# 01

# Project Background and Definition

# 1. Background and Definition

- Pertussis, or Whooping cough, is a highly contagious lung infection
- Two vaccines: whole-cellular (wP) and acellular (aP)
- Challenges of balancing vaccine safety and efficacy
- Importance of ongoing research and monitoring of vaccine effectiveness over time
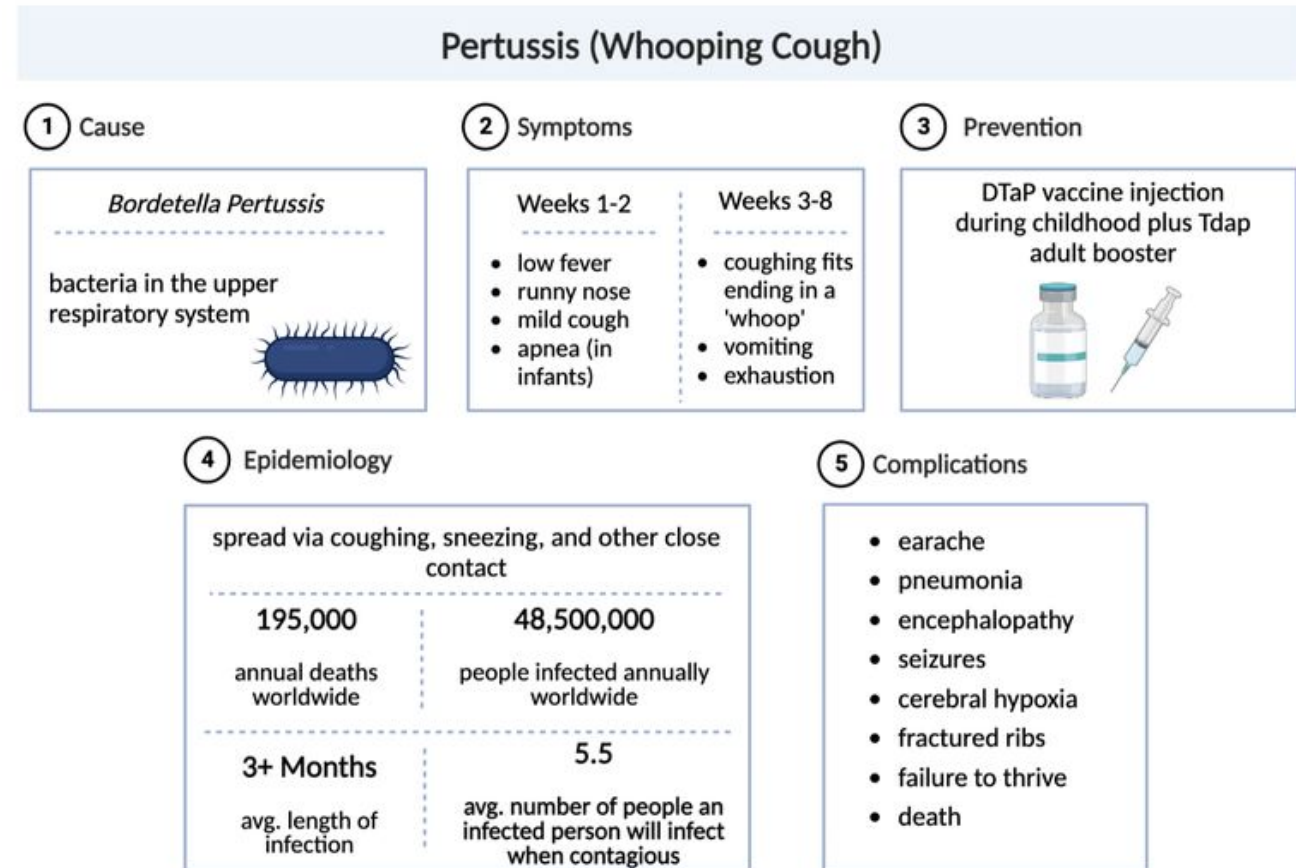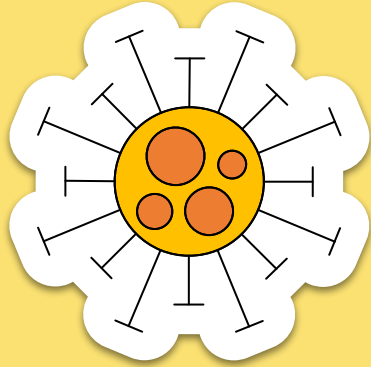- How does vaccine induced immunity change over time per person?



## Pertussis (Whooping Cough)

**① Cause**

*Bordetella Pertussis*

bacteria in the upper respiratory system

**② Symptoms**

Weeks 1-2
- low fever
- runny nose
- mild cough
- apnea (in infants)

Weeks 3-8
- coughing fits ending in a 'whoop'
- vomiting
- exhaustion

**③ Prevention**

DTaP vaccine injection during childhood plus Tdap adult booster

**④ Epidemiology**

spread via coughing, sneezing, and other close contact

| 195,000 | 48,500,000 |
|---------|-----------|
| annual deaths worldwide | people infected annually worldwide |
| 3+ Months | 5.5 |
| avg. length of infection | avg. number of people an infected person will infect when contagious |

**⑤ Complications**

- earache
- pneumonia
- encephalopathy
- seizures
- cerebral hypoxia
- fractured ribs
- failure to thrive
- death

Sources: Centers for Disease Control, World Health Organization, PLOS Medicine, PubMed

# 1. Background and Definition

## Goals

- Help save lives from preventable pertussis cases

- Share academic research with growing community

- Advance the understanding of immunology and use models in the real world to improve vaccine effectiveness
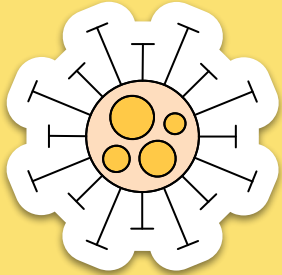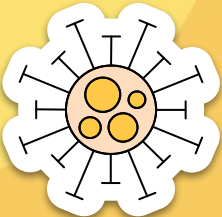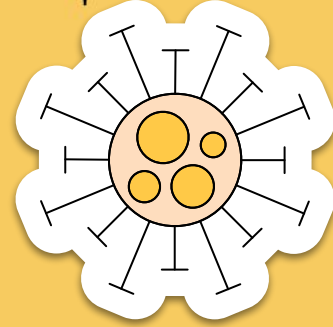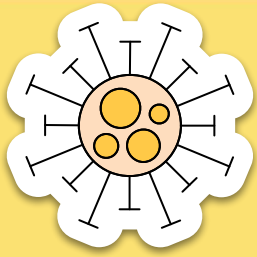
UC San Diego

**JACOBS SCHOOL OF ENGINEERING**

**MAS Data Science and Engineering**

CMI-PB
COMPUTATIONAL MODELS OF IMMUNITY
— PERTUSSIS BOOST —

# 02
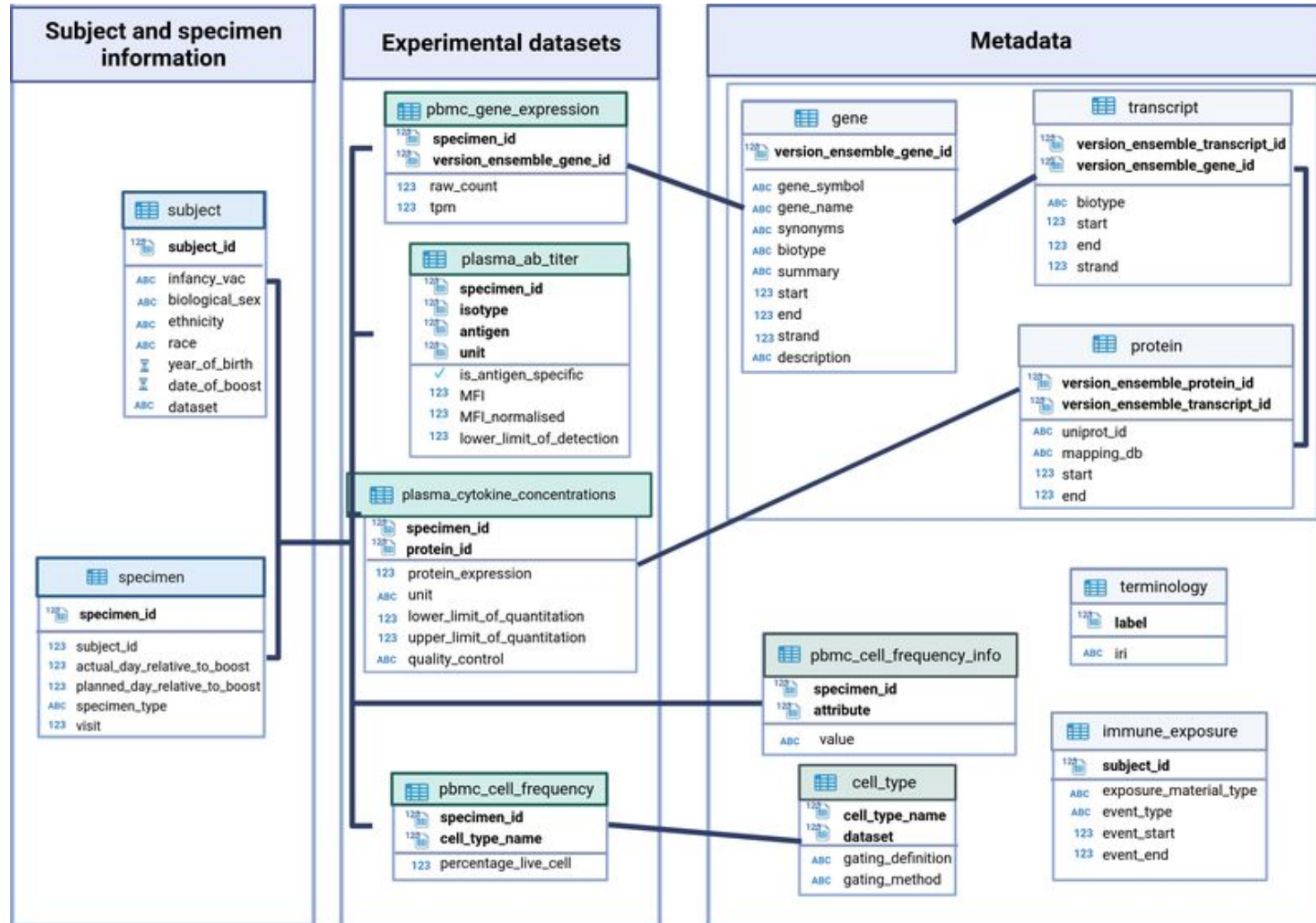# Datasets Overview

# 2. Overview of Data Sets

**Objective:** Analyze immune responses post-Tdap booster vaccination.

**Subjects:** 118 individuals contributing 500+ blood specimens.
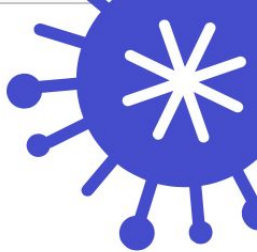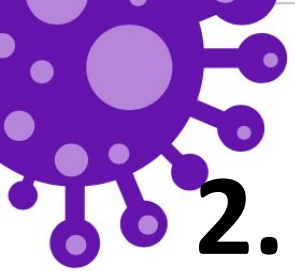
**Timeframe:** Pre- and post-vaccination (days 1, 3, 7, and 14).
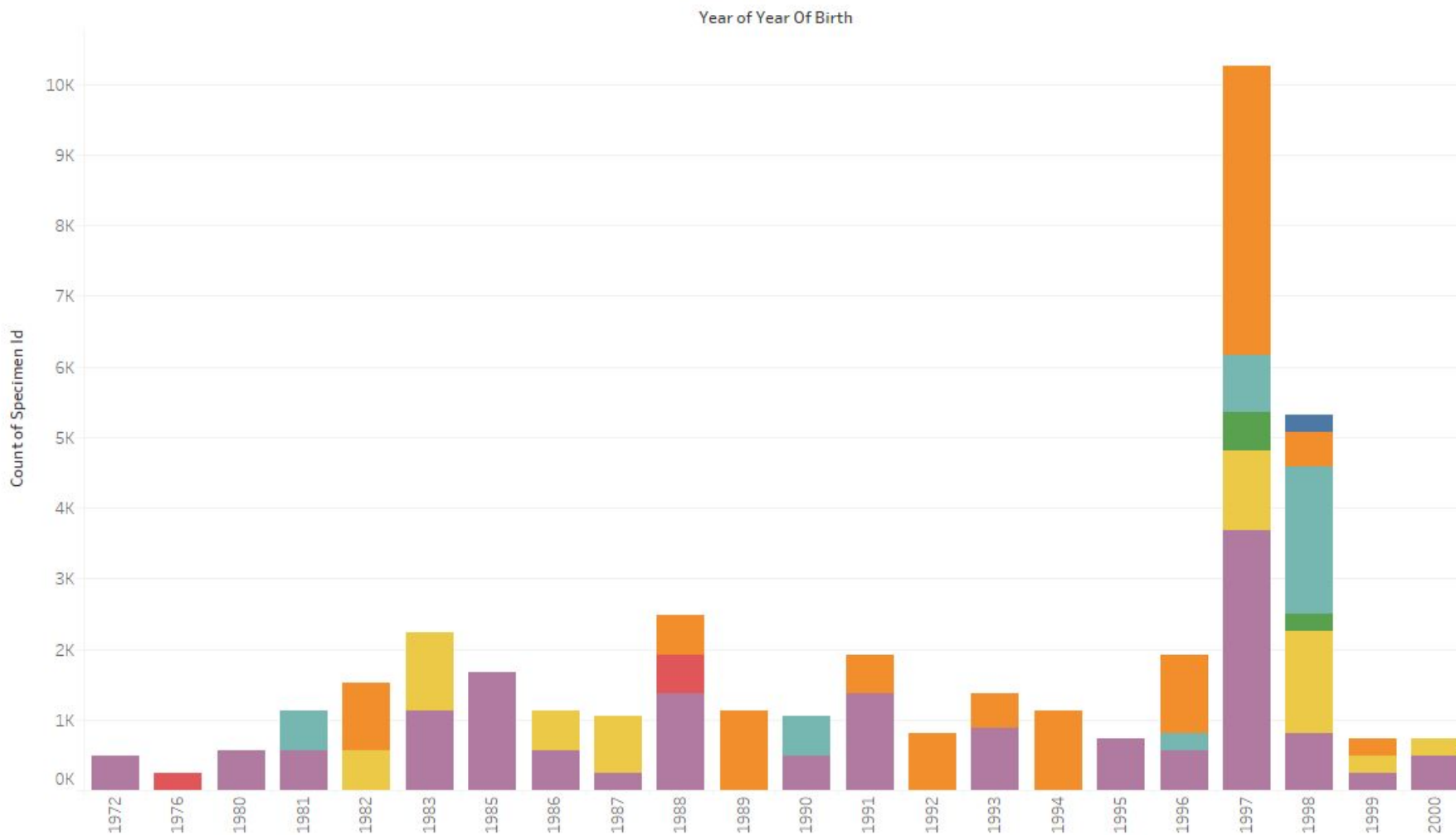
# 2. Overview of Data Sets



Source: CMI-PB

# 2. Overview of Data Sets

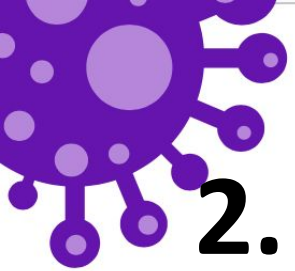# 2. Overview of Data Sets

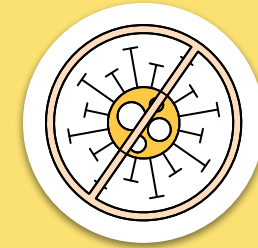# 2. Overview of Data Sets



Frequency of Monocytes

# 2. Overview of Data Sets

DNA → mRNA(Gene Expression) → Protein Expression → CCL3

# 03
# Project Objectives

# 3. Project Objectives

- **Building Computational Models:** To predict vaccination outcomes for newly tested individuals.

- **Participation in the Prediction Challenge:** Showcase intuition and analytical skills by engaging in the community prediction challenge
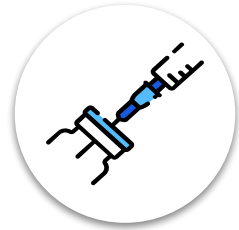
| | Annual prediction challenge title | Contestants | Number of subjects | | Current status |
|---|---|---|---|---|---|
| | | | Training dataset | Test dataset | |
| 1 | **First Challenge:** Internal dry run | CMI-PB consortium | 60 (28 aP + 32 wP) | 36 (19 aP + 17 wP) | Concluded in May 2022 |
| 2 | **Second Challenge:** Invited challenge | Invited contestants | 96 (47 aP + 49 wP) | 22 (13 aP + 9 wP) | Will be announced in September 2023 |
| 3 | **Third Challenge:** Open Challenge 1 | Public | 118 (60 aP + 58 wP) | 32 (16 aP + 16 wP) | Will be announced in April 2024 |
| 4 | **Fourth Challenge:** Open Challenge 2 | Public | 150 (76 aP + 74 wP) | 32 (16 aP + 16 wP)* | Will be announced in December 2024 |

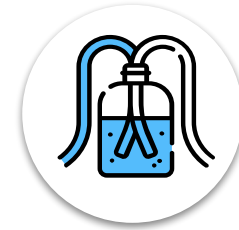*Goal

# 3. Project Objectives

- Accurate prediction of immune responses
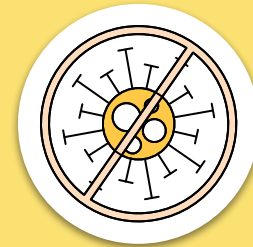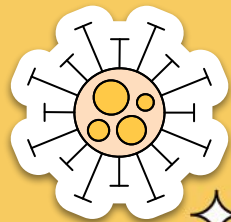
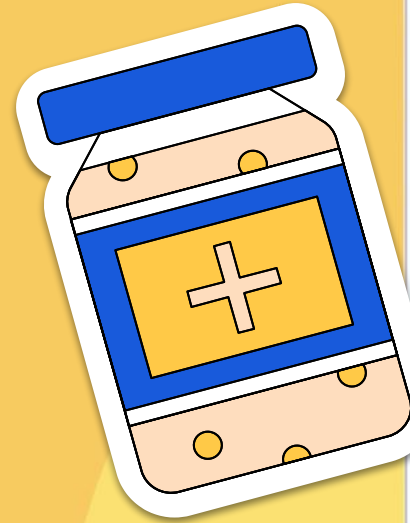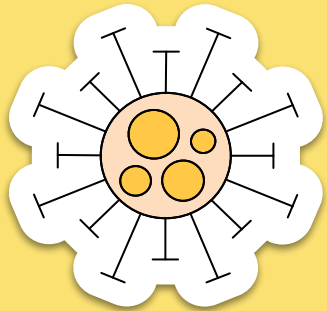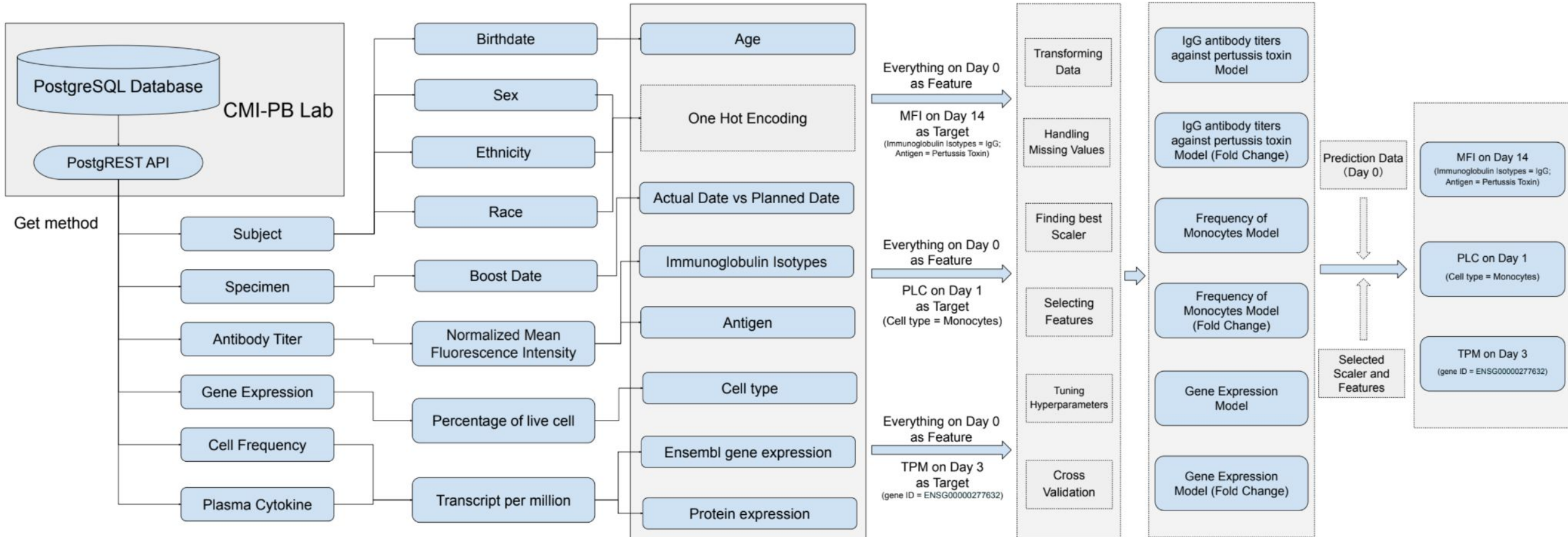- Identifying what variables induce a strong response

- Real-time access to dynamic immune response data

- Efficient handling of diverse features

# 04

# Methods
# and Techniques

# 4. Methods and Techniques
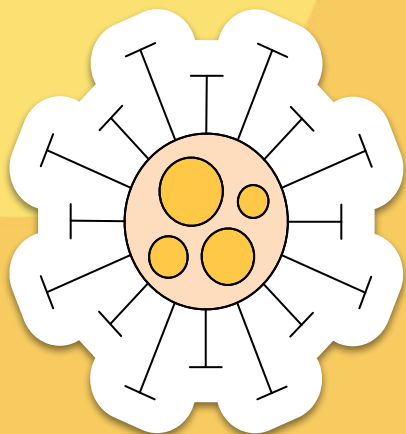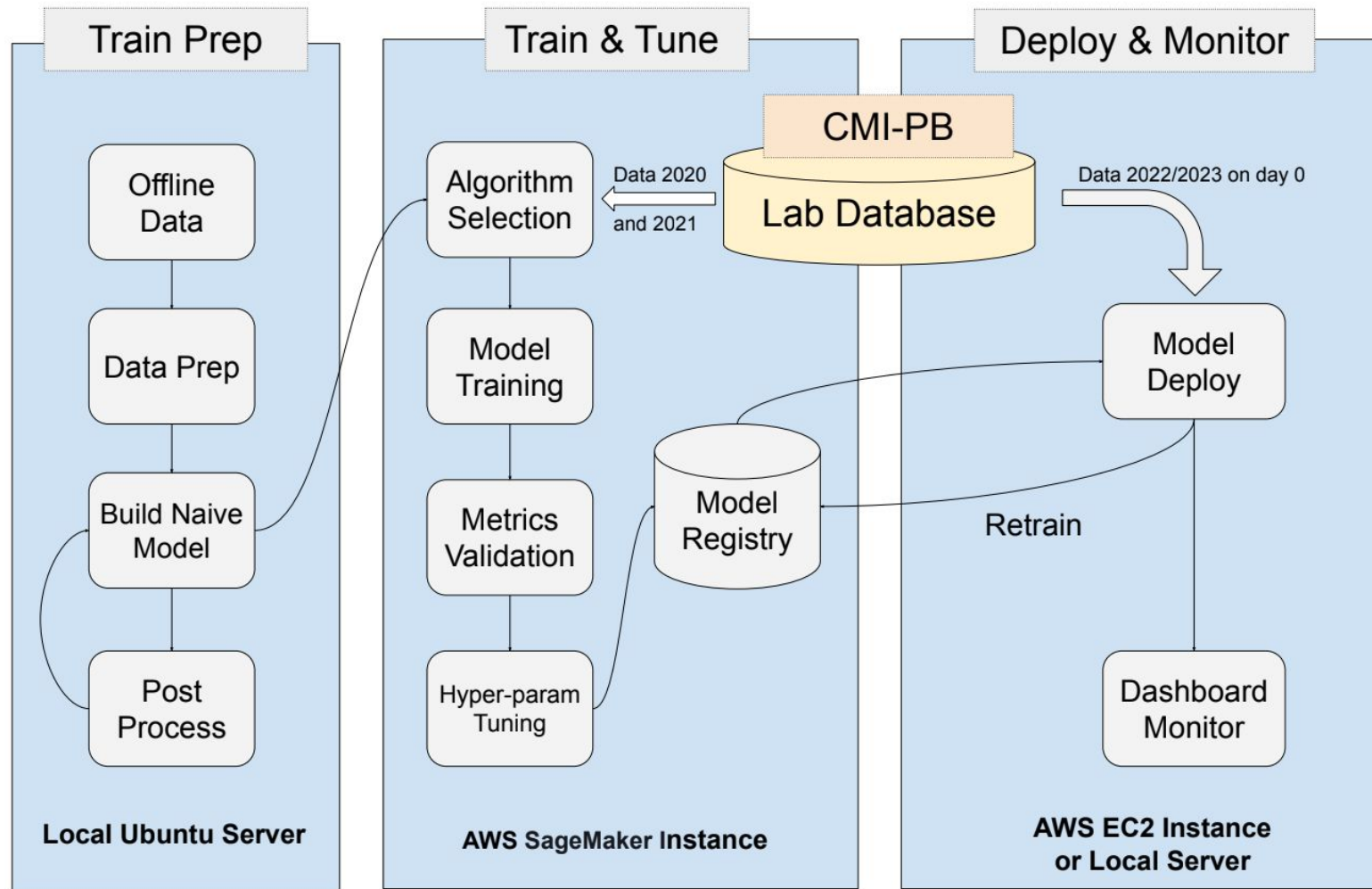
# 05

# Solution Architecture

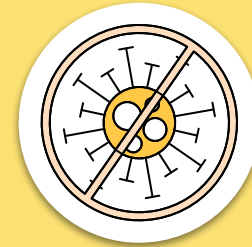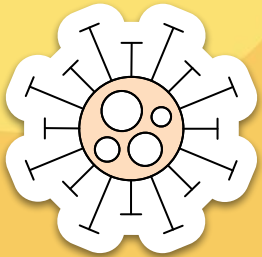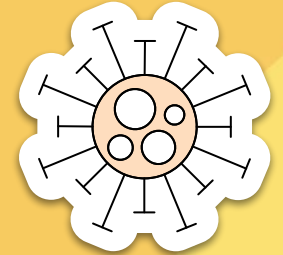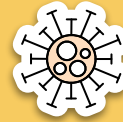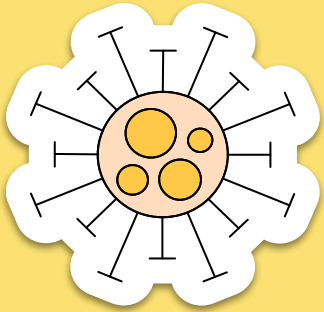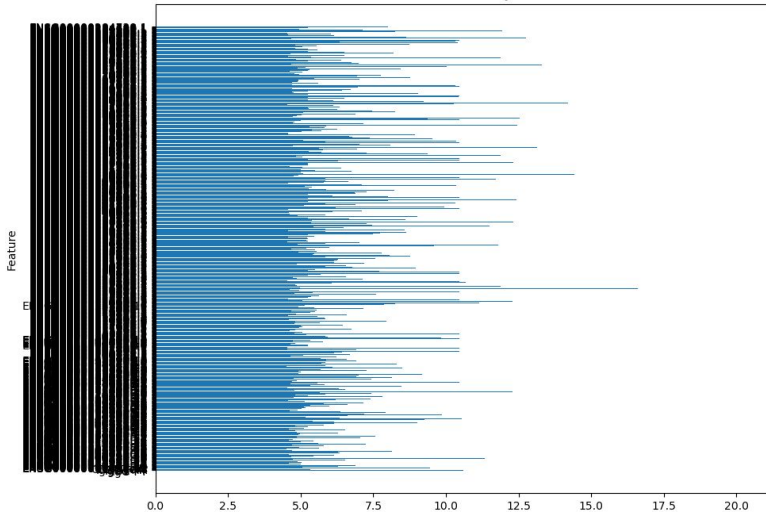# 5. Solution Architecture
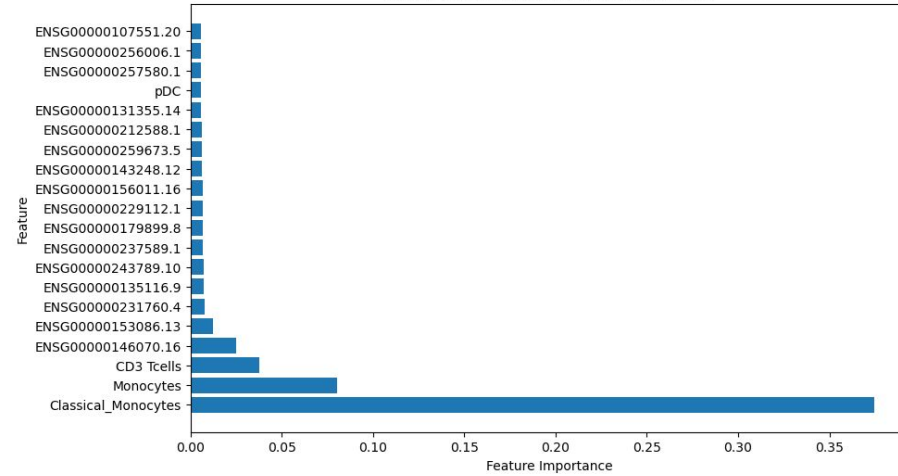
# 06
# Analysis of Results
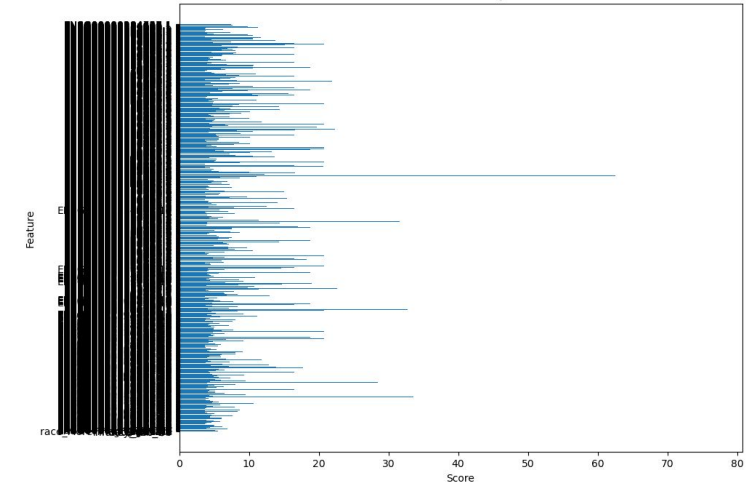
# 6. Analysis of Results
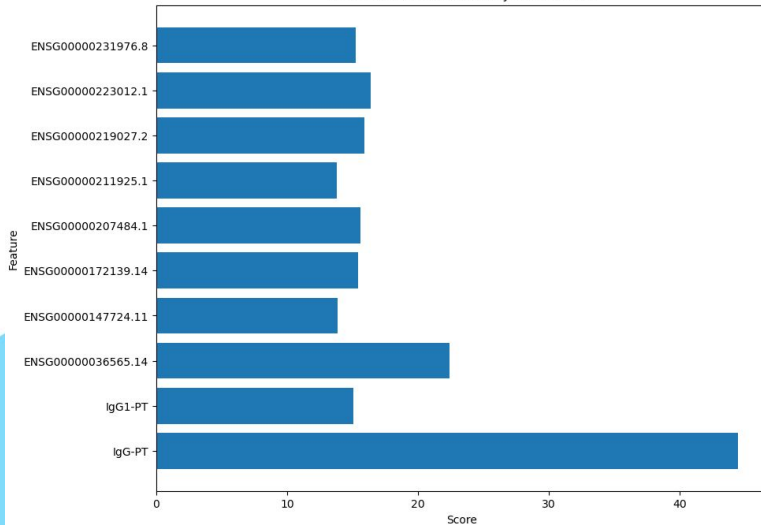


task11 Feature Selected by SelectKBest

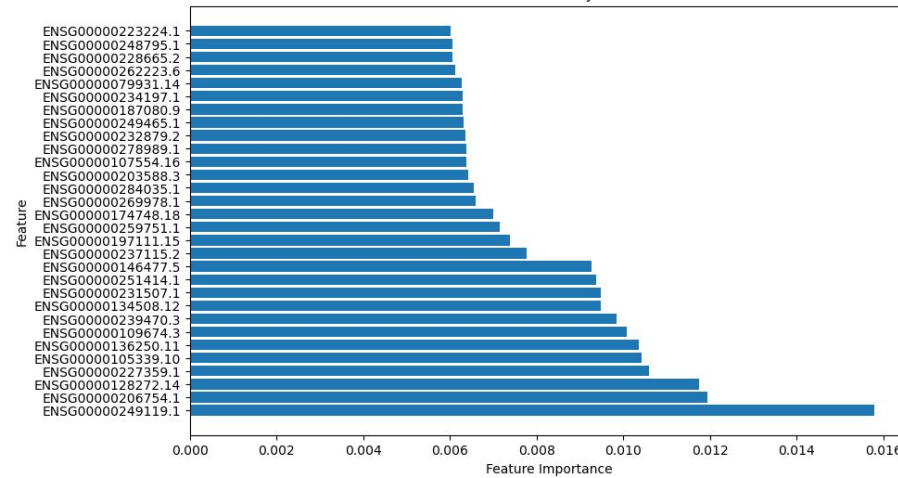task21 Feature Selected by Random Forest
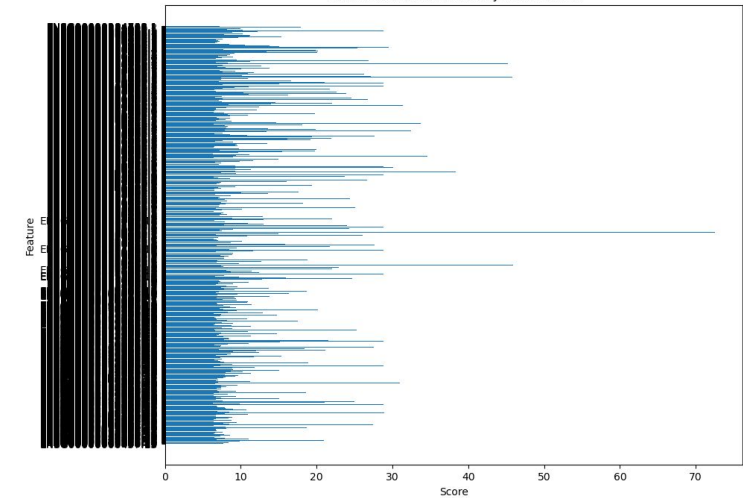
task31 Feature Selected by SelectKBest

task12 Feature Selected by SelectKBest
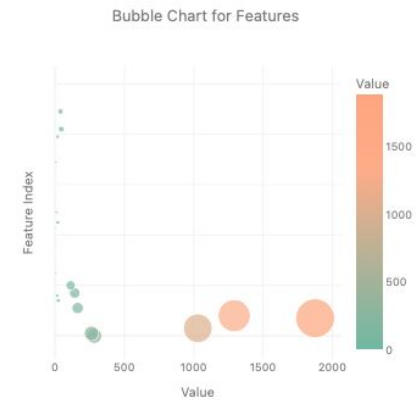
task22 Feature Selected by Random Forest

task32 Feature Selected by SelectKBest
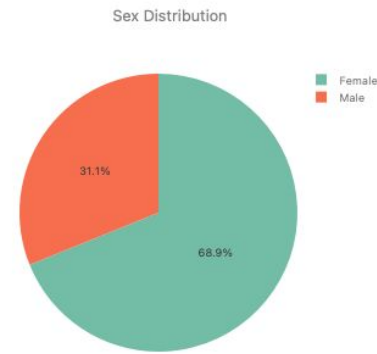
# 6. Dashboard

# 6. Dashboard - Demo

# 6. Analysis of Results
– Predict IgG antibody titers against pertussis toxin on day 14



R-squared Model Evaluation for 1.1



R-squared Model Evaluation for 1.2

Take Log2 Fold Change
the Target

$$\log_2 \left( \frac{MFI \ on \ day \ 14}{MFI \ on \ day \ 0} - 1 \right)$$

# 6. Analysis of Results
– Predict frequency of Monocytes on day 1

R-squared Model Evaluation for 2.1

R-squared Model Evaluation for 2.2

Take Log2 Fold Change the Target

$$\log_2 \left( \frac{PLC \text{ on } day\ 1}{PLC \text{ on } day\ 0} - 1 \right)$$

# 6. Analysis of Results

– Predict gene expression of CCL3 on day 3



R-squared Model Evaluation for 3.1

R-squared Model Evaluation for 3.2

Take Log2 Fold Change the Target

$$\log_2 \left( \frac{TPM \ on \ day \ 3}{TPM \ on \ day \ 0} - 1 \right)$$

# 07

# Insights

# 7. Insights

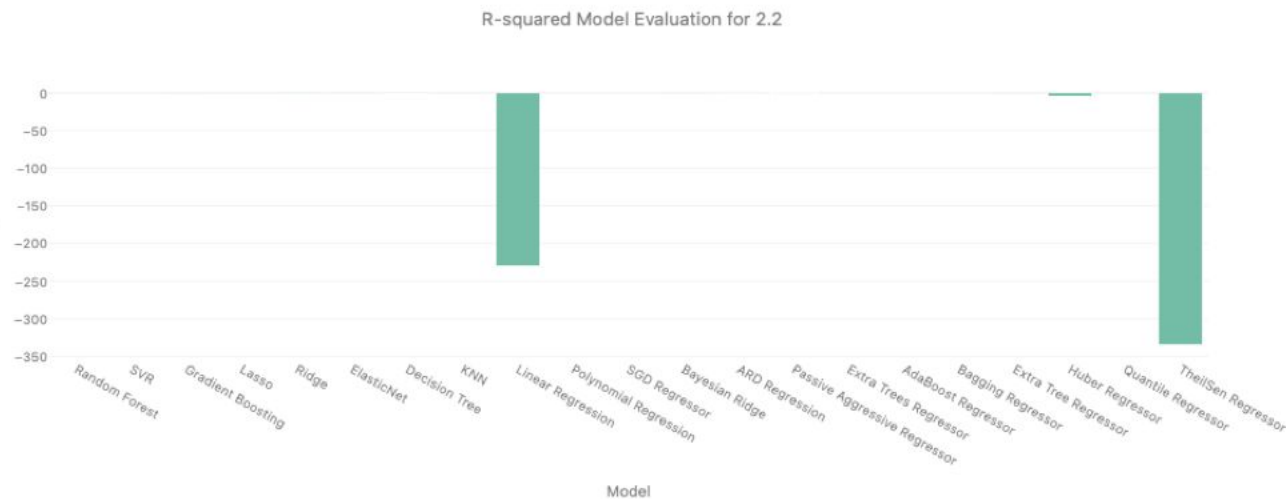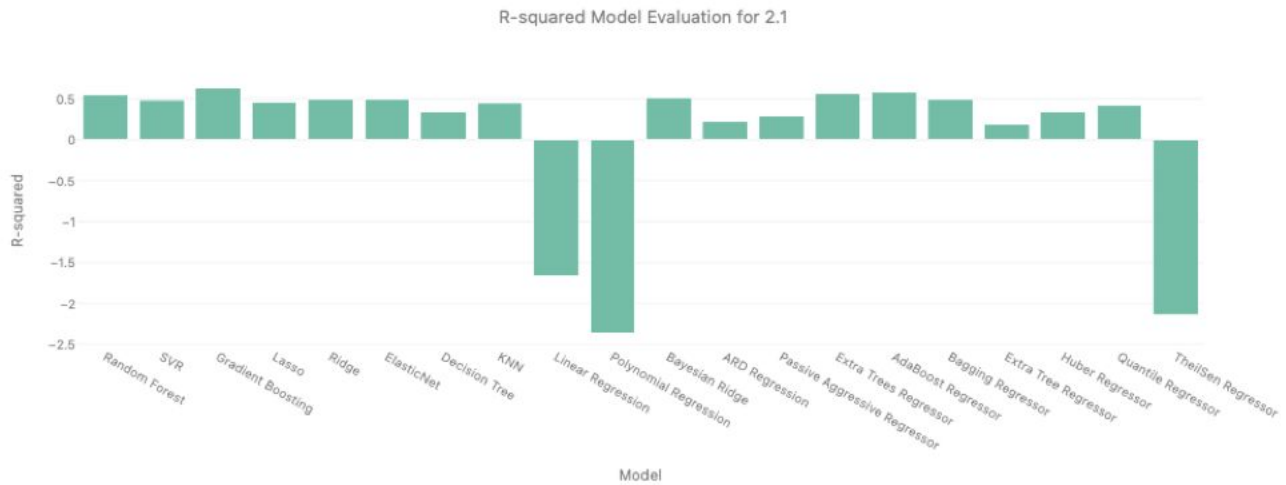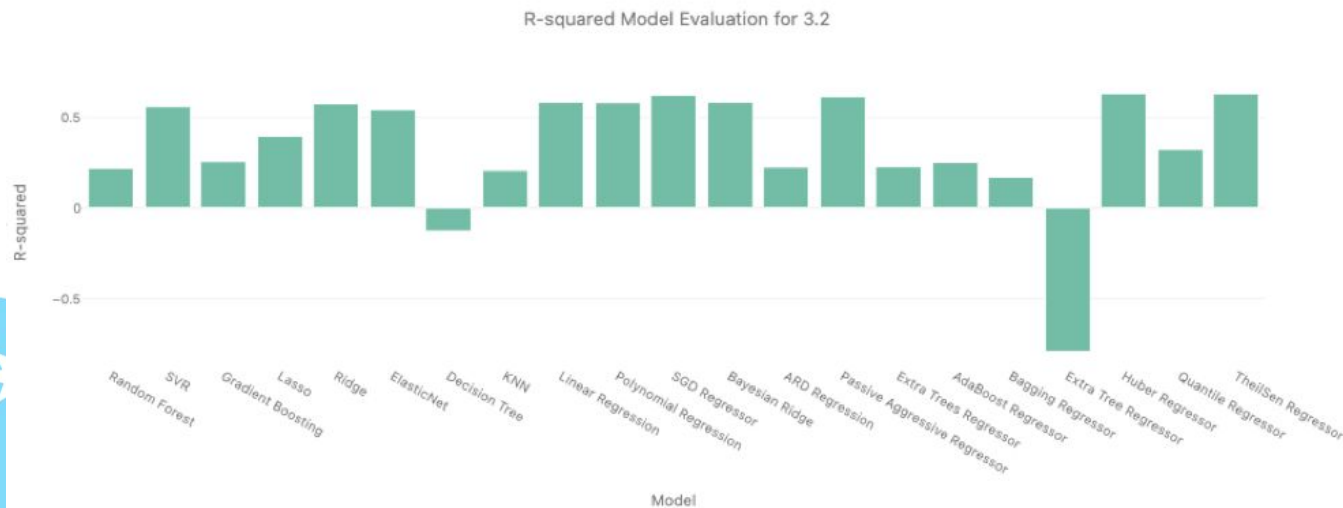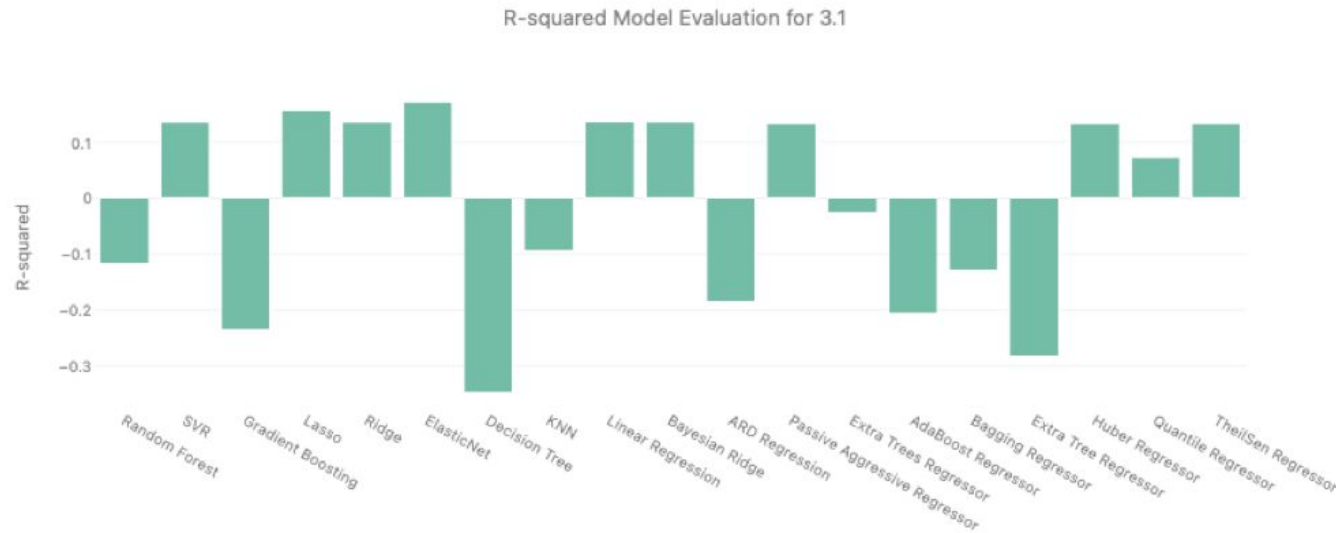- *Support Vector Regression (SVR)*: The model performed well in predicting antibody titer levels. SVR is effective in high-dimensional spaces and can handle nonlinear relationships, making it suitable for handling our complex biological data.

- *Extra Tree Regressor and Gradient Boosting*: These ensemble methods outperformed other methods in predicting the Log2 fold change of antibody levels at day 14 relative to day 0. They were able to reduce overfitting through bagging and boosting techniques, making them a good fit for our dataset.

# 7. Insights

- *Gradient Boosting for Monocyte Frequency*: Gradient Boosting showed the best performance in predicting monocyte frequency at day 1. The approach of this model helps minimize the prediction error, making it applicable to a variety of tasks.

- *ElasticNet for Gene Expression Levels*: Initially, ElasticNet performed best in predicting CCL3 gene expression levels at day 3. However, simpler models such as Stochastic Gradient Descent(SGD) Regressor and TheilSen Regressor outperformed ElasticNet when predicting the Log2 fold change of the target variable. This shift highlights the importance of model simplicity and regularization in dealing with large amounts of gene expression data.

08

Next Steps Plan

# 8. Next Steps Plan

## Model performance metrics:

- *Regular re-evaluation*: Continuously monitor performance metrics (MSE, MAE, $R^2$) to ensure the model remains accurate as new data becomes available.

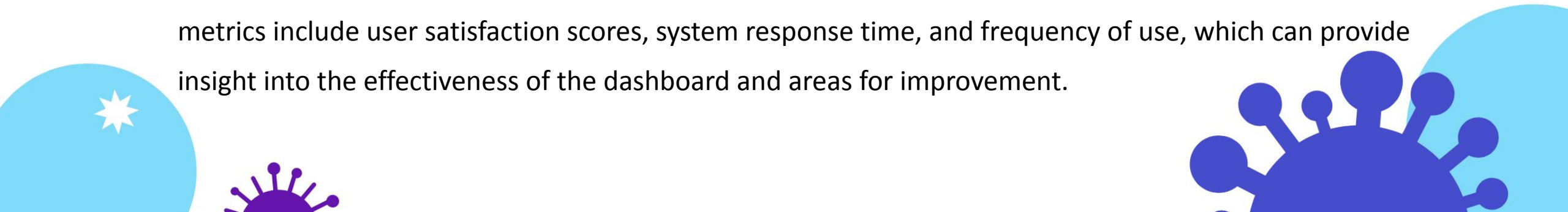- *Model retraining*: Regularly retrain the model with updated datasets to maintain prediction accuracy and adapt to any changes in data patterns.

## Operational metrics:

- *Data quality:* This involves regular data cleansing and validation processes. High-quality data is critical to maintaining the integrity of the model and ensuring reliable predictions.

- *System availability*: Monitor user feedback on dashboard usability for iterative improvements. Usability metrics include user satisfaction scores, system response time, and frequency of use, which can provide insight into the effectiveness of the dashboard and areas for improvement.

# Thank you!