

Comprehensive Guide to Building Foundational LLMs Like ChatGPT From Scratch for Enterprise Deployment

Authored by: A PhD Researcher specializing in Large Language Models and Systems Engineering

Introduction

The advent of Large Language Models (LLMs) like ChatGPT has marked a significant inflection point in artificial intelligence, demonstrating remarkable capabilities in natural language understanding and generation. These models, primarily based on the Transformer architecture¹, have shown proficiency across a wide array of tasks, from text summarization and translation to code generation and complex reasoning. Building such foundational models from scratch, however, is a monumental undertaking, requiring deep expertise across mathematics, data science, software engineering, and large-scale systems infrastructure. This report provides a comprehensive, step-by-step technical guide for constructing GPT-style Transformer-based LLMs, suitable for enterprise-scale deployment. It delves into the mathematical underpinnings, data curation processes, implementation details, infrastructure requirements, training dynamics, optimization techniques, and deployment strategies necessary for building and operationalizing these powerful models.

1. Mathematical Foundations and Model Architecture

The foundation of modern LLMs like ChatGPT is the Transformer architecture, introduced by Vaswani et al. in their seminal 2017 paper, "Attention Is All You Need".¹ This architecture revolutionized sequence transduction tasks by relying entirely on attention mechanisms, eschewing the recurrent or convolutional layers prevalent in earlier models.² This shift enabled significantly greater parallelization during training, a key factor in scaling models to billions or trillions of parameters.¹

1.1 The Original Transformer Architecture (Encoder-Decoder)

The original Transformer was designed for sequence-to-sequence tasks like machine translation and comprised two main components: an encoder stack and a decoder stack.²

- **Encoder Stack:** Composed of N (typically 6) identical layers. Each layer contains two sub-layers:
 1. A multi-head self-attention mechanism.
 2. A position-wise fully connected feed-forward network (FFN).

- **Decoder Stack:** Also composed of N (typically 6) identical layers. Each layer has three sub-layers:
 1. A *masked* multi-head self-attention mechanism (to prevent attending to future tokens).
 2. A multi-head attention mechanism over the encoder's output.
 3. A position-wise fully connected feed-forward network (FFN).

Crucially, residual connections² are employed around each sub-layer, followed by layer normalization.² The output of each sub-layer is thus computed as $\text{LayerNorm}(x + \text{Sublayer}(x))$, facilitating gradient flow in deep networks. Some later work suggests that applying normalization *before* the sub-layer (Pre-LN) can improve training stability.⁴

1.2 Self-Attention: The Core Mechanism

The heart of the Transformer is the self-attention mechanism, specifically Scaled Dot-Product Attention.¹ It allows the model to weigh the importance of different tokens in the input sequence when computing the representation for a specific token.

The attention function maps a query (Q) and a set of key-value (K,V) pairs to an output. The query, keys, and values are vectors derived from the input token embeddings. The output is a weighted sum of the values, where the weight for each value is computed based on the compatibility (dot product) of the query with the corresponding key.²

The mathematical formulation is:

$$\text{Attention}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$

Where:

- Q, K, V are matrices packing the query, key, and value vectors, respectively.
- d_k is the dimension of the key vectors.
- The scaling factor $\frac{1}{\sqrt{d_k}}$ is crucial. It prevents the dot products from growing too large (especially for high d_k), which could push the softmax function into regions with extremely small gradients, hindering learning.¹

In *self-attention* (or intra-attention), Q, K, and V are derived from the same source sequence, allowing the model to relate different positions within that sequence.¹

1.3 Multi-Head Attention

Instead of performing a single attention calculation, the Transformer uses multi-head attention. The queries, keys, and values are linearly projected h times (the number of heads, typically 8⁶) into lower-dimensional subspaces. Scaled Dot-Product Attention

is then applied in parallel to each of these projected versions. The outputs from the h heads are concatenated and projected again with a final linear layer.²

$\text{MultiHead}(Q,K,V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h) W_O$

where $\text{head}_i = \text{Attention}(Q W_{iQ}, K W_{iK}, V W_{iV})$

Here, W_{iQ}, W_{iK}, W_{iV} are the projection matrices for the i -th head, and W_O is the output projection matrix. This allows the model to jointly attend to information from different representation subspaces at different positions simultaneously.²

1.4 Positional Encoding

Since self-attention itself is permutation-invariant (it doesn't inherently know the order of tokens), positional information must be explicitly injected.² The original Transformer added positional encodings to the input embeddings. These encodings use sine and cosine functions of different frequencies:

$\text{PE}(\text{pos}, 2i) = \sin(\text{pos}/10000^{2i/d_{\text{model}}})$

$\text{PE}(\text{pos}, 2i+1) = \cos(\text{pos}/10000^{2i/d_{\text{model}}})$

where pos is the position and i is the dimension index.² This scheme allows the model to learn relative positions, as the encoding for any position $\text{pos}+k$ can be represented as a linear function of the encoding at position pos .

1.5 Position-wise Feed-Forward Networks (FFN)

Each layer in the encoder and decoder includes a fully connected feed-forward network applied independently to each position.² It typically consists of two linear transformations with a ReLU activation function in between:

$\text{FFN}(x) = \max(0, x W_1 + b_1) W_2 + b_2$

The dimensionality is often expanded in the inner layer (e.g., from $d_{\text{model}}=512$ to $d_{\text{ff}}=2048$) before being projected back.

1.6 GPT-Style Decoder-Only Architecture

For generative tasks like language modeling, decoder-only architectures, popularized by the GPT series, have become standard.¹¹ These models essentially use only the decoder stack of the original Transformer.

A key component is **Masked Multi-Head Self-Attention**. During training and generation, a mask is applied to the attention scores matrix before the softmax step. This mask prevents any token at position i from attending to tokens at positions $j > i$.¹¹ This is typically achieved by setting the attention scores corresponding to future positions to negative infinity.¹² This ensures the autoregressive property: the prediction for the next token depends only on the preceding tokens.¹¹

1.7 Training Objective: Causal Language Modeling (CLM)

Decoder-only models are typically trained using a Causal Language Modeling (CLM) objective.¹⁴ The model learns to predict the next token in a sequence given all the preceding

tokens.¹¹ Mathematically, the objective is to maximize the log-likelihood of the sequence:
$$\text{LLM}(x) = \sum_i \log P(x_i | x_{<i}; \theta)$$

This is typically achieved by minimizing the Cross-Entropy Loss between the model's predicted probability distribution for the next token (obtained by applying softmax to the output logits) and the actual next token in the training data.¹⁴ During loss calculation, padding tokens are usually ignored (e.g., by setting their label index to -100 in PyTorch).¹⁸

1.8 Optimizers and Learning Rate Schedulers

Training large models requires sophisticated optimization algorithms and learning rate schedules.

- **Optimizers:** Adam¹⁹ and AdamW¹⁹ are the most common choices. AdamW is generally preferred for training Transformers.¹⁹ The key difference lies in the implementation of weight decay (L2 regularization). Adam incorporates the L2 penalty into the gradient estimate, which can interact poorly with the adaptive learning rates.²² AdamW decouples weight decay, applying it directly to the weights after the gradient-based update step²², leading to more stable training and often better generalization.²² Memory-efficient alternatives like Lion²¹ and Adafactor¹⁹ also exist, reducing the memory overhead associated with optimizer states.²¹
- **Learning Rate Schedulers:** A fixed learning rate is rarely optimal. Common schedules include:
 - **Warmup:** Starting with a small learning rate and gradually increasing it over the first few thousand steps helps stabilize training early on.³⁰
 - **Cosine Decay:** After warmup, the learning rate decreases following a cosine curve, often decaying to 10% of the peak learning rate.³⁰ This is a very common schedule for LLM pretraining.³¹
 - **Linear Decay:** An alternative where the learning rate decreases linearly, potentially to zero (D2Z). Some studies suggest linear D2Z can outperform cosine decay, especially for compute-optimal training regimes, potentially saving significant compute.³⁰
 - **Warmup-Stable-Decay (WSD):** A newer schedule involving a warmup, a long phase with a constant high learning rate, and a final rapid decay phase.³¹

1.9 Architectural Variations and Improvements

The basic Transformer architecture has undergone numerous refinements:

- **GPT Evolution (GPT-2/3/4):** Key changes involve massive increases in model size (parameters), expansion of the context window (from 512/1024 in early models to 128k+ in GPT-4), and the introduction of multimodality (GPT-4 accepts image inputs).³³ GPT-3.5 introduced instruction tuning and RLHF.³⁵

- Rotary Positional Embeddings (RoPE):** Instead of adding absolute positional encodings, RoPE modifies the query and key vectors through rotation matrices whose angles depend on the token's absolute position.³⁷ This elegantly encodes relative positional information directly into the attention mechanism (qmTkn becomes dependent on $m-n$).³⁷ RoPE has shown strong performance, especially for long sequences⁴², and is used in models like LLaMA and PaLM.³⁸ Potential inefficiencies for long-distance retrieval have been noted.⁴⁴
- Attention with Linear Biases (ALiBi):** An alternative relative positional encoding method that doesn't modify queries/keys but adds a static, non-learned bias penalty to the attention scores based on the distance between tokens.⁴⁶ Closer tokens receive less penalty. ALiBi has demonstrated good extrapolation capabilities beyond the training sequence length.⁴⁶ It can be combined with other techniques, e.g., in SWAT⁴⁷ or FlexAttention.⁴⁹ The "attention sink" phenomenon, where initial tokens receive high attention regardless of semantic importance (potentially due to softmax normalization), is a related observation often discussed alongside windowed attention and methods like ALiBi.⁴⁶
- Gated Linear Units (GLU) in FFN:** The FFN sub-layer has seen alternatives to ReLU/GELU. GLU variants like SwiGLU (Sigmoid Gated Linear Unit) and GeGLU (GELU Gated Linear Unit) have become popular.⁴ SwiGLU, used in models like LLaMA and PaLM, computes $\text{FFN}(x, W, V) = (\text{Swish}(xW) \odot xV)W_2$, where $\text{Swish}(x) = x \cdot \sigma(\beta x)$ and \odot is element-wise multiplication. These gated activations can potentially offer better performance.⁵⁰ ReLU² is another variant explored for sparsity benefits.⁵¹
- Root Mean Square Normalization (RMSNorm):** A simplification of Layer Normalization that omits the re-centering (mean subtraction) step, only performing re-scaling based on the root mean square of the input.⁵³ The formula is $\text{RMSNorm}(x) = \frac{1}{\sqrt{n}} \sum_{i=1}^n x_i^2 + \epsilon \cdot \gamma$. It is computationally cheaper than LayerNorm⁵³ while often achieving comparable performance.⁵⁴ It's used in models like LLaMA⁵³ and Mistral.⁵³ Some research suggests the re-centering in LayerNorm might be redundant for LLMs.⁵⁶
- FlashAttention:** An optimized implementation of the attention mechanism designed to be IO-aware, significantly reducing memory reads/writes between the GPU's slow High Bandwidth Memory (HBM) and fast on-chip SRAM.⁵⁹ It achieves this using techniques like tiling (processing queries, keys, and values in blocks that fit into SRAM) and recomputation, avoiding the materialization of the full $N \times N$ attention matrix in HBM.⁶⁰ This results in substantial speedups (e.g., 3x for GPT-2⁶⁶) and linear memory complexity with respect to sequence length for the attention computation, enabling longer context windows.⁶¹ FlashAttention computes the exact attention result, not an approximation.⁶² Distributed versions

like DistFlashAttn exist for multi-GPU settings.⁶⁴

The progression from absolute positional encodings like sinusoidal to relative ones like RoPE and ALiBi, coupled with the development of memory-efficient attention mechanisms like FlashAttention, directly addresses the challenge of scaling Transformers to handle ever-longer context lengths. This capability is crucial for enhancing LLM performance on tasks requiring understanding of extensive background information or long conversational histories. Simultaneously, the adoption of RMSNorm over LayerNorm and the exploration of different FFN activations (like SwiGLU) reflect a growing emphasis on computational efficiency alongside model accuracy. These architectural choices represent a balancing act, optimizing performance within the constraints of available hardware and training budgets. The shift from Adam to AdamW also underscores the refinement of optimization strategies specifically for the nuances of large Transformer training, prioritizing stable convergence and generalization.

Table 1.1: Comparison of Positional Encoding Methods

Method	Encoding Type	Mechanism	Extrapolation Capability	Computational Cost	Key Papers/Models
Sinusoidal	Absolute	Adds fixed sine/cosine waves based on position and dimension	Poor	Low (fixed, no params)	Vaswani et al. 2017 (Original Transformer) ²
Learned Absolute	Absolute	Learns a unique embedding vector for each position	Poor	Moderate (adds parameters)	Devlin et al. 2018 (BERT) ³⁷
Rotary (RoPE)	Relative (Implicit)	Rotates query/key vectors based on	Good	Low (no extra params)	Su et al. 2021 (RoFormer), LLaMA,

		absolute position; relative in dot product			PaLM ³⁸
ALiBi	Relative (Explicit)	Adds a distance-based bias penalty directly to attention scores	Good	Very Low (no params)	Press et al. 2022, MPT ⁴⁶

Table 1.2: Comparison of Normalization Layers

Method	Formula (Simplified)	Key Difference	Computational Cost	Stability Notes	Common Models Using It
LayerNorm	$\sigma^2 + \epsilon x - \mu \gamma + \beta$	Re-centering	Higher	Standard, Pre-LN may improve	BERT, GPT series ⁵³
RMSNorm	$\text{RMS}(x)^2 + \epsilon xy$	No Re-centering	Lower	Simpler, often comparable perf.	LLaMA, Mistral ⁵³

Table 1.3: Comparison of FFN Activation Functions

Activation	Formula/Mechanism	Properties	Computational Cost	Common Models Using It
ReLU	$\max(0, x)$	Sparse, Simple, Non-smooth	Low	Original Transformer ⁴
GELU	$x\Phi(x)$ (approx. $0.5x(1+\tanh[2/\pi(x+0.044715x^3)])$)	Smooth approx. of ReLU	Moderate	BERT, GPT-2/3

SwiGLU	$(\text{Swish}(xW1) \odot xV)W2$ where $\text{Swish} = x\sigma(\beta x)$	Gated, Smooth, Potentially Better Perf.	Higher	PaLM, LLaMA ⁵⁵
GeGLU	$(\text{GELU}(xW1) \odot xV)W2$	Gated, Smooth	Higher	T5 v1.1
SquaredReLU	$(\max(0, xW1 + b1))^2$	Non-gated, Used in some variants	Moderate	PaLM (some layers) ⁶⁸

Table 1.4: Adam vs. AdamW

Optimizer	Weight Decay Implementation	Interaction with Adaptive LR	Typical Use Case/Recommendation for LLMs
Adam	L2 Regularization (penalty added to loss/gradient) ²⁴	Coupled, can interfere	Less common for modern LLMs
AdamW	Decoupled Weight Decay (applied directly to weights after update) ²²	Decoupled, more stable	Preferred for LLM training ¹⁹

2. Data Collection and Preprocessing

The quality and scale of the training data are paramount to the performance of foundational LLMs. Building a suitable dataset involves collecting vast amounts of text from diverse sources and applying rigorous cleaning, filtering, and preprocessing steps.

2.1 Types of Data Required

A typical pretraining corpus for a general-purpose LLM aims for diversity in topic, style, and domain. Common sources include:

- Public Web Data:** Large web crawls like Common Crawl are a primary source, providing petabytes of raw web text.⁶⁹ However, raw crawls require extensive cleaning. Datasets like C4 (Colossal Clean Crawled Corpus)⁷¹ and RefinedWeb⁷¹ represent filtered subsets of Common Crawl.

- **Encyclopedic Knowledge:** Wikipedia is frequently included for its structured, factual content.⁷²
- **Books:** Corpora like BooksCorpus⁷² and Books3⁷⁵ provide long-form narrative and diverse writing styles. However, licensing and copyright issues are significant concerns, particularly with datasets like Books3, which was derived from pirated sources.⁷⁵
- **Code:** Source code from repositories like GitHub⁷⁷ and Q&A sites like Stack Exchange¹⁵ are included to imbue models with coding and technical reasoning abilities.
- **Academic & Scientific Texts:** Sources like arXiv¹⁵ and PubMed Central⁷⁶ provide formal language and domain-specific knowledge.
- **Conversation/Dialogue Data:** Datasets derived from sources like Reddit¹⁵ can help improve conversational abilities, although over-reliance can cause issues.¹⁵
- **Multilingual Data:** For models intended to support multiple languages, datasets like mC4⁷⁷, CC100⁷⁷, OSCAR⁷⁷, CulturaX⁷⁷, and MADLAD-400⁷⁷ are used.

Aggregated datasets like The Pile⁷⁶ and RedPajama⁷² attempt to create balanced, high-quality mixtures from many of these sources. The specific blend and weighting of these sources (data mixture) is a critical, often proprietary, aspect of LLM training.⁷³

2.2 Data Collection and Initial Filtering

Raw data, especially from web crawls, requires initial processing:

- **Crawling:** While large crawls like Common Crawl are available, specialized crawling strategies like Crawl4LLM aim to prioritize pages based on their likely value for LLM pretraining, potentially reducing wasted effort.⁶⁹
- **Language Identification:** Raw web data contains multiple languages. Tools like fastText⁸¹ or langdetect⁸² are used to identify and filter documents based on language (e.g., selecting only English documents).⁷⁰
- **Basic Cleaning:** Removing HTML tags, boilerplate text, navigation menus, and other non-prose content is essential.⁷⁶ Tools like trafilatura⁷¹ are often used for robust text extraction from web pages.

2.3 Data Cleaning and Quality Filtering

Further filtering is applied to enhance data quality:

- **Heuristic Filtering:** Simple rules are often used, such as filtering documents based on length, ratio of symbols to words, presence of "bad words", repetition levels, or perplexity scores against a smaller language model.⁷⁰
- **Model-Based Filtering:** More sophisticated approaches use machine learning

models to classify document quality. This can involve training classifiers on datasets known to be high-quality (e.g., DCLM using OpenHermes/ELI5 data ⁷¹) or using powerful LLMs to score documents based on criteria like educational value (e.g., FineWeb-Edu ⁷¹) or general quality (e.g., Nemotron-CC using classifier ensembles ⁷⁰).

- **Toxicity Filtering:** Explicit filtering is often applied to remove or mitigate toxic, offensive, or harmful content.⁷⁴ This can involve using pre-built classifiers (like Google's Perspective API ⁸³ mentioned in related contexts ⁸⁴) or custom models trained to detect hate speech, harassment, etc..⁸⁴ Defining and detecting toxicity robustly, especially subtle forms, remains a challenge.⁸⁵

2.4 Data Deduplication

Removing duplicate or near-duplicate documents is crucial to prevent models from overfitting to repeated content and to ensure fair evaluation.⁷⁴

- **Exact Deduplication:** Simple methods involve checking for exact matches of documents or large substrings, often using hashing techniques like SHA256.⁸⁸
- **Near-Duplicate Detection:** Identifying documents that are substantially similar but not identical is more complex. Common techniques include:
 - **MinHash + Locality Sensitive Hashing (LSH):** This is a standard approach for large-scale near-duplicate detection.⁸⁶
 1. **Shingling:** Documents are broken down into overlapping sets of n-grams (shingles) of characters or words.⁸⁷
 2. **MinHashing:** Each document's shingle set is converted into a compact signature (a vector of hash values) using the MinHash algorithm. The probability of two signatures matching in a given position approximates the Jaccard similarity of the original shingle sets.⁹¹
 3. **LSH Banding:** The MinHash signatures are divided into multiple "bands". Documents are hashed into buckets based on the hash of each band. Documents that hash to the same bucket in at least one band are considered candidate pairs.⁹¹ This avoids comparing all N² pairs, making the process scalable.⁹⁴ The choice of the number of hash functions, bands, and rows per band involves trade-offs between accuracy, speed, and memory usage.⁹¹
 - **Bloom Filters:** These are space-efficient probabilistic data structures used for set membership testing.⁹⁶ While primarily for exact matches, they can be incorporated into near-duplicate detection workflows, for instance, to quickly check if a MinHash band has been seen before.⁸⁶ LSHBloom is a technique combining LSH with Bloom filters to improve space and potentially runtime

efficiency compared to standard MinHashLSH, at the cost of a small, controllable increase in false positives.⁹⁰

The data processing pipeline is complex and computationally intensive, often requiring distributed processing frameworks. The quality of the final dataset is a direct result of careful filtering and deduplication, balancing the removal of low-quality or redundant data with the preservation of diversity.

2.5 Tokenization

The final preprocessing step is tokenization, which converts the cleaned text into a sequence of integer IDs that the model can process.⁷⁴

- **Byte Pair Encoding (BPE):** A widely used algorithm, popularized by GPT-2.¹⁰⁰ It starts with individual bytes (covering all possible characters in UTF-8) and iteratively merges the most frequent adjacent pair of tokens in the corpus to form new vocabulary entries, until a predefined vocabulary size is reached.⁷⁴ Implementations like minbpe¹⁰⁰ and OpenAI's tiktoken¹⁰³ exist. A common refinement used in GPT models involves using regex patterns to pre-split the text by character categories (letters, numbers, punctuation) before applying BPE merges, preventing merges across these boundaries.¹⁰⁰
- **SentencePiece:** An unsupervised library developed by Google¹⁰⁵ that supports both BPE and the Unigram Language Model algorithm.⁹⁹ Key advantages include training directly from raw sentences without requiring pre-tokenization and handling whitespace consistently by treating it as a regular character (escaped with a meta symbol like U+2581), ensuring lossless detokenization.¹⁰⁵
- **Unigram Language Model:** An alternative subword algorithm (supported by SentencePiece) that initializes with a large vocabulary and iteratively removes tokens that least increase the likelihood of the training data according to a unigram language model, until the target vocabulary size is reached.¹⁰⁵

The choice of tokenization algorithm and vocabulary size affects the average number of tokens per word (compression rate) and can impact downstream performance, especially for multilingual models.⁷⁴

2.6 Tools and Frameworks

Handling terabyte-scale datasets requires specialized tools:

- **Apache Spark:** A distributed computing framework well-suited for large-scale data processing, including cleaning, filtering, deduplication, and tokenization pipelines.⁸⁸ Libraries like Spark NLP provide additional NLP capabilities.¹⁰⁶

- **Petastorm:** An open-source library from Uber for loading large datasets, particularly in Apache Parquet format or from Spark DataFrames, into deep learning frameworks like TensorFlow or PyTorch.¹⁰⁹ It includes a Spark converter API.¹⁰⁹
- **Hugging Face Datasets:** A popular library providing easy access to thousands of datasets and efficient tools for loading, processing, mapping, shuffling, and streaming data in various formats.¹¹¹ It integrates well with Arrow for memory efficiency.
- **MosaicML StreamingDataset (MDS):** A format developed by MosaicML (now Databricks) designed for high-throughput streaming of large datasets during distributed training, offering efficient shuffling and random access.¹¹³ Tools exist to convert data into MDS format and integrate with Databricks Unity Catalog.¹¹³

The increasing sophistication of data processing pipelines, involving multi-stage filtering, advanced near-duplicate detection, and careful tokenization, underscores the critical role data quality plays in building high-performance LLMs. Techniques like MinHashLSH and Bloom Filters offer scalable solutions for deduplication, each presenting different trade-offs in accuracy, speed, and resource usage. Similarly, the choice of tokenization algorithm impacts model efficiency and cross-lingual performance, making it an integral design consideration.

Table 2.1: Common LLM Pretraining Data Sources

Data Source	Description	Typical Size Scale	Content Type	Key Considerations
Common Crawl	Public web crawl data ⁶⁹	Petabytes (Raw)	Web	Needs extensive cleaning/filtering; Permissive Use
The Pile	Curated mix of 22 datasets (web, books, code, academic) ⁷⁶	~825 GiB	Mixed	Diverse, documented; contains Books3 (copyright issue)
C4 /	Cleaned/filtered subsets of	TBs	Web	Higher quality

RefinedWeb	Common Crawl ⁷¹			than raw crawl
Wikipedia	Online encyclopedia ⁷²	~20 GB (English)	Encyclopedic	Factual, structured; relatively small
Books3	Collection of ~170k books from Bibliotik ⁷⁵	~100 GiB	Books	Copyright issues (pirated source) ⁷⁵
GitHub	Public code repositories ⁷⁷	Terabytes	Code	Permissive licenses vary; code quality varies
Stack Exchange	Q&A site data ¹⁵	~70 GB (Pile)	Q&A, Code	Conversational, technical content
arXiv / PubMed	Academic preprints / biomedical literature ¹⁵	~60 GB / ~100 GB	Academic	Formal language, domain knowledge
mC4 / OSCAR	Multilingual web corpora from Common Crawl ⁷⁷	Terabytes	Multilingual Web	Broad language coverage but quality varies
RedPajama	Open replication of LLaMA dataset sources ⁷²	~1.2T Tokens	Mixed	Aims for transparency in data sourcing

Table 2.2: Data Deduplication Techniques

Technique	Mechanism	Detects	Pros	Cons
Exact Hash	Compare cryptographic	Exact Duplicates	Simple, fast for	Misses

	hashes (e.g., SHA256) of documents ⁸⁸		exact matches	near-duplicates
MinHash + LSH	Approximates Jaccard similarity using MinHash signatures; LSH bands reduce comparisons ⁸⁶	Near Duplicates	Scalable, good accuracy for near-duplicates, tunable threshold	More complex, memory/compute intensive than exact hash ⁸⁹
Bloom Filter	Probabilistic check if element (e.g., shingle, hash) is in a set ⁹⁶	Exact Membership (probabilistic)	Very space/time efficient for membership check	False positives possible; not directly for near-duplicates alone ⁹⁷
LSHBloom	Uses Bloom filters within LSH framework to approximate MinHashLSH index ⁹⁰	Near Duplicates	Faster & much less disk space than MinHashLSH ⁹⁰	Small increase in false positives vs MinHashLSH ⁹⁰

Table 2.3: Tokenization Algorithms

Algorithm	Key Principle	Handling of Unknowns/Whitespace	Pros	Cons	Common Models Using It
BPE	Iteratively merge most frequent adjacent token pairs ¹⁰⁰	Can pre-split by regex; whitespace handling depends on pre-tokenization	Simple, effective for common languages	Can be suboptimal for some languages/code; not reversible if pre-tokenized	GPT-2/3/4 ¹⁰⁰
SentencePie	BPE implementation	Treats whitespace	Trains on raw text,		LLaMA,

ce (BPE)	on within SentencePie ce framework ¹⁰⁵	as normal symbol (reversible); handles UTF-8 directly	language-ag nostic, reversible		Mistral
SentencePie ce (Unigram)	Iteratively removes tokens to maximize corpus likelihood under unigram model ¹⁰⁵	Treats whitespace as normal symbol (reversible); handles UTF-8 directly	Trains on raw text, language-ag nostic, reversible, multiple segmentatio ns possible	Can be slower to train than BPE	T5, ALBERT
TikToken	OpenAI's implementati on, likely BPE-based ¹⁰³	Handles special tokens; likely uses regex splitting	Optimized for OpenAI models	Less documented internals	OpenAI GPT models ¹⁰⁴

3. Model Implementation (Code-Level)

Implementing a foundational LLM requires translating the theoretical architecture into functional code using deep learning frameworks. This section focuses on framework choices, a minimal GPT implementation, and memory efficiency techniques crucial for handling large models.

3.1 Frameworks Overview

Several frameworks are commonly used for LLM development, each with its strengths and ecosystem:

- **PyTorch:** A widely adopted framework known for its Pythonic interface, dynamic computation graph (facilitating easier debugging and flexibility), and extensive community support.¹¹⁴ It forms the basis for many higher-level LLM libraries and is often preferred for research due to its ease of iteration.¹¹⁴
- **DeepSpeed:** Developed by Microsoft, DeepSpeed is a library built *on top of* PyTorch, specifically designed to optimize the training and inference of very large models.¹¹⁶ It integrates advanced parallelism techniques (like ZeRO) and memory optimization strategies with minimal code changes required in the base PyTorch model.¹¹⁷

- **Megatron-LM:** An NVIDIA-developed framework, also often used with PyTorch, specializing in efficient Tensor Parallelism (TP) and Pipeline Parallelism (PP) for training massive Transformer models.¹²⁰ It provides highly optimized implementations for NVIDIA GPUs and is often used for state-of-the-art pretraining.¹²⁰ NVIDIA's NeMo framework often utilizes Megatron-LM components.¹²⁰
- **Hugging Face Transformers:** A high-level library providing standardized implementations of numerous Transformer architectures (including GPT variants), pre-trained model weights, tokenizers, and utilities like the Trainer class for simplifying training and fine-tuning workflows.¹²⁵ It abstracts away much of the underlying complexity of PyTorch or TensorFlow.
- **JAX / Flax:** JAX is a Google-developed library for high-performance numerical computation and machine learning research, known for its composable function transformations (grad, jit, vmap, pmap).¹¹⁵ Flax is a popular neural network library built on JAX.¹¹⁵ JAX's strengths lie in its Just-In-Time (JIT) compilation via XLA and its design facilitating efficient parallelization across multiple accelerators (GPUs/TPUs).¹²⁸ While potentially offering superior performance at scale, it has a steeper learning curve due to its functional programming paradigm and explicit state management.¹¹⁴

The choice between these frameworks often involves navigating a trade-off between ease of use and control over low-level optimizations. Hugging Face offers the highest abstraction, PyTorch provides flexibility, JAX offers potential performance advantages at scale (especially on TPUs), and libraries like DeepSpeed and Megatron-LM add specialized scaling capabilities on top of PyTorch.

3.2 Minimal GPT-like Model from Scratch (PyTorch)

Building a minimal GPT model provides foundational understanding.¹³⁰ Using PyTorch, the implementation involves defining modules for each architectural component:

1. **Token Embeddings:** An `nn.Embedding` layer maps input token IDs (integers) to dense vectors of dimension `dmodel`.⁸
2. **Positional Encoding:** A mechanism (either fixed sinusoidal or a learned `nn.Embedding`) adds positional information to the token embeddings.⁸ The sum of token and positional embeddings forms the input to the Transformer blocks.
3. **Transformer Decoder Block (`nn.Module`):** This is the core repeating unit.
 - **Masked Multi-Head Self-Attention:**
 - Define linear layers (`nn.Linear`) to project the input `x` into Query (Q), Key (K), and Value (V) matrices.⁸
 - Split Q,K,V across the attention heads dimension (e.g., reshape from

- (B,T,dmodel) to (B,h,T,dk) where $dk=dmodel/h$).⁸
- Compute attention scores: $S=QKT/dk$.⁸
 - **Apply Causal Mask:** Create a lower-triangular mask (e.g., using `torch.tril`) where future positions are masked out. Apply this mask to the scores, typically by adding $-\infty$ (or a very large negative number) to the masked positions before the softmax.¹² PyTorch's `nn.MultiheadAttention` or `nn.TransformerDecoderLayer` handle this internally via the `attn_mask` or `tgt_mask/tgt_is_causal` arguments.¹³
 - Apply softmax to get attention weights: $P=\text{softmax}(S+\text{mask})$.⁸
 - Compute the output of the attention head: $O=PV$.⁸
 - Concatenate outputs from all heads and pass through a final output linear layer (`nn.Linear`).⁸
- **Residual Connection & Layer Normalization 1:** Add the input x to the attention output and apply `nn.LayerNorm`.²
 - **Feed-Forward Network (FFN):** Implement the two linear layers with an intermediate activation function (e.g., `nn.ReLU`, `nn.GELU`, or a `SwiGLU` implementation).⁴
 - **Residual Connection & Layer Normalization 2:** Add the input to the FFN sub-layer to its output and apply `nn.LayerNorm`.²
4. **Model Body:** Stack L instances of the Transformer Decoder Block sequentially (`nn.Sequential` or a custom loop).¹³¹
 5. **Language Model Head:** Apply a final `nn.LayerNorm` (optional, sometimes done within the block) and a final `nn.Linear` layer that projects the output dimension $dmodel$ back to the vocabulary size V , producing the final logits.¹¹

The forward method orchestrates the flow of input token IDs through these layers to produce the output logits for predicting the next token. Educational implementations like `minbpe`¹⁰⁰ and `nanoGPT`¹³¹ provide excellent references.

Python

Simplified PyTorch Attention Block Example (Single Head, No Masking for brevity)

```
import torch
import torch.nn as nn
import torch.nn.functional as F
import math
```

```

class SimpleSelfAttention(nn.Module):
    def __init__(self, embed_dim, head_dim):
        super().__init__()
        self.head_dim = head_dim
        self.q_proj = nn.Linear(embed_dim, head_dim, bias=False)
        self.k_proj = nn.Linear(embed_dim, head_dim, bias=False)
        self.v_proj = nn.Linear(embed_dim, head_dim, bias=False)
        # Output projection omitted for simplicity

    def forward(self, x):
        B, T, E = x.shape # Batch, Time, Embedding dim
        q = self.q_proj(x) # (B, T, H) H=head_dim
        k = self.k_proj(x) # (B, T, H)
        v = self.v_proj(x) # (B, T, H)

        # Attention scores
        scores = q @ k.transpose(-2, -1) / math.sqrt(self.head_dim) # (B, T, H) @ (B, H, T) -> (B,
T, T)

        # --- Causal Mask would be applied here ---
        # Example:
        # mask = torch.tril(torch.ones(T, T, device=x.device)).unsqueeze(0)
        # scores = scores.masked_fill(mask == 0, float('-inf'))

        attn_weights = F.softmax(scores, dim=-1) # (B, T, T)
        output = attn_weights @ v # (B, T, T) @ (B, T, H) -> (B, T, H)
        return output

# --- Minimal Decoder Block would combine Attention, FFN, LayerNorm, Residuals ---

```

3.3 GPU Tensor Parallelism and Sharding (Conceptual Introduction)

Training models with hundreds of billions or trillions of parameters necessitates distributing the model itself across multiple GPUs, as even a single layer's weights might exceed the memory of one device. **Tensor Parallelism (TP)**, pioneered by Megatron-LM¹²², addresses this by splitting large weight matrices (e.g., in `nn.Linear` layers or attention projections) across GPUs.¹²³ For example, a weight matrix might be split column-wise across GPUs, and the input activation is broadcasted. Each GPU computes a partial result, and communication (e.g., all-reduce) is used to combine these results.¹²² This allows computation on tensor shards in parallel but requires

high-speed interconnects like NVLink.¹²³ Frameworks like Megatron-LM and DeepSpeed handle the complexities of this sharding and communication. (Detailed strategies are covered in Section 4).

3.4 Techniques for Memory Efficiency

Beyond parallelism, several techniques are vital for managing memory consumption during training:

- **Gradient Checkpointing (Activation Checkpointing):** Standard backpropagation requires storing activations from the forward pass to compute gradients in the backward pass. For deep networks, this activation memory can be prohibitive. Gradient checkpointing trades computation for memory by *not* storing intermediate activations for designated "checkpointed" layers during the forward pass.¹³⁶ During the backward pass, when these activations are needed, they are recomputed on-the-fly.¹³⁶ This significantly reduces activation memory usage, allowing larger models or batch sizes to fit, at the cost of increased training time (typically around 20% slower).¹³⁶ PyTorch and other frameworks provide utilities to easily enable gradient checkpointing.
- **Mixed Precision Training (FP16/BF16):** This involves using 16-bit floating-point formats (FP16 or BF16) for storing weights, activations, and performing computations, instead of the standard 32-bit (FP32).¹³⁹ This halves the memory required for these components and can significantly accelerate computation on hardware with specialized units like NVIDIA's Tensor Cores.¹³⁹
 - **FP16:** Offers higher precision than BF16 but has a much smaller dynamic range. This necessitates **Loss Scaling**, where the loss value is multiplied by a scaling factor before backpropagation to prevent small gradient values from becoming zero (underflowing) in FP16, and gradients are unscaled before the weight update.¹³⁹
 - **BF16 (Brain Float):** Has the same dynamic range as FP32 but lower precision than FP16.¹⁴⁰ Its wider range often eliminates the need for loss scaling, simplifying training.¹⁴⁰ Conversion between BF16 and FP32 is also computationally cheaper.¹⁴⁰ BF16 is generally preferred on hardware that supports it (e.g., NVIDIA A100/H100).
 - Frameworks like PyTorch (torch.cuda.amp) and DeepSpeed provide automatic mixed precision (AMP) capabilities, managing the casting between precisions and loss scaling automatically.¹³⁹
- **FlashAttention:** As discussed in Section 1, FlashAttention provides memory savings by avoiding the materialization of the $N \times N$ attention matrix, making memory usage linear instead of quadratic in sequence length.⁶¹ This is crucial for

enabling long context lengths during training and inference.

- **LoRA (Low-Rank Adaptation):** While primarily a fine-tuning technique, LoRA drastically reduces memory requirements during adaptation.¹⁴² By freezing the large base model and only training small low-rank matrices (A and B such that the update $\Delta W=BA$), the number of trainable parameters, gradients, and optimizer states is significantly reduced.¹⁴³ This allows fine-tuning very large models on consumer-grade GPUs.

These memory-saving techniques are often essential, not merely optimizations, for training state-of-the-art LLMs. The sheer size of model parameters, gradients, and optimizer states (especially for optimizers like AdamW) quickly exceeds the memory capacity of even high-end accelerators.¹⁴⁴ Techniques like mixed precision and gradient checkpointing directly address these bottlenecks, making large-scale training feasible.

Table 3.1: Comparison of Deep Learning Frameworks for LLM Training

Framework	Key Features	Programming Paradigm	Ecosystem/ Tooling	Performance Profile	Common Use Cases
PyTorch	Dynamic graph, Pythonic API, Large community, AMP ¹¹⁴	Imperative	Extensive (Hugging Face, DeepSpeed, TorchServe)	Fast iteration, Scalable	Research, Production
TensorFlow	Static/Dynamic graph (TF2), Keras API, Production focus (TF Serving, TFLite) ¹¹⁴	Declarative/Imperative	Mature deployment tools, TensorBoard	Robust production	Production Deployment
JAX / Flax	Function transforms (jit, grad, vmap,	Functional	Growing (Optax, Orbax), Strong TPU	High performance at scale	Research, Large-scale HPC

	pmap), XLA compilation ¹²⁸		support		
DeepSpeed	ZeRO, 3D Parallelism, MII (Inference), Built on PyTorch ¹¹⁶	Imperative (via PyTorch)	Training/Inference Optimization Library	Scalability focus	Very Large Model Training
Megatron-LM	Tensor/Pipeline Parallelism, Optimized kernels, NVIDIA focus ¹²⁰	Imperative (via PyTorch)	Specialized Training Framework (often via NeMo)	Scalability focus	SOTA Pretraining
HF Transformers	Model Hub, Standardized APIs, Trainer class, Tokenizers ¹²⁵	High-Level API	Large ecosystem of models, datasets, tools	Ease of use, Fine-tuning	Fine-tuning, Prototyping

Table 3.2: Memory Efficiency Techniques

Technique	Mechanism	Primary Benefit	Trade-offs
Gradient Checkpointing	Recompute activations during backward pass instead of storing all ¹³⁶	Activation Memory	Slower training (~20%) ¹³⁶
Mixed Precision (FP16)	Use 16-bit floats for weights/activations; requires loss scaling ¹³⁹	Memory & Compute	Potential underflow/overflow, requires scaling

Mixed Precision (BF16)	Use 16-bit brain floats (FP32 range, FP16 precision) ¹⁴⁰	Memory & Compute	Lower precision than FP16, avoids loss scaling
FlashAttention	IO-aware attention with tiling, avoids materializing N x N matrix ⁶⁰	Attention Memory & Speed	Exact attention, requires specific hardware support
LoRA (Fine-tuning)	Train low-rank updates instead of full weights ¹⁴²	Parameter/Optimizer Memory (Fine-tuning)	Reduced trainable params, slight accuracy trade-off

4. Training Infrastructure & Systems Engineering

Pretraining foundational LLMs demands cutting-edge infrastructure capable of handling immense computational loads, massive datasets, and prolonged training runs. This section details the hardware, software, and systems engineering principles required.

4.1 Compute Architecture for Pretraining

- **GPU Clusters:** State-of-the-art LLM training relies heavily on large clusters of high-performance GPUs. NVIDIA's A100 and H100 Tensor Core GPUs are the current industry standards, often deployed in nodes containing 8 or 16 GPUs each.¹⁴⁶ Future architectures like Blackwell (B200, GB200) promise further performance gains.¹⁴⁸ These GPUs provide the massive parallel processing power needed for matrix multiplications and other core deep learning operations.
- **Interconnects:** Efficient communication between GPUs (intra-node) and between nodes (inter-node) is critical for distributed training performance. Bottlenecks here can severely limit overall throughput.¹⁴⁶
 - **Intra-Node (GPU-to-GPU):** NVIDIA NVLink and NVSwitch provide high-bandwidth, low-latency direct connections between GPUs within a server.¹⁴⁸ Bandwidths can reach 900 GB/s per GPU or higher ¹⁵², essential for Tensor Parallelism strategies.
 - **Inter-Node (Server-to-Server):** High-speed networking fabrics like InfiniBand or Ethernet with RDMA (Remote Direct Memory Access) capabilities (e.g., RoCEv2) are required for efficient data transfer between nodes in the cluster.¹⁴⁶ RDMA allows GPUs in one node to directly access memory in another node, bypassing the CPU and reducing latency.
- **Integrated Systems (NVIDIA DGX SuperPOD):** NVIDIA offers turnkey AI

supercomputing solutions like the DGX SuperPOD.¹⁴⁹ These systems integrate optimized compute (DGX servers with A100, H100, or B200 GPUs), high-speed networking (NVLink/NVSwitch and InfiniBand), high-performance storage (often Lustre-based¹⁴⁶), and a software stack including cluster management and AI libraries.¹⁴⁸ DGX Cloud provides similar capabilities as a managed service in partnership with cloud providers.¹⁴⁶ Deploying DGX systems on-premise often involves Kubernetes for orchestration.¹⁵⁴

4.2 Orchestration and Job Scheduling

Managing large GPU clusters requires robust orchestration and scheduling systems:

- **Slurm:** A widely used workload manager in High-Performance Computing (HPC) environments.¹⁵⁸ Slurm excels at managing batch jobs, allocating resources (nodes, GPUs, memory), and handling job queues.¹⁵⁸ It integrates well with MPI and is often used in DGX SuperPOD deployments.¹⁴⁷ However, it's traditionally less suited for dynamic scaling or managing containerized, service-oriented workloads compared to Kubernetes.¹⁶⁰
- **Kubernetes (K8s):** The standard for container orchestration, offering flexibility, dynamic scaling, and self-healing capabilities.¹⁵⁸ While not originally designed for HPC batch jobs, it's increasingly used for AI/ML workloads, including on-premise DGX clusters.¹⁵⁴ Managing stateful applications and specialized hardware like GPUs requires extensions such as the NVIDIA GPU Operator and potentially specialized schedulers (e.g., Volcano).¹⁵⁴
- **Slurm vs. Kubernetes:** The choice often depends on the environment and workload type. Slurm is mature and optimized for large, long-running batch training jobs typical in HPC. Kubernetes offers more flexibility for dynamic workloads, containerization benefits, and integration with cloud-native tooling, but may require more effort to configure optimally for large-scale distributed training.¹⁵⁸

4.3 Distributed Training Strategies

No single GPU can hold or train the largest LLMs. Distributed training strategies parallelize the workload across multiple GPUs and nodes.

- **Data Parallelism (DP / DDP):** The simplest approach. The model is replicated on each GPU (worker). The global data batch is split, with each worker processing a mini-batch.¹¹⁶ Gradients are computed locally and then aggregated (e.g., via all-reduce) across all workers before updating the model weights synchronously.¹¹⁶ PyTorch's DistributedDataParallel (DDP) is the standard implementation.¹²⁴ **Limitation:** The entire model, gradients, and optimizer states

must fit in the memory of each individual GPU.¹¹⁶

- **Tensor Parallelism (TP):** Splits individual layers or tensors *within* the model across multiple GPUs, typically within a single node.¹²³ For a linear layer $Y=XA$, the weight matrix A might be split column-wise $[A1,A2]$, with X broadcasted. Each GPU computes YAi , and results are gathered.¹²³ Requires frequent communication between participating GPUs, necessitating high-speed interconnects like NVLink.¹²³ Megatron-LM provides a well-known implementation.¹²²
- **Pipeline Parallelism (PP):** Splits the model *vertically* across layers, assigning consecutive layer chunks to different GPUs (stages), potentially spanning multiple nodes.¹²³ The input batch is divided into micro-batches, which are processed sequentially through the pipeline stages.¹²² Activations are passed from one stage to the next. This reduces memory per GPU but introduces potential "pipeline bubbles" (idle time) that micro-batching aims to minimize.¹²² Communication occurs between stages.
- **ZeRO (Zero Redundancy Optimizer):** A DeepSpeed technique that enhances Data Parallelism by partitioning model states (optimizer states, gradients, and optionally parameters) across data-parallel workers.¹¹⁶
 - Stage 1: Partitions optimizer states.
 - Stage 2: Partitions optimizer states and gradients.
 - Stage 3: Partitions optimizer states, gradients, and model parameters. Each GPU only holds a fraction of the partitioned state, drastically reducing memory per GPU.¹¹⁹ Stage 3 allows training models whose parameters alone exceed single-GPU memory.¹¹⁶ ZeRO introduces communication steps to gather the necessary states/parameters when needed. ZeRO-Infinity extends this by offloading to CPU/NVMe.¹⁴¹
- **Fully Sharded Data Parallelism (FSDP):** PyTorch's native implementation of sharded data parallelism, conceptually similar to ZeRO Stage 3.¹¹⁶ It shards model parameters, gradients, and optimizer states across data-parallel ranks.¹⁶⁴ It offers various sharding strategies (e.g., FULL_SHARD, SHARD_GRAD_OP, HYBRID_SHARD) and supports CPU offloading.¹⁶⁴
- **Sequence Parallelism (SP):** An optimization used in conjunction with TP or PP. It splits the activations along the sequence length dimension across the TP or PP group, reducing activation memory further, especially beneficial for long sequences.¹³⁵ Requires communication (e.g., all-gather) before operations that need the full sequence context within a layer.

These strategies are often combined (e.g., "3D Parallelism" = DP + PP + TP) to train the largest models across massive clusters.¹¹⁹ The optimal combination depends heavily on the model architecture, scale, and the specific hardware interconnect

topology. Infrastructure limitations, particularly network bandwidth and latency, are critical; a slow network can negate the benefits of aggressive sharding strategies like ZeRO-3 or FSDP, making simpler DDP potentially faster in some scenarios.¹¹⁶

4.4 Storage and Logging

- **Storage Systems:** Training requires fast access to petabyte-scale datasets and frequent checkpointing.
 - **Object Storage (S3, GCS, Azure Blob):** Suitable for storing raw datasets, intermediate processed data, and final model artifacts due to its scalability and durability.¹⁶⁵ S3-compatible solutions like MinIO offer on-premise/local alternatives.¹⁶⁵
 - **Parallel File Systems (PFS):** High-performance file systems like Lustre¹⁶⁸ and IBM Spectrum Scale (GPFS)¹⁶⁸ are common in HPC and DGX SuperPOD environments.¹⁴⁶ They provide high-bandwidth, parallel access needed during training but can face challenges with metadata operations or small files.¹⁷⁰ Modern alternatives like BeeGFS¹⁶⁸ and WEKA¹⁷⁰ also exist.
- **Dataset Sharding:** Large datasets are typically split into smaller files or shards (e.g., 1GB-10GB chunks).¹⁷² This allows different nodes/workers in distributed training to read different parts of the dataset concurrently, improving I/O efficiency.¹⁷³ Formats like MosaicML's StreamingDataset (MDS) are designed for this.¹¹³
- **Logging and Monitoring:** Tracking experiments is crucial for reproducibility and debugging. Tools like Weights & Biases (W&B)¹⁷⁴, MLflow¹²⁶, and TensorBoard¹⁷⁴ are used to log metrics (loss, accuracy, perplexity), hyperparameters, system resource utilization (GPU memory/compute), visualize training progress, and store model artifacts.¹⁷⁵

4.5 Cost Estimation

Estimating the cost of pretraining is complex but crucial for planning.

- **Key Factors:** Compute cost (GPU hours * cost per hour), dataset size (tokens processed), model size (parameters), hardware efficiency (MFU - Model FLOPs Utilization), training stability, and infrastructure choices (cloud vs. on-prem) all influence the final cost.¹⁷⁹
- **Estimates:** Public estimates vary widely. MosaicML suggested a compute-optimal 30B parameter model could be trained for <\$500k in 2023.¹⁸³ Other estimates place a 7B model around \$85k.¹⁸⁴ Training GPT-3 likely cost millions (\$4.6M - \$12M)¹⁸⁰, and GPT-4 significantly more.¹⁸⁰ Costs generally scale with model size and the number of tokens processed, following patterns related to scaling laws

(see Section 9).

The complexity of LLM training necessitates a holistic view of infrastructure. The interplay between compute power, high-speed communication fabrics, performant storage, and sophisticated distributed training software (like DeepSpeed or Megatron-LM integrated with PyTorch/JAX) is critical. Achieving high MFU (e.g., the 55.2% reported by MegaScale¹⁸⁵) requires careful co-design and optimization across this entire stack. The choice of distributed strategy (DDP, ZeRO, FSDP, TP, PP) involves balancing memory savings against communication costs, heavily influenced by the underlying hardware capabilities, particularly interconnect bandwidth.

Table 4.1: Comparison of Distributed Training Strategies

Strategy	What is Sharded	Key Benefit	Communication Overhead	Implementation Complexity	Typical Use Case
DDP	Data (Gradients aggregated)	Simple, Widely supported ¹⁶²	Low (Gradients only)	Low	Model fits on single GPU
ZeRO-1	Optimizer States	Reduces optimizer memory ¹¹⁶	Low	Moderate (via DeepSpeed)	Modest memory saving
ZeRO-2	Optimizer States, Gradients	Reduces optimizer & gradient memory ¹¹⁶	Moderate (ReduceScatter)	Moderate (via DeepSpeed)	Significant memory saving (model fits GPU)
ZeRO-3 / FSDP	Params, Gradients, Optimizer States ¹¹⁶	Fits models > GPU memory ¹¹⁶	High (AllGather)	High (DeepSpeed/PyTorch)	Very large models
Tensor (TP)	Model Weights/Activations	Reduces memory per layer, faster	High (Intra-node)	High (Megatron/D)	Very large layers, needs

	(Intra-layer) ¹³⁵	compute		DeepSpeed)	NVLink
Pipeline (PP)	Model Layers (Inter-layer) ¹³⁵	Fits models > node memory	Moderate (Inter-stage)	High (Megatron/DeepSpeed)	Very deep models, spans nodes

Table 4.2: LLM Training Infrastructure Components

Component	Examples	Key Role in LLM Training
Compute	NVIDIA H100/A100/B200 GPUs, DGX Systems ¹⁴⁸	Massively parallel computation for training steps
Intra-Node Interconnect	NVLink, NVSwitch ¹⁴⁸	High-speed GPU-to-GPU communication (for TP, ZeRO/FSDP)
Inter-Node Interconnect	InfiniBand, Ethernet (RoCE) ¹⁴⁶	High-speed node-to-node communication (for DP, PP, ZeRO/FSDP)
Storage	Parallel FS (Lustre, GPFS), Object Storage (S3) ¹⁶⁵	Storing massive datasets & checkpoints, high I/O throughput
Orchestration	Slurm, Kubernetes ¹⁵⁸	Resource management, job scheduling, cluster coordination
Monitoring/Logging	Weights & Biases, MLflow, TensorBoard ¹⁷⁴	Tracking experiments, metrics, debugging, visualization

Table 4.3: Estimated Training Cost Ranges (Illustrative, Highly Variable)

Model Size (Params)	Estimated Cost Range (USD)	Key Assumptions / Notes
---------------------	----------------------------	-------------------------

~1B	\$10k - \$50k+	Depends heavily on tokens, efficiency. MosaicML estimate for 1.3B was in 'thousands'. ¹⁸³
~7B	\$50k - \$200k+	HN estimate ~\$85k. ¹⁸⁴ MosaicML 6.7B in 'tens of thousands'. ¹⁸³
~13B	\$100k - \$500k+	MosaicML estimate in 'hundreds of thousands'. ¹⁸³
~30B	\$400k - \$1M+	MosaicML estimate <\$500k for compute-optimal 30B. ¹⁸³
~70B	\$500k - \$2M+	MosaicML estimate in 'hundreds of thousands'. ¹⁸³ Llama 2 70B likely cost millions.
175B+ (GPT-3 scale)	\$4M - \$15M+	GPT-3 estimated \$4.6M-\$12M. ¹⁸⁰ Costs depend heavily on compute-optimality.

Note: Costs are highly dependent on hardware generation (A100 vs H100 vs B200), cloud provider pricing vs on-prem amortization, training efficiency (MFU), dataset size/quality, and training stability.

5. Training the Model

The process of training an LLM involves iterating through the dataset, computing losses, updating model weights, and carefully monitoring convergence over potentially weeks or months. Stability and efficiency are paramount.

5.1 Curriculum Learning Strategies

Instead of randomly sampling data throughout training, curriculum learning organizes the training data, typically progressing from easier to more complex examples.¹⁸⁶ The motivation is to improve convergence speed and generalization by mimicking human learning patterns.¹⁸⁷ For LLM pretraining, various strategies exist, though their effectiveness can be context-dependent¹⁸⁹:

- **Vocabulary Curriculum:** Start training with a small vocabulary (e.g., characters or bytes) and progressively expand it by merging predictable token sequences based on model entropy or other metrics.¹⁸⁶ This allows the model to focus compute on harder-to-predict parts of the sequence.
- **Data Complexity/Difficulty Ordering:** Sequence data based on metrics like document length, vocabulary rarity¹⁸⁸, perplexity difference between strong/weak models¹⁹¹, gradient norms¹⁹², or model uncertainty (least certainty).¹⁹² Frameworks like MATES use auxiliary models to predict data influence and select the most beneficial samples for the current training stage.¹⁹²
- **Task/Domain Sequencing:** In continual pretraining or domain adaptation scenarios, models might first be trained on general data before moving to specialized domains (e.g., legal¹⁸⁷, medical) or tasks of increasing complexity.¹⁹³

While promising, designing effective curricula for large-scale pretraining remains an active research area.¹⁸⁹

5.2 The Training Loop

The core training loop iterates over the dataset, performing the following steps for each batch (or accumulation cycle):

1. **Data Loading:** Fetch a batch of tokenized sequences from the sharded dataset.
2. **Forward Pass:** Pass the input batch through the model to obtain logits for each token position.¹⁹⁴
3. **Loss Calculation:** Compute the cross-entropy loss between the predicted logits and the target next tokens.¹⁹⁴ If using gradient accumulation, normalize the loss by the number of accumulation steps.¹⁹⁵
4. **Backward Pass:** Compute gradients of the loss with respect to the model parameters using backpropagation.¹⁹⁴ If using gradient accumulation, these gradients are summed with those from previous micro-batches.¹⁹⁵
5. **Optimizer Step:** After the required number of accumulation steps, the optimizer (e.g., AdamW) updates the model weights using the accumulated gradients and its internal states (momentum, variance).¹⁹⁵ Optimizer states are also updated.¹⁹⁷
6. **Zero Gradients:** Reset the gradients before the next accumulation cycle.¹⁹⁵
7. **Learning Rate Update:** Adjust the learning rate according to the chosen schedule (e.g., cosine decay after warmup).

Key metrics tracked per step or globally include **token throughput** (tokens processed per second per GPU), which measures training speed¹⁹⁵, and the number of **steps per epoch** or total training steps/tokens processed.¹⁹⁵

5.3 Batch Size Scaling and Convergence

The choice of batch size significantly impacts training dynamics, convergence speed, and stability.¹⁹⁹

- **Effective Batch Size:** The global batch size (mini-batch size per GPU * number of GPUs * gradient accumulation steps) determines how much data informs each weight update.
- **Learning Rate Interaction:** The optimal learning rate is strongly coupled with the batch size.¹⁹⁹ While SGD often follows a linear scaling rule (larger batch -> proportionally larger LR), the relationship for Adam-style optimizers used in LLMs is more complex. Theoretical and empirical studies suggest the optimal LR for Adam might first increase and then *decrease* as batch size grows beyond a certain point, before potentially converging to a stable value for very large batches.¹⁹⁹ Tuning AdamW's β_2 hyperparameter can also be important, especially at smaller batch sizes.²⁰²
- **Large Batch Training:** Essential for efficiency in large-scale distributed settings²⁰¹, but requires careful tuning. Very large batches can sometimes lead to poorer generalization (converging to sharper minima) compared to smaller batches which introduce more noise into the gradient estimates.²⁰⁰ They can also increase the risk of instability if the learning rate is not adjusted appropriately.⁶⁸ Gradient accumulation allows simulating large batch sizes without exceeding GPU memory limits.¹⁹⁵

5.4 Convergence Monitoring

Continuously monitoring training progress is vital to ensure the model is learning effectively and to detect issues early.

- **Loss Curves:** Plotting training loss and validation loss over training steps is fundamental.²⁰³ A decreasing trend indicates learning. Signs of problems include:
 - *Plateauing:* Loss stops decreasing, suggesting the learning rate might be too small or the model has converged.
 - *Increasing Validation Loss:* Indicates overfitting; the model is memorizing the training data but not generalizing.²⁰³
 - *Loss Spikes:* Sudden jumps in loss can indicate instability, often related to gradients or learning rate.²⁰⁷
 - *Divergence:* Loss increases uncontrollably.²⁰³
- **Perplexity:** Often used alongside loss, perplexity ($PPL = \exp(\text{CrossEntropyLoss})$) measures how well the model predicts the validation set.¹⁷ Lower perplexity is better.⁸³ While useful, it may not perfectly correlate with performance on specific downstream tasks.²⁰⁹

- **Evaluation Sets:** Periodically evaluating the model on held-out validation sets or standard benchmarks (e.g., subsets of The Pile⁷⁶, HELM²¹⁰, MMLU²¹⁰, BIG-bench²¹⁰) provides insights into generalization capabilities.²⁰³
- **Early Stopping:** A common regularization technique where training is halted if performance on a validation set (measured by loss or perplexity) stops improving or starts worsening for a predefined number of steps (patience).²⁰³ This prevents overfitting and saves computational resources.²¹³ Monitoring post-hoc metrics might inform better stopping points than raw validation loss alone.²¹⁴

Monitoring requires a combination of metrics. While loss and perplexity track the primary training objective, periodic evaluation on downstream tasks or specific capabilities provides a more holistic picture of model progress, as pretraining loss doesn't always perfectly predict final task performance.²⁰⁹

5.5 Common Pitfalls and Stability Issues

Training LLMs at scale is prone to various issues:

- **Divergence / Instability:** The training loss increases indefinitely, often due to excessively high learning rates, large batch sizes without proper LR scaling, or numerical errors.⁶⁸ Requires careful hyperparameter tuning.
- **NaNs (Not-a-Number):** Loss or gradients become NaN, halting training.²¹⁷ This can result from numerical overflow (e.g., with exploding gradients or in mixed precision without proper scaling), or invalid operations like $\log(0)$ or $0/0$.
- **Vanishing Gradients:** Gradients become near-zero during backpropagation, especially in very deep networks, preventing weight updates in earlier layers.²¹⁷ Mitigated by ReLU activations, normalization layers (LayerNorm/RMSNorm), residual connections, and careful initialization.²¹⁹
- **Exploding Gradients:** Gradients become excessively large, causing unstable weight updates and divergence.²¹⁸ Often caused by large weights or high learning rates.²¹⁷ Mitigated by **Gradient Clipping** (rescaling gradients if their norm exceeds a threshold)²⁰⁶, weight regularization (AdamW), normalization, and proper initialization.²¹⁸
- **Loss Spikes:** Sudden, sharp increases in the loss value during otherwise stable training.²⁰⁷ Can sometimes be recovered from, but may indicate underlying numerical instability or issues with specific data batches. Mitigation might involve temporarily reducing the learning rate, skipping problematic data batches, or restarting from a recent checkpoint.⁶⁸

Successfully navigating LLM training requires robust monitoring, careful hyperparameter tuning (especially learning rate and batch size, considering their

complex interplay with AdamW ¹⁹⁹), and employing techniques like gradient clipping and appropriate initialization to maintain stability over long training runs.¹⁸⁵

Table 5.1: Curriculum Learning Strategies for Pretraining

Strategy	Description	Potential Benefits	Challenges/Status
Vocabulary Curriculum	Start with small vocab, expand based on model predictability/entropy ¹⁸⁶	Improved efficiency, better scaling	Active research, optimal expansion strategy?
Data Complexity Metric	Order data by length, perplexity difference, influence score, etc. ¹⁸⁸	Faster convergence, better generalization	Defining effective complexity metrics is hard
Task/Domain Progression	Start general, move to specific domains/tasks ¹⁸⁷	Better adaptation to target domain	Primarily for continual learning/fine-tuning stages
Certainty/Gradient Norm Based	Prioritize data based on model uncertainty or gradient magnitude ¹⁹²	Focus on 'learnable' examples	Effectiveness for LLM pretraining debated

Table 5.2: Common LLM Training Pitfalls

Pitfall	Symptoms	Common Causes	Mitigation Techniques
Divergence	Loss consistently increases ²¹⁶	High LR, large batch size, numerical instability	Reduce LR, tune batch size/LR scaling, gradient clipping ²¹⁸
NaNs	Loss/Gradients become NaN ²¹⁷	Exploding gradients, numerical overflow	Gradient clipping, check mixed

		(mixed precision), invalid ops	precision scaling, check data/ops
Vanishing Gradients	Slow/stalled learning, early layer weights unchanged ²¹⁷	Deep network, sigmoid/tanh activations, poor initialization	ReLU/variants, Norm layers, ResNets, proper initialization ²¹⁹
Exploding Gradients	Oscillating/NaN loss, large weight updates ²¹⁷	High LR, large weights, recurrent computations	Gradient Clipping ²¹⁹ , Weight Regularization (AdamW), Norm layers, Initialization ²¹⁸
Loss Spikes	Sudden, temporary large increase in loss ²⁰⁷	Data issues, numerical instability, LR schedule interaction	Reduce LR temporarily, skip data batch, restart from checkpoint ⁶⁸

6. Fine-tuning, RLHF, and Safety Layers (Optional)

While pretraining endows LLMs with broad knowledge and capabilities, further steps are often necessary to align their behavior with specific tasks, human preferences, and safety requirements.

6.1 Instruction Tuning

Instruction tuning adapts a pre-trained LLM to better follow natural language instructions and perform well on unseen tasks in a zero-shot manner.²²² It involves supervised fine-tuning on a dataset composed of examples formatted as (instruction, input (optional), output) triplets.²²²

- **Datasets:** Numerous datasets have been created for instruction tuning, varying in source, size, language, and format:
 - *Based on Existing NLP Tasks:* Collections like Natural Instructions ²²⁴, P3 (Public Pool of Prompts) ²²⁴, and FLAN (Fine-tuned Language Net) ²²² reformat existing NLP datasets into instruction formats.
 - *LLM-Generated Instructions:* Techniques like Self-Instruct ²²⁵ use a powerful LLM (e.g., GPT-3/4) to generate new instruction-output pairs, often starting from a small set of human-written seed examples. Alpaca ²²³ is a well-known example generated using text-davinci-003.
 - *Human-Generated/Curated Instructions:* Datasets like Dolly ²²⁵ (created by Databricks employees) and OpenAssistant ²²⁵ (crowd-sourced) contain

instructions and responses written entirely by humans, often resulting in higher quality but at greater cost. ShareGPT data consists of real user conversations scraped from a website, offering multi-turn interactions but with potential privacy and licensing concerns.²²⁵

- **Key Considerations:** The choice of instruction dataset significantly impacts the resulting model's capabilities. Factors include data source (human vs. AI-generated, influencing quality and potential bias), data quality (filtering and curation efforts), license (some datasets derived from proprietary models inherit restrictive licenses), and whether the data includes multi-turn dialogues (important for chatbot training).²²⁵

6.2 Reinforcement Learning from Human Feedback (RLHF)

RLHF is a technique used to further align LLMs with complex human preferences, values, and desired behaviors (e.g., helpfulness, honesty, harmlessness) that are difficult to capture solely through supervised instruction tuning.²²⁸ It was notably used for models like InstructGPT/ChatGPT²²⁹ and Claude.²³² The process typically involves three stages:

1. **Supervised Fine-Tuning (SFT) (Optional but Recommended):** Start with a pre-trained LLM and fine-tune it on a high-quality instruction dataset (as described above). This provides a good initial policy that can already follow instructions reasonably well.²²⁸
2. **Reward Model (RM) Training:**
 - *Collect Preference Data:* Generate multiple (typically two) responses from the SFT model for a diverse set of prompts. Human labelers then compare these responses and indicate which one they prefer based on predefined criteria (e.g., helpfulness, harmlessness).²³¹ This creates a dataset of preference pairs (prompt,yw,yl), where yw is the preferred (winning) response and yl is the rejected (losing) response.
 - *Train Reward Model:* Train a separate model (the RM), often initialized from the SFT model, to predict a scalar reward score reflecting human preference.²²⁸ The RM takes a prompt and a response as input and outputs a reward. It's trained to assign a higher reward to yw than to yl for a given prompt. A common loss function is based on the Bradley-Terry model, maximizing the log-likelihood of the preferences, often expressed as minimizing $-\log(\sigma(r(x,yw)-r(x,yl)))$, where r is the reward model's score and σ is the sigmoid function.²³⁷
3. **RL Fine-Tuning (Policy Optimization):**
 - *RL Setup:* Treat the LLM (initialized from the SFT model) as the RL policy

agent. The "action" is generating a response token by token. The "reward" for a completed response is given by the trained RM.²³¹

- *Optimization:* Use an RL algorithm, most commonly Proximal Policy Optimization (PPO)²³¹, to fine-tune the LLM policy. The LLM generates responses to prompts, the RM scores these responses, and PPO updates the LLM's weights to maximize the expected reward from the RM.²³¹
- *PPO Details:* PPO is an on-policy algorithm that optimizes the policy (LLM) using a clipped surrogate objective function. This clipping mechanism prevents the policy from changing too drastically in one update step, ensuring more stable training compared to simpler policy gradient methods.²³⁸ PPO typically involves estimating an advantage function (how much better an action is than average) often using a separate value function (critic).²⁴⁰ A KL divergence penalty term is usually added to the PPO objective to keep the updated policy close to the original SFT policy, preventing the model from deviating too much and losing its general capabilities.²³²
- *Libraries:* The Hugging Face TRL (Transformer Reinforcement Learning) library provides implementations for RLHF, including the PPOTrainer.²⁴²

RLHF allows for fine-grained control over model behavior based on nuanced human judgments, but it is complex, computationally intensive, and sensitive to the quality of human feedback and the reward model.²³⁵

6.3 Safety Layers and Mechanisms

Ensuring LLMs behave safely and reliably requires multiple layers of defense, applied during training and inference.

- **Alignment During Training:** Techniques like RLHF and Constitutional AI are fundamental for instilling desired behaviors (e.g., refusing harmful requests) during the training process itself.²³² Research also investigates identifying and potentially reinforcing specific "safety layers" within the model's architecture.²⁴⁸
- **Moderation Filters:** Input prompts and model outputs can be passed through separate moderation models or filters designed to detect and block harmful, toxic, or policy-violating content.⁸⁴ These can be based on classifiers or keyword/pattern matching.⁸⁴
- **Red Teaming:** Proactive adversarial testing where humans or automated methods attempt to find vulnerabilities or elicit undesirable behavior from the model.²⁴⁹ This helps identify weaknesses (e.g., susceptibility to jailbreaks) that may not be apparent through standard evaluations.²⁴⁹
- **Jailbreaking and Prompt Injection Defenses:** These are specific attack vectors red teaming aims to uncover.²⁵⁰

- *Jailbreaking*: Techniques trying to bypass the model's safety constraints to generate forbidden content.²⁵⁰
- *Prompt Injection*: Manipulating the model by embedding malicious instructions within the prompt, often hidden or obfuscated.²⁵⁰ Techniques include role-playing ("Act as..."), prefix injection, obfuscation (using different languages, encodings), multi-turn manipulation, and indirect injection (via retrieved documents or external tools).²⁵¹
- *Defenses*: Include strict input validation and sanitization (checking formatting, normalizing characters, filtering harmful patterns)²⁵³, clearly separating trusted instructions from untrusted user input, monitoring conversation context for suspicious shifts²⁵⁴, and output filtering.²⁴⁹
- **Constitutional AI (CAI)**: Anthropic's approach uses a predefined "constitution" – a set of explicit principles (e.g., "Choose the response that is least harmful") – to guide model behavior.²⁴⁷ The model learns to critique and revise its own responses based on the constitution (Supervised phase) and is then fine-tuned using RL where the reward signal comes from ranking responses against the constitution (RL phase).²⁴⁷ This makes the alignment principles more transparent.²⁴⁷

The overall alignment process reflects a progression from general language understanding (pretraining) to task-specific competence (SFT) and finally to nuanced behavioral shaping based on human values (RLHF/CAI). The quality of data at each stage, particularly the human preference data for RLHF, is critical for the final model's alignment and performance. Furthermore, safety cannot rely solely on training; a defense-in-depth strategy combining robust training-time alignment with proactive testing (red teaming) and runtime safeguards (filters, input sanitization) is necessary to mitigate the risks associated with powerful LLMs.

Table 6.1: Overview of Instruction Tuning Datasets

Dataset	Source	Approx. Size	Language(s)	License Notes	Key Characteristics
FLAN	NLP Tasks + Templates ²²⁵	N/A	EN, Multi	Varies by source dataset	Zero/Few-shot/CoT templates ²²⁵

P3	NLP Tasks + Prompts ²²⁴	N/A	EN	Varies by source dataset	Crowd-sourced prompts, shorter ²²⁴
Alpaca	LLM (Davinci-003) Generated ²²⁵	52K	EN	CC BY-NC 4.0 (Non-Commercial) ²²⁶	Self-Instruct method, single-turn
Dolly	Human Generated ²²⁵	15K	EN	Apache 2.0 (Commercial Use) ²²⁵	Databricks employees, 7 use cases
OpenAssistant	Human Generated (Crowd) ²²⁵	11K convos	Multi	Apache 2.0 (Commercial Use) ²²⁵	Multi-turn dialogues, quality ratings
ShareGPT	User Conversations (Scraped) ²²⁵	70K-90K convos	Multi (mainly EN)	Unclear/Problematic (User data) ²²⁵	Real multi-turn conversations, potentially noisy

Table 6.2: RLHF vs. Constitutional AI

Method	Alignment Signal Source	Key Training Process Steps	Pros	Cons	Key Proponents
RLHF	Human Preferences (Pairwise Comparisons) ²³¹	1. (SFT) 2. Collect Human Preferences 3. Train Reward Model (RM) 4. Fine-tune LLM Policy via RL (e.g., PPO) using	Captures nuanced preferences, empirically effective	Complex, data-intensive (human labeling), sensitive to RM quality ²³⁵	OpenAI

		RM ²³¹			
Constitutional AI (CAI)	Written Principles (Constitution) ²⁵⁵	1. (SFT) 2. Supervised Phase (Model critiques/revises based on constitution) 3. RL Phase (Fine-tune using constitution-based rewards) ²⁴⁷	Transparent principles, less reliance on direct human labeling for RL phase	Constitution design is crucial, may not capture all nuances as well as RLHF	Anthropic ²⁵⁵

Table 6.3: LLM Safety Mechanisms

Mechanism	Description	Stage Applied	Target Threats
Alignment Training (RLHF/CAI)	Train model to prefer helpful/harmless/honest outputs based on feedback/principles ²³²	Training	Generating harmful content, bias, misalignment
Moderation Filters	Classify/block harmful or forbidden content in inputs/outputs ⁸⁴	Inference (Input/Output)	Toxicity, hate speech, policy violations
Red Teaming	Adversarial testing to find vulnerabilities and failure modes ²⁴⁹	Testing/Evaluation	Jailbreaks, prompt injection, unforeseen harmful behaviors
Input Sanitization	Clean/validate user input to remove malicious instructions or formatting ²⁵³	Inference (Input)	Prompt injection, code injection, obfuscation
Output Filtering	Validate/sanitize	Inference (Output)	XSS, sensitive data

	model output before displaying to user or passing to other systems ²⁵⁸		leakage, harmful content generation (secondary defense)
Access Control	Limit model's access to external tools, APIs, or sensitive data based on least privilege ²⁵⁹	Deployment/Runtime	Unauthorized actions, data exfiltration via tool use

7. Inference and Optimization

Once an LLM is trained, deploying it efficiently for inference (generating text for user prompts) presents a different set of challenges compared to training. Optimization focuses on minimizing latency and maximizing throughput while managing resource consumption, particularly GPU memory.

7.1 Inference vs. Training: The Role of KV Cache

Autoregressive generation, the standard for LLMs like GPT, involves generating tokens one by one, with each new token depending on all previously generated tokens.²⁶⁰ A naive implementation would require recomputing the attention scores over the entire sequence for every new token, leading to quadratic complexity in sequence length for each step.

The **Key-Value (KV) Cache** is a critical optimization to avoid this.²⁶¹ During the attention calculation in each Transformer layer, Key (K) and Value (V) vectors are computed for each token. In inference, once the K and V vectors for a token at position t are computed, they are stored (cached) in GPU memory.²⁶³ When generating the token at position $t+1$, the model only needs to compute the Query (Q), Key (K), and Value (V) vectors for the *new* token. It then retrieves the cached K and V vectors for all preceding tokens ($1...t$) and computes attention using the new Q vector against all keys (new and cached).²⁶²

- Impact:** This drastically reduces computation, making generation time roughly linear (instead of quadratic) with respect to the sequence length after the initial prompt processing (prefill phase).²⁶² However, the KV cache itself consumes significant GPU memory, scaling linearly with sequence length, batch size, model size (layers, heads, dimensions), and precision.²⁶¹ This memory consumption often becomes the limiting factor for inference throughput.²⁶³

7.2 Decoding Strategies

Given the model's output logits (a probability distribution over the vocabulary for the next token), a decoding strategy selects the actual token to generate.²⁶⁰ Different strategies offer trade-offs between output quality, diversity, and computational cost²⁶⁰.

- **Greedy Search:** Always selects the single token with the highest probability.²⁶⁰ It's deterministic and computationally cheap but often produces repetitive, dull, or suboptimal sequences.²⁶⁰
- **Beam Search:** Maintains the k (beam width) most probable partial sequences at each step. It explores more of the search space than greedy search, often leading to higher-quality (higher overall probability) sequences.²⁶⁰ However, it is computationally more expensive (proportional to k) and can still lack diversity, sometimes favoring generic outputs.²⁶⁰
- **Sampling Strategies:** Introduce randomness to improve diversity and creativity.
 - **Temperature Sampling:** Rescales the logits before applying softmax ($\text{softmax}(\text{logits}/T)$). Temperature $T > 1$ flattens the distribution, increasing randomness; $T < 1$ sharpens it, making it closer to greedy.
 - **Top-k Sampling:** Restricts sampling to the k most probable tokens. The probabilities of these k tokens are renormalized, and one is sampled.²⁶⁰ It prevents sampling very low-probability tokens while allowing diversity among the top candidates.²⁶⁵
 - **Top-p (Nucleus) Sampling:** Samples from the smallest set of tokens whose cumulative probability mass exceeds a threshold p.²⁶⁶ The size of the sampling set (the nucleus) adapts dynamically based on the distribution's shape – narrow for high-confidence predictions, wider for uncertain ones.²⁶⁵ Often considered to provide a good balance between coherence and diversity.²⁶⁵
- **Combined Strategies:** Top-k and Top-p are often used together, applying Top-k first and then Top-p to potentially further restrict the sampling pool.²⁶⁶

7.3 Efficient Inference Libraries and Techniques

Specialized libraries and techniques optimize LLM inference execution:

- **ONNX (Open Neural Network Exchange):** An open standard format for representing ML models.²⁶⁸ Exporting a trained model (e.g., from PyTorch) to ONNX allows it to be run using optimized runtimes like **ONNX Runtime**.²⁶⁹ ONNX Runtime provides cross-platform compatibility and hardware-accelerated execution.²⁶⁹ Tools like Hugging Face Optimum facilitate export.²⁷⁰
- **NVIDIA TensorRT / TensorRT-LLM:** TensorRT is NVIDIA's SDK for

high-performance deep learning inference optimization and deployment.²⁷²

TensorRT-LLM is a specialized open-source library built on TensorRT, providing highly optimized kernels and runtime components specifically for LLMs.²⁷² Key optimizations include:

- *Kernel Fusion*: Merging multiple operations (e.g., within attention or FFN blocks) into single CUDA kernels to reduce memory access and launch overhead.²⁷²
- *Quantization*: Using lower precision (INT8, FP8) for weights and/or activations to reduce memory footprint and accelerate computation on supported hardware (like Hopper GPUs with FP8).²⁷²
- *Paged Attention*: An efficient memory management technique for the KV cache, similar to virtual memory paging, reducing fragmentation and waste.²⁷²
- *In-Flight Batching (Continuous Batching)*: Dynamically managing batches of requests, allowing new requests to start processing as soon as resources become available, rather than waiting for the entire current batch to finish. This improves GPU utilization and overall throughput.²⁷²
- *Multi-GPU / Multi-Node Support*: Efficiently scales inference across multiple GPUs using tensor parallelism.²⁷⁴
- **vLLM**: An open-source library from UC Berkeley specifically designed for fast LLM inference and serving.²⁷⁶ Its core innovations are:
 - *PagedAttention*: Manages the KV cache in non-contiguous memory blocks (pages), similar to OS virtual memory, significantly reducing internal fragmentation and wasted memory compared to traditional KV cache allocation.²⁷⁷ This allows for larger batch sizes and longer sequences.
 - *Continuous Batching*: Processes incoming requests in continuous batches rather than static ones, maximizing GPU utilization and throughput, especially with variable input/output lengths.²⁷⁶
- **NVIDIA FasterTransformer**: An earlier NVIDIA library providing optimized C++/CUDA implementations for Transformer layers.²⁷⁹ It includes layer fusion, activation caching, low-precision support (FP16, INT8), and GEMM auto-tuning.²⁷⁹ It was often used with Triton Inference Server.²⁸⁰ While still functional, TensorRT-LLM is generally the recommended NVIDIA solution for newer models and features.²⁷⁹
- **DeepSpeed-MII**: Part of the DeepSpeed ecosystem, MII (Model Implementations for Inference) provides a Python library for low-latency, high-throughput inference, leveraging DeepSpeed-Inference optimizations like DeepFusion, tensor slicing, ZeRO-Inference, and optimized CUDA kernels.¹¹⁷ It aims for ease of use, allowing deployment with minimal code changes.¹¹⁷ DeepSpeed-FastGen is a recent enhancement focusing on high-throughput text generation using

techniques like continuous batching and dynamic SplitFuse.²⁸³

- **Hugging Face Text Generation Inference (TGI):** A production-ready toolkit from Hugging Face for serving LLMs.²⁸⁵ It incorporates many state-of-the-art optimizations, including tensor parallelism, continuous batching, Flash Attention, Paged Attention, quantization (bitsandbytes, GPTQ), safetensors loading, and watermarking.²⁸⁶

7.4 Model Compression Techniques

To further reduce model size and accelerate inference, especially for deployment on resource-constrained devices:

- **Quantization:** Reducing the numerical precision of model weights and/or activations (e.g., from FP32/FP16 to INT8 or INT4) [²⁸⁸]

Works cited

1. Attention Is All You Need - Wikipedia, accessed May 5, 2025, https://en.wikipedia.org/wiki/Attention_Is_All_You_Need
2. papers.neurips.cc, accessed May 5, 2025, <https://papers.neurips.cc/paper/7181-attention-is-all-you-need.pdf>
3. Attention is All you Need - NIPS papers, accessed May 5, 2025, <https://papers.nips.cc/paper/7181-attention-is-all-you-need>
4. Transformer (deep learning architecture) - Wikipedia, accessed May 5, 2025, [https://en.wikipedia.org/wiki/Transformer_\(deep_learning_architecture\)](https://en.wikipedia.org/wiki/Transformer_(deep_learning_architecture))
5. How Transformers Work: A Detailed Exploration of Transformer Architecture - DataCamp, accessed May 5, 2025, <https://www.datacamp.com/tutorial/how-transformers-work>
6. Attention Is All You Need (Vaswani et al., ArXiv 2017) | Jonathan K. Kummerfeld, accessed May 5, 2025, https://jkk.name/reading-notes/old-blog/2017-10-20_onlyattention/
7. What changed in the Transformer architecture - Hugging Face, accessed May 5, 2025, <https://huggingface.co/blog/rishiraj/what-changed-in-the-transformer-architecture>
8. Transformer from scratch using pytorch - Kaggle, accessed May 5, 2025, <https://www.kaggle.com/code/arunmohan003/transformer-from-scratch-using-pytorch>
9. What is a Transformer Model? - IBM, accessed May 5, 2025, <https://www.ibm.com/think/topics/transformer-model>
10. Transformer: PyTorch Implementation of "Attention Is All You Need" - GitHub, accessed May 5, 2025, <https://github.com/hyunwoongko/transformer>
11. The Illustrated GPT-2 (Visualizing Transformer Language Models ..., accessed May 5, 2025, <https://jalammar.github.io/illustrated-gpt2/>
12. Masking in Transformer Encoder/Decoder Models - Sanjaya's Blog, accessed May

- 5, 2025, <https://sanjayasubedi.com.np/deeplearning/masking-in-attention/>
13. TransformerDecoderLayer — PyTorch 2.7 documentation, accessed May 5, 2025, <https://pytorch.org/docs/stable/generated/torch.nn.TransformerDecoderLayer.html>
 14. 11.9. Large-Scale Pretraining with Transformers - Dive into Deep Learning, accessed May 5, 2025, http://d2l.ai/chapter_attention-mechanisms-and-transformers/large-pretraining-transformers.html?highlight=transformer
 15. Pretraining of Large Language Models - GitHub Gist, accessed May 5, 2025, <https://gist.github.com/ritwikraha/77e79990992043f60a9588610b2781c5>
 16. How Transformers solve tasks - Hugging Face, accessed May 5, 2025, https://huggingface.co/docs/transformers/tasks_explained
 17. The Ultimate Guide to Fine-Tuning LLMs from Basics to Breakthroughs: An Exhaustive Review of Technologies, Research, Best Practices, Applied Research Challenges and Opportunities (Version 1.0) - arXiv, accessed May 5, 2025, <https://arxiv.org/html/2408.13296v1>
 18. Transformers cross-entropy loss masked label issue - Stack Overflow, accessed May 5, 2025, <https://stackoverflow.com/questions/77628127/transformers-cross-entropy-loss-masked-label-issue>
 19. Day: 25 Optimizer Algorithms for Large Language Models (LLMs) - DEV Community, accessed May 5, 2025, <https://dev.to/nareshnishad/day-25-optimizer-algorithms-for-large-language-models-llms-1p4>
 20. Session 5: LLM training and fine-tuning - MLSys 2025, accessed May 5, 2025, <https://mlsys.org/virtual/2025/session/3148>
 21. Anything but SGD: Evaluating Optimizers for LLM Training ..., accessed May 5, 2025, <https://kempnerinstitute.harvard.edu/research/deeper-learning/anything-but-sgd-evaluating-optimizers-for-llm-training/>
 22. AdamW Optimizer in PyTorch Tutorial - DataCamp, accessed May 5, 2025, <https://www.datacamp.com/tutorial/adamw-optimizer-in-pytorch>
 23. AdamW and Super-convergence is now the fastest way to train neural nets - Fast.ai, accessed May 5, 2025, <https://www.fast.ai/posts/2018-07-02-adam-weight-decay.html>
 24. PyTorch Optimizer: AdamW and Adam with weight decay - Stack Overflow, accessed May 5, 2025, <https://stackoverflow.com/questions/64621585/pytorch-optimizer-adamw-and-adam-with-weight-decay>
 25. What are the key differences between Adam and AdamW optimizers in the context of transformer-based language models? - Massed Compute, accessed May 5, 2025, <https://massedcompute.com/faq-answers/?question=What%20are%20the%20key%20differences%20between%20Adam%20and%20AdamW%20optimizers%20in%20the%20context%20of%20transformer-based%20language%20models?>

26. Mastering AdamW Optimizer: Enhancing Deep Learning Models with Superior Regularization - LUNARTECH, accessed May 5, 2025, <https://www.lunartech.ai/blog/mastering-adamw-optimizer-enhancing-deep-learning-models-with-superior-regularization>
27. Understanding L2 regularization, Weight decay and AdamW | Another Deep-Learning Blog, accessed May 5, 2025, <https://benihime91.github.io/blog/machinelearning/deeplearning/python3.x/tensorflow2.x/2020/10/08/adamW.html>
28. [R] Why is AdamW often superior to Adam with L2-Regularization in practice? The answer may lie in how weight decay balances updates across layers. - Reddit, accessed May 5, 2025, https://www.reddit.com/r/MachineLearning/comments/1731pcg/r_why_is_adamw_often_superior_to_adam_with/
29. [D] Lion , An Optimizer That Outperforms Adam - Symbolic Discovery of Optimization Algorithms : r/MachineLearning - Reddit, accessed May 5, 2025, https://www.reddit.com/r/MachineLearning/comments/1138jpp/d_lion_an_optimizer_that_outperforms_adam/
30. Straight to Zero: Why Linearly Decaying the Learning Rate to Zero ..., accessed May 5, 2025, <https://openreview.net/forum?id=hrOIBgHsMI>
31. Understanding Warmup-Stable-Decay Learning Rates: A River Valley Loss Landscape Perspective - arXiv, accessed May 5, 2025, <https://arxiv.org/html/2410.05192v1>
32. Glossary: LLM fine-tuning hyperparameters | Modal Blog, accessed May 5, 2025, <https://modal.com/blog/fine-tuning-llms-hyperparameters-glossary-article>
33. GPT-3 vs GPT-4 | What's the difference? - Botpress, accessed May 5, 2025, <https://botpress.com/blog/gpt-3-vs-gpt-4-whats-the-difference>
34. Models - OpenAI API, accessed May 5, 2025, <https://platform.openai.com/docs/models>
35. ChatGPT and OpenAI's Evolution: From GPT-2 to GPT-4 - LLM Directory, accessed May 5, 2025, <https://llmmodels.org/blog/chatgpt-and-openais-evolution-from-gpt-2-to-gpt-4/>
36. Generative pre-trained transformer - Wikipedia, accessed May 5, 2025, https://en.wikipedia.org/wiki/Generative_pre-trained_transformer
37. RoPE: Rotary Positional Embedding - Sushant Kumar, accessed May 5, 2025, <https://sushant-kumar.com/blog/rope-rotary-positional-embedding>
38. Extending the RoPE - EleutherAI Blog, accessed May 5, 2025, <https://blog.eleuther.ai/yarn/>
39. Decoding Llama3: Part 4 - Rotary Positional Embeddings, accessed May 5, 2025, <https://hasgeek.com/simrathanspal/the-llama3-guide/sub/decoding-llama3-part-4-rotary-positional-embedding-3K8ZHpdLi6E56N8ejnaWzm>
40. RoPE - Rotary Positional Embedding - After Hours Research, accessed May 5, 2025, <https://afterhoursresearch.hashnode.dev/rope-rotary-positional-embedding>
41. Rotary Embeddings: A Relative Revolution | EleutherAI Blog, accessed May 5, 2025, <https://blog.eleuther.ai/rotary-embeddings/>

42. Rotary Position Embedding (RoPE) - AI Resources - Modular, accessed May 5, 2025, <https://www.modular.com/ai-resources/rotary-position-embedding-rope>
43. Rotatory Position Embedding (RoPE) - Karthick Panner Selvam, accessed May 5, 2025, [https://karthick.ai/blog/2024/Rotatory-Position-Embedding-\(RoPE\)/](https://karthick.ai/blog/2024/Rotatory-Position-Embedding-(RoPE)/)
44. [2502.11276] The Rotary Position Embedding May Cause Dimension Inefficiency in Attention Heads for Long-Distance Retrieval - arXiv, accessed May 5, 2025, <https://arxiv.org/abs/2502.11276>
45. Round and Round We Go! What makes Rotary Positional Encodings useful? - arXiv, accessed May 5, 2025, <https://arxiv.org/abs/2410.06205>
46. Efficient Streaming Language Models with Attention Sinks - arXiv, accessed May 5, 2025, <https://arxiv.org/html/2309.17453v3>
47. arxiv.org, accessed May 5, 2025, <https://arxiv.org/abs/2502.18845>
48. [R] Sliding Window Attention Training for Efficient LLMs : r/MachineLearning - Reddit, accessed May 5, 2025, https://www.reddit.com/r/MachineLearning/comments/1j1ckye/r_sliding_window_attention_training_for_efficient/
49. 1 Introduction - arXiv, accessed May 5, 2025, <https://arxiv.org/html/2412.05496v1>
50. Polynomial Composition Activations: Unleashing the Dynamics of Large Language Models, accessed May 5, 2025, <https://arxiv.org/html/2411.03884v3>
51. arxiv.org, accessed May 5, 2025, <https://arxiv.org/abs/2402.03804>
52. ReLU2 Wins: Discovering Efficient Activation Functions for Sparse LLMs - arXiv, accessed May 5, 2025, <https://arxiv.org/pdf/2402.03804>
53. HAAN: A Holistic Approach for Accelerating Normalization Operations in Large Language Models - arXiv, accessed May 5, 2025, <https://arxiv.org/html/2502.11832v1>
54. Paper page - Root Mean Square Layer Normalization - Hugging Face, accessed May 5, 2025, <https://huggingface.co/papers/1910.07467>
55. Understanding Llama2: KV Cache, Grouped Query Attention, Rotary Embedding and More, accessed May 5, 2025, <https://plainenglish.io/community/understanding-llama2-kv-cache-grouped-query-attention-rotary-embedding-and-more-9a79bd>
56. Geometric Interpretation of Layer Normalization and a Comparative Analysis with RMSNorm, accessed May 5, 2025, <https://arxiv.org/html/2409.12951v2>
57. Batch Normalization, Layer Normalization and Root Mean Square ..., accessed May 5, 2025, <https://afterhoursresearch.hashnode.dev/batch-normalization-layer-normalization-and-root-mean-square-layer-normalization-a-comprehensive-guide-with-python-implementations>
58. Re-Introducing LayerNorm: Geometric Meaning, Irreversibility and a Comparative Study with RMSNorm - arXiv, accessed May 5, 2025, <https://arxiv.org/html/2409.12951v1>
59. Progressive Sparse Attention: Algorithm and System Co-design for Efficient Attention in LLM Serving - arXiv, accessed May 5, 2025, <https://arxiv.org/html/2503.00392v1>
60. Designing Hardware-Aware Algorithms: FlashAttention - DigitalOcean, accessed

- May 5, 2025, <https://www.digitalocean.com/community/tutorials/flashattention>
61. Paper Summary #8 - FlashAttention: Fast and Memory-Efficient Exact Attention with IO-Awareness | Shreyansh Singh, accessed May 5, 2025, https://shreyansh26.github.io/post/2023-03-26_flash-attention/
 62. FlashAttention: Fast and memory-efficient exact attention with IO-Awareness - Together AI, accessed May 5, 2025, <https://www.together.ai/blog/flashattentionfandm>
 63. FlashAttention: Fast and Memory-Efficient Exact Attention with IO-Awareness | OpenReview, accessed May 5, 2025, <https://openreview.net/forum?id=H4DqfPSibmx>
 64. DistFlashAttn: Distributed Memory-efficient Attention for Long-context LLMs Training - arXiv, accessed May 5, 2025, <https://arxiv.org/html/2310.03294v2>
 65. Flash attention(Fast and Memory-Efficient Exact Attention with IO ..., accessed May 5, 2025, <https://towardsdatascience.com/flash-attention-fast-and-memory-efficient-exact-attention-with-io-awareness-a-deep-dive-724af489997b>
 66. Flash attention(Fast and Memory-Efficient Exact Attention with IO-Awareness): A deep dive, accessed May 5, 2025, <https://towardsdatascience.com/flash-attention-fast-and-memory-efficient-exact-attention-with-io-awareness-a-deep-dive-724af489997b/>
 67. [2310.03294] DISTFLASHATTN: Distributed Memory-efficient Attention for Long-context LLMs Training - arXiv, accessed May 5, 2025, <https://arxiv.org/abs/2310.03294>
 68. Methods of improving LLM training stability - arXiv, accessed May 5, 2025, <https://arxiv.org/html/2410.16682v1>
 69. Crawl4LLM: Efficient Web Crawling for LLM Pretraining - arXiv, accessed May 5, 2025, <https://arxiv.org/html/2502.13347v1>
 70. Announcing Nemotron-CC: A Trillion-Token English Language Dataset for LLM Pretraining, accessed May 5, 2025, <https://developer.nvidia.com/blog/announcing-nemotron-cc-a-trillion-token-english-language-dataset-for-llm-pretraining/>
 71. Technical Deep-Dive: Curating Our Way to a State-of-the-Art Text Dataset - Blog, accessed May 5, 2025, <https://blog.datologyai.com/technical-deep-dive-curating-our-way-to-a-state-of-the-art-text-dataset/>
 72. Open-Sourced Training Datasets for Large Language Models (LLMs) - Kili Technology, accessed May 5, 2025, <https://kili-technology.com/large-language-models-llms/9-open-sourced-datasets-for-training-large-language-models>
 73. Recipe for Success: Blending Data for Better LLM Pretraining - Snowflake, accessed May 5, 2025, <https://www.snowflake.com/en/engineering-blog/blending-data-for-better-llm-pretraining/>
 74. Large language model - Wikipedia, accessed May 5, 2025, https://en.wikipedia.org/wiki/Large_language_model

75. Towards a Books Data Commons for AI Training - Open Future Foundation, accessed May 5, 2025, https://openfuture.eu/wp-content/uploads/2024/04/240404Towards_a_Books_Data_Commons_for_AI_Training.pdf
76. The Pile (dataset) - Wikipedia, accessed May 5, 2025, [https://en.wikipedia.org/wiki/The_Pile_\(dataset\)](https://en.wikipedia.org/wiki/The_Pile_(dataset))
77. Immlzn/Awesome-LLMs-Datasets: Summarize existing ... - GitHub, accessed May 5, 2025, <https://github.com/Immlzn/Awesome-LLMs-Datasets>
78. DCAD-2000: A Multilingual Dataset across 2000+ Languages with Data Cleaning as Anomaly Detection - arXiv, accessed May 5, 2025, <https://arxiv.org/html/2502.11546v1>
79. Unveiling CulturaX: An Enormous, Multilingual Dataset for LLMs in 167 Languages, accessed May 5, 2025, <https://www.e2enetworks.com/blog/unveiling-culturax-an-enormous-multilingual-dataset-for-llms-in-167-languages>
80. The Pile, accessed May 5, 2025, <https://pile.eleuther.ai/>
81. facebook/fasttext-language-identification - Hugging Face, accessed May 5, 2025, <https://huggingface.co/facebook/fasttext-language-identification>
82. LlmKira/fast-langdetect: ⚡ 80x faster Fasttext language detection out of the box - GitHub, accessed May 5, 2025, <https://github.com/LlmKira/fast-langdetect>
83. LLM Evaluation: Best Metrics & Tools - UBAI, accessed May 5, 2025, <https://ubai.tools/llm-evaluation-best-metrics-tools/>
84. Concept | Guardrails against risks from Generative AI and LLMs - Dataiku Knowledge Base, accessed May 5, 2025, <https://knowledge.dataiku.com/latest/gen-ai/admin/concept-genai-guardrails.html>
85. Combating Toxic Language: A Review of LLM-Based Strategies for Software Engineering, accessed May 5, 2025, <https://arxiv.org/html/2504.15439v1>
86. Evaluation of Document Deduplication Algorithms for Large Text Corpora - Fraunhofer-Publica, accessed May 5, 2025, <https://publica.fraunhofer.de/entities/publication/631a494a-83a1-456b-94e3-518221257087>
87. Large-scale Near-deduplication Behind BigCode - Hugging Face, accessed May 5, 2025, <https://huggingface.co/blog/dedup>
88. LLM training datasets - Glenn K. Lockwood, accessed May 5, 2025, <https://www.glennklockwood.com/garden/LLM-training-datasets>
89. FED: Fast and Efficient Dataset Deduplication Framework with GPU Acceleration - arXiv, accessed May 5, 2025, <https://arxiv.org/html/2501.01046v2>
90. 1 Introduction - arXiv, accessed May 5, 2025, <https://arxiv.org/html/2411.04257v1>
91. LSH/examples/Introduction.ipynb at master · mattilyra/LSH - GitHub, accessed May 5, 2025, <https://github.com/mattilyra/LSH/blob/master/examples/Introduction.ipynb>
92. Minhash LSH Implementation Walkthrough: Deduplication - DZone, accessed May 5, 2025, <https://dzone.com/articles/minhash-lsh-implementation-walkthrough>
93. Large scale document similarity search with LSH and MinHash ..., accessed May 5,

- 2025, <https://mrhasankthse.github.io/riz/2020/03/19/Minhash-and-LSH.html>
94. LSH methods for data deduplication in a Wikipedia artificial dataset - arXiv, accessed May 5, 2025, <https://arxiv.org/pdf/2112.11478>
 95. Finding near-duplicates with Jaccard similarity and MinHash - Made of Bugs, accessed May 5, 2025, <https://blog.nelhage.com/post/fuzzy-dedup/>
 96. Bloom filter - Wikipedia, accessed May 5, 2025, https://en.wikipedia.org/wiki/Bloom_filter
 97. A Gentle Introduction to Bloom Filter | iunera, accessed May 5, 2025, <https://www.iunera.com/kraken/fabric/a-gentle-introduction-to-bloom-filter/>
 98. Personalized version of the Bloom Filter Algorithm for duplicate detection of massive datasets of passwords. - GitHub, accessed May 5, 2025, <https://github.com/paologentleman/Bloom-Filter-Algorithm>
 99. Introduction to LLM Tokenization - Airbyte, accessed May 5, 2025, <https://airbyte.com/data-engineering-resources/llm-tokenization>
 100. karpathy/minbpe: Minimal, clean code for the Byte Pair ... - GitHub, accessed May 5, 2025, <https://github.com/karpathy/minbpe>
 101. So many tokens, so little time: Introducing a faster, more flexible byte-pair tokenizer, accessed May 5, 2025, <https://github.blog/ai-and-ml/llms/so-many-tokens-so-little-time-introducing-a-faster-more-flexible-byte-pair-tokenizer/>
 102. Implementing A Byte Pair Encoding (BPE) Tokenizer From Scratch - Sebastian Raschka, accessed May 5, 2025, <https://sebastianraschka.com/blog/2025/bpe-from-scratch.html>
 103. Tiktoken Tokenizer for GPT-4o, GPT-4, and o1 OpenAI models - Pub.dev, accessed May 5, 2025, https://pub.dev/documentation/tiktoken_tokenizer_gpt4o_o1/latest/
 104. Tokenizer - OpenAI API, accessed May 5, 2025, <https://platform.openai.com/tokenizer>
 105. google/sentencepiece: Unsupervised text tokenizer for ... - GitHub, accessed May 5, 2025, <https://github.com/google/sentencepiece>
 106. Spark NLP - State of the Art NLP Library for Large Language Models ..., accessed May 5, 2025, <https://sparknlp.org/>
 107. How to effectively use Spark for natural language processing and text mining? - HopHR, accessed May 5, 2025, <https://www.hophr.com/tutorial-page/effectively-use-spark-for-natural-language-processing-and-text-mining>
 108. Deep Learning with Apache Spark and NetApp AI - Financial Sentiment Analysis, accessed May 5, 2025, <https://www.netapp.com/blog/deep-learning-with-apache-spark-and-netapp-ai-financial-sentiment-analysis/>
 109. Load data using Petastorm - Azure Databricks | Azure Docs, accessed May 5, 2025, <https://docs.azure.cn/en-us/databricks/archive/machine-learning/petastorm>
 110. Petastorm | AI and Machine Learning - Howdy, accessed May 5, 2025, <https://www.howdy.com/glossary/petastorm>
 111. Load Dataset Guide - Hugging Face, accessed May 5, 2025,

- <https://huggingface.co/docs/datasets/loading>
112. Loading methods - Hugging Face, accessed May 5, 2025,
https://huggingface.co/docs/datasets/package_reference/loading_methods
 113. LLM Training on Unity Catalog data with MosaicML Streaming Dataset - Databricks, accessed May 5, 2025,
<https://www.databricks.com/blog/llm-training-unity-catalog-data-mosaicml-streaming-dataset>
 114. TensorFlow vs PyTorch vs JAX: Performance Benchmark, accessed May 5, 2025,
<https://www.apxml.com/posts/tensorflow-vs-pytorch-vs-jax-performance-benchmark>
 115. ML Engineer comparison of Pytorch, TensorFlow, JAX, and Flax - SoftwareMill, accessed May 5, 2025,
<https://softwaremill.com/ml-engineer-comparison-of-pytorch-tensorflow-jax-and-flax/>
 116. Deep Learning Model Multi-Node, Distributed Training Strategies (Primer) - Mark III Systems, accessed May 5, 2025,
<https://www.markiiisys.com/blog/ml-dl-model-multi-node-distributed-training-strategies-primer/>
 117. DeepSpeed is a deep learning optimization library that makes distributed training and inference easy, efficient, and effective. - GitHub, accessed May 5, 2025, <https://github.com/deepspeedai/DeepSpeed>
 118. DeepSpeed vs FSDP: A Comprehensive Comparison - BytePlus, accessed May 5, 2025, <https://www.byteplus.com/en/topic/499106>
 119. Training Overview and Features - DeepSpeed, accessed May 5, 2025,
<https://www.deepspeed.ai/training/>
 120. Megatron-LM - NVIDIA NeMo Framework User Guide, accessed May 5, 2025,
<https://docs.nvidia.com/nemo-framework/user-guide/24.12/nemotoolkit/nlp/megatron.html>
 121. Exploring the Exciting Possibilities of NVIDIA Megatron LM: A Fun and Friendly Code Walkthrough with PyTorch & NVIDIA Apex! - DEV Community, accessed May 5, 2025,
https://dev.to/hassan_sherwani_9dd766c43/exploring-the-exciting-possibilities-of-nvidia-megatron-lm-a-fun-and-friendly-code-walkthrough-with-pytorch-nvidia-apex-n3d
 122. Megatron-LM - Hugging Face, accessed May 5, 2025,
https://huggingface.co/docs/accelerate/usage_guides/megatron_lm
 123. Model Parallelism - Hugging Face, accessed May 5, 2025,
<https://huggingface.co/docs/transformers/v4.13.0/parallelism>
 124. Efficient Training on Multiple GPUs - Hugging Face, accessed May 5, 2025,
https://huggingface.co/docs/transformers/v4.28.1/perf_train_gpu_many
 125. How to Fine-Tune an LLM from Hugging Face - GeeksforGeeks, accessed May 5, 2025,
<https://www.geeksforgeeks.org/how-to-fine-tune-an-llm-from-hugging-face/>
 126. Fine-tune Hugging Face models for a single GPU - Databricks Documentation,

- accessed May 5, 2025,
<https://docs.databricks.com/aws/en/machine-learning/train-model/huggingface/finetune-model>
127. Fine-tune a pretrained model - Hugging Face, accessed May 5, 2025,
<https://huggingface.co/docs/transformers/v4.37.1/training>
 128. JAX vs PyTorch: Comparing Two Deep Learning Frameworks - New Horizons, accessed May 5, 2025,
<https://www.newhorizons.com/resources/blog/jax-vs-pytorch-comparing-two-deep-learning-frameworks>
 129. A guide to JAX for PyTorch developers | Google Cloud Blog, accessed May 5, 2025,
<https://cloud.google.com/blog/products/ai-machine-learning/guide-to-jax-for-pytorch-developers>
 130. kevinpdev/gpt-from-scratch: Educational implementation of ... - GitHub, accessed May 5, 2025, <https://github.com/kevinpdev/gpt-from-scratch>
 131. GPT-2 style transformer implementation from scratch : r/learnmachinelearning - Reddit, accessed May 5, 2025,
https://www.reddit.com/r/learnmachinelearning/comments/1k138c1/gpt2_style_transformer_implementation_from_scratch/
 132. Implementation of self-attention in Transformers.jl? - Machine Learning - Julia Discourse, accessed May 5, 2025,
<https://discourse.julialang.org/t/implementation-of-self-attention-in-transformers-jl/94732>
 133. pytorch - Training torch.TransformerDecoder with causal mask - Stack Overflow, accessed May 5, 2025,
<https://stackoverflow.com/questions/77923078/training-torch-transformerdecoder-with-causal-mask>
 134. SmolGPT: A minimal PyTorch implementation for training a small LLM from scratch | Hacker News, accessed May 5, 2025,
<https://news.ycombinator.com/item?id=42868770>
 135. Paradigms of Parallelism | Colossal-AI, accessed May 5, 2025,
https://colossalai.org/docs/concepts/paradigms_of_parallelism/
 136. Aman's AI Journal • Gradient Accumulation and Checkpointing, accessed May 5, 2025, <https://aman.ai/primers/ai/grad-accum-checkpoint/>
 137. Gradient Checkpointing | Towards AI, accessed May 5, 2025,
<https://towardsai.net/p//gradient-checkpointing>
 138. Messing around with fine-tuning LLMs, part 9 -- gradient checkpointing - Giles' blog, accessed May 5, 2025,
<https://www.gilesthomas.com/2024/09/fine-tuning-9>
 139. DeepSpeed Mixed Precision Training - Tutorialspoint, accessed May 5, 2025,
<https://www.tutorialspoint.com/deepspeed/deepspeed-mixed-precision-training.htm>
 140. train - Understanding the advantages of BF16 vs. FP16 in mixed ..., accessed May 5, 2025,
<https://stats.stackexchange.com/questions/637988/understanding-the-advantage>

- [s-of-bf16-vs-fp16-in-mixed-precision-training](#)
141. DeepSpeed - Hugging Face, accessed May 5, 2025, <https://huggingface.co/docs/transformers/deepspeed>
 142. RoRA: Efficient Fine-Tuning of LLM with Reliability Optimization for Rank Adaptation - arXiv, accessed May 5, 2025, <https://arxiv.org/abs/2501.04315>
 143. arxiv.org, accessed May 5, 2025, <https://arxiv.org/abs/2402.12354>
 144. Breaking Memory Limits: Gradient Wavelet Transform Enhances LLMs Training - arXiv, accessed May 5, 2025, <https://arxiv.org/html/2501.07237v1>
 145. GaLore: Memory-Efficient LLM Training by Gradient Low-Rank Projection - arXiv, accessed May 5, 2025, <https://arxiv.org/html/2403.03507v1>
 146. Continued Pretraining of State-of-the-Art LLMs for Sovereign AI and Regulated Industries with iGenius and NVIDIA DGX Cloud, accessed May 5, 2025, <https://developer.nvidia.com/blog/continued-pretraining-of-state-of-the-art-llms-for-sovereign-ai-and-regulated-industries-with-igenius-and-nvidia-dgx-cloud/>
 147. NeMo Framework Foundation Model Pre-training - NVIDIA Docs Hub, accessed May 5, 2025, <https://docs.nvidia.com/nemo-framework/user-guide/24.07/playbooks/pretraining.html>
 148. NVIDIA DGX Systems | DGX SuperPOD, DGX GB200, B200, H100 | Lambda, accessed May 5, 2025, <https://lambdalabs.com/nvidia/dgx-systems>
 149. DGX SuperPOD: AI Infrastructure for Enterprise Deployments | NVIDIA, accessed May 5, 2025, <https://www.nvidia.com/en-us/data-center/dgx-superpod/>
 150. Interconnect Needs for LLM Inference to Drive Networking ..., accessed May 5, 2025, <https://650group.com/blog/interconnect-needs-for-llm-inference-to-drive-networking-bandwidth/>
 151. Multi-GPU Systems and GPU Clusters (NVLink, UALink, InfiniBand) Essentials - Tonex Training, accessed May 5, 2025, <https://www.tonex.com/training-courses/multi-gpu-systems-and-gpu-clusters-nvlink-ualink-infiniband-essentials/>
 152. NVIDIA GB200 NVL72 Delivers Trillion-Parameter LLM Training and Real-Time Inference, accessed May 5, 2025, <https://developer.nvidia.com/blog/nvidia-gb200-nvl72-delivers-trillion-parameter-llm-training-and-real-time-inference/>
 153. Training LLMs: An efficient GPU traffic routing mechanism within AI/ML cluster with rail-only connections - Outshift | Cisco, accessed May 5, 2025, <https://outshift.cisco.com/blog/training-llms-efficient-gpu-traffic-routing>
 154. Deploy Kubernetes — NVIDIA DGX BasePOD, accessed May 5, 2025, <https://docs.nvidia.com/dgx-basepod/deployment-guide-multi-cloud-arch-aws/la test/deploy-kubernetes.html>
 155. GPU nodes on-premise : r/kubernetes - Reddit, accessed May 5, 2025, https://www.reddit.com/r/kubernetes/comments/1isk9ne/gpu_nodes_onpremise/
 156. Building RAG Applications with NVIDIA NIM and Haystack on K8s, accessed May 5, 2025, <https://haystack.deepset.ai/blog/haystack-nvidia-nim-rag-guide>
 157. Maximizing NVIDIA DGX with Kubernetes | NVIDIA Technical Blog, accessed

- May 5, 2025,
<https://developer.nvidia.com/blog/maximizing-nvidia-dgx-kubernetes/>
158. Choosing the Right Orchestration Tool for ML Workloads: Slurm vs ...,
accessed May 5, 2025,
<https://www.nscale.com/blog/choosing-the-right-orchestration-tool-for-ml-workloads-slurm-vs-kubernetes>
 159. Slurm vs Kubernetes: Which to choose for model training - Nebius, accessed
May 5, 2025, <https://nebius.com/blog/posts/model-pre-training/slurm-vs-k8s>
 160. Understanding Slurm for AI/ML Workloads - WhiteFiber, accessed May 5,
2025,
<https://www.whitefiber.com/blog/understanding-slurm-for-ai-ml-workloads>
 161. Is Kubernetes or Slurm the best orchestrator for 512+ GPU jobs. - Fluidstack,
accessed May 5, 2025,
<https://www.fluidstack.io/post/is-kubernetes-or-slurm-the-best-orchestrator-for-512-gpu-jobs>
 162. Scaling model training - ROCm Documentation - AMD, accessed May 5, 2025,
<https://rocm.docs.amd.com/en/latest/how-to/rocm-for-ai/training/scale-model-training.html>
 163. How to decide the distributed inference strategy? - vLLM, accessed May 5,
2025, https://docs.vllm.ai/en/v0.5.2/serving/distributed_serving.html
 164. FullyShardedDataParallel - Hugging Face, accessed May 5, 2025,
<https://huggingface.co/docs/transformers/fsdp>
 165. S3 / Blob Storage (self-hosted) - Langfuse, accessed May 5, 2025,
<https://langfuse.com/self-hosting/infrastructure/blobstorage>
 166. Object storage for data lakes: S3, GCS, Azure | AI Data Analytics - Starburst,
accessed May 5, 2025, <https://www.starburst.io/data-glossary/object-storage/>
 167. Compare Storage Services on Azure and AWS - Azure Architecture Center |
Microsoft Learn, accessed May 5, 2025,
<https://learn.microsoft.com/en-us/azure/architecture/aws-professional/storage>
 168. Parallel File Systems - Advanced HPC, accessed May 5, 2025,
<https://www.advancedhpc.com/pages/parallel-file-systems>
 169. Parallel File Systems - | HPC @ LLNL - Lawrence Livermore National
Laboratory, accessed May 5, 2025,
<https://hpc.llnl.gov/hardware/file-systems/parallel-file-systems>
 170. Lustre File System Explained | Key Components, Architecture & More,
accessed May 5, 2025,
<https://www.weka.io/learn/glossary/file-storage/lustre-file-system-explained/>
 171. High Performance File Systems for AI/ML - WWT, accessed May 5, 2025,
<https://www.wwt.com/article/high-performance-file-systems-for-aiml>
 172. Handling Large Datasets in LLM Training: Distributed Training Architectures
and Techniques - Just Total Tech, accessed May 5, 2025,
<https://justtotaltech.com/llm-training-distributed-training-architectures/>
 173. Strategies for distributed training - Amazon SageMaker AI, accessed May 5,
2025,
<https://docs.aws.amazon.com/sagemaker/latest/dg/distributed-training-strategies>

- [.html](#)
174. Weights & Biases: The AI Developer Platform, accessed May 5, 2025, <https://wandb.ai/site/>
 175. Machine Learning Experiment Tracking with Weights & Biases - Wandb, accessed May 5, 2025, <https://wandb.ai/site/experiment-tracking/>
 176. Intro to MLOps: Machine Learning Experiment Tracking - Weights & Biases - Wandb, accessed May 5, 2025, <https://wandb.ai/site/articles/intro-to-mlops-machine-learning-experiment-tracking/>
 177. LLM Monitoring Sign Up - Weights & Biases LLM - Wandb, accessed May 5, 2025, <https://wandb.ai/site/llm-monitoring-sign-up/>
 178. Machine Learning Experiment Tracking with Weights & Biases, accessed May 5, 2025, <https://wandb.ai/site/experiment-tracking>
 179. LLM Training Hardware API Pricing - BytePlus, accessed May 5, 2025, <https://www.byteplus.com/en/topic/382277>
 180. What is the Cost of Training LLM Models? A Comprehensive Guide for AI Professionals, accessed May 5, 2025, <https://www.galileo.ai/blog/llm-model-training-cost>
 181. Lenovo LLM Sizing Guide, accessed May 5, 2025, <https://lenovopress.lenovo.com/lp2130.pdf>
 182. What is the Cost of Training LLM Models? Key Factors Explained, accessed May 5, 2025, <https://botpenguin.com/blogs/what-is-the-cost-of-training-llm-models>
 183. Mosaic LLMs: GPT-3 quality for - Databricks, accessed May 5, 2025, <https://www.databricks.com/blog/gpt-3-quality-for-500k>
 184. If you read around, training a 7B model costs on the order of \$85000 - Hacker News, accessed May 5, 2025, <https://news.ycombinator.com/item?id=39224534>
 185. arXiv:2402.15627v1 [cs.LG] 23 Feb 2024, accessed May 5, 2025, <https://mirrors.qiql.net/books/MegaScale%20Scaling%20Large%20Language%20Model%20Training%20to%20More%20Than%2010%2C000%20GPUs.pdf>
 186. Scaling LLM Pre-training with Vocabulary Curriculum - arXiv, accessed May 5, 2025, <https://arxiv.org/html/2502.17910v1>
 187. SynLexLM: Scaling Legal LLMs with Synthetic Data and Curriculum Learning - arXiv, accessed May 5, 2025, <https://arxiv.org/html/2504.18762v1>
 188. Strategic Data Ordering: Enhancing Large Language Model Performance through Curriculum Learning - arXiv, accessed May 5, 2025, <https://arxiv.org/html/2405.07490v1>
 189. Scaling LLM Pre-training with Vocabulary Curriculum - arXiv, accessed May 5, 2025, <https://arxiv.org/pdf/2502.17910>
 190. Fine-tuning Large Language Models with Human-inspired Learning Strategies in Medical Question Answering - arXiv, accessed May 5, 2025, <https://arxiv.org/html/2408.07888v1>
 191. Preference Curriculum: LLMs Should Always Be Pretrained on Their Preferred Data - arXiv, accessed May 5, 2025, <https://arxiv.org/html/2501.13126v1>
 192. arxiv.org, accessed May 5, 2025, <https://arxiv.org/pdf/2406.06046>

193. Continual Learning for Large Language Models: A Survey - arXiv, accessed May 5, 2025, <https://arxiv.org/html/2402.01364v1>
194. Machine Learning Glossary - Google for Developers, accessed May 5, 2025, <https://developers.google.com/machine-learning/glossary>
195. Gradient Accumulation in LLM Training - Algos, accessed May 5, 2025, <https://algos-ai.com/gradient-accumulation-in-llm-training/>
196. The Novice's LLM Training Guide - GitHub Gist, accessed May 5, 2025, <https://gist.github.com/btbytes/cf845f9ade1cb34348110c14c8c49cea>
197. LLM training - Glenn K. Lockwood, accessed May 5, 2025, <https://www.glennklockwood.com/garden/LLM-training>
198. How to Fine-tune an LLM Part 3: The HuggingFace Trainer | alpaca_ft - Wandb, accessed May 5, 2025, https://wandb.ai/capecape/alpaca_ft/reports/How-to-Fine-tune-an-LLM-Part-3-The-HuggingFace-Trainer--Vmlldzo1OTEyNjMy
199. Surge Phenomenon in Optimal Learning Rate and Batch Size Scaling - NIPS papers, accessed May 5, 2025, https://proceedings.neurips.cc/paper_files/paper/2024/file/ef74413c7bf1d915c3e45c72e19a5d32-Paper-Conference.pdf
200. How does batch size affect convergence of SGD and why? - Stats Stackexchange, accessed May 5, 2025, <https://stats.stackexchange.com/questions/316464/how-does-batch-size-affect-convergence-of-sgd-and-why>
201. Surge Phenomenon in Optimal Learning Rate and Batch Size Scaling - arXiv, accessed May 5, 2025, <https://arxiv.org/html/2405.14578v5>
202. Resolving Discrepancies in Compute-Optimal Scaling of Language Models - NIPS papers, accessed May 5, 2025, https://proceedings.neurips.cc/paper_files/paper/2024/file/b6341525cd84f3be0ef203e4d7cd8556-Paper-Conference.pdf
203. llm-research-summaries/training/ultimate-guide-fine-tuning-llm_parthasarathy-2408.13296.md at main - GitHub, accessed May 5, 2025, https://github.com/cognitivetech/llm-research-summaries/blob/main/training/ultimate-guide-fine-tuning-llm_parthasarathy-2408.13296.md
204. Machine Learning Model Training: Pro Tips & Tools - Label Your Data, accessed May 5, 2025, <https://labelyourdata.com/articles/model-training>
205. Introduction to Fine-tuning Large Language Models - Stephen Diehl, accessed May 5, 2025, https://www.stephendiehl.com/posts/training_llms/
206. Fine-Tuning Large Language Models: A Comprehensive Guide, accessed May 5, 2025, <https://machinelearningartificialintelligence.com/blog/finetuningllms.html>
207. Spike No More: Stabilizing the Pre-training of Large Language Models - arXiv, accessed May 5, 2025, <https://arxiv.org/html/2312.16903v2>
208. Which LLM to Play? Convergence-Aware Online Model Selection with Time-Increasing Bandits - arXiv, accessed May 5, 2025, <https://arxiv.org/html/2403.07213v1>
209. Deciphering Cross-Modal Alignment in Large Vision-Language Models with Modality Integration Rate - arXiv, accessed May 5, 2025,

- <https://arxiv.org/html/2410.07167v1>
210. Evaluating Large Language Models with Enterprise Benchmarks - ACL Anthology, accessed May 5, 2025, <https://aclanthology.org/2025.naacl-industry.40.pdf>
 211. Enterprise Benchmarks for Large Language Model Evaluation - arXiv, accessed May 5, 2025, <https://arxiv.org/html/2410.12857v1>
 212. 10 Must-Know LLM Benchmarks for Comprehensive Analysis, accessed May 5, 2025, <https://datasciencedojo.com/blog/llm-benchmarks-for-evaluation/>
 213. Fine-Tuning LLMs: Expert Guide to Task-Specific AI Models - Rapid Innovation, accessed May 5, 2025, <https://www.rapidinnovation.io/post/for-developers-step-by-step-guide-to-fine-tuning-llms-for-specific-tasks>
 214. NeurIPS Poster Post-Hoc Reversal: Are We Selecting Models Prematurely?, accessed May 5, 2025, <https://neurips.cc/virtual/2024/poster/96745>
 215. Language models scale reliably with over-training and on downstream tasks - OpenReview, accessed May 5, 2025, <https://openreview.net/forum?id=5tOVh81aze-eld=2u3btz122m>
 216. Troubleshooting PPO Training Instability - ApX Machine Learning, accessed May 5, 2025, <https://apxml.com/courses/rllhf-reinforcement-learning-human-feedback/chapter-4-rl-ppo-fine-tuning/troubleshooting-ppo-instability>
 217. Vanishing and Exploding Gradients in Neural Network Models: Debugging, Monitoring, and Fixing - Neptune.ai, accessed May 5, 2025, <https://neptune.ai/blog/vanishing-and-exploding-gradients-debugging-monitoring-fixing>
 218. Exploding Gradient Explained: How To Detect & Overcome It [6 Best Practices], accessed May 5, 2025, <https://spotintelligence.com/2023/12/06/exploding-gradient-problem/>
 219. Vanishing and Exploding Gradients Problems in Deep Learning ..., accessed May 5, 2025, <https://www.geeksforgeeks.org/vanishing-and-exploding-gradients-problems-in-deep-learning/>
 220. Vanishing and Exploding Gradients in Deep Neural Networks - Analytics Vidhya, accessed May 5, 2025, <https://www.analyticsvidhya.com/blog/2021/06/the-challenge-of-vanishing-exploding-gradients-in-deep-neural-networks/>
 221. Understanding Gradient Clipping (and How It Can Fix Exploding Gradients Problem), accessed May 5, 2025, <https://neptune.ai/blog/understanding-gradient-clipping-and-how-it-can-fix-exploding-gradients-problem>
 222. What Is Instruction Tuning? | IBM, accessed May 5, 2025, <https://www.ibm.com/think/topics/instruction-tuning>
 223. Instruction Tuning: What is fine-tuning? - DataScientest, accessed May 5, 2025, <https://datascientest.com/en/instruction-tuning-what-is-fine-tuning>
 224. Instruction Tuning for Large Language Models: A Survey - arXiv, accessed May

- 5, 2025, <https://arxiv.org/html/2308.10792v5>
225. An Overview of Instruction Tuning Data - rudr.io, accessed May 5, 2025, <https://www.rudr.io/an-overview-of-instruction-tuning-data/>
226. yaodongC/awesome-instruction-dataset: A collection of open-source dataset to train instruction-following LLMs (ChatGPT, LLaMA, Alpaca) - GitHub, accessed May 5, 2025, <https://github.com/yaodongC/awesome-instruction-dataset>
227. Large-Scale Data Selection for Instruction Tuning - arXiv, accessed May 5, 2025, <https://arxiv.org/html/2503.01807v1>
228. What is RLHF? - Reinforcement Learning from Human Feedback Explained - AWS, accessed May 5, 2025, <https://aws.amazon.com/what-is/reinforcement-learning-from-human-feedback/>
229. Understanding Reinforcement Learning from Human Feedback (RLHF) in LLMs - Turing, accessed May 5, 2025, <https://www.turing.com/resources/rlhf-in-llms>
230. Reinforcement learning with human feedback (RLHF) for LLMs - SuperAnnotate, accessed May 5, 2025, <https://www.superannotate.com/blog/rlhf-for-llm>
231. Reinforcement Learning From Human Feedback (RLHF) For LLMs, accessed May 5, 2025, <https://neptune.ai/blog/reinforcement-learning-from-human-feedback-for-llms>
232. Reinforcement learning from human feedback - Wikipedia, accessed May 5, 2025, https://en.wikipedia.org/wiki/Reinforcement_learning_from_human_feedback
233. How to Implement Reinforcement Learning from Human Feedback (RLHF) - Labelbox, accessed May 5, 2025, <https://labelbox.com/guides/how-to-implement-reinforcement-learning-from-human-feedback-rlhf/>
234. The 5 Steps of Reinforcement Learning with Human Feedback - Appen, accessed May 5, 2025, <https://www.appen.com/blog/the-5-steps-of-reinforcement-learning-with-human-feedback>
235. Secrets of RLHF in Large Language Models Part II: Reward Modeling - arXiv, accessed May 5, 2025, <https://arxiv.org/html/2401.06080v2>
236. Train a reward model for RLHF - Argilla 1.11 documentation, accessed May 5, 2025, <https://docs.v1.argilla.io/en/v1.11.0/guides/llms/examples/train-reward-model-rlhf.html>
237. Reward Modeling | RLHF Book by Nathan Lambert, accessed May 5, 2025, <https://rlhfbook.com/c/07-reward-models.html>
238. PPO (Proximal Policy Optimization) - AI Engineering Academy, accessed May 5, 2025, <https://aiengineering.academy/LLM/TheoryBehindFinetuning/PPO/>
239. Pairwise Proximal Policy Optimization: Harnessing Relative Feedback for LLM Alignment, accessed May 5, 2025, <https://openreview.net/forum?id=JzAuFCKioV>
240. The State of Reinforcement Learning for LLM Reasoning - Sebastian Raschka, accessed May 5, 2025, <https://sebastianraschka.com/blog/2025/the-state-of-reinforcement-learning-for>

[-llm-reasoning.html](#)

241. A vision researcher's guide to some RL stuff: PPO & GRPO - Yuge (Jimmy) Shi, accessed May 5, 2025, <https://yugeten.github.io/posts/2025/01/ppogrp/>
242. TRL + PPO + Using Conditioned Reference Model - Intermediate - Hugging Face Forums, accessed May 5, 2025, <https://discuss.huggingface.co/t/trl-ppo-using-conditioned-reference-model/136361>
243. TRL - Transformer Reinforcement Learning - Hugging Face, accessed May 5, 2025, <https://huggingface.co/docs/trl/index>
244. Quickstart - Hugging Face, accessed May 5, 2025, <https://huggingface.co/docs/trl/quickstart>
245. PPO Trainer - Hugging Face, accessed May 5, 2025, https://huggingface.co/docs/trl/v0.7.4/ppo_trainer
246. PPO Trainer - Hugging Face, accessed May 5, 2025, https://huggingface.co/docs/trl/ppo_trainer
247. Claude's Constitution - Anthropic, accessed May 5, 2025, <https://www.anthropic.com/news/claudes-constitution>
248. Safety Layers in Aligned Large Language Models: The Key to LLM Security - arXiv, accessed May 5, 2025, <https://arxiv.org/html/2408.17003v5>
249. Responsible AI in action: How Data Reply red teaming supports generative AI safety on AWS | AWS Machine Learning Blog, accessed May 5, 2025, <https://aws.amazon.com/blogs/machine-learning/responsible-ai-in-action-how-data-reply-red-teaming-supports-generative-ai-safety-on-aws/>
250. The Ultimate Guide to Red Teaming LLMs and Adversarial Prompts ..., accessed May 5, 2025, <https://kili-technology.com/large-language-models-llms/red-teaming-llms-and-adversarial-prompts>
251. Prompt Injection Attacks on LLMs - HiddenLayer, accessed May 5, 2025, <https://hiddenlayer.com/innovation-hub/prompt-injection-attacks-on-llms/>
252. What Is a Prompt Injection Attack? [Examples & Prevention] - Palo Alto Networks, accessed May 5, 2025, <https://www.paloaltonetworks.com/cyberpedia/what-is-a-prompt-injection-attack>
253. Prompt Injection & the Rise of Prompt Attacks: All You Need to Know ..., accessed May 5, 2025, <https://www.lakera.ai/blog/guide-to-prompt-injection>
254. Jailbreaking LLMs: A Comprehensive Guide (With Examples) - Promptfoo, accessed May 5, 2025, <https://www.promptfoo.dev/blog/how-to-jailbreak-llms/>
255. Collective Constitutional AI: Aligning a Language Model with Public ..., accessed May 5, 2025, <https://www.anthropic.com/research/collective-constitutional-ai-aligning-a-language-model-with-public-input>
256. Claude AI's Ethical Framework | What is Constitutional AI?, accessed May 5, 2025, <https://claudeaihub.com/constitutional-ai/>
257. Constitutional AI, accessed May 5, 2025, <https://www.constitutional.ai/>
258. Quick Guide to OWASP Top 10 LLM: Threats, Examples & Prevention -

- Tigera.io, accessed May 5, 2025,
<https://www.tigera.io/learn/guides/llm-security/owasp-top-10-llm/>
259. OWASP Top 10 LLM, Updated 2025: Examples & Mitigation Strategies, accessed May 5, 2025,
<https://www.oligo.security/academy/owasp-top-10-llm-updated-2025-examples-and-mitigation-strategies>
260. Decoding Strategies: How LLMs Choose The Next Word - AssemblyAI, accessed May 5, 2025,
<https://www.assemblyai.com/blog/decoding-strategies-how-llms-choose-the-next-word>
261. Transformers Key-Value Caching Explained - neptune.ai, accessed May 5, 2025, <https://neptune.ai/blog/transformers-key-value-caching>
262. KV Caching Explained: Optimizing Transformer Inference Efficiency - Hugging Face, accessed May 5, 2025, <https://huggingface.co/blog/not-lain/kv-caching>
263. Key Value Cache | Continuum Labs, accessed May 5, 2025,
<https://training.continuumlabs.ai/inference/why-is-inference-important/key-value-cache>
264. LLM Inference with KV Cache Explained - Gala codes, accessed May 5, 2025,
<https://galacodes.hashnode.dev/introduction-to-llm-inferencing>
265. Decoding Strategies in Large Language Models - Hugging Face, accessed May 5, 2025, <https://huggingface.co/blog/mlabonne/decoding-strategies>
266. Token Sampling Methods - aman.ai, accessed May 5, 2025,
<https://aman.ai/primers/ai/token-sampling/>
267. Top-k and Top-p Decoding - Aussie AI, accessed May 5, 2025,
<https://www.aussieai.com/research/top-k-decoding>
268. ONNX | Home, accessed May 5, 2025, <https://onnx.ai/>
269. ONNX Runtime and models - Azure Machine Learning | Microsoft ..., accessed May 5, 2025,
<https://learn.microsoft.com/en-us/azure/machine-learning/concept-onnx?view=azureml-api-2>
270. Export to ONNX - Hugging Face, accessed May 5, 2025,
<https://huggingface.co/docs/transformers/v4.43.2/serialization>
271. Optimizing Large Language Model Performance with ONNX on DataRobot MLOps, accessed May 5, 2025,
<https://www.datarobot.com/blog/optimizing-large-language-model-performance-with-onnx-on-datarobot-mlops/>
272. Accelerate AI & Machine Learning Workflows | NVIDIA Run:ai, accessed May 5, 2025, <https://www.run.ai/guides/generative-ai/tensorrt-llm>
273. TensorRT-LLM: A Comprehensive Guide to Optimizing Large ..., accessed May 5, 2025,
<https://www.unite.ai/tensorrt-llm-a-comprehensive-guide-to-optimizing-large-language-model-inference-for-maximum-performance/>
274. Overview — TensorRT-LLM - GitHub Pages, accessed May 5, 2025,
<https://nvidia.github.io/TensorRT-LLM/overview.html>
275. NVIDIA/TensorRT-LLM: TensorRT-LLM provides users with an easy-to-use

Python API to define Large Language Models (LLMs) and support state-of-the-art optimizations to perform inference efficiently on NVIDIA GPUs. TensorRT-LLM also contains components to create Python and C++ runtimes that orchestrate the - GitHub, accessed May 5, 2025,

<https://github.com/NVIDIA/TensorRT-LLM>

276. LLM inference: vLLM - BentoML Documentation, accessed May 5, 2025,

<https://docs.bentoml.com/en/latest/examples/vllm.html>

277. What is vLLM: A Guide to Quick Inference - Hyperstack, accessed May 5, 2025,

<https://www.hyperstack.cloud/blog/case-study/what-is-vllm-a-guide-to-quick-inference>

278. Meet vLLM: For faster, more efficient LLM inference and serving - Red Hat, accessed May 5, 2025,

<https://www.redhat.com/en/blog/meet-vllm-faster-more-efficient-llm-inference-and-serving>

279. FasterTransformer: The Basics and a Quick Tutorial - Run:ai, accessed May 5, 2025, <https://www.run.ai/guides/ai-open-source-projects/fastertransformer>

280. Accelerated Inference for Large Transformer Models Using NVIDIA ..., accessed May 5, 2025,

<https://developer.nvidia.com/blog/accelerated-inference-for-large-transformer-models-using-nvidia-fastertransformer-and-nvidia-triton-inference-server/>

281. Increasing Inference Acceleration of KoGPT with NVIDIA FasterTransformer, accessed May 5, 2025,

<https://developer.nvidia.com/blog/increasing-inference-acceleration-of-kogpt-with-fastertransformer/>

282. What is DeepSpeed-MII? Features & Getting Started - Deepchecks, accessed May 5, 2025, <https://www.deepchecks.com/llm-tools/deepspeed-mii/>

283. deepspeedai/DeepSpeed-MII: MII makes low-latency and ... - GitHub, accessed May 5, 2025, <https://github.com/deepspeedai/DeepSpeed-MII>

284. DeepSpeed MII vs. TensorRT LLM: Choosing the Best Inference Library for Large Language Models - Inferless, accessed May 5, 2025,

<https://www.inferless.com/learn/deepspeed-mii-vs-tensorrt-llm-a-complete-guide-to-optimized-large-language-model-inference>

285. Hugging Face's Text Generation Inference Toolkit for LLMs - A Game Changer in AI, accessed May 5, 2025,

<https://www.datacamp.com/tutorial/hugging-faces-text-generation-inference-toolkit-for-llms>

286. Text Generation Inference - Hugging Face, accessed May 5, 2025, <https://huggingface.co/text-generation-inference>

287. Text Generation Inference - Hugging Face, accessed May 5, 2025, <https://huggingface.co/docs/text-generation-inference/index>

288. A Guide to Quantization in LLMs | Syml.ai, accessed May 5, 2025, <https://syml.ai/developers/blog/a-guide-to-quantization-in-llms/>

289. Quantization - Hugging Face, accessed May 5, 2025,

https://huggingface.co/docs/transformers/main_classes/quantization

290. The Ultimate Handbook for LLM Quantization | Towards Data Science,
accessed May 5, 2025,
<https://towardsdatascience.com/the-ultimate-handbook-for-llm-quantization-88bb7cb0d9d7/>