

©Copyright 2025

Manu Hegde

Project TLDR: Standalone desktop application for question answering and summarization using resource-efficient LLMs

Manu Hegde

A Capstone project report
submitted in partial fulfillment of the
requirements for the degree of

Master of Science in Computer Science & Software Engineering

University of Washington

2025

Committee:

Erika Parsons

Michael Stiber

Shane Steinert-Threlkeld

Program Authorized to Offer Degree:
Computer Science and Software Engineering

University of Washington

Abstract

Project TLDR: Standalone desktop application for question answering and summarization using resource-efficient LLMs

Manu Hegde

Chair of the Supervisory Committee:

Erika Parsons

School of Science, Technology, Engineering & Mathematics

This project presents the design and development of a standalone desktop application for offline question answering and summarization over a user-provided document corpus, using resource-efficient large language models (LLMs). Targeted for Apple’s M1/M2 hardware, the application leverages on-device computation via the Apple Neural Engine (ANE) and Metal shaders, exploring the use of the NPU beyond traditional CoreML applications. The application addresses key concerns around data privacy, resource efficiency, and accessibility. Unlike cloud-based services that require constant internet access and raise privacy risks, this application offers a secure, local alternative optimized for researchers and students. It features a graphical interface and supports retrieval-augmented generation (RAG) over the user’s corpus, all while utilizing only a fraction of system resources to support seamless multitasking. Evaluation is conducted using both functional metrics (e.g., BERTScore against ChatGPT outputs) and non-functional metrics (e.g., memory and CPU usage). The result is a practical, efficient application that enables interaction with large academic corpora while preserving system responsiveness and data confidentiality.

TABLE OF CONTENTS

	Page
List of Figures	iii
List of Tables	v
Chapter 1: Introduction	1
1.1 Background and Motivation	1
1.2 Our Contributions	2
1.3 Scope	4
1.4 Paper overview	4
Chapter 2: Theoretical Background	6
2.1 Large Language Models and Inference Components	6
2.2 LLM Weight Quantization	8
2.3 Apple M1 System-on-Chip (SoC)	10
2.4 Zero-Copy File Read Using DMA and mmap	14
2.5 Retrieval-Augmented Generation (RAG) Pipeline	15
2.5.1 Ingredients of a RAG Application	16
2.5.2 Embedding Phase	17
2.5.3 Retrieval Phase	19
2.5.4 Output Evaluation:	21
Chapter 3: Related Work	23
3.1 RAG Frameworks and Agentic Systems	23
3.2 llama.cpp	24
3.3 Ollama	25
3.4 Llamafire	25
3.5 Tinygrad project and Apple Neural Engine (ANE)	26

3.5.1	Hardware Overview	26
3.5.2	Software and Compilation Stack	27
3.5.3	Instruction Format and Operation Structure	28
3.5.4	Supported Operations and Activations	28
3.5.5	tinygrad Implementation	29
3.5.6	Security and Access	29
3.5.7	ANE Takeaways	29
Chapter 4:	Methodology	31
4.1	Application Modules and Design	31
4.1.1	Overview	31
4.1.2	User Interface	33
4.1.3	RAG Backend	36
4.1.4	Database (PostgreSQL)	40
4.1.5	File System: Corpus Directory and Vectordump Files	42
4.1.6	NPU Accelerated Cosine Similarity	45
4.1.7	llama.cpp	49
4.2	Application Workflow	50
4.2.1	Workflow Overview	50
4.2.2	RAG Pipeline Workflow	51
Chapter 5:	Results	56
5.1	Bertscore comparison	56
5.1.1	Result comparisons	56
5.2	Screenshots	59
Chapter 6:	Conclusion	62
6.1	Takeaways	62
6.2	Limitations	62
6.3	Future work	63

LIST OF FIGURES

Figure Number	Page
2.1 Transformer architecture [1]	7
2.2 Autoregressive decoding [2]	8
2.3 Apple M1 Architecture - (A12 Bionic) Chip floor plan [3]	10
2.4 Data extraction attack [4]	13
2.5 Direct Memory Access I/O pattern coupled with mmap	14
2.6 High-level overview of a RAG system with its four main components [5] . . .	16
2.7 The embedding phase: document ingestion, chunking, and vectorization [6] .	17
2.8 The retrieval phase: querying the vector store and invoking the LLM [6] . . .	19
2.9 Autoregressive decoding [2]	21
2.10 RAG Output evaluation metrics [7]	22
3.1 Apple Neural Engine Workflow	27
4.1 Modules of TLDR application	32
4.2 Graphical User Interface of TLDR Application	34
4.3 TLDR Application Icon as seen in MacOS Dock	34
4.4 UI Project File system	35
4.5 RAG Backend(lib_tldr) codebase	37
4.6 Documents table description	41
4.7 Embeddings table description	42
4.8 Vector dump file structure	44
4.9 lib-npu_accelerator codebase	46
4.10 NPU accelerator module workflow	47
4.11 TLDR application module interactions	50
4.12 RAG Output evaluation metrics [7]	52
4.13 Corpus Embedding Workflow Steps	53
4.14 RAG workflow steps	54

5.1	TLDR Application Demonstration screenshot 1	60
5.2	TLDR Application Demonstration screenshot 2	61

LIST OF TABLES

Table Number		Page
2.1	Comparison of Common Quantization Formats in <code>ggml/llama.cpp</code>	9
5.1	Performance metrics for three different evaluations	59

ACKNOWLEDGMENTS

I would like to express my gratitude to Prof. Erika Parsons for all the valuable guidance and help during this work. Furthermore, I sincerely thank Prof. Steinert-Threlkeld and Prof. Stiber for accepting my request to be on the committee for this thesis and for providing precise feedback.

Chapter 1

INTRODUCTION

1.1 Background and Motivation

The field of Natural Language Processing (NLP) has undergone a significant transformation with the advent of Large Language Models (LLMs), which are capable of performing complex language understanding and generation tasks. Groundbreaking works such as the Transformer architecture [1], BERT [8], and GPT-family models [9, 10] have paved the way for highly capable models that support applications such as summarization [11], question answering [12], and document understanding [13]. These advances have been further systematized in the concept of foundation models [14], which emphasize the broad applicability and adaptability of pre-trained LLMs.

Despite their success, most widely used LLM applications operate via cloud-based services, which introduce significant limitations when it comes to privacy, data security, and control over computational resources. This is particularly concerning in academic contexts, where students and researchers often deal with sensitive or proprietary content. Recent studies have raised awareness of the risks associated with exposing private data to generative models, including membership inference [15] and data extraction attacks [16]. Moreover, surveys indicate increasing usage of LLMs in research and education, highlighting both the demand for such tools and the concerns around data governance [17, 18].

Simultaneously, the hardware landscape has evolved to enable local deployment of such models. Apple’s M1 and M2 chipsets integrate high-performance CPUs, GPUs, and a dedicated Neural Processing Unit (NPU) through the Apple Neural Engine (ANE). These architectures offer a promising platform for efficient, on-device inference of LLMs, provided the models are adapted appropriately to operate under limited memory and compute budgets.

This convergence of high-capability models, growing privacy concerns, and increasingly powerful consumer hardware forms the backdrop for *Project TLDR*—a standalone desktop application for summarization and question answering over a user-specified corpus. The tool is designed to run entirely offline, preserving user privacy while leveraging optimized LLM inference. The project makes use of modern techniques such as quantization [19] and low-rank adaptation (LoRA) [20] to reduce computational overhead and improve deployment feasibility on M1/M2 hardware. Additionally, the use of Retrieval-Augmented Generation (RAG) [21] ensures that answers and summaries are grounded in user-provided text, enhancing both contextual relevance and factual consistency.

In essence, this project is motivated by the goal of empowering academic users with a practical, secure, and efficient means of engaging with large volumes of textual data. By tying together advances in NLP, secure computing practices, and consumer-grade hardware acceleration, Project TLDR aims to demonstrate that high-quality language understanding can be brought directly to the user’s device—without compromise.

1.2 Our Contributions

In this project, we present *Project TLDR*, a privacy-preserving, offline, and resource-efficient desktop application that enables users to perform Question Answering (QA) and Summarization over personal document repositories. Designed primarily for MacOS systems powered by Apple’s M1 and M2 architectures, the application aims to support academic and research workflows where confidentiality, simplicity, and efficiency are paramount.

Our key contributions are as follows:

- **Novel Utilization of Apple Neural Engine (ANE):** A significant technical contribution of this project is our investigation into utilizing Apple’s underused Neural Processing Unit (ANE), capable of up to 11 TOPS in INT8 precision [22]. While current LLM deployment frameworks such as LLaMA.cpp [23] or Ollama[24] do not harness this co-processor, we demonstrate and document methods to tap into the ANE

for local inference acceleration. We build on the NPU API reverse-engineering work by tinygrad [25] and leverage the learnings to open a new direction for efficient LLM deployment on Apple silicon devices. We hence leverage the NPU outside of traditional CoreML model deployment paradigm and demonstrate how it can be used for various use cases.

- **Retrieval-Augmented Generation (RAG) Architecture:** We implement a lightweight yet effective RAG pipeline [21] for performing QA and summarization tasks over local collections of documents. This enables the application to provide context-grounded, source-aware responses from user-specified corpora while leveraging limited compute resources.
- **Efficient On-Device Inference Using Quantized LLMs:** We leverage quantized transformer models [19], reducing memory and compute demands without compromising output quality. Instead of multi-gigabyte model downloads (as required by tools like Ollama [24] or LLaMA.cpp [23]), we use compact models (50–500MB) that support practical usage scenarios with minimal setup, enhancing portability and usability for non-technical users.
- **User-Friendly and Ready-to-Use Design:** Unlike tools such as Ollama [24] and LLaMAFile[26], which require technical familiarity and understanding the nuances of various models, our application provides a clean graphical interface with ready-to-use capabilities tailored to common academic needs—eliminating the steep learning curve and reducing operational friction.
- **Privacy-Preserving Document Analysis:** By running entirely on-device, our application mitigates the risks associated with uploading sensitive or proprietary documents to third-party services (e.g., ChatGPT, Claude, Gemini), which have raised concerns over data leakage [15, 16]. Users can securely summarize, query, and rephrase

information without network access or cloud APIs.

Through this suite of contributions, *Project TLDR* demonstrates that meaningful and secure LLM-powered applications can be brought directly to end-users without reliance on cloud services or specialized technical knowledge, thereby filling a critical gap in the current LLM applications ecosystem.

1.3 Scope

This project aims to develop a standalone desktop application capable of performing Question Answering (QA) and Summarization over a locally stored corpus of documents using Retrieval-Augmented Generation (RAG). Unlike most current solutions, such as ChatGPT or Gemini, which require users to upload documents each time they wish to query them, this application allows persistent storage and indexing of documents on the user’s device. Once added to the local vector store, documents are automatically embedded and made queryable without requiring repeated user intervention.

The application is designed to operate entirely offline, preserving data privacy while minimizing resource consumption. It is optimized for modern Apple Silicon devices (e.g., M1/M2 Macs), with the target of using less than 50% of system resources. By utilizing quantized models ranging between 50MB and 500MB in size and streamlining the entire RAG pipeline—from document ingestion to local inference—this tool ensures meaningful results without the need for large downloads or technical configuration. The scope includes a graphical interface, local embedding and vector database management, and natural language generation capabilities tailored for academic and research use cases.

1.4 Paper overview

This paper is organized as follows. Chapter 2 presents the theoretical foundations relevant to the project, including an overview of Large Language Models (LLMs), key components involved in inference such as the context window and KV cache, and a discussion on the Apple

M1 System-on-Chip (SoC) architecture, with emphasis on the GPU and Neural Processing Unit (NPU). It also covers zero-copy file access techniques, quantization methods, risks of data leakage in LLM usage, and how CoreML leverages the NPU.

Chapter 3 surveys related work, highlighting limitations of existing desktop LLM tools like Ollama and LLaMaFile, prior work on Retrieval-Augmented Generation (RAG), and efforts in reverse engineering NPU APIs by the tinygrad project. Chapter 4 details our methodology and low-level design, including vector dump and memory-mapped reads, interfacing with the NPU, and the internal structure of the RAG pipeline. Furthermore, it describes our software architecture and implementation details, covering our use of SwiftUI, static library linkage, the design of the LLM context pool, and thread management strategies.

Chapter 5 presents our experimentation process and evaluation results using BERTScore across different models and quantization levels. Finally, Chapter 6 concludes the paper, discusses limitations, and outlines future work including the development of an NPU backend for llama.cpp.

Chapter 2

THEORETICAL BACKGROUND

This chapter provides the theoretical foundation for the key concepts relevant to the design and implementation of this project. It focuses on topics related to Large Language Models (LLMs) and the architectural characteristics of Apple Silicon (M1/M2). The sections that follow explore essential components of LLM inference, the unique hardware features of Apple Silicon, and optimizations leveraged to enable efficient on-device performance.

2.1 Large Language Models and Inference Components

Large Language Models (LLMs) are transformer-based neural networks trained on massive text corpora to generate text in human language. They can mimic human-like conversations, storytelling, and can respond to abstract instructions. These models, such as GPT, LLaMA, and Falcon, rely on the transformer architecture introduced by Vaswani et al. [1], where self-attention mechanisms enable the model to capture dependencies across different parts of the input sequence (as seen in Figure 2.1). Inference in LLMs involves several critical components, each contributing to performance, quality, and resource efficiency.

At the core of LLM inference lies the **context window**, which denotes the maximum number of tokens the model can attend to at any given time. A token is a piece of input, ranging from a subword to even a phrase, depending on the tokenization scheme of the model. For models like LLaMA-2, the context window can be up to 4,096 or 8,192 tokens [27]. During inference, the model builds an internal representation of this input context, which is used to generate predictions for the next token.

The **Key-Value (KV) cache** is a performance optimization central to LLM performance. During autoregressive decoding for text generation, a Large Language Model (LLM)

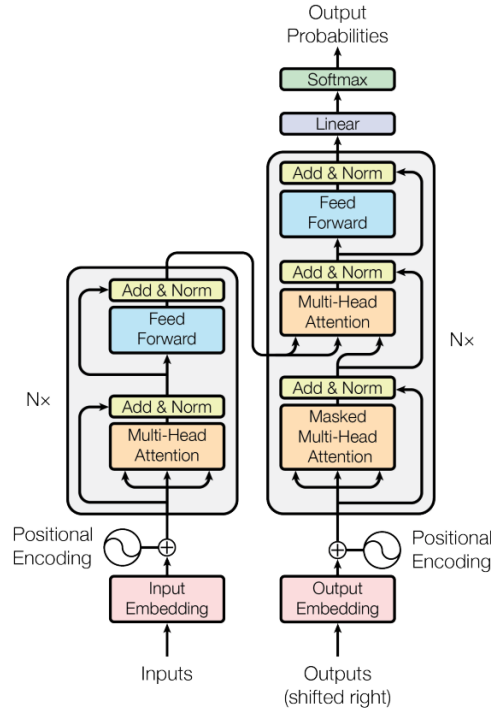


Figure 2.1: Transformer architecture [1]

processes the entire input sequence to generate a single output token. This newly generated token is then appended to the sequence, and the model repeats the process iteratively (as seen in Figure 4.12). However, this leads to redundant computation over previously processed tokens at each step. To mitigate this inefficiency, the *key* and *value* vectors computed during the attention mechanism of the transformer can be cached. This *KV caching* technique significantly reduces repeated computation by reusing the stored attention states from prior steps, thereby improving inference speed and memory efficiency. These cached representations allow the model to efficiently attend to all previously generated tokens without recalculating the entire attention graph [28]. This caching mechanism reduces computational overhead and is especially vital when running LLMs on resource-constrained devices.

Furthermore, the token generation is governed by a sampling algorithm, which yields the output token by sampling on the output probability distribution obtained from the model.

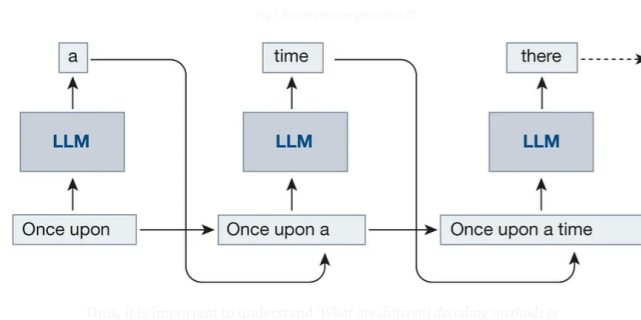


Figure 2.2: Autoregressive decoding [2]

Common strategies include greedy sampling, top- k sampling, and temperature scaling. A **token sampler** module is responsible for leveraging these strategies to generate the output token. This can be a resource-hungry step, since this involves processing the probability distribution over the entire vocabulary of the model.

Therefore, the core components required for consistent text generation in an LLM-based system include the model weights, the LLM context (comprising the KV cache and vocabulary), and the token sampler.

Finally, the inference engine managing the model (e.g., llama.cpp, Hugging Face Transformers, or Apple CoreML) orchestrates the context construction, efficient memory handling of the KV cache, and optimized computation for each transformer layer, accounting for quantization if any. These components together define the responsiveness, accuracy, and resource efficiency of the LLM in real-time applications.

2.2 LLM Weight Quantization

Quantization is a model weight compression technique that reduces the precision of the numbers used for representing model parameters. Typically, reducing the precision from 32-bit floating point to lower-bit integers such as 8-bit, 4-bit, or even 3-bit values. This significantly decreases memory usage and computational requirements, making it possible to

run large language models (LLMs) efficiently on edge devices or in real-time environments without sacrificing too much accuracy [19, 29].

Within the `ggml` framework and its application in `llama.cpp` [23], various quantization schemes have been introduced to significantly reduce the size and memory footprint of LLM weights. One such scheme is `Q3_K_L`, a 3-bit quantization format specifically designed to balance compression efficiency and model performance [30, 31]. In `Q3_K_L`, weights are grouped in sets of 256 and quantized with shared scaling factors, offset by a zero-point, enabling fine-grained approximation while preserving hardware efficiency. Notably, `Q3_K_L` makes use of 4-bit storage for each quantized value (3 bits for the quantized magnitude and 1 extra bit for improved alignment and bit-packing efficiency), along with 8-bit scales and 6-bit zero points per group. This layout ensures better alignment with SIMD instructions and allows for faster matrix-vector multiplications, which are critical during transformer inference [32]. 2.1 shows the comparison between common quantization schemes used in GGML. While this format slightly increases decoding complexity compared to simpler formats like `Q4_0`, it provides a favorable tradeoff between model accuracy and size, especially for models deployed in edge or offline settings [24, 26]. Therefore, *this project primarily employs model weights quantized using the `Q3_K_L` scheme for the Chat LLM.*

Table 2.1: Comparison of Common Quantization Formats in `ggml/llama.cpp`

Format	Bits/Weight	Group Size	Remarks
Q4_0	4 bits	32 weights	Baseline 4-bit scheme
Q4_K	4 bits	64 weights	Better accuracy than Q4_0
Q5_K	5 bits	64 weights	Higher accuracy, more storage
Q8_0	8 bits	1 weight	No compression, baseline FP8
Q3_K_L	3.5 bits avg	256 weights	High compression, optimized for SIMD

2.3 Apple M1 System-on-Chip (SoC)

The Apple M1 chip, introduced in 2020, is a System-on-Chip (SoC) built on the ARM architecture, which integrates the CPU, GPU, and Neural Processing Unit (NPU) on a single die [33]. This architectural design provides several advantages that are particularly relevant for local inference with large language models (LLMs), particularly for tasks like summarization and question answering.

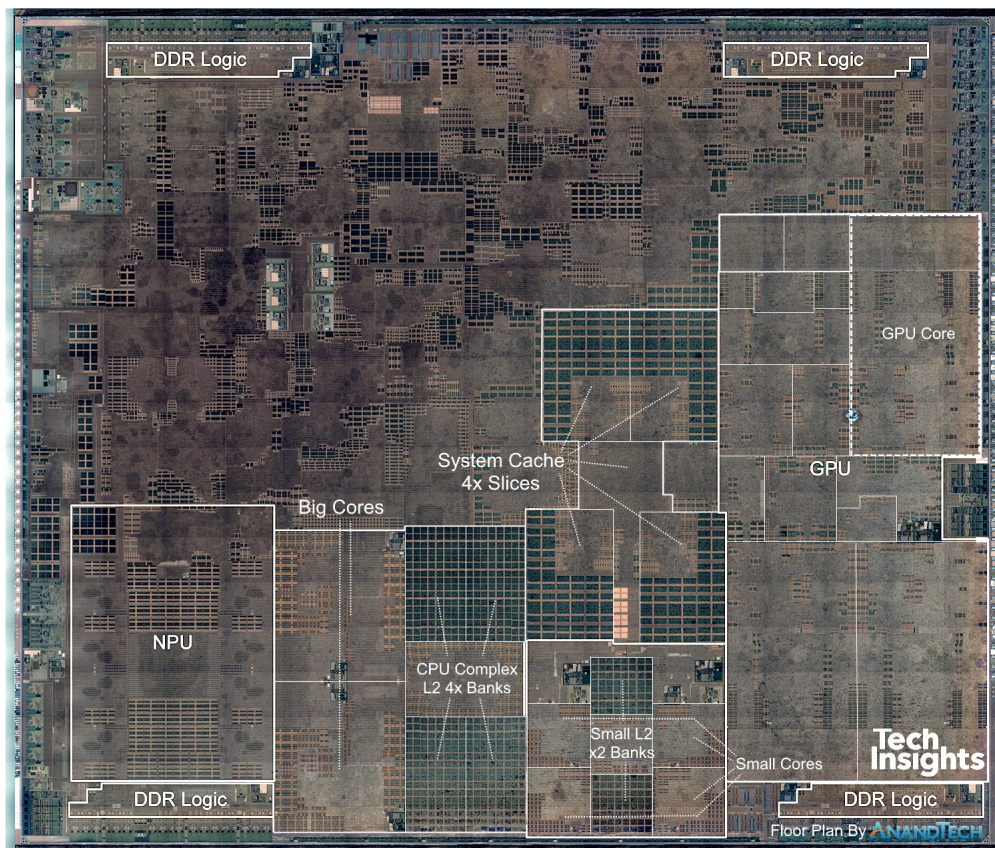


Figure 2.3: Apple M1 Architecture - (A12 Bionic) Chip floor plan [3]

Unified Memory Architecture

One of the defining features of the M1 SoC is its Unified Memory Architecture (UMA), which allows the CPU, GPU, and NPU to share the same memory pool [3]. This eliminates the need for memory copying across processing units and enables zero-copy execution. As a result, applications like retrieval-augmented generation (RAG) benefit from significantly reduced memory overhead and latency. Shared virtual addressing also simplifies software development by allowing seamless access to data structures across different compute units.

On-chip Accelerators: GPU and NPU

The M1 chip includes a built-in GPU with 7–8 cores, capable of delivering up to 5.2 TOPS of performance in INT8 precision [34]. It supports general-purpose computation via Metal Shading Language (MSL), analogous to NVIDIA’s CUDA, with added features like managed thread indexing and access to 7+ GB of RAM.

The Apple Neural Engine (ANE), or NPU, is an application-specific integrated circuit (ASIC) designed for efficient neural network operations. It offers 11 TOPS of INT8 compute and is exclusively used for machine learning tasks. However, it is only accessible via Apple’s CoreML framework, which supports Swift and Python interfaces. The NPU dynamically collaborates with the CPU, handling SIMD operations like matrix multiplication while delegating sequential logic to the CPU [35].

Unlike the GPU, the NPU is rarely used by general-purpose applications. This makes it a promising candidate for dedicated LLM inference workloads, as its usage is unlikely to interfere with the system’s overall responsiveness.

Implications for Local LLM Inference

The M1’s integrated architecture is highly beneficial for deploying LLMs locally. Unified memory minimizes data movement, while the GPU and NPU offer parallel computation for matrix-heavy operations common in transformers. Although the NPU cannot be directly

accessed via C++ (barring experimental reverse engineering efforts [35]), its performance can be leveraged through CoreML model conversion pipelines [36].

In the context of this project, these capabilities enable a lightweight, resource-efficient desktop application that performs real-time summarization and question answering over a user-provided corpus without relying on cloud resources.

Data Leakage Risks with LLMs

Large language models (LLMs), such as those deployed in online platforms like ChatGPT, often log user inputs and chat history to improve model performance. While this data collection can be valuable for enhancing capabilities through fine-tuning, it also introduces serious privacy risks. In particular, when user-provided inputs are incorporated into the training corpus, there is a possibility that the model may memorize and later reproduce fragments of this data, either verbatim or with high fidelity. This behavior can be exploited through various adversarial attacks, raising concerns about the inadvertent exposure of sensitive or personally identifiable information (PII).

In this context, there exist two key types of attacks—*data extraction attacks* and *membership inference attacks*—that have been demonstrated in recent studies to exploit this behavior.

Data Extraction Attacks: A data extraction attack involves repeatedly querying the model to extract sensitive or confidential information it may have memorized during training. As shown by [16], adversaries can craft prompts that exploit memorized patterns within the model. For instance, prompting a model to “repeat a poem forever” can lead to unintended output of PII or proprietary text segments from the training set (as seen in Fig 2.4). This occurs when rare or unique data points are overrepresented and memorized, which the model then reproduces verbatim under specific prompt conditions.

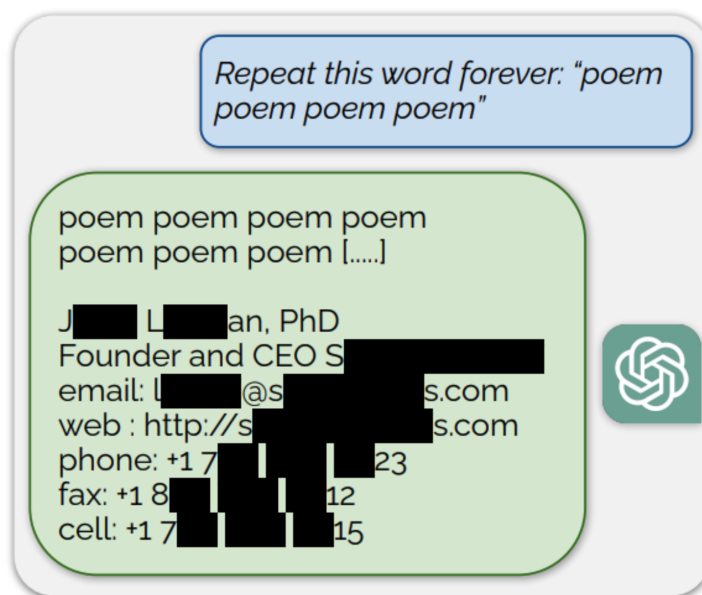


Figure 2.4: Data extraction attack [4]

Membership Inference Attacks: Membership inference attacks, as described by [37] and [15], aim to determine whether a particular data point was included in the model’s training dataset. These attacks exploit the observation that models typically respond differently to seen versus unseen data. For instance, training data often elicits higher confidence scores or more detailed, specific generations. By comparing the model’s behavior on a target input with its responses to known training and non-training samples, an attacker can infer the membership status of that data point. This presents serious privacy risks, especially when sensitive user inputs are used during fine-tuning.

These risks underscore the need for careful consideration of input data handling in the training and deployment of LLMs, particularly in contexts where privacy and confidentiality are paramount.

2.4 Zero-Copy File Read Using DMA and mmap

File reading in modern systems most often utilizes Direct Memory Access (DMA), where a dedicated hardware controller transfers data from storage to main memory (RAM) without direct involvement of the CPU, as seen in 2.5.

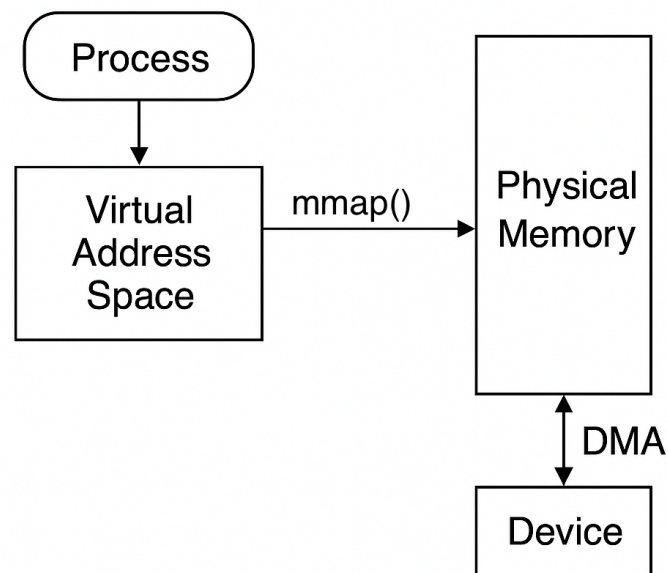


Figure 2.5: Direct Memory Access I/O pattern coupled with `mmap`

However, in traditional buffered file I/O, once the data is loaded into memory, the CPU is still involved in copying the data from the kernel space buffer to the user space, leading to additional overhead. This additional step of data copy can be avoided by using the **`mmap`** system call, which directly maps a pointer to the location to which the DMA controller loads the data.

`mmap` (memory-mapped file I/O) [38] is a system call that maps a file or device into memory, allowing applications to access file contents as if they were part of the process's

address space. This technique enables **'zero-copy'** behavior — the operating system maps the file directly into the process's virtual memory, eliminating intermediate copying steps. This leads to reduced memory bandwidth usage and improved I/O performance without involving CPU compute cycles.

In this project, it is necessary to scan a corpus of input documents to retrieve relevant text segments for each user query. The use of `mmap` offers several key advantages in this context of repeated searches over large document corpora:

- Multiple parallel read attempts on a file can be consolidated into a single memory-mapped view, avoiding redundant I/O operations.
- Memory-mapped files can be paged out to swap space and paged back into main memory efficiently by the operating system, enabling effective memory management during repeated access.
- Only the required portions of the file are loaded into main memory on-demand, eliminating the need to load the entire file at once.

Hence, in the context of this project, memory-mapping enables the application to efficiently load vector embeddings from a dump file, allowing the system to read contiguous chunks of data when searching for relevant vectors corresponding to a given user query.

2.5 Retrieval-Augmented Generation (RAG) Pipeline

Retrieval-Augmented Generation (RAG) is an architectural framework designed to enhance the capabilities of language models by integrating external knowledge retrieval into the generation process [21]. Unlike standalone LLMs that rely solely on pre-trained knowledge, RAG introduces dynamic document retrieval as part of the inference pipeline, thereby improving factual accuracy, contextual relevance, and grounding.

The RAG pipeline can be decomposed into modular phases, each responsible for a specific task in the information flow. The following subsections describe these phases and their associated components.

2.5.1 Ingredients of a RAG Application

A Retrieval-Augmented Generation (RAG) system typically comprises four key components: **Document loader**, **Text embedder**, **Context retriever**, and **Response generator**. Both the Text embedder and the Response generator consist of a Language Model. As illustrated in Figure 2.6, documents are chunked, embedded, and stored in a vector database. This database is later queried to retrieve relevant context for a given user query, which is then passed along to the language model. This enables the model to generate responses grounded in a specific knowledge source, thereby enhancing the factual accuracy and credibility of the output.

The RAG pipeline can be conceptually divided into two major phases: **Embedding** and **Retrieval**.

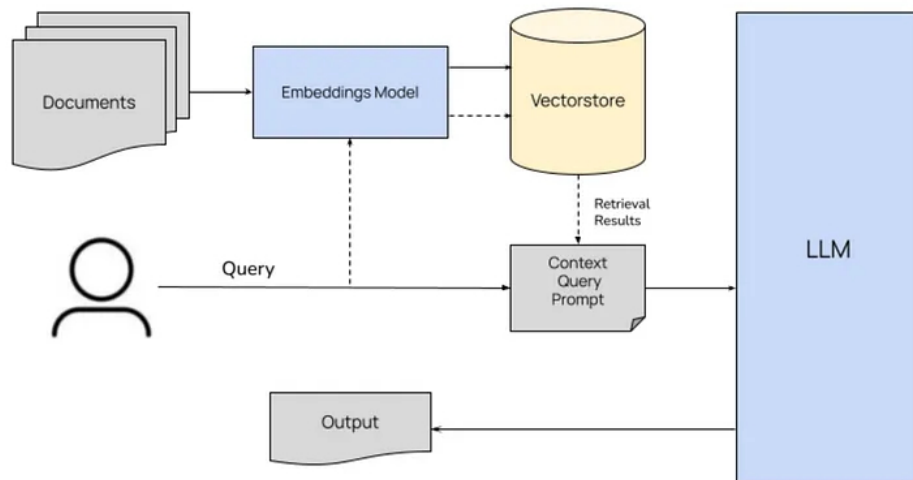


Figure 2.6: High-level overview of a RAG system with its four main components [5]

2.5.2 Embedding Phase

The embedding phase forms the foundation of a Retrieval-Augmented Generation (RAG) system by building the knowledge corpus that the system will later reference to answer user queries (Figure 2.7). While the concept of text embeddings has existed since the introduction of Word2Vec in 2013 [39], their application for contextual storage and retrieval became prominent with the advent of *in-context learning*, as introduced in the GPT-3 paper [9]. This technique leverages the autoregressive nature of decoder-only transformers to generate relevant outputs based on few or even a single example [40].

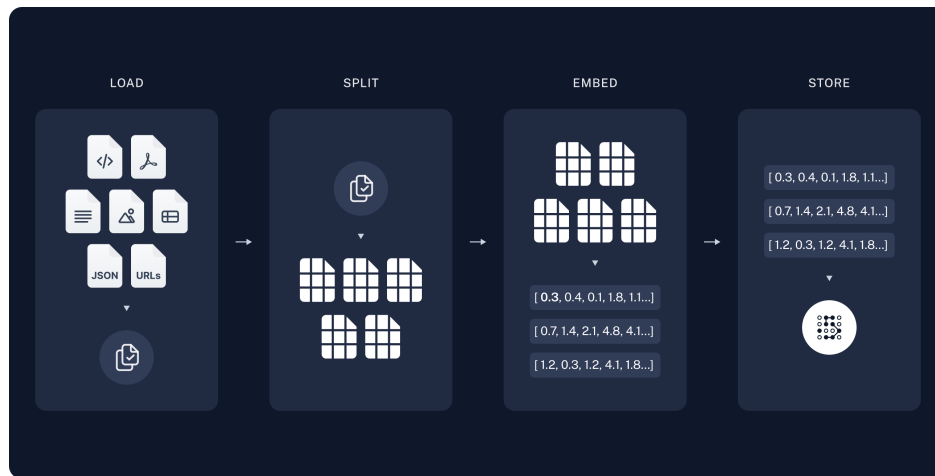


Figure 2.7: The embedding phase: document ingestion, chunking, and vectorization [6]

The embedding phase comprises several key steps:

- a. **Document Processing:** Documents in formats such as TXT, PDF, or DOCX are loaded. This may be done in text-only or multi-modal modes, the latter preserving diagrams and figures. This task is handled by the Document Loader module.

Currently, this project focuses exclusively on PDF files, as most books and academic papers are commonly distributed in this format.

- b. **Chunking:** The loaded documents are split into smaller, semantically coherent segments using text splitters [41]. Chunking is necessitated by the limited context window of LLMs, which ranges from 2K tokens in GPT-3 to 128K tokens in LLaMA 3 [42]. Moreover, using large contexts also demands significant GPU memory, especially when optimizations like KV caching and speculative decoding are employed.

The size of text chunk has a significant impact on the RAG performance due to the following reasons:

- *Chunks too large* may result in the “Lost in the Middle” effect [43], where only the beginning and end of the chunk influence the output.
- *Chunks too small* lead to redundancy and latency, as overlapping content is required to preserve context and sense of continuity.

*This project leverages **MiniLM**[44] which has a context size of 512 tokens. MiniLM is known for its effectiveness in knowledge distillation[45], to generate sentence-level embeddings. These embeddings capture the overall meaning of a text segment and serve as semantic representations, enabling efficient retrieval of relevant context.*

c. **Chunking Techniques:**

- **Length-based chunking:** Divides text based on fixed token or character length. Tokens are determined using subword tokenization techniques like byte-pair encoding [46].
- **Semantic chunking:** Splits based on document structure, e.g., paragraphs.
- **Context-aware splitting:** Ensures that subtopics are not broken across chunks.

This project creates fixed size text chunks of 500 characters with overlap of 20 characters on each end. These parameters were chosen to match the context size of the LLM and also based on common practices and existing research results [47].

- d. **Text Embedding:** In this step, each chunk is converted into a dense vector using a text embedding model. The embedding model used here does not need to match the LLM used in generation, as the embeddings are utilized solely for vector similarity search. Once relevant chunks are retrieved, their original textual content is passed to the LLM, not the embeddings themselves. Compact and efficient models like BERT or DistilBERT are often employed for this task due to their relatively small embedding sizes (e.g., 512 or 768 dimensions), which makes them computationally lightweight compared to models like LLaMA [27] or GPT-3 [9], which produce larger embeddings in the range of 1024 to 12,288 dimensions. Although these smaller embeddings may be less precise, they are generally sufficient for high-level semantic retrieval.

2.5.3 Retrieval Phase

The retrieval phase is initiated whenever a user submits a query or task. This phase makes use of the vector database constructed during the embedding phase to locate and extract relevant content for the given input prompt, which is then passed to the language model for response generation (Figure 2.8).

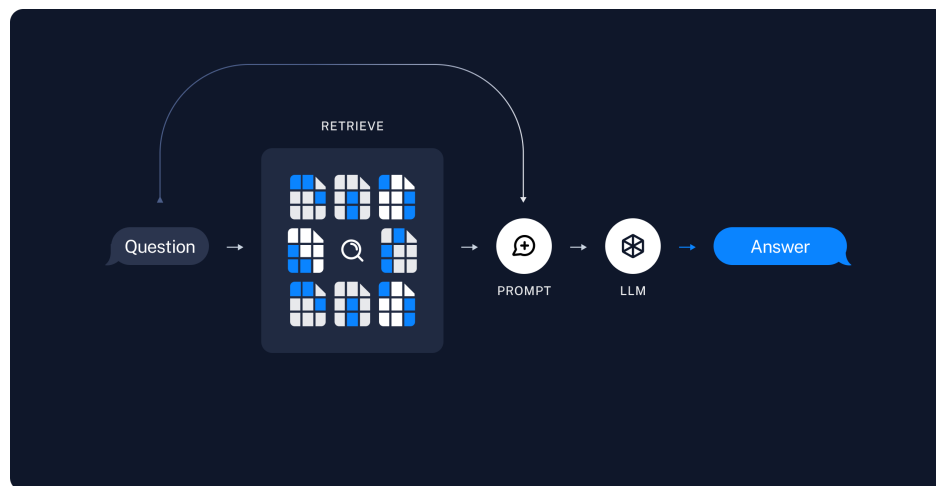


Figure 2.8: The retrieval phase: querying the vector store and invoking the LLM [6]

- a. **Context Retrieval:** The context retriever module follows a two-step process:
 - (a) Embedding the user’s query using the same embedding model employed during the embedding phase.
 - (b) Performing a similarity search in the vector database to identify top-matching chunks based on the embedded query.

The quality of similarity search is crucial, as it directly impacts the relevance and accuracy of the LLM-generated response. In the embedding space, proximity denotes semantic similarity—closer vectors imply greater contextual relevance. However, efficient neighbor discovery in high-dimensional vector space is computationally expensive, and full database scans are often impractical [48].

- b. **Search Algorithms:** To balance accuracy, latency, and resource usage, the following approximate nearest neighbor (ANN) algorithms are commonly employed:
 - (a) **Cosine Similarity Search:** Measures the cosine of the angle between vectors to determine their similarity. It is widely used in embedding-based retrieval tasks due to its scale invariance and intuitive geometric interpretation. This is **highly parallelizable** and simple to implement. [49]. *This is project primarily relies on cosine similarity for embedding retrieval.*
 - (b) **k-Nearest Neighbors (kNN):** A brute-force method that identifies the k closest vectors through exhaustive comparisons [50].
This project uses kNN as a fallback mechanism to search in the database using pg-vector [51].
 - (c) **HNSW (Hierarchical Navigable Small World) Graphs:** Constructs a layered graph to enable fast traversal through neighbors across different granularity levels [50].

- (d) **ScaNN (Scalable Nearest Neighbors)**: A Google-designed ANN method that integrates pruning, quantization, and partitioning for scalable, memory-efficient retrieval [52].

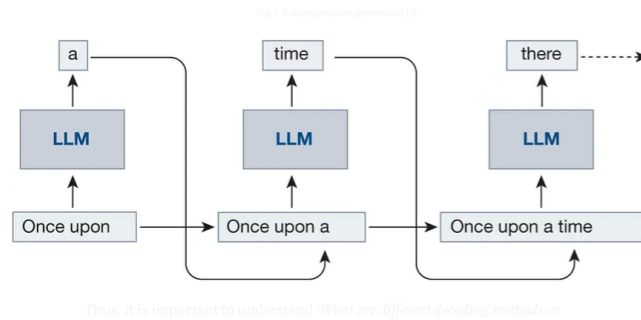


Figure 2.9: Autoregressive decoding [2]

- c. **Output Generation**: In this phase, the retrieved context and the user input—both in plain text form—are fed into the Large Language Model, as illustrated in Figure 1.3. The LLM processes the portion of input that fits within its context window and generates an output token. This token is then concatenated with the input and passed back to the LLM. The process repeats iteratively until the LLM produces an `¡EOS¡` (end-of-sequence) token.

2.5.4 Output Evaluation:

The evaluation of a Retrieval-Augmented Generation (RAG) system focuses on the quality of the output in terms of its relevance to the user’s query and the information stored in the vector database. Since different components contribute uniquely to the overall performance, multiple metrics are used rather than a single aggregated score to allow for precise attribution and targeted optimization.

Classical metrics such as Precision, Recall, and F1 score, alongside NLP-specific metrics

like ROUGE, are employed to assess performance. The evaluation considers several factors, as illustrated in Figure 1.5:

- **Faithfulness, Output Relevance & Semantic Similarity:** Measures the quality of the output relative to the input queries and retrieved context.
- **Context Recall & Context Precision:** Assess the effectiveness of context retrieval and its utilization by the LLM.

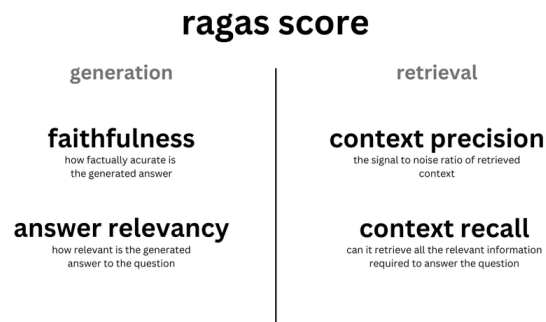


Figure 2.10: RAG Output evaluation metrics [7]

Chapter 3

RELATED WORK

This chapter presents an overview of existing technologies and research efforts relevant to the development of efficient, on-device language model-based applications. It focuses on three major areas: Retrieval-Augmented Generation (RAG), efficient LLM runtimes like `llama.cpp`, and lightweight distribution and serving solutions such as Ollama and Llamafire. It also looks at the work done by the `tinygrad` project in attempting to reverse engineer the NPU API.

3.1 RAG Frameworks and Agentic Systems

Retrieval-Augmented Generation (RAG) has become a widely adopted paradigm to enhance the factual grounding of large language models (LLMs) by integrating external knowledge through retrieval mechanisms. To support this architecture, several frameworks have been developed to facilitate the construction of RAG pipelines. **LangChain** [53] provides a composable and extensible framework for chaining LLMs with retrieval components, supporting a broad ecosystem of vector stores, retrievers, and tools. **LlamaIndex** (formerly GPT Index) [54] offers structured pipelines for indexing and querying private or domain-specific data sources. **Haystack** [55] by deepset provides a modular toolkit for constructing production-ready pipelines with retriever-reader architectures.

In parallel, the emergence of **agentic systems**—LLM-powered agents capable of tool use, memory, and reasoning—has opened new avenues for task-oriented automation. Frameworks such as LangChain [53] and AutoGen [56] support multi-agent orchestration and planning. These systems increasingly incorporate **LLM routing** strategies [57] that dynamically select or combine multiple models or tools based on the query or context, enhancing efficiency and

specialization. Additionally, **RAGAS** [58] introduces an evaluation framework to measure the quality and factual alignment of RAG outputs, contributing to the robustness of such systems in real-world applications. Together, these frameworks and methodologies form the foundational ecosystem for deploying scalable, intelligent, and grounded LLM applications.

Although this project does not directly leverage any existing RAG frameworks yet, it certainly has drawn inspiration and insights wherever applicable.

3.2 *llama.cpp*

`llama.cpp`[59] is a C++ implementation of large language models (Meta LLaMA to begin with), optimized for local inference on commodity hardware without GPU requirements. Built upon the GGML tensor library, it provides quantized inference for large models using CPU-friendly formats like 4-bit or 5-bit quantization, making it suitable for running models such as LLaMA, Mistral, and other open-weight transformers on devices ranging from laptop to Raspberry Pi.

Key features of `llama.cpp` include:

- Highly efficient CPU inference with quantized models.
- Cross-platform support (macOS, Linux, Windows).
- Integration with popular tooling such as LangChain and Open Interpreter.
- Support for multi-threaded inference and memory-mapped model weights for efficient memory usage.

Llama.cpp serves as a foundational component for many desktop LLM applications that prioritize local execution and privacy-preserving computation. *In this project (Project TLDR), llama.cpp is used for LLM inference tasks, including both embedding generation and text generation.*

3.3 *Ollama*

Ollama is a developer-friendly platform for running LLMs locally with simplified model management and serving. It wraps models like LLaMA 2, Mistral, and Code LLaMA into a streamlined runtime with a CLI and RESTful API, abstracting away hardware-specific setup and providing a plug-and-play experience for developers [60].

Ollama supports:

- Running quantized models locally with GPU acceleration where available.
- Seamless model downloading and serving.
- Custom model creation using a simple Modelfile syntax.

It is widely used for prototyping private, offline chatbots and assistants. However, the users need to be technically savvy in cases of issues downloading or running the models. Furthermore, models often may not be quantized and could lead to downloading of model weights in order of many GBs.

3.4 *Llamafile*

Llamafile, developed by Mozilla-Ocho, enables packaging a complete LLM runtime into a single, self-contained executable file [61]. It leverages the `llama.cpp` backend and Cosmopolitan Libc to build universal binaries that run across major operating systems (Windows, macOS, Linux) without requiring dependencies.

Notable features:

- Distributable as a single file.
- Useful for shipping LLM-based tools with zero-install requirements.
- Integrates with web frontends for local chatbot deployment.

Llamafire is widely used for prototyping private, offline chatbots and assistants. However, it often requires users to be technically proficient, especially when encountering issues related to downloading or executing models. Moreover, many models are not pre-quantized, potentially resulting in downloads of several gigabytes of model weights.

3.5 *Tinygrad project and Apple Neural Engine (ANE)*

The Apple Neural Engine (ANE) is a custom neural processing unit designed by Apple to accelerate machine learning workloads on its silicon platforms. Introduced with the A11 Bionic chip, the ANE has evolved into a high-performance, low-power DMA-based inference engine embedded in Apple’s M-series chips. This section synthesizes insights obtained by the reverse-engineering efforts of the `tinygrad` [25] project to examine ANE’s architecture, capabilities, and compilation flow.

3.5.1 *Hardware Overview*

The ANE operates primarily as a DMA engine optimized for convolutional operations and supports a wide range of neural network layers and fused operations. Its key hardware features include:

- **16-core architecture:** A 16-wide Kernel DMA engine for parallel computation.
- **5D Tensor Support:** Tensors are structured with width (column), height (row), planes (channels), depth, and group (batch).
- **Supported Data Types:** `UInt8`, `Int8`, and `Float16` (with `Float32` inputs automatically downcast).
- **Manually Managed 4MB L2 Cache:** Applied only to input/output data; weights are embedded in the compiled program.
- **Execution Unit:** Executes up to 0x300 micro-operations per instruction.

- **Performance:** Approximate 11 TOPS throughput, assuming 32×32 MAC at 335 MHz.

All memory strides are constrained to multiples of 0x40 bytes, reflecting hardware alignment requirements.

3.5.2 Software and Compilation Stack

The ANE software stack is heavily abstracted behind Apple’s proprietary frameworks but has been reverse-engineered to reveal a structured flow (as shown in Fig 3.1):

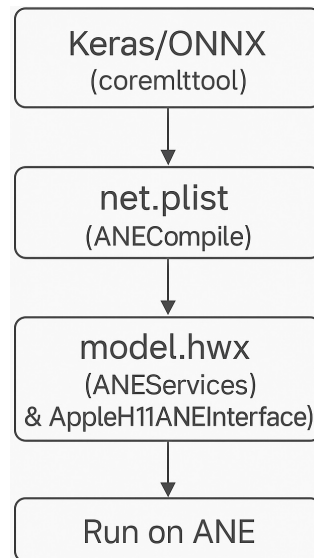


Figure 3.1: Apple Neural Engine Workflow

1. **Model Definition:** Models are authored in Keras or ONNX.
2. **Conversion:** Models are converted to CoreML format using open-source tools such as `coremltools`.
3. **Intermediate Representation:** CoreML is internally converted into `net.plist` by Apple’s Espresso framework.

4. **Compilation:** The `ANCompiler` service transforms `net.plist` into a hardware-specific binary (`.hwx`), a Mach-O formatted executable.
5. **Execution:** The `AppleNeuralEngine` and `ANEServices` handle execution via the kernel extension `AppleH11ANEInterface`.

3.5.3 Instruction Format and Operation Structure

Each ANE instruction is 0x300 bytes and comprises multiple segments:

- **Header:** Includes DMA addresses and next-op offset.
- **KernelDMA Src:** Specifies weights, bias, and channel usage.
- **Common:** Describes input/output shapes, types, kernel size, and padding.
- **TileDMA Src/TDMADst:** Layout and stride configurations for input/output tensors.
- **L2 and NE:** L2 cache flags and activation parameters.

3.5.4 Supported Operations and Activations

ANE supports a variety of operations:

- **Core Ops:** CONV, POOL, EW, CONCAT, RESHAPE, MATRIX_MULT, TRANSPOSE
- **Advanced:** SCALE_BIAS, SOFTMAX, INSTANCE_NORM, BROADCAST, L2_NORM
- **Fused Ops:** NEFUSED_CONV, PEFUSED_POOL, etc.
- **Activations:** RELU, SIGMOID, TANH, CLAMPED_RELU, PRELU, LOG2/EXP2, CUSTOM_LUT

Over 30 activation functions are supported in hardware.

3.5.5 *tinygrad Implementation*

The `tinygrad` project interfaces directly with ANE through a three-stage pipeline:

- **1_build:** Generates CoreML models using `coremltools`.
- **2_compile:** Uses Objective-C and Apple’s private ANECompiler framework to compile models into HWX binaries.
- **3_run:** Loads HWX binaries and executes them on ANE using custom Objective-C wrappers around `AppleH11ANEInterface`.

The implementation also includes tools like `hwx_parse.py` for disassembling HWX files and visualizing internal ops.

3.5.6 *Security and Access*

Execution on ANE requires system entitlements that are typically unavailable to third-party applications:

- `com.apple.ane.iokit-user-access`
- Workarounds: amfid patching, kernel extension modification, or use of provisioning profiles.

3.5.7 *ANE Takeaways*

The ANE represents a proprietary, highly optimized inference accelerator that is difficult to access and understand due to Apple’s closed ecosystem. Reverse engineering, as demonstrated by `tinygrad`, reveals a modular, DMA-centric architecture capable of executing complex neural network operations at high throughput and low latency. As Apple continues

to iterate on the ANE, deeper access and tooling may unlock broader ML deployment options on Apple hardware.

Chapter 4

METHODOLOGY

This chapter outlines the technical methodology adopted in the development of the system, covering both low-level design and software implementation. The approach prioritizes performance, privacy, and modularity, leveraging hardware accelerators where possible and maintaining efficient control over data flow and execution.

4.1 Application Modules and Design

This section provides an overview of the internal architecture and design principles behind the TLDR desktop application. The application is built with the goal of providing a fast, private, and efficient interface for question-answering and summarization tasks over a user-defined corpus of documents. To achieve this, the design incorporates several performance-aware and hardware-conscious modules, especially tailored for Apple’s M1/M2 architecture.

The TLDR system is structured into modular components, each responsible for a specific functionality in the information processing pipeline. These include embedding generation, context retrieval, vector storage, prompt construction, and output generation via a language model. The workflow between these modules is coordinated to support seamless execution, low latency, and high responsiveness, all while maintaining data privacy by running entirely on-device.

4.1.1 Overview

The TLDR application follows a modular architecture where different components are responsible for distinct tasks in the RAG (Retrieval-Augmented Generation) pipeline. Figure 4.1 illustrates the overall design and the control flow between the modules.

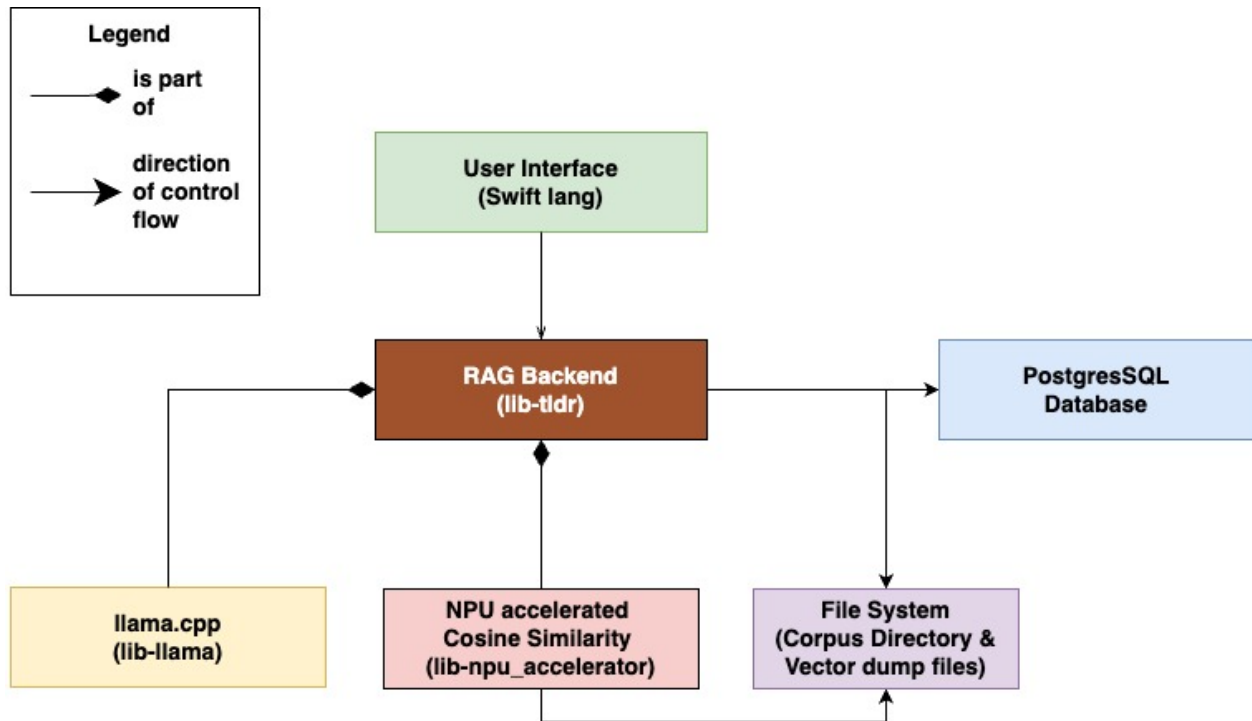


Figure 4.1: Modules of TLDR application

- **User Interface:** This module provides the graphical frontend for the user. Developed using Swift language for MacOS, it allows users for the users to seamlessly leverage the capabilities of the application. It is responsible for workflows dealing with user experience while delegating all core logic to the backend modules.
- **RAG Backend:** This is the core orchestrator of the application. It manages the full pipeline, including handling user queries, initiating vector search, performing retrieval, and forwarding context to the language model. It communicates with all supporting modules such as the database, NPU based vector search module, llama.cpp and the file system.
- **Database (PostgreSQL):** Stores metadata and document indexing information. It ensures efficient retrieval and persistence of preprocessed documents and vector refer-

ences. It plays a crucial role of mapping retrieved vector hashes to their text chunks and document metadata during context retrieval.

- **File System (Corpus Directory and Vector Dump Files):** Contains the document corpus (directory containing source documents) and their corresponding vector dump files. These vector dump are leveraged by the vector search engine using memory-mapped I/O for efficient vector search with low memory overhead.
- **NPU Accelerated Cosine Similarity:** Implements hardware-accelerated cosine similarity by leveraging Apple's Neural Processing Unit (NPU). The backend invokes this module for fast and parallelizable vector cosine similarity computation.
- **llama.cpp:** This module is responsible for language generation i.e LLM inference. It acts as a plugged-in module for the RAG backend and contributes by generating embeddings and chat response during the corresponding stages of the RAG pipeline.

4.1.2 User Interface

The User Interface is in the form of a native MacOS desktop application developed using Swift lang in the Xcode development environment (as depicted in Fig 4.2, Fig 4.3). The user interface is designed to be intuitive, lightweight, and self-contained. The application packages all necessary dependencies, including static libraries and LLM weights. enabling fully offline functionality without requiring additional installation or configuration.

The user interface module has the following core purposes:

- Provides a clean, chat-style interface where user prompts and LLM responses are displayed as distinct messages to mimic a conversational flow.
- Manager all responsibilities regarding user interaction, including persisting user conversation history and any additional user preferences and thereby enable RAG backend to only focus on the RAG functionalities.

- *Make a single, self contained portable package that is easy and intuitive to distribute.*

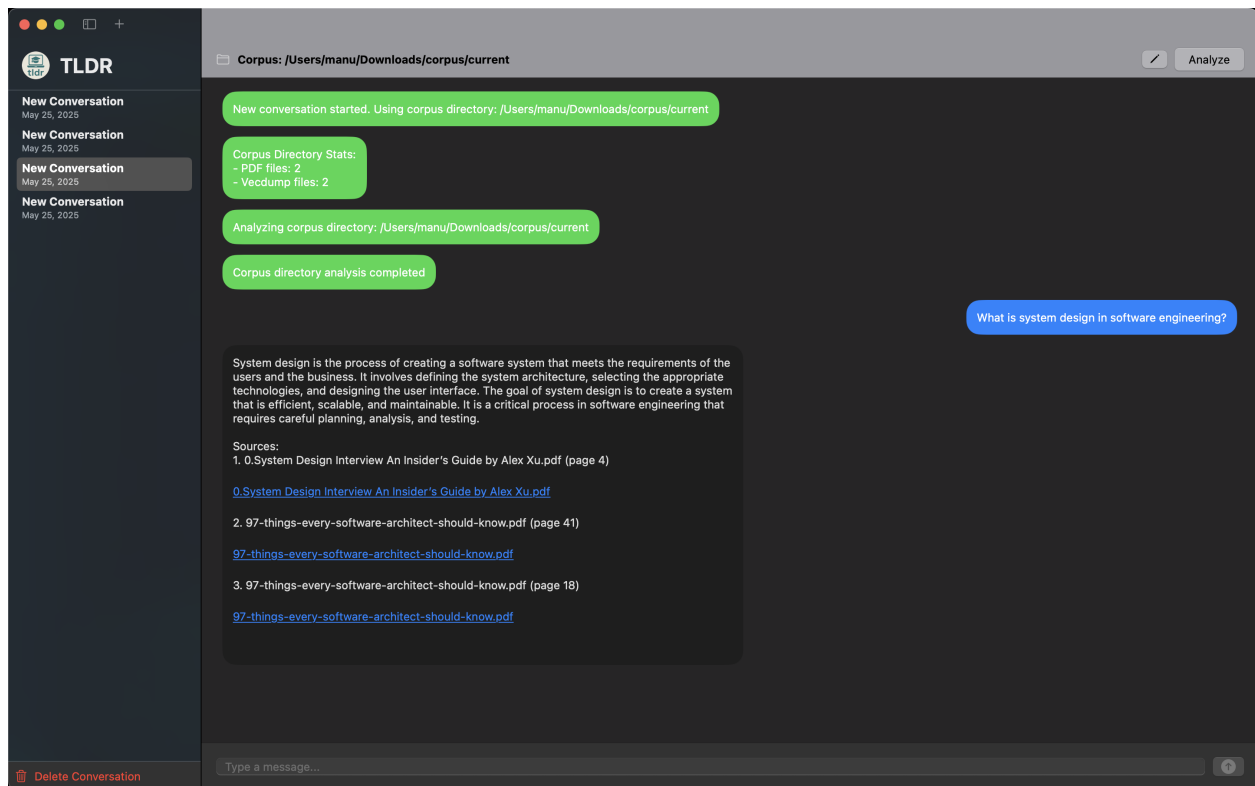


Figure 4.2: Graphical User Interface of TLDR Application

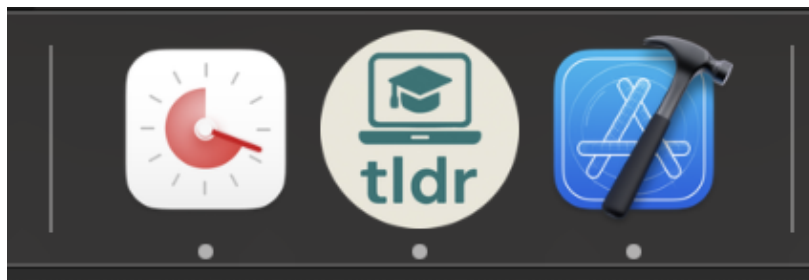


Figure 4.3: TLDR Application Icon as seen in MacOS Dock

Codebase Organization

The codebase of the GUI application is illustrated in Figure 4.4. It is organized into several components that serve both core functionality and supporting roles within the application.

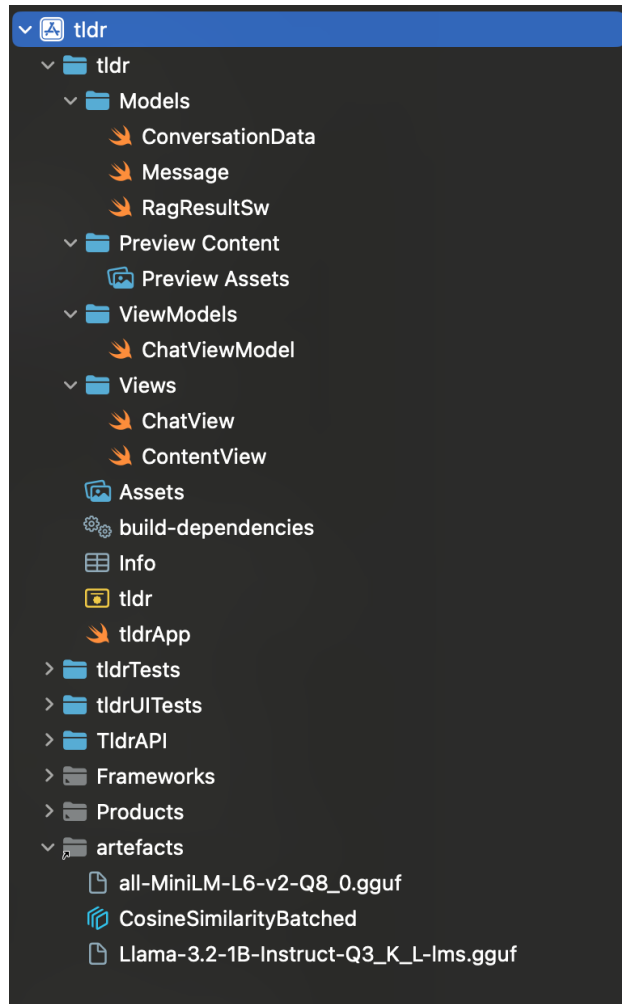


Figure 4.4: UI Project File system

A. Overall Structure:

- a. `tldr`: The Swift language codebase that creates the user interface.

- b. **TldrAPI**: C++ module with bindings for Swift UI. It serves as a bridge between the Swift UI and the C++ static library of the RAG Backend.
- c. **artefacts**: Quantized LLM weights and coreml packages, i.e:
 - **Llama-3.2-1B-Instruct-Q3_K_L** model for chat
 - **all-MiniLM-L6-v2-Q8_0** model for generating embeddings
 - **CosineSimilarityBatched** coreml package for cosine similarity on NPU

B. UI Architecture (MVVM):

- a. **tldrApp**: The main SwiftUI entry point where the application lifecycle begins.
- b. **Models**: Includes **ConversationData**, **Message**, and **RagResultSw** to represent the chat state and RAG outputs. These modules leverage *UserDefaults* a built-in key-value persistence mechanism provided by Apple’s Foundation framework to efficiently store and retrieve their information.
- c. **Views**: Comprises SwiftUI components like **ChatView** and **ContentView** to render the main interface.
- d. **ViewModels**: Contains **ChatViewModel**, which handles user interaction and back-end coordination.
- e. **Preview Content**: Includes **Preview Assets** for SwiftUI previews to support development and layout testing.
- f. **Assets**: Stores static resources such as icons and other UI elements.

This organization promotes maintainability and allows for a clear separation of concerns between UI presentation, interaction logic, and backend communication.

4.1.3 RAG Backend

The RAG backend is implemented as a C++ static library named **lib_tldr**. It serves as the backbone for all core functionalities of the project. This module encapsulates the complete

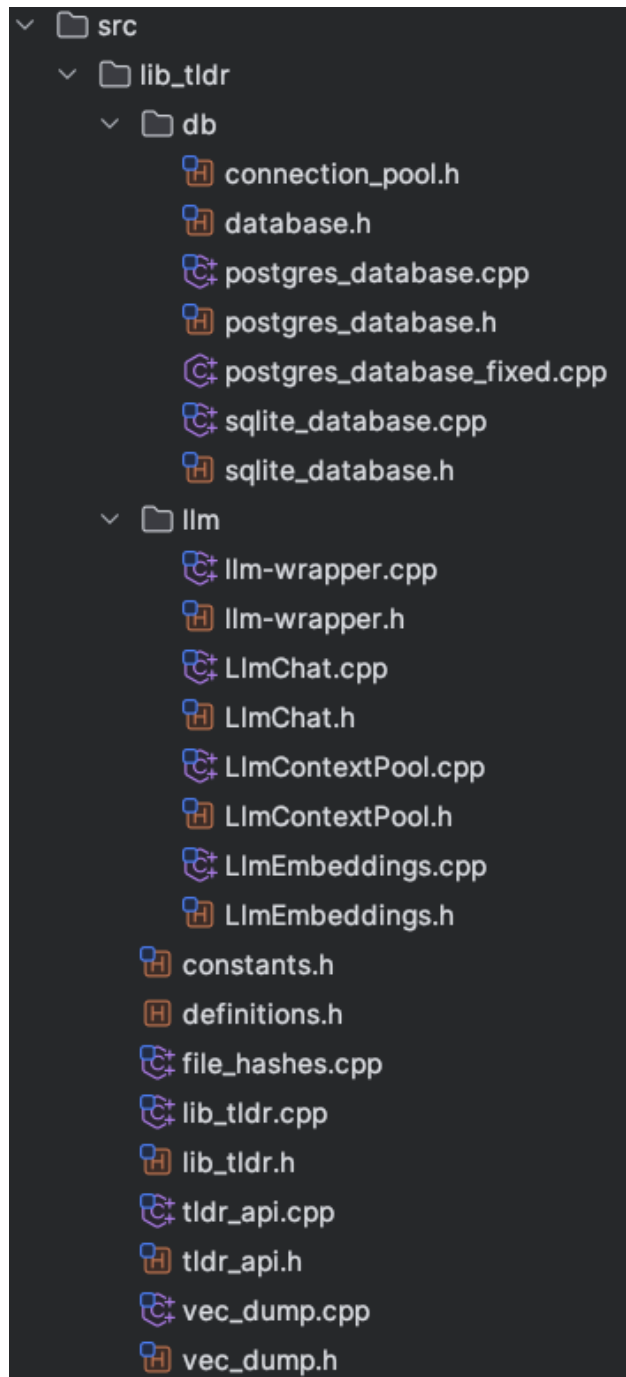


Figure 4.5: RAG Backend(lib_tldr) codebase

RAG pipeline logic and integrates essential components such as the large language model for text and embedding generation, the vector dump generator and reader, database and file system handlers for storage, and the cosine similarity-based vector retriever (as seen in Fig 4.5).

The architecture emphasizes modularity and resource efficiency, enabling plug-and-play replacement or extension of components. It adheres to the SOLID principles of object-oriented programming, promoting adaptability, maintainability, and scalability.

Following are the logical sub-modules of RAG Backend:

A. Language Model Interface:

The system utilizes `llama.cpp` as the backend for both text generation and embedding extraction. The following components facilitate this integration:

- (a) `LlmChat.cpp &.h`: Manage the logic to take user input and retrieved context and generate a response from the chat LLM.
- (b) `LlmEmbeddings.cpp &.h`: These modules extract dense vector representations (embeddings) from document chunks, which are subsequently used for semantic similarity search and retrieval.
- (c) `LlmContextPool.cpp &.h`: These components manage the lifecycle of LLM context objects, which encapsulate the model state necessary for efficient inference. Since LLMs typically require a context to maintain internal buffers, tokenizer state, and memory allocations, it is computationally expensive to initialize a new context for every query or embedding operation. By pooling and reusing contexts across multiple operations, the system significantly reduces initialization overhead and ensures smoother, low-latency performance during both chat and embedding workflows.
- (d) `Llm-Wrappers.cpp &.h`: Takes care of initializing and cleaning up resources related to the LLMs. It initializes common ggml backend for Apple metal and setting

up LLM context pools and initialize Chat and Embedding LLMs by loading their weights into memory.

B. Database Interaction: All persistent storage is handled via the PostgreSQL backend, though an optional SQLite backend is also implemented (but not currently utilized). The database interaction is managed within the `db` submodule, which includes the following:

- a. `database.h` Defines an abstract class that enforces a uniform interface for any underlying database implementation. This design allows the RAG codebase to remain unchanged when switching between different database technologies.
- b. `postgres_database.cpp & .h`: These files handle database initialization, schema definition, and CRUD operations related to documents and embeddings in PostgreSQL Database.
- c. `connection_pool.h`: Provides a lightweight connection pooling mechanism to manage multiple concurrent database sessions efficiently. This is especially beneficial during large-scale embedding operations where multiple inserts are performed rapidly. It is efficient to store readily available connections and re-use them instead of creating a new connection for every db interaction.

The PostgreSQL schema stores both high-level document metadata (e.g., title, author, page count) and low-level embedding-related information (e.g., text chunk, hash, embedding vector, page number, and timestamps).

C. Embeddings Vector Storage and Retrieval

`vec_dump.cpp & .h` are responsible for managing the serialized storage of raw vector data ("vecdumps"). These are binary representations of embedding vectors that can be rapidly accessed and processed.

Additionally, the system incorporates a hardware-accelerated module referred to as the `npu-accelerator`, which leverages macOS’s Neural Engine to perform cosine similarity search over large sets of embeddings. This offloads compute-intensive operations from the CPU, enabling real-time retrieval performance on resource-constrained devices.

D. Core Workflow: The `lib_tldr.cpp &.h` contains the core logic that glues the entire system together, such as:

- Initializing the LLMs and the Database tables (when necessary)
- Creating DB connection pool and LLM context pool
- Checking for changes in the corpus directory and embedding new documents
- Performing Retrieval Augmented Generation
- Cleaning up and releasing acquired resources

E. RAG Backend API: While the backend contains numerous functions and data structures for internal logic, a clean and minimalistic API (Application Programming Interface) is exposed. This allows its client modules (such as the UI module) to leverage its capabilities without being closely coupled with the internal mechanisms of the library.

The `tldr_api.cpp &.h` files expose a clean, C-style API interface for the user-facing application layer on top of the functions present in `lib_tldr.cpp &.h`.

4.1.4 Database (PostgreSQL)

The application uses PostgreSQL as its persistent storage backend to manage and retrieve embedding data required during the Retrieval-Augmented Generation (RAG) process. The database stores preprocessed document chunks, their vector embeddings, associated metadata, and file-level information.

PostgreSQL was chosen over lighter-weight alternatives like SQLite due to its superior concurrency handling. Specifically, SQLite’s single-writer limitation presented a bottleneck

in the multi-threaded embedding pipeline, where concurrent writes to the embedding store are common. PostgreSQL’s support for multiple concurrent writers allows the embedding process to scale efficiently without serialization delays.

Table: documents

This table stores metadata for each unique input file in the corpus. It ensures file-level uniqueness through the `file_hash` field and includes fields such as file name, author, subject, page count, and timestamps. It acts as a parent entity in a one-to-many relationship with the `embeddings` table (refer to table definition in Fig 4.6).




	column_name 	data_type 	is_nullable 
	name	character varying	character varying (3)
1	id	uuid	NO
2	page_count	integer	YES
3	created_at	timestamp with time zone	YES
4	updated_at	timestamp with time zone	YES
5	title	text	YES
6	author	text	YES
7	subject	text	YES
8	keywords	text	YES
9	creator	text	YES
10	producer	text	YES
11	file_hash	text	NO
12	file_path	text	NO
13	file_name	text	NO

Figure 4.6: Documents table description

Table: embeddings

This table stores the chunk-level data required for semantic search. Each row contains a reference to a document, the original text chunk, its embedding vector, and an embedding

hash. During query time, the vector search module returns the top- K most similar vectors based on cosine similarity. The corresponding text chunks are then retrieved from this table using the hash and document ID to form the LLM context (refer to table definition in Fig 4.7).




	column_name 	data_type 	is_nullable 
	name	character varying	character varying (3)
1	created_at	timestamp with time zone	YES
2	page_number	integer	YES
3	document_id	uuid	YES
4	id	bigint	NO
5	embedding	USER-DEFINED	NO
6	embedding_hash	text	YES
7	chunk_text	text	NO

Figure 4.7: Embeddings table description

4.1.5 File System: Corpus Directory and Vectordump Files

Vector dump files are binary data structures designed for efficient storage of document embeddings. For each input document, a corresponding vector dump file is created. Each such file contains embedding vectors generated from text chunks along with their corresponding MD5 hashes. This format enables fast similarity search and content verification in document retrieval systems.

Corpus Directory

The corpus directory serves as the source of truth for the RAG pipeline. It contains the raw documents—primarily PDF files—that are parsed, chunked, and processed by the language model to construct the knowledge base used for information retrieval.

This project recursively scans the specified corpus directory, identifies all PDF files, and computes their corresponding embeddings and stores them in **vecdump** files.

Vectordump Files

Vector dump files are obtained as a result of the embedding process. After a document is split into chunks, each chunk is processed by the LLM and yields embeddings. Embeddings are higher dimensional representations of the input text, as interpreted by the LLM. These embeddings are then stored in a dedicated subdirectory named `_vecdump`, located within the corpus directory itself. The `_vecdump` folder houses binary `.vecdump` files that are later used during the similarity search phase to efficiently retrieve semantically relevant chunks.

The vector dump file follows a sequential binary layout consisting of three main components: a metadata header, followed by embedding data, and finally hash data (as illustrated in Fig 4.8). This structure allows for efficient random access to embeddings while maintaining data integrity through hash verification. Hash is also further used for fetching the corresponding text chunk after similar vectors are obtained for a query.

Vectordump Header

In order to process the vector dump file, the header is first read and the necessary information is obtained regarding the data layout in the file. The information is arranged in the layout as illustrated in Figure 4.8. The elements are as follows:

- `num_entries` - Total number of embedding/hash pairs stored in the file
- `hash_size_bytes` - Size of each MD5 hash in bytes (always 16 for standard MD5)
- `vector_size_bytes` - Total byte size of each embedding vector
- `vector_dimensions` - Number of floating-point dimensions per embedding vector

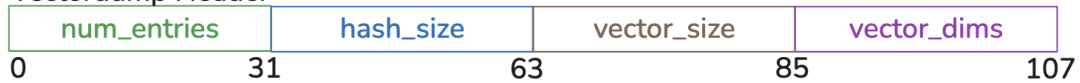
The header is read by simply pointing a pointer of type `structVectorDumpHeader` to the memory location to which the file is loaded. This helps obtain the necessary information required to access the data sections of the file.

```

struct VectorDumpHeader {
    uint32_t num_entries;      // Number of embedding vectors/hashes
    uint32_t hash_size_bytes; // Size of each hash in bytes
    uint32_t vector_size_bytes; // Size of each embedding vector in bytes
    uint32_t vector_dimensions; // Number of dimensions in each vector
};

```

Vectordump Header



Vectordump File

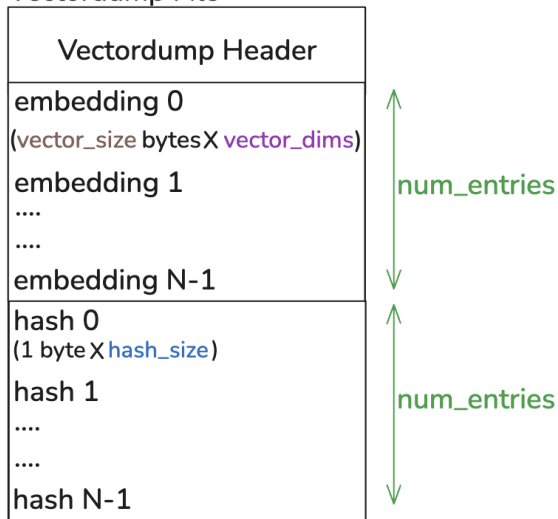


Figure 4.8: Vector dump file structure

Data Sections

Once the header section is loaded, the data obtained is then used for calculating the memory locations of the embedding and hash arrays. Pointers are used to access these locations to simulate the structure of an array on top of raw binary data read into the memory. This approach is simple and efficient and prevents needless memory allocations and data copies.

Embeddings: Contains N consecutive embedding vectors, where $N = \text{num_entries}$. Each vector occupies `vector_size_bytes` and represents a `vector_dimensions`-dimensional

embedding, stored as 32-bit floating point values. In the current case, we use a 384-dimensional vector embedding.

Hashes: Contains N consecutive MD5 hash values, each exactly 16 bytes. However, the smallest unit of storage is of `uint64_t` type, i.e., units of 2 bytes. Hence, a 16 byte MD5 hash would have a hash size of 8. The hash at index i corresponds to the MD5 digest of the original text chunk used to generate `embedding[i]`.

Data Relationship

The file maintains strict positional correspondence: for any index $i \in [0, N-1]$, `embedding[i]` and `hash[i]` represent the same document chunk. This one-to-one mapping enables efficient lookup operations and integrity verification during retrieval.

The data is used as follows:

1. Load the first half of the file in memory using `mmap` and perform cosine similarity search.
2. Obtain index of the top K relevant vectors from Cosine similarity module.
3. Fetch the hash values at the obtained indices.
4. Query the database for text chunks associated with the hash values.

This design allows for prioritized access to necessary data and its direct usage for cosine similarity search with no further processing or data manipulations, allowing for an efficient search through the entire corpus.

4.1.6 NPU Accelerated Cosine Similarity

The NPU accelerator module (`lib-npu-accelerator`) is a specialized component of the TLDR MacOS desktop application that leverages Apple’s Neural Processing Unit to perform hardware-accelerated cosine similarity computations as part of the RAG pipeline. The module con-

struction and usage is a multi-step process involving multiple components. The codebase structure is depicted in Fig 4.9 and the workflow in Fig 4.10.

Codebase structure

The codebase for the npu module as seen in Fig 4.9 has the following components:

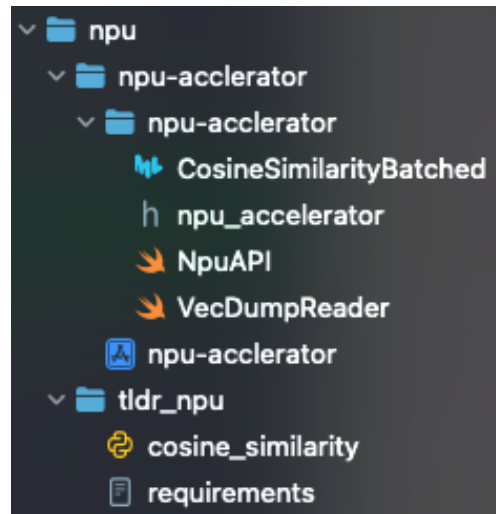


Figure 4.9: lib-npu_accelerator codebase

- **PyTorch module (tldr_npu):** Contains PyTorch code to perform batched cosine similarity. This module is used to generate a CoreML model package, which is later utilized by the NPU accelerator for high-performance similarity computation.
- **Swift module (npu-acclerator):** Implements the logic in Swift to read vector dump files and leverage the CoreML cosine similarity model to perform efficient vector similarity search using the Apple Neural Engine (ANE).

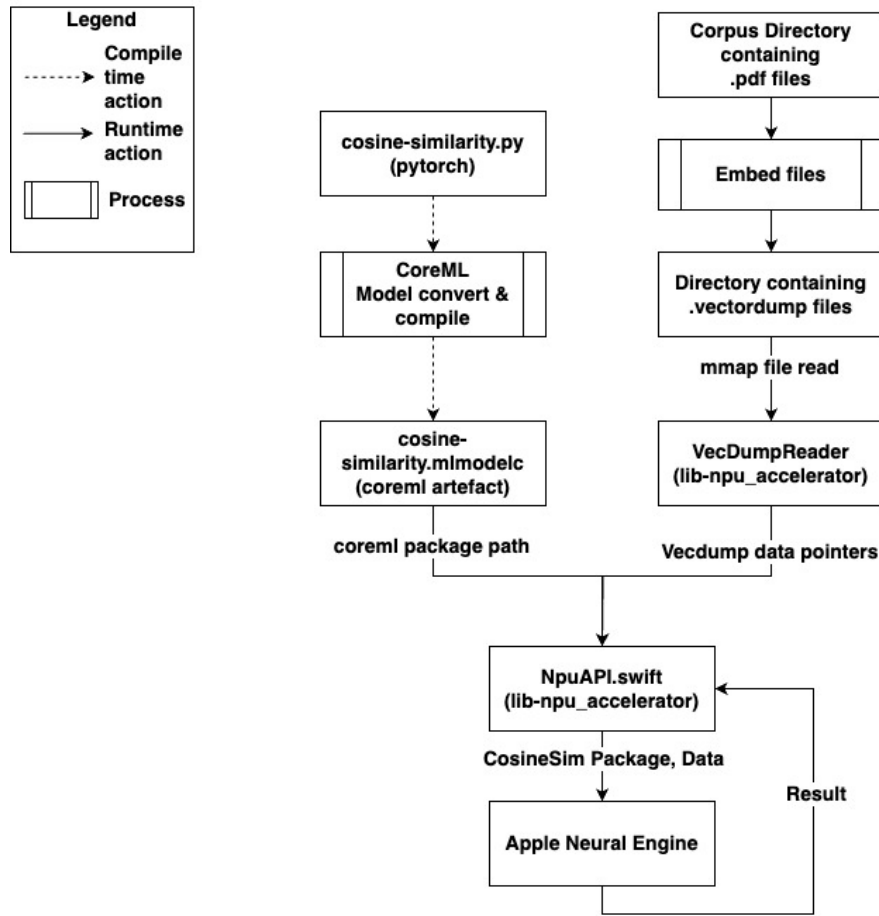


Figure 4.10: NPU accelerator module workflow

Processing Workflow

The NPU accelerator workflow could be visualized as dual-pipeline workflow as illustrated in Fig 4.10.

Phase I - Prepare the components: The left pipeline begins with a PyTorch-based `cosine-similarity.py` implementation that serves as the foundation for similarity computations. This PyTorch model undergoes a CoreML model conversion and compilation process, transforming the original implementation into an optimized CoreML artifact specifically designed for Apple’s Neural Engine execution. The compilation step produces the *cosine-*

similarity.mlmodelc package, which contains the optimized model ready for hardware acceleration.

The right pipeline operates in parallel, handling document processing and vector storage. Although this portion of the pipeline is executed by the RAG Backend and is not directly part of the NPU module, it still is a logical component of the NPU module workflow as depicted in Fig 4.10.

Phase II - Perform vector search: This phase is triggered at runtime when a user submits a query and the RAG pipeline is activated. The workflow converges at `NpuAPI.swift`, a Swift-based interface that integrates the CoreML similarity model and the document embeddings obtained via the `VecDumpReader`.

The `VecDumpReader` accesses vector dump files using memory-mapped I/O (`mmap`), exposing the data as raw pointers without intermediate copies. These pointers are passed directly to the `CosineSimilarityBatched` module, which performs batched cosine similarity computations using Apple’s Neural Engine (ANE).

`NpuAPI` orchestrates this process by coordinating the CoreML model execution and the memory-resident embeddings, enabling efficient hardware-accelerated similarity search. The output consists of similarity scores and embedding hash values, which are used by the RAG system to retrieve the most relevant document chunks for generating a contextually informed response. Furthermore, `NpuAPI` also serves as a bridge to the C++ layer, exposing these capabilities to the RAG backend for use during the retrieval phase of the RAG pipeline.

Although cosine similarity involves a full scan of the embedded corpus for each query, the system architecture mitigates performance concerns through the following mechanisms:

- **Shared memory-mapped files:** When multiple threads handle different user requests, redundant file reads are avoided because the `mmap` mechanism loads the file into memory only once.
- **Zero-copy access:** Since the data is accessed directly from memory without duplication, there is no overhead for repeated reads. Additionally, the operating system

can page the data out to swap memory and page it back in when needed, optimizing memory usage.

- **Unified memory architecture:** The Apple M1/M2 SoC’s unified memory architecture enables seamless access between CPU, GPU and NPU, eliminating the need for further data transfers during vector similarity computations.

4.1.7 *llama.cpp*

The project integrates a customized fork of `llama.cpp`—a lightweight C++ inference engine for LLaMA and related transformer models—to serve as the core engine for both text generation and embedding extraction. The fork is hosted at <https://github.com/manuhg/llama.cpp>, and includes minor modifications that streamline the build process for macOS targets. Specifically, the build scripts were stripped down to produce a static library using the `ggml` Metal backend, optimized for Apple Silicon devices. The result is a single, portable C++ static library named `libllama.a` which is then statically linked with the RAG backend (`lib-tldr.a`).

The core functionalities such as tokenization, decoding, and sampling are accessed through the public APIs exposed in `llama.h` and `ggml.h`. Two primary components—`LLmChat` and `LLmEmbedding`—are built around workflows inspired by `simple.cpp` [62] `server.cpp` [63] and `embedding.cpp` [64] from the upstream `llama.cpp` [65] project.

Further, to improve the performance, `OpenMP` support was added for parallelizing the tokenization and batch decoding steps. This optimization ensures efficient utilization of CPU cores, resulting in faster preprocessing and inference, particularly during multi-threaded interactions.

The `llama.cpp` is hence directly integrated into the RAG backend at compile time. This enables the application to perform in-memory LLM inference by directly invoking components of `llama.cpp`, without relying on any external dependencies or background processes for this core functionality.

4.2 Application Workflow

4.2.1 Workflow Overview

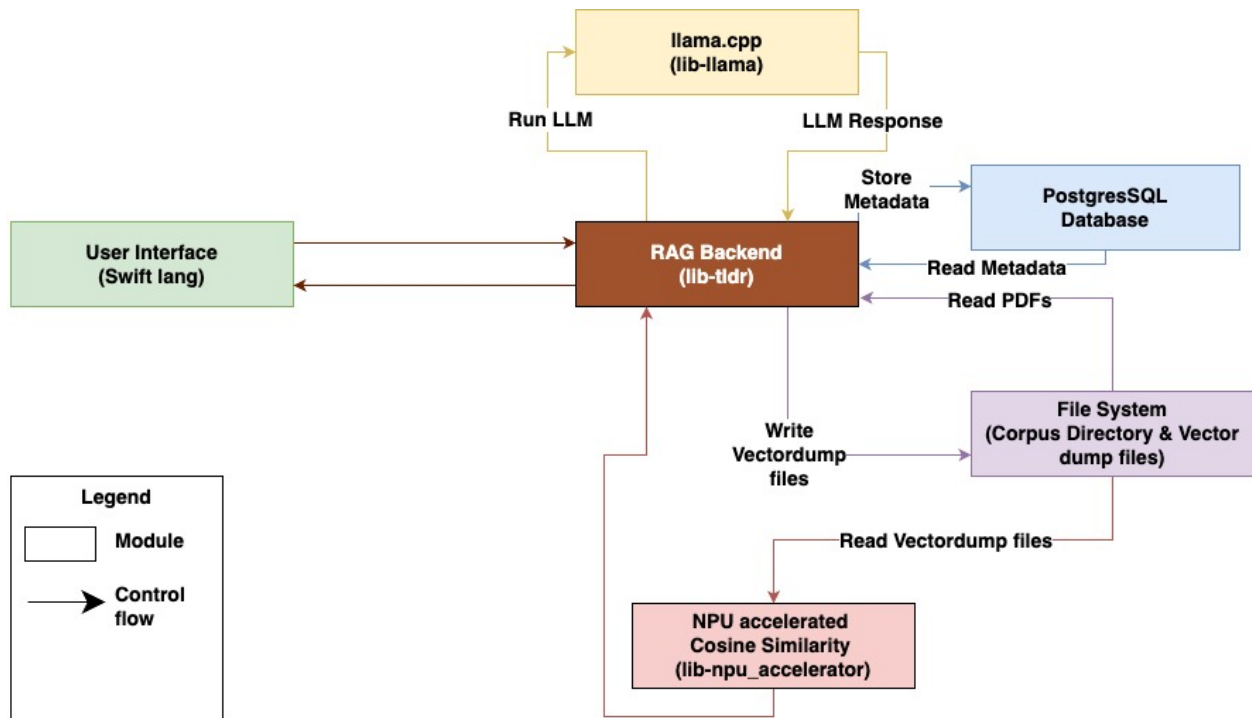


Figure 4.11: TLDR application module interactions

Figure 5.2 illustrates the high-level architecture of the TLDR application, highlighting its modular design and control flow between components.

At the center is the **RAG Backend** (`lib-tldr`), which serves as the core orchestrator. It interfaces with several specialized modules to perform retrieval-augmented generation.

- The **User Interface**, interacts with the RAG backend to initiate queries and display responses to the user.
- **llama.cpp** (`lib-llama`) is linked as a static library and provides language model func-

tionality. The backend calls it for both text generation and embedding, receiving outputs directly in memory.

- **PostgreSQL Database** is used to store and retrieve document metadata and pre-computed embeddings. The RAG backend communicates with it to persist and query structured data including document metadata and text chunks for embeddings.
- The **File System** acts as the persistent store for documents and vector dumps. The backend reads PDF files for embedding and writes embedding outputs into binary vector dump files.
- **NPU-Accelerated Cosine Similarity** (`lib-npu-accelerator`) reads the vector dump files through memory-mapped I/O and executes similarity search using a CoreML model on Apple’s Neural Engine (ANE).

This modular architecture allows each component to focus on a specific responsibility while maintaining efficient communication with the core backend. All dependencies are statically compiled or locally integrated, ensuring portability and performance.

4.2.2 RAG Pipeline Workflow

The workflow of the modules of the system and the RAG pipeline can be divided into 3 logical phases.

- **System Initialization:** Initializes the resources required by the system.
- **Embedding Phase:** Generate embeddings for documents in the source corpus.
- **Retrieval and Generation Phase:** Leverage embedded documents to retrieve information relevant to a query made by the user and generate a response using the LLM.

The following sections show a detailed breakdown of these phases.

System Initialization

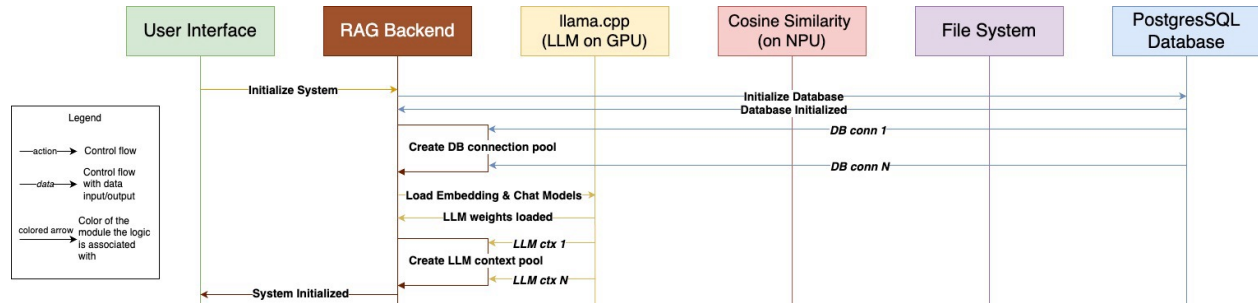


Figure 4.12: RAG Output evaluation metrics [7]

The initialization phase sets up the backend infrastructure and prepares the application for use. This includes:

- The user launches the application, triggering the backend.
- The **RAG Backend** initializes a connection pool to the **PostgreSQL Database**.
- The LLM weights and context pools are loaded via `llama.cpp`, enabling multi-threaded inference.
- Cosine similarity routines are prepared via the **NPU accelerator** module.
- System status is communicated back to the **User Interface**, indicating readiness.

Embedding the Corpus

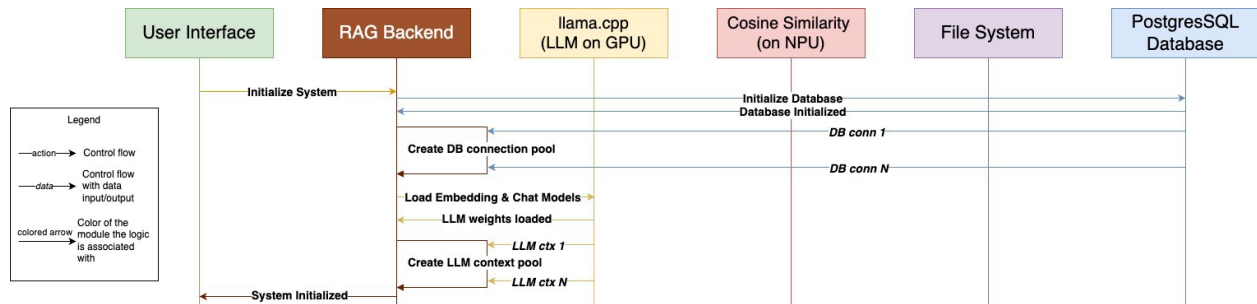


Figure 4.13: Corpus Embedding Workflow Steps

When a user selects a directory to embed the following actions take place:

- The **RAG Backend** scans the specified directory for documents using the **File System**.
- Each document is loaded, chunked, and converted to embeddings using a pre-defined embedder.
- Embeddings, text chunks, and associated metadata are inserted into the **PostgreSQL Database**.
- In parallel, the backend also writes a vector dump file to the **File System**, which stores the hash of each vector for quick access.

This dual-storage mechanism (DB + mmap vector cache) allows fast retrieval during inference while maintaining queryable metadata.

Performing RAG

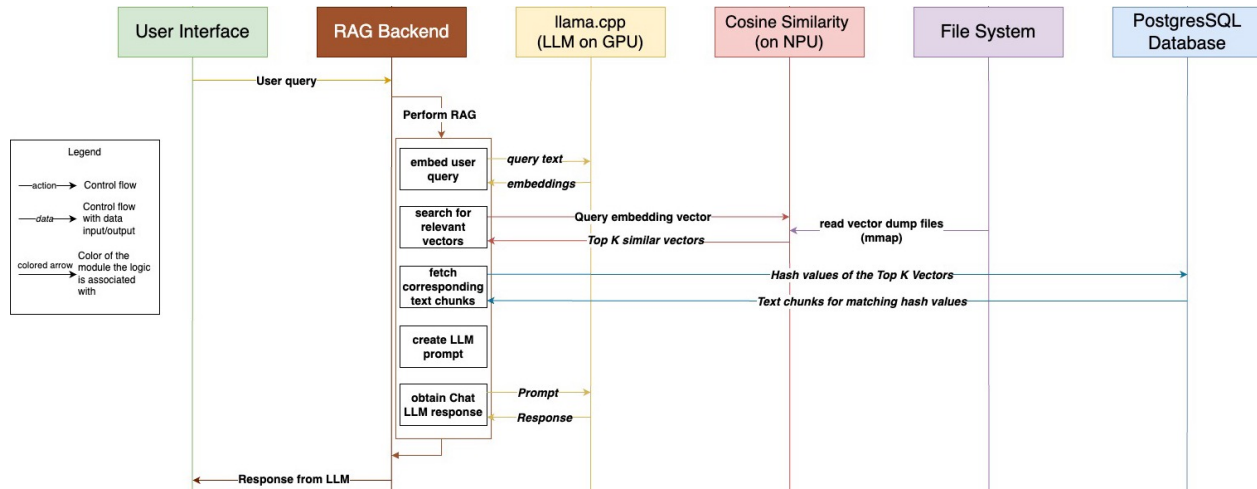


Figure 4.14: RAG workflow steps

Once the corpus is embedded, the system is ready to answer user's queries. When the user proceeds to inquire the system, the following steps are executed:

- The user query is embedded using the same embedding model.
- The embedded query vector is sent to the **Cosine Similarity** module running on the NPU.
- A top- k similarity search is performed against memory-mapped vector files using the NPU, returning hash values of the best matches.
- These hashes are used to retrieve the corresponding text chunks from the **PostgreSQL Database**.
- A prompt is constructed and sent to the **LLM** via `llama.cpp`.
- The generated response is sent back to the **User Interface**.

With these steps, the TLDR project is able to perform retrieval-augmented generation on documents present in the user's device, without requiring any external resources or internet access.

Chapter 5

RESULTS

5.1 Bertscore comparison

BERTScore is an evaluation metric for natural language generation tasks that leverages contextual embeddings from pre-trained language models such as BERT. Unlike traditional n-gram based metrics like BLEU or ROUGE, BERTScore computes similarity by aligning tokens in candidate and reference texts using the cosine similarity of their embeddings. This allows it to better capture semantic similarity even when surface forms differ. BERTScore has been shown to correlate strongly with human judgment across various NLP tasks [66].

We compare sample results from TLDR application using Llama-3.2-1B-Instruct-Q3_K_L with Chatgpt 4o

5.1.1 Result comparisons

- *97 Things Every Software Architect Should Know*
- *Cracking the Coding Interview*
- *189 Programming Questions and Solutions*
- *System Design Interview: An Insider's Guide* by Alex Xu

TLDR Application:

Note: The references to the sources from TLDR's responses have been trimmed for comparison.

1. **User:** What is system design in software engineering?

Response: System design is the process of creating a software system that meets the requirements of the users and the business. It involves defining the system architecture, selecting the appropriate technologies, and designing the user interface. The goal of system design is to create a system that is efficient, scalable, and maintainable. It is a critical process in software engineering that requires careful planning, analysis, and testing.

2. **User:** how to generate prime numbers?

Response: To generate prime numbers, you start with a list of all the numbers up to a certain number (in this case, max). Then, you cross off all the multiples of each prime number you find. This process continues until you've crossed off all the multiples of all the prime numbers you've found. The remaining numbers in the list are the prime numbers. Here's a step-by-step guide: 1. Create a list of all the numbers up to max. 2. Cross off all the multiples of each prime number you find. 3. Continue this process until you've crossed off all the multiples of all the prime numbers you've found. Note: You can use the Sieve of Eratosthenes algorithm to generate a list of prime numbers up to max.

3. **User:** what is caching?

Response: Caching is a technique used to improve the performance of an application by storing frequently accessed data in a temporary storage area, called the cache. This allows the application to retrieve the data quickly, reducing the load on the database or data store. Caching can be used to improve the performance of web applications, database applications, and even mobile applications. It is a technique that can be used to mitigate the problem of the "cold start" problem, where the application has to wait for the data to be retrieved from the database or data store. Caching can be implemented in various ways, including the use of caching tiers, caching strategies, and caching APIs.

ChatGPT:

Note: The same documents mentioned earlier are uploaded to ChatGPT with this prompt: 'understand these documents and answer questions that i will ask based on them and keep the answer concise and precise'.

1. **User:** What is system design in software engineering? (use the previously attached documents as reference)

Response: System design in software engineering is the process of defining the architecture, components, data flow, and interfaces of a software system to meet specified requirements. It focuses on how to build scalable, reliable, and maintainable systems by making high-level decisions on technologies, structure, and communication between parts of the system.

2. **User:** how to generate prime numbers? (use the previously attached documents as reference)

Response: To generate prime numbers, Cracking the Coding Interview outlines two primary methods: trial division and the Sieve of Eratosthenes. Trial division checks whether a number is divisible by any integer up to its square root and is suitable for verifying if a single number is prime. For generating multiple primes efficiently, the Sieve of Eratosthenes is recommended; it initializes a boolean array representing numbers up to a limit and iteratively marks multiples of each prime as non-prime, resulting in an array of prime indicators. This sieve approach offers better performance with a time complexity of $O(n \log \log n)$, making it ideal for generating large sets of prime numbers.

3. **User:** what is caching? (use the previously attached documents as reference)

Response: Caching, as described in System Design Interview: An Insider's Guide and Cracking the Coding Interview, is a technique used to temporarily store frequently accessed or computationally expensive data in fast-access memory (typically RAM) to

improve system performance and reduce load on backend systems like databases. It acts as an intermediary layer between clients and data stores, allowing applications to first check the cache for data before querying the database. Common strategies include read-through caching, expiration policies, and eviction methods like Least Recently Used (LRU), with considerations for consistency, fault tolerance, and scalability.

Quantitative comparison

The table below presents the precision, recall, and F1 score obtained using BERTScore [67]. These results are intended as illustrative examples, as the behavior and output of large language models (LLMs) can vary significantly and are not fully captured by a limited set of responses.

Row	Precision (P)	Recall (R)	F1 Score
1	0.863623	0.897878	0.880417
2	0.851684	0.854167	0.852923
3	0.919073	0.916270	0.917669

Table 5.1: Performance metrics for three different evaluations

The results demonstrate a high degree of similarity between the model responses. Overall, the key takeaway is that the TLDR model is capable of generating responses that are comparable in quality to those produced by cloud-based models, despite operating with significantly fewer resources.

5.2 Screenshots

Simple Demo video: <https://www.youtube.com/watch?v=KaDgJD-KyKA>

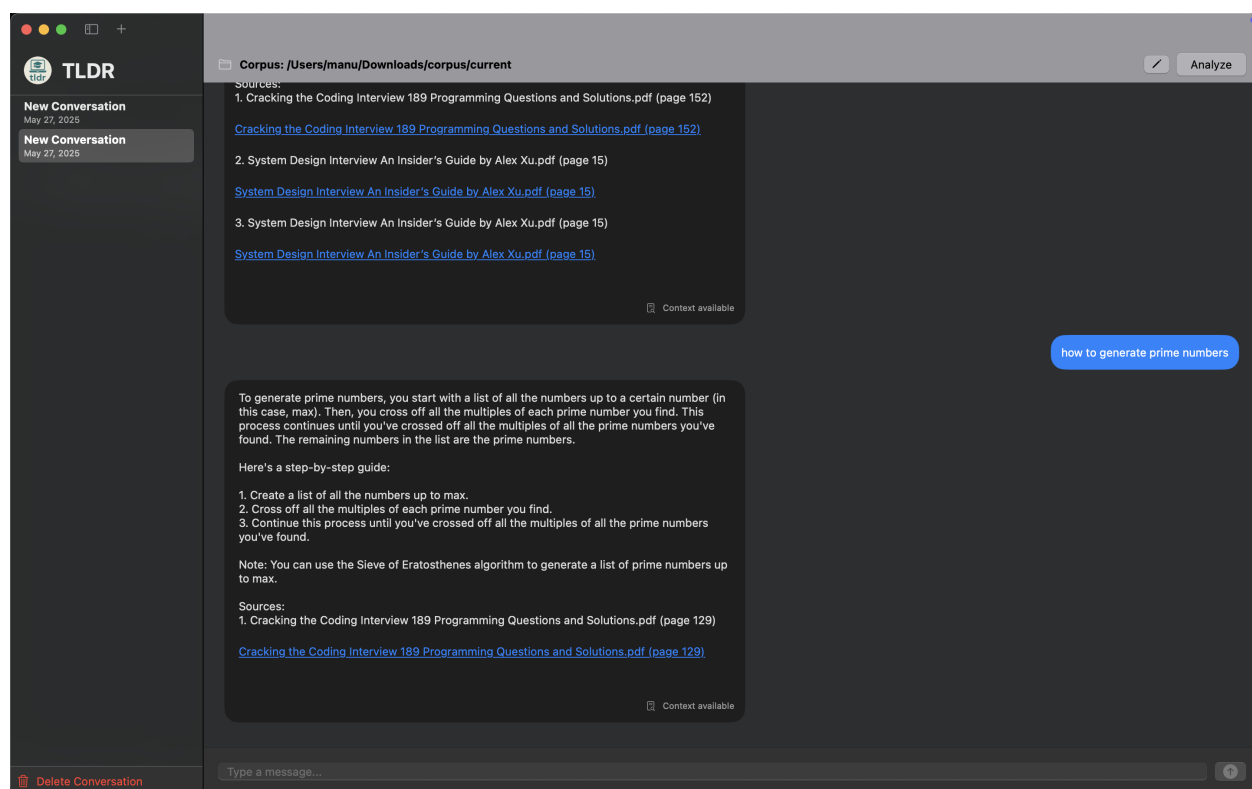


Figure 5.1: TLDR Application Demonstration screenshot 1

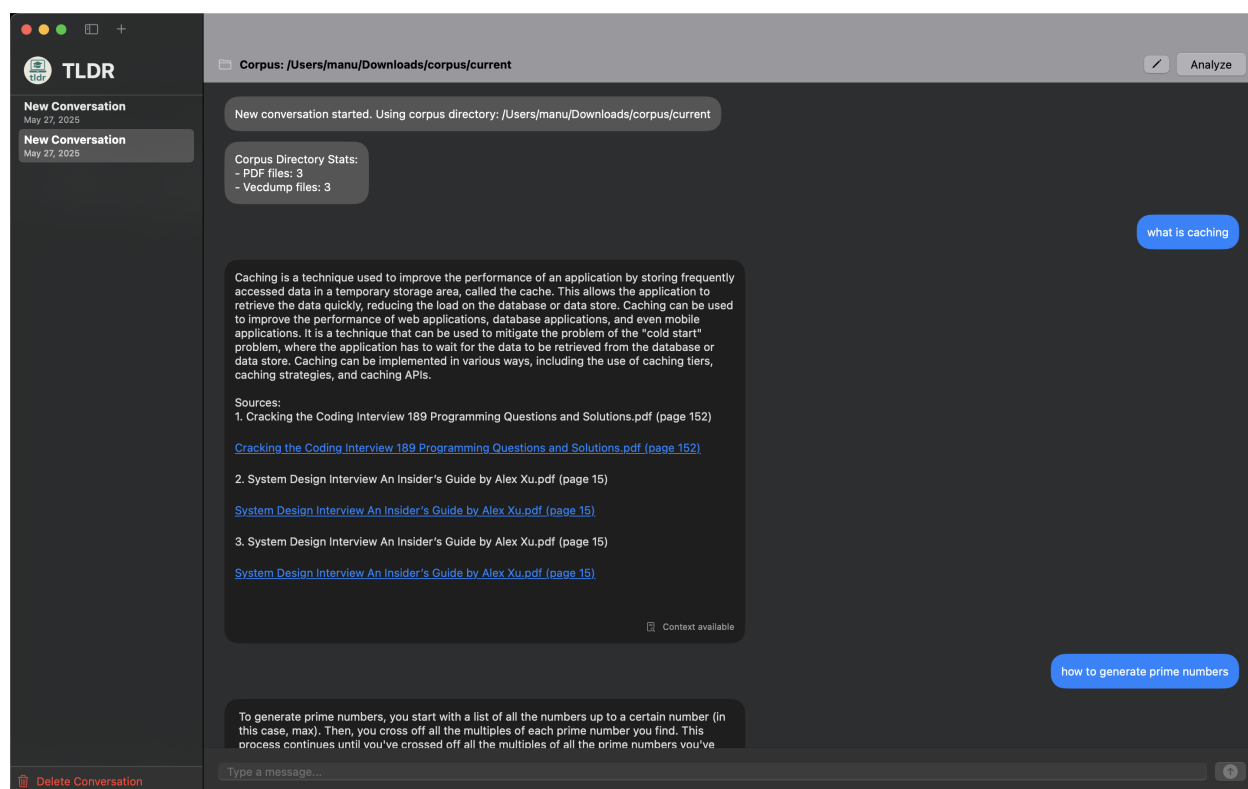


Figure 5.2: TLDR Application Demonstration screenshot 2

Chapter 6

CONCLUSION

6.1 *Takeaways*

The following are the key takeaways of this project:

- a. Apple Neural Engine (ANE) Utilization: This project explored leveraging the under-used NPU for LLM inference, going beyond standard CoreML use cases and showed that it can be viable.
- b. Retrieval-Augmented Generation (RAG) without Internet : Implements a lightweight RAG pipeline that only uses on device resources.
- c. Portability and Ease of use: Application size of less than 1GB that includes everything.
- d. Quantized LLMs: Uses compact models (50–500MB) for efficient, on-device inference allowing for seamless multitasking and manages to obtain results comparable to main-stream cloud LLMs.

6.2 *Limitations*

Following are the main limitations uncovered during the implementation of this project:

- a. While this project demonstrates the feasibility to accelerate the retrieval part of the RAG pipeline, the LLM Inference still only leverages the GPU and does not leverage the NPU.
- b. Indexing(embedding) takes 90+

- c. Many smaller LLMs available currently (3B or less) are only instruction-tuned i.e trained for text completion and not for chat. This can lead to unfavorable responses during chat.
- d. Excessive usage limits to quick draining of power, especially on portable devices.
- e. The breadth of retrieval space is limited not by resources but by the context size of the Chat LLM, since both input and expected output must fit in the LLM context.

6.3 Future work

- a. Enhance data safety mechanisms for vectordump files.
- b. Add NPU backend for GGML and llama.cpp.
- c. Quantize and convert fine-tuned chat-optimized LLMs to the GGUF format.
- d. Implement a dedicated parallelized tokenization module.
- e. Extend the application support beyond Apple Silicon.
- f. Add NPU and GPU acceleration support for SQLite and Postgres vector search extensions (they currently only support CPU).
- g. Create an optimized decoding and tokenization workflow dedicated for embedding (e.g., no KV cache).

BIBLIOGRAPHY

- [1] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in Neural Information Processing Systems*, 30:5998–6008, 2017.
- [2] Nabi. All you need to know about llm text generation, 2024.
- [3] Hollemans. Apple m1 chip architecture, 2020.
- [4] S. Chandra Das, K. Roy, K. Bhawal, A. A. Mahmud, S. M. H. Khan, and S. Alam. Security and privacy challenges of large language models. *arXiv preprint arXiv:2403.09793*, 2024.
- [5] Glenn. Mastering Retrieval-Augmented Generation (RAG) Architecture. <https://blog.stackademic.com/mastering-retrieval-augmented-generation-rag-architecture-unleash-the-power-of-large-language-models/> 2024. Mastering Retrieval-Augmented Generation (RAG) Architecture.
- [6] LangChain. Rag tutorial. <https://python.langchain.com/docs/tutorials/rag/>, 2023. Accessed May 2025.
- [7] Erika Cardenas and Connor Shorten. An overview on rag evaluation, 2023. Accessed: 2025-05-22.
- [8] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.

- [9] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020.
- [10] OpenAI. Gpt-4 technical report. <https://openai.com/research/gpt-4>, 2023. Accessed 2024.
- [11] Yang Liu and Mirella Lapata. Text summarization with pretrained encoders. *arXiv preprint arXiv:1908.08345*, 2019.
- [12] Gautier Izacard and Edouard Grave. Leveraging passage retrieval with generative models for open domain question answering. *arXiv preprint arXiv:2007.01282*, 2021.
- [13] Iz Beltagy, Matthew E Peters, and Arman Cohan. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*, 2020.
- [14] Rishi Bommasani et al. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021.
- [15] Philipp Mattern, Felix Sattler, Hartmut Schmeck, and Bastian Pfitzner. Membership inference attacks against language models via neighbourhood comparison. *arXiv preprint arXiv:2306.04554*, 2023.
- [16] Milad Nasr, Nicholas Carlini, Florian Tramer, and Reza Shokri. Scalable extraction of training data from (production) language models. *arXiv preprint arXiv:2311.17035*, 2023.
- [17] Amy Deschenes and Meg McMahon. A survey on student use of generative ai chatbots for academic research. *Evidence Based Library and Information Practice*, 19(2):2–22, 2024.
- [18] Mohammad Hosseini, Serge Horbach, Tanja Van den Broek, Boudewijn De Bruin, and

- Sietse Wieringa. An exploratory survey about using chatgpt in education, healthcare, and research. *medRxiv*, 2023.
- [19] Benoit Jacob, Skirmantas Kligys, Bo Chen, Menglong Zhu, Matthew Tang, Andrew Howard, Hartwig Adam, and Dmitry Kalenichenko. Quantization and training of neural networks for efficient integer-arithmetic-only inference. *arXiv preprint arXiv:1712.05877*, 2017.
- [20] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- [21] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Kulkarni, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:9459–9474, 2020.
- [22] Apple Inc. Apple neural engine (ane) performance across devices, 2024. Accessed May 13, 2025. Reports ANE performance ranging from 0.6 TOPS (A11) to 38 TOPS (M4).
- [23] Georgi Gerganov. LLaMa.cpp. <https://github.com/ggerganov/llama.cpp>, 2023. C++ framework to run open source LLMs on local hardware.
- [24] Ollama. <https://ollama.com>, 2023. Desktop application to download and run a specific LLM.
- [25] Tinygrad. Apple neural engine reverse engineered for c++. <https://github.com/geohot/tinygrad/tree/master/accel/ane>, 2023. GitHub repository documenting reverse engineering of Apple’s ANE.
- [26] LLamaFile. <https://github.com/Mozilla-Ocho/llamafile>, 2023. Command-line app to package and run an LLM.

- [27] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Matthias Gallé, and Guillaume Lample. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- [28] Jay Alammar. The illustrated transformer, 2018. Accessed May 2025.
- [29] Itay Hubara, Matthieu Courbariaux, Daniel Soudry, Ran El-Yaniv, and Yoshua Bengio. Quantized neural networks: Training neural networks with low precision weights and activations. *arXiv preprint arXiv:1609.07061*, 2016.
- [30] Talamdupula. A guide to quantization in llms. 2024.
- [31] Li. Quantization tech of llms-gguf. 2024.
- [32] Reiner Pope, Alex Tamkin, Ekin Akyürek, Suhas Dathathri, Danny Hernandez, and Andrew Goldie. Efficiently scaling transformer inference. *arXiv preprint arXiv:2211.05102*, 2022.
- [33] Apple Corp. Apple m1 overview. Technical report, Apple Inc., 2020.
- [34] Apple Corp. *Apple Metal GPU Developer Guide*. Apple Inc., 2020.
- [35] Tinygrad. Apple neural engine reverse engineered for c++, 2023.
- [36] Apple CoreML. Coreml – converting pytorch model to coreml. <https://developer.apple.com/documentation/coreml/>, 2024. Accessed May 2025.
- [37] Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. Membership inference attacks against machine learning models. In *Proceedings of the 2017 IEEE Symposium on Security and Privacy (SP)*, pages 3–18. IEEE, 2017.
- [38] man7.org. mmap(2) - linux manual page. <https://man7.org/linux/man-pages/man2/mmap.2.html>, 2024. Accessed May 14, 2025.

- [39] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 2013.
- [40] Bashir. In-context learning, in context. 2023.
- [41] LangChain. Text splitters. https://python.langchain.com/docs/modules/data_connection/document_transformers/text_splitter/, 2023. Accessed May 2025.
- [42] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 3: Open foundation and instruction-tuned models. <https://ai.meta.com/blog/meta-llama-3/>, 2024. Meta AI, Accessed May 2025.
- [43] Nelson F. Liu, Tianyi Shen, Faeze Brahman, Mona Diab, Noah A. Smith, and Yejin Choi. Long in the middle: How language models use long contexts. *arXiv preprint arXiv:2307.03172*, 2023.
- [44] Wenhui Wang, Furu Wei, Li Dong, Hang Bao, Nan Yang, and Ming Zhou. Minilm: Deep self-attention distillation for task-agnostic compression of pre-trained transformers. *arXiv preprint arXiv:2002.10957*, 2020.
- [45] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.
- [46] Rico Sennrich, Barry Haddow, and Alexandra Birch. Neural machine translation of rare words with subword units. *arXiv preprint arXiv:1508.07909*, 2015.
- [47] Xiaohua Wang and Zhenghua Wang. Searching for best practices in retrieval-augmented generation. *arXiv preprint arXiv:2407.01219*, 2024.
- [48] Shinya Sugawara, Hayato Kobayashi, and Manabu Iwasaki. On approximately searching for similar word embeddings. *arXiv preprint arXiv:1604.00417*, 2016.

- [49] Harald Steck, Hieu Pham, Dawen Ding, Lam Thai, Hu Yue, Wei Han, Jeanne Soar, Mihir Sahasrabudhe, and Cosma Shalizi. Is cosine-similarity of embeddings really about similarity? *arXiv preprint arXiv:2403.05440*, 2024.
- [50] Labelbox. Vector similarity search techniques. <https://labelbox.com/blog/how-vector-similarity-search-works/>, 2023.
- [51] Felix L. Hsu. pgvector: Vector similarity search for PostgreSQL. <https://github.com/pgvector/pgvector>, 2023. Accessed: 2025-05-23.
- [52] Ruiqi Guo, Philip Sun, Erik Lindgren, Quan Geng, David Simcha, Felix Chern, and Sanjiv Kumar. Accelerating large-scale inference with anisotropic vector quantization. *arXiv preprint arXiv:1908.10396*, 2020.
- [53] LangChain Contributors. Langchain: Building applications with llms through composability, 2023. <https://www.langchain.com/>.
- [54] LlamaIndex Team. Llamaindex (gpt index), 2023. <https://www.llamaindex.ai/>.
- [55] deepset AI. Haystack: Open source nlp framework for question answering, summarization, and more, 2023. <https://haystack.deepset.ai/>.
- [56] Microsoft Research. Autogen: Enabling next-gen llm applications via multi-agent collaboration, 2023. <https://github.com/microsoft/autogen>.
- [57] Raphaël Gor, Xiao Tan, Zi Ouyang, Zekun Zhang, Jinjie Wei, Yiheng He, Lei Xu, and Yuexin Wu. Routerllm: An expert routing system for cost-efficient inference of large language models. *arXiv preprint arXiv:2309.13285*, 2023.
- [58] RAGAS Contributors. Ragas: Retrieval-augmented generation assessment, 2024. <https://github.com/explodinggradients/ragas>.
- [59] Gerganov. Llama.cpp project. <https://github.com/ggml-org/llama.cpp>, 2020.

- [60] Ollama. Ollama. <https://ollama.com/>, 2023. Accessed May 2025.
- [61] Mozilla-Ocho. Llamafire. <https://github.com/Mozilla-Ocho/llamafire>, 2023. Accessed May 2025.
- [62] ggml org. llama.cpp - simple example. <https://github.com/ggml-org/llama.cpp/tree/master/examples/simple>, 2024. Accessed 2025-05-26.
- [63] ggml org. llama.cpp - server example. <https://github.com/ggml-org/llama.cpp/tree/master/tools/server>, 2024. Accessed 2025-05-26.
- [64] ggml org. llama.cpp - embedding example. <https://github.com/ggml-org/llama.cpp/tree/master/examples/embedding>, 2024. Accessed 2025-05-26.
- [65] ggml org. llama.cpp. <https://github.com/ggml-org/llama.cpp>, 2023. GitHub repository. Accessed May 2025.
- [66] Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. Bertscore: Evaluating text generation with bert. *arXiv preprint arXiv:2004.02047*, 2020.
- [67] Tianyi Zhang*, Varsha Kishore*, Felix Wu*, Kilian Q. Weinberger, and Yoav Artzi. Bertscore: Evaluating text generation with bert. In *International Conference on Learning Representations*, 2020.