

# Action spotting on SoccerNet Challenge

MSc Artificial Intelligence for Science and  
Technology

Advanced techniques

Submitted by **Daniele Cecca Mat.914358**





# Introduction



This project focuses on the development and implementation of **classification system for ball events in a football match**, also called **ball action spotting**.

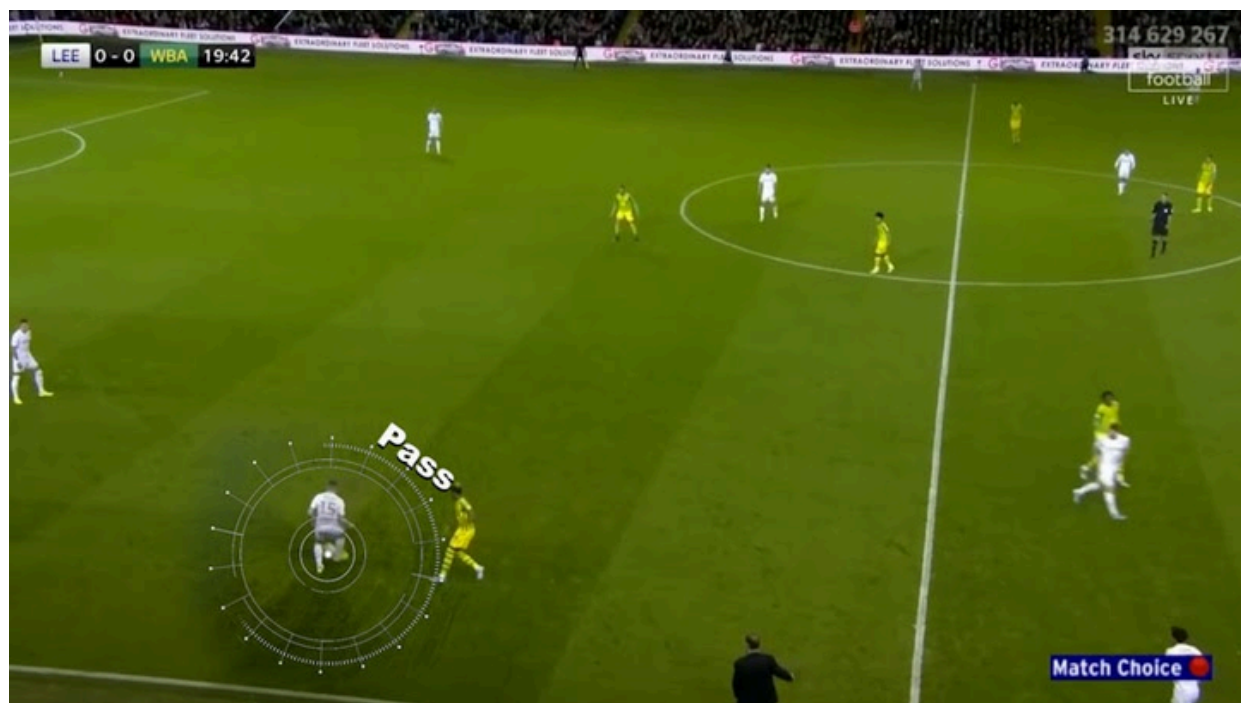
It utilizes the **SoccerNet dataset** and draws inspiration from previously proposed solutions

# PROJECT STRUCTURE



1. Load create and explore the **dataset**
2. Create the **ActionRecognitionDataset**
3. Design and develop the **architecture of the network.**
4. Train the network using different **hyperparameters.**
5. Select the **best-performing network** and train it for additional epochs.
6. **Evaluate** the network's performance.
7. Implementation of a web app with **Streamlit**

# SoccerNet Dataset



The dataset is composed of **7 videos of English Football League games**, and to each video a **JSON** with the **timestamp** and the **label** is associated.

In total we have have **12 different types of action**

Pass	Drive	Header	High Pass
Out	Cross	Throw in	Shot
Ball Player Block	Player Successful Tackle	Free Kick	Goal



# SoccerNet Dataset

## CREATION OF THE DATASET



01

Divide the videos matches into segments

Since defining precise **temporal boundaries** for actions is challenging because it's hard to fix the exact **start** and **end** times, and knowing what occurs after the action can be beneficial, I **extend** the interval by one **second beyond** the defined action times in the JSON file.

02

Create the dataframe

I will apply label encoding on the label

	clip_filename	label	clip_duration
0	/content/data/new_val_data/2019-10-01 - Middle...	PASS	1
1	/content/data/new_val_data/2019-10-01 - Middle...	DRIVE	2
2	/content/data/new_val_data/2019-10-01 - Middle...	PASS	3
3	/content/data/new_val_data/2019-10-01 - Middle...	DRIVE	4
4	/content/data/new_val_data/2019-10-01 - Middle...	HIGH PASS	7

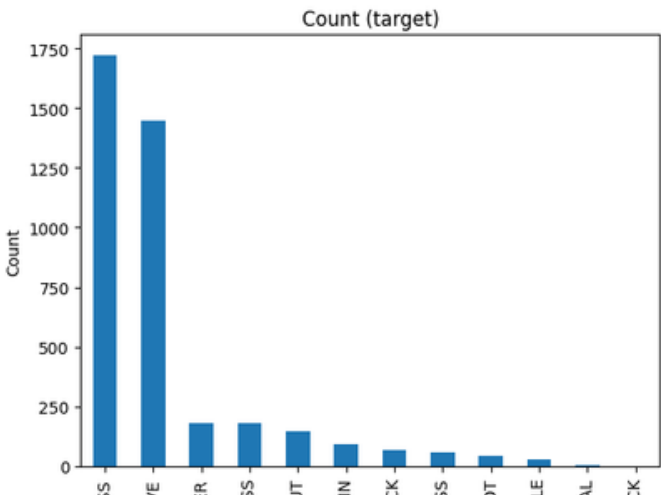
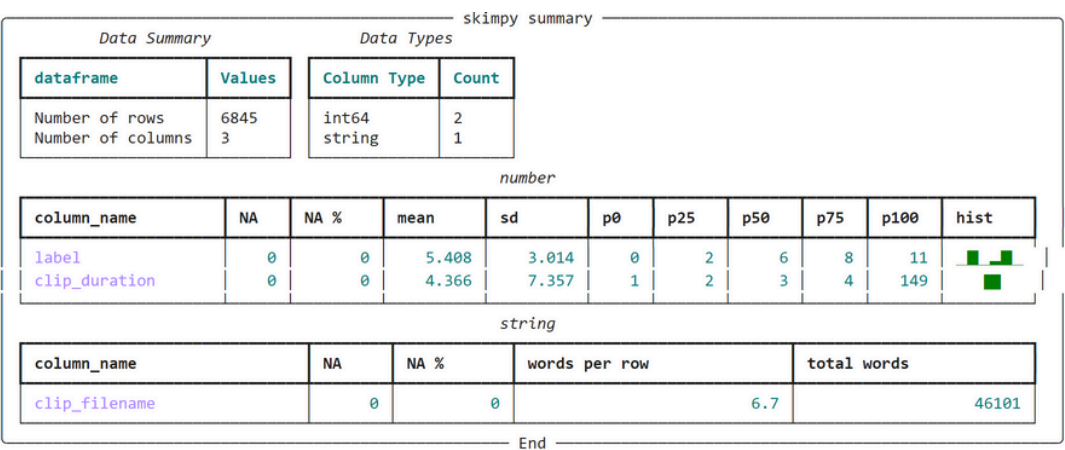


# SoccerNet Dataset

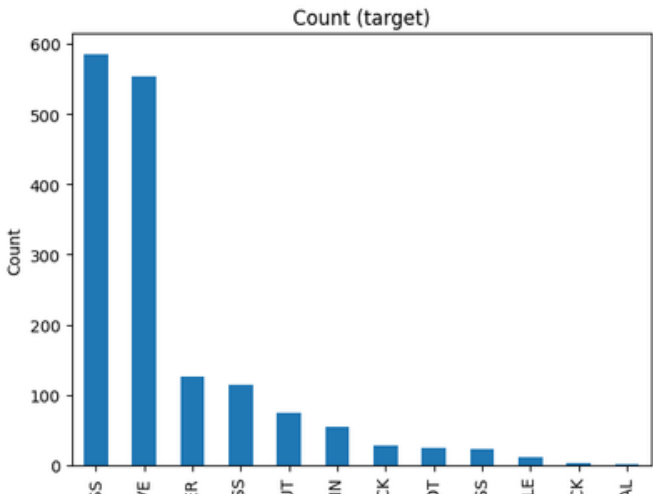
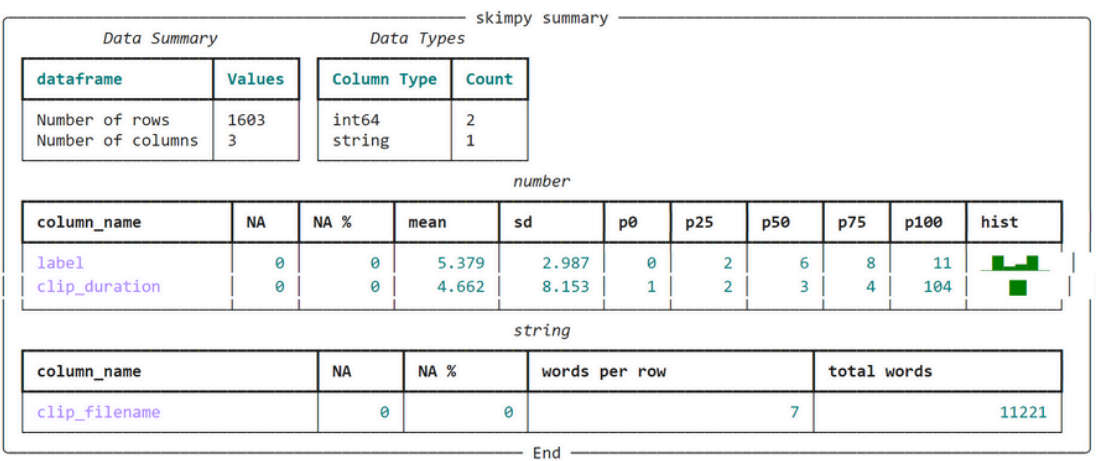
## EXPLORATION OF THE DATASET



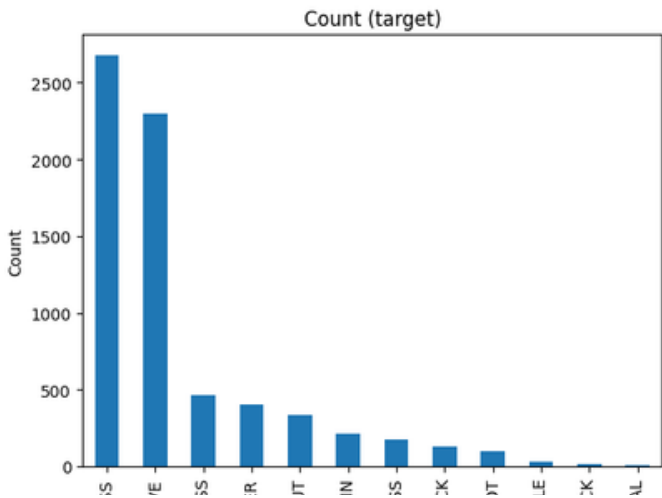
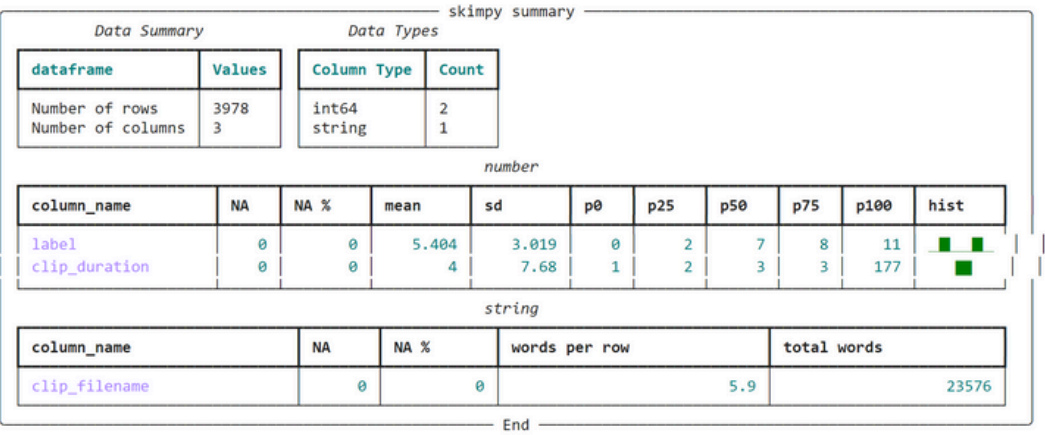
### Training Set



### Validation set



### Test set





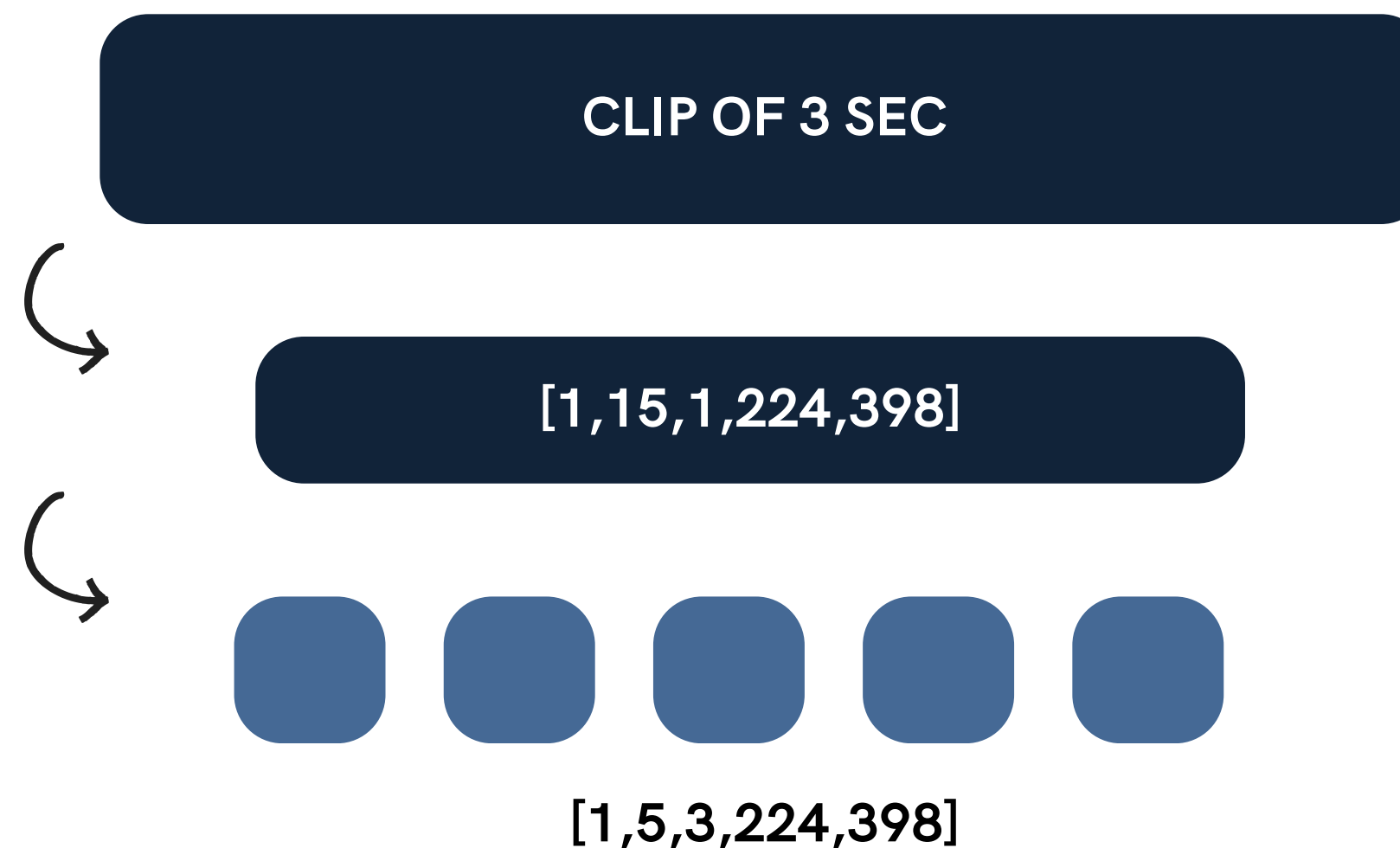
# ActionRecognition Dataset class

The main idea is to take for each clip **15 frames**, one every 12 frames (because the mean length is 4), and then divide these frames into **stacks of 3 frames**.

So in the end we will have a vector of **[1, 5, 3, 224, 398]** where **[B, T, C, H, W]**.

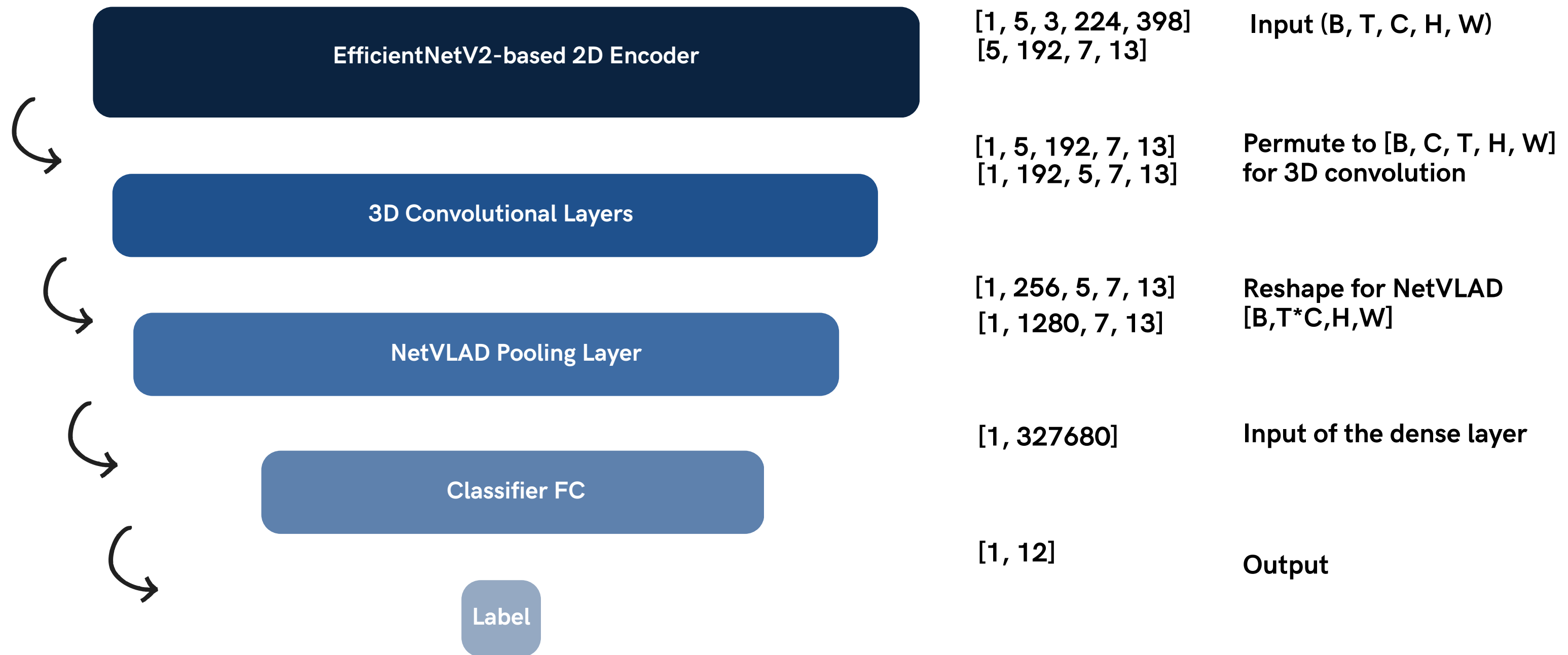
This will be useful to stack temporal information that will be processed by the network.

Videos will be transformed into **grayscale**





# Network architecture



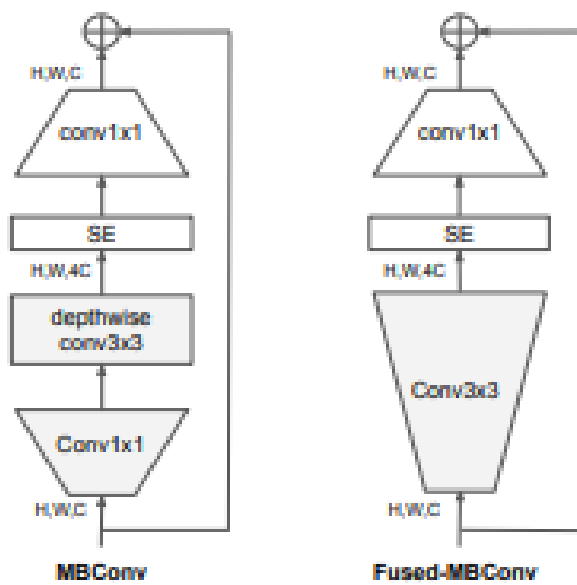


# Building blocks

## NETWORK

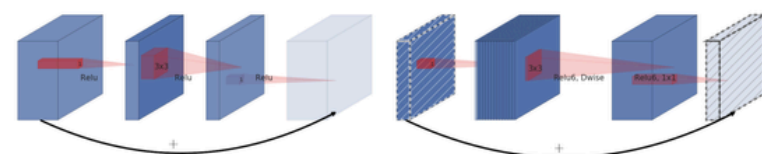
### EfficientNetV2

Inverted residual network



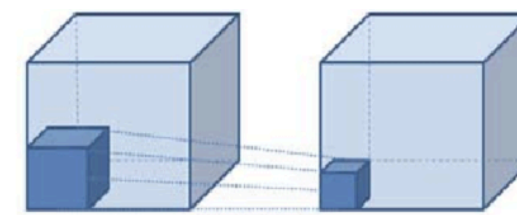
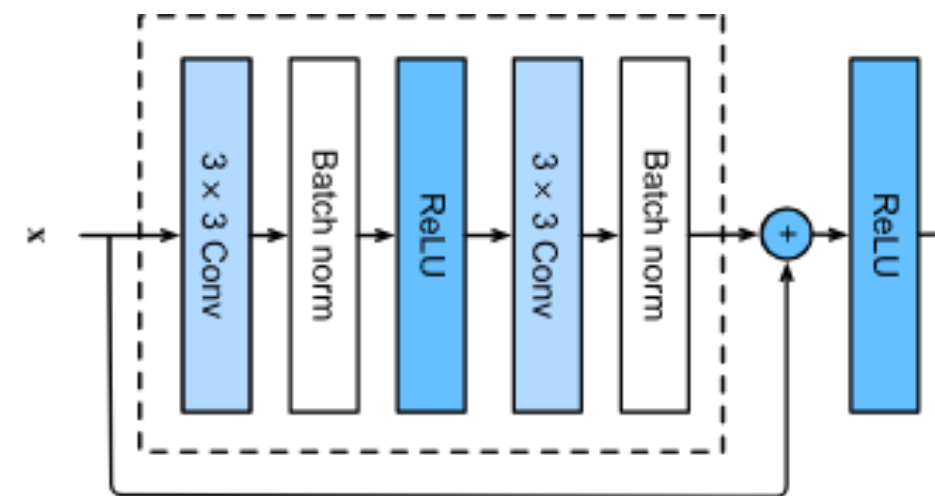
(a) Residual block

(b) Inverted residual block



### Residual block

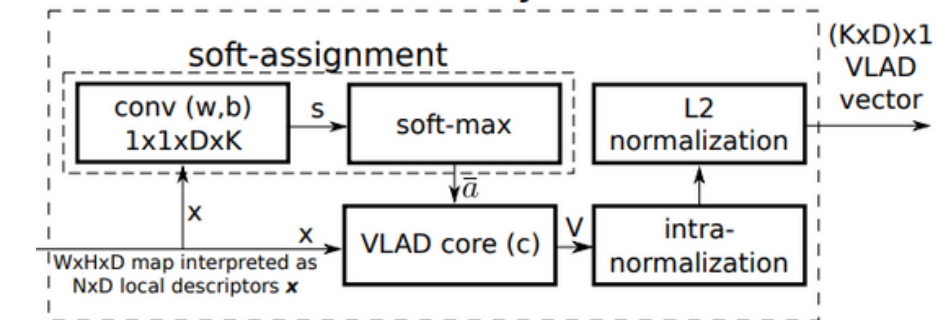
3D Convolution



(b) 3D convolution

### NetVLAD

NetVLAD layer



$$V(j, k) = \sum_{i=1}^N a_k(\mathbf{x}_i) (x_i(j) - c_k(j)),$$

$$V(j, k) = \sum_{i=1}^N \frac{e^{\mathbf{w}_k^T \mathbf{x}_i + b_k}}{\sum_{k'} e^{\mathbf{w}_{k'}^T \mathbf{x}_i + b_{k'}}} (x_i(j) - c_k(j))$$

# HyperParameters Tuning

## NETWORK

I create a configuration that will be used by the **Weights and Biases agent** to set the different **hyperparameters** during various experiments. I used the simplest agent which chose the combination of the **parameters randomly**.

Parameters	Values
epochs	5
learning rate	[0.01, 0.001, 0.0001]
dropout	[0.4, 0.5]
Batch Size	[1,3,5]
loss function	focal loss
optimizer	adam

# Focal Loss

## NETWORK

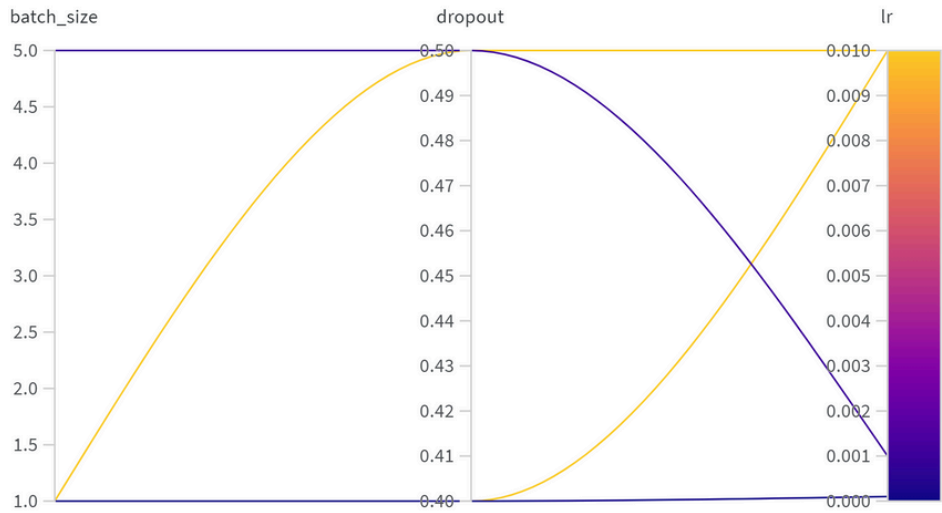
**Focal Loss** is a loss function designed to address the challenge of class imbalance in tasks such as object detection.

Focal Loss is a modified version of the standard **Cross-Entropy Loss** that down-weights the contribution of easy-to-classify examples and focuses more on hard-to-classify examples. It is defined as

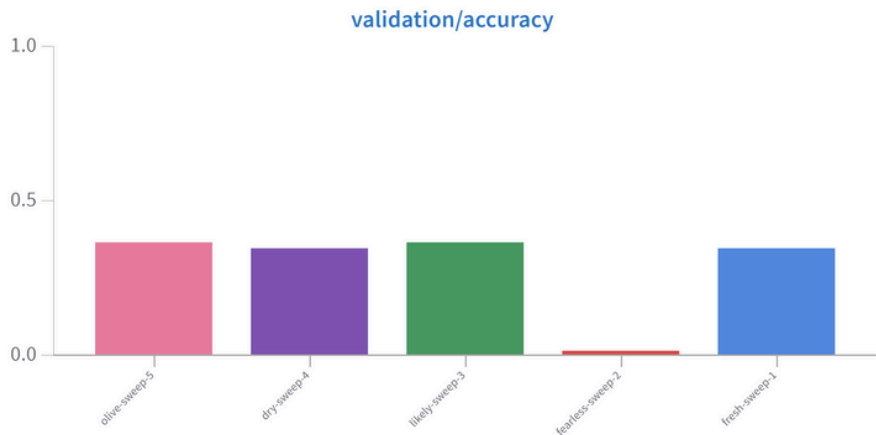
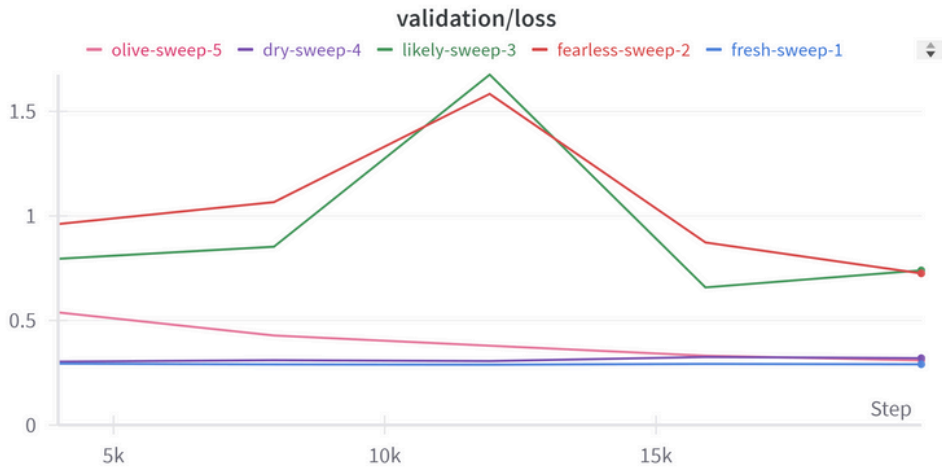
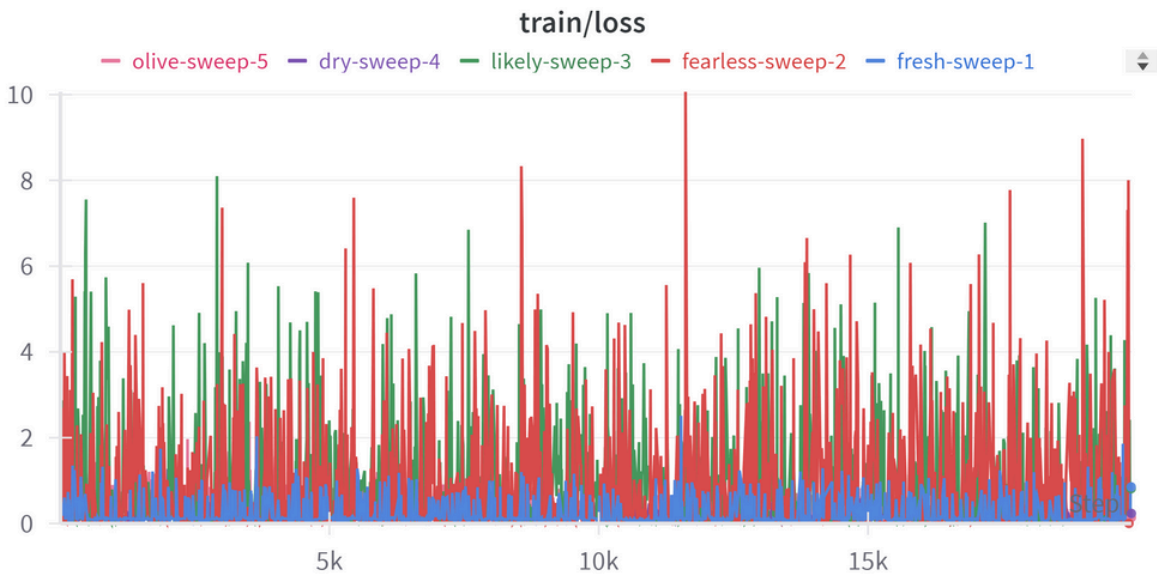
$$FL(p_t) = -\alpha_t (1 - p_t)^\gamma \log(p_t)$$

- $p_t$  is the model's estimated probability for the true class.
- $\alpha_t$  is a weighting factor for the class, balancing the importance of positive/negative example
- $\gamma$  is a focusing parameter that adjusts the rate at which easy examples are down-weighted.

# Training NETWORK



Parameter	Value
Batch Size	1
Drop Out	0.4
Epochs	5
Learning Rate	0.0001



# Training the best model

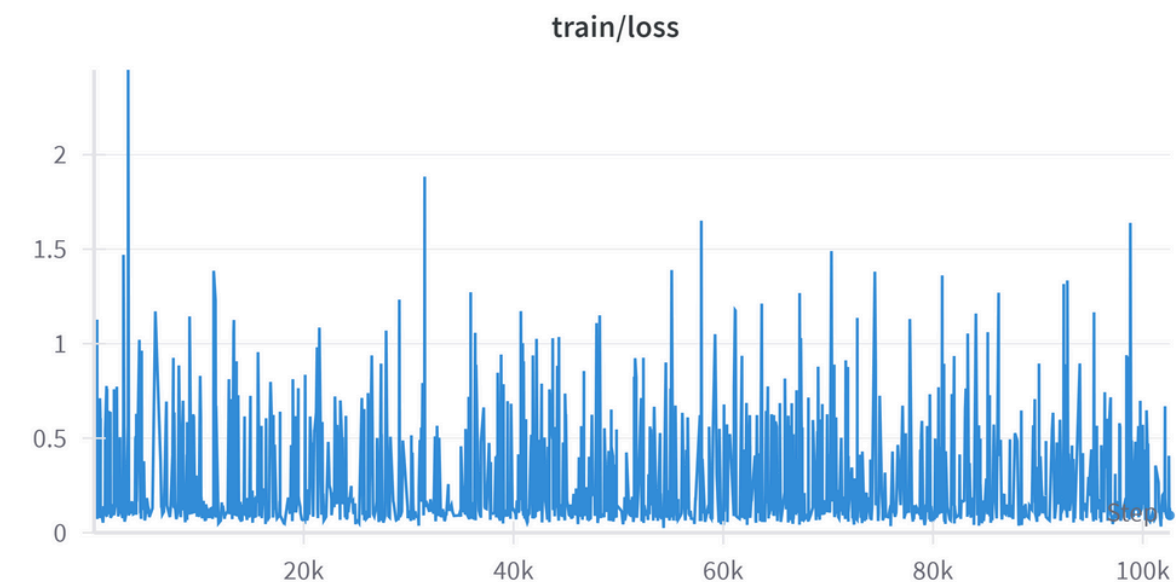
## NETWORK

The behavior of the model is a little bit strange because both the **training loss** and the **validation loss** are **volatile**.

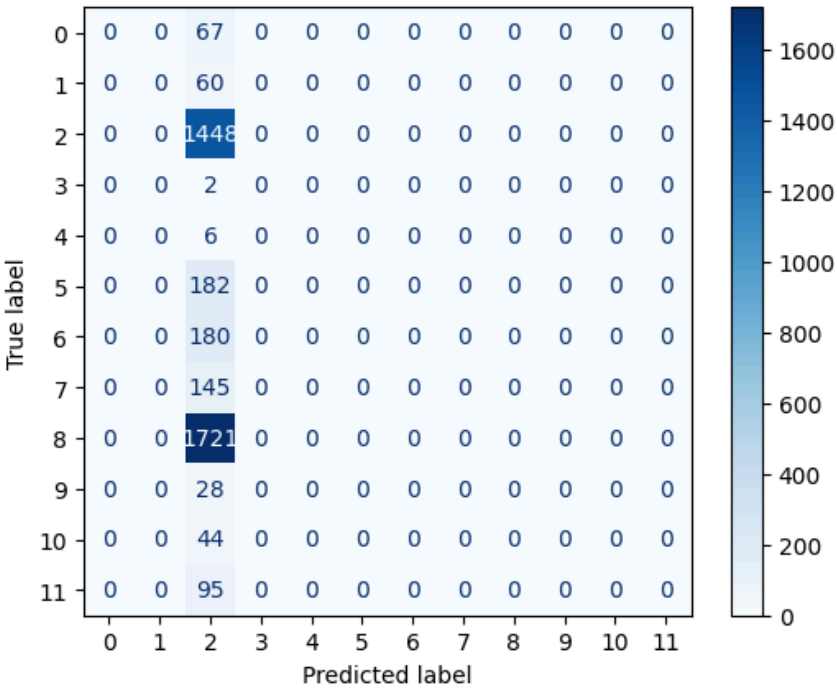
Despite it all I think the model is becoming better at generalizing because both validation and training losses are **decreasing on average**.

So in this case, it is difficult to talk about **overfitting** and **underfitting**.

Config parameter	Importance ① ↓	Correlation
batch_size	<div></div>	<div></div>
dropout	<div></div>	<div></div>
lr	<div></div>	<div></div>

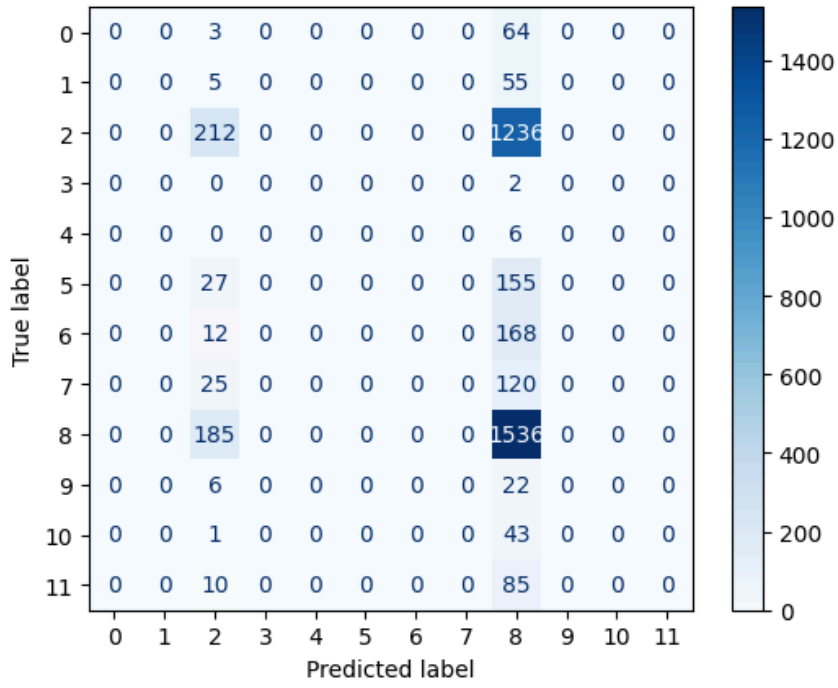


5 epochs model



- accuracy micro: 0.36
- precision micro: 0.36
- recall micro: 0.3
- f1 micro: 0.36
- accuracy macro: 0.36
- precision macro: 0.03
- recall macro: 0.08
- f1 macro: 0.045

20 epochs model



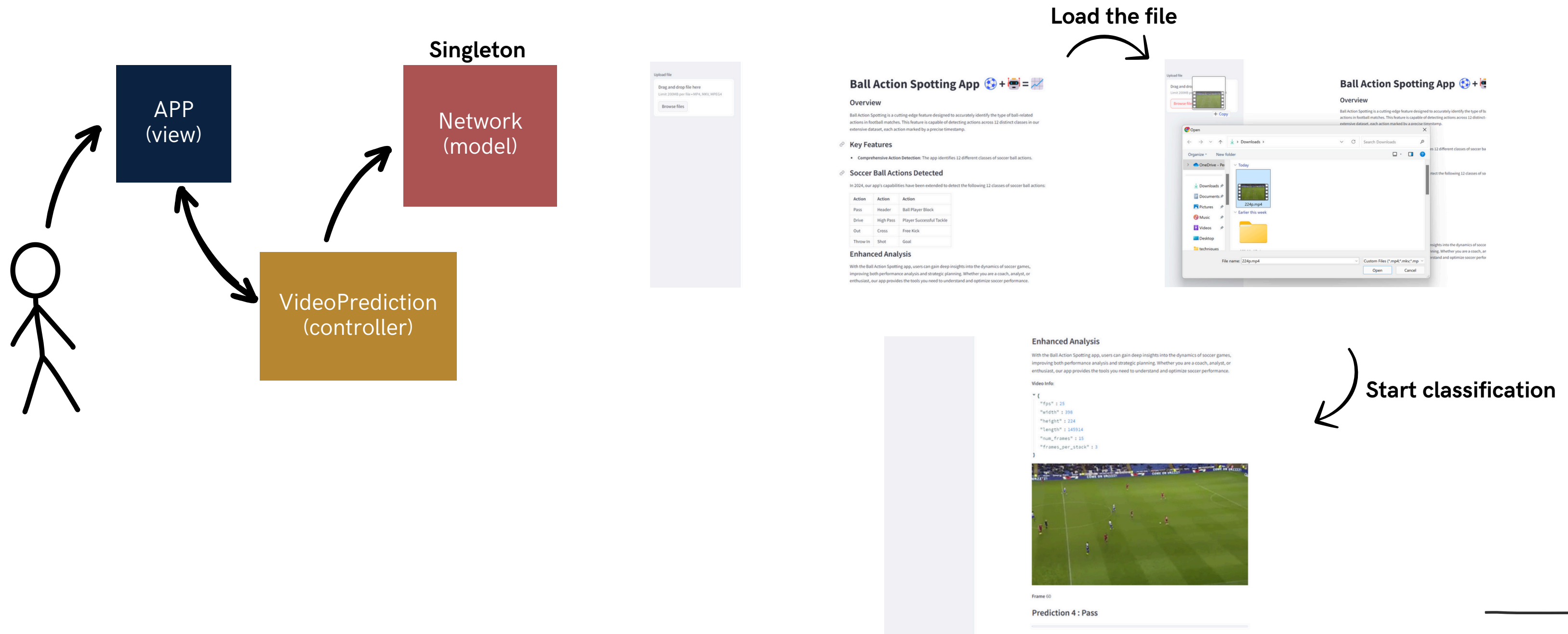
- accuracy macro : 0.44
- precision macro : 0.08
- recall macro : 0.09
- f1 macro: 0.07
- accuracy micro: 0.44
- precision micro: 0.44
- recall micro: 0.44
- f1 micro: 0.44

Evaluation NETWORK

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$
$$\text{Recall} = \frac{TP}{TP + FN}$$

$$\text{Precision} = \frac{TP}{TP + FP}$$
$$\text{F1 Score} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

# → Streamlit app





# → FUTURE WORKS

## Network

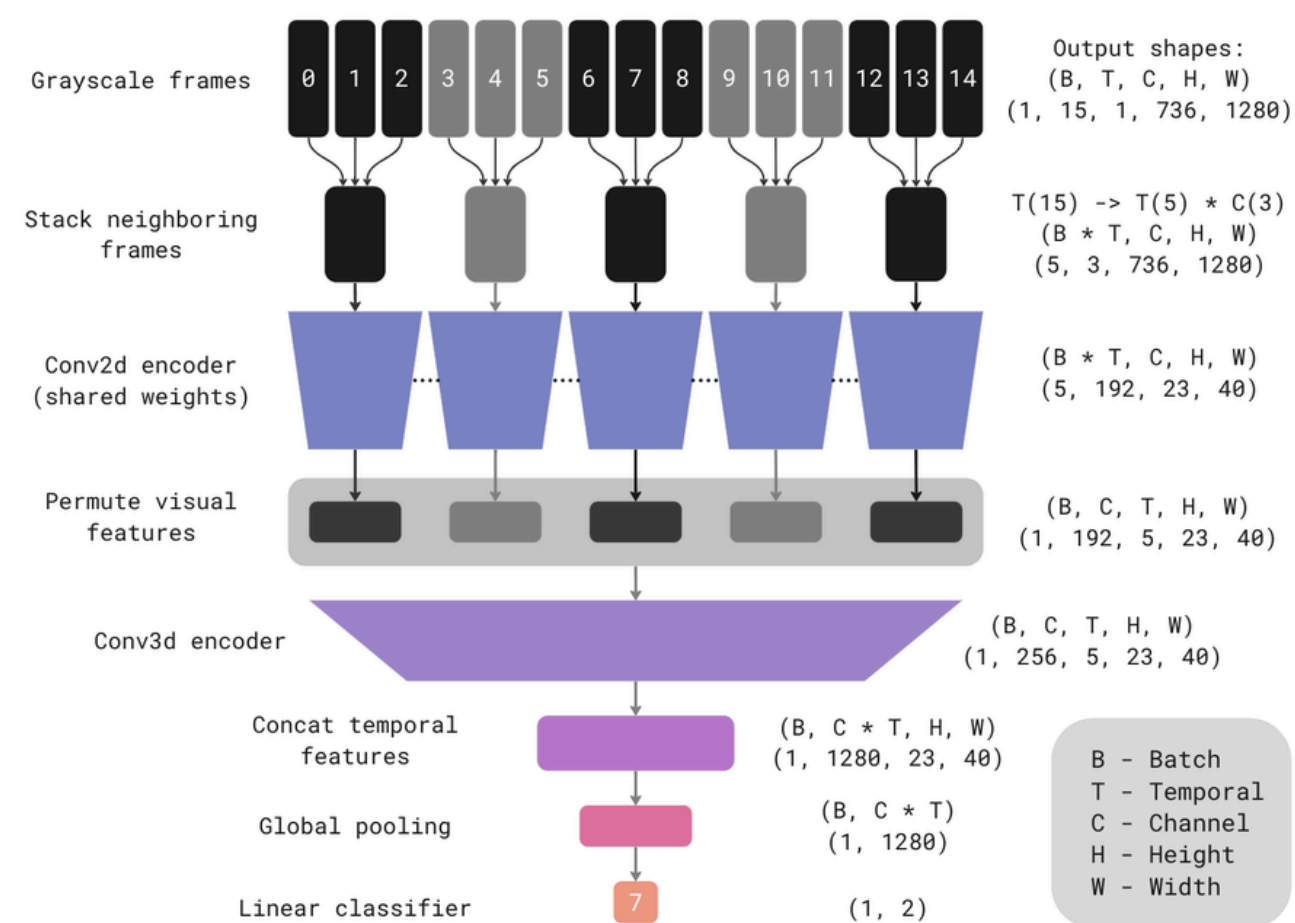
- Increase the dataset
- FineTune the model on the dataset of the task of action spotting
- Train for more epochs
- Try other hyperparameters
- Use high quality video
- Use RGB video
- Try other architectures

## Streamlit

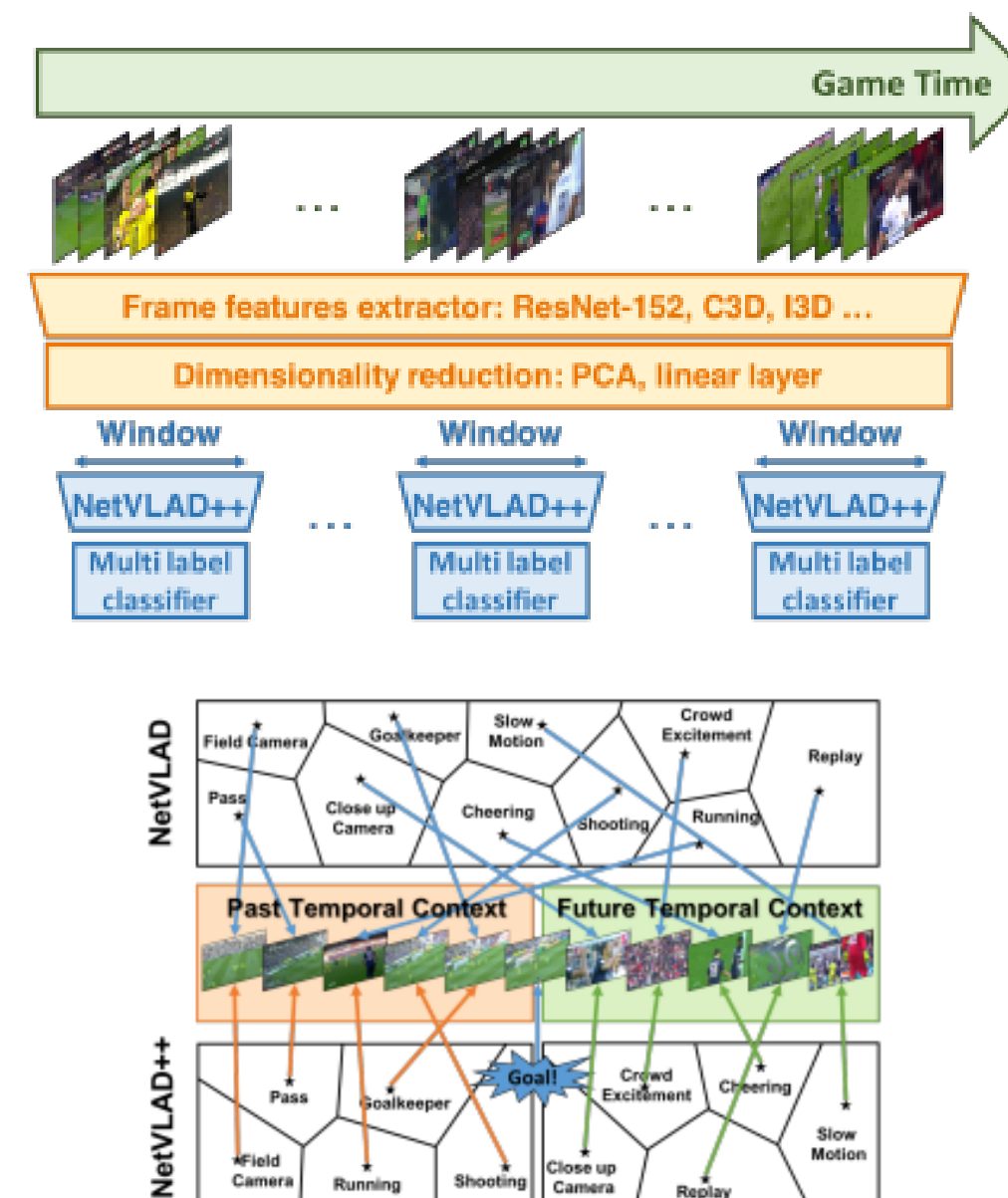
- Add a trigger to start the video when the classification begins.
- Use YOLO or another object detection model to assign the event to each player and compute some statistics relevant to the players.

# → State of art

## INSPIRATION



+



# State of art

## INSPIRATION-BIBLIOGRAPHY

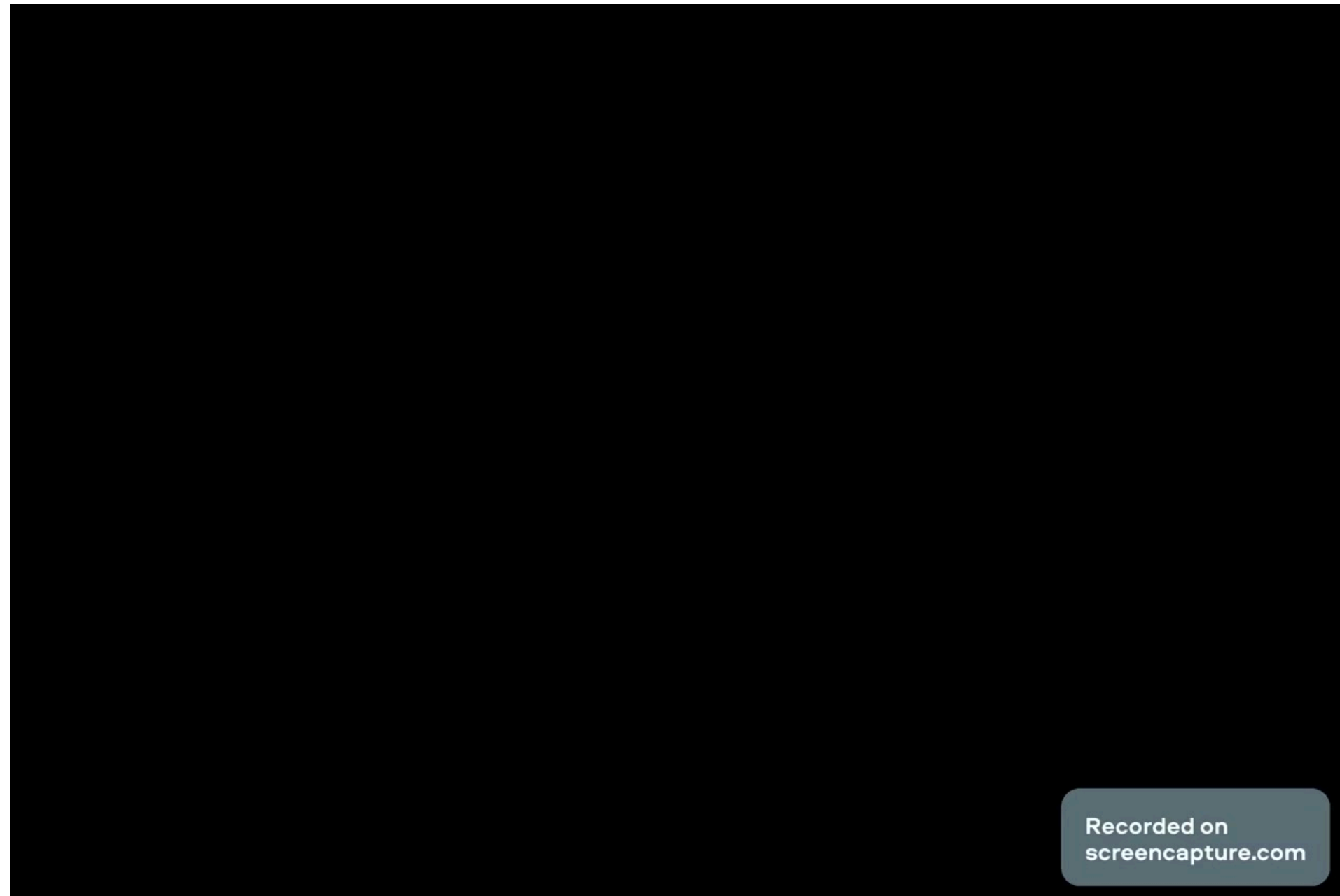
- NetVlad ++
- Encoder2D3DLateFusion
- EfficientNetV2
- SoccerNet
- NetVlad
- Vlad
- DatasetSoccerNet
- OverviewTrends
- InvertedResidualBlocks
- LibraryActionSpotting
- FocalLoss

# THANKS FOR ATTENTION

DANIELE CECCA MAT 914358

Submitted by Daniele Cecca

# DEMO



# DEMO

DANIELE CECCA MAT 914358

Advanced techniques  
2023-2024

Presentation

