



Final Project  
Course- Information Architecture  
(DAV-6100)

**Presented By - Group 2**



Yeshiva University | KATZ SCHOOL

# CONTENTS



1

Overview

2

About Data

3

Architecture

4

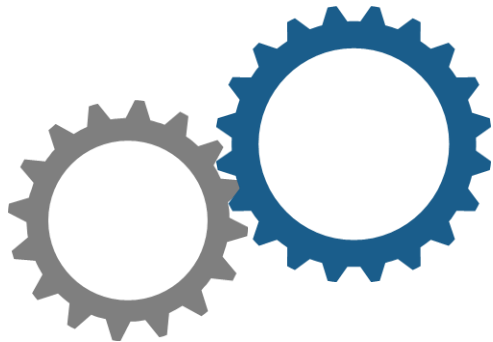
OLAP & OLTP

5

Reporting & Analysis

6

Challenges



# Stack Overflow Data Warehouse Solution

## Project Overview

- Our project is focused around building a data warehouse solution for the Stack Overflow platform
- The data warehouse will help Stack Overflow achieve their analytical needs like user interests, membership info, active user info & trending technologies/topics etc.
- Stack Overflow can improve their advertising and marketing revenue model by analyzing these users activities, interest and engagement metrics
- This can also help Stack Overflow with providing better job recommendations to end users who are actively looking for change. This improved jobs listing model can help them to increase ad revenue
- The data warehouse will provide a more comprehensive view of Stack Overflow and will provide data which will (in theory) be used to improve site performance



# Roles and Responsibilities

## **Surbhi - AWS Architect:**

Project lead and system architect.

## **Nosson - Data Engineer:**

Collect relevant data according to the needs of the system and users.

## **Yihang - Database Administrator:**

Operate, maintain and manage database management systems.

## **Yifeng - Database developer:**

Design and develop database management systems.

## **Qianwen - Data Analyst:**

Analyze and visualize data, and make industry assessments and forecasts based on data.



# CONTENTS



1

Overview

2

About Data

3

Architecture

4

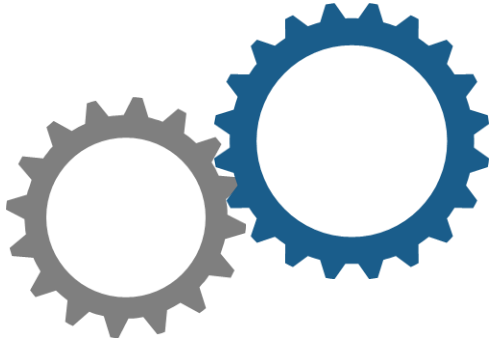
OLAP & OLTP

5

Reporting & Analysis

6

Challenges



# About Data

**Stack Overflow** - Stack Overflow is a question and answer website for professional and enthusiast programmers. It features wide range of topics in the world of computer programming. They also provide membership options along with job portal features for their customers.



**Data Source** - <https://relational.fit.cvut.cz/dataset/Stats>

We are acquiring datasets from **RELATIONAL DATASET REPOSITORY** (An anonymized dump of all user-contributed content on the Stats Stack Overflow network)



# Summary of Datasets

## Entities in Data Warehouse-

- Posts
- Comments
- Tags
- Users
- Posts History
- Post Links
- Votes
- Badges

## Method of access-

For extracting data from the internet open source we are using mysql connector in python environment and loading all the relevant data into the AWS S3 bucket for storage/Staging.

## Data Quality-

- There are some missing values in the datasets which will deal according to the requirement.
- Basic Exploratory analysis (EDA) is being performed in python environment.

# Data Profile

## Dataset 1 Summary

|                         |   |
|-------------------------|---|
| Source of information   | <a href="https://relational.fit.cvut.cz/dataset/Stats">https://relational.fit.cvut.cz/dataset/Stats</a> |
| Number of Table         | 8   |
| Number of Records       | 1,027,838 Rows<br>71 Columns  |
| Data type and structure | Numeric, String, Temporal   |
| Data Acquisiton Method  | MySQL connector   |

## Dataset 2 Summary

|                         |   |
|-------------------------|---|
| Source of information   | <a href="https://content.techgig.com/technology/24-highest-paying-programming-languages-for-developers/articleshow/76243434.cms">https://content.techgig.com/technology/24-highest-paying-programming-languages-for-developers/articleshow/76243434.cms</a> |
| Number of Table         | 1   |
| Number of Records       | 24 Rows   |
| Data type and structure | Numeric, String   |
| Data Acquisiton Method  | Web Scraping using python   |





# Timeline of the Project

## Waterfall Model

|   |   |                |                |                 |                 |                |                |                 |                 |                |               |                |  |
|---|---|----------------|----------------|-----------------|-----------------|----------------|----------------|-----------------|-----------------|----------------|---------------|----------------|--|
| Team members -                                | Surbhi Nayak, Yifeng, Yihang, Nossou, Qianwen |                |                |                 |                 |                |                |                 |                 |                |               |                |  |
|   | Week -1                                       | Week - 2       | Week - 3       | Week - 4        | Week - 5        | Week - 6       | Week - 7       | Week - 8        | Week - 9        | Week - 10      | Week - 12     | Week - 13      |  |
|   | Feb 21 - Feb 27                               | Feb 28 - Mar 6 | Mar 7 - Mar 13 | Mar 14 - Mar 20 | Mar 21 - Mar 27 | Mar 28 - Apr 3 | Apr 4 - Apr 10 | Apr 11 - Apr 17 | Apr 18 - Apr 24 | Apr 25 - May 1 | May 2 - May 8 | May 9 - May 12 |  |
| Bus Matrices                                  |   |                |                |                 |                 |                |                |                 |                 |                |               |                |  |
| Basic Use Cases                               |   |                |                |                 |                 |                |                |                 |                 |                |               |                |  |
| Dimensional Model with ETL Instructions       |   |                |                |                 |                 |                |                |                 |                 |                |               |                |  |
| Conceptual Model                              |   |                |                |                 |                 |                |                |                 |                 |                |               |                |  |
| Logical Model                                 |   |                |                |                 |                 |                |                |                 |                 |                |               |                |  |
| Physical Model                                |   |                |                |                 |                 |                |                |                 |                 |                |               |                |  |
| Data Dictionaries                             |   |                |                |                 |                 |                |                |                 |                 |                |               |                |  |
| Architecture Diagram with Network Information |   |                |                |                 |                 |                |                |                 |                 |                |               |                |  |
| Waterfall Plan                                |   |                |                |                 |                 |                |                |                 |                 |                |               |                |  |
| ETL/ELT Code                                  |   |                |                |                 |                 |                |                |                 |                 |                |               |                |  |
| Sample Queries                                |   |                |                |                 |                 |                |                |                 |                 |                |               |                |  |
| Visualization                                 |   |                |                |                 |                 |                |                |                 |                 |                |               |                |  |
| Presentation                                  |   |                |                |                 |                 |                |                |                 |                 |                |               |                |  |



# CONTENTS



1

Overview

2

About Data

3

Architecture

4

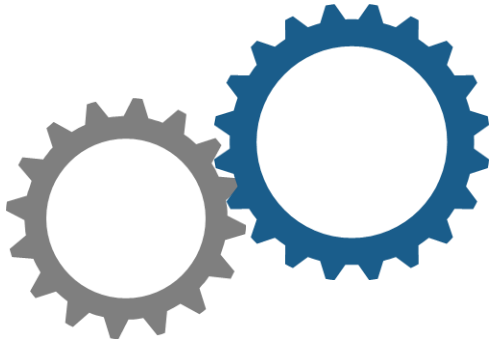
OLAP & OLTP

5

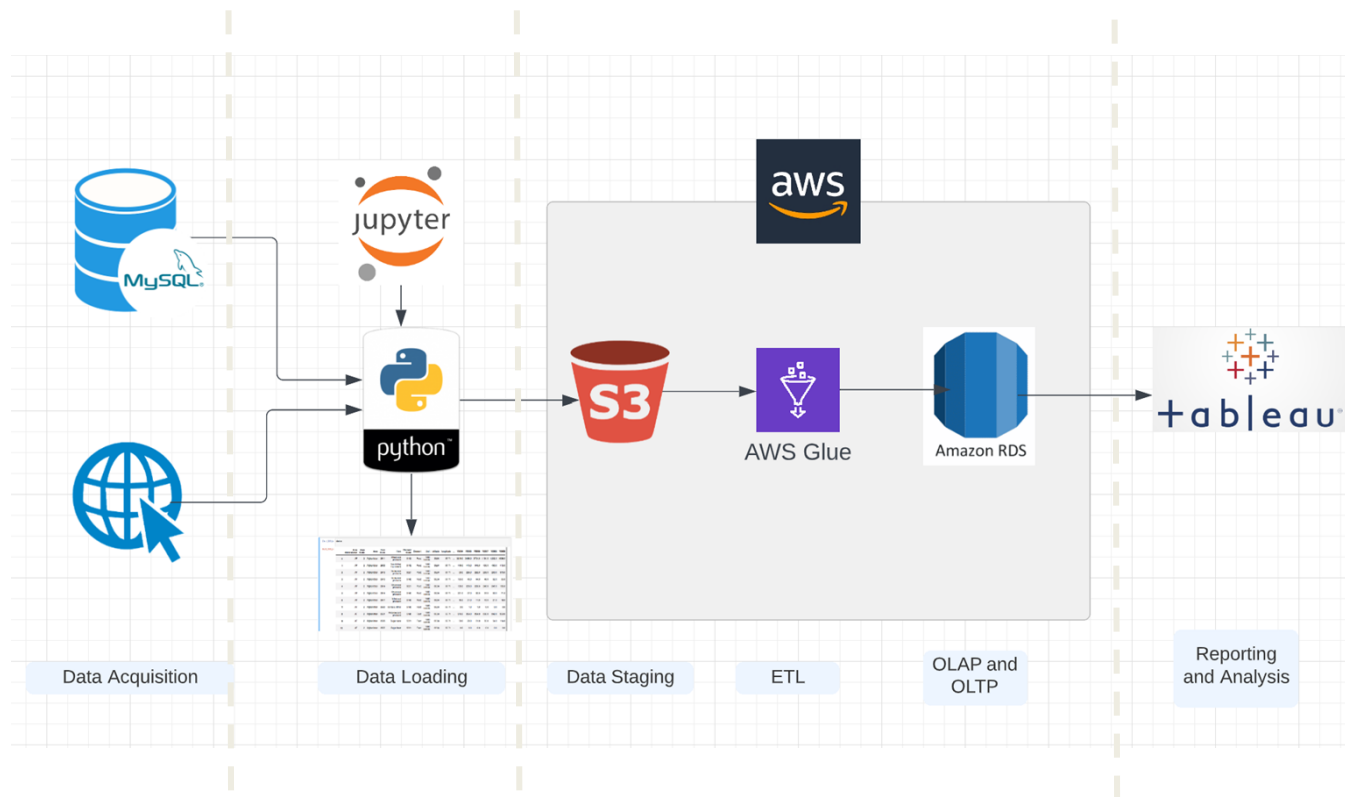
Reporting & Analysis

6

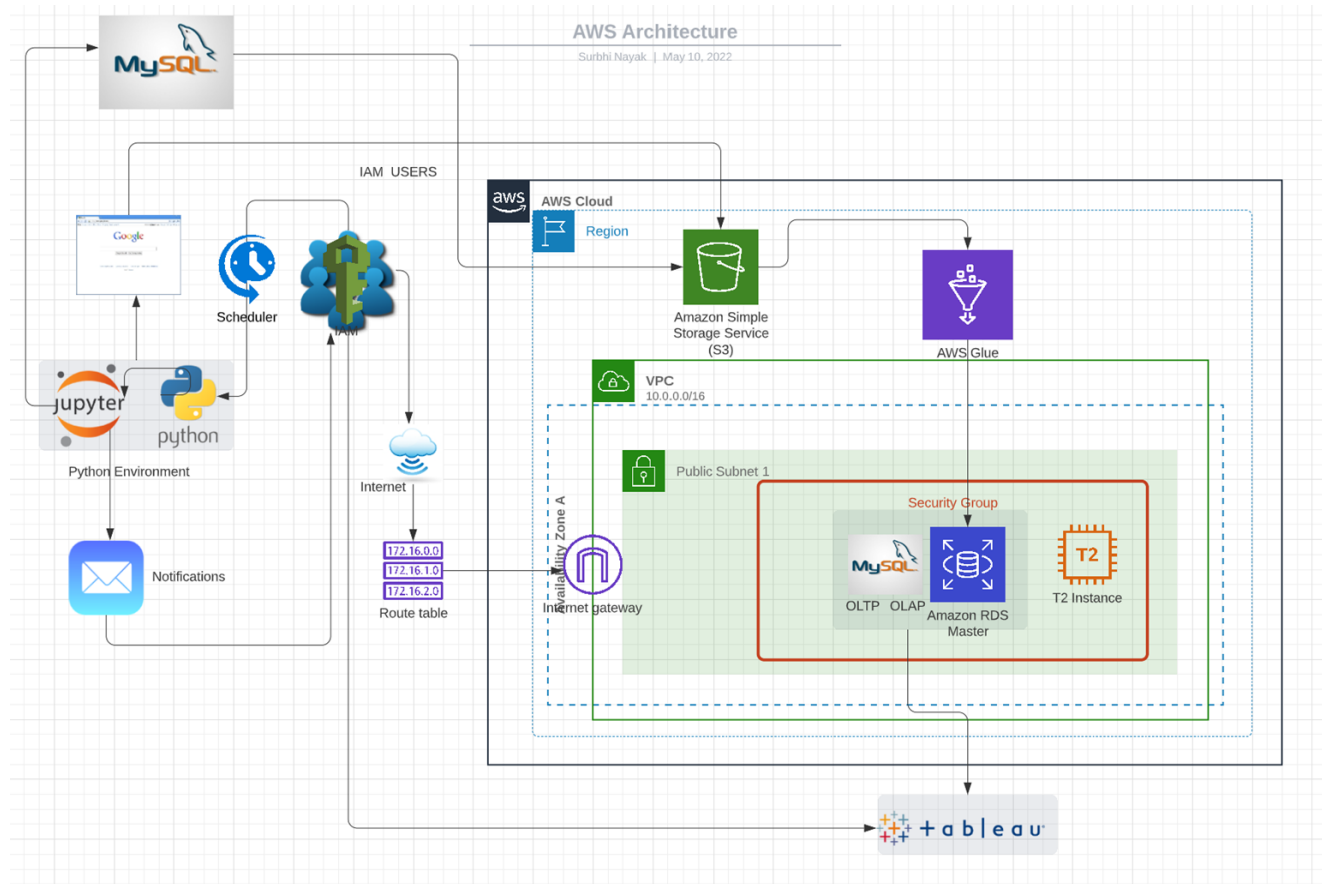
Challenges



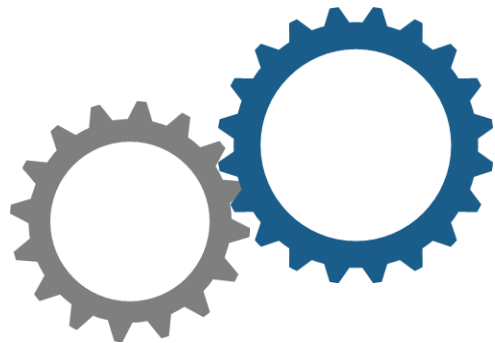
# Work Flow



# AWS Architecture



# CONTENTS



1

Overview

2

About Data

3

Architecture

4

OLAP & OLTP

5

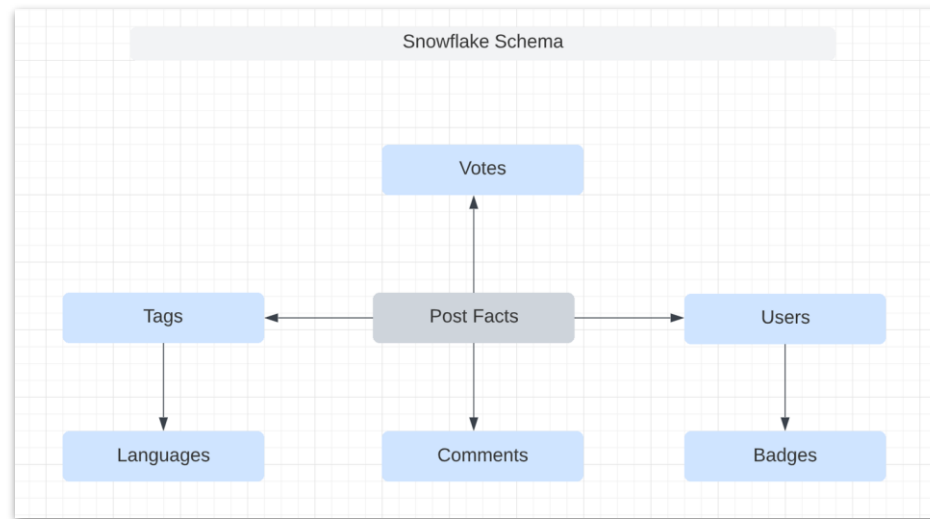
Reporting & Analysis

6

Challenges



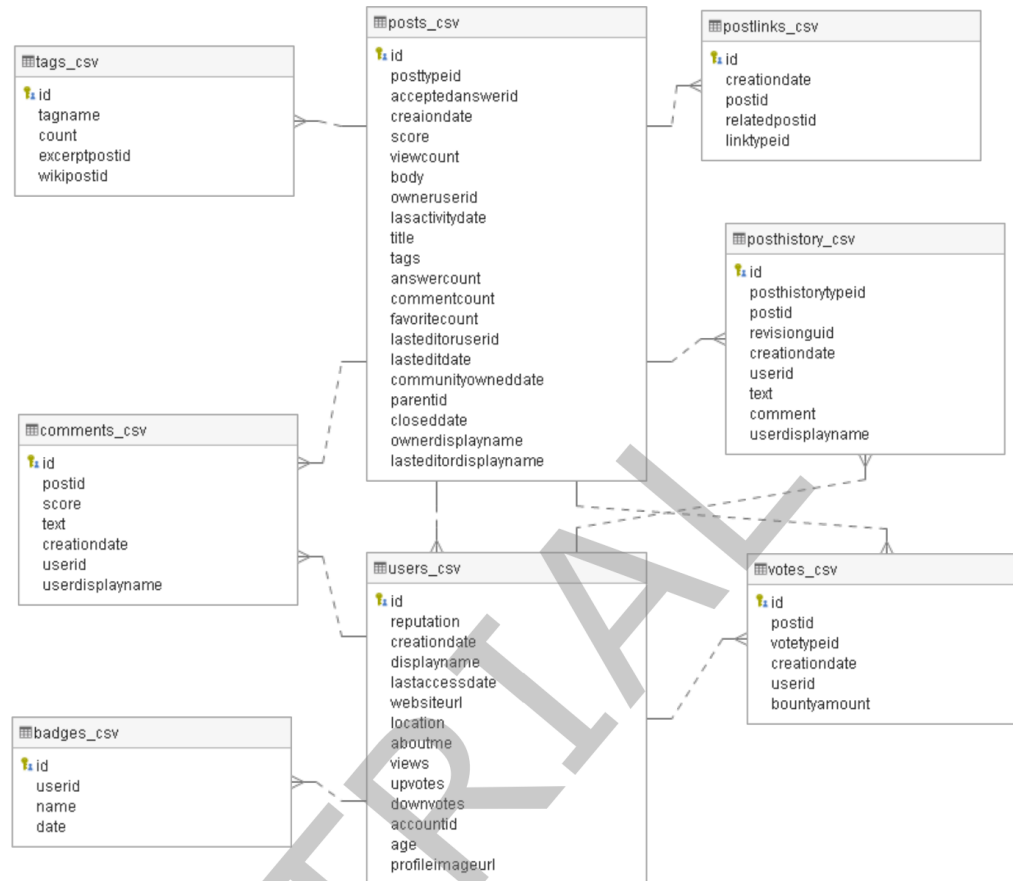
# Conceptual Model



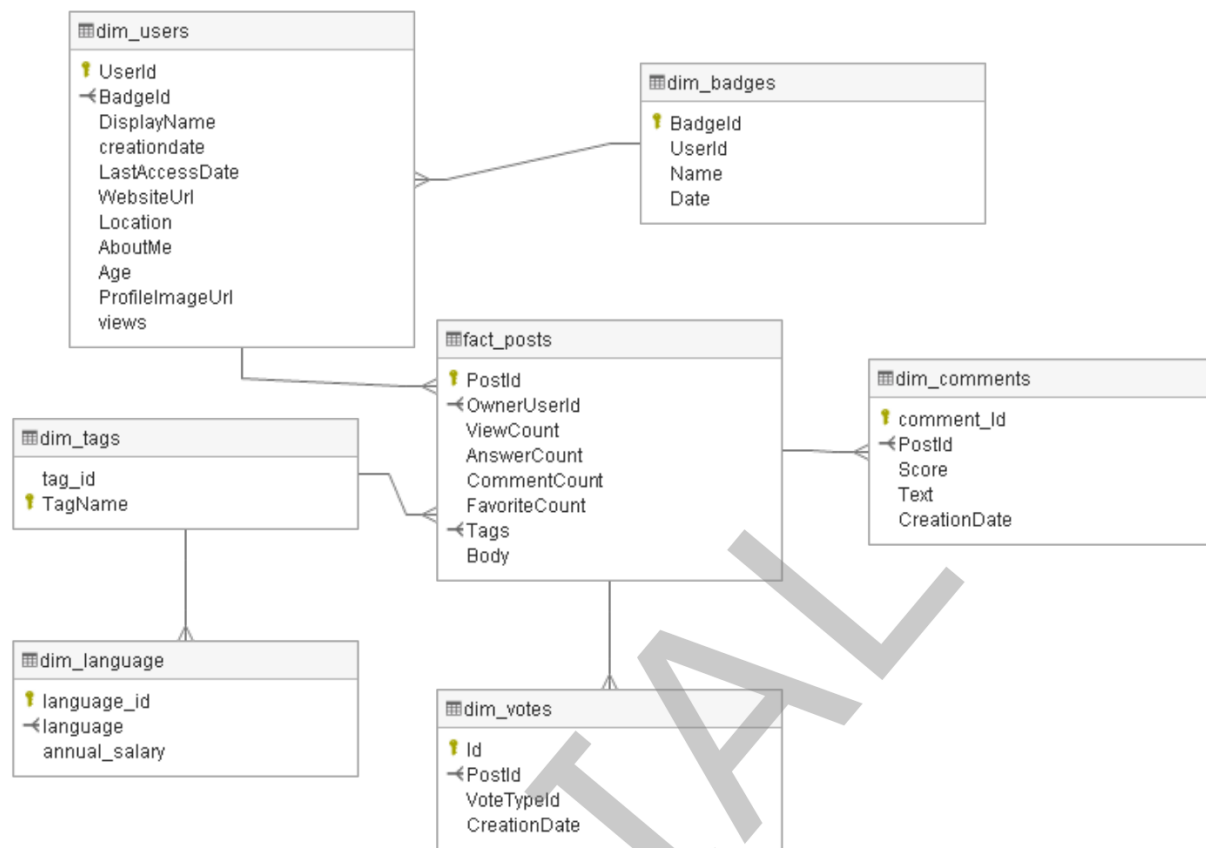
## Bus Matrix

|                            | Dimensions |      |       |       |      |
|----------------------------|------------|------|-------|-------|------|
| Facts (Business Processes) | date       | text | score | count | link |
| access                     |            |      |       |       |      |
| creation                   |            |      |       |       |      |
| answer                     |            |      |       |       |      |
| comment                    |            |      |       |       |      |
| edit                       |            |      |       |       |      |
| accept                     |            |      |       |       |      |
| vote                       |            |      |       |       |      |
| view                       |            |      |       |       |      |
| close                      |            |      |       |       |      |

# OLTP ERD

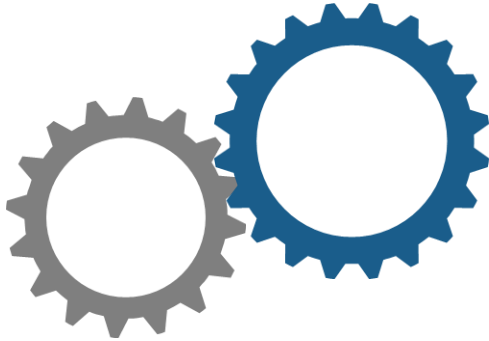


# OLAP ERD





# CONTENTS



1

Overview

2

About Data

3

Architecture

4

OLAP & OLTP

5

Reporting & Analysis

6

Challenges



# Connection in Tableau



☒ Live ☐ Extract

## Connections

Add

final.c7klak...mazonaws.com  
MySQL

## Database

Stack\_overflow\_OLAP

## Table



- dim\_badges
- dim\_comments
- dim\_language
- dim\_tags
- dim\_users
- dim\_votes
- fact\_posts

fact\_posts

dim\_comments

dim\_tags

dim\_users

dim\_votes

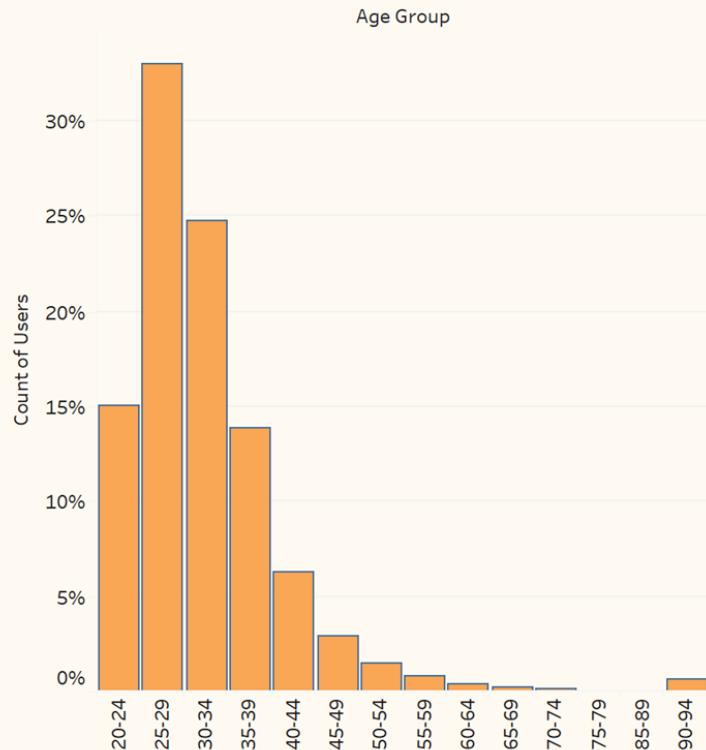
dim\_language

dim\_badges

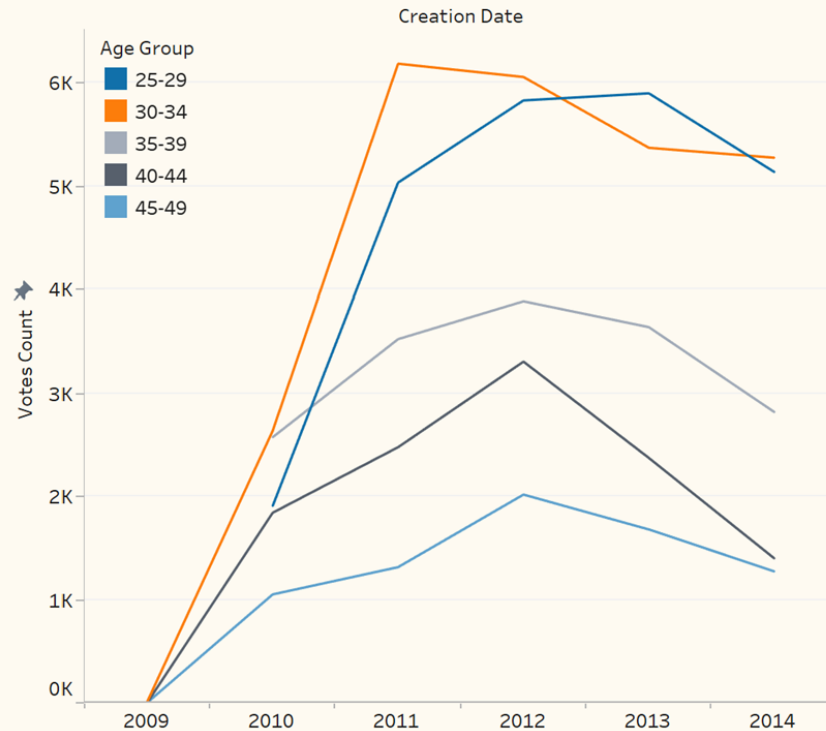


By understanding age groups and increase in user base participation over the years can help understand platform engagement rates and popularity.

## Age groups & count of users

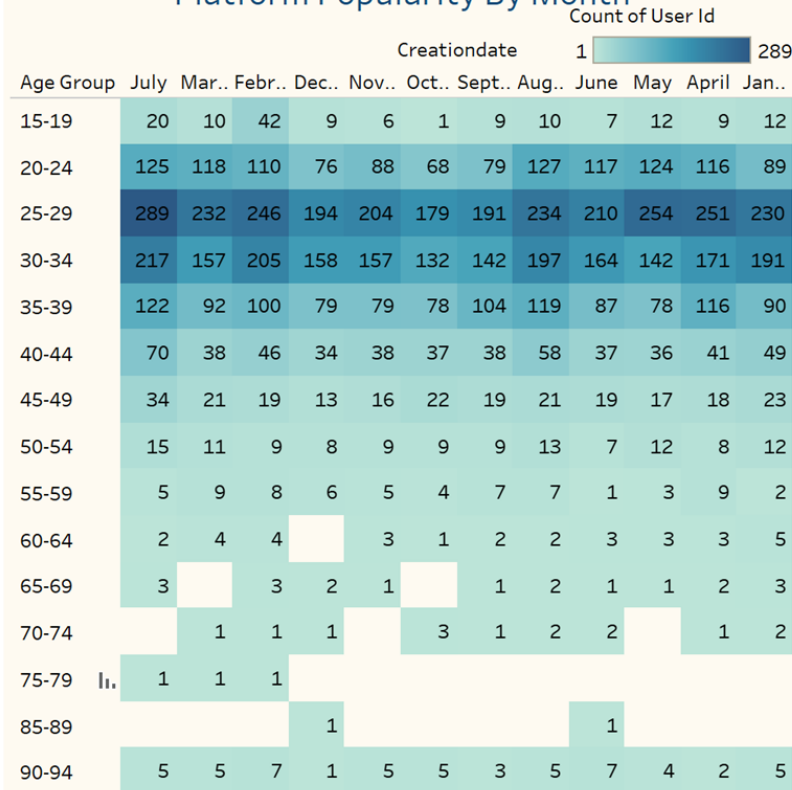


## Age groups participation

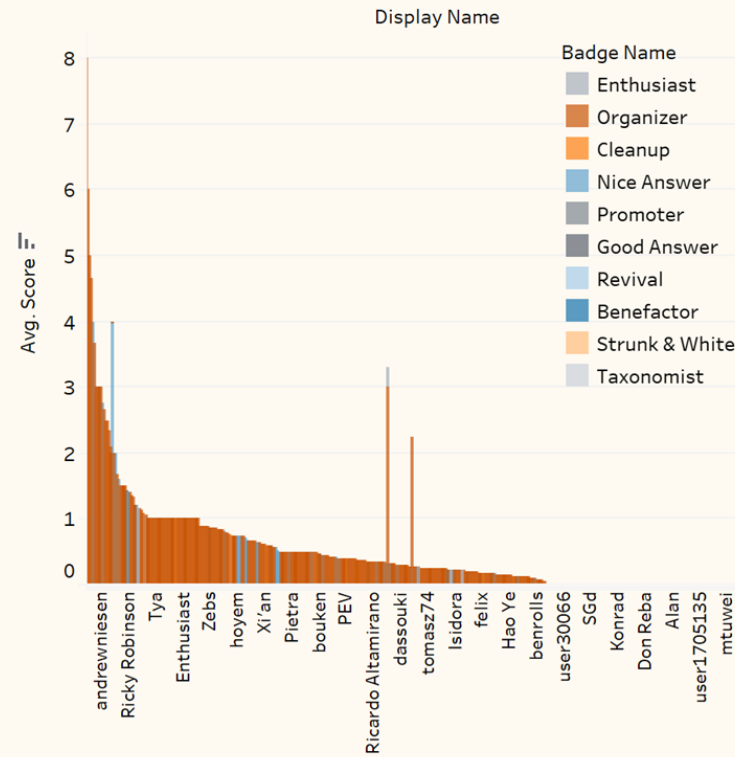


Month wise users count can help to understand user's participation pattern on the stack overflow platform. Also we can analyze badges popularity among top reputed users.

## Platform Popularity By Month

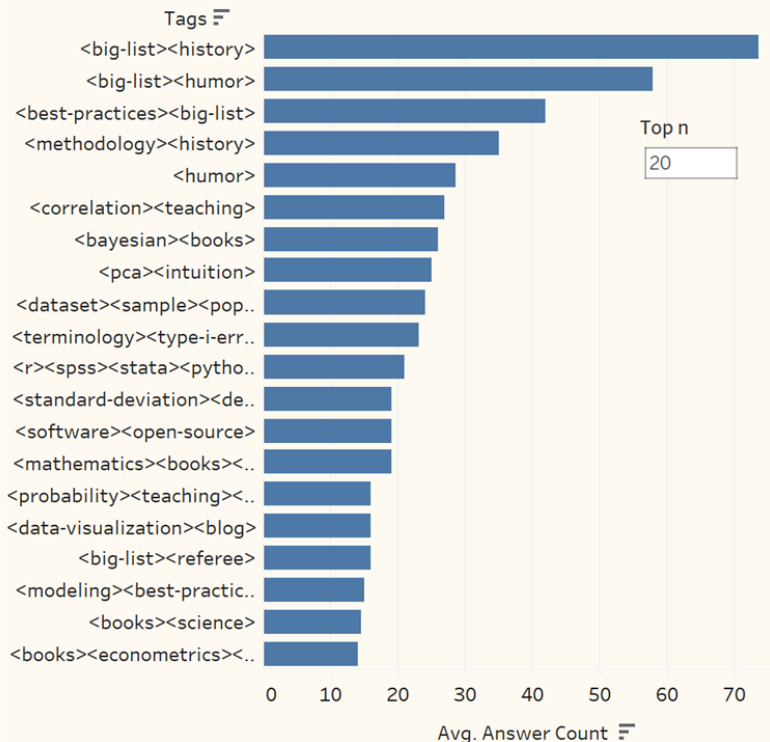


## Top reputed users

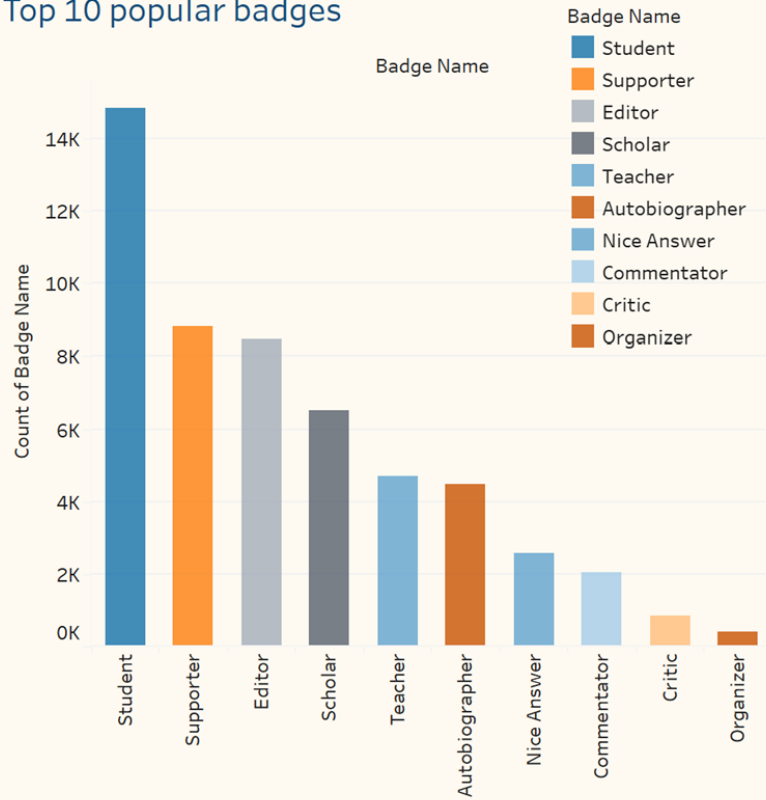


We can analyze various tags associated with the posts to help understand latest market trends/topics. Also, badges assigned to users can help further explain the engagement rate on platform.

## Top 20 Most Answered Tags

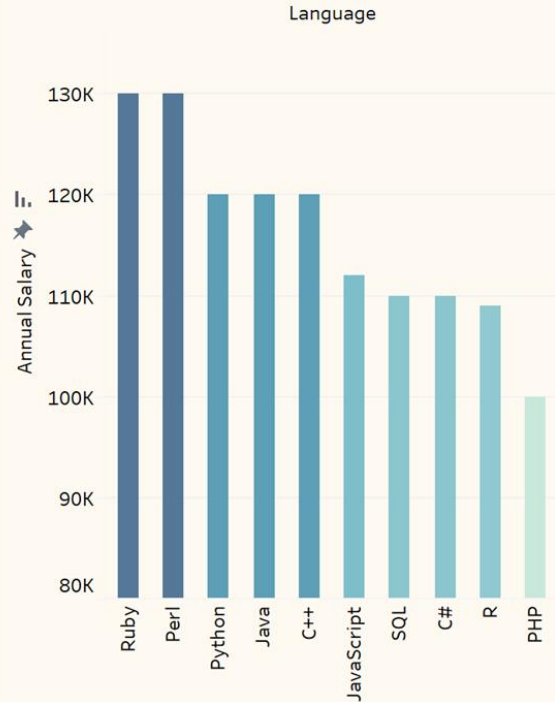


## Top 10 popular badges

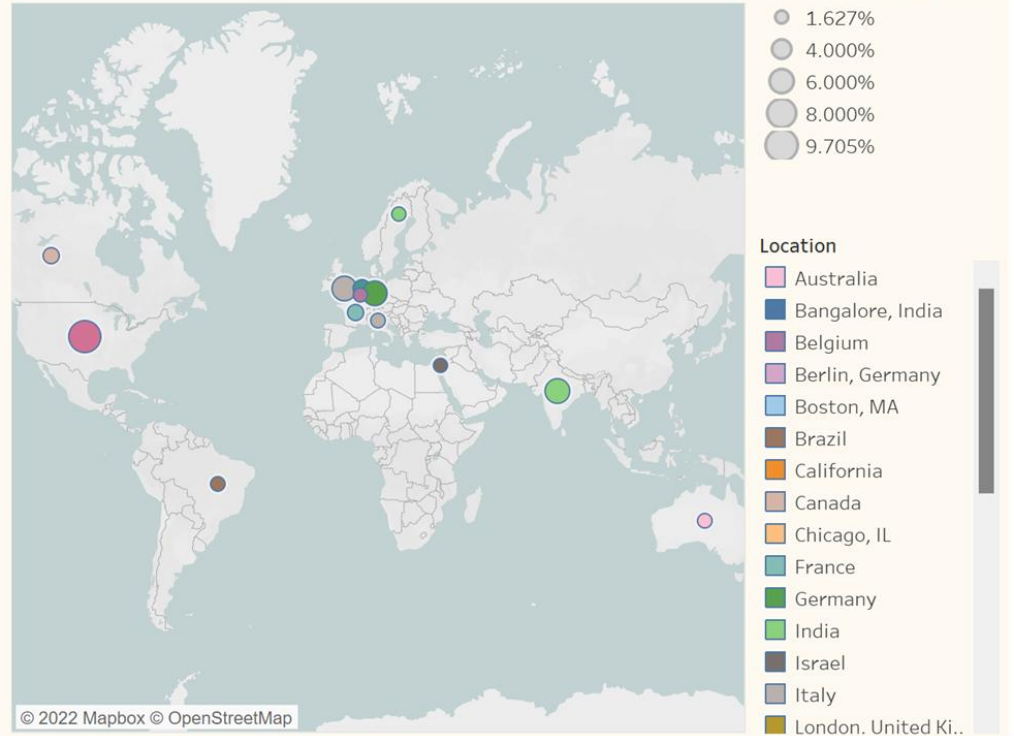


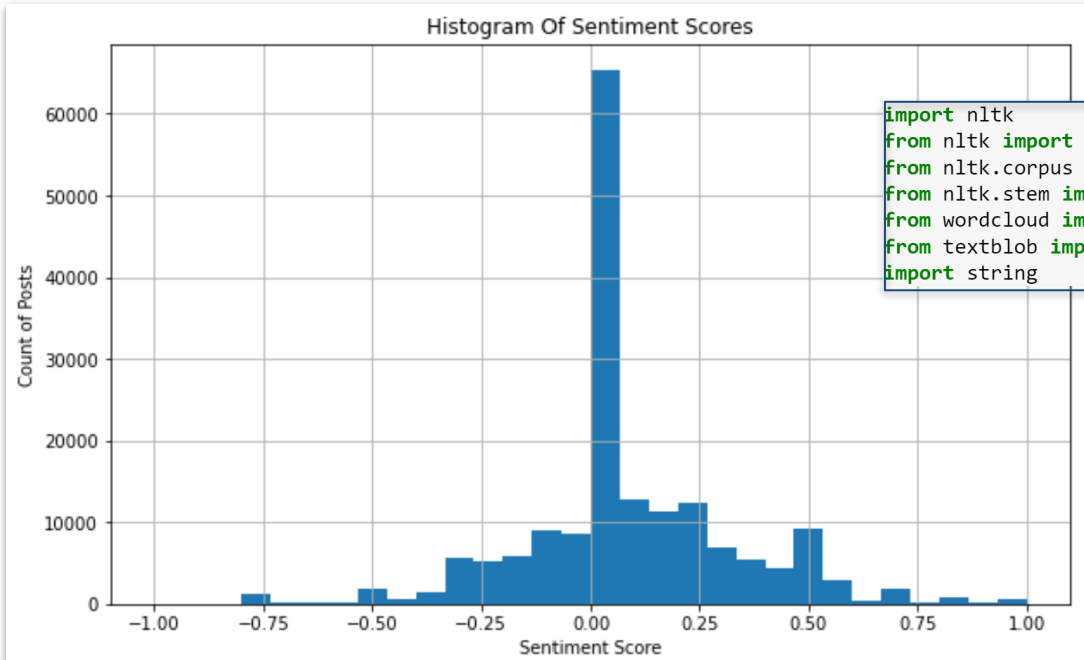


## Language & average annual salary



## Country & Users





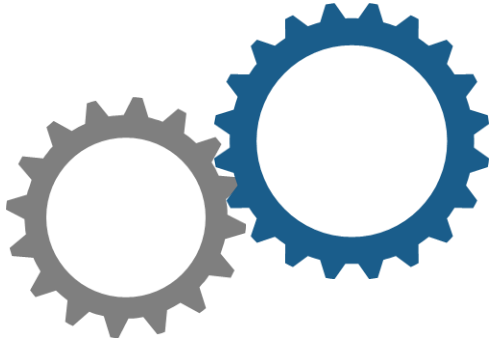
```
import nltk
from nltk import word_tokenize, sent_tokenize
from nltk.corpus import stopwords
from nltk.stem import LancasterStemmer, WordNetLemmatizer, PorterStemmer
from wordcloud import WordCloud, STOPWORDS
from textblob import TextBlob
import string
```

| Comment   | senti_score_polarity |
|---|----------------------|
| could poster child fo argumentative subjective least neec | 1                    |
| yes r nice 'valuable'                                     | 1                    |
| would convince boss use say excel                         | 1                    |
| mature well supported standard within certain scientific  | 1                    |

| Comment  | senti_score_polarity |
|--|----------------------|
| syntax would be: xtmixed studentscore lagstudentscore l  | -1                   |
| maybe create time series difference two variables        | -1                   |
| er "if *link* longer works" keyboard misbehaving droppi  | -1                   |
| use second model get get coefficient estimate effect sch | -1                   |



# CONTENTS



1

Overview

2

About Data

3

Architecture

4

OLAP & OLTP

5

Reporting & Analysis

6

Challenges





# Challenges

01

Creating timeline for project

02

Task allocation among the team

03

Finding the data sources and acquisition methods



Thank You!

Any Questions ?

