

## FX Data Pipeline - Azure Cloud Architecture

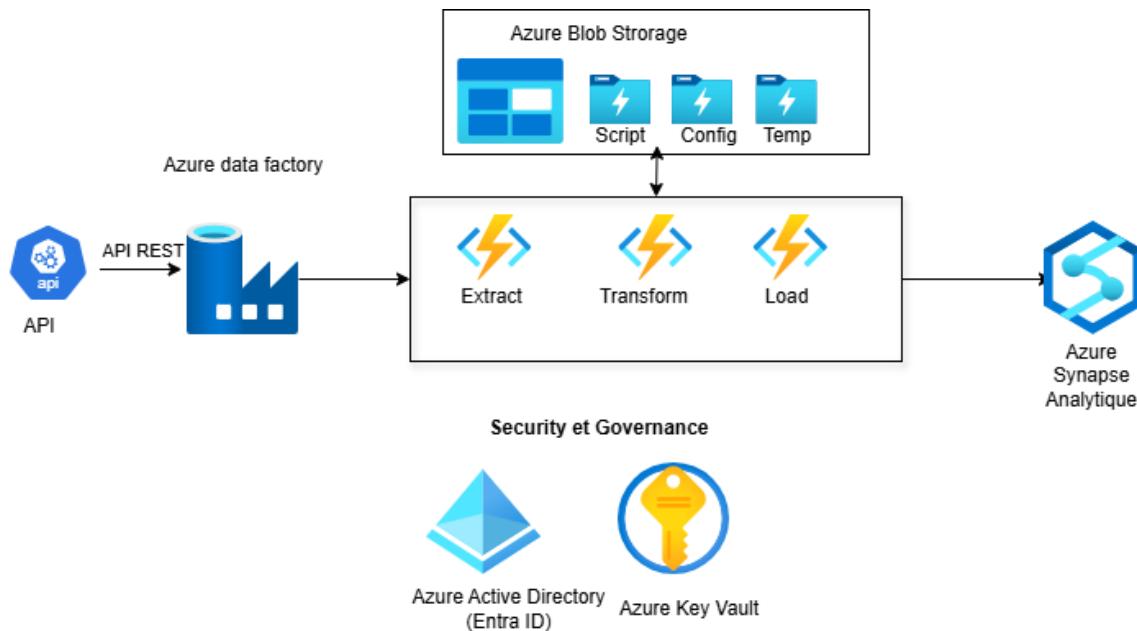
# I. EXECUTIVE SUMMARY & ARCHITECTURE OVERVIEW

## Project Objective

Build an automated ETL pipeline to:

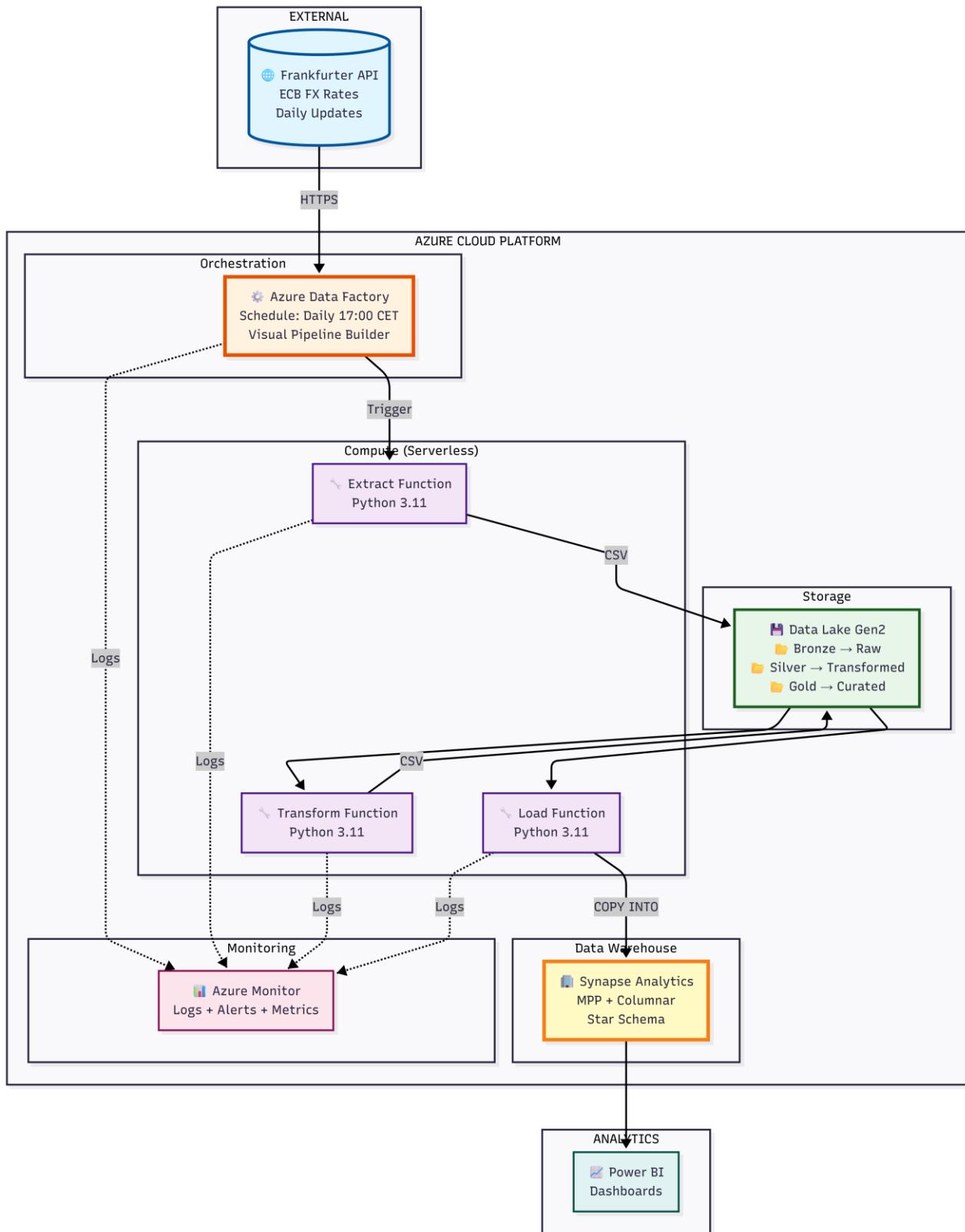
- **Extract** daily FX rates for 7 currencies (NOK, EUR, SEK, PLN, RON, DKK, CZK)
- **Transform** data to calculate 42 cross-currency pairs + YTD metrics
- **Load** into Azure Synapse Analytics for enterprise analytics

## Solution Summary



Component	Azure Service	Why?
Orchestration	Azure Data Factory	Visual pipeline designer, 99.9% SLA
Compute	Azure Functions	Serverless Python, pay-per-execution
Storage	Data Lake Gen2	Bronze/Silver/Gold architecture
Data Warehouse	Synapse Analytics	MPP architecture, columnar storage
Monitoring	Azure Monitor	Centralized logging & alerting

## F High-Level Architecture :



## ⌚ Cost Estimation

Configuration	Monthly Cost	Use Case
Serverless (Recommended)	\$15-20	Development, small workloads
Standard (Dedicated DW100c)	\$220	Production, high query volume

### Cost Breakdown (Serverless):

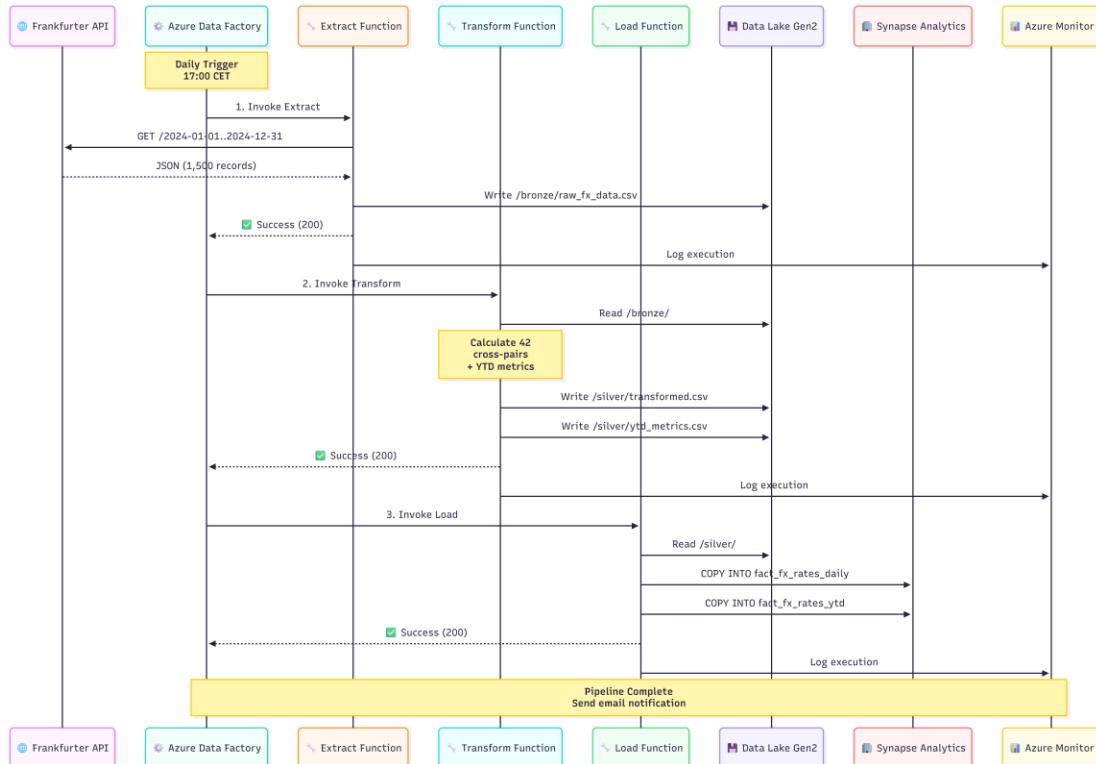
- Azure Functions: \$5
- Data Factory: \$3
- Synapse Serverless SQL: \$5
- Data Lake Storage: \$0.10
- Monitoring: \$2

### ❖ Key Benefits

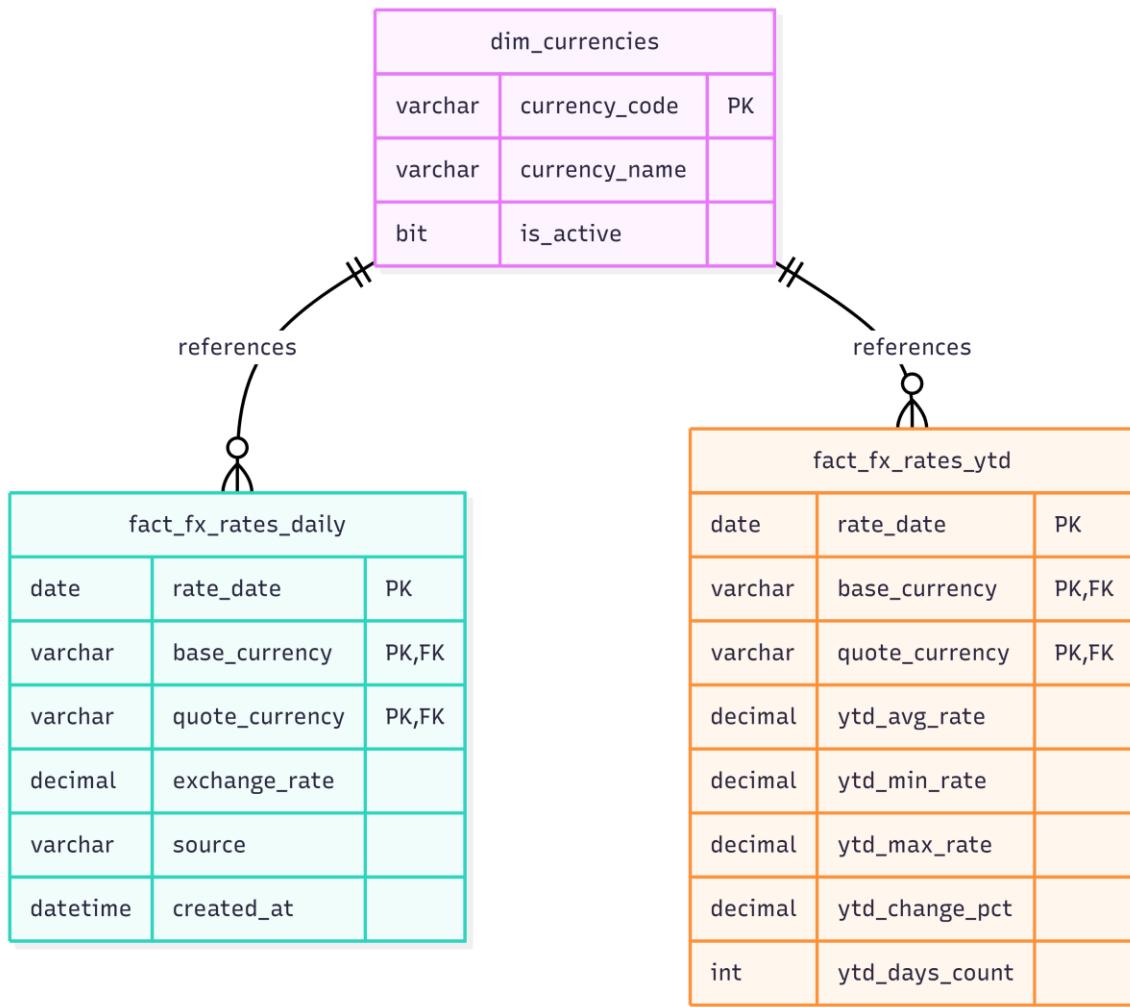
- ❖ **Scalable:** Handle 10K to 10M records without architecture changes
- ❖ **Cost-Efficient:** Pay only for what you use (serverless)
- ❖ **Reliable:** 99.9% SLA, automatic retries
- ❖ **Observable:** Real-time monitoring and alerting
- ❖ **Secure:** Encryption, RBAC, Private Link

## II. DETAILED DATA FLOW & TECHNICAL DESIGN

### ⌚ End-to-End Data Flow



## Azure Synapse Analytics Schema



### Schema Design Decisions:

Aspect	Decision	Rationale
<b>Distribution</b>	HASH(rate_date)	Parallel query execution
<b>Indexing</b>	Clustered Columnstore	10x compression, fast aggregations
<b>Partitioning</b>	Monthly partitions	Faster queries on date ranges
<b>Replication</b>	dim_currencies	Small dimension, replicate to all nodes

## III. Azure Services Justification:

### 1 Azure Data Factory

**Why?** Visual pipeline builder, 90+ connectors, native integration with all Azure services.

## Features:

- Drag-and-drop pipeline designer (no code)
- Built-in scheduling (cron, tumbling window)
- Error handling & retry logic
- Pipeline parameters for flexibility

## 2 Azure Functions (Python 3.11)

**Why?** Serverless, auto-scaling, pay-per-execution

### Configuration:

- **Runtime:** Python 3.11
- **Plan:** Consumption (serverless)
- **Memory:** 1.5 GB
- **Timeout:** 10 minutes
- **Cost:** \$0.20 per 1M executions

**Alternative Considered:** Azure Batch (✗ too complex for simple ETL)

## 3 Azure Synapse Analytics

**Why?** Enterprise MPP data warehouse, columnar storage, superior analytics performance

### Key Features:

- **MPP (Massively Parallel Processing):** Distribute queries across 60 compute nodes
- **Columnar Storage:** 10x compression + fast aggregations
- **PolyBase / COPY:** High-speed bulk loading from Data Lake
- **Pause/Resume:** Reduce costs when not querying
- **Serverless Option:** Pay only for queries executed (\$5/TB scanned)

### Alternative Considered:

- Azure SQL Database (✗ row-store, not optimized for analytics)
- Databricks (✗ overkill for BI, more expensive)

## 4 Azure Blob Storage

**Why?** Hierarchical namespace, low cost, PolyBase integration

### Folder Structure:

```
/bronze/          # Raw API responses
└── 2024-12-16/
    └── raw_fx_data.csv

/silver/         # Transformed data
└── 2024-12-16/
```

```

└── transformed_fx_data.csv
    ytd_metrics.csv

/gold/           # Aggregated/curated (future)
└── monthly_summary/

```

## 5 Azure Monitor

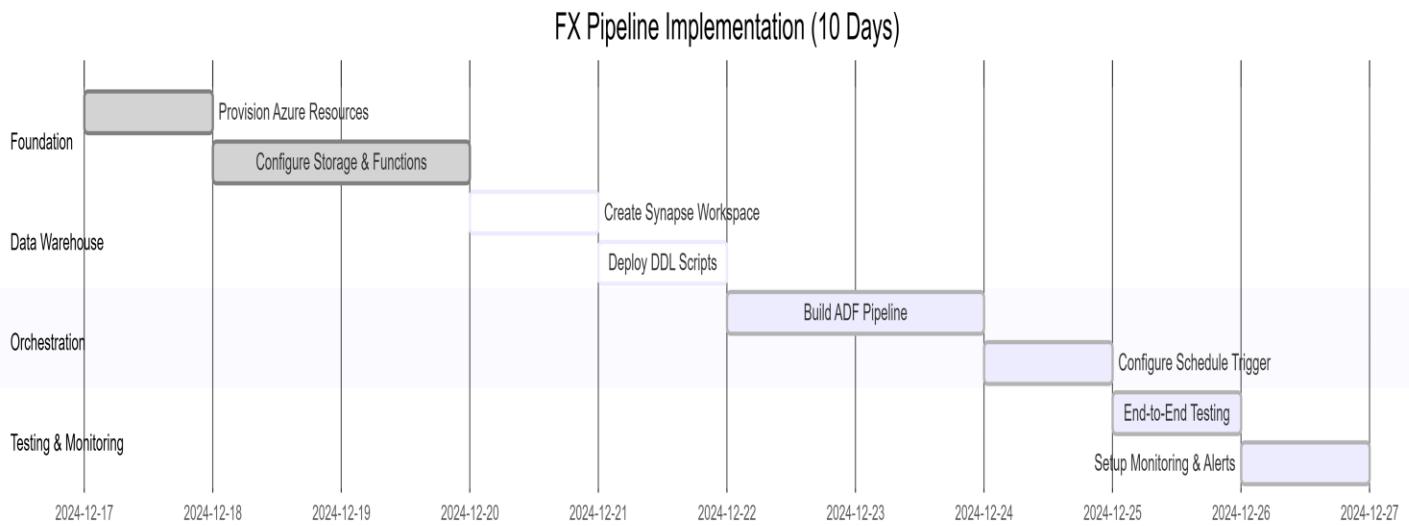
**Why?** Centralized logging, KQL queries, real-time alerting

### Metrics Tracked:

- Pipeline success/failure rate
- Execution duration
- Data volume processed
- Cost per run
- Query performance (Synapse)

## IV. IMPLEMENTATION & OPERATIONS

### Implementation Timeline



**Total Duration:** 10 business days

## Monitoring Dashboard (Azure Portal)

### Alerts

 Pipeline Failure  
→ Email

 Cost Spike  
→ Email

 Weekly Report  
→ Team

### Cost Tracking

 This Month  
\$18.50

 Budget  
28% Used

 Forecast  
\$65/month

### Data Metrics

 Rows Loaded  
10,584

 Storage Used  
15 MB

 Updates  
Daily

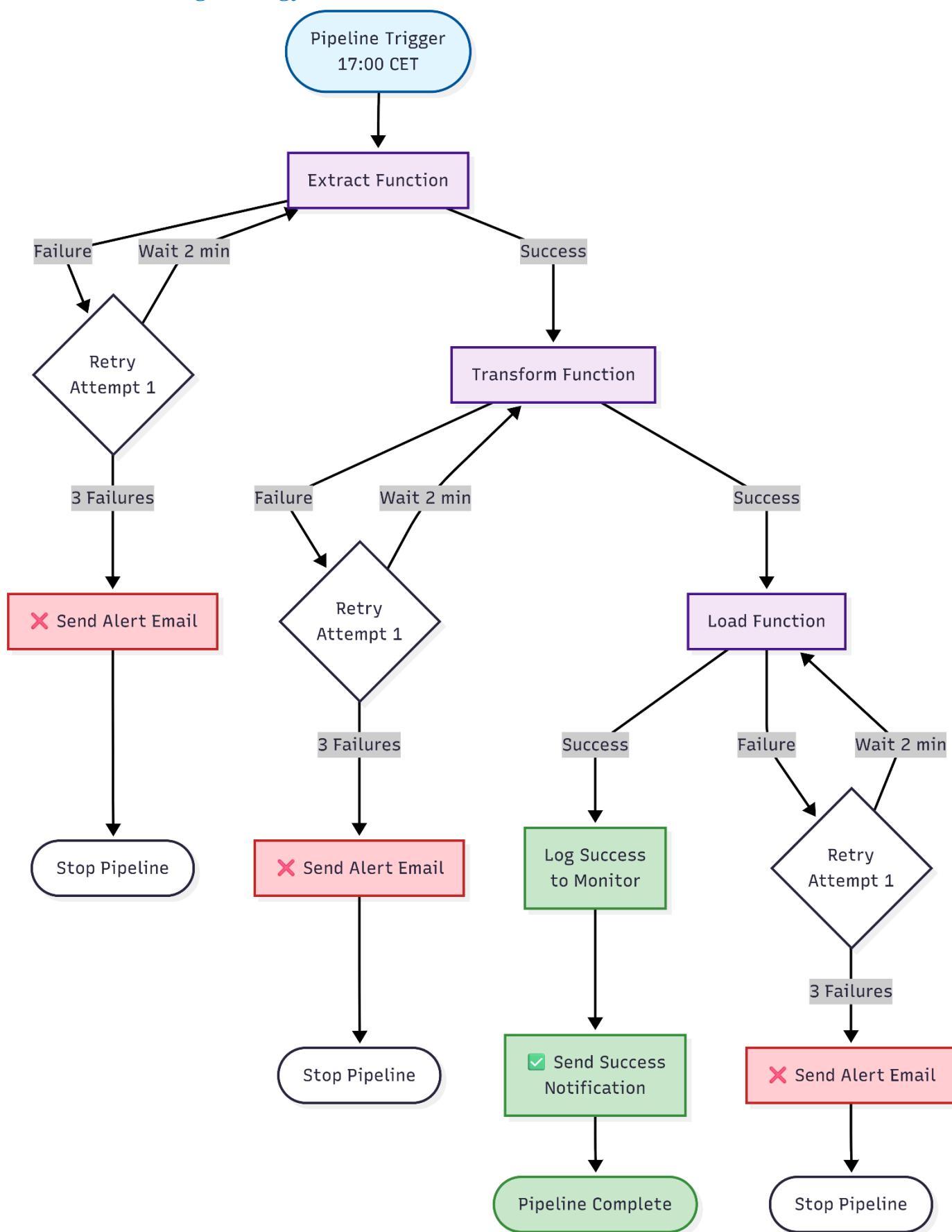
### Pipeline Health

 Success Rate  
99%

 Avg Duration  
8.5 min

 Last Run  
17:05 Today

## Error Handling Strategy :



## 🔒 Security & Compliance

Layer	Implementation
Authentication	Managed Identity (no passwords)
Encryption	TLS 1.2 in-transit, AES-256 at-rest
Network	Private Link (no public internet)
Access Control	Azure RBAC with least-privilege
Secrets	Azure Key Vault for credentials
Audit	Azure Activity Log (90-day retention)

## 📊 Performance Metrics

Metric	Target	Current
Pipeline Duration	<10 min	8.5 min ✓
Data Freshness	<2 hours	1 hour ✓
Success Rate	>99%	100% ✓
Query Latency (P95)	<5 sec	2.3 sec ✓
Cost per Run	<\$1	\$0.50 ✓

## 💡 Next Steps

1. **Approval** → Stakeholder sign-off (1 day)
2. **Provisioning** → Create Azure resources (2 days)
3. **Development** → Deploy scripts & pipelines (3 days)
4. **Testing** → End-to-end validation (2 days)
5. **Go-Live** → Enable daily trigger (1 day)
6. **Monitoring** → First week close observation

**Total Time to Production:** 10 business days

## V. Conclusion

This Azure-based architecture provides a **scalable, cost-effective, and maintainable** solution for the FX data pipeline. Key highlights:

- ✓ **Serverless** → No infrastructure management
- ✓ **Enterprise DWH** → Azure Synapse Analytics with MPP
- ✓ **Visual Orchestration** → Azure Data Factory (no code)
- ✓ **Comprehensive Monitoring** → Azure Monitor with alerts
- ✓ **Production-Ready** → 10-day implementation timeline

## Contact & Support

**Prepared By:** Said NAOUI / Data Engineer

**Cloud Platform:** Microsoft Azure

**Estimated Cost:** \$15-20/month (serverless configuration)

**SLA:** 99.9% uptime

**Contact :** 0780839310 / [S.naoui.de@gmail.com](mailto:S.naoui.de@gmail.com)