# 3D Scene Understanding, Generation, and Interaction for Embodied AI

Baoxiong Jia

General Vision Lab, BIGAI

# Perception

# Embodied? ☹


Robot vs. Door, DARPA Challenge (2015)


Fu et al., Mobile-ALOHA (2024)

## Moravec's Paradox

*It's comparatively **easy** to make computers exhibit adult level performance on **intelligence tests or playing checkers**, and **difficult or impossible** to give them skills of a one-year old when it comes to **perception and mobility**.*

Hans Moravec, Mind Children, 1988

# Embodied AI

*"The embodiment hypothesis is the idea that **intelligence emerges in the interaction of an agent with an environment** and as a result of sensorimotor activity"*

*Smith & Gasser, The Development of Embodied Cognition: Six Lessons from Babies, 2005*



### Manipulation & Locomotion

RL / Imitation learning / MPC on **specific scenes or skills**

*Boston Dynamics, Atlas | Partners in Parkour, 2022*
*https://www.youtube.com/watch?v=tF4DML7FIWk*

### Interaction with scenes in daily life

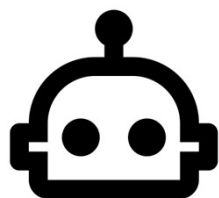**Various** object attributes and **diverse** scene configurations

**Long-horizon interaction** with scenes

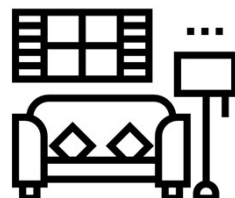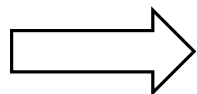*Damen et al., Scaling Egocentric Vision: The Epic-Kitchens Dataset, 2018*
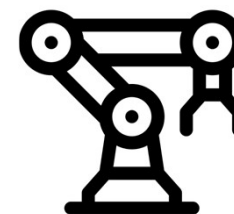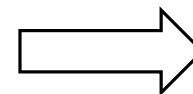
# From the scene perspective

Perception        Grounding        Action

- 3D or even 4D data capture
- Representation efficiency
- …

- Spatial relationships in situations
- Affordance & functionality
- …

- Object Geometry / Physics
- Embodiment gap
- …

# 3D scene understanding for EAI



| Dataset | 3D Data | | Language | | Total |
|---|---|---|---|---|---|
| | Scene | Object | Anno. | Syn. | |
| ScanRefer | | | 52K | - | 52K |
| ReferIt3D | | | 42K | 200K | 242K |
| ScanQA | 1.5K | 33K | 27K | - | 27K |
| SQA3D | | | - | 33K | 33K |
| Multi3DRefer | | | 52K | 10K | 62K |
| Cap3D | - | 666K | 58K | 666K | 724K |
| ScanScribe | 3K | 56K | 94K | 184K | 278K |

| Dataset | 2D Image-text pairs |
|---|---|
| MS-COCO | 330K |
| Visual Genome | 5.4M |
| WIT | 5.5M |
| Conceptual Captions-12M | 12M |
| YFCC100M | 100M |
| LAION-5B | 2.3B |

# Scaling 3D-VL with SceneVerse



*Jia et al., SceneVerse: Scaling 3D Vision-Language Learning for Grounded Scene Understanding, ECCV 2024*

# Uniting scene representations



*Zhu et al., Unifying 3D Vision-Language Learning via Promptable Queries, ECCV 2024*

# Findings & Takeaways

- Language
  - Relatively easy to scale
  - Quality of language matters

- Scene
  - Imbalanced classes
  - Domain gap between synthetic and real data
  - Domain gap between real-world datasets



*"In the corner of the room are boxes. the first two book shelves in the corner to the right of the boxes are the bookshelves we are looking for."*



*Lamp*



*Plant*

| Real | Synthetic | SCENEVERSE-val | S3D | ProcTHOR |
|------|-----------|----------------|------|----------|
| All | ✗ | 64.8 | 37.1 | 43.4 |
| ✗ | S3D | 7.0 | 85.1 | 16.1 |
| ✗ | ProcTHOR | 4.2 | 16.3 | 91.0 |

# Can scene generation or reconstruction help?



*Fu et al., 3D-FRONT: 3D Furnished Rooms with layOuts and semaNTics, ICCV 2021*

Collisions between objects

Objects outside of floor plan

Areas unreachable to agents

*Yang et al., PhyScene: Physically Interactable 3D Scene Synthesis for Embodied AI, CVPR 2024 (Highlight)*

# Reconstruction?



**Monocular Cues**

**Previous Method**

In Isaac Gym

video

t=0

Final state

# Reconstructing scenes with physical constraints



**Monocular Cues**

**Physical Simulator**

t=0 ... Final state

**Previous Method**

In Isaac Gym
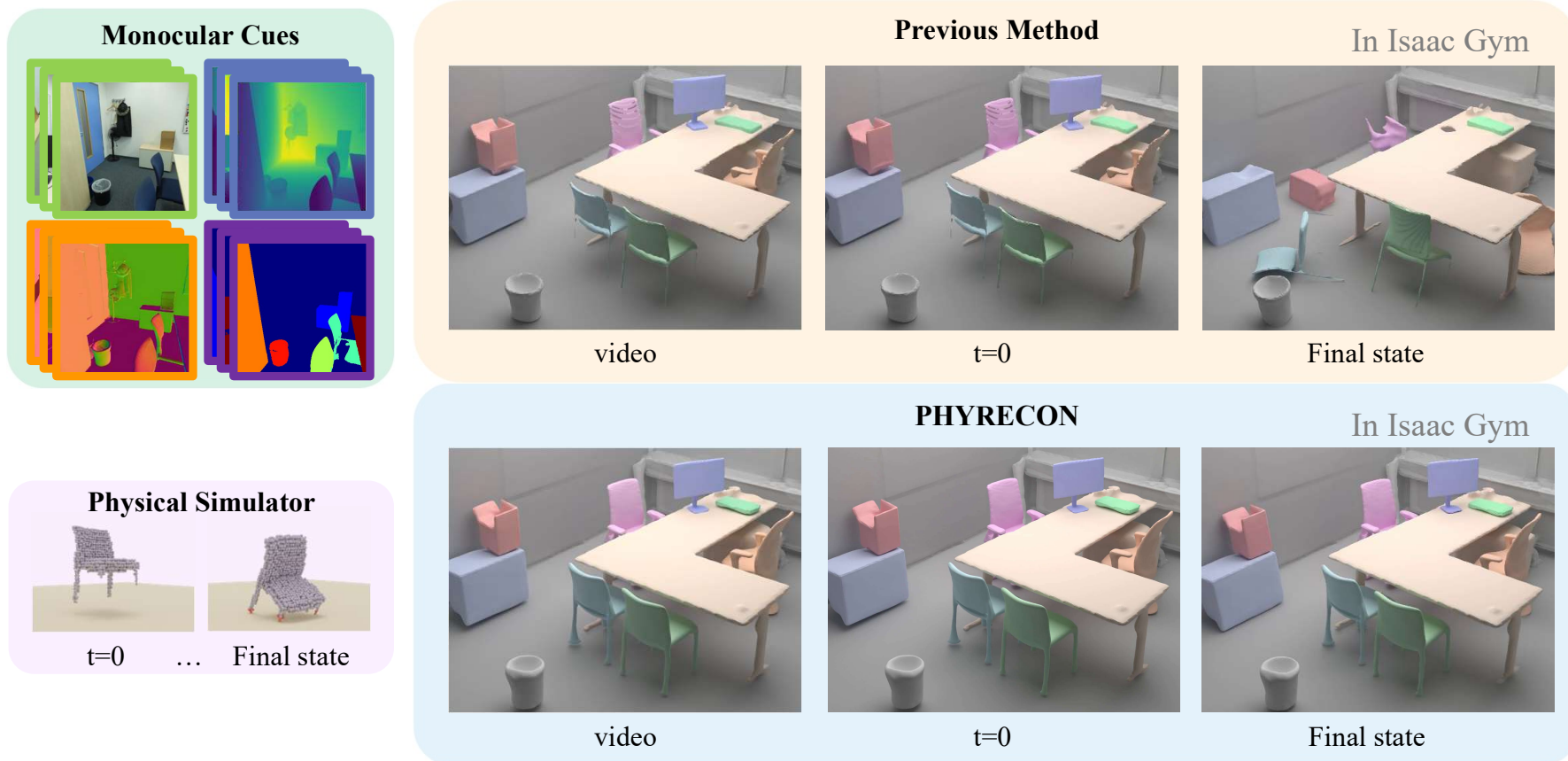
video          t=0          Final state

**PHYRECON**

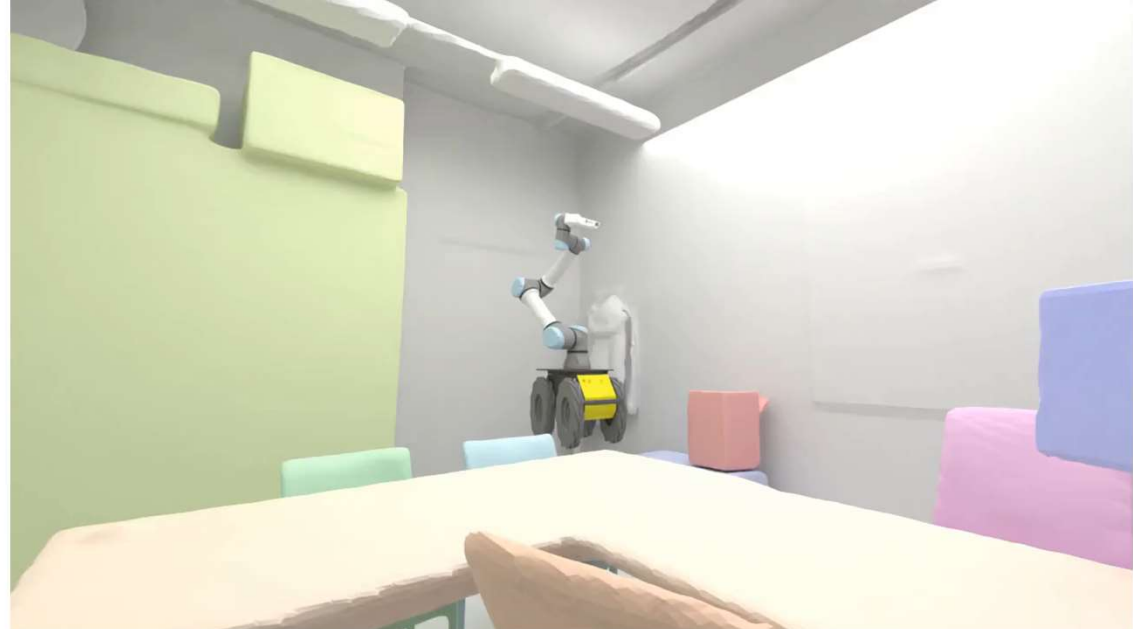In Isaac Gym

video          t=0          Final state

*Ni et al., PhyRecon: Physically Plausible Neural Scene Reconstruction, NeurIPS 2024*

- Generation
  - Insufficient training scale & diversity
  - Limited diversity both assets and layout

- Reconstruction
  - No object articulation, pick & place only
  - Limited efficiency and scaling potential

- Articulated asset reconstruction

- Retrieval augmented reconstruction

General Vision Lab, BIGAI

*Afford-motion, CVPR 2024 Highlight*

The man walks to the chair in a curve.

*TRUMANS, CVPR 2024 Highlight*

*LingoMotions, SIGGRAPH ASIA 2024*

Qualitative results of our synthesized motions

Text Instruction

Go to the sink to wash the hands

I am hungry. Could you give me some food? And pass me a cup of juice.

15x

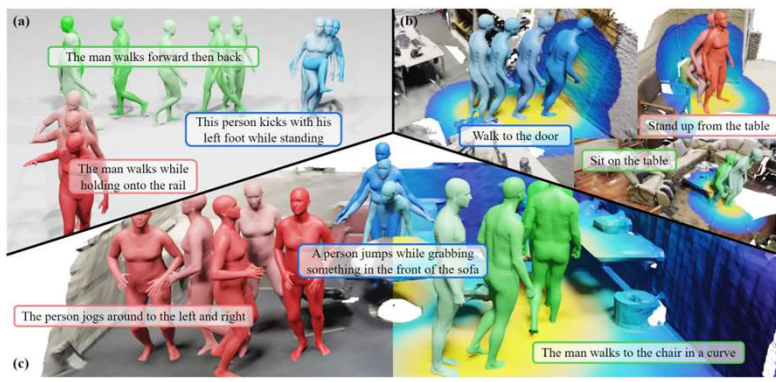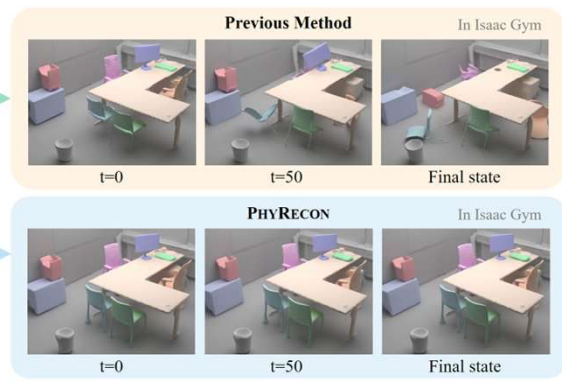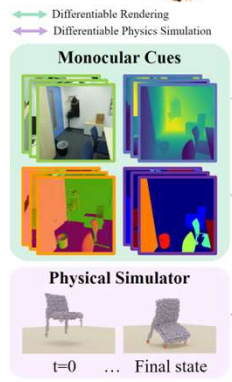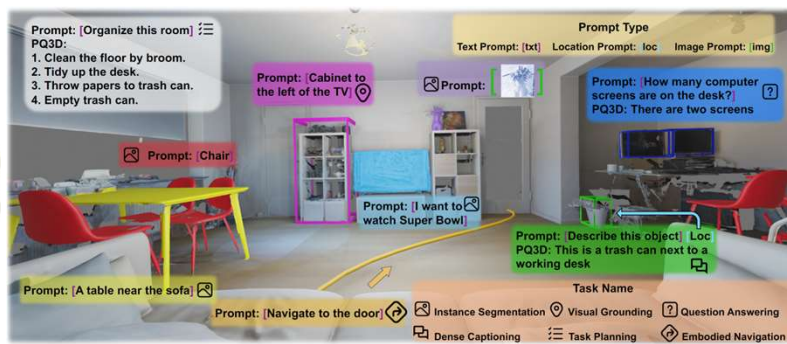Active Perception

COME-Robot, ArXiv 2024

# Above all

- Scaling works in 3D scene understanding
  - Unifying different domains and situation modeling

- Scene curation for embodied AI is still challenging
  - Ensuring naturalness and physics while maintaining diversity

- Interaction data from human-scene interaction, egocentric videos
  - Robust robot system for coordinating different modules for real-world applications
  - Transfer from motion to robots, agent agnostic policy learning

# More to come at BIGAI



Thank you!