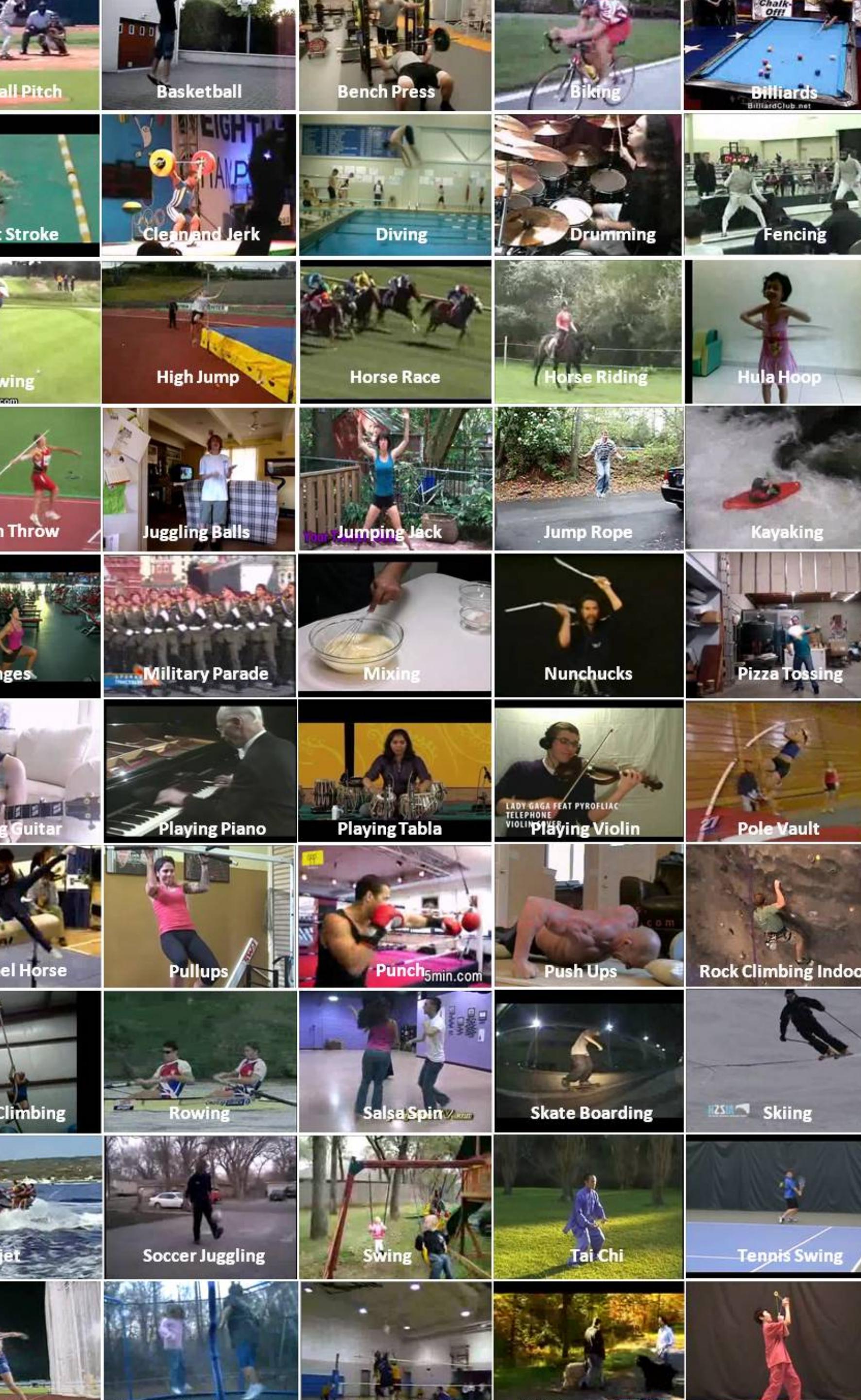


Enhancing Action Recognition & Advanced Frames Selection

Innopolis University, Computer Vision Course 2024



Project Overview

Problem Statement:

Video action recognition is computationally demanding and resource-intensive.

Objective:

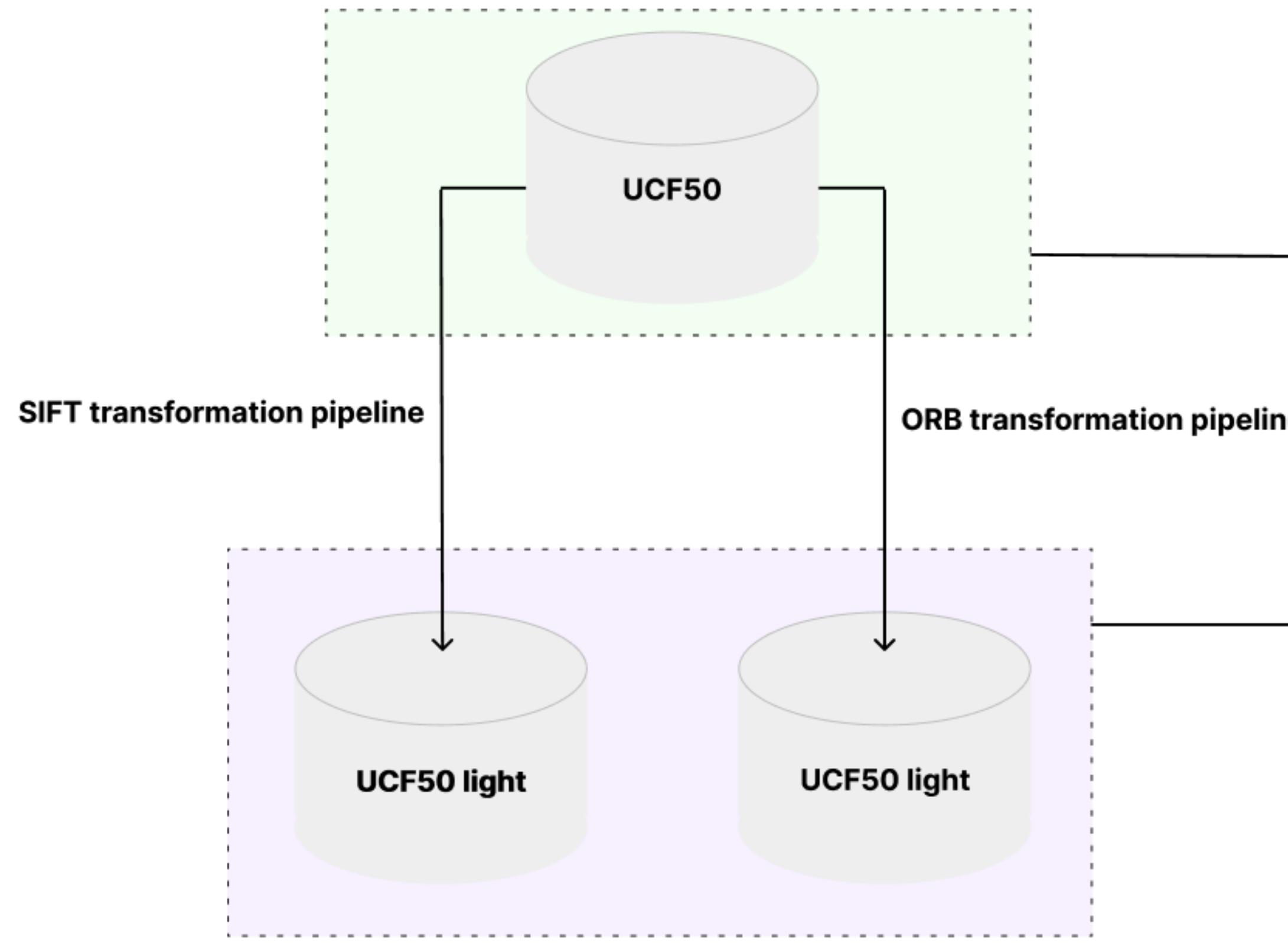
Develop a lightweight video action recognition pipeline that reduces data size while maintaining competitive accuracy.

Key Performance Indicators:

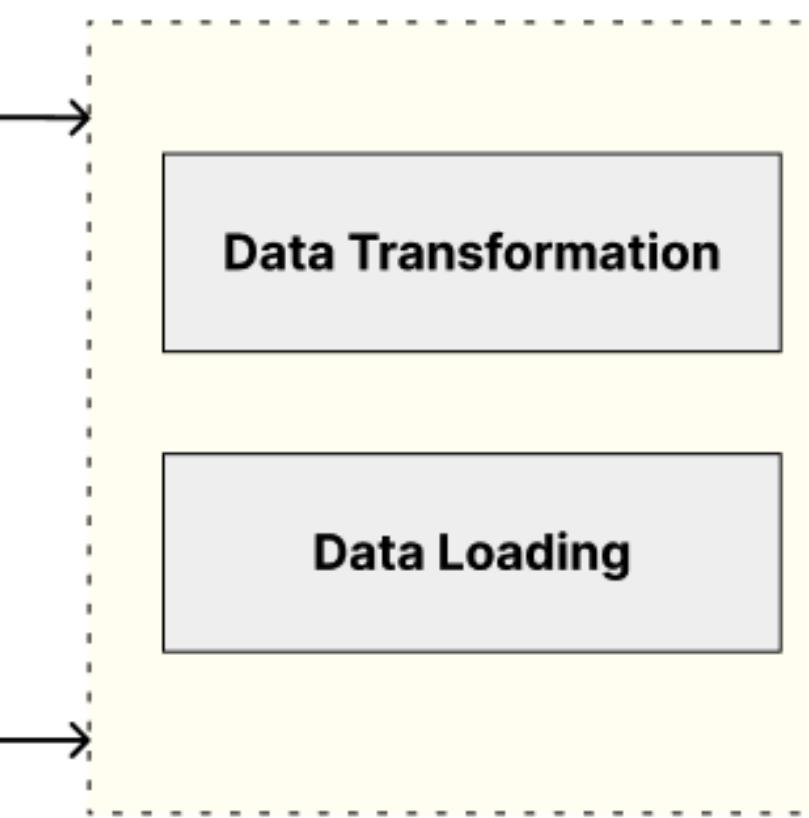
Reduction in data size, training and inference time, and maintenance of classification metrics (F1, accuracy, precision, recall, ROC-AUC).

Project architecture

Data Source



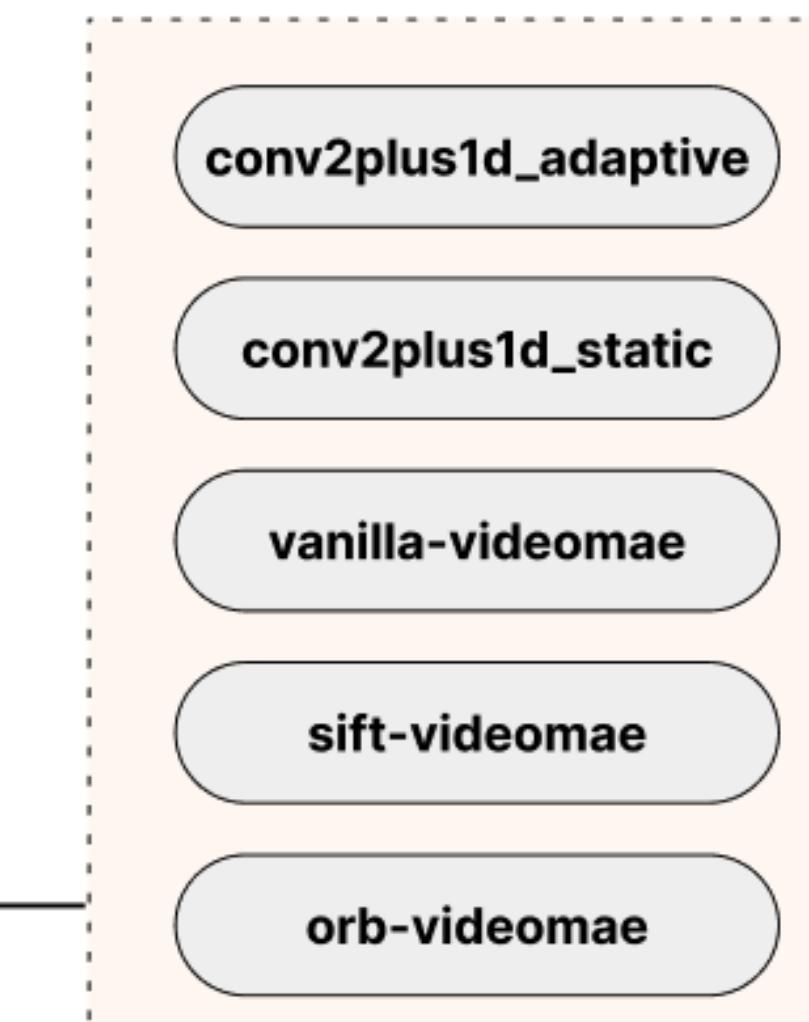
Experiment Setup



Train and Evaluation



Models store

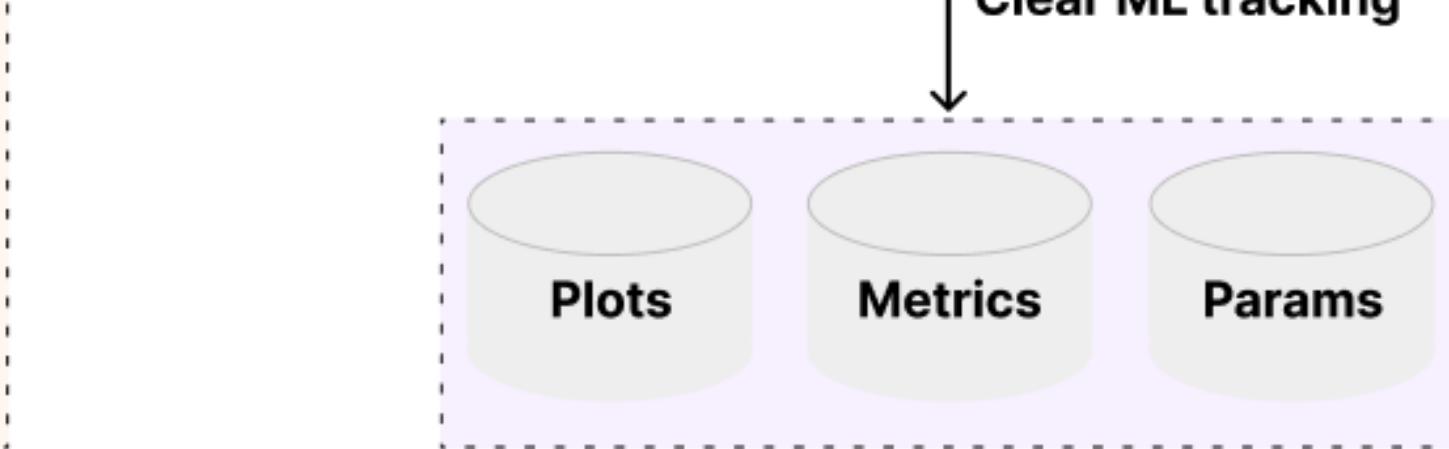
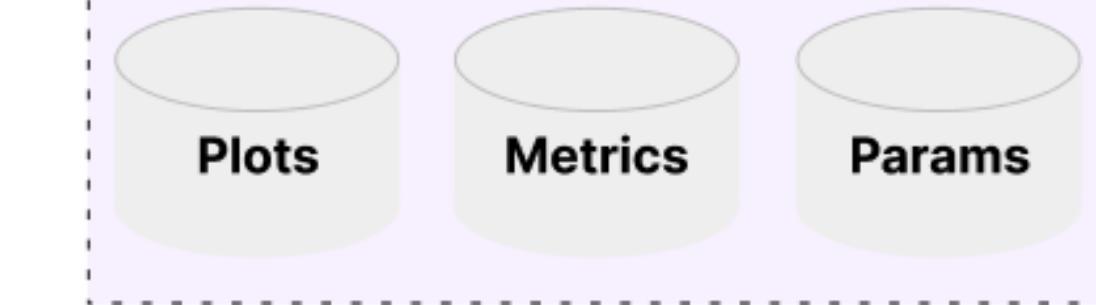


Data Artifacts Store



Deployment

Experiments Store



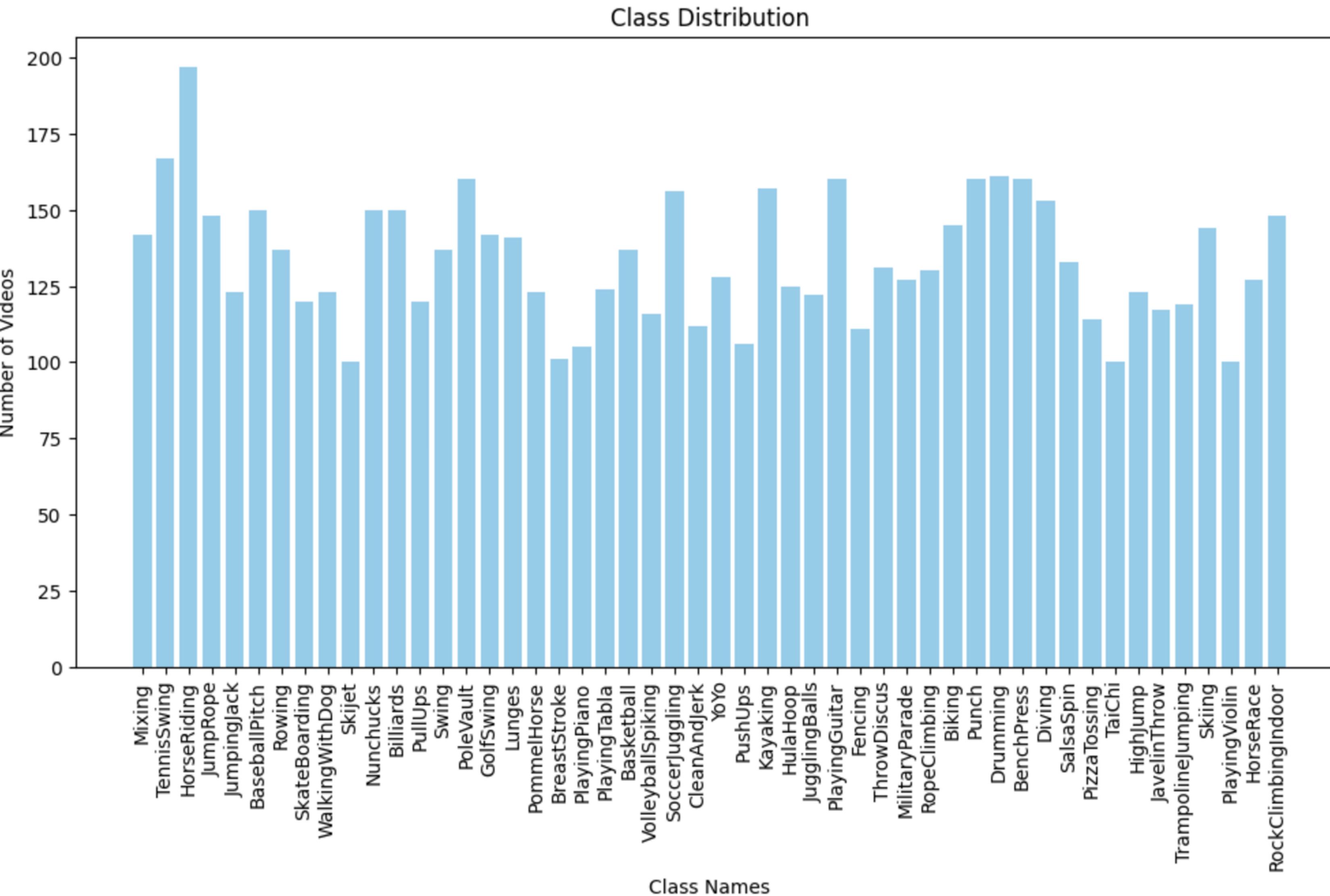
Research design

UCF 50 Dataset

6,618 clips

50 categories

3.23 GB size

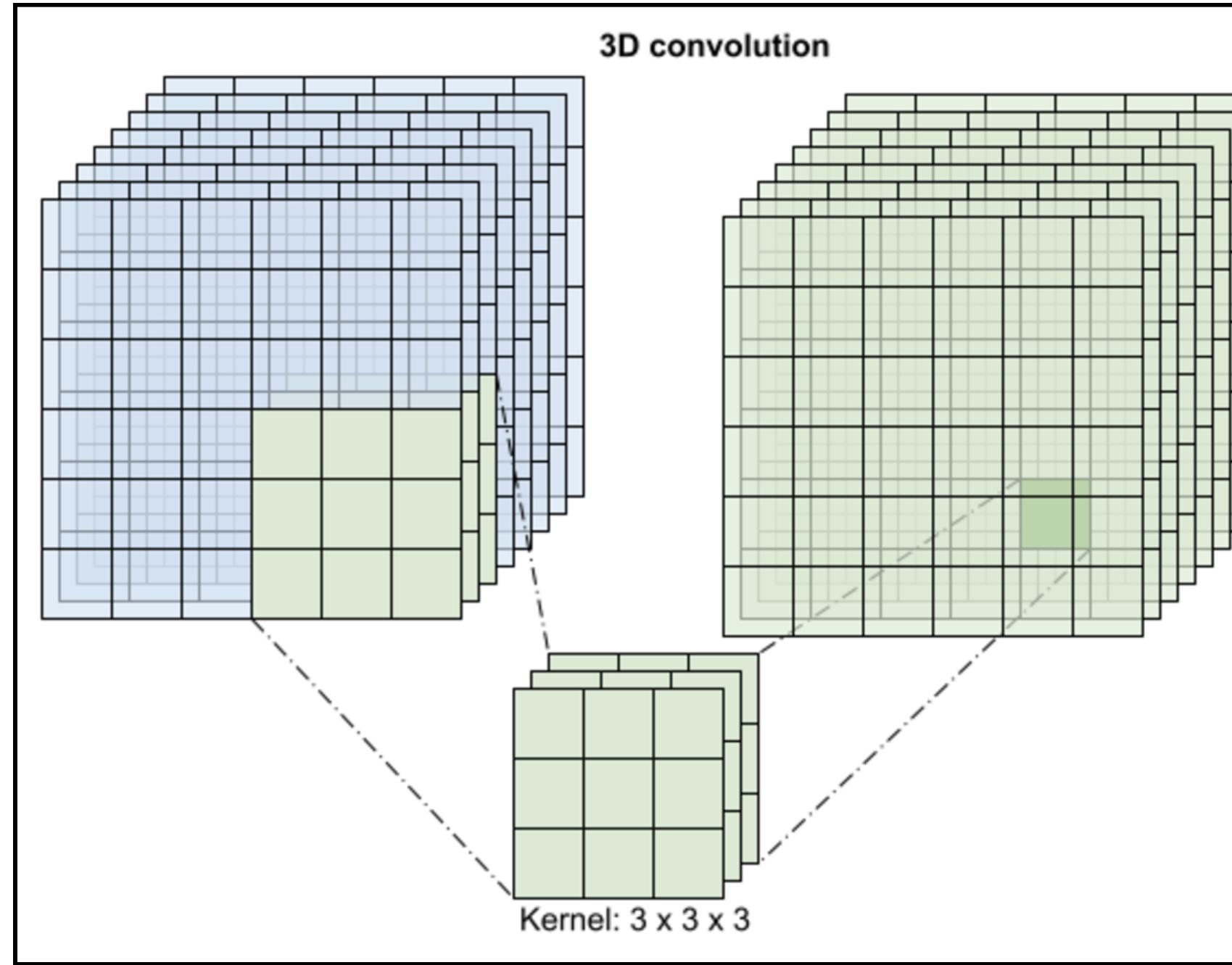


(2+1)D Convolutions

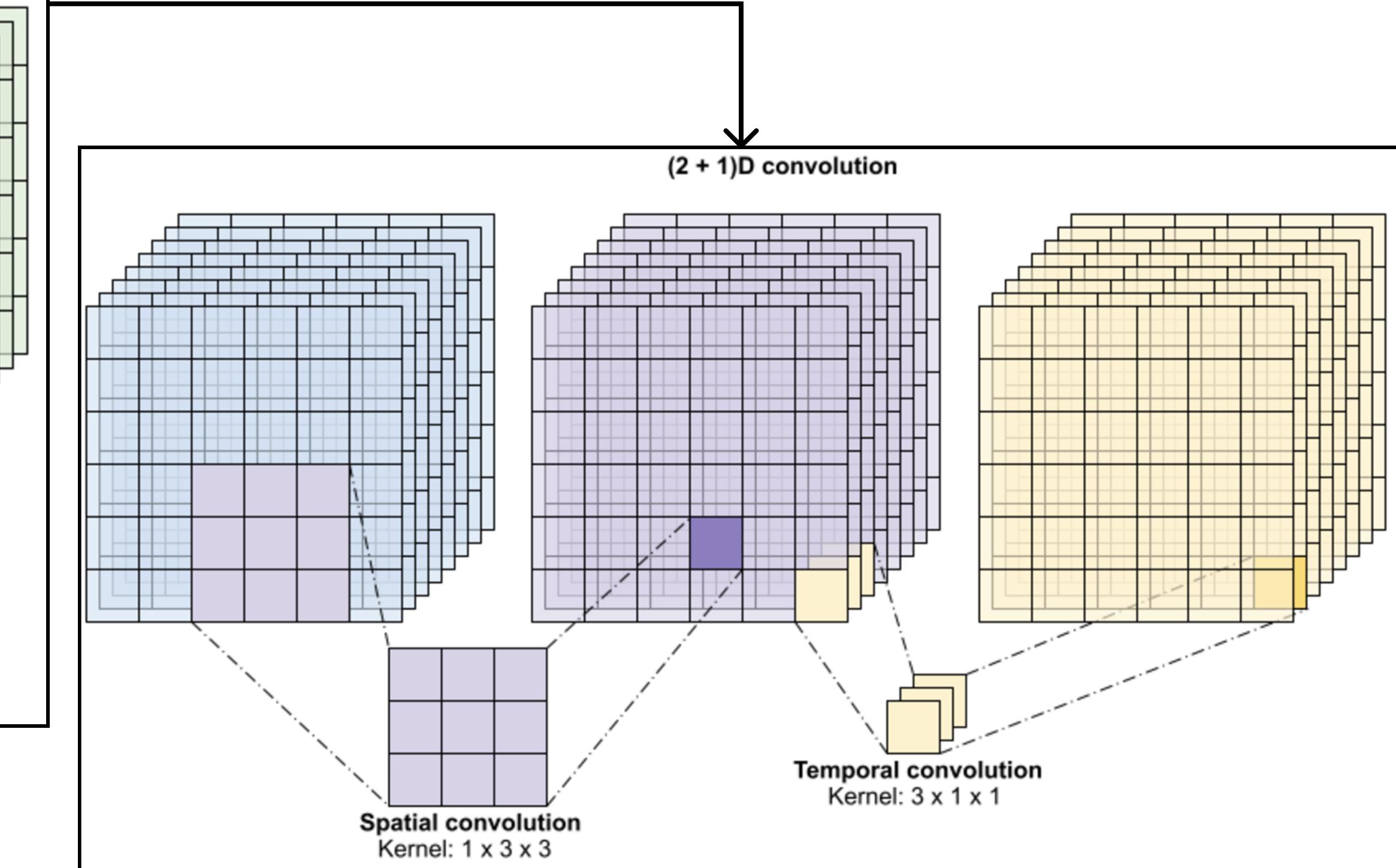
Basic idea

Models

param#3D = $(3 \times 3 \times 3) \times \text{in} \times \text{out} + \text{out}$



$$\begin{aligned} \text{param\#(2+1)D} = & (1 \times 3 \times 3) \times \text{in} \times \text{out} + \text{out} + \\ & + (3 \times 1 \times 1) \times \text{out} \times \text{out} + \text{out} \end{aligned}$$



(2+1)D Convolutions

Results

Models

Static Method: $\text{step_frame} = \text{const}$

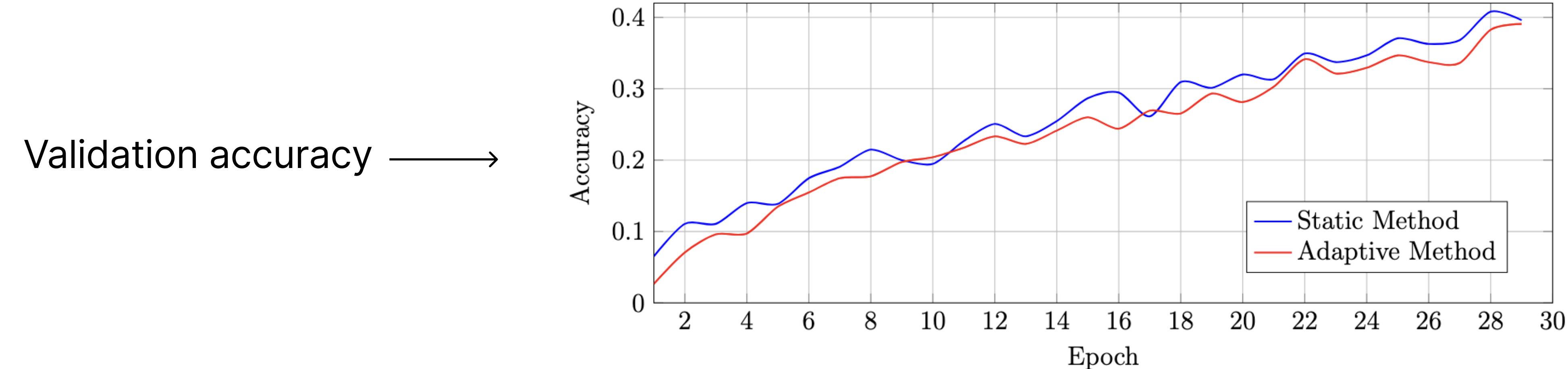
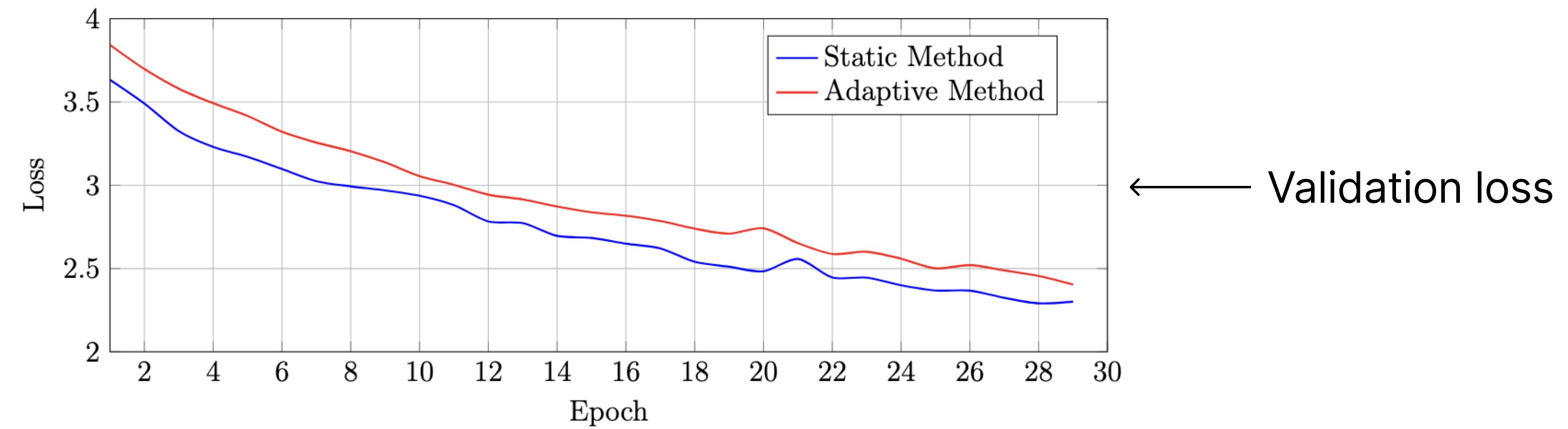
Adaptive Method: $\text{step_frame} = \text{video_length/n_frames}$

Method	Test Loss	Test Accuracy	Macro-average ROC-AUC
Static Method	2.288	0.376	0.93
Adaptive Method	2.401	0.364	0.92

(2+1)D Convolutions

Results

Models



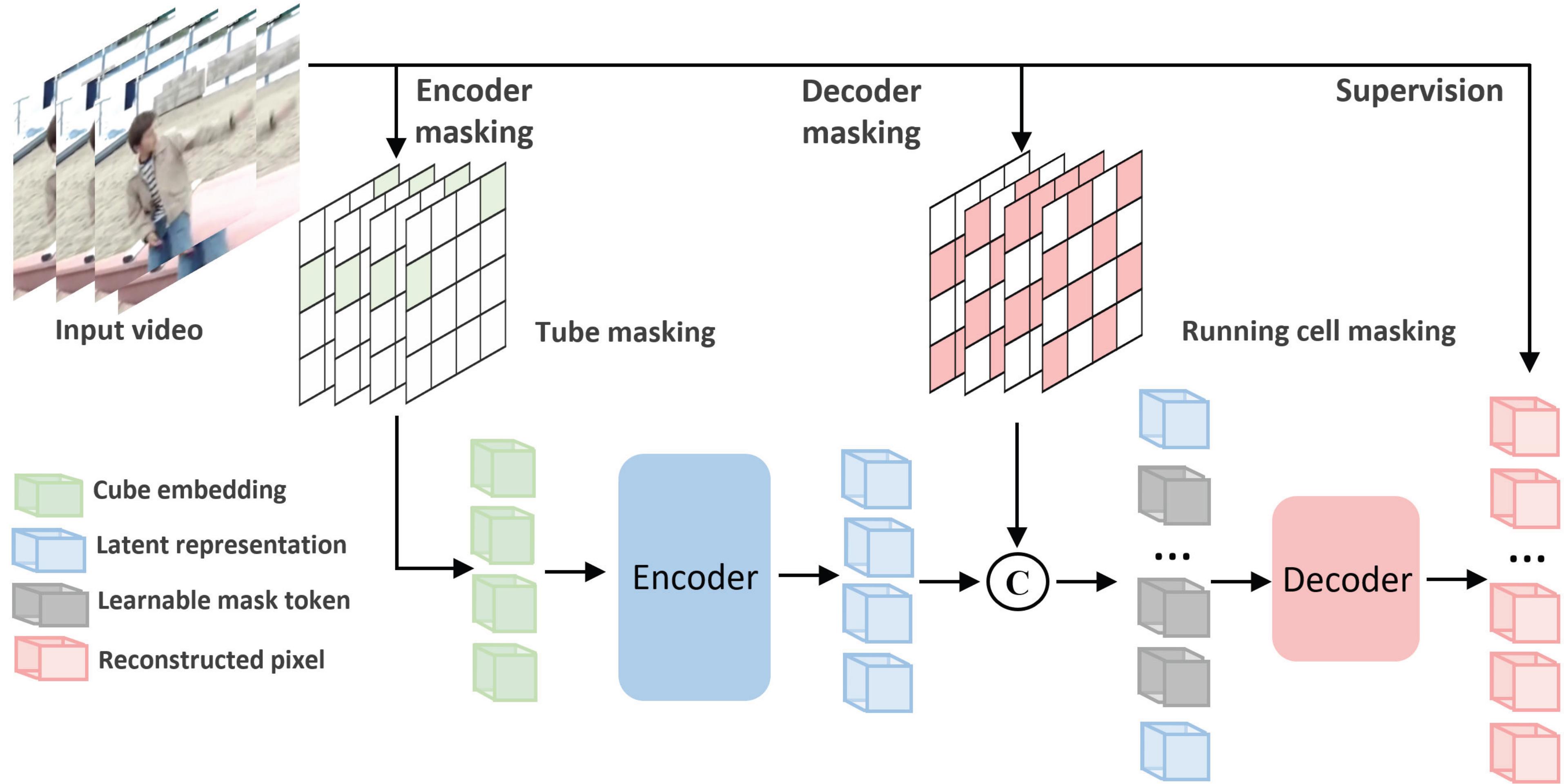
Validation accuracy →

← Validation loss

VideoMAE

Basic idea

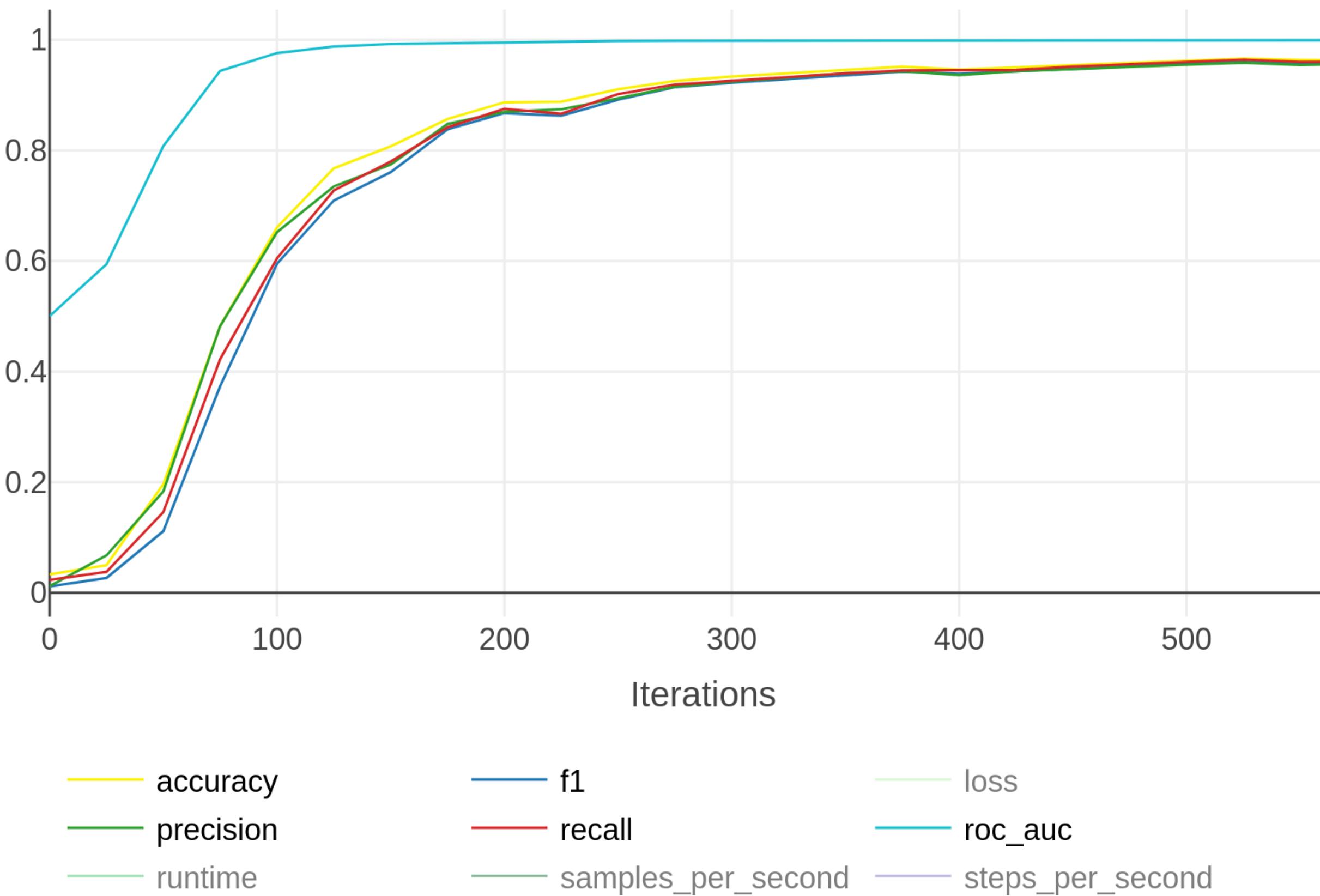
Models



VideoMAE

Results

Models



Models

Model Comparison

	Parameters	Size (.safetensors)	Pretrained	Framework
(2+1)D Convolutions	0.4M	1.8MB	-	
VideoMAE	94.2M	330MB	Kinetics-400	

Frames Selection

ORB

(Oriented FAST and Rotated BRIEF)

Identify keypoints and compute their descriptors

```
keypoints, descriptors = orb.detectAndCompute(processed_frame, None)
```

Compare descriptors
using Hamming distance

```
bf = cv2.BFMatcher(cv2.NORM_HAMMING, crossCheck=True)  
matches = bf.match(prev_descriptors, descriptors)
```

Calculate distances

```
distances = [m.distance for m in matches]  
avg_distance = np.mean(distances)
```

Select frames based on a distance threshold

```
if avg_distance > threshold:  
    relevant_frames.append(frame)
```

SIFT

Scale-Invariant Feature Transform

Frames Selection

Identify keypoints and compute their descriptors

```
keypoints, descriptors = sift.detectAndCompute(processed_frame, None)
```

Find potential matches

```
bf = cv2.BFMatcher()
matches = bf.knnMatch(prev_descriptors, descriptors, k=2)
```

Apply Lowe's ratio test to filter matches

```
good_matches = []
for match in matches:
    if len(match) == 2:
        m, n = match
        if m.distance < 0.75 * n.distance:
            good_matches.append(m)
```

Maintain a record of matches

```
total_frames += 1
all_frames.append((total_frames, len(good_matches), processed_frame))
```

SIFT

Scale-Invariant Feature Transform

Frames Selection

Sort all frames by the number of good matches (ascending)

```
all_frames.sort(key=lambda x: x[1])
```

Choose the top frames

```
selected_frames = all_frames[:min(len(all_frames), 2 * number_of_frames)]
```

Reorder the selected frames to match their original sequence

```
selected_frames.sort(key=lambda x: x[0])
```

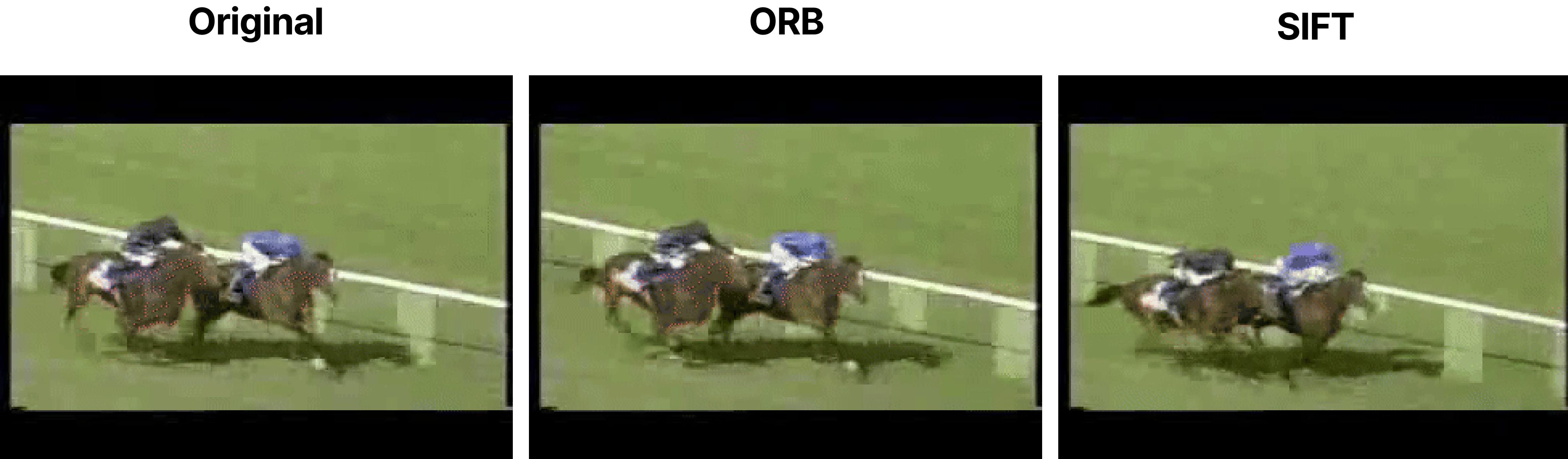
Save only the contents of the frame

```
relevant_frames = [frame[2] for frame in selected_frames]
```

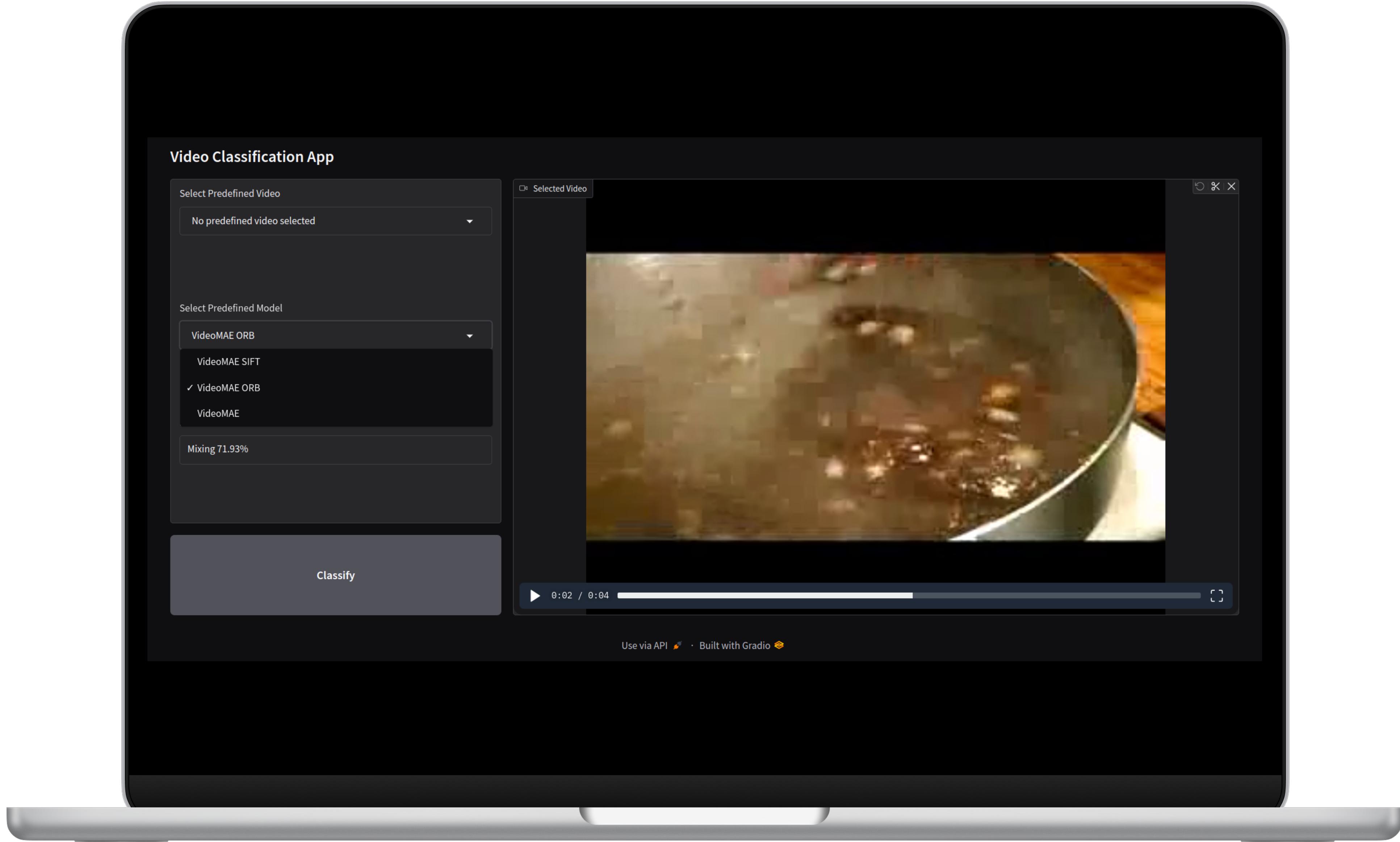
Dataset Size Comparison

Results

Dataset	UCF50	UCF50 ORB	UCF50 SIFT
Size of zip file	3.23 GB	0,87 GB	1.62 GB

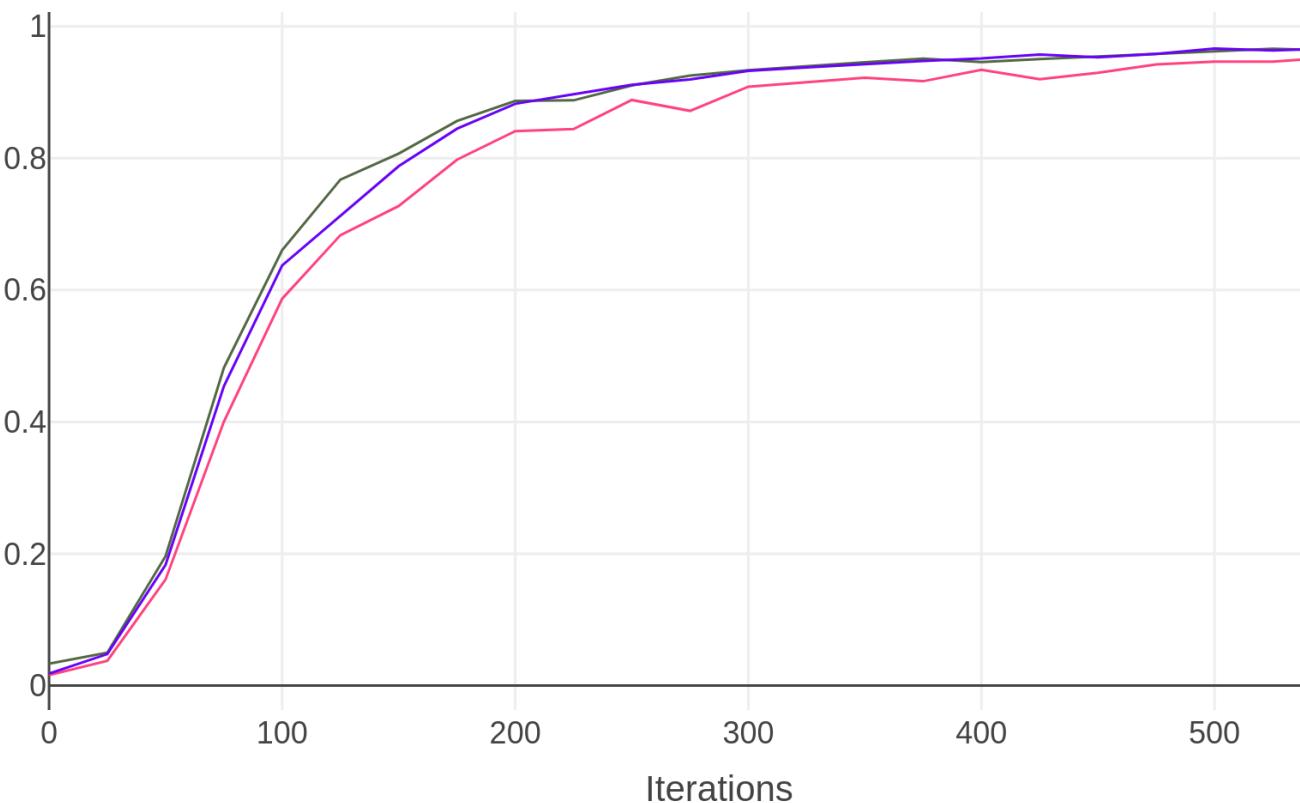


Let's see the demo

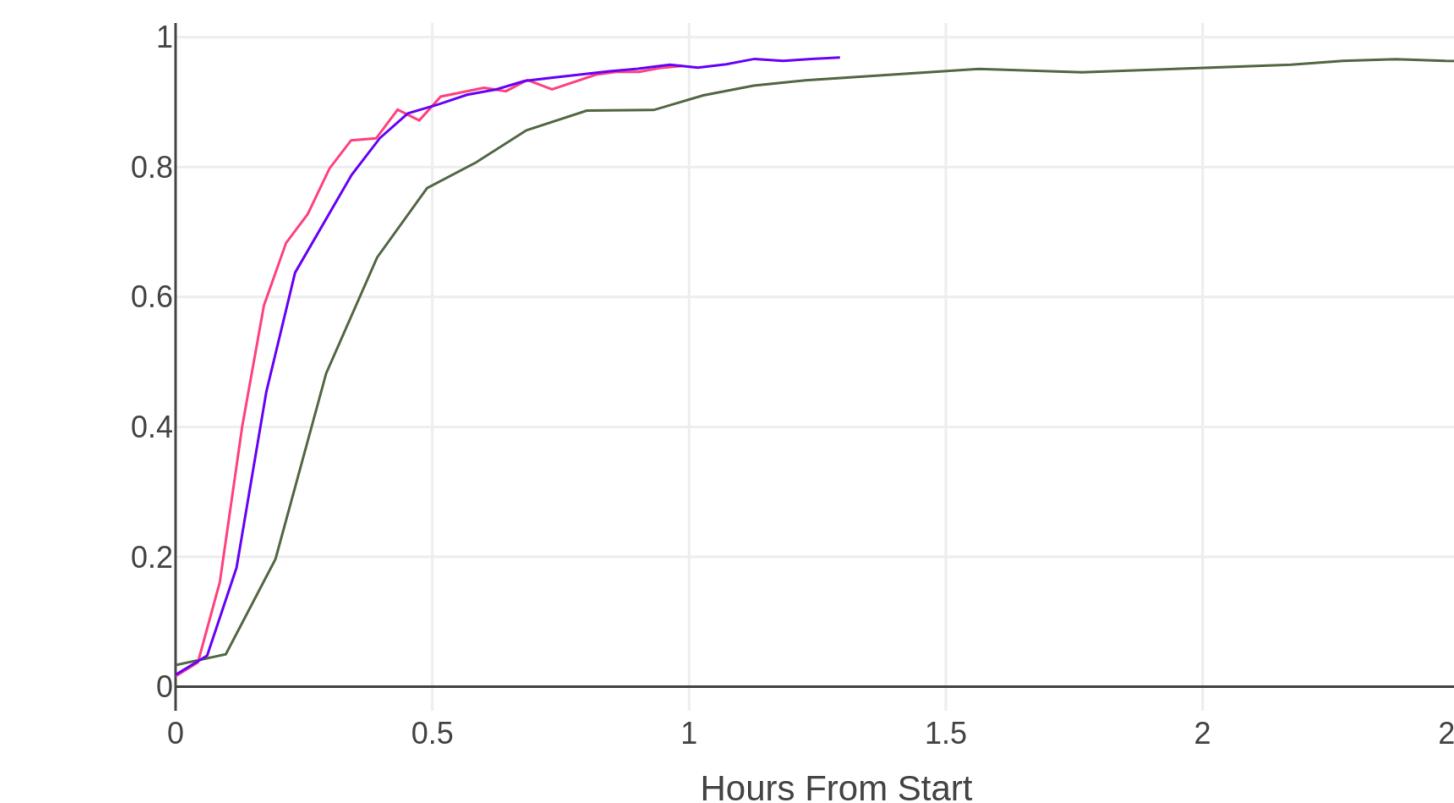


Results

VideoMAE Accuracy Comparison

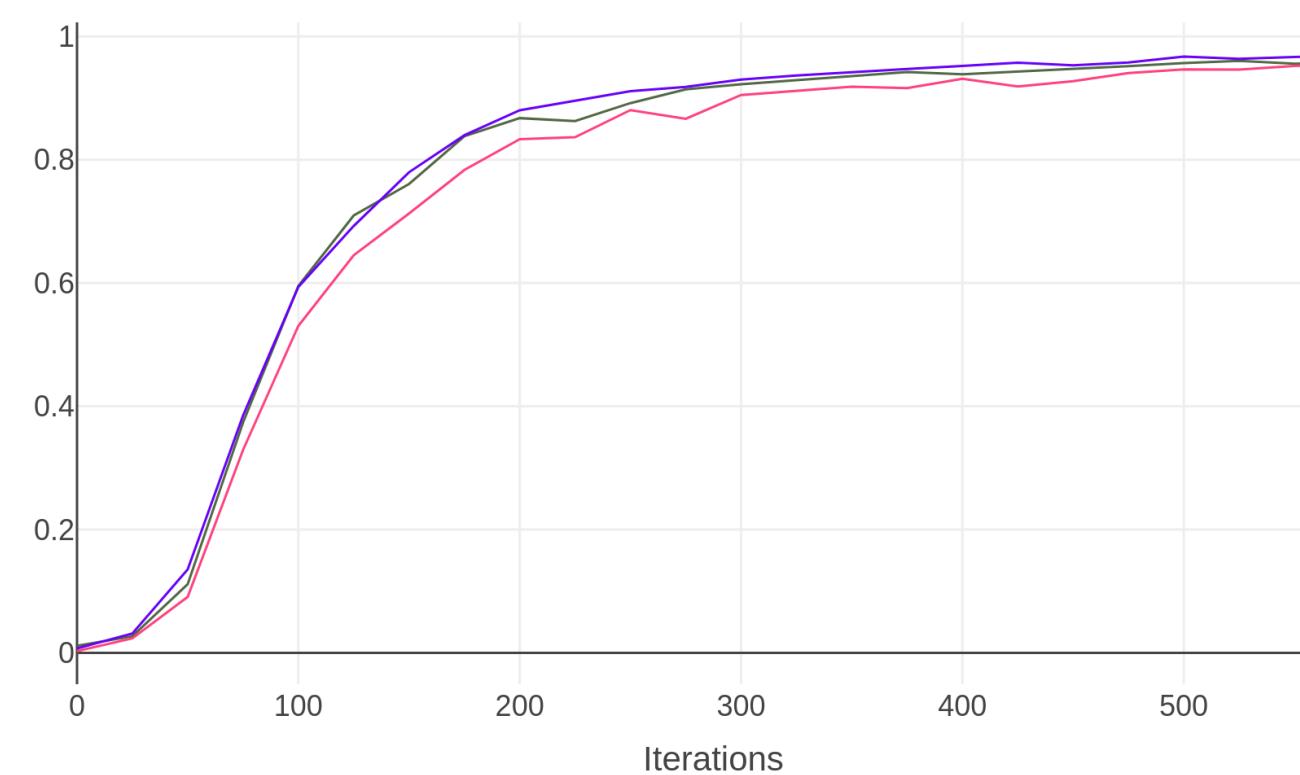


— VideoMAE Vanilla — VideoMAE ORB — VideoMAE SIFT

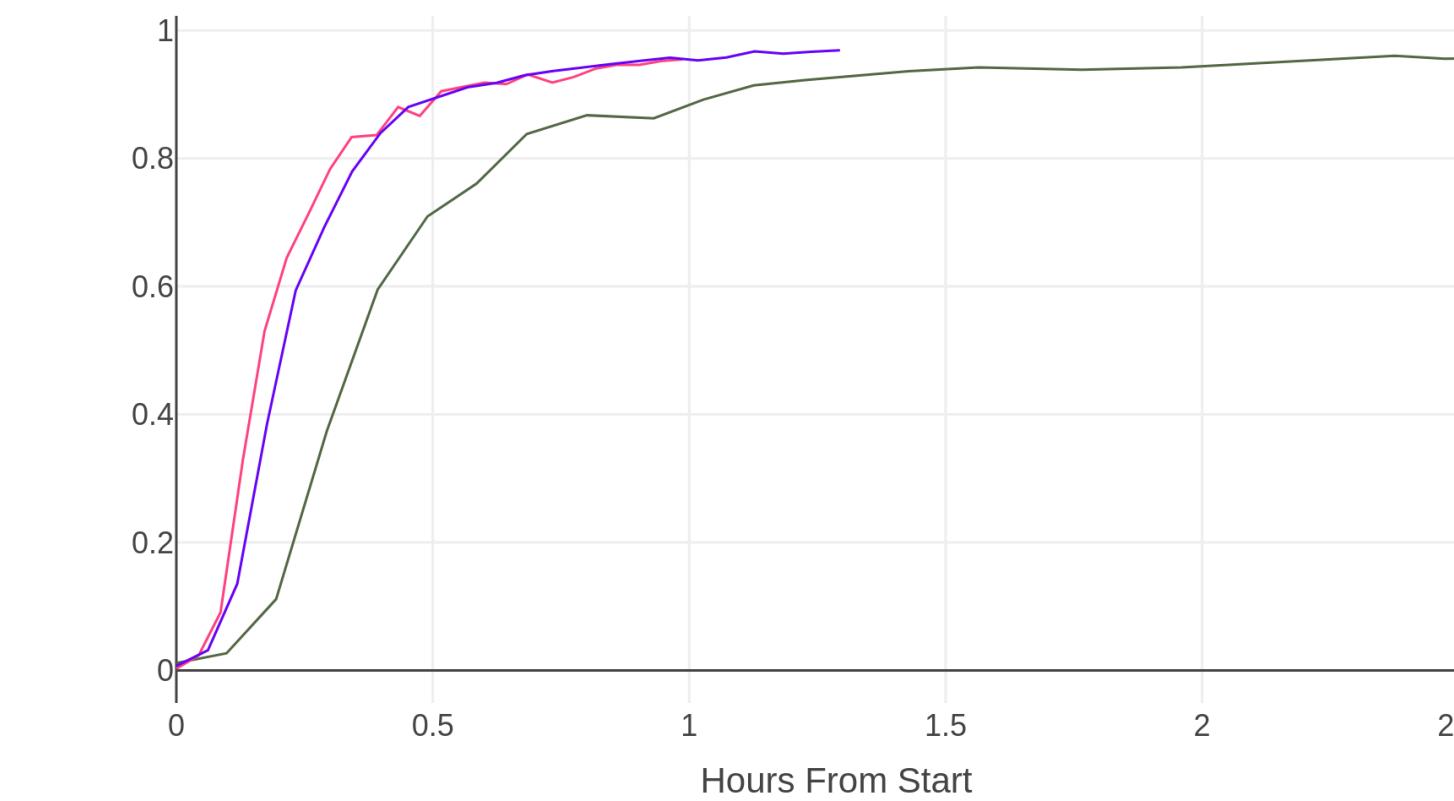


— VideoMAE Vanilla — VideoMAE ORB — VideoMAE SIFT

VideoMAE F1 Comparison



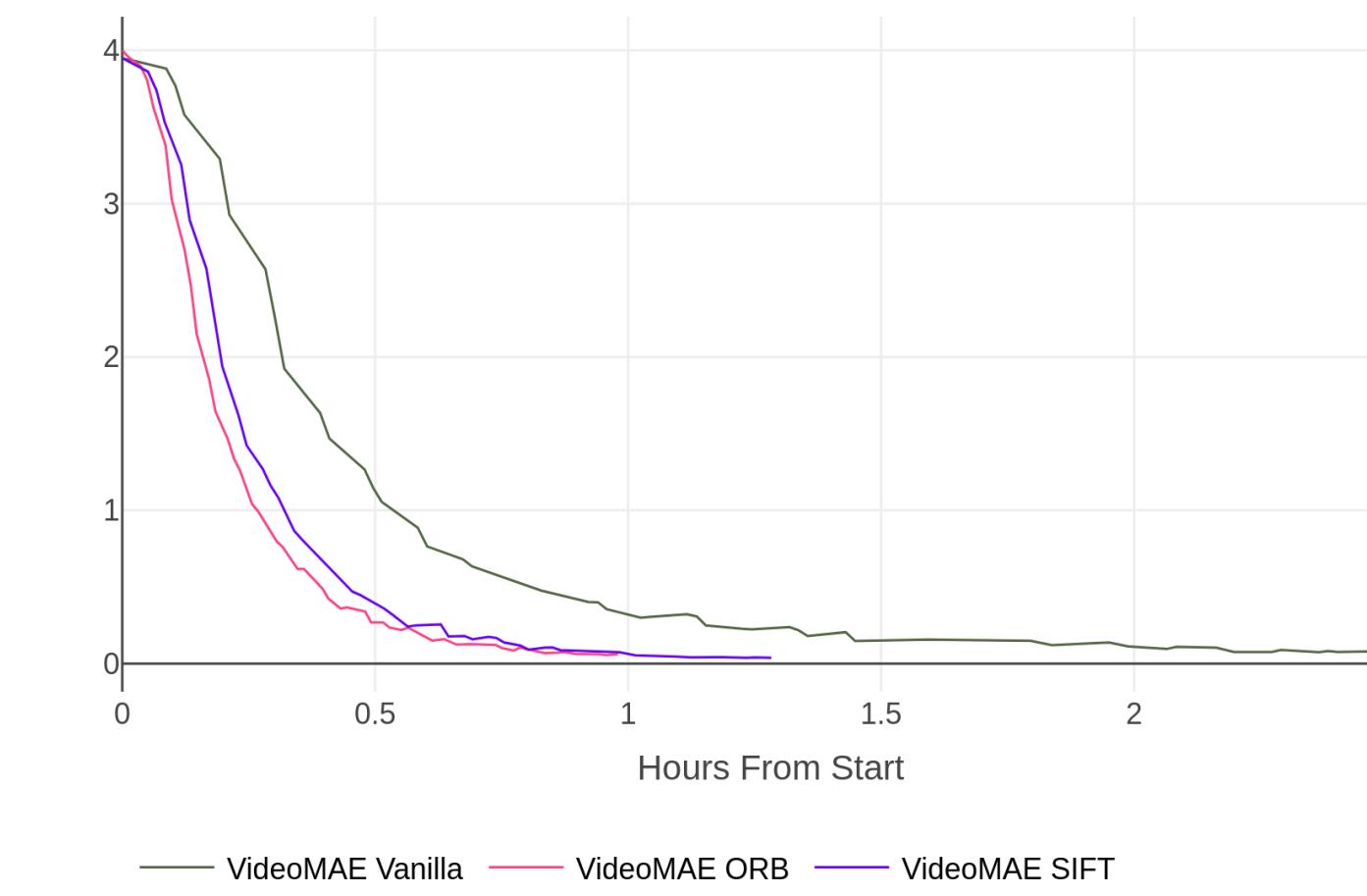
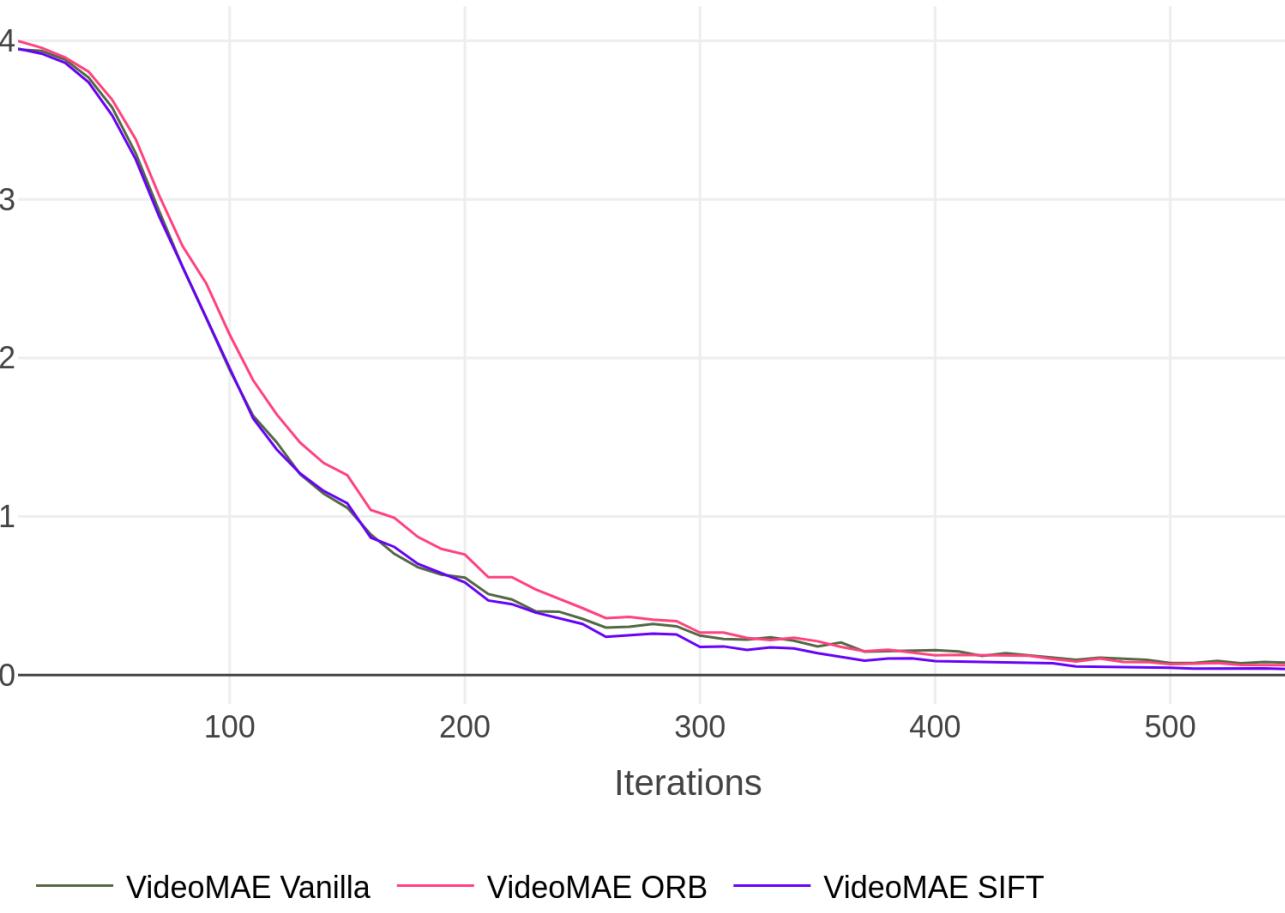
— VideoMAE Vanilla — VideoMAE ORB — VideoMAE SIFT



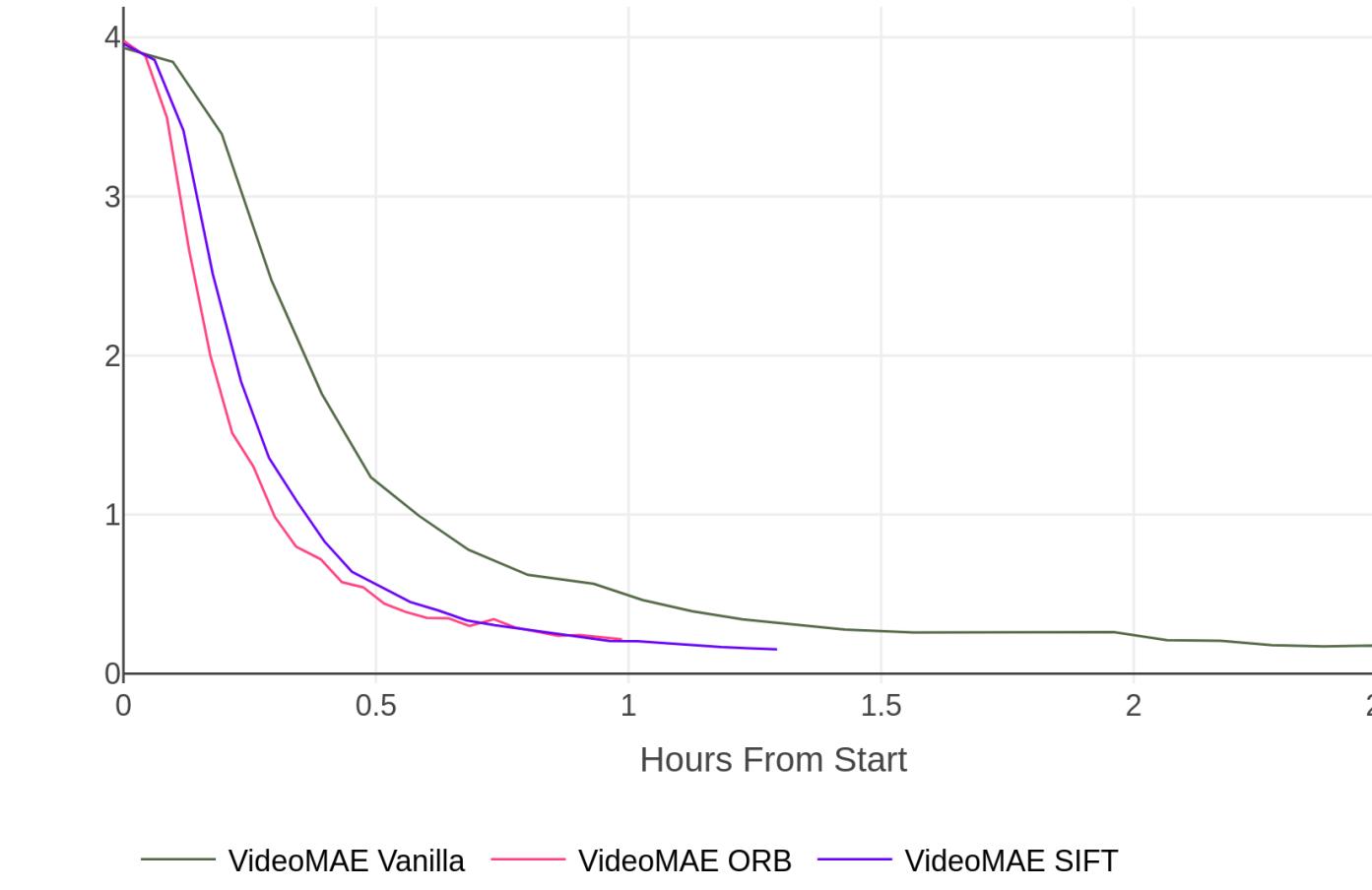
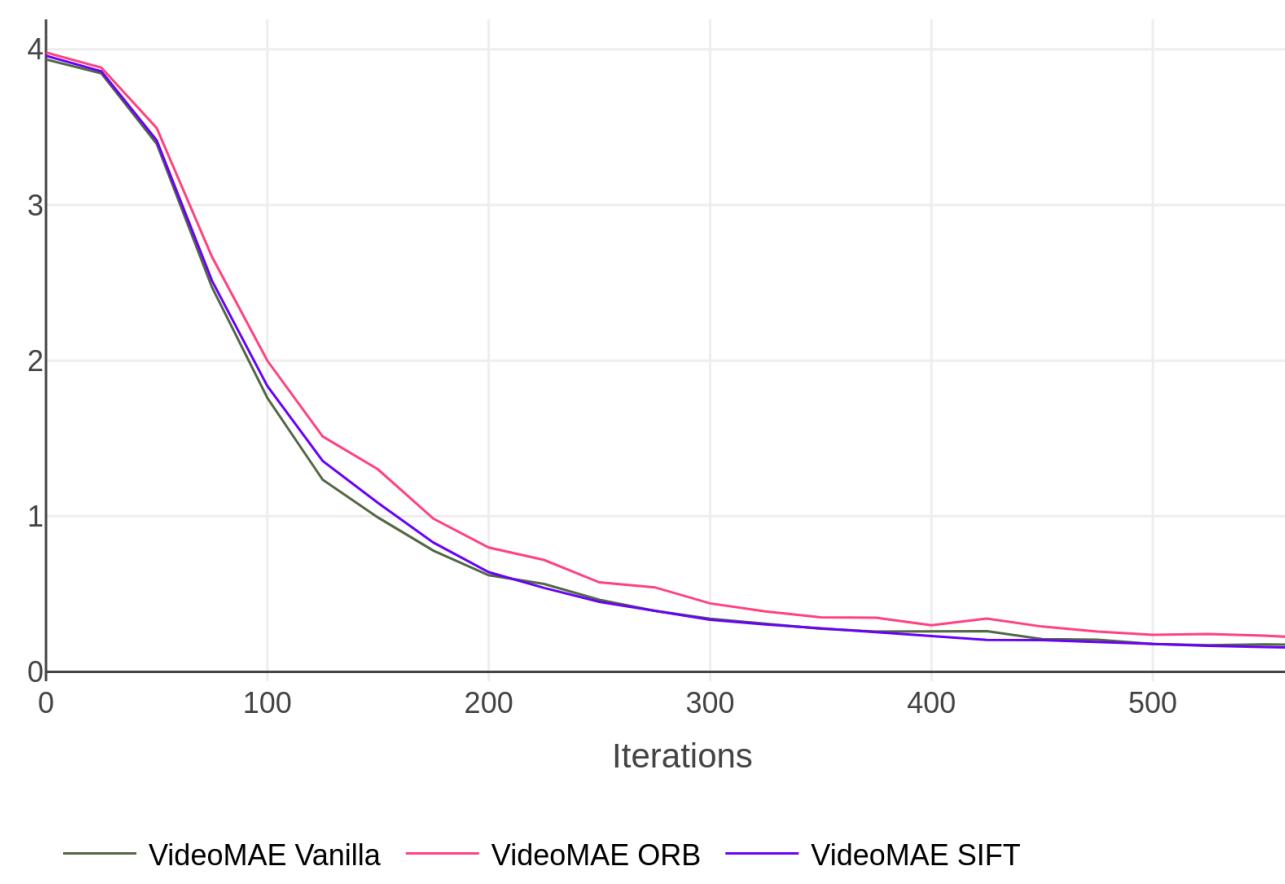
— VideoMAE Vanilla — VideoMAE ORB — VideoMAE SIFT

Results

VideoMAE Train Loss



VideoMAE Eval Loss



Results

Key Idea: selecting the most informative frames reduces dataset size and processing time without significantly degrading performance.

Findings

- ✓ A model with reduced dataset achieves almost the same accuracy as a model with the full dataset.
- ✓ Models reach optimal performance faster, reducing training time.
- ✓ Less computational resources (memory, processing time) are used.

Project Team



Ivan Golov

Teamlead, A100 hunter, ORB master



Anatoly Soldatov

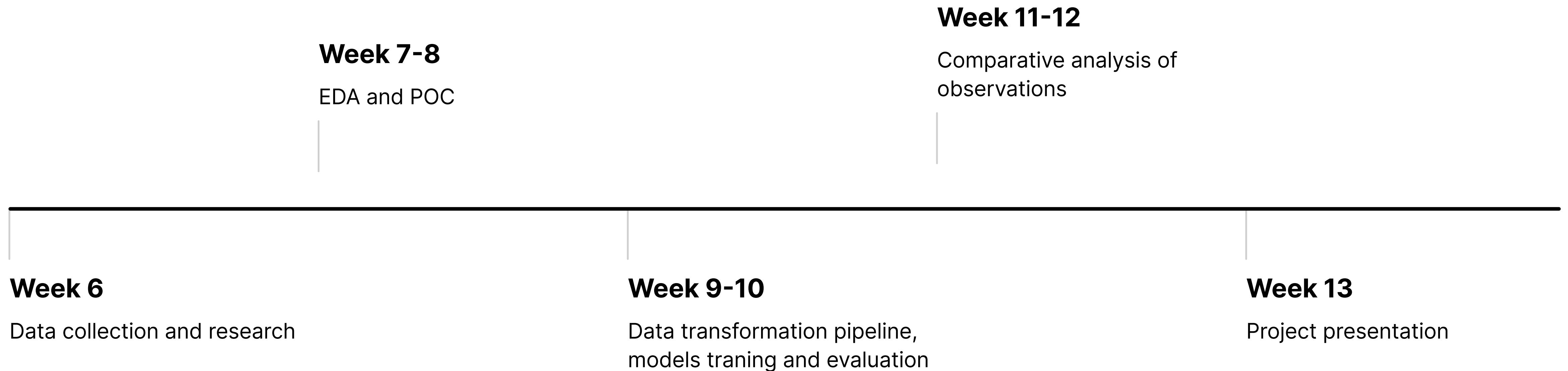
Transformer enjoyer, MLOps conqueror



Rufina Gafiiatullina

Convolutional majesty, SIFT guru

Project timeline



Q/A section