

Data selection 1

IMPROVING PRETRAINING DATA USING PERPLEXITY CORRELATIONS

By Tristan Thrush, Christopher Potts, Tatsunori Hashimoto

Setup

Goal: select pre-training data in order to optimize performance on a fixed benchmark

Often: train small model on a bunch of mixtures of domains (data sets), check performance. Computationally expensive!

Here: only use “observational data” to find optimal data mix, no LLM training required.

Input: publicly available models + data sets + their performance in benchmarks

Intuition

If for a fixed data set...

If a LLM has low loss on validation data -> high score on target benchmark

and

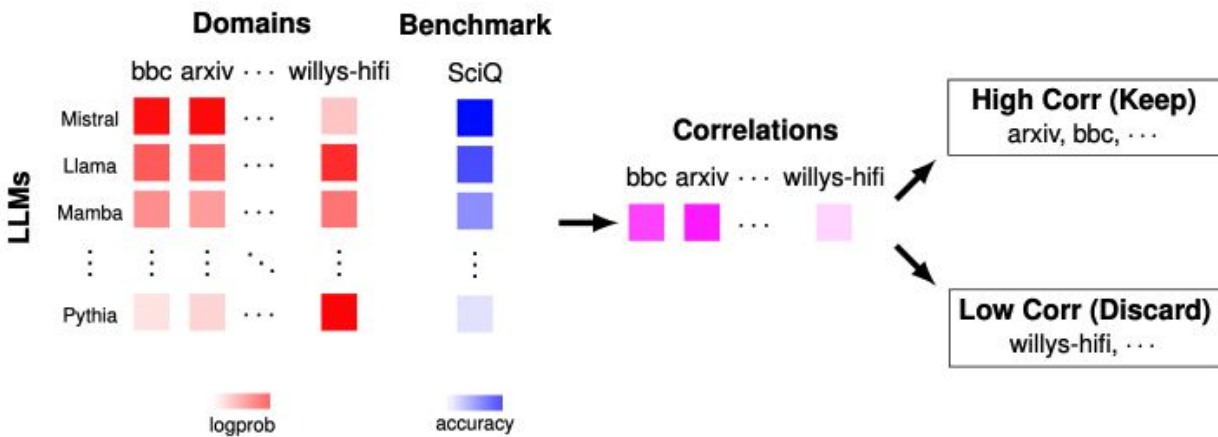
If a LLM has large loss on validation data -> low score on target benchmark

We would want to use this data set for pre-training (?)

Proposal

Run a regression of benchmark accuracies on domains (data set)

Each row is a LLM, each column is a domain (data set)



LLMs: 90 taken from Open LLM Leaderboard,

Domains: 9841 domains from sample-100B RedPajama V2

Seems weird, but at second glance it's not that unusual

Similar ideas have appeared in the causal inference literature.

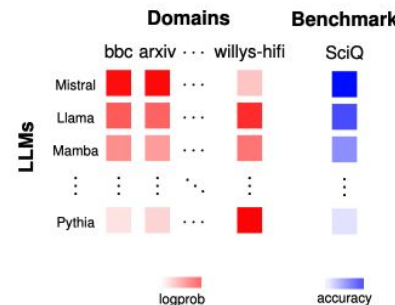
In the synthetic control approach, we regress a target data set (California) on several source data sets (Texas, Washington,...)

Columns: data sets (here: data sets)

Rows: pre-treatment variables (here: LLMs)

Entries: means of pre-treatment variables (here: log-likelihoods)

Target vector: mean of post-treatment outcome (here: accuracy)



Assumption

The perplexity-performance hypothesis. We formulate the task of predicting errors y_i from negative log-probabilities \mathbf{x}_i as a single-index model (SIM),

$$y_i = f(\langle \boldsymbol{\theta}^*, \mathbf{x}_i \rangle + \epsilon_i) \quad (1)$$

where $f : \mathbb{R} \mapsto \mathbb{R}$ is some unknown monotonically increasing univariate function, ϵ_i is zero-mean noise which is independent of \mathbf{x} , and $\boldsymbol{\theta}^* \in \mathbb{R}^D$ are unknown weights over D domains.

Outcomes y_i : accuracy of LLM i on benchmark data set

Covariates x_{ij} : negative log-likelihood of LLM i on document j

Parameter theta: best prediction

This is just a low-dimensional single index model that we can learn

Assumption: single index model holds not only for log-likelihoods of trained models in leaderboard, but also for the log-likelihoods when training a new model

-> To optimize benchmark performance (left hand side), we can just optimize weighted log-likelihoods (right hand side)

Estimating theta (idea from previous papers)

If X follow a standard Gaussian distribution, using integration by parts

$$E[X_j Y] = E[X_j f(\langle \theta^*, X \rangle + \epsilon)] = E[\theta_j f'(\langle \theta^*, X \rangle + \epsilon)] = \theta_j^* C$$

where $C > 0$. We are averaging over LLMs.

So, we can just compute the correlation between X and Y to estimate theta (up to a constant).

Estimation

The actual proposed algorithm avoids estimating f .

Distribution of loglikelihoods has heavy tails (solution: rank correlations, transformation)

Solutions θ may be negative and not sum up to one (solution: project θ on probability simplex - set of nonnegative vectors that sum up to one)

Algorithm 1 Perplexity Correlation Based Data Selection

Input: Benchmark error vector $\mathbf{y} \in [0, 1]^N$, log-loss matrix normalized as bits-per-byte $\mathbf{X} \in \mathbb{R}_0^{+N \times D}$, available tokens per domain $\mathbf{a} \in \mathbb{N}^D$, and pretraining token target $b \in \mathbb{N}$.

Output: Target token counts per domain $\mathbf{t} \in \mathbb{N}_0^D$, a fastText classifier to filter pretraining data.

Initialize: $\gamma \leftarrow \mathbf{0} \in \mathbb{R}^D$, $\mathbf{t} \leftarrow [0 \dots] \in \mathbb{N}_0^D$, counter $\leftarrow 0$.

$\mathbf{r}_0, \mathbf{r}_1, \dots, \mathbf{r}_N \leftarrow \text{rank}(\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_N)$ ▷ 1. Compute the γ correlation coefficient

for $i, j \in 0$ **to** N **do**

$\gamma \leftarrow \gamma + \text{sign}(y_i - y_j) \cdot (\mathbf{r}_i - \mathbf{r}_j)$

for $i \in \text{ArgSort}(\gamma, \text{descending}=\text{True})$ **do** ▷ 2. Select most to least correlated domains

$t_i \leftarrow \min(a_i, b - \text{counter})$

counter \leftarrow counter + a_i

if counter $\geq b$ **then**

Break

classifier = trainFastText(positive = $1_{t>0}$, negative = $1_{t=0}$)

Return \mathbf{t} , classifier

In practice...

- 1) No weights (domain included yes/no)
- 2) Train FastText classifier on selected domains to be able to select documents, not just domains

Evaluation

Fix pre-training architecture, change data sets.

Pre-training on Pythia 160M LLM configuration from Biderman et al. (2023)

Hyperparameters including learning rate, weight decay, and warmup are optimized to minimize loss on the uniform sampling (no selection algorithm) baseline.

Benchmarks

SciQ - scientific questions

ARC easy - grade school level questions

PIQA - commonsense reasoning, physics questions

Lambada - recover missing word from corpus of text (multilingual)

Accuracy
Higher is better

Log perplexity
Lower is better

LLM	Accuracy			Log perplexity							Method
	ARC Easy	PIQA	SciQ	LAMBADA	LAMBADA DE	LAMBADA ES	LAMBADA FR	LAMBADA IT	Perplexity Correlations		
Uniform RPJv2	42.21	58.92	63.20	10.77	23.95	29.04	22.90	26.95	None		
EN Filtered	43.48	57.94	64.90	9.78	72.52	69.89	59.09	67.10	equal weights		
DE Filtered	26.60	52.39	30.10	57.72	24.63	90.78	70.91	78.24	Lang Filt		
ES Filtered	28.79	51.52	32.80	55.86	86.75	31.48	67.79	69.91	language filter		
FR Filtered	29.08	52.56	34.60	53.67	80.39	80.22	28.30	70.79	Domain selection via importance sampling		
IT Filtered	27.95	51.52	29.10	70.29	70.24	80.90	74.83	37.09			
ARC Easy*	43.01	57.40	59.40	15.86	24.60	31.00	25.01	29.05	DSIR		
PIQA*	42.30	56.96	60.40	15.40	24.55	30.93	24.67	28.93			
SciQ*	42.17	57.94	59.60	15.62	24.51	32.17	24.56	28.95	Handcrafted fastText		
LAMBADA*	41.33	57.83	56.70	12.11	23.42	29.87	23.32	27.51			
LAMBADA DE*	40.28	57.13	57.60	17.67	22.32	30.40	23.40	27.60			
LAMBADA ES*	39.56	55.71	58.80	17.74	23.72	29.33	24.23	28.07			
LAMBADA FR*	40.15	56.20	57.30	18.20	23.94	28.64	23.41	26.70			
LAMBADA IT*	38.97	56.04	59.90	18.61	23.76	29.18	24.02	27.62			
fastText, EN Filt	45.50	60.12	63.90	9.06	71.80	70.94	58.27	66.16	Perplexity Correlations		
fastText	42.72	57.56	60.30	14.34	22.04	27.79	21.74	24.86			
ARC Easy'	44.15	58.49	61.20	11.64	83.06	78.03	64.99	70.96			
PIQA'	44.02	58.98	59.00	11.69	84.22	78.76	65.10	69.84			
SciQ'	46.63	58.65	65.70	12.36	75.96	73.09	61.24	67.77			
LAMBADA'	42.13	58.92	61.40	9.42	79.48	76.47	62.71	69.34			
LAMBADA DE'	42.51	56.80	62.00	12.94	20.88	26.97	21.29	24.85			
LAMBADA ES'	42.97	55.77	62.20	15.62	21.85	28.59	22.96	26.31			
LAMBADA FR'	43.18	56.75	63.80	12.92	55.07	26.98	22.15	25.64			
LAMBADA IT'	41.54	56.75	57.40	13.49	21.10	27.28	21.33	25.24			

Data selection method

you get performance

TLDR: fastText is great, perplexity pretty good



Aside: what is fastText (SOTA)?

<https://huggingface.co/mlfoundations/fasttext-oh-eli5>

Text classifier

Filter data by sorting pages by “high quality” score, until we had reached 3.2B unique tokens.

Comments

- Should we do data selection at the domain level or document level (or both?)
- Can we improve performance by using weights (instead of include data yes/no?). The greedy strategy should be suboptimal
- We should have a principled approach for choosing between data collection strategies based on what type of distribution shift we have
- We probably want to find weights that do well universally across tasks
- Rank correlations are not very satisfying since they hide whether the changes are large

Awesome!

(...) we use this paper also as a preregistration for further pretraining experiments using different data sources and evaluation benchmarks at the 1.4B model scale.

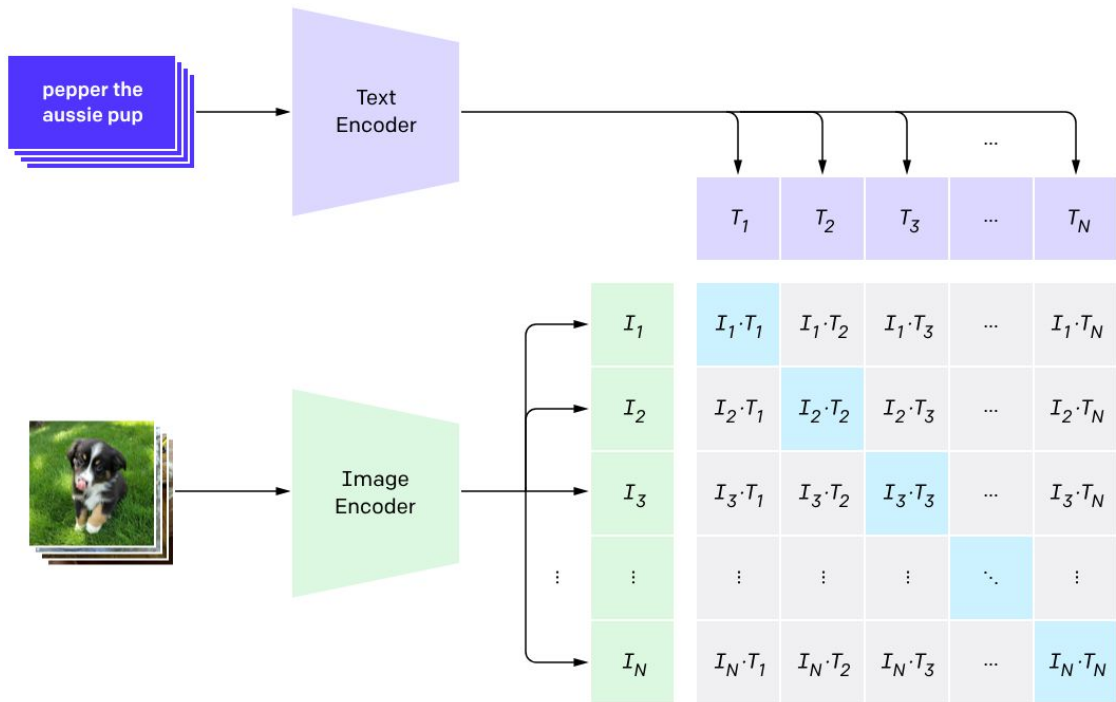
We commit to updating the arXiv version of our manuscript with both positive and negative results on these preregistered validation experiments.

SCALING LAWS FOR DATA FILTERING – DATA CURATION CAN NOT BE COMPUTE AGNOSTIC

Sachin Goyal, Pratyush Maini
Zachary Lipton, Aditi Raghunathan, Zico Kolter

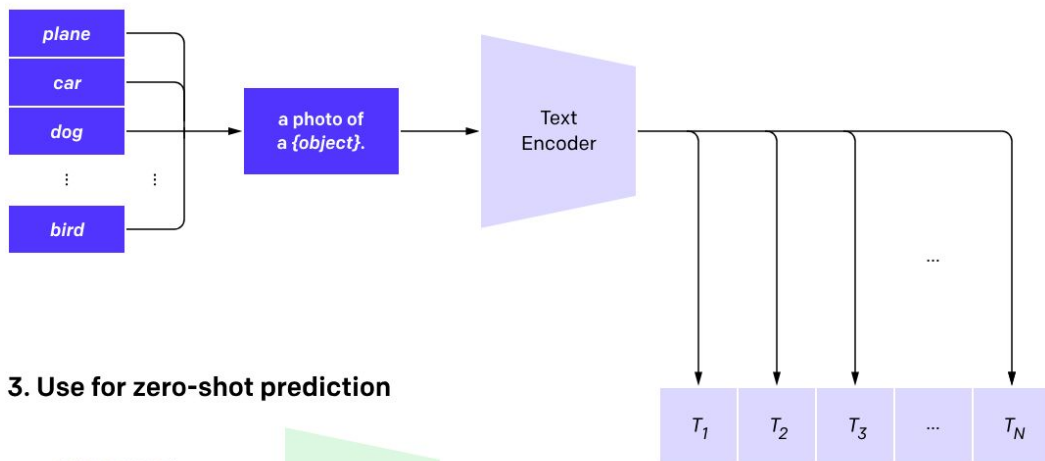
CLIP (Contrastive Language-Image Pretraining)

1. Contrastive pre-training

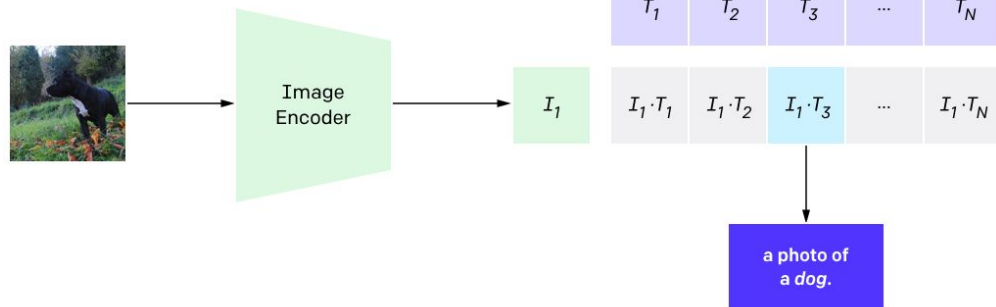


CLIP (Contrastive Language-Image Pretraining)

2. Create dataset classifier from label text



3. Use for zero-shot prediction



Setting

128M image-caption pairs, taken from DataComp ([Gadre et al, 2023](#)) medium

Split into 10 buckets of 12.8M images either by

- CLIP Score - cosine similarity between text and image embeddings
- T-MARS ([Maini et al, 2023](#)) - CLIP score after masking out text in images

Setting

128M image-caption pairs, taken from DataComp ([Gadre et al, 2023](#)) medium

Split into 10 buckets of 12.8M images either by

- CLIP Score - cosine similarity between text and image embeddings
- T-MARS ([Maini et al, 2023](#)) - CLIP score after masking out text in images

For a sufficiently high compute budget, there is an inherent tradeoff between seeing new samples or repeating existing high quality samples

Optimal filtering depends on compute budget

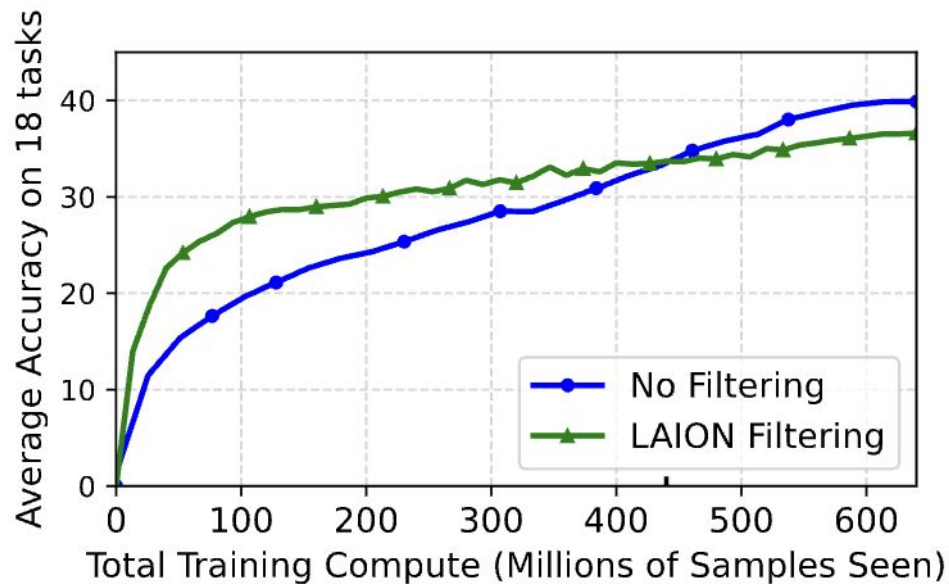


Figure 2. Given an initial data pool of 128M samples, we train ViT-B/32 CLIP models for a total of 640M samples. As we train for longer, the accuracy gains using the LAION-filtering subset that filters the common crawl to 10% of its initial size plateau. Surprisingly, even no-filtering of the common crawl is better than the popular LAION-filtering after seeing more than 450M samples.

Optimal filtering depends on compute budget

As you scale the compute budget, diversity becomes more important than quality

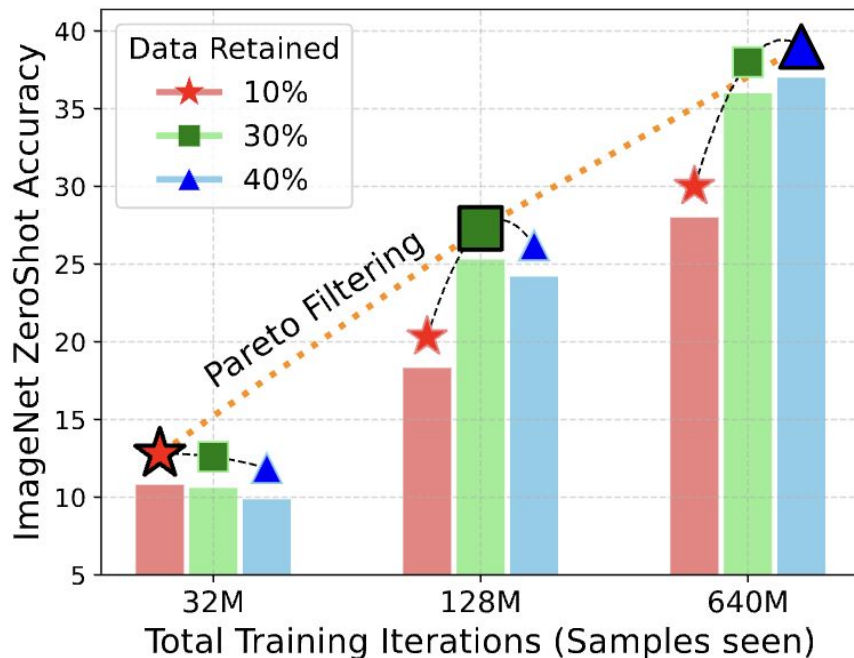


Figure 3. We vary the CLIP filtering threshold after ranking the data by their metric. While the original paper proposed retaining 30% of the data, our results show that depending on the ratio of compute to data pool size, we must adaptively make the filtering less (or more) aggressive to account for the diminishing utility of good data with repetitions. Results are presented on an average of 18 visual understanding tasks with a global data pool size of 128M samples, and varying compute scales.

Can we predict the downstream
performance?

Scaling laws for a single bucket, single epoch

y = downstream performance

n = number of samples

a, d = problem-dependent constants

b = sample utility

$$y = a \cdot n^b + d$$

Scaling laws for a single bucket, single epoch

y = downstream performance

n = number of samples

a, d = problem-dependent constants

b = sample utility

$$y = a \cdot n^b + d$$

$$\frac{dy}{dn} = b \cdot \frac{y}{n}$$

Scaling laws for a single bucket, multiple epochs

$$y = a \cdot n^b + d$$

(1 epoch)

Scaling laws for a single bucket, multiple epochs

$$y = a \cdot n^b + d \quad (1 \text{ epoch})$$

$$y = a \cdot n^b \left(\frac{2n}{n}\right)^b \dots \left(\frac{kn}{(k-1)n}\right)^b + d \quad (k \text{ epochs with no utility decrease})$$

Scaling laws for a single bucket, multiple epochs

$$y = a \cdot n^b + d \quad (1 \text{ epoch})$$

$$y = a \cdot n^b \left(\frac{2n}{n}\right)^b \cdots \left(\frac{kn}{(k-1)n}\right)^b + d \quad (k \text{ epochs with no utility decrease})$$

$$y = a \cdot n^{b^{(1)}} \left(\frac{2n}{n}\right)^{b^{(2)}} \cdots \left(\frac{kn}{(k-1)n}\right)^{b^{(k)}} + d \quad (k \text{ epochs with different utilities})$$

Scaling laws for a single bucket, multiple epochs

$$y = a \cdot n^b + d \quad (1 \text{ epoch})$$

$$y = a \cdot n^b \left(\frac{2n}{n}\right)^b \cdots \left(\frac{kn}{(k-1)n}\right)^b + d \quad (k \text{ epochs with no utility decrease})$$

$$y = a \cdot n^{b^{(1)}} \left(\frac{2n}{n}\right)^{b^{(2)}} \cdots \left(\frac{kn}{(k-1)n}\right)^{b^{(k)}} + d \quad (k \text{ epochs with different utilities})$$

$$y = a \cdot n^{b \cdot \delta^0} \left(\frac{2n}{n}\right)^{b \cdot \delta^1} \cdots \left(\frac{kn}{(k-1)n}\right)^{b \cdot \delta^{k-1}} + d \quad (\text{assume utility exponentially decays})$$

$$b^{(k)} = b \cdot \delta^{k-1}$$

Scaling laws for multiple buckets, single epoch

Initialize model state at performance y_0 , trained on n_0 samples

Make a mixed dataset of two samples with utility and decay (b_1, δ_1) and (b_2, δ_2)

What is the effective utility of a sample in this new dataset?

Scaling laws for multiple buckets, single epoch

Initialize model state at performance y_0 , trained on n_0 samples

Make a mixed dataset of two samples with utility and decay (b_1, δ_1) and (b_2, δ_2)

What is the effective utility of a sample in this new dataset?

$$\frac{y_2}{y_0} = \left(\frac{n_0 + 1}{n_0} \right)^{b_1} \left(\frac{n_0 + 2}{n_0 + 1} \right)^{b_2} \quad (\text{two gradient steps})$$

Scaling laws for multiple buckets, single epoch

Initialize model state at performance y_0 , trained on n_0 samples

Make a mixed dataset of two samples with utility and decay (b_1, δ_1) and (b_2, δ_2)

What is the effective utility of a sample in this new dataset?

$$\frac{y_2}{y_0} = \left(\frac{n_0 + 1}{n_0}\right)^{b_1} \left(\frac{n_0 + 2}{n_0 + 1}\right)^{b_2} \quad (\text{two gradient steps})$$

$$\approx \left(\frac{n_0 + 1}{n_0}\right)^{b_1} \left(\frac{n_0 + 1}{n_0}\right)^{b_2} \quad (\text{assume large } n_0)$$

Scaling laws for multiple buckets, single epoch

Initialize model state at performance y_0 , trained on n_0 samples

Make a mixed dataset of two samples with utility and decay (b_1, δ_1) and (b_2, δ_2)

What is the effective utility of a sample in this new dataset?

$$\frac{y_2}{y_0} = \left(\frac{n_0 + 1}{n_0}\right)^{b_1} \left(\frac{n_0 + 2}{n_0 + 1}\right)^{b_2} \quad (\text{two gradient steps})$$

$$\approx \left(\frac{n_0 + 1}{n_0}\right)^{b_1} \left(\frac{n_0 + 1}{n_0}\right)^{b_2} \quad (\text{assume large } n_0)$$

$$= \left(\frac{n_0 + 1}{n_0}\right)^{b_1 + b_2} \quad (\text{collapse into single term})$$

Scaling laws for multiple buckets, single epoch

Initialize model state at performance y_0 , trained on n_0 samples

Make a mixed dataset of two samples with utility and decay (b_1, δ_1) and (b_2, δ_2)

What is the effective utility of a sample in this new dataset?

$$\frac{y_2}{y_0} = \left(\frac{n_0 + 1}{n_0}\right)^{b_1} \left(\frac{n_0 + 2}{n_0 + 1}\right)^{b_2} \quad (\text{two gradient steps})$$

$$\approx \left(\frac{n_0 + 1}{n_0}\right)^{b_1} \left(\frac{n_0 + 1}{n_0}\right)^{b_2} \quad (\text{assume large } n_0)$$

$$= \left(\frac{n_0 + 1}{n_0}\right)^{b_1 + b_2} \quad (\text{collapse into single term})$$

$$\approx \left(\frac{n_0 + 2}{n_0}\right)^{\frac{b_1 + b_2}{2}} \quad (\text{split back into two data points})$$

Scaling laws for multiple buckets, single epoch

Initialize model state at performance y_0 , trained on n_0 samples

Make a mixed dataset of two samples with utility and decay (b_1, δ_1) and (b_2, δ_2)

What is the effective utility of a sample in this new dataset?

$$\frac{y_2}{y_0} = \left(\frac{n_0 + 1}{n_0}\right)^{b_1} \left(\frac{n_0 + 2}{n_0 + 1}\right)^{b_2} \quad (\text{two gradient steps})$$

$$\approx \left(\frac{n_0 + 1}{n_0}\right)^{b_1} \left(\frac{n_0 + 1}{n_0}\right)^{b_2} \quad (\text{assume large } n_0)$$

$$= \left(\frac{n_0 + 1}{n_0}\right)^{b_1 + b_2} \quad (\text{collapse into single term})$$

$$\approx \left(\frac{n_0 + 2}{n_0}\right)^{\frac{b_1 + b_2}{2}} \quad (\text{split back into two data points})$$

$$b_{\text{eff}} = \frac{b_1 + b_2}{2}$$

Scaling laws for multiple buckets, multiple epochs

$$b_{\text{eff}} = \frac{b_1 + b_2}{2}$$

(two buckets, one epoch)

Scaling laws for multiple buckets, multiple epochs

$$b_{\text{eff}} = \frac{b_1 + b_2}{2} \quad (\text{two buckets, one epoch})$$

$$b_{\text{eff}} = \frac{1}{p} \sum_{i \in [p]} b_i \quad (p \text{ buckets, one epoch})$$

Scaling laws for multiple buckets, multiple epochs

$$b_{\text{eff}} = \frac{b_1 + b_2}{2} \quad (\text{two buckets, one epoch})$$

$$b_{\text{eff}} = \frac{1}{p} \sum_{i \in [p]} b_i \quad (p \text{ buckets, one epoch})$$

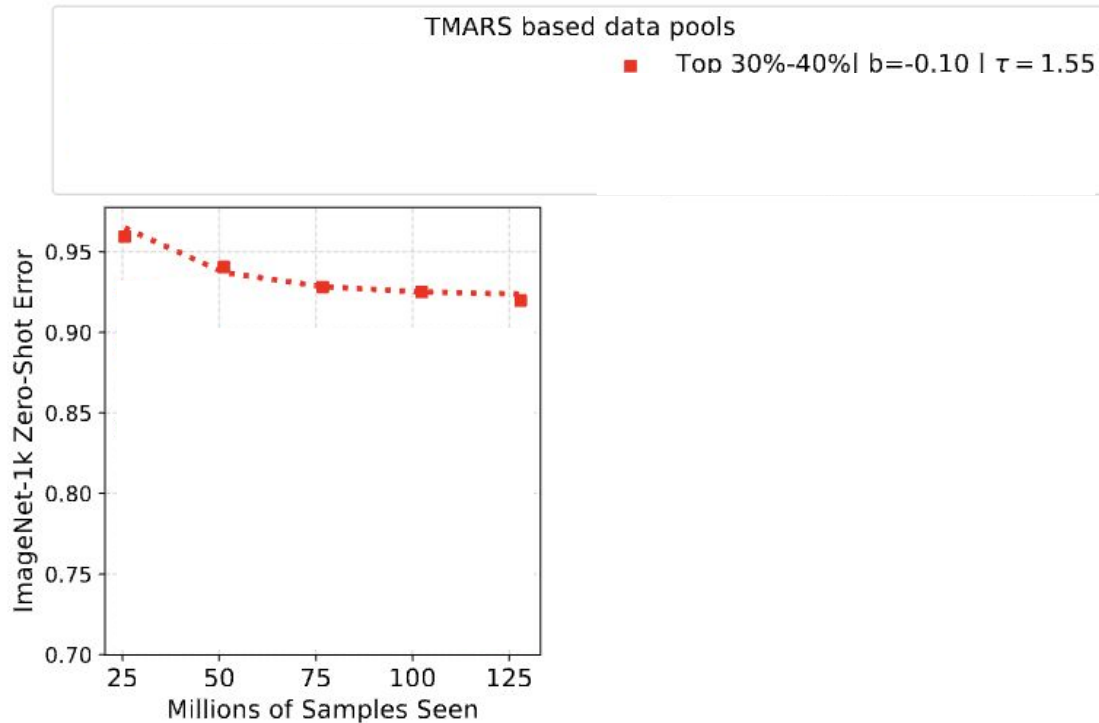
$$b_{\text{eff}}^{(k)} = \frac{1}{p} \sum_{i \in [p]} b_i \delta_i^{k-1} \quad (p \text{ buckets, } k \text{ epochs})$$

Alternative formulations, not tested

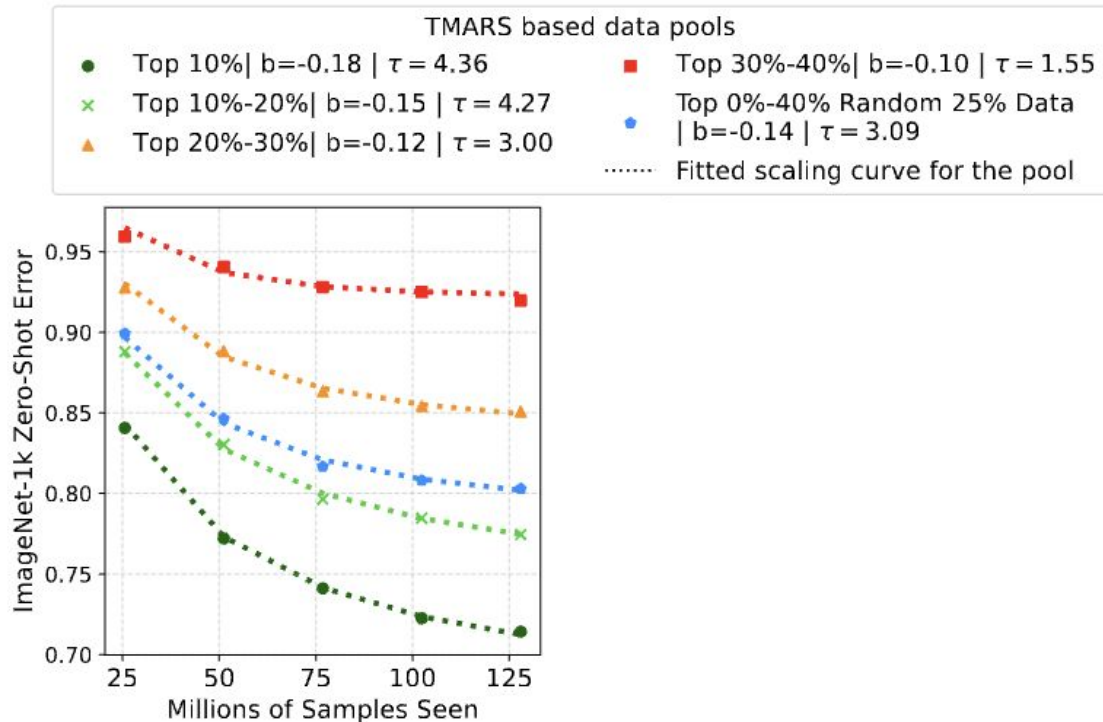
Utility does not have to decay exponentially

Instead of decreasing effective utility, you could decrease effective sample count ([Muenninghoff et al, 2023](#))

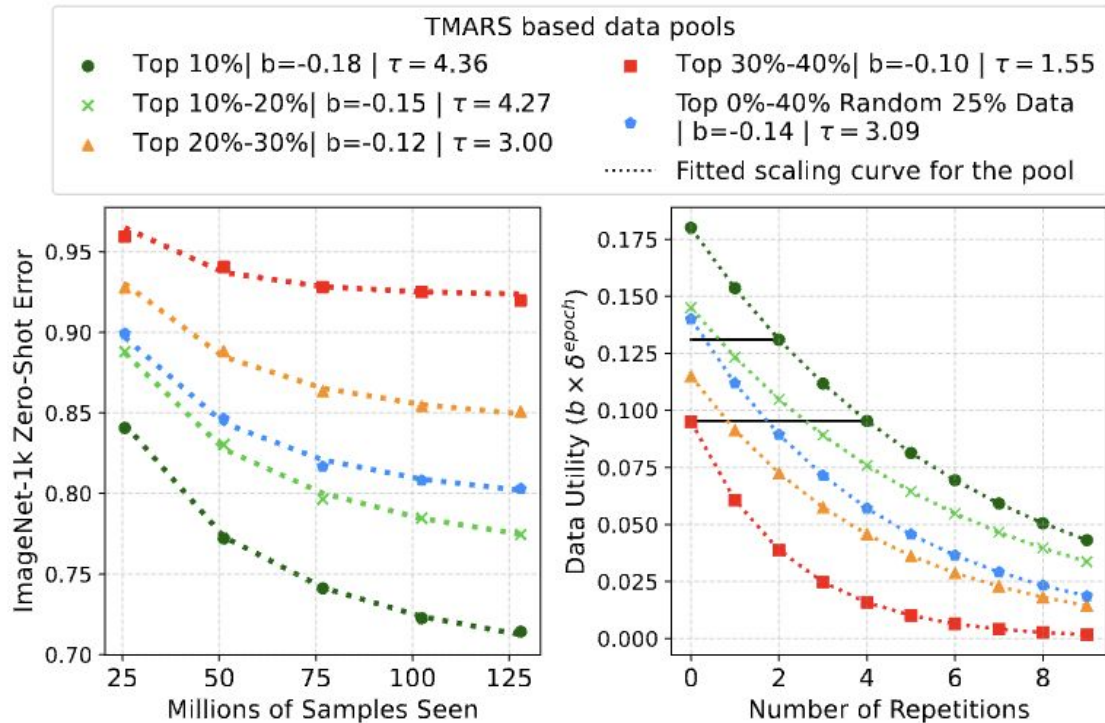
Fitting scaling laws to CLIP buckets



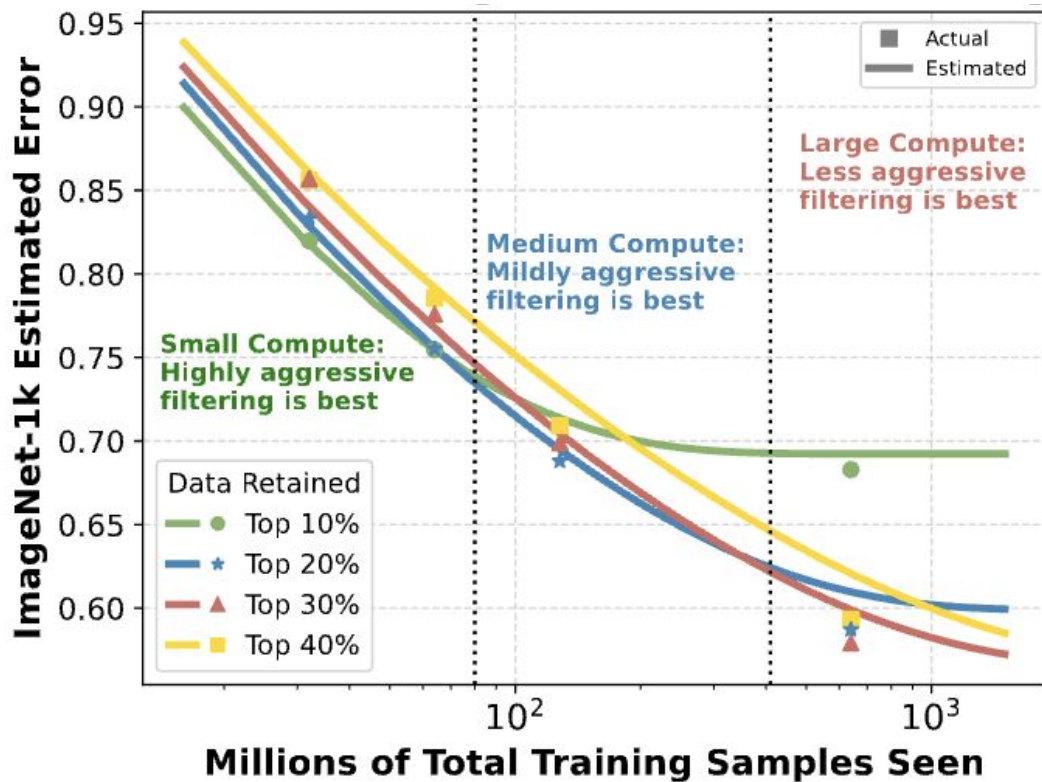
Fitting scaling laws to CLIP buckets



Fitting scaling laws for single bucket, multi epoch



Testing scaling laws for multiple buckets, multi epoch



Conclusion

In settings where repetition reduces the value of a data point, there is a tradeoff between same high quality data or more low quality data

The optimal decision depends on the amount of available compute

Scaling laws (may be able to) predict the optimal mixture

Open Problems

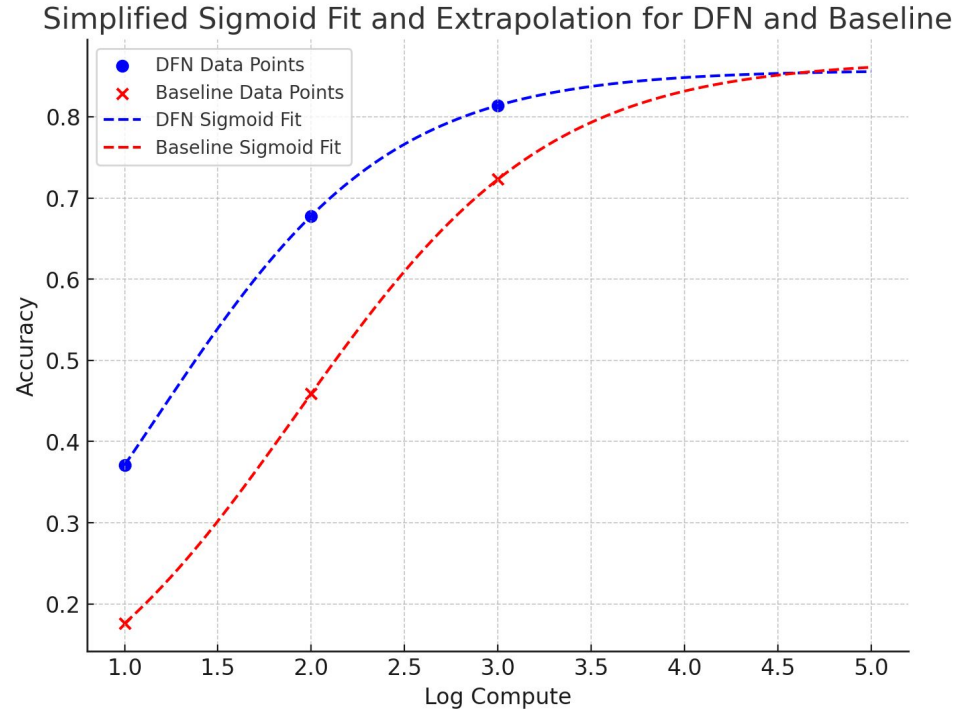
1. Importance of filtering in the limit

If compute increases and data stays fixed, data filtering will be less relevant. How much data filtering would one like to do in the infinite compute limit?

1. Importance of filtering in the limit

If compute increases and data stays fixed, data filtering will be less relevant. How much data filtering would one like to do in the infinite compute limit?

Open problem: build scaling laws for the effective benefit of data as compute increases



2. Distributionally aware filtering

Current data filtering strategies consider each domain to have additive impact. However, the composition of your data matters: the benefit of code data may depend on the amount of math data. With unlimited data, today's strategies may just pick a single domain, which is known to be suboptimal

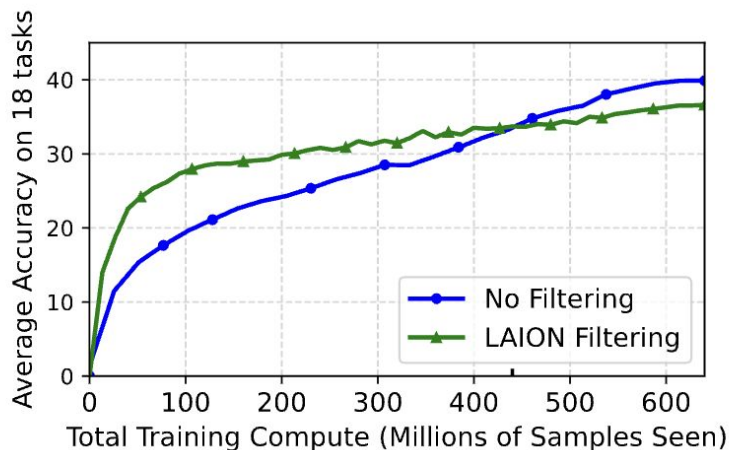
$$y_i = f(\langle \boldsymbol{\theta}^*, \mathbf{x}_i \rangle + \epsilon_i)$$

$$y = a \cdot n^b + d$$

This is really important for synthetic data, since one possible future has infinite data from each domain with strong mode collapse concerns

3. Curriculum: does varying mixture over time help?

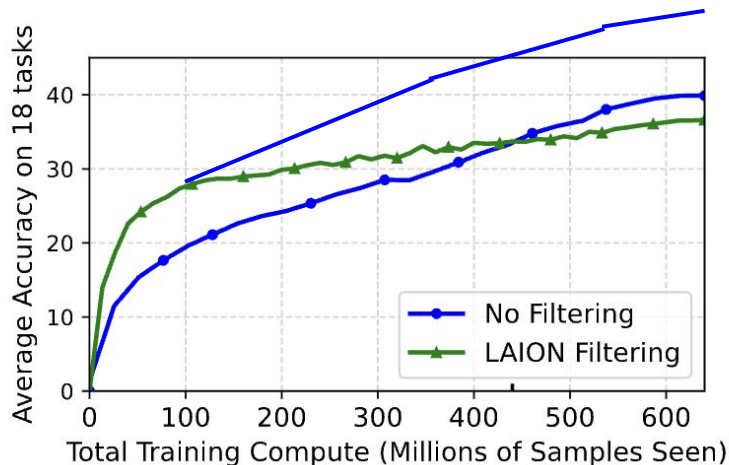
Given domains and a target validation set, we know that the optimal mixture can be better than the mixture of the validation set.



3. Curriculum: does varying mixture over time help?

Given domains and a target validation set, we know that the optimal mixture can be better than the mixture of the validation set.

Can varying the mixture over time be much better than the optimal mixture?



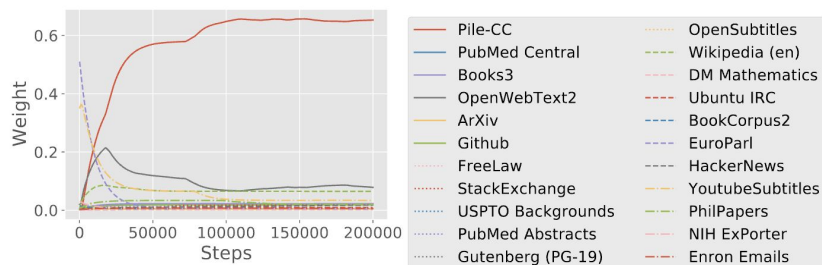
3. Curriculum: does varying mixture over time help?

Given domains and a target validation set, we know that the optimal mixture can be better than the mixture of the validation set.

Can varying the mixture over time be much better than the optimal mixture?

One Evidence Against

Doremi mixture stabilizes quickly



One Evidence For

Industry models mix high quality data at the end

“We used curriculum learning for pretraining, changing the data mix during training in ways we found to substantially improve model quality.” - DBRX

“We stage training to alter the mixture composition during training – increasing the weight of domain-relevant data towards the end of training” - Gemini