

# DataComp

Creating large public Datasets for  
the AI open-source community



Institute for Foundations of  
**MACHINE LEARNING**

Alex Dimakis  
UT Austin  
Bespoke Labs

Datacomp is a collaboration with Ludwig Schmidt and the DataComp team

# Talk overview

- **Datacomp**: A scientific framework for dataset curation.
- 1. Datacomp for multimodal data (Neurips 23)  
2. Datacomp for Language Model Pre-training (DCLM) (Neurips 24)  
3. Ongoing now: Datacomp for Post-Training
- **Philosophical ponderings**: AI systems and the role of synthetic data curation.

# Talk overview

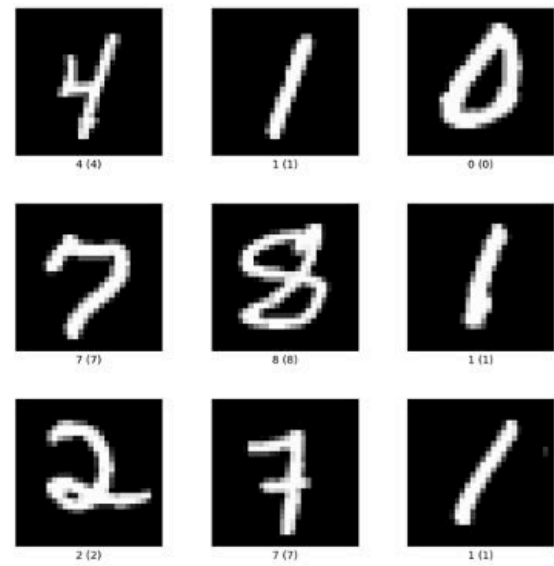
Ludwig Schmidt (in this building) told me about LAION. A dataset of billions of images and captions created from Common Crawl.

- LAION has been used to train Stable diffusion, and many other multimodal models. Also used for most open CLIP model pipelines.
- Idea: Make a better version of the LAION dataset. But also:  
Open-source the **tooling** for dataset curation.  
Create **scientific standardized benchmarks** for dataset curation  
Create a **community**.

# Datasets are the foundation of progress in AI

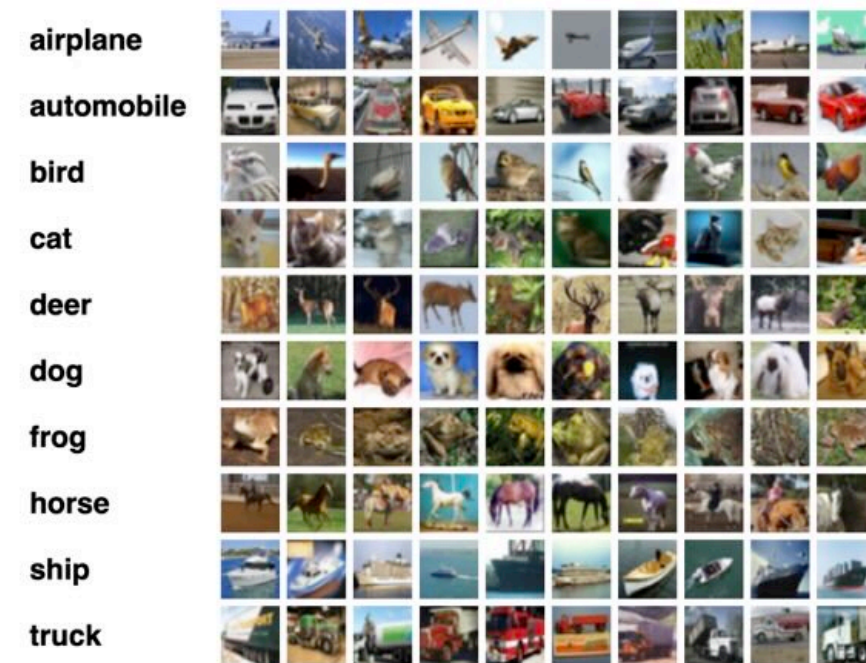
- For text:
- GPT-1 (2018): 3 B Tokens
- GPT-2 (2019): 30 B Tokens
- GPT-3 (2020): 300 B Tokens
- GPT-4 (2013): 3000? B Tokens
- 1000x growth in 5 years
- For images:
- Imagenet (2009): 1 Million images
- LAION-5B (2022): 5 Billion Images
- 5000x growth in 5 years

# ML Discoveries **enabled** by datasets



## MNIST (1994)

Convolutional  
neural networks



CIFAR-10 (2009)



## ImageNet (2012)

Deep learning  
resurgence, ResNets,  
transfer learning, etc.



## WebImageText (2021)

Zero-shot classification


# Data is a poorly understood (yet crucial) component of LLMs

LLMs rely on **trillions of tokens** of training data crawled from the Internet.

Details about training sets have become sparse for state of the art models.

Amongst open-source models, significant gap between **closed** v.s. **open** dataset models.

Average of 22 standard downstream evaluations



Model	Params	Tokens	Open dataset?	CORE	MMLU
<b>Open weights, closed datasets</b>					
Llama2	7B	2T	X	49.2	45.8
DeepSeek	7B	2T	X	50.7	48.5
Mistral-0.3	7B	?	X	57.0	62.7
QWEN-2	7B	?	X	57.5	<b>71.9</b> ←
Llama3	8B	15T	X	57.6	66.2
Gemma	8B	6T	X	57.8	64.3
Phi-3	7B	?	X	<b>61.0</b>	69.9
<b>Open weights, open datasets</b>					
Falcon	7B	1T	✓	44.1	27.4
OLMo-1.7	7B	2.1T	✓	47.0	54.0
MAP-Neo	7B	4.5T	✓	<b>50.2</b>	<b>57.1</b> ←


# Data is a poorly understood (yet crucial) component of LLMs

LLMs rely on **trillions of tokens** of training data crawled from the Internet.



Details about training sets have become sparse for state of the art models.

Amongst open-source models, significant gap between **closed** v.s. **open** dataset models.

Average of 22 standard downstream evaluations



Model	Params	Tokens	Open dataset?	CORE	MMLU
<b>Open weights, closed datasets</b>					
Llama2	7B	2T	X	49.2	45.8
DeepSeek	7B	2T	X	50.7	48.5
Mistral-0.3	7B	?	X	57.0	62.7
QWEN-2	7B	?	X	57.5	<b>71.9</b>
Llama3	8B	15T	X	57.6	66.2
Gemma	8B	6T	X	57.8	64.3
Phi-3	7B	?	X	<b>61.0</b>	69.9
<b>Open weights, open datasets</b>					
Falcon	7B	1T	✓	44.1	27.4
OLMo-1.7	7B	2.1T	✓	47.0	54.0
MAP-Neo	7B	4.5T	✓	<b>50.2</b>	<b>57.1</b>



Spoiler:

We got an open data model with MMLU 64

# A problem with previous work on dataset curation

**Scaling Language Models & Insights from Training DeepMind GPT-4**  
Jack W. Rae, Sebastian Borgeaud, Trevor Cai, Katie Mi Sarah Henderson, Roman Ring, Susannah Young, Eliza

**SlimPajama-DC: Understanding Data Combinations for LLM Training**  
Zhiqiang Shen<sup>†</sup> Tianhua Tao<sup>†,‡</sup> Liqun Ma<sup>†</sup> Willie Neiswanger<sup>§</sup>

**How to Train Data-Efficient LLMs**  
Kang<sup>1</sup> Jianmo Ni<sup>1</sup> Lichan Hong<sup>1</sup> Ed H. Chi<sup>1</sup>  
<sup>2</sup> Derek Zhiyuan Cheng<sup>1</sup>

**QuRating: Selecting High-Quality Data for Training Language Models**  
Alexander Wettig<sup>1</sup> Aatmik Gupta<sup>1</sup> Saumya Malik<sup>1</sup> Danqi Chen<sup>1</sup>

**The RefinedWeb Dataset for Falcon LLM: Training Curated Corpora with Web Data, and Web Data Only**  
The Falcon LLM team  
Guilherme Penedo<sup>1</sup> Quentin Malartic<sup>2</sup>  
<sup>1</sup> Ruxandra Cojocaru<sup>2</sup> Alessandro Cappelli<sup>1</sup> Hamza Alobeidli<sup>2</sup> Baptiste Pannier<sup>1</sup>  
Ebtessam Almazrouei<sup>2</sup> Julien Launay<sup>1,3</sup>

**FineWeb-Edu**  
The finest collection of educational content the web has to offer

**d4m: an Open Corpus of Text for Language Model Pretraining**  
Luca Soldaini<sup>α</sup> Rodney Kinney<sup>α</sup> Akshita Bh  
David Atkinson<sup>α</sup> Russell Authur<sup>α</sup> Ben Bogir  
Jennifer Dumas<sup>α</sup> Yanai Elazar<sup>α,ω</sup> Valentin Hofn  
Sachin Kumar<sup>α</sup> Li Lucy<sup>β</sup> Xinxu Lyu<sup>ω</sup> Nathan I  
Jacob Morrison<sup>α</sup> Niklas Muennighoff<sup>α</sup> Aakank  
Matthew E. Peters<sup>σ</sup> Abhilasha Ravichander<sup>α</sup> Kyle  
Emma Strubell<sup>x,α</sup> Nishant Subramani<sup>x,α</sup> Oyv  
Luke Zettlemoyer<sup>ω</sup> Noah A. Smith<sup>α,ω</sup> Ha  
Iz Beltagy<sup>α</sup> Dirk Groeneveld<sup>α</sup> J  
Kyle Lo<sup>α</sup>

Travis Hoppe Charles Foster Jason Phang Horace He ret Mitchell Matt Gardner



# A problem with previous work on dataset curation

The image shows a collage of research papers related to LLM dataset curation. A central blue callout box contains the text: "Different evaluation and training protocols result in apples-to-oranges comparisons".

**DeepMind**  
Scaling Language Models & Insights from Training GPT-4  
Jack W. Rae, Sebastian Borgeaud, Trevor Cai, Katie Mirowski, Sarah Henderson, Roman Ring, Susannah Young, Eliza S. Rutherford, Eric Nijkamp, David Shih, Paul Davidsson, Taehoon Kim, Shreya Pathak, Sander Stachiw, Nemanja Djokic, Markus B. Duwig, Joelle Barral, David Dale, Pieter H. Schaul, Ed H. Chi

**SlimPajama-DC: Understanding Data Combinations for LLM Training**  
Zhiqiang Shen<sup>†</sup> Tianhua Tao<sup>†,‡</sup> Liqun Ma<sup>†</sup> Willie Neiswanger<sup>§</sup>

**The F**  
Leo  
Travis Hoppe

**Ed H. Chi<sup>1</sup>**

**doqa: an Open Dataset for Language Model Pretraining**  
Luca Soldaini<sup>α</sup> Rodney Kinney<sup>α</sup> Akshita Bhosale<sup>α</sup>  
David Atkinson<sup>α</sup> Russell Authur<sup>α</sup> Ben Bogin<sup>α</sup>  
Jennifer Dumas<sup>α</sup> Yanai Elazar<sup>α,ω</sup> Valentin Hofmann<sup>α</sup>  
Sachin Kumar<sup>α</sup> Li Lucy<sup>β</sup> Xinxin Lyu<sup>ω</sup> Nathan Lambert<sup>α</sup>  
Jacob Morrison<sup>α</sup> Niklas Muennighoff<sup>α</sup> Aakanksha Mateika<sup>α</sup>  
Matthew E. Peters<sup>σ</sup> Abhilasha Ravichander<sup>α</sup> Kyle Richardson<sup>α</sup>  
Emma Strubell<sup>α,α</sup> Nishant Subramani<sup>α,α</sup> Oyvind Tafjord<sup>α</sup>  
Luke Zettlemoyer<sup>ω</sup> Noah A. Smith<sup>α,ω</sup> Hanukh Weininger<sup>α</sup>  
Iz Beltagy<sup>α</sup> Dirk Groeneveld<sup>α</sup> Jesse Alpert<sup>α</sup>  
Kyle Lo<sup>α</sup>

**FineWeb-Edu**  
The finest collection of educational content the web has to offer

**Diversification**  
The RefinedWeb Dataset for Falcon LLM: Mining Curated Corpora with Web Data, and Web Data Only  
The Falcon LLM team  
Guilherme Penedo<sup>1</sup> Quentin Malartic<sup>2</sup>  
<sup>1</sup> Ruxandra Cojocaru<sup>2</sup> Alessandro Cappelli<sup>1</sup> Hamza Alobeidli<sup>2</sup> Baptiste Pannier<sup>1</sup>  
Ebtessam Almazrouei<sup>2</sup> Julien Launay<sup>1,3</sup>

# The DataComp idea: Shift to Data-Centric AI

Traditional ML Benchmark: Dataset fixed (e.g. Imagenet),  
improve the model (Model Centric)

DataComp: Model training + Evaluation fixed.  
Improve the dataset. (Data Centric)

# DATAComp:

## In search of the next generation of multimodal datasets

Samir Yitzhak Gadre\*<sup>2</sup> Gabriel Ilharco\*<sup>1</sup> Alex Fang\*<sup>1</sup> Jonathan Hayase<sup>1</sup> Georgios Smyrnis<sup>5</sup>  
Thao Nguyen<sup>1</sup> Ryan Marten<sup>7,10</sup> Mitchell Wortsman<sup>1</sup> Dhruva Ghosh<sup>1</sup> Jieyu Zhang<sup>1</sup> Eyal Orgad<sup>3</sup>  
Rahim Entezari<sup>11</sup> Giannis Daras<sup>5</sup> Sarah Pratt<sup>1</sup> Vivek Ramanujan<sup>1</sup> Yonatan Bitton<sup>12</sup>  
Kalyani Marathe<sup>1</sup> Stephen Mussmann<sup>1</sup> Richard Vencu<sup>6</sup> Mehdi Cherti<sup>6,8,9</sup> Ranjay Krishna<sup>1</sup>  
Pang Wei Koh<sup>1</sup> Olga Saukh<sup>11</sup> Alexander Ratner<sup>1</sup> Shuran Song<sup>2</sup> Hannaneh Hajishirzi<sup>1,7</sup>  
Ali Farhadi<sup>1</sup> Romain Beaumont<sup>6</sup> Sewoong Oh<sup>1</sup> Alexandros G. Dimakis<sup>5</sup> Jenia Jitsev<sup>6,8,9</sup>  
Yair Carmon<sup>3</sup> Vaishaal Shankar<sup>4</sup> Ludwig Schmidt<sup>1,6,7</sup>

### Abstract

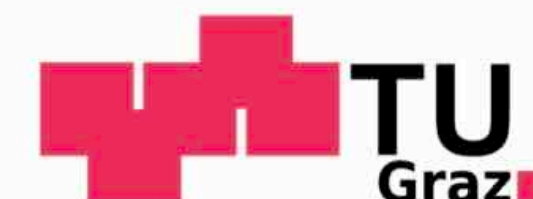
Large multimodal datasets have been instrumental in multiple breakthroughs like CLIP, DALL-E, Stable Diffusion, Flamingo and GPT-4, yet datasets rarely receive the same research attention as model architectures or training algorithms. To address this shortcoming in the machine learning ecosystem, we introduce DATAComp, a participatory benchmark where the training code is fixed and researchers innovate by proposing new training sets. Concretely, we provide an experimental testbed centered around a new



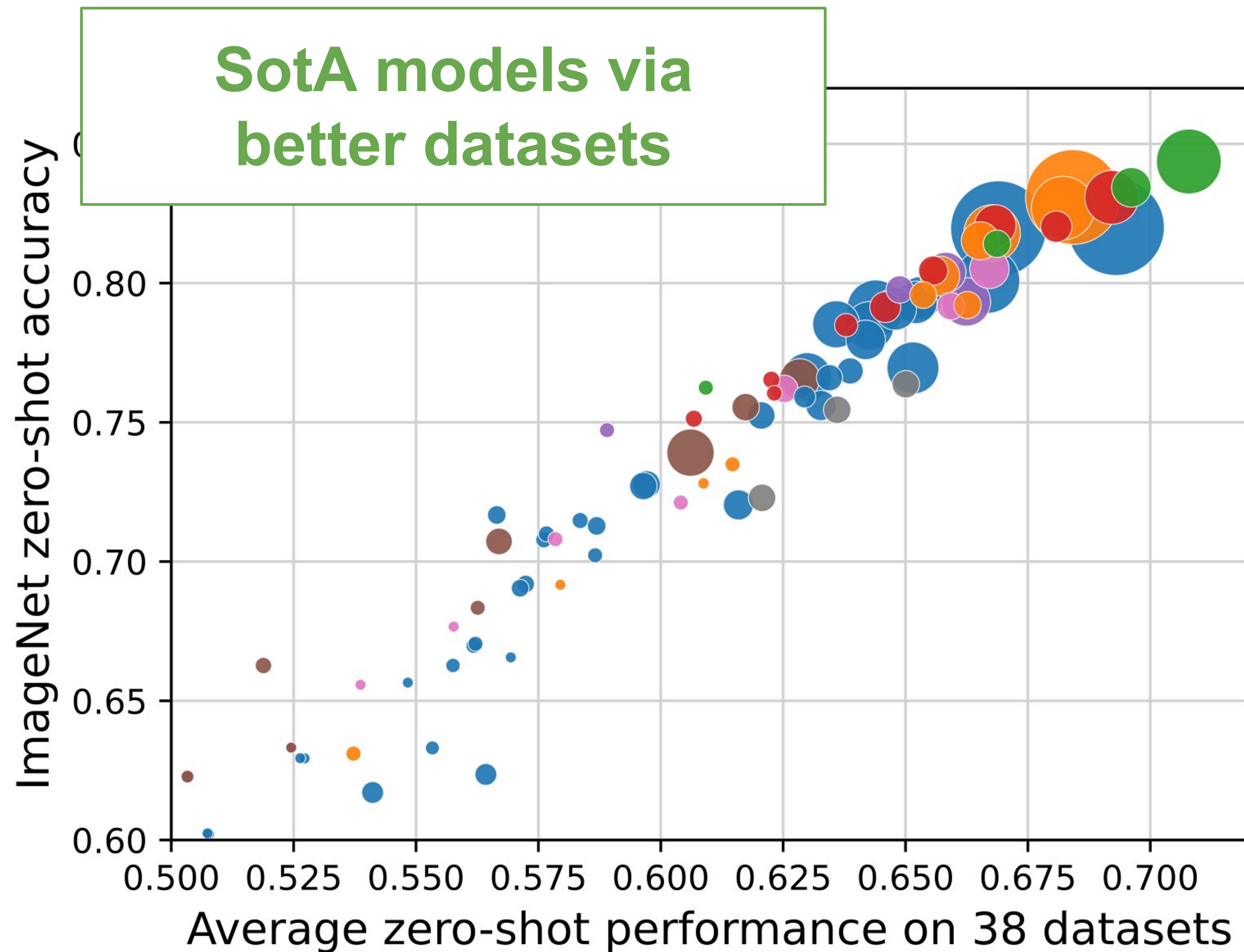
האוניברסיטה העברית בירושלים  
THE HEBREW UNIVERSITY OF JERUSALEM



UNIVERSITY OF  
ILLINOIS  
URBANA-CHAMPAIGN



# Impact of DataComp



## Dataset

- LAION
- DataComp
- WebLI
- LAION+COYO
- OpenAI WIT
- MetaCLIP
- CommonPool
- DFN

## FLOPs (B)

- 400
- 800
- 1200
- 1600
- 2000

# Impact of DataComp

huggingface.co/datasets/mlfoundations/datacomp\_pools

Google Scholar Outlook mail The Information Hacker News TechCrunch Join a Meeting - Z... https://voyaretire... Reddit Machine Le... DataTau reddit Austin, Texas All Bookmarks

Hugging Face

Search models, datasets, users...

Models

Datasets

Spaces

Posts

Docs

Pricing

Datasets: mlfoundations/datacomp\_pools

like 15

Following ML Foundations 80

Modalities: Image

License: cc-by-4.0

Dataset card

Viewer

Files and versions

Community 2

Dataset Preview

API

Embed

Full Screen Viewer

Split (1)  
train

The full dataset viewer is not available (click to read why). Only showing a preview of the rows.

uid  
string

url  
string

text  
string

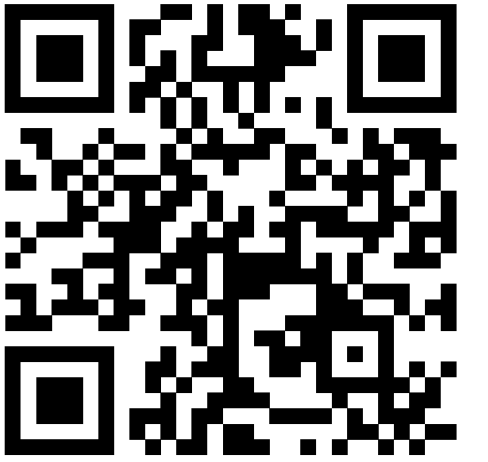
Downloads last month

826,258

Edit dataset card

Data Sourcing report

powered by Spawning.ai

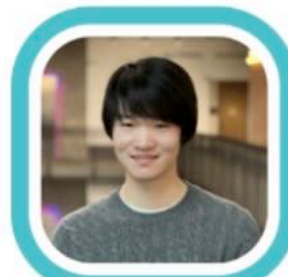


# DataComp - LM

In search of the next generation of training sets for language models



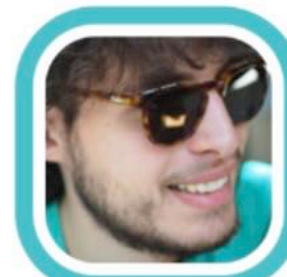
# The DCLM Team



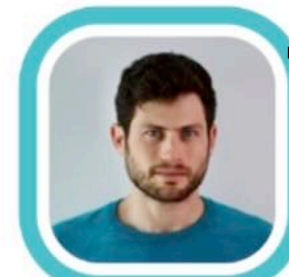
**Alex Fang**  
University of Washington



**Jeffrey Li**  
University of Washington



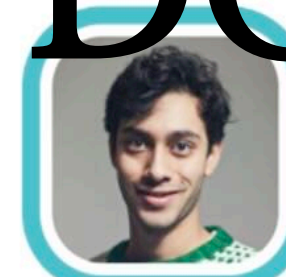
**Georgios Smyrnis**  
UT Austin



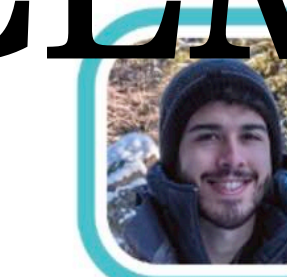
**Maor Ivgi**  
Tel Aviv University



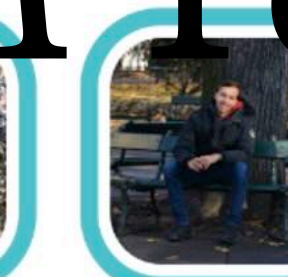
**Matt Jordan**  
UT Austin



**Samir Gadre**  
Columbia University



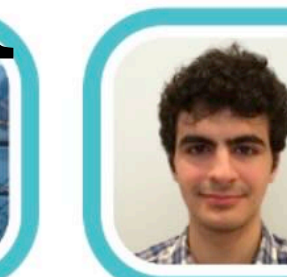
**Gabriel Ilharco**  
University of Washington



**Giannis Daras**  
UT Austin



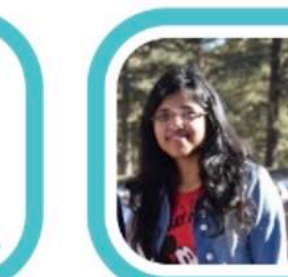
**Kalyani Marathe**  
University of Washington



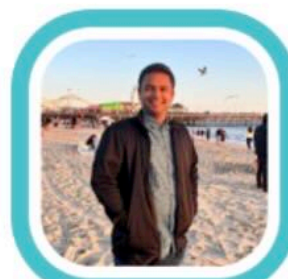
**Aaron Gokaslan**  
Cornell University



**Jieyu Zhang**  
University of Washington



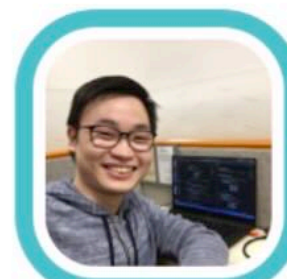
**Khyathi Chandu**  
AI2



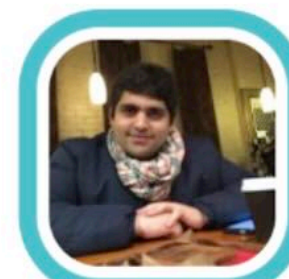
**Hritik Bansal**  
UCLA



**Etash Guha**  
University of Washington



**Sedrick Keh**  
Toyota Research Institute



**Kushal Arora**  
Toyota Research Institute



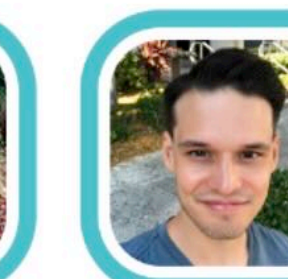
**Saurabh Garg**  
Carnegie Mellon University



**Rui Xin**  
University of Washington



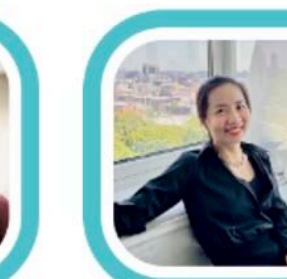
**Thao Nguyen**  
University of Washington



**Igor Vasiljevic**  
Toyota Research Institute



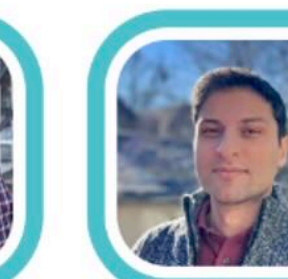
**Sham Kakade**  
Harvard University



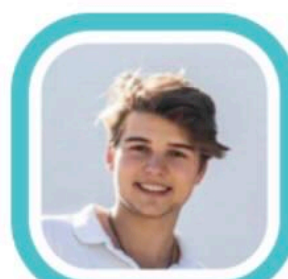
**Shuran Song**  
Stanford University



**Sujay Sanghavi**  
UT Austin



**Fartash Faghri**  
Apple



**Niklas Muennighoff**  
Contextual AI



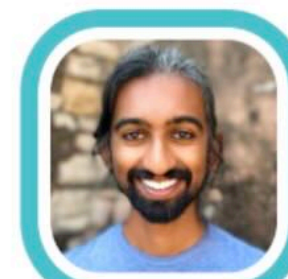
**Reinhard Heckel**  
Technical University of Munich



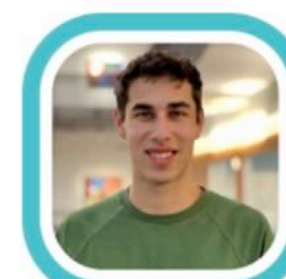
**Jean Mercat**  
Toyota Research Institute



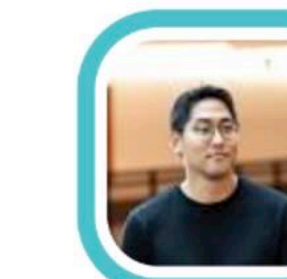
**Mayee Chen**  
Stanford University



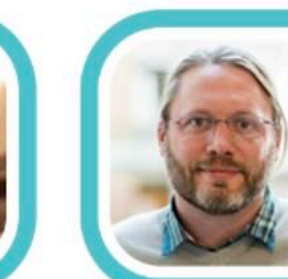
**Suchin Gururangan**  
University of Washington



**Mitchell Wortsman**  
University of Washington



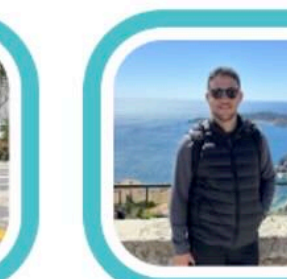
**Sewoong Oh**  
University of Washington



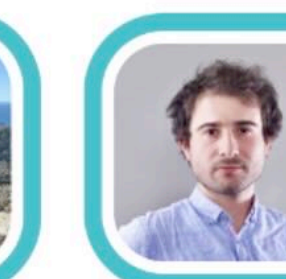
**Luke Zettlemoyer**  
University of Washington



**Kyle Lo**  
AI2



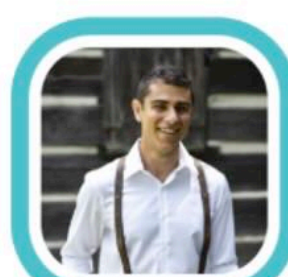
**Alaa El-Nouby**  
Apple



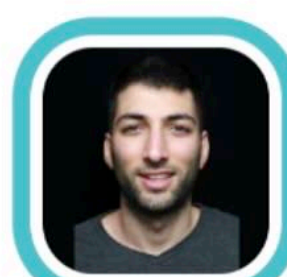
**Hadi Pour Ansari**  
Apple



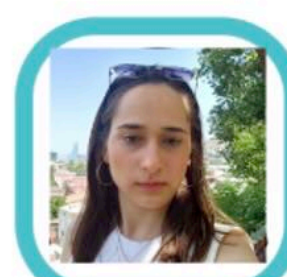
**Alexander Toshev**  
Apple



**Alon Albalak**  
UCSB



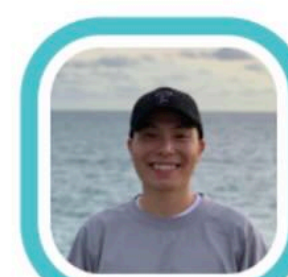
**Yonatan Bitton**  
Hebrew University



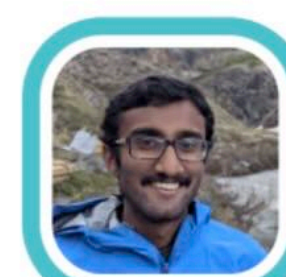
**Marianna Nezhurina**  
Research Center Juelich & JSC & LAION



**Amro Abbas**  
Datology AI



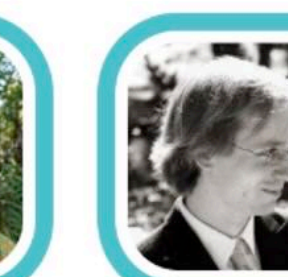
**Cheng-Yu Hsieh**  
University of Washington



**Dhruba Ghosh**  
University of Washington



**Stephanie Wang**  
University of Washington



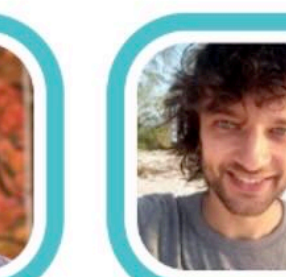
**Dirk Groeneveld**  
AI2



**Luca Soldaini**  
AI2



**Pang Wei Koh**  
University of Washington



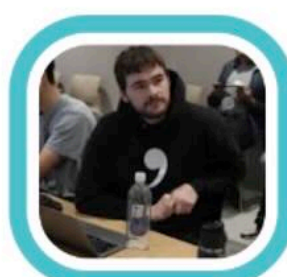
**Jenia Jitsev**  
Research Center Juelich & JSC & LAION



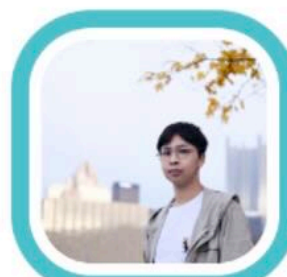
**Thomas Kollar**  
Toyota Research Institute



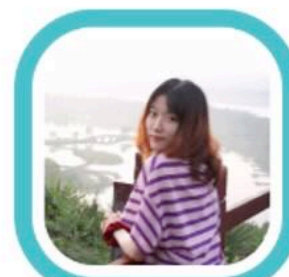
**Joshua Gardner**  
University of Washington



**Maciej Kilian**  
USC



**Hanlin Zhang**  
Harvard University



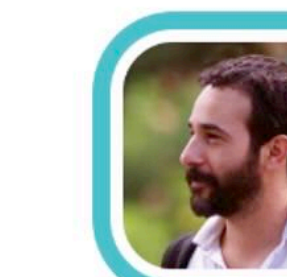
**Rulin Shao**  
University of Washington



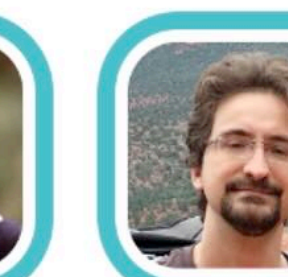
**Sarah Pratt**  
University of Washington



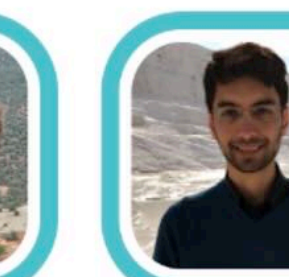
**Sunny Sanyal**  
UT Austin



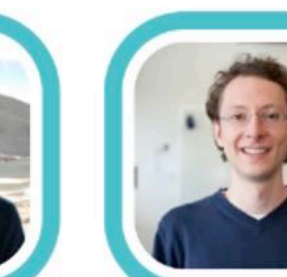
**Alex Dimakis**  
UT Austin BespokeLabs.ai



**Yair Carmon**  
Tel Aviv University



**Achal Dave**  
Toyota Research Institute



**Ludwig Schmidt**  
University of Washington, Stanford



**Vaishaal Shankar**  
Apple

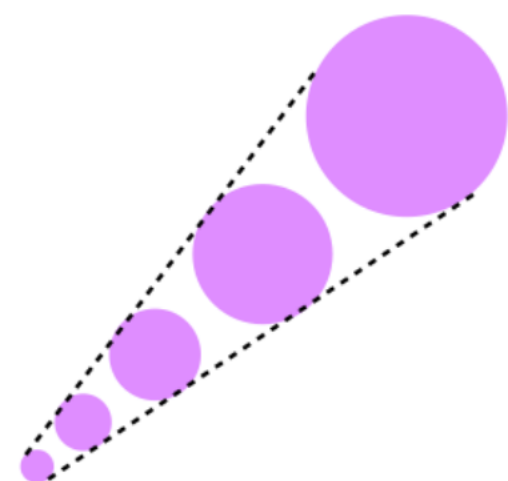
# DCLM Contributions

- A benchmark for language model data curation
- 416 experiments to identify the key parts of an effective data curation pipeline
- A state-of-the-art public training set for 7B parameter models
- A leaderboard for the community to try different data curation methods
- Our best model, DCLM 7B achieves 64 MMLU, trained only on 2T tokens.
- This is better than Llama2 models while Llama3 8B was trained on 6x times more tokens (i.e. 6 times more compute efficient, due to better data curation!).



# DCLM benchmark

## A. Select a scale



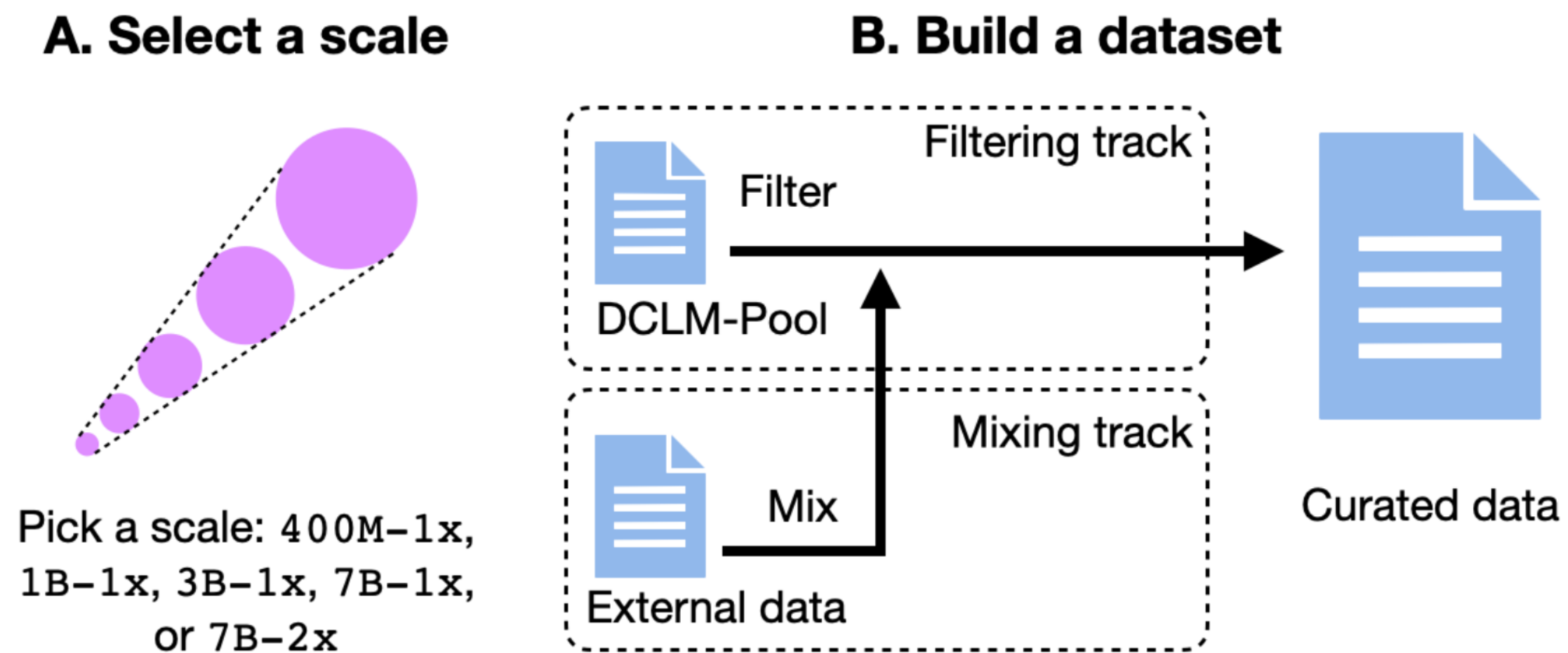
Pick a scale: 400M-1x,  
1B-1x, 3B-1x, 7B-1x,  
or 7B-2x

Scale	Model parameters	Train tokens	Train FLOPs	Train H100 hours	Pool size
400M-1x	412M	8.2B	2.0e19	26	469B
1B-1x	1.4B	28.8B	2.4e20	240	1.64T
3B-1x	2.8B	55.9B	9.4e20	740	3.18T
7B-1x	6.9B	138B	5.7e21	3,700	7.85T
7B-2x	6.9B	276B	1.1e22	7,300	15.7T

We fix **initial data pools**, **training recipes**, and **evaluations**.

Participants submit **datasets** that can then be fairly compared to one another.

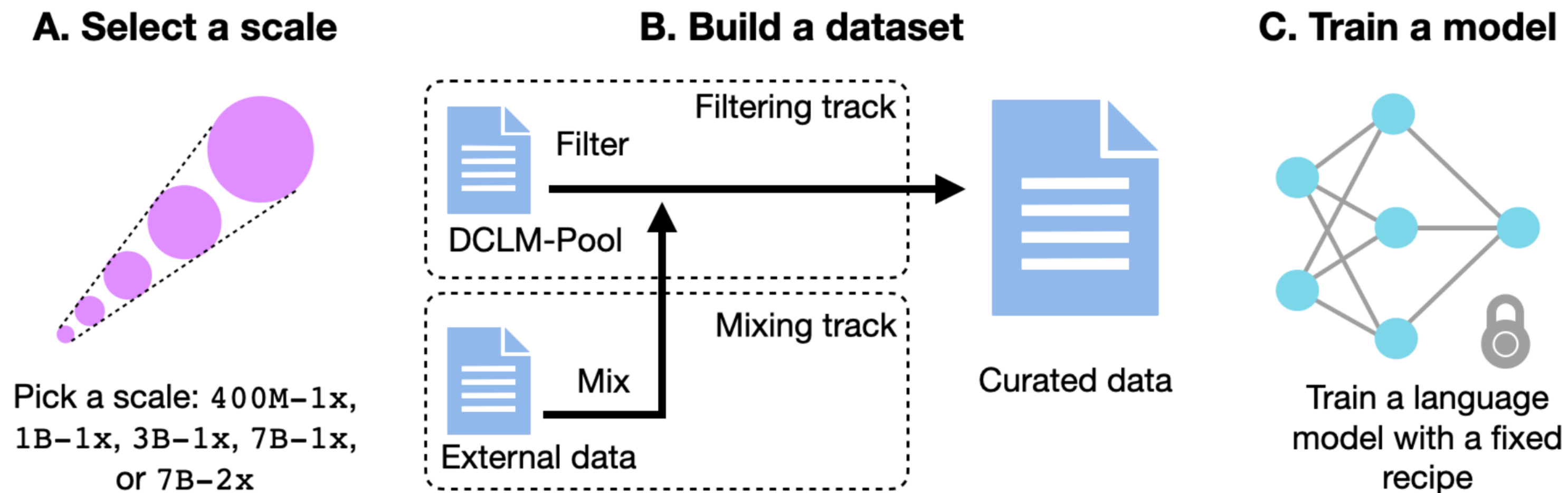
# DCLM benchmark



We fix **initial data pools**, **training recipes**, and **evaluations**.

Participants submit **datasets** that can then be fairly compared to one another.

# DCLM benchmark

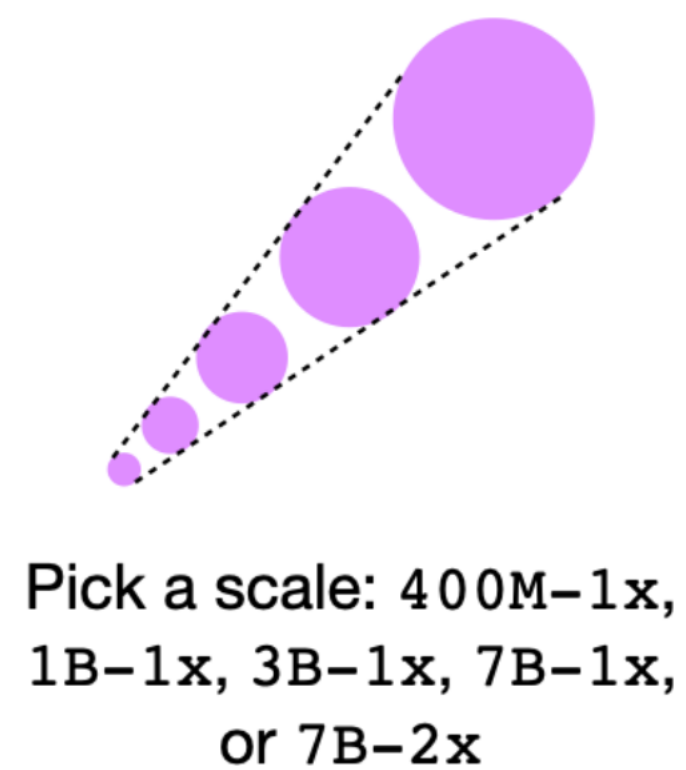


We fix **initial data pools**, **training recipes**, and **evaluations**.

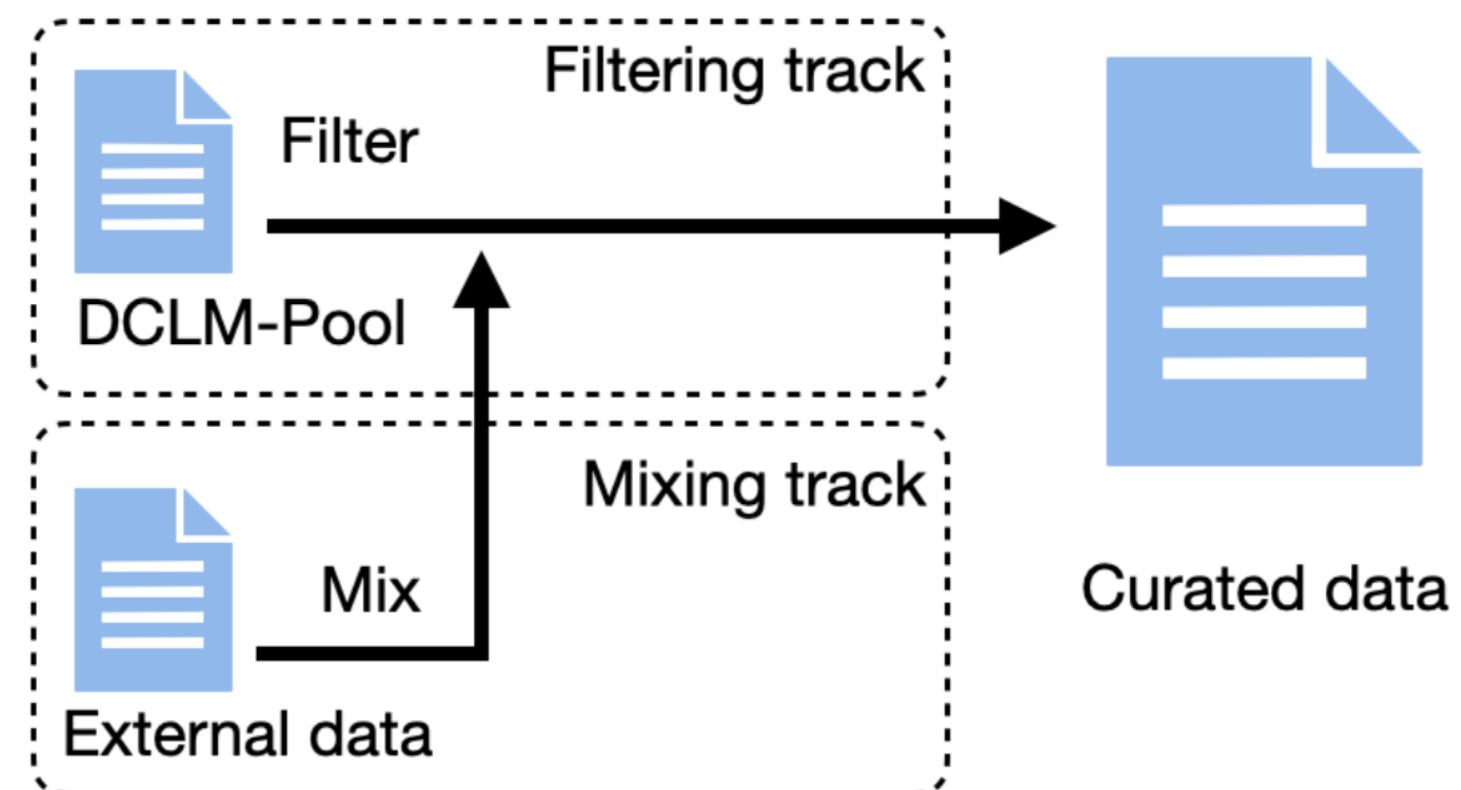
Participants submit **datasets** that can then be fairly compared to one another.

# DCLM benchmark

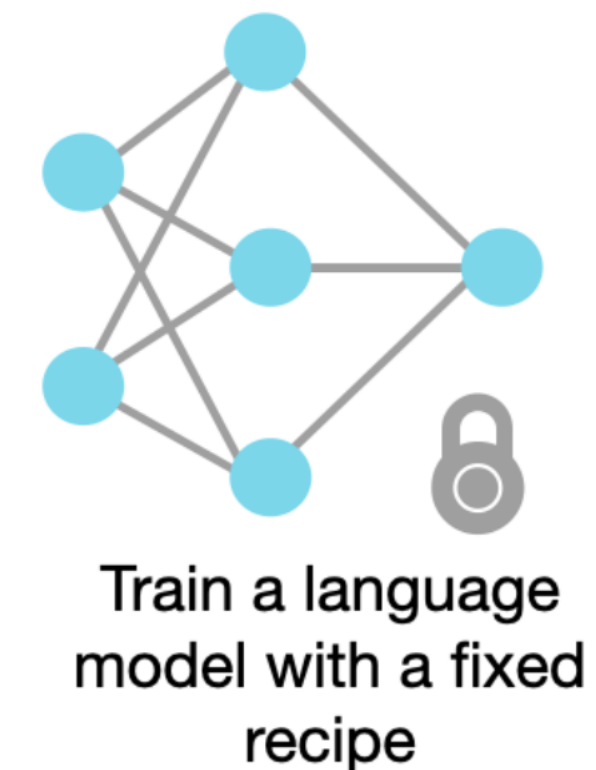
## A. Select a scale



## B. Build a dataset



## C. Train a model



## D. Evaluate



We fix **initial data pools**, **training recipes**, and **evaluations**.

Participants submit **datasets** that can then be fairly compared to one another.

# Lessons from DCLM

1. Data curation algorithms that work at 400M scale **can indeed predict performance at bigger sizes** (at least up to 7B).
2. That means we can do data curation science by developing data curation algorithms at small scale and extrapolating.
3. You don't have to build the whole ship and throw it in the ocean to check if it floats. Build a miniature ship and test in the bathtub.

# Lessons from DCLM

1. Deduplication and data cleaning are very important.  
(Bloom filters, Min-hash, a lot more to do here)
2. Humans are terrible at deciding which paragraphs will be good vs bad pretraining data. We found this very surprising.
3. The best data filter we found is a FastText Classifier trained on Instruction tuning vs CommonCrawl.  
This outperformed everything else we tried!
4. Once you have created a good-quality filtered dataset, no mixing methods seemed to help further which was surprising.  
(Right ways to Ensemble datasets?)

# Effective data filtering

We use a combination of the RefinedWeb heuristic filters + a trained classifier.

[Penedo et al., 2023]

Filter	CORE	EXTENDED
RefinedWeb reproduction	27.5	14.6
Top 20% by Pagerank	26.1	12.9
SemDedup [1]	27.1	13.8
Classifier on BGE features [176]	27.2	14.0
AskLLM [139]	28.6	14.3
Perplexity filtering	29.0	15.0
Top-k average logits	29.2	14.7
→ fastText [81] OH-2.5 +ELI5	<b>30.2</b>	<b>15.4</b>

Results for 1B-1x scale (1.4B model, ~28B tokens)

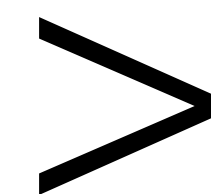
→ A fastText bigram classifier **combined with the right training data** does best

# Training data is key for the data filtering model

Dataset	Threshold	CORE	MMLU
→ OH-2.5 + ELI5	10%	<b>41.0</b>	<b>29.2</b>
Wikipedia	10%	35.7	27.0
OpenWebText2	10%	34.7	25.0
GPT-3 Approx	10%	37.5	24.4
OH-2.5 + ELI5	15%	39.8	27.2
OH-2.5 + ELI5	20%	38.7	24.2

Results for 7B-1x scale (7B model, ~140B tokens)

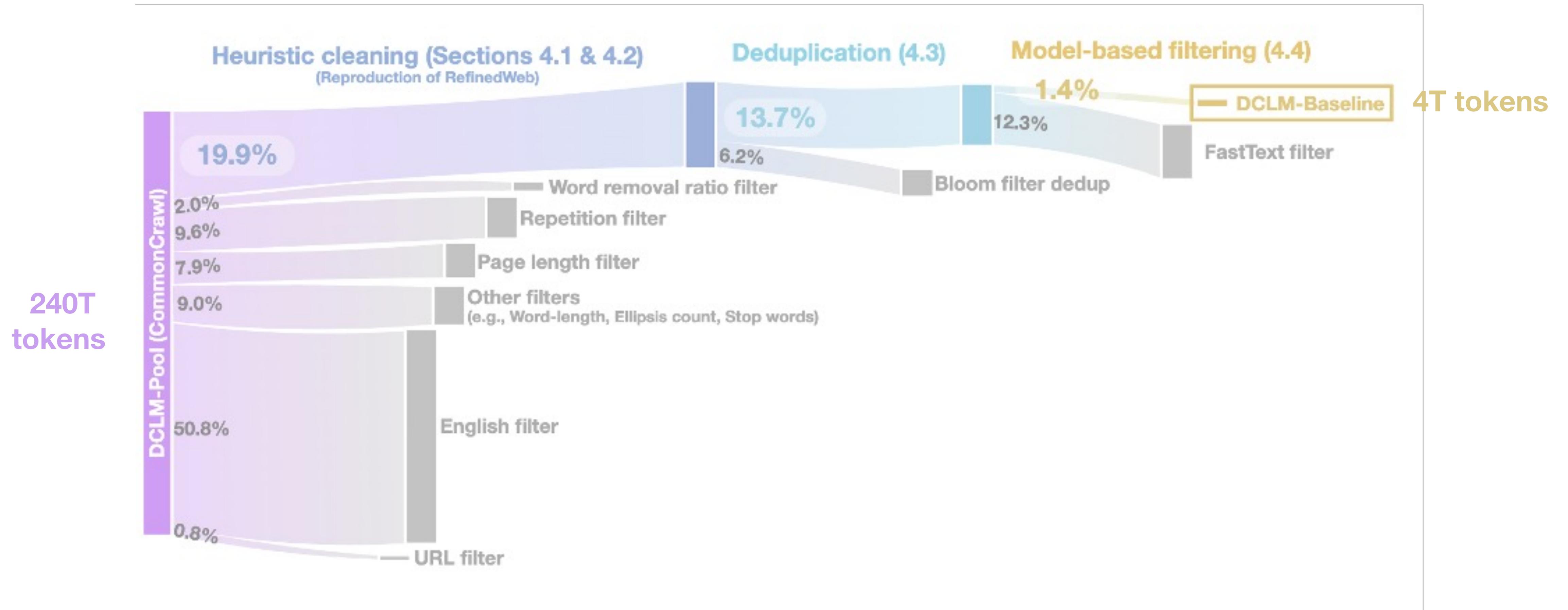
- **OpenHermes-2.5:** a SoTA instruction tuning dataset (mix of ~15 sources)
- **ELI5:** pairs of questions and top answers from the /r/explainlikeimfive subreddit

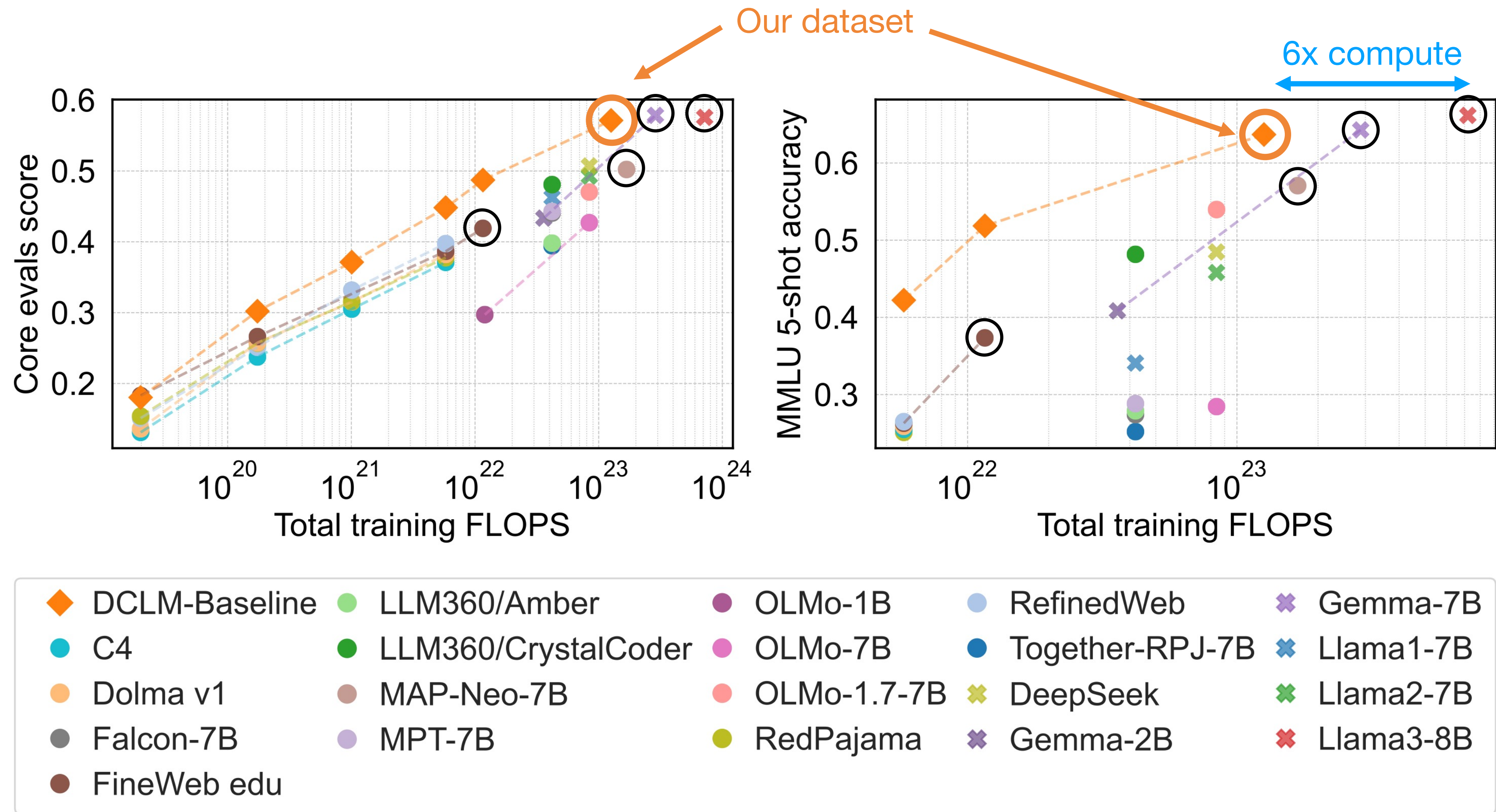


**Previous sources (GPT-3 Approx):**  
Wiki + Books + OpenWebText2



# Putting it all together: DCLM-Baseline





DCLM-Baseline outperforms leading open-source datasets like **FineWeb-Edu** and **MAP-Neo** and closes gap to **Mistral-7B**, **Gemma-7B**, and **Llama 3-8B** (with 6x less compute).

# Check out the DCLM paper for more findings!

## DataComp-LM: In search of the next generation of training sets for language models

Jeffrey Li\*<sup>1,2</sup> Alex Fang\*<sup>1,2</sup> Georgios Smyrnis\*<sup>4</sup> Maor Ivgi\*<sup>5</sup>  
Matt Jordan<sup>4</sup> Samir Gadre<sup>3,6</sup> Hritik Bansal<sup>8</sup> Etash Guha<sup>1,15</sup> Sedrick Keh<sup>3</sup> Kushal Arora<sup>3</sup>  
Saurabh Garg<sup>13</sup> Rui Xin<sup>1</sup> Niklas Muennighoff<sup>22</sup> Reinhard Heckel<sup>12</sup> Jean Mercat<sup>3</sup> Mayee  
Chen<sup>7</sup> Suchin Gururangan<sup>1</sup> Mitchell Wortsman<sup>1</sup> Alon Albalak<sup>19,20</sup> Yonatan Bitton<sup>14</sup>  
Marianna Nezhurina<sup>9,10</sup> Amro Abbas<sup>23</sup> Cheng-Yu Hsieh<sup>1</sup> Dhruva Ghosh<sup>1</sup> Josh Gardner<sup>1</sup>  
Maciej Kilian<sup>17</sup> Hanlin Zhang<sup>18</sup> Rulin Shao<sup>1</sup> Sarah Pratt<sup>1</sup> Sunny Sanyal<sup>4</sup> Gabriel Ilharco<sup>1</sup>  
Giannis Daras<sup>4</sup> Kalyani Marathe<sup>1</sup> Aaron Gokaslan<sup>16</sup> Jieyu Zhang<sup>1</sup> Khyathi Chandu<sup>11</sup>  
Thao Nguyen<sup>1</sup> Igor Vasiljevic<sup>3</sup> Sham Kakade<sup>18</sup> Shuran Song<sup>6,7</sup> Sujay Sanghavi<sup>4</sup> Fartash  
Faghri<sup>2</sup> Sewoong Oh<sup>1</sup> Luke Zettlemoyer<sup>1</sup> Kyle Lo<sup>11</sup> Alaaeldin El-Nouby<sup>2</sup> Hadi Pouransari<sup>2</sup>  
Alexander Toshev<sup>2</sup> Stephanie Wang<sup>1</sup> Dirk Groeneveld<sup>11</sup> Luca Soldaini<sup>11</sup> Pang Wei Koh<sup>1</sup>  
Jenia Jitsev<sup>9,10</sup> Thomas Kollar<sup>3</sup> Alexandros G. Dimakis<sup>4,21</sup>  
Yair Carmon<sup>5</sup> Achal Dave<sup>13</sup> Ludwig Schmidt<sup>11,7</sup> Vaishal Shankar<sup>12</sup>

<sup>1</sup>University of Washington, <sup>2</sup>Apple, <sup>3</sup>Toyota Research Institute, <sup>4</sup>UT Austin,  
<sup>5</sup>Tel Aviv University, <sup>6</sup>Columbia University, <sup>7</sup>Stanford, <sup>8</sup>UCLA, <sup>9</sup>JSC, <sup>10</sup>LAION, <sup>11</sup>AI2,  
<sup>12</sup>TUM, <sup>13</sup>CMU, <sup>14</sup>Hebrew University, <sup>15</sup>SambaNova, <sup>16</sup>Cornell, <sup>17</sup>USC, <sup>18</sup>Harvard,  
<sup>19</sup>UCSB, <sup>20</sup>SynthLabs, <sup>21</sup>Bespokelabs.AI, <sup>22</sup>Contextual AI, <sup>23</sup>DatologyAI

contact@datacomp.ai


### Abstract

We introduce DataComp for Language Models (DCLM), a testbed for controlled dataset experiments with the goal of improving language models. As part of DCLM, we provide a standardized corpus of 240T tokens extracted from Common Crawl, effective pretraining recipes based on the OpenLM framework, and a broad suite of 53 downstream evaluations. Participants in the DCLM benchmark can experiment with data curation strategies such as deduplication, filtering, and data mixing at model scales ranging from 412M to 7B parameters. As a baseline for DCLM, we conduct extensive experiments and find that model-based filtering is key to assembling a high-quality training set. The resulting dataset, DCLM-BASELINE, enables training a 7B parameter language model from scratch to 64% 5-shot accuracy on MMLU with 2.6T training tokens. Compared to MAP-Neo, the previous state-of-the-art in open-data language models, DCLM-BASELINE represents a 6.6 percentage point improvement on MMLU while being trained with 40% less compute. Our baseline model is also comparable to Mistral-7B-v0.3 and Llama 3 8B on MMLU (63% & 66%), and performs similarly on an average of 53 natural language understanding tasks while being trained with 6.6× less compute than Llama 3 8B. Our results highlight the importance of dataset design for training language models and offer a starting point for further research on data curation. We release the DCLM benchmark, framework, models, and datasets at <https://datacomp.ai/dclm>.

- Resiliparse and Bloom Filters are effective for extraction and deduplication
- Dataset rankings are consistent across scales and hyperparameters
- Adding traditional “high-quality” sources (e.g., Wikipedia) doesn’t always help
- Contamination is unlikely to explain the strong performance of DCLM-Baseline
- Agreement with human labels isn’t a good predictor among high-quality filters

# Let's all improve open-source datasets!

DataComp [Home](#) [FAQs](#) [Team](#) [Leaderboard](#)




## DataComp - LM

**Welcome to DataComp**, the machine learning benchmark where the models are fixed and the challenge is to find the best possible data!

Prior competitions in machine learning have focused on finding the best model, with a fixed set of training and test data. However, many recent advances (GPT-4, Gemini, LLAMA, Mistral) are due in part to large and diverse language datasets. DataComp centers the role that data plays by fixing the training code, and encouraging researchers to innovate by proposing new training sets. We provide an experimental testbed centered around a new candidate pool of 240T tokens from Common Crawl. Participants in our benchmark design new filtering techniques or curate new data sources and then evaluate them by running our standardized language model training code followed by an evaluation on 53 downstream datasets. Our benchmark consists of multiple scales, with various candidate pool sizes and associated compute budgets ranging from 412M to 7B parameters. This multi-scale design facilitates the study of scaling trends and makes the benchmark accessible to researchers with varying resources.

[Paper](#) [Code](#) [Data](#)

Rank	Created	Submission	Core	Extended	MMLU	Authors	Writeup
1	06-19-2024	Baseline: DCLM-Baseline	0.44823	0.28776	0.42228	DCLM team	<a href="https://arxiv.org/abs/2406.19001">https://arxiv.org/abs/2406.19001</a>
2*	06-19-2024	External: RefinedWeb	0.39749	0.21665	0.24765	The Falcon LLM team	<a href="https://arxiv.org/abs/2406.19001">https://arxiv.org/abs/2406.19001</a>
3*	06-19-2024	External: Fineweb-Edu	0.38658	0.22062	0.26294	HuggingFace	<a href="https://huggingface.co/fineweb-edu">https://huggingface.co/fineweb-edu</a>
4*	06-19-2024	External: Dolma v1	0.38220	0.19573	0.25890	AI2	<a href="https://arxiv.org/abs/2406.19001">https://arxiv.org/abs/2406.19001</a>
5*	06-19-2024	External: RPJ	0.37752	0.20614	0.25111	Together AI	<a href="https://www.together.ai">https://www.together.ai</a>
6*	06-19-2024	External: C4	0.37103	0.19840	0.25505	Google	<a href="https://arxiv.org/abs/1909.09555">https://arxiv.org/abs/1909.09555</a>

 **dclm** Public Edit Pins Watch 38 Fork 104 Starred 1.2k



Dataset Name	Size	Number of Objects	Description
<a href="#">DCLM-pool</a>	279.6 TiB	5,047,684	Raw extracted text
<a href="#">DCLM-refinedweb</a>	41.6 TiB	357,163	Raw text subsequently processed with a pipeline similar to <a href="#">RefinedWeb</a>
<a href="#">DCLM-baseline</a>	6.6 TiB	27,938	RefinedWeb filtered using ML

<https://www.datacomp.ai/dclm/>



# How this project incorporates foundational work

The idea of data centric AI is a new paradigm.

We need to design new AI models for filtering. We need theoretical principles for data quality filtering, data subset selection etc.

Evals are hard. Just released Evalchemy, an LLM evaluation platform today.

A lot of on-going interest on training with synthetic data. Several new foundational questions arise and Datacomp can be a benchmark for synthetic datasets.

On-going: Datacomp for post-training: Create a summarization dataset by prompting GPT. Post-train a small (8B Llama) models on synthetic instruction datasets.

Theoretical principles for foundation model specialization, using GPT4o as an Oracle.

# Small Specialized Models through Synthetic Data curation

**Paradigm 1:** The GPT 5 AGI monolith.

A gigantic model that knows everything.

You prompt it to summarize papers.

Students prompt it to ask questions about your lecture notes etc.

# Small Specialized Models through Synthetic Data curation

**Paradigm 1:** The GPT 5 AGI monolith.  
A gigantic model that knows everything.  
You prompt it to summarize papers.  
Students prompt it to ask questions about your lecture notes etc.

**Paradigm 2:** Small specialized models  
Can be created by synthetic dataset curation:  
1. You prompt the big model to create an instruction dataset:  
1. Summarize these 10k Neurips papers.  
2. Check if this synthetic dataset is truthful, faithful, helpful etc.  
3. Distill a small LLM to make a Neurips paper summarizer  
4. Evaluate that this is working well

The small model can be better than its teacher (e.g. as shown in Bespoke-Minicheck 7B)  
(if the prompting pipeline uses external knowledge, external validators or competition).

- Fin