# Action Recognition in Videos using Deep Learning

Deep Learning course

MSc in AI

Georgios Batsis

# About project

**Action Recognition:**
- Classifying the activity being performed by a human
- We need a set of evidence to recognize an action → Video classification

**Video:**
- A signal which combined spatial and temporal information
- Sequence of images-frames

# MODELS

# Baseline

Step 1 → Feature Extraction of each frame: Texture information → GLCM features:

- o Contrast
- o Dissimilarity
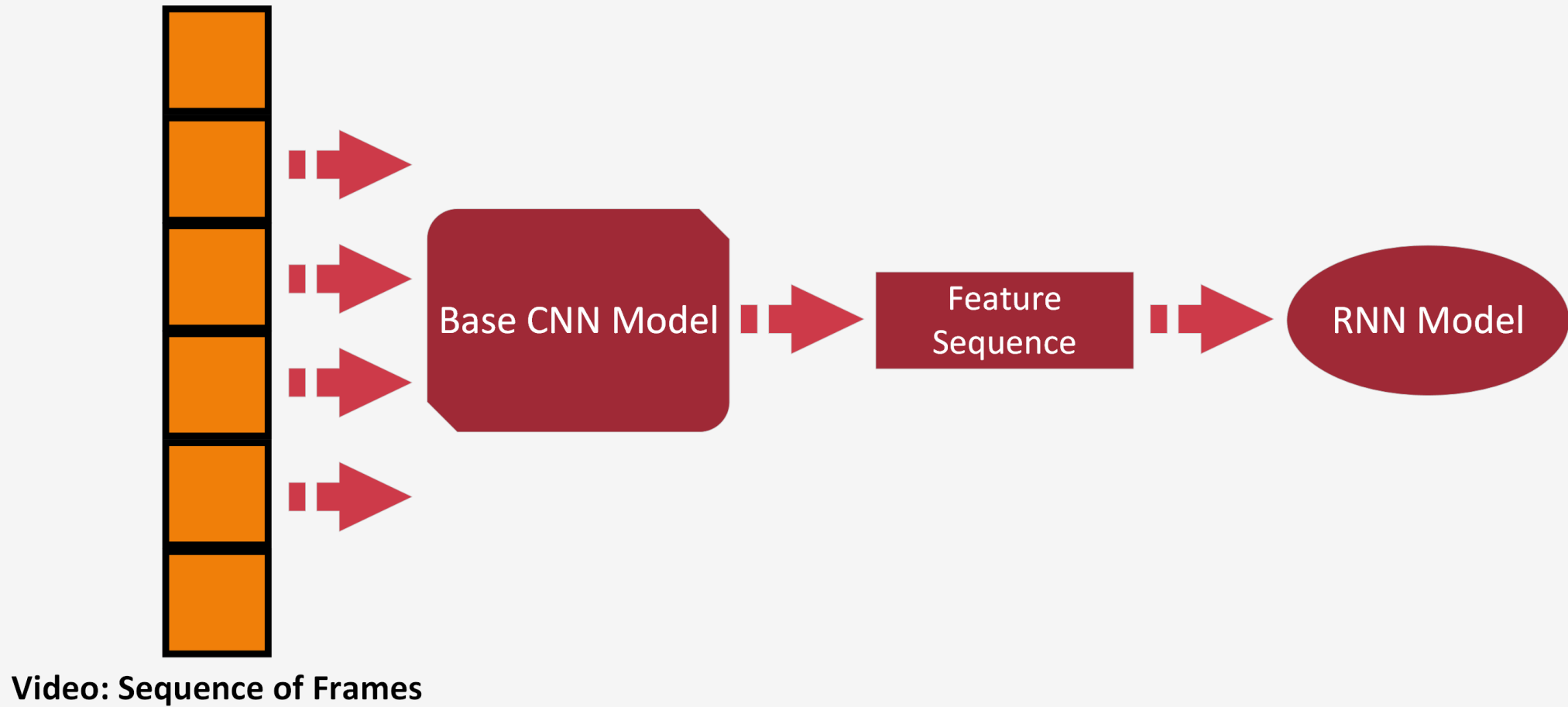- o Homogeneity
- o ASM
- o Energy
- o Correlation

Step 2 → Temporal Aggregation: Statistics → Mean & Std

Step 3 → Definition of a classifier pipeline:

- o Standard Scaler
- o SVM with RBF Kernel

# Going Deeper...



**Video: Sequence of Frames**

# Experiments: Step 1 → Choose Base CNN Model

**Most popular pretrained models:**

- VGG

- ResNet

- MobileNet.

**Choose which layers will be left frozen:**

*CNNs consist of:*

- o Convolutional Block → Convolutional and Pooling Layers

- o Classifier → Fully Connected Layers

*Freezing options:*

- o Last (or other) Convolutional Layer

- o Penultimate Layer of classifier

# Experiments: Step 2 ➔ Deep Model Construction
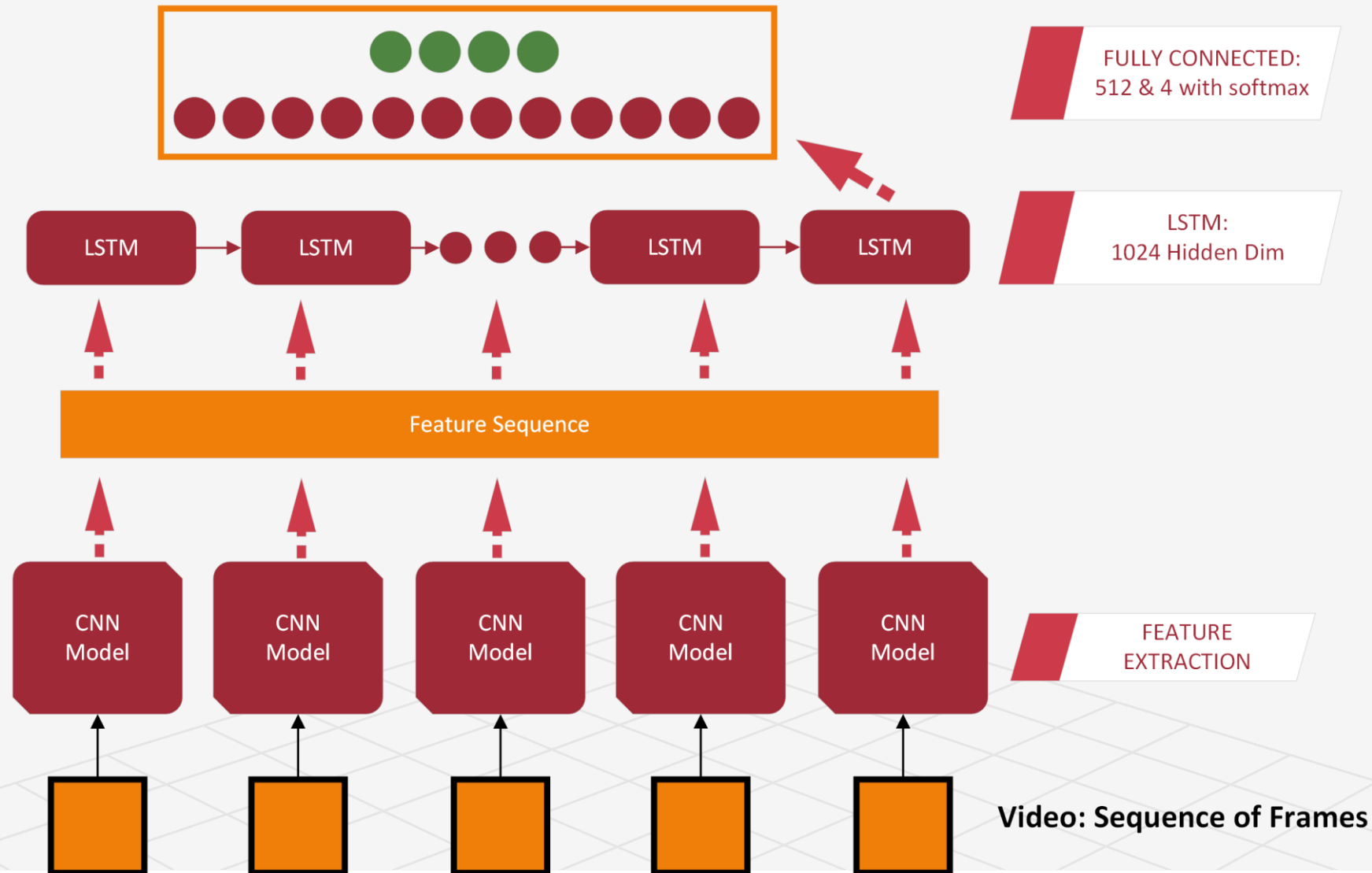
**After pre-trained CNN selection, decide if:**

- Use the extracted features directly

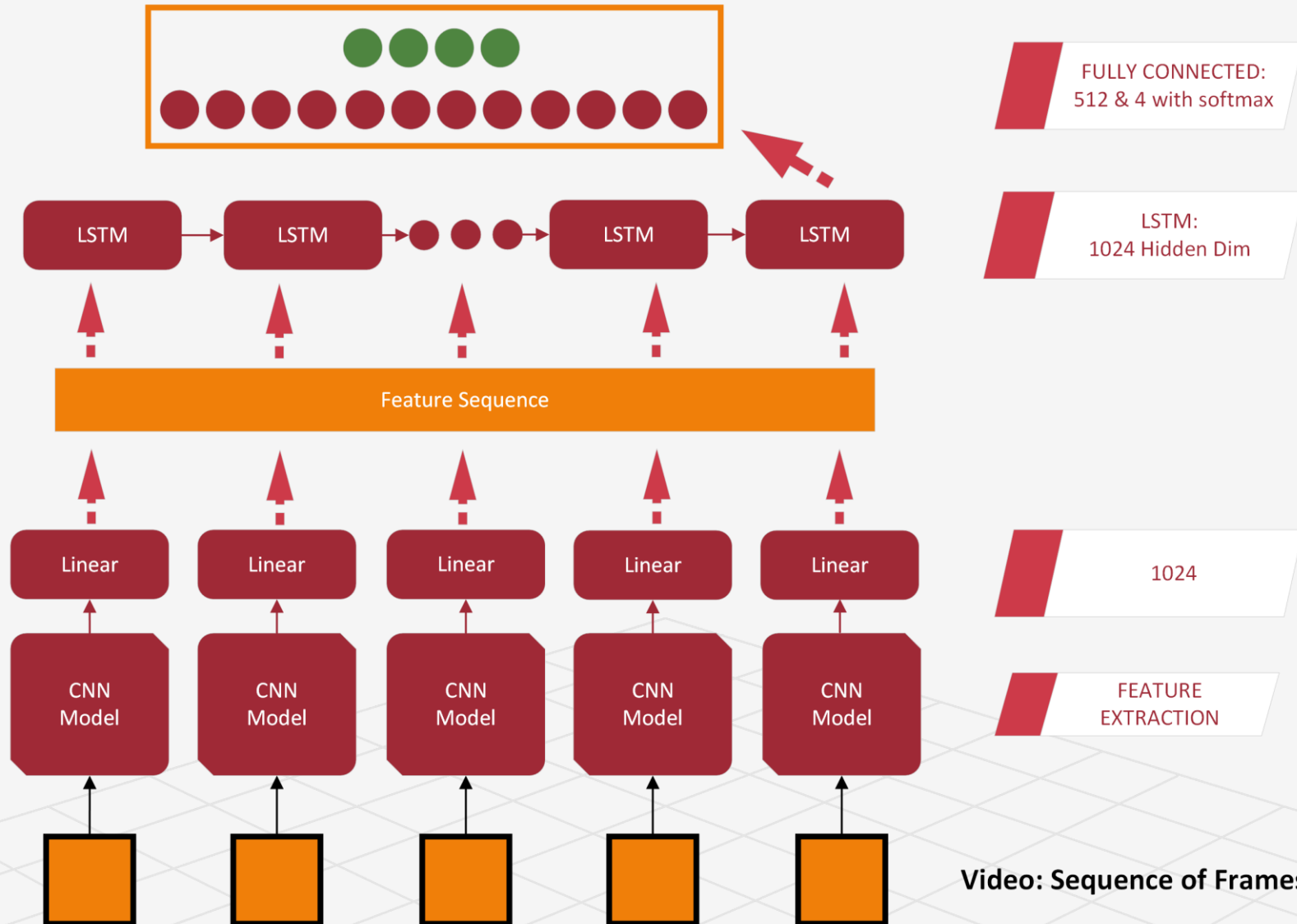- Add trainable layer(s) after feature extraction

**Define the RNN-based part:**

- Model (e.g. LSTM)

  o Hidden state dimension

  o Number of layers

  o Direction

- Final classifier

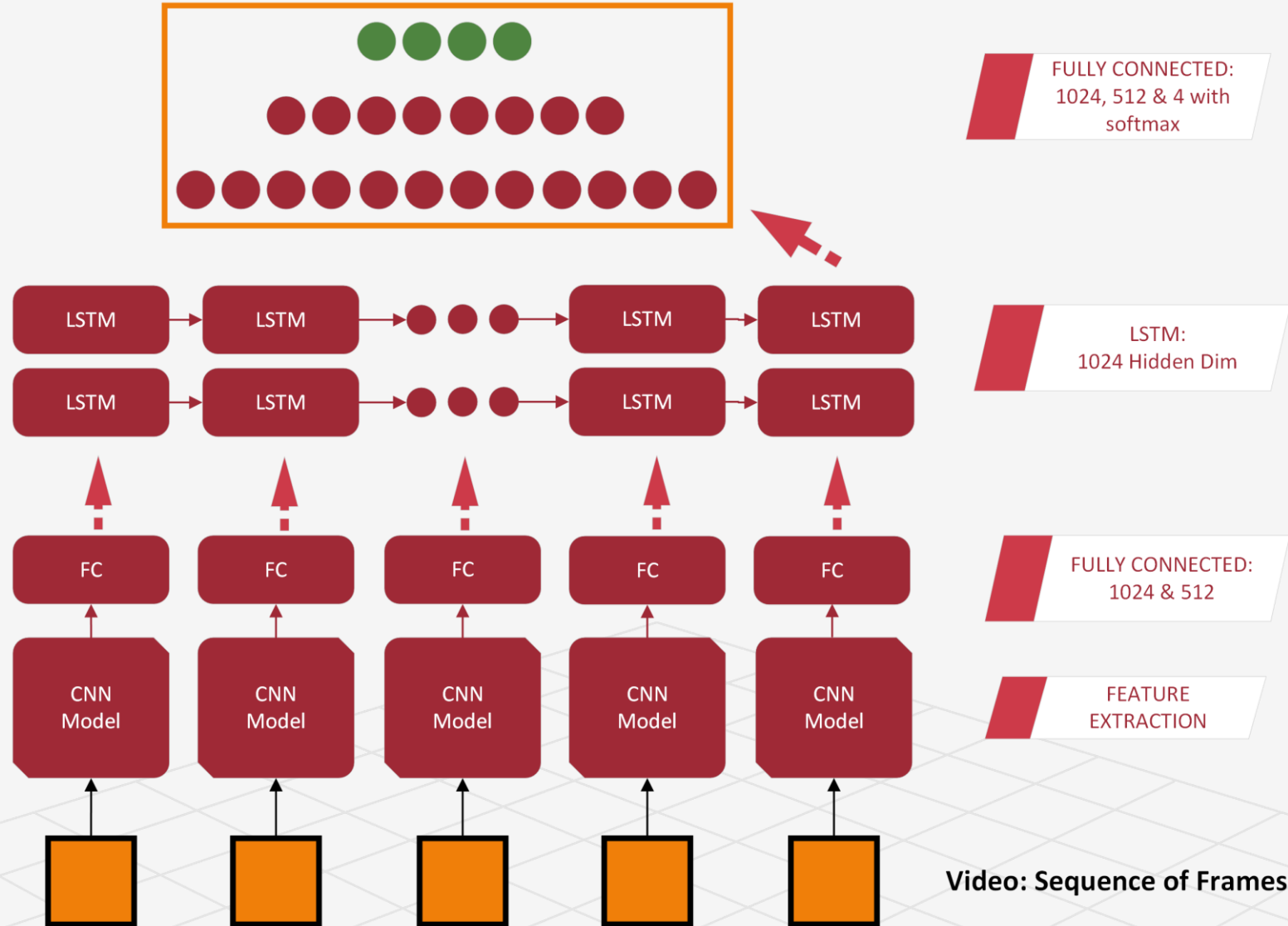# Approach A



FULLY CONNECTED:
512 & 4 with softmax

LSTM:
1024 Hidden Dim

FEATURE
EXTRACTION

LSTM

LSTM

LSTM

LSTM

Feature Sequence

CNN Model

CNN Model

CNN Model

CNN Model

CNN Model

Video: Sequence of Frames

# Approach B



FULLY CONNECTED:
512 & 4 with softmax

LSTM:
1024 Hidden Dim

1024

FEATURE
EXTRACTION

**Video: Sequence of Frames**

# Approach C



FULLY CONNECTED:
1024, 512 & 4 with softmax

LSTM:
1024 Hidden Dim

FULLY CONNECTED:
1024 & 512

FEATURE EXTRACTION

**Video: Sequence of Frames**

# Approach C'



FULLY CONNECTED:
**1024 * 2**, 512 & 4
with softmax

LSTM: 2 Bidirectional
Layers with
1024 Hidden Dim

FULLY CONNECTED:
1024 & 512

FEATURE
EXTRACTION

**Video: Sequence of Frames**

# Approach D



FULLY CONNECTED:
1024 * **2**, 512 & 4
with softmax

LSTM: 3 Bidirectional
Layers with
1024 Hidden Dim

FULLY CONNECTED:
1024 - 512 - 256

FEATURE
EXTRACTION

**Video: Sequence of Frames**
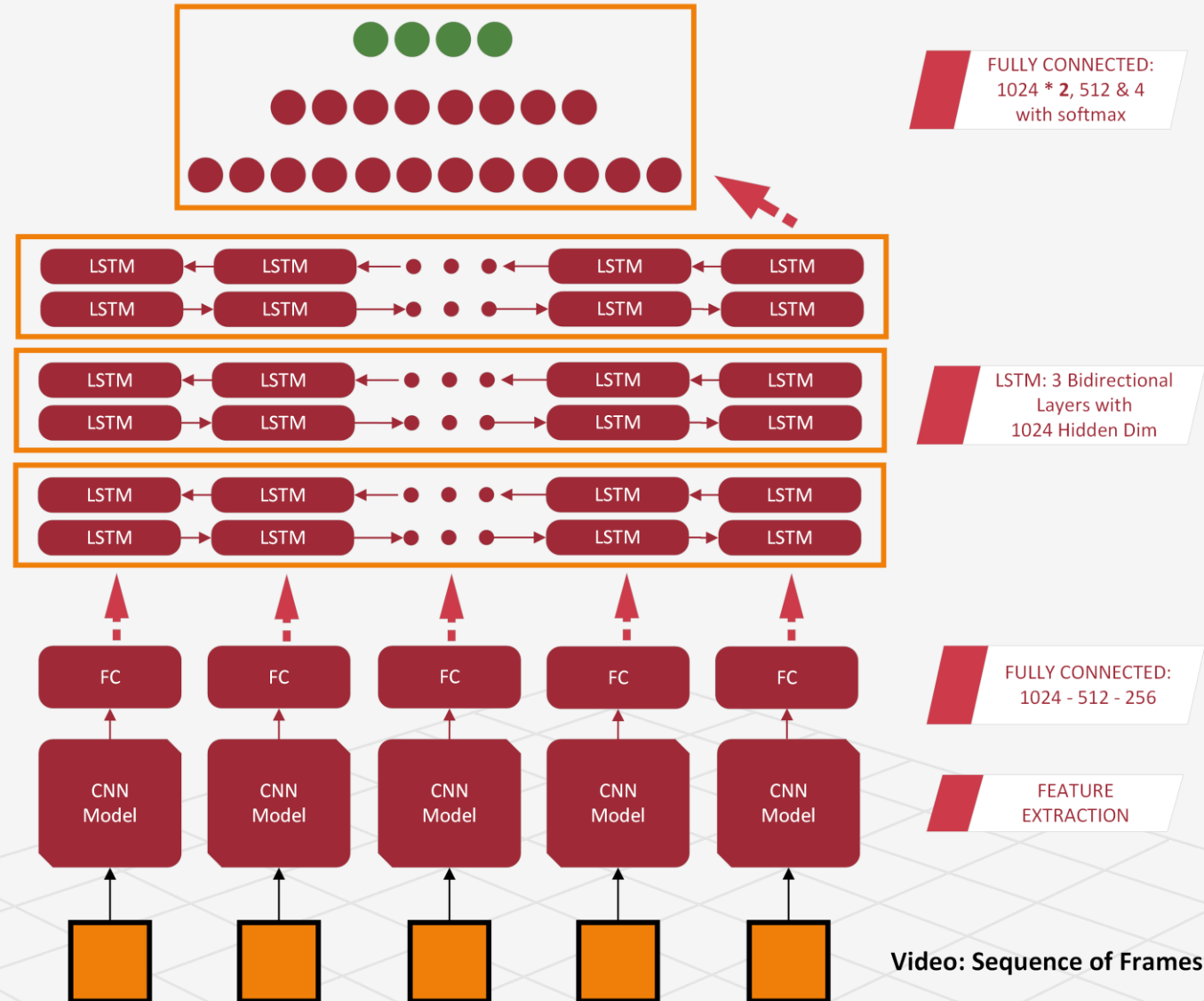
# RESULTS

# Dataset

## UCF101 Human Actions dataset:

A small subset was used in this project

Official train-test splitting

Classes:

0: Playing Guitar

1: Rock Climbing Indoor

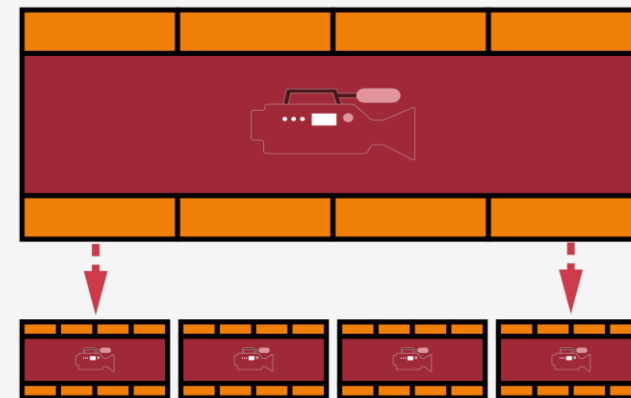2: Soccer Juggling

3: Band Marching

# Preprocessing

## Sequence level transformation:

➢ Problem: Differences in the total number of frames per video & the Fps

➢ Solution: Extract video segments ➔ Shorter videos with fixed number of frames

## Frame-level transformation:

➢ Resize (e.g. 224x224)

➢ Normalize ➔ *mean=[0.485, 0.456, 0.406], std=[0.229, 0.224, 0.225]*

# Results: Baseline

| True\Predicted | 0 | 1 | 2 | 4 |
|---|---|---|---|---|
| 0 | 54 | 28 | 82 | 4 |
| 1 | 17 | 63 | 12 | 0 |
| 2 | 59 | 8 | 135 | 0 |
| 3 | 4 | 0 | 18 | 112 |

```
               precision    recall  f1-score   support

           0        0.41      0.32      0.36       168
           1        0.64      0.70      0.67        92
           2        0.55      0.67      0.60       202
           3        0.97      0.84      0.90       134

    accuracy                            0.61       596
   macro avg        0.64      0.63      0.63       596
weighted avg        0.62      0.61      0.61       596
```

# Results

## Pre-trained CNN selection:
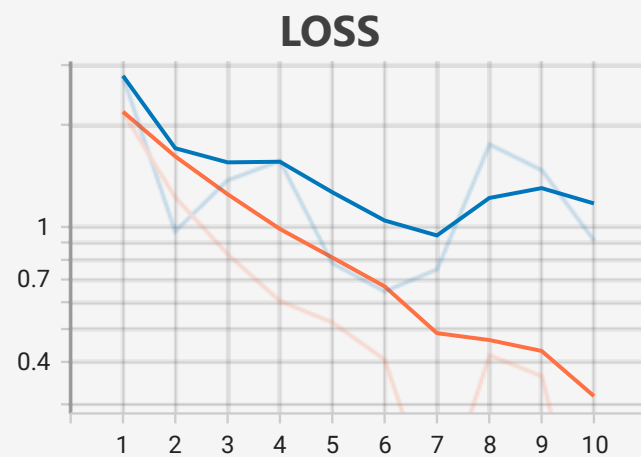
*All models tested using the DNN of **Approach A***

➤ ResNet – 125 freezed at the last Convolutional Layer (we receive output of Average Pooling) outperforms VGG 19

   freezed  at the penultimate Layer of classifier [87.8% vs 83.4% F1 Score in validation dataset].

➤ If we freeze the last Convolutional Layer of VGG the model can't learn…

➤ MobileNet overfits…

# Results

| Model Name | F1 Score | Loss |
|---|---|---|
| Approach A | 87.8 | 2.9 |
| Approach B | 97.7 | 0.64 |
| Approach C | 98.4 | 0.44 |
| Approach C (Bid) | 97.7 | 0.57 |
| Approach D | 98.6 | 0.29 |

# Results



**LOSS**

| Model Name | F1 Score | Loss |
| --- | --- | --- |
| Approach A | 87.8 | 2.9 |
| Approach B | 97.7 | 0.64 |
| Approach C | 98.4 | 0.44 |
| Approach C (Bid) | 97.7 | 0.57 |
| Approach D | 98.6 | 0.29 |

# Results

## LOSS



## F1



| Model Name | F1 Score | Loss |
|---|---|---|
| Approach A | 87.8 | 2.9 |
| Approach B | 97.7 | 0.64 |
| Approach C | 98.4 | 0.44 |
| Approach C (Bid) | 97.7 | 0.57 |
| Approach D | 98.6 | 0.29 |

# Results

## LOSS



## F1



| Model Name | F1 Score | Loss |
|---|---|---|
| Approach A | 87.8 | 2.9 |
| Approach B | 97.7 | 0.64 |
| Approach C | 98.4 | 0.44 |
| Approach C (Bid) | 97.7 | 0.57 |
| Approach D | 98.6 | 0.29 |

# Results

**LOSS**



**F1**



| Model Name | F1 Score | Loss |
|---|---|---|
| Approach A | 87.8 | 2.9 |
| Approach B | 97.7 | 0.64 |
| Approach C | 98.4 | 0.44 |
| Approach C (Bid) | 97.7 | 0.57 |
| Approach D | 98.6 | 0.29 |

# Results: Approach C & Test Dataset

| True\Predicted | 0 | 1 | 2 | 4 |
|---|---|---|---|---|
| 0 | 43 | 0 | 0 | 0 |
| 1 | 0 | 39 | 2 | 0 |
| 2 | 0 | 0 | 38 | 1 |
| 3 | 0 | 0 | 0 | 43 |

```
              precision    recall  f1-score   support

           0       1.00      1.00      1.00        43
           1       1.00      0.95      0.97        41
           2       0.95      0.97      0.96        39
           3       0.98      1.00      0.99        43

    accuracy                           0.98       166
   macro avg       0.98      0.98      0.98       166
weighted avg       0.98      0.98      0.98       166
```



FULLY CONNECTED:
1024, 512 & 4 with softmax

LSTM:
1024 Hidden Dim

FULLY CONNECTED:
1024 & 512

FEATURE EXTRACTION

Video: Sequence of Frames

Video level classification → Combine predictions of all segments.

# DEMO

# Thank you!