

# LLMs for every language

A how-to guide

**BigScience**

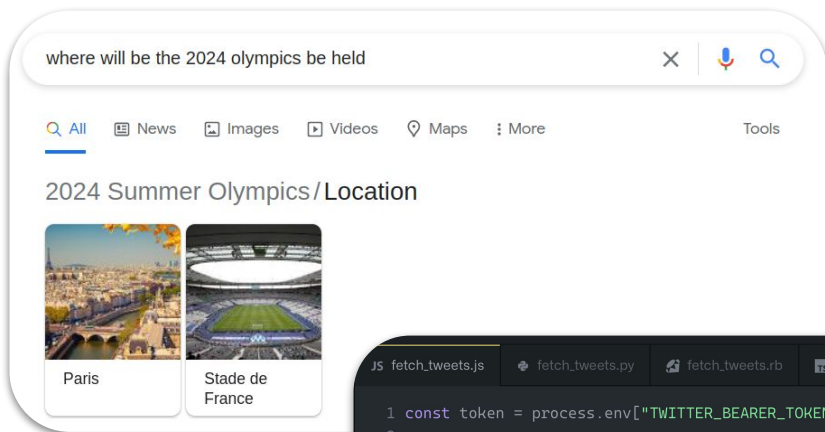


# Collaboratively training a large multilingual language model

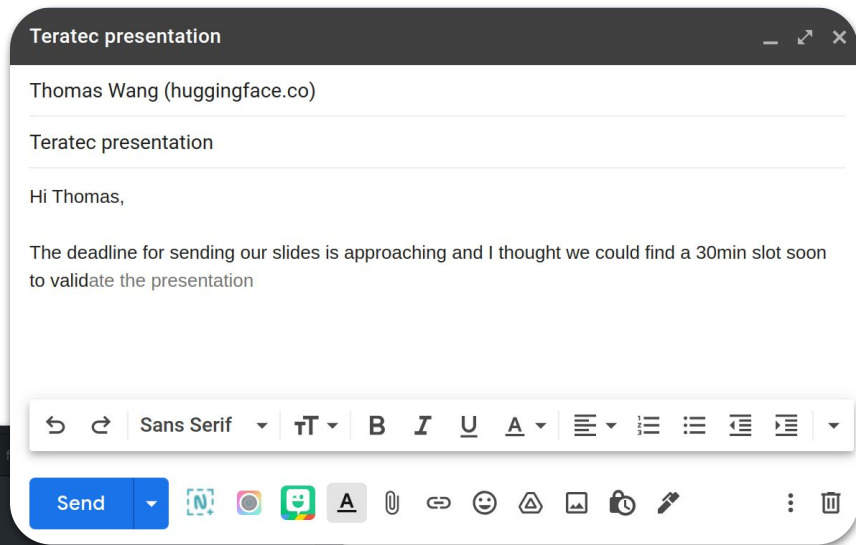
Teven Le Scao

What motivated us to do BigScience?

# Language models



```
Js fetch_tweets.js  fetch_tweets.py  fetch_tweets.rb  fetch_tweets.ts  -90-
1 const token = process.env["TWITTER_BEARER_TOKEN"]
2
3 const fetchTweetsFromUser = async (screenName, count) => {
4   const response = await fetch(
5     `https://api.twitter.com/1.1/statuses/user_timeline.json?screen_name=${screenName}&count=${count}`
6   )
7   const headers = {
8     Authorization: `Bearer ${token}`,
9   }
10  }
11 }
12 const json = await response.json()
13 return json
14 }
```



# Language models

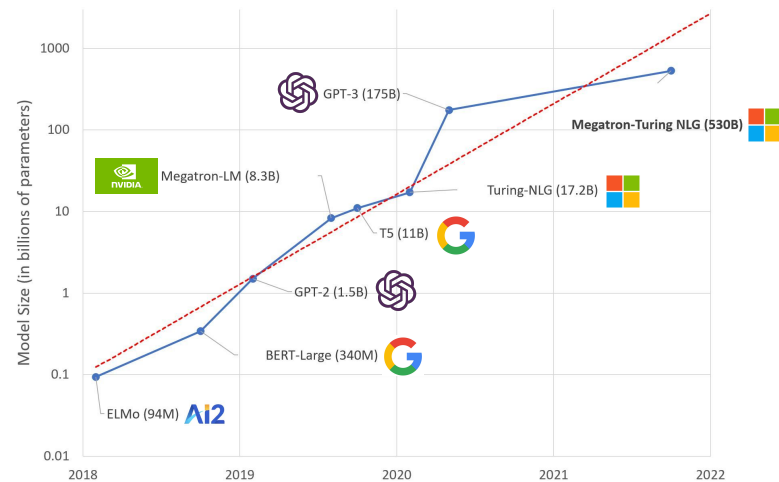
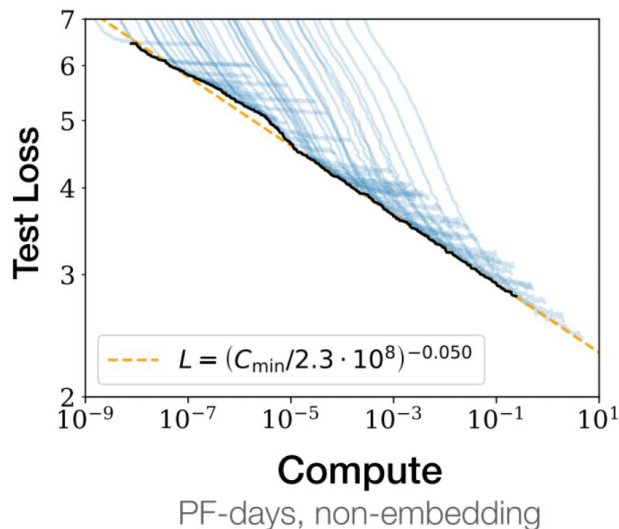
TE

Are language models useful?



Yes, language models are useful for a variety of tasks in natural language processing (NLP), such as machine translation, speech recognition, and text summarization. These models are trained to understand and generate natural language, which allows them to perform tasks that involve understanding and generating human-like text. For example, a language model could be used to automatically generate human-like text responses to customer inquiries, or to help a machine translation system produce accurate and fluent translations. Additionally, language models can be used to improve the performance of other NLP models by providing them with a better understanding of the structure and nuances of natural language.

# Scaling



## GPT-3's generation example:

[...]

To do a "farduddle" means to jump up and down really fast. An example of a sentence that uses the word farduddle is:

One day when I was playing tag with my little sister, she got really excited and she started doing these crazy farduddles.

# Access

## Training cost

- typically \$2-5M
- million of gpu hours

## Closed access for most of them

**VB** VentureBeat

**Naver trained a 'GPT-3-like' Korean language model**

Naver claims the system learned 6,500 times more Korean data than OpenAI's ... Some experts believe that while HyperCLOVA, GPT-3, PanGu-a, ...  
1 Jun 2021



**TC** TechCrunch

**Anthropic is the new AI research outfit from OpenAI's Dario Amodei, and it has \$124M to burn**

Anthropic, as it's called, was founded with his sister Daniela and its goal is to create "large-scale AI systems that are steerable, ...  
28 May 2021



**VB** VentureBeat

**AI21 Labs trains a massive language model to rival OpenAI's GPT-3**

"AI21 Labs was founded to fundamentally change and improve the way people read and write. Pushing the frontier of language-based AI requires ...  
1 month ago



**FC** Fast Company

**Ex-Googlers raise \$40 million to democratize language AI**

This story has been updated with more information about Cohere's approach to responsible AI. About the author. Fast Company Senior Writer Mark ...  
2 days ago



# BigScience

“During **one-year**, from May 2021 to May 2022, 1000+ researchers from 60 countries and more than 250 institutions are **creating together a very large multilingual neural network language model** and a **very large multilingual text dataset** on the 28 petaflops Jean Zay (IDRIS) supercomputer located near Paris, France.

During the workshop, the participants plan to investigate the dataset and the model from all angles: bias, social impact, capabilities, limitations, ethics, potential improvements, specific domain performances, carbon impact, general AI/cognitive research landscape.”



# Spirit of the project

## **Make LLM research accessible**

- Open-source, open-access
- Organize research around the models
- Make compute accessible

## **Create and share knowledge around the process**

- Train in the open
- Freely discuss engineering problems and solutions


# Artifacts: what came out of BigScience? (0/4)

```
from transformers import AutoModel, AutoTokenizer

model_name = "bigscience/bloom"
tokenizer = AutoTokenizer.from_pretrained(model_name)
model = AutoModel.from_pretrained(model_name)
```

<https://huggingface.co/bigscience/bloom>

# Artifacts: what came out of BigScience? (1/4)




## BigScience Workshop

Research workshop on large language models - The Summer of Language Models 21

<https://bigscience.huggingface.co> [@BigScienceW](#) [bigscience-contact@googlegroups.com](mailto:bigscience-contact@googlegroups.com)

[Overview](#) [Repositories 29](#) [Projects 4](#) [Packages](#) [People 22](#)



### Popular repositories

#### promptsource

Toolkit for creating, sharing and using natural language prompts.

Python 660 144

Public

#### bigscience

Central place for the engineering/scaling WG: documentation, SLURM scripts and logs, compute environment and data.

Shell 277 26

Public

#### Megatron-DeepSpeed

Ongoing research training transformer language models at scale, including: BERT & GPT-2

Python 145 34

Public

#### t-zero

Reproduce results and replicate training for T0 (Multitask Prompted Training Enables Zero-Shot Task Generalization)

Python 143 26

Public

#### biomedical

Tools for curating biomedical training data for large-scale language modeling

Python 125 69

Public


#### evaluation

Code and Data for Evaluation WG

Python 36 23

Public

### People



[View all](#)

### Top languages

Python Jupyter Notebook HTML  
Shell Makefile

### Most used topics

machine-learning nlp

# Artifacts: what came out of BigScience? (2/4)

## Masader: Metadata Sourcing for Arabic Text and Speech Data Resources

Zaid Alyafei<sup>1</sup>, Maraim Masoud<sup>2</sup>, Mustafa Ghaleb<sup>1</sup>, and Maged S. Al-shaibani<sup>1</sup>

<sup>1</sup> King Fahd University of Petroleum and Minerals, Dhahran, Saudi Arabia  
<sup>2</sup> Independent Researcher

### Abstract

The NLP pipeline has evolved dramatically in the last few years. The first step in the pipeline is to find suitable annotated datasets to evaluate the tasks we are trying to solve. Unfortunately, most of the published datasets lack metadata annotations that describe their attributes. Not to mention, the absence of a public catalogue that indexes all the publicly available datasets related to specific regions or languages. When we consider low-resource dialectal languages, for example, this issue becomes more prominent. In this paper we create *Masader*, the largest public catalogue for Arabic NLP datasets, which consists of 200 datasets annotated with 25 attributes. Further-

and so on. This study attempts to identify the publicly available Arabic NLP datasets and to provide a catalogue of Arabic datasets to researchers. The catalogue will increase the discoverability and provide some key metadata that will help researchers identify the most suitable dataset for their research questions.

We

## Between words and characters: A Brief History of Open-Vocabulary Modeling and Tokenization in NLP

Sabrina J. Mielke<sup>1,2</sup> Zaid Alyafei<sup>3</sup> Elizabeth Salesky<sup>1</sup>  
 Colin Raffel<sup>2</sup> Manan Dey<sup>4</sup> Matthias Gallé<sup>5</sup> Arun Raja<sup>6</sup>  
 Chenglei Si<sup>7</sup> Wilson Y. Lee<sup>8</sup> Benoît Sagot<sup>9\*</sup> Samson Tan<sup>10\*\*</sup>

*BigScience Workshop Tokenization Working Group*

<sup>1</sup>Johns Hopkins University <sup>2</sup>HuggingFace <sup>3</sup>King Fahd University of Petroleum and Minerals <sup>4</sup>SAP  
<sup>5</sup>Naver Labs Europe <sup>6</sup>Institute for Infocomm Research, A\*STAR Singapore <sup>7</sup>University of Maryland  
<sup>8</sup>BigScience Workshop <sup>9</sup>Inria Paris <sup>10</sup>Salesforce Research Asia & National University of Singapore  
 \*s.jm@s.jm.i.e.k.e. com \*\*s.tan@salesforce.com

### Abstract

What are the units of text that we want to model? From bytes to multi-word expressions, text can be analyzed and generated at many granularities. Until recently, most natural language processing (NLP) models operated over words, treating those as discrete and atomic tokens, but starting with byte-pair encoding (BPE), subword-based approaches have become dominant in many areas, enabling small vocabularies while still allowing for fast inference. Is the end of the road character-level model or byte-level processing? In this survey, we connect several lines of work from the pre-neural and neural era, by showing how hybrid approaches

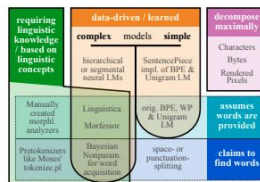


Figure 1: A taxonomy of segmentation and tokenization algorithms and research directions

## MULTITASK PROMPTED TRAINING ENABLES ZERO-SHOT TASK GENERALIZATION

Victor Sanh<sup>\*</sup> Hugging Face Albert Webson<sup>\*</sup> Brown University Colin Raffel<sup>\*</sup> Hugging Face Stephen H. Bach<sup>\*</sup> Brown & Snorkel AI

Lintang Sutowika<sup>\*</sup> BigScience Zaid Alyafei<sup>\*</sup> KFUPM Antoine Chaffin<sup>\*</sup> IRISA & IMATAG Arnaud Stiegler<sup>\*</sup> Hyperscience Teven Le Scao<sup>\*</sup> Hugging Face Arun Raja<sup>\*</sup> Manan Dey<sup>\*</sup> M Saiful Bari<sup>\*</sup> M Taewoon Kim<sup>\*</sup> Gunjan Chhabhani<sup>\*</sup> Nihal V. Nayak<sup>\*</sup> Urmish Thakker<sup>\*</sup> Harshit Pandey<sup>\*</sup> BigScience Michael McKenna<sup>\*</sup> Parity Rachel Bawden<sup>\*</sup> Inria, France Thomas Wang<sup>\*</sup> Inria, France Trishala Neeraj<sup>\*</sup> BigScience Jos Rozen<sup>\*</sup> Naver Labs Europe Abheesh Sharma<sup>\*</sup> BITS Pilani, India Andrea Santilli<sup>\*</sup> University of Rome

lan Fries & Snorkel AI Ryan Teehan Charles River Analytics Tali Bers Brown University

Leo Gao EleutherAI Thomas Wolf Hugging Face Alexander M. Rush Hugging Face

### ABSTRACT

It has recently been shown to attain reasonable zero-shot performance set of tasks (Brown et al., 2020). It has been hypothesized that implicit multitask learning in language models (Sanh et al., 2019). Can zero-shot generalization instead be directly attributed to multitask learning? To test this question at scale, we develop a framework for training language models on a large set of natural language tasks into a human-readable

And many mores..

## What Language Model to Train if You Have One Million GPU Hours?

### The BigScience Architecture & Scaling Group

Teven Le Scao<sup>1\*</sup> Thomas Wang<sup>1\*</sup> Daniel Hesse<sup>2\*</sup> Lucile Saulnier<sup>1\*</sup> Stas Bekman<sup>1\*</sup>  
 M Saiful Bari<sup>3</sup> Stella Biderman<sup>4,5</sup> Hady Elsahar<sup>6</sup> Jason Phang<sup>7</sup> Ofir Press<sup>7</sup> Colin Raffel<sup>1</sup>  
 Victor Sanh<sup>1</sup> Sheng Shen<sup>9</sup> Lintang Sutowika<sup>10</sup> Jaesung Tae<sup>1</sup> Zheng Xin Yong<sup>11</sup>

Julien Launay<sup>2,12†</sup> Iz Beltagy<sup>13†</sup>

<sup>1</sup>Hugging Face <sup>2</sup>LightOn <sup>3</sup>NTU, Singapore <sup>4</sup>Booz Allen <sup>5</sup>EleutherAI <sup>6</sup>Naver Labs Europe <sup>7</sup>New York University  
<sup>8</sup>University of Washington <sup>9</sup>Berkeley University <sup>10</sup>Big Science <sup>11</sup>Brown University <sup>12</sup>LPENS <sup>13</sup>Allen Institute for AI

Modeling methods have been a well-motivated transfer across impact of modeling the emergence of parameters models, reusing expert-train. Notably, how modeling capabilities, use mainly from singular language scale, our goal training setup

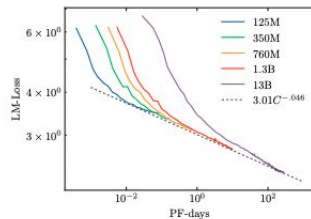


Figure 1: Smooth scaling of language modeling loss as compute budget and model size increase. We observe a power-law coefficient  $\alpha_C \sim 0.046$ , in-line with pre-dict

## What Language Model Architecture and Pretraining Objective Work Best for Zero-Shot Generalization?

### The BigScience Architecture & Scaling Group

Thomas Wang<sup>1\*</sup> Adam Roberts<sup>2\*</sup>  
 Daniel Hesse<sup>3</sup> Teven Le Scao<sup>1</sup> Hyung Won Chung<sup>2</sup>  
 Iz Beltagy<sup>4</sup> Julien Launay<sup>3,5†</sup> Colin Raffel<sup>1†</sup>

<sup>1</sup>Hugging Face <sup>2</sup>Google <sup>3</sup>LightOn

<sup>4</sup>Allen Institute for AI <sup>5</sup>LPENS, École Normale Supérieure

### Abstract

Large pretrained Transformer language models have been shown to exhibit zero-shot generalization, i.e. they can perform a wide variety of tasks that they were not explicitly trained on. However, the architectures and pretraining objectives used across state-of-the-art models differ significantly, and there has been limited systematic comparison of these factors. In this work, we present a large-scale evaluation of modeling choices and their impact on zero-shot generalization. In particular, we focus on text-to-text models and experiment with three model architectures (causal/non-causal decoder-only and encoder-decoder), trained with two different pretraining objectives (autoregressive and masked language modeling), and evaluated with and without multitask prompted finetuning. We train

## Artifacts: what came out of BigScience? (3/4)



### **BigScience**

#### **BigScience RAIL License v1.0**

This is the home of the BigScience RAIL License v1.0. If you would like to download the license you can get it as [.txt](#), [.docx](#), or [.html](#) file.

<https://huggingface.co/spaces/bigscience/license>

# Artifacts: what came out of BigScience? (4/4)



Stas Bekman  
@StasBekman

The embed matrix with 250k multi-lingual vocab is on par in size with the transformer block, so rebalancing the pipeline to count embedding matrices as transformer blocks leads to even faster throughput and less memory usage on ranks 0 and -1

Benchmarks:

bigscience-workshop/  
**bigscience**



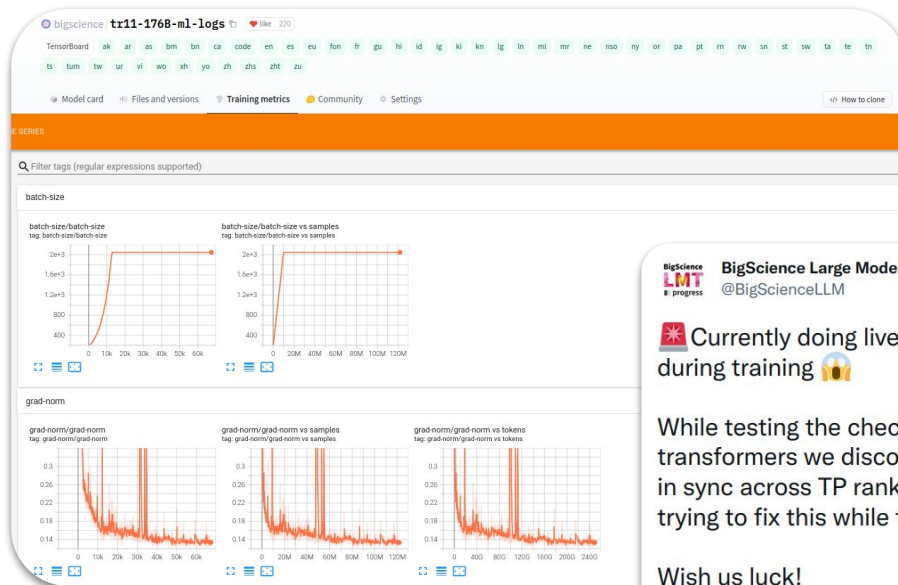
Central place for the engineering/scaling WG:  
documentation, SLURM scripts and logs, compute  
environment and data.

13 Contributors   2 Issues   272 Stars   24 Forks

github.com

bigscience/chronicles-prequel.md at master · bigscience-workshop/bigscience  
Central place for the engineering/scaling WG: documentation, SLURM scripts  
and logs, compute environment and data. - bigscience/chronicles-prequel.md a...

6:34 AM · Mar 4, 2022 · Twitter Web App



BigScience Large Model Training  
@BigScienceLLM

Currently doing live surgery on the 176B model during training 🏠

While testing the checkpoint weights for integration in transformers we discovered that layer norms were not in sync across TP ranks contrary to what expected - trying to fix this while training 🙌

Wish us luck!

12:38 PM · Mar 25, 2022 · Twitter Web App

<https://huggingface.co/bigscience/tr11-176B-ml-logs>

<https://github.com/bigscience-workshop/bigscience/blob/master/train/lessons-learned.md>

# Training a 176B model

# Jean Zay

✿ This work was granted access to the HPC resources of *Institut du développement et des ressources en informatique scientifique* (IDRIS) du *Centre national de la recherche scientifique* (CNRS) under the allocation 2021-A0101012475 made by *Grand équipement national de calcul intensif* (GENCI) - Thank you!

✿ Compute grant:

- 2.5M V100 hours
- 1.25M A100 hours: a reserved allocation of 416 A100 (80GB)
- and a ton of CPU

✿ Technical support - Thanks Rémi Lacroix!

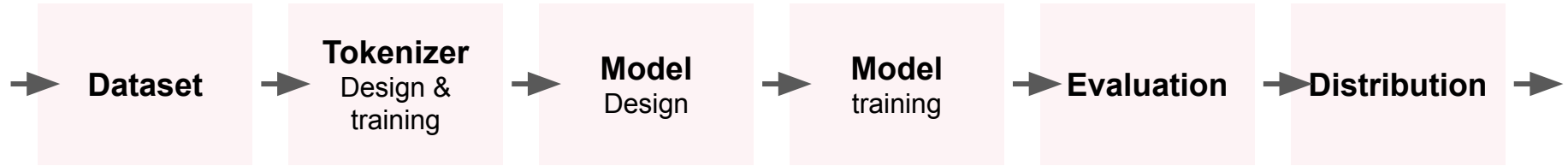


INSTITUT DU  
DÉVELOPPEMENT ET DES  
RESSOURCES EN  
INFORMATIQUE  
SCIENTIFIQUE

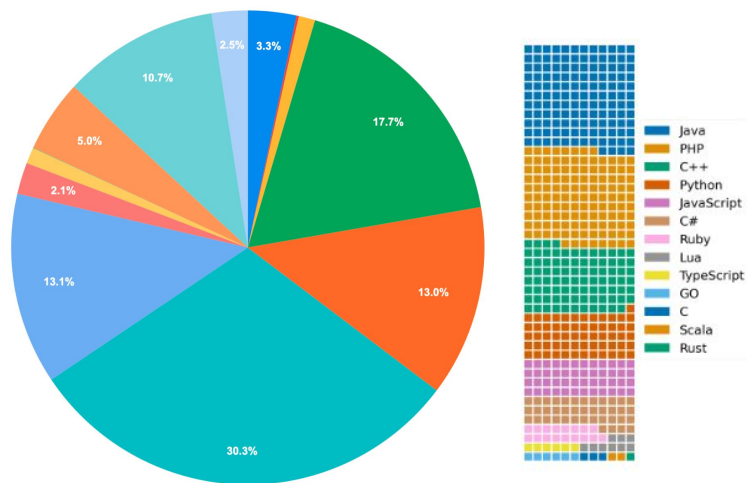




# Pipeline



# Dataset: 1.6T multilingual



- Arabic (3,3%)
- Basque (0,2%)
- Catalan (1,1%)
- Chinese (17,7%)
- Code (13%)
- English (30,3%)
- French (13,1%)
- Indic (2,1%)
- Indonesian (1,1%)
- Niger Congo (0,03%)
- Portuguese (5%)
- Spanish (10,7%)
- Vietnamese (2,5%)

<https://openreview.net/forum?id=UoEw6KigkUn>

# Sources: catalogue

Crowdsourced multilingual datasets from BigScience participants

~60% of data in tokens

Lessons:

- Some filtering and deduplication required, lots of templates to remove
- A lot of those are *not* clean
- LMs require specific corpora: unsupervised, diverse text with long documents

# Web crawl

Our filtered and deduped version of OSCAR-v1

~35% of data in tokens

Lessons:

- Heavy filtering needed, but you can find good data
- Hard to assess deduplication at TB-scale

# Web crawl filtering

## 7 simple filters

AR	EU	BN	CA	ZH	EN	FR	HI	ID	PT	UR	VI	ES
20.3	5.2	48.8	21.1	23.1	17.2	17.0	25.7	10.4	12.6	15.8	21.3	16.9

Table 1: Percentage of documents removed by the filtering per language (ISO 639-1 code).

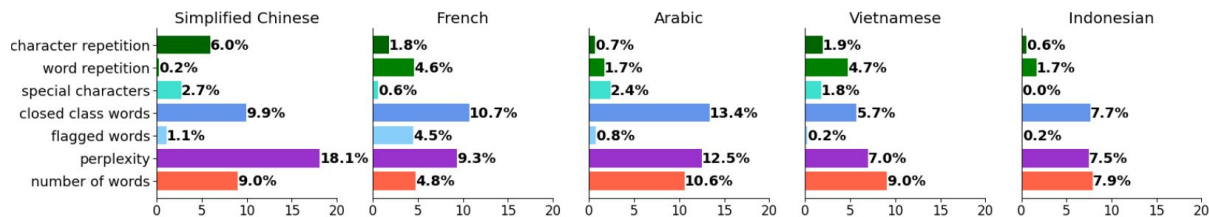


Figure 3: Percentage of documents discarded by each filter independently for 5 languages

# Sources: pseudo-crawl

List of domains from native speakers, corresponding Common Crawl WARC files

5% of data

## Lessons:

- Doing your own WARC parsing is rewarding but a project in itself
- We probably should have crawled the sites directly
- Tons of processing/deduplication/template removal needed

# Deduplication

For catalogue and pseudo-crawl

- Line by line within every document
- Across all documents

And for the web crawl:

- MinHash near-dedup (0.7% of data removed, wary of false +)
- Suffix tree exact dedup over long documents (21.67% duplication, stringent)

# Tokenizer choices

Metric: **fertility**

- Minimize vocab size
- Minimize tokens per byte of text



# Tokenizer choices

Objective: no more than 10%+ fertility compared to monolingual tokenizers

Reached at **250k tokens**

Tokenizer	fr	en	es	zh	hi	ar
Monolingual	1.30	1.15	1.12	1.50	1.07	1.16
BLOOM	1.17 (-11%)	1.15 (+0%)	1.16 (+3%)	1.58 (+5%)	1.18 (+9%)	1.34 (+13%)

Table 2: Fertilities obtained on Universal Dependencies treebanks on languages with existing monolingual tokenizers. The monolingual tokenizers we used were the ones from CamemBERT (Martin et al., 2020), GPT-2 (Radford et al., 2019), `DeepESP/gpt2-spanish`, `bert-base-chinese`, `monsoon-nlp/hindi-bert` and Arabic BERT (Safaya et al., 2020), all available on the HuggingFace Hub.

Lessons:

- Weird tokens (e.g. URLs) are good indicators of training data pathologies
- Massively multilingual vocabs are hard... Still messed up for Devanagari

# GPT-style autoregressive



The latest from Google Research



Language modelling at scale:  
Gopher, ethical considerations,  
and retrieval

December 8, 2021

RESEARCH

## Democratizing access to large-scale language models with OPT-175B

May 3, 2022

## Pathways Language Model (PaLM): Scaling to 540 Billion Parameters for Breakthrough Performance

Monday, April 4, 2022

---

## YUAN 1.0: LARGE-SCALE PRE-TRAINED LANGUAGE MODEL IN ZERO-SHOT AND FEW-SHOT LEARNING

---

Shaohua Wu*	Xudong Zhao	Tong Yu	
Rongguo Zhang	Chong Shen	Hongli Liu	Feng Li
Hong Zhu	Jiangang Luo	Liang Xu	Xuanwei Zhang



An empirical analysis of compute-optimal large language model training

Download

View publication

## Announcing AI21 Studio and Jurassic-1 Language Models


AI21 Labs' new developer platform offers instant access to our 178B-parameter language model, to help you build sophisticated text-based AI applications at scale

## Announcing GPT-NeoX-20B

Announcing GPT-NeoX-20B, a 20 billion parameter model trained in collaboration with CoreWeave.  
February 2, 2022 - Connor Leahy

# Code stack

```
~checkpoints/tr11-176B-ml/checkpoints/main> du -h global_step63600/  
2.3T  global_step63600/
```

 Including optimizer states and checkpoints

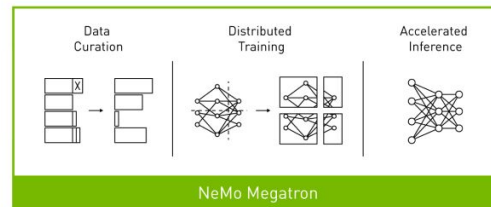


 [microsoft / Megatron-DeepSpeed](#) Public

forked from [NVIDIA/Megatron-LM](#)

# Code stacks I would use now

- If you have high interconnect
- If not
- If you have TPUs



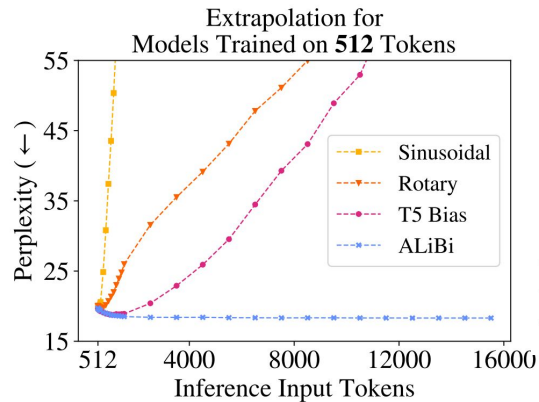
google-research/  
**t5x**



70 Contributors   6 Used by   8 Discussions   1k Stars   160 Forks

# Modeling adjustments

- **ALiBi** positional embeddings which allow long-sequence extrapolation



Positional Embedding	Average EAI Results
None	41.23
Learned	41.71
Rotary	41.46
ALiBi	<b>43.70</b>

Table 2: **ALiBi significantly outperforms other embeddings for zero-shot generalization.** All models are trained on the OSCAR dataset for 112 billion tokens.

- **Embed LayerNorm** for stability (there's a perf tradeoff though!)

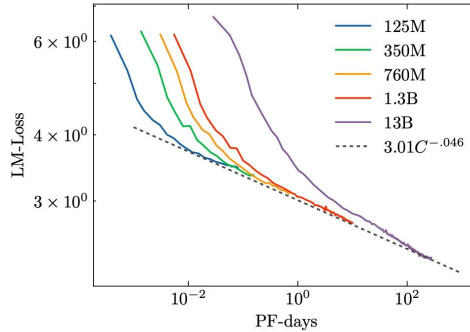
# Scaling

Compute (total operations) =  $k * \text{Data} * \text{Parameters}$

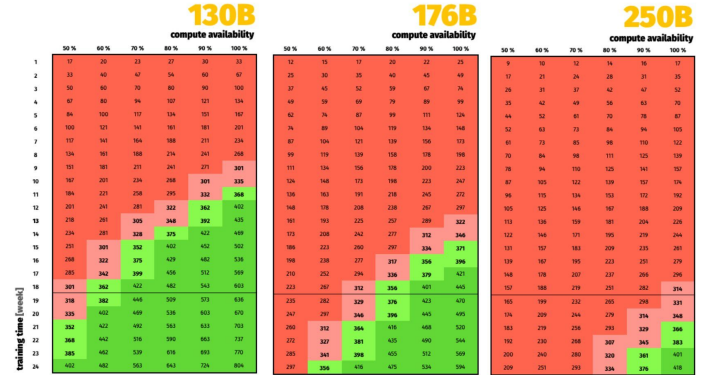
As C (compute) goes up, how much should go into D (data size) vs N (params)?

# Scaling

## Scaling laws on English



Training matrix [tokens]



Fixed budget scenarios for different model sizes

Model	Size [Bparams.]	Pretraining [Btokens]	Budget [PF-days]	Layers	Hidden dim.	Attention heads num.	Attention heads dim.
LaMDA (Thoppilan et al., 2022)	137	432	4,106	64	8,192	128	64
GPT-3 (Brown et al., 2020)	175	300	3,646	96	12,288	96	128
J1-Jumbo (Lieber et al., 2021)	178	300	3,708	76	13,824	96	144
PanGu- $\alpha$ (Zeng et al., 2021)	207	42	604	64	16,384	128	128
Yuan (Wu et al., 2021)	245	180	3,063	76	16,384		
Gopher (Rae et al., 2021)	280	300	4,313	80	16,384	128	128
MT-530B (Smith et al., 2022)	530	270	9,938	105	20,480	128	160

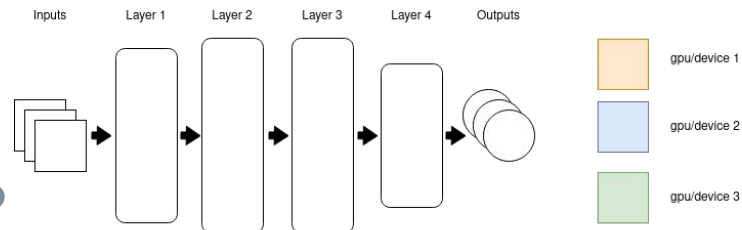
Model	Size [params.]	Layers	Hidden dim.	Attention heads num.	Attention heads dim.	Memory [GB]	Performance [sec/iter.]	Performance [TFLOPs]
(1)	178	82	13,312	64	208	63	104	152
(2)	178	82	13,312	128	104	60	109	146
(3)	176	70	14,336	112	128	59	105	150

[arxiv.org/abs/2001.08361](https://arxiv.org/abs/2001.08361)

[arxiv.org/abs/2006.12467](https://arxiv.org/abs/2006.12467)

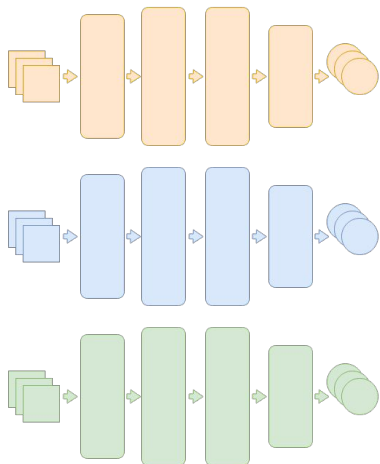
# Parallelism:

## how to use a cluster wisely for DL?



### Data parallelism

to accelerate the training speed

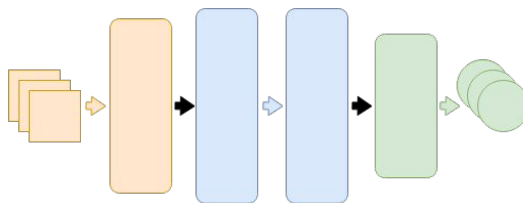


Each device has a replica of the model and receives a different batch of training data on which it performs a forward and backward pass

### Model parallelism

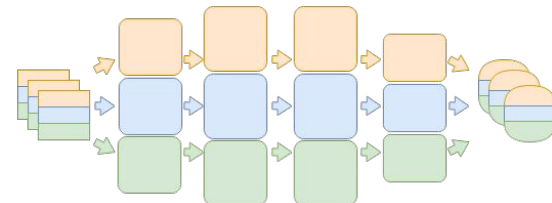
to train models that don't fit in the memory of one device

#### Pipeline parallelism



Only one or several consecutive layers of the model are placed on a single GPU

#### Tensor parallelism



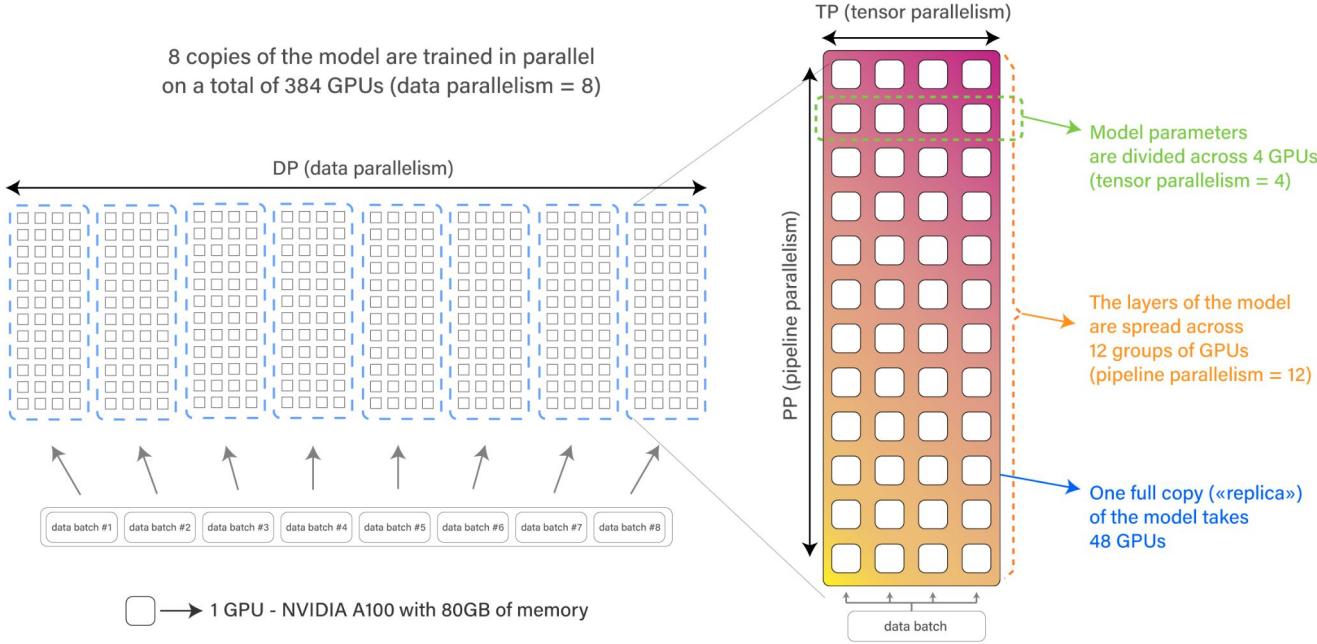
Each tensor is divided into several pieces so that instead of having the whole tensor residing on a single GPU each piece of the tensor resides on a different GPU



# Parallelism:

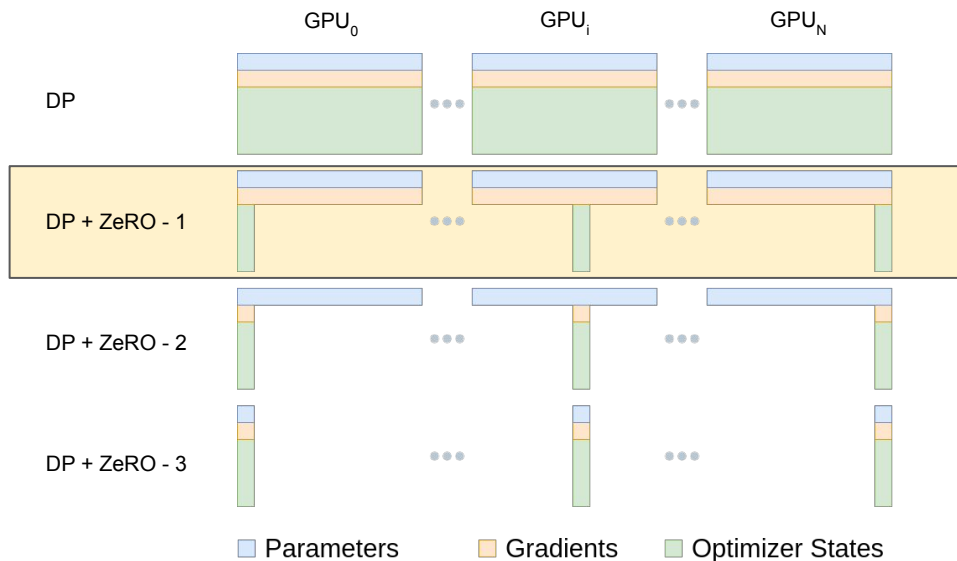
DP?  
TP?  
PP?

All 3 techniques were used!



# ZeRO data parallelism

- instead of replicating everything each GPU stores only a slice of it
- free the GPUs for larger batch sizes or more layers



Memory Consumption		Comm Volume
Formulation	Specific Example $K=12 \Psi=7.5B N_d=64$	
$(2 + 2 + K) * \Psi$	120GB	1x
$2\Psi + 2\Psi + \frac{K * \Psi}{N_d}$	31.4GB	1x
$2\Psi + \frac{(2+K) * \Psi}{N_d}$	16.6GB	1x
$\frac{(2 + 2 + K) * \Psi}{N_d}$	1.9GB	1.5x

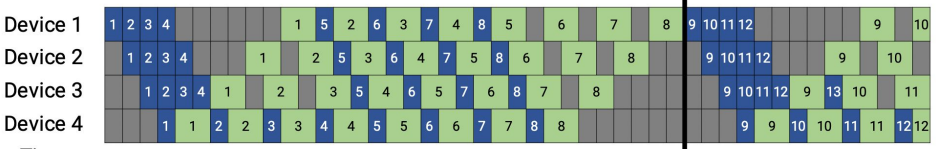
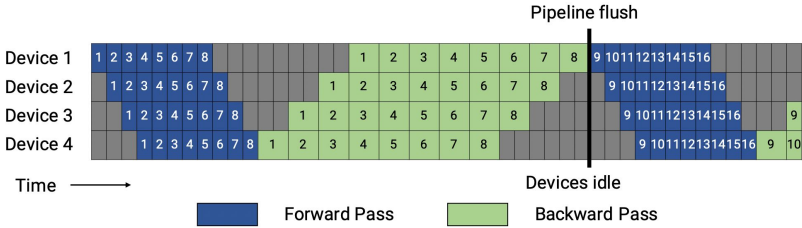
# Pipeline scheduling

👉 The one we used

All forward, all backward

Reduce memory →

One forward, one backward (1f1b)



↓ Reduce bubble at the cost of communication

Interleaved 1f1b



# BF16



Source: <https://mooabolic.medium.com/f684-fb32-fb16-bfloat16-f32-and-other-members-of-the-zoo-a1ca7897d407>



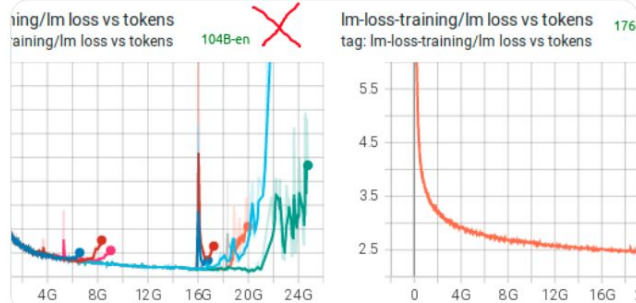
Stas Bekman  
@StasBekman

[1/2] What makes the @BigScienceLLM 176B-ml training so stable?

The 176B-ml succeeded to cross the 24B-token barrier whereas 104B-en failed.


We would love to hear your speculative and experiential reasoning for why this is so!


Following are the main candidates:





6:53 PM · Mar 20, 2022 · Twitter Web App


# Evaluation: is hard

 **Extrinsic Evaluation:** Focus on downstream, user-facing tasks

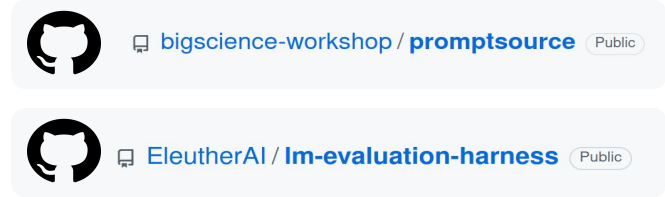
 **Intrinsic Evaluation:** Focus on encoding of linguistic and world knowledge

 **Bias/Social Impact:** Quantify encoding of stereotypes and risk of user harm

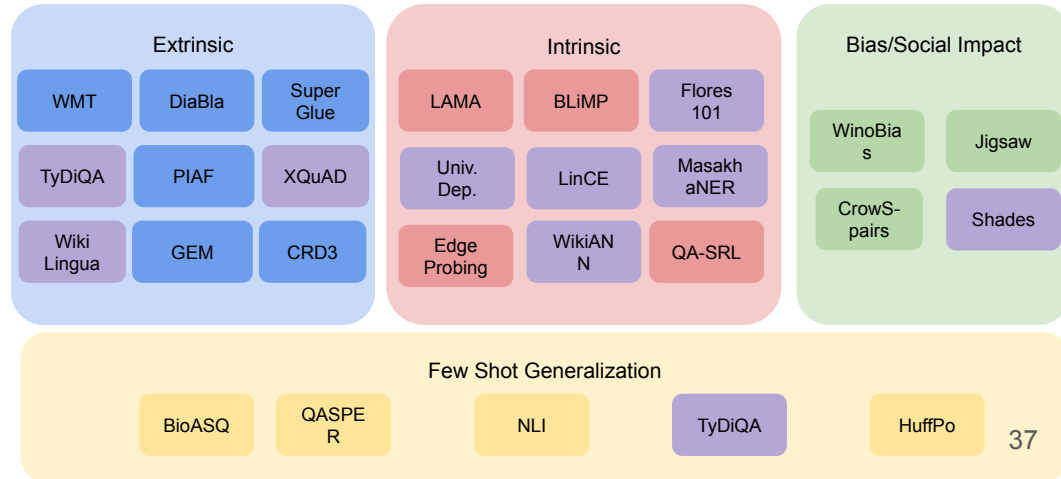
 **Multilingualism:** Ensure coverage of training and unseen language in all evaluations

 **Few-Shot Generalization:** Focus on evaluation on distributions not seen in pretraining

## Code Bases

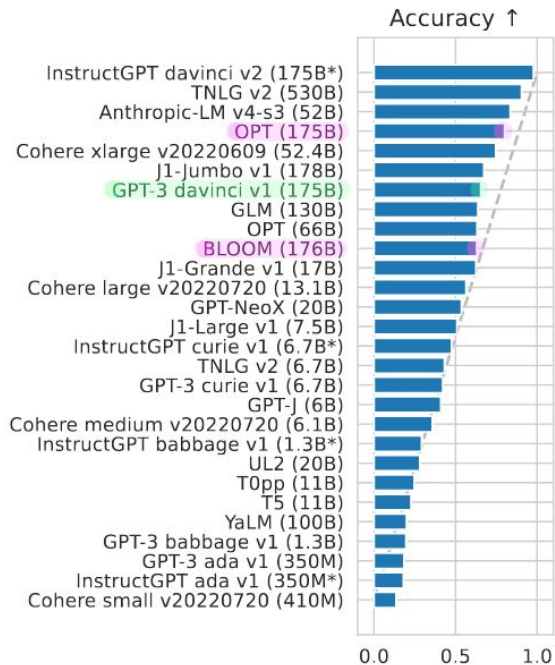


## Benchmark



After training

# HELM



In 2022, open source was ~1.5 years behind closed-source

# Zero-shot

Checkpoint: 65k steps (240B tokens)

To do a "farduddle" means to jump up and down really fast. An example of a sentence that uses the word farduddle is:

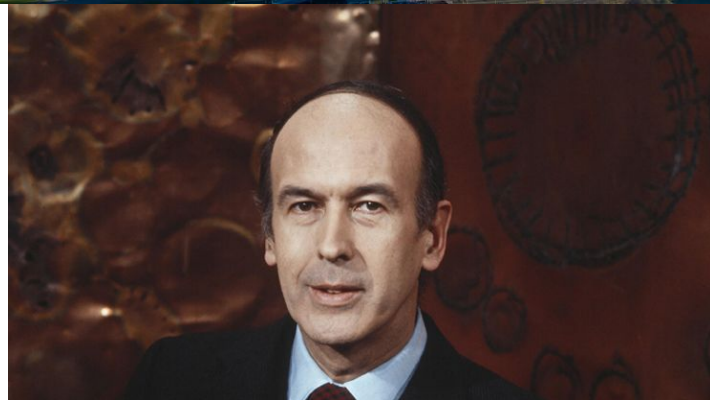
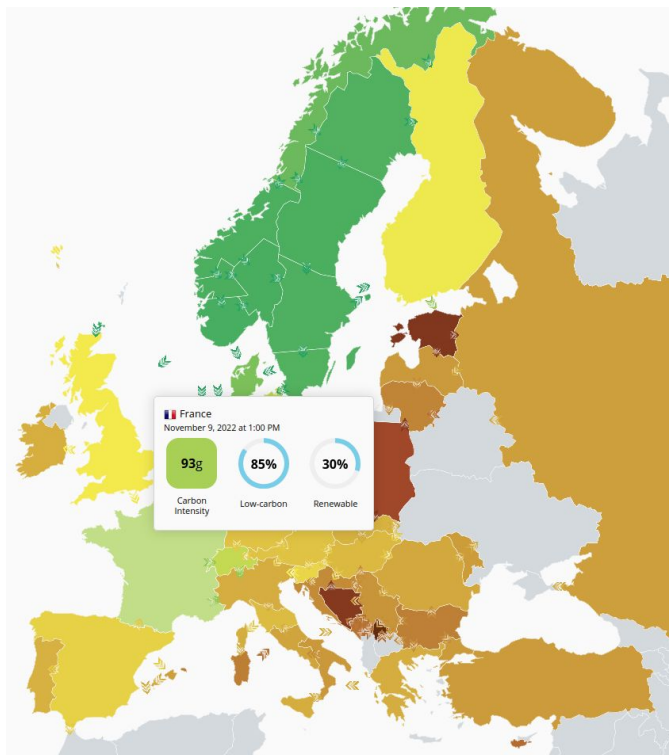
We were on a Kindrex toyshop hop and we got to see if Kindrex toys were built to withstand constant jumping, lots of jumping! To do this part, I had to farduddle to simulate jumping.

Prompt

Generated



# Carbon footprint



# Carbon footprint

Model name	Number of parameters	Datacenter PUE	Carbon intensity of grid used	Power consumption	CO <sub>2</sub> eq emissions	CO <sub>2</sub> eq emissions × PUE
GPT-3	175B	1.1	429 gCO <sub>2</sub> eq/kWh	1,287 MWh	502 tonnes	552 tonnes
Gopher	280B	1.08	330 gCO <sub>2</sub> eq/kWh	1,066 MWh	352 tonnes	380 tonnes
OPT	175B	1.09 <sup>2</sup>	231gCO <sub>2</sub> eq/kWh	324 MWh	70 tonnes	76.3 tonnes <sup>3</sup>
BLOOM	176B	1.2	57 gCO <sub>2</sub> eq/kWh	433 MWh	25 tonnes	30 tonnes

~ 20 Paris <-> NYC return trips

# Example generations

A poem about the Large Hadron Collider by Alfred Tennyson.

Title: The Particle at the End of the Universe

The **mindless swarm** that fills the void,  
The blind, unthinking, purposeless crowd  
Of atoms, blindly rushing to their doom,  
With never a thought of what they are or why,  
Or whence they came, or whither they must go;  
This is the end of all we see and know,

**Hark!** from the abyss there comes a sound  
As of a mighty wind that shakes the ground,  
A sound of rending, tearing, crashing sound,  
As if the solid earth were rent in twain,  
And all the elements in wild uproar  
Were hurled together in a chaotic storm;

And then a pause, and then again  
The sound as of ten thousand thunders rolled  
Upon the startled ear of night,  
Then silence

استخراج العدد العائلي في لغة بايثون:

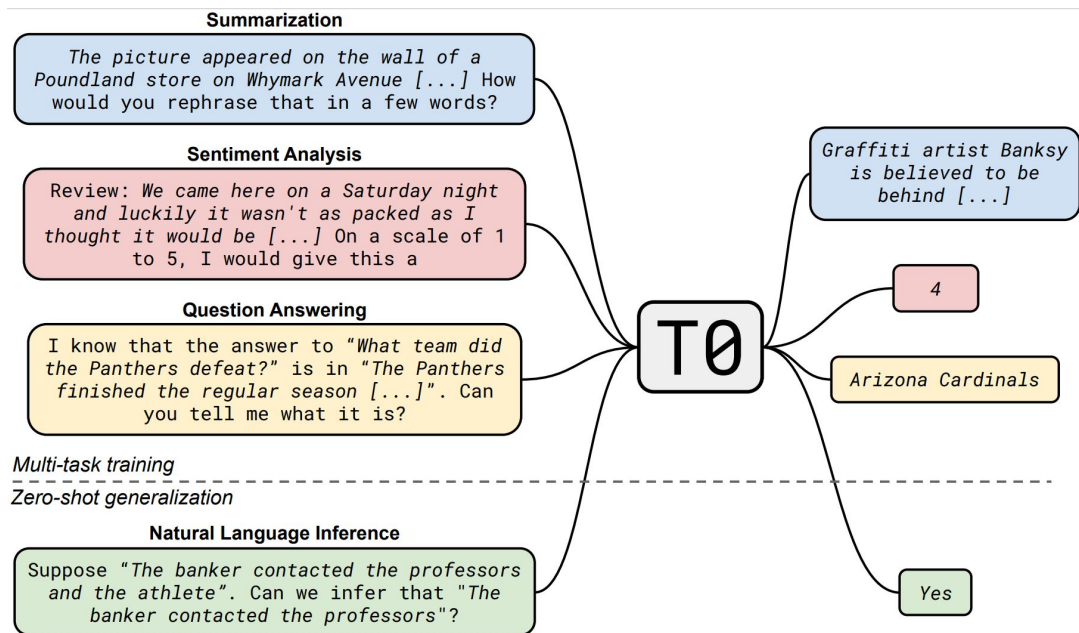
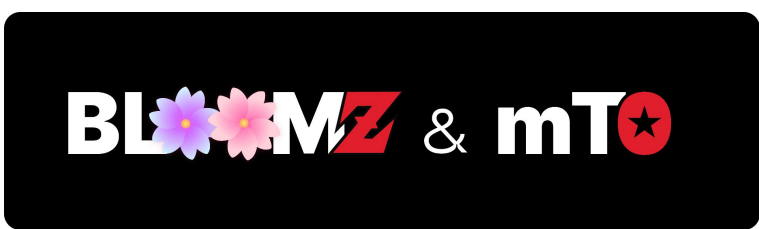
```
def factorial(n):  
    if n == 0:  
        return 1  
    else:  
        result = 1  
        for i in range(1, n + 1):  
            result *= i  
        return result
```

Traduce español de España a español de Argentina  
El coche es rojo - el auto es rojo  
El ordenador es nuevo - la computadora es nueva  
el bolígrafo es negro - lapicera es negra  
la nevera - heladera  
los zapatos - zapatillas  
las gafas - anteojos

**BLM**  
a BigScience initiative

Input 176B params  
Output 59 languages  
Open-access

# Multi-task fine-tuning



# Code models

**BigCode is an open scientific collaboration working on the responsible development of large language models for code**

Learn more...

Supported by:

**servicenow**



**Hugging Face**

<https://www.bigcode-project.org/>

Monolingual EU LMs?

# Scaling

Compute (total operations) =  $k * \text{Data} * \text{Parameters}$

As C (compute) goes up, how much should go into D (data size) vs N (params)?

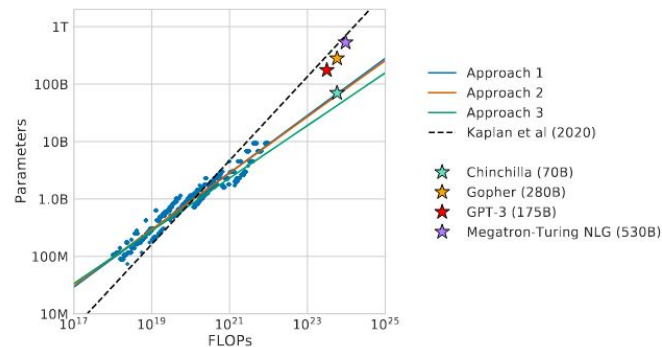
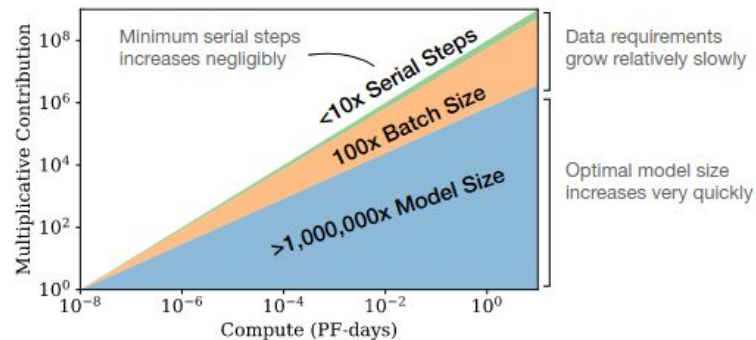
# The scaling controversy

Kaplan '20:

$$D \sim C^{0.27}, N \sim C^{0.73}$$

Hoffman '22:

$$D \sim C^{0.50}, N \sim C^{0.50}$$





# Data, data, data

BLOOM: 384 A100s for 111 days

->  $5.5E+23$  operations

Chinchilla optimality:

59B parameters, for 1,18T tokens

BLOOM:

176B parameters, for 366B tokens

# Data, data, data

So where do you find 1T tokens? ~4TB text, 800B words.

# Data, data, data

So where do you find 1T tokens? ~4TB text, 800B words.

- Even in English, finding 1T quality tokens is non-trivial
- Norwegian OSCAR: 2.8GB

# Data, data, data

So where do you find 1T tokens? ~4TB text, 800B words.

- Even in English, finding 1T quality tokens is non-trivial
- Norwegian OSCAR: 2.8GB

Libraries, silver bullet? The French National Library contains roughly 800B words

# Our current work



Niklas Muennighoff, HF



Nouamane Tazi, HF



Sampo Pyysalo, Turku



Thom Wolf, HF



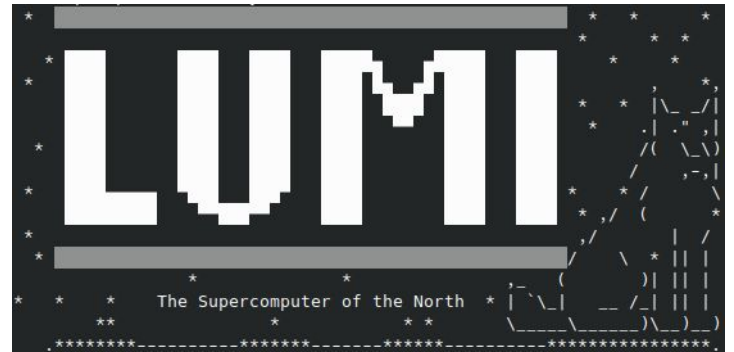
Ola Piktus, HF



Teven Le Scao, HF

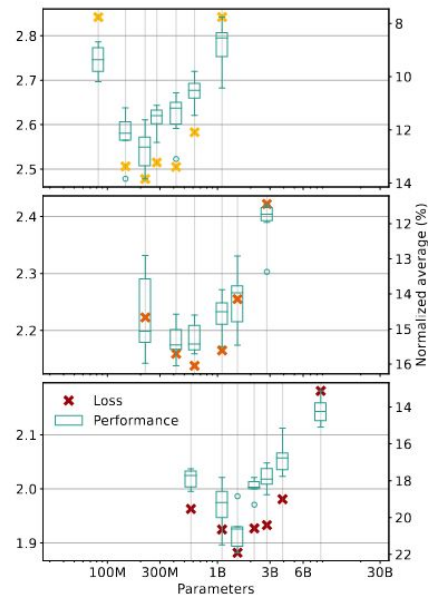


**UNIVERSITY  
OF TURKU**



# Val loss vs downstream zero-shot perf

Chinchilla-optimal models for loss are also optimal for 0-shot performance.



# Is repeating data bad?

Validation loss: yes

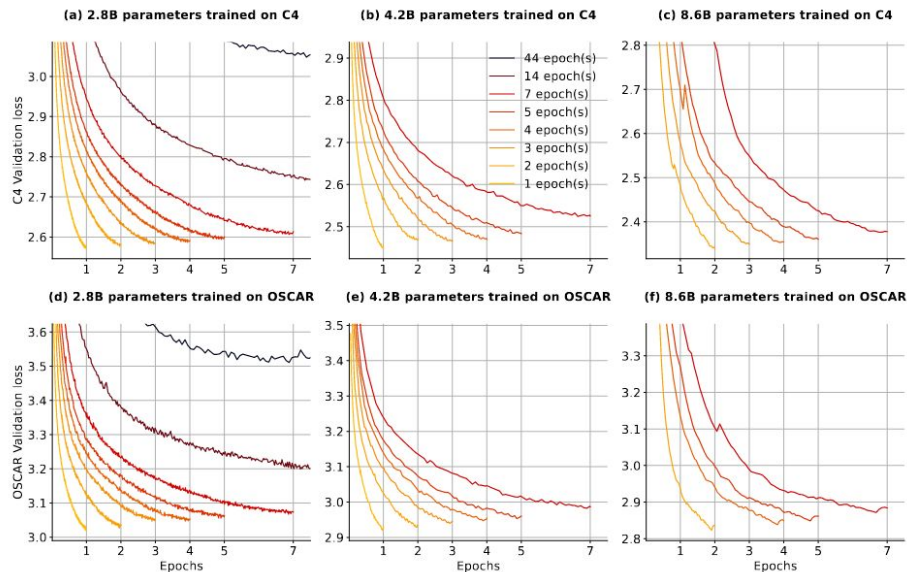


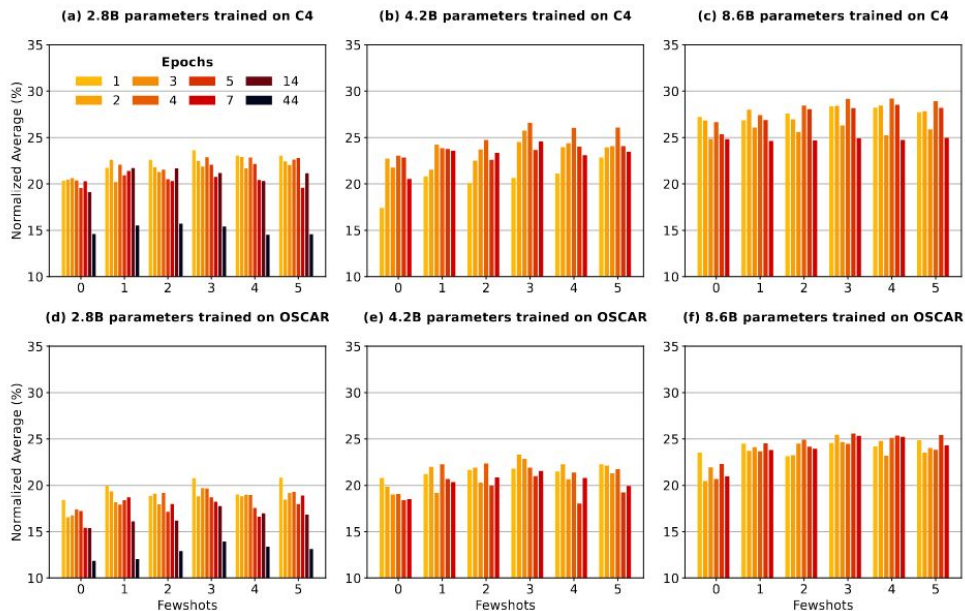
Figure 3: Epochs

# Is repeating data bad?

Zero- or few-shot downstream perf:

Not really.....

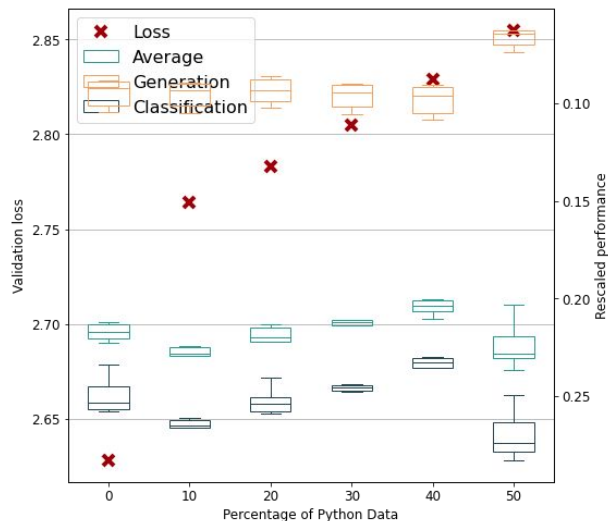
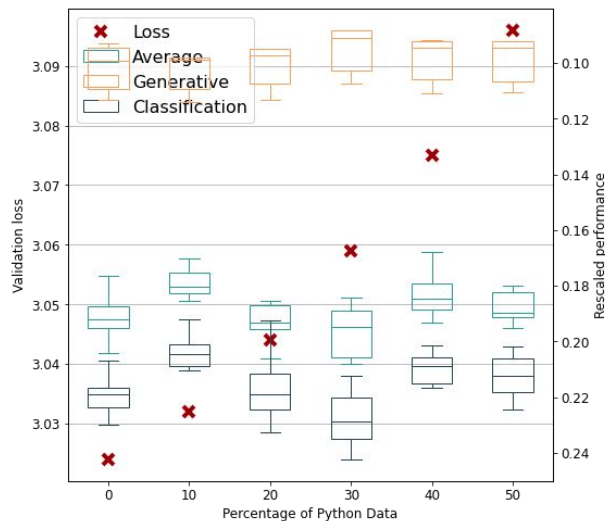
You can probably get away with x5-8





# Adding code data

50% Python data: no difference



Would it work for non-en languages?

# Multilinguality: when does it help?

## SCALING LAWS FOR GENERATIVE MIXED-MODAL LANGUAGE MODELS

Armen Aghajanyan<sup>\*†</sup>, Lili Yu<sup>\*†</sup>, Alexis Conneau<sup>†</sup>, Wei-Ning Hsu<sup>†</sup>

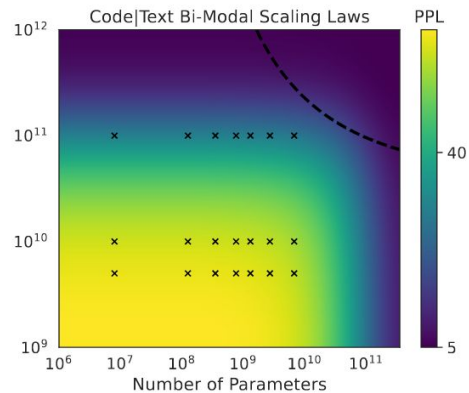
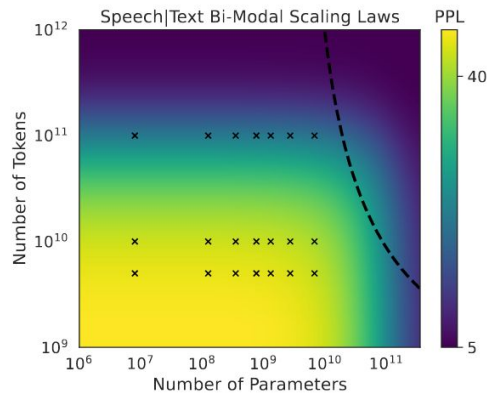
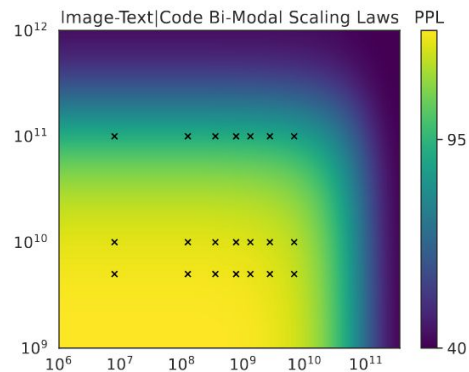
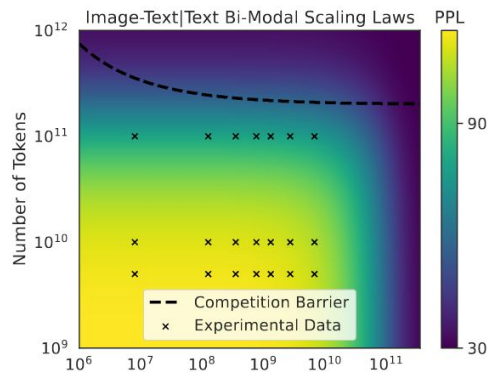
Karen Hambardzumyan<sup>◇</sup>, Susan Zhang<sup>†</sup>, Stephen Roller<sup>†</sup>, Naman Goyal<sup>†</sup>

Omer Levy<sup>†</sup> & Luke Zettlemoyer<sup>†,◇</sup>

FAIR<sup>†</sup>, University of Washington<sup>◇</sup>, YerevanN<sup>◇</sup>  
armenag@meta.com

### ABSTRACT

Generative language models define distributions over sequences of tokens that can represent essentially any combination of data modalities (e.g., any permutation of image tokens from VQ-VAEs, speech tokens from HuBERT, BPE tokens for language or code, and so on). To better understand the scaling properties of such mixed-modal models, we conducted over 250 experiments using seven different modalities and model sizes ranging from 8 million to 30 billion, trained on 5-100 billion tokens. We report new mixed-modal scaling laws that unify the contributions of individual modalities and the interactions between them. Specifically, we explicitly model the optimal synergy and competition due to data and model size as an additive term to previous uni-modal scaling laws. We also find four empirical phenomena observed during the training, such as emergent coordinate-ascent style training that naturally alternates between modalities, guidelines for selecting critical hyper-parameters, and connections between mixed-modal competition and training stability. Finally, we test our scaling law by training a 30B speech-text model, which significantly outperforms the corresponding unimodal models. Overall, our research provides valuable insights into the design and training of mixed-modal generative models, an important new class of unified models that have unique distributional properties.



# Multilinguality: when does it help?

Different languages ~ very close modalities

Find close languages, and compute competition barrier

Or even start from a code or other language model, then fine-tune it, the same way ChatGPT is descended from a Python LM

# RLHF

Instruct, Sparrow: ~50k human annotations

As easy in any language as in English



Questions?