# CUSTOMER
# DATA INGESTION

# TABLE OF CONTENTS

# INTRODUCTION
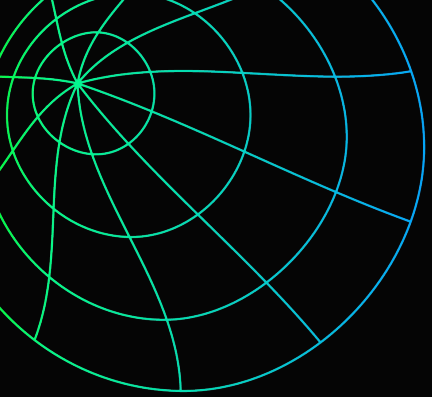
This is a data management project that involves collecting, processing, and storing data from various sources in a centralized location for further analysis.
The Goal of this project is to make data readily available for analysis by extracting data from variety of sources and integrating it into a single repository, such as a data warehouse
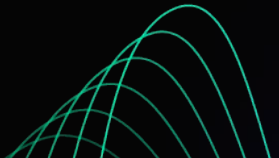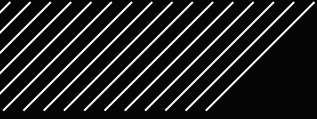
# 01

## DATA STORAGE
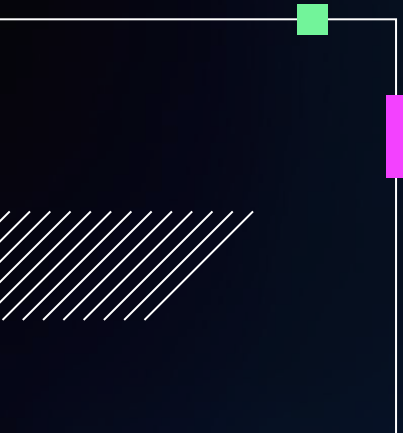
# STORAGE USED

## AWS S3

AWS S3 (Amazon Simple Storage Service) is a a highly scalable ,secure, and durable cloud storage service.
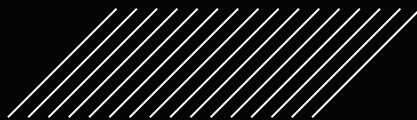
## MongoDB

A popular document-oriented NoSQL database that allows for flexible and scalable data management

# DATA SECURITY
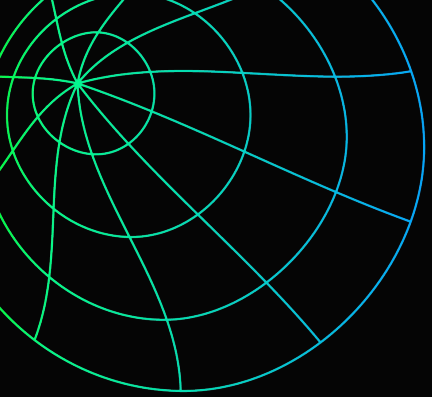
02

# SECURITY OF DATA

## VPC

Virtual Private Cloud is a customizable and secure virtual network within AWS that allows users to launch AWS resources in a defined virtual network environment which is isolated

## Authenticated Users

Authenticated users are individuals or entities who have successfully proven their identity and been granted access to a system or application with the appropriate permissions.
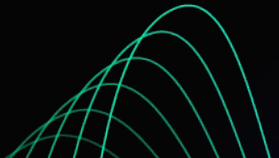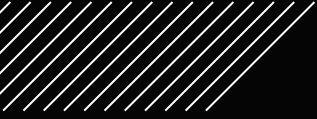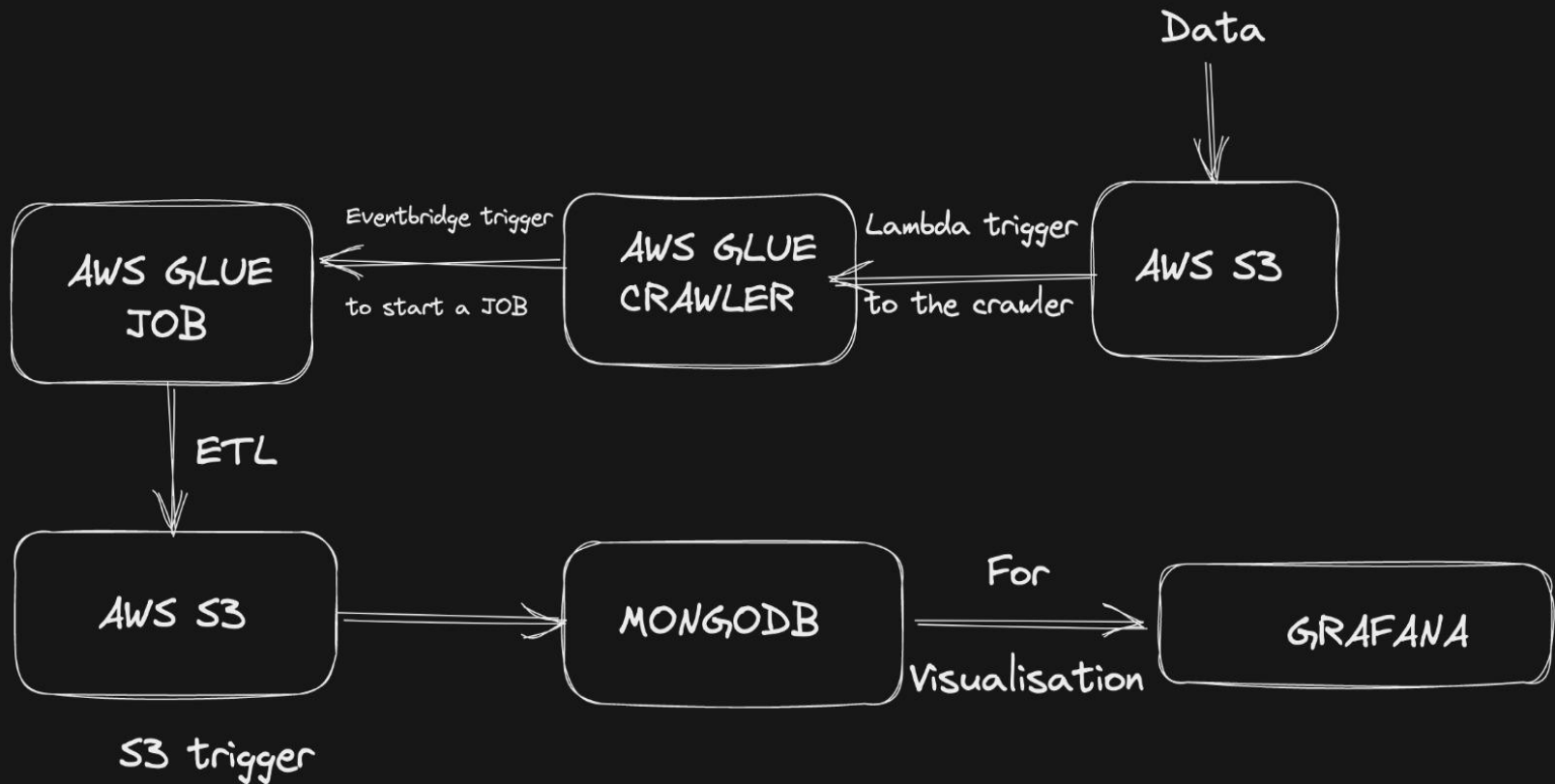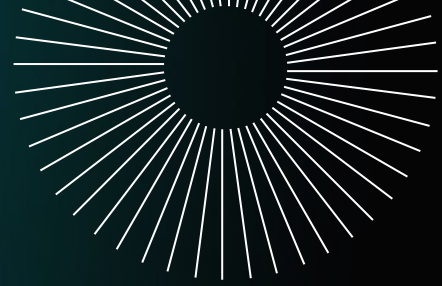
# 03

## DATA INGESTION

# INGESTION FLOW

# GET THE API KEY

## Customer Data Ingestion Pipeline

☑ API Key Copied ✕

**Email :**

mayank@gmail.com

**Name :**

mayankjha

**Organization :**

India

**User Authenticated**

### File upload status

1. S3 Status:

2. Glue Status:

3. Database Status:

# UPLOAD FILE FROM POSTMAN

AuthServer / Bucket Upload

Save ▾   •••   ✏️ 💬

POST ▾   https://orguser.singhharshit.me/api/bucketUpload   **Send** ▾

Params   Authorization   Headers (9)   **Body** ●   Pre-request Script   Tests   Settings   Cookies

○ none   ● form-data   ○ x-www-form-urlencoded   ○ raw   ○ binary   ○ GraphQL
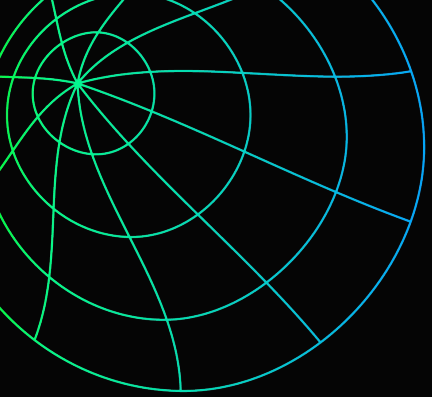
| | Key | Value | Description | ••• Bulk Edit |
|---|---|---|---|---|
| ☑ | file | datatesting.csv ✕ | | |
| | Key | Value | Description | |

Body   Cookies   Headers (18)   Test Results   🌐   Status: 200 OK   Time: 725 ms   Size: 889 B   💾 Save as Example   •••
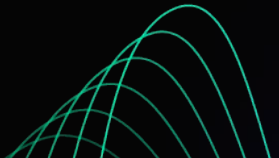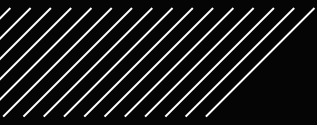
Pretty   Raw   Preview   Visualize   HTML ▾   ⇥

```
1  File uploaded successfully to S3
```
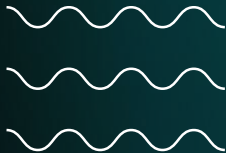
# 04

## DATA PROCESSING

# APACHE SPARK AS ETL FRAMEWORK

Apache Spark is a powerful ETL (Extract, Transform, Load) framework that can efficiently process large volumes of data from various sources and transform it into a desired format for analysis and reporting.

# CRAWLER RUNS

Objects (1)

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

| | Name | ▲ | Type | ▽ | Last modified | ▽ | Size | ▽ | Storage class | ▽ |
|---|---|---|---|---|---|---|---|---|---|---|
| | parsed.json | | json | | April 1, 2023, 10:25:20 (UTC+05:30) | | 392.8 KB | | Standard | |

Copy S3 URI | Copy URL | Download | Open | Delete | Actions | Create folder | Upload

Find objects by prefix

1

---

customerpipeline.customerData

998 DOCUMENTS    1 INDEXES

Documents    Aggregations    Schema    Explain Plan    Indexes    Validation

Filter    Type a query: { field: 'value' }    Reset    Find    More Options

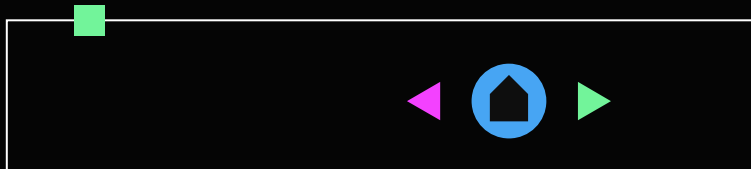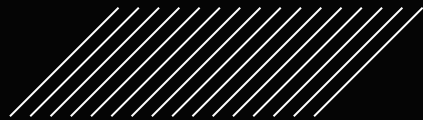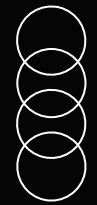ADD DATA    EXPORT COLLECTION    1 – 20 of 998

```
_id: ObjectId('6427ad9a53860be70603d1c5')
CustomerName: "Yash"
CustomerID: "C5841053"
CustGender: "F"
CustomerDOB: "10-01-1994"
Martial Status: "married"
CustLocation: "JAMSHEDPUR"
Job: "management"
Education: "tertiary"
loan: "no"
CustAccountBalance: "17819.05"
Last TransactionDate: "02-08-2016"
Last TransactionTime: "143207"
LastTransactionDetails: "to JHAANU@okicici via mayankoksbi"
```
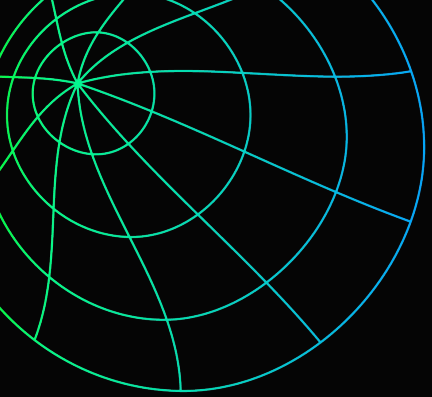
# DATA SCALABILITY

05

# PROPERTIES FOR SCALABILITY

Serverless

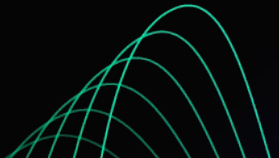Event Driven

Big Data Processing

Resilient

# 06

## DATA ANALYSIS

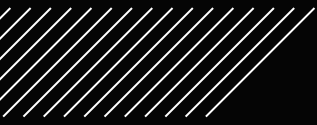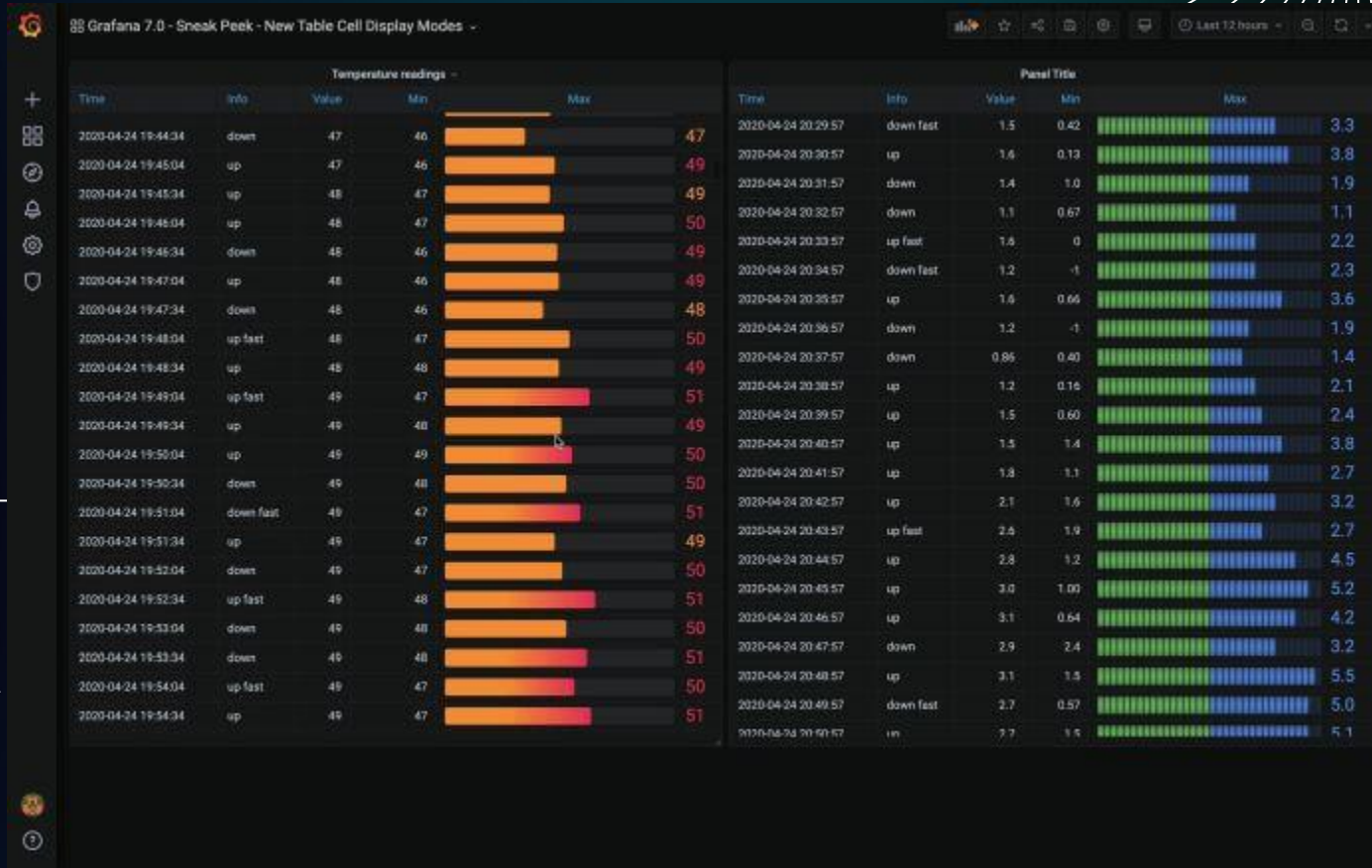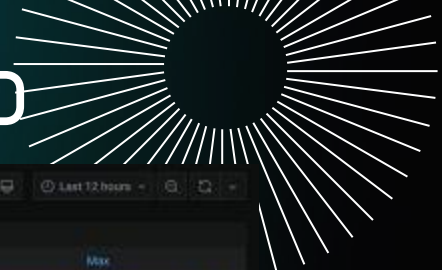# GRAFANA TO BE IMPLEMENTED

# OUR TEAM

Anupama Jha

Harshit Singh

Mayank Jha

Shashank Kumar