

Day 20

Speaker: Prof Dhruv Batra, FAIR & Georgia Tech

Title: From Disembodied to Embodied AI.

## From Disembodied to Embodied AI

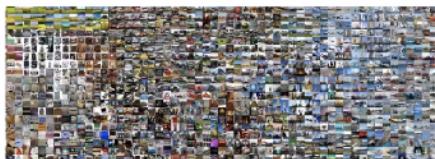
Dhruv Batra



## What is Embodied AI?

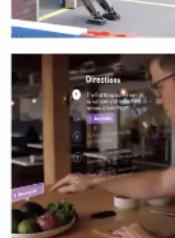
### Disembodied AI

Static datasets



### Embodied AI

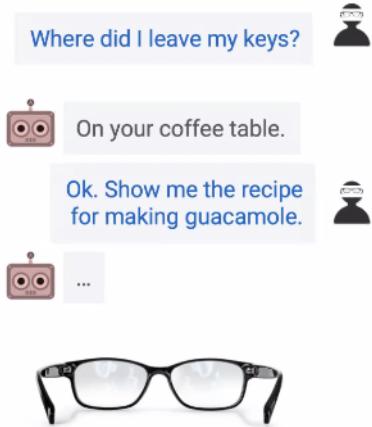
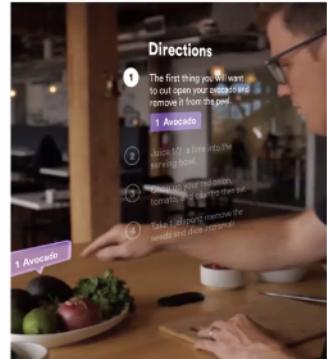
Agents acting in environments



- Is there smoke in any room around you? Yes, in one room.
- Go there and look for people.
- ...
- Where did I leave my keys?
- On your coffee table.
- Show me the recipe for making guacamole.
- ...

## Embodied Agents in Real & Virtual Environments

Virtual agent  
on smart glasses  
taking virtual actions  
talking to humans  
in natural language



12

## Plan for this talk

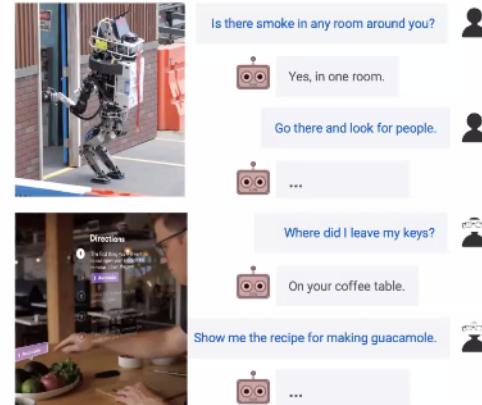
- Simulation as a test-bed
  - Habitat: A Platform for Embodied AI [ICCV 19]
  - Habitat 2.0: Training Home Assistants to Rearrange their Habitat [Arxiv 21]
- Emergence of intelligent behavior at scale
  - DD-PPO: Learning Near-Perfect PointGoal Navigators from 2.5 Billion Frames [ICLR 20]

15

# The Challenge of Embodied AI

- Hardware doesn't exist, or
- Reality is
  - Slow
  - Dangerous
  - Expensive
  - Not easy reproducible/controllable

## Agents acting in environments



17

## Common Solution: Sim2Real

Simulation



Reality



18

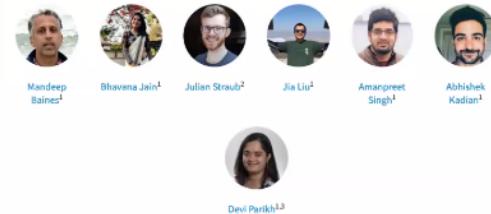


## Habitat: A Platform for Embodied AI Research

### Team: Current Contributors



### Team: Past Contributors



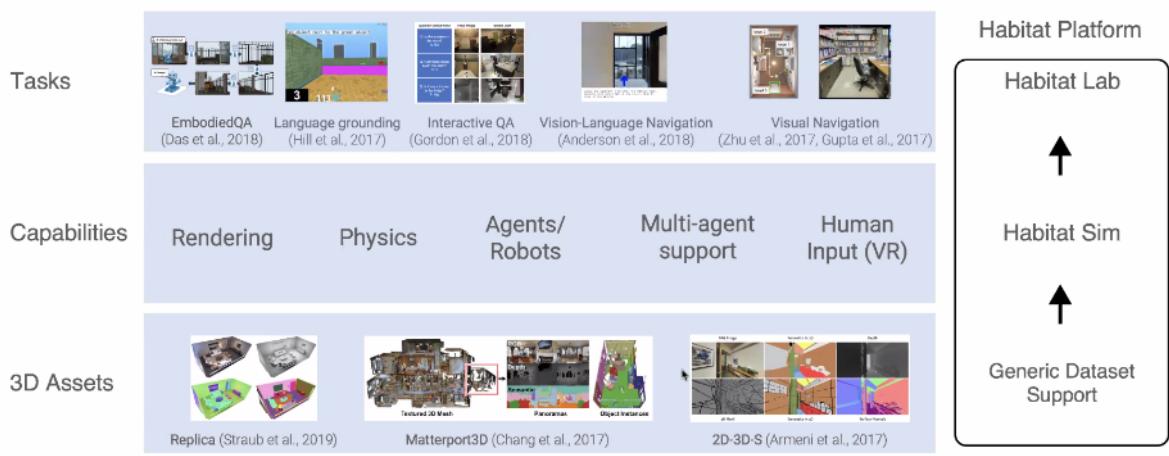
# Our Vision

- Create the ImageNet/COCO/VQA of Embodied AI
  - Dataset → Simulator → Task → Benchmark Challenge



20

## Habitat: A Simulation Platform for Embodied AI



21

## FRL Surreal Team: high quality 3D reconstructions



The Replica Dataset: A Digital Replica of Indoor Spaces [Straub et al. 2019]

## FRL Surreal Team: high quality 3D reconstructions



## Bring Your Own Scan: Virtualizing Reality



## Object datasets: YCB, Ikea



Yale-CMU-Berkeley dataset for robotic manipulation research  
Berk Cali<sup>1</sup>, Arjen Stigle<sup>2</sup>, James Bran<sup>3</sup>, Azrael Walorski<sup>4</sup>, Kurt Konstig<sup>5</sup>,  
Siddhartha Srinivasa<sup>6</sup>, Pieter Abbeel<sup>7</sup> and Asenka M Dular<sup>8</sup>



## Instruction Following



## Embodied Question Answering



## Navigation

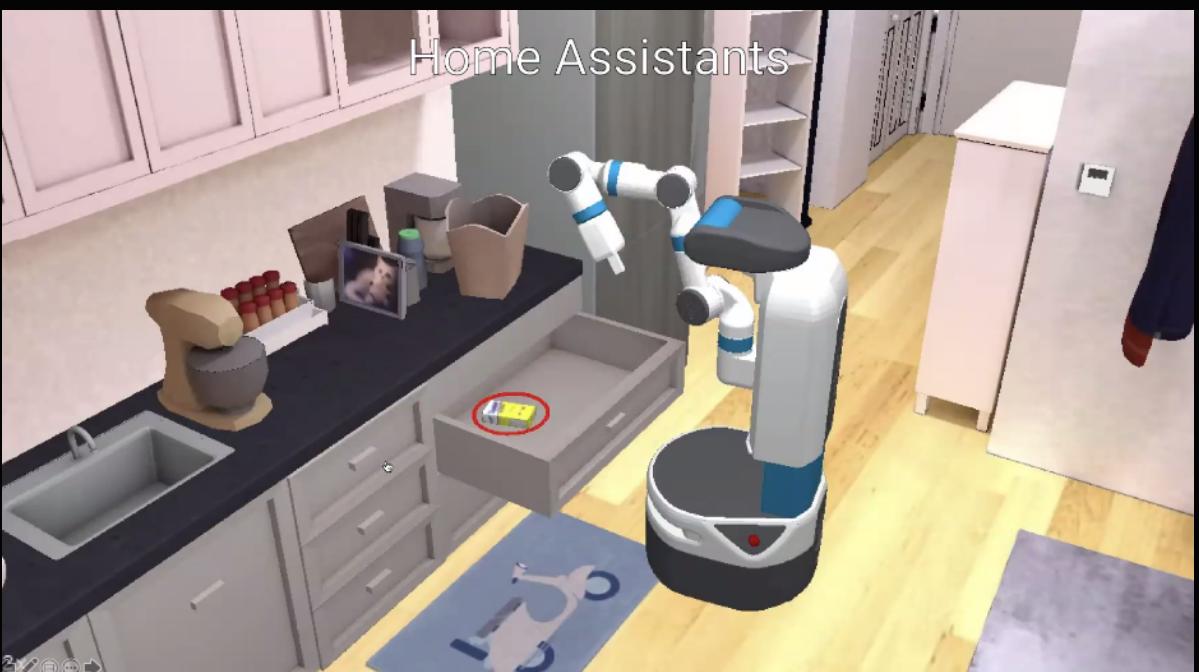
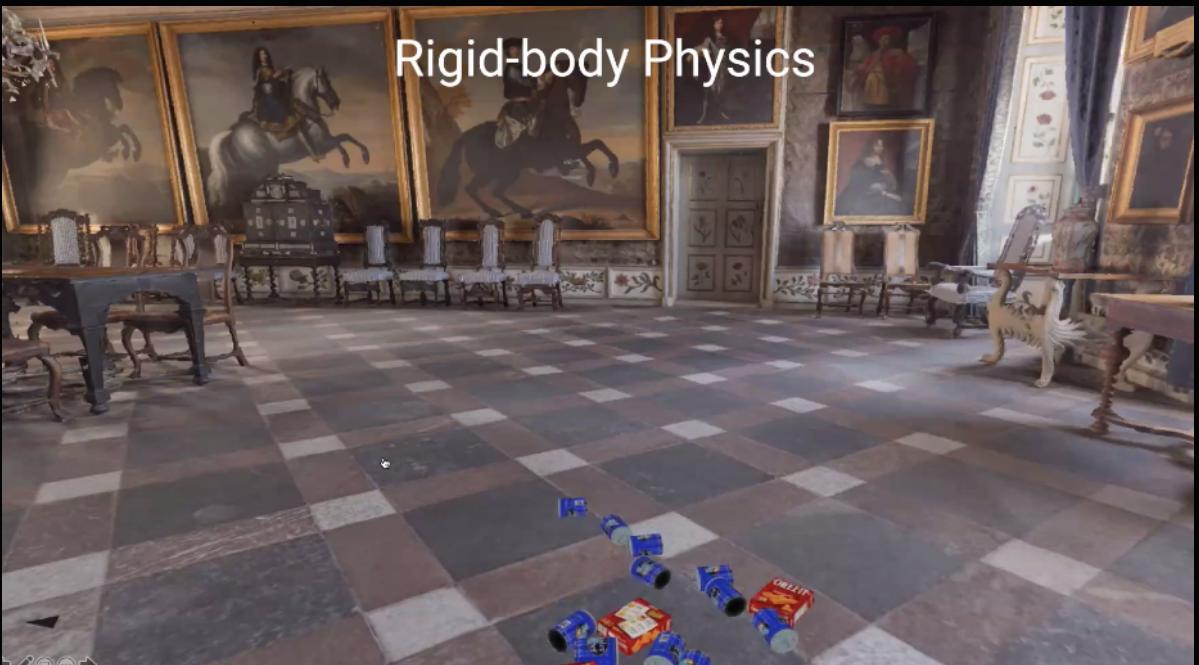


Habitat 1.0  
[ICCV '19]

## Mobile Manipulation



Habitat 2.0



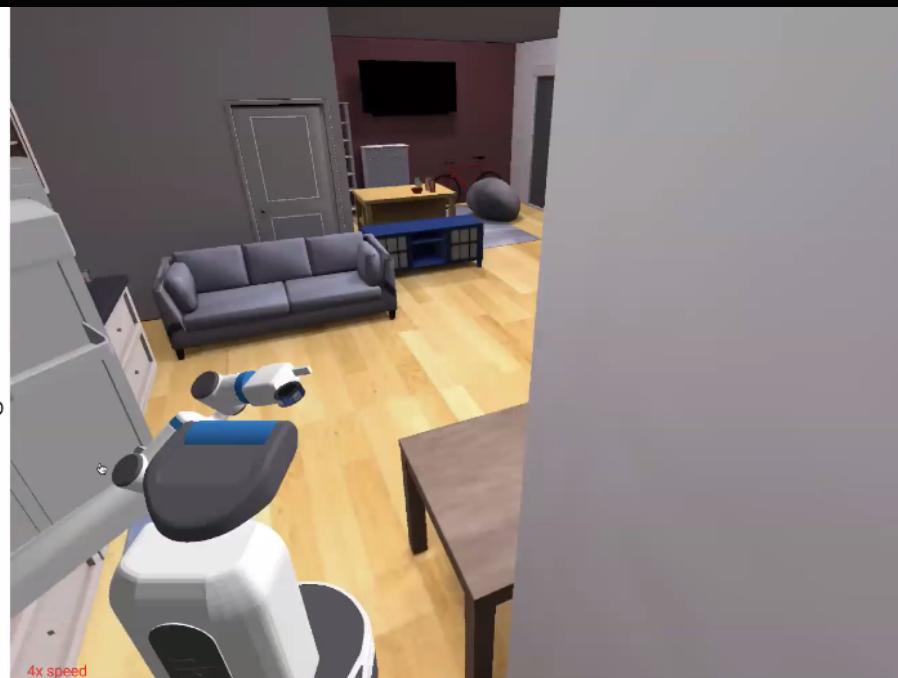
## Tidy House

- Move 5 objects from start to goal location



## Prepare Groceries

- Object 1 from fridge to counter
- Object 2 from fridge to table
- Object 3 from counter to fridge



## Set Table

- Bowl from drawer to table
- Apple from fridge into bowl
- (Fridge closed)
- (Drawer closed)





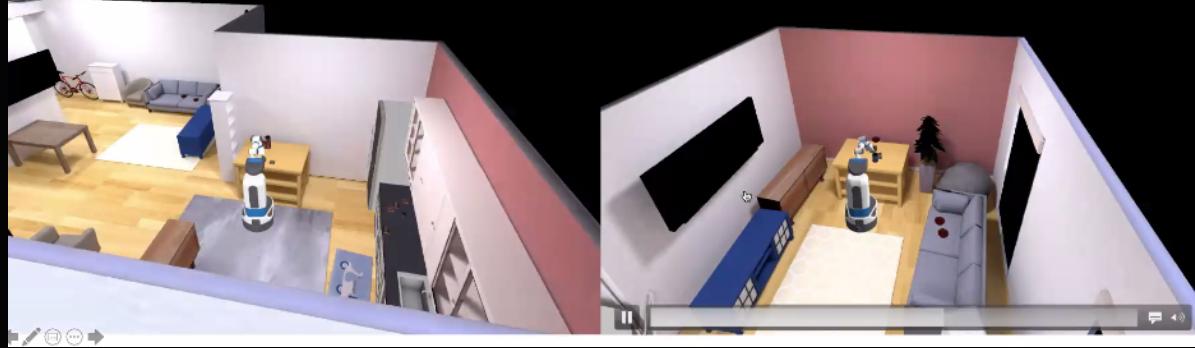
## Generalization to New Objects



## Generalization to New Layouts

**Training Layout**

**Testing Layout**



# Why not use a video game engine?



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

56

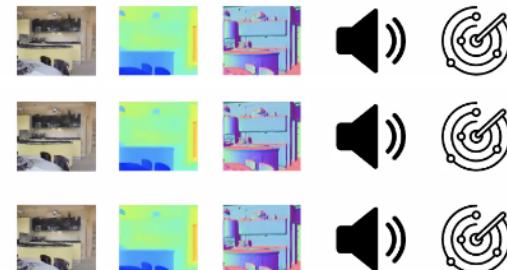
## Human needs vs AI needs



Human: 1080p @ 60Hz



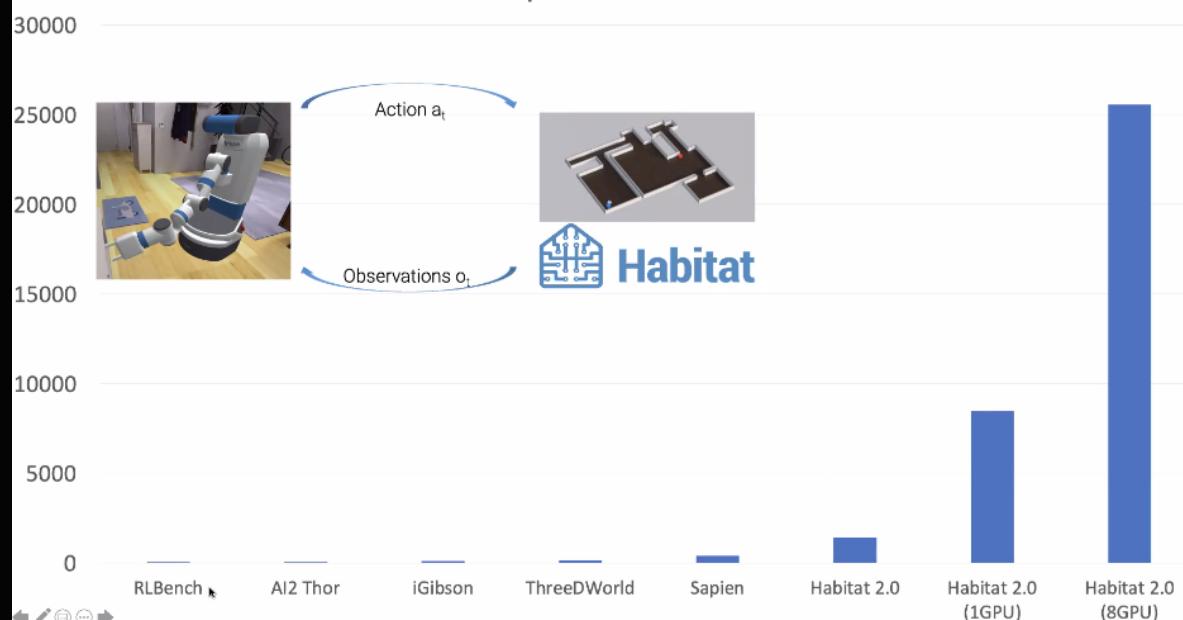
AI: 84x84 @ 1000+ Hz

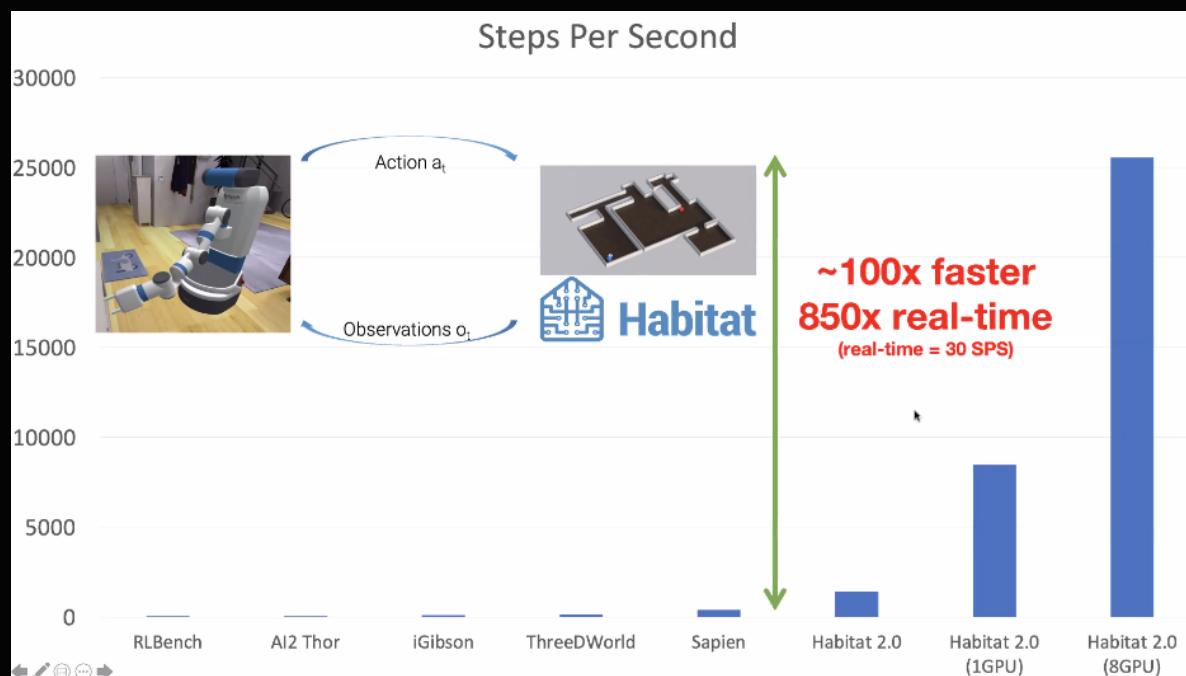


⋮

58

## Steps Per Second





## Why does speed matter?

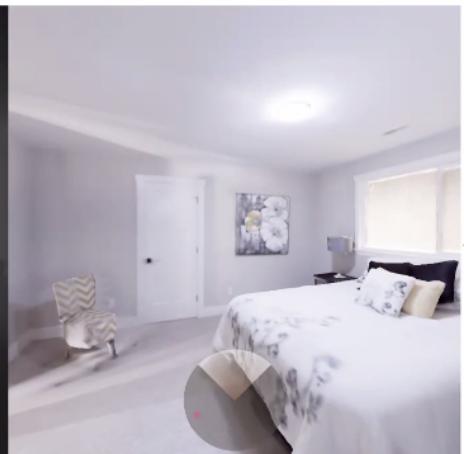
Because you can now run experiments you couldn't before.

## Plan for this talk

- Simulation as a test-bed
  - Habitat: A Platform for Embodied AI [ICCV 19]
  - Habitat 2.0: Training Home Assistants to Rearrange their Habitat [Arxiv 21]
- Emergence of intelligent behavior at scale
  - DD-PPO: Learning Near-Perfect PointGoal Navigators from 2.5 Billion Frames [ICLR 20]



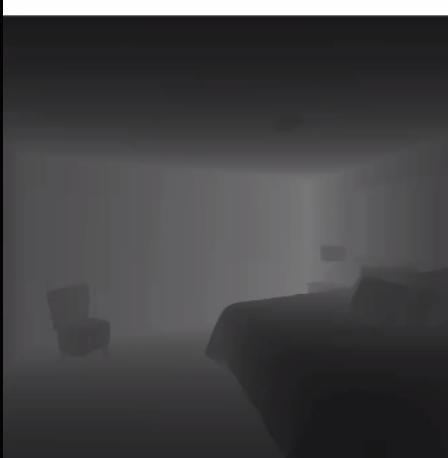
Depth



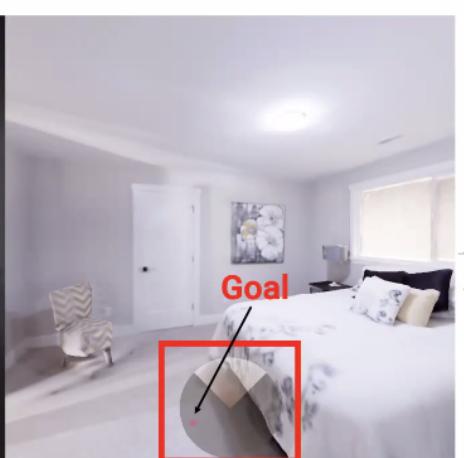
RGB and GPS+Compass



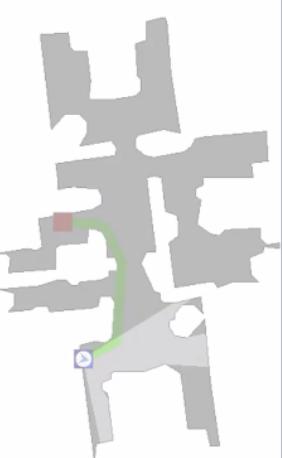
Top Down Map



Depth



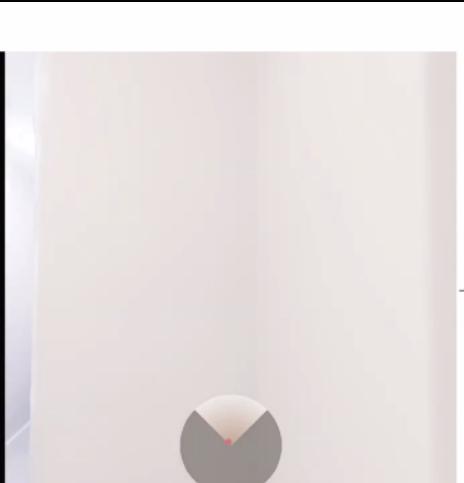
RGB and GPS+Compass



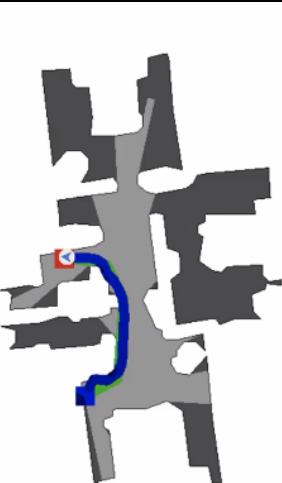
Top Down Map



Depth



RGB and GPS+Compass



Top Down Map

# Classical Robots Pipeline

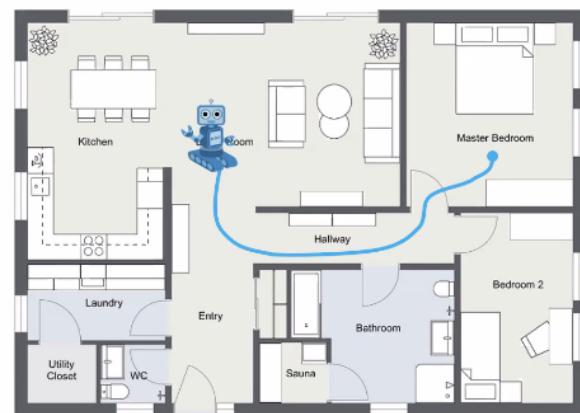
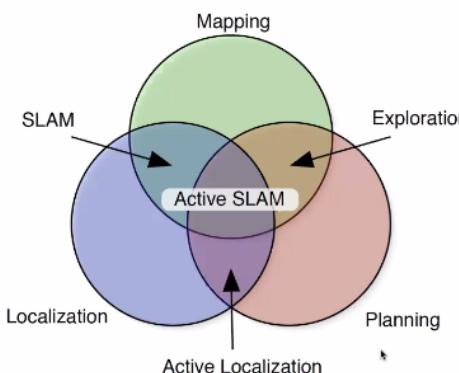


Image Credit: Nathaniel Fairfield, PhD Thesis 2009

82

## Decentralized Distributed PPO: Learning Near-Perfect PointGoal Navigators from 2.5 Billion Frames



Erik Wijmans



Abhishek Kadian



Ari Morcos



Stefan Lee



Irfan Essa



Devi Parikh



Manolis Savva



Dhruv Batra

facebook

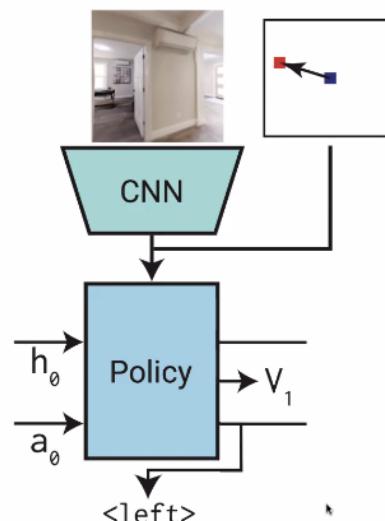
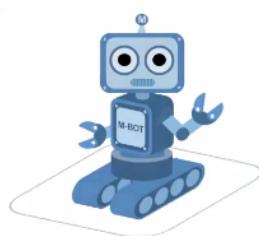
Artificial Intelligence Research

Georgia Tech

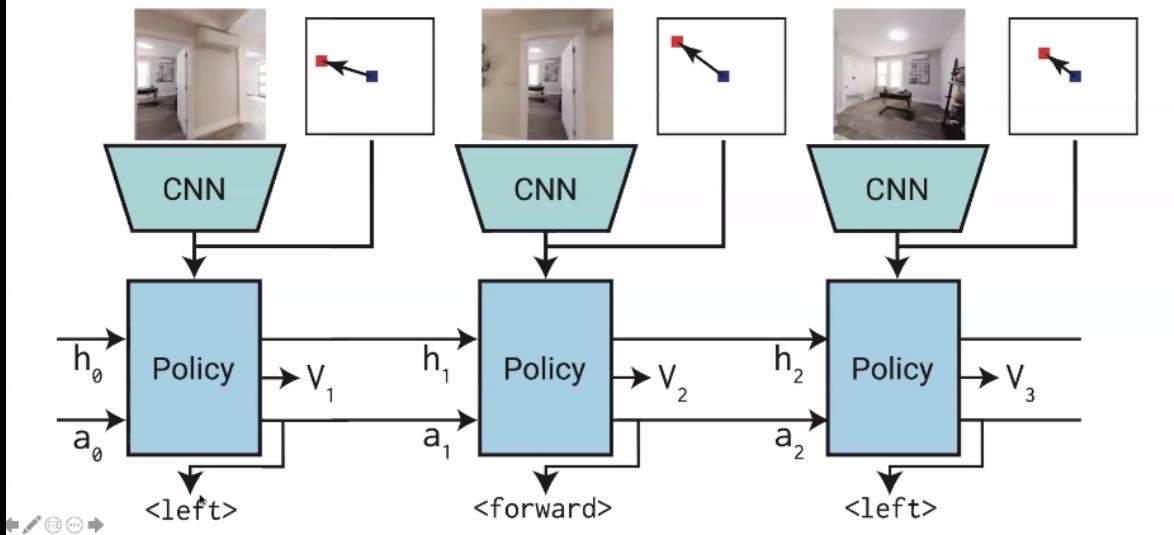
SFU

Oregon State University

## Agent and Model Design

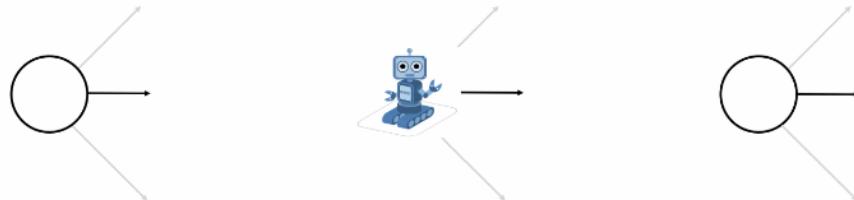


## Agent and Model Design



## Agent and Model Design

$$r_t = \begin{cases} s^* & \text{if goal is reached} \\ \underbrace{(d_{t-1} - d_t)}_{\text{reward shaping}} - \underbrace{\lambda}_{\text{efficiency}} & \text{otherwise} \end{cases}$$



## Radical Empiricism

- No task-specific modules
  - No SLAM, no mapping, no path planning
- No domain-specific inductive biases
  - No spatial memory, knowledge of projective geometry or 3D
- No additional learning signals
  - No mapping supervision, auxiliary tasks
  - No expert demonstrations or imitation learning
  - No pre-training of representations
  - No look-ahead search trees
- No learning tricks
  - No curriculum, replay buffer, reset to partially-successful-states

# *Radical Empiricism*

- On-policy Episodic Reinforcement Learning (on steroids)
  - Dense rewards

95

How far can we **scale** model-free RL  
for embodied visual navigation?

96

## Decentralized Distributed PPO

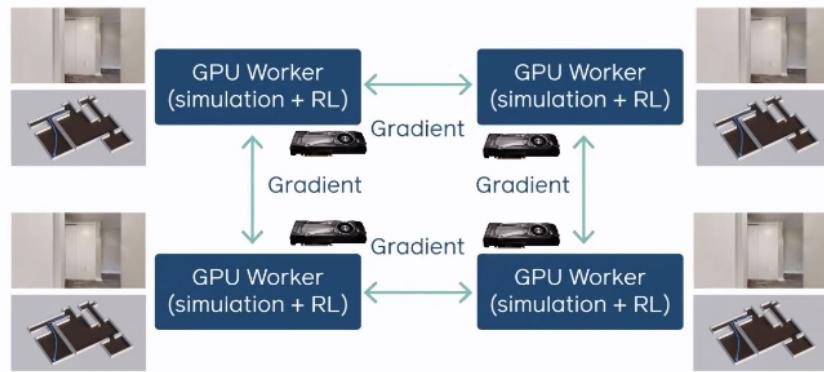
Collect experience



103

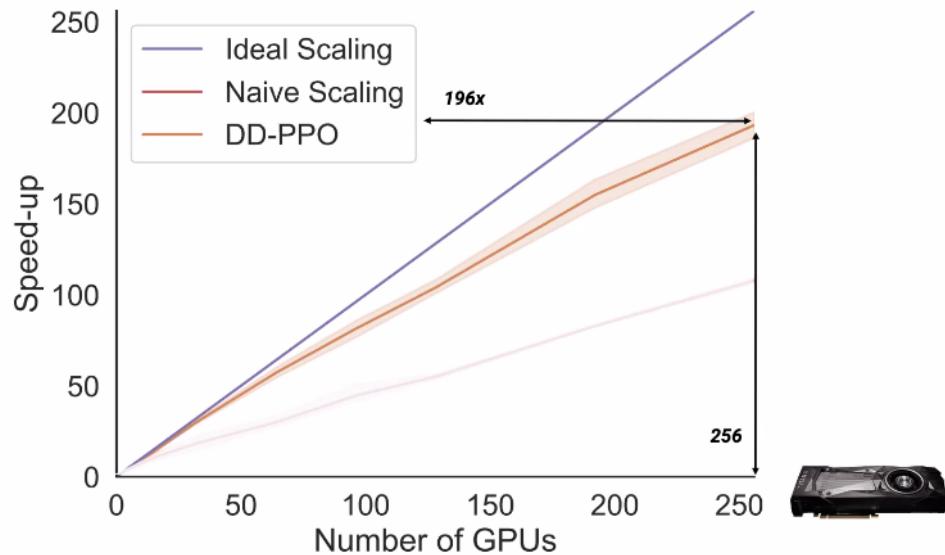
# Decentralized Distributed PPO

Synchronously optimize model

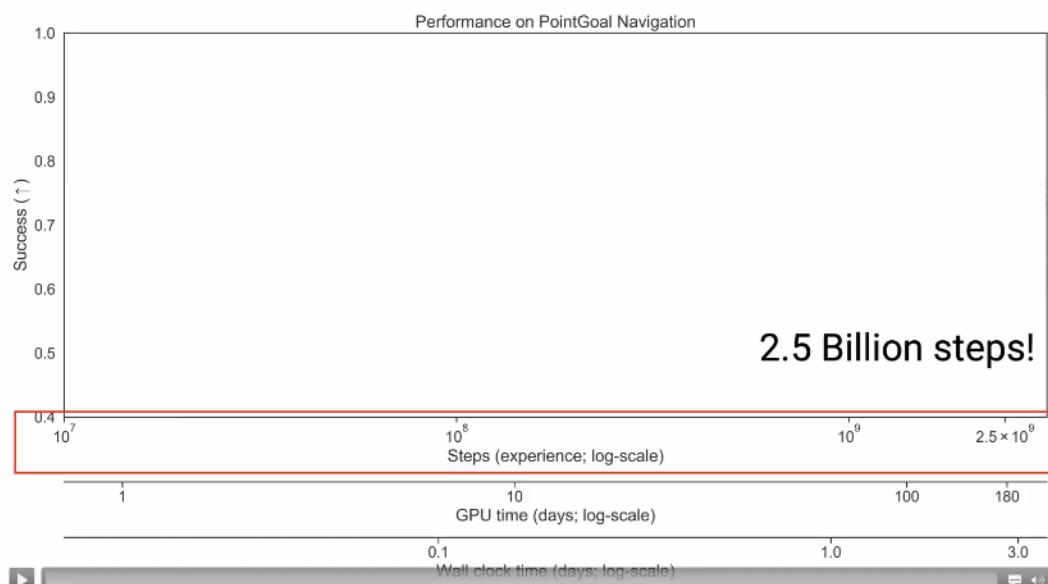


104

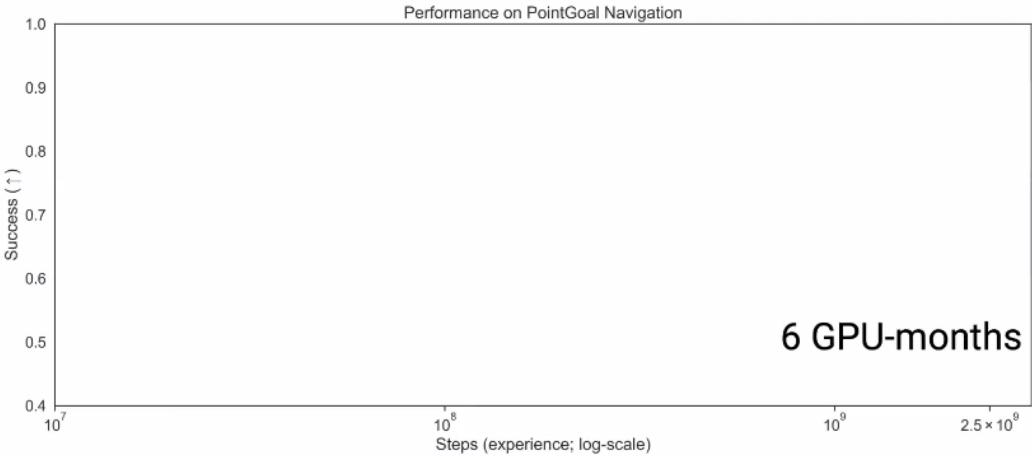
How does DD-PPO scale? Near-linearly!



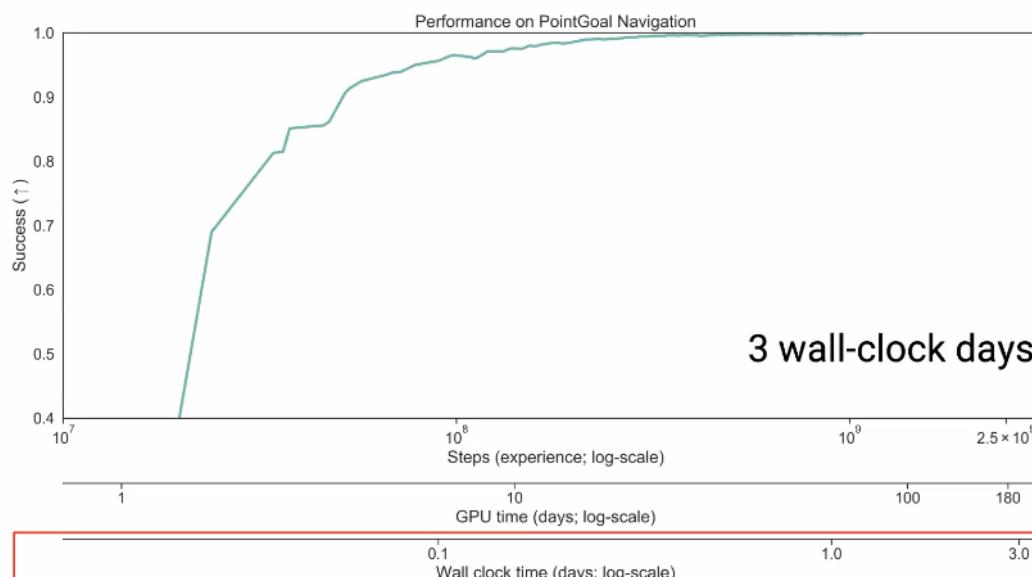
110



111



114

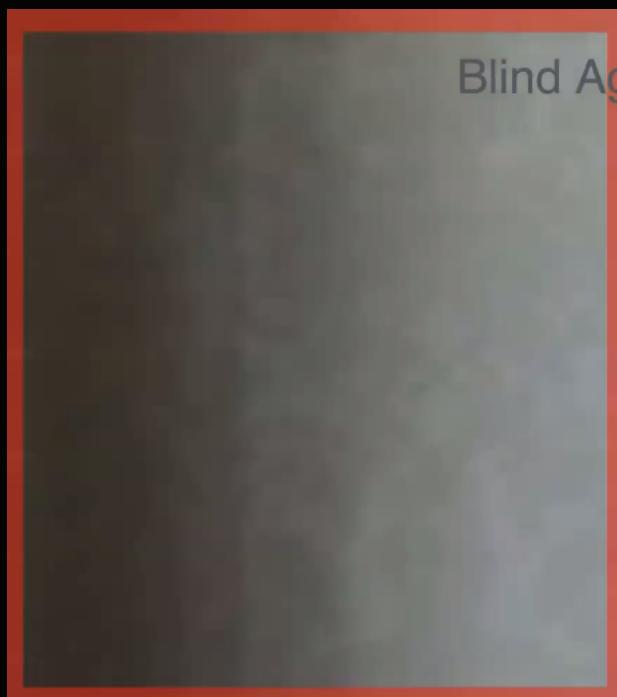


115





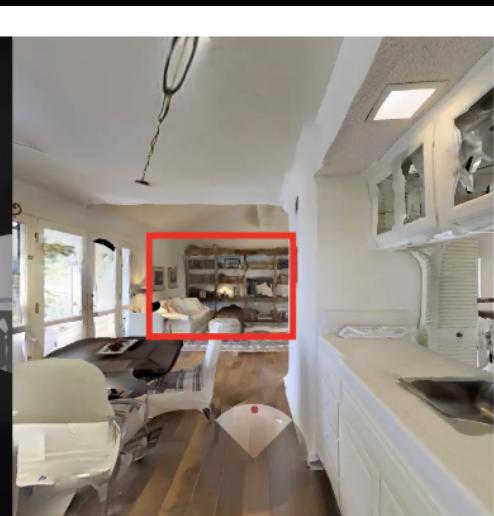
Depth Agent (RL)



Blind Agent (RL)



Depth

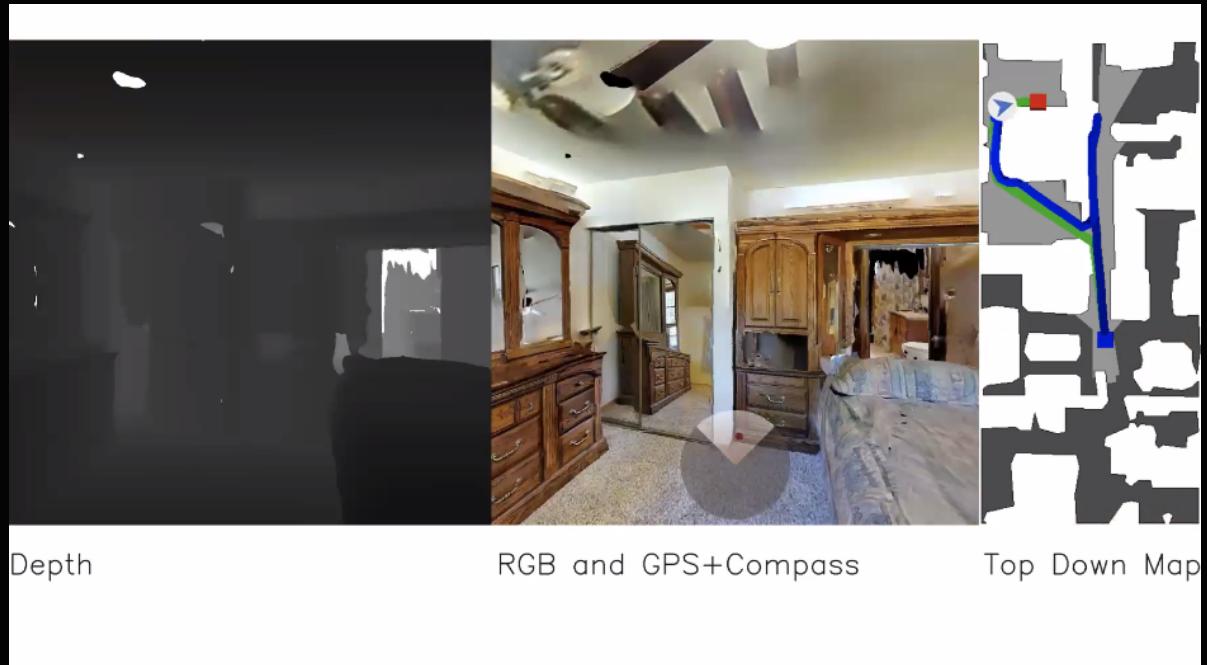
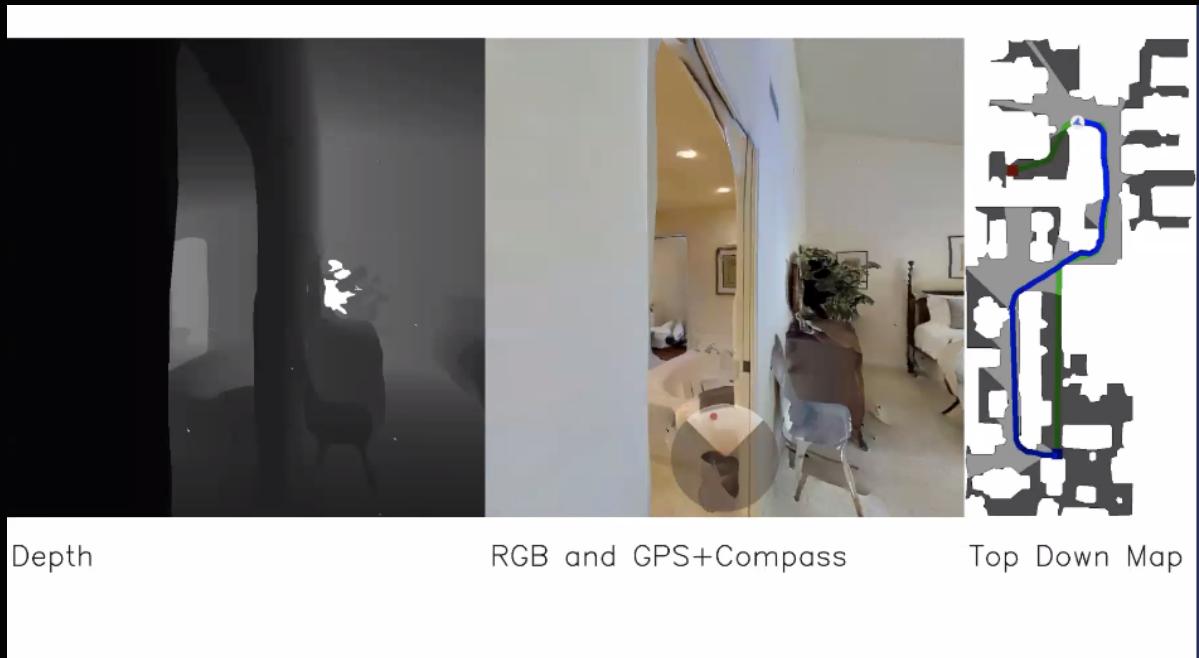


RGB and GPS+Compass



Top Down Map

However, there is wall straight



# How far can we scale model-free RL for embodied visual navigation?

- Surprisingly far!
    - PointNav w/ RGBD + GPS+Compass: Near-perfect
  - Not far enough yet (but rapidly improving)

	2019	2020	2021
◦ PointNav w/ RGBD (no localization):	15%	→ 28%	→ 96%

# How far can we scale model-free RL for embodied visual navigation?

- Surprisingly far!

- PointNav w/ RGBD + GPS+Compass: Near-perfect

- Not far enough yet (but rapidly improving) 2019 2020 2021

- PointNav w/ RGBD (no localization): 15% → 28% → 96%

- Object-Goal Navigation: 3-5% → 21% → 30%