

# 3D Graphics System Challenges for Simulation: Lessons from AI Habitat

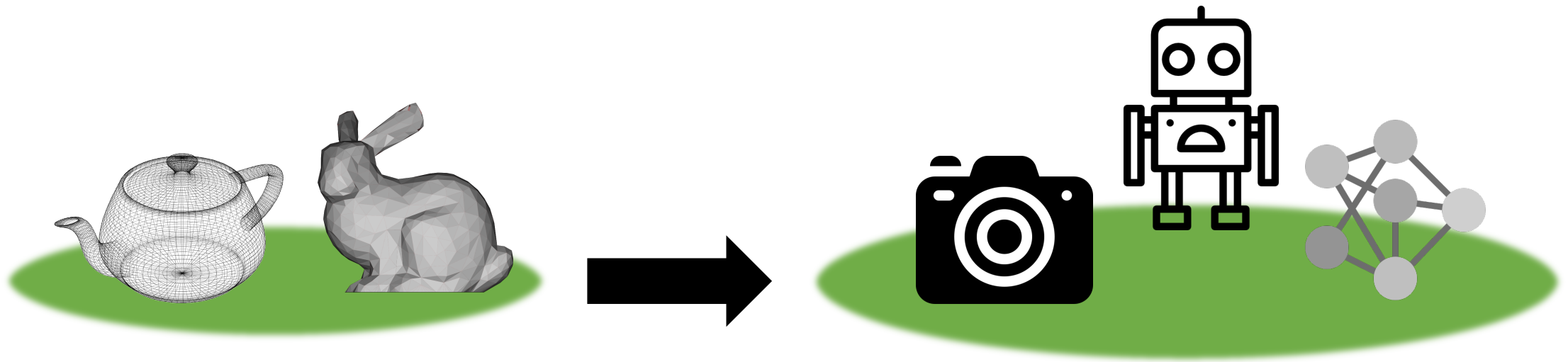
Manolis Savva

HPG 2020

2020-07-16



Preface: “thoughts from a graphics expat”



“Simulation”?

# Terminology: 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.”*

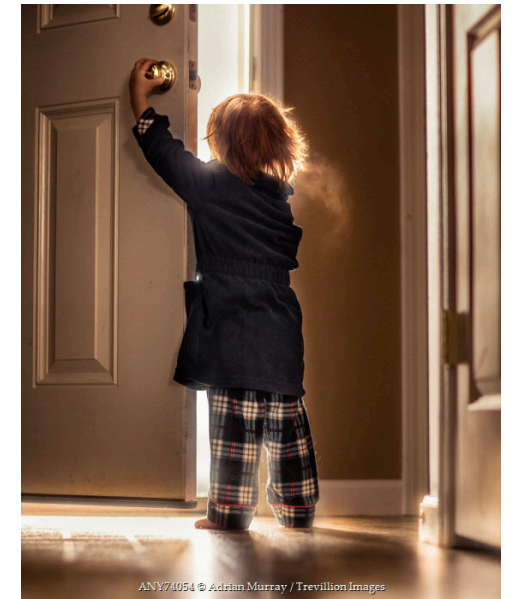
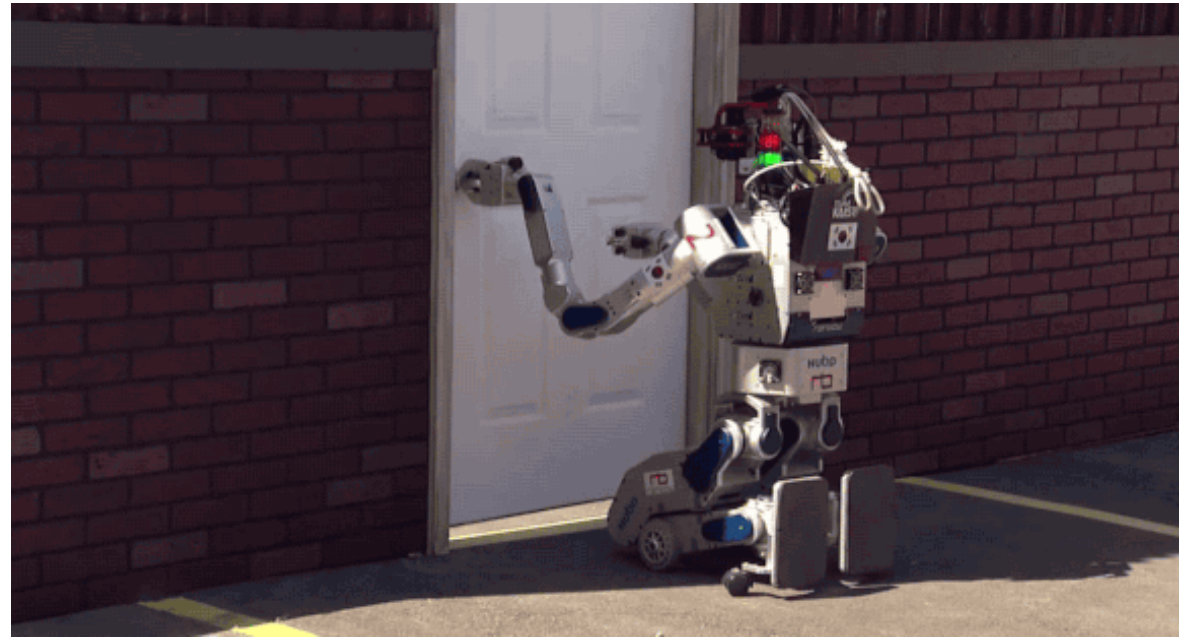
The Development of Embodied Cognition: Six Lessons from Babies  
[Smith & Gasser 2005]

# Embodied Agents

Physically embodied agents  
taking actions in the world

= Human-like AI

- Active perception
- Long-term planning
- Learning by interaction



# Simulation for embodied AI

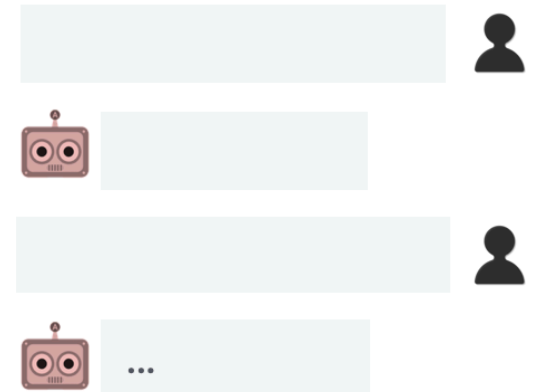
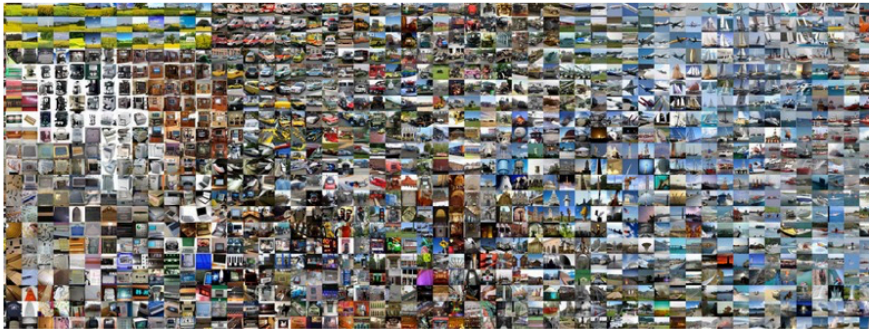
Physically embodied agents  
taking actions in the world

*Virtual* embodied agents  
taking actions in a *virtual* world

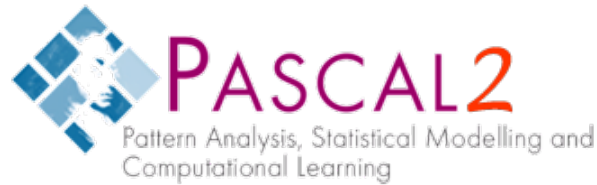
# Internet AI



# Embodied AI



# From internet image datasets to 3D simulators

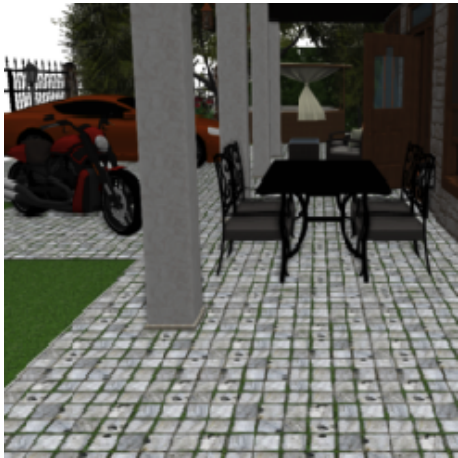


Dataset → Simulator → Task → Benchmark



Year 2017: exciting times!

# 3D simulators galore!



*HoME Platform*  
[Brodeur et al. 2017]



*House3D*  
[Wu et al. 2017]



*MINOS*  
[Savva et al. 2017]



*AI2-THOR*  
[Kolve et al. 2017]



*Matterport3D Simulator*  
[Anderson et al. 2018]



*Gibson Environment*  
[Zamir et al. 2018]



*InteriorNet / ViSim*  
[Li et al. 2018]

...

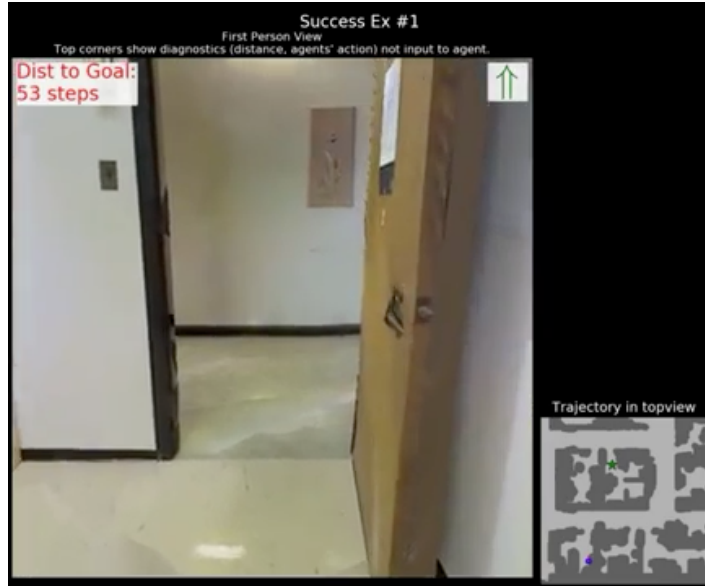
# 3D simulators galore!

Environment	3D	Large-Scale	Customizable	Physics	Photorealistic	Actionable
Atari						
OpenAI Universe	✓	✓	✓			
Malmo	✓	✓	✓			
DeepMind Lab	✓		✓			
VizDoom	✓		✓			
Matterport3D	✓				✓	
MINOS (Matterport3D)	✓				✓	
House3D	✓	✓	✓			
MINOS (SUNCG)	✓	✓	✓			
HoME	✓	✓	✓	✓		
AI2-THOR	✓		✓	✓	✓	✓

Table from AI2-THOR [Kolve et al. 2017]

# Impact: research tasks and communities

## Visual navigation



[Gupta et al. 2017]

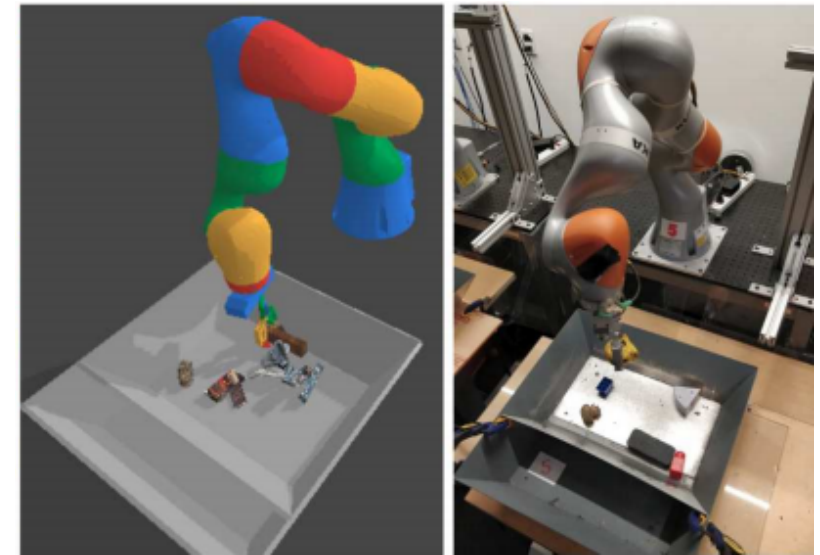
## Instruction following



Leave the bedroom, and enter the kitchen. Walk forward, and take a left at the couch. Stop in front of the window.

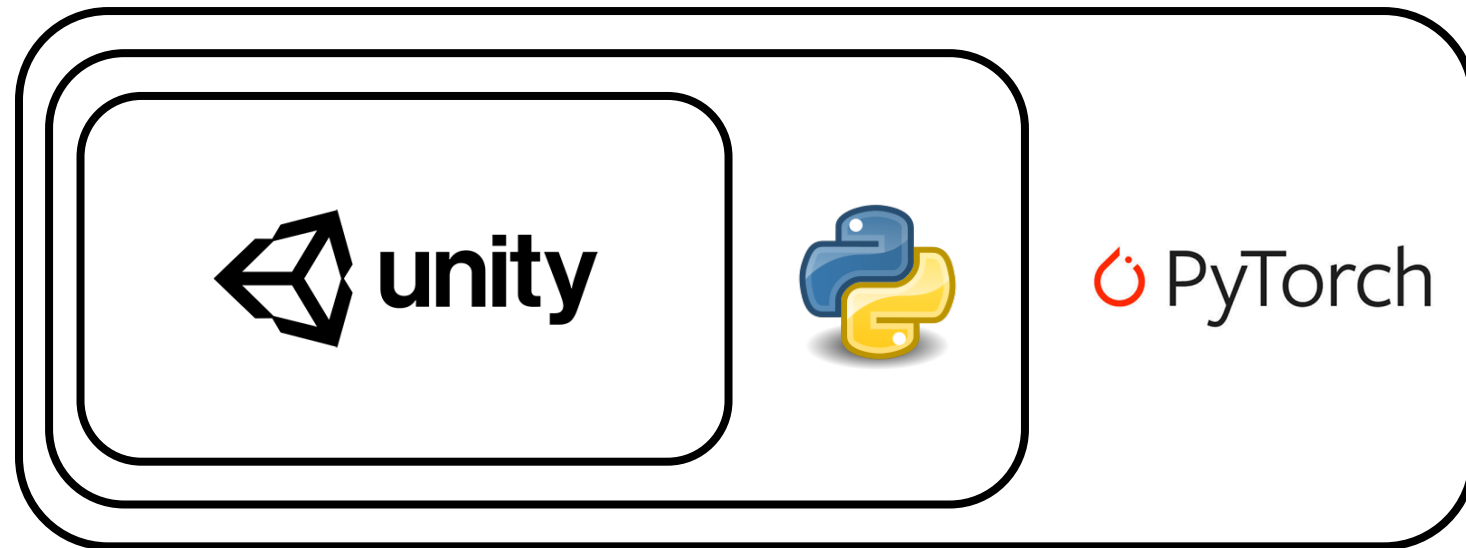
[Anderson et al. 2018]

## Robotic manipulation



[James et al. 2019]

# Common: black-boxed 3D game engine binary

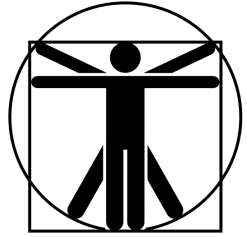


AI2-THOR [Kolve et al. 2017]  
architecture example sketch

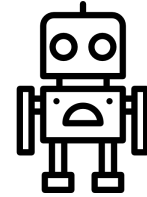


10 – 60 FPS

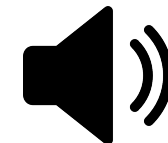
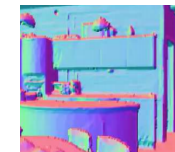
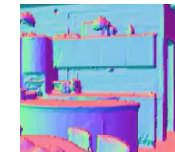
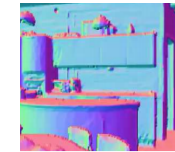
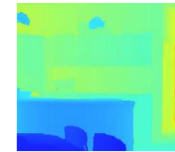
# However: not for human eyeballs!



**Human:** 1080p @ 60Hz



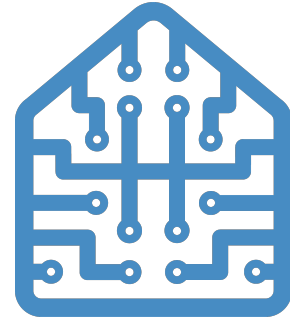
**RL:** 84x84 @ 1000+ Hz



⋮

Can we do better?

# Habitat: A Platform for Embodied AI Research



[aihabitat.org](http://aihabitat.org)



Manolis Savva\*



Abhishek Kadian\*



Oleksandr Maksymets\*



Yili Zhao



Erik Wijmans



Bhavana Jain



Julian Straub



Jia Liu



Vladlen Koltun



Jitendra Malik



Devi Parikh



Dhruv Batra

**facebook**  
Artificial Intelligence Research

**facebook**  
Reality Labs

**Georgia  
Tech**

**SFU**

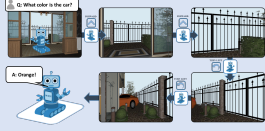

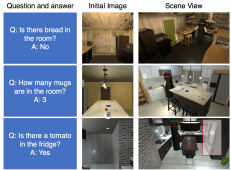


**intel**

**Berkeley**  
UNIVERSITY OF CALIFORNIA



# Habitat: standardizing the Embodied AI “software stack”

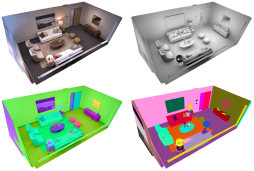
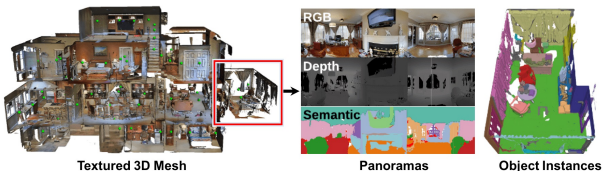
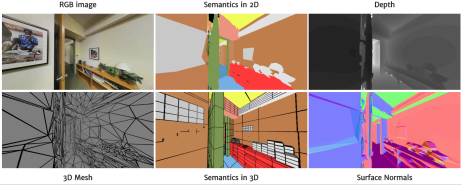
Tasks

 <p><b>EmbodiedQA</b> (Das et al., 2018)</p>	 <p><b>Language grounding</b> (Hill et al., 2017)</p>	 <p><b>Interactive QA</b> (Gordon et al., 2018)</p>	 <p><b>Instruction following</b> (Anderson et al., 2018)</p>	 <p><b>Visual Navigation</b> (Zhu et al., 2017, Gupta et al., 2017)</p>
---	--	---	---	--

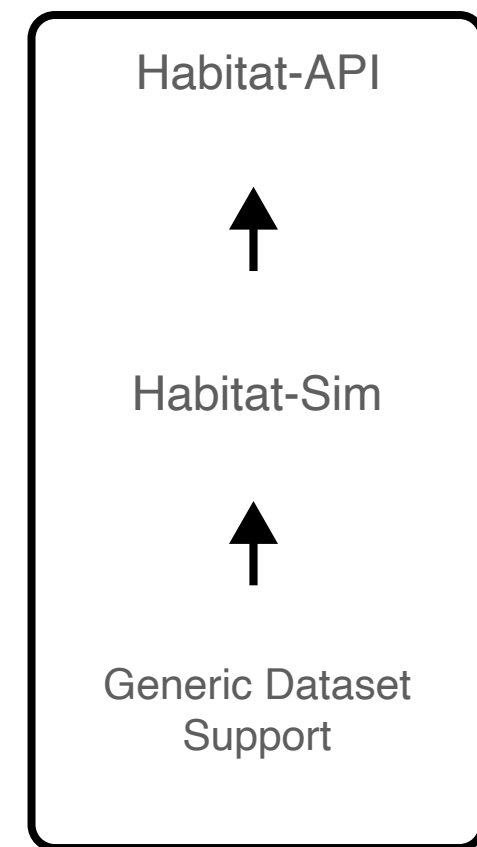
Simulators

 <p><b>House3D</b> (Wu et al., 2017)</p>	 <p><b>AI2-THOR</b> (Kolve et al., 2017)</p>	 <p><b>MINOS</b> (Savva et al., 2017)</p>	 <p><b>Gibson</b> (Zamir et al., 2018)</p>	 <p><b>CHALET</b> (Yan et al., 2018)</p>
---	---	--	---	---

Datasets

 <p><b>Replica</b> (Straub et al., 2019)</p>	 <p><b>Matterport3D</b> (Chang et al., 2017)</p>	 <p><b>2D-3D-S</b> (Armeni et al., 2017)</p>
---	--	---

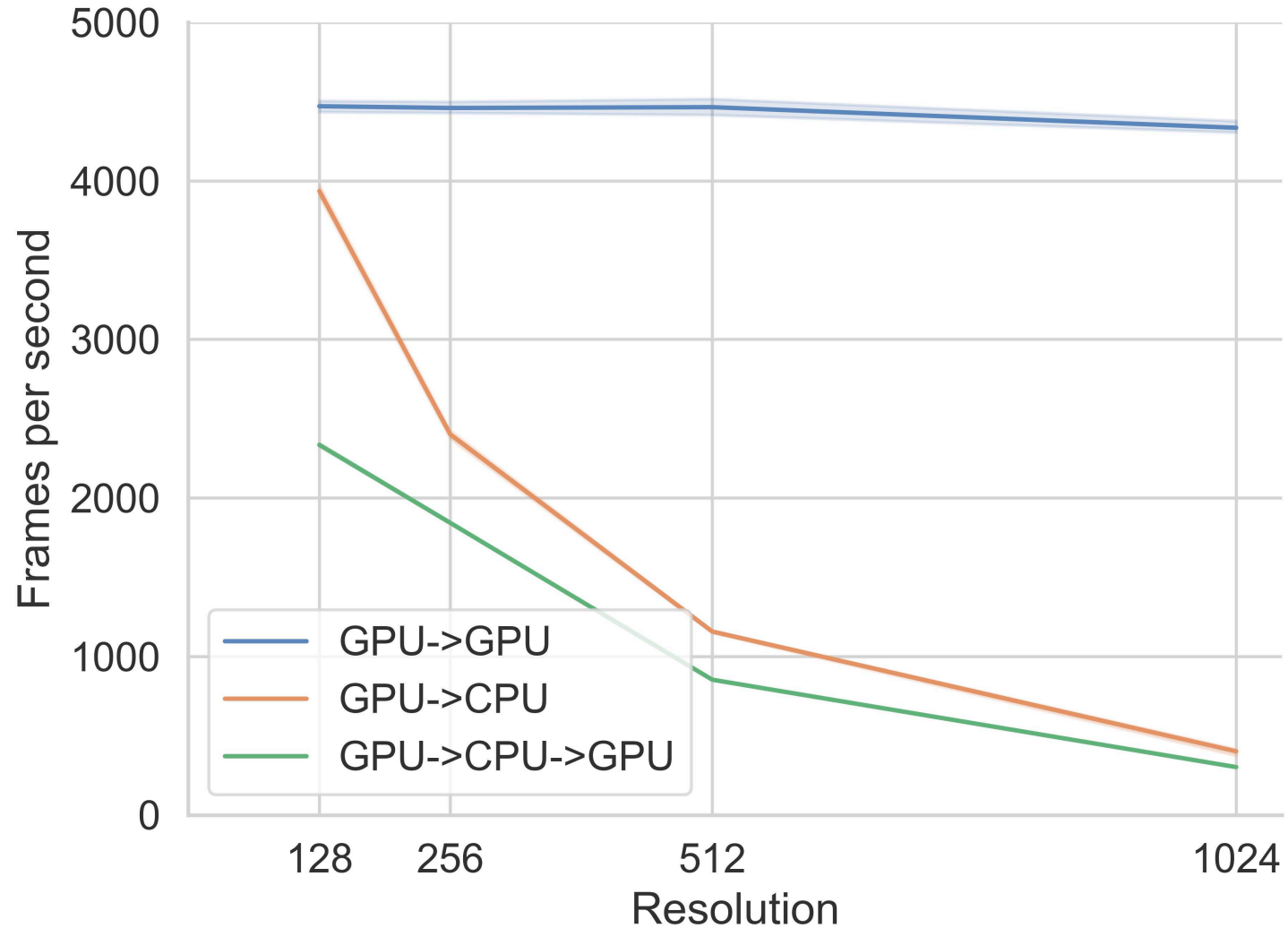
Habitat Platform





Matterport3D

# Attention to speed



Did speed matter?

# Learned vs classical navigation agents

## To Learn or Not to Learn: Analyzing the Role of Learning for Navigation in Virtual Environments

Noriyuki Kojima  
University of Michigan  
2260 Hayward St, Ann Arbor, MI 48109  
kojimana@umich.edu

Jia Deng  
Princeton University  
35 Olden St 423, Princeton, NJ 08540  
jiadeng@cs.princeton.edu

### Abstract

*In this paper we focus on the task of geometric navigation, i.e. navigation when ground-truth 3D information is available. Specifically, we explore the dichotomy between "learning" and "coding" for this task. We construct a hand-coded navigating agent, and demonstrate that it outperforms state-of-the-art learning based agents on two popular benchmarks, MINOS [37] and Stanford large-scale 3D Indoor Spaces (S3DIS) [2]. We also observe that as the environment becomes more challenging, the performance gap between learning-based and hand coded-agent increases.*

ods. Therefore, in the context of geometric navigation, the strengths and weaknesses of "learning" over "coding" are not clear. In this paper, we attempt to clarify this so that intelligent choices can be made while developing real-world systems.

We construct a hand-coded agent for the task of geometric navigation and compare its performance with state-of-the-art learning based methods on two challenging benchmarks: S3DIS [2] and MINOS [37]. On MINOS, the UNREAL agent [37] (which is based on deep reinforcement learning) and on S3DIS, the CMP agent [14] (which uses imitation learning to jointly train a mapper and plan-

## Benchmarking Classic and Learned Navigation in Complex 3D Environments

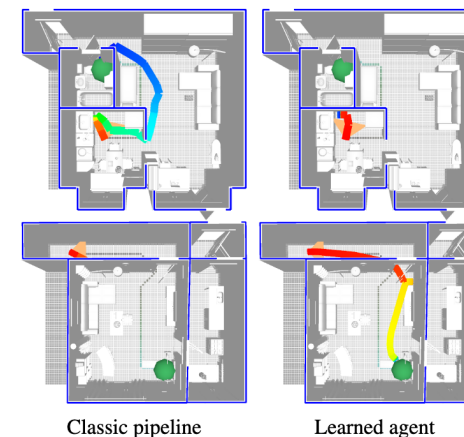
Dmytro Mishkin\*  
Czech Technical University

Alexey Dosovitskiy  
Intel Labs

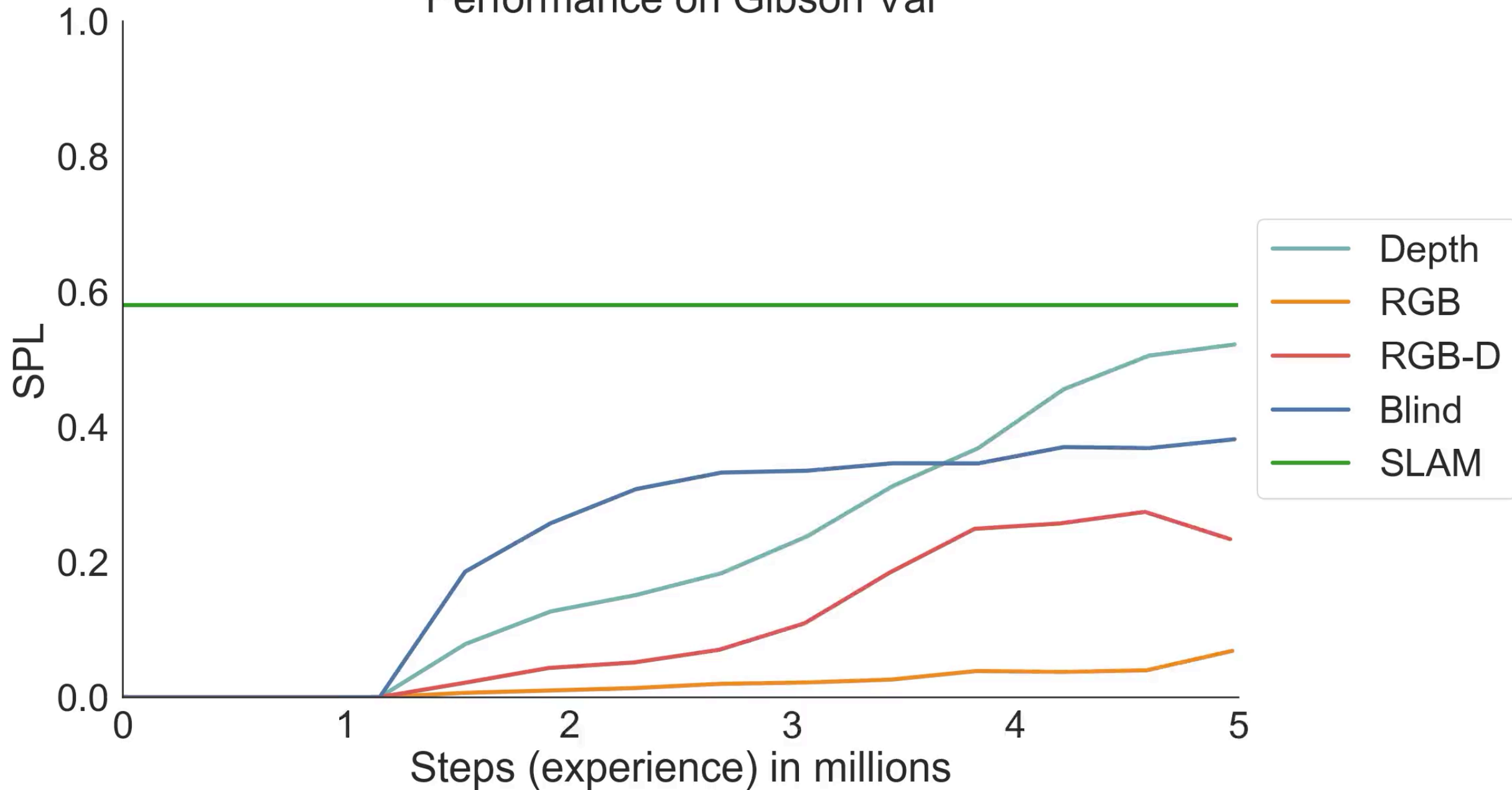
Vladlen Koltun  
Intel Labs

### Abstract

*Navigation research is attracting renewed interest with the advent of learning-based methods. However, this new line of work is largely disconnected from well-established classic navigation approaches. In this paper, we take a step towards coordinating these two directions of research. We set up classic and learning-based navigation systems in common simulated environments and thoroughly evaluate them in indoor spaces of varying complexity, with access to different sensory modalities. Additionally, we measure human performance in the same environments. We find that a classic pipeline, when properly tuned, can perform very well in complex cluttered environments. On the other hand, learned systems can operate more robustly with a limited sensor suite. Both approaches are still far from human-level performance.*

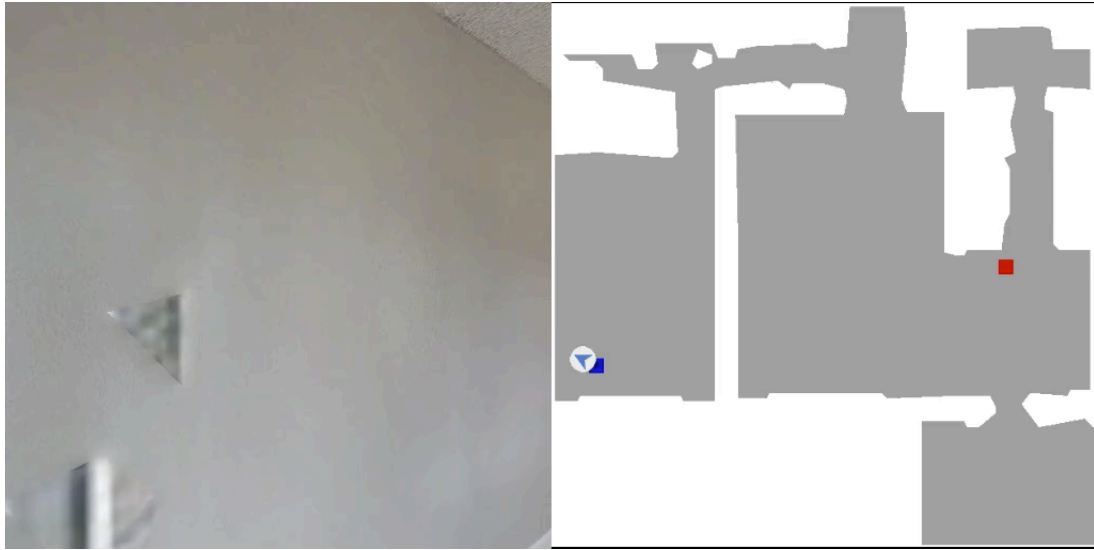


# Performance on Gibson Val

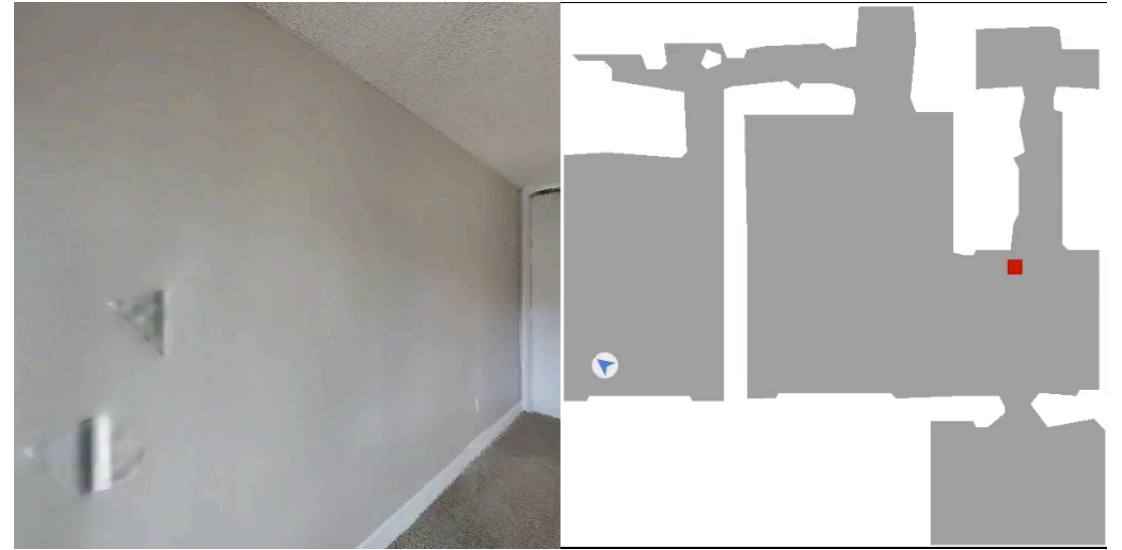


# Example navigation episodes

Blind Agent



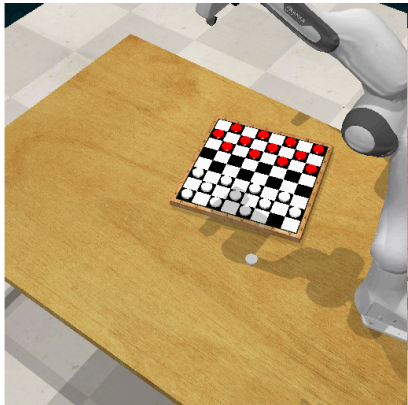
Depth Agent



# Back to today: simulators galore part 2!

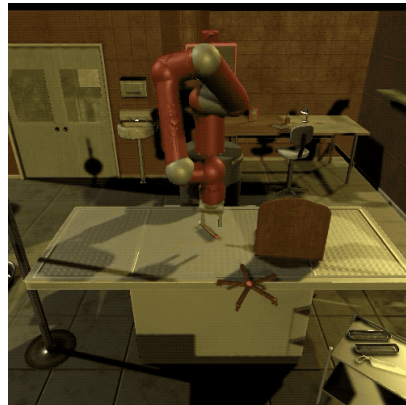
*RLBench*

[James et al. 2019]



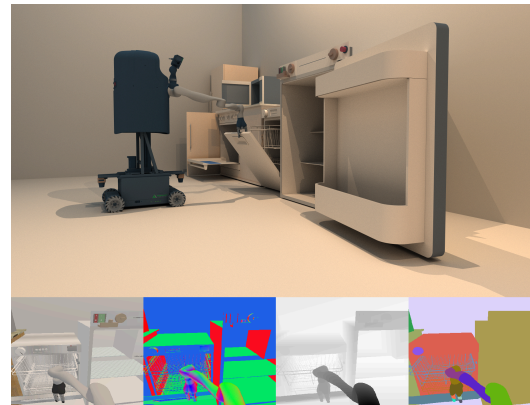
*IKEA Furniture Assembly*

[Lee et al. 2019]



*SAPIEN*

[Xiang et al. 2020]



*iGibson*

[Xia et al. 2020]



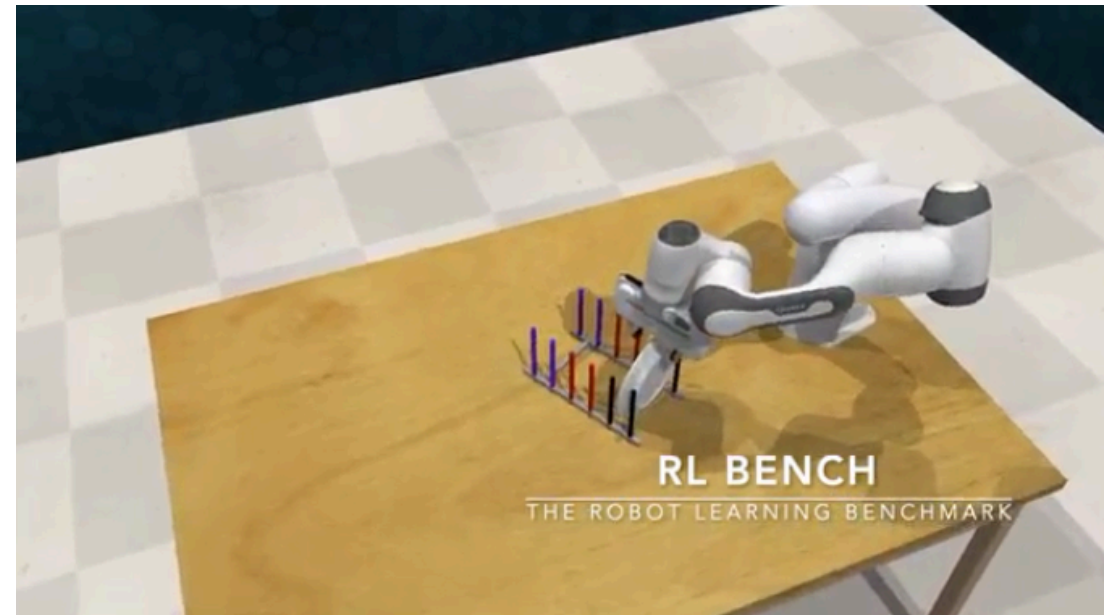


Emerging trends

# Emerging trends: interaction

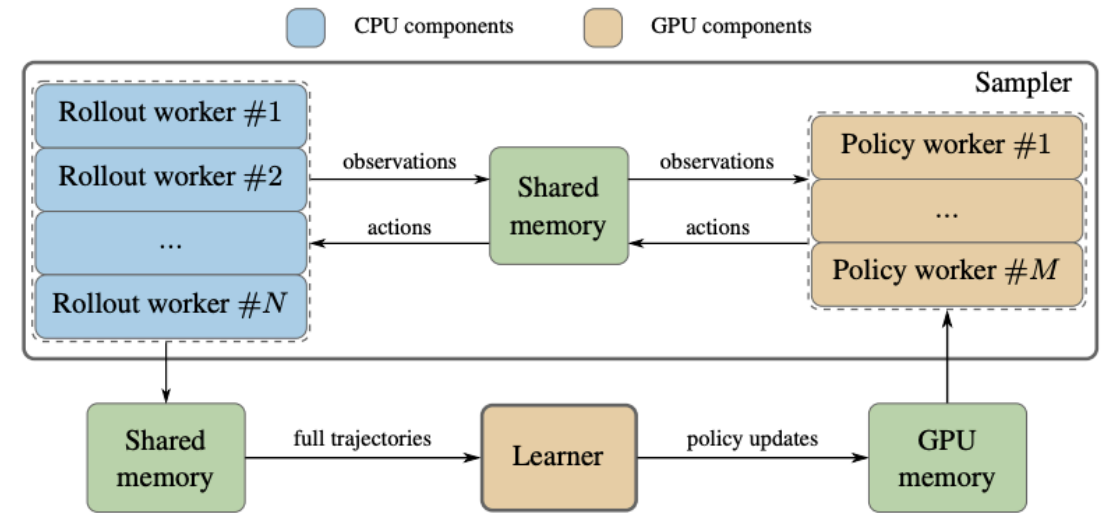
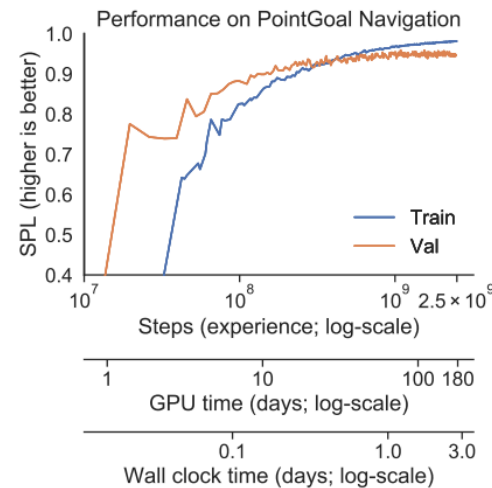
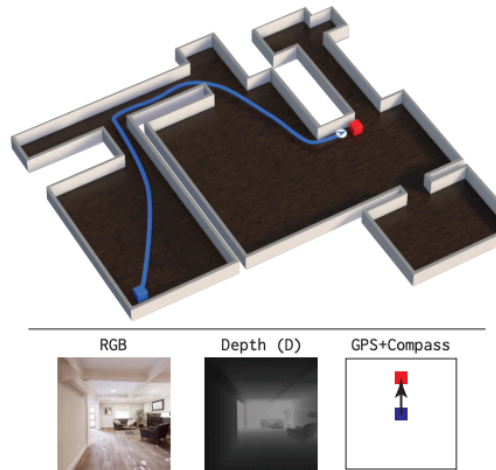


iGibson [Xia et al. 2020]



RLBench [James et al. 2019]

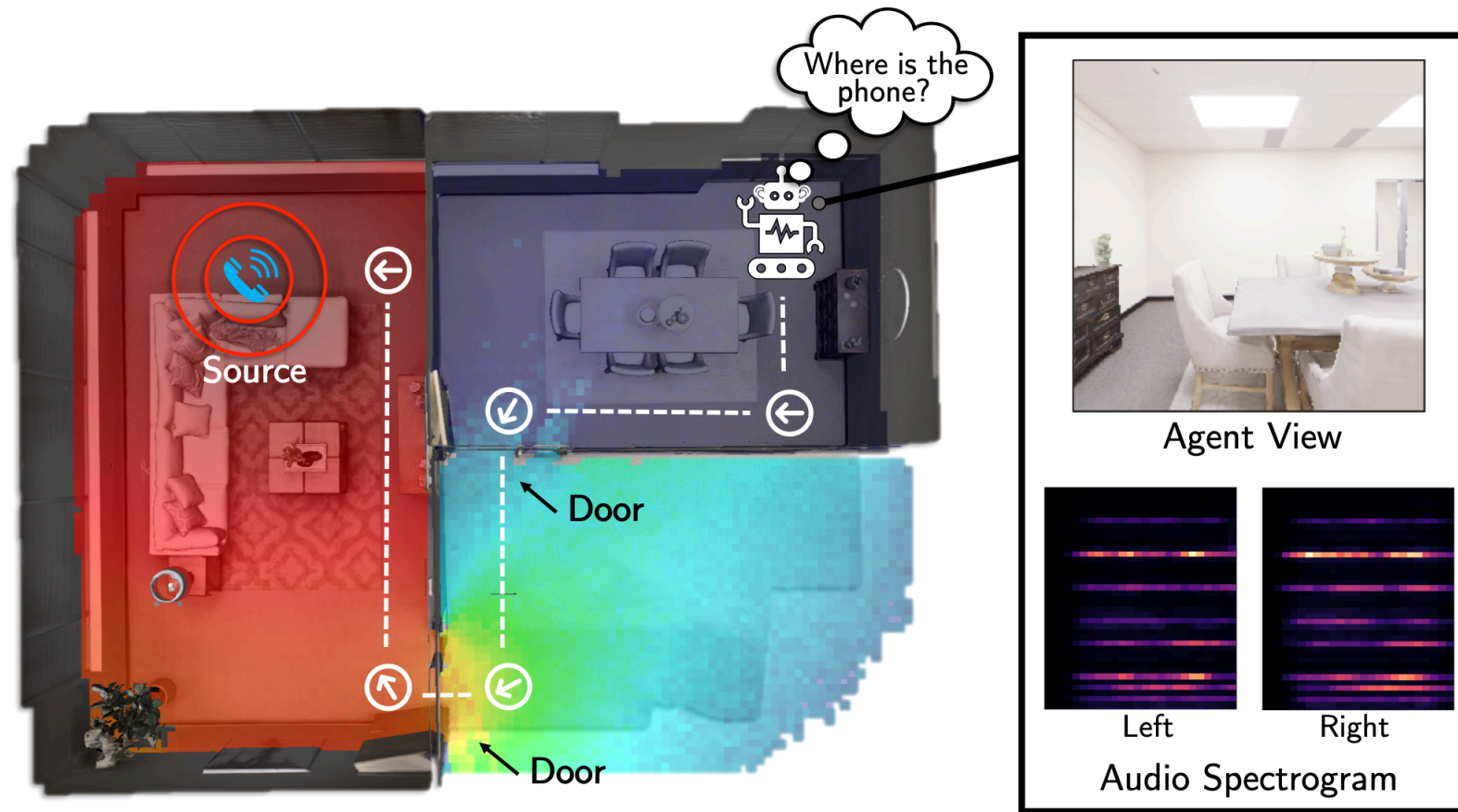
# Emerging trends: scale (& more speed)



DD-PPO: Learning Near-Perfect PointGoal Navigators from 2.5 Billion Frames [Wijmans et al. 2020]

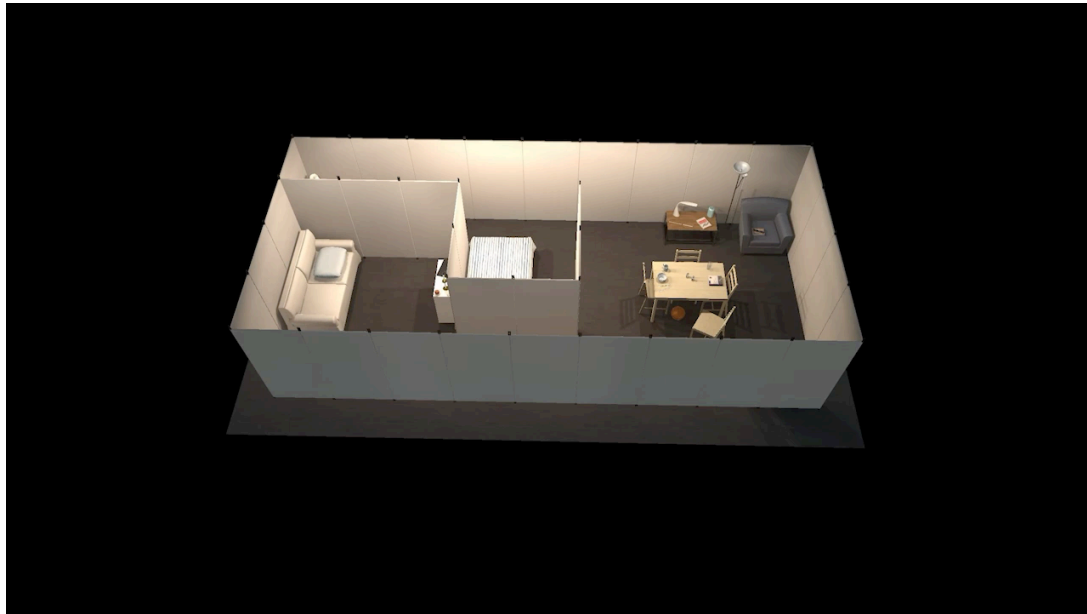
Sample Factory: Egocentric 3D Control from Pixels at 100000 FPS with Asynchronous Reinforcement Learning [Petrenko et al. 2020]

# Emerging trends: multimodality



Audio-Visual Embodied Navigation [Chen et al. 2020]

# Emerging trends: Sim2Real



RoboTHOR [Deitke et al. 2020]

Are We Making Real Progress in Simulated Environments?  
Measuring the Sim2Real Gap in Embodied Visual Navigation



Abhishek Kadian<sup>1\*</sup>



Joanne Truong<sup>2\*</sup>



Aaron Gokaslan<sup>1\*</sup>



Alex Clegg<sup>1\*</sup>



Erik Wijmans<sup>1,2</sup>



Stefan Lee<sup>2</sup>



Manolis Savva<sup>1,4</sup>



Sonia Chernova<sup>1,2</sup>



Dhruv Batra<sup>1,2</sup>

\* denotes equal contribution

facebook  
Artificial Intelligence Research

1

Georgia  
Tech

2

Oregon State  
University

3

SFU

4

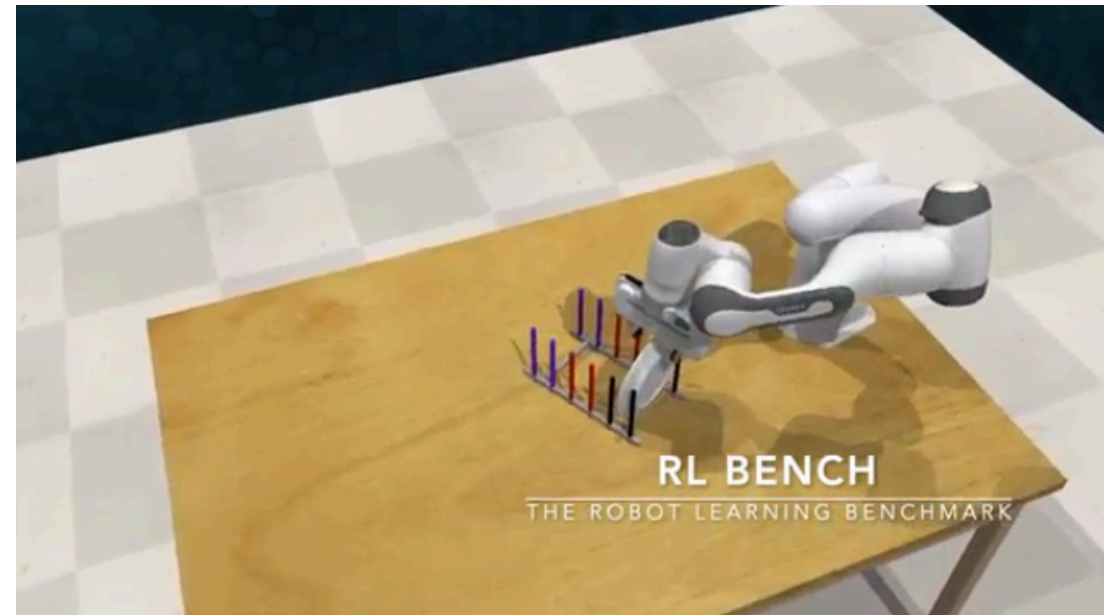
Sim2Real Coefficient [Kadian et al. 2020]

Graphics system challenges

# Challenge: “fast physics”



iGibson [Xia et al. 2020]



RLBench [James et al. 2019]

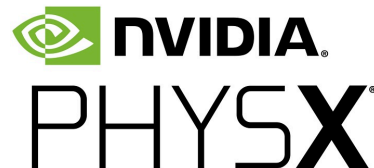
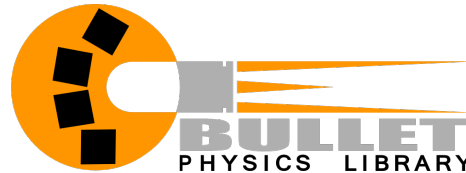
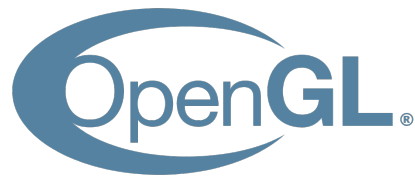
# Challenge: “GPU cohabitation”



Rendering

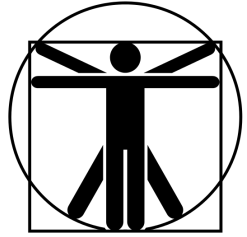
Physics

Learning

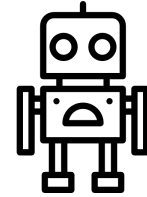




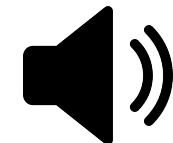
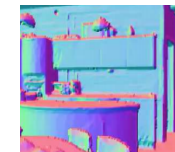
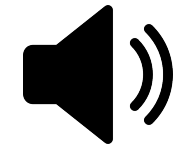
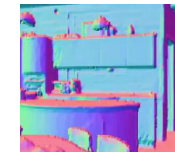
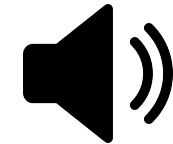
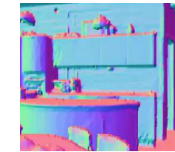
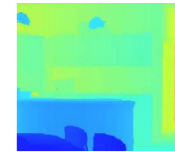
# Challenge: “not for eyeballs”



**Human:** 1080p @ 60Hz



**RL:** 84x84 @ 1000+ Hz



⋮

# Challenge: “asset soup”

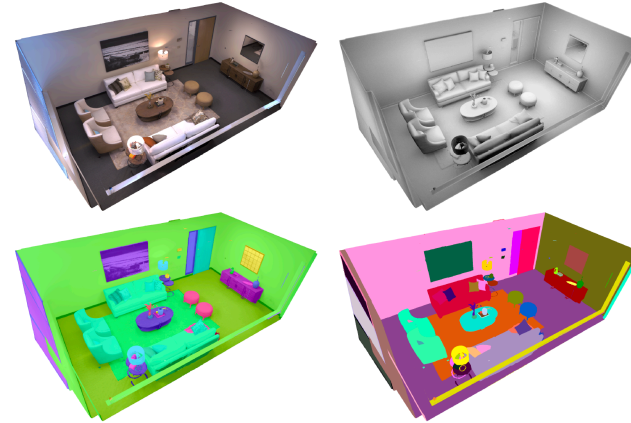
## AI2-THOR

120 virtual rooms



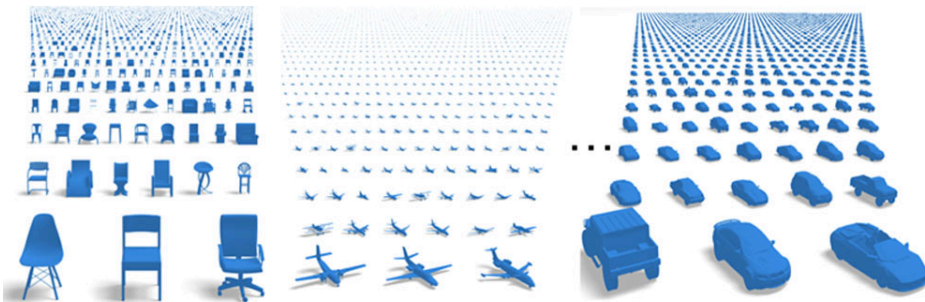
## Replica

18 near-photorealistic rooms



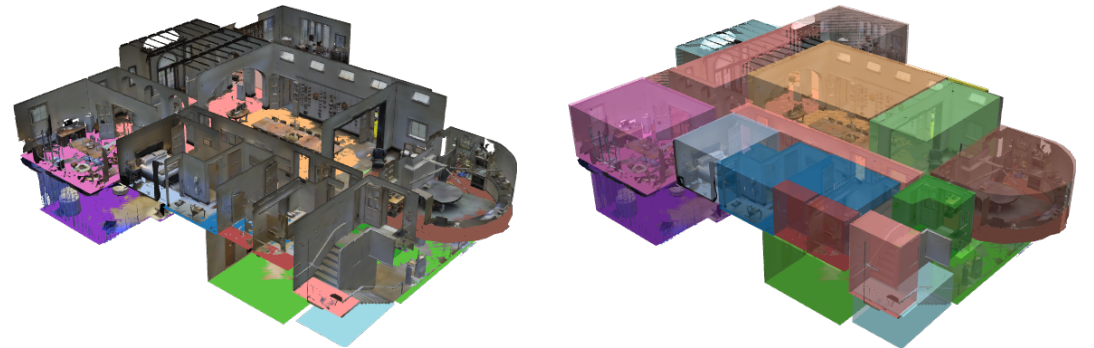
## ShapeNet

65K virtual objects



## Matterport3D

90 multi-floor house reconstructions



# Summary

## Trends

- Interaction
- Scale & more speed
- Multimodality
- Sim2Real

## Challenges

- “Fast physics”
- “GPU cohabitation”
- “Not for eyeballs”
- “Asset soup”

# Takeaway messages

- Growing interest in embodied AI
- Simulation for embodied AI: new frontiers for GFX-ML systems
- Opportunities for broad impact!

# Visual Computing @ Simon Fraser University

We're hiring at all levels! MSc, PhD, postdocs, researchers, faculty 😊



SFU campus over Metro Vancouver



Greg Mori



Ping Tan



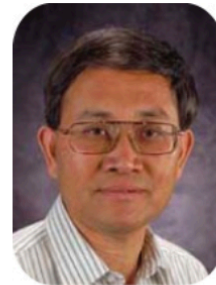
Manolis Savva



Angel Chang



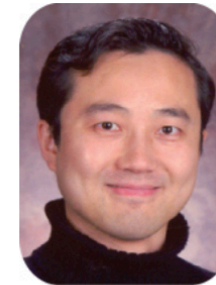
Yağız Aksoy



Ze-Nian Li



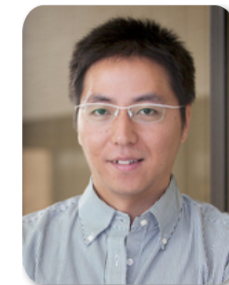
KangKang Yin



Richard Zhang



Eugene Fiume



Yasutaka Furukawa

SFU VC Faculty

[msavva@sfu.ca](mailto:msavva@sfu.ca)

Thank you!

