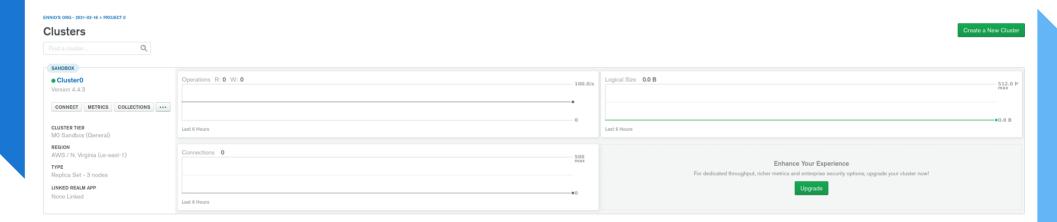# BTS – MBDS
# Big Data Infrastructure

## A2: MongoDB and DataBricks

Ennio Maldonado
18 February, 2021

# MongoDB

MongoDB cluster creation
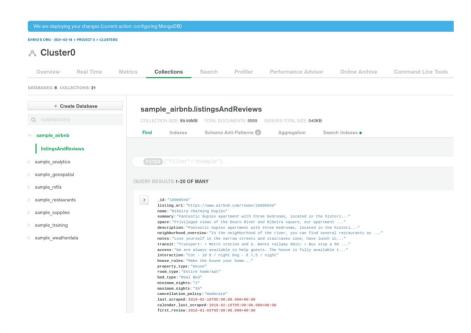
# MongoDB

Installed mongodb shell tool

# MongoDB

Loaded sample data_set from mongodb for testing connection.

# MongoDB on DataBricks

Connect MongoDB to DataBricks followed this guide:

https://docs.databricks.com/data/data-sources/mongodb.html