

# Leveraging Temporal Context in Low Representational Power Regimes

WED-AM-235

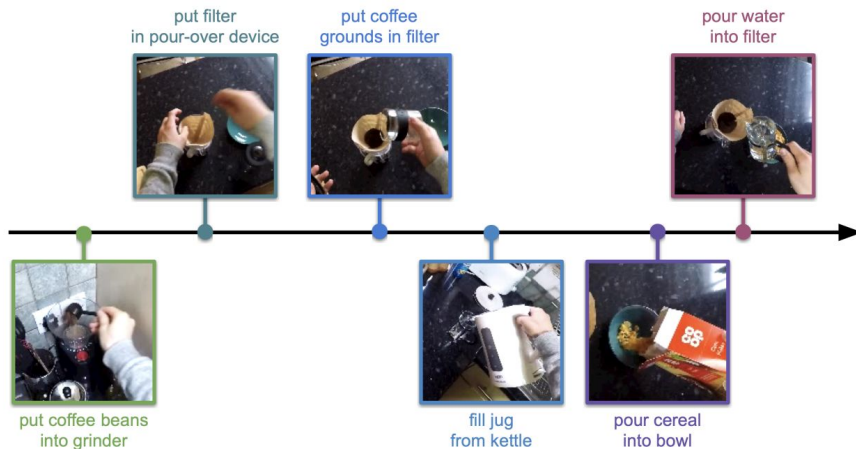
Camilo Fosco, SouYoung Jin, Emilie Josephs, Aude Oliva  
**Olivalab**, MIT CSAIL



# The problem

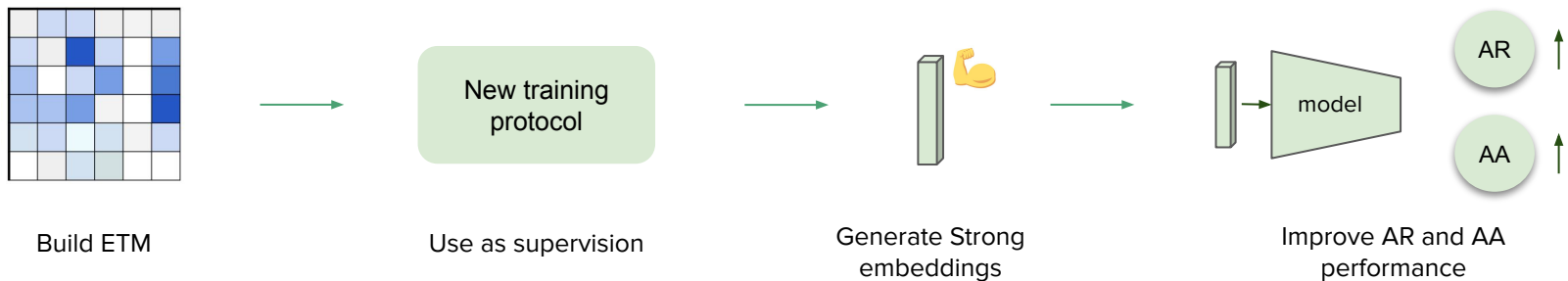
Using temporal information is crucial to understand videos. Yet, current models don't explicitly attempt to leverage temporal regularities in datasets with long videos.

- Can we leverage the **statistics in temporal sequences** of video datasets to **improve performance** in downstream tasks?
- Can we build richer embeddings with this information?
- Where does this type of information **help the most**?



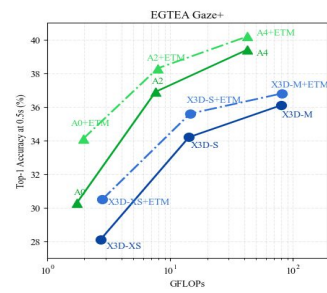
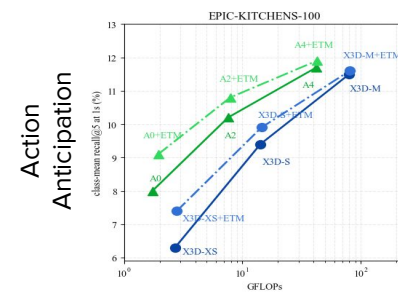
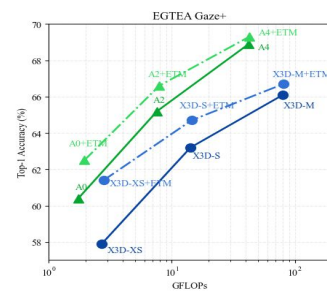
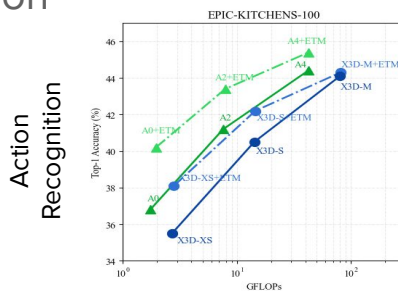
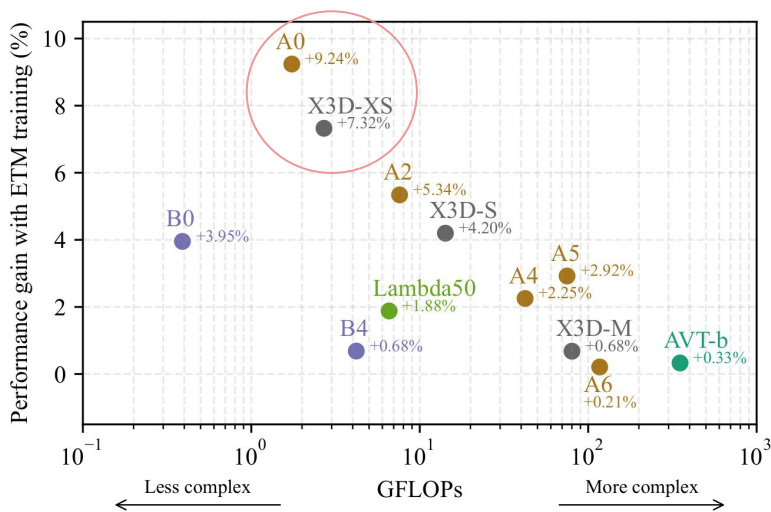
# The idea

- We propose to build an **Event Transition Matrix**: a representation that captures typical transition probabilities between actions in long video sequences
- We use this matrix as supervision in a new training protocol to generate **strong embeddings for video snippets**
- We leverage these embeddings to **improve action recognition and action anticipation performance**, especially on low complexity models.



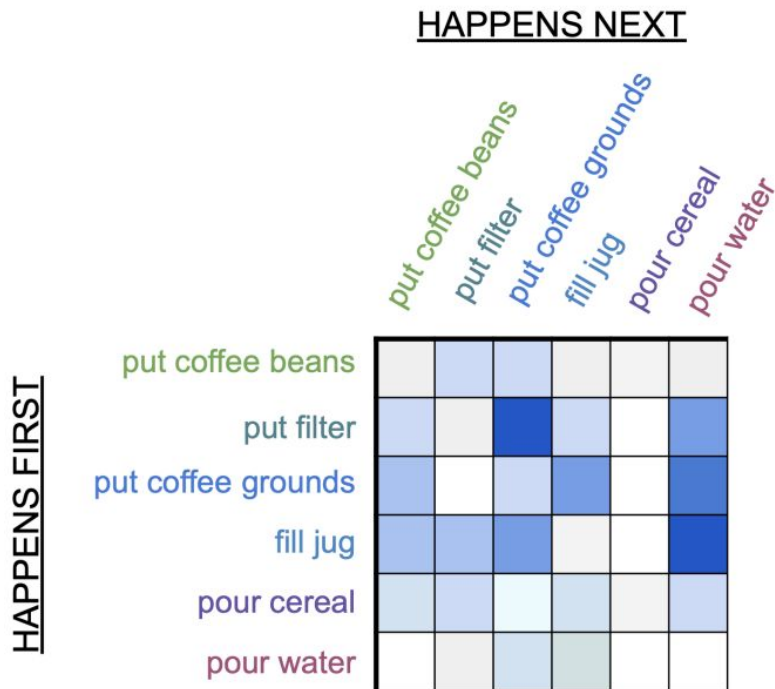
# Key results

Our model-agnostic framework helps **low complexity models** **improve performance** on action recognition and action anticipation across 3 datasets.



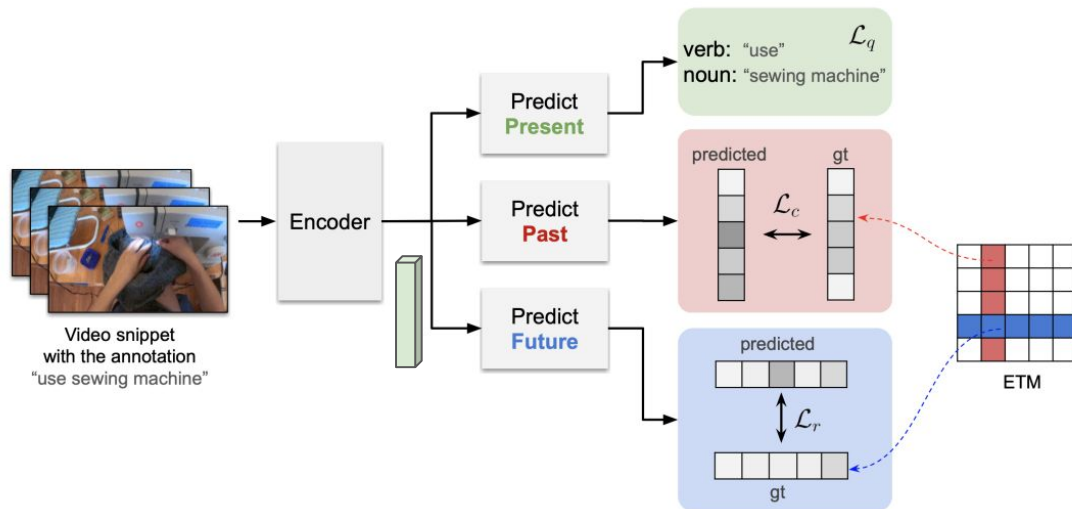
# The Event Transition Matrix

- Computed by looking at **all actions** happening after a given action, weighted by a decay function
- Square matrix, **not symmetric**
- Several postprocessing steps:
  - Dimensionality reduction
  - Decay definition
  - Normalization



# How can we leverage this ETM?

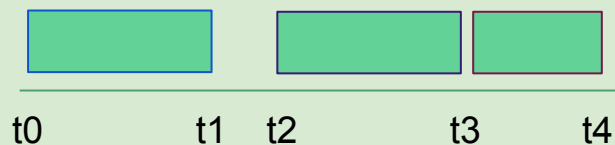
- We propose to use the rows and columns of the matrix as **targets in a regression problem**
- An encoder is tasked to **generate an embedding** that can predict the action + regress the past and future





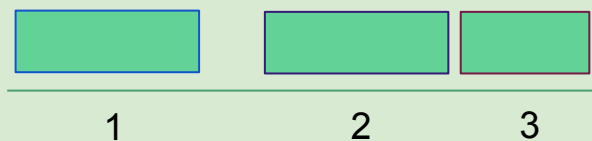
# Types of distance metrics

## Distance in frames/seconds between actions



- Takes into account temporal difference
- Can differentiate between end/start (but which one to choose?)
- But dependent on length of actions

## Distance as index difference in ordinal sequence



- Only considers ordering
- Might be more adapted to causality concepts
- Doesn't depend on length of actions



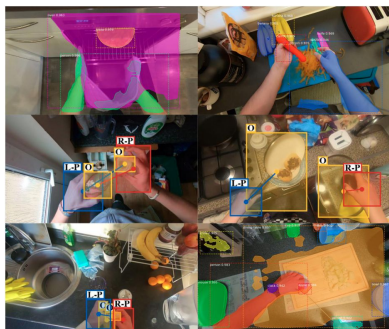
# Testing ETM design choices

We test several design choices, including different decay functions, ETM sizes and distance metrics.

Size	Decay	Temp. Metric	<b>Present</b> (top-1 accuracy)		
			Verb	Noun	Action
13k	linear	time	0.551	0.462	0.288
13k	exponential	time	0.556	0.477	0.291
2.5k	exponential	time	0.586	0.488	0.313
2.5k	no decay	-	0.581	0.480	0.305
2.5k	linear	index	0.601	0.493	0.319
2.5k	exponential	index	<b>0.603</b>	<b>0.503</b>	<b>0.324</b>

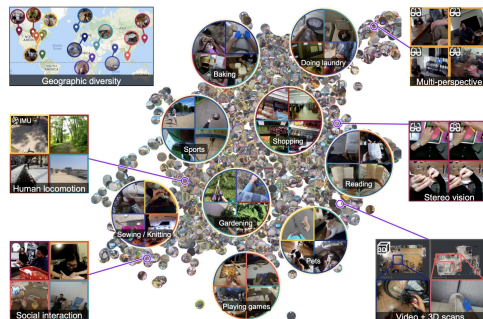
# Datasets used

## EPIC-KITCHENS-100



700 videos depicting cooking actions, totalling 100 hours.

## EGO4D LTA



3670 hours of video from 71 different participants.

## EGTEA Gaze+



10k segments annotated with 19 verbs, 51 nouns and 106 unique actions.

# Tasks where we test our embeddings

## Action Recognition

Receive a snippet, predict class label (here, verb and noun)



Roll Dough

## Action Anticipation

Receive a collection of snippets, predict next action



Put Down Dough

# Experimental results - Action recognition

## Cross Dataset

Results with MoViNet A0

Dataset	Model	Present		
		Verb	Noun	Action
EK100 [14]	Baseline	64.8	47.4	36.8
	ETM(Ours)	<b>67.9</b>	<b>51.2</b>	<b>40.2</b>
EGO4D	Baseline	32.3	23.5	21.1
LTA [16]	ETM(Ours)	<b>32.9</b>	<b>24.2</b>	<b>22.0</b>
EGTEA	Baseline	81.2	71.7	60.4
Gaze+ [28]	ETM(Ours)	<b>83.4</b>	<b>72.9</b>	<b>62.5</b>

## Cross Model

Results on EK100 across a wide variety of models

Model	w/o ETM	w/ ETM
MoviNet A0 [24]	36.8	<b>40.2</b>
MoviNet A2 [24]	41.2	<b>43.4</b>
X3D-XS [11]	35.5	<b>38.1</b>
X3D-S [11]	40.5	<b>42.2</b>
ConvNeXt-S 224 [31]	20.1	<b>32.4</b>
LambdaResNet-50 [4]	26.6	<b>27.1</b>
EfficientNet-B0 [57]	25.3	<b>26.3</b>
EfficientNet-B4 [11]	29.2	<b>29.4</b>
AVT-b [14]	30.4	<b>30.7</b>

# Experimental results - Action Anticipation

## Cross Dataset

Results with MoViNet A0

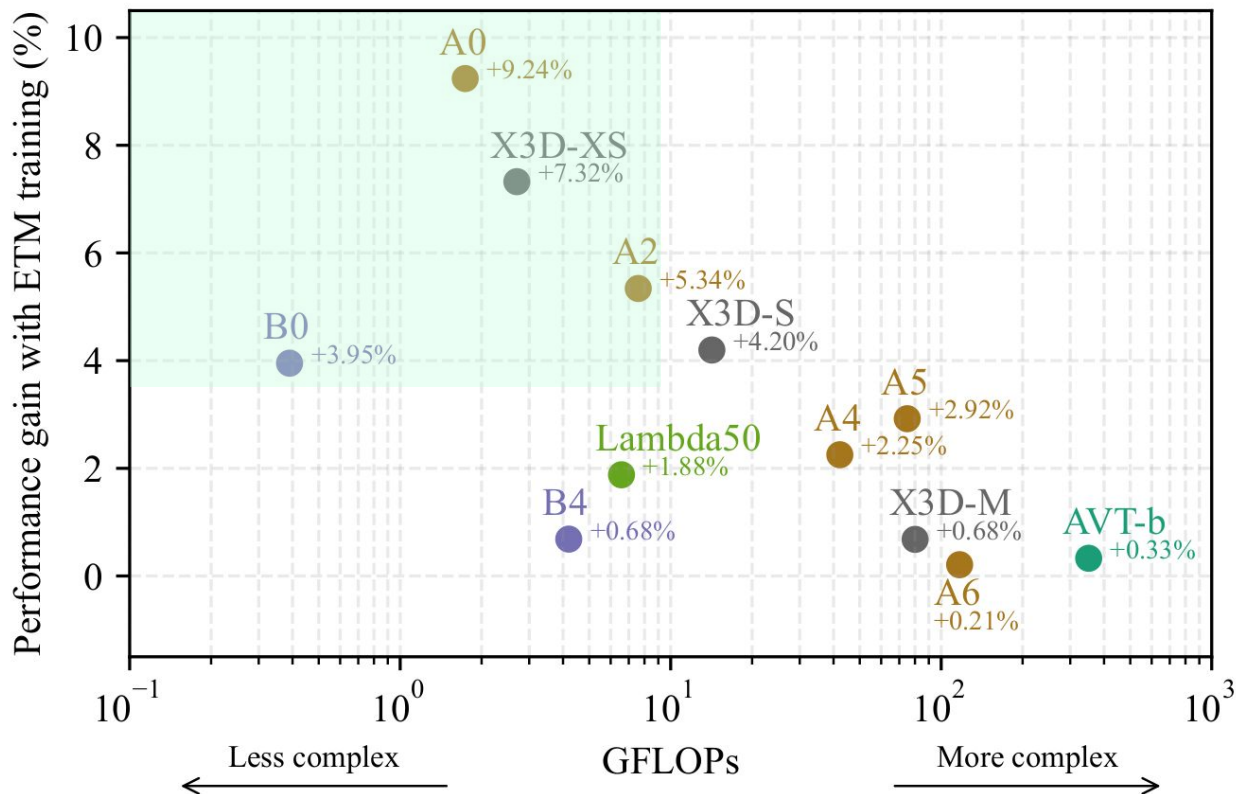
Dataset	Frozen Encoder?	Baseline			ETM (Ours)		
		Verb ↑	Noun ↑	Action ↑	Verb ↑	Noun ↑	Action ↑
EK100	✓	19.9	20.4	7.2	<b>21.5</b>	<b>20.5</b>	<b>8.1</b>
		20.8	21.3	8.0	<b>22.4</b>	<b>22.7</b>	<b>9.1</b>
EGO4D LTA	✓	17.1	16.6	10.3	<b>18.1</b>	<b>17.8</b>	<b>11.4</b>
		18.2	17.5	11.1	<b>19.9</b>	<b>19.1</b>	<b>12.9</b>
EGTEA Gaze+	✓	42.1	37.6	28.9	<b>43.4</b>	<b>38.9</b>	<b>31.3</b>
		43.5	38.5	30.3	<b>46.5</b>	<b>40.7</b>	<b>34.1</b>

## Cross Models

Results on EK100

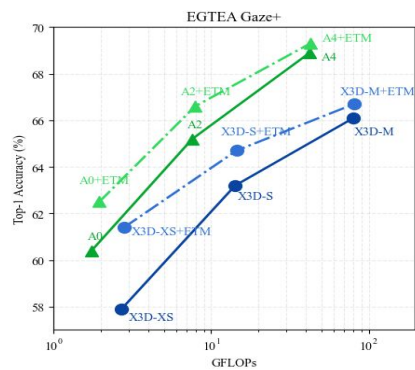
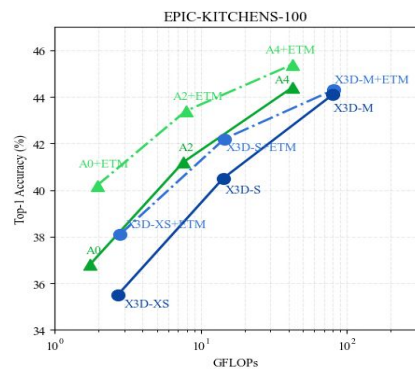
Encoder	w/o ETM	with ETM
MoViNet A0 [24]	8.0	<b>9.1</b>
MoViNet A2 [24]	10.2	<b>10.8</b>
X3D-XS [11]	6.3	<b>7.4</b>
X3D-S [11]	9.4	<b>9.9</b>
ConvNeXt-S 224 [31]	4.1	<b>5.0</b>
EfficientNet B0 [57]	7.2	<b>8.0</b>
EfficientNet B4 [57]	9.4	<b>10.1</b>
AVT-b [14]	13.4	<b>13.5</b>

# Larger gains on smaller models!

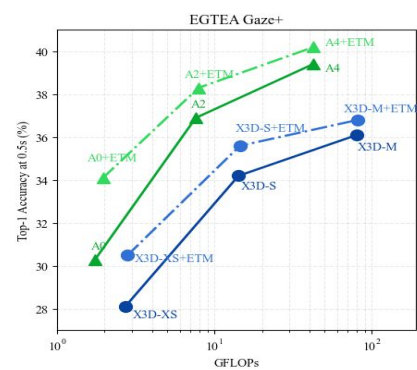
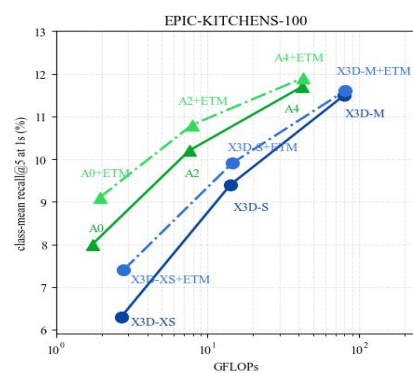


# Performance on architecture families

## Action Recognition



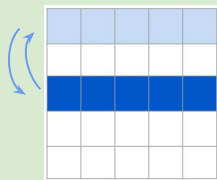
## Action Anticipation



MoViNet+ETM (—▲); MoViNet (—▲); X3D+ETM (—●); X3D (—●).

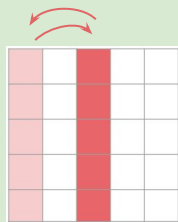
# Baseline tests and ablations

## Shuffle Rows



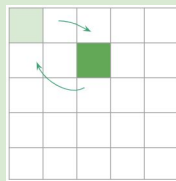
Row shuffling in ETM matrix: Distribution of **future actions** doesn't match the action index at a given **row**.

## Shuffle Columns



Row shuffling in ETM matrix: Distribution of **past actions** doesn't match the action index at a given **column**.

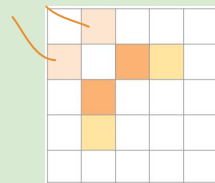
## Full Shuffle



Full shuffling in ETM matrix: no transition probability estimation matches its original action pair.

## Co-occurrence

Co-occurrence frequency between a1 and a2



Using a co-occurrence matrix: cells correspond to **co-occurrence frequencies** instead of transition probabilities

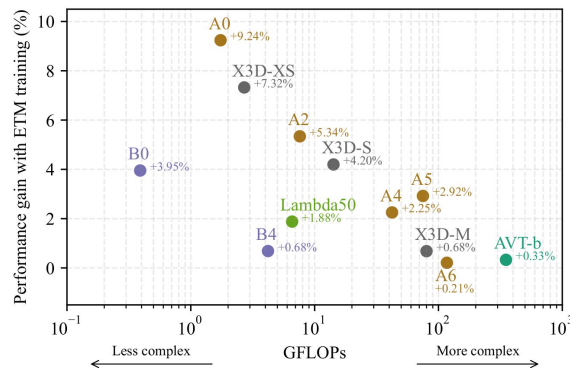
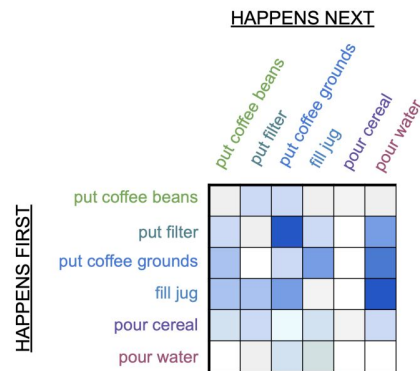


# Baseline tests and ablations

Model	Present			MAE on	MAE on
	Verb $\uparrow$	Noun $\uparrow$	Action $\uparrow$	<b>Past</b> $\downarrow$	<b>Future</b> $\downarrow$
Baseline	64.8	47.4	36.8	-	-
Full shuffle	64.1	47.2	36.3	4.117	4.012
Columns/rows shuffle	64.7	47.6	36.7	3.254	3.101
Co-occurrence	65.3	49.0	37.9	1.211	1.115
Only past vector	65.7	49.3	38.2	0.901	-
Only future vector	65.5	49.8	38.3	-	0.898
<b>ETM (Ours)</b>	<b>67.9</b>	<b>51.2</b>	<b>40.2</b>	<b>0.882</b>	<b>0.859</b>

# Conclusions

- We introduce a new training regime that uses **external temporal regularities** to boost video understanding.
- Using our ETM as a training target enhances action recognition and anticipation, particularly on **low representational power models**.
- Our ETM protocol's key benefits: flexibility, simplicity, cost-effectiveness, and easy integration.



# Thank you!

## Leveraging Temporal Context in Low Representational Power Regimes

Camilo Fosco, SouYoung Jin, Emilie Josephs, Aude Oliva

Project page:

[camilofosco.com/etm\\_website](https://camilofosco.com/etm_website)

Contact: [camilolu@mit.edu](mailto:camilolu@mit.edu)

