

# Developing an Automated IT Ticket Resolution System with Retrieval-Augmented Generation and Word Embeddings

Taeyang Kim*	Jeremy Mumford
Pattern	Pattern
taeyang.kim@pattern.com	jeremy.mumford@pattern.com
Andrew Marquez	Jacob Miller
Pattern	Pattern
andrew.marquez@pattern.com	jacob@pattern.com

December 2023

## Abstract

This research investigates the application of Retrieval-Augmented Generation (RAG) and word embeddings to automate the ticket resolution system within Pattern, a company specializing in supply chain management software. Leveraging a combination of word embedding techniques and advanced language models, our methodology integrates a vector database to store and retrieve relevant information, enhancing the efficiency and accuracy of ticket resolution. Real support tickets are employed for evaluation, and the results demonstrate that using RAG and word embeddings provides improved performance over traditional fine-tuning approaches. The study contributes to the automation of complex tasks in the supply chain domain and explores the challenges and opportunities inherent in integrating these advanced NLP technologies.

## 1 Introduction

Supply chain management serves as the backbone of contemporary business operations, demanding precision and efficiency to navigate the complexities of global commerce. Within this intricate framework, the resolution of support tickets emerges

---

\* First author.

as a critical component, directly influencing the overall agility and reliability of supply chain processes. This research delves into the realm of cutting-edge natural language processing (NLP) technologies to explore their transformative potential in automating the ticket resolution system of Pattern, a prominent player in the supply chain management software landscape.

Traditional approaches to automating ticket resolution often rely on fine-tuning large language models (LLMs) using question-answer pairs. However, this method has limitations in terms of accuracy and practicality [2]. Recent advancements suggest that integrating Retrieval-Augmented Generation (RAG) with word embeddings can significantly enhance performance. By establishing a maintained knowledge base stored in a vector database, relevant information can be retrieved and supplied to the LLMs to generate more accurate and contextually appropriate responses [5].

In this study, we leverage RAG and word embeddings to bridge the gap between human-language understanding and automated resolution processes. Our approach involves creating a "canon" of truth maintained by process owners, which the system uses to retrieve relevant information for ticket resolution. This integration aims to overcome the limitations of direct fine-tuning, offering a more robust and scalable solution for automating support ticket responses in the supply chain domain.

## **2 Related Works**

### **2.1 Limitations of Direct Fine-Tuning for Ticket Resolution**

Prior studies have indicated that directly fine-tuning language models with question-answer pairs may not yield sufficient accuracy for practical implementation [1, 9]. These approaches often lack the ability to generalize beyond the provided data and may not maintain consistency over time due to the static nature of the training data.

### **2.2 Retrieval-Augmented Generation (RAG)**

Retrieval-Augmented Generation combines traditional LLMs with a retrieval mechanism that sources relevant documents from a database. This method has been shown to improve the factual accuracy of generated responses by grounding them in a dynamic knowledge base [2].

### **2.3 Word Embeddings and Vector Databases**

Word embeddings like Doc2Vec provide a way to represent textual data in a numerical format that preserves semantic relationships [4, 5]. When stored in a vector

database, these embeddings enable efficient retrieval of relevant documents based on similarity measures.

## 2.4 Applications in Supply Chain Management

The integration of NLP techniques in supply chain management has gained attention for its potential to automate and enhance various processes [3]. However, the application of RAG and word embeddings specifically for ticket resolution remains an emerging area of study.

# 3 Methodology

## 3.1 Overview

Our methodology involves integrating Retrieval-Augmented Generation with word embeddings to automate IT ticket resolution. This approach utilizes a vector database to store embeddings of existing ticket resolutions, allowing the system to retrieve relevant information dynamically.

## 3.2 Data Preparation

We curated a dataset comprising real support tickets and their corresponding resolutions. Each ticket and resolution pair was transformed into vector representations using Doc2Vec [5]. This dataset forms the knowledge base for the retrieval component.

## 3.3 Vector Database Implementation

All vector representations were stored in a vector database implemented using FAISS. This enables efficient similarity search to retrieve the most relevant documents based on the input ticket.

## 3.4 Retrieval-Augmented Generation

When a new ticket is submitted, the system performs the following steps:

1. **Embedding Generation:** The new ticket is converted into a vector representation using the same embedding model.
2. **Similarity Search:** The system queries the vector database to retrieve the top  $k$  most similar ticket-resolution pairs.

3. **Context Formation:** The retrieved resolutions are compiled into a context that is supplied to the language model.
4. **Response Generation:** The language model, such as GPT-3.5 Turbo, uses this context to generate a response to the ticket.

### 3.5 Language Model Configuration

We utilized GPT-3.5 Turbo via the OpenAI API for response generation. The model parameters were configured to accept the retrieved context within its input prompt, thereby grounding its response in the relevant information.

## 4 Evaluation

### 4.1 Experimental Setup

We conducted experiments to evaluate the performance of the RAG-based approach compared to the traditional fine-tuning method. The evaluation metrics included accuracy, response quality, computational cost, and time efficiency.

### 4.2 Results

Method	Accuracy (%)	Cost (\$)	Time (s)
Fine-Tuning	57.14	0.45	240
RAG + Embeddings	<b>75.6</b>	0.2	<b>1</b>

Table 1: Performance comparison between fine-tuning and RAG with embeddings

### 4.3 Analysis

The RAG approach achieved an accuracy of 75.6%, significantly higher than the fine-tuning method. The cost and time efficiency were also improved, demonstrating the benefits of utilizing a retrieval mechanism with embeddings.

### 4.4 Discussion

The RAG method effectively addresses the limitations of direct fine-tuning by dynamically incorporating relevant information into the response generation process.

This results in higher accuracy and more contextually appropriate responses, making it more suitable for practical implementation in a corporate environment.

## 5 Conclusion

This study demonstrates that integrating Retrieval-Augmented Generation with word embeddings significantly improves the performance of automated IT ticket resolution systems. By maintaining a dynamic knowledge base and leveraging efficient retrieval mechanisms, the system can generate accurate and contextually relevant responses. This approach overcomes the limitations of traditional fine-tuning methods and offers a scalable solution for automating support processes in supply chain management.

## Acknowledgments

We would like to thank Andrew Marquez for leading the technical efforts in implementing the RAG infrastructure within our project. We also acknowledge Jeremy Mumford and Jacob Miller for their valuable contributions to the development and maintenance of the vector database and embedding models.

## References

- [1] Smith, J., & others. (2019). Automating Supply Chain Management: A Rule-Based Approach. *Journal of Supply Chain Management*, 55(2), 35–50.
- [2] Zhang, Y., & Chen, X. (2020). Enhancing Ticket Resolution with AI: A Case Study. *International Journal of Logistics Management*, 31(1), 120–138.
- [3] Smith, L., & Jones, M. (2018). Domain Adaptation of Transformer Models in Supply Chain Analytics. *Journal of Business Analytics*, 2(1), 22–37.
- [4] Wang, H., & others. (2017). Doc2Vec for Document Representation in Classification Tasks. In *Proceedings of the International Conference on Data Mining and Big Data*, 456–465.
- [5] Li, X., & Zhang, Y. (2019). Leveraging Doc2Vec for NLP Applications. *Journal of Artificial Intelligence Research*, 64, 789–806.
- [6] Brown, T. B., & others. (2021). GPT-3: Advancements and Applications in Natural Language Processing. *AI Magazine*, 42(2), 58–67.

- [7] Chen, L., & others. (2022). Utilizing OpenAI API in Business Analytics. *Journal of Applied Business Research*, 38(4), 145–154.
- [8] Garcia, R., & Kim, T. (2018). Evaluating NLP Systems in Real-World Business Cases. *Journal of Information Technology*, 33(2), 150–163.
- [9] Patel, S., & others. (2020). Assessing the Accuracy of Automated Solutions in IT Ticketing. *Journal of Business and IT*, 15(3), 244–260.