

Decoupling Exploration and Exploitation for Meta-Reinforcement Learning without Sacrifices



Evan Zheran Liu



Aditi Raghunathan



Percy Liang



Chelsea Finn

The world constantly **changes**

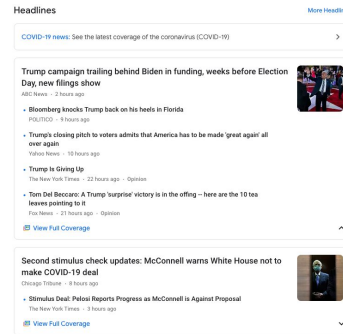
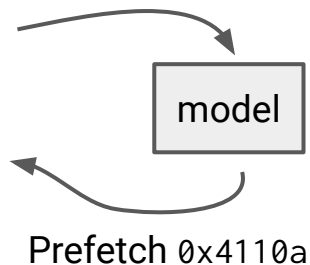


New kitchens / recipes

```
add $1, %rdi
mov %edx, %eax
shr $8, %rdx
xor -1(%rdi), %al
movzx %al, %eax
xor 0x4110a(, %rax, 8), %rdx
cmp %rcx, %rdi
```

Code

New programs



New users / preferences

The world constantly **changes**

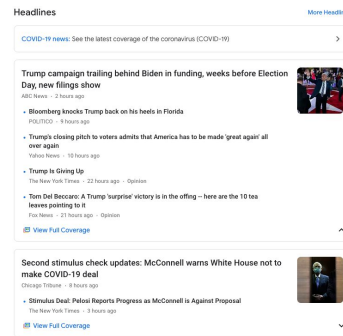
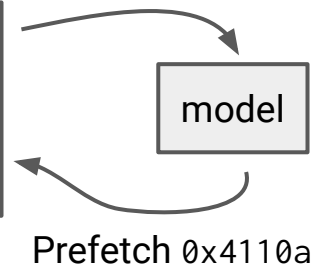


New kitchens / recipes

```
add $1, %rdi
mov %edx, %eax
shr $8, %rdx
xor -1(%rdi), %al
movzx %al, %eax
xor 0x4110a(, %rax, 8), %rdx
cmp %rcx, %rdi
```

Code

New programs



New users / preferences

Re-learning from scratch is **expensive** and **sample inefficient**

Goal: quickly learn new tasks

How? Leverage *prior experience*

Goal: quickly learn new tasks



meta-training

How? Leverage *prior experience*

Goal: quickly learn new tasks



meta-training

How? Leverage *prior experience*



meta-testing

Goal: quickly learn new tasks



meta-training

How? Leverage *prior experience*



explore



Goal: quickly learn new tasks



meta-training

How? Leverage *prior experience*



exploit

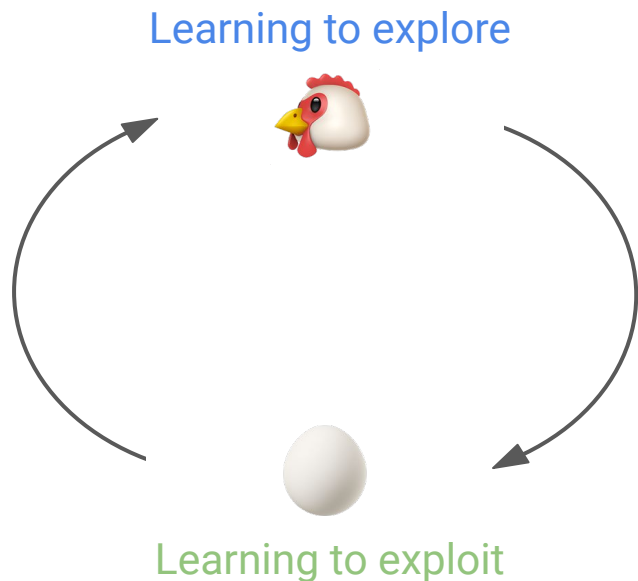


End-to-End Meta-Reinforcement Learning

Natural approach: Optimize **exploration** and **exploitation** **end-to-end** to maximize returns

End-to-End Meta-Reinforcement Learning

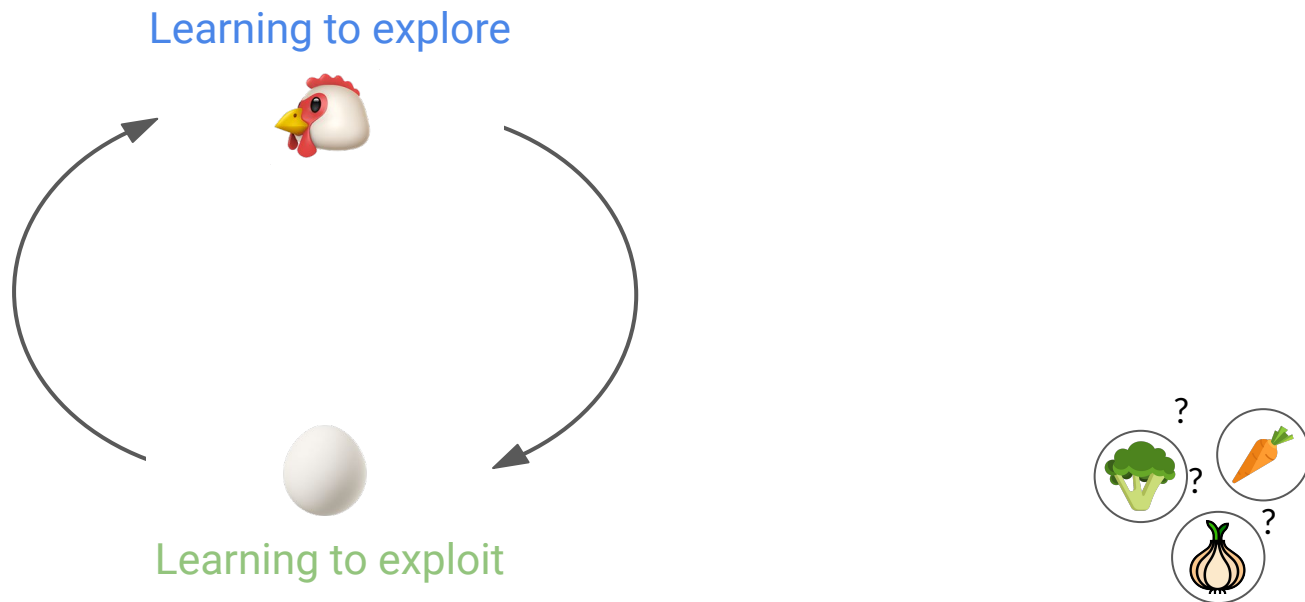
Natural approach: Optimize **exploration** and **exploitation** **end-to-end** to maximize returns



Coupling problem: learning exploration and exploitation depend on each other

End-to-End Meta-Reinforcement Learning

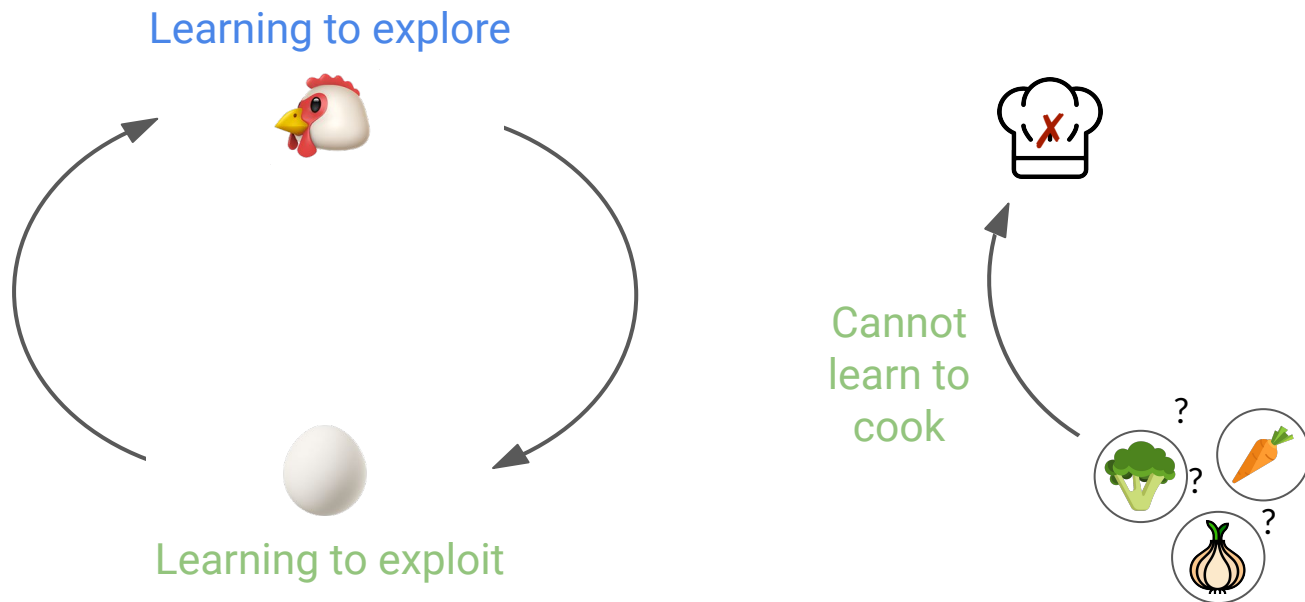
Natural approach: Optimize **exploration** and **exploitation** **end-to-end** to maximize returns



Coupling problem: learning exploration and exploitation depend on each other

End-to-End Meta-Reinforcement Learning

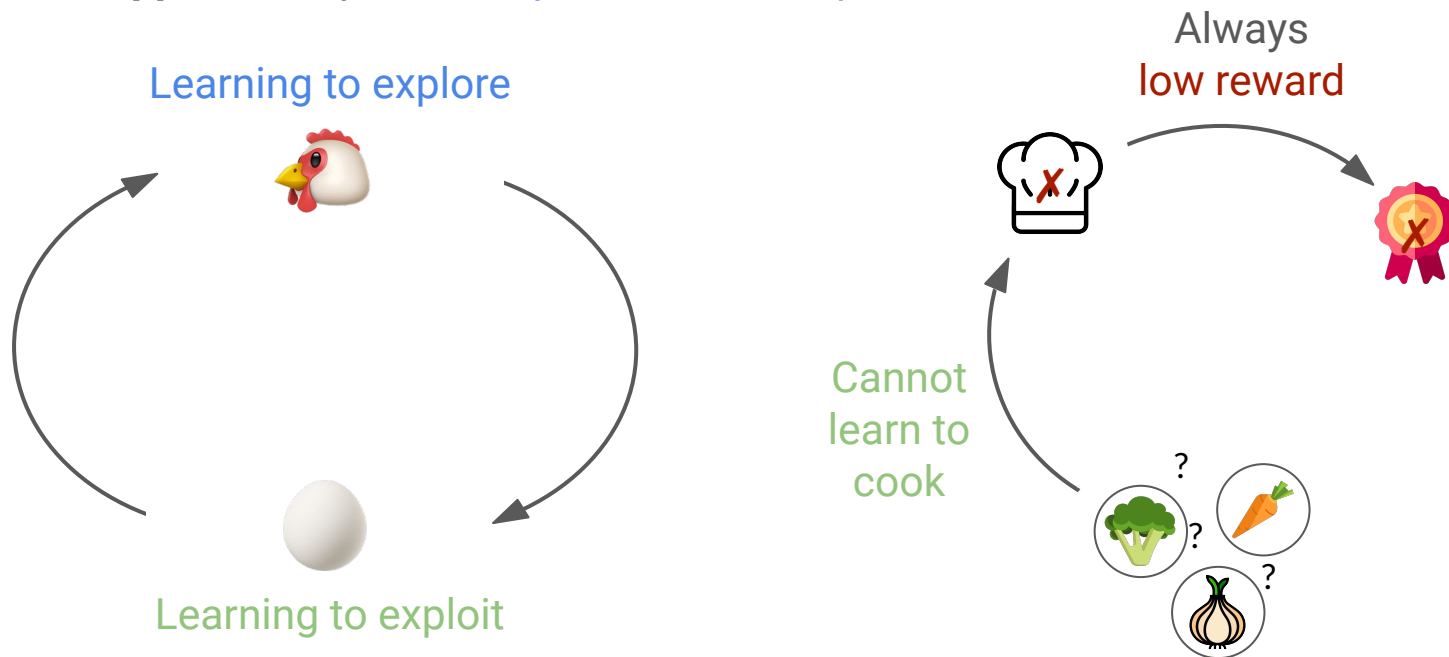
Natural approach: Optimize **exploration** and **exploitation** **end-to-end** to maximize returns



Coupling problem: learning exploration and exploitation depend on each other

End-to-End Meta-Reinforcement Learning

Natural approach: Optimize **exploration** and **exploitation** **end-to-end** to maximize returns

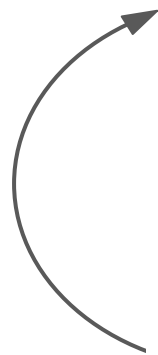


Coupling problem: learning exploration and exploitation depend on each other

End-to-End Meta-Reinforcement Learning

Natural approach: Optimize **exploration** and **exploitation** **end-to-end** to maximize returns

Learning to explore



Learning to exploit



Always
low reward



Cannot
learn to
cook

Cannot learn
to find
ingredients

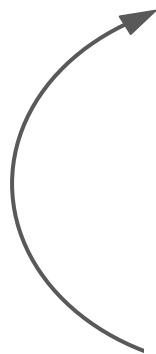


Coupling problem: learning exploration and exploitation depend on each other

End-to-End Meta-Reinforcement Learning

Natural approach: Optimize **exploration** and **exploitation** **end-to-end** to maximize returns

Learning to explore



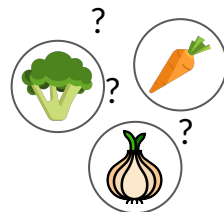
Learning to exploit

Always
low reward



Cannot
learn to
cook

Cannot learn
to find
ingredients



Coupling problem: learning exploration and exploitation depend on each other

... leads to **local optima** and **poor sample complexity**

DREAM Overview

Goal: Create **exploration** objective to *all* and *only* recover **task-relevant** information

DREAM Overview

Goal: Create **exploration** objective to *all* and *only* recover **task-relevant** information

Key (mild) assumption: can distinguish all meta-training tasks from each other with a unique problem ID

DREAM Overview

Goal: Create **exploration** objective to *all* and *only* recover **task-relevant** information

Meta-training

1) Learn **exploitation** & identify task-relevant information

|
|
|
|
|
|
|
|
|

DREAM Overview

Goal: Create **exploration** objective to *all* and *only* recover **task-relevant** information

Meta-training

1) Learn **exploitation** & identify task-relevant information

2) Learn to **explore** by recovering that information

DREAM Overview

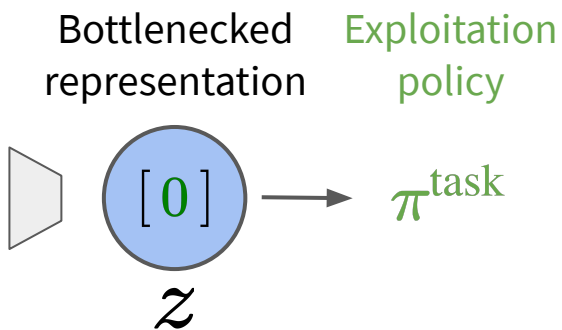
Goal: Create **exploration** objective to *all* and *only* recover **task-relevant** information

Meta-training

1) Learn **exploitation** & identify task-relevant information

Problem ID μ

wall color	2
ingredients	0
decorations	1



2) Learn to **explore** by recovering that information

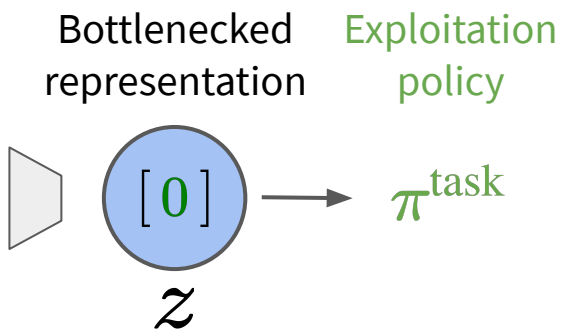
DREAM Overview

Goal: Create **exploration** objective to *all* and *only* recover **task-relevant** information

Meta-training

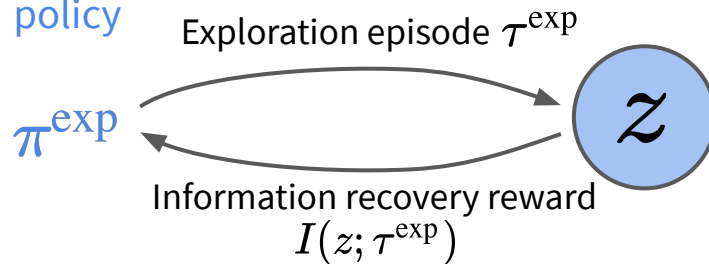
1) Learn **exploitation** & identify task-relevant information

Problem ID μ
wall color $\begin{bmatrix} 2 \\ 0 \\ 1 \end{bmatrix}$
ingredients
decorations

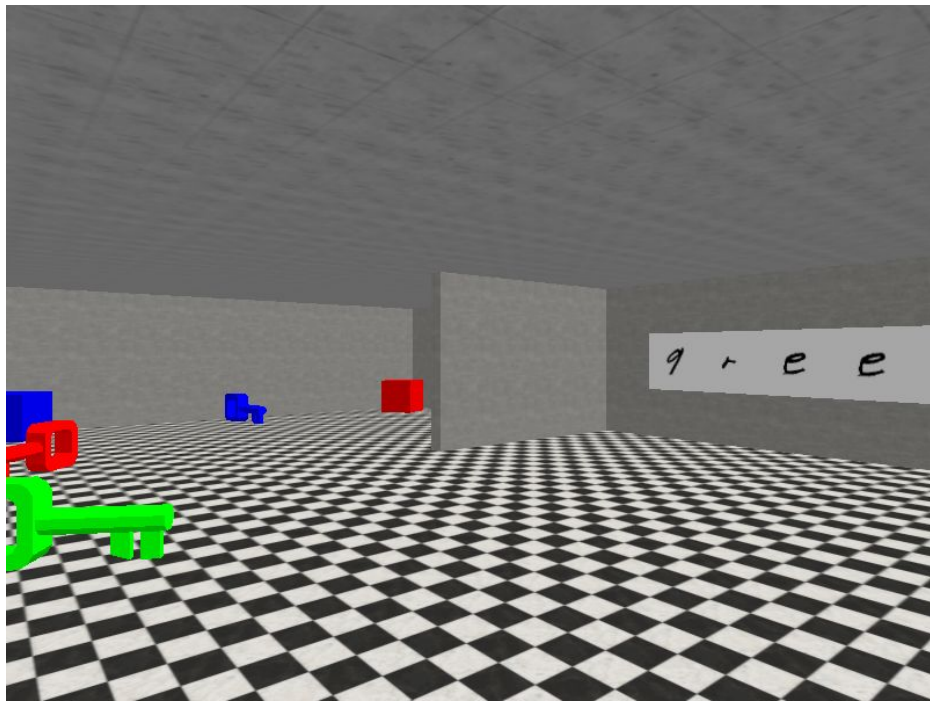


2) Learn to **explore** by recovering that information

Exploration policy π^{exp}



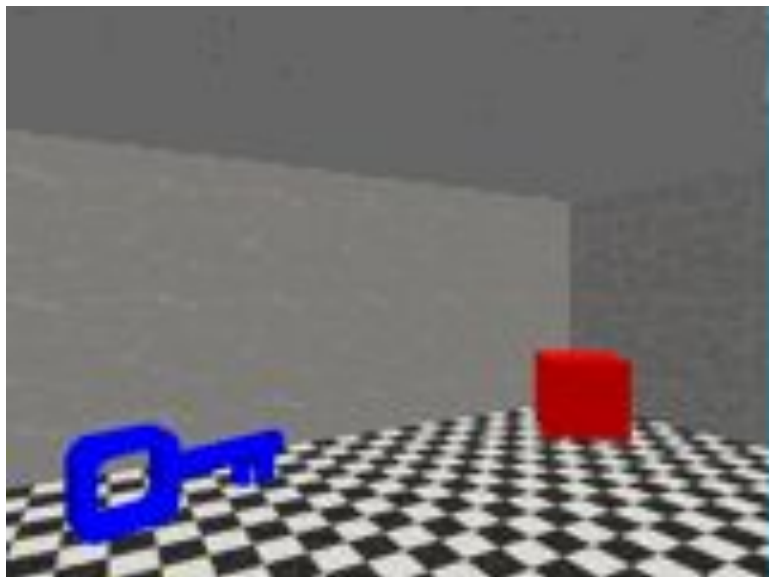
Experiments: Sparse Reward 3D Visual Navigation



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

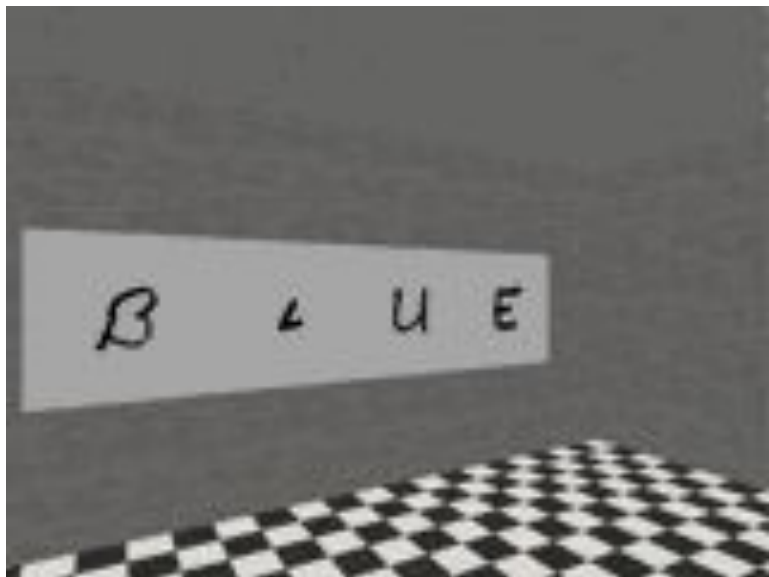
- Task: go to the *goal* = key / block, color specified by the sign
- Agent starts on other side of barrier and must walk around to read the sign
- Pixels observations (80 x 60 RGB)
- Sparse binary reward
- Existing benchmarks don't typically use **pixel observations** and **sparse rewards**

Experiments: Qualitative Results for DREAM

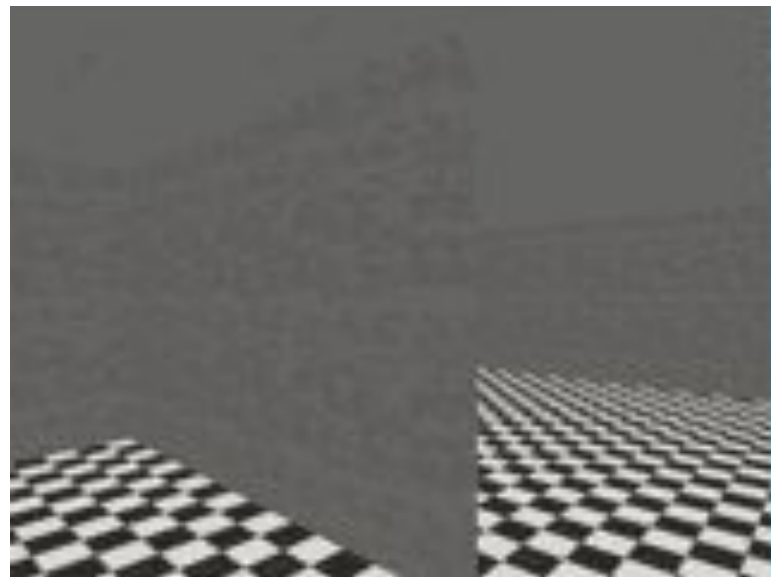


DREAM exploration

Experiments: Qualitative Results for DREAM

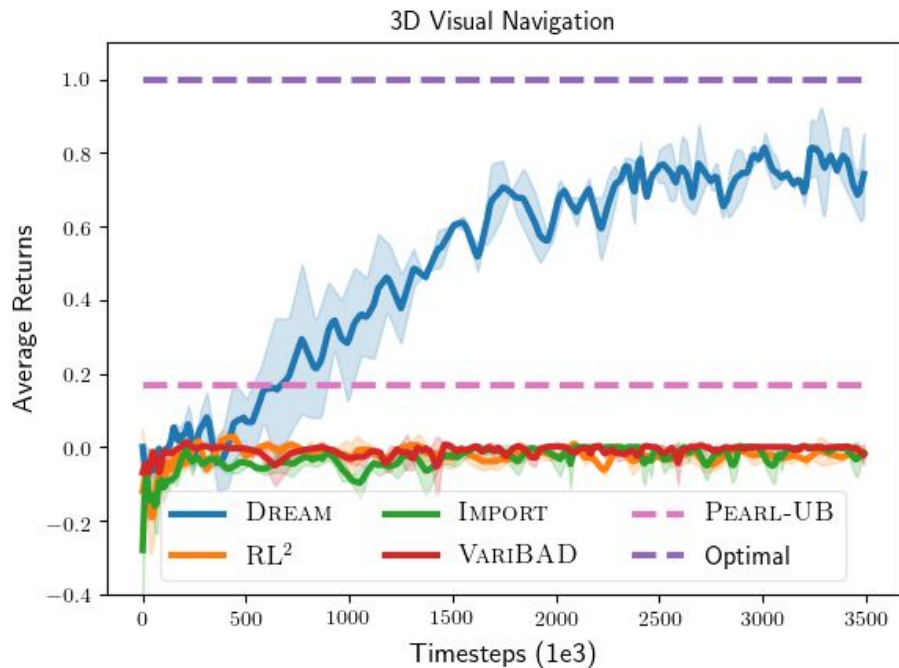


DREAM exploration



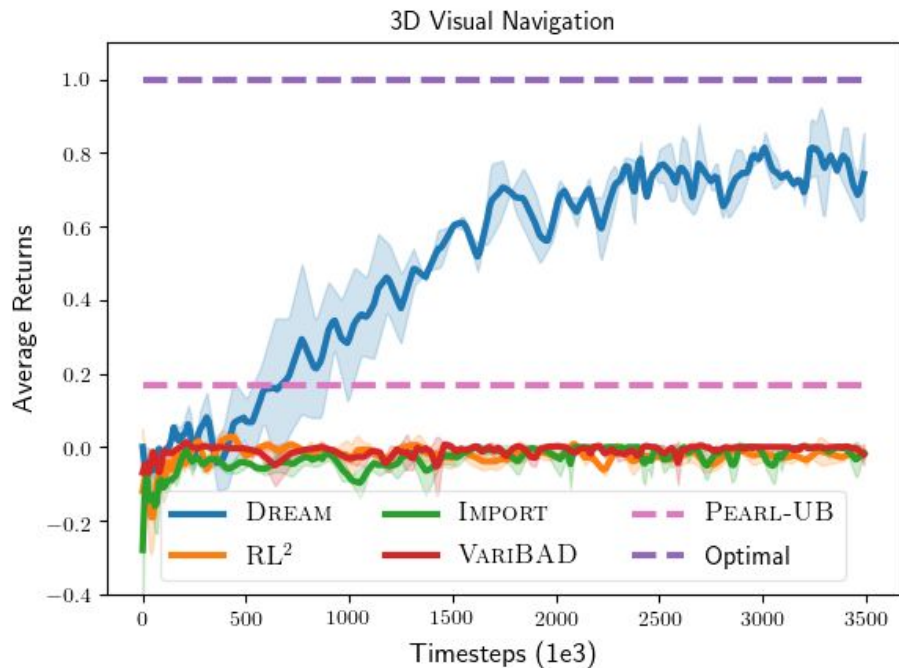
DREAM exploitation: Go to key

Experiments: Quantitative Results



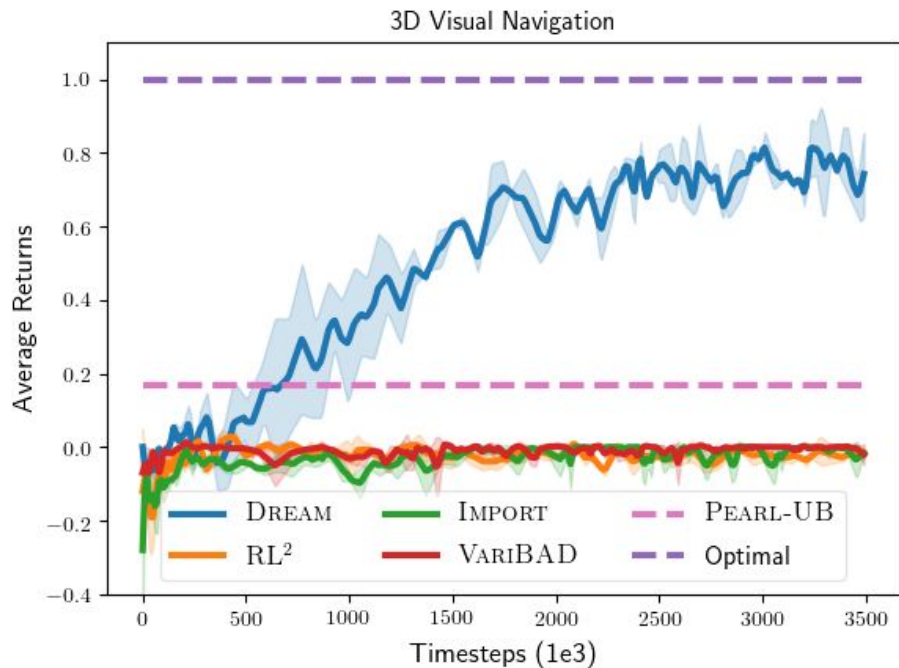
- Only DREAM scales to high-dimensional states and sparse rewards

Experiments: Quantitative Results



- Only DREAM scales to high-dimensional states and sparse rewards
- End-to-end approaches achieve zero returns due to **coupling problem**

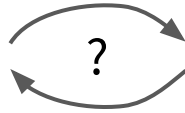
Experiments: Quantitative Results



- Only DREAM scales to high-dimensional states and sparse rewards
- End-to-end approaches achieve *zero* returns due to **coupling problem**
- Decoupled approaches, e.g., Thompson Sampling do not learn the optimal exploration strategy

Takeaways

- I. **Coupling** between exploration and exploitation prevents existing **end-to-end** methods from solving tasks with challenging exploration



Takeaways

- I. **Coupling** between exploration and exploitation prevents existing **end-to-end** methods from solving tasks with challenging exploration



- II. **DREAM** provides separate exploration and exploitation objectives that avoid the **coupling** problem

