



Vaasan yliopisto
UNIVERSITY OF VAASA

Fahad Ibne Fahian

Modular Software Implementation of an Edge Voice Chatbot

Solving Latency and Privacy through Concurrency and Adaptive Listening

School of Technology and Innovations
Master's thesis in
Computing Sciences

Vaasa 2026

UNIVERSITY OF VAASA

School of Technology and Innovations

Author: Fahad Ibne Fahian**Title of the thesis:** Modular Software Implementation of an Edge Voice Chatbot: Solving Latency and Privacy through Concurrency and Adaptive Listening**Degree:** Master of Science in Technology**Degree Programme:** Sustainable and Autonomous Systems**Supervisor:** Jani Boutellier**Year:** 2026 **Pages:** 115

ABSTRACT:

Historically, conversational artificial intelligence systems primarily relied on continuous, cloud-based processing architectures to handle intensive computational workloads. While effective, these traditional paradigms frequently introduced substantial communication latency and raised critical concerns regarding data privacy and security. In recent years, a paradigm shift toward local edge computing emerged as a viable alternative. By deploying advanced language models and voice processing applications directly onto consumer-grade hardware, these privacy and latency issues were significantly mitigated. Consequently, the development of natural, highly responsive, and secure spoken communication with machines was facilitated.

The primary objective of this research was the development and implementation of a fully local, modular voice chatbot application utilizing the Python programming language. The proposed system integrated a sophisticated pipeline comprising Speech-to-Text transcription, a Large Language Model, and Text-to-Speech synthesis. To ensure practical viability, the architecture was implemented and evaluated on consumer-grade edge hardware. Specifically, an Nvidia RTX 3060 GPU equipped with six gigabytes of video random access memory served as the practical minimum baseline to run this local pipeline simultaneously without memory failure. During the initial phase of the research, a baseline sequential pipeline utilized a rigid five-second recording loop, inadvertently forcing the transcription model to process ambient room silence. This generated severe artificial intelligence hallucinations, characterized by the transcription of nonsensical text and repeating numerical sequences. Consequently, this continuous processing of silence severely degraded the language model's context window and precipitated severe Out-of-Memory (OOM) crashes, critically undermining conversational naturalness.

To resolve these system failures, the architecture evolved through a modular approach. Phase 2 implemented adaptive listening by integrating a lightweight Voice Activity Detection module. Executing entirely on the central processing unit to preserve GPU memory, this module dynamically terminated recording after approximately 640 milliseconds of consecutive human silence. In Phase 3, the system was upgraded to a concurrent, multithreaded pipeline. The monolithic loop was decoupled into isolated, asynchronous threads communicating via maxsize queues to eliminate sequential processing bottlenecks.

Ultimately, the integration of these methodologies yielded highly successful results, eliminating the silence-induced hallucinations and memory crashes. The final multithreaded deployment operated stably with a steady-state memory footprint of approximately 3.3 gigabytes of VRAM, successfully proving the viability of deploying highly responsive, natural, and low-latency conversational artificial intelligence locally on constrained edge hardware.

KEYWORDS: Edge Computing, Voice Chatbot, Voice Activity Detection, Large Language Models, Conversational Artificial Intelligence, Low-Latency Processing, Modular Architecture

Contents

1	Introduction	7
1.1	Background	7
1.2	Problem Statement	9
1.3	Research Motivation	11
1.4	Research Objectives	12
1.5	Research Questions	13
1.6	Scope and Delimitations	14
1.7	Thesis Structure	14
2	Literature Review	16
2.1	Evolution of Conversational Agents	16
2.2	Evolution of Neural Language Models	18
2.3	Transformer Architectures and Attention Mechanisms	20
2.4	Edge Computing and Edge AI Deployment Challenges	22
2.5	Speech AI Foundations	23
2.6	Practical Local Conversational AI Systems	24
3	Theoretical Foundations	27
3.1	Transformer Architecture	27
3.1.1	Input Representation and Positional Encoding	28
3.1.2	Network Topology: Encoder-Decoder Structure	29
3.1.3	Autoregressive Language Modeling and Decoder-Only Architectures	30
3.1.4	Inference-Stage Generation Principles	31
3.1.5	Scaling Behavior and Model Alignment	32
3.2	Self-Attention Mechanisms	33
3.2.1	The Theoretical Genesis of Attention Mechanisms	34
3.2.2	Information Retrieval Paradigm: Queries, Keys, and Values	35
3.2.3	The Mathematics of Scaled Dot-Product Attention	36
3.2.4	Multi-Head Attention and Subspace Representations	37
3.2.5	Masked Self-Attention for Autoregressive Modeling	38

3.2.6	Computational Complexity and Path Length Optimization	39
3.3	Quantization Theory	41
3.3.1	Fundamentals of Neural Network Quantization	41
3.3.2	The Outlier Phenomenon and Information Loss	43
3.3.3	Mixed-Precision Decomposition (LLM.int8())	44
3.3.4	Equivalent Smoothing Transformations (SmoothQuant)	44
3.3.5	Second-Order Weight Quantization (GPTQ)	45
3.3.6	Activation-Aware Weight Saliency (AWQ)	46
3.3.7	Implications for Edge Intelligence	48
3.4	Speech Processing Theory	48
3.4.1	Acoustic Signals and Speech Representation	49
3.4.2	Automatic Speech Recognition (ASR) Theory	50
3.4.3	Voice Activity Detection (VAD) Theory	51
3.4.4	Text-to-Speech (TTS) Synthesis Theory	52
3.4.5	Challenges of Low-Latency Voice Interaction	54
3.5	Concurrent Systems Theory	55
3.5.1	Sequential Execution Models and Latency Propagation	55
3.5.2	Asynchronous Processing and Decoupled Pipelines	56
3.5.3	The Producer-Consumer Paradigm and Buffering	57
3.5.4	Backpressure Theory in Low-Latency Environments	58
3.5.5	Synchronization, Race Conditions, and Liveness Properties	59
3.5.6	Complexity Trade-offs in Edge Architectures	61
4	Implementation and Results	62
4.1	Hardware Environment and Constraints	62
4.2	Software Implementation and Configuration	65
4.2.1	Speech Recognition Component	65
4.2.2	Language Model Component	67
4.2.3	Text-to-Speech Component	68
4.2.4	Voice Activity Detection Component	69
4.3	Architectural Evolution	71

4.3.1	Phase 1 - Sequential Baseline	71
4.3.2	Phase 2 - Adaptive Listening	74
4.3.3	Phase 3 - Concurrent Multithreaded Pipeline	81
4.4	Performance Evaluation	86
4.4.1	Memory Utilization	87
4.4.2	Latency Analysis	89
4.4.3	Subjective Testing	94
4.5	Summary of Findings	97
5	Discussion and Future Work	100
5.1	Discussion of Research Findings	100
5.2	Limitations and Hardware Nuances	103
5.3	Future Work	105
6	Conclusions	108
	References	112
	Implementation / Software References	115

Figures

- Figure 1. The Phase-1 sequential processing pipeline. The stop-and-wait architecture creates a bottleneck because each stage must be completed before the next begins. 72
- Figure 2. The Phase-2 adaptive listening architecture. The CPU-based VAD module dynamically gates the audio capture, passing data to the inference models only after detecting the end of user speech. 78
- Figure 3. The Phase-3 concurrent multithreaded architecture. Independent processing threads communicate through bounded queues while interruption signals enable responsive conversational interaction. 85
- Figure 4. Per-turn processing times recorded during steady-state Phase 3 execution. Speech-to-text processing, language model generation, and text-to-speech synthesis remained consistently below 2.5 seconds total across multiple conversational turns. 92

Tables

- Table 1. Steady-State VRAM Utilization of Core Conversational Pipeline Components 88
- Table 2. Comparison of Average Processing Latency 91

1 Introduction

1.1 Background

The ambition to develop artificial intelligence capable of engaging in natural, seamless communication with humans has driven decades of computing research, originating from early philosophical frameworks that proposed indistinguishability in conversation as a benchmark for machine intelligence (Turing, 1950). In recent years, this field has experienced a profound acceleration fueled by the availability of massive datasets and the development of highly complex neural architectures. Specifically, the scaling of Large Language Models (LLMs) has fundamentally transformed the capabilities of natural language processing. By training massive autoregressive networks on vast corpora of text, these models have emerged as powerful few-shot and zero-shot learners, capable of generating coherent, context-aware text across a wide array of topics (Brown et al., 2020). Through subsequent alignment techniques, such as instruction tuning and reinforcement learning from human feedback, these models have been explicitly optimized to follow user intent safely and helpfully, establishing the modern standard for conversational artificial intelligence (Ouyang et al., 2022).

While text-based chatbots have achieved remarkable proficiency, human communication is inherently multimodal, with spoken language representing the most natural and efficient interface. Transitioning advanced language models into fully realized voice conversational agents requires the integration of advanced acoustic modeling. This multimodal pipeline necessitates robust automatic speech recognition to transcribe raw audio waveforms into text (Hinton et al., 2012; Radford et al., 2023), voice activity detection to intelligently segment human speech from ambient background noise (Ramírez et al., 2004), and end-to-end text-to-speech synthesis to generate expressive, naturalistic audio responses (Kim et al., 2021).

Historically, deploying this computationally intensive combination of speech and language models has relied almost exclusively on continuous, cloud-based computing

paradigms. In the traditional cloud architecture, edge devices act primarily as passive sensors and data transmitters, continuously streaming captured audio payloads over a wide area network to centralized, remote data centers where the heavy neural network inference is performed (Shi et al., 2016). While the cloud offers vast, centralized computational throughput, it introduces critical systemic limitations that inherently degrade the conversational experience.

The most immediate and disruptive limitation of cloud-reliant artificial intelligence is end-to-end communication latency. Because large-scale data centers are frequently geographically distant from the end user, every spoken interaction incurs inevitable network propagation delays, multi-hop routing, and potential bandwidth congestion (Satyanarayanan, 2017). Interactive human communication is exceptionally sensitive to timing; cognitive and perceptual thresholds require end-to-end delays to be sufficiently low to maintain a natural, overlapping interaction flow. Consequently, relying on a distant cloud server is fundamentally misaligned with the strict temporal requirements of rapid conversational turn-taking.

Furthermore, continuous cloud connectivity raises severe privacy and security vulnerabilities. Voice data inherently contains highly sensitive, biometric, and contextual personal information. Transmitting unencrypted or raw audio sensor data across the internet to a third-party hub forces users to relinquish control over their privacy and creates a vulnerability to interception or centralized data breaches (Satyanarayanan, 2017; Zhou et al., 2019). Additionally, the cloud-based paradigm suffers from strict internet dependency, rendering voice interfaces practically useless in offline, remote, or unstable network environments.

To resolve these inherent limitations, edge computing has emerged as a critical architectural paradigm. Edge computing pushes cloud services and computing tasks away from the network core directly to the network edges, in closer physical proximity to the data sources (Shi et al., 2016). By offloading neural network inference onto local

consumer devices, significant privacy preservation can be achieved, as voice processing occurs entirely on-device and sensitive data is shielded from external cloud infrastructure (Zhou et al., 2019). Equally important, physical proximity to the computational hardware bypasses wide area network delays, promising the highly responsive, reduced latency required for true conversational usability (Satyanarayanan, 2017).

However, transferring state-of-the-art conversational artificial intelligence to the edge ecosystem is non-trivial. Edge devices are constrained by strict physical limitations regarding processing capability, memory footprint, and energy consumption (Wang et al., 2020). Deploying massive language models and acoustic processing pipelines on constrained hardware necessitates sophisticated software execution strategies and architectural optimizations to prevent system instability and resource exhaustion (Lin et al., 2024; Xiao et al., 2023). Thus, a fundamental motivation exists to conceptualize and engineer efficient, fully localized software architectures capable of hosting highly responsive voice chatbots autonomously on the edge. This research was conducted in conjunction with the Human–AI Interaction and Innovation Nexus at Vaasa.

1.2 Problem Statement

Conventional conversational artificial intelligence systems have historically relied upon continuous, cloud-based computing paradigms to handle the substantial computational workloads required by deep neural networks. In this traditional architecture, end devices function primarily as passive data consumers and raw data producers, continuously streaming captured audio payloads over a wide area network to centralized, remote data centers for processing (Shi et al., 2016). While this paradigm exploits economies of scale and vast centralized computing power, it introduces critical systemic limitations regarding latency, privacy, network dependency, and conversational responsiveness.

A primary limitation of cloud-reliant architectures is end-to-end communication latency. Because large-scale data centers are frequently consolidated and geographically distant

from the end user, every interaction incurs inevitable propagation delays, multi-hop queuing, and routing congestion (Satyanarayanan, 2017). Interactive human communication is exceptionally sensitive to timing; cognitive and perceptual tasks require end-to-end delays to be tightly controlled to less than a few tens of milliseconds to maintain a natural interaction flow (Satyanarayanan, 2017). Consequently, relying on a distant cloud data center is fundamentally misaligned with the strict temporal requirements of natural, overlapping human conversation.

Furthermore, continuous cloud connectivity raises severe privacy and security concerns. Voice data contains highly sensitive, biometric, and contextual personal information. Transmitting unencrypted or raw audio sensor data across the internet to a third-party hub inherently shifts control over privacy policies away from the user (Satyanarayanan, 2017). The centralization of such data creates a vulnerability to interception, unauthorized deep analytics, and major data breaches (Zhou et al., 2019). Beyond latency and privacy, the cloud-based paradigm suffers from strict internet dependency and scalability constraints. Applications relying on persistent connectivity are highly vulnerable to network failures, denial-of-service attacks, and bandwidth degradation (Satyanarayanan, 2017). As billions of smart interfaces and internet-of-things devices are deployed, the cumulative ingress bandwidth required to stream continuous, high-fidelity audio threatens to saturate metropolitan area networks, making centralized processing structurally unsustainable (Shi et al., 2016).

Edge computing has emerged as a promising architectural solution to these challenges by pushing cloud services and computing tasks from the network core to the network edges, in closer physical proximity to the data sources (Shi et al., 2016; Satyanarayanan, 2017). However, transferring artificial intelligence capabilities to the edge ecosystem is highly non-trivial due to severe resource constraints. Edge devices are typically limited by computing capability, memory footprint, and energy consumption budgets (Wang et al., 2020). Deploying state-of-the-art, large-scale deep learning models, such as

transformer-based language generators and acoustic synthesizers onto these constrained devices introduces a significant processing bottleneck.

Without rigorous architectural redesign, simply transferring cloud-based sequential processing pipelines to edge hardware inevitably results in system instability, memory exhaustion, and unacceptably slow processing times. When speech recognition, language generation, and audio synthesis are executed in a rigid, linear sequence on constrained devices, computational resources are locked in waiting states, dramatically increasing the time to response. Thus, a fundamental problem remains in engineering modular, localized software architectures capable of bridging the gap between heavy neural networks and constrained edge environments while preserving the speed required for natural conversation.

1.3 Research Motivation

The imperative to resolve the inherent limitations of cloud-based conversational agents drives the motivation for edge-based local artificial intelligence. By offloading complex neural network inference directly onto local edge devices, a host of structural, temporal, and security advantages can be realized, rendering local conversational artificial intelligence highly desirable (Zhou et al., 2019).

A key motivation for shifting inference to the edge is enhanced privacy preservation. When an edge device serves as the primary and only point of contact in the infrastructure, all voice processing occurs entirely on-device (Satyanarayanan, 2017). Because voice processing occurs locally rather than being transmitted to external cloud infrastructure, privacy exposure can be substantially reduced. This isolation enables supports stronger user control over privacy and mitigates the risks associated with centralized data storage and external interception (Zhou et al., 2019).

Equally critical is the achievement of highly responsive, reduced latency. Physical proximity to the computational hardware entirely bypasses wide area network

transmission delays, bandwidth throttling, and external server routing (Shi et al., 2016). For voice conversational artificial intelligence, minimizing processing latency is paramount to replicating the rapid turn-taking and spontaneous flow characteristic of natural human speech. A localized architecture holds the promise of delivering response times that support true conversational usability. Furthermore, operating conversational artificial intelligence locally guarantees offline autonomy. Consumer devices equipped with localized inference capabilities are completely insulated from internet outages, variable signal strength, and remote server downtime, ensuring that natural language interfaces remain reliable and accessible in disconnected or unstable environments (Satyanarayanan, 2017).

Democratizing this technology, however, requires making it accessible on standard, consumer-grade hardware rather than relying on enterprise-grade server clusters or specialized micro-datacenters. Deploying large-scale language models and complex acoustic processing pipelines on constrained hardware necessitates sophisticated software execution strategies (Wang et al., 2020; Lin et al., 2024). Consequently, there is a strong motivation to innovate beyond traditional linear software design. Exploring architectural optimization strategies for efficient localized inference is vital to preventing resource exhaustion (Xiao et al., 2023). Developing an architecture that intelligently orchestrates speech recognition, and language generation is a fundamental prerequisite for establishing a seamless, fully localized conversational agent capable of running autonomously on the edge.

1.4 Research Objectives

The primary objective of this research is the conceptualization, development, and evaluation of a modular, fully localized software architecture capable of supporting a highly responsive edge voice chatbot. The research aims to shift the conversational artificial intelligence paradigm away from cloud dependency by investigating whether sophisticated natural language understanding and acoustic synthesis can be reliably

hosted on consumer-grade edge hardware without sacrificing conversational naturalness or system stability.

A central objective is the seamless architectural integration of three distinct neural network modalities: automatic speech recognition, large language model generation, and text-to-speech synthesis, into a unified, locally executing framework. To ensure this integration is practically viable on devices with strictly constrained memory boundaries, the research must identify and apply optimal structural configurations, efficient model deployment strategies, and parameter limitations that prevent memory exhaustion and system failure.

Furthermore, a critical objective is to overcome the severe latency limitations associated with traditional sequential inference processing. The research aims to minimize end-to-end conversational latency by exploring architectural optimization approaches that decouple and streamline the processing loop. By exploring more efficient architectural processing strategies that can manage audio ingestion, transcription, and response generation dynamically, the research seeks to neutralize inherent computational bottlenecks and establish a robust theoretical framework for natural, localized human-machine interaction.

1.5 Research Questions

To fulfill the stated objectives and systematically address the theoretical and architectural challenges of edge-based conversational artificial intelligence, this research is guided by the following core questions:

1. How can a modular, fully localized voice chatbot architecture be designed to support natural spoken interaction on consumer-grade edge hardware?
2. How can speech recognition, language generation, and speech synthesis components be integrated efficiently under constrained memory resources without triggering system instability?

3. What architectural optimization strategies can reduce end-to-end conversational latency while preserving system stability in localized conversational AI?

1.6 Scope and Delimitations

The scope of this research is confined to the architectural design, software structure, and theoretical optimization of a fully localized conversational artificial intelligence pipeline. The study focuses on integrating pre-existing, state-of-the-art neural network models for transcription, language generation, and speech synthesis into a cohesive software framework operating entirely at the network edge.

Several delimitations have been established to maintain the theoretical and academic focus of the thesis. First, the research exclusively targets inference-stage deployment and structural optimization; the fundamental pre-training, structural modification, or foundational fine-tuning of the underlying neural architectures is strictly outside the scope of this work. Second, the study explicitly excludes any reliance on cloud-based Application Programming Interfaces (APIs) or distributed network offloading, maintaining a pure focus on autonomous local execution.

Additionally, while the architectural design adheres to strict memory budgets typical of consumer-grade environments, the thesis does not encompass custom hardware design, microchip fabrication, or low-level systems programming. The evaluation is focused on systemic stability, memory preservation through architectural design, and reductions in conversational latency achieved via software orchestration.

1.7 Thesis Structure

The remainder of this thesis is logically structured to progress from theoretical foundations to architectural design and evaluation, maintaining a rigorous separation between scientific theory and practical implementation.

Chapter 2 provides a comprehensive literature review, examining the academic context of edge computing, the evolution of local language models, and previous architectures utilized in conversational systems. It contextualizes the proposed research against existing paradigms and identifies the gaps in current localized inference strategies.

Chapter 3 explores the theoretical framework underlying the system. This encompasses the scientific principles of language model neural architectures, sequence-to-sequence transcription, the mathematical mechanics of attention and transformers, and the theories of model quantization and acoustic modeling that make edge deployment viable.

Chapter 4 details the practical implementation of the proposed software architecture. It outlines the specific integration of the chatbot pipeline, the deployment constraints, and the structural methodologies applied to achieve latency reduction and memory stability on consumer hardware.

Chapter 5 offers a discussion of the project's successes, analyzes identified architectural limitations, and proposes directions for future research. Finally, Chapter 6 summarizes the core findings and concludes the thesis.

2 Literature Review

This chapter presents a literature review examining the scientific evolution of conversational artificial intelligence, language model architectures, and edge-oriented deployment strategies relevant to localized voice chatbot systems. The chapter first reviews the historical progression from symbolic conversational agents to modern Transformer-based large language models, followed by practical deployment paradigms involving speech processing, model quantization, and edge inference optimization. Finally, the chapter contextualizes the present research within the broader academic and technological landscape by contrasting localized modular conversational architectures against conventional cloud-dependent conversational AI systems.

2.1 Evolution of Conversational Agents

The historical evolution of conversational artificial intelligence spans over seven decades, characterized by a profound transition from early rule-based symbolic manipulation to highly complex, data-driven neural architectures. The theoretical foundation for machine conversation was formally established when the question of whether machines could think was substituted with a practical indistinguishability test known as the Imitation Game (Turing, 1950). This foundational framework proposed that if a human interrogator could not reliably distinguish between the conversational text responses of a machine and those of a human, the machine could be considered to exhibit intelligent behavior (Turing, 1950). This paradigm intentionally separated physical characteristics from intellectual capacities, establishing an enduring academic benchmark for evaluating natural language communication between humans and computers.

Early attempts to construct conversational agents relied exclusively on symbolic pattern matching and syntactic transformation, eschewing any genuine semantic understanding. The most iconic of these early systems was ELIZA, a program designed to mimic a Rogerian psychotherapist (Weizenbaum, 1966). ELIZA operated by scanning input text

for specific keywords, isolating minimal contextual information, and applying predefined decomposition and reassembly rules to generate responses (Weizenbaum, 1966). Despite its reliance on simple text manipulation and pronoun swapping, ELIZA demonstrated that the illusion of human comprehension could be effectively maintained, highlighting the human psychological tendency to anthropomorphize machines that mirror user inputs (Weizenbaum, 1966).

Following ELIZA, the PARRY system introduced a psychological dimension to conversational agents by simulating a paranoid belief system (Colby et al., 1971). Unlike ELIZA's purely linguistic transformations, PARRY integrated internal affective variables, specifically fear, anger, and mistrust; which dynamically influenced the model's output strategies (Colby et al., 1971). By tracking the context of the interview and modifying its affective states in response to the human interlocutor, PARRY represented a significant advancement in modeling the internal psychological state of an artificial conversationalist (Colby et al., 1971).

The progression of symbolic, rule-based systems culminated in more sophisticated architectures such as the Artificial Linguistic Internet Computer Entity (ALICE), which utilized the Artificial Intelligence Markup Language (AIML) (Wallace, 2009). ALICE expanded the stimulus-response paradigm by deploying tens of thousands of conversational categories managed by algorithms that facilitated highly efficient pattern matching. Furthermore, ALICE demonstrated the value of supervised refinement through backward-looking log file analysis, leveraging the Zipf-like distribution of natural language inputs to systematically author new rules for frequently occurring queries (Wallace, 2009).

However, despite these structural advancements, symbolic chatbots suffered from severe systemic limitations. As the demand for open-domain conversation grew, the inability of rule-based systems to scale without exhaustive, manual engineering became a critical bottleneck. Because it is practically impossible to manually script responses for

every conceivable linguistic variation, symbolic systems remained exceptionally brittle when confronted with novel, open-domain dialogue. These inherent limitations catalyzed a paradigm shift toward neural conversational artificial intelligence, wherein dialogue generation is framed as a sequential, data-driven decision-making process optimized via deep learning (Gao et al., 2018).

2.2 Evolution of Neural Language Models

In this modern era, Recurrent Neural Networks (RNNs) emerged as the natural computational architecture for sequential natural language problems. However, early traditional RNNs struggled fundamentally with long-term dependencies due to the mathematical challenge of vanishing and exploding error backflow during the backpropagation through time process (Hochreiter & Schmidhuber, 1997). This architectural limitation prevented networks from retaining context over extended conversational sequences. The problem was elegantly solved by the introduction of the Long Short-Term Memory (LSTM) architecture, which enforced constant error flow through internal memory cells utilizing specialized multiplicative gate units; specifically input, output, and forget gates, to protect memory contents from irrelevant perturbations (Hochreiter & Schmidhuber, 1997).

The successful application of LSTMs directly enabled the Sequence to Sequence (Seq2Seq) learning framework, which revolutionized machine translation and end-to-end conversational modeling (Sutskever et al., 2014). The Seq2Seq architecture eschewed hand-crafted linguistic rules entirely, utilizing a multilayered LSTM encoder to map a variable-length input sequence into a fixed-dimensional continuous vector representation, which a secondary LSTM decoder then conditioned upon to sequentially generate target text (Sutskever et al., 2014). This framework established the foundation for fully data-driven generation of conversational responses (Gao et al., 2018).

While the Seq2Seq architecture represented a monumental leap forward, compressing an entire sentence into a single fixed-length vector created a severe informational

bottleneck that degraded performance on longer sequences (Bahdanau et al., 2015). This structural limitation was resolved by the introduction of the attention mechanism, which allowed the decoder to adaptively soft-search a sequence of annotations computed by the encoder for each generated word (Bahdanau et al., 2015). By calculating a dynamic context vector based on alignment scores, the model learned to autonomously determine which parts of the source sequence to pay attention to during generation, freeing neural models from the fixed-length vector constraint.

The introduction of attention mechanisms established the theoretical foundation for the Transformer architecture, which subsequently became the dominant paradigm in modern conversational artificial intelligence. However, as these models proliferated, a critical systemic flaw emerged: optimizing massive networks purely for a next-token prediction objective resulted in models that were fundamentally misaligned with human conversational intent (Ouyang et al., 2022). Predicting the next word on a diverse internet corpus does not inherently teach a model to be a helpful conversational assistant; consequently, base language models frequently generated outputs that were untruthful, toxic, biased, or unhelpful (Brown et al., 2020; Ouyang et al., 2022).

To rectify this misalignment and establish modern neural conversational modeling, the paradigm of instruction tuning and Reinforcement Learning from Human Feedback (RLHF) was established (Ouyang et al., 2022). Researchers developed a multi-stage fine-tuning process where models were first supervised via human-written demonstrations of desired conversational behavior. Next, a reward model was trained on vast datasets of human-labeled comparisons, learning to mathematically predict which generated response a human judge would prefer. Finally, the base language model's policy was optimized against this reward model using proximal policy optimization reinforcement learning (Ouyang et al., 2022). RLHF proved that aligning language models to human intent yielded significantly higher conversational quality and safety, definitively establishing the theoretical architecture of the modern conversational assistant (Ouyang et al., 2022).

2.3 Transformer Architectures and Attention Mechanisms

The profound success of the attention mechanism paved the way for the Transformer architecture, proposed by Vaswani et al. (2017), which eschewed complex recurrence and convolutions entirely. The Transformer is based solely on attention mechanisms, drawing global dependencies between inputs and outputs without relying on sequence-aligned RNNs. At the core of this architecture is self-attention, an attention mechanism relating different positions of a single sequence in order to compute a rich, contextual representation of that sequence (Vaswani et al., 2017).

The specific attention function employed by the Transformer is termed "Scaled Dot-Product Attention." This mathematical operation maps a query and a set of key-value pairs to an output, computed as a weighted sum of the values. To counteract the effect of dot products growing massive in magnitude; which pushes the softmax function into regions with vanishing gradients, the Transformer scales the dot products by the inverse square root of the key dimension, ensuring stable gradient flow during backpropagation (Vaswani et al., 2017).

Building upon this, the architecture utilizes Multi-Head Attention to dramatically enhance the model's representational capacity. Instead of performing a single attention function, multi-head attention linearly projects the queries, keys, and values multiple times with different learned projections. This parallel attention structure allows the model to jointly attend to information from different representation subspaces at different positions simultaneously (Vaswani et al., 2017). Because the Transformer dispenses with recurrence entirely, it employs Positional Encodings, utilizing sine and cosine functions of different frequencies added to the input embeddings to inject information regarding the relative or absolute position of the tokens in the sequence (Vaswani et al., 2017).

The original Transformer retains an encoder-decoder structure. The encoder is composed of a stack of identical layers containing multi-head self-attention and position-wise feed-forward networks, utilizing residual connections and layer normalization (Vaswani et al., 2017). The decoder follows a similar architecture but inserts an additional sub-layer to perform multi-head attention over the output of the encoder stack. Furthermore, the self-attention sub-layer within the decoder incorporates a masking mechanism, preventing leftward information flow to preserve the fundamental auto-regressive property required for sequence generation (Vaswani et al., 2017).

The architectural shift from recurrence to self-attention provides massive computational and parallelization advantages. A self-attention layer connects all positions within a sequence using a constant number of sequentially executed operations, reducing the maximum path length between any two input and output positions to $O(1)$, compared to the $O(n)$ sequential operations required by recurrent networks (Vaswani et al., 2017). The elimination of sequential computation allows matrix operations to be highly parallelized across modern hardware, drastically reducing training time.

This parallelizability sparked an exponential scaling of language models, splintering the Transformer into distinct architectural sub-families. Deep bidirectional models, such as BERT (Bidirectional Encoder Representations from Transformers), utilized the Transformer's encoder stack and a masked language modeling objective to pre-train semantic representations by jointly conditioning on both left and right context, establishing deep natural language understanding (Devlin et al., 2019). Conversely, the advancement of modern conversational artificial intelligence has been predominantly driven by the scaling of autoregressive, decoder-only Transformer models. Generative Pre-trained Transformers (GPT) leveraged this decoder-only architecture, demonstrating that exponentially scaling parameters and training data could result in powerful few-shot and zero-shot learners capable of complex text synthesis purely via in-context learning (Brown et al., 2020). The transition toward this highly parallelizable, autoregressive

Transformer architecture forms the definitive theoretical bedrock upon which all contemporary state-of-the-art conversational language models operate.

2.4 Edge Computing and Edge AI Deployment Challenges

The theoretical evolution of conversational artificial intelligence has culminated in models capable of highly sophisticated natural language understanding. However, transitioning these deep learning architectures from theoretical frameworks into practical, deployable systems introduces substantial engineering challenges. Historically, the deployment of large language models and complex acoustic pipelines relied almost exclusively on cloud computing infrastructures. Cloud-based paradigms centralized the massive computational workloads required for inference, leveraging vast data centers to process user requests (Shi et al., 2016). While this approach provided necessary computational throughput, it intrinsically suffered from critical systemic limitations.

A primary limitation of cloud-reliant architectures is end-to-end communication latency. Because large-scale data centers are frequently geographically distant from the end user, every interaction incurs inevitable propagation delays and wide-area network routing congestion (Satyanarayanan, 2017). Interactive human communication is exceptionally sensitive to timing; natural conversation demands tightly controlled latency. Consequently, the delays inherent to cloud routing become a critical bottleneck for voice chatbots. Furthermore, continuous cloud connectivity raises severe vulnerabilities regarding data privacy. Transmitting unencrypted or raw audio sensor data across the internet to a third-party hub inherently shifts control over privacy policies away from the user (Satyanarayanan, 2017). Finally, the cloud-based paradigm suffers from strict internet dependency, rendering systems useless in offline environments (Wang et al., 2020).

To overcome these limitations, the research community has increasingly focused on edge computing, a paradigm that pushes computational services away from the network core directly to the proximity of the data source (Zhou et al., 2019). Edge intelligence

facilitates the local execution of machine learning models on consumer-grade hardware, ensuring privacy preservation by keeping sensitive audio and text data on-device while simultaneously reducing network-induced latency (Wang et al., 2020).

However, deploying conversational artificial intelligence at the edge is highly non-trivial. Edge devices are constrained by strict physical limitations, primarily regarding processing capability, thermal dissipation, and Video Random Access Memory (VRAM) capacity. In transformer-based language models, the autoregressive generation phase processes one token at a time. This decoding stage is fundamentally memory-bound because the entire massive weight matrix of the model must be loaded from memory to the processing cores for every single generated token (Lin et al., 2024). Consumer-grade hardware features restricted VRAM capacities and relatively narrow memory bandwidths compared to enterprise data centers. Attempting to load an uncompressed multi-billion parameter model alongside acoustic processing models simultaneously into local VRAM inevitably triggers system memory exhaustion.

2.5 Speech AI Foundations

A functional voice conversational agent requires a comprehensive pipeline integrating speech recognition and acoustic synthesis. For local Automatic Speech Recognition (ASR), the Whisper architecture has established the contemporary standard. Whisper utilizes a sequence-to-sequence Transformer trained via large-scale weak supervision on 680,000 hours of multilingual audio (Radford et al., 2023). By training on a vastly diverse dataset without relying exclusively on gold-standard human transcripts, the model generalized effectively to real-world acoustic variations, achieving zero-shot transfer capabilities that remove the need for dataset-specific fine-tuning (Radford et al., 2023). The release of its underlying software ecosystem has allowed researchers to integrate highly robust, offline transcription directly into edge pipelines.

To ensure continuous, real-time audio capture does not overwhelm computational resources or degrade the language model's context window, practical systems integrate

Voice Activity Detection (VAD). VAD algorithms dynamically segment temporal speech features, isolating human speech from ambient background noise (Ramirez et al., 2004). By actively gating the audio stream and truncating periods of silence, VAD modules prevent the ASR from hallucinating transcriptions based on ambient noise, thereby stabilizing the downstream language model inference and conserving hardware resources.

To complete the conversational loop, Text-to-Speech (TTS) systems are deployed. Contemporary TTS models utilize conditional variational autoencoders coupled with adversarial learning networks to synthesize raw speech waveforms end-to-end, producing highly naturalistic voice outputs (Kim et al., 2021). These architectures allow for the modular integration of expressive voice synthesis directly alongside local language generation, forming the necessary acoustic foundation for edge-based conversational agents.

2.6 Practical Local Conversational AI Systems

Overcoming the hardware constraints and memory bottlenecks discussed previously requires an ecosystem of highly optimized, open-weight models and advanced mathematical quantization strategies. A primary enabler of practical edge deployment has been the democratization of open-weight large language models. The Qwen family of models, particularly the Qwen2 series, was developed to push the boundaries of accessible language generation (Yang et al., 2024). Crucially for edge deployment, smaller dense models within this series, such as the 1.5 billion parameter variants, were explicitly designed for portable devices and consumer-grade hardware. These models have been extensively instruction-tuned using RLHF to align with human conversational preferences, rendering them highly adept at dialogue generation without requiring cloud resources (Yang et al., 2024).

To fit these models into local VRAM budgets, deployment relies heavily on mathematical model quantization. Quantization systematically maps high-precision floating-point

weights into discrete, lower-bit integer representations. However, conventional quantization methods suffer from severe accuracy degradation when applied to large language models due to the spontaneous emergence of systematic large-magnitude outlier features in the network's activations (Dettmers et al., 2022). Naive quantization fails because these extreme outliers stretch the quantization grid, effectively extinguishing the model's learned information (Dettmers et al., 2022).

To combat the outlier phenomenon, several sophisticated post-training quantization algorithms have been developed. The LLM.int8() method addressed this by mathematically isolating outlier dimensions, retaining them in high-precision 16-bit floating-point formats, while aggressively quantizing the remaining 99.9% of regular weights into 8-bit matrices (Dettmers et al., 2022). An alternative approach, SmoothQuant, solves the outlier problem by mathematically migrating the quantization difficulty from the activations to the weights prior to execution. By dividing activation channels by a smoothing factor and scaling weights in the reverse direction, SmoothQuant enables efficient 8-bit weight and activation (W8A8) quantization without runtime mixed-precision overhead (Xiao et al., 2023).

To achieve even higher compression rates, Generative Pre-trained Transformer Quantization (GPTQ) introduced an approximate second-order optimization method for weight-only compression. GPTQ utilizes inverse Hessian information to iteratively adjust unquantized weights to compensate for the mathematical error introduced by rounding, successfully compressing networks down to 3 or 4 bits per weight with negligible accuracy degradation (Frantar et al., 2023). Furthermore, Activation-aware Weight Quantization (AWQ) was introduced to optimize execution based on the insight that protecting just 1% of the most salient weights prevents almost all quantization degradation. AWQ profiles the distribution of activations to identify these salient weight channels and applies per-channel scaling to minimize relative quantization error, making it highly effective for instruction-tuned conversational models (Lin et al., 2024).

These theoretical advancements in quantization are operationalized through highly optimized software inference engines. The llama.cpp project established a revolutionary inference framework explicitly designed to execute quantized large language models with minimal overhead. To facilitate this, the Generalized Generic Universal Form (GGUF) file format was standardized, allowing models to be seamlessly packaged alongside their tokenizers and quantized weight tensors into a single distributable binary. By utilizing unified memory architectures, these frameworks enable intelligent offloading between CPUs and GPUs. Together, these methodologies and deployment frameworks collectively establish the practical feasibility of hosting highly responsive, private, and autonomous conversational artificial intelligence directly at the network edge.

3 Theoretical Foundations

This chapter details the theoretical foundations underlying modern edge-based conversational artificial intelligence systems. While Chapter 2 examined the historical evolution and literature surrounding conversational AI and edge computing, this chapter focuses specifically on the mathematical and architectural principles that enable contemporary large language models, speech processing systems, and efficient edge deployment. The chapter first examines Transformer architectures and autoregressive language modeling principles, followed by self-attention mechanisms, model quantization theory, speech processing foundations, and concurrent systems theory relevant to low-latency conversational architectures.

3.1 Transformer Architecture

The fundamental objective of sequence modeling and transduction in natural language processing is to map an input sequence of discrete symbol representations into an output sequence. Prior to the introduction of the Transformer, the dominant theoretical paradigms for these sequence-to-sequence mapping tasks were based on complex recurrent or convolutional neural networks (Vaswani et al., 2017). To rigorously understand the architectural innovations of the Transformer, it is first necessary to examine the mathematical and structural limitations of the recurrent architectures it replaced.

Recurrent Neural Networks (RNNs) inherently factor computation along the symbol positions of the input and output sequences. By aligning these positions to discrete steps in computation time, an RNN generates a sequence of hidden states, denoted as h_t , as a deterministic function of the previous hidden state h_{t-1} and the input for position t (Vaswani et al., 2017). This purely sequential nature imposes two severe theoretical constraints. First, the step-by-step processing requirement strictly precludes computational parallelization within training examples, which creates a massive computational bottleneck as sequence lengths increase (Vaswani et al., 2017). Second,

traditional recurrent networks suffer from the phenomenon of vanishing and exploding error backflow; during the backpropagation through time process, the temporal evolution of error signals depends exponentially on the size of the network's weights, preventing the model from effectively learning long-range dependencies (Hochreiter & Schmidhuber, 1997).

While architectures such as Long Short-Term Memory (LSTM) mitigated vanishing gradients by enforcing constant error flow through internal memory cells (Hochreiter & Schmidhuber, 1997), and Sequence to Sequence (Seq2Seq) frameworks allowed mapping variable-length inputs to outputs (Sutskever et al., 2014), the fundamental sequential constraint remained. Furthermore, traditional Seq2Seq models forced the encoder to squash all semantic information from a potentially long source sentence into a single, fixed-length continuous vector, creating a severe informational bottleneck (Bahdanau et al., 2015).

The Transformer architecture, introduced by Vaswani et al. (2017), provided a definitive theoretical solution to these limitations by eschewing recurrence and convolutions entirely. Instead, the Transformer relies exclusively on an attention mechanism to draw global dependencies between inputs and outputs. By doing so, the Transformer connects all positions within a sequence using a constant number of sequentially executed operations. In computational complexity theory, this reduces the maximum path length that forward and backward signals must traverse between any two arbitrary input and output positions to an $O(1)$ operation, compared to the $O(n)$ sequential operations required by recurrent networks (Vaswani et al., 2017). This $O(1)$ path length is the foundational architectural property that allows Transformers to seamlessly learn long-range linguistic dependencies while enabling massive hardware parallelization.

3.1.1 Input Representation and Positional Encoding

To process discrete natural language within a continuous neural network, the input text must first be tokenized into subword units and mapped into a continuous vector space. In the Transformer architecture, learned embeddings are utilized to convert the input

tokens into vectors of a fixed dimensionality, denoted as d_{model} . In the original formulation, d_{model} is uniformly set to 512 across all layers (Vaswani et al., 2017).

Because the Transformer dispenses with recurrence, the architecture inherently lacks a built-in mechanism to recognize the sequential order of the tokens. To enable the model to utilize sequence order, information regarding the relative or absolute position of the tokens must be explicitly injected into the input representation. The Transformer achieves this through Positional Encodings, which are added directly to the input embeddings at the bottom of the network (Vaswani et al., 2017).

The positional encodings share the same d_{model} dimensionality as the embeddings, allowing the two vectors to be summed. The theoretical formulation of these encodings relies on sine and cosine functions of different frequencies: $PE_{(pos,2i)} = \sin(pos/10000^{2i/d_{model}})$ $PE_{(pos,2i+1)} = \cos(pos/10000^{2i/d_{model}})$ where pos represents the absolute position of the token in the sequence and i represents the feature dimension (Vaswani et al., 2017). Each dimension of the positional encoding corresponds to a sinusoid. This specific mathematical formulation was selected because it allows the model to easily learn to attend to relative positions; for any fixed offset k , PE_{pos+k} can be represented as a linear function of PE_{pos} (Vaswani et al., 2017). Furthermore, this continuous functional representation allows the model to theoretically extrapolate to sequence lengths longer than those encountered during the training phase.

3.1.2 Network Topology: Encoder-Decoder Structure

The overarching topology of the original Transformer follows an encoder-decoder structure. The encoder's theoretical objective is to map an input sequence of symbol representations (x_1, \dots, x_n) into a sequence of continuous contextual representations $z = (z_1, \dots, z_n)$. Given the continuous representations z , the decoder then generates an output sequence of symbols (y_1, \dots, y_m) one element at a time (Vaswani et al., 2017).

The encoder is composed of a stack of N identical layers (originally $N = 6$). Each layer consists of two distinct sub-layers. The first sub-layer is a multi-head self-attention

mechanism, the mathematical specifics of which will be detailed in the subsequent section. The second sub-layer is a fully connected Position-wise Feed-Forward Network (FFN). This FFN is applied to each position in the sequence separately and identically, consisting of two linear transformations with a Rectified Linear Unit (ReLU) activation function applied in between: $FFN(x) = \max(0, xW_1 + b_1) W_2 + b_2$ where W_1 and W_2 represent the learned weight matrices, and b_1 and b_2 represent the learned biases (Vaswani et al., 2017). While the linear transformations are identical across different sequence positions, they utilize entirely different parameters from layer to layer within the stack.

To facilitate stable gradient flow during deep network training, the Transformer employs a residual connection around each of the two sub-layers, followed immediately by layer normalization. Mathematically, the output of each sub-layer is defined as $LayerNorm(x + Sublayer(x))$, where $Sublayer(x)$ is the function implemented by the sub-layer itself (Vaswani et al., 2017). This architectural requirement dictates that all sub-layers in the model, including the embedding layers, produce outputs of the exact same dimension d_{model} .

The decoder is similarly composed of a stack of N identical layers, but it introduces a third sub-layer designed to perform multi-head attention over the output z generated by the encoder stack (Vaswani et al., 2017). Crucially, the self-attention sub-layer within the decoder is structurally modified to prevent positions from attending to subsequent positions. This masking mechanism, combined with the fact that output embeddings are offset by one position, guarantees that the prediction for position i can only depend on the known, previously generated outputs at positions less than i (Vaswani et al., 2017).

3.1.3 Autoregressive Language Modeling and Decoder-Only Architectures

While the original Transformer utilized a bidirectional encoder paired with a decoder for translation tasks, the evolution of modern conversational artificial intelligence has been heavily driven by decoder-only Transformer architectures. Massive language models, such as the Generative Pre-trained Transformer (GPT) series, strip away the encoder

entirely, relying exclusively on a deep stack of masked decoder blocks to perform autoregressive language modeling (Brown et al., 2020).

Autoregressive models process sequences by decomposing the joint probability of a sequence of tokens into the product of ordered conditional probabilities. Given a sequence $y = (y_1, \dots, y_T)$, the joint probability is defined as: $p(y) = \prod_{t=1}^T p(y_t | y_1, \dots, y_{t-1})$ (Sutskever et al., 2014). At each computational step, the model is strictly auto-regressive, consuming the previously generated symbols as additional input when calculating the probability distribution for the next symbol (Vaswani et al., 2017).

In a decoder-only architecture, the masking mechanism inside the self-attention layers is the absolute structural requirement that enforces this autoregressive property. By setting all values in the input of the attention softmax function that correspond to illegal (future) connections to $-\infty$, the model ensures strict left-to-right information flow (Vaswani et al., 2017). During the pre-training phase, these models are optimized to maximize the log-probability of predicting the correct next token across vast corpora of text.

3.1.4 Inference-Stage Generation Principles

The theoretical mechanics of training a Transformer differ significantly from how the architecture behaves during inference. During training, the mathematical objective is to maximize the likelihood of the training data through gradient descent and backpropagation. During inference, however, the model must autonomously generate novel sequences.

To convert the continuous output vectors of the final decoder layer into discrete token predictions, the model applies a learned linear transformation followed by a softmax function. This projects the d_{model} -dimensional output into a massive vector whose length equals the entire vocabulary size, outputting a normalized probability distribution over all possible next tokens (Vaswani et al., 2017).

Selecting the optimal sequence of tokens from these probabilities requires specialized decoding algorithms. A purely greedy decoding approach simply selects the single token with the highest probability at each time step. However, greedy search can easily trap the generation process in suboptimal local maxima. Consequently, advanced inference engines frequently utilize beam search. A beam search decoder maintains a small number, B , of the most likely partial hypotheses at each time step (Sutskever et al., 2014). At each generative step, every partial hypothesis in the beam is expanded with every possible token in the vocabulary. The system evaluates the cumulative log-probability of these expanded sequences, discarding all but the top B hypotheses (Sutskever et al., 2014). Additionally, generation parameters such as length penalties and sampling temperatures are applied to modify the mathematical sharpness of the softmax distribution, ensuring that the generated conversational responses remain both coherent and highly diverse (Brown et al., 2020).

3.1.5 Scaling Behavior and Model Alignment

The theoretical elegance of the purely attention-based, autoregressive Transformer is that its performance scales predictably with increased computational resources. Research into the scaling laws of language models demonstrates that exponentially increasing the number of network parameters, the size of the pre-training dataset, and the computational throughput results in smooth, predictable power-law improvements in language modeling loss (Brown et al., 2020).

As decoder-only Transformers scale into the hundreds of billions of parameters, they undergo a theoretical phase shift, evolving into powerful zero-shot and few-shot learners. At this immense scale, models are capable of performing complex natural language tasks, such as logical reasoning, translation, and conversational synthesis, purely via in-context learning. Instead of requiring task-specific architectural modifications or gradient updates, the model infers the desired task directly from the natural language prompt provided in its context window (Brown et al., 2020).

However, optimizing massive networks purely for a mathematical next-token prediction objective results in models that are fundamentally misaligned with human conversational intent. Predicting the statistically most likely next word on a diverse internet corpus does not inherently teach a model to be a helpful or harmless conversational assistant; base models frequently generate outputs that are untruthful, toxic, or unhelpful (Ouyang et al., 2022).

To bridge the gap between the autoregressive pre-training objective and safe conversational utility, modern Transformers undergo an alignment process utilizing Reinforcement Learning from Human Feedback (RLHF) (Ouyang et al., 2022). In this theoretical framework, a secondary reward model is trained on a dataset of human-labeled comparisons to mathematically predict which generated response a human judge would prefer. The base language model's generative policy is then explicitly fine-tuned to maximize this reward signal using the Proximal Policy Optimization (PPO) reinforcement learning algorithm (Ouyang et al., 2022). This structural alignment ensures that the immense representational power of the Transformer architecture is safely constrained and optimally steered for human-machine conversational interaction.

3.2 Self-Attention Mechanisms

Self-attention mechanisms constitute the mathematical core of modern Transformer architectures and large language models. Unlike recurrent neural networks, which process tokens sequentially, self-attention enables models to dynamically compute contextual relationships between all tokens within a sequence simultaneously. This section examines the theoretical foundations, mathematical formulations, and computational properties of self-attention mechanisms that enable efficient large-scale language modeling. The foundational architecture of the Transformer relies entirely on the mechanism of self-attention to draw global dependencies between inputs and outputs, dispensing with recurrence and convolutions entirely (Vaswani et al., 2017). To understand the theoretical power of modern large language models, a rigorous mathematical examination of the self-attention mechanism is required. This section

explores the theoretical genesis of attention, the mathematical formulations of scaled dot-product and multi-head attention, and the computational complexity paradigms that afford self-attention its superiority in long-range dependency modeling.

3.2.1 The Theoretical Genesis of Attention Mechanisms

The theoretical necessity for attention mechanisms originated from the mathematical limitations of early neural sequence-to-sequence (Seq2Seq) architectures. In a conventional recurrent Seq2Seq model, an encoder processes an input sequence of vectors step-by-step to compute a hidden state (Sutskever et al., 2014). The fundamental flaw in this traditional approach is the requirement that the encoder compress all semantic and syntactic information of a variable-length source sentence into a single, fixed-dimensional continuous vector (Bahdanau et al., 2015). From an information-theoretic perspective, squashing a potentially massive amount of information into a rigid vector space creates a severe informational bottleneck, causing the predictive performance of the neural network to degrade precipitously as sequence length and complexity increase (Bahdanau et al., 2015).

To resolve this bottleneck, the concept of attention was introduced to allow a model to automatically soft-search for a sequence of annotations computed by the encoder, evaluating them dynamically for each target word generated (Bahdanau et al., 2015). In this precursor to self-attention, often termed additive attention; the context vector c_i for generating the i -th target word is computed as a weighted sum of the encoder annotations h_j : $c_i = \sum_{j=1}^{T_x} \alpha_{ij} h_j$ (Bahdanau et al., 2015). The weight α_{ij} of each annotation h_j is determined by a softmax-normalized alignment score, which quantifies how well the inputs around position j and the output at position i match mathematically.

While this formulation successfully freed neural models from the fixed-length vector constraint, early attention mechanisms were still strictly utilized in conjunction with a recurrent neural network to maintain sequence alignment (Vaswani et al., 2017). The

profound theoretical leap of the Transformer was the realization that attention mechanisms could be generalized to relate different positions of a single sequence in order to compute a rich representation of that sequence entirely independently of recurrence. This paradigm is defined as self-attention, or intra-attention (Vaswani et al., 2017).

3.2.2 Information Retrieval Paradigm: Queries, Keys, and Values

Self-attention fundamentally operates as a differentiable information retrieval system mapped onto a continuous vector space. An attention function can be described mathematically as mapping a query and a set of key-value pairs to an output, where the query, keys, values, and the resulting output are all represented as continuous high-dimensional vectors (Vaswani et al., 2017).

In the context of self-attention, the queries, keys, and values are not derived from distinct input and target sequences; rather, they all originate from the exact same place; the output of the previous layer within the network (Vaswani et al., 2017). Given an input sequence of continuous representations, the self-attention layer applies three distinct, learned linear transformations to project the input sequence into three separate matrices: the Query matrix (Q), the Key matrix (K), and the Value matrix (V).

The theoretical abstraction here is that the Query matrix represents the current token's search intent, what information the token is looking for. The Key matrix represents the addressing system; what information each token in the sequence holds. Finally, the Value matrix represents the actual semantic content that will be aggregated and passed forward through the network. The output of the self-attention layer is computed as a mathematically weighted sum of the Values, where the specific weight assigned to each individual Value is determined by a compatibility function computed between the Query and the corresponding Key (Vaswani et al., 2017).

3.2.3 The Mathematics of Scaled Dot-Product Attention

The specific attention mechanism utilized by the Transformer architecture is formally defined as "Scaled Dot-Product Attention" (Vaswani et al., 2017). The inputs to this function consist of queries and keys of dimension d_k , and values of dimension d_v . To compute the compatibility score between a query and a key, the algorithm utilizes a simple dot product (multiplicative attention), which is highly space-efficient and exceptionally fast in practice due to its compatibility with highly optimized matrix multiplication algorithms (Vaswani et al., 2017).

To process entire sequences simultaneously, the queries are packed together into a matrix Q , while the keys and values are packed into matrices K and V , respectively. The unnormalized compatibility scores for all queries and keys are computed via the matrix multiplication QK^T . A softmax function is then applied to these scores to yield a normalized probability distribution of attention weights, ensuring that the weights sum to 1.0. These normalized weights are finally multiplied by the Value matrix V . The complete scaled dot-product attention function is mathematically defined as:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (\text{Vaswani et al., 2017}).$$

The inclusion of the scaling factor $\frac{1}{\sqrt{d_k}}$ is a critical mathematical necessity for preserving the stability of the neural network during training. While additive attention and dot-product attention perform similarly well for small representation dimensions, standard unscaled dot-product attention drastically underperforms as the dimensionality d_k grows large (Vaswani et al., 2017). The theoretical explanation for this degradation lies in the variance of the dot products. Assuming that the individual components of the query vector q and the key vector k are independent random variables with a mean of 0 and a variance of 1, their dot product $q \cdot k = \sum_{i=1}^{d_k} q_i k_i$ will have a mean of 0 but a variance of d_k (Vaswani et al., 2017).

As the dimension d_k scales (e.g., $d_k = 64$ in standard Transformer heads), the variance of the dot product grows proportionally large. Consequently, the dot products generate scalar values that are extremely large in magnitude. When these large magnitude values are passed into the softmax function, they push the softmax outputs into regions where the gradient is infinitesimally small (Vaswani et al., 2017). Because the derivative of the softmax function approaches zero at its extremes, learning ceases, as the vanishing gradients prevent backpropagation from successfully updating the network's weights. To rigorously counteract this variance inflation, the dot products are divided by the standard deviation $\sqrt{d_k}$. This scaling mathematically restricts the variance of the logits back to 1, ensuring that the softmax function operates in its dynamic middle region where gradients flow optimally (Vaswani et al., 2017).

3.2.4 Multi-Head Attention and Subspace Representations

While a single self-attention operation is powerful, it is mathematically limited by the averaging effect of the softmax distribution. If a single attention head is forced to attend to a sequence, the probability mass is distributed across the keys, which inhibits the model's ability to jointly focus on multiple distinct, highly relevant representation subspaces simultaneously (Vaswani et al., 2017). Natural language is inherently multifaceted; a single token must simultaneously track syntactic relationships (such as subject-verb agreement), semantic meanings, and positional proximity.

To overcome the limitations of a single attention distribution, the architecture expands scaled dot-product attention into Multi-Head Attention. Instead of performing a single attention function using d_{model} -dimensional keys, values, and queries, the theoretical approach is to linearly project the queries, keys, and values h times using different, learned linear projection matrices (Vaswani et al., 2017). Specifically, they are projected to d_k , d_k , and d_v dimensions, respectively.

On each of these h projected versions of queries, keys, and values, the scaled dot-product attention function is performed in parallel. This yields h distinct d_v -dimensional

output value vectors. These individual head outputs are then mathematically concatenated along the feature dimension and once again linearly projected via a final weight matrix W^O to produce the final output of the multi-head attention layer. The operation is defined as: $MultiHead(Q, K, V) = Concat(head_1, \dots, head_h)W^O$ where each individual head is computed as: $head_i = Attention(QW_i^Q, KW_i^K, VW_i^V)$ (Vaswani et al., 2017).

The parameter matrices are defined as $W_i^Q \in \mathbb{R}^{d_{model} \times d_k}$, $W_i^K \in \mathbb{R}^{d_{model} \times d_k}$, $W_i^V \in \mathbb{R}^{d_{model} \times d_v}$, and $W^O \in \mathbb{R}^{hd_v \times d_{model}}$ (Vaswani et al., 2017). By reducing the dimension of each head (e.g., $d_k = d_v = d_{model}/h$), the total computational cost of multi-head attention is kept roughly equivalent to that of single-head attention with full dimensionality (Vaswani et al., 2017). Crucially, this architectural division allows the model to jointly attend to information from different representation subspaces at different positions. One head may learn to map pronouns to their antecedents, while another simultaneously tracks local grammatical structures, drastically enhancing the expressivity of the contextual representation (Devlin et al., 2019; Vaswani et al., 2017).

3.2.5 Masked Self-Attention for Autoregressive Modeling

The standard self-attention mechanism allows every position in the input sequence to attend to every other position, generating a deep bidirectional representation. While this bidirectional conditioning is ideal for natural language understanding tasks (Devlin et al., 2019), it is strictly prohibited in the context of sequence generation. Large language models operate on the principle of autoregressive language modeling, decomposing the joint probability of a sequence into a product of conditional probabilities where each token is predicted solely based on the preceding tokens (Brown et al., 2020; Sutskever et al., 2014).

If bidirectional self-attention were utilized during autoregressive generation, the model would be able to trivially "see the future." The query for position t would compute high compatibility scores with the keys for positions $t + 1, t + 2$, etc., rendering the next-

token prediction objective mathematically invalid. To strictly preserve the autoregressive property, the self-attention sub-layers within the decoder are structurally modified into Masked Self-Attention layers (Vaswani et al., 2017).

Masking is implemented directly inside the scaled dot-product attention function prior to the application of the softmax normalization. A strictly lower-triangular mask matrix is constructed, where all elements above the diagonal are set to $-\infty$, and all elements on and below the diagonal are set to 0. This mask is added to the unnormalized attention logits matrix $\frac{QK^T}{\sqrt{d_k}}$. Because the exponential function evaluates $e^{-\infty}$ as 0, passing the masked matrix through the softmax function mathematically forces the attention weights for all future illegal connections to exactly 0.0. Consequently, the prediction for position i can only depend on the known outputs at positions strictly less than i (Vaswani et al., 2017). This specific structural alteration ensures that leftward information flow is strictly maintained, enabling the Transformer to function as an autoregressive generative engine.

3.2.6 Computational Complexity and Path Length Optimization

To theoretically justify the absolute dominance of self-attention mechanisms in modern large language models, their computational complexity and sequence processing properties must be analyzed against the traditional recurrent and convolutional layers they replaced. The theoretical desiderata for mapping a variable-length sequence of symbol representations (x_1, \dots, x_n) to another sequence of equal length (z_1, \dots, z_n) revolves around three primary metrics: total computational complexity per layer, the amount of computation that can be parallelized (measured by the minimum number of sequential operations required), and the path length between long-range dependencies in the network (Vaswani et al., 2017).

In terms of algorithmic complexity, a standard recurrent neural network layer processes an input sequence using $O(n)$ sequential operations, resulting in an overall

computational complexity of $O(n \cdot d^2)$ per layer, where n is the sequence length and d is the representation dimensionality (Vaswani et al., 2017). The sequential nature of the $O(n)$ operations fundamentally precludes parallelization across the temporal dimension of a sequence. A hidden state h_t simply cannot be computed until h_{t-1} is fully resolved. Conversely, a self-attention layer connects all positions within a sequence using a constant, $O(1)$ number of sequentially executed operations (Vaswani et al., 2017). Because the compatibility between every query and every key is calculated via dense matrix multiplication simultaneously, self-attention allows for total parallelization across the sequence. The computational complexity of a self-attention layer is $O(n^2 \cdot d)$ (Vaswani et al., 2017). From a theoretical standpoint, self-attention is definitively faster and more computationally efficient than recurrent layers whenever the sequence length n is smaller than the representation dimensionality d (Vaswani et al., 2017). Given that state-of-the-art language models utilize representation dimensions extending into the thousands (e.g., $d_{model} = 4096$ or higher), self-attention maintains a vast computational advantage for standard sequence lengths.

Beyond raw computational speed, the most critical theoretical advantage of self-attention lies in its routing of gradient flow and its ability to learn long-range dependencies. In deep learning architectures, learning dependencies between distant positions in a sequence is a key challenge (Hochreiter & Schmidhuber, 1997). The difficulty of learning such dependencies is mathematically proportional to the length of the paths that forward and backward signals must traverse through the network structure (Vaswani et al., 2017). In recurrent networks, the maximum path length between any arbitrary input position and any output position is $O(n)$, as the gradient signal must be backpropagated step-by-step through the unfolded time sequence, risking exponential decay (Hochreiter & Schmidhuber, 1997). In convolutional networks utilizing contiguous kernels, achieving a global receptive field requires a stack of $O(n/k)$ convolutional layers, while dilated convolutions require $O(\log_k(n))$ layers, scaling the path length accordingly (Vaswani et al., 2017).

Self-attention eradicates this path length dependency entirely. Because the scaled dot-product attention computes direct, pairwise alignments between every single token in the input sequence regardless of their positional distance, the maximum path length between any two arbitrary positions in the sequence is strictly reduced to $O(1)$ (Vaswani et al., 2017). The shorter these paths are, the easier it is to backpropagate error signals without vanishing or exploding gradients. By guaranteeing an $O(1)$ path length, the self-attention mechanism provides the optimal mathematical framework for capturing the profound, long-range semantic correlations inherent in human language, establishing it as the definitive theoretical core of modern natural language processing.

3.3 Quantization Theory

The exponential scaling of large language models has fundamentally transformed the capabilities of artificial intelligence, yet this scaling introduces profound computational and memory challenges. The Autoregressive Transformer models inherently require sequential token generation during inference, creating substantial computational and memory bandwidth demands (Lin et al., 2024). Consequently, the decoding stage is fundamentally memory-bound; the entire massive weight matrix of the neural network must be loaded from memory to the processing cores for every single generated token (Lin et al., 2024). When executed using standard high-precision representations, the memory footprint and the bandwidth required to fetch these weights heavily exceed the capacities of localized and edge-based computing environments (Wang et al., 2020; Zhou et al., 2019). To bridge the gap between massive neural architectures and constrained deployment environments, the mathematical framework of model quantization is required.

3.3.1 Fundamentals of Neural Network Quantization

In deep learning, neural network parameters, specifically weights and activations are conventionally stored and computed using high-precision floating-point representations, most commonly 32-bit (FP32) or 16-bit (FP16) formats. Floating-point formats allocate

specific bits to a sign, an exponent, and a fraction (or mantissa), granting the neural network a vast dynamic range capable of representing both infinitesimally small gradients and massive activation values with high precision. However, this precision comes at the cost of significant storage requirements and high computational energy per operation (Wang et al., 2020).

Quantization is the theoretical process of systematically mapping these high-precision, continuous floating-point values into a set of discrete, lower-bit integer representations, such as 8-bit (INT8), 4-bit (INT4), or even 3-bit arrays. By representing numbers with fewer bits, quantization drastically shrinks the physical memory footprint of the model, reduces the required memory bandwidth for tensor fetching, and allows the system to utilize highly efficient integer arithmetic logic units during matrix multiplication (Zhou et al., 2019).

The most prevalent mathematical approach to this mapping is uniform absolute maximum (absmax) quantization. Absmax quantization linearly scales floating-point inputs into a symmetric integer range, such as $[-127, 127]$ for INT8. Given an FP16 input matrix X_{FP16} , the discrete mapping is achieved by dividing the tensor by its absolute maximum magnitude (equivalent to its infinity norm), multiplying by the maximum integer range, and applying a rounding function: $X_{INT8} = \lfloor 127 \cdot \frac{X_{FP16}}{\max(|X_{FP16}|)} \rfloor$ (Dettmers et al., 2022; Xiao et al., 2023). To recover the approximated floating-point values during or after matrix multiplication, the system applies a dequantization step by multiplying the integer outputs by the inverse scaling constant. While computationally efficient, this mapping is inherently lossy. The rounding function $\lfloor \cdot \rfloor$ introduces a strict quantization error, and the magnitude of this error is directly governed by the scaling step size, which relies entirely on the maximum absolute value within the tensor (Xiao et al., 2023).

3.3.2 The Outlier Phenomenon and Information Loss

While quantization has been successfully applied to smaller convolutional networks and early language models, applying naive absmax quantization to billion-parameter language models results in catastrophic accuracy degradation. The theoretical reason for this failure lies in the spontaneous emergence of highly systematic, extreme-magnitude features within the network's hidden states (Dettmers et al., 2022).

Research reveals that as autoregressive language models are scaled beyond approximately 6.7 billion parameters, a phase shift occurs wherein massive outlier features emerge in the activation layers (Dettmers et al., 2022). These activation outliers are exceptionally sparse; representing approximately 0.1% of all feature dimensions, yet their magnitudes can be up to 20 times larger than the values in the remaining regular dimensions (Dettmers et al., 2022). Furthermore, these outliers are highly systematic; if an outlier appears in a specific channel, it persistently appears across all tokens for that specific channel, concentrating the variance (Xiao et al., 2023).

When naive per-tensor absmax quantization is applied to a matrix containing these outliers, the extreme magnitudes stretch the quantization grid entirely out of proportion. Because the scaling denominator $\max(|X_{FP16}|)$ is dominated by the outliers, the vast majority of regular, non-outlier features are compressed into an exceptionally narrow band of discrete bins, frequently rounding to zero (Dettmers et al., 2022; Xiao et al., 2023). This phenomenon effectively extinguishes the model's learned information, severing gradient flow and decimating the network's predictive perplexity. To address model compression without the prohibitive computational costs of Quantization-Aware Training (QAT) - which requires full backpropagation and retraining on massive datasets, the field shifted entirely toward Post-Training Quantization (PTQ) paradigms (Frantar et al., 2023; Lin et al., 2024). Advanced PTQ theories focus on mathematically neutralizing the outlier phenomenon using either structural decomposition, mathematically equivalent transformations, or second-order error compensation.

3.3.3 Mixed-Precision Decomposition (LLM.int8())

To support effective INT8 quantization in the presence of extreme activation outliers, the LLM.int8() framework proposed a mixed-precision decomposition paradigm. The theoretical basis of this method relies on the observation that because outliers are restricted to specific, systematic feature dimensions, they can be mathematically isolated from the rest of the computation (Dettmers et al., 2022).

Given an input matrix X and a weight matrix W , mixed-precision decomposition separates the outlier feature dimensions into a distinct set O , which contains all dimensions that possess at least one value exceeding a defined magnitude threshold. The matrix multiplication is therefore bifurcated. The feature dimensions belonging to set O are extracted and multiplied in strict 16-bit floating-point precision, ensuring that the critical structural outliers suffer zero quantization degradation (Dettmers et al., 2022). Simultaneously, the remaining 99.9% of the regular dimensions are subjected to vector-wise quantization, where a separate quantization normalization constant is assigned to each inner product row and column and multiplied in 8-bit precision (Dettmers et al., 2022). The resulting INT32 outputs of the 8-bit multiplication are dequantized and mathematically accumulated with the 16-bit outlier outputs. This theoretical decomposition proves that preserving the exact magnitude of sparse outliers is sufficient to retain full 16-bit predictive accuracy at scales up to 175 billion parameters, though it inherently introduces computational scheduling complexities.

3.3.4 Equivalent Smoothing Transformations (SmoothQuant)

While decomposing matrices effectively isolates outliers, it complicates hardware execution. The SmoothQuant methodology provides an alternative theoretical framework by recognizing an asymmetry in transformer layers: while activations are incredibly difficult to quantize due to outliers, the static model weights are generally uniform, flat, and exceptionally easy to quantize (Xiao et al., 2023).

To execute fully quantized 8-bit weight and 8-bit activation (W8A8) matrix multiplications without mixed-precision overhead, SmoothQuant mathematically migrates the quantization difficulty from the dynamic activations to the static weights offline. This relies on the mathematical equivalence of linear layers. For a linear operation $Y = X \cdot W$, the computation can be rewritten by introducing a per-channel smoothing factor, a continuous vector s : $Y = (X \text{diag}(s)^{-1}) \cdot (\text{diag}(s)W) = \hat{X}\hat{W}$ (Xiao et al., 2023).

By dividing the activation channels by the smoothing factor s and multiplying the corresponding weights by the same factor in the reverse direction, the exact output Y is perfectly preserved. The objective of this transformation is to choose s such that the variance of the new activations \hat{X} is suppressed, thereby eliminating the outlier magnitude spikes. To evenly balance the quantization difficulty between the weights and activations, SmoothQuant introduces a migration strength hyperparameter α , formally defining the smoothing factor as: $s_j = \frac{\max(|X_j|)^\alpha}{\max(|W_j|)^{1-\alpha}}$ (Xiao et al., 2023). When α is tuned appropriately (typically around 0.5), the outlier channels in the activation matrix are flattened, while the corresponding channels in the weight matrix are expanded (Xiao et al., 2023). By absorbing the extreme variance into the static weights prior to runtime, both matrices become highly quantization-friendly, permitting uniform PTQ to succeed where naive methods fail.

3.3.5 Second-Order Weight Quantization (GPTQ)

As memory constraints become more severe at the network edge, quantization theories must compress networks beyond 8-bit, targeting 4-bit or 3-bit representations. At these extreme low-bit regimes, activation quantization is frequently abandoned in favor of weight-only quantization to preserve accuracy. The Generative Pre-trained Transformer Quantization (GPTQ) algorithm approaches weight-only compression as a layer-wise optimization problem, seeking to find a matrix of discrete quantized weights \hat{W} that

minimizes the squared reconstruction error relative to the full-precision output:

$$\operatorname{argmin}_{\widehat{W}} \left\| WX - \widehat{W}X \right\|_2^2 \text{ (Frantar et al., 2023).}$$

To solve this objective efficiently across billions of parameters, GPTQ utilizes approximate second-order information, generalizing the Optimal Brain Quantization (OBQ) framework. The algorithm quantizes weights sequentially while continuously updating the remaining unquantized weights to explicitly compensate for the mathematical error introduced by rounding the current weight. If a continuous weight w_q is rounded to a discrete grid point, the optimal update vector δ_F to apply to the remaining flat weights is determined by the inverse Hessian matrix H^{-1} of the layer's

$$\text{activation inputs: } \delta_F = -\frac{w_q - \text{quant}(w_q)}{[H^{-1}]_{qq}} \cdot (H^{-1})q, \text{ (Frantar et al., 2023).}$$

To scale this computationally expensive operation to massive transformer blocks, GPTQ processes weights in consecutive column blocks rather than individually. Furthermore, it applies a Cholesky decomposition to the inverse Hessian matrix upfront (Frantar et al., 2023). This mathematically stabilizes the updating process against numerical inaccuracies and indefinite matrices, allowing the algorithm to compress models containing hundreds of billions of parameters to 3 or 4 bits per weight with negligible accuracy degradation, fundamentally solving the high-density storage problem for inference.

3.3.6 Activation-Aware Weight Saliency (AWQ)

An alternative to second-order reconstruction is the Activation-aware Weight Quantization (AWQ) framework, which optimizes low-bit weight quantization without relying on backpropagation or heavy Hessian reconstructions, thereby preventing overfitting to the calibration data. The theoretical foundation of AWQ is the insight that not all weights in a language model contribute equally to the final perplexity; retaining the precision of merely 1% of the most salient weights can prevent almost all quantization degradation (Lin et al., 2024).

Crucially, AWQ demonstrates that weight saliency cannot be accurately identified by examining the magnitude or $L2$ -norm of the weights themselves. Instead, weight saliency is strictly determined by the activation distribution; weight channels that correspond to large activation magnitudes process the most critical temporal features and are therefore deemed salient (Lin et al., 2024).

However, explicitly keeping 1% of salient weights in high precision introduces the same hardware inefficiencies observed in mixed-precision models. To circumvent this, AWQ scales the salient weight channels to mathematically minimize their relative quantization error. Let the quantization function be $Q(w) = \Delta \cdot \text{Round}(w/\Delta)$. If a salient weight w is multiplied by a scaling factor $s > 1$, and the corresponding activation x is inversely scaled by $1/s$, the operation becomes $Q(w \cdot s)(x/s)$ (Lin et al., 2024). The quantization error for this scaled calculation is derived as: $\text{Err}(Q(w \cdot s)(x/s)) = \Delta' \cdot \text{RoundErr}\left(\frac{w \cdot s}{\Delta'}\right) \cdot x \cdot \frac{1}{s}$ (Lin et al., 2024).

Because the expected rounding error is fundamentally bounded by a uniform distribution, and multiplying a single channel by s rarely alters the overall maximum scalar Δ' of the entire weight group, the expression simplifies. The presence of the $1/s$ term dictates that as the salient weight is scaled up prior to quantization, its relative quantization error decreases inversely (Lin et al., 2024). AWQ defines an activation-aware search space, $s = s_x$, where s_x is the average magnitude of the activation channel. By optimizing α , AWQ identifies the precise mathematical scaling vector that minimizes the overall quantization loss across the layer. By embedding this scale into the static weights, AWQ achieves robust, generalization-preserving 4-bit weight compression entirely optimized for localized hardware execution.

3.3.7 Implications for Edge Intelligence

The theoretical principles of quantization, spanning mixed-precision decomposition, offline equivalent transformations, second-order error compensation, and activation-aware scaling; are paramount for the deployment of conversational artificial intelligence outside of centralized cloud infrastructures. In the context of edge computing, where devices are constrained by fixed processing power and strict thermal limitations, the memory bottleneck poses the most absolute barrier to local inference (Zhou et al., 2019). By compressing 16-bit architectures into stable 8-bit or 4-bit configurations, the required memory bandwidth for the autoregressive decoding stage is slashed proportionally, directly increasing the theoretical maximum arithmetic intensity and multiplying token generation throughput (Lin et al., 2024). Consequently, these rigorous post-training quantization algorithms represent the definitive mathematical methodology required to decouple natural language generation from the cloud, enabling the deployment of autonomous, privacy-preserving, and low-latency artificial intelligence directly at the network edge.

3.4 Speech Processing Theory

Human communication is intrinsically multimodal, with spoken language representing the most natural and efficient interface. To elevate a language model from a text-based processor to a fully realized conversational agent, it must be deeply integrated with acoustic modeling pipelines. Speech processing theory encompasses the mathematical and architectural frameworks required to translate continuous, noisy acoustic waveforms into discrete semantic tokens, and conversely, to synthesize expressive acoustic signals from generated text. This section rigorously examines the theoretical foundations of acoustic representations, automatic speech recognition (ASR), voice activity detection (VAD), and text-to-speech (TTS) synthesis.

3.4.1 Acoustic Signals and Speech Representation

Raw audio waveforms are high-dimensional, continuous temporal signals that contain vast amounts of information, much of which is irrelevant to semantic linguistic comprehension, such as ambient background noise, acoustic reverberation, and speaker-specific biological traits. Processing these raw waveforms directly is computationally prohibitive and mathematically unstable for most traditional learning algorithms (Hinton et al., 2012). Consequently, the foundational step in speech processing involves transforming the raw waveform into a lower-dimensional, discriminative feature representation.

Historically, this preprocessing was highly engineered. Early speech recognition systems relied heavily on Mel-frequency cepstral coefficients (MFCCs) and perceptual linear predictive (PLP) coefficients (Hinton et al., 2012). These representations were explicitly designed to decouple the highly correlated components of the acoustic signal. By transforming the signal into a domain where individual components are roughly statistically independent, MFCCs satisfied the strict mathematical assumptions required by Gaussian Mixture Models (GMMs) with diagonal covariance matrices, which were the bedrock of early acoustic modeling (Hinton et al., 2012).

However, with the advent of deep learning, the necessity for strictly uncorrelated features diminished. Modern neural network architectures are highly capable of modeling complex, correlated data representations. Consequently, contemporary speech processing predominantly utilizes log-magnitude Mel spectrograms (Radford et al., 2023). To compute a Mel spectrogram, the continuous audio signal is first digitized and segmented into short, overlapping temporal windows (e.g., 25-millisecond windows with a 10-millisecond stride). A Short-Time Fourier Transform (STFT) is applied to each window to map the time-domain signal into the frequency domain, capturing the magnitude of different frequency bins (Kim et al., 2021). Because human auditory perception is non-linear; specifically, humans are highly sensitive to variations in low frequencies but less sensitive to high frequencies, the linear frequency spectrum is

subsequently warped onto the Mel scale using a specialized filterbank (Radford et al., 2023). Taking the logarithm of these filtered magnitudes compresses the dynamic range of the audio, resulting in a dense, two-dimensional time-frequency matrix that serves as the standard input for deep neural acoustic models (Radford et al., 2023).

3.4.2 Automatic Speech Recognition (ASR) Theory

Automatic speech recognition is the theoretical problem of mapping an acoustic observation sequence to a sequence of discrete linguistic tokens. For decades, the dominant theoretical paradigm for ASR relied upon Hidden Markov Models (HMMs) combined with Gaussian Mixture Models (GMMs) (Hinton et al., 2012). In this generative framework, the HMMs dealt with the temporal variability of speech by modeling the sequence of phonemes as a Markov process, while the GMMs evaluated how well each hidden state of the HMM fit a specific frame of acoustic coefficients (Hinton et al., 2012).

The paradigm shifted profoundly with the application of Deep Neural Networks (DNNs) to acoustic modeling. Researchers demonstrated that a deep feed-forward neural network, taking a contextual window of several contiguous acoustic frames as input, could produce vastly superior posterior probability estimates over HMM states compared to traditional GMMs (Hinton et al., 2012). This hybrid DNN-HMM architecture benefited from the non-linear, hierarchical feature extraction capabilities of multiple hidden layers, establishing new state-of-the-art results across various large-vocabulary continuous speech recognition benchmarks (Hinton et al., 2012).

Despite the success of hybrid models, the complex, multi-component nature of DNN-HMM pipelines, which required separate acoustic models, pronunciation lexicons, and language models drove the evolution toward fully end-to-end Sequence-to-Sequence (Seq2Seq) architectures. Modern ASR systems utilize attention-based Transformer models to consume the entire sequence of Mel spectrogram features and autoregressively predict the corresponding text transcript directly (Radford et al., 2023).

A prominent contemporary example of this theoretical shift is the Whisper architecture. Whisper employs an encoder-decoder Transformer structure that processes 80-channel log-magnitude Mel spectrograms (Radford et al., 2023). The encoder maps the continuous acoustic features into a sequence of contextual hidden states, and the decoder applies cross-attention over these states to generate a sequence of byte-pair encoded text tokens.

Crucially, modern ASR theory has expanded beyond purely architectural innovations into the realm of data-driven robustness. Historically, neural ASR systems were trained on small, highly curated, gold-standard datasets, resulting in brittle models that suffered catastrophic performance degradation when exposed to out-of-distribution noise or diverse accents (Radford et al., 2023). The contemporary solution relies on large-scale weak supervision. By training encoder-decoder architectures on hundreds of thousands of hours of noisy, multilingual, web-scraped audio-transcript pairs, the neural network is forced to learn highly generalized, robust acoustic representations (Radford et al., 2023). This weakly supervised pre-training paradigm explicitly removes the need for dataset-specific fine-tuning, allowing the ASR model to achieve zero-shot transfer capabilities that approach human-level transcription accuracy across diverse, real-world acoustic environments (Radford et al., 2023).

3.4.3 Voice Activity Detection (VAD) Theory

In real-world conversational systems, continuous audio streaming captures vast amounts of ambient background noise. Feeding unstructured silence or environmental noise into an autoregressive ASR model frequently triggers artificial intelligence hallucinations, generating catastrophic insertion errors that destabilize downstream language processing (Ramírez et al., 2004). To enforce temporal segmentation and protect the processing pipeline, Voice Activity Detection (VAD) is required.

VAD is fundamentally a statistical hypothesis testing problem aimed at determining the presence or absence of human speech in a given signal. Let \mathbf{x} represent a feature vector extracted from an audio frame. Assuming that speech and noise are additive, the VAD

module must decide between two hypotheses: H_0 (the frame contains only noise) and H_1 (the frame contains speech and noise) (Ramírez et al., 2004). The optimal decision rule that minimizes the error probability is the Bayes classifier, which selects the class with the highest posterior probability $P(H_i|\mathbf{x})$. Applying Bayes' rule yields the statistical likelihood ratio test (LRT) (Ramírez et al., 2004).

If the Discrete Fourier Transform (DFT) coefficients of the clean speech and the noise are modeled as asymptotically independent Gaussian random variables, the decision rule can be mathematically formulated in terms of a priori and a posteriori Signal-to-Noise Ratios (SNRs) (Ramírez et al., 2004). An observation vector is classified as speech if the calculated likelihood ratio exceeds an empirically determined threshold η .

While instantaneous spectral analysis is effective in high-SNR conditions, it rapidly loses discriminative power when background noise increases. Advanced VAD algorithms mitigate this by incorporating long-term speech information (Ramírez et al., 2004). Because the most significant information for detecting voice activity remains in the temporal variations of the spectral envelope, tracking a multi-frame Long-Term Spectral Envelope (LTSE) drastically improves robustness against non-stationary noise. The decision rule is subsequently defined using the Long-Term Spectral Divergence (LTSD), which measures the deviation of the LTSE with respect to the average noise spectrum magnitude (Ramírez et al., 2004).

By applying these mathematical decision rules, VAD algorithms dynamically segment continuous audio streams, providing accurate endpoint detection. This active gating mechanism explicitly drops non-speech frames from the processing pipeline, conserving computational resources and preventing the ASR model from misinterpreting ambient noise as linguistic input (Ramírez et al., 2004).

3.4.4 Text-to-Speech (TTS) Synthesis Theory

To complete the conversational loop, an artificial intelligence system must synthesize natural-sounding acoustic waveforms from the text generated by the language model.

Traditional TTS systems relied on complex, two-stage pipelines: a front-end model converted preprocessed text into intermediate speech representations, such as linguistic features or Mel spectrograms, and a back-end vocoder subsequently synthesized raw waveforms conditioned on those intermediate representations (Kim et al., 2021). While capable of generating intelligible speech, these decoupled stages suffered from sequential training inefficiencies and error propagation.

Modern TTS theory has advanced toward parallel, end-to-end architectures that leverage deep generative modeling to map text directly to waveforms. A highly successful theoretical framework for this is the Conditional Variational Autoencoder (CVAE) augmented with normalizing flows and adversarial learning (Kim et al., 2021).

In a VAE-based TTS framework, the complex, one-to-many relationship of human speech must be mathematically addressed. A single text sequence can be spoken in infinite variations, with different pitches, intonations, and speaking rates. To model this uncertainty, the CVAE introduces latent variables z . During training, a posterior encoder takes the high-resolution linear spectrogram of the target speech as input and maps it to a latent distribution $q_{\phi}(z|x_{lin})$. Simultaneously, a prior encoder processes the input text phonemes c_{text} and an alignment A to produce a prior distribution $p_{\theta}(z|c_{text}, A)$ (Kim et al., 2021). The model is optimized by maximizing the Evidence Lower Bound (ELBO), which includes a reconstruction loss (ensuring the decoder accurately synthesizes the audio from the latent variables) and a Kullback-Leibler (KL) divergence loss (which minimizes the mathematical distance between the approximate posterior distribution and the text-conditioned prior) (Kim et al., 2021).

To increase the expressive power and flexibility of the prior distribution beyond a simple Gaussian, normalizing flows are applied to the prior encoder. Normalizing flows are sequences of invertible, non-linear transformations that can map simple distributions into highly complex, multimodal distributions perfectly suited for modeling the diverse nuances of human prosody (Kim et al., 2021).

A critical theoretical challenge in end-to-end TTS is estimating the monotonic alignment A between the discrete text phonemes and the continuous acoustic frames. Because ground truth alignment labels do not exist, systems utilize Monotonic Alignment Search (MAS), a dynamic programming algorithm designed to find the specific alignment that maximizes the log-likelihood of the latent variables z parameterized by the normalizing flow (Kim et al., 2021).

Once aligned, the system must determine how long each phoneme should be sounded out. Rather than predicting a rigid, deterministic duration, state-of-the-art TTS architectures employ a stochastic duration predictor (Kim et al., 2021). Using variational dequantization and variational data augmentation, the stochastic duration predictor models the distribution of phoneme lengths as a continuous variable. During inference, phoneme durations are randomly sampled from this distribution via the inverse transformation of the flow, allowing the model to generate highly diverse, human-like speech rhythms from the exact same text input (Kim et al., 2021).

Finally, to synthesize high-fidelity raw audio waveforms that bypass the robotic artifacts typical of older vocoders, the CVAE incorporates Generative Adversarial Networks (GANs). A sophisticated decoder acts as the generator, while a multi-period discriminator evaluates periodic patterns in the output to distinguish between generated waveforms and ground truth human speech (Kim et al., 2021). By combining least-squares adversarial loss with feature-matching loss, the generator is mathematically forced to produce waveforms that are perceptually indistinguishable from real human speech. This integration of variational inference, stochastic duration prediction, and adversarial learning forms the theoretical apex of contemporary Text-to-Speech synthesis (Kim et al., 2021).

3.4.5 Challenges of Low-Latency Voice Interaction

The theoretical mechanics of ASR, VAD, and TTS collectively establish the operational loop of a voice-based conversational agent. However, combining these complex neural models introduces profound latency challenges. Sequential dependencies in traditional

pipelines, where the ASR must transcribe an entire audio sequence before the language model can begin processing, and the language model must output a full sentence before TTS synthesis begins; create severe communication delays. Furthermore, large sequence-to-sequence ASR models require sufficient audio context to perform global attention, making continuous, real-time streaming mathematically difficult without significant architectural trade-offs (Radford et al., 2023). Addressing these latency barriers necessitates the exploration of highly concurrent execution strategies and decoupled processing frameworks.

3.5 Concurrent Systems Theory

The theoretical evolution of conversational artificial intelligence demonstrates that integrating advanced neural networks, specifically automatic speech recognition, large language models, and text-to-speech synthesis is necessary to facilitate natural human-machine interaction. However, hosting these computationally intensive models on edge environments introduces profound temporal challenges. Human cognitive and perceptual systems are acutely sensitive to interaction delays; interactive communication demands that end-to-end latency be tightly controlled, often to within a few tens to hundreds of milliseconds, to maintain a natural, overlapping conversational flow (Satyanarayanan, 2017). When multiple complex neural networks are deployed in a localized environment, achieving this strict real-time responsiveness requires more than optimizing the individual models; it necessitates the theoretical restructuring of the software execution paradigm. This section explores the computer science foundations of concurrent systems, asynchronous processing, and backpressure theory, establishing the theoretical framework required to overcome the latency bottlenecks inherent in sequential conversational pipelines.

3.5.1 Sequential Execution Models and Latency Propagation

To understand the necessity of concurrency, the mathematical and structural limitations of sequential execution models must first be analyzed. In a strictly sequential (or synchronous) processing architecture, a predefined set of computational tasks must be

executed in a rigid, linear order. For a given sequence of operations T_1, T_2, \dots, T_n , the fundamental constraint of sequential computation dictates that operation T_{i+1} cannot initiate until operation T_i has completely resolved (Vaswani et al., 2017).

In the context of a voice conversational agent, this sequential dependency creates a monolithic control flow. The system must record a complete audio utterance (T_1), wait for the transcription model to process the audio (T_2), wait for the language model to generate a complete text response (T_3), and finally wait for the acoustic synthesizer to render the audio waveform (T_4). From a theoretical standpoint, the total end-to-end processing latency L_{total} of a sequential system is strictly additive, defined as the sum of the individual latencies of each component: $L_{total} = \sum_{i=1}^n L_i$

This additive latency propagation presents a massive theoretical barrier to real-time responsiveness. Furthermore, the individual stages within an artificial intelligence pipeline possess vastly different computational complexities and arithmetic intensities. For example, the autoregressive generation stage of a large language model is heavily memory-bound and computationally intensive compared to the relatively lightweight processing required for voice activity detection or audio ingestion (Lin et al., 2024). In a sequential architecture, these asymmetric workloads result in highly inefficient resource utilization. The central processing unit (CPU) or graphics processing unit (GPU) is frequently forced into an idle, blocking state, waiting for memory-bound tasks to complete before the next stage can begin. Because sequential processing strictly precludes parallelization within a given interaction cycle (Vaswani et al., 2017), the overall throughput of the system is artificially bottlenecked by the summation of all operational delays, rendering it mathematically incapable of supporting rapid conversational turn-taking.

3.5.2 Asynchronous Processing and Decoupled Pipelines

To neutralize the severe latency accumulation of sequential models, concurrent systems theory proposes the transition toward asynchronous processing and decoupled pipelines. Concurrency is a theoretical software execution paradigm wherein multiple

computational tasks make progress overlapping in time, independent of a strictly defined sequential order. Unlike true parallelism, which requires multiple physical processing cores to execute instructions simultaneously at the exact same physical instant concurrency is a property of the software architecture that deals with the management and scheduling of independent execution streams.

By decoupling a monolithic conversational pipeline into independent, concurrent processes, the strict mathematical dependency $L_{total} = \sum L_i$ is broken. In a concurrent architecture, the audio ingestion, natural language understanding, and acoustic synthesis stages are isolated into distinct computational boundaries that execute asynchronously. When a system extracts as much concurrency as possible from its workloads, it drastically improves resource utilization and minimizes waiting times (Zhou et al., 2019).

In a decoupled pipeline, the latency is no longer the sum of all tasks, but is instead constrained primarily by the processing time of the single slowest stage, the system bottleneck; plus the minimal propagation delay between stages. Furthermore, asynchronous processing enables temporal overlapping. The system can capture a new audio utterance at the exact same time it is synthesizing the speech for the previous response. By establishing a software-defined computing stream where data is passed between distributed, independent functions (Shi et al., 2016), the conversational agent transitions from a stop-and-wait operation to a continuous, fluid processing loop, effectively mirroring the overlapping nature of human dialogue.

3.5.3 The Producer-Consumer Paradigm and Buffering

The theoretical foundation for managing data flow between decoupled, asynchronous processes is the producer-consumer paradigm. In this classical computer science model, a "producer" process generates data entities and places them into a shared buffer, while a "consumer" process removes the data entities from the buffer and processes them. Within a modular conversational artificial intelligence system, this paradigm exists at multiple intersections: the microphone capture stream produces raw audio that the

transcription module consumes; the transcription module produces text that the language model consumes; and the language model produces sentences that the speech synthesis module consumes.

Because these independent processes execute asynchronously and possess inherently different, often stochastic, execution times, they rarely operate at the exact same speed. An acoustic sensor may capture audio data at a highly consistent, deterministic rate, whereas a generative language model processes data at a variable rate depending on the complexity of the input prompt and the resulting generation length. If a producer and consumer were tightly coupled without an intermediary structure, the faster process would constantly be forced to halt and wait for the slower process, reverting the system back to the inefficient sequential timing model.

To preserve the independence and efficiency of both processes, the architecture utilizes a buffer, a designated region of memory structured as a First-In-First-Out (FIFO) queue. The buffer mathematically absorbs the variance in processing speeds. When the producer is faster, the generated data accumulates in the buffer, allowing the producer to continue executing without blocking. When the consumer is ready, it retrieves the data from the buffer at its own maximum operational speed. This buffering mechanism is the definitive structural requirement for enabling asynchronous communication between neural network modules operating under different computational complexities.

3.5.4 Backpressure Theory in Low-Latency Environments

While buffering resolves the immediate timing discrepancies between asynchronous processes, it introduces a severe theoretical vulnerability regarding memory stability and real-time responsiveness if the queue capacities are left unregulated. If an asynchronous queue is theoretically unbounded (infinite capacity), a producer that consistently operates faster than a consumer will continuously deposit data into the buffer. Over time, the buffer will grow infinitely, inevitably leading to catastrophic memory exhaustion, a primary cause of system failure on constrained edge devices.

More critically for conversational systems, unbounded queues destroy real-time responsiveness through latency bloat. If a user speaks continuously, generating numerous audio segments faster than the language model can evaluate them, an unbounded queue will store all historical audio. Consequently, the consumer will eventually process data that is seconds or minutes old. This "stale" data processing completely violates the strict end-to-end temporal requirements of interactive human communication (Satyanarayanan, 2017).

To mathematically guarantee system stability and preserve low-latency execution, concurrent architectures must implement bounded queues, which enforce a strict maximum capacity on the buffer. The introduction of bounded buffering creates the theoretical mechanism of backpressure. Backpressure is a regulatory feedback signal that propagates from the consumer backward to the producer.

When a bounded queue reaches its maximum capacity, the producer is theoretically prohibited from adding new data. The producer must either block (suspend its execution) until the consumer removes an item, or it must drop the new data entirely. In the context of a real-time edge voice chatbot, dropping the newest data or blocking the audio stream enforces a strict temporal alignment. By limiting the queue size, often to a theoretical maximum of exactly one element; the architecture mathematically guarantees that the language and acoustic models only ever process the single most recent, up-to-date utterance. Backpressure neutralizes the accumulation of stale data, bounding the maximum possible latency of the pipeline to the processing time of the current active queue elements, thereby ensuring the system remains anchored in the present conversational moment.

3.5.5 Synchronization, Race Conditions, and Liveness Properties

While the transition from sequential to concurrent execution provides massive theoretical reductions in latency, it introduces profound complexities regarding system state management and data integrity. In a concurrent pipeline, multiple independent processes operate simultaneously within the same overarching application, frequently

requiring access to shared computational resources, such as memory addresses, hardware sensors, or neural network weights.

Because the execution scheduling of asynchronous processes is fundamentally non-deterministic governed by the operating system rather than a fixed program sequence; the exact timing of resource access cannot be predicted. If two or more concurrent processes attempt to read and modify a shared resource simultaneously without regulation, the system encounters a race condition. In a race condition, the final state of the data becomes dependent on the arbitrary timing of the execution threads, leading to unpredictable, mathematically invalid states and catastrophic software crashes.

To prevent race conditions, concurrent systems theory dictates the use of synchronization primitives. Synchronization is the theoretical mechanism of restricting access to shared resources to enforce mutual exclusion. When a process requires access to a shared variable or a critical section of the pipeline, it must acquire a mathematical lock. While the lock is held, all other concurrent processes attempting to access that specific resource are suspended in a blocked state until the lock is released. As edge computing relies on complex, distributed computing streams, addressing these collaboration and synchronization issues across multiple operational layers is a mandatory requirement for system stability (Shi et al., 2016).

However, the application of synchronization introduces theoretical risks to the system's liveness properties; the guarantee that the software will continue to make computational progress. The most severe liveness failure is deadlock, a state wherein two or more concurrent processes are indefinitely blocked, each waiting for a resource that the other process currently holds. Similarly, starvation can occur if a specific process is perpetually denied access to a necessary resource due to the scheduling priorities of other threads. Consequently, engineering a low-latency concurrent pipeline requires meticulous architectural design to ensure that locks are acquired and released in a mathematically provable, consistent order, preventing deadlocks while minimizing the time resources remain locked.

3.5.6 Complexity Trade-offs in Edge Architectures

The application of concurrent systems theory to localized artificial intelligence involves a deliberate trade-off between architectural complexity and execution efficiency. Sequential systems are theoretically trivial to implement, inherently safe from race conditions, and mathematically deterministic. However, their inability to parallelize independent tasks and their strict additive latency make them fundamentally inadequate for natural voice conversational interfaces.

Conversely, concurrent and asynchronous architectures are highly non-deterministic and require sophisticated mechanisms for state migration, backpressure regulation, and resource synchronization (Shi et al., 2016). They introduce computational overhead through context switching, the processing time required by the hardware to save the state of one task and load the state of another, and lock contention. Nevertheless, in the domain of edge computing and artificial intelligence, where heavy neural networks process asymmetric workloads, the latency reductions achieved by decoupling the pipeline vastly outweigh the theoretical overhead of synchronization.

By applying the principles of producer-consumer buffering, bounded backpressure, and asynchronous scheduling, the fundamental computational bottlenecks of the system are isolated and neutralized. This concurrent theoretical framework represents the definitive software engineering requirement for translating complex, heavy artificial intelligence models into highly responsive, interactive, and autonomous conversational agents operating natively at the network edge.

4 Implementation and Results

This chapter presents the practical implementation and evaluation of the proposed edge-based voice chatbot architecture. Following the theoretical foundations established in the preceding chapters, the focus shifts to the engineering decisions, deployment constraints, architectural evolution, and empirical observations obtained during system development. The chapter first describes the hardware and software environment used for implementation, followed by the iterative progression from a sequential baseline architecture toward a concurrent multithreaded design. Finally, the chapter evaluates the resulting system in terms of memory utilization, conversational latency, and overall operational feasibility within a consumer-grade edge computing environment.

4.1 Hardware Environment and Constraints

The empirical development and evaluation of the fully localized conversational artificial intelligence pipeline were conducted within a strictly defined hardware environment intended to represent standard, consumer-grade computing capabilities. The primary deployment target was a computing system equipped with an Nvidia RTX 3060 graphical processing unit (GPU). This specific hardware was selected because it imposes a strict physical limitation of six gigabytes of Video Random Access Memory (VRAM). The underlying operating system environment was configured to support cross-platform execution, explicitly accommodating both Windows 10/11 and Linux (Ubuntu) distributions. The runtime environment was established using Python version 3.12.9, supported by PyTorch version 2.8.0 and Compute Unified Device Architecture (CUDA) versions 12.1 and 12.2, to ensure that the system could leverage full hardware acceleration for deep learning inference.

In the context of local conversational artificial intelligence, stringent memory limitations represent the most critical deployment constraint. Traditional cloud-based voice assistants mitigate computational bottlenecks by offloading processing to massive

external server farms equipped with virtually unlimited memory bandwidth and storage. In stark contrast, deploying a responsive voice chatbot entirely at the network edge requires that the automatic speech recognition, large language model, and text-to-speech synthesis components all reside simultaneously within the limited physical memory of a single local device. If the combined memory footprint of these neural networks exceeds the available VRAM capacity, the system experiences severe Out-of-Memory crashes, entirely halting execution and rendering autonomous conversation impossible. Therefore, the 6 GB VRAM constraint served as the absolute upper boundary for the entire project, directly dictating the selection and configuration of all subsequent machine learning models.

To successfully execute within this hardware boundary, highly specific model selections were mandated. For the audio transcription component, the system utilized the 'base' configuration of the Whisper Speech-to-Text model. While larger variants of the Whisper architecture provide higher transcription accuracy for complex audio inputs, they require significantly more memory. The 'base' model was strategically selected because it consumes approximately 145 megabytes of VRAM, prioritizing a low memory footprint and rapid processing speed over absolute transcription precision.

The Qwen 2.5 1.5B-Instruct model was selected as the language generation component and contributed approximately 3.0 GB of VRAM consumption during inference. The detailed rationale behind the model selection and deployment configuration is discussed in Section 4.2.2.

For the final stage of the conversational loop, the architecture utilized the Coqui Text-to-Speech synthesis engine. This specific module was integrated to provide high-quality, localized voice generation and requires approximately 150 megabytes of VRAM during execution.

When combined, these strategic model selections yield a total, concurrent steady-state VRAM footprint of approximately 3.3 gigabytes. Achieving this specific memory footprint is highly significant relative to the available 6.0 GB budget, as it consumes roughly 55% of the total hardware capacity. Maintaining a steady-state footprint of 3.3 gigabytes deliberately leaves a 45% safety reserve within the GPU memory. This safety margin is a critical engineering requirement for edge deployment. During steady-state execution, the language model requires dynamic memory allocation to manage expanding context windows and conversational histories, which temporarily inflate VRAM usage during generation. Furthermore, the local operating system requires dedicated memory to perform background tasks and render graphical user interfaces. Without this 45% reserve, unexpected memory spikes during concurrent processing could exceed the hardware limits and trigger fatal memory failures.

The rigid 6 GB memory constraint not only dictated the selection of the primary generative models but also profoundly influenced the engineering of auxiliary system components. Because the core conversational models consumed over half of the available GPU memory, the architecture could not afford to allocate additional VRAM to secondary, non-generative tasks. For example, implementing continuous background noise filtering or acoustic monitoring natively on the GPU would have risked crossing the 6 GB threshold. To circumvent this, the system's Voice Activity Detection module, a lightweight neural network of approximately 1.5 megabytes, was explicitly configured to execute entirely on the central processing unit. This deliberate offloading strategy ensured that the auxiliary monitoring tasks did not compete with the primary neural networks for memory, preserving the limited GPU VRAM exclusively for the intensive transcription and language generation operations.

Ultimately, validating that the integrated speech and language models could operate safely within the 3.3 gigabyte footprint established the foundational viability of the system. By satisfying the strict physical limitations of consumer-grade hardware, the research proved that the required machine learning assets could be hosted locally. With

the hardware environment stabilized and the memory constraints successfully navigated, the research subsequently focused on configuring the specific software parameters and synchronization mechanisms required to orchestrate these models into a fluid, conversational pipeline. This software implementation is detailed in the following section.

4.2 Software Implementation and Configuration

The fully localized edge voice chatbot relies on a sophisticated, integrated software architecture designed to replicate the capabilities of cloud-based conversational artificial intelligence without transmitting sensitive data to external servers. To accomplish this, the system is constructed from four primary software components, each responsible for a distinct stage of the multimodal conversational pipeline. These components include an automatic speech recognition module to transcribe incoming audio, a large language model to perform natural language understanding and response generation, a text-to-speech synthesis engine to render acoustic waveforms, and a voice activity detection module to intelligently monitor and segment the audio stream. The selection, configuration, and integration of these specific software libraries were strictly dictated by the necessity to operate entirely on-device, prioritize cross-platform compatibility across Windows and Linux operating systems, and conform to the rigid physical memory constraints of consumer-grade graphical processing units.

4.2.1 Speech Recognition Component

The initial stage of the conversational artificial intelligence pipeline is the automatic speech recognition component, which is responsible for ingesting raw acoustic data and converting it into discrete text tokens that the language model can process. To fulfill this role, the architecture integrates the Whisper speech-to-text framework developed by OpenAI. Whisper was selected because it provides highly accurate, offline transcription capabilities that have been proven to generalize robustly across diverse acoustic environments without requiring domain-specific fine-tuning or cloud connectivity.

Within the Whisper model family, the specific configuration selected for the edge deployment was the 'base' model. The engineering rationale for selecting the 'base' variant over the 'tiny', 'small', or 'medium' architectures was fundamentally driven by the need to balance transcription precision with the strict computational limits of the target hardware. The 'base' model operates with a highly efficient memory footprint, consuming approximately 145 megabytes of Video Random Access Memory (VRAM) when loaded onto the graphical processing unit. This low memory allocation was critical to preserving the majority of the 6-gigabyte VRAM budget for the computationally dominant natural language generation engine. While deployment evidence indicates that switching to larger variants, such as the 'small' model (which requires an additional 320 megabytes of VRAM) or the 'medium' model (which requires an additional 1.5 gigabytes of VRAM), yields noticeable improvements in absolute transcription accuracy, scaling up the acoustic model introduces unacceptable increases in processing latency and memory consumption that threaten the stability of the entire pipeline. Therefore, the 'base' model was established as the optimal baseline for achieving rapid, localized transcription within the defined hardware constraints.

In terms of software integration, the speech recognition component utilizes the official Python implementation of Whisper, supported by the PyTorch deep learning framework. To accurately process user input without introducing artificial padding or audio truncation errors, the component invokes the specific transcription methods native to the Whisper library, which are explicitly designed to handle variable-length audio arrays. Operating directly on the graphical processing unit, this component continuously receives floating-point audio arrays captured from the user's microphone, processes the log-magnitude Mel spectrograms, and outputs the decoded text strings. By performing this transcription entirely locally, the component eliminates the variable network propagation delays inherent in traditional cloud-based speech-to-text application programming interfaces, forming the foundational input layer of the offline conversational loop.

4.2.2 Language Model Component

The core reasoning and natural language generation engine of the localized chatbot is powered by the Qwen 2.5 1.5B-Instruct large language model. Serving as the central node of the pipeline, this component consumes the transcribed text provided by the speech recognition module, evaluates the conversational context, and autoregressively generates the text of the chatbot's response. Because language modeling represents the most computationally intensive phase of the system, the selection and deployment strategy for this specific component required meticulous engineering to navigate severe hardware and software compatibility limitations.

The integration of the language model was achieved utilizing the Hugging Face Transformers library, relying on the native PyTorch framework to load the model tensors directly into the graphical processing unit's memory. The model was initialized using automatic device mapping algorithms to optimize tensor placement and was explicitly configured to operate in 16-bit floating-point (fp16) precision. By operating in fp16 precision, the 1.5 billion parameter model consumes approximately 3.0 gigabytes of VRAM during steady-state execution. This footprint represents the single largest memory allocation within the entire software architecture, occupying half of the available 6-gigabyte hardware budget.

During development, larger quantized model variants, including Qwen 2.5 3B-Instruct-AWQ, were evaluated as potential alternatives. However, deployment of AWQ models required software dependencies based on the Triton compiler framework, which is primarily supported on Linux environments and lacks straightforward compatibility with Windows systems. Because a key objective of the project was to maintain cross-platform portability, the quantized 3 billion parameter configuration was not adopted. Consequently, the standard Qwen 2.5 1.5B-Instruct model was selected as a practical compromise between computational efficiency, memory consumption, and deployment flexibility.

To optimize the selected language model for interactive spoken dialogue, strict generation parameters were embedded into the software configuration. The system utilizes specific sampling techniques, applying a temperature of 0.3 and a top-p threshold of 0.7 to ensure the generated text remains coherent and focused. Furthermore, a strict maximum generation limit of 150 new tokens is enforced, accompanied by a targeted system prompt that explicitly instructs the model to answer conversationally in one to two short sentences. This parameterization is essential for edge deployment, as it mathematically prevents the language model from generating excessively verbose responses that would needlessly inflate processing latency, delay the subsequent acoustic synthesis, and consume excess dynamic memory.

4.2.3 Text-to-Speech Component

The final generative phase of the conversational pipeline relies on the text-to-speech component, which transforms the text responses generated by the language model back into audible, naturalistic human speech. To accomplish this, the architecture integrates the Coqui Text-to-Speech (TTS) engine, utilizing the specific version 0.27.2 library to ensure strict compatibility with the Python 3.12 runtime environment.

Within the Coqui framework, the system is configured to deploy a multilingual, multi-dataset model architecture capable of high-fidelity voice cloning and synthesis. The primary role of this component is to provide the chatbot with an expressive, localized voice that mimics the cadence and intonation of a natural human speaker. Upon receiving a generated text string from the language model, the text-to-speech component processes the linguistic features, synthesizes the raw audio waveform, and subsequently interfaces with the local machine's audio hardware to play the resulting speech through the connected speakers. To ensure the generated output is parsed correctly, the software implementation includes a specialized preprocessing function that intercepts the text from the language model and expands numerical digits into their corresponding written English words, preventing the synthesis engine from silently dropping characters or generating decoding warnings during playback.

From a resource perspective, the Coqui TTS engine is highly optimized for localized execution. The synthesis model is loaded directly into the graphical processing unit alongside the speech recognition and language models, enabling accelerated acoustic rendering. Despite generating high-quality audio, the specific configuration utilized in this architecture demands a remarkably low memory footprint, consuming approximately 150 megabytes of VRAM. This minimal resource requirement is highly advantageous, as it allows the final output stage of the chatbot to operate efficiently without threatening the memory limitations established by the heavier language model. By executing the acoustic synthesis entirely on-device, this component ensures that the conversational loop is completed without the need to transmit text payloads to external cloud-based synthesis application programming interfaces, thereby preserving user privacy.

4.2.4 Voice Activity Detection Component

To govern the ingestion of raw audio and protect the generative neural networks from processing continuous ambient noise, the architecture relies on a dedicated voice activity detection component. This component serves as the intelligent gatekeeper for the entire conversational pipeline, actively monitoring the incoming audio stream from the microphone to identify the precise boundaries of human speech.

The software implementation utilizes the Silero Voice Activity Detection (VAD) library, a robust, pre-trained neural network specifically engineered for rapid speech segmentation. Unlike the speech recognition, language generation, and text-to-speech models, which all require heavy hardware acceleration on the graphical processing unit, the Silero VAD module is deployed using a strict central processing unit (CPU) execution strategy. The model itself is extremely lightweight, occupying approximately 1.5 megabytes of storage. By explicitly confining the execution of this monitoring neural network to the system's central processor, the architecture deliberately isolates the continuous audio scanning workload from the graphical processing unit. This strategic engineering decision guarantees that the persistent background task of listening for

human speech never competes for the strictly limited 6-gigabyte VRAM budget, effectively reserving the maximum possible graphical memory for the computationally expensive transcription and language modeling phases.

The operational role of the voice activity detection component is critical for maintaining conversational control and system stability. The software is configured to sample the live microphone input in microscopic temporal segments, evaluating the statistical probability that each segment contains human speech. It operates using strict, configurable thresholds to filter out transient background noises, such as heating and ventilation systems or computer fan interference, ensuring that only sustained vocalization is passed forward.

Most importantly, the component is responsible for detecting the cessation of speech. By continuously tracking the duration of absolute silence following a vocal input, the software can autonomously determine when a user has finished their conversational turn. Upon detecting a sustained period of silence, the voice activity detection component dynamically halts the audio capture and immediately releases the compiled audio data to the subsequent speech recognition module. This intelligent segmentation ensures that the core artificial intelligence models process only validated, variable-length human utterances, thereby preventing the ingestion of unstructured silence and reducing transcription errors, unnecessary processing loops, and inefficient resource utilization that may contribute to hallucinations and memory overload issues. Ultimately, this lightweight, CPU-bound software component establishes the foundational timing mechanics that allow the localized pipeline to mirror the natural, adaptive flow of human conversation.

Collectively, these software components established the functional foundation of the localized conversational artificial intelligence system. However, achieving stable operation within the hardware limitations discussed previously required several architectural revisions throughout the development process. The following section

therefore examines the evolution of the system architecture, beginning with an initial sequential baseline implementation and progressing toward the final adaptive and concurrent pipeline design.

4.3 Architectural Evolution

The successful deployment of a fully localized conversational artificial intelligence pipeline on consumer-grade hardware was not achieved through a single, static software iteration. Instead, the software architecture evolved systematically through multiple distinct development phases. Each developmental phase was specifically engineered to identify, isolate, and resolve the structural bottlenecks, hardware limitations, and conversational shortcomings exposed by the preceding design. The overall progression moved from a simple sequential processing pipeline toward an adaptive, low-latency concurrent architecture capable of facilitating more natural human-machine interaction. The purpose of this section is to document the engineering evolution of the system, providing an academic rationale for the architectural decisions, trade-offs, and refinements implemented throughout the development process.

4.3.1 Phase 1 - Sequential Baseline

The first functional prototype of the edge voice chatbot was implemented as a monolithic, single-threaded application. This initial Phase 1 architecture established a foundational sequential baseline operating under a strict stop-and-wait execution model. In this configuration, the entire conversational pipeline was executed sequentially within a single primary application thread, meaning that each computational stage had to complete before the next stage could begin.

The execution flow of the baseline architecture followed a strictly linear processing sequence. The conversational cycle commenced with the audio capture module, which was configured to record audio using a fixed-duration recording window. The system activated the microphone and captured acoustic input for exactly five seconds before

any subsequent processing could occur. During this recording period, the main application thread remained blocked while waiting for the recording buffer to fill. Once the recording interval had elapsed, the captured audio data was transferred to the speech recognition module, where the Whisper transcription model converted the acoustic signal into textual form.

Following transcription, the recognized text was appended to the conversational history and forwarded to the Qwen large language model. The language model processed the user input and conversational context to generate a textual response. This generated response was then passed to the Coqui text-to-speech engine, which synthesized the corresponding speech waveform and played it through the system's audio output. Only after audio playback had completely finished did the application return to the beginning of the cycle and initiate a new five-second recording window.

This highly structured sequential design was intentionally selected as the initial implementation strategy in order to validate the fundamental technological components of the project. Deploying speech recognition, language generation, and speech synthesis models simultaneously on a consumer-grade graphics processing unit introduced significant uncertainty regarding memory consumption and system stability. A single-threaded execution model provided a predictable and observable environment in which the interoperability of the selected software components could be verified. Furthermore, by eliminating the complexities associated with asynchronous execution, thread synchronization, race conditions, and shared resources, the sequential baseline substantially simplified debugging and functional validation during the early stages of development.

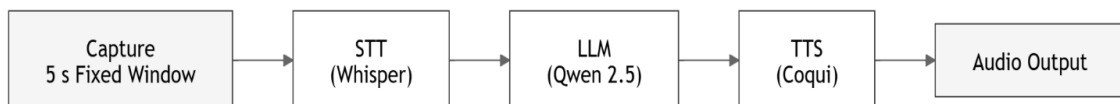


Figure 1. The Phase-1 sequential processing pipeline. The stop-and-wait architecture creates a bottleneck because each stage must be completed before the next begins.

While the Phase 1 architecture successfully demonstrated the feasibility of integrating the selected machine learning models into a single localized pipeline, practical testing quickly exposed several significant limitations. The most immediate drawback of the design was the accumulation of processing latency. Because each stage was required to wait for the previous stage to finish, the total response time became the sum of the recording duration, transcription time, language generation time, and speech synthesis time. This resulted in substantial delays between user input and system response, creating an interaction style that felt noticeably slower than natural human conversation.

The single-threaded design also prevented any form of conversational interruption. Once the system transitioned from audio capture to transcription, language generation, or speech synthesis, the microphone was no longer actively monitored. Consequently, users could not interrupt the chatbot while it was processing information or speaking. Any attempt to interject, correct a statement, or redirect the conversation was ignored until the entire processing cycle had completed. This imposed a rigid turn-taking structure that reduced conversational flexibility and negatively affected the perceived naturalness of the interaction.

Another significant limitation originated from the fixed-duration audio capture strategy. The recording mechanism operated independently of the actual duration of the user's speech. If a user issued a brief command, the system continued recording until the full five-second interval had elapsed, capturing unnecessary silence and environmental noise. Conversely, longer utterances could be truncated if they exceeded the predetermined recording window. As a result, the recording strategy frequently failed to align with natural speech behavior.

This inefficient handling of silence introduced significant reliability issues within the conversational pipeline. Because the entire recording buffer was forwarded to the speech recognition component regardless of speech content, the Whisper model was

frequently required to process extended periods of silence or ambient background noise. In some cases, these conditions produced inaccurate transcriptions consisting of irrelevant words, characters, or nonsensical text fragments.

These transcription errors were subsequently forwarded to the language model as valid conversational input. As the accumulated conversational context became increasingly polluted with irrelevant or nonsensical content, the overall stability and responsiveness of the system deteriorated. The observations obtained during testing demonstrated that the fixed recording strategy was unsuitable for practical conversational interaction and highlighted the need for more intelligent speech boundary detection.

Collectively, these limitations demonstrated that merely deploying conversational artificial intelligence models on local hardware was insufficient to achieve a natural user experience. The combination of accumulated latency, rigid turn-taking behavior, and inefficient audio capture mechanisms revealed the need for a more adaptive architecture. Addressing these shortcomings directly motivated the development of the second architectural phase, which introduced adaptive speech detection mechanisms to improve conversational responsiveness and recording efficiency.

4.3.2 Phase 2 - Adaptive Listening

The empirical observations gathered during the evaluation of the Phase 1 baseline architecture definitively established that a fixed-duration recording strategy was fundamentally unsuitable for interactive conversational applications. The utilization of a rigid, hardcoded five-second listening loop operated completely independently of the user's actual speech patterns. This static execution forced the system to ingest prolonged segments of ambient room noise whenever a user provided a brief utterance, and arbitrarily truncated sentences if the user spoke beyond the allotted timeframe. Most critically, forcing the speech recognition model to process extended periods of silence frequently resulted in inaccurate transcriptions and reduced overall system reliability. To salvage system stability and establish reliable speech boundaries, a fundamental

redesign of the audio capture module was strictly required. Phase 2 of the architectural evolution was specifically engineered to resolve these flaws by replacing the static recording loop with an intelligent, dynamic speech detection mechanism.

To achieve adaptive listening without violating the stringent hardware constraints of the deployment environment, the architecture integrated the Silero Voice Activity Detection module. This specific module was selected because it utilizes a highly optimized, pre-trained neural network engineered for rapid and precise speech segmentation. A critical implementation decision was made regarding the hardware allocation for this new component. Because the core conversational models; specifically the Whisper transcription network and the Qwen large language model already consumed the vast majority of the six-gigabyte Video Random Access Memory budget on the target graphical processing unit, allocating additional video memory to a continuous background monitoring task would have unnecessarily increased pressure on the limited GPU memory. Consequently, the voice activity detection network, which featured a highly lightweight memory footprint of approximately 1.5 megabytes, was explicitly configured to execute entirely on the system's central processing unit. This deliberate offloading strategy ensured that continuous audio scanning could occur asynchronously without competing for the strictly limited graphical memory required for intensive generative inference.

With the monitoring neural network successfully isolated on the central processing unit, the audio capture module was redesigned to perform continuous microphone monitoring. Instead of blocking the application thread for a predetermined duration, the system was configured to stream raw audio directly from the local microphone in microscopic, discrete data chunks. Each chunk consisted of exactly 512 audio samples, equating to approximately 32 milliseconds of acoustic data at the required 16 kilohertz sample rate. As each audio chunk was captured, it was immediately converted into a floating-point tensor and evaluated by the voice activity detection model. The neural

network computed a statistical probability indicating the likelihood that the specific 32-millisecond segment contained human vocalization.

This continuous monitoring enabled the system to perform dynamic speech start detection. The software evaluated the output of the neural network against a strict, predefined probability threshold to filter out transient background noises, such as computer fan interference or environmental mechanical sounds. Once the acoustic probability exceeded this threshold, the system registered the onset of human speech and began accumulating the audio chunks into a temporary recording buffer.

Equally important to detecting the start of a conversation was the ability to perform dynamic speech endpoint detection. To mirror the natural flow of human dialogue, the system needed to accurately distinguish between a brief pause for breath and the actual completion of a spoken sentence. To accomplish this, the architecture implemented a continuous silence-tracking mechanism. Once active speech was detected, the system monitored all subsequent incoming audio chunks. If the neural network determined that a chunk lacked human speech, a silence counter was incremented. To accommodate natural mid-sentence pauses without prematurely cutting off the user, the software was configured to require twenty consecutive silent chunks before concluding the conversational turn. This equated to approximately 640 milliseconds of sustained absolute silence. The moment this 640-millisecond threshold was reached, the dynamic endpoint was triggered, and the microphone recording loop was instantly terminated. As a secondary architectural safety net, a hard timeout limit of fifteen seconds was also implemented to guarantee that the recording loop would eventually release system resources in the event of continuous, uninterrupted background noise mimicking human speech.

Because the length of the captured audio was no longer a predictable five-second constant, the interface between the audio capture module and the speech recognition component required structural modification. The architecture shifted to a variable-

length transcription methodology. The specific transcription functions invoked within the Whisper software library were updated to natively process arbitrary-length floating-point audio arrays, completely bypassing legacy methods that artificially padded or trimmed acoustic data to fit fixed temporal constraints.

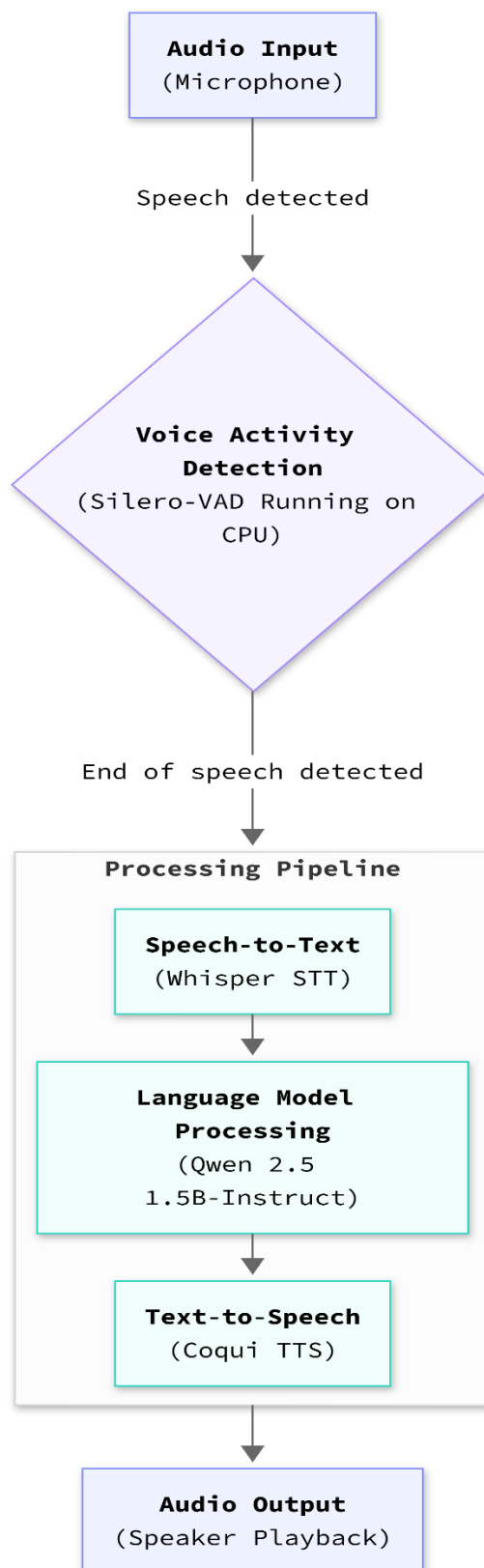


Figure 2. The Phase-2 adaptive listening architecture. The CPU-based VAD module dynamically gates the audio capture, passing data to the inference models only after detecting the end of user speech.

The implementation of this adaptive, variable-length capture architecture yielded immediate and profound improvements to the overall stability and reliability of the chatbot. By dynamically terminating the recording exactly 640 milliseconds after the user stopped speaking, the system successfully excluded empty room silence from the transcription payload. Consequently, the Whisper model was spared from processing ambient noise, which eliminated the severe artificial intelligence hallucinations observed in the Phase 1 baseline. Because the speech recognition module no longer generated spurious characters or repeating numerical sequences, the language model exclusively received validated, coherent human input. This protected the language model's conversational context from degradation, preventing the uncontrolled dynamic memory spikes during text generation that previously caused out-of-memory (OOM) crashes. Through intelligent speech segmentation, the Phase 2 architecture permanently stabilized the memory footprint of the system.

In addition to resolving the most significant reliability issues, the adaptive listening architecture significantly improved conversational responsiveness. Unnecessary processing overhead was significantly reduced because the Whisper model was no longer forced to analyze seconds of empty audio buffer. Furthermore, the overall turnaround time for brief interactions was substantially accelerated. If a user provided a short, one-second command, the system terminated the recording and immediately initiated transcription roughly 640 milliseconds later, rather than forcing the user to wait out the remainder of a rigid five-second timer.

Despite these massive improvements in system stability and the elimination of silence-induced errors, empirical evaluation revealed that the Phase 2 architecture remained constrained by significant conversational limitations. While the audio capture phase had been rendered dynamic and highly efficient, the overarching software architecture remained fundamentally sequential. The execution flow was still governed by a single, monolithic application thread operating under a strict stop-and-wait paradigm.

Because the architecture was sequential, considerable latency still existed within the conversational loop. The system successfully minimized the time spent in the recording phase, but it still required the transcription, language generation, and text-to-speech synthesis stages to execute to absolute completion, one after another, before any audio was played back to the user. This additive latency accumulation meant that while the system was highly reliable, it still lacked the rapid, overlapping responsiveness characteristic of natural human-to-human dialogue.

Furthermore, the single-threaded nature of the Phase 2 pipeline meant that the inability to interrupt audio playback remained entirely unresolved. The microphone was only active during the dedicated audio capture phase. The moment the 640-millisecond silence threshold was reached, and the audio array was passed to the speech recognition module, the microphone stream was completely deactivated. During the subsequent inference processing and the entirety of the synthesized audio playback, the system remained completely deaf to its environment. If the chatbot generated an unusually long response, the user was forced to listen to the entire acoustic output without any physical mechanism to interject, correct the chatbot, or redirect the conversation.

Ultimately, the Phase-2 architecture successfully engineered a highly stable, resource-efficient, and adaptive acoustic capture mechanism that solved the most critical memory and hallucination failures of the baseline system. However, the persistence of sequential processing bottlenecks and the rigid, unbreakable turn-taking structure demonstrated that optimizing the individual components was insufficient to achieve true conversational naturalness. Overcoming the additive latency of the sequential pipeline and enabling interactive voice interruptions necessitated a complete structural departure from the single-threaded paradigm, directly motivating the development of a fully decoupled, concurrent architecture in the final phase of the research.

4.3.3 Phase 3 - Concurrent Multithreaded Pipeline

Although the adaptive listening architecture developed in Phase 2 successfully resolved the most significant reliability issues associated with silence processing and inaccurate transcription caused by processing ambient silence, empirical observations revealed that the system remained fundamentally constrained by its sequential execution model. The rigid "stop-and-wait" paradigm mandated that each discrete conversational stage audio capture, speech recognition, language generation, and acoustic synthesis had to execute to absolute completion before the subsequent stage could commence. This strict linearity resulted in a significant accumulation of processing latency, as the total waiting time between a user concluding an utterance and the chatbot initiating a response was the arithmetic sum of all individual processing delays. Furthermore, the sequential flow inherently prevented the system from listening while speaking. Once the audio capture phase concluded and the inference operations began, the microphone was entirely deactivated. This completely neutralized the user's ability to naturally interrupt the chatbot or redirect the conversation during audio playback, enforcing a highly mechanical and unnatural turn-taking structure. To achieve a truly interactive and highly responsive conversational agent, the pipeline required a structural paradigm shift away from sequential processing.

Phase-3 of the architectural evolution was explicitly engineered to eliminate these sequential bottlenecks through the implementation of a fully decoupled, concurrent multithreaded architecture. The monolithic processing loop was decomposed into four isolated, asynchronous threads, each dedicated to a specific operational domain within the conversational pipeline. The first thread, designated as the capture thread, was solely responsible for continuous microphone streaming and voice activity detection. The second thread, the inference thread, managed the heavy computational inference workloads, sequentially executing the Whisper speech-to-text transcription and the Qwen large language model generation. The third thread, the output thread, handled the Coqui text-to-speech acoustic rendering and local speaker playback. Finally, a fourth control thread was implemented to monitor cross-platform local keyboard interrupts. By

decoupling these stages, the system transitioned to an asynchronous operational model where multiple processing tasks could overlap in time. For instance, the capture thread could immediately resume listening for new acoustic input the moment a finalized utterance was passed forward, operating simultaneously while the inference thread analyzed the previous input and the output thread synthesized the response.

To safely orchestrate the flow of data between these asynchronous threads without introducing race conditions or data corruption, the architecture implemented thread-safe communication pipelines utilizing bounded queues. Specifically, an audio queue bridged the capture thread and the inference thread, while a text queue connected the inference thread and the output thread. Crucially, these queues were strictly initialized with a maximum size limit of exactly one element. The implementation of this strict capacity boundary was a deliberate engineering decision designed to enforce backpressure throughout the concurrent pipeline. In a low-latency conversational system, utilizing unbounded queues introduces a severe risk of latency bloat and stale audio accumulation. If a user speaks multiple rapid sentences before the language model can generate responses, an unbounded queue would indiscriminately accumulate historical audio inputs. The inference models would eventually process outdated utterances, destroying the real-time relevancy of the conversation.

By constraining the communication queues to a single element, the architecture effectively ensures that the generative neural networks only ever process the single most recent, up-to-date user utterance. If the inference thread is actively processing a prompt, the full audio queue forces the capture thread to block or discard new inputs until the system is ready, naturally synchronizing the asymmetric workloads of the pipeline. Furthermore, this strict buffering mechanism provides critical memory protection benefits. It prevents unchecked data accumulation from inflating dynamic memory usage, ensuring that the heavy neural network models and the data pipelines remain securely within the strict memory limits of the deployment hardware.

The transition to a concurrent architecture uniquely enabled the integration of dynamic interruption support. Because the capture thread operates continuously and asynchronously from the output thread, the system possesses the theoretical capacity to monitor the microphone simultaneously while synthesizing audio playback. The architecture leverages this capability to implement a natural voice barge-in behavior. If the chatbot is actively speaking and the user decides to interject, the capture thread's voice activity detection module actively evaluates the incoming acoustic stream. To prevent transient noises from erroneously cutting off the chatbot, the system requires a sustained period of approximately 640 milliseconds, equivalent to twenty consecutive processing chunks of validated human speech to trigger an interruption. Once this sustained speech threshold is detected, the capture thread broadcasts a global interrupt signal. This signal instantly terminates the active audio playback, flushes the remaining text from the processing queues, and forces the system to immediately focus on the newly captured user's utterance.

While the voice barge-in mechanism proved effective during testing; however, practical deployment on consumer-grade hardware revealed a critical physical limitation regarding acoustic echo. Standard computing environments, such as open laptops, lack dedicated hardware-based Acoustic Echo Cancellation. Consequently, when the chatbot played its synthesized response through open speakers, the active microphone captured the chatbot's own voice. The voice activity detection module would interpret this acoustic bleed as human speech, inadvertently causing the chatbot to trigger an interruption against itself. To resolve this hardware-induced self-interruption, a specific software toggle, denoted as the headset mode, was introduced into the capture thread logic. When this mode is disabled, the architecture enforces a strict microphone muting protocol during active speaker playback, effectively suppressing the acoustic echo but temporarily sacrificing the user's ability to perform voice interruptions. Conversely, when the mode is enabled, operating under the assumption that the user is wearing headphones, which physically isolates the acoustic output from the microphone, the system keeps the microphone continuously active. This headset configuration unlocks

the full potential of the concurrent architecture, permitting true voice barge-in and highly overlapping, natural conversational dynamics.

To supplement the voice barge-in capabilities, particularly when the system operates without headphones, the dedicated control thread provides robust keyboard interruption handling. The software monitors raw input streams across both Windows and Unix-based operating systems without blocking the main execution, utilizing different underlying system libraries to ensure seamless cross-platform compatibility. Pressing the spacebar dispatches an immediate halt signal to the output thread, instantly silencing the current speech playback and allowing the user to seamlessly take over the conversation, while a separate key command initiates a safe, synchronized shutdown of all active threads.

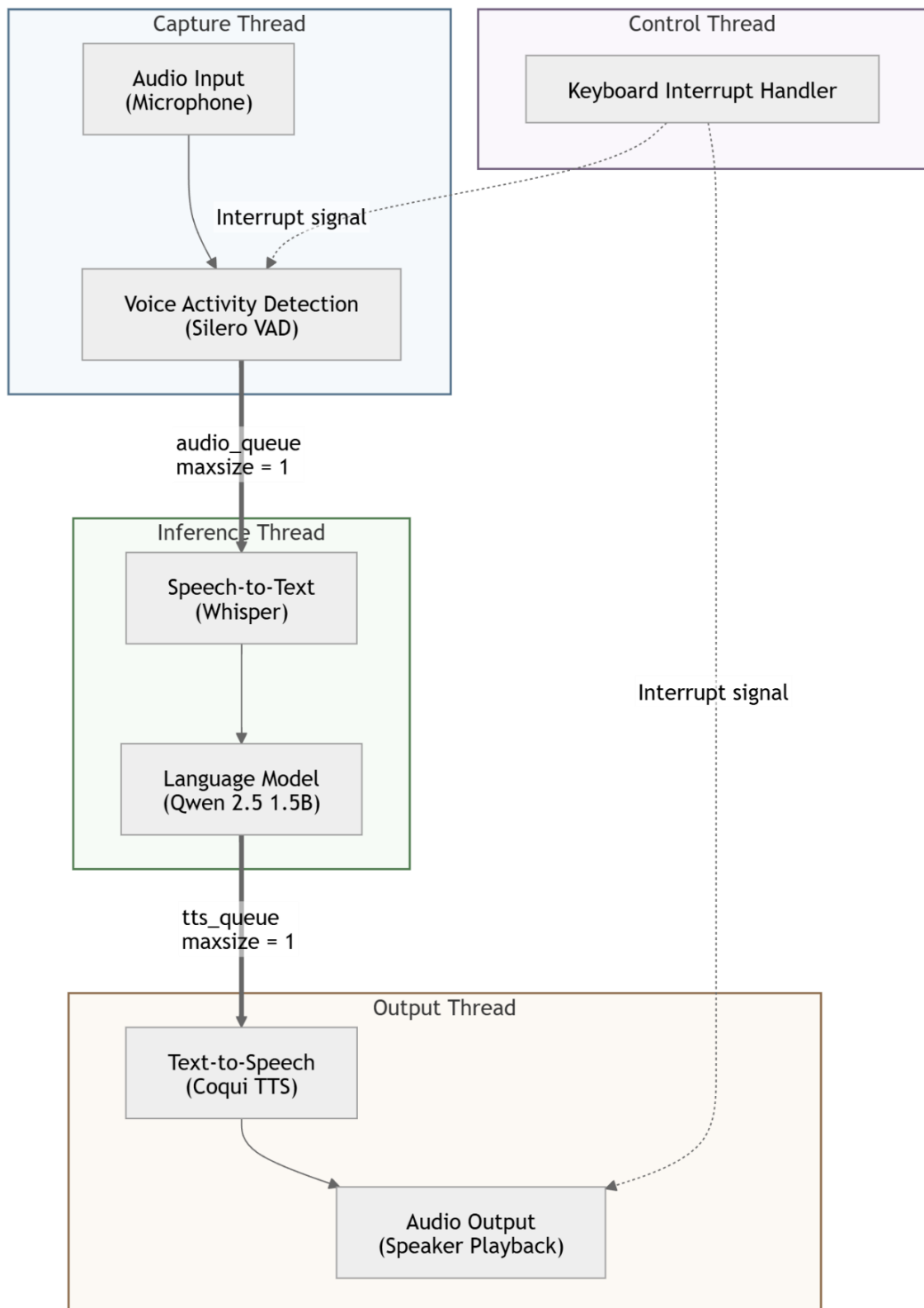


Figure 3. The Phase-3 concurrent multithreaded architecture. Independent processing threads communicate through bounded queues while interruption signals enable responsive conversational interaction.

The implementation of this concurrent multithreaded architecture fundamentally transformed the operational capabilities of the localized chatbot. By decoupling the monolithic loop, the system achieved a profound improvement in responsiveness. The total waiting time between a user completing an utterance and the initiation of the chatbot's vocal response was significantly reduced, as the staggered, asynchronous processing eliminated the mandatory idle states inherent in the sequential baseline. Furthermore, the successful integration of both hardware-based keyboard interrupts and true acoustic voice barge-in established a highly dynamic interaction model. Users were no longer constrained by rigid, unbreakable turn-taking, but could instead interject, correct, and converse with the artificial intelligence utilizing the same overlapping, fluid patterns characteristic of natural human dialogue. The architectural evolution from a simple sequential script into an adaptive, multithreaded pipeline successfully neutralized the primary latency and rigidity bottlenecks, fulfilling the software engineering prerequisites for edge-based conversational artificial intelligence.

With the architectural framework stabilized and the concurrent processing pipeline fully established, the subsequent phase of the research focused on empirical validation. The following section presents the comprehensive performance evaluation of the final Phase-3 architecture, including memory utilization, latency characteristics, and user evaluation results.

4.4 Performance Evaluation

The transition from a theoretical framework to a practical deployment necessitates rigorous empirical validation to confirm the system's overall operational viability. To ensure the localized chatbot met the foundational requirements for edge-based conversational artificial intelligence, the implemented architecture was evaluated across three distinct, complementary dimensions. First, memory utilization was systematically profiled to verify that the integrated neural networks could operate safely within the strict physical boundaries of consumer-grade hardware without triggering memory

exhaustion. Second, a quantitative latency analysis was performed utilizing precise runtime measurements to objectively demonstrate the temporal acceleration achieved by the concurrent pipeline. Finally, qualitative user observations were gathered through human-in-the-loop testing to evaluate the subjective conversational naturalness, interactive flow, and the efficacy of the interruption mechanics.

4.4.1 Memory Utilization

The most inflexible physical limitation governing the deployment of localized conversational artificial intelligence is the maximum memory capacity of the target hardware. The primary environment for this research was constrained to an Nvidia RTX 3060 graphics processing unit, which possesses a strict upper boundary of six gigabytes of Video Random Access Memory (VRAM). In a multi-model pipeline, the automatic speech recognition, natural language generation, and acoustic synthesis models must all reside in the graphical memory simultaneously to enable real-time, overlapping execution. Exceeding this six-gigabyte boundary may result in out-of-memory errors, potentially causing application instability or termination. Therefore, the foundational metric for edge deployment feasibility was validating the system's steady-state memory footprint.

Empirical measurements recorded during the continuous execution of the concurrent pipeline demonstrated that the system successfully adhered to this rigorous hardware constraint. The total steady-state VRAM utilization of the active architecture stabilized at approximately 3.3 gigabytes. This overall memory footprint is unevenly distributed among the generative neural networks required for the conversational loop. The vast majority of the graphical memory is consumed by the natural language generation engine, the Qwen 2.5 1.5B-Instruct model. Operating in 16-bit floating-point precision, this language model requires an allocation of approximately 3.0 gigabytes of VRAM. The automatic speech recognition component, utilizing the Whisper 'base' architecture, demands a highly efficient memory footprint of approximately 145 megabytes. For the final acoustic output, the Coqui Text-to-Speech synthesis engine requires an allocation

of approximately 150 megabytes to render the waveform responses. Furthermore, to protect the critical graphical memory budget, the Silero Voice Activity Detection neural network, which features a footprint of approximately 1.5 megabytes, was successfully routed to execute entirely on the central processing unit rather than the graphical processing unit.

Table 1. Steady-State VRAM Utilization of Core Conversational Pipeline Components

Component	Function	Approximate Memory Usage
Whisper Base	Speech-to-text processing	145 MB
Qwen 2.5 1.5B-Instruct	Language generation	3000 MB
Coqui TTS	Speech synthesis	150 MB
Silero VAD	Voice activity detection	1.5 MB
Total	Entire active pipeline	~3300 MB

The analysis of these memory measurements confirms the practical feasibility of hosting sophisticated, multi-modal artificial intelligence natively on accessible edge hardware. By consuming roughly 3.3 gigabytes of graphical memory during steady-state execution, the concurrent architecture utilizes approximately 55 percent of the available hardware capacity. This configuration intentionally leaves a substantial, unallocated hardware reserve of approximately 45 percent. Maintaining this significant margin of unused VRAM capacity is a paramount requirement for both immediate system stability and long-term deployment viability.

During live conversational interactions, large language models do not maintain a completely static memory footprint. The generation process requires dynamic memory allocation to manage accumulating context windows, expanding conversation histories, and the caching of self-attention keys and values. As the dialogue lengthens over a session, the temporary memory required to generate new tokens inherently expands. If an architecture's static footprint consumes the entirety of the hardware budget, these necessary dynamic memory spikes will breach the hardware limitations and trigger fatal process crashes. The 45 percent safety reserve provides sufficient capacity to absorb

these dynamic memory fluctuations during extended conversational turns without threatening core stability.

Furthermore, this substantial hardware reserve accommodates the operational realities of consumer-grade computing environments. Unlike dedicated, headless enterprise servers, local edge devices rely heavily on their graphical processing units to manage underlying operating system tasks, render graphical user interfaces, and execute concurrent background applications. The verified 45 percent reserve ensures that the voice chatbot can function as an unobtrusive, stable background process without aggressively competing with the host operating system for critical graphical resources. Finally, this margin of unused capacity provides essential headroom for future architectural scalability. The landscape of open-weight language and speech models evolves rapidly. The demonstrated memory buffer ensures that the current system architecture possesses the fundamental flexibility to integrate larger, more capable generative models or more sophisticated software-based acoustic echo cancellation modules in subsequent iterations without requiring immediate, prohibitive upgrades to the underlying edge hardware.

4.4.2 Latency Analysis

In the domain of conversational artificial intelligence, end-to-end processing latency is an extraordinarily critical factor determining the practical viability of a system. Human cognitive and perceptual systems are highly sensitive to interaction delays during verbal communication. Natural conversational dynamics rely on rapid turn-taking, and excessive computational delays actively disrupt this flow, forcing human users into rigid, unnatural communication patterns. If a voice chatbot requires multiple seconds to process an input and generate a response, the interaction inevitably feels mechanical, unresponsive, and severely degraded, regardless of the linguistic sophistication of the underlying language model. Consequently, evaluating the temporal acceleration achieved by the software implementation is as paramount as verifying its memory stability. To objectively quantify these temporal improvements, a rigorous comparative

analysis was conducted between the initial Phase 1 sequential baseline and the final Phase 3 concurrent multithreaded architecture.

The quantitative evaluation of conversational latency relies on exact per-turn timing measurements extracted directly from the system's runtime logs during steady-state execution. The analysis dissects the total system turnaround latency into four distinct operational stages: the audio capture delay, the speech recognition latency, the language model latency, and the speech synthesis latency. In the Phase 1 sequential baseline architecture, the total system turnaround time averaged a highly disruptive 7.00 seconds. This massive delay was primarily dominated by the audio capture module, which operated on a strictly hardcoded recording loop, introducing an inflexible 5.00-second audio capture delay for every interaction. Following this mandatory capture period, the Whisper speech recognition latency averaged approximately 0.55 seconds to transcribe the audio. The Qwen large language model latency accounted for approximately 1.25 seconds of text generation time, and the Coqui speech synthesis latency added a final 0.20 seconds to render the acoustic waveform. Because the baseline architecture was strictly sequential, these individual processing delays were entirely additive, directly resulting in the unacceptable 7.00-second total turnaround latency.

The transition to the Phase 3 concurrent multithreaded architecture yielded a profound acceleration in system responsiveness. Empirical runtime measurements of the final deployment demonstrate that the total system turnaround latency decreased to just 1.87 seconds. This represents a massive net latency reduction of 5.13 seconds, equating to a 73 percent decrease in total conversational delay compared to the sequential baseline. When analyzing the per-stage latency breakdown of the final architecture, the audio capture delay was drastically reduced to an average of approximately 1.83 seconds. In contrast, the processing times for the generative neural networks exhibited only marginal, fractional improvements. The speech recognition latency improved slightly from approximately 0.55 seconds to 0.49 seconds. The language model latency

remained practically static, shifting from 1.25 seconds to 1.23 seconds, and the speech synthesis latency similarly experienced only a negligible reduction from 0.20 seconds to 0.15 seconds.

This precise distribution of latency reductions highlights a fundamental conclusion regarding the software implementation. The speech recognition, language model, and speech synthesis components remained largely unchanged across the developmental phases, utilizing the exact same neural network weights, hardware allocations, and generative parameters. Because the inference processing times of these core artificial intelligence models were virtually identical in both the slowest and fastest versions of the chatbot, it is evident that the overwhelming majority of the 73 percent total latency reduction was achieved strictly through software architectural redesign. Specifically, the largest single latency improvement was directly generated by eliminating the fixed 5.00-second recording window and replacing it with the dynamic Voice Activity Detection module. By terminating the audio capture immediately after detecting approximately 640 milliseconds of human silence, the system entirely bypassed the necessity of recording and processing empty ambient noise, cutting over three seconds of artificial delay from the pipeline before the inference stages even commenced.

Table 2. Comparison of Average Processing Latency

Processing Stage	Phase 1 (Sequential Baseline)	Phase 3 (Concurrent Multithreaded)	Net Latency Change
Audio Capture	5.00 s (<i>Fixed Window</i>)	~ 1.83 s (<i>Dynamic VAD</i>)	- 3.17 s
Whisper STT	~ 0.55 s	0.49 s	- 0.06 s
Qwen LLM Generation	~ 1.25 s	1.23 s	- 0.02 s
Coqui TTS Synthesis	~ 0.20 s	0.15 s	- 0.05 s
Total System Turnaround	~ 7.00 s	~ 1.87 s	- 5.13 s

These latency reductions are significant in the context of edge-based conversational artificial intelligence systems. By dropping the total system turnaround time from 7.00 seconds down to a highly responsive 1.87 seconds, the architecture fundamentally transforms the user experience from an asynchronous, stop-and-wait interface into a dynamic conversational agent. The transition from sequential blocking to concurrent, asynchronous thread management ensures that the central processing unit and graphical processing unit are utilized optimally, eliminating the idle waiting states that previously bottlenecked the conversational loop. Achieving sub-two-second latency strictly on consumer-grade hardware proves that local computing environments can successfully achieve the responsive conversational interaction traditionally associated with massive, cloud-reliant data centers.

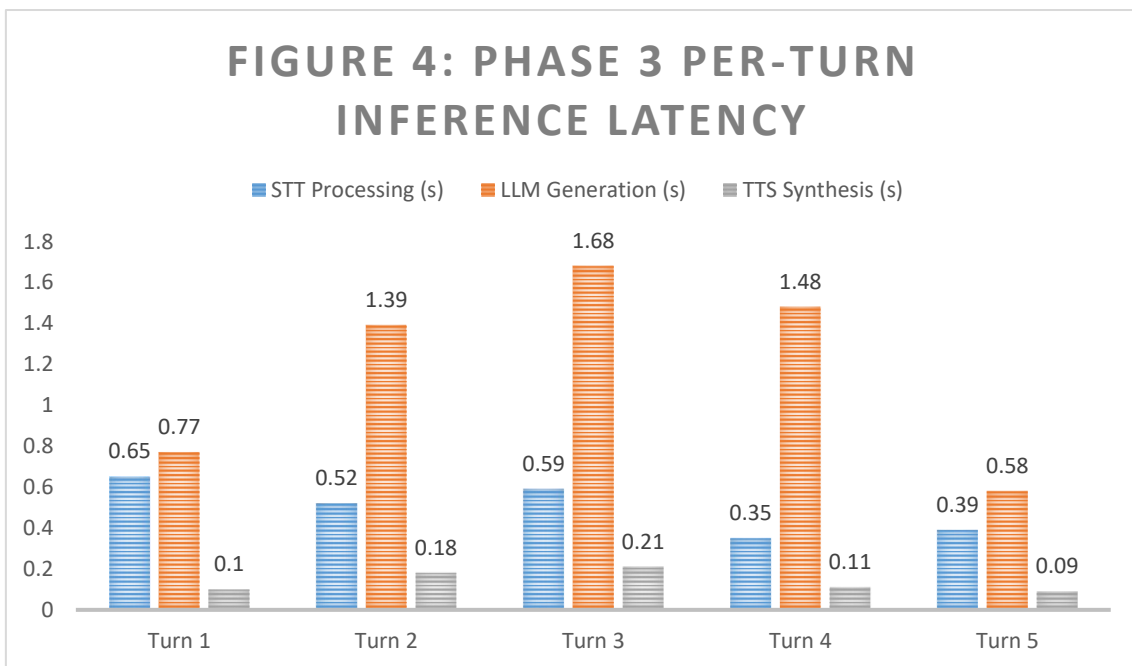


Figure 4. Per-turn processing times recorded during steady-state Phase 3 execution. Speech-to-text processing, language model generation, and text-to-speech synthesis remained consistently below 2.5 seconds total across multiple conversational turns.

An analysis of the per-turn processing times recorded during steady-state Phase 3 execution, as illustrated in Figure 4, further validates the robustness of the concurrent architecture. The empirical logs demonstrate a high degree of consistency in processing

times across varied and continuous interactions. Despite the variable nature of user speech lengths and the dynamic token generation of the language model, the combined inference operations; specifically the speech-to-text and language model threads, alongside the text-to-speech output thread, consistently execute in under 2.5 seconds total. This remarkable stability across multiple conversational turns proves that the implementation of thread-safe, bounded communication queues effectively regulates the flow of data through the system.

By enforcing backpressure through these maximum-size queues, the architecture entirely prevents the accumulation of stale audio arrays or unprocessed text tokens, which would otherwise cause latency to bloat uncontrollably during extended dialogue sessions. Instead, the processing times remain anchored and consistent, guaranteeing that the system performance does not degrade over time. The ability to maintain this level of temporal consistency ensures a near-real-time interaction capability, establishing a highly reliable and predictable foundation for continuous, autonomous human-machine dialogue at the network edge.

While the quantitative latency measurements and runtime logs definitively prove the temporal acceleration and computational efficiency of the concurrent software architecture, evaluating a voice-based artificial intelligence system necessitates more than mathematical benchmarking. The ultimate goal of conversational artificial intelligence is to facilitate seamless interactions with human users. To fully understand how these objective latency reductions and architectural modifications translate into practical usability, the system must be evaluated based on human perception. Therefore, following the quantitative runtime analysis, the architecture was subjected to qualitative human-in-the-loop evaluations. The following section details the setup and results of this subjective testing, focusing on how the mathematical improvements directly impacted the perceived conversational flow.

4.4.3 Subjective Testing

While quantitative metrics such as steady-state memory utilization and per-turn processing latency provide objective validation of the system's operational efficiency, evaluating a conversational artificial intelligence system inherently requires assessing the quality of human-machine interaction. Conducting a formal, large-scale user evaluation study with statistically significant sample sizes and quantitative questionnaires was outside the scope of this research. However, the software was successfully demonstrated during the thesis defense seminar, where attendees were invited to observe and interact with the operational system in real time. Although no formal participant feedback or numerical survey scores were collected during this public demonstration, comprehensive qualitative observations were systematically gathered throughout the development lifecycle. These subjective evaluations were conducted through iterative developer testing, cross-device hardware testing, supervisor evaluation, and informal peer testing involving a small group of users interacting with the various developmental phases.

The informal peer testing sessions were specifically designed to assess the qualitative improvements introduced by the architectural evolution. Evaluators observed the system's behavior based on three primary subjective metrics: conversational naturalness, interaction flow, and interruption handling capabilities. During the evaluation of the Phase 1 sequential baseline, qualitative feedback confirmed that the system felt highly mechanical and unresponsive. Because the baseline architecture utilized a fixed five-second recording window and lacked interruption capabilities, testers were forced into a rigid, unnatural turn-taking structure. Conversely, when evaluating the final concurrent multithreaded architecture, testers reported a noticeable improvement in interaction flow. The dynamic speech segmentation and asynchronous processing introduced in Phases 2 and 3 translated into a noticeably more responsive conversational experience, allowing users to communicate using more natural cadences and variable-length sentences with fewer perceived processing delays or silence-induced hallucinations.

A critical component of the subjective evaluation involved rigorous testing of the system's newly introduced interruption mechanics, which generated vital feedback regarding cross-platform software behavior. During a dedicated evaluation session conducted by the thesis supervisor, a specific functional anomaly was identified within the voice barge-in architecture. The supervisor observed that when the user attempted to interrupt the chatbot mid-sentence utilizing voice commands, the command-line interface successfully registered the interruption signal and immediately halted the underlying text generation processes. However, the active audio playback from the text-to-speech engine failed to stop simultaneously, allowing the chatbot to finish speaking its previously queued sentence despite the successful computational interrupt.

This specific supervisor feedback directly guided subsequent cross-device developer testing, which revealed a nuanced operating system-level audio buffering limitation. The software had been engineered to support cross-platform execution across both Windows and Unix-based environments. While the audio playback interruption behaved seamlessly on the Windows development environment, testing revealed that Linux audio backend servers, such as Advanced Linux Sound Architecture (ALSA) or PulseAudio, struggled to flush the hardware audio buffer instantly upon receiving a software stop command. This buffering lag caused the text-to-speech engine to temporarily block the thread during synthesis on Linux distributions, creating a discrepancy between the logical interruption and the acoustic output. Furthermore, this cross-platform evaluation validated the implementation of the hardware keyboard interrupts, confirming that the system successfully managed the differing underlying input libraries without failing, utilizing standard input monitoring on Windows and terminal input/output interfaces on Linux to provide a reliable manual stopping mechanism.

Beyond operating system nuances, interactive testing exposed significant physical hardware challenges related to localized acoustic feedback. The concurrent architecture theoretically permitted the microphone to remain continuously active, listening for new

input while the chatbot simultaneously played synthesized audio responses. However, practical observations demonstrated that standard consumer-grade computing environments, such as open laptops, inherently lack dedicated hardware-based Acoustic Echo Cancellation. Without this hardware isolation, the active microphone captured the chatbot's own synthesized speech bleeding from the external speakers. The voice activity detection module logically interpreted this acoustic bleed as human speech, inadvertently causing the chatbot to trigger the voice barge-in protocol against its own output, trapping the system in a loop of self-interruption.

To resolve this hardware limitation while preserving the concurrent functionality, the qualitative evaluation phase prompted the implementation of a dedicated software toggle designated as the headset mode. Testing confirmed that this toggle provided a necessary practical compromise for edge deployment. When the system was operated with the mode disabled, the architecture enforced a strict microphone muting protocol during active speaker playback. This effectively suppressed the acoustic echo and stabilized the conversation, though it temporarily sacrificed the user's ability to perform voice interruptions. Conversely, when the mode was enabled, under the assumption that the user was wearing headphones to physically isolate the acoustic output from the microphone array, the system permitted true voice barge-in. Evaluators observed that this headset configuration unlocked the full potential of the multithreaded architecture, enabling highly overlapping, dynamic, and natural conversational interruptions that closely mirrored human-to-human dialogue.

Ultimately, the combination of quantitative hardware benchmarking and qualitative interaction testing confirmed that the system successfully fulfilled its primary objectives. The architectural evolution from a rigid sequential script to a concurrent, adaptive pipeline not only addressed the primary memory-management and latency challenges identified during development, but also yielded a robust, highly interactive conversational agent. The following section provides a comprehensive summary of these

evaluation findings, synthesizing the performance metrics and subjective observations to establish the overall viability of the implemented edge voice chatbot.

4.5 Summary of Findings

The implementation and evaluation of the modular edge voice chatbot successfully demonstrated the practical feasibility of hosting highly responsive, privacy-preserving conversational artificial intelligence on consumer-grade hardware. By systematically transitioning from a basic sequential script to a robust, concurrent pipeline, the developed system achieved the primary architectural and operational objectives established at the onset of the research. Furthermore, by operating entirely offline without relying on external cloud-based application programming interfaces, the software architecture significantly improved data privacy by eliminating the need to transmit user speech to external cloud services, addressing a fundamental vulnerability inherent in traditional voice assistants.

Empirical hardware validation proved that complex, multi-modal neural networks could operate safely within strictly constrained deployment environments. The strategic integration of the Whisper speech recognition model, the Qwen 2.5 large language model, and the Coqui speech synthesis engine resulted in a steady-state graphical memory footprint of approximately 3.3 gigabytes on the target Nvidia RTX 3060 hardware. By explicitly routing the lightweight voice activity detection network to execute on the central processing unit, the architecture deliberately preserved critical graphical memory. This orchestration maintained a substantial safety reserve of nearly half the total available video memory, effectively eliminating the out-of-memory crashes observed during early development and improving operational stability during extended interactions and ensuring enduring operational stability during extended, dynamic dialogue generation.

The quantitative temporal evaluation confirmed that while hardware constraints necessitated the use of specific, highly optimized generative models, the most profound

performance enhancements were achieved through architectural redesign rather than modifications to the underlying neural network weights. The evolution from the Phase 1 sequential baseline to the Phase 2 adaptive listening architecture successfully resolved critical vulnerabilities associated with processing ambient silence. By dynamically terminating audio capture based on intelligent speech boundary detection, the system eradicated silence-induced artificial intelligence hallucinations. The subsequent transition to the Phase 3 concurrent multithreaded architecture fundamentally resolved the additive latency bottlenecks inherent in sequential processing. By decoupling audio capture, inference, and acoustic synthesis into isolated, asynchronous threads communicating via bounded queues, the system achieved a substantial reduction in total conversational turnaround time, dropping from a highly disruptive delay to near-real-time responsiveness. Because the core inference speeds of the underlying machine learning models remained largely unchanged throughout development, the observed latency improvements can be primarily attributed to the architectural redesign and concurrent software orchestration.

Qualitative evaluations and interactive observations corroborated the quantitative improvements, revealing a substantially enhanced conversational experience. The implementation of thread-safe asynchronous processing not only accelerated response times but also enabled highly dynamic interaction mechanics. Evaluators observed that the integration of voice barge-in capabilities, facilitated by continuous microphone monitoring and intelligent acoustic echo suppression via the headset mode toggle, allowed users to seamlessly interrupt and redirect the chatbot mid-sentence. Supplemented by cross-platform hardware keyboard interrupts, these capabilities shattered the rigid turn-taking structure of earlier software iterations, cultivating an adaptive conversational flow that closely mirrored natural human-to-human interaction.

Ultimately, the empirical evidence gathered throughout the implementation and evaluation phases confirms that localized, multimodal conversational artificial intelligence can be viably and efficiently deployed at the network edge. The successful

orchestration of these software components comprehensively resolves the latency and privacy limitations traditionally associated with cloud-reliant architectures. Having established the functional success of the system, the subsequent chapter will provide a broader discussion of these findings, examining the overall implications of the research, analyzing specific hardware nuances and operational limitations, and proposing strategic directions for future developmental work.

5 Discussion and Future Work

The implementation and evaluation presented in the preceding chapter demonstrated the practical feasibility of deploying a fully localized conversational artificial intelligence system on consumer-grade edge hardware. However, beyond the quantitative performance metrics and implementation outcomes, it is necessary to examine the broader implications of these findings in relation to the research objectives and the existing challenges of edge-based conversational systems. This chapter discusses the significance of the results, evaluates the limitations encountered during development and deployment, and outlines potential directions for future research and system enhancement.

5.1 Discussion of Research Findings

The overarching ambition of this research was to address the inherent systemic limitations of traditional, cloud-reliant conversational artificial intelligence by engineering a highly responsive, privacy-preserving, and fully localized voice chatbot. Historically, the continuous streaming of sensitive audio data to remote server clusters introduced severe end-to-end communication latencies and profound vulnerabilities regarding user privacy. By successfully shifting the computational workload entirely to the network edge, this thesis demonstrated that sophisticated natural language understanding and acoustic synthesis can be reliably hosted on standard, consumer-grade hardware. The empirical results, obtained through the iterative development of a modular software architecture, directly address the foundational research questions and validate the theoretical advantages of edge computing.

The first research question (RQ1) explored how a modular, fully localized voice chatbot architecture could be designed to support natural spoken interaction on consumer-grade edge hardware. The findings of this study demonstrate that achieving natural interaction necessitates abandoning monolithic software designs in favor of decoupled, task-specific operational boundaries. The successful design was realized by integrating four distinct

components: the Whisper 'base' model for automatic speech recognition, the Qwen 2.5 1.5B-Instruct large language model for natural language generation, the Coqui engine for text-to-speech synthesis, and the Silero Voice Activity Detection module for acoustic segmentation. By isolating these components into a modular framework, the architecture successfully replicated the multimodal capabilities of enterprise-grade cloud assistants while operating entirely offline. The significance of this offline operation lies in its ability to reduce reliance on external cloud infrastructure while improving user privacy. The system was designed so that user data remains on the physical device, successfully answering the imperative for secure, autonomous conversational artificial intelligence.

The second research question addressed the highly challenging problem of integrating these computationally heavy speech recognition, language generation, and speech synthesis components efficiently under constrained memory resources without triggering system instability. The deployment environment, restricted to an Nvidia RTX 3060 graphics processing unit with a strict six-gigabyte Video Random Access Memory (VRAM) boundary, represented a formidable physical limitation. The findings indicate that system stability under such constraints requires highly strategic model selection, asymmetric hardware allocation, and strict data flow regulation. By utilizing the 16-bit floating-point precision Qwen 1.5B model alongside the lightweight Whisper 'base' and Coqui synthesis models, the architecture achieved a steady-state graphical memory footprint of approximately 3.3 gigabytes. The significance of this footprint lies not merely in fitting within the hardware boundary, but in deliberately maintaining a substantial unallocated memory reserve. This reserve proved to be a stable deployment under constrained hardware conditions. During natural conversation, large language models inherently demand dynamic memory expansion to manage accumulating context windows and cache self-attention keys. The reserved capacity successfully absorbed these dynamic fluctuations, eliminating the out-of-memory issues encountered during early development that plagued early developmental phases.

Furthermore, the resolution to Research Question 2 heavily relied on the implementation of strict data flow regulation through concurrent systems theory. To prevent memory exhaustion caused by the accumulation of unprocessed data, the architecture implemented bounded communication queues with a maximum capacity of exactly one element. This design introduced necessary backpressure into the pipeline. By constraining the queues, the system mathematically guaranteed that the heavy generative models only processed the single most recent user utterance. This prevented the unchecked buffering of historical audio, ensuring that the hardware was never burdened by stale data processing. Additionally, the strategic offloading of the Silero Voice Activity Detection module to the central processing unit served as a critical optimization. By confining continuous background acoustic monitoring to the central processor, the architecture shielded the scarce graphical memory, reserving it exclusively for the most intensive neural network inference tasks.

The third research question sought to identify the architectural optimization strategies that could reduce end-to-end conversational latency while preserving system stability. The evolutionary transition from the Phase 1 sequential baseline to the Phase 3 concurrent multithreaded architecture provided the definitive answer to this question. The initial sequential pipeline yielded a highly disruptive average turnaround latency of 7.00 seconds. However, empirical runtime analysis revealed that the core inference speeds of the neural networks remained largely static throughout all developmental phases. Therefore, the substantial reduction in total latency to just 1.87 seconds in the final architecture was achieved strictly through software-based architectural orchestration.

The most profound optimization strategy for reducing latency was the implementation of adaptive listening. By abandoning the rigid, fixed-duration five-second recording window and integrating dynamic speech boundary detection, the system immediately eliminated over three seconds of artificial delay. Terminating the audio capture exactly 640 milliseconds after detecting human silence not only accelerated the conversational

loop but also eradicated silence-induced artificial intelligence hallucinations, preventing the Whisper model from transcribing ambient noise into polluting text.

Equally significant was the transition to a fully concurrent, asynchronous processing pipeline. By decomposing the monolithic loop into isolated capture, inference, and output threads, the architecture neutralized the additive latency of sequential blocking. Asynchronous processing allowed the system to overlap operational tasks in time; the capture thread could immediately resume listening for subsequent input while the inference and output threads finalized the previous response. This decoupling proved essential for enabling natural conversational dynamics. Because the microphone remained active independently of the generative models, the system successfully supported voice barge-in capabilities. Users were empowered to interrupt the chatbot mid-sentence with vocal commands or keyboard interrupts, reducing the rigid turn-taking behavior observed in the baseline implementation. Ultimately, the successful deployment of concurrent thread management, bounded queues, and adaptive listening proved that intelligent software orchestration can mitigate many of the latency challenges associated with edge deployment, delivering a highly responsive and fluid conversational agent.

5.2 Limitations and Hardware Nuances

While the implemented architecture successfully demonstrated the viability of localized edge voice chatbots, the system remains subject to several critical limitations and hardware-induced nuances. A comprehensive academic evaluation necessitates a balanced discussion of these constraints, which primarily stem from the reliance on consumer-grade hardware, environmental variability, and the inherent trade-offs between processing speed and generative capability.

The most prominent limitation of the system is its absolute dependence on the computational boundaries of consumer-grade graphical processing units. The strict six-gigabyte VRAM constraint of the target Nvidia RTX 3060 mandated the use of relatively

small neural network architectures. Consequently, the pipeline utilized the 1.5 billion parameter Qwen model and the Whisper 'base' transcription model. While these models are highly efficient, they possess fundamentally lower representational capacities than the massive, hundred-billion-parameter models deployed in cloud-scale systems. The localized language model, while capable of maintaining coherent conversational context, lacks the deep, generalized world knowledge, advanced logical reasoning, and complex problem-solving proficiencies characteristic of enterprise-grade artificial intelligence. Similarly, the Whisper 'base' model sacrifices absolute transcription precision for execution speed, leading to occasional transcription errors when processing highly complex vocabulary or thick regional accents. Therefore, a persistent trade-off exists between the autonomy and privacy of edge execution and the generalized accuracy of the deployed artificial intelligence.

In addition to constraints on model size, the system demonstrated significant sensitivity to acoustic environmental variability. The efficacy of the Silero Voice Activity Detection module, which dictates the dynamic segmentation of the entire pipeline, relies on distinguishing human speech from ambient noise. In testing environments characterized by high levels of non-stationary background noise, overlapping conversations, or mechanical interference, the voice activity detection threshold occasionally struggled to accurately identify the cessation of speech. This environmental interference could lead to either premature truncation of user utterances or delayed recording termination, inadvertently inflating response latency. The reliance on physical microphone quality further compounded this limitation, as lower-fidelity hardware introduced acoustic artifacts that degraded both speech segmentation and subsequent transcription accuracy.

A highly specific and disruptive hardware nuance observed during evaluation was the phenomenon of localized acoustic echo. Because the concurrent architecture was designed to keep the microphone active to support voice interruptions, utilizing standard, open laptop speakers created a detrimental physical feedback loop. The active

microphone captured the chatbot's own synthesized speech, which the voice activity detection module erroneously interpreted as new user input. This acoustic bleed trapped the system in a continuous cycle of self-interruption. To mitigate this hardware deficiency, the architecture required the implementation of a "headset mode" software toggle. While enforcing strict microphone muting during open-speaker playback successfully stabilized the system, it temporarily revoked the user's ability to perform voice interruptions. The necessity of physical acoustic isolation via headphones to unlock the full concurrent potential of the system highlights a significant gap in standard edge computing hardware, which typically lacks dedicated, hardware-level Acoustic Echo Cancellation.

Finally, the evaluation methodology employed during this research presents certain academic limitations. While the quantitative measurements of memory utilization and processing latency were extracted using precise runtime logging, the subjective evaluation of conversational naturalness relied on informal peer testing, supervisor observation, and public demonstration. The research did not include a large-scale, statistically significant human-computer interaction study utilizing standardized usability questionnaires. Consequently, while the qualitative feedback overwhelmingly confirmed improvements in interaction flow and responsiveness, the subjective findings lack the rigorous statistical validation required to definitively quantify the user experience across a diverse, multi-demographic participant pool.

5.3 Future Work

The successful implementation of the modular, concurrent edge voice chatbot establishes a robust foundational architecture that invites numerous avenues for future research and development. Addressing the identified limitations and further optimizing the system for edge environments provides a clear trajectory for advancing localized conversational artificial intelligence.

A primary direction for future work is the integration of advanced model quantization methodologies to safely deploy larger, more capable language models within strict memory constraints. The current implementation utilizes 16-bit floating-point precision. As explored in the theoretical foundations of this thesis, advanced post-training quantization techniques, such as Activation-aware Weight Quantization (AWQ) or Generative Pre-trained Transformer Quantization (GPTQ), can compress neural networks down to 4-bit integer representations with negligible accuracy degradation. Implementing these quantization frameworks would theoretically allow a consumer-grade GPU with six gigabytes of VRAM to host a significantly more intelligent 7-billion parameter language model, directly closing the reasoning and knowledge gap between edge-based and cloud-based systems without inducing memory exhaustion.

Furthermore, to relentlessly pursue latency reduction, future iterations of the architecture should explore fully streaming inference pipelines. The current Phase 3 architecture, while highly concurrent, still processes discrete blocks of data; the transcription model waits for the full audio chunk to finish recording, and the speech synthesis model waits for the language model to generate a complete text sentence. Transitioning to streaming speech recognition and streaming text generation would allow the pipeline to operate on a rolling basis. By feeding partial text tokens to a streaming text-to-speech engine the moment they are generated by the language model, the system could initiate acoustic playback virtually instantaneously, fundamentally blurring the lines of processing delay and achieving true real-time conversational overlap.

To resolve the hardware-induced limitations regarding acoustic echo and self-interruption, significant future engineering effort must be directed toward implementing robust Acoustic Echo Cancellation (AEC). Integrating sophisticated, software-based adaptive filtering algorithms into the audio capture thread could mathematically subtract the synthesized output waveform from the incoming microphone signal. A successful AEC implementation would permanently eliminate the necessity of the

"headset mode" compromise, allowing users to seamlessly interrupt the chatbot while utilizing open, ambient speaker configurations, thereby drastically improving the accessibility and naturalness of the interface.

Expanding the knowledge and utility of the localized chatbot without scaling the parameter count of the language model could be achieved through the integration of Retrieval-Augmented Generation (RAG). Future architectures could implement localized vector databases to allow the chatbot to search and retrieve relevant external documents, personal user notes, or domain-specific manuals during the generation phase. This would empower the edge agent to provide highly accurate, fact-based responses and maintain long-term personalized conversational memory, all while remaining strictly offline and preserving user privacy. Additionally, integrating multilingual support across the transcription, generation, and synthesis modules would broaden the global applicability of the system.

Finally, to validate the conversational efficacy of these proposed advancements, future research should encompass comprehensive, large-scale user evaluations. Conducting formal human-computer interaction studies with standardized metrics, such as the Mean Opinion Score for conversational naturalness and subjective latency perception, would provide the necessary statistical rigor to benchmark the localized architecture against commercial cloud alternatives. Ultimately, pursuing these developmental directions will continue to push the boundaries of what is computationally achievable at the network edge, advancing the democratization of autonomous, private, and highly intelligent conversational systems.

6 Conclusions

The proliferation of conversational artificial intelligence has fundamentally altered the landscape of human-computer interaction, enabling machines to understand and generate natural language with unprecedented proficiency. Historically, the deployment of these sophisticated systems has relied almost entirely on continuous, cloud-based computing paradigms. While this centralization leverages the immense computational power of remote data centers, it inherently introduces profound systemic vulnerabilities. The continuous transmission of captured audio payloads over wide area networks produces inevitable communication latencies, multi-hop routing delays, and bandwidth congestion that severely disrupt the rapid turn-taking required for natural human conversation. Furthermore, the mandatory streaming of raw, biometric voice data to third-party servers creates severe privacy and security risks, stripping users of control over their personal information and establishing a strict dependency on stable internet connectivity. Driven by the imperative to resolve these critical limitations, this research aimed to conceptualize, engineer, and evaluate a fully localized, edge-based conversational artificial intelligence system capable of operating autonomously on constrained consumer-grade hardware.

The primary outcome of this thesis is the successful development and implementation of a highly responsive, modular voice chatbot that executes entirely at the network edge. The proposed software architecture seamlessly integrates four state-of-the-art neural network components: the Whisper base model for automatic speech recognition, the Qwen 2.5 1.5B-Instruct large language model for natural language understanding and generation, the Coqui Text-to-Speech engine for acoustic waveform synthesis, and the Silero Voice Activity Detection module for intelligent audio segmentation. By deploying this comprehensive multimodal pipeline on a standard consumer-grade Nvidia RTX 3060 graphics processing unit, the research demonstrates that enterprise-level conversational capabilities can be democratized and hosted locally. Because all speech transcription, language reasoning, and voice synthesis operations occur exclusively on the physical

device, the architecture significantly improves privacy by ensuring that speech processing and language generation occur locally on the user's device. The system completely eliminates the transmission of sensitive acoustic data to external cloud infrastructures, reducing exposure of user data to external cloud services while providing robust offline autonomy.

Through the systematic evaluation of this architecture, the research provides definitive answers to the core questions guiding the study. The first research question investigated how a modular, fully localized voice chatbot architecture could be designed to support natural spoken interaction on consumer-grade edge hardware. The findings conclude that engineering such a system requires a fundamental departure from monolithic software designs in favor of decoupled, task-specific operational boundaries. The successful architecture isolated the acoustic input, inference processing, and acoustic output into a cohesive but modular framework. By replicating the functional pipeline of cloud-based assistants without the associated network dependencies, the architecture proved that local deployment is not only viable but highly effective. Natural spoken interaction was ultimately supported by ensuring that the decoupled modules could communicate efficiently, allowing the system to mirror the interactive, multimodal capabilities of human dialogue while preserving the autonomy of the edge environment.

The second research question addressed the profound challenge of integrating computationally intensive speech recognition, language generation, and speech synthesis components efficiently under constrained memory resources without triggering system instability. Operating within the strict six-gigabyte Video Random Access Memory boundary of the target graphics processing unit, the research concludes that system stability relies on strategic model selection, asymmetric hardware allocation, and rigorous data flow regulation. The successful integration utilized a 16-bit floating-point precision language model alongside highly efficient, lightweight acoustic models to establish a stable steady-state graphical memory footprint of approximately 3.3 gigabytes. Crucially, the research demonstrates that maintaining a substantial

unallocated memory reserve proved beneficial for preserving system stability. This substantial buffer successfully absorbed the dynamic memory spikes inherent in large language model generation; such as the expansion of context windows and self-attention caching, eliminating the out-of-memory issues observed during early development. Furthermore, system stability was improved by offloading the continuous background execution of the voice activity detection network to the central processing unit, thereby shielding the scarce graphical memory. To regulate data flow and prevent memory exhaustion from unprocessed inputs, the architecture implemented bounded communication queues with a maximum capacity of exactly one element. This application of backpressure theory ensured that the heavy generative models only processed the single most recent user utterance, preventing the accumulation of stale audio and synchronizing the asymmetric workloads of the pipeline.

The third research question sought to identify the specific architectural optimization strategies capable of reducing end-to-end conversational latency while preserving system stability in localized environments. The research concludes that effective latency reduction requires a combination of adaptive listening mechanics and concurrent multithreaded software orchestration. The evaluation of the initial sequential baseline architecture revealed a highly disruptive average turnaround latency of 7.00 seconds, exacerbated by a rigid, fixed-duration recording window. The primary optimization strategy that resolved this bottleneck was the implementation of dynamic speech segmentation. By utilizing the central processing unit-bound voice activity detection module to actively monitor the acoustic stream, the system was engineered to terminate audio capture exactly 640 milliseconds after detecting human silence. This adaptive listening strategy instantly eliminated over three seconds of artificial delay and eliminated the severe artificial intelligence hallucinations caused by forcing the transcription model to process ambient room noise.

In addition to adaptive listening, the transition to a fully concurrent, multithreaded pipeline was identified as the primary architectural strategy for minimizing latency. By

decomposing the pipeline into isolated, asynchronous capture, inference, and output threads, the architecture neutralized the additive latency of sequential blocking. This decoupling allowed computational tasks to overlap in time, reducing the total system turnaround latency to an average of just 1.87 seconds. Furthermore, concurrent execution permitted the microphone to remain actively monitored while the chatbot synthesized and played acoustic responses. This capability unlocked natural interruption handling, enabling users to seamlessly interject and redirect the chatbot mid-sentence through true voice barge-in and keyboard interrupts. The synthesis of these optimization strategies transformed the system from a mechanical, stop-and-wait program into a highly responsive conversational agent.

The primary scientific and practical contribution of this research lies in the successful synthesis of adaptive listening, concurrent multithreaded execution, and bounded queue architectures for edge-based artificial intelligence. By explicitly navigating the severe physical limitations of consumer-grade hardware, the study provides a replicable, optimized software framework that addresses key challenges related to latency, hallucination, and memory management. The empirical findings validate the feasibility of edge deployment, proving that sophisticated conversational systems do not inherently require the massive computational throughput of enterprise server farms to achieve real-time responsiveness.

Ultimately, the development and optimization of this fully localized voice chatbot underscore the broader significance of edge intelligence in modern computing. By effectively decoupling natural language generation and acoustic processing from the cloud, the implemented architecture demonstrates a practical approach for privacy-preserving, autonomous, and highly responsive human-computer interaction. As the demand for intelligent, conversational interfaces continues to expand into sensitive, real-world applications, the methodologies validated within this thesis provide a critical foundation for ensuring that these technologies remain secure, accessible, and naturally interactive directly at the network edge.

References

- Bahdanau, D., Cho, K., & Bengio, Y. (2015). Neural machine translation by jointly learning to align and translate. In *Proceedings of the 3rd International Conference on Learning Representations (ICLR 2015)*.
- Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D., Wu, J., Winter, C., ... Amodei, D. (2020). Language models are few-shot learners. *Advances in Neural Information Processing Systems*, 33, 1877–1901.
- Colby, K. M., Weber, S., & Hilf, F. D. (1971). Artificial paranoia. *Artificial Intelligence*, 2(1), 1–25.
- Dettmers, T., Lewis, M., Belkada, Y., & Zettlemoyer, L. (2022). LLM.int8(): 8-bit matrix multiplication for transformers at scale. *Advances in Neural Information Processing Systems*, 35, 30318–30332.
- Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT)* (pp. 4171–4186).
- Frantar, E., Ashkboos, S., Hoefler, T., & Alistarh, D. (2023). GPTQ: Accurate post-training quantization for generative pre-trained transformers. In *Proceedings of the 11th International Conference on Learning Representations (ICLR 2023)*.
- Gao, J., Galley, M., & Li, L. (2018). Neural approaches to conversational AI. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics: Tutorial Abstracts* (pp. 27–31).
- Hevner, A. R., March, S. T., Park, J., & Ram, S. (2004). Design science in information systems research. *MIS Quarterly*, 28(1), 75–105.
- Hinton, G., Deng, L., Yu, D., Dahl, G. E., Mohamed, A., Jaitly, N., Senior, A., Vanhoucke, V., Nguyen, P., Sainath, T. N., & Kingsbury, B. (2012). Deep neural networks for

- acoustic modeling in speech recognition: The shared views of four research groups. *IEEE Signal Processing Magazine*, 29(6), 82–97.
- Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural Computation*, 9(8), 1735–1780.
- Kim, J., Kong, J., & Son, J. (2021). Conditional variational autoencoder with adversarial learning for end-to-end text-to-speech. In *Proceedings of the 38th International Conference on Machine Learning* (pp. 5530–5540). PMLR.
- Lin, J., Tang, J., Tang, H., Yang, S., Dang, X., & Han, S. (2024). AWQ: Activation-aware weight quantization for on-device LLM compression and acceleration. In *Proceedings of the 7th MLSys Conference*.
- Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C. L., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., Schulman, J., Hilton, J., Kelton, F., Miller, L., Simens, M., Askell, A., Welinder, P., Christiano, P., Leike, J., & Lowe, R. (2022). Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35, 27730–27744.
- Peppers, K., Tuunanen, T., Rothenberger, M. A., & Chatterjee, S. (2007). A design science research methodology for information systems research. *Journal of Management Information Systems*, 24(3), 45–77.
- Radford, A., Kim, J. W., Xu, T., Brockman, G., McLeavey, C., & Sutskever, I. (2023). Robust speech recognition via large-scale weak supervision. In *Proceedings of the 40th International Conference on Machine Learning* (pp. 28492–28518). PMLR.
- Ramírez, J., Segura, J. C., Benítez, C., de la Torre, A., & Rubio, A. (2004). Efficient voice activity detection algorithms using long-term speech information. *Speech Communication*, 42(3–4), 271–287.
- Satyanarayanan, M. (2017). The emergence of edge computing. *Computer*, 50(1), 30–39.
- Shi, W., Cao, J., Zhang, Q., Li, Y., & Xu, L. (2016). Edge computing: Vision and challenges. *IEEE Internet of Things Journal*, 3(5), 637–646.
- Sutskever, I., Vinyals, O., & Le, Q. V. (2014). Sequence to sequence learning with neural networks. *Advances in Neural Information Processing Systems*, 27, 3104–3112.
- Turing, A. M. (1950). Computing machinery and intelligence. *Mind*, 59(236), 433–460.

- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., & Polosukhin, I. (2017). Attention is all you need. *Advances in Neural Information Processing Systems*, 30, 5998–6008.
- Wallace, R. S. (2009). The anatomy of A.L.I.C.E. In R. Epstein, G. Roberts, & G. Beber (Eds.), *Parsing the Turing test* (pp. 181–210). Springer.
- Wang, X., Han, Y., Wang, C., Zhao, Q., Chen, X., & Chen, M. (2020). Convergence of edge computing and deep learning: A comprehensive survey. *IEEE Communications Surveys & Tutorials*, 22(2), 869–904.
- Weizenbaum, J. (1966). ELIZA—A computer program for the study of natural language communication between man and machine. *Communications of the ACM*, 9(1), 36–45.
- Xiao, G., Lin, J., Seznec, M., Wu, H., Demouth, J., & Han, S. (2023). SmoothQuant: Accurate and efficient post-training quantization for large language models. In *Proceedings of the 40th International Conference on Machine Learning* (pp. 38087–38099). PMLR.
- Yang, A., Yang, B., Hui, B., Zheng, B., Yu, B., Zhou, C., Li, C., Li, C., Liu, D., Huang, F., Dong, G., Wei, H., Lin, H., Tang, J., Wang, J., Yang, J., Tu, J., Zhang, J., Ma, J., ... Fan, Z. (2024). Qwen2 technical report. *arXiv*. <https://arxiv.org/abs/2407.10671>
- Zhou, Z., Chen, X., Li, E., Zeng, L., Luo, K., & Zhang, J. (2019). Edge intelligence: Paving the last mile of artificial intelligence with edge computing. *Proceedings of the IEEE*, 107(8), 1738–1762.

Implementation / Software References

- Alibaba Cloud. (2024). *Qwen2.5-1.5B-Instruct* [Large language model]. Hugging Face.
<https://huggingface.co/Qwen/Qwen2.5-1.5B-Instruct>
- Coqui AI. (2023). *TTS* (Version 0.27.2) [Computer software]. GitHub.
<https://github.com/coqui-ai/TTS>
- Gerganov, G., & contributors. (2023). *llama.cpp* [Computer software]. GitHub.
<https://github.com/ggerganov/llama.cpp>
- Hugging Face. (2024). *Transformers* [Computer software]. GitHub.
<https://github.com/huggingface/transformers>
- OpenAI. (2025). *Whisper* [Computer software]. GitHub.
<https://github.com/openai/whisper>
- Python Software Foundation. (2025). *Python* (Version 3.x) [Computer software].
<https://www.python.org/>
- PyTorch Contributors. (2025). *PyTorch* [Computer software]. <https://pytorch.org/>
- Silero Team. (2024). *Silero VAD* [Computer software]. GitHub.
<https://github.com/snakers4/silero-vad>
- SoundDevice Contributors. (2025). *python-sounddevice* [Computer software]. GitHub.
<https://github.com/spatialaudio/python-sounddevice>