

# Yixuan Li

CONTACT INFORMATION	163 Xianlin Road, Qixia District Nanjing, Jiangsu, P.R. China, 210023	(+86) 18120190595 <a href="mailto:yixuanli@smail.nju.edu.cn">yixuanli@smail.nju.edu.cn</a>
HOME PAGE	<a href="https://yixuanli98.github.io/">https://yixuanli98.github.io/</a>	
RESEARCH INTERESTS	<b>Computer Vision:</b> spatio-temporal action detection, action recognition.	
EDUCATION	<b>Department of Computer Science and Technology, Nanjing University</b> M.Sc. Candidate in <b>MCG Lab</b> Supervisor: <b>Prof. Limin Wang</b> <b>Kuang Yaming Honors School, Nanjing University</b> B.Sc., Major in Computer Science (GPA: 86.2/100) Supervisor: <b>Prof. Gangshan Wu</b>	Nanjing, China August 2019 – Present  Nanjing, China August 2015 – June 2019
PUBLICATION	<b>Yixuan Li*</b> , Zixu Wang*, Limin Wang, Gangshan Wu. Actions as Moving Points. European Conference on Computer Vision (ECCV'20), Glasgow, United Kingdom, 2020. <b>Yixuan Li</b> , Lei Chen, Runyu He, Zhenzhi Wang, Gangshan Wu, Limin Wang. MultiSports: A Multi-Person Video Dataset of Spatio-Temporally Localized Sports Actions. International Conference on Computer Vision (ICCV'21), 2021.	
RESEARCH EXPERIENCE	<b>MultiSports: A Multi-Person Video Dataset of Spatio-Temporally Localized Sports Actions.</b> Advisor: <b>Prof. Limin Wang</b> HomePage: <a href="#">MultiSports Dataset</a> <ul style="list-style-type: none"><li>As the first author, presented a large-scale, fine-grained, multi-person, and untrimmed spatio-temporal action detection dataset with well-defined temporal boundaries, <i>MultiSports</i>. Besides, adapted several representative methods to it and gave in-depth analysis to inspire new advances in this field.</li><li><i>MultiSports</i> contained 66 fine-grained action categories from 4 different sports, where we collected 3200 video clips and annotated around 37790 action instances with 907k bounding boxes. Our datasets had more fine-grained action categories (66 vs. 21 or 24), more instances per video clip (11.8 vs. 1 or 1.4) and much more instances (37790 vs. 928 or 4458) than the existing datasets JHMDB and UCF101-24.</li><li>Existing methods achieved satisfactory performance on JHMDB and UCF101-24 but obtained low performance on <i>MultiSports</i> (video-mAP@0.2 of 77.3 or 82.8 vs. 12.88 for MOC).</li></ul> <b>Actions as Moving Points</b> Advisor: <b>Prof. Limin Wang</b> <ul style="list-style-type: none"><li>As the first author, presented an conceptually simple, computationally efficient, and more precise spatial-temporal action detection framework, MOC-detector, which would recognize all the action instances present in a video and localize them in both space and time.</li><li>MOC outperformed the existing state-of-the-art methods under the same setting on the JHMDB and UCF101-24 datasets. The code is available at <a href="https://github.com/MCG-NJU/MOC-Detector">https://github.com/MCG-NJU/MOC-Detector</a>.</li><li>MOC could handle online real-time video stream and reach 53 fps with only RGB as input.</li></ul> <b>MR2Flow: Efficient Motion Representations for Real-time Video Recognition</b> Advisor: <b>Prof. Limin Wang</b> <ul style="list-style-type: none"><li>As the first author, presented an efficient motion representation by enhancing the discriminative power of motion vector for real-time video recognition, termed as MR2Flow.</li><li>The whole pipeline achieved 94.0% with 100 fps on UCF101 dataset, where the accuracy rate of previous method was 95.8% with 12 fps.</li></ul>	
ACADEMIC SERVICE	Track organizer of ICCV2021 Workshop <a href="#">DeeperAction</a> on localized-and-detailed understanding of human actions in videos. Our track, <i>MultiSports</i> , focused on localizing all action instances with spatio-temporal tubes and recognizing their labels from untrimmed and multi-person videos.	
CONTESTS	<b>Human-centric Spatio-Temporal Video Grounding Challenge.</b> In CVPR2021 Workshop <a href="#">Person in Context</a> . <ul style="list-style-type: none"><li>We got the <b>1st place</b>. First, we extracted tube-level features by SlowFast and CSN on linked tubes based on person boxes predicted by Faster R-CNN. Then we used a 2d-map proposal representation like 2D-TAN and enhanced the feature representation to be more discriminative by multi-modal contrastive learning.</li></ul>	

- Contribution: I generated the person boxes for single frame, then linked the boxes into the tubes and finally extracted the visual features of the tubes.

#### HONORS AND AWARDS

• <b>National Scholarship</b> (9/231)	Ministry of Education	2021
• <b>National Scholarship</b> (2/231)	Ministry of Education	2020
• <b>Outstanding Graduate Students, awardee</b> (2/231)	Nanjing University	2020
• <b>Final Round, Google Girl Hackathon 2020</b> (18/94)	Google	2020
• <b>1st Award, Scholarship for Graduate Students</b> (20%)	Nanjing University	2019&2020
• <b>2rd Award, People's Scholarship</b>	Nanjing University	2018
• <b>3rd Place, ROBOMASTER 2017 (Eastern Division)</b>	Da-Jiang Innovations	2017
• <b>3rd Award, Academic Excellent Scholarship</b>	Nanjing University	2016
• <b>3rd Award, People's Scholarship</b>	Nanjing University	2016

#### SKILLS

- Programming: Python, PyTorch, Matlab, C, Latex,
- Languages: Mandarin, English