



Kinetics-TPS Track on Part-level Action Parsing and Action Recognition Tech Report

Zheming Yu¹ Lin Li² Jietian Guo³

Hikvision Research Institute



Challenge introduction



belly dancing



hopscotch







clean_and_jerk



lunge



push_up



deadlifting



punching_bag



throwing_discus

needs to predict:

- human location
- body part location
- part state in the frame level
- human action in the video level

ICCV 2021 Workshop

Deeper Action

Dataset Statistics

Kinetics-TPS contains 4741 videos:

- 1. Bounding boxes of human instances: 1.6 M
- 2. Bounding boxes of body parts: 7.9 M
- 3. Part state tags of each annotated part: 7.9 M
- 4. Bounding boxes and tags of objects: 0.5 M
- 5. 'body part, part state' pairs of four exemplar classes in the





ICCV 2021

Workshop







• Our Method

This Challenge aims at locating the human location, body part location and their part states simultaneously in the frame level. There are four modules to predict the results.



The framework of the proposed method

Deeper Action

 \succ



Best result is cascade-RCNN with Swin-L backbone and

ICCV 2021 Workshop

		method	Backbone	Dataset size	mAP(%)	
		Retinanet	Resnet-50	5k	25	
n	mAP(%)	GFL	ResNeXt-101	5k	54	
	93	GFL	ResNeXt-101	50k	56.6	
		Cascade-rcnn	Swin-L	50k	57.1	

Action recognition

Human detection

Method

Cascade-RCNN

Method	Pre-train	ACC(%)
Video Swin Transformer	SSv2	93

Backbone

Resnet-50

Pre-trai

COCO

- Human Body Parts Detection
- We adopts a top-down strategy to decompose multiperson problems into single-person problems;

The risk of pseudo-label

trained on larger dataset.

Method	Backbone	pseudo-label	mAP(%)
Cascade-rcnn	Swin-L	×	57.1
Cascade-rcnn	Swin-L	\checkmark	56.7

Deeper

ICCV 2021 Workshop

- Module results
- Human Part State Recognition



Human part state recognition structure

- Use the pre-trained Transformer network as the backbone, and perform ROL Pooling to obtain part features;
- Connate human body feature on each part specific to obtain global context;
- Use Multi-Head Attention to fuse information from body parts;
- Add action category information to the network.

Augme ntation	Action info	ROL Pooling	Multi- Head attention	Top1 (%)
\checkmark				72.4
\checkmark	\checkmark			73.6
	\checkmark	\checkmark		74.2
	\checkmark	\checkmark	\checkmark	74.9



- Prospects
 - Incorporating middle level information to improve the action recognition performance.
 - Incorporating stable and accurate temporal information to improve the performance of human body part localization and human body part state prediction.
 - Incorporating pre-trained contrastive Language-Video Model, like ActionCLIP, and open concepts.





Thanks