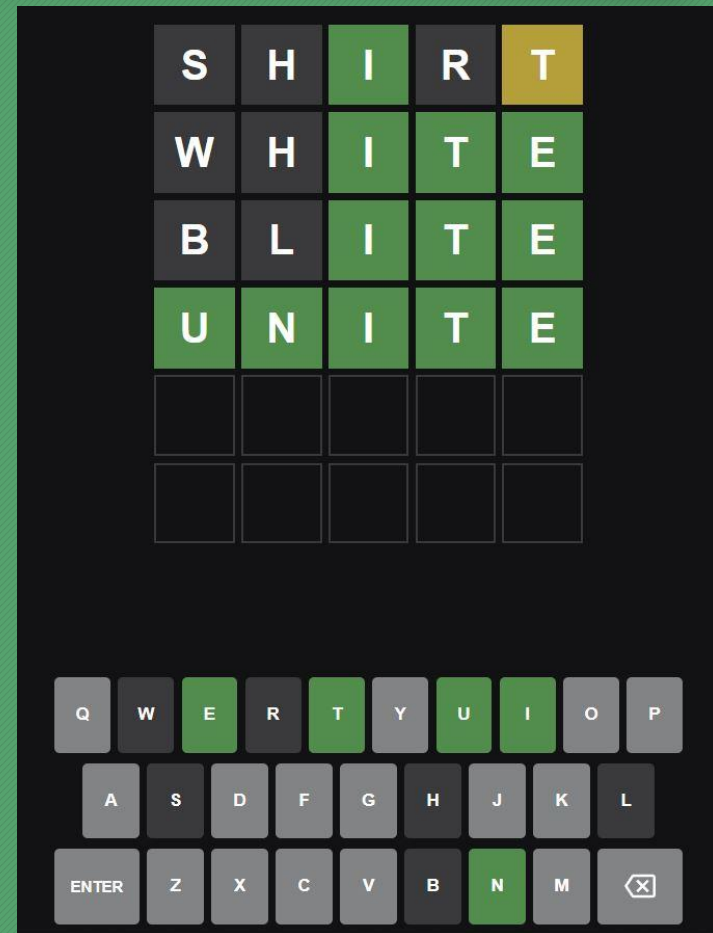# Wordle Match

Team 1 – Ayush Venkatesh, Olasunkanmi Olayinka,
Sai Sumana Puppala, Jared Girouard

# Wordle

- 5-letter word
- 6 guesses
- No color – letter not in word
- Yellow – letter in word but wrong position
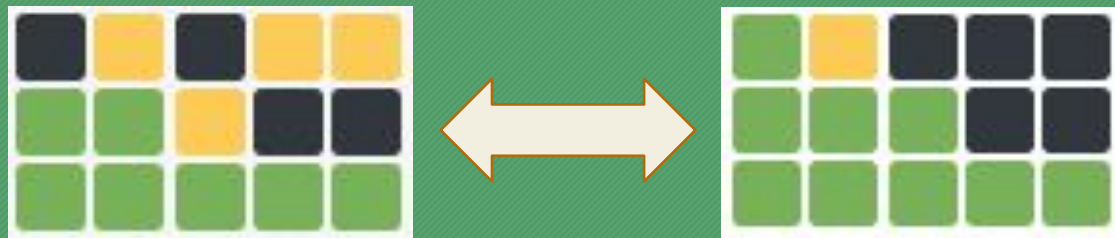- Green – letter in word in correct position

# Use Cases

- System queries wordle tweets once each day.
- Top 10 User Search
  - Given a user, provide top 10 users with similar game pattern results for a specific day
- On a given day, compute user rankings based on scores.

# Methodology

- Data Collection and Analysis
  - Extract wordle game results of users from Twitter API, over a period.
  - Clean and transform data to get 5*6 pattern matrices of game results.

- User Match
  - Compute user similarity scores among user game result matrices using a Euclidean distance metric.

- Dashboards
  - Build dashboard to display top 10 similar users.

# Dashboard

# Architecture

# Data Sources

- Data was extracted via twint, a twitter tool, each day.
  - User who posted the tweet
  - Created date
  - Wordle score
  - Date/Time
- Live stream of data preprocessed before storing it in mysql database.
- Extracted data for the last 8 days, accounting to 117k user distances.

# Data Sources: Raw Stream

# Data Transformation

- Data extracted via twint is extremely raw; unicode formatted data. A series of transformations done to compute wordle sequence vectors.
  - Regex transformations: series of regex extract and replacements to decode unicode characters.
  - Formatting the extracted string to the right format to computer vectors; fill a partial wordle result to a full flattened matrix (5*6): vector size 30
  - Computing a vector of doubles for distance calculation from a string sequence representing wordle result.
- Compute all pair user distances for users each day before writing to database

# Snapshot of transformed data

# Tests: Parse Tweets

# Tests: Ingest Tweets

# Milestones

- Sprint 1 (11/6–11/12): Define requirements and create repository
- Sprint 2 (11/13 – 11/19): Data collection
- Sprint 3 (11/20 – 11/26): Data collection and data cleaning
- Sprint 4 (11/27 – 12/3): Data visualization
- Sprint 5 (12/4 – 12/8): Finalize Data visualization

# Use of Scala

- Spark consumer gets data from Kafka
- Scala processes and transforms before writing it to the database
- Spark parallelly computes user similarity scores using Euclidean distance (best metric evaluated).
- Databricks visualizes data with Scala.

# Clustering of users - Additional

- Given the users, and the distance between each pair of users, use the k-means clustering algorithm to find clusters.

- Use the in-built mllib.clustering library in Databricks.

- Idea is to find the optimal number of centroids by iterating through a range, computing the error (SSE).

- Find the number of centroids which produce the least error.

- Using this, visualize the clusters graphically.

# Goals

- Demonstrate uses of Scala and Spark by creating an end-to-end project that extracts, transforms and visualizes data from an API.
- Create a similarity matrix among all users that shows how similar a player's wordle attempts were to other players.
- Visualize the top 10 performing players for a given user.
- (If time permits) Create graphs of commonly occurring words. (Did not achieve)

# Acceptance Criteria

- 1 million rows are extracted. (Actual: 2500 per day, more than 117k rows after distance calculations)

- User search response in 5 seconds. (Actual: 4.73 seconds on Databricks community)

- Distance metric gives a similarity score of at least 75% in the range of a base metric. (Actual: N/A since we chose to use Euclidean distance)

- Data is extracted twice a day. (Actual: Data is extracted once a day, user distances calculated eod)