



Group 3 Final Project: Neptune Graph Database

SP 2022 DAV 6100: Information Architectures

By

Atreish Ramlakhan
Mahlet Melese
Shichao Zhou
Yuxiao Shen



Katz

Katz School
of Science and Health

Group Members



Atreish Ramlakhan

M.S. Artificial
Intelligence



Mahlet Melese

M.S. in Data Analytics
& Visualization



Shichao Zhou

M.S. in Data Analytics
& Visualization



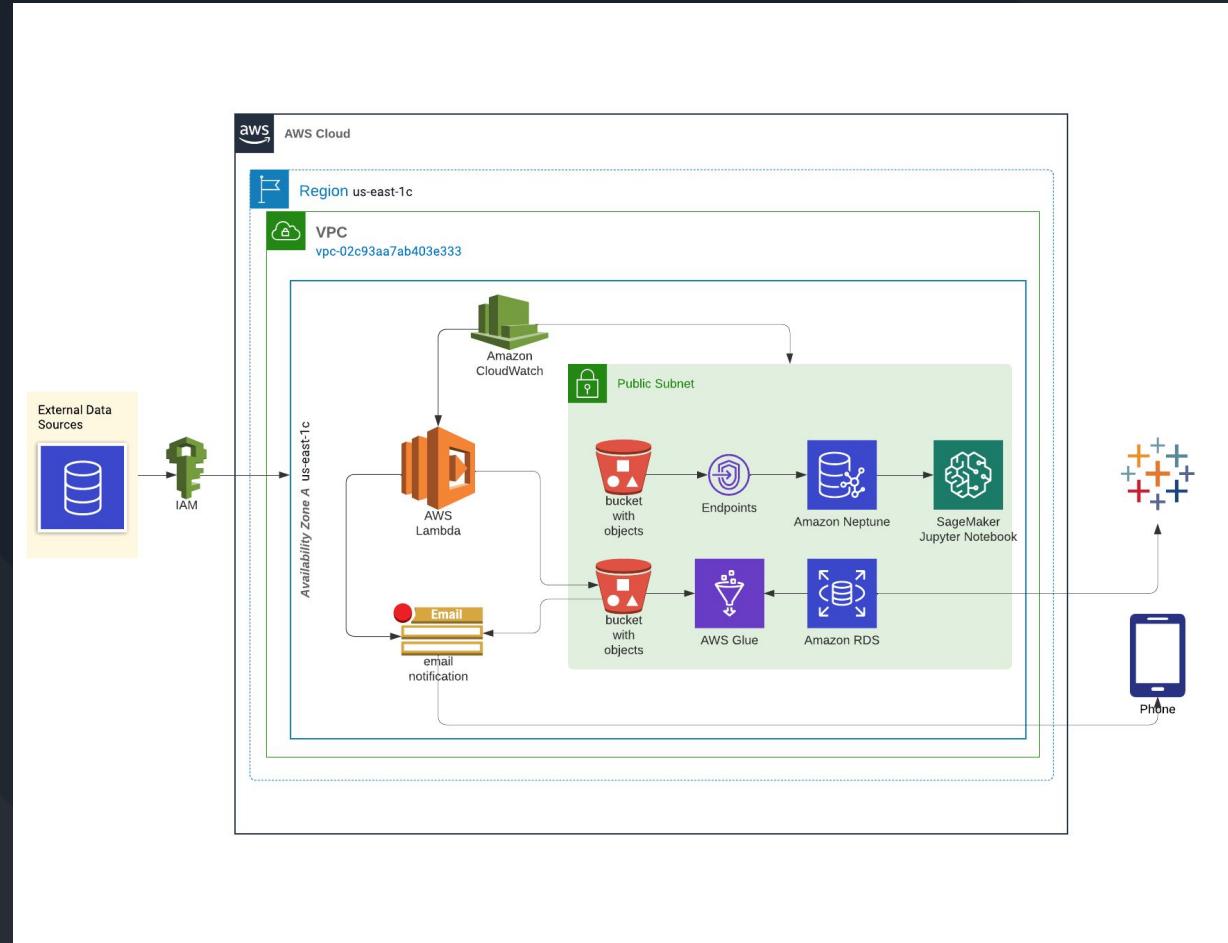
Yuxiao (Henry) Shen

M.S. in Data Analytics
& Visualization

Introduction



AWS Architecture



“Original” Data Structures

Student Focused Structure

Student ID	Course1	Course2	Course3	Course4	Course5	Course6	Course7	Course8	Course9	Course1
1168340921	Spring-2016: DAV 5200 Visual Design and Storytelling	Spring-2016: DAV 6100 Information Architectures	Spring-2016: DAV 5300 Computational Math and Statistics	Fall-2016: MAN 5580 Project Management	Fall-2016: DAV 6000 Talent Analytics	Fall-2016: Natural Language Processing	Spring-2017: DAV 5000 Business Modeling and Data Analysis	Spring-2017: DAV 6500 Capstone	Spring-2017: DAV 5100 Structured Data Management	Fall-2019: Spec Topic
1191512830	Fall-2018: Predictive Models	Fall-2018: Neural Networks and Deep Learning	Fall-2018: Machine Learning	Spring-2019: Data Acquisition and Management	Spring-2019: Natural Language Processing	Spring-2019: AI Capstone: R&D Experience	Fall-2019: Numerical Methods	Fall-2019: Artificial Intelligence	Fall-2019: Computational Statistics and Probability	Spring-2022: Bayesia Method
1232473375	Spring-2019: Data Visualization	Spring-2019: DAV 5400 Analytics Programming	Spring-2019: DAV 5200 Visual Design and Storytelling	Fall-2019: DAV 5000 Business Modeling and Data Analysis	Fall-2019: DAV 6000 Talent Analytics	Fall-2019: DAV 6200 Data Product Design	Spring-2020: DAV 5580 Capstone	Spring-2020: DAV 5100 Structured Data Management	Spring-2020: DAV 6500 Project Management	Fall-2021: Spec Topic
1313020897	Fall-2016: Advanced Data Engineering	Fall-2016: DAV 5100 Structured Data Management	Fall-2016: DAV 5000 Business Modeling and Data Analysis	Spring-2017: DAV 5200 Analytics Programming	Spring-2017: DAV 6400 Internship	Fall-2017: DAV 6450 Independent Study	Fall-2017: Independent Study	Fall-2017: DAV 5300 Computational Math and Statistics	Spring-2017: DAV 6500 Capstone	Fall-2021: DAV 6000 Project
1364243602	Spring-2020: MAN 5590 Project	Spring-2020: DAV 6000 Talent	Spring-2020: DAV 5200 Visual Design and	Fall-2020: DAV 6500 Capstone	Fall-2020: DAV 5100 Internship	Fall-2020: Structured Data	Spring-2021: DAV 5400 Analytics	Spring-2021: Independent Study	Spring-2021: Natural Language	Spring-2021: DAV 5000 Business Modelin

Course Focused Structure

Course Name	Student ID	First Name	Last Name	Katz School Major	GPA	Graduation Semester	Country of Origin	Languages	Undergraduate Major	Age at Graduation	Years of Experience	Location	Dept
Fall-2015: AI Capstone: R&D Experience	4784813239	Calahan	Martin	Artificial Intelligence	2.95	Spring-2017	USA	English	Computer Science	27	5.5	New York, NY	Engineering
	4800891893	Johnson	Angela	Artificial Intelligence	3.51	Spring-2017	China	Mandarin	Mathematics	24	2.5	New York, NY	Engineering
Fall-2015: AI Product Studio	8479614053	Pockrus	Clifford	Data Analytics and Visualization	4.00	Spring-2017	China	Mandarin	Other	29	8.0	New York, NY	Business
	9290159667	Potter	Eugene	Artificial Intelligence	3.63	Spring-2017	Pakistan	Sindhi	Mathematics	27	6.0	New York, NY	Arts
Fall-2015: Advanced Data Engineering	4991566949	Berry	Susan	Data Analytics and Visualization	3.55	Spring-2017	China	Mandarin	Computer Science	25	3.0	New York, NY	Business

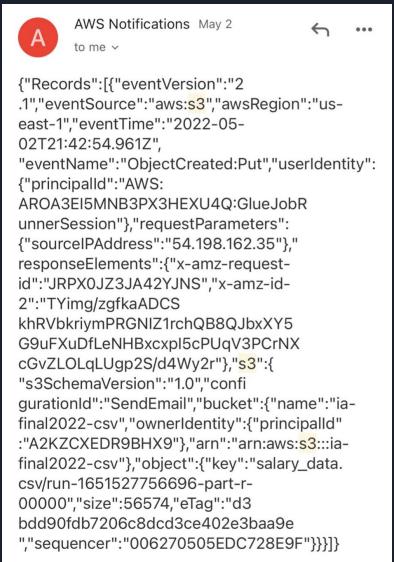
Spring-2022: MAN 5580				Data Analytics		Spring-			Computer			Los Angeles	Design

Data Profile: Student Dummy Dataset

Dataset Summary	
Source of Information	Randomly generated dataset
Number of Records	1319
Frequency of updates	Yearly
Data type and structure	
Number of columns	14
Granularity	Age

Lambda Functions and Layers

Email notification



The screenshot shows the AWS Lambda Function Overview page for 'ia-final2022'. The function was last modified 'last month' and has a Function ARN of 'arn:aws:lambda:us-east-1:765102680182:function:ia-final2022'. The 'Code source' tab is selected, showing the Python code for generating student data. The 'Layers' tab shows two layers: 'ia-final2022-layer1' and 'ia-final2022-layer2'. The 'Runtime settings' tab shows the runtime as Python 3.8, handler as 'StudentDataGenerator.main', and architecture as x86_64. The 'Layers' table lists the ARNs for each layer.

Merge order	Name	Layer version	Compatible runtimes	Compatible architectures	Version ARN
1	ia-final2022-layer1	1	-	-	arn:aws:lambda:us-east-1:765102680182:layer:ia-final2022-layer1:1
2	ia-final2022-layer2	1	-	-	arn:aws:lambda:us-east-1:765102680182:layer:ia-final2022-layer2:1

S3 Bucket

Python packages & 'names' generator package

The screenshot shows the AWS S3 console interface. The left sidebar includes links for Services, S3, Lambda, CloudWatch, and Neptune. The main content area displays the 'ia-final2022' bucket. A blue banner at the top says, "We're continuing to improve the S3 console to make it faster and easier to use. If you have feedback on the updated experience, choose [Provide feedback](#)". The 'Objects' tab is selected, showing two items:

Name	Type	Last modified	Size	Storage class
python.zip	zip	April 4, 2022, 20:54:50 (UTC-04:00)	47.0 MB	Standard
studentGenerator_layer.zip	zip	April 4, 2022, 20:54:50 (UTC-04:00)	782.9 KB	Standard

S3 Bucket

Dummy, LinkedIn, Salary.com

AWS Services Search for services, features, blogs, docs, and more [Option+S] Global final-project @ 7651-0268-0182 ▾

53 Lambda CloudWatch Neptune

Amazon S3

We're continuing to improve the S3 console to make it faster and easier to use. If you have feedback on the updated experience, choose Provide feedback. Provide feedback X ⓘ

Amazon S3 > Buckets > ia-final2022-csv

ia-final2022-csv [Info](#)

Objects Properties Permissions Metrics Management Access Points

Objects (4)

Objects are the fundamental entities stored in Amazon S3. You can use [Amazon S3 inventory](#) to get a list of all objects in your bucket. For others to access your objects, you'll need to explicitly grant them permissions. [Learn more](#)

[Copy S3 URI](#) [Copy URL](#) [Download](#) [Open](#) [Delete](#) Actions Create folder Upload

Find objects by prefix < 1 > ⚙

	Name	Type	Last modified	Size	Storage class
linkedin_info.csv	csv	April 11, 2022, 23:22:45 (UTC-04:00)	672.0 B	Standard	
salary_data.csv	csv	April 4, 2022, 21:18:23 (UTC-04:00)	55.1 KB	Standard	
salary_data.csv/	Folder	-	-	-	
student_data.csv	csv	April 12, 2022, 22:42:24 (UTC-04:00)	1.3 MB	Standard	

Block Public Access settings for this account

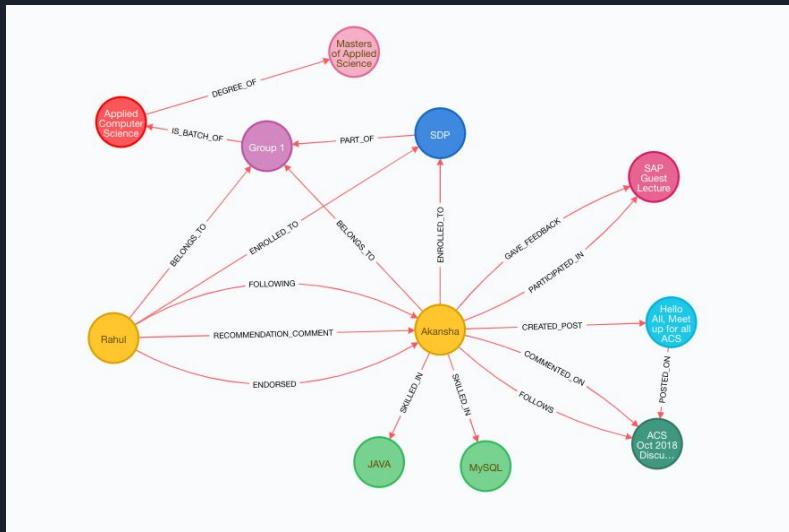
Storage Lens Dashboards AWS Organizations settings

Feature spotlight 3

AWS Marketplace for S3

Graph Databases(Neptune)

- Used to create relationship between data.
- Store data as Vertices of graph and the relationships as edge.



Graph Database

Several nodes exist in the graph database which allow for edges or relationships across students. For example, Rahul is a student who is following Akansha. Both Rahul and Akansha belong to Group 1 as part of the Applied Computer Science Class, which is part of the Master of Applied Science Program.

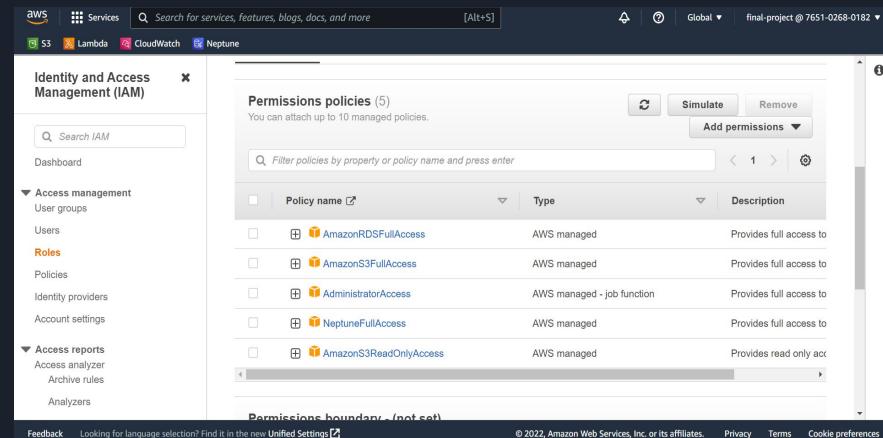
Neptune

A screenshot of Microsoft Excel showing a table of data. The table has columns labeled "id", "from", "to", and "label". Rows 9 through 18 are selected, indicated by a light blue background. The status bar at the bottom shows "edge" and "Sheet1".

Create Neptune cluster

Creating new database cluster

- Setup VPC where the database cluster is located.
- Setup the security group
- Add IAM role



Neptune Database

The screenshot shows the AWS Neptune console interface. On the left, there's a sidebar with links for Neptune, Databases, Snapshots, Notebooks, Subnet groups, Parameter groups, Events, and Event subscriptions. The main area is titled "Neptune > Databases". It displays a table of databases with columns for DB identifier, Role, Engine, Region & AZ, and Size. Two entries are visible: "database-2" (Cluster, Neptune, us-east-1) and "database-2-instance-1" (Writer, Neptune, us-east-1a, db.t3.medium). Above the table are buttons for "Group resources", "Modify", "Actions", and "Create database". A search bar at the top of the table allows filtering by database identifier.

Databases

Group resources Modify Actions ▾ Create database

Filter databases

DB identifier	Role	Engine	Region & AZ	Size
database-2	Cluster	Neptune	us-east-1	-
database-2-instance-1	Writer	Neptune	us-east-1a	db.t3.medium

Connectivity & security

Endpoint & port	Networking	Security
Endpoint database-4-1.cmln4r0t1v.us-east-1.neptune.amazonaws.com	Availability zone us-east-1a	VPC security groups neptune-sg-05770898a0327562 (active)
Port 8112	VPC vpc-0f1d7fa8a2e6	Public accessibility No
Subnet group default-vpc-0f407a2a5a9e	Subnet subnet-04691006a10a645 subnet-07156217a202536 subnet-0355294a045421 subnet-0963520730a16 subnet-05d62ca059a43a2 subnet-06c703a1757a880	

Load data from S3 to Neptune

The screenshot shows the AWS Management Console with three main windows:

- S3 Service Dashboard:** Shows two objects: `edge.csv` and `nodemahi.csv`. `edge.csv` was modified on May 11, 2022, at 23:34:47 (UTC-04:00) and is 2.4 KB in size. `nodemahi.csv` was modified on May 12, 2022, at 01:08:30 (UTC-04:00) and is 23.9 KB in size.
- VPC Endpoints:** Displays a list of endpoints, including one for the Neptune service.
- Neptune Service Dashboard:** Shows configuration options for loading data from S3. The "Source" is set to `s3:neptuneDb`, "Format" to "csv", "AWS Region" to "us-west-2", and "Load ARN" to `arn:aws:lambda:216153009192:ne`. The "Mode" is set to "AUTO".

Jupyter Notebook (Left): A Python 3 notebook cell containing the following code to configure the Neptune client:

```
openCypher: {'version': 'Neptune-9.0.20190305-1.0'},  
labMode: {'objectIndexing': 'disabled',  
          'queryExecution': 'enabled'},  
features: {'resultCache': {'status': 'disabled'},  
           'IAMAuthentication': 'disabled',  
           'Streams': 'disabled',  
           'Auditing': 'disabled'},  
settings: {'clusterQueryTimeoutInMs': '120000'}
```

Jupyter Notebook (Right): A Python 3 notebook cell showing the results of a Gremlin query:

```
In [1]: %load  
Total execution time: 6 seconds  
Done.  
  
In [2]: %%gremlin  
g.V().limit(20).valueMap()  
Console  
Show 25 entries  
Search: _____  
  
1  ('name': ['Fall2015: AI Capstone: R&D Experience', 'Calahan Martin', 'marko'], 'program': ['Artificial Intel  
2  ('name': ['Fall2015: AI Capstone: R&D Experience', 'Calahan Martin', 'lop'], 'program': ['Artificial Intel  
3  ('undergrad': ['Computer Science'], 'yearsofexperience': [6], 'graduationsemester': ['Spring-2017']), 'name': ['C  
4  ('name': ['Fall2015: AI Capstone: R&D Experience'], 'program': ['Artificial Intelligence'])  
5  ('undergrad': ['Mathematics'], 'yearsofexperience': [3], 'graduationsemester': ['Spring-2017']), 'name': ['Johnso  
6  ('name': ['Fall2015: AI Product Studio'], 'program': ['Artificial Intelligence'])  
7  ('undergrad': ['Other'], 'yearsofexperience': [8], 'graduationsemester': ['Spring-2017']), 'name': ['Pockrus Clif  
8  ('name': ['Fall2015: AI Product Studio'], 'program': ['Data Analytics and Visualization'])  
9  ('undergrad': ['Mathematics'], 'yearsofexperience': [6], 'graduationsemester': ['Spring-2017']), 'name': ['Potter
```

Install the Gremlin Console and Connect to Neptune

- ❑ Create EC2 instance with key pair
- ❑ Install Java 8 on EC2 instance
- ❑ wget`https://archive.apache.org/dist/tinkerpop/3.5.2/apache-tinkerpop-gremlin-console-3.5.2-bin.zip`

```
[ec2-user@ip-172-30-0-62: ~]$ login -f ec2-user
[ec2-user@ip-172-30-0-62: ~]$ Authenticating with public key "imported-openssh-key"
Last login: Thu May 12 18:12:00 2022 from 69.143.190.85
[ec2-user@ip-172-30-0-62: ~]$ Amazon Linux 2 AMI
[ec2-user@ip-172-30-0-62: ~]$ https://aws.amazon.com/amazon-linux-2/
[ec2-user@ip-172-30-0-62: ~]$ wget https://www.amazontrust.com/repository/SFSRoot
[ec2-user@ip-172-30-0-62: ~]$ curl -O https://www.amazontrust.com/repository/SFSRootCA02.cer
Resolving www.amazontrust.com (www.amazontrust.com)... 52.85.130.97, 52.85.130.1
[ec2-user@ip-172-30-0-62: ~]$ curl -O https://www.amazontrust.com/repository/SFSRootCA02.cer
Connecting to www.amazontrust.com (www.amazontrust.com)[52.85.130.97]:443...
connected.
HTTP request sent, awaiting response... 200 OK
Length: 1011 [application/xkie-certs]
Saving to: 'SFSRootCA02.cer' [100%] 1,011 --.-K/s in 0s
2022-05-12 18:15:18 (137 MB/s) - 'SFSRootCA02.cer.lz' saved [101/1011]

[ec2-user@ip-172-30-0-62: ~]$ mkdir /tmp/certs/
[ec2-user@ip-172-30-0-62: ~]$ mv SFSRootCA02.cer /tmp/certs/
[ec2-user@ip-172-30-0-62: ~]$ cp /tmp/certs/SFSRootCA02.cer /etc/pki/tls/certs/cacerts
cp: cannot stat '/tmp/certs/SFSRootCA02.cer': No such file or directory
[ec2-user@ip-172-30-0-62: ~]$ java -version
openjdk version "1.8.0_312-8u312-b07-1~rhel8.0.312~"
OpenJDK Runtime Environment (build 1.8.0_312-b07)
OpenJDK 64-Bit Server VM (build 25.312-b07, mixed mode)
[ec2-user@ip-172-30-0-62: ~]$ which java
/usr/bin/java
[ec2-user@ip-172-30-0-62: ~]$ file $(which java)
/usr/bin/java: ELF 64-bit LSB shared object, x86_64, version 1.0.0, dynamically linked, Build ID: 1.8.0_312_b07-1_amzn2.0.2_x86_64/jre/bin/java
[ec2-user@ip-172-30-0-62: ~]$ rm /etc/alternatives/java
[ec2-user@ip-172-30-0-62: ~]$ ln -s $(which java) /etc/alternatives/java
[ec2-user@ip-172-30-0-62: ~]$ rm /etc/alternatives/java
[ec2-user@ip-172-30-0-62: ~]$ cp /usr/lib/jvm/java-1.8.0-openjdk-1.8.0.312.b07-1_amzn2.0.2_x86_64/jre/bin/java*
[ec2-user@ip-172-30-0-62: ~]$ cp /usr/lib/jvm/java-1.8.0-openjdk-1.8.0.312.b07-1_amzn2.0.2_x86_64/jre*
cp: missing destination file operand after '/usr/lib/jvm/java-1.8.0-openjdk-1.8.0.312.b07-1_amzn2.0.2_x86_64/jre'
Try 'cp --help' for more information.
```

```
\.,.,/
(o o)
-----o00o-(3)-o00o-----
plugin activated: tinkerpop.server
plugin activated: tinkerpop.utilities
plugin activated: tinkerpop.tinkergraph
gremlin>
```

Create Vertex and Edge for Gremlin

```
File Edit Format View Help

v1 = g.addV('person')\
    .property(id, 1)\
    .property('name', 'marko')\
    .property('age', 29)\
    .next()

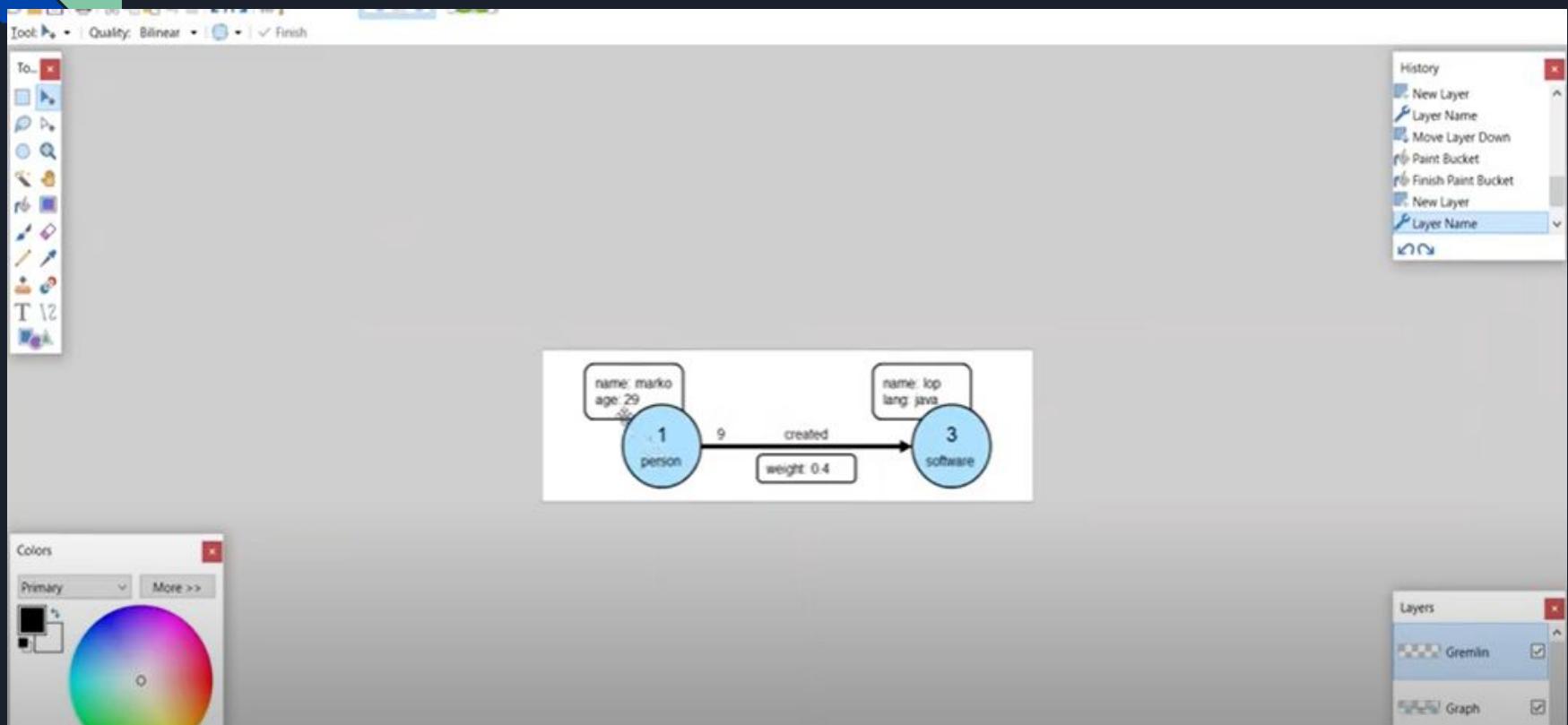
// Create the "software" vertex and assign it to "v2"

v2 = g.addV('software')\
    .property(id, 3)\
    .property('name', 'lop')\
    .property('lang', 'java')\
    .next()

// Create the "created" edge from "v1" to "v2"

g.AddE('created')\
    .from(v1)\
    .to(v2) \
    .property(id, 9)\
```

Con't...



Glue

The screenshot shows the AWS Glue Data Catalog interface. On the left, there's a sidebar with navigation links: Data catalog, Databases, Tables, Connections, Crawlers, Classifiers, Schema registries, Schemas, Settings, ETL, Security, Tutorials, and Security configurations. The 'Classifiers' link is highlighted. The main area shows a breadcrumb path: AWS Glue > Databases > db-final-project. Below the path are two buttons: 'Edit database' and 'Delete database'. To the right, detailed information about the database is displayed:

Name	db-final-project
Description	
Location	db-final-project.cmlm4mbt7riv.us-east-1.rds.amazonaws.com

The image contains three side-by-side screenshots of the AWS Glue Classifier configuration interface.

Screenshot 1 (Left): Edit classifier - classifier-csv

This screenshot shows the 'Edit classifier' dialog for 'classifier-csv'. It includes fields for 'Classifier name' (classifier-csv), 'Classifier type' (CSV), 'Column delimiter' (Comma), 'Quote symbol' (Double-quote), 'Column headings' (Detect headings), and 'Processing options' (Allow files with single column, Trim whitespace before identifying column values). A note at the top states: "If you change a classifier definition, any data that was previously crawled using the classifier is not re-crawled. New data is classified with the updated classifier which might result in an updated schema. Learn more."

Screenshot 2 (Middle): Crawlers > crawler-student-data

This screenshot shows the 'Edit' view for the crawler 'crawler-student-data'. It displays the following configuration:

Name	crawler-student-data
Description	Create a single schema for each S3 path
Security configuration	false
Tags	-
State	Ready
Schedule	Last updated: Tue May 03 20:44:17 GMT-400 2022 Date created: Tue May 03 20:44:17 GMT-400 2022
Database	db-final-project
Table level	
Service role	iam_glue

Screenshot 3 (Right): Crawlers > crawler-salary-data

This screenshot shows the 'Edit' view for the crawler 'crawler-salary-data'. It displays the following configuration:

Name	crawler-salary-data
Description	Create a single schema for each S3 path
Security configuration	false
Tags	-
State	Ready
Schedule	Last updated: Tue May 03 20:36:56 GMT-400 2022 Date created: Tue May 03 20:28:40 GMT-400 2022
Database	db-final-project
Table level	
Service role	iam_glue

Bottom Navigation Bar

The bottom navigation bar includes links for AWS Glue Studio, Jobs (New), Jobs (legacy), ML Transforms, Blueprints, Workflows, Triggers, Dev endpoints, Notebooks, Interactive Session - New, Configuration options, Schema updates in the data store, Object deletion in the data store, and Update the table definition in the data catalog.

Glue

AWS Glue Tables > student_data.csv

Last updated 3 May 2022 08:45 PM Table Version (Current version)

Data catalog

Databases
Tables
Connections
Crawlers
Classifiers
Schema registries
Schemas
Settings

ETL

AWS Glue Studio (2)
Jobs (2) New
Jobs (legacy)
ML Transformations
Blueprints
Workflows
Triggers
Dev endpoints
Notebooks
Interactive Session - New

Schemas

student_data.csv

Name: student_data.csv
Description: db-final-project
Database: db-final-project
Classification: csv
Location: s3://fe-fns2022-csv/student_data.csv
Connection: Connection
Dependencies: None
Last updated: Tue May 03 20:45:22 GMT+00:00 2022
Input format: org.apache.hadoop.mapred.TextInputFormat
Output format: org.apache.hadoop.hive.ql.io.HiveIgnoreKeyTextOutputFormat
Serde parameters: fileFormat:csv
Table properties:
recordCount: 2897 averageRecordSize: 519 CrawlerSchemaDeserializerVersion: 1.0 compressionType: none columnsOrdered: true
areColumnQuoted: false delimiter: , typeOfData: file

Columns

Column name	Data type	Partition key	Comment
1 student_id	bigint		
2 first_name	string		
3 last_name	string		
4 katz school major	string		
5 gpa	double		
6 graduation semester	string		
7 country of origin	string		
8 languages	string		
9 undergraduate major	string		
10 age at graduation	bigint		
11 years of experience	double		
12 location	string		
13 zip_code	string		

Showing: 1 - 25 of 25

AWS Glue Tables > salary_data.csv

Last updated 3 May 2022 08:40 PM Table Version (Current version)

Data catalog

Databases
Tables
Connections
Crawlers
Classifiers
Schema registries
Schemas
Settings

ETL

AWS Glue Studio (2)
Jobs (2) New
Jobs (legacy)
ML Transformations
Blueprints
Workflows
Triggers
Dev endpoints
Notebooks
Interactive Session - New

Schemas

salary_data.csv

Name: salary_data.csv
Description: db-final-project
Database: db-final-project
Classification: csv
Location: s3://fe-fns2022-csv/salary_data.csv
Connection: Connection
Dependencies: None
Last updated: Tue May 03 20:45:22 GMT+00:00 2022
Input format: org.apache.hadoop.hive.ql.io.HiveIgnoreKeyTextInputFormat
Output format: org.apache.hadoop.hive.serde2.lazy.LazySimpleSerDe
Serde parameters: field.delim: ,
Table properties:
recordCount: 91 averageRecordSize: 616 CrawlerSchemaDeserializerVersion: 1.0 compressionType: none columnsOrdered: true
areColumnQuoted: false delimiter: , typeOfData: file

Columns

Column name	Data type	Partition key	Comment
1 title	string		
2 location	string		
3 description	string		
4 mtc10	bigint		
5 mtc25	bigint		
6 mt650	bigint		
7 mt675	bigint		
8 mt690	bigint		

Showing: 1 - 8 of 8

AWS Glue Jobs (2)

Jobs: A job is your business logic required to perform extract, transform and load (ETL) work. Job runs are initiated by triggers which can be scheduled or driven by events.

User preferences

Add job Action (2) Filter by job and attributes

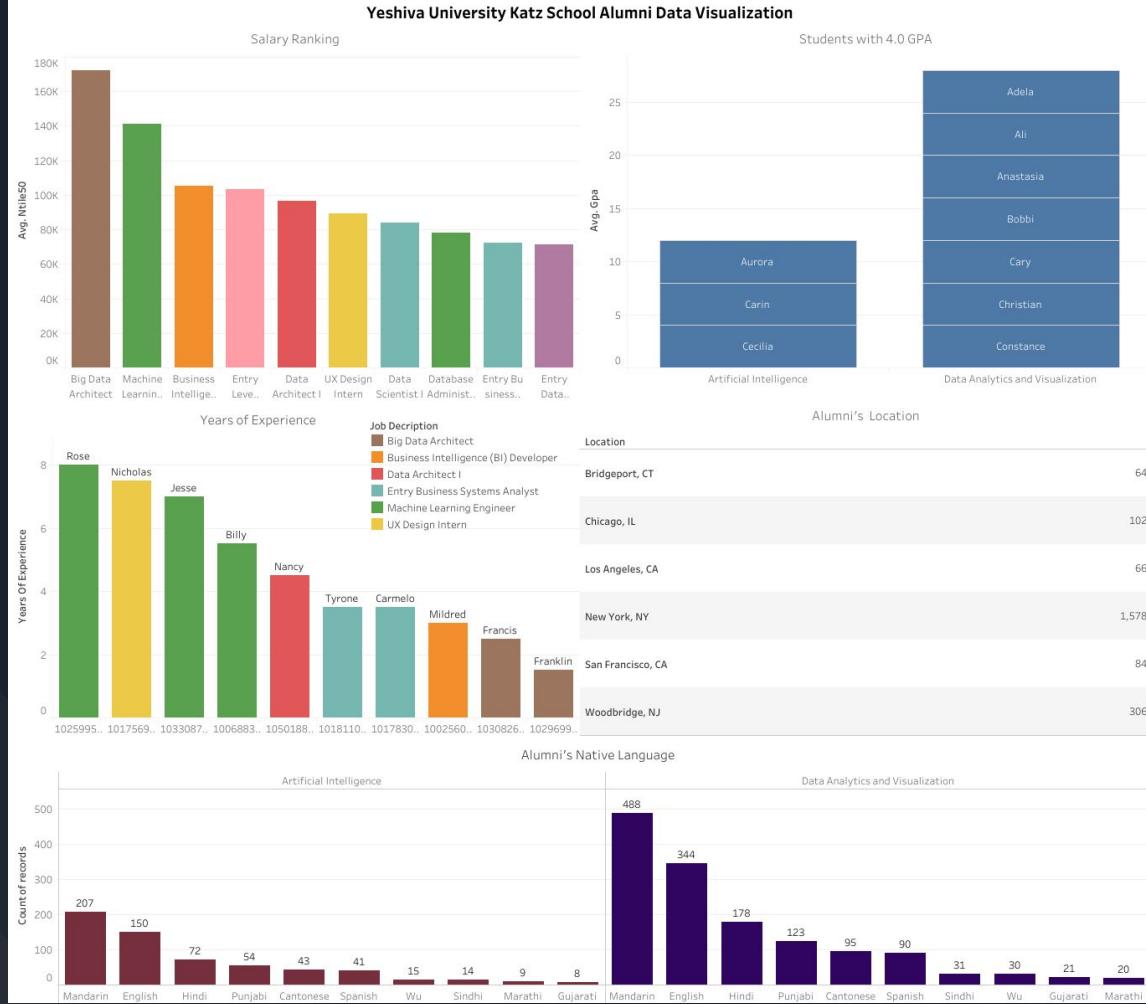
Name	Type	ETL language	Script location	Last modified	Job bookmark
job-salary-data	Spark	python	s3://aws-glue-s.../3 May 2022 08:40 PM UTC-4	Disable	
job-student-data	Spark	python	s3://aws-glue-s.../3 May 2022 08:40 PM UTC-4	Disable	

History

View run metrics Refresh job bookmarks

Showing: 1 - 1															
Run ID	Retry attempt	Status	Error	Output	Logs	Error logs	Glue version	Maximum capacity	Triggered by	Start time	End time	Start up time	Execution time	Timeout delay	Job run input
1_30fb0022fed...	-	Success	-	Log	Error Logs	2.0	10		3 May... 3 May...	1 min	2860 millis	3 May...	1 min	2860 millis	s3://aws-glue-t...

Tableau Dashboard





Challenges and Solutions

Challenges	Solutions
<p>Lambda function execution</p> <ul style="list-style-type: none">• Timeout Limit• Package Size Limit• Incompatible Dependencies <p>Neptune graph database</p> <ul style="list-style-type: none">• Node and edge selection• The connection between s3 and neptune database and notebook.• Dataset loading into JupyterLab notebook• CA Certificate failure	<p>Lambda function execution</p> <ul style="list-style-type: none">• Extend the timeout to 15 mins• Clean the script code• Run time setting from 3.9 to 3.8 <p>Neptune graph database</p> <ul style="list-style-type: none">• Identify entities and relationships• Create endpoint, add “all traffic” to inbound and outbound and connect to relevant VPC• Convert dataset format to Gremlin csv• Create the connection with gremlin.bat console.

Thanks for listening!

Any comments/questions?