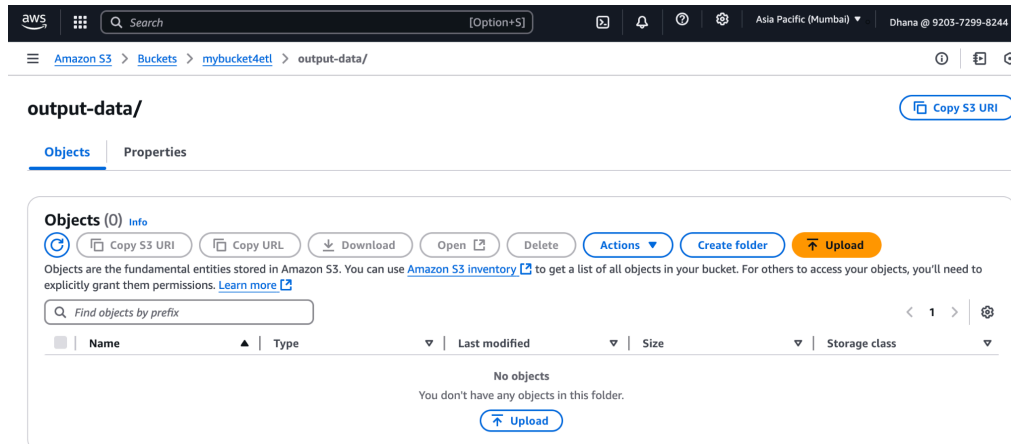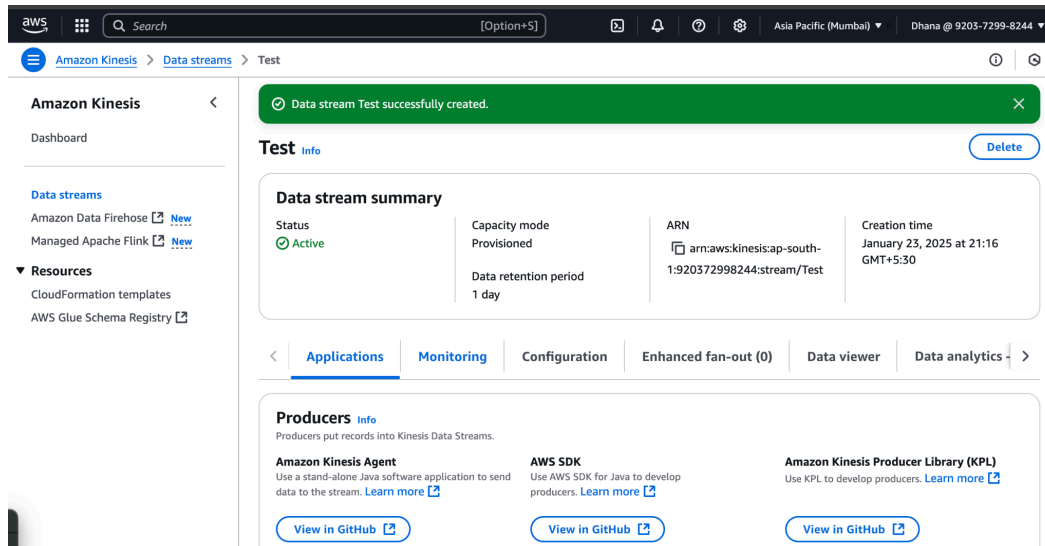For real-time streaming, we are using **Kinesis Data Generator**.

First, we will create an S3 folder named **output-data**, as shown below. Next, we will create an IAM role with full access to **AWS Glue**, **Kinesis**, **Athena**, and **S3**.
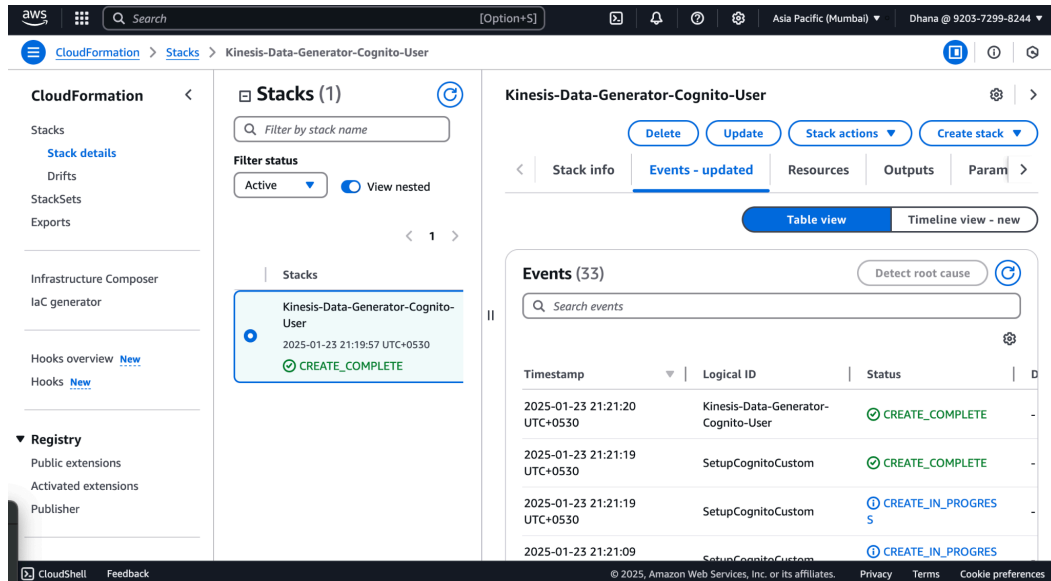


We will create a **Kinesis Data Stream** for real-time data streaming.



Next, we will create a **Cognito user** for **Kinesis Data Generator** using a CloudFormation stack. You can use the following link to create it:
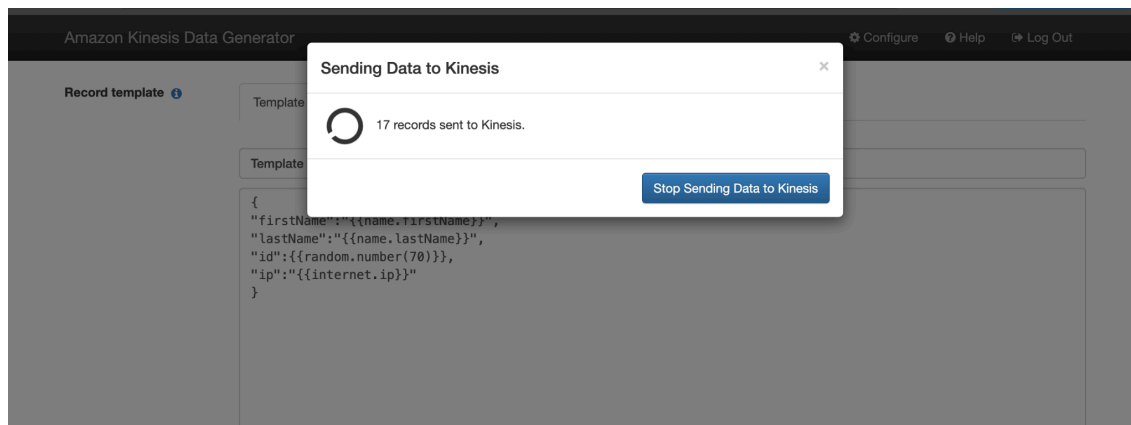https://awslabs.github.io/amazon-kinesis-data-generator/web/help.html

Once the stack creation is complete, log in to **Kinesis Data Generator using the link from Resources tab in CloudFormation** and start streaming data.

Use the following code for data generation:

```
{
"firstName":"{{name.firstName}}",
"lastName":"{{name.lastName}}",
"id":{{random.number(70)}},
"ip":"{{internet.ip}}"
}
```



Now, create a **Glue job** with the **source** set to the Kinesis Data Stream we created, and the **target** set to the S3 location.

**Edit** job Details and **Run** the job.



Once job is completed we should see the files placed in S3 from streaming data.

Next, we will create a **database** in **Glue** and then **create Crawler** and **run the crawler for table creation.**

And then Use **Athena** to query the data