




# SIPENG ZHENG

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

I am a research scientist at Beijing Academy of Artificial Intelligence (BAAI). I received my PhD and bachelor degrees from Renmin University of China in 2023 and 2018 respectively under the supervision of Prof. Qin Jin. My research interest lies in human behavior understanding, vision-and-language learning, and embodied artificial intelligence. I am now working towards the general-purpose humanoid robot.

## EDUCATION

---

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

## WORK EXPERIENCE

---

<b>Beijing Academy of Artificial Intelligence</b> <i>Research scientist</i> <ul style="list-style-type: none"><li>• Large multi-modal pre-training for open-world agents (e.g., humanoid robot).</li></ul>	<b>Beijing, China</b> <i>2023.07 - Present</i>
<i>Research intern</i> <ul style="list-style-type: none"><li>• Multi-lingual language-vision-audio pre-training.</li></ul>	<i>2021.09 - 2022.02</i>
<b>Microsoft Research Asia</b> <i>Research intern</i> <ul style="list-style-type: none"><li>• Temporal sentence grounding for long-term videos.</li></ul>	<b>Beijing, China</b> <i>2022.04 - 2022.10</i>

## PUBLICATION

---

- Ye Wang, Yuting Mei, **Sipeng Zheng**, and Qin Jin. QuadrupeDgpt: Towards a versatile quadruped agent in open-ended worlds. Under review 2024
- Boshen Xu, Ziheng Wang, Yang Du, Zhinan Song, **Sipeng Zheng**, and Qin Jin. Egonce++: Do egocentric video-language models really understand hand-object interactions? Under review 2024
- Boshen Xu, **Sipeng Zheng**, and Qin Jin. Spaformer: Sequential 3d part assembly with transformers. arxiv 2024
- Sipeng Zheng**, Bohan Zhou, Yicheng Feng, Ye Wang, and Zongqing Lu. Unicode: Learning a unified codebook for multimodal large language models. ECCV 2024
- BAAI Multimodal Interaction Group. Towards general computer control: A multimodal agent for red dead redemption ii as a case study. ICLR 2024 workshop
- Sipeng Zheng**, Jiazheng Liu, Yicheng Feng, and Zongqing Lu. Steve-eye: Equipped llm-based embodied agents with visual perception in open worlds. ICLR 2024
- Yicheng Feng, Yuxuan Wang, Jiazheng Liu, **Sipeng Zheng**, and Zongqing Lu. Llama rider: Spurring large language models to explore the open worlds. NAACL 2024
- Qi Zhang, **Sipeng Zheng**, and Qin Jin. No-frills temporal video grounding: Multi-scale neighboring attention and zoom-in boundary detection. arxiv 2023
- Sipeng Zheng**, Boshen Xu, and Qin Jin. Open-category human-object interaction pre-training via language modeling framework. In *CVPR*, 2023

10. 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
11. Ludan Ruan, Anwen Hu, Yuqing Song, Liang Zhang, **Sipeng Zheng**, and Qin Jin. Accommodating audio modality in clip for multimodal processing. In *AAAI*, 2023
12. **Sipeng Zheng**, Shizhe Chen, and Qin Jin. Few-shot action recognition with hierarchical matching and contrastive learning. In *ECCV*, 2022
13. **Sipeng Zheng**, Shizhe Chen, and Qin Jin. Vrdformer: End-to-end video visual relation detection with transformers. In *CVPR*, 2022 (Oral)
14. **Sipeng Zheng**, Qi Zhang, and Qin Jin. Exploring anchor-based detection for ego4d natural language query. In *CVPR Ego4D workshop*, 2022
15. Bei Liu, **Sipeng Zheng**, Jianlong Fu, and Wen-Huang Cheng. Anchor-based detection for natural language localization in ego-centric videos. In *IEEC*, 2022
16. **Sipeng Zheng**, Shizhe Chen, and Qin Jin. Skeleton-based interactive graph network for human object interaction detection. In *ICME*, 2020
17. **Sipeng Zheng**, Shizhe Chen, and Qin Jin. Visual relation detection with multi-level attention. In *ACM MM*, 2019
18. **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.