

SIPENG ZHENG

✉ zhengsipeng27@gmail.com 🏠 [zhengsipeng.github.io](https://github.com/zhengsipeng) 🌐 Sipeng Zheng ☎ +86 15905058181

I am a research scientist at the Beijing Academy of Artificial Intelligence (BAAI), working with Prof. Zongqing Lu. I received both my PhD and bachelor degrees from Renmin University of China in 2023 and 2018, respectively, under the supervision of Prof. Qin Jin. My research interests focus on human behavior learning, human motion understanding, vision-and-language pretraining, and embodied artificial intelligence. Currently, I am working on developing general-purpose embodied agents, including game agents and humanoid robots.

EDUCATION

Renmin University of China Ph.D in Computer Science	<i>2018.09 - 2023.06</i> Supervisor: Qin Jin
Renmin University of China B.S in Computer Science	<i>2014.09 - 2018.06</i> Supervisor: Qin Jin

WORK EXPERIENCE

Beijing Academy of Artificial Intelligence <i>Research scientist</i>	Beijing, China <i>2023.07 - Present</i>
<ul style="list-style-type: none">• Experience in human motion understanding and generation, as well as human-object interaction.• Skilled in pretraining, data curation, and deploying large multimodal models on embodied agents, including humanoid robots, human-computer interaction systems, and open-world agents (e.g., Minecraft).• Experience in large-scale pretraining on 1,000+ GPU clusters with billion-level datasets.• Hands-on experience with humanoid robots, including models such as H1, H1-b, GR1, GO1, GO2, and others.	
Beijing Academy of Artificial Intelligence <i>Research intern</i>	Beijing, China <i>2021.09 - 2022.02</i>
<ul style="list-style-type: none">• Multi-lingual language-vision-audio pre-training.	
Microsoft Research Asia <i>Research intern</i>	Beijing, China <i>2022.04 - 2022.10</i>
<ul style="list-style-type: none">• Temporal sentence grounding for long-term videos.	

PUBLICATION

1. Boshen Xu, Yuting Mei, **Sipeng Zheng**, and Qin Jin. Learning 3d-aware representation for ego-centric video-language models. CVPR 2025 Submission
2. Yicheng Feng, Yijiang Li, Wanpeng Zhang, **Sipeng Zheng**, and Zongqing Lu. Videorion: Tokenizing object dynamics in videos. CVPR 2025 Submission
3. Jiazheng Liu, Börje F. Karlsson, **Sipeng Zheng**, and Zongqing Lu. Taking notes brings focus? towards multi-turn multimodal dialogue learning. CVPR 2025 Submission
4. Wanpeng Zhang, Zilong Xie, Yicheng Feng, Yijiang Li, Xingrun Xing, **Sipeng Zheng**, and Zongqing Lu. From pixels to tokens: Byte-pair encoding on quantized visual modalities. ICLR 2025 Submission
5. **Sipeng Zheng**, Ye Wang, Bin Cao, Qianshan Wei, and Zongqing Lu. Quo vadis, motion generation? from large language models to large motion models. ICLR 2025 Submission
6. Ye Wang, Yuting Mei, **Sipeng Zheng**, and Qin Jin. Quadrupedgpt: Towards a versatile quadruped agent in open-ended worlds. ICRA 2025 Submission

7. Boshen Xu, Ziheng Wang, Yang Du, Zhinan Song, **Sipeng Zheng**, and Qin Jin. Egocent++: Do egocentric video-language models really understand hand-object interactions? ICLR 2025 Submission
8. Boshen Xu, **Sipeng Zheng**, and Qin Jin. Spaformer: Sequential 3d part assembly with transformers. 3DV 2025
9. **Sipeng Zheng**, Bohan Zhou, Yicheng Feng, Ye Wang, and Zongqing Lu. Unicode: Learning a unified codebook for multimodal large language models. ECCV 2024
10. BAAI Multimodal Interaction Group. Towards general computer control: A multimodal agent for red dead redemption ii as a case study. ICLR 2024 workshop
11. **Sipeng Zheng**, Jiazheng Liu, Yicheng Feng, and Zongqing Lu. Steve-eye: Equipped llm-based embodied agents with visual perception in open worlds. ICLR 2024
12. Yicheng Feng, Yuxuan Wang, Jiazheng Liu, **Sipeng Zheng**, and Zongqing Lu. Llama rider: Spurring large language models to explore the open worlds. NAACL 2024
13. Qi Zhang, **Sipeng Zheng**, and Qin Jin. No-frills temporal video grounding: Multi-scale neighboring attention and zoom-in boundary detection. arxiv 2023
14. **Sipeng Zheng**, Boshen Xu, and Qin Jin. Open-category human-object interaction pre-training via language modeling framework. In *CVPR*, 2023
15. Boshen Xu, **Sipeng Zheng**, and Qin Jin. Pov: Prompt-oriented view-agnostic learning for egocentric hand-object interaction in the multi-view world. In *ACM MM*, 2023
16. Ludan Ruan, Anwen Hu, Yuqing Song, Liang Zhang, **Sipeng Zheng**, and Qin Jin. Accommodating audio modality in clip for multimodal processing. In *AAAI*, 2023
17. **Sipeng Zheng**, Shizhe Chen, and Qin Jin. Few-shot action recognition with hierarchical matching and contrastive learning. In *ECCV*, 2022
18. **Sipeng Zheng**, Shizhe Chen, and Qin Jin. Vrdformer: End-to-end video visual relation detection with transformers. In *CVPR*, 2022 (Oral)
19. **Sipeng Zheng**, Qi Zhang, and Qin Jin. Exploring anchor-based detection for ego4d natural language query. In *CVPR Ego4D workshop*, 2022
20. Bei Liu, **Sipeng Zheng**, Jianlong Fu, and Wen-Huang Cheng. Anchor-based detection for natural language localization in ego-centric videos. In *IEEC*, 2022
21. **Sipeng Zheng**, Shizhe Chen, and Qin Jin. Skeleton-based interactive graph network for human object interaction detection. In *ICME*, 2020
22. **Sipeng Zheng**, Shizhe Chen, and Qin Jin. Visual relation detection with multi-level attention. In *ACM MM*, 2019
23. **Sipeng Zheng**, Xiangyu Chen, Shizhe Chen, and Qin Jin. Relation understanding in videos. In *ACM MM*, 2019

AWARDS

-
- | | |
|--|------|
| ★ National Scholarship for Ph.D Students. | 2022 |
| ★ Ranked 3th in Facebook CVPR 2022 Ego4D Natural Language Query Challenge. | 2022 |
| ★ Ranked 3th, NIST TRECVID 2021 Ad-hoc Video Search (AVS) Challenge. | 2021 |

- ★ Ranked 4th in CVPR 2021 HOMAGE Scene-graph Generation Challenge. 2021
- ★ Ranked 2nd in ACM MM 2020 Video Relationship Understanding Grand Challenge. 2020
- ★ Ranked 2nd in ACM MM 2019 Video Relationship Understanding Grand Challenge. 2019
- ★ Best Method Prize in ACM MM 2019 Grand Challenge 2019
- ★ First Class Scholarship for Ph.D Students. 2018-2021
- ★ First Prize in National University Mathematical Modeling Competition. 2015

PROFESSIONAL ACTIVITIES

- ★ Conference Reviewer for CVPR, ICCV, ECCV, NeurIPS, ICLR, AAAI, ACM MM.
- ★ Journal Reviewer for IJCV, TCSVT, TMM.