

Australian National University

Understanding the Behaviors of Modern Video Models Through Signal Perception



Author: Yifan Chen; Supervisor: Dr Lei Wang





Backbone I3D TimeSformer (Divided) TimeSformer (Joint) VideoMAE (V1)	Pretraining		Experiments			
I3D TimeSformer (Divided) TimeSformer (Joint) VideoMAE (V1)		Fine-tuning/testing	Top-1	Top-5	Reference Top-1	Reference Top-5
TimeSformer (Divided) TimeSformer (Joint) VideoMAE (V1)		HMDB51 UCE101	71.44 94 92	93.20 99.50	74.30	_
TimeSformer (Joint) VideoMAE (V1)	- - Kinetics-400 -	HMDB51	71.37	99.30	-	
TimeSformer (Joint) VideoMAE (V1)		UCF101 HMDB51	<u>95.82</u> 68.17	99.68		
VideoMAE (V1)		UCF101	93.71	99.58	-	_
		UCF101	72.22 95.32	92.48 99.37	73.30 96.10	_
VideoMAE (V2)	Kinetics-710	HMDB51 UCE101	79.35	96.86 00.02	_	_
UniFormer (V1)	Kinetics-400	HMDB51	76.41	99.92	_	
	CLIP-400M	UCF101 HMDB51	97.17	<u>99.71</u> 95.42		
UniFormer (V2)	+ Kinetics-710	UCF101	96.51	99.76	-	_
Swin	Kinetics-400	HMDB51 UCF101	75.16 96.11	93.46 99.50	_	
³ -1.0 -1.5 -2.0 -2.0 -2.5 -2.5 -3.0 0.0 0.2 0.4 0.2 0.4 0.4 0.2 0.4 Normalized Data (a) Figure 1: Temporal F he HMDB51 datase outputs from all MSA he CNN-based mode nodels, and from all MLP blocks of the tra- Noise Ratio pinology	Divided) Joint) Joint) Sourier analysis t, with models A and MLP blocks of ansformer-base (%) 60 Yourger analysis t, with models A and MLP blocks of ansformer-base	-1.0 -1.5 -1.5 -2.0 -2.5 -2.5 -3.0 -1.5 -2.5 -3.0 -1.5 -2.5 -3.0 -1.5 -2.5 -3.0 -1.5 -2.5 -1.5 -2.5 -1.5 -2.5 -1.5 -2.5 -1.5 -1.5 -2.5 -1.5 -2.5 -1.5	er (Divided) er (Joint) (V1) (V2) (V1) (V2) 0.6 0.8 0 Depth 0.6 0.8 0 Depth 0	nediate om left sed mod ll MSA (c) the Conv blo	Part of the transformation of the transforma	Sformer (Divided) Sformer (Joint) OMAE (V1) OMAE (V2) Ormer (V1) Ormer (V2) 4 0.6 0.8 Depth (C) during testing (e analysis of the ll Conv blocks ansformer-base outputs from a N-based model Baseline Top-1 = 71.3 Baseline Top-5 = 93.8 2 Layer Ratio = 50 71.6 93.8 54.1 93.8 54.8 54.8 54.8 55.8
	ussian noise w	ith a consistent la	ayer rati	o in rom	a activaly applie	
Figure 2: Random ga Veight Matrix of the eft to right.	spatial attention 80 % 60 Ser Layer Noise S	JI Tayers III Divit	ded Spa	o.10,15,20,25,00	e TimeSformer	ed to Q, K, and visualized from Baseline Top-1 = 71.3 Baseline Top-5 = 93.8 Layer Ratio = 50 71.4 54.1 $\frac{1}{2}$ 36.9 $\frac{7}{2}$ 19.6 $\frac{1}{2}$ 32.9 12.6

Figure 4: Random Gaussian noise with a consistent noise standard deviation is applied to the spatial attention layers of Divided Space-Time TimeSformer. The model's layers are sequentially divided into three equal parts (front, middle, and rear), and selected from each part respectively, visualized from left to right.









