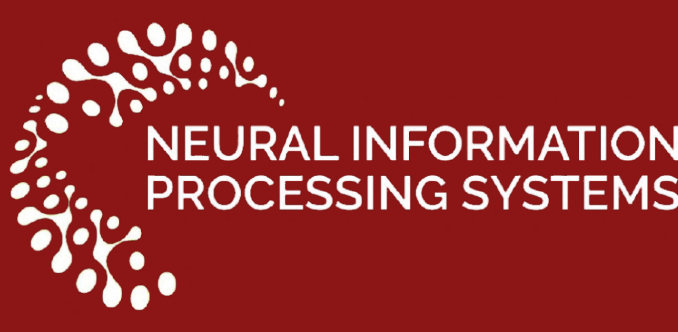


IKEA Manuals at Work: 4D Grounding of Assembly Instructions on Internet Videos

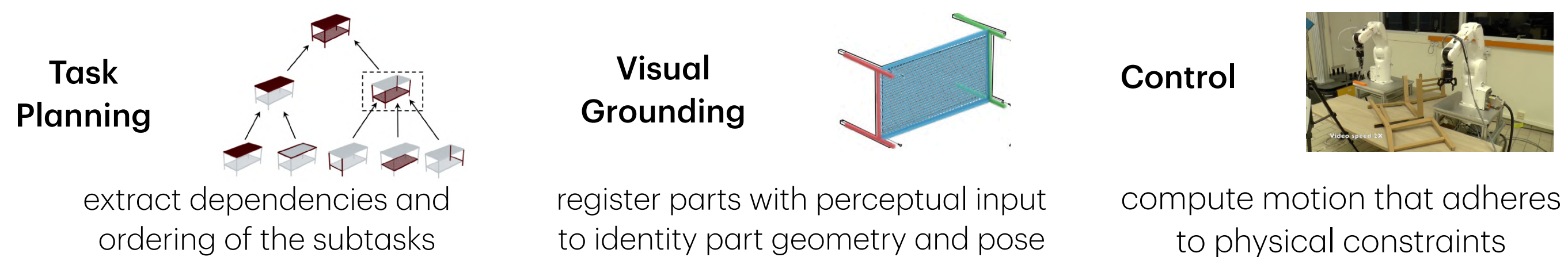
Yunong Liu¹, Cristobal Eyzaguirre¹, Manling Li¹, Shubh Khanna¹, Juan Carlos Niebles¹, Vineeth Ravi², Saumitra Mishra², Weiyu Liu^{*1}, Jiajun Wu^{*1} (*Equal advising)



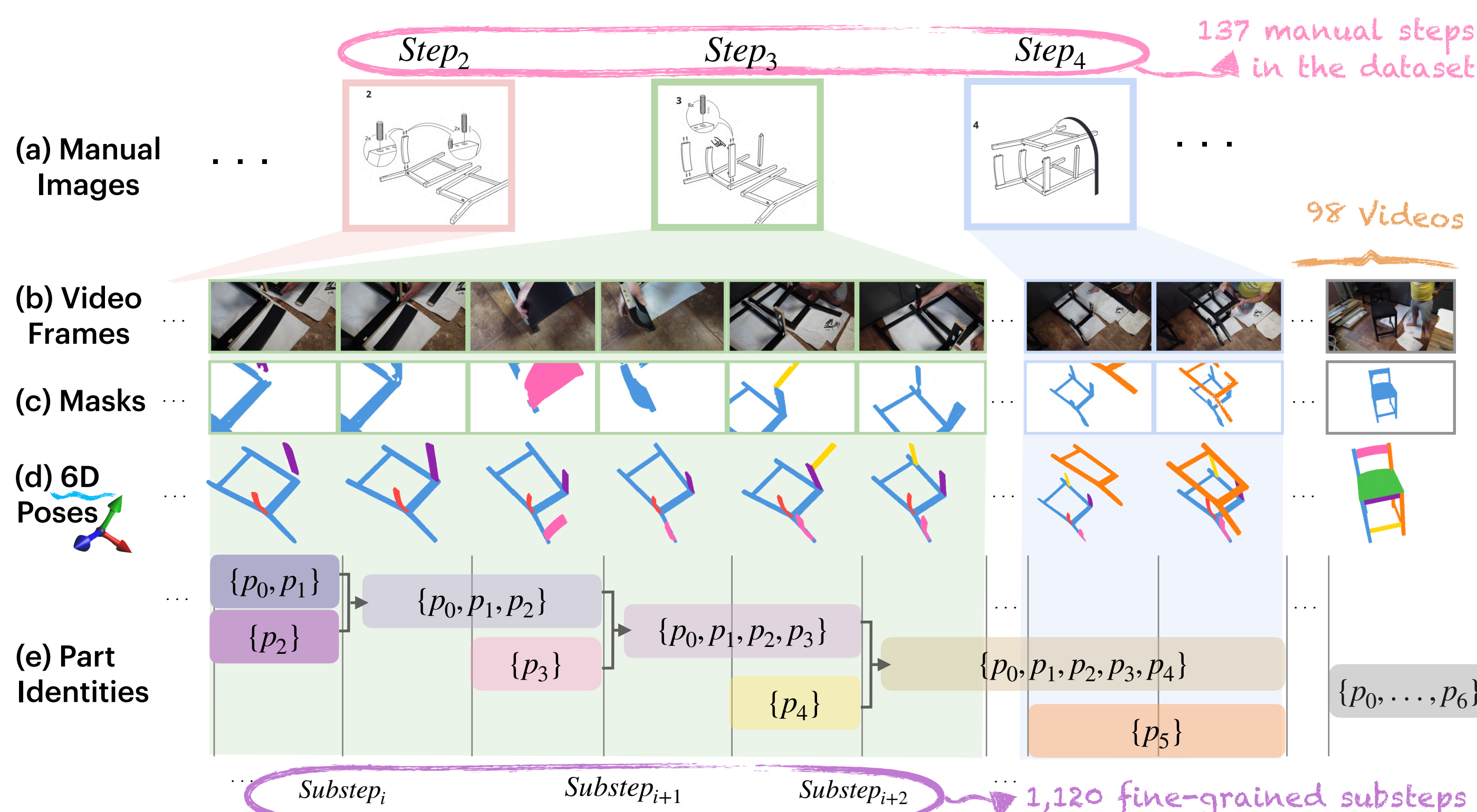
→ I'm looking for 2025 Fall PhD/Internship 🙋🏻🙋🏻¹Stanford University ²J.P. Morgan AI Research

Why it matters — Teaching AI to interpret assembly instructions

- Autonomous 3D shape assembly requires task planning, visual grounding, and control.
- Utilizing instructions, including manuals and how-to videos, requires grounding them in 3D over time.



Our Dataset: IKEA Manuals at Work



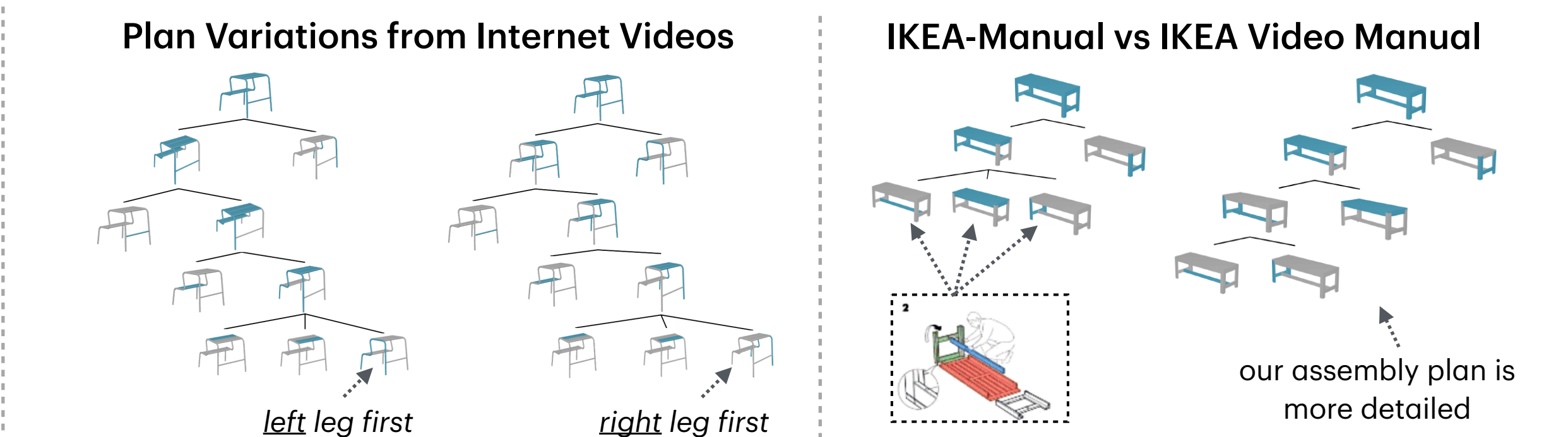
- The first multimodal dataset with extensive 4D groundings of assembly in internet videos.
- Reveal how furniture parts connect in 2D and 3D with parts' 6-DoF poses and segmentation masks.
- Provide fine-grained temporal alignments (steps → substeps → frames) from instruction manuals to internet videos to 3D models.

A Diverse Collection of 36 IKEA Furniture Pieces



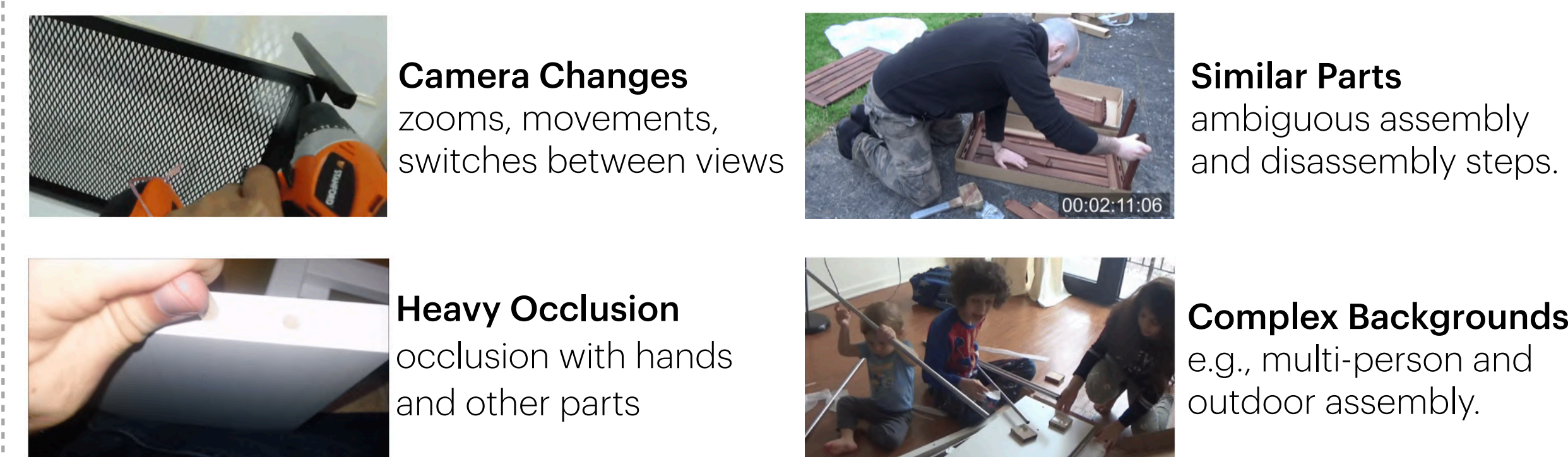
- 36 furniture items provide diversity in 3D geometry, structure, material, and dimension.

Detailed Step-by-Step Instructions



- 25% of furniture objects have more than one unique assembly sequences.

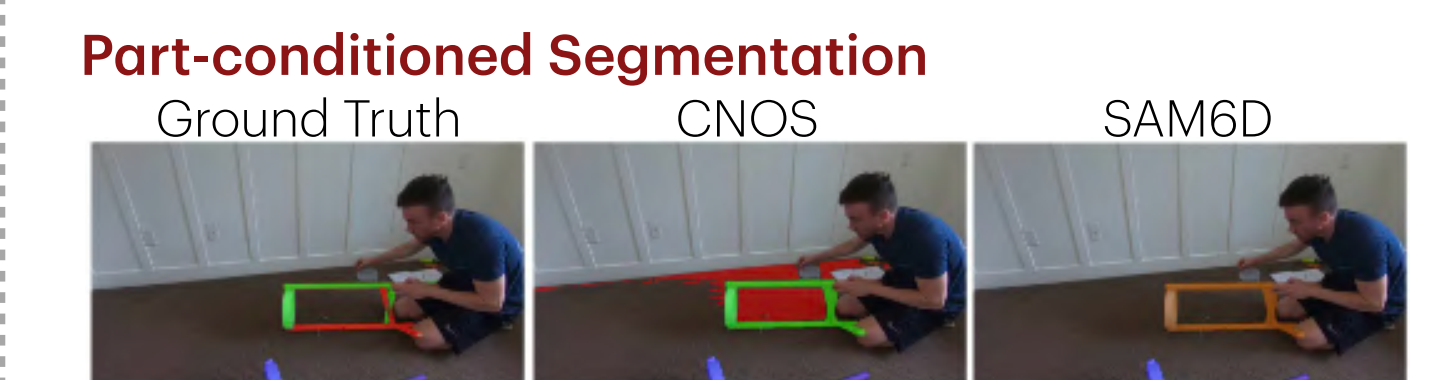
Complexities in Real-World Videos



- Instruction videos from the internet bring challenges in visual grounding.

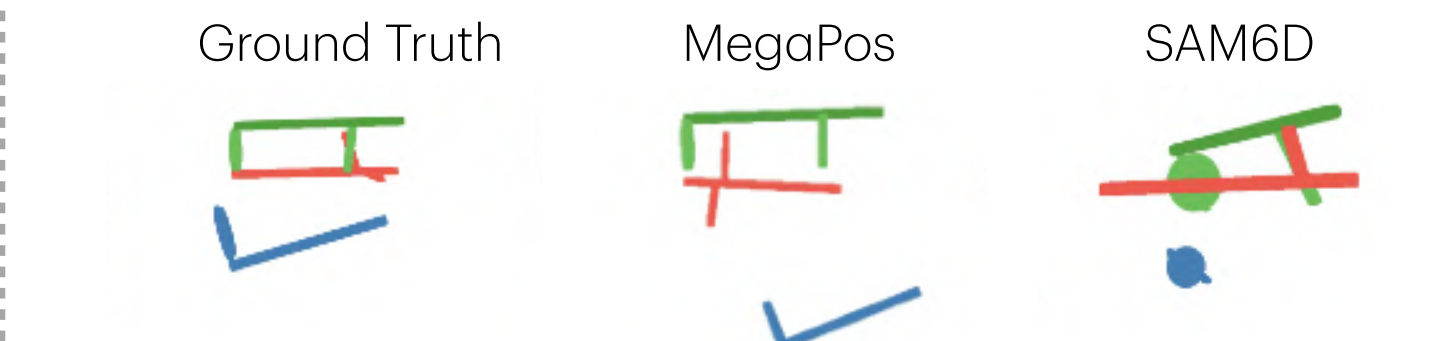
Applications

Break down 4D grounding of videos into sub-tasks.



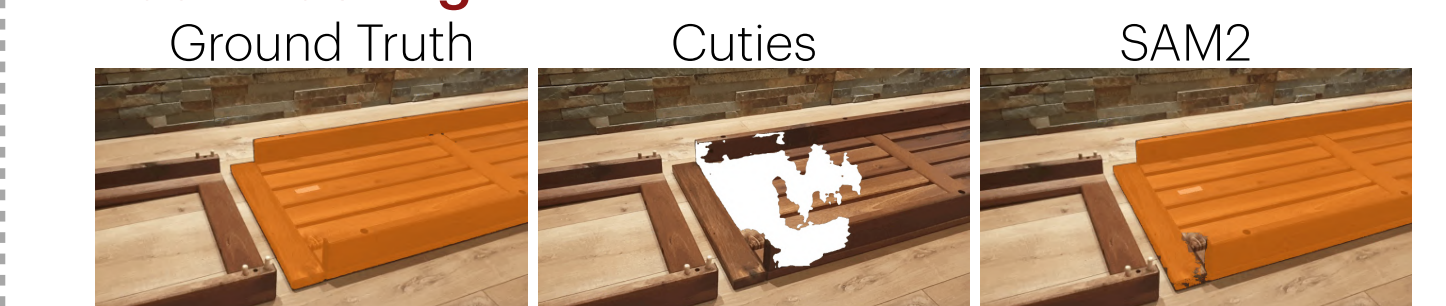
Challenges from visual complexities.

Part conditioned Pose Estimation



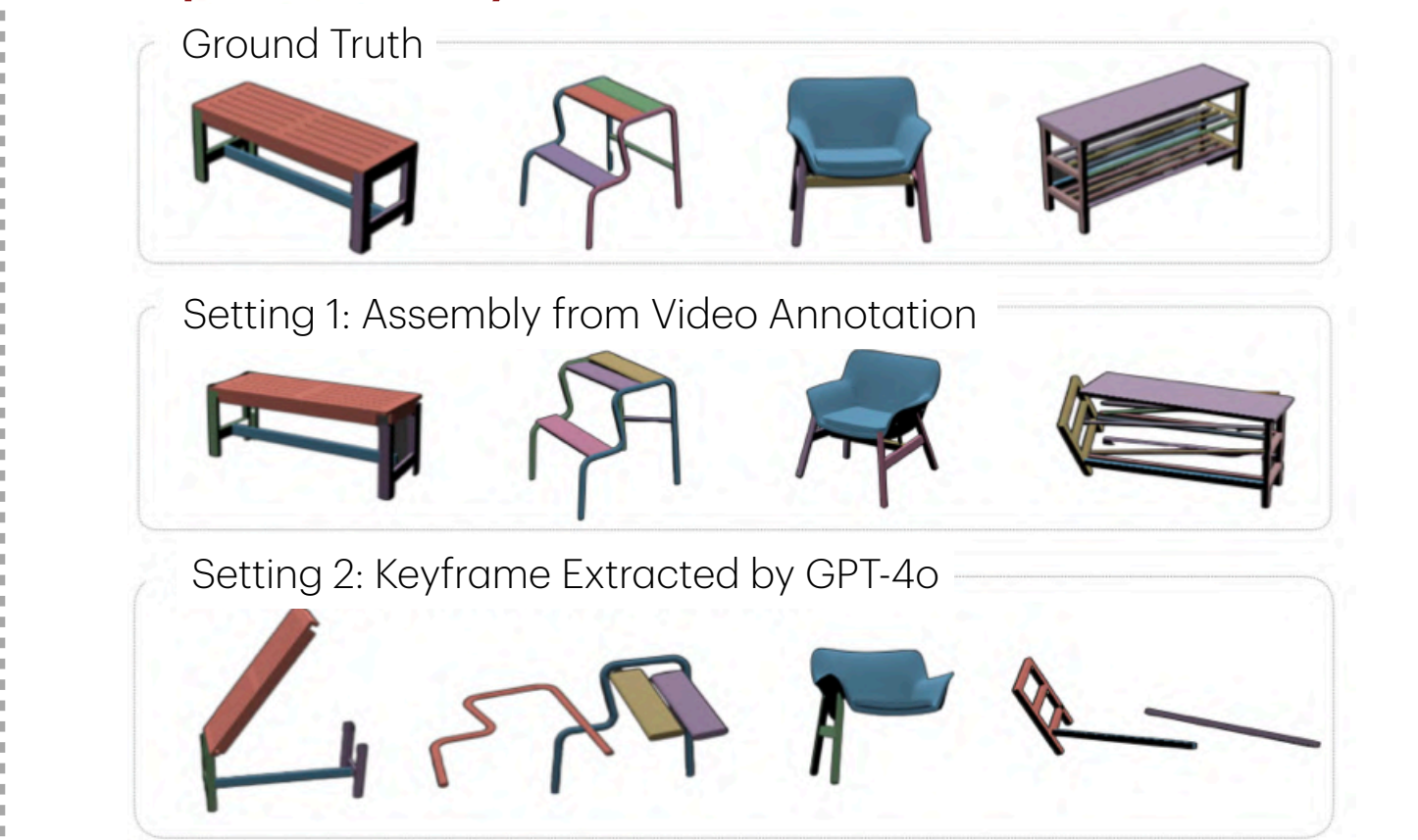
Inaccurate estimated depth can weaken the results.

Mask Tracking



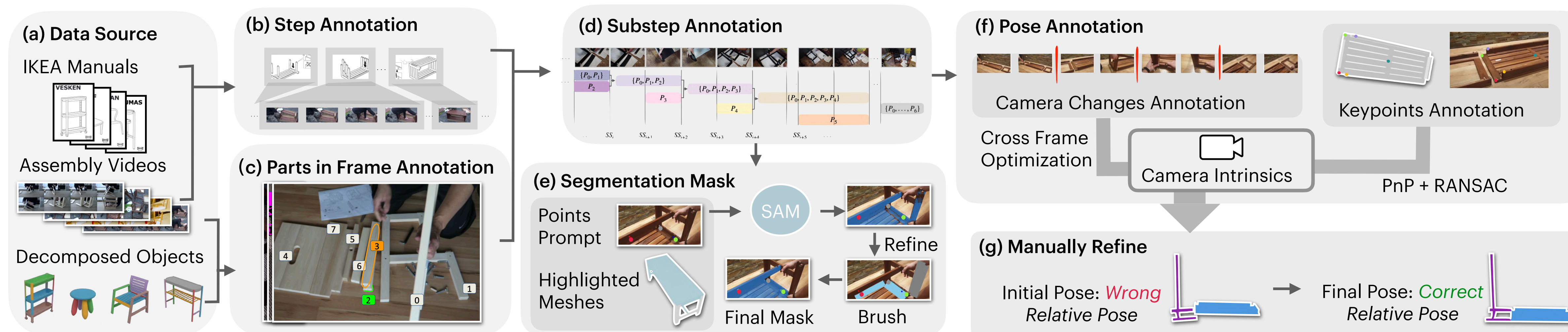
Harder than existing benchmarks for VOS.

Shape Assembly with Instruction Videos

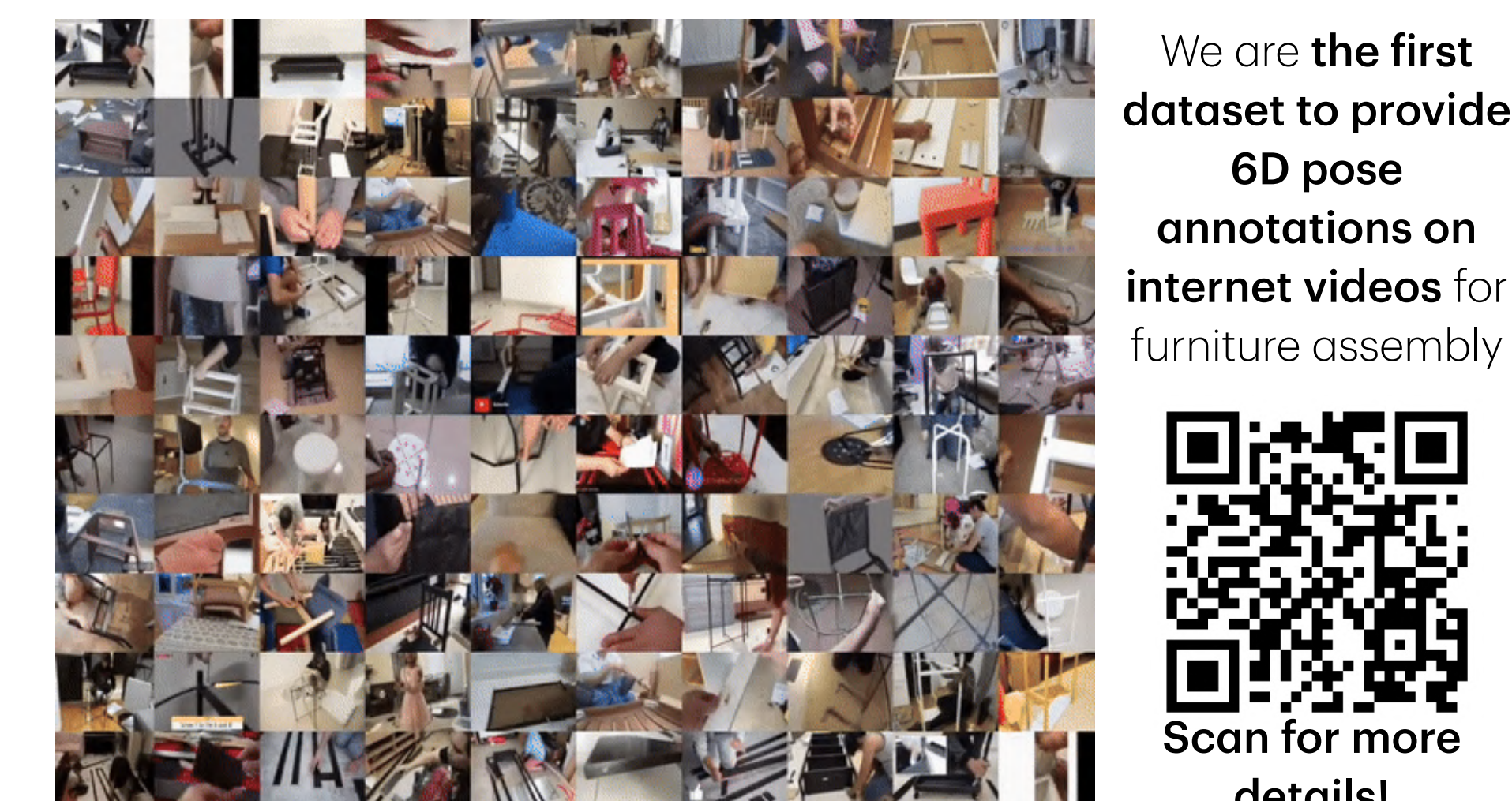


VLMs face challenges in robustly detecting steps and identifying involved parts in videos.

How we got the high-quality spatial-temporal aligned annotations? A semi-automatic pipeline



More details on the website!



We are the first dataset to provide 6D pose annotations on internet videos for furniture assembly



Scan for more details!