



# DailyDVS-200: A Comprehensive Benchmark Dataset for Event-Based Action Recognition

Qi Wang<sup>1\*</sup>, Zhou Xu<sup>1\*</sup>, Yuming Lin<sup>1</sup>, Jingtao Ye<sup>1</sup>, Hongsheng Li<sup>1</sup>, Guangming Zhu<sup>1</sup>,  
Syed Afaq Ali Shah<sup>2</sup>, Mohammed Bennamoun<sup>3</sup>, Liang Zhang<sup>1</sup>✉

<sup>1</sup>Xidian University, <sup>2</sup>Edith Cowan University, <sup>3</sup>University of Western Australia



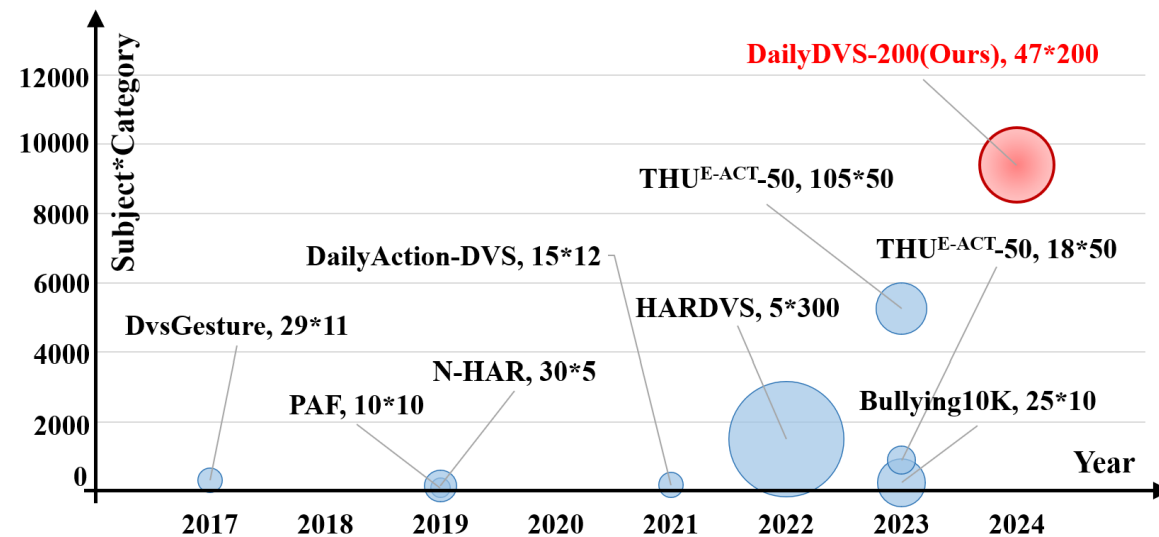
EUROPEAN CONFERENCE ON COMPUTER VISION

M I L A N O  
2 0 2 4

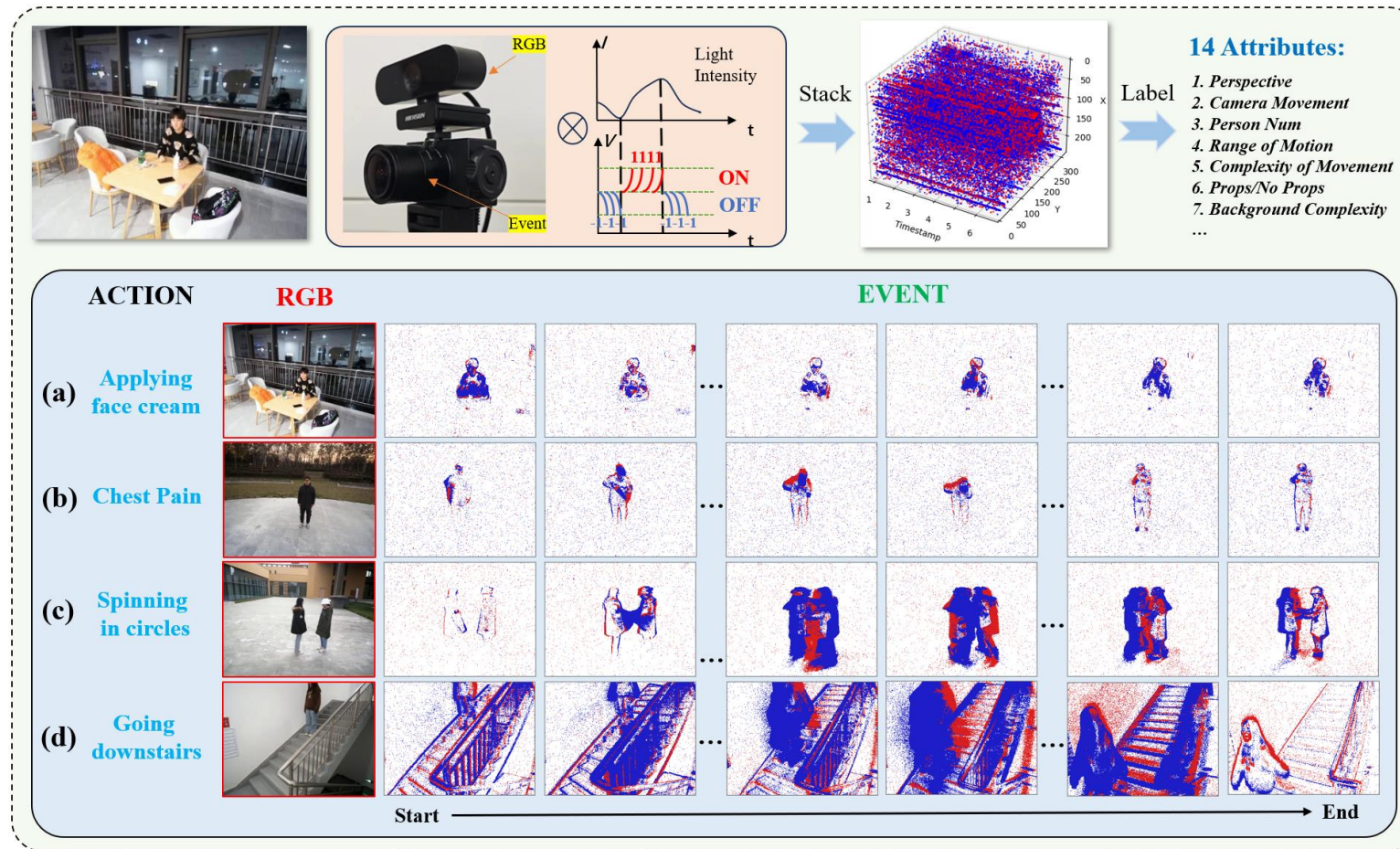
# Contributions

- We introduce DailyDVS-200, a pioneering large-scale neuromorphic dataset for action recognition, featuring over 22,000 samples across 200 categories with real-world challenges and detailed attribute annotations.
- We've set up benchmarks for diverse recognition models on this dataset, establishing a solid foundation for future performance comparisons in the field of neuromorphic action recognition.
- Our group evaluation and analysis, based on attribute annotations, not only assesses the influence of various attributes on event-based models but also paves the way for new research directions in neuromorphic datasets.

DailyDVS-200					
Scale	22046	14 attributes:			
Class	200	PR	Props	PE	Perspective
Subject	47	PO	Posture	DI	Diurnality
Sensors	DVXplorer Lite	DR	Duration	LO	Location
Resolution	320×240	AR	Action Range	DT	Distance
Duration	1-20s	PN	Person Num	HE	Height
EventCount	18 billion	CM	Camera Motion	SH	Shadow
EventTime	24 hour	IL	Illumination Direction	BC	Background Complexity



# Data acquisition process



- As we move forward, you'll see how these meticulously captured and annotated data points contribute to the robustness of our dataset. They form the backbone of our benchmarks and are instrumental in pushing the boundaries of neuromorphic action recognition.

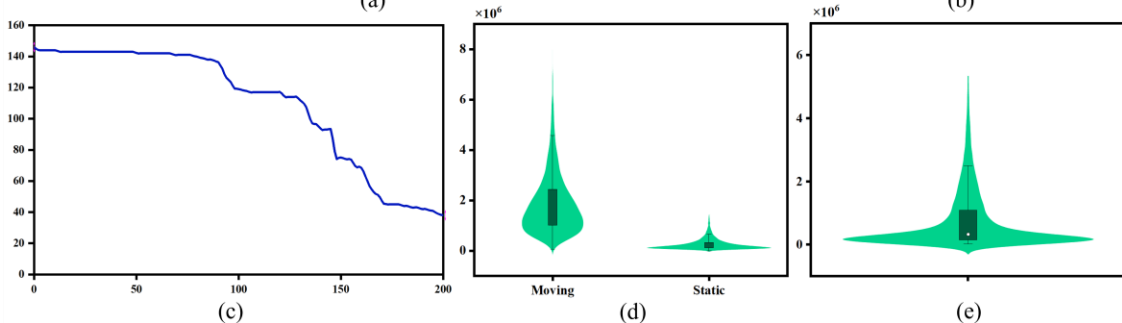
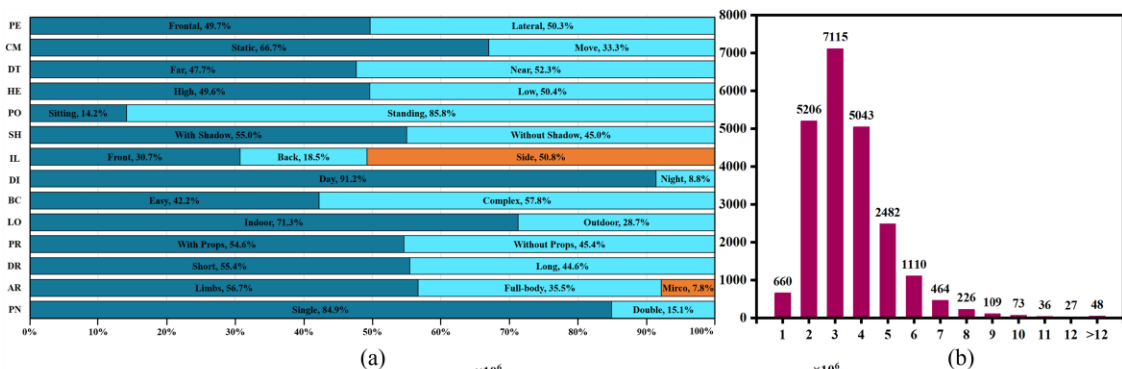
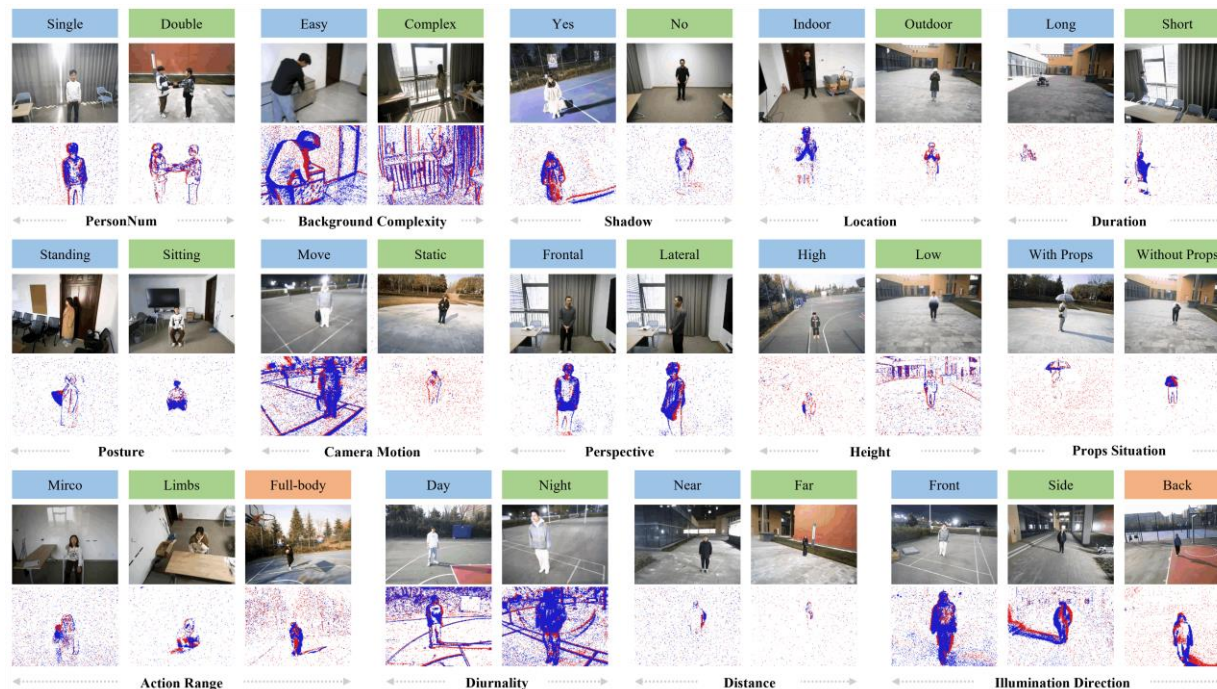
# DailyDVS-200: Event-Based Action Recognition Dataset

Dataset	Year	Sensors	Object	Scale	Class	Sub	Real	MA	AA	DR
ASLAN-DVS	2011	DAVIS240c	Action	3,697	432	-	✗	-	-	-
MNISTDVS	2013	DAVIS128	Image	30,000	10	-	✗	-	-	-
N-Caltech101	2015	ATIS	Image	8,709	101	-	✗	-	-	0.3s
N-MNIST	2015	ATIS	Image	70,000	10	-	✗	-	-	0.3s
CIFAR10-DVS	2017	DAVIS128	Image	10,000	10	-	✗	-	-	1.2s
HMDB-DVS	2019	DAVIS240c	Action	6,766	51	-	✗	-	-	19s
UCF-DVS	2019	DAVIS240c	Action	13,320	101	-	✗	-	-	25s
N-ImageNet	2021	Samsung-Gen3	Image	1,781,167	1,000	-	✗	-	-	-
ES-ImageNet	2021	-	Image	1,306,916	1,000	-	✗	-	-	-
DvsGesture	2017	DAVIS128	Action	1,342	11	29	✓	✗	✗	6s
N-CARS	2018	ATIS	Car	24,029	2	-	✓	✗	✗	0.1s
ASL-DVS	2019	DAVIS240	Hand	100,800	24	5	✓	✗	✗	0.1s
PAF	2019	DAVIS346	Action	450	10	10	✓	✗	✗	5s
DailyAction	2021	DAVIS346	Action	1,440	12	15	✓	✗	✗	5s
HARDVS	2022	DAVIS346	Action	107,646	300	5	✓	✓	✗	5s
THU <sup>E-ACT</sup> -50	2023	CeleX-V	Action	10,500	50	105	✓	✗	✗	2-5s
THU <sup>E--ACT</sup> -50-CHL	2023	DAVIS346	Action	2,330	50	18	✓	✗	✗	2-5s
Bullying10K	2023	DAVIS346	Action	10,000	10	25	✓	✗	✗	2-20s
DailyDVS-200 (Ours)	2024	DVXplorer Lite	Action	22,046	200	47	✓	✓	✓	1-20s

- Our proposed dataset consists of 200 distinct action categories, collected from 47 subjects. The dataset we provide has the richest number of subject-category combinations. The detailed comparison with existing benchmark datasets can be found in table.

# DailyDVS-200: Event-Based Action Recognition Dataset

- A preview of our proposed DailyDVS-200 dataset and examples of our attribute annotations.
- Statistical data and detailed analysis of DailyDVS-200.
- DailyDVS-200 dataset exemplifies diversity in action recognition.



## Various scene:

- Household Activities
- Office Tasks
- Sports and Physical Activities
- Health-related Activities
- Interactions
- Bullying and Violence
- Transportation-related Activities

## Various setting & Various action characteristics:

- Different Perspectives
- Different Distances
- Lighting Conditions
- Different Camera Movements
- Fine-grained Micro, Limb, and Whole-body Movements
- Short-duration and Long-duration
- Actions with and without Props Interaction

# Experimental Results

Methods	Year	Input Type	Backbone	top-1 acc.(%)	top-5 acc.(%)
C3D	2015	Frame	3D CNN	21.99	45.81
I3D	2017	Frame	ResNet50	32.30	59.05
R2Plus1D	2018	Frame	ResNet34	36.06	63.67
SlowFast	2019	Frame	ResNet50	<b>41.49</b>	68.19
TSM	2019	Frame	ResNet50	40.87	<b>71.46</b>
EST	2019	Learnable	ResNet34	<b>32.23</b>	<b>59.66</b>
TimeSformer	2021	Token	Transformer	44.25	74.03
Swin-T	2022	Token	Transformer	<b>48.06</b>	<b>74.47</b>
ESTF	2022	Token	ResNet18	24.68	50.18
GET	2023	Token	Transformer	37.28	61.59
Spikformer	2022	Spike	Transformer	<b>36.94</b>	<b>62.37</b>
SDT	2024	Spike	Transformer	35.43	58.81

- Evaluation of different methods on our dataset. Swin Transformer achieved the highest accuracy in both top-1 and top-5, with 48.06% and 74.47% respectively, and also achieved the highest accuracy among all methods. In the spike-based methods, Spikformer achieved the highest top-1 and top-5 accuracy of 36.94% and 62.37%, respectively.

# Experimental Results

Methods	1. BC		2. DR		3. IL			4. DT		5. HE		6. LO		7. CM	
	Complex	Easy	Long	Short	Back	Front	Side	Far	Near	High	Low	Indoor	Outdoor	Move	Static
SlowFast	<u>45.98</u>	39.20	<u>40.94</u>	42.18	<u>44.05</u>	34.52	46.13	<u>41.62</u>	41.64	38.34	45.36	43.21	39.08	21.05	<u>51.81</u>
TSM	45.37	<u>40.84</u>	39.68	<u>44.69</u>	42.91	<u>36.60</u>	<u>47.16</u>	40.71	<u>44.09</u>	<u>38.90</u>	<u>46.51</u>	<u>43.41</u>	<u>40.93</u>	<u>26.88</u>	50.16
EST	35.22	30.59	29.42	34.52	32.89	24.71	38.26	28.93	35.28	27.45	37.64	33.76	29.82	13.52	41.51
Spikformer	38.96	33.45	33.86	36.68	36.77	30.09	39.11	31.57	39.00	31.82	39.52	38.54	30.33	13.74	46.15
Swin-T	<b>52.86</b>	<b>45.37</b>	<b>48.57</b>	<b>47.64</b>	<b>51.61</b>	<b>39.11</b>	<b>53.39</b>	<b>47.16</b>	<b>48.89</b>	<b>44.46</b>	<b>52.14</b>	<b>50.15</b>	<b>44.70</b>	<b>27.84</b>	<b>58.05</b>

Methods	8. PN		9. PE		10. PR		11. AR			12. SH		13. PO		14. DI	
	One	Two	Frontal	Lateral	No	Yes	Full-body	Limbs	Micro	No	Yes	Stand	Sit	Day	Night
SlowFast	39.15	55.57	44.52	39.04	41.81	<u>41.48</u>	52.05	<u>37.59</u>	25.16	37.07	<u>42.73</u>	42.93	35.92	43.10	<u>33.65</u>
TSM	<u>40.07</u>	<u>55.90</u>	<u>46.13</u>	<u>39.17</u>	<u>46.66</u>	38.98	<u>53.31</u>	37.21	<u>33.02</u>	<u>42.75</u>	42.39	<u>43.11</u>	<u>39.61</u>	<u>44.20</u>	33.02
EST	29.88	45.56	33.42	31.15	34.81	30.08	41.18	28.55	19.81	35.18	31.55	32.70	30.13	34.74	18.55
Spikformer	33.88	44.10	37.09	33.94	37.39	33.80	43.86	31.95	23.27	39.47	34.45	35.01	37.24	38.30	19.81
Swin-T	<b>44.70</b>	<b>66.88</b>	<b>50.98</b>	<b>45.43</b>	<b>51.35</b>	<b>45.33</b>	<b>59.24</b>	<b>43.15</b>	<b>34.59</b>	<b>46.78</b>	<b>48.36</b>	<b>48.96</b>	<b>44.08</b>	<b>50.22</b>	<b>36.32</b>

Methods	0.5s			0.25s			0.125s		
	gap0	gap2	gap4	gap0	gap2	gap4	gap0	gap2	gap4
C3D [53]	21.99	14.98	11.70	31.10	24.82	22.82	44.81	32.62	27.75
I3D [9]	32.30	22.94	20.82	45.39	29.54	29.42	59.10	36.70	37.50
R2Plus1D [54]	36.06	<u>26.39</u>	<b>24.97</b>	<u>49.65</u>	<u>36.62</u>	<u>32.10</u>	58.88	<u>48.06</u>	40.29
SlowFast [17]	<b>41.49</b>	<b>26.90</b>	<u>24.55</u>	<b>52.16</b>	33.28	25.43	<b>64.09</b>	44.81	<u>44.64</u>
TSM [31]	<u>40.87</u>	23.67	22.97	49.55	<b>37.94</b>	<b>32.37</b>	<u>61.76</u>	<b>51.48</b>	<b>48.77</b>

- The fine-grained group testing results on models trained on our dataset.

- Evaluation of existing models with different frame settings on our dataset.

# Experimental Results

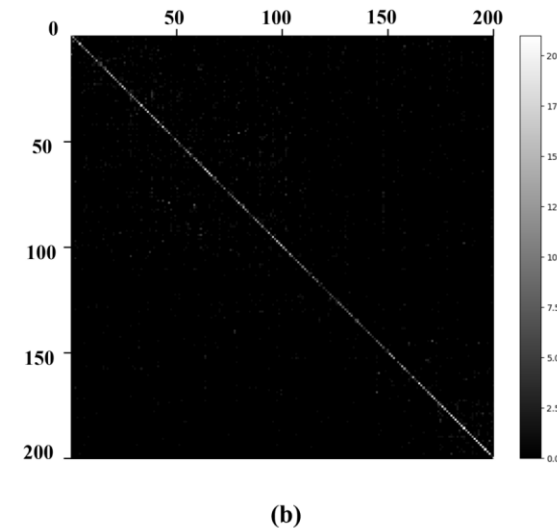
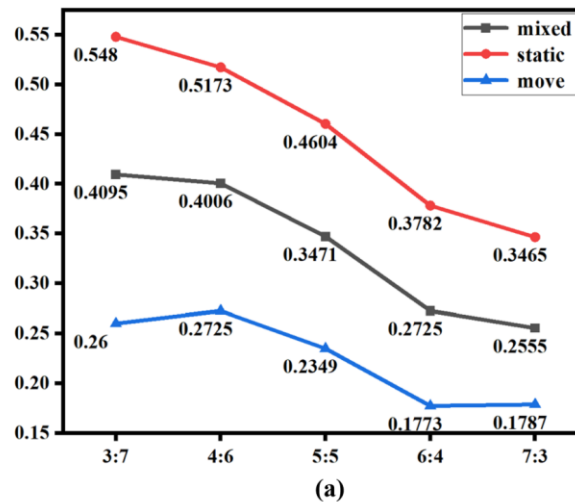
Methods	Top-10 accurate actions	Top-10 incorrect actions	Methods	Top-10 accurate actions	Top-10 incorrect actions
Slowfast	<ol style="list-style-type: none"> <li>pull out the chair</li> <li>open curtains</li> <li><b>hand in hand circling</b></li> <li>close the door</li> <li>wash towels</li> <li>fall down</li> <li>turn the light on</li> <li><b>wipe the table</b></li> <li>cross your legs</li> <li><b>push-up</b></li> </ol>	<ol style="list-style-type: none"> <li>hammer table</li> <li><b>play table tennis</b></li> <li>open window</li> <li>clean windows</li> <li><b>pitch</b></li> <li>charge a phone</li> <li>arrang cards</li> <li>V sign</li> <li>use a tablet</li> <li>kick a ball</li> </ol>	EST	<ol style="list-style-type: none"> <li>tie shoelaces</li> <li>mutual bow</li> <li>turn the light on</li> <li><b>push-up</b></li> <li>turn the light off</li> <li>fall down</li> <li>lie on the table</li> <li>close window</li> <li>cheers</li> <li><b>hand in hand circling</b></li> </ol>	<ol style="list-style-type: none"> <li><b>play table tennis</b></li> <li>open window</li> <li>close curtains</li> <li>slap the table</li> <li><b>pitch</b></li> <li>OK sign</li> <li>make paper cuttings</li> <li>charge a phone</li> <li>hit people with things</li> <li>headache</li> </ol>
Swin-T	<ol style="list-style-type: none"> <li><b>hand in hand circling</b></li> <li>arm wrestling</li> <li><b>push-up</b></li> <li>close curtains</li> <li>open curtains</li> <li>close the door</li> <li>sit-up</li> <li><b>wipe the table</b></li> <li>moves heavy objects</li> <li>cross legs</li> </ol>	<ol style="list-style-type: none"> <li>hammer table</li> <li><b>pitch</b></li> <li>take off headphones</li> <li>write</li> <li>blow nose</li> <li>V sign</li> <li>crush paper into a ball</li> <li>take something from bag</li> <li>chest pain</li> <li>open the bottle</li> </ol>	Spikformer	<ol style="list-style-type: none"> <li><b>wipe the table</b></li> <li>close curtains</li> <li>stand up</li> <li><b>hand in hand circling</b></li> <li>lie on the table</li> <li>mutual bow</li> <li>cross legs</li> <li>clean the windows</li> <li><b>push-up</b></li> <li>go upstairs</li> </ol>	<ol style="list-style-type: none"> <li><b>pitch</b></li> <li><b>play table tennis</b></li> <li>plug in the power strip</li> <li>take off shoes</li> <li>trim nails</li> <li>stomachache</li> <li>play with hair</li> <li>roll up sleeves</li> <li>backache</li> <li>headache</li> </ol>

- Top-10 accurate and top-10 incorrect actions of different methods. The same action is marked with the same color.



# Experimental Results

- Evaluation of using different sizes of Moving camera set for action recognition and Confusion matrix of swin transformer.



- Comparison with other large-scale datasets under the same training configuration.

Methods	TSM	ESTF
THU <sup>E-ACT</sup> -50	95.60 / 98.75 <sup>†</sup>	95.25
THU <sup>E-ACT</sup> -50-CHL	49.07 / 83.83 <sup>†</sup>	49.50
Hardvs	97.33 / 98.55 <sup>†</sup>	96.67
Bullying10K	74.22 / 91.90 <sup>†</sup>	84.72
DailyDVS-200( <b>Ours</b> )	<b>36.05 / 65.90<sup>†</sup></b>	<b>31.29</b>

# Thanks



**Code & Dataset:**

[github.com/QiWang233/DailyDVS-200](https://github.com/QiWang233/DailyDVS-200)