

Lecture 17:

High-Throughput World Simulation for Agent Training

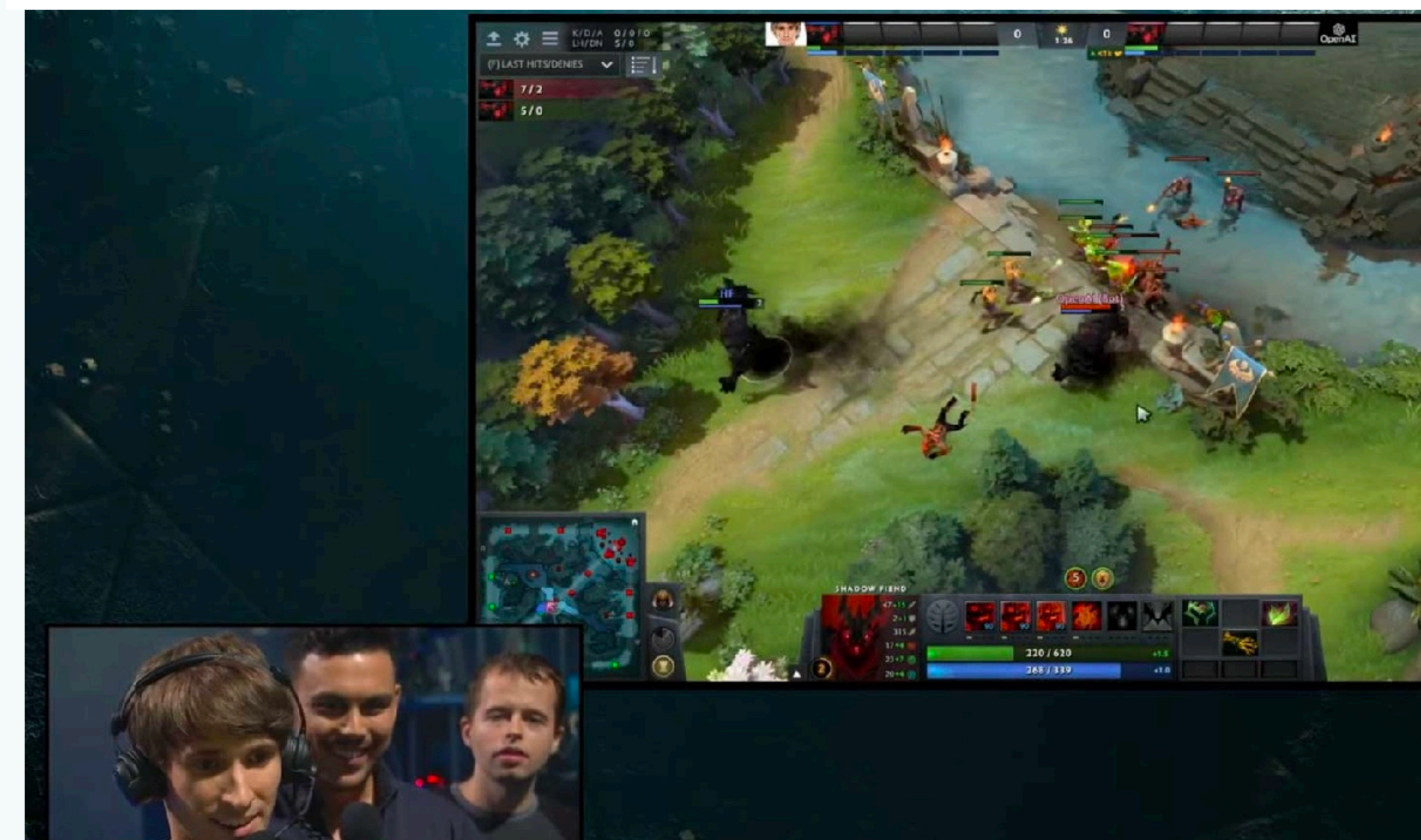
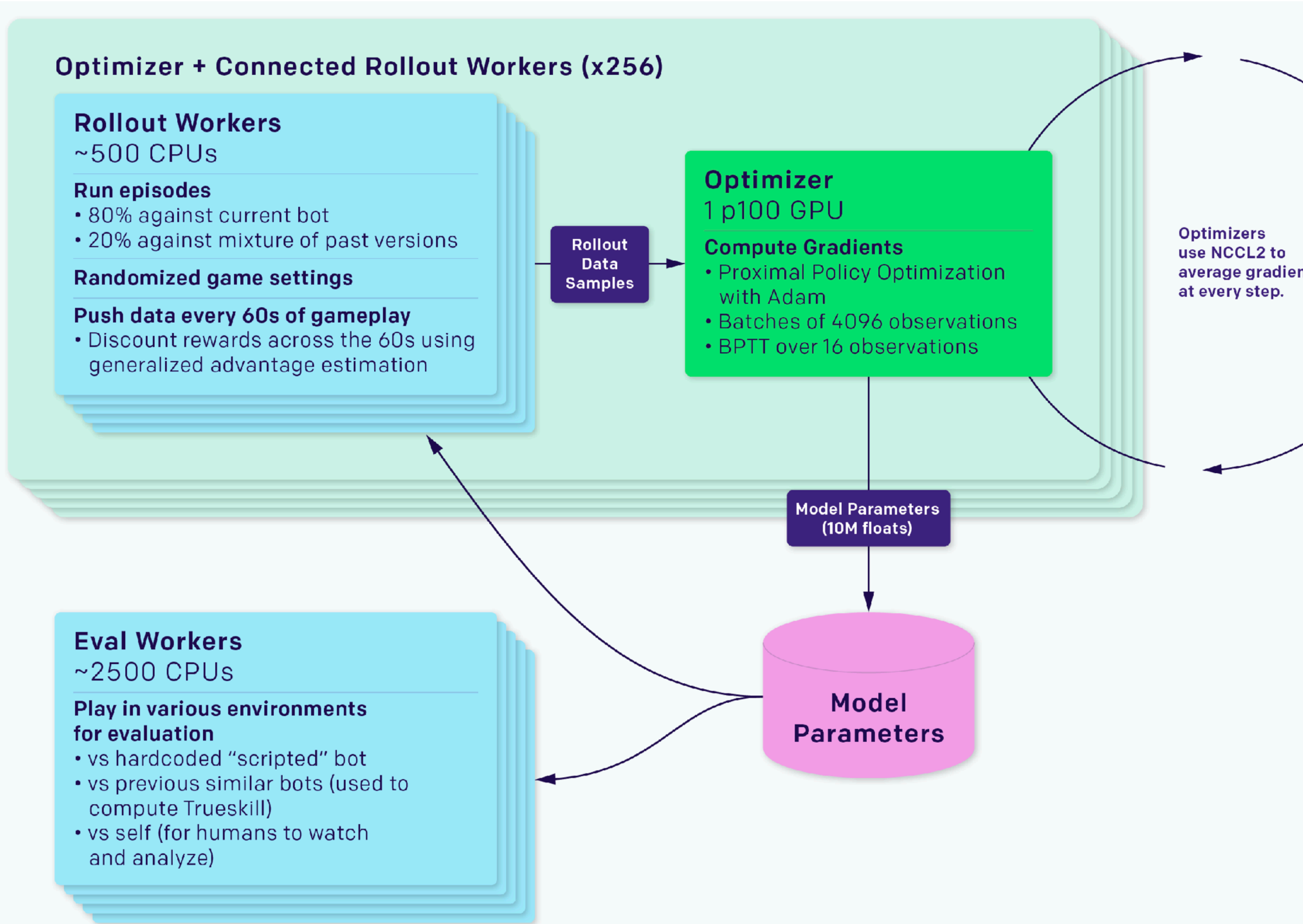
**Visual Computing Systems
Stanford CS348K, Spring 2026**

Experience collection for Dota 2 bots

OpenAI Five stats

OPENAI FIVE

CPUs	128,000 preemptible CPU cores on GCP
GPUs	256 P100 GPUs on GCP
Experience collected	~180 years per day (~900 years per day counting each hero separately)
Size of observation	~36.8 kB
Observations per second of gameplay	7.5
Batch size	1,048,576 observations
Batches per minute	~60



Modern AI “bots”

Large-scale RL



4. Mass Scale Training Infrastructure

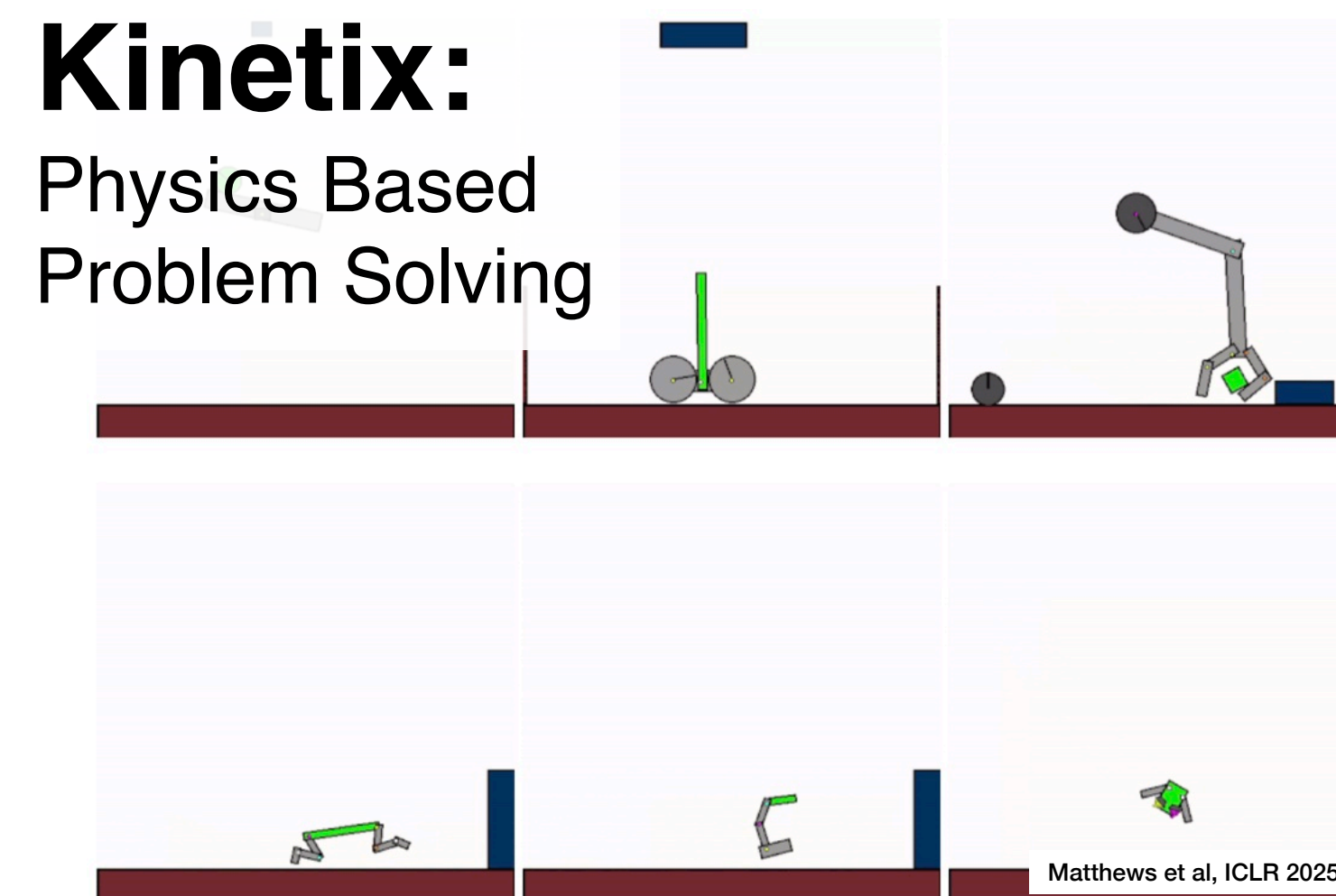
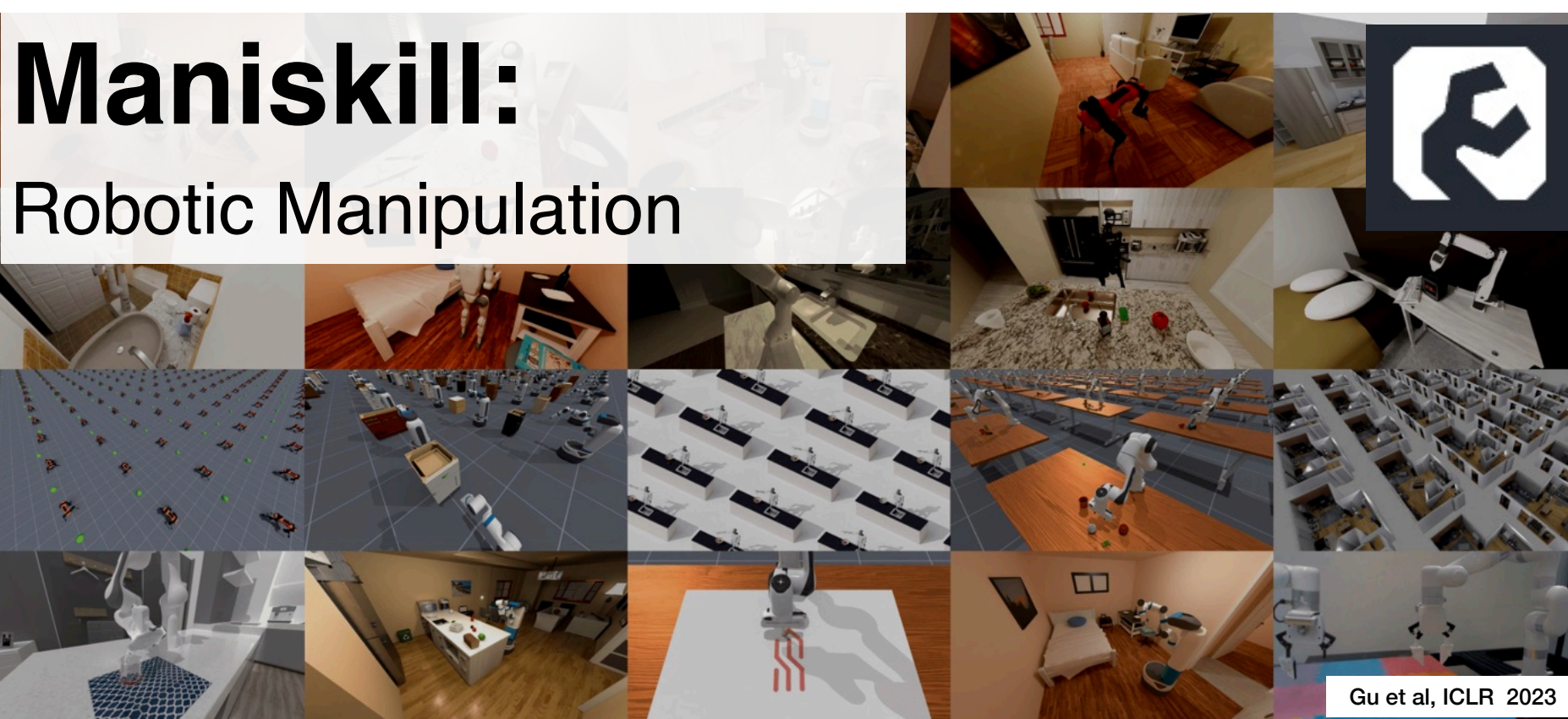
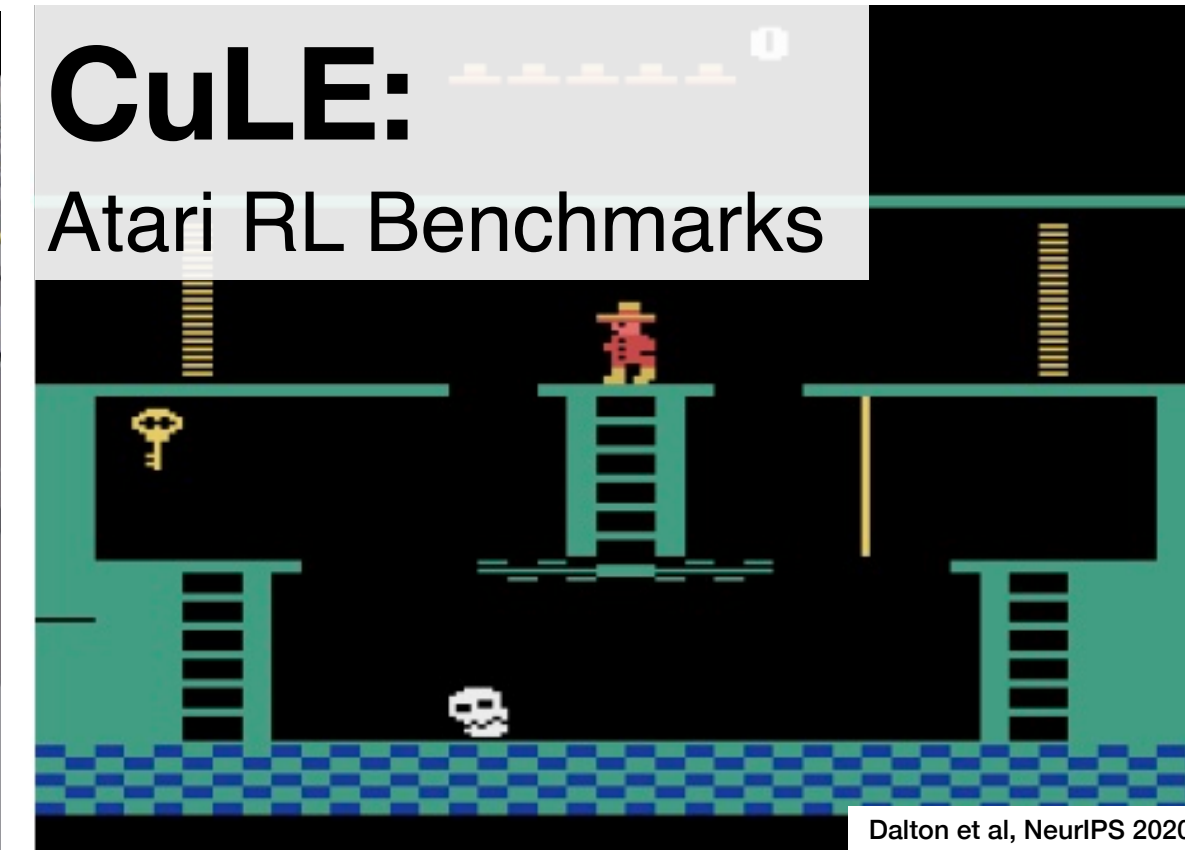
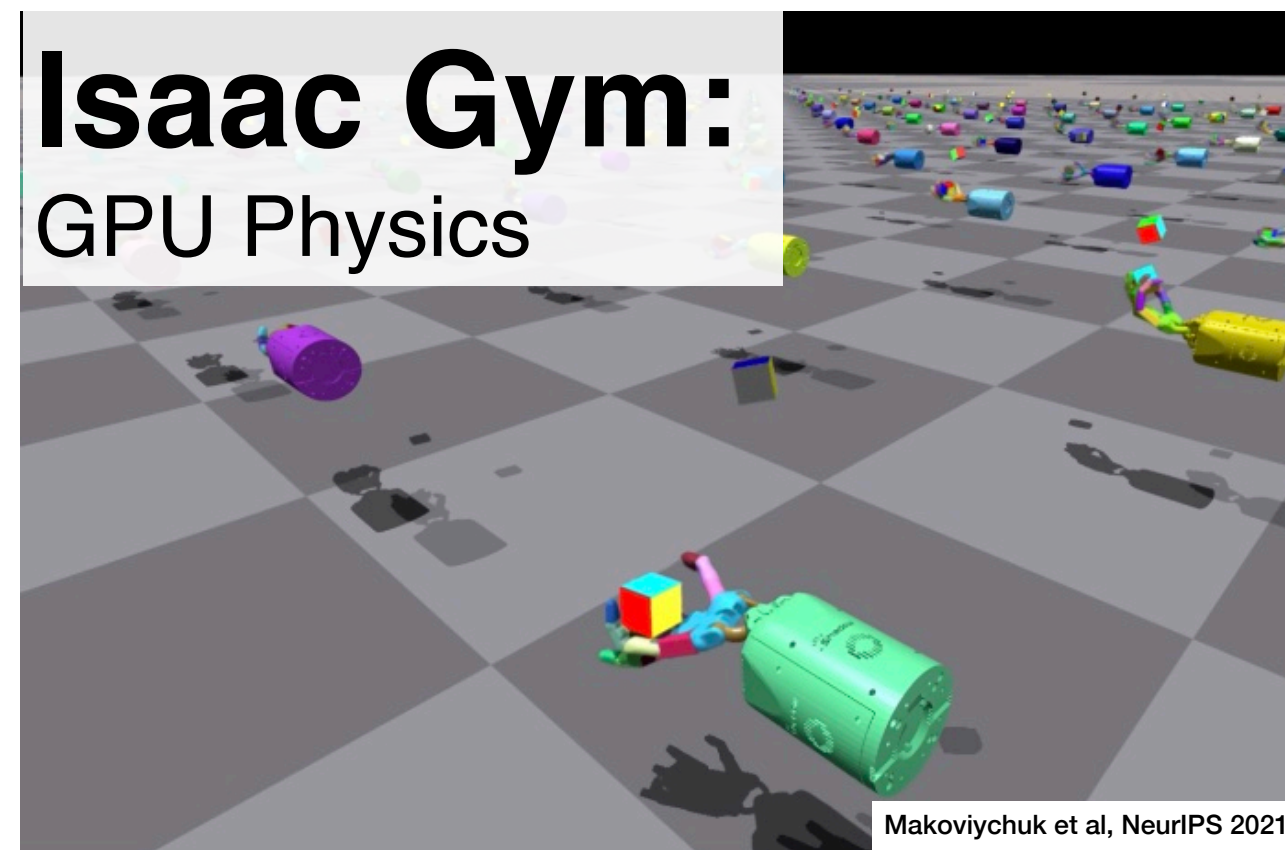
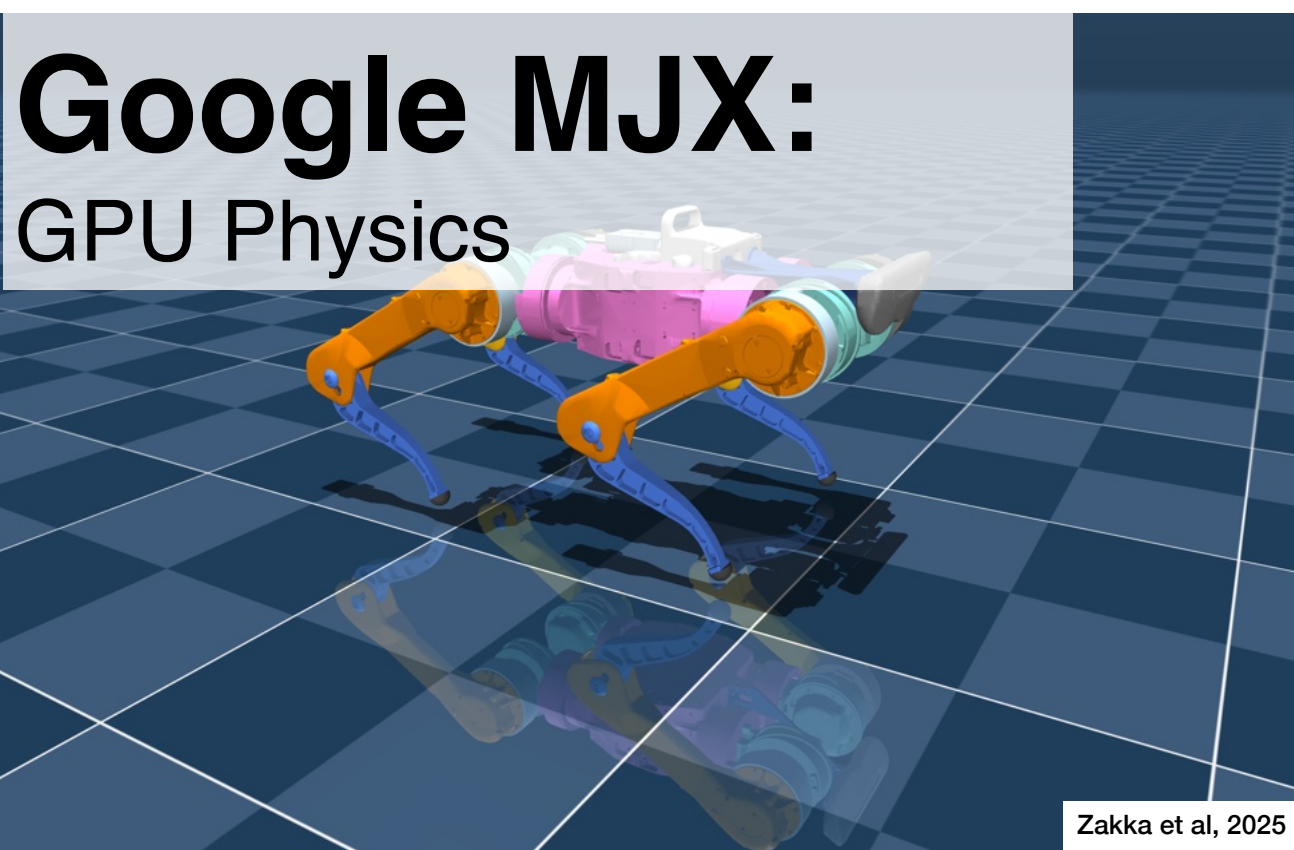
The DART platform had access to over 1,000 PlayStation 4 (PS4) consoles. Each was used to collect data for training GT Sophy or evaluate a trained version. The platform consisted of the necessary computing components (GPUs, CPUs) to interact with a large number of PS4s and support large scale training over an extended period of time.



for over 1500 epochs⁵ from scratch, which takes approximately ten days to train on a single A100 GPU and a dozen Playstation 4 systems. For our experiments, we reduced the training time of each iteration to 300 epochs, and only train the policy from the final iteration for 1500 epochs without a secondary replay buffer (to compare with GT Sophy) for evaluation. From now on, we refer to

Large-Batch World Simulation

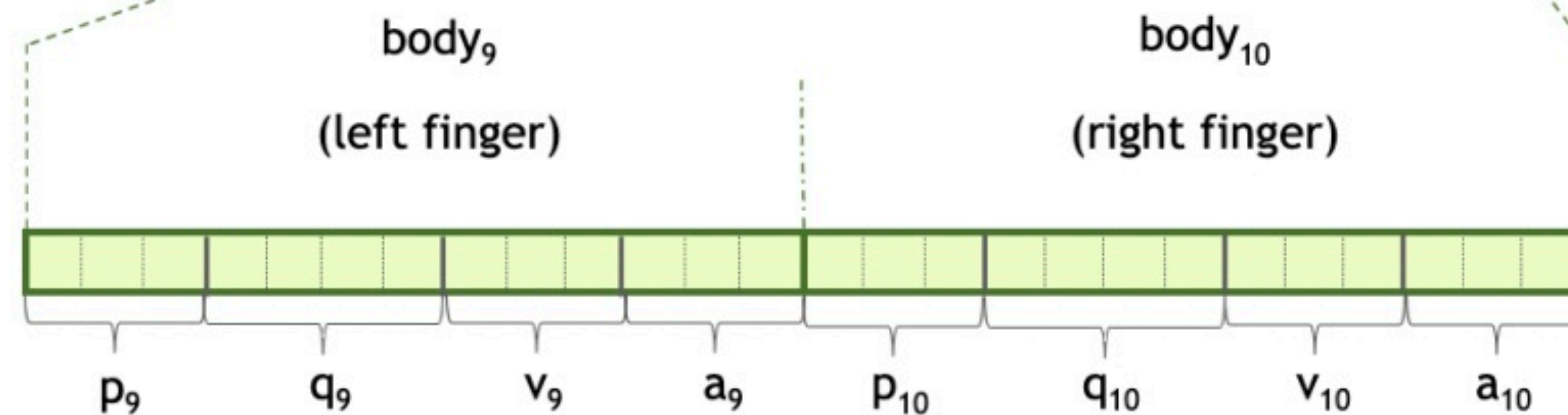
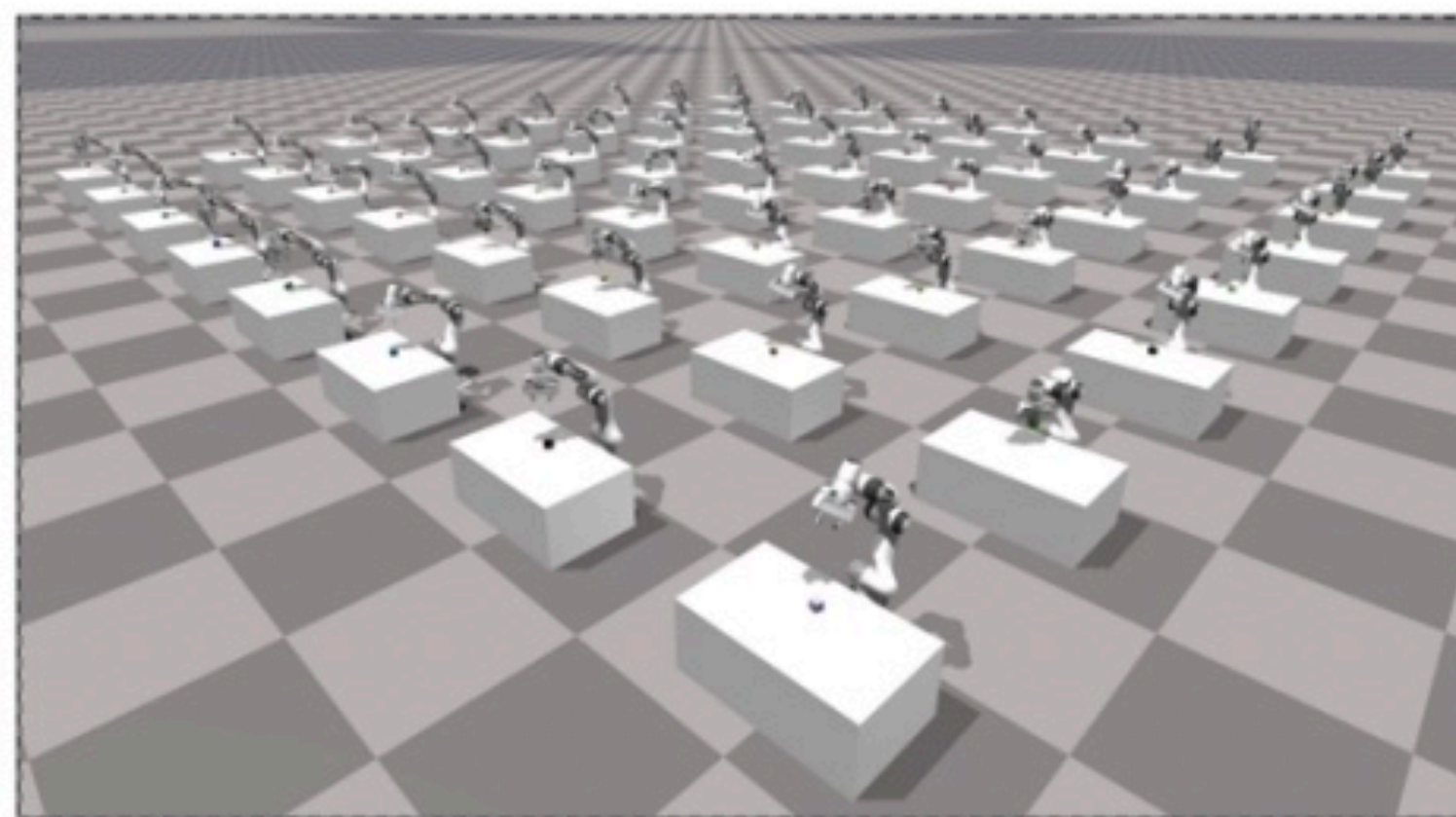
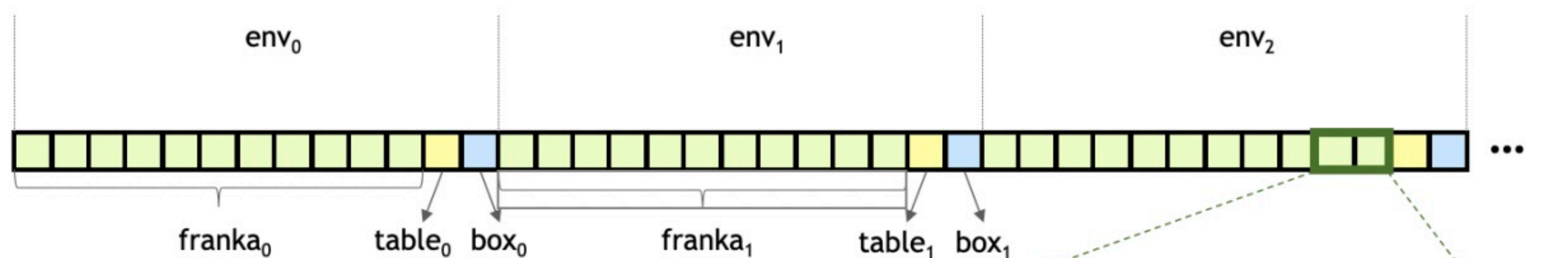
2022-2025 “batch simulators” that achieve millions of steps/sec by executing thousands of environments in parallel on a single GPU



But... requires a simulator rewrite for the GPU

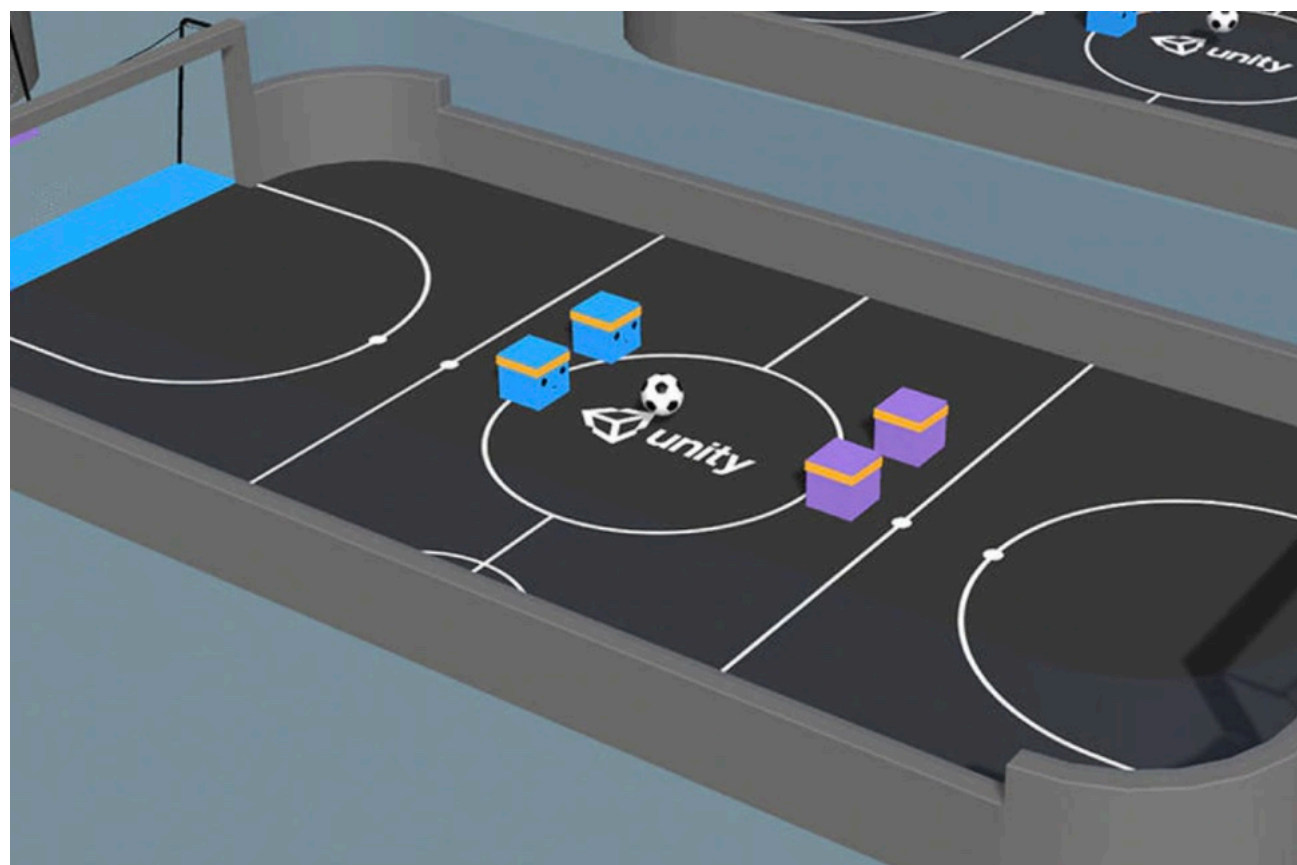
Example: Issac Gym

- Batched many-environment execution applied to rigid body physics sim
- Simulate 1000's of world environments simultaneously on the GPU

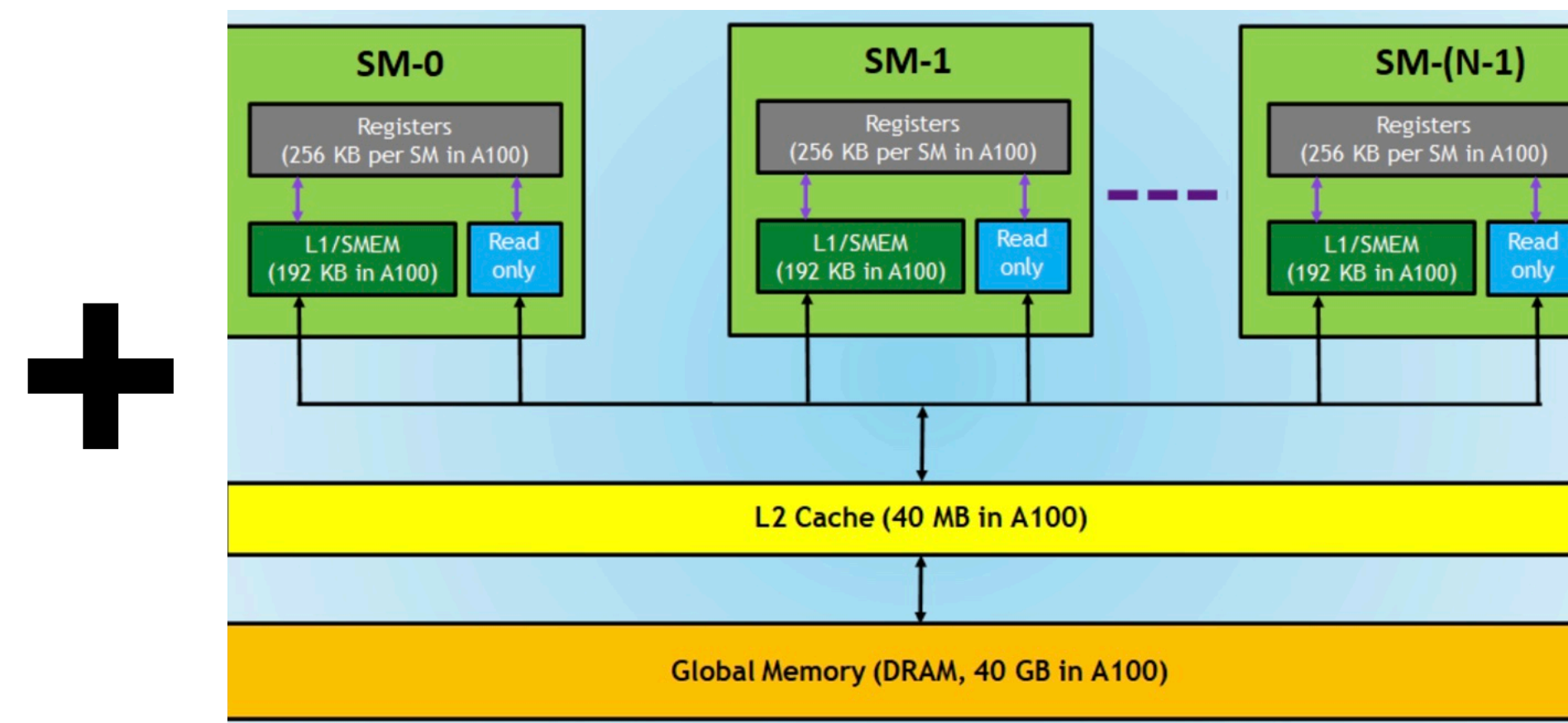


Significant engineering skill and effort to build new GPU-accelerated batch simulator from scratch!

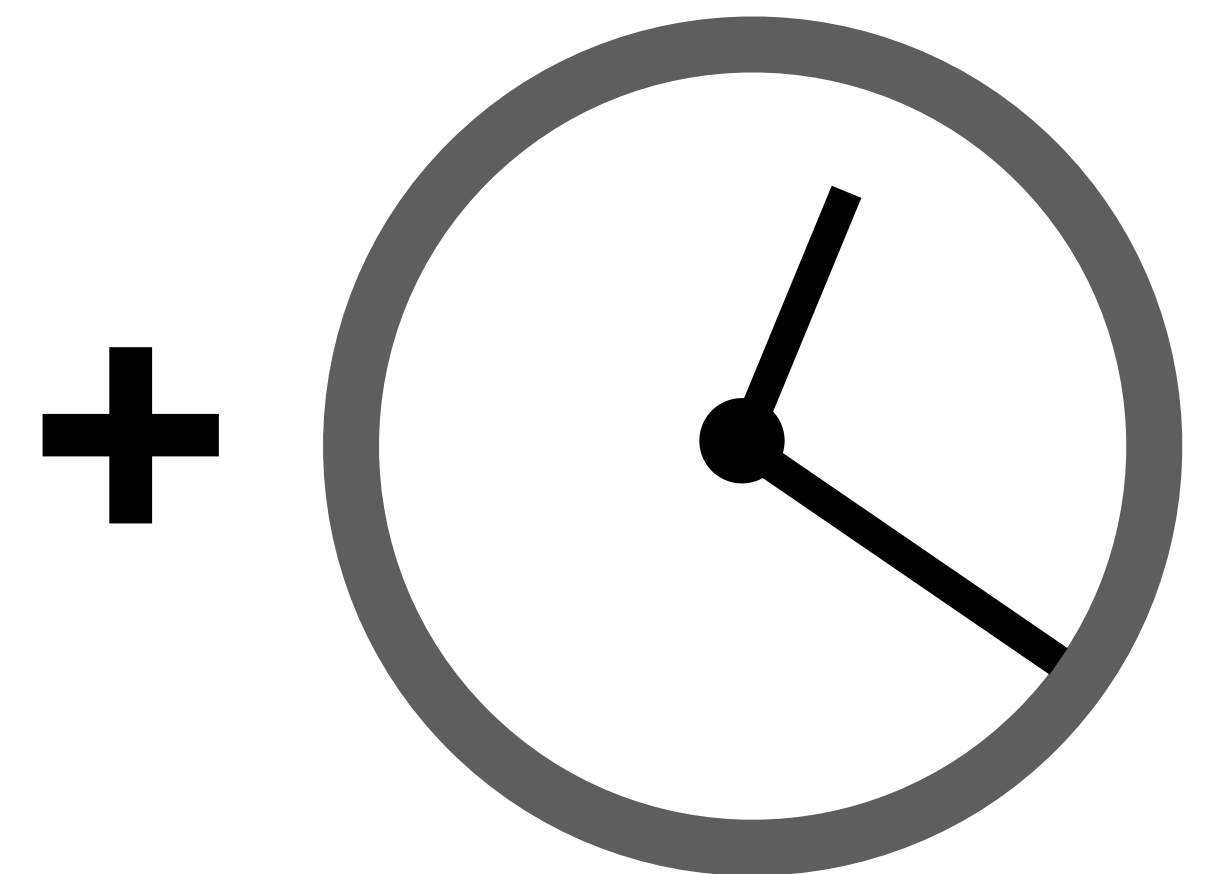
Task Domain Knowledge



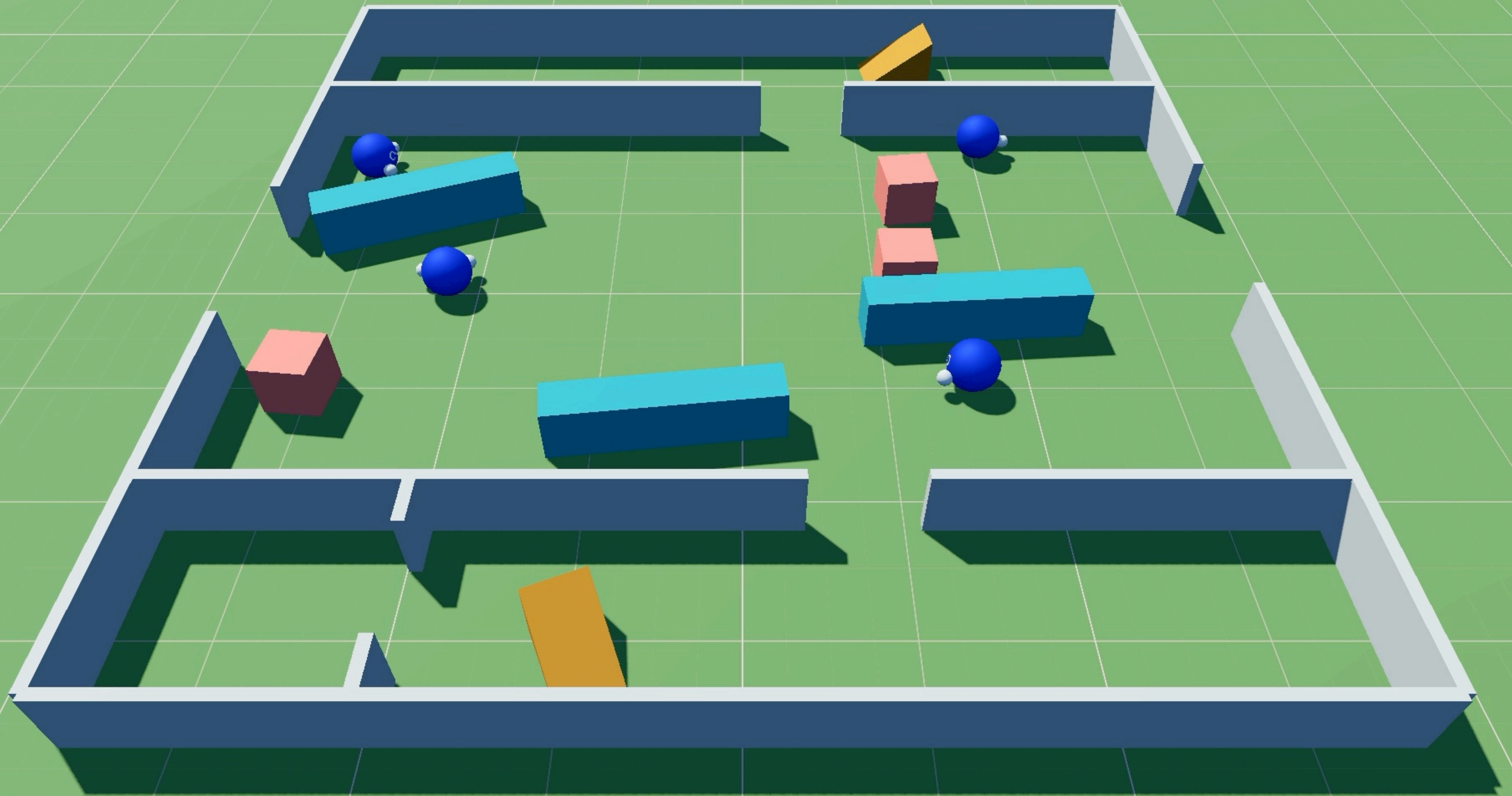
GPU Programming Skill

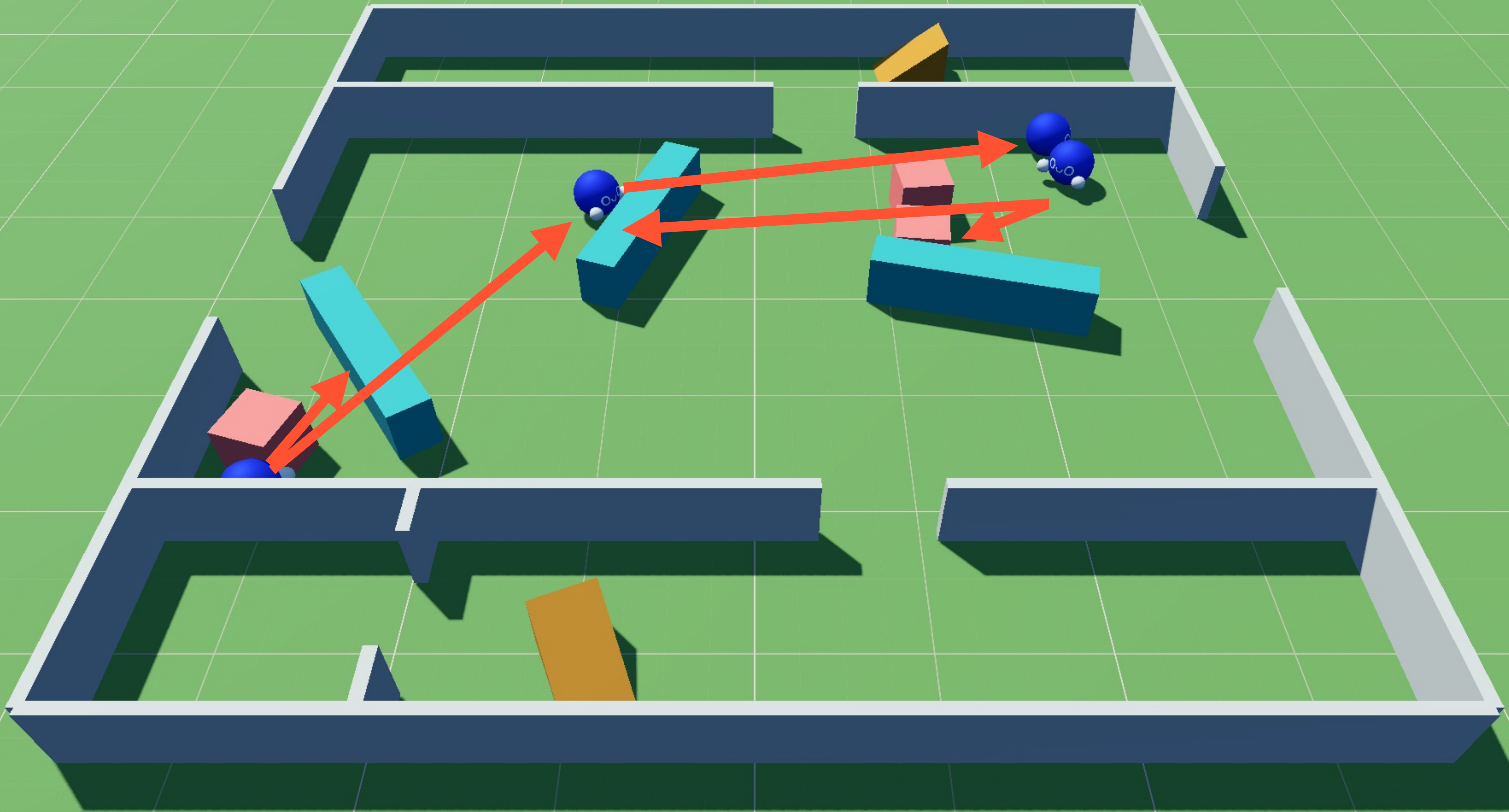


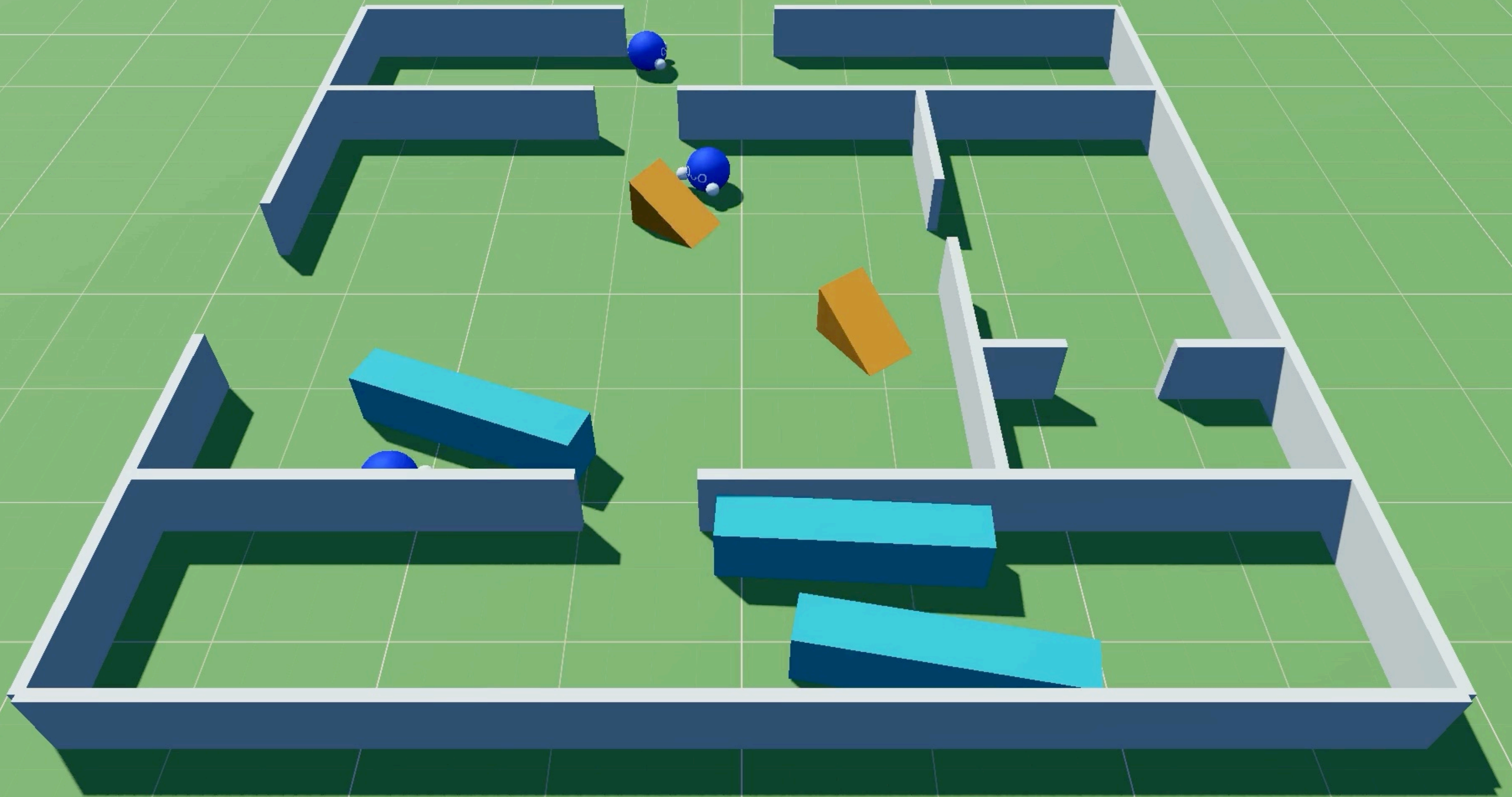
Engineering Time



Why isn't there a "game engine" to make this task easier?

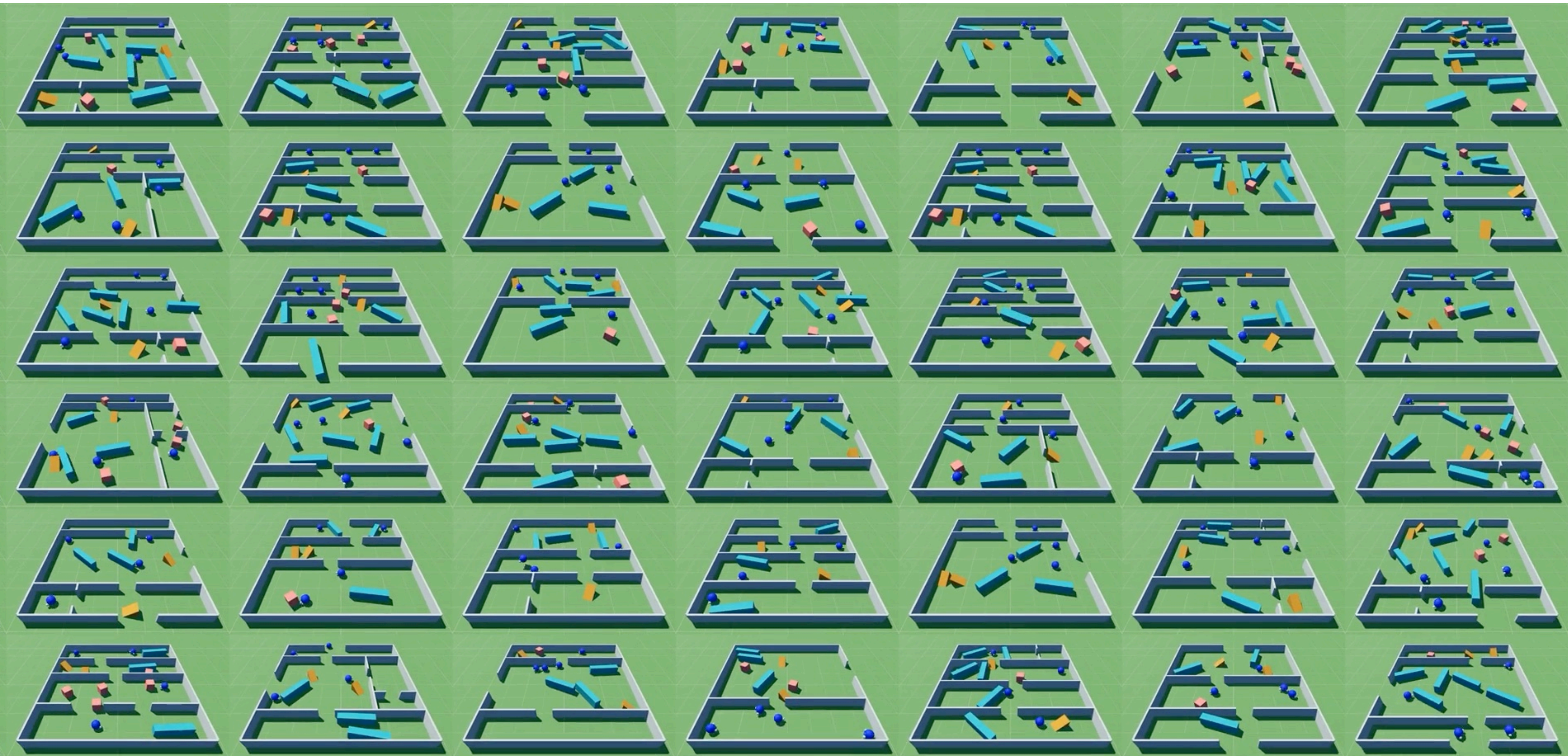




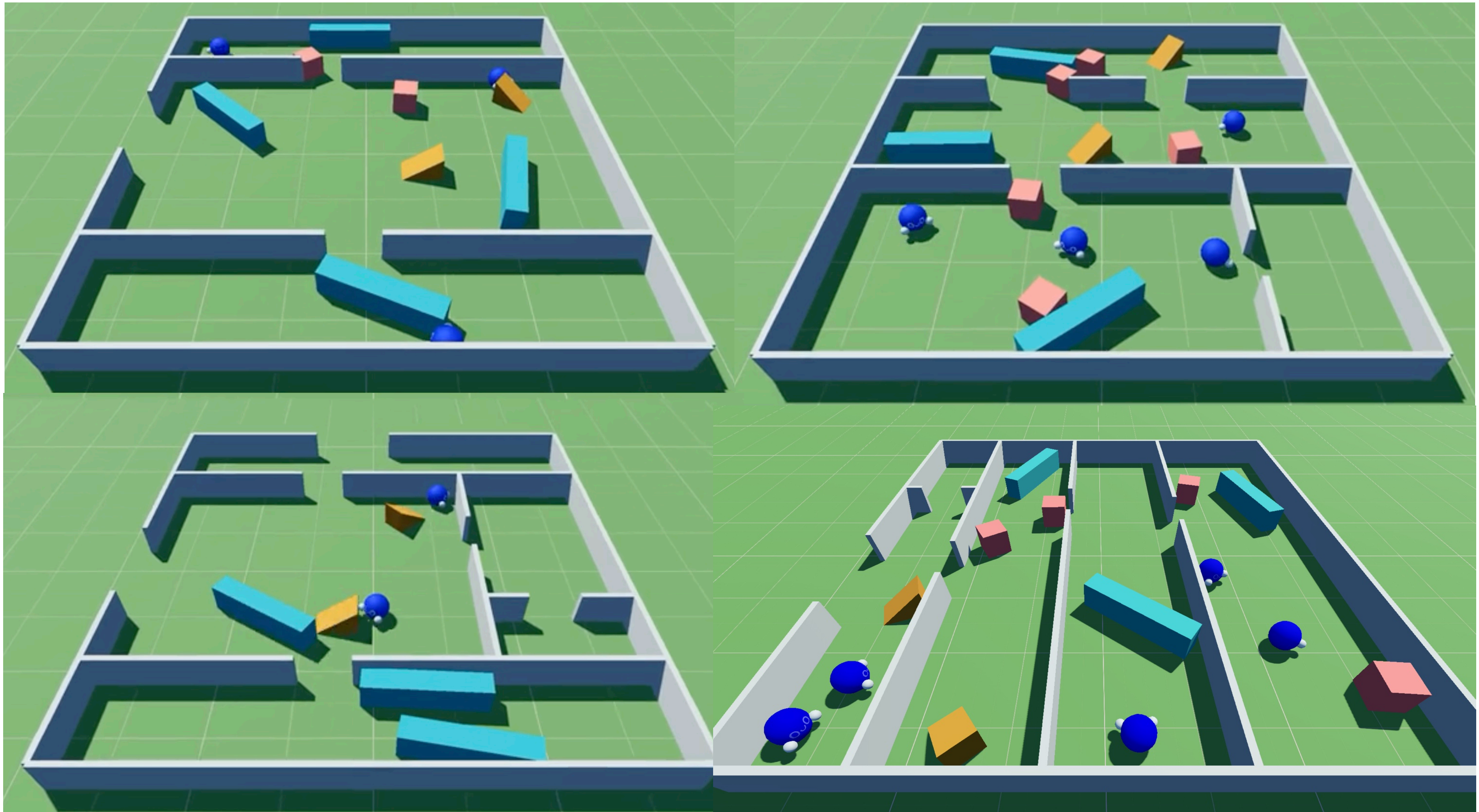


What are the parallel systems programming requirements?

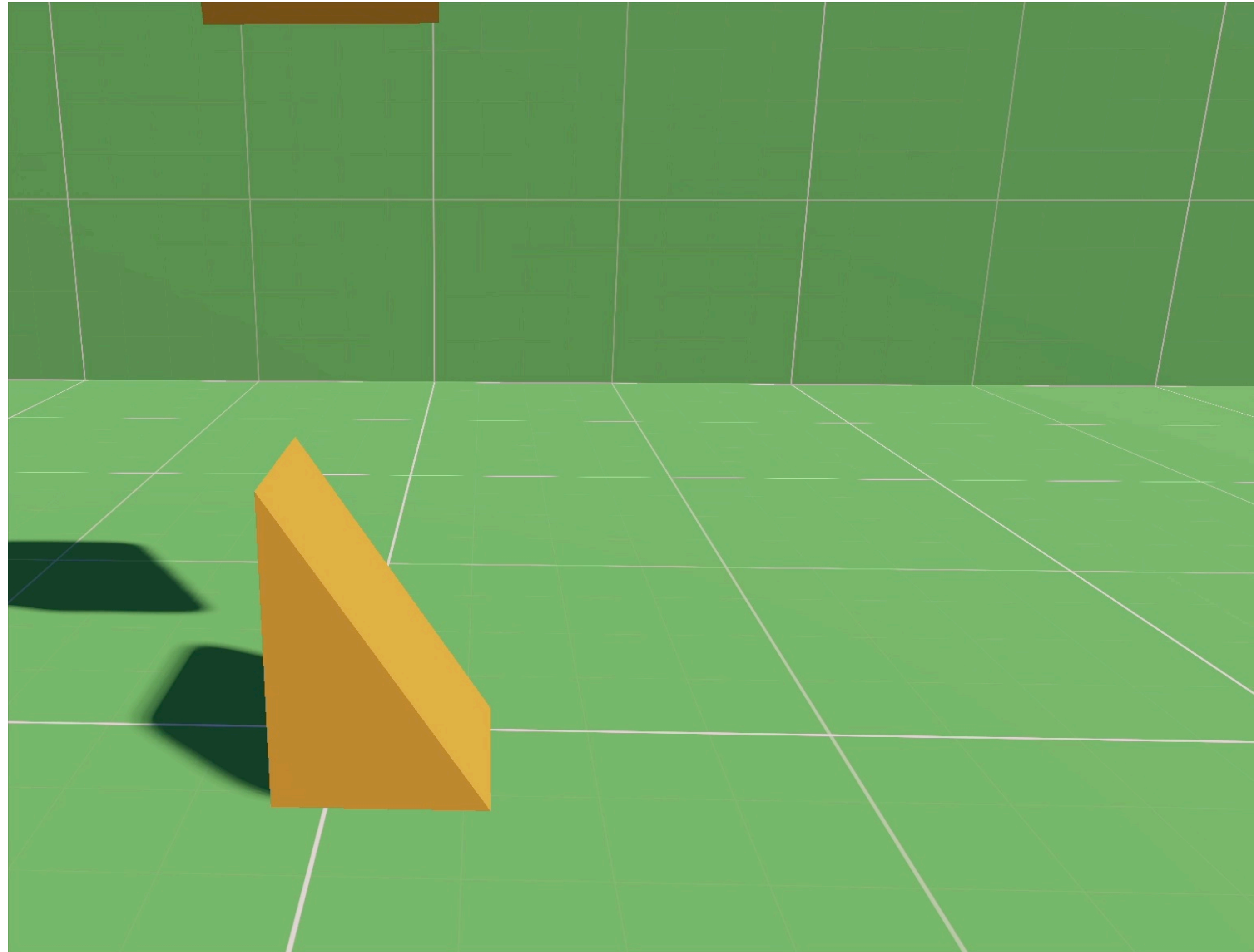
1. Nested Parallelism: Task logic for each object in each world



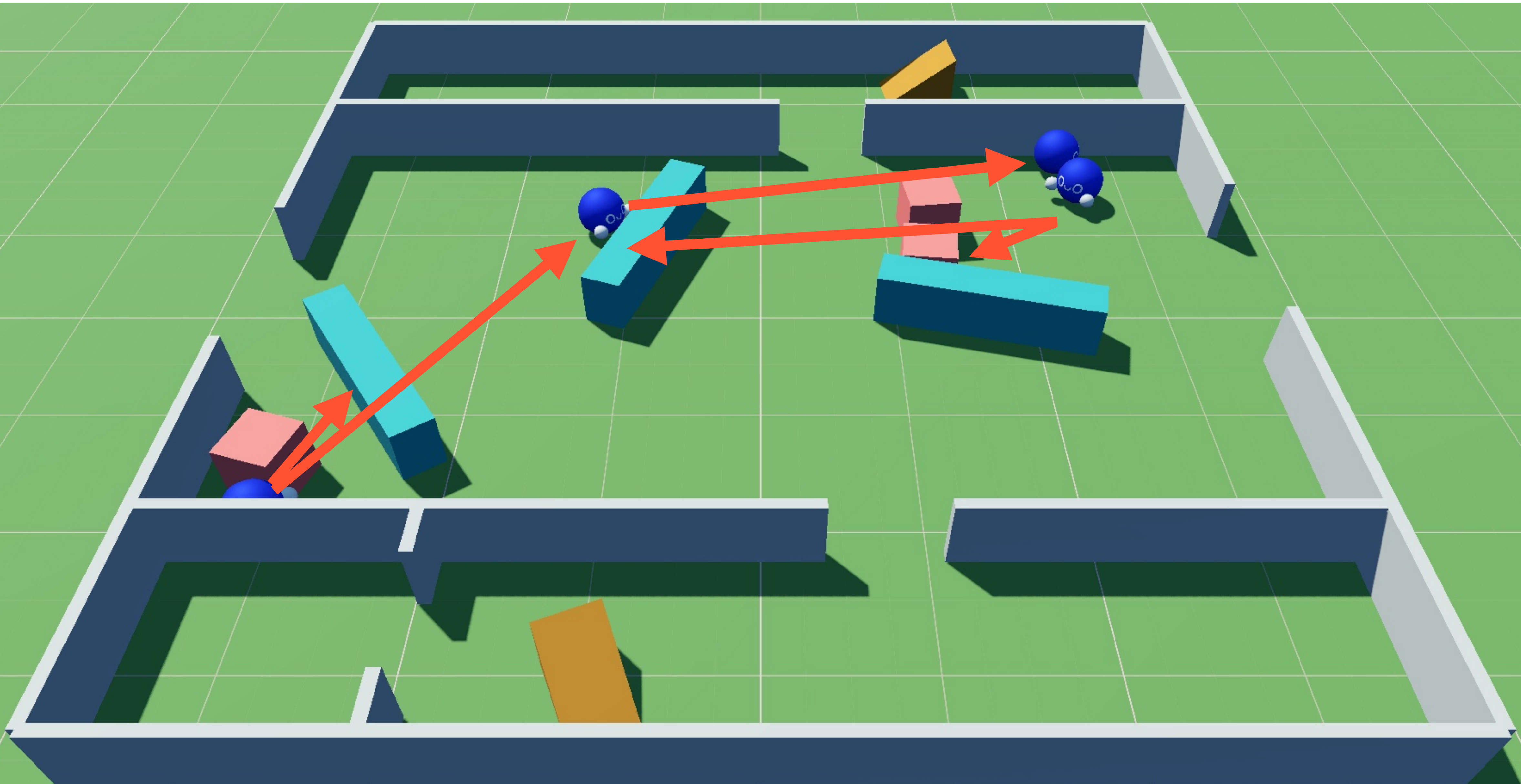
2. Irregular nested parallelism and irregular per-world state: worlds with varying numbers of objects



3. Dynamic allocation/deallocation of environment state **New entities, physics contacts come and go, sparse events**



4. Must perform complex spatial queries (aka “joins”): Ray casting, 3D collisions, physical proximity...



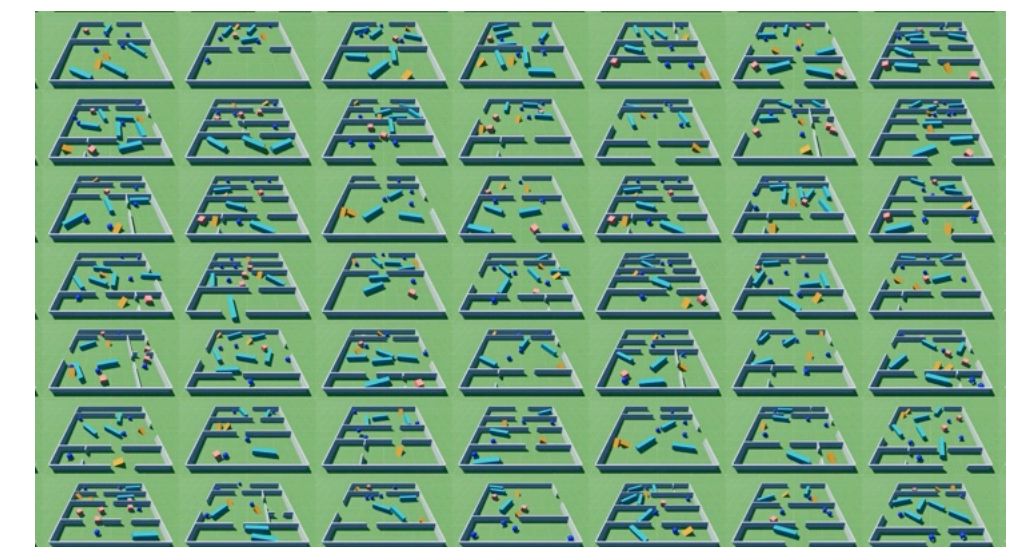
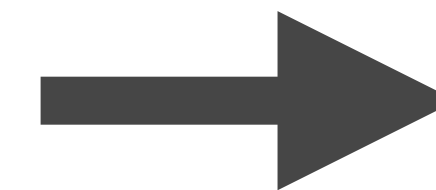
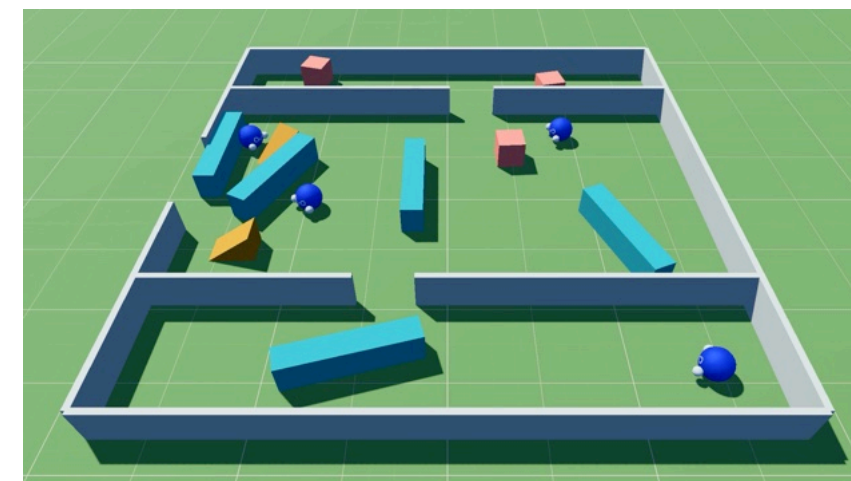
Usability: easy scripting of task-specific logic

```
def process_action(agent_position, action):  
    if action.type == MOVE_CHARACTER:  
        force = computeMovementForce(action.dir)  
    if action.type == LOCK_OBJECT_IN_PLACE:  
        hit_obj = raycastForward(agent_position)  
        if hit_obj:  
            lockObject(hit_obj)  
    ...
```

Usability: easy scripting of task-specific logic

```
def process_action(agent_position, action):  
    if action.type == MOVE_CHARACTER:  
        force = computeMovementForce(action.dir)  
    if action.type == LOCK_OBJECT_IN_PLACE:  
        hit_obj = raycastForward(agent_position)  
        if hit_obj:  
            lockObject(hit_obj)  
    ...
```

5. SPMD Implicit parallelism: logic written in terms of one environment, automatically batched



We need to create a "game engine" (that meets these requirements) for building batch simulators.

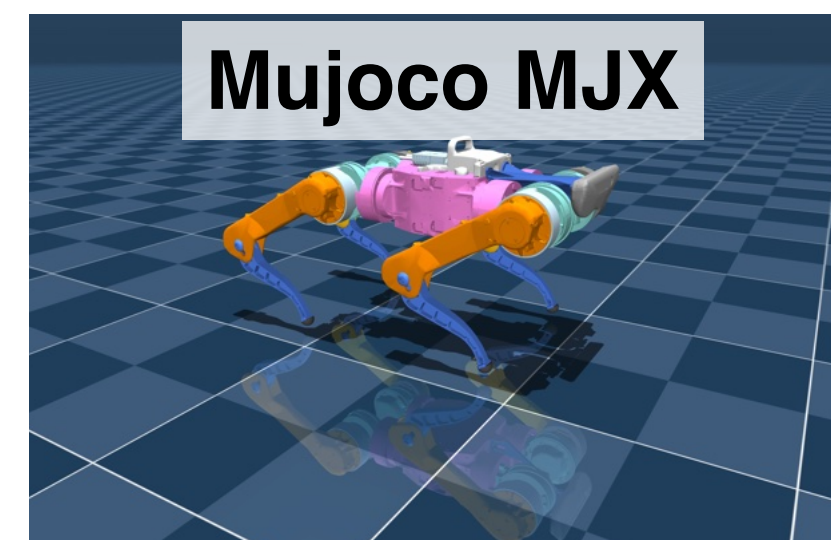
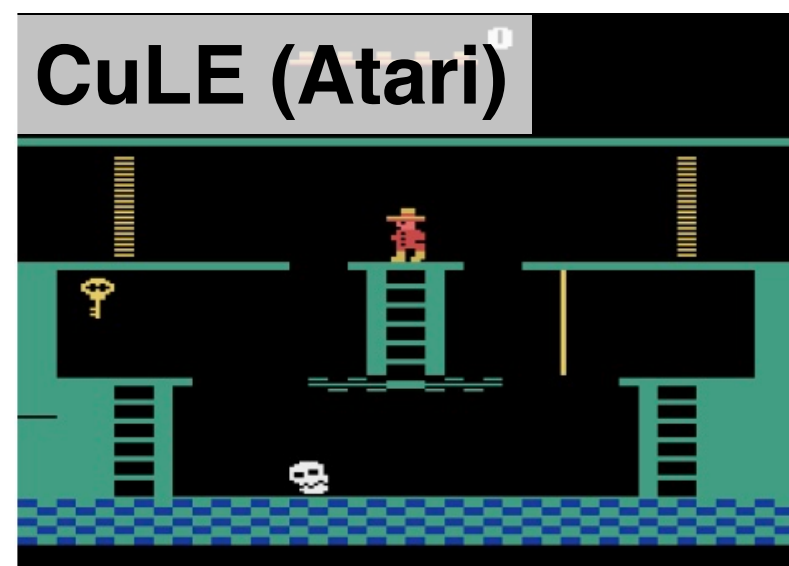
1. Nested Parallelism
2. Irregular Parallelism & Collection Sizes
3. Dynamic GPU-Controlled Allocation
4. Complex Spatial Joins
5. SPMD-Style Control Flow, Implicit Parallelism

What about existing systems / frameworks that have been used to build batch simulators?

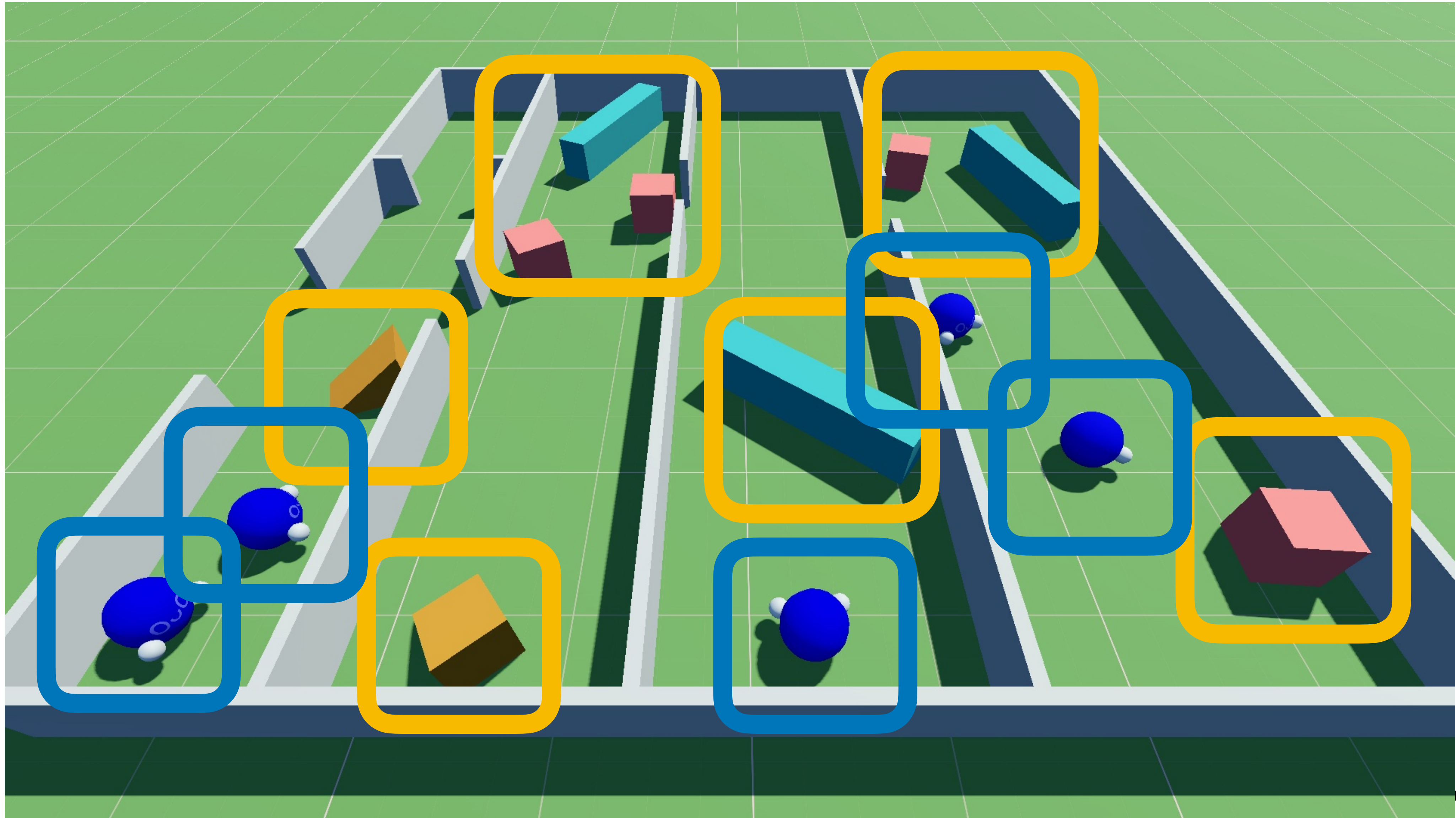
CUDA C++

NVIDIA Warp

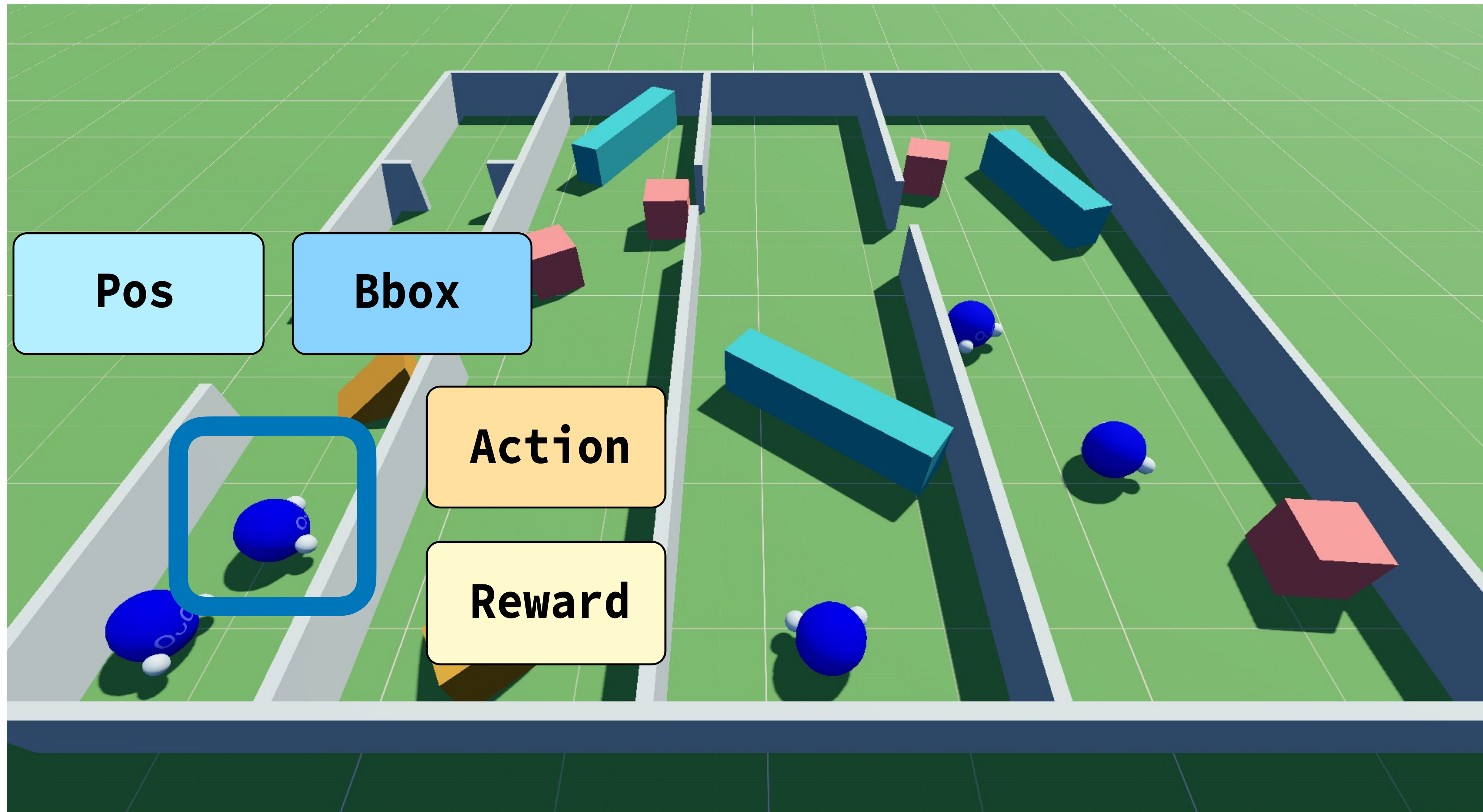

 PyTorch

Concept 1: ECS Entities



Concept 2: ECS Components



Components represent simulation state

Agents						
Id	EnvID	Pos	Bbox	Action	Reward	
12	0	[0,0,.5]	{min...	LEFT	0.1	...
32	0	[2,1,0]	{min...	FWD	-0.1	...
51	0	[1.5,0,1]	{min...	FWD	-2.5	...

Obstacles				
Id	EnvID	Pos	Bbox	
13	0	[0.5,0,.5]	{min...	...
72	0	[0,1,3]	{min...	...
61	0	[1,1,2]	{min...	...

Store data across **all environments in a batch** in columnar tables in GPU memory

Agents						
Id	EnvID	Pos	Bbox	Action	Reward	
12	0	[0,0,.5]	{min...	LEFT	0.1	...
32	0	[2,1,0]	{min...	FWD	-0.1	...
51	0	[1.5,0,1]	{min...	FWD	-2.5	...
62	1	[1.5,0.5,1]	{min...	RIGHT	0.3	...
65	1	[-0.5,0,0]	{min...	RIGHT	0.5	...
20	2	[0.5,1,1]	{min...	LEFT	1.5	...
23	2	[-1.5,0,0]	{min...	LEFT	10.5	...
41	2	[0.5,1,0]	{min...	BACK	-10	...

Obstacles				
Id	EnvID	Pos	Bbox	
13	0	[0.5,0,.5]	{min...	...
72	0	[0,1,3]	{min...	...
61	0	[1,1,2]	{min...	...
49	1	[0.5,1,2]	{min...	...
70	1	[-0.5,1,0]	{min...	...
33	2	[1.5,0,2.5]	{min...	...
81	3	[2.5,0.5,2]	{min...	...
11	3	[1.5,0.5,2]	{min...	...

Unified table storage also enables throughput-oriented dynamic memory allocation

Agents						
Id	EnvID	Pos	Bbox	Action	Reward	
12	0	[0,0,.5]	{min...	LEFT	0.1	...
32	0	[2,1,0]	{min...	FWD	-0.1	...
51	X	[1.5,0,1]	{min...	FWD	2.5	...
62	1	[1.5,0.5,1]	{min...	RIGHT	0.3	...
65	1	[-0.5,0,0]	{min...	RIGHT	0.5	...
20	2	[0.5,1,1]	{min...	LEFT	1.5	...
23	2	[-1.5,0,0]	{min...	LEFT	10.5	...
41	2	[0.5,1,0]	{min...	BACK	-10	...
51	1	[0.5,1,1]	{min...	BACK	-10	...

Obstacles				
Id	EnvID	Pos	Bbox	
13	0	[0.5,0,.5]	{min...	...
72	0	[0,1,3]	{min...	...
61	0	[1,1,2]	{min...	...
49	X	[0.5,1,2]	{min...	...
70	1	[-0.5,1,0]	{min...	...
33	2	[1.5,0,2.5]	{min...	...
81	3	[2.5,0.5,2]	{min...	...
11	X	[1.5,0.5,2]	{min...	...
15	2	[1.5,1.5,2]	{min...	...
12	0	[2.5,0.5,3]	{min...	...

Concept 3: Define functions to operate on these tables

Agents						
Id	EnvID	Pos	Bbox	Action	Reward	
12	0	[0,0,.5]	{min...	LEFT	0.1	...
32	1	[2,1,0]	{min...	FWD	-0.1	...
51	2	[1.5,0,1]	{min...	FWD	2.5	...
22	2	[2.5,0,1.5]	{min...	RIGHT	1.1	...

Obstacles				
Id	EnvID	Pos	Bbox	
13	0	[0.5,0,.5]	{min...	...
72	1	[0,1,3]	{min...	...
61	1	[1,1,2]	{min...	...
25	2	[1.5,1,2.5]	{min...	...

ProcessActions

Pos, Action

Concept 3: Functions define logic on these tables

Agents						
Id	EnvID	Pos	Bbox	Action	Reward	
12	0	[0,0,.5]	{min...	LEFT	0.1	...
32	1	[2,1,0]	{min...	FWD	-0.1	...
51	2	[1.5,0,1]	{min...	FWD	2.5	...
22	2	[2.5,0,1.5]	{min...	RIGHT	1.1	...

Obstacles				
Id	EnvID	Pos	Bbox	
13	0	[0.5,0,.5]	{min...	...
72	1	[0,1,3]	{min...	...
61	1	[1,1,2]	{min...	...
25	2	[1.5,1,2.5]	{min...	...

ProcessActions
Pos, Action



Written as straight-line, per-entity logic, including conditional behavior

Agents					
Id	EnvID	Pos	Bbox	Action	Reward
12	0	[0,0,.5]	{min...	LEFT	0.1
32	1	[2,1,0]	{min...	FWD	-0.1
51	2	[1.5,0,1]	{min...	FWD	2.5
22	2	[2.5,0,1.5]	{min...	RIGHT	1.1

```
def process_action(agent_position, action):  
    if action.type == MOVE:  
        force = computeMovementForce(action.dir)  
  
    if action.type == LOCK:  
        hit_obj = raycastForward(agent_position)  
        if hit_obj:  
            lockObject(hit_obj)  
  
    ...
```

ProcessActions

Pos, Action

Easy implement of 'foreach': parallel GPU threads execute system logic over each table row

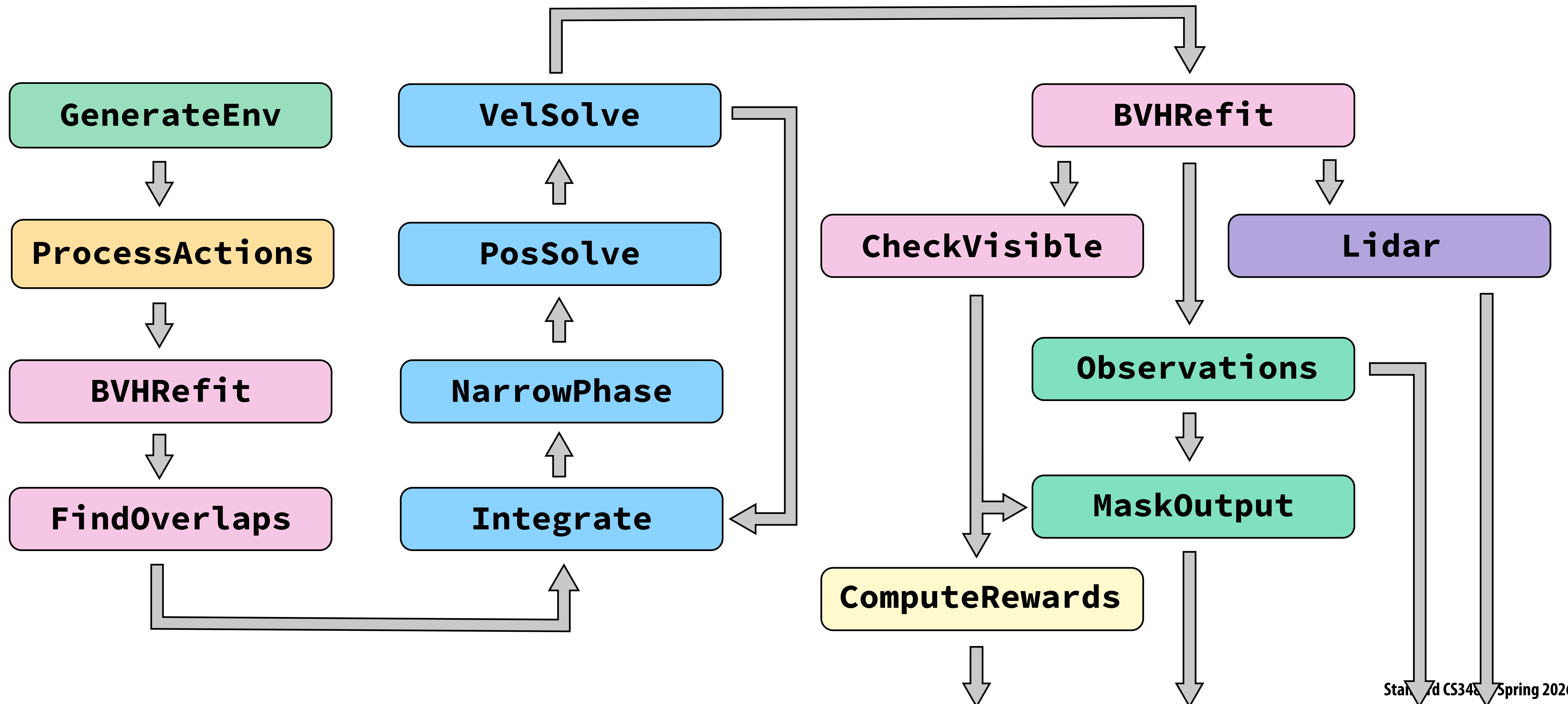
Agents					
Id	EnvID	Pos	Bbox	Action	Reward
GPU Thread	0	[0,0,.5]	{min...	LEFT	0.1
GPU Thread	1	[2,1,0]	{min...	FWD	-0.1
GPU Thread	2	[1.5,0,1]	{min...	FWD	2.5
GPU Thread	2	[2.5,0,1.5]	{min...	RIGHT	1.1

```
def process_action(agent_position, action):  
    if action.type == MOVE:  
        force = computeMovementForce(action.dir)  
  
    if action.type == LOCK:  
        hit_obj = raycastForward(agent_position)  
        if hit_obj:  
            lockObject(hit_obj)  
  
    ...
```

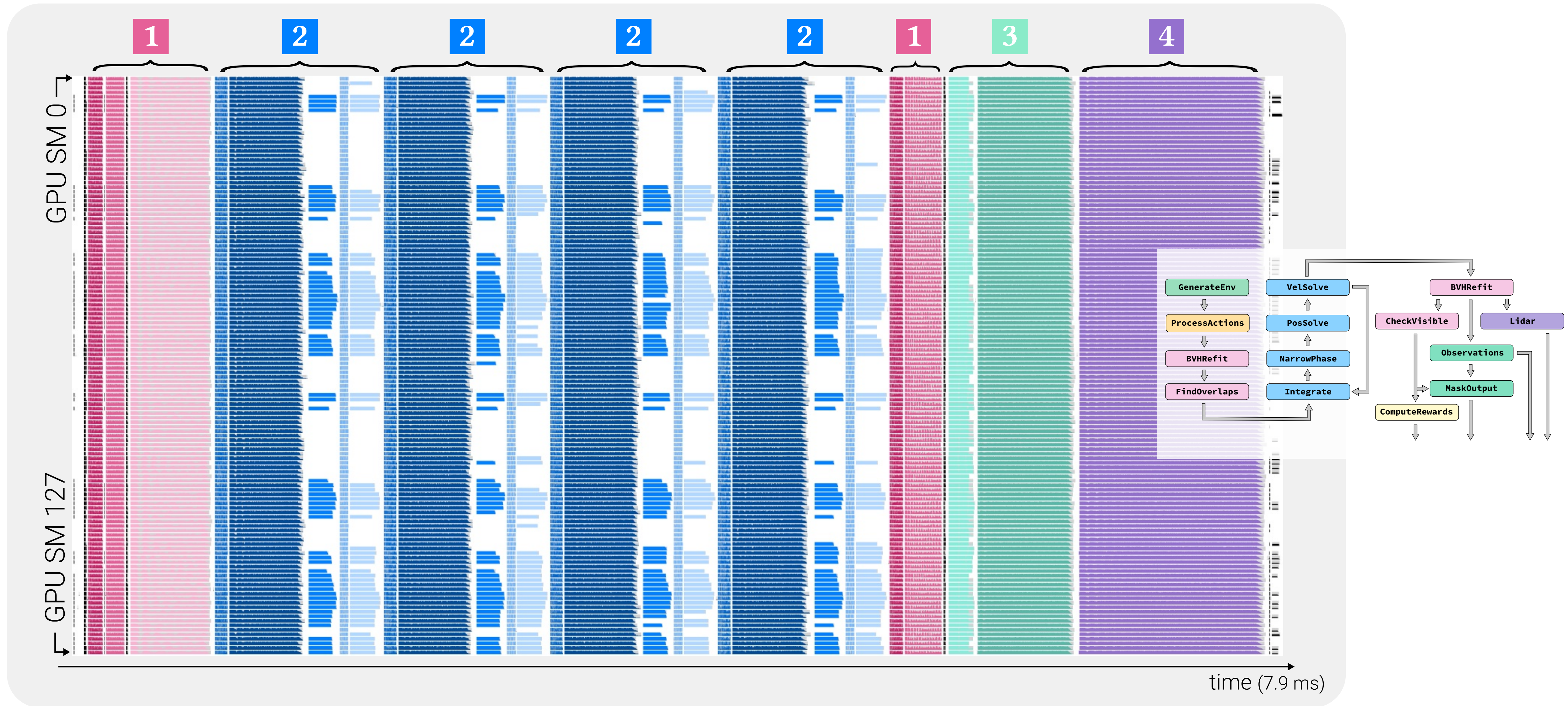
ProcessActions

Pos, Action

Create a task graph of what needs to be done per sim step...

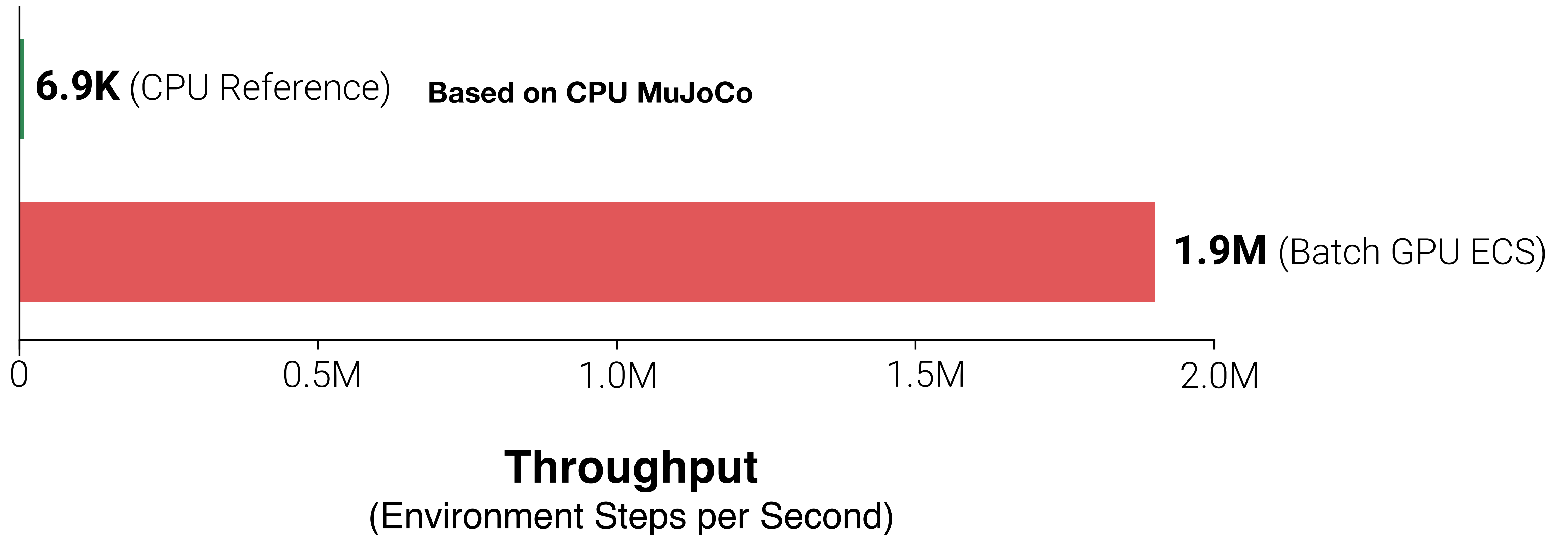


Batch of 8192 worlds running on the GPU

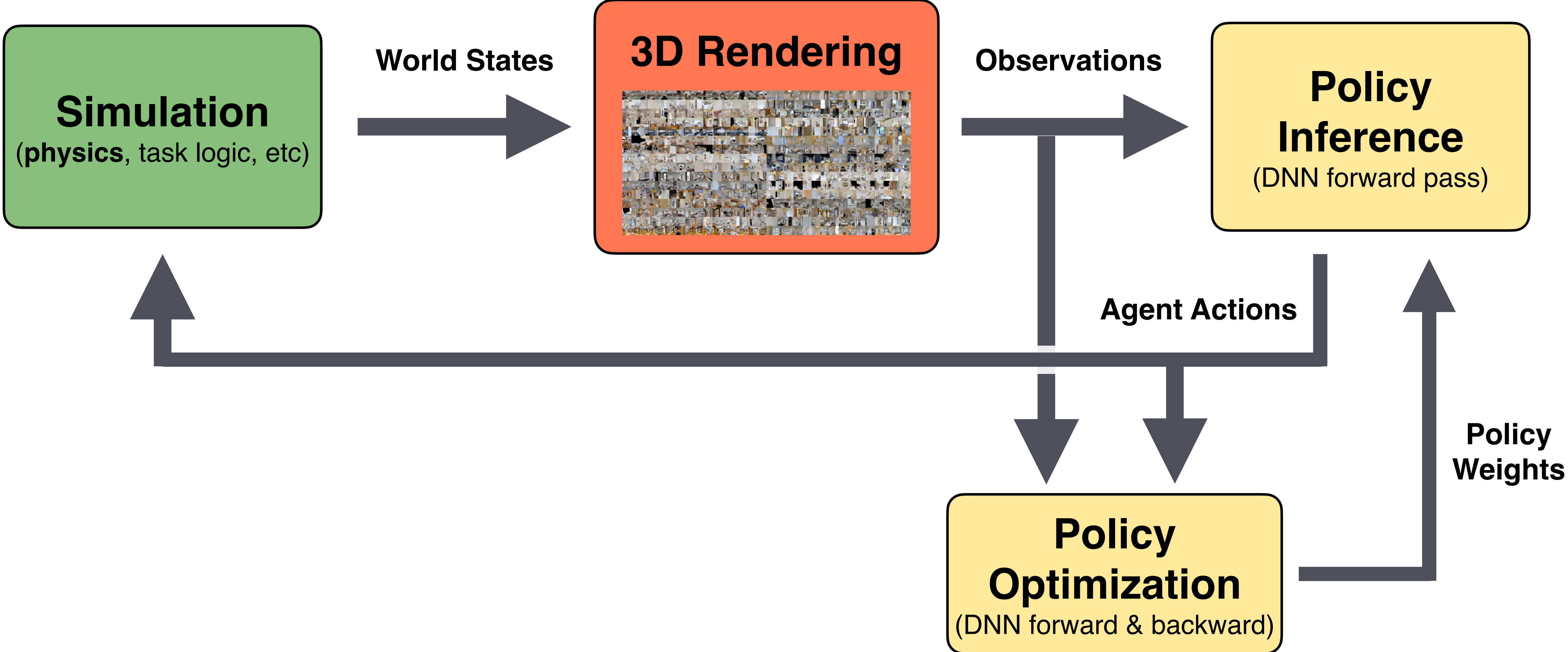


1 BVH & Broad Phase **2** Physics Sub-Step **3** Agent Observations & Rewards **4** LIDAR

Our Hide-and-Seek batch simulator is orders of magnitude faster than original open source reference



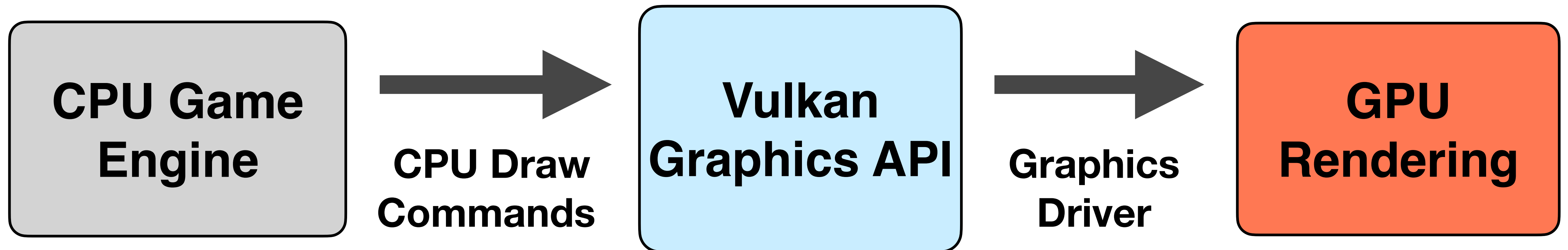
Key subsystems in an RL agent training loop: If the subsystems don't communicate efficiently, the whole loop will be slow



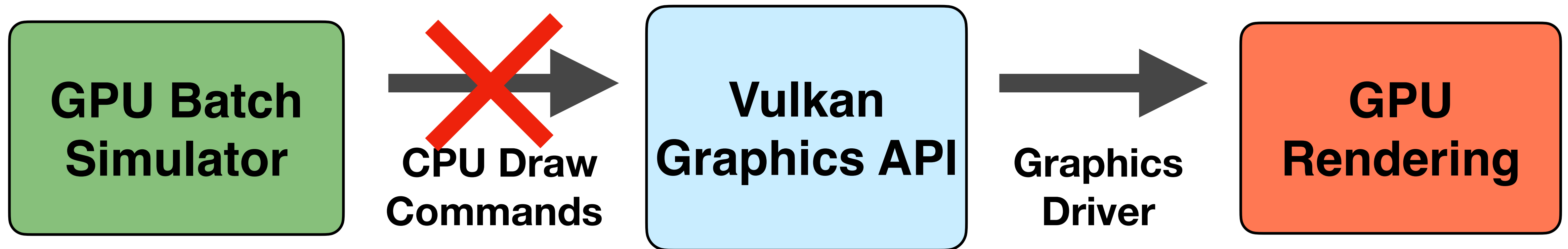


How to interface with a 1 million+ FPS game engine on the GPU?

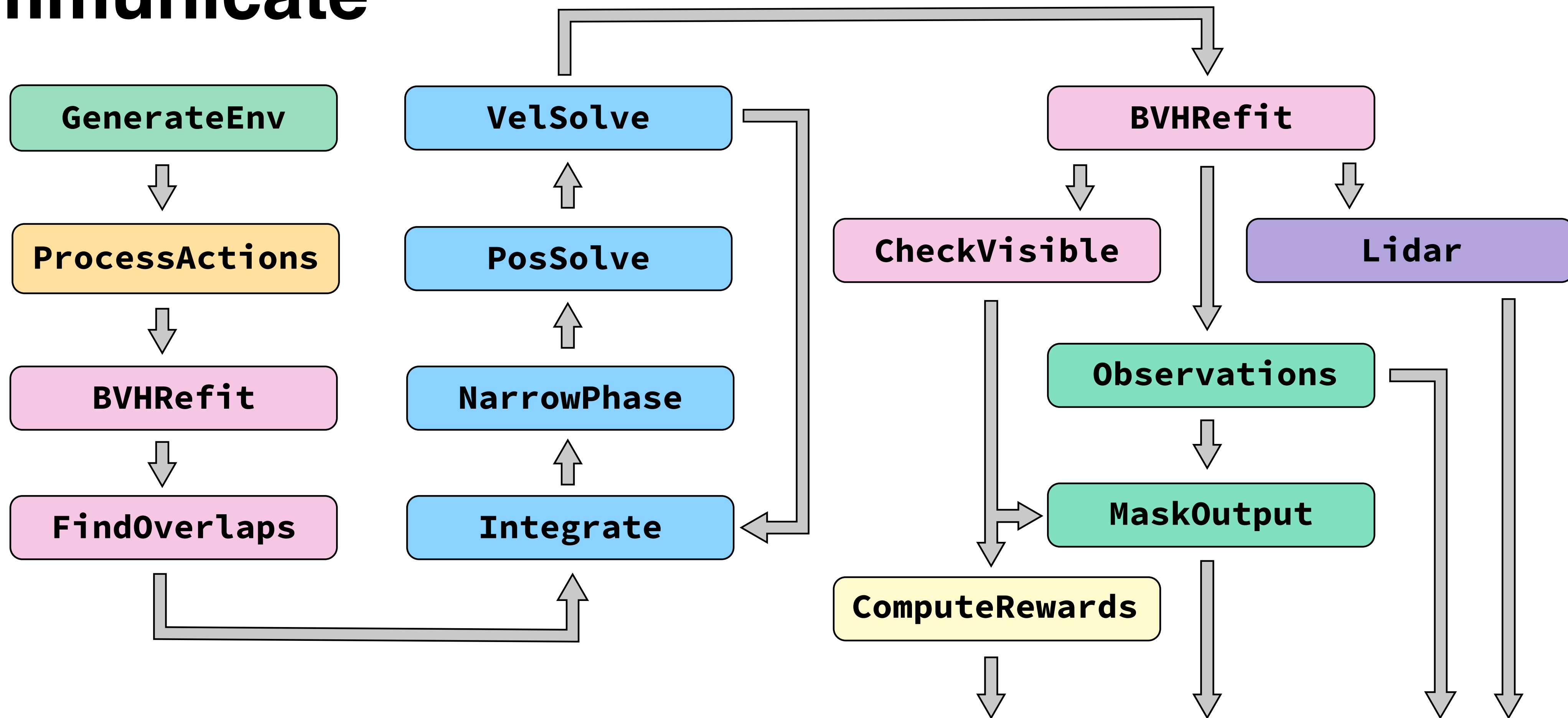
Problem: existing graphics hardware APIs and renderers only provide CPU-side interfaces
“They expect the CPU to tell the GPU what to do!”)



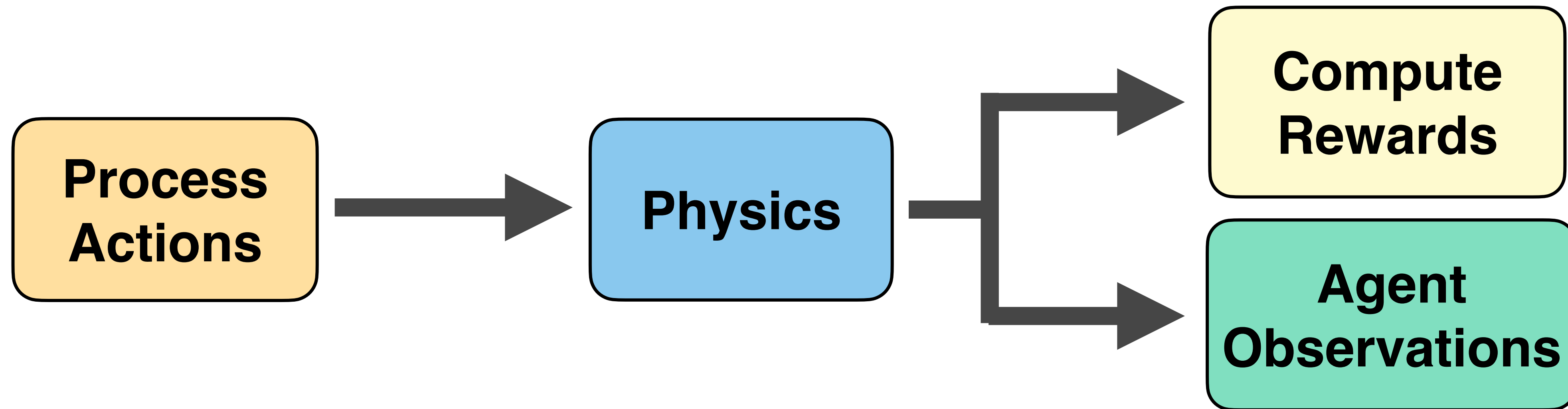
Problem: our engine runs entirely on the GPU



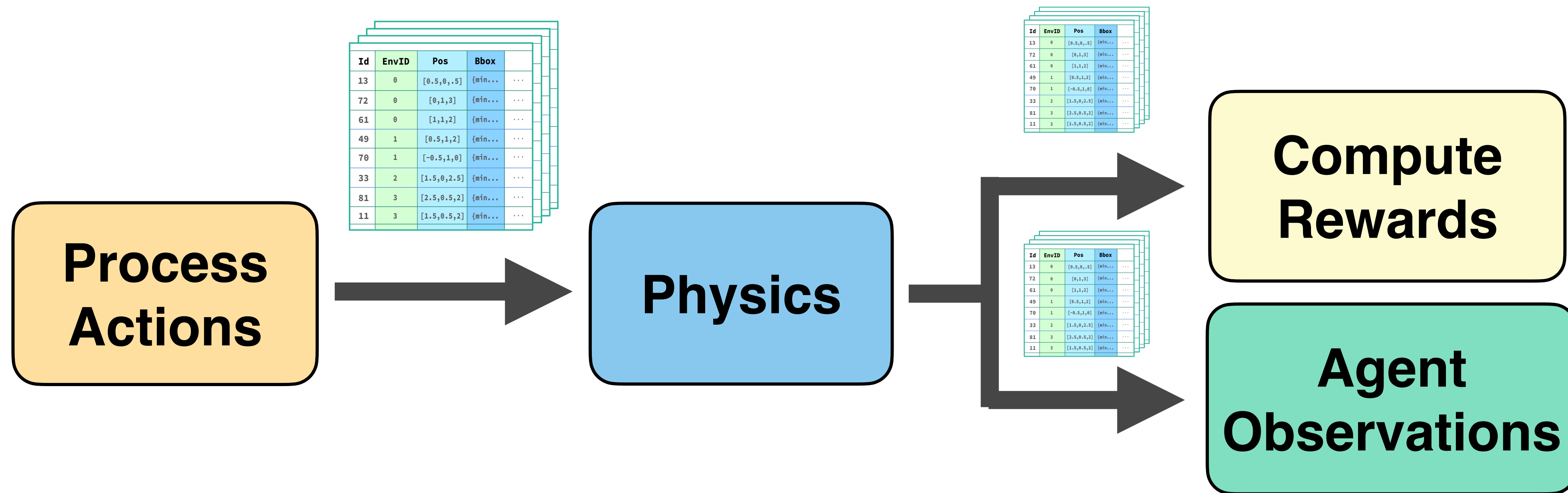
Recall: multi-world ECS tables provide high-throughput interface for ECS systems to communicate



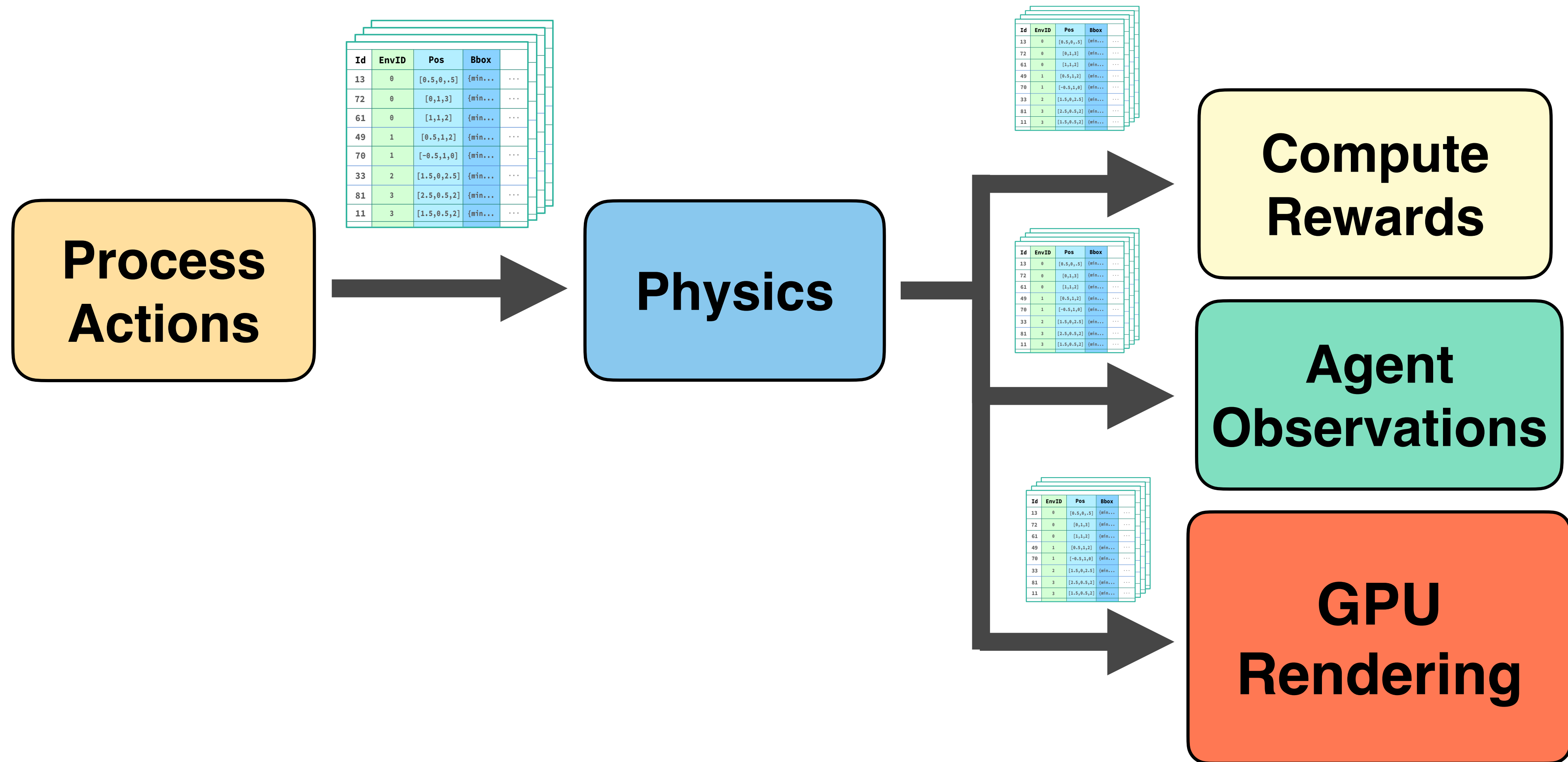
Multi-world ECS tables provide high-throughput interface for ECS systems to communicate



Multi-world ECS tables provide high-throughput interface for ECS systems to communicate

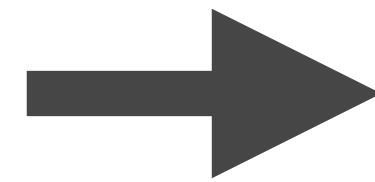


Idea: use multi-world ECS tables as renderer interface



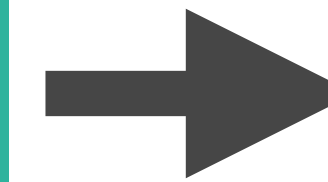
Batch simulator and batch renderer communicate through reads/writes to shared ECS table: no CPU intervention!

GPU Batch Simulator



Renderable Entities

ID	EnvID	ObjID	Transform
100	0	0	[[0.5, 0.0, ...]]
99	0	5	[[0.2, 0.1, ...]]
102	1	0	[[-0.5, 1.0, ...]]
105	2	2	[[1.5, 2.0, ...]]
50	2	3	[[0.2, -1.5, ...]]
200	3	2	[[1.0, 0.0, ...]]
330	3	5	[[1.5, 2.0, ...]]
⋮	⋮	⋮	⋮



GPU Rendering

Worlds can easily update positions & other properties of renderable objects in parallel

GPU Batch Simulator

Renderable Entities

GPU Rendering

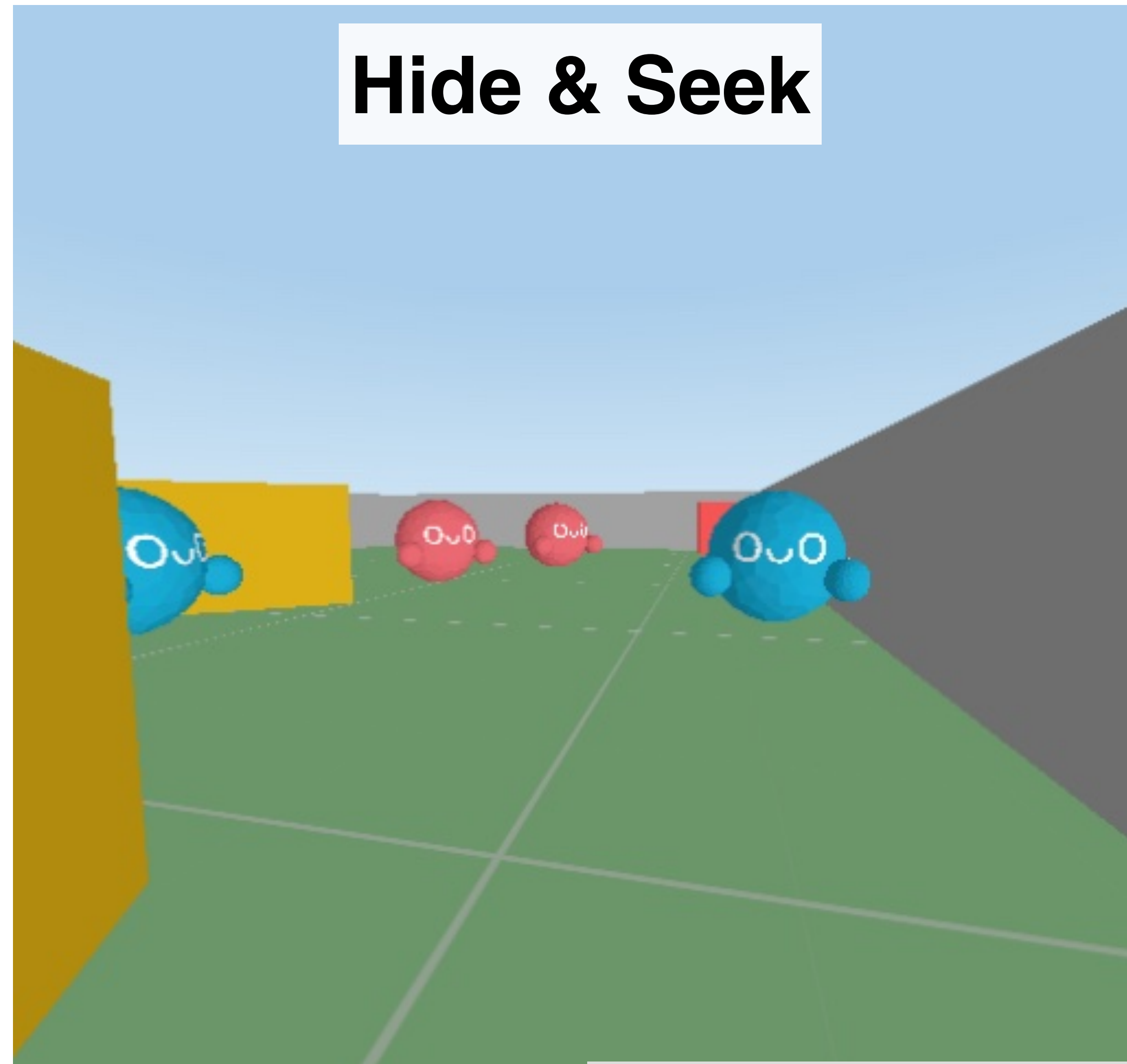
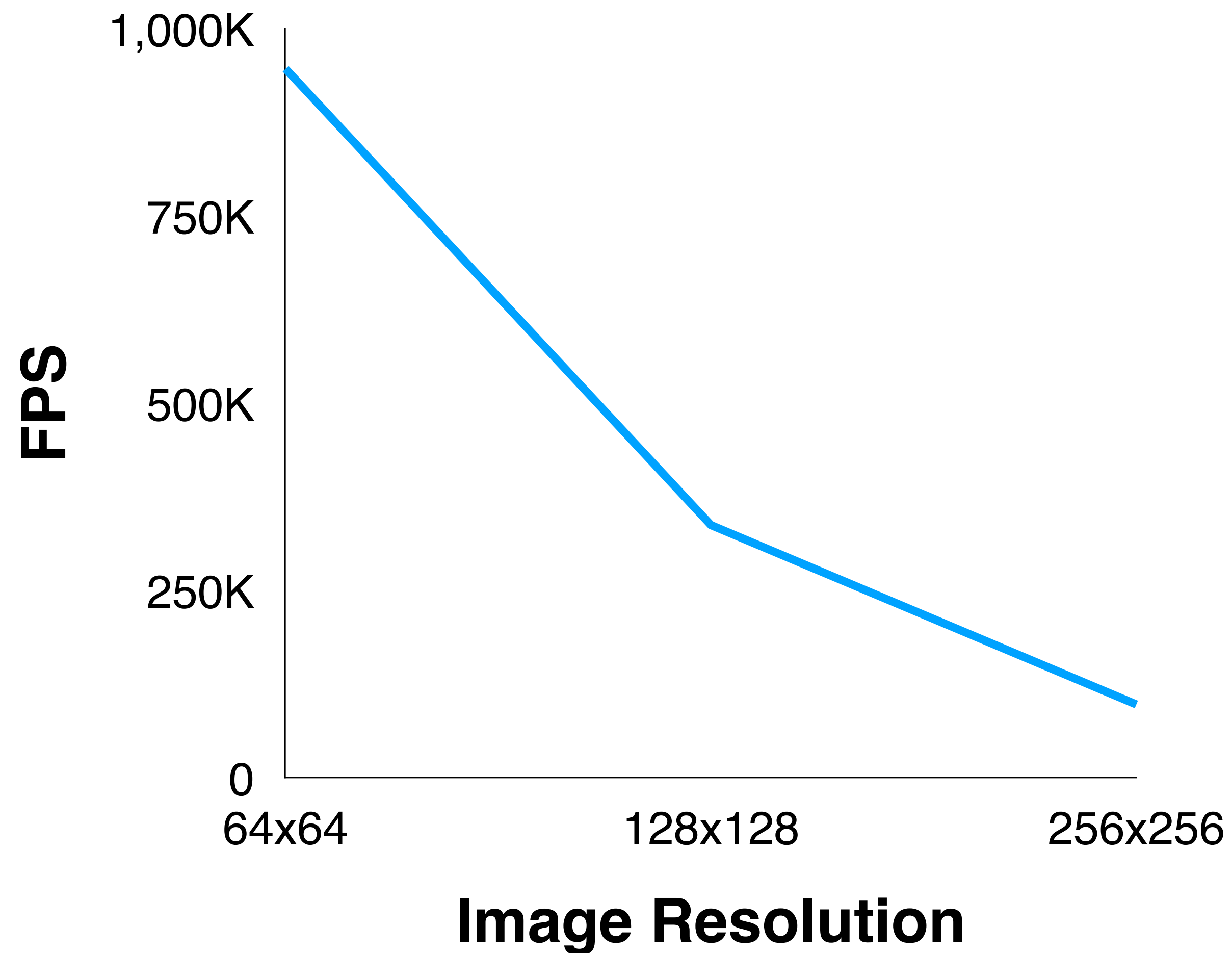
update_position(id)

update_position(id)

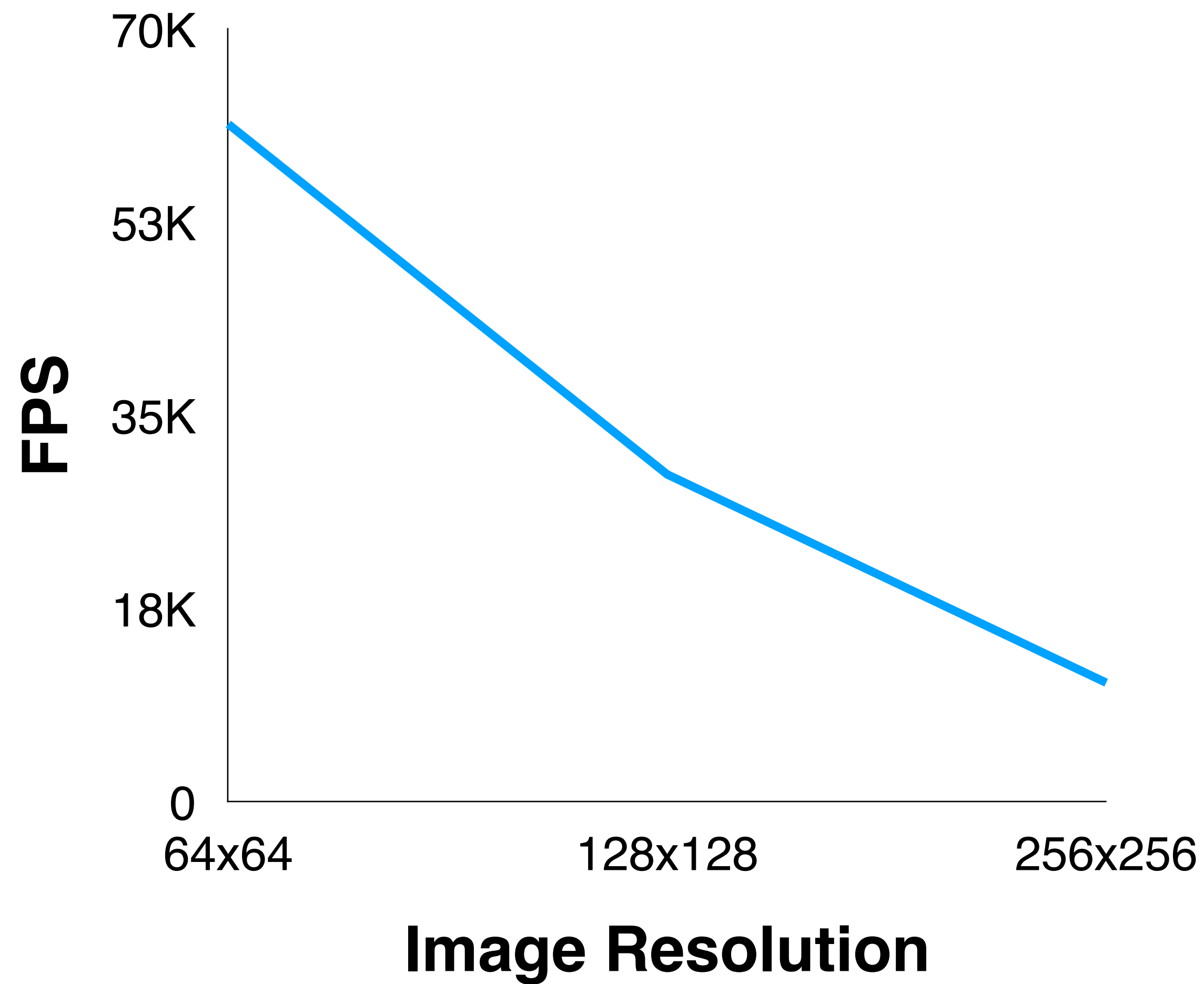
update_position(id)

ID	EnvID	ObjID	Transform
100	0	0	[[0.5, 0.0, ...]]
99	0	5	[[0.2, 0.1, ...]]
102	1	0	[[-0.5, 1.0, ...]]
105	2	2	[[1.5, 2.0, ...]]
50	2	3	[[0.2, -1.5, ...]]
200	3	2	[[1.0, 0.0, ...]]
330	3	5	[[1.5, 2.0, ...]]
⋮	⋮	⋮	⋮

>100K FPS on low complexity scenes, including cost of batch simulator => batch renderer interface



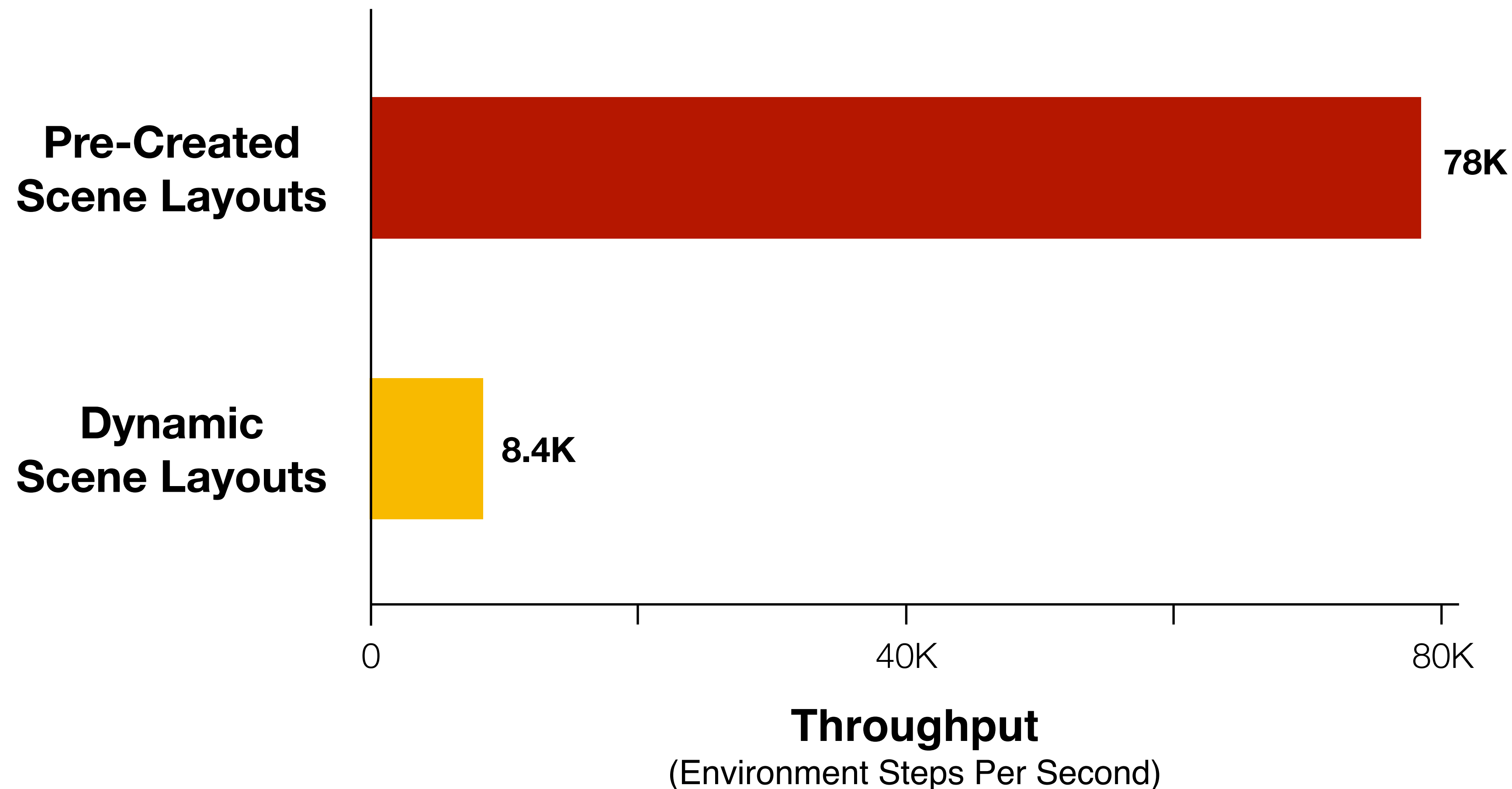
>10K FPS on high geometric complexity Robotics scenes



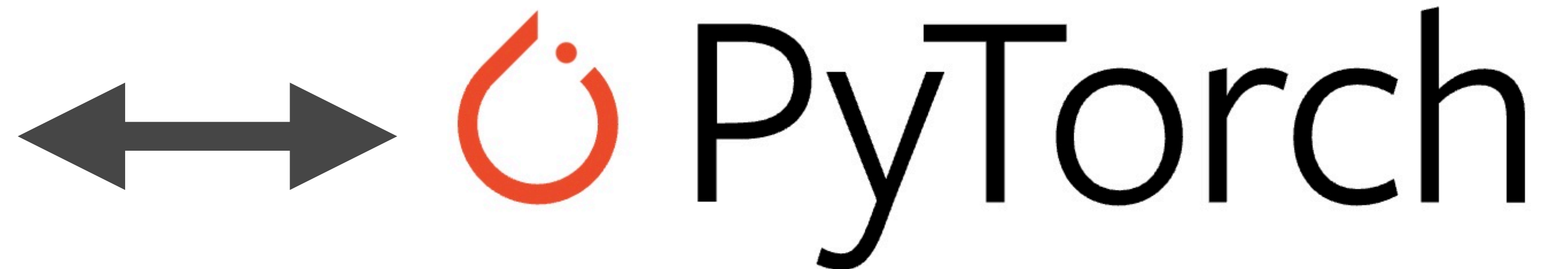
**Similar CPU \Rightarrow GPU Interface problems
exist in off-the-shelf GPU physics libraries!**

PhysX requires CPU to “setup the physics scene”: a bottleneck for procedural generation during training

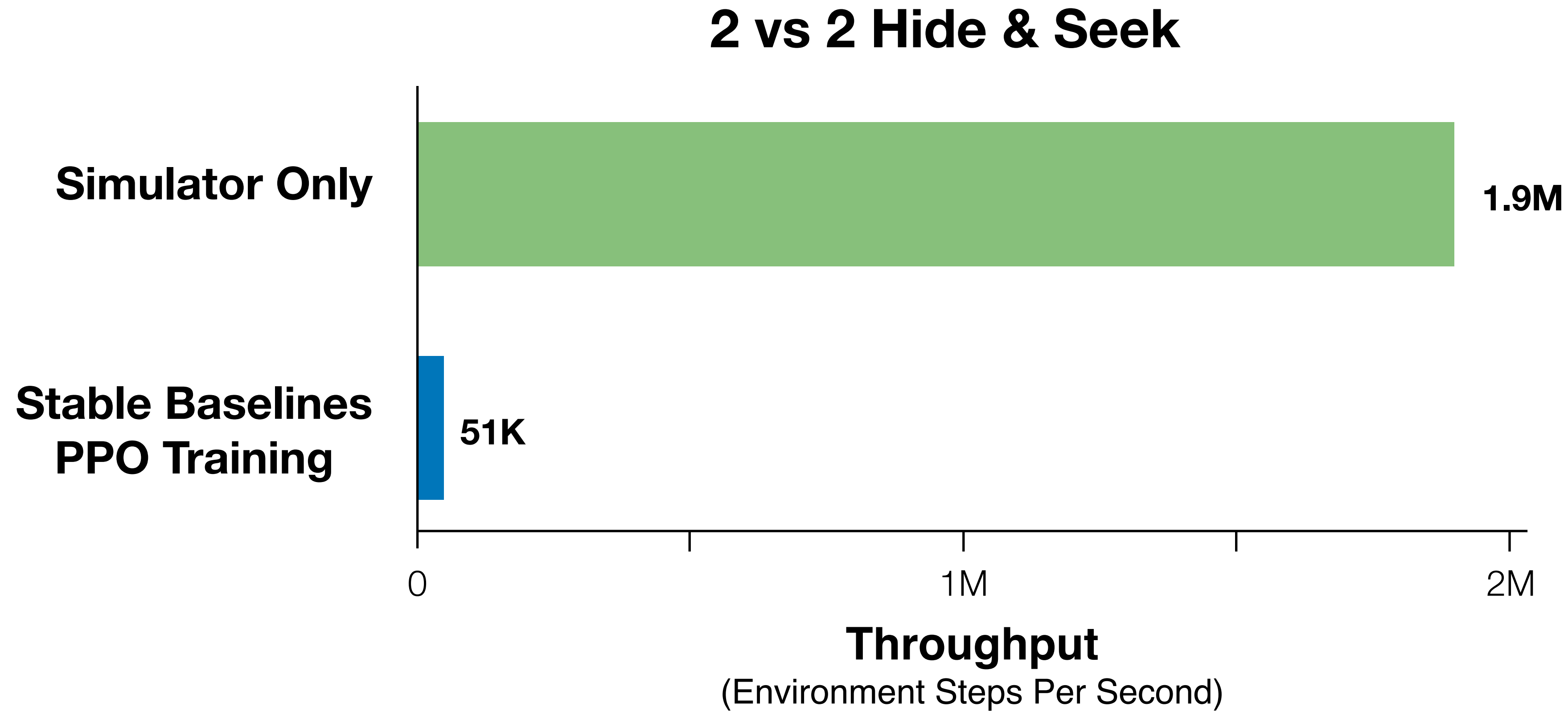
Maniskill Open Cabinet Drawer Task
200 Step Episodes



Interfacing batch simulators with deep learning frameworks

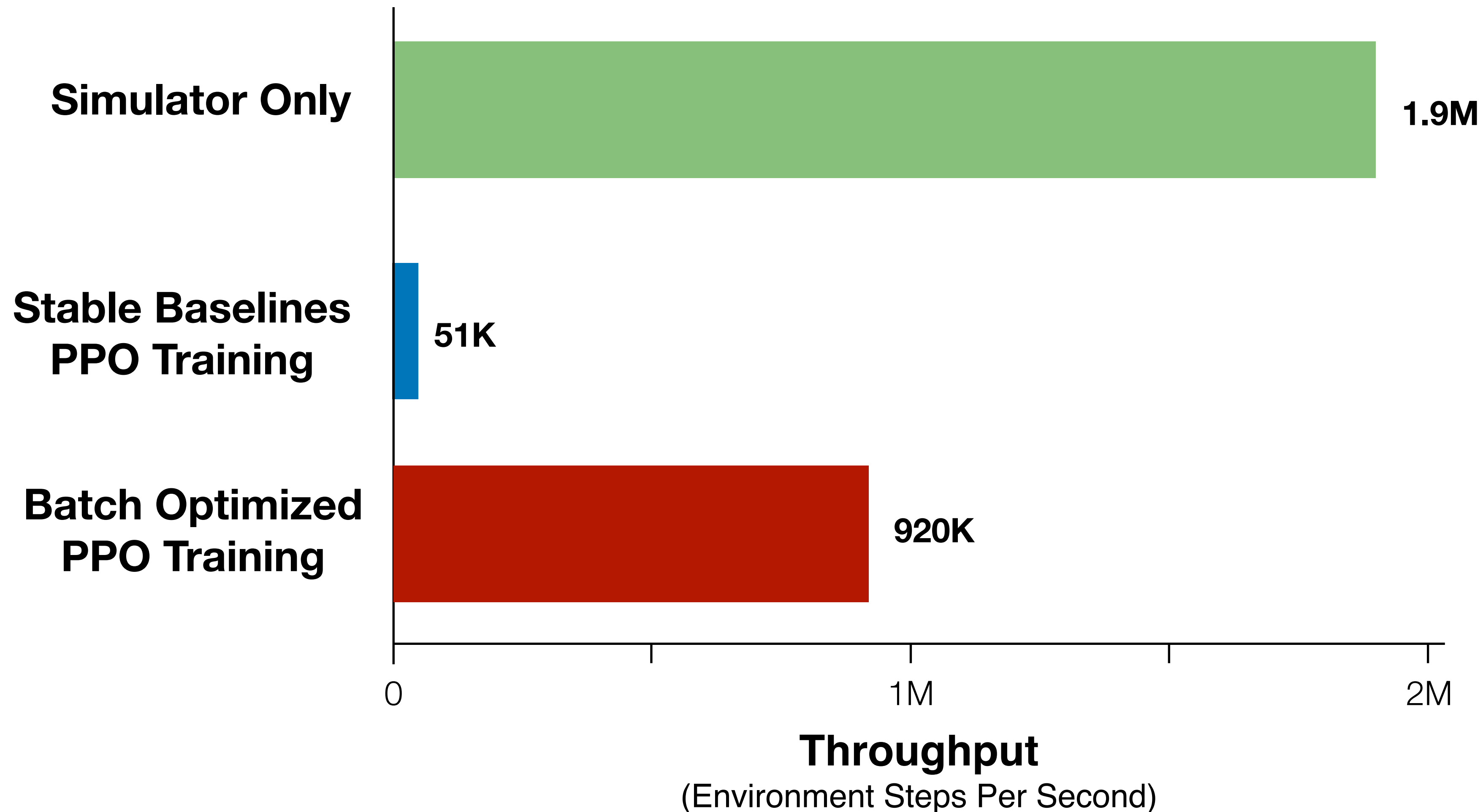


Many off-the-shelf RL libraries do not have efficient interfaces to batch simulators

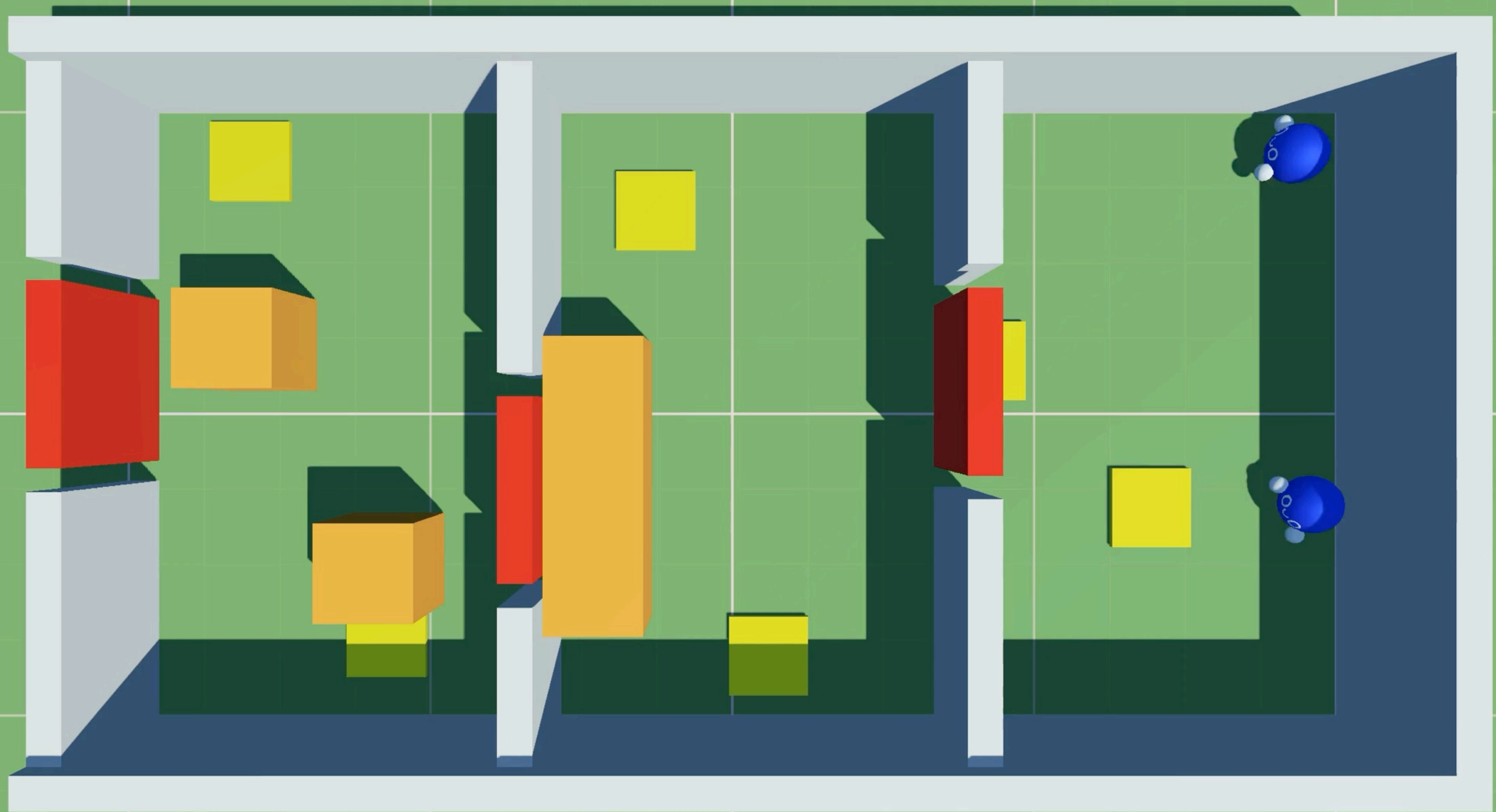


Building new training codebase around ECS-exported tensors achieves ~20x end-to-end speedup

2 vs 2 Hide & Seek



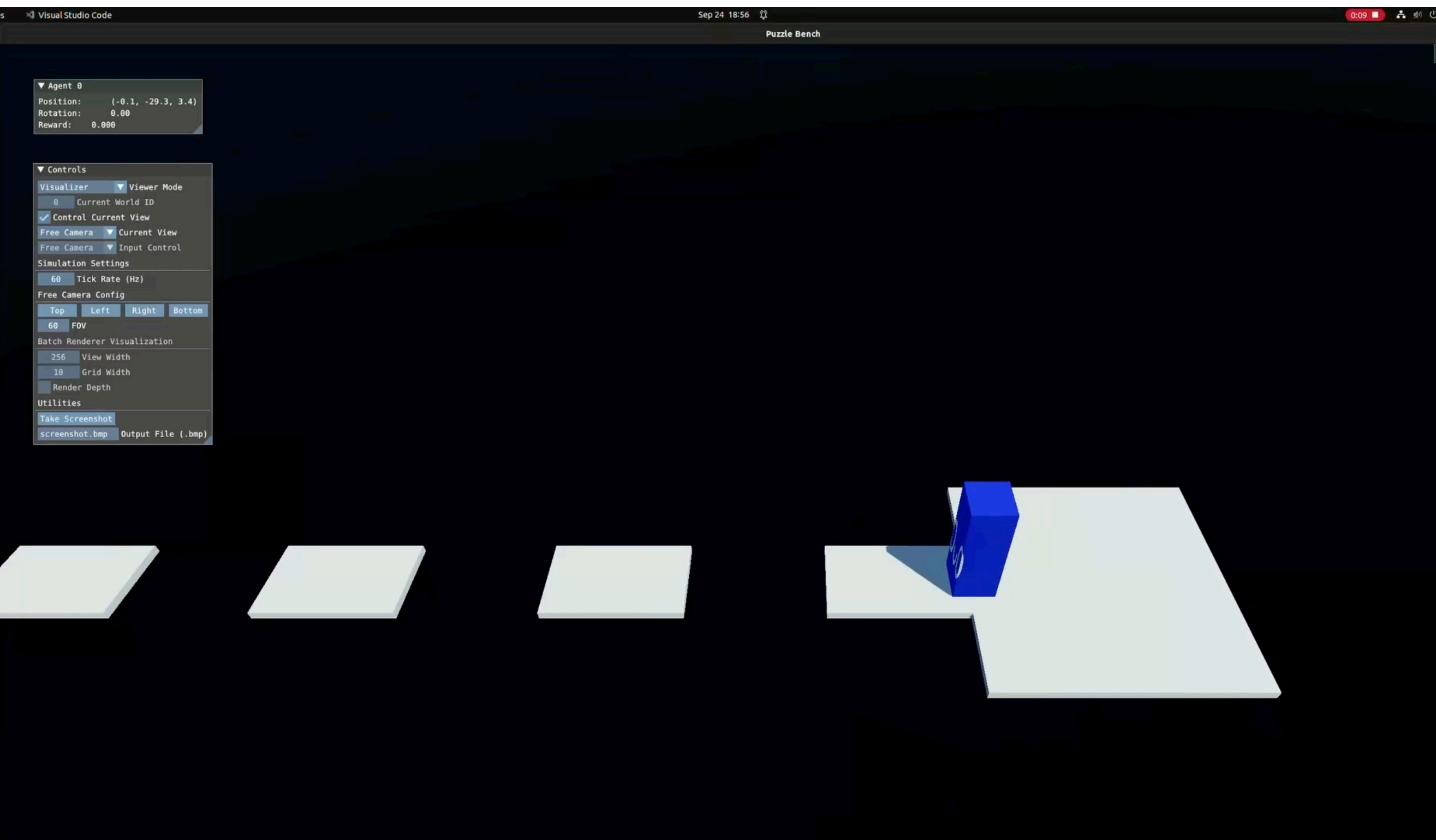
Training cooperative agents to exist an escape room



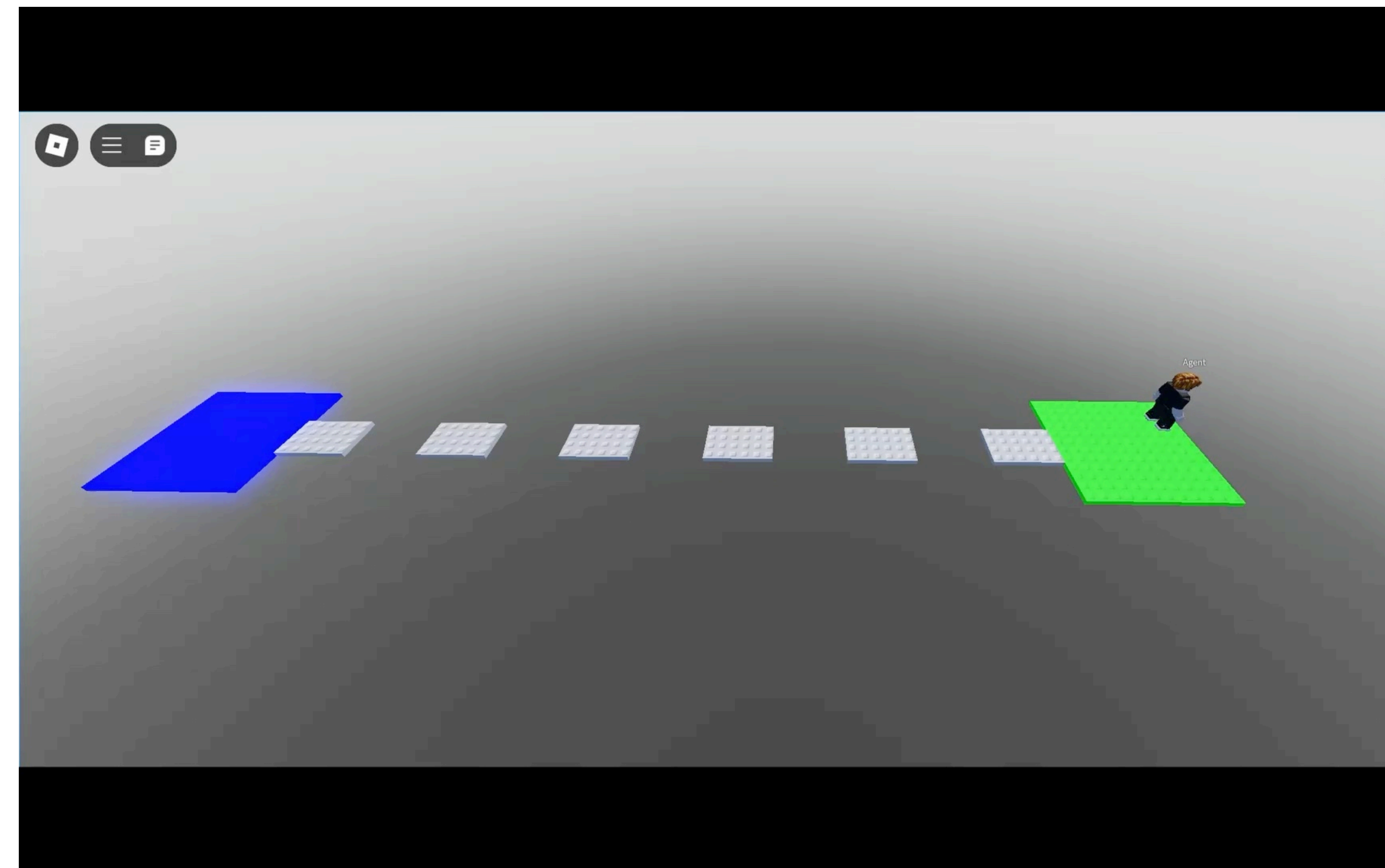
3.2B steps of experience in 1.5 hours on V100

Implement simple proxies of real games in fast-simulator: sim-to-sim transfer

1M steps/sec simulation in low-fi simulator



Same policy running in 60Hz Roblox engine

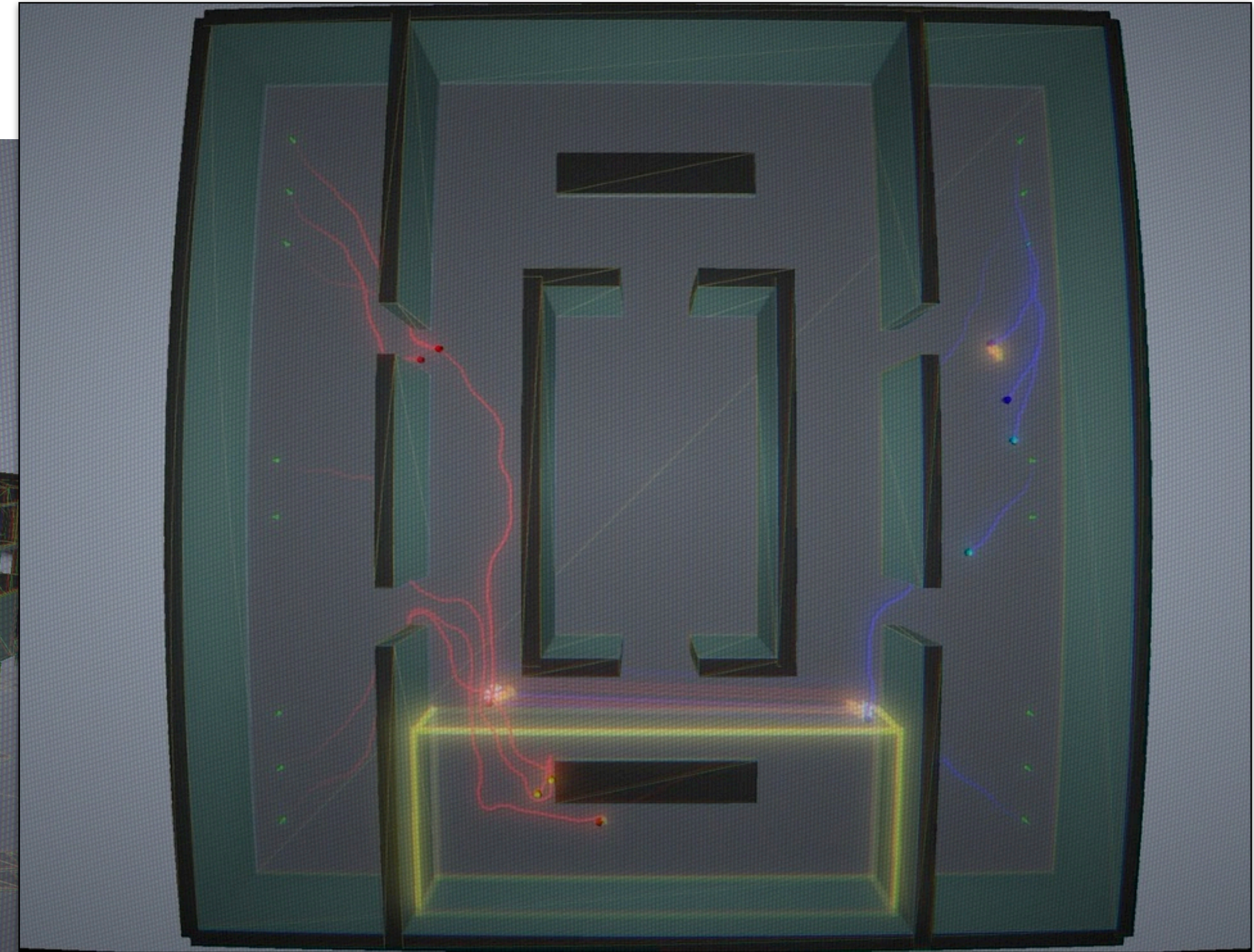
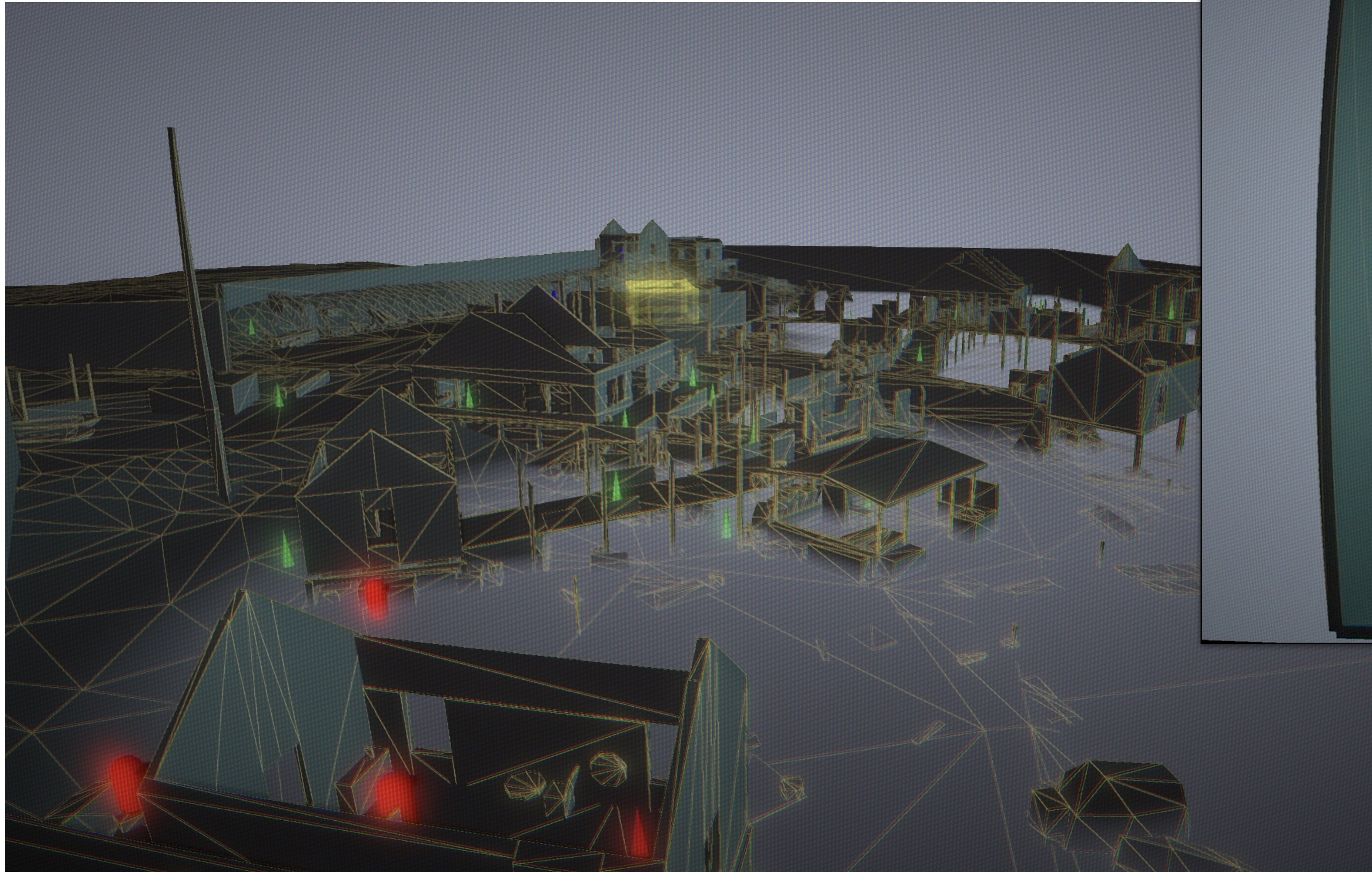


Sim-to-Sim: policies trained in coarse sim in a day can transfer directly to unseen obstacles in a higher-fidelity simulator (Roblox)

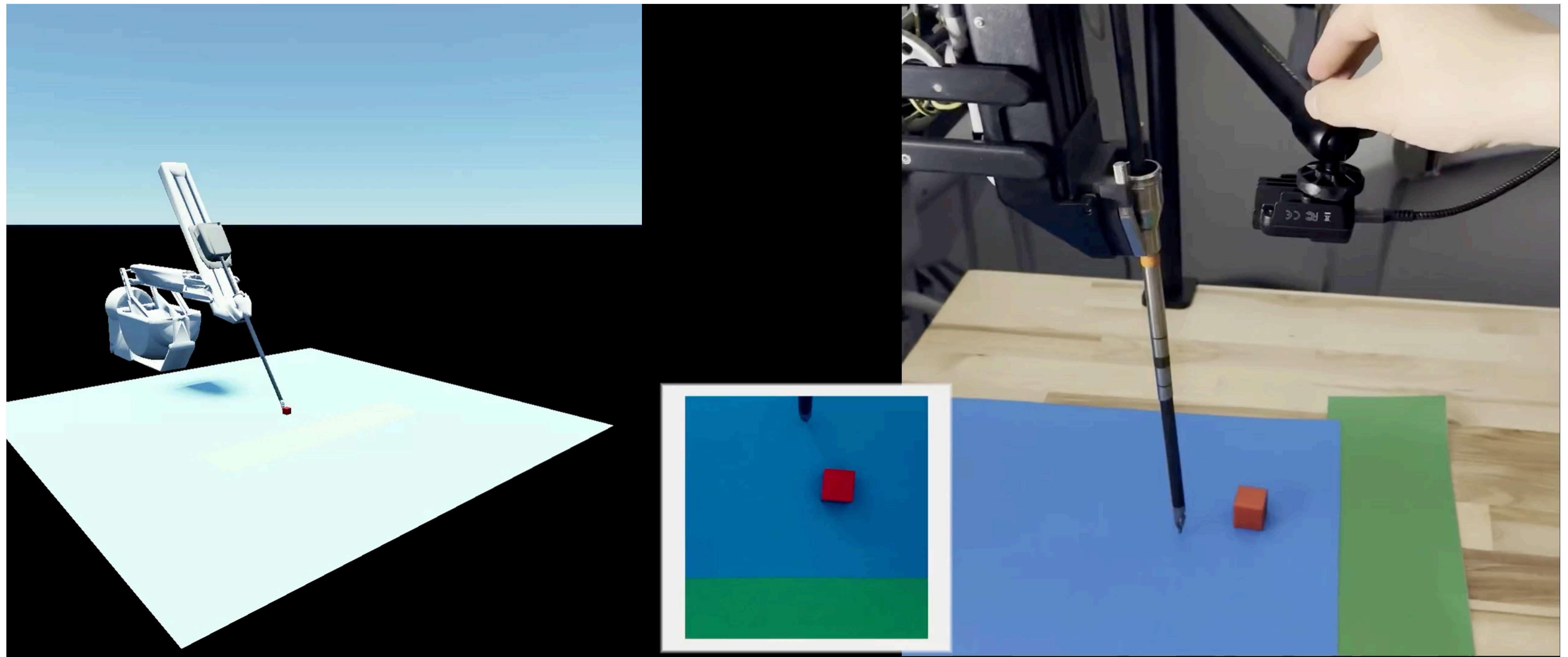


We tried learning competitive strategies for 6 vs 6 FPS play

6 x 6 FPS simulator (collaboration with Activision)

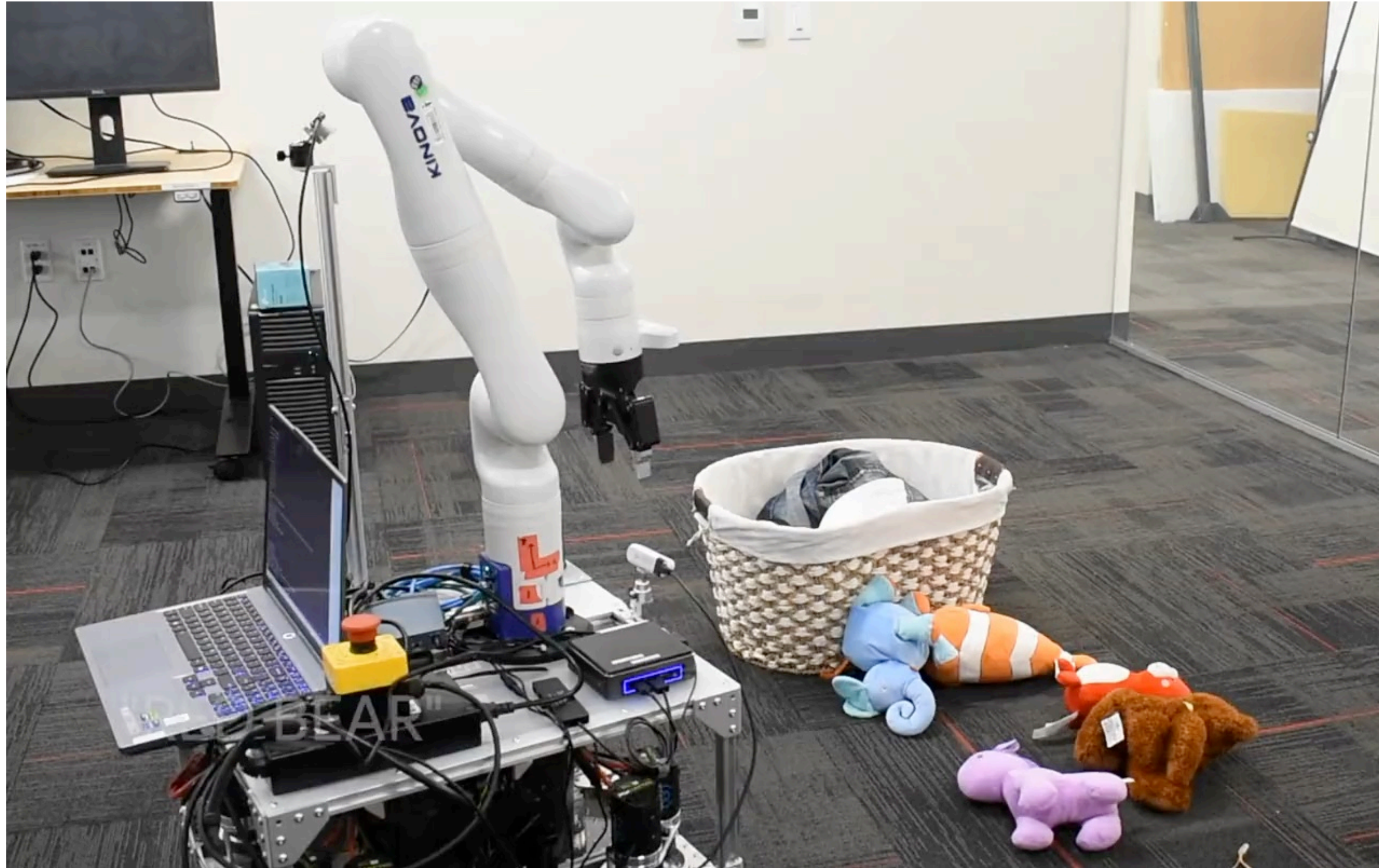


**Leveraging system efficiency:
simple sim-to-real transfer via domain randomization
(100x more samples = a lot of domain randomization)**



**But how do we make sufficiently
accurate simulations to train useful
policies?**

Consider challenges of simulating these scenarios



Consider arguments for video-world model-based simulation.



Speed = level-of-detail = reasoning at the right level of abstraction

How can we identify and generate a world simulation at the right level of abstraction for a given learning task?

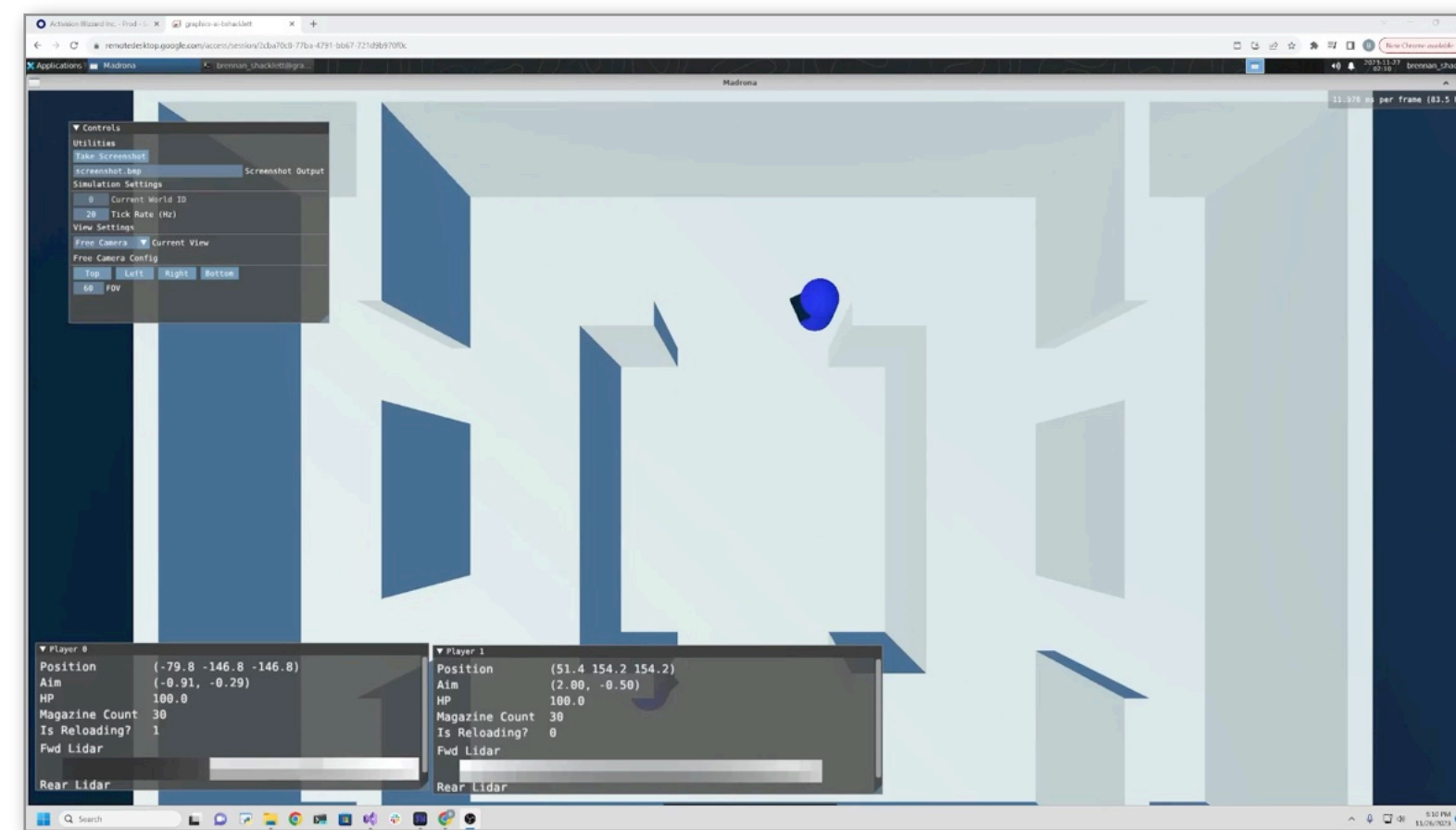
What aspects of problem solving should be done at each level of abstraction?

Text-based game

Low-fidelity sim

High-fidelity game

“You are a character in a battle Royale FPS game, you are standing near a corner and you see an enemy peak out...”



LLM policy action:
“You should take cover”

Low-fi agent action:
Move to (x,y)

Game agent policy action:
Game controller input

Project preso tips

Project presentation logistics

- **Next Tuesday June 2 (in CoDa)**
 - **We'll send out a sign up sheet for slots**
 - **Everyone in the class is welcome to attend any preso**
- **Each group has 10 minute "slot" to talk**
 - **Aim for at most 8 minutes of speaking**
- **I would like all team members to participate in delivering part of the presentation**
- **There is no required format or template, but here are a few tips**
- **Note: your writeup/handin is due 11:59pm on Tuesday**

Communicating like an architect (Aka. project presentation tips)

A few course themes

- ***Thinking like a systems architect***
 - **What are inputs/outputs, constraints, and goals?**
 - **What are the “services” the system should perform (what is hard for a user to do in a world without the system)**
 - **Once you can establish answers to these questions, you can consider your solution options**
- **Knowledge of applications and systems is necessary to choose efficient solutions**
 - **Algorithms folks use their hammer to reduce cost of algorithms (smaller DNN models, new optimization hyperparameters, use different data structures, etc.)**
 - **HW designers use their hammer (custom accelerators, new interconnects) to execute a given workload faster**
 - **SW systems folks use their hammer (parallel/distributed computing, workload-specific scheduling, high level programming abstractions, etc.)**
- **The best solutions architects pick the right mixture of hammers for the job**

A few course themes

- **Knowledge of applications and systems is necessary to do meaningful evaluation**
 - **Is this the right workload to evaluate?**
 - **Pitfall: measuring speedup on a part of the workload that is not the most significant**
 - **Does dataset I'm using have the right data distribution?**
 - **Am I measuring cost in FLOPs, but increasing data movement?**
 - **Am I optimizing an algorithm that is not the right choice of algorithm for this problem?**

Why get good at communication?

- **Systems architects need to be some of the best communicators in an organization**
- **Must be able to:**
 - **Communicate with users to understand workloads and constraints**
 - **Communicate with the individual contributors/engineers to:**
 - **Understand emerging problems/constraints**
 - **Understand how to evolve/extending existing designs to enable new functionality**
 - **Must constantly be communicating to various parties why:**
 - **Their desired features are not being added**
 - **They must do the same work with 50% of the resources**
 - **Communicate to everyone a strategic vision for a system**
 - **Communicate to executives/management/funders how goals are being met.**

My motivation

- **I have found I give nearly the same feedback over and over to students making talks**
 - **It is not profound feedback, it is just application of a simple set of techniques and principles that are consistently useful when making talks**

- **I am hoping these slides serve as a useful checklist you can refer to vet your own project presentation talks before giving them next Tuesday**
 - **Don't worry: I still make these mistakes all the time when creating first drafts of talks**

Who painted this painting?



Salvador Dali (age 22)

My point: learn the basic principles before you consciously choose to break them



Put yourself in your audience's shoes

Tip 1

This is a major challenge for most technical speakers. (including professors)

(Tip: recite a sentence out loud to yourself. * Do you really expect someone who has not been working with you everyday on the project to understand what you just said?)

*** I'm not kidding. Say it out loud. I find hearing myself say something out loud makes it easier to parse it from an audience's perspective.**

Consider your audience

- **Everyone in the audience knows about course readings/topics**
 - Terminology/concepts we all know about need not defined (just say “remember we talked about X”)
- **Most of the audience knows little-to-nothing about the specific application domain or problem you are trying to solve**
 - Application-specific terminology should be defined or avoided
 - What they can get their head around is inputs/outputs and goals/constraints
- **Everyone wants to know the “most interesting” thing that you found out or accomplished (your job is to define most interesting for them)**

Tip 2

**A good principle for any talk (or paper):
“Every sentence matters”**

**What are you trying to say?
What technical story are you trying to tell?
What is point you are trying to make?**

You might be trying to define a hypothesis.

Or define goals

Or establish inputs/outputs

**Or show data that suggests you
were successful.**

Is what you just said making that point? (If not, remove it)

If you can't justify how it will help the listener understand the point, take it out.

If it's not likely to be clear, take it out.

Pick a focus

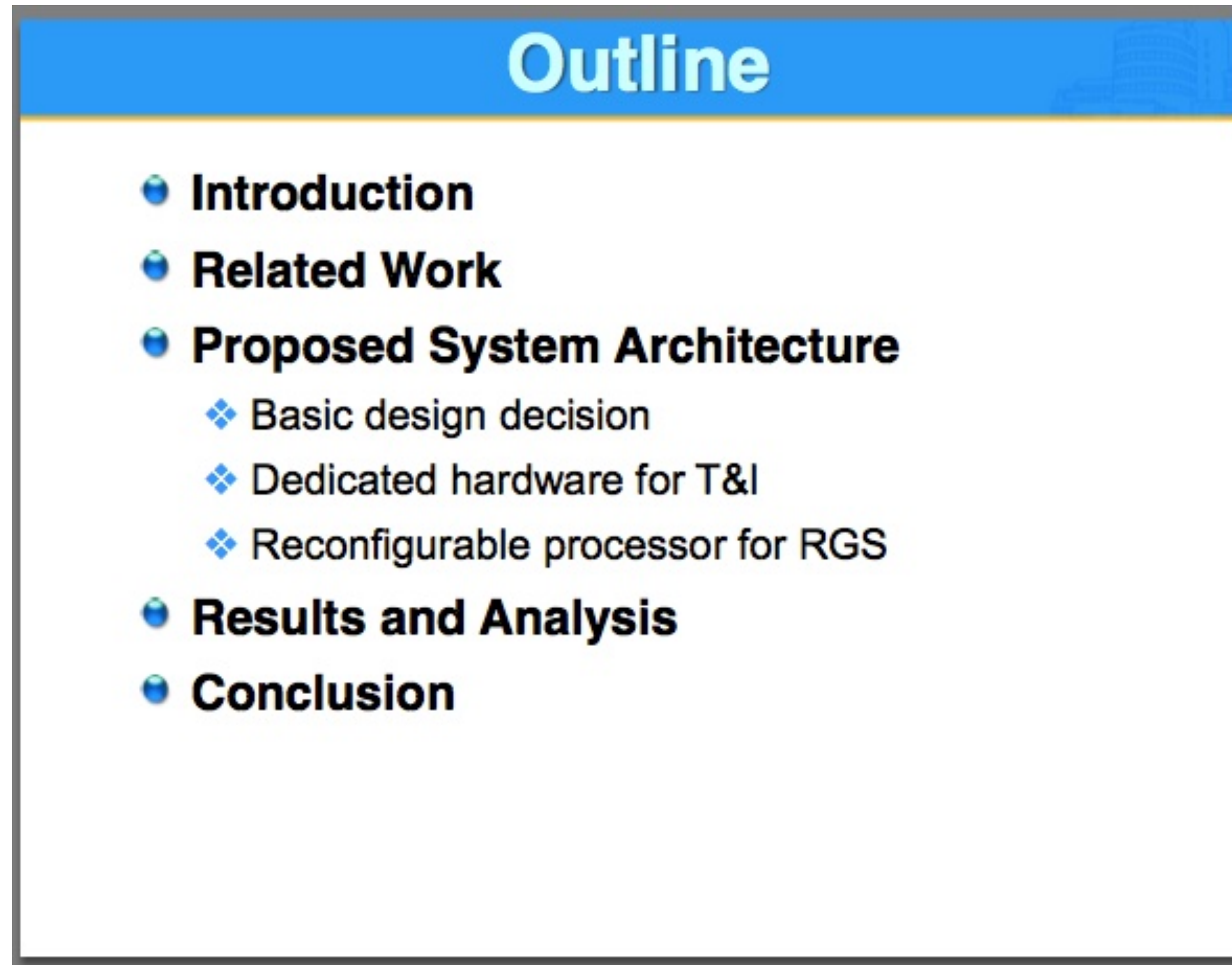
- In this class, different projects should stress different results
- Some projects may wish to show a flashy demo and describe how it works (proof by “it works”)
- Other projects may wish to show a sequence of graphs (path of progressive optimization) and describe the optimization that took system from performance A to B to C
- Other projects may wish to clearly contrast parallel CPU vs. parallel GPU performance for a workload

Your job is not to explain what you did, but to explain what you think we should know.

**And really the most important thing we want to know is (1) what was your goal? and
(2) what's the evidence you have that you were successful?**

Ignoring every sentence matters

Never ever, ever, ever do this!



Outline

- **Introduction**
- **Related Work**
- **Proposed System Architecture**
 - ❖ Basic design decision
 - ❖ Dedicated hardware for T&I
 - ❖ Reconfigurable processor for RGS
- **Results and Analysis**
- **Conclusion**

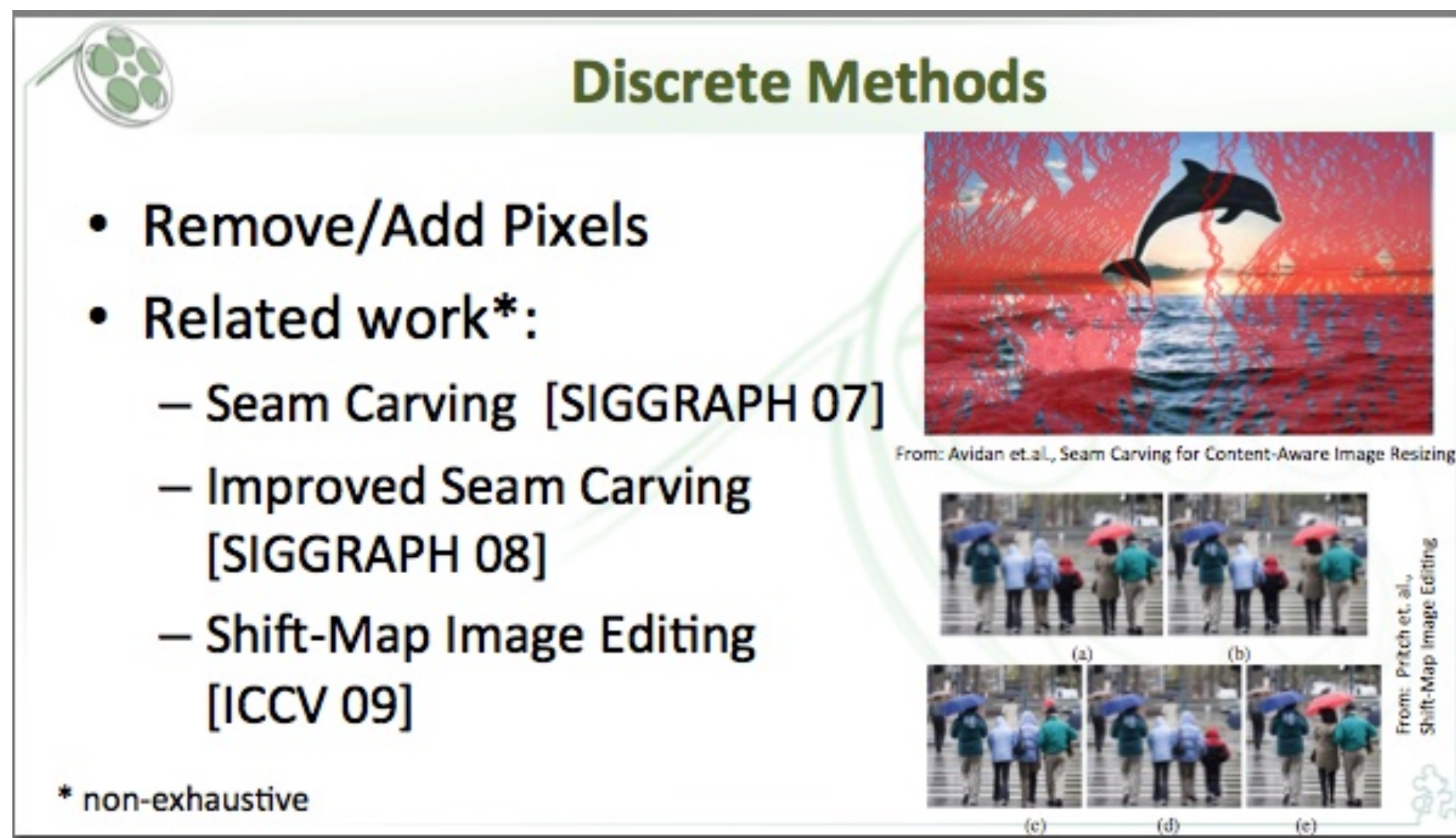
Bad example 2

■ Who is the audience for this? (how does this benefit them?)

Discrete Methods

- Remove/Add Pixels
- Related work*:
 - Seam Carving [SIGGRAPH 07]
 - Improved Seam Carving [SIGGRAPH 08]
 - Shift-Map Image Editing [ICCV 09]

* non-exhaustive



From: Avidan et al., Seam Carving for Content-Aware Image Resizing

From: Pritch et al., Shift-Map Image Editing

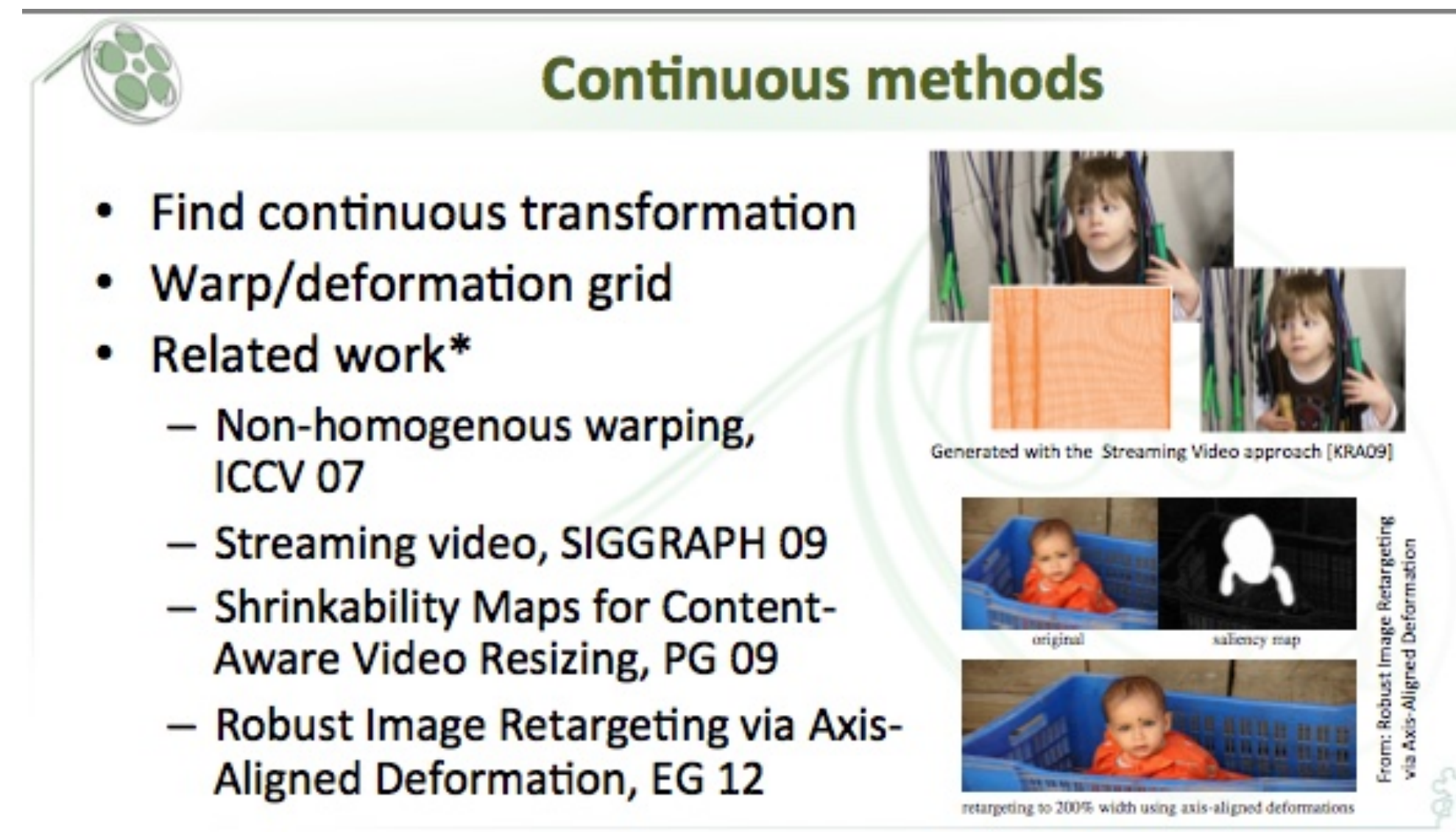
This type of related work section says little more than “others have worked in this area before”.

- I suspect your audience assumes this is the case.
- Every sentence matters: if it doesn't provide value, take it out (or replace it with comments that do provide value)

- **Experts?**
 - They likely know these papers exist. These slides don't tell them what about these papers is most relevant to this talk
- **Non-experts?**
 - They won't learn the related work from these two slides

Continuous methods

- Find continuous transformation
- Warp/deformation grid
- Related work*
 - Non-homogenous warping, ICCV 07
 - Streaming video, SIGGRAPH 09
 - Shrinkability Maps for Content-Aware Video Resizing, PG 09
 - Robust Image Retargeting via Axis-Aligned Deformation, EG 12



Generated with the Streaming Video approach [KRA09]

original saliency map

retargeting to 200% width using axis-aligned deformations

From: Robust Image Retargeting via Axis-Aligned Deformation

Tip 3

The audience prefers not to think (much)

The audience has a finite supply of mental effort

- **The audience does not want to burn mental effort about things you know and can just tell them.**
 - They want to be led by hand through the major steps of your story
 - They do not want to interpret any of your figures or graphs, they want to be directly told how to interpret them (e.g., what to look for in a graph).
 - They want to be told about your key assumptions
- **The audience does want to spend their energy thinking about:**
 - Potential problems with what you did (did you consider all edge cases? Is your evaluation methodology sound? Is this a good platform for this workload?)
 - Implications of your approach to other things
 - Connections to their own project or interests

Tip 4

Set up the problem.

**Establish inputs, outputs, and constraints
(goals and assumptions)**

Basics of problem setup

- **What is the computation performed (or system built)?**
 - **What are the inputs? What are the outputs?**
- **Why does this problem stand to benefit from optimization?**
 - **“Real-time performance could be achieved”**
 - **“Researchers could run many more trials, changing how science is done”**
 - **“It is 90% of the execution time in this particular system”**
- **Why is it hard? (What made your project interesting? What should we reward you for?)**
 - **What turned out to be the hardest part of the problem?**
 - **Optimization projects: e.g., overcoming SIMD divergence, increasing arithmetic intensity**
 - **Applications projects: coming up with an algorithm to estimate depth, dealing with motion blur**
 - **Abstraction/API projects: deciding between stateless or stateful interfaces**

Example: 3D rendering problem

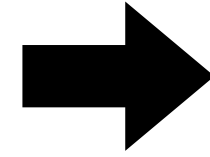
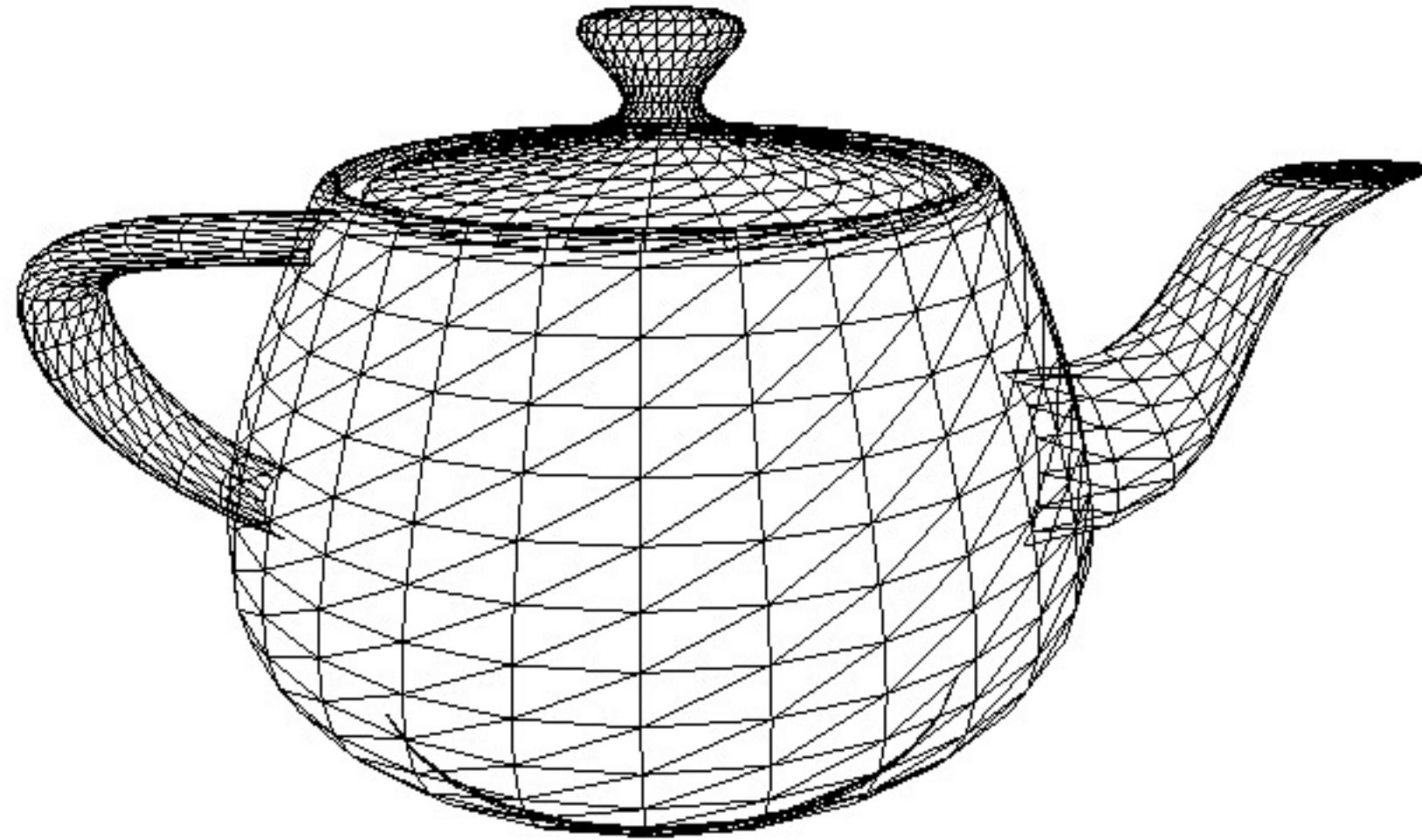


Image credit: Henrik Wann Jensen

Input: description of a scene:

3D surface geometry (e.g., triangle mesh)
surface materials, lights, camera, etc.

Output: image of the scene

Simple definition of rendering task: computing how each triangle in 3D mesh contributes to appearance of each pixel in the image?

Why is knowing the goals and constraints important?

Your contribution is typically a system or algorithm that meets the stated goals under the stated constraints.

Understanding whether a solution is “good” requires having this problem context.

Tip 5

How to describe a system

How to describe a system

■ Start with the nouns (the key boxes in a diagram)

- Major components (processors, memories, interconnects, etc.)
- Major entities (particles, neighbor lists, pixels, pixel tiles, features, etc.)
- What is state in the system?

■ Then describe the verbs

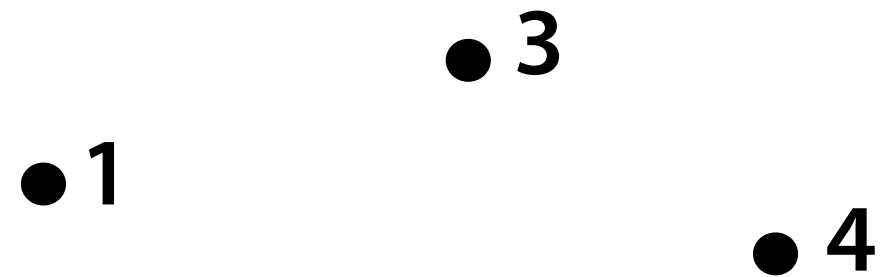
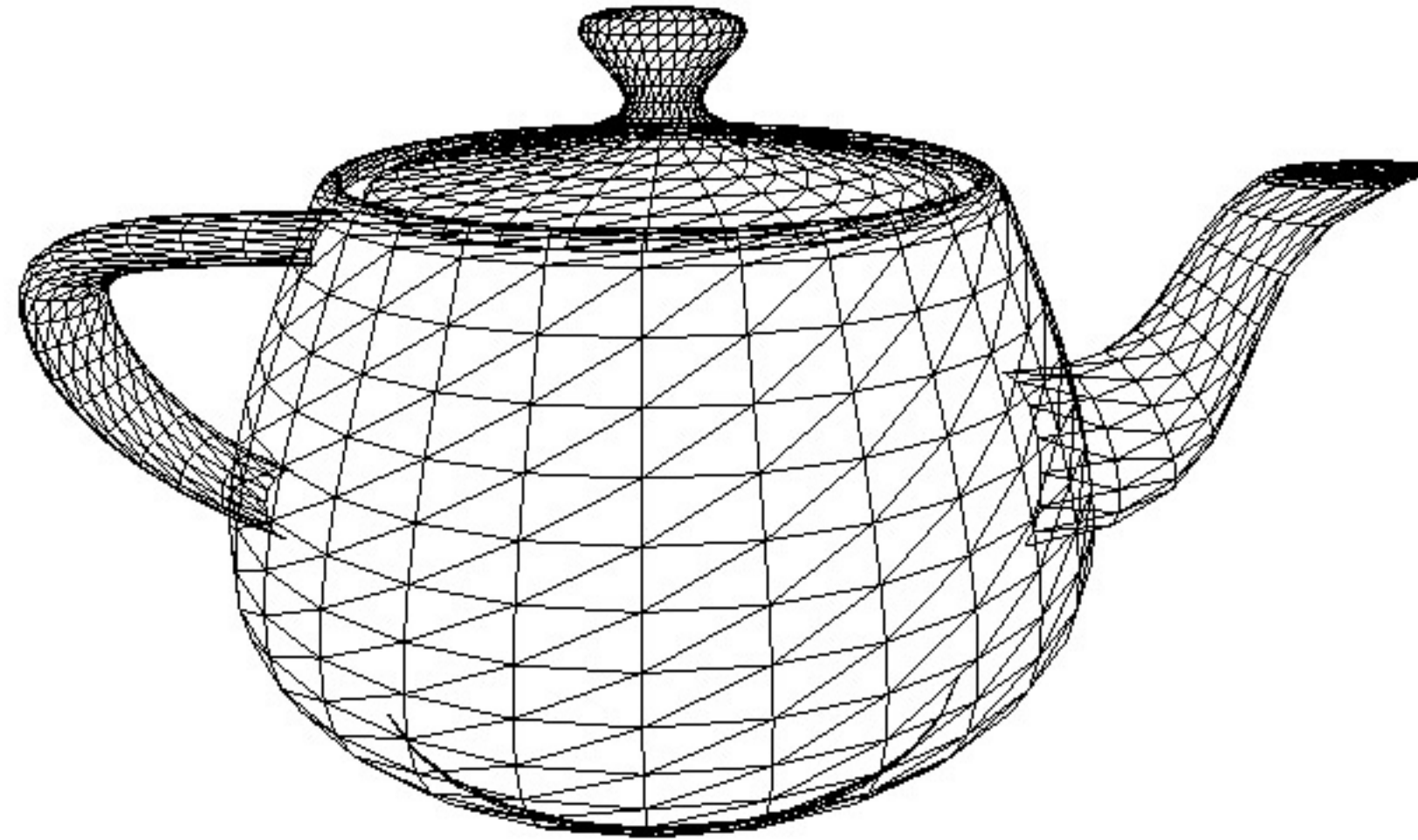
- Operations that can be performed on the state (update particle positions, compute gradient of pixels, traverse graph, etc.)
- Operations produce, consume, or transform entities

Tip: how to explain “a system”

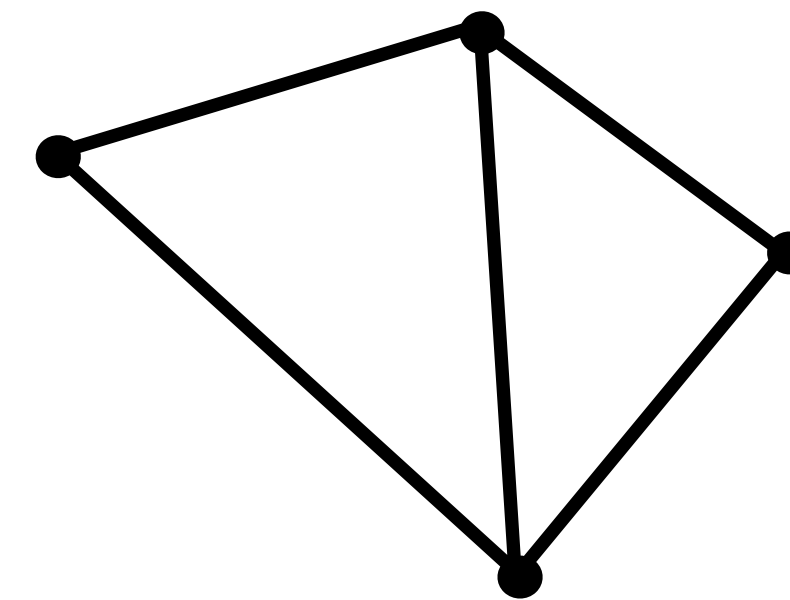
- Step 1: describe the things (key entities) that are manipulated
 - The nouns

Example: real-time graphics primitives (entities)

Represent surface as a 3D triangle mesh

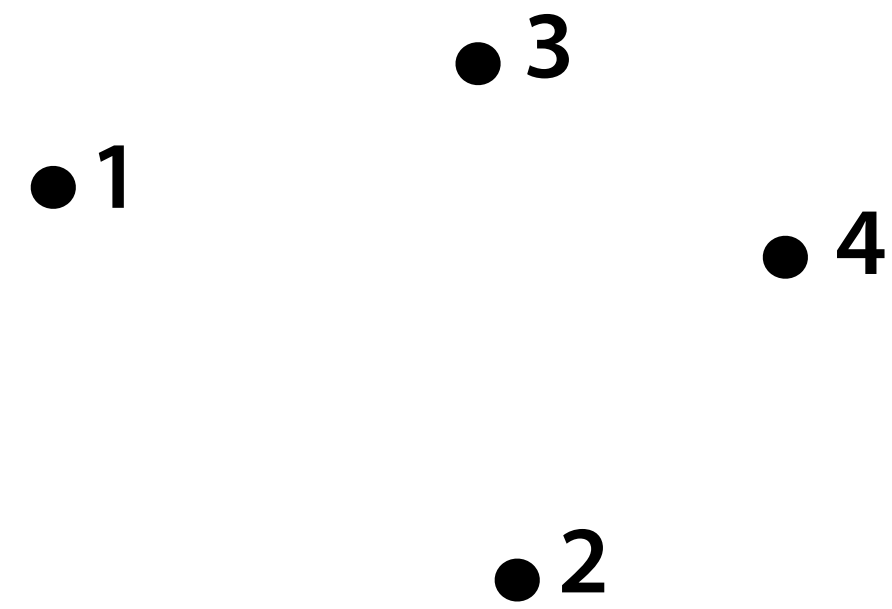


Vertices
(points in space)

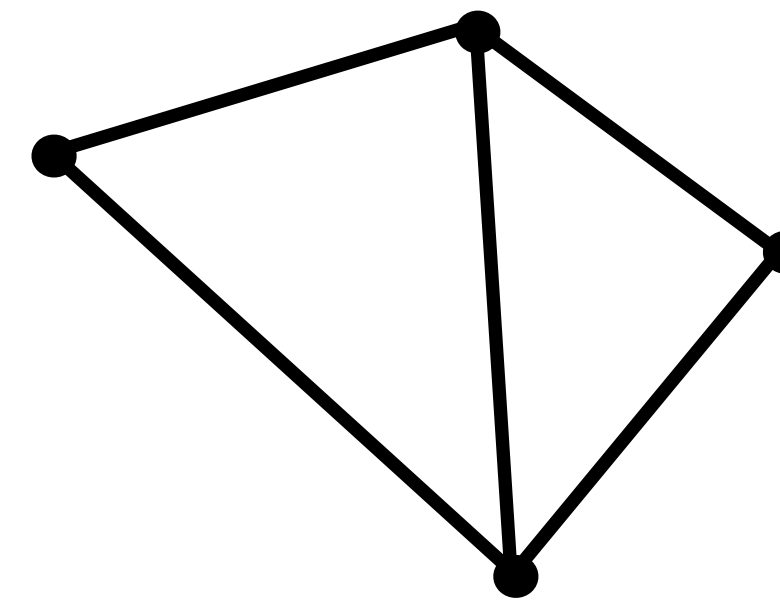


Primitives
(e.g., triangles, points, lines)

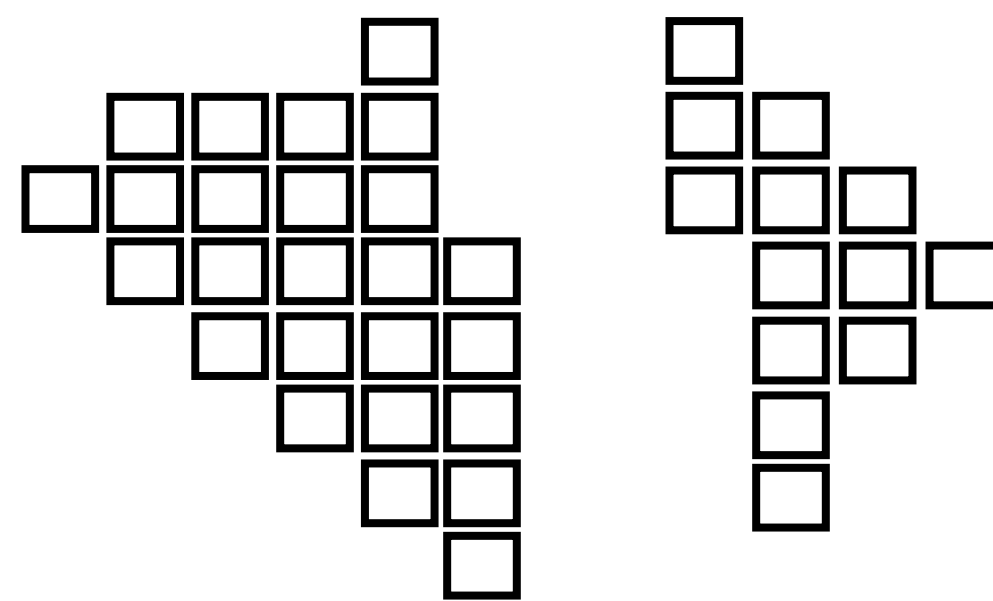
Real-time graphics primitives (entities)



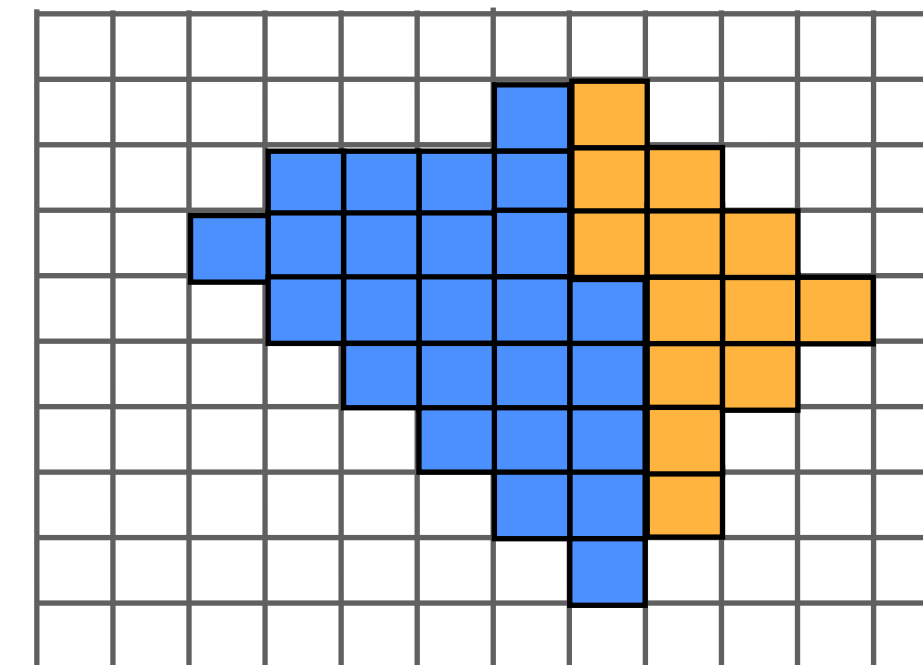
Vertices
(points in space)



Primitives
(e.g., triangles, points, lines)



Fragments



Pixels (in an image)

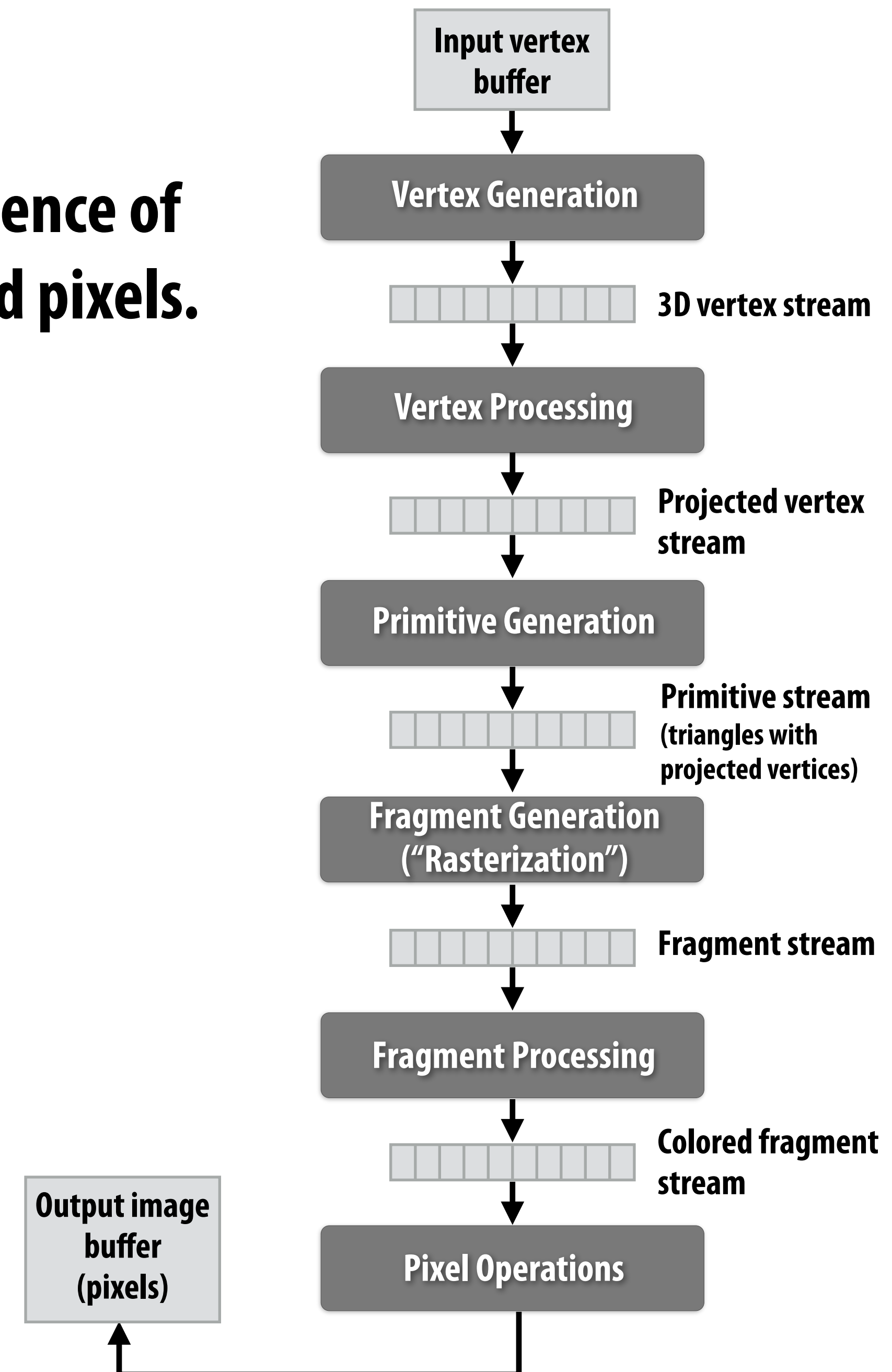
How to explain “a system”

- **Step 1: describe the things (key entities) that are manipulated**
 - **The nouns**

- **Step 2: describe the operations the system performs on these entities**
 - **The verbs**

Real-time graphics pipeline

- Abstracts process of rendering a picture as a sequence of operations on vertices, primitives, fragments, and pixels.



Tip 6

Surprises* are almost always bad:

Say where you are going and why you must go there before you say what you did.

*** I am referring to surprises in talk narrative and/or exposition. A surprising result is great.**

Give the why before the what

■ Why provides the listener context for...

- **Compartmentalizing: assessing how hard they should pay attention (is this a critical idea, or just an implementation detail?). Especially useful if they are getting lost.**
- **Understanding how parts of the talk relate (“Why is the speaker now introducing a new optimization framework?”)**

■ In the algorithm description:

- **“We need to first establish some terminology”**
- **“Even given X , the problem we still haven’t solved is...”**
- **“Now that we have defined a cost metric we need a method to minimize it...”**

■ In the results/evaluation:

- **Speaker: “Key questions to ask about our approach are...”**
- **Audience: “Thanks! I agree, those are good questions. Let’s see what the results say!”**

Two key questions:

- How much does SRDH improve traversal cost when perfect information about shadow rays is present?
- How does the benefit of the SRDH decrease as less shadow ray information is known a priori? (Is a practical implementation possible?)

Big surprises in a narrative are a bad sign

- **Ideally, you want the audience to always be able to anticipate* what you are about to say**
 - **This means: your story is so clear it's obvious!**
 - **It also means the talk is really easy to present without notes or text on slides (it just flows)**

- **If you are practicing your talk, and you keep forgetting what's coming on the next slide (that is, you can't anticipate it)...**
 - **This means: you probably need to restructure your talk because a clear narrative is not there.**
 - **It's not even obvious to you! Ouch!**

Tip 7

Show, don't tell

**It's much easier to communicate with
figures/images than text**

(And it saves the speaker a lot of work explaining... you can just describe the picture)

Example:

- **In a recent project, we asked the question... given enough video of tennis matches of a professional athlete, could we come up with an algorithm for turning all this input video into a controllable video game character?**

Compare the description above to the following sequence...

Here's an example of that source video

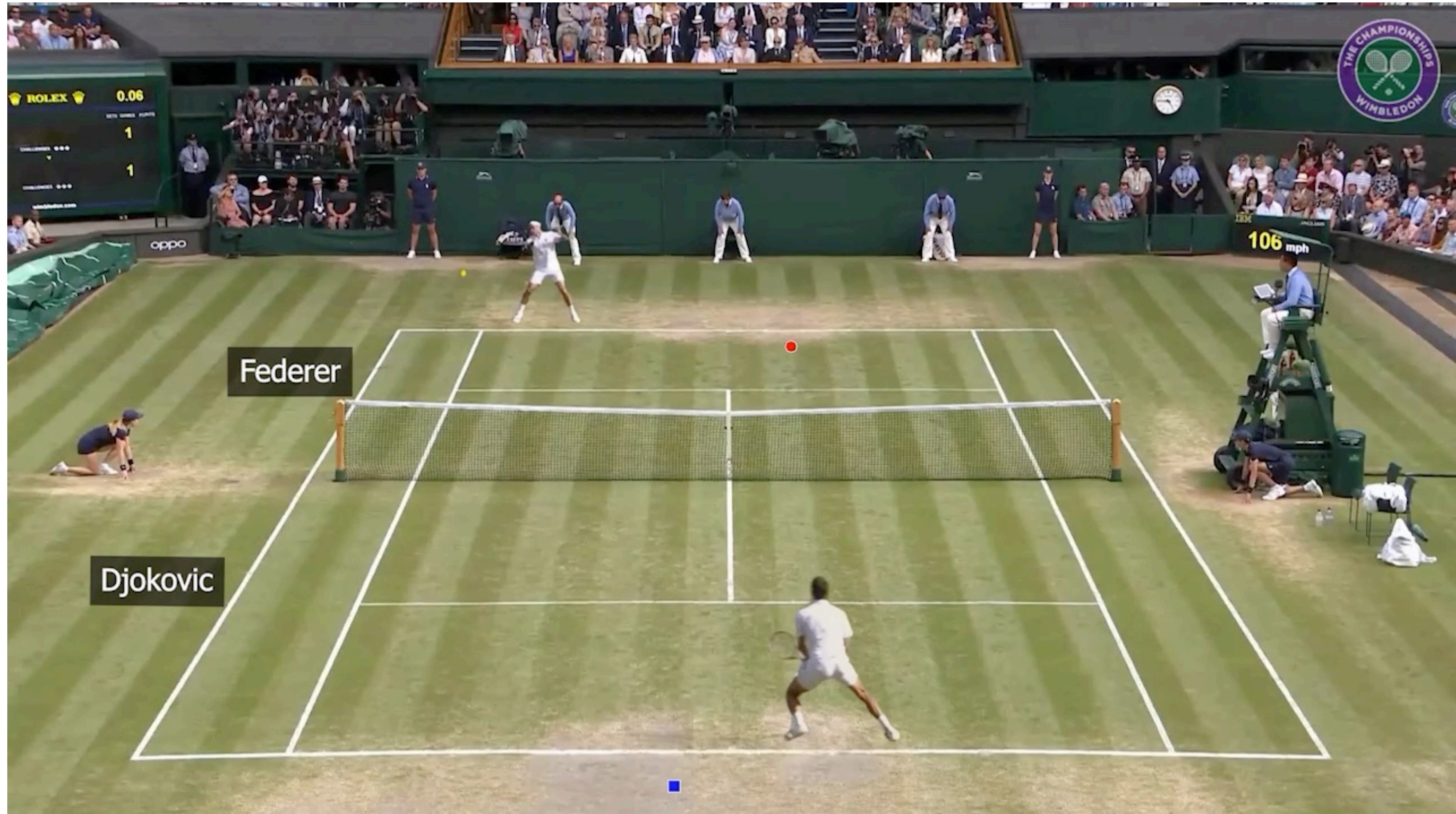


The best way to describe the input data is just show it! ("This is what the input looks like!")

And there's a lot of it out there!



And here's an example of controllable output



The best way to describe the output we seek is just show the result of the system!
("We click to specify a target ball location, and the player hits the incoming ball back to the red dot")

Another example:

The problem (lighting differences)



After the fix



Another example: we recently created a renderer that achieved high frame rates by rendering many views of the scene at the same time



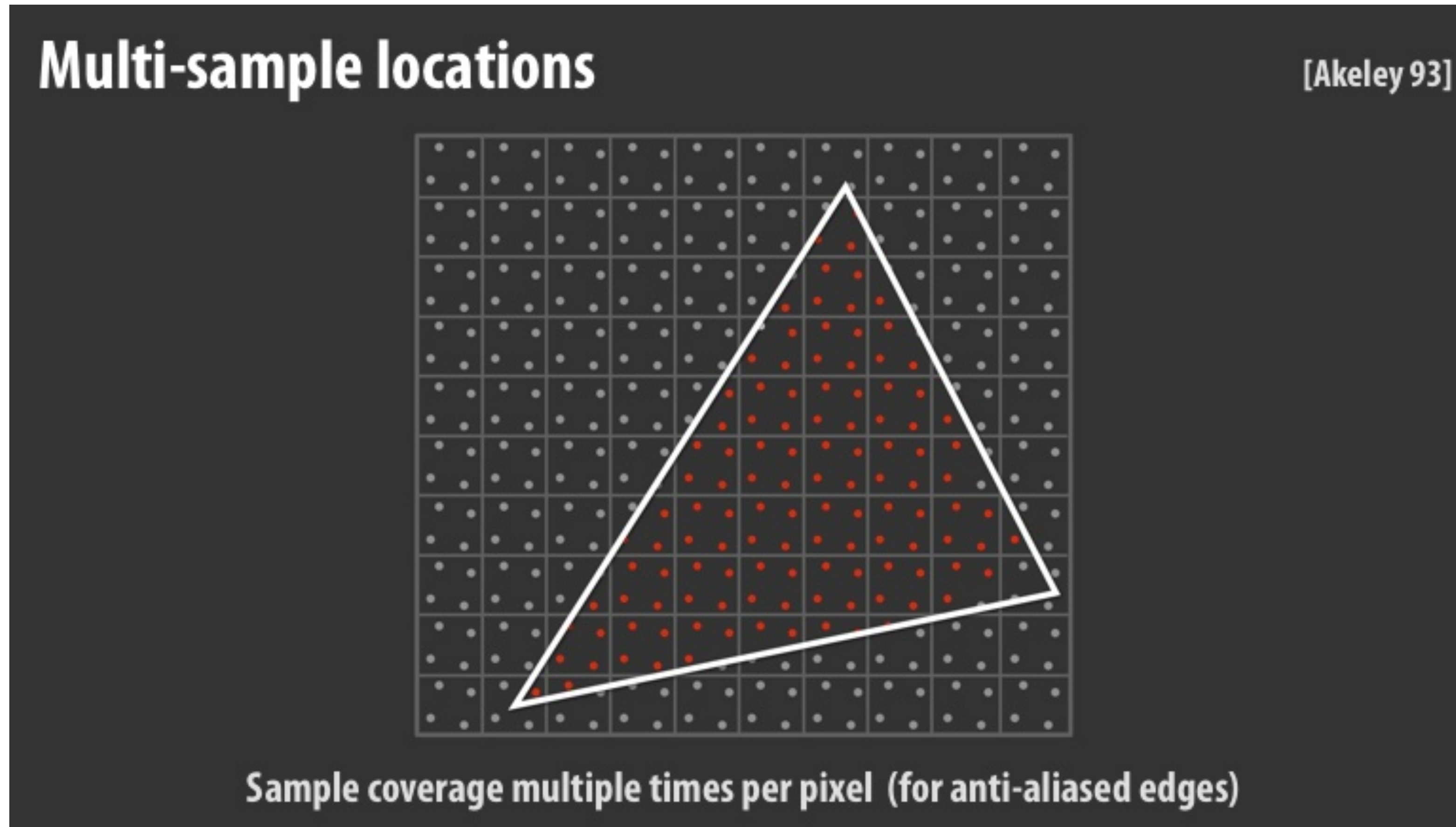
Tip 8

Always, always, always
explain any figure or graph

(remember, the audience does not want to think about things you can tell them)

Explain every figure

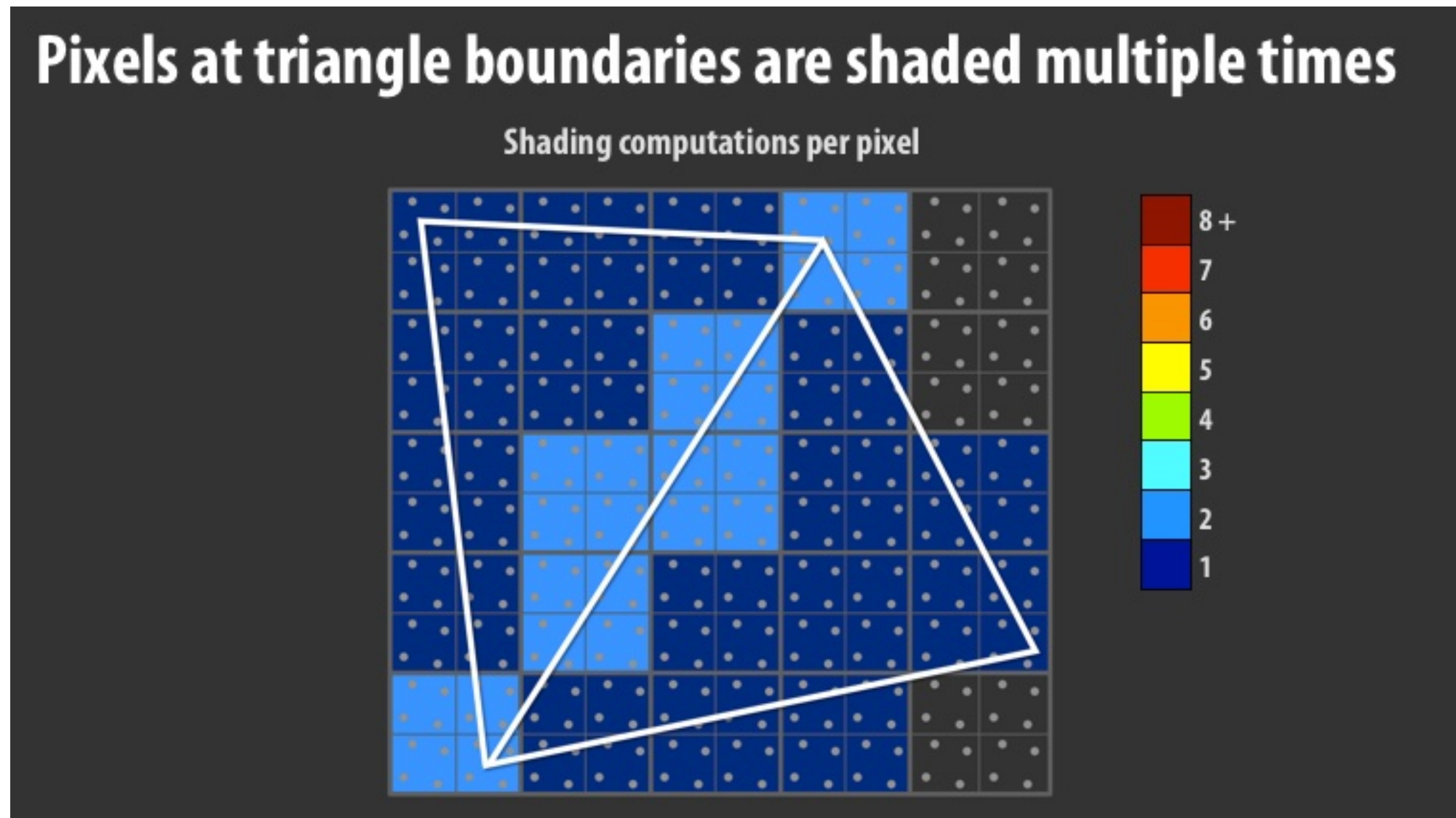
- Explain every visual element in the figure (never make the audience decode a figure)
- Refer to highlight colors explicitly (explain why the visual element is highlighted)



Example voice over: "Here I'm showing you a pixel grid, a projected triangle, and the location of four sample points at each pixel. Sample points falling within the triangle are colored red."

Explain every figure

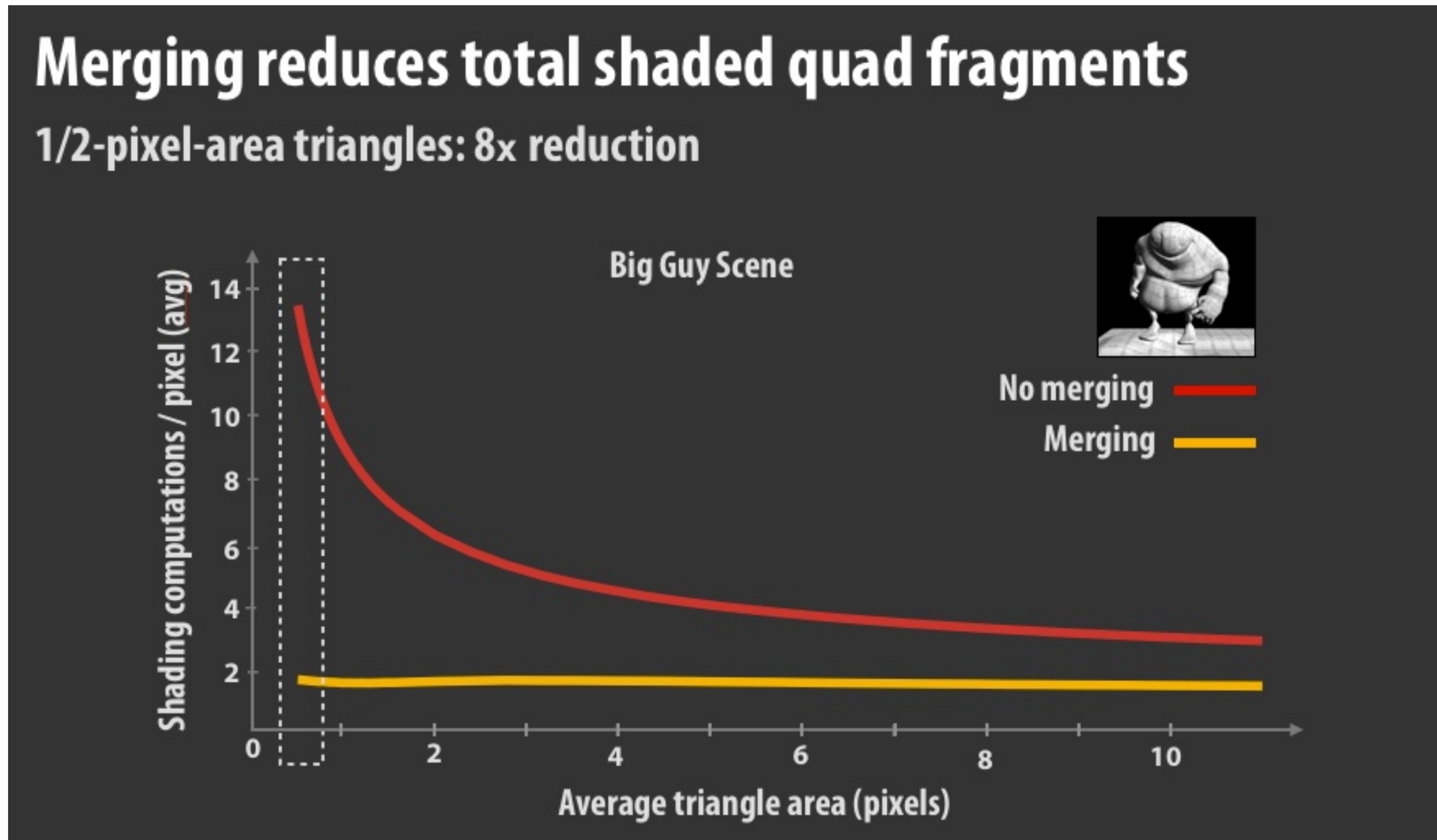
- Lead the listener through the key points of the figure
- Useful phrase: “As you can see...”
 - It’s like verbal eye contact. It keeps the listener engaged and makes the listener happy... “Oh yeah, I can see that! I am following this talk!”



Example voice over: “Now I’m showing you two adjacent triangles, and I’m coloring pixels according to the number of shading computations that occur at each pixel as a result of rendering these two triangles. As you can see from the light blue region, pixels near the boundary of the two triangles get shaded twice.

Explain every results graph

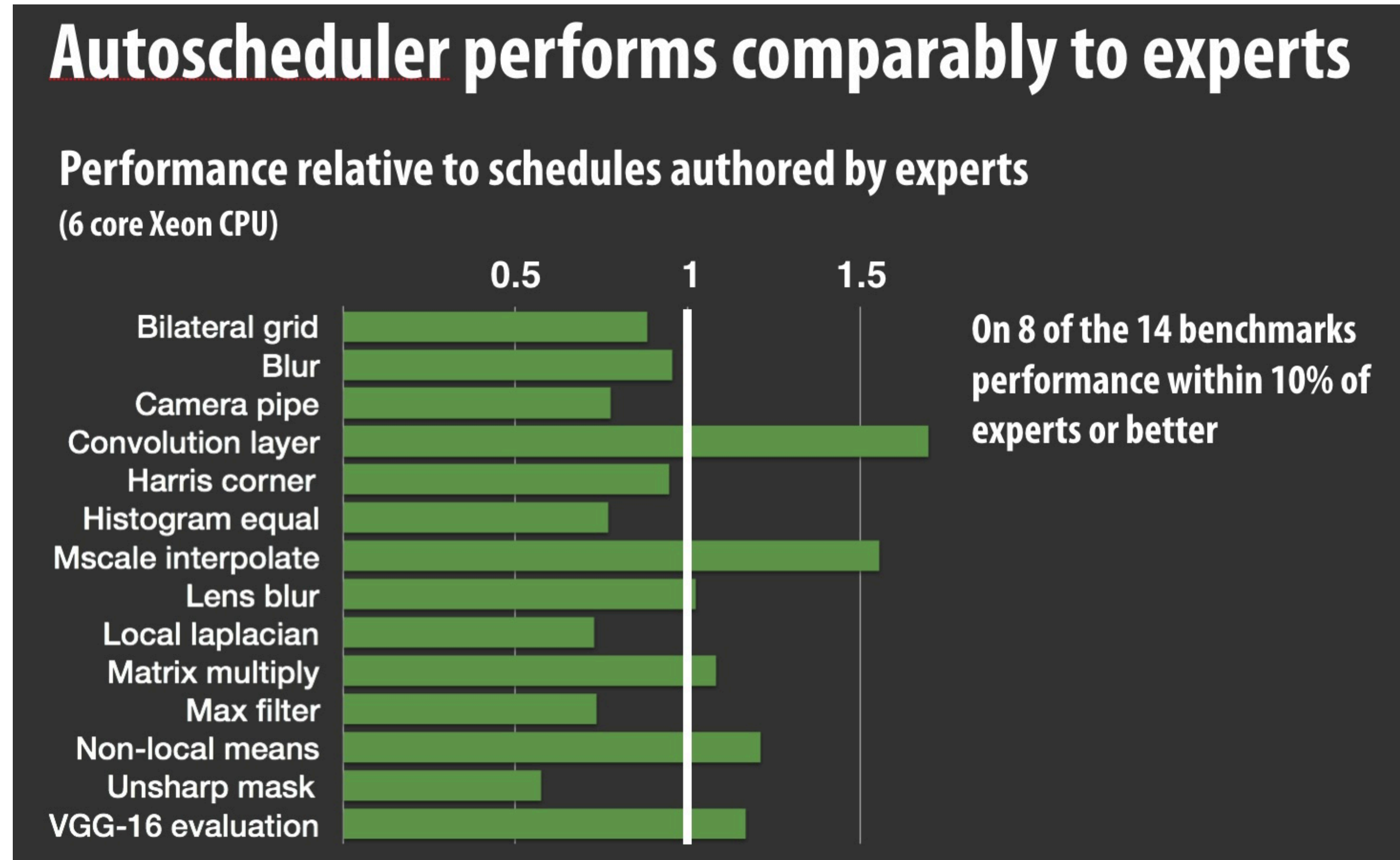
- May start with a general intro of what the graph will address (so audience “anticipates” the result)
- Then describe the axes (and your axes better have labels!)
- Then describe the one point that you wish to make with this results slide (more on this later!)



Example voice over: “Our first questions were about performance: how much did the algorithm reduce the number of the shading computations? And we found out that the answer is a lot. This figure plots the number of shading computations per pixel when rendering different tessellations of the big guy scene. X-axis gives triangle size. If you look at the left side of the graph, which corresponds to a high-resolution micropolygon mesh, you can see that merging, shown by yellow line, shades over eight times less than the convention pipeline.

Explain every results graph

- May start with a general intro of what the graph will address.
- Then describe the axes (your axes better have labels!)
- Then describe the one point that you wish to make with this results slide (more on this later!)



Example voice over: "Our first question was about performance: how fast is the auto scheduler compared to experts? And we found out that it's quite good. This figure plots the performance of the autoscheduler compared to that of expert code. So expert code is 1. Faster code is to the right. As you can see, the auto scheduler is within 10% of the performance of the experts in many cases, and always within a factor of 2."

Tip 9

In the results section:

One point per slide!

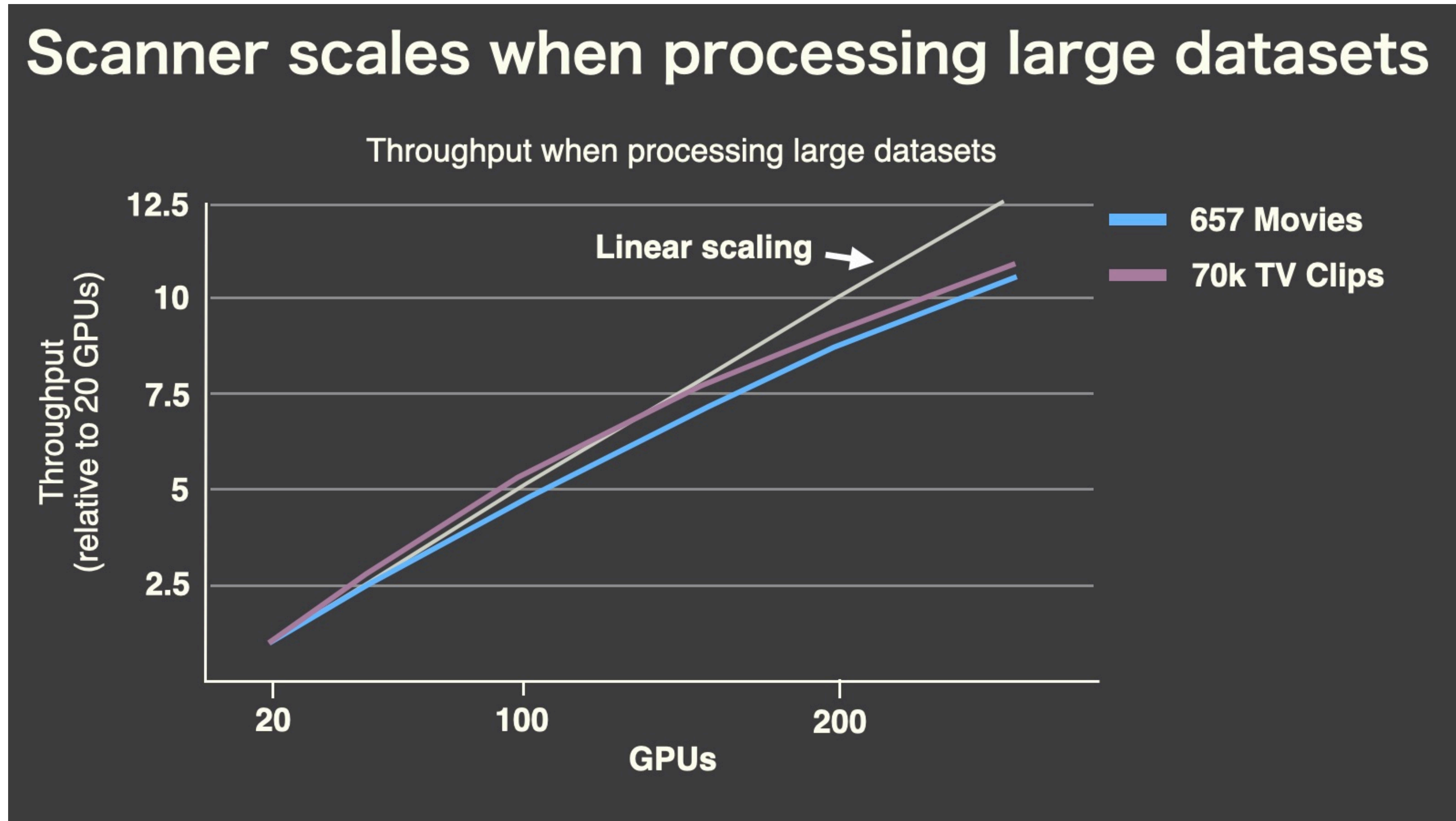
One point per slide!

One point per slide!

(and the point is the title of the slide!!!)

■ **Make the point of the graph the slide's title:**

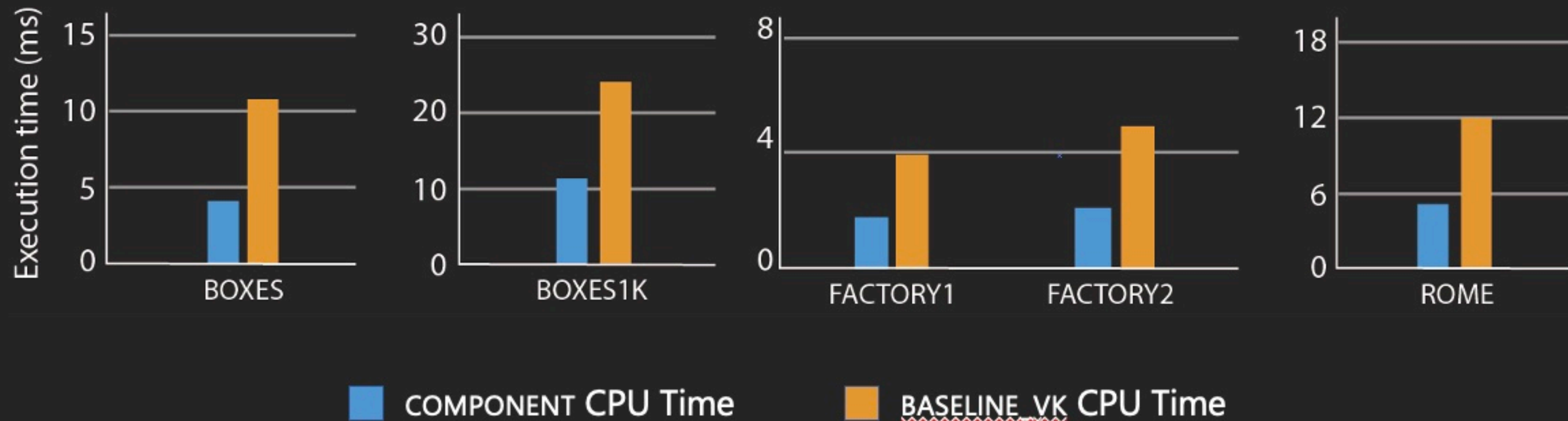
- **It provides audience context for interpreting the graph (“Let me see if I can verify that point in the graph to check my understanding”)**
- **Another example of the “audience prefers not to think” principle**



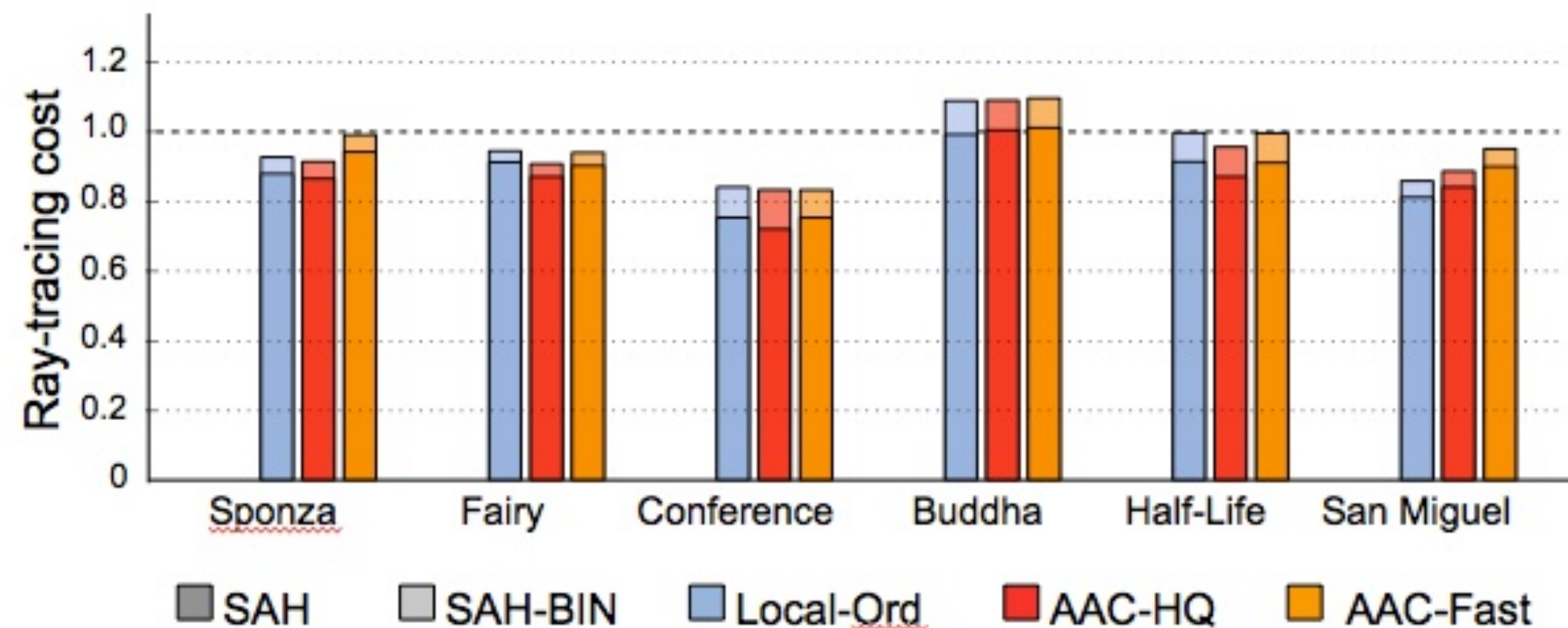
Another example:

The COMPONENTS renderer uses 2x less CPU time than BASELINE_VK

CPU Performance Comparison (single core)

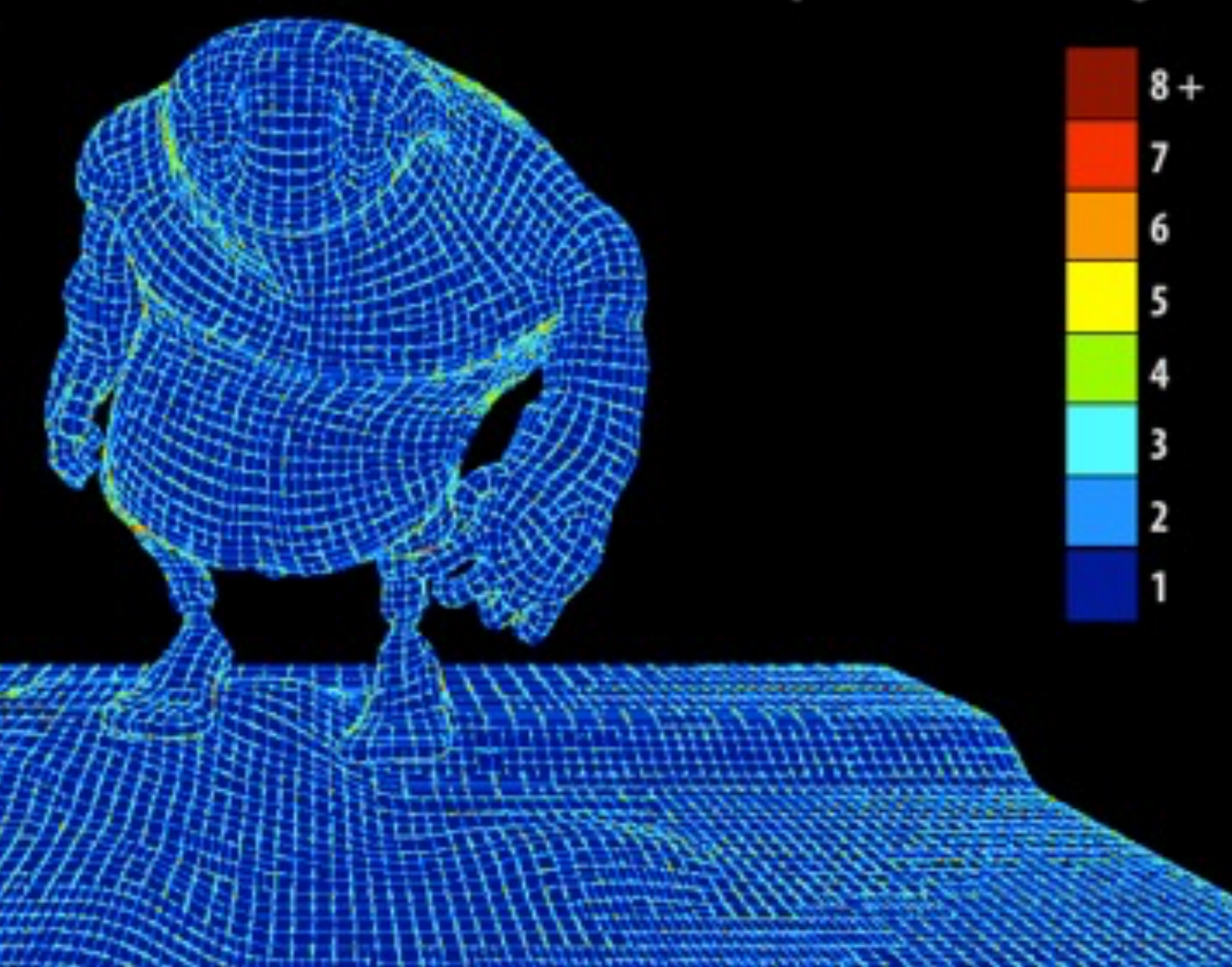


AAC-Fast produces BVHs with equal or lower cost than the **full sweep build** in all cases except Buddha.

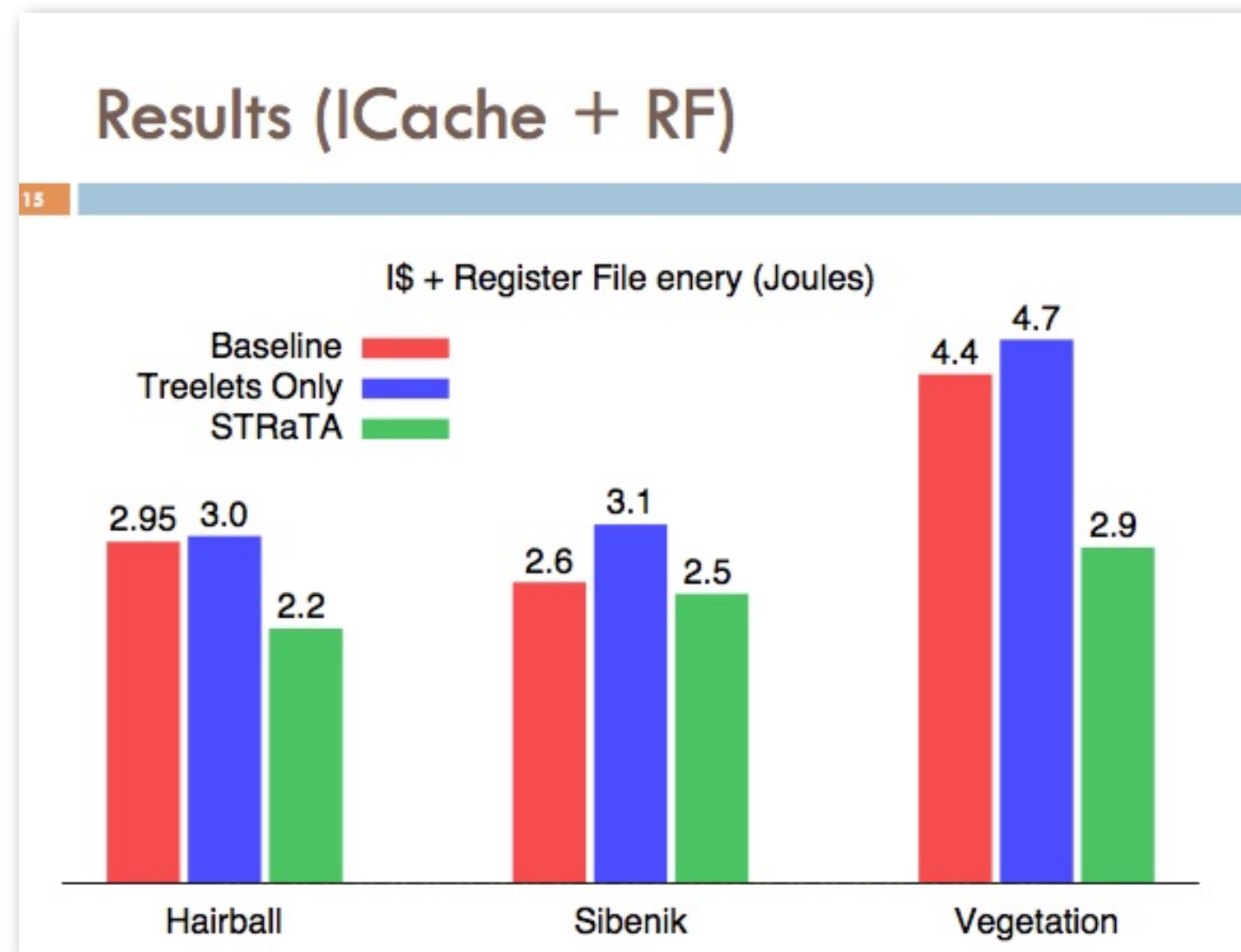


More examples

Extra shading occurs at merging window boundaries



Bad examples of results slides



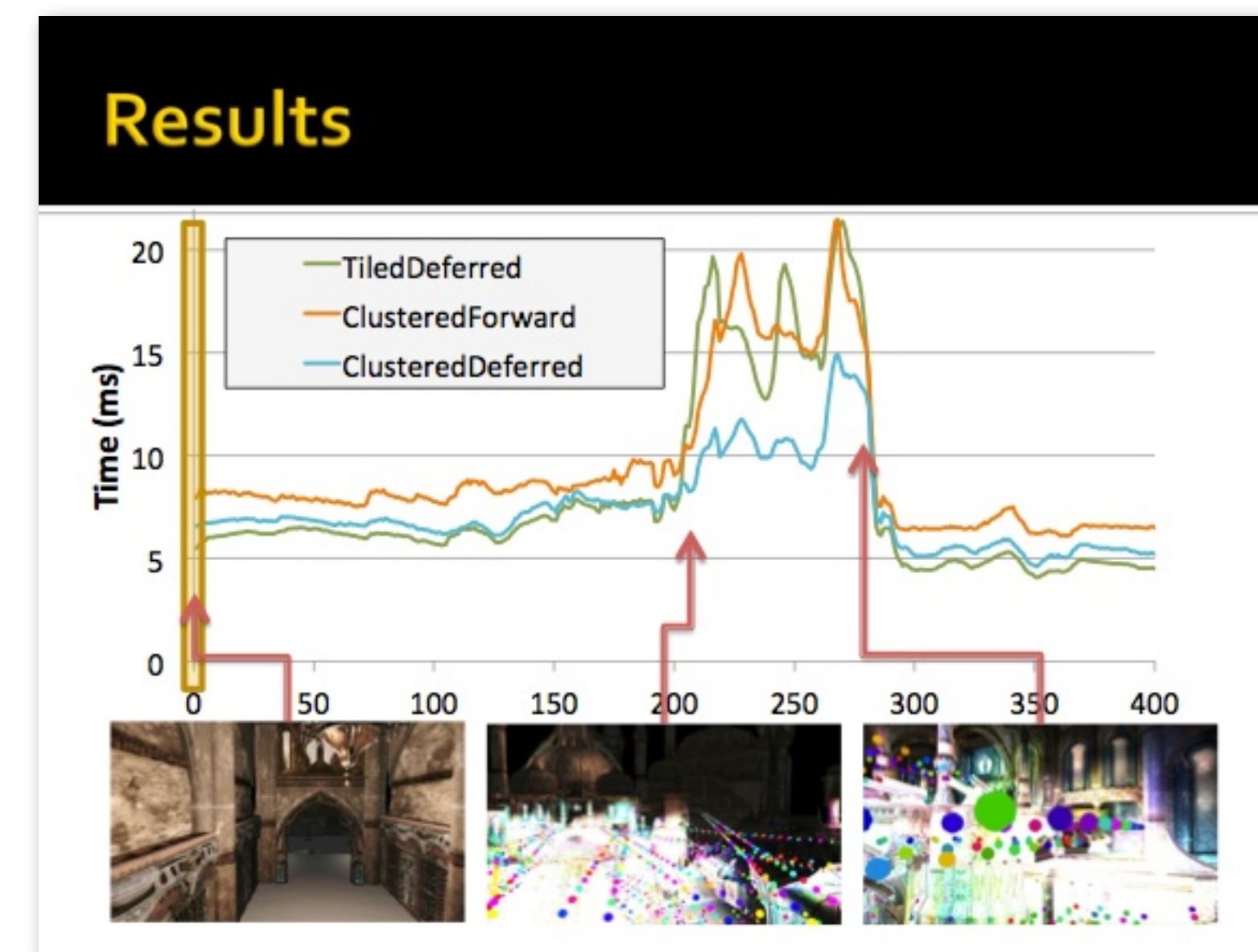
- Notice how you (as an audience member) are working hard to interpret the trends in these graphs
 - You are asking: what do these results say?
 - What am I supposed to be concluding?
- The audience just wants to be told what to look for!
 - They are reading the graphs to verify the main point, not determine the main point.

Simulation Results : RGS

● RGS Performance

- ❖ 147-198 Mray/sec
- ❖ Texture cache concerns : Mip-mapping & Compression

Test scene	Ray type	Cache hit rate (%)		Bandwidth (GB/s)	Performance (Mrays/sec)
		Texture	Data		
Sibenik (80K tri.)	Primary	-	96.76	0.5	182.11
	FSR	-	91.24	1.9	172.25
Fairy (179K tri.)	Primary	93.25	96.87	0.8	175.66
	FSR	81.49	94.91	1.9	147.45
Ferrari (210K tri.)	Primary	86.12	98.09	0.6	183.28
	FSR	75.95	95.71	2.0	163.67
Conference (282K tri.)	Primary	-	98.44	0.2	198.32
	FSR	-	95.72	0.8	158.79



Tip 10

Titles matter

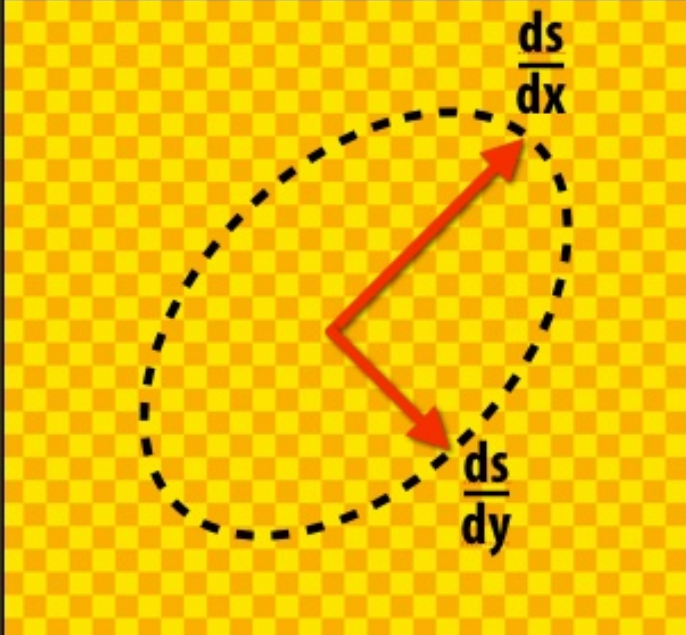
If you read the titles of your talk all the way through, it should be a great summary of the talk.

(basically, this is “one-point-per-slide” for the whole talk)

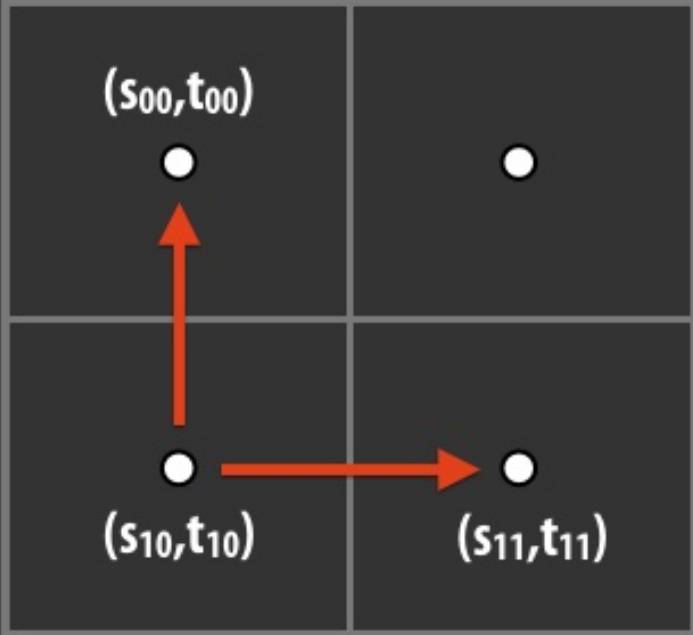
Examples of descriptive slide titles

GPUs shade quad fragments (2x2 pixel blocks)

Texture data



Quad fragment



use differences between neighboring texture coordinates to estimate derivatives

Greedy SRDH build optimizes over partitions and traversal policies

SAH:
`forall(partitions in set-of-partitions)
...evaluate SAH and pick min...`

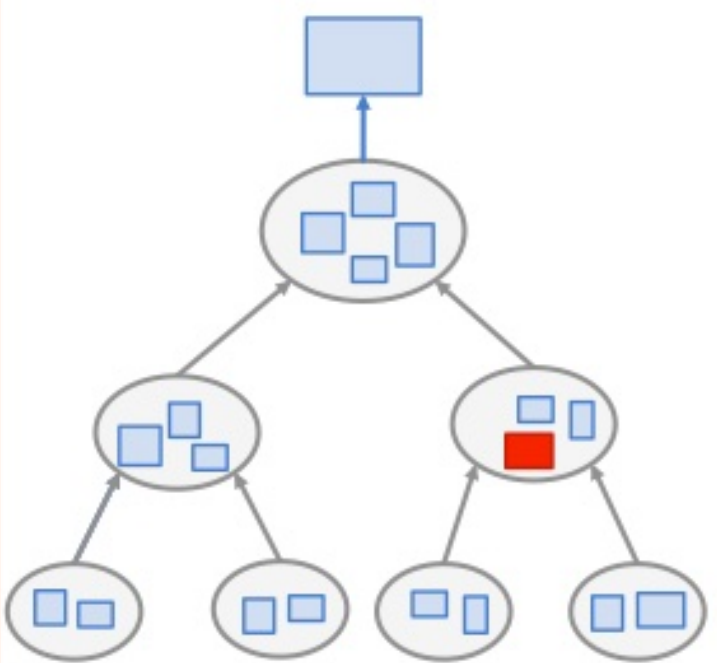
SRDH:
`forall(partitions in set-of-partitions)
forall(traversalKernels in set-of-kernels)
...evaluate SRDH and pick min...`

$$\text{SRDH}(R,L,\kappa,r) = (1 - \kappa(r)H(L,r))|R| + (1 - \kappa(r)H(R,r))|L|$$

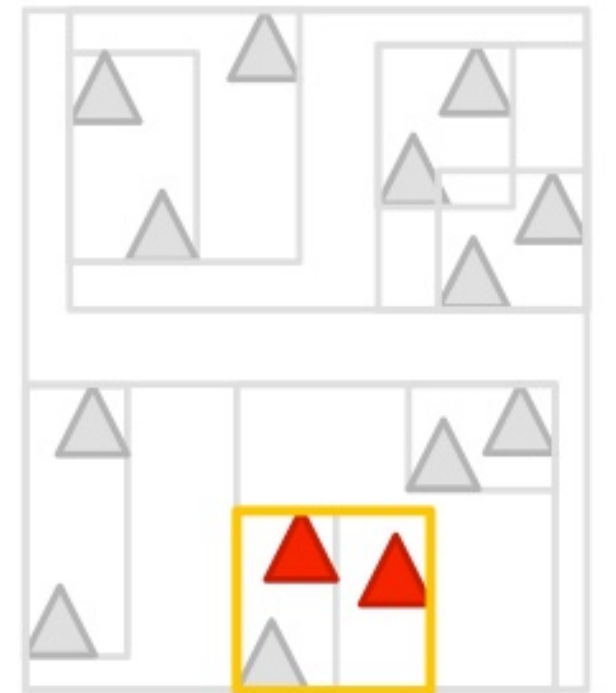
51

AAC IS AN APPROXIMATION TO THE TRUE AGGLOMERATIVE CLUSTERING SOLUTION.

Computation graph:



Primitive partitioning:



The reason for meaningful slide titles is convenience and clarity for the audience

“Why is the speaker telling me this again?”

(Recall “why before what”)

Read your slide titles in thumbnail view

Do they make all the points of the story you are trying to tell?

1 Reducing Shading on GPUs using Quad-Fragment Merging
Rayson Fatahalhan, Solomon Bealos, James Hegarty, Harry Morston, Kurt Akeley, Pat Hanrahan, William R. Mark

2 High-resolution meshes are appearing in games
Low detail

3 High-resolution meshes are appearing in games

4 PROBLEM
Current GPUs shade small triangles inefficiently

5 Multi-sample locations
Triangle coverage multiple times per pixel (for each aligned edge)

6 Shading sample locations
Triangle shading once per pixel

7 Surface derivatives are needed for texture filtering

8 GPUs shade quad fragments (2x2 pixel blocks)
Quad fragment

9 Shaded quad fragments

10 Final pixel values

11 Pixels at triangle boundaries are shaded multiple times
Shading computations per pixel

12 Pixels at triangle boundaries are shaded multiple times
Shading computations per pixel

13 Pixels at triangle boundaries are shaded multiple times
Shading computations per pixel

14 Small triangles result in extra shading
Unmerged small triangles, Merged small triangles, 7 pixel area triangles

15 Goal: Shade high-resolution meshes (not individual triangles) approximately once per pixel
Approach: Evolve GPU's quad-fragment shading system (Provide smooth evolution from status quo)

16 QUAD-FRAGMENT MERGING

17 GPU pipeline (with tessellation)

18 Rasterized quad-fragment

19 Rasterized quad-fragment
multi-sample coverage mask

20 Rasterized quad fragments
shader inputs (interpolated from triangle vertices)

21 GPU pipeline: triangle connectivity is known
Triangle connectivity is known

22 Pipeline with quad-fragment merging

23 Pipeline with quad-fragment merging

24 Two key merging operations
1. Identifying when quad fragments can be merged
2. Constructing a merged quad fragment

25 Merging quad fragments
Step 1: aggregate coverage

26 Merging quad fragments
Step 2: sample shading inputs

27 Merging quad fragments
Step 2: sample shading inputs

28 Two key merging operations
1. Identifying when quad fragments can be merged
2. Constructing a merged quad fragment

29 Challenge
Avoiding merges that introduce visual artifacts

30 Example: surface with a silhouette
anti-aliased silhouette

31 Naive merging results in aliasing

32 Avoid merging across discontinuities

33 Conditions required to merge quad fragments
1. Same screen location
2. Same sidedness (triangles front-facing or back-facing)
3. Source triangles are adjacent in the mesh

34 High-frequency geometric detail may cause aliasing
Our merging rules are designed for real-time performance

35 Implementation: the cost of merging is low
Merging operations are cheap
Merge buffer is small
Expectation: quad-fragment merging can be encapsulated in fixed-function hardware

36 EVALUATION

Tip 11

Practice the presentation

Practice the presentation

- **Given the time constraints, you'll need to be smooth to say everything you want to say**
- **To be smooth you'll have to practice**
- **Rehearse your presentation several times the night before (in front of a partner or friend)**
 - **It's only a short presentation, so a couple of practice runs are possible in a small amount of time**