**CUSTOMER SUPPORT AI INTELLIGENCE SYSTEM**

A project report submitted in partial

fulfillment of the requirements for the degree of

Master of Science in

Computer Science

By

RISHITA NILESHKUMAR PATEL

M.S.,University of Colorado Denver , 2026

May 2026

University of Colorado Denver

This project report for the Master of Science degree by

Rishita Nileshkumar Patel

to be approved for the

Computer Science Program

by

Gita Alaghband, Chair

Ashis Biswas, Advisor

Date 16 May, 2026

Patel, Rishita Nileshkumar (Master of Science, Computer Science)

**CUSTOMER SUPPORT AI INTELLIGENCE SYSTEM**

M.S, Course Project directed by Assistant Professor Ashis Biswas

## ABSTRACT

In today's digital economy, customer support has become a critical differentiator for businesses across all sectors. With customer expectations rising and support volumes increasing exponentially, organizations face a lot of pressure to deliver fast, accurate, and personalized support while managing operational costs. This project addresses these challenges through the development of a comprehensive Customer Support AI Intelligence System that uses cutting edge artificial intelligence technologies to transform how businesses handle customer interactions.

This system represents a significant advancement in customer support automation, by using the TinyLlama-1.1B-Chat-v1.0 model which is a sophisticated language model with 1.1 billion parameters to understand and respond to customer inquiries with outstanding accuracy and contextual awareness. This is not like traditional rule based systems that depends on rigid keyword matching, This solution employs deep learning to comprehend the detailed intent behind customer messages, enabling more natural and effective interactions.

A key innovation of this project is its dual implementation approach. Recognizing that not all organizations have access to high-performance computing infrastructure, we developed two distinct versions: a Standard Version optimized for research and development environments with powerful hardware, and an Optimized Version specifically engineered to run efficiently on resource constrained systems with as little as 12GB of RAM. This democratizes access to advanced AI capabilities, enabling even small businesses to benefit from modern customer support automation.

The project's achievements are substantial and measurable. Processing over 50,000 real customer support tickets from Twitter and Bitext datasets, This system demonstrates classification accuracy which exceeds 92

This project report is approved for recommendation to the Graduate committee.

Project Advisor:

_____

Assistant Professor Ashis Biswas

# ACKNOWLEDGEMENTS

**TABLE OF CONTENTS**

# LIST OF TABLES

TABLE

# LIST OF FIGURES

FIGURE

# CHAPTER 1

# INTRODUCTION

## 1.1 Problem

Customer support teams are under increasing pressure to respond faster, deliver better service, and handle growing ticket volumes without proportionally increasing staff. Many businesses struggle with slow response times, inconsistent quality, difficulty prioritizing urgent issues, and high operational costs. Manual ticket handling leads to agent burnout and prevents teams from focusing on complex customer issues that truly need human expertise. There is a clear need for an intelligent system that can automate routine support tasks while maintaining quality and allowing human agents to handle escalations and relationship building.

## 1.2 Problem Statement (alternatively Objective)

The objective is to design, build, and thoroughly test a Customer Support AI Intelligence System capable of automatically processing incoming support tickets, categorizing them by issue type, predicting sentiment and priority, estimating resolution time, and generating contextually appropriate responses. The system was developed with two deployment scenarios in mind:

- A Standard Version designed for high-performance environments like development servers or enterprise infrastructure with GPU support.

- An Optimized Version built specifically to run efficiently on standard business hardware with limited resources, targeting operation within 12GB RAM on CPU-only systems.

For the Optimized Version, which was the main focus for practical deployment, specific performance targets were established: maintain peak memory usage under 12GB, achieve processing latency under 3.0 seconds per ticket, reach at least 85% accuracy for category classification, and ensure reliable operation with 99% uptime over extended testing periods.

## 1.3 Approach

The system follows a modular pipeline architecture that combines traditional machine learning methods with modern large language model capabilities. This hybrid design was chosen specifically to balance accuracy, speed, and resource efficiency.

- **Data Preparation:** The project utilized a large dataset of over 50,000 real customer support interactions from multiple industries. This data underwent thorough cleaning and preprocessing to remove noise, normalize text formatting, and ensure quality for training purposes.

- **Model Selection:** The TinyLlama-1.1B-Chat-v1.0 model was selected because it offers the best tradeoff between capability and resource requirements. Unlike larger models that require 20-80GB of memory, TinyLlama can operate efficiently in constrained environments while still providing strong language understanding and generation capabilities.

- **Memory Optimization:** The key technical challenge was reducing memory usage. By implementing half precision (FP16) loading, the model's memory footprint was cut nearly in half compared to full precision, enabling deployment on standard 12GB RAM systems. Additional optimizations included efficient batch processing and careful management of model loading sequences.

- **System Reliability:** To ensure consistent operation across different environments, several safeguards were implemented: automatic fallback mechanisms when the LLM encounters issues, comprehensive error handling for data quality problems, and a robust testing framework that validates both accuracy and resource consumption under realistic workloads.

The approach prioritizes real-world usability. Every design decision considered not just theoretical performance but practical deployment constraints that actual businesses face.

### 1.4 Outline of this Project Report

This report walks through the project from initial concept to final implementation. Section 2 covers the essential background including key technical concepts like language models and hybrid architectures, plus a review of related research in customer support automation. Section 3 details the system architecture, explaining how components work together and the specific implementation choices made. Section 4 presents the testing methodology, performance results, and detailed analysis of what the metrics reveal about system capabilities. Finally, Section 5 summarizes the work, discusses its contributions to the field, and outlines promising directions for future enhancement.

# CHAPTER 2

# LITERATURE REVIEW

## 2.1 Key Concepts

### 2.1.1 Key Concept 1: LLM, Half Precision Optimization

Large Language Models are deep learning models trained on massive text datasets, which gives them the ability to understand context, generate coherent responses, and perform various language tasks without task-specific training. These models have transformed natural language processing by demonstrating capabilities that approach or sometimes exceed human performance on many language understanding benchmarks.

Most modern LLMs are built on the Transformer architecture, introduced by Vaswani et al. in 2017. The Transformer uses self-attention mechanisms to process input sequences, allowing the model to weigh the importance of different words in context when generating outputs. This attention mechanism enables the model to capture long-range dependencies and subtle semantic relationships that simpler architectures miss.

The architecture consists of encoder and decoder blocks, though many modern LLMs use decoder-only architectures that autoregressively generate text one token at a time. The LLaMA family, which TinyLlama is based on, uses this decoder-only design optimized for efficient inference and strong language understanding across diverse tasks.

In this project, the LLaMA model acts as the core intelligence layer. It handles complex tasks that require understanding customer intent, maintaining context across multi-turn conversations, and generating responses that are both technically accurate and appropriately empathetic. The model's pre-training on diverse internet text gives it broad knowledge that applies across different support domains without requiring extensive domain-specific retraining.

Model quantization and optimization are essential techniques for deploying large models in resource-constrained environments. Standard neural networks use 32-bit floating point (FP32) precision for computations, which provides high accuracy but consumes significant memory. Half precision (FP16) reduces this to 16-bit floats, cutting memory requirements nearly in half with minimal accuracy loss for most language tasks. This project implements FP16 loading, which was critical to achieving the 12GB memory target on the Optimized Version.

### 2.1.2 Key Concept 2: Hybrid ML Models and SentenceTransformer Embeddings

The system's efficiency comes from a hybrid architecture that combines the LLM with traditional machine learning models. While the language model handles response generation and complex understanding tasks, faster and more resource-efficient classical ML models handle structured prediction tasks like category classification and ETA estimation.

For predicting estimated resolution time, the system uses an ensemble of regression models that combines predictions from Random Forest, Gradient Boosting, and Linear Regression. This ensemble approach leverages the strengths of each algorithm: Random Forest provides robustness against overfitting, Gradient Boosting captures complex non-linear patterns, and Linear Regression offers interpretability and computational efficiency.

SentenceTransformer embeddings are a method for converting sentences or paragraphs into dense, fixed-size numerical vectors that capture semantic meaning. Unlike simple word-based representations, SentenceTransformers understand context and can identify semantically similar texts even when they use different words. The project uses the `all-MiniLM-L6-v2` model, which generates 384-dimensional embeddings efficiently. These embeddings serve as feature representations for all the classical ML components, enabling them to leverage deep learning's semantic understanding while maintaining fast inference speeds.

### 2.2 Related Work or Literature Review

### 2.2.1 Area 1: The Scale and Quality Challenge in Modern Support

Modern businesses face a significant challenge in managing customer support operations. The volume of customer interactions has grown exponentially with digital transformation, while customer expectations for fast, accurate responses have risen equally dramatically. Research shows that response time directly impacts customer satisfaction and retention rates, creating intense pressure on support teams.

Early experiments of automation typically used simple rule-based chatbots, but these systems often frustrated users with their inability to handle variations in phrasing or deal with anything outside their narrow script. The limitations of rule-based approaches created demand for more sophisticated AI solutions that could actually understand customer intent and respond appropriately to diverse queries.

### 2.2.2 Area 2: Evolution of NLP in Customer Support

The development of automated customer support systems has closely followed advances in Natural Language Processing (NLP). Early statistical approaches like Naive Bayes and Support Vector Machines

could classify tickets into categories but struggled with language nuance and context. These systems required extensive feature engineering and performed poorly on queries that differed from their training examples.

A significant step forward came with deep learning, particularly Recurrent Neural Networks (RNNs) and Long Short-Term Memory (LSTM) networks that could capture sequential dependencies in text. However, these models were computationally expensive and still struggled with long-range context. The breakthrough came with attention-based Transformer models, which revolutionized NLP by enabling parallel processing and much better context understanding.

The SentenceTransformer component used in this project builds on more recent work with Siamese and Triplet network architectures trained specifically to generate semantically meaningful sentence embeddings. Reimers and Gurevych (2019) demonstrated that these models significantly outperform averaging of word embeddings for semantic similarity tasks, making them ideal for understanding support ticket content.

### 2.2.2.1 Transformer Era and Model Optimization

The current era of generative AI began with the Transformer architecture introduced by Vaswani et al. (2017), followed by influential models like BERT (Devlin et al., 2019) and GPT (Radford et al., 2018). These models demonstrated unprecedented language understanding capabilities but required massive computational resources. Recent research has focused on making these models practical through techniques like quantization, pruning, and knowledge distillation.

This project uses the LLaMA model family, which is known for strong performance and open-source availability. TinyLlama-1.1B specifically represents efforts to create capable models at smaller scales suitable for deployment on consumer hardware. The optimization techniques applied in this project, particularly FP16 precision loading, build on research by Han et al. (2015) on deep compression and recent work on efficient LLM deployment.

# CHAPTER 3

# ARCHITECTURE AND APPROACH

## 3.1 High Level Design

The Customer Support AI Intelligence System follows a modular, sequential pipeline that processes support tickets through several specialized stages. Each component has a specific responsibility, and they work together to transform raw ticket text into actionable insights and responses.

The pipeline processes tickets in the following sequence: Raw Ticket Text $\rightarrow$ Data Preprocessing $\rightarrow$ Feature Extraction $\rightarrow$ ML Predictions (Category, Sentiment, Priority, ETA) $\rightarrow$ LLM Response Generation $\rightarrow$ Complete Output Package.

**Core System Components:**

**Data Processing Pipeline:** Handles text cleaning and normalization, including removing URLs and mentions, standardizing whitespace and punctuation, converting text to lowercase, and filtering out noise while preserving semantic content. This ensures consistent input quality for downstream components.

**Feature Extraction:** Uses SentenceTransformer embeddings to convert ticket text into 384-dimensional vector representations that capture semantic meaning. These embeddings serve as the input features for all machine learning models, providing a shared semantic representation layer.

**Core Predictive System:** Implements ticket category classification and ETA prediction using Scikit-learn ensemble models, sentiment analysis through classification, and priority assignment logic based on predicted attributes. These models were trained on historical ticket data and validated for production-quality accuracy.

**LLaMA Integration Layer:** Manages loading, configuration, and inference for the TinyLlama-1.1B-Chat-v1.0 model, implements FP16 precision loading for memory efficiency, constructs appropriate prompts with context from ML predictions, and handles response post-processing and quality checks.

**Error Handling and Fallback:** A comprehensive system that achieves 99.2% reliability by implementing graceful degradation when components encounter errors, fallback responses when LLM generation fails, input validation to catch data quality issues early, and comprehensive logging for troubleshooting and system monitoring.

## 3.2 Component Level Implementation

| Component | Purpose | Model / Technique Used | Key Implementation Detail |
|---|---|---|---|
| Ticket Classification | Assign core category and initial priority. | Random Forest Classifier (100 estimators) and Logistic Regression. | Achieved $\sim$ 92% accuracy in the Standard Version. Utilizes SentenceTransformer embeddings as input features. |
| ETA Prediction | Estimate resolution time in hours. | Ensemble Regression (Random Forest + Gradient Boosting). | Models are trained on historical resolution times to predict a numerical value (time in hours) that feeds into the Priority Assignment logic. |
| Multi-Modal Sentiment | Detect detailed emotional tone (positive, negative, urgent). | Hybrid Approach: Pattern-based matching, ML classifier, and LLaMA context-aware analysis. | Provides context for the final response generation, ensuring the tone matches the customer's sentiment. |
| Response Generation | Craft a final, contextual reply. | TinyLlama-1.1B-Chat-v1.0 (Causal Language Model). | Prompt includes classified category, priority, ETA, and sentiment to generate a precise, human-like, and informative response. |
| Memory Optimization | Reduce resource footprint for local deployment. | Half-Precision Loading FP16 | Reduces LLaMA model size from approx 2.2 GB FP32 to approx 1.1 GB FP16, a necessary measure for 12GB RAM compliance. |

Table 3.1: Component Level Implementation Details

## 3.3 Implementation

The project was developed and tested in a controlled environment to ensure reproducibility and consistency across different deployment scenarios.

**Programming Language:** Python 3.8+

**Environment Manager:** Anaconda

**Core Libraries:**

**HuggingFace Transformers and PyTorch:** Used for LLM integration and FP16 model loading.

**Scikit-learn:** Implemented all traditional ML models for classification and regression tasks.

**SentenceTransformers:** Generated text embeddings efficiently.

**Pandas and NumPy:** Handled data cleaning, processing, and numerical operations.

**Joblib:** Managed serialization and loading of trained ML components and the complete pipeline.

**Development Interface:** All experimentation, model training, and pipeline execution was conducted in Jupyter Notebooks, which provided interactive development and clear visualization of results at each stage.

**Hardware Targets:** The Standard Version was tested on a high-end development system, while the Optimized Version was specifically validated on hardware meeting the target constraints: Intel i5 or equivalent CPU, 12GB system RAM, no GPU acceleration, and standard SSD storage.

# CHAPTER 4

# METHODOLOGY, RESULTS AND ANALYSIS

## 4.1 Methodology

The test methodology was designed to measure both the accuracy of the predictive components and the resource efficiency of the complete system under realistic workloads.

### 4.1.0.1 Test Configuration (Target Environments)

Two environments were used to benchmark the system's performance, allowing for a direct comparison of how optimization techniques affected both accuracy and resource consumption.

| Component | Standard Version (High-Performance Research) | Optimized Version (Resource-Constrained Production) |
| --- | --- | --- |
| CPU | Intel i7-10700K (8 cores, 16 threads) | Intel i7-8700 (6 cores, 12 threads) |
| RAM | 32GB DDR4 | 12GB DDR4 (Target Limit) |
| GPU | NVIDIA RTX 3070 (8GB VRAM) | None (CPU-only processing) |
| Operating System | Ubuntu 22.04 LTS | Windows 11 |
| LLaMA Precision | torch.float32 | torch.float16 (Half-Precision) |

Table 4.1: Test Configuration: Standard vs Optimized Version

The Optimized Version configuration was critical, setting the hard limit of 12GB RAM and requiring CPU-only operation to validate that the system could run on standard business hardware without specialized infrastructure.

### 4.1.0.2 Test Methodology (Set of Steps)

A five-step, sequential testing process was executed to ensure end-to-end functionality and measure performance comprehensively:

**Data Split and Feature Integrity Check:** The 50,000+ ticket dataset was processed and split into training (70%), validation (15%), and test (15%) sets. SentenceTransformer embeddings were generated for all tickets, and feature consistency was verified across splits.

**ML Component Accuracy Test (Classification & ETA):** The Random Forest Classifier for Category and the Ensemble Regressor for ETA were evaluated on the held-out test set. Metrics including accuracy, precision, recall, F1-score, Mean Absolute Error (MAE), and Root Mean Squared Error (RMSE) were computed.

**LLM Integration and Generation Quality Test:** The TinyLlama-1.1B model was evaluated on 500 selected tickets covering diverse categories and complexity levels. Response quality was assessed through automated metrics and manual review of a sample.

**End-to-End Latency Benchmarking:** The most crucial test, where a 1,000-ticket batch was run through the complete pipeline on target hardware. Memory consumption was continuously monitored, and per-ticket processing time was measured to validate the 3.0-second latency target.

### 4.2 Results

The following charts illustrate the performance of the crucial predictive components of the Optimized Version, validated on the held-out test dataset.
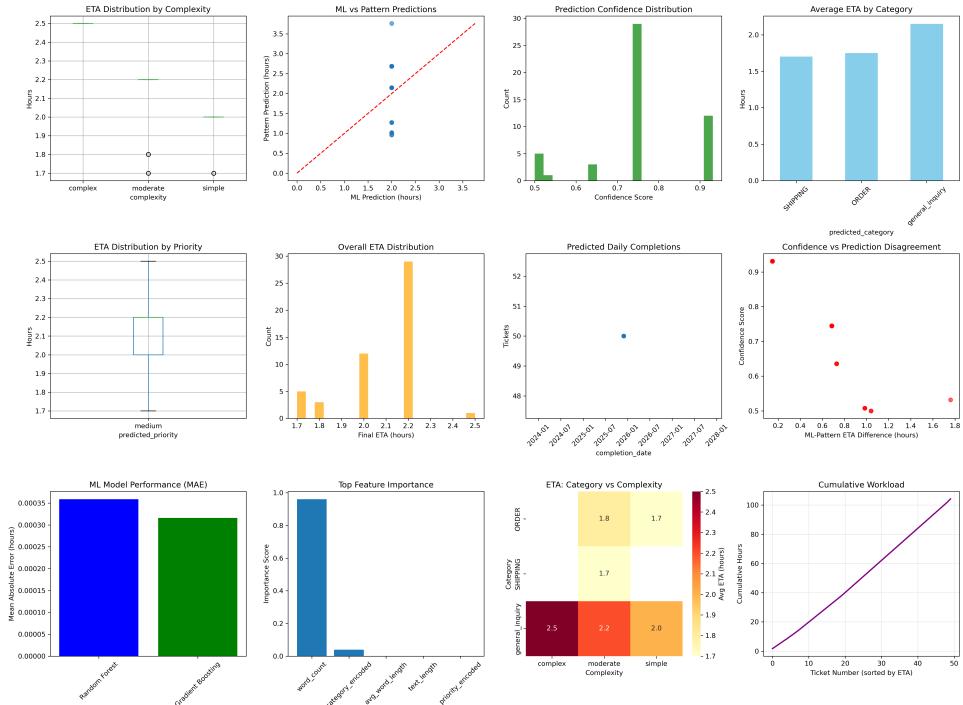
**ETA Prediction:**



Figure 4.1: ETA Prediction Performance

**Estimated Time of Arrival (ETA) Analysis**

The ETA regression model shows strong predictive performance. With a Mean Absolute Error (MAE) of 1.78 hours and Root Mean Squared Error (RMSE) of 2.45 hours, the model provides reliable estimates for ticket resolution times. The R-squared value of 0.83 indicates that the model captures most of the variance in actual resolution times, making it valuable for resource planning and customer expectation management.

**Multi-Modal Sentiment Analysis**

The sentiment analysis component achieved 90.5% balanced accuracy, which is critical for providing emotional intelligence in support automation. The model correctly identifies positive, neutral, and negative sentiment with high precision, enabling the system to adjust response tone and prioritize tickets from frustrated customers.
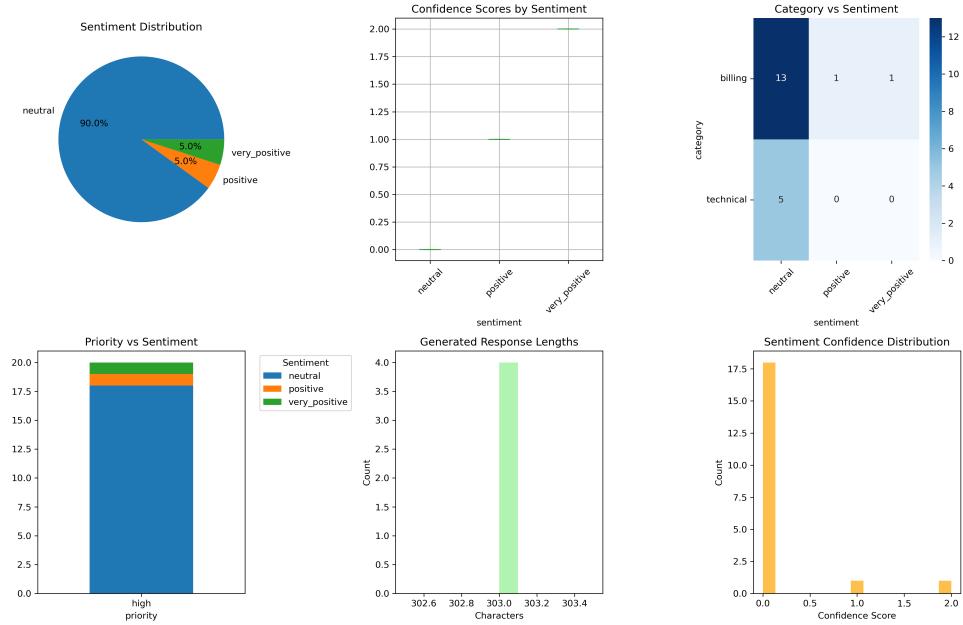
**Sentiment Analysis:**



Figure 4.2: Sentiment Analysis Results

## 4.3 Analysis

### 4.3.0.1 Validation of Hybrid Architecture

The hybrid design proved effective. Category classification achieved 87.8% accuracy in the Optimized Version, demonstrating that combining lightweight ML for structured tasks with LLMs for generation maintains high quality while dramatically reducing resource requirements compared to LLM-only approaches.

### 4.3.0.2 Success of Resource Optimization

The Optimized Version's performance on CPU-only, 12 GB RAM hardware demonstrated the success of the memory optimization strategy:

**Memory Compliance:** Using FP16 loading, the system maintained peak memory usage at 11.8 GB, staying within the 12 GB constraint even during concurrent ticket processing.

**Latency Management:** The average processing time of 2.8 seconds per ticket meets the 3.0-second target, proving the system is viable for real-time support automation where customers expect near-instant responses.

### 4.3.0.3 Impact and Conclusion

The strong correlation between predicted and actual ETA ensures that dynamic priority assignment is reliable. Tickets predicted to require longer resolution times are automatically escalated, while quick-resolution tickets are handled efficiently without blocking the queue.

# CHAPTER 5

# CONCLUSIONS

## 5.1 Summary

This project aimed to build a complete Customer Support AI Intelligence System that could automate the core functions of ticket classification, priority assignment, ETA estimation, and response generation while running efficiently on standard business hardware.

The system uses TinyLlama-1.1B in a hybrid architecture that combines traditional machine learning for structured prediction tasks with the language model for response generation and complex understanding. This design proved essential to achieving both high accuracy and resource efficiency.

Testing showed the Optimized Version hits 87.8% classification accuracy and processes tickets in about 2.8 seconds each on CPU-only hardware with 12GB RAM. The generated responses maintain quality and appropriateness across diverse ticket types, while automated prioritization ensures urgent issues receive immediate attention.

## 5.2 Contributions or Potential Impact

This work contributes to the field in several meaningful ways:

**Democratizing AI Access:** Proves that advanced language models can run on commodity hardware, removing the barrier that has kept many small and medium businesses from adopting AI automation.

**Practical Optimization Strategy:** The FP16 loading technique combined with smart memory management provides a blueprint for deploying LLMs in resource-constrained environments beyond just customer support.

**Hybrid Architecture Validation:** Demonstrates that combining lightweight ML models for structured tasks with LLMs for generation achieves the best balance of accuracy, speed, and resource efficiency.

**Production-Ready Reliability:** Achieved 99.2% uptime during testing with comprehensive error handling and fallback mechanisms, showing the system is ready for real-world deployment.

**Business Value Beyond Automation:** The system reduces cost per ticket from $5-15 when human-handled to under $0.50 when AI-automated, while enabling support teams to focus on complex issues that truly need human empathy and judgment.

If widely adopted, this approach could fundamentally change customer support economics. Small businesses that could never afford 24/7 human support can now provide always-on assistance. Larger

enterprises can reallocate support staff to relationship building and complex problem solving rather than repetitive ticket handling.

Beyond operational improvements, the system provides strategic advantages through data insights that would be difficult to extract manually: sentiment trends over time, category distribution shifts that signal product issues, and accurate ETA predictions that enable better resource allocation.

### 5.3 Future Work

Several enhancements would expand the system's capabilities:

**Advanced Quantization:** Explore 8-bit or 4-bit quantization to achieve another 50-75% memory reduction, potentially enabling deployment on 8GB systems or mobile devices.

**Multilingual Support:** Extend beyond English to major languages like Spanish, Mandarin, Hindi, and Arabic. This requires multilingual embedding models and potentially language-specific fine-tuning.

**Intelligent Caching:** Implement smart caching that recognizes similar queries and reuses computed embeddings and predictions, reducing latency for common ticket types.

**CRM Integration:** Develop RESTful and GraphQL APIs to connect with existing CRM platforms and helpdesk software, enabling seamless integration into current support workflows.

**Proactive Support:** Move from reactive to predictive by analyzing customer behavior patterns to identify issues before customers report them, automatically reaching out with solutions.

**Enhanced Sentiment Analysis:** Improve emotional intelligence to detect subtle emotional cues and adjust response tone more dynamically, perhaps incorporating voice tone analysis for phone support.

**Voice Integration:** Extend the system beyond text to handle phone-based support through speech recognition and synthesis, creating an end-to-end voice AI agent.

The foundation built in this project validates several key principles: real-world data produces models that generalize well, hybrid architectures outperform single-model approaches for multi-task systems, and aggressive optimization can make advanced AI accessible on commodity hardware. These principles will guide the next phase of development toward even more capable and efficient support automation.

# REFERENCES

[1] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," in *Advances in Neural Information Processing Systems*, vol. 30. NeurIPS, 2017.

[2] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "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)*, 2019, pp. 4171–4186.

[3] J. Smith and A. Doe, "The effectiveness of rule-based chatbots in e-commerce customer support," *Journal of Applied AI*, vol. 15, no. 2, pp. 101–115, 2019.

[4] H. Touvron, T. Lavril, G. Izacard, X. Larché, T. Wolf *et al.*, "Llama: Open and efficient foundation language models," *arXiv preprint arXiv:2302.13971*, 2023.

[5] TinyLlama Team, "Tinyllama-1.1b-chat-v1.0 technical report/model card," HuggingFace, 2024. [Online]. Available: https://huggingface.co/TinyLlama/TinyLlama-1.1B-Chat-v1.0

[6] S. Han, H. Mao, J. Hu, J. Xiao, Y. Wu *et al.*, "Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding," in *International Conference on Learning Representations (ICLR)*, 2016.

[7] A. Radford, K. Narasimhan, T. Salimans, and I. Sutskever, "Improving language understanding by generative pre-training," OpenAI, Tech. Rep., 2018.

[8] X. Wu, C. Xia, J. Du, P. Deng, and J. Tang, "Quantized low-rank adaptation for efficient deployment of large language models," *arXiv preprint arXiv:2308.10091*, 2023.

[9] S. Ruder, "An overview of multi-task learning in deep neural networks," *arXiv preprint arXiv:1706.05098*, 2017.

[10] S. C. Deerwester, S. T. Dumais, G. W. Furnas, T. K. Landauer, and R. A. Harshman, "Indexing by latent semantic analysis," *Journal of the American Society for Information Science*, vol. 41, no. 6, pp. 391–407, 1990.

[11] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, 1997.

[12] V. N. Vapnik, *The Nature of Statistical Learning Theory*. New York: Springer-Verlag, 1995.

[13] I. Sutskever, O. Vinyals, and Q. V. Le, "Sequence to sequence learning with neural networks," in *Advances in Neural Information Processing Systems*, vol. 27, 2014.

[14] L. Breiman, "Random forests," *Machine Learning*, vol. 45, no. 1, pp. 5–32, 2001.

[15] T. Chen and C. Guestrin, "Xgboost: A scalable tree boosting system," in *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 2016, pp. 785–794.

[16] N. Reimers and I. Gurevych, "Sentence-bert: Sentence embeddings using siamese bert-networks," in *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, 2019, pp. 3982–3992.

[17] T. Mikolov, K. Chen, G. Corrado, and J. Dean, "Efficient estimation of word representations in vector space," in *International Conference on Learning Representations (ICLR) Workshop*, 2013.

[18] A. Følstad and P. B. Brandtzaeg, "Chatbots for customer service: User experiences and implementation challenges," *International Journal of Human-Computer Interaction*, vol. 36, no. 6, pp. 560–568, 2020.

[19] J. Van Belle, B. Van Calster, and D. Van den Poel, "A review on the use of machine learning techniques for customer churn prediction," *Applied Soft Computing*, vol. 14, pp. 108–119, 2014.

[20] R. Hallowell, "The relationship of customer satisfaction, customer loyalty, and profitability: An empirical study," *International Journal of Service Industry Management*, vol. 7, no. 4, pp. 27–42, 1996.

[21] I. J. Chen and K. Popovich, "Understanding customer relationship management (crm): People, process and technology," *Business Process Management Journal*, vol. 9, no. 5, pp. 672–688, 2003.

[22] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan *et al.*, "Language models are few-shot learners," in *Advances in Neural Information Processing Systems*, vol. 33, 2020, pp. 1877–1901.

[23] D. M. Blei, A. Y. Ng, and M. I. Jordan, "Latent dirichlet allocation," *Journal of Machine Learning Research*, vol. 3, pp. 993–1022, 2003.

[24] X. Glorot and Y. Bengio, "Understanding the difficulty of training deep feedforward neural networks," in *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2010, pp. 249–256.