# Natural Language Processing

## Team

# ⚗ Catalyst

### ~ Corporate Joker ~

| Data Collection | <ol><li>https://www.reddit.com/r/Jokes/</li><li>https://bestlifeonline.com/what-do-you-call-jokes/</li><li>https://www.kaggle.com/code/alohahejahe/what-when-why-how/data</li><li>We scraped different websites and scraped data from reddit by using their API Calls.</li></ol> |
|---|---|
| Steps of Data Wrangling | 1. Discovery<br>    a. We compiled the Disparate, Siloed data sources and configure each of them so they can be understood and examined to find patterns and trends in the data |
| Model Training | We used the **Sec2Sec Deep learning model** because of variable length input and output to train our Prediction System. We also combined an attention module to make our model smarter. |
| Model Evaluation | We followed **92/4/4 distribution** to train our model on test data sets. Although we faced a lot of difficulties in model evaluation, we finally managed to increase the perplexity of our model by 15%. |

| Model Deployment | The *Source code of the Project is hosted on github [here](here) and the Project is deployed on to streamlit.io [here](here)* |
|---|---|
| **MLOPS Practices** | |

| MLOps Best Practices | Data | ML Model | Code |
|---|---|---|---|
| Documentation | 1) Data sources<br>2) Decisions, how/where to get data<br>3) Labelling methods | 1) Model selection criteria<br>2) Design of experiments<br>3) Model pseudo-code | 1) Deployment process<br>2) How to run locally |
| Project Structure | 1) Data folder for raw and processed data<br>2) A folder for data engineering pipeline<br>3) Test folder for data engineering methods | 1) A folder that contains the trained model<br>2) A folder for notebooks<br>3) A folder for feature engineering<br>4)A folder for ML model engineering | 1) A folder for bash/shell scripts<br>2) A folder for tests<br>3) A folder for deployment files (e.g Docker files) |