# MLQDA

Zsolt Takacs – 2472886T

# Developing a Machine Learning webapp for Qualitative Data Analysis

# Overview

- Overview of presentation
- Background
- Requirements
- Design
- Technologies
- Demonstration of finished system
- Testing
- Evaluation
- Future Work

# Background

**How do Machine Learning and Qualitative Data Analysis meet?**

- Qualitative techniques are often described as subjective, unreliable and extremely time consuming

- Machine Learning techniques are becoming accessible rapidly

- Using machine learning to support qualitative data analysis can help overcome the mentioned issues
  - Less subjective coding
  - More reliable and reproducible results
  - Much faster initial coding process

- This has been proposed by multiple researchers previously and have been used in research already

- Most of the prominent applications to support qualitative data analysis offer similar solutions too
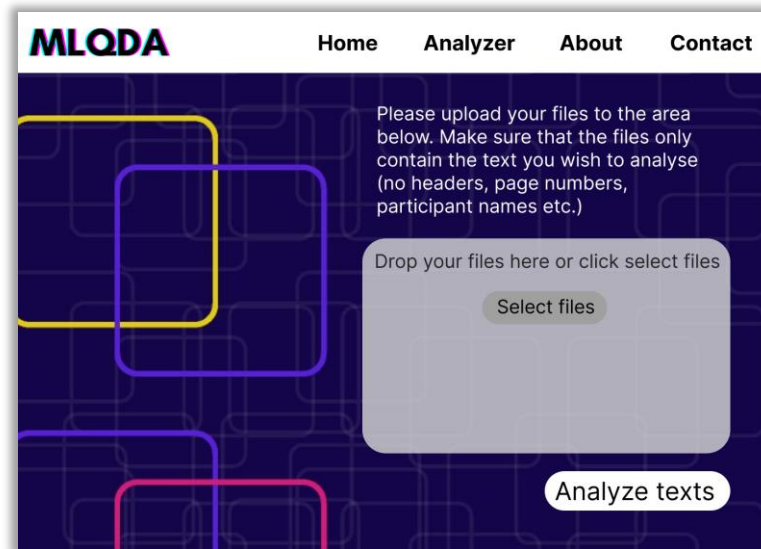
# Requirements

- Functional requirements prioritised using the **MoSCoW** technique
  - Container requirements:
    - Users must be able to upload a text for processing
    - Users should receive valuable information about how the product processes their texts
    - The entire system could compile a zip file of all the result files to avoid cluttering and
  - ML script requirements:
    - The Machine learning script must be able to work with different file extensions
    - The Machine learning script should be able to automatically calculate the optimal parameters for the machine learning models of their results
    - The Machine learning script could provide visualization of the most important words and their contribution to the topics
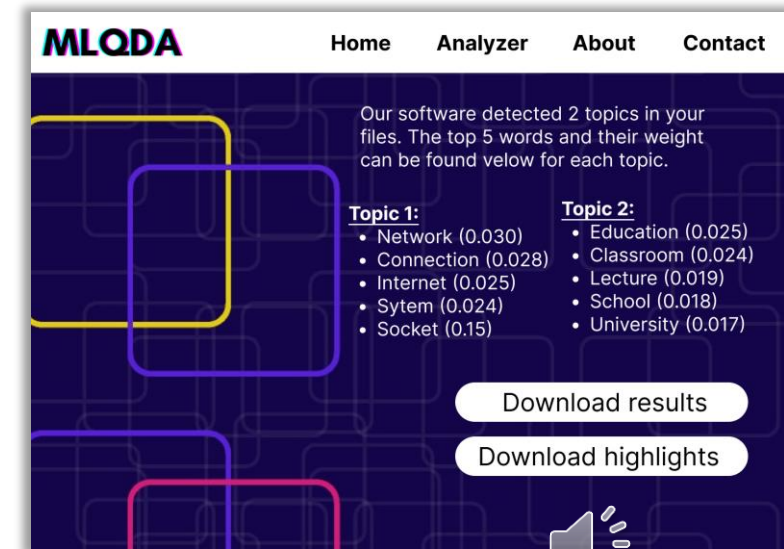- Non-functional requirements

# Design

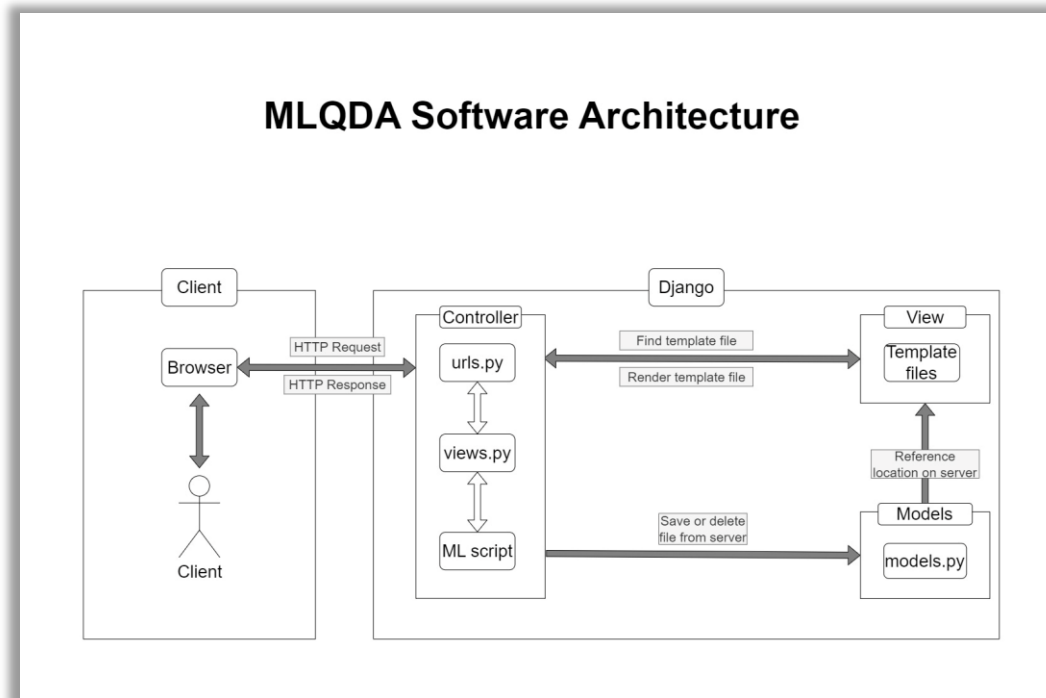Based on requirements Figma wireframes were developed and converted into a realistic prototype
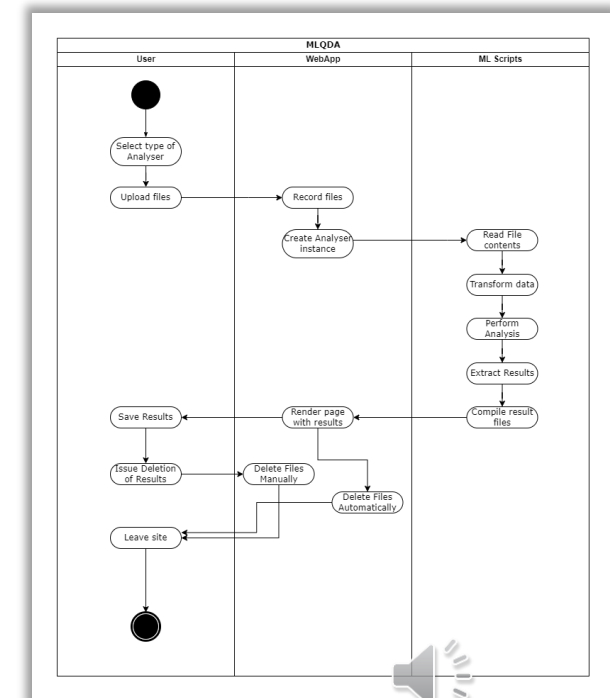


*Analyser start page*



*Analyser results page*

# Design

Based on requirements a system architecture diagram and a user activity diagram was sketched out to represent the inner logic of the proposed system.



*Proposed Software Architecture*



*Activity Diagram*

MLQDA

Technologies

# Demonstration of the finished system

# Testing

Two major type of testing – Manual Testing and Unit Testing:
1. Manual Testing table filled out after every deployment
2. Unit tests developed concurrently to code base and run as part of the CI pipeline



*Manual Testing Table*



*Unit Test Coverage Report*

# Evaluation

The evaluation of the system took place in two steps
1. Pilot user evaluation by an qualitative analysis expert to gather general feedback
2. Open user evaluation to quantify system usability and gather more feedback

# Future Work

- Include alternative machine learning techniques
- Improve Topic Modelling visualisation by updating pyLDAvis
- Scale up Topic Modelling to allow a more dynamic way of hyperparameter usage
- Automate highlighting after manual changes in topic words
- Include more instructions to support first-time users

# References

- Blei, D. M., Ng, A. Y., & Jordan, M. I. (2003). Latent dirichlet allocation. Journal of machine Learning research, 3(Jan), 993-1022.
- Braun, V., & Clarke, V. (2006). Using thematic analysis in psychology. Qualitative research in psychology, 3(2), 77-101.
- Chen, N. C., Drouhard, M., Kocielnik, R., Suh, J., & Aragon, C. R. (2018). Using machine learning to support qualitative coding in social science: Shifting the focus to ambiguity. ACM Transactions on Interactive Intelligent Systems (TiiS), 8(2), 1-20
- Crowston, K., Allen, E. E., & Heckman, R. (2012). Using natural language processing technology for qualitative data analysis. International Journal of Social Research Methodology, 15(6), 523-543.
- Janasik, N., Honkela, T., & Bruun, H. (2009). Text mining in qualitative research: Application of an unsupervised learning method. Organizational Research Methods, 12(3), 436-460.
- Parks, L., & Peters, W. (2022). Natural Language Processing in Mixed-methods Text Analysis: A Workflow Approach. International Journal of Social Research Methodology, 1-13.