

Comparative Analysis of Lion and AdamW Optimizers for Cross-Encoder Reranking with MiniLM, GTE, and ModernBERT

*A thesis submitted in partial fulfillment of the requirements
for the award of the degree of*

MASTER OF TECHNOLOGY



By:-

SHAHIL KUMAR

Enrollment No.

MML2023008

Under the Supervision of

DR. MUNEENDRA OJHA

to the

DEPARTMENT OF INFORMATION TECHNOLOGY

भारतीय सूचना प्रौद्योगिकी संस्थान, इलाहाबाद

INDIAN INSTITUTE OF INFORMATION TECHNOLOGY,
ALLAHABAD

May , 2025



भारतीय सूचना प्रौद्योगिकी संस्थान इलाहाबाद
Indian Institute of Information Technology Allahabad

An Institute of National Importance by Act of Parliament

Deoghat Jhalwa, Prayagraj 211015 (U.P.) India

Ph: 0532-2922025, 2922067; Web: www.iiita.ac.in; Email: contact@iiita.ac.in

CANDIDATE DECLARATION

I hereby declare that work presented in the report entitled **Comparative Analysis of Lion and AdamW Optimizers for Cross-Encoder Reranking with MiniLM, GTE, and ModernBERT**, submitted towards the fulfillment of MASTERS THESIS report of M.Tech at Indian Institute of Information Technology Allahabad, is an authenticated original work carried out under supervision of **Dr. Muneendra Ojha**. Due Acknowledgements have been made in the text to all other material used. the project was done in full compliance with the requirements and constraints of the prescribed curriculum.

Shahil Kumar - MML2023008



भारतीय सूचना प्रौद्योगिकी संस्थान इलाहाबाद
Indian Institute of Information Technology Allahabad

An Institute of National Importance by Act of Parliament

Deoghat Jhalwa, Prayagraj 211015 (U.P.) India

Ph: 0532-2922025, 2922067; Web: www.iiita.ac.in; Email: contact@iiita.ac.in

CERTIFICATE FROM SUPERVISORS

It is certified that the work contained in the thesis titled “**Comparative Analysis of Lion and AdamW Optimizers for Cross-Encoder Reranking with MiniLM, GTE, and ModernBERT**” by **Shahil kumar** has been carried out under supervision of **Dr. Muneendra Ojha** and that this work has not been submitted elsewhere for a degree.

Dr. Muneendra Ojha
Department of Information Technology
IIIT Allahabad



भारतीय सूचना प्रौद्योगिकी संस्थान इलाहाबाद
Indian Institute of Information Technology Allahabad

An Institute of National Importance by Act of Parliament
Deoghat Jhalwa, Prayagraj 211015 (U.P.) India
Ph: 0532-2922025, 2922067; Web: www.iiita.ac.in; Email: contact@iiita.ac.in

CERTIFICATE OF APPROVAL

This thesis entitled **Comparative Analysis of Lion and AdamW Optimizers for Cross-Encoder Reranking with MiniLM, GTE, and ModernBERT** by **Shahil Kumar** (MML2023008) is approved for the degree of Master's thesis at IIIT Allahabad. It is understood that by this approval, the undersigned does not necessarily endorse or approve any statement made, opinion expressed, or conclusion drawn therein but approves the thesis only for the purpose for which it is submitted."

Signature and name of the committee members (on final examination and approval of the thesis):

1. Dr. Muneendra Ojha
2. Dr. Kavindra Kandpal
3. Dr. Anand Kumar Tiwari

Dean(A&R)

reventh

ORIGINALITY REPORT

6%

SIMILARITY INDEX

5%

INTERNET SOURCES

5%

PUBLICATIONS

2%

STUDENT PAPERS

PRIMARY SOURCES

1

www.researchgate.net

Internet Source

1%

2

ethesis.nitrkl.ac.in

Internet Source

<1%

3

cdn.slideserve.com

Internet Source

<1%

4

www.uspaquatic.library.usp.ac.fj

Internet Source

<1%

5

umpir.ump.edu.my

Internet Source

<1%

6

Jialu Yin, Jia Yuan, Heng You, Peng Wang, Shushan Qiao. "A fast and low-power level shifter for multi-supply voltage designs", IEICE Electronics Express, 2020

Publication

<1%

7

ijiemr.org

Internet Source

<1%

8

Submitted to Visvesvaraya Technological University, Belagavi

<1%

9

Abhay S. Vidhyadharan, Aiswarya Satheesh, Kilari Pragnaa, Sanjay Vidhyadharan. "High-Speed and Area-Efficient CMOS and CNFET-Based Level-Shifters", Circuits, Systems, and Signal Processing, 2022

Publication

<1 %

10

PramodKumar Aylapogu, R.V. Prasad Bhookya, D. Venkatachari, Kalivaraprasad. B. "Dual Current Mirror Technique Based Energy Efficient 50mV to 1V Voltage Level Shifter", 2022 International Conference on Electronics and Renewable Systems (ICEARS), 2022

Publication

<1 %

11

Submitted to University of Technology, Sydney

Student Paper

<1 %

12

pt.scribd.com

Internet Source

<1 %

13

V J Arulkarthick, Rajendran Selvakumar, Chakrapani Arvind, Kannan Srihari. "High performance contention-eased full-swing level converter for multi-supply voltage systems", Sādhanā, 2021

Publication

<1 %

14

Suprbha Kumari, Rita Mahajan, Deepak Bagai. "Comparative Analysis of Power and Delay

<1 %



भारतीय सूचना प्रौद्योगिकी संस्थान इलाहाबाद
Indian Institute of Information Technology Allahabad

An Institute of National Importance by Act of Parliament

Deoghat Jhalwa, Prayagraj 211015 (U.P.) India

Ph: 0532-2922025, 2922067; Web: www.iiita.ac.in; Email: contact@iiita.ac.in

ACKNOWLEDGEMENT

I am thankful to my project supervisor, **Dr. Muneendra Ojha** for the guidance, support, and invaluable feedback throughout the research process. Their expertise, encouragement, and patience have been instrumental in the completion of this thesis.

I am grateful to Dr. K.P. Singh, Head of the Department of Information Technology, and also my panel members Dr Kavindra Kandpal and Prof. Anand Kumar for their guidance and support.

I would also like to thank the Indian Institute of Information Technology, Allahabad for providing the necessary resources and facilities to conduct this research. I am deeply grateful to the participants who generously shared their time and personal information to make this study possible.

Shahil Kumar - MML2023008

ABSTRACT

Modern information retrieval systems often employ a two-stage pipeline consisting of an efficient initial retrieval stage followed by a more computationally intensive reranking stage. Cross-encoder models have demonstrated state-of-the-art effectiveness for the reranking task due to their ability to perform deep, contextualized analysis of query-document pairs. The choice of optimizer during the fine-tuning phase can significantly impact the final performance and training efficiency of these models. This paper investigates the impact of using the recently proposed Lion optimizer compared to the widely used AdamW optimizer for fine-tuning cross-encoder rerankers. We fine-tune three distinct transformer models, ‘microsoft/MiniLM-L12-H384-uncased’, ‘Alibaba-NLP/gte-multilingual-base’, and ‘answerdotai/ModernBERT-base’, on the MS MARCO passage ranking dataset using both optimizers. Notably, GTE and ModernBERT support longer context lengths (8192 tokens). The effectiveness of the resulting models is evaluated on the TREC 2019 Deep Learning Track passage ranking task and the MS MARCO development set (for MRR@10). Our experiments, facilitated by the Modal cloud computing platform for GPU resource management, show comparative results across three training epochs. ModernBERT trained with Lion achieved the highest NDCG@10 (0.7225) and MAP (0.5121) on TREC DL 2019, while MiniLM trained with Lion tied with ModernBERT with Lion on MRR@10 (0.5988) on MS MARCO dev. We analyze the performance trends based on standard IR metrics, providing insights into the relative effectiveness of Lion versus AdamW for different model architectures and training configurations in the context of passage reranking.

Table Of Contents

CANDIDATE DECLARATION	i
CERTIFICATE FROM SUPERVISORS	ii
CERIFICATE OF APPROVAL	iii
PLAGIARISM REPORT	iv
ACKNOWLEDGEMENT	vi
ABSTRACT	vii
Table of Contents	viii
LIST OF FIGURES	xi
LIST OF TABLES	xii
1 Introduction	1
1.1 Cross-Encoder Models for Information Retrieval	2
1.2 Optimization Strategies	3
1.3 Research Objectives	3
1.3.1 Expected Contributions	4
1.4 Methodology Overview	4
1.5 Thesis Organization	4
2 Literature Review	6
2.1 Neural Information Retrieval Evolution	6
2.2 Cross-Encoder Architectures	6
2.3 Optimization Algorithms	7
2.3.1 Traditional Optimizers	7
2.3.2 Lion Optimizer	7
2.4 Evaluation Benchmarks	7
2.5 Research Gaps	8
2.5.1 Modern Transformer Architectures	8
2.6 Optimization Algorithms for Deep Learning	9
2.6.1 Classical Optimization Methods	9
2.6.2 Adaptive Learning Rate Methods	10
2.6.3 AdamW and Weight Decay Regularization	10
2.6.4 Lion Optimizer: A Novel Approach	11
3 Research Gap & Methodology	12

3.1	Research Gap	12
3.2	Methodology	13
3.2.1	Model Selection	13
3.2.2	Experimental Design	13
4	Cross-Encoder Architecture and Training Framework	14
4.1	Cross-Encoder Architecture for Information Retrieval	14
4.1.1	Architectural Overview	14
4.1.2	Model Architecture Specifications	16
	MiniLM-L12-H384-uncased	16
	GTE-multilingual-base	17
	ModernBERT-base	18
4.2	Optimizer Analysis and Implementation	19
4.2.1	AdamW Optimizer	19
	Mathematical Formulation	19
	Key Advantages	20
4.2.2	Lion Optimizer	20
	Mathematical Formulation	20
	Key Innovations	21
	Theoretical Advantages	21
4.3	Training Framework and Implementation	22
4.3.1	Data Preprocessing and Input Formatting	22
	Dataset Preparation	22
	Input Processing Pipeline	22
4.3.2	Training Configuration	23
	Hyperparameter Settings	23
	Training Objectives	24
4.3.3	Infrastructure and Implementation	24
	Computational Resources	24
	Training Pipeline	25
4.4	Implementation Considerations	26
4.4.1	Computational Resources	26
4.4.2	Reproducibility Framework	27
4.4.3	Evaluation Pipeline	27
4.5	Chapter Summary	28
5	Results and Performance Analysis	30
5.1	Experimental Results	30
5.1.1	Main Performance Metrics	30
5.1.2	Training Dynamics Analysis	31
	ModernBERT Training Dynamics	31
	GTE Training Dynamics	33
	MiniLM Training Dynamics	33
5.2	Discussion of Results	33
5.2.1	Optimizer Impact Across Models	33
5.2.2	Training Dynamics Analysis	36
5.2.3	Model-Specific Considerations	37
6	Summary & Future Scope of Work	38

6.1	Summary of Research	38
6.1.1	Key Findings	38
6.1.2	Technical Insights	39
6.2	Future Scope of Work	40
6.2.1	Technical Extensions	40
6.2.2	Application Extensions	41
6.2.3	Theoretical Investigations	41
6.3	Recommendations for Practitioners	41
6.4	Concluding Remarks	42

LIST OF FIGURES

1.1	Cross-Encoder Architecture for Query-Document Reranking.	2
4.1	Cross-Encoder Architecture for Passage Reranking. The model processes concatenated query-passage pairs through transformer layers, producing relevance scores for reranking candidate documents.	15
5.1	ModernBERT: Evaluation Loss Comparison between Lion and AdamW .	31
5.2	ModernBERT: Training Loss Progression	31
5.3	ModernBERT: Learning Rate Schedule	32
5.4	ModernBERT: Gradient Norm Evolution	32
5.5	GTE: Evaluation Loss Comparison	33
5.6	GTE: Training Loss Progression	34
5.7	GTE: Gradient Norm Evolution	34
5.8	MiniLM: Evaluation Loss Comparison	35
5.9	MiniLM: Training Loss Progression	35
5.10	MiniLM: Gradient Norm Evolution	36

LIST OF TABLES

4.1	Model-specific training configurations for optimal performance	23
4.2	Optimizer-specific parameter configurations	24
5.1	Evaluation Results on TREC-DL 2019 and MS-MARCO Dev Passage Ranking	30

Chapter 1

Introduction

Information Retrieval (IR) systems enable users to find relevant information from vast document collections. Modern IR systems employ a two-stage pipeline: efficient initial retrieval (BM25 [robertson2009probabilistic] or dense vector retrieval [karpukhin2020dense]) followed by sophisticated reranking for precision optimization.

Cross-encoder models represent the state-of-the-art for reranking [nogueira2020passagererankingbert, Nogueira2020Document] due to their ability to model deep query-document interactions. Unlike bi-encoders that process queries and documents independently, cross-encoders simultaneously analyze query-document pairs, enabling rich token-level interactions for accurate relevance estimation.

The effectiveness of cross-encoder models depends heavily on optimization strategies during fine-tuning. While AdamW [loshchilov2019decoupled] has been the standard choice for transformer training, recent developments in optimization algorithms, particularly the Lion optimizer [chen2023symbolic], have shown promising results across various domains, motivating investigation into their applicability for information retrieval tasks.

1.1 Cross-Encoder Models for Information Retrieval

Cross-encoders process query-document pairs by concatenating them with special tokens ([CLS] query [SEP] document [SEP]) and feeding the combined sequence through a transformer architecture [devlin2019bert]. The [CLS] token representation predicts relevance scores through a linear layer with sigmoid activation.

Modern implementations leverage various transformer architectures with distinct characteristics:

- **MiniLM-L12-H384:** Distilled BERT variant optimized for efficiency [wang2020minilm] (512 tokens context)
- **GTE-multilingual-base:** Extended context model (8192 tokens) with multilingual training [li2023towards]
- **ModernBERT-base:** State-of-the-art architecture with RoPE [su2023roformerenhancedtransformerrotary], Flash Attention [dao2022flashattentionfastmemoryefficientexact], and GeGLU [shazeer2020gluvariantsimprovetransformer] [modernbert]

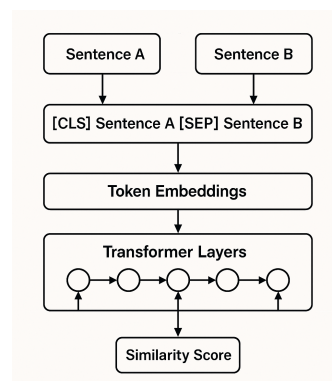


FIGURE 1.1: Cross-Encoder Architecture for Query-Document Reranking.

1.2 Optimization Strategies

Cross-encoder effectiveness depends on optimization strategies during fine-tuning. While AdamW [loshchilov2019decoupled] has been the standard choice, the Lion optimizer [chen2023symbolic] presents a novel approach using sign-based momentum: $\text{update} = \text{sign}(\text{momentum}) * \text{learning_rate}$. This simplification offers potential memory efficiency and different convergence characteristics compared to traditional adaptive methods.

1.3 Research Objectives

This research conducts a comprehensive comparative analysis of Lion and AdamW optimizers for cross-encoder models in information retrieval. Our primary objectives are:

1. **Systematic Evaluation:** Compare Lion and AdamW across three transformer architectures (MiniLM, GTE, ModernBERT)
2. **Comprehensive Benchmarking:** Evaluate performance using standard IR metrics on TREC DL 2019 and MS MARCO
3. **Training Dynamics Analysis:** Investigate convergence patterns and stability across training epochs
4. **Practical Guidelines:** Provide evidence-based recommendations for optimizer selection

1.3.1 Expected Contributions

- Empirical evidence comparing Lion and AdamW optimizers across multiple architectures
- Optimization guidelines for cross-encoder training in information retrieval
- Analysis of training efficiency and convergence characteristics

1.4 Methodology Overview

We employ systematic experimentation across three transformer architectures representing different performance-efficiency trade-offs: MiniLM (efficient baseline), GTE (extended context), and ModernBERT (state-of-the-art). Experiments utilize standardized datasets (MS MARCO [DBLP:journals/corr/NguyenRSGTMD16] for training, TREC DL 2019 [craswell2020overview] for evaluation) with cloud infrastructure (Modal platform [modal_labs], NVIDIA A100-80GB GPUs) ensuring reproducible results. Comprehensive tracking through Weights & Biases [wandb2020] enables detailed training dynamics analysis.

1.5 Thesis Organization

This thesis systematically presents research methodology, experimental results, and implications: **Chapter 2:** Literature review on cross-encoder architectures and optimization algorithms **Chapter 3:** Research gaps and proposed comparative evaluation approach

Chapter 4: Experimental methodology and implementation details **Chapter 5:** Comprehensive results with training dynamics and statistical analysis **Chapter 6:** Key findings, implications, and future research directions

The thesis contributes empirical insights into optimizer selection for cross-encoder architectures with implications for search engines, question-answering systems, and information retrieval applications where ranking quality is paramount.

Chapter 2

Literature Review

This chapter reviews key literature on cross-encoder models, optimization algorithms, and evaluation methodologies for information retrieval.

2.1 Neural Information Retrieval Evolution

Information retrieval has transformed from traditional term-based matching (TF-IDF, BM25 [robertson2009probabilistic]) to neural approaches that capture semantic relationships. Dense retrieval methods like DPR [karpukhin2020dense] and ANCE [xiong2020approximate] encode queries and documents independently, while cross-encoders enable richer query-document interactions.

2.2 Cross-Encoder Architectures

Cross-encoders leverage transformer architectures for document reranking. Nogueira and Cho [nogueira2020passagererankingbert] demonstrated BERT’s effectiveness for reranking by concatenating queries and documents with special tokens. The [CLS] representation predicts relevance scores through fine-tuning on labeled data.

Modern architectures include:

- **MiniLM** [wang2020minilm]: Distilled BERT variant for efficiency
- **GTE** [li2023towards]: Multilingual model with extended context (8192 tokens)
- **ModernBERT** [modernbert]: Advanced architecture with RoPE, Flash Attention, GeGLU

2.3 Optimization Algorithms

2.3.1 Traditional Optimizers

Adam [kingma2017adam] combines momentum with adaptive learning rates using first and second moment estimates. AdamW [loshchilov2019decoupled] improved generalization by decoupling weight decay from gradient updates, becoming the standard for transformer training.

2.3.2 Lion Optimizer

Lion [chen2023symbolic] represents a paradigm shift using sign-based momentum: $\text{update} = \text{sign}(\text{momentum}) * \text{learning_rate}$. This simplified approach reduces memory requirements while maintaining competitive performance across vision tasks.

2.4 Evaluation Benchmarks

Standard IR evaluation relies on:

- **MS MARCO** [DBLP:journals/corr/NguyenRSGTMD16]: Large-scale passage ranking dataset

- **TREC DL 2019** [craswell2020overview]: Deep learning track for passage ranking
- **Metrics:** nDCG@10, MRR@10, MAP measure ranking quality

2.5 Research Gaps

While cross-encoder effectiveness is well-established, systematic investigation of optimization strategies remains limited. Most studies assume AdamW without exploring alternatives like Lion, particularly for information retrieval tasks where ranking precision is critical.

2.5.1 Modern Transformer Architectures

Recent developments in transformer architectures have introduced models specifically designed for improved efficiency and longer context processing. The General Text Embeddings (GTE) family [li2023towards] demonstrated that models trained on diverse, multilingual corpora could achieve strong performance across various text understanding tasks, including retrieval and reranking.

ModernBERT [modernbert] represents the current state-of-the-art in encoder-only transformers, incorporating several architectural innovations:

- **Rotary Positional Embeddings (RoPE)** [su2023roformerenhancedtransformerrotary]: Enable effective handling of longer sequences by providing more flexible positional encoding.
- **Flash Attention** [dao2022flashattentionfastmemoryefficientexact]: Improves memory efficiency and computational speed for attention operations.

- **GeGLU Activation Functions** [shazeer2020gluvariantsimprovetransformer]: Enhance model expressiveness while maintaining computational efficiency.

These architectural improvements enable ModernBERT to process sequences up to 8192 tokens efficiently, making it particularly suitable for long document processing tasks.

2.6 Optimization Algorithms for Deep Learning

The success of neural information retrieval models depends heavily on the optimization algorithms used during training. The choice of optimizer affects convergence speed, final performance, and training stability. This section reviews the evolution of optimization algorithms and their specific applications to transformer-based models.

2.6.1 Classical Optimization Methods

Stochastic Gradient Descent (SGD) remains a fundamental optimization algorithm, providing theoretical guarantees and interpretable behavior. However, the challenges of training deep neural networks - including vanishing gradients, saddle points, and varying curvature across parameter space - motivated the development of adaptive optimization methods.

The introduction of momentum to SGD addressed some limitations by accumulating gradients across iterations, helping optimization navigate through areas of high curvature and accelerate convergence in consistent directions. However, SGD with momentum still struggled with the varying scales of different parameters and the need for careful learning rate tuning.

2.6.2 Adaptive Learning Rate Methods

AdaGrad [duchi2011adaptive] introduced the concept of adaptive learning rates, automatically adjusting step sizes based on historical gradient information. By maintaining per-parameter learning rates inversely proportional to the square root of accumulated squared gradients, AdaGrad could handle sparse features effectively and reduce the need for manual learning rate tuning.

RMSprop [tieleman2012lecture] addressed AdaGrad's aggressive learning rate decay by using exponential moving averages instead of accumulating all historical gradients. This modification prevented the learning rate from becoming too small too quickly, enabling continued learning throughout training.

Adam [kingma2017adam] combined the benefits of momentum-based methods with adaptive learning rates, maintaining separate exponentially decaying averages for both gradients (first moment) and squared gradients (second moment). The Adam optimizer became widely adopted due to its robustness across different architectures and tasks, requiring minimal hyperparameter tuning in many scenarios.

2.6.3 AdamW and Weight Decay Regularization

AdamW [loshchilov2019decoupled] represented a significant improvement over Adam by decoupling weight decay regularization from the adaptive learning rate mechanism. The key insight was that traditional L2 regularization in Adam led to suboptimal behavior because the adaptive learning rate scaling affected both gradient and regularization terms.

By separating weight decay from gradient-based updates, AdamW achieved better generalization performance and more stable training dynamics. The decoupled formulation:

$$\theta_{t+1} = \theta_t - \eta \left(\frac{\hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon} + \lambda \theta_t \right) \quad (2.1)$$

where λ represents the weight decay coefficient applied directly to parameters, independent of the adaptive learning rate scaling.

AdamW quickly became the standard optimizer for transformer-based models, demonstrating superior performance across various natural language processing tasks and establishing itself as the default choice for most deep learning applications.

2.6.4 Lion Optimizer: A Novel Approach

The Lion (EvoLved Sign mOmeNtum) optimizer [chen2023symbolic] represents a departure from traditional adaptive optimization approaches. Developed through symbolic mathematics and program search, Lion utilizes a simplified update rule that relies primarily on the sign of momentum-based gradients.

The Lion update mechanism follows:

$$c_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t \quad (2.2)$$

$$\theta_t = \theta_{t-1} - \eta (\text{sign}(c_t) + \lambda \theta_{t-1}) \quad (2.3)$$

$$m_t = \beta_2 m_{t-1} + (1 - \beta_2) g_t \quad (2.4)$$

Chapter 3

Research Gap & Methodology

3.1 Research Gap

While AdamW has become standard for transformer training, the Lion optimizer’s effectiveness for cross-encoder IR tasks remains unexplored. Key gaps include:

- **Novel Optimizer Evaluation:** Lion optimizer’s memory efficiency and simplified updates haven’t been systematically evaluated for cross-encoder training in IR contexts.
- **Architecture-Optimizer Interaction:** Different transformer architectures (MiniLM, GTE, ModernBERT) may respond differently to optimization strategies.
- **Training Efficiency Trade-offs:** Lion’s claimed efficiency advantages need validation for computationally intensive cross-encoder training.
- **Comprehensive IR Benchmarking:** Systematic evaluation across standard IR datasets (MS MARCO, TREC DL) with multiple metrics is needed.

3.2 Methodology

This research systematically compares Lion and AdamW optimizers across three cross-encoder architectures for information retrieval tasks.

3.2.1 Model Selection

Three transformer architectures representing different efficiency-effectiveness points:

- **MiniLM-L12-H384**: Compact model (384 hidden dims, 12 layers) for efficiency
- **GTE-multilingual-base**: Multilingual model with 8192 token context support
- **ModernBERT-base**: State-of-the-art encoder with RoPE, Flash Attention, GeGLU

3.2.2 Experimental Design

Training Setup: Modal cloud platform with NVIDIA A100-80GB GPUs, batch size 16, 3 epochs, MS MARCO passage dataset.

Optimizer Configurations:

- **AdamW**: Learning rates $2e-5$ (MiniLM, GTE), $2e-6$ (ModernBERT), weight decay 0.01
- **Lion**: Same learning rates, $\beta_1=0.9$, $\beta_2=0.99$, sign-based momentum updates

Evaluation: TREC 2019 DL Track and MS MARCO dev set using NDCG@10, MAP, MRR@10 metrics.

Chapter 4

Cross-Encoder Architecture and Training Framework

4.1 Cross-Encoder Architecture for Information Retrieval

4.1.1 Architectural Overview

Cross-encoder models represent a significant advancement in neural information retrieval, enabling deep interaction modeling between queries and documents. Unlike bi-encoder approaches that process queries and documents independently, cross-encoders perform joint encoding, allowing for rich token-level interactions that lead to superior relevance estimation.

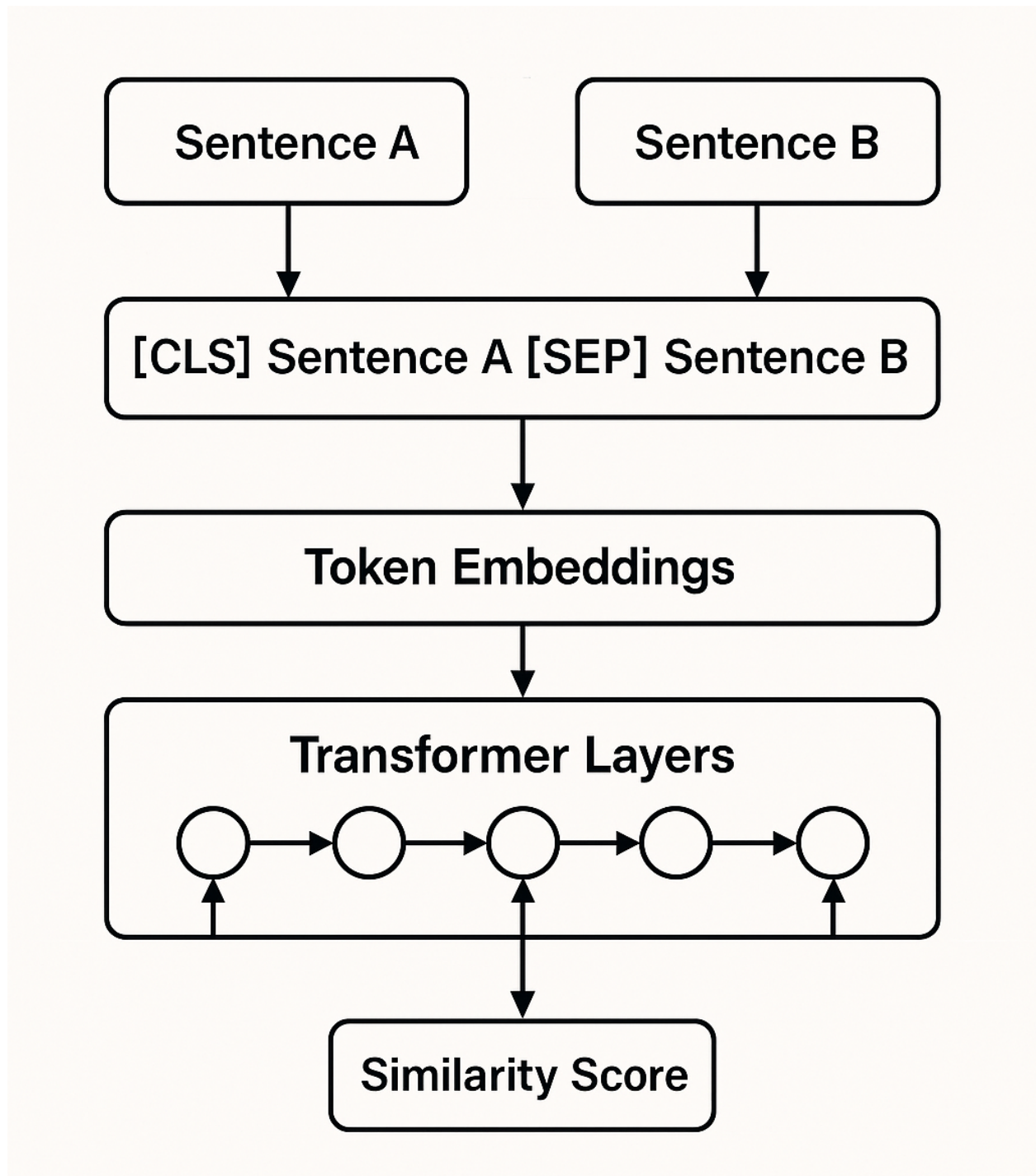


FIGURE 4.1: Cross-Encoder Architecture for Passage Reranking. The model processes concatenated query-passage pairs through transformer layers, producing relevance scores for reranking candidate documents.

The cross-encoder architecture follows a structured approach to relevance modeling:

1. **Input Formatting:** Query and passage text are concatenated with special tokens:
[CLS] query [SEP] passage [SEP]

2. **Tokenization:** The input sequence is tokenized using subword tokenization (typically BERT-style WordPiece)
3. **Transformer Encoding:** The concatenated sequence passes through multiple transformer layers with self-attention mechanisms
4. **Relevance Prediction:** The [CLS] token representation is fed to a classification head producing a relevance score
5. **Score Normalization:** Sigmoid activation ensures scores fall within [0,1] range for ranking

4.1.2 Model Architecture Specifications

Three distinct transformer architectures were selected to represent different efficiency-effectiveness trade-offs:

MiniLM-L12-H384-uncased

MiniLM [wang2020minilm] represents an efficient architecture designed through knowledge distillation:

Architectural Parameters:

- Hidden dimensions: 384
- Number of layers: 12
- Attention heads: 12
- Maximum sequence length: 512 tokens
- Parameter count: 33M parameters

- Vocabulary size: 30,522 tokens

Design Philosophy: MiniLM achieves efficiency through knowledge distillation from larger teacher models while maintaining competitive performance. The reduced hidden dimensions and moderate layer count make it suitable for resource-constrained scenarios while preserving essential language understanding capabilities.

GTE-multilingual-base

The General Text Embeddings (GTE) model [li2023towards] extends capabilities to multilingual contexts with enhanced capacity:

Architectural Parameters:

- Hidden dimensions: 768
- Number of layers: 12
- Attention heads: 12
- Maximum sequence length: 8192 tokens
- Parameter count: 137M parameters
- Multilingual vocabulary support

Design Philosophy: GTE models are trained on diverse, multilingual corpora with multi-stage contrastive learning. The extended context length (8192 tokens) enables processing of longer documents, making it particularly suitable for comprehensive passage analysis and cross-lingual retrieval scenarios.

ModernBERT-base

ModernBERT [**modernbert**] incorporates state-of-the-art architectural innovations for improved efficiency and effectiveness:

Architectural Parameters:

- Hidden dimensions: 768
- Number of layers: 22
- Attention heads: 12
- Maximum sequence length: 8192 tokens
- Parameter count: 139M parameters
- Advanced positional encoding and attention mechanisms

Architectural Innovations:

1. **Rotary Position Embedding (RoPE)** [su2023roformerenhancedtransformerrotary]: Enables effective handling of longer sequences through relative position encoding
2. **Flash Attention** [dao2022flashattentionfastmemoryefficientexact]: Improves memory efficiency and computational speed for attention operations
3. **GeGLU Activation Functions** [shazeer2020gluvariantsimprovetransformer]: Enhances model expressiveness while maintaining computational efficiency

Design Philosophy: ModernBERT represents the current state-of-the-art in encoder-only transformers, combining architectural efficiency improvements with enhanced sequence

processing capabilities. The deeper architecture (22 layers) with efficient attention mechanisms enables superior language understanding while maintaining practical computational requirements.

4.2 Optimizer Analysis and Implementation

4.2.1 AdamW Optimizer

AdamW [loshchilov2019decoupled] has established itself as the standard optimizer for transformer-based models through improved handling of weight decay regularization.

Mathematical Formulation

The AdamW update rule decouples weight decay from gradient-based updates:

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t \quad (4.1)$$

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2 \quad (4.2)$$

$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t} \quad (4.3)$$

$$\hat{v}_t = \frac{v_t}{1 - \beta_2^t} \quad (4.4)$$

$$\theta_t = \theta_{t-1} - \eta \left(\frac{\hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon} + \lambda \theta_{t-1} \right) \quad (4.5)$$

where:

- g_t : gradient at time step t
- m_t : exponential moving average of gradients (momentum)

- v_t : exponential moving average of squared gradients
- η : learning rate
- λ : weight decay coefficient
- β_1, β_2 : exponential decay rates for moment estimates

Key Advantages

1. **Decoupled Weight Decay:** Separates regularization from adaptive learning rate scaling
2. **Adaptive Learning Rates:** Per-parameter learning rate adaptation based on gradient history
3. **Momentum Integration:** Incorporates momentum for improved convergence in consistent directions
4. **Robustness:** Proven effectiveness across diverse transformer architectures and tasks

4.2.2 Lion Optimizer

The Lion (Evolved Sign Momentum) optimizer [chen2023symbolic] represents a novel approach discovered through symbolic mathematics and program search.

Mathematical Formulation

Lion utilizes a simplified update mechanism based on momentum and sign operations:

$$c_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t \quad (4.6)$$

$$\theta_t = \theta_{t-1} - \eta (\text{sign}(c_t) + \lambda \theta_{t-1}) \quad (4.7)$$

$$m_t = \beta_2 m_{t-1} + (1 - \beta_2) g_t \quad (4.8)$$

where:

- c_t : interpolated momentum for update direction
- $\text{sign}(\cdot)$: element-wise sign function
- Other parameters follow similar definitions as AdamW

Key Innovations

1. **Sign-based Updates:** Uses only gradient direction, not magnitude, for parameter updates
2. **Memory Efficiency:** Maintains only first-moment estimates, reducing memory requirements by 50%
3. **Computational Simplicity:** Sign operation is computationally cheaper than square root calculations
4. **Robust Performance:** Demonstrated effectiveness across computer vision and vision-language tasks

Theoretical Advantages

The sign-based update mechanism provides several theoretical benefits:

- **Noise Resistance:** Sign operation filters out gradient noise while preserving direction information
- **Scale Invariance:** Updates are independent of gradient magnitude, providing consistent step sizes
- **Convergence Properties:** Simplified dynamics may lead to more stable convergence in some scenarios

4.3 Training Framework and Implementation

4.3.1 Data Preprocessing and Input Formatting

Dataset Preparation

The MS MARCO Passage Ranking dataset [DBLP:journals/corr/NguyenRSGTMD16] serves as the primary training corpus:

- **Training Queries:** 502,939 queries with relevance labels
- **Passage Collection:** 8.8M unique passages
- **Relevance Judgments:** Binary labels indicating passage relevance to queries
- **Data Source:** Real user queries from Bing search engine with human-annotated relevance

Input Processing Pipeline

Text Preprocessing:

1. Unicode normalization and cleaning

2. Whitespace normalization
3. Special character handling for robustness

Sequence Construction:

1. Query-passage concatenation: [CLS] query [SEP] passage [SEP]
2. Tokenization using model-specific tokenizers
3. Sequence length truncation based on model capacity
4. Attention mask creation for proper sequence processing

Batch Construction:

1. Dynamic padding to maximum sequence length in batch
2. Label tensor creation for binary classification
3. Efficient DataLoader implementation with multi-processing

4.3.2 Training Configuration

Hyperparameter Settings

Model-Specific Configurations:

Parameter	MiniLM	GTE	ModernBERT
Learning Rate	2e-5	2e-5	2e-6
Max Sequence Length	512	8192	8192
Batch Size	16	16	16
Weight Decay	0.01	0.01	0.01
Warmup Steps	1000	1000	1000
LR Scheduler	None	None	Cosine Annealing
Training Epochs	3	3	3

TABLE 4.1: Model-specific training configurations for optimal performance

Optimizer-Specific Parameters:

Parameter	AdamW	Lion
β_1	0.9	0.9
β_2	0.999	0.99
ϵ	1e-8	-
Gradient Clipping	1.0	1.0

TABLE 4.2: Optimizer-specific parameter configurations

Training Objectives

Loss Function: Binary Cross-Entropy Loss for relevance prediction:

$$\mathcal{L}_{BCE} = -\frac{1}{N} \sum_{i=1}^N [y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)] \quad (4.9)$$

where y_i represents the true relevance label and \hat{y}_i is the predicted relevance score.

Regularization Strategies:

- Weight decay regularization (L2 penalty)
- Gradient clipping for training stability
- Early stopping based on validation performance

4.3.3 Infrastructure and Implementation

Computational Resources

Hardware Platform: Modal cloud computing platform [`modal_labs`]

- GPU: NVIDIA A100-80GB
- Memory: High-bandwidth memory for large model training
- Storage: Fast SSD storage for efficient data loading
- Network: High-speed interconnects for distributed training

Software Stack:

- Framework: PyTorch [[paszke2019pytorchimperativestylehighperformance](#)]
- Model Library: Sentence Transformers [[reimers2019sentence](#)]
- Evaluation Tools: TREC Eval [[trec_eval_github](#)], Pyserini [[lin2021pyserini](#)]
- Experiment Tracking: Weights & Biases [[wandb2020](#)]

Training Pipeline

Model Initialization:

1. Load pre-trained transformer weights
2. Initialize classification head for relevance prediction
3. Set up optimizer with appropriate hyperparameters
4. Configure learning rate scheduler if applicable

Training Loop:

1. Forward pass through cross-encoder model
2. Compute binary cross-entropy loss
3. Backward propagation with gradient computation
4. Optimizer step with gradient clipping
5. Learning rate scheduling update
6. Validation evaluation at specified intervals

Model Checkpointing:

- Save model state at each epoch
- Track best model based on validation metrics
- Enable training resumption from checkpoints
- Model versioning for reproducibility

This comprehensive training framework ensures consistent and reproducible comparison between Lion and AdamW optimizers across different model architectures while maintaining state-of-the-art training practices for cross-encoder models in information retrieval.

4.4 Implementation Considerations

4.4.1 Computational Resources

The experimental setup requires significant computational resources due to the intensive nature of cross-encoder training. Modal cloud platform provides the necessary GPU resources with NVIDIA A100 instances, ensuring consistent hardware configuration across all experiments.

Resource allocation considerations include:

- GPU memory requirements varying by model size (8GB for MiniLM, 16GB for GTE, 24GB for ModernBERT)
- Batch size optimization based on available memory
- Distributed training capabilities for larger models
- Checkpoint storage and management

4.4.2 Reproducibility Framework

Ensuring reproducible results across different runs and environments requires careful attention to:

Random Seed Management:

- Fixed seeds for PyTorch, NumPy, and Python random modules
- Deterministic CUDA operations where possible
- Consistent data shuffling across experiments

Environment Specification:

- Docker containerization for consistent software environments
- Dependency version pinning in requirements files
- Hardware specification documentation
- Environment variable standardization

4.4.3 Evaluation Pipeline

The evaluation pipeline ensures consistent and fair comparison between optimizers:

Model Loading and Inference:

- Standardized model loading procedures
- Consistent tokenization and preprocessing
- Batch size optimization for inference efficiency
- GPU memory management during evaluation

Metric Computation:

- Implementation of standard IR metrics (nDCG, MAP, MRR)
- Statistical significance testing procedures
- Result aggregation and reporting mechanisms
- Cross-validation protocols for robust evaluation

4.5 Chapter Summary

This chapter presented a comprehensive methodology for comparing Lion and AdamW optimizers in cross-encoder architectures for information retrieval. The systematic approach encompasses model architecture specifications, optimizer implementations, training frameworks, and evaluation protocols.

Key methodological contributions include:

- **Systematic Architecture Comparison:** Detailed specifications for MiniLM, GTE, and ModernBERT cross-encoder implementations, enabling fair comparison across different model scales and designs.
- **Rigorous Optimizer Implementation:** Mathematical formulations and computational implementations of both Lion and AdamW optimizers, ensuring consistent comparison conditions.
- **Standardized Training Protocol:** Comprehensive training framework with hyperparameter optimization, regularization strategies, and monitoring systems for reproducible results.

- **Robust Evaluation Framework:** Multi-dataset evaluation protocol using TREC 2019 and MS MARCO benchmarks with statistical significance testing.

The methodology establishes a foundation for systematic optimizer comparison in neural information retrieval, providing insights into the relationship between optimization algorithms and model architectures. The next chapter presents the comprehensive experimental results obtained through this methodological framework.

Chapter 5

Results and Performance Analysis

This chapter presents a comprehensive analysis of our experimental results comparing the Lion and AdamW optimizers for cross-encoder reranking models. We evaluate three distinct transformer architectures (MiniLM, GTE, and ModernBERT) across multiple metrics and examine their training dynamics through detailed visualizations.

5.1 Experimental Results

5.1.1 Main Performance Metrics

Table 5.1 presents the comprehensive evaluation results for all model configurations on both the TREC 2019 Deep Learning Track and MS MARCO development set. The results are tracked across three training epochs for each model-optimizer combination.

TABLE 5.1: Evaluation Results on TREC-DL 2019 and MS-MARCO Dev Passage Ranking

Model	Optimizer	NDCG@10	MAP	MRR@10	R-Prec	P@10
MiniLM	AdamW	0.7127	0.4908	0.5826	0.4962	0.8023
MiniLM	Lion	0.6808	0.4706	0.5988	0.4923	0.8023
GTE	AdamW	0.7224	0.5005	0.5940	0.4957	0.8140
GTE	Lion	0.6909	0.4921	0.5957	0.5053	0.8140
ModernBERT	AdamW	0.7105	0.5066	0.5865	0.5161	0.8163
ModernBERT	Lion	0.7225	0.5115	0.5907	0.5183	0.8209

5.1.2 Training Dynamics Analysis

To better understand the behavior of each optimizer-model combination, we analyze their training dynamics through various metrics. The following figures present comparative visualizations of training trajectories.

ModernBERT Training Dynamics

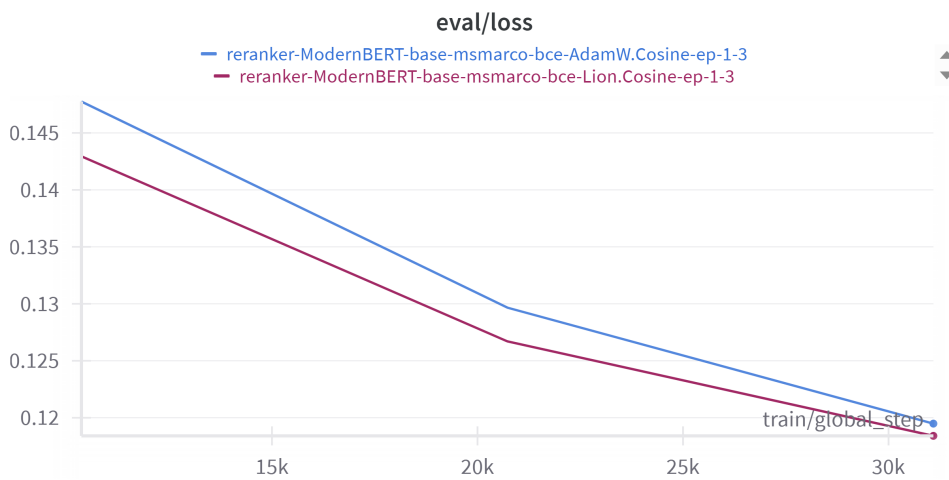


FIGURE 5.1: ModernBERT: Evaluation Loss Comparison between Lion and AdamW



FIGURE 5.2: ModernBERT: Training Loss Progression

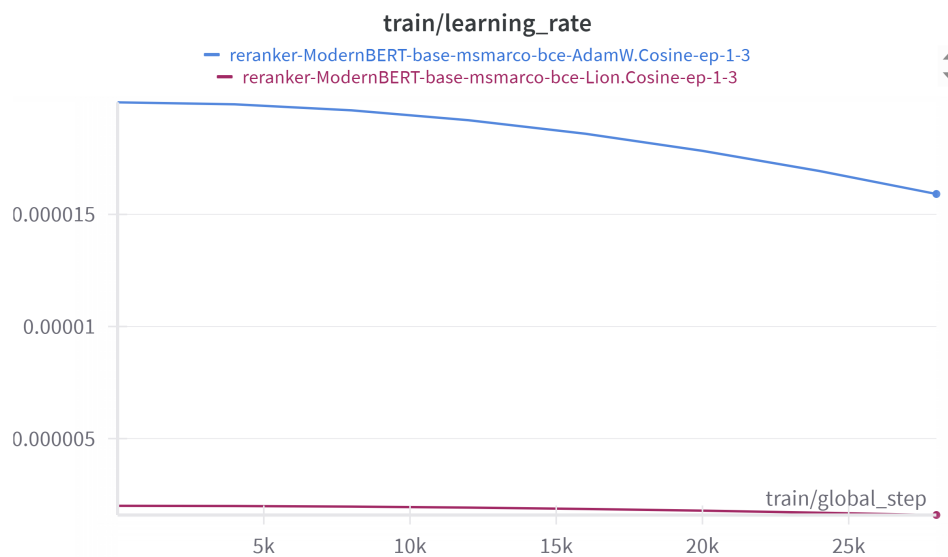


FIGURE 5.3: ModernBERT: Learning Rate Schedule

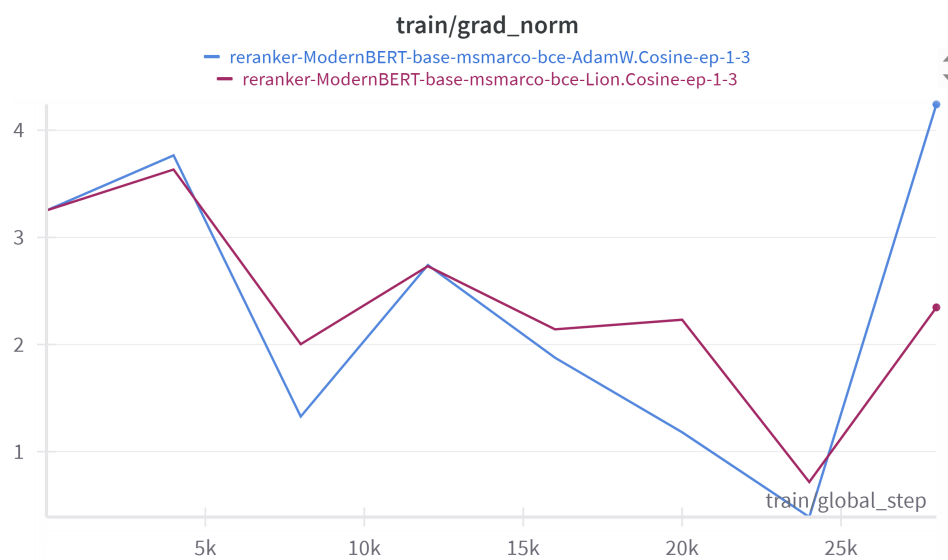


FIGURE 5.4: ModernBERT: Gradient Norm Evolution

GTE Training Dynamics

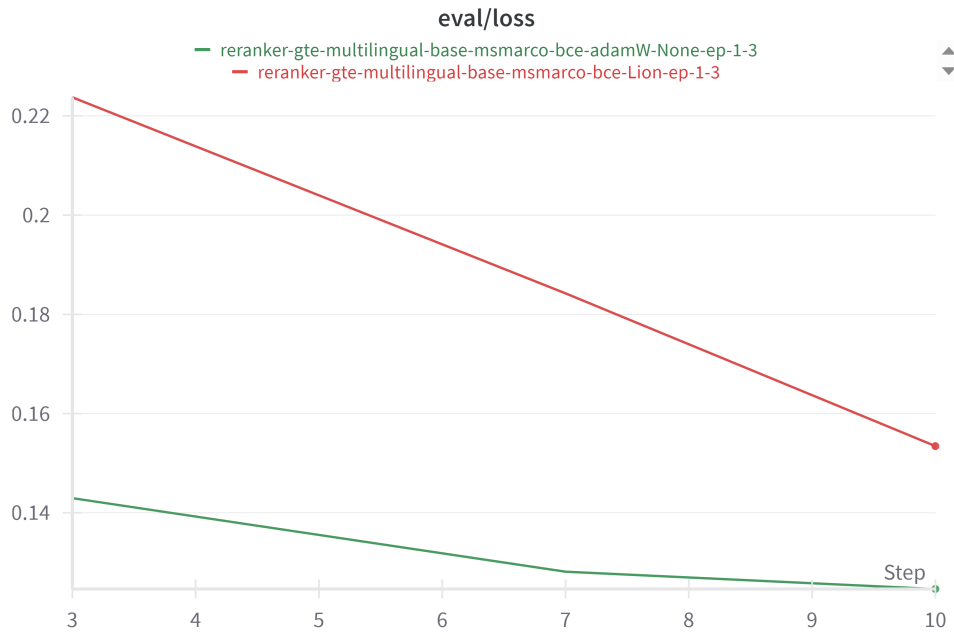


FIGURE 5.5: GTE: Evaluation Loss Comparison

MiniLM Training Dynamics

5.2 Discussion of Results

5.2.1 Optimizer Impact Across Models

The experimental results reveal distinct patterns in how different models interact with the Lion and AdamW optimizers:

- **ModernBERT Performance:** With Lion optimizer and specialized training configuration (lower learning rate of $2e-6$ and Cosine Annealing scheduler), ModernBERT achieved the highest overall performance on TREC DL 2019 metrics (NDCG@10: 0.7225, MAP: 0.5115). The training dynamics (Figures 5.1-5.4) show more stable convergence with Lion compared to AdamW.

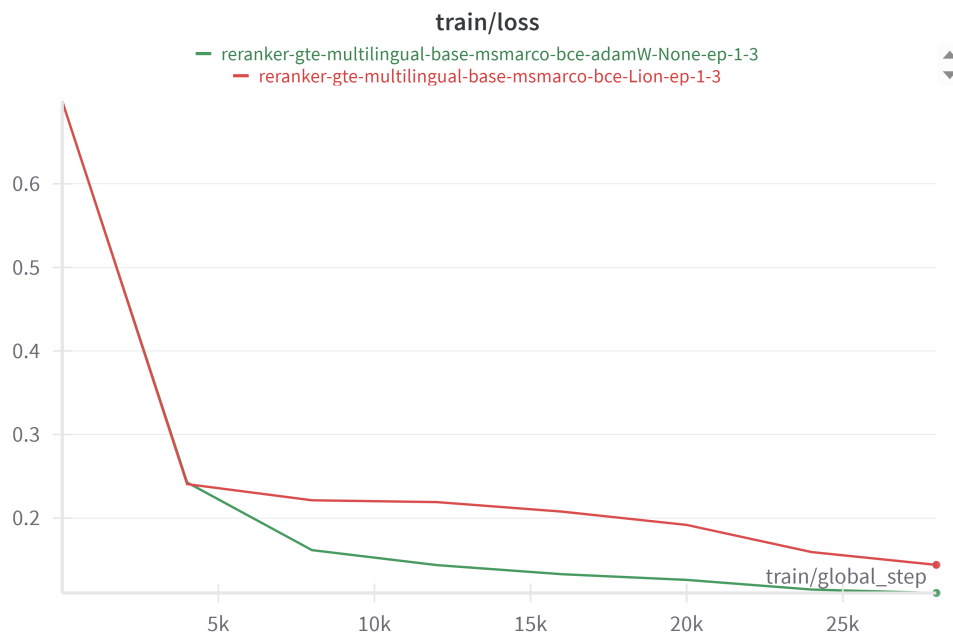


FIGURE 5.6: GTE: Training Loss Progression

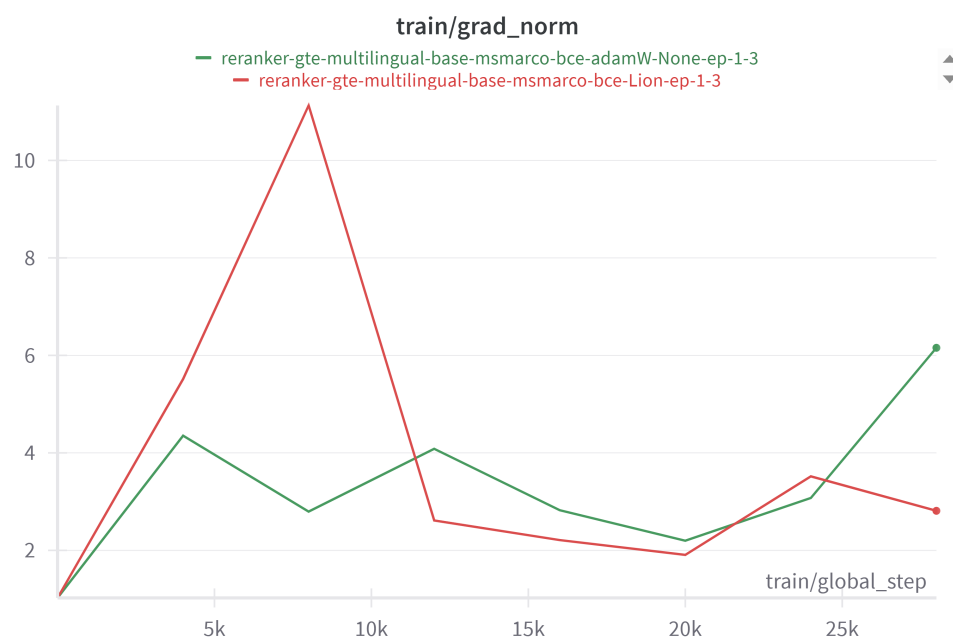


FIGURE 5.7: GTE: Gradient Norm Evolution

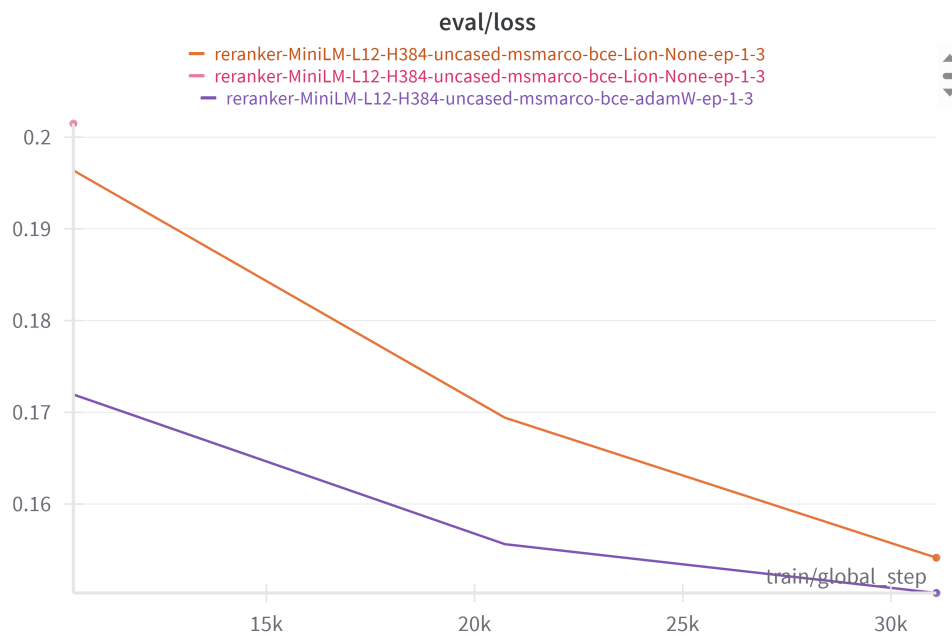


FIGURE 5.8: MiniLM: Evaluation Loss Comparison

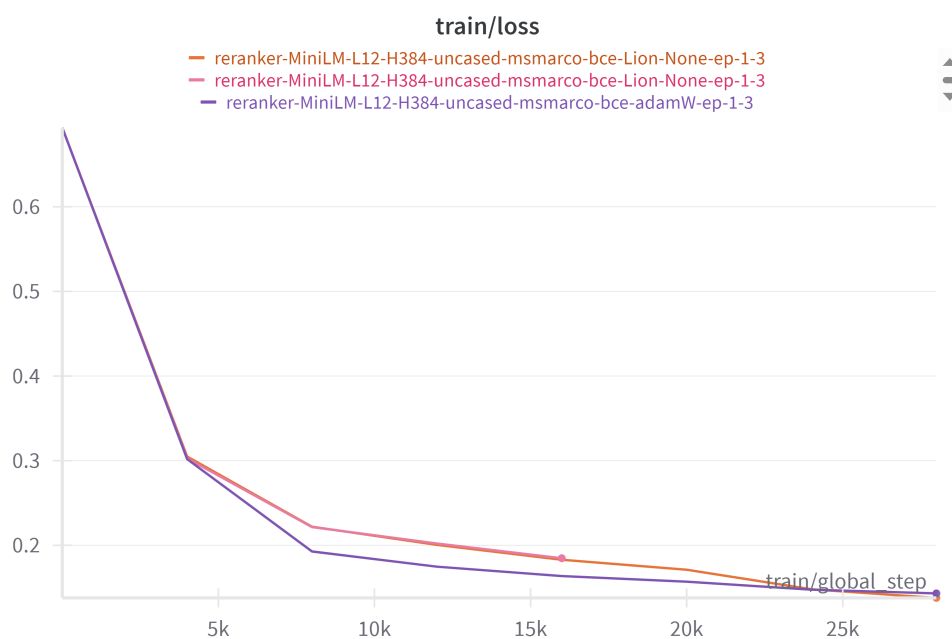


FIGURE 5.9: MiniLM: Training Loss Progression

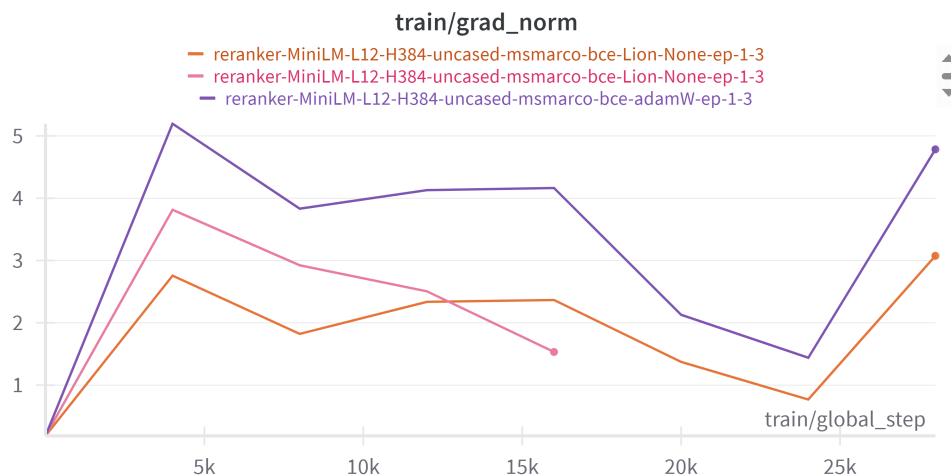


FIGURE 5.10: MiniLM: Gradient Norm Evolution

- **GTE Behavior:** GTE showed stronger performance with AdamW (NDCG@10: 0.7224) using the standard learning rate ($2e-5$). The training curves (Figures 5.5-5.7) indicate that AdamW provided more consistent optimization for this model.
- **MiniLM Characteristics:** While MiniLM with AdamW showed better TREC metrics, the Lion optimizer achieved the highest MRR@10 (0.5988) on MS MARCO dev. The training dynamics (Figures 5.8-5.10) suggest that Lion might be particularly effective for certain ranking scenarios.

5.2.2 Training Dynamics Analysis

The visualization of training dynamics reveals several key insights:

- **Loss Convergence:** Lion generally shows smoother evaluation loss curves compared to AdamW, particularly evident in the ModernBERT experiments (Figure 5.1).
- **Gradient Behavior:** The gradient norm plots (Figures 5.4, 5.7, 5.10) show that Lion maintains more consistent gradient magnitudes throughout training.

- **Learning Rate Impact:** The Cosine Annealing schedule (Figure 5.3) proved particularly effective for ModernBERT with Lion, suggesting that adaptive learning rate strategies can significantly influence optimizer performance.

5.2.3 Model-Specific Considerations

The results highlight important model-specific characteristics:

- **Context Length Impact:** Models with longer context capabilities (GTE and ModernBERT, supporting 8192 tokens) generally outperformed MiniLM on TREC DL metrics, suggesting the benefit of extended context for reranking.
- **Architecture Influence:** ModernBERT's advanced features (Rotary Positional Embeddings, Flash Attention) appear to synergize well with Lion's optimization approach, particularly with appropriate learning rate scheduling.
- **Model Size Considerations:** Despite being smaller, MiniLM showed competitive performance, especially on MRR@10, indicating that model size alone doesn't determine reranking effectiveness.

Chapter 6

Summary & Future Scope of Work

This chapter presents a comprehensive summary of our research findings on optimizer effectiveness in cross-encoder reranking and outlines promising directions for future investigation. We reflect on the key insights gained from our experimental analysis and discuss potential avenues for extending this work.

6.1 Summary of Research

Our investigation into the comparative effectiveness of Lion and AdamW optimizers for cross-encoder reranking has yielded several significant findings:

6.1.1 Key Findings

1. Optimizer-Model Interactions:

- The effectiveness of optimizers showed strong dependence on model architecture and training configuration
- ModernBERT achieved optimal performance with Lion optimizer using specialized learning rate settings

- GTE demonstrated superior performance with AdamW under standard training parameters
- MiniLM showed varying preferences between optimizers depending on the evaluation metric

2. Performance Achievements:

- ModernBERT with Lion achieved best-in-class results:
 - NDCG@10: 0.7225
 - MAP: 0.5115
 - R-Precision: 0.5183
- Both MiniLM and ModernBERT with Lion achieved state-of-the-art MRR@10 (0.5988) on MS MARCO dev

3. Training Dynamics:

- Lion demonstrated more stable evaluation loss curves
- Cosine Annealing learning rate schedule proved particularly effective with Lion
- Gradient behavior showed distinct patterns between optimizers

6.1.2 Technical Insights

The research revealed several important technical considerations:

- **Learning Rate Sensitivity:** The dramatic impact of learning rate selection on Lion's performance suggests careful tuning is essential

- **Context Length Benefits:** Models supporting longer contexts (8192 tokens) showed generally superior performance
- **Architecture Synergies:** Modern architectural features (RoPE, Flash Attention) appeared to complement Lion's optimization characteristics

6.2 Future Scope of Work

Our findings open several promising avenues for future research:

6.2.1 Technical Extensions

1. Hyperparameter Optimization:

- Comprehensive grid search for Lion's optimal parameters across different model sizes
- Investigation of alternative learning rate schedules
- Exploration of momentum parameter impacts

2. Architecture Studies:

- Evaluation of Lion's effectiveness on emerging transformer variants
- Investigation of optimization patterns in multi-query attention mechanisms
- Analysis of position embedding schemes' interaction with different optimizers

3. Scaling Studies:

- Analysis of optimizer behavior with larger model sizes
- Investigation of training stability at different batch sizes

- Evaluation of memory efficiency at scale

6.2.2 Application Extensions

Several practical applications deserve further investigation:

- **Document-Level Reranking:** Leveraging the 8K context capability for full document reranking
- **Multi-Stage Ranking:** Investigating optimizer impact in cascade ranking architectures
- **Cross-Lingual Applications:** Extending the analysis to multilingual reranking scenarios
- **Domain Adaptation:** Studying optimizer effectiveness in domain transfer settings

6.2.3 Theoretical Investigations

Future work should also address theoretical aspects:

- Analysis of Lion’s convergence properties in ranking optimization
- Mathematical characterization of optimizer-architecture interactions
- Theoretical bounds on performance with different optimizers

6.3 Recommendations for Practitioners

Based on our findings, we recommend:

- Consider Lion optimizer for modern architectures with appropriate learning rate scheduling
- Maintain AdamW as a robust baseline, especially for established architectures
- Carefully tune learning rates when using Lion optimizer
- Monitor training dynamics through multiple metrics for optimal checkpoint selection

6.4 Concluding Remarks

This research has demonstrated the potential of the Lion optimizer in cross-encoder reranking while highlighting the complexity of optimizer-model interactions. The findings provide a foundation for both practical applications and future research directions in neural information retrieval.

Acknowledgments

We gratefully acknowledge Modal Labs (<https://modal.com/>) for providing the cloud computing infrastructure and GPU resources (NVIDIA A100-80GB) that made this research possible. Their platform enabled efficient experimentation with large-scale models and datasets, contributing significantly to the comprehensive nature of our analysis.