# Proposed Solution

## Task 1: Document Classification System - Use Case

### Context

A bank seeks to automate the classification of customer-submitted PDF documents, like statements, contracts, and applications, by detecting their type and storing them correctly.

### System Requirements

#### Input (1.1)

Documents arrive via:

- Email.
- Web portal.
- Mobile app.
- Third-party API.

#### Process (1.2)

- **Extract content**: Read text and analyze images.
- **Classify**: Use AI to identify document type.
- **Store**: Save in the right repository with metadata.

#### Output (1.3)
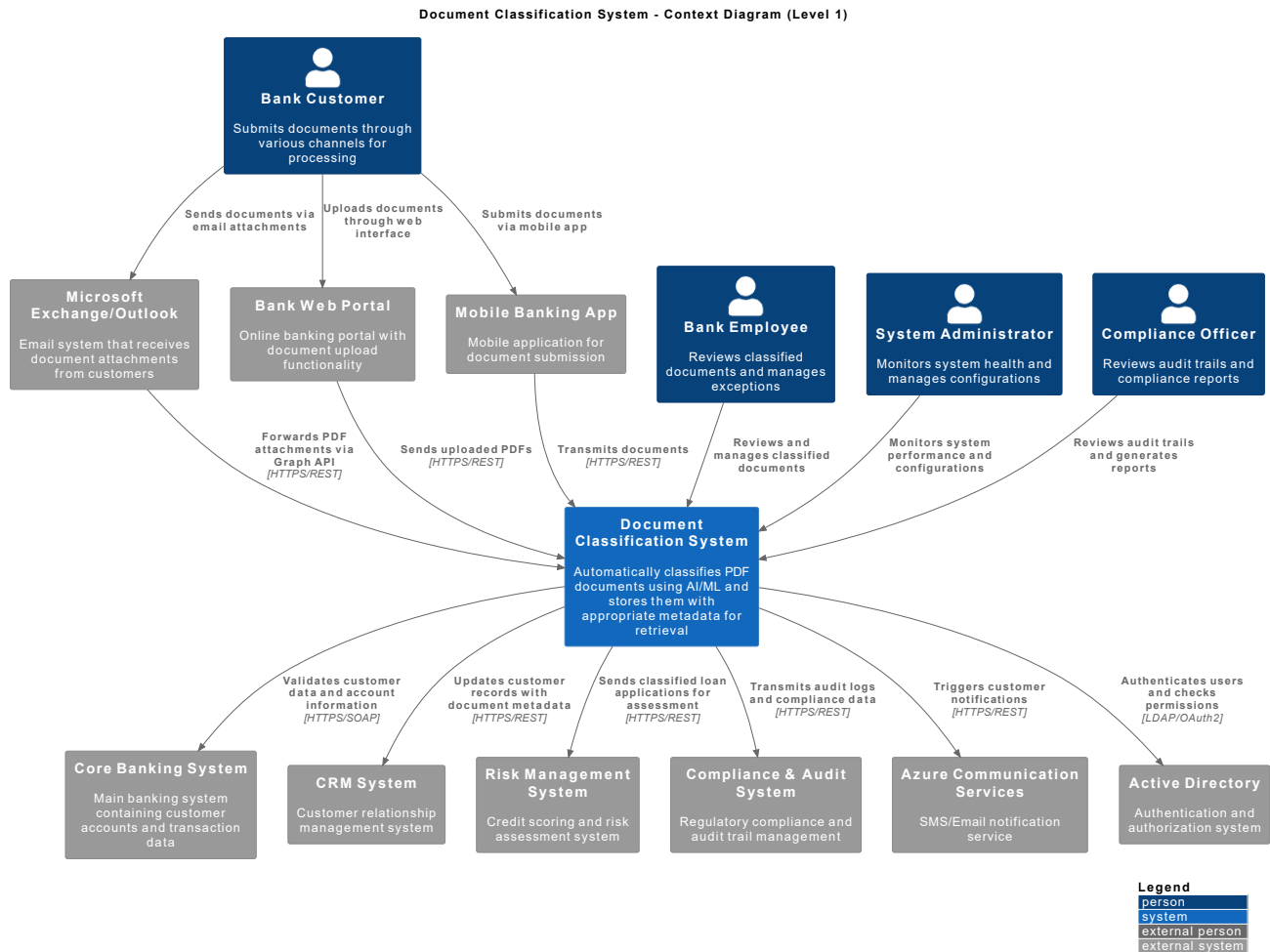
- Tagged by type (e.g., Contract, Statement)
- Indexed for search
- Logged for compliance

# C4 Model Architecture

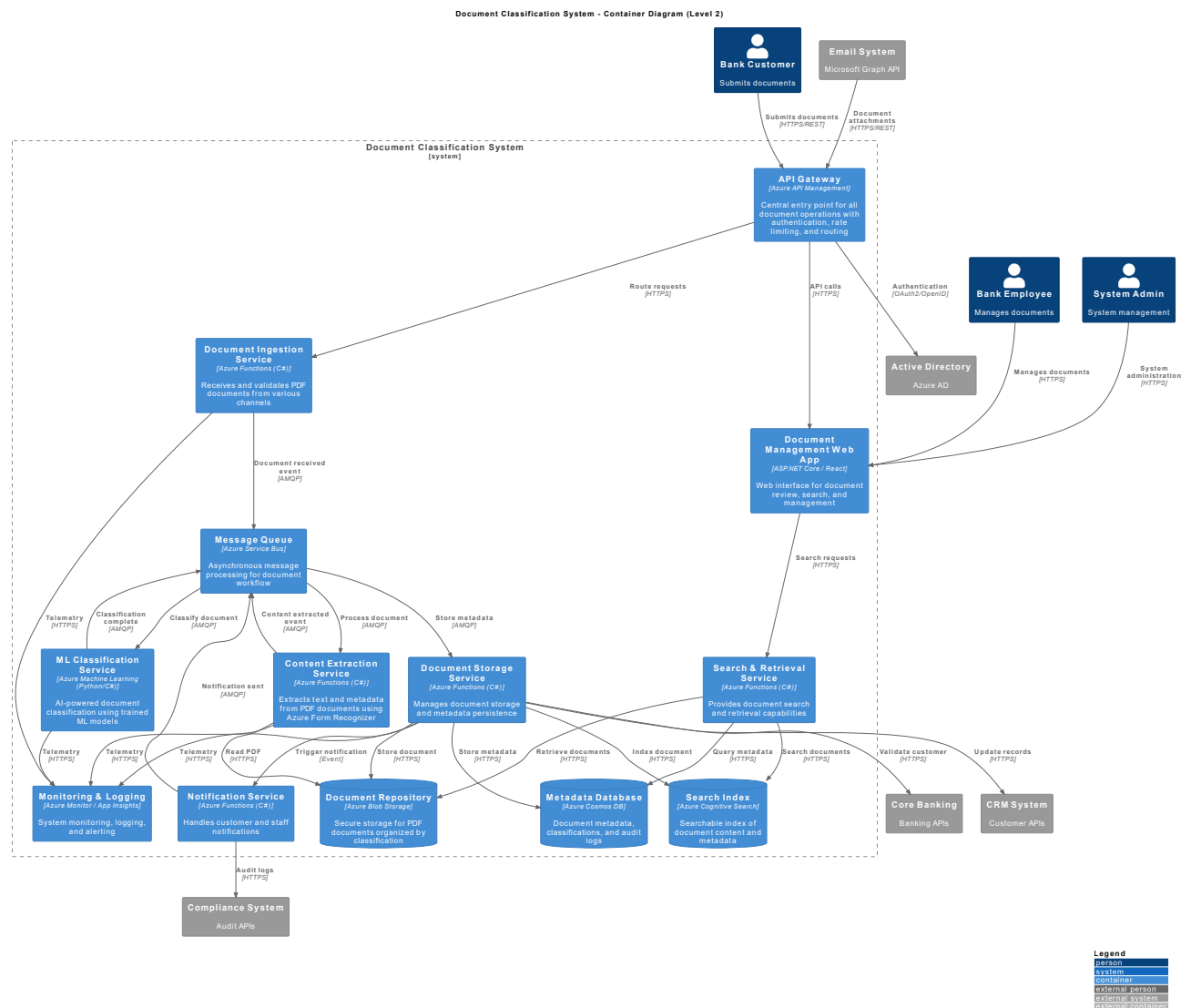This solution uses the **C4 Model** with multiple layers to explain all system interactions:

## 1. Context Diagram (Level 1)

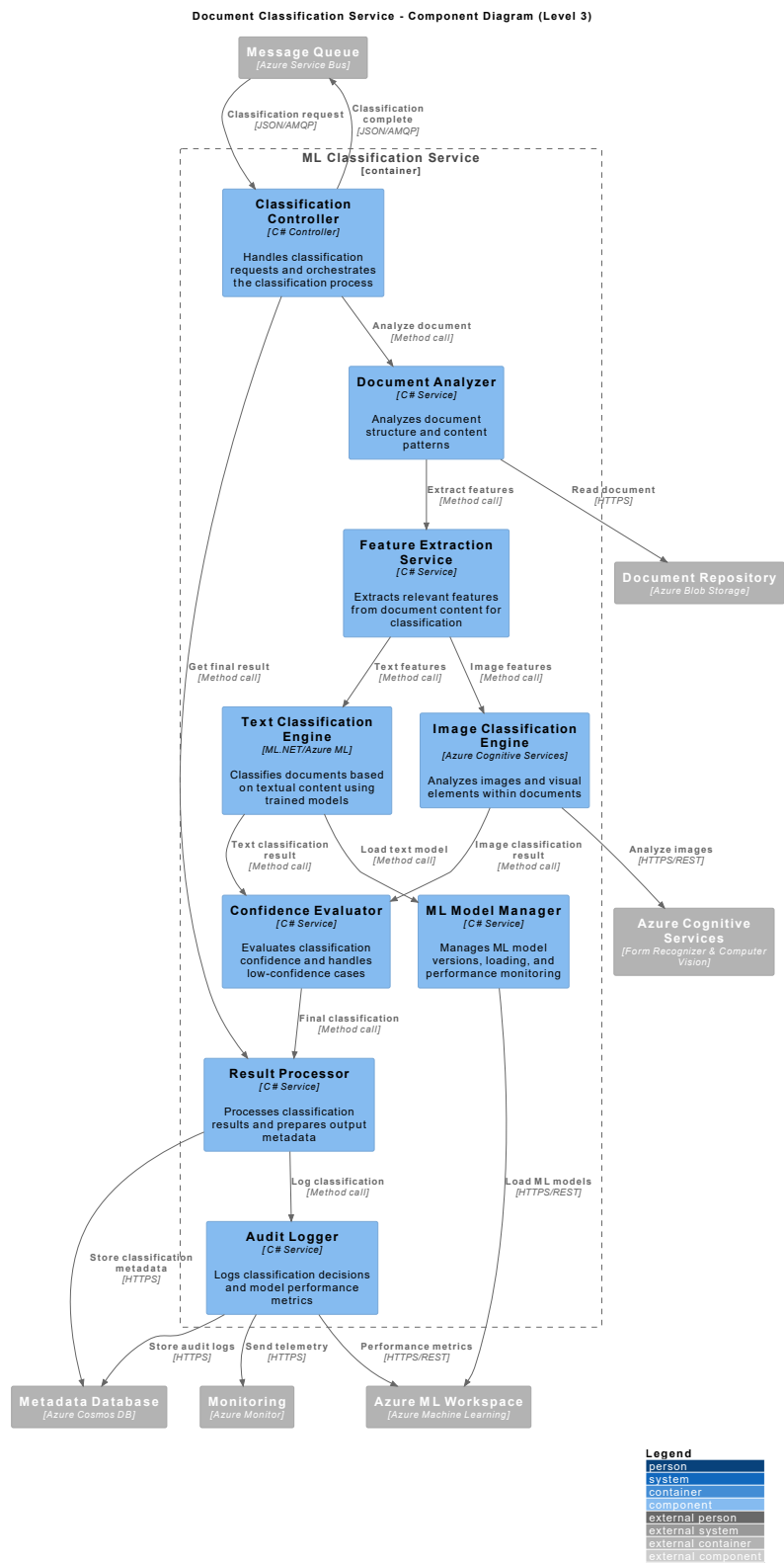Shows the system in its environment with external actors and systems.



Document Classification System - Context Diagram (Level 1)

# 2. Container Diagram (Level 2)

Shows high-level technical building blocks and their interactions.

Document Classification System - Container Diagram (Level 2)

**Bank Customer**
Submits documents

**Email System**
Microsoft Graph API

Submits documents
[HTTPS/REST]

Document attachments
[HTTPS/REST]

Document Classification System
[system]

**API Gateway**
[Azure API Management]
Central entry point for all document operations with authentication, rate limiting, and routing

API calls
[HTTPS]

Authentication
[OAuth2/OpenID]

**Bank Employee**
Manages documents

**System Admin**
System management

Route requests
[HTTPS]

**Document Ingestion Service**
[Azure Functions (C#)]
Receives and validates PDF documents from various channels

**Active Directory**
Azure AD

Manages documents
[HTTPS]

System administration
[HTTPS]

Document received event
[AMQP]

**Document Management Web App**
[ASP.NET Core / React]
Web interface for document review, search, and management

**Message Queue**
[Azure Service Bus]
Asynchronous message processing for document workflow

Search requests
[HTTPS]

Telemetry
[HTTPS]

Classification complete
[AMQP]

Classify document
[AMQP]

Content extracted event
[AMQP]

Process document
[AMQP]

Store metadata
[AMQP]

**ML Classification Service**
[Azure Machine Learning (Python/C#)]
AI-powered document classification using trained ML models

Notification sent
[AMQP]

**Content Extraction Service**
[Azure Functions (C#)]
Extracts text and metadata from PDF documents using Azure Form Recognizer

**Document Storage Service**
[Azure Functions (C#)]
Manages document storage and metadata persistence

**Search & Retrieval Service**
[Azure Functions (C#)]
Provides document search and retrieval capabilities

Telemetry
[HTTPS]

Telemetry
[HTTPS]

Telemetry
[HTTPS]

Read PDF
[HTTPS]

Trigger notification
[Event]

Store document
[HTTPS]

Store metadata
[HTTPS]

Retrieve documents
[HTTPS]

Index document
[HTTPS]

Query metadata
[HTTPS]

Search documents
[HTTPS]

Validate customer
[HTTPS]

Update records
[HTTPS]

**Monitoring & Logging**
[Azure Monitor / App Insights]
System monitoring, logging, and alerting

**Notification Service**
[Azure Functions (C#)]
Handles customer and staff notifications

**Document Repository**
[Azure Blob Storage]
Secure storage for PDF documents organized by classification

**Metadata Database**
[Azure Cosmos DB]
Document metadata, classifications, and audit logs

**Search Index**
[Azure Cognitive Search]
Searchable index of document content and metadata

**Core Banking**
Banking APIs

**CRM System**
Customer APIs

Audit logs
[HTTPS]

**Compliance System**
Audit APIs

Legend
person
system
container
external person
external system
external container

# 3. Component Diagram (Level 3)

Shows internal components within key containers.

**Document Classification Service - Component Diagram (Level 3)**

**Message Queue**
*[Azure Service Bus]*

**Classification request**
*[JSON/AMQP]*

**Classification complete**
*[JSON/AMQP]*

**ML Classification Service**
*[container]*

**Classification Controller**
*[C# Controller]*

Handles classification requests and orchestrates the classification process

**Analyze document**
*[Method call]*

**Document Analyzer**
*[C# Service]*

Analyzes document structure and content patterns

**Extract features**
*[Method call]*

**Read document**
*[HTTPS]*

**Document Repository**
*[Azure Blob Storage]*

**Feature Extraction Service**
*[C# Service]*

Extracts relevant features from document content for classification

**Get final result**
*[Method call]*

**Text features**
*[Method call]*

**Image features**
*[Method call]*

**Text Classification Engine**
*[ML.NET/Azure ML]*

Classifies documents based on textual content using trained models

**Image Classification Engine**
*[Azure Cognitive Services]*

Analyzes images and visual elements within documents

**Analyze images**
*[HTTPS/REST]*

**Text classification result**
*[Method call]*

**Load text model**
*[Method call]*

**Image classification result**
*[Method call]*

**Azure Cognitive Services**
*[Form Recognizer & Computer Vision]*

**Confidence Evaluator**
*[C# Service]*

Evaluates classification confidence and handles low-confidence cases

**ML Model Manager**
*[C# Service]*

Manages ML model versions, loading, and performance monitoring

**Final classification**
*[Method call]*

**Result Processor**
*[C# Service]*

Processes classification results and prepares output metadata

**Log classification**
*[Method call]*

**Load ML models**
*[HTTPS/REST]*

**Store classification metadata**
*[HTTPS]*

**Audit Logger**
*[C# Service]*

Logs classification decisions and model performance metrics

**Store audit logs**
*[HTTPS]*

**Send telemetry**
*[HTTPS]*

**Performance metrics**
*[HTTPS/REST]*

**Metadata Database**
*[Azure Cosmos DB]*

**Monitoring**
*[Azure Monitor]*

**Azure ML Workspace**
*[Azure Machine Learning]*

**Legend**
person
system
container
component
external person
external system
external container
external component

# 4. Deployment Diagram

Shows infrastructure and deployment architecture on Microsoft Azure.

# Key Components & Integrations

## Core System Components

The solution is built from modular services that work together to classify and store documents efficiently:

- **API Gateway**: Entry point for document submissions.
- **Ingestion Service**: Receives and pre-processes PDFs.
- **ML Engine**: Uses AI to classify documents.
- **Document Repository**: Secure archive for storing files.
- **Metadata Manager**: Handles document tags and indexing.
- **Search Service**: Enables quick and accurate retrieval.

## External Integrations

To connect with key systems and ensure smooth operations, the platform integrates with:

- **Email API**: Microsoft Graph for processing attachments.
- **Core Banking API**: Syncs with internal bank systems.
- **CRM API**: Connects customer data for context.
- **Compliance API**: Supports audits and reporting.
- **Notifications**: Azure Communication for user alerts.
- **Cognitive Services**: Powers AI document analysis.

## Azure Technology Stack

The solution is cloud-native, built on Microsoft Azure components:

- API gateway via **Azure API Management**.
- Scalable compute with **Azure Functions**.
- AI-powered insights from **Cognitive Services**.
- Custom ML models with **Azure Machine Learning**.
- Secure storage in **Azure Blob Storage**.
- Metadata housed in **Azure Cosmos DB**.
- Search powered by **Azure Search**.
- Messaging via **Azure Service Bus**.
- Secrets managed with **Azure Key Vault**.
- Real-time monitoring using **Azure Monitor**.

# Task 2: FinHub Data Processing Solution - From HLD to LLD

## Overview

FinHub needs a data platform to collect financial data from multiple APIs (BCRA, ARCA, etc.), run ETL workflows, and deliver analytics to business leaders.

## 2.1 High-Level Design (HLD)

### Architecture Overview

### 1. Ingestion

- Central API Gateway.
- Connectors for BCRA, ARCA, market data.
- Scheduling service and data validation.

### 2. Processing

- ETL Orchestrator.
- Data cleaning and enrichment.
- Business rule application.
- Quality checks and alerts.

### 3. Storage

- Data Lake for raw info.
- Data Warehouse for processed results.
- Metadata and archive layers.

### 4. Analytics & Reports

- Analytics Engine.
- Scheduled & on-demand reports.
- Live dashboards.
- Exposed via REST APIs.

### 5. Security & Compliance

- User authentication.
- Audit logs.
- Encryption for data at rest and in transit.
- GDPR/SOX compliance monitoring.

### 6. Operations

- Monitoring and alerts.
- Central config.

- Backup and recovery.

## Data Sources & Flow

### Sources

- Central banks, tax agencies, market feeds.
- Internal systems and credit bureaus.

### Flow
Ingest → Validate → Transform → Store → Analyze → Report.

## Design Considerations

### Scalability

- Microservices and auto-scaling.
- Redis caching and load balancing.

### Security

- Zero Trust model.
- API protection via OAuth2.
- Secure networking and secret storage.

### Compliance

- Governance, retention policies.
- Full audit trails.
- GDPR data privacy tools.

## 2.2 Low-Level Design (LLD)

### Selected Components - Detailed View

### Data Ingestion Service

Designed with modularity and flexibility in mind, this service uses factory and strategy patterns to handle multiple data sources.

### Main Classes

- `DataIngestionController` : Receives API requests.
- `DataOrchestrator` : Oversees the ingestion flow.
- `DataConnectorFactory` : Instantiates connectors based on source type.
- Specialized connectors: `BCRAConnector` , `ARCAConnector` , `MarketDataConnector` .
- `DataValidator` : Checks data integrity and applies rules.
- `AzureStorageService` : Manages cloud storage.
- `ServiceBusMessageQueue` : Publishes events asynchronously.

### Design Patterns Used

- Factory: For connector creation.
- Strategy: Custom validation per source.
- Dependency Injection: Improves modularity and testability.
- Repository: Abstracts data access.

### Class Interactions

- Controller depends on `Orchestrator` and `Validator` .
- Orchestrator manages flow using `Factory` , `Validator` , `Storage` , `Queue`
- Connectors implement a shared `IDataConnector` interface.

### Data Processing Engine

A robust ETL pipeline, built to support transformation, rule application, and quality control.

### Main Components

- `ETLOrchestrator` : Manages full ETL lifecycle.
- `TransformationService` : Performs data cleanup and enhancement.
- `TransformationEngine` : Executes transformation rules
- `BusinessRulesEngine` : Applies financial logic.
- `FinancialCalculationService` : Handles calculations.
- `DataQualityMonitor` : Flags anomalies.
- `SynapseDataWarehouseService` : Stores results.
- `ProcessingJobManager` : Manages ETL jobs.

**Patterns Applied**

- Chain of Responsibility: Rule processing flow.
- Strategy: Dynamic rule categories.
- Observer: Real-time quality alerts.
- Command: Task execution and tracking.
- Template Method: Structured ETL steps.

**Workflow Summary**

1. Ingestion request hits `ETLOrchestrator`.
2. Data gets validated.
3. Transformation and enrichment are applied.
4. Business rules and calculations executed.
5. Quality checks run.
6. Processed data stored with full lineage.

## 2.3 Architectural Enabler Features

### 1. Azure Data Lake Storage Gen2 Integration

**Purpose**
To handle large volumes of financial data with scalable, secure storage optimized for analytics and regulatory compliance.

**Design Highlights**

- **Hierarchical Namespace**: Faster directory-level operations for analytics.
- **Multi-Protocol**: Compatible with Blob APIs and HDFS.
- **Tiered Storage**: Lifecycle automation across Hot, Cool, Archive.
- **Security**: Azure AD integration, encryption, private endpoints.
- **Performance**: Supports parallel access and large file handling.

**Implementation Phases**

1. Partitioned setup.
2. Lifecycle and tiering automation.
3. Security and compliance hardening.
4. Performance tuning and monitoring.

**Key Features**

- Auto-partitioning by date/source.
- Smart tiering based on access patterns.
- Synapse integration for analytics.
- Audit logs and granular access control.

### 2. Azure Service Bus Message Queue

**Purpose**
To ensure reliable, ordered delivery of financial messages across decoupled services.

**Design Highlights**

- **Topics/Subs**: Multi-consumer message routing.
- **Message Sessions**: Maintain transaction order.
- **Dead Letter Queues**: Retry failed messages.
- **Duplicate Detection**: Ensure data consistency.
- **Auto-Scaling**: Adjusts to queue load dynamically.

**Implementation Phases**

1. Core messaging setup.
2. Advanced features (DLQ, sessions, deduplication).
3. Monitoring & alerting.

4. Performance and scaling optimization.

**Key Features**

- At-least-once delivery.
- Retry logic with exponential backoff.
- Correlated messaging for data tracing.
- Full observability via Azure Monitor.

## 3. Monitoring with Azure Monitor

**Purpose**

To ensure visibility into system health, performance, and data quality, with alerts and automation for proactive issue resolution.

**Design Highlights**

- **Multi-Layer Monitoring**: Infra, app, business metrics.
- **Real-Time Alerts**: Issues surfaced immediately.
- **Distributed Tracing**: Full visibility across services.
- **Custom Dashboards**: KPIs tailored to business needs.
- **Automation**: Enables self-healing workflows.

**Implementation Phases**

1. Basic system monitoring.
2. Business metrics and quality tracking.
3. Predictive alerts and analytics.
4. Automated remediation.

**Key Features**

- Telemetry for data pipelines.
- Anomaly detection on quality signals.
- KPI dashboards for executives.
- Integration with incident management tools.

# Diagrams

## High-Level Design architecture diagram



FinHub Data Processing Solution - High Level Design (HLD)

# Data Ingestion Service class diagram


Data Ingestion Service – Class Diagram (LLD)

# Data Processing Engine class diagram



Data Processing Engine – Class Diagram (LLD)

# Data ingestion sequence diagram

**Data Ingestion Process - Sequence Diagram**

| External API (BCRA/ARCA) | API Gateway | Data Ingestion Controller | Data Orchestrator | Data Connector Factory | BCRA Connector | Data Validator | Azure Storage Service | Service Bus Message Queue | ETL Orchestrator |
|---|---|---|---|---|---|---|---|---|---|

**Data Ingestion Request**

POST /api/ingest/bcra {scheduleId, parameters}

ProcessDataAsync(request)

CreateConnector(BCRA)

new BCRAConnector()

BCRAConnector instance

**Data Fetching**

ConnectAsync()

Authenticate & Connect

Connection established

true

FetchDataAsync(parameters)

GET /exchange-rates?date=today

Exchange rate data (JSON)

RawData object

**Data Validation**

ValidateAsync(rawData)

ValidateSchemaAsync(rawData, schema) — Check JSON structure, required fields, data types

ValidateBusinessRulesAsync(rawData) — Check exchange rate ranges, currency codes, timestamps

ValidateDataQuality(rawData) — Check completeness, accuracy, consistency

ValidationResult(isValid: true)

**Data Storage**

StoreRawDataAsync(rawData)

GenerateStoragePath(rawData) — Path: /raw/bcra/2025/01/20/exchange-rates-{timestamp}.json

Azure Blob Storage Upload raw data

Azure Cosmos DB Store metadata

dataId: "bcra-20250120-001"

**Event Publishing**

PublishAsync(IngestionCompleteEvent)

Event: {
jobId: guid,
dataId: "bcra-20250120-001",
sourceType: "BCRA",
status: "Completed"
}

Message published

IngestionResult(success, dataId)

HTTP 200 OK {jobId, status, dataId}

**ETL Processing Trigger**

IngestionCompleteEvent received — Async processing starts in background

ProcessDataAsync(dataId) — Starts transformation and business rules pipeline

**Error Handling Scenario**

Alternative Flow: Validation Failure

Controller -> Orchestrator : ProcessDataAsync(request)
Orchestrator -> Validator : ValidateAsync(invalidData)
Validator --> Orchestrator : ValidationResult(isValid: false, errors)

Orchestrator -> Queue : PublishAsync(IngestionFailedEvent)
note right : Event: {\n jobId: guid,\n status: "Failed",\n errors: [validation errors]\n}

Orchestrator --> Controller : IngestionResult(failed, errors)
Controller --> Gateway : HTTP 400 Bad Request\n{errors, jobId}

**Monitoring and Logging**

Cross-cutting Concerns

Controller -> Controller : Log request received
Orchestrator -> Orchestrator : Log processing steps
Validator -> Validator : Log validation results
Storage -> Storage : Log storage operations
Queue -> Queue : Log message publishing

All components send telemetry to Azure Monitor

| External API (BCRA/ARCA) | API Gateway | Data Ingestion Controller | Data Orchestrator | Data Connector Factory | BCRA Connector | Data Validator | Azure Storage Service | Service Bus Message Queue | ETL Orchestrator |
|---|---|---|---|---|---|---|---|---|---|

# Data flow diagram

**FinHub Data Processing Solution - Data Flow Diagram**



*«external»*
**External Data Sources**

- BCRA API Exchange Rates
- ARCA API Settlement Data
- Market Data APIs Stock Prices
- Credit Bureau APIs Risk Data
- Internal Systems Transactions

1. Exchange rate data (JSON, REST API)
2. Settlement data (JSON, REST API)
3. Market prices (JSON, REST API)
4. Credit scores (JSON, REST API)
5. Transaction data (JSON, REST API)

*«ingestion»*
**Data Ingestion Layer**

- API Gateway
- Ingestion Scheduler
- Data Connectors
- Data Validator

6. Route to appropriate connector
7. Trigger scheduled collections
8. Raw data validation (Schema, business rules)
9. Store validated raw data
10. Publish ingestion complete event

**Message Queue**

- Azure Service Bus Topics & Subscriptions

11. Trigger ETL processing
Processing status & error events
12. Data transformation (Clean, normalize)

*«processing»*
**Data Processing Engine**

- ETL Orchestrator
- Business Rules Engine
- Data Transformation Service
- Data Quality Monitor

**Event-Driven Architecture**
- Loose coupling between services
- Guaranteed message delivery
- Dead letter queue handling
- Scalable message processing

13. Apply business rules & calculations
14. Data quality monitoring
15. Store processed data (partitioned)
16. Store lineage & metadata
Transformation metrics
Quality alerts & notifications
Quality metrics for validation tuning

17. Archive old data (compliance)
Performance feedback for rule optimization

*«storage»*
**Data Storage Layer**

- Metadata Repository (Azure SQL Database)
- Data Warehouse (Azure Synapse)
- Archive Storage (Azure Blob)
- Raw Data Lake (Azure Data Lake Gen2)

**Data Warehouse**
- Star/snowflake schema
- Columnar storage
- Automated indexing
- Query optimization

**Raw Data Storage**
- Immutable source data
- Partitioned by date/source
- Multiple formats supported
- Retention: 7 years

18. Query processed data for analysis

*«analytics»*
**Analytics & Reporting**

- Analytics Engine (Synapse Analytics)
- Report Generator (Power BI Service)
- Dashboard Service (Power BI)
- Analytics APIs (Azure Functions)

**Analytics Capabilities**
- Real-time processing
- Complex financial calculations
- Machine learning integration
- Regulatory compliance reporting

19. Generate scheduled reports
20. Real-time dashboard data feed
21. Expose analytics via REST APIs
22. Management reports & insights
23. Interactive analytics dashboards
24. System monitoring & administration

**Management Users**

- Management Team
- Financial Analysts
- System Administrators

# System Architecture Design Approach and Decisions

## Design Summary

I built both solutions to be cloud-first, scalable, and secure—leaning heavily on Microsoft Azure's native tools. For the document classifier, I used C4 modeling to break down the architecture. For the FinHub pipeline, I went with a clean, layered approach powered by microservices and event-driven logic.

## Architectural Highlights

- **Cloud-Native by Default**: Azure Functions and Synapse reduce overhead and handle scale with ease.
- **Event-Driven Messaging**: Service Bus keeps components loosely coupled and responsive.
- **Built-In Security**: Active Directory, Key Vault, and API Management ensure protection and compliance.

## Technical Decisions

I went with C# for strong tooling and Azure integration. Architecture follows microservices with patterns like CQRS and event sourcing. AI features combine Azure Cognitive Services with custom models for flexibility and fast delivery.