

The most expensive part of an LLM should be its training data

Colin Raffel

(presenting work primarily by Nikhil Kandpal)

Language model training costs have increased ~6.5x each year



From "Data on Notable AI Models" by EpochAI

Increasing model sizes leads to increasing training dataset sizes



From "Data on Notable AI Models" by EpochAI

Typical training datasets comprise many sources

LLaMa's Training Data Sources



Writing text involves human labor

Writing Medium	Typical Length (words)	Estimated Writing Time
Blog Post	2,000	1.1 hours
Academic Paper	5,000	2.8 hours
Novel	70,000	1.6 days
Textbook	300,000	1 week
Encyclopedia Britannica	40,000,000	2.5 years

Human labor has a cost



From https://www.visualcapitalist.com/minimum-wage-around-the-world/

Training data would cost much more than training if annotators were paid



From "The Most Expensive Part of an LLM should be its Training Data" by Kandpal and Raffel

Copyright-related lawsuits against OpenAI/Microsoft as of March 2024

Coders

1. Joseph Saveri Firm: overview, complaint

Writers

- 2. Joseph Saveri Firm: overview, complaint
- 3. Authors Guild & Alter: overview, complaint
- 4. Nicholas Gage: overview & complaint

Media

- 5. New York Times: overview, complaint
- 6. Intercept Media: overview, complaint
- 7. Raw Story & Alternet: overview, complaint
- 8. Denver Post & seven others: overview, complaint
- 9. Center for Investigative Reporting: overview, complaint

From https://forum.effectivealtruism.org/posts/EdBkBBFkCaHrw2iwL/twelve-lawsuits-against-openai

Idea 1:

Train only on "permissively licensed" text

"... it would be impossible to train today's leading AI models without using copyrighted materials."

- OpenAI

... let's try!

What is "permissively licensed" (to us)?

- I'm not a lawyer, but I know people who know lawyers
- Public domain/CC0
 - Mostly very old and/or governmental text
- Creative Commons-Attribution (CC-BY, CC-BY-SA)
 - Non-commercial (NC) considered non-permissive
 - "No derivative works" is ambiguous
 - (so is attribution, sort of...)
- Blue Oak Council gold/silver/bronze licenses
 - BSD, MIT, Apache, etc.
- Licenses "equivalent" to the above
- <u>https://github.com/r-three/common-pile/blob/main/licensed_pile/licenses.py</u>

Common Pile raw source sizes (in UTF-8 bytes)



Proportion of licence types



Proportion of source types



How it's going



Try it out!

😕 Hugging Face	Search models, dataset	s, users	Models = Data	asets 📑 Spaces	Posts Docs Docs Pricing	~≡ 🛛 👰
 Datasets: nkand Modalities: Text Form 	pa2/ common-pile ats: (+) json Size: [1B	e-filtered 🗅	⊘ like 0 Datasets Ø Dask ♥ Croi	issant		
Dataset card Hi	ewer 🛛 🗏 Files and vers	sions 🥚 Commun	ity 2			
This dataset has 1 file scale	nned as unsafe. Show files				Downloads last month	1,150
Dataset Viewer (First 5	GB) (1)	S Auto-converted	to Parquet 🛷 API 🖥 Embed	Full Screen Viewer	↔ Use this dataset <	1
default · ~1.18B rows (showi	ng the first 2.53M)	✓ train · ~1.1	8B rows (showing the first 2.53M)	SQL Console	Size of the auto-converted Parquet files (First 5GB per spl 48.6 GB	lit):
added string · lengths	<pre>created</pre>	id string · lengths	metadata 💠	source string · <i>classes</i>	Number of rows (First 5GB per split): 17,386,843	
19 26 2024-08-12T18:16:26.58537	70 2007-04-02T19:18:42	9 16	{ "authors": "C. Bal	2 values arxiv-abstracts	Estimated number of rows: 1,365,167,755	
2024-08-12T18:16:26.58557	2007-03-31T02:26:18	0704.0002	{ "authors": "Ileana Streinu and Louis…	arxiv-abstracts		
2024-08-12T18:16:26.58564	4 2007-04-01T20:46:54	0704.0003	{ "authors": "Hongjun Pan",…	arxiv-abstracts		
2024-08-12T18:16:26.58569	2007-03-31T03:16:14	0704.0004	{ "authors": "David Callan",	arxiv-abstracts		
2024-08-12T18:16:26.58572	25 2007-04-02T18:09:58	0704.0005	<pre>{ "authors": "Wael Abu- Shammala and Alberto</pre>	arxiv-abstracts		
2024-08-12T18:16:26.58576	0 2007-03-31T04:24:59	0704.0006 2 3 25	and C. K. Law",	arxiv-abstracts		

https://huggingface.co/datasets/nkandpa2/common-pile-filtered

Idea 2: Attribute predictions back to training data

Context attribution finds influential spans in a model's input

Context 📚

- <solar_eclipse_2024.pdf>
- weekly_weather_forecast.pdf>

Query 🔎

I live in Boston, MA. When and where should I go to see the eclipse?



Generated response 💡

You can drive to Burlington, Vermont to see the eclipse on April 8, 2024. The weather forecast for Burlington is sunny with temperature in the 60s^[1]...

Context attribution

solar_eclipse_2024.pdf:

A solar eclipse will be visible in the continental United States on April 8, 2024. The path of totality arches from Mexico to Texas to Maine ...

weekly_weather_forecast.pdf:

... [1] The weather in Burlington should be sunny, with mostly clear skies and temperatures ranging from mid to high 60s ...

From "ContextCite: Attributing Model Generation to Context" by Cohen-Wang et al.

Observation 1: Smaller models provide similar attributions



From "AttriBoT: A Bag of Tricks for Efficiently Approximating Leave-One-Out Context Attribution" by Liu et al.

Observation 2: Attributions combine linearly



From "AttriBoT: A Bag of Tricks for Efficiently Approximating Leave-One-Out Context Attribution" by Liu et al.

AttriBoT is pareto-optimal for context attribution



From "AttriBoT: A Bag of Tricks for Efficiently Approximating Leave-One-Out Context Attribution" by Liu et al.

Try it out!

KattriBoT Public	S Edit Pins	▼ O Watch O ▼	양 Fork 0 ▼ ☆ Star 5 ▼	r
	Go to file	+ <> Code •	About ध्	Ķ
FY-Liu Fixed gradnorm device iss	sue. e3f	fcea · 3 months ago 🕚	AttriBoT: A Bag of Tricks for Efficiently Approximating Leave-	-
context_attribution	Fixed gradnorm device issue.	3 months ago		
context_attribution_simple.e	Initial upload.	3 months ago	-^- Activity	
example	Initial upload.	3 months ago	Custom properties	
🗋 README.md	Initial upload.	3 months ago	☆ 5 stars ⊙ 0 watching	
🗋 requirements.txt	Initial upload.	3 months ago	양 0 forks	
🗅 setup.py	Initial upload.	3 months ago	Report repository	

Measuring training data influence



scoring rule ϕ assigns a fair share of u to each player

From https://ml-data-tutorial.org/

AirRep is pareto-optimal for group data influence



From "Enhancing Training Data Attribution with Representational Optimization" by Sun et al.

Thanks. Please give me feedback: <u>http://bit.ly/colin-talk-feedback</u> <u>craffel@qmail.com</u>