

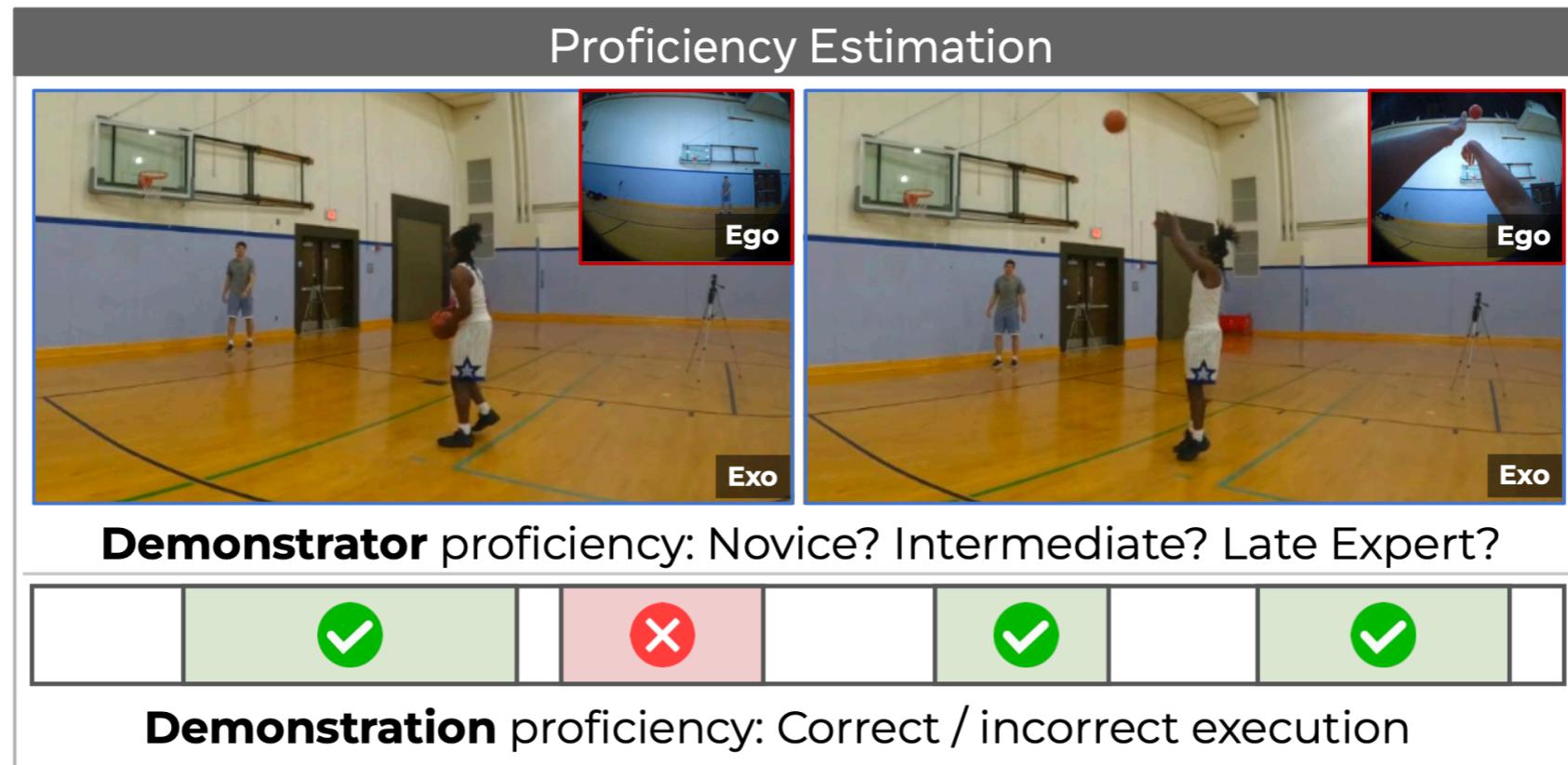
# Proficiency Estimation Final Project Update

Noah Frahm, Andrea Dunn Beltran

790: 3D Generative Models Project Update

# Understanding Proficiency Estimation

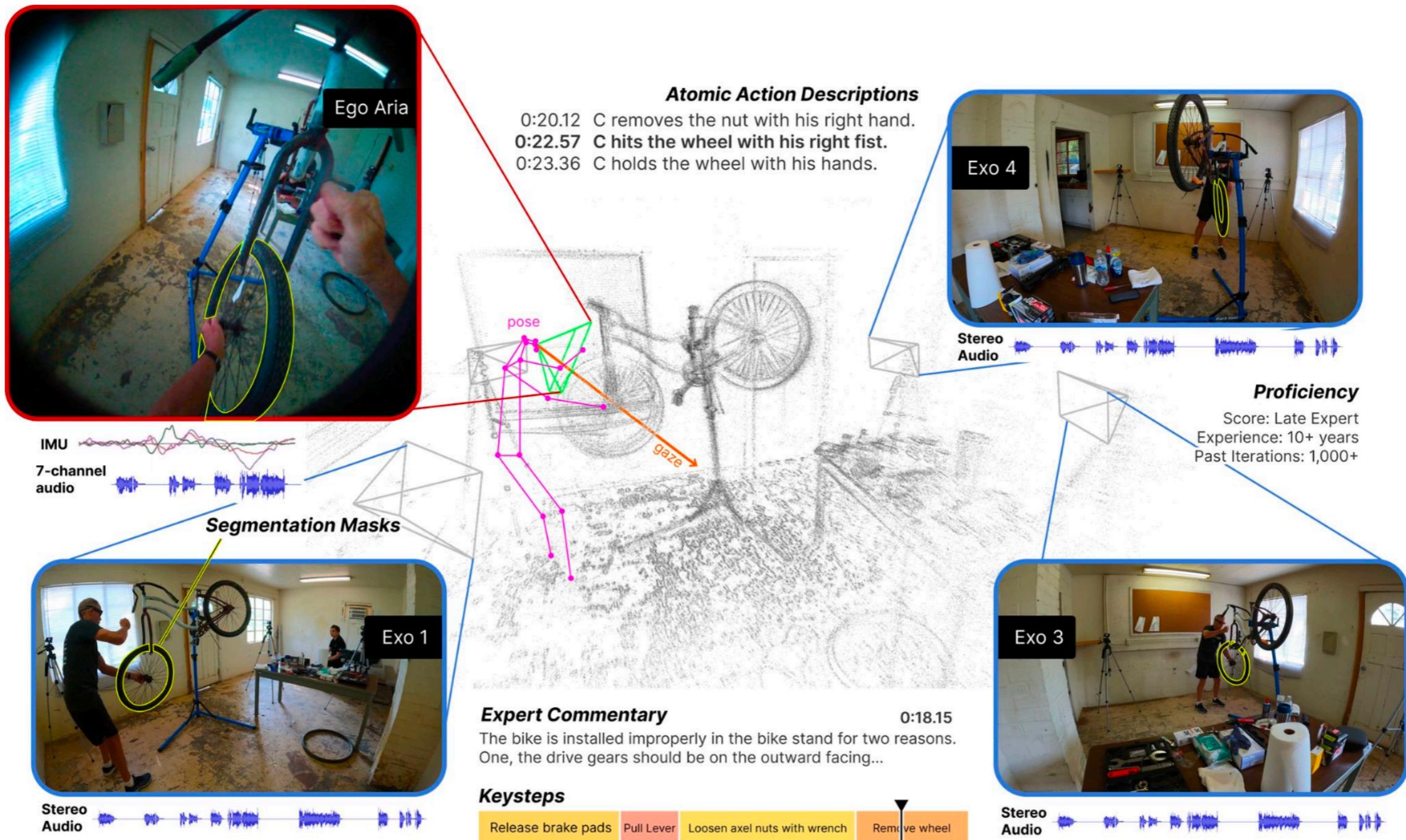
## Analyzing Video for Skill Assessment



- Labeling/recognizing key moments for skill level assessment
- Analyzing video footage to determine skill level

# Ego-Exo 4D Dataset

Diverse, large-scale multi-modal, multi-view, video dataset



# Building on Baselines

## Combining video and 3D features

1

**Classification using  
Video Features**

2

**Classification using  
3D Features**

3

**Classification using  
Video Features &  
3D Features**

# Project Progress

## Part 1



### Project Pitch

Invested EgoExo4D dataset and identified missing information

### Project Update 1

Recreated TimeSFormer baselines and tested different data splits

### Project Update 2

Extracted 3D features using MotionAGFormer with aria validation

Constructed Transformer

# Project Progress

## Part 2



### Final Project Update

Ran classification transformer on pose only data

Experimented with different tokenization techniques  
and tuned hyper parameters

Investigated different proficiency classes splits

# EgoExo4D Benchmarks

## Demonstrator Proficiency

### EgoExo4D

Method	Pretrain	Accuracy		
		Ego	Exos	Ego + Exos
Random	-	25.7	25.7	25.7
Majority-class	-	43.0	43.0	43.0
TimeSFormer [13]	-	<b>43.3</b>	40.5	40.8
TimeSFormer [13]	K400	41.9	45.2	45.2
TimeSFormer [13]	EgoVLP	42.7	<b>52.3</b>	<b>51.8</b>
Inference with multiple takes per demonstrator				
TimeSFormer [13]	EgoVLP	46.3	51.0	52.6

Table 20. **Demonstrator proficiency estimation benchmark.**  
We report top-1 accuracies for various baselines on the demonstrator proficiency estimation task.

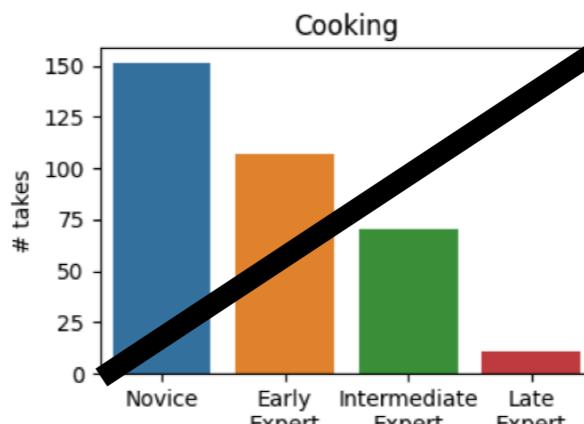
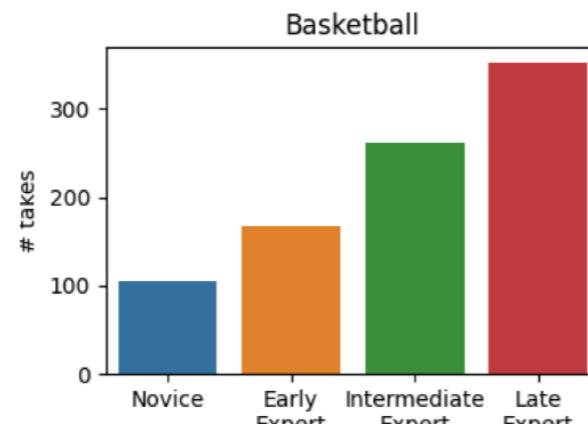
Scenario	Random	Majority-class	TimeSFormer
All Scenarios	25.25	30.25	37.21
Dance	23.96	48.04	43.55
Rock Climbing	25.60	33.06	-
Soccer	29.04	69.12	68.75
Basketball	25.56	42.88	47.29

# Dwindling Dataset

## Data availability

### Basketball

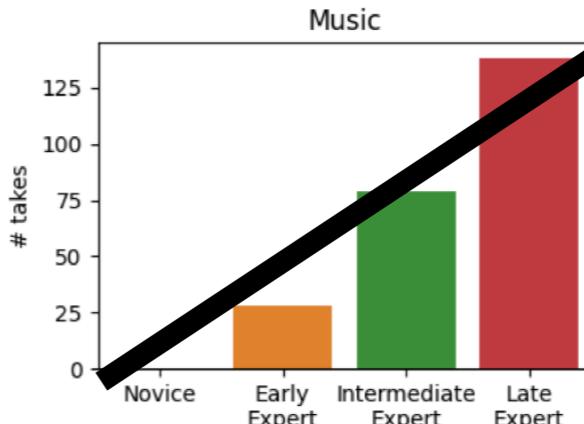
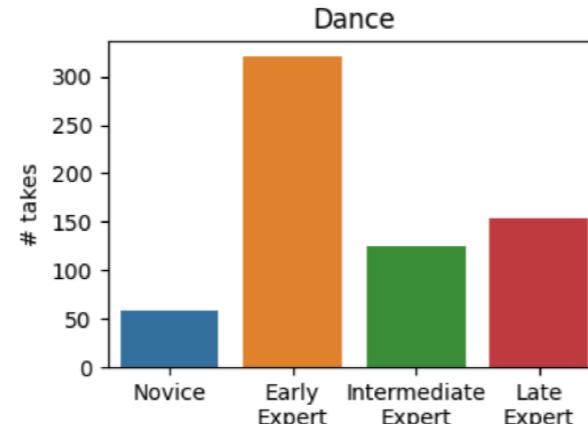
- Good proficiency class spread
- Large dataset



Not Physical Activities

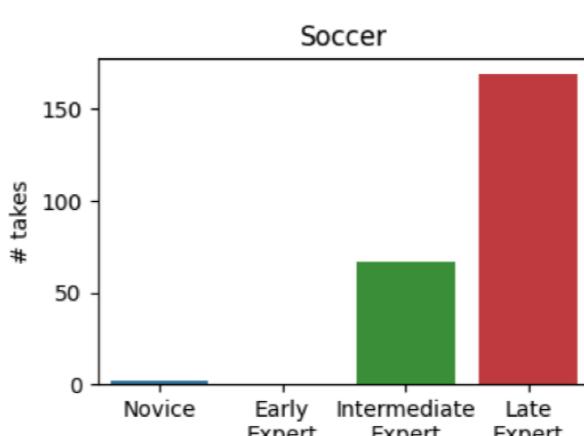
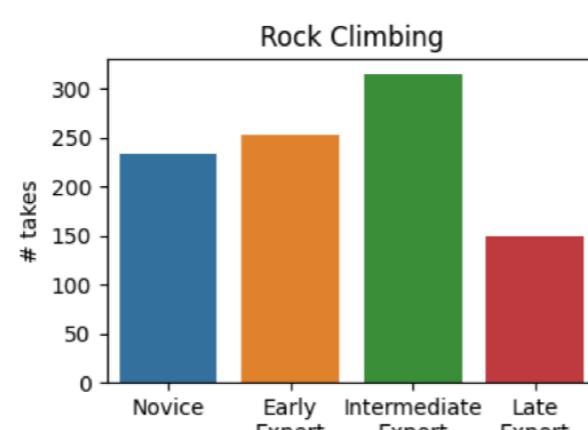
### Dance

- Good proficiency class spread
- Large dataset



### Rock Climbing

- Poor pose extraction due to heavy occlusions and camera angle



### Soccer

- Only has 2 proficiency classes
- Smallest task dataset

# Classification: 3D Features

## Old Pipeline

**Per Frame:**

Pose feature:  $17 \times 128$

3D pose:  $17 \times 3$

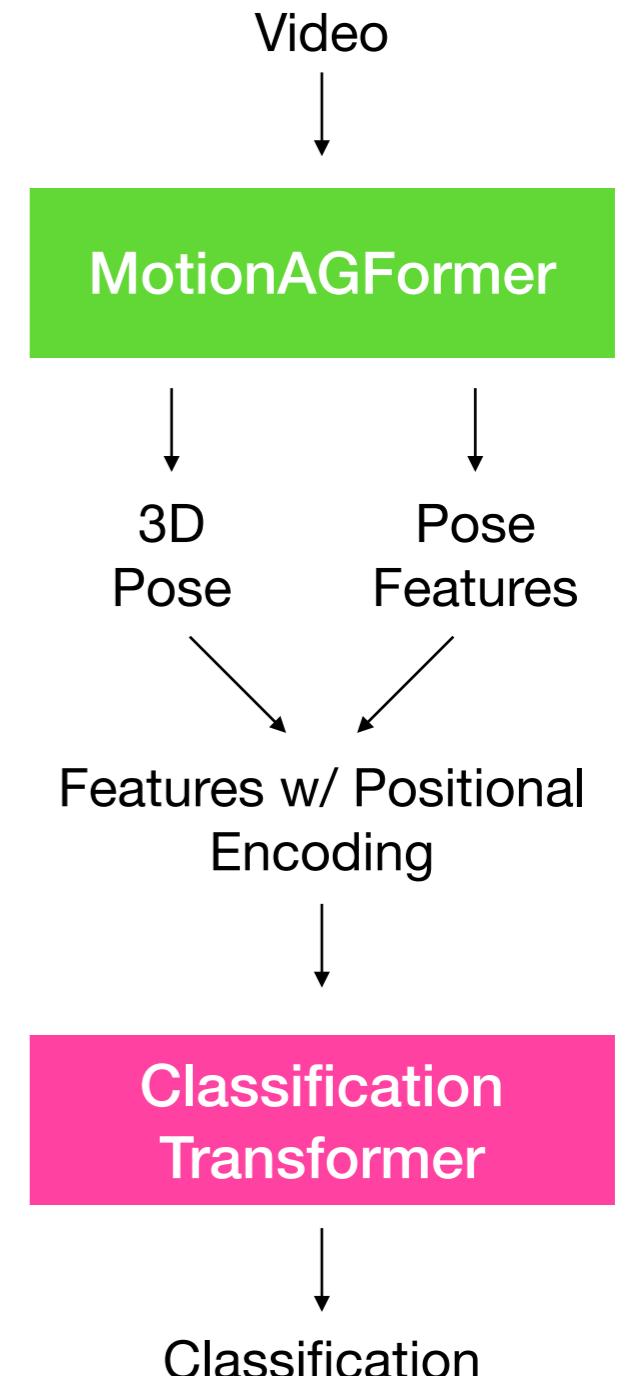
Encoded 3D Pose:  $17 \times 60$

Token size:  $17 \times 188$

Flattened:  $1 \times 3196$

**Input:** 100 frame/video

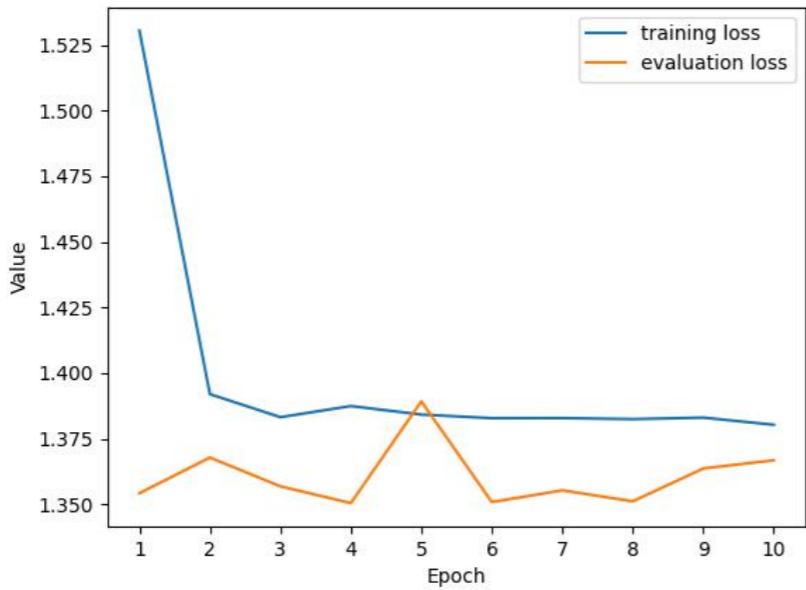
$100 \times 3196$



# Results

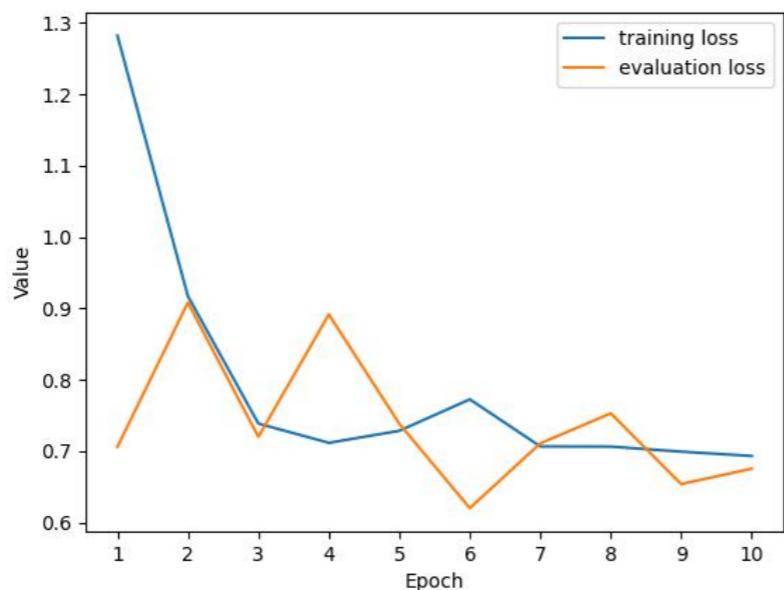
## Old Pipeline

Soccer



Top1 Acc: 80.70%  
Majority: 80.70%

Dance

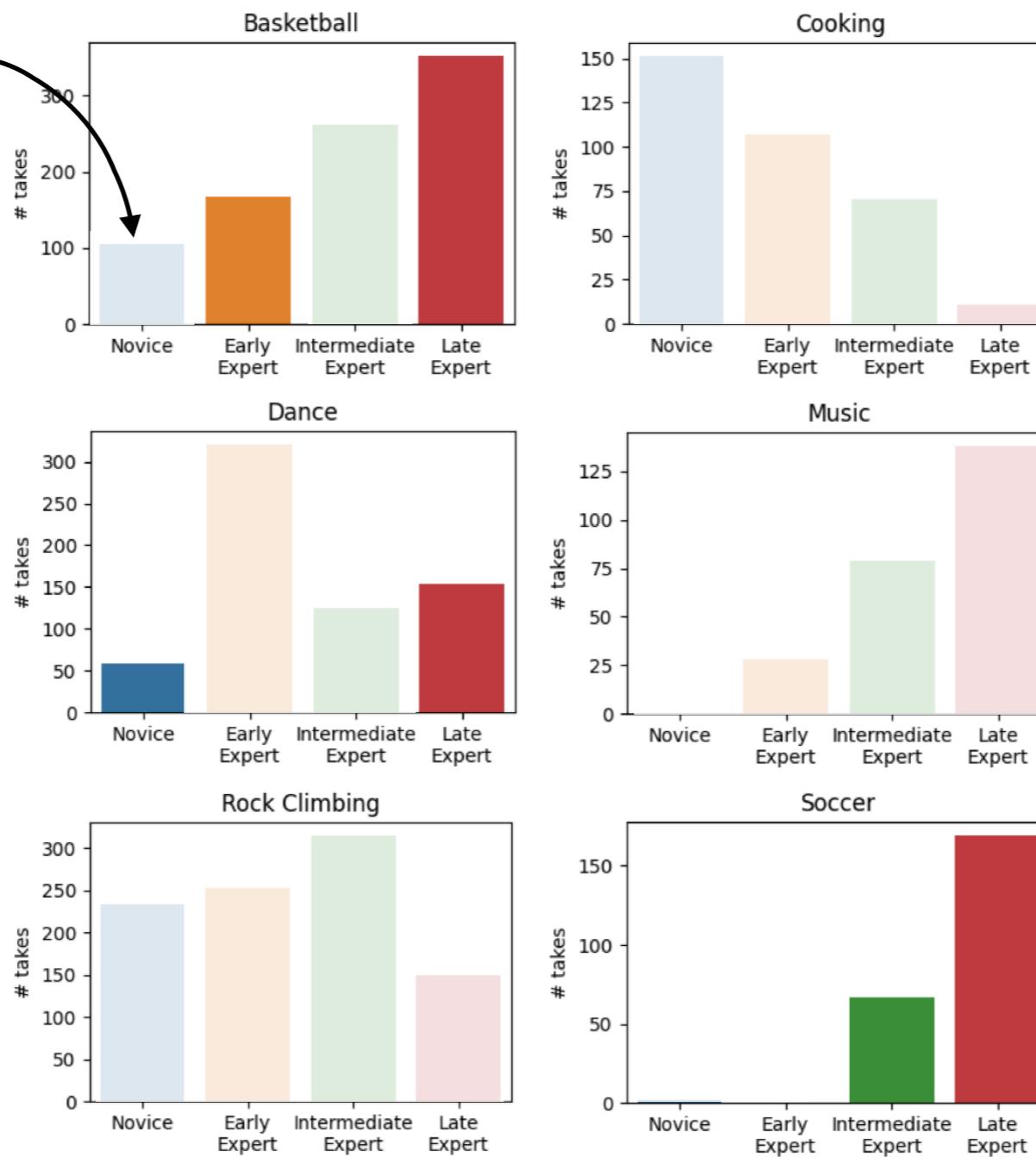


Top1 Acc: 62.20%  
Majority: 62.20%

# Binary vs Multi-class

## Simplifying the Problem

*not present in  
previous release*



# TimeSformer Results

## Majority, Multiclass, Binary

Tasks	Number classes	TimeSFormer (video only)	
		Majority of Dataset	Top_1
Soccer	2 classes	0.625	0.6875
Dance	4 classes	0.516	0.4355
	2 classes	0.75	0.9688
Basketball	3 classes	0.434	0.4729
	2 classes	0.746	0.5479

# Experimenting with Tokenization

## Changing Transformer Inputs

Old:

$100 \text{ fs} \times (17 \text{ kpts} \times (60 \text{ encoded pose} + 128 \text{ features})) \rightarrow 100 \times (17 \times 188) = 100 \times (3,196)$

Body\_Position:

$17 \text{ kpts} \times (100 \text{ fs} \times 3 \text{ pose})$   
 $17 \times (300)$

Body\_Enposition:

$17 \text{ kpts} \times (100 \text{ fs} \times 60 \text{ encoded pose})$   
 $17 \times (600)$

Frame\_Position:

$100 \text{ fs} \times (17 \text{ kpts} \times 60 \text{ encoded pose})$   
 $100 \times (1021)$

Frame\_Features:

$100 \text{ fs} \times (17 \text{ kpts} \times (3 \text{ pose} + 128 \text{ features}))$   
 $100 \times (17 \times 131) = 100 \times (2,227)$

# Other Ablations

## Changing Transformer Inputs and Architecture

- **Linear Layer** vs no linear layer
  - Linear Layer + ReLU -> 512 dim
- **Normalized** vs denormalized pose
- Encoded pose (60 dim) vs **non-encoded** pose (3 dim)
- Multi-class vs **binary**

Metrics we considered: F1, AUC, Top-1

# Classification: 3D Features

## New Pipeline

### Per Frame:

Pose feature:  $17 \times 128$

3D pose:  $17 \times 3$

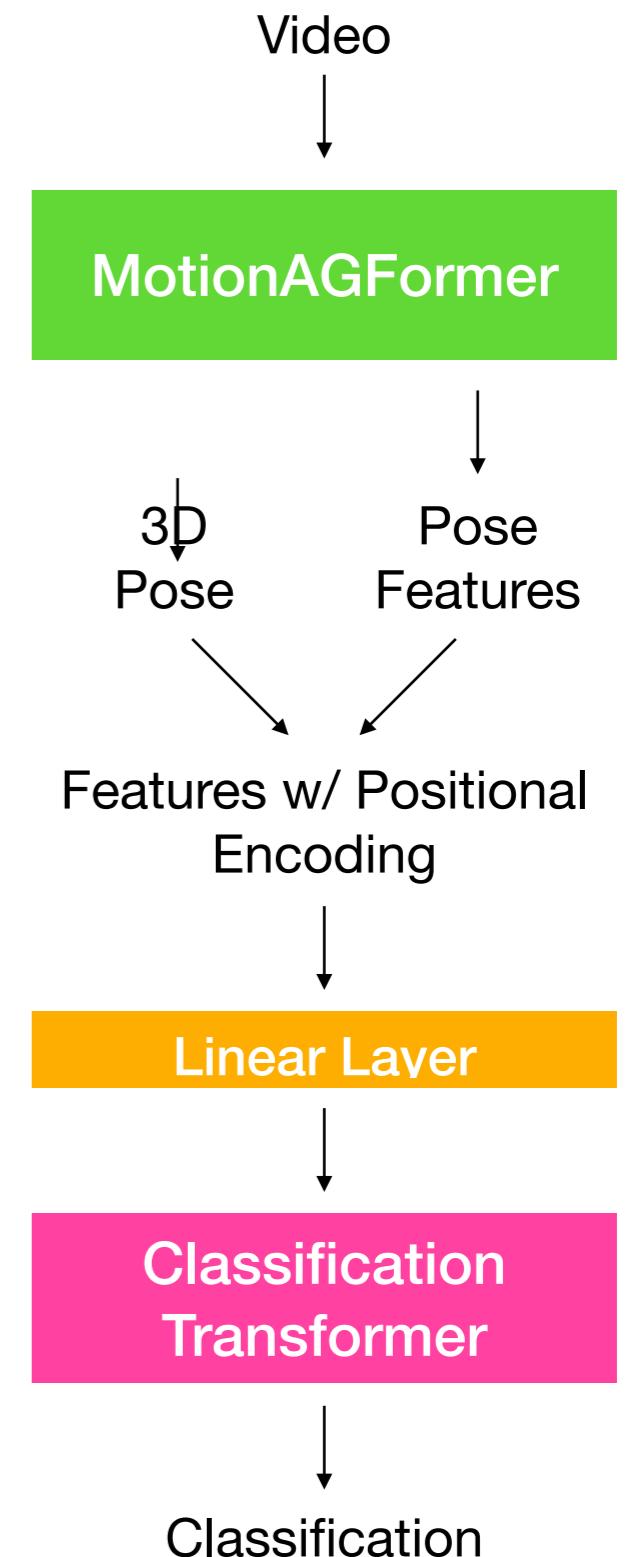
Token size:  $17 \times 131$

Flattened:  $1 \times 2227$

Projected:  $1 \times 512$

**Input:** 100 frame/video

$100 \times 512$



# Classification: 3D Features

## New Pipeline

**Per Frame:**

Flow Feature:

$128 \times 64 \times 120$

Flattened:

$128 \times 1920$

Token size:

$1 \times 1920$

Video

Unimatch

Flow  
Features

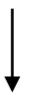
Optical  
Flow

AutoEncoder

Flow Features Token

Classification  
Transformer

Classification



# Pose Only Classification

## 3D feature Classification

Tasks	Number classes	Majority Guess			Pose Only Clasification		
		Top_1	F1	AUC	Top_1	F1	AUC
Soccer	2 classes	0.807	0.4464	0.5	0.679	0.413	0.313
Dance	4 classes	0.622	0.1916	0.5	0.452	0.453	0.64
	2 classes	0.756	0.4305	0.5	0.806	0.752	0.818
Basketball	3 classes	0.43	0.2004	0.5	0.494	0.319	0.654
	2 classes	0.672	0.4018	0.5	0.682	0.637	0.717

Model Details:

- NORMALIZED Features
- NO ENCODING on 3D pose
- Linear Layer

# Soccer Results

## Majority, Multiclass, Binary

Soccer	Top_1	F1	AUC
Majority Guess	0.807	0.4464	0.5
Pose Only Clasification	0.679	0.413	0.313

# Dance Results

## Majority, Multiclass, Binary

Dance 2 classes	Top_1	F1	AUC
Majority Guess	0.756	0.4305	0.5
Pose Only Clasification	<b>0.806</b>	<b>0.752</b>	<b>0.818</b>

Dance 4 classes	Top_1	F1	AUC
Majority Guess	<b>0.622</b>	0.1916	0.5
Pose Only Clasification	0.452	<b>0.453</b>	<b>0.64</b>

# Basketball Results

## Majority, Multiclass, Binary

<b>Basketball 2 classes</b>	<b>Top_1</b>	<b>F1</b>	<b>AUC</b>
<b>Majority Guess</b>	0.672	0.4018	0.5
<b>Pose Only Clasification</b>	<b>0.682</b>	<b>0.637</b>	<b>0.717</b>

<b>Basketball 3 classes</b>	<b>Top_1</b>	<b>F1</b>	<b>AUC</b>
<b>Majority Guess</b>	0.43	0.2004	0.5
<b>Pose Only Clasification</b>	<b>0.494</b>	<b>0.319</b>	<b>0.654</b>

# TimeSformer Results

## Majority, Multiclass, Binary

Tasks	Number classes	TimeSFormer (video only)	
		Majority of Dataset	Top_1
Soccer	2 classes	0.625	0.6875
Dance	4 classes	0.516	0.4355
	2 classes	0.75	0.9688
Basketball	3 classes	0.434	0.4729
	2 classes	0.746	0.5479

# Next Steps

## Tasks and Modalities

- Investigate other modalities
  - Optical Flow
  - Monocular Depth
- Fusing different modalities for tasks
- Expand to more tasks