

Motion meets Attention: Video Motion Prompts

Qixiang Chen¹, Lei Wang^{1,2}, Piotr Koniusz^{2,1}, Tom Gedeon³

¹Australian National University ²Data61/CSIRO ³Curtin University

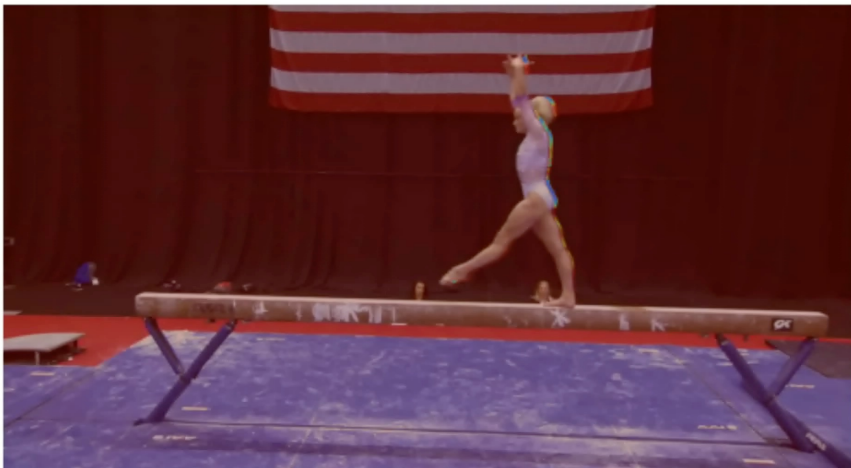
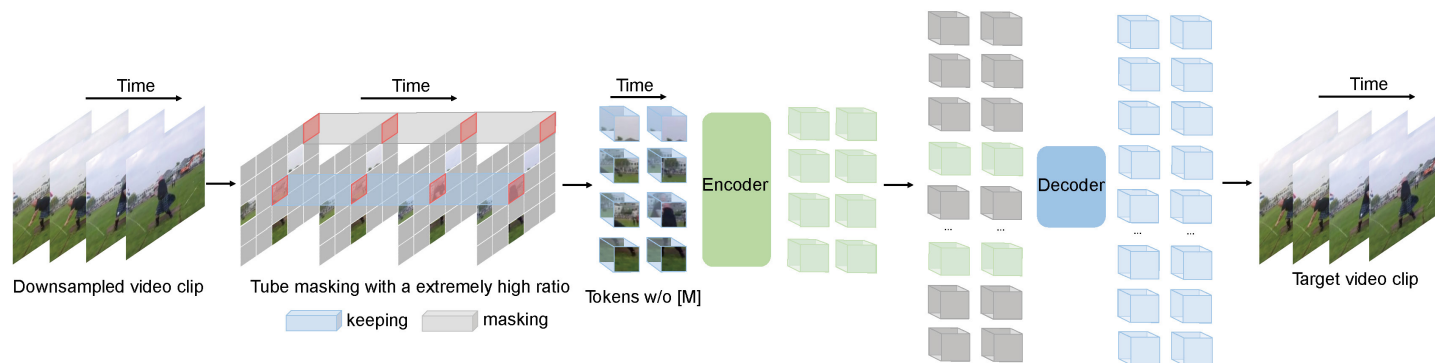


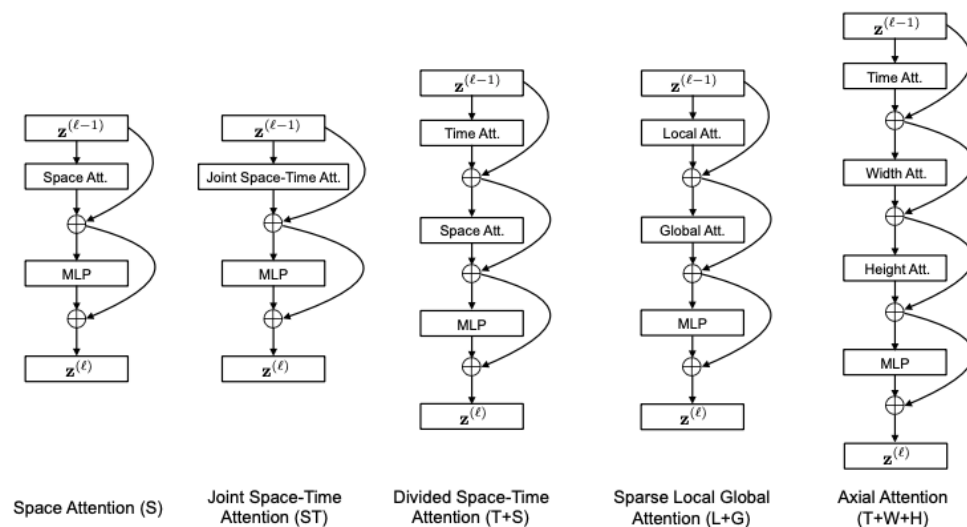
Table of content

- Background & Motivation
- Our Method
- Experiments & Discussions
- Conclusion & Future Work

Background & Motivation



VideoMAE¹



TimeSformer²

Background & Motivation

UCF-Crime: Fighting



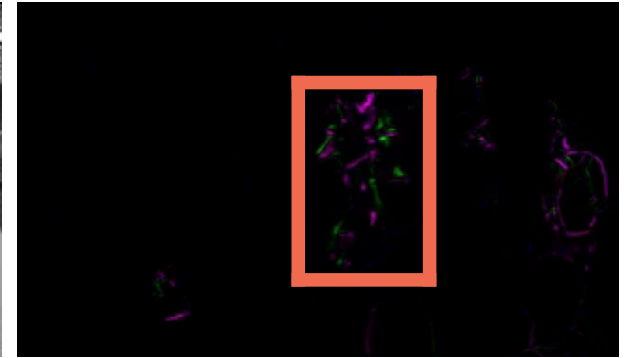
Original video



Normalized frame differencing map

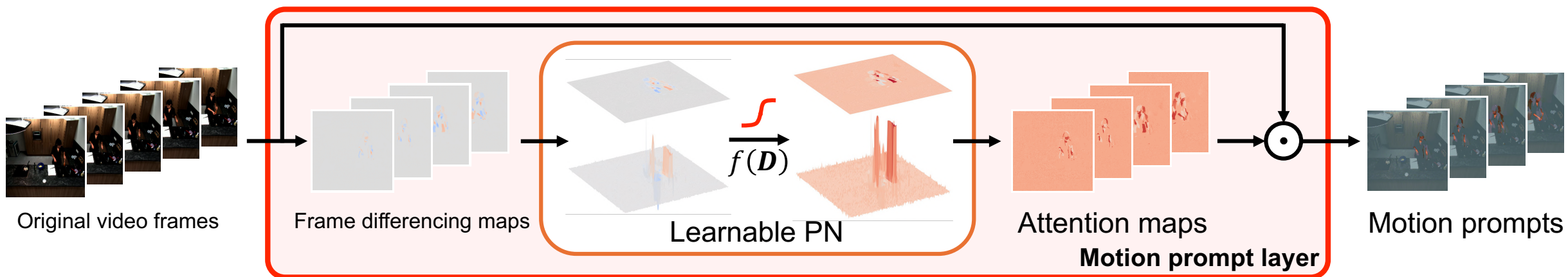


Time-color reordering frame



Taylor video frame

Method

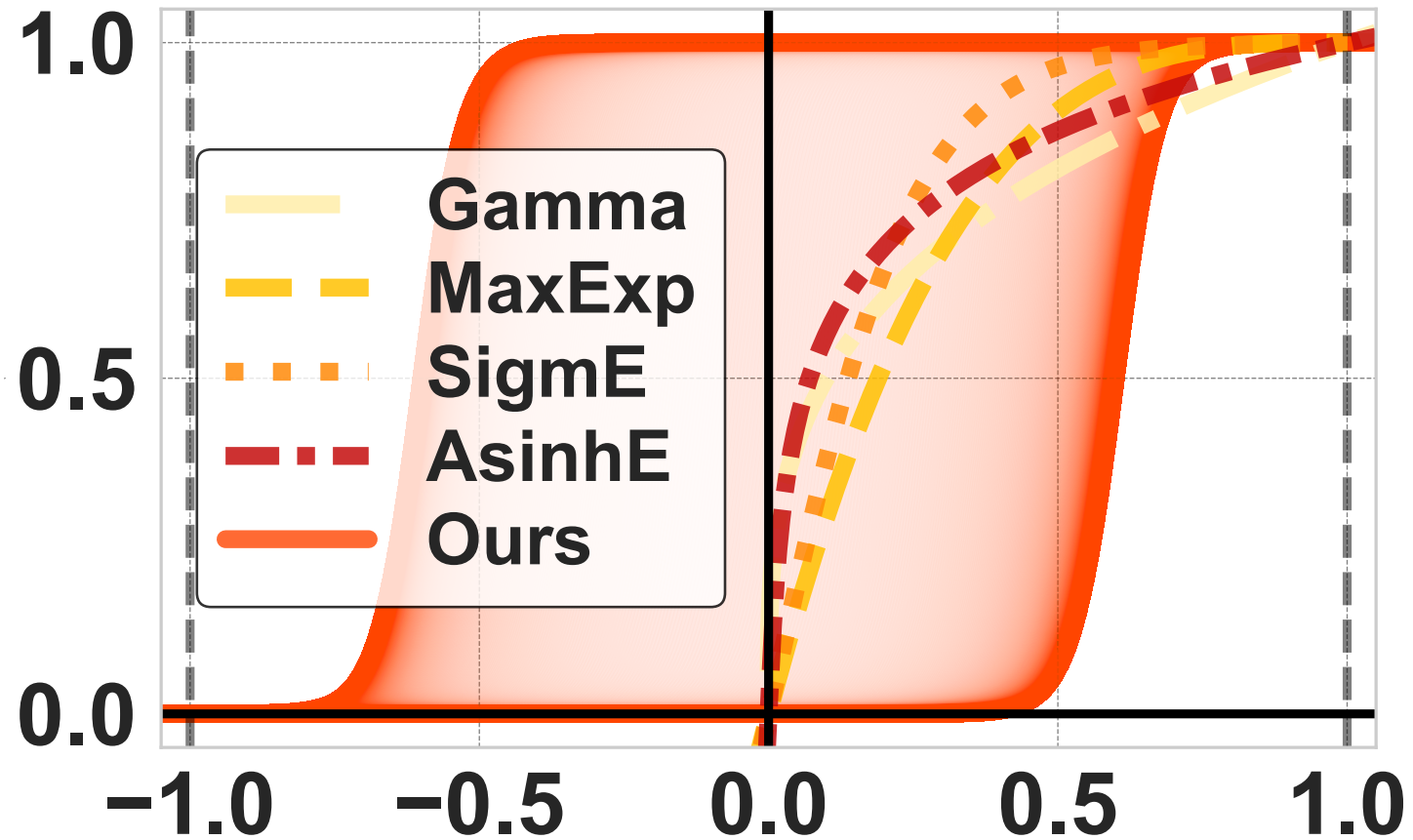


Method

$$f(\mathbf{D}) = \frac{1}{1 + e^{-a(\mathbf{D}-b)}}$$

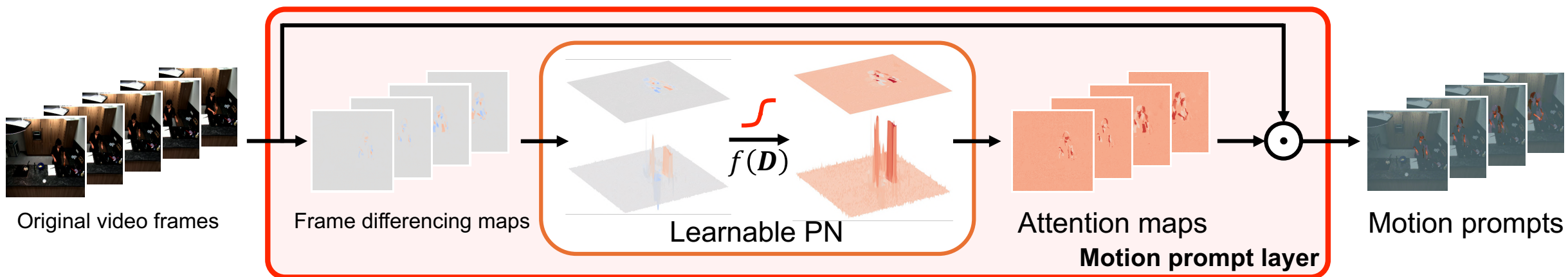
$$\begin{cases} a(m) = \frac{\alpha}{\beta |\tanh(m)| + \epsilon} \\ b(n) = \gamma \tanh(n) \end{cases}$$

Method



*Comparison of existing well-behaved Power Normalization (PN) functions Ko- niusz and Zhang (2021) and our learnable PN function

Method



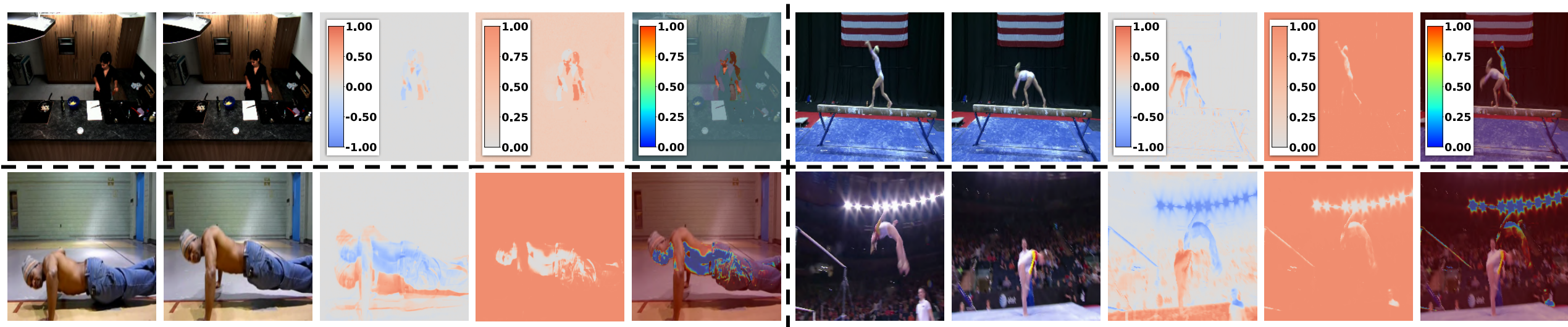
Method

Loss function: $\mathcal{L} = \mathcal{L}_{\text{ori}} + \lambda \mathcal{V},$

Temporal attention variation regularization:

$$\mathcal{V} = \frac{1}{T-2} \sum_{t=1}^{T-2} \|f(\mathbf{D}_{t+1}) - f(\mathbf{D}_t)\|_F^2,$$

Experiments & Discussion

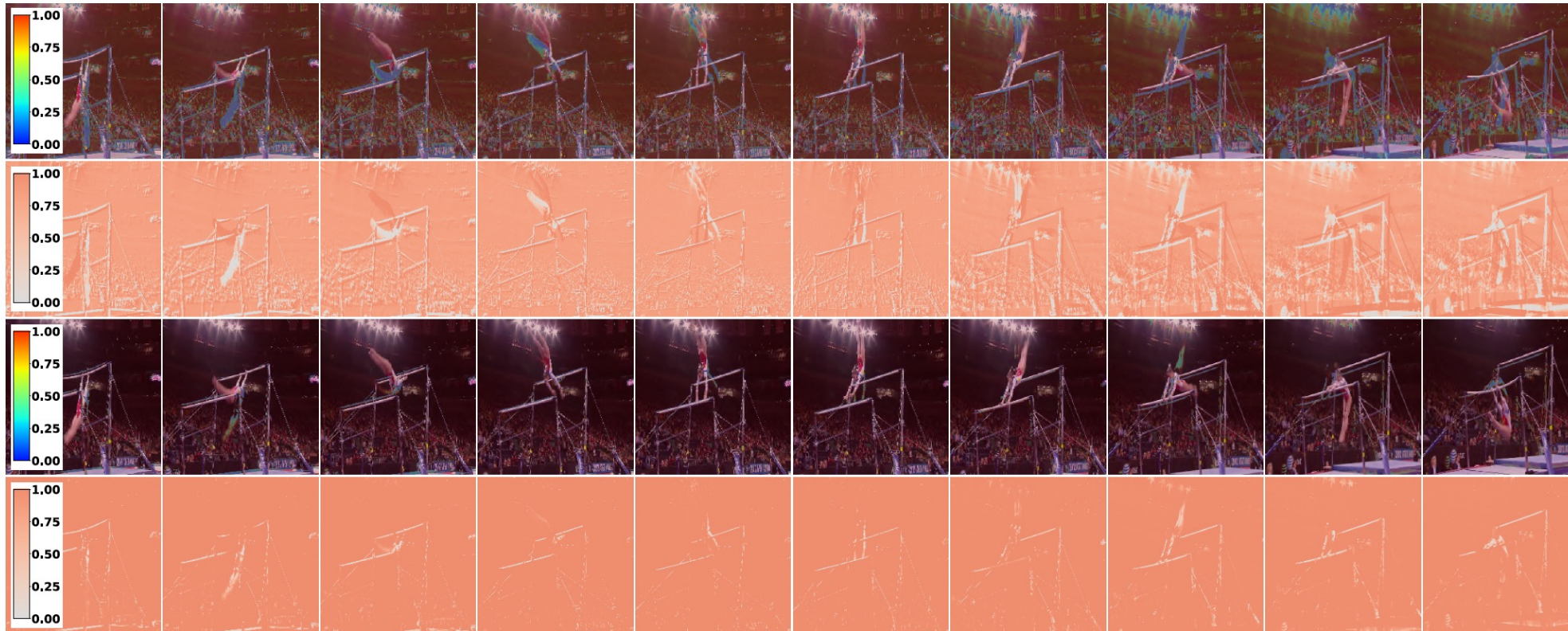


Experiments & Discussion

	[Pretrained❄]	[Pretrained❄ +VMPs🌟]	Baseline ([Pretrained❄🌟])	[Pretrained❄🌟 +VMPs🌟]	[Pretrained❄🌟❄ +VMPs🌟]	[Pretrained❄🌟❄🌟 +VMPs🌟]
Top-1	36.6	37.1 ↑0.5	50.6	56.6 ↑6.0	56.2 ↑5.6	57.1 ↑6.5
Top-5	66.9	66.2↓0.7	81.8	84.4 ↑2.6	84.3 ↑2.5	83.7 ↑1.9

Table1: Variant study of finetuning on MPlI Cooking 2 using TimeSformer.

Experiments & Discussion



$\lambda = 0$

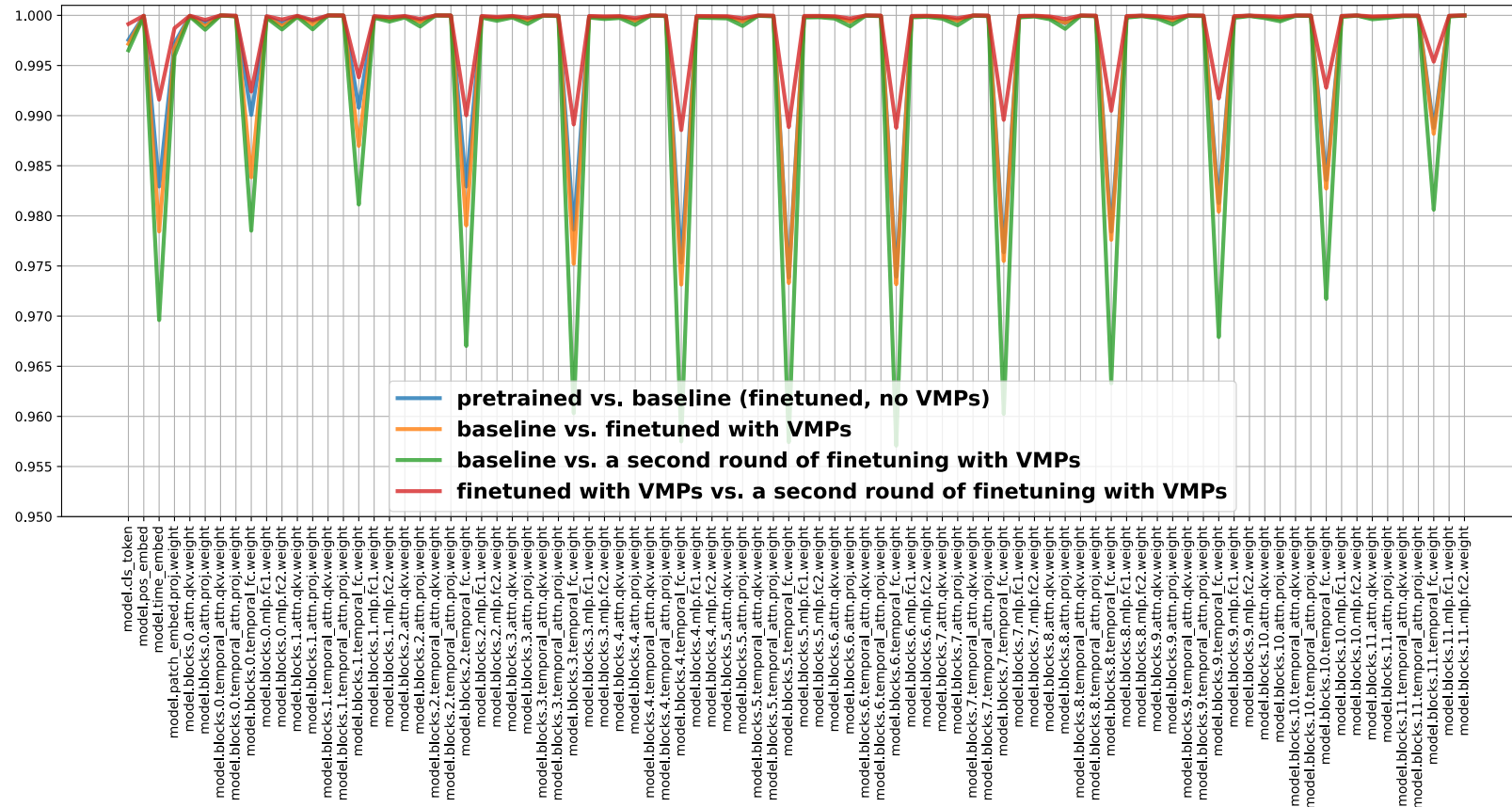
$\lambda = 2.5$

Experiments & Discussion

Table 2: Evaluations are conducted on (*left*) HMDB-51, and (*right*) FineGym, MPII Cooking 2, using SlowFast, X3D and TimeSformer as backbones. For SlowFast, we explore three variants by adding motion prompts into the slow-only stream, fast-only stream, and both slow and fast streams. We highlight improvements in red.

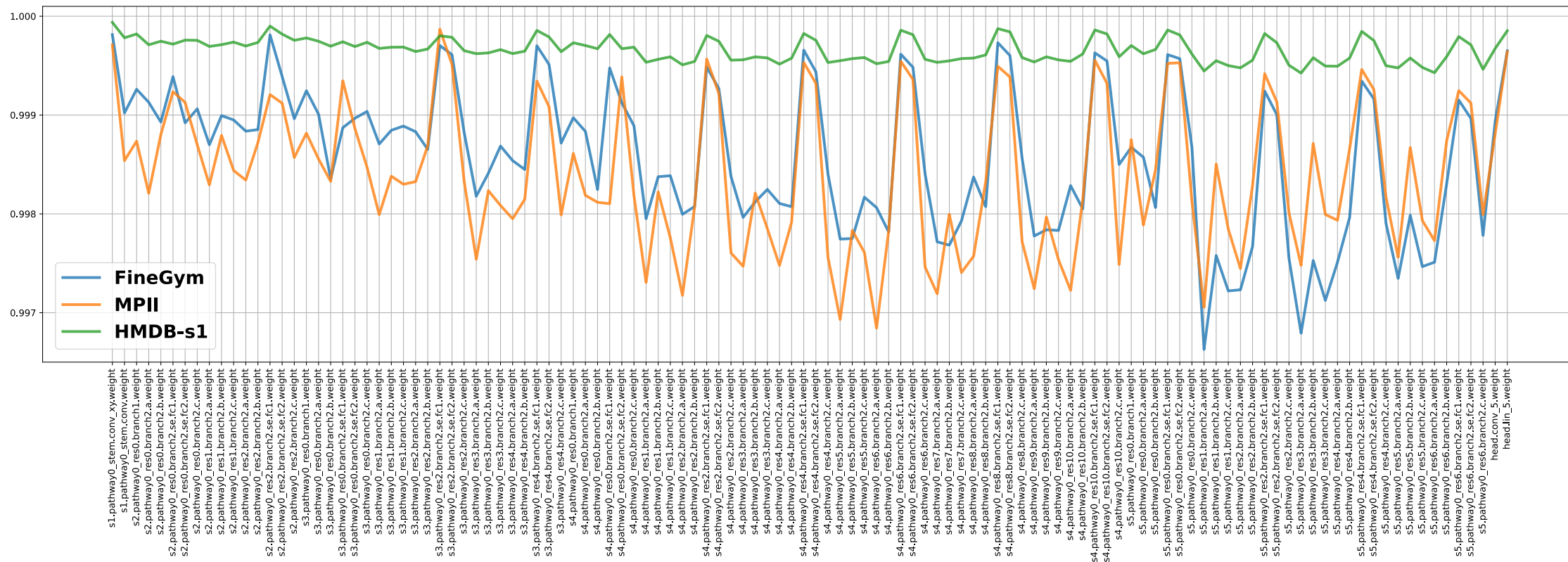
Model	HMDB-51				Mean	Model	FineGym		MPII Cooking 2	
	Split 1	Split 2	Split 3	Top-1			Top-5	Top-1	Top-5	
SlowFast	75.4	76.2	76.9	76.2	SlowFast	89.8	99.2	52.9	86.1	
+VMPs (slow-only)	76.8 ^{↑1.4}	77.0 ^{↑0.8}	77.3 ^{↑0.4}	77.0 ^{↑0.8}	+VMPs (slow-only)	89.7 ^{↓0.1}	99.2	55.5 ^{↑2.6}	84.5 ^{↓1.6}	
+VMPs (fast-only)	76.5 ^{↑1.1}	77.4 ^{↑1.2}	77.1 ^{↑0.2}	77.0 ^{↑0.8}	+VMPs (fast-only)	90.3 ^{↑0.5}	99.3 ^{↑0.1}	55.2 ^{↑2.3}	84.0 ^{↓2.1}	
+VMPs (slow&fast)	76.2 ^{↑0.8}	76.7 ^{↑0.5}	77.1 ^{↑0.2}	76.6 ^{↑0.4}	+VMPs (slow&fast)	90.1 ^{↑0.3}	99.3 ^{↑0.1}	56.8 ^{↑3.9}	86.6 ^{↑0.5}	
X3D	75.0	72.6	73.4	73.7	X3D	83.0	98.4	48.4	80.8	
+VMPs	75.8 ^{↑0.8}	73.2 ^{↑0.6}	73.6 ^{↑0.2}	74.2 ^{↑0.5}	+VMPs	83.8 ^{↑0.8}	98.6 ^{↑0.2}	49.1 ^{↑0.7}	80.6 ^{↓0.2}	
TimeSformer	72.7	73.1	72.2	72.7	TimeSformer	83.6	98.7	50.6	81.5	
+VMPs	74.2 ^{↑1.5}	74.3 ^{↑1.2}	72.9 ^{↑0.7}	73.8 ^{↑1.1}	+VMPs	84.4 ^{↑0.8}	98.5 ^{↓0.2}	56.6 ^{↑6.0}	84.4 ^{↑2.9}	

Experiments & Discussion



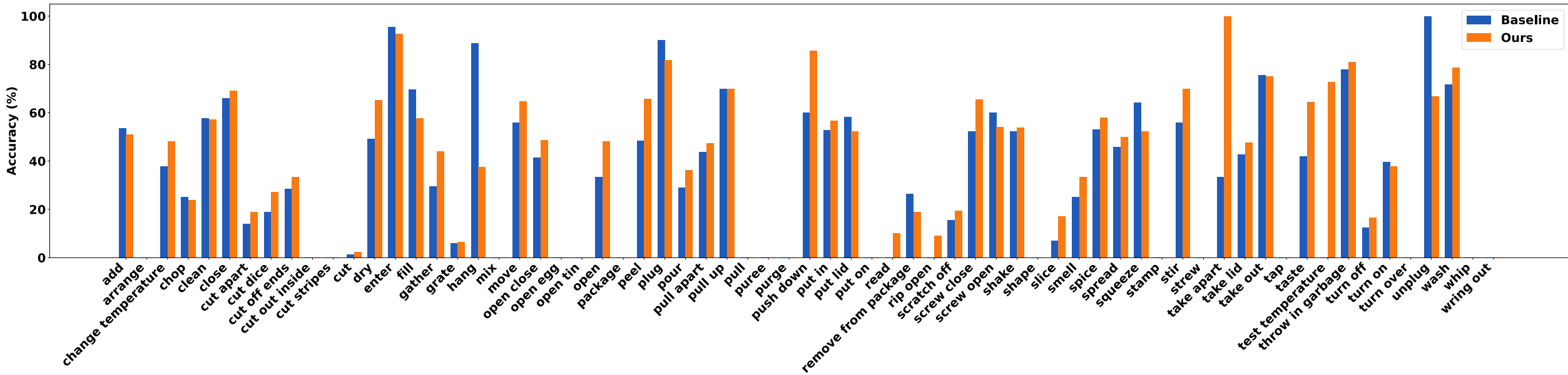
*Roles of VMPs in model finetuning via per-layer weight similarity comparison. We use TimeSformer pretrained on Kinetics-600 as the backbone, and finetuned on MPlI Cooking 2 with or without VMPs.

Experiments & Discussion



*Roles of VMPs in model finetuning via per-layer weight similarity comparison. We use X3D pretrained on Kinetics-600 as the backbone and finetune it on FineGym, MPII Cooking 2 (MPII), and HMDB-51 split 1 (HMDB-s1).

Experiments & Discussion



*Per-class accuracy comparison is conducted between the baseline model (pretrained on Kinetics-600 and then finetuned on MPII Cooking 2, without VMPs) and our VMP-enhanced model on MPII Cooking 2, using TimeSformer as the backbone.

Conclusion & Future Work