

Leveraging Temporal Context in Low Representational Power Regimes Camilo Fosco, SouYoung Jin, Emilie Josephs, Aude Oliva MIT CSAIL

Problem

Can the performance of low-complexity video understanding models be improved by supplementing training with additional information about temporal regularities in a dataset?

We explore this in the domains of action recognition and action anticipation, over egocentric videos.

- We introduce The Event Transition Matrix (ETM), computed from action labels in an untrimmed video dataset, which captures the temporal context of a given action.
- We show that including ETM information during training improves action recognition and anticipation performance on various egocentric video datasets.



The advantage

Training with the ETM approach increases performance on action recognition. Interestingly, lower complexity models appear to benefit more from this approach.



How do we build the ETM?

We combine information from all previous and subsequent events, weighted by their temporal distance from the queried action, capturing long-range relationships among events. This temporal distance can be measured in two ways: time or number of events, with different tradeoffs.





Ablation studies

To show that these effects are due to our specific protocol, we train models with simple baselines and compare AR performance.

Model	Present			MAE on	MAE on
WIGGET	Verb ↑	Noun ↑	Action ↑	Past ↓	Future \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
ETM (Ours)	67.9	51.2	40.2	0.882	0.859

Proposed Approach

A MoViNet A0 pre-t supervision shows i performance over In AA, We observe 12.5%, 10.7% and 8 EGO4D and EGTEr



Results

					20
trained with ETM	Dataset	Model	Present		
incrosed	Dataset	wiodei	Verb	Noun	Action
	FK100	Baseline	64.8	47.4	36.8
3 different datasets.	LIXIOU	ETM(Ours)	67.9	51.2	40.2
improvements of	EGO4D	Baseline	32.3	23.5	21.1
1 8 3% on EK100	LTA	ETM(Ours)	32.9	24.2	22.0
$\mathbf{O} \cdot \mathbf{O} \cdot $	EGTEA	Baseline	81.2	71.7	60.4
A Gaze+, respectively.	Gaze+	ETM(Ours)	83.4	72.9	62.5

Encoder?	Baseline			ETM (Ours)		
	Verb ↑	Noun ↑	Action ↑	Verb ↑	Noun ↑	Action ↑
1	19.9	20.4	7.2	21.5	20.5	8.1
	20.8	21.3	8.0	22.4	22.7	9.1
~	17.1	16.6	10.3	18.1	17.8	11.4
	18.2	17.5	11.1	19.9	19.1	12.9
~	42.1	37.6	28.9	43.4	38.9	31.3
	43.5	38.5	30.3	46.5	40.7	34.1