

# **DeepFake Detection**

**Dataset creation, feature engineering and classification.**

**[cheliotisnick@gmail.com](mailto:cheliotisnick@gmail.com)**

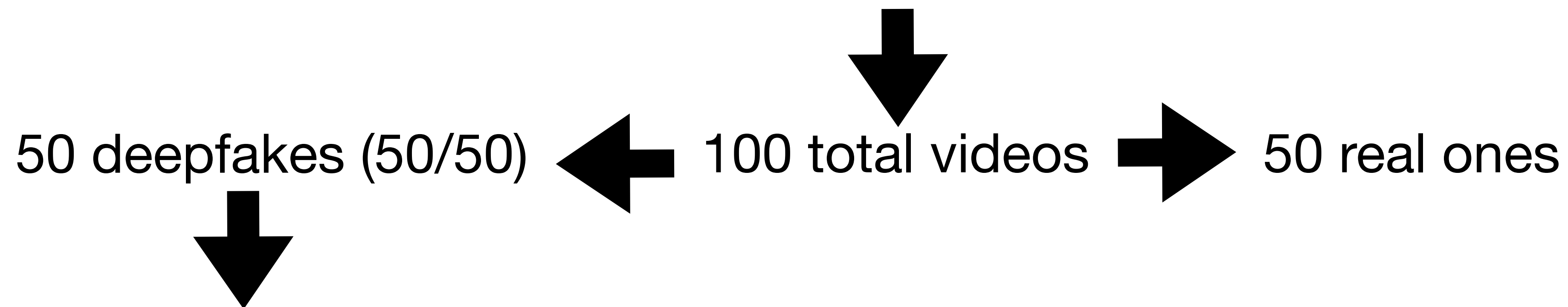
# Dataset creation

# Dataset creation

## Video collection

We collected short, face-focused videos by searching YouTube Shorts. To ensure quality each video had to meet specific criteria:

1. Only one person
2. No background music
3. No face occlusions
4. No sudden camera/angle changes
5. Exactly 10 seconds

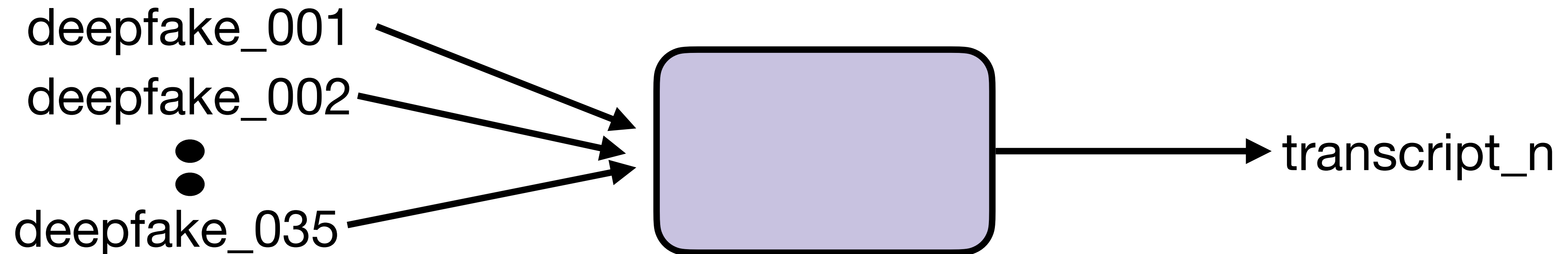


20 video+audio, 15 video and 15 audio.

# Dataset creation

## Deepfake Audio Generation

We used OpenAi's **Whisper**\* model to convert audio into text.



---

We used ElevenLabs\*\* to generate DeepFake voices. A variety of synthetic voice profiles were used to simulate different speech patterns.

\*<https://github.com/openai/whisper>

\*\*<https://elevenlabs.io/>



# Dataset creation

## Deepfake audio insertion

We used **Wav2Lip**\* to sync the deepfake voices with the speaker's lip movements



\*<https://github.com/Rudrabha/Wav2Lip>



# Dataset creation

## DeepFake video generation

We collected 35 unique female/  
male faces.

We used **online face-swapping  
models** to generate deepfake  
videos by replacing the original  
face with synthetic ones.



# Feature engineering

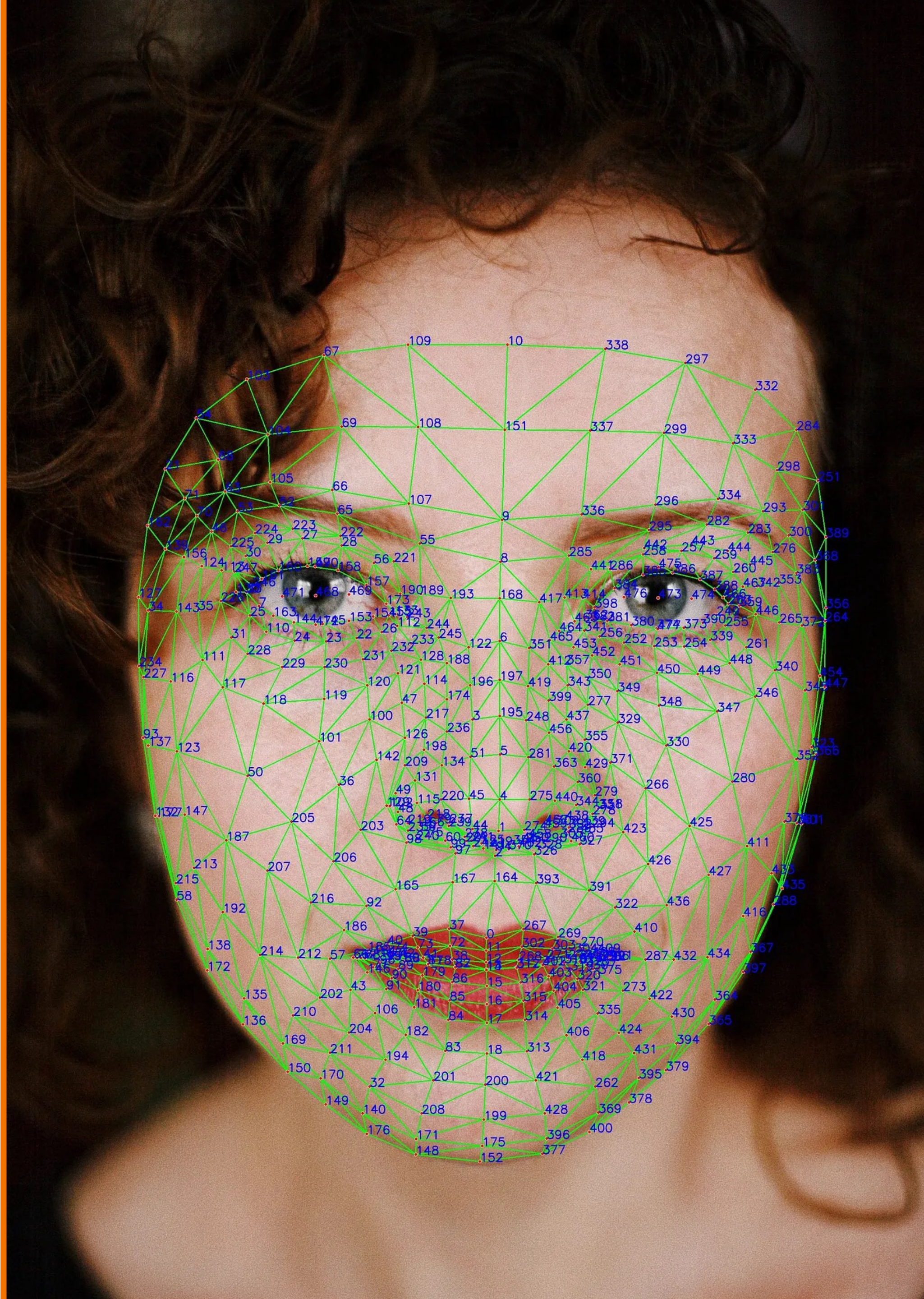


# Feature engineering

## Visual features

To extract visual features we used MediaPipe, a library that detects facial landmarks. From these landmarks we applied custom logic to compute behavioural metrics.

We calculated **blink rate**, **mouth open ratio**, **head motion** (yaw and pitch) and **expression entropy** (using DeepLab)





# Feature engineering

## Visual features - Blink rate

1. Input frame
2. Calculate Eye Aspect Ratio (EAR)
3. If  $EAR < 0.21 \Rightarrow$  potential blink
4. Get next frame
5. If EAR goes back up next frame  $\Rightarrow$  blink
6. Divide Blinks/seconds





# Feature engineering

## Visual features - Mouth open ratio

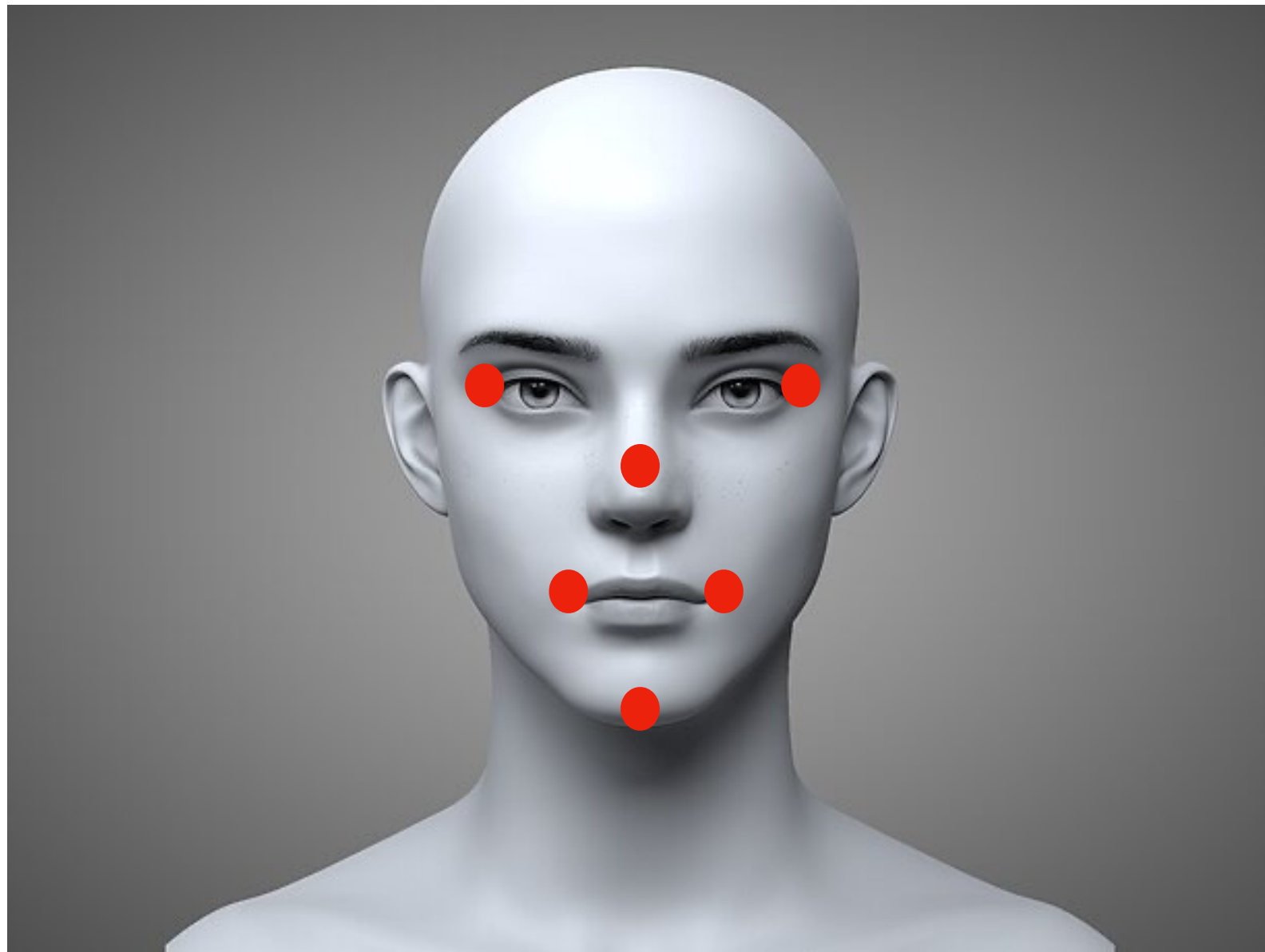
1. Input frame
2. Calculate Mouth Open Ratio (MOR)
3. Store it
4. Calculate the **mean** (or **variance**).



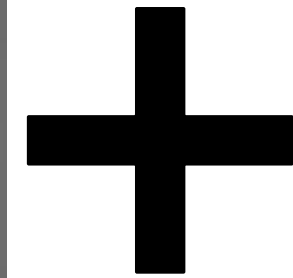


# Feature engineering

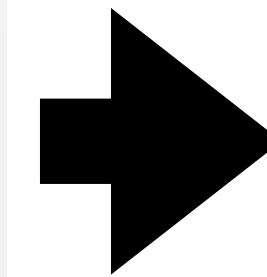
## Visual features - Head motion (left-right and up-down)



3D Face



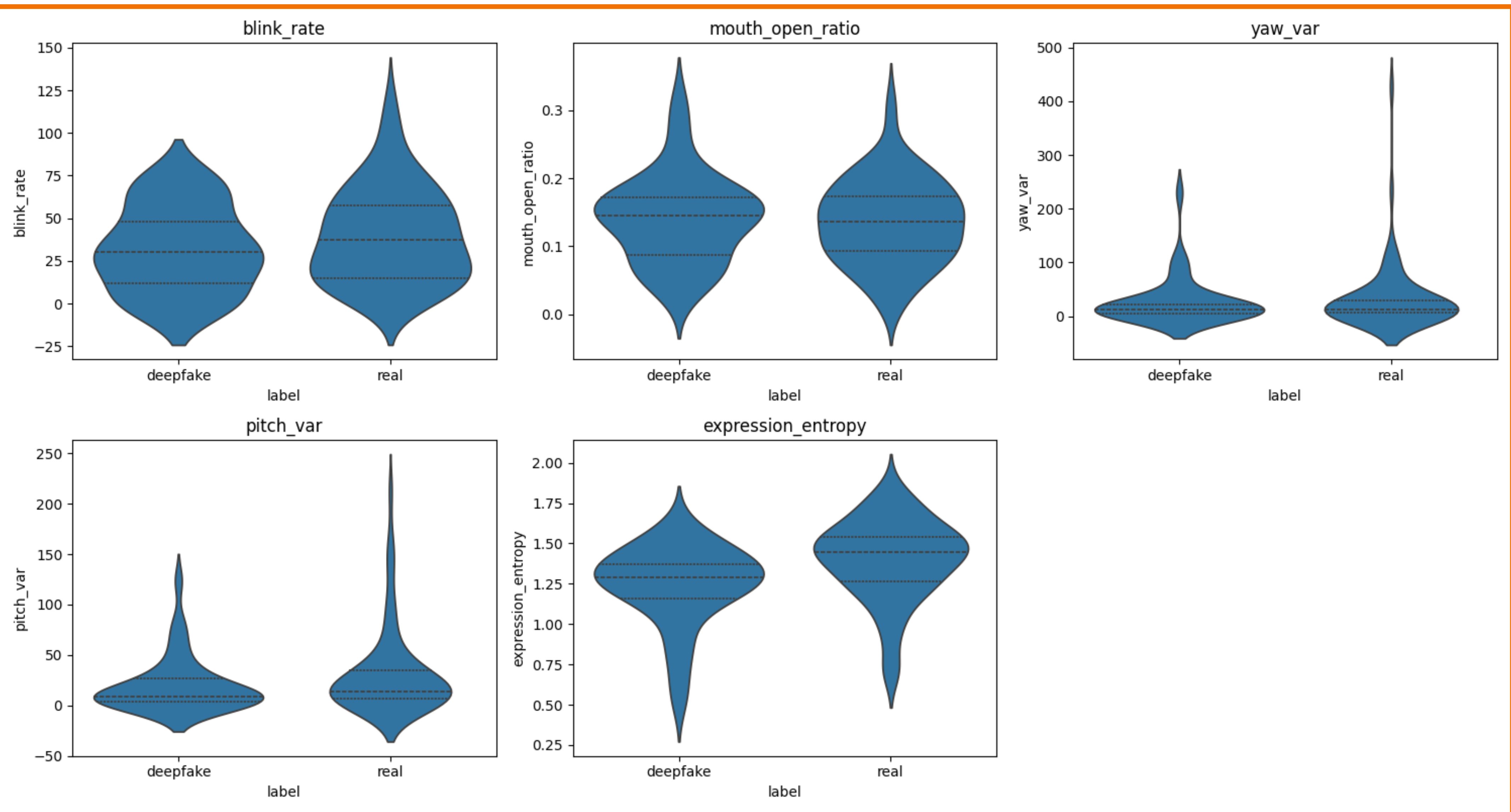
2D Frame



1. Rotate (yaw and pitch) the 3D face so it matches 2D frame.
2. Calculate the variance of yaw and pitch across all frames.

# Feature engineering

## Visual features - Visual feature plots



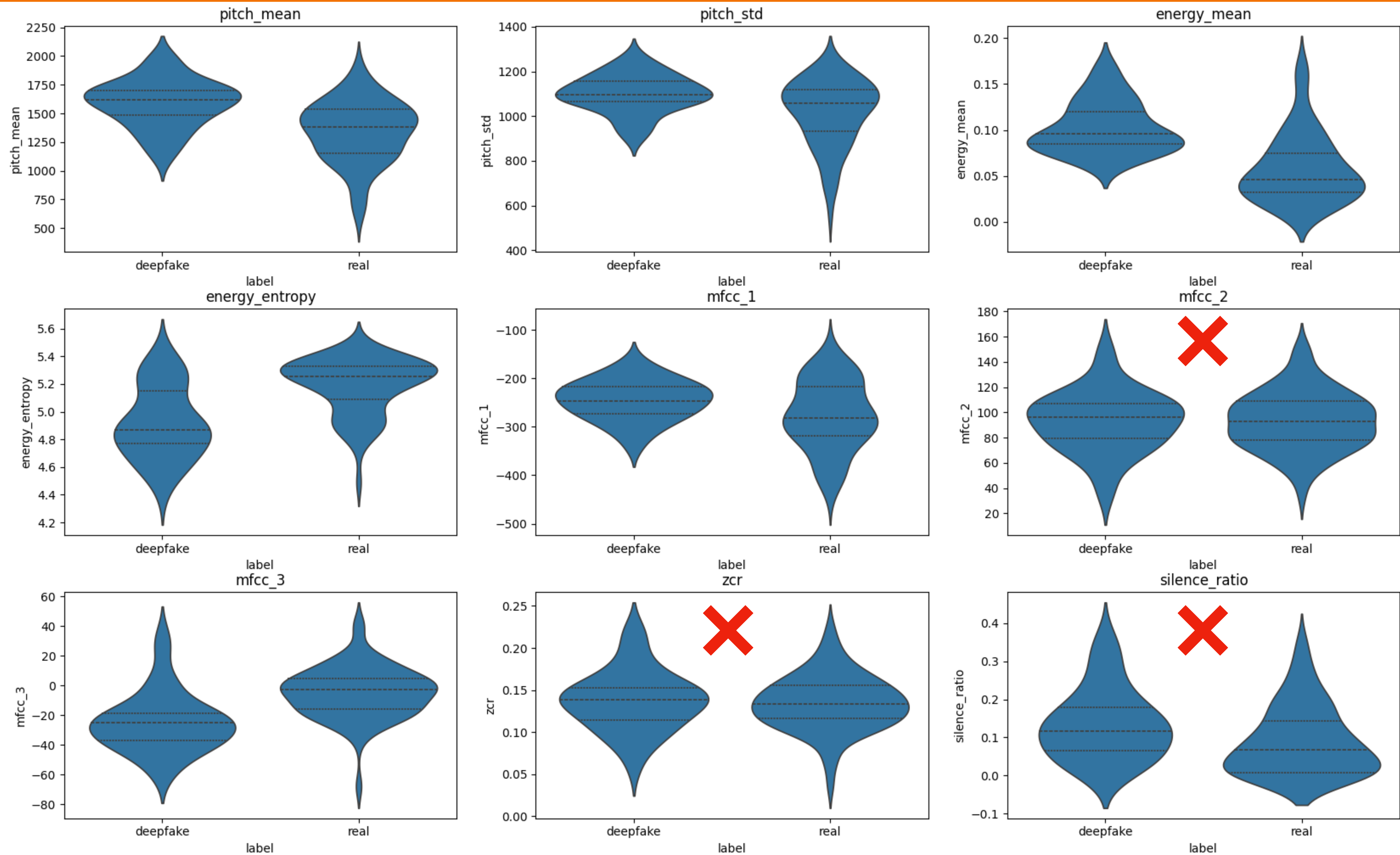


# Feature engineering

## Audio features

We used Librosa to extract **pitch** (mean and std) ,**energy** (mean and entropy), **MFCC** (first 3), **zero cross rate** and **silence ratio**.







# Model training

# Model training

## Models

We trained and evaluated two classification models on our Dataset.Svm and Logistic Regression.Each model was tasked with a binary classification problem, predict whether a video is real or deepfake.

SVM

	Precision	Recall	F1-score	Support
Real	<b>0.88</b>	0.64	0.74	11
Deepfake	0.67	<b>0.89</b>	0.76	9

Accuracy:75%

Logistic Regression

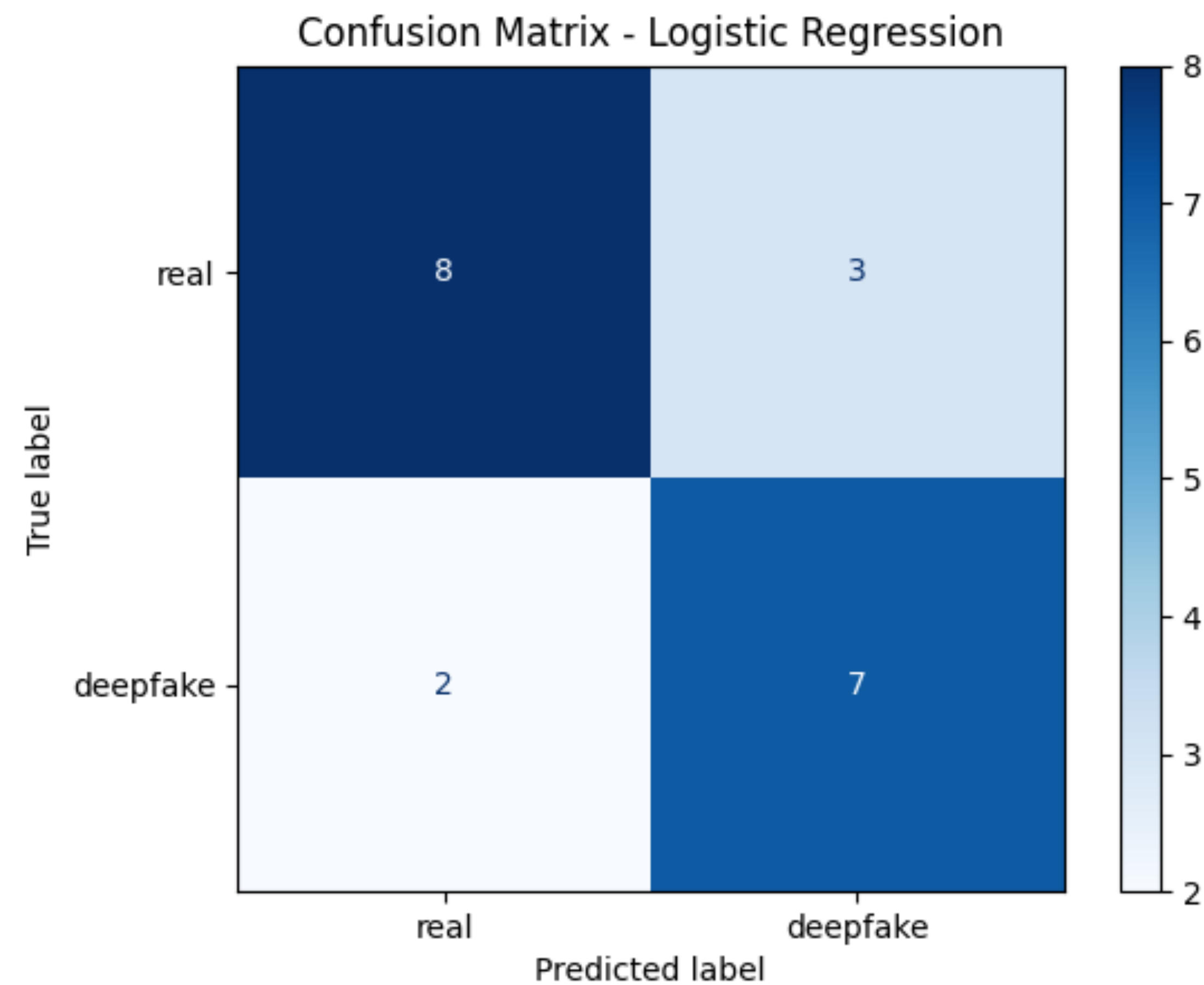
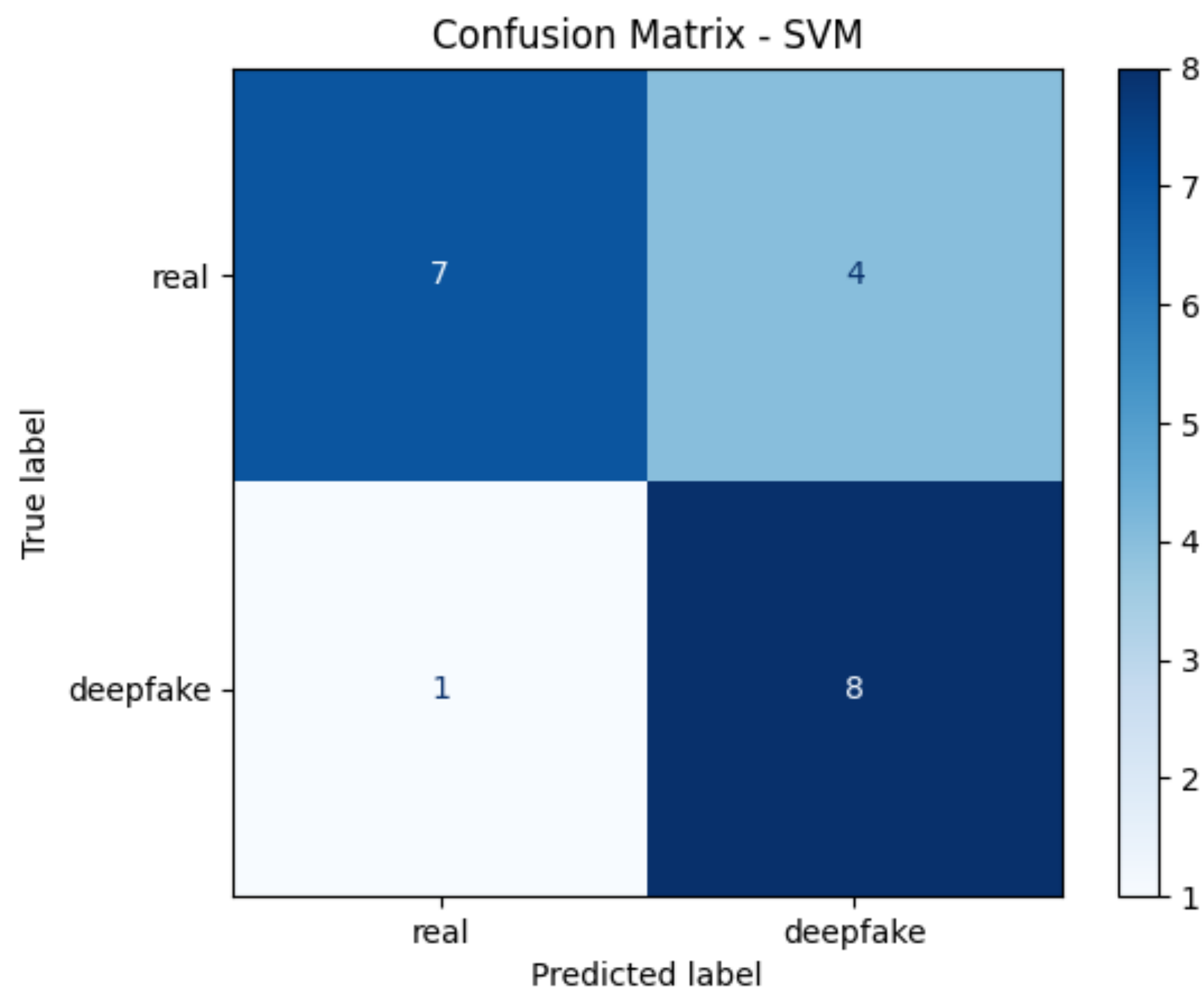
	Precision	Recall	F1-score	Support
Real	0.80	<b>0.73</b>	0.76	11
Deepfake	<b>0.70</b>	0.78	0.74	9

Accuracy:75%



# Model training

## Graphs

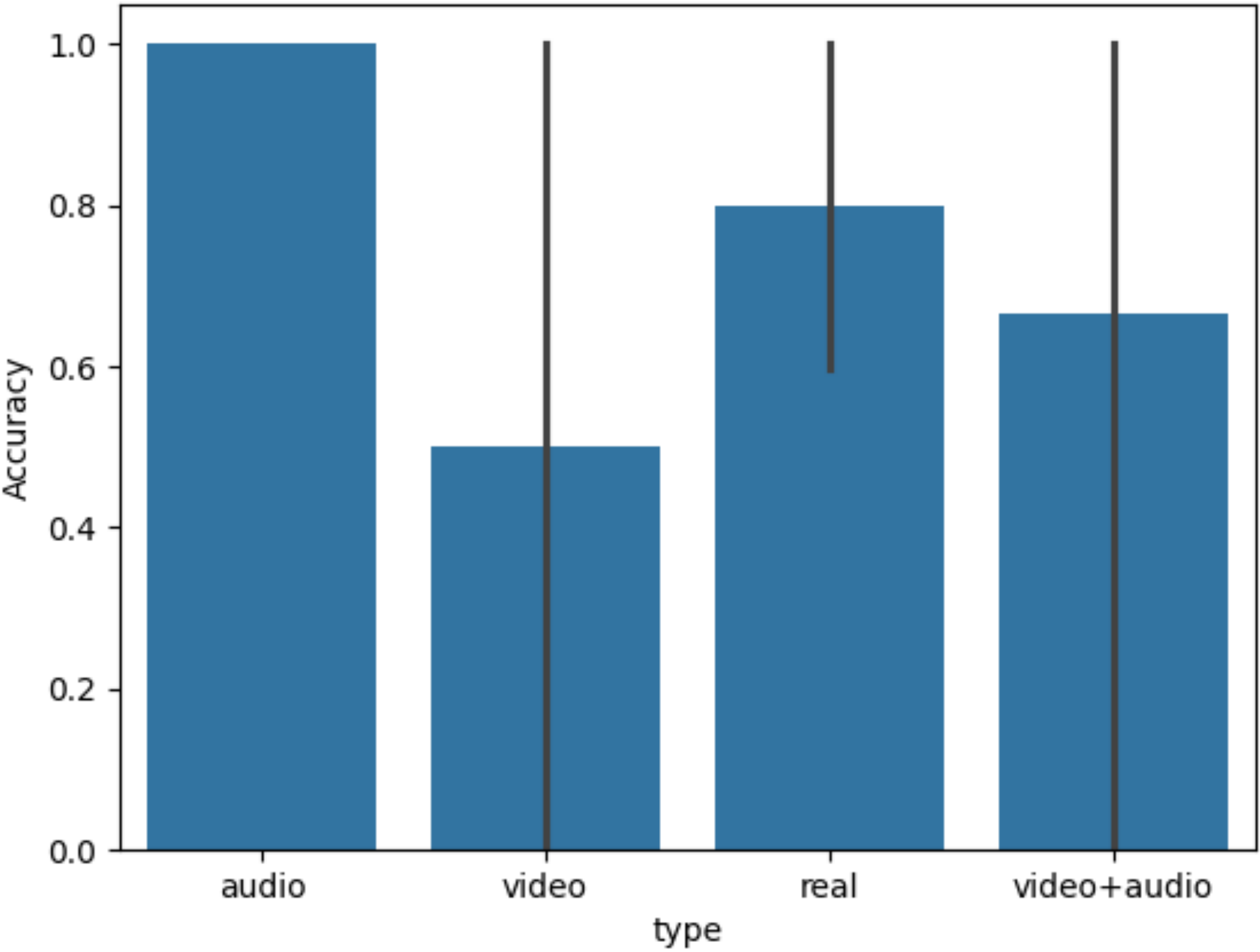


# Model training

## Graphs

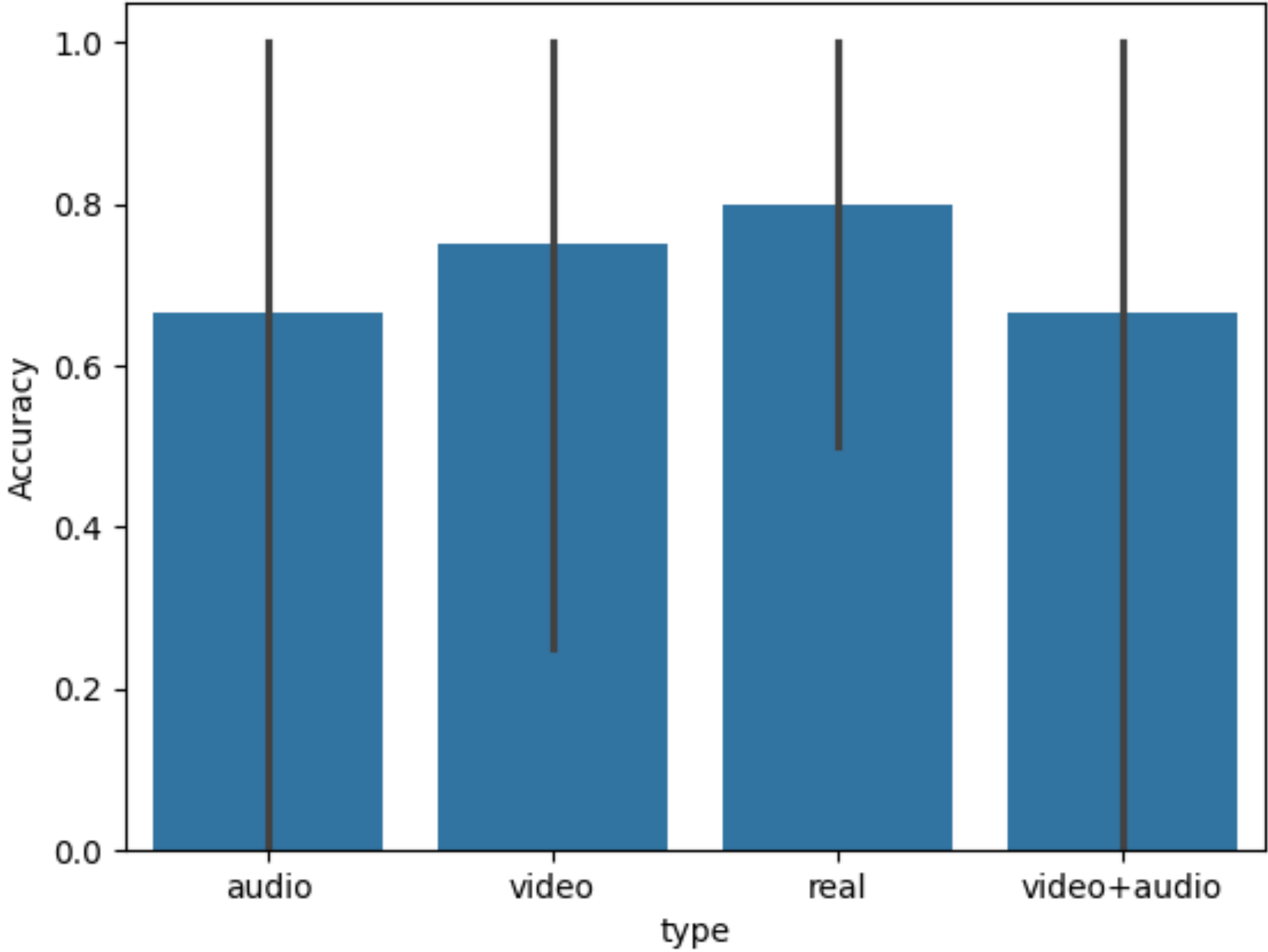
### SVM

Model Accuracy by Deepfake Type



### Logistic Regression

Model Accuracy by Deepfake Type





# How could we do better?

1. More data (quantity and diversity) to improve generalisation.
2. Better visual features.
3. Integrate cross modal consistency features, such as verifying whether mouth movements are synchronized with the spoken words.

**Thank you!**