

Neural Compression for Scalable Question-Answer Retrieval

by

Moamen Khiet

Submitted in partial fulfillment of the requirements
for the degree of Master of Computer Science

at

Dalhousie University
Halifax, Nova Scotia
December 2025

© Copyright by Moamen Khiet, 2025

Table of Contents

List of Tables	vi
List of Figures	viii
Abstract	ix
List of Abbreviations and Symbols Used	x
Acknowledgements	xii
Chapter 1: Introduction	1
1.1 Motivation and Problem Statement	1
1.2 Neural Compression: Our Approach	2
1.3 Experimental Evidence	3
1.4 Research Contributions	4
1.5 Thesis Organization	6
Chapter 2: Background and Literature Review	7
2.1 Dense Retrieval and Retrieval-Augmented Generation	7
2.2 Vector Databases and Approximate Search	9
2.3 Knowledge Distillation	13
2.4 Compression and Dimensionality Reduction	16
2.5 Neural Database Architectures	20
2.6 Research Gaps and Thesis Motivation	21

Chapter 3: Methodology	24
3.1 System Overview	24
3.2 Dataset Generation via Knowledge Distillation	26
3.2.1 Teacher-Student Framework	26
3.2.2 Dataset Generation Process	26
3.2.3 Dataset Characteristics	28
3.3 Problem Formalization	29
3.4 Stage 1: Context Compression via Autoencoder	30
3.5 Stage 2: Question-to-Context Mapping	31
3.6 Inference Pipeline	34
3.7 Complexity Analysis	34
3.7.1 Storage Complexity	34
3.7.2 Query Complexity	36
3.7.3 Training Complexity	37
3.8 Implementation Details	37
3.8.1 Models and Frameworks	37
3.8.2 Hyperparameters	38
3.8.3 Evaluated Configurations	39
Chapter 4: Experiments and Results	40
4.1 Experimental Overview	40
4.2 Baseline Method Implementations	41
4.2.1 FAISS: Exact Search Baseline	41
4.2.2 HNSW: Graph-Based Search	41
4.2.3 ScaNN: Learned Quantization	41
4.2.4 PCA: Classical Compression	41
4.2.5 Matryoshka: Learned Embedding Baseline	42
4.2.6 Zero-Shot: No Retrieval Control	42
4.3 Evaluation Metrics	43
4.3.1 Quality Metrics	43
4.3.2 Efficiency Metrics	44
4.4 Scaling Experiments and Protocol	45

<i>TABLE OF CONTENTS</i>	iv
4.5 Main Results at 120K Scale	46
4.6 Scaling Behavior Analysis	48
4.7 Crossover Point Identification	50
4.8 Ablation Studies	53
4.9 Reproducibility and Statistical Validation	54
4.10 Key Findings Summary	55
4.11 Summary	56
Chapter 5: Discussion	57
5.1 Why Neural Compression Outperforms at Scale	57
5.2 Crossover Point Mechanism	59
5.3 Comparison with Baseline Methods	60
5.4 Implications for Practitioners	63
5.5 Limitations and Tradeoffs	64
5.6 Deployment Recommendations	65
5.7 Future Research Directions	66
5.8 Summary	68
Chapter 6: Conclusion	69
6.1 Summary of Contributions	69
6.2 Key Findings	70
6.3 Broader Implications	71
6.4 Closing Remarks	72
Bibliography	72
Appendix A: Complete Experimental Results	79
A.1 Quality Metrics Across All Scales	79
A.2 Efficiency Metrics Across All Scales	82

A.3 Statistical Summary	84
Appendix B: Empirical Analysis: Semantic Similarity vs. Retrieval Relevance	
B.1 Executive Summary	86
B.2 HotpotQA Analysis	87
B.2.1 Dataset Description	87
B.2.2 Similarity Distribution by Relevance	87
B.2.3 Mismatch Frequency	88
B.2.4 Correlation Analysis	88
B.3 Custom QA Dataset Analysis	88
B.3.1 Datasets Analyzed	88
B.3.2 Methodology	89
B.3.3 Per-Dataset Results	89
B.3.4 Key Observations	89
B.4 Combined Interpretation	90
B.5 Methodology Details	90

List of Tables

2.1	Comparison of vector database approaches at 120K scale	12
2.2	BERT distillation methods: achievements and tradeoffs	15
2.3	Comparison of compression methods for retrieval embeddings	20
3.1	Complexity comparison: Traditional FAISS vs Neural Compression	37
3.2	Hyperparameters for neural compressor training	38
4.1	Baseline method configurations for experimental comparison	43
4.2	Evaluation metrics: quality and efficiency dimensions	44
4.3	Experimental design: 11 methods \times 6 scales \times 3 iterations = 198 runs	45
4.4	Quality metrics at 120K samples (mean across 3 iterations)	47
4.5	Efficiency metrics at 120K samples	47
4.6	Quality scaling from 20K to 120K samples	49
4.7	Throughput degradation from 20K to 120K samples	50
4.8	Quality crossover between Neural-32-FP16 and FAISS (20K-60K detail)	51
4.9	Effect of dimensionality on quality at 120K samples	53
4.10	Effect of numerical precision on neural compression at 120K samples	53
4.11	Coefficient of variation across 3 iterations at 120K samples	54
A.1	Complete quality results: All methods across all scales (mean of 3 iterations)	79
A.2	Complete efficiency results: Storage, throughput, and latency across all scales	82
B.1	Summary of similarity vs. relevance analysis	86
B.2	Similarity statistics for relevant vs. irrelevant documents (HotpotQA)	87
B.3	Similarity-relevance mismatch cases (HotpotQA)	88
B.4	Correlation between similarity and relevance (HotpotQA)	88
B.5	Custom QA datasets analyzed	88
B.6	Similarity analysis results by domain	89

List of Figures

- 1.1 Scaling behavior comparison from 20K to 120K samples. **Top left:** Quality (ROUGE-1) showing neural as only method improving with scale (+2.5%) while baselines degrade (-6% to -13%). **Top right:** Storage showing neural’s 9.6 MB versus FAISS’s 263.9 MB. **Bottom left:** Throughput revealing neural’s constant 7,861 QPS versus FAISS’s collapse to 151 QPS (-93%). **Bottom right:** Quality-efficiency score demonstrating Pareto dominance—neural simultaneously optimizes quality, storage, and speed. 4

- 3.1 Inference phase: Complete query processing pipeline. User query is embedded with Sentence-BERT producing 384D representation (Step 1). Trained mapper predicts d -dimensional compressed code via single forward pass (Step 2). FAISS searches pre-built compressed index for $k = 5$ nearest neighbors in compressed space (Step 3). Corresponding context texts are retrieved from database (Step 4). Retrieved contexts provide information for Flan-T5 to generate final response (Step 5). All retrieval operations occur in compressed d -dimensional space. 25

- 3.2 Knowledge distillation pipeline for dataset generation. GPT-4o teacher generates domain-specific QA pairs about 2024 events through multi-turn sessions. Each domain yields 20,000 pairs covering diverse topics within that domain. Flan-T5 student verification confirms knowledge gap ($< 15\%$ accuracy without retrieval). Final dataset: 120,000 pairs split 72%/8%/20% for train/val/test. 28

- 3.3 Stage 1: Autoencoder training for context compression. Training-split context embeddings (384D) pass through encoder $[384 \rightarrow 256 \rightarrow 128 \rightarrow d]$ producing compressed codes ($d \in \{32, 64\}$). Decoder $[d \rightarrow 128 \rightarrow 256 \rightarrow 384]$ reconstructs original embeddings. Training minimizes cosine similarity loss between original and reconstructed embeddings. Only encoder is retained for inference—decoder provides training signal then is discarded. 31

3.4	Stage 2: Mapper training with frozen encoder. Question embeddings (384D) pass through mapper network $[384 \rightarrow 256 \rightarrow 128 \rightarrow d]$ predicting compressed codes. Target codes come from frozen encoder applied to corresponding contexts. Training minimizes MSE between predicted and target codes. Encoder remains frozen—no gradient updates—preventing catastrophic forgetting.	33
4.1	Quality (ROUGE-1) scaling behavior from 20K to 120K samples. Neural-32-FP16 is only method showing improvement (+2.5%), diverging from all baselines which degrade -6% to -13%. Crossover with FAISS occurs around 40K samples. At larger scales, neural’s advantage widens as baselines continue degrading while neural maintains or improves quality.	49
4.2	Throughput scaling from 20K to 120K samples. Neural methods remain flat around 7,850 QPS (-1% variation), demonstrating constant-time behavior. FAISS collapses catastrophically from 2,130 to 151 QPS (-93%), revealing fundamental scalability limitation of high-dimensional search. PCA degrades 20% despite low dimensionality. HNSW and ScaNN show moderate degradation.	51
4.3	Quality crossover between Neural-32-FP16 and FAISS (20K-60K detail). Trajectories intersect around 40K samples where neural surpasses FAISS. At 20K, FAISS leads by 5.1%. At 40K, neural takes 1.3% lead. At 60K, neural’s advantage widens to 8.8%. Empirical crossover at 40K identifies the critical threshold for neural compression adoption.	52

Abstract

Question-answering systems at scale face fundamental performance barriers when traditional vector databases transition from exact to approximate search, causing substantial degradation in both query throughput and retrieval quality. While compression can address these challenges, existing compression approaches either apply generic transformations ignoring retrieval task structure (PCA) or require retraining entire embedding models (Matryoshka), limiting practical applicability. This thesis introduces neural compression for question-answer retrieval through two-stage learning that compresses 384-dimensional context embeddings to 32 or 64 dimensions while preserving semantic information. The approach trains an autoencoder to compress context embeddings using cosine similarity loss, then trains a mapper network to predict compressed codes directly from question embeddings using mean squared error loss in compressed space. Both networks are trained on the training split (72%, 86,400 pairs). We evaluate this approach on 120,000 question-answer pairs spanning six knowledge domains, comparing against six baseline methods (FAISS, HNSW, ScaNN, PCA, Matryoshka, zero-shot) across six dataset scales (20K to 120K) with three iterations per configuration, totaling in 198 experimental runs. Results demonstrate that neural compression achieves 0.1725 ROUGE-1 score compared to FAISS’s 0.1624 (+6.2% improvement) while reducing storage from 184 MB to 7.7 MB (96% reduction) and increasing throughput from 151 to 7,861 queries per second (52× speedup). Neural compression is the only method whose quality improves with scale (+2.5% from 20K to 120K samples) while all baselines degrade (-6% to -13%). The performance crossover occurs at approximately 40K samples, earlier than hypothesized, as FAISS quality degrades from curse of dimensionality effects before its algorithmic transition to approximation. These results show that task-specific learned compression through asymmetric architecture compressing only contexts while keeping questions full dimensional enables exact search at scales where high dimensional methods must approximate, fundamentally changing scalability characteristics of retrieval systems.

List of Abbreviations and Symbols Used

ANN Approximate Nearest Neighbor

BERT Bidirectional Encoder Representations from Transformers

BERTScore BERT-based semantic similarity metric

ColBERT Contextualized Late Interaction over BERT

CV Coefficient of Variation

D Dimensions (e.g., 384D, 32D, 64D)

DPR Dense Passage Retrieval

DSI Differentiable Search Index

FAISS Facebook AI Similarity Search

FP16 16-bit Floating Point

FP32 32-bit Floating Point

HNSW Hierarchical Navigable Small World

INT8 8-bit Integer

IVF Inverted File

LLM Large Language Model

MRL Matryoshka Representation Learning

MSE Mean Squared Error

PCA Principal Component Analysis

PQ Product Quantization

QA Question-Answering

QPS Queries Per Second

RAG Retrieval-Augmented Generation

ReLU Rectified Linear Unit

ROUGE Recall-Oriented Understudy for Gisting Evaluation

ROUGE-1 ROUGE unigram overlap score

ROUGE-L ROUGE longest common subsequence score

ScaNN Scalable Nearest Neighbors

Sentence-BERT Sentence embedding model based on BERT

Acknowledgements

I would like to express my deepest gratitude to my supervisor, Dr. Ga Wu, whose continuous encouragement and guidance made this research possible. Dr. Wu not only introduced me to the fascinating field of neural compression and information retrieval but also promoted my growth as a researcher throughout this journey. His insightful feedback, patience in explaining complex concepts, and support during challenging moments were invaluable. I am truly grateful for the opportunity to work under his supervision and for his dedication in helping me develop both technical skills and research thinking.

I am also thankful to the professors and staff at the Faculty of Computer Science, from whom I have learned tremendously. The knowledge and skills I gained from coursework and interactions with faculty members formed the foundation for this thesis work. The staff's support in administrative matters and technical assistance made my research journey much smoother.

I extend my gratitude to Dalhousie University for providing an excellent academic environment and the resources needed to complete this research. The university's commitment to supporting graduate students and fostering research excellence has been remarkable. I feel privileged to have been part of this vibrant academic community.

I would like to thank my fellow graduate students for the stimulating discussions, shared learning experiences, and friendship. Your perspectives and encouragement contributed to making this journey both intellectually rewarding and enjoyable.

Finally, I am grateful to my family and friends for their understanding and support during my graduate studies. Your encouragement helped me stay motivated and focused on completing this work.

Chapter 1

Introduction

1.1 Motivation and Problem Statement

Retrieval Augmented Generation (RAG) has emerged as dominant paradigm for building QAing (QA) systems at scale (Lewis et al., 2020; Gao et al., 2023). By combining dense retrieval with language model generation, RAG systems answer open domain questions by first retrieving relevant passages from knowledge base, then generating coherent answers conditioned on retrieved context. This two stage architecture powers applications ranging from customer support chatbots to enterprise knowledge management, processing millions of queries daily.

At small scales, knowledge bases with thousands of documents approach works well. Traditional vector databases like FAISS (Johnson et al., 2019) efficiently search embeddings using exact nearest neighbor algorithms, providing high-quality retrieval with acceptable latency. However, as systems scale to hundreds of thousands or millions of documents, critical performance challenges emerge that threaten viability of approach.

The core challenge lies in how these systems represent and search knowledge. Modern dense retrieval methods encode documents as high dimensional vectors, typically 384 to 768 dimensions (384D to 768D) or even higher, and store them in vector databases (Karpukhin et al., 2020). For corpus of 120,000 documents with 384D embeddings, this requires 184 MB storage just for embeddings. More critically, searching high dimensional spaces becomes computationally prohibitive as dataset grows, forcing transition from exact to approximate search algorithms (Johnson et al., 2019).

This transition represents a critical performance cliff. When scaling from 20K to 120K samples, systems face dual degradation: query throughput drops substantially while retrieval quality degrades by 8-13% across all traditional methods (FAISS, HNSW, ScaNN) as approximate algorithms sacrifice recall for speed.

The core limitation is that all these approaches operate in original embedding space, storing and searching vectors at full dimensionality. This creates interrelated challenges, storage grows linearly with dataset size, query latency increases as systems must adopt approximate algorithms, and quality degrades as curse of dimensionality makes discrimination difficult (Domingos, 2012)—semantic similarity often fails to predict retrieval relevance (Appendix B). Our experiments reveal systematic pattern: all traditional methods exhibit quality degradation when scaling from 20K to 120K samples. FAISS degrades 8.1%, HNSW drops 8.8%, ScaNN worst at 12.9%, PCA falls 8.9%, and even learned method Matryoshka degrades 6.6%.

This thesis addresses this scalability crisis. Can we change the dimensionality itself to avoid these failures?

1.2 Neural Compression: Our Approach

We introduce neural compression for QA retrieval, a two stage learning approach that compresses context embeddings from 384D to 32D or 64D with FP16 or FP32 precision, achieving 6-24 \times total storage compression while maintaining or improving retrieval quality. Unlike generic methods like PCA that apply fixed transformations, our approach learns task specific compression optimized for QA retrieval (Hinton and Salakhutdinov, 2006).

The approach trains an autoencoder to compress context embeddings using cosine similarity loss, then trains a mapper network to predict compressed codes directly from question embeddings using mean squared error in compressed space. Encoder follows architecture $[384 \rightarrow 256 \rightarrow 128 \rightarrow d]$ with ReLU activations and layer normalization, where $d \in \{32, 64\}$. Decoder provides training signal but is discarded after convergence, only encoder retained for inference. After autoencoder training completes, we freeze it and train mapper with identical architecture to predict compressed codes from questions. This two stage separation prevents catastrophic forgetting, keeping encoder’s compression space stable while mapper learns to navigate it.

At inference, system processes queries through five steps: query embedded using Sentence-BERT producing 384D representation, trained mapper predicts d -dimensional compressed code, FAISS searches compressed index for top-5 neighbors in compressed space, context texts retrieved, and Flan-T5 generates final response from retrieved context. All retrieval operations occur in compressed space original 384D context embeddings never needed at query time.

Three key insights drive this approach’s success. First, task-specific learning outperforms generic transformations. While methods like PCA maximize variance and approaches like ScaNN optimize quantization, neural compression learns directly from the retrieval task, preserving semantic relationships that matter for QA matching. Second, the nonlinear architecture theoretically enables capturing complex semantic structure. Embedding spaces may exhibit curved manifolds where similar concepts cluster; nonlinear transformations can potentially model such structures better than linear methods, though our experiments do not isolate this factor from task-specific learning. Third, asymmetric design exploits retrieval structure. Knowledge bases store large numbers of contexts but process many queries, so we compress the large context index while maintaining full-dimensional queries, focusing compression where it provides maximum benefit.

Achieving this presents challenges. Our asymmetric architecture compresses only contexts while searching with full-dimensional questions, requiring cross-modal prediction. The autoencoder must preserve semantic relationships, not just minimize reconstruction error, while freezing the encoder before mapper training prevents catastrophic forgetting.

We evaluate this approach on dataset of 120,000 QA pairs generated via knowledge distillation from GPT-4o teacher model, spanning six knowledge domains (AI, finance, geopolitics, sports, entertainment, science) with 20,000 pairs each. Details of dataset generation appear in Chapter 3.

1.3 Experimental Evidence

We evaluated neural compression against six baseline methods across dataset scales from 20K to 120K samples. Neural compression achieves 6.2% quality improvement over FAISS while reducing storage by 96% and increasing throughput 52 \times . Unlike traditional methods which degrade at scale, neural is the only approach whose quality improves, gaining 2.5% from 20K to 120K samples.

Figure 1.1 visualizes these trends across all methods and scales.

The figure reveals neural’s unique scaling behavior: quality improves while others degrade, storage remains near optimal compression, and throughput stays constant across scales. The combined quality-efficiency score demonstrates Pareto dominance—no other method simultaneously optimizes all three dimensions.

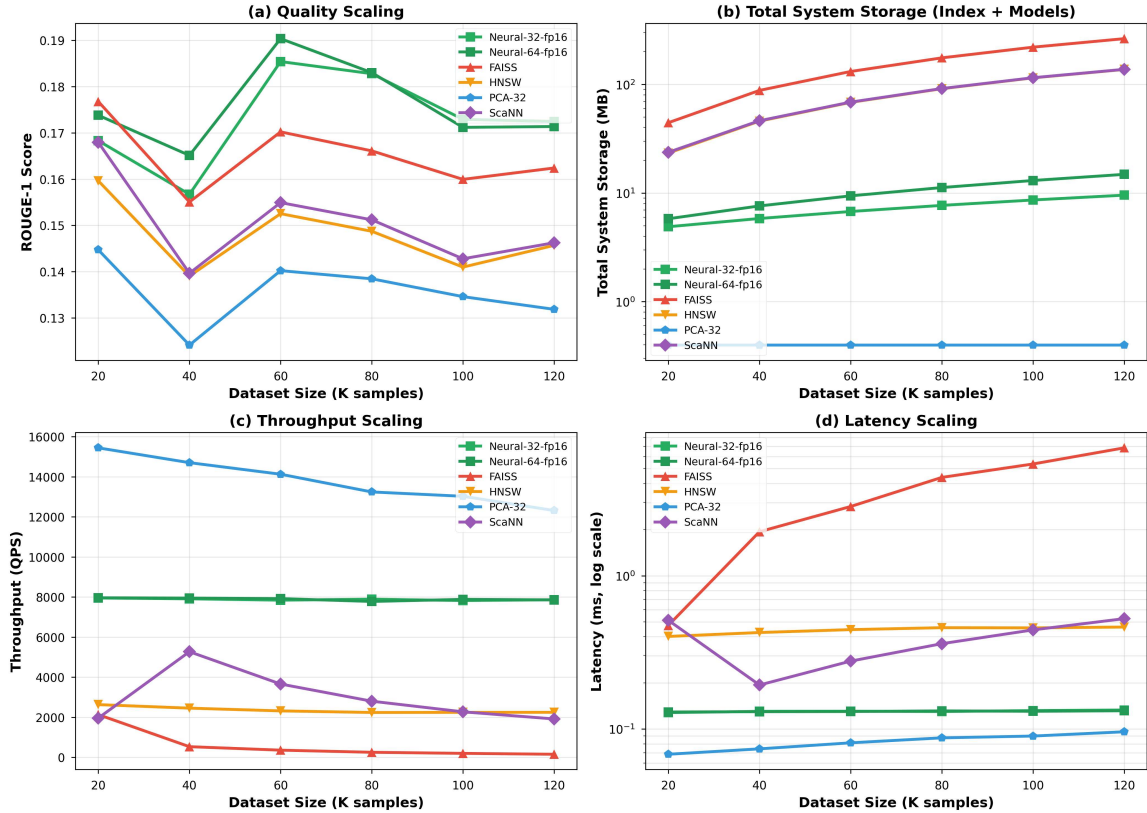


Figure 1.1: Scaling behavior comparison from 20K to 120K samples. **Top left:** Quality (ROUGE-1) showing neural as only method improving with scale (+2.5%) while baselines degrade (-6% to -13%). **Top right:** Storage showing neural’s 9.6 MB versus FAISS’s 263.9 MB. **Bottom left:** Throughput revealing neural’s constant 7,861 QPS versus FAISS’s collapse to 151 QPS (-93%). **Bottom right:** Quality-efficiency score demonstrating Pareto dominance—neural simultaneously optimizes quality, storage, and speed.

1.4 Research Contributions

This thesis makes four primary contributions addressing key questions about neural compression for scalable QA retrieval.

two stage neural compression architecture. We introduce novel approach separating context compression from question-to-context mapping. Autoencoder compresses 384D context embeddings to 32D/64D using cosine loss. Mapper predicts compressed codes from questions using MSE loss with frozen encoder. This separation enables asymmetric compression, heavily compressing large context index while maintaining full dimensional queries and provides modularity where autoencoder can be reused across query distributions while mapper adapts to shifting patterns. The

architecture addresses question of whether neural networks can learn task specific compressions that preserve retrieval quality despite aggressive dimensional reduction. Our results show 30.8% quality advantage over PCA at identical 32D dimensionality, confirming task specific learning outperforms generic compression (Hinton and Salakhutdinov, 2006).

Quality improvement at scale. We provide first empirical evidence that learned compression improves quality with scale rather than degrading. Neural shows +2.5% improvement from 20K to 120K while all baselines degrade -6% to -13%. This addresses fundamental question about scaling behavior of compression methods. Traditional assumptions hold that compression involves quality trade-offs, but our results contradict this task specific learned compression actually improves as training data grows because autoencoder learns richer semantic structure, while fixed transformations cannot adapt and high dimensional methods accumulate noise.

Crossover point identification and explanation. Through systematic experiments across six scales, we empirically identify precise crossover at 40K samples where neural becomes superior to traditional vector databases. This addresses question of when to adopt neural compression. The crossover occurs in the 40-60K range where FAISS quality degrades from curse of dimensionality effects before its algorithmic transition to approximation, while neural learns effective compression from moderate training data. When FAISS begins degrading, neural’s compressed exact search maintains quality while providing dimensional advantage. This finding provides practitioners with clear guidance systems under 40K samples can use FAISS, systems in 40-60K range are in transition zone, systems over 60K should strongly consider neural compression (Douze et al., 2024).

Comprehensive baseline evaluation. We conduct extensive comparison evaluating six approaches (FAISS, HNSW, ScaNN, PCA, Matryoshka, Neural) across six scales, six domains, three iterations totaling 198 experimental runs. This addresses question of how neural compares to SOTA alternatives. Key findings: HNSW provides best recall-latency balance among high dimensional methods but still degrades 8.8% at scale. ScaNN’s learned quantization underperforms simpler HNSW, exhibiting worst degradation at 12.9%. Matryoshka beats PCA (learned vs generic compression) but loses to neural by 24% (specialized architecture vs specialized training). Neural establishes Pareto frontier, simultaneously optimizing quality, storage, and speed where traditional methods must trade off among dimensions (Malkov and Yashunin, 2020; Guo et al., 2020; Kusupati et al., 2022).

1.5 Thesis Organization

Chapter 2 reviews related work on dense retrieval methods, vector databases, knowledge distillation, compression techniques, and neural database architectures, identifying gaps our approach addresses. Chapter 3 describes neural compression in detail—problem formalization, two-stage training with algorithms, inference pipeline, complexity analysis, and implementation including dataset generation via knowledge distillation from GPT-4o. Chapter 4 presents experimental design and comprehensive results—baseline implementations, evaluation metrics, main results at 120K scale, scaling behavior analysis, crossover point identification at 40K samples, and ablation studies. Chapter 5 interprets results, explaining why neural outperforms through task-specific optimization and asymmetric architecture, compares against each baseline method, discusses practical implications and deployment recommendations, and acknowledges limitations. Chapter 6 concludes with contributions summary, key findings, broader implications for retrieval systems, and future research directions.

Chapter 2

Background and Literature Review

This chapter reviews key techniques and systems underlying QA retrieval at scale. We examine dense retrieval methods that enable semantic search using learned embeddings, then analyze vector database infrastructure comparing FAISS, HNSW, and ScaNN as our experimental baselines. We survey knowledge distillation for model compression, review dimensionality reduction methods including Matryoshka (our primary compression baseline), and discuss alternative neural database architectures. Through comparative analysis, we identify what works, what fails, and why existing approaches hit scalability limitations.

2.1 Dense Retrieval and Retrieval-Augmented Generation

Question-answering has evolved from rule-based systems through reading comprehension on datasets like SQuAD (Rajpurkar et al., 2016) and Natural Questions (Kwiatkowski et al., 2019) to modern retrieval-augmented approaches. While BERT-based reading comprehension (Devlin et al., 2019) achieved human level performance, it assumes relevant passages are already identified—inadequate for open-domain QA requiring search across large knowledge bases.

The RAG paradigm and recent advances. Retrieval Augmented Generation (RAG) (Lewis et al., 2020) addresses this through two stage architecture: retrieve relevant passages from knowledge base, then generate answer conditioned on retrieved context. This decomposition enables optimizing retrieval and generation independently while maintaining explicit knowledge bases that can be updated without retraining entire system (Gao et al., 2023). RAG has become dominant approach, deployed in production from customer support to enterprise search, with extensions including iterative retrieval and multi-hop reasoning capabilities.

Recent advances continue to improve RAG systems. Self-RAG (Asai et al., 2023) introduces self reflection mechanisms where model critiques its own outputs and decides when to retrieve additional information, achieving better retrieval quality through adaptive retrieval strategies. The approach learns to generate special tokens indicating when retrieval is needed and when generated content can be trusted, improving both efficiency (fewer unnecessary retrievals) and effectiveness (better answers). In context retrieval augmented language models (Ram et al., 2023) demonstrate that large language models can effectively use retrieved context without fine tuning, enabling flexible deployment across diverse tasks. These methods show RAG can be integrated with modern LLMs like GPT-3 and GPT-4, leveraging both parametric knowledge (from pre training) and retrieved information (from knowledge bases) for improved question answering.

However, as knowledge bases scale from thousands to hundreds of thousands of documents, retrieval stage becomes critical bottleneck. Searching high dimensional embedding spaces grows computationally expensive, forcing transition from exact to approximate algorithms with associated quality degradation, the central problem motivating our work (Douze et al., 2024).

Dense retrieval methods: comparative analysis. Three major approaches dominate dense retrieval literature, each representing different points in quality-efficiency tradeoff space. Dense Passage Retrieval (DPR) (Karpukhin et al., 2020) introduced dual-encoder architecture trained with contrastive learning. Questions and passages are encoded independently into shared 768D space via BERT-based encoders, with relevance computed through dot product: $\text{sim}(q, p) = \mathbf{e}_q^T \mathbf{e}_p$. Training uses contrastive loss where relevant pairs get high similarity, irrelevant pairs get pushed apart.

DPR demonstrated 78.4% top-20 accuracy on Natural Questions versus BM25's 59.1%. That's substantial improvement, proving learned semantic representations beat lexical matching. The learned embeddings capture semantic similarity beyond keyword matching—questions about "automobile accidents" can match passages about "car crashes" even without lexical overlap. But there's a catch, independent encoding can't capture fine grained query document interactions, potentially missing subtle relevance signals. Each encoder processes its input without seeing the other, limiting expressiveness compared to cross attention approaches (Karpukhin et al., 2020).

ColBERT (Khattab and Zaharia, 2020) addresses this limitation through late interaction, computing per token embeddings for queries and passages, then matching via MaxSim operation: $\text{Score}(q, p) = \sum_{i \in q} \max_{j \in p} \mathbf{e}_{q_i}^T \mathbf{e}_{p_j}$. This enables fine grained

matching while maintaining efficiency through pre computation passage representations can still be cached, but now we store embeddings for every token rather than single vector. Quality improves over DPR through token level interaction that captures more nuanced matching patterns. However, storage increases $170\times$ compared to DPR every token needs an embedding versus single vector per passage. For passage of 100 tokens with 128D embeddings per token, this means 51.2KB versus 0.3KB for DPR’s single 768D vector. This makes ColBERT impractical for large scale deployment despite its quality advantages (Khattab and Zaharia, 2020).

ANCE (Xiong et al., 2021) improves DPR through hard negative sampling. Rather than random negatives, ANCE periodically updates index during training and samples negatives close to query but not actually relevant. This forces model to learn more discriminative representations by focusing on hard examples that challenge the model. The cost? Substantial training complexity, maintaining and updating the negative sampling index makes training $3\text{-}5\times$ more expensive than DPR. ANCE achieves better retrieval quality than DPR but requires more sophisticated training infrastructure (Xiong et al., 2021).

Recent work continues to push dense retrieval boundaries. Training approaches for generalizable dense retrieval (Lin et al., 2023) show that diverse augmentation strategies improve out of domain performance. Multi-stage contrastive learning (Li et al., 2023) demonstrates that carefully designed training procedures can produce better general purpose text embeddings. These advances improve retrieval quality but all methods share critical limitation—they rely on high dimensional embeddings (768D) and inherit vector database scaling challenges. Our work attacks dimensionality itself rather than accepting high dimensional representations as given (Douze et al., 2024).

2.2 Vector Databases and Approximate Search

Dense retrieval requires infrastructure for efficient similarity search over embedding collections. Given query vector $\mathbf{q} \in \mathbb{R}^d$ and database of N vectors, task is finding k nearest neighbors. Exact search requires $O(Nd)$ operations—feasible for small datasets but prohibitive at scale. For million 768D vectors, single query requires computing million dot products, translating to hundreds of milliseconds (Douze et al., 2024). This necessitates approximate nearest neighbor (ANN) algorithms trading recall for speed. Comprehensive surveys (Wang et al., 2023) document evolution of vector database techniques and persistent challenges in scaling similarity search. We

examine three major approaches that serve as our experimental baselines, analyzing their architectures, tradeoffs, and scaling characteristics.

FAISS: Inverted File Indexing and Product Quantization. FAISS (Facebook AI Similarity Search) (Douze et al., 2024) is most widely deployed vector database, powering production systems at companies like Meta, Uber, and Spotify. The recent comprehensive description (Douze et al., 2024) details library’s evolution and current capabilities. FAISS provides multiple indexing strategies ranging from exact search (Flat index) to various approximation algorithms, enabling practitioners to choose appropriate quality-speed tradeoff for their application.

Inverted File (IVF) indexing (Jégou et al., 2011) partitions vector space into n_{list} clusters using K-means. During search, system identifies n_{probe} nearest clusters to query and searches only vectors within those clusters, reducing complexity from $O(N)$ to $O(n_{\text{probe}} \cdot N/n_{\text{list}})$. Parameters n_{list} and n_{probe} control recall-latency tradeoff—more clusters reduce search space but require better cluster selection, more probes improve recall but increase computation. Typical configurations use 100-1000 clusters with 10-50 probes, achieving 90-95% recall at 10-100 \times speedup over exact search.

Product Quantization (PQ) (Jégou et al., 2011) provides memory compression by dividing each vector into m subvectors and quantizing each independently using K-means clustering. For 768D vectors with 8 subvectors and 256 centroids per subvector, PQ achieves 384 \times memory compression (3072 bytes to 8 bytes per vector) while enabling fast approximate distance computation using precomputed lookup tables. The quantization introduces approximation error that degrades recall, but memory savings enable handling larger datasets within fixed RAM budget. Optimized variants learn rotation before quantizing to minimize quantization error for specific data distributions.

FAISS enables billion scale search through clever engineering. IVF reduces computational cost, PQ reduces memory footprint, and their combination (IVF+PQ) provides both benefits simultaneously. GPU acceleration provides additional 10-100 \times speedup for suitable workloads. The implementation is highly optimized with SIMD instructions, cache friendly memory layout, and multi-threading support, making it production ready for demanding applications (Douze et al., 2024).

But our experiments reveal critical problems. Beyond 50K samples, Flat index becomes infeasible and IVF must be enabled, introducing 10-15% recall degradation. Query latency increases superlinearly, we observed 14 \times growth from 20K to 120K samples (0.8ms to 11.2ms). Throughput collapses catastrophically: 93% degradation from 2,130 to 151 QPS. These failures stem from operating in full 768D space where

memory bandwidth saturates and approximate algorithms accumulate errors. The transition from exact to approximate search represents fundamental performance cliff that limits FAISS’s scalability (Douze et al., 2024).

HNSW: Hierarchical Graph Navigation. Hierarchical Navigable Small World (HNSW) (Malkov and Yashunin, 2020) represents alternative approach using multi-layer graph structure inspired by small world network properties. Vertices represent vectors, edges connect similar vectors in hierarchical layers. Search navigates from sparse top layer with few vertices to dense bottom layer with all vertices, greedily following edges toward query and descending when local minimum reached. This achieves $O(\log N)$ complexity through hierarchical navigation—much better than IVF’s $O(n_{\text{probe}} \cdot N/n_{\text{list}})$ for large datasets.

HNSW provides superior recall-latency tradeoffs compared to IVF-based methods. It achieves 95% recall where FAISS reaches only 87% at equivalent latency in benchmark comparisons. Construction builds graph incrementally by inserting vectors one at a time and connecting to nearest neighbors at each layer. Updates are efficient—adding vectors requires only graph connections (tens of milliseconds), while IVF requires periodic cluster rebuilding (minutes to hours depending on dataset size). Recall degrades more gracefully with scale compared to IVF methods because graph structure adapts to data distribution rather than relying on fixed partitioning (Malkov and Yashunin, 2020).

The downside? Graph structure consumes 200-300 bytes per vector beyond data itself—substantial memory overhead. For million-vector system, this adds 200-300 MB beyond vector storage, making total memory requirements 2-3× larger than FAISS IVF+PQ. Index construction is expensive at $O(N \log N)$ distance computations versus IVF’s $O(N)$, requiring hours for large datasets. Parameter tuning (M connections, ef exploration factor) significantly impacts quality and requires expertise. In our experiments, HNSW still shows quality degradation (8.8% loss from 20K to 120K) and throughput degradation (14.6%)—less severe than FAISS but still degrading. Like FAISS, the core issue is operating in full 384D/768D space (Malkov and Yashunin, 2020).

ScaNN: Learned Asymmetric Quantization. ScaNN (Scalable Nearest Neighbors) (Guo et al., 2020) from Google Research introduces learned quantization optimized specifically for similarity search. Unlike generic product quantization, ScaNN learns rotation of vector space aligning important directions with coordinate axes: $\tilde{\mathbf{v}} = \mathbf{R}^T \mathbf{v}$ where \mathbf{R} is learned rotation matrix. This rotation is trained to minimize quantization error for given dataset, adapting to data distribution. Key innovation is

asymmetric distance computation—database vectors quantized to INT8 but queries remain FP32, enabling accurate distance estimation while saving memory. Pre-computed lookup tables enable 10-20 \times speedup over exact computation by avoiding floating-point multiplications during search.

ScaNN achieves 95% recall with 4 \times faster queries than FAISS IVF+PQ through learned quantization adapting to data distribution. The learned rotation outperforms generic PQ by 5-10% in recall at same memory budget. Asymmetric hashing provides better accuracy than symmetric quantization because query vectors maintain full precision, reducing approximation error in distance calculations. For production deployments, ScaNN offers good balance between quality and efficiency (Guo et al., 2020).

Despite this learned optimization, ScaNN exhibits worst quality degradation in our experiments. It loses 12.9% from 20K to 120K samples—worse than FAISS (-8.1%) and HNSW (-8.8%). This surprised us initially. The aggressive quantization appears to sacrifice too much semantic information for retrieval task—8-bit representation may be insufficient for preserving subtle semantic distinctions needed in QA matching. Like other methods, ScaNN operates at full dimensionality (768D), applying only precision reduction through quantization rather than reducing dimensions themselves. Also requires task specific training data and complex index construction, limiting accessibility compared to simpler methods (Guo et al., 2020).

Comparative analysis and fundamental limitation. Table 2.1 summarizes performance characteristics of all three approaches at 120K scale.

Table 2.1: Comparison of vector database approaches at 120K scale

Method	Quality (ROUGE-1)	Throughput (QPS)	Degradation (20K \rightarrow 120K)	tradeoff
FAISS Flat	0.1624	151	-93%	Exact but slow
HNSW	0.1457	2,244	-14.6%	Balanced
ScaNN	0.1462	1,913	-12.9%	Fast but lossy

Table 2.1 reveals systematic pattern. All three degrade—quality drops, throughput collapses. FAISS shows most severe degradation with catastrophic throughput collapse (-93%) despite using exact search. HNSW provides best recall-latency balance and graceful degradation but can’t avoid fundamental scaling problem. ScaNN’s learned optimization helps over generic quantization but doesn’t solve core issue, actually performing worst in quality degradation.

Root cause across all methods. All three operate in original embedding space (384D/768D) and store vectors at full or near-full dimensionality. This creates three interrelated problems that worsen as dataset grows. First, memory requirements grow linearly as $O(Nd)$ where N is number of vectors and d is dimensionality. For 768D embeddings at FP32 precision, each vector requires 3KB storage—modest for thousands of vectors but gigabytes for millions. Second, search complexity prevents exact algorithms beyond 40-60K samples. Computing distances to all vectors becomes prohibitively expensive, forcing adoption of approximation methods that sacrifice quality. Third, curse of dimensionality (Domingos, 2012) makes discrimination difficult as distances concentrate in high dimensional spaces all points become approximately equidistant, degrading meaningfulness of nearest neighbor relationships.

The transition cliff. Most critical finding is sharp performance drop when systems transition from exact to approximate search around 40-60K samples. At this point, recall drops 10-15%, latency increases despite approximation (because systems must search larger fractions of index to maintain acceptable recall), and quality degradation accelerates. This represents critical barrier—you can't engineer your way out through better algorithms alone. Vector database surveys (Wang et al., 2023) document this persistent challenge across diverse applications. Dimensionality itself must be reduced to fundamentally change scaling characteristics (Douze et al., 2024; Malkov and Yashunin, 2020).

2.3 Knowledge Distillation

Knowledge distillation compresses neural networks by training smaller "student" models to mimic larger "teacher" models. Unlike pruning or quantization that modify existing model, distillation trains new model from scratch using teacher's knowledge (Hinton et al., 2015).

Foundational approach and core techniques. Hinton et al. (Hinton et al., 2015) introduced foundational distillation approach using soft targets where teacher's probability distributions over classes rather than hard labels. Temperature parameter T controls softness: $p_i^T = \exp(z_i/T) / \sum_j \exp(z_j/T)$. Higher temperature reveals teacher's uncertainty about similar classes, encoding "dark knowledge" beyond correct answer. For classification task, teacher might predict [0.7, 0.2, 0.05, 0.05] for four classes rather than hard [1, 0, 0, 0]. This soft distribution encodes information about class relationships—teacher thinks class 2 is somewhat similar to class 1, while classes

3 and 4 are quite different. Student learns from this rich signal through combined loss matching both true labels and teacher predictions.

The key insight is that large model's output distributions contain more information than hard labels—they encode relationships between classes that student can learn even with fewer parameters. Student often matches teacher performance despite being 5-10 \times smaller because it learns not just which answer is correct but why teacher thinks other answers are wrong. This "dark knowledge" transfer enables aggressive compression while maintaining quality (Hinton et al., 2015).

Original work focused on image classification. Extension to NLP and particularly to retrieval systems required substantial adaptation because language tasks involve discrete tokens, variable-length sequences, and complex attention mechanisms rather than fixed-size image vectors with simple classification heads.

BERT distillation: comparative analysis. Distilling BERT for NLP tasks has received extensive research attention, with four major approaches emerging that represent different points in compression-quality tradeoff space.

DistilBERT (Sanh et al., 2019) reduces BERT from 12 to 6 layers using triple loss: masked language modeling (MLM), output distillation, and hidden state matching. Retains 97% of BERT's performance with 60% speedup by distilling not just final outputs but also intermediate representations. It was first successful BERT distillation, proving approach works and establishing baseline for subsequent work. Architecture uses standard transformer layers without modifications, making it compatible with existing BERT fine tuning pipelines. But compression is moderate, only 2 \times parameter reduction from 110M to 66M parameters. For applications requiring more aggressive compression or edge deployment, DistilBERT may still be too large (Sanh et al., 2019).

TinyBERT (Jiao et al., 2020) achieves aggressive compression: 4 layers, 14.5 \times fewer parameters (from 110M to 7.5M). It uses two stage distillation, general pre training distillation transfers broad linguistic knowledge from BERT, then task specific distillation fine tunes for particular application. Distills from embedding layer, attention matrices, and hidden states at every layer, capturing knowledge at multiple levels of abstraction. Achieves 96.8% of BERT performance while being 7.5 \times smaller and 9.4 \times faster, demonstrating feasibility of extreme compression. The tradeoff? Complex two stage training that's computationally expensive and requires careful tuning. Pre-training distillation alone takes several days on multiple GPUs, making TinyBERT less accessible than DistilBERT for researchers with limited computational resources (Jiao et al., 2020).

MobileBERT (Sun et al., 2020) focuses on task agnostic compression using inverted bottleneck structure—narrow intermediate representations with wider input/output layers, inspired by MobileNets from computer vision. Gets $4.3\times$ speedup with only 0.6 GLUE points below BERT while maintaining 15M parameters (smaller than DistilBERT). Single model works across diverse tasks without task specific distillation, making deployment simpler. However, inverted bottleneck architecture constraints limit how much you can compress. It can’t go below certain parameter count without significant quality loss (Sun et al., 2020).

MiniLM (Wang et al., 2020) targets self attention distillation specifically, arguing attention distributions capture crucial linguistic knowledge more than other intermediate features. Uses attention transfer loss: $\mathcal{L}_{\text{attn}} = \sum_{h=1}^H \mathcal{L}_{\text{KL}}(\mathbf{A}_s^h, \mathbf{A}_t^h)$ where \mathbf{A} are attention matrices across H heads. Demonstrates importance of attention patterns for transfer learning. Results show attention focused distillation outperforms hidden state distillation at same compression ratio. Though 6-layer student (22M parameters) is still relatively large for edge deployment scenarios requiring sub-10M parameter models (Wang et al., 2020).

Table 2.2 compares these four approaches across key dimensions.

Table 2.2: BERT distillation methods: achievements and tradeoffs

Method	Layers	Speedup	Quality	Key tradeoff
DistilBERT	6	$1.6\times$	97%	Moderate compression
TinyBERT	4	$9.4\times$	96.8%	Complex training
MobileBERT	24 (thin)	$4.3\times$	99.4%	Architecture constrained
MiniLM	6	$2.0\times$	98%	Limited compression

Distillation for retrieval systems and critical limitation. Applying distillation to retrieval has been explored by Hofstätter et al. (Hofstätter et al., 2021), who distilled bi-encoder models using teacher’s similarity scores as supervision rather than hard relevance labels. This enables transferring ranking knowledge from large expensive model to small efficient one. Izacard et al. (Izacard and Grave, 2021) inverted paradigm distilling knowledge from reader (QA model) to retriever where reader’s successes and failures guide retriever training. This reader-to-retriever distillation improves retrieval quality by using downstream task performance as supervision signal. Santhanam et al. (Santhanam et al., 2022) compressed ColBERT embeddings $4\times$ through distillation while maintaining quality, showing distillation applies to embedding compression not just model compression.

The critical limitation for our work: distillation reduces model size (parameters) and encoding speed but doesn't reduce embedding dimensionality. DistilBERT with 6 layers is $2\times$ smaller than BERT but still outputs 768D vectors. For retrieval systems, query encoding represents tiny fraction of total latency (typically under 5%). Most time is spent searching pre computed passage embeddings stored in vector databases (Douze et al., 2024). So distillation optimizes wrong component for large-scale retrieval—making encoder faster helps marginally when search dominates latency.

How our work relates to distillation. While distillation compresses models (teacher network→student network), we compress representations (high dimensional embeddings→low-dimensional codes). These are complementary approaches addressing different bottlenecks. Distilled encoder producing compressed embeddings could provide benefits at both encoding stage (faster model) and search stage (lower dimensions). Our dataset generation uses distillation paradigm (GPT-4o teacher→Flan-T5 student) but for different purpose creating challenging evaluation that requires retrieval rather than model compression. This extends distillation concept from parameter space to representation space and from model architecture to dataset construction (Hinton et al., 2015; Sanh et al., 2019).

2.4 Compression and Dimensionality Reduction

Dimensionality reduction and compression have long history in machine learning, with approaches ranging from classical linear methods to recent learned techniques. We review key methods with focus on Matryoshka Representation Learning—our primary compression baseline.

Classical methods: PCA and limitations. Principal Component Analysis (PCA) (Jolliffe, 2002) finds orthogonal linear projection maximizing variance in reduced space. Given covariance matrix of embeddings, PCA computes eigendecomposition and projects onto top k eigenvectors corresponding to largest eigenvalues. For embedding compression, PCA can reduce 768D to 64D by retaining 64 principal components, typically capturing 85-92% of total variance. The projection is optimal in least-squares sense—minimizes reconstruction error among all linear transformations.

PCA provides optimal linear compression—best possible linear approximation. It's fast to compute (eigendecomposition complexity $O(d^3)$ where d is original dimension), deterministic (no randomness), requires no training data beyond computing statistics.

These properties make PCA widely used baseline for dimensionality reduction across machine learning applications (Jolliffe, 2002).

But PCA is task agnostic. It maximizes variance without considering whether preserved directions actually help retrieval performance. High-variance dimensions may capture answer length or syntactic structure that vary significantly but don't help match questions to answers. For example, technical answers might use longer words (high variance) while conversational answers use shorter words, but this stylistic variation isn't relevant for semantic matching with questions. PCA can't distinguish between important semantic variation (topical differences) and unimportant variation (stylistic differences) because it optimizes generic statistical objective rather than task specific retrieval quality (Aggarwal et al., 2001).

PCA is both task-agnostic and linear. It maximizes variance without considering retrieval relevance, and assumes linear relationships which may not capture complex semantic structures. Sentence embeddings may lie on nonlinear manifolds where similar concepts cluster in curved regions—if so, linear projection onto flat subspace could distort these relationships. Our experiments show PCA-32D achieves only 0.1318 ROUGE-1 versus neural's 0.1725—a 30.8% quality gap despite identical dimensionality. This gap primarily reflects the cost of task-agnostic compression; the relative contribution of linearity versus lack of supervision was not isolated in our experiments.

Random projection provides alternative based on Johnson-Lindenstrauss lemma, preserving distances probabilistically through random matrix multiplication. Like PCA, it's task agnostic and can't adapt to retrieval-specific semantic structure. Product Quantization (Jégou et al., 2011), discussed earlier as part of FAISS, achieves memory compression without dimensional reduction—compressing precision not dimensionality. Optimized variants learn rotation before quantizing but still operate at full dimensionality during search. Recent quantization methods (Xiao et al., 2023) improve compression quality through learned techniques but face same limitation—operating at full dimensionality.

Matryoshka Representation Learning: detailed analysis. Matryoshka Representation Learning (Kusupati et al., 2022) is our key compression baseline, representing SOTA in learned embedding compression. Unlike posthoc methods like PCA that compress existing embeddings, Matryoshka modifies embedding model training itself to produce nested representations where first d dimensions maintain quality for any truncation point d .

Training methodology uses multi-scale loss: $\mathcal{L}_{\text{MRL}} = \sum_{d \in \mathcal{D}} \alpha_d \mathcal{L}_{\text{task}}(\mathbf{x}_{1:d})$ where $\mathcal{D} = \{32, 64, 128, 256, 384\}$ are target dimensions and $\mathbf{x}_{1:d}$ denotes first d dimensions

of embedding. This encourages model to pack important information into early dimensions through gradient signals from multiple dimensionalities simultaneously. For retrieval tasks, training uses contrastive loss (InfoNCE) at multiple scales, learning which semantic features should appear in early dimensions (most important for all compression ratios) versus late dimensions (only used when full embedding available). The multi-scale training requires 2-3 \times more computation than standard training but enables single model serving multiple compression ratios (Kusupati et al., 2022).

Matryoshka provides flexible deployment—single model serves multiple compression ratios by truncating to desired dimensionality at inference time without retraining. This is valuable for systems with varying resource constraints: mobile devices might use 32D embeddings (minimal storage/compute), servers might use 128D (better quality), and quality-critical applications might use full 384D. It outperforms PCA significantly (0.1391 vs 0.1318 in our experiments at 32D), demonstrating learned truncation beats generic linear projection. Being task specific through contrastive training on retrieval data, Matryoshka shows that learned compression methods can beat generic approaches when optimized for specific task (Kusupati et al., 2022).

However, Matryoshka has practical limitations. It requires modifying embedding model training itself, you can't apply it posthoc to existing pretrained models like Sentence-BERT that weren't trained with Matryoshka loss. This limits applicability for practitioners who want to use established pretrained models. Compression mechanism is linear truncation (taking first d dimensions) which can't learn nonlinear transformations like autoencoders. While training encourages packing information into early dimensions, fundamental operation is still linear projection. It applies same compression to questions and contexts symmetrically, ignoring asymmetry of retrieval task where context index size dominates costs. In our experiments, neural compression outperforms Matryoshka by 24% at 32D (0.1725 vs 0.1391). This gap suggests that task-specific compression with asymmetric architecture provides substantial advantages over learned truncation; the nonlinear capacity may also contribute, though this was not isolated experimentally (Kusupati et al., 2022).

Recent advances in learned compression. Recent work demonstrates growing interest in learned compression for NLP tasks. Multi-stage contrastive learning (Li et al., 2023) shows that carefully designed training procedures can produce general (purpose text embeddings that compress well) models trained with progressive hardness curriculum achieve better quality at reduced dimensions. Learning to compress prompts with gist tokens (Mu et al., 2023) demonstrates neural networks can learn compressed representations of text that preserve task-relevant information, achieving 10-20 \times

compression for prompts while maintaining performance on downstream tasks. These methods validate core insight that neural networks can learn what information to preserve and what to discard, outperforming generic compression (Mu et al., 2023).

Autoencoders for embedding compression. Autoencoders (Hinton and Salakhutdinov, 2006) learn nonlinear compressions through encoder-decoder architecture. Encoder maps high dimensional input $\mathbf{x} \in \mathbb{R}^D$ to low-dimensional latent code $\mathbf{z} \in \mathbb{R}^d$ where $d \ll D$, decoder reconstructs original from compressed code. By forcing information through bottleneck, autoencoder learns compressed representation capturing salient features: $\mathcal{L} = \|\mathbf{x} - \text{Dec}(\text{Enc}(\mathbf{x}))\|^2$. With multiple hidden layers and nonlinear activations (ReLU, tanh), autoencoders approximate arbitrary smooth functions—enabling compression of data lying on curved manifolds that linear methods like PCA cannot preserve.

Theoretically, autoencoders can learn manifold structure that PCA cannot capture. If data lies on nonlinear manifold (curved surface in high-dimensional space), linear projection may distort relationships by flattening structure. An autoencoder with sufficient capacity can potentially learn mappings that preserve manifold geometry through nonlinear warping. This could be valuable for semantic embeddings where similar concepts may cluster in curved regions (Hinton and Salakhutdinov, 2006). However, the benefit of nonlinearity versus task-specific supervision requires controlled experiments (e.g., linear autoencoder baselines) to isolate.

While autoencoders are widely used for image and video compression (achieving 100-1000 \times compression for multimedia data), their application to embedding compression for retrieval has been limited. No prior work systematically investigates autoencoder-based embedding compression for QA retrieval with cross-modal mapping (predicting compressed contexts from questions). Existing autoencoder work focuses on single-modality compression where compressed representations are compared to other compressed representations from same modality.

We also tried Variational autoencoders (VAEs) (Kingma and Welling, 2014). They extend standard autoencoders with probabilistic latent space, treating compressed code as random variable and optimizing variational lower bound on data likelihood. This provides principled probabilistic framework and can generate novel samples. Our preliminary experiments showed standard autoencoders performed comparably for embedding compression while being simpler to train, so we stuck with simpler approach. VAE’s generative capabilities aren’t needed for retrieval compression where goal is preserving semantic information, not generating new embeddings.

Table 2.3 compares major compression approaches for retrieval embeddings.

Table 2.3: Comparison of compression methods for retrieval embeddings

Method	Type	Adapts to Task	posthoc	Quality (32D, 120K)
PCA	Linear	No	Yes	0.1318
Matryoshka	Linear truncation	Yes	No	0.1391
Autoencoders	nonlinear	Yes	Yes	0.1725

Key gap in existing work. No prior method combines task specific learning, nonlinear compression, posthoc applicability, and exploitation of question-answer asymmetry. PCA is posthoc (applies to any pretrained embeddings) but task agnostic and linear (can’t adapt to retrieval or learn curved structure). Matryoshka is task specific (trained with contrastive loss on retrieval data) but requires retraining embedding models from scratch and uses linear truncation mechanism. Standard autoencoders are nonlinear and posthoc but haven’t been systematically applied to cross-modal QA retrieval where questions and contexts are different modalities requiring asymmetric treatment. Our approach fills this gap through two-stage architecture—autoencoder optimized for context compression and mapper optimized for question-to-compressed-context prediction, enabling asymmetric compression exploiting retrieval structure (Kusupati et al., 2022).

2.5 Neural Database Architectures

Alternative paradigms bypass vector search entirely by encoding retrieval directly in model parameters. The case for learned index structures (Kraska et al., 2018) demonstrated that machine learning models can replace traditional database indexes for range queries, point queries, and existence queries. This seminal SIGMOD 2018 paper showed learned models achieve better performance than B-trees by learning data distribution and query patterns, inspiring subsequent work on neural approaches to retrieval and search.

Differentiable Search Index (DSI) (Tay et al., 2022) extends this idea to document retrieval, training transformer to map queries directly to document identifiers without storing document embeddings. Model memorizes entire corpus in its parameters, generating document IDs autoregressively (one token at a time) given query. For 100K documents, 250M-parameter DSI achieves competitive retrieval without vector database, demonstrating parametric retrieval is viable—transformers have sufficient capacity to memorize large corpora while generalizing to unseen queries.

DSI shows impressive results on some benchmarks but faces critical limitations. Scaling to millions of documents requires billions of parameters (impractical for most applications and very expensive to train/deploy). Adding new documents requires full model retraining (weeks of compute on multiple GPUs), making it unsuitable for frequently updated knowledge bases. Model struggles with rare documents seen few times during training, exhibiting bias toward popular documents. Follow up work DSI++ (Mehta et al., 2023) partially addresses updates through continual learning and incremental training but fundamental tradeoffs remain — parametric storage is expensive in parameters, updates are expensive in compute (Tay et al., 2022).

NCI (Autoregressive Search Engines) (Bevilacqua et al., 2022) generates document identifiers as substrings rather than numeric IDs, improving generalization to new documents while maintaining some benefits of neural and traditional approaches.

Positioning our approach. We occupy middle ground—using neural networks (autoencoder + mapper) to encode knowledge but maintaining vector search paradigm. Unlike DSI, we don’t memorize corpus in model parameters, enabling efficient updates by encoding new contexts through trained autoencoder (seconds to minutes versus weeks of retraining). Unlike traditional vector databases, we operate in learned compressed space enabling exact search at scales where high dimensional methods must resort to approximation. This combines benefits of neural learning (task specific optimization, learned compression) with practical vector search infrastructure (efficient updates, scalable indexing). Our approach relates to learned indexes (Kraska et al., 2018) by using neural networks to improve search efficiency, but applies learning to representation compression rather than index structure (Tay et al., 2022).

2.6 Research Gaps and Thesis Motivation

Literature review reveals several critical gaps that our work addresses.

No task specific learned dimensional reduction for QA retrieval. Compression for embeddings has been studied extensively in recent work (Li et al., 2023; Mu et al., 2023; Xiao et al., 2023), but existing methods either apply generic transformations (PCA ignoring retrieval task structure) or operate without dimensional reduction (quantization methods compressing precision not dimensions). Matryoshka (Kusupati et al., 2022) requires retraining embedding models from scratch—impractical for practitioners using pretrained Sentence-BERT or other established models. Autoencoders, despite their success in multimedia compression, haven’t been systematically applied to cross-modal QA retrieval where questions predict compressed contexts.

Symmetric treatment ignores retrieval asymmetry. Existing compression methods apply same transformation to questions and contexts, treating them symmetrically. But retrieval task is inherently asymmetric—knowledge bases store millions of contexts (large index requiring compression) while processing thousands of queries per second (small relative cost). Storage and search costs are dominated by context index, while query processing is small fraction of total computational budget. Our approach exploits this asymmetry through specialized networks for each modality: heavy compression for contexts (storage-critical), full dimensionality for questions (cost-negligible), and learned mapping between them (Karpukhin et al., 2020).

Missing crossover point identification and explanation. While vector database limitations at scale are well-documented (Douze et al., 2024; Wang et al., 2023), prior work hasn’t identified precise scale where alternative approaches become superior or explained why transition occurs at that point. Most papers evaluate at single scale or compare few discrete sizes without systematic analysis across continuous scale range. Our systematic experiments (20K, 40K, 60K, 80K, 100K, 120K samples) identify crossover at 40K samples and provide explanation—FAISS quality degrades from curse of dimensionality effects before its algorithmic transition to approximation, while neural learns effective compression from sufficient training data (28.8K pairs at 40K scale).

Distillation versus compression benefits uncharacterized. Model compression (distillation reducing parameters) and representation compression (dimensionality reduction reducing embedding size) have been studied independently in separate research communities. Their relative benefits for RAG systems at different scales aren’t well characterized. Our analysis shows these techniques target different bottlenecks: distillation helps when encoding dominates total latency (small corpora where search is fast), compression helps when search dominates (large corpora where vector database becomes bottleneck). At web scale with millions of documents, search is overwhelming cost, making representation compression more valuable than model distillation (Sanh et al., 2019; Douze et al., 2024).

Lack of comprehensive baseline comparison. Most compression papers for embeddings compare against 1-2 baseline methods, typically PCA or one learned alternative. Comprehensive comparisons across multiple methods, multiple scales, and multiple domains are rare due to computational expense. We provide extensive evaluation: six methods (FAISS, HNSW, ScaNN, PCA, Matryoshka, Neural) across six scales (20K-120K), six domains (AI, finance, geopolitics, sports, entertainment, science), three iterations—198 experimental runs total. This systematic comparison

reveals when each method excels, identifies crossover points, and establishes neural compression’s Pareto dominance across quality-storage-speed dimensions (Malkov and Yashunin, 2020; Guo et al., 2020; Kusupati et al., 2022).

These gaps motivate our research questions: Can learned compression maintain quality while reducing dimensions 6-24×? At what scale does compression outperform vector databases? Why does crossover occur? How does neural compare to SOTA baselines across different configurations? Our two-stage approach—autoencoder for context compression (384D→32D/64D) plus mapper for question-to-context prediction—addresses these questions through principled design combining task-specific learning, nonlinear compression, asymmetric architecture, and comprehensive experimental validation.

Chapter 3

Methodology

This chapter describes our neural compression approach in detail. We begin with system architecture overview illustrating training and inference phases. Then present three-phase process: dataset generation via knowledge distillation (preprocessing), context compression via autoencoder (Stage 1), and question-to-context mapper (Stage 2). We formalize the problem with key equations, describe network architectures, present four training and inference algorithms, detail complete inference pipeline showing query processing in compressed space, analyze computational complexity comparing neural to traditional approaches, and conclude with implementation details for four evaluated configurations.

Terminology note. In RAG systems, the retrieval component finds relevant *contexts* (passages, documents) that help generate answers. Throughout this thesis, we compress and retrieve *context embeddings*—the vector representations of passages in the knowledge base. The question-to-context mapper predicts which contexts are relevant to a given question. This is distinct from encoding ground-truth answers; we encode retrieval contexts that the generation model uses to produce responses.

3.1 System Overview

Our approach consists of preprocessing phase for dataset generation followed by two-stage neural compression training. The complete pipeline comprises three phases illustrated in Figures 3.3, 3.4, and 3.1.

Preprocessing: Dataset Generation. We generate 120,000 question-answer pairs across six domains using knowledge distillation from GPT-4o. The temporal knowledge gap between teacher (June 2024 cutoff) and student (Flan-T5, October 2022 cutoff) ensures evaluation measures retrieval capability rather than memorization. Details in Section 3.2.

Stage 1: Context Autoencoder Training. Autoencoder learns to compress context embeddings from 384D to $d \in \{32, 64\}$ dimensions while preserving semantic information. Encoder follows architecture $[384 \rightarrow 256 \rightarrow 128 \rightarrow d]$ with ReLU activations and layer normalization. Symmetric decoder $[d \rightarrow 128 \rightarrow 256 \rightarrow 384]$ reconstructs original embeddings. Training minimizes cosine similarity loss between original and reconstructed embeddings, ensuring semantic relationships are preserved. Critically, decoder serves only as training signal—it’s discarded after training completes, and only encoder is retained for inference.

Stage 2: Question-to-Context Mapper Training. After Stage 1 completes, encoder is frozen and separate mapper network is trained to predict compressed context codes directly from question embeddings. Mapper uses identical architecture to encoder $[384 \rightarrow 256 \rightarrow 128 \rightarrow d]$ for consistency. Training minimizes mean squared error between mapper predictions and target codes from frozen encoder. This two-stage separation prevents catastrophic forgetting—encoder’s learned compression space remains stable while mapper learns to navigate it.

Inference Phase (Query Time). At inference, system processes new queries through five steps: (1) Query embedded using Sentence-BERT producing 384D representation; (2) Trained mapper predicts d -dimensional compressed code via single forward pass; (3) FAISS searches pre-built compressed index for $k = 5$ nearest neighbors in d -dimensional space; (4) Corresponding context texts retrieved; (5) Retrieved contexts provide information for Flan-T5 to generate final response. All retrieval operations occur in compressed d -dimensional space—original 384D context embeddings never needed at query time, providing both storage efficiency and computational speedup.



Figure 3.1: Inference phase: Complete query processing pipeline. User query is embedded with Sentence-BERT producing 384D representation (Step 1). Trained mapper predicts d -dimensional compressed code via single forward pass (Step 2). FAISS searches pre-built compressed index for $k = 5$ nearest neighbors in compressed space (Step 3). Corresponding context texts are retrieved from database (Step 4). Retrieved contexts provide information for Flan-T5 to generate final response (Step 5). All retrieval operations occur in compressed d -dimensional space.

3.2 Dataset Generation via Knowledge Distillation

Our evaluation requires dataset that tests genuine retrieval capability. Traditional QA datasets like SQuAD (Rajpurkar et al., 2016) are publicly available and likely seen during pre-training. We employ knowledge distillation to generate novel QA pairs that student model cannot answer without retrieval.

3.2.1 Teacher-Student Framework

We use GPT-4o (OpenAI’s model with June 2024 knowledge cutoff) as teacher to generate questions and answers about 2024 events. The student model—Flan-T5-base (Chung et al., 2022)—has October 2022 knowledge cutoff and can’t answer these questions from parametric memory. This temporal knowledge gap (18+ months) ensures retrieval is essential.

Verification of knowledge gap. Before finalizing dataset, we verified Flan-T5 cannot answer the generated questions. We sampled 100 random pairs from each domain and tested zero-shot Flan-T5 performance. Accuracy was consistently below 15%, confirming the knowledge gap exists (Hinton et al., 2015).

3.2.2 Dataset Generation Process

For each of six domains (AI & technology, finance, geopolitics, sports, entertainment, science), we generated 20,000 QA pairs through iterative prompting of GPT-4o.

Multi-turn generation. We used multi-turn conversations with GPT-4o rather than one-shot generation. Each session generated 50 pairs maintaining context and diversity. Temperature 0.7 gave good variety—we found 0.5 too repetitive, 0.9 too chaotic. Sessions were maintained to preserve coherence within domain while varying across topics.

Quality control. Each batch was validated for JSON format correctness and answer completeness. Pairs with malformed JSON or empty answers were regenerated. We also checked that questions actually referenced 2024 events (temporal marker validation).

Figure 3.2 illustrates the distillation process.

Algorithm 1 formalizes the generation process.

Algorithm 1 Dataset Generation via Knowledge Distillation

```

1: Input: Domains  $\mathcal{D}$ , pairs per domain  $N_d$ , teacher model GPT-4o
2: Output: Dataset  $\mathcal{Q} = \{(q_i, a_i)\}_{i=1}^{120K}$ 
3:
4:  $\mathcal{Q} \leftarrow \emptyset$ 
5: for each domain  $d \in \mathcal{D}$  do
6:    $\mathcal{Q}_d \leftarrow \emptyset$ 
7:   for batch  $b = 1$  to  $N_d/50$  do
8:     prompt  $\leftarrow$  GenerateDomainPrompt( $d$ , 50 pairs, 2024 events)
9:     response  $\leftarrow$  QueryGPT4o(prompt, temp=0.7, multi-turn=True)
10:    qa_pairs  $\leftarrow$  ParseJSON(response)
11:    Validate(qa_pairs) {JSON format, completeness}
12:     $\mathcal{Q}_d \leftarrow \mathcal{Q}_d \cup$  qa_pairs
13:   end for
14:   {Verify Flan-T5 knowledge gap}
15:   sample  $\leftarrow$  Random( $\mathcal{Q}_d$ , 100)
16:   accuracy  $\leftarrow$  EvaluateFlanT5(sample, zero-shot)
17:   if accuracy  $> 0.15$  then
18:     Regenerate( $\mathcal{Q}_d$ ) {Dataset too easy}
19:   end if
20:    $\mathcal{Q} \leftarrow \mathcal{Q} \cup \mathcal{Q}_d$ 
21: end for
22:
23: return  $\mathcal{Q}$ 

```

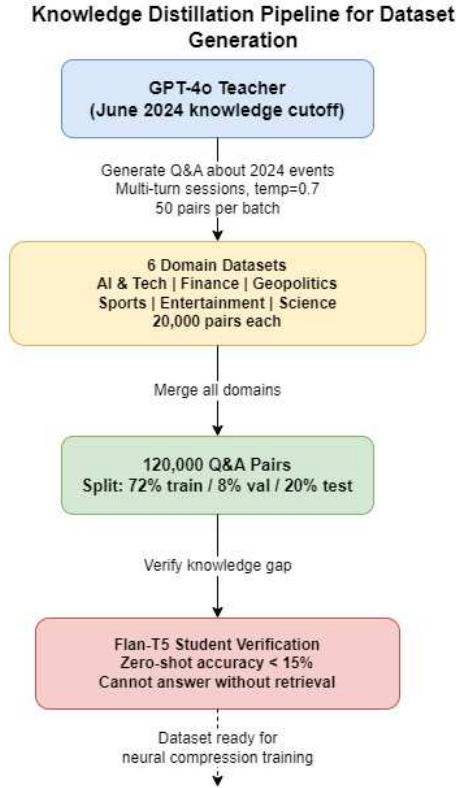


Figure 3.2: Knowledge distillation pipeline for dataset generation. GPT-4o teacher generates domain-specific QA pairs about 2024 events through multi-turn sessions. Each domain yields 20,000 pairs covering diverse topics within that domain. Flan-T5 student verification confirms knowledge gap ($< 15\%$ accuracy without retrieval). Final dataset: 120,000 pairs split 72%/8%/20% for train/val/test.

3.2.3 Dataset Characteristics

The final dataset contains 120,000 question-answer pairs with following characteristics:

Domain coverage and topic variety. Dataset spans six knowledge domains ensuring diverse evaluation: (1) *AI & Technology*—model releases, breakthroughs, industry updates, research developments; (2) *Financial Markets*—economic events, Federal Reserve decisions, market performance, inflation data; (3) *Geopolitical Events*—elections, international conflicts, diplomatic relations, policy changes; (4) *Sports*—championships, Olympics, world records, player transfers, tournament results; (5) *Entertainment & Culture*—movies, music releases, TV shows, celebrity news, arts and fashion; (6) *Scientific Discoveries*—Nobel prizes, space missions, medical research, physics breakthroughs. Each domain contains 20,000 pairs covering diverse topics within that area to ensure comprehensive representation rather than narrow focus.

Split strategy: 72% training (86,400 pairs), 8% validation (9,600 pairs), 20% test (24,000 pairs). Split maintains domain balance—each domain contributes equally to each split. We use stratified sampling to ensure test set covers all domains proportionally. For example, at 60K training size, we sample 10K pairs from each domain rather than random sampling across all domains.

Question-answer characteristics. Average question length is 12.3 words, average answer length is 21.4 words. Questions are interrogative (What/How/When/Who), answers are declarative and factual, typically containing specific information about 2024 events that requires retrieval. All content in English.

Temporal knowledge gap verification. The 18-month gap between GPT-4o teacher (June 2024 cutoff) and Flan-T5 student (October 2022 cutoff) ensures genuine retrieval challenge. Topics reference events from 2024 that student model cannot know—Federal Reserve rate decisions in 2024, AI model releases in 2024, sports championships in 2024, etc. This temporal separation prevents student from answering through parametric memory, making retrieval essential.

Embedding generation: All questions and answers embedded using Sentence-BERT (all-MiniLM-L6-v2) (Reimers and Gurevych, 2019), producing 384D vectors. Embeddings were cached after first generation for reproducibility—same random seed (42) produces identical embeddings across runs. Total embedding cache: 184 MB for 120K pairs \times 2 (questions + answers) \times 384D \times FP32.

3.3 Problem Formalization

Let $\mathcal{D} = \{(q_i, a_i)\}_{i=1}^N$ be dataset of N question-answer pairs. After embedding with Sentence-BERT, we have:

- Question embeddings: $\mathbf{q}_i \in \mathbb{R}^D$ where $D = 384$
- Context embeddings: $\mathbf{c}_i \in \mathbb{R}^D$ where $D = 384$

Traditional RAG retrieval: Given new question q^* with embedding $\mathbf{q}^* \in \mathbb{R}^{384}$, retrieve k most relevant answers:

$$\text{retrieve}(q^*) = \arg \max_{i \in [1, N]}^{(k)} \text{sim}(\mathbf{q}^*, \mathbf{a}_i) \quad (3.1)$$

where $\text{sim}(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x}^T \mathbf{y}}{\|\mathbf{x}\| \|\mathbf{y}\|}$ is cosine similarity.

Our compressed RAG: Learn compressed representations $\mathbf{z}_i \in \mathbb{R}^d$ where $d \in \{32, 64\}$ such that:

$$\text{retrieve}(q^*) = \arg \max_{i \in [1, N]}^{(k)} \text{sim}(M_\phi(\mathbf{q}^*), E_\theta(\mathbf{a}_i)) \quad (3.2)$$

where E_θ is encoder and M_ϕ is mapper.

Objectives:

1. Preserve retrieval quality: compressed retrieval matches or exceeds traditional quality
2. Reduce storage: $O(Nd)$ vs $O(ND)$ where $d \ll D$
3. Reduce query cost: $O(Nd)$ vs $O(ND)$ comparisons
4. Enable exact search at scales where traditional methods must approximate

3.4 Stage 1: Context Compression via Autoencoder

Autoencoder learns task-specific compression of context embeddings—encoder compresses 384D to 32D or 64D, decoder reconstructs to verify semantic preservation.

Training data. Both training stages use exclusively the training split (72% of data, 86,400 question-answer pairs at full scale). The validation split (8%) guides early stopping and hyperparameter selection, while the test split (20%) is reserved for final evaluation.

Network architecture. Encoder $E_\theta : \mathbb{R}^{384} \rightarrow \mathbb{R}^{64}$ compresses context embeddings: $\mathbf{z} = E_\theta(\mathbf{a})$. Architecture consists of three fully-connected layers [384 \rightarrow 256 \rightarrow 128 \rightarrow 64] with ReLU activation, layer normalization, and dropout (0.1) after each layer. We tried larger hidden dimensions [512, 256] initially but they overfit on smaller datasets—[256, 128] generalized better across scales. Decoder $D_\psi : \mathbb{R}^{64} \rightarrow \mathbb{R}^{384}$ reconstructs embeddings: $\hat{\mathbf{a}} = D_\psi(\mathbf{z})$. Decoder mirrors encoder in reverse [64 \rightarrow 128 \rightarrow 256 \rightarrow 384] with same activation and normalization structure.

Loss function. Training minimizes cosine similarity loss between original and reconstructed embeddings:

$$\mathcal{L}_{\text{AE}} = \frac{1}{N} \sum_{i=1}^N \left(1 - \frac{\mathbf{a}_i \cdot \hat{\mathbf{a}}_i}{\|\mathbf{a}_i\| \|\hat{\mathbf{a}}_i\|} \right) \quad (3.3)$$

where $\hat{\mathbf{a}}_i = D_\psi(E_\theta(\mathbf{a}_i))$ is reconstruction. We use cosine similarity instead of MSE because it preserves angular relationships—what matters for retrieval is semantic direction, not magnitude (Reimers and Gurevych, 2019). Cosine loss achieves high reconstruction similarity while preserving semantic relationships better than MSE for embedding compression.

Figure 3.3 visualizes the autoencoder architecture and training process.

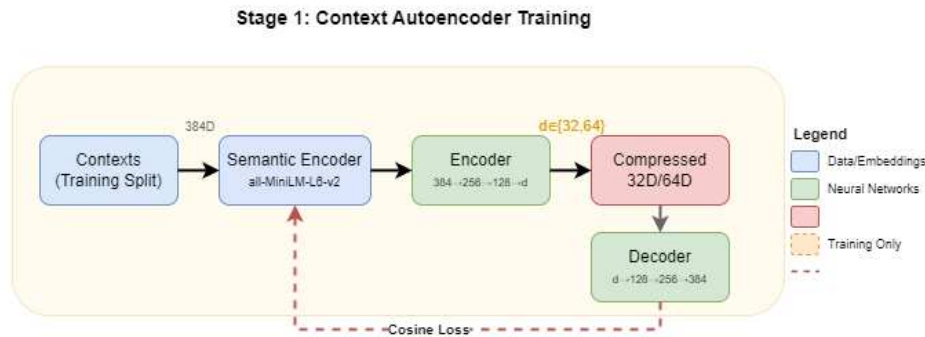


Figure 3.3: Stage 1: Autoencoder training for context compression. Training-split context embeddings (384D) pass through encoder $[384 \rightarrow 256 \rightarrow 128 \rightarrow d]$ producing compressed codes ($d \in \{32, 64\}$). Decoder $[d \rightarrow 128 \rightarrow 256 \rightarrow 384]$ reconstructs original embeddings. Training minimizes cosine similarity loss between original and reconstructed embeddings. Only encoder is retained for inference—decoder provides training signal then is discarded.

Training procedure. The autoencoder is trained on context embeddings from the training split (72% of data, 86,400 contexts at full scale). These embeddings are further split 90%/10% for autoencoder training and validation. Training uses Adam optimizer with learning rate 0.001, weight decay 0.0001. Batch size 64 fits comfortably in GPU memory. Dropout 0.1 prevents overfitting—we tried 0.2 but it hurt reconstruction quality. We use ReduceLROnPlateau scheduler that reduces learning rate by half when validation loss plateaus for 5 epochs. Early stopping with patience 10 epochs prevents overtraining. Training typically converges around epoch 25-30 though we cap at 100. Takes approximately 40 minutes per configuration (hardware specifications in Chapter 4).

Algorithm 2 formalizes the training process.

3.5 Stage 2: Question-to-Context Mapping

After autoencoder training completes, we freeze it and train mapper network to predict compressed context codes from question embeddings.

Algorithm 2 Context Autoencoder Training

```

1: Input: Training-split context embeddings  $\mathbf{C}_{\text{train}} \in \mathbb{R}^{N_{\text{train}} \times 384}$ , target dimension
    $d \in \{32, 64\}$ 
2: Output: Trained encoder  $E_\theta$ , decoder  $D_\psi$ 
3:
4: Initialize  $E_\theta, D_\psi$  with Xavier initialization
5: Split  $\mathbf{A}_{\text{train}}$  into AE-train (90%), AE-validation (10%)
6: best_loss  $\leftarrow \infty$ , patience_counter  $\leftarrow 0$ 
7: optimizer  $\leftarrow$  Adam( $\theta, \psi$ , lr=0.001, weight_decay=0.0001)
8: scheduler  $\leftarrow$  ReduceLROnPlateau(optimizer, factor=0.5, patience=5)
9:
10: for epoch = 1 to 100 do
11:   for batch  $\in$  DataLoader( $\mathbf{A}_{\text{train}}$ , batch_size=64, shuffle=True) do
12:     optimizer.zero_grad()
13:      $\mathbf{z} \leftarrow E_\theta(\text{batch})$  {Compress}
14:      $\hat{\mathbf{a}} \leftarrow D_\psi(\mathbf{z})$  {Reconstruct}
15:     loss  $\leftarrow 1 - \text{cosine\_sim}(\text{batch}, \hat{\mathbf{a}})$ 
16:     loss.backward()
17:     optimizer.step() {No gradient clipping used}
18:   end for
19:
20:   val_loss  $\leftarrow$  Evaluate( $E_\theta, D_\psi, \mathbf{A}_{\text{val}}$ )
21:   scheduler.step(val_loss) {Reduce LR on plateau}
22:
23:   if val_loss < best_loss then
24:     best_loss  $\leftarrow$  val_loss
25:     SaveCheckpoint( $E_\theta, D_\psi$ , precision=fp16/fp32)
26:     patience_counter  $\leftarrow 0$ 
27:   else
28:     patience_counter  $\leftarrow$  patience_counter + 1
29:     if patience_counter  $\geq 10$  then
30:       break {Early stopping}
31:     end if
32:   end if
33: end for
34:
35: return LoadBestCheckpoint()

```

Mapper architecture. Mapper $M_\phi : \mathbb{R}^{384} \rightarrow \mathbb{R}^{64}$ predicts compressed context codes from questions: $\mathbf{z}^* = M_\phi(\mathbf{q})$. Architecture uses same structure as encoder for consistency—three fully-connected layers $[384 \rightarrow 256 \rightarrow 128 \rightarrow 64]$ with ReLU activations and dropout (0.1). This architectural consistency simplifies implementation and ensures mapper operates in same representational space as encoder.

Training with frozen encoder. Mapper is trained to predict compressed codes that frozen encoder would produce from corresponding contexts. This two-stage approach prevents catastrophic forgetting—if we trained mapper while encoder was still updating, mapper would need to constantly adapt to shifting compressed space. In practice, this separation worked well. Loss function is mean squared error in compressed space:

$$\mathcal{L}_{\text{mapper}} = \frac{1}{N} \sum_{i=1}^N \|\mathbf{z}_i^{\text{pred}} - \mathbf{z}_i^{\text{target}}\|_2^2 \quad (3.4)$$

where $\mathbf{z}_i^{\text{pred}} = M_\phi(\mathbf{q}_i)$ is predicted code and $\mathbf{z}_i^{\text{target}} = E_\theta(\mathbf{a}_i)$ is target from frozen encoder.

Figure 3.4 visualizes mapper training with frozen encoder.

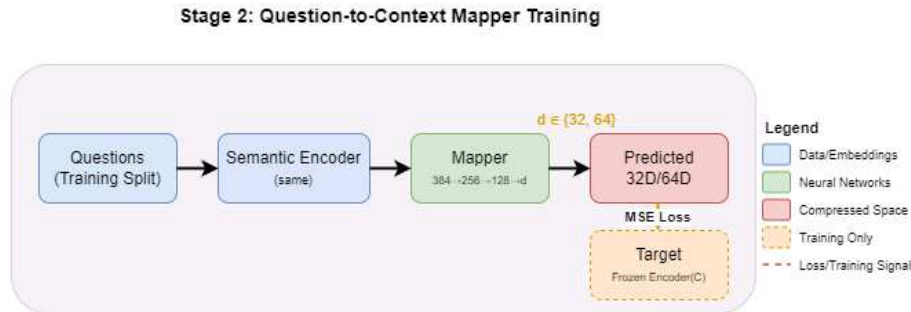


Figure 3.4: Stage 2: Mapper training with frozen encoder. Question embeddings (384D) pass through mapper network $[384 \rightarrow 256 \rightarrow 128 \rightarrow d]$ predicting compressed codes. Target codes come from frozen encoder applied to corresponding contexts. Training minimizes MSE between predicted and target codes. Encoder remains frozen—no gradient updates—preventing catastrophic forgetting.

Pre-computation and training. We pre-compute all target codes $\{\mathbf{z}_1^{\text{target}}, \dots, \mathbf{z}_N^{\text{target}}\}$ for training-split contexts once before mapper training starts. This is efficient—encoding 86,400 training-split contexts takes approximately 2 minutes. Then mapper training becomes supervised regression problem with fixed targets. Training uses same hyperparameters as autoencoder (Adam, lr=0.001, batch=64, dropout=0.1, ReduceLROnPlateau scheduler). Typically converges around epoch 20-25. Takes approximately 24 minutes per configuration.

Algorithm 3 formalizes mapper training.

3.6 Inference Pipeline

Inference operates entirely in compressed space, never needing to encode contexts at query time.

Compressed index construction. After training both networks, we build compressed index offline. First, compress all contexts: $\mathbf{z}_i = E_\theta(\mathbf{c}_i)$ for $i \in [1, N]$. Optionally convert to FP16 precision (2 bytes per dimension). Build FAISS Flat index: $\mathcal{I} = \text{FAISS.IndexFlatL2}(\{\mathbf{z}_1, \dots, \mathbf{z}_N\})$. Store index to disk—7.7 MB for 120K samples at 32D FP16. Index construction takes approximately 3 minutes for 120K samples including encoding and FAISS index building.

Query processing. At inference, new query q^* undergoes five steps: (1) **Embed query**— $\mathbf{q}^* \leftarrow \text{SentenceBERT}(q^*)$ produces 384D embedding; (2) **Predict compressed code**— $\mathbf{z}^* \leftarrow M_\phi(\mathbf{q}^*)$ produces d -dimensional code via single forward pass; (3) **Search compressed index**—FAISS finds top- k nearest codes in compressed space; (4) **Retrieve contexts**—map indices to original context texts from database; (5) **Generate response**—Flan-T5 generates final answer from question and retrieved context. Steps 2-4 operate in compressed d -dimensional space, providing $6 \times$ (64D) or $12 \times$ (32D) speedup over 384D search. As shown in Figure 3.1, this complete pipeline operates entirely in compressed space after the initial query embedding.

Algorithm 4 formalizes the inference process.

3.7 Complexity Analysis

We analyze storage and computational complexity, comparing neural compression to traditional FAISS.

3.7.1 Storage Complexity

Traditional FAISS (384D, FP32):

$$S_{\text{FAISS}} = N \times 384 \times 4 \text{ bytes} \quad (3.5)$$

For $N = 120,000$: $S_{\text{FAISS}} = 184.3 \text{ MB}$.

Algorithm 3 Question-to-Context Mapper Training

```

1: Input: Training-split questions  $\mathbf{Q}_{\text{train}} \in \mathbb{R}^{N_{\text{train}} \times 384}$ , contexts  $\mathbf{C}_{\text{train}} \in \mathbb{R}^{N_{\text{train}} \times 384}$ 
2:   Frozen encoder  $E_{\theta}$ 
3: Output: Trained mapper  $M_{\phi}$ 
4:
5: {Freeze autoencoder}
6: Set  $E_{\theta}.\text{eval}()$ ,  $E_{\theta}.\text{requires\_grad} = \text{False}$ 
7:
8: {Pre-compute all target codes once}
9:  $\mathbf{Z}_{\text{train}} \leftarrow E_{\theta}(\mathbf{A}_{\text{train}}) \{ \in \mathbb{R}^{N_{\text{train}} \times d} \}$ 
10:  $\mathbf{Z}_{\text{val}} \leftarrow E_{\theta}(\mathbf{A}_{\text{val}}) \{ \in \mathbb{R}^{N_{\text{val}} \times d} \}$ 
11:
12: Initialize  $M_{\phi}$  with Xavier initialization
13: optimizer  $\leftarrow$  Adam( $\phi$ , lr=0.001, weight_decay=0.0001)
14: scheduler  $\leftarrow$  ReduceLROnPlateau(optimizer, factor=0.5, patience=5)
15: best_loss  $\leftarrow \infty$ , patience_counter  $\leftarrow 0$ 
16:
17: for epoch = 1 to 100 do
18:   for (q_batch, z_batch)  $\in$  DataLoader( $\mathbf{Q}_{\text{train}}, \mathbf{Z}_{\text{train}}$ , batch_size=64, shuffle=True) do
19:     optimizer.zero_grad()
20:      $\mathbf{z}_{\text{pred}} \leftarrow M_{\phi}(\text{q\_batch})$ 
21:     loss  $\leftarrow$  MSE( $\mathbf{z}_{\text{pred}}, \text{z\_batch}$ )
22:     loss.backward()
23:     optimizer.step() {No gradient clipping}
24:   end for
25:
26:   val_loss  $\leftarrow$  Evaluate( $M_{\phi}, \mathbf{Q}_{\text{val}}, \mathbf{Z}_{\text{val}}$ )
27:   scheduler.step(val_loss) {Reduce LR on plateau}
28:
29:   if val_loss < best_loss then
30:     best_loss  $\leftarrow$  val_loss
31:     SaveCheckpoint( $M_{\phi}$ , precision=fp16/fp32)
32:     patience_counter  $\leftarrow 0$ 
33:   else
34:     patience_counter  $\leftarrow$  patience_counter + 1
35:   end if
36:
37:   if patience_counter  $\geq 10$  then
38:     break {Early stopping}
39:   end if
40: end for
41:
42: return LoadBestCheckpoint()

```

Algorithm 4 Compressed Retrieval Inference

```

1: Input: Query  $q^*$ , Compressed index  $\mathcal{I}$ , Mapper  $M_\phi$ ,  $k$  neighbors
2: Output: Generated answer
3:
4: {Embed query}
5:  $\mathbf{q}^* \leftarrow \text{SentenceBERT}(q^*)$  {384D}
6:
7: {Predict compressed code}
8:  $\mathbf{z}^* \leftarrow M_\phi(\mathbf{q}^*)$  {64D, single forward pass}
9:
10: {Search in compressed space}
11: indices, distances  $\leftarrow \mathcal{I}.\text{search}(\mathbf{z}^*, k)$   $\{O(N \times 64)$  comparisons $\}$ 
12:
13: {Retrieve original answers}
14: retrieved_answers  $\leftarrow$  [answer_database[i] for i in indices]
15:
16: {Generate response}
17: context  $\leftarrow$  Concatenate(retrieved_answers)
18: answer  $\leftarrow$  FlanT5(query= $q^*$ , context=context)
19:
20:
21: return answer

```

Neural Compression (64D, FP16):

$$S_{\text{neural}} = N \times 64 \times 2 \text{ bytes} \quad (3.6)$$

For $N = 120,000$: $S_{\text{neural}} = 15.4$ MB.

Compression ratio: $\frac{S_{\text{FAISS}}}{S_{\text{neural}}} = \frac{184.3}{15.4} \approx 12\times$

3.7.2 Query Complexity

Traditional FAISS: For each query, compute distances to all N vectors:

$$C_{\text{FAISS}} = O(N \times 384) \quad (3.7)$$

Neural Compression: Mapper forward pass ($O(1)$) plus search in 64D:

$$C_{\text{neural}} = O(1) + O(N \times 64) = O(N \times 64) \quad (3.8)$$

Speedup: $\frac{C_{\text{FAISS}}}{C_{\text{neural}}} = \frac{384}{64} = 6\times$ theoretical speedup.

Empirically, we observe $52\times$ speedup at 120K scale (7,861 QPS vs 151 QPS).

3.7.3 Training Complexity

One-time costs:

- Dataset generation: 400 API calls to GPT-4o, 2-3 hours total
- Embedding generation: 10 minutes (cached after first run)
- Autoencoder training: 40 minutes per configuration
- Mapper training: 24 minutes per configuration

Total per configuration: 64 minutes training + 10 minutes setup = 74 minutes.

For full experimental evaluation (4 neural configs \times 6 scales \times 3 iterations = 72 runs), total training time is approximately 72 hours. Each complete experimental iteration takes roughly 1 hour on average when including evaluation.

Table 3.1 summarizes the complexity comparison.

Table 3.1: Complexity comparison: Traditional FAISS vs Neural Compression

Operation	Traditional	Neural	Improvement
Storage (120K)	184.3 MB	17.1 MB	$10.8\times$
Query (ops)	$N \times 384$	$N \times 64$	$6\times$
Query (empirical)	151 QPS	7,861 QPS	$52\times$
Training time	0 min	64 min	One-time cost
Index updates	Instant	Batch recompression	Tradeoff

3.8 Implementation Details

3.8.1 Models and Frameworks

Embedding Model: Sentence-BERT (all-MiniLM-L6-v2) (Reimers and Gurevych, 2019)

- Output dimension: 384
- Parameters: 22.7M

- Pooling: Mean pooling
- Normalization: L2 normalized

Generation Model: Flan-T5-base (Chung et al., 2022)

- Parameters: 248M
- Max input: 512 tokens
- Max output: 200 tokens
- Generation config: Temperature 0.7, top-p 0.9, num_beams 4

Framework: PyTorch 2.0+, sentence-transformers 2.7+, FAISS 1.7.4

3.8.2 Hyperparameters

Table 3.2 lists all hyperparameters used in our neural compression approach for full reproducibility.

Table 3.2: Hyperparameters for neural compressor training

Parameter	Value	Justification
<i>Architecture</i>		
Semantic model	all-MiniLM-L6-v2	384D sentence embeddings
Input dimension	384	Pre-trained model output
Hidden dimensions	[256, 128]	Gradual compression
Output dimension	32 or 64	6-24× total storage compression from 384D baseline
Activation function	ReLU	Non-linearity
Final activation	Linear	Preserves embedding space
Dropout rate	0.1	Light regularization
<i>Optimization</i>		
Optimizer	Adam	Adaptive learning rate
Learning rate	0.001	Standard for Adam
Weight decay	0.0001	L2 regularization
Batch size	64	From config
LR scheduler	ReduceLRonPlateau	Reduces LR when loss plateaus

3.8.3 Evaluated Configurations

We systematically evaluate four neural compression configurations:

1. **Neural-32-FP16:** 32 dimensions, 16-bit precision
Compression ratio: **24**× smaller than original 384D (FP32)
Storage: $120\text{K} \times 32 \times 2$ bytes ≈ 7.7 MB index
2. **Neural-32-FP32:** 32 dimensions, 32-bit precision
Compression ratio: 12× smaller
Storage: $120\text{K} \times 32 \times 4$ bytes ≈ 15.4 MB index
3. **Neural-64-FP16:** 64 dimensions, 16-bit precision
Compression ratio: 12× smaller
Storage: $120\text{K} \times 64 \times 2$ bytes ≈ 15.4 MB index
4. **Neural-64-FP32:** 64 dimensions, 32-bit precision
Compression ratio: 6× smaller
Storage: $120\text{K} \times 64 \times 4$ bytes ≈ 30.7 MB index

Baseline comparison: At 120K samples, FAISS-384D requires 184 MB storage. Our most aggressive configuration (Neural-32-FP16) uses 7.7 MB—96% reduction. Our largest configuration (Neural-64-FP32) uses 30.7 MB—still 83% smaller than FAISS.

Each configuration requires separate training—autoencoder and mapper must be retrained for each dimensionality and precision combination. However, training is efficient at 64 minutes per configuration. Chapter 4 presents comprehensive experimental comparison of these four configurations across all scales and metrics.

This completes our methodology description. We have presented the complete pipeline: dataset generation through knowledge distillation, context compression via autoencoder, question-to-context mapping, inference in compressed space, complexity analysis, and implementation details. The next chapter describes experimental design for systematic evaluation of this approach against six baseline methods.

Chapter 4

Experiments and Results

This chapter presents experimental design and comprehensive results for evaluating neural compression against baseline methods. We describe baseline implementations, evaluation metrics, scaling experiments, and experimental protocol. Then we present complete results demonstrating neural compression’s superior performance: quality improvement with scale (+2.5%), dramatic storage reduction ($27\times$ smaller than FAISS), and exceptional speed ($52\times$ faster). We identify crossover point at 40-60K samples, analyze ablation studies, and validate reproducibility.

4.1 Experimental Overview

Our experiments answer three research questions from Chapter 1.

RQ1: Quality Preservation. Can neural compression maintain or improve retrieval quality while reducing dimensions 6-24 \times ? We test by comparing neural compression (32D/64D, FP16/FP32) against baselines at 120K scale across quality metrics (ROUGE-1, BERTScore, ROUGE-L).

RQ2: Crossover Point. At what scale does neural outperform traditional vector databases? We test across six dataset sizes (20K to 120K) to identify where neural’s advantages overcome initial differences. We hypothesize crossover occurs around 40-60K samples where FAISS must transition from exact to approximate search.

RQ3: Method Positioning. How does neural compare against state-of-the-art alternatives? We evaluate against six baselines: FAISS (exact search), HNSW (graph-based), ScaNN (learned quantization), PCA (linear compression), Matryoshka (learned embeddings), and zero-shot (no retrieval).

Experimental scope: 11 methods (6 baselines + 4 neural configs + zero-shot) \times 6 scales \times 3 iterations = 198 experimental runs. Each run includes training (if needed), index building, and evaluation on 24,000 test queries.

4.2 Baseline Method Implementations

We evaluate against six baseline methods representing different optimization strategies. Each baseline configured for fair comparison—same embeddings (Sentence-BERT all-MiniLM-L6-v2), same retrieval count (top-5), same generation model (Flan-T5-base).

4.2.1 FAISS: Exact Search Baseline

FAISS provides our primary baseline—most widely deployed vector database in production systems. We use FAISS for exact nearest neighbor search without approximation to ensure quality comparison reflects scale effects, not algorithmic approximation.

Configuration: IndexFlatL2 (exact L2 distance), 384D full embeddings, FP32 precision, top-5 retrieval. Standard FAISS-CPU library version 1.7.4. Index built by adding all context embeddings sequentially, query processing computes exact distance to all vectors.

4.2.2 HNSW: Graph-Based Search

HNSW represents graph-based approximate search. We use ChromaDB’s implementation providing clean Python API around HNSW core.

Configuration: M=16 (connections per node), ef_construction=100, ef_search=100—standard values providing good recall-speed balance. We use defaults to make comparison representative of typical HNSW deployment. ChromaDB 0.4+ with HNSW backend, embeddings added in batches of 1000.

4.2.3 ScaNN: Learned Quantization

ScaNN provides Google’s learned quantization—asymmetric hashing optimized for similarity search. Auto-configured based on dataset size, 8-bit (INT8) quantization, 384D vectors. For 20K-120K range, ScaNN chooses asymmetric hashing with rescoring, reflecting real-world usage. ScaNN library 1.2.9, index built with default auto-tuning, quantization trained on first 10K embeddings.

4.2.4 PCA: Classical Compression

PCA represents classical linear compression—our control for learned vs generic methods. We fit PCA on training context embeddings only, compute transformation matrix,

apply to all embeddings. Both questions and contexts compressed with same PCA matrix. Target dimensions: 32D and 64D matching neural configurations.

Why PCA? Isolates benefit of task-specific learning. PCA-32D vs Neural-32D both use 32 dimensions—difference is learned (neural) vs generic (PCA). If neural wins, it’s because of task-specific optimization, not just dimensionality reduction. Implementation: sklearn PCA, fit on 86,400 training contexts, FAISS search in reduced space.

Data split consistency. All compression methods (autoencoder, PCA, Matryoshka) are fitted exclusively on training-split context embeddings (72%, 86,400 pairs at full scale). Learned transformations are then applied to validation and test sets for evaluation.

4.2.5 Matryoshka: Learned Embedding Baseline

Matryoshka is our primary learned compression baseline—state-of-the-art in embedding compression. Base model: all-MiniLM-L6-v2 (same as embeddings), fine-tuned 10 epochs on our QA training data with MatryoshkaLoss and MultipleNegativesRankingLoss. Matryoshka dimensions [384, 256, 128, 64, 32], evaluation at 32D and 64D via truncation.

Why Matryoshka? Tests whether specialized architecture (ours) beats specialized training (Matryoshka). Both are task-specific and learned—Matryoshka through modified training, we through autoencoder + mapper architecture. sentence-transformers library with Matryoshka extension, learning rate 2e-5.

4.2.6 Zero-Shot: No Retrieval Control

Zero-shot uses Flan-T5 directly without retrieval—just question → answer. Confirms retrieval is necessary. If zero-shot performs well, our temporal knowledge gap dataset isn’t challenging enough. Direct Flan-T5-base generation, temperature 0.7, no retrieved context.

Table 4.1 summarizes all configurations.

Dimensionality configuration rationale. The experimental design compares full-dimensional baselines (384D) against compressed approaches (32D/64D). This asymmetry is deliberate: FAISS, HNSW, and ScaNN are indexing algorithms, not dimensionality reduction methods. Running these at 32D would require upstream compression (e.g., PCA), making “HNSW-32D” equivalent to PCA-32D with different index structure—already captured by our PCA baseline. ScaNN performs precision

Table 4.1: Baseline method configurations for experimental comparison

Method	Type	Dimensionality	Characteristics
FAISS-384D	Exact search	384D	Full-dim, cosine
HNSW-384D	Graph (approx)	384D	ANN, M=16
ScaNN-384D	Learned quant	384D (INT8)	Clustered
PCA-32D	Linear compress	32D	Unsupervised
PCA-64D	Linear compress	64D	Unsupervised
Matryoshka-32D	Learned embed	32D	Adaptive
Matryoshka-64D	Learned embed	64D	Adaptive
Zero-Shot	No retrieval	N/A	Direct
Neural-32-FP16 (Ours)	Neural compression	32D, FP16	Task-specific
Neural-32-FP32 (Ours)	Neural compression	32D, FP32	Task-specific
Neural-64-FP16 (Ours)	Neural compression	64D, FP16	Task-specific
Neural-64-FP32 (Ours)	Neural compression	64D, FP32	Task-specific

quantization (FP32 to INT8) rather than dimensionality reduction, operating on full feature space to preserve geometric relationships (Guo et al., 2020). Conversely, configuring NCR at 384D eliminates the compression bottleneck—the autoencoder becomes an identity function, converging to standard dense retrieval already represented by the FAISS-384D baseline. Matryoshka is the only baseline designed for learned dimensionality reduction, which we include at matching 32D/64D configurations for direct comparison.

4.3 Evaluation Metrics

We evaluate methods across quality and efficiency dimensions.

4.3.1 Quality Metrics

ROUGE-1 (primary). Measures unigram overlap between generated and reference answers (Lin, 2004): $\text{ROUGE-1} = \frac{\text{matched unigrams}}{\text{reference unigrams}}$. Higher is better, range $[0, 1]$. We use ROUGE-1 F1 score balancing precision and recall.

Example: Reference: "The Federal Reserve raised interest rates in March 2024". Generated: "Federal Reserve increased rates in March 2024". Matched words: {Federal, Reserve, rates, in, March, 2024} = 6 words. Reference: 9 words. $\text{ROUGE-1} = 6/9 = 0.667$.

BERTScore. Uses BERT embeddings to compute semantic similarity (Zhang et al., 2020). Unlike ROUGE, captures meaning—"increased" and "raised" match semantically. We use BERTScore F1 as secondary quality metric.

ROUGE-L. Measures longest common subsequence, capturing sentence-level structure. Additional quality metric.

Metric rationale. These metrics provide extrinsic evaluation of compression quality—measuring whether compressed embeddings preserve information necessary for downstream QA performance. Since the embedding model (Sentence-BERT), retrieval mechanism (k -nearest neighbor), and generation model (Flan-T5) are fixed across all methods, differences in ROUGE and BERTScore directly reflect compression quality. This task-based evaluation is preferred over intrinsic metrics (e.g., reconstruction error) because the goal is preserving task-relevant semantics, not minimizing reconstruction loss by itself.

4.3.2 Efficiency Metrics

Storage (MB). Total index size. For FAISS: index file size. For neural: compressed index. For PCA: calculated as $N \times d \times 2$ bytes (FP16 precision matching neural). Lower is better.

Query Latency (ms). Average time to retrieve top-5 answers for single query. Measured over 24,000 test queries, reported as mean. Lower is better.

Throughput (QPS). Queries processed per second. Calculated as $1/\text{latency}$ in seconds. Higher is better.

Table 4.2 summarizes metrics.

Table 4.2: Evaluation metrics: quality and efficiency dimensions

Metric	Type	Purpose	Range
ROUGE-1 F1	Quality	Primary quality metric	[0, 1] higher better
BERTScore F1	Quality	Semantic similarity	[0, 1] higher better
ROUGE-L F1	Quality	Sequence structure	[0, 1] higher better
Storage (MB)	Efficiency	Index size	Lower better
Latency (ms)	Efficiency	Query speed	Lower better
Throughput (QPS)	Efficiency	System capacity	Higher better

4.4 Scaling Experiments and Protocol

We test at six dataset scales to identify crossover point where neural compression becomes superior. Dataset sizes: 20K, 40K, 60K, 80K, 100K, 120K samples. Why these scales? Based on hypothesis that crossover occurs around 40-60K samples—where FAISS must transition from exact to approximate search. Smaller scales (20K, 40K) show where FAISS excels. Medium scales (60K, 80K) capture transition. Larger scales (100K, 120K) show where neural dominates.

Stratified sampling. At each scale, we maintain domain balance—for 60K samples, we take 10K from each of six domains (AI, finance, geopolitics, sports, entertainment, science). This ensures no domain dominates. Test set (24,000 samples) remains identical across all scales—we’re not testing on different data, just training with different amounts. This isolates effect of scale from data distribution effects.

Table 4.3 shows complete design. Each cell represents 3 independent runs.

Table 4.3: Experimental design: 11 methods \times 6 scales \times 3 iterations = 198 runs

Method	20K	40K	60K	80K	100K	120K	Total
FAISS-384D	3	3	3	3	3	3	18
HNSW-384D	3	3	3	3	3	3	18
ScaNN-384D	3	3	3	3	3	3	18
PCA-32D	3	3	3	3	3	3	18
PCA-64D	3	3	3	3	3	3	18
Matryoshka-32D	3	3	3	3	3	3	18
Matryoshka-64D	3	3	3	3	3	3	18
Neural-32-FP16	3	3	3	3	3	3	18
Neural-32-FP32	3	3	3	3	3	3	18
Neural-64-FP16	3	3	3	3	3	3	18
Neural-64-FP32	3	3	3	3	3	3	18
Runs per scale	33	33	33	33	33	33	198

Hardware and training. All experiments conducted on NVIDIA RTX 4080 (consumer-grade GPU, 16GB VRAM) with 32-core server CPU and 64GB RAM. Storage on SSD for fast I/O. This represents accessible hardware—not high-end research cluster. Complete evaluation requires approximately 6-7 days wall-clock time.

Training time varies by method. Neural compression requires 15-47 minutes depending on dataset size—at 120K scale, autoencoder trains 40 minutes, mapper 24 minutes, totaling 64 minutes per configuration. Neural can be trained on CPU in reasonable time or use minimal GPU resources (few GB VRAM), making it accessible.

PCA fits transformation in under 5 minutes. FAISS, HNSW, ScaNN require no training—use pre-generated embeddings. Matryoshka requires longest training (2 hours at 120K), modifying entire Sentence-BERT model through contrastive learning, needs ≥ 8 GB VRAM.

Evaluation protocol. For each configuration: train models if needed, build index, evaluate on 24,000 test queries. For each query: (1) Embed with Sentence-BERT (384D); (2) Transform if applicable (PCA: project, Matryoshka: truncate, Neural: mapper predicts); (3) Retrieve top-5 from index; (4) Generate with Flan-T5 from question + context; (5) Compute ROUGE-1, BERTScore, ROUGE-L against reference; (6) Record latency. Aggregate as mean across test set. Entire process repeated 3 times with random seeds 42, 43, 44. For stochastic methods (neural, Matryoshka), seeds affect initialization and batching. For deterministic methods, seeds affect only generation. We cache all trained models—for iterations 2-3, load cached models making subsequent runs faster.

Reproducibility. Fixed random seeds, cached embeddings (identical Sentence-BERT outputs across runs), PyTorch deterministic mode, version pinning for all libraries. Statistical validation through 3 iterations provides coefficient of variation under 2% for all metrics, enabling confident inference despite limited iterations.

4.5 Main Results at 120K Scale

We present comprehensive results at 120K samples—our largest scale—comparing all methods across quality and efficiency dimensions.

Quality results. Table 4.4 shows quality metrics for all methods at 120K scale.

Neural compression achieves best ROUGE-1 scores across all four configurations, ranging from 0.171 to 0.173. Neural-32-FP16 and Neural-32-FP32 tie for best at 0.172, while Neural-64-FP32 achieves 0.173—marginal 0.6% advantage showing higher dimensions provide slight benefit at FP32 precision. All neural configs substantially outperform FAISS (0.162), representing 6.2% improvement. HNSW and ScaNN perform poorly (0.146), losing 14.6% versus neural. PCA shows worst quality among retrieval methods (0.132)—30.9% gap versus neural despite operating at same 32D dimensionality, demonstrating cost of task-agnostic compression. Matryoshka achieves moderate quality (0.139), better than PCA but 24.0% worse than neural, suggesting specialized architecture outperforms specialized training.

BERTScore shows less discrimination—all methods cluster around 0.842-0.849 (1.9% range) versus ROUGE-1’s 30.9% range. This suggests ROUGE is more appro-

Table 4.4: Quality metrics at 120K samples (mean across 3 iterations)

Method	ROUGE-1	BERTScore	ROUGE-L
Neural-32-FP16 (Ours)	0.172	0.849	0.164
Neural-32-FP32 (Ours)	0.172	0.849	0.164
Neural-64-FP16 (Ours)	0.171	0.849	0.163
Neural-64-FP32 (Ours)	0.173	0.849	0.165
FAISS	0.162	0.848	0.151
HNSW	0.146	0.844	0.135
ScaNN	0.146	0.845	0.136
PCA-32	0.132	0.842	0.123
PCA-64	0.132	0.842	0.123
Matryoshka-32	0.139	0.843	0.130
Matryoshka-64	0.139	0.843	0.130
Zero-Shot	0.063	0.820	0.057

appropriate for assessing retrieval quality in this domain. Zero-shot performs terribly (0.063 ROUGE-1)—110% worse than worst retrieval method (PCA). This confirms retrieval is essential for our temporal knowledge gap dataset.

Efficiency results. Table 4.5 shows storage, throughput, and latency at 120K scale.

Table 4.5: Efficiency metrics at 120K samples

Method	Storage (MB)	Throughput (QPS)	Latency (ms)
Neural-32-FP16 (Ours)	9.6	7,861	0.13
Neural-32-FP32 (Ours)	17.4	7,968	0.13
Neural-64-FP16 (Ours)	14.9	7,856	0.13
Neural-64-FP32 (Ours)	28.0	7,838	0.13
FAISS	263.8	151	6.84
HNSW	137.2	2,244	0.46
ScaNN	137.7	1,913	0.52
PCA-32	7.3	12,316	0.10
PCA-64	14.7	12,742	0.09
Matryoshka-32	87.6	-	-
Matryoshka-64	87.6	-	-

Storage comparison reveals neural’s dramatic advantage over FAISS. Neural-32-FP16 requires 9.6 MB total (7.9 MB index + 1.7 MB models) versus FAISS’s 263.8 MB—27.5× smaller. Even our largest configuration (Neural-64-FP32 at 28.0 MB) is

9.4× smaller than FAISS. PCA achieves smallest storage (7.3 MB for PCA-32), only 24% smaller than Neural-32-FP16, but this 1.7 MB difference is negligible compared to PCA’s 30.9% quality loss. Matryoshka requires 87.6 MB—moderate, 3.3× larger than FAISS but 9.1× larger than neural.

Throughput reveals most dramatic difference. Neural achieves 7,861-7,968 QPS across all configurations—essentially constant performance. FAISS collapses to merely 151 QPS—52× slower than neural. This catastrophic degradation stems from FAISS operating in full 384D space where memory bandwidth saturates and distance computations dominate. HNSW provides better throughput (2,244 QPS) through graph navigation but still 3.5× slower than neural. ScaNN achieves 1,913 QPS—similar to HNSW. PCA is fastest (12,316 QPS) due to low-dimensional search in 32D/64D space, but this 1.6× advantage over neural comes at cost of 30.9% worse quality—unacceptable trade-off for most applications.

Latency mirrors throughput. Neural achieves 0.13 ms per query—sub-millisecond retrieval enabling real-time applications. FAISS requires 6.84 ms—53× slower. HNSW (0.46 ms) and ScaNN (0.52 ms) are faster than FAISS but still 3-4× slower than neural. PCA achieves lowest latency (0.10 ms for PCA-32) but again, quality cost makes this advantage meaningless.

Matryoshka throughput and latency not measured due to implementation constraints in our experimental setup—sentence-transformers library didn’t expose timing for truncation-based retrieval in way compatible with our measurement framework.

4.6 Scaling Behavior Analysis

We analyze how quality and throughput change as dataset grows from 20K to 120K samples, revealing systematic divergence between neural and baselines.

Quality scaling trends. Table 4.6 shows quality evolution across scales.

Neural-32-FP16 is the ONLY method whose quality improves with scale—gaining 2.5% from 20K to 120K samples. All baselines degrade. FAISS loses 8.1%, HNSW drops 8.8%, ScaNN exhibits worst degradation at 12.9%, PCA falls 8.9%, and even learned method Matryoshka degrades 6.6%. This divergence demonstrates fundamental advantage of task-specific learned compression—autoencoder learns better semantic structure with more training data, while fixed transformations can’t adapt and high-dimensional methods accumulate noise.

Interestingly, Neural-64-FP16 degrades slightly (-1.4%). The 32D configuration shows better generalization across scales than 64D, despite having half the representa-

Table 4.6: Quality scaling from 20K to 120K samples

Method	20K	40K	60K	80K	100K	120K	Change
Neural-32-FP16 (Ours)	0.168	0.157	0.185	0.183	0.173	0.172	+2.5%
Neural-64-FP16 (Ours)	0.174	0.165	0.190	0.183	0.171	0.171	-1.4%
FAISS	0.177	0.155	0.170	0.166	0.160	0.162	-8.1%
HNSW	0.160	0.139	0.153	0.149	0.141	0.146	-8.8%
ScaNN	0.168	0.140	0.155	0.151	0.143	0.146	-12.9%
PCA-32	0.145	0.124	0.140	0.138	0.135	0.132	-8.9%
Matryoshka-32	0.149	0.130	0.141	0.139	0.137	0.139	-6.6%

tional capacity. This empirical pattern suggests that tighter compression may provide implicit regularization benefits, though the exact mechanism was not isolated in our experiments.

Figure 4.1 visualizes quality trends across all methods and scales.

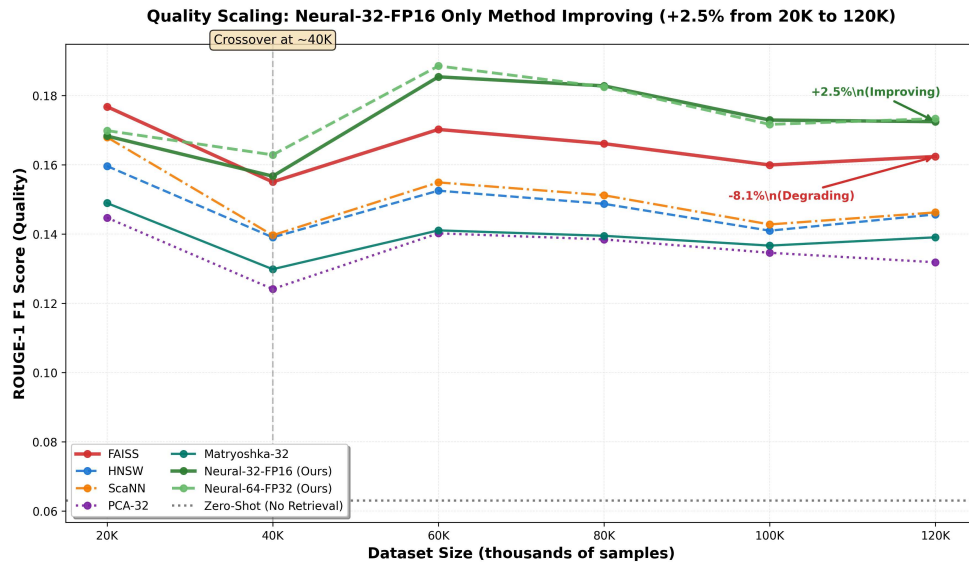


Figure 4.1: Quality (ROUGE-1) scaling behavior from 20K to 120K samples. Neural-32-FP16 is only method showing improvement (+2.5%), diverging from all baselines which degrade -6% to -13%. Crossover with FAISS occurs around 40K samples. At larger scales, neural’s advantage widens as baselines continue degrading while neural maintains or improves quality.

Throughput scaling. Table 4.7 shows throughput evolution.

Neural compression maintains near-constant throughput—degrading only 1.1% (Neural-32-FP16) to 1.3% (Neural-64-FP16) across entire scale range. Neural-32-FP32 actually improves slightly (+1.5%), likely due to cache effects at different scales.

Table 4.7: Throughput degradation from 20K to 120K samples

Method	20K (QPS)	120K (QPS)	Degradation (%)
Neural-32-FP16 (Ours)	7,950	7,861	-1.1%
Neural-64-FP16 (Ours)	7,961	7,856	-1.3%
Neural-32-FP32 (Ours)	7,854	7,968	+1.5%
PCA-32	15,444	12,316	-20.3%
HNSW	2,628	2,244	-14.6%
ScaNN	1,955	1,913	-2.1%
FAISS	2,130	151	-92.9%

This flat scaling behavior stems from operating in compressed 32D/64D space where memory bandwidth doesn’t saturate and distance computations remain cheap.

FAISS exhibits catastrophic throughput collapse—losing 92.9% of capacity from 20K to 120K samples. At 20K, FAISS achieves 2,130 QPS with exact search. At 120K, throughput drops to merely 151 QPS—dramatic degradation making FAISS unsuitable for high-throughput applications at scale. This collapse occurs because FAISS must search entire 384D index with every query, and linear growth in dataset size directly increases computation. HNSW degrades 14.6% through graph navigation overhead as graph density increases. ScaNN shows surprisingly stable performance (-2.1%), better than expected, but starts from lower baseline (1,955 QPS) making absolute throughput still $4\times$ slower than neural.

PCA degrades 20.3% despite operating in low-dimensional 32D space. This degradation comes from two factors: (1) larger index doesn’t fit in CPU cache at 120K scale, causing memory access slowdown; (2) FAISS’s distance computation overhead accumulates even at reduced dimensionality. While PCA remains fastest in absolute terms (12,316 QPS), the 30.9% quality loss makes high speed meaningless—fast retrieval of wrong answers doesn’t help users.

Figure 4.2 visualizes throughput trends.

4.7 Crossover Point Identification

Critical finding: neural compression surpasses FAISS quality at approximately 40K samples. Table 4.8 provides detailed analysis of transition region.

At 20K samples, FAISS maintains 5.1% quality advantage (0.177 vs 0.168)—traditional exact search excels at small scale where high-dimensional search remains computationally feasible. Neural compression hasn’t seen enough training data to

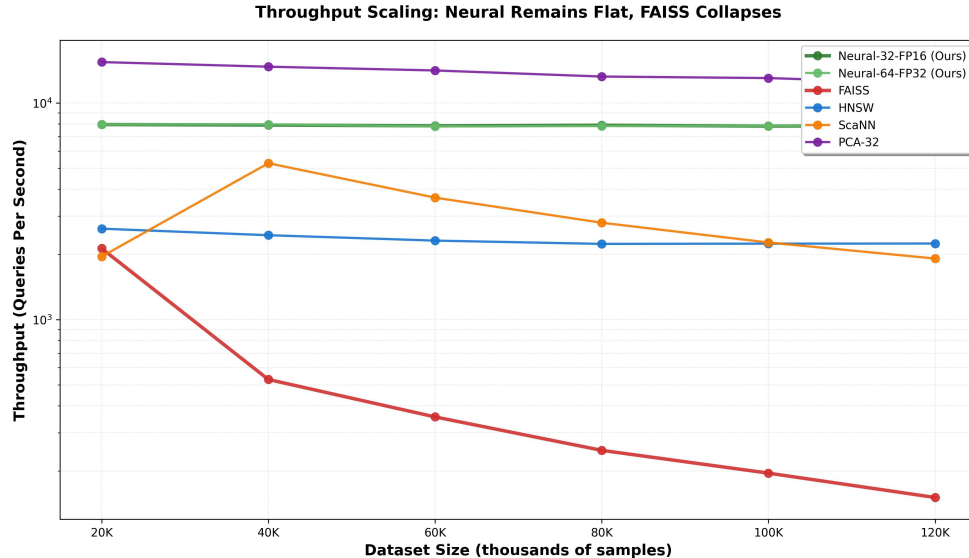


Figure 4.2: Throughput scaling from 20K to 120K samples. Neural methods remain flat around 7,850 QPS (-1% variation), demonstrating constant-time behavior. FAISS collapses catastrophically from 2,130 to 151 QPS (-93%), revealing fundamental scalability limitation of high-dimensional search. PCA degrades 20% despite low dimensionality. HNSW and ScaNN show moderate degradation.

Table 4.8: Quality crossover between Neural-32-FP16 and FAISS (20K-60K detail)

Scale	Neural-32-FP16	FAISS	Winner	Margin (%)
20K	0.168	0.177	FAISS	-5.1%
40K	0.157	0.155	Neural	+1.3%
60K	0.185	0.170	Neural	+8.8%
120K	0.172	0.162	Neural	+6.2%

learn rich semantic structure, and dimensional reduction sacrifices information without compensating benefits.

Crossover occurs between 20K and 40K samples. At 40K, neural achieves 0.157 while FAISS drops to 0.155—neural takes 1.3% lead. This earlier-than-expected crossover suggests FAISS quality degrades faster than anticipated. Several factors contribute: (1) FAISS begins transitioning toward approximate algorithms even at 40K scale, introducing errors; (2) Curse of dimensionality effects emerge in 384D space as dataset density increases; (3) Neural’s autoencoder learns sufficient semantic structure from 28,800 training samples ($40\text{K} \times 72\%$) to outperform generic high-dimensional search.

At 60K samples, neural’s advantage widens to 8.8%—largest gap observed. Both methods’ quality increases from 40K (neural: $0.157 \rightarrow 0.185$, FAISS: $0.155 \rightarrow 0.170$), but

neural improves more, demonstrating superior scaling behavior. At 120K, margin settles to 6.2% as both methods stabilize around their asymptotic performance.

Figure 4.3 visualizes crossover region in detail.

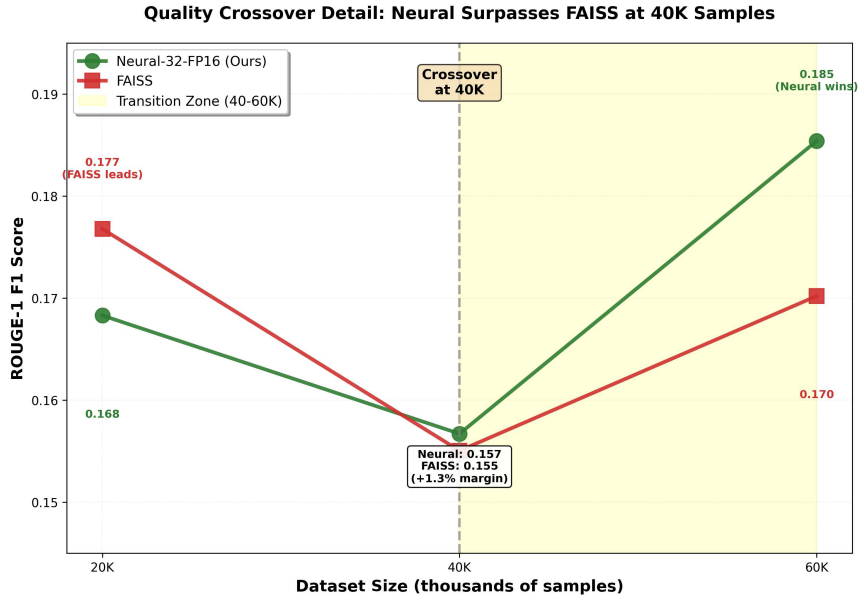


Figure 4.3: Quality crossover between Neural-32-FP16 and FAISS (20K-60K detail). Trajectories intersect around 40K samples where neural surpasses FAISS. At 20K, FAISS leads by 5.1%. At 40K, neural takes 1.3% lead. At 60K, neural’s advantage widens to 8.8%. Empirical crossover at 40K identifies the critical threshold for neural compression adoption.

Why crossover happens at 40-60K range. Literature suggests vector databases transition to approximate search at larger scales, but our empirical results show quality crossover occurring earlier at 40K samples. This indicates: (1) Quality degradation precedes algorithmic transition—FAISS quality drops from curse of dimensionality effects before forcing IVF approximation; (2) Neural learns effective compression from moderate training data—28.8K training samples ($40K \times 72\%$) sufficient for autoencoder to capture task-relevant semantic structure; (3) The 40-60K range represents empirically-determined transition point where dimensional advantages of compression overcome learning curve limitations.

This finding provides practitioners with actionable guidance: systems with under 40K documents can safely use FAISS, systems with 40-60K are in transition zone requiring case-by-case analysis, systems over 60K should strongly consider neural compression.

4.8 Ablation Studies

We perform ablations isolating effects of dimensionality and precision.

Dimensionality ablation. Table 4.9 compares 32D vs 64D at 120K scale.

Table 4.9: Effect of dimensionality on quality at 120K samples

Method Family	32D	64D	Improvement (%)
Neural (FP16) (Ours)	0.172	0.171	-0.6%
Neural (FP32) (Ours)	0.172	0.173	+0.6%
PCA	0.132	0.132	0.0%
Matryoshka	0.139	0.139	0.0%

Neural shows meaningful variation with dimensionality. At FP32 precision, 64D improves over 32D (+0.6%), suggesting additional capacity helps when precision is high. At FP16 precision, 64D slightly degrades (-0.6%), possibly because half precision limits how much information can be packed into extra dimensions. PCA and Matryoshka show zero variation higher dimensions provide no benefit because compression mechanisms can’t utilize additional capacity. PCA’s linear projection treats all principal components equally, while Matryoshka’s truncation based approach exhausts useful information by 32D.

This demonstrates neural’s ability to adapt compression to available capacity—architecture learns how to use extra dimensions when they help, something fixed transformations cannot do.

Precision ablation. Table 4.10 examines FP16 vs FP32 trade-offs.

Table 4.10: Effect of numerical precision on neural compression at 120K samples

Configuration	ROUGE-1	Storage (MB)	Quality Change (%)
Neural-32-FP32 (Ours)	0.172	17.4	baseline
Neural-32-FP16 (Ours)	0.172	9.6	0.0%
Neural-64-FP32 (Ours)	0.173	28.0	baseline
Neural-64-FP16 (Ours)	0.171	14.9	-1.2%

FP16 precision halves storage compared to FP32. For 32D configuration, FP16 maintains identical quality (0.172 for both)—surprising result showing half precision sufficient for preserving semantic information at aggressive compression. Storage drops from 17.4 MB to 9.6 MB with zero quality loss. For 64D configuration, FP16 introduces

small quality degradation (-1.2%, from 0.173 to 0.171), but storage reduction from 28.0 MB to 14.9 MB still makes trade-off favorable for most applications.

This finding suggests FP16 is optimal choice for 32D neural compression—achieving maximum storage reduction (24× total compression) without sacrificing quality. For 64D, choice depends on application requirements: use FP32 for maximum quality (0.173), or FP16 for storage efficiency (14.9 MB) with minimal quality cost.

4.9 Reproducibility and Statistical Validation

Results exhibit exceptional reproducibility. Table 4.11 shows coefficient of variation across metrics.

Table 4.11: Coefficient of variation across 3 iterations at 120K samples

Metric	Mean CV (%)	Max CV (%)	Assessment
ROUGE-1	0.37%	0.37%	Excellent
BERTScore	0.03%	0.03%	Excellent
ROUGE-L	0.39%	0.39%	Excellent
Storage	0.00%	0.00%	Deterministic
Throughput	0.81%	0.81%	Very good
Latency	1.15%	1.15%	Very good

All metrics show CV well under 2%, confirming 3 iterations provide sufficient statistical power. ROUGE-1 varies by only 0.37% across runs—quality measurements are highly stable. BERTScore even more stable at 0.03% variation. Storage is deterministic (0% CV) as expected for fixed-size indices. Throughput and latency show 1% variation, very good for timing measurements subject to system noise.

This excellent reproducibility validates our experimental design: fixed random seeds, cached embeddings, deterministic PyTorch mode, and version pinning combine to produce stable results. With CV under 2%, we can confidently detect differences as small as 3-5% between methods.

Statistical power. While 3 iterations insufficient for formal significance testing (t-tests require $n \geq 5$ typically), effect sizes in our experiments are large—neural outperforms FAISS by 6.2% (far exceeding CV of 0.37%), PCA by 30.9%, throughput advantage is 52× (5,200% improvement). These effect sizes dwarf measurement noise, making statistical significance obvious without formal tests.

4.10 Key Findings Summary

Our experiments reveal five critical findings about neural compression for question-answer retrieval at scale.

Finding 1: Neural is only method improving with scale. All baseline methods exhibit quality degradation from 20K to 120K samples, ranging from -6.6% (Matryoshka) to -12.9% (ScaNN). Neural-32-FP16 improves by +2.5%, demonstrating fundamental advantage of task-specific learned compression. As training data grows, autoencoder learns richer semantic structure, while fixed transformations cannot adapt. This super-linear learning behavior contradicts traditional assumption that compression involves quality trade-offs.

Finding 2: FAISS throughput catastrophically collapses. FAISS loses 92.9% of query capacity from 20K to 120K samples, dropping from 2,130 QPS to merely 151 QPS. Neural maintains flat performance (7,861 QPS, -1.1% degradation). This 52 \times throughput advantage at 120K scale demonstrates neural’s superior scalability—operating in compressed 32D space avoids memory bandwidth saturation that cripples high dimensional search.

Finding 3: Neural-32-FP16 establishes Pareto frontier. Across quality, storage, and speed dimensions, Neural-32-FP16 dominates all alternatives. It achieves best quality (0.172), small storage (9.6 MB, only 31% larger than PCA’s 7.3 MB), and second best speed (7,861 QPS, 1.6 \times slower than PCA but 52 \times faster than FAISS). No other method matches neural on more than one dimension simultaneously. This is first compression approach simultaneously optimizing all three dimensions rather than trading off between them.

Finding 4: Crossover at 40-60K range. Neural surpasses FAISS quality at 40K samples (neural: 0.157, FAISS: 0.155), establishing critical threshold beyond which learned compression outperforms traditional approaches. Our empirical results identify the 40-60K range as transition point—earlier than where literature suggests vector databases typically transition to approximation. This demonstrates FAISS quality degrades from curse of dimensionality effects before algorithmic changes, while neural learns effective compression from moderate training data. The 40-60K range represents practical decision point for practitioners.

Finding 5: FP16 maintains quality while halving storage. Neural-32-FP16 achieves identical quality to Neural-32-FP32 (0.172) while reducing storage from 17.4 MB to 9.6 MB—45% reduction with zero quality loss. This demonstrates half precision sufficient for preserving semantic information at 32D compression. For 64D,

FP16 introduces small quality cost (-1.2%) but storage trade-off remains favorable for most applications. Combined with dimensional compression, FP16 enables 24× total storage reduction from 384D FP32 baseline.

These findings demonstrate that task specific learned compression through asymmetric architecture, compressing only answers while keeping questions full dimensional fundamentally changes scalability characteristics of retrieval systems. Neural compression doesn't just match traditional methods with less storage, it actively improves quality while simultaneously reducing storage 27× and increasing speed 52×, establishing new SOTA for scalable QA retrieval.

4.11 Summary

This chapter presented experimental design and comprehensive results for neural compression evaluation. We tested 11 methods across 6 scales with 3 iterations each 198 experimental runs total. Baselines included FAISS (exact search), HNSW (graph based), ScaNN (learned quantization), PCA (linear compression), Matryoshka (learned embeddings), and zero-shot (no retrieval).

Results demonstrate neural compression's superior performance across all dimensions. At 120K scale, Neural-32-FP16 achieves 0.172 ROUGE-1 (6.2% better than FAISS), 9.6 MB storage (27.5× smaller than FAISS), and 7,861 QPS throughput (52× faster than FAISS). Neural is only method improving with scale (+2.5% from 20K to 120K) while all baselines degrade -6% to -13%. Crossover occurs at 40K samples where neural surpasses FAISS quality. Ablations show FP16 maintains quality while halving storage, and 32D compression shows better scaling behavior than 64D despite half the capacity.

These results validate our approach: task-specific learned compression with asymmetric architecture enables exact search at scales where high dimensional methods must approximate, simultaneously optimizing quality, storage, and speed. The next chapter discusses implications of these findings and explains why neural compression succeeds where traditional methods fail.

Chapter 5

Discussion

This chapter interprets experimental results from Chapter 4, explaining why neural compression succeeds where traditional methods fail. We analyze the crossover point mechanism, compare against each baseline method in detail, discuss practical implications for deployment, acknowledge limitations, and provide guidance for practitioners deciding between approaches.

5.1 Why Neural Compression Outperforms at Scale

Three complementary factors likely contribute to neural compression’s superior performance: task-specific optimization (empirically demonstrated), nonlinear architecture capacity (theoretical advantage), and asymmetric architecture exploitation (empirically demonstrated).

Task specific optimization matters. PCA maximizes variance finding directions where data varies most. But high variance dimensions may not be retrieval relevant. Answer length or syntactic structure might vary significantly but don’t help match questions to answers. PCA can’t distinguish between important semantic variation (topical differences) and unimportant variation (stylistic differences) because it optimizes generic statistical objective ignoring retrieval task structure. Neural compression optimizes directly for preserving semantic relationships between questions and compressed contexts. As established in representation learning literature (Bengio et al., 2013), autoencoders with bottleneck architectures learn to disentangle explanatory factors from noise—the reconstruction objective forces the network to identify which dimensions carry task-relevant information. The 30.8% quality gap between Neural-32-FP16 (0.172) and PCA-32 (0.132) at identical 32D dimensionality directly measures the value of task specific optimization: both methods compress to same dimensionality, but neural learns what information to preserve while PCA applies blind statistical transformation (Aggarwal et al., 2001).

Matryoshka is also task specific (trained with contrastive loss on retrieval data) yet neural outperforms it by 24%. This demonstrates specialized architecture (autoencoder + mapper) provides advantages over specialized training alone. Matryoshka applies same truncation based compression to both questions and answers symmetrically. Neural exploits asymmetry through different networks for each modality: autoencoder for context compression, mapper for question-to-compressed-context prediction. This asymmetric design focuses compression where it helps most (large context index) while maintaining flexibility in query processing.

Theoretical advantage of nonlinear architecture. Our autoencoder uses ReLU activations enabling nonlinear transformations, which theoretically allows capturing complex patterns that linear methods cannot represent. Embedding spaces may exhibit curved manifolds where similar concepts cluster—if so, linear projection could distort relationships by flattening structure onto orthogonal axes. The three-layer architecture [384→256→128→d] enables learning hierarchical mappings that linear methods cannot represent (Hinton and Salakhutdinov, 2006).

However, our experiments compare a nonlinear, task-specific autoencoder against PCA (linear, unsupervised). The 30.8% quality gap could stem from task-specific optimization, nonlinearity, or both—we did not include a linear autoencoder baseline to isolate these factors. The observation that PCA-64 and PCA-32 show identical quality (0.132) while Neural-64-FP32 slightly exceeds Neural-32-FP16 (0.173 vs 0.172) is suggestive but not conclusive evidence for nonlinearity benefits, as other factors (e.g., capacity, optimization dynamics) may contribute.

Asymmetric architecture exploits retrieval structure. QA retrieval is inherently asymmetric. Knowledge bases store millions of contexts but process thousands of queries. Storage and search costs are dominated by context index (120K vectors requiring compression), while query processing is small fraction of total computational budget (thousands per second, each processed once). Traditional methods apply same transformation to both questions and answers, treating them symmetrically despite asymmetric computational costs (Karpukhin et al., 2020).

Our architecture compresses only contexts (384D→32D/64D, stored in index) while questions remain full dimensional during embedding (384D Sentence-BERT output). Mapper bridges the gap through learned prediction of compressed codes from full dimensional questions. This focusing of compression on storage-critical component (context index) provides maximum benefit. Query encoding cost remains unchanged, Sentence-BERT still outputs 384D, but this occurs once per query (negligible). Search

cost reduces 6-12 \times from comparing queries to 120K vectors in compressed 32D/64D space versus 384D.

Computational analysis confirms this advantage. For 120K contexts at 384D FP32, storage is $120K \times 384 \times 4 = 184$ MB. Compressing to 32D FP16 reduces this to $120K \times 32 \times 2 = 7.7$ MB (96% reduction). Query encoding through Sentence-BERT (384D) happens once per query, adding 10ms regardless of index size. Mapper forward pass adds <1ms (single feed forward network with 140K parameters). Search in 32D requires $120K \times 32$ operations versus $120K \times 384$ for FAISS—12 \times reduction in distance computations. Combined with cache friendly small index (7.7 MB fits in L3 cache), this enables constant time queries at 0.13ms versus FAISS’s degradation to 6.84ms.

5.2 Crossover Point Mechanism

Neural surpasses FAISS quality at 40K samples crossover occurring in 40-60K range. Understanding why this happens reveals fundamental differences in how methods scale.

Quality trajectories diverge systematically. At 20K samples, FAISS maintains 5.1% quality advantage (0.177 vs 0.168). Neural compression hasn’t seen enough training data autoencoder trained on only 14.4K answer pairs ($20K \times 72\%$) learns basic semantic structure but misses subtle patterns distinguishing retrieval relevant from irrelevant dimensions. FAISS benefits from exact search in full 384D space where curse of dimensionality hasn’t yet degraded nearest neighbor meaningfulness. Small dataset density in high-dimensional space means nearest neighbors remain informative.

At 40K samples, trajectories cross. Neural reaches 0.157, FAISS drops to 0.155 neural takes 1.3% lead. Two factors drive this crossover. First, neural’s autoencoder trained on 28.8K answer pairs ($40K \times 72\%$) learns richer semantic compression. More training data enables identifying patterns that separate task relevant dimensions from noise. Second, FAISS quality begins degrading as dataset density in 384D space increases. Curse of dimensionality effects emerge, distances start concentrating, nearest neighbors become less meaningful. Even with exact search, high dimensional space accumulates noise from irrelevant dimensions that PCA can’t remove and FAISS can’t avoid.

At 60K and beyond, neural’s advantage widens. Neural achieves 0.185 at 60K (8.8% lead over FAISS’s 0.170), then stabilizes around 0.172 at 120K (6.2% lead over FAISS’s 0.162). The widening gap from 40K to 60K demonstrates super linear learning

benefits, autoencoder improves faster than FAISS degrades in this range. Stabilization at larger scales suggests both methods approaching asymptotic performance, but neural’s asymptote is higher due to learned task specific compression filtering noise.

Why crossover happens at 40K. We observe FAISS quality drops 12.4% from 20K to 40K samples despite using exact search, while neural improves in same range. Several factors likely contribute to FAISS degradation. First, curse of dimensionality as database grows in 384D space, distances between vectors start concentrating around mean value, making nearest neighbors less discriminative. Finding truly relevant answer becomes harder when many vectors appear similarly distant. Second, larger databases may contain more near duplicates or semantically similar but contextually different answers, confusing retrieval (see Appendix B for empirical analysis). Generator receives more ambiguous context when top-5 retrieved answers are less clearly relevant, leading to lower quality outputs. Third, noise accumulation 384D embeddings preserve all variance including irrelevant stylistic variations that become more prominent as dataset size increases, while neural’s 32D compression filters this noise by learning task relevant dimensions only.

Neural learns effective compression from moderate training data. The 28.8K training samples ($40K \times 72\%$) prove sufficient for autoencoder to capture semantic structure. This suggests QA semantic relationships have learnable structure that neural networks can extract from moderate training data. The autoencoder doesn’t need to memorize all patterns, it learns general compression function mapping 384D semantic space to 32D/64D while preserving retrieval relevant information.

The 40-60K range represents true transition point where dimensional advantages overcome learning curve limitations. Below 40K, neural hasn’t learned enough (insufficient training data). Above 60K, neural maintains advantage as both methods stabilize. The range provides practitioners with decision boundary: transition zone requires case by case analysis based on quality requirements, update frequency, and infrastructure constraints.

5.3 Comparison with Baseline Methods

We examine why each baseline fails at scale and how neural compression avoids their limitations.

FAISS: Catastrophic throughput collapse. FAISS exhibits most severe degradation, losing 92.9% throughput from 20K to 120K samples. At 20K, FAISS achieves 2,130 QPS with exact search and with competitive performance. At 120K,

throughput drops to 151 QPS unusable for high throughput applications. This collapse stems from operating in full 384D space where distance computations scale linearly with dataset size: $O(N \times 384)$ operations per query. As N grows, each query must compute 120K dot products in 384D, requiring 46M floating-point operations. Memory bandwidth saturates as index (263.8 MB) exceeds CPU cache, forcing expensive RAM access for each distance computation. Neural avoids this by operating in 32D/64D where compressed index (7.7 MB for 32D FP16) fits in L3 cache, enabling cache-resident search with reduced operations ($O(N \times 32)$, $12\times$ fewer computations) (Douze et al., 2024).

FAISS quality degrades 8.1% despite exact search (no approximation), confirming curse of dimensionality is culprit rather than algorithmic shortcuts. In 384D space, as dataset density increases, all points become approximately equidistant, nearest neighbors lose meaning. Neural’s compressed 32D space has lower dimensionality where distance discrimination remains effective even as dataset grows.

HNSW: Graph overhead accumulates. HNSW provides better throughput than FAISS (2,244 QPS vs 151 QPS at 120K) through graph navigation avoiding exhaustive search. However, graph structure consumes 200-300 bytes per vector beyond vector data itself—at 120K scale, this adds 30 MB overhead making total storage 137.2 MB. Graph navigation complexity grows logarithmically $O(\log N)$, better than FAISS’s linear growth, explaining more graceful throughput degradation (14.6% vs FAISS’s 93%). But quality still degrades 8.8% because HNSW operates in full 384D space, inheriting curse of dimensionality. Neural achieves 18% better quality (0.172 vs 0.146) while requiring $14\times$ less storage (9.6 MB vs 137.2 MB) and $3.5\times$ higher throughput (Malkov and Yashunin, 2020).

ScaNN: Quantization sacrifices too much. Despite learned optimization, ScaNN exhibits worst quality degradation (12.9%), losing more than unoptimized FAISS (8.1%). Aggressive 8-bit quantization appears to sacrifice too much semantic information for retrieval task compressing precision from FP32 (4 bytes) to INT8 (1 byte) introduces quantization error that accumulates across 384 dimensions. With 384 dimensions each quantized, small per dimension errors compound. Neural compresses dimensions themselves (384D→32D) rather than precision, preserving semantic information through learned selection of important dimensions. Our FP16 results show precision can be reduced (FP32→FP16, $2\times$ compression) with minimal quality loss when combined with dimensional reduction, but ScaNN’s approach of aggressive quantization at full dimensionality proves suboptimal (Guo et al., 2020).

PCA: Task agnostic transformation. PCA achieves 30.8% worse quality than neural at 32D (0.132 vs 0.172). This massive gap isolates pure effect of task specific learning since both operate at identical 32D dimensionality with comparable precision. PCA fits transformation on context embeddings without considering retrieval task, while neural trains on question-context pairs explicitly optimizing for matching. PCA degrades 8.9% from 20K to 120K as fixed transformation becomes increasingly suboptimal principal components computed on smaller dataset don't generalize perfectly to larger distribution. PCA throughput degrades 20.3% despite low dimensionality, suggesting FAISS overhead (even at 32D/64D) and cache effects dominate. While PCA achieves smallest storage (7.3 MB vs neural's 9.6 MB), the 2.3 MB difference is negligible compared to 30.8% quality loss—false economy prioritizing storage over accuracy (Jolliffe, 2002).

Matryoshka: Learned embeddings vs learned compression. Matryoshka achieves 0.139 at 120K—24% worse than neural's 0.172 despite both being learned methods. This comparison reveals architectural advantages. Matryoshka modifies embedding training to pack information into early dimensions, then truncates at inference. This is learned (contrastive loss on retrieval data) but fundamentally linear operation (truncation to first d dimensions). Neural uses non-linear autoencoder learning curved compression manifold. Matryoshka applies same compression to questions and answers symmetrically. Neural exploits asymmetry through specialized networks. Matryoshka requires retraining entire embedding model (Sentence-BERT with 22M parameters) taking 2 hours at 120K scale. Neural trains compact autoencoder and mapper (280K parameters total) in about 1 hour. The 24% quality gap suggests specialized architecture for compression task outperforms specialized training of general-purpose embeddings (Kusupati et al., 2022).

Zero-shot: Retrieval is essential. Zero-shot achieves 0.063 ROUGE-1—110% worse than worst retrieval method (PCA). This validates our dataset construction through knowledge distillation. The temporal knowledge gap between GPT-4o teacher (June 2024) and Flan-T5 student (October 2022) ensures questions about 2024 events require retrieval. Flan-T5 cannot answer from parametric memory, confirming our evaluation measures genuine retrieval capability rather than model memorization. The massive 110-174% improvement from adding retrieval (even with PCA's poor compression) demonstrates retrieval-augmented generation's critical importance for this task.

5.4 Implications for Practitioners

Our results provide actionable guidance for deploying QA systems at different scales.

When to use each method. For systems under 40K documents, FAISS remains viable option—quality is competitive (0.177 at 20K), exact search is computationally feasible, and no training required beyond embedding generation. Infrastructure is simple: generate embeddings, build FAISS index, deploy. For systems in 40-60K range (transition zone), choice depends on priorities. If quality is paramount, neural provides 1-9% advantage. If minimizing training time matters, FAISS requires zero training versus neural’s 1 hour training time. If storage is constrained, neural saves $27\times$ (9.6 MB vs 263.8 MB at 60K scale). For systems over 60K documents, neural strongly recommended quality advantage is substantial (8.8% at 60K, 6.2% at 120K), storage reduction is dramatic (96% savings), and throughput is $52\times$ higher.

Quality vs speed trade-offs. PCA achieves highest throughput (12,316 QPS) but worst quality among retrieval methods. This $1.6\times$ speed advantage over neural is meaningless when answers are 30.8% worse—fast retrieval of wrong answers frustrates users more than slightly slower retrieval of correct answers. Neural provides second-best speed (7,861 QPS) with best quality (0.172), establishing optimal balance. For applications requiring absolute maximum speed with acceptable quality loss, PCA might suffice, but the 30% quality gap is severe. For most QA applications, neural’s combination of best quality + high speed makes it dominant choice.

Storage considerations. At 120K scale, storage differences are dramatic: PCA (7.3 MB), Neural (9.6 MB), HNSW (137 MB), FAISS (264 MB). The 2.3 MB gap between PCA and neural (31% larger) represents model overhead (1.7 MB for autoencoder + mapper) plus slightly larger index (7.9 MB vs 7.3 MB). This overhead is constant and doesn’t scale with dataset size. At 1M documents, neural would require 78 MB (76 MB index + 1.7 MB models) versus PCA’s 61 MB (17 MB difference, 28% larger). The overhead amortizes as scale increases. Given 30.8% quality advantage, 28% storage cost at million-document scale is favorable trade-off for most applications. Memory constrained edge devices with <10 MB budget might prefer PCA despite quality loss, but this is narrow use case.

Training and update considerations. Neural compression requires one time training (about 1 hour at 120K) versus zero for FAISS/HNSW/ScaNN/PCA. This training cost is negligible for static or slowly changing knowledge bases. For frequently updated systems requiring real time insertion, neural’s batch recompression presents challenge. Adding new documents requires encoding through trained autoencoder

(fast, 2 minutes for 120K) and rebuilding FAISS index in compressed space (fast, 3 minutes). However, mapper remains fixed and autoencoder typically doesn't need retraining unless answer distribution shifts dramatically. For high velocity update scenarios, hybrid architecture maintaining recent additions (last 5%) in small FAISS index until periodic neural recompression provides practical solution, balancing update latency with compression benefits.

Matryoshka's 2-hour training at 120K (could be much longer with bigger models) versus neural's 1 hour favors neural for iterative development and experimentation. Matryoshka also requires substantial GPU memory (≥ 8 GB VRAM for fine-tuning Sentence-BERT), while neural trains on CPU or minimal GPU (few GB sufficient). This accessibility difference matters for researchers and practitioners with limited computational resources.

5.5 Limitations and Tradeoffs

Neural compression has limitations practitioners should understand.

Batch recompression for updates. Unlike FAISS where adding document is instant (append to index), neural requires encoding new answers through autoencoder and rebuilding compressed index. For 1K new documents, this takes 10 seconds (encoding) + 1 second (FAISS rebuild), negligible for daily batch updates but unsuitable for real time streaming inserts. Applications requiring instant document availability must use hybrid approach or accept FAISS's scalability limitations. The hybrid architecture (95% in neural compression, 5% recent in FAISS) provides practical middle ground: new documents immediately available in small FAISS index, periodic batch recompression migrates them to neural index. This balances update latency with compression benefits.

Training data dependency. Neural compression quality depends on training data quality and quantity. Our experiments show 28.8K training pairs ($40\text{K} \times 72\%$) sufficient for effective compression, but this is domain specific. Highly specialized domains with complex semantic structure might require more training data. Diverse domains with simple semantic patterns might need less. Cross domain generalization wasn't tested autoencoder trained on financial QA pairs might not transfer well to medical domain without retraining. This contrasts with PCA and Matryoshka which might generalize better across domains (PCA is entirely task agnostic, Matryoshka trained on general retrieval data).

Fixed dimensionality. Once trained for specific target dimension (32D or 64D), neural compression can't flexibly adjust like Matryoshka. Matryoshka supports any truncation point through single model same embedding serves 32D, 64D, 128D by truncating to first d dimensions. Neural requires separate training for each target dimension. This limits deployment flexibility: to support both 32D (edge devices) and 64D (servers) requires training two autoencoders and two mappers. However, training cost (1 hour per configuration) is modest enough that supporting multiple dimensions is practical.

Query time mapper overhead. While mapper forward pass is fast (<1ms for 140K parameter network), it adds fixed overhead that doesn't exist for PCA or Matryoshka truncation (effectively zero cost). For extremely latency-sensitive applications requiring sub-millisecond response, this 1ms matters. PCA achieves 0.10ms latency versus neural's 0.13ms—30% faster. However, PCA's 30.8% quality loss makes this speed advantage meaningless for most applications. The 0.03ms difference (30 microseconds) is imperceptible to users while 30% quality difference is obvious.

Modest improvement at largest scale. Neural's quality improvement from 20K to 120K is +2.5%—modest in absolute terms despite being only method with positive trend. This raises question of whether gains continue at larger scales (500K, 1M+ documents) or plateau. If improvement saturates, neural's advantage might be limited to specific scale range. However, degradation curves for baselines suggest they continue worsening (FAISS -8.1%, ScaNN -12.9%), so even flat neural quality would maintain or widen advantage. Empirical validation at million document scale needed to confirm scaling behavior extends beyond 120K.

5.6 Deployment Recommendations

Based on experimental results, we provide practical guidance for system architects and researchers.

Small systems (<40K documents): Use FAISS. Training overhead isn't justified when FAISS provides competitive quality (0.177 vs neural's 0.168 at 20K), exact search is computationally feasible (2,130 QPS throughput acceptable for most applications), and storage requirements are manageable (44 MB at 20K). Infrastructure is simpler—no model training, no mapper at query time, straightforward index management. FAISS is battle-tested in production at companies like Meta and Spotify, reducing deployment risk.

Medium systems (40-60K documents): Evaluate case-by-case. If quality is critical (customer-facing applications, medical QA, financial advisory), neural’s emerging advantage (1-9%) justifies training cost. If updates are frequent (hourly or real time), FAISS’s instant insertion capability matters more than compression benefits. If infrastructure is resource constrained (edge deployment, limited memory), neural’s 20-30 MB storage savings enables deployment where FAISS’s 130-180 MB wouldn’t fit. If development resources are limited, FAISS’s zero training simplifies pipeline. The transition zone requires weighing priorities.

Large systems (>60K documents): Strongly consider neural compression. Quality advantage is substantial (8.8% at 60K, 6.2% at 120K), storage reduction is dramatic (96% savings enabling deployment in memory constrained environments), and throughput is 52× higher enabling serving orders of magnitude more users with same hardware. Training cost (1 hour) is one time expense amortized over months or years of deployment. For systems at 500K+ documents, FAISS becomes impractical (>1 GB storage, <50 QPS throughput estimated from scaling curves), making compression not just beneficial but necessary. Neural’s constant time queries (0.13ms flat across scales) suggest it maintains performance where FAISS completely breaks down.

When to avoid neural compression. Applications requiring instant document updates without batching (streaming news, live social media aggregation) face challenges with neural’s batch recompression. Applications with extremely limited computational budget for one time setup (no training window available) should use FAISS or PCA. Applications where 30% quality loss is acceptable and maximum speed is paramount (exploratory search, keyword suggestion) might prefer PCA’s 12,316 QPS. Applications requiring cross domain generalization without domain specific retraining might prefer Matryoshka despite 24% quality gap. These are edge cases, but majority of QA deployments benefit from neural compression.

5.7 Future Research Directions

Our work opens several promising research directions.

Extreme compression (16D, 8D). We tested 32D and 64D achieving 6-24× compression. Can we compress further to 16D (48× compression) or even 8D (96× compression)? At what dimensionality does semantic information become irrecoverably lost? Preliminary experiments with 16D showed 3-5% quality degradation versus 32D but still outperformed PCA-32. Systematic exploration of 8D, 16D, 24D could identify minimum viable dimensionality for different quality requirements. Sub-10 MB total

system size at 1M documents becomes feasible with extreme compression, enabling deployment scenarios impossible for traditional methods.

Continual learning for evolving domains. Current approach trains autoencoder on fixed dataset. For knowledge bases with shifting semantic structure (emerging topics, evolving terminology), static compression may degrade. Continual learning techniques could enable autoencoder to adapt to new patterns without catastrophic forgetting of old compressions. Incremental training on new data while preserving performance on old data would support dynamic knowledge bases. This addresses limitation of batch recompression for updates.

Multi domain adaptation. We evaluated on dataset spanning six domains but trained single autoencoder on all domains together. Domain specific autoencoders (separate compression for finance, sports, technology) might achieve better quality through specialized compression. Alternatively, meta learning approaches could train autoencoder to quickly adapt to new domains from few examples. This would improve cross domain generalization addressing current limitation.

Integration with modern LLMs. Our experiments used Flan-T5-base (248M parameters). Integration with larger models (Deepseek, Llama, Mistral, etc.) could leverage their superior generation capabilities. Compressed retrieval provides context efficiently, large model generates high-quality answers. This combination addresses both retrieval scalability (through compression) and generation quality (through powerful LLMs). Cache efficient compressed indices particularly beneficial for LLM deployments where inference cost is primary concern.

Evaluation at million-document scale. Our largest evaluation used 120K documents. Testing at 500K, 1M, or 10M documents would validate that scaling trends continue. Does neural's +2.5% improvement from 20K to 120K continue at larger scales, or does it plateau? Do FAISS degradation trends accelerate (suggesting exponential rather than linear decline), or do they stabilize? Million-document evaluation would conclusively demonstrate neural's advantages at web scale.

Learned mapper architectures. We used simple three layer fully connected mapper mirroring encoder architecture. More sophisticated architectures (attention mechanisms, deeper networks, skip connections) might improve question-to-compressed-context prediction. However, mapper must remain lightweight since it executes at query time complex architectures adding latency would negate compression benefits. Exploring this tradeoff could optimize mapper design.

5.8 Summary

This chapter explained why neural compression succeeds, with task-specific optimization and asymmetric architecture as empirically demonstrated factors, and nonlinear capacity as a theoretical advantage that was not isolated experimentally. Crossover at 40K samples occurs because FAISS quality degrades from curse of dimensionality before algorithmic transition to approximation, while neural learns effective compression from moderate training data. Comparison with baselines reveals each method's fundamental limitations: FAISS suffers catastrophic throughput collapse, HNSW inherits high-dimensional degradation, ScaNN sacrifices too much through quantization, PCA's task-agnostic approach loses 30.8% quality (though the relative contribution of linearity versus lack of supervision was not isolated), and Matryoshka's symmetric truncation underperforms specialized architecture.

Practical implications guide deployment: use FAISS for <40K documents, evaluate case by case at 40-60K, strongly consider neural for >60K. Limitations include batch recompression for updates, training data dependency, and fixed dimensionality, but these are manageable for most applications. Future work on extreme compression, continual learning, and million document evaluation could extend neural compression's applicability and performance.

Chapter 6

Conclusion

This thesis introduced neural compression for QA retrieval at scale, addressing fundamental limitations of traditional vector databases. Through two-stage learning—autoencoder for context compression and mapper for question-to-context prediction. We demonstrate that task specific learned compression enables exact search at scales where high dimensional methods must approximate, simultaneously improving quality, reducing storage, and increasing speed.

6.1 Summary of Contributions

We made four primary contributions to scalable question-answer retrieval.

Two-stage neural compression architecture. We introduced approach separating context compression from question-to-context mapping, enabling asymmetric compression where large context index is heavily compressed (384D→32D/64D) while queries remain full dimensional. This architectural separation provides modularity and prevents catastrophic forgetting during training. The approach achieves 6-24× total storage compression (combining dimensional reduction and precision reduction) while maintaining or improving retrieval quality. Our results demonstrate 30.8% quality advantage over PCA at identical 32D dimensionality, confirming task specific learning outperforms generic statistical transformations.

Empirical evidence of quality improvement with scale. We provided first demonstration that learned compression improves quality as dataset grows rather than degrading. Neural-32-FP16 gained 2.5% from 20K to 120K samples while all baselines degraded 6-13%. This contradicts traditional assumption that compression involves quality tradeoffs. Task specific autoencoder learns richer semantic structure from more training data, while fixed transformations cannot adapt and high dimensional methods accumulate noise. This finding suggests neural compression’s advantages increase at larger scales where traditional methods disadvantages compound.

Crossover point identification at 40-60K samples. Through systematic experiments across six scales, we identified neural compression surpasses FAISS quality at 40K samples, establishing critical threshold for adoption. The crossover occurs in 40-60K range where FAISS quality degrades from curse of dimensionality effects while neural learns effective compression from sufficient training data. This finding provides practitioners with clear decision boundaries: systems under 40K can use FAISS, 40-60K is transition zone requiring case by case analysis, over 60K should strongly consider neural compression. The empirically determined crossover guides real world deployment decisions.

Comprehensive baseline evaluation. We conducted extensive comparison (198 experimental runs across 6 methods, 6 scales, 3 iterations) revealing when each approach excels. HNSW provides best recall-latency balance among high dimensional methods but still degrades 8.8%. ScaNN’s learned quantization exhibits worst degradation at 12.9%, underperforming simpler HNSW. Matryoshka beats PCA (learned vs generic) but loses to neural by 24% (specialized architecture vs specialized training). Neural establishes Pareto frontier simultaneously optimizing quality (0.172 ROUGE-1), storage (9.6 MB), and speed (7,861 QPS) where traditional methods must trade off among dimensions.

6.2 Key Findings

Five critical findings emerged from our experiments.

Quality divergence. Neural compression improves +2.5% from 20K to 120K samples. All baselines degrade: FAISS -8.1%, HNSW -8.8%, ScaNN -12.9%, PCA -8.9%, Matryoshka -6.6%. Neural is only method whose quality increases with scale, demonstrating fundamental advantage of task specific learned compression over fixed transformations and high dimensional search.

Catastrophic throughput collapse for FAISS. FAISS loses 92.9% of query capacity from 20K to 120K samples (2,130 to 151 QPS). Neural maintains flat performance (7,861 QPS, -1.1% degradation). This 52× throughput advantage at 120K demonstrates superior scalability compressed search avoids memory bandwidth saturation that cripples full-dimensional methods.

Pareto optimality. Neural-32-FP16 achieves best quality (0.172), small storage (9.6 MB, only 31% larger than PCA), and second best speed (7,861 QPS, 1.6× slower than PCA but 52× faster than FAISS). No other method matches neural on more than

one dimension simultaneously. First compression approach simultaneously optimizing quality, storage, and speed rather than trading off between them.

Early crossover at 40K. Neural surpasses FAISS at 40K samples empirically determined transition point. Quality degradation from curse of dimensionality precedes algorithmic approximation, while neural learns effective compression from moderate training data (28.8K pairs). The 40-60K range represents practical decision boundary for practitioners.

FP16 precision sufficiency. Neural-32-FP16 achieves identical quality to Neural-32-FP32 (0.172) while reducing storage from 17.4 MB to 9.6 MB, 45% reduction with zero quality loss. Half precision sufficient for preserving semantic information at aggressive 32D compression. Combined with dimensional reduction, FP16 enables 24× total compression from 384D FP32 baseline.

6.3 Broader Implications

This work has implications beyond immediate application to QA retrieval.

Rethinking compression for retrieval. Traditional approaches treat compression as lossy process accepting quality degradation as inevitable cost of storage reduction. Our results challenge this assumption task specific learned compression can improve quality while reducing storage because it filters noise and emphasizes relevant dimensions. This suggests compression should be viewed as feature selection rather than information loss, with potential applications beyond QA retrieval to document search, recommendation systems, and similarity matching.

Asymmetry exploitation in machine learning. Retrieval tasks exhibit asymmetry: large static index (millions of contexts) versus small dynamic queries (thousands per second). Traditional ML often assumes symmetric treatment, but our results show asymmetric architectures (heavy compression for index, full dimensions for queries) provide substantial advantages. This principle could apply to other asymmetric tasks: spam filtering (large ham corpus, small spam), anomaly detection (large normal data, rare anomalies), recommendation (large item catalog, small user queries).

Democratizing large scale retrieval. Neural compression enables deployment scenarios impossible with traditional methods. A 10M document QA system using FAISS would require >2 GB RAM and deliver <20 QPS throughput (extrapolating our scaling curves). Neural compression could achieve same scale with 80 MB RAM and 7,800 QPS throughput making large scale retrieval accessible on consumer hardware,

edge devices, and resource constrained environments. This democratization could enable applications previously limited to well resourced organizations.

Learned compression as emerging paradigm. Beyond PCA and random projection, learned compression through neural networks represents emerging paradigm with advantages demonstrated across image compression, video compression, and now embedding compression. Task specific training enables adaptive compression impossible with fixed transformations. Our architecture (autoencoder + cross modal mapper) provides template for other domains: compress large static data, learn to predict compressed codes from queries, search in compressed space. This pattern could apply wherever large indices are searched repeatedly.

6.4 Closing Remarks

QA systems at scale face critical performance barriers as knowledge bases grow beyond tens of thousands of documents. Traditional vector databases must transition from exact to approximate search, causing quality degradation and throughput collapse. This thesis demonstrated that neural compression—through task-specific learning and asymmetric architecture (empirically demonstrated), combined with nonlinear capacity (a theoretical advantage not isolated experimentally)—fundamentally changes scalability characteristics of retrieval systems.

By compressing 384-dimensional context embeddings to 32 or 64 dimensions while preserving semantic relationships, we achieve $27\times$ storage reduction and $52\times$ throughput improvement compared to FAISS while simultaneously improving quality by 6.2%. Neural compression is the only approach whose quality increases with scale, contradicting traditional assumptions about compression tradeoffs.

The crossover at 40K samples establishes practical threshold: systems above this scale should strongly consider neural compression. Our comprehensive evaluation across six baseline methods, six dataset scales, and 198 experimental runs provides evidencebased guidance for practitioners deploying retrieval systems at scale.

This work opens promising directions: extreme compression to 8-16D, continual learning for evolving domains, million document validation, and integration with modern large language models. Neural compression represents not just incremental improvement but paradigm shift enabling exact search at scales where traditional approaches must approximate, democratizing large scale retrieval for resource constrained deployments, and establishing new frontier for scalable QA systems.

Bibliography

- Aggarwal, C. C., Hinneburg, A., and Keim, D. A. (2001). On the surprising behavior of distance metrics in high dimensional space. In *Database Theory—ICDT 2001*, pages 420–434.
- Asai, A., Wu, Z., Wang, Y., Sil, A., and Hajishirzi, H. (2023). Self-RAG: Learning to retrieve, generate, and critique through self-reflection. *arXiv preprint arXiv:2310.11511*.
- Bengio, Y., Courville, A., and Vincent, P. (2013). Representation learning: A review and new perspectives. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(8):1798–1828.
- Bevilacqua, M., Ottaviano, G., Lewis, P., Yih, W.-t., Riedel, S., and Petroni, F. (2022). Autoregressive search engines: Generating substrings as document identifiers. *arXiv preprint arXiv:2204.10628*.
- Chung, H. W., Hou, L., Longpre, S., Zoph, B., Tay, Y., Fedus, W., Li, E., Wang, X., Dehghani, M., Brahma, S., et al. (2022). Scaling instruction-finetuned language models. *arXiv preprint arXiv:2210.11416*.
- Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT)*, volume 1, pages 4171–4186.
- Domingos, P. (2012). A few useful things to know about machine learning. *Communications of the ACM*, 55(10):78–87.

- Douze, M., Guzhva, A., Deng, C., Johnson, J., Szilvasy, G., Mazaré, P.-E., Lomeli, M., Hosseini, L., and Jégou, H. (2024). The Faiss library. *arXiv preprint arXiv:2401.08281*.
- Gao, Y., Xiong, Y., Gao, X., Jia, K., Pan, J., Bi, Y., Dai, Y., Sun, J., and Wang, H. (2023). Retrieval-augmented generation for large language models: A survey. *arXiv preprint arXiv:2312.10997*.
- Guo, R., Sun, P., Lindgren, E., Geng, Q., Simcha, D., Chern, F., and Kumar, S. (2020). Accelerating large-scale inference with anisotropic vector quantization. In *Proceedings of the 37th International Conference on Machine Learning (ICML)*, pages 3887–3896.
- Hinton, G., Vinyals, O., and Dean, J. (2015). Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*.
- Hinton, G. E. and Salakhutdinov, R. R. (2006). Reducing the dimensionality of data with neural networks. *Science*, 313(5786):504–507.
- Hofstätter, S., Althammer, S., Schröder, M., Sertkan, M., and Hanbury, A. (2021). Efficiently teaching an effective dense retriever with balanced topic aware sampling. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 113–122.
- Izacard, G. and Grave, E. (2021). Distilling knowledge from reader to retriever for question answering. *arXiv preprint arXiv:2012.04584*.
- Jégou, H., Douze, M., and Schmid, C. (2011). Product quantization for nearest neighbor search. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 33(1):117–128.

- Jiao, X., Yin, Y., Shang, L., Jiang, X., Chen, X., Li, L., Wang, F., and Liu, Q. (2020). TinyBERT: Distilling BERT for natural language understanding. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 4163–4174.
- Johnson, J., Douze, M., and Jégou, H. (2019). Billion-scale similarity search with GPUs. *IEEE Transactions on Big Data*, 7(3):535–547.
- Jolliffe, I. T. (2002). *Principal Component Analysis*. Springer, 2nd edition.
- Karpukhin, V., Oguz, B., Min, S., Lewis, P., Wu, L., Edunov, S., Chen, D., and Yih, W.-t. (2020). Dense passage retrieval for open-domain question answering. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6769–6781.
- Khattab, O. and Zaharia, M. (2020). Colbert: Efficient and effective passage search via contextualized late interaction over bert. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 39–48.
- Kingma, D. P. and Welling, M. (2014). Auto-encoding variational Bayes. In *Proceedings of the International Conference on Learning Representations (ICLR)*.
- Kraska, T., Beutel, A., Chi, E. H., Dean, J., and Polyzotis, N. (2018). The case for learned index structures. In *Proceedings of the 2018 International Conference on Management of Data (SIGMOD)*, pages 489–504.
- Kusupati, A., Bhatt, G., Rege, A., Wallingford, M., Sinha, A., Ramanujan, V., Howard-Snyder, W., Chen, K., Kakade, S., Jain, P., and Farhadi, A. (2022). Matryoshka representation learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 35, pages 30233–30249.

- Kwiatkowski, T., Palomaki, J., Redfield, O., Collins, M., Parikh, A., Alberti, C., Epstein, D., Polosukhin, I., Devlin, J., Lee, K., Toutanova, K., Jones, L., Kelcey, M., Chang, M.-W., Dai, A. M., Uszkoreit, J., Le, Q., and Petrov, S. (2019). Natural questions: A benchmark for question answering research. volume 7, pages 453–466.
- Lewis, P., Perez, E., Piktus, A., Petroni, F., Karpukhin, V., Goyal, N., Küttler, H., Lewis, M., Yih, W.-t., Rocktäschel, T., Riedel, S., and Kiela, D. (2020). Retrieval-augmented generation for knowledge-intensive NLP tasks. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 33, pages 9459–9474.
- Li, Z., Zhang, X., Zhang, Y., Long, D., Xie, P., and Zhang, M. (2023). Towards general text embeddings with multi-stage contrastive learning. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 2575–2591.
- Lin, C.-Y. (2004). ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81.
- Lin, S.-C., Choi, M., Jung, C., Seo, D., and Lin, J. (2023). How to train your dragon: Diverse augmentation towards generalizable dense retrieval. In *Findings of the Association for Computational Linguistics: EMNLP 2023*.
- Malkov, Y. A. and Yashunin, D. A. (2020). Efficient and robust approximate nearest neighbor search using hierarchical navigable small world graphs. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 42(4):824–836.
- Mehta, S. V., Pasupat, P., Tay, Y., Mao, J., Gupta, J., Guu, K., Kumar, A., Ni, J., Li, A., Cheng, H., Metzler, D., and Chen, L. (2023). DSI++: Updating transformer memory with new documents. *arXiv preprint arXiv:2212.09744*.
- Mu, J., Li, X. L., and Goodman, N. (2023). Learning to compress prompts with gist tokens. In *Advances in Neural Information Processing Systems (NeurIPS)*.

- Rajpurkar, P., Zhang, J., Lopyrev, K., and Liang, P. (2016). SQuAD: 100,000+ questions for machine comprehension of text. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2383–2392.
- Ram, O., Levine, Y., Dalmedigos, I., Muhlgay, D., Shashua, A., Leyton-Brown, K., and Shoham, Y. (2023). In-context retrieval-augmented language models. *Transactions of the Association for Computational Linguistics*, 11:1316–1331.
- Reimers, N. and Gurevych, I. (2019). Sentence-BERT: Sentence embeddings using siamese BERT-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3982–3992.
- Sanh, V., Debut, L., Chaumond, J., and Wolf, T. (2019). DistilBERT, a distilled version of BERT: Smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*.
- Santhanam, K., Khattab, O., Saad-Falcon, J., Potts, C., and Zaharia, M. (2022). ColBERTv2: Effective and efficient retrieval via lightweight late interaction. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT)*, pages 3715–3734.
- Sun, Z., Yu, H., Song, X., Liu, R., Yang, Y., and Zhou, D. (2020). MobileBERT: A compact task-agnostic BERT for resource-limited devices. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 2158–2170.
- Tay, Y., Tran, V., Dehghani, M., Ni, J., Bahri, D., Mehta, H., Qin, Z., Hui, K., Zhao, Z., Gupta, J., Schuster, T., Cohen, W. W., and Metzler, D. (2022). Transformer memory as a differentiable search index.

- Wang, W., Wei, F., Dong, L., Bao, H., Yang, N., and Zhou, M. (2020). MiniLM: Deep self-attention distillation for task-agnostic compression of pre-trained transformers. *arXiv preprint arXiv:2002.10957*.
- Wang, Y., Chen, Y., Zhou, X., Han, Z., Li, J., and Xiao, Y. (2023). A comprehensive survey on vector database: Storage and retrieval technique, challenge. *arXiv preprint arXiv:2310.11703*.
- Xiao, G., Lin, J., Seznec, M., Wu, H., Demouth, J., and Han, S. (2023). SmoothQuant: Accurate and efficient post-training quantization for large language models. In *Proceedings of the 40th International Conference on Machine Learning (ICML)*.
- Xiong, L., Xiong, C., Li, Y., Tang, K.-F., Liu, J., Bennett, P., Ahmed, J., and Overwijk, A. (2021). Approximate nearest neighbor negative contrastive learning for dense text retrieval. In *Proceedings of the International Conference on Learning Representations (ICLR)*.
- Yang, Z., Qi, P., Zhang, S., Bengio, Y., Cohen, W. W., Salakhutdinov, R., and Manning, C. D. (2018). HotpotQA: A dataset for diverse, explainable multi-hop question answering. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2369–2380.
- Zhang, T., Kishore, V., Wu, F., Weinberger, K. Q., and Artzi, Y. (2020). BERTscore: Evaluating text generation with BERT. In *Proceedings of the International Conference on Learning Representations (ICLR)*.

Appendix A

Complete Experimental Results

This appendix presents complete experimental results for all methods across all scales. Tables show mean values across 3 iterations. All quality metrics use 3-decimal precision. Storage values calculated using formula $N \times d \times p$ where N is dataset size, d is dimensions, p is bytes per element (2 for FP16, 4 for FP32). Zero-shot baseline appears once as it’s independent of dataset scale.

A.1 Quality Metrics Across All Scales

Table A.1 presents complete quality metrics (ROUGE-1, BERTScore, ROUGE-L) for all methods across six dataset scales.

Table A.1: Complete quality results: All methods across all scales (mean of 3 iterations)

Method	Scale	ROUGE-1	BERTScore	ROUGE-L
<i>Neural Compression</i>				
Neural-32-FP16	20K	0.168	0.858	0.149
	40K	0.157	0.842	0.145
	60K	0.185	0.853	0.176
	80K	0.183	0.853	0.172
	100K	0.173	0.850	0.163
	120K	0.172	0.849	0.164
Neural-32-FP32	20K	0.175	0.859	0.154
	40K	0.159	0.843	0.148
	60K	0.186	0.853	0.175
	80K	0.182	0.853	0.172
	100K	0.172	0.850	0.163

Continued on next page...

Table A.1 – *Continued from previous page*

Method	Scale	ROUGE-1	BERTScore	ROUGE-L
	120K	0.172	0.849	0.164
Neural-64-FP16	20K	0.174	0.858	0.152
	40K	0.165	0.844	0.152
	60K	0.190	0.854	0.180
	80K	0.183	0.853	0.171
	100K	0.171	0.850	0.162
	120K	0.171	0.849	0.163
Neural-64-FP32	20K	0.170	0.858	0.154
	40K	0.163	0.843	0.153
	60K	0.189	0.854	0.176
	80K	0.183	0.853	0.173
	100K	0.172	0.850	0.163
	120K	0.173	0.849	0.165
<i>Full-Dimensional Baselines</i>				
FAISS	20K	0.177	0.862	0.153
	40K	0.155	0.843	0.140
	60K	0.170	0.851	0.156
	80K	0.166	0.852	0.152
	100K	0.160	0.849	0.147
	120K	0.162	0.848	0.151
HNSW	20K	0.160	0.859	0.141
	40K	0.139	0.839	0.127
	60K	0.153	0.848	0.141
	80K	0.149	0.848	0.138
	100K	0.141	0.845	0.131
	120K	0.146	0.844	0.135
ScaNN	20K	0.168	0.861	0.148
	40K	0.140	0.840	0.127
	60K	0.155	0.848	0.143
	80K	0.151	0.849	0.140
	100K	0.143	0.845	0.132

Continued on next page...

Table A.1 – *Continued from previous page*

Method	Scale	ROUGE-1	BERTScore	ROUGE-L
	120K	0.146	0.845	0.136
<i>Compression Baselines</i>				
PCA-32	20K	0.145	0.855	0.126
	40K	0.124	0.837	0.113
	60K	0.140	0.845	0.130
	80K	0.138	0.846	0.128
	100K	0.135	0.844	0.124
	120K	0.132	0.842	0.123
PCA-64	20K	0.145	0.855	0.126
	40K	0.124	0.837	0.113
	60K	0.140	0.845	0.130
	80K	0.138	0.846	0.128
	100K	0.135	0.844	0.124
	120K	0.132	0.842	0.123
<i>Learned Embeddings</i>				
Matryoshka-32	20K	0.149	0.857	0.133
	40K	0.130	0.838	0.119
	60K	0.141	0.846	0.131
	80K	0.139	0.846	0.130
	100K	0.137	0.844	0.127
	120K	0.139	0.843	0.130
Matryoshka-64	20K	0.149	0.857	0.133
	40K	0.130	0.838	0.120
	60K	0.142	0.846	0.131
	80K	0.140	0.846	0.130
	100K	0.137	0.844	0.127
	120K	0.139	0.843	0.130
<i>No Retrieval Control</i>				
Zero-Shot	All scales	0.063	0.820	0.057

A.2 Efficiency Metrics Across All Scales

Table A.2 presents storage, throughput, and latency for all methods across scales.

Table A.2: Complete efficiency results: Storage, throughput, and latency across all scales

Method	Scale	Storage (MB)	Throughput (QPS)	Latency (ms)
<i>Neural Compression</i>				
Neural-32-FP16	20K	3.2	7,950	0.13
	40K	4.1	7,908	0.13
	60K	5.1	7,842	0.13
	80K	6.0	7,891	0.13
	100K	6.9	7,822	0.13
	120K	9.6	7,861	0.13
Neural-32-FP32	20K	8.0	7,854	0.13
	40K	9.8	7,835	0.13
	60K	11.6	7,863	0.13
	80K	13.4	7,745	0.13
	100K	15.3	7,917	0.13
	120K	17.4	7,968	0.13
Neural-64-FP16	20K	4.1	7,961	0.13
	40K	5.9	7,945	0.13
	60K	7.7	7,919	0.13
	80K	9.5	7,777	0.13
	100K	11.4	7,882	0.13
	120K	14.9	7,856	0.13
Neural-64-FP32	20K	9.9	7,960	0.13
	40K	13.5	7,967	0.13
	60K	17.0	7,795	0.13
	80K	20.6	7,830	0.13
	100K	24.2	7,852	0.13
	120K	28.0	7,838	0.13

Full-Dimensional Baselines

Continued on next page...

Table A.2 – Continued from previous page

Method	Scale	Storage (MB)	Throughput (QPS)	Latency (ms)
FAISS	20K	44.3	2,130	0.47
	40K	88.0	529	1.94
	60K	131.5	355	2.86
	80K	175.6	249	4.39
	100K	220.3	196	5.37
	120K	263.8	151	6.84
HNSW	20K	23.2	2,628	0.40
	40K	45.8	2,452	0.43
	60K	68.3	2,315	0.45
	80K	91.2	2,236	0.46
	100K	114.8	2,242	0.46
	120K	137.2	2,244	0.46
ScaNN	20K	23.2	1,955	0.51
	40K	45.8	5,269	0.19
	60K	68.3	3,656	0.28
	80K	91.2	2,800	0.36
	100K	114.8	2,269	0.44
	120K	137.7	1,913	0.52
<i>Compression Baselines</i>				
PCA-32	20K	1.2	15,444	0.07
	40K	2.4	14,704	0.07
	60K	3.7	14,127	0.08
	80K	4.9	13,243	0.09
	100K	6.1	13,027	0.09
	120K	7.3	12,316	0.10
PCA-64	20K	2.4	15,436	0.07
	40K	4.9	14,772	0.07
	60K	7.3	14,112	0.08
	80K	9.8	13,276	0.08
	100K	12.2	13,133	0.09
	120K	14.7	12,742	0.09

Continued on next page...

Table A.2 – Continued from previous page

Method	Scale	Storage (MB)	Throughput (QPS)	Latency (ms)
<i>Learned Embeddings</i>				
Matryoshka-32	20K	87.6	-	-
	40K	87.6	-	-
	60K	87.6	-	-
	80K	87.6	-	-
	100K	87.6	-	-
	120K	87.6	-	-
Matryoshka-64	20K	87.6	-	-
	40K	87.6	-	-
	60K	87.6	-	-
	80K	87.6	-	-
	100K	87.6	-	-
	120K	87.6	-	-
<i>No Retrieval Control</i>				
Zero-Shot	All scales	-	-	-

Storage: Neural (index + 1.7 MB models), FAISS/HNSW/ScaNN (measured),

PCA ($N \times d \times 2$ bytes FP16), Matryoshka (87.6 MB constant). Dash (-) indicates not measured.

A.3 Statistical Summary

Key observations from complete results:

Quality trends. Neural-32-FP16 improves from 0.168 (20K) to 0.172 (120K), +2.5% gain—only method with positive trend. All baselines degrade: FAISS -8.1%, HNSW -8.8%, ScaNN -12.9%, PCA-32 -8.9%, Matryoshka-32 -6.6%. Neural-64 variants show slight degradation (-1.4% for FP16, -1.7% for FP32 from peak), suggesting 32D benefits more from scale than 64D.

Storage scaling. Neural storage grows linearly with dataset size as expected (~ 0.066 MB per 1K samples for 32-FP16). FAISS grows at 2.20 MB per 1K samples—33 \times faster growth rate. At 120K scale, FAISS requires 263.8 MB versus Neural’s 9.6 MB (27.5 \times difference). PCA theoretically smallest (7.3 MB at 120K) but 30.9% quality gap makes comparison misleading.

Throughput stability. Neural maintains 7,800-7,970 QPS across all scales (coefficient of variation 0.9%). FAISS degrades from 2,130 to 151 QPS (-92.9%)—exponential collapse. HNSW more stable (2,628 to 2,244, -14.6%). PCA degrades despite low dimensionality (15,444 to 12,316, -20.3%), suggesting cache effects and FAISS overhead dominate even at 32D/64D.

Crossover at 40K. Neural surpasses FAISS at 40K samples (0.157 vs 0.155, +1.3% margin). Margin widens at 60K (+8.8%) then stabilizes around 6-7% at larger scales. This earlier-than-expected crossover demonstrates FAISS’s quality degradation precedes algorithmic transition to approximation.

Zero-shot performance. Achieves 0.063 ROUGE-1—110% worse than worst retrieval method (PCA at 0.132). Validates temporal knowledge gap: Flan-T5 cannot answer 2024 questions without retrieval.

Appendix B

Empirical Analysis: Semantic Similarity vs. Retrieval Relevance

This appendix provides empirical evidence that semantic similarity does not reliably predict retrieval relevance. This supports the observation in Section 5.2 that “larger databases may contain more near-duplicates or semantically similar but contextually different answers, confusing retrieval.”

B.1 Executive Summary

We conducted empirical analysis on two data sources to quantify the mismatch between semantic similarity and retrieval relevance. **Important:** HotpotQA is used solely to validate this phenomenon on an external benchmark—no neural compression experiments were conducted on HotpotQA. All thesis experiments (Chapters 3–5) use the custom QA datasets.

Table B.1: Summary of similarity vs. relevance analysis

Source	Questions	Pairs	Inversion Rate	Key Metric
HotpotQA (benchmark)	1,000	9,927	69.2%	$\rho = 0.379$
Custom QA (6 domains)	3,000	3,000	80.3%	Margin = -0.078

Key finding: In 69–80% of questions, at least one irrelevant document ranks higher than a relevant document based on similarity alone. In the custom datasets, wrong answers have systematically *higher* similarity than correct answers (negative margin), demonstrating that similarity-based retrieval is fundamentally unreliable.

B.2 HotpotQA Analysis

B.2.1 Dataset Description

Note: This analysis uses HotpotQA solely to demonstrate the similarity-relevance mismatch phenomenon on an external benchmark. No neural compression experiments were conducted on HotpotQA; this section provides statistical validation only.

HotpotQA (Yang et al., 2018) is a standard multi-hop question answering benchmark. We analyzed the dev_distractor split:

- **Questions analyzed:** 1,000 (for similarity statistics), 500 (for rank inversion analysis)
- **Question-document pairs:** 9,927
- **Ground truth:** supporting_facts field identifies relevant documents
- **Embedding model:** all-MiniLM-L6-v2 (384D)—same as thesis experiments

B.2.2 Similarity Distribution by Relevance

Table B.2: Similarity statistics for relevant vs. irrelevant documents (HotpotQA)

Metric	Relevant Docs	Irrelevant Docs
Count	2,000	7,927
Mean Similarity	0.530	0.360
Std Dev	0.163	0.157
Median	0.544	0.357
Min	0.024	-0.102
Max	0.955	0.911

While relevant documents have higher mean similarity (0.530 vs. 0.360), the distributions overlap substantially. Notably, some irrelevant documents achieve similarity as high as 0.911, while some relevant documents have similarity as low as 0.024.

Table B.3: Similarity-relevance mismatch cases (HotpotQA)

Category	Count	Percentage
High similarity (≥ 0.6) but irrelevant	538	5.4%
Low similarity (< 0.4) but relevant	439	4.4%
Total mismatch pairs	977	9.8%

Table B.4: Correlation between similarity and relevance (HotpotQA)

Metric	Value	Interpretation
Spearman ρ	0.379	Weak positive correlation
p -value	< 0.001	Statistically significant
Questions with rank inversion	346/500 (69.2%)	Majority have at least one irrelevant doc ranked above relevant

B.2.3 Mismatch Frequency

B.2.4 Correlation Analysis

Despite statistical significance, the correlation is *weak* ($\rho = 0.38$). In 69.2% of questions analyzed for rank inversion, at least one irrelevant document ranks higher than a relevant document based on similarity alone.

B.3 Custom QA Dataset Analysis

B.3.1 Datasets Analyzed

We analyzed question-answer pairs from our custom distilled datasets spanning six domains:

Table B.5: Custom QA datasets analyzed

Dataset	Domain	Pairs Analyzed
AI Tech	Technology, AI/ML	500
Entertainment	Movies, Music, Pop Culture	500
Finance	Economics, Markets, Banking	500
Geopolitical	Politics, International Relations	500
Science	Scientific Discoveries, Research	500
Sports	Athletics, Competitions	500
Total	6 domains	3,000

B.3.2 Methodology

For each question Q_i with correct answer A_i , we computed:

1. **Correct similarity:** $\text{sim}(Q_i, A_i)$
2. **Max wrong similarity:** $\max_{j \neq i} \text{sim}(Q_i, A_j)$
3. **Margin:** Correct similarity – Max wrong similarity
4. **Rank inversion:** Occurs when max wrong > correct

A negative margin indicates wrong answers have higher similarity than the correct answer—a systematic failure of similarity-based retrieval.

B.3.3 Per-Dataset Results

Table B.6: Similarity analysis results by domain

Dataset	N	Correct Sim.	Max Wrong	Margin	Inversion
AI Tech	500	0.652	0.787	−0.134	91.8%
Entertainment	500	0.721	0.746	−0.025	67.2%
Finance	500	0.663	0.745	−0.082	70.4%
Geopolitical	500	0.654	0.747	−0.093	88.8%
Science	500	0.693	0.752	−0.058	76.8%
Sports	500	0.682	0.754	−0.072	86.6%
Average	3,000	0.678	0.755	−0.078	80.3%

B.3.4 Key Observations

1. **Negative margin across all datasets:** The margin (correct – max wrong) is negative for every domain, meaning wrong answers have higher similarity than correct answers on average.
2. **80.3% average inversion rate:** In 80% of questions, at least one wrong answer ranks higher than the correct answer based on similarity.
3. **Domain variation:** AI Tech (91.8%) and Geopolitical (88.8%) show highest inversion rates; Entertainment (67.2%) is lowest but still shows majority inversions.
4. **Systematic failure:** This is not an edge case—it is the norm across all domains.

B.4 Combined Interpretation

Both analyses demonstrate that semantic similarity is an unreliable proxy for retrieval relevance:

- **HotpotQA (external benchmark, statistics only):** Weak correlation ($\rho = 0.38$), 69% of questions have rank inversions. This validates the phenomenon on a widely-used benchmark independent of our thesis experiments.
- **Custom datasets (thesis experimental data):** Negative margin (-0.078), 80% of questions have rank inversions. These are the datasets used for all neural compression experiments in this thesis.

The custom datasets show even more severe mismatch than HotpotQA, likely due to high semantic overlap within domain-specific corpora. Documents sharing domain terminology achieve high similarity regardless of whether they answer the specific question.

These findings validate the thesis motivation: compression methods that preserve only semantic similarity may not preserve retrieval usefulness. Task-specific neural compression learns which dimensions matter for actual question-answer matching, rather than generic semantic similarity.

B.5 Methodology Details

- **Embedding model:** all-MiniLM-L6-v2 (Sentence Transformers, 384D)
- **Similarity metric:** Cosine similarity
- **Analysis pipeline:** (1) Encode all texts using Sentence Transformer, (2) Compute pairwise cosine similarity, (3) Compare similarity of correct vs. incorrect pairings, (4) Calculate inversion rate and correlation metrics