# Problem Definition

## 1. Introduction and Background

In modern customer service operations, support teams are inundated with large volumes of incoming emails—ranging from product inquiries and technical issues to feedback and general requests. Manual handling of each email imposes significant labor costs, introduces response delays, and increases the risk of inconsistent or unhelpful replies. At the same time, prompt, accurate, and personalized responses directly impact customer satisfaction, brand reputation, and operational efficiency.

## 2. Real-World Relevance and Motivation

- **Operational Scalability**: As businesses grow, support email volume often scales faster than staffing. Hiring and training new agents is time-consuming and costly. Automated systems can help absorb surges and maintain service levels without linear headcount increases.
- **Customer Expectations**: Research shows that 78% of consumers expect a response from customer support within 24 hours, and 32% expect a response within one hour. Failing to meet these expectations leads to reduced customer loyalty and lost revenue.
- **Consistency and Quality**: Even experienced agents vary in tone and content. Automated, model-driven replies ensure a consistent brand voice while embedding best practices (e.g., politeness, clarity, next steps).

## 3. Problem Statement

**Objective**: Develop an end-to-end email auto-reply system that: 1) classifies incoming support emails into business-relevant categories, and 2) generates concise, accurate, and contextually appropriate reply drafts.

**Key Requirements**:

1. **Classification Accuracy**: Automatically assign each email to one of a finite set of categories (e.g., inquiry, issue, suggestion, other) with high precision and recall, to drive downstream templates or specialized handling.
2. **Reply Generation Quality**: Produce reply drafts that address the customer's question or issue, using professional tone and relevant information, minimizing manual editing.
3. **Operational Efficiency**: Integrate into existing workflows via a command-line interface (CLI) and RESTful API, enabling rapid iteration and deployment without extensive engineering overhead.

## 4. Limitations of Existing Solutions

- **Rule-Based Autoresponders**: Traditional autoresponder systems rely on keyword matching and static templates. They often misclassify nuanced requests, fail to handle edge cases, and produce generic, unhelpful responses.
- **Vanilla Prompt-Engineering**: Using a large language model (LLM) out-of-the-box with ad-hoc prompts can yield impressive text, but lacks domain adaptation. Without fine-tuning, responses may ignore company policies, misuse terminology, or hallucinate unsupported details.
- **SaaS Email Bots**: Commercial email automation platforms (e.g., Zendesk triggers, Intercom bots) provide template chaining and simple ML classification, but they typically require significant manual rule-crafting, lack customization for specific products, and incur ongoing licensing costs.

## 5. Proposed Approach and Scope

I propose a two-stage, fine-tuning-based solution:

1. **Stage 1 – Classification**: Fine-tune a DistilBERT model on a labeled dataset of customer emails. This lightweight transformer yields fast inference and high accuracy on support-specific categories.
2. **Stage 2 – Reply Generation**: Fine-tune a T5-style sequence-to-sequence transformer on paired (email, actual reply) data from historical tickets. We incorporate an explicit task prefix ("generate reply:") to improve the model's understanding of the generation objective.

Both stages are orchestrated via a Python CLI and exposed through a FastAPI service. The system reads training CSVs, persists models in a standardized directory layout, and exposes /classify and /reply endpoints for integration.

## 6. Why This Matters

- **Reduced Response Time**: Automated drafts cut average first-response time by 50–80%, enabling agents to focus on complex cases.
- **Cost Savings**: By deflecting repetitive inquiries, businesses can save on hourly support costs and reduce turnover due to burnout.
- **Data-Driven Improvement**: Metrics from classification accuracy and generation quality feed back into continuous model retraining, driving ongoing performance gains.