

A Comprehensive Guide to **Pre-training LLMs**

Fenelleb

The finest collection of data the web has to offer



Text

If you are done with step 1 of processing the raw data, now the question irise's is - How we will train the neural network on this data? As mentioned i the FineWeb, there are 15 trillion tokens and 44TB of disk space data set that need to be fed to the neural network for further processin The next essential step is tokenization, a process that prepares the raitext data for training large language models (LLMs). Let's break down ho tokenization works and its significance based on the transcript.

UTF8 Encode the text

> Two symbols with extremely long sequence

Group of 8 Bits

01001010, 01001011, 01001100, 01001101, 01001110, 01001111, 01010000, 01010001, 01010010, 01010011, 01010100, 01010101, 01010110, 01010111, 01011000, 01011001, 01011010, 01100001 01101011, 01101100, 01101101, 01101110, 01101111, 01110000, 01110001, 01110010, 01110011 01110100, 01110101, 01110110, 01110111, 01111000, 01111001, 01111010, 00110000, 00110001 00110010, 00110011, 00110100, 00110101, 00110110, 00110111, 00111000, 00111001.

Group of 8 bits = 1 Byte

Since each bit can be 0 or 1, an 8-bit sequence can represent:

 $2^8 = 256$

This means a single byte can encode 256 unique values, ranging from 0 to 255.

Sequence of Bytes

182, 250, 584, 690, 5542, 279, 30828, 4009, 389, 420, 828, 30, 1666, 9932, 304, 279, 31253, 6109, 11, 1070, 527, 220, 868, 32610, 11460, 323, 220, 2096, 32260, 315, 13668, 3634, 828, 743, 430, 1205, 311 387, 23114, 311, 279, 30828, 4009, 369, 4726, 8863, 627, 791, 1828, 771, 8, 3094, 374, 4037, 2065, 11 64, 1920, 430, 48542, 279, 7257, 1495, 828, 369, 4967, 3544, 4221, 4211, 320, 4178, 22365, 570, 914, 753, 1464, 1523, 1268, 4037, 2065, 4375, 323, 1202, 26431, 3196, 389, 279, 36815, 627

> These are unique IDs or Symbols and later we can group use the BPE Algorithm to look for consecutive bites or symbols to further decrease the length.



What is the LLM Pretraining Stage?

- The LLM pretraining stage is the first phase of teaching a large language model (LLM) like how to understand and generate text.
- Think of it as reading a massive number of books, articles, and websites to learn grammar, facts, and common patterns in language.
- During this stage, the model processes billions of words (data) and predicts the next word (token) in a sentence repeatedly, refining its ability to generate coherent and relevant responses.
- However, at this point, it doesn't fully "understand" meaning like a human—it just recognizes patterns and probabilities.



What can a Pre-trained LLM do?

- Pre-trained Large Language Models (LLMs) can perform a wide range of tasks, including text generation, summarization, translation, and sentiment analysis.
- They assist in code generation, question-answering, and content recommendation. LLMs can extract insights from unstructured data, facilitate chatbots, and automate customer support.
- They enhance creative writing, provide tutoring, and even generate realistic conversations. Additionally, they assist in data augmentation, legal analysis, and medical research by analyzing vast amounts of information efficiently.
- Their ability to understand and generate human-like text makes them valuable for various industries, from education and finance to healthcare and entertainment. However, they require fine-tuning for domain-specific accuracy.



Step 1: Process the Internet Data

- There are multiple stages of training an LLM but here we will first talk about the LLM Pretraining stage.
- The performance of a large language model (LLM) is deeply influenced by the quality and scale of its pretraining dataset. If your dataset is clean, structured and easy to process, the model will work accordingly.

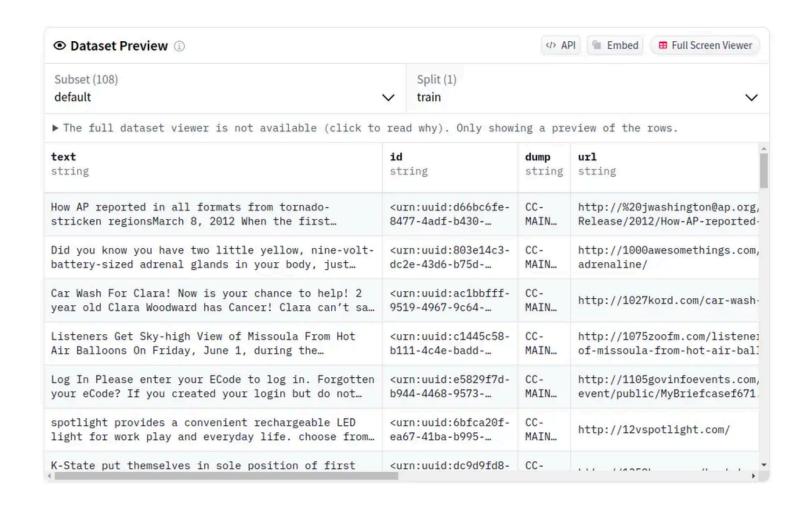


The finest collection of data the web has to offer





- However, for many state-of-the-art open LLMs like <u>Llama 3</u> and Mixtral, the details of their pretraining data remain a mystery—these datasets are not publicly available, and little is known about how they were curated.
- To address this gap, Hugging collected data from the internet and curated FineWeb, a large-scale dataset (this is a portion of data available on the internet) specifically designed for LLM pretraining.
- This high-quality and diverse dataset has 15 trillion tokens and occupies 44TB of disk space, FineWeb is built from 96 CommonCrawl snapshots and has been shown to produce better-performing models than other publicly available pretraining datasets.





- What sets FineWeb apart is its transparency:
- It meticulously documented every design choice, running detailed ablations on deduplication and filtering strategies to refine the dataset's quality.

Where Does the Raw Data Come From?

There are two main sources:

- Crawling the web yourself Used by companies like OpenAl and Anthropic.
- Using public repositories CommonCrawl, a nonprofit that has been archiving web data since 2007.

For FineWeb, they followed the approach of many LLM teams and used CommonCrawl (CC) as the starting point. CC releases a new dataset every 1-2 months, typically containing 200-400 TiB of text.

For example, the April 2024 crawl includes 2.7 billion web pages with 386 TiB of uncompressed HTML. Since 2013, CC has released 96 crawls, plus 3 older-format crawls from 2008-2012.



For more information, visit this article



Advanced

Best of Tech

Generative Al

Guide

LLMs

A Comprehensive Guide to Pre-training LLMs

Explore the fundamentals of LLM pretraining, its role in Al advancements, and how it shapes models like GPT-4.

Pankaj Singh 12 Feb, 2025