# Individual Report

*AWS Infrastructure Setup:*
In the AWS infrastructure setup, I played a role that complemented our team's efforts. My focus was on configuring EC2 instances and S3 buckets, ensuring they met our project's requirements. This task was done in tandem with the team's strategy on IAM roles and user management. Together, we navigated through security challenges, analyzing access needs and establishing a secure yet user-friendly setup. My involvement was a blend of contributing individual insights and working within the team's framework.
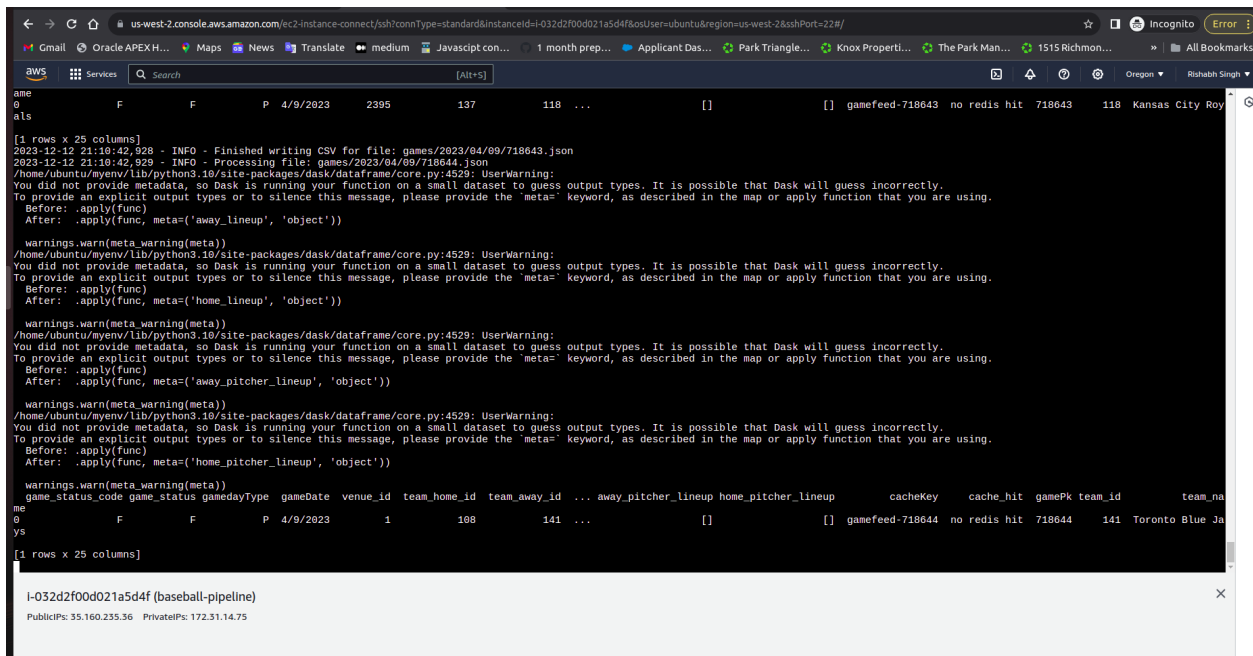
*Data Processing with PySpark and Dask:*
For data processing, I was involved in integrating PySpark and Dask with our existing processes. My work centered around refining data transformation scripts, making sure they were aligned with the project's scale and complexity. This involved a deep understanding of both the tools and the data, ensuring our methods were efficient and in line with the team's strategy.

The decision to use Dask was a collective one, and I contributed by helping adapt this technology into our workflow. Implementing Dask's parallel processing capabilities helped improve our data handling efficiency, a task that required both individual technical skills and a collaborative approach.

There was analysis done through google colab by reading in the cleaned data for player engagement for fans. Used the python visualization tools for easy understanding of the results and analysis derived

Here is one of the screenshots of the logs: