

training data : raw text

span a broad range of domains, genres, languages

Web: Google Search Index 100 PB

WalMart 2.5 PB / hour

Representation:

Representation. Despite the richness of web data, it has been noted in [Bender et al, 2021](#) that:

- Despite the size, large-scale data still has **uneven representation** over the population.
- Internet data overrepresents younger users from developed countries.
- GPT-2's training data is based on Reddit, which according to Pew Internet Research's 2016 survey, 67% of Reddit users in the US are men, 64% between ages 18 and 29.
- 8.8-15% of Wikipedians are female.
- Harassment on Internet could turn away certain people (trans, queer, neurodivergent people).
- Filtering "bad words" could further marginalize certain populations (e.g., LGBT+).

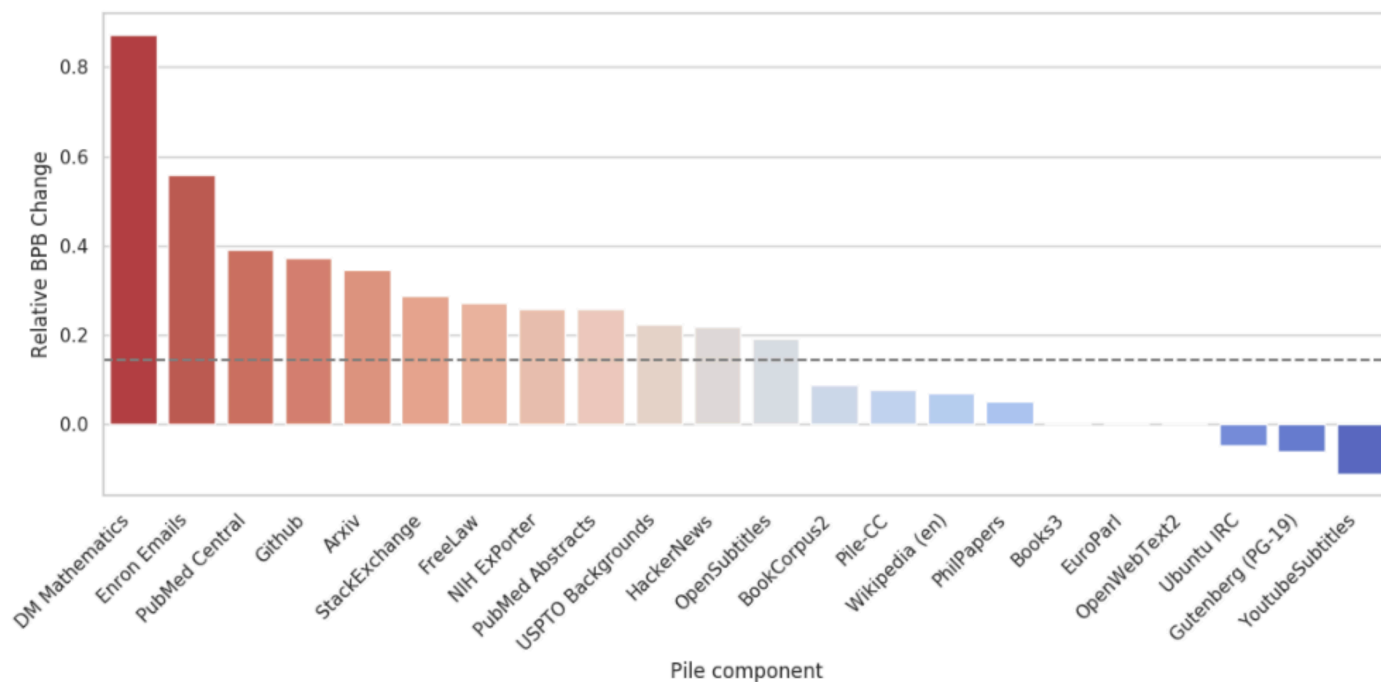
Takeaway: it is crucial to understand and document the composition of the datasets used to train large language models.

WebText - 40 GB text - GPT-2

OpenWebText - 38 GB

Colossal Clean Crawled Corpus - 806 GB text
(156 billion tokens)
- TS

Component	Raw Size	Weight	Epochs	Effective Size	Mean Document Size
Pile-CC	227.12 GiB	18.11%	1.0	227.12 GiB	4.33 KiB
PubMed Central	90.27 GiB	14.40%	2.0	180.55 GiB	30.55 KiB
Books3 [†]	100.96 GiB	12.07%	1.5	151.44 GiB	538.36 KiB
OpenWebText2	62.77 GiB	10.01%	2.0	125.54 GiB	3.85 KiB
ArXiv	56.21 GiB	8.96%	2.0	112.42 GiB	46.61 KiB
Github	95.16 GiB	7.59%	1.0	95.16 GiB	5.25 KiB
FreeLaw	51.15 GiB	6.12%	1.5	76.73 GiB	15.06 KiB
Stack Exchange	32.20 GiB	5.13%	2.0	64.39 GiB	2.16 KiB
USPTO Backgrounds	22.90 GiB	3.65%	2.0	45.81 GiB	4.08 KiB
PubMed Abstracts	19.26 GiB	3.07%	2.0	38.53 GiB	1.30 KiB
Gutenberg (PG-19) [†]	10.88 GiB	2.17%	2.5	27.19 GiB	398.73 KiB
OpenSubtitles [†]	12.98 GiB	1.55%	1.5	19.47 GiB	30.48 KiB
Wikipedia (en) [†]	6.38 GiB	1.53%	3.0	19.13 GiB	1.11 KiB
DM Mathematics [†]	7.75 GiB	1.24%	2.0	15.49 GiB	8.00 KiB
Ubuntu IRC	5.52 GiB	0.88%	2.0	11.03 GiB	545.48 KiB
BookCorpus2	6.30 GiB	0.75%	1.5	9.45 GiB	369.87 KiB
EuroParl [†]	4.59 GiB	0.73%	2.0	9.17 GiB	68.87 KiB
HackerNews	3.90 GiB	0.62%	2.0	7.80 GiB	4.92 KiB
YoutubeSubtitles	3.73 GiB	0.60%	2.0	7.47 GiB	22.55 KiB
PhilPapers	2.38 GiB	0.38%	2.0	4.76 GiB	73.37 KiB
NIH ExPorter	1.89 GiB	0.30%	2.0	3.79 GiB	2.11 KiB
Enron Emails [†]	0.88 GiB	0.14%	2.0	1.76 GiB	1.78 KiB
The Pile	825.18 GiB			1254.20 GiB	5.91 KiB



Takeaway: The Pile contains a lot of information that's not well covered by GPT-3's dataset.