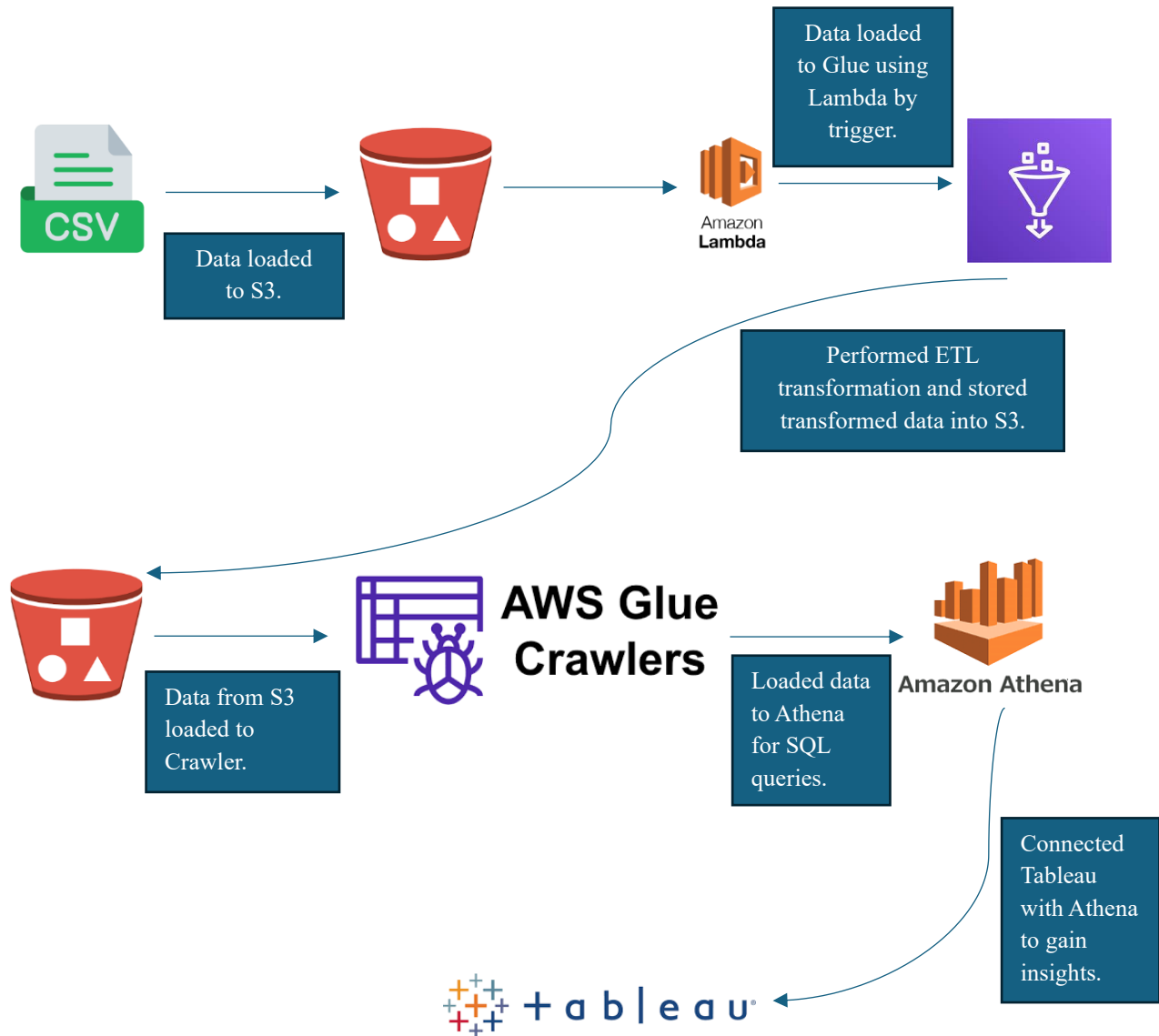# Milestone 6

## Data Architecture



## Data Ingestion

- Prepared structured data in CSV format.
- Created an S3 bucket, load csv file to S3.
- Used AWS lambda to fetch data from S3 and pass it to AWS glue.
- Used AWS glue for ETL to transform data.
- After transformation, sent data to the new S3 bucket.
- Loaded the transformed data to AWS Athena by AWS Glue Crawler.
- Used Athena to run queries and get granular insights from the data.

- Connected Athena with Tableau to visually get some insights.

us-east-1.console.aws.amazon.com/athena/home?region=us-east-1#/query-editor/history/1da0b954-f5b7-4415-b0be-c6c112278067

Apps ▶ YouTube 🔗 LinkedIn 📧 N Canvas ⬤ Meet Vasani ▶ Nuworks ▮ dataviz | 🗀 All Bookmarks

aws ⠿ Services | Q Search [Alt+S] | N. Virginia ▾ | airbnb-project @ 6374-2326-7574 ▾

Q Filter tables and views

▼ Tables (1) ‹ 1 ›

⊞ transformedairbnbdata  Partitioned ⋮

▶ Views (0) ‹ 1 ›

SQL  Ln 1, Col 1

Run again | Explain ⬈ | Cancel | Clear | Create ▾  ⬤ Reuse query results
up to 60 minutes ago ✎

Query results | Query stats

⊘ Completed | Time in queue: 74 ms  Run time: 1.012 sec  Data scanned: 336.09 KB

Results (3,585) | 📋 Copy | Download results

🔍 Search rows  ‹ 1 … › ⚙

| # ▽ | id ▽ | listing_url | scrape_id ▽ | last_scraped ▽ | name ▽ | experiences_o |
|---|---|---|---|---|---|---|
| 1 | 0 | https://www.airbnb.com/rooms/12147973 | 2.02E+13 | 9/7/2016 | Sunny Bungalow in the City | none |
| 2 | 1 | https://www.airbnb.com/rooms/3075044 | 2.02E+13 | 9/7/2016 | Charming room in pet friendly apt | none |
| 3 | 2 | https://www.airbnb.com/rooms/6976 | 2.02E+13 | 9/7/2016 | Mexican Folk Art Haven in Boston | none |
| 4 | 3 | https://www.airbnb.com/rooms/1436513 | 2.02E+13 | 9/7/2016 | Spacious Sunny Bedroom Suite in Historic Home | none |
| 5 | 4 | https://www.airbnb.com/rooms/7651065 | 2.02E+13 | 9/7/2016 | Come Home to Boston | none |
| 6 | 5 | https://www.airbnb.com/rooms/12386020 | 2.02E+13 | 9/7/2016 | Private Bedroom + Great Coffee | none |
| 7 | 6 | https://www.airbnb.com/rooms/5706985 | 2.02E+13 | 9/7/2016 | New Lrg Studio apt 15 min to Boston | none |
| 8 | 7 | https://www.airbnb.com/rooms/2843445 | 2.02E+13 | 9/7/2016 | Tranquility on Top of the Hill | none |
| 9 | 8 | https://www.airbnb.com/rooms/753446 | 2.02E+13 | 9/7/2016 | 6 miles away from downtown Boston! | none |

⬜ CloudShell  Feedback | © 2024, Amazon Web Services, Inc. or its affiliates.  Privacy  Terms  Cookie preferences