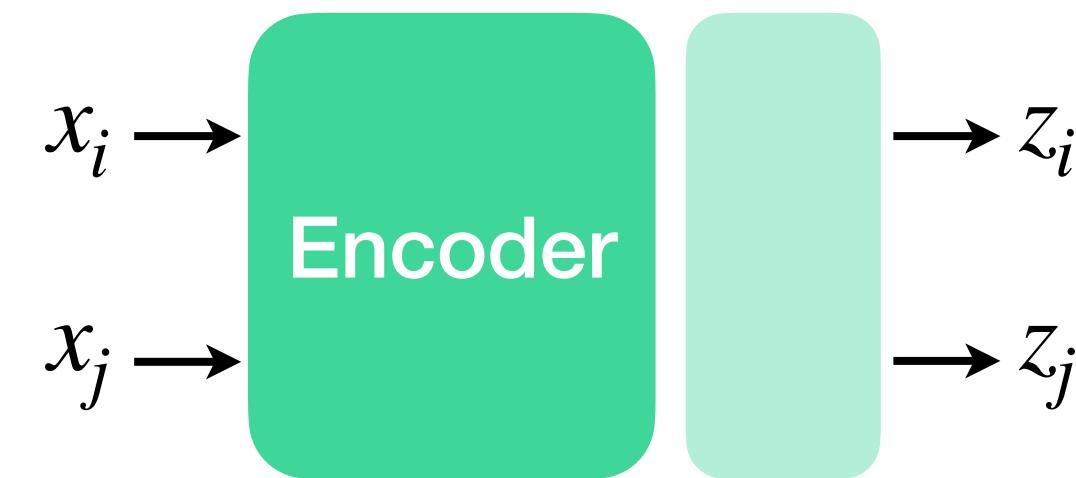
Example from Vision: SimCLR Framework

... framework for contrastive learning of visual representations by maximizing agreement between differently augmented views of the same image

 **Code Notebook 1 · Task A**

Strategy: Data Augmentation as Supervision

batch of N ,
data augmentation to $2N$
with $2(N - 1)$ negative examples

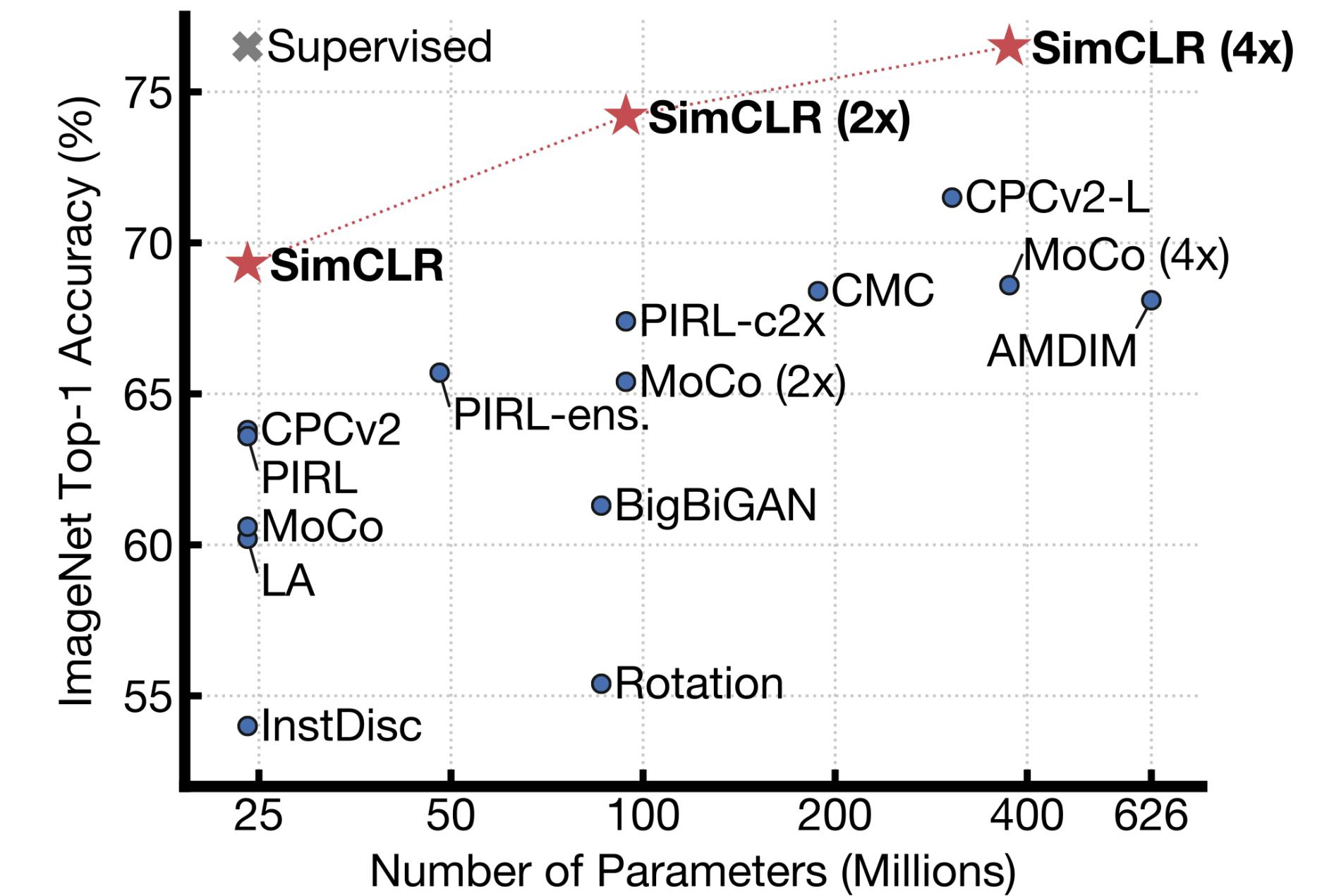


Loss: Symmetric version of InfoNCE

$$\mathcal{L}_{\text{SimCLR}}^{(i,j)} = - \log \frac{\exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_j) / \tau)}{\sum_{k=1}^{2N} \mathbb{1}_{[k \neq i]} \exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_k) / \tau)}$$

Chen et al., (2018)

Key to Success: Large Batch Sizes



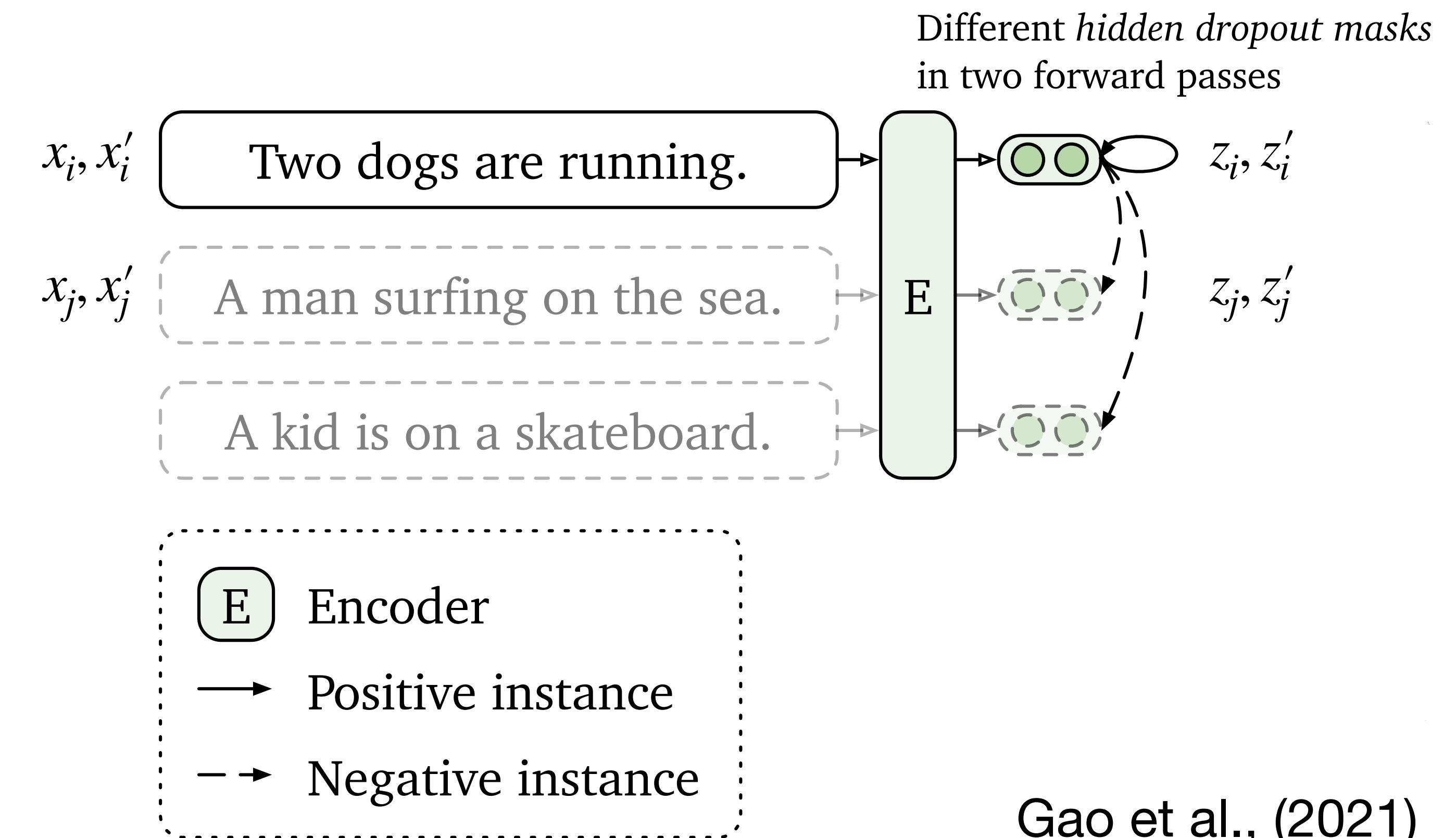
Self-Supervised Learning: Contrastive Learning

Example from Language: SimCSE Framework

... applies contrastive learning to create high-quality sentence embeddings using **dropout as minimal augmentation**

Loss: Symmetric version of InfoNCE

Here, z_i, z'_i is the positive pair.



Self-Supervised Learning: Masked Learning

Principle: Hide parts of the data and train the model to predict what is missing.

This forces the model to understand the structure and relationships within the data.

Data (x_1, x_2, \dots, x_n) and mask $M \subset \{1, \dots, m\}$.

Predict masked parts from unmasked parts $p(x_M | x_{\setminus M})$.

e.g., in language

Randomly masked A quick [MASK] fox jumps over the [MASK] dog

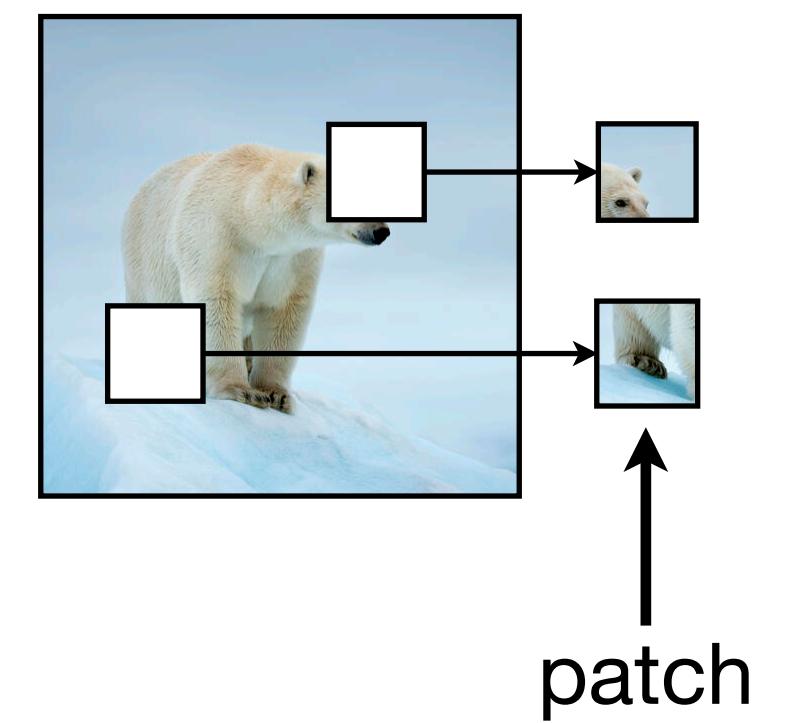
Predict A quick brown fox jumps over the lazy dog

special token

Design choices between *random structured adversarial curriculum*
masking

simple and unbiased
mask semantically meaningful units
mask the most informative parts
start with easy masking, increase difficulty
→ 15-75% depending on domain.

e.g., in vision



Self-Supervised Learning: Masked Learning

General Formulation

Masking Loss

$$\mathcal{L}_{\text{mask}} = \mathbb{E}_{x \sim \mathcal{D}} \mathbb{E}_M \left[\frac{1}{|M|} \sum_{i \in M} \ell(x_i, \hat{x}_i) \right]$$

position-wise loss function
dependent on domain

$\hat{x}_i = f_\theta(x_{\setminus M})_i$ is the model's prediction based on
unmasked parts for position i

For Language:

$$\mathcal{L}_{\text{MLM}} = -\mathbb{E}_{x \sim \mathcal{D}} \mathbb{E}_M \left[\frac{1}{|M|} \sum_{i \in M} \log p_\theta(x_i | x_{\setminus M}) \right]$$

cross-entropy between true distribution δ_{x_i} (one-hot at true token)
and predicted distribution $p_\theta(\cdot | x_{\setminus M})$

For Vision:

$$\mathcal{L}_{\text{MAE}} = \mathbb{E}_{x \sim \mathcal{D}} \mathbb{E}_M \left[\frac{1}{|M|} \sum_{i \in M} \|x_i - f_\theta(x_{\setminus M})_i\|^2 \right]$$

pixel/patch
reconstruction

Self-Supervised Learning: Masked Learning

Self-Supervised

Example from Vision:
Masked Autoencoder (MAE)

Lecture 6: Vision Foundation Models

Example from Language:
Bidirectional Encoder Representations from Transformers (BERT)

Lecture 6: Language Foundation Models

Self-Supervised Learning: Autoregressive Learning

Principle: Learn by **predicting next** elements in sequence.

Autoregressive Loss

$$\mathcal{L}_{\text{AR}} = -\mathbb{E}_{x \sim \mathcal{D}} \left[\sum_{t=1}^T \log p_{\theta}(x_t | x_{<t}) \right]$$

Any joint distribution can be factored as a product of conditionals:

chain rule of probability

$$\frac{\text{prefix/history}}{x_{<t} = (x_1, \dots, x_{t-1})}$$

$p(x_1, \dots, x_T) = \prod_{t=1}^T p(x_t | x_{<t})$

- Characteristics:**
1. **Exact factorization**, i.e., no approximation, allows to objectively compare models, detect anomalies, etc.
 2. **Causal structure** that enforces temporal/sequential dependencies, i.e., foundation to generate sequences, one token at a time, left to right.
 3. Deep **connection to information theory**, i.e., ideal AR model achieves Shannon-optimal compression.

Excursion: Perplexity

Perplexity (PPL) is the exponential of the average negative log-likelihood per token:

Perplexity

$$\text{PPL} = \exp \left(-\frac{1}{T} \sum_{t=1}^T \log p(x_t | x_{<t}) \right) = p(x_1, x_2, \dots, x_N)^{-1/T}$$

Intuition: How "surprised" the model is on average

- PPL = 10 → Model is as confused as choosing uniformly from 10 options
- PPL = 1 → Perfect prediction (no surprise)

Historical Origin:

- **Information Theory** (1948): Claude Shannon introduced as $2^{H(X)}$ (exponential of entropy H)
- **Language Modeling** (1970s-80s): Frederick Jelinek and IBM team popularized for speech recognition
- **Modern NLP**: Standard metric since statistical language models replaced rule-based systems

Excursion: Pseudo-Likelihood

Autoregressive strategies directly learn the joint probability through the chain rule, allowing to compute the exact likelihood *and* perplexity → direct way to evaluate models!

What probability distribution do approaches based on masking learn?

→ Masked language models (MLM) optimize something called the **pseudo-likelihood**!

What is the difference? True joint probability: $p(x_1, x_2, x_3) = p(x_1) \cdot p(x_2 | x_1) \cdot p(x_3 | x_1, x_2)$

Pseudo joint probability: $p_{\text{pseudo}}(x_1, x_2, x_3) = p(x_1 | x_2, x_3) \cdot p(x_2 | x_1, x_3) \cdot p(x_3 | x_1, x_2)$

Pseudo-Likelihood

$$p_{\text{pseudo}}(x_1, \dots, x_T) = \prod_{t=1}^T p(x_t | x_{\setminus t})$$

→ approximation to the joint probability distribution

Excursion: Pseudo-Perplexity

 Exercise 1 · Task 2

Masked language models learn the pseudo-likelihood

Negative Pseudo-Log-Likelihood

$$\mathcal{L}_{\text{MLM}} = q \cdot \text{NPLL}(\theta) = -q \cdot \mathbb{E}_{x \sim \mathcal{D}} \left[\sum_{t=1}^T \log p_\theta(x_t | x_{\setminus t}) \right]$$

Following the same structure as perplexity:

Pseudo-Perplexity

$$\text{Pseudo-PPL} = \exp \left(-\frac{1}{T} \sum_{t=1}^T \log p(x_t | x_{\setminus t}) \right)$$

Note: These metrics are not comparable to each other despite similar names and formulas.

Remark: For continuous data (e.g., image patches), pseudo-likelihood extends by assuming Gaussian conditionals $p(x_t | x_{\setminus t}) = \mathcal{N}(f_\theta(x_{\setminus t}), \sigma^2)$, making MSE minimization equivalent to maximizing pseudo-likelihood.

Self-Supervised Learning: Autoregressive Learning

Self-Supervised

Example from Language:
Generative Pretrained Transformer (GPT)

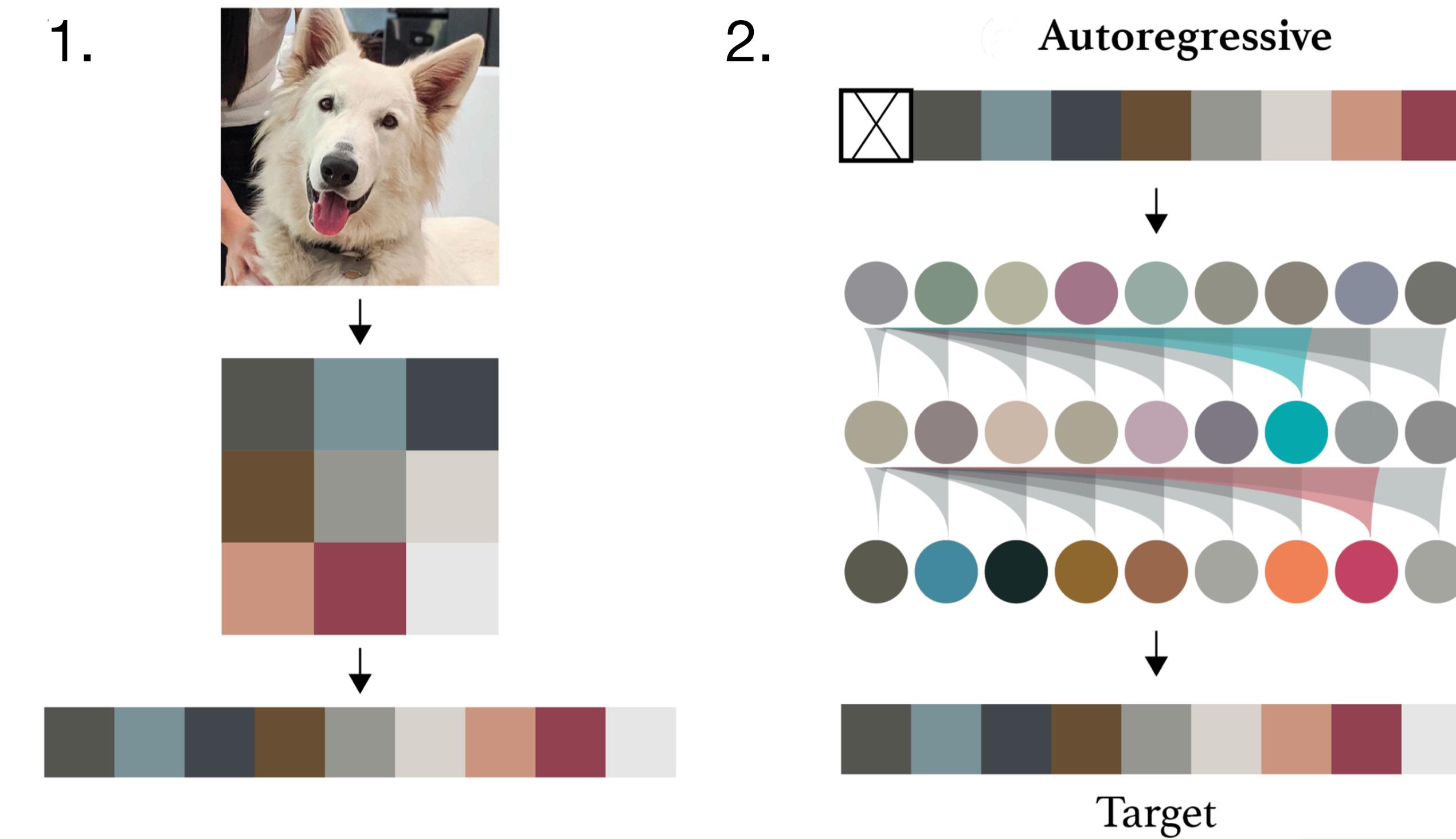
Radford et al., (2018)

Lecture 6: Language Foundation Models

Example from Vision:
Image GPT (iGPT)

Chen et al., (2020)

1. Pre-process raw images by resizing to a low resolution and reshaping into a 1D sequence.
2. Train via auto-regressive next pixel prediction.



Self-Supervised Learning: Fusions of Learning Principles

Self-Supervised

Combining multiple self-supervised learning strategies leverages their complementary strengths to learn richer, more versatile representations than any single approach alone.

e.g., **ELECTRA**

Combines: Masked Learning + Contrastive Learning

How: Generator creates replacements for masked tokens, discriminator distinguishes real vs. replaced tokens.

Benefit: More sample-efficient than pure masked learning.

Clark et al., (2020)

Masked language models:
learn from 15% of tokens

A quick [MASK] fox jumps over the [MASK] dog



A quick brown fox jumps over the lazy dog

original or
replacement?

ELECTRA:
learns from 100% of tokens

A quick [MASK] fox jumps over the [MASK] dog

etc.



A quick brown fox jumps over the lazy dog

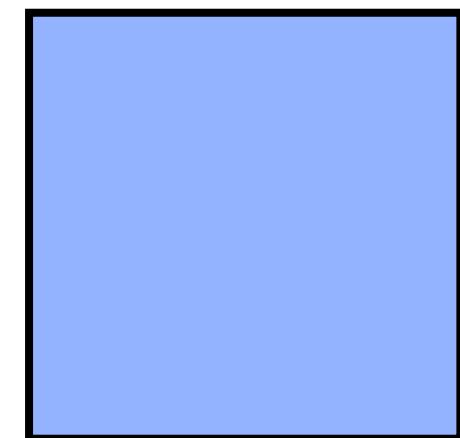
From Learning Principle to Architecture

Important!

The **self-supervision objective** is *not* a training detail but the fundamental constraint that **shapes many architectural choices** from attention patterns to encoder-decoder asymmetry.

1. Contrastive → Must compare → Needs batch/memory

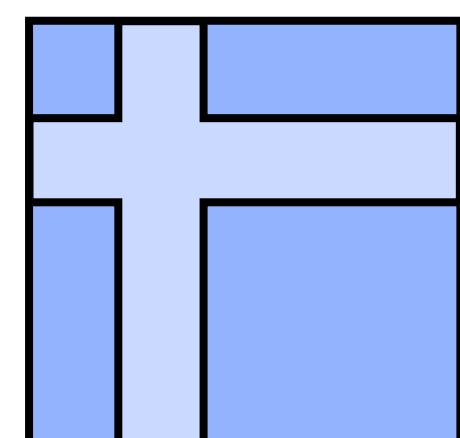
- Requires denominator with many negatives: $\exp(\text{sim}^+)/\sum \exp(\text{sim})$
- **Solution:** Large batches (SimCLR: 4096) or memory banks (MoCo: 65K queue)



full
attention

2. Masked → Must see all → Needs bidirectional

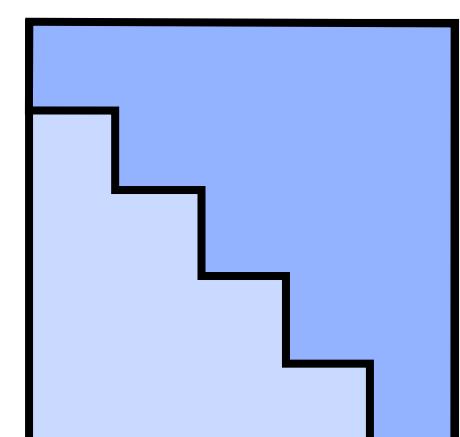
- Predict from full context: $p(x_{\text{mask}} | x_{\text{left}}, x_{\text{right}})$
- **Solution:** Full attention patterns



bidirectional
masked
attention

3. Autoregressive → Must generate → Needs causality

- Factorization chain: $p(x) = \prod p(x_t | x_{<t})$
- **Solution:** Causal self-attention to avoid future leakage



causal
attention

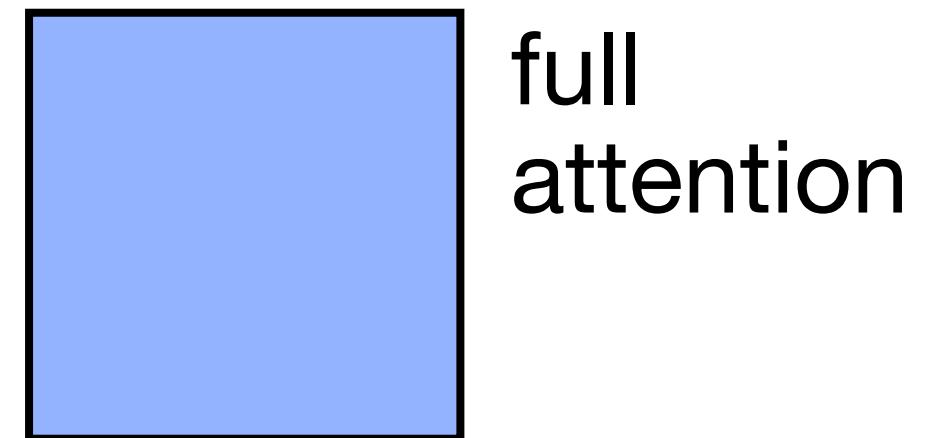
From Learning Principle to Architecture

Important!

The **self-supervision objective** is *not* a training detail but the fundamental constraint that **shapes many architectural choices** from attention patterns to encoder-decoder asymmetry.

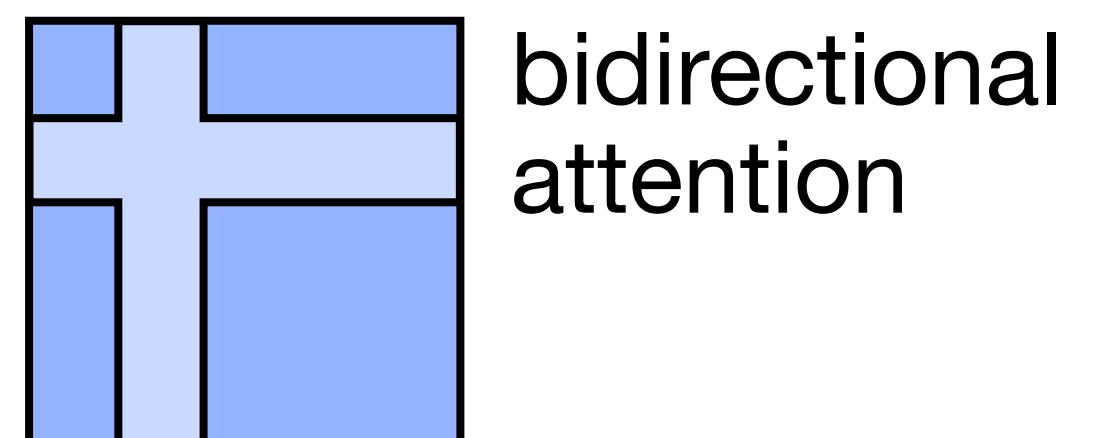
1. Contrastive

```
python
1 attn = (Q @ K.T) / sqrt(d_k) # no mask needed
2 attn = softmax(attn)          # full attention
3 out = attn @ V              # then global pool
```



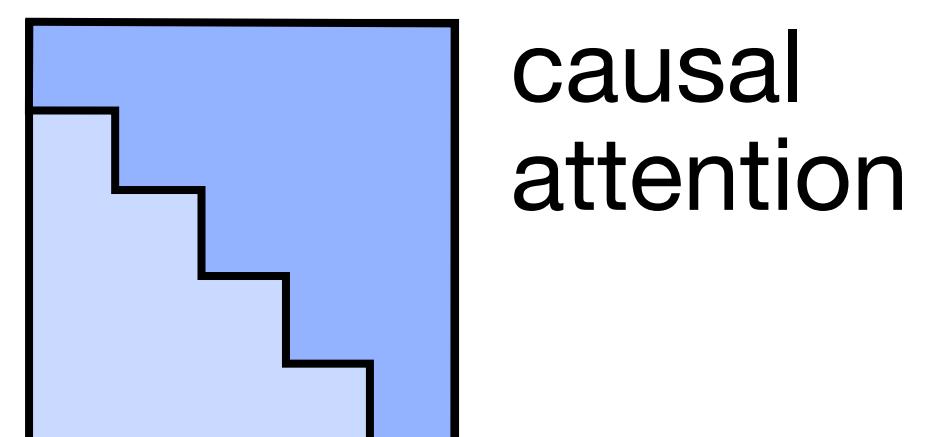
2. Masked

```
python
1 attn = Q @ K.T           # full attention
2 attn[mask_idx, :] = -inf # mask positions cannot attend
3 attn[:, mask_idx] = -inf # nothing attends to masked
```



3. Autoregressive

```
python
1 mask = torch.tril(torch.ones(n, n))
2 attn = (Q @ K.T).masked_fill(mask == 0, -inf)
```

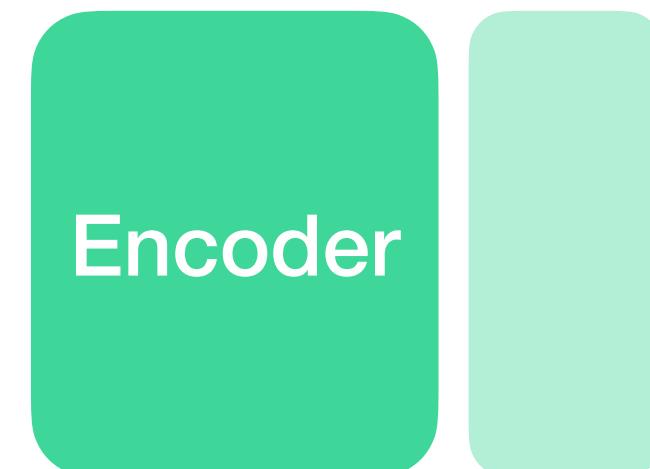


From Learning Principle to Architecture

Important!

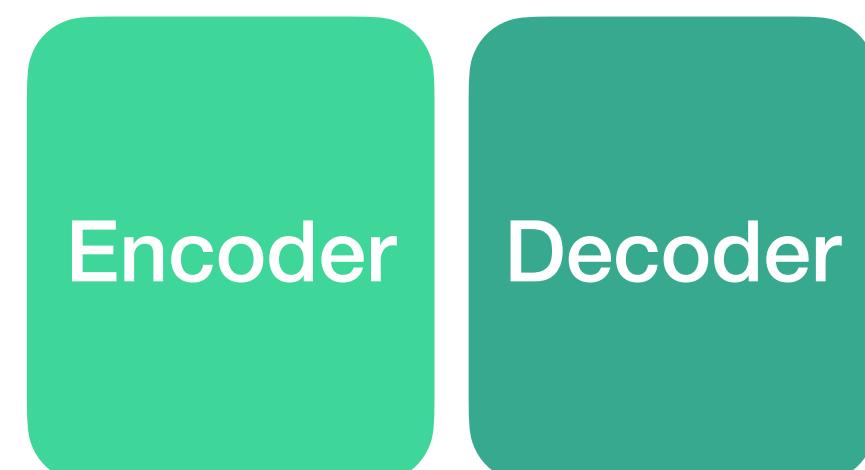
The **self-supervision objective** is *not* a training detail but the fundamental constraint that **shapes many architectural choices** from attention patterns to encoder-decoder asymmetry.

1. Contrastive



Encoder + projection head
for comparing global features

2. Masked



Encoder-decoder split that
can be asymmetric

3. Autoregressive



Unified decoder-only architecture
with causal masking



Lecture 5-7: Architectures

Model Size Evolution of GPT

+ sparse attention in alternating layers

•

GPT-1

117M

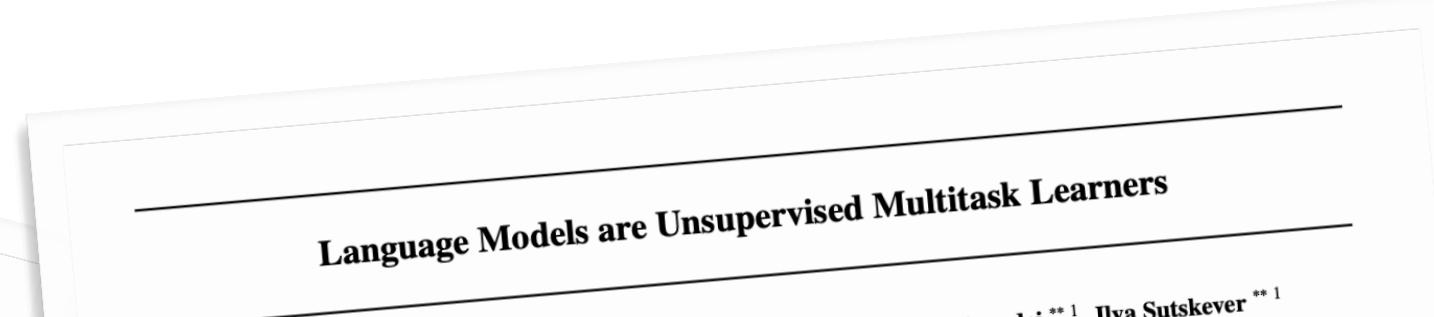
(2018)



GPT-2

1.5B

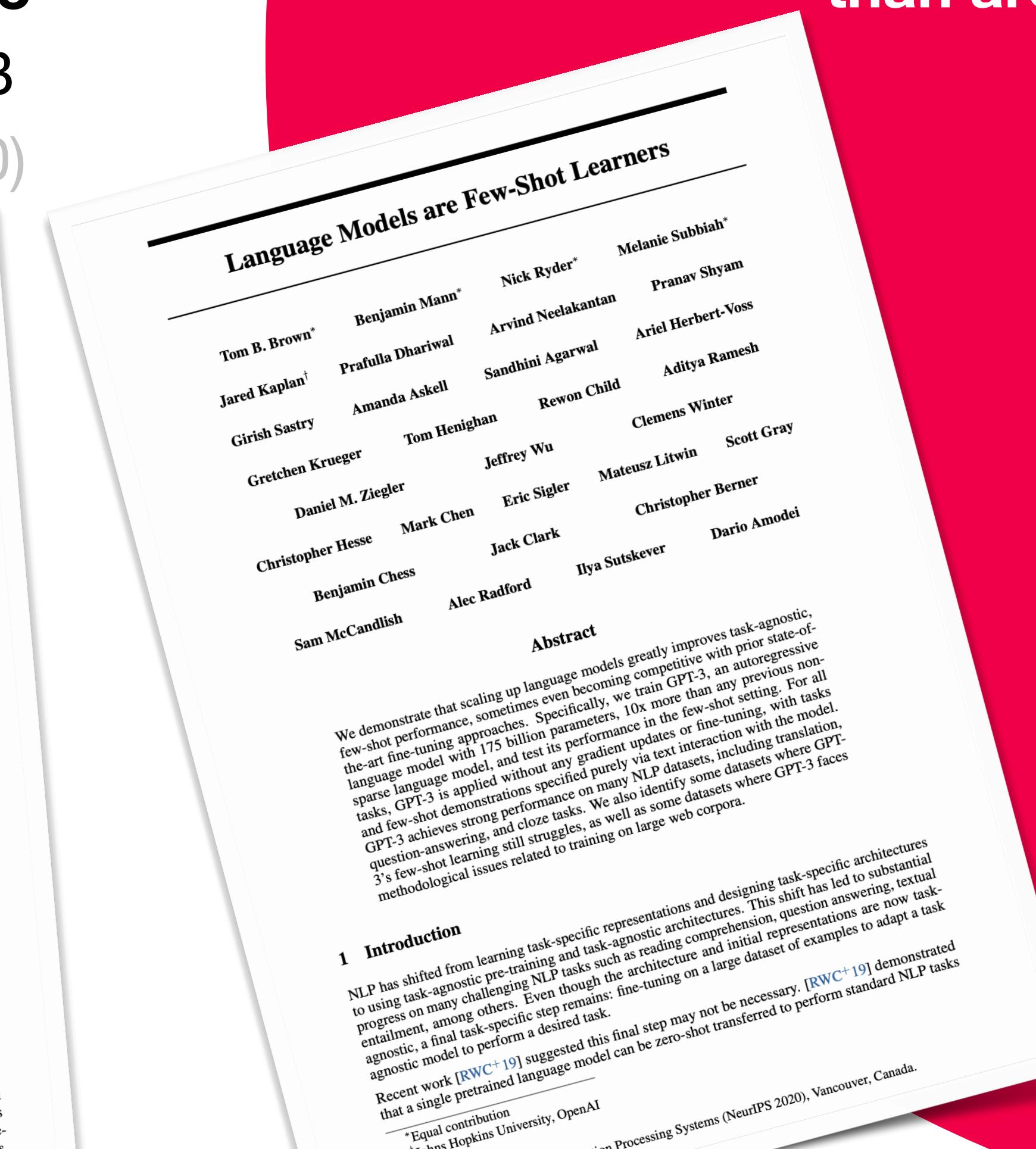
(2019)



GPT-3

175B

(2020)



→ architecture changed minimally

→ scale can be more important than architecture

and Genera

1 Introduction

The ability to learn effectively from raw text is crucial to alleviating the dependence on supervised learning in natural language processing (NLP). Most deep learning methods require substantial amounts of manually labeled data, which restricts their applicability in many domains that suffer from a dearth of annotated resources [6]. In these situations, models that can leverage linguistic information from unlabeled data provide a valuable alternative to gathering more annotation, which can be time-consuming and expensive. Further, even in cases where considerable supervision is available, learning good representations in an unsupervised fashion can provide a significant performance boost. The most compelling evidence for this so far has been the extensive use of pre-trained word embeddings [10][39][42] to improve performance on a range of NLP tasks [8][11][26][45]. Leverage more than word-level information from unlabeled text, however, is challenging for two main reasons. First, it is unclear what type of optimization objectives are most effective at learning text representations that are useful for transfer. Recent research has looked at various methods outperforming the others learned representations to the best of our knowledge. A creative way to transfer these learned representations to the best of our knowledge is to combine the use of task-specific changes to the training objective with fine-tuning schemes [21] and adding auxiliary components to the model to encourage it to learn features that are useful for specific tasks [12].

1. Introduction

Machine learning systems now excel (in expectation) at tasks they are trained for by using a combination of large datasets, high-capacity models, and supervised learning (Krizhevsky et al., 2012) (Sutskever et al., 2014) (Amodei et al., 2016). Yet these systems are brittle and sensitive to slight changes in the data distribution (Recht et al., 2018) and overfitting (Kirkpatrick et al., 2017). Current systems are trained to learn a single task at a time rather than multiple tasks simultaneously.

Multitask learning (Caruana, 1997) is a promising framework for improving general performance. However, multitask training in NLP is still nascent. Recent work reports modest performance improvements (Yogatama et al., 2019) and the two most ambitious efforts to date have trained on a total of 10 and 17 (dataset, objective) pairs respectively (McCann et al., 2018) (Bowman et al., 2018). From a meta-learning perspective, each (dataset, objective) pair is a single training example sampled from the distribution of datasets and objectives. Current ML systems need hundreds to thousands of examples to induce functions which generalize well. This suggests that multitask training may need just as many effective training pairs to realize its promise with current approaches. It will be very difficult to continue to scale the creation of datasets and the design of objectives to the degree that may be required to achieve state-of-the-art performance. This shift has led to substantial progress on many challenging NLP tasks such as reading comprehension, question answering, textual entailment, among others. Even though the architecture and initial representations are now task-agnostic, a final task-specific step remains: fine-tuning on a large dataset of examples to adapt a task-agnostic model to perform a desired task.

1 Introduction

NLP has shifted from learning task-specific representations and designing task-specific architectures to using task-agnostic pre-training and task-agnostic architectures. This shift has led to substantial progress on many challenging NLP tasks such as reading comprehension, question answering, textual entailment, among others. Even though the architecture and initial representations are now task-agnostic, a final task-specific step remains: fine-tuning on a large dataset of examples to adapt a task-agnostic model to perform a desired task. Recent work [RWC⁺19] suggested this final step may not be necessary. [RWC⁺19] demonstrated that a single pretrained language model can be zero-shot transferred to perform standard NLP tasks

*Equal contribution

[†]Johns Hopkins University, OpenAI

Neural Information Processing Systems (NeurIPS 2020), Vancouver, Canada.

The Scale Breakthrough

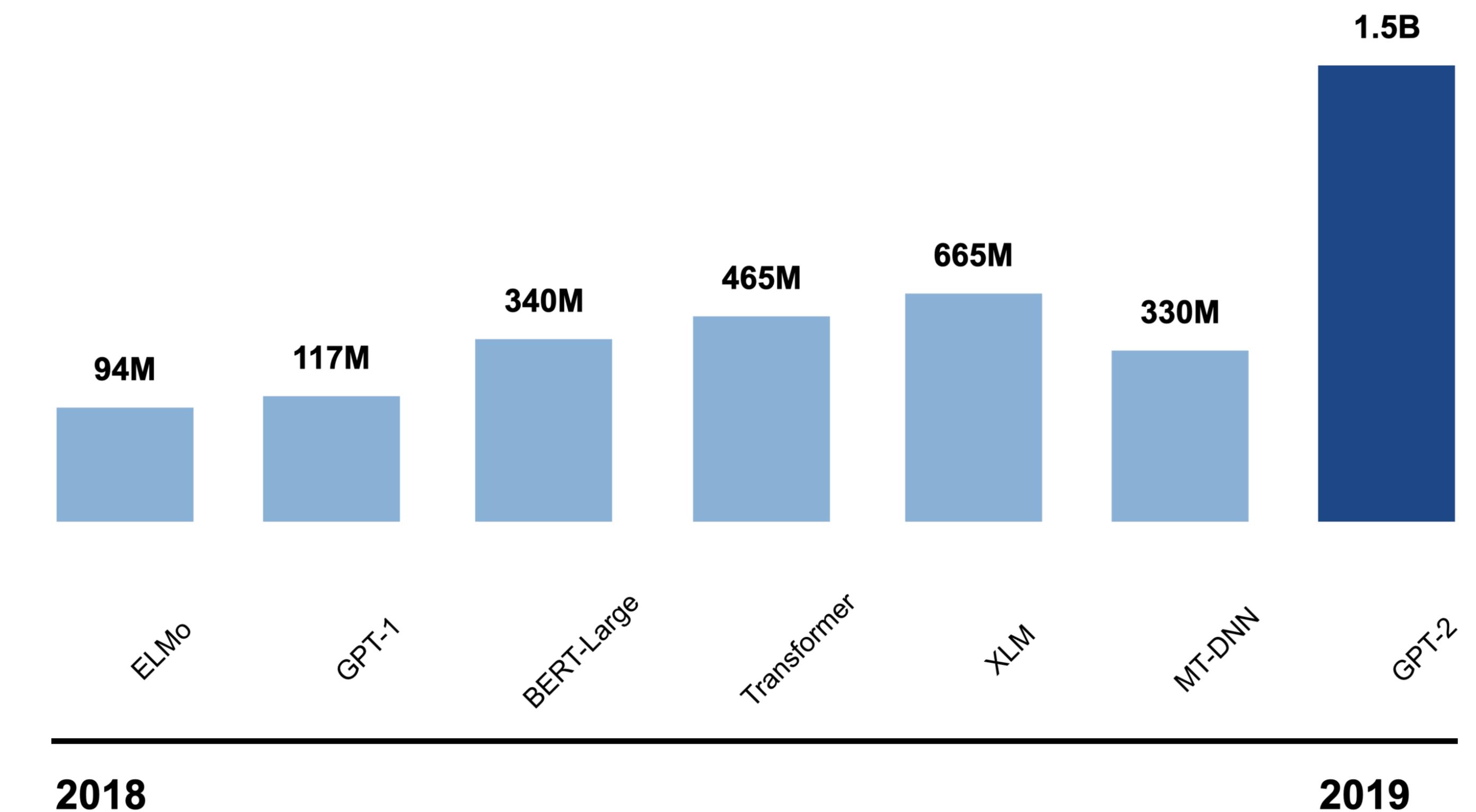
GPT-2's 1.5B parameters in 2019 marked a dramatic **2-3x leap in model scale** compared to previous architectures, signaling the start of the large language model era.

How much data do we need?

How big should our model be?

How much compute should we use?

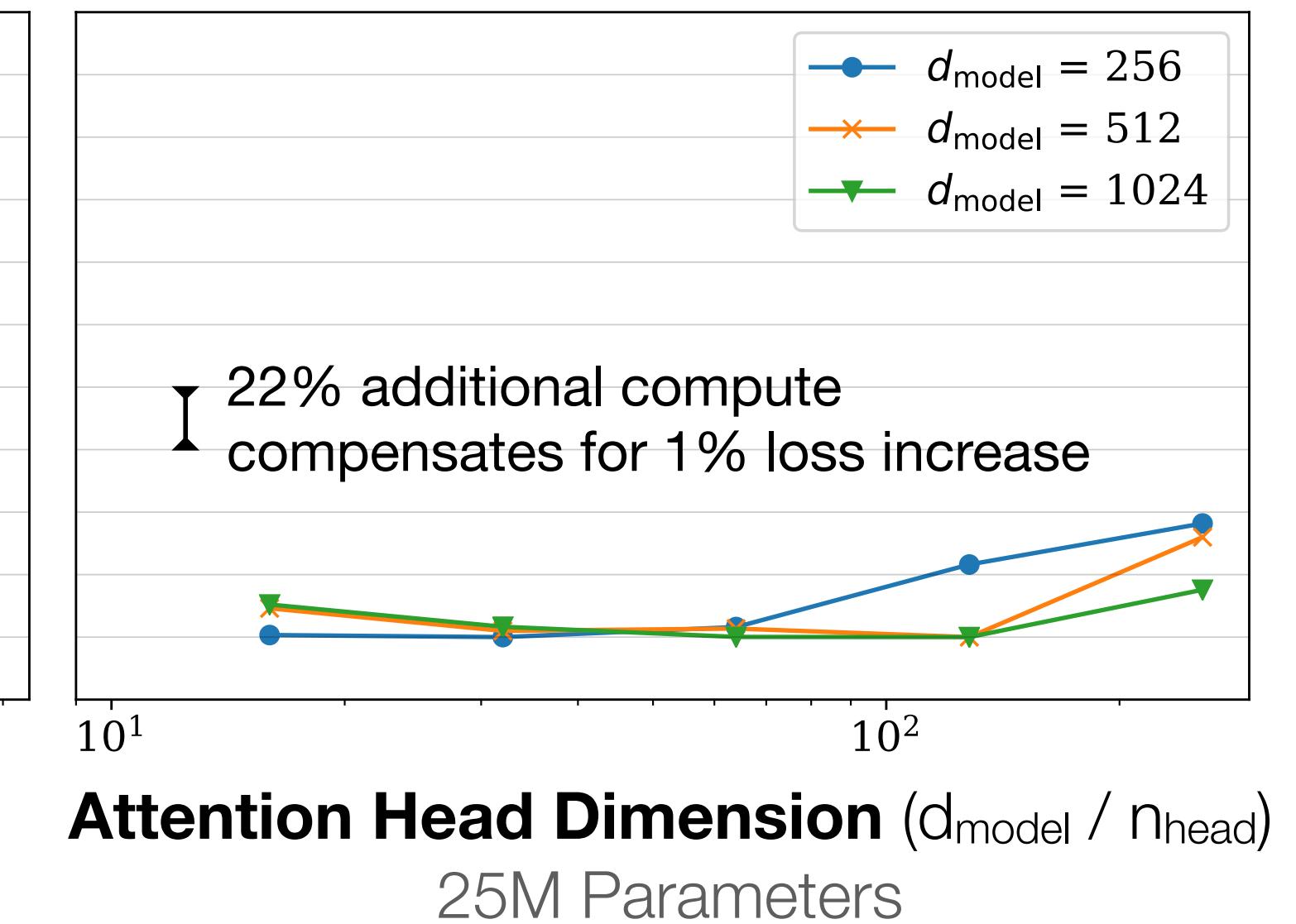
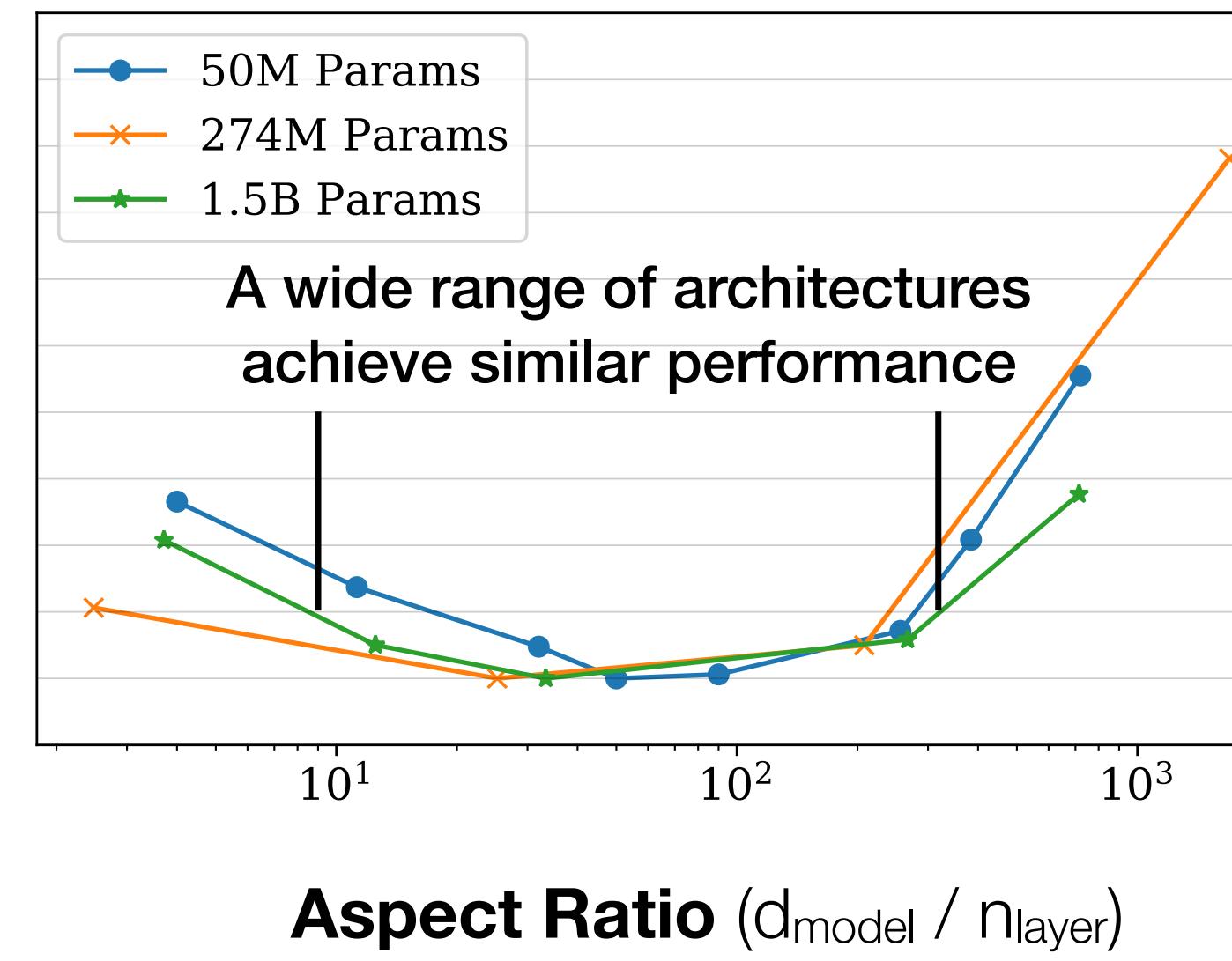
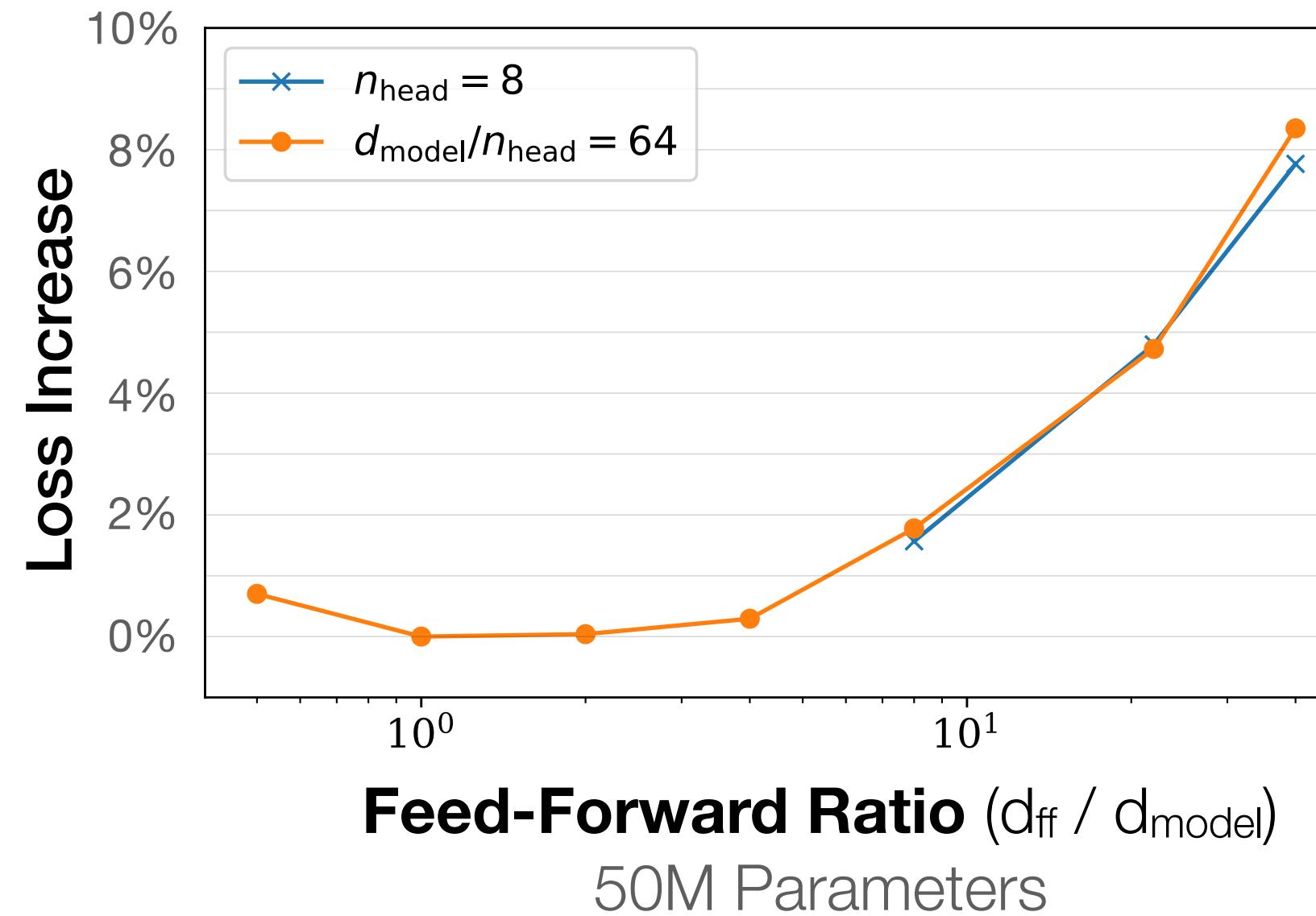
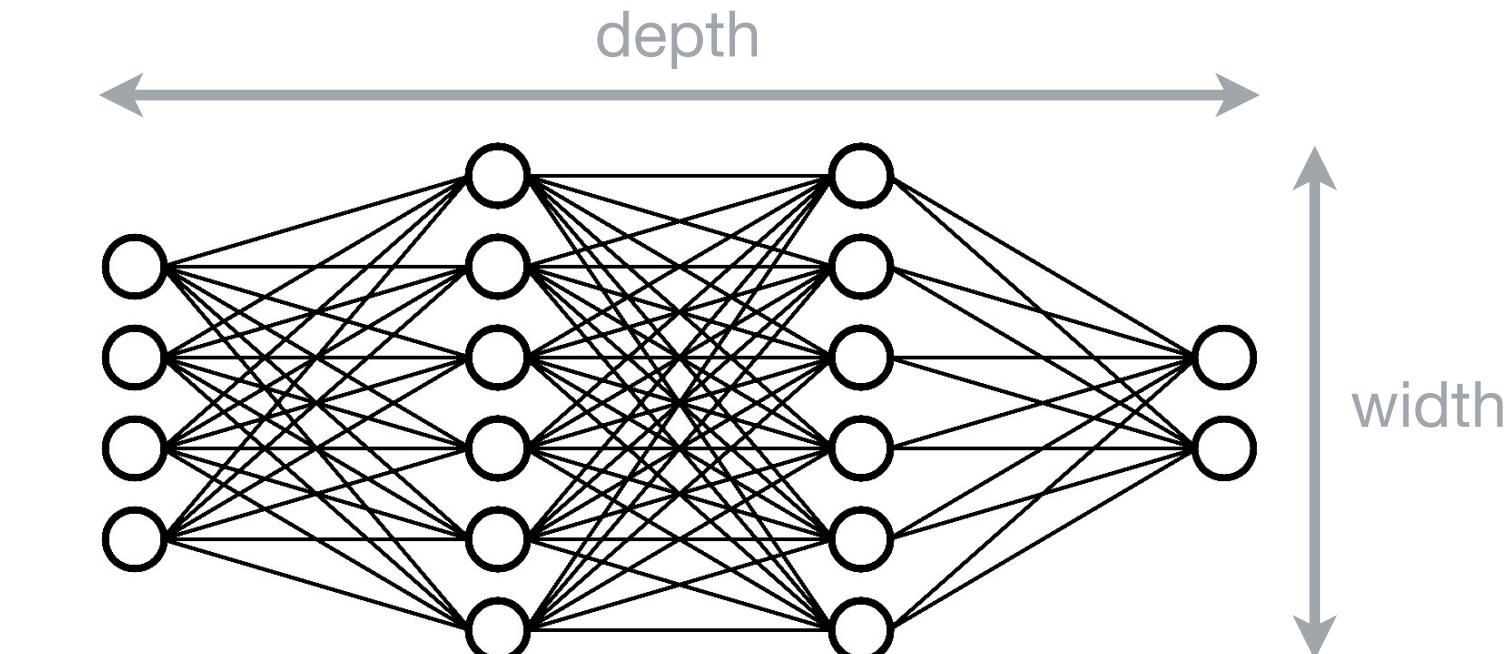
... in a principled way.



Hyperparameter Tuning vs. Scaling

Classic approach: Hyperparameter tuning on big models.

Wouldn't it be great to tune hyperparameters on small models and extrapolate to large ones?



Sample Complexity and Rates

Learning Theory Fundamentals: Generalization Bounds

$$\text{True Error} \leq \text{Training Error} + \text{Complexity}$$
$$\mathbb{E}[\mathcal{L}(f)] \leq \hat{\mathcal{L}}(f) + O\left(\sqrt{\frac{d}{n}}\right)$$

↑
loss of
model f ↑
number of
training samples ↑
VC dimension d ,
i.e., model capacity

VC Dimension (d):

- Measures model capacity
- Linear: $d = n_{\text{params}} + 1$
- Neural nets: $O(W \log W)$

Rademacher Complexity:

- Expected correlation with random labels
- Data-dependent measure

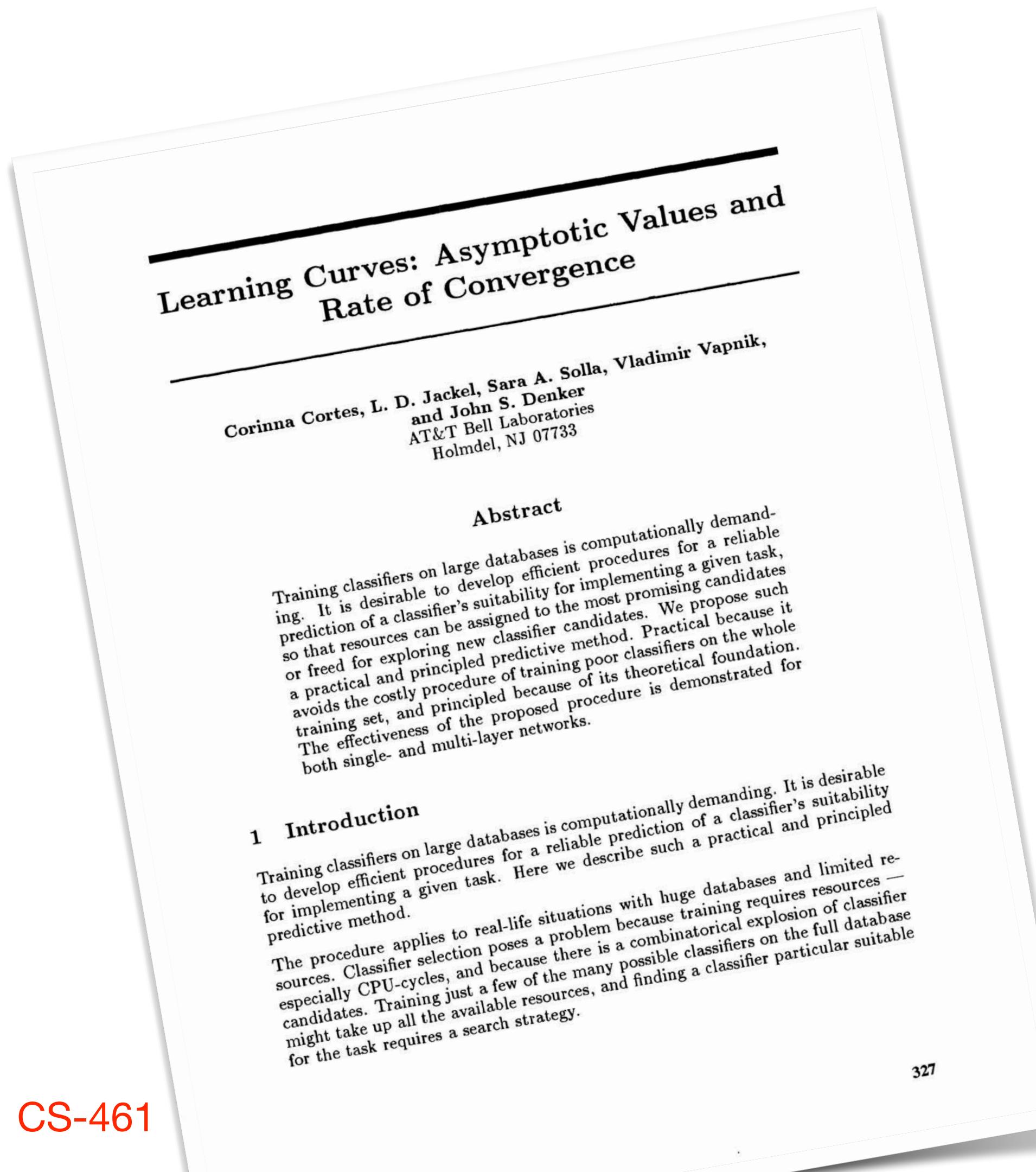
Key Insights: More data → Less overfitting

But fixed model → Performance plateau

Earliest Papers on Scaling Laws

1993

Theory meets reality!

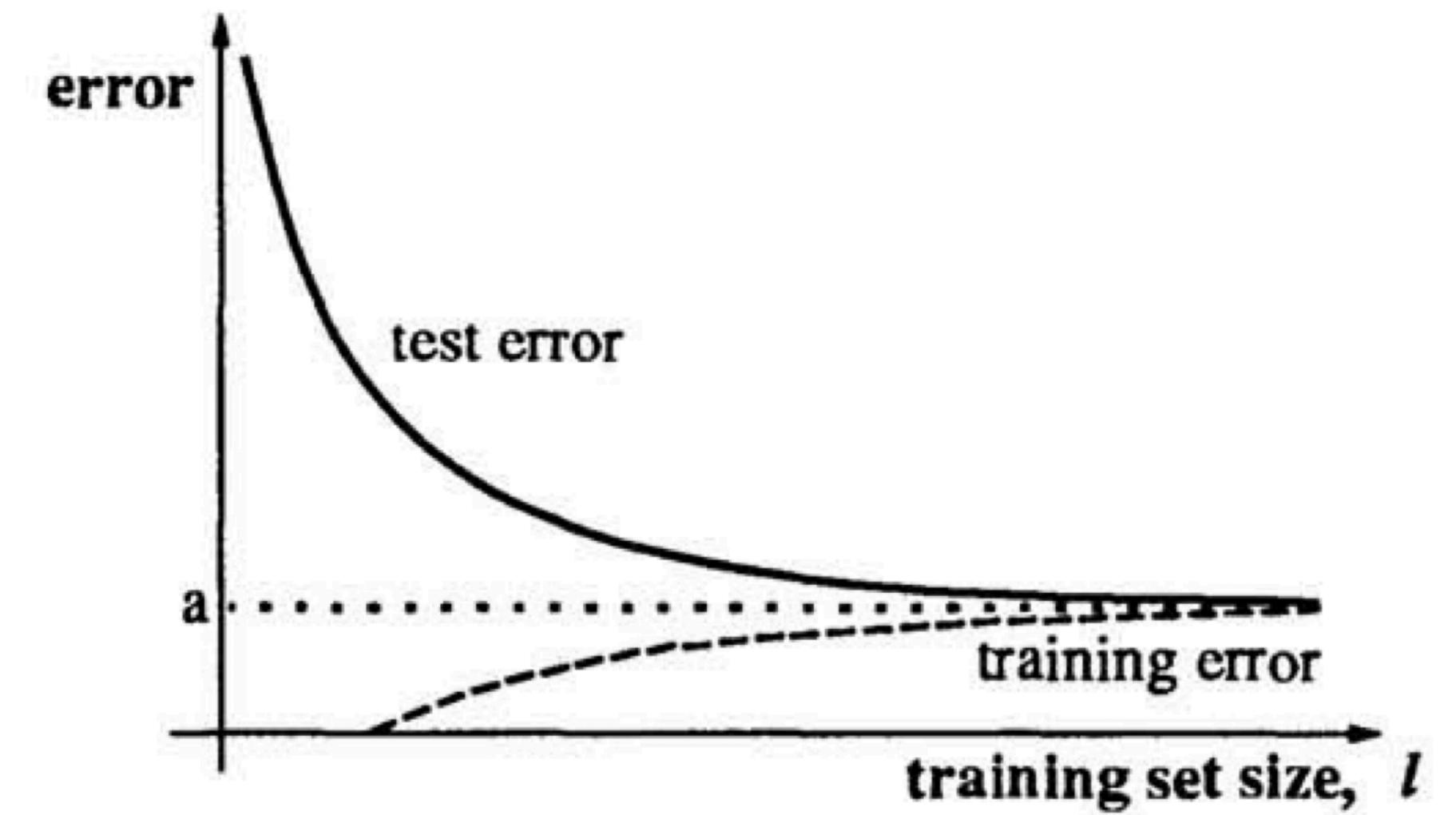


A typical example of learning curves is shown in Fig. 2. The test error is always larger than the training error, but asymptotically they reach a common value, a . We model the errors for large sizes of the training set as power-law decays to the asymptotic error value, a :

$$\mathcal{E}_{\text{test}} = a + \frac{b}{l^\alpha} \quad \text{and} \quad \mathcal{E}_{\text{train}} = a - \frac{c}{l^\beta}$$

where l is the size of the training set, and α and β are positive exponents.

Figure 2



Earliest Papers on Scaling Laws

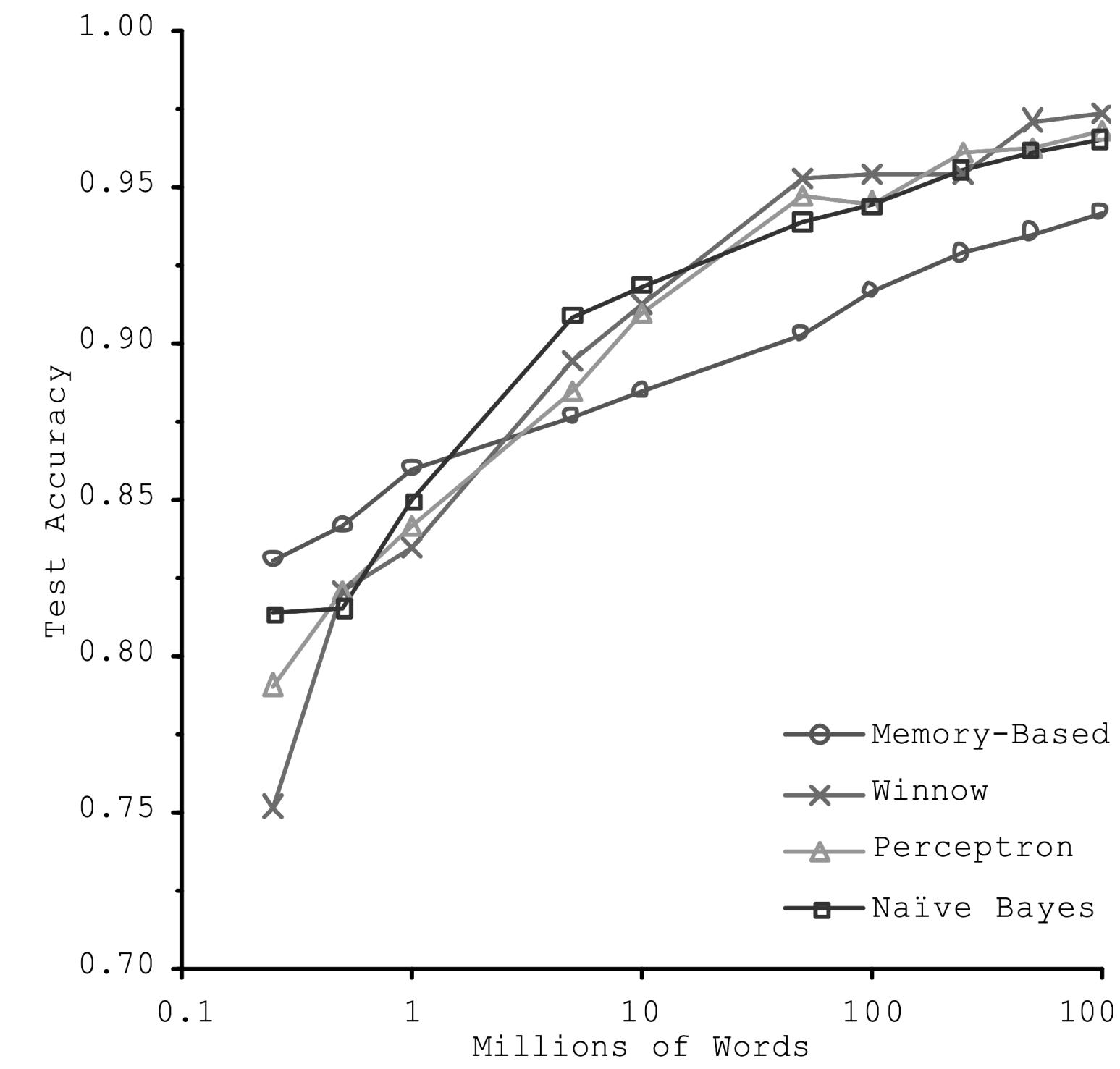
2001

Data scaling!



“ ... these results suggest that we may want to reconsider the trade-off between spending time and money on algorithm development versus spending it on corpus development/

**Log-linear scaling
with growing
amount of data**

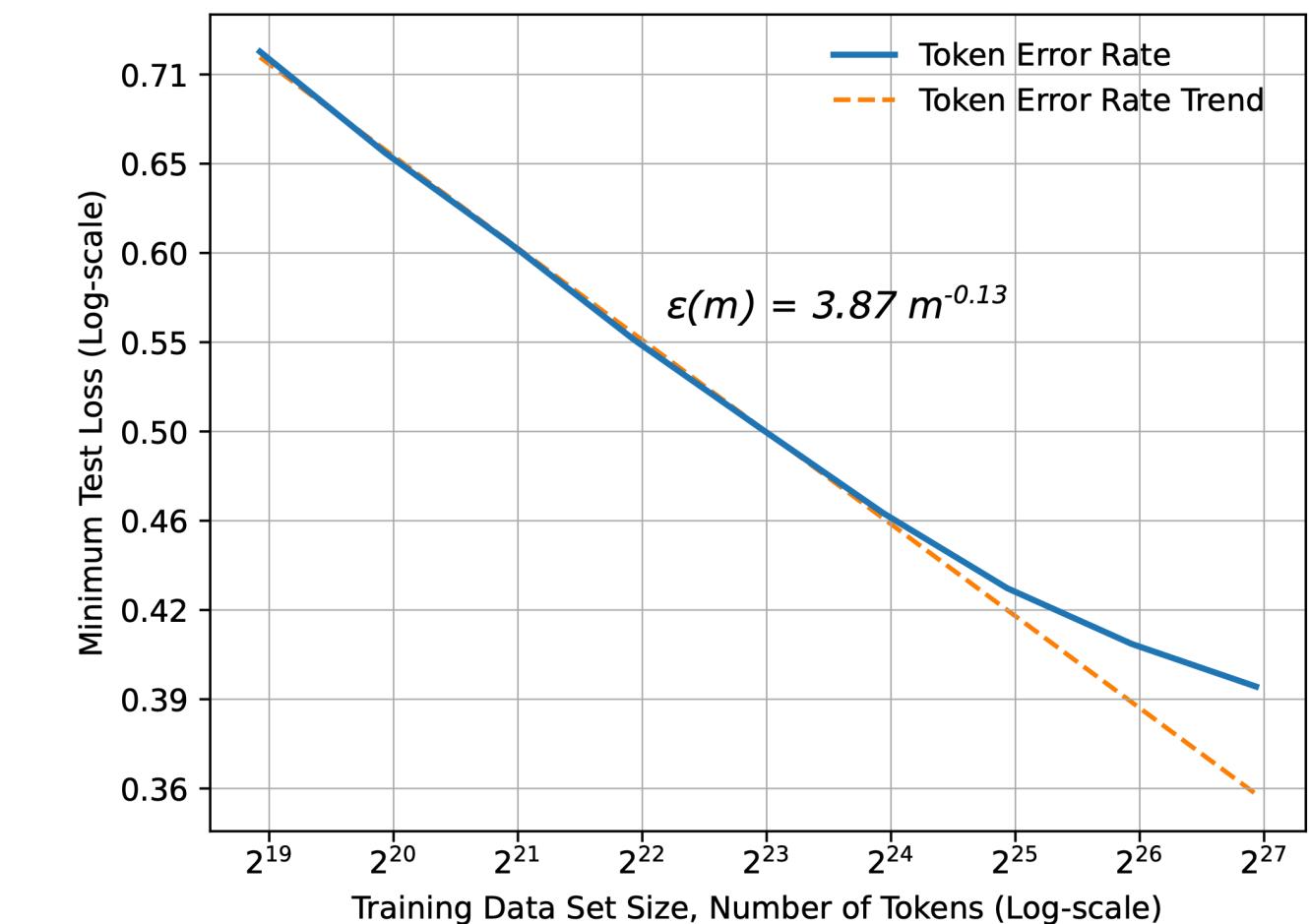
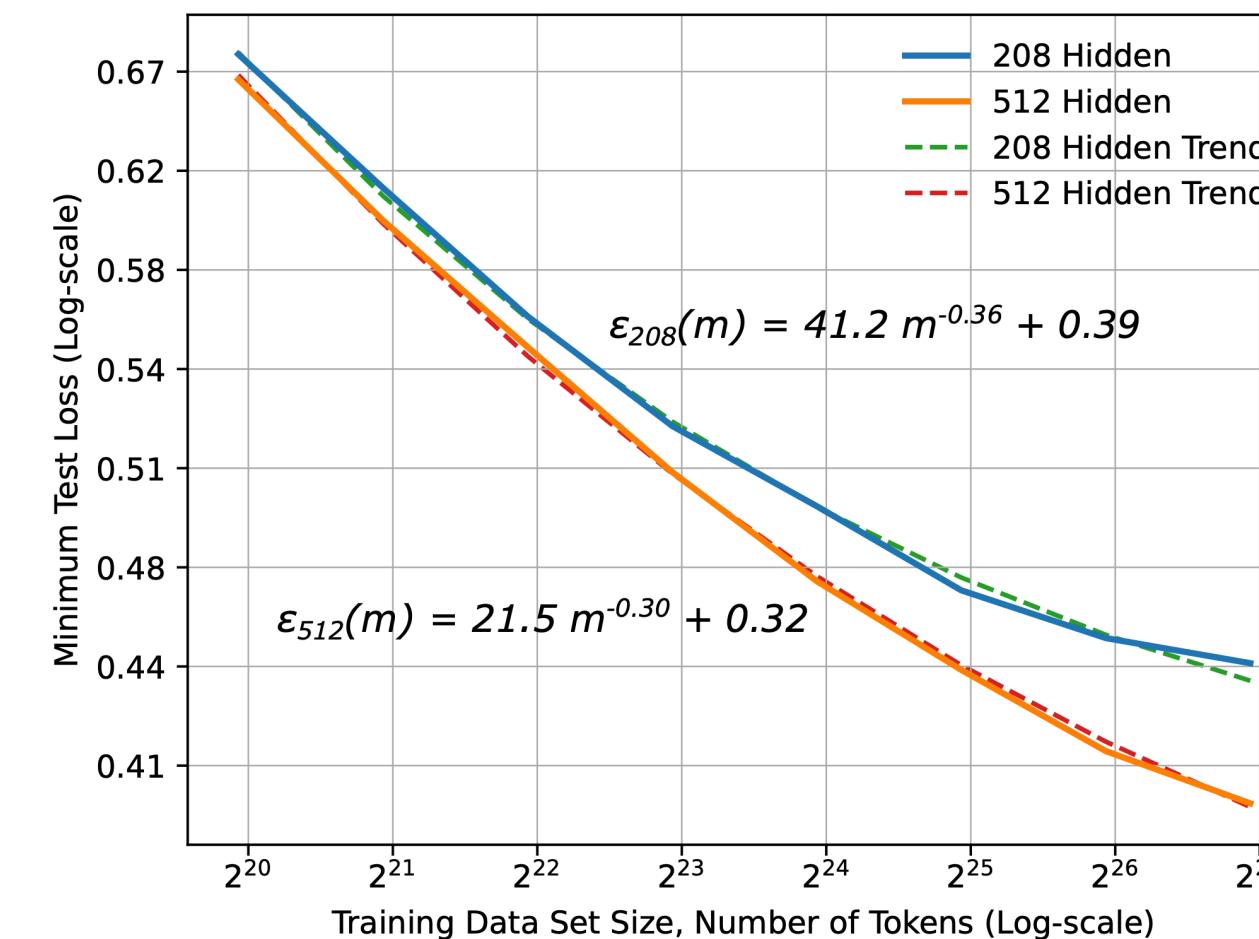


Charlotte Bunne

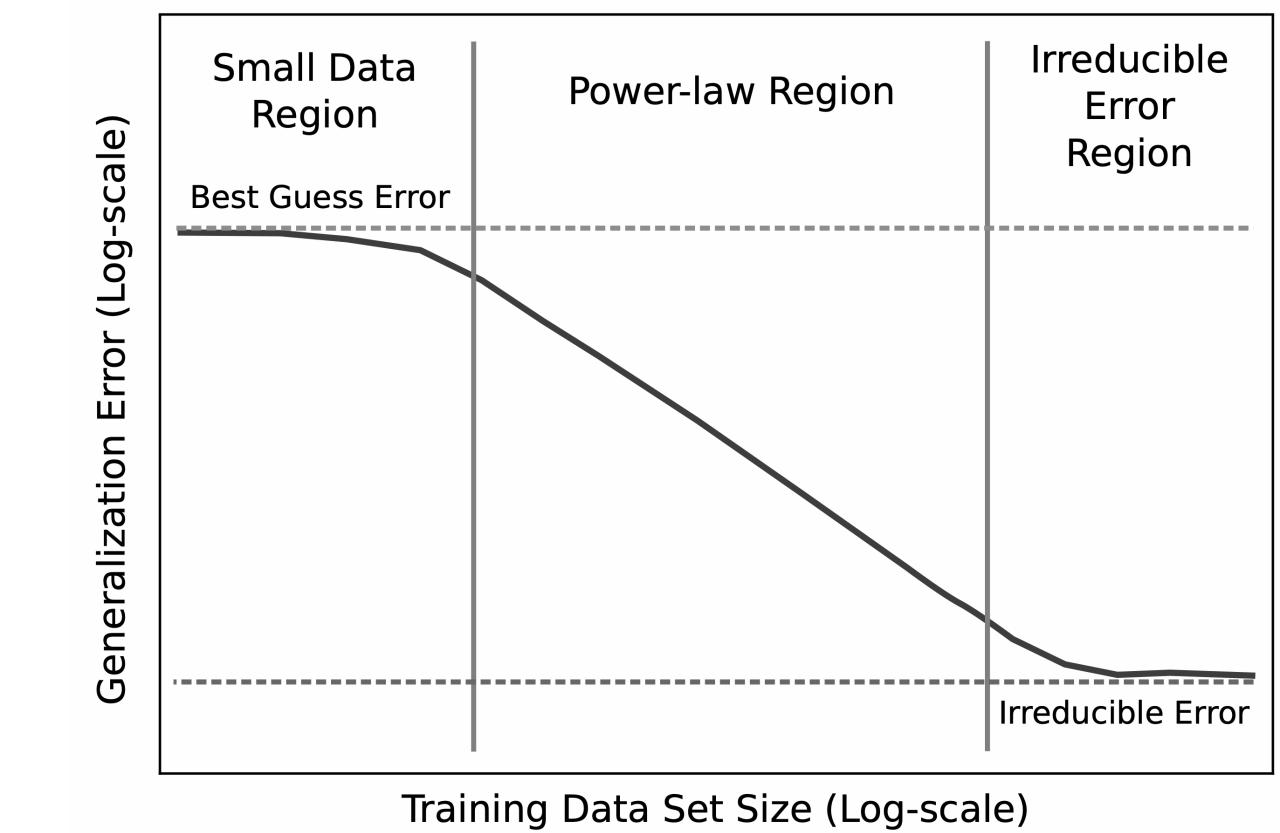
Earliest Papers on Scaling Laws

2017

Earliest large-scale neural' scaling work. Very ahead of their time!



different regions of model behavior



Earliest Papers on Scaling Laws

2017

Earliest large-scale neural' scaling work. Very ahead of their time!



“ ... it can be difficult to ensure that training data is large enough to see the power-law learning curve region. We have found that models ... where accuracy is only as good as best guessing, but the model trains on enough data to be in the power-law region.

“ ... model size curves may offer a way to project the compute requirements to reach a particular accuracy level. The compute requirements could inform decisions about how to scale computational capacity to unlock these compute-limited applications.

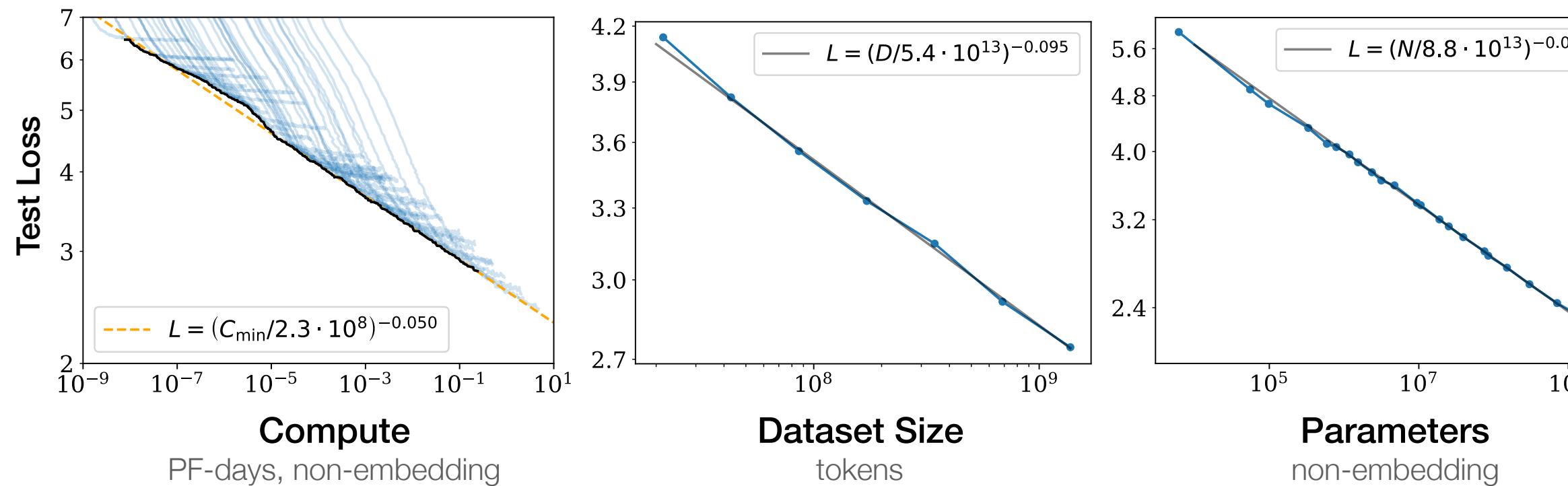
→ foreshadowed the idea of **scaling by compute**.

Scaling Laws for Neural Networks

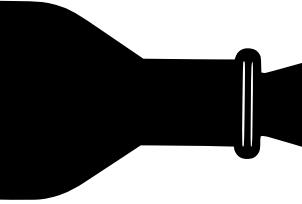
Neural scaling laws establish a mathematical relationship between the following variables in training a model.

Performance has a **power-law relationship** with each of the three scale factors N, D, C when not bottlenecked by the other two.

- \mathcal{L} prediction loss
- N number of model parameters
- D size of the dataset used to train the model
- C total compute used to train the model



Kaplan et al., (2020)

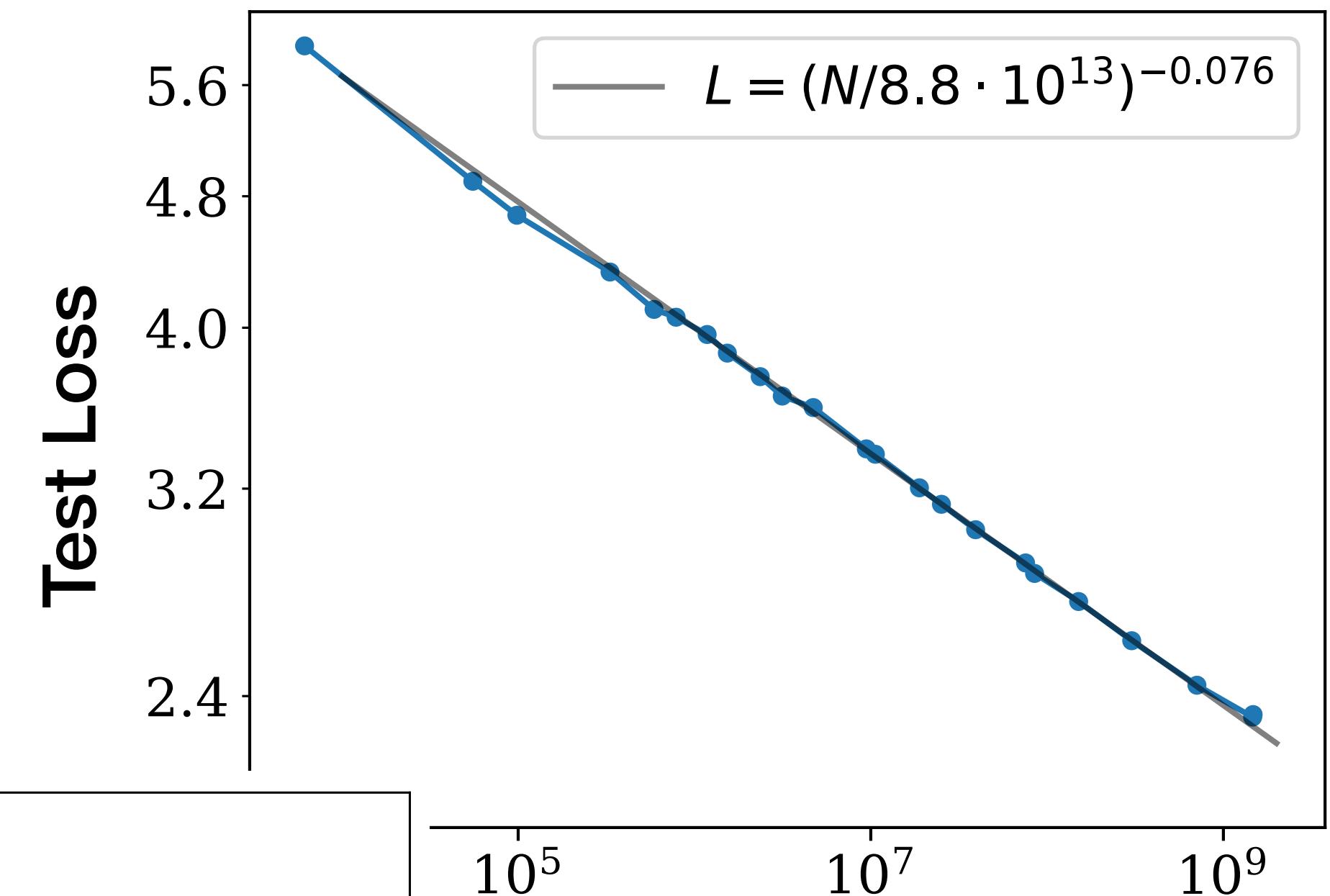
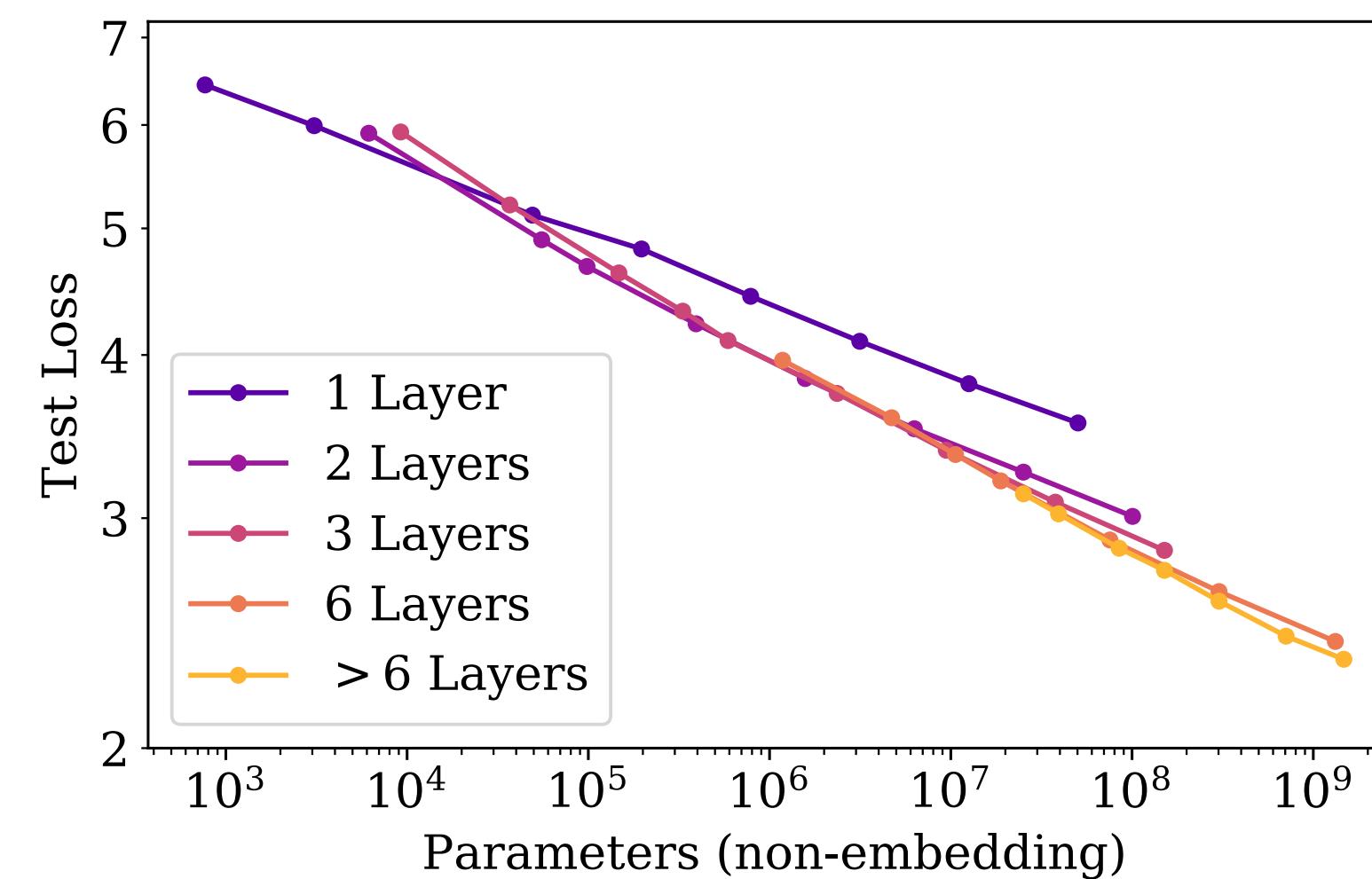


Bottleneck 1: Number of Parameters

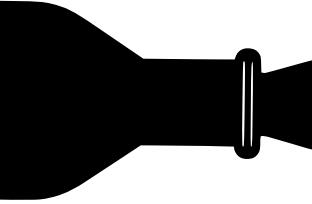
1. Fix the number of parameters N . Dataset size D and amount of compute C are not considered a bottleneck.
2. Train model until convergence.
3. Plot test loss \mathcal{L} as function of the number of parameters N .

$$\rightarrow \mathcal{L} \approx \left(\frac{N_c}{N} \right)^{\alpha_N}$$

... specifying the number of parameters does not uniquely specify model architecture but **model size matters more than model shape**



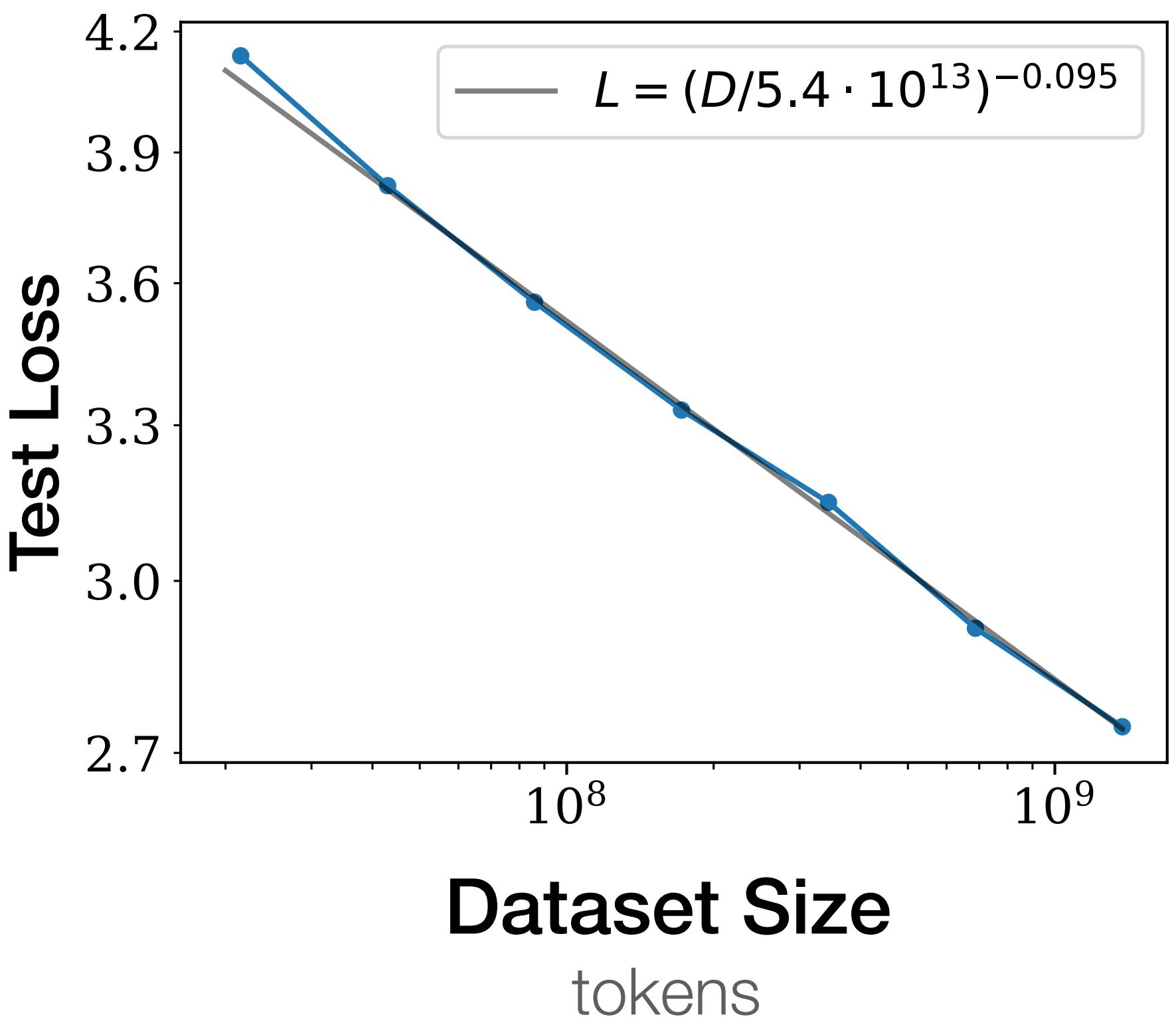
Parameters
non-embedding

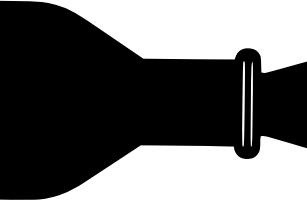


Bottleneck 2: Dataset Size

1. Fix the number of training examples D . Number of parameters N and amount of compute C are not considered a bottleneck.
2. Train a model until test loss stops decreasing.
3. Plot test loss \mathcal{L} as function of the number of training examples D .

$$\longrightarrow \mathcal{L} \approx \left(\frac{D_c}{D} \right)^{\alpha_D}$$





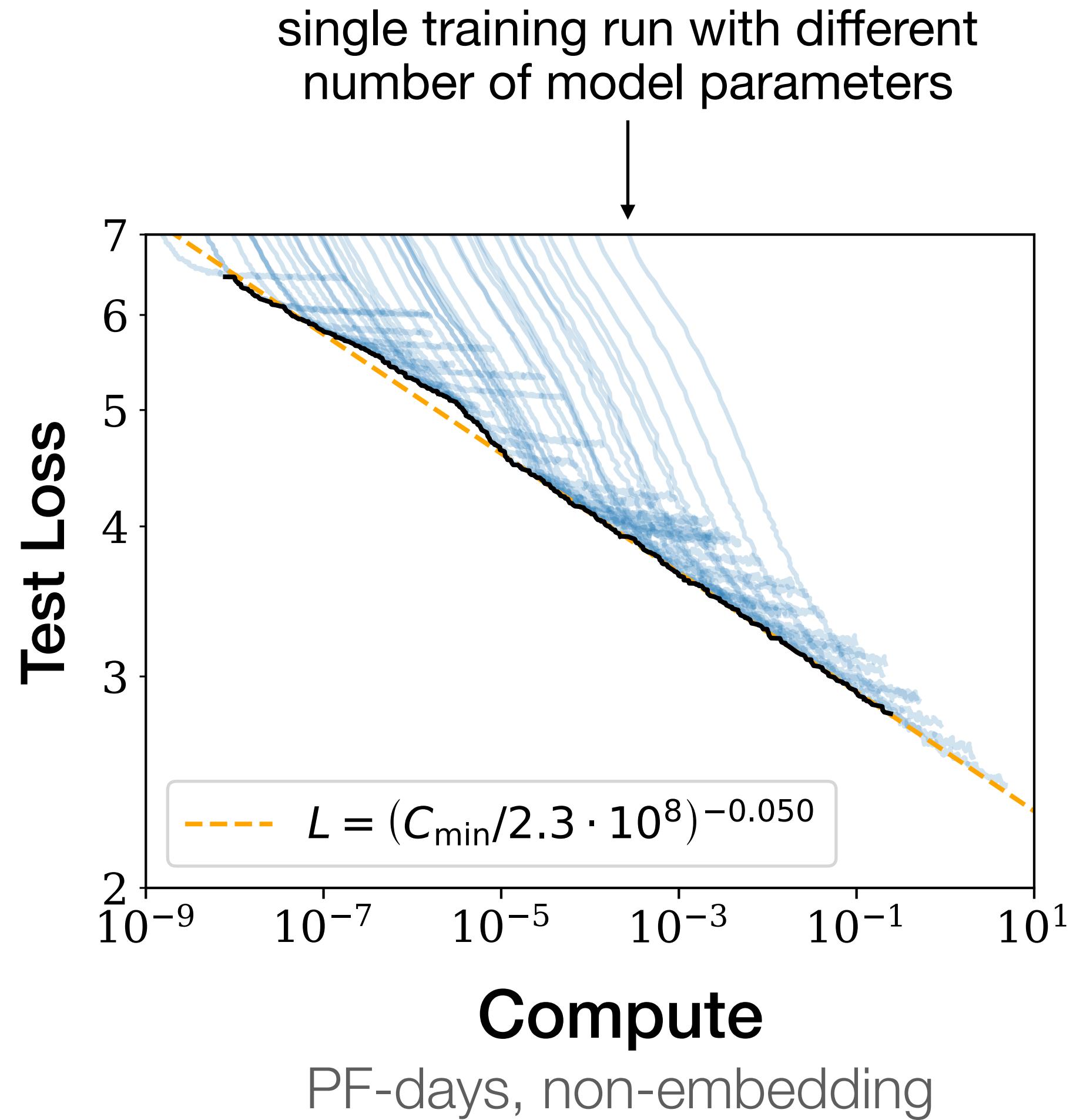
Bottleneck 3: Amount of Compute

1. Train models of various sizes N with a very large dataset of size D .
2. For each model, plot the test loss as function of the compute used.
3. Plot minimum test loss \mathcal{L} as function of the amount of compute C used.

$$\longrightarrow \mathcal{L} \approx \left(\frac{C_c}{C} \right)^{\alpha_C}$$

Measuring Compute (in a Transformer)

- Forward compute $\approx 2ND$
- Backward compute $\approx 4ND$



Language Model Sizes

State of AI, 'Report 2020'.
Published online at '<https://www.stateof.ai/2020>'.

Huge models and massive training costs dominate the hottest area of AI.

2018 (left) through 2019 (right)

>_ Code Notebook 1 · Task B

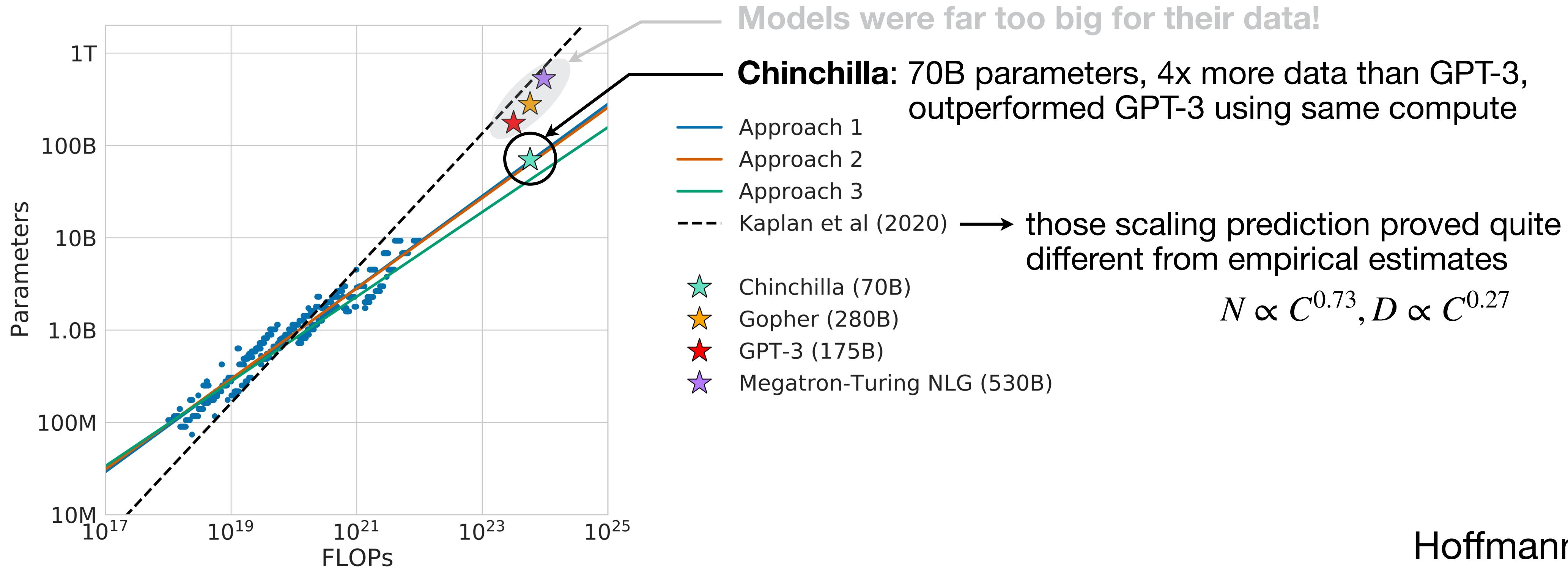


Caution: Model Size Is Not Everything

Researchers noticed giant models seemed undertrained ...

... Kaplan et al., (2020) had fixed one entity and varied the others.

But what if we systematically vary *both* model size and dataset size?



Estimating Optimal Parameter/Training Tokens Allocation

Approach 1: Fix model sizes, vary training tokens

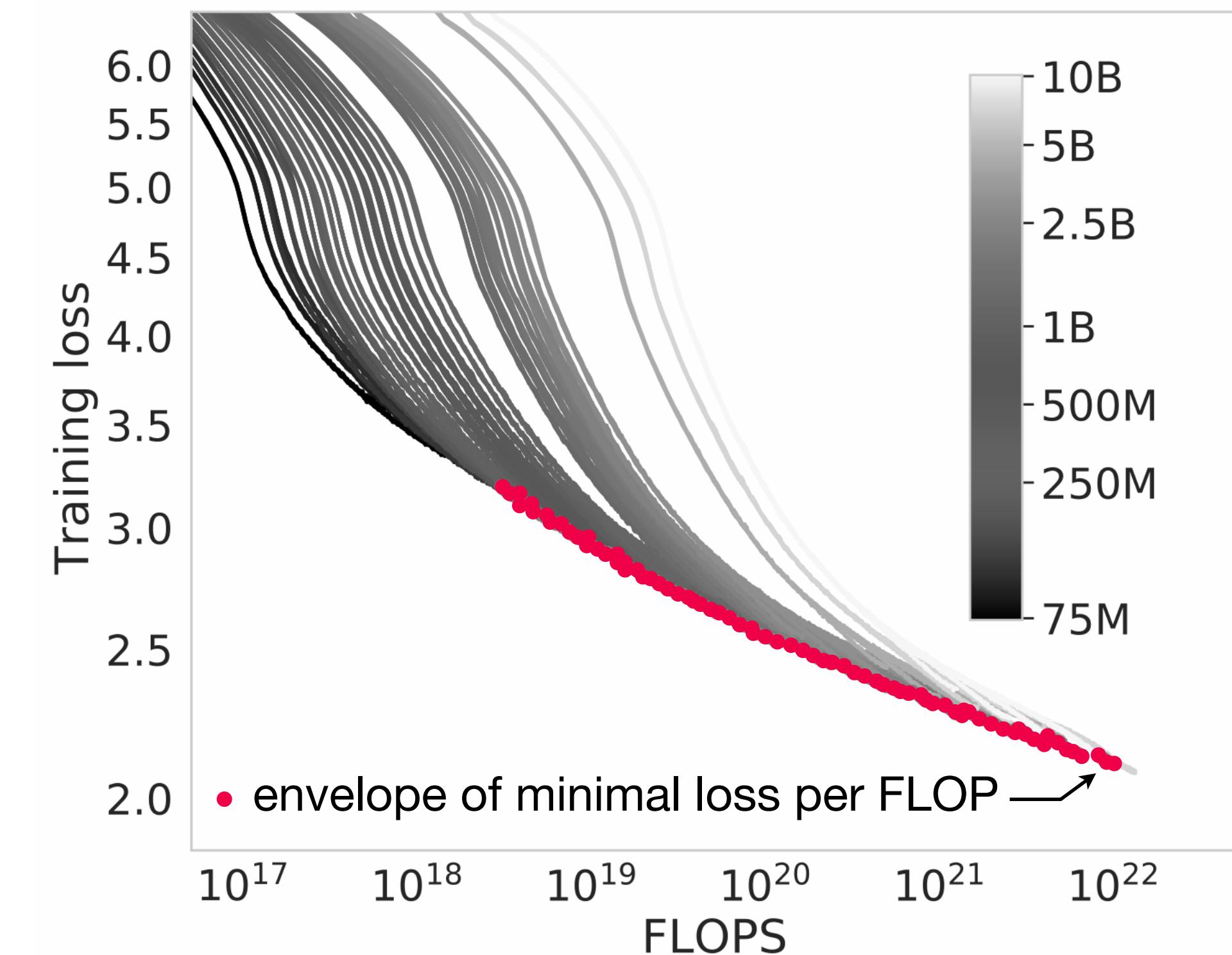
- Trained models (70M to 10B) for different durations
- Found training duration that minimized loss per FLOP
- Extracted the "envelope" of best performance at each compute
- **Finding:** $N \propto C^{0.50}$, $D \propto C^{0.50}$

Approach 2: Fixed compute budgets

- Fixed 9 compute budgets, varied model size at each
- Found the "valley": optimal model size that minimizes loss
- Clear parabolas showing sweet spot between "too small/large"
- **Finding:** $N \propto C^{0.49}$, $D \propto C^{0.51}$

Approach 3: Parametric fitting

- Fit equation $L(N,D)$ to >400 training runs
- Minimized loss equation under compute constraints
- Predicted optimal allocation mathematically
- **Finding:** $N \propto C^{0.46}$, $D \propto C^{0.54}$



• use envelope points to estimate
the optimal model size

Estimating Optimal Parameter/Training Tokens Allocation

Approach 1: Fix model sizes, vary training tokens

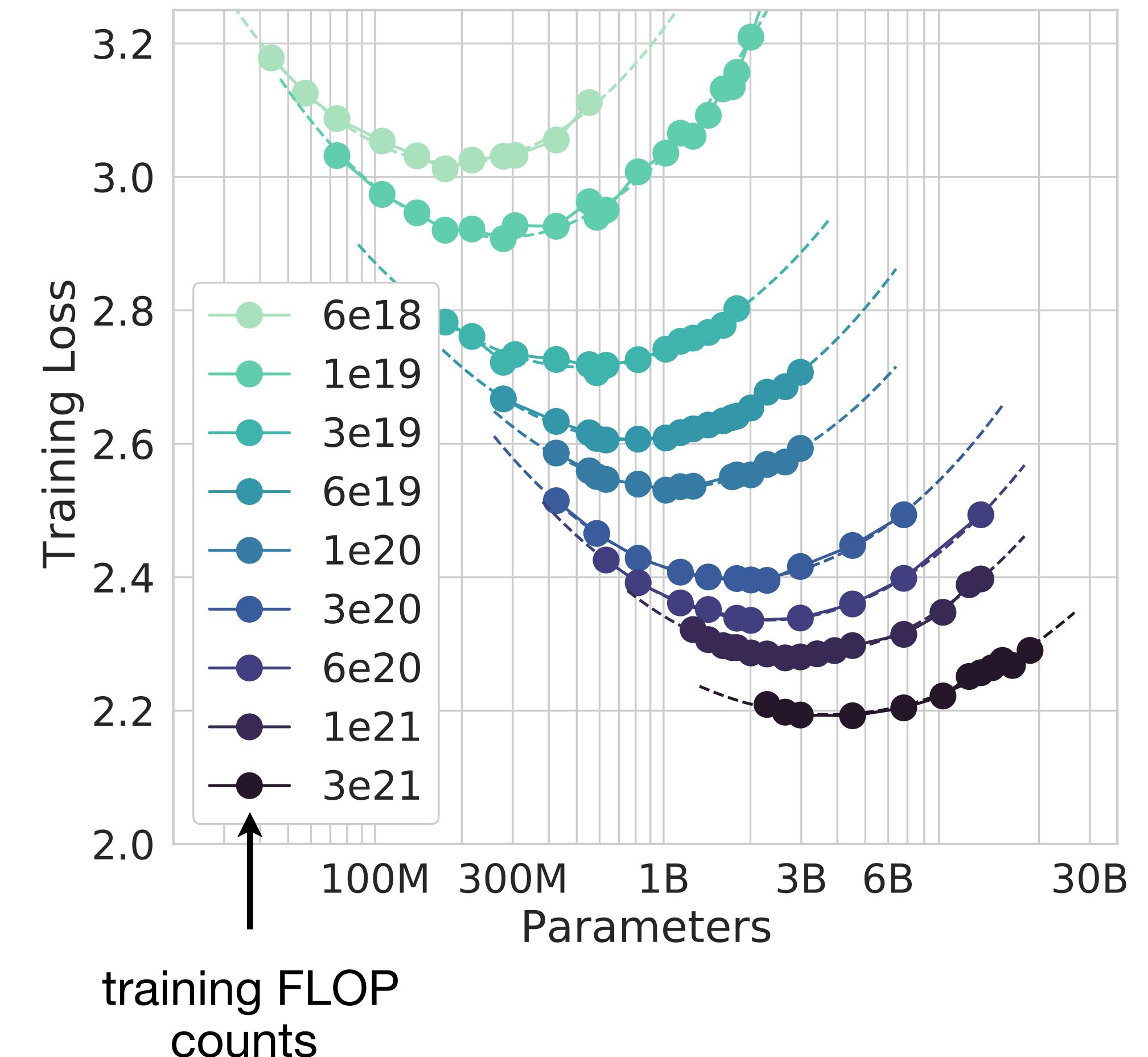
- Trained models (70M to 10B) for different durations
- Found training duration that minimized loss per FLOP
- Extracted the "envelope" of best performance at each compute
- **Finding:** $N \propto C^{0.50}$, $D \propto C^{0.50}$

Approach 2: Fixed compute budgets

- Fixed 9 compute budgets, varied model size at each
- Found the "valley": optimal model size that minimizes loss
- Clear parabolas showing sweet spot between "too small/large"
- **Finding:** $N \propto C^{0.49}$, $D \propto C^{0.51}$

Approach 3: Parametric fitting

- Fit equation $L(N,D)$ to >400 training runs
- Minimized loss equation under compute constraints
- Predicted optimal allocation mathematically
- **Finding:** $N \propto C^{0.46}$, $D \propto C^{0.54}$



Estimating Optimal Parameter/Training Tokens Allocation

Approach 1: Fix model sizes, vary training tokens

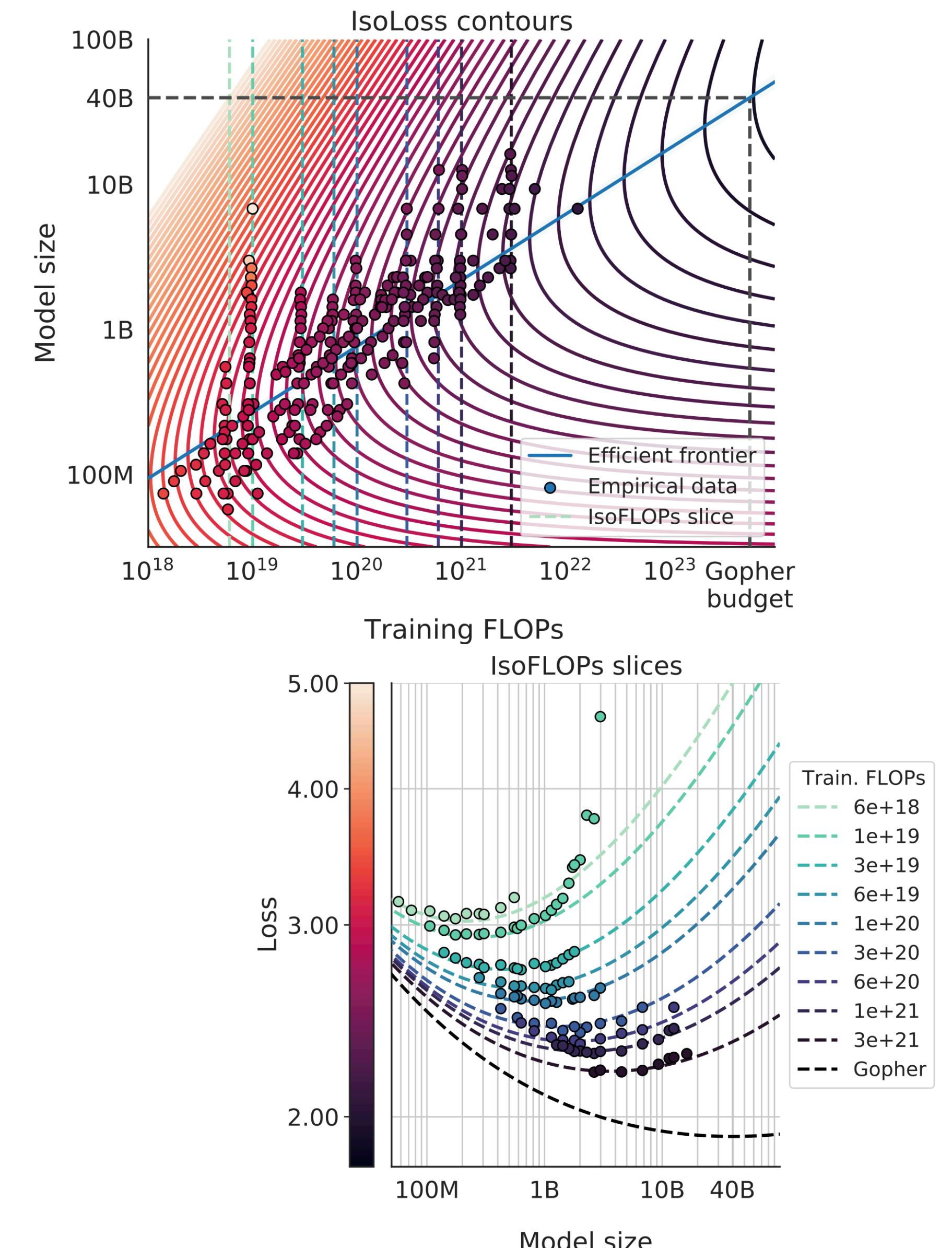
- Trained models (70M to 10B) for different durations
- Found training duration that minimized loss per FLOP
- Extracted the envelope of best performance at each compute
- **Finding:** $N \propto C^{0.50}, D \propto C^{0.50}$

Approach 2: Fixed compute budgets

- Fixed 9 compute budgets, varied model size at each
- Found the “valley”: optimal model size that minimizes loss
- Clear parabolas showing sweet spot between “too small/large”
- **Finding:** $N \propto C^{0.49}, D \propto C^{0.51}$

Approach 3: Parametric fitting

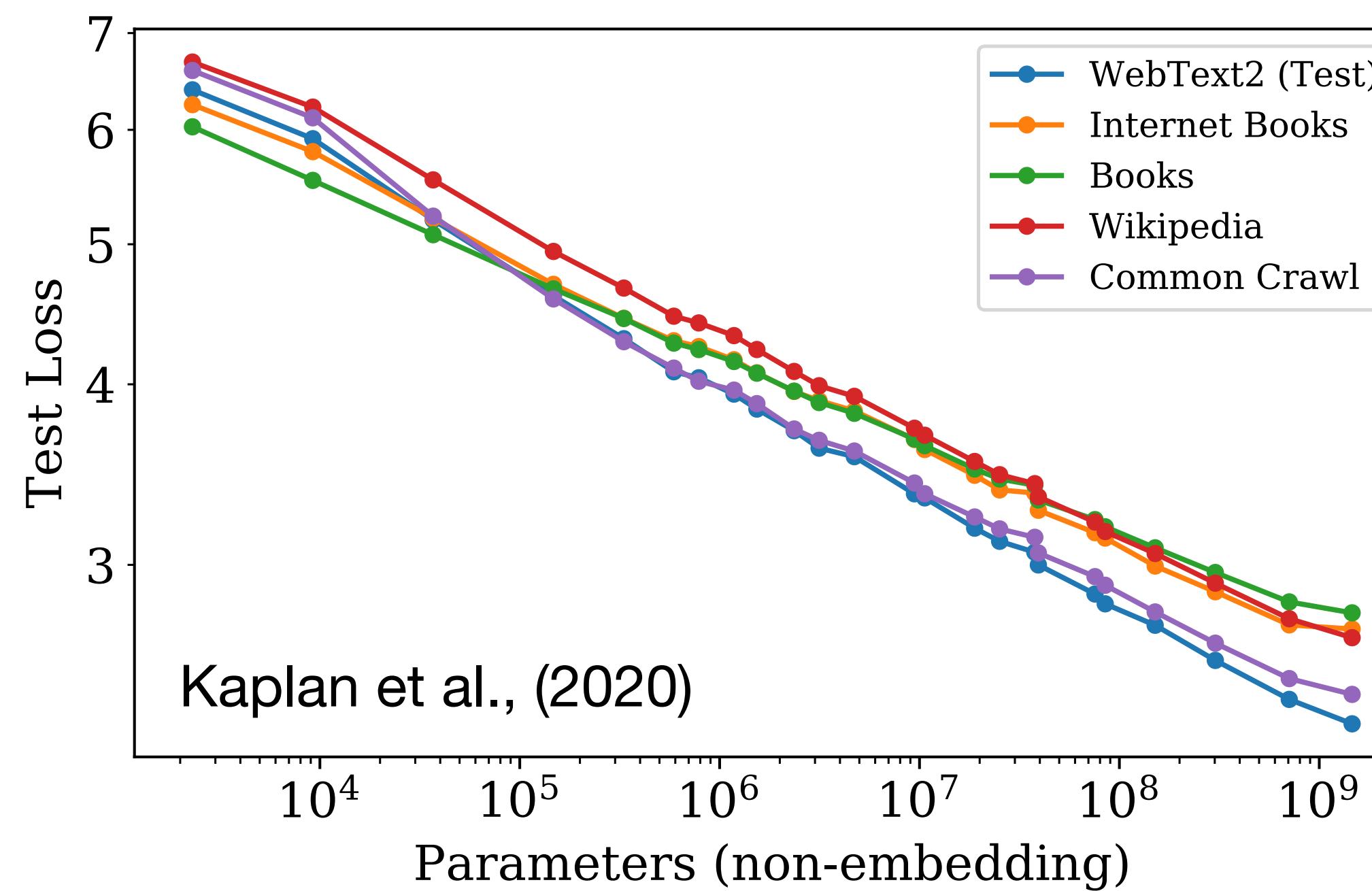
- Fit equation $L(N,D)$ to >400 training runs
- Minimized loss equation under compute constraints
- Predicted optimal allocation mathematically
- **Finding:** $N \propto C^{0.46}, D \propto C^{0.54}$



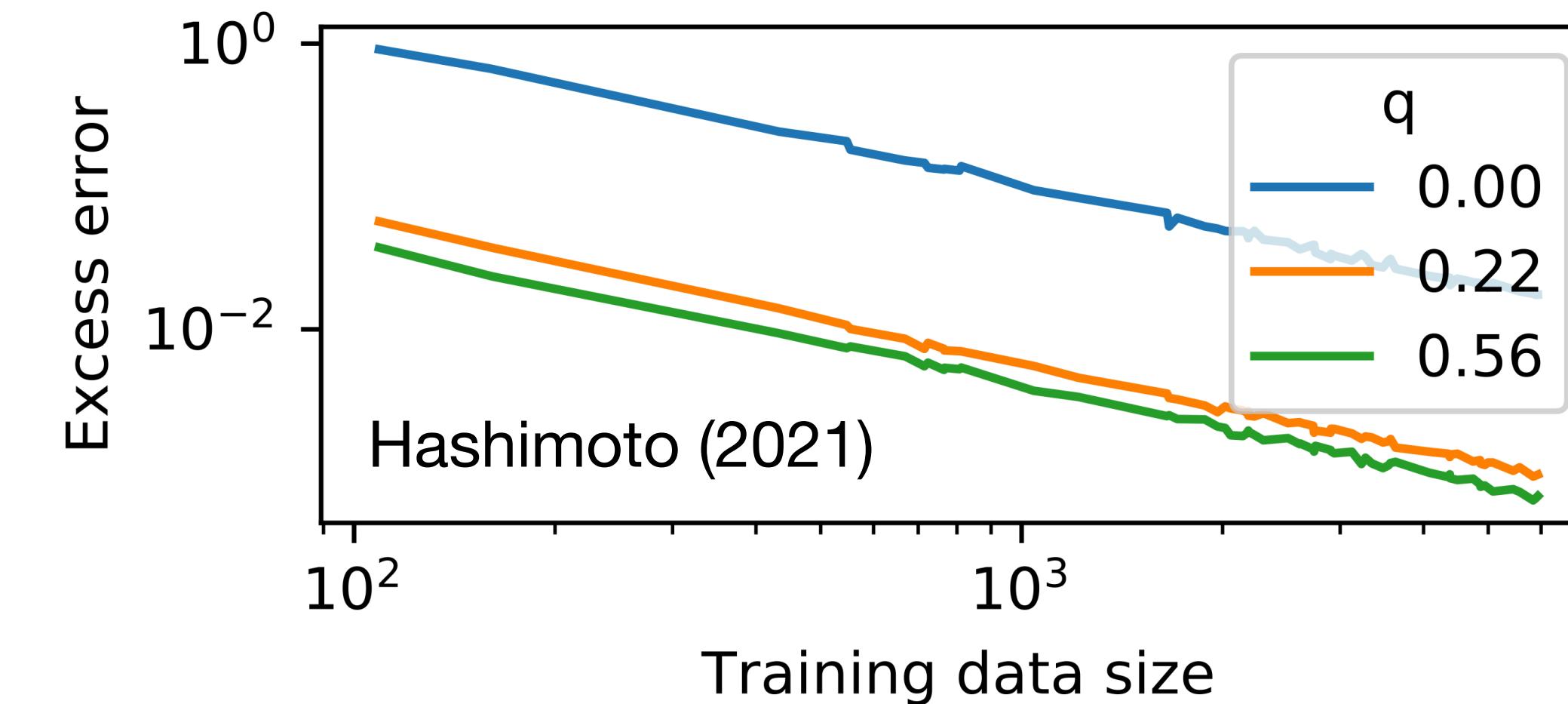
It's Also About Data

Scaling Laws under Distribution Shifts

- Different datasets follow the same power law slope
- Only the vertical offset changes between datasets
- Models can transfer across domains with performance penalties



- Data mixture composition affects performance non-linearly
- Optimal performance often requires diverse data sources
- As q changes from pure ($q = 0$) to mixed datasets ($q > 0$), performance improves



Similar Patterns for Beyond Language

... for generative modeling



... for vision transformers



... for molecular representation learning



1 Introduction

Large language models lie at the center of recent advances in artificial intelligence. Shared across nearly all such language models is a common recipe: learn a model that maximizes data likelihoods using an autoregressive, left-to-right factorization. Maximum-likelihood pretraining has been a remarkably successful paradigm, leading to models that perform well on a range of downstream tasks and display complex behaviors like in-context learning [1, 28].

Thus far, autoregressive modeling has been a core part of this process due to its computational efficiency and empirical performance. However, this choice carries drawbacks. Autoregressive models generate tokens one at a time, making it difficult to perform long-range planning or controllable generation [17, 19, 21]. In addition, certain sequence distributions may be fundamentally more difficult to model autoregressively [22].

Given the importance of language modeling, these potential drawbacks motivate us to explore alternatives to the autoregressive approach. As a promising candidate, we turn to continuous diffusion models [32, 12], which have achieved state-of-the-art results in image modeling [5, 30, 31]. In language, prior works on diffusion models exist (e.g. [21, 10, 6]), but these optimize non-likelihood-based objectives. Without the ability to use standard likelihood-based benchmarks [25, 14, 26], it is difficult to say precisely how these models compare to autoregressive models (see Section 7 for a discussion). Somewhat concerning, there is no work showing that it is possible for diffusion language models to achieve any nontrivial likelihoods on standard benchmarks.

In this work, we explore the limits of likelihood-based diffusion language models. Our goal is to train and release a diffusion model which achieves better likelihoods than GPT-2 124M [29], which

1. Introduction

Attention-based Transformer architectures [45] have taken computer vision domain by storm [8, 16] and are becoming an increasingly popular choice in research and practice. Previously, Transformers have been widely adopted in the natural language processing (NLP) domain [7, 15]. Optimal scaling of Transformers in NLP was carefully studied in [22], with the main conclusion that large models not only perform better, but do use large computational budgets more efficiently. However, it remains unclear to what extent these findings transfer to the vision domain, which has several important differences. For example, the most successful pre-training schemes in vision are supervised, as opposed to unsupervised pre-training in the NLP domain.

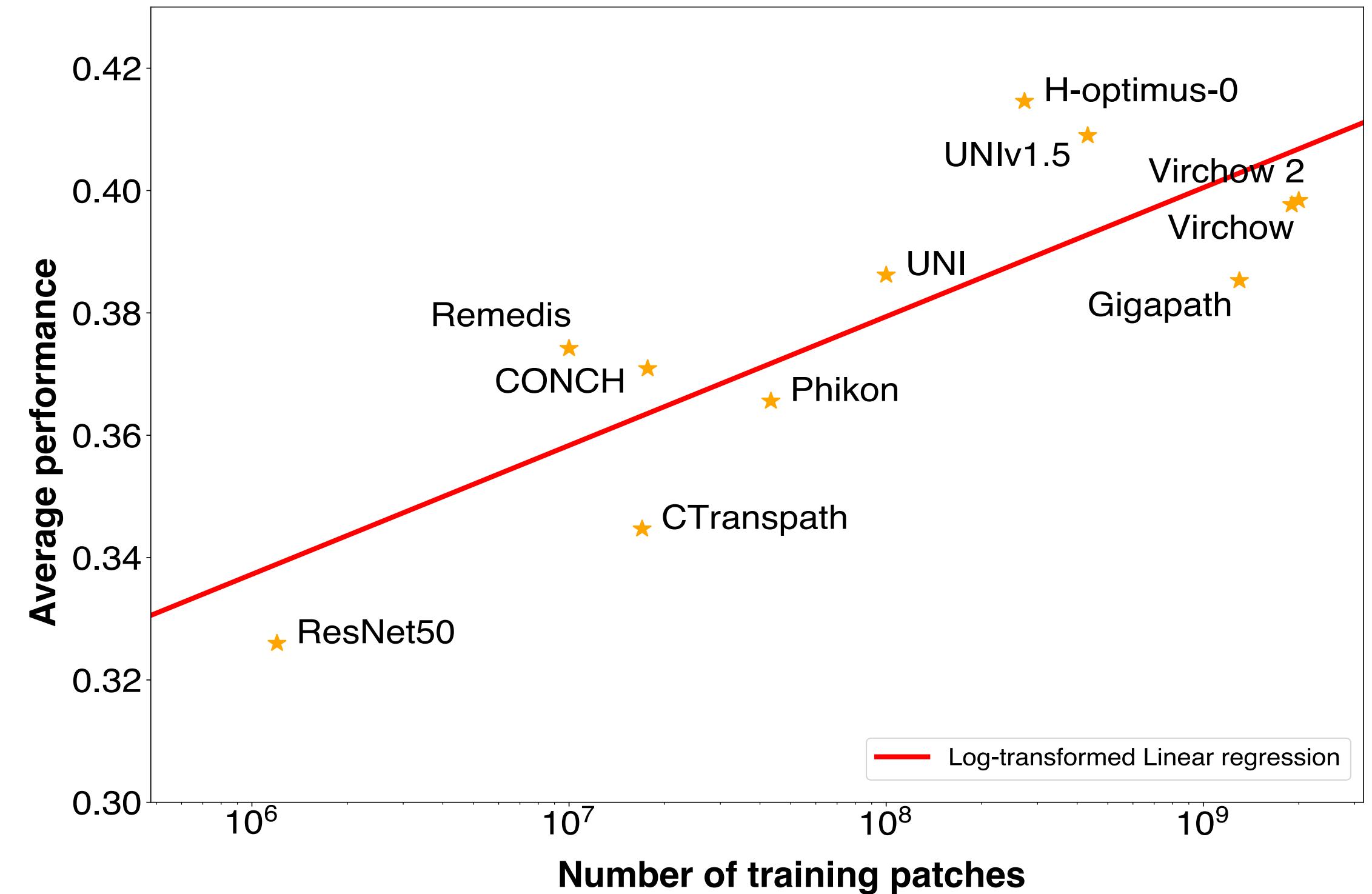
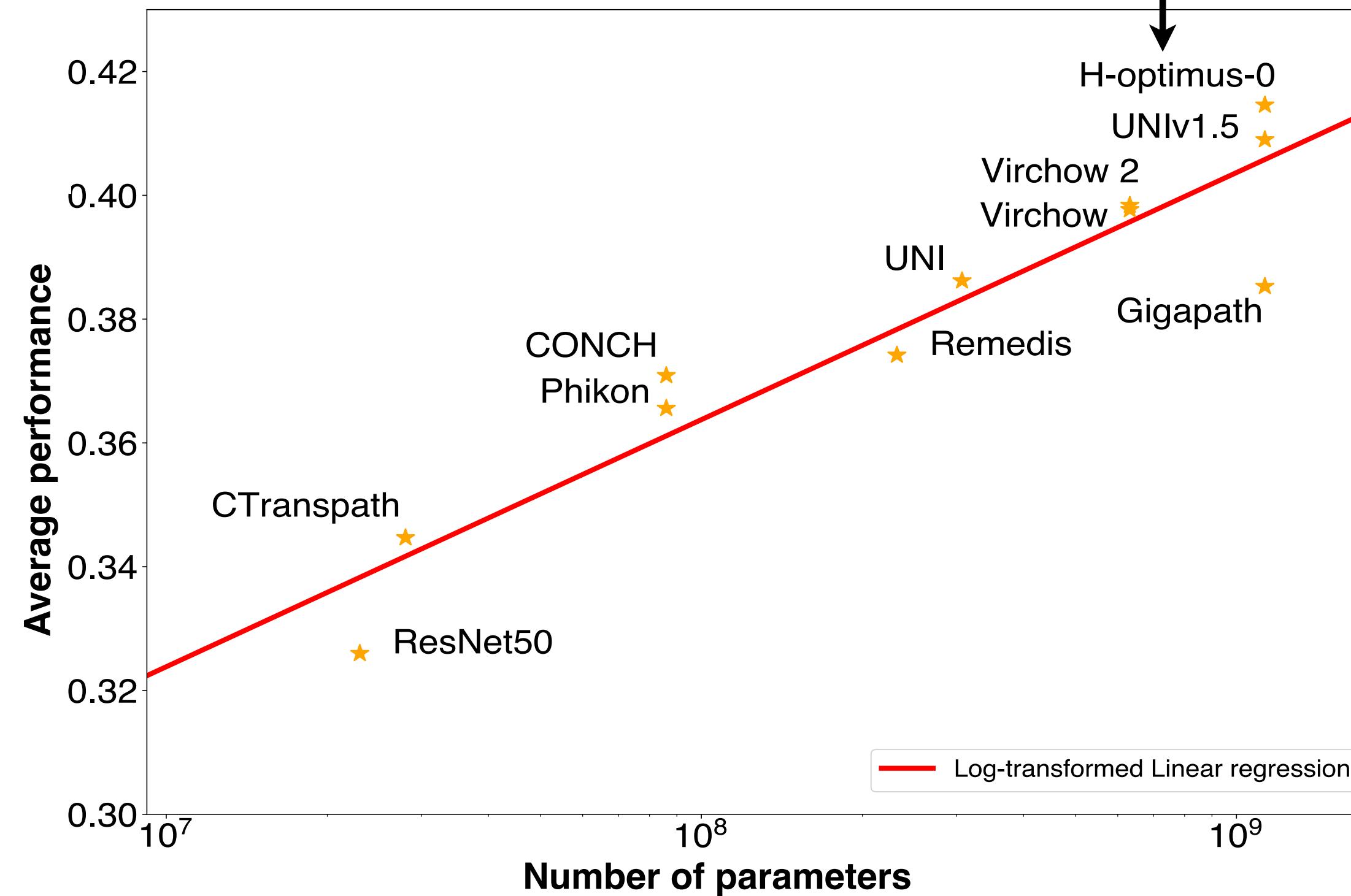
In this paper we concentrate on scaling laws for transfer performance of ViT models pre-trained on image classification datasets. We achieve this by hardware-specific architecture changes and a different optimizer. As a result, we train a model with two billion parameters and attain a new state-of-the-art 90.45% accuracy on ImageNet.

* equal contribution

[†]These authors made equal contribution to this research.
[†]Corresponding author: Shu Wu (shu.wu@nlpri.ac.cn).

Scaling Laws in Histopathology Foundation Models

Lecture 9: Guest Lecture



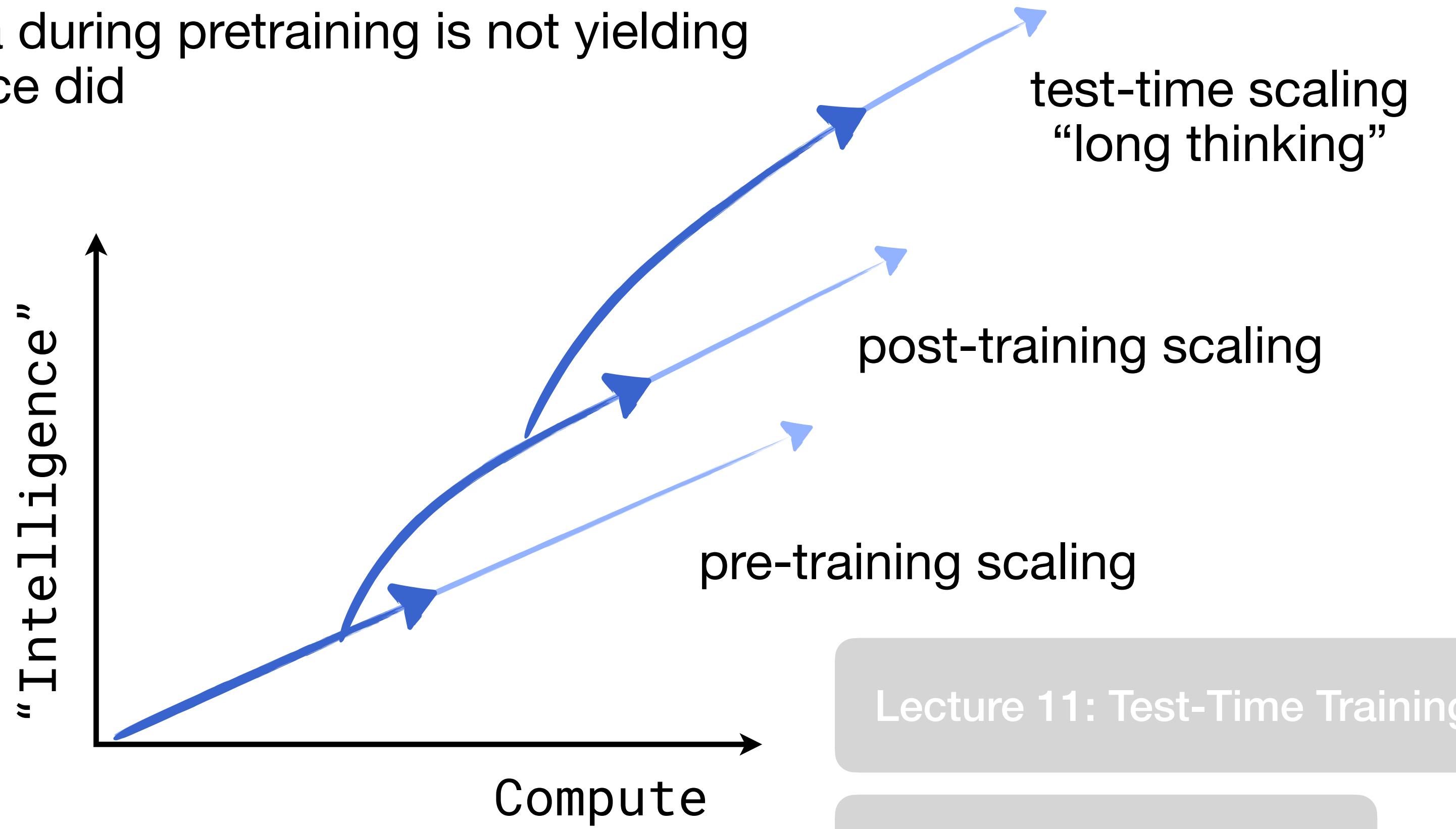
e.g., H-optimus-0 is trained on **1 million+ whole-slide images** from **800,000+ patients** with over **1 billion parameters**

Scale is not Intelligence: The end of scaling laws?

AI scaling laws are showing signs of diminishing returns

- Leading AI labs are seeing models improve more slowly than before
- Simply adding more GPUs and data during pretraining is not yielding the exponential improvements it once did

From One to Three Scaling Laws



Lecture 11: Test-Time Training

Lecture 13: Reasoning

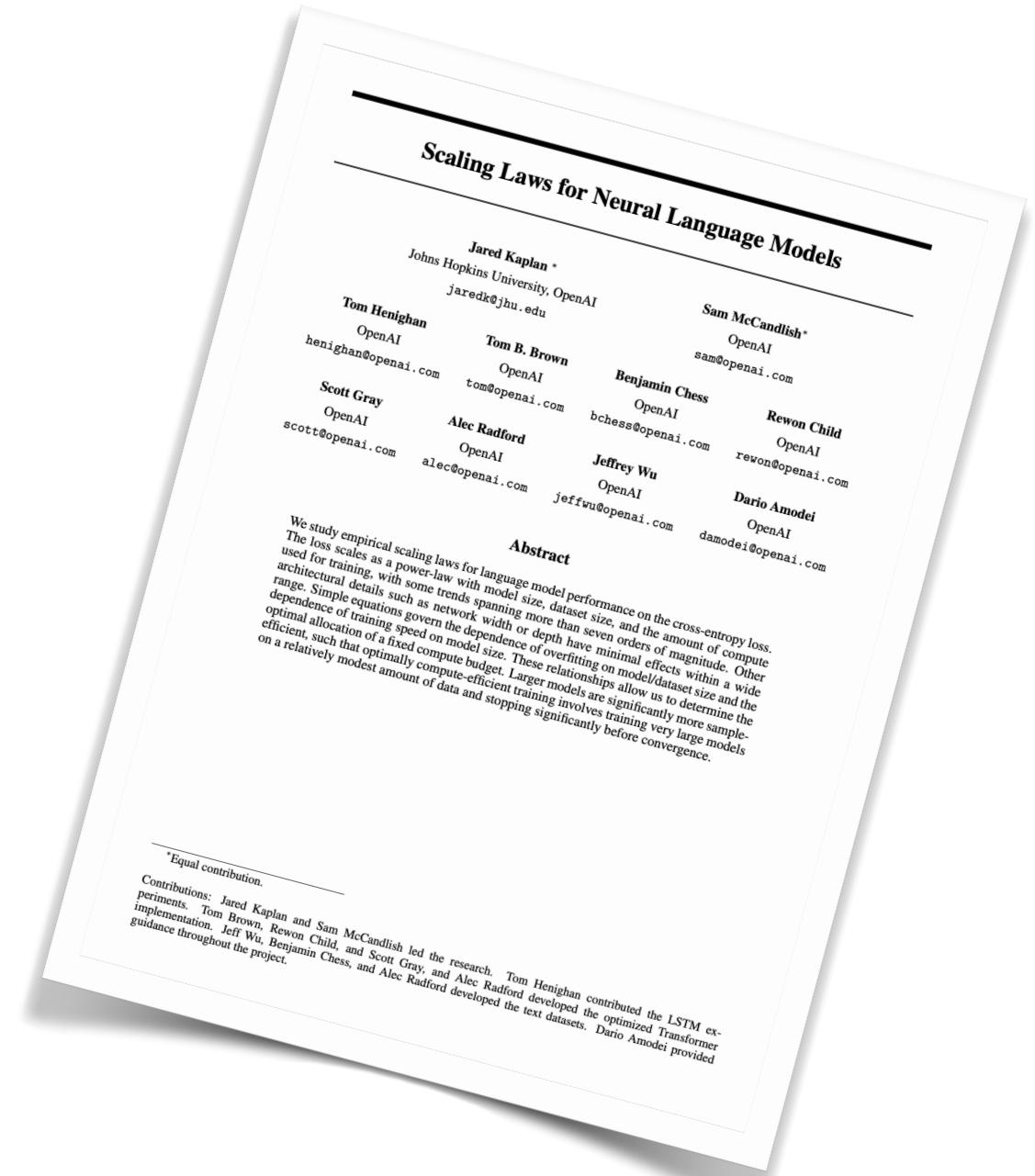
This Week's Papers



Papers are linked in Moodle.



Kaplan et al., "Scaling Laws for Neural Language Models." arXiv preprint arXiv:2001.08361 (2020).



Hoffmann et al., "Training Compute-Optimal Large Language Models." arXiv preprint arXiv:2203.15556 (2022).

... in addition to papers mentioned on exercise sheet.

This Week's Exercise Sheet



Exercise 1 · Task 1

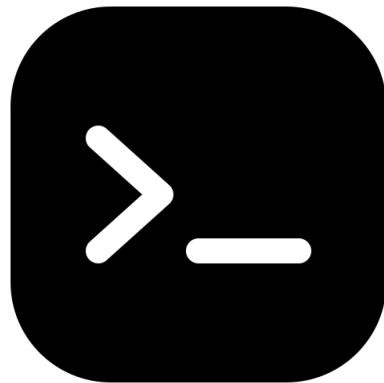
InfoNCE in Contrastive Learning: Explores the InfoNCE loss from a probabilistic perspective, connecting it to cross-entropy loss and NCE, while analyzing how the number of negative samples K affects mutual information learning.



Exercise 1 · Task 2

Masked Language Modeling as Pseudo-Likelihood: Establishes the mathematical relationship between MLM training objectives and pseudo-likelihood, including connections to pseudo-perplexity and comparisons with autoregressive models.

This Week's Code Demonstration



 Code Notebook 1 · Task A

Exploring Contrastive Learning with SimCLR: Hands-on implementation and exploration of the SimCLR framework for self-supervised visual representation learning.
→ Jupyter notebook exercise

 Code Notebook 1 · Task B

Exploring Scaling Behaviour of LMs with Pythia Models: Empirical investigation of how language model performance scales with model size using the Pythia model series.
→ Jupyter notebook exercise

CS-461

Foundation Models and Generative AI

Have a great week!