# Video Perception Analyzer
## Project Advisor: Vijay Eranti

**Anandram,Swathi (MS Software Engineering)**
**Jadon, Aryan (MS Software Engineering)**
**Nalam, Harika (MS Software Engineering)**
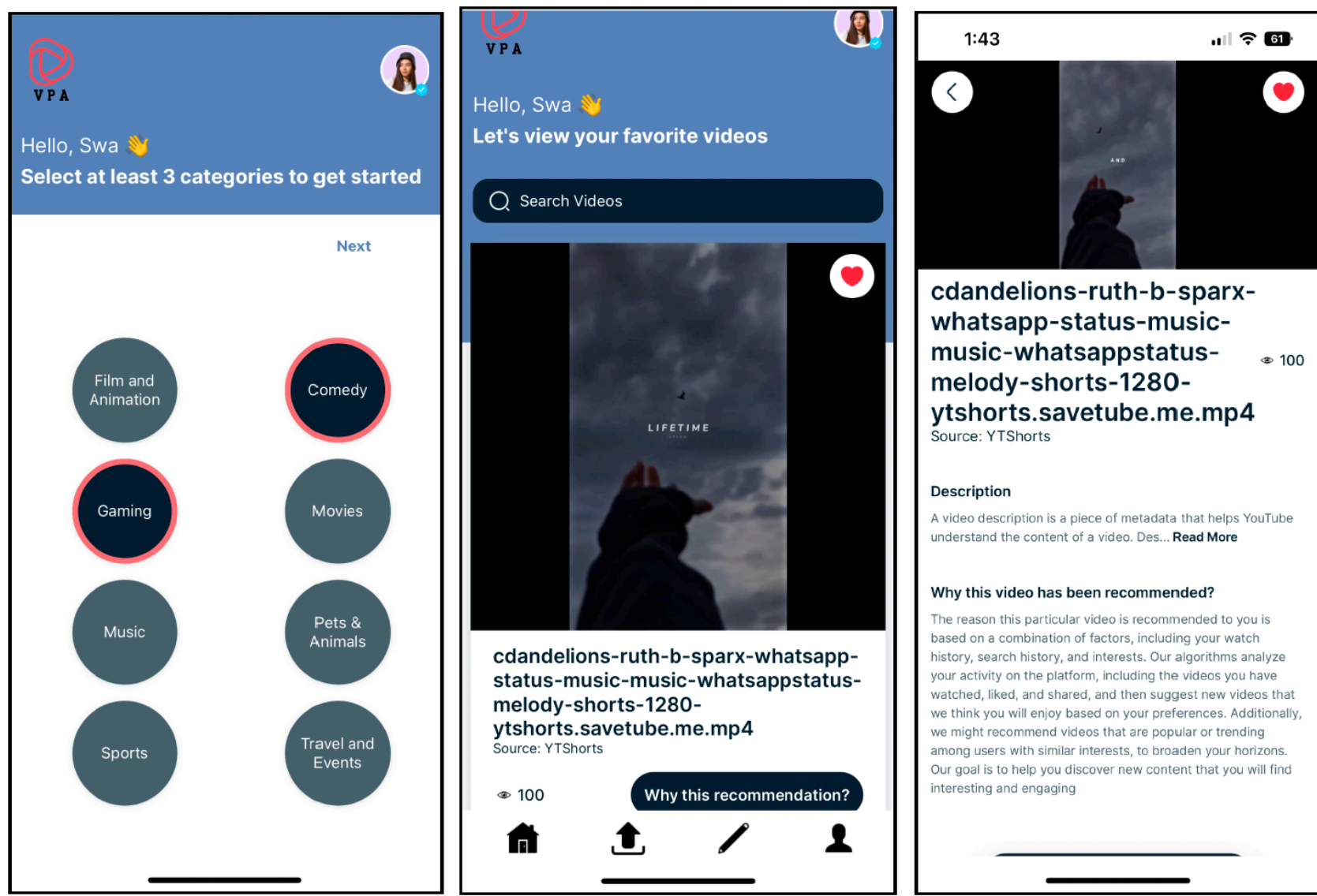**Nimbhorkar, Shreya (MS Software Engineering)**

## Introduction

In recent news, California legislators have passed a groundbreaking bill that compels social media platforms to publicly disclose their content screening and recommendation policies. This mandate aims to combat the alarming spread of hate speech and violence online.

However, numerous platforms have fallen short of the government's requirement to filter restricted content based on user age. In response, we present an innovative, state-of-the-art solution designed to uphold these essential government regulations.

Our project proposes the development of a cutting-edge tool, the Video Perception Analyzer, capable of automatically analyzing video content using a multimodal approach. By leveraging the power of multiple Deep Learning (DL) and Machine Learning (ML) models, this system will extract valuable insights from video data, including images, audio, and video metadata. The primary goal is to create a safer online environment for users.



## Methodology

### UI implementation

The UI contains these components :

1. User Interface
2. User authentication
3. Video Uploading and Storage
4. Video Player
5. Video feed and recommendation system



## Methodology

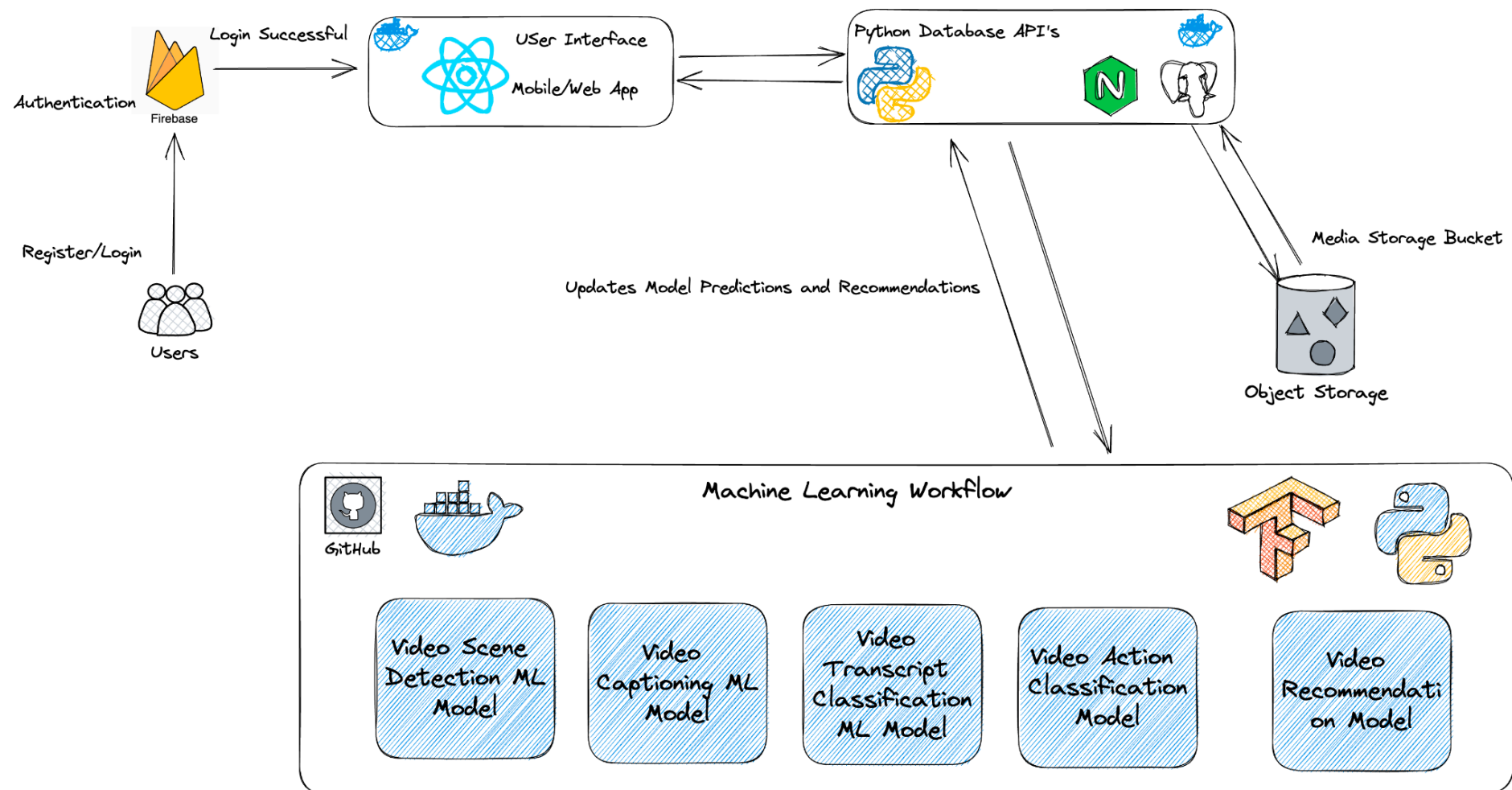### Middle Tier implementation

The Middle Tier Implementation have these components :

1. Environment Setup and Dependency Installation
2. User Profile Creation and Editing API
3. User Details and Preferences Retrieval API
4. Admin Dashboard and User Videos API
5. Video Recommendation API and ML Agent Interaction

### Machine Learning implementation

The Machine Learning implementation have these components:

1. Data Pipeline
2. REST API Development for Inference
3. Inference using API's POST Requests
4. Status Request using API's GET Requests
5. Storage of Inference Results in Postgres Database



### Data Tier implementation

The Data Tier Implementation have these components:

1. Digital Ocean Spaces Buckets
2. Digital Ocean Droplets
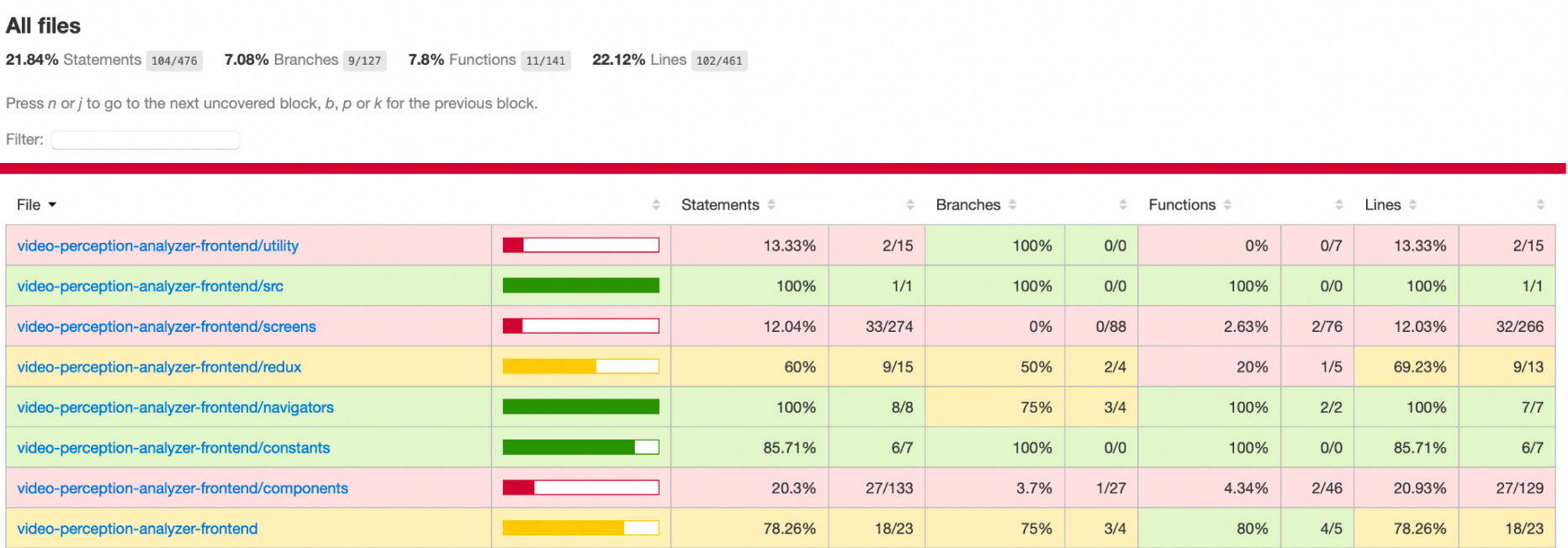3. Postgres Relational Database



Hence, the project implementation emphasizes on security, reliability, compatibility, availability, scalability, maintainability, localization, and usability. The application is designed to provide a secure and reliable video streaming experience while being compatible with different screen sizes, easy to maintain, and user-friendly.

## Analysis and Results

The testing process involves using carefully curated datasets that represent a diverse range of scenarios and cover different aspects of the problem domain. Performance metrics such as accuracy, precision, recall, F1-score, or mean Average Precision (mAP) are computed to assess the models' predictive abilities. Additionally, techniques such as cross-validation or holdout validation are employed to validate the models' performance and identify any potential overfitting or under-fitting issues.

The testing phase also involves rigorous error analysis to understand the models' limitations and identify areas for improvement. Throughout the testing process, the ML models are fine-tuned and iteratively refined to achieve optimal performance. By conducting comprehensive testing, the project ensures the reliability and effectiveness of the ML models in providing accurate insights and supporting the application's video understanding capabilities.



The Image above shows the code coverage report of the react-native application. This report shows how much of the code in each folder of our project is covered by the test file of the project. The table shows that almost 78% of the code is covered by tests, which leaves less room for errors in the project.



The table above shows the results of various models on Flicker30K dataset using zero-shot image-text retrieval. These results show that the BLIP model outperforms existing methods by a large margin. The video perception analyzer uses the BLIP model based on the above results.

The Image above shows the code coverage report of the react-native application. This report shows how much of the code in each folder of our project is covered by the test file of the project. The table shows that almost 78% of the code is covered by tests, which leaves less room for errors in the project.

This application has been tested using the Jest library, which is a JavaScript testing framework. Basic functionalities of the application like login, page rendering, correct functioning of the button have been tested successfully, which made the application ready for deployment into production

| Test Case | Test Description | Expected Result | Test Result |
|---|---|---|---|
| 1. | Has 1 child | The component has 1 child component. | Pass |
| 2. | Renders email and password input fields | The component renders the email and password input fields. | Pass |
| 3. | Navigates to "Home" screen on button press | The component navigates to the "Home" screen when the "Sign up" button is pressed. | Pass |
| 4. | Renders correctly | The component renders with the expected test ID. | Pass |
| 5. | Login renders the correct | The login component renders consistently over time. | Pass |

## Summary/Conclusions

The Video Perception Analyzer application represents a successful integration of various technologies to create a robust video analysis and recommendation system. By utilizing React-Native, Python, Django, and other cutting-edge tools, the application enables seamless video browsing, searching, and uploading. The integration of machine learning models for video classification, action detection, and content verification ensures accurate analysis and enhanced user safety. With personalized recommendations based on user preferences, the application further enhances user engagement and satisfaction. The Video Perception Analyzer application sets the stage for a sophisticated and enjoyable video browsing experience while leveraging the power of machine learning for accurate analysis and recommendation.

## Key References

[1] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," 2017. [Online]. Available: https://arxiv.org/abs/1706.037622]

[2] M. D. Zeiler and R. Fergus, "Visualizing and understanding convolutional networks," 2013. [Online]. Available: https://arxiv.org/abs/1311.2901

[3] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," 2014.

[4] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," 2015.

[5] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," 2015.

[6] J. Dai, Y. Li, K. He, and J. Sun, "R-fcn: Object detection via region-based fully convolutional networks," 2016. [Online]. Available: https://arxiv.org/abs/1605.06409

## Acknowledgements