# A Survey of Mathematical Optimization Methods for Large Language Model Training

**Samrat Kar (BL.SC.R4CSE24007) | Amrita Vishwa Vidyapeetham**

## Abstract

The unprecedented scale and capabilities of modern Large Language Models (LLMs) are underpinned by sophisticated mathematical optimization techniques. Training these models, which can have hundreds of billions or even trillions of parameters, presents a formidable challenge in navigating a high-dimensional, non-convex loss landscape. This paper provides a comprehensive survey of the pivotal optimization algorithms that have enabled the LLM revolution. We begin with the foundational principles of Stochastic Gradient Descent (SGD) and its variants before delving into the adaptive methods, particularly Adam and its successor, AdamW, which have become the de-facto standards for LLM training. We analyze the theoretical underpinnings of these methods, their practical advantages, and the specific reasons for their widespread adoption. Furthermore, we explore the cutting edge of optimization research, documenting emerging techniques designed to enhance training efficiency, reduce memory footprints, and improve model generalization. These new frontiers include 8-bit optimizers, approximations of second-order methods, and sharpness-aware minimization, all of which are poised to define the next generation of LLM training paradigms.

In the area of aviation application of LLMs, the 8-bit optimization solutions can be utilized to drastically reduce the memory footprint upto 75%. This makes it possible for onboarded LLMs into the embedded system of avaiation flight deck, facilitating Edge AI by reducing size, weight and power.

## 1. Introduction

The last five years have witnessed a paradigm shift in artificial intelligence, driven largely by the advent of Large Language Models (LLMs) such as GPT-3/4 (Brown et al., 2020), PaLM (Chowdhery et al., 2022), and Llama (Touvron et al., 2023). These models have demonstrated remarkable abilities in natural language understanding, generation, and reasoning. However, their power comes at an immense computational cost. Training an LLM involves minimizing a loss function over a parameter space of staggering dimensionality, a task that is both computationally intensive and algorithmically complex.

The choice of optimization algorithm is not merely a detail but a critical component that dictates the feasibility, speed, and ultimate performance of the trained model. A well-chosen optimizer can navigate the treacherous, non-convex loss landscape to find a "good" local minimum—one that generalizes well to unseen data—within a practical timeframe. A poor choice can lead to slow convergence, divergence, or convergence to a sharp minimum with poor generalization properties.

This survey aims to provide a clear and detailed overview of the mathematical optimization methods that form the bedrock of LLM training. We will trace the evolution of these methods, from the classic SGD to the current industry-standard AdamW, explaining the mathematical intuition and practical significance of each step. Finally, we will cast our gaze toward the future, outlining active research areas that promise to make LLM training more accurate, efficient, and accessible.

# 2. Introduction to Gradient Descent

Gradient descent is a first-order iterative optimization algorithm used to find the minimum of a function. In machine learning, we use gradient descent to minimize a loss function by iteratively moving in the direction of steepest descent as defined by the negative of the gradient.

## 2.1 The Basic Concept

At its core, gradient descent works by:

1. Starting at an initial point
2. Computing the gradient (direction of steepest ascent) of the loss function
3. Taking a step in the opposite direction (steepest descent)
4. Repeating until convergence

Mathematically, we update parameters θ using the rule:

$$\theta_{t+1} = \theta_t - \eta \nabla_\theta J(\theta_t)$$

Where:

- $\theta_t$ is the parameter vector at iteration t
- $\eta$ is the learning rate (step size)
- $\nabla_\theta J(\theta_t)$ is the gradient of the objective function J with respect to parameters θ

## 2.2 Types of Gradient Descent

There are three main variants of gradient descent:

**Batch Gradient Descent**: Computes the gradient using the entire dataset.

- Pros: Stable convergence; theoretically can reach global minimum for convex functions
- Cons: Computationally expensive for large datasets; can be slow

**Stochastic Gradient Descent (SGD)**: Updates parameters using only one sample at a time.

- Pros: Fast; can escape local minima due to noise
- Cons: High variance in parameter updates; may never converge exactly

**Mini-batch Gradient Descent**: Updates parameters using a small batch of samples.

- Pros: Balances computational efficiency with update stability
- Cons: Requires tuning of batch size

## 2.3 Challenges with Basic Gradient Descent

Despite its simplicity, gradient descent faces several challenges:

1. **Learning Rate Selection**: Too small leads to slow convergence; too large can cause divergence

2. **Saddle Points**: Areas where the gradient is zero but not a minimum
3. **Plateaus**: Regions with very small gradients that slow learning
4. **Ravines**: Areas where the surface curves more steeply in one dimension than others
5. **Local Minima**: Points where the algorithm may get stuck

These challenges led to the development of more sophisticated optimization algorithms, including ADAM.

# 3. Advanced Gradient Descent Techniques

Before diving into ADAM, let's explore some key improvements to the basic gradient descent algorithm.

## 3.1 Momentum

Momentum adds a fraction of the previous update vector to the current update:

$$v_t = \gamma v_{t-1} + \eta \nabla_\theta J(\theta_t)$$

$$\theta_{t+1} = \theta_t - v_t$$

Where:

- $v_t$ is the velocity vector (momentum)
- $\gamma$ is the momentum coefficient (typically 0.9)

**Benefits**:

- Accelerates convergence
- Helps overcome local minima and plateaus
- Reduces oscillations in ravines

## 3.2 RMSprop

RMSprop (Root Mean Square Propagation) adapts the learning rate for each parameter based on the historical gradient:

$$s_t = \beta s_{t-1} + (1 - \beta)(\nabla_\theta J(\theta_t))^2$$

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{s_t + \epsilon}} \nabla_\theta J(\theta_t)$$

Where:

- $s_t$ is the exponentially decaying average of squared gradients
- $\beta$ is the decay rate (typically 0.9)
- $\epsilon$ is a small constant to prevent division by zero

**Benefits**:

- Adaptive learning rates for each parameter
- Handles different scales of gradients effectively
- Mitigates the effects of vanishing/exploding gradients

# 4. Foundational Methods: From SGD to Momentum

The core of training deep neural networks is to find the set of parameters $\theta$ that minimizes a loss function $L(\theta)$, typically averaged over a large dataset. The primary tool for this is gradient descent.

## 4.1. Stochastic Gradient Descent (SGD)

Given the massive size of modern datasets, computing the true gradient over the entire dataset for each parameter update is infeasible. Stochastic Gradient Descent (SGD) addresses this by approximating the gradient using a small, random subset of the data called a mini-batch. The parameter update rule is:

$$\theta_{t+1} = \theta_t - \eta \nabla L(\theta_t; x^{(i:i+n)}; y^{(i:i+n)})$$

where:

- $\theta_t$ is the parameter vector at step $t$.
- $\eta$ is the learning rate, a crucial hyperparameter.
- $\nabla L(\theta_t; \dots)$ is the gradient of the loss with respect to $\theta_t$ computed on a mini-batch of size $n$.

While simple and computationally light, vanilla SGD suffers from high variance in its updates, leading to noisy convergence paths. It can struggle in regions with steep ravines, where it tends to oscillate across the ravine instead of moving along its bottom, slowing down progress.

## 4.2. SGD with Momentum

To mitigate the issues of SGD, the Momentum method was introduced (Polyak, 1964). It adds a fraction $\beta$ of the previous update vector to the current one, acting as a velocity term that accumulates over steps.

$$v_{t+1} = \beta v_t + \eta \nabla L(\theta_t)$$
$$\theta_{t+1} = \theta_t - v_{t+1}$$

The momentum term $v_t$ helps to dampen oscillations in directions of high curvature and accelerate convergence in consistent directions. For LLMs, while not the primary choice, the concept of momentum is a critical building block for more advanced optimizers.

# 5. The Era of Adaptive Methods: Adam and AdamW

A key limitation of SGD and its momentum variant is the use of a single, global learning rate $\eta$ for all parameters. In a model as complex as an LLM, some parameters (e.g., those associated with rare tokens) may require larger updates, while others need finer adjustments. This observation led to the development of adaptive learning rate methods.

## 5.1. Precursors: AdaGrad and RMSProp

- **AdaGrad (Adaptive Gradient):** (Duchi et al., 2011) adapts the learning rate for each parameter by dividing it by the square root of the sum of all past squared gradients. This works well for sparse data but has a major flaw: the accumulated sum grows indefinitely, causing the learning rate to shrink to near zero, prematurely halting training.
- **RMSProp (Root Mean Square Propagation):** (Hinton, unpublished) fixes AdaGrad's issue by using an exponentially decaying moving average of squared gradients, preventing the denominator from growing monotonically.

## 5.2. Adam: Adaptive Moment Estimation

Adam (Kingma & Ba, 2015) is arguably the most significant optimization algorithm in the deep learning era. It elegantly combines the ideas of **Momentum** (first-order moment estimation) and **RMSProp** (second-order moment estimation).

At each step $t$, Adam computes:

1. **Decaying average of past gradients (First Moment, Momentum):**
$m_t = \beta_1 m_{t-1} + (1 - \beta_1)\nabla L(\theta_t)$
2. **Decaying average of past squared gradients (Second Moment, Uncentered Variance):**
$v_t = \beta_2 v_{t-1} + (1 - \beta_2)(\nabla L(\theta_t))^2$

Since $m_t$ and $v_t$ are initialized to zero, they are biased towards zero, especially during the initial steps. Adam corrects for this bias:

3. **Bias Correction:**
$\hat{m}_t = \frac{m_t}{1 - \beta_1^t}$
$\hat{v}_t = \frac{v_t}{1 - \beta_2^t}$

Finally, the parameters are updated using these corrected estimates:

4. **Parameter Update:**
$\theta_{t+1} = \theta_t - \eta \frac{\hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon}$

Here, $\beta_1$ and $\beta_2$ are decay rates (typically ~0.9 and ~0.999), and $\epsilon$ is a small constant for numerical stability. The term $\sqrt{\hat{v}_t}$ provides a per-parameter adaptive learning rate, scaling down updates for parameters with consistently large gradients and vice versa. Adam's robustness and fast initial convergence made it the go-to optimizer for years.

ADAM (Adaptive Moment Estimation) combines the benefits of both momentum and RMSprop, making it one of the most popular optimization algorithms in deep learning.

## 5.3 The ADAM Algorithm

ADAM maintains two moving averages:

1. First moment (mean) of gradients - similar to momentum
2. Second moment (uncentered variance) of gradients - similar to RMSprop

The update rules are:

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1)\nabla_\theta J(\theta_t)$$

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2)(\nabla_\theta J(\theta_t))^2$$

To correct the bias in these estimates (which start at zero):
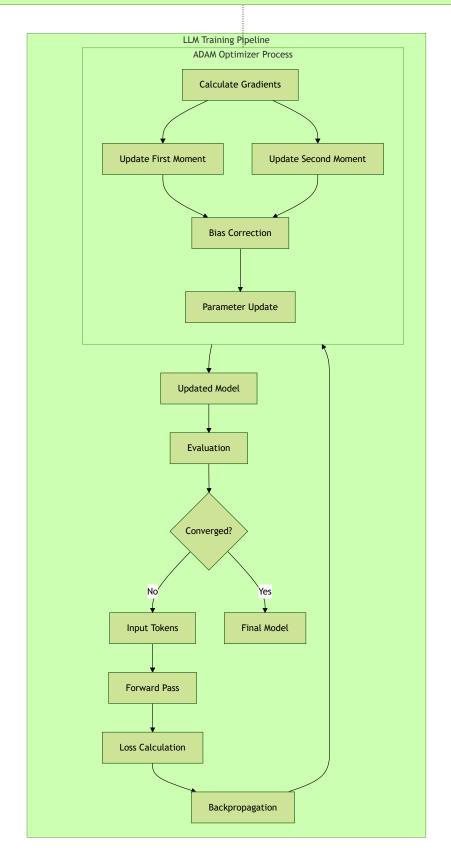
$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t}$$

$$\hat{v}_t = \frac{v_t}{1 - \beta_2^t}$$

And finally, the parameter update:

$$\theta_{t+1} = \theta_t - \frac{\eta \hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon}$$

Where:

- $m_t$ and $v_t$ are the first and second moment estimates
- $\beta_1$ and $\beta_2$ are the decay rates (typically 0.9 and 0.999)
- $\hat{m}_t$ and $\hat{v}_t$ are the bias-corrected moment estimates
- $\epsilon$ is a small constant (typically $10^{-8}$)

## LLM-Specific Enhancements

| Learning Rate Scheduling | Weight Decay | Gradient Clipping | Mixed Precision | Gradient Accumulation |
|---|---|---|---|---|

## LLM Training Pipeline

### ADAM Optimizer Process

Calculate Gradients

→ Update First Moment

→ Update Second Moment

Update First Moment → Bias Correction

Update Second Moment → Bias Correction

Bias Correction → Parameter Update

Parameter Update → Updated Model

Updated Model → Evaluation

Evaluation → Converged?

Converged? —No→ Input Tokens

Converged? —Yes→ Final Model

Input Tokens → Forward Pass

Forward Pass → Loss Calculation

Loss Calculation → Backpropagation

Backpropagation → (Calculate Gradients)

## 5.4 Key Benefits of ADAM

1. **Adaptive Learning Rates**: Different parameters have different effective learning rates
2. **Momentum**: Captures the concept of momentum for faster convergence
3. **Bias Correction**: Corrects the bias in moment estimates
4. **Robust to Noisy Gradients**: Works well with sparse gradients and non-stationary objectives
5. **Minimal Hyperparameter Tuning**: Often works well with default values

## 5.5 AdamW: Decoupled Weight Decay

LLMs are massively over-parameterized and thus highly prone to overfitting. A standard regularization technique is L2 regularization (also known as weight decay), which adds a penalty term $\frac{\lambda}{2}||\theta||^2$ to the loss function. When applied to Adam, the gradient of this term, $\lambda\theta$, gets incorporated into the adaptive learning rate mechanism.

Loshchilov & Hutter (2019) observed that this coupling is suboptimal. In Adam, the weight decay for a parameter $\theta_i$ is scaled by $1/\sqrt{\hat{v}_{t,i}}$. This means parameters with large historical gradients (large $\hat{v}_{t,i}$) receive a smaller effective weight decay, which is counter-intuitive. Large weights, which we often want to decay the most, might not be decayed effectively if they also have large gradients.

**AdamW (Adam with Decoupled Weight Decay)** solves this by decoupling the weight decay from the gradient update. The update rule is modified as follows:

1. Perform the standard Adam update based only on the loss gradient $\nabla L(\theta_t)$:
   $\theta'_{t+1} = \theta_t - \eta \frac{\hat{m}_t}{\sqrt{\hat{v}_t}+\epsilon}$
2. Apply the weight decay directly to the parameters *after* the update:
   $\theta_{t+1} = \theta'_{t+1} - \eta\lambda\theta_t$

This seemingly small change ensures that weight decay acts as intended—shrinking all weights towards zero at a rate proportional to their magnitude, independent of their gradient history. This improved regularization has proven critical for achieving state-of-the-art results in LLM training, making **AdamW the current de-facto standard optimizer.**

## 5.6 ADAM Variants

Several variants of ADAM have been proposed to address specific challenges:

**AdamW**: Decouples weight decay from the gradient update to improve generalization.

**AMSGrad**: Maintains the maximum of all past squared gradients to address potential convergence issues.

**AdaBelief**: Incorporates "belief" in the gradient update to improve stability and generalization.

# 6. ADAM in Large Language Models (LLMs)

Large Language Models like GPT, BERT, and Claude present unique optimization challenges due to their massive parameter space and complex loss landscapes.

## 6.1 Why ADAM is Preferred for LLMs

1. **Parameter Scale**: LLMs contain billions of parameters that benefit from adaptive learning rates
2. **Training Stability**: ADAM's bias correction and momentum properties help maintain stable updates
3. **Convergence Speed**: Faster convergence is critical when training models that might take weeks or months
4. **Different Learning Dynamics**: Different layers in LLMs often need different effective learning rates
5. **Sparse Gradients**: Many parameters in LLMs receive sparse gradient updates, which ADAM handles well

## 6.2 Implementation Considerations in LLMs

When applying ADAM to LLMs, several practical considerations come into play:

**Learning Rate Scheduling**: Often combined with:

- Warmup periods (gradually increasing learning rate)
- Decay schedules (gradually decreasing learning rate)
- Cosine annealing schedules

**Weight Decay**: Usually implemented as AdamW to properly separate weight decay from gradient updates.

**Mixed Precision Training**: To save memory, gradients are often computed in lower precision (FP16) while maintaining optimizer states in higher precision (FP32).

**Gradient Accumulation**: Updates are often performed after accumulating gradients from multiple mini-batches to simulate larger batch sizes.

**Gradient Clipping**: To prevent exploding gradients, the norm of gradients is often clipped to a maximum value.

## 6.3 Hyperparameter Selection for LLMs

Typical hyperparameters for ADAM when training LLMs:

- $\beta_1$: 0.9 (first moment decay rate)
- $\beta_2$: 0.999 (second moment decay rate)
- $\epsilon$: 1e-8 (numerical stability constant)
- Learning rate: 1e-4 to 5e-5 (model-specific)
- Weight decay: 0.01 to 0.1 (model-specific)

## 6.4 ADAM's Impact on LLM Training

The adoption of ADAM has been crucial for enabling:

1. Training larger models without excessive hyperparameter tuning
2. Faster convergence, reducing the computational resources needed
3. Better handling of the complex loss landscapes in LLMs
4. More stable training across different architectural choices
5. Effective fine-tuning of pre-trained models

# 7. Case Studies and Practical Examples

## 7.1 Training GPT Models with ADAM

OpenAI's GPT models rely heavily on ADAM optimization. For GPT-3, researchers used:

- ADAM optimizer with $\beta_1$=0.9, $\beta_2$=0.95, $\varepsilon$=1e-8
- Learning rate: 6e-5 with cosine decay to zero
- Batch size of 3.2M tokens
- Gradient clipping at 1.0
- Weight decay of 0.1

## 7.2 BERT Training Approach

Google's BERT used:

- ADAM optimizer with $\beta_1$=0.9, $\beta_2$=0.999, $\varepsilon$=1e-8
- Learning rate: 1e-4 with linear decay
- Warmup over first 10,000 steps
- L2 weight decay of 0.01

## 7.3 LLaMA Training

Meta's LLaMA models used:

- AdamW optimizer
- Learning rate: 3e-4 with cosine decay to 1e-5
- Weight decay of 0.1
- Batch size of 4M tokens

# 8. Limitations and Challenges

Despite its popularity, ADAM has some limitations when applied to LLMs:

1. **Memory Requirements**: Maintaining moment estimates doubles the memory needed compared to SGD
2. **Generalization Gap**: Some research suggests ADAM may generalize slightly worse than SGD in some cases
3. **Sensitivity to Learning Rate**: Still requires careful tuning of initial learning rate
4. **Computational Overhead**: Slightly more computation per step than simpler optimizers
5. **Distributed Training Complexity**: Requires careful implementation in distributed settings

# 9. Future Directions

While ADAM has become the standard optimizer for training LLMs, active research continues to push the boundaries of what's possible in optimization. As models grow larger and training becomes more resource-intensive, several promising research directions may shape the next generation of optimization algorithms. This expanded section explores these future directions in greater detail.

Research continues on improving optimization for LLMs:

1. **Memory-Efficient Optimizers**: Reducing the memory footprint without sacrificing performance
2. **Second-Order Methods**: Incorporating curvature information more effectively
3. **Neural Optimizers**: Using neural networks to learn optimal update rules
4. **Sparsity-Aware Optimizers**: Better handling of the sparse nature of updates in large models
5. **Scale-Invariant Optimizers**: Automatically adapting to different parameter scales

# 10. New Frontiers and Research Directions

While AdamW is a robust workhorse, the scale of future models demands further innovation. Research is now focused on several key areas to improve efficiency and accuracy.

## 10.1. Memory-Efficient Optimizers: 8-bit Adam

A significant bottleneck in training large models is the memory footprint of the optimizer states. For every model parameter stored in 32-bit floating-point (FP32), a standard Adam optimizer requires two additional parameters for its states ($m_t$ and $v_t$), also in FP32. For a 175B parameter model, the optimizer states alone can consume over 1.4 Terabytes of GPU memory.

**8-bit Optimizers** (Dettmers et al., 2022) address this by quantizing the optimizer states. The core idea is to store the 32-bit momentum and variance states in 8-bit format. This is non-trivial, as naive quantization would lead to significant precision loss and training instability. The key innovations include:

- **Block-wise Quantization:** Instead of quantizing the entire tensor with one set of scaling factors, the tensor is divided into smaller blocks. Each block is quantized independently, preserving the relative magnitudes of values within the block.
- **Dynamic Quantization:** A special non-uniform 8-bit data type (e.g., quantile quantization) is used, which allocates more precision to values that occur more frequently in the optimizer state distribution.

These techniques can reduce the optimizer's memory footprint by a factor of 4x with negligible impact on model performance, making it possible to train larger models on existing hardware or fine-tune models on consumer-grade GPUs.

## 10.2. Approximating Second-Order Information: Shampoo

First-order methods like Adam only use the gradient (slope). Second-order methods, which use the Hessian matrix (curvature of the loss landscape), can converge much faster in theory. However, for a model with $N$ parameters, computing and inverting the Hessian is an $O(N^3)$ operation, which is utterly infeasible for LLMs where $N$ is in the billions.

**Shampoo (Shampoo: Preconditioned Stochastic Tensor Optimization)** (Gupta et al., 2018; Anil et al., 2020) is a practical quasi-Newton method that approximates the Hessian. It works by creating block-diagonal or block-tridiagonal preconditioners that capture local curvature information without forming the full Hessian. These preconditioners are themselves learned and updated efficiently.

For a weight matrix $W \in \mathbb{R}^{m \times n}$, Shampoo maintains separate preconditioners for the input and output dimensions. This drastically reduces the complexity from inverting an $(mn \times mn)$ matrix to inverting two smaller matrices ($m \times m$ and $n \times n$). Recent work has shown that Shampoo can outperform AdamW in terms of final model quality and convergence speed on large-scale models, although it can have higher computational overhead per step.

## 10.3. Improving Generalization: Sharpness-Aware Minimization (SAM)

It has been observed that the flatness of the loss minimum found by an optimizer correlates strongly with a model's generalization ability. Models that converge to sharp, narrow minima are often brittle and perform poorly on out-of-distribution data.

**Sharpness-Aware Minimization (SAM)** (Foret et al., 2021) is an optimization procedure that explicitly seeks out flat minima. Instead of just minimizing the loss $L(\theta)$, SAM minimizes the loss in a neighborhood around $\theta$. The objective is:

$$\min_\theta \max_{||\epsilon||_p \leq \rho} L(\theta + \epsilon)$$

This is a minimax problem. SAM solves it with a two-step update:

1. **Ascent Step:** First, find the perturbation $\epsilon$ within a radius $\rho$ that maximally increases the loss. This is approximated by taking a single gradient ascent step:
   $\hat{\epsilon}(\theta) = \rho \frac{\nabla L(\theta)}{||\nabla L(\theta)||_2}$
2. **Descent Step:** Then, update the original parameters $\theta$ using the gradient computed at the perturbed point $\theta + \hat{\epsilon}(\theta)$:
   $\theta_{t+1} = \theta_t - \eta \nabla L(\theta_t + \hat{\epsilon}(\theta_t))$

By doing so, SAM forces the optimizer to find a parameter configuration where not only the loss is low, but the loss also remains low in a surrounding neighborhood. The primary drawback of SAM is that it requires two forward-backward passes per update (one to find $\epsilon$, one to compute the final gradient), doubling the computational cost. Research is ongoing to develop more efficient variants of SAM.

# 10.4 Second-Order Methods

## 10.4.1 Limitations of First-Order Methods

ADAM and other popular optimizers are first-order methods, using only gradient information. They don't account for the curvature of the loss landscape, which can lead to slower convergence in highly curved regions.

## 10.4.2 Advanced Approaches

K-FAC (Kronecker-Factored Approximate Curvature)
K-FAC approximates the Fisher information matrix using Kronecker products, making second-order optimization practical for large neural networks. Research has shown that K-FAC can significantly accelerate training while maintaining or improving final model quality.
Shampoo Optimizer
Shampoo approximates a full-matrix adaptive algorithm with computational efficiency comparable to ADAM. It uses matrix roots to compute preconditioning matrices separately for each tensor dimension, effectively incorporating second-order information.
Quasi-Newton Methods
Limited-memory BFGS (L-BFGS) and other quasi-Newton methods approximate second-order information without explicitly computing the Hessian matrix. Recent research is making these methods more practical for large-scale deep learning.

## 10.5 Neural Optimizers

### 10.5.1 Learning to Optimize

Rather than hand-designing optimization algorithms, neural optimizers use machine learning to discover optimal update rules. This meta-learning approach could potentially discover optimization strategies that outperform human-designed algorithms.

### 10.5.2 Promising Directions

Learned Update Rules
Systems like "Learned Optimizer" use RNNs to learn update rules from data. These learned optimizers can adapt to specific problem domains and loss landscapes, potentially outperforming general-purpose optimizers like ADAM.
Meta-Learning for Optimization
Meta-learning approaches train on collections of similar tasks to learn optimization strategies that transfer well to new tasks. For LLMs, this could mean learning specialized optimizers for language model pre-training and fine-tuning.
Self-Tuning Optimizers
Recent research explores optimizers that can automatically adapt their hyperparameters during training, reducing the need for extensive hyperparameter tuning experiments.

## 10.6 Sparsity-Aware Optimizers

### 10.6.1 The Sparsity Challenge

LLM training often involves sparse gradients, where only a small subset of parameters receive meaningful updates in each step. Standard optimizers don't fully exploit this sparsity pattern.

### 10.6.2 Advanced Approaches

#### Sparse-to-Dense Training

Rather than training the full model from the beginning, sparse-to-dense approaches start by training a smaller subset of parameters and gradually increase model capacity. This can both accelerate training and improve final model quality.

#### Selective Backpropagation

Methods like "Top-K" training only update the parameters with the largest gradients, reducing computational overhead while maintaining model quality. This approach naturally adapts to the sparse nature of updates in large models.
Dynamic Sparse Training
Instead of using a fixed architecture, dynamic sparse training methods adaptively determine which parameters to update during training, focusing computational resources on the most important parameters at each stage.

## 10.7 Scale-Invariant Optimizers

### 10.7.1 The Scale Challenge

Different layers in LLMs may operate at vastly different scales, making it difficult to find a single learning rate that works well for all parameters. While ADAM partially addresses this through adaptive learning rates, further improvements are possible.

### 10.7.2 Advanced Approaches

Layer-wise Adaptive Rate Scaling (LARS)
LARS adjusts learning rates based on the ratio of parameter norm to gradient norm, enabling stable training with much larger batch sizes. This approach has shown promise for extremely large models.
LAMB (Layer-wise Adaptive Moments optimizer for Batch training)
LAMB combines the benefits of LARS with ADAM, providing layer-wise adaptation along with momentum and adaptive learning rates. This has proven effective for pre-training large language models like BERT.
Normalized Direction-Preserving Adam (ND-Adam)
ND-Adam normalizes gradients to preserve their direction while adapting their magnitude, addressing scale variation issues across layers more effectively than standard ADAM.

## 10.8 Federated and Distributed Optimization

### 10.8.1 The Distribution Challenge

Training very large LLMs requires distributed computing across many devices. However, standard optimization algorithms were not designed with massive parallelism in mind.

### 10.8.2 Advanced Approaches

Federated Averaging
This approach trains local models on decentralized data and periodically averages parameters. Recent research is making federated learning more efficient and robust, potentially enabling collaborative training of LLMs across organizations without sharing raw data.
Decentralized Optimization
Rather than using a parameter server architecture, decentralized optimization methods like D-PSGD (Decentralized Parallel Stochastic Gradient Descent) allow each worker to communicate only with its neighbors, reducing communication bottlenecks.
Compressed Communication
Methods like PowerSGD and QSGD reduce communication overhead by compressing gradients before sharing them between workers. This is particularly important for large LLMs where communication can become the primary bottleneck.

## 10.9 Hybrid and Composite Optimizers

### 10.9.1 No One-Size-Fits-All Solution

Different phases of training (early vs. late) and different parts of the model (embeddings vs. attention layers) may benefit from different optimization strategies.

### 10.9.2 Advanced Approaches

Phase-Dependent Optimization
Using different optimizers or hyperparameters at different phases of training. For example, using higher learning rates with ADAM early in training, then switching to SGD with momentum for better generalization in later stages.
Layer-Specific Optimizers
Applying different optimization algorithms to different components of the model. For example, using ADAM for attention mechanisms and a more memory-efficient optimizer for embedding layers.

Composite Optimization

Rather than choosing a single algorithm, composite optimizers combine multiple update rules, selecting the most appropriate one for each parameter based on its recent update history.

## 10.10 Energy-Efficient Optimization

### 10.10.1 The Energy Challenge

Training large LLMs consumes enormous amounts of energy. There's growing interest in optimization methods that reduce energy consumption without sacrificing model quality.

### 10.10.2 Advanced Approaches

Mixed Precision with Dynamic Loss Scaling

Training models using lower precision (FP16 or bfloat16) dramatically reduces memory and computation requirements. Advanced loss scaling techniques maintain numerical stability while benefiting from the efficiency of lower precision.

Efficient Attention Optimizers

Specialized optimizers for attention mechanisms that exploit the unique structure of self-attention to reduce computational overhead during training.

Early Stopping with Uncertainty Estimation

Advanced techniques for determining the optimal point to stop training based on uncertainty estimates, avoiding wasted computation once improvements become marginal.

# 10.11 Hardware-Aware Optimization

## 10.11.1 The Hardware Gap

Standard optimization algorithms don't account for the specific capabilities and limitations of modern AI accelerators like GPUs, TPUs, and specialized ASIC chips.

## 10.11.2 Advanced Approaches

Hardware-Aware Training (HAT)

Optimization methods that adapt to the specific characteristics of the underlying hardware, such as memory hierarchy, compute units, and communication bandwidth.

Training-Aware Chip Design

Co-designing optimization algorithms alongside specialized hardware. For example, chips with dedicated circuits for efficiently computing ADAM updates.

Optimizer-in-Memory Architectures

Novel computing architectures that perform optimizer updates directly in memory, reducing data movement and dramatically improving energy efficiency.

# 10.12 Curriculum and Progressive Optimization

## 10.12.1 The Learning Efficiency Challenge

Training LLMs by randomly sampling from the entire dataset may not be the most efficient approach. There's growing interest in methods that optimize the order and difficulty of training examples.

### 10.12.2 Advanced Approaches

Difficulty-Based Curriculum Learning
Progressively increasing the difficulty of training examples based on the model's current capabilities. This can lead to faster convergence and better final performance.
Importance Sampling
Dynamically adjusting the sampling probability of different examples based on their current difficulty for the model. Examples that produce larger gradients may be sampled more frequently.
Progressive Growing of LLMs
Starting with a smaller model and progressively growing it during training, similar to techniques used in image generation models like GANs. This approach can stabilize early training and reduce overall computational requirements.

## 10.13 Conclusion: Towards More Efficient LLM Training

The future of optimization for LLMs likely involves combinations of these approaches, tailored to specific model architectures, hardware configurations, and training objectives. As models continue to grow in size and complexity, innovations in optimization will be crucial for making training more efficient, accessible, and environmentally sustainable. The most promising direction may be adaptive, composable optimization frameworks that can automatically select and configure the most appropriate techniques based on the specific characteristics of the model and training data. This would reduce the need for extensive hyperparameter tuning and make advanced optimization accessible to a wider range of researchers and practitioners.
As these methods mature, we can expect significant reductions in the computational resources required to train state-of-the-art LLMs, potentially democratizing access to these powerful models and enabling new applications that were previously infeasible due to resource constraints.

# 10. Technical Deep Dive 4.1.1: Mathematical and Implementational Details of 8-bit Optimizers

The high-level motivation for 8-bit optimizers is memory savings. To fully appreciate the ingenuity of these methods, we must analyze the mathematical challenges involved in reducing the precision of optimizer states and the specific techniques developed to overcome them.

## 10.1 The Memory Imperative: A Quantitative Look

A standard AdamW optimizer maintains two state vectors for each model parameter $\theta$:

1. The first moment (momentum) vector, $m_t$.
2. The second moment (uncentered variance) vector, $v_t$.

In a typical training setup, the model parameters ($\theta$), gradients ($\nabla L(\theta)$), and optimizer states ($m_t, v_t$) are all stored in 32-bit floating-point (FP32) format. Therefore, for every single parameter in the model, we store:

- 1 parameter (FP32) = 4 bytes
- 1 gradient (FP32) = 4 bytes (transient, but needed for the update)
- 1 momentum state (FP32) = 4 bytes
- 1 variance state (FP32) = 4 bytes

This means the optimizer states alone consume **8 bytes for every 4-byte parameter**, effectively tripling the base memory requirement for the model's trainable weights. For a 1-billion parameter model, this translates to:

- Parameters: 1B * 4 bytes = 4 GB
- Optimizer States: 1B * 8 bytes = 8 GB
- Total: 12 GB (before accounting for activations, gradients, etc.)

This memory burden is a primary constraint on the maximum model size and batch size that can be used on a given hardware setup.

## 10.2 The Challenge: The Perils of Naive Quantization

The simplest approach would be to cast the 32-bit float values of $m_t$ and $v_t$ directly to 8-bit integers (INT8). An INT8 value can represent integers from -128 to 127. This naive approach fails catastrophically for two main reasons:

1. **Low Dynamic Range:** Optimizer states, especially the variance term $v_t$, can have a massive dynamic range. Some values might be very close to zero, while others (corresponding to parameters with large gradients) can be very large. Squeezing this entire range into just 256 discrete integer levels would lead to massive quantization error. Most values would be clipped to zero or the maximum value, effectively erasing all useful information.
2. **Destructive Precision Loss:** The distribution of values within optimizer state tensors is highly non-uniform. It's typically a bell-shaped curve centered near zero, with long tails. A linear quantization scheme would allocate its 256 levels evenly across the entire range. This means the vast majority of values, which are concentrated in a small region near the mean, would be mapped to only a few integer levels. This "rounding off" of subtle but important differences in momentum and variance would destabilize training.

## 10.3 The Solution: Block-wise Dynamic Quantization

The work by Dettmers et al. (2022) introduced a robust method that successfully overcomes these challenges. It is built on two core mathematical principles: **Block-wise Quantization** and **Dynamic (Quantile) Quantization**.

### 10.3.1. Block-wise Quantization

Instead of quantizing the entire tensor (e.g., the millions of values in the momentum state for a single layer) with a single scaling factor, the tensor is first partitioned into smaller, independent blocks or chunks. A typical block size might be 2048 elements.

For each block, we perform quantization independently. This allows the process to adapt to local variations in the magnitude of the state values. An outlier in one block will not affect the precision of quantization in another.

**Mathematical Formulation:**

Let $X$ be a tensor of optimizer states (e.g., $m_t$ or $v_t$) that we want to quantize.

1. Partition $X$ into blocks $X_1, X_2, \ldots, X_k$.
2. For each block $X_b$, find the absolute maximum value:
   $c_b = \max(|X_b|)$
   This value, $c_b$, is the **quantization constant** for block $b$. It is stored in a higher precision format (e.g., FP32).

3. **Quantize** the block $X_b$ to an 8-bit integer tensor $X_b^{INT8}$ using the formula:

$$X_b^{INT8} = \text{round}\left(\frac{X_b}{c_b} \times 127\right)$$

This scales the values in the block to the range [-127, 127].

4. To use the value in a computation, we **dequantize** it back to 32-bit precision:

$$X_b^{FP32} = \frac{X_b^{INT8}}{127} \times c_b$$

By storing one FP32 constant ($c_b$) for every 2048 INT8 values, the memory overhead is minimal, but the fidelity of the restored values is dramatically improved.

## 10.3.2. Dynamic (Quantile) Quantization

While block-wise processing solves the dynamic range problem, linear scaling still does not optimally handle the non-uniform distribution of values. To address this, a non-linear, dynamic quantization scheme is used. **Quantile quantization** is a powerful method for this.

Instead of mapping the range $\left[-c_b, c_b\right]$ linearly to [-127, 127], we build a "codebook" that is tailored to the specific distribution of the data.

**Process:**

1. For a given block of data $X_b$, determine its empirical cumulative distribution function (CDF).
2. Find the values that correspond to $2^k$ (e.g., 256 for 8-bit) evenly spaced quantiles. For 8-bit, we would find the values at the 0-th percentile, the (100/255)-th percentile, the (200/255)-th percentile, and so on, up to the 100th percentile.
3. These 256 values form a non-linear codebook (or lookup table) $Q$. The values in $Q$ are not evenly spaced; they are densely packed where the data is dense and sparse where the data is sparse.
4. **Quantization:** For each element $x \in X_b$, find the closest value in the codebook $Q$ and store its index (an 8-bit integer).
5. **Dequantization:** To restore the value, simply use the 8-bit integer index to look up the corresponding FP32 value from the codebook $Q$.

This ensures that our 256 available "precision levels" are spent wisely, capturing the structure of the dense part of the distribution with high fidelity.

## 10.3.3 The Full 8-bit AdamW Algorithm in Practice

Here is a step-by-step breakdown of how these concepts are integrated into the training loop for a single parameter update.

Let $\theta_t$ be the model parameters (FP16/BF16), $m_{t-1}^{INT8}$ and $v_{t-1}^{INT8}$ be the 8-bit optimizer states from the previous step, and $C_{m,t-1}, C_{v,t-1}$ be their corresponding FP32 block-wise constants.

1. **Forward & Backward Pass:**
   - Compute the loss $L(\theta_t)$.
   - Compute the gradient $\nabla L(\theta_t)$. This gradient is typically in FP16 or BF16. For the optimizer update, it is cast to FP32 for precision.
2. **Dequantize Optimizer States:**
   - For each block in $m_{t-1}^{INT8}$ and $v_{t-1}^{INT8}$:
     $$m_{t-1}^{FP32} = \text{dequantize}(m_{t-1}^{INT8}, C_{m,t-1})$$
     $$v_{t-1}^{FP32} = \text{dequantize}(v_{t-1}^{INT8}, C_{v,t-1})$$

3. **Perform AdamW Update in Full Precision (FP32):**
   - Update the first moment:
   $$m_t^{FP32} = \beta_1 m_{t-1}^{FP32} + (1 - \beta_1)\nabla L(\theta_t)^{FP32}$$
   - Update the second moment:
   $$v_t^{FP32} = \beta_2 v_{t-1}^{FP32} + (1 - \beta_2)(\nabla L(\theta_t)^{FP32})^2$$

4. **Quantize New Optimizer States for Storage:**
   - For each block in the newly computed $m_t^{FP32}$ and $v_t^{FP32}$:
   $$(m_t^{INT8}, C_{m,t}) = \text{quantize}(m_t^{FP32})$$
   $$(v_t^{INT8}, C_{v,t}) = \text{quantize}(v_t^{FP32})$$
   - These INT8 tensors and FP32 constants are stored, overwriting the previous states.

5. **Update Model Parameters:**
   - The parameter update itself uses the high-precision (FP32) intermediate values calculated in Step 3, before quantization. This is critical for stability.
   - First, perform bias correction on the FP32 states. Let's call the corrected states $\hat{m}_t^{FP32}$ and $\hat{v}_t^{FP32}$.
   - Update the parameters (which are often in BF16/FP16):
   $$\theta_{t+1} = \theta_t - \eta \left( \frac{\hat{m}_t^{FP32}}{\sqrt{\hat{v}_t^{FP32}+\epsilon}} + \lambda\theta_t \right)$$

This cycle ensures that the **storage** is memory-efficient (8-bit), but the **computation** of the update retains the numerical precision of FP32, preventing the accumulation of quantization errors from destabilizing the training process. The small computational overhead of the quantize/dequantize steps is far outweighed by the 4x reduction in optimizer memory, which enables larger batch sizes or entirely larger models to be trained on the same hardware.

Excellent question. The connection between a low-level mathematical technique like an 8-bit optimizer and a high-stakes, physically-grounded domain like aviation is not immediately obvious, but it is incredibly significant. The value lies not in the optimizer directly affecting an aircraft's aerodynamics, but in its role as an **enabling technology** that makes the development and deployment of crucial AI systems in aviation more feasible, efficient, and effective.

Here is a detailed breakdown of how 8-bit optimizers are helpful in the aviation domain, framed for industry stakeholders like aerospace engineers, safety officers, and operations managers.

# 11 Executive Summary: The Bridge from Optimization Theory to Aviation Practice

In aviation, the adoption of new technology is governed by strict requirements for **safety, reliability, certification, and operational efficiency**. Advanced AI and Large Language Models (LLMs) offer transformative potential, from predictive maintenance to intelligent pilot assistance. However, these models are notoriously large and resource-intensive.

**The core problem:** The immense computational and memory requirements of training and fine-tuning state-of-the-art AI models create a barrier to their practical use in the aviation ecosystem.

**The 8-bit optimizer solution:** By drastically reducing the memory footprint of the AI development process (by up to 75%), 8-bit optimizers act as a catalyst. They do not change *what* the AI model does, but they fundamentally change *how* it is built and deployed. This makes advanced AI more accessible, affordable, and adaptable for aviation-specific needs, directly impacting safety, cost, and operational readiness.

# 12 Specific Applications and Benefits in the Aviation Domain

Let's explore concrete use cases where the benefits of 8-bit optimizers become tangible.

## 1. Onboard / Edge AI Systems: Reducing SWaP (Size, Weight, and Power)

Aircraft are environments with extreme constraints on Size, Weight, and Power (SWaP). Any onboard computational system must be small, light, and energy-efficient.

- **Application: Predictive Maintenance and Real-time Anomaly Detection**
  - **Challenge:** An AI model running on an avionics unit could analyze real-time sensor data from engines, control surfaces, and hydraulic systems to predict component failure *before* it occurs. However, the onboard hardware is not a data center.
  - **How 8-bit Optimizers Help:** They enable the **fine-tuning** of a powerful, pre-trained predictive model on a specific aircraft's data stream using the limited computational resources available *on-premise* or even on a specialized ground-support computer. A maintenance organization can take a general "engine health" model and rapidly adapt it for a specific engine variant (e.g., a CFM56-7B vs. a GEnx) without needing a massive cloud computing budget. The resulting smaller, specialized model is easier to deploy onto a certified, low-SWaP onboard computer.
- **Application: Intelligent Pilot and Crew Assistance**
  - **Challenge:** An LLM in the cockpit could interpret pilot commands, automatically cross-reference the flight manual for emergency procedures, or monitor crew speech patterns for signs of fatigue or cognitive overload. This requires a sophisticated model running locally for real-time response and data privacy.
  - **How 8-bit Optimizers Help:** Using techniques like **QLoRA** (Quantized Low-Rank Adaptation), which leverages 8-bit optimizers, a general-purpose LLM can be efficiently fine-tuned on aviation-specific data (manuals, checklists, terminology). This fine-tuning can be done on a single GPU. The resulting highly specialized, yet compact, model can then be deployed within the cockpit's hardware constraints, providing real-time, reliable assistance.

## 2. Ground-Based Operations: Increasing Adaptability and Reducing Costs

Airlines and MRO (Maintenance, Repair, and Overhaul) centers operate on tight margins. The ability to develop and maintain proprietary AI models in-house is a significant competitive advantage.

- **Application: Fuel Optimization and Flight Planning**
  - **Challenge:** Fuel is a primary operational cost. AI models can analyze vast datasets—historical flight paths, real-time weather, air traffic control directives, and aircraft performance data—to recommend the most fuel-efficient routes. These models must be constantly updated to remain accurate.
  - **How 8-bit Optimizers Help:** They dramatically lower the hardware barrier for retraining. An airline's operations team can take a base model and **retrain it daily or even hourly** on the latest data using their own on-premise servers. Without 8-bit optimizers, this level of frequent, agile retraining would require a massive and expensive GPU cluster, making it unaffordable for many operators. This leads to more accurate, timely, and cost-saving flight plans.
- **Application: Safety Data Analysis**
  - **Challenge:** The NTSB and NASA's Aviation Safety Reporting System (ASRS) contain millions of unstructured text reports from pilots and controllers. Manually analyzing these to find emerging safety trends is a monumental task. An LLM can be used to "read" and categorize these reports to identify latent risks.

- **How 8-bit Optimizers Help:** A safety department can fine-tune a powerful open-source LLM (like Llama) on the specific lexicon and context of aviation safety reports. The use of 8-bit optimizers makes this process feasible for a data science team with access to just one or two commercial-grade GPUs. This **democratizes access to cutting-edge NLP**, allowing organizations to proactively identify safety trends instead of reacting to incidents.

## 3. Training and Simulation: Enhancing Realism and Development Speed

- **Application: Dynamic and Adaptive Flight Simulators**
  - **Challenge:** Creating realistic and unpredictable scenarios for pilot training is crucial. An AI "Director" could control virtual air traffic, inject unexpected (but plausible) system failures, or adapt the scenario based on the trainee's performance. Developing this AI is computationally intensive.
  - **How 8-bit Optimizers Help:** The development and iteration cycle for this AI Director model is significantly accelerated. Game developers already use similar techniques. In aviation simulation, this means a new training module or a more intelligent virtual instructor can be developed and refined much faster and on less hardware, reducing the cost and time-to-deployment for new training programs.

## Summary of Key Benefits for Aviation

| Benefit | Technical Enabler (8-bit Optimizer) | Impact on Aviation |
|---|---|---|
| **Increased Accessibility** | Reduces memory/hardware requirements for AI model training and fine-tuning. | Allows airlines, MROs, and safety boards to develop and customize their own AI solutions in-house without massive capital expenditure. |
| **Reduced SWaP** | Enables the creation of smaller, more efficient models suitable for deployment on constrained hardware. | Makes powerful real-time AI feasible for onboard avionics, improving predictive maintenance and pilot assistance. |
| **Faster Iteration & Adaptation** | Speeds up the retraining cycle by allowing larger batch sizes or running on existing hardware. | Models for fuel optimization or supply chain management can be updated constantly with fresh data, improving accuracy and saving costs. |
| **Enhanced Feasibility for Certification** | While not a direct solution, smaller, more contained models are often easier to subject to the rigorous Verification & Validation (V&V) required for certification. The development process becomes more manageable. | Accelerates the path from AI research to certified, flight-worthy application. |

In conclusion, 8-bit optimizers are a foundational technology. They are the "lightweight composite material" of the AI software world—they don't change the laws of physics, but by making the components lighter and more efficient, they enable the construction of far more advanced and capable systems that can operate effectively within the demanding constraints of the aviation industry.

# 13. Conclusion

ADAM has become the de facto standard optimizer for training Large Language Models due to its robust performance, adaptive learning rates, and relatively low hyperparameter sensitivity. While newer techniques continue to emerge, ADAM's combination of momentum and per-parameter adaptive learning rates has proven particularly effective for navigating the complex loss landscapes of modern LLMs.

The success of ADAM in LLM training demonstrates how critical optimization algorithm choice is for pushing the boundaries of what's possible in deep learning. As models continue to grow in size and complexity, further innovations in optimization will likely play a key role in enabling the next generation of language models.

The journey of optimization for LLMs is a story of continuous refinement and adaptation. We have moved from the simple, noisy updates of SGD to the robust, adaptive, and well-regularized updates of AdamW, which currently reigns as the dominant algorithm for training state-of-the-art models.

However, the field is far from static. The future of LLM optimization lies in tackling the tripartite challenges of **memory efficiency**, **computational speed**, and **generalization**. Innovations like 8-bit optimizers are making training more accessible, advanced methods like Shampoo are re-introducing the power of curvature information in a tractable way, and novel procedures like SAM are redefining the optimization objective itself to prioritize generalization. The co-evolution of these optimization algorithms with model architectures and hardware will continue to be the primary engine driving progress in the remarkable field of generative AI.

Of course. Let's transition from the survey-level overview to a more focused and detailed technical exposition on 8-bit optimizers. This section is structured as a deeper follow-up, suitable for a technical appendix or a specialized chapter in our ongoing research documentation.

# 14. References

- Anil, R., Gupta, V., Koren, T., & Singer, Y. (2020). *Second Order Optimization for Deep Learning*.
- Brown, T. B., Mann, B., Ryder, N., et al. (2020). *Language Models are Few-Shot Learners*. In Advances in Neural Information Processing Systems (NeurIPS).
- Chowdhery, A., Narang, S., Devlin, J., et al. (2022). *PaLM: Scaling Language Modeling with Pathways*. Journal of Machine Learning Research.
- Dettmers, T., Lewis, M., Belkada, Y., & Zettlemoyer, L. (2022). *LLM.int8(): 8-bit Matrix Multiplication for Transformers at Scale*.
- Duchi, J., Hazan, E., & Singer, Y. (2011). *Adaptive Subgradient Methods for Online Learning and Stochastic Optimization*. Journal of Machine Learning Research.
- Foret, P., Kleiner, A., Mobahi, H., & Neyshabur, B. (2021). *Sharpness-Aware Minimization for Efficiently Improving Generalization*. In International Conference on Learning Representations (ICLR).
- Gupta, V., Anil, R., et al. (2018). *Shampoo: Preconditioned Stochastic Tensor Optimization*. In International Conference on Machine Learning (ICML).
- Kingma, D. P., & Ba, J. (2015). *Adam: A Method for Stochastic Optimization*. In International Conference on Learning Representations (ICLR).
- Loshchilov, I., & Hutter, F. (2019). *Decoupled Weight Decay Regularization*. In International Conference on Learning Representations (ICLR).

- Polyak, B. T. (1964). *Some methods of speeding up the convergence of iteration methods*. USSR Computational Mathematics and Mathematical Physics.
- Touvron, H., Lavril, T., Izacard, G., et al. (2023). *LLaMA: Open and Efficient Foundation Language Models*. arXiv preprint arXiv:2302.13971.