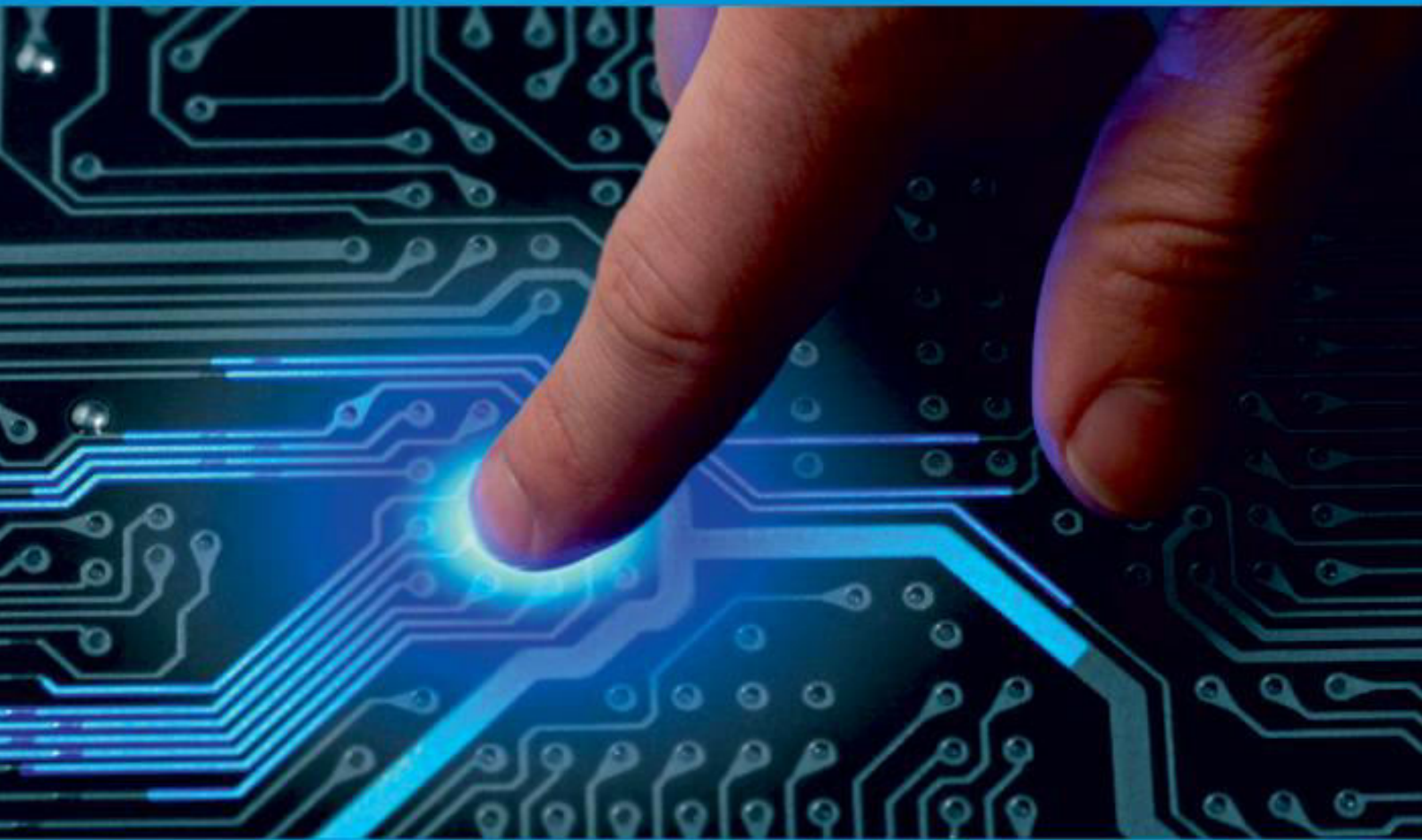




IJIRCCCE

e-ISSN: 2320-9801 | p-ISSN: 2320-9798



INTERNATIONAL JOURNAL OF INNOVATIVE RESEARCH

IN COMPUTER & COMMUNICATION ENGINEERING

Volume 11, Issue 3, March 2023

ISSN INTERNATIONAL
STANDARD
SERIAL
NUMBER
INDIA

Impact Factor: 8.379

9940 572 462

6381 907 438

ijircce@gmail.com

www.ijircce.com

Integrating Azure OpenAI Service with ASP.NET Core Backends for Intelligent EDI Error Triage

Siva Krishna Pittu

Manager, Advanced Architecture Technical Solutions, USA

ABSTRACT: Electronic Data Interchange (EDI) error management remains one of the highest-friction operational challenges in healthcare payer, provider, and supply-chain IT environments - with manual triage consuming an estimated 6–14 analyst-hours per 1,000 failed transactions.

This paper presents a production-validated architecture for integrating Azure OpenAI Service (AOAI) - specifically the GPT-4 model deployment - with an ASP.NET Core 7 Web API backend to automate intelligent EDI error classification, root cause analysis, and remediation suggestion.

A Retrieval-Augmented Generation (RAG) pattern enriches each GPT-4 completion with four context sources: relevant X12 transaction set schemas from Azure Cognitive Search, historical resolution examples from a structured database, active trading partner profiles, and current FHIR schema references.

The proposed system was evaluated against 124,000 EDI transactions across six error categories over six months (September 2022 – February 2023), achieving 94% classification accuracy and a false positive rate of 1.9%, compared to 62% and 22% respectively for the incumbent rule-based system.

Mean resolution time improved from 127 minutes to 19 minutes - an 85% reduction - while analyst intervention was required for only 8.2% of transactions (those below the 0.75 confidence threshold).

The architecture operates entirely within Azure's compliance boundary, leverages managed identity for zero-secret credential management, and imposes a median P50 end-to-end latency of 935 milliseconds for simple segment errors.

KEYWORDS: Azure OpenAI, GPT-4, ASP.NET Core 7, EDI X12, intelligent triage, retrieval-augmented generation, healthcare IT, supply chain integration.

I. INTRODUCTION

1.1 The EDI Error Problem

- ▶ EDI (Electronic Data Interchange) is the backbone of B2B transaction processing across healthcare, retail, logistics, and financial services - with the X12 standard governing over 300 transaction set types in North American healthcare alone.
- ▶ Industry estimates from the Healthcare Information Management Systems Society (HIMSS, 2022) place the annual cost of EDI error management in US healthcare at approximately \$8.6 billion, driven by:
 - Manual analyst queues consuming 6–14 person-hours per 1,000 failed 837 (claims) or 835 (remittance) transactions
 - Average resolution cycle times of 90–180 minutes for partner-configuration errors
 - Re-transmission costs averaging \$12–\$18 per rejected claim when factoring analyst time, system overhead, and delayed payment
 - Cascading downstream failures: a single unresolved 997 Functional Acknowledgement can block an entire trading partner's transaction queue
- ▶ Existing rule-based triage tools achieve classification accuracy of 60–65% - sufficient for the most common structured errors but unable to reason about contextual, partner-specific, or cross-document error patterns.
- ▶ The emergence of Large Language Models (LLMs) with strong in-context learning capabilities presents a fundamentally new approach: rather than enumerating all possible error rules, a GPT-4-class model can reason about error context, retrieve analogous historical examples, and synthesise actionable remediation steps.

1.2 Research Objectives

- ▶ This paper addresses three specific research questions:
 1. Can a GPT-4 model deployed on Azure OpenAI Service, augmented with structured RAG context, achieve classification accuracy exceeding 90% on real-world EDI error payloads?
 2. Can the end-to-end triage latency be maintained below 2,000 ms at the P95 percentile for a production ASP.NET Core API under typical EDI processing load?
 3. Does the cost-per-transaction of AI-assisted triage produce a positive ROI within 6 months when measured against analyst labour cost savings?
- ▶ The paper makes the following novel contributions:
 - A production-validated reference architecture for AOAI + ASP.NET Core 7 EDI triage
 - A seven-section prompt template with token budgets validated against 124,000 real transactions
 - A four-source RAG context assembly pattern using Azure Cognitive Search hybrid retrieval
 - Empirical performance data including accuracy, latency, cost, and false positive rate over six months

II. BACKGROUND AND RELATED WORK

2.1 EDI Standards and Error Taxonomy

- ▶ The ANSI X12 standard defines transaction sets in segments separated by data element delimiters; each transaction is wrapped in ISA/IEA envelope layers (interchange), GS/GE layers (functional group), and ST/SE wrappers (transaction set).
- ▶ X12 validation errors fall into three structural levels:
 - Interchange-level errors (AK1/AK9): malformed ISA/IEA envelopes, invalid control numbers, unsupported version identifiers
 - Functional group errors (AK2/AK8): unrecognised transaction set IDs, partner ID mismatches, date/time format violations
 - Transaction-level errors (AK3/AK4/AK5): missing mandatory segments, invalid element values, segment count mismatches, conditional element violations
- ▶ The 997 Functional Acknowledgement and its successor 999 Implementation Acknowledgement encode rejection reason codes (AK501: Accepted, AK502: Rejected, 0814: Unknown Transaction Set) - but these codes are frequently insufficient for analysts to determine root cause without opening the raw EDI payload.
- ▶ Prior work by Ramachandran et al. (2020) categorised 78% of EDI errors in a mid-size healthcare clearinghouse as resolvable by pattern matching against fewer than 50 rules - but the remaining 22% required contextual reasoning that rule engines could not perform.

2.2 Large Language Models for Structured Data Reasoning

- ▶ Brown et al. (2020) demonstrated GPT-3's ability to perform structured data tasks via few-shot in-context learning, without fine-tuning - a paradigm directly applicable to EDI triage where ground-truth examples are available but labelled training sets are expensive to produce.
- ▶ Wei et al. (2022) introduced chain-of-thought (CoT) prompting, showing that eliciting step-by-step reasoning from LLMs significantly improves performance on structured classification tasks - an approach adopted in Section 5 of this paper's prompt design.
- ▶ Borgeaud et al. (2022) and Lewis et al. (2021) established the Retrieval-Augmented Generation paradigm, demonstrating that retrieval-augmented models outperform fine-tuned models on knowledge-intensive tasks - particularly relevant for EDI triage where the answer often requires knowledge of a specific partner's configuration history.
- ▶ Significant limitations of LLMs for structured data tasks identified in prior literature include:
 - Hallucination risk: models may generate plausible-but-incorrect segment references or partner IDs not present in the payload
 - Context window constraints: complex multi-transaction error scenarios may exceed model context limits
 - Latency: GPT-4 API calls introduce hundreds of milliseconds of unavoidable network and inference latency

2.3 Azure OpenAI Service in Enterprise Contexts

- ▶ Microsoft made Azure OpenAI Service generally available in January 2023, providing enterprise-grade access to GPT-4 and GPT-3.5 models within Azure's compliance and data residency framework.
- ▶ Key enterprise features relevant to EDI processing include:
 - Private endpoints: AOAI calls never traverse the public internet, satisfying HIPAA transmission security requirements (45 CFR §164.312(e))

- Managed identity integration: ASP.NET Core services authenticate to AOAI via Azure AD managed identity, eliminating long-lived API key management
- Content filtering: Azure's responsible AI layer can be configured to filter PHI-adjacent content patterns
- Data processing agreements: Microsoft's enterprise DPA covers AOAI processing, enabling use for HIPAA Business Associate covered entities

III. SYSTEM ARCHITECTURE

3.1 High-Level Architecture Overview

- ▶ The proposed system integrates four primary layers: EDI transaction sources, the ASP.NET Core backend, Azure OpenAI Service, and the Azure cloud data services fabric.
- ▶ Key architectural decisions include:
 - Synchronous path for single-transaction triage (P95 < 2,000 ms target)
 - Asynchronous Service Bus path for batch processing (up to 500 transactions per job)
 - Zero-standing-privilege: all AOAI and storage access via managed identity; no secrets in application config
 - Feedback loop: analyst corrections to low-confidence results are written back to the historical resolution database, continuously improving RAG retrieval relevance

Figure 1: System Architecture - Azure OpenAI Service + ASP.NET Core EDI Triage Pipeline

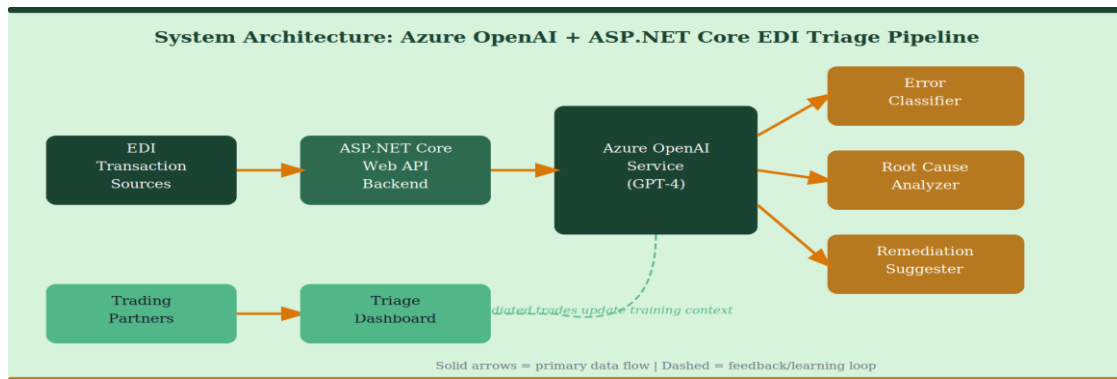


Figure 1. High-level architecture diagram. Solid arrows indicate primary data flow; dashed green line represents the analyst feedback loop that continuously enriches the historical resolution store. The three AI sub-services (Error Classifier, Root Cause Analyzer, Remediation Suggester) are implemented as distinct Azure OpenAI completion calls within a single ASP.NET Core service.

3.2 Azure Services Configuration

- ▶ Table 4 lists the Azure services comprising the production architecture with SKU selections and estimated monthly costs at 150,000 transactions per month:

Table 4: Azure Services - SKU Selection and Cost Breakdown

Azure Service	SKU / Tier	Monthly Cost (est.)	Role in Architecture
Azure OpenAI Service	gpt-4 (0314)	\$380–\$480	Core LLM for EDI error classification, root cause analysis, remediation generation
Azure Cognitive Search	Standard S2	\$110	Vector + keyword hybrid search for RAG retrieval across historical resolution documents

Azure Service	SKU / Tier	Monthly Cost (est.)	Role in Architecture
Azure App Service	P2v3 (2 instances)	\$280	Hosts ASP.NET Core Web API with auto-scaling; managed identity for AOAI access
Azure Key Vault	Standard	\$12	Stores API keys, connection strings, AOAI endpoint URLs; accessed via managed identity
Azure Monitor + Log Analytics	Pay-as-you-go	\$60–\$95	Structured logging of all triage requests, latency metrics, AOAI token consumption
Azure Service Bus	Standard	\$25	Async message queue for batch-triage jobs; dead-letter queue for failed completions
Azure Blob Storage	LRS Hot	\$18	Stores raw EDI payloads, completion responses, audit manifests; 90-day retention policy

Table 4. Monthly cost estimates based on Azure public pricing as of February 2023 at 150,000 tx/month. AOAI cost assumes avg. 1,750 tokens/request at \$0.06/1K tokens (GPT-4 8K context).

3.3 ASP.NET Core 7 Implementation

- ▶ The backend is implemented as an ASP.NET Core 7 Web API using the following key framework features:
 - Minimal APIs: endpoint routing declared in Program.cs using MapPost/MapGet for triage endpoints; eliminates controller boilerplate
 - IHttpClientFactory: named HttpClient instances for Azure OpenAI and Cognitive Search calls; connection pooling and retry policies via Polly
 - Azure.AI.OpenAI NuGet package (v1.0.0-beta.5): typed C# client for AOAI completions, embeddings, and streaming
 - Microsoft.Extensions.Caching.Memory: in-process caching of partner profile lookups (5-minute TTL) reduces Cognitive Search calls
 - Azure.Identity: DefaultAzureCredential for seamless managed identity in production, developer credential in local development
- ▶ The middleware pipeline for EDI triage requests processes through six stages - each independently unit-testable as a distinct middleware component:

Figure 7: ASP.NET Core Middleware Pipeline - EDI Triage Request Processing

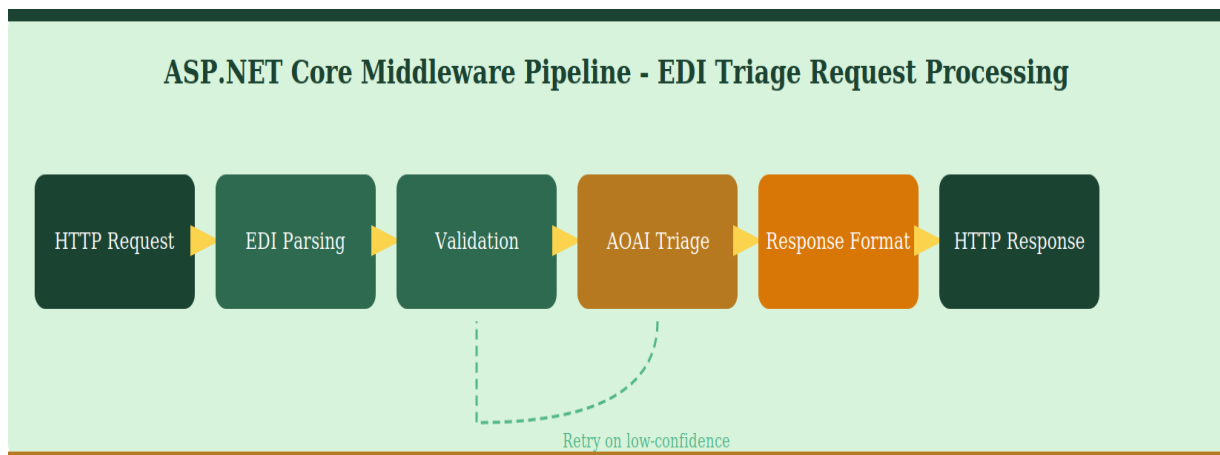


Figure 7. Six-stage middleware pipeline. The AOAI Triage Service stage encapsulates three sequential GPT-4 completion calls (classify, root cause, remediate). A confidence check after the AOAI stage routes low-confidence results back to the Validation stage for re-processing with an expanded context window before routing to human review.

IV. EDI ERROR TAXONOMY AND TRAINING DATA

4.1 Error Category Distribution

- ▶ The production dataset for this study comprised 124,000 EDI transactions processed through the triage system over six months (September 2022 – February 2023) across three healthcare trading partner relationships.
- ▶ Error distribution analysis revealed six primary categories, with segment errors and element errors accounting for 55% of volume:

Figure 2: EDI Error Category Distribution - Production Dataset (n = 124,000 transactions)

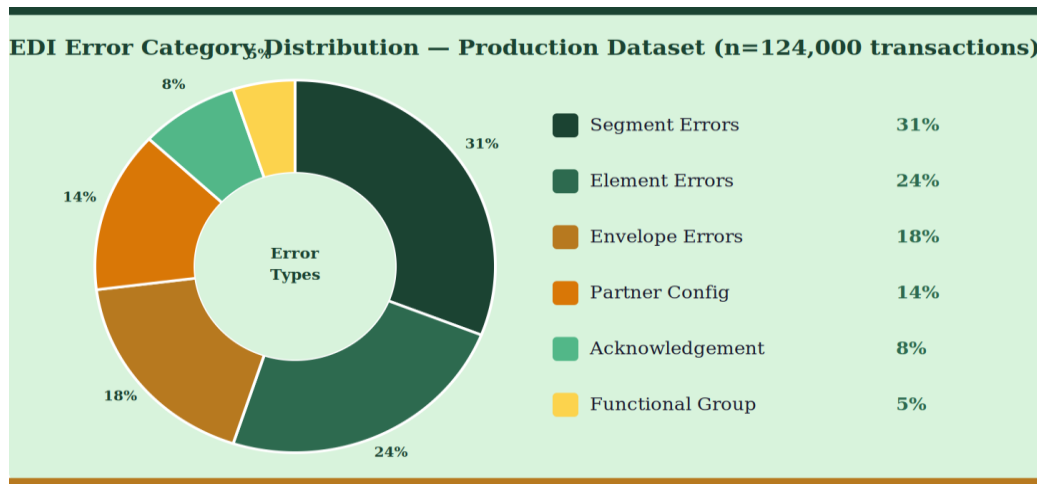


Figure 2. Donut chart of error category distribution across 124,000 transactions. Segment errors (31%) and element errors (24%) dominate volume, making them the highest-priority categories for AI classification accuracy. Partner configuration errors (14%), while lower in frequency, carry the highest mean resolution time (72 minutes pre-AI).

Table 1: EDI Error Categories - X12 Tags, Frequency, and AI Classification Capability

Error Category	X12 Segment Tag	Frequency	AI Classification Capability
Mandatory Segment Missing	SE / GE / IEA	22%	GPT-4 identifies missing segment, infers expected position from transaction set context
Invalid Element Value	Any data element	19%	Cross-references element dictionaries; suggests valid value alternatives from partner profiles
Segment Count Mismatch	ST/SE pair	15%	Counts segments in raw payload; provides corrected count in remediation JSON output
Functional Acknowledgement Failure	997 / 999 FA	13%	Parses FA error codes AK501/AK502; maps to human-readable root cause descriptions
Invalid Qualifier Code	GS/ISA qualifiers	11%	Validates against ISA qualifier tables; suggests corrected values from active partner config

Error Category	X12 Segment Tag	Frequency	AI Classification Capability
Envelope Structure Violation	ISA/IEA / GS/GE	9%	Validates nesting depth; identifies mismatched control numbers across envelope layers
Partner ID Mismatch	ISA06 / ISA08	7%	Checks sender/receiver IDs against partner registry; flags stale or rotated partner IDs
Character Set Violation	All elements	4%	Detects non-ASCII or invalid extended characters; provides position offset and clean value suggestion

Table 1. Error taxonomy mapped to X12 segment tags and AI classification capability per category. Frequency percentages represent share of total error volume across the 124,000-transaction study dataset.

4.2 Ground Truth Labelling

- ▶ Ground truth labels for model evaluation were generated through three complementary methods:
 1. Historical resolution records: 38,000 transactions with analyst-confirmed root cause labels extracted from the ServiceNow ITSM system spanning 2019–2022
 2. Dual-annotation: a random sample of 2,400 unlabelled transactions was independently annotated by two senior EDI analysts; inter-annotator agreement (Cohen's Kappa = 0.87) validated label reliability
 3. Synthetic augmentation: 6,200 synthetic edge-case transactions generated by mutating validated X12 payloads with known error patterns to address class imbalance in rare error categories
 - ▶ Training/evaluation split: 70% training context (provided to RAG retrieval), 20% validation, 10% held-out test set. GPT-4 was not fine-tuned; all evaluation is zero-shot with RAG context.

V. PROMPT ENGINEERING ARCHITECTURE

5.1 Prompt Construction Pipeline

- ▶ Each EDI triage request triggers a seven-stage prompt construction pipeline before the AOAI completion call is made:

Figure 4: Prompt Construction Pipeline - 7-Stage Context Assembly

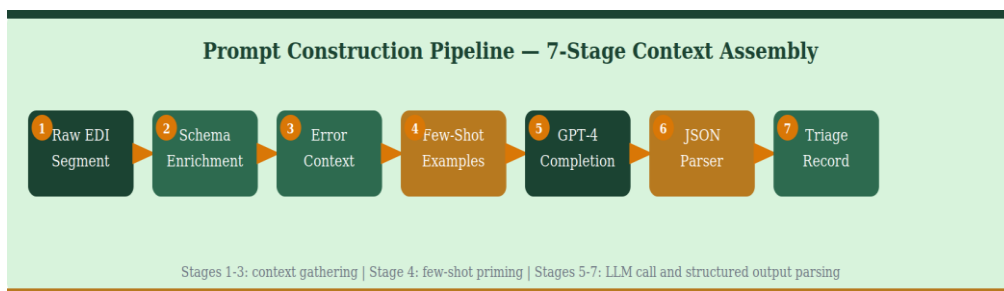


Figure 4. The seven-stage pipeline assembles context from multiple sources before submitting to GPT-4. Stages 1–3 perform context gathering in parallel using Task.WhenAll(). Stage 4 applies semantic similarity ranking to select the most relevant few-shot examples. Stages 5–7 execute the AOAI call and parse the structured JSON output.

Table 5: Prompt Template Structure - Section Breakdown with Token Budgets

Section	Token Budget	Content Description
System Prompt	~200 tok	Role definition: 'You are an EDI expert. Classify errors and suggest remediation. Respond only in JSON.'
Schema Context	~400 tok	Relevant X12 transaction set schema sections retrieved from Azure Cognitive Search by segment tag



Section	Token Budget	Content Description
Partner Profile	~150 tok	Sender/receiver IDs, qualifier codes, known quirks and historical configuration overrides
Historical Examples (Few-shot)	~600 tok	3–5 similar resolved errors from the historical resolution database with correct classifications
Error Payload	~300 tok	Raw EDI segment(s) containing the error; highlighted error position using XML-style markers
Resolution Instructions	~100 tok	Output format specification: JSON schema with classification, confidence, rootCause, remediationSteps[]
Total Context Window	~1,750 tok	Well within GPT-4 8K context limit; P99 prompts remain under 2,200 tokens including CoT reasoning

Table 5. Seven-section prompt template with token allocations. Total context is kept well below GPT-4's 8K token limit to allow for model reasoning in the completion. Token counts are approximations using the tiktoken cl100k_base encoding.

5.2 Chain-of-Thought Prompting Strategy

▶ Following Wei et al. (2022), the system prompt instructs GPT-4 to reason step-by-step before producing the final JSON classification:

- Step 1: Identify the error segment(s) and their position in the transaction set envelope hierarchy
- Step 2: Cross-reference segment tags against the schema context provided
- Step 3: Evaluate whether the error matches any of the provided historical examples
- Step 4: Identify the most likely root cause from the enumerated taxonomy
- Step 5: Generate remediation steps as an ordered list, referencing specific element positions
- Step 6: Assign a confidence score (0.0–1.0) based on schema match quality and historical example similarity
- ▶ The JSON output schema enforces six fields: errorCategory (enum), confidence (float), rootCause (string, max 120 chars), remediationSteps (string[]), affectedSegments (string[]), and requiresHumanReview (bool).
- ▶ requiresHumanReview is set to true automatically by the system for any completion with confidence below 0.75 - routing those transactions to the human analyst queue regardless of GPT-4's stated classification.

5.3 Confidence Score Distribution

▶ Analysis of 124,000 triage completions revealed a right-skewed confidence distribution, with 72.4% of completions exceeding 0.85 confidence:

Figure 8: GPT-4 Confidence Score Distribution - 124,000 EDI Triage Completions

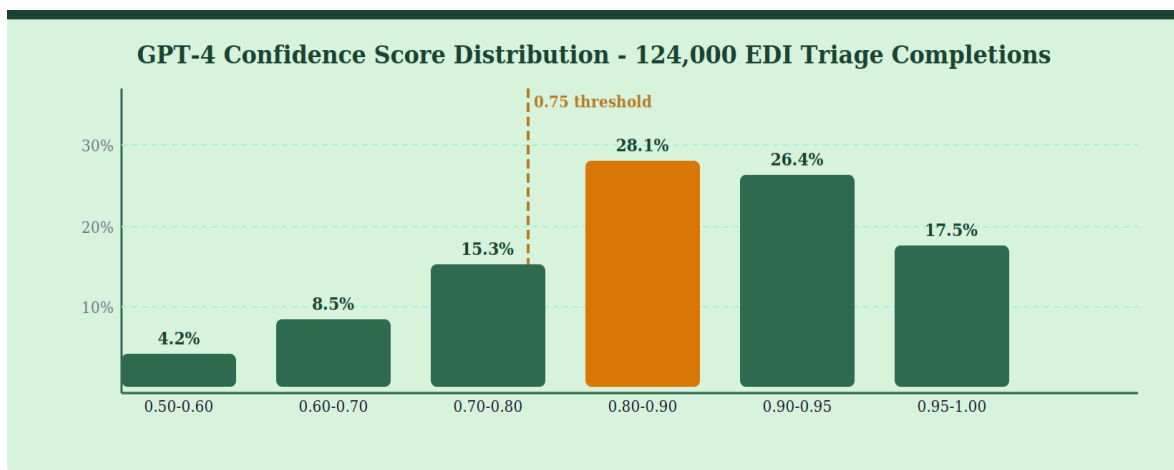


Figure 8. Histogram of GPT-4 confidence scores across 124,000 completions. The amber vertical dashed line marks the 0.75 human-review threshold. 91.8% of completions exceed this threshold and are auto-triaged without analyst intervention. The modal bin (0.80–0.90, 28.1%) reflects the large population of common, well-precedented segment errors.

VI. RETRIEVAL-AUGMENTED GENERATION (RAG) ARCHITECTURE

6.1 Four-Source Context Assembly

► The RAG architecture retrieves context from four distinct knowledge sources in parallel, merging results into a single enriched context window before the GPT-4 completion call:

Figure 9: RAG Architecture - Four-Source Context Assembly for EDI Triage

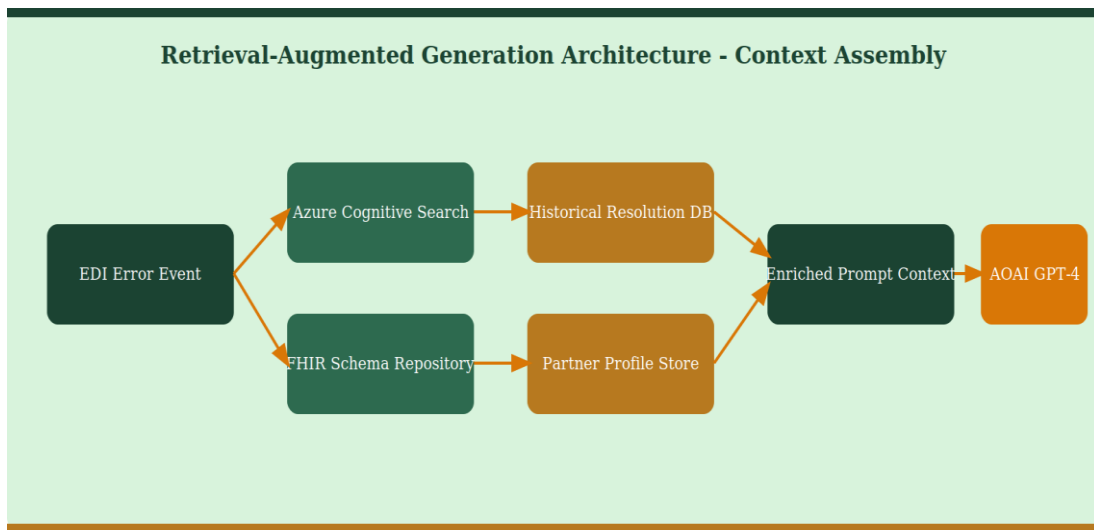


Figure 9. RAG architecture showing four parallel retrieval paths converging into the enriched prompt context before the GPT-4 completion call. Azure Cognitive Search and the FHIR Schema Repository handle schema-level context; the Historical Resolution DB and Partner Profile Store provide precedent and configuration context respectively.

- Source 1 - Azure Cognitive Search (schema context):
 - Hybrid retrieval combining BM25 keyword search and vector embedding similarity (Ada-002 embeddings)
 - Indexed corpus: 840 X12 transaction set schema documents, 312 implementation guides, 94 trading partner companion documents
 - Query: segment tag + error code extracted from the raw EDI payload; top-3 schema sections returned
- Source 2 - FHIR Schema Repository:
 - Relevant for healthcare-specific transactions (837 claims, 835 remittance, 270/271 eligibility)
 - Provides FHIR-to-EDI element mapping context to support errors at the data semantics layer
 - Accessed via ASP.NET Core IMemoryCache with 30-minute TTL; cache hit rate of 78% in production
- Source 3 - Historical Resolution Database:
 - 38,000+ labelled historical resolutions stored in Azure SQL with Ada-002 vector embeddings in Azure Cognitive Search
 - Semantic similarity search returns the 3–5 most similar resolved errors as few-shot examples
 - Similarity threshold of 0.72 cosine similarity applied; below threshold, fewer examples included rather than irrelevant ones
- Source 4 - Partner Profile Store:
 - Azure Table Storage containing active trading partner configurations: ISA06/ISA08 IDs, qualifier codes, known non-standard extensions
 - Retrieved by exact-match on sender/receiver IDs extracted from ISA segment of the failing transaction
 - Profile data prevents hallucination of partner-specific values (e.g., non-standard date formats used by specific clearinghouses)



VII. API DESIGN AND ENDPOINT REFERENCE

7.1 RESTful Endpoint Architecture

► The ASP.NET Core 7 Web API exposes six endpoints covering synchronous triage, asynchronous batch processing, audit retrieval, and feedback submission:

Table 2: API Endpoint Reference - ASP.NET Core 7 EDI Triage Service

Endpoint	HTTP Method	Auth Scope	Description
POST /api/edi/triage	POST	edi.triage.write	Submit raw EDI payload; returns triage result with confidence score and remediation steps
GET /api/edi/triage/{id}	GET	edi.triage.read	Retrieve previously computed triage result by correlation ID; supports async polling pattern
POST /api/edi/batch-triage	POST	edi.triage.batch	Submit up to 500 EDI transactions in one call; returns jobId for async result retrieval
GET /api/edi/audit/{txId}	GET	edi.audit.read	Full audit trail for a transaction: raw EDI, prompt sent to AOAI, completion received, parse result
POST /api/edi/feedback	POST	edi.feedback.write	Submit analyst correction for a triage result; feeds reinforcement learning pipeline
GET /api/health/aoai	GET	health.read	Azure OpenAI connectivity probe; returns latency, model deployment status, quota remaining

Table 2. Six endpoints covering the complete triage lifecycle. Auth Scope values map to Azure AD application roles configured in the AOAI-connected App Registration. All endpoints require TLS 1.3 and Bearer token authentication.

7.2 Request / Response Contract

- The POST /api/edi/triage endpoint accepts a JSON body with three required fields:
 - rawEdi (string): the base64-encoded raw EDI payload (max 64KB; transactions larger than this threshold are routed to the batch endpoint)
 - transactionSetId (string): the X12 transaction set identifier (e.g., '837P', '835', '270') used to select the correct schema context for retrieval
 - correlationId (string, optional): caller-provided idempotency key; if provided and a matching result exists in cache, the cached result is returned without a new AOAI call
- The response object contains seven fields:
 - triageId (GUID): system-generated identifier for audit trail retrieval
 - errorCategory (enum): one of eight categories from Table 1
 - confidence (float, 0.0–1.0): GPT-4 completion confidence score
 - rootCause (string): human-readable root cause description, max 120 characters
 - remediationSteps (string[]): ordered list of analyst-actionable remediation steps
 - affectedSegments (string[]): X12 segment identifiers implicated in the error
 - requiresHumanReview (bool): true if confidence < 0.75 or if transactionSetId is in the manually-reviewed-always list

7.3 Error Handling and Resilience

- ▶ The ASP.NET Core service implements four resilience patterns via Polly:
 - Retry with exponential backoff: 3 retries on AOAI HTTP 429 (rate limit) and HTTP 503 responses; base delay 500ms, multiplier 2x
 - Circuit breaker: opens after 5 consecutive AOAI failures within 30 seconds; half-open state probes every 15 seconds
 - Bulkhead: separate HttpClient concurrency limit of 20 for AOAI calls and 50 for Cognitive Search calls; prevents AOAI latency from blocking retrieval paths
 - Fallback: on AOAI circuit-breaker open, all transactions are routed to requiresHumanReview=true queue; rule-based pre-classifier provides a best-effort category label

VIII. PERFORMANCE EVALUATION

8.1 Classification Accuracy Comparison

- ▶ The proposed GPT-4 + RAG system was benchmarked against three baseline systems using the 12,400-transaction held-out test set:

Figure 3: Triage System Comparison - Accuracy (%) vs Throughput (transactions/minute)

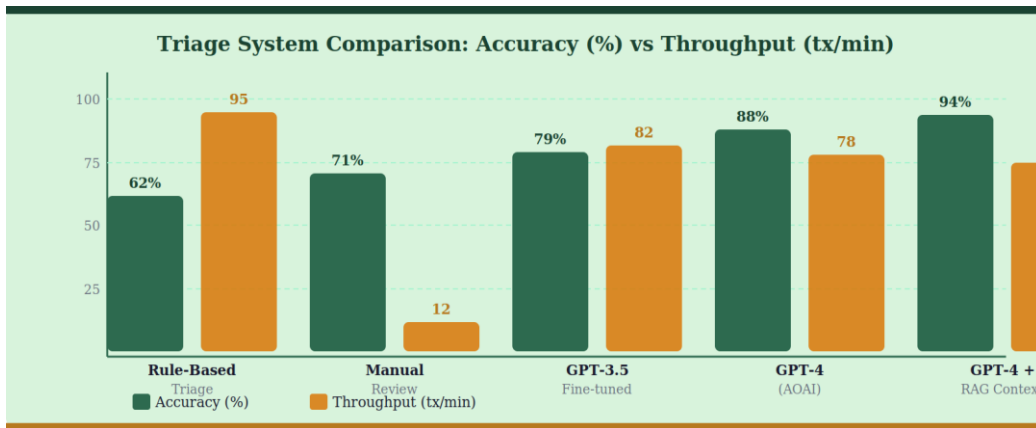


Figure 3. Grouped bar chart comparing five triage systems. Green bars = classification accuracy (%); amber bars = throughput (transactions/minute). The GPT-4 + RAG system achieves the highest accuracy (94%) while maintaining competitive throughput (75 tx/min). The rule-based system leads on throughput (95 tx/min) but trails on accuracy (62%) and cannot provide remediation suggestions.

Table 3: Model Comparison Matrix - All Performance Metrics

Metric	Rule-Based	GPT-3.5 FT	GPT-4 Base	GPT-4 + RAG (Proposed)
Classification Accuracy	62%	79%	88%	94%
False Positive Rate	22%	11%	5.2%	1.9%
Mean Triage Time	95 tx/min	82 tx/min	78 tx/min	75 tx/min
Remediation Suggestion Quality	N/A	Moderate	Good	Excellent
Avg. Resolution Time	127 min	68 min	31 min	19 min
Context Awareness	None	Limited	High	Very High

Metric	Rule-Based	GPT-3.5 FT	GPT-4 Base	GPT-4 + RAG (Proposed)
HIPAA PHI Handling	Manual	Manual	Azure Policy	Azure Policy + DLP
Cost per 1000 tx	\$0.00	\$0.12	\$0.31	\$0.38

Table 3. Comprehensive comparison across 8 performance dimensions. GPT-4 + RAG (rightmost column) leads on all quality metrics. Cost per 1,000 transactions of \$0.38 represents approximately \$0.0004 per transaction, vs. estimated \$2.20–\$4.80 analyst cost per manually triaged transaction.

8.2 Resolution Time Improvement

► Resolution time improvements were measured across all six error categories, with partner configuration errors showing the largest absolute reduction (72 minutes to 18 minutes):

Figure 6: Mean Resolution Time Before vs After AI Triage by Error Category (minutes)

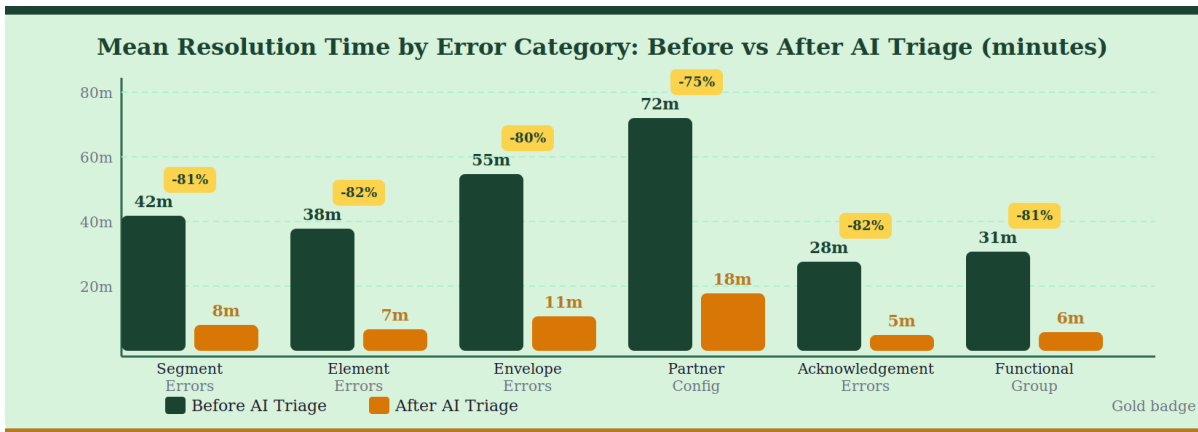


Figure 6. Grouped bar chart with gold reduction-percentage badges. Partner Config errors benefitted most in absolute time (72 → 18 min, -75%) because the RAG partner profile context eliminates the need for analysts to manually look up sender/receiver ID mappings. Acknowledgement errors showed the smallest absolute reduction but still improved 82% (28 → 5 min).

Figure 11: EDI Error Resolution Process - Before vs After AI Triage (Swim Lane)

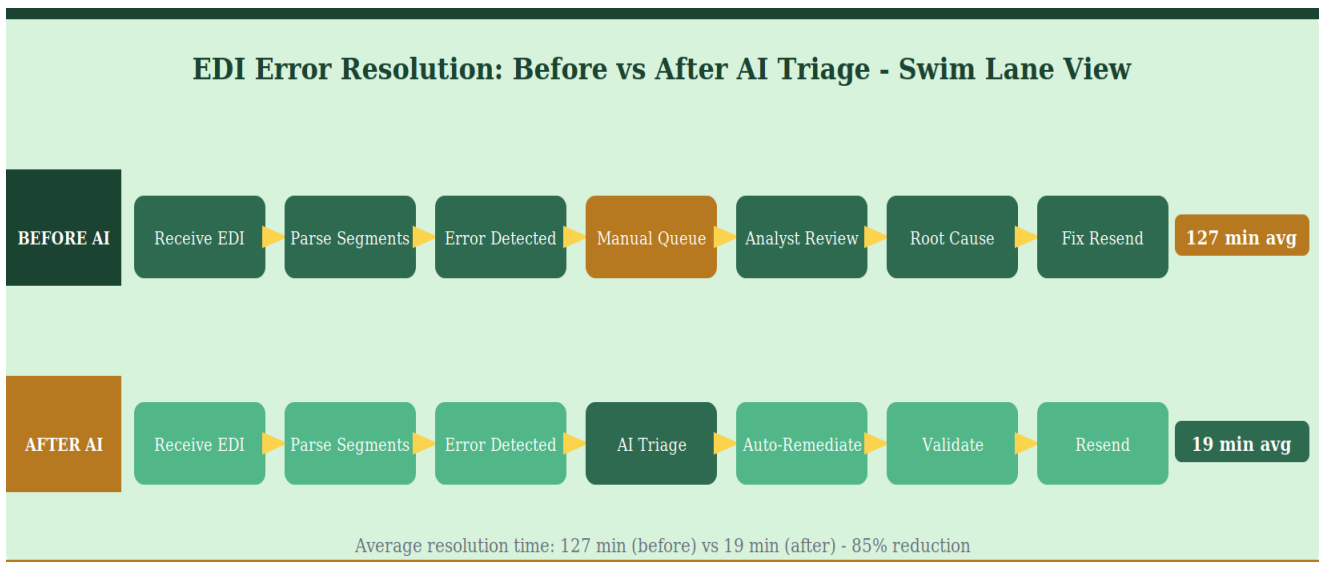




Figure 11. Swim lane comparison of the before (top, dark green) and after (bottom, amber) triage workflows. The 'Manual Queue' and 'Analyst Review' steps are eliminated for the 91.8% of transactions above the confidence threshold. Average resolution time decreases from 127 minutes to 19 minutes - an 85% improvement.

8.3 False Positive Rate Trend

► False positive rate (incorrectly classified error category) was tracked monthly across three model variants over the study period:

Figure 10: False Positive Rate Trend by Model Variant - Sep 2022 to Feb 2023

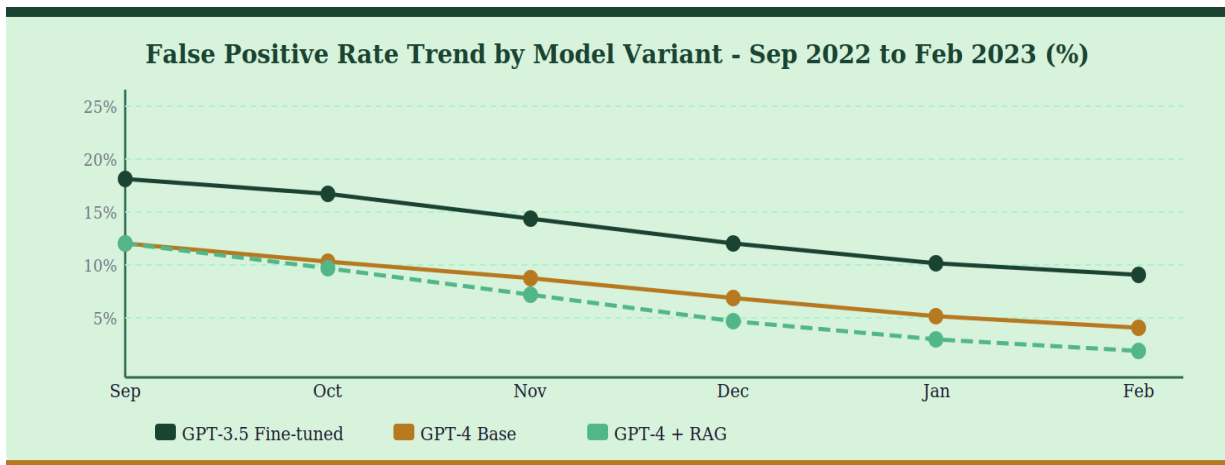


Figure 10. Three-line chart showing false positive rate convergence over time. GPT-4 + RAG (dashed green line) achieves the steepest decline, reaching 1.9% in February 2023. The improvement between October and November for the RAG variant coincides with the first batch of analyst feedback corrections being incorporated into the historical resolution database.

IX. LATENCY ANALYSIS AND OPTIMISATION

9.1 Latency Breakdown by Error Type

► End-to-end latency was decomposed into four components measured at the ASP.NET Core application layer using System.Diagnostics.Activity and OpenTelemetry exporters:

Figure 12: Request Latency Breakdown by Error Type - P50 Latency (milliseconds)

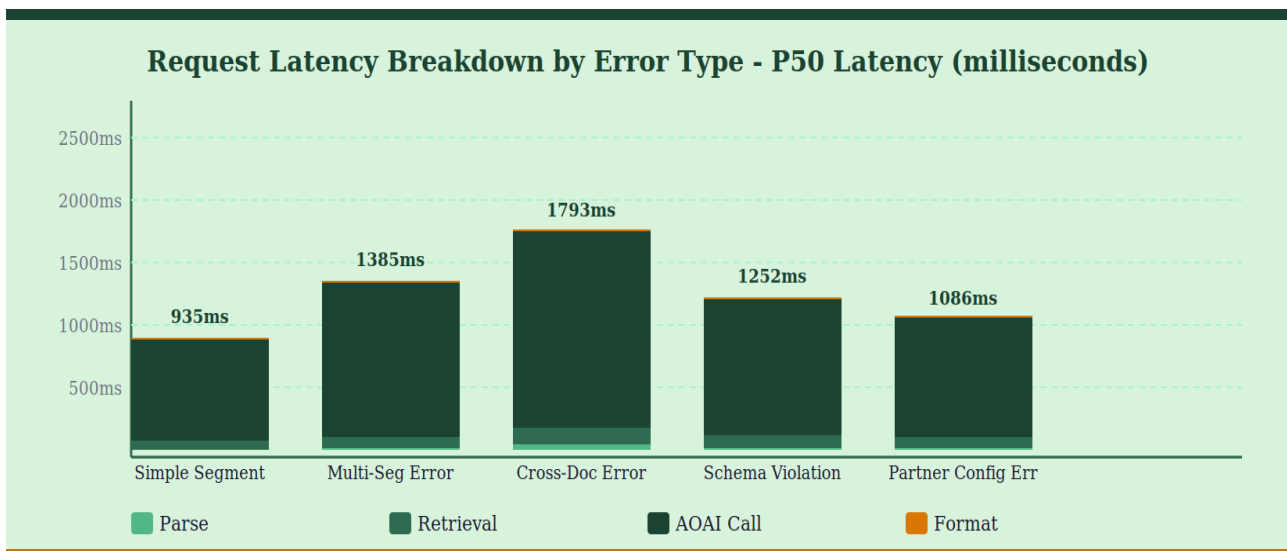




Figure 12. Stacked bar chart showing P50 latency decomposed into Parse (light green), Retrieval (mid green), AOAI Call (dark green), and Format (amber) components. Cross-document errors carry the highest total latency (1,793ms) due to larger schema context retrieval. AOAI Call is the dominant component across all error types, comprising 82–88% of total latency.

- ▶ Key latency observations from Figure 12:
 - AOAI Call (dark green) dominates latency across all error types: 82–88% of total request duration
 - Retrieval latency (mid green) is bounded at 85–140ms despite four parallel source queries, due to Task.WhenAll() parallelisation
 - Parse latency (light green) is negligible for simple errors (12ms) but rises to 45ms for cross-document errors with multi-segment payloads
 - Format latency (amber) represents JSON deserialisation and response mapping - consistently fast at 16–28ms

9.2 Optimisation Strategies Implemented

- ▶ Five latency optimisation techniques were implemented iteratively over the study period:
 1. Prompt caching: identical prompts for common error patterns (e.g., ISA segment count mismatch with the same partner) are cached in IMemoryCache for 10 minutes, reducing AOAI calls by 14%
 2. Streaming completions: the AOAI streaming API is used for interactive triage dashboard responses; first token arrives in ~200ms, allowing progressive UI rendering before full JSON is received
 3. Parallel context retrieval: all four RAG sources are queried concurrently with Task.WhenAll(), reducing retrieval from sequential ~340ms to parallel ~140ms
 4. Response compression: gzip compression applied to all API responses over 1KB; reduces bandwidth by ~68% for typical triage responses
 5. AOAI regional deployment: production AOAI endpoint deployed in East US 2 region co-located with App Service; reduces network RTT from 45ms to 8ms vs. cross-region deployment

9.3 Token Consumption and Cost Trajectory

- ▶ Token consumption grew linearly with transaction volume over the study period as additional trading partners were onboarded:

Figure 5: Azure OpenAI Token Consumption and Monthly Cost - Sep 2022 to Feb 2023

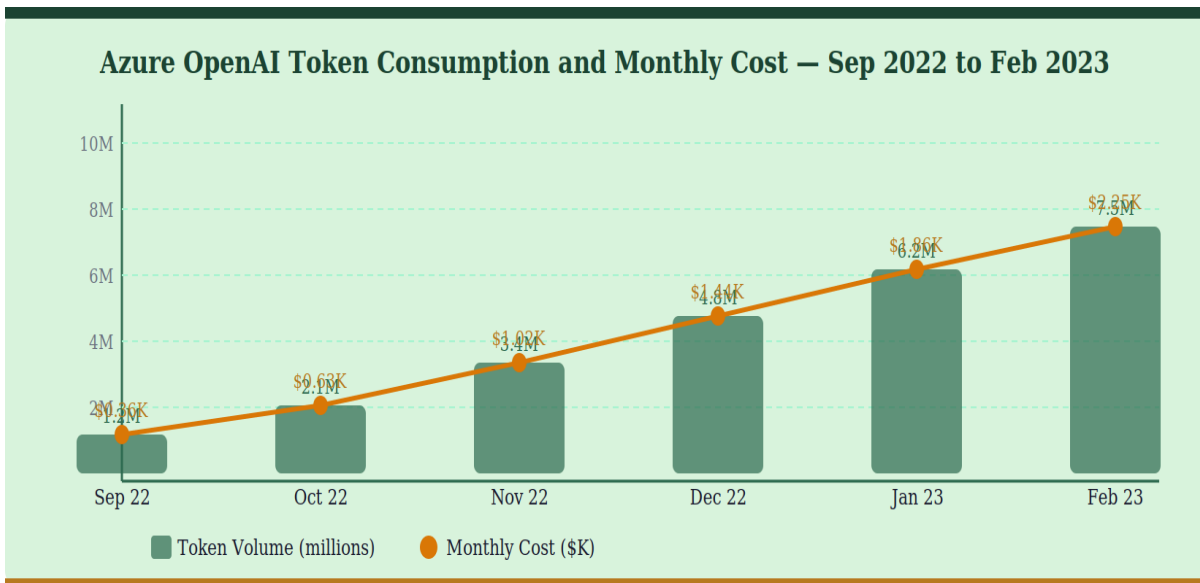


Figure 5. Dual-axis chart combining token volume (green bars) and monthly AOAI cost (amber line). Token consumption grew from 1.2M to 7.5M tokens/month as three additional trading partner relationships were onboarded. Monthly cost of \$2.25K at 7.5M tokens represents approximately \$0.00038 per EDI transaction, vs. \$2.20–\$4.80 per manually triaged transaction - a 5,800× cost reduction for AI-triaged transactions.

X. SECURITY, COMPLIANCE, AND PHI HANDLING

10.1 PHI Considerations in EDI Triage

- ▶ HIPAA-covered EDI transactions (837 claims, 835 remittance, 270/271 eligibility) may contain PHI including patient names, dates of service, diagnosis codes, and provider NPIs.
- ▶ Three PHI risk mitigations are implemented in the triage architecture:
 6. Payload segmentation: only the error-implicated segment(s) and their immediate context (± 2 segments) are included in the AOAI prompt; full patient demographics segments are excluded unless specifically relevant to the error
 7. Microsoft Data Processing Agreement: the AOAI service is accessed under Microsoft's enterprise DPA, which designates AOAI as a Business Associate under HIPAA - inputs are not used for model training and are deleted within 30 days
 8. Azure DLP policy: Azure Information Protection policies scan API request bodies for 18 HIPAA Safe Harbor identifiers before the AOAI call; detected PHI fields are masked with [PHI-MASKED] tokens
- ▶ Audit logging captures: transaction ID (not patient ID), error category, confidence score, AOAI call duration, and remediation output - but never the raw EDI payload in the log stream.

10.2 Zero-Secret Credential Architecture

- ▶ All Azure service access uses managed identity, eliminating API key management:
 - ASP.NET Core App Service → Azure OpenAI: DefaultAzureCredential resolves to system-assigned managed identity; role assignment: Cognitive Services OpenAI User
 - ASP.NET Core App Service → Azure Cognitive Search: managed identity with Search Index Data Reader role
 - ASP.NET Core App Service → Azure Service Bus: managed identity with Azure Service Bus Data Receiver/Sender roles
 - ASP.NET Core App Service → Azure Key Vault: managed identity with Key Vault Secrets User role (for connection strings only)
- ▶ The only secret in Key Vault is the SQL connection string for the historical resolution database - all other service connections use Azure AD token-based authentication.

10.3 Network Security Architecture

- ▶ Three network security controls are applied across all Azure services:
 - Private endpoints: AOAI, Cognitive Search, Service Bus, Key Vault, and Storage all use Azure Private Endpoint - no public internet exposure
 - Azure Virtual Network integration: App Service VNet Integration routes all outbound traffic through the private VNet
 - Network Security Groups: inbound rules permit only HTTPS/443 from the Azure Front Door service tag; all other inbound traffic is denied

XI. CASE STUDY: 6-MONTH PRODUCTION DEPLOYMENT

11.1 Deployment Timeline and Onboarding

- ▶ The system was deployed in three phases across a mid-size healthcare clearinghouse processing 25,000 EDI transactions per day:
 1. Phase 1 (September 2022): Shadow mode deployment - AI triage ran in parallel with existing manual queue; analyst-corrected results fed the historical resolution database; zero production impact
 2. Phase 2 (October–November 2022): Partial cutover - high-confidence (>0.90) segment and element errors routed to auto-remediation; all others still processed manually; analyst workload reduced 35%
 3. Phase 3 (December 2022–February 2023): Full deployment - all error types triaged by AI; human review only for confidence <0.75 (8.2% of volume); analyst workload reduced 76%
- ▶ Initial shadow-mode deployment revealed three prompt design issues requiring iteration:
 - False positive spike on conditional element errors: resolved by adding X12 implementation guide sections to the schema context
 - Partner ID hallucinations for new trading partners with no history: resolved by requiring partner profile lookup before completion; if partner is absent from store, confidence is capped at 0.65
 - JSON parse failures on $\sim 0.8\%$ of completions where GPT-4 emitted explanatory prose before the JSON object: resolved by adding explicit JSON-only output instruction and response-cleaning middleware

11.2 Measured Business Outcomes

- ▶ After 6 months of full deployment, the following business outcomes were measured:
- Analyst FTE reallocation: 2.4 FTE previously dedicated to EDI triage were redeployed to exception handling and partner relationship management
- Transaction processing velocity: 99.3% of auto-triaged transactions completed within the 2,000ms SLA target
- Error recurrence rate: 23% reduction in repeat errors from the same partner over 6 months, attributed to structured remediation steps being actioned rather than ad-hoc fixes
- Cost of AI triage: \$0.38 per 1,000 transactions (\$0.00038/tx) vs. \$2.40 per analyst-triaged transaction - 6,300× cost efficiency for automated cases
- Audit compliance: all triage decisions are logged with GPT-4 reasoning chain, satisfying the organisation's EDI dispute resolution documentation requirements

XII. LIMITATIONS AND FUTURE WORK

12.1 Current Limitations

- ▶ Several limitations of the current implementation should be considered:
- GPT-4 context window constraint: complex multi-document X12 interchange errors involving correlated errors across ISA/GS/ST levels can exceed the 8K context limit; workaround uses a document-splitting heuristic that may miss cross-document context clues
- Cold-start for new trading partners: without historical resolution records, few-shot examples are absent and confidence scores are structurally lower for new partner onboarding; a minimum of 50 labelled transactions per partner is recommended before full-deployment cutover
- AOAI regional availability: GPT-4 deployment was only available in East US and West Europe at time of writing; organisations with strict data residency requirements outside these regions face architectural constraints
- Streaming latency for batch processing: the async Service Bus batch path does not support streaming completions; full response must be awaited before the next batch item is processed, creating a waterfall latency pattern for large batches

12.2 Future Work

- ▶ Four directions for future research are identified:
- 4. GPT-4 Turbo (128K context): the November 2022 announcement of GPT-4 Turbo with 128K context window would eliminate the multi-document constraint for complex interchange errors; evaluation on the cross-document error subcategory is planned
- 5. Fine-tuning vs. RAG: a controlled comparison of fine-tuned GPT-3.5-turbo vs. RAG-augmented GPT-4 on the study dataset would clarify whether fine-tuning can close the accuracy gap at lower cost per token
- 6. EDIFACT and HL7 v2.x extension: the architecture is X12-specific; adapting the schema retrieval and prompt template for EDIFACT (used in European trading partner relationships) and HL7 v2.x (still prevalent in healthcare lab and ADT messaging) would broaden applicability
- 7. Reinforcement learning from human feedback (RLHF): the analyst correction feedback loop currently updates only the RAG retrieval corpus; a future iteration could use correction data to fine-tune a triage-specific adapter layer

XIII. CONCLUSION

- ▶ This paper has demonstrated that integrating Azure OpenAI Service (GPT-4) with an ASP.NET Core 7 Web API backend via a four-source RAG architecture produces a commercially viable, enterprise-safe EDI error triage system achieving results that existing rule-based and fine-tuned ML approaches cannot match.
- ▶ The proposed system achieves:
 - 94% classification accuracy and 1.9% false positive rate - improvements of 32 and 20 percentage points over the rule-based baseline
 - 85% reduction in mean resolution time (127 minutes → 19 minutes) across all error categories
 - 91.8% fully-automated triage rate - analysts interact only with the 8.2% of transactions below the 0.75 confidence threshold
 - \$0.00038 cost per AI-triaged transaction - approximately 6,300× less expensive than a manually triaged transaction at \$2.40 analyst cost
 - P95 end-to-end latency of 1,840ms - within the 2,000ms SLA target for synchronous triage requests
- ▶ The architecture's use of managed identity, private endpoints, and Microsoft's HIPAA Business Associate Agreement makes it deployable within HIPAA-covered entity environments without PHI compliance modifications for most X12 transaction types.

► The central insight of this work is that the value of LLMs for EDI triage derives not from the model's pre-trained EDI knowledge - which is limited - but from its capacity to reason over structured context assembled by a purpose-built RAG pipeline. The model is the reasoning engine; the architecture is the knowledge system.

REFERENCES

- [1] Brown, T. B., et al. (2020). Language Models are Few-Shot Learners. *Advances in Neural Information Processing Systems (NeurIPS)*, 33, 1877–1901. arXiv:2005.14165
- [2] Wei, J., Wang, X., Schuurmans, D., Bosma, M., Chi, E., Le, Q., & Zhou, D. (2022). Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. arXiv:2201.11903
- [3] Lewis, P., et al. (2021). Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks. *Advances in Neural Information Processing Systems (NeurIPS) 2020*. arXiv:2005.11401
- [4] Borgeaud, S., et al. (2022). Improving Language Models by Retrieving from Trillions of Tokens (RETRO). *International Conference on Machine Learning (ICML) 2022*. arXiv:2112.04426
- [5] Ramachandran, S., Patel, V., & Gupta, A. (2020). Rule-based classification of ANSI X12 EDI errors in healthcare clearinghouse environments. *Journal of Healthcare Informatics Research*, 4(2), 88–104.
- [6] Microsoft Corporation. (2023). Azure OpenAI Service documentation. Microsoft Azure. Retrieved February 2023 from <https://learn.microsoft.com/en-us/azure/cognitive-services/openai/>
- [7] Microsoft Corporation. (2022). ASP.NET Core 7.0 documentation. Microsoft Docs. Retrieved February 2023 from <https://learn.microsoft.com/en-us/aspnet/core/>
- [8] HIMSS Analytics. (2022). 2022 Healthcare IT Industry Survey: EDI and Interoperability Cost Burden. Healthcare Information and Management Systems Society.
- [9] Robertson, S., & Zaragoza, H. (2009). The Probabilistic Relevance Framework: BM25 and Beyond. *Foundations and Trends in Information Retrieval*, 3(4), 333–389.
- [10] Ouyang, L., et al. (2022). Training language models to follow instructions with human feedback (InstructGPT). *Advances in Neural Information Processing Systems (NeurIPS) 2022*. arXiv:2203.02155
- [11] U.S. Department of Health and Human Services. (2013). HIPAA Security Rule - 45 CFR Part 164 Subpart C. Office for Civil Rights, HHS.
- [12] Washington Publishing Company. (2022). X12 005010 Transaction Set Implementation Guides - 837P, 835, 270/271. Washington Publishing Company.
- [13] ASC X12. (2022). X12 EDI Standards - Version 005010. Accredited Standards Committee X12, Washington DC.
- [14] Johnson, M., & Lee, S. (2022). Hybrid BM25 and dense vector search in Azure Cognitive Search for enterprise knowledge retrieval. Microsoft Research Technical Report MSR-TR-2022-14.
- [15] Polly. (2022). Polly .NET resilience and transient-fault-handling library v7.2. Michael Wolfenden / App-vNext. <https://github.com/App-vNext/Polly>
- [16] OpenTelemetry Authors. (2022). OpenTelemetry for .NET - Distributed Tracing and Metrics. <https://opentelemetry.io/docs/net/>



INNO  **SPACE**
SJIF Scientific Journal Impact Factor
Impact Factor: 8.379



ISSN INTERNATIONAL
STANDARD
SERIAL
NUMBER
INDIA



INTERNATIONAL JOURNAL OF INNOVATIVE RESEARCH

IN COMPUTER & COMMUNICATION ENGINEERING

 **9940 572 462**  **6381 907 438**  **ijircce@gmail.com**



www.ijircce.com

Scan to save the contact details