

大模型预训练

通用人工智能与数字经济创新团队

西南财经大学

本章内容

- **8.1 概述**
- **8.2 预训练数据工程**
 - 8.2.1 预训练数据源
 - 8.2.2 多模态数据集
 - 8.2.3 数据处理
 - 8.2.4 模型性能关系
- **8.3 预训练方法**
 - 8.3.1 预训练任务
 - 8.3.2 优化参数设置
 - 8.3.3 可扩展训练技术
- **8.4 讨论**

8.1 概述

大模型的预训练训练，主要涉及到**数据源**和**分布式训练**两个关键方面

数据源

- ◆ 需要庞大的标注或者未标注数据集进行自监督学习，种类包括**文本、图像、声音**等
- ◆ 并需要经过**预处理**和**清洗**步骤。数据的质量和多样性直接影响模型的性能和泛化能力
- ◆ 为了获取全面的数据，有时需要聚合多个子数据集，会导致**数据冗余**和**不平衡**等问题

分布式训练

- 【解决问题】
- ◆ 由于单一计算设备（GPU）的计算能力和内存有限
- 【主要方法】
- ◆ 大模型通常需要分布式架构进行训练，通过**数据并行**、**模型并行**或**张量并行**等多种方式来实现

本章内容

- **8.1 概述**
- **8.2 预训练数据工程**
 - 8.2.1 预训练数据源
 - 8.2.2 多模态数据集
 - 8.2.3 数据处理
 - 8.2.4 模型性能关系
- **8.3 预训练方法**
 - 8.3.1 预训练任务
 - 8.3.2 优化参数设置
 - 8.3.3 可扩展训练技术
- **8.4 讨论**

8.2.1 预训练数据源

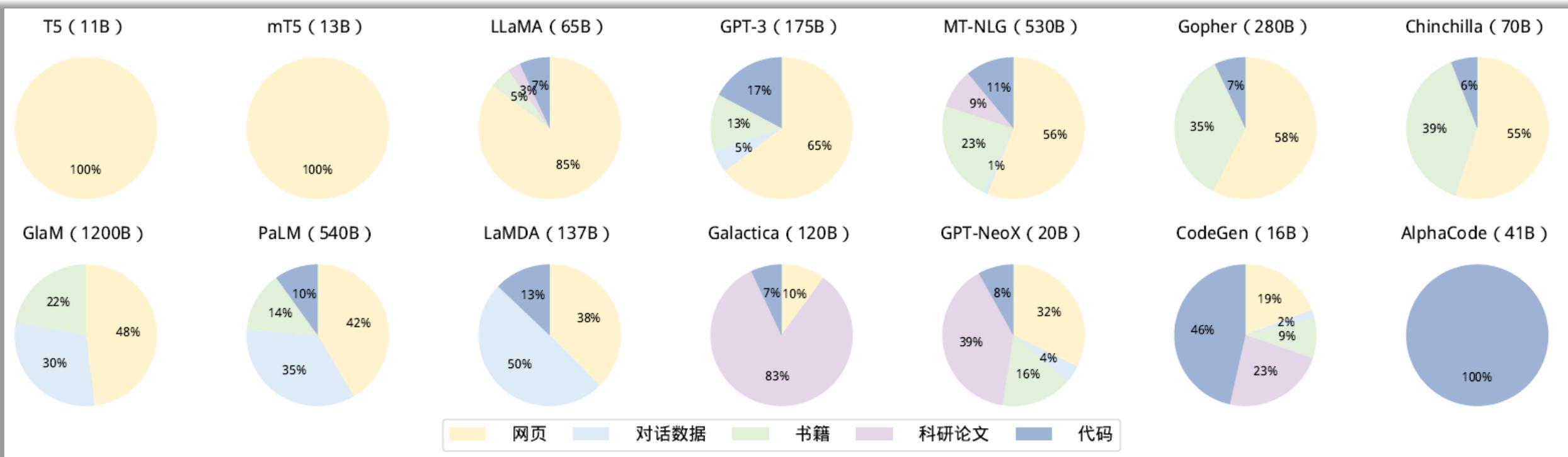
大模型预训练的性能高度依赖两个主要因素：1) 预训练语料库的质量与规模；2) 数据预处理的策略与方法。

通用数据

网页、书籍和对话文本等

专业数据

语言数据、科学论文和编程代码等

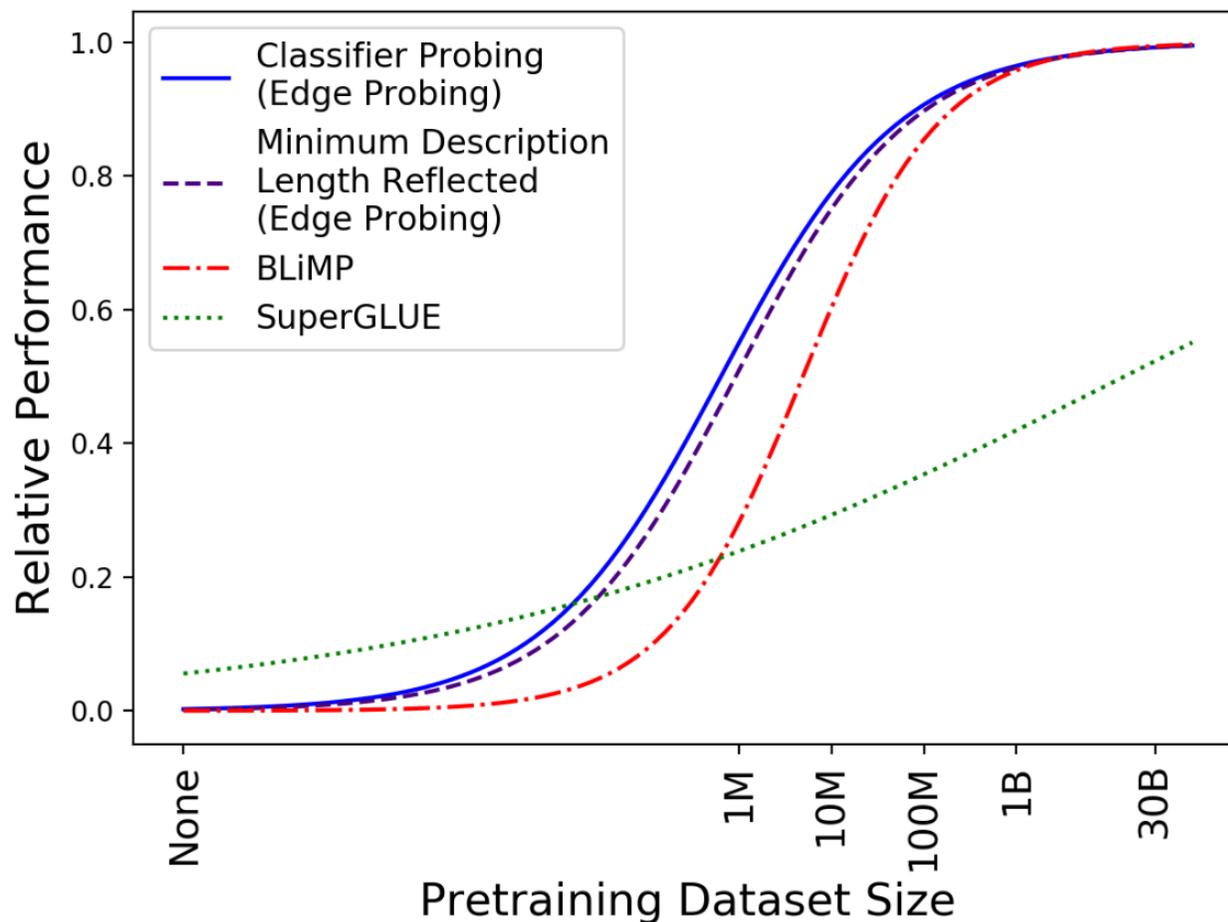


需要多少语料库才能学会“文字接龙”？

1. 语言知识

“我马上就”后面可以接“写”等动词，而不能接“猫”等名词

When Do You Need Billions of Words of Pretraining Data?



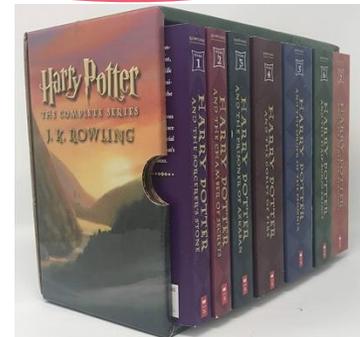
需要多少语料库才能学会“文字接龙”？

2. 世界知识

“中国的首都是”
后面接“北京”，
而不能接“上海”

模型	年份	数据规模
GPT-1	2018	7000本书
GPT-2	2019	40GB文本
GPT-3	2020	从45TB数据清洗得到570GB文本
Llama 3	2024	15TB文本
Qwen 2.5	2024	大概60TB数据

大概 300B
tokens, 相对于
《哈利·波特》
全集的30万遍



In terms of Qwen2.5, the language models, all models are pretrained on our latest large-scale dataset, encompassing up to **18 trillion tokens**.

8.2.1 预训练数据源

■ 通用数据源

网页数据源 互联网作为一个庞大的数据源，为语言模型提供了丰富的文本材料，具备规模大、动态、多语言和主题丰富等特点，是目前LLMs中使用最广泛的数据源。



代表性网页预训练语料库

数据集	发布者	发布时间	规模	特点
CC-Stories	Google Brain	2018-7	31 GB	基于 Common Crawl, 英文
RealNews	华盛顿大学等	2019-5	120 GB	基于 Common Crawl, 英文
C4	Google Research	2019-10	12.68TB	基于 Common Crawl, 英文
CLUECorpus2020	CLUE	2020-3	100 GB	基于 Common Crawl, 中文
CC100	Facebook AI	2020-7	2.5 TB	基于 Common Crawl, 100 种语言
WuDaoCorpora-Text	北京智源人工智能研究院	2021-6	5 TB	中文, 开源 200 GB
mC4	Google Research	2021-6	251 GB	基于 Common Crawl, 108 种语言
OSCAR 22.01	Inria	2022-1	8.41 TB	基于 Common Crawl, 151 种语言
MNBVC	里屋	2023-1	20.3 TB	中文, 包括网页、书籍、论文等
RefinedWeb	Falcon	2023-6	4.88 TB	基于 Common Crawl, 英文
Wan.Juan.Text	上海 AI 实验室	2023-8	1 TB	中文, 包括网页、书籍等

8.2.1 预训练数据源

■ 所有网上文本都能拿来训练吗？

The Times Sues OpenAI and Microsoft Over A.I. Use of Copyrighted Work

Millions of articles from The New York Times were used to train chatbots that now compete with it, the la

程序员起诉 Copilot

2021 年 6 月，GitHub 和 OpenAI 发布了 Copilot，可以“通过使用人工智能提供或填充代码块来帮助软件编码人员”。2021 年 8 月，OpenAI 又发布了 Codex，“可将自然语言转换为代码并集成到 Copilot 中”。GitHub 用户每月支付 10 美元或每年 100 美元才能访问 Copilot。Codex 和 Copilot 接受了“数十亿行”公开可用代码的训练，包括来自公共 GitHub 存储库的代码，诉讼由此而起。2023 年 5 月 11 日，美国加利福尼亚州北区地方法院针对 J. DOE 1 等诉 GitHub 等案做出了部分允许并部分拒绝驳回动议的裁定。该案的被告包括 GitHub、微软、OpenAI 等。

<https://www.nytimes.com/2023/12/27/business/media/new-york-times-open-ai-microsoft-lawsuit.html>

<https://36kr.com/p/2609988349031944>

8.2.1 预训练数据源

通用数据源

书籍数据集 作为长形式文本的优质来源，有助于语言模型在理解复杂语句结构和生成连贯文本方面的训练。

Smashwords

Project Gutenberg

Anna's Archive

The screenshot shows the Smashwords website with a search bar and navigation menu. It features a statistics section with the following data:

Words Published:	34.79 billion
Books Published:	901,708
Free Books:	92,456
Books on Sale:	13,020

Below the statistics are filters for price and length, and a list of categories including Adventure, African American fiction, and more.

The screenshot shows the Project Gutenberg website with a search bar and navigation menu. It features a welcome message and a list of featured eBooks:

Travels in Eastern Africa, volume 2 (of 2)	Travels in Eastern Africa, volume 1 (of 2)	Portraits of women by Gamaliel	My autobiography by Benito	Heath's French and English dictionary by	Elämän pyöreissä by Gudda	Loved and lost by Bertha M. Clay	The twelve best short stories in the
--	--	--------------------------------	----------------------------	--	---------------------------	----------------------------------	--------------------------------------

Additional text on the page includes: "Choose among free epub and Kindle eBooks, download them or read them online. You will find the world's great literature here, with focus on older works for which U.S. copyright has expired. Thousands of volunteers digitized and diligently proofread the eBooks, for you to enjoy."

The screenshot shows the Anna's Archive website with a search bar and navigation menu. It features a mission statement and a table of data sources:

Our mission is to archive all the books in the world (as well as papers, magazines, etc), and make them widely accessible. We believe that all books should be mirrored far and wide, to ensure redundancy and resiliency. This is why we're pooling together files from a variety of sources. Some sources are completely open and can be mirrored in bulk (such as Sci-Hub). Others are closed and protective, so we try to scrape them in order to "liberate" their books. Yet others fall somewhere in between.

All our data can be [torrented](#), and all our metadata can be [generated](#) or [downloaded](#) as Elasticsearch and MariaDB databases. The raw data can be manually explored through JSON files such as [this](#).

Source	Size	% mirrored by AA / torrents available	Last updated
Libgen.rs [lgrs] Non-Fiction and Fiction	7,379,910 files 83.0 TB	100% / 100%	2024-09-30
Sci-Hub [scihub] Via Libgen.II "scimag"	101,004,457 files 96.0 TB	87.052% / 87.052%	Sci-Hub: frozen since 2021; most available through torrents Libgen.II: minor additions since then
Libgen.II [lgli] Excluding "scimag"	19,457,185 files 281.4 TB	86.528% / 84.751%	2024-09-01 Fiction torrents are behind (though IDs ~4-6M not torrented since they overlap)

8.2.1 预训练数据源

■ 专业数据源

科学文本 这类数据集通常包括学术论文、专利和其他类型的专业文献。由于科学文本常常涉及专业术语、复杂的数据结构和公式，因此在预处理时需要采用特殊方法

arXiv is a free distribution service and an open-access archive for nearly 2.4 million scholarly articles in the fields of physics, mathematics, computer science, quantitative biology, quantitative finance, statistics, electrical engineering and systems science, and economics. Materials on this site are not peer-reviewed by arXiv.

Subject search and browse:
Physics Search Form Interface Catchup

Physics

- Astrophysics (**astro-ph** new, recent, search) Astrophysics of Galaxies; Cosmology and Nongalactic Astrophysics; Earth and Planetary Astrophysics; High Energy Astrophysical Phenomena; Instrumentation and Methods for Astrophysics; Solar and Stellar Astrophysics
- Condensed Matter (**cond-mat** new, recent, search) Disordered Systems and Neural Networks; Materials Science; Mesoscale and Nanoscale Physics; Other Condensed Matter; Quantum Gases; Soft Condensed Matter; Statistical Mechanics; Strongly Correlated Electrons; Superconductivity
- General Relativity and Quantum Cosmology (**gr-qc** new, recent, search)
- High Energy Physics - Experiment (**hep-ex** new, recent, search)
- High Energy Physics - Lattice (**hep-lat** new, recent, search)
- High Energy Physics - Phenomenology (**hep-ph** new, recent, search)
- High Energy Physics - Theory (**hep-th** new, recent, search)
- Mathematical Physics (**math-ph** new, recent, search)
- Nonlinear Sciences (**nlin** new, recent, search) includes: Adaptation and Self-Organizing Systems; Cellular Automata and Lattice Gases; Chaotic Dynamics; Exactly Solvable and Integrable Systems; Pattern Formation and Solitons
- Nuclear Experiment (**nucl-ex** new, recent, search)
- Nuclear Theory (**nucl-th** new, recent, search)
- Physics (**physics** new, recent, search) includes: Accelerator Physics; Applied Physics; Atmospheric and Oceanic Physics; Atomic and Molecular Clusters; Atomic Physics; Biological Physics; Chemical Physics; Classical Physics; Computational Physics; Data Analysis, Statistics and Probability; Fluid Dynamics; General Physics; Geophysics; History and Philosophy of Physics; Instrumentation and Detectors; Medical Physics; Optics; Physics and Society; Physics Education; Plasma Physics; Popular Physics; Space Physics
- Quantum Physics (**quant-ph** new, recent, search)

Mathematics

- Mathematics (**math** new, recent, search) includes: (see detailed description): Algebraic Geometry; Algebraic Topology;

National Library of Medicine
National Center for Biotechnology Information

PMC PubMed Central®

Search PMC Full-Text Archive Search in PMC

Advanced | Journal List

PubMed Central® (PMC) is a free full-text archive of biomedical and life sciences journal literature at the U.S. National Institutes of Health's National Library of Medicine (NIH/NLM)

- About PMC**
Discover a digital archive of scholarly articles, spanning centuries of scientific research.
- User Guide**
Learn how to find and read articles of interest to you.
- Collections**
Browse the PMC Journal List or learn about some of PMC's unique collections.
- For Authors**
Navigate the PMC submission methods to
- For Publishers**
Learn about deposit options for journals and
- For Developers**
Find tools for bulk download, text mining,

sci-hub

数据库 关于 Elbakyan 统计 捐助

未来我们将会加入浏览和搜索Sci-Hub数据库内容的选项，请持续关注更新

Sci-Hub数据库收集了 **88,343,822** 份研究文件，全部免费下载。大约80%的文件是发表在期刊上的研究论文 6%是会议论文集 (conference proceedings) 中的论文, 5%是书籍的章节 剩下的就是其他种类的文件了。Sci-Hub上能查询到的文件中, 77%发表于1980年到2020年 36%发表于2010年到2020年。所有主要科学出版社的覆盖率都在95%以上。Sci-Hub数据库的总规模约为100TB。

学科	文件数量
医学	24,557,530
化学	16,460,921
生物学	15,499,507
人文学科	12,592,316
物理学	8,658,518
工程学	6,892,853
数学	2,676,780
生态学	2,676,780
计算机科学	2,768,241
经济学	2,572,842
地理学	2,571,177

8.2.1 预训练数据源

■ 专业数据源

代码数据集 在程序生成方面，代码作为一种专业数据类型已经得到了广泛的研究关注。模型在大规模代码库上的预训练能显著提升代码生成质量。

The Stack

6 TB of permissive code data



@BigCodeProject
https://www.bigcode-project.org/
contact@bigcode-project.org

Dataset Collection

GH Archive query → 220 M repo names → git clone → 137 M repos, 52 B files, 102 TB of data → selecting file extensions → 49 TB of data → license filtering → 6.4 TB of data → near-deduplication → 2.9 TB of data

Licensing + Governance

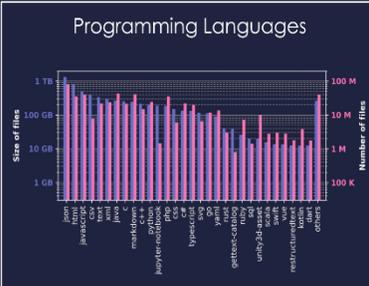
Raw dataset: No license, MIT, Apache 2.0, BSD-3-Clause, CC0-1.0, BSD-2-Clause, WTFPL, RSA-MD, MIT-0, Others

Permissive: MIT, Apache 2.0, BSD-3-Clause, CC0-1.0, BSD-2-Clause, WTFPL, RSA-MD, MIT-0, Others

Opt-out: If users would like to exclude their code from the corpus we have an opt-out mechanism. Visit: <https://www.bigcode-project.org/stack/about/the-stack/>

Permissive license distribution of licenses used to filter the dataset:
MIT (67.7%) | Apache-2.0 (19.1%) | BSD-3-Clause (3.9%) | Unlicense (2.0%) | CC0-1.0 (1.5%) | BSD-2-Clause (1.2%) | CC-BY-4.0 (1.1%) | CC-BY-3.0 (0.7%) | 0BSD (0.4%) | RSA-MD (0.3%) | WTFPL (0.2%) | MIT-0 (0.2%) | Others (166) (2.2%)

Programming Languages



Language	Size of files (TB)	Number of files
python	~1.0	~100M
java	~0.8	~80M
javascript	~0.7	~70M
cpp	~0.6	~60M
rust	~0.5	~50M
go	~0.4	~40M
php	~0.3	~30M
scala	~0.2	~20M
typescript	~0.1	~10M
others	~0.1	~10M

Evaluation

We trained several **GPT-2 models (350M parameters)** on different parts of the dataset both with and without near-deduplication. The models trained on the Python subset of The Stack performed on par with CodeX and CodeGen of similar size when using near-deduplication.

Dataset	Filtering	pass@1	pass@10	pass@100
CodeX (300M)	unknown	13.17	20.17	34.27
CodeGen (350M)	unknown	12.76	23.11	35.19
Python all-license	None	13.11	21.77	34.67
	Near-dedup	17.34	27.44	45.52
	Permissive	10.99	15.94	27.21
Python permissive-license	None	10.99	15.94	27.21
	Near-dedup	12.89	22.26	34.01

*results obtained with The Stack v1.0

Build and ship software on a single, collaborative platform

Join the world's most widely adopted AI-powered developer platform where millions of developers, businesses, and the largest open source community build software that advances humanity.

Sign up for GitHub | Try GitHub Copilot (30 days free)

Gitee 最有价值开源项目

Project	Stars	Language
kwdb	1.1K	Python
GSTVideoPlayer	1K	C++
SpireCV	157	C++

加入 GVP 计划 | 查看全部项目

AntV | apollo | 滴滴 | 开源 | DolphinScheduler | dotNET

stackoverflow

Every **developer** has a tab open to Stack Overflow.

For over 15 years we've been the Q&A platform of choice that millions of people visit every month to ask questions, learn, and share technical knowledge.

Sign up | Visit the community >



Services for companies of all shapes & sizes

Full business solutions

8.2.1 预训练数据源

■ 专业数据源

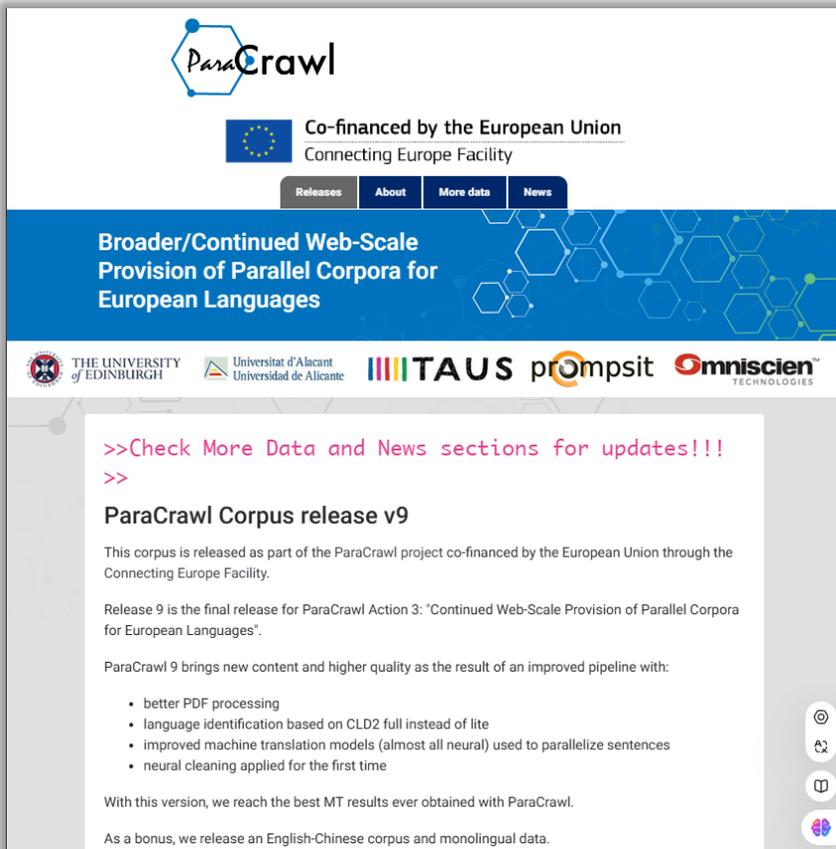
Python是代码数据集中最常见的程序语言。如 GPT-3 (175B) 的训练数据主要由 Common Crawl、BooksCorpus 和代码数据集构成，而 LaMDA (540B) 采用对话文本、代码数据和 Common Crawl。

数据集	数据源	程序语言	规模 (GB)	类型
CodeSearchNet (2019)	GitHub	Go, Java, JS, PHP, Python, Ruby PL	17	NL-PL
CodeNet (2021)	AIZU, AtCoder	-	8	NL-PL
THEPILE (2021)	GitHub, ArXiv,...	-	825	NL-PL
thestack	GitHub	-	3100	PL
BigQuery	GitHub	C/C++, Go, Java, JS, Python	340	PL
BIGPYTHON (2022)	GitHub	Python	217	PL
CodeParrot	GitHub	Python	180	PL
GCPY (2022)	GitHub	Python	-	PL

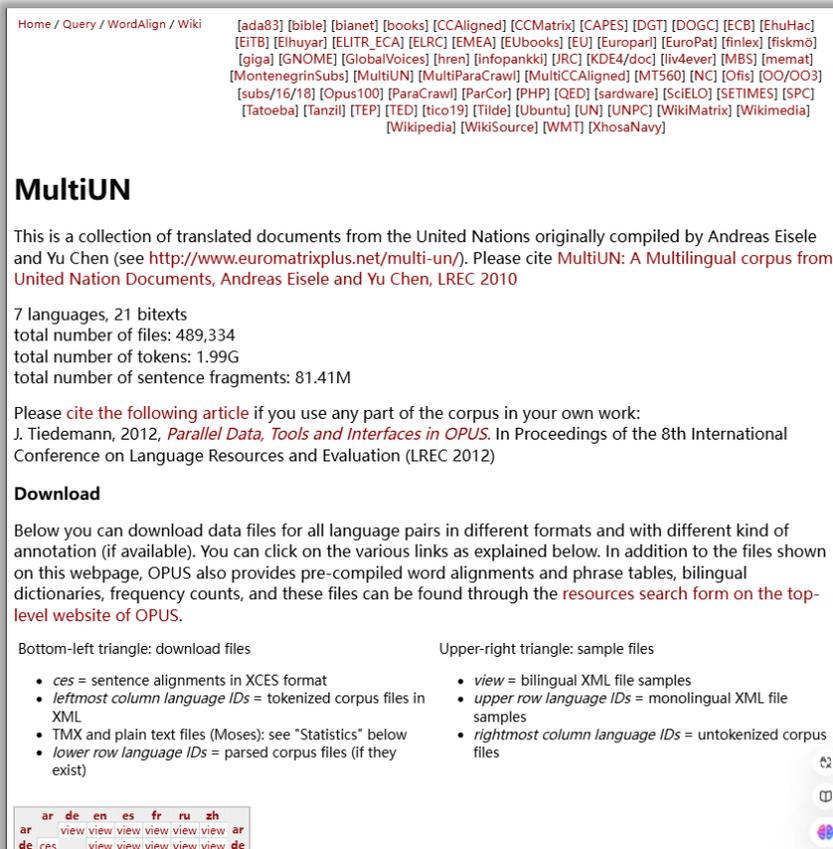
8.2.1 预训练数据源

■ 专业数据源

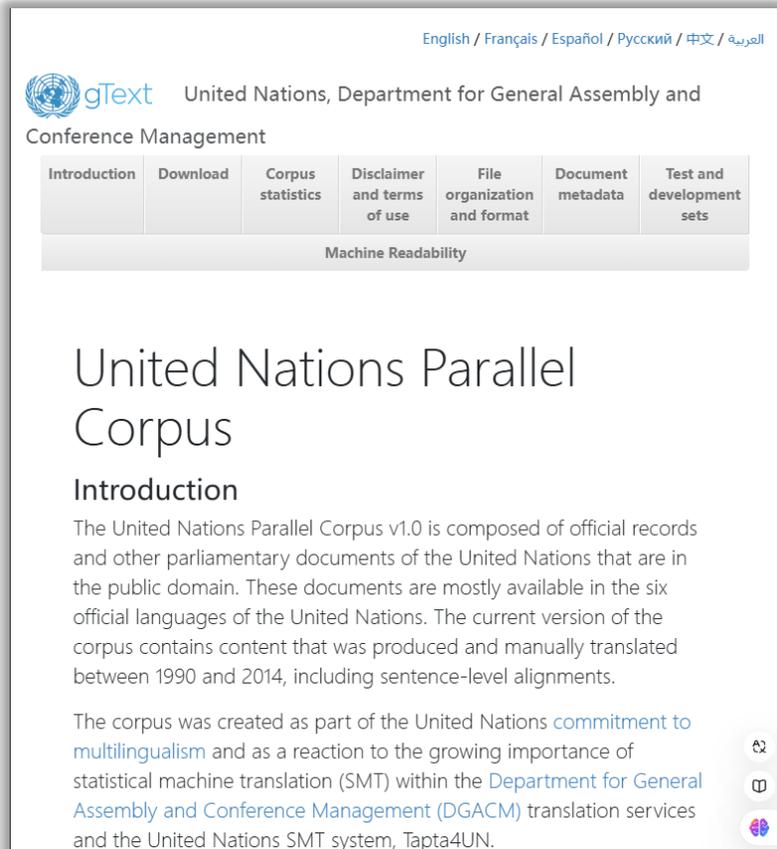
多语言语料 这类数据集通常包括多种语言的相似或相同的内容，例如双语或多语言新闻文章、翻译文本等，也被称为平行语料库。



The screenshot shows the ParaCrawl website. At the top, it features the ParaCrawl logo and a banner stating "Co-financed by the European Union Connecting Europe Facility". Below this, there are navigation tabs for "Releases", "About", "More data", and "News". The main heading reads "Broader/Continued Web-Scale Provision of Parallel Corpora for European Languages". Logos for partner institutions like The University of Edinburgh, Universitat d'Alicante, TAUS, procompit, and Omnicien are displayed. A pink text box says ">>Check More Data and News sections for updates!!!>>". The main content area is titled "ParaCrawl Corpus release v9" and describes the project's goals and improvements. A list of features includes better PDF processing, language identification, improved machine translation models, and neural cleaning. It also mentions that with this version, they reached the best MT results ever obtained with ParaCrawl. At the bottom, it notes that as a bonus, they release an English-Chinese corpus and monolingual data.



The screenshot shows the MultiUN website. At the top, there is a navigation bar with links for "Home / Query / WordAlign / Wiki" and a list of language codes. The main heading is "MultiUN". Below this, a paragraph describes the corpus as a collection of translated documents from the United Nations, compiled by Andreas Eisele and Yu Chen. It provides a citation: "MultiUN: A Multilingual corpus from United Nation Documents, Andreas Eisele and Yu Chen, LREC 2010". Statistics are listed: "7 languages, 21 bitexts", "total number of files: 489,334", "total number of tokens: 1.99G", and "total number of sentence fragments: 81.41M". A section titled "Please cite the following article if you use any part of the corpus in your own work:" lists a paper by J. Tiedemann from 2012. A "Download" section explains that users can download data files in various formats and provides instructions on how to find resources on the OPUS website. It also includes a "Bottom-left triangle: download files" section with a list of file types and formats, and an "Upper-right triangle: sample files" section with a list of file types and formats. At the bottom, there is a small table with language codes and links to view samples.

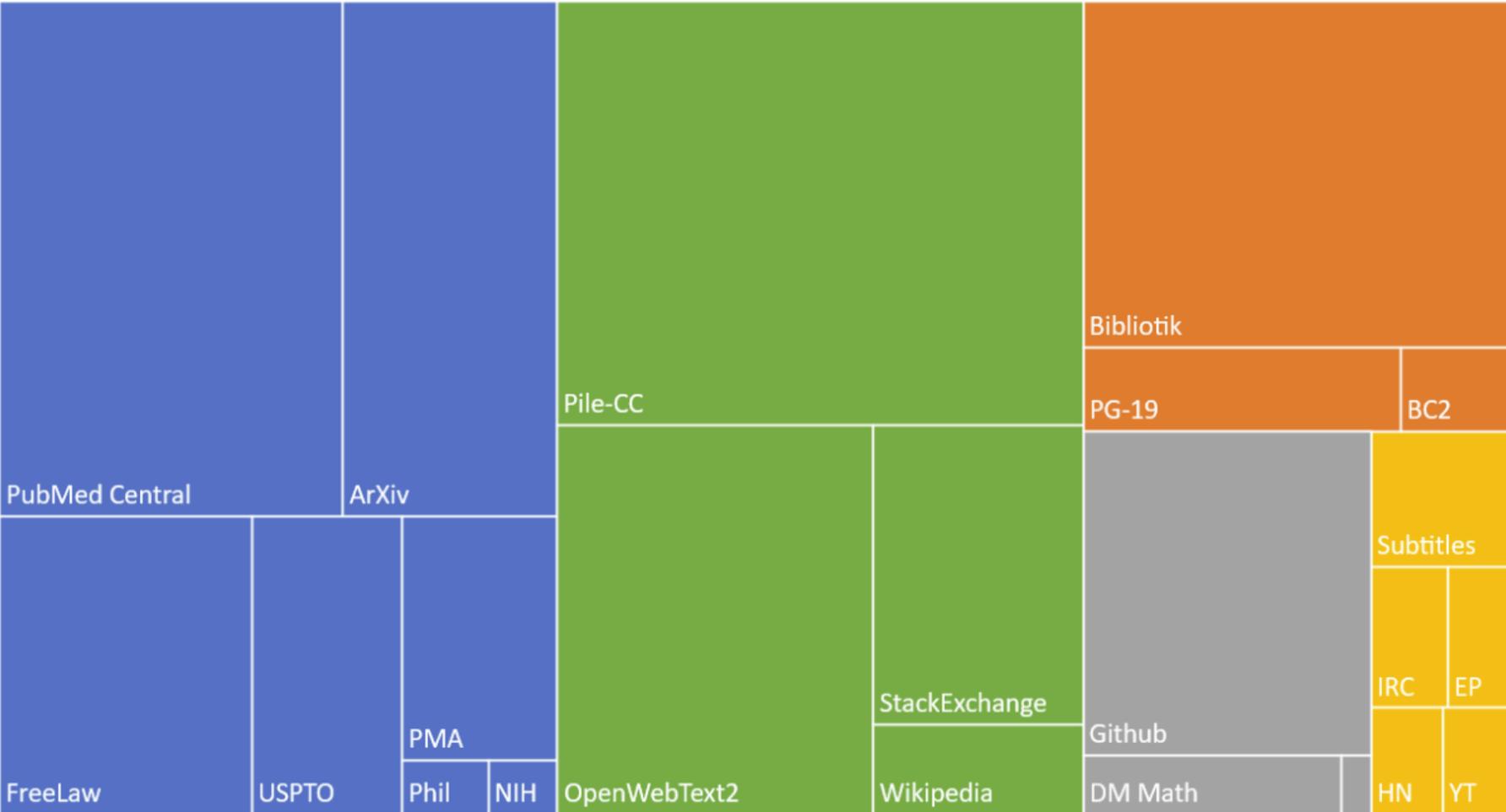


The screenshot shows the United Nations Parallel Corpus website. At the top, there is a navigation bar with links for "English / Français / Español / Русский / 中文 / العربية". The main heading is "United Nations Parallel Corpus". Below this, there is a section titled "Introduction" which describes the corpus as a collection of official records and other parliamentary documents of the United Nations. It mentions that the current version of the corpus contains content that was produced and manually translated between 1990 and 2014, including sentence-level alignments. At the bottom, it states that the corpus was created as part of the United Nations commitment to multilingualism and as a reaction to the growing importance of statistical machine translation (SMT) within the Department for General Assembly and Conference Management (DGACM) translation services and the United Nations SMT system, Tapta4UN.

8.2.1 预训练数据源

多类别语料

Pile: 2020年, 825G的语料 (英文数据集)

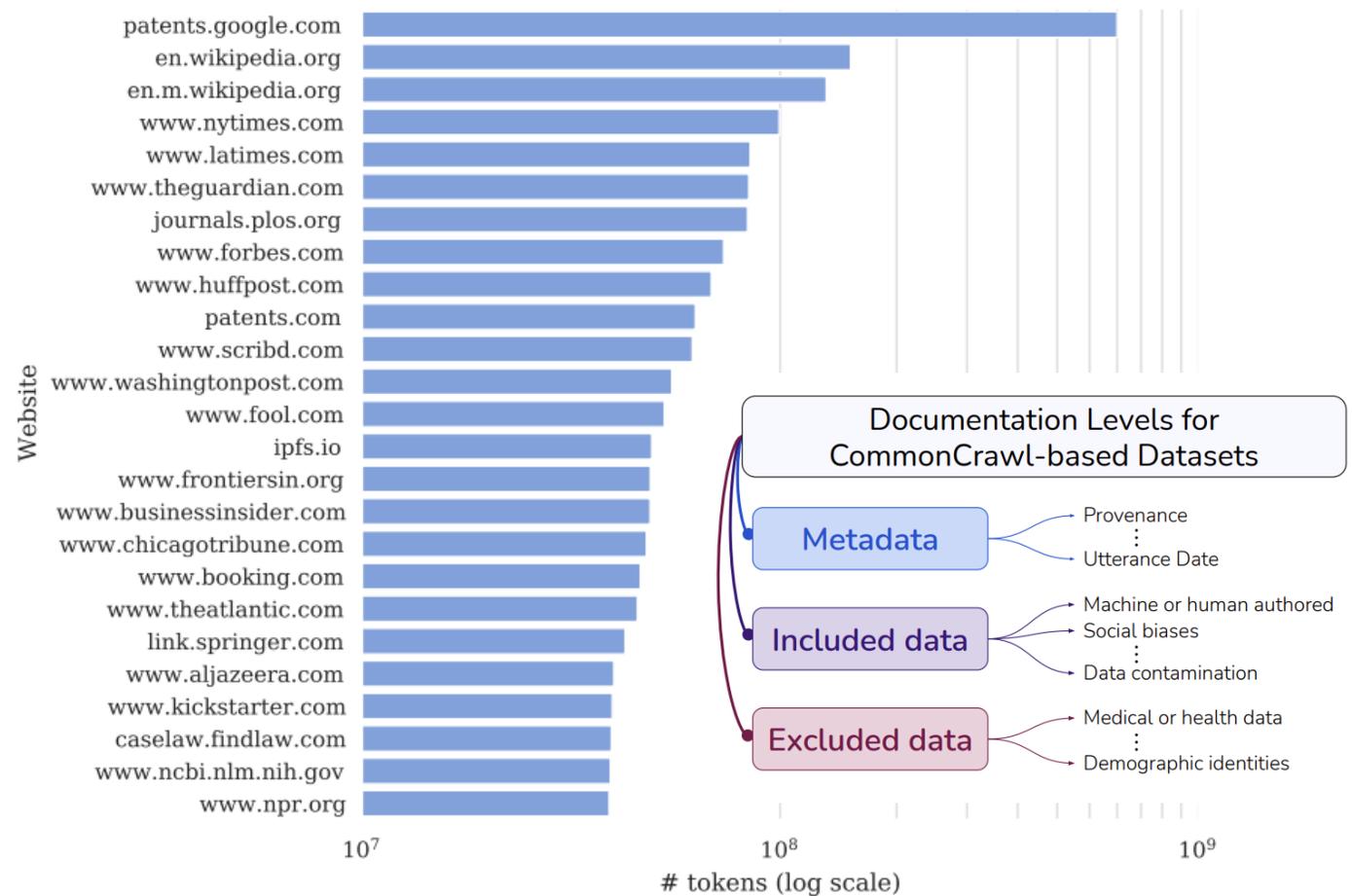
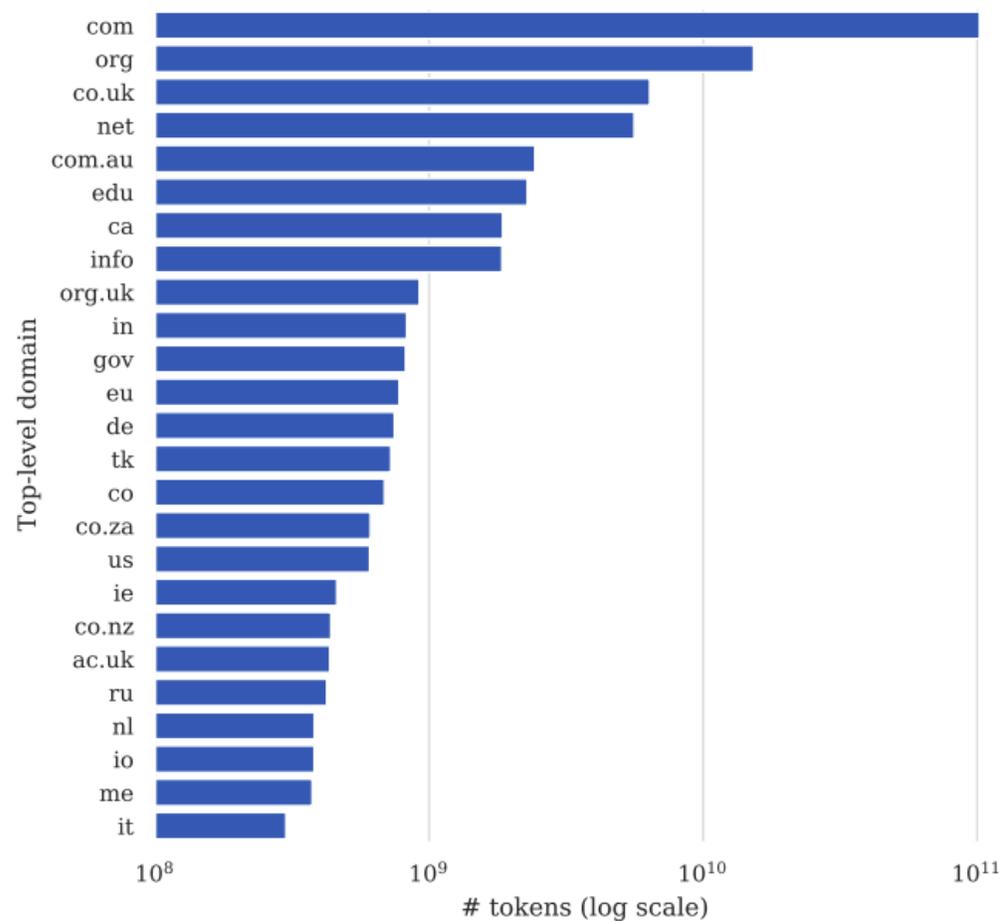


Component	Raw Size	Weight
Pile-CC	227.12 GiB	18.11%
PubMed Central	90.27 GiB	14.40%
Books3 [†]	100.96 GiB	12.07%
OpenWebText2	62.77 GiB	10.01%
ArXiv	56.21 GiB	8.96%
Github	95.16 GiB	7.59%
FreeLaw	51.15 GiB	6.12%
Stack Exchange	32.20 GiB	5.13%
USPTO Backgrounds	22.90 GiB	3.65%
PubMed Abstracts	19.26 GiB	3.07%
Gutenberg (PG-19) [†]	10.88 GiB	2.17%
OpenSubtitles [†]	12.98 GiB	1.55%
Wikipedia (en) [†]	6.38 GiB	1.53%
DM Mathematics [†]	7.75 GiB	1.24%
Ubuntu IRC	5.52 GiB	0.88%
BookCorpus2	6.30 GiB	0.75%
EuroParl [†]	4.59 GiB	0.73%
HackerNews	3.90 GiB	0.62%
YoutubeSubtitles	3.73 GiB	0.60%
PhilPapers	2.38 GiB	0.38%
NIH ExPorter	1.89 GiB	0.30%
Enron Emails [†]	0.88 GiB	0.14%
The Pile	825.18 GiB	

8.2.1 预训练数据源

多类别语料

C4, T5的训练语料, 2021 EMNLP (英文数据集)



8.2.1 预训练数据源

■ 多类别语料

WuDaoCorpora是北京智源研究院最新构建的高质量数据集，由全球最大的**纯文本数据集**、全球最大的**多模态图文数据集**和全球最大的**中文对话数据集**三部分构成。包含**3TB 文本数据**、**90TB 图文数据**和**181GB 对话数据**。

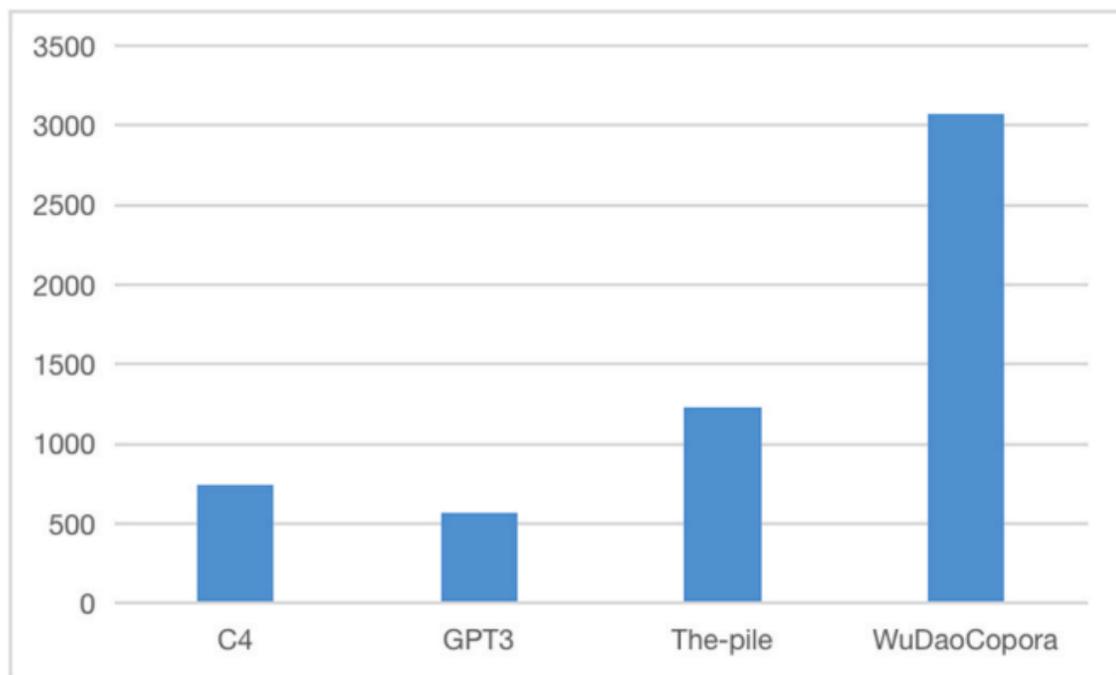


Fig. 1. Comparison with other Chinese corpora.

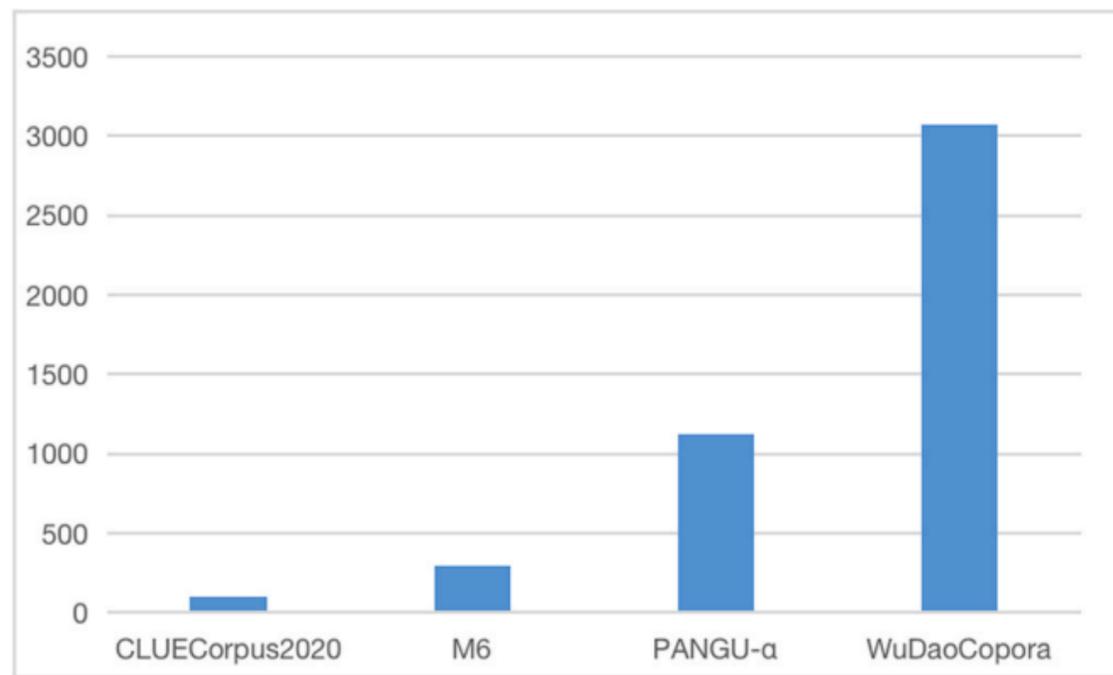


Fig. 2. Comparison with other language corpora.

8.2.1 预训练数据源

多类别语料

表 8.2: 多类别语料的主要数据分布

数据集	规模	网页	代码	书籍	科学文本
Pile	825 GB	36.2%	7.6%	15%	38.1%
TigerBot_pretrain_en	51 GB	33.9%	30.2%	35.9%	-
TigerBot_pretrain_zh	55 GB	50.3%	-	25.9%	-
WanJuanText	1 TB	96.8%	-	0.07%	2.1%

表 8.3: 垂直领域语料库

数据集	分类	规模	数据源
BBT-FinCorpus	金融	256 GB	公司公告、研究报告、金融新闻、社交媒体
FinCorpus	金融	60 GB	公司公告、金融新闻、金融考试题目
FinGLM	金融	69 GB	上市公司年报
Medical-pt	医疗	632 MB	医学百科、教科书
Proof-Pie-2	数学	55 B 词元	ArXiv、OpenWebMath
TigberBot-law	法律	29.9 MB	法律条款
TransGPT-pt	交通	35.8 MB	技术文档、工程施工信息

金融NLP数据集

SmoothNLP

CCKS

OpenKG

Minds14

FIQA

FPB

...

无标签金融文本

金融新闻

咨询函

金融研报

网络金融论坛

...

开源数据集

金融数据集

3.1 中文

- Huatu-26M
 - 地址: <https://github.com/FreedomIntelligence/Huatu-26M> Stars 218
 - 简介: Huatu-26M 是迄今为止最大的中医问答数据集。
- 中文医疗对话数据集
 - 地址: <https://github.com/Toyhom/Chinese-medical-dialogue-data> Stars 1.2k
 - 简介: 包含六个科室的医学问答数据
- CBLUE
 - 地址: <https://github.com/CBLUEbenchmark/CBLUE> Stars 727
 - 简介: 涵盖了医学文本信息抽取 (实体识别、关系抽取)
- cMedQA2 (108K)
 - 地址: <https://github.com/zhangsheng93/cMedQA2> Stars 309
 - 简介: 中文医药方面的问答数据集, 超过10万条
- xywy-KG(294K三元组)
 - 地址: <https://github.com/baiyang2464/chatbot-base-on-Knowledge-Graph> Stars
 - 简介: 44.1K实体 294.1K 三元组
- 39Health-KG (210K三元组)
 - 地址: <https://github.com/zhiaoh>
 - 简介: 包括15项信息, 其中7类实体

3.2 英文

- MedMentions
 - 地址: <https://github.com/chanzuckerberg/MedMentions> Stars 312
 - 简介: 基于PubMed摘要的生物医学实体链接数据集
- webMedQA
 - 地址: <https://github.com/hejunjing/webMedQA> Stars 60
 - 简介: 医疗问答
- COMETA
 - 地址: <https://www.siphs.org/>
 - 简介: 社交媒体中的医疗实体链接数据, 发表于EMNLP2020
- PubMedQA
 - 地址: <https://arxiv.org/abs/1909.06146>
 - 简介: 基于PubMed提取的医学问答数据集
- MediQA
 - 地址: <https://sites.google.com/view/mediqa2021>
 - 简介: 文本概括

中英文医疗数据集

Datasets: TigerResearch tigerbot-law-plugin like 25 Follow Tiger Research 65

Split (1)
train - 55.9k rows

Search this dataset

type	title	chapter1	content
string - cClasses	string - lengths	string - lengths	string - lengths
1.2 values	4	74	31.4k
7	0	0	0

type	title	chapter1	content
宪法	中华人民共和国宪法		1982年12月4日 第五届全国人民代表大会第五次会议通过 1982年12月4日 全国
宪法	中华人民共和国宪法	序言	中国是世界上历史最悠久的国家之一。中国各族人民共同创造了光辉灿烂的文化。
宪法	中华人民共和国宪法	第一章 总纲	第一条 中华人民共和国是工人阶级领导的、以工农联盟为基础的人民民主专政的社
宪法	中华人民共和国宪法	第一章 总纲	第二条 中华人民共和国的一切权力属于人民。人民行使国家权力的机关是全国人
宪法	中华人民共和国宪法	第一章 总纲	第三条 中华人民共和国的国家机构实行民主集中制的原则。全国人民代表大会和
宪法	中华人民共和国宪法	第一章 总纲	第四条 中华人民共和国各民族一律平等。国家保障少数民族合法的权利和利
宪法	中华人民共和国宪法	第一章 总纲	第五条 中华人民共和国实行依法治国，建设社会主义法治国家。国家维护社会主
宪法	中华人民共和国宪法	第一章 总纲	第六条 中华人民共和国的社会主义经济制度的基础是生产资料的社会主义公有制，
宪法	中华人民共和国宪法	第一章 总纲	第七条 国有经济，即社会主义全民所有制经济，是国民经济中的主导力量。国家
宪法	中华人民共和国宪法	第一章 总纲	第八条 农村集体经济组织实行家庭联产承包责任制，统分结合的双层经营体制。农
宪法	中华人民共和国宪法	第一章 总纲	第九条 矿藏、水流、森林、山岭、草原、荒地、滩涂等自然资源，都属于国家所有

法律大模型数据集

8.2.2 多模态数据集

多模态数据集是大规模数据集的重要分支，其独特之处在于能够同时包含并融合多种格式的数据资源，为模型预训练提供了更加多样化的输入

数据集	发布时间	规模	特点
Flickr [202]	2014	3 万张图片，每张 5 条描述	英文，人工标注
COCO [23]	2014	33 万张图片，每张 5 条描述	英文，2022 年发布了 7.47 亿张图片的 COCO-700M
Conceptual Caption ⁹	2018	30 万张图片，每张 5 条描述	英文，2021 年发布 1200 万张图片的 Conceptual12M
WIT [157]	2021	3700 万条图文对	多语言，来自维基百科
悟空 [58]	2022	1 亿条图文对	中文
LAION-5B [152]	2022	58 亿条图文对	多语言，提供 LAION2B-en（英文描述）等子集
WuDaoMM ¹⁰	2022	6.5 亿条图文对	中文

一些代表性的图文语料库

多模态数据集能促进大模型预训练，让模型具备更强的泛化能力

8.2.2 多模态数据集

多模态数据集被广泛用于多模态匹配（检索）任务。该任务是指通过处理和比较来自不同模态的数据，实现高效且精确的跨模态信息检索。

M5PRODUCT M5产品

About 大约

The M5Product dataset is a large-scale multi-modal pre-training dataset with coarse and fine-grained annotations for E-products. M5Product 数据集是一个大规模的多模态预训练数据集，具有对电子产品的粗粒度和细粒度注释。

- 6 Million multi-modal samples, 5k properties with 24 Million values
- 600 万个多模态样本，5k 个属性，2400 万个值
- 5 modalities - image text table video audio
- 5 种模态 - 图像文本表视频音频
- 6 Million category annotations with 6k classes
- 600 万个类别注释，具有 6k 个类
- Wide data source (1 Million merchants provide)
- 广泛的数据源 (100 万商家提供)

Sampler 采样

The data acquisition page is shown as follows. 数据获取页面如下所示。



Dataset 数据

TVQA is a large-scale video QA dataset based on 6 popular TV shows (*Friends*, *The Big Bang Theory*, *How I Met Your Mother*, *House M.D.*, *Grey's Anatomy*, *Castle*). It consists of 152.5K QA pairs from 21.8K video clips, spanning over 460 hours of video. The questions are designed to be compositional, requiring systems to jointly localize relevant moments within a clip, comprehend subtitles-based dialogue, and recognize relevant visual concepts. Download TVQA data from [/data](#).

TVQA 是基于 6 个热门电视节目的大规模视频 QA 数据集（《老友记》、《生活大爆炸》、《我是如何遇见玛妈的》、《医学博士》、《实习医生格蕾》、《城堡》）。它由来自 21.8K 视频剪辑的 152.5K QA 对组成，视频时长超过 460 小时。这些问题被设计为组合式的，需要系统共同定位剪辑中的相关时刻，理解基于字幕的对话，并识别相关的视觉概念。从 [/data](#) 下载 TVQA 数据。

- QA example QA 示例



See examples in video: [click here](#)
观看视频中的示例: [单击此处](#)

HowTo100M

video clips run in real time, this demo uses a lighter (and less accurate) version of the model than the one described in the paper.

HowTo100M 搜索演示 | 统计 | 下载 | 团队

Dataset statistics



HowTo100M
23611 tasks
136.6M clips

OpenDataLab

YFCC100M

OpenDataLab/YFCC100M

Video Image CV

Large Model Pre-training Image Classification

Image Retrieval MIT

1.3KB 9.4k 2

OpenDataLab 2024.08.23

Introduction Download

数据集介绍

简介

YFCC100M 是一个包含总共 1 亿个媒体对象的数据集，其中大约 9920 万个照片和 80 万个视频，所有这些都带有知识共享许可证。数据集集中的每个媒体对象都由几条元数据表示，例如 Flickr 标识符、所有者名称、相机、标题、标签、地理位置、媒体来源。该集合提供了从 Flickr 于 2004 年成立以来到 2014 年初这些年来如何拍摄、描述和共享照片和视频的全面快照。

类定义

8.2.3 数据处理

数据处理的主要目标是剔除数据集中的**噪声**、**冗余信息**、**无关数据**，以及**潜在有害内容**。不合适的数据可能对语言模型的训练效果产生不利影响。

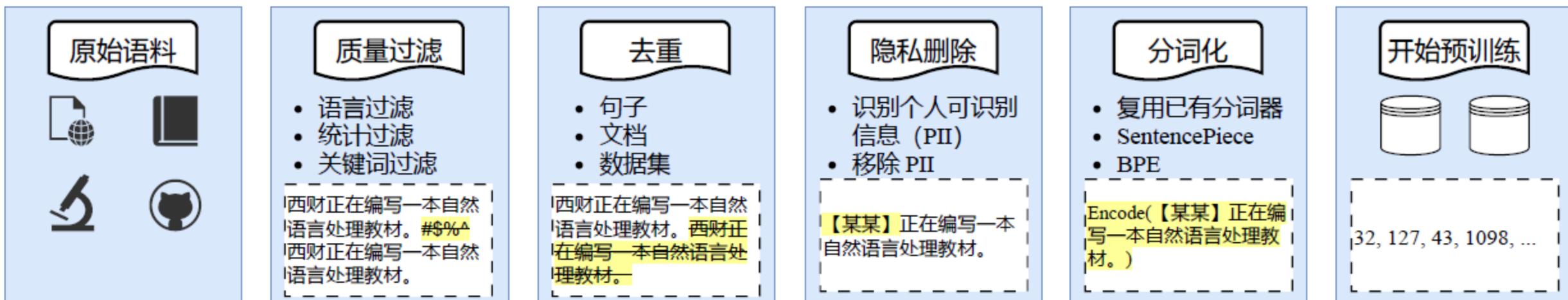


图 8.2: 典型的大型语言模型预训练的数据预处理流程的示意图

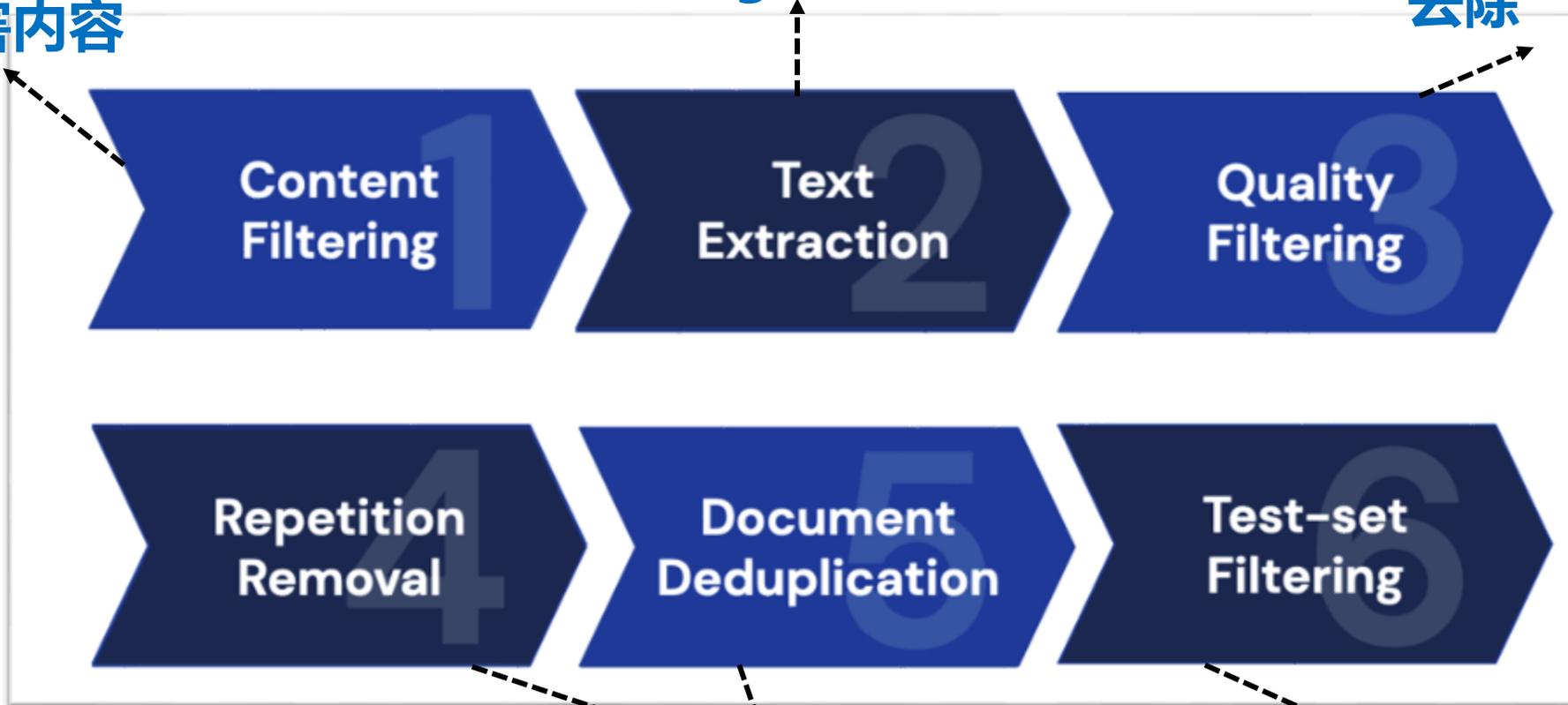
主要涵盖**质量过滤 (Quality Filtering)**、**去重 (De-duplication)**、**隐私删除 (Privacy Redaction)** 及**分词化 (Tokenization)** 四个步骤。

Scaling Language Models: Methods, Analysis & Insights from Training *Gopher*

去掉HTML Tag, 而保留内容

过滤有害内容

去除“低质”文本



去除重复文本

为了实验的严谨

8.2.3 数据处理 — 质量过滤

■ 质量过滤常用方法

分类器法 训练一个小型分类器判断文本质量

规则判断 人为设置规则过滤低质量数据

指标阈值 设置特定指标选择阈值范围内数据

聚类方法 聚类过滤，筛选指定类别的数据

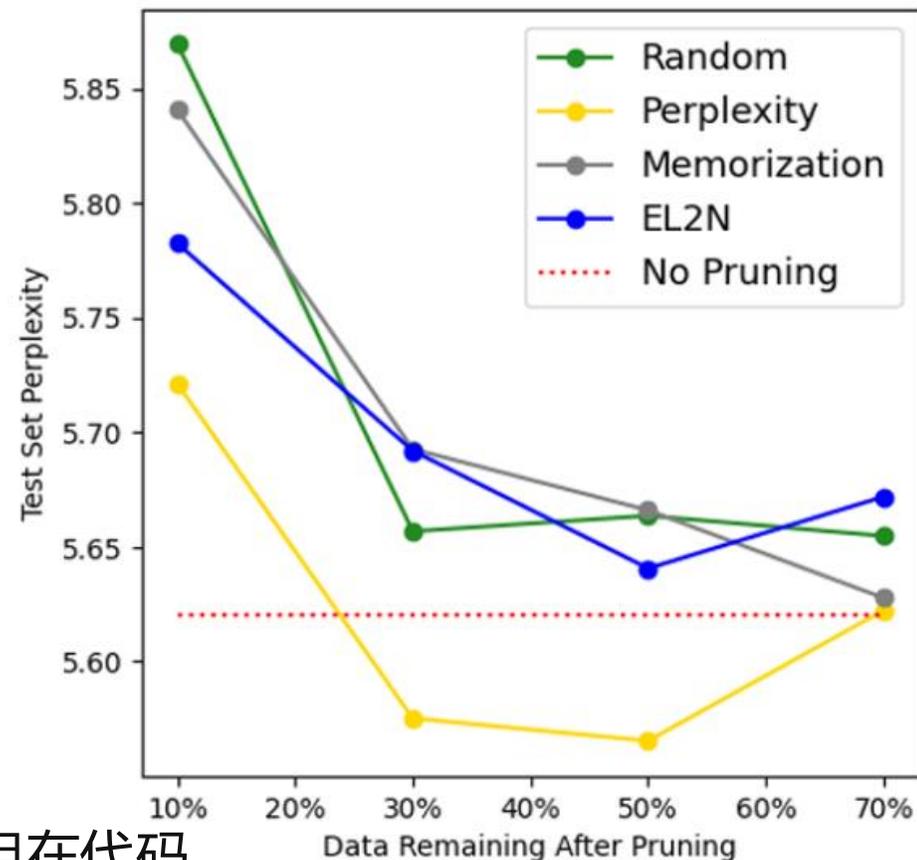
■ 使用小规模高质量样本训练轻量级模型

□ Microsoft phi-1/phi-1.5/phi-2仅有1.3B-2.7B参数量，但在代码和推理任务上表现良好

■ 过于激进的过滤可能会导致模型性能下降[1]

□ 选定的过滤指标不能很好地衡量数据的质量

■ Perplexity作为简单的衡量指标效果最好[2]



[1] Leo Gao. 2021. An empirical exploration in quality filtering of text data. arXiv preprint arXiv:2109.00698.

[2] Max Marion et al. 2023. When less is more: Investigating data pruning for pretraining llms at scale. arXiv preprint arXiv:2309.04564.

Perplexity 困惑度

A better model is better at predicting upcoming words, and so it will be less surprised by (i.e., assign a higher probability to) each word when it occurs in the test set.

$$\begin{aligned}\text{perplexity}(W) &= P(w_1 w_2 \dots w_N)^{-\frac{1}{N}} \\ &= \sqrt[N]{\frac{1}{P(w_1 w_2 \dots w_N)}}\end{aligned}$$

困惑度越小表明模型能力越强

$$\text{perplexity}(W) = \sqrt[N]{\prod_{i=1}^N \frac{1}{P(w_i | w_1 \dots w_{i-1})}}$$

8.2.3 数据处理 — 数据去重

■ 数据去重的作用

- 提高训练效率[1]
- 减少训练集和测试集的重叠[1]
- 减轻模型的记忆现象 - 降低隐私攻击的成功率[2]

¹“by combining fantastic ideas, interesting arrangements, and follow the current trends in the field of that make you more inspired and give artistic touches. We’d be honored if you can apply some or all of these design in your wedding. believe me, brilliant ideas would be perfect if it can be applied in real and make the people around you amazed!”

**这段话在语料库中重复了
61,036次!**

[1] Katherine Lee et al., 2022. Deduplicating training data makes language models better. ACL.

[2] Nikhil Kandpal et al., 2022. Deduplicating training data mitigates privacy risks in language models. ICML.

[3] Amro Abbas et al., 2023. Semdedup: Data-efficient learning at web-scale through semantic deduplication. arXiv preprint arXiv:2303.09540.

8.2.3 数据处理 — 数据去重

■ 数据去重的常用方法

- N-gram-and-hashing 是最为常用的方法
 - Line-level / document-level / both
- 神经网络模型
- 语义聚类方法：去除语义重复数据[3]

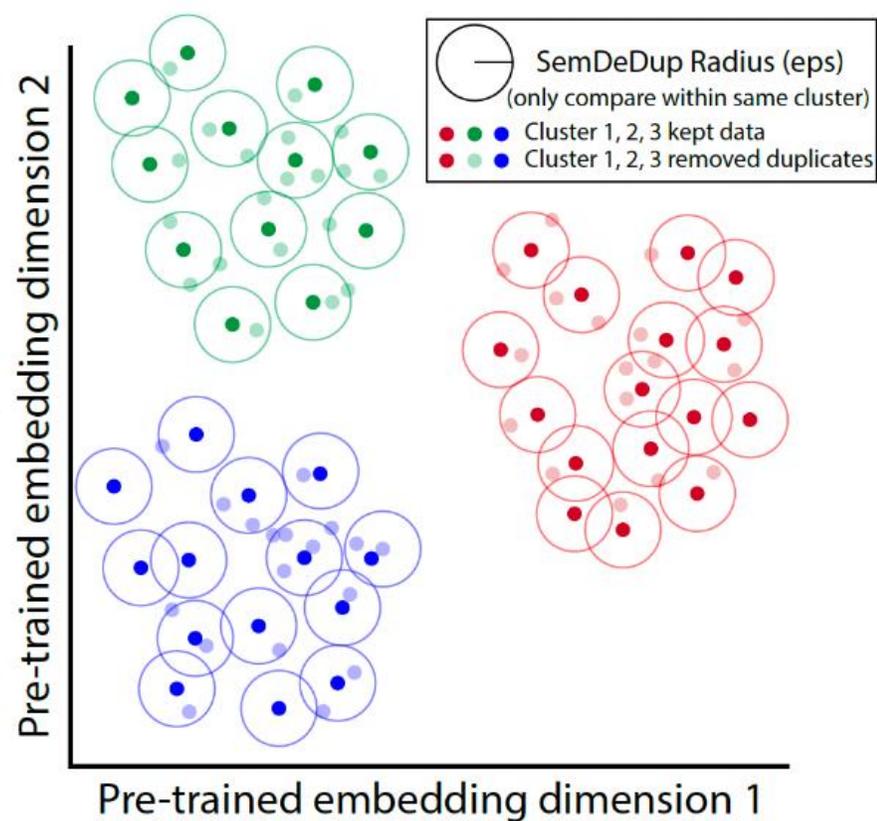


MinHash 文A 4 languages

Article [Talk](#) [Read](#) [Edit](#) [View history](#) [Tools](#)

From Wikipedia, the free encyclopedia

In [computer science](#) and [data mining](#), **MinHash** (or the **min-wise independent permutations locality sensitive hashing scheme**) is a technique for quickly estimating how [similar](#) two sets are. The scheme was published by [Andrei Broder](#) in a 1997 conference,^[1] and initially used in the [AltaVista](#) search engine to detect duplicate web pages and eliminate them from search results.^[2] It has also been applied in large-scale [clustering](#) problems, such as [clustering documents](#) by the similarity of their sets of words.^[1]



[1] Katherine Lee et al., 2022. Deduplicating training data makes language models better. ACL.

[2] Nikhil Kandpal et al., 2022. Deduplicating training data mitigates privacy risks in language models. ICML.

[3] Amro Abbas et al., 2023. Semdedup: Data-efficient learning at web-scale through semantic deduplication. arXiv preprint arXiv:2303.09540.

8.2.3 数据处理 — 隐私删除

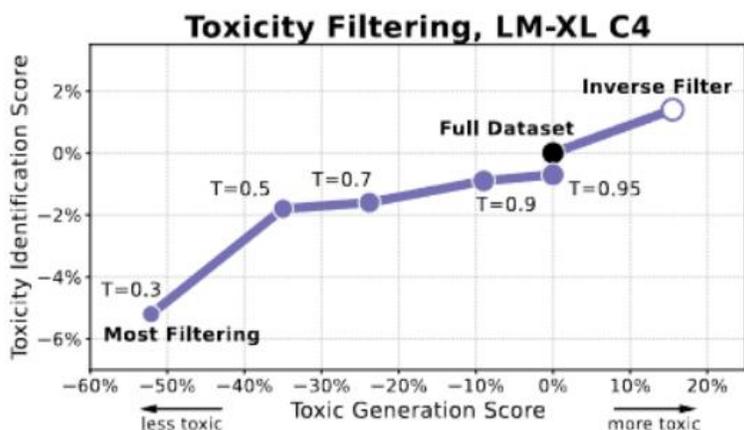
常用的过滤删除方法：

分类器法 训练一个小型分类器判断文本质量

规则判断 人为设置规则过滤低质量数据

隐私删除的类似任务 — 质量过滤

质量过滤会增强模型的泛化性和对有害信息判别的能力 [1]



	Wiki	Web	Books	Biomed	Academic	Common Sense	Contrast Sets	Average
Inverse T=0.5 (73%)	-5.0	-4.5	2.1	-2.2	-2.7	1.2	-6.4	-3.1
Full Dataset (100%)	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
T=0.975 (91%)	1.2	0.7	-2.2	6.1	6.4	4.7	6.1	2.5
T=0.95 (84%)	-1.2	1.0	-4.0	3.7	-0.3	3.2	4.9	1.0
T=0.9 (73%)	-0.3	0.8	-3.5	1.8	1.0	1.9	6.8	1.2
T=0.7 (46%)	-1.2	0.8	-6.7	1.7	0.8	2.0	4.2	0.7

思考：隐私删除是否会增强模型泛化性和对是否是隐私信息判别的能力？

数据信息的过滤和删除可能会导致少数群体被边缘化的现象 [2][3]

[1] S. Longpreet et al. 2021A Pretrainer's Guide to Training Data: Measuring the Effects of Data Age, Domain Coverage, Quality, & Toxicity. 2023.

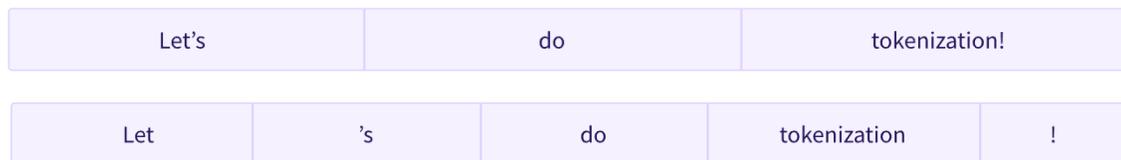
[2] Suchin Gururangan et al. 2022. Whose language counts as high quality? measuring language ideologies in text data selection. EMNLP 2021.

[3] Shangbin Feng et al. 2023. From pretraining data to language models to downstream tasks: Tracking the trails of political biases leading to unfair nlp models. ACL 2023.

8.2.3 数据处理

分词化

基于单词的分词器 (Word-based Tokenizer)



基于字符的分词器 (Character-based Tokenizer)



基于子词的分词器 (Subword-based Tokenizer)



“unfortunately” = “un” + “for” + “tun” + “ate” + “ly”

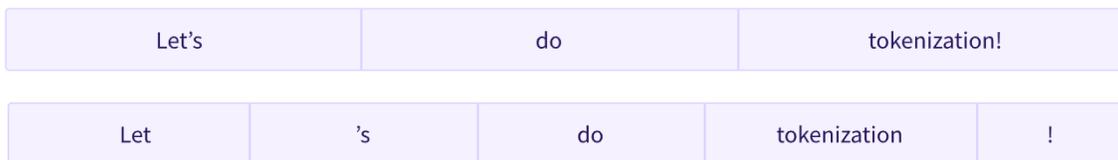
子词分词算法依赖于这样一个原则，即不应将常用词拆分为更小的子词，而应将稀有词分解为有意义的子词。

例如，“annoyingly”可能被认为是一个罕见的词，可以分解为“annoying”和“ly”。这两者都可能作为独立的子词出现得更频繁，同时“annoyingly”的含义由“annoying”和“ly”的复合含义保持。

8.2.3 数据处理

分词化

基于单词的分词器 (Word-based Tokenizer)



基于子词的分词器 (Subword-based Tokenizer)



“*unfortunately*” = “*un*” + “*for*” + “*tun*” + “*ate*” + “*ly*”

WordPiece Tokenizer

原始文本: I have a new GPU.

标记文本: ['i', 'have', 'a', 'new', 'gp', '## u', '.']

基于字符的分词器 (Character-based Tokenizer)



字节对编码标记化(Byte-level BPE)

Byte Pair Encoding Data Compression Example

aaabdaaabc

aaabdaaabc Replace Z = aa

ZabdZabc Replace Y = ab

ZYdZYac Replace X = ZY

XdXac Final compressed string

Replacement Table

Byte pair	Replacement
X	ZY
ab	Y
aa	Z

Unigram Tokenizer (SentencePiece or Unigram)

8.2.3 数据处理

■ BPE

```
import tiktoken

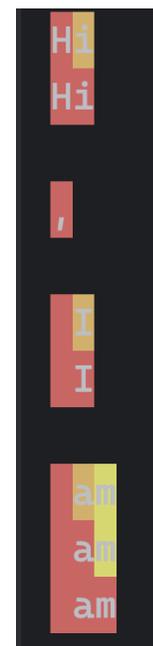
enc = tiktoken.encoding_for_model("gpt-4o")

print(enc.encode("你好, 我正在学习NLP"))
```

[177519, 40824, 70104, 64550, 45, 19318]

```
from tiktoken._educational import *

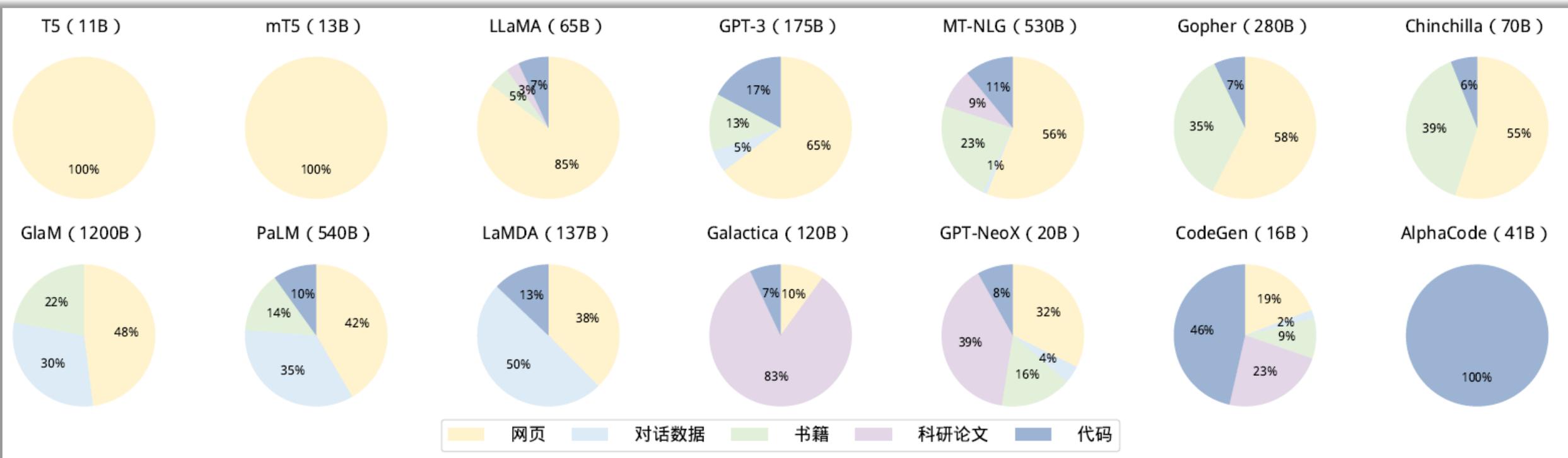
enc = SimpleBytePairEncoding.from_tiktoken("cl100k_base")
tokens = enc.encode("Hi, I am learning NLP")
```



8.2.4 模型性能关系

■ 数据源的多样性

数据分布对模型性能具有影响。Gopher 团队进行了一系列消融实验 [138]，发现增加**书籍的数据比例**能够提高模型捕捉文本**中长期依赖关系**的能力 [127]。当某一**特定领域**的训练数据**过多时**，也可能会**降低** LLMs 在其他领域的**泛化能力** [164, 138]。



8.2.4 模型性能关系

■ 预训练数据量 — Scaling Laws

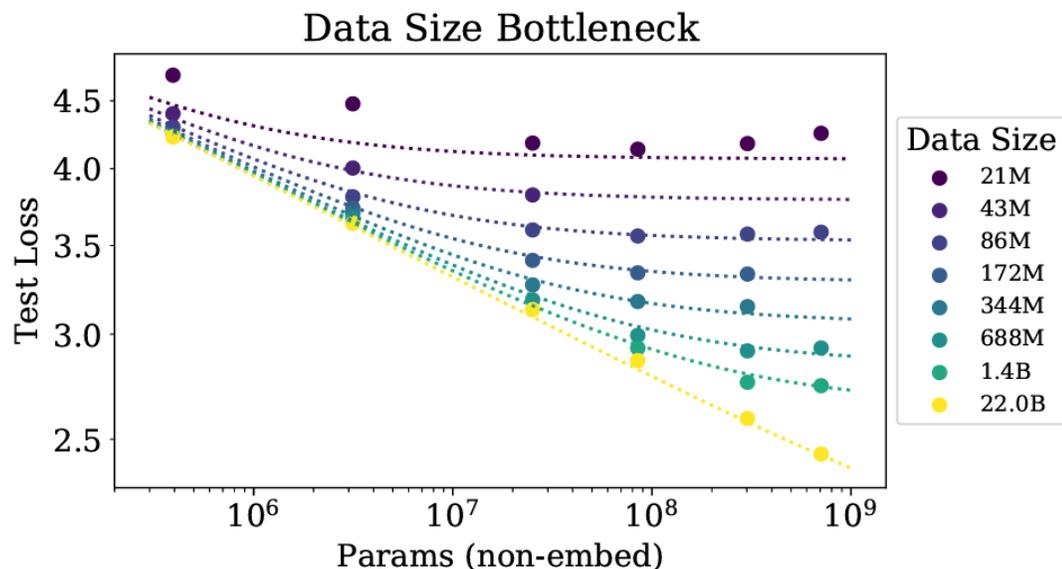
OpenAI Kaplan's Scaling Laws

- 在给定训练计算量的情况下，模型性能和数据数量符合幂律关系

$$L(N) = (N_c/N)^{\alpha_N}; \quad \alpha_N \sim 0.076, \quad N_c \sim 8.8 \times 10^{13} \text{ (non-embedding parameters)}$$

$$L(D) = (D_c/D)^{\alpha_D}; \quad \alpha_D \sim 0.095, \quad D_c \sim 5.4 \times 10^{13} \text{ (tokens)}$$

- 随着训练计算量的增加，数据数量和模型大小应该同步增长，但是模型大小的增速应该更快

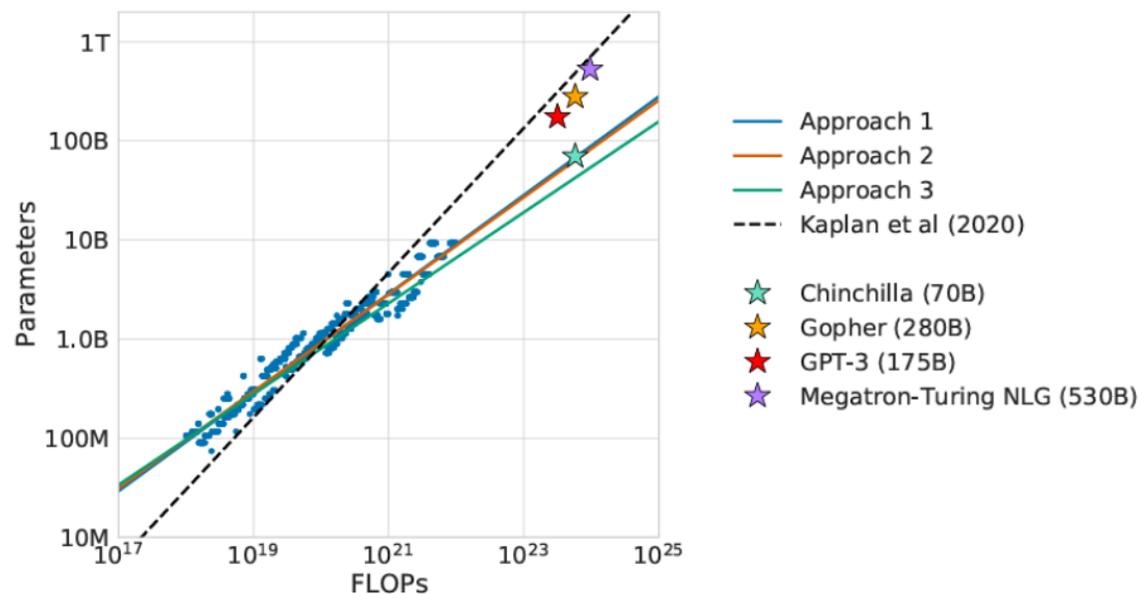


DeepMind Chinchilla Scaling Laws

- 新的 scaling law

$$\hat{L}(N, D) \triangleq E + \frac{A}{N^\alpha} + \frac{B}{D^\beta}$$

- 随着训练计算量增长，数据数量和模型大小应该以几乎一致的速度增长 $N_{opt} \propto C^a$ and $D_{opt} \propto C^b$



8.2.4 模型性能关系

■ 预训练数据质量

高质量的预训练数据是优化 LLMs 性能的关键要素之一

■ 数据低质量表现形式

噪声数据 对模型训练和性能产生负面影响

有毒内容 生成毒性文本的能力下降，而识别毒性的能力提高

重复内容 导致模型性能在初始阶段下降，影响模型在上下文信息中进行有效信息抽取的能力

思考1 什么样的数据集才是高质量的预训练数据？

数据源1: <https://blog.csdn.net/u011559552/article/details/142217358>

数据源2: <https://github.com/lonePatient/awesome-pretrained-chinese-nlp-models>

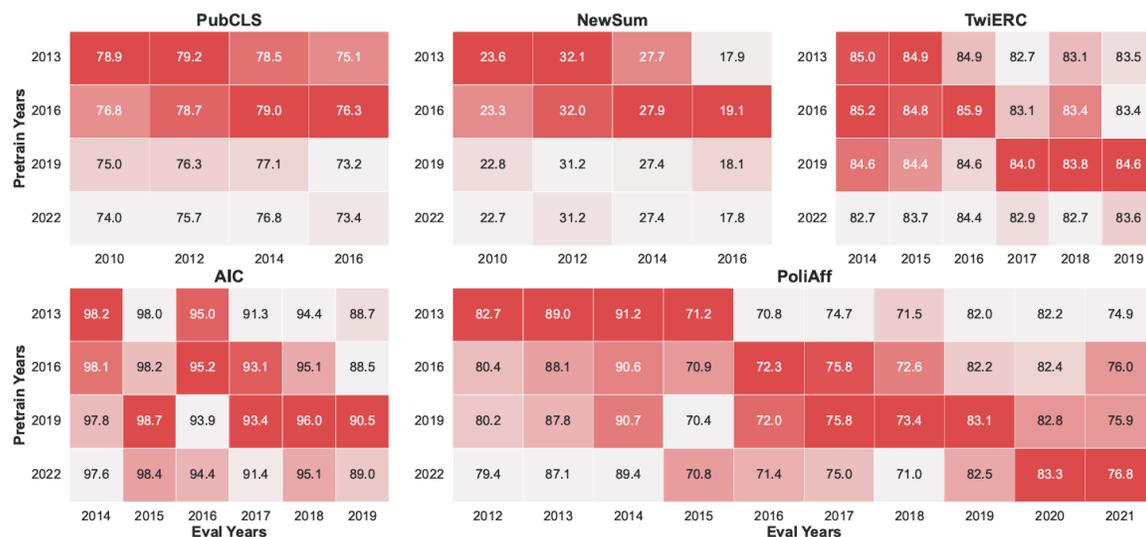
数据源3: <https://gitee.com/oschina/awesome-llm>

思考2 还有什么因素会影响大模型预训练性能？

8.2.4 模型性能关系

■ 预训练数据质量 — 数据年龄（时效性）

- 预训练数据集和评估数据集之间的时间不对齐会导致模型性能下降
- 时间不对齐问题不容易用微调来解决
- 时间不对齐导致的性能下降在更大的模型上更显著



DOMAIN	TASK	FINETUNING				PRETRAINING			
		LM-SMALL		LM-XL		LM-SMALL		LM-XL	
		TD	<i>r</i>	TD	<i>r</i>	TD	<i>r</i>	TD	<i>r</i>
NEWS	PUBCLS	5.82	0.84	5.63	0.80	0.02	0.01 [†]	0.59	0.67
	NEWSUM	0.80	0.82	2.91	0.92	-0.31	-0.29	0.73	0.45
TWITTER	POLI AFF	3.74	0.84	4.93	0.89	0.50	0.21	0.28	0.56
	TwiERC	0.49	0.73	0.53	0.82	0.05	0.27	0.23	0.72
SCIENCE	AIC	0.94	0.83	0.24	0.36	0.11	0.18 [†]	0.23	0.66
	MEAN	2.36	0.81	2.84	0.76	0.08	0.07	0.41	0.61

综上所述，预训练数据的来源、数量、质量和时效都是影响 LLMs 性能的重要因素。需要进行细致的优化，从而提升 LLMs 在泛化场景性能

本章内容

- 8.1 概述
- 8.2 预训练数据工程
 - 8.2.1 预训练数据源
 - 8.2.2 多模态数据集
 - 8.2.3 数据处理
 - 8.2.4 模型性能关系
- 8.3 预训练方法
 - 8.3.1 预训练任务
 - 8.3.2 优化参数设置
 - 8.3.3 可扩展训练技术
- 8.4 讨论

8.3.1 预训练任务



8.3.1 预训练任务



8.3.1 预训练任务



所以这就是预训练阶段的作用。但是

8.3.1 预训练任务

■ 语言建模

Next-token prediction

$$\mathcal{L}_{LM}(\mathbf{x}) = \sum_{i=1}^n \log P(x_i | x_{<i})$$

这样的思想事实上是在传达只要模型足够大，学到的知识足够多，任何有监督任务都可以通过无监督的方式来完成，即任何任务都可以视作生成任务。

8.3.1 预训练任务

■ 去噪自编码

通过将输入数据添加噪声，然后训练模型还原出原始的、无噪声的输入数据，去噪自编码器能够学习到数据的鲁棒表示。

$$\mathcal{L}_{DAE}(\mathbf{x}) = \log P(\tilde{\mathbf{x}} \mid \mathbf{x} \setminus \tilde{\mathbf{x}})$$

假设输入是 $\mathbf{x} \setminus \tilde{\mathbf{x}}$ ，其中 $\tilde{\mathbf{x}}$ 是带有随机替换片段的损坏文本，语言模型被训练用来恢复被替换的词元 $\tilde{\mathbf{x}}$ 。DAE 的训练目标是使重建的输出 $\tilde{\mathbf{x}}$ 尽可能接近原始的无噪声输入 \mathbf{x} 。形式上，DAE 的训练目标表示如下：

然而，相比于 LM 任务，DAE 任务在实现上似乎更加复杂 [214]。采用 DAE 作为预训练目标的 LLMs 包括 T5 和 GLM-130B。

8.3.2 优化参数设置

批量训练

- 影响**训练速度、收敛性、泛化能力和资源使用效率**。大批量可**减少梯度噪声**，使训练过程更加**稳定，加快收敛**
- 在训练过程中**动态增加批量大小**，能够有效地**提升训练稳定性**

学习率

- 为平衡**训练速度和稳定性**，现有的大模型通常在预训练期间采用**动态的学习率调度策略**，包括**预热 (Warm-up)** 和**衰减 (Learning Rate Decay)** 策略

优化器

- SGD**计算效率高、收敛速度较慢**
- RMSprop自适应调整，**收敛快、稳定性高**
- Adam 通过计算梯度的一/二阶矩估计来自动调整每个参数的学习率
- AdamW (带权重衰减 Adam) **收敛快、对超参数不敏感、对复杂损失函数适应性好**

模型	批次大小	学习率 (预热 → 峰值 → 衰减)	优化器	精度 类型	权重 衰减	梯度 裁剪
GPT-3	32K → 3.2M	预热 → 6×10^{-5} → 余弦	Adam	FP16	0.1	1.0
PanGu- α	-	2×10^{-5}	Adam	-	0.1	-
OPT	2M	预热 → 1.2×10^{-4} → 手动	AdamW	FP16	0.1	-
PaLM	1M → 4M	1×10^{-2} → 平方根倒数	AdaFactor	BF16	lr^2	1.0
BLOOM	4M	预热 → 6×10^{-5} → 余弦	Adam	BF16	0.1	1.0
MT-NLG	64K → 3.75M	预热 → 5×10^{-5} → 余弦	Adam	BF16	0.1	1.0
Gopher	3M → 6M	预热 → 4×10^{-5} → 余弦	Adam	BF16	-	1.0
Chinchilla	1.5M → 3M	预热 → 1×10^{-4} → 余弦	AdamW	BF16	-	-
Galactica	2M	预热 → 7×10^{-6} → 余弦	AdamW	-	0.1	1.0
LaMDA	256K	-	-	BF16	-	-
Jurassic-1	32k → 3.2M	预热 → 6×10^{-5}	-	-	-	-
LLaMA-2	4M	预热 → 1.5×10^{-4} → 余弦	AdamW	-	0.1	1.0
Pythia	2M	预热 → 1.4×10^{-4} → 余弦	Adam	FP16	0.01	1.0
Baichuan-2	-	预热 → 1.5×10^{-4} → 余弦	AdamW	BF16	0.1	0.5
Qwen-1.5	4M	预热 → 3×10^{-4} → 余弦	AdamW	BF16	0.1	1.0
InternLM-2	5M	预热 → 3×10^{-4} → 余弦	AdamW	-	0.1	-
Falcon	预热 → 4M	预热 → 1.25×10^{-4} → 余弦	AdamW	BF16	0.1	0.4
DeepSeek	18M	预热 → 3.2×10^{-4} → 余弦	AdamW	BF16	0.1	1.0
Yi	256K	1×10^{-5}	AdamW	BF16	0.1	1.0
YuLan	4.5M	预热 → 3×10^{-4} → 余弦	Adam	BF16	0.1	1.0
GLM-130B	0.4M → 8.25M	预热 → 8×10^{-5} → 余弦	AdamW	FP16	0.1	1.0
T5	64K	1×10^{-2} → 平方根倒数	AdaFactor	-	-	-

8.3.2 优化参数设置

■ 稳定训练

权重衰减 (Weight Decay)

通过在损失函数中添加与模型权重平方成正比的项，减少模型复杂度并防止过拟合。

鼓励模型在训练过程中保持较小的权重值，有助于提高模型的泛化能力，并增强训练稳定性。

梯度裁剪 (Gradient Clipping)

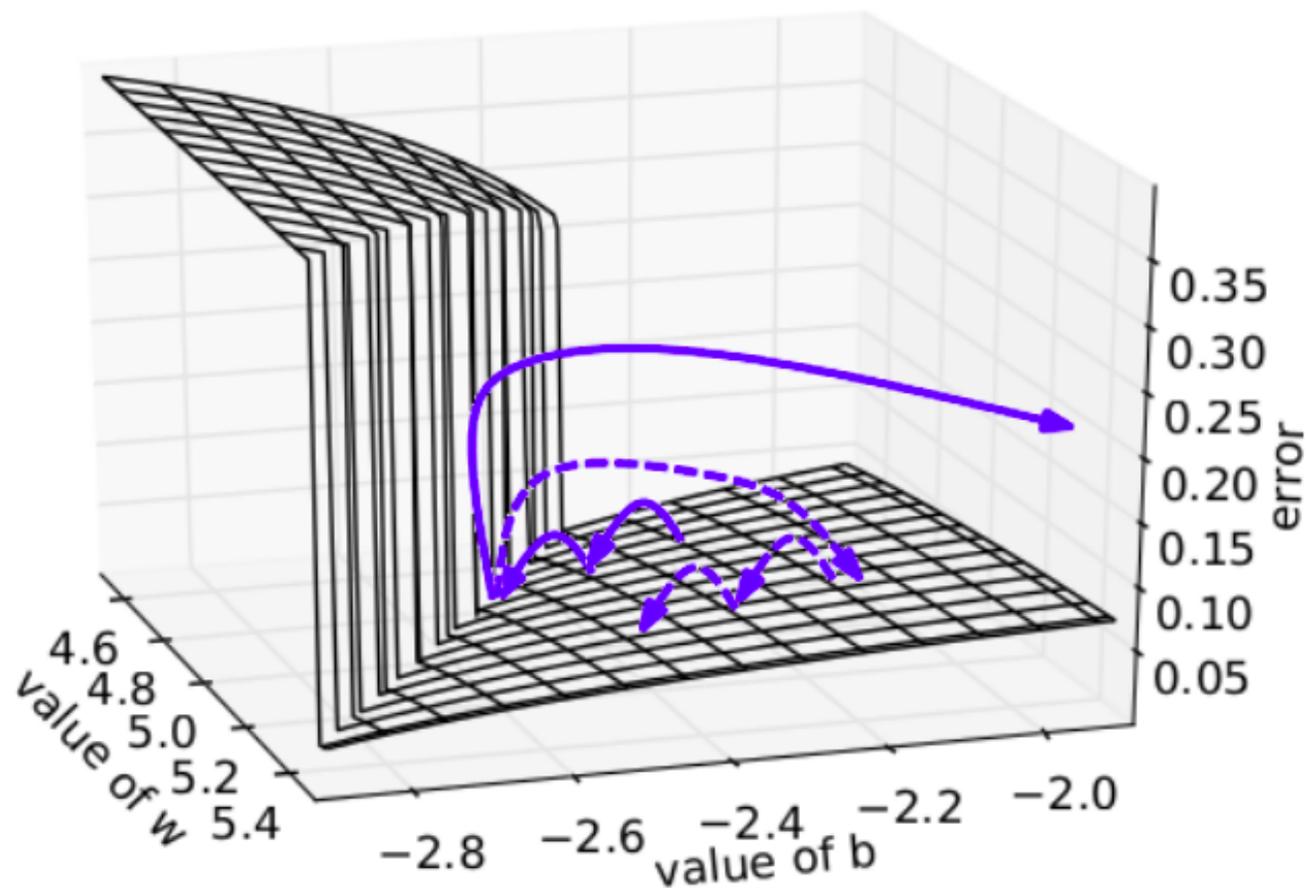
梯度裁剪则用于防止梯度爆炸。在训练过程中，由于梯度累积，可能会出现异常大的梯度，导致权重更新时步长过大，破坏训练稳定性。通过设置一个阈值限制梯度大小，当梯度范数超过该阈值时进行缩放，避免极端步长变化，从而防止训练中的不稳定现象。

其它稳定训练的策略

随着 LLMs 规模扩大，训练中损失突变的情况变得更加普遍，加剧了训练的不稳定性。为了解决该问题，PaLM 和 OPT 模型采用了一种简单但有效的策略，即在损失突增发生之前，从较早的检查点重新开始训练，并跳过可能导致问题的数据。该策略提高了训练稳定性。

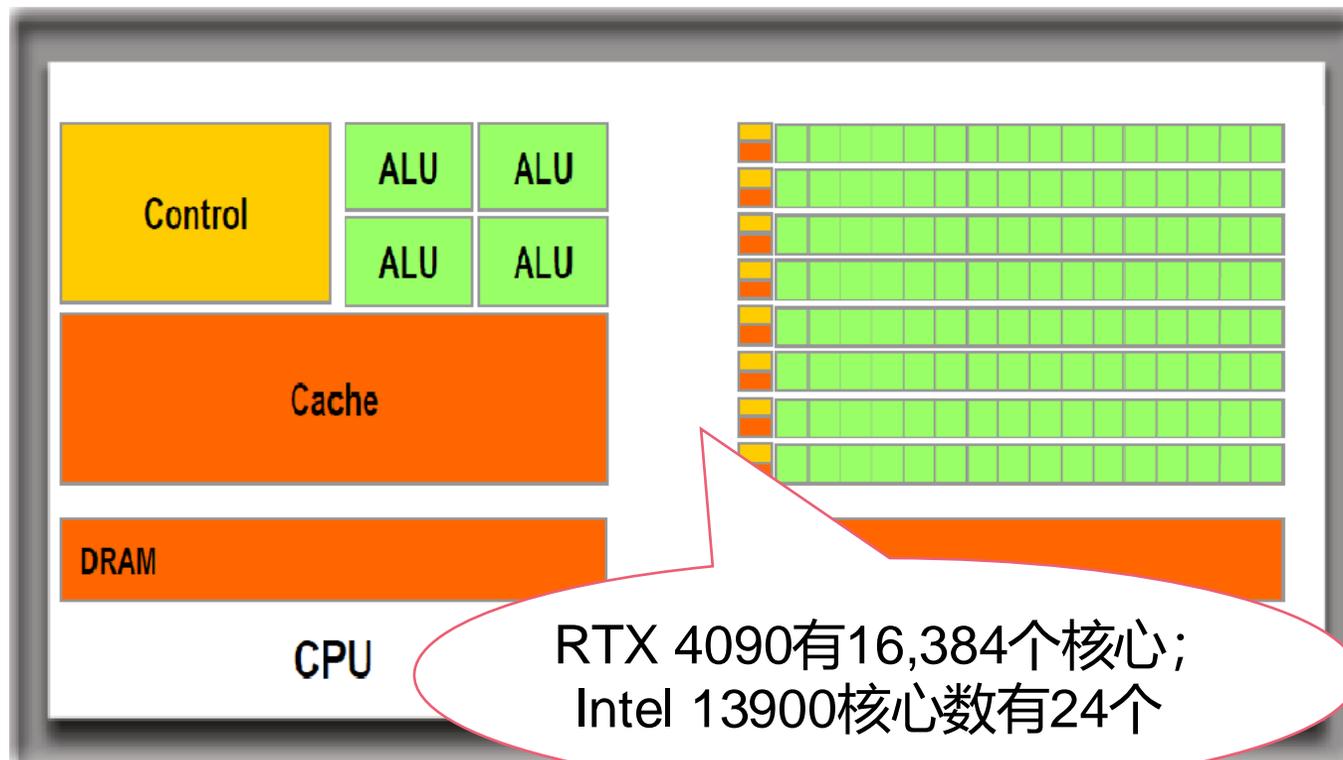
8.3.2 优化参数设置

■ 稳定训练：梯度裁剪



8.3.2 优化参数设置

CPU 与 GPU 架构对比



CPU:

- 由控制单元 (Control)、运算单元 (ALU)、缓存 (Cache) 等组成;
- 对比GPU, CPU的控制单元数量较小
- 适合处理复杂性任务

GPU:

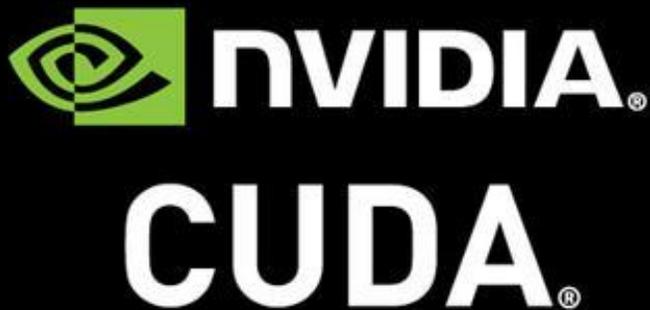
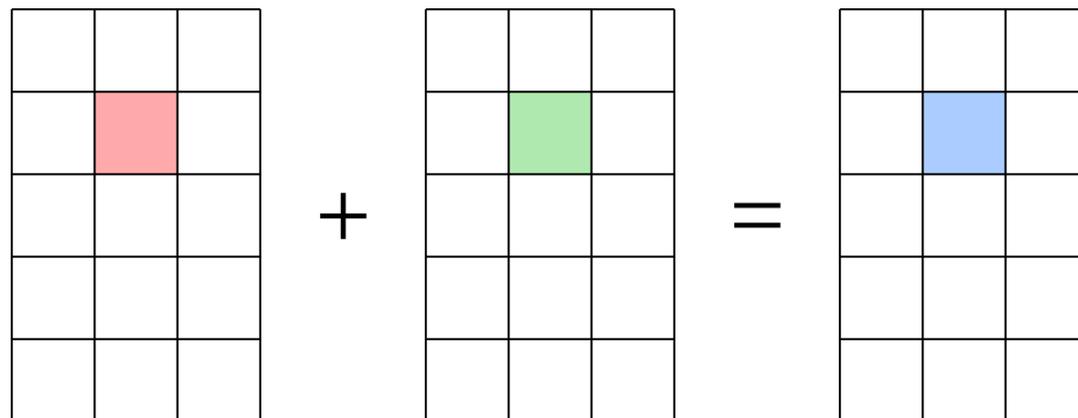
- 拥有大量运算单元
- 具有强大的并行计算能力
- 简化的控制器, 无法处理复杂性较高的任务

深度神经网络训练涉及大量的矩阵 (张量) 运算。由于图形处理器 (GPU) 架构包含大量的计算单元, 这些单元能够并行处理矩阵运算, 相较于中央处理器 (CPU), GPU提供了更高的计算效率。

8.3.2 优化参数设置

■ CUDA编程示例

```
void vecAdd(float *a, float *b, int n) {  
    for (int i = 0; i < n; i++) {  
        b[i] = a[i] + b[i];  
    }  
}
```



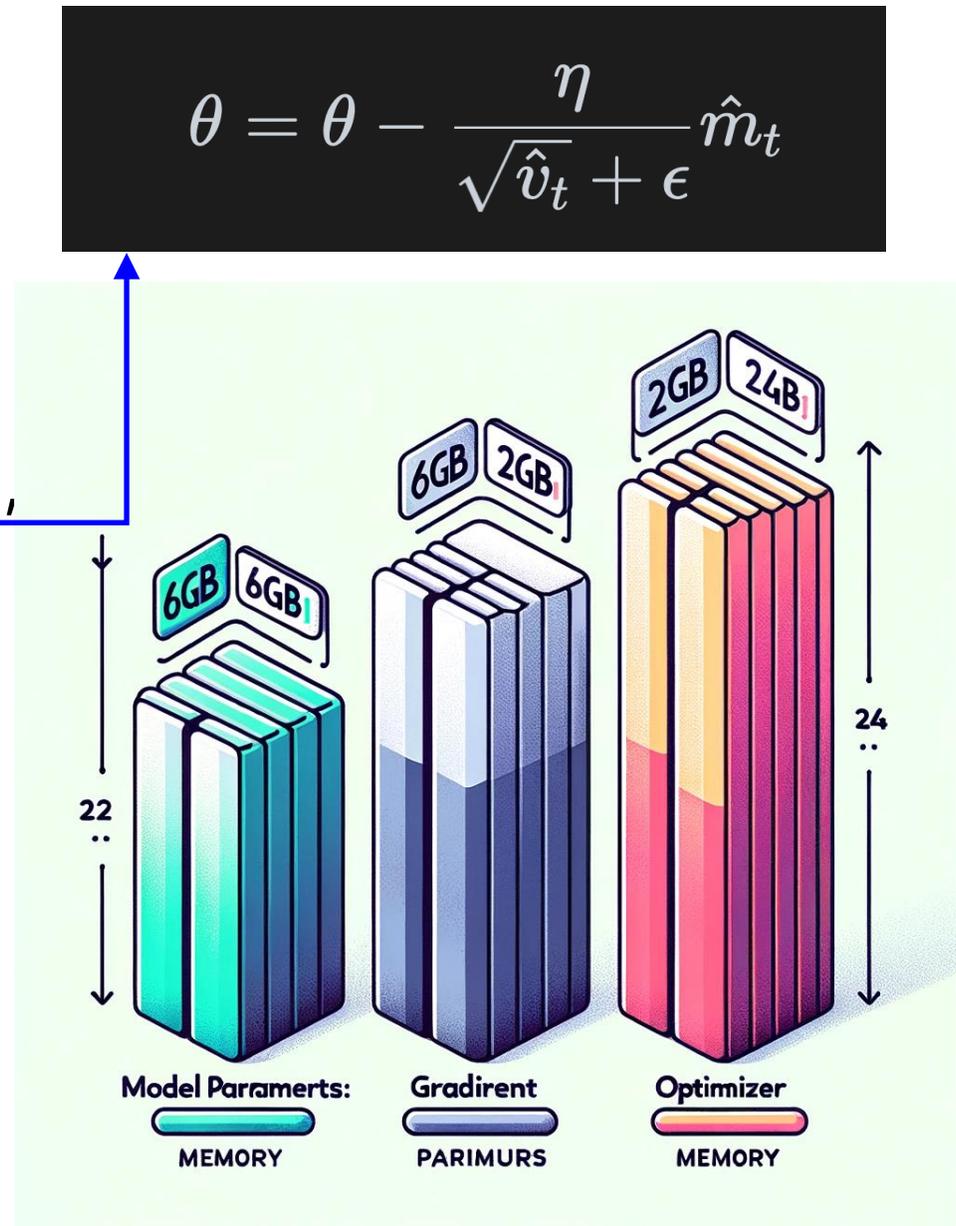
```
7 // Kernel function to add the elements of two arrays  
8 __global__ void add(int n, float *x, float *y) {  
9     for (int i = 0; i < n; i++)  
10         y[i] = x[i] + y[i];  
11 }
```

8.3.3 可扩展训练技术

□ ChatGLM-6B 模型参数:

- ChatGLM-6B架构: 隐层-4096, 中间层-11008, block数-32, 数据类型: Int8
- 模型参数所占内存: 6B x 1 bytes = 6GB
- 梯度所占内存: 6B x 1 bytes = 6GB
- 优化器 (AdamW) 参数: 2倍模型参数, 6GB x 2 = 12GB
- 总计: 约 24GB 显存

思考题 如果训练130B的模型呢?



8.3.3 可扩展训练技术

■ 量化 (quantization)

Quantization is a generic method that refers to the compression of data into a smaller space. 比如 32-bit float 量化得到 int-8。

$$\text{scaled_value} = \frac{(\text{float_value} - \text{min_float})}{(\text{max_float} - \text{min_float})} \times 255$$

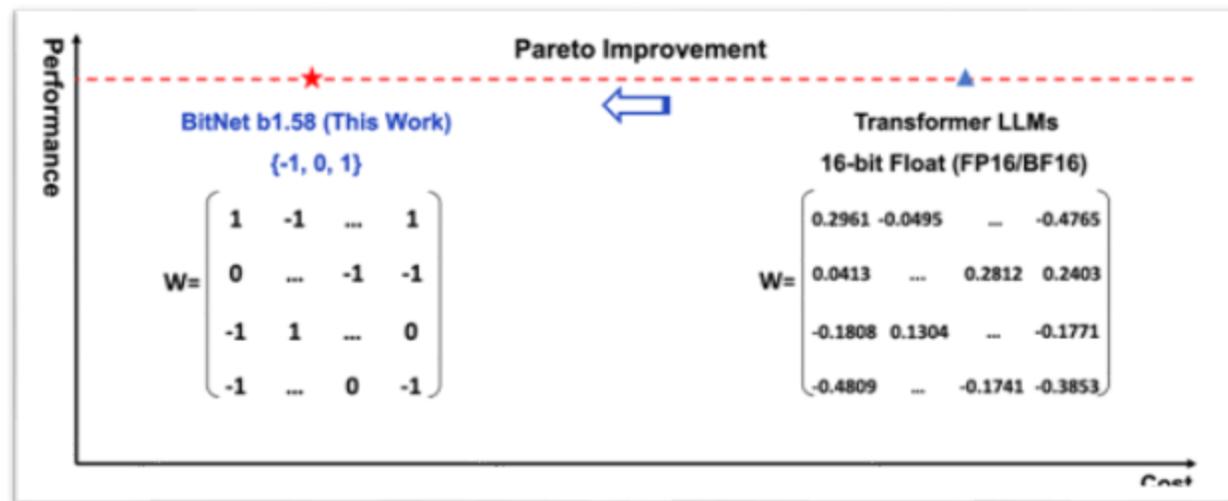
Qwen / Qwen-7B-Chat-Int8 like 8 Follow Qw

Text Generation Transformers Safetensors Chinese

The Era of 1-bit LLMs: All Large Language Models are in 1.58 Bits

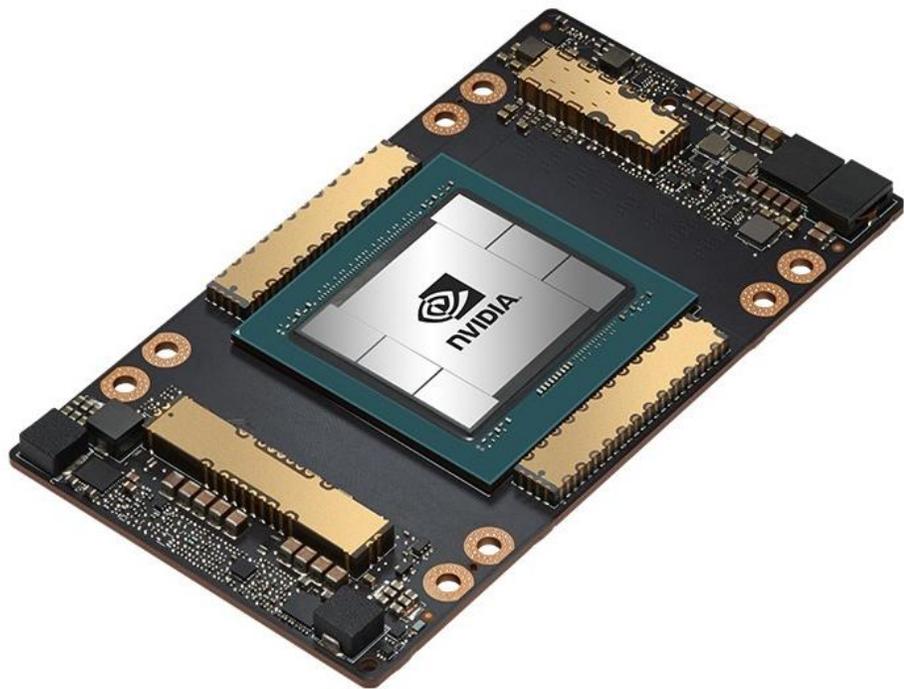
Shuming Ma* Hongyu Wang* Lingxiao Ma Lei Wang Wenhui Wang
Shaohan Huang Li Dong Ruiping Wang Jilong Xue Furu Wei^o
<https://aka.ms/GeneralAI>

Abstract



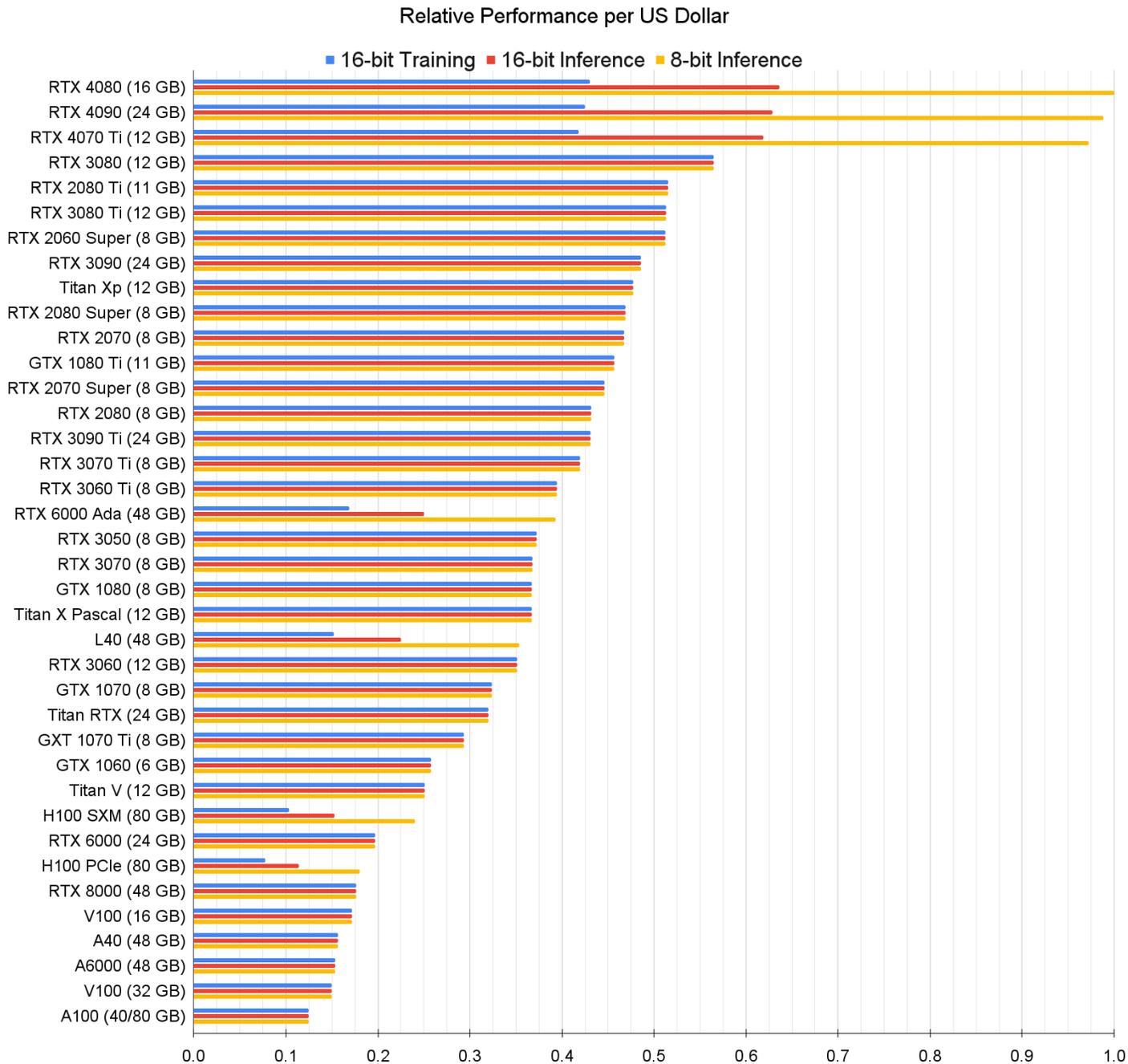
8.3.3 可扩展训练技术

3D 并行



Nvidia-A100: 40GB 或 80 GB显存，但仍无法满足使用一张卡满足训练大模型

因此需要可扩展训练技术

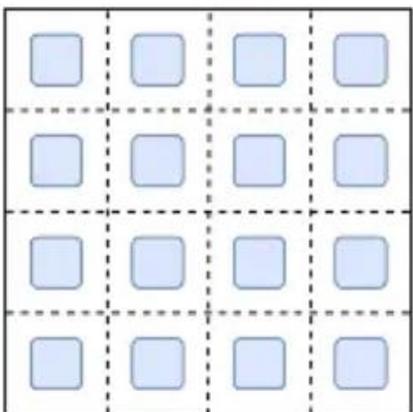


8.3.3 可扩展训练技术

■ 3D 并行

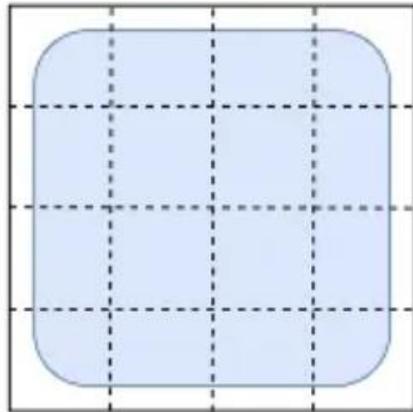
大模型突破上百亿规模参数，传统的单机单卡模式已经无法满足超大模型的训练需求。需要单机多卡、甚至是多机多卡进行分布式大模型训练

Data Parallelism



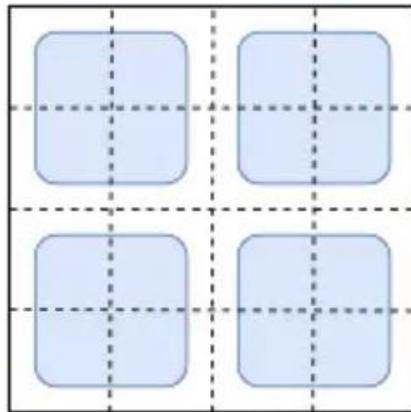
数据并行

Model Parallelism



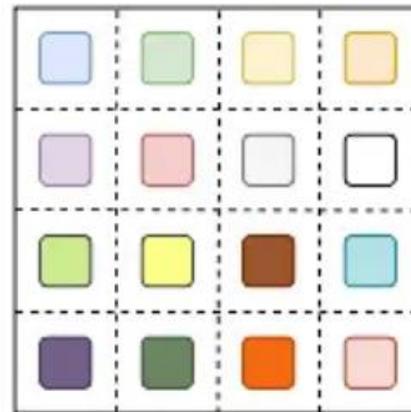
模型并行

Model and Data Parallelism



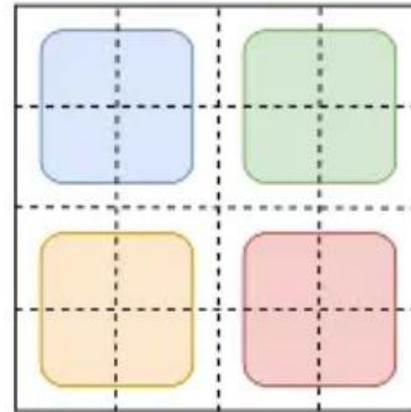
数据模型并行

Expert and Data Parallelism



专家模型

Expert, Model and Data Parallelism



混合专家

利用AI集群，一般需要根据**硬件资源与数据/模型规模**的匹配情况，考虑对**计算任务、训练数据和模型**进行划分，从而进行分布式存储和分布式训练。

8.3.3 可扩展训练技术

■ 数据并行

■ 大语言模型并行训练方法—数据并行

数据并行的特点:

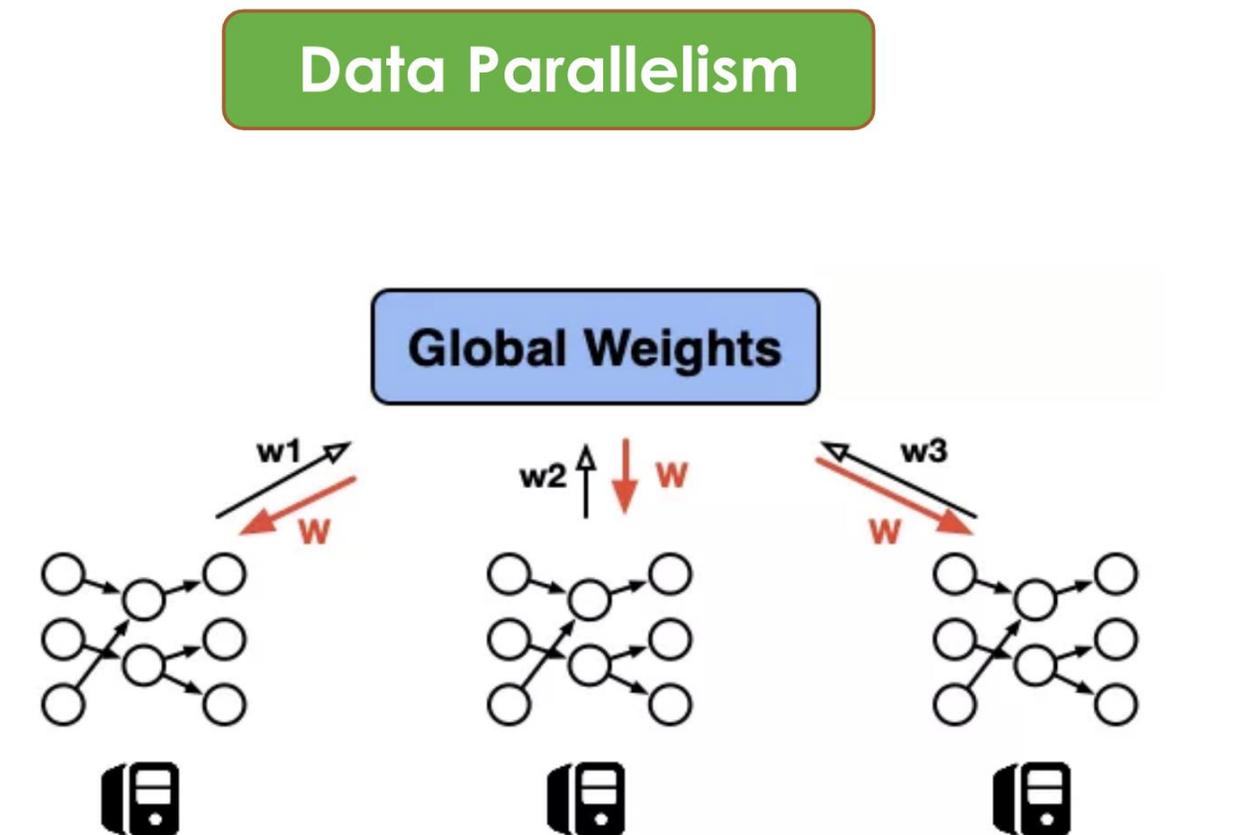
■ 如果**模型太大**无法嵌入到**一台机器**, 就将暂时未使用的参数**卸载回CPU**。

■ 数据交换传输通常在后端进行 (不干扰训练计算), 每个Mini-batch计算结束后**worker需要同步梯度或权重**, 以保证学习效率。

批量同步并行 (BSP): worker在每个Mini-batch结束时同步数据

—优点: 保证了模型**权重传递的及时性**

—缺点: 每台机器都必须排队**等待**其他机器发送**梯度**



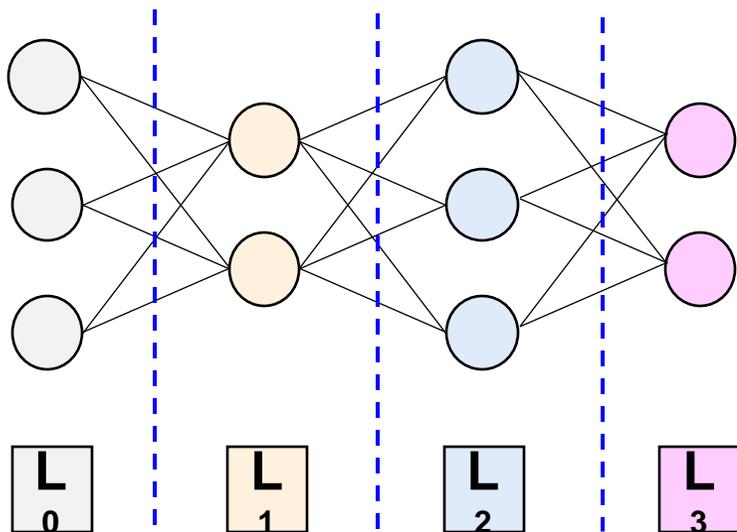
异步并行 (ASP): 每个GPU采用异步方式处理数据, 异步更新模型

—优点: **避免了**异构机器间的**相互等待**

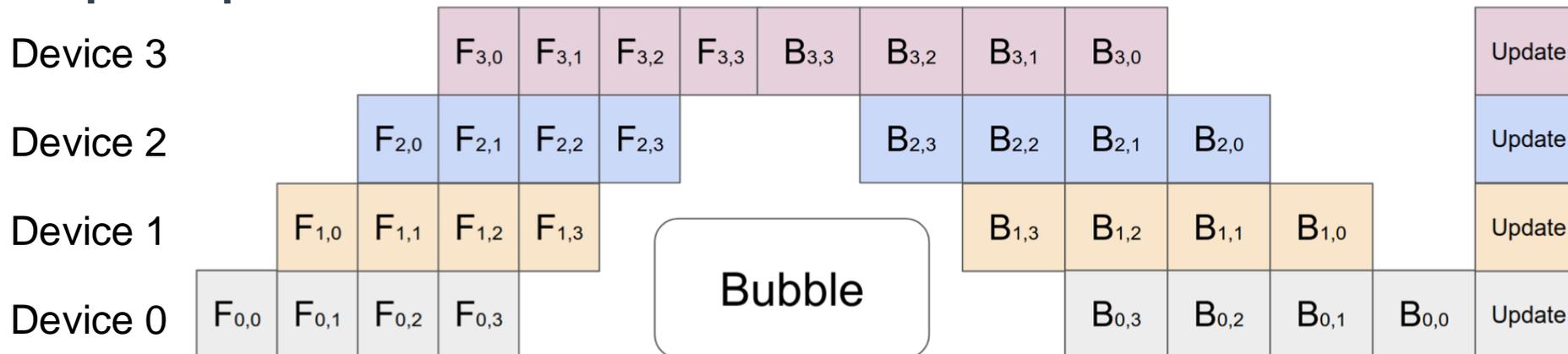
—缺点: 影响了**权重传递的时效**, **降低**了统计**学习效率**

8.3.3 可扩展训练技术

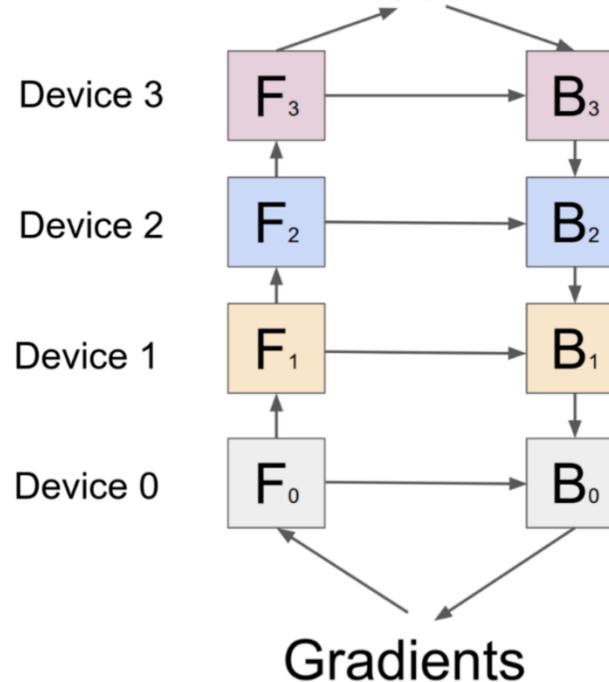
模型并行



Pipeline parallelism

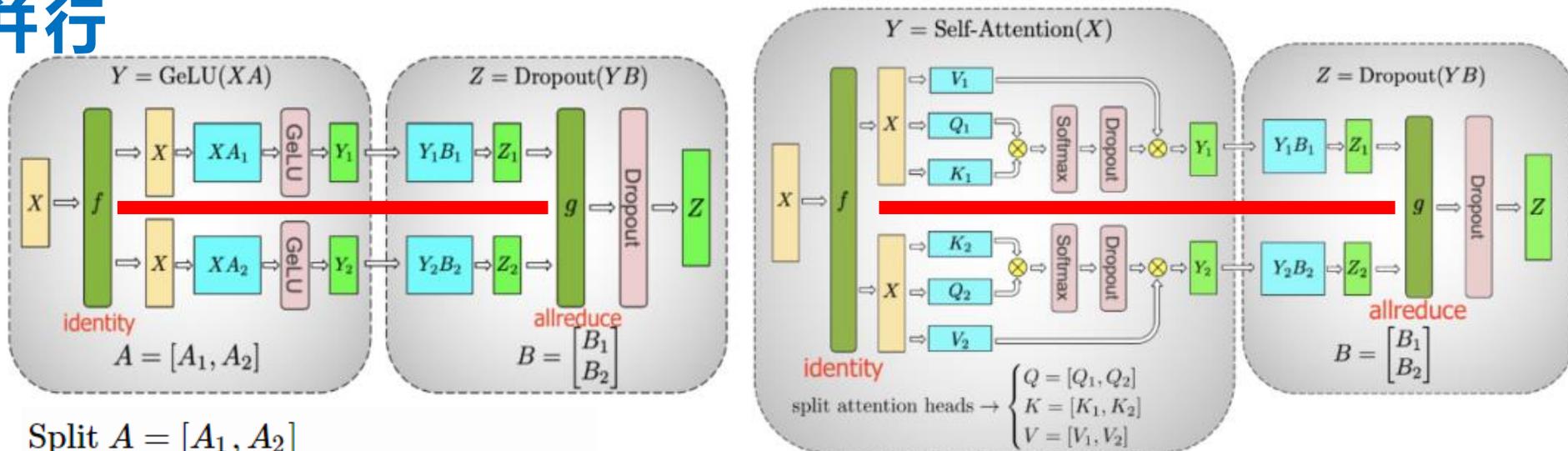


Workflow Loss



8.3.3 可扩展训练技术

张量并行

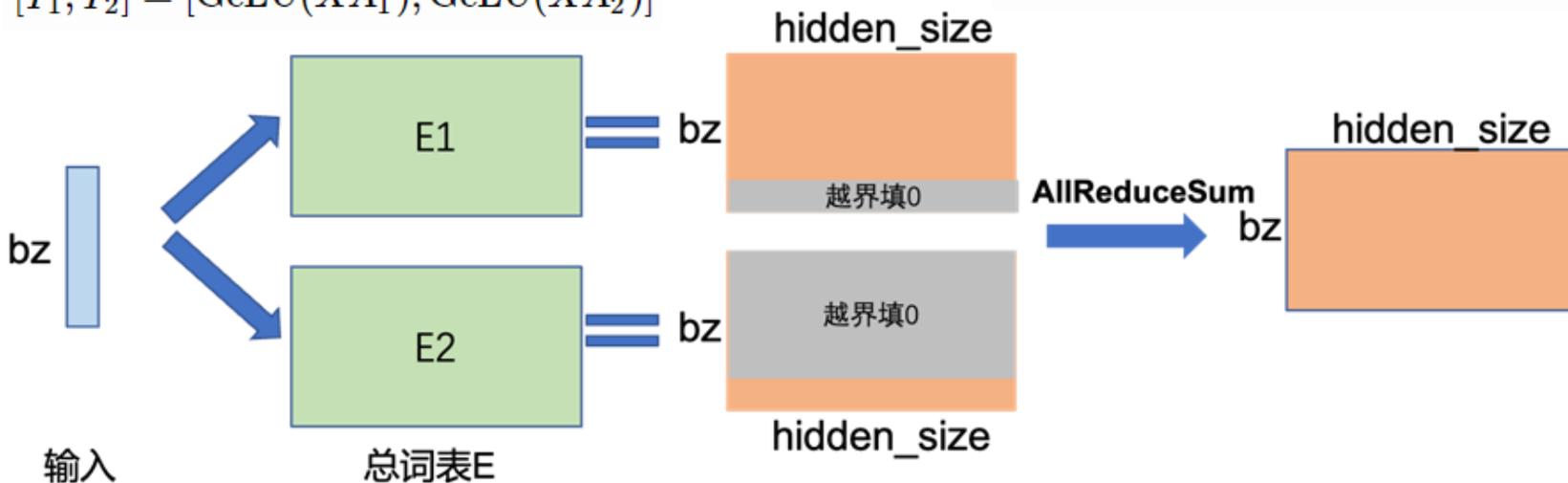


Split $A = [A_1, A_2]$

$Y = \text{GeLU}(XA)$

$[Y_1, Y_2] = [\text{GeLU}(XA_1), \text{GeLU}(XA_2)]$

$$\text{Attention}(X, Q, K, V) = \text{softmax}\left(\frac{(XQ)(XK)^T}{\sqrt{d_k}}\right)XV$$



8.3.3 可扩展训练技术

数据并行+模型并行+张量并行 = NVIDIA + Microsoft
PTD-P解决方案 (APEX)

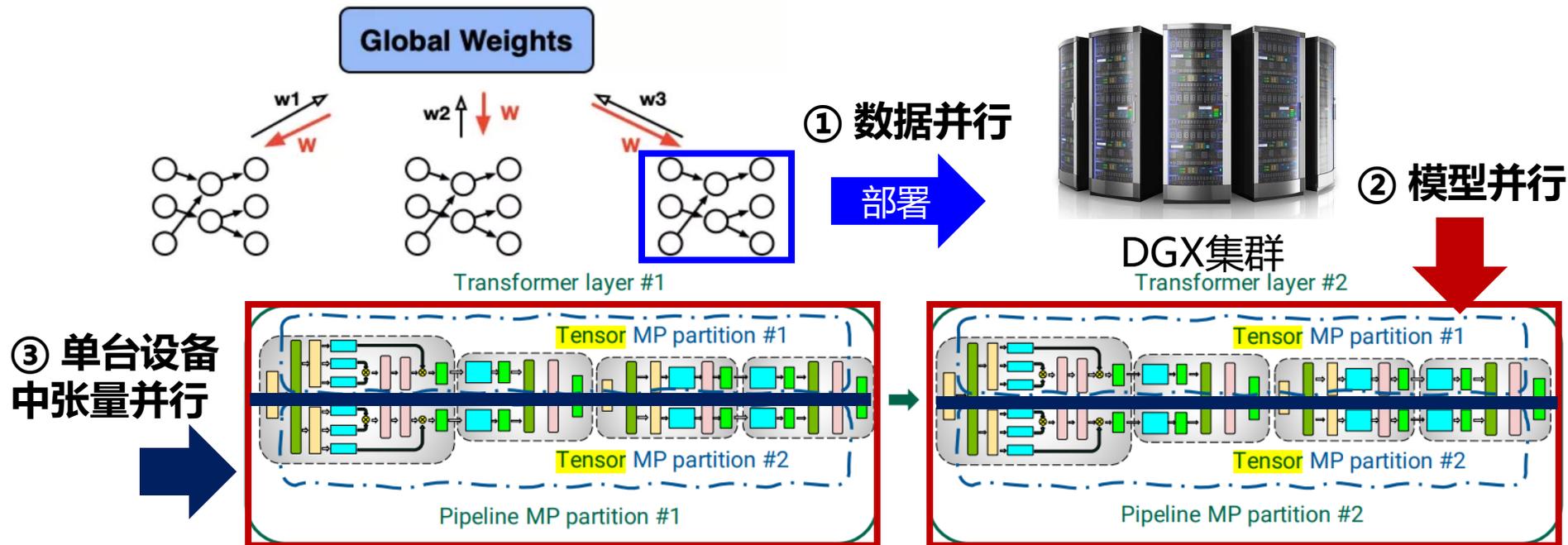
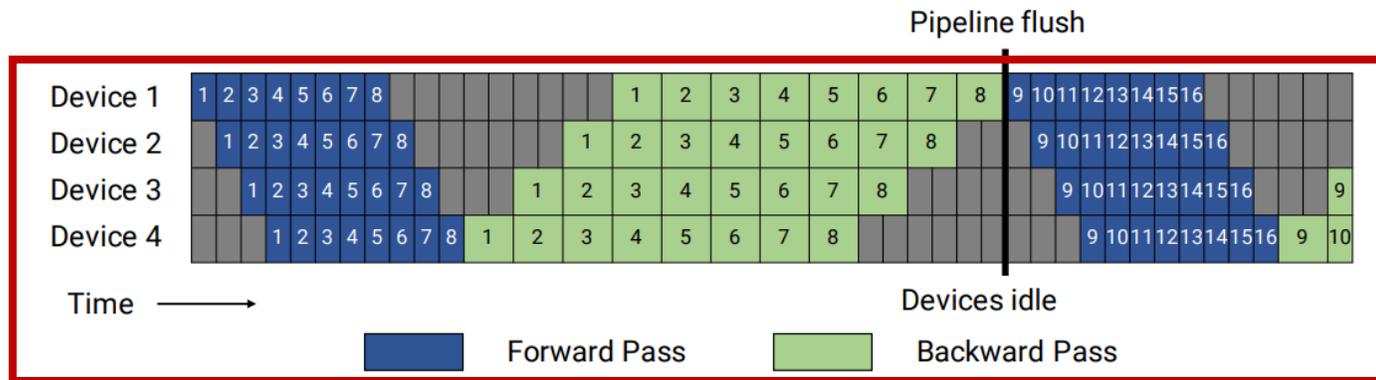


Figure 2: Combination of **tensor** and pipeline model parallelism (MP) used in this work for transformer-based models.

② 拆分为8个
Transformer+
MLP 结构, 4台
设备构建处理流



□ Narayanan D, Shoeybi M, Casper J, et al. Efficient large-scale language model training on gpu clusters using megatron-lm[C] //Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis. 2021: 1-15.

8.3.3 可扩展训练技术

■ ZeRO-1

训练过程中需要维护的大量中间变量，如动量、平方梯度等。这些状态通常会占用大量内存

方案 通过将优化器参数分布式存储在多个设备，使得单设备只保存一部分参数，从而显著减少了内存的使用

■ ZeRO-2

训练大规模模型时，梯度计算作为反向传播算法的核心，其数据规模非常庞大，会占用大量内存

方案 梯度也进行拆分，并将其分布到多个 GPU 上。每个 GPU 只需存储和处理一部分梯度

■ ZeRO-3

大模型的参数规模也非常庞大

方案 除了拆分优化器状态和梯度外，还通过拆分模型参数来最大化内存效率。

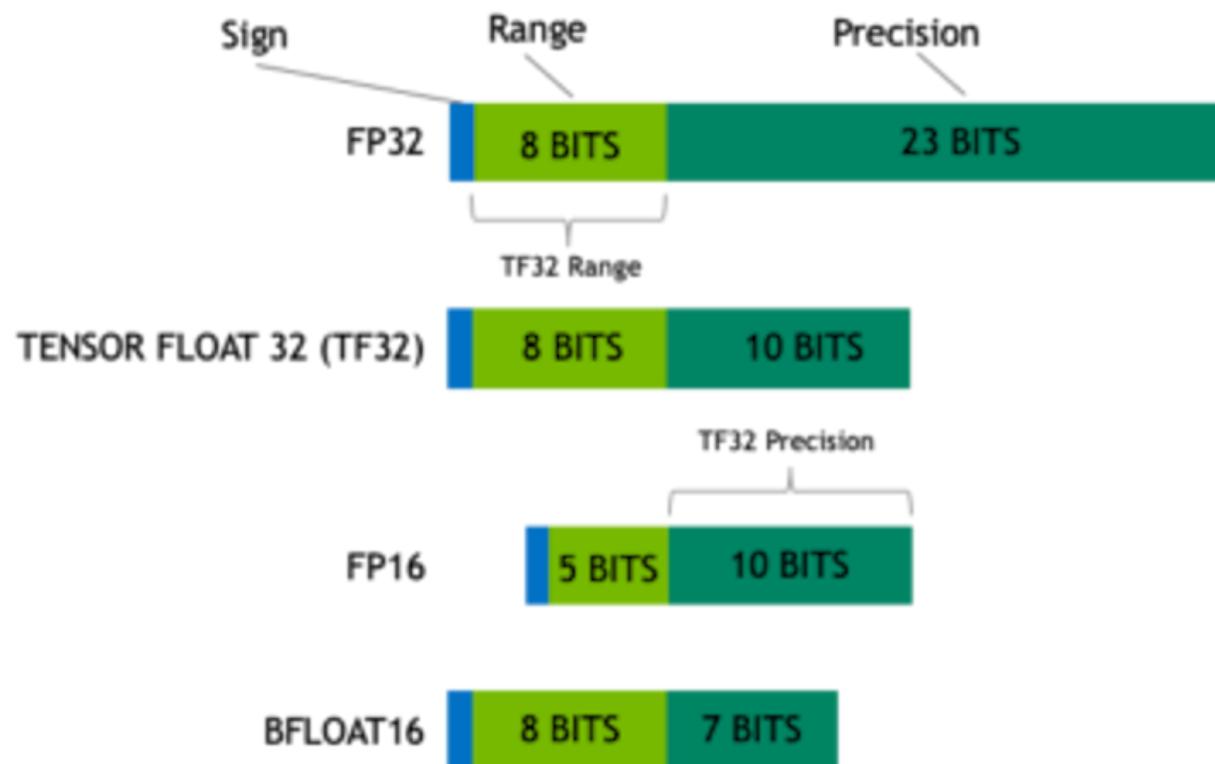
分布式存储和按需加载的机制有效减少了内存占用。ZeRO技术适用于在显存资源极其有限的情况下，仍然需要训练超大规模模型的场景。

8.3.3 可扩展训练技术

■ 混合精度训练

□ 大模型训练常用数据类型：

- FP32 (Full Precise Float 32) 单精度：单个参数需要 32 bits
- TF32 (Tensor Float 32) 单精度：由NVIDIA提出的单精度
- FP16 半精度：单个参数需要 16 bits
- INT8 整型：单个参数需要 8 bits
- BF16 (Brain Floating Point) 半精度：单个参数需要16 bits，由google提出



一些研究工作尝试开发更激进的 INT4 量化方法。LLaMA 和 通义千问 等开源模型也均提供了 INT4 量化版本的模型副本

浮点数(Floating Point)的表示范围

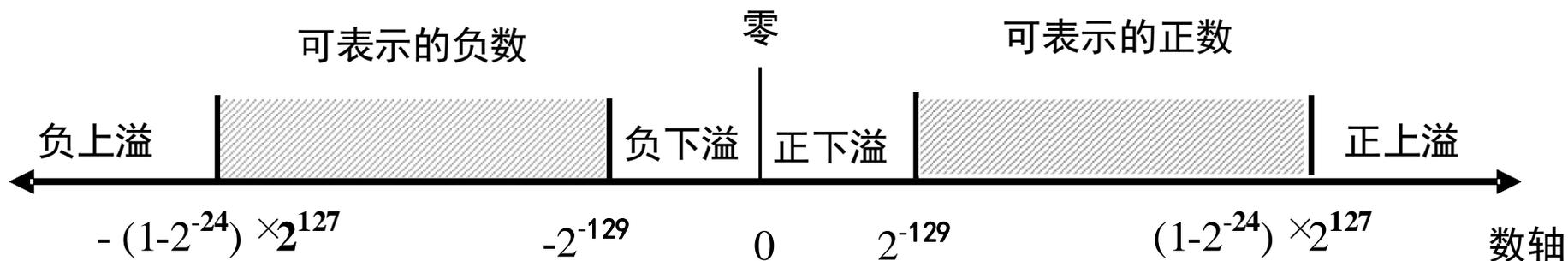
例：画出下述32位浮点数格式的规格化数的表示范围。



第0位数符S；第1~8位为8位移码表示阶码E（偏置常数为128）；第9~31位为24位二进制原码小数表示的尾数数值部分M。规格化尾数的小数点后第一位总是1，故规定第一位默认的“1”不明显表示出来。这样可用23个数位表示24位尾数。

最大正数： $0.11...1 \times 2^{11...1} = (1-2^{-24}) \times 2^{127}$ 最小正数： $0.10...0 \times 2^{00...0} = (1/2) \times 2^{-128}$

因为原码是对称的，所以其表示范围关于原点对称。



机器0：尾数为0 或 落在下溢区中的数

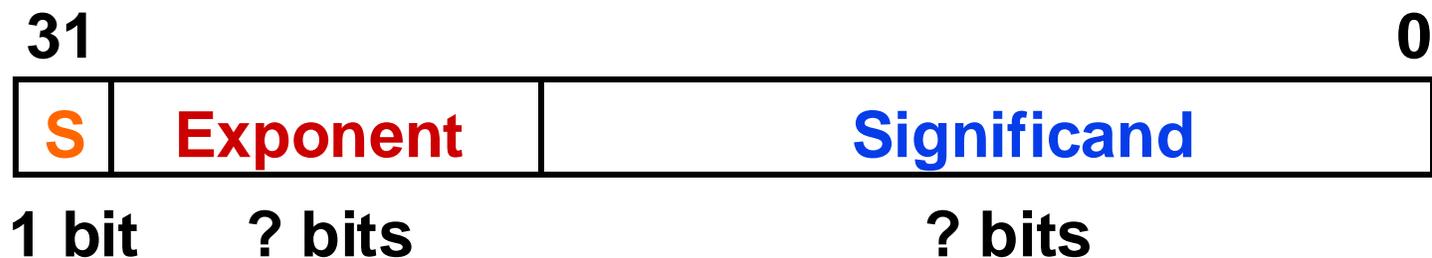
浮点数范围比定点数大，但数的个数没变多，故数之间更稀疏，且不均匀

浮点数的表示

- Normal format（规格化数形式）：

$$\text{+/-}1.\text{xxxxxxxxxxx} \times R^{\text{Exponent}}$$

- 32-bit 规格化数：



S 是符号位（Sign）

Exponent 用移码（增码）来表示

Significand 表示 **xxxxxxxxxxxx**，尾数部分

（基可以是 2 / 4 / 8 / 16，约定信息，无需显式表示）

- 早期的计算机，各自定义自己的浮点数格式

问题：浮点数表示不统一会带来什么问题？

本章内容

- 8.1 概述
- 8.2 预训练数据工程
 - 8.2.1 预训练数据源
 - 8.2.2 多模态数据集
 - 8.2.3 数据处理
 - 8.2.4 模型性能关系
- 8.3 预训练方法
 - 8.3.1 预训练任务
 - 8.3.2 优化参数设置
 - 8.3.3 可扩展训练技术
- 8.4 讨论

8.4 讨论

思考题1 如何以较小的代价修正大语言模型存储的知识？

思考题2 讨论在大规模语言模型预训练中常用的优化技巧，如学习率调度、混合精度训练、分布式训练等。