

A Pretrainer's Guide to Training Data: Measuring the Effects of Data Age, Domain Coverage, Quality, & Toxicity

Shayne Longpre¹, Gregory Yauney², Emily Reif³, Katherine Lee^{2,3}, Adam Roberts³,
Barret Zoph⁴, Denny Zhou³, Jason Wei⁴, Kevin Robinson³, David Mimno², Daphne Ippolito^{3,5}
¹ MIT ² Cornell University ³ Google Research ⁴ OpenAI ⁵ Carnegie Mellon University

All pretraining data is curated, but data curation decisions are not always disclosed.

1. Practitioners are guided by intuition.
2. Experiments are frequently repeated because results are not disclosed.
3. Data curation has large impact because pretrained models are reused.

We pretrain 28 LMs at the 1.5B-parameter scale on differently-curated pretraining datasets in order to measure the effects of curation choices.

Compute is expensive! But so is dark data & documentation debt.

Datasets: C4 and the Pile

Models: 1.5B-parameter decoder-only autoregressive transformers

Setting: Pretrain, then finetune on downstream tasks individually

Takeaways

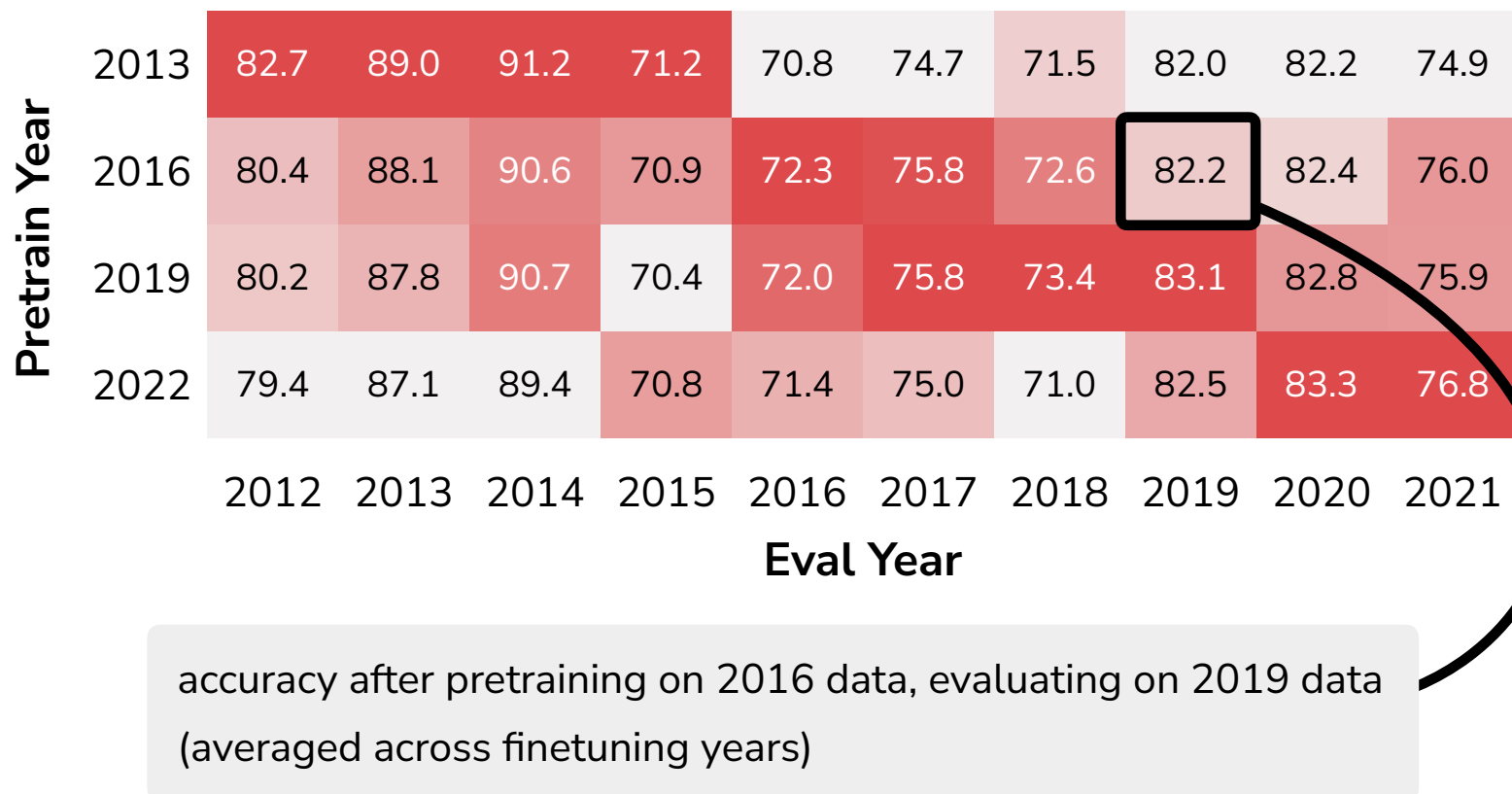
1. Stale pretraining data matters and is not overcome by finetuning!
2. Temporal misalignment effects grow with model size.
3. “Quality” filters boost performance, even while reducing training data.
4. Toxicity filters hurt. Inverse toxicity filters can help a lot for some tasks.
5. Data heterogeneity and quantity matter most, especially web and books data.

Data Age

Mismatch in data age between pretraining and evaluation data causes performance degradation.

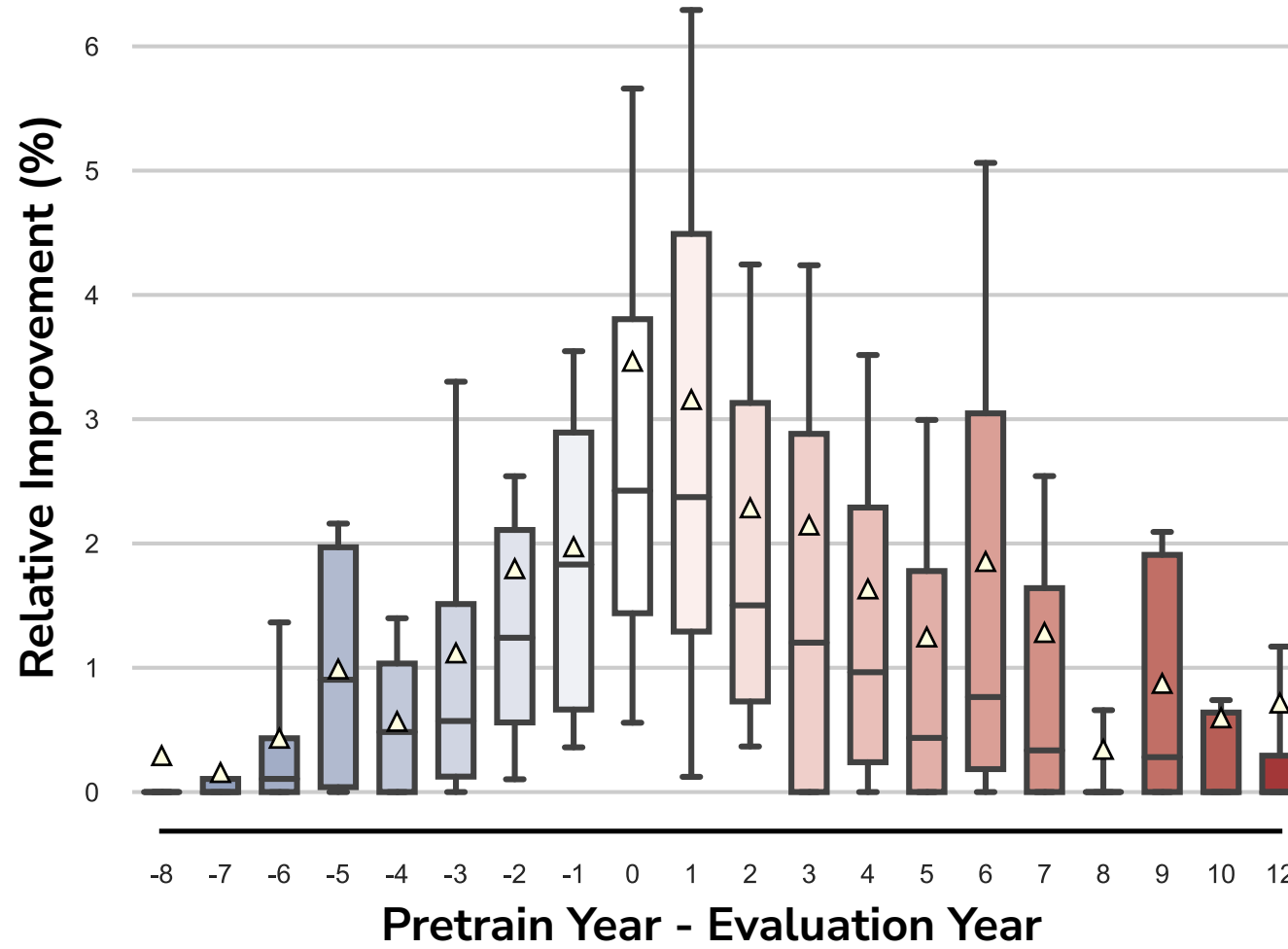
1. Less impact than finetuning mismatch, but adds up.
2. Release age distributions for pretraining data.

Example dataset: PoliAff



Accuracy is higher when pretraining and eval year are closer in time, even after finetuning

Pretrained models become stale



Temporal degradation happens faster when evaluating old models on new benchmarks.

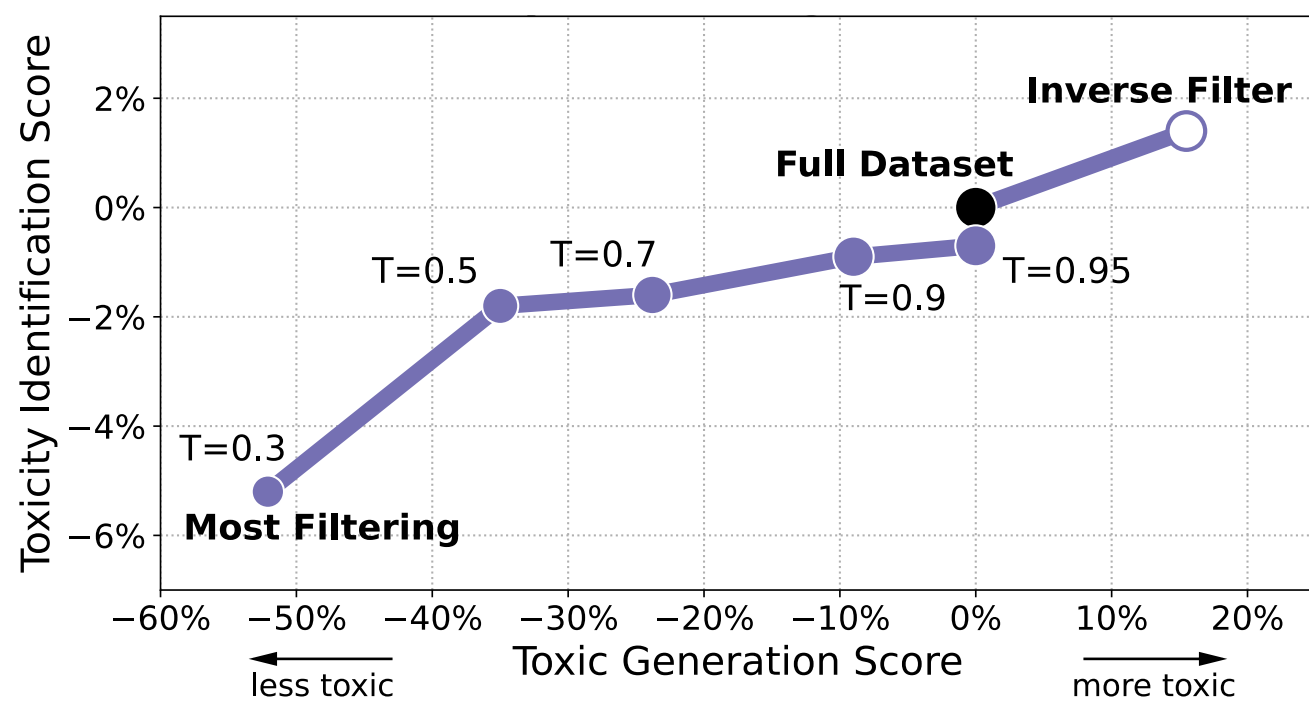
Toxicity and Quality

Toxicity: Perspective API, a classifier that assigns every document a score from 0 (nontoxic) to 1 (toxic)

Quality: GLaM/PaLM classifier, Wikipedia + books are “high quality”, every document gets a score from 0 (high quality) to 1 (low quality)

- Scalable, consistent with current practice
- Many downsides
- Lots of open questions!

Toxicity filtering induces a tradeoff: reduces toxic generation at the cost of decreased toxicity identification.



If the goal is to identify toxic text, then training on toxic data is more effective

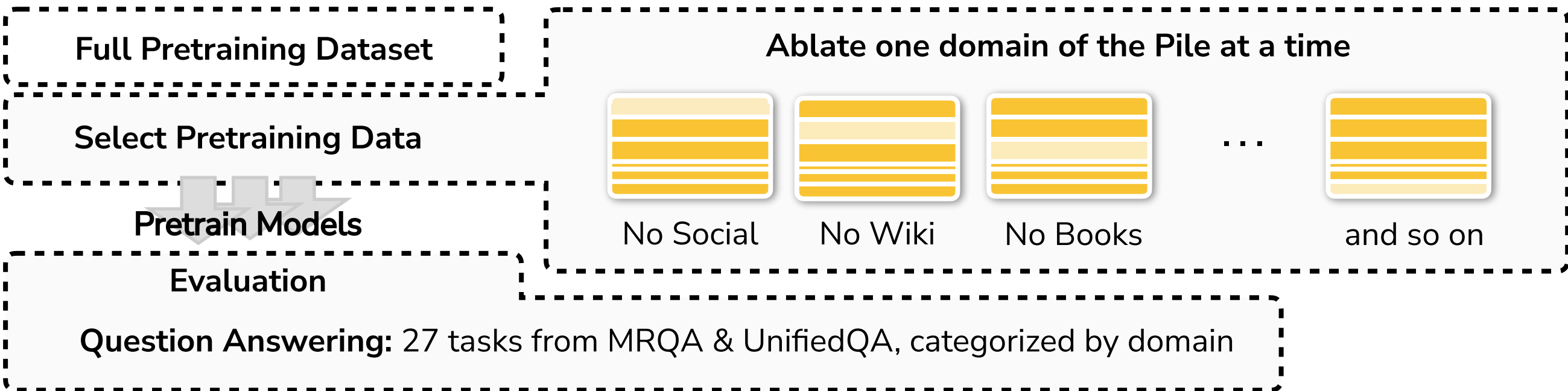
Content filtering impacts downstream QA performance

	Filter	Data	QA domain				Mean
			Wiki	Web	Acad	CS	
Baseline	Full Data	100%	0	0	0	0	0
Toxicity	Light	95%	-2.2	-1.1	0.2	0.2	-0.7
	Heavy	76%	-4.2	-2.4	-1.1	-3.5	-2.7
	Inverse	92%	0.4	-1.4	4.9	2.7	1.7
Quality	Light	91%	1.2	0.7	6.4	6.1	2.5
	Heavy	73%	-0.3	0.8	0.8	6.8	1.2
	Inverse	73%	-5.0	-4.5	-2.7	-6.4	-3.1

1. Toxicity filtering hurts QA performance across domains.

2. Quality filtering improves performance across most domains, despite removing data.

Domain Coverage



Heterogeneous domains have biggest effect on QA performance

	Wiki	Web	Biomed	Academic	Common Sense	Contrast Sets	Average
Full Dataset (100%)	0.0	0.0	0.0	0.0	0.0	0.0	0.0
No Social (99%)	-0.8	-3.7	0.1	3.5	-3.5	3.5	0.3
No Wiki (98%)	-1.3	-5.3	0.2	0.9	-4.4	7.2	-0.4
No Books (93%)	-3.5	-6.3	0.0	-1.6	-6.5	-4.4	-2.8
No OpenWeb (93%)	-2.0	-4.1	-1.0	0.6	-5.8	-2.9	-1.4
No Legal (91%)	-2.7	-2.9	0.4	0.8	-2.6	-0.4	-0.7
No Academic (87%)	-0.3	-2.5	-0.9	2.2	-1.1	4.3	0.2
No Pubmed (85%)	-0.3	-3.0	-5.8	-1.5	-5.9	3.9	-1.4
No Code (81%)	-0.5	-3.1	-1.2	1.2	-5.8	4.4	-0.2
No CC (73%)	-3.2	-6.2	-4.6	-5.9	-8.0	-5.2	-4.9

1. Removing Books and Common Crawl domains hurt downstream performance most.

2. Targeted data helps for targeted evaluations.

Paper: arxiv.org/pdf/2305.13169.pdf