# DATA X: FUZZY JOINS

Project 38, Group 2

# MEET THE TEAM

**Jake Mainwaring**

**Chase Smith**

**Cyril Tamraz**

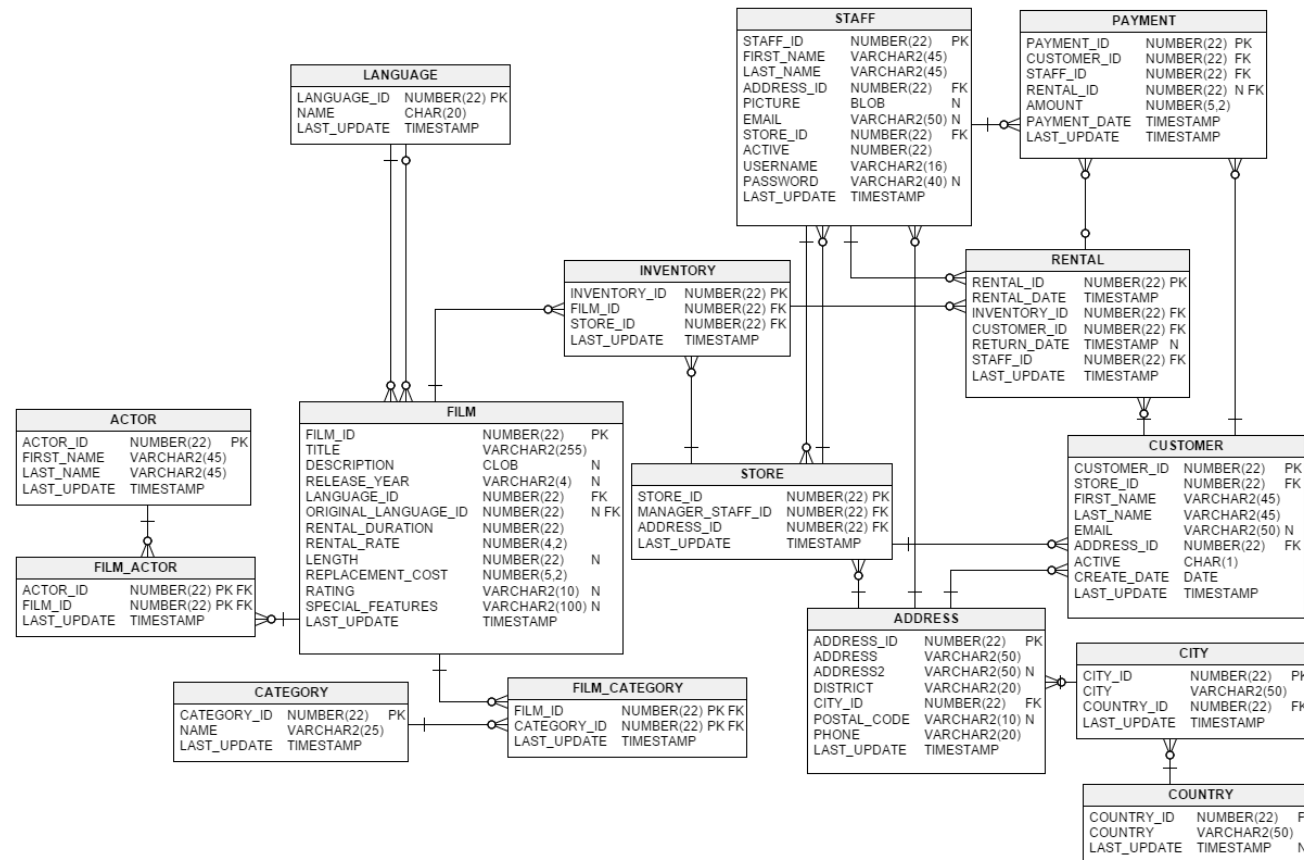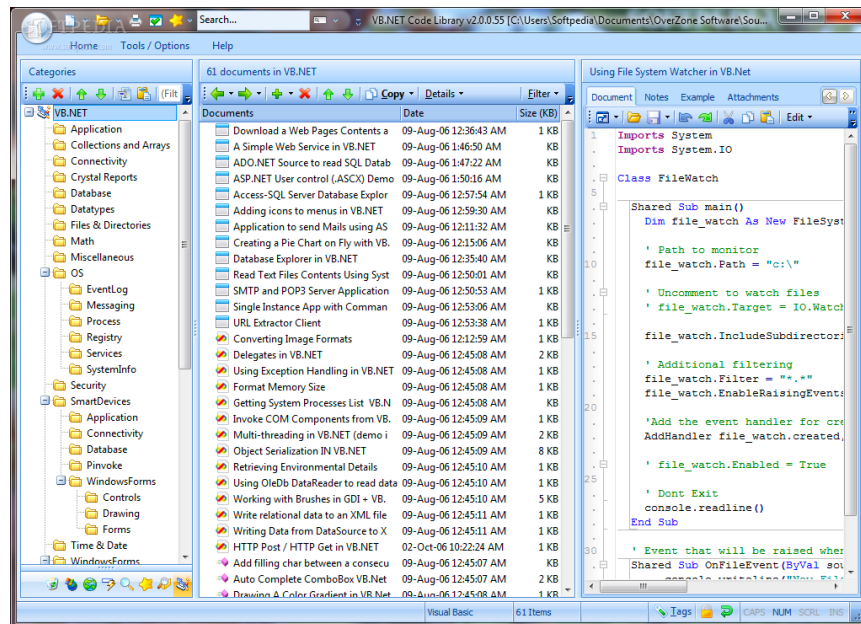**Kochise Bennett**

# HIGH LEVEL DESCRIPTION



**Current problem** – companies rely on a large number of databases to store information, but the majority of time is spent cleaning and aggregating rather than conducting analysis
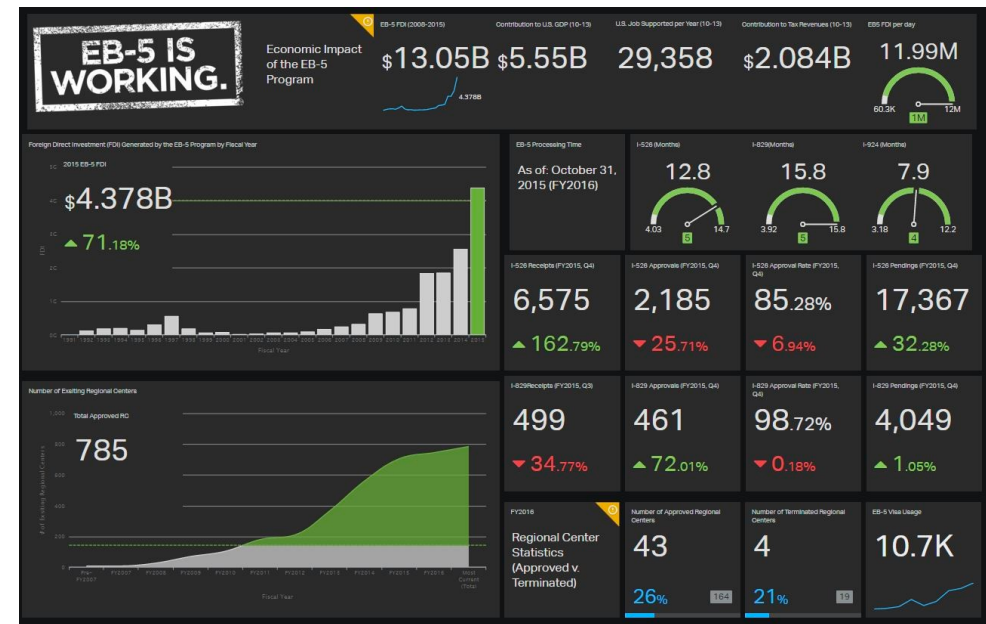
**Our solution** – create a tool that helps users standardize and aggregate data based on fuzzy comparisons. This could include minor differences in column name (e.g. "User_ID" vs. "UserID") specific values (e.g. "123 MLK Jr. Way" vs. "123 Martin Luther King Jr. Way"), or other comparisons

# USER PERSPECTIVE

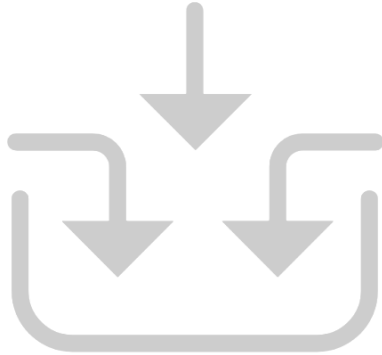*Option 1 (library):*

*Option 2 (dashboard):*

# TECHNICAL COMPONENTS

- Development of fuzzy join/inexact matching criteria:
  - Column name comparison
  - String values within columns
  - Dictionary based on standardized column names, where users can input "User-id", "User-ID", etc. and it maps to the standard name for that column
- Joining data tables on cleaned column names and values
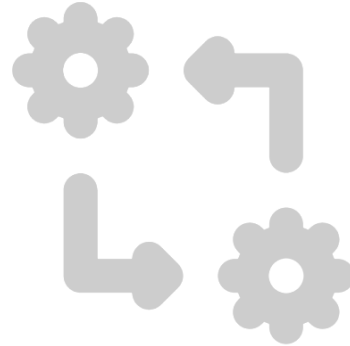- UI-database connector
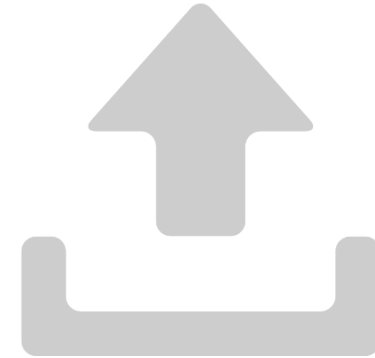
# SAMPLE ARCHITECTURE/DATA MODEL

*Input*

1. Collection of two or more data tables

*Process*

1. Fuzzy match foreign keys with corresponding primary key column names
2. Fuzzy match column names between primary key tables with overlapping column names
3. Fuzzy match specific values based on similarity metrics

*Output*

1. Updated data tables with dictionary of column names

# NEXT STEPS