

**CONTINO**

# **Mainframe & Serverless Integration**

**How to Liberate the Data  
and Stay Competitive in the New World!**

Federico Fregosi - VP of Engineering

**CONTINO**

# **Mainframe & Serverless Integration**

**Come liberare i dati e rimanere competitivi  
nel nuovo mondo!**

Federico Fregosi - VP of Engineering

# #WhoAml



## Federico Fregosi

Current position: **VP of Engineering, EMEA**

Past:



[federico.fregosi@gmail.com](mailto:federico.fregosi@gmail.com)

<https://www.linkedin.com/in/federico-fregosi/>

# Large Payments Providers Are at a Crossroads

From the 90s, the payments industry made billions (\$\$\$!) off the back of first-class reliability and resilience.

But modern tech is pushing the source of value beyond just reliability and resilience.

But there is a solution! Payment providers are sitting on a **DATA GOLD MINE** in their mainframe.

This data contains insights that will power their product and customer strategies for years to come.

Getting to that data isn't easy.



# Payment Providers Need to Evolve!

You need to shift your business from delivering resilient, reliable transactions....

...to providing real-time integrated data-driven services.



You can process and stream millions of transactions to the cloud, where they will then be available for analysis! It's **low risk** and **high value**.

## Stream & Use

1. **Replicate** the information on the mainframe in the cloud
2. **Deploy** cloud-native, real time streaming platforms (Managed Kafka, Data Streams)
3. **Define** pre or post arrival triggers
4. **Separate** batch processing where appropriate
5. **Make** available to use by your teams

By **connecting** your cloud, data and product programs, you can start delivering business value from real-time transactions in the cloud in around **6-12 months**.

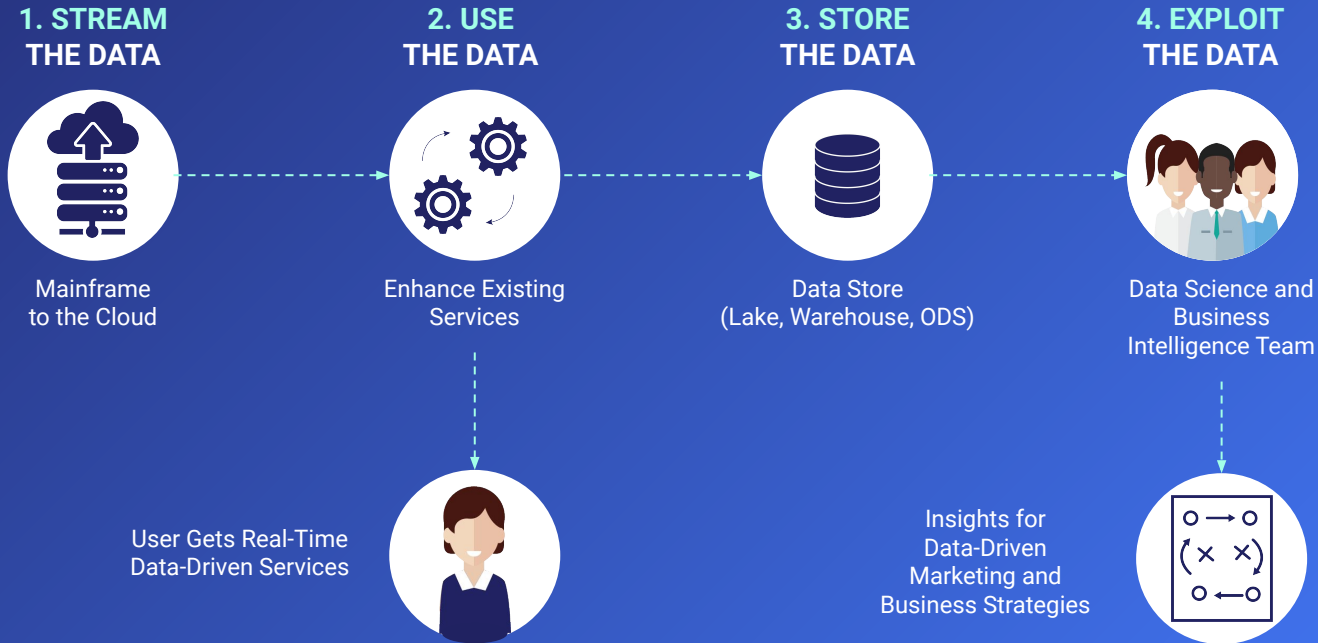
## Store & Exploit

1. **Ensure:** the data is retained in a secure and flexible store
2. **Use** the flexibility of the cloud to meet your requirements
3. **Democratize** access to your data
4. **Generate** business value by providing data-driven evidence
5. **Turn** data into knowledge and actions

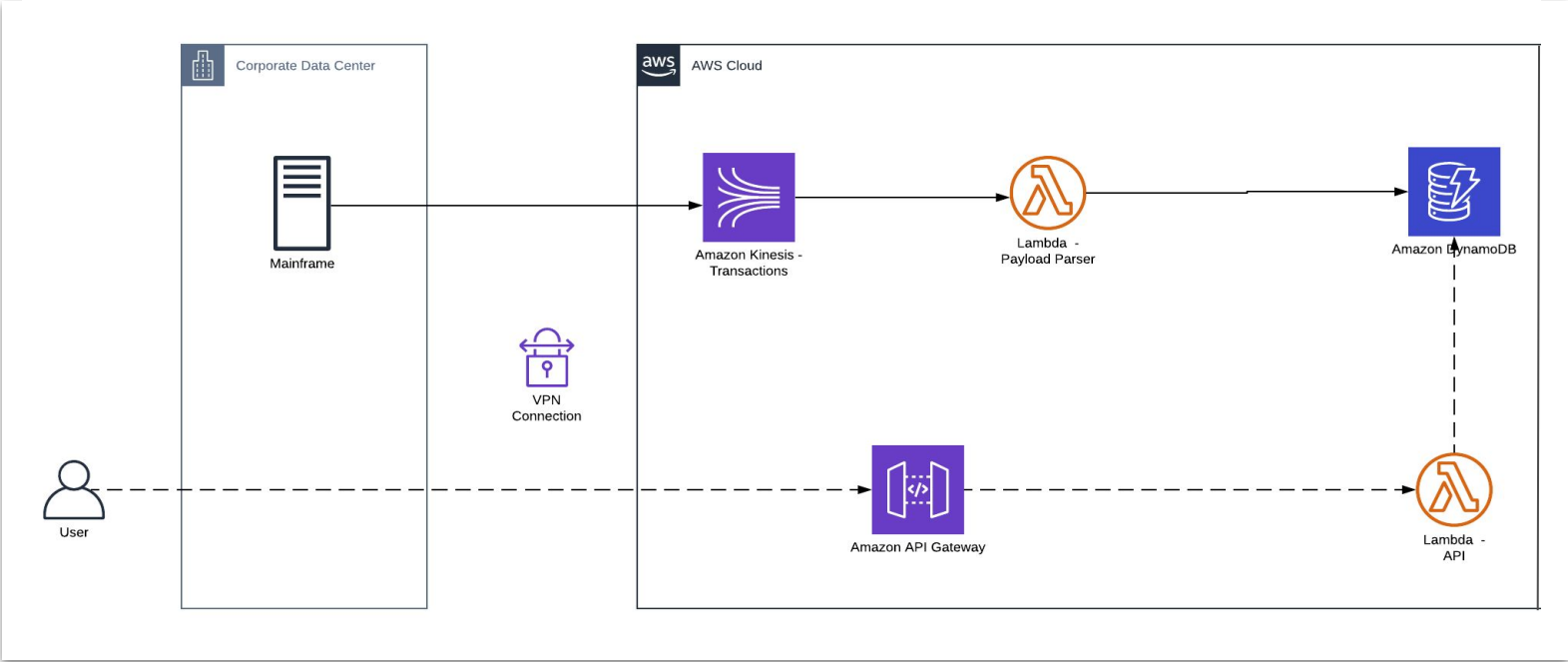
# The Architecture



# Live-Streaming Mainframe Data to the Cloud



# The Architecture (Extremely Simplified)



# Quick Primer on Some Core AWS Services



**AWS Lambda**



**Amazon API Gateway**



**Amazon Kinesis Data Firehose**



**Amazon Simple Notification Service**



**AWS Step Functions**



**AWS Snowball**



**AWS Fargate**



**Amazon Kinesis**



**Amazon Elasticsearch Service**



**Amazon Simple Queue Service**



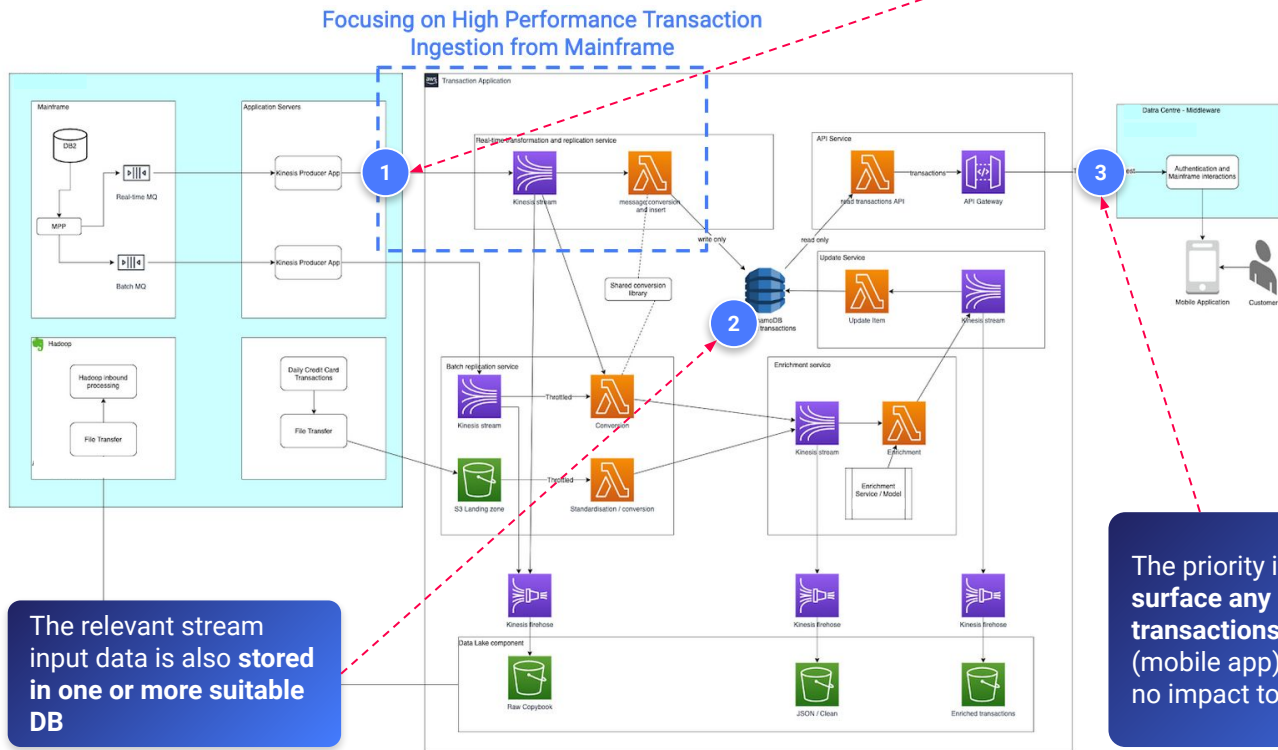
**Amazon RDS**



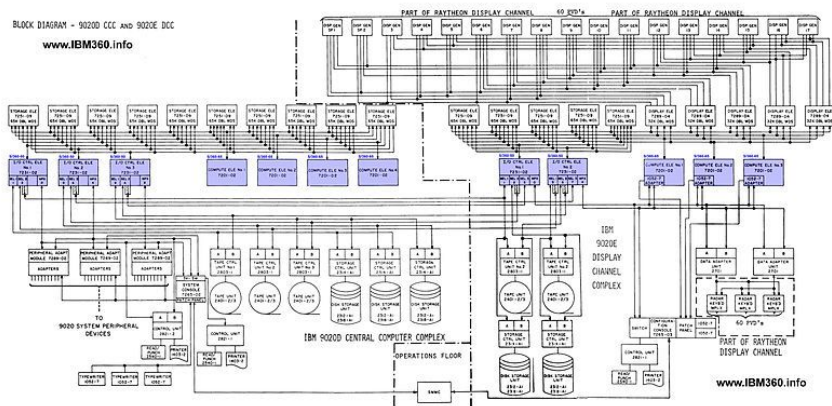
**Amazon DynamoDB**

# The Architecture (Not as Simple)

The focus is on providing **performance and accuracy** of the transactions stream initiated by a mainframe process



# Architecture at a Glance



- Serverless/Microservices architecture
- Lambda 80 %, Fargate 20%
- 3 Teams, multiple AWS accounts
- 9 Microservices, tens of functions
- Shared Kernel (lightly versioned / library)
- 120 rps API, latency <100 ms 95th percentile
- Event sourcing pattern with Data lake integration
- CI/CD with AWS services & Serverless Framework



The Actual  
Architecture... [here...](#)

# Microservices Boundaries



- First - Domain Model.  
Start identifying the aggregates
  - Second - Operational needs.  
Who's going to run this?  
It had already been decided!
- (Inverse Conway's Manoeuvre!)*

A composite image showing the Earth's horizon from space. The bottom half of the image shows the curved horizon of the Earth, with a bright orange and yellow glow on the left side, suggesting a sunrise or sunset. The top half of the image shows the dark, starry sky with the Milky Way galaxy visible as a dense band of stars and dust stretching across the frame.

# Architectural Highlights



# Ensuring Correctness

It's a distributed systems problem. This is financial transaction data; end users need correct values.

- Acceptable levels of service as SLOs (accuracy, speed, availability)
- Live End-To-End Testing System
- In-band and Out-of-Band Transaction Reconciliation System (TRS)



# Acceptable Levels of Service (NFRs)



1

Define expected **correctness**: “99% of Get API calls for transactions, over a minute, will return correct data at a “certain” moment in time”

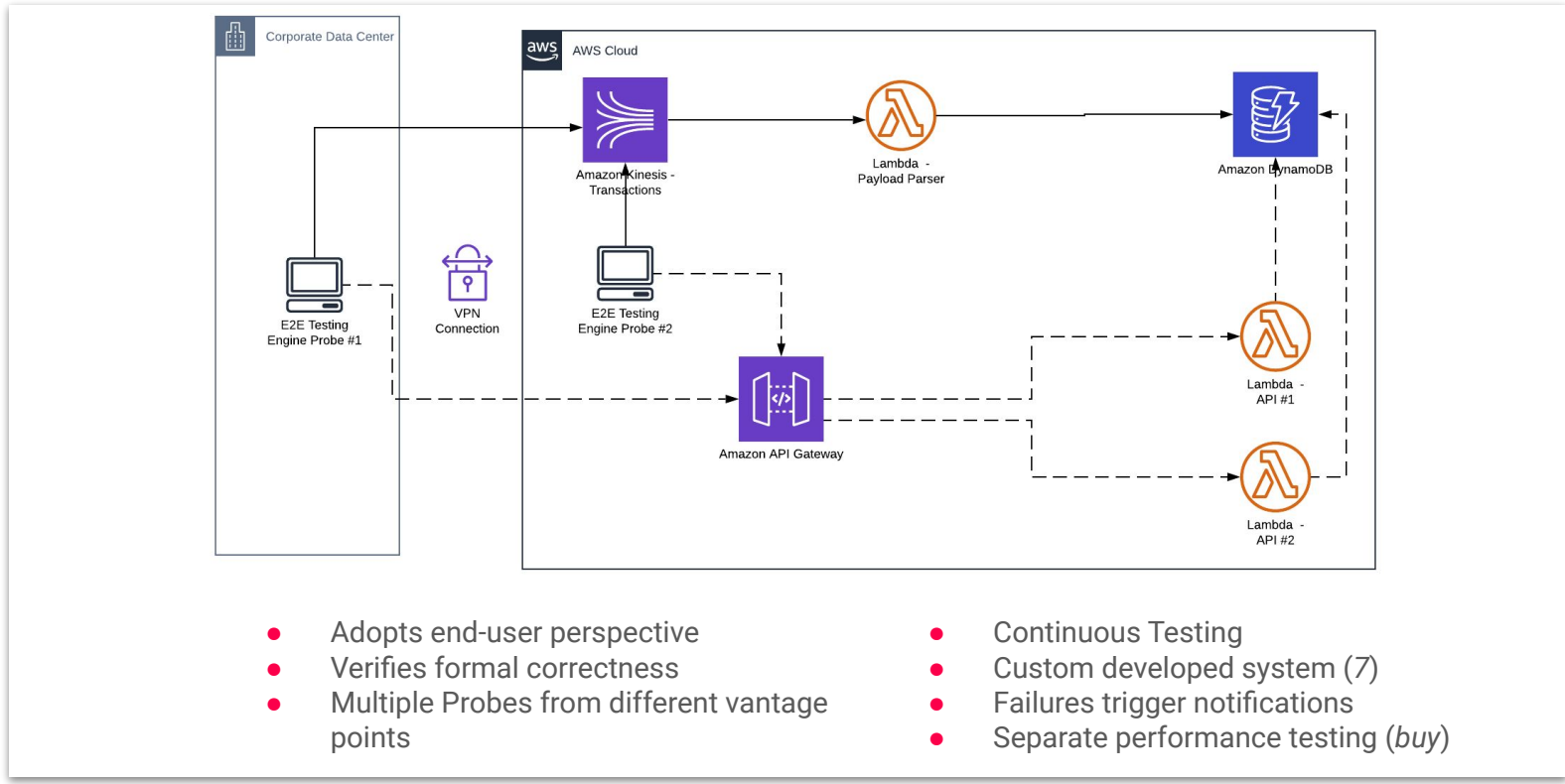
2

Define expected **consistency**: “99.99% of Get API calls for transactions, over a minute, will return correct data not-older-than 30 secs”

3

Define expected **latency**: “99% of Get API calls for transactions, over a minute, will return the most up-to-date data in less than 100ms”

# End-To-End Testing System



# Two Transaction Reconciliation Systems

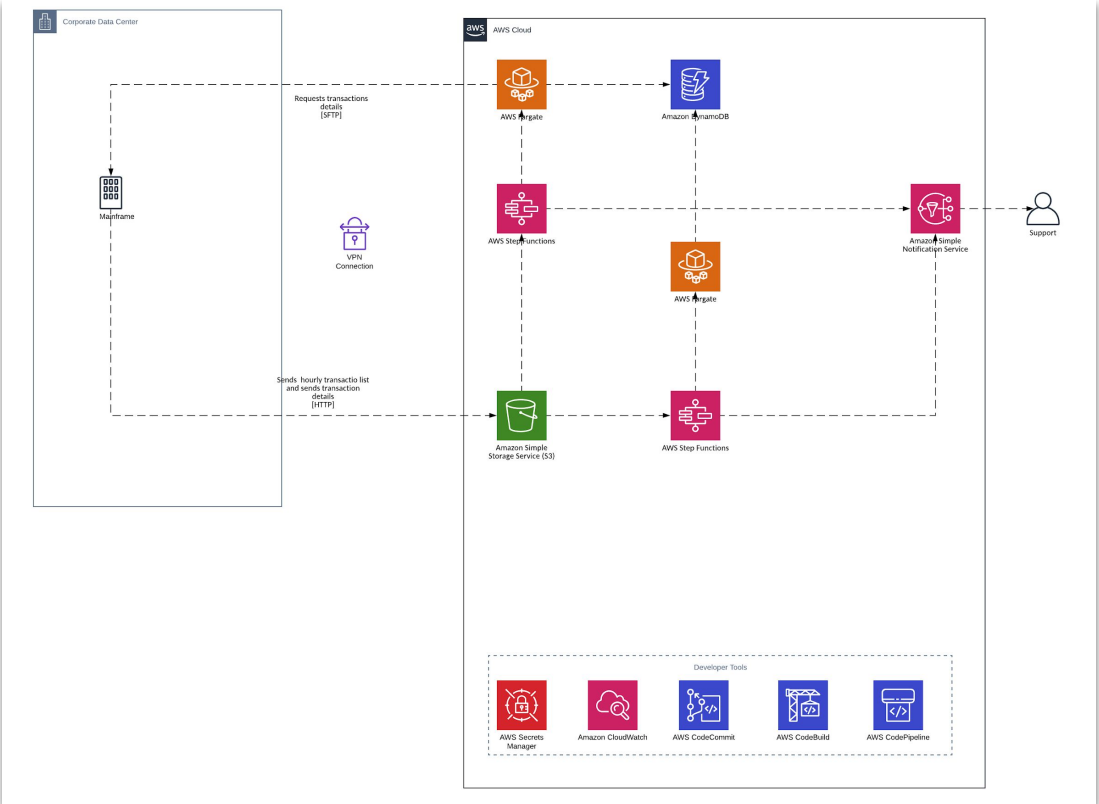
## In-Band

- Uses streaming transactional data
- Idempotency tags
- Retries on AWS side
- Repeats on the mainframe
- Monitoring and alerting enabled
- With High-Water mark

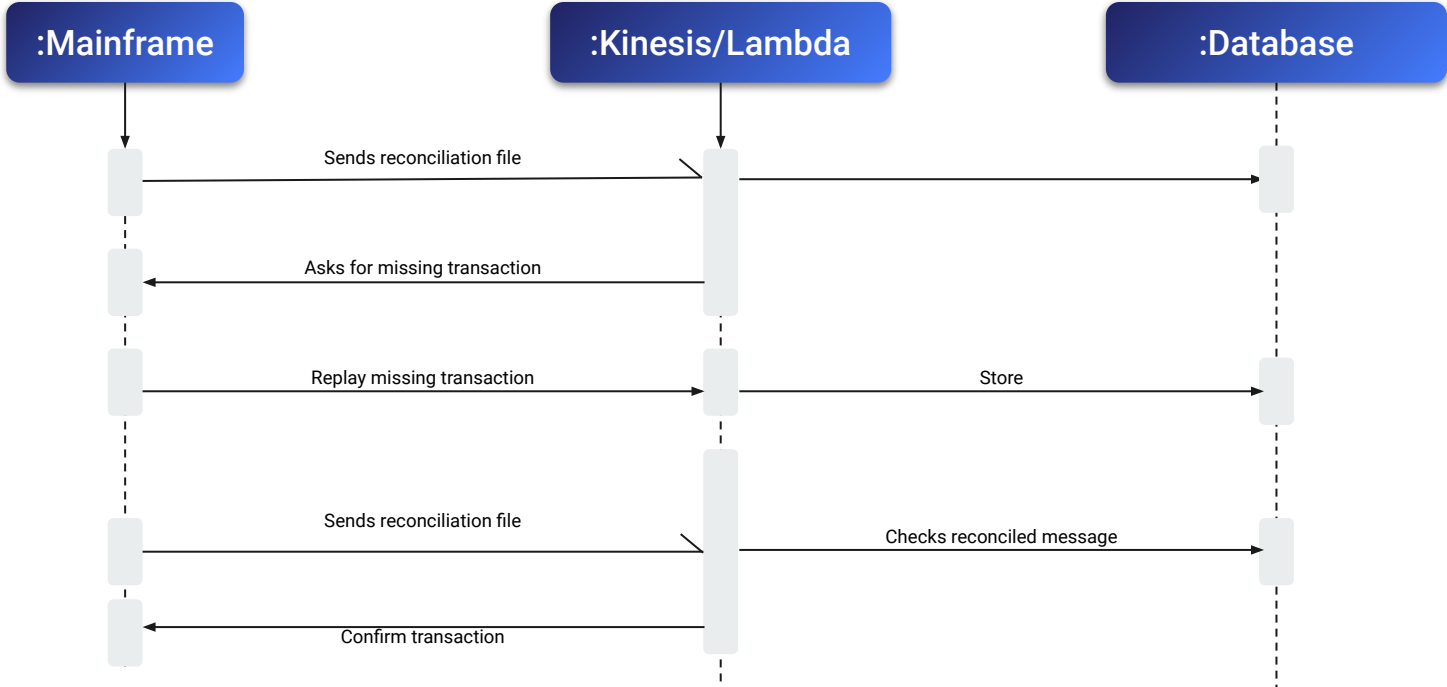
## Out-Of-Band

- Complex system with Fargate and AWS Step Functions
- Uses a different datasource (batch) and processes transactions in a different way
- With automated and manual resolution
- Monitoring and alerting enabled
- With High-Water mark

# Transaction Reconciliation System

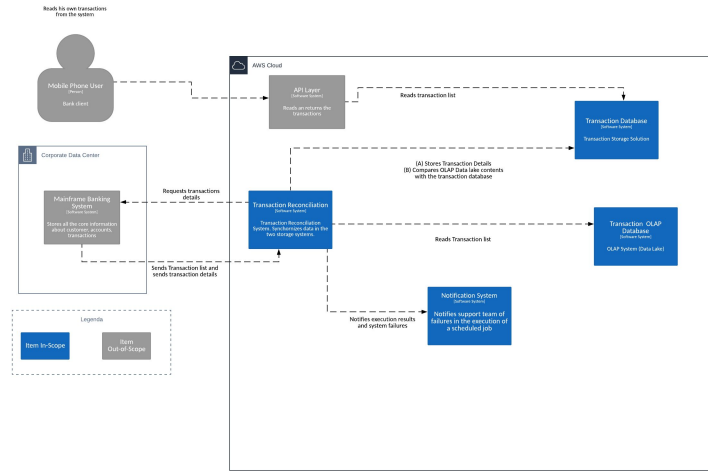


# Transaction Reconciliation System

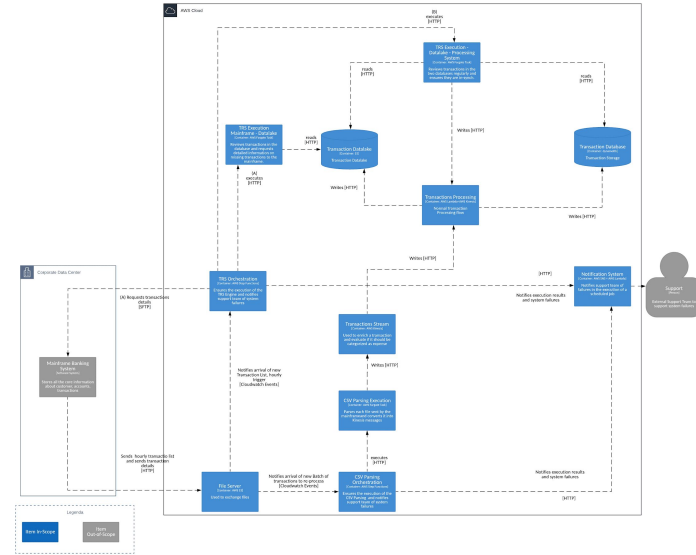


UML Sequence Diagrams

# Transaction Reconciliation System - C4 Model

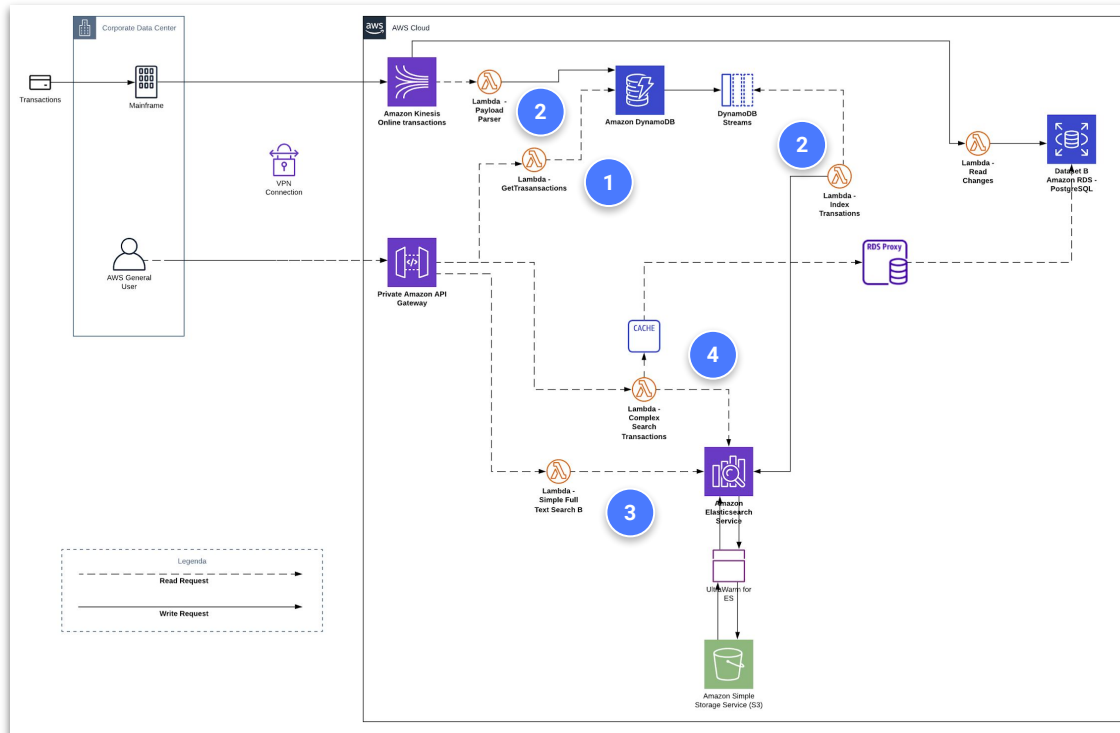


System Context



Container Model

# Full-Text and Advanced Search

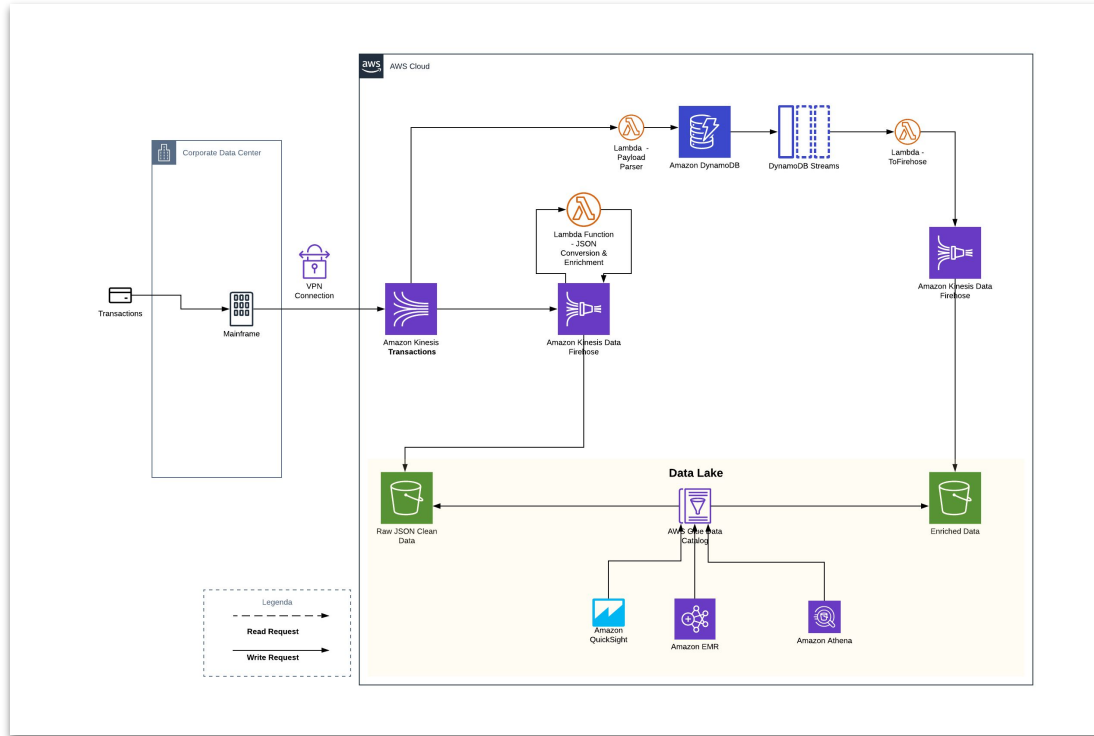


- Enables users to search their transactions
- Multi-Database : DynamoDB, RDS, Amazon ElasticSearch
- Uses DynamoDB Streams to replicate data to ES
- Ultrawarm for low-frequency data
- With RDS Proxy and Simple Cache
- Explicitly disregards API Gateway Velocity Templates

*NOTE: Simplified diagram*



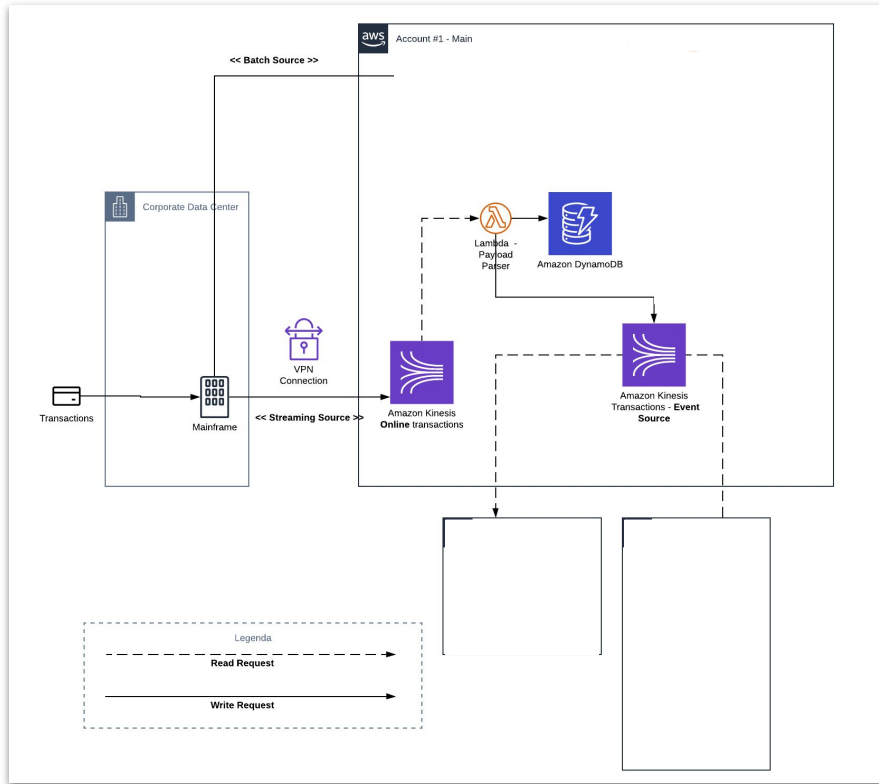
# Data Lake Integration



- Where value lies: liberating the data for wider consumption
- Easy “in team” data exploration
- Used for analytics
- Standardized out-of-team consumption (data stewardship and governance external)
- Raw and Enriched (clean) data available

**NOTE: Simplified diagram**

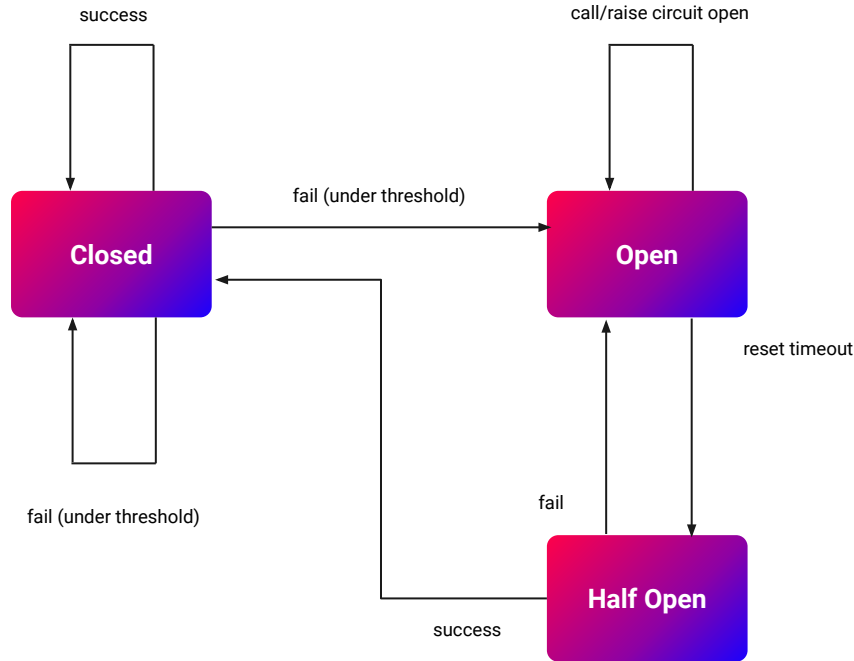
# Event Source



- Kinesis-based event source
- Follows financial transaction lifecycle
- Strong decoupling
- Useful to other live systems
- With idempotency tag
- Does not allow to replay history
- Contract testing with Pact

# Serverless Patterns

# Circuit Breaker Pattern



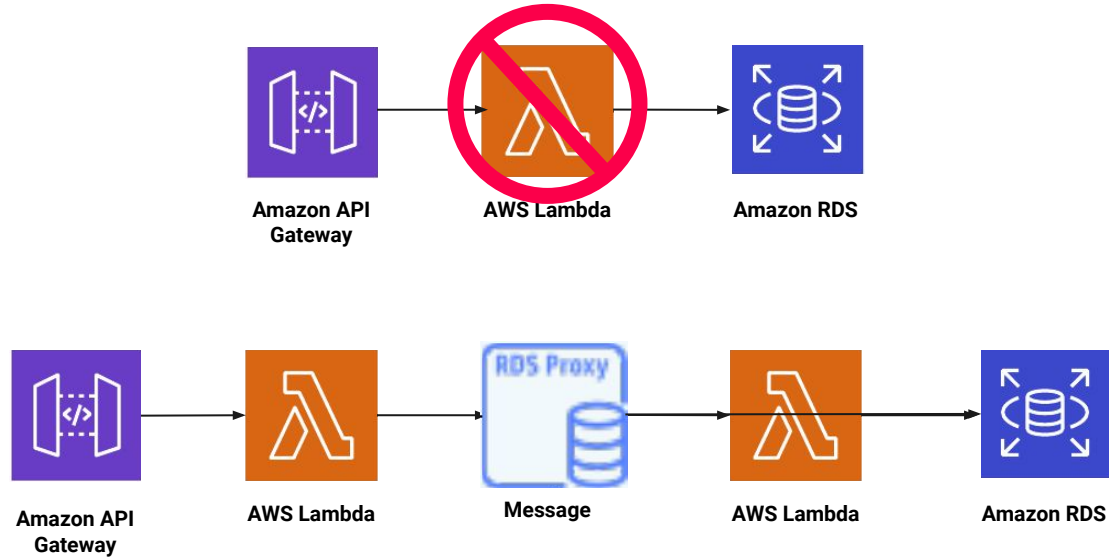
- Strong requirements for transaction consistency and availability
- Automated switch-off on high number of errors
- Mainframe system as fallback
- Prevents downstream service flooding
- Reduces resource consumption

# Controlling Function Scaling and Throttling



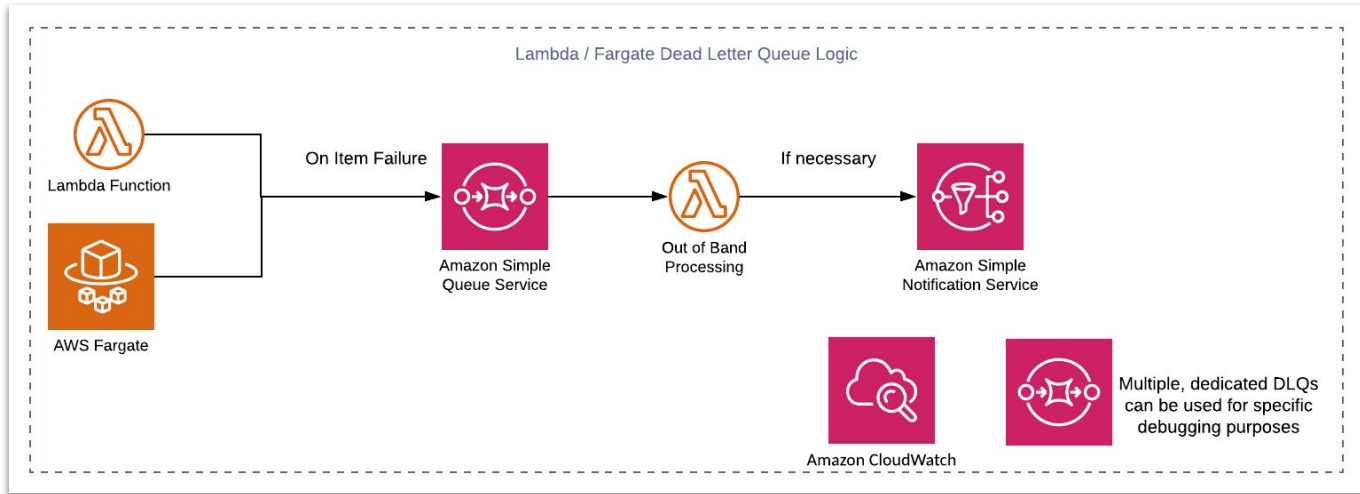
- Make sure to monitor account quotas
- Consider Provisioned Concurrency
- Set-up extensive monitoring in Cloudwatch (*Throttles, Invocations, ExecutionTime*)
- Look for tail-latency events

# Impact on Downstream Systems



Set Reserved Concurrent Executions

# Dead-Letter-Queue Pattern on Lambda



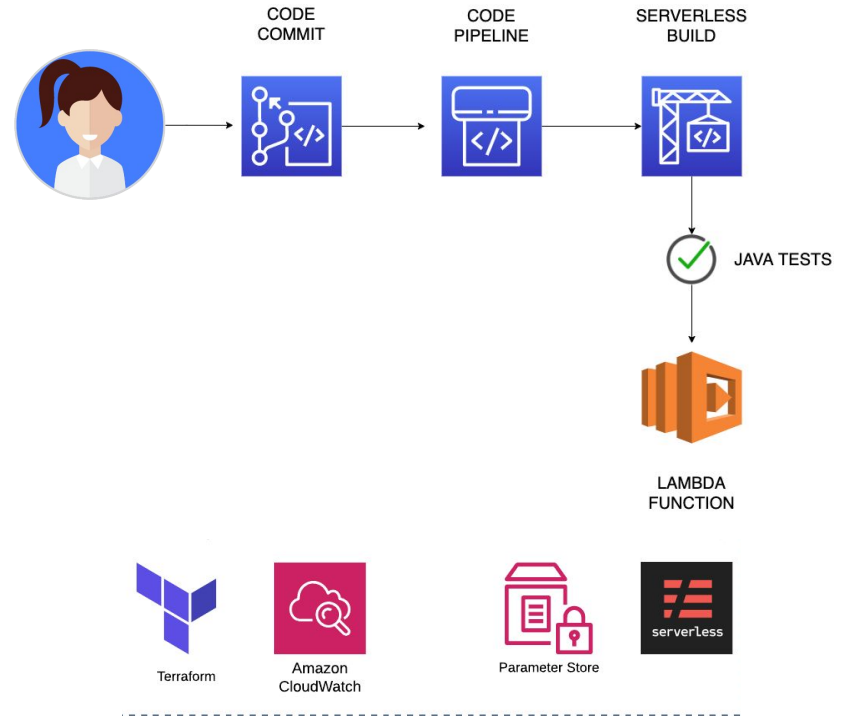
- Kinesis specific logic: Bisect-on-failure and multiple retries
- Aware of Lambda context
- Can have performance impact on the overall system
- Failures trigger notifications
- Lambda Destinations unsuitable

# Developer Experience



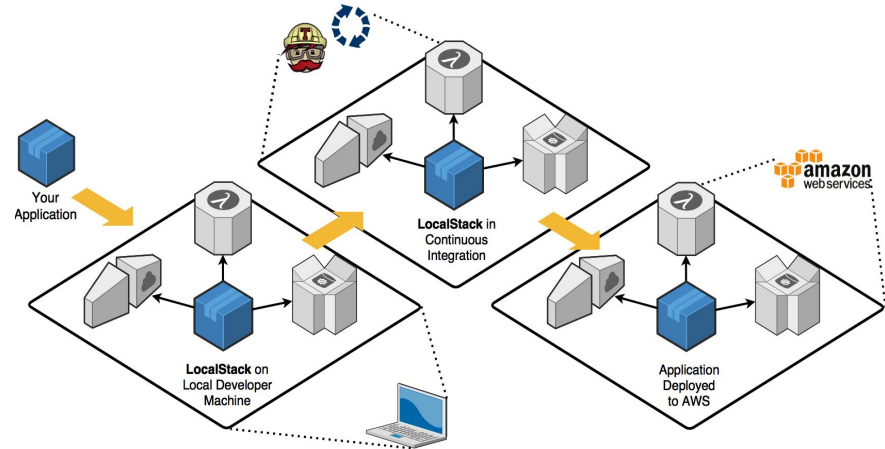
# CI/CD

- Using AWS Developer Tools
- Monorepo for functions, continuous deployment multi-branch pipeline
- Backed by build automation tools and integrated with existing on-prem solutions
- **Serverless Framework** core tool for Lambda Functions
- Terraform for infra resources
- Unit & Integrations Tests running in the pipeline with artefact promotion
- Canary Deployment with gradual traffic rollout
- Supports manual approvals for compliance, where needed



# Testing & Debugging Serverless

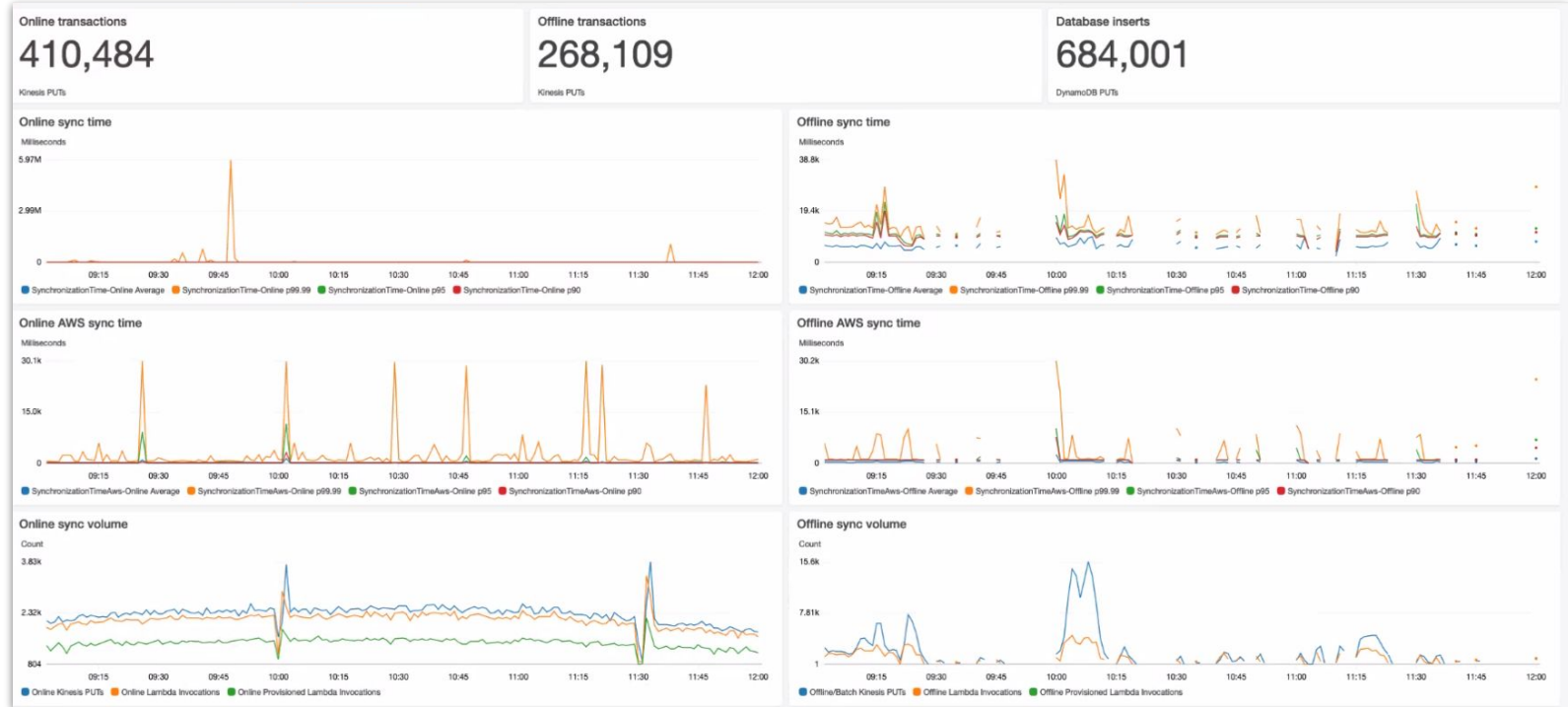
- Used [Localstack](#),
- Integrated with IntelliJ using AWS toolkit to test/debug Lambda's
- Make ample use of AWS X-Ray & AWS Cloudwatch
- Testing in production with feature flags and Canaries (Serverless framework)



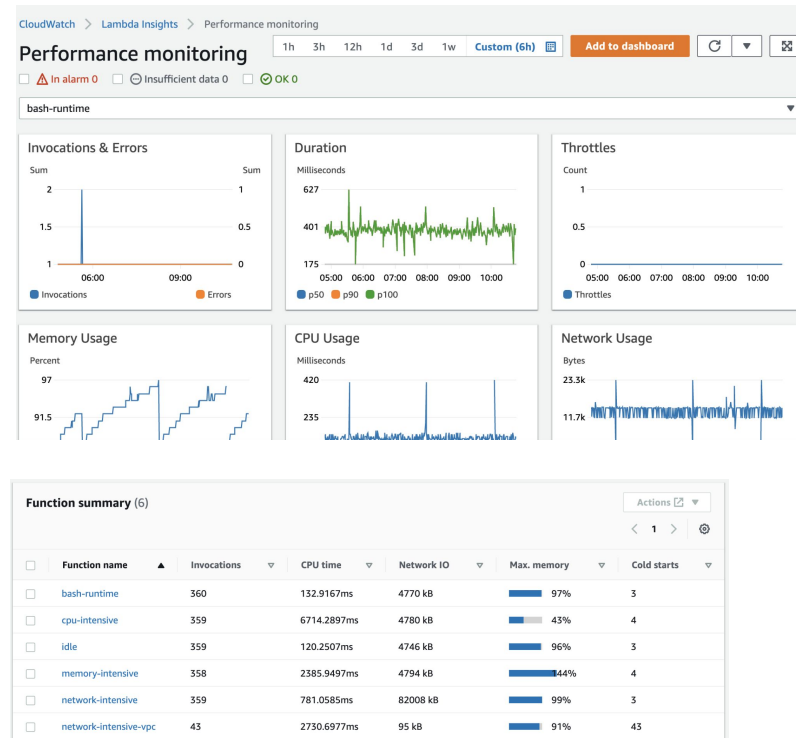
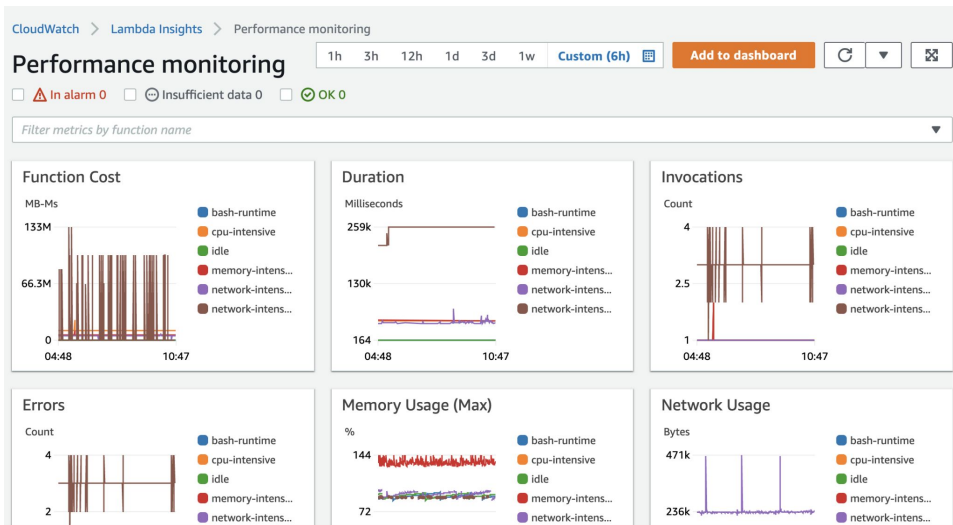
# Performances



# Results



# From numbers to percentages

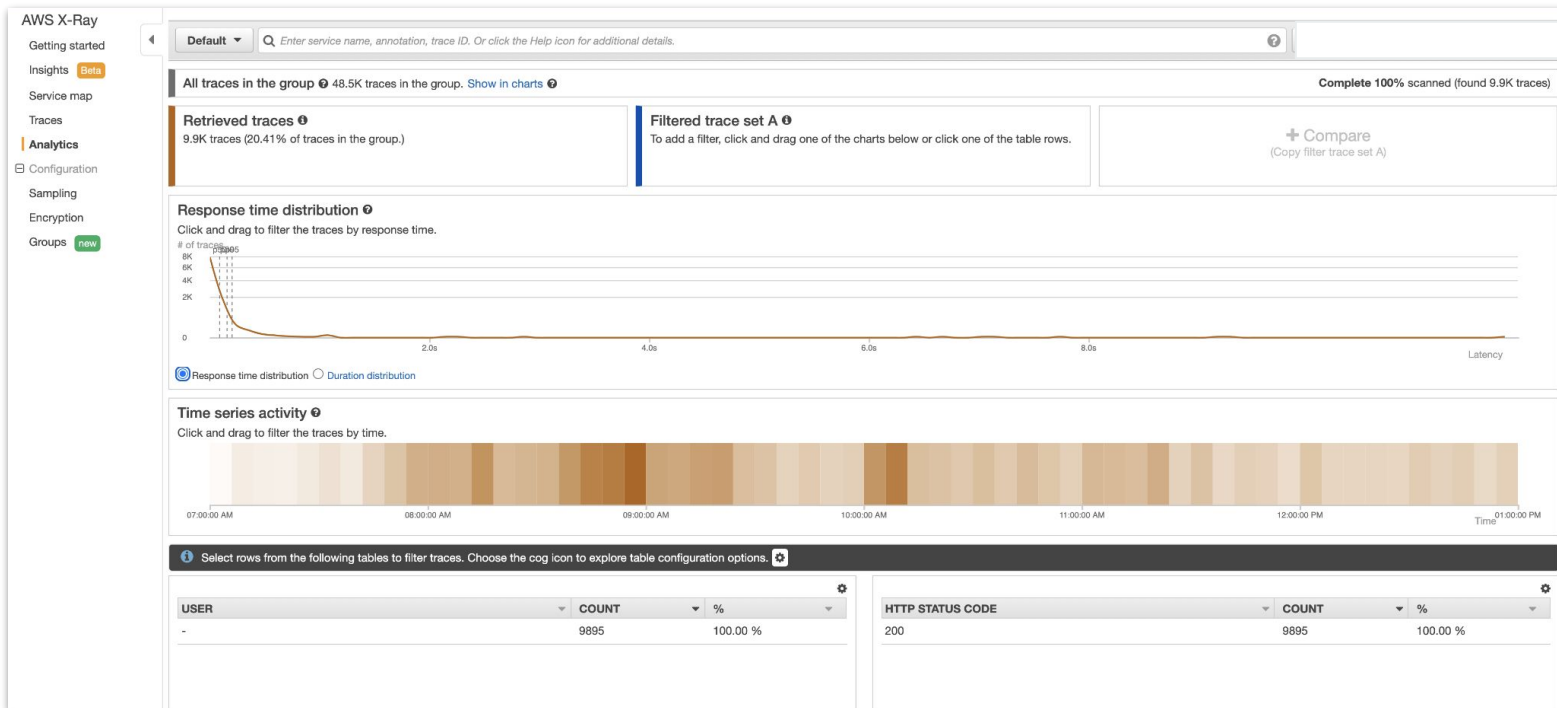


# Measuring for real

- Distributed systems hard to tackle issue
- Lambda's NTP servers cannot be configured
- Lambda **should run on** Stratum 3 NTP servers
- Run a small test using Lambda & Python
- With DC<->AWS RTT at 40 ms, clock offset < 10ms
- No guarantees are made



# Response Time



Response time for the system p95 < 100ms , p50 ~60 ms

# Example Java Optimizations on

## Cold-Start

- **Optimize** libraries import
- Use **Provisioned Concurrency** (> cost, >complexity, >deployment time)
- X-Ray **adds** non-negligible **latency** (Cloudwatch Embedded Metrics)
- **Reduce** artifact Size
- Monitor the **p99**

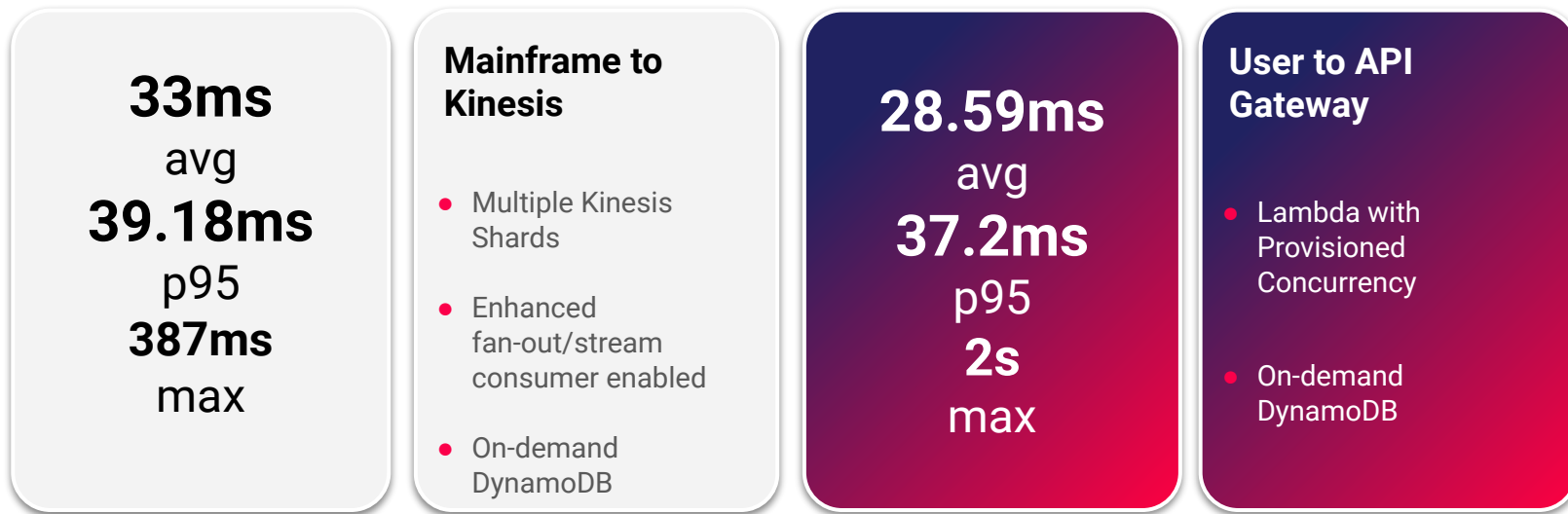
## Execution

- Tune to the **right** memory allocation
- Move as much as possible in the init phase (DynamoDB or 3rd party client)
- Fetch secrets at **init** time and handle failures
- Watch your **framework**
- Monitor the **average** and tail latency



# Performance samples...

Below are the results from our load testing activity, along with the adjustments we made in order to keep response times as low as possible.



# Takeaway System Optimizations

1 Java **can** be optimized for low latency application on Lambda

2 **Enhanced Fan-Out** increases throughput and **reduces latency**

3 Use **Lambda Destinations** (Bisect on Function Error)

4 **Provisioned Concurrency** is a bit expensive but significantly helps controlling the cold-start problem

5 Kinesis batch size tuning **had more impact on avg latency** than JVM optimizations

# Cloud Operating Model

# Serverless



**On-Premises**



**IaaS**  
Infrastructure as a Service



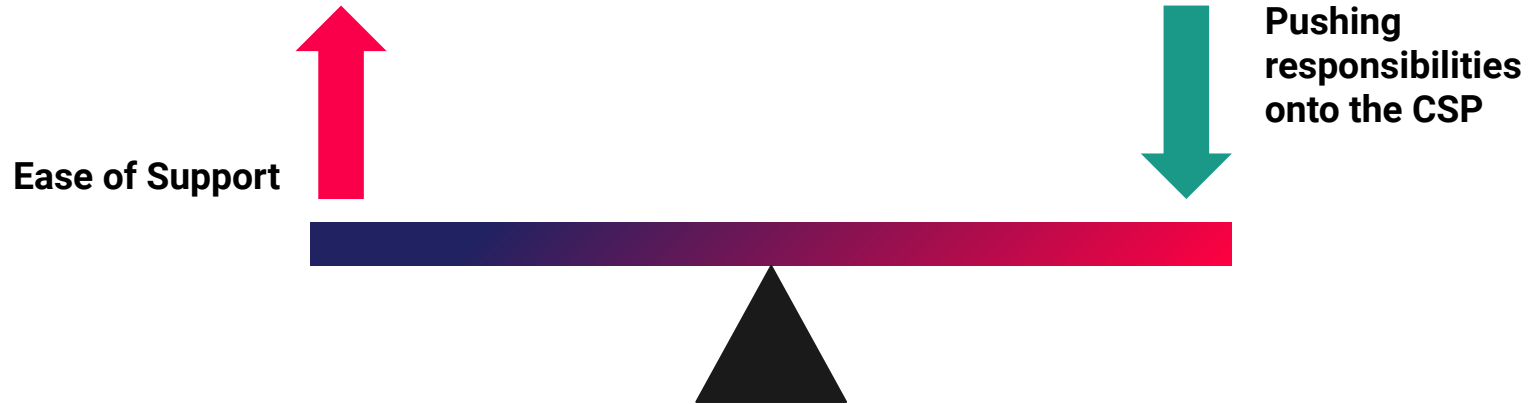
**PaaS**  
Platform as a Service



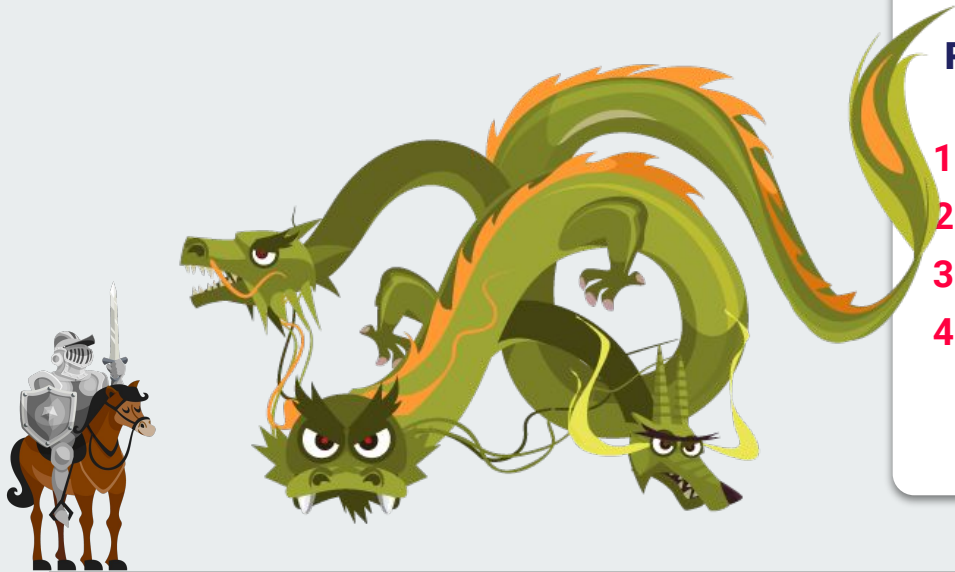
**SaaS**  
Software as a Service

Application	Application	Application	Application	You Manage
Data	Data	Data	Data	Others Manage
Runtime	Runtime	Runtime	Runtime	
Middleware	Middleware	Middleware	Middleware	
O/S	O/S	O/S	O/S	
Virtualization	Virtualization	Virtualization	Virtualization	
Servers	Servers	Servers	Servers	
Storage	Storage	Storage	Storage	
Networking	Networking	Networking	Networking	

# Fully Serverless



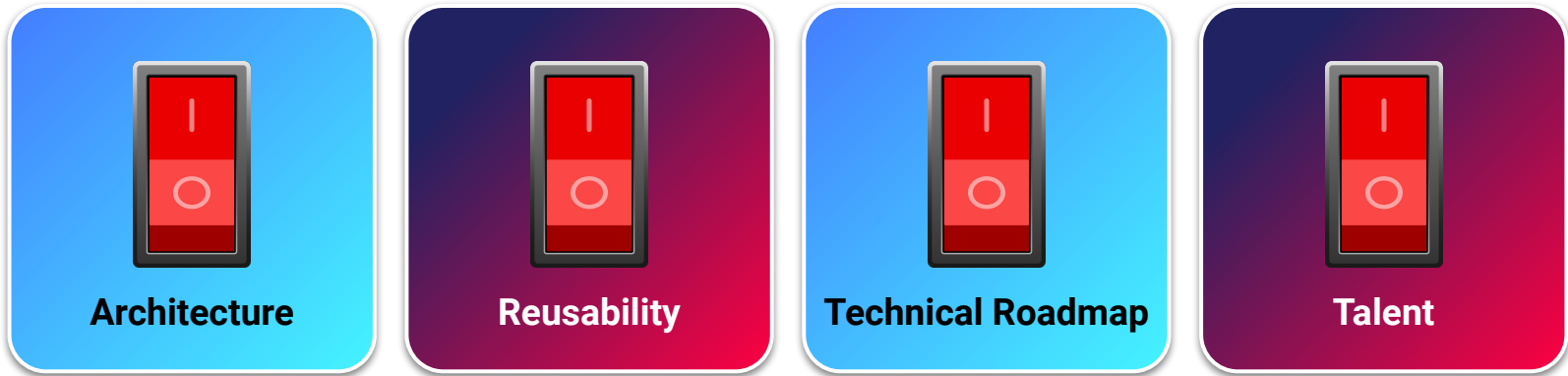
Leapfrogging cloud transformation, neglecting the cultural and organizational changes results in **higher risk** and **more complexity**.



### Running before walking

1. Completely **different** operating model
2. **Low levels** of engineering maturity
3. **Complex** technical solution
4. **No man's land** of compliance and internal processes

# Switches



At low levels of engineering maturity we can't break the link between cloud architecture and operating model



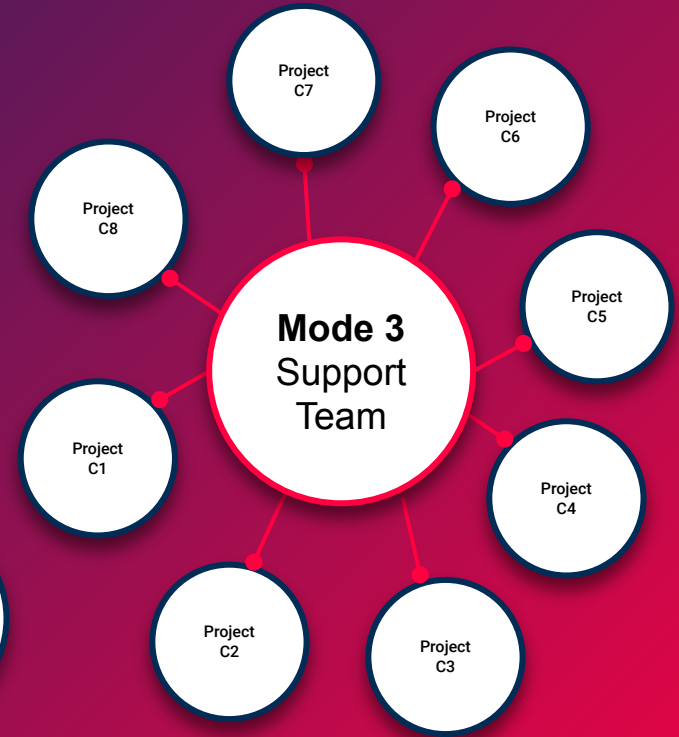
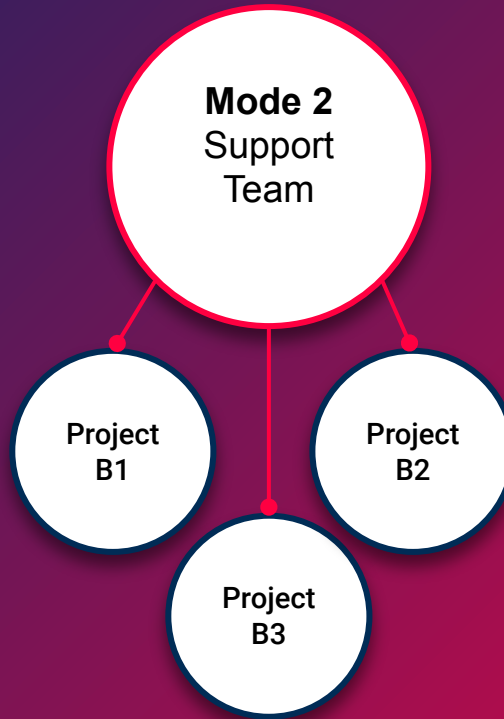
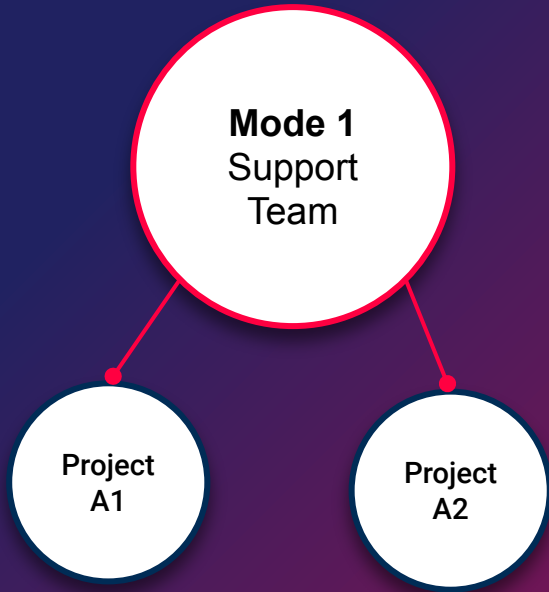


# Two Key Questions

*“What is our level of engineering maturity ?”*

*“What operating model are we going to use?”*

# Limiting the Support ratio...



# CONTINO

# Thank You

**London**

[london@contino.io](mailto:london@contino.io)

**New York**

[newyork@contino.io](mailto:newyork@contino.io)

**Atlanta**

[atlanta@contino.io](mailto:atlanta@contino.io)

**Melbourne**

[melbourne@contino.io](mailto:melbourne@contino.io)

**Sydney**

[sydney@contino.io](mailto:sydney@contino.io)

**Brisbane**

[brisbane@contino.io](mailto:brisbane@contino.io)



[www.contino.io](http://www.contino.io)



[continohq](https://twitter.com/continohq)



[contino](https://www.linkedin.com/company/contino)