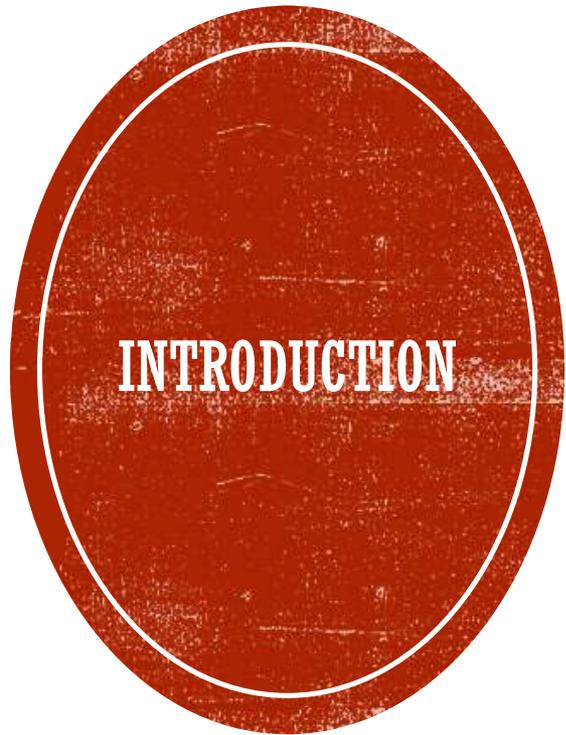




# **AWS Data Lake Proof Of Concept**



---

**Project outline  
and business case**

---

**Legacy  
Technology Stack**

---

**AWS Proof of  
Concept**

---

**Future Roadmap**



**INCREASING  
BUSINESS DEMAND**



**COST**



**FUNCTIONALITY**



**PERFORMANCE**

# **PROJECT OUTLINE AND BUSINESS CASE**

# LEGACY TECHNOLOGY STACK

- Changing business requirements
- SAS Resourcing challenges
- Too costly to maintain and difficult to scale infrastructures
- Complex legacy business rules
- Lack of Documentation



# **AWS PROOF OF CONCEPT**

- Asked for 6 weeks POC using AWS but delivered within 4 weeks.
- Serverless data lake
- 5 new dashboards
- Full automation



Find the right technology partner/consultants



Make sure you have IT supports



Discuss with Stakeholder early



Need data SME (Subject Matter Expert)

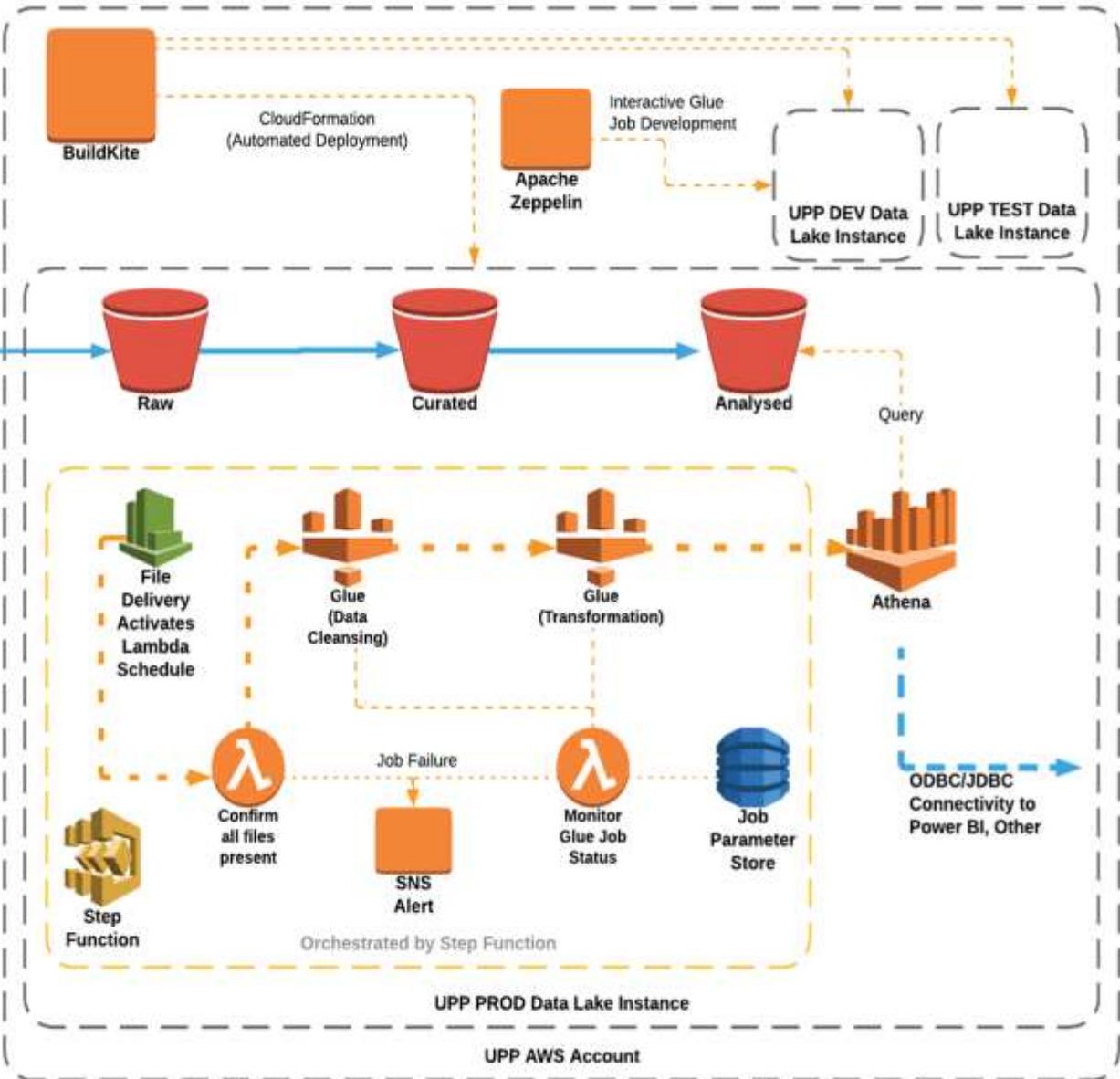
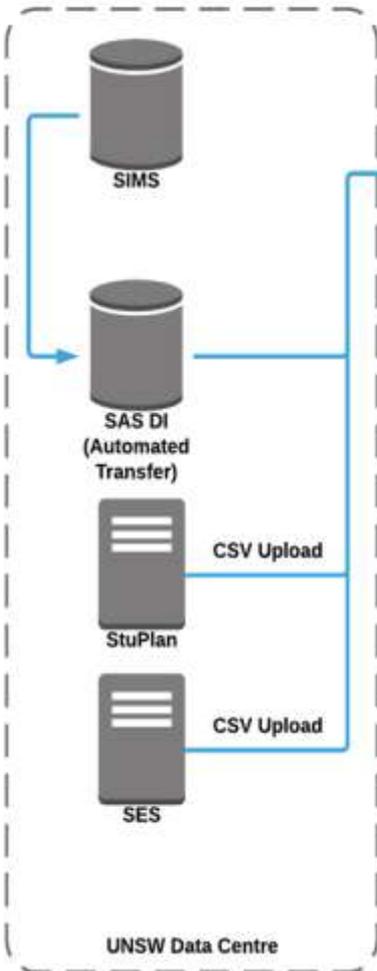


Document your success and share it

## LESSONS LEARNED

# UNSW Planning & Performance AWS Data Lake

Data Flow      Processing Flow



# SOLUTION COMPONENT BREAKDOWN

## ■ AWS Native

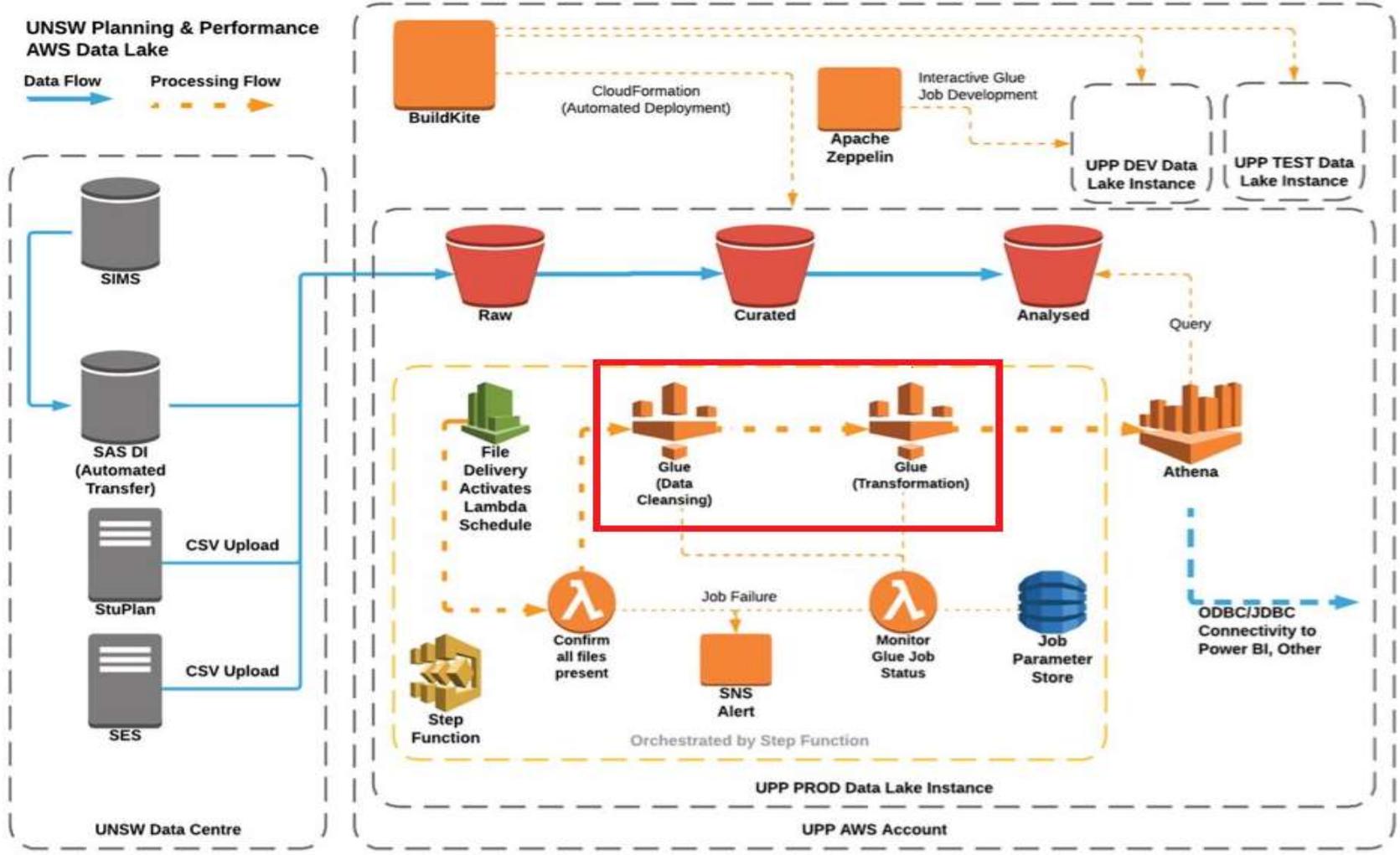
- Glue
- Athena
- S3
- Lambda
- Step Function
- DynamoDB
- SNS
- CloudWatch
- CloudFormation

## ■ 3rd Party Tooling

- BitBucket
- BuildKite
- Zeppelin Notebook



# GLUE



# GLUE



- Glue is a serverless ETL service for batch workload.

A screenshot of the AWS Glue console showing a list of tables. The table has columns for Name, Database, Location, Classification, and Last updated. The data is as follows:

Name	Database	Location	Classification	Last updated
glue_catalog_table	glue_catalog_db	s3://aws-logs-912461231913-us-east-1-logs-1/2018-11-15/	Log	17 November 2018 10:26:00 PT
glue_catalog_table	glue_catalog_db	s3://aws-logs-912461231913-us-east-1-logs-1/2018-11-15/	Log	17 November 2018 10:26:00 PT
glue_catalog_table	glue_catalog_db	s3://aws-logs-912461231913-us-east-1-logs-1/2018-11-15/	Log	17 November 2018 10:26:00 PT
glue_catalog_table	glue_catalog_db	s3://aws-logs-912461231913-us-east-1-logs-1/2018-11-15/	Log	17 November 2018 10:26:00 PT

## Glue Catalog

Replaced Hive Metastore, stores your data's metadata

A screenshot of the AWS Glue console showing a list of databases. The table has columns for Name, Location, Size, Type, Last updated, and Status. The data is as follows:

Name	Location	Size	Type	Last updated	Status
glue_catalog_db	s3://aws-logs-912461231913-us-east-1-logs-1/2018-11-15/	10 MB	Log	17 November 2018 10:26:00 PT	Active

## Glue Crawler

Crawls new file that comes in and creates table/partitions in the glue catalog

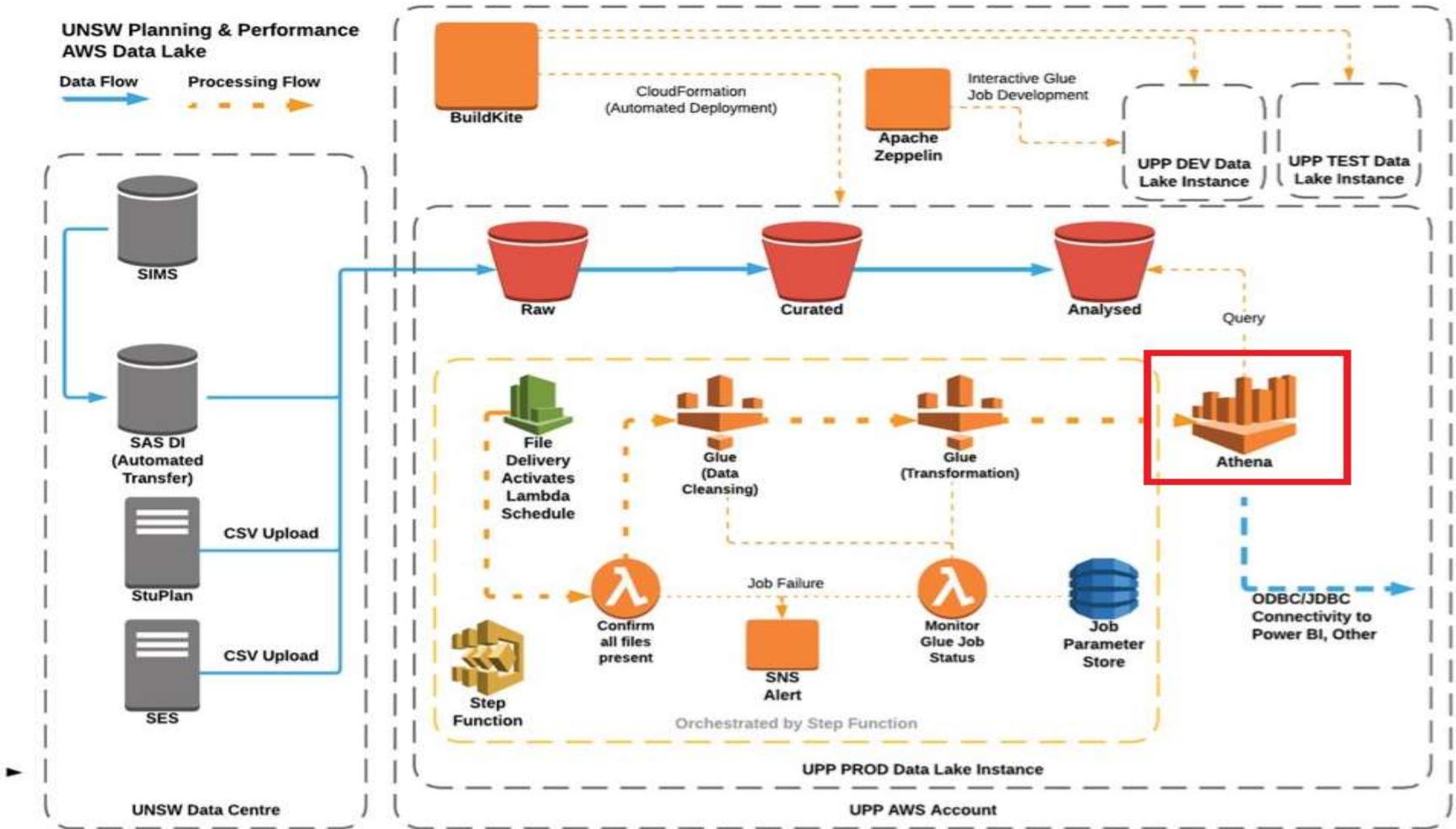
A screenshot of the AWS Glue console showing a list of crawlers. The table has columns for Name, Status, Type, Last updated, and Last run. The data is as follows:

Name	Status	Type	Last updated	Last run
glue_crawler	Running	Log	17 November 2018 10:26:00 PT	17 November 2018 10:26:00 PT

## Glue Job

Runs Spark Code in either Scala or Python

# ATHENA



# ATHENA

Interactive query service that makes it easy to analyze data in Amazon S3 using standard SQL.

Query using AWS console or your favourite thick client using JDBC/ODBC connection

Build in integration with popular BI tools like Tableau/PowerBI

Serverless

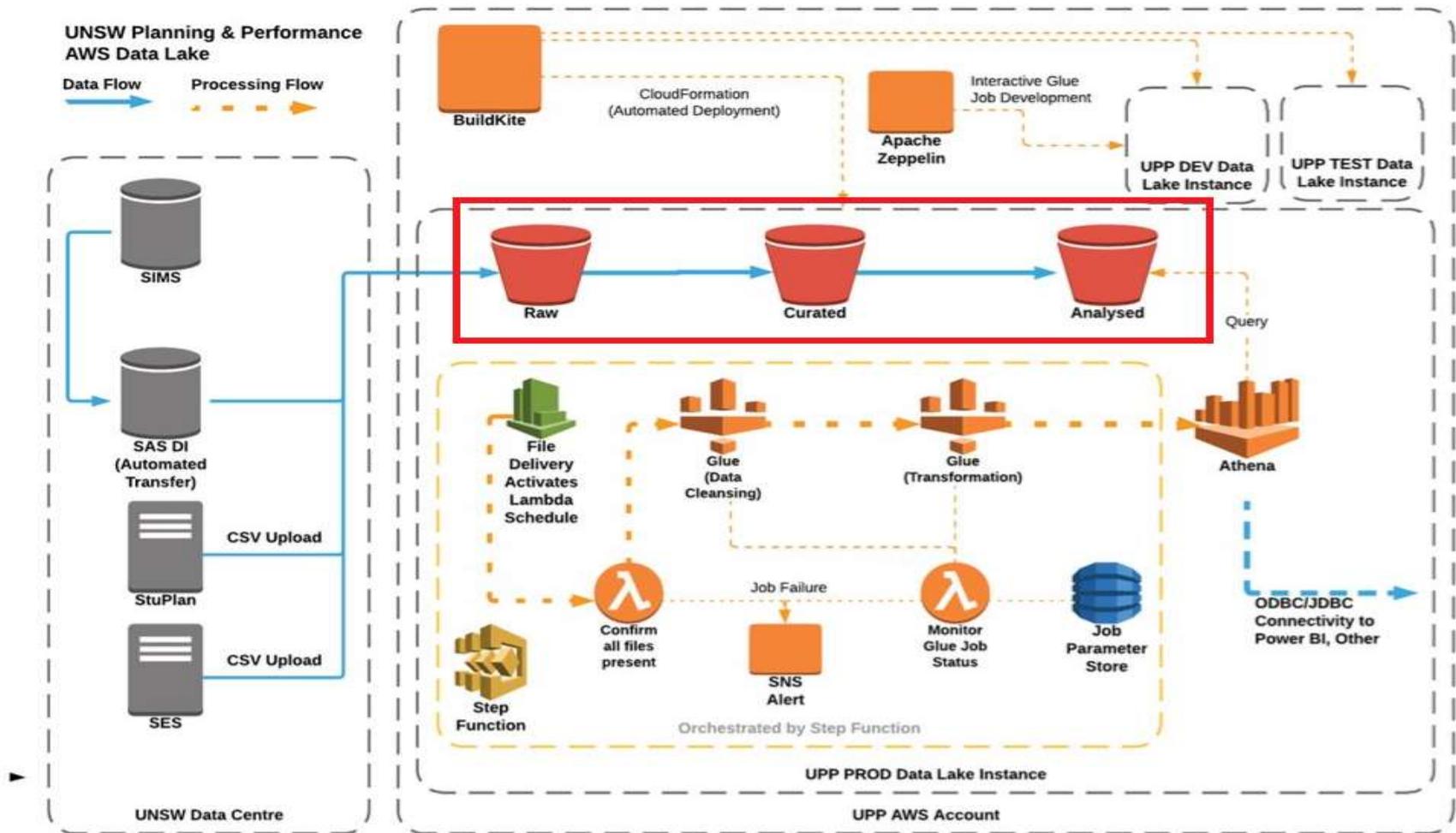


The screenshot displays the AWS Athena console interface. At the top, there are tabs for "New query 1" and "New query 2". The main area contains a SQL query with various clauses including "with", "select", "sum", and "case when". Below the query, there are buttons for "Run query", "Save as", "Create", and "Format query", along with a "Clear" button. The "Run query" button is highlighted in blue. Below the buttons, the console shows the execution time and data scanned: "(Run time: 27.02 seconds, Data scanned: 101.86 MB)".

The "Results" section at the bottom shows a table with the following columns: event\_date, term, inst\_doe, faculty\_description, career, degree\_type, school, program\_id, program\_description, and application\_stage\_group. Two rows of data are visible:

	event_date	term	inst_doe	faculty_description	career	degree_type	school	program_id	program_description	application_stage_group
1	2017-05-17 00:00:00.000	Semester 2 2017	Domestic	UNSW Business School	UGRD	Bachelor Degree (Pass)	UNSW Business School	3379	Information Systems	Application
2	2016-01-18 00:00:00.000	Semester 1 2016	Domestic	UNSW Business School	UGRD	Bachelor Degree (Pass)	UNSW Business School	3367	Commerce / Computer Science	Application

# SIMPLE STORAGE SERVICE (S3)



# SIMPLE STORAGE SERVICE (S3)

- ⑩ Object Storage
- ⑩ Unlimited on size
- ⑩ Serverless
- ⑩ Cheap
- ⑩ Resilient (99.999999999% Durability of objects)

## Storage pricing

Region: Asia Pacific (Sydney) ▾

	Pricing
<b>S3 Standard Storage</b>	
First 50 TB / Month	\$0.025 per GB
Next 450 TB / Month	\$0.024 per GB
Over 500 TB / Month	\$0.023 per GB
<b>S3 Standard-Infrequent Access (S3 Standard-IA) Storage</b>	
All storage / Month	\$0.019 per GB
<b>S3 One Zone-Infrequent Access (S3 One Zone-IA) Storage</b>	
All storage / Month	\$0.0152 per GB
<b>S3 Glacier Storage</b>	
All storage / Month	\$0.005 per GB

- Allows event driven workflow by S3 event.

# DYNAMODB

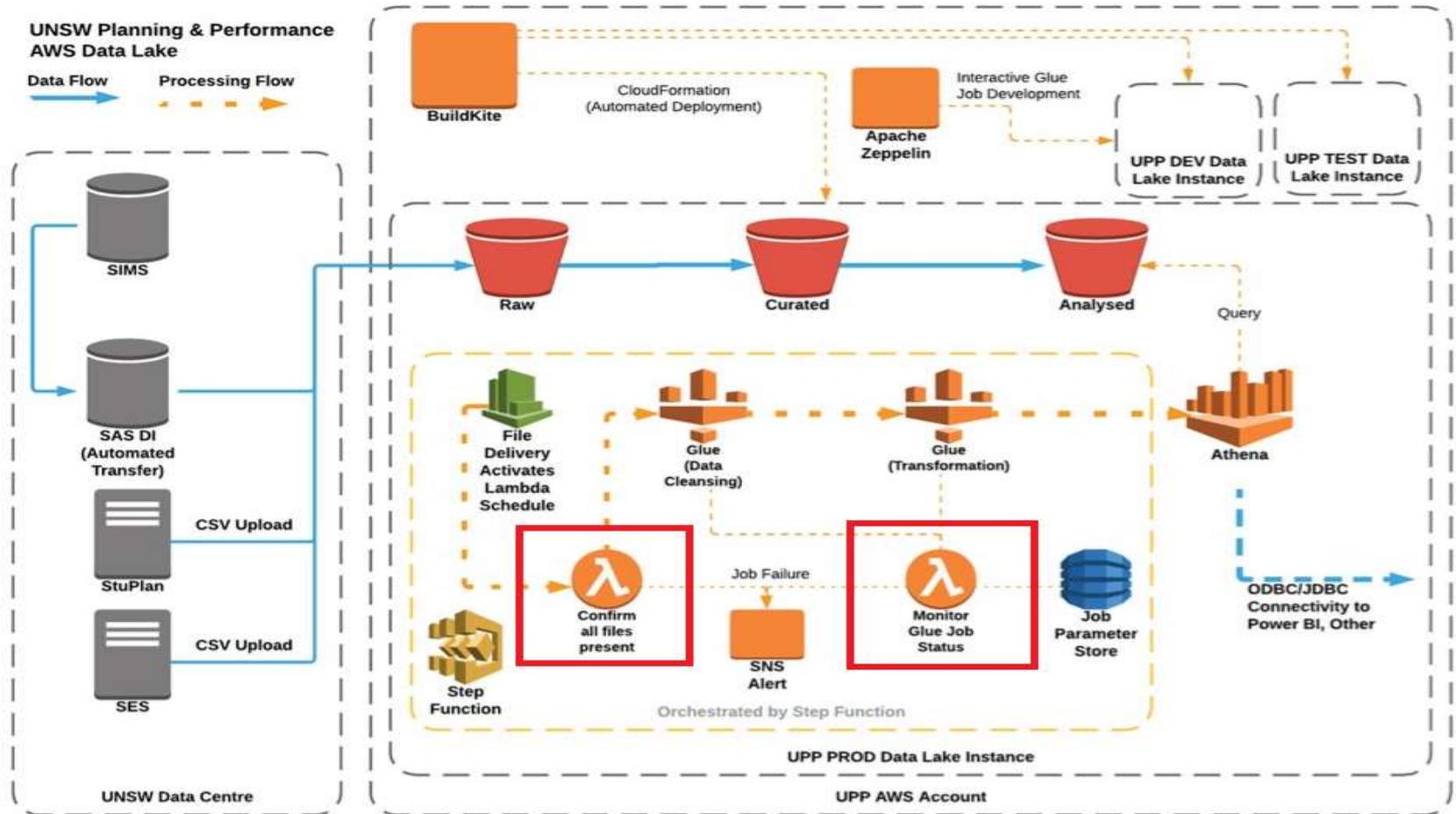
Serverless No-SQL DB, key-value storage

We use this for state storage for our serverless services



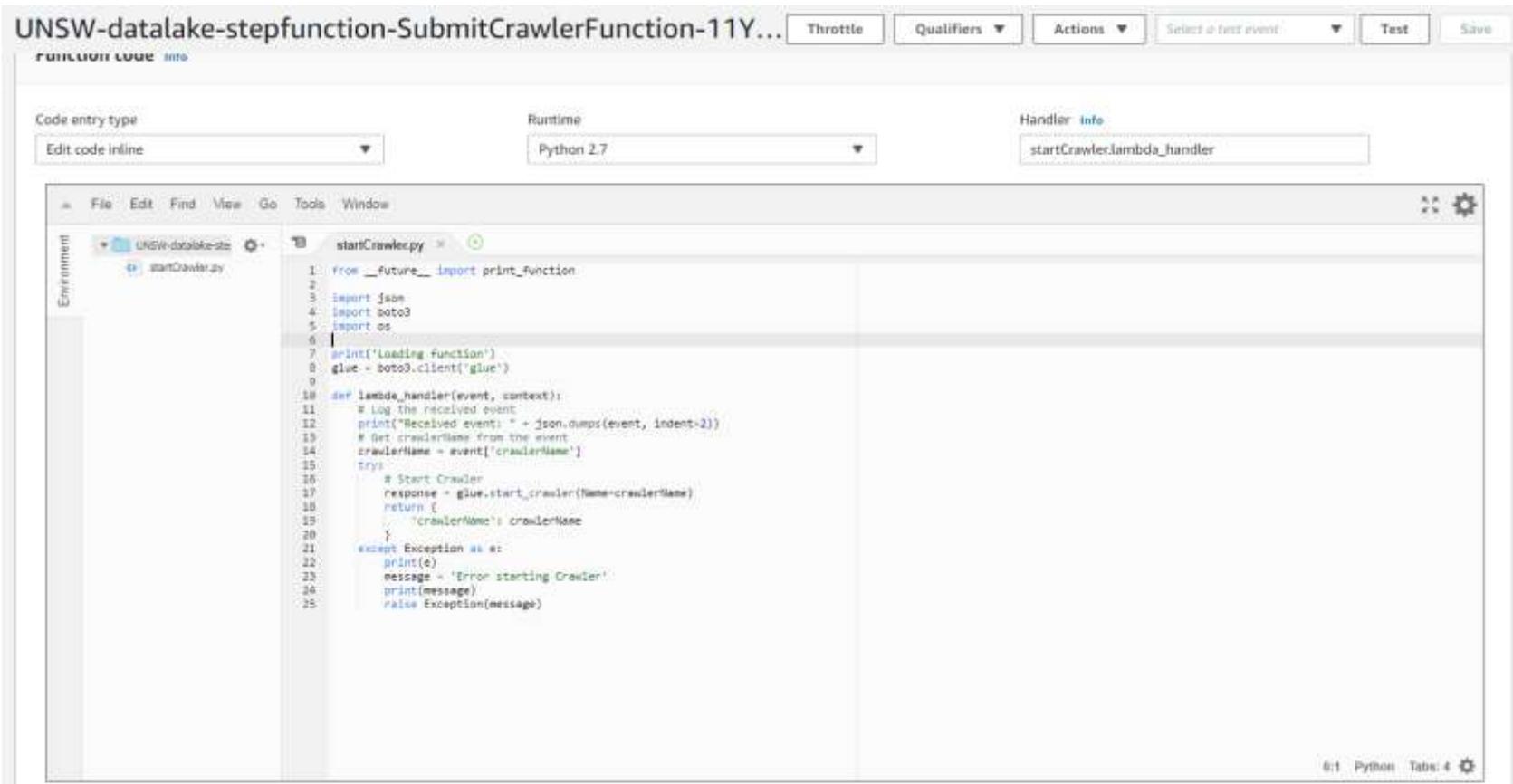
**amazon**  
DynamoDB

# LAMBDA



# LAMBDA

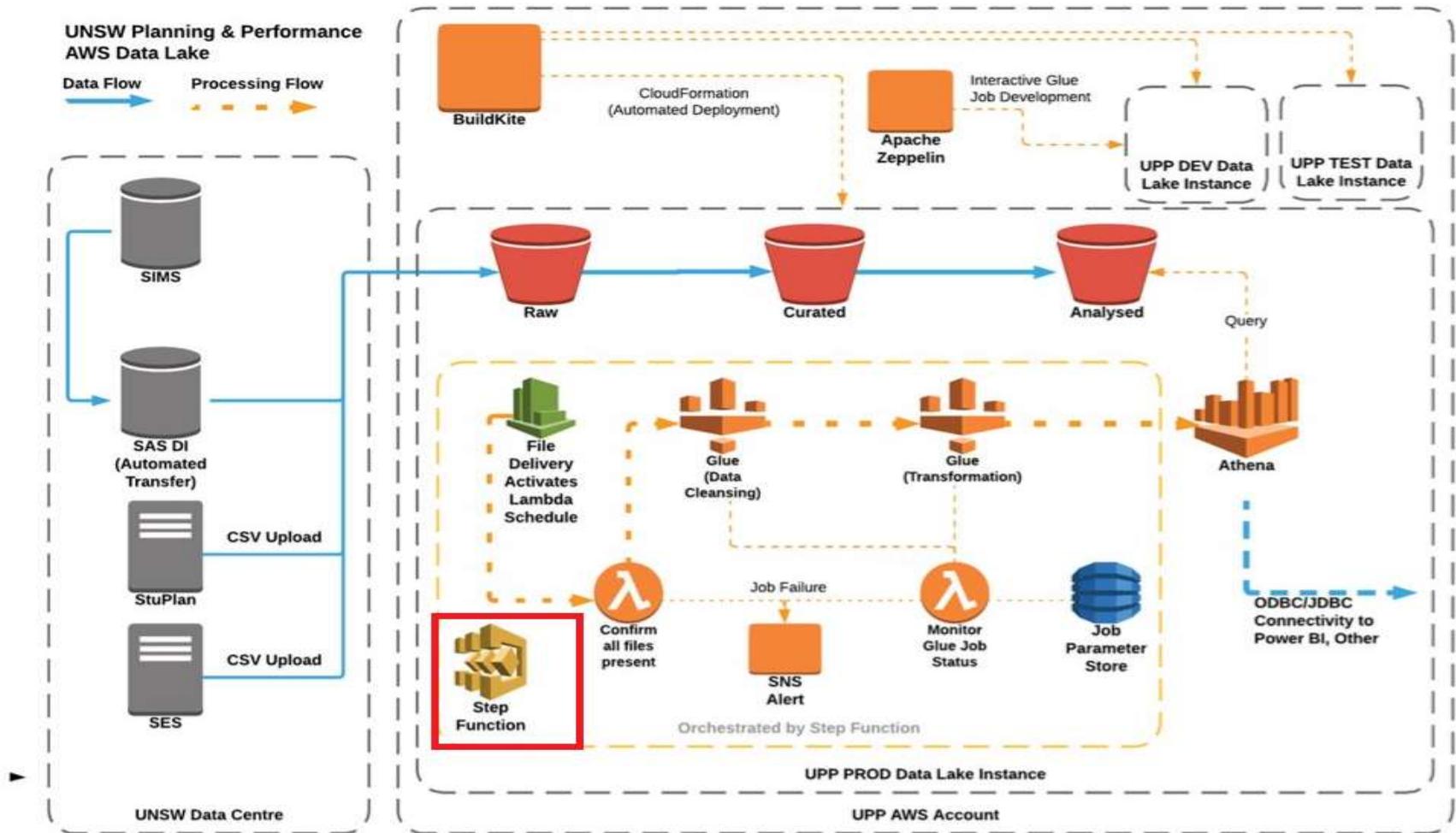
Serverless function, per second billing. Can integrate with different AWS service e.g. S3, Step Function, SNS



The screenshot shows the AWS Lambda console interface for a function named 'UNSW-datalake-stepfunction-SubmitCrawlerFunction-11Y...'. The function is configured with Python 2.7 as the runtime and 'startCrawler.lambda\_handler' as the handler. The code editor displays the following Python code:

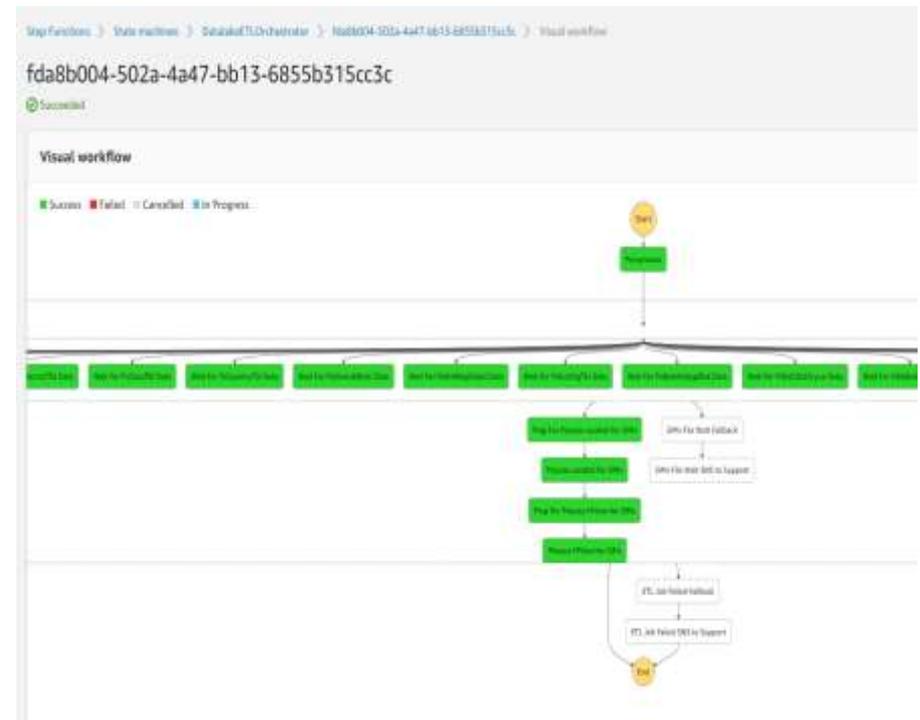
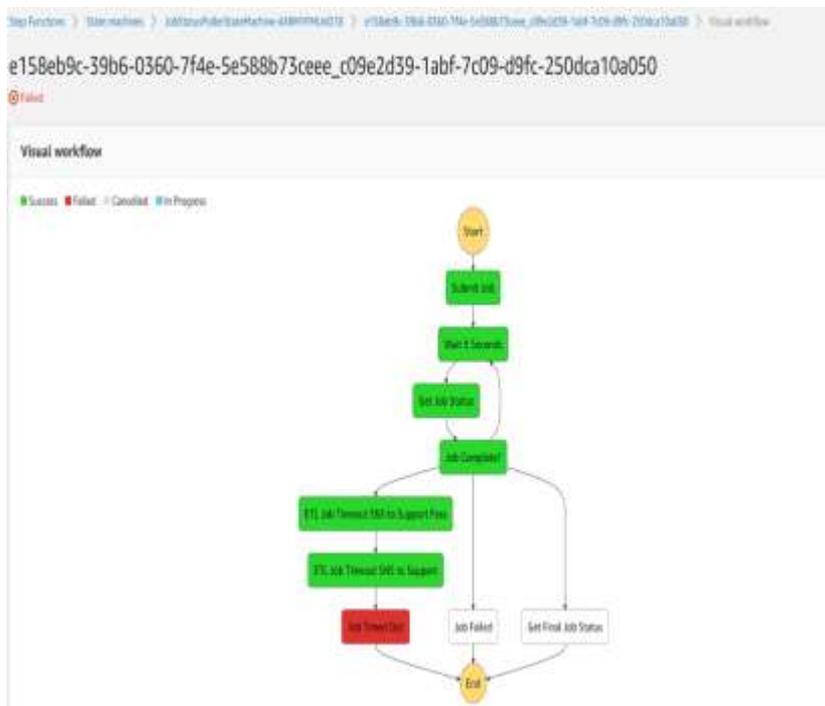
```
1: from __future__ import print_function
2:
3: import json
4: import boto3
5: import os
6:
7: print('Loading function')
8: glue = boto3.client('glue')
9:
10: def lambda_handler(event, context):
11:     # Log the received event
12:     print("Received event: " + json.dumps(event, indent=2))
13:     # Get crawlerName from the event
14:     crawlerName = event['crawlerName']
15:     try:
16:         # Start Crawler
17:         response = glue.start_crawler(Name=crawlerName)
18:         return {
19:             'crawlerName': crawlerName
20:         }
21:     except Exception as e:
22:         print(e)
23:         message = 'Error starting Crawler'
24:         print(message)
25:         raise Exception(message)
```

# STEP FUNCTION

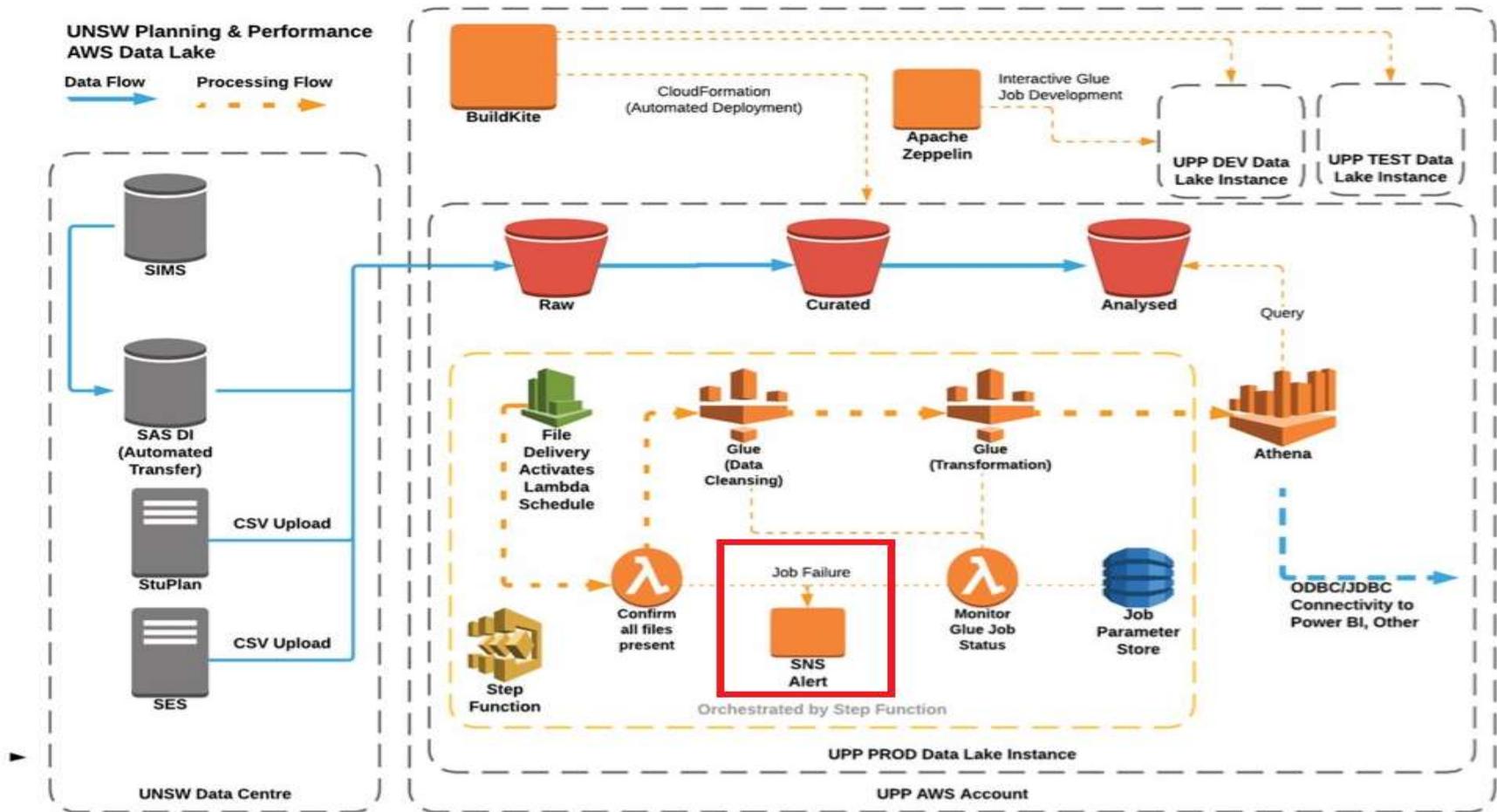


# STEP FUNCTION

AWS Step Functions lets you coordinate multiple AWS services into serverless workflows. We use step function in combination with S3 event to kick off and chain glue job, also with notification on time-out and failure.



# SIMPLE NOTIFICATION SERVICE (SNS)



# SIMPLE NOTIFICATION SERVICE (SNS)

## Failure Notifications

Amazon SNS > Topics > UNSW-datalake-stepfunction-UNSWDataLakeJobFailureSNSTopic-1X6WJPA50BRWT

### UNSW-datalake-stepfunction-UNSWDataLakeJobFailureSNSTopic-1X6WJPA50BRWT

Edit Delete Publish message

#### Details

Name	UNSW-datalake-stepfunction-UNSWDataLakeJobFailureSNSTopic-1X6WJPA50BRWT	Display name	UNSWDataLake-Job-Failure-Topic
ARN	arn:aws:sns:ap-southeast-2:386858679224:UNSW-datalake-stepfunction-UNSWDataLakeJobFailureSNSTopic-1X6WJPA50BRWT	Topic owner	386858679224

**Subscriptions** | Access policy | Delivery retry policy (HTTP/5) | Delivery status logging | Encryption | Tags

#### Subscriptions (1)

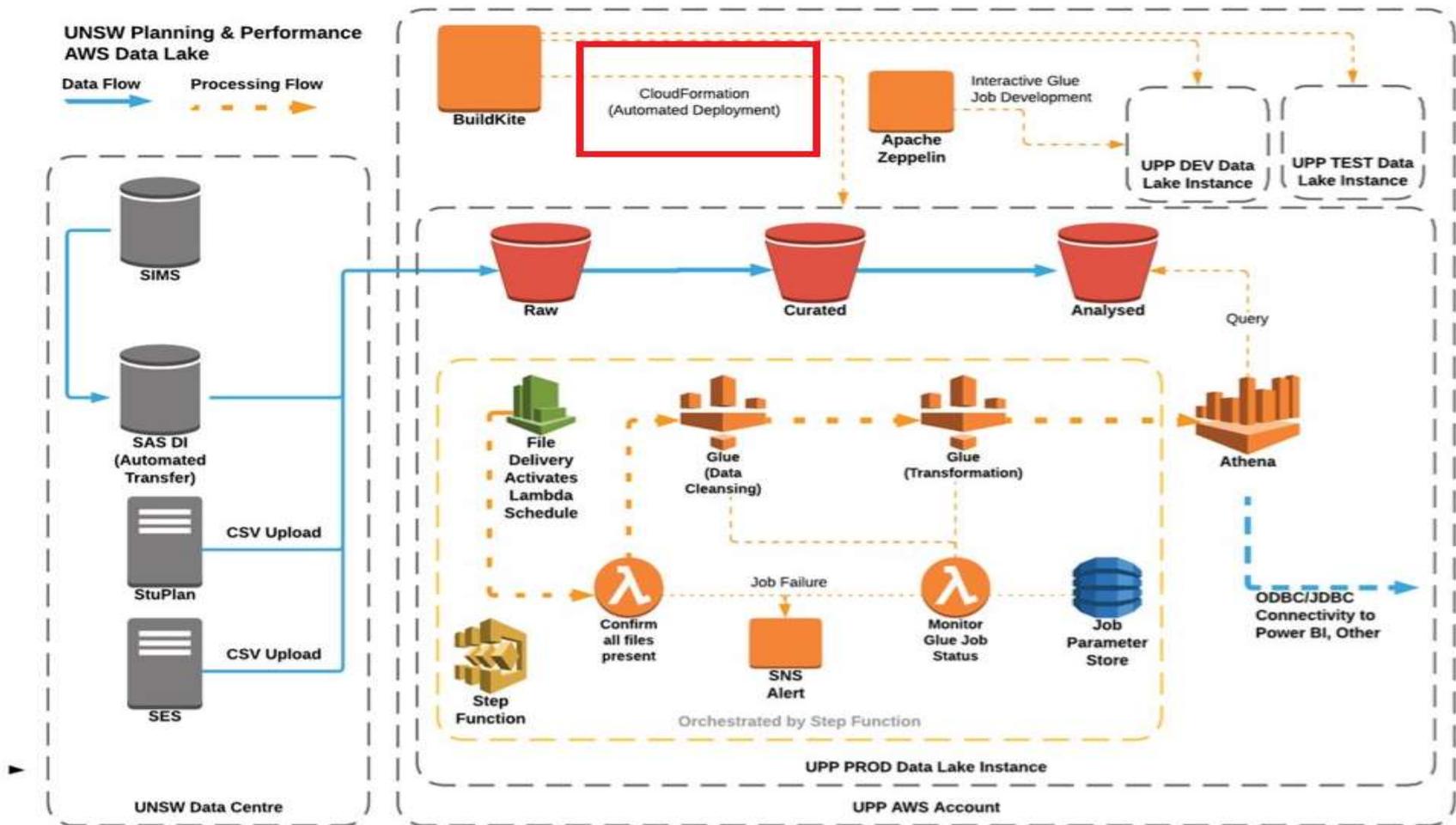
Edit Delete Request confirmation Confirm subscription **Create subscription**

< 1 > ⚙

ID	Endpoint	Status	Protocol
39875e76-9b8d-4139-bc02-1e0a76bfc06	████████@unsw.edu.au	Confirmed	EMAIL



# CLOUDFORMATION



# CLOUDFORMATION

## Code Templates for Deployment, Easily Control and Track Changes

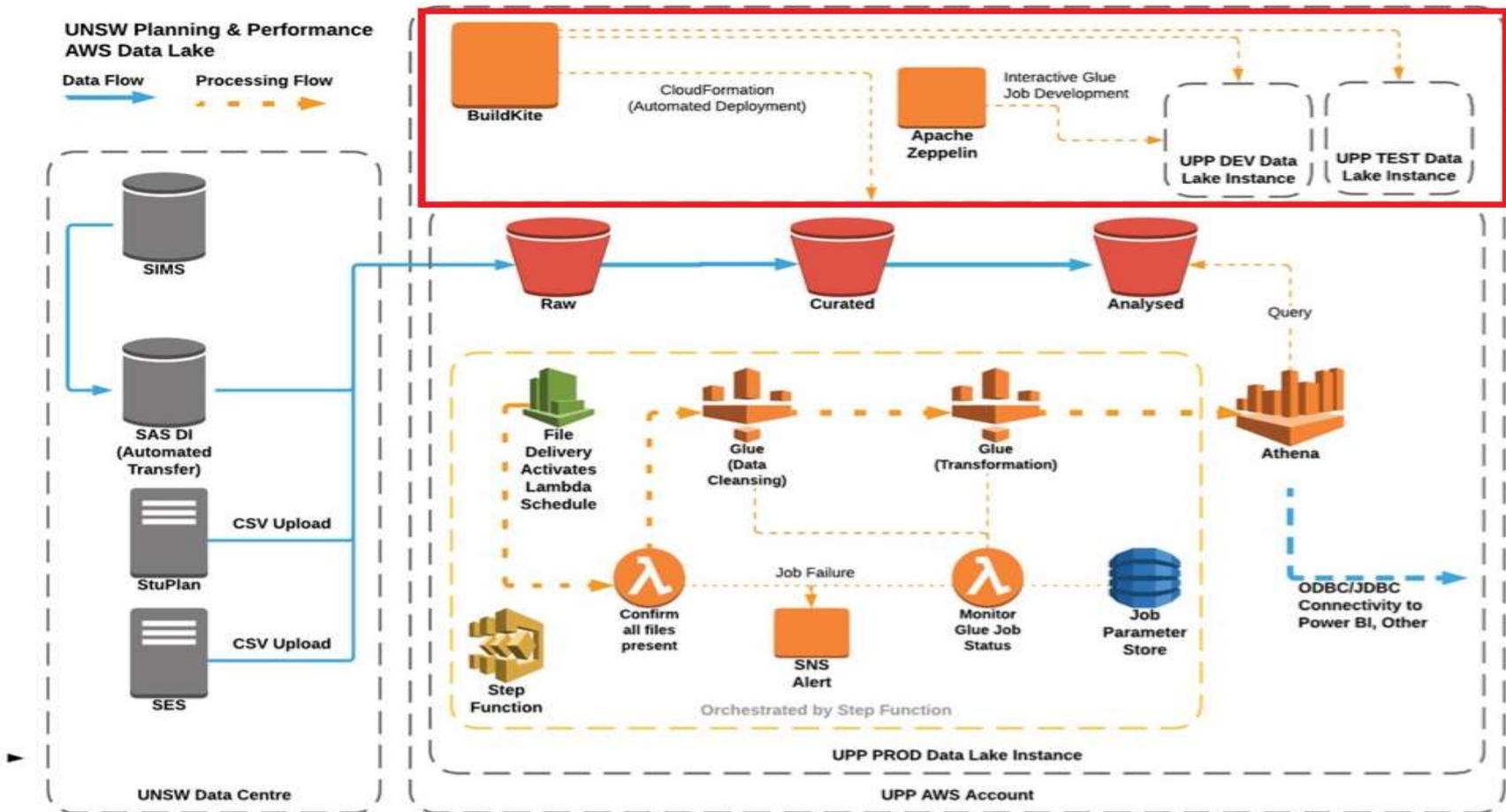
CloudFormation ▾ Stacks

Create Stack ▾ Actions ▾ Design template

Filter: Active ▾

Stack Name	Created Time	Status	Drift Status	Description
 UNSW-datalake-glue-SIMsRawDatabaseTableStack2-1EHQ10YLZ282T <span>NESTED</span>	2019-02-08 14:01:22 UTC+1100	UPDATE_COMPLETE	NOT_CHECKED	UNSW Datalake Glue
 UNSW-datalake-glue-CuratedDatabaseTableStack-1CM03388WZALT <span>NESTED</span>	2019-02-08 14:01:22 UTC+1100	UPDATE_COMPLETE	NOT_CHECKED	UNSW Datalake Glue
 UNSW-datalake-glue-SIMsRawDatabaseTableStack1-CAVFKQM5XRNL <span>NESTED</span>	2019-02-08 14:01:22 UTC+1100	UPDATE_COMPLETE	NOT_CHECKED	UNSW Datalake Glue
 UNSW-datalake-stepfunction	2019-02-01 13:42:46 UTC+1100	UPDATE_COMPLETE	NOT_CHECKED	UNSW Datalake Step Functions
 UNSW-datalake-azure-refresh	2019-01-15 09:43:13 UTC+1100	UPDATE_COMPLETE	NOT_CHECKED	Lambda function for Azure data extract from U
 UNSW-datalake-glue	2018-12-11 11:29:58 UTC+1100	UPDATE_COMPLETE	NOT_CHECKED	UNSW Datalake Glue
 gluerunner-lambda	2018-11-15 14:13:37 UTC+1100	UPDATE_COMPLETE	NOT_CHECKED	The CloudFormation template for AWS resour
 UNSW-buildkite-Agents	2018-09-04 10:32:08 UTC+1000	UPDATE_COMPLETE	NOT_CHECKED	Buildkite Elastic CI
 UNSW-buildkite-buckets	2018-09-04 10:21:09 UTC+1000	CREATE_COMPLETE	NOT_CHECKED	UNSW Datalake buildkite Buckets
 UNSW-datalake-s3	2018-09-03 10:22:02 UTC+1000	UPDATE_COMPLETE	NOT_CHECKED	UNSW Datalake Buckets

# OTHER SERVICES AND TOOLS



# OTHER SERVICES AND TOOLS

- 🔗 BitBucket
- 🔗 CI/CD
- 🔗 BuildKite
- 🔗 Zeppelin Notebook (Spark Script development and interact with Glue Dev Endpoint)



# ARCHITECTURAL CONSIDERATIONS

## ▪ **EMR vs Glue + Athena + Step Function + Lambda**

- ⌘ UNSW has minimal admin staff and does not want to maintain a hadoop cluster
- ⌘ Lower cost

## ▪ **Glue Job Batching**

- ⌘ Glue has minimum cost of 10 mins, some small job runs less than 1 min. Batched up to save cost.

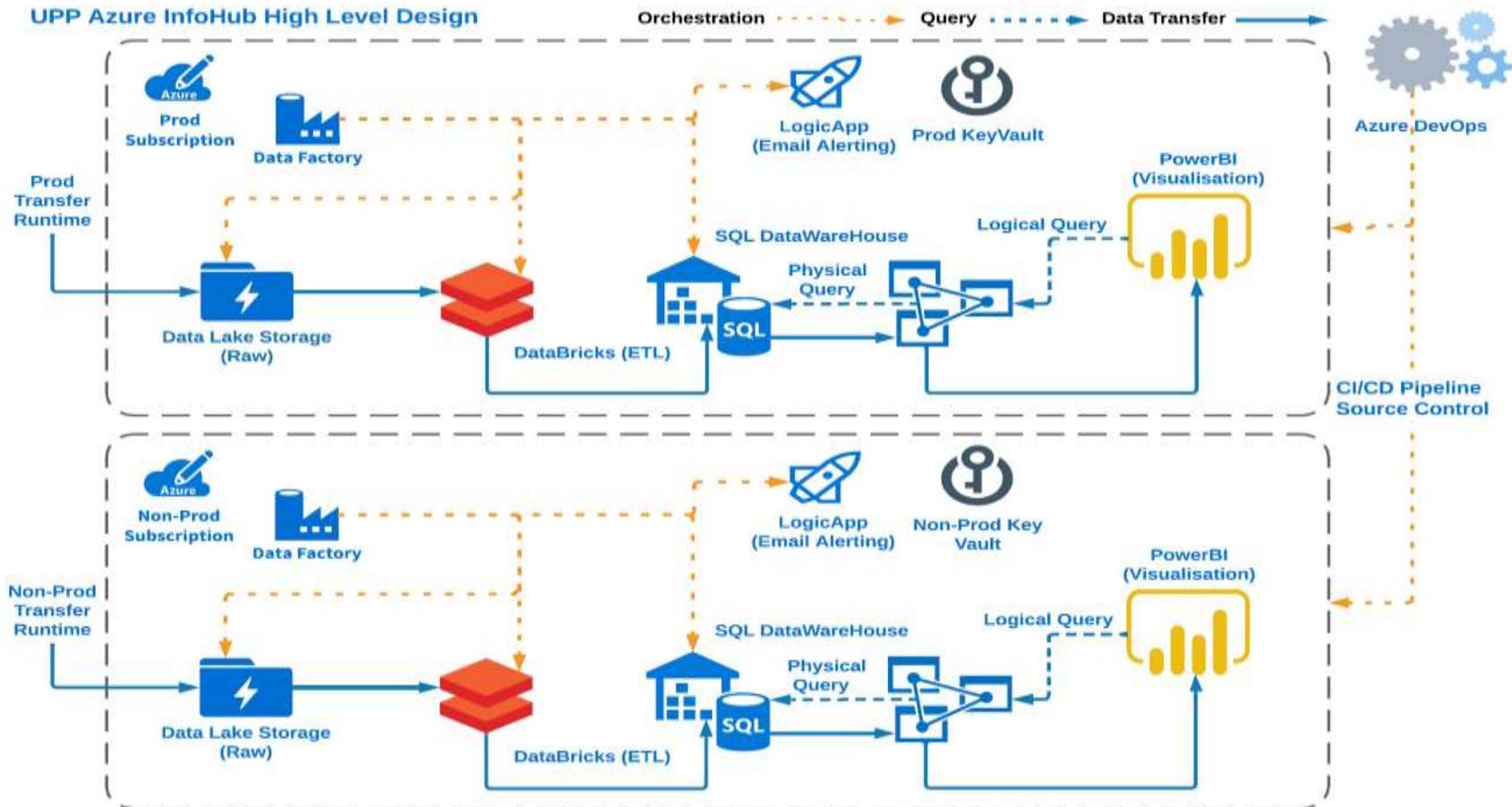
## ▪ **Visualization Tool**

- ⌘ Quicksight - lacks of Important features such as
  - Code promotion
  - Object ownership
  - Visualisation Features

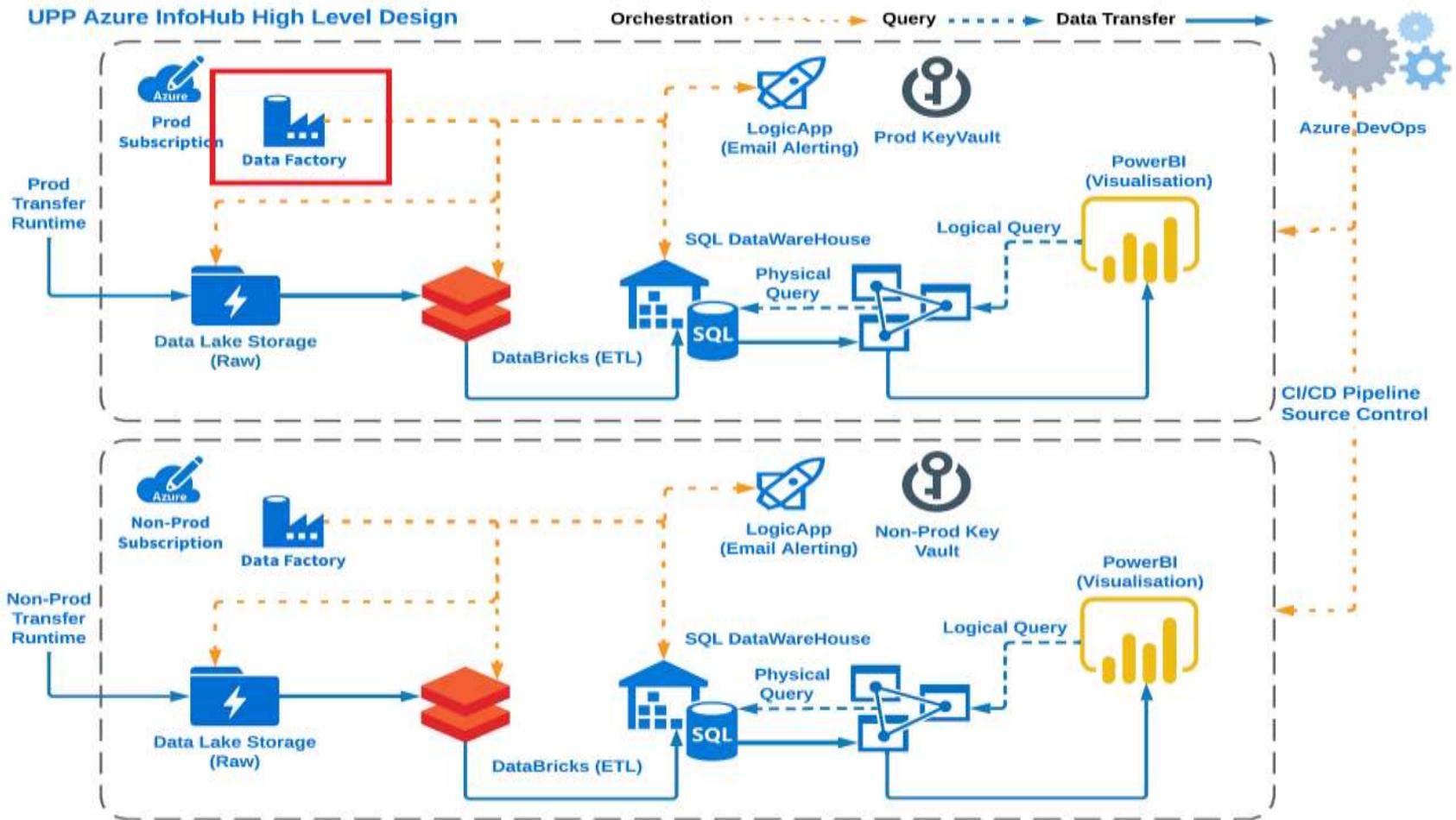
# KEY ACHIEVEMENTS

- ❧ Created the first **visual pipeline** of Students into the university, breaking down applicants by country of origin, age, degree and other data points through the Uni Application process as far back as the first student enrolled in the university in 1950 - a Mining Engineering Major.
- ❧ Established first true **data lake** for UPP storing raw, curated and analysed data from both source systems and legacy data marts in a central pipeline for easy consumption as a service.
- ❧ Purely **serverless** implementation of a data lake with all functions controlled through **event-driven orchestration**, ensuring lowest cost per workload.
- ❧ Implemented **CI/CD tooling** and **SCM** for the data lake with BuildKite, allowing one-click deployments
- ❧ Illustrated an example of **serverless cost-of-exit** - with a new DWH being built, data extracts could be provided in a matter of days via ODBC connection to Azure from Athena, facilitating transition.
- ❧ UPP team **uplifted skills** in Python and SQL etc.

# FUTURE ROADMAP - UPP AZURE INFOHUB

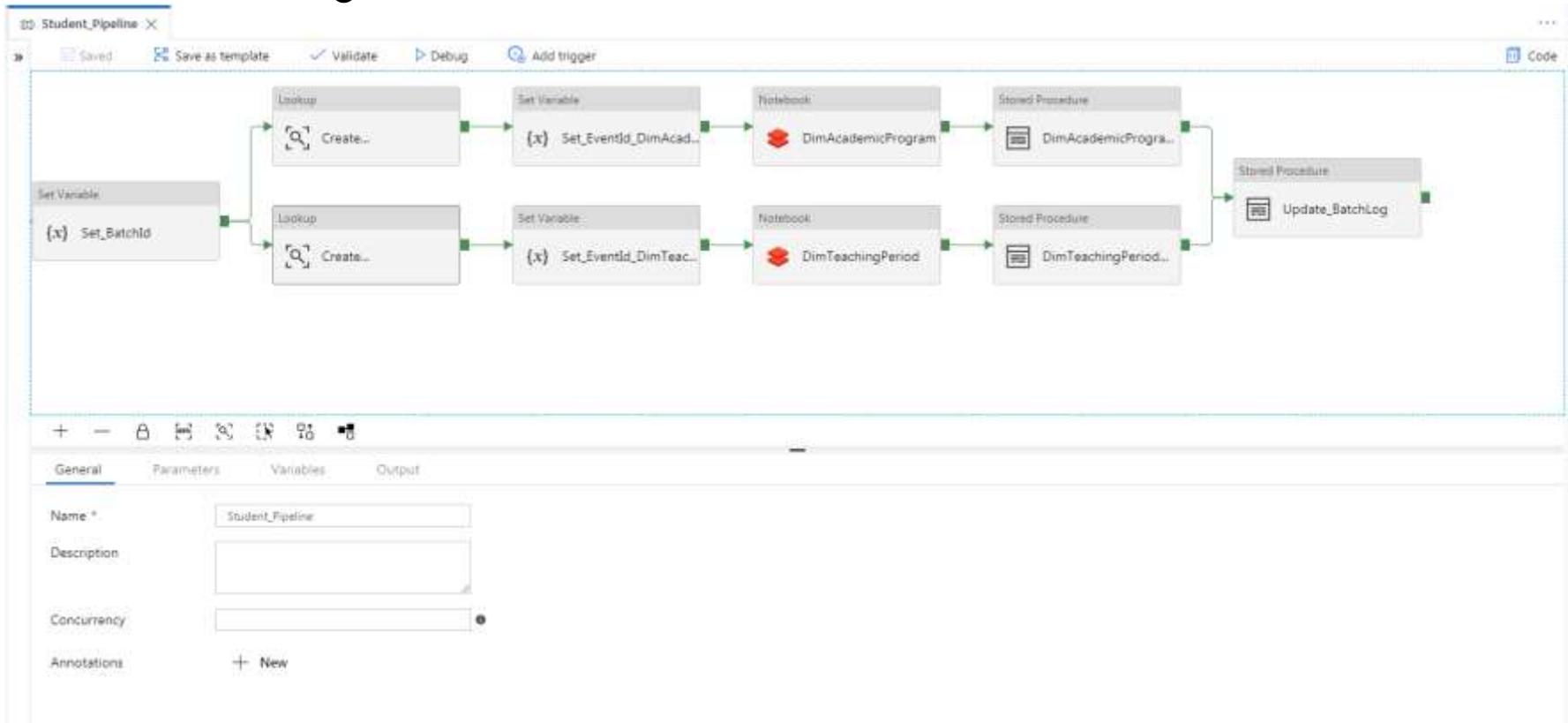


# AZURE DATA FACTORY

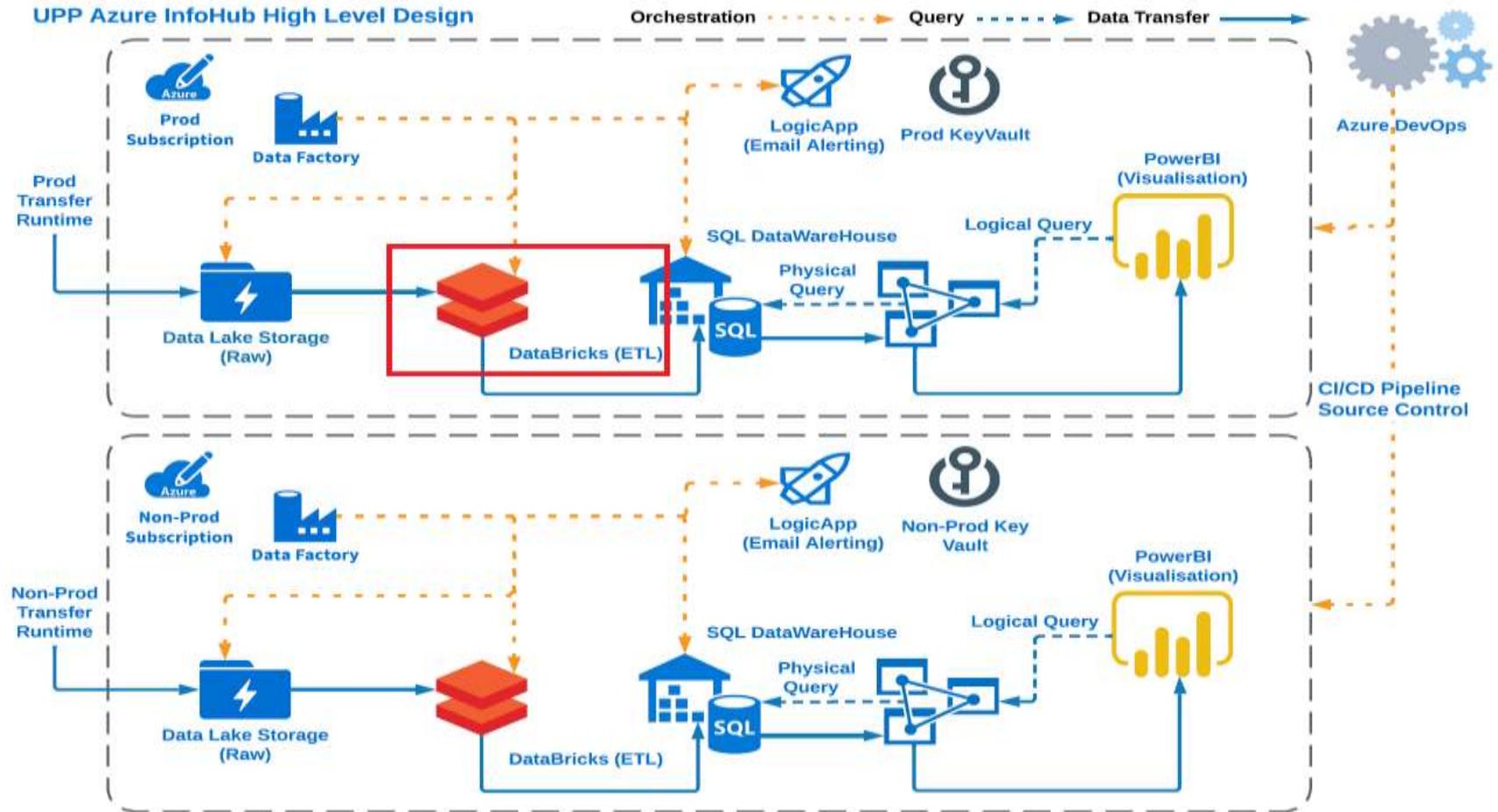


# AZURE DATA FACTORY

- ⑩ ETL and ELT processes code-free or write your own code
- ⑩ Natively built and maintenance-free connectors
- ⑩ serverless integration service

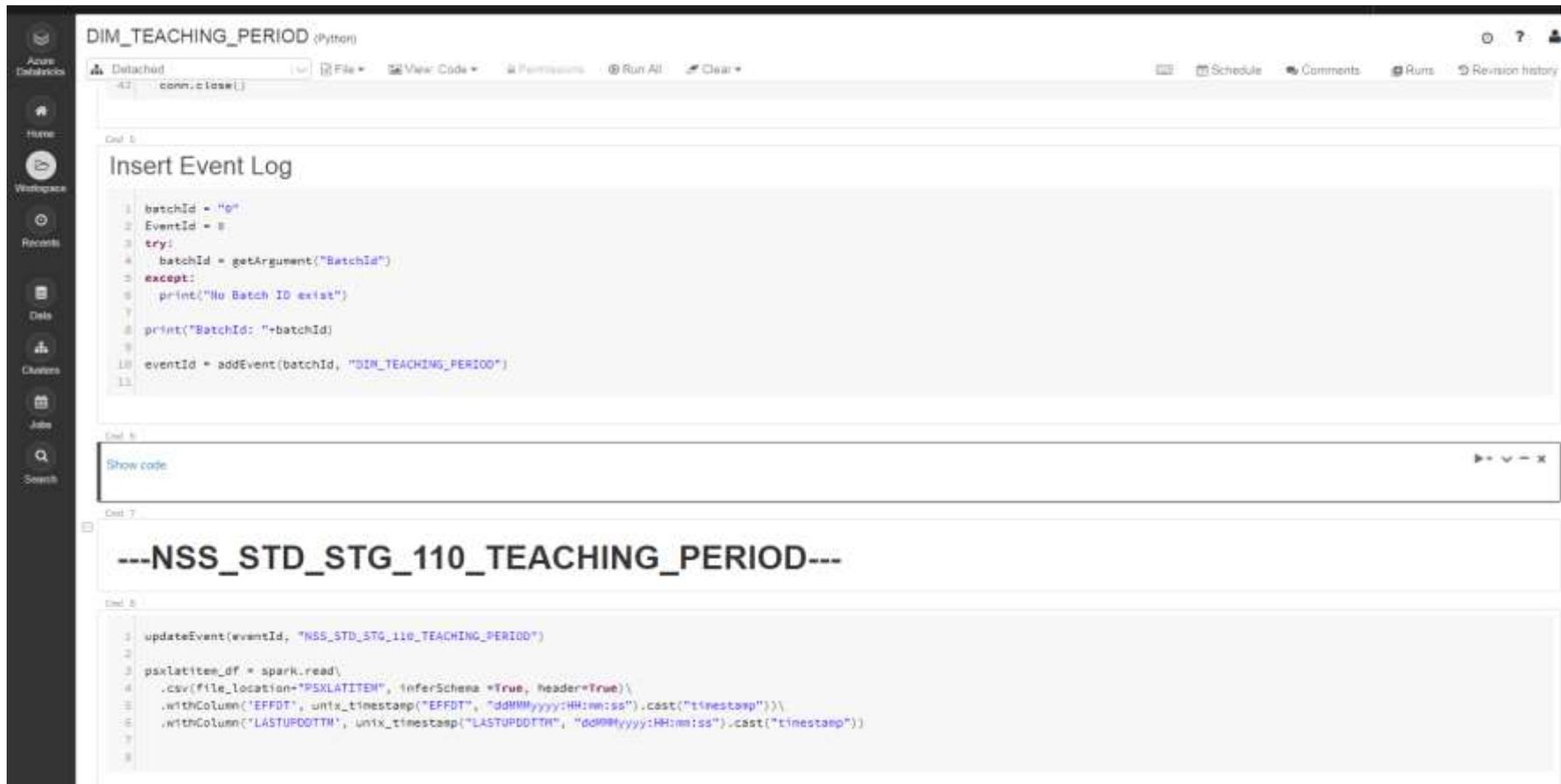


# AZURE DATABRICKS



# AZURE DATABRICKS

## Spark-based analytics service



The screenshot displays the Azure Databricks workspace interface. The top navigation bar includes options like 'Home', 'Workspace', 'Recent', 'Data', 'Clusters', 'Jobs', and 'Search'. The main area shows a notebook titled 'DIM\_TEACHING\_PERIOD' with a Python script. The script is divided into three sections: 'Insert Event Log', a 'Show code' section, and a section for updating an event.

```
conn.close()

Out 5:

Insert Event Log

1 batchId = "6"
2 eventId = 8
3 try:
4     batchId = getArguments("BatchId")
5 except:
6     print("No Batch ID exist")
7
8 print("BatchId: "+batchId)
9
10 eventId = addEvent(batchId, "DIM_TEACHING_PERIOD")
11

Out 6:

Show code

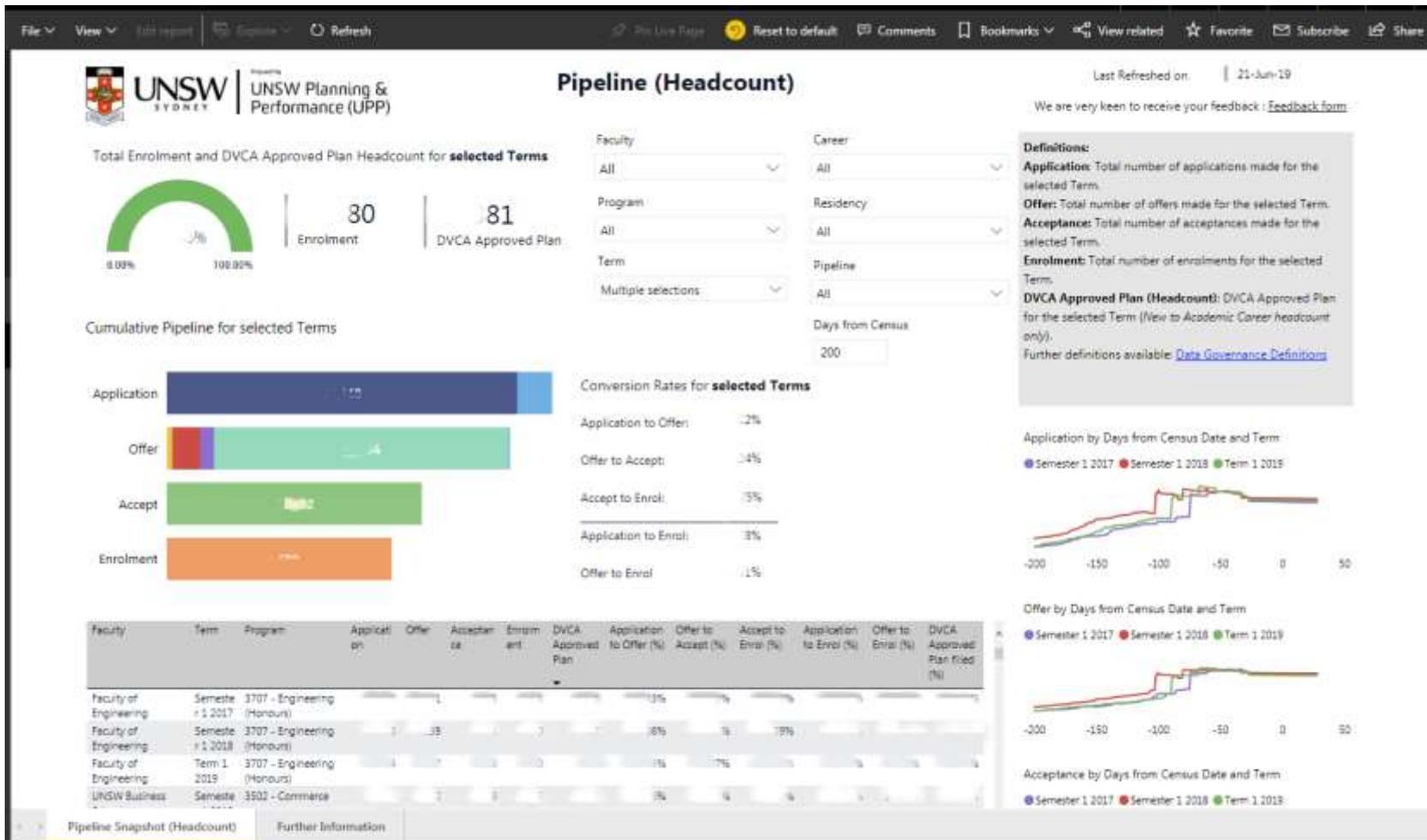
Out 7:

---NSS_STD_STG_110_TEACHING_PERIOD---
```

```
Incl 8:

1 updateEvent(eventId, "NSS_STD_STG_110_TEACHING_PERIOD")
2
3 psxlatites_df = spark.read
4     .csv(file_location+"PSXLATITEM", inferSchema =True, header=True))
5     .withColumn('EFFDT', unix_timestamp("EFFDT", "ddMM/yyyy:HH:mm:ss").cast("timestamp"))
6     .withColumn('LASTUPDDTTM', unix_timestamp("LASTUPDDTTM", "ddMM/yyyy:HH:mm:ss").cast("timestamp"))
7
8
```

# POWERBI



# AWS VS AZURE

AWS	Azure	Description
Step Function, AWS Glue	Data Factory 	Move data between storage. Schedule and orchestrate
Amazon Athena 	Azure data Lake Analytics	Serverless interactive query service
Glue	Azure Databricks 	Spark-based analytics
SNS	LogicApp	Sent out alerting email
CloudWatch	Azure Monitor	Track performance and create logs
Parameter Store	Key Vault	Store parameters
Simple Storage Service(S3)	Azure Blob storage	Object storage
QuickSight	PowerBI 	BI visualisation tools

