

I worked on an end-to-end healthcare data engineering pipeline using AWS services to simulate a real-world scenario. I started by downloading a synthetic healthcare dataset from Kaggle, which contained patient demographics, medical history, hospital details, doctor, insurance provider, medication, admission types, billing amount, and test results. I manually uploaded this CSV dataset to Amazon S3, which served as my raw data storage layer. From there, I created an AWS Glue Crawler to scan the S3 folder and populate the Glue Data Catalog. This automatically inferred the schema from the CSV and allowed me to query the raw data using Amazon Athena. Athena converted the CSV into structured, queryable tables, which helped me inspect data quality, identify nulls and datatype mismatches, and finalize how I wanted to structure the data downstream.

Once the structure was validated, I created a Glue ETL job using PySpark. This job cleaned and transformed the raw dataset into a dimensional star schema. I extracted key columns to build four dimension tables: `dim_patient`, `dim_doctor`, `dim_hospital`, and `dim_insurance`. I used `.dropDuplicates()` and assigned clean surrogate keys where needed. Then I created the `fact_admissions` table, which included foreign keys to each of the dimensions along with transactional fields like room number, admission type, medication, test results, and billing amount. After transformation, I wrote each table back to S3 in **Parquet** format — a compressed, columnar format that's optimized for analytics.

Next, I moved to Redshift Serverless to load the transformed data warehouse. I created a new namespace and workgroup and then wrote SQL `CREATE TABLE` scripts that exactly matched the structure of the Parquet files. I used the `COPY` command to load each table from S3 into Redshift. But this step came with a few challenges. First, Redshift `COPY` failed until I correctly associated the IAM role with both the namespace and the workgroup, and set it as default. I also had to make the workgroup publicly accessible and allow inbound traffic on port 5439 via the associated security group to connect external tools like Power BI. Another challenge was Spark versioning — Glue was running Spark 3, which threw strict datetime parsing errors on some fields. I solved that by dropping or cleaning problematic rows. Redshift also threw schema mismatch errors during `COPY` — for example, my table had 5 columns but the Parquet had 6. I resolved that by using Athena and Glue Catalog to confirm the Parquet schema before recreating the table with exactly matching columns.

Once all the tables were successfully loaded into Redshift, I queried the schema using Query Editor v2 and verified joins between the fact and dimension tables. After validating the model, I connected Redshift to Power BI using the Redshift workgroup endpoint, the database name, and the admin user credentials. Once connected, I was able to use the structured tables from Redshift to build live visualizations in Power BI based on the star schema.

Blog reference: [End to End Data Engineering Project Using AWS Services | by Gemma ✨ | Medium](#)



the pipeline

Amazon S3

healthcare-data-project-rupak

General purpose buckets

Directory buckets

Table buckets

Access Grants

Access Points

Object Lambda Access Points

Multi-Region Access Points

Batch Operations

IAM Access Analyzer for S3

Block Public Access settings for this account

Storage Lens

Dashboards

Storage Lens groups

AWS Organizations settings

Feature spotlight 11

healthcare-data-project-rupak

Objects

Metadata

Properties

Permissions

Metrics

Management

Access Points

Objects (4)

Copy S3 URI

Copy URL

Download

Open

Delete

Actions

Create folder

Upload

Find objects by prefix

	Name	Type	Last modified	Size	Storage class
<input type="checkbox"/>	athena-results/	Folder	-	-	-
<input type="checkbox"/>	healthcare_dataset.csv	csv	April 14, 2025, 16:49:21 (UTC-04:00)	8.0 MB	Standard
<input type="checkbox"/>	transformed_\$(folder\$	-	April 15, 2025, 12:29:03 (UTC-04:00)	0 B	Standard
<input type="checkbox"/>	transformed/	Folder	-	-	-

The screenshot shows the AWS Glue console for a crawler named 'healthcare-crawler'. The crawler is in a 'READY' state. The configuration includes an IAM role of 'AWSGlueServiceRole-s3accessrupak', a database named 'healthcare\_db', and a table prefix. The crawler has one run completed on April 14, 2025, at 21:07:46, with a duration of 01 min 16 s and 0.039 DPU hours. The run resulted in 1 table change and 0 partition changes.

**healthcare-crawler** Last updated (UTC) April 16, 2025 at 21:10:42 [Run crawler](#) [Edit](#) [Delete](#)

**Crawler properties**

<b>Name</b> healthcare-crawler	<b>IAM role</b> <a href="#">AWSGlueServiceRole-s3accessrupak</a>	<b>Database</b> healthcare_db	<b>State</b> READY
<b>Description</b> -	<b>Security configuration</b> -	<b>Lake Formation configuration</b> -	<b>Table prefix</b> -
<b>Maximum table threshold</b> -			

[Advanced settings](#)

**Crawler runs (1)** [Stop run](#) [View CloudWatch logs](#) [View run details](#)

The list of crawler runs for this crawler.

	Start time (UTC)	End time (UTC)	Current/last duration	Status	DPU hours	Table changes
<input type="radio"/>	April 14, 2025 at 21:06:30	April 14, 2025 at 21:07:46	01 min 16 s	Completed	0.039	1 table change, 0 partition changes

The screenshot shows the Amazon Athena console's query editor. A SQL query is entered: 'SELECT \* FROM "healthcare\_db"."patient\_data" LIMIT 10;'. The left sidebar shows the data source as 'AwsDataCatalog', the catalog as 'None', and the database as 'healthcare\_db'. The top navigation bar includes 'Editor', 'Recent queries', 'Saved queries', and 'Settings'. The bottom status bar shows 'CloudShell' and 'Feedback'.

**Amazon Athena** > Query editor

[Introducing a new tutorial](#) Learn how to create a table from an S3 table bucket and run queries. [View tutorials](#)

[Editor](#) Recent queries Saved queries Settings **Workgroup** primary

[Athena now supports typeahead code suggestions to speed up SQL query development](#) Typeahead suggestions are turned on by default. You can change this setting in query editor preferences. [Edit preferences](#)

**Data**

**Data source**  
AwsDataCatalog

**Catalog**  
None

**Database**  
healthcare\_db

**Tables and views** [Create](#)

**Query 1**

```
1 SELECT * FROM "healthcare_db"."patient_data" LIMIT 10;
```

## Glue script

```
import sys
import hashlib
from pyspark.sql.functions import col, udf
from pyspark.sql.types import StringType
from awsglue.context import GlueContext
from awsglue.utils import getResolvedOptions
from awsglue.job import Job
from pyspark.context import SparkContext
```

```
# Glue job setup
args = getResolvedOptions(sys.argv, ['JOB_NAME'])
sc = SparkContext()
glueContext = GlueContext(sc)
spark = glueContext.spark_session
job = Job(glueContext)
```

```

job.init(args['JOB_NAME'], args)

# Load raw data from Glue Catalog
df = glueContext.create_dynamic_frame.from_catalog(
    database="healthcare_db",
    table_name="healthcare_data_project_rupak"
).toDF()

# Normalize column names
df = df.select([col(c).alias(c.strip().lower().replace(" ", "_")) for c in df.columns])

# Drop rows missing critical fields
df = df.dropna(subset=["name", "hospital"])

# ✅ Drop date columns to avoid Spark 3.x issues
df = df.drop("date_of_admission", "discharge_date")

# Create surrogate key generator
def gen_id(val):
    return hashlib.md5(str(val).encode('utf-8')).hexdigest()
hash_udf = udf(gen_id, StringType())

# Add surrogate keys
df = df.withColumn("patient_id", hash_udf(col("name") + col("gender") + col("age").cast("string")))
df = df.withColumn("doctor_id", hash_udf(col("doctor")))
df = df.withColumn("hospital_id", hash_udf(col("hospital")))
df = df.withColumn("insurance_id", hash_udf(col("insurance_provider")))
df = df.withColumn("admission_id", hash_udf(col("name") + col("room_number").cast("string")))

# Build dimension tables
dim_patient = df.select("patient_id", "name", "age", "gender", "blood_type", "medical_condition").dropDuplicates()
dim_hospital = df.select("hospital_id", col("hospital").alias("hospital_name")).dropDuplicates()
dim_doctor = df.select("doctor_id", col("doctor").alias("doctor_name")).dropDuplicates()
dim_insurance = df.select("insurance_id", col("insurance_provider").alias("provider")).dropDuplicates()

# Build fact table (no admission_date/discharge_date)
fact_admissions = df.select(
    "admission_id", "patient_id", "doctor_id", "hospital_id", "insurance_id",
    "admission_type", "room_number", "billing_amount", "test_results", "medication"
)

# Write transformed tables to S3 in Parquet format
output_path = "s3://healthcare-data-project-rupak/transformed/"
dim_patient.write.mode("overwrite").parquet(output_path + "dim_patient/")
dim_hospital.write.mode("overwrite").parquet(output_path + "dim_hospital/")
dim_doctor.write.mode("overwrite").parquet(output_path + "dim_doctor/")
dim_insurance.write.mode("overwrite").parquet(output_path + "dim_insurance/")
fact_admissions.write.mode("overwrite").parquet(output_path + "fact_admissions/")

# Finish the job
job.commit()

```

aws

Search

[Alt+S]

United States (N. Virginia)

rupakkulkarni17

AWS Glue

Getting started

ETL jobs

Visual ETL

Notebooks

Job run monitoring

Data Catalog tables

Data connections

Workflows (orchestration)

Zero-ETL integrations

Data Catalog

Databases

Tables

Stream schema registries

Schemas

Connections

Crawlers

Classifiers

Catalog settings

Data Integration and ETL

Legacy pages

healthcare-etl-job

Last modified on 4/15/2025, 1:04:23 PM

Actions

Save

Run

Script

Job details

Runs

Data quality

Schedules

Version Control

Job runs (1/8)

Info

Last updated (UTC)

April 16, 2025 at 21:13:18

View details

Stop job run

Troubleshoot with AI

Table View

Card View

Filter job runs by property

Run status

Retries

Start time (Local)

End time (Local)

Duration

Capacity (...)

Worker type

Glue

Succeeded

0

04/15/2025 13:04:25

04/15/2025 13:06:04

1 m 32 s

3 DPUs

G.1X

5.0

Failed

0

04/15/2025 12:58:06

04/15/2025 13:00:04

1 m 51 s

3 DPUs

G.1X

5.0

Run details

Input arguments (10)

Continuous logs

Run insights

Metrics

Troubleshooting analysis - preview

Spark UI

Job name

healthcare-etl-job

Id

jr\_9d550d600cd6ef1c3372a69475313cc2aa82ba9de2570d00e87513d756622a4c

Run status

Succeeded

Start time (Local)

04/15/2025 13:04:25

End time (Local)

04/15/2025 13:06:04

Start-up time

6 seconds

Glue version

5.0

Worker type

G.1X

Max capacity

3 DPUs

Last modified on (Local)

04/15/2025 13:06:04

Log group name

/aws-glue/jobs

Number of workers

3

aws

Services

Search

[Alt+S]

United States (N. Virginia)

rupakkulkarni17

You can now register your Redshift table definitions to AWS Glue Data Catalog, making them accessible from other query engines such as EMR Serverless, Athena and other Redshift data warehouses. To get started registering a namespace, select a Redshift serverless namespace from the namespaces, then select "Register with AWS Glue Data Catalog" from the Actions menu.

Learn more

Amazon Redshift Serverless

Namespace configuration

test2

test2

Info

Actions

Query data

General information

Namespace

test2

Namespace ID

36b5d0c9-2d34-4235-8cd9-8e587f458765

Namespace ARN

arn:aws:redshift-serverless:us-east-1:633263914408:namespace/36b5d0c9-2d34-4235-8cd9-8e587f458765

Namespace register status

Deregistered

Status

Available

Date created

April 15, 2025, 14:43 (UTC-04:00)

Total used storage

2.1 GB

Admin user name

admin

Database name

dev

Total table count

5

aws

Services

Search

[Alt+S]

United States (N. Virginia)

rupakkulkarni17

Redshift query editor v2

Create

Load data

Filter resources

Serverless: healthcare-rupak

Serverless: test2

native databases (2)

external databases (1)

Serverless: default-workgroup

Untitled 1

Run

Limit 100

Explain

Isolated session

Serverless: test2

dev

Schedule

1

SELECT \* FROM dim\_patient LIMIT 10;

Result 1 (10)

Export

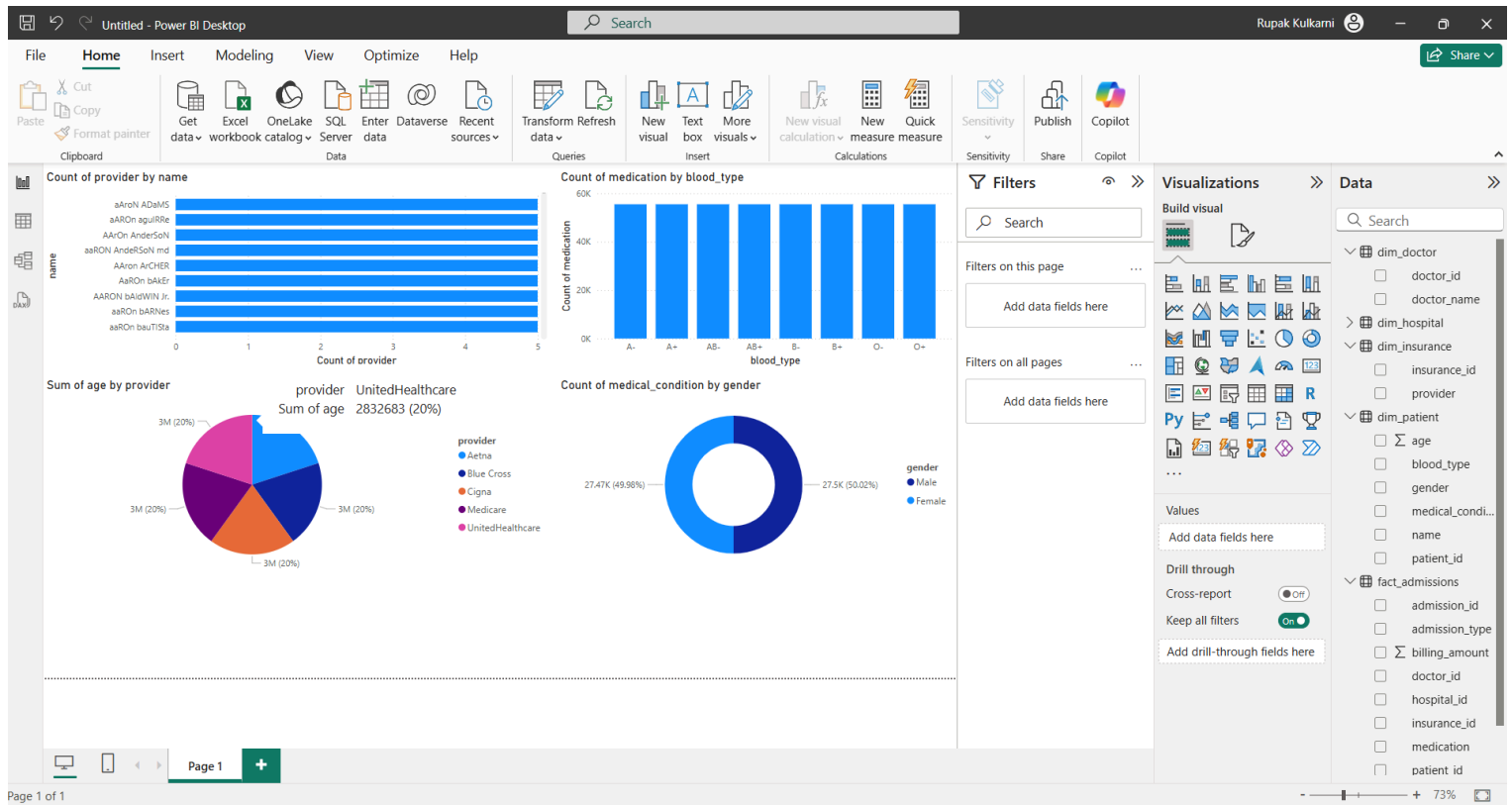
Chart

patient_id	name	age	gender	blood_type	medical_conditi
6ad97f83acf6453d4a6a4...	cAniY acosta	18	Female	AB-	Diabetes
6ad97f83acf6453d4a6a4...	JAsmlne Wiley	68	Female	A-	Obesity
6ad97f83acf6453d4a6a4...	dAKOTA weBsTeR	30	Male	AB-	Obesity
6ad97f83acf6453d4a6a4...	stEpHanIE MURILlo	20	Male	AB-	Arthritis
6ad97f83acf6453d4a6a4...	PATRicIA potTeR	38	Male	O-	Hypertension
6ad97f83acf6453d4a6a4...	LindA bARNETt	21	Female	A+	Arthritis
6ad97f83acf6453d4a6a4...	JimMY MathiS	59	Female	AB+	Hypertension

Query ID 6943

Elapsed time: 99 ms

Total rows: 10



aws

Search

[Alt+S]

Global rupakkulkarni17

Identity and Access Management (IAM)

Search IAM

Dashboard

Access management

- User groups
- Users
- Roles**
- Policies
- Identity providers
- Account settings
- Root access management [New](#)

Access reports

- Access Analyzer
- Archive rules
- Analizers
- Settings

### Roles (7) Info

An IAM role is an identity you can create that has specific permissions with credentials that are valid for short durations. Roles can be assumed by entities that you trust.

Search

<input type="checkbox"/>	Role name	Trusted entities	Last activity
<input type="checkbox"/>	<a href="#">AmazonRedshift-CommandsAccessRole-20250415T134644</a>	AWS Service: redshift-serverless, <a href="#">anc</a>	-
<input type="checkbox"/>	<a href="#">AmazonRedshift-CommandsAccessRole-20250415T134659</a>	AWS Service: redshift-serverless, <a href="#">anc</a>	Yesterday
<input type="checkbox"/>	<a href="#">AWSGlueServiceRole-s3accessrupak</a>	AWS Service: glue	Yesterday
<input type="checkbox"/>	<a href="#">AWSGlueServiceRole-s3csvaccess</a>	AWS Service: glue	-
<input type="checkbox"/>	<a href="#">AWSServiceRoleForRedshift</a>	AWS Service: redshift ( <a href="#">Service-Linker</a> )	18 minutes ago
<input type="checkbox"/>	<a href="#">AWSServiceRoleForSupport</a>	AWS Service: support ( <a href="#">Service-Linker</a> )	-
<input type="checkbox"/>	<a href="#">AWSServiceRoleForTrustedAdvisor</a>	AWS Service: trustedadvisor ( <a href="#">Service</a> )	-

### Roles Anywhere Info

Authenticate your non AWS workloads and securely provide access to AWS services.

[Manage](#)

CloudShell Feedback

© 2025, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences