

What do We Need to Learn Before Burning the Next One Million GPU Hours?

Iz Beltagy

Allen Institute for AI

Julien Launay

LightOn

What we did for the past year in architecture & scaling...



**What Language Model to Train if
You Have One Million GPU
Hours?
Le Scao et al. (2022).**



**What Language Model Architecture and
Pretraining Objective Work Best for Zero-Shot
Generalization?
Wang et al. (ICML 2022).**

**...but this talk is about what we
learned and what open questions we
still need to answer?**

Overview

Evaluation

Architecture and Pretraining Objective

Scaling

Datasets

Engineering

Efficient Pretraining

Overview

Evaluation

Architecture and Pretraining Objective

Scaling

Datasets

Engineering

Efficient Pretraining

Evaluation

What evaluation setup to use to make modeling decisions?

Many settings

- LM loss vs. downstream
- Zero-shot vs. few-shot
- Prompting vs. finetuning
- Parameter-efficient vs. full finetuning
- With/without multi-task finetuning

Evaluation

What evaluation setup to use to make modeling decisions?

Many settings **and results don't necessarily transfer**

- LM loss vs. downstream (Tay et al., 2021, Abnar et. al., 2021)
- Zero-shot vs. few-shot
- Prompting vs. finetuning (Wang et al., 2022)
- Parameter-efficient vs. full finetuning
- With/without multi-task finetuning (Wang et al., 2022)

Tay et. al., 2021: Scale Efficiently: Insights From Pre-trained and Fine-tuned Transformers

Abnar et. al., 2021: Exploring the Limits of Large Scale Pre-training

Wang et. al., 2022: What Language Model Architecture and Pretraining Objective Work Best for Zero-Shot Generalization?

Evaluation

What evaluation setup to use to make modeling decisions?

Many settings and results don't necessarily transfer

- LM loss vs. **downstream** (Tay et al., 2022, Abnar et. al., 2021)
- **Zero-shot** vs. few-shot
- **Prompting** vs. finetuning (Wang et al., 2022)
- ~~Parameter-efficient vs. full finetuning~~
- **With/without multi-task finetuning** (Wang et al., 2022)

Tay et. al., 2021: Scale Efficiently: Insights From Pre-trained and Fine-tuned Transformers

Abnar et. al., 2021: Exploring the Limits of Large Scale Pre-training

Wang et. al., 2022: What Language Model Architecture and Pretraining Objective Work Best for Zero-Shot Generalization?

Evaluation

What evaluation setup to use to make modeling decisions?

Many settings and results don't necessarily transfer

- LM loss vs. **downstream** (Tay et al., 2022, Abnar et. al., 2021)
 - **Zero-shot** vs. few-shot
 - **Prompting** vs. finetuning (Wang et al., 2022)
 - ~~Parameter-efficient vs. full finetuning~~
 - **With/without multi-task finetuning** (Wang et al., 2022)
- Fast and easy to run
 - Representative of how the model will be used in practice

Tay et. al., 2021: Scale Efficiently: Insights From Pre-trained and Fine-tuned Transformers

Abnar et. al., 2021: Exploring the Limits of Large Scale Pre-training

Wang et. al., 2022: What Language Model Architecture and Pretraining Objective Work Best for Zero-Shot Generalization?

What evaluation setup to use to make modeling decisions?

Open question

- Better understanding of the relation between different setup, and why the results differ

Overview

Evaluation

Architecture and Pretraining Objective

Scaling

Datasets

Engineering

Efficient Pretraining

Architecture and Pretraining Objective

Architecture and Pretraining Objective

What is a model architecture and pretraining objective that work for all settings?

Architecture and Pretraining Objective

What is a model architecture and pretraining objective that work for all settings?

Architecture and Pretraining Objective

What is a model architecture and pretraining objective that work for all settings?

Wang et al., 2022: zero-shot generalization

Architecture and Pretraining Objective

What is a model architecture and pretraining objective that work for all settings?

Wang et al., 2022: zero-shot generalization

- GPT2-style: after pretraining

Architecture and Pretraining Objective

What is a model architecture and pretraining objective that work for all settings?

Wang et al., 2022: zero-shot generalization

- GPT2-style: after pretraining
- T5-style: after pretraining + multitask finetuning (e.g. T0 adaptation)

Architecture and Pretraining Objective

What is a model architecture and pretraining objective that work for all settings?

Wang et al., 2022: zero-shot generalization

- GPT2-style: after pretraining
- T5-style: after pretraining + multitask finetuning (e.g. T0 adaptation)

Architecture and Pretraining Objective

What is a model architecture and pretraining objective that work for all settings?

Wang et al., 2022: zero-shot generalization

- GPT2-style: after pretraining
- T5-style: after pretraining + multitask finetuning (e.g. T0 adaptation)

Architecture and Pretraining Objective

What is a model architecture and pretraining objective that work for all settings?

Wang et al., 2022: zero-shot generalization

- GPT2-style: after pretraining
- T5-style: after pretraining + multitask finetuning (e.g. T0 adaptation)

Architecture and Pretraining Objective

What is a model architecture and pretraining objective that work for all settings?

Wang et al., 2022: zero-shot generalization

- GPT2-style: after pretraining
- T5-style: after pretraining + multitask finetuning (e.g. T0 adaptation)

Architecture and Pretraining Objective

What is a model architecture and pretraining objective that work for all settings?

Wang et al., 2022: zero-shot generalization

- GPT2-style: after pretraining
- T5-style: after pretraining + multitask finetuning (e.g. T0 adaptation)

Wang et. al., 2022: What Language Model Architecture and Pretraining Objective Work Best for Zero-Shot Generalization?

Architectures

Architectures

Architectures:

Architectures

Architectures:

- Vanilla transformer

Architectures

Architectures:

- Vanilla transformer
 - Decoder-only

Architectures

Architectures:

- Vanilla transformer
 - Decoder-only
 - Encoder-decoder

Architectures

Architectures:

- Vanilla transformer
 - Decoder-only
 - Encoder-decoder
- Transformer variants

Architectures

Architectures:

- Vanilla transformer
 - Decoder-only
 - Encoder-decoder
- Transformer variants

Architectures

Architectures:

- Vanilla transformer
 - Decoder-only
 - Encoder-decoder
- Transformer variants

Narang et. al.,: most transformer variants are comparable to vanilla transformers.

Architectures

Architectures:

- Vanilla transformer
 - Decoder-only
 - Encoder-decoder
- Transformer variants

Narang et. al.,: most transformer variants are comparable to vanilla transformers.

Architectures

Architectures:

- Vanilla transformer
 - Decoder-only
 - Encoder-decoder
- Transformer variants

Narang et. al.,: most transformer variants are comparable to vanilla transformers.

Mixture of Experts (**MoEs**) or Conditionally-Activated Transformers:

Architectures

Architectures:

- Vanilla transformer
 - Decoder-only
 - Encoder-decoder
- Transformer variants

Narang et. al.,: most transformer variants are comparable to vanilla transformers.

Mixture of Experts (**MoEs**) or Conditionally-Activated Transformers:

- Different examples activate different parts (“experts”) of the model

Architectures

Architectures:

- Vanilla transformer
 - Decoder-only
 - Encoder-decoder
- Transformer variants

Narang et. al.,: most transformer variants are comparable to vanilla transformers.

Mixture of Experts (**MoEs**) or Conditionally-Activated Transformers:

- Different examples activate different parts (“experts”) of the model
- Examples: Switch Transformer, Base-layer, DEMix Layers, GLaM, LaMDa

Architectures

Architectures:

- Vanilla transformer
 - Decoder-only
 - Encoder-decoder
- Transformer variants

Narang et. al.,: most transformer variants are comparable to vanilla transformers.

Mixture of Experts (**MoEs**) or Conditionally-Activated Transformers:

- Different examples activate different parts (“experts”) of the model
- Examples: Switch Transformer, Base-layer, DEMix Layers, GLaM, LaMDa
- Interesting but haven’t seen the same adoption as dense models

Architectures

Architectures:

- Vanilla transformer
 - Decoder-only
 - Encoder-decoder
- Transformer variants

Narang et. al.,: most transformer variants are comparable to vanilla transformers.

Mixture of Experts (**MoEs**) or Conditionally-Activated Transformers:

- Different examples activate different parts (“experts”) of the model
- Examples: Switch Transformer, Base-layer, DEMix Layers, GLaM, LaMDa
- Interesting but haven’t seen the same adoption as dense models

Architectures

Architectures:

- Vanilla transformer
 - Decoder-only
 - Encoder-decoder
- Transformer variants

Narang et. al.,: most transformer variants are comparable to vanilla transformers.

Mixture of Experts (**MoEs**) or Conditionally-Activated Transformers:

- Different examples activate different parts (“experts”) of the model
- Examples: Switch Transformer, Base-layer, DEMix Layers, GLaM, LaMDa
- Interesting but haven’t seen the same adoption as dense models

Narang et. al., 2021: Do Transformer Modifications Transfer Across Implementations and Applications?

Pretraining Objective

Pretraining Objective

Mixture of objectives

Pretraining Objective

Mixture of objectives

- UL2: Tay et. al., 2022

Pretraining Objective

Mixture of objectives

- UL2: Tay et. al., 2022
- CM3: Aghajanyan et. al., 2022

Pretraining Objective

Mixture of objectives

- UL2: Tay et. al., 2022
- CM3: Aghajanyan et. al., 2022
- Argue that the objectives are complementary \Rightarrow train on a mix of them

Pretraining Objective

Mixture of objectives

- UL2: Tay et. al., 2022
- CM3: Aghajanyan et. al., 2022
- Argue that the objectives are complementary \Rightarrow train on a mix of them

Pretraining Objective

Mixture of objectives

- UL2: Tay et. al., 2022
- CM3: Aghajanyan et. al., 2022
- Argue that the objectives are complementary \Rightarrow train on a mix of them

Pretraining Objective

Mixture of objectives

- UL2: Tay et. al., 2022
- CM3: Aghajanyan et. al., 2022
- Argue that the objectives are complementary \Rightarrow train on a mix of them

Pretraining Objective

Mixture of objectives

- UL2: Tay et. al., 2022
- CM3: Aghajanyan et. al., 2022
- Argue that the objectives are complementary \Rightarrow train on a mix of them

Tay et. al., 2022: Unifying Language Learning Paradigms

Pretraining Objective

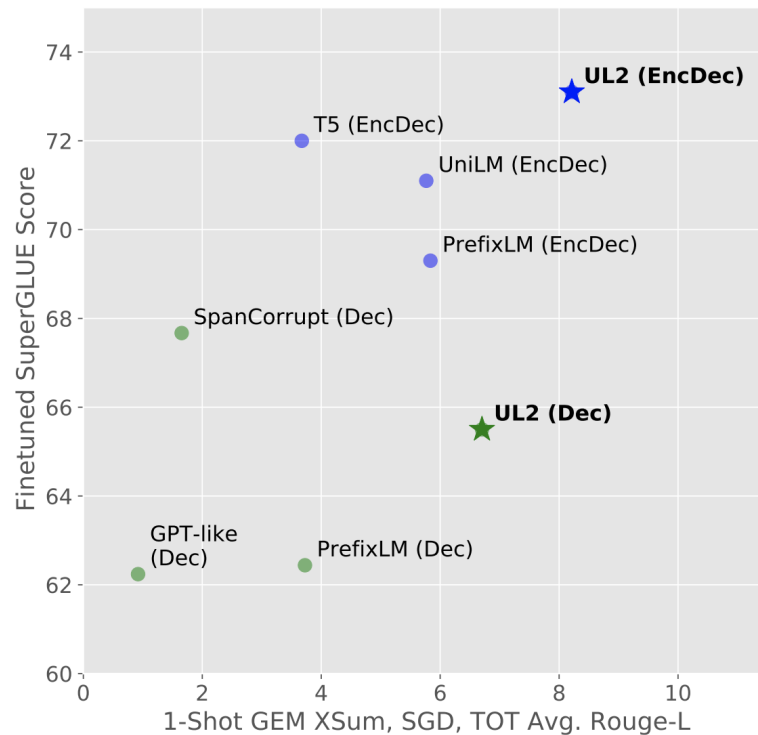
Mixture of objectives

- UL2: Tay et. al., 2022
- CM3: Aghajanyan et. al., 2022
- Argue that the objectives are complementary \Rightarrow train on a mix of them

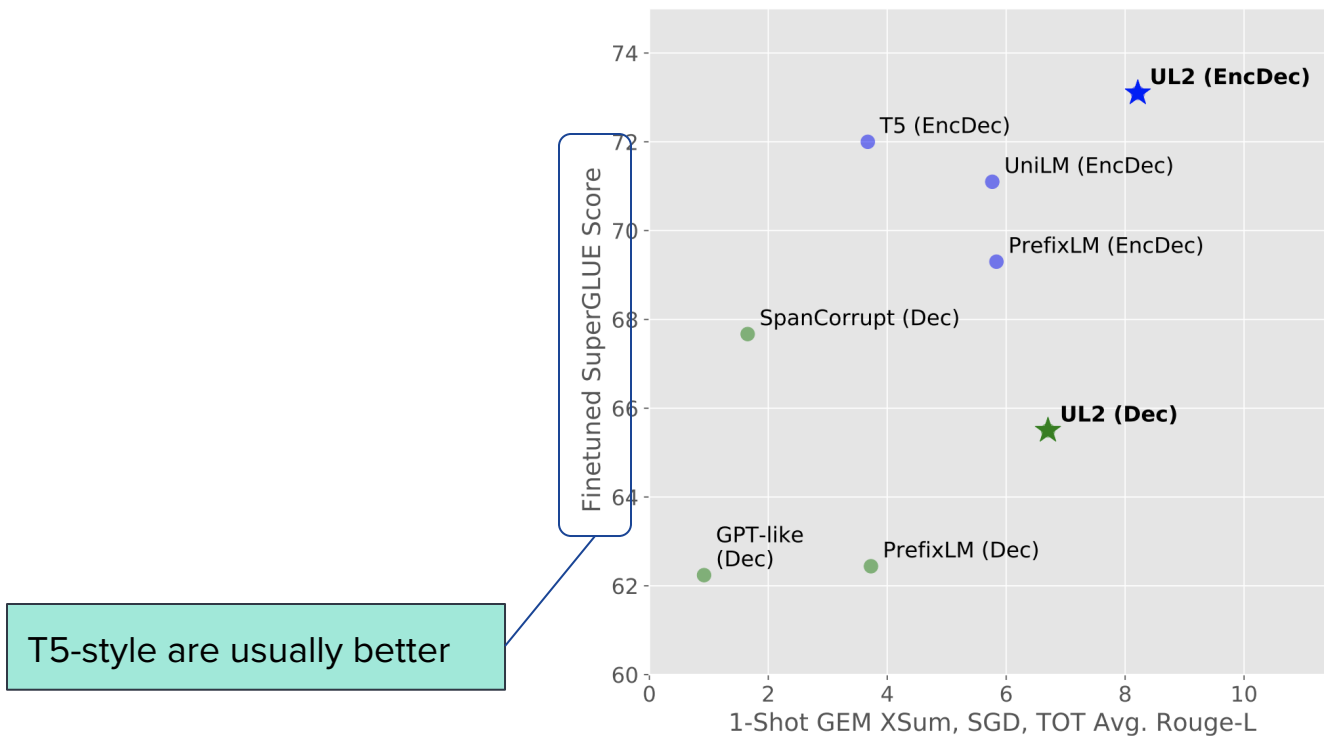
Tay et. al., 2022: Unifying Language Learning Paradigms

Aghajanyan et. al., 2022: CM3: A Causal Masked Multimodal Model of the Internet

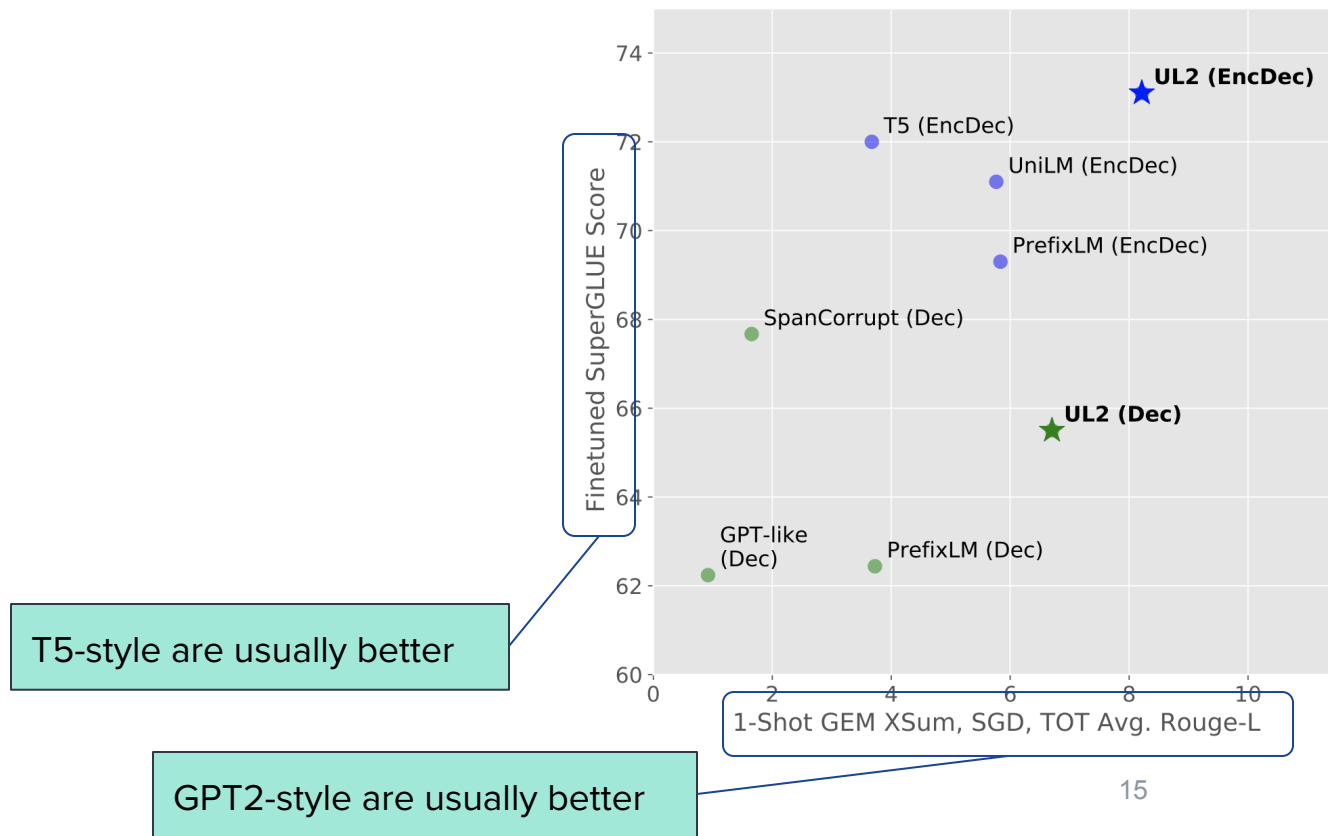
Tay et. al., 2022: Unifying Language Learning Paradigms



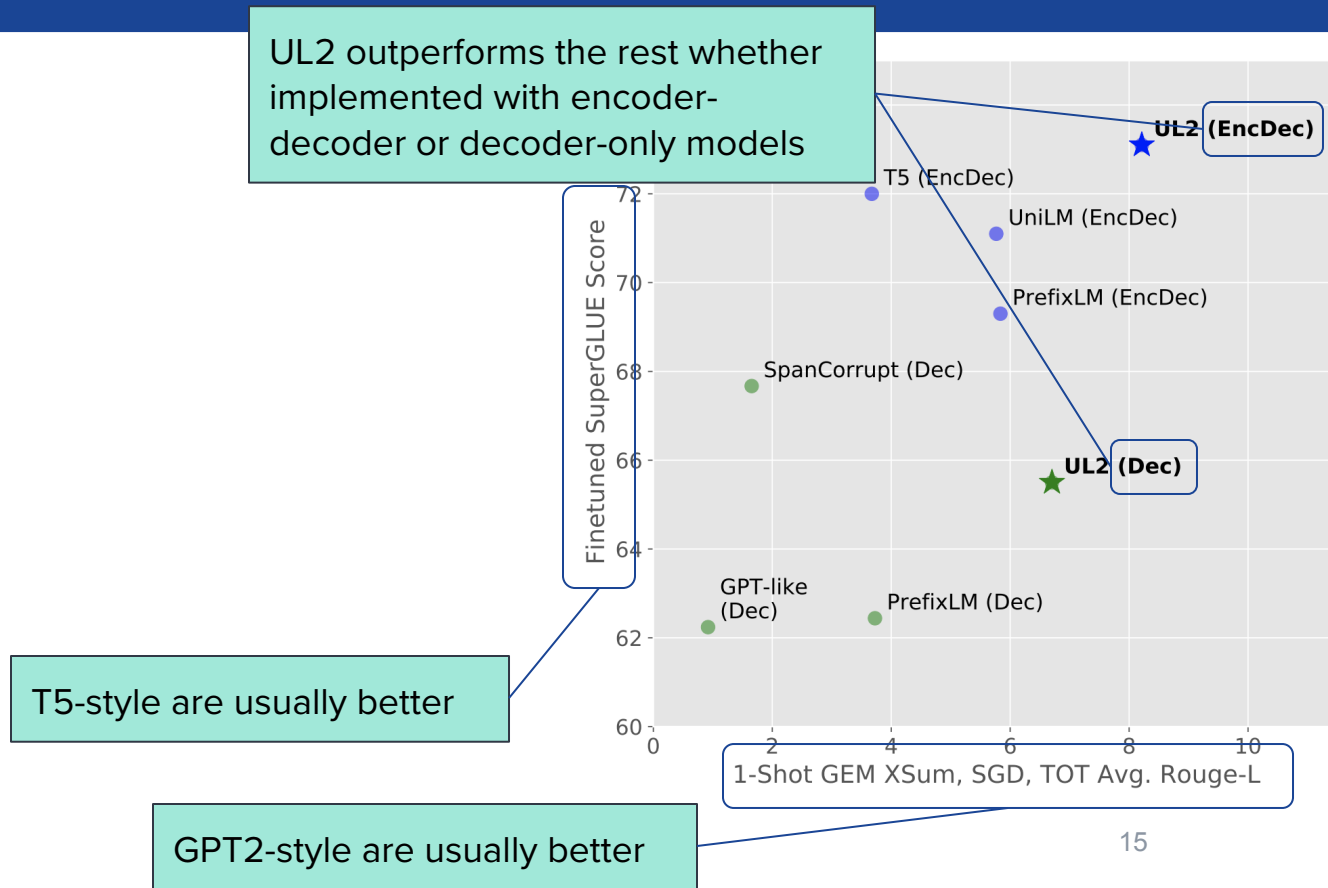
Tay et. al., 2022: Unifying Language Learning Paradigms



Tay et. al., 2022: Unifying Language Learning Paradigms



Tay et. al., 2022: Unifying Language Learning Paradigms



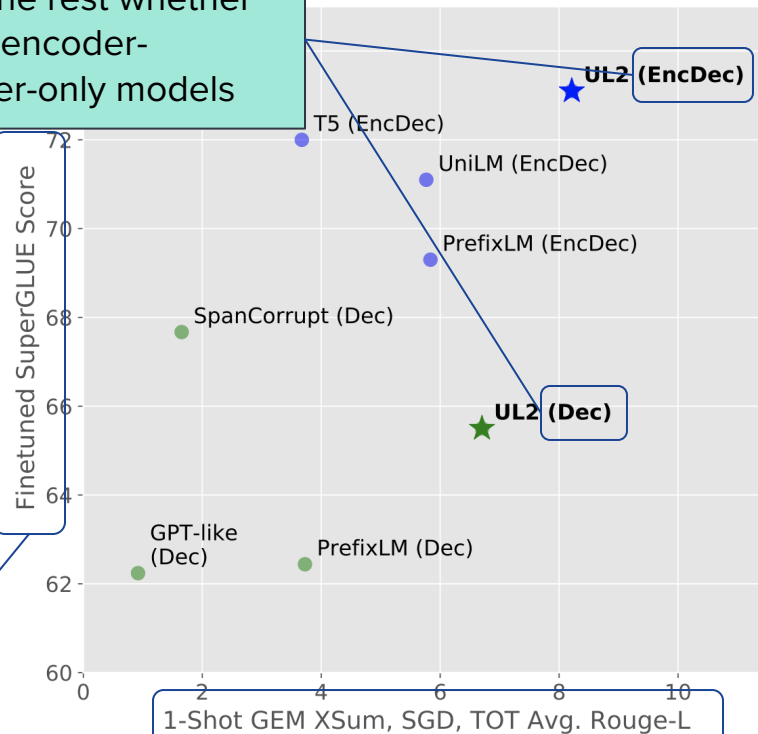
Tay et. al., 2022: Unifying Language Learning Paradigms

Limitations:

UL2 outperforms the rest whether implemented with encoder-decoder or decoder-only models

T5-style are usually better

GPT2-style are usually better



Tay et. al., 2022: Unifying Language Learning Paradigms

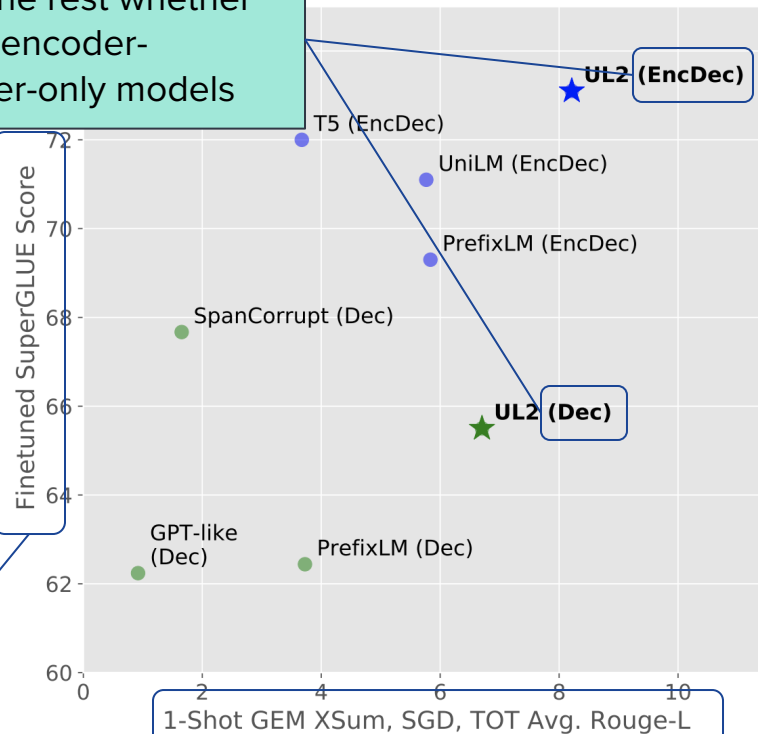
Limitations:

- Missing the zero-shot evaluation with/without MT-F [as in Wang et. al., discussed earlier]

UL2 outperforms the rest whether implemented with encoder-decoder or decoder-only models

T5-style are usually better

GPT2-style are usually better



Tay et. al., 2022: Unifying Language Learning Paradigms

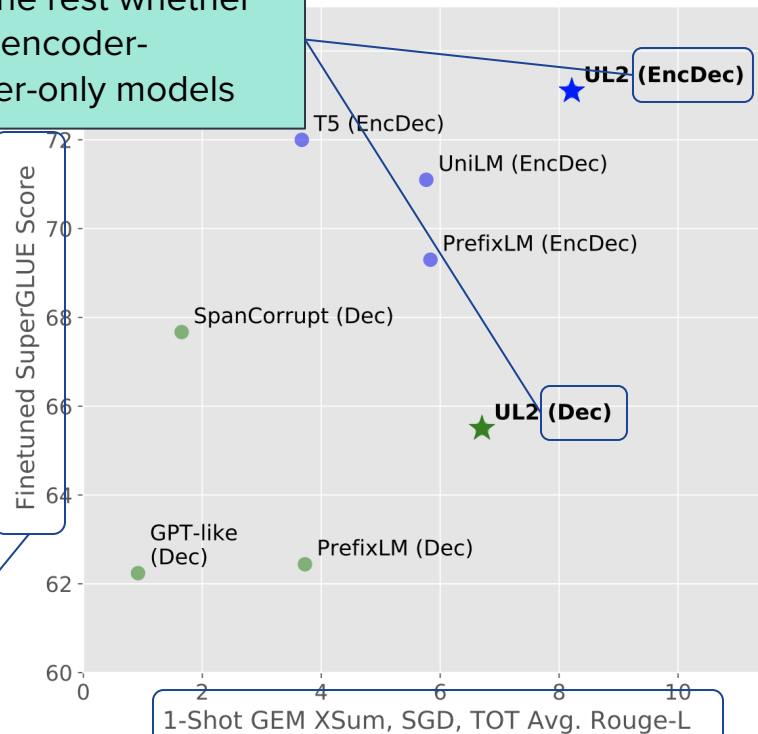
Limitations:

- Missing the zero-shot evaluation with/without MT-F [as in Wang et. al., discussed earlier]
- Ignored causal LM
 - Claimed that Prefix LM is generally better

UL2 outperforms the rest whether implemented with encoder-decoder or decoder-only models

T5-style are usually better

GPT2-style are usually better



Architecture and Pretraining Objective

Architecture and Pretraining Objective

What is a model architecture and pretraining objective that work for all settings?

Architecture and Pretraining Objective

What is a model architecture and pretraining objective that work for all settings?

Architecture and Pretraining Objective

What is a model architecture and pretraining objective that work for all settings?

We don't know, but we are getting closer

Overview

Evaluation

Architecture and Pretraining Objective

Scaling

Datasets

Engineering

Efficient Pretraining

Scaling - Kaplan et. al., 2020 vs. Hoffmann et. al., 2022

Scaling - Kaplan et. al., 2020 vs. Hoffmann et. al., 2022

Scaling - Kaplan et. al., 2020 vs. Hoffmann et. al., 2022

Kaplan et. al.: grow model size **much faster** than tokens (e.g, GPT3, OPT, PaLM)

Scaling - Kaplan et. al., 2020 vs. Hoffmann et. al., 2022

Kaplan et. al.: grow model size **much faster** than tokens (e.g, GPT3, OPT, PaLM)

- Given 10x compute, increase N by 5.5x, and D by 1.8x

Scaling - Kaplan et. al., 2020 vs. Hoffmann et. al., 2022

Kaplan et. al.: grow model size **much faster** than tokens (e.g, GPT3, OPT, PaLM)

- Given 10x compute, increase N by 5.5x, and D by 1.8x

Scaling - Kaplan et. al., 2020 vs. Hoffmann et. al., 2022

Kaplan et. al.,: grow model size **much faster** than tokens (e.g, GPT3, OPT, PaLM)

- Given 10x compute, increase N by 5.5x, and D by 1.8x

Hoffmann et. al.,: grow model size and number of token **at the same rate** (e.g, Chinchilla)

Scaling - Kaplan et. al., 2020 vs. Hoffmann et. al., 2022

Kaplan et. al.,: grow model size **much faster** than tokens (e.g, GPT3, OPT, PaLM)

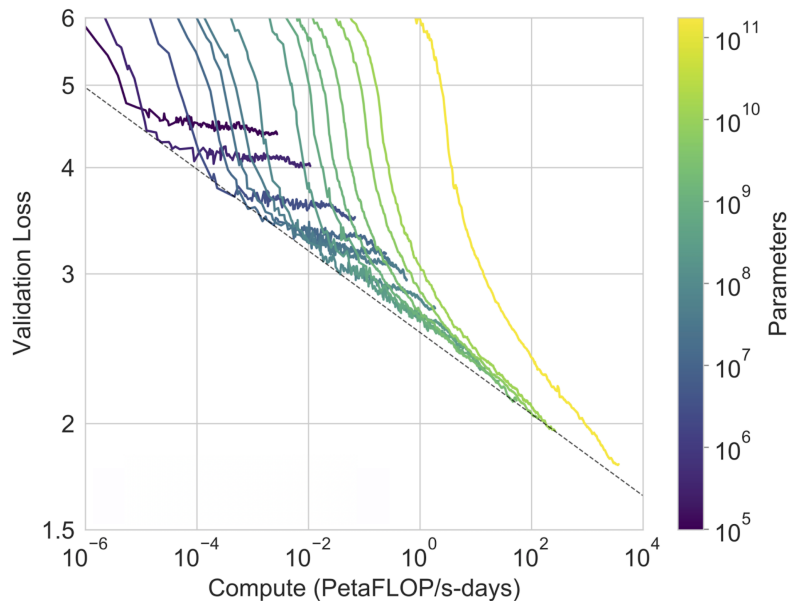
- Given 10x compute, increase N by 5.5x, and D by 1.8x

Hoffmann et. al.,: grow model size and number of token **at the same rate** (e.g, Chinchilla)

- Given 10x compute, grow N by 3.2x and D by 3.2x

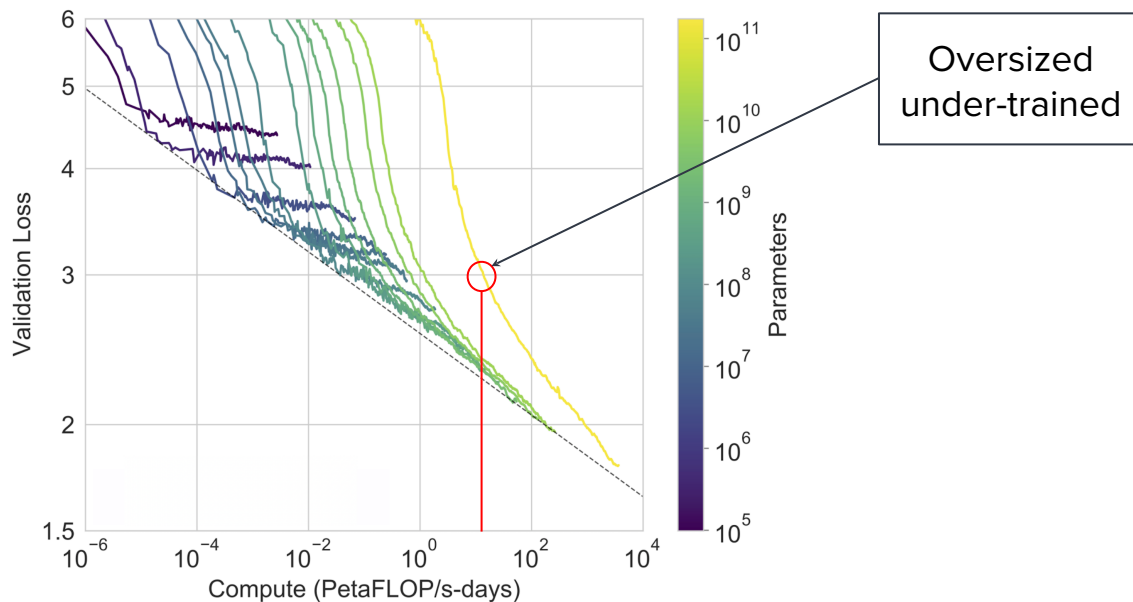
Scaling - Kaplan et. al., 2020 vs. Hoffmann et. al., 2022

- We trained (and are training) “oversized” and “under-trained” models



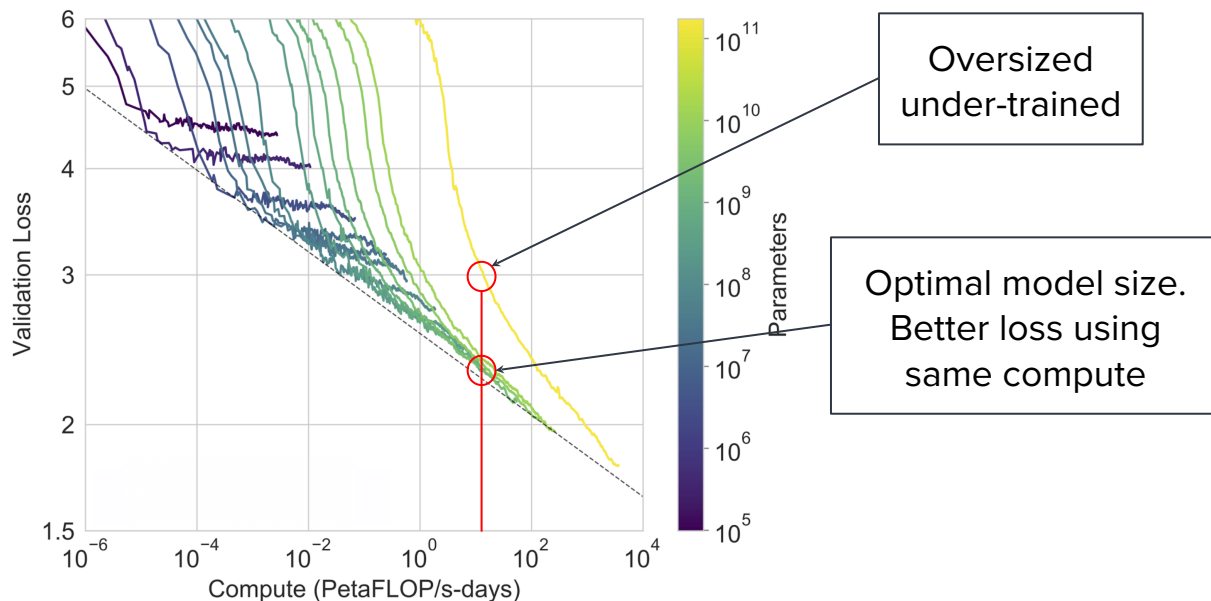
Scaling - Kaplan et. al., 2020 vs. Hoffmann et. al., 2022

- We trained (and are training) “oversized” and “under-trained” models



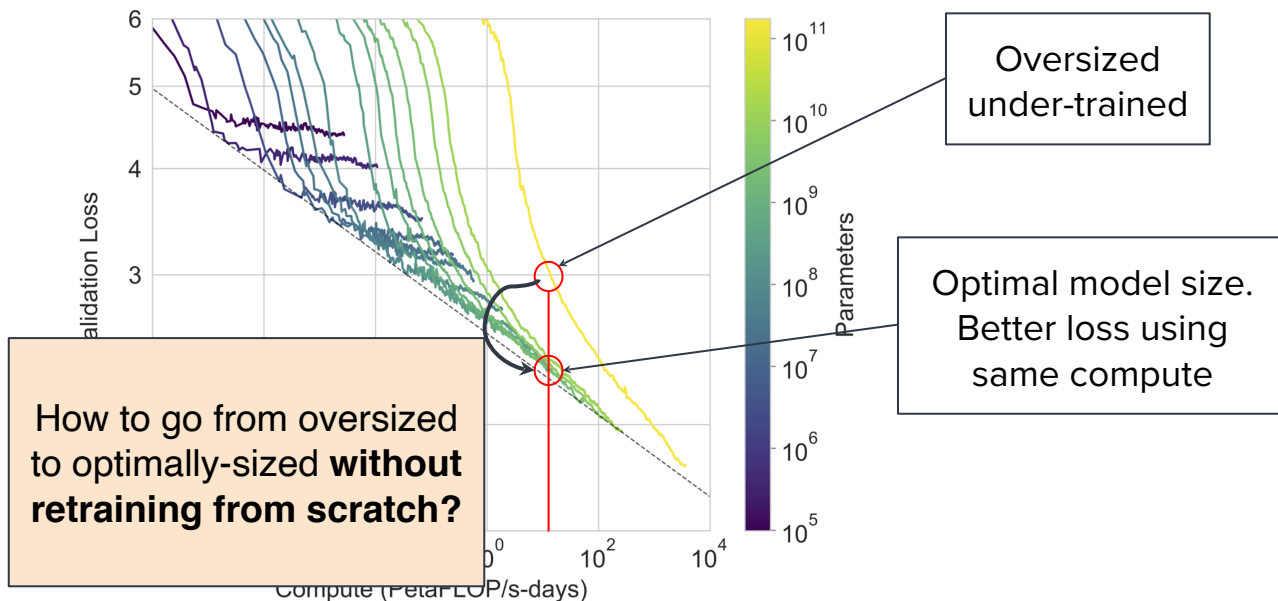
Scaling - Kaplan et. al., 2020 vs. Hoffmann et. al., 2022

- We trained (and are training) “oversized” and “under-trained” models



Scaling - Kaplan et. al., 2020 vs. Hoffmann et. al., 2022

- We trained (and are training) “oversized” and “under-trained” models



Chonk' me up Scotty

For the next generation of LLMs, we will need to scale...

Overview

Evaluation

Architecture and Pretraining Objective

Scaling

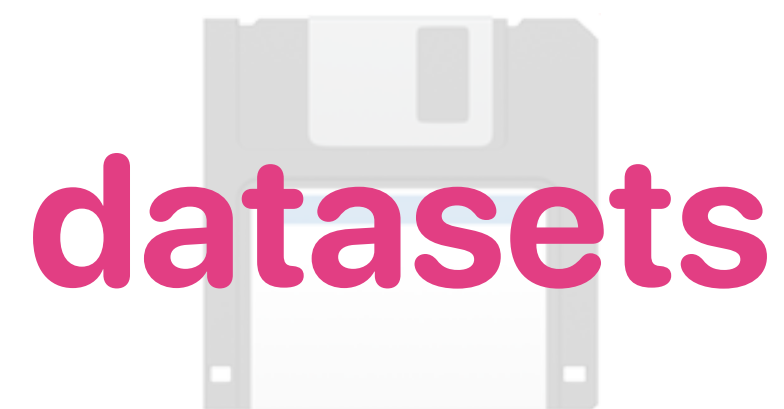
Datasets

Engineering

Efficient Pretraining

Chonk' me up Scotty

For the next generation of LLMs, we will need to scale...



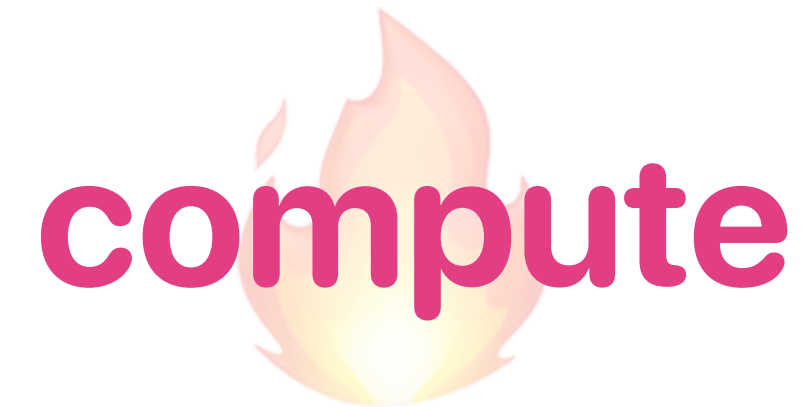
quality at scale

Chonk' me up Scotty

For the next generation of LLMs, we will need to scale...



quality at scale



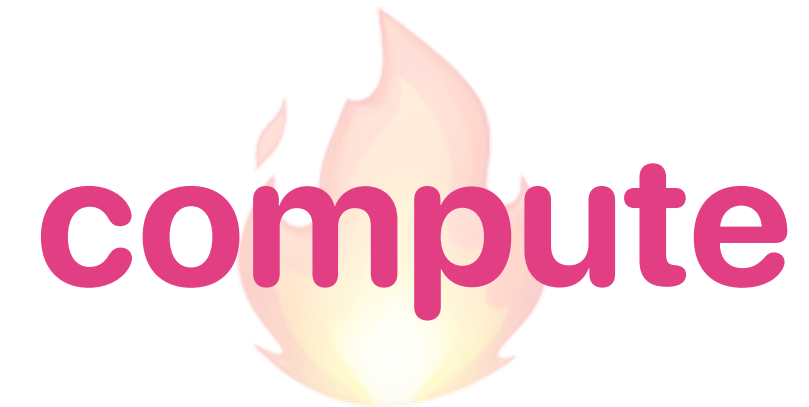
engineering challenges

Chonk' me up Scotty

For the next generation of LLMs, we will need to scale...



quality at scale



engineering challenges



accelerate scaling

Model quality is all about **data** quality

Training data matters a lot!
(more than most modeling choices?)

Model quality is all about data quality

Training data matters a lot!
(more than most modeling choices?)

Aggregated performance on EAI harness

| Model | Parameters | Pretraining tokens | | | |
|----------------------|------------|--------------------|--------------|-------|--------------|
| | | Dataset | 112B | 250B | 300B |
| OpenAI — Curie | 6.7B | | | | <u>49.28</u> |
| OpenAI — Babbage | 1.3B | | | | 45.30 |
| EleutherAI — GPT-Neo | 1.3B | The Pile | | | 42.94 |
| Ours | 13B | OSCAR | | | 47.09 |
| | 1.3B | The Pile | 42.79 | 43.12 | 43.46 |
| | 1.3B | C4 | 42.77 | | |
| | 1.3B | OSCAR | 41.72 | | |

Le Scao et al.,2022.

Model quality is all about data quality

Training data matters a lot!
(more than most modeling choices?)

Aggregated performance on EAI harness

| Model | Parameters | Pretraining tokens | | | |
|----------------------|------------|--------------------|--------------|-------|--------------|
| | | Dataset | 112B | 250B | 300B |
| OpenAI — Curie | 6.7B | | | | <u>49.28</u> |
| OpenAI — Babbage | 1.3B | | | | 45.30 |
| EleutherAI — GPT-Neo | 1.3B | The Pile | | | 42.94 |
| Ours | 13B | OSCAR | | | 47.09 |
| | 1.3B | The Pile | 42.79 | 43.12 | 43.46 |
| | 1.3B | C4 | 42.77 | | |
| | 1.3B | OSCAR | 41.72 | | |

Same architecture, different data:

Le Scao et al.,2022.

Model quality is all about data quality

Training data matters a lot!
(more than most modeling choices?)

Aggregated performance on EAI harness

| Model | Parameters | Pretraining tokens | | | |
|----------------------|------------|--------------------|-------|-------|-------|
| | | Dataset | 112B | 250B | 300B |
| OpenAI — Curie | 6.7B | | | | 49.28 |
| OpenAI — Babbage | 1.3B | | | | 45.30 |
| EleutherAI — GPT-Neo | 1.3B | The Pile | | | 42.94 |
| Ours | 13B | OSCAR | | | 47.09 |
| | 1.3B | The Pile | 42.79 | 43.12 | 43.46 |
| | 1.3B | C4 | 42.77 | | |
| | 1.3B | OSCAR | 41.72 | | |

Le Scao et al.,2022.

Same architecture, different data:

45.30%

OpenAI-Babbage(1.3B)

Model quality is all about data quality

Training data matters a lot!
(more than most modeling choices?)

Aggregated performance on EAI harness

| Model | Parameters | Pretraining tokens | | | |
|----------------------|------------|--------------------|-------|-------|-------|
| | | Dataset | 112B | 250B | 300B |
| OpenAI — Curie | 6.7B | | | | 49.28 |
| OpenAI — Babbage | 1.3B | | | | 45.30 |
| EleutherAI — GPT-Neo | 1.3B | The Pile | | | 42.94 |
| | 13B | OSCAR | | | 47.09 |
| Ours | 1.3B | The Pile | 42.79 | 43.12 | 43.46 |
| | 1.3B | C4 | 42.77 | | |
| | 1.3B | OSCAR | 41.72 | | |

Le Scao et al.,2022.

Same architecture, different data:

45.30%

OpenAI-Babbage(1.3B)

43.46%

Ours-1.3B@The Pile

Model quality is all about data quality

Training data matters a lot!
(more than most modeling choices?)

Aggregated performance on EAI harness

| Model | Parameters | Pretraining tokens | | | |
|----------------------|------------|--------------------|--------------|-------|--------------|
| | | Dataset | 112B | 250B | 300B |
| OpenAI — Curie | 6.7B | | | | <u>49.28</u> |
| OpenAI — Babbage | 1.3B | | | | 45.30 |
| EleutherAI — GPT-Neo | 1.3B | The Pile | | | 42.94 |
| Ours | 13B | OSCAR | | | 47.09 |
| | 1.3B | The Pile | 42.79 | 43.12 | 43.46 |
| | 1.3B | C4 | 42.77 | | |
| | 1.3B | OSCAR | 41.72 | | |

Scale can't compensate for bad data:

Le Scao et al.,2022.

Model quality is all about data quality

Training data matters a lot!
(more than most modeling choices?)

Aggregated performance on EAI harness

| Model | Parameters | Pretraining tokens | | | |
|----------------------|------------|--------------------|-------|-------|-------|
| | | Dataset | 112B | 250B | 300B |
| OpenAI — Curie | 6.7B | | | | 49.28 |
| OpenAI — Babbage | 1.3B | | | | 45.30 |
| EleutherAI — GPT-Neo | 1.3B | The Pile | | | 42.94 |
| Ours | 13B | OSCAR | | | 47.09 |
| | 1.3B | The Pile | 42.79 | 43.12 | 43.46 |
| | 1.3B | C4 | 42.77 | | |
| | 1.3B | OSCAR | 41.72 | | |

Le Scao et al.,2022.

Scale can't compensate for bad data:

49.28%

OpenAI-Curie(6.7B)

Model quality is all about data quality

Training data matters a lot!
(more than most modeling choices?)

Aggregated performance on EAI harness

| Model | Parameters | Pretraining tokens | | | |
|----------------------|------------|--------------------|-------|-------|-------|
| | | Dataset | 112B | 250B | 300B |
| OpenAI — Curie | 6.7B | | | | 49.28 |
| OpenAI — Babbage | 1.3B | | | | 45.30 |
| EleutherAI — GPT-Neo | 1.3B | The Pile | | | 42.94 |
| | 13B | OSCAR | | | 47.09 |
| Ours | 1.3B | The Pile | 42.79 | 43.12 | 43.46 |
| | 1.3B | C4 | 42.77 | | |
| | 1.3B | OSCAR | 41.72 | | |

Le Scao et al.,2022.

Scale can't compensate for bad data:

49.28%

OpenAI-Curie(6.7B)

47.09%

Ours-13B@OSCAR

We are gonna need a bigger dataset!

Bad news: we need a lot more data than expected...



We are gonna need a bigger dataset!

Bad news: we need a lot more data than expected...



Previously... [Kaplan et al., 2020](#)

176B parameters → 300B tokens

We are gonna need a **bigger** dataset!

Bad news: we need a lot more data than expected...



Previously... [Kaplan et al., 2020](#)

176B parameters → 300B tokens

Now... [Hoffmann et al., 2020](#)

isoFLOPs 50B parameters → 1000B tokens

We are gonna need a **bigger** dataset!

Bad news: we need a lot more data than expected...



Previously... [Kaplan et al., 2020](#)

176B parameters → 300B tokens

Now... [Hoffmann et al., 2020](#)

isoFLOPs 50B parameters → 1000B tokens

isoparams 176B parameters → 3700B tokens

[~1 year of CC English](#)

We are gonna need a **bigger** dataset!

Bad news: we need a lot more data than expected...



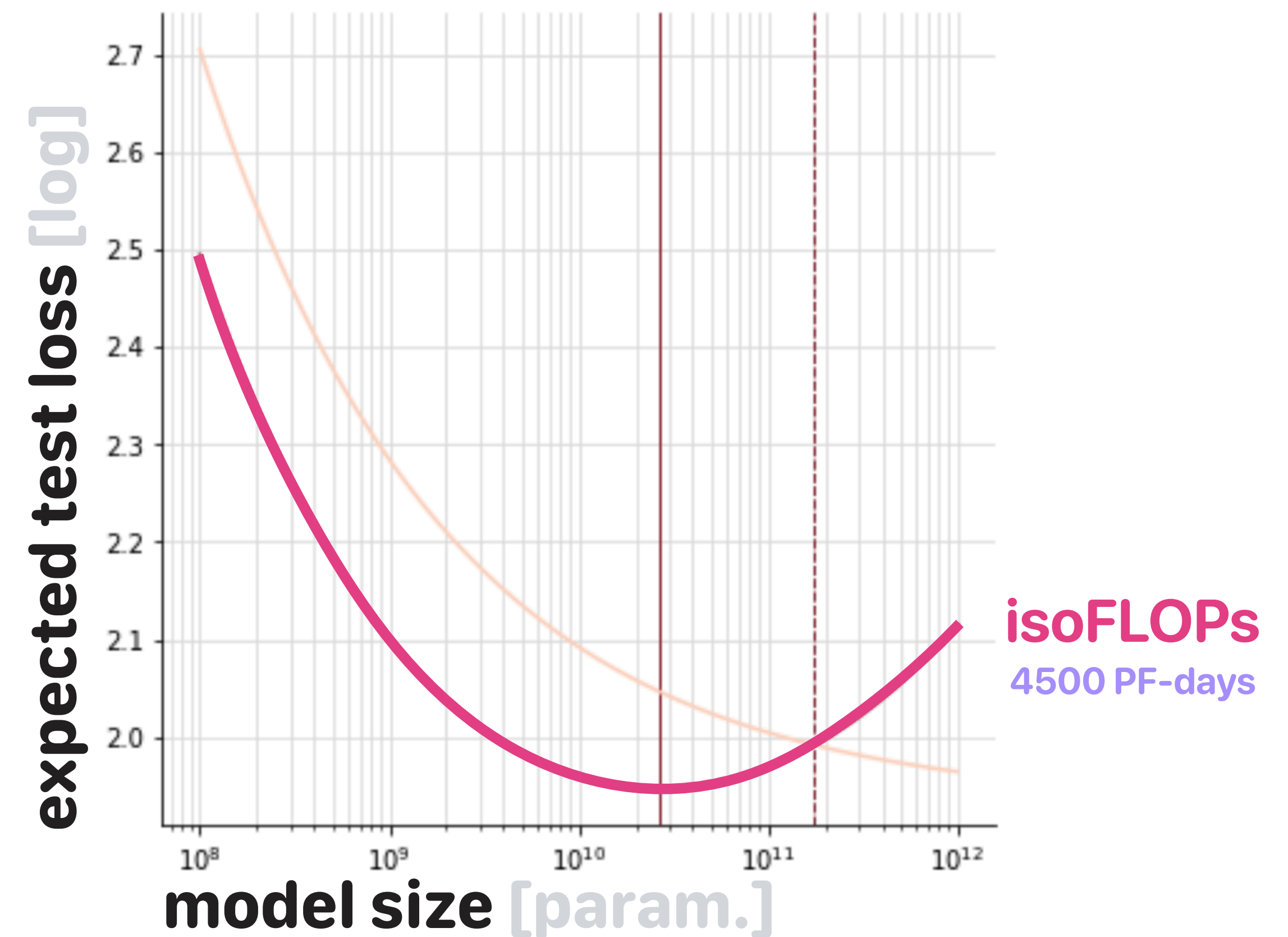
Previously... [Kaplan et al., 2020](#)

176B parameters → 300B tokens

Now... [Hoffmann et al., 2020](#)

isoFLOPs 50B parameters → 1000B tokens

isoparams 176B parameters → 3700B tokens
~1 year of CC English



We are gonna need a **bigger** dataset!

Bad news: we need a lot more data than expected...



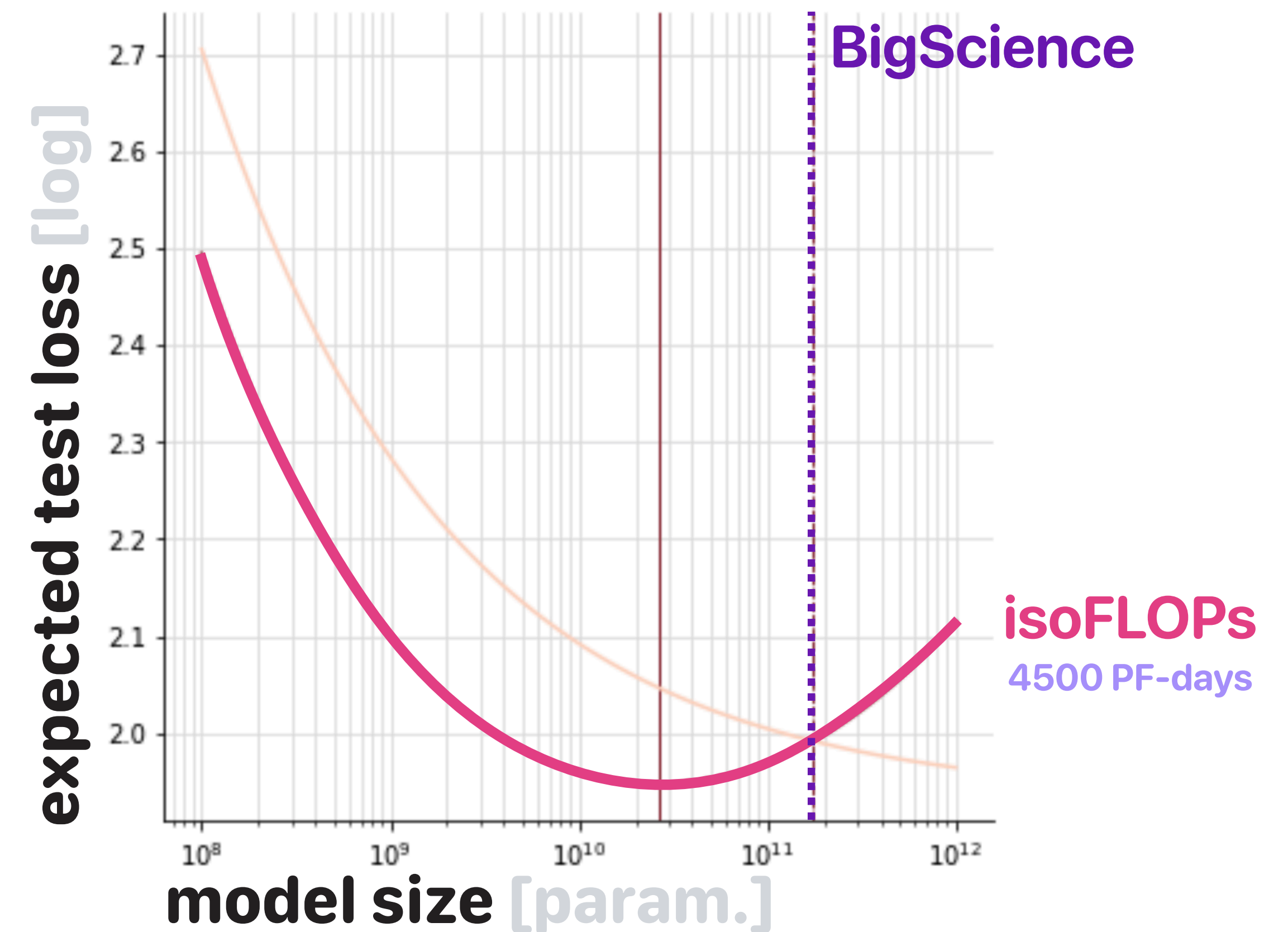
Previously... [Kaplan et al., 2020](#)

176B parameters → 300B tokens

Now... [Hoffmann et al., 2020](#)

isoFLOPs 50B parameters → 1000B tokens

isoparams 176B parameters → 3700B tokens
~1 year of CC English



We are gonna need a **bigger** dataset!

Bad news: we need a lot more data than expected...



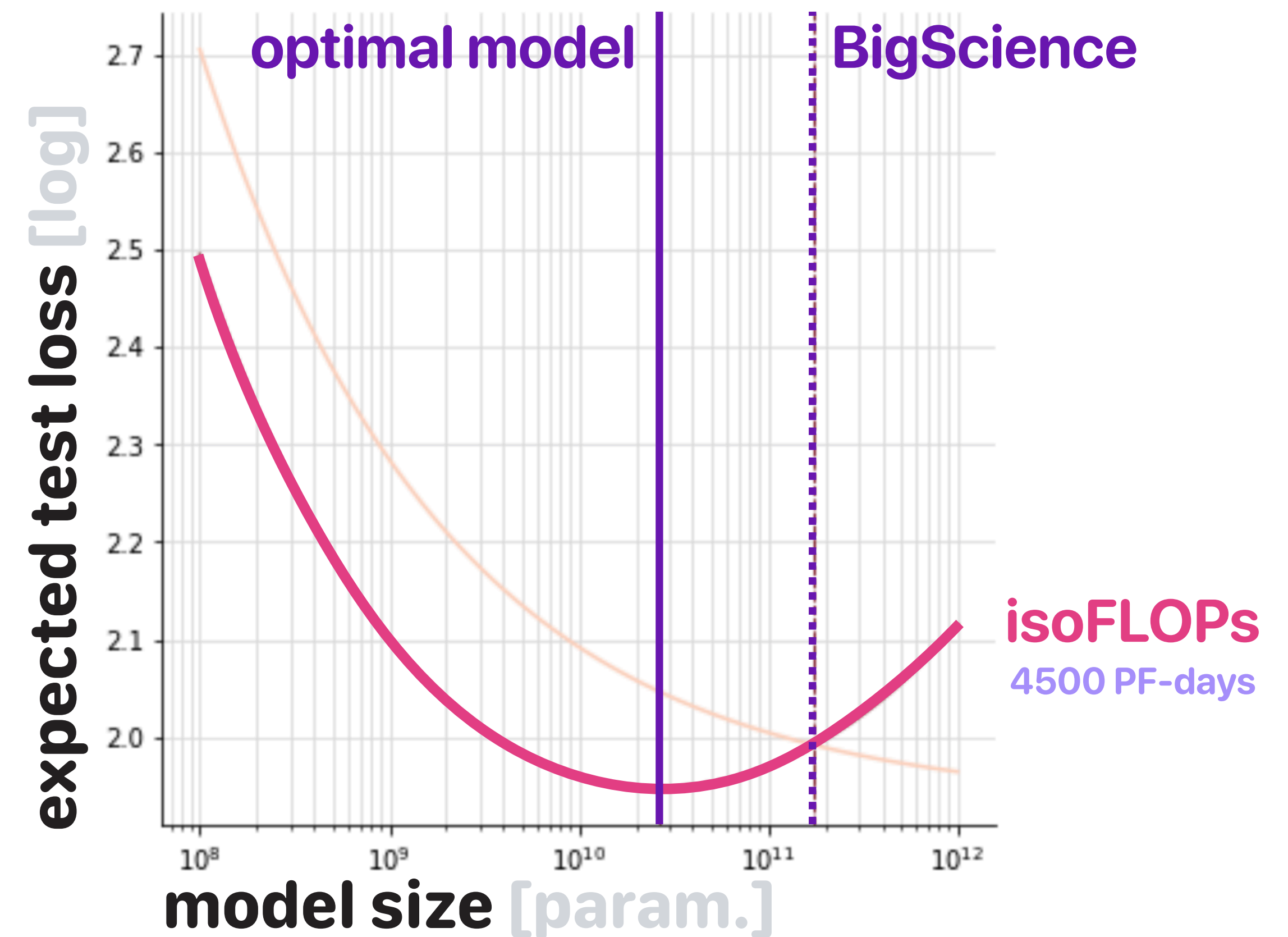
Previously... [Kaplan et al., 2020](#)

176B parameters → 300B tokens

Now... [Hoffmann et al., 2020](#)

isoFLOPs 50B parameters → 1000B tokens

isoparams 176B parameters → 3700B tokens
~1 year of CC English



We are gonna need a **bigger** dataset!

Bad news: we need a lot more data than expected...



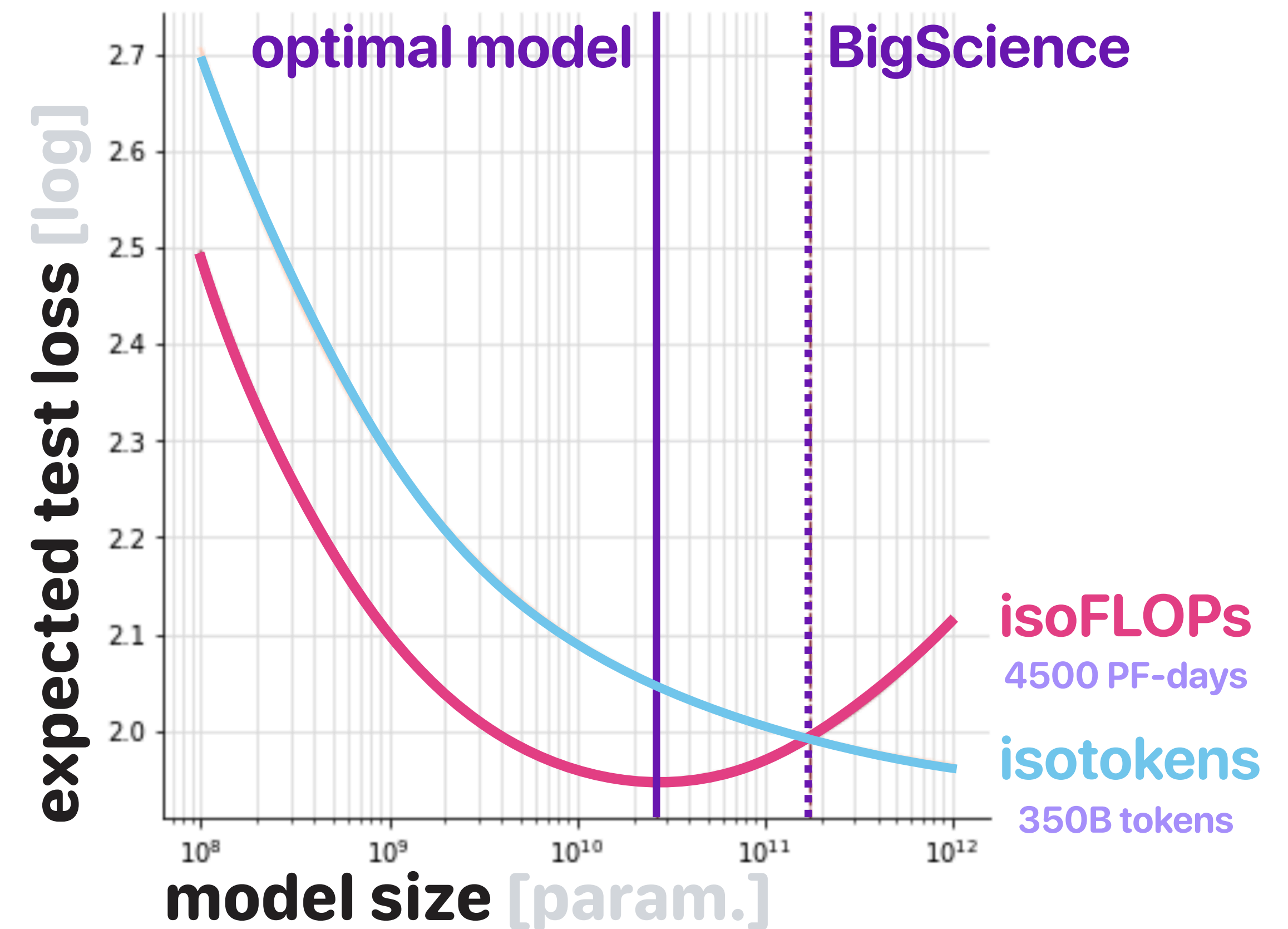
Previously... [Kaplan et al., 2020](#)

176B parameters → 300B tokens

Now... [Hoffmann et al., 2020](#)

isoFLOPs 50B parameters → 1000B tokens

isoparams 176B parameters → 3700B tokens
~1 year of CC English



We are gonna need a **bigger** dataset!

Bad news: we need a lot more data than expected...



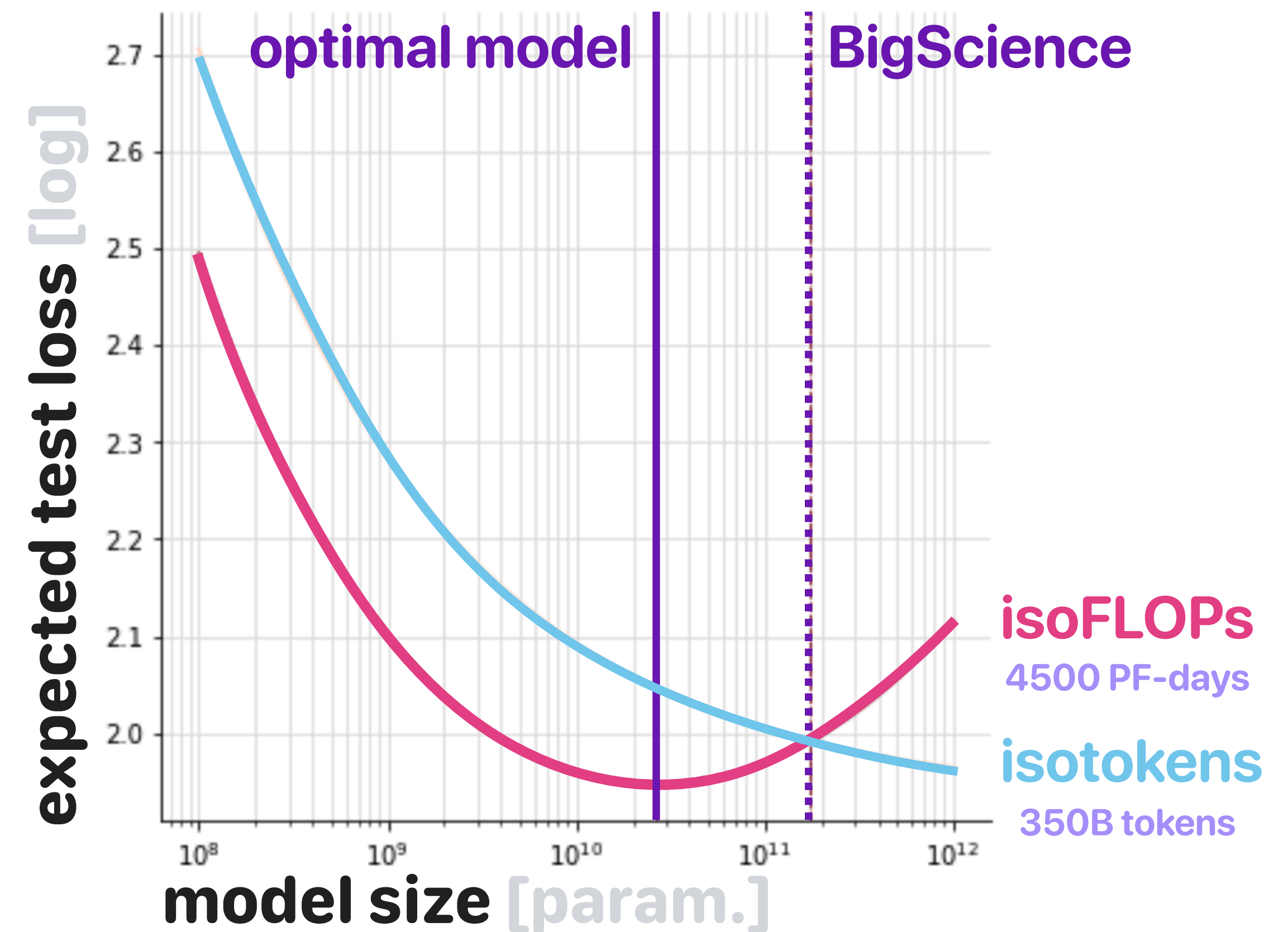
Previously... [Kaplan et al., 2020](#)

176B parameters → 300B tokens

Now... [Hoffmann et al., 2020](#)

isoFLOPs 50B parameters → 1000B tokens

isoparams 176B parameters → 3700B tokens
~1 year of CC English



Will we be data-bound instead of compute-bound?

Fantastic **training data** and where to find it



What even is **high-quality** data?

Fantastic **training data** and where to find it



What even is **high-quality** data? **technical filtering** deduplication, lack of artefacts, etc.

Fantastic training data and where to find it



What even is high-quality data?

technical filtering

deduplication, lack of artefacts, etc.

curation

diverse, cross-domain, etc.

Fantastic training data and where to find it



What even is **high-quality** data?

- technical filtering** deduplication, lack of artefacts, etc.
- curation** diverse, cross-domain, etc.

"social media conversations"

| Total dataset size = 780 billion tokens | |
|---|--------------------|
| Data source | Proportion of data |
| Social media conversations (multilingual) | 50% |
| Filtered webpages (multilingual) | 27% |
| Books (English) | 13% |
| GitHub (code) | 5% |
| Wikipedia (multilingual) | 4% |
| News (English) | 1% |

Chowdhery et al.,2022.

Fantastic training data and where to find it



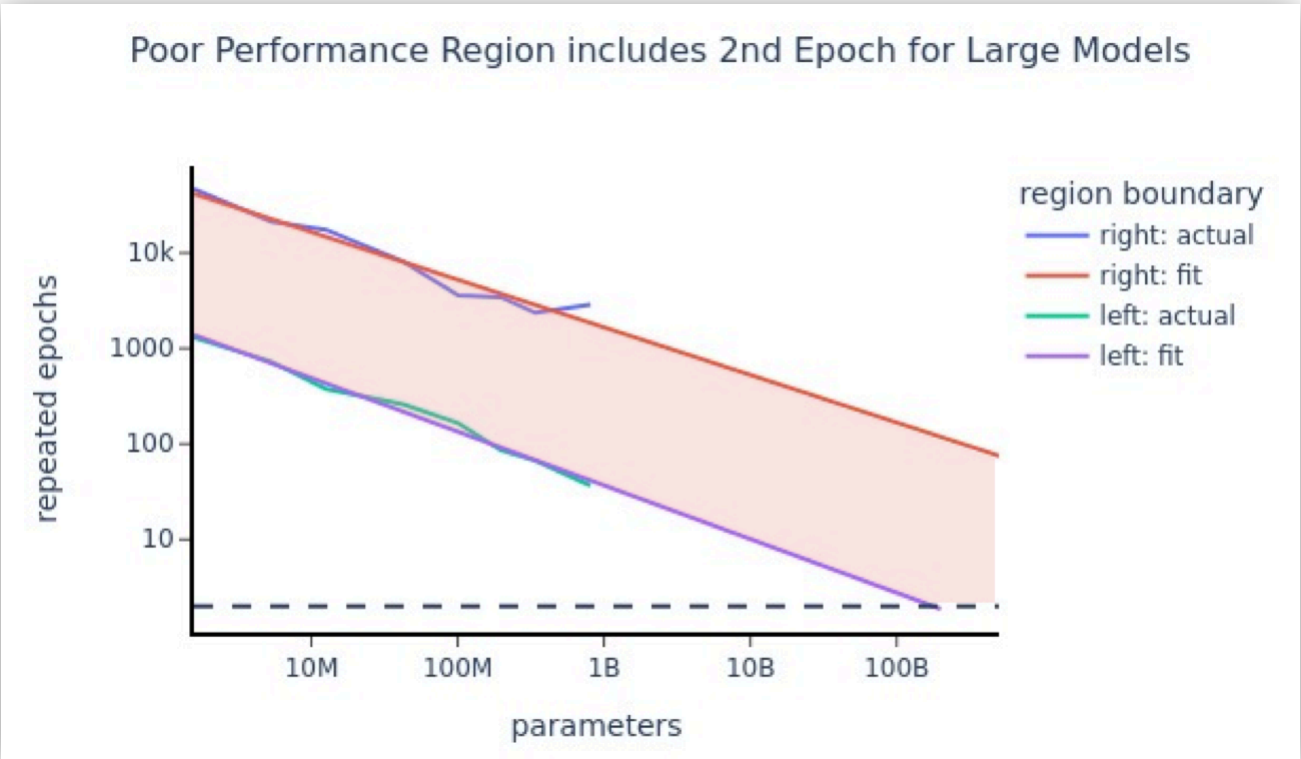
What even is **high-quality** data? **technical filtering** deduplication, lack of artefacts, etc.
curation diverse, cross-domain, etc.

"social media conversations"

| Total dataset size = 780 billion tokens | |
|---|--------------------|
| Data source | Proportion of data |
| Social media conversations (multilingual) | 50% |
| Filtered webpages (multilingual) | 27% |
| Books (English) | 13% |
| GitHub (code) | 5% |
| Wikipedia (multilingual) | 4% |
| News (English) | 1% |

Chowdhery et al.,2022.

double descent for duplication?



Hernandez et al.,2022.

Fantastic training data and where to find it

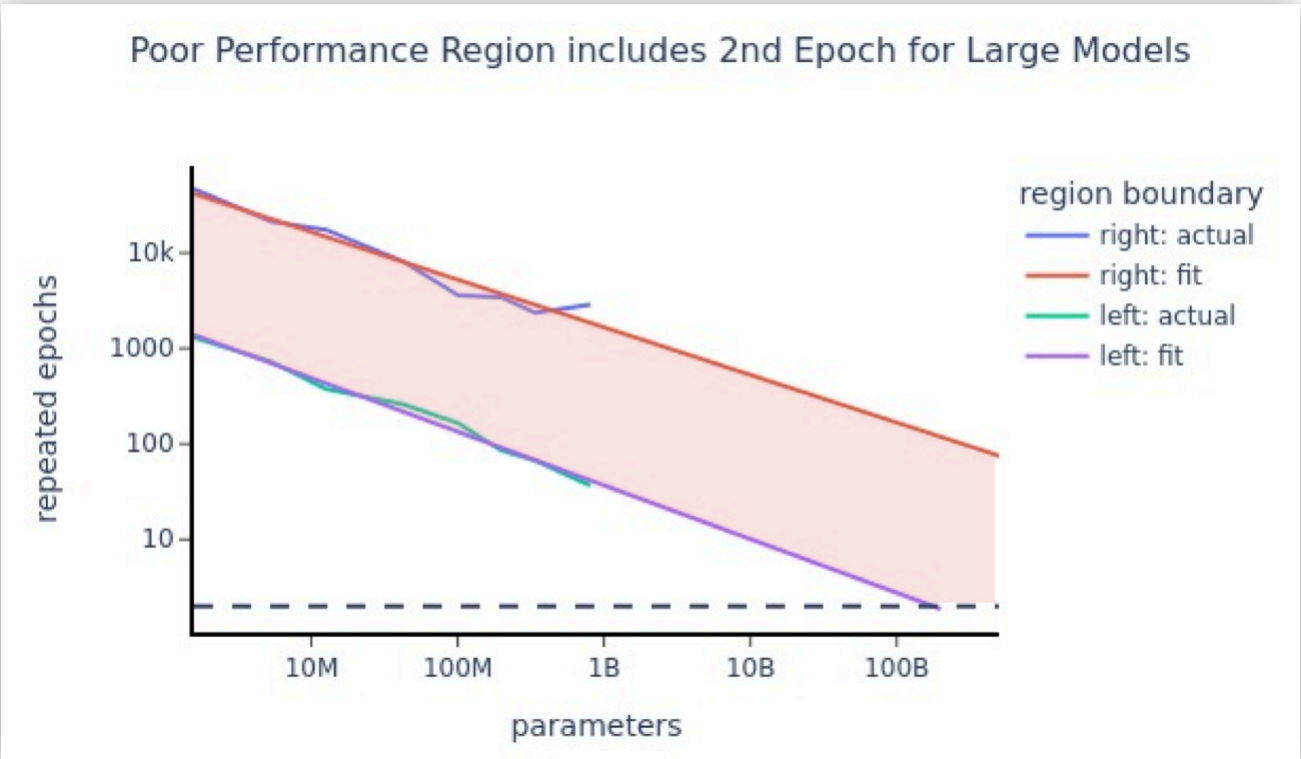
🔬 What even is **high-quality** data? **technical filtering** deduplication, lack of artefacts, etc.
curation diverse, cross-domain, etc.

“social media conversations”

| Total dataset size = 780 billion tokens | |
|---|--------------------|
| Data source | Proportion of data |
| Social media conversations (multilingual) | 50% |
| Filtered webpages (multilingual) | 27% |
| Books (English) | 13% |
| GitHub (code) | 5% |
| Wikipedia (multilingual) | 4% |
| News (English) | 1% |

Chowdhery et al.,2022.

double descent for duplication?



Hernandez et al.,2022.

💡 Currently, dataset construction is more akin to magic... Need **principled methods**!

Fantastic training data and where to find it

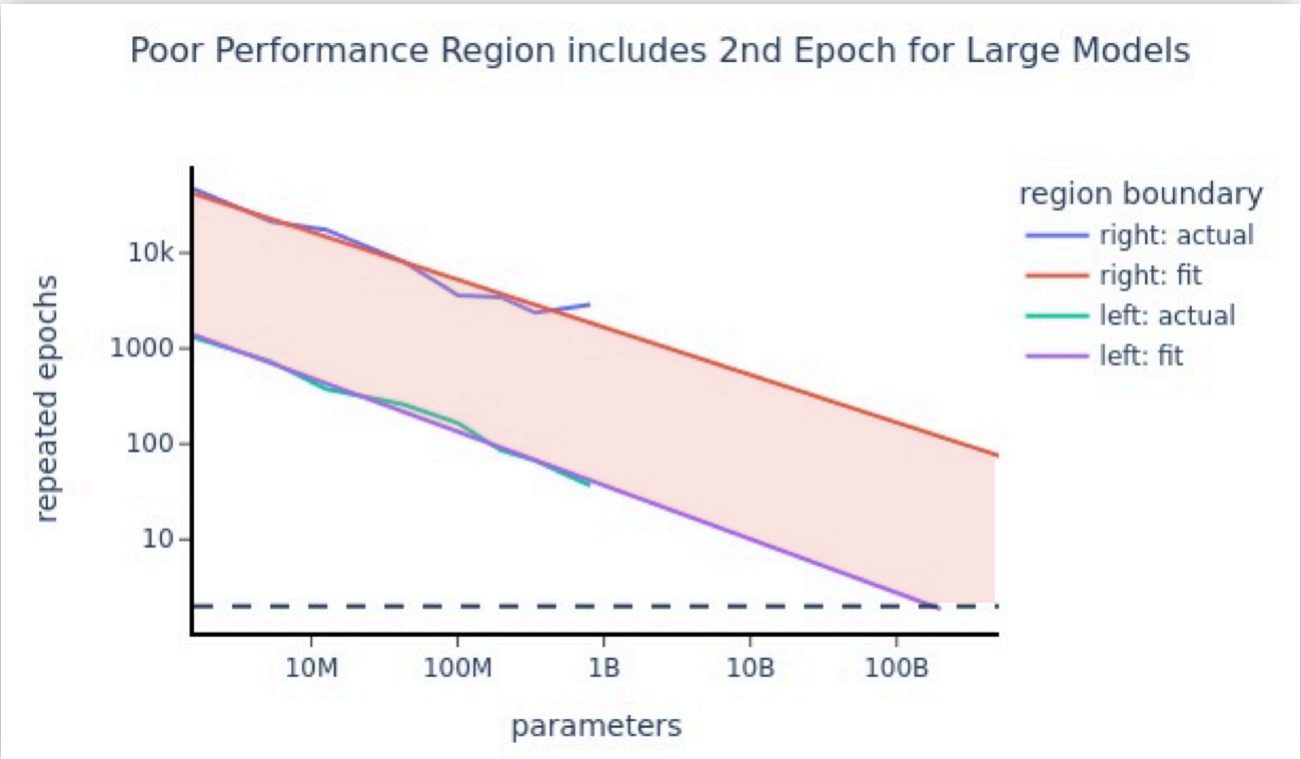
🔍 What even is **high-quality** data? **technical filtering** deduplication, lack of artefacts, etc.
curation diverse, cross-domain, etc.

“social media conversations”

| Total dataset size = 780 billion tokens | |
|---|--------------------|
| Data source | Proportion of data |
| Social media conversations (multilingual) | 50% |
| Filtered webpages (multilingual) | 27% |
| Books (English) | 13% |
| GitHub (code) | 5% |
| Wikipedia (multilingual) | 4% |
| News (English) | 1% |

Chowdhery et al.,2022.

double descent for duplication?



Hernandez et al.,2022.

💡 Currently, dataset construction is more akin to magic... Need **principled methods**!

⚠️ Emergence of **data moats** which could stand in the way of research.

Fantastic training data and where to find it

Oh and by the way...

We need this in >100 languages!

We are doing Big Science, and this comes with challenges...

We are doing **Big Science**, and this comes with challenges...

 LLMs are a true **big science** and require significant engineering efforts...
state-of-the-art HPC challenges

We are doing **Big Science**, and this comes with challenges...

 LLMs are a true **big science** and require significant engineering efforts...
state-of-the-art HPC challenges

 **Principled approaches** are very much needed:

We are doing **Big Science**, and this comes with challenges...

 LLMs are a true **big science** and require significant engineering efforts...
state-of-the-art HPC challenges

 **Principled approaches** are very much needed: tested and validated frameworks

We are doing **Big Science**, and this comes with challenges...

 LLMs are a true **big science** and require significant engineering efforts...
state-of-the-art HPC challenges

 **Principled approaches** are very much needed: tested and validated frameworks
expert HPC/software engineering knowledge

We are doing **Big Science**, and this comes with challenges...

 LLMs are a true **big science** and require significant engineering efforts...
state-of-the-art HPC challenges

 **Principled approaches** are very much needed: tested and validated frameworks
expert HPC/software engineering knowledge
performance tuning is magic currently
e.g. tile/wave quantization, distributed hyperparameters, etc.

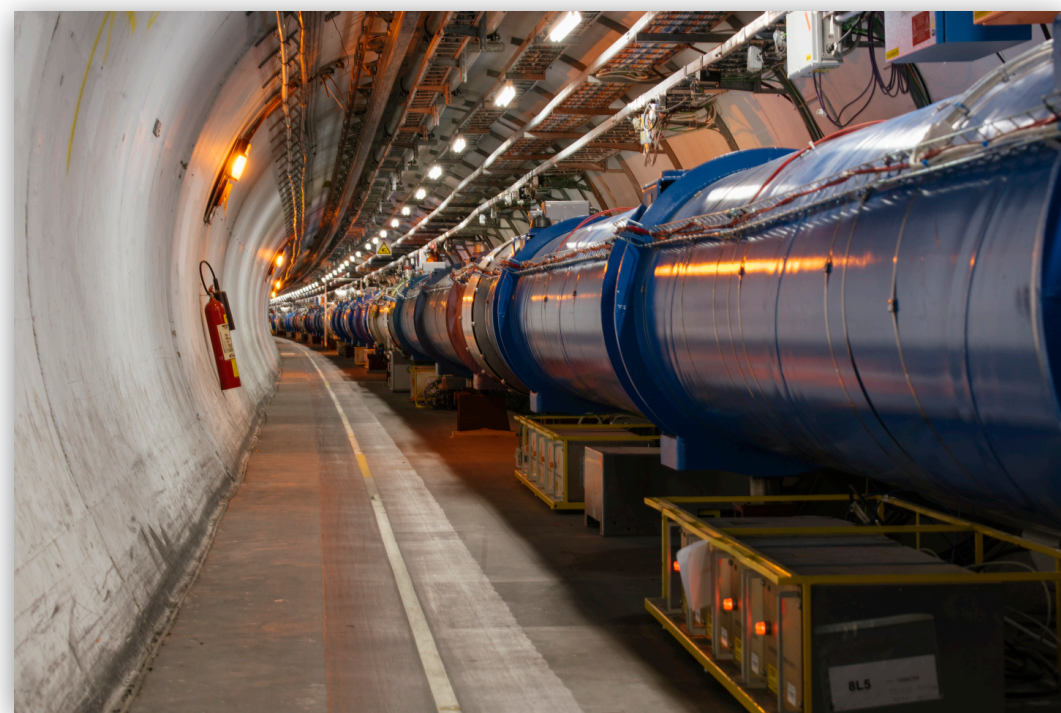
BLOOM: >100 configurations tested!

We are doing **Big Science**, and this comes with challenges...

🚀 LLMs are a true **big science** and require significant engineering efforts...
state-of-the-art HPC challenges

📁 **Principled approaches** are very much needed: tested and validated frameworks
expert HPC/software engineering knowledge
performance tuning is magic currently
e.g. tile/wave quantization, distributed hyperparameters, etc.

BLOOM: >100 configurations tested!



We are doing **Big Science**, and this comes with challenges...

🚀 LLMs are a true **big science** and require significant engineering efforts...
state-of-the-art HPC challenges

📁 **Principled approaches** are very much needed: tested and validated frameworks
expert HPC/software engineering knowledge
performance tuning is magic currently
e.g. tile/wave quantization, distributed hyperparameters, etc.

BLOOM: >100 configurations tested!



(let's avoid this)

Case-study of how hard it can get: Meta’s OPT



OPT: Open Pre-Trained Transformer Language Models

Zhang et al., 2022

Case-study of how hard it can get: Meta's OPT



OPT: Open Pre-Trained Transformer Language Models Zhang et al., 2022

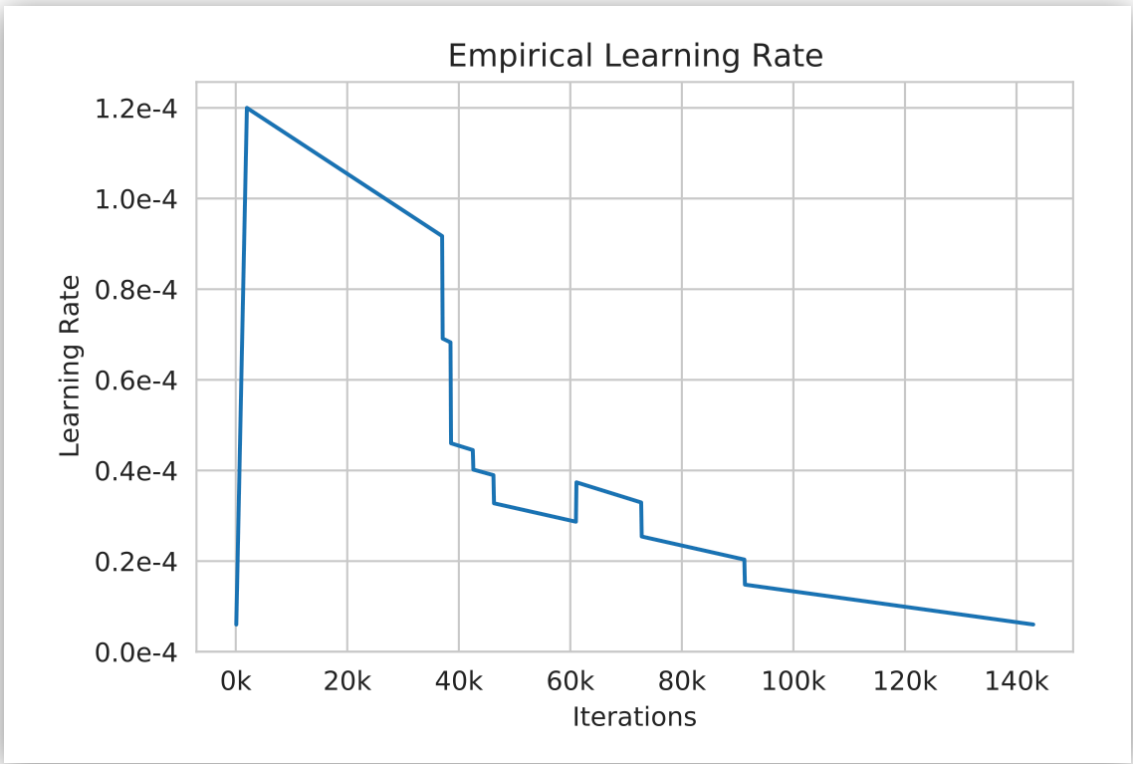
😄 Meta's open "reproduction" of GPT-3 was... a **challenging** experience!

Case-study of how hard it can get: Meta's OPT



OPT: Open Pre-Trained Transformer Language Models Zhang et al., 2022

😄 Meta's open "reproduction" of GPT-3 was... a challenging experience!



manually tuned learning rate

hundreds of restarts, spikes, etc.

Case-study of how hard it can get: Meta's OPT



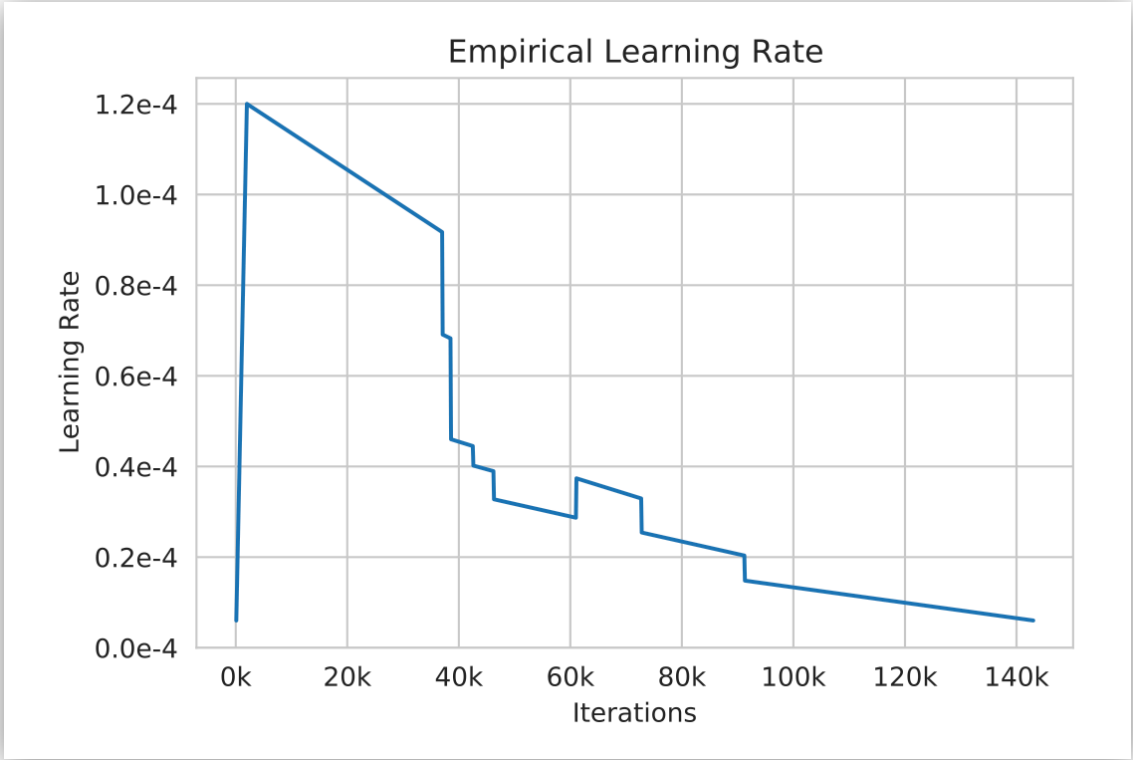
OPT: Open Pre-Trained Transformer Language Models Zhang et al., 2022

😅 Meta's open "reproduction" of GPT-3 was... a challenging experience!

But why?

FP16

BF16



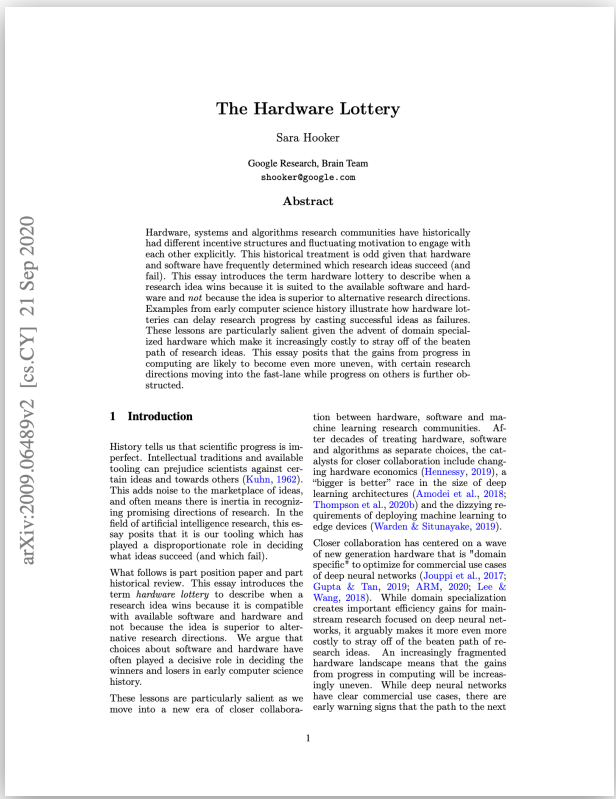
manually tuned learning rate

hundreds of restarts, spikes, etc.



template: Karpathy, 2020

He who controls the chips controls the LLMs

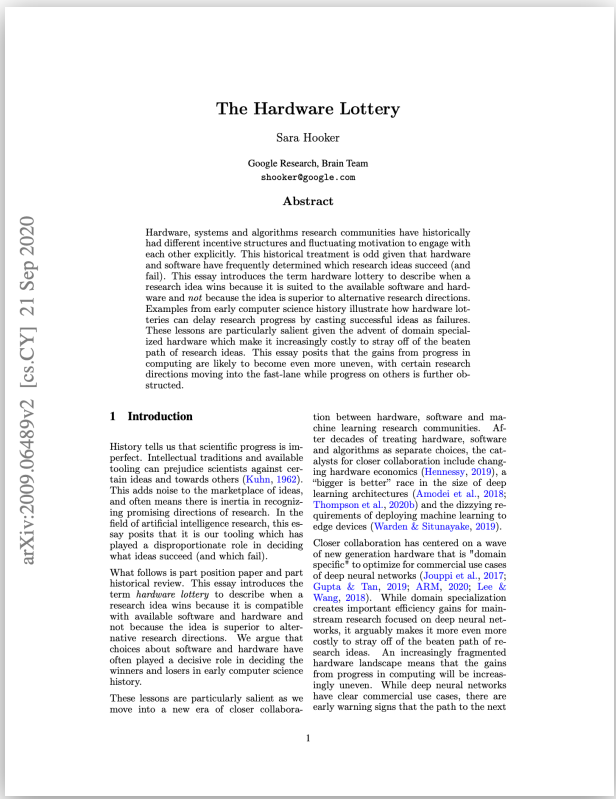


Hardware progress is secretly shaping machine learning

The Hardware Lottery

Sara Hooker, 2020

He who controls the chips controls the LLMs



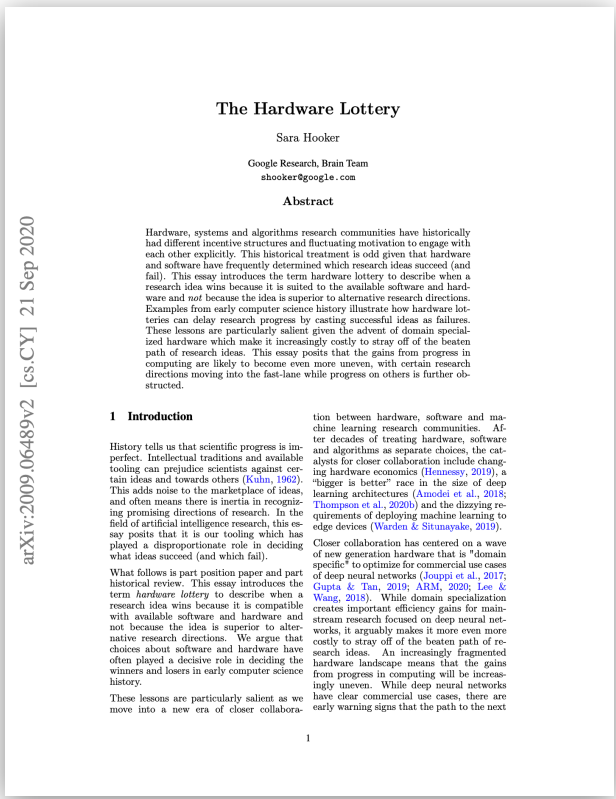
Hardware progress is secretly shaping machine learning

The Hardware Lottery Sara Hooker, 2020

TPU

pure data/model parallelism
uniform platform experience

He who controls the chips controls the LLMs



Hardware progress is secretly shaping machine learning

The Hardware Lottery Sara Hooker, 2020

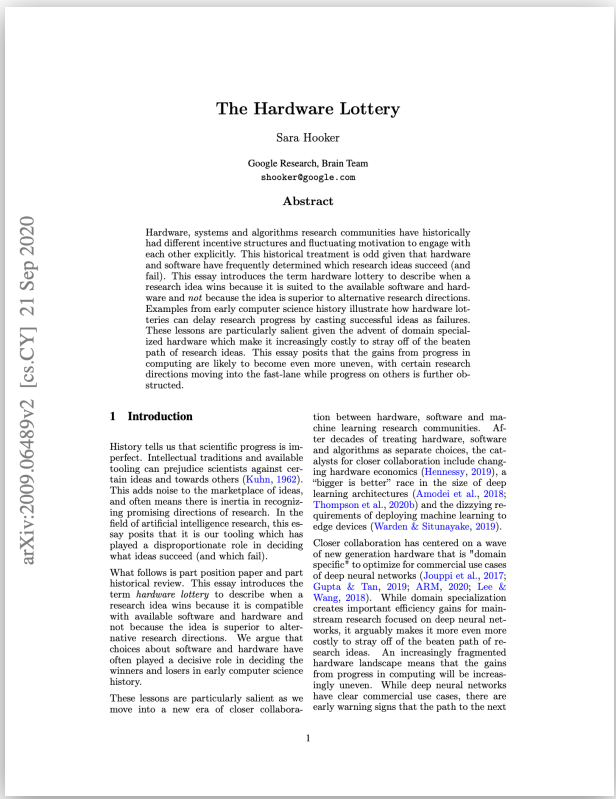
TPU

- pure data/model parallelism
- uniform platform experience

GPU

- data/model/pipeline/sequence parallelism
- diversity in HPC platforms
- network topology, etc.

He who controls the chips controls the LLMs



Hardware progress is secretly shaping machine learning

The Hardware Lottery
Sara Hooker, 2020

TPU

pure data/model parallelism
uniform platform experience

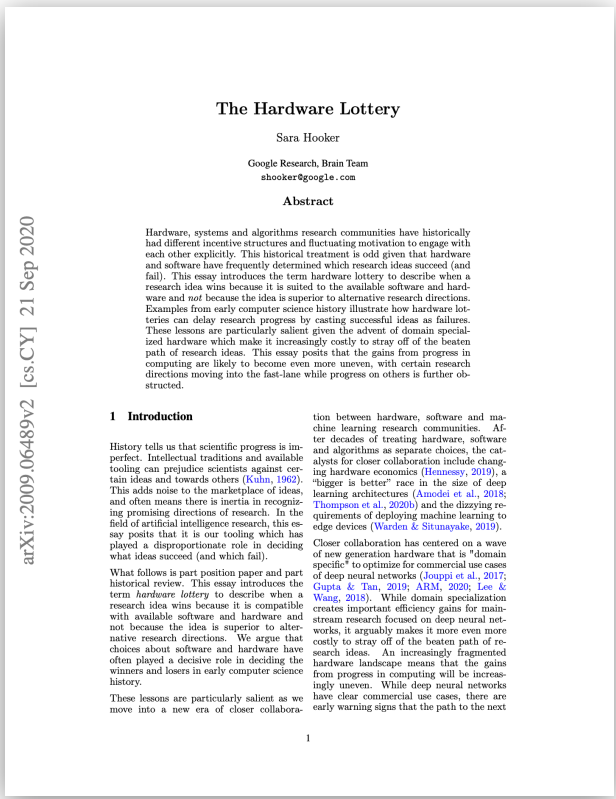


simpler experience?

GPU

data/model/pipeline/sequence parallelism
diversity in HPC platforms
network topology, etc.

He who controls the chips controls the LLMs



Hardware progress is secretly shaping machine learning

The Hardware Lottery
Sara Hooker, 2020

TPU

pure data/model parallelism
uniform platform experience



simpler experience?

GPU

data/model/pipeline/sequence parallelism
diversity in HPC platforms
network topology, etc.

Google, Facebook, Tesla, Amazon are all making their own chips!

Can better **modeling** & more efficient **pretraining** change the playing field?

💨 We can gain in **efficiency**...

Can better **modeling** & more efficient **pretraining** change the playing field?

💡 We can gain in **efficiency**... current approaches, **~50% GPU FLOPs usage**

Can better **modeling** & more efficient **pretraining** change the playing field?

💨 We can gain in **efficiency**... current approaches, **~50% GPU FLOPs usage**
reduced numerical precision: down to **int8**
see Transformer engine in H100

Can better **modeling** & more efficient **pretraining** change the playing field?



We can gain in **efficiency**...

current approaches, **~50% GPU FLOPs usage**

reduced numerical precision: down to **int8**

see Transformer engine in H100

reduce number of **computations**

efficient attention, etc.

Can better **modeling** & more efficient **pretraining** change the playing field?

💨 We can gain in **efficiency**... current approaches, **~50% GPU FLOPs usage**
reduced numerical precision: down to **int8**
see Transformer engine in H100
reduce number of **computations**
efficient attention, etc.

But can we also fundamentally change scaling behaviour?

Can better **modeling** & more efficient **pretraining** change the playing field?

💨 We can gain in **efficiency**... current approaches, **~50% GPU FLOPs usage**
reduced numerical precision: down to **int8**
see Transformer engine in H100
reduce number of **computations**
efficient attention, etc.

But can we also fundamentally change scaling behaviour?

First, **optimise** pretraining:

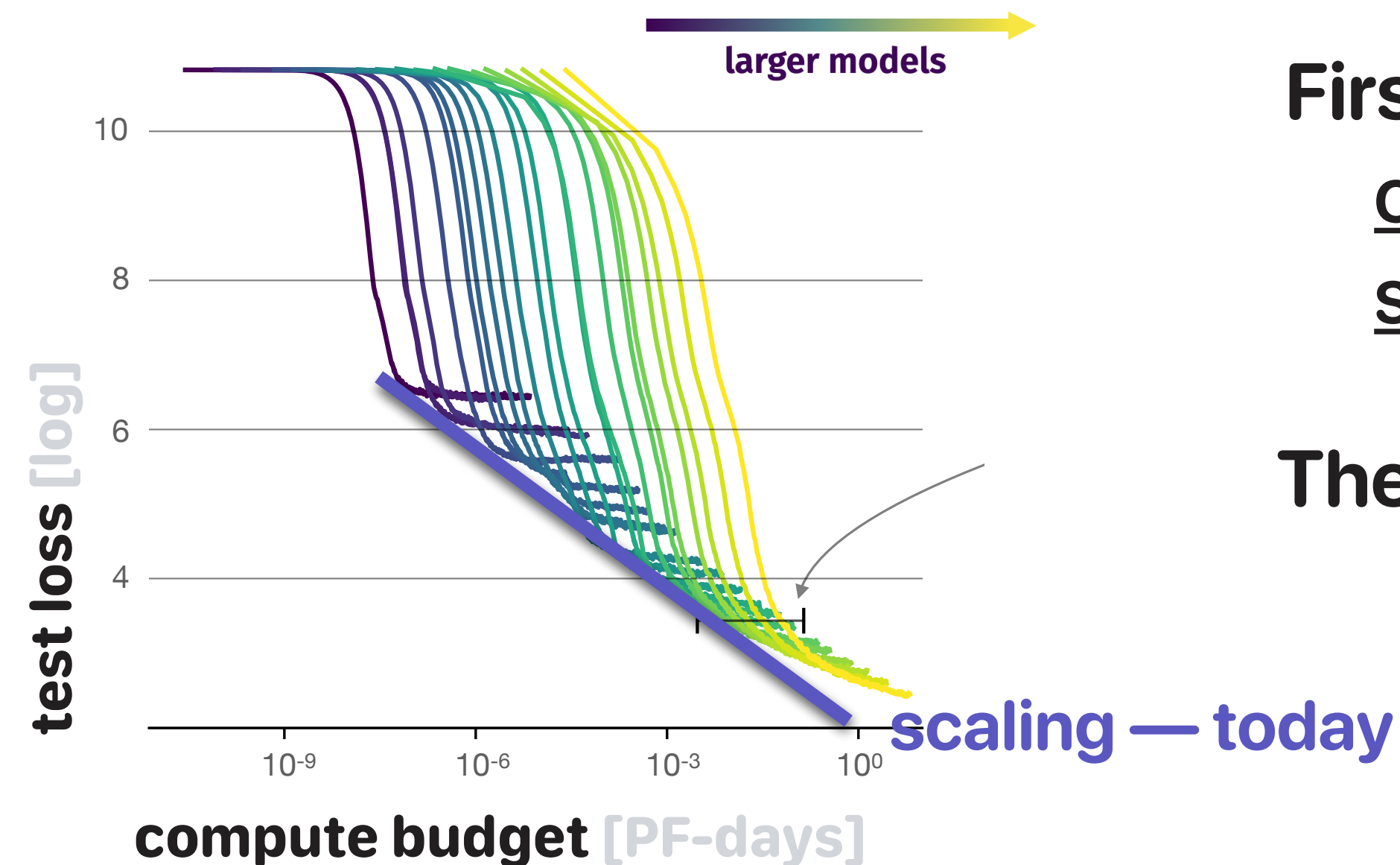
Curriculum learning, grow sequence length Li et al., 2021

Staged training, progressively grow model Shen et al., 2022

Can better **modeling** & more efficient **pretraining** change the playing field?

- 💨 We can gain in **efficiency**...
- current approaches, **~50% GPU FLOPs usage**
 - reduced numerical precision: down to **int8**
see Transformer engine in H100
 - reduce number of **computations**
efficient attention, etc.

But can we also fundamentally change scaling behaviour?



First, **optimise** pretraining:

Curriculum learning, grow sequence length

Li et al., 2021

Staged training, progressively grow model

Shen et al., 2022

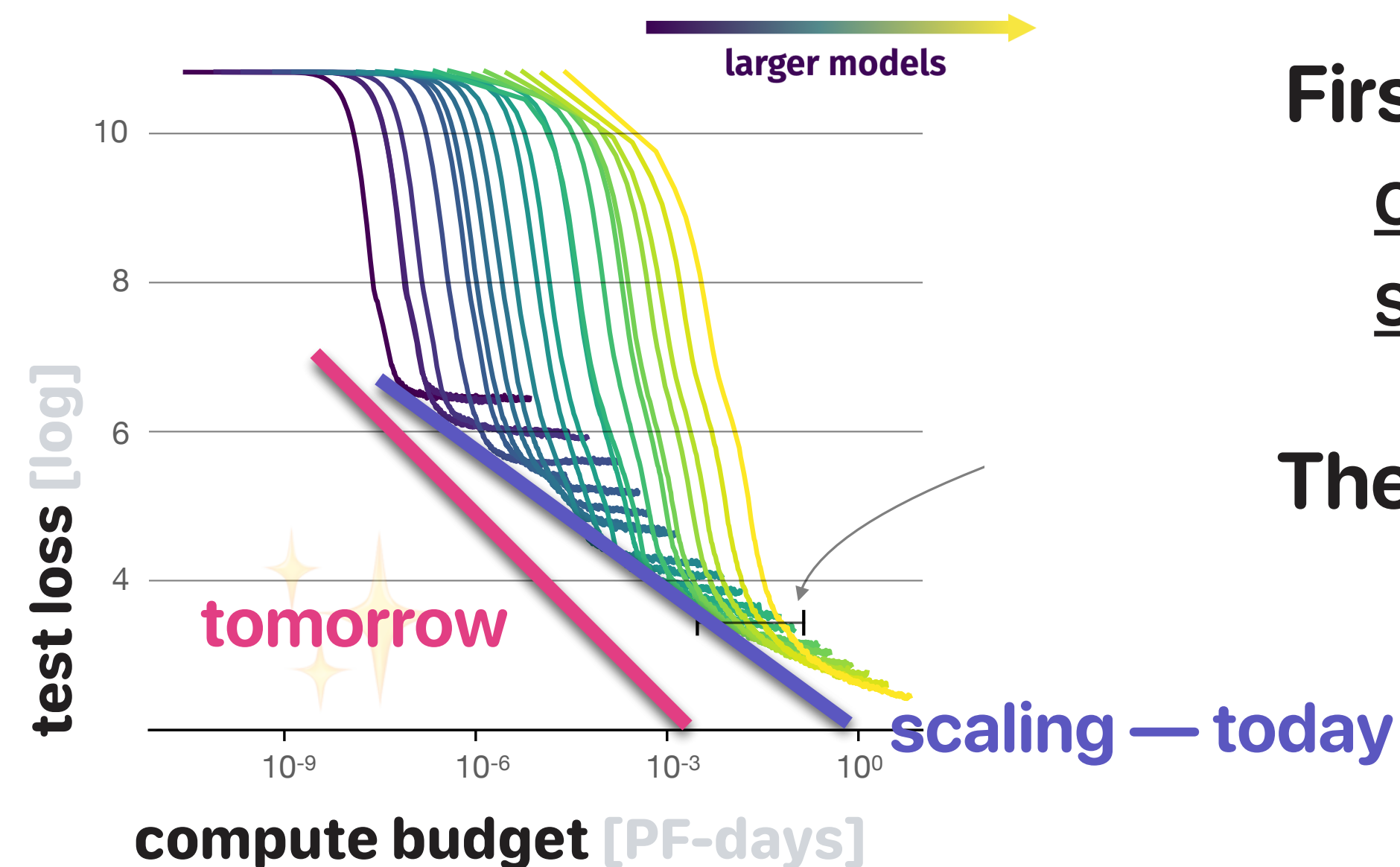
Then, can we get **better** scaling?

🙏 please report scaling laws in your modeling work!

Can better **modeling** & more efficient **pretraining** change the playing field?

- 🧠 We can gain in **efficiency**...
 - current approaches, **~50% GPU FLOPs usage**
 - reduced numerical precision: down to **int8**
 - see Transformer engine in H100
 - reduce number of **computations**
 - efficient attention, etc.

But can we also fundamentally change scaling behaviour?



First, **optimise** pretraining:

Curriculum learning, grow sequence length

Li et al., 2021

Staged training, progressively grow model

Shen et al., 2022

Then, can we get **better** scaling?

🙏 please report scaling laws in your modeling work!

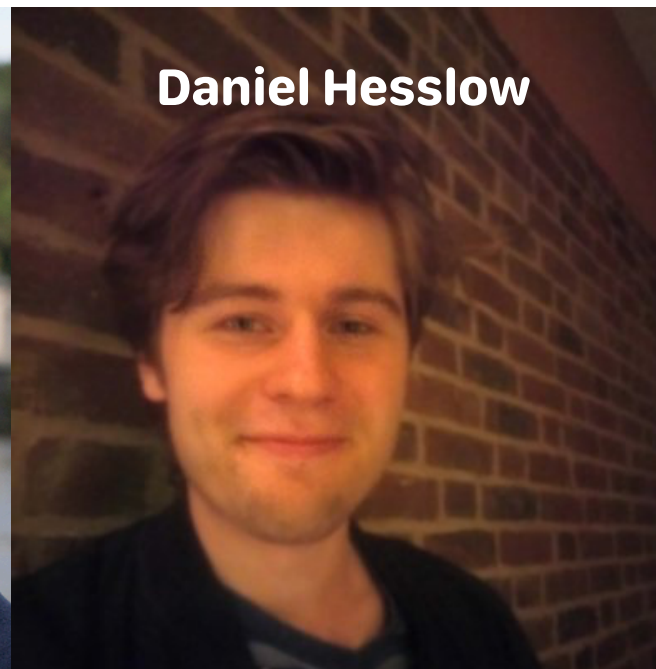
🙌 Thank you to all contributors!



Teven Le Scao



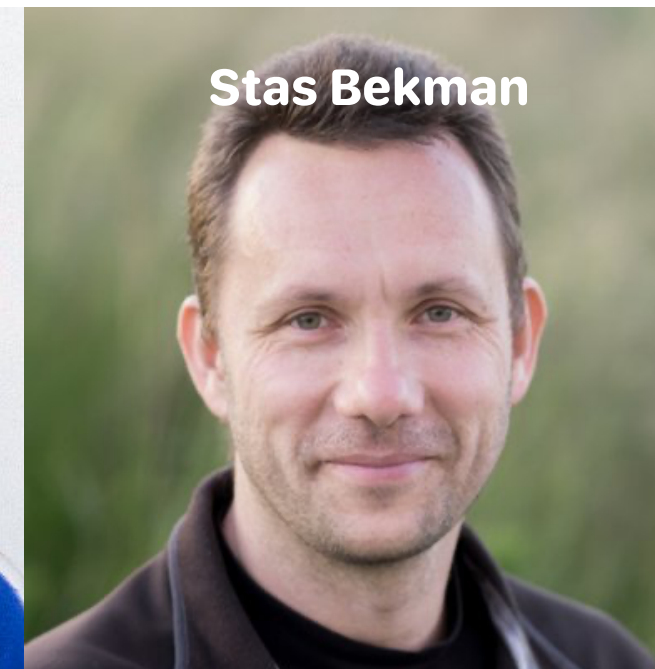
Thomas Wang



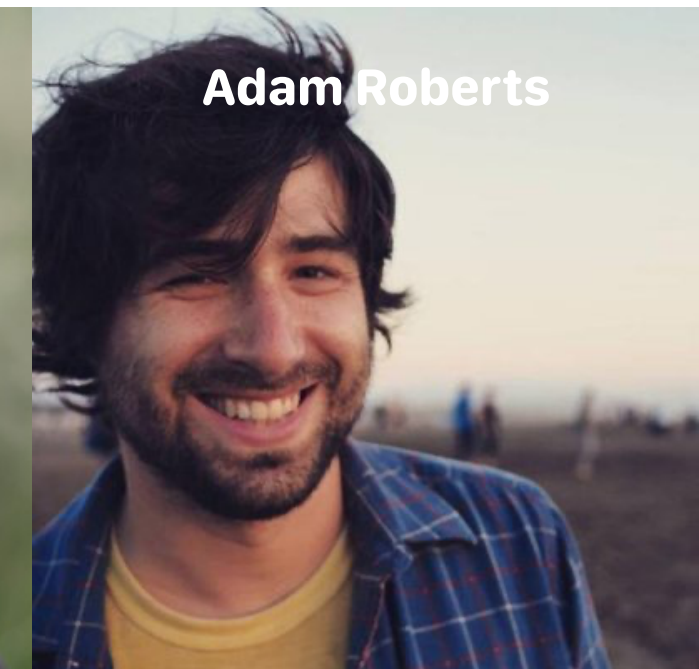
Daniel Hesslow



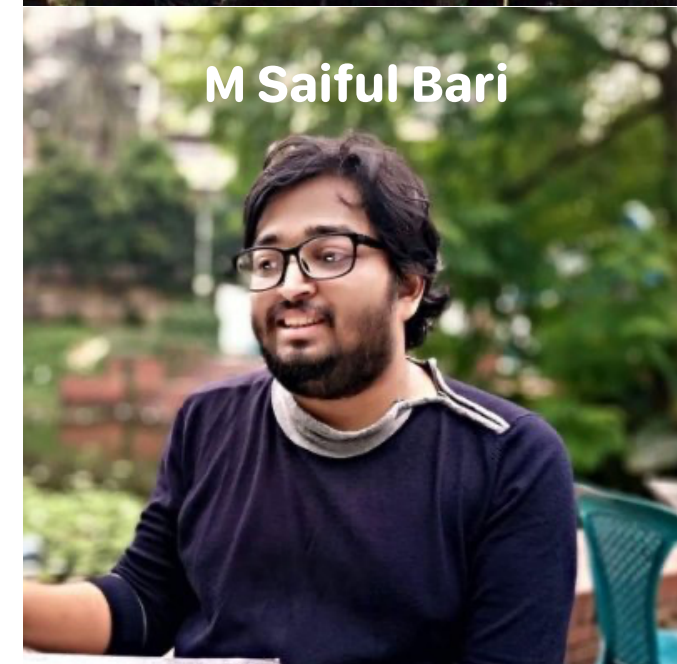
Lucile Saulnier



Stas Bekman



Adam Roberts



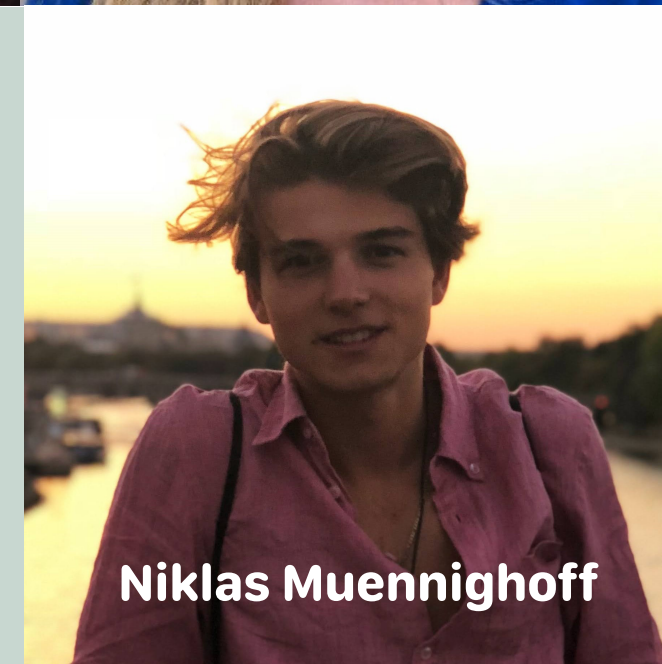
M Saiful Bari



Stella Biderman



Hady Elsahar



Niklas Muennighoff



Jason Phang



Ofir Press



Colin Raffel



Victor Sanh



Sheng Shen



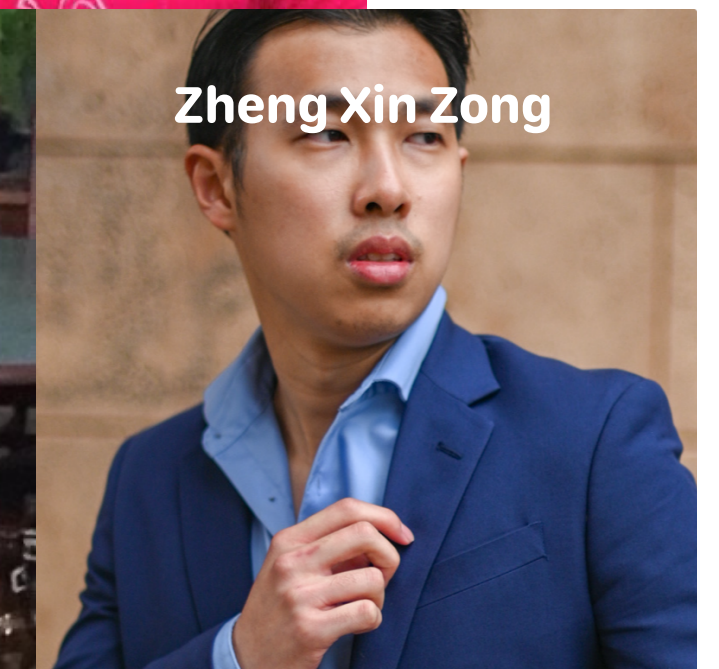
Lintang Sutawika



Jaesung Tae



Hyung Won Chung



Zheng Xin Zong