Hierarchical Temporal Transformer for 3D Hand Pose Estimation and Action Recognition from Egocentric RGB Videos

Task

Unified framework for both 3D hand pose estimation and a egocentric RGB video.

Key Designs



Leverage the temporal information for both pose an

- frequent self-occlusions between hands and objects.
- severe ambiguity of action types judged from individual fra

Build a hierarchical temporal transformer with two

- leverage different time spans for pose and action estimation
- model the semantic correlation by deriving the high-lev hand motion and manipulated object label.

Result on H2O

Action Recognition							3D Hand Pose Estimation (Camera Space)					
	C2D	חכו	SlowFast	H+O	H2O w/ ST-GCN	H2O w/ TA-GCN	Ours	MEPE (in <i>mm</i>)	H+O	LPC	H2O	Ours
	620	עטו						Left	41.42	39.56	41.45	35.02
Acc.	70.66	75.21	77.69	68.88	73.86	79.25	86.36	Right	38.86	41.87	37.21	35.63

Result on FPHA

Action Recognition						
	Joule- color	Two- Stream	H+O	Collaborative	Ours	
Acc.	66.78	75.30	82.43	85.22	94.09	

[†]Work partially done during internships with Microsoft Research Asia.

Yilin Wen^{1,†}, Hao Pan², Lei Yang^{3,1}, Jia Pan¹, Taku Komura¹, Wenping Wang⁴ ¹The University of Hong Kong ²Microsoft Research Asia ⁴Texas A&M University ³TransGP

	Framework				
action recognition from the	Hand Pose Estimation with Sh				
	Pose block P focuses on a narrow manipulated object label.				
pla	Input video is divided into consec segments are processed by P in				
Ce m	Action Recognition with Long				
	> Action block <i>A</i> uses the full video				
nd action. to cope with:	\succ The input of <i>A</i> leverages the per-				
	Segmentation Strategy to Div				
ames.	Videos longer than T frames are of with a window size T:				
nation.	In testing stage, the hand pos among segmented clips.				
el action from the low-level	To augment training data, the				



hort-Term Temporal Cue

wer temporal receptive field to output per-frame 3D hand pose and

Hand Pose

cutive segments by a shifting window strategy with window size t, parallel.

J-Term Temporal Cue

to predict the action label.

frame predicted hand pose, object label and image feature.

ide Long Videos into HTT Inputs

divided into consecutive clips, by adopting the shifting window strategy

se are estimated by **P**, the action category is voted from the predictions

e starting frame for shifting window is offset to each of the first t frames



