# DataComp: In search of the next generation of multimodal datasets

Samir Yitzhak Gadre

COLUMBIA UNIVERSITY
IN THE CITY OF NEW YORK

# DATACOMP:
# In search of the next generation of multimodal datasets

Samir Yitzhak Gadre[*,2]  Gabriel Ilharco[*,1]  Alex Fang[*,1]  Jonathan Hayase[1]  Georgios Smyrnis[5]
Thao Nguyen[1]  Ryan Marten[7,9]  Mitchell Wortsman[1]  Dhruba Ghosh[1]  Jieyu Zhang[1]
Eyal Orgad[3]  Rahim Entezari[10]  Giannis Daras[5]  Sarah Pratt[1]  Vivek Ramanujan[1]
Yonatan Bitton[11]  Kalyani Marathe[1]  Stephen Mussmann[1]  Richard Vencu[6]
Mehdi Cherti[6,8]  Ranjay Krishna[1]  Pang Wei Koh[1,12]  Olga Saukh[10]  Alexander Ratner[1,13]
Shuran Song[2]  Hannaneh Hajishirzi[1,7]  Ali Farhadi[1]  Romain Beaumont[6]
Sewoong Oh[1]  Alexandros G. Dimakis[5]  Jenia Jitsev[6,8]
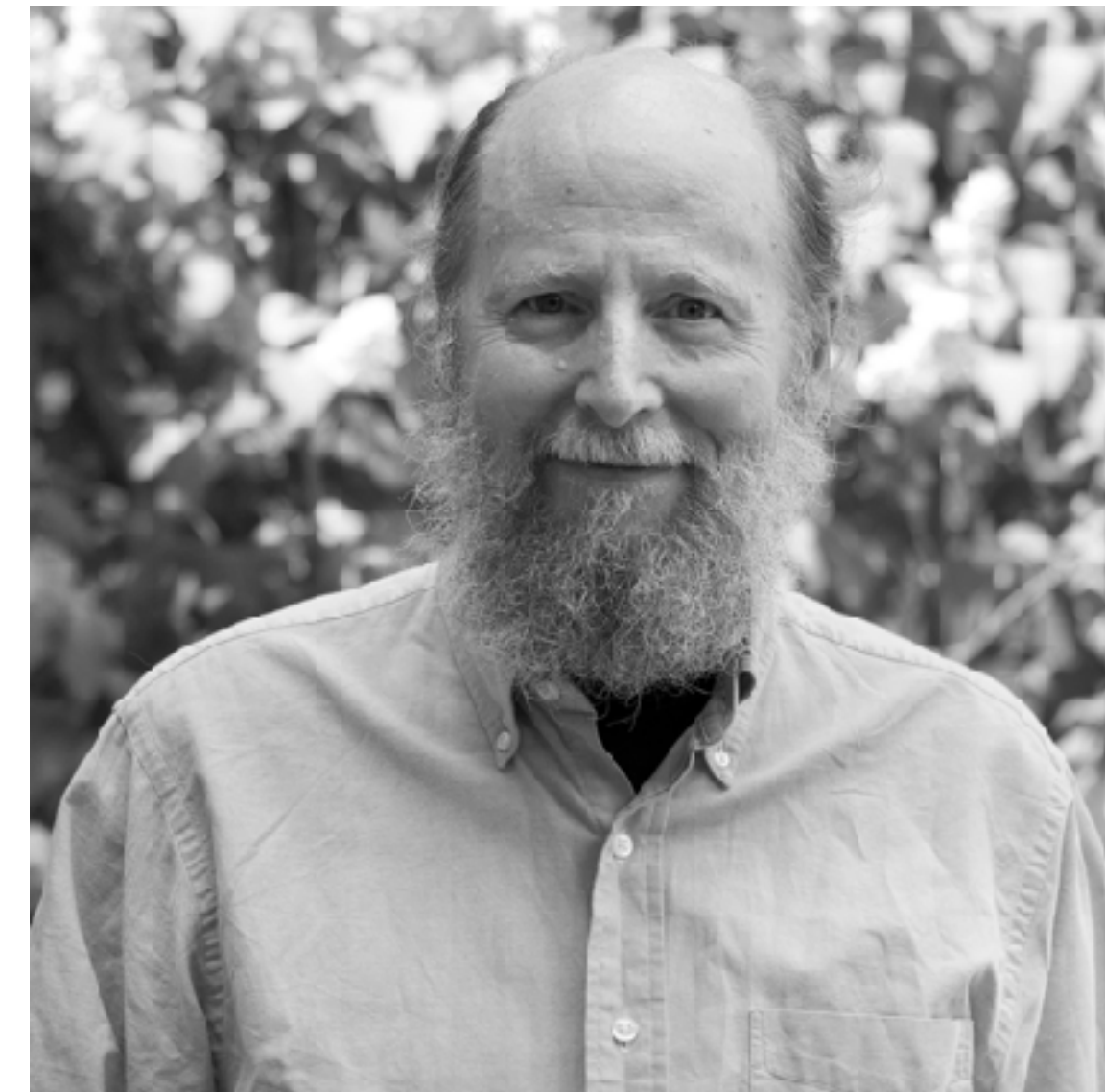Yair Carmon[3]  Vaishaal Shankar[4]  Ludwig Schmidt[1,6,7]

## Abstract

Multimodal datasets are a critical component in recent breakthroughs such as Stable Diffusion and GPT-4, yet their design does not receive the same research attention as model architectures or training algorithms. To address this shortcoming in the ML ecosystem, we introduce DATACOMP, a testbed for dataset experiments centered around a new candidate pool of 12.8 billion image-text pairs from Common Crawl. Participants in our benchmark design new filtering techniques or curate new data sources and then evaluate their new dataset by running our standardized CLIP training code and testing the resulting model on 38 downstream test sets. Our benchmark consists of multiple compute scales spanning four orders of magnitude, which enables the study of scaling trends and makes the benchmark accessible to researchers with varying resources. Our baseline experiments show that the DATACOMP workflow leads to better training sets. In particular, our best baseline, DATACOMP-1B, enables training a CLIP ViT-L/14 from scratch to 79.2% zero-shot accuracy on ImageNet, outperforming OpenAI's CLIP ViT-L/14 by 3.7 percentage points while using the same training procedure and compute. We release DATACOMP and all accompanying code at www.datacomp.ai.

*Gadre et al. DataComp: In search of the next generation of multimodal datasets. NeurIPS D&B 2023.*
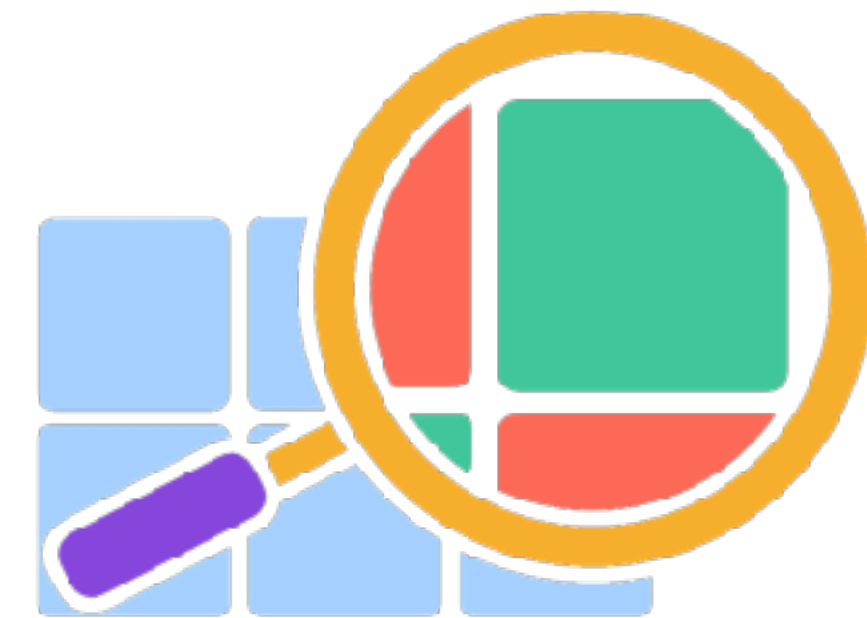
# The bitter lesson

"The biggest lesson that can be read from 70 years of AI research is that general methods that leverage computation are ultimately the most effective, and by a large margin."
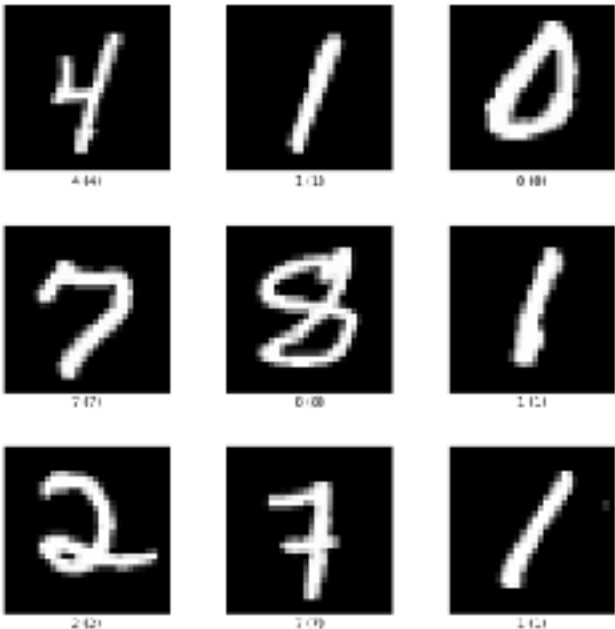
*Rich Sutton. The bitter lesson. 2019.*

# The data lesson addendum?

"The biggest lesson that can be read from … AI research is that general methods that leverage computation are ultimately the most effective," **especially when applied in conjunction with rigorous dataset construction.**
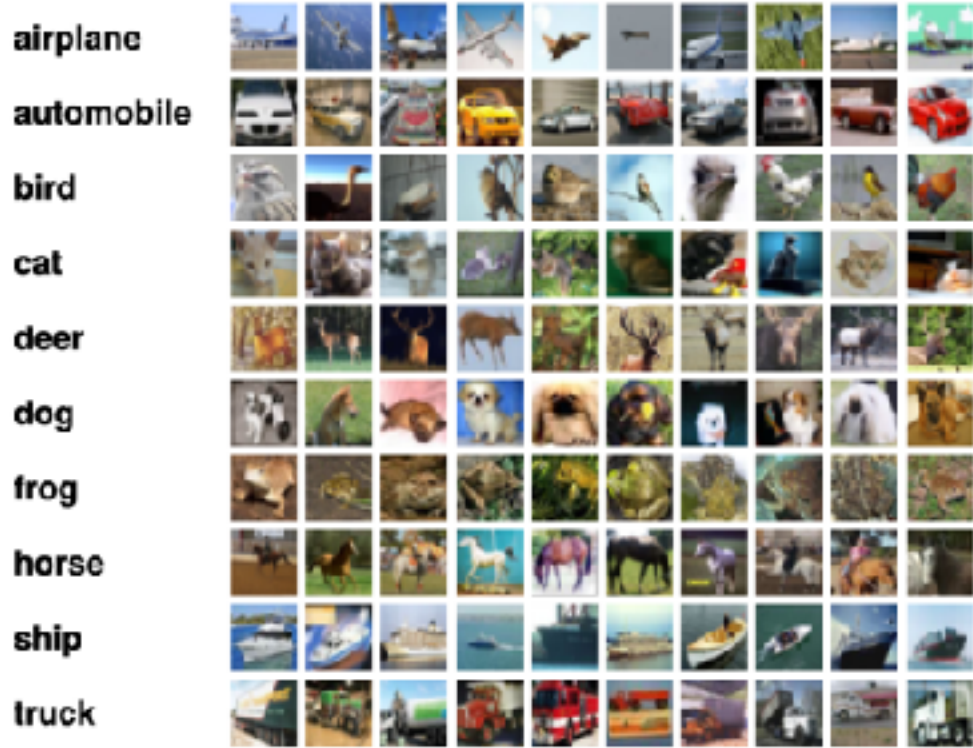
# Datasets are the foundation of progress in ML



**ImageNet (2012)**

Deep learning resurgence, ResNets, transfer learning, etc.

**MNIST (1994)**

Convolutional neural networks

**CIFAR-10 (2009)**
Training on GPUs

**WebImageText (2021)**

Zero-shot classification (CLIP), text-guided image generation (DALL-E)

# The standard ML research pipeline



**A. Select datasets** 🔒   **B. Train** 🔄   **C. Evaluate** 🔒

# This pipeline has produced better models

- Architectures

- Optimizers

- Normalization

- Tuned hyperparameters

- Activation functions

- Weight initialization schemes

- Stable training tricks

But how much performance are we leaving on the table by fixing datasets?

# Dataset iteration for better models?

- Diversity?

- Hard vs. easy examples?

- Class distributions?

- Label quality?

- Scale?



**Clearly this is a platypus says ImageNet**

# DataComp is a benchmark for dataset development

# Enter DataComp

**A. Select datasets** ↻

**B. Train** 🔒

**C. Evaluate** 🔒

# Interlude: DataComp targets CLIP training



This doggie is adorable!

Lol that's a weird frog

Embedding Space

*Radford et al. Learning Transferable Vision Models From Natural Language Supervision. ICML 2021.*

# Interlude: CLIP for zero-shot inference

- With vision-language features we can create *arbitrary* image classifiers.

**Input image**       **Prompts to create classifier**       **Similarity scores give class label**



**A photo of a dog.**

**A photo of a frog.**

*Radford et al. Learning Transferable Vision Models From Natural Language Supervision. ICML 2021.*

# Interlude: Why CLIP?



| Dataset | | ImageNet ResNet101 | CLIP ViT-L |
|---|---|---|---|
| ImageNet | | 76.2% | 76.2% |
| ImageNet V2 | | 64.3% | 70.1% |
| ImageNet Rendition | | 37.7% | 88.9% |
| ObjectNet | | 32.6% | 72.3% |
| ImageNet Sketch | | 25.2% | 60.2% |
| ImageNet Adversarial | | 2.7% | 77.1% |

*https://openai.com/research/clip*

# Interlude: Why CLIP?

- Many model vision models utilize CLIP backbones for V&L tasks, segmentation, detection, image generation, embodied tasks, etc.

- Reasonable signal that improving CLIP models also leads to downstream model gains



Figure 3: **Flamingo architecture overview.** Flamingo is a family of visual language models (VLMs) that take as input visual data interleaved with text and produce free-form text as output.

*Alayrac et al. Flamingo: A visual language model for few-shot learning. NeurIPS 2022.*

# Re-enter DataComp



A. Choose Scale  B. Select Data  C. Train  D. Evaluate  E. Submit

# Choosing a scale



A. Choose Scale  B. Select Data  C. Train  D. Evaluate  E. Submit

# DataComp is compute accessible

- Academics may have less resources and usually can't train many FLOPs

- Industry labs may not want to participate unless DataComp can produce SOTA models

- Solution: different compute scales for participants

**samir gadre**
@sy_gadre

academia          industry
🤝
gossip

2:51 AM · Jun 25, 2023

A. Choose Scale

# Scale configurations

|  | small | medium | large | xlarge |
|---|---|---|---|---|
| samples seen | | | | |
| model | | | | |
| training A100 hours | | | | |
| compute analogy | | | | |

# Scale configurations

|  | small | medium | large | xlarge |
|---|---|---|---|---|
| samples seen | 12.8M |  |  |  |
| model | ViT-B/32 |  |  |  |
| training A100 hours | 8 |  |  |  |
| compute analogy | fine-tune IN-1k |  |  |  |

A. Choose Scale

# Scale configurations

| | small | medium | large | xlarge |
|---|---|---|---|---|
| samples seen | 12.8M | 128M | | |
| model | ViT-B/32 | ViT-B/32 | | |
| training A100 hours | 8 | 80 | | |
| compute analogy | fine-tune IN-1k | training IN-1k | | |

A. Choose Scale

# Scale configurations

|  | small | medium | large | xlarge |
|---|---|---|---|---|
| samples seen | 12.8M | 128M | 1.28B |  |
| model | ViT-B/32 | ViT-B/32 | ViT-B/16 |  |
| training A100 hours | 8 | 80 | 1,000 |  |
| compute analogy | fine-tune IN-1k | training IN-1k | training IN-21k |  |

A. Choose Scale

# Scale configurations

|  | small | medium | large | xlarge |
|---|---|---|---|---|
| samples seen | 12.8M | 128M | 1.28B | 12.8B |
| model | ViT-B/32 | ViT-B/32 | ViT-B/16 | ViT-L/14 |
| training A100 hours | 8 | 80 | 1,000 | 40,000 |
| compute analogy | fine-tune IN-1k | training IN-1k | training IN-21k | training OAI CLIP |

A. Choose Scale

# Scale configurations

|  | small | medium | large | xlarge |
|---|---|---|---|---|
| samples seen | 12.8M | 128M | 1.28B | 12.8B |
| model | ViT-B/32 | ViT-B/32 | ViT-B/16 | ViT-L/14 |
| training A100 hours | 8 | 80 | 1,000 | 40,000 |
| compute analogy | fine-tune IN-1k | training IN-1k | training IN-21k | training OAI CLIP |

**No constraint on dataset size! Real constraints are pool size and compute.**

A. Choose Scale

# Example of samples seen

- Participating at the medium scale

  (128M samples seen)

- Filter a dataset to 64M samples

- Each sample will then get seen twice

  (in expectation) during training

**Initial pool**

**After filtering**

**Medium scale training**

A. Choose Scale
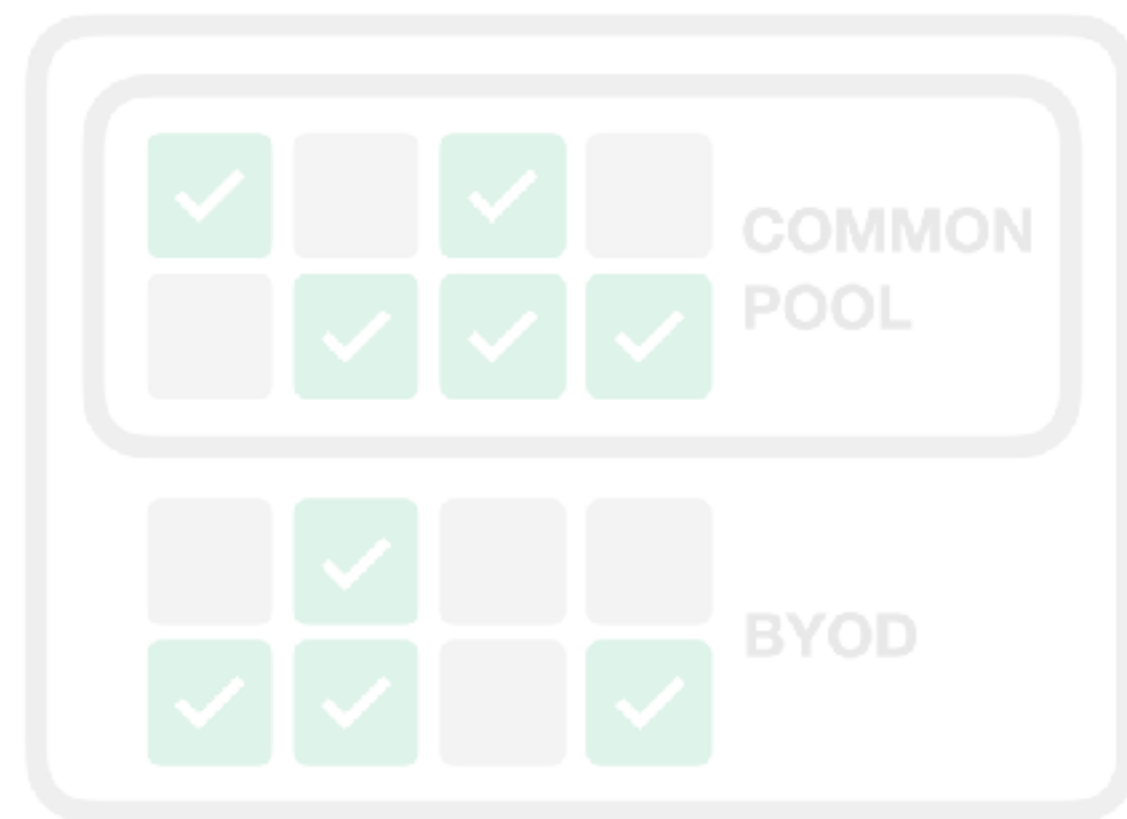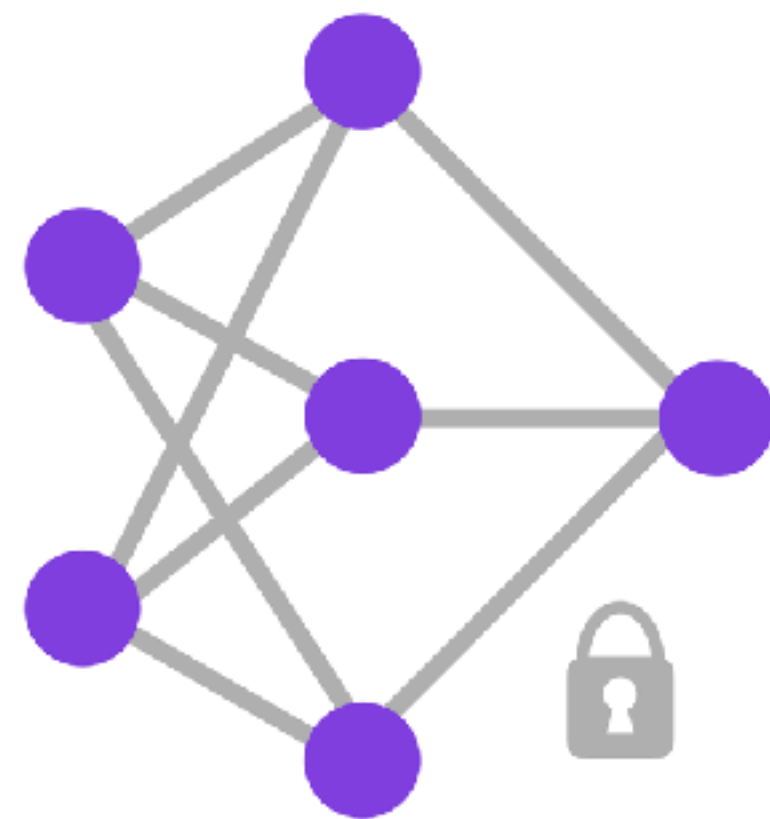
# Selecting data
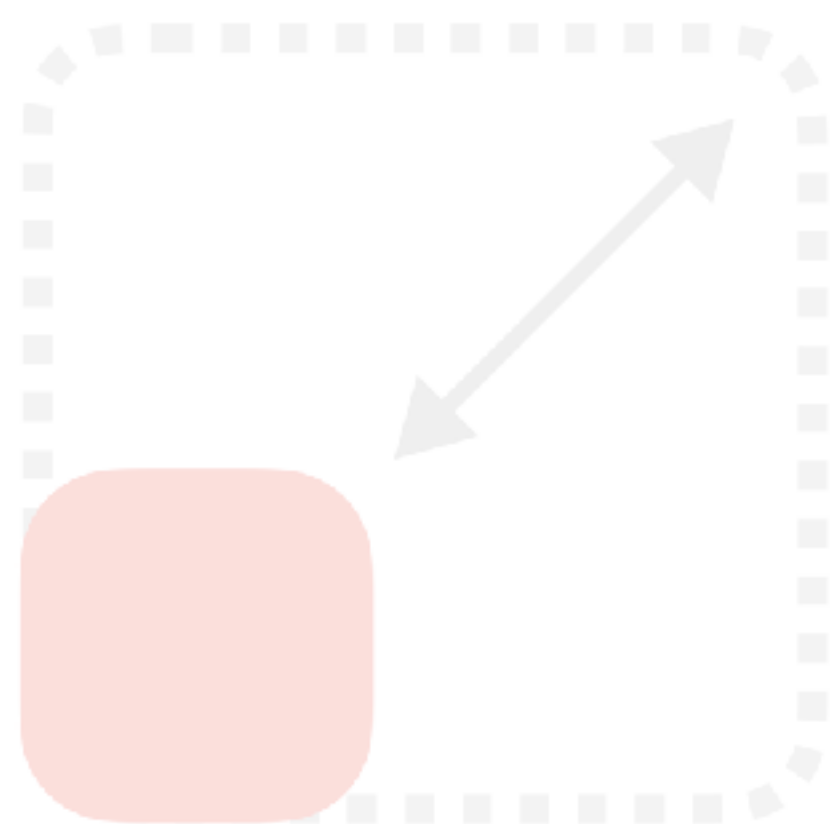


A. Choose Scale
B. Select Data
C. Train
D. Evaluate
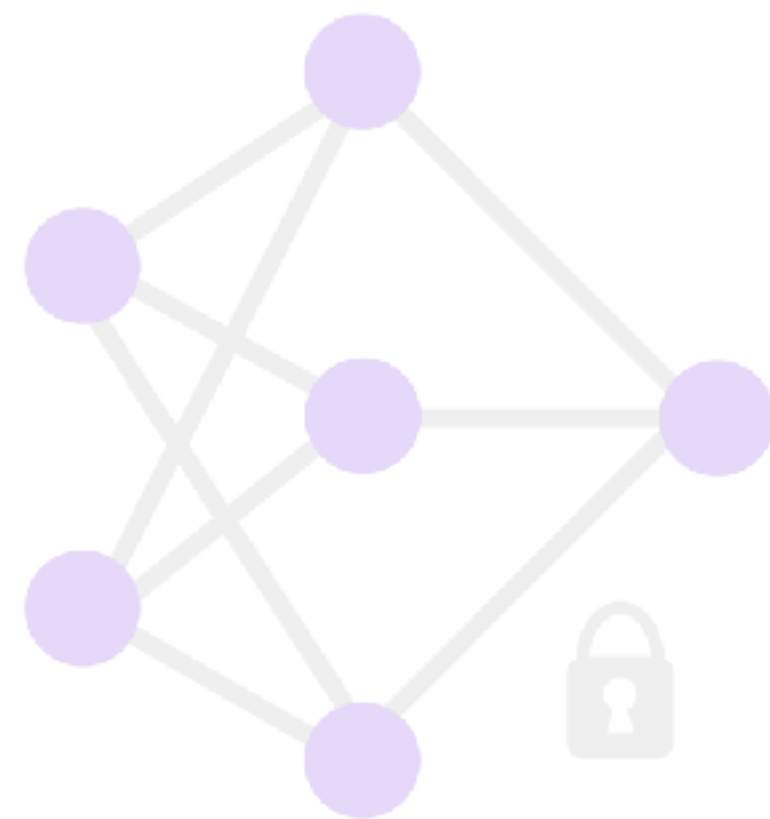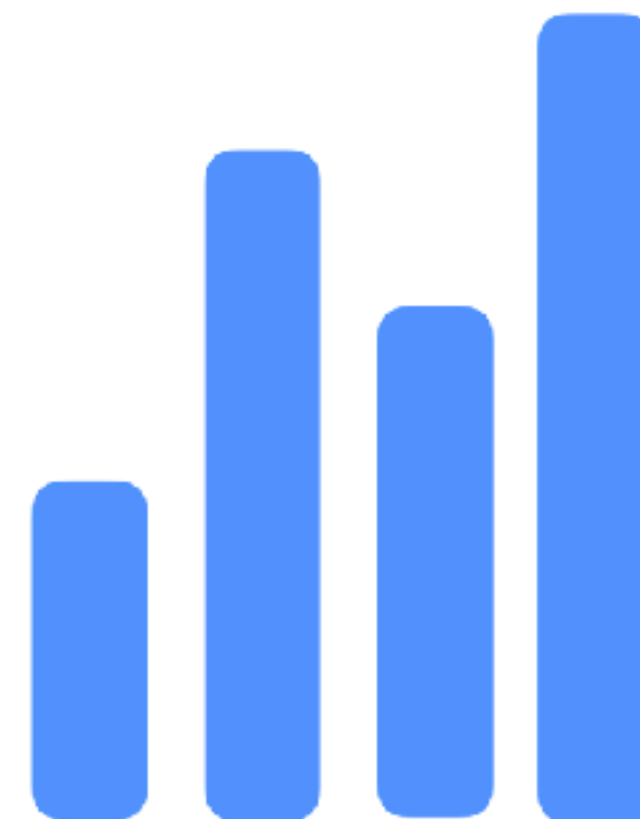E. Submit

# Two tracks: Filtering and BYOD

**Filtering**

**Bring your own data (BYOD)**

Pare down a noisy source

Cobble together many sources

COMMON POOL

BYOD

B. Select Data

# CommonPool to facilitate the filtering track

- 88B url-(alt)text pairs from CommonCrawl

- 40B attempted image downloads

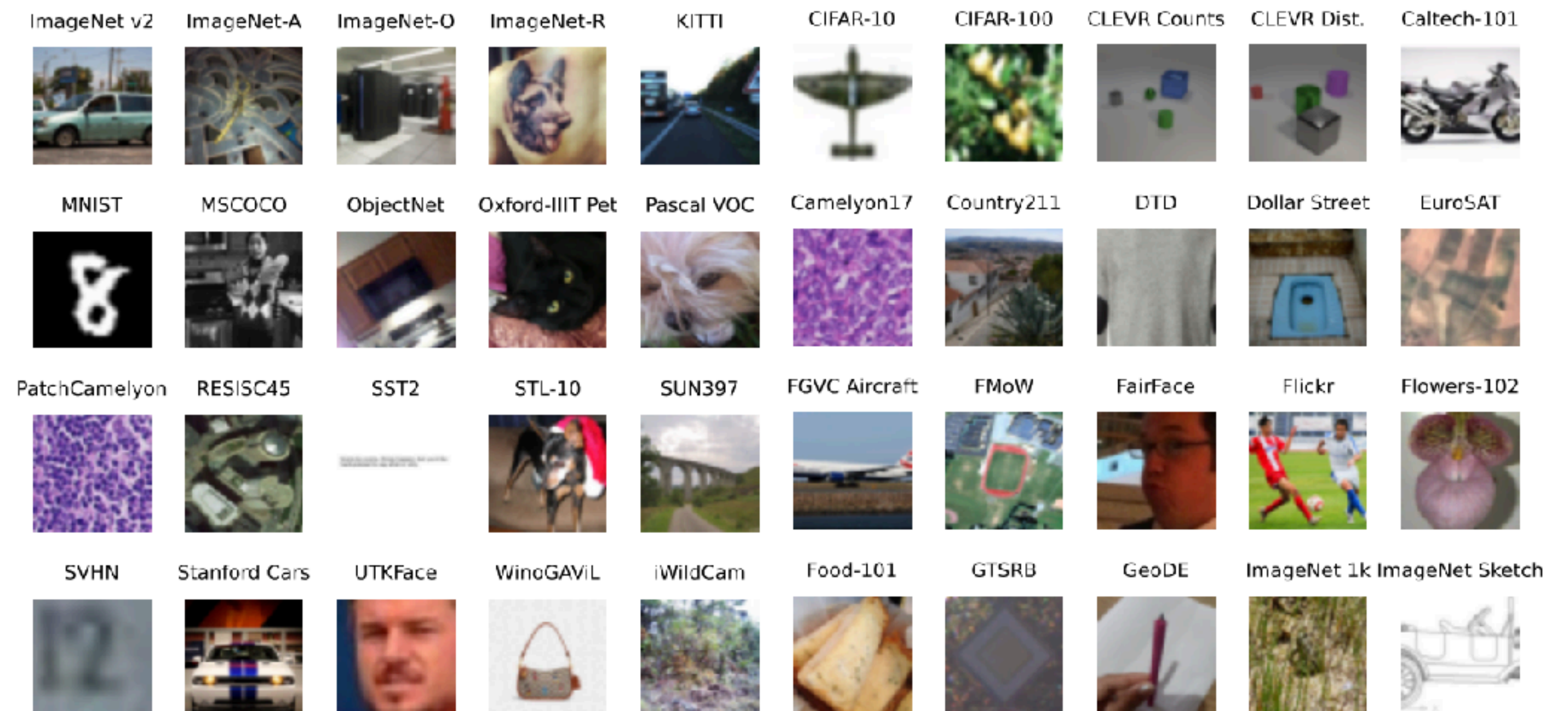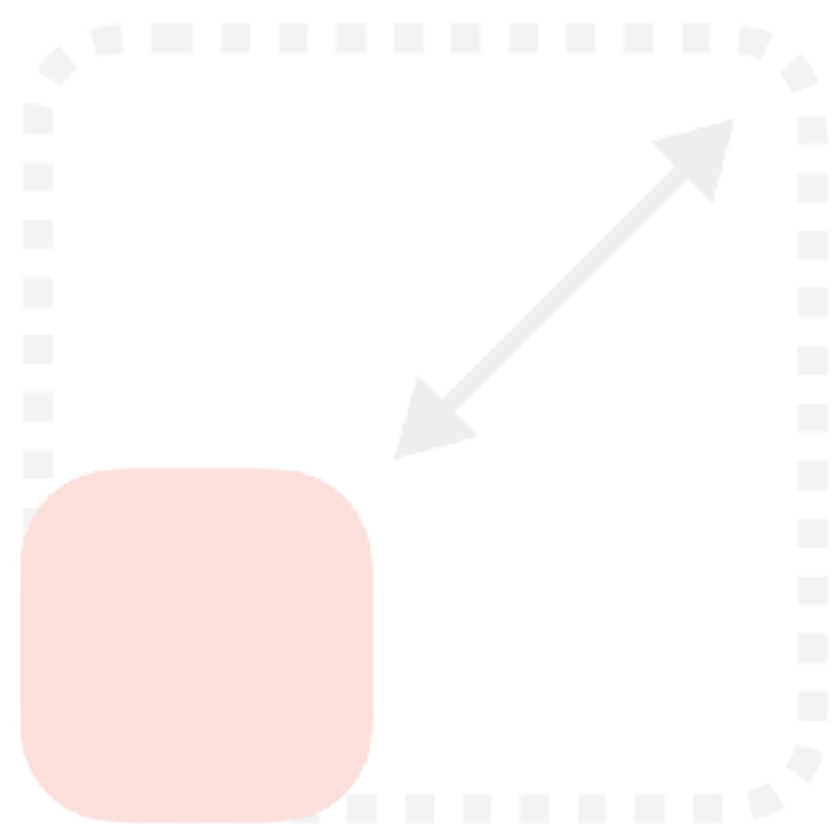- 16.8B successfully downloaded

- 13.1B retained after pre-processing

- 12.8B sampled for the xlarge pool

**40B potential candidates**

CommonPool data funnel

**12.8B xlarge candidate pool**

COMMON POOL

BYOD

**B. Select Data**

# Pre-filtering for safety and eval decontamination

- Near deduplication against downstream evaluation images

- NSFW image removal

- NSFW text removal

- Face blurring automatically in download tooling

- Notably, not pre-processing for "quality"

- Dataset safety is an active area of research!



B. Select Data

# Metadata

- Original width/height

- Caption

- Image sha256

- CLIP features (B/32 and L/14)

- CLIP scores

- Face bounding boxes (for automatic blurring)

# Bring your own data (BYOD)

- Filtering is only one way to curate datasets

- Combine other data sources (e.g., YFCC-15M, CC12M, RedCaps, etc.)

- CommonPool filtering ++

- The BYOD track allows this flexibility



COMMON POOL

BYOD

**B. Select Data**

# Train



A. Choose Scale     B. Select Data     C. Train     D. Evaluate     E. Submit

# Fixed training configurations

- Hyperparameters based on OpenAI, LAION (open_clip) runs

- Fixed, so participants cannot modify

- Ablations on architecture, batch size, etc. show relatively consistent trends, suggesting dataset and modeling choices can be considered independently

```
"medium": {
    "batch_size": 4096,
    "learning_rate": 5e-4,
    "train_num_samples": 128_000_000,
    "warmup": 500,
    "model": "ViT-B-32",
    "beta2": None,
},
```

C. Train

# Evaluate



A. Choose Scale   B. Select Data   C. Train   D. Evaluate   E. Submit

# Downstream eval sets

- 38 core classification and retrieval tasks

- Evaluations are zero-shot (no fine-tuning)

- We look at both ImageNet and average acc.



D. Evaluate

# Submit



A. Choose Scale    B. Select Data    C. Train    D. Evaluate    E. Submit

# datacomp.ai

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

E. Submit

# A unified leaderboard

**Select the track and scale**

| Filtering track | BYOD track |
|:---:|:---:|

| small | medium | large | xlarge |
|:---:|:---:|:---:|:---:|

**Leaderboard**

| Rank | Created | Submission | ImageNet acc. | Average perf. |
|---|---|---|---|---|
| 1 | 10-02-2023 | Data Filtering Networks | 0.371 | 0.373 |
| 2 | 09-08-2023 | The Devil Is in the Details | 0.320 | 0.371 |
| 3 | 08-17-2023 | T-MARS: Improving Visual Representations by Circumventing Text Feature Learning | 0.330 | 0.361 |

E. Submit

# DataComp leads to better models

# Teaser results

| Dataset | Dataset size | Samples seen | Architecture | ImageNet-1k accuracy |
|---------|--------------|--------------|--------------|----------------------|
| OpenAI WIT | 0.4B | 13B | ViT-L/14 | 75.5 |
| | | | | |
| | | | | |
| | | | | |
| | | | | |

# Teaser results

| Dataset | Dataset size | Samples seen | Architecture | ImageNet-1k accuracy |
|---|---|---|---|---|
| OpenAI WIT | 0.4B | 13B | ViT-L/14 | 75.5 |
| LAION-400M | 0.4B | 13B | ViT-L/14 | 72.8 |
| | | | | |
| | | | | |
| | | | | |
| | | | | |

# Teaser results

| Dataset | Dataset size | Samples seen | Architecture | ImageNet-1k accuracy |
|---|---|---|---|---|
| OpenAI WIT | 0.4B | 13B | ViT-L/14 | 75.5 |
| LAION-400M | 0.4B | 13B | ViT-L/14 | 72.8 |
| LAION-2B | 2.3B | 13B | ViT-L/14 | 73.1 |
| LAION-2B | 2.3B | 34B | ViT-H/14 | 78.0 |
| LAION-2B | 2.3B | 34B | ViT-g/14 | 78.5 |
| | | | | |

# Teaser results

| Dataset | Dataset size | Samples seen | Architecture | ImageNet-1k accuracy |
|---------|-------------|--------------|--------------|---------------------|
| OpenAI WIT | 0.4B | 13B | ViT-L/14 | 75.5 |
| LAION-400M | 0.4B | 13B | ViT-L/14 | 72.8 |
| LAION-2B | 2.3B | 13B | ViT-L/14 | 73.1 |
| LAION-2B | 2.3B | 34B | ViT-H/14 | 78.0 |
| LAION-2B | 2.3B | 34B | ViT-g/14 | 78.5 |
| DataComp-1B | 1.4B | 13B | ViT-L/14 | **79.2** |

# Teaser results

| Dataset | Dataset size | Samples seen | Architecture | ImageNet-1k accuracy |
|---------|--------------|--------------|--------------|----------------------|
| OpenAI WIT | 0.4B | 13B | ViT-L/14 | 75.5 |
| LAION-400M | 0.4B | 13B | ViT-L/14 | 72.8 |
| LAION-2B | 2.3B | 13B | ViT-L/14 | 73.1 |
| LAION-2B | 2.3B | 34B | ViT-H/14 | 78.0 |
| LAION-2B | 2.3B | 34B | ViT-g/14 | 78.5 |
| DataComp-1B | 1.4B | 13B | ViT-L/14 | **79.2** |

+3.7pp

# Teaser results

| Dataset | Dataset size | Samples seen | Architecture | ImageNet-1k accuracy |
|---------|--------------|--------------|--------------|----------------------|
| OpenAI WIT | 0.4B | 13B | ViT-L/14 | 75.5 |
| LAION-400M | 0.4B | 13B | ViT-L/14 | 72.8 |
| LAION-2B | 2.3B | 13B | ViT-L/14 | 73.1 |
| LAION-2B | 2.3B | 34B | ViT-H/14 | 78.0 |
| LAION-2B | 2.3B | 34B | ViT-g/14 | 78.5 |
| DataComp-1B | 1.4B | 13B | ViT-L/14 | **79.2** |

**9x compute savings**

# How did we get there? Baselines!

- No filtering

- CLIP-score filtering

- Basic: filtering based on aspect ratio, caption length, etc.

- Image-based filtering: clustering against ImageNet-1k train

- Text-based filtering: looking for ImageNet-1k synsets



IMG_2187.jpg

**No filtering**



Porsche Cayman S

**CLIP filtering (pool top 30%)**

# How did we get there? Baselines!

# How did we get there? Baselines!



CLIP ViT-L/14 similarity score
mean: 0.208, median: 0.203, min: -0.114, max: 0.524

# How did we get there? Baselines!



CLIP ViT-L/14 similarity score
mean: 0.208, median: 0.203, min: -0.114, max: 0.524

# How did we get there? Baselines!



CLIP ViT-L/14 30% filter

# How did we get there? Workflow!

- Ran many experiments at small and medium scale

- Best methods we run at large scale

- Best at large (often same as at medium), we run at xlarge

# DataComp-1B

- Combination of 2 baseline strategies (CLIP-filter ∩ image-based)

- First (public) training set that's better than OpenAI's.

# Dataset size is not the full story



ImageNet

Average over 38 datasets

- small scale
- medium scale
- large scale
- CLIP score (L/14)
- CLIP score (B/32)
- Rand. subset

# Dataset size is not the full story

# Consistent ordering across scales

# Consistent ordering across scales

# BYOD mixing can help

- BYOD seems to help up to large scale

- Currently diminishing returns at xlarge

- Lots left to explore here (rich, scalable sources beyond CommonCrawl?)

|  | large | xlarge |
|---|---|---|
| CLIP-filter ∩ image-based | 63.1 | 79.2 |
| + 4 BYOD sources | 65.6 | 79.2 |

# Future directions

- Curating more dataset sources

- Improved filtering and dataset balancing methods

- Further (weak) supervision signals

- Additional modalities

- Broader evaluation on vision, vision-language, robotics tasks, and model bias

# People are already iterating on DataComp

# On the horizon

- DataComp NLP (DCNLP)

- Larger pools (100B candidate image-text pairs)

- DataComp for image generation

- Multimodal DataComp?

- Interleaved DataComp?

# Everything is open source

- Central webpage: datacomp.ai

- Main repo: github.com/mlfoundations/datacomp

- CLIP training code: github.com/mlfoundations/open_clip

- Downloading billions of image-text pairs: github.com/rom1504/img2dataset

- Processing metadata for billions of image-text pairs: github.com/mlfoundations/dataset2metadata