

Meta-Reinforcement Learning: Learning to Explore

CS224R



Reminders

Homework 3 due **tonight**
(and HW4 out today)

Project milestone due **next Wednesday**

Following up on high-res feedback:

- Wanting homeworks to require more conceptual understanding
- Request for summary table of approaches
- Unofficial lecture notes

Why meta-reinforcement learning?

Why are humans good at RL?

Our RL agents start tabula rasa.



People have previous experience.

They have developed representations that facilitate exploration & learning.

Can we allow RL agents to leverage prior experience?

Should we be using the same exploration algorithm for:

- Learning to **navigate an environment**
- Learning to **make recommendations** to users
- Learning a policy for **computer system caching**
- Learning to **physically operate a new tool or machine**

This is how we currently approach exploration.

Today's Lecture

Can we *learn exploration strategies* based on experience from other tasks in that domain?

Outline

Brief Recap on Meta-RL

Algorithms for Learning to Explore

End-to-End Optimization of Exploration Strategies

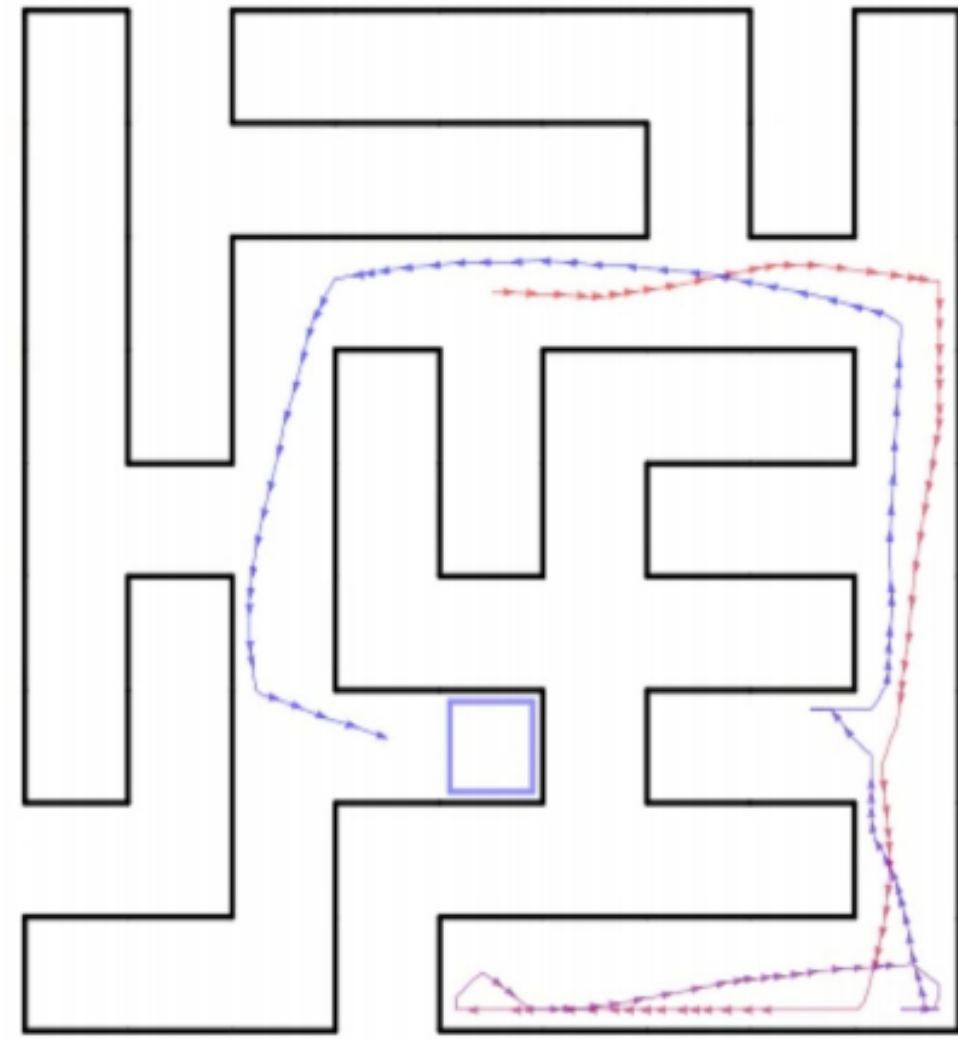
Alternative Decoupled Exploration Strategies

Decoupled but Consistent Exploration & Exploitation

Case Study: Applying Meta-RL to CS Education

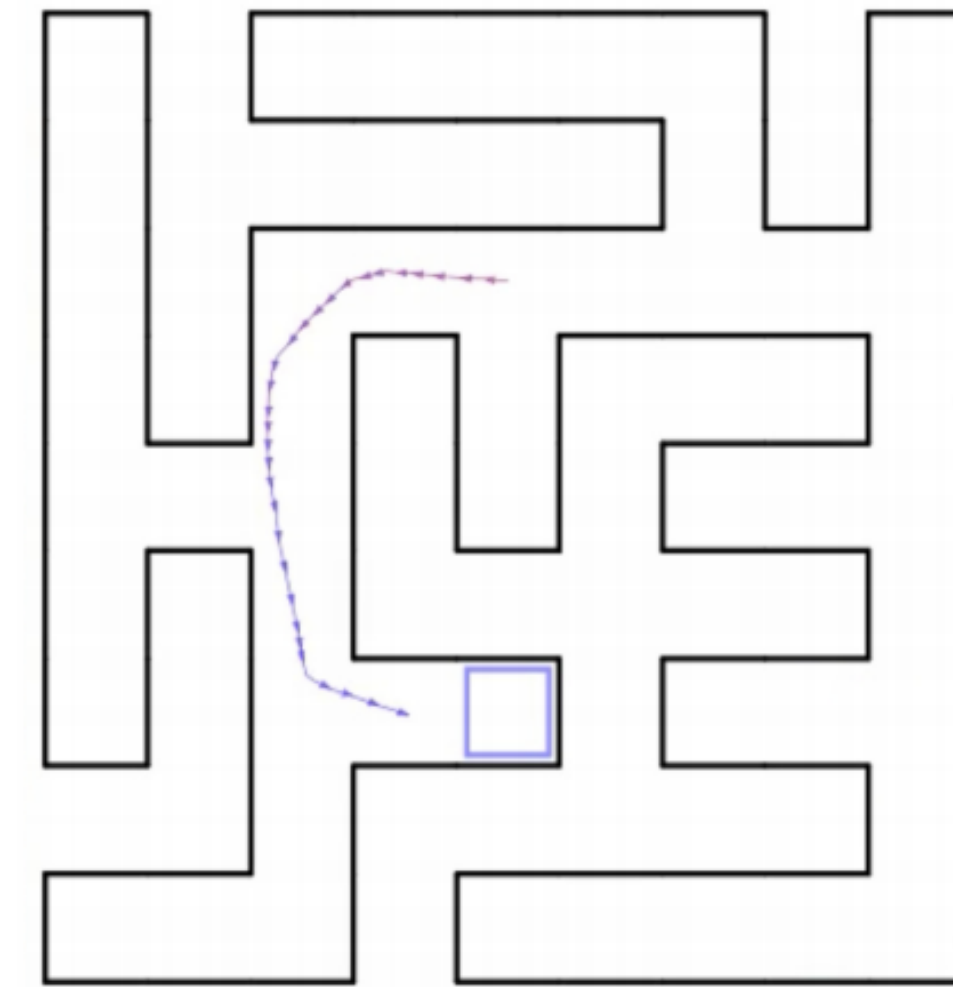
A brief recap of meta-reinforcement learning

Collect small amount of
experience in new MDP



Collect $\mathcal{D}_{\text{tr}} \sim \pi^{\text{exp}}$

Learn policy that
solves that MDP



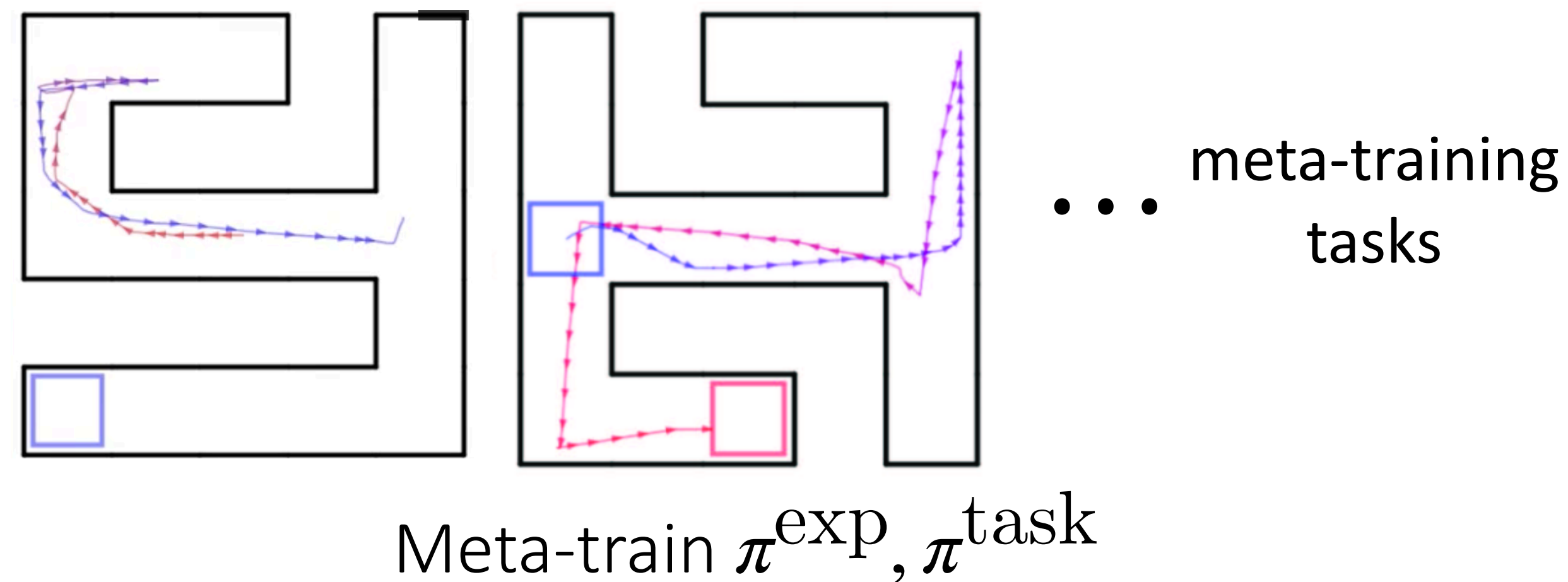
$\mathcal{D}_{\text{tr}} \rightarrow \pi^{\text{task}}$

Goal:

A brief recap of meta-reinforcement learning

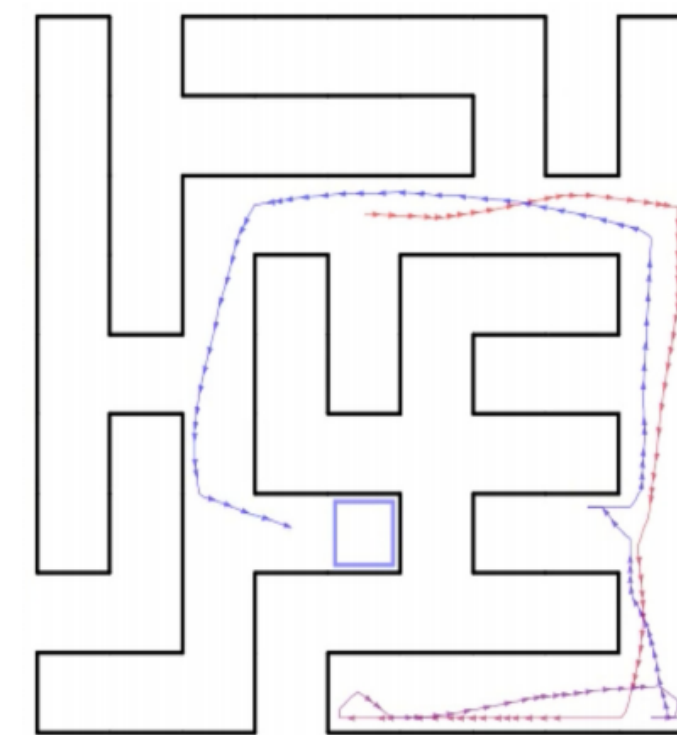
Meta-Train Time:

Learn how to efficiently explore & solve many MDPs:

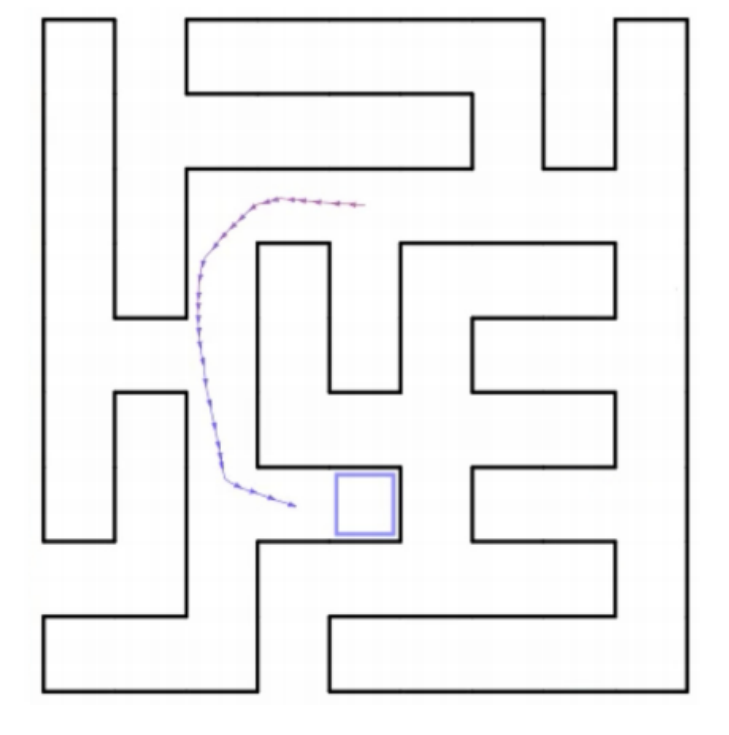


Meta-Test Time:

Collect small amount of experience in new MDP



Learn policy that solves that MDP

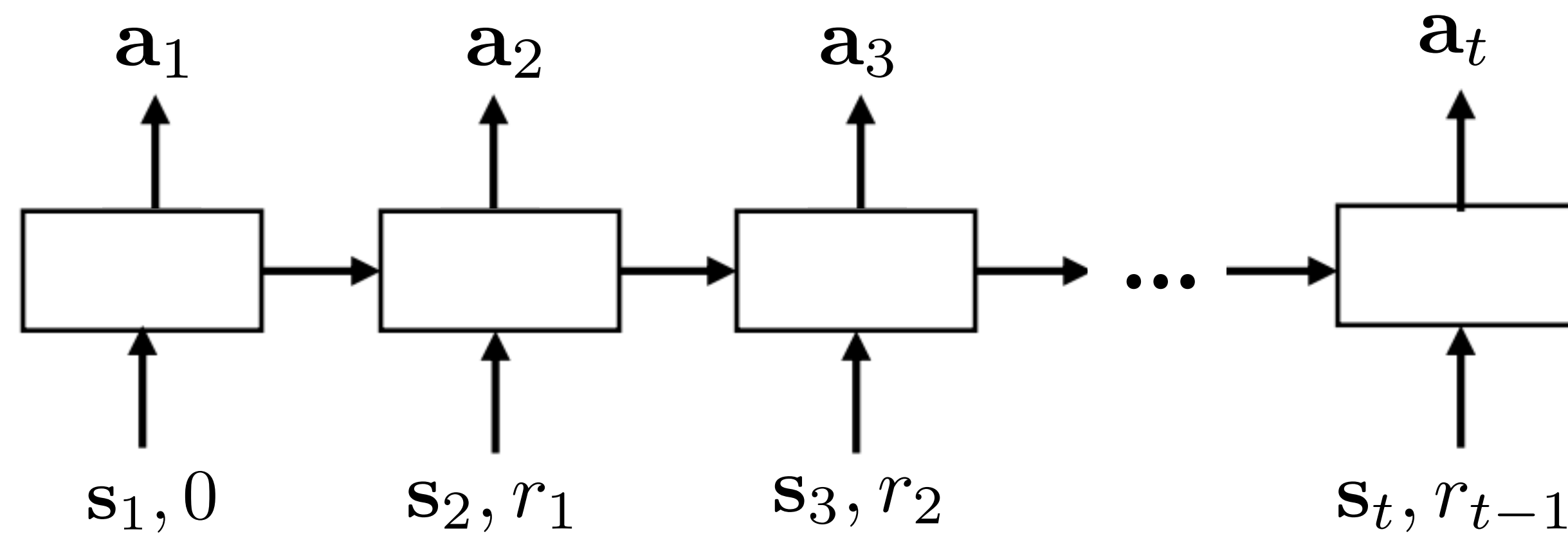


Key assumption: Meta-training & meta-testing MDPs come from same distribution.

(so that we can expect generalization)

A brief recap of meta-reinforcement learning

Common approach: Implement the learning procedure with a recurrent network.



Is this just a recurrent policy?

Hidden state maintained
across episodes within a task!

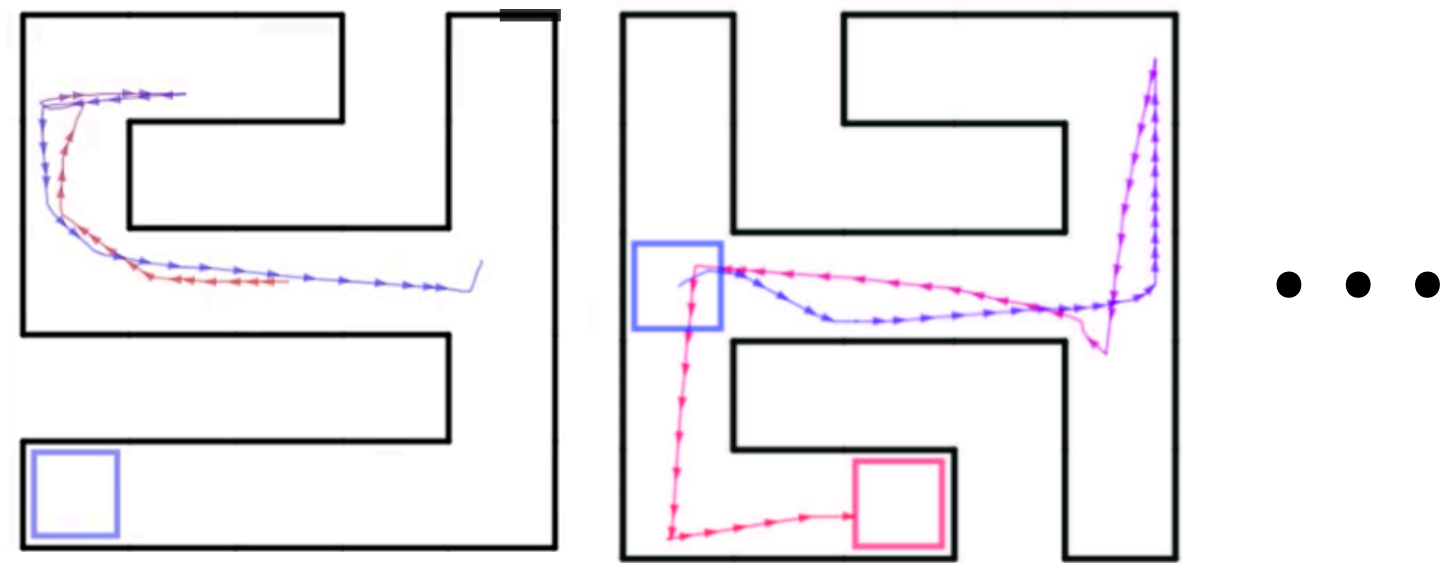
Trained across a *family of MDPs*
with varying dynamics, rewards.

RL² with Policy Gradients:
$$\nabla_{\theta} J(\theta) = E_{\tau \sim \pi_{\theta}, \mathcal{T}_i} \left[\left(\sum_t \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t, \mathcal{D}_i^{\text{tr}}) \right) \left(\sum_t r_i(\mathbf{s}_t, \mathbf{a}_t) \right) \right]$$

A brief recap of meta-reinforcement learning

Examples of meta-RL tasks

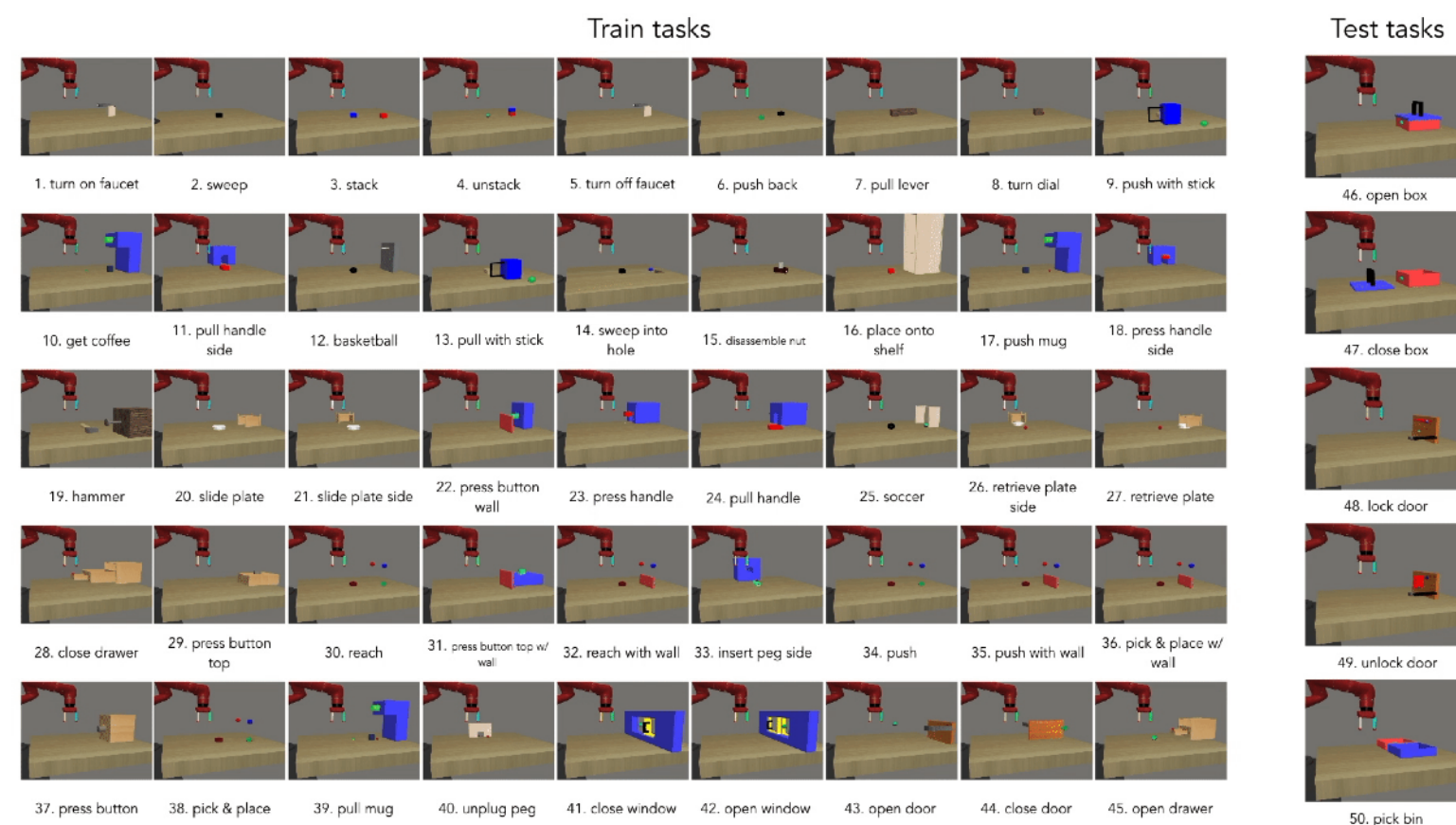
Navigation through different mazes



Locomotion on different terrains, slopes



Object manipulation with different objects, goals



Dialog with different users w/ different preferences



Outline

Brief Recap on Meta-RL

Algorithms for Learning to Explore

End-to-End Optimization of Exploration Strategies

Alternative Decoupled Exploration Strategies

Decoupled but Consistent Exploration & Exploitation

Case Study: Applying Meta-RL to CS Education

How Do We Learn to Explore?

Solution #1: Optimize for Exploration & Exploitation *End-to-End* w.r.t. Reward

(Duan et al., 2016, Wang et al., 2016, Mishra et al., 2017, Stadie et al., 2018, Zintgraf et al., 2019, Kamienny et al., 2020)

- + simple
- + leads to optimal strategy in principle
- challenging optimization when exploration is hard

A simple, running example

Hallway 1

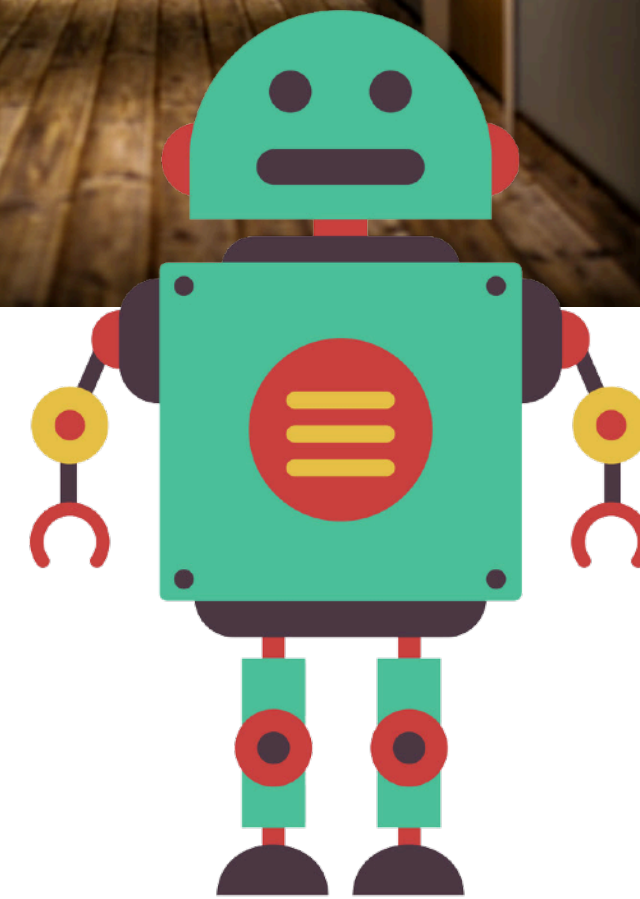


Hallway 2



...

Hallway N



agent



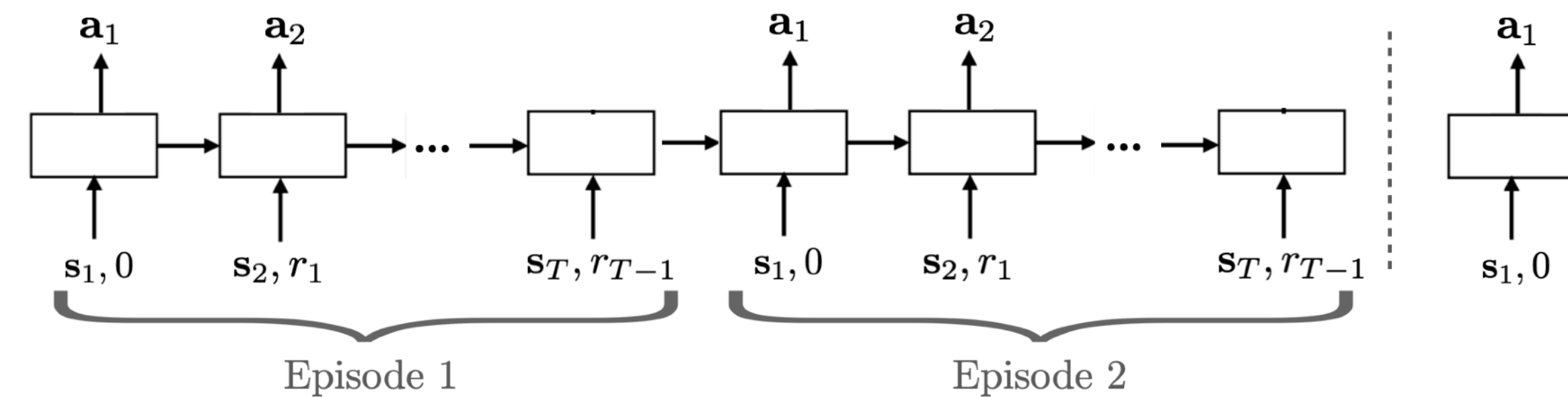
information on
where to go

Different tasks: navigating to
the ends of different hallways

How Do We Learn to Explore?

Solution #1: Optimize for Exploration & Exploitation *End-to-End* w.r.t. Task Reward

(Duan et al., 2016, Wang et al., 2016, Mishra et al., 2017, Stadie et al., 2018, Zintgraf et al., 2019, Kamienny et al., 2020)



Example episodes during meta-training:

agent goes to the end of the correct hallway

agent goes to wrong hallway then correct hallway

agent looks at the instructions

- gets positive reward for current task,
but $\mathcal{D}_i^{\text{tr}}$ won't be different than for any other task

+/- provides signal on a **suboptimal**
exploration + exploitation strategy

- good exploratory behavior, but won't
get any reward for this behavior

It's hard to learn exploration & exploitation at the same time!

Another Example of a Hard Exploration Meta-RL Problem

Learned cooking tasks in previous kitchens



meta-training

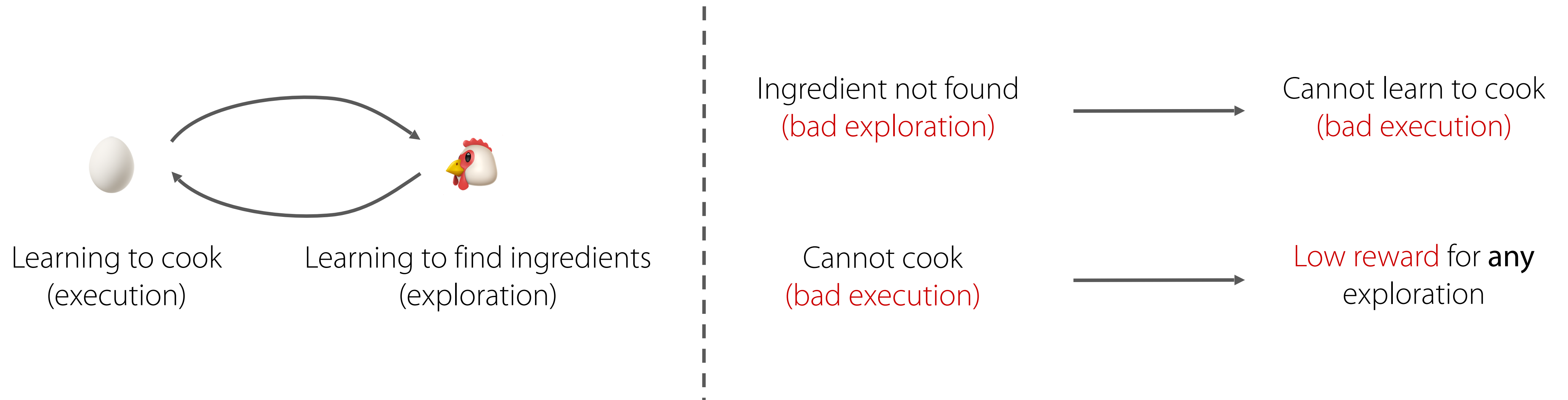
Goal: Quickly learn tasks in a new kitchen.



meta-testing

Why is End-to-End Training Hard in This Example?

End-to-end approach: optimize exploration and execution episode behaviors end-to-end to maximize reward of execution



Coupling problem: learning exploration and execution depend on each other

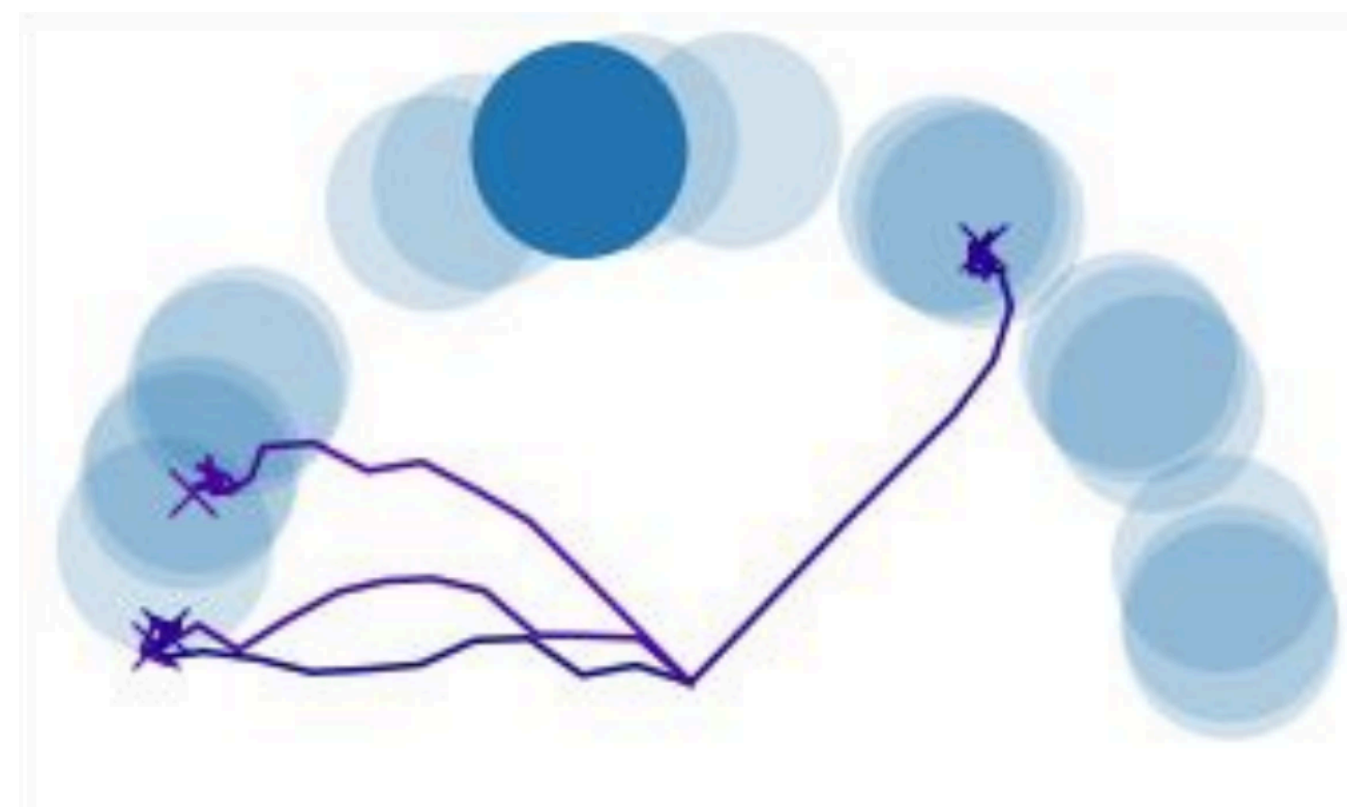
—> can lead to poor local optima, poor sample efficiency

Solution #2: Leverage Alternative Exploration Strategies

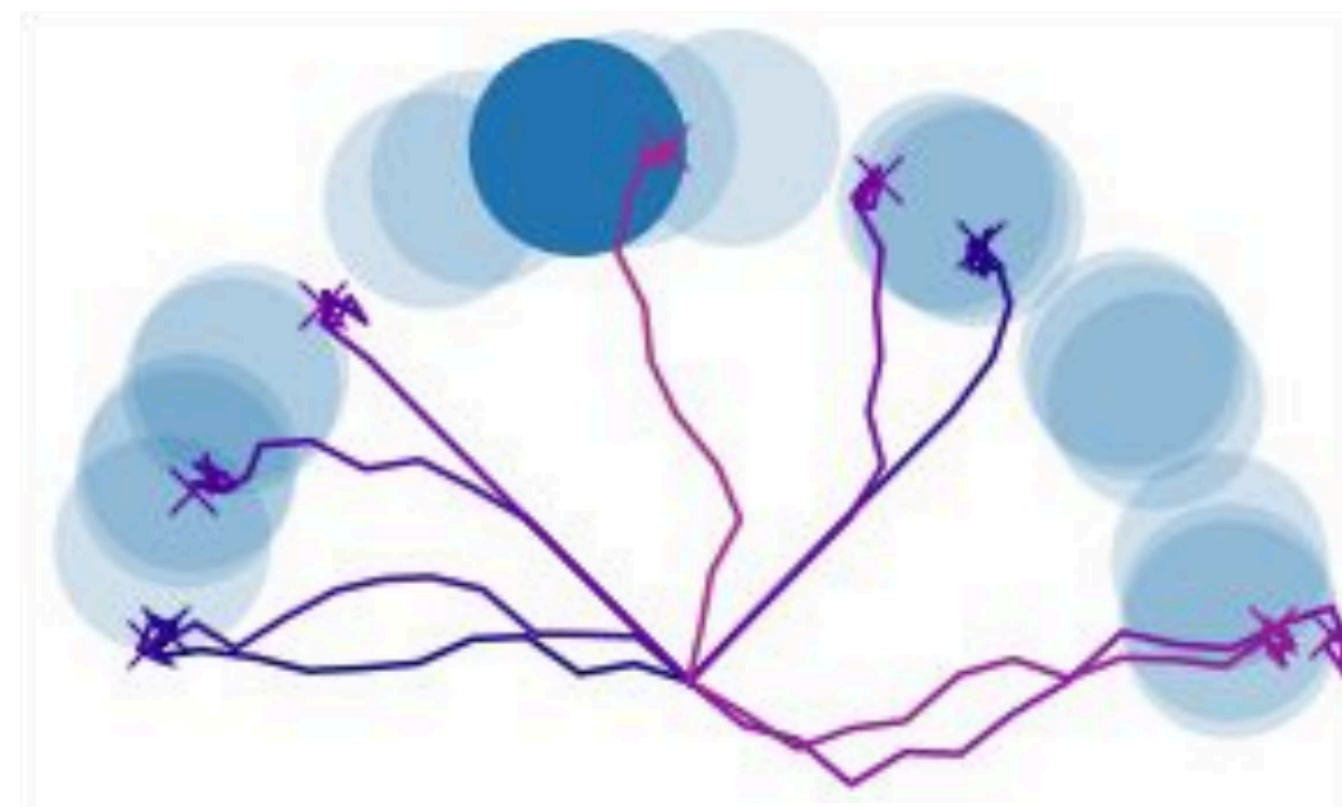
2a. Use posterior sampling
(also called Thompson sampling)

PEARL (Rakelly, Zhou, Quillen, Finn, Levine. ICML '19)

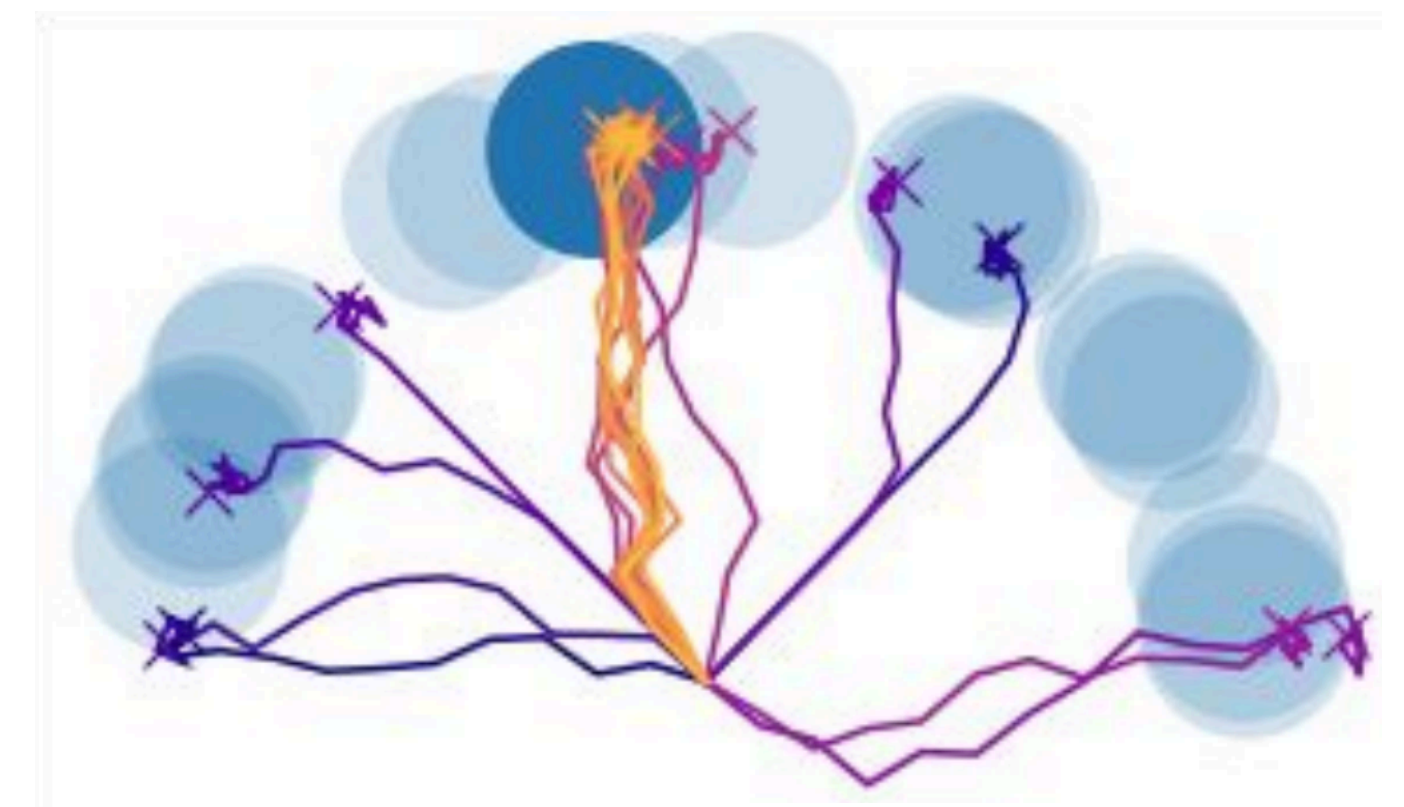
- i. Learn distribution over latent task variable $p(\mathbf{z})$, $q(\mathbf{z} | \mathcal{D}_{\text{tr}})$ and corresponding task policies $\pi(\mathbf{a} | \mathbf{s}, \mathbf{z})$
- ii. Sample \mathbf{z} from current *posterior* and sample from policy $\pi(\mathbf{a} | \mathbf{s}, \mathbf{z})$



$$\mathbf{z} \sim p(\mathbf{z})$$



$$\mathbf{z} \sim q_{\phi}(\mathbf{z} | c_{1:10})$$



$$\mathbf{z} \sim q_{\phi}(\mathbf{z} | c_{1:30})$$

When might posterior sampling be bad? Eg. Goals far away & sign on wall that tells you the correct goal.

Solution #2: Leverage Alternative Exploration Strategies

2a. Use posterior sampling
(also called Thompson sampling)

PEARL (Rakelly, Zhou, Quillen, Finn, Levine. ICML '19)

- i. Learn distribution over latent task variable $p(\mathbf{z}), q(\mathbf{z} | \mathcal{D}_{\text{tr}})$ and corresponding task policies $\pi(\mathbf{a} | \mathbf{s}, \mathbf{z})$
- ii. Sample \mathbf{z} from current *posterior* and sample from policy $\pi(\mathbf{a} | \mathbf{s}, \mathbf{z})$

2b. Use intrinsic rewards

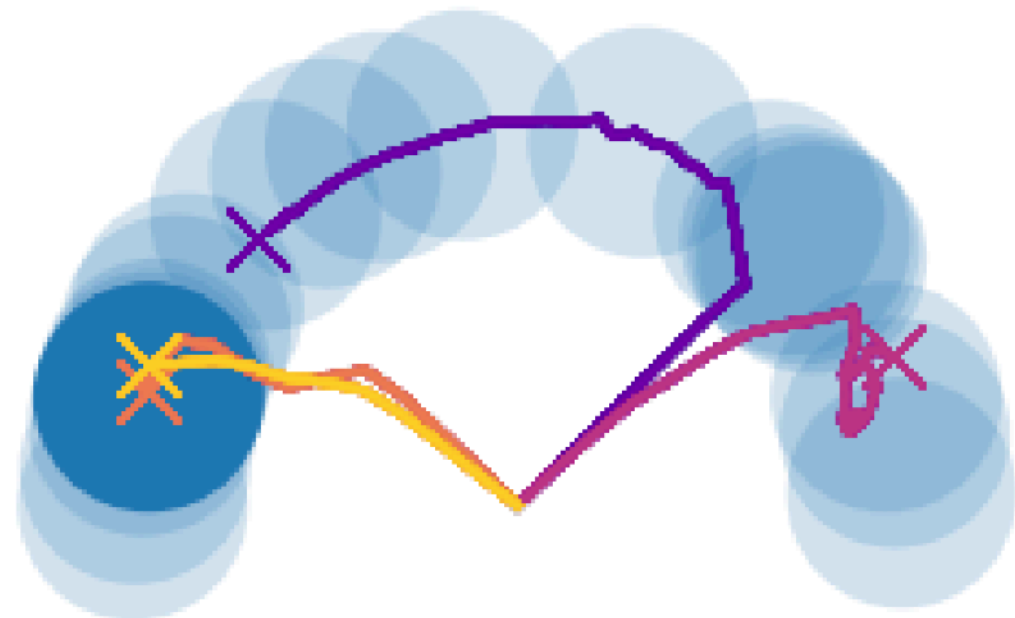
MAME (Gurumurthy, Kumar, Sycara. CoRL '19)

2c. Task dynamics & reward prediction

MetaCURE (Zhang, Wang, Hu, Chen, Fan, Zhang. '20)

- i. Train model $f(\mathbf{s}', r | \mathbf{s}, \mathbf{a}, \mathcal{D}_{\text{train}})$

- ii. Collect $\mathcal{D}_{\text{train}}$ so that model is accurate.



When might this be bad?

Lots of distractors,
or complex, high-dim state dynamics

Solution #2: Leverage Alternative Exploration Strategies

- 2a. Use posterior sampling (also called Thompson sampling) PEARL (Rakelly, Zhou, Quillen, Finn, Levine. ICML '19)
- Learn distribution over latent task variable $p(\mathbf{z}), q(\mathbf{z} | \mathcal{D}_{\text{tr}})$ and corresponding task policies $\pi(\mathbf{a} | \mathbf{s}, \mathbf{z})$
 - Sample \mathbf{z} from current *posterior* and sample from policy $\pi(\mathbf{a} | \mathbf{s}, \mathbf{z})$
- 2b. Use intrinsic rewards MAME (Gurumurthy, Kumar, Sycara. CoRL '19)
- 2c. Task dynamics & reward prediction MetaCURE (Zhang, Wang, Hu, Chen, Fan, Zhang. '20)
- Train model $f(\mathbf{s}', r | \mathbf{s}, \mathbf{a}, \mathcal{D}_{\text{train}})$
 - Collect $\mathcal{D}_{\text{train}}$ so that model is accurate.

- + easy to optimize
- + many based on principled strategies

- suboptimal by arbitrarily large amount in some environments.

Can we avoid the chicken-and-egg problem without sacrificing optimality?

(best of both worlds?)

Yes!

Solution #3

Idea from solution #2b: Train model $f(\mathbf{s}', r | \mathbf{s}, \mathbf{a}, \mathcal{D}_{\text{tr}})$ & collect \mathcal{D}_{tr} so that model is accurate.

Do we have to learn a *full dynamics & reward model*?

Idea 3.0: Label each training task with a **unique ID μ**

Meta training

Exploration policy: train policy $\pi^{\text{exp}}(\mathbf{a} | \mathbf{s})$ and task identification model $q(\mu | \mathcal{D}_{\text{tr}})$
such that $\mathcal{D}_{\text{tr}} \sim \pi^{\text{exp}}$ allows accurate task prediction from f

Execution policy: train ID-conditioned policy $\pi^{\text{exec}}(\mathbf{a} | \mathbf{s}, \mu_i)$

Meta testing

Explore: $\mathcal{D}_{\text{tr}} \sim \pi^{\text{exp}}(\mathbf{a} | \mathbf{s})$ Infer task: $\hat{\mu} \sim q(\mu | \mathcal{D}_{\text{tr}})$ Perform task: $\pi^{\text{exec}}(\mathbf{a} | \mathbf{s}, \hat{\mu})$

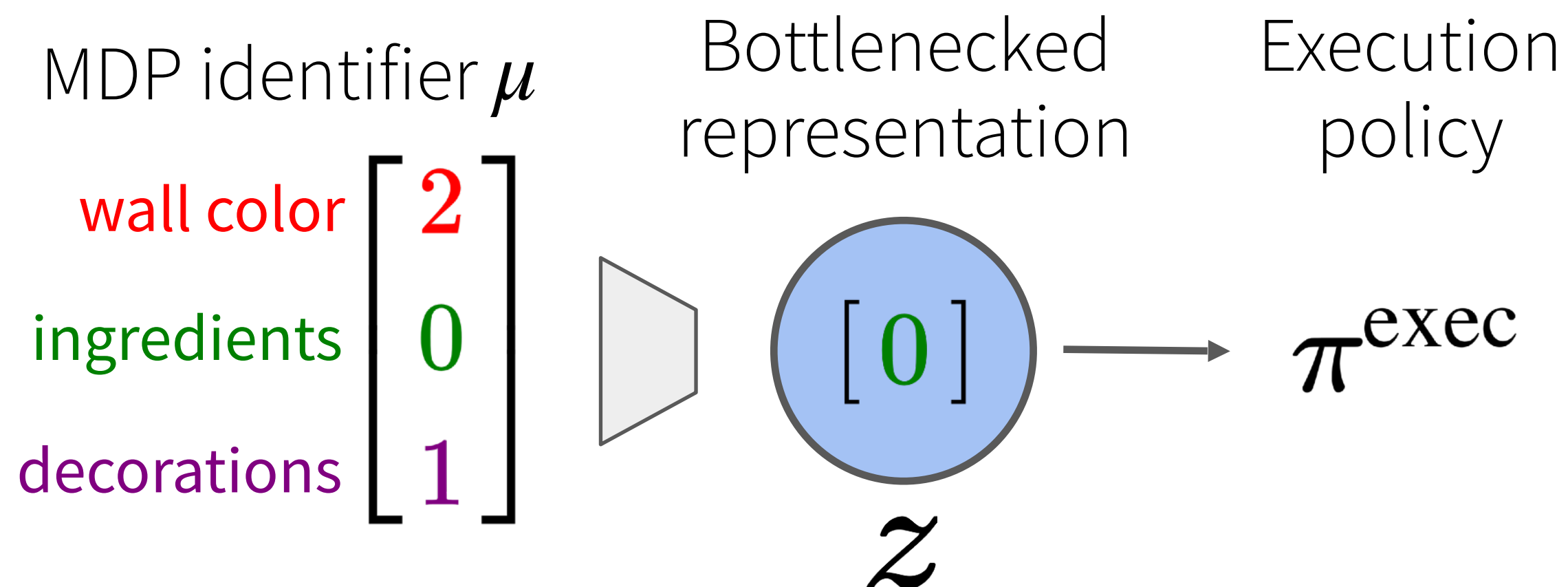
+ no longer need to model dynamics, rewards

— may not generalize well for one-hot μ

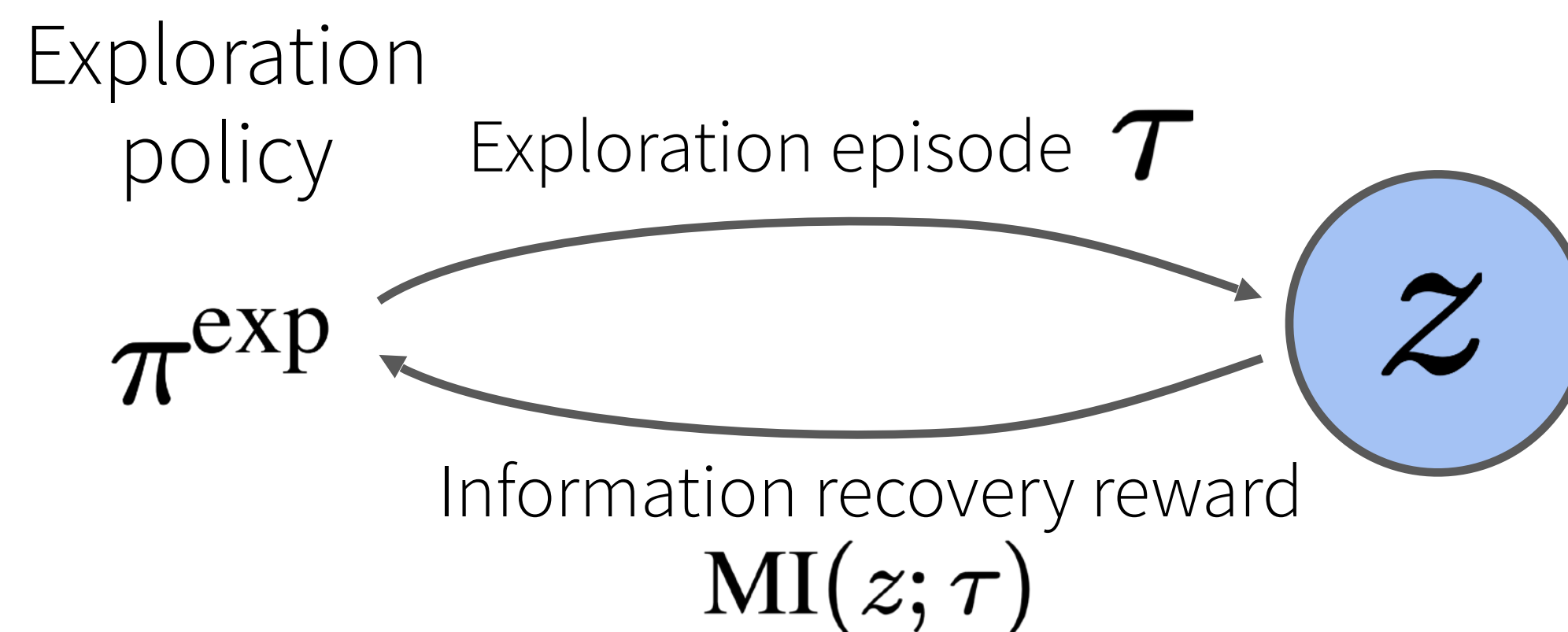
Solution #3: **Decouple** by acquiring representation of task relevant information

Meta-training

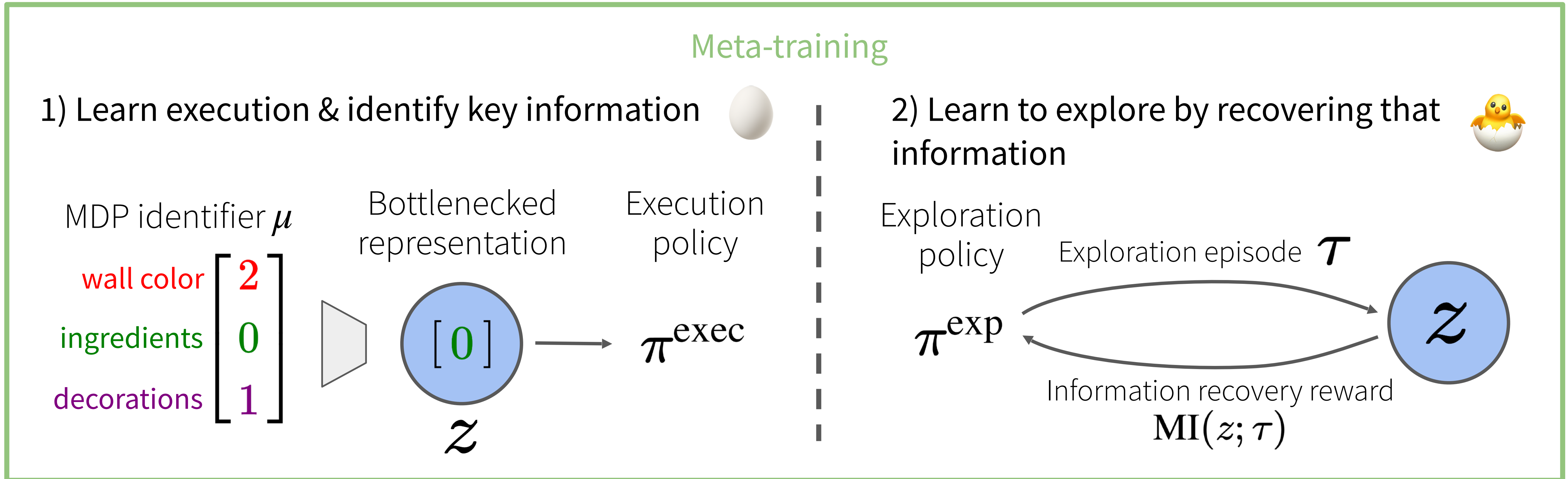
1) Learn execution & identify key information 



2) Learn to explore by recovering that information 



Solution #3: **Decouple** by acquiring representation of task relevant information



Train $\pi^{\text{exec}}(\mathbf{a} | \mathbf{s}, z_i)$ and encoder $F(z_i | \mu_i)$ to:

$$\max \sum_i \mathbb{E}_{\pi^{\text{exec}}} [r_i] - D_{\text{KL}} (F(z_i | \mu_i) || \mathcal{N}(0, 1))$$

Train π^{exp} such that collected \mathcal{D}_{tr} is predictive of z_i .

In practice: (1) and (2) can be trained simultaneously.

Solution #3: **Decouple** by acquiring representation of task relevant information

Meta-training

1) Learn execution & identify key information 

2) Learn to explore by recovering that information 

Train $\pi^{\text{exec}}(\mathbf{a} | \mathbf{s}, z_i)$ and encoder $F(z_i | \mu_i)$ to:

$$\max \sum_i \mathbb{E}_{\pi^{\text{exec}}} [r_i] - D_{\text{KL}} (F(z_i | \mu_i) || \mathcal{N}(0, 1))$$

Train π^{exp} such that collected \mathcal{D}_{tr} is predictive of z_i .

How to formulate the *reward function* for π^{exp} ?

(a) Train model $q(z_i | \mathcal{D}_{\text{tr}})$ (b) $r_t =$ per-step information gain

$r_t =$ prediction error from $\tau_{1:t-1}$ — prediction error from $\tau_{1:t}$

Decoupled Reward-free Exploration and Execution in Meta-Reinforcement Learning (DREAM)

Aside: How can we bottleneck the information in a neural net's representation?

V0: Add noise the representation.

$$\epsilon \sim \mathcal{N}(0, I) \quad \bar{\mathbf{z}} = \mathbf{z} + \epsilon \quad + \text{ will discard information } \text{😊}$$

- if done at test time, my discard good info

- if done during training, model can increase magnitude of \mathbf{z}

- Key ideas:
1. Add Gaussian noise during training
 2. Prevent the model from increasing magnitude

V1: Variational information bottleneck

Add noise before passing representation to next layer: $\epsilon \sim \mathcal{N}(0, I) \quad \bar{\mathbf{z}} = \mathbf{z} + \epsilon$ Modify loss term:

$$L_{\text{tr}} + \|\mathbf{z}\|^2$$

-> equivalent to $D_{KL} (F(\mathbf{z} | \mu_i) \| \mathcal{N}(0, 1))$.

Solution #3: **Decouple** by acquiring representation of task relevant information

(Informal) Theoretical Analysis

(1) **DREAM** objective is **consistent** with **end-to-end optimization**.

[under mild assumptions]

-> can in principle recover the optimal exploration strategy

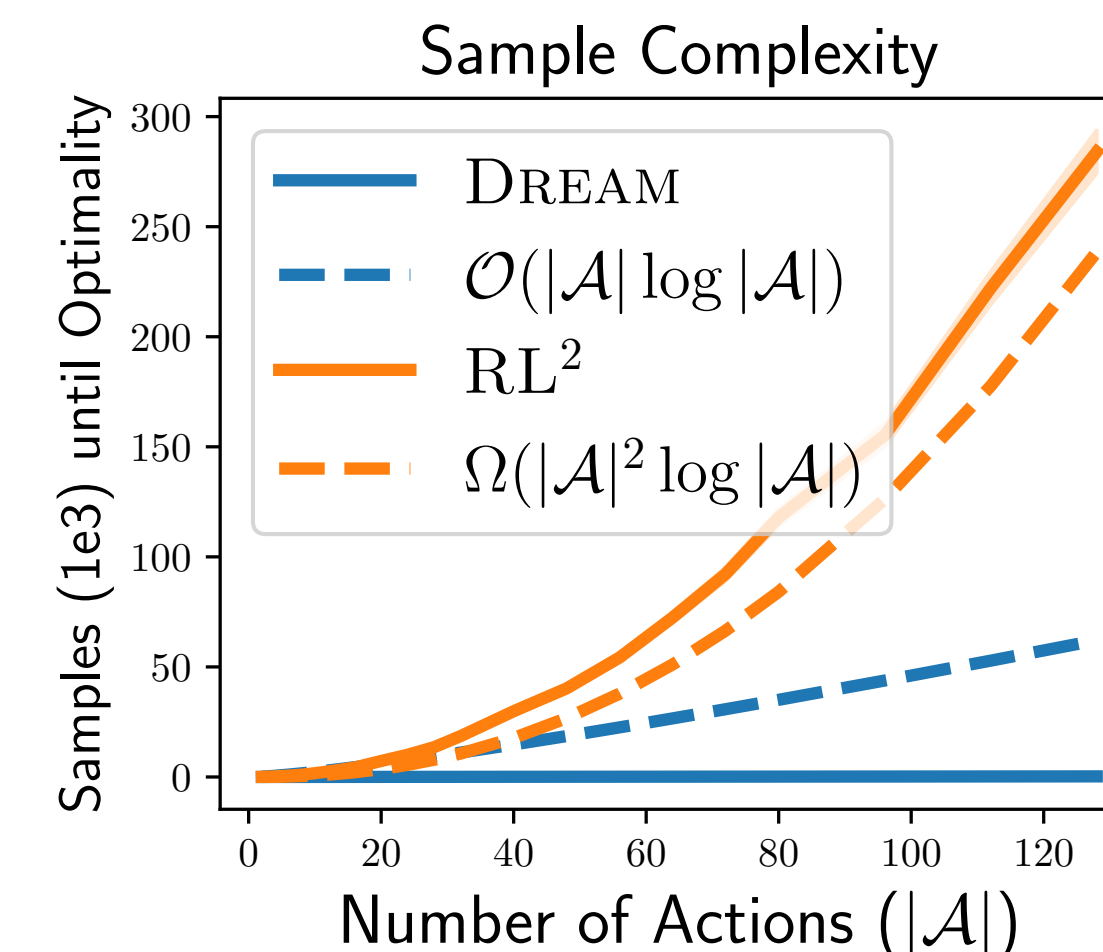
(2) Consider a bandit-like setting with $|\mathcal{A}|$ arms.

In MDP i , arm i yields reward. In all MDPs, arm 0 reveals the rewarding arm.

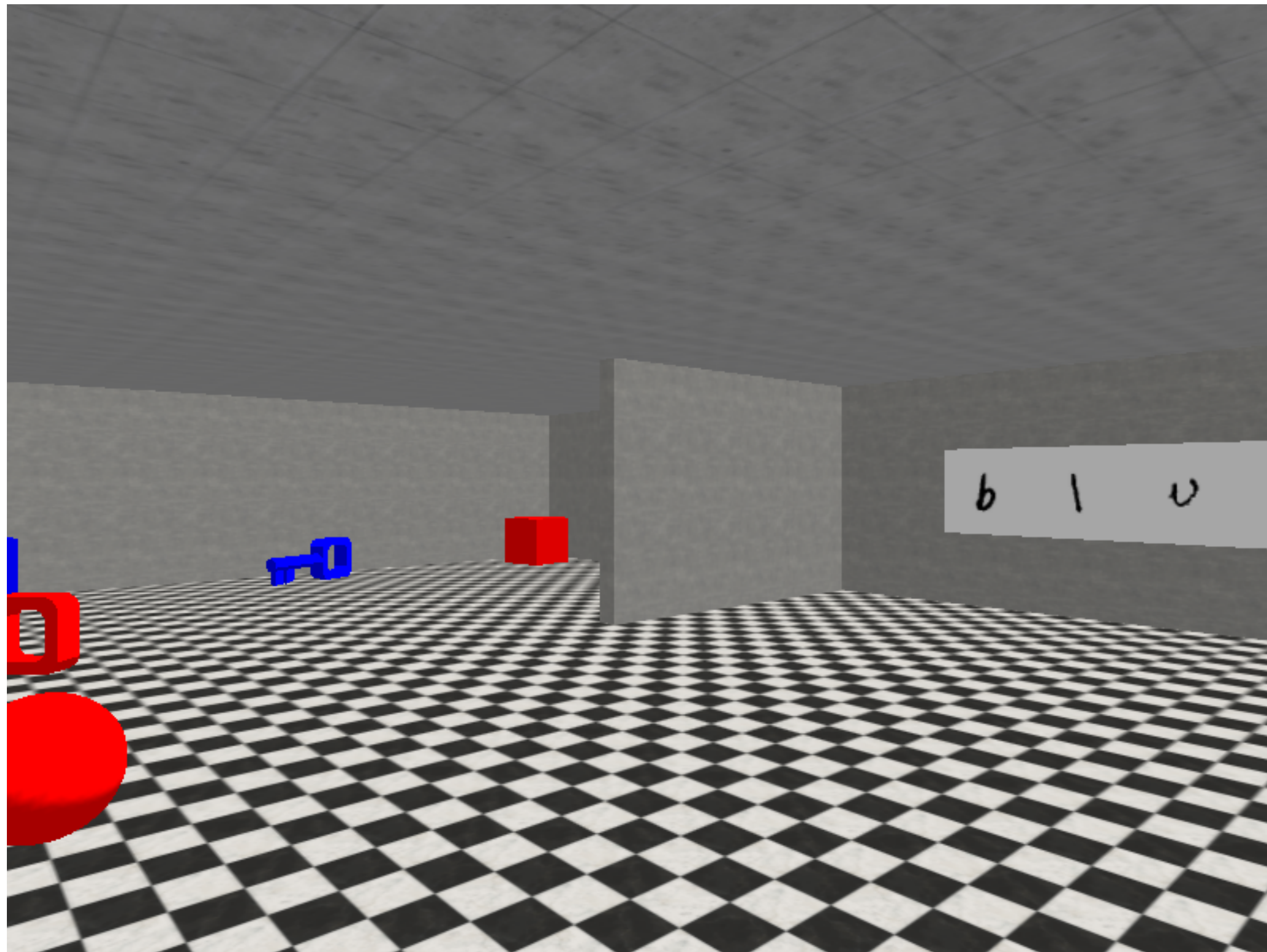
RL² requires $\Omega(|\mathcal{A}|^2 \log |\mathcal{A}|)$ samples for meta-optimization.

DREAM requires $\mathcal{O}(|\mathcal{A}| \log |\mathcal{A}|)$ samples for meta-optimization.

[assuming Q-learning with uniform outer-loop exploration]



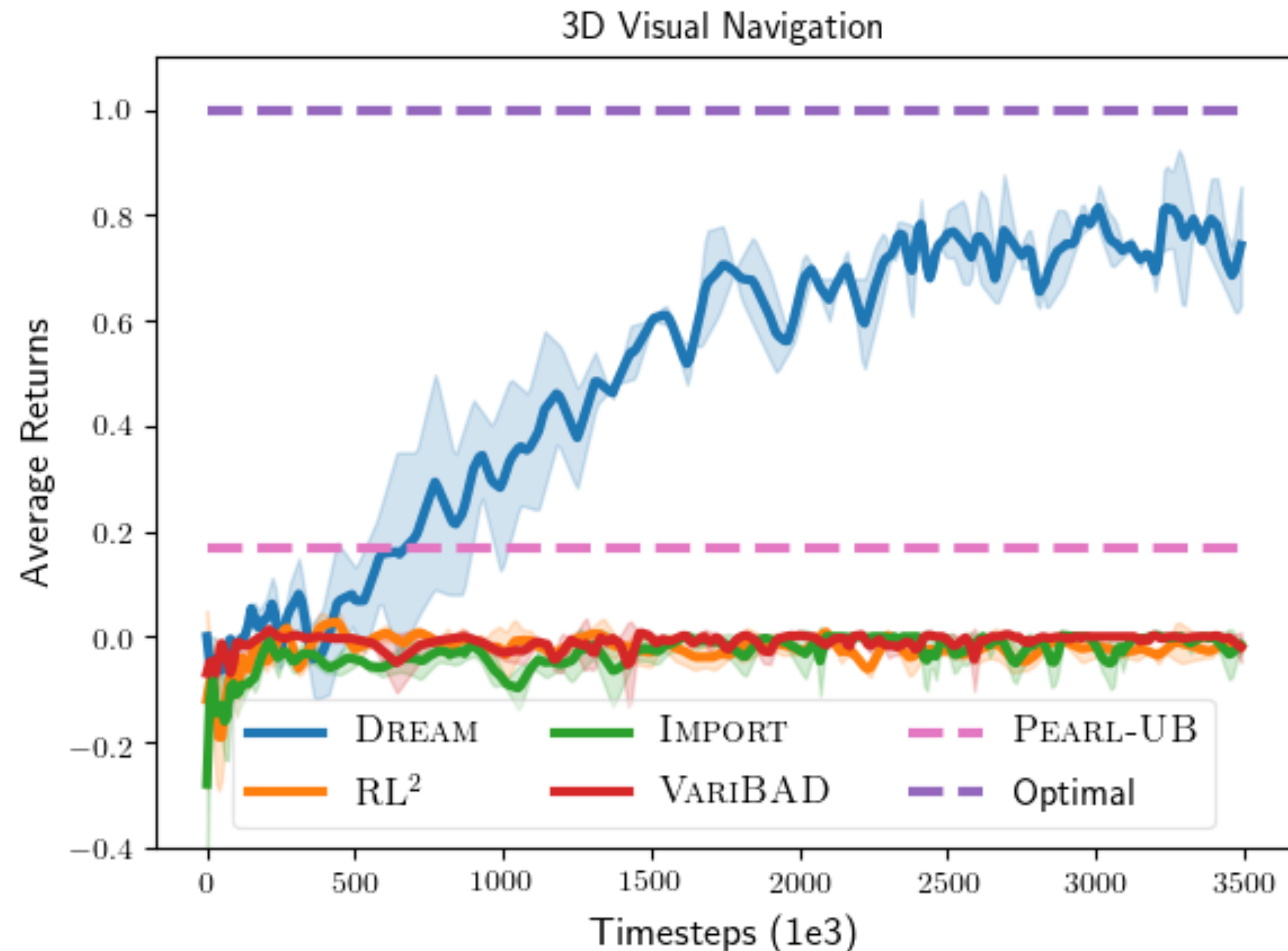
Empirical Comparison: Sparse Reward 3D Visual Navigation Problem



More challenging variant of task from Kamienny et al., 2020

- Task: go to the (key / block / ball), color specified by the sign
- Agent starts on other side of barrier, must walk around to read the sign
- Pixels observations (80 x 60 RGB)
- Sparse binary reward

Quantitative Comparison

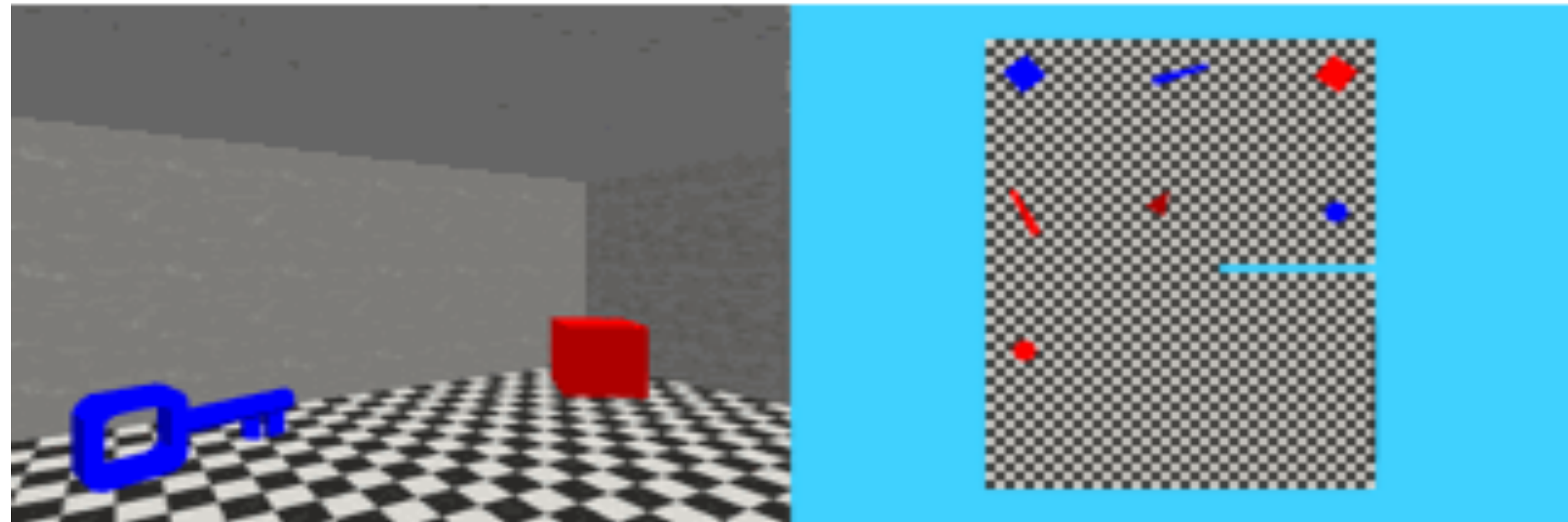


RL² (Duan et al., 2016), IMPORT (Kamienny et al., 2020), VARIBAD (Zintgraf et al., 2019), PEARL (Rakelly, et. al., 2019), Thompson, 1933

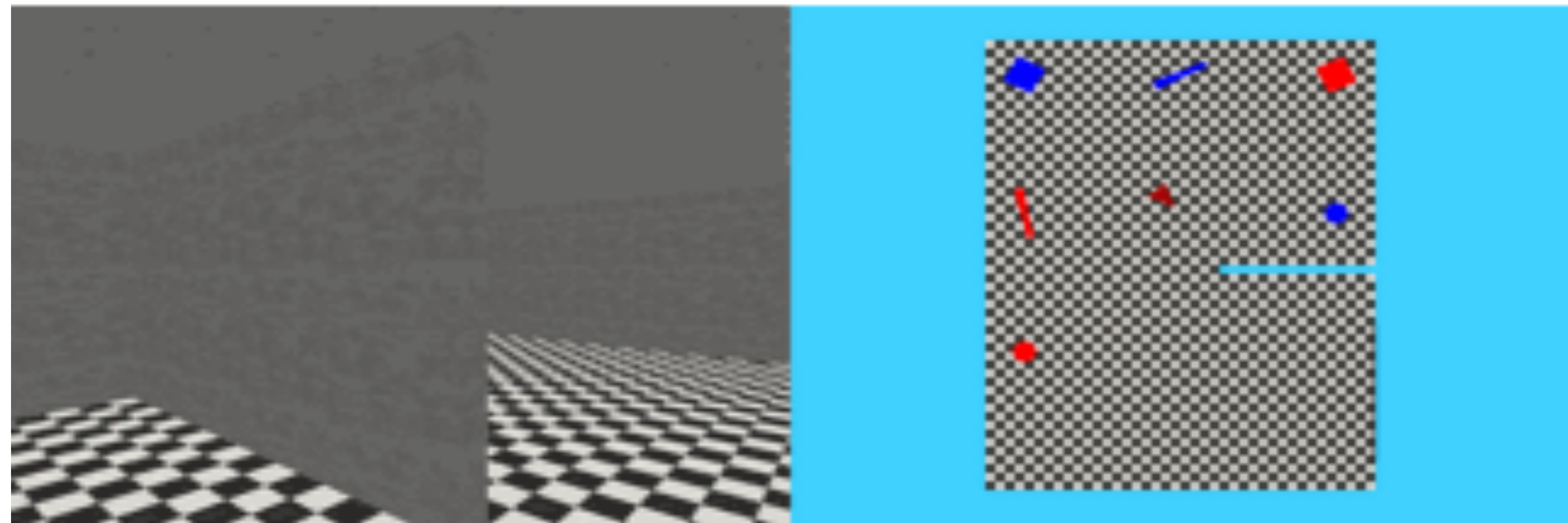
- End-to-end algorithms (RL², IMPORT, VARIBAD) perform poorly due to **coupling**
- PEARL-UB: Upper-bound on PEARL: optimal policy and Thompson-Sampling exploration, does not learn the optimal exploration strategy
- DREAM achieves near-optimal reward

Qualitative Results for DREAM

Exploration episode



Execution episode
Goal: Go to key



How Do We Learn to Explore?

End-to-End

Alternative Strategies

Decoupled Exploration & Execution

+ leads to optimal strategy in principle

+ easy to optimize

+ leads to optimal strategy in principle

+ many based on principled strategies

+ easy to optimize in practice

-- challenging optimization when exploration is hard

-- suboptimal by arbitrarily large amount in some environments.

-- requires task identifier

Outline

Brief Primer on Meta-RL

Algorithms for Learning to Explore

End-to-End Optimization of Exploration Strategies

Alternative Decoupled Exploration Strategies

Decoupled but Consistent Exploration & Exploitation

Case Study: Applying Meta-RL to CS Education

Problem: Providing Feedback on Interactive Software

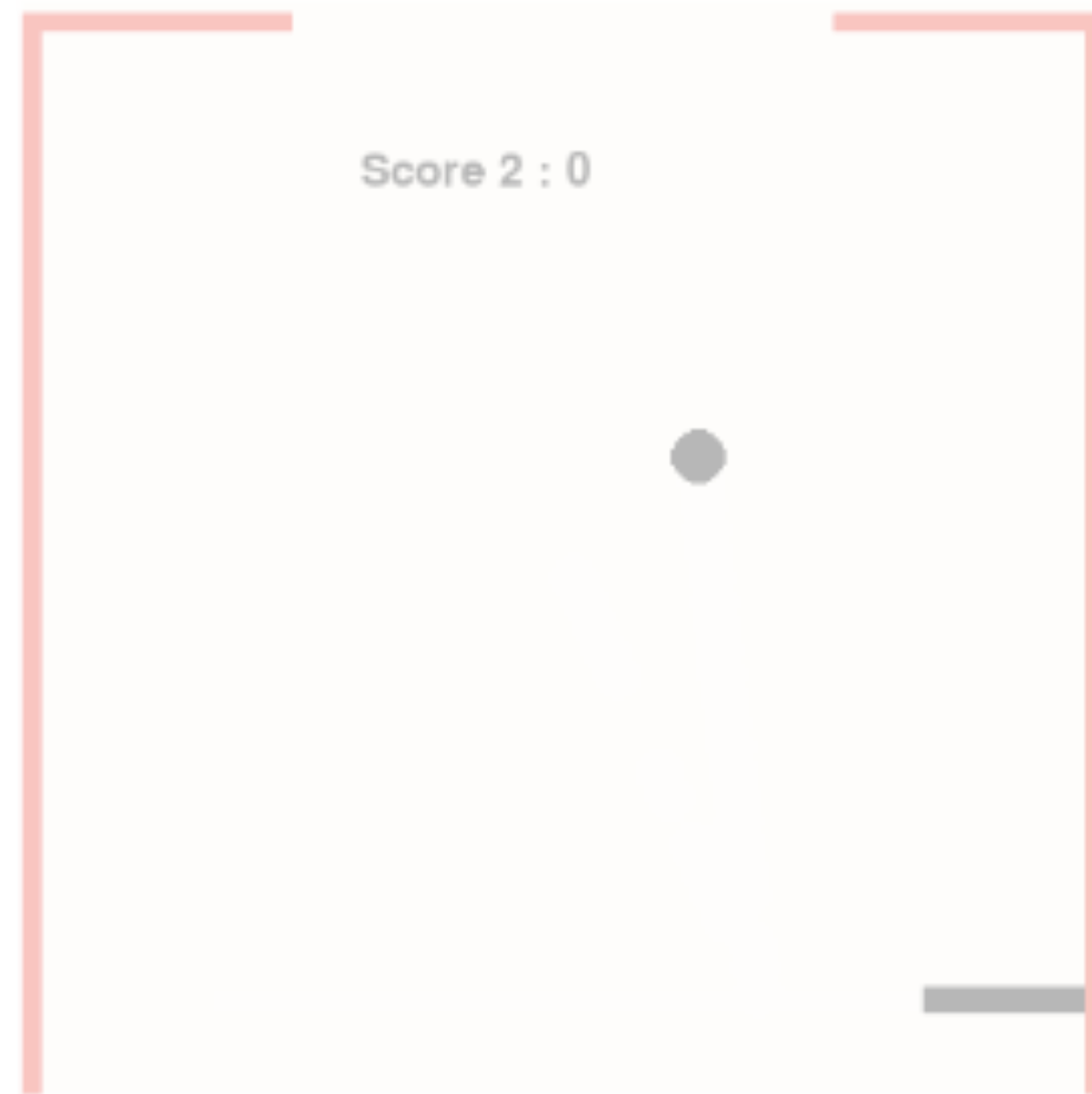
Common CS assignment: **interactive software**

```
Underlying env ID: 7340  
Env ID: 1  
Label: [1 1 0 0 0 0 1 0 0 1 0 1 0 1 1]  
Binary label: whenGoal-noBallLaunch  
Action: None  
Reward: 0  
Timestep: 0  
Exploration reward: 0.020  
Prob: 0.456
```

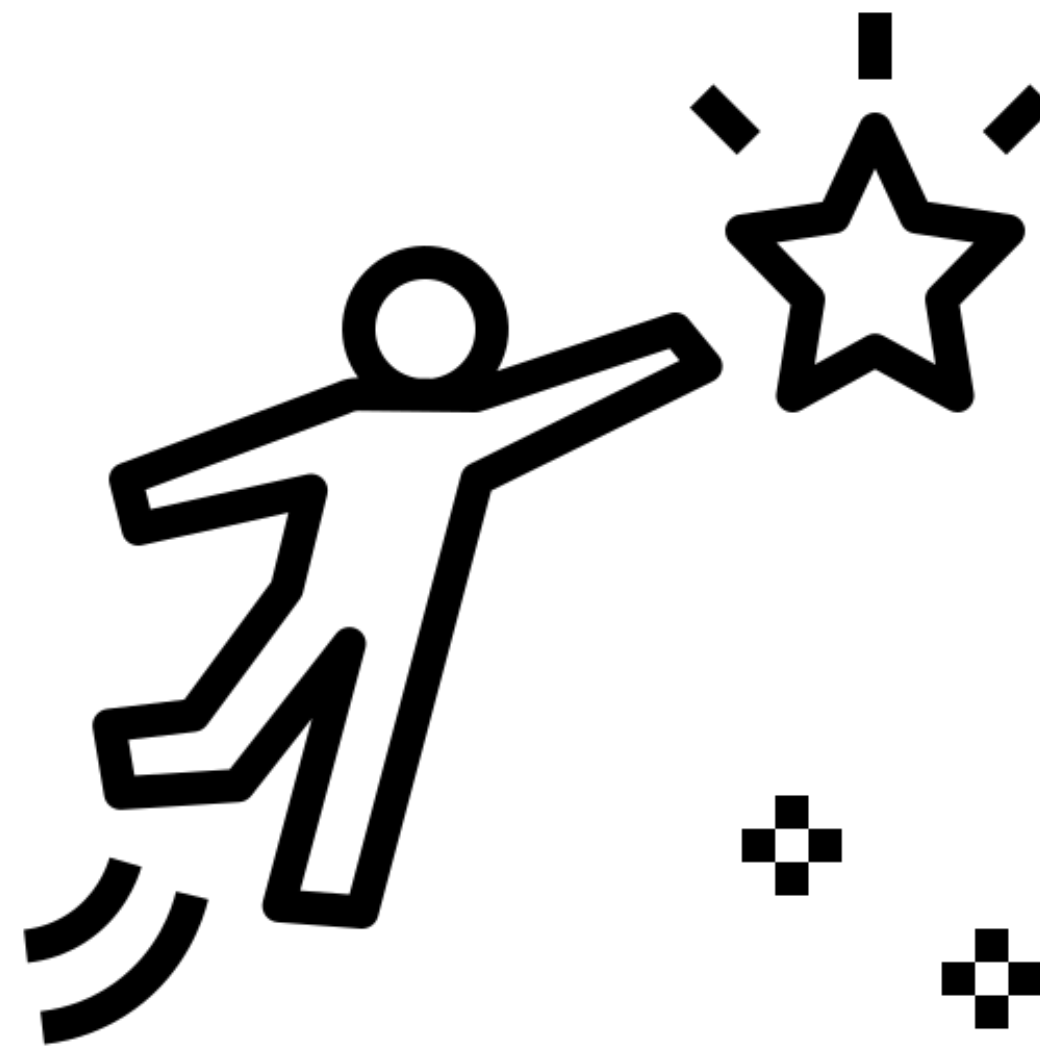
Bounce assignment
(Code.org)

Problem: Providing Feedback on Interactive Software

Common CS assignment: **interactive software**



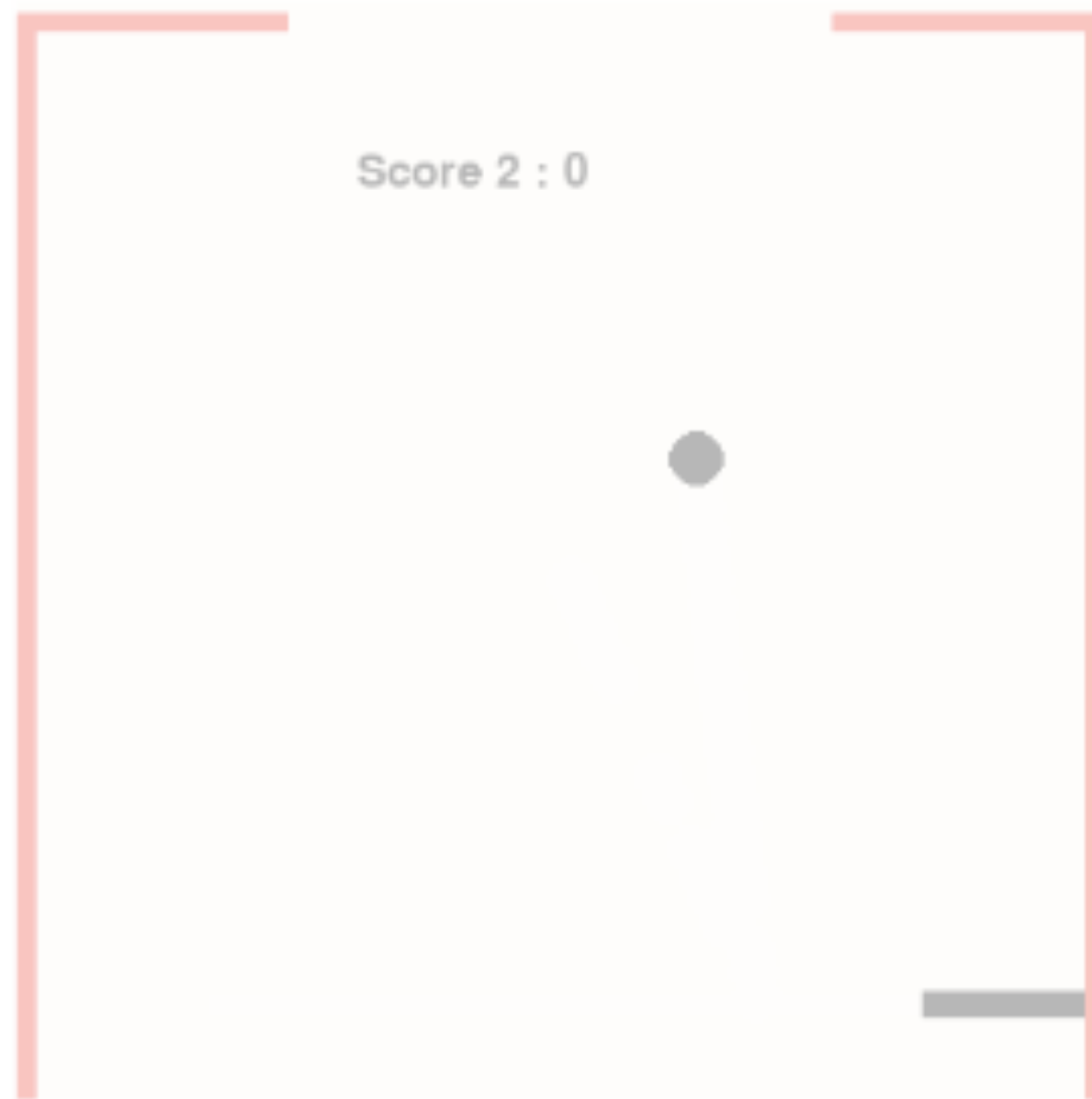
Bounce assignment
(Code.org)



Motivating and engaging (fun!)
⇒ can enrich learning (Pfaffman et al., 2003)

Problem: Providing Feedback on Interactive Software

Common CS assignment: **interactive software**



Bounce assignment
(Code.org)



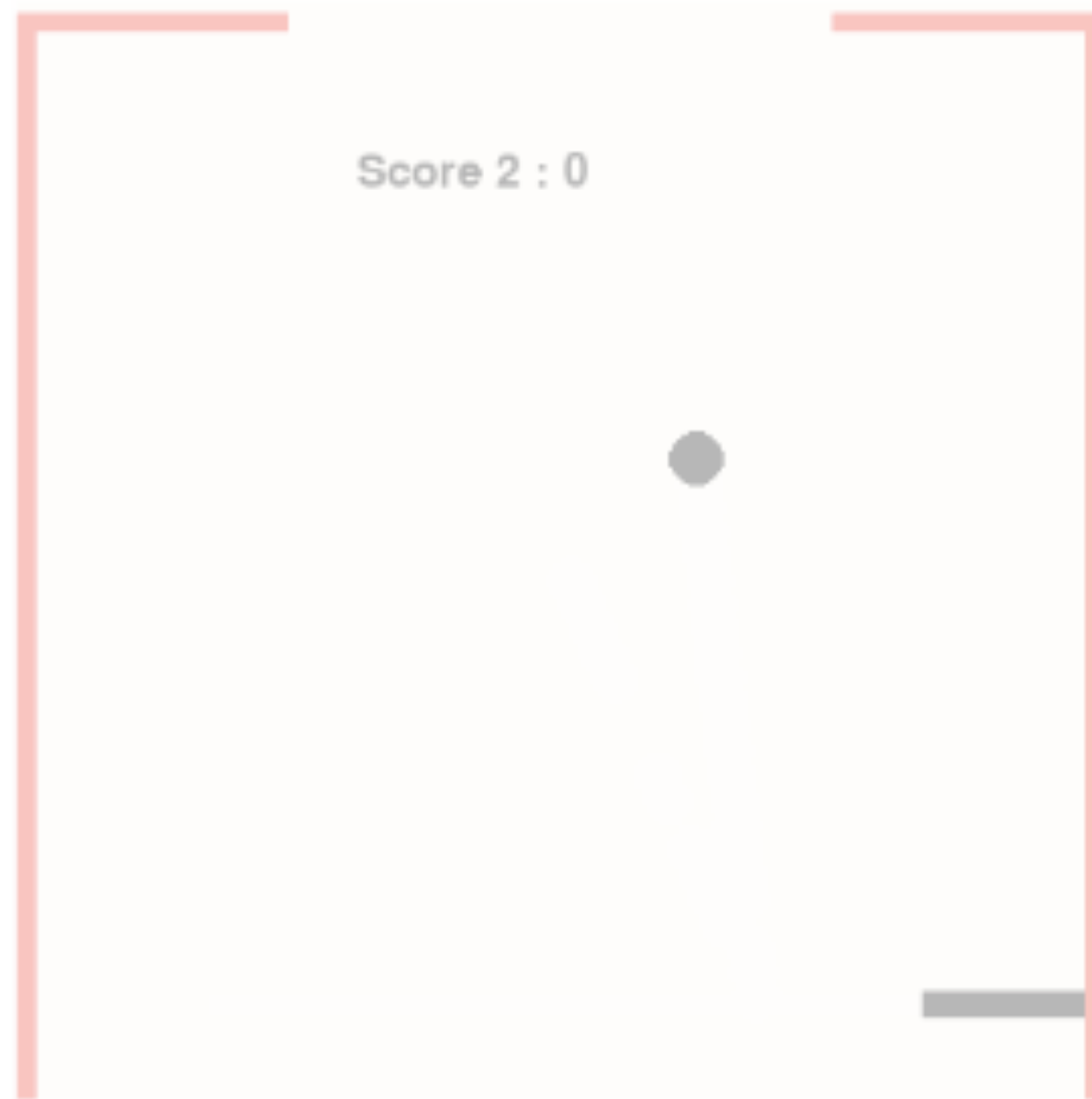
Motivating and engaging (fun!)
⇒ can enrich learning (Pfaffman et al., 2003)

Harvard CS50
UC Berkeley CS61B
UCLA CS32
Stanford CS106A
Code.org
Camp K12
Tynker
Google Applied CS Skills
...

Increasingly found
everywhere

Problem: Providing Feedback on Interactive Software

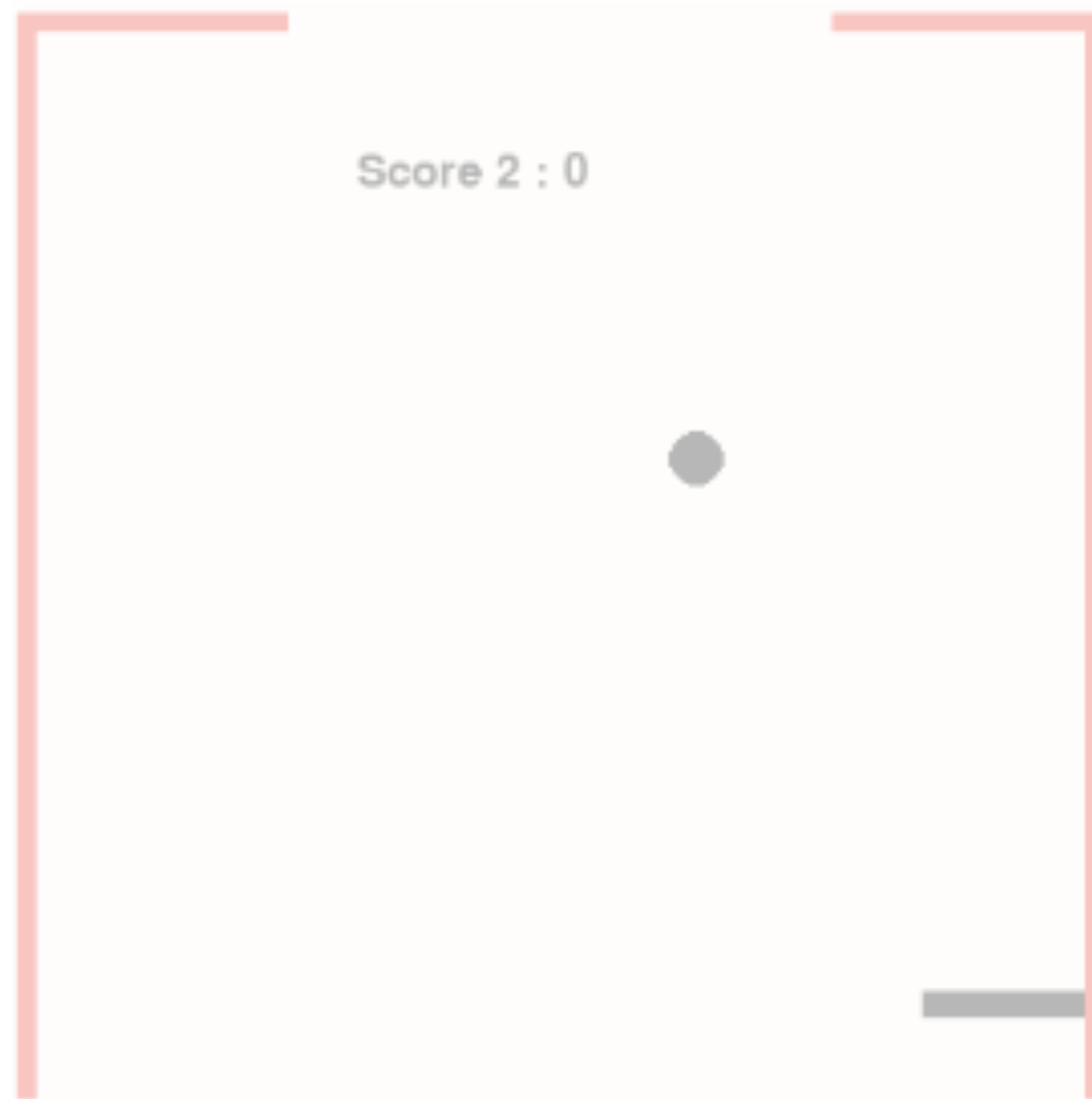
Providing feedback / grading is hard



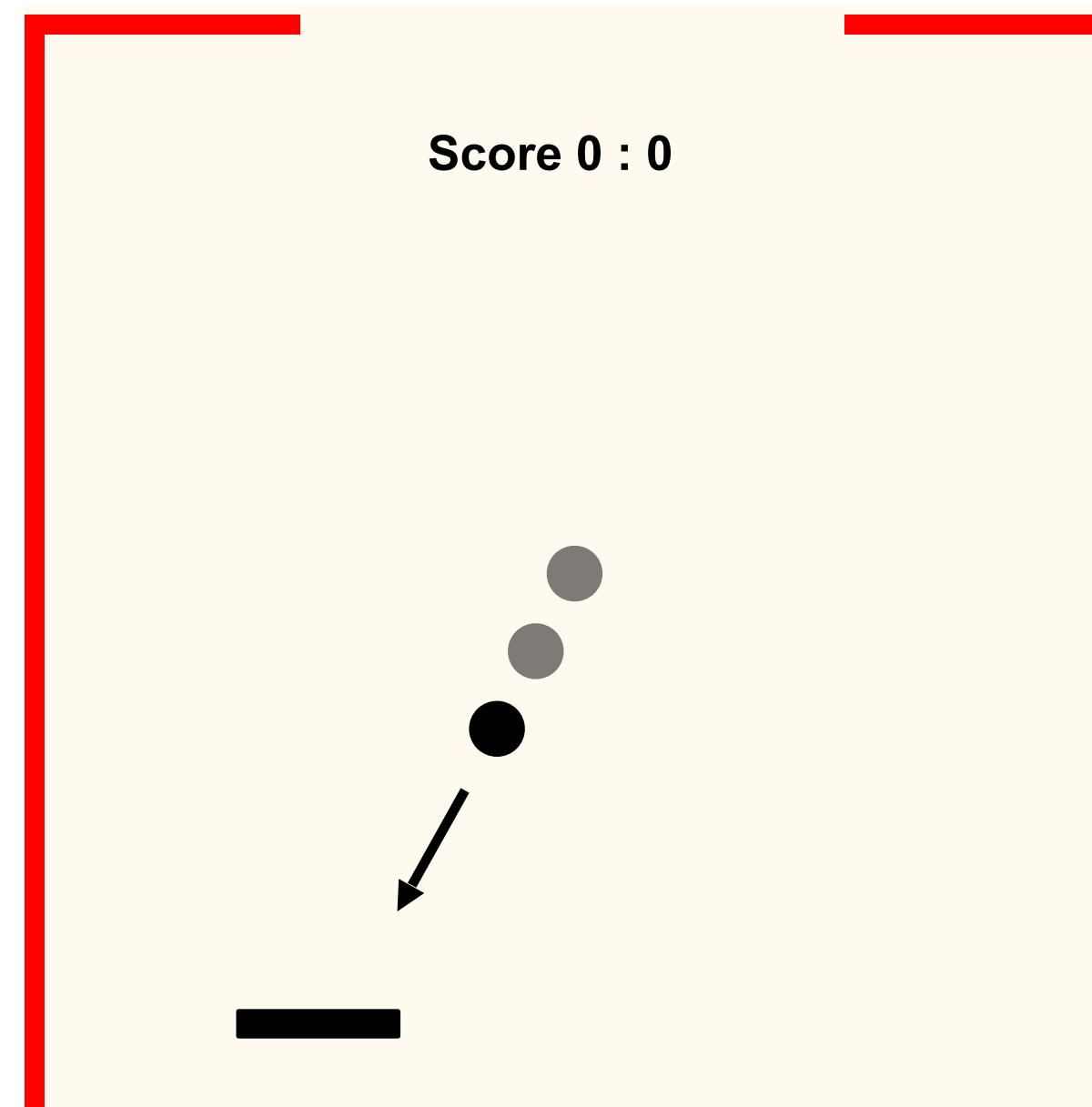
Bounce assignment
(Code.org)

Problem: Providing Feedback on Interactive Software

Providing feedback / grading is hard



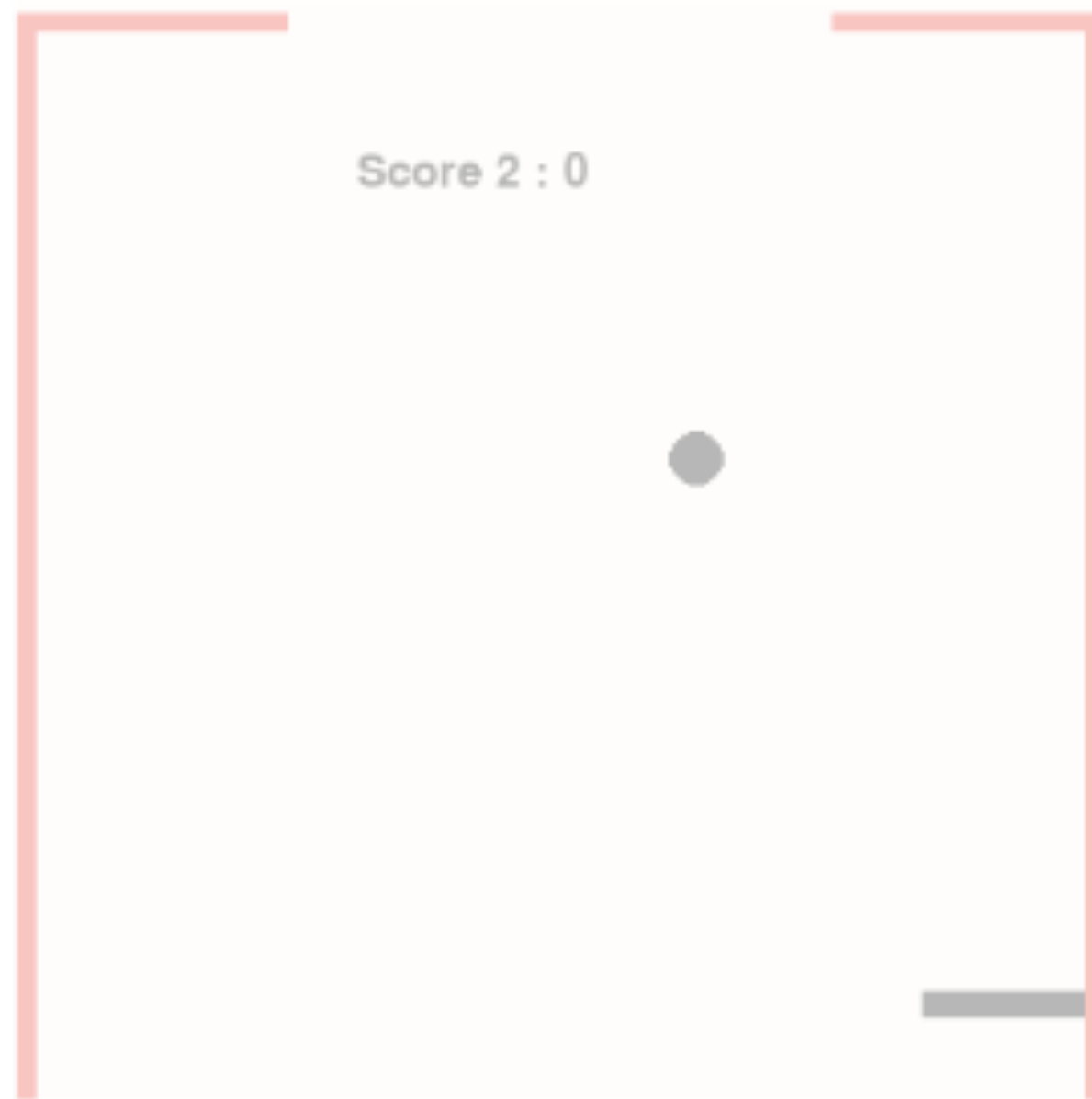
Bounce assignment
(Code.org)



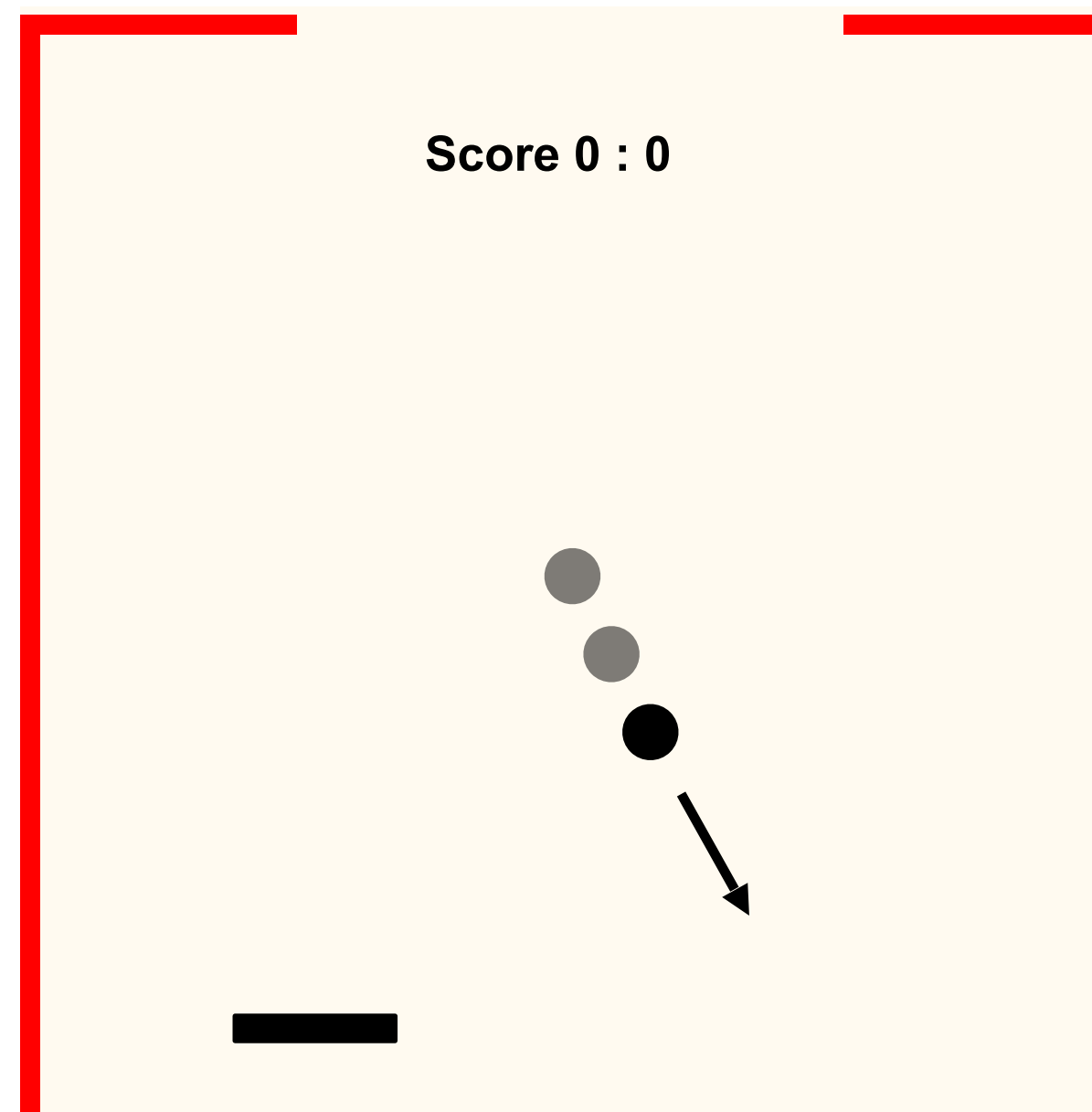
Stochasticity

Problem: Providing Feedback on Interactive Software

Providing feedback / grading is hard



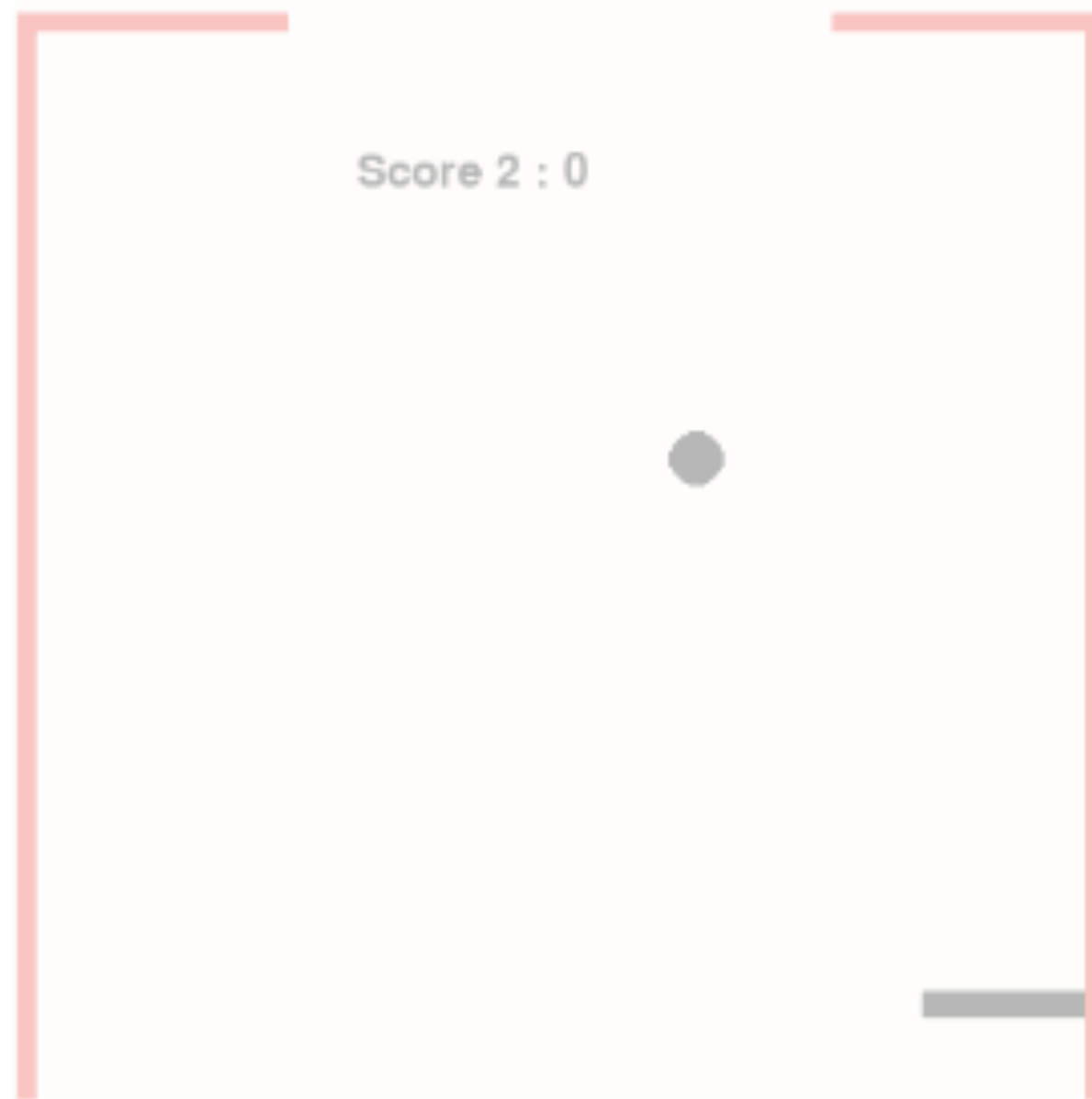
Bounce assignment
(Code.org)



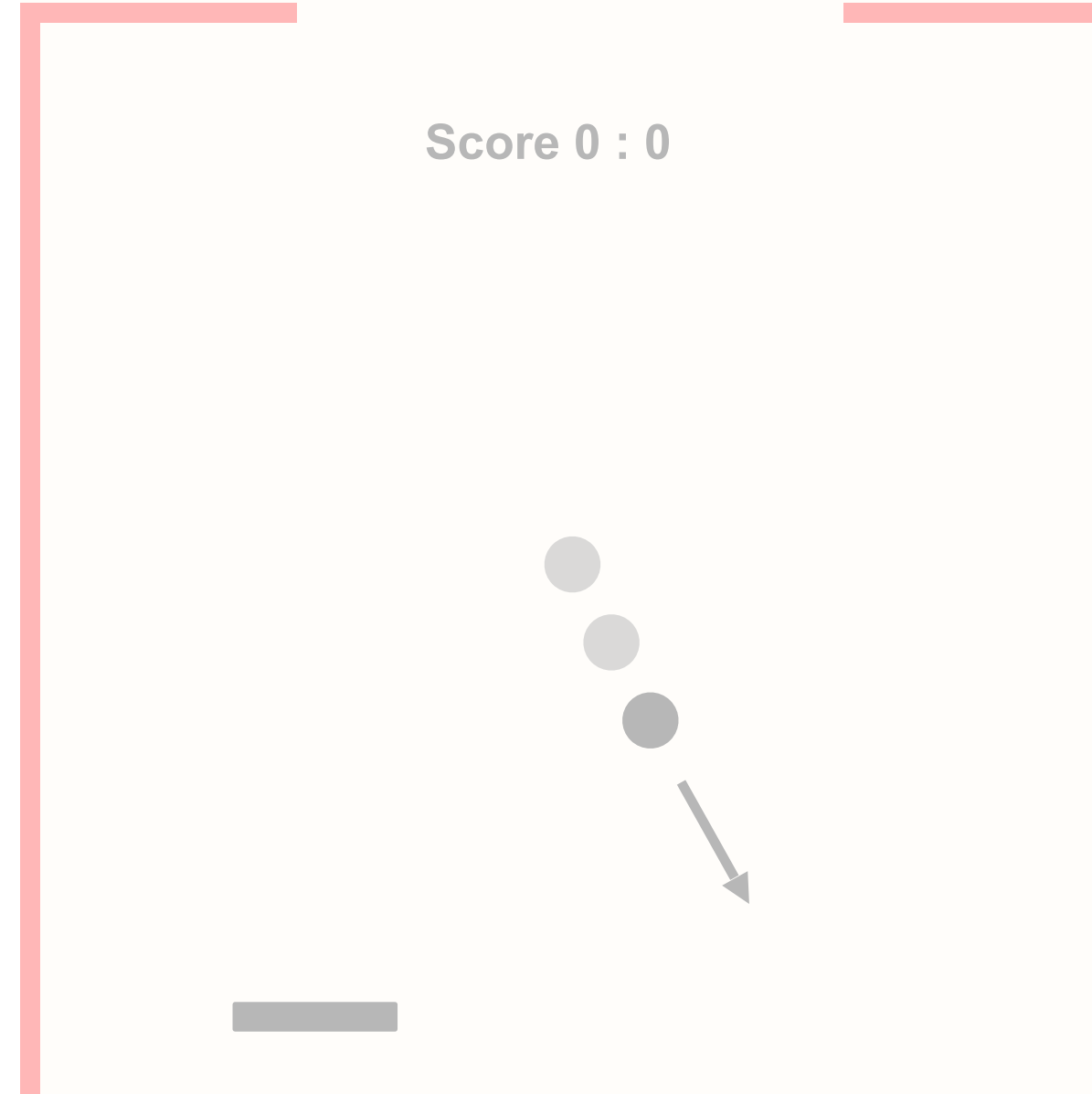
Stochasticity

Problem: Providing Feedback on Interactive Software

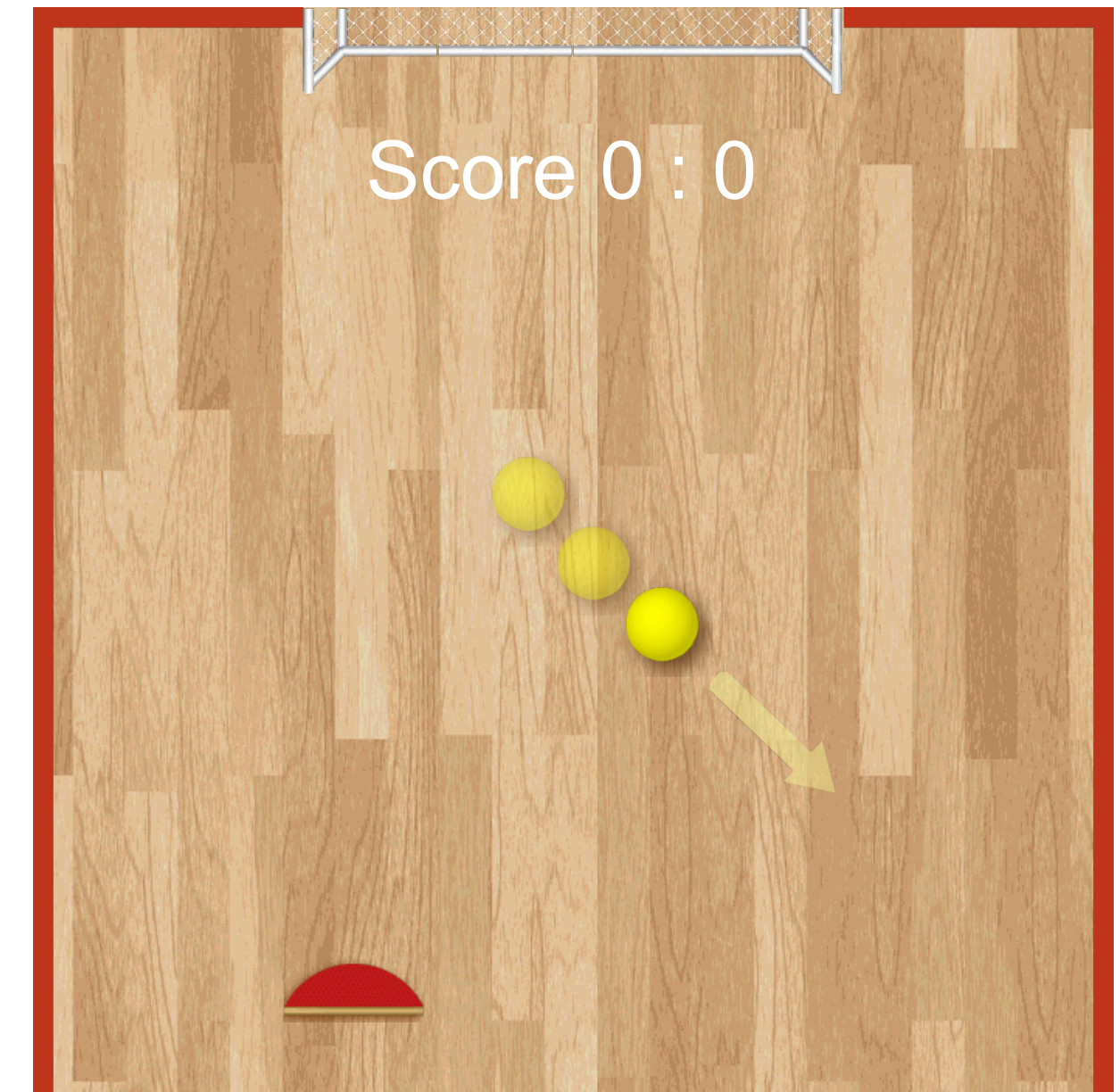
Providing feedback / grading is hard



Bounce assignment
(Code.org)

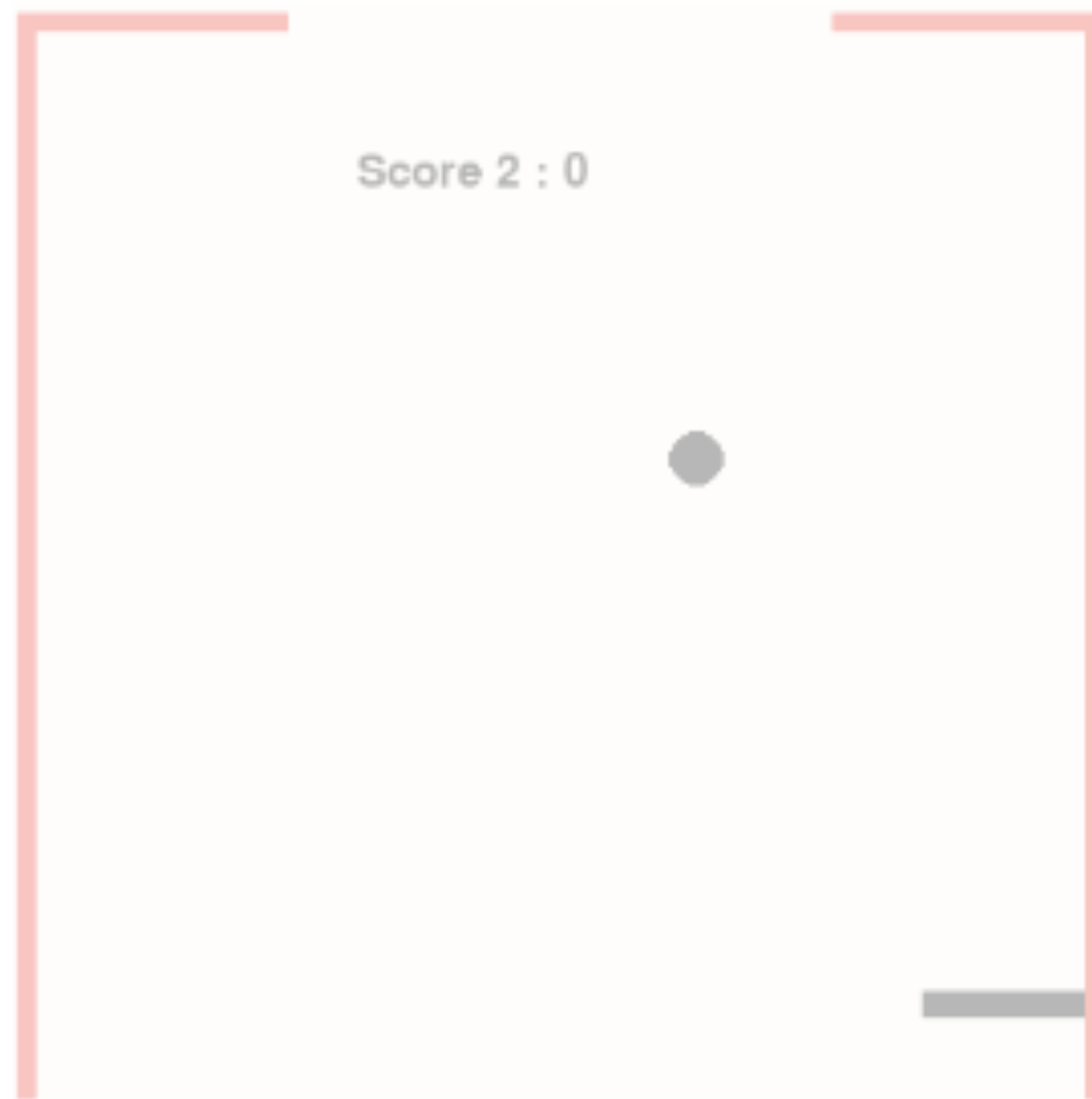


Stochasticity

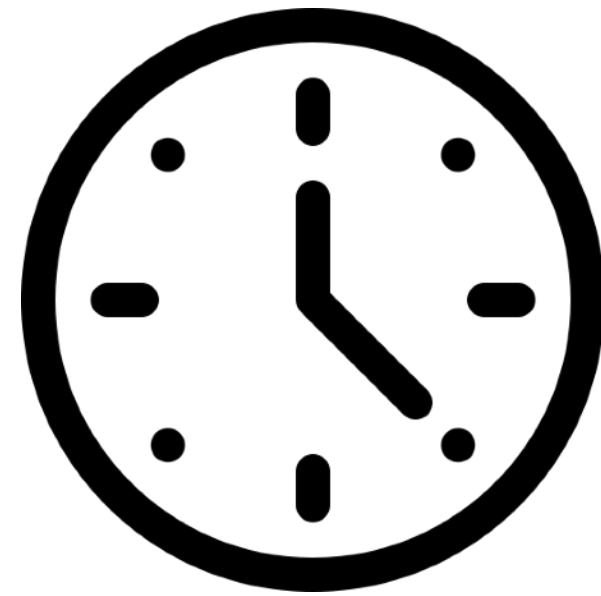


Student creativity

Problem: Providing Feedback on Interactive Software



Bounce assignment
(Code.org)



3+ min manual grading
per assignment

×



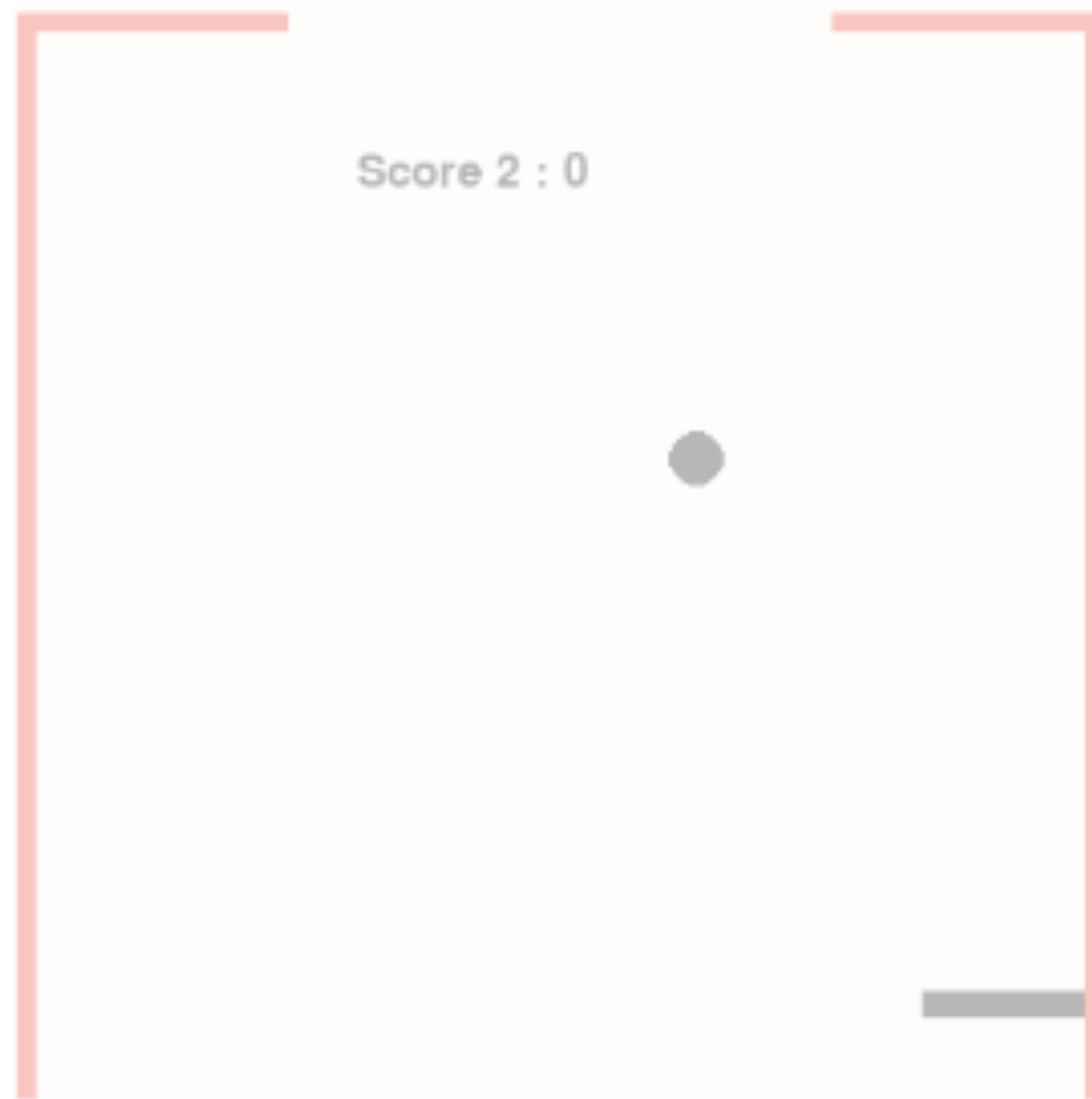
70M+
enrolled students

=



300+ years
of grading work

Problem: Providing Feedback on Interactive Software



Bounce assignment
(Code.org)



3+ min manual grading
per assignment



70M+
enrolled students



300+ years
of grading work

Our goal: Automatically provide feedback



Reduce enormous human grading burden



Provide *faster* and *iterative* feedback

Setting

Rubric: List of possible errors



moveError

whenWall-newBallError

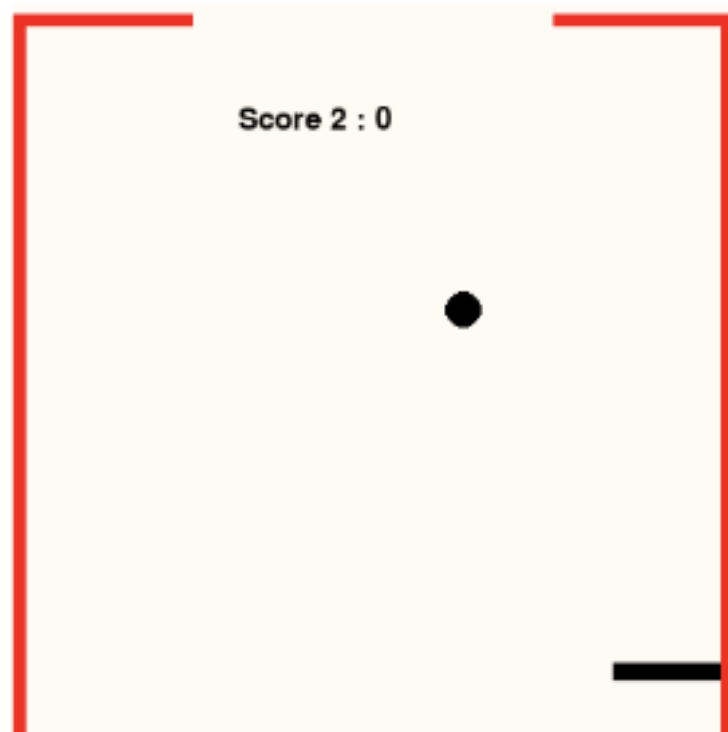
whenGoal-scoreError

...

Setting

Training (~3500 labeled programs)

Program μ



Label y : Subset of rubric items present in program

moveError

~~whenGoal-scoreError~~

...

Rubric: List of possible errors



moveError

whenWall-newBallError

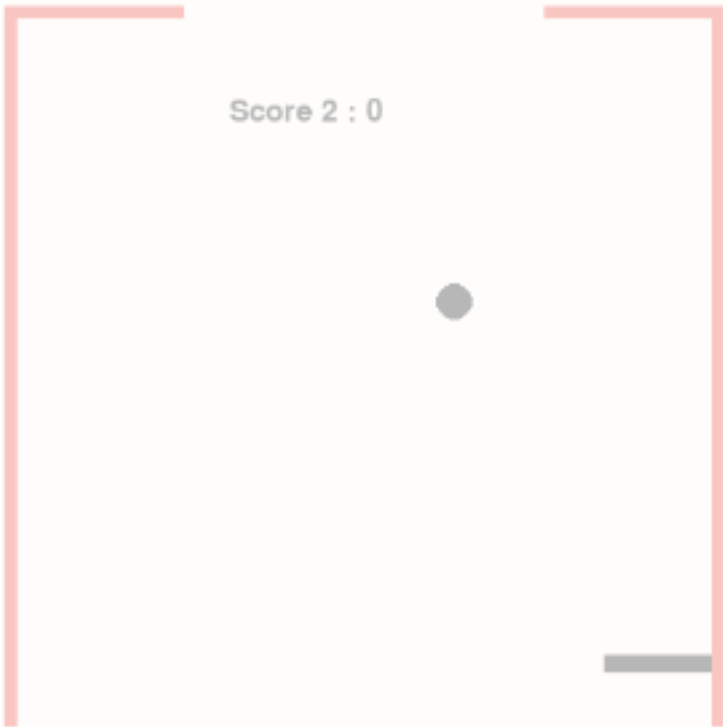
whenGoal-scoreError

...

Setting

Training (~3500 labeled programs)

Program μ



Score 2 : 0

Label y : Subset of rubric items present in program

moveError

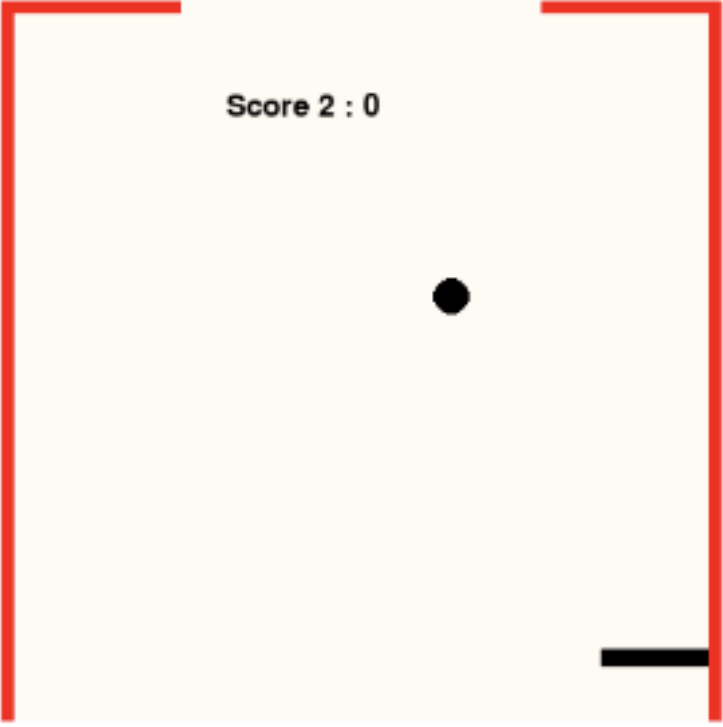
~~whenGoal scoreError~~

...

The training data consists of multiple instances of a soccer game interface. Each instance is a rectangular frame containing a soccer field. The field has a yellow background, a black center circle, and a black goal line on the right. A score 'Score 2 : 0' is displayed at the top left. A red border highlights the field area. To the right of the field, a list of rubric items is shown: 'moveError' in a green box, and 'whenGoal scoreError' which is crossed out with a red line. Ellipses indicate other items.

Testing

New student program

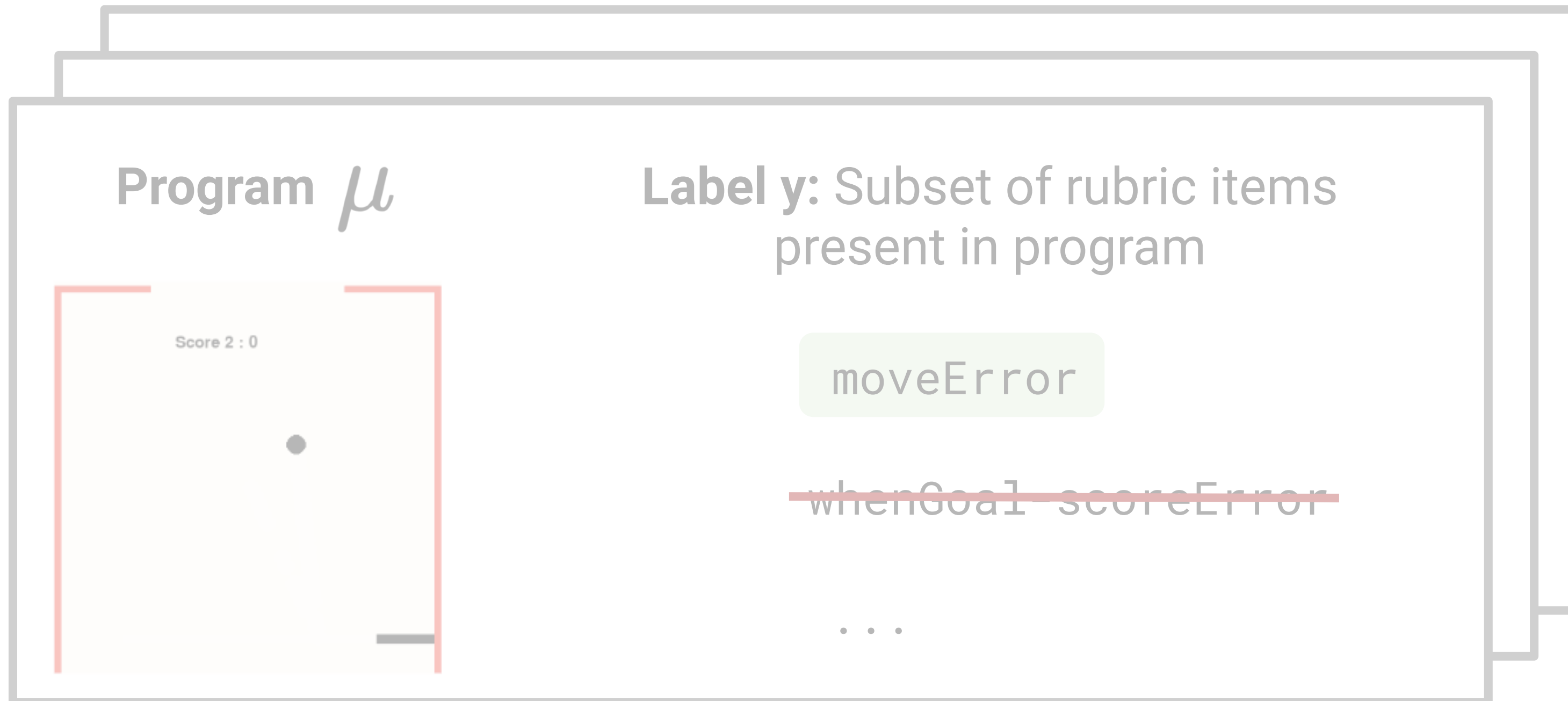


Score 2 : 0

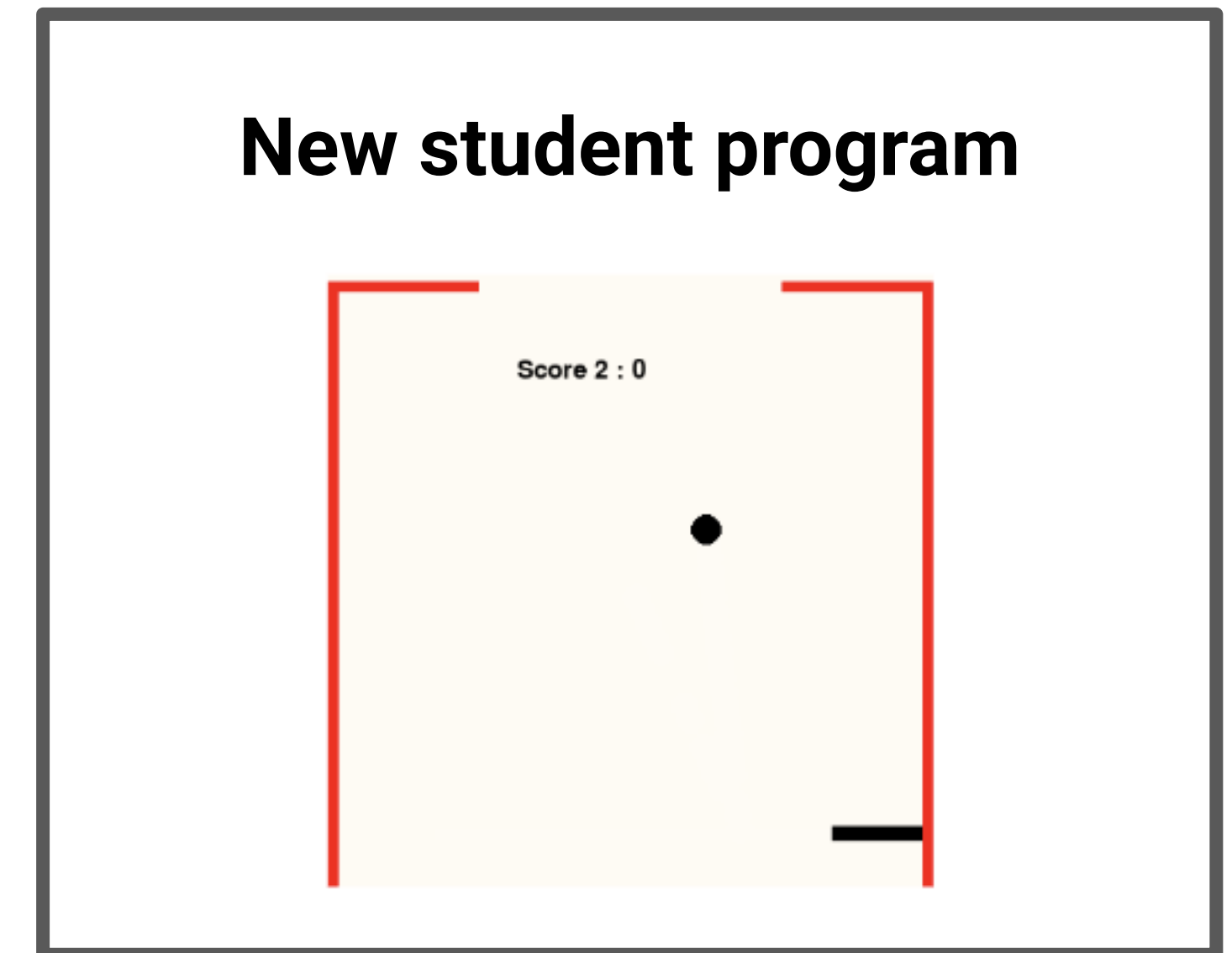
The testing phase shows a single instance of a soccer game interface, labeled 'New student program'. It is a rectangular frame containing a soccer field. The field has a yellow background, a black center circle, and a black goal line on the right. A score 'Score 2 : 0' is displayed at the top left. A red border highlights the field area.

Setting

Training (~3500 labeled programs)



Testing



Goal: Output which bugs are in the program (i.e., predict the label)

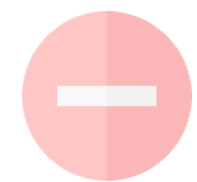
Related Work: Two Paradigms for Automated Feedback

Analyze program *code*

(Singh et al., '13, Piech et al., '15, Bhatia et al., '16, , Rivers et al., '17, Paaßen et al., '17, Wang et al., '17, Malik et al., '19, Wu et al., '19, Wu et al., '21)



Works well for shorter programs (e.g., <50 lines of code)



Existing methods struggle to scale to longer programs

Analyze program *behavior*

(King et al., '76, Godefroid et al., '08, Zheng et al., '19, Nie et al., '21, Gordillo et al., '21)



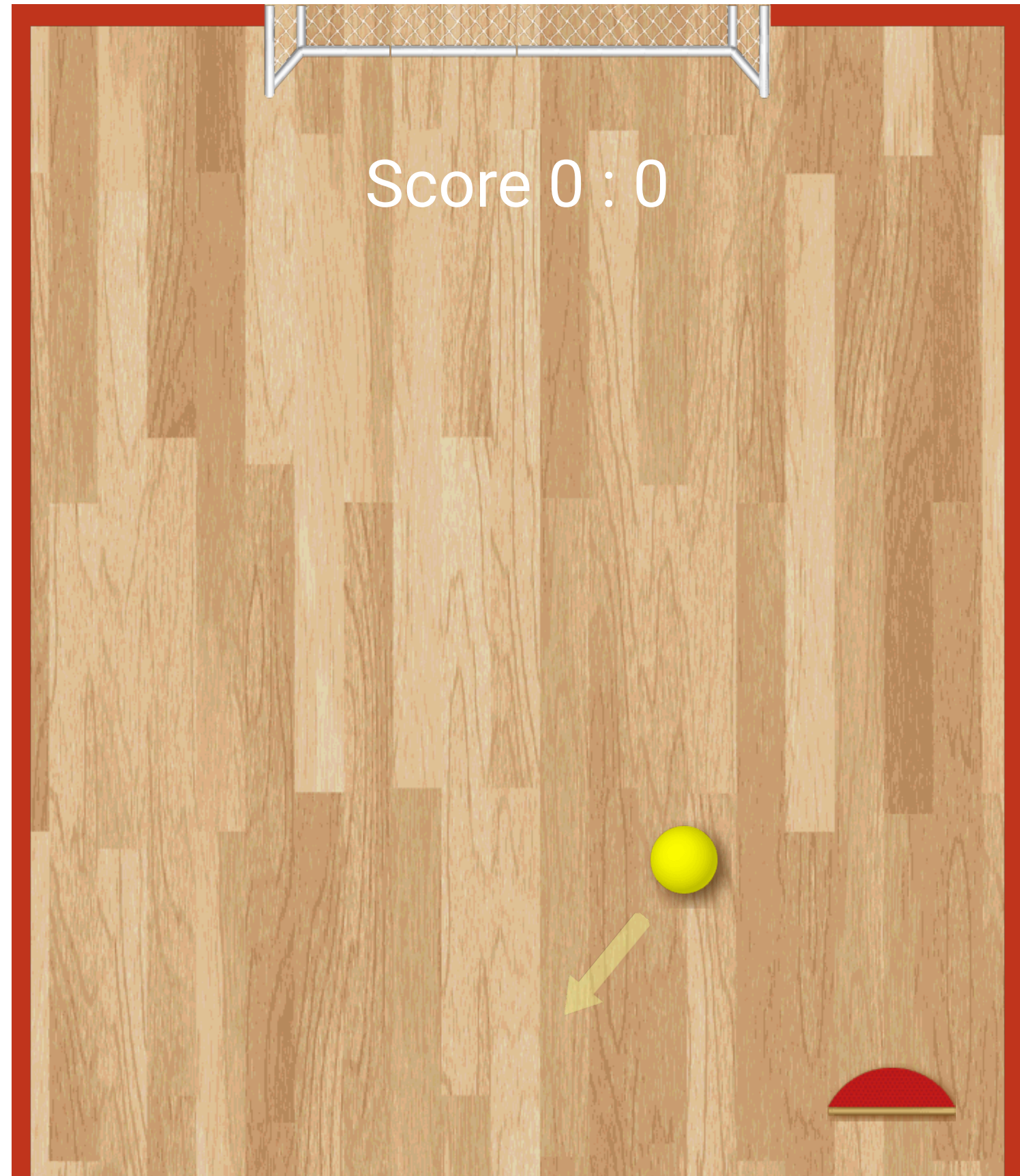
Independent of program length



Assumes that the program can compile and run

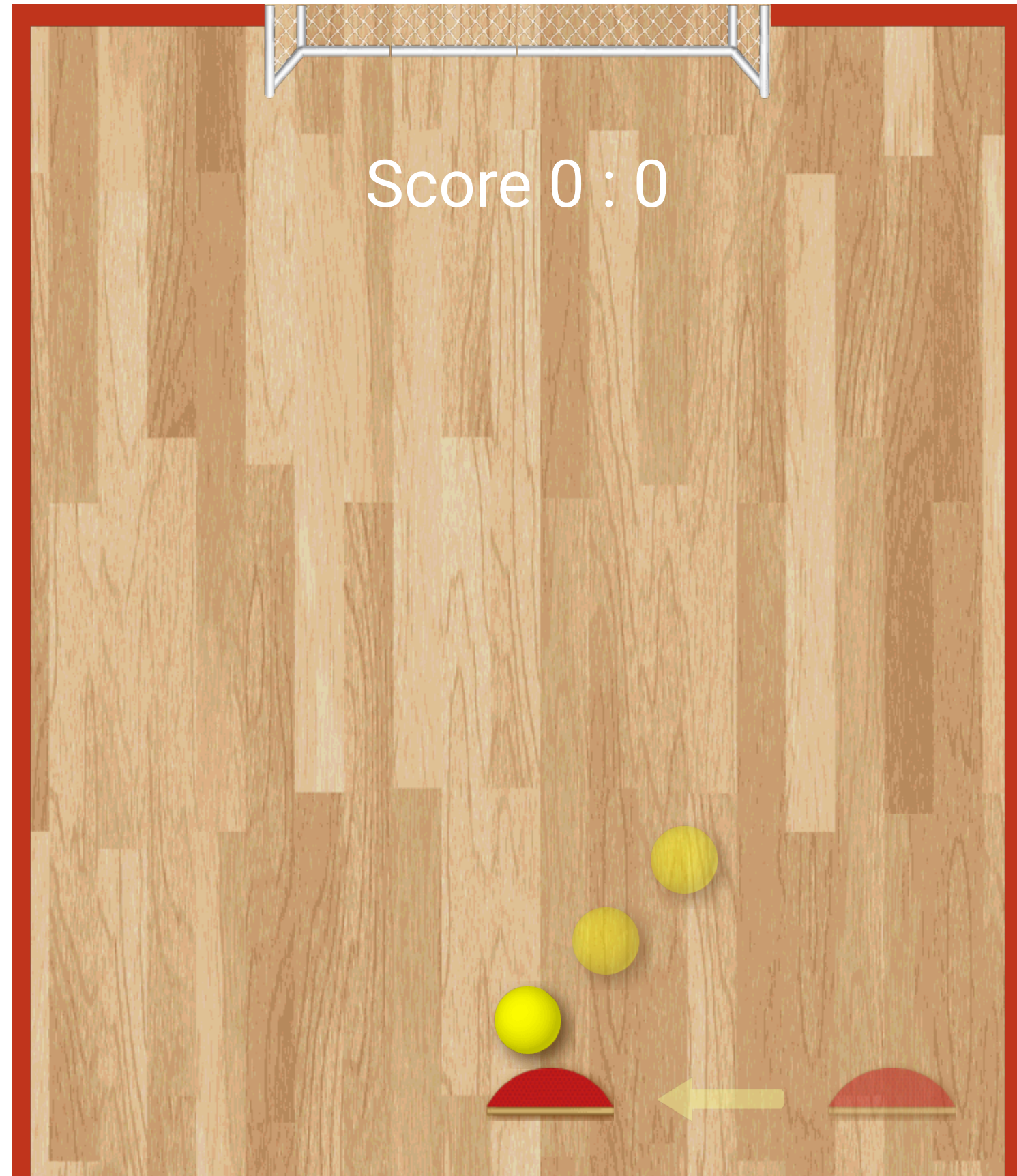
We opt for this approach

The Play-to-Grade Paradigm



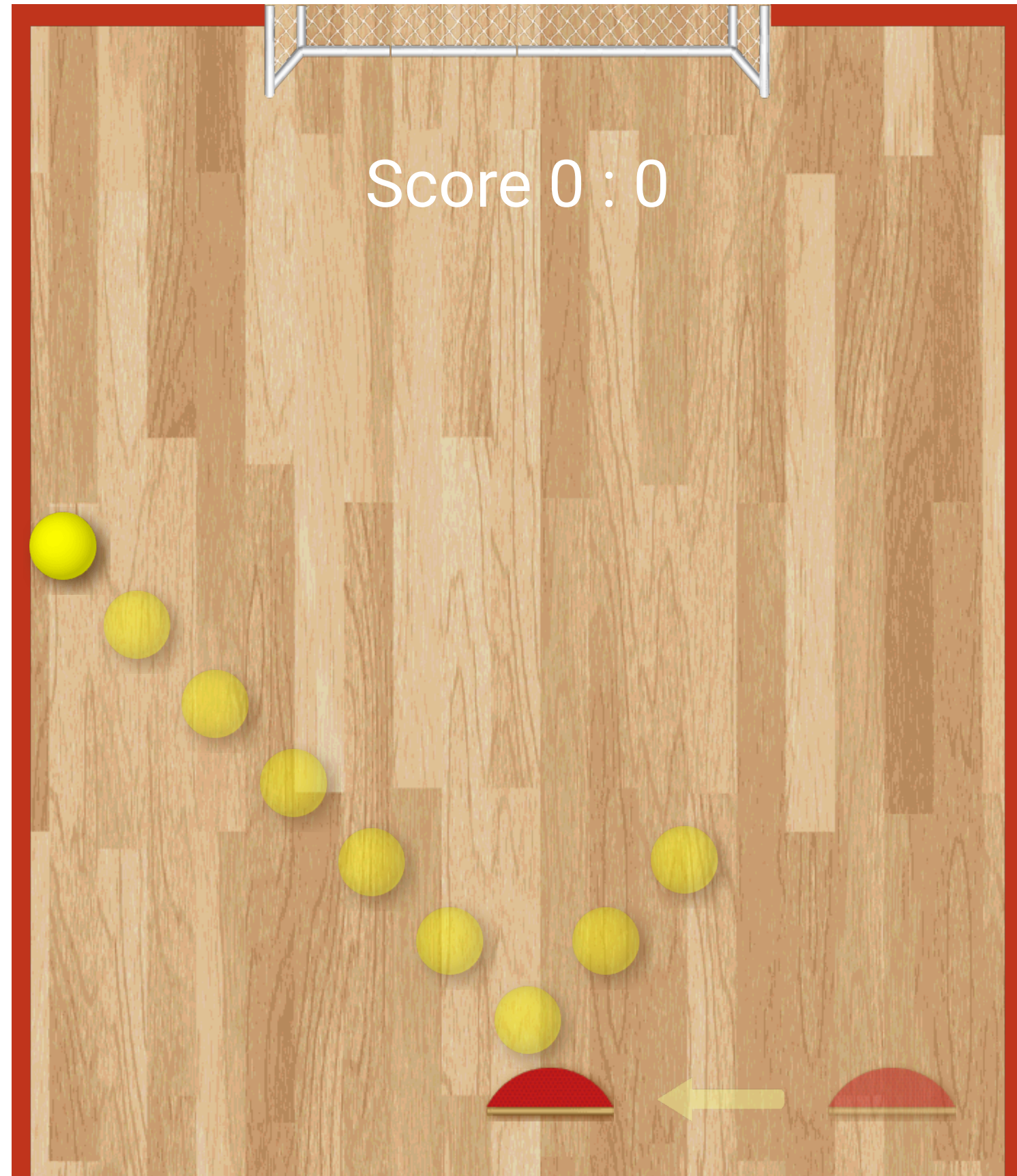
Agent interacts with program like human

The Play-to-Grade Paradigm



Agent interacts with program like human

The Play-to-Grade Paradigm



Agent interacts with program like human

The Play-to-Grade Paradigm



Agent interacts with program like human

The Play-to-Grade Paradigm



Agent interacts with program like human

The Play-to-Grade Paradigm



Agent interacts with program like human

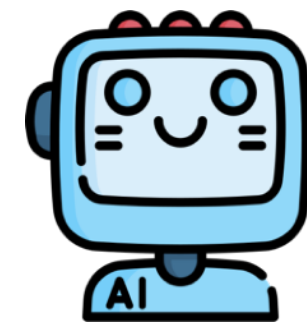
Agent outputs feedback

The Play-to-Grade Paradigm



Agent interacts with program like human

Existing work (Nie et al., '21):
Coarse binary feedback



The program is
incorrect

Agent outputs feedback

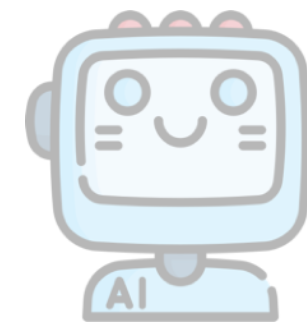
Not **specific** enough for
student to learn and
correct mistakes

The Play-to-Grade Paradigm



Agent interacts with program like human

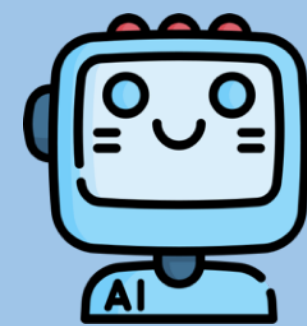
Existing work (Nie et al., '21):
Coarse binary feedback



The program is incorrect



Our goal: **Fine-grained** feedback

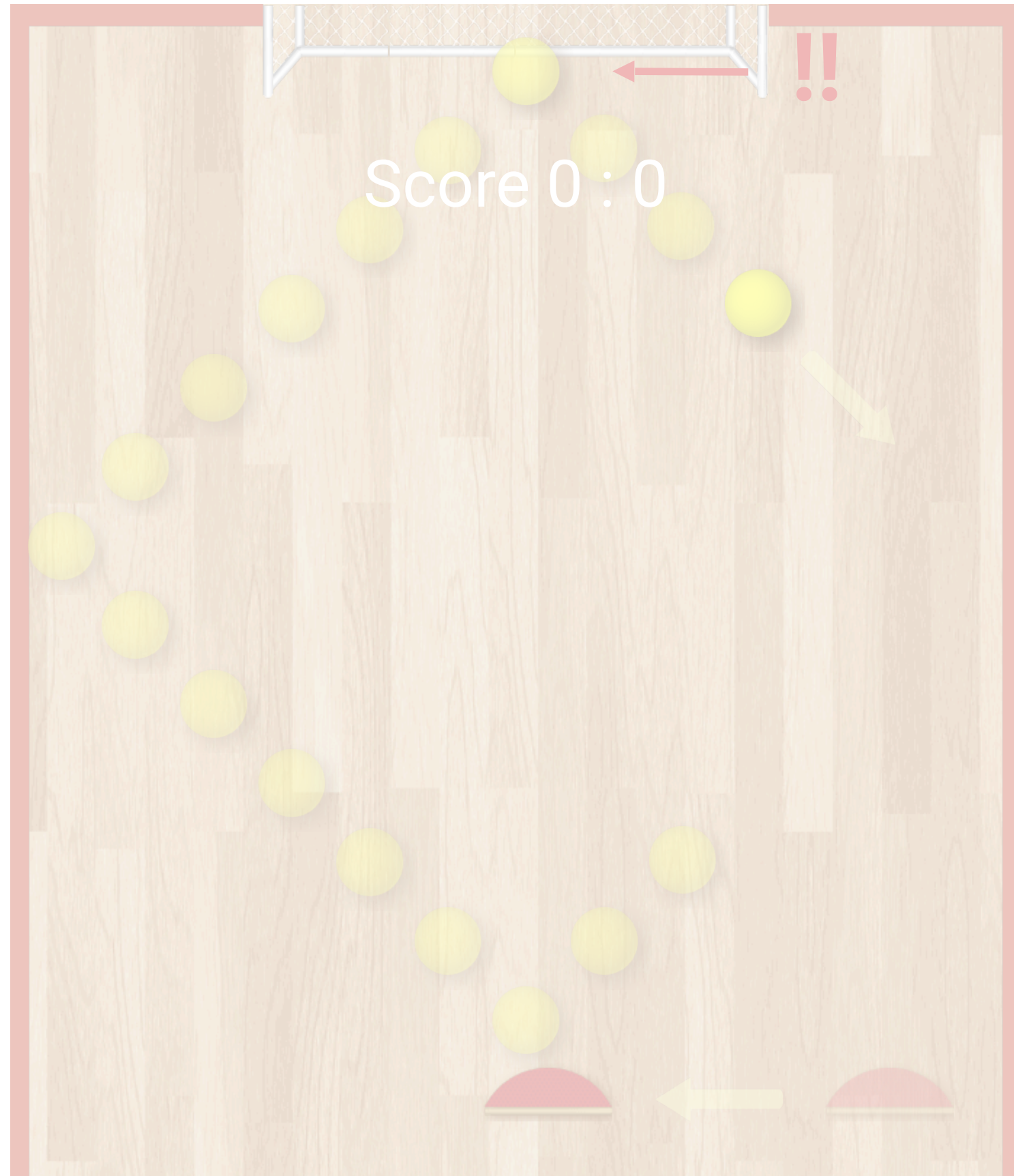


The ball incorrectly bounces off the goal

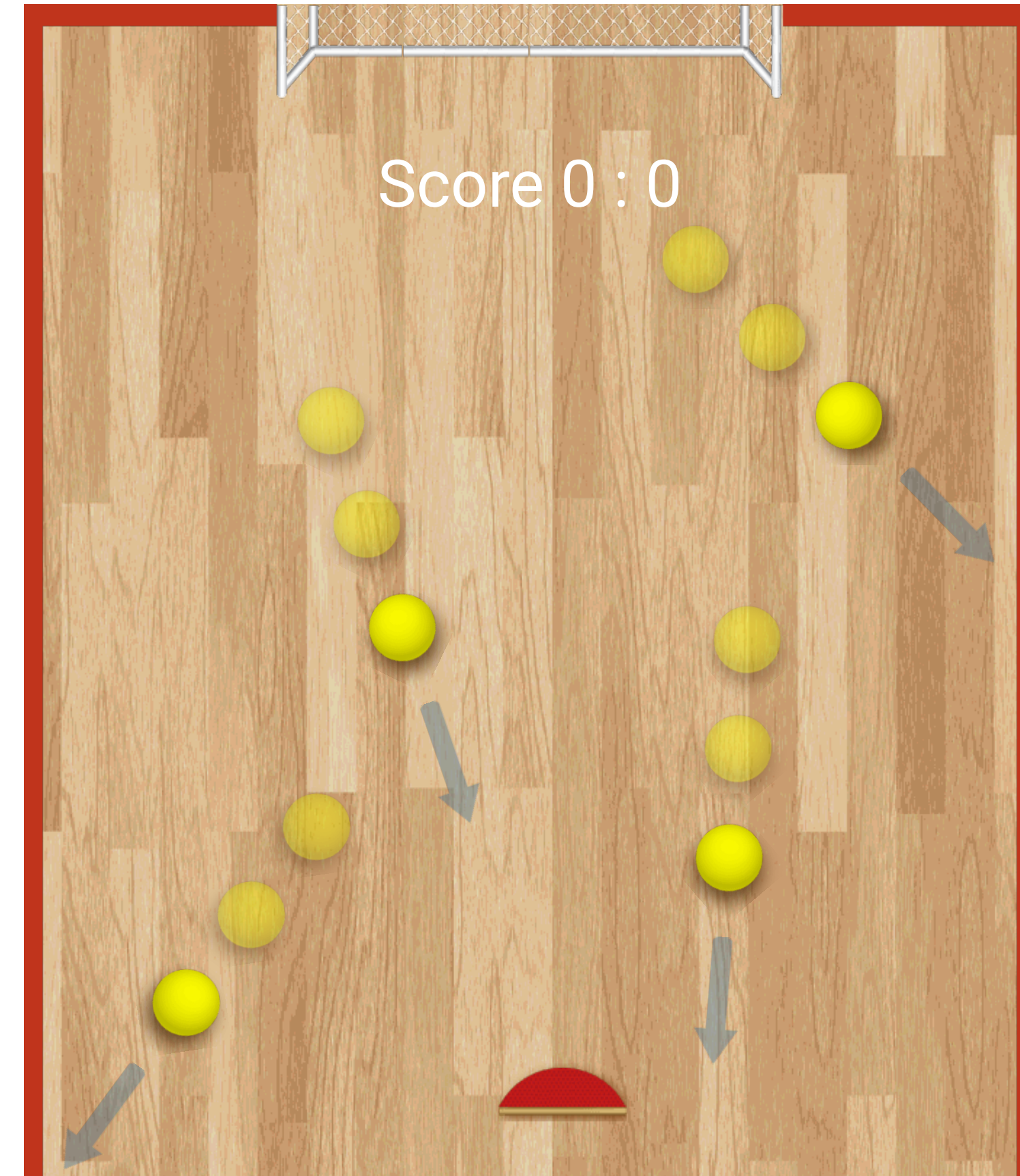
Agent outputs feedback

Not **specific** enough for student to learn and correct mistakes

What makes providing feedback hard?



Targeted exploration

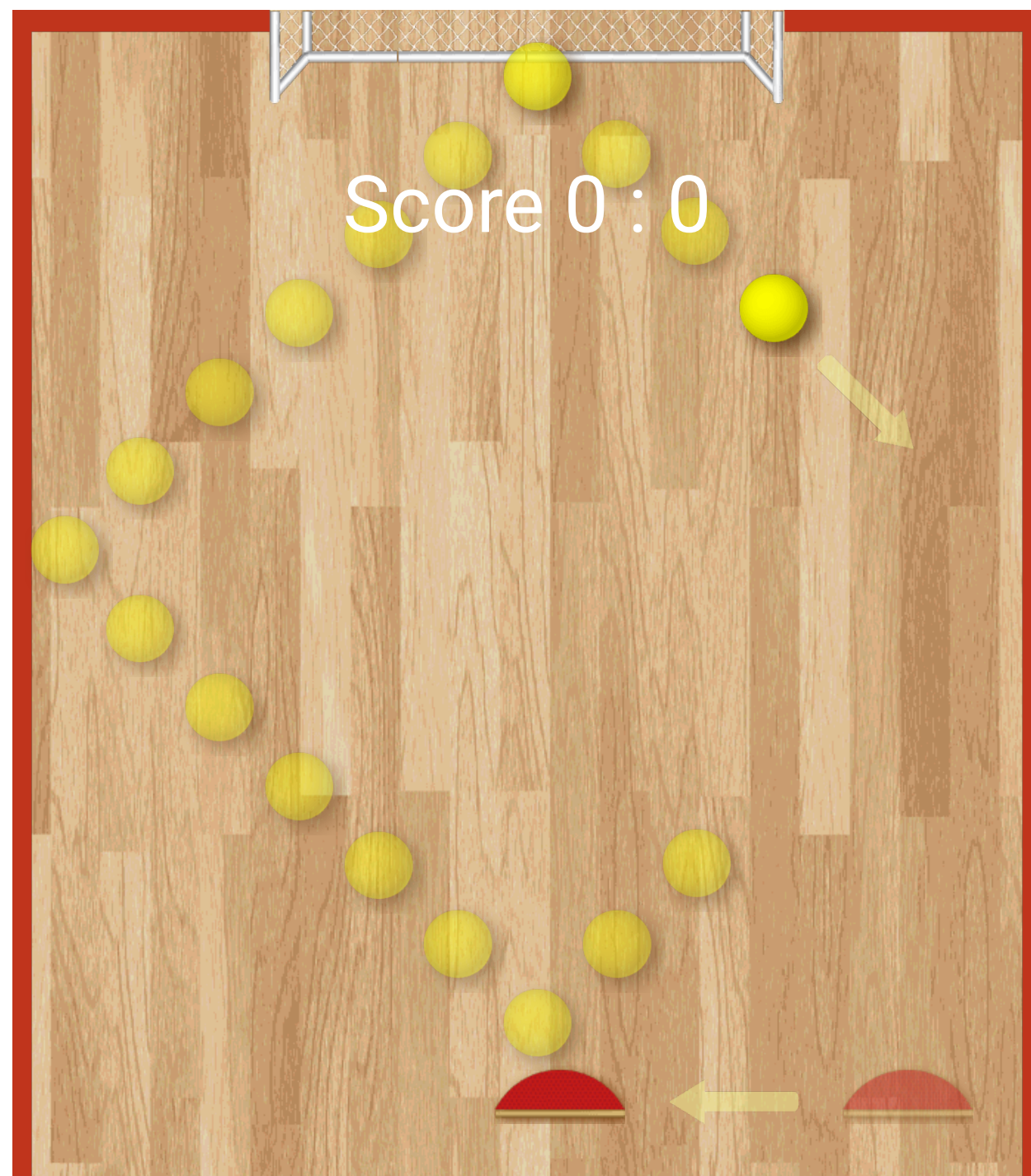


Adaptive exploration

Approach: DREAMGRADER

Exploration policy π :

Takes program μ and produces trajectories \mathcal{T}



Feedback classifier g :

Takes trajectories \mathcal{T} and predicts label y

whenWall-newBallError

whenWall-scoreError

...

Approach: DREAMGRADER

Exploration policy π :

Takes program μ and produces trajectories τ

Feedback classifier g :

Takes trajectories τ and predicts label y

Maximize probability of correct label

$$\mathcal{J}(\pi, g) = \mathbb{E}_{\mu \sim p(\mu), \tau \sim \pi(\mu)} [g(y \mid \tau)]$$

Sample a program and roll out
exploration policy

Approach: DREAMGRADER

Exploration policy π :

Takes program μ and produces trajectories τ

Feedback classifier g :

Takes trajectories τ and predicts label y

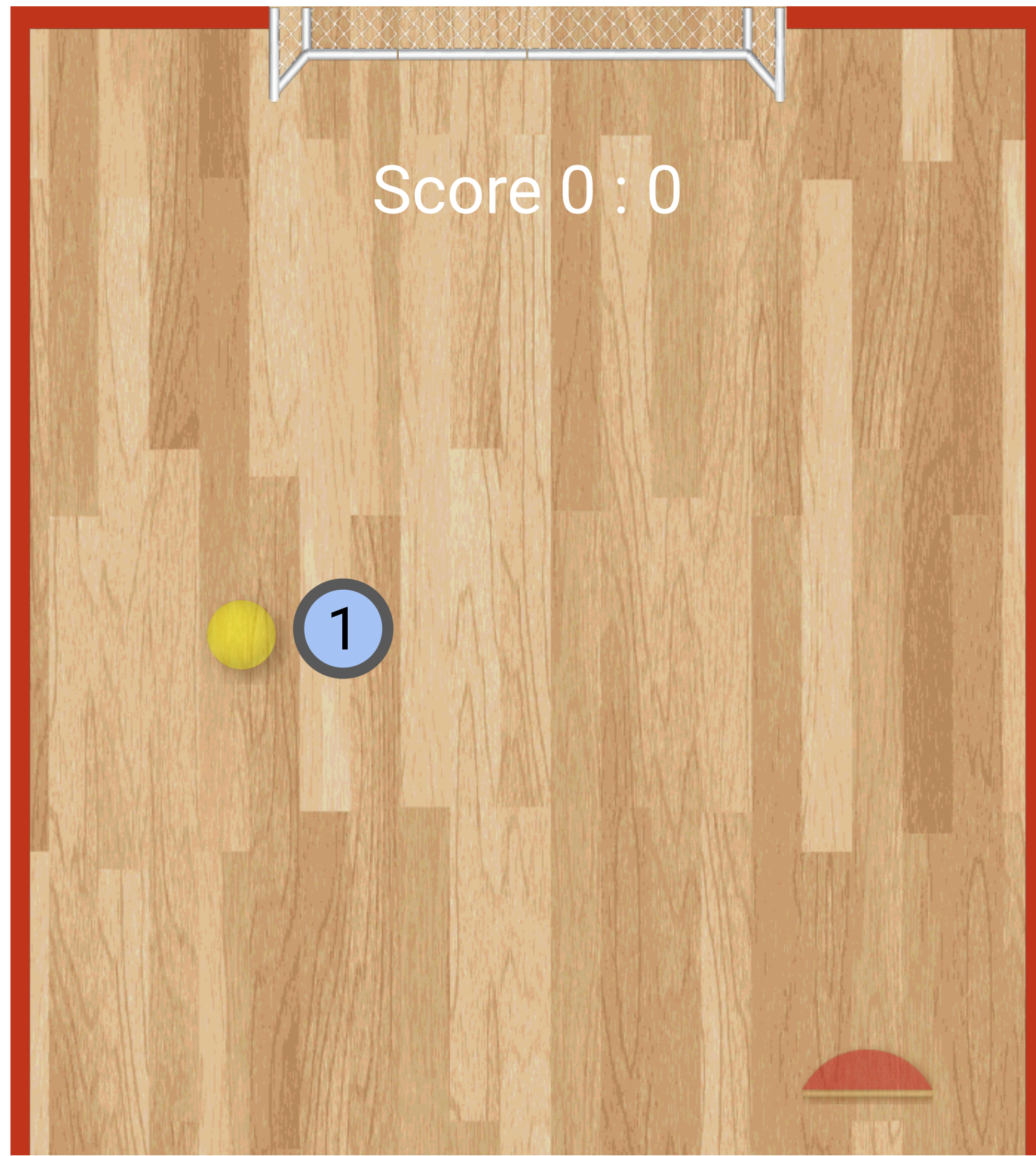
Naive approach:

Treat this as end-of-episode reward

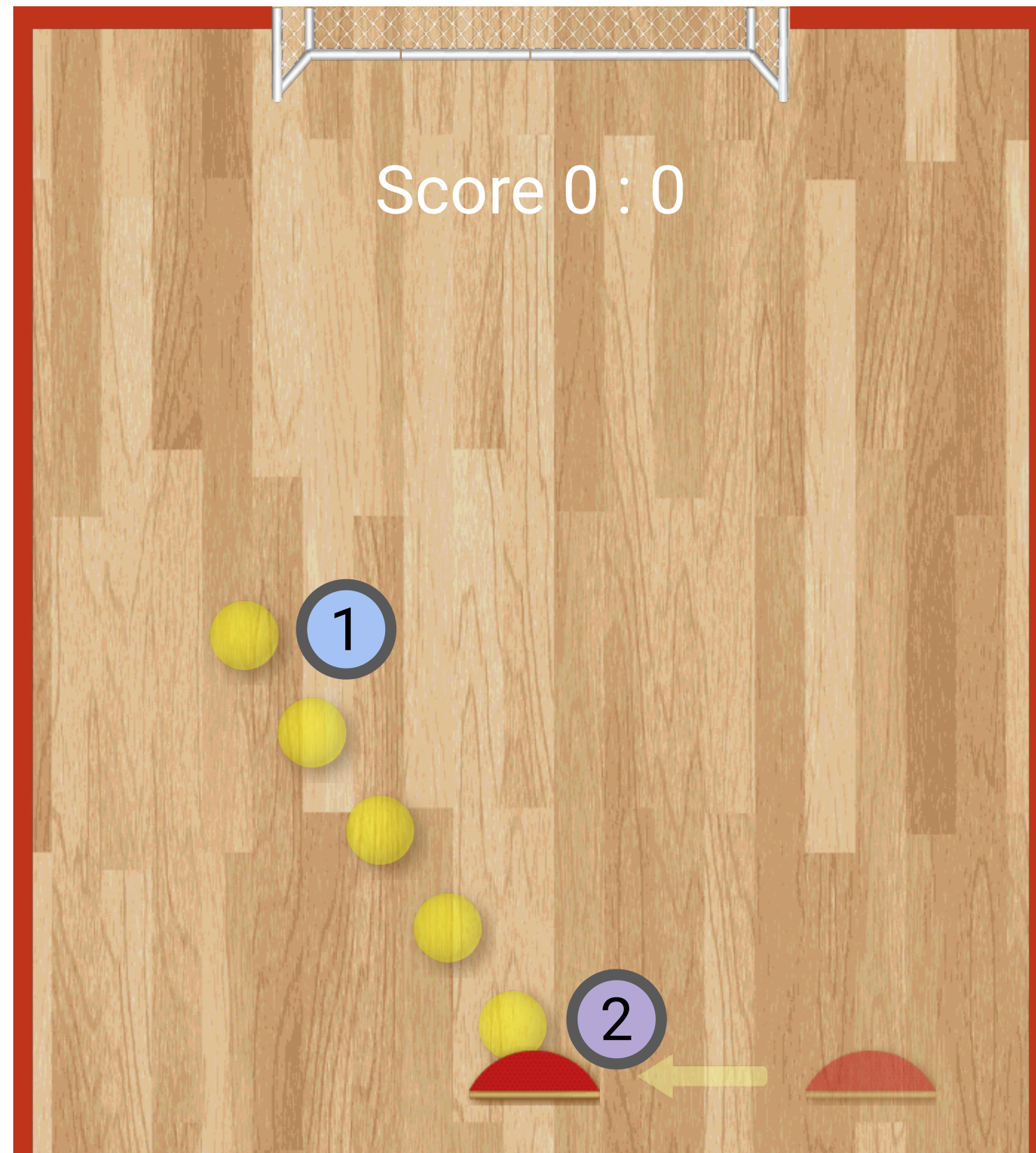
$$\mathcal{J}(\pi, g) = \mathbb{E}_{\mu \sim p(\mu), \tau \sim \pi(\mu)} [g(y \mid \tau)]$$

Sample a program and roll out
exploration policy

Approach: DREAMGRADER



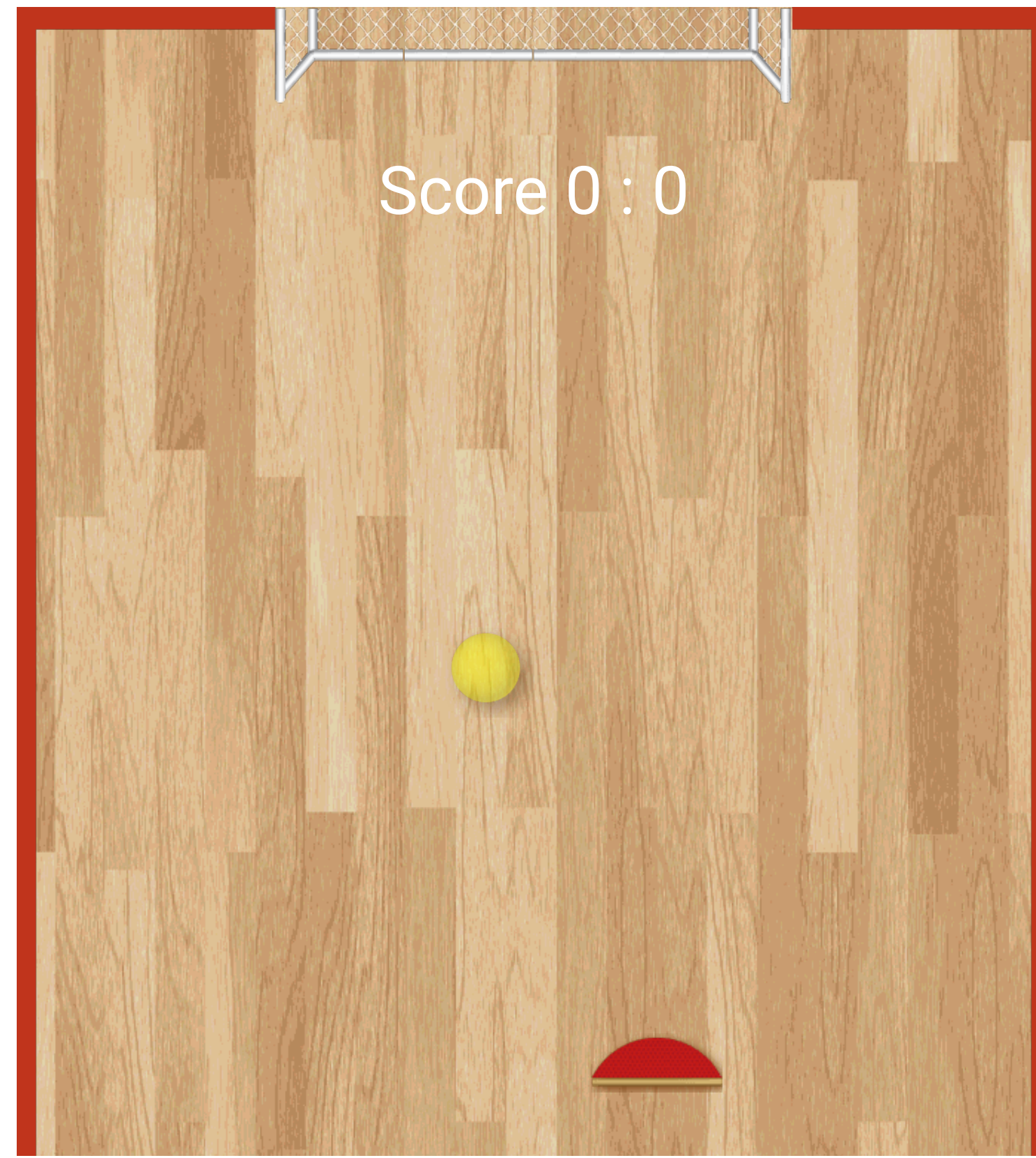
Approach: DREAMGRADER



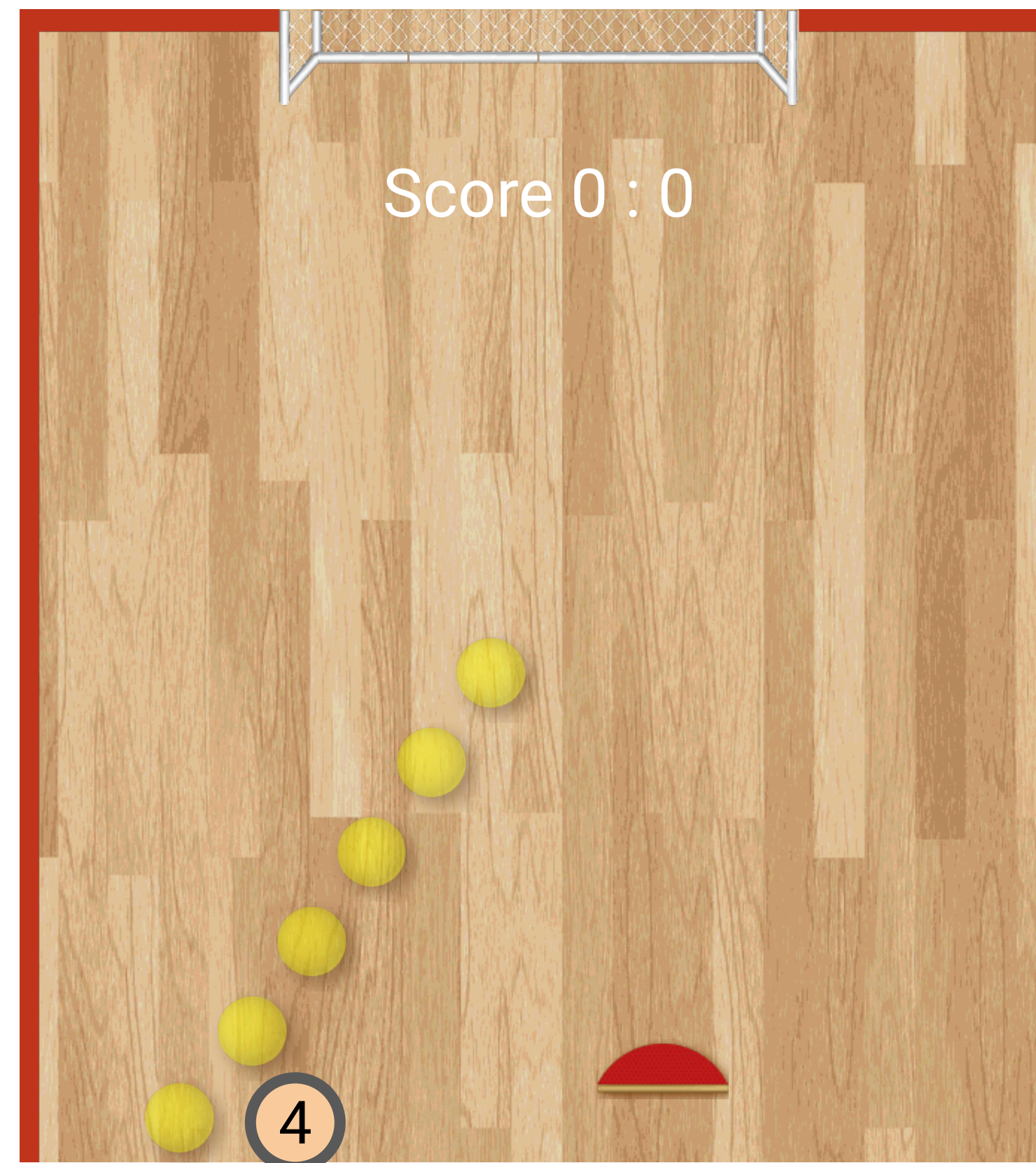
Approach: DREAMGRADER



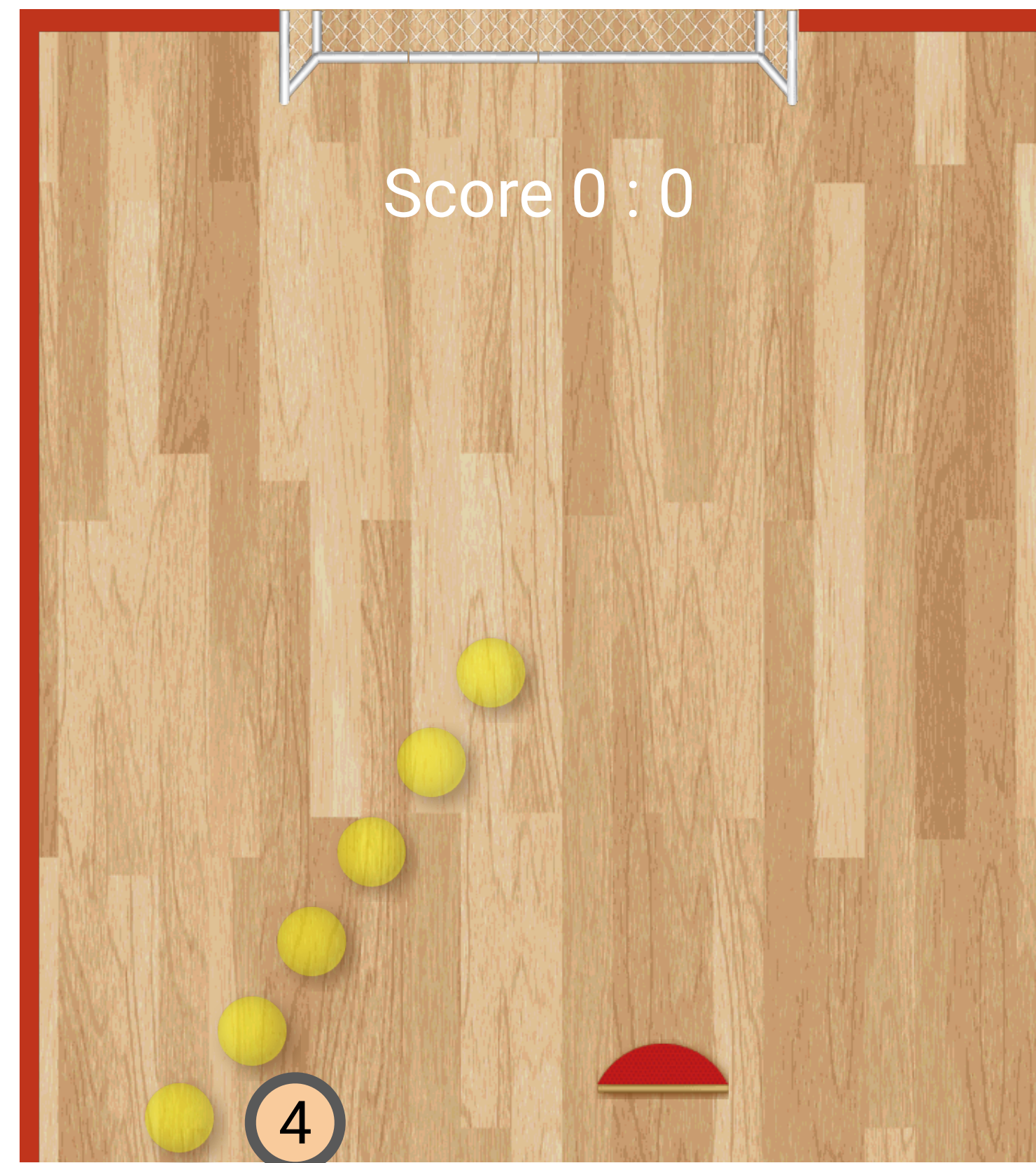
Approach: DREAMGRADER



Approach: DREAMGRADER



Approach: DREAMGRADER

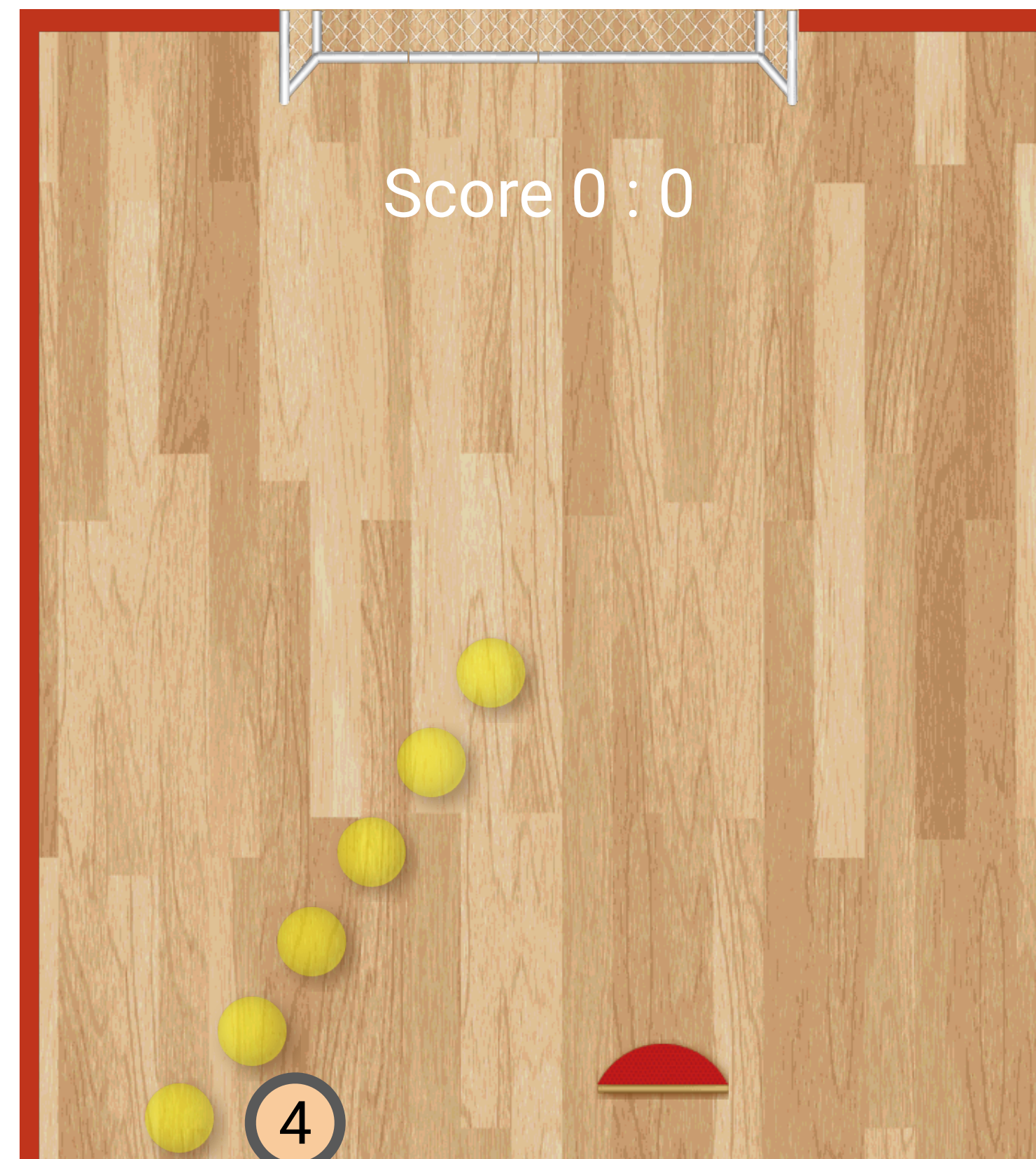


Naive approach

$$\mathcal{J}(\pi, g) = \mathbb{E}_{\mu \sim p(\mu), \tau \sim \pi(\mu)} [\underline{g(y | \tau)}]$$

End-of-episode reward

Approach: DREAMGRADER



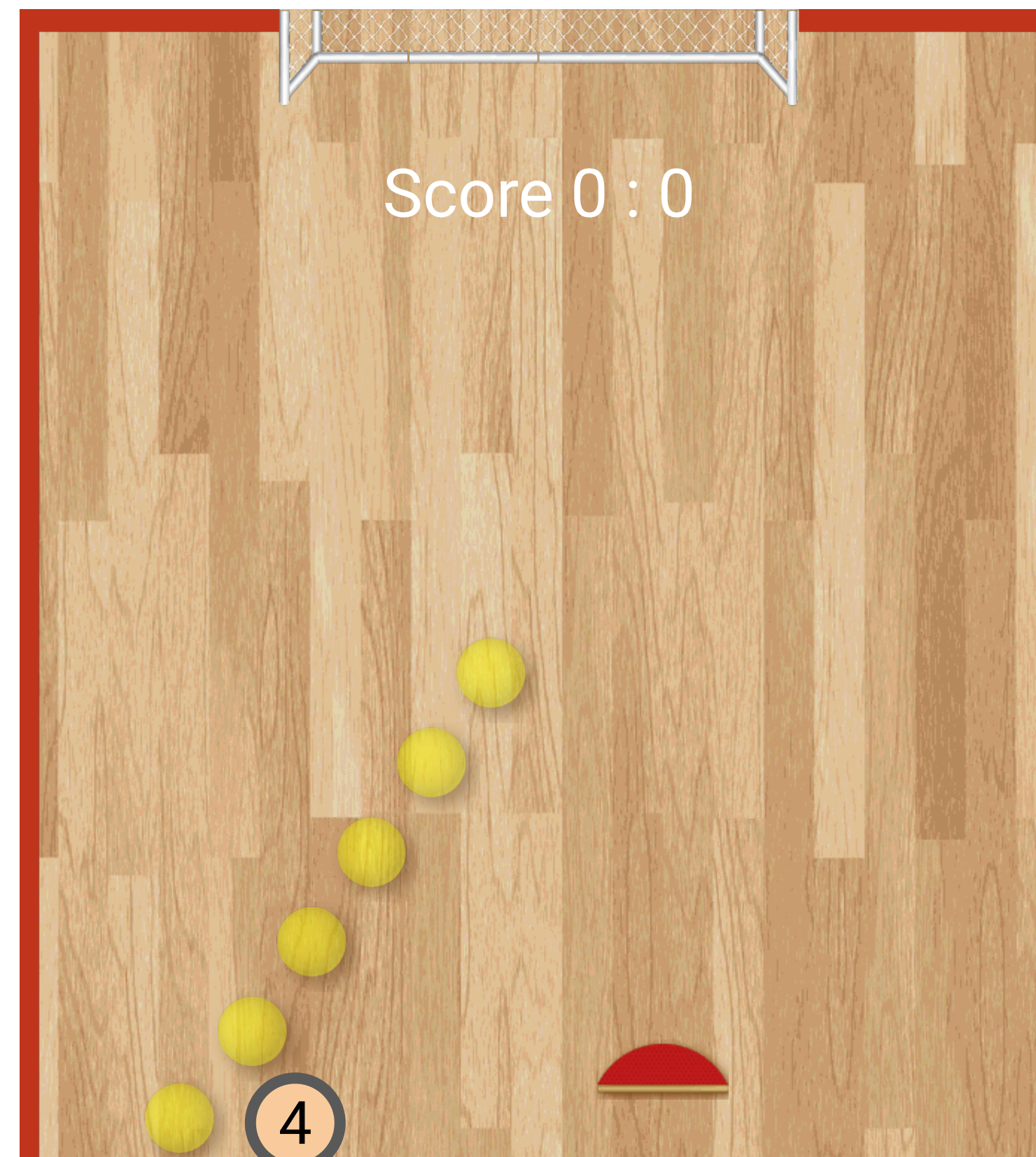
Naive approach

$$\mathcal{J}(\pi, g) = \mathbb{E}_{\mu \sim p(\mu), \tau \sim \pi(\mu)} [g(\mathbf{y} \mid \tau)]$$

End-of-episode reward

⊖ Reward given at 4 but bug discovered at 3

Approach: DREAMGRADER



Naive approach

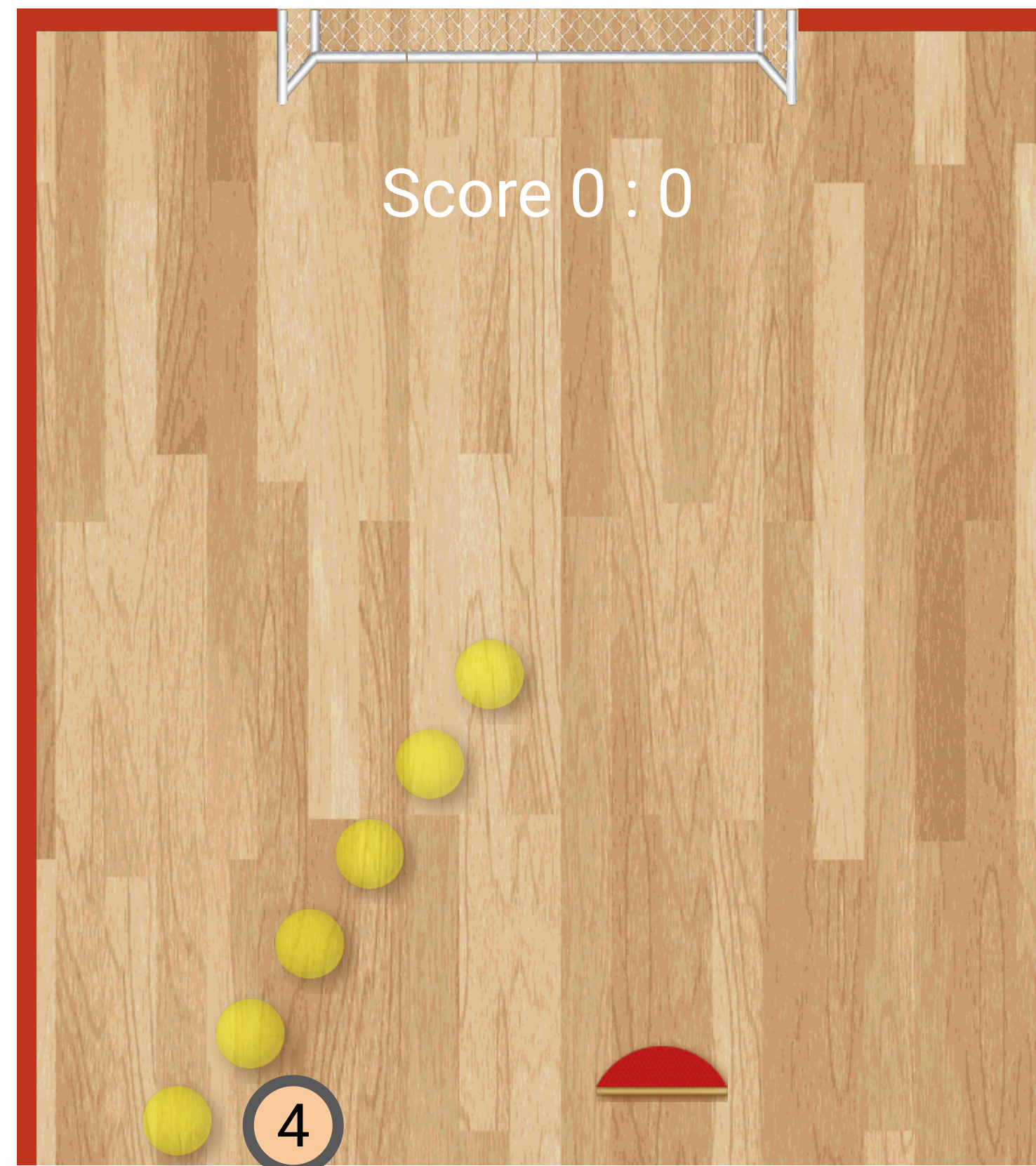
$$\mathcal{J}(\pi, g) = \mathbb{E}_{\mu \sim p(\mu), \tau \sim \pi(\mu)} [g(y | \tau)]$$

End-of-episode reward

— Reward given at 4 but bug discovered at 3

Instead, use DREAM (Liu et al., '21) to provide credit at 3

Approach: DREAMGRADER



Naive approach

$$\mathcal{J}(\pi, g) = \mathbb{E}_{\mu \sim p(\mu), \tau \sim \pi(\mu)} [g(\mathbf{y} \mid \tau)]$$

End-of-episode reward

— Reward given at 4 but bug discovered at 3

Instead, use DREAM (Liu et al., '21) to provide credit at 3

Intuition: maximize information gain

$$r_t = \log \frac{g(\mathbf{y} \mid \tau_{:t+1})}{g(\mathbf{y} \mid \tau_{:t})}$$

Approach: DREAMGRADER

Why does the DREAM meta-RL algorithm apply here?

Few-shot meta-RL:



1) Agent is given new task



2) Agent gets to explore for a few episodes



3) Agent uses exploration to maximize returns on new episode

Approach: DREAMGRADER

Why does the DREAM meta-RL algorithm apply here?

Few-shot meta-RL:



1) Agent is given new task

Agent is given new program



2) Agent gets to explore for a few episodes



3) Agent uses exploration to maximize returns on new episode

Approach: DREAMGRADER

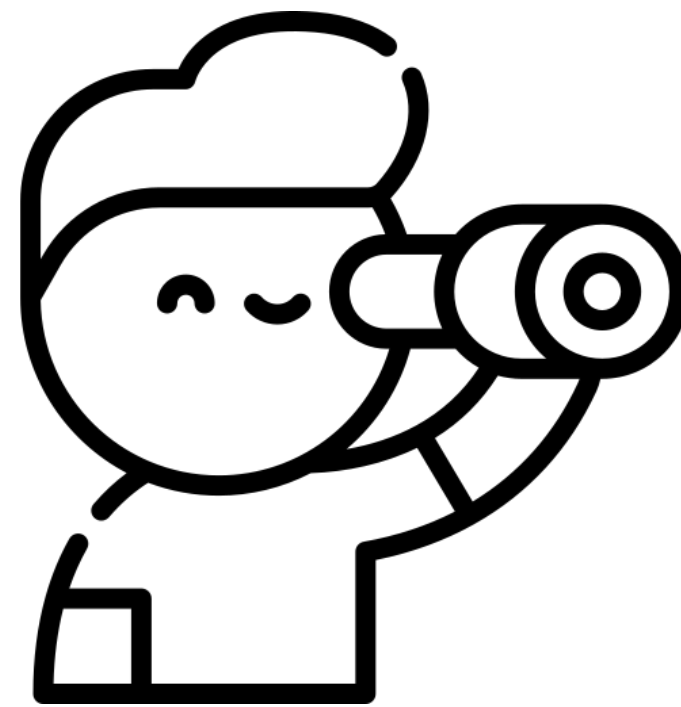
Why does the DREAM meta-RL algorithm apply here?

Few-shot meta-RL:



1) Agent is given new task

Agent is given new program



2) Agent gets to explore for a few episodes

Exploration policy runs to find bugs



3) Agent uses exploration to maximize returns on new episode

Approach: DREAMGRADER

Why does the DREAM meta-RL algorithm apply here?

Few-shot meta-RL:



1) Agent is given new task

Agent is given new program



2) Agent gets to explore for a few episodes

Exploration policy runs to find bugs



3) Agent uses exploration to maximize returns on new episode

Feedback classifier uses exploration to predict label

Experiments: Questions

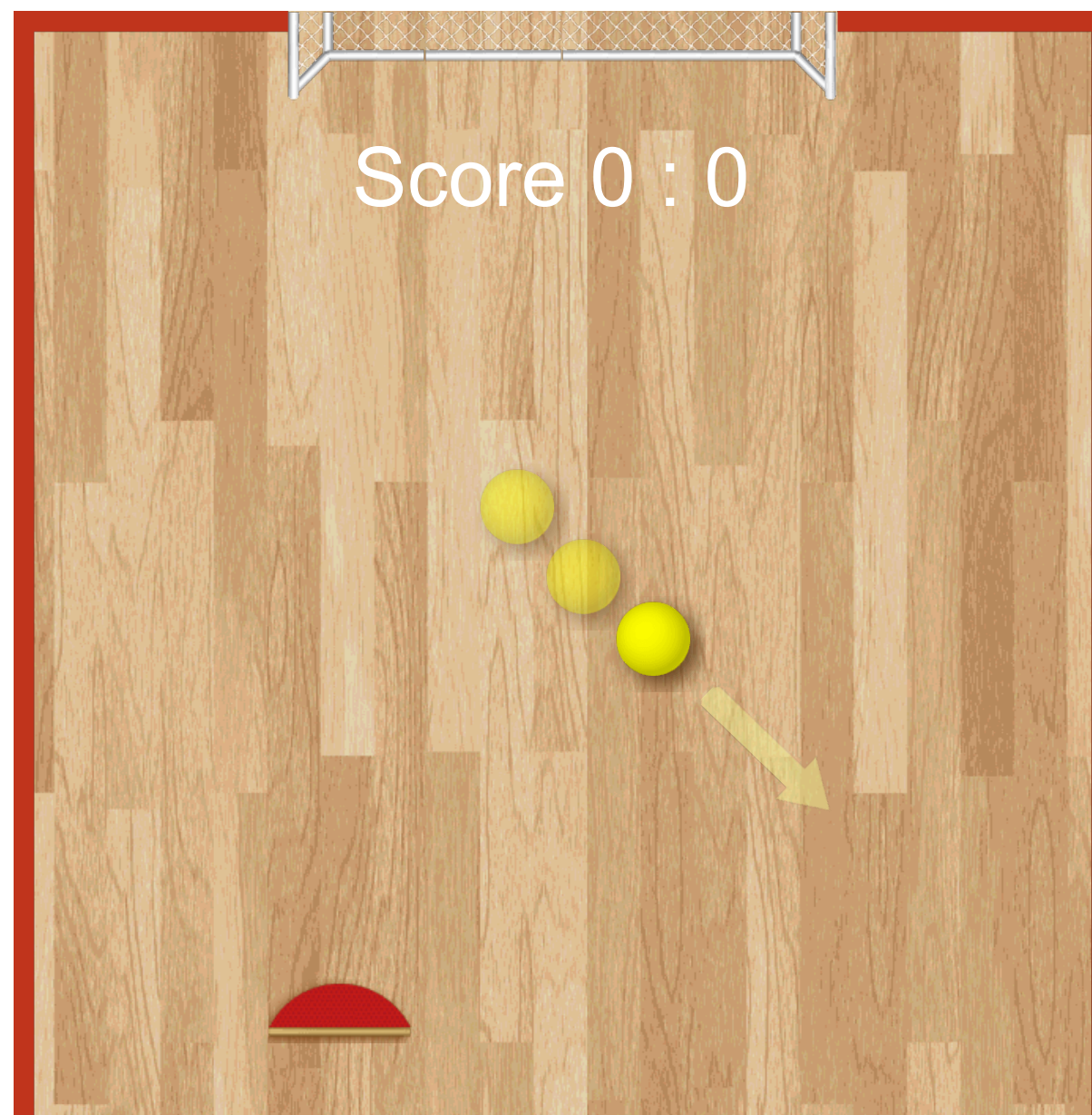
How feasible is automated feedback generation?

Can automated feedback generation handle student creativity?

Experiments: Questions

How feasible is automated feedback generation?

Can automated feedback generation handle student creativity?



Bounce programming assignment from Code.org

Dataset of ~700K real student submissions, released by Nie et al., '21

Train systems on 3,500 programs – hold out the rest

Experiments: Questions

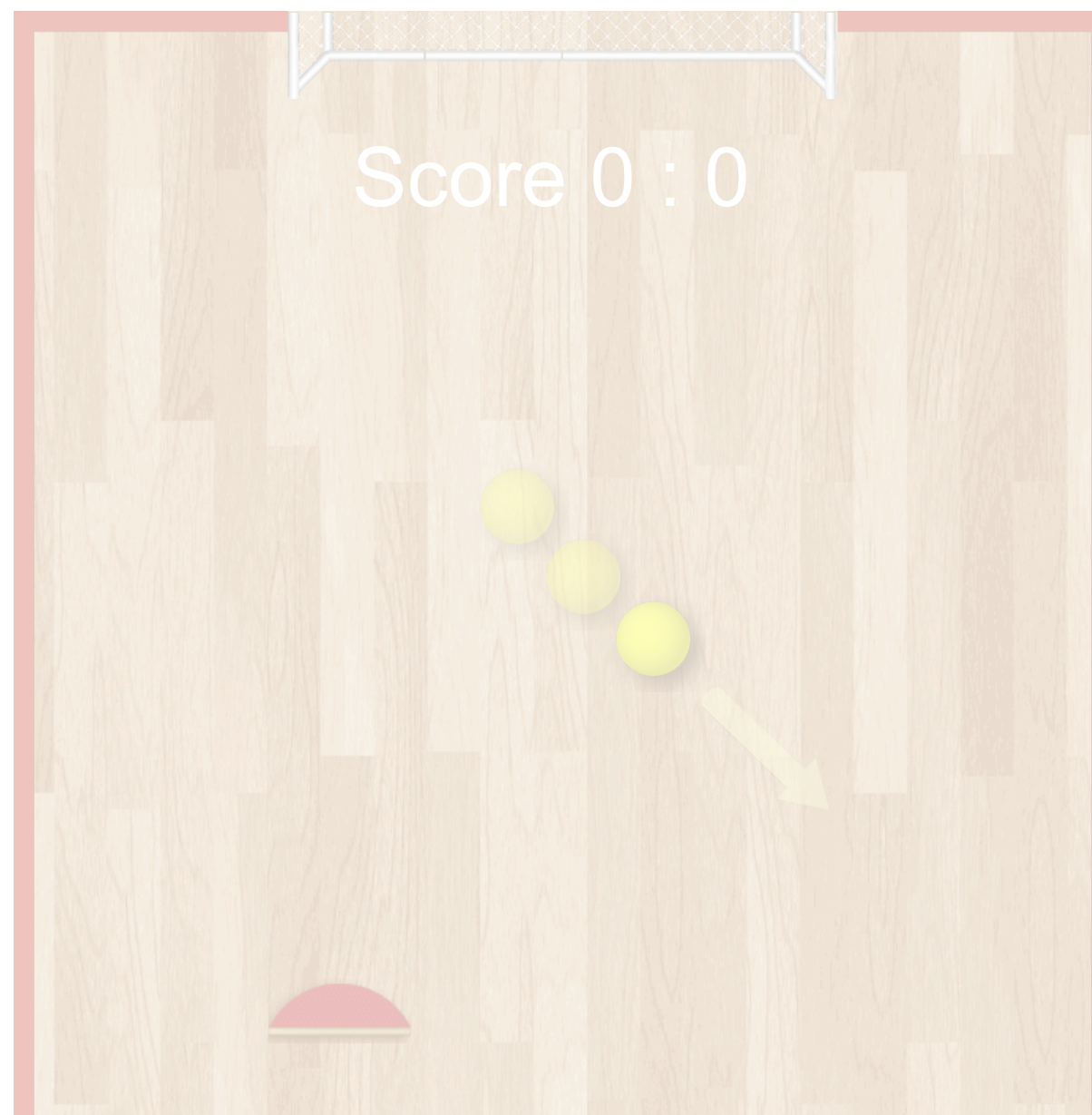
How feasible is automated feedback generation?

Humans

Naive approach of direct maximization

Existing state-of-the-art approach (Nie et al., '21)

Can automated feedback generation handle student creativity?

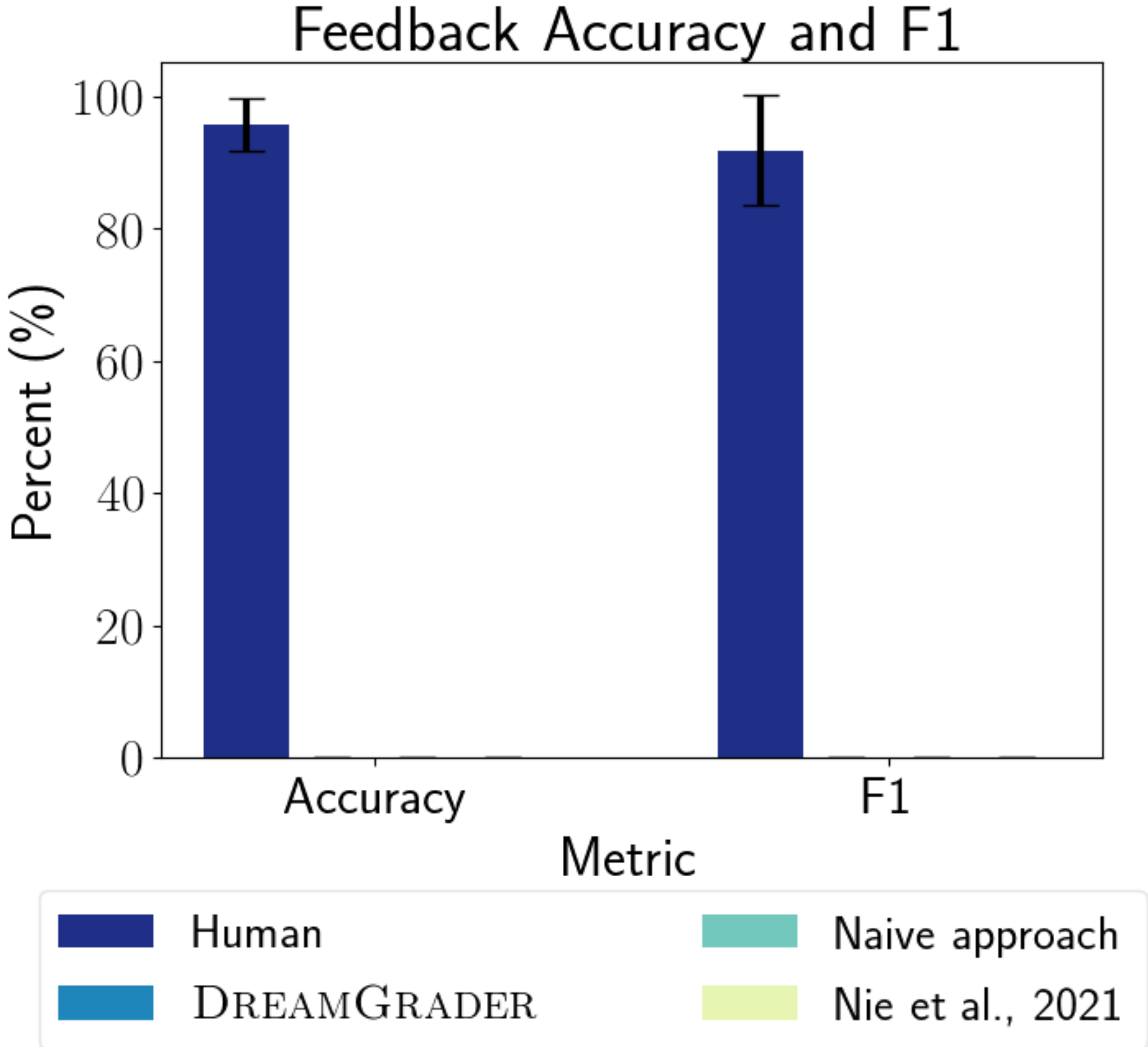


Bounce programming assignment from Code.org

Dataset of ~700K real student submissions, released by Nie et al., '21

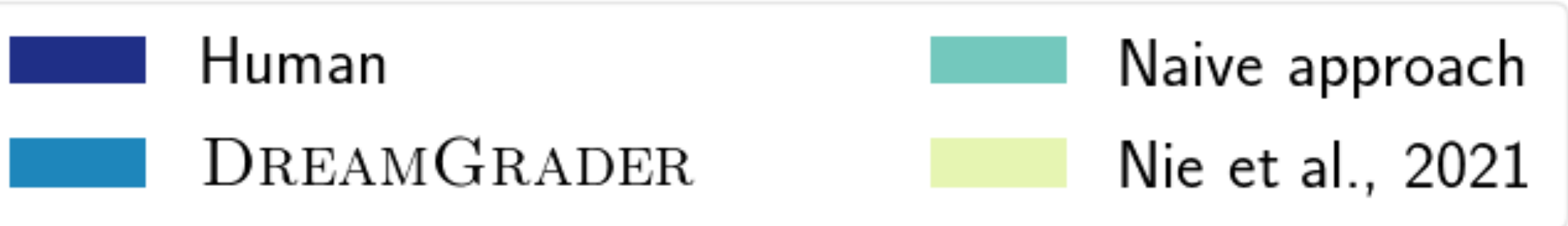
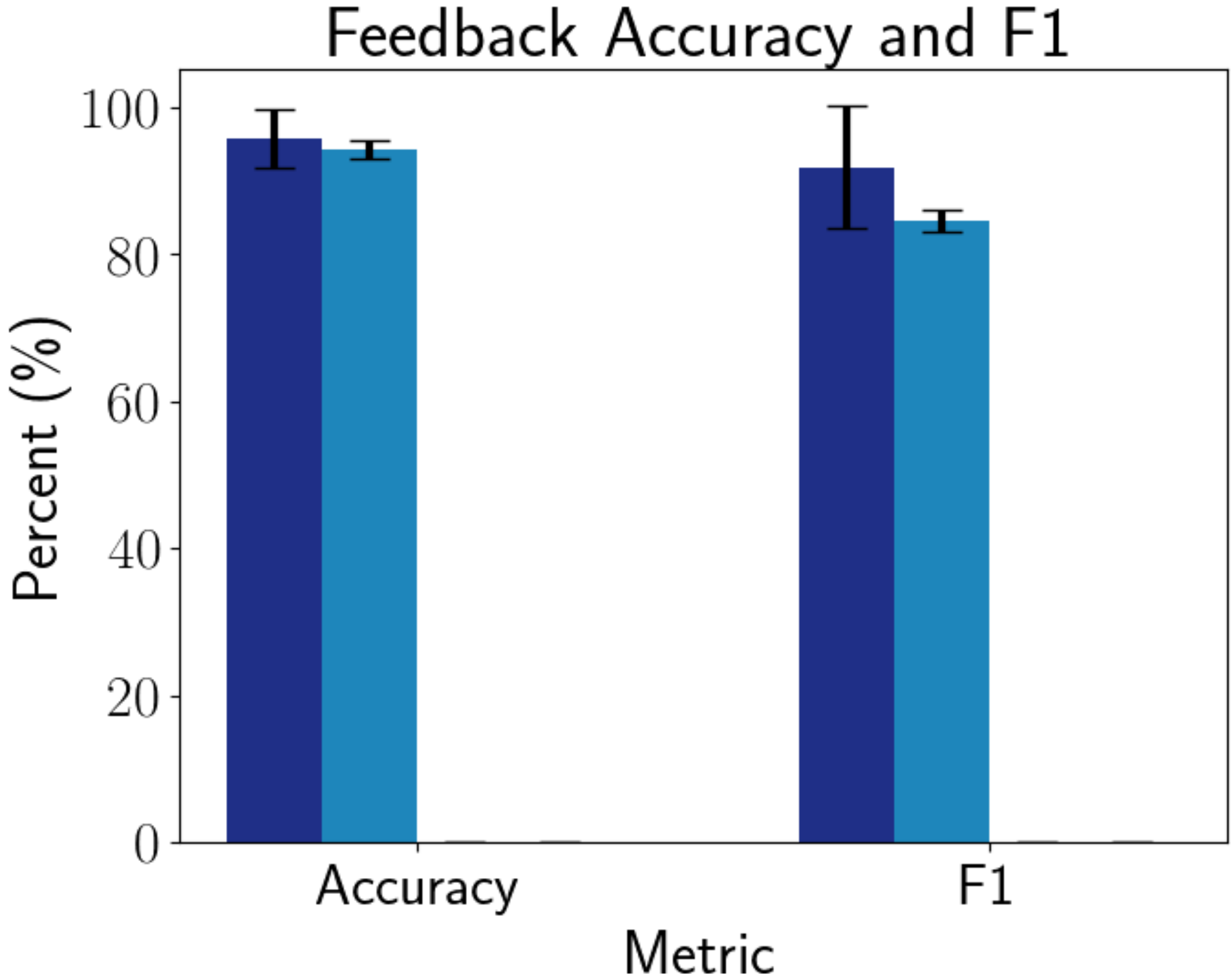
Train systems on 3,500 programs – hold out the rest

Experiments: How Feasible is Automated Feedback?

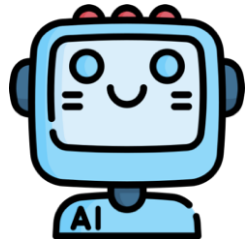


Humans are accurate, but **infeasible**:
Requires ~4 years to grade the dataset

Experiments: How Feasible is Automated Feedback?

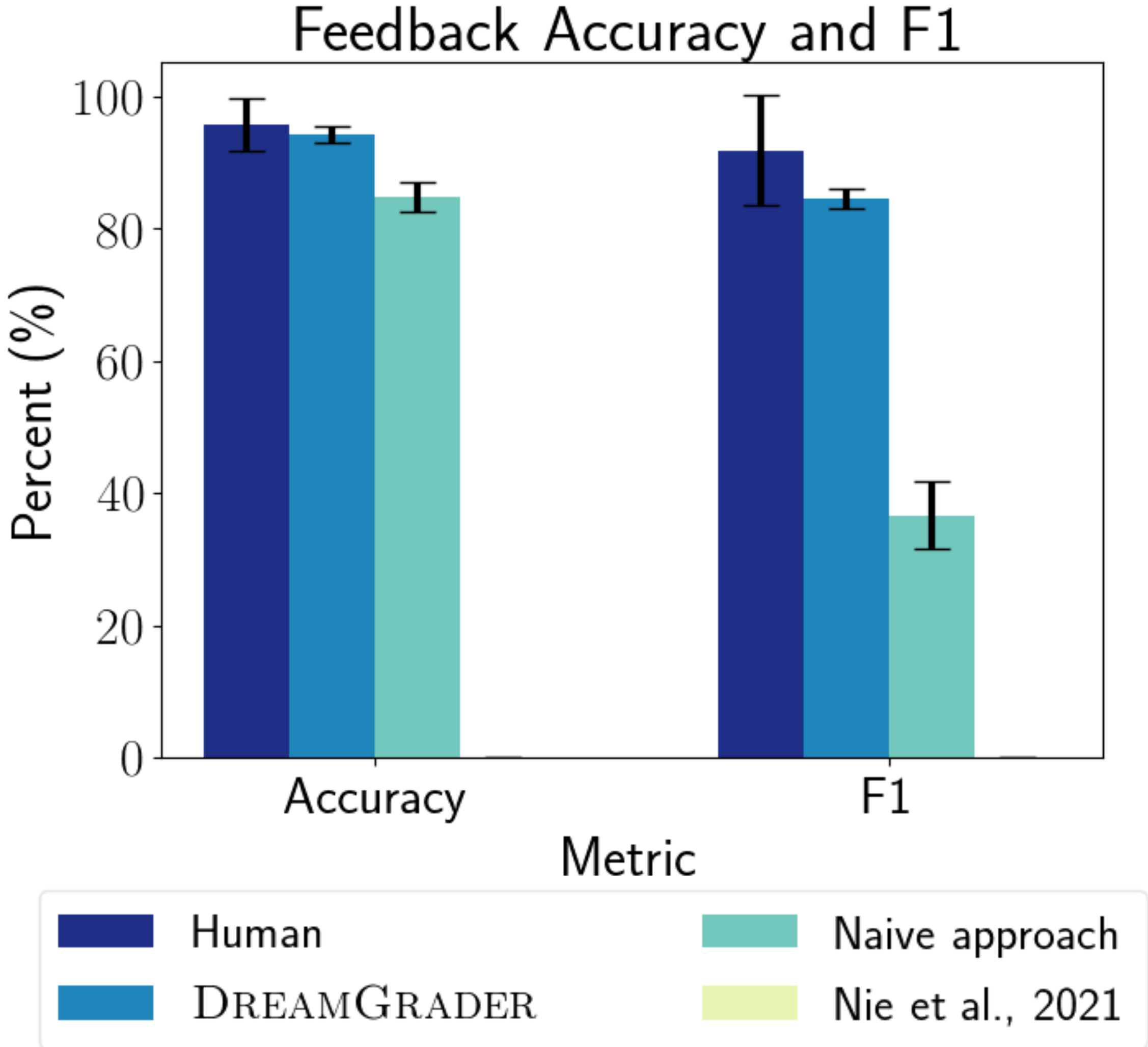


Humans are accurate, but **infeasible**:
Requires ~4 years to grade the dataset

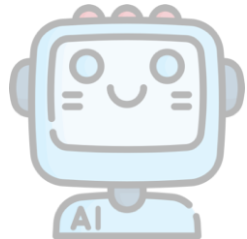


DREAMGRADER achieves within **1.5%** of
human accuracy

Experiments: How Feasible is Automated Feedback?



Humans are accurate, but **infeasible**:
Requires ~4 years to grade the dataset

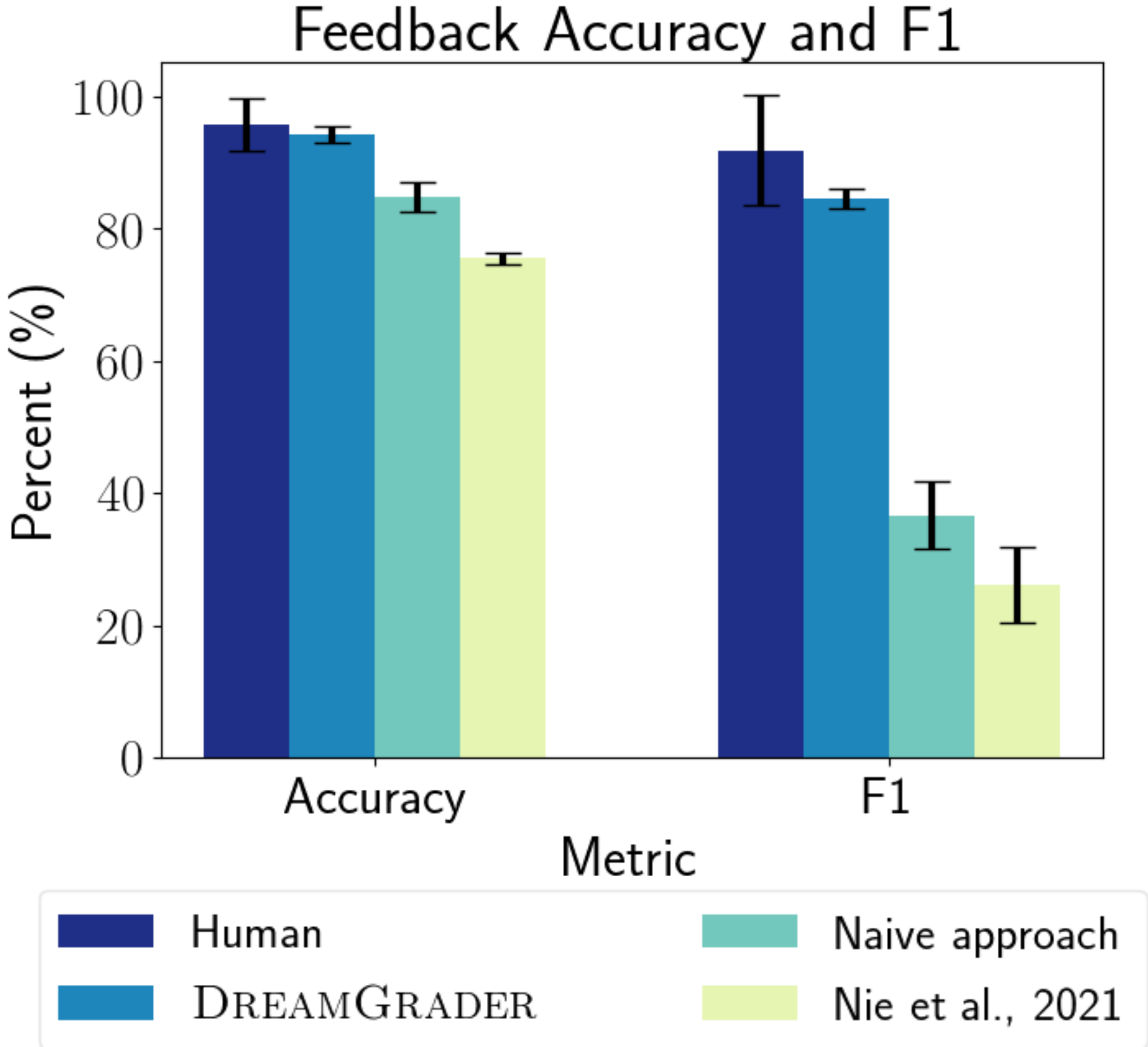


DREAMGRADER achieves within **1.5%** of
human accuracy

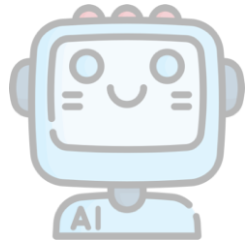


Appropriate credit assignment is **critical** for
learning effective exploration

Experiments: How Feasible is Automated Feedback?



Humans are accurate, but **infeasible**:
Requires ~4 years to grade the dataset



DREAMGRADER achieves within **1.5%** of human accuracy



Appropriate credit assignment is **critical** for learning effective exploration



Improves on existing methods by **18.8%**

Experiments: Learned Exploration Behavior

```
Underlying env ID: 7340  
Env ID: 1  
Label: [1 1 0 0 0 0 1 0 0 1 0 1 0 1 1]  
Binary label: whenGoal-noBallLaunch  
Action: None  
Reward: 0  
Timestep: 0  
Exploration reward: 0.020  
Prob: 0.456
```

```
Underlying env ID: 4843  
Env ID: 0  
Label: [0 1 0 0 0 0 0 0 0 0 0 0 0 1 1]  
Binary label: whenMiss-noBallLaunch  
Action: None  
Reward: 0  
Timestep: 0  
Exploration reward: 0.005  
Prob: 0.507
```

```
Underlying env ID: 2732  
Env ID: 1  
Label: [0 1 0 1 1 0 0 1 0 0 1 1 0 0 1]  
Binary label: whenWall-illegal-moveRight  
Action: None  
Reward: 0  
Timestep: 0  
Exploration reward: 0.079  
Prob: 0.331
```

What happens when...

the ball hits the goal?

the ball hits the floor?

the ball hits the wall?

Main gap with humans appears in these sorts of programs with many balls

Experiments: Can We Handle Some Student Creativity?

One type of student creativity in the dataset: ball and paddle speed

- Test handling student creativity by evaluating on held out ball and paddle speeds

	Both held out	Held out ball speed	Held out paddle speed	Neither held out
Accuracy	88.0%	88.8%	88.2%	88.4%
Precision	38.8%	41.6%	44.9%	38.6%
Recall	82.1%	87.2%	91.4%	85.6%
F1	52.8%	56.3%	60.2%	53.2%



Performance on held out speeds roughly matches speeds seen during training

Bonus Experiment: Beyond Code.org bounce game?

- Stanford CS106A: Students program Breakout in homework assignment
- Ball “skewering” bug: common mistake, most difficult to detect/grade

Learned exploration policy



Outline

Brief Recap on Meta-RL

Algorithms for Learning to Explore

End-to-End Optimization of Exploration Strategies

Alternative Decoupled Exploration Strategies

Decoupled but Consistent Exploration & Exploitation

Case Study: Applying Meta-RL to CS Education

Reminders

Homework 3 due **tonight**
(and HW4 out today)

Project milestone due **next Wednesday**

Next week: Can we make reinforcement learning more autonomous?
Can RL agents discover skills themselves?
Can we do hierarchical RL?