# 2025-09-08 Meeting Agenda

## 🗓️ Date

08 Sep 2025

## 👥 Participants

- RAIL PG-2 project team

  o   Tao Xu                a1937511
  o   Sheng Wang            a1903948
  o   Jinchao Yuan          a1936476
  o   Zilun Ma              a1915860
  o   Di Zhu                a1919727
  o   Xin Wei               a1912958
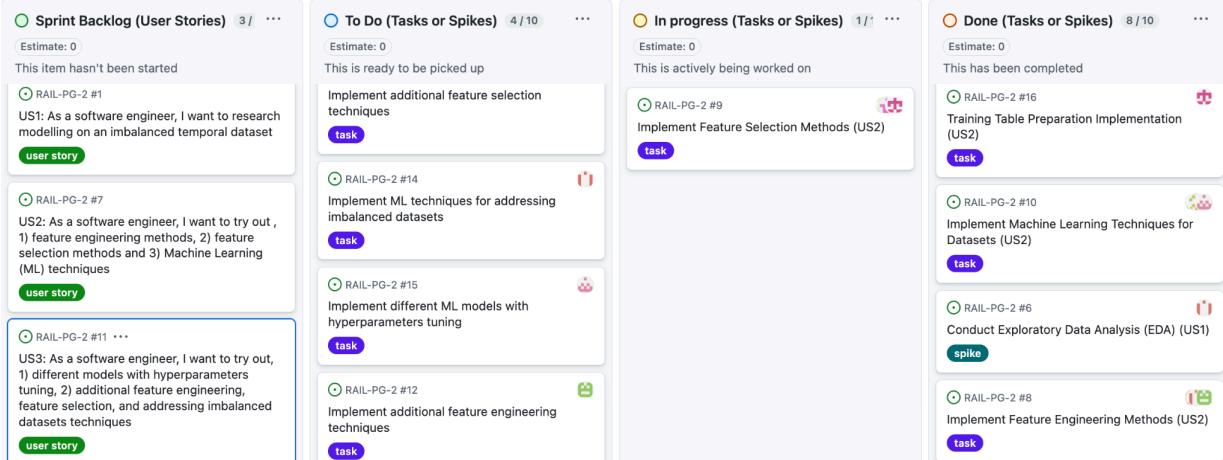  o   Yifan Gu              a1909803
  o   Tianhua Zhang         a1915934
  o   Zihan Luo             a1916700

- Murtaza (Proxy Client)

## 📕 Goals

- Github backlog overview

- Report progress

- Submission results

- Blockers

- Next steps

- QA

## 🗣️ Discussion topics

### 1. GitHub backlog overview

## 2. Report progress

### 2.1 EDA

- Apply visualization techniques (e.g heapmap, histogram, scatter plot, and time series plot) to analyze data

### 2.2 Feature selection

- Joining and integrating trainingcontext, wagondata, and tonnagedata table into a unified total_training_table
- Preprocessing the total training table
- Implement a basic version of the Transformer model with Group Lasso and REF using LightGBM, conducting comparative testing to evaluate their performance.

### 2.3 Feature engineering

- Sensor features and threshold design
- Baseline features, fourier transform, and trend features

### 2.4 ML model training

- Training SVM, DNN, and transformer models and submit

### 2.5 Build production line

## 3. Submission results

- **Transformer**

| 53fb232c | RAIL-PG-2 | Completed | 16 hours ago | 68bae9b17021.csv | Competition 1 - Legolas | ★ **Accuracy**: 66.53%, **AUC_PR**: 40.24%, **F1_Score**: 14.97% |
|---|---|---|---|---|---|---|

- **DNN**

| 3d12fb91 | RAIL-PG-2 | Completed | 2 days ago | 68b832e14464.csv | Competition 1 - Legolas | **Accuracy**: 65.16%, **AUC_PR**: 36.33%, **F1_Score**: 18.27% |
|---|---|---|---|---|---|---|

- **SVM**

| 4eeac20f | RAIL-PG-2 | Completed | 2 days ago | 68b838b91784.csv | Competition 1 - Legolas | Accuracy: 43.37%, AUC_PR: 61.49%, F1_Score: 51.88% |

## 4. Blockers

- **Data availability**: At the beginning of the sprint 2, no datasets were provided. Therefore, we started to implement this project a little late.
- **Limited training time:** Team members need time to select features, preprocessing data, and create training tables for training models. Therefore, time for building and training model was quite limited.
- **Unstable platform runtime**: After submitting notebooks, the runtime duration on the IF platform was unstable. For example, sometimes the execution took half an hour, while other times it finished within 4 minutes. The instability increases uncertainty for model training and inference creation.

## 5. Next steps

**Goal:** In the sprint 3, we will improve models overall performance and achieve an F1 score over 55%.

- Implement at least 2 methods of feature engineering, feature selection, and handling imbalance datasets.
- Tuning hyperparameters of 3 models, feature selection, feature engineering, and imbalanced dataset handling technique.

## 6. QA

Q: For addressing imbalanced datasets, do we need to create a balanced table by preprocessing techniques, or can we directly apply handling imbalanced datasets techniques when training model (e.g. class weights)?