

Learning to Synthesize Programs as Interpretable and Generalizable Policies

Dweep Trivedi*, Jesse Zhang*, Shao-Hua Sun*, Joseph J Lim

University of Southern California



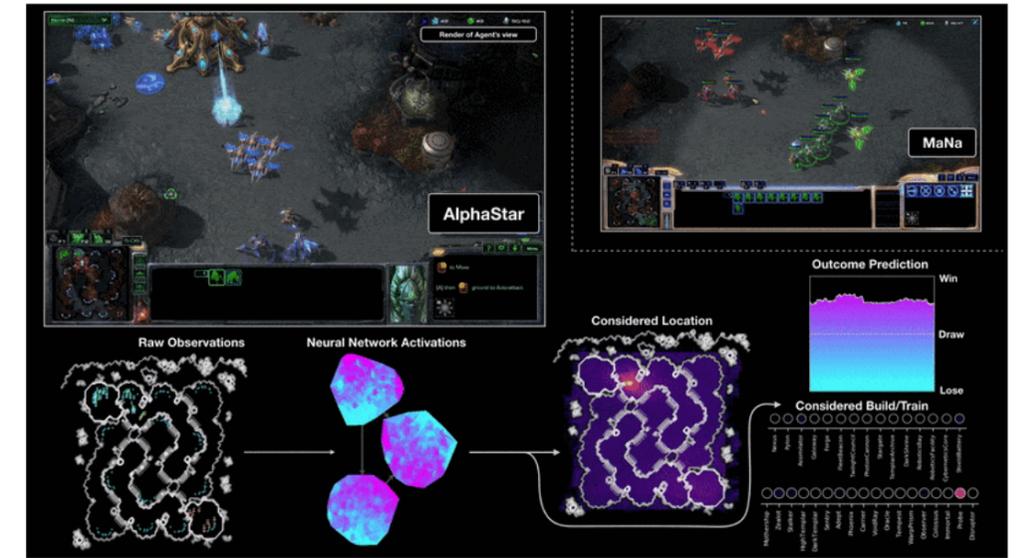
Advances in Deep Reinforcement Learning



Autonomous Driving

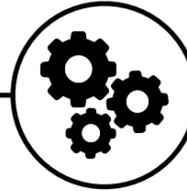


Robotics



Game AI

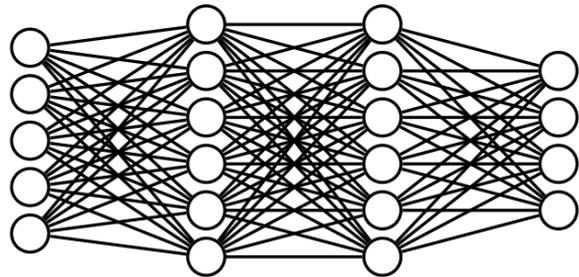
Execute



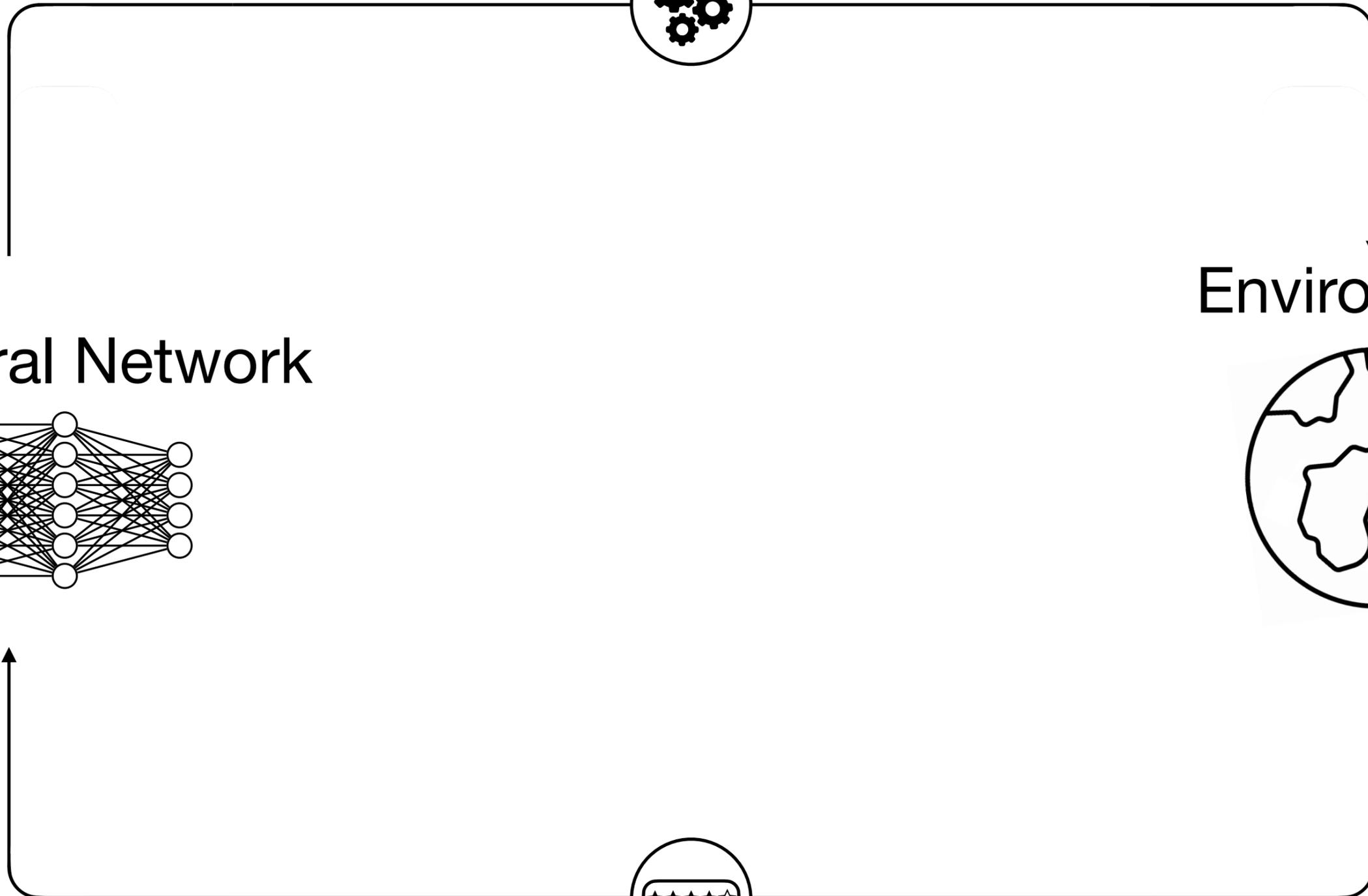
Environment



Deep Neural Network

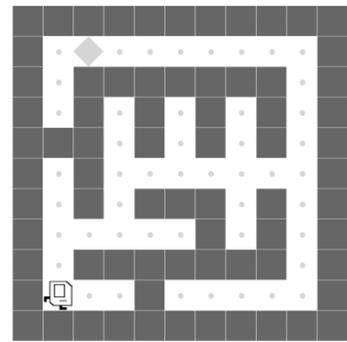


Reward

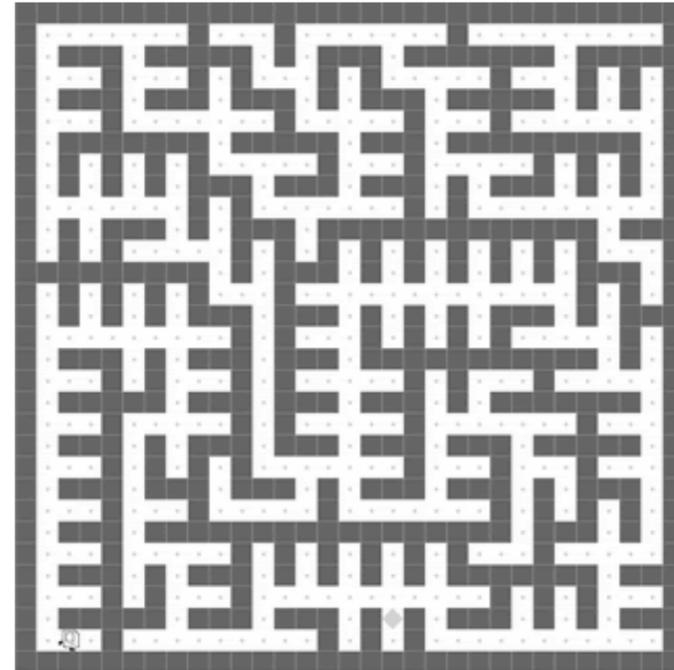


Issues with Deep Reinforcement Learning (DRL)

Generalization



Simple task



Complex task

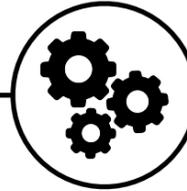
Interpretability

Trust

Safety

Contestability

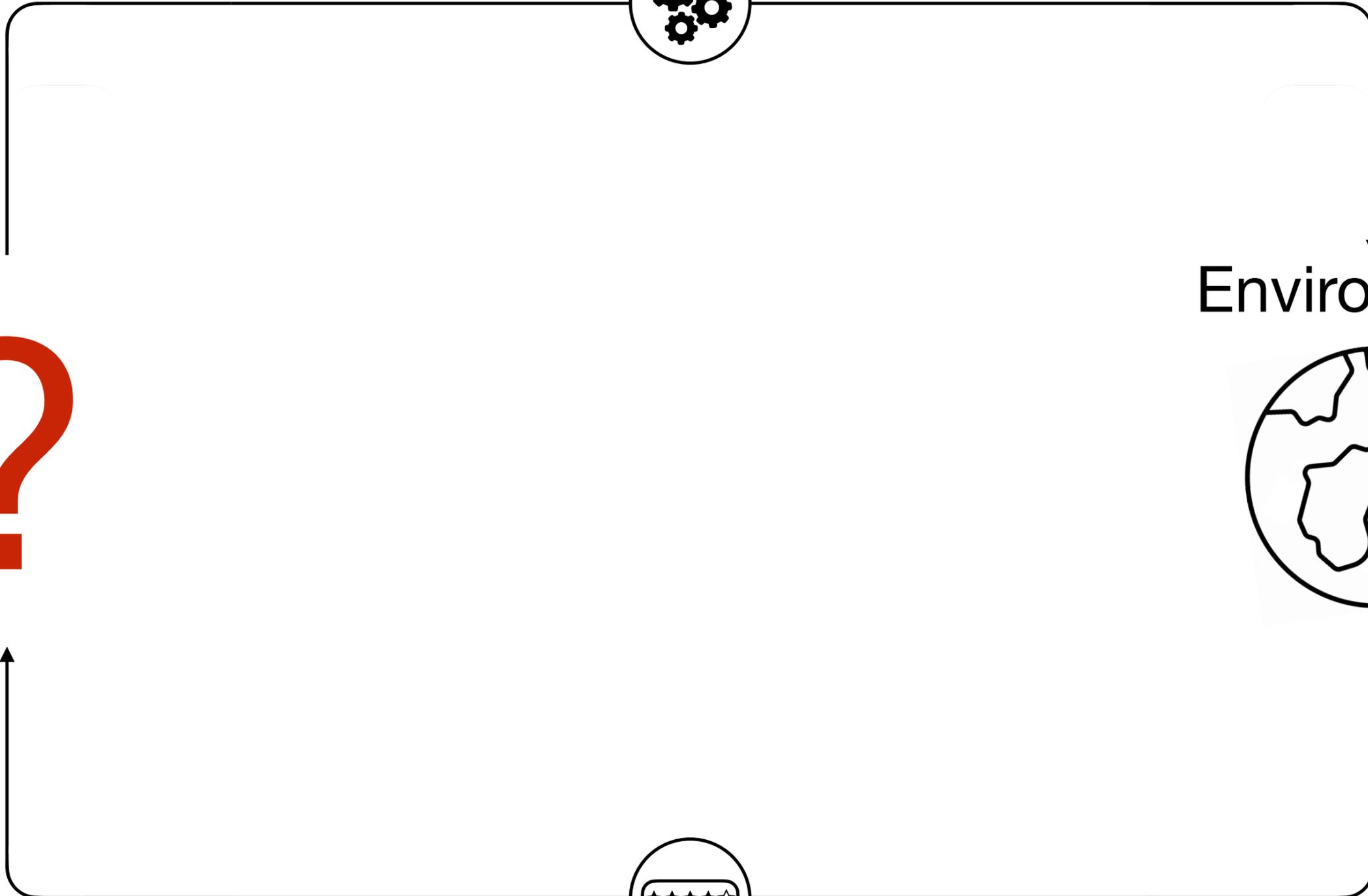
Execute



Environment

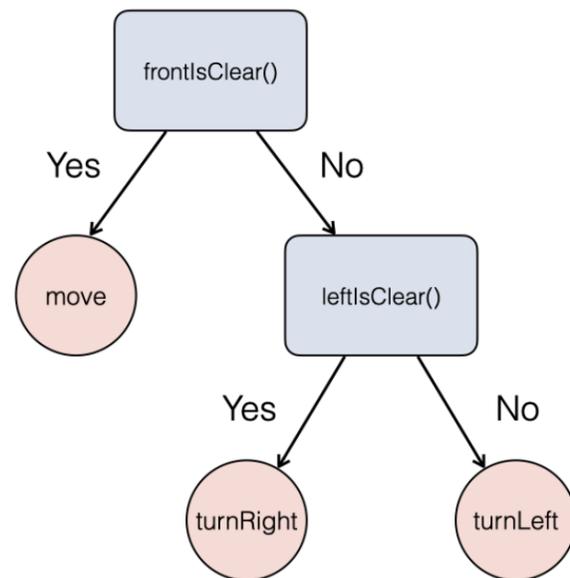


Reward



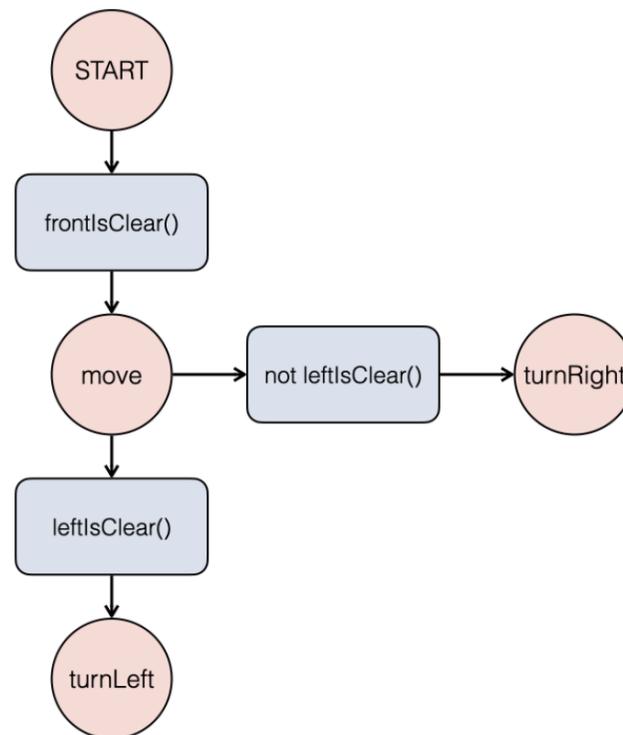
Programmatic Reinforcement Learning

Decision Trees



- Hard to represent repetitive behaviors

State Machines



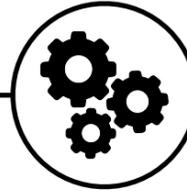
- Difficult to scale

Programs

```
DEF run()
  IF frontIsClear()
    move
  ELSE
    IF frontIsClear()
      turnLeft
    ELSE
      turnRight
```

- Flexible
- Human readable
- Hard to synthesize

Execute



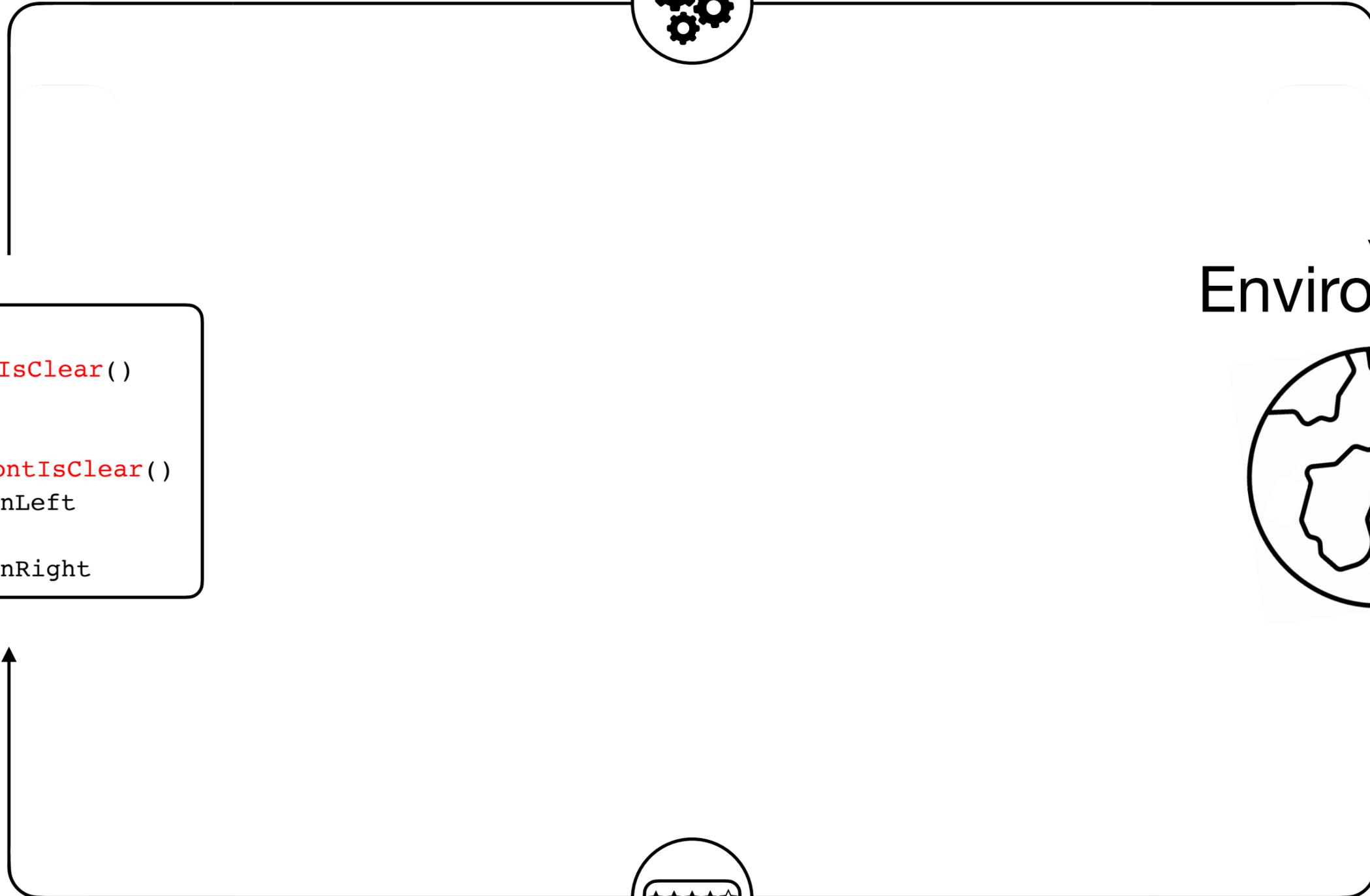
Environment

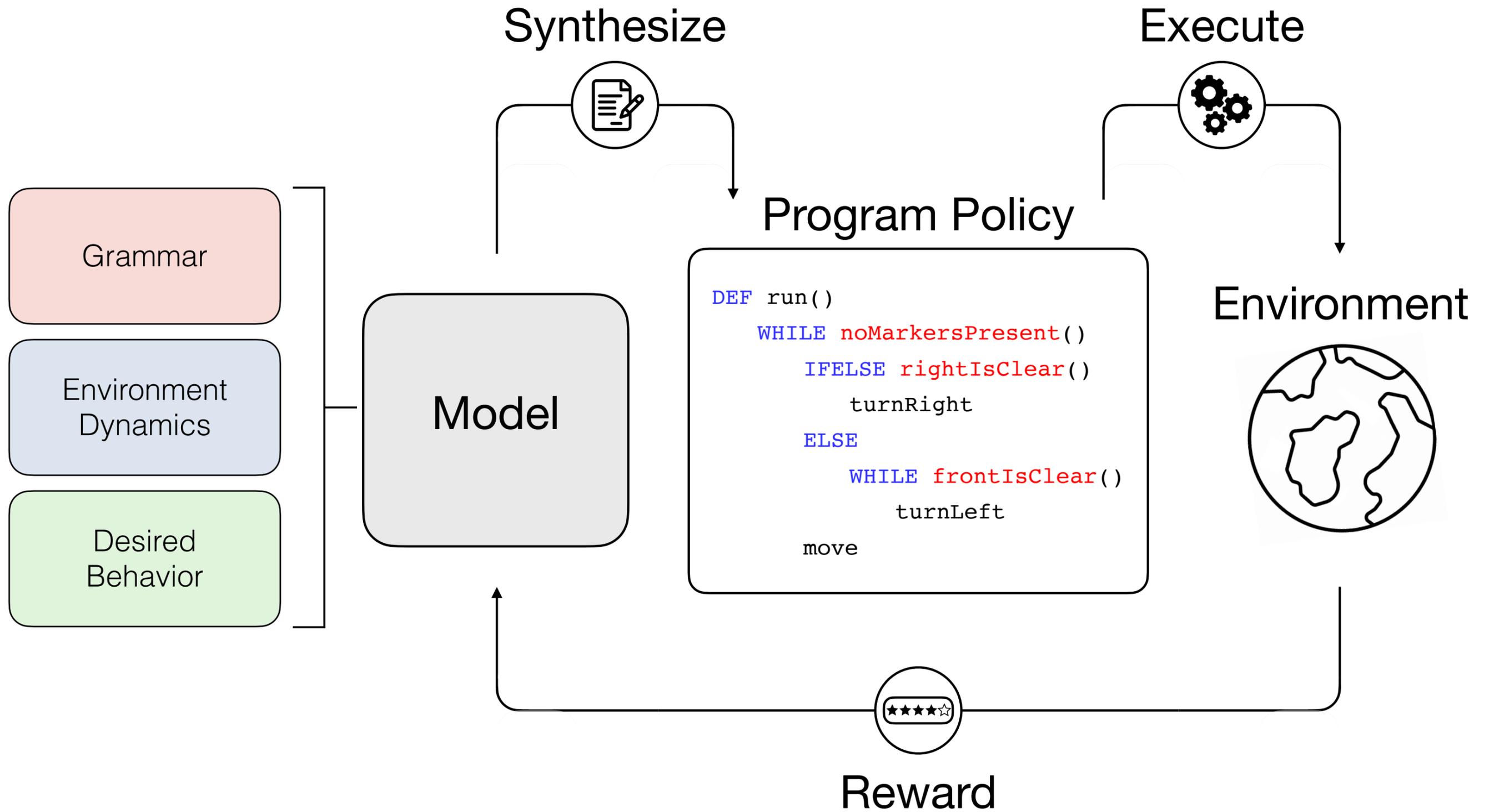


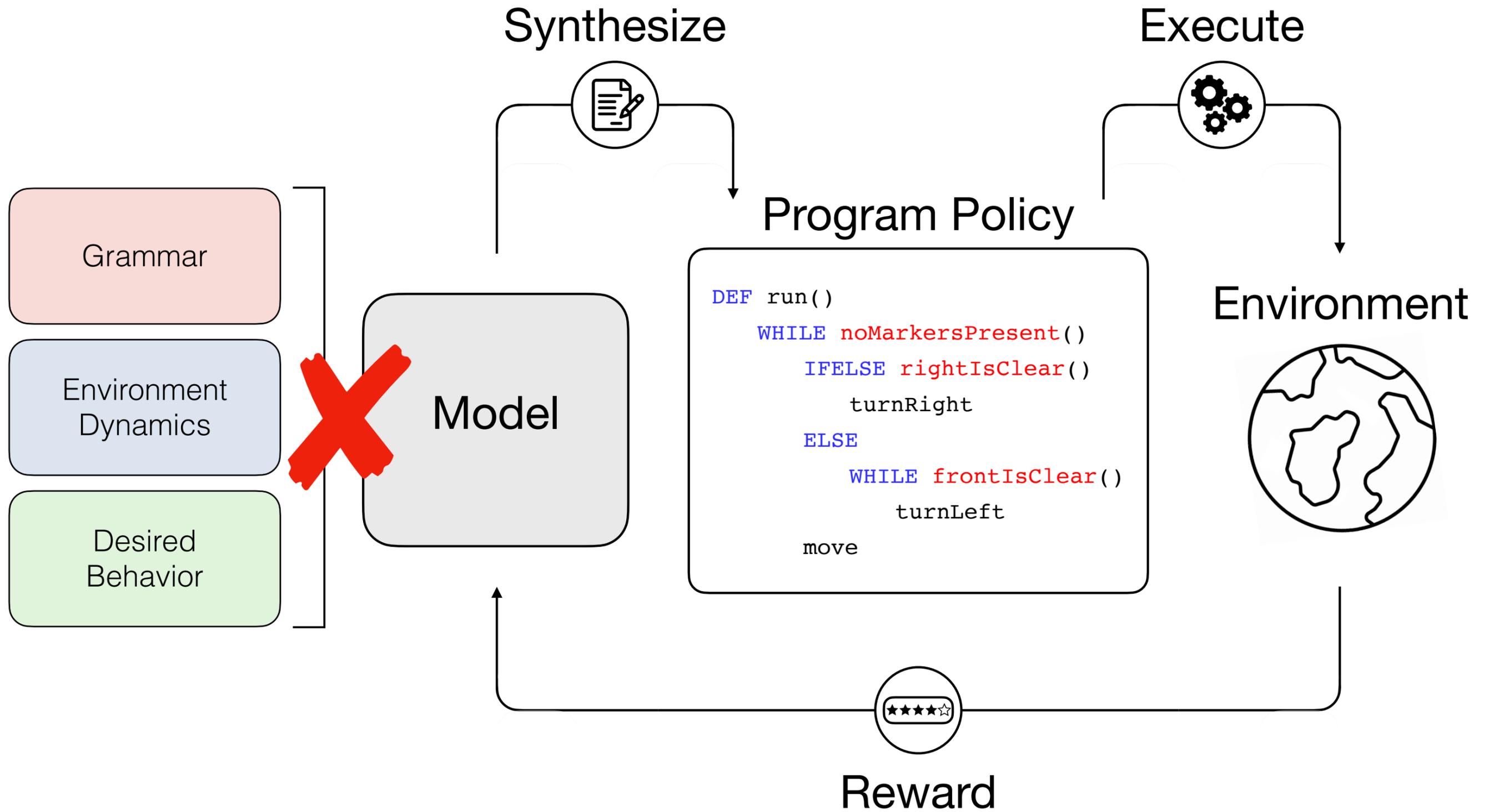
Reward



```
DEF run()  
  IF frontIsClear()  
    move  
  ELSE  
    IF frontIsClear()  
      turnLeft  
    ELSE  
      turnRight
```







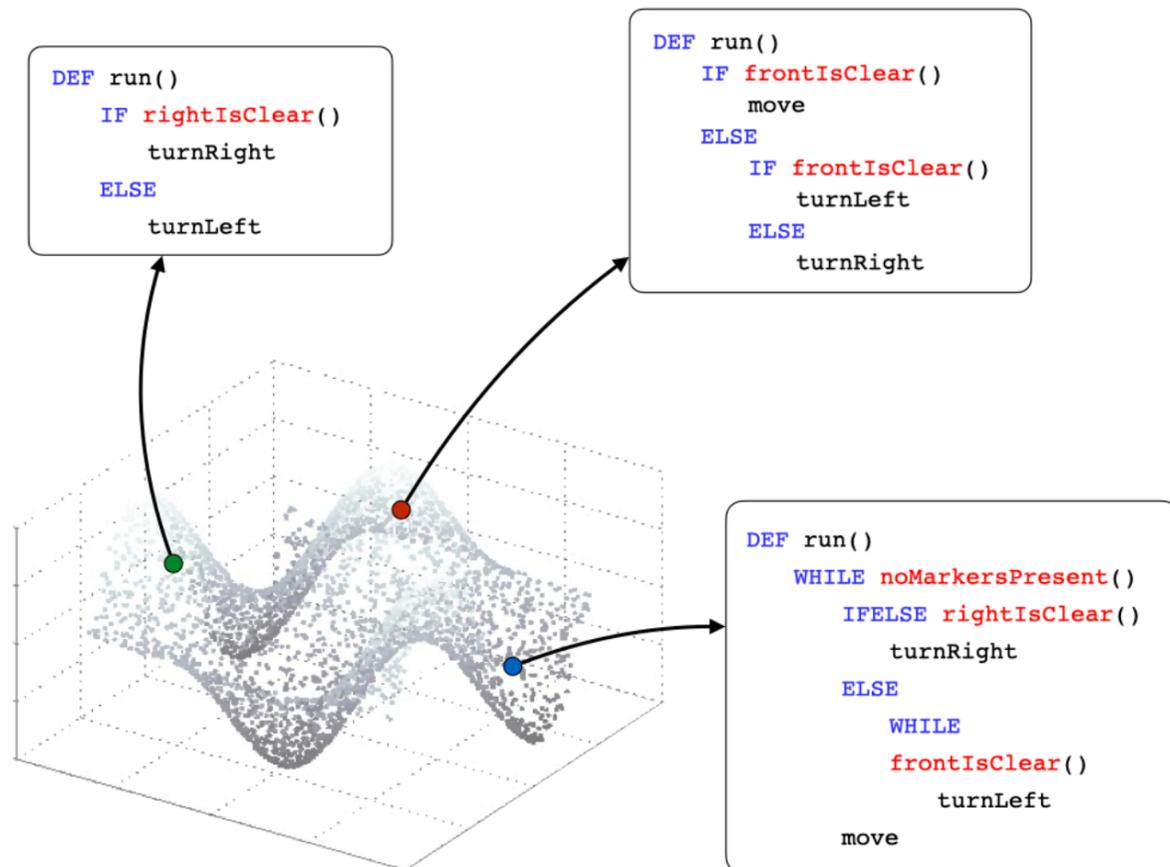
LEAPS: Learning Embeddings for Latent Program Synthesis

Stage 1

Learning a program embedding space from a set of randomly generated programs

Grammar

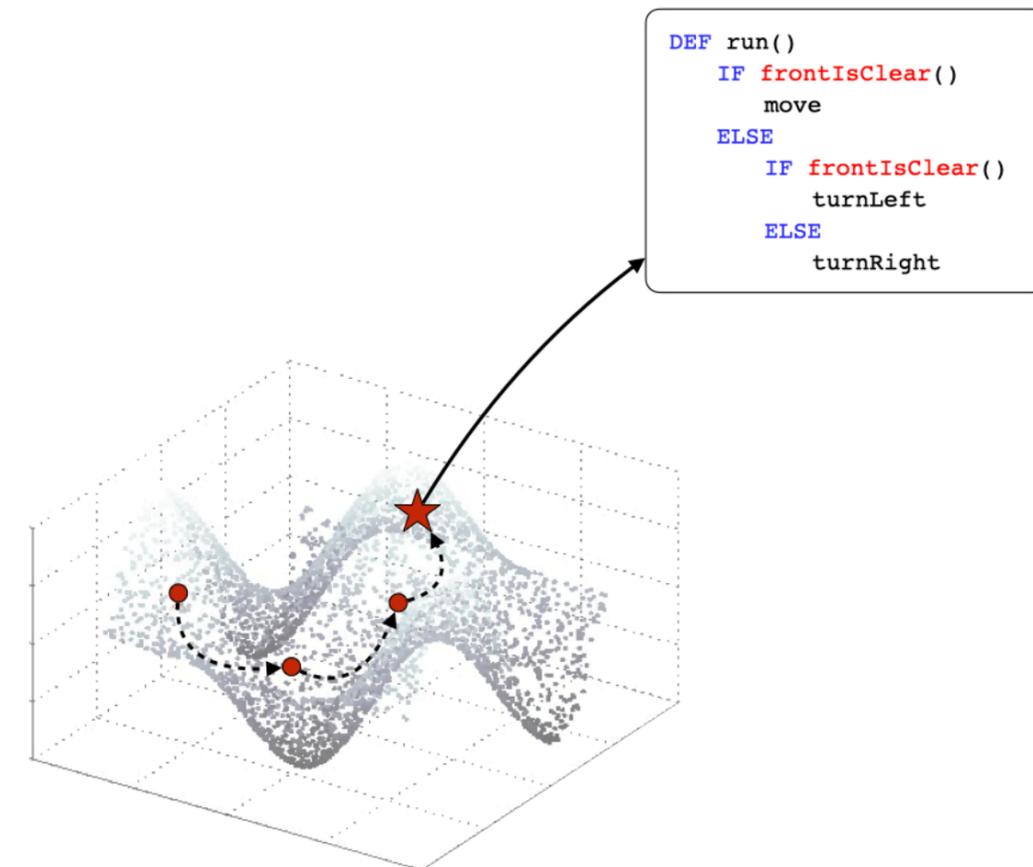
Environment Dynamics



Stage 2

Searching for a task-solving program

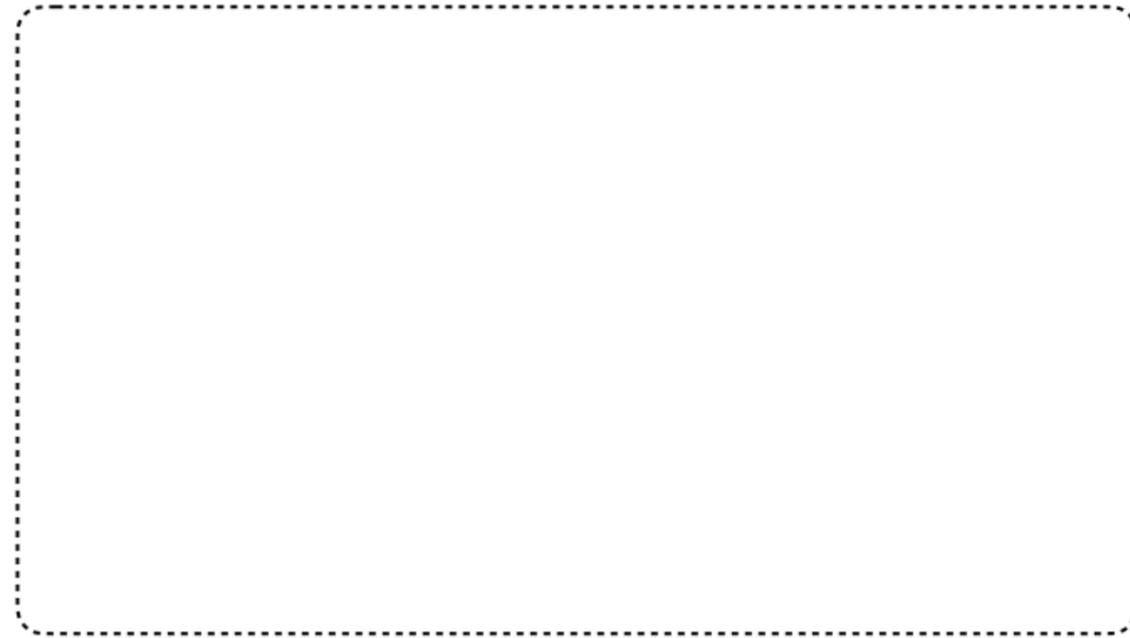
Desired Behavior



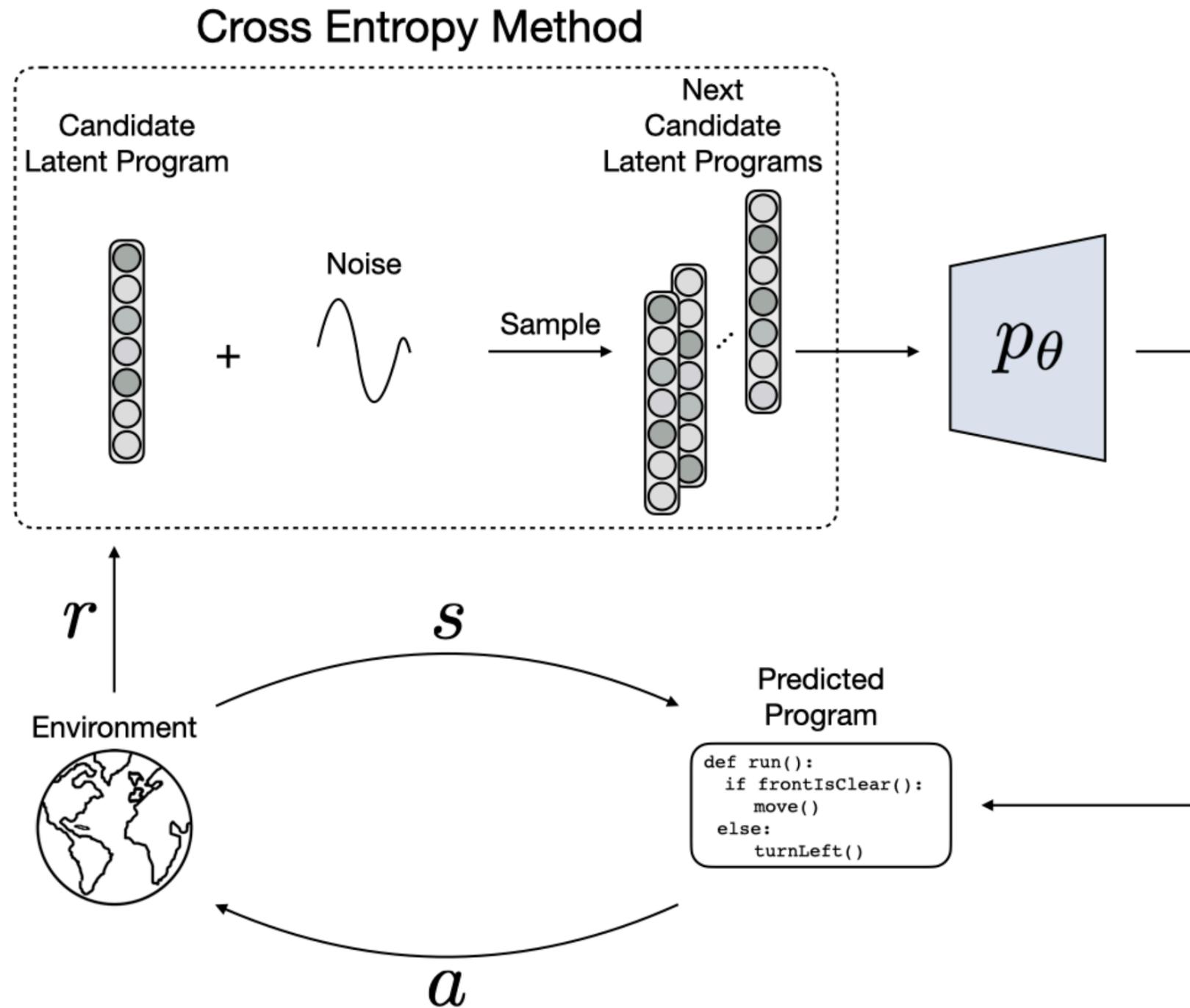
Learning a Program Embedding Space

Latent Program Search with Cross-Entropy Method

Cross Entropy Method

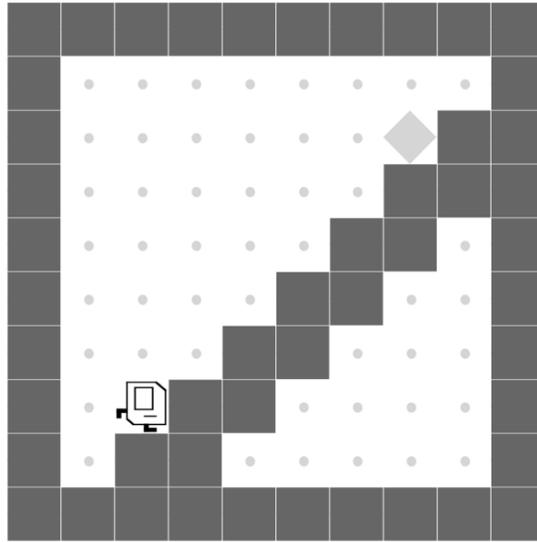


Latent Program Search with Cross-Entropy Method

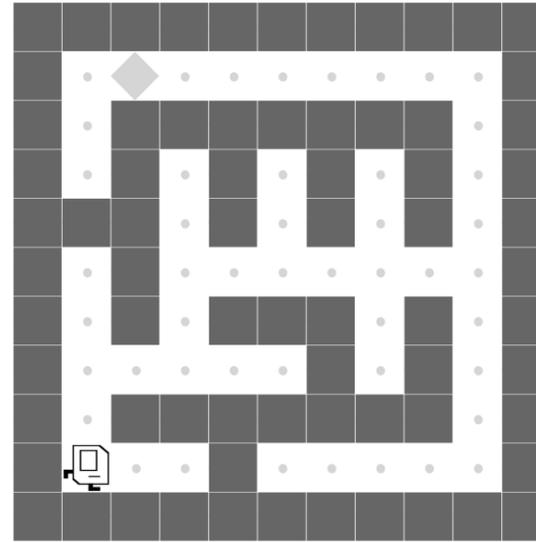


Karel Tasks

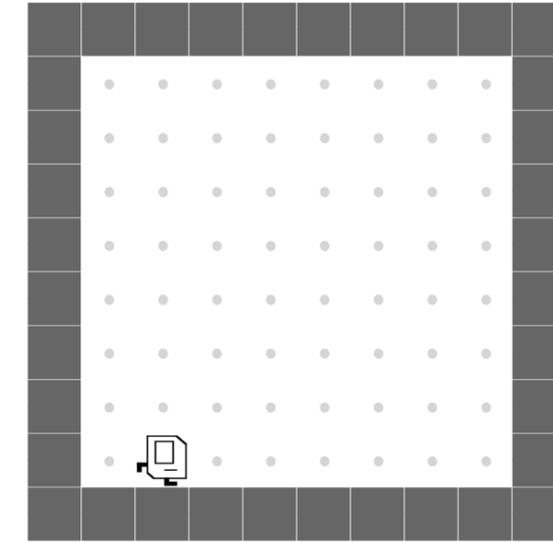
StairClimber



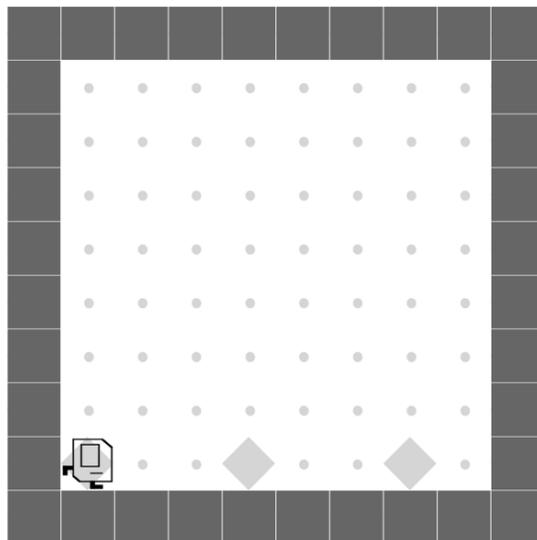
Maze



FourCorners

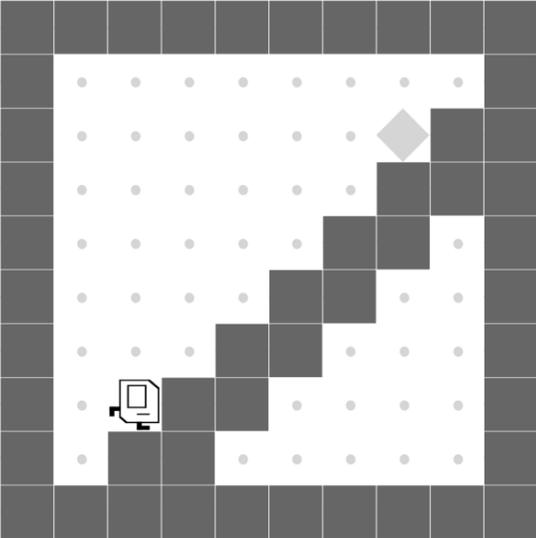


TopOff

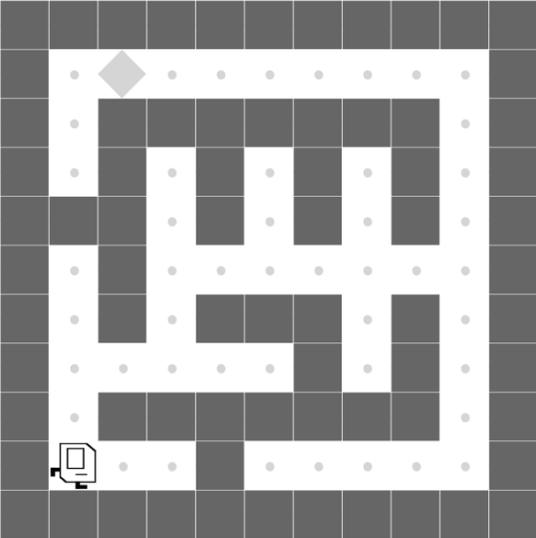


Karel Tasks

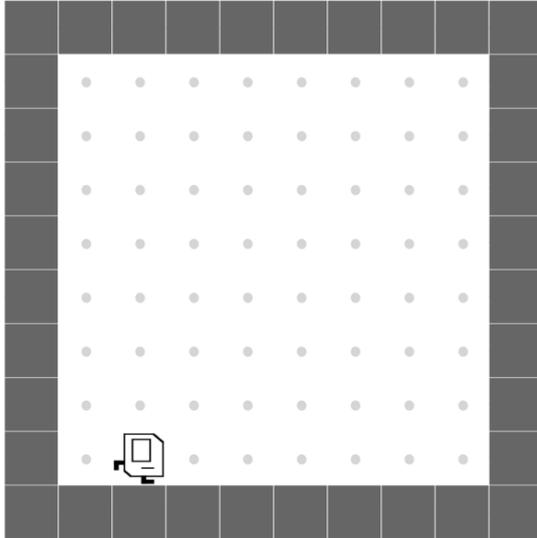
StairClimber



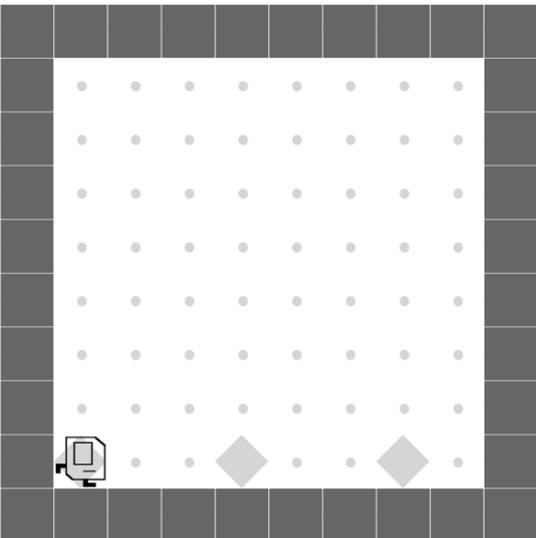
Maze



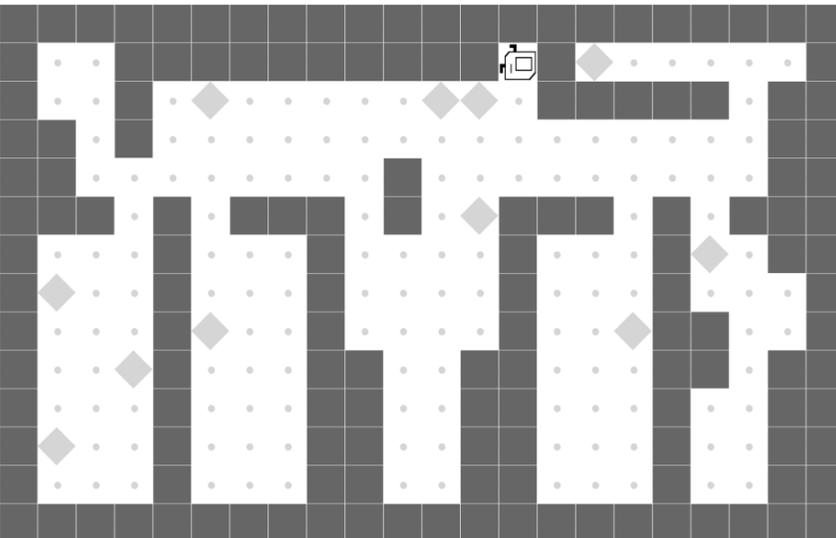
FourCorners



TopOff

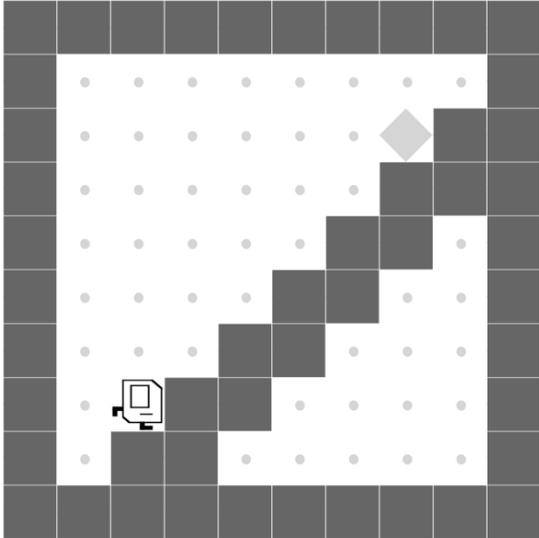


CleanHouse

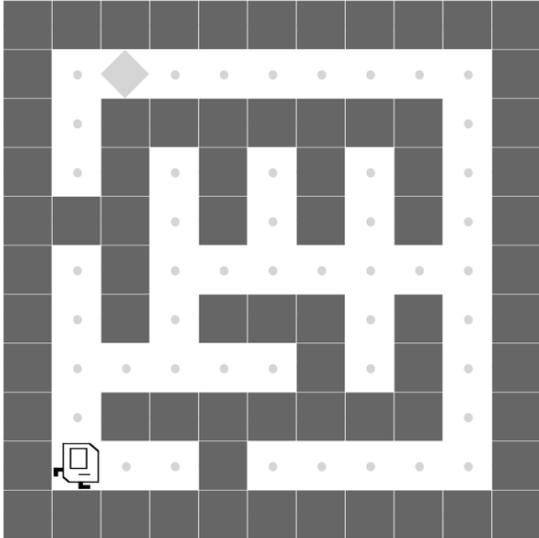


Karel Tasks

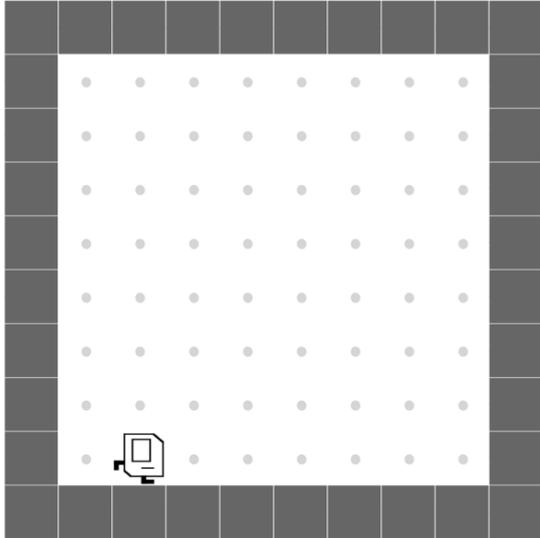
StairClimber



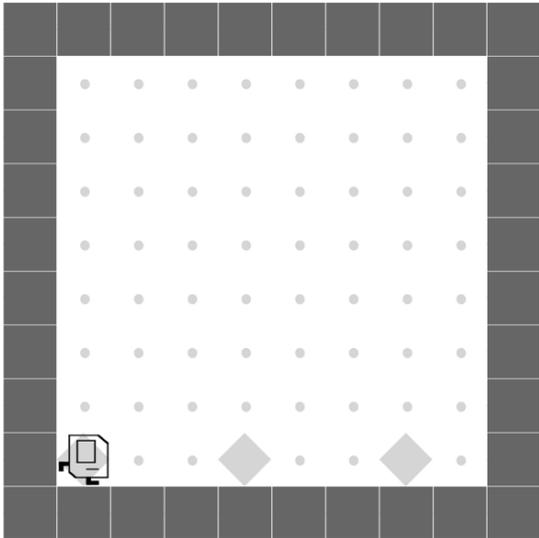
Maze



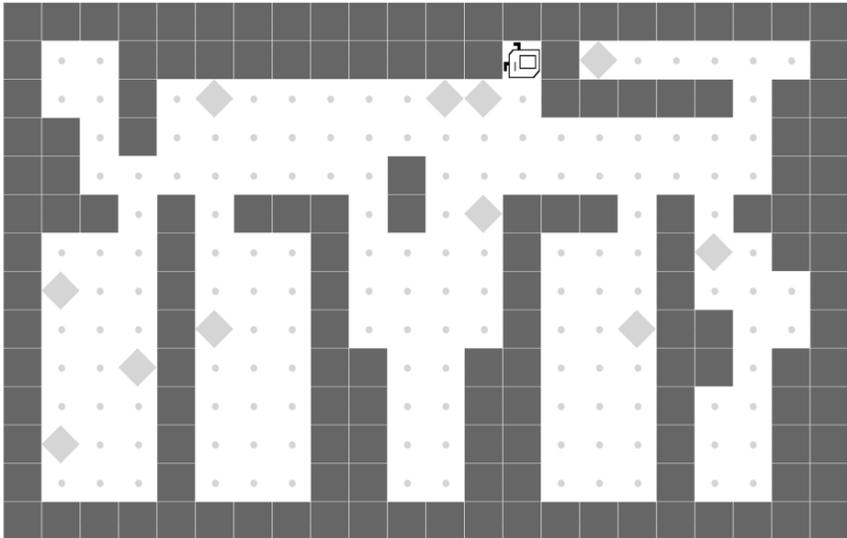
FourCorners



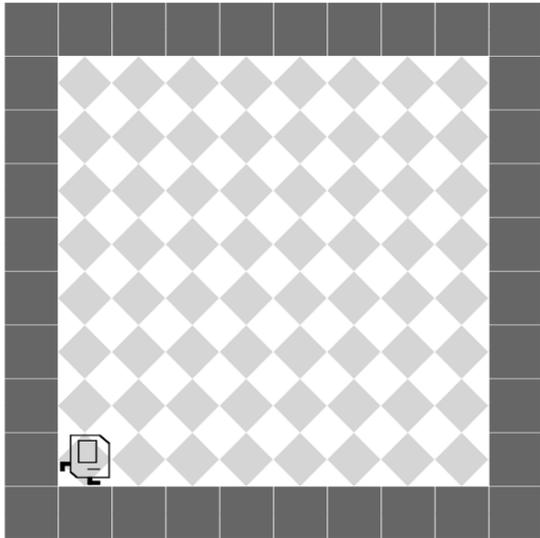
TopOff



CleanHouse



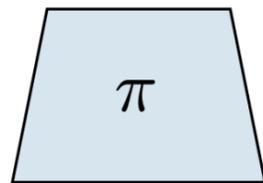
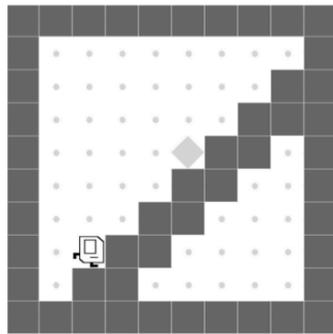
Harvester



Baselines

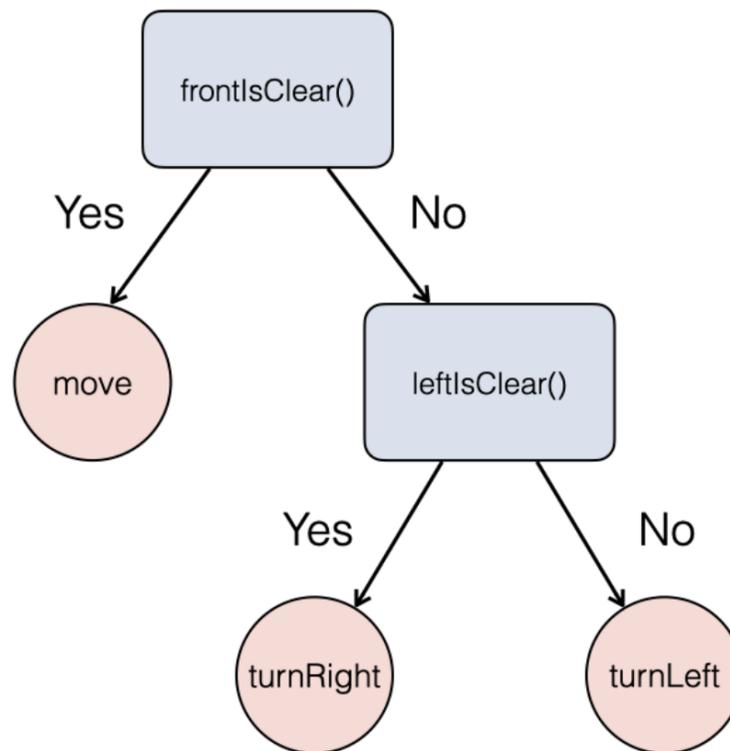
DRL

Raw State



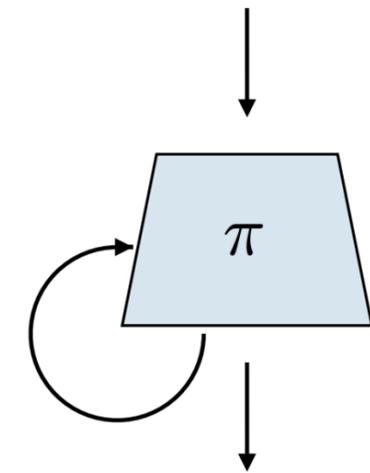
a

VIPER (Decision Tree)



Naive Program Synthesis

startToken



Program Token Generated at t

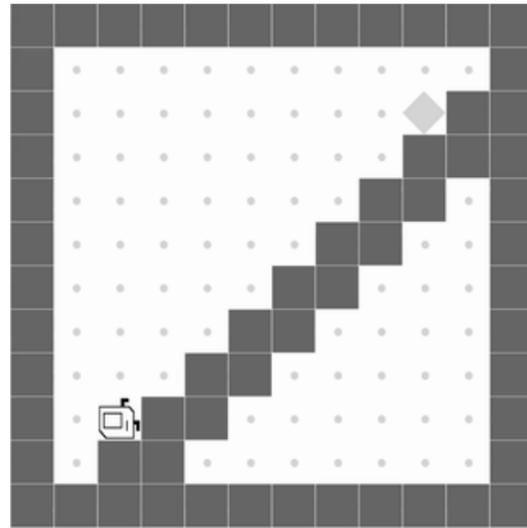
turnLeft()

Program Synthesized So Far

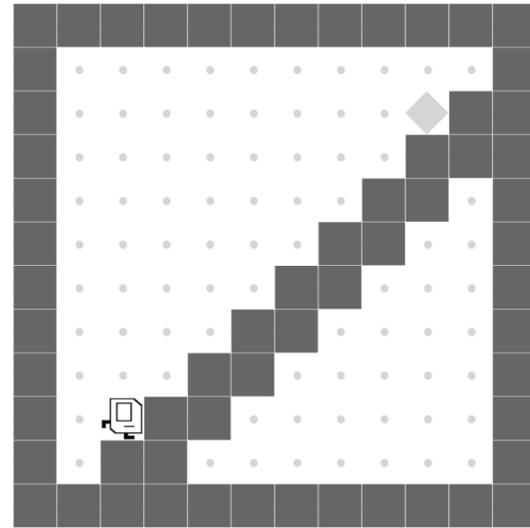
```
def run():
  if frontIsClear():
    move()
  else:
    turnLeft()
    ⋮
```

Results

StairClimber



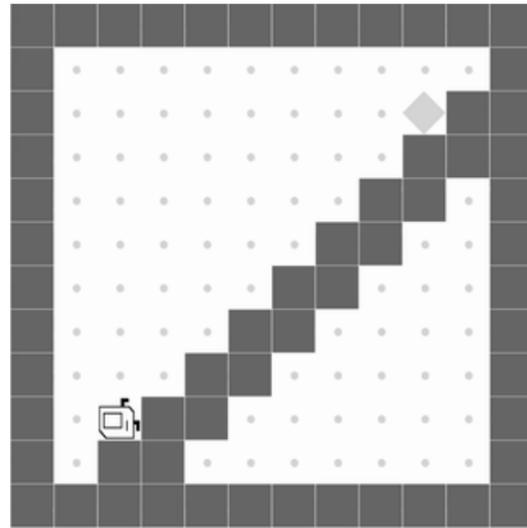
DRL



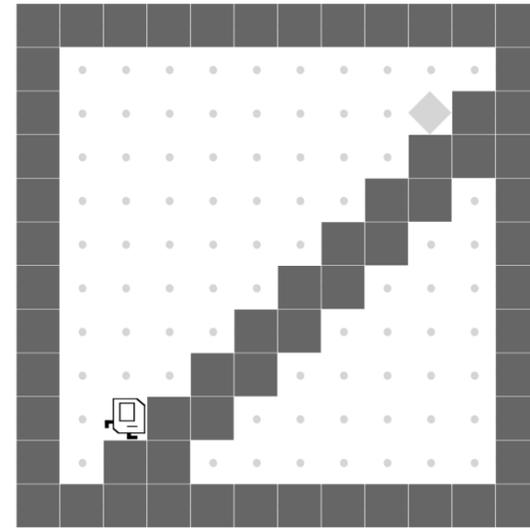
LEAPS

Results

StairClimber

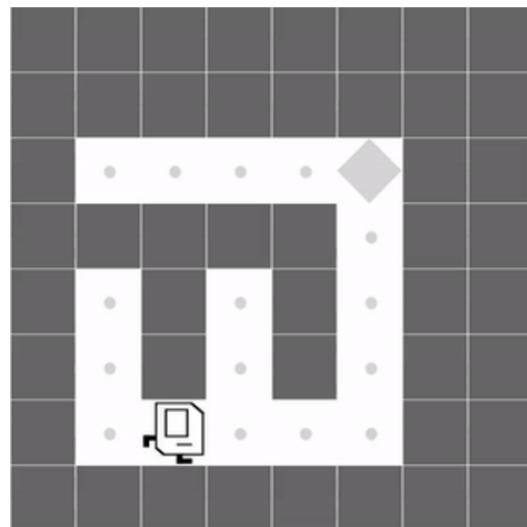


DRL

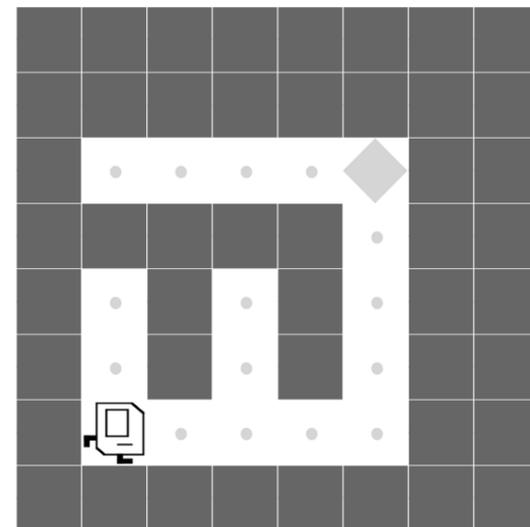


LEAPS

Maze



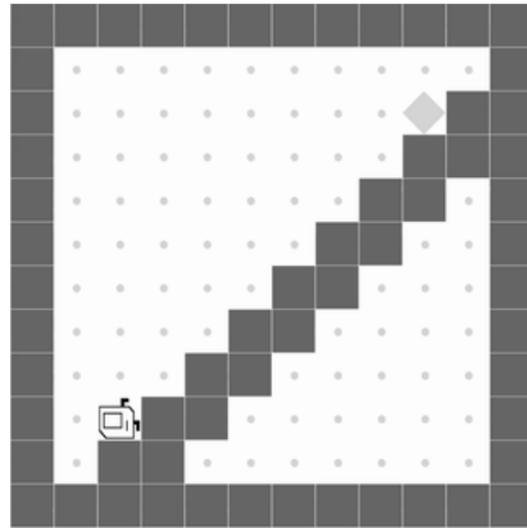
DRL



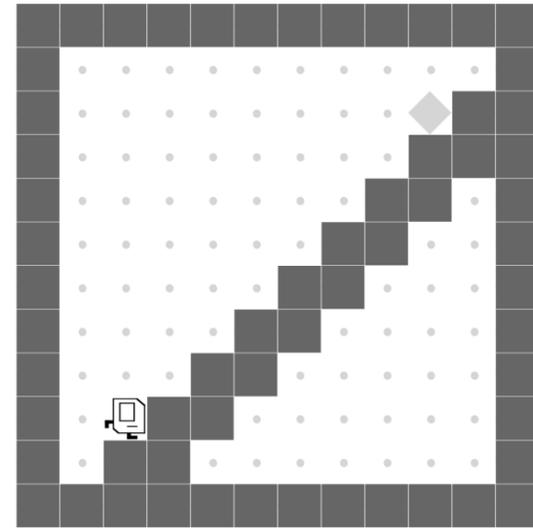
LEAPS

Results

StairClimber

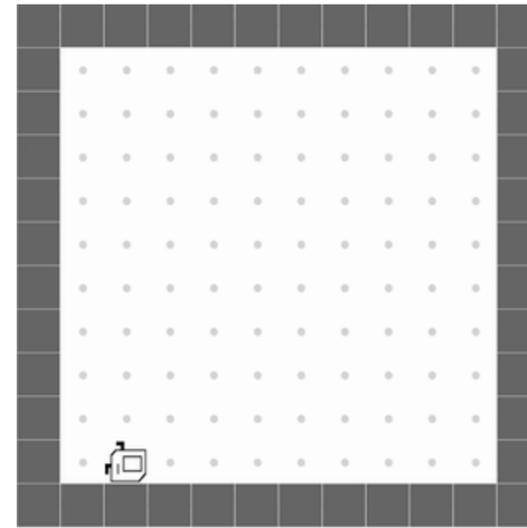


DRL

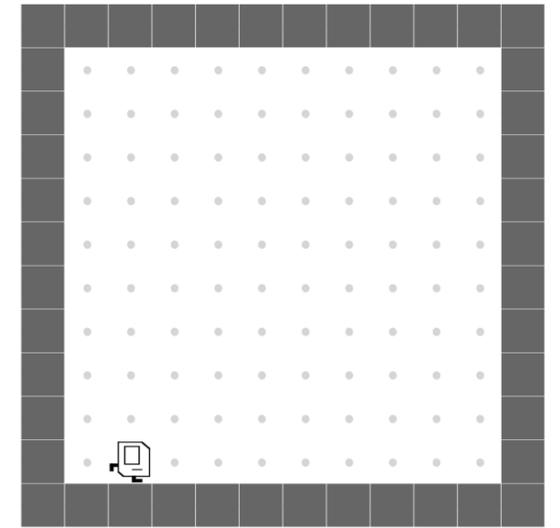


LEAPS

FourCorners

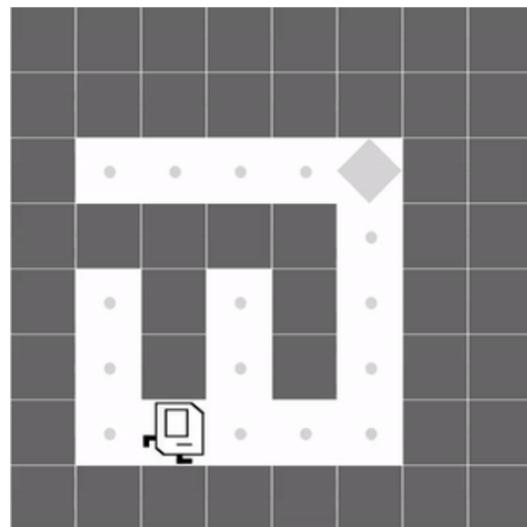


DRL

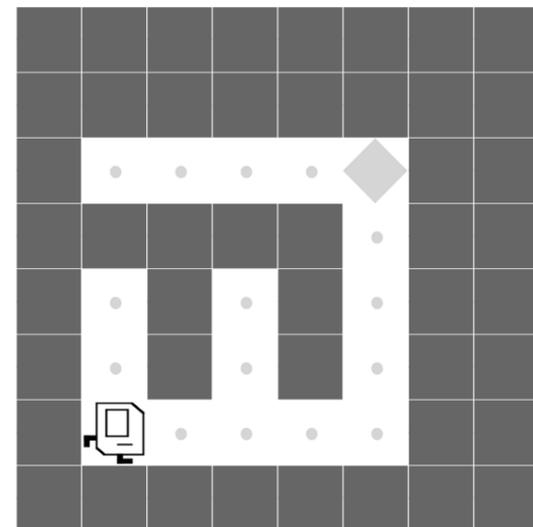


LEAPS

Maze



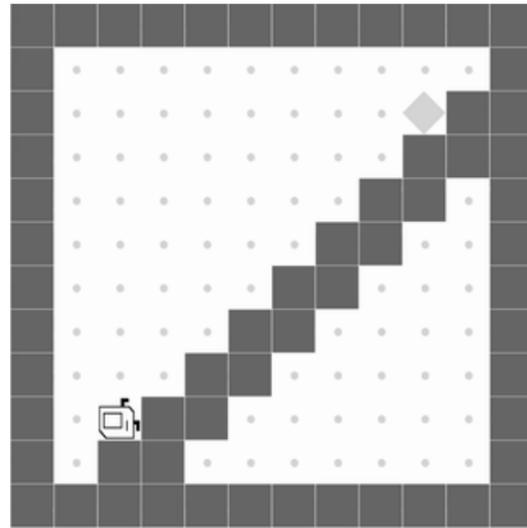
DRL



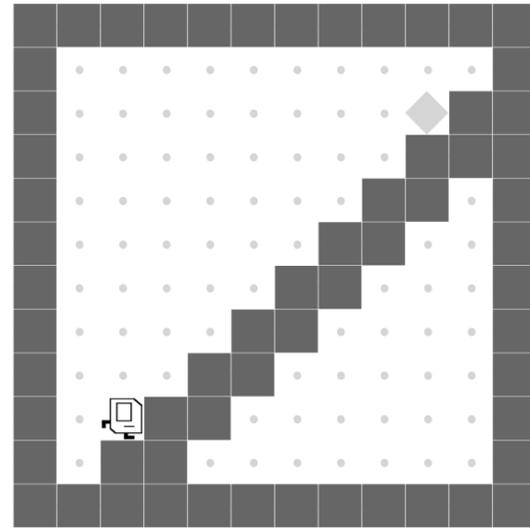
LEAPS

Results

StairClimber

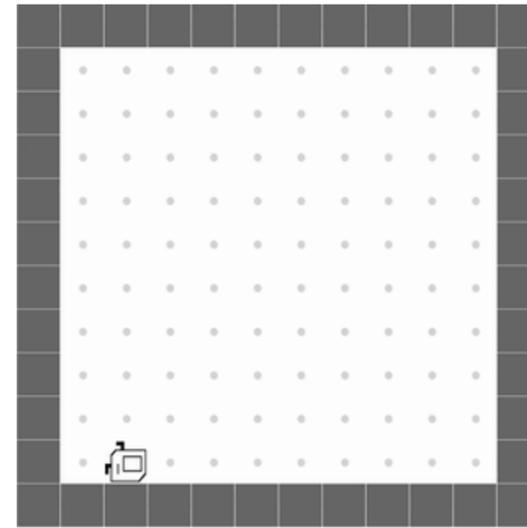


DRL

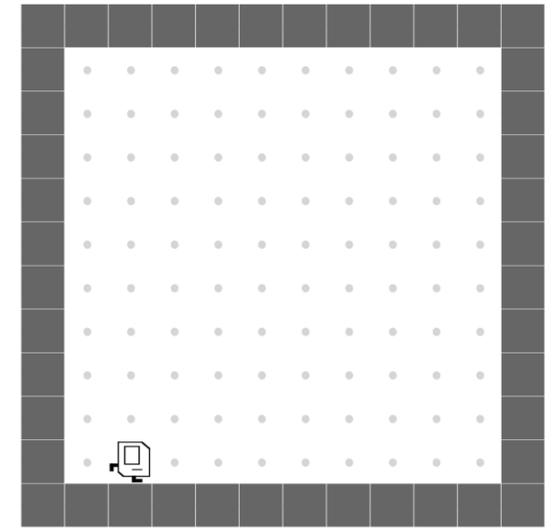


LEAPS

FourCorners

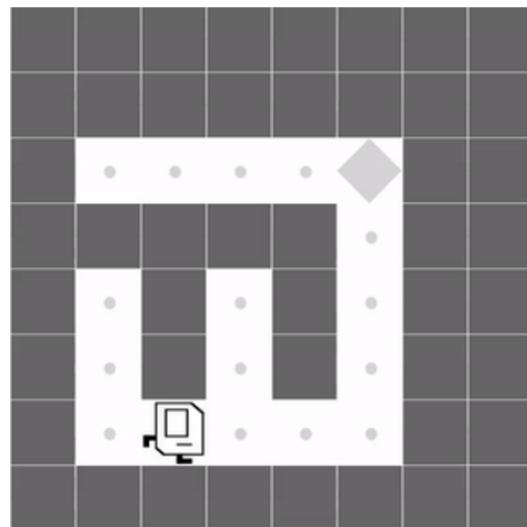


DRL

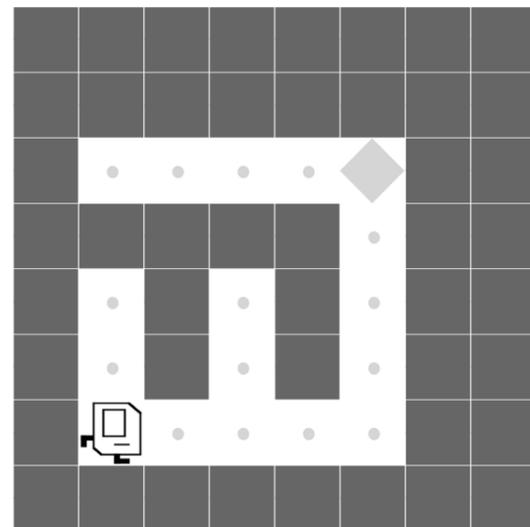


LEAPS

Maze

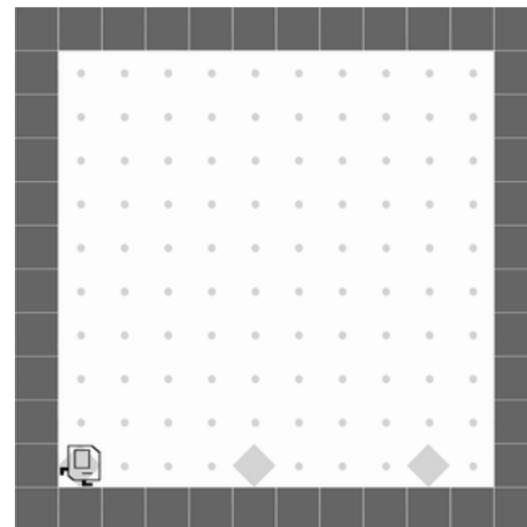


DRL

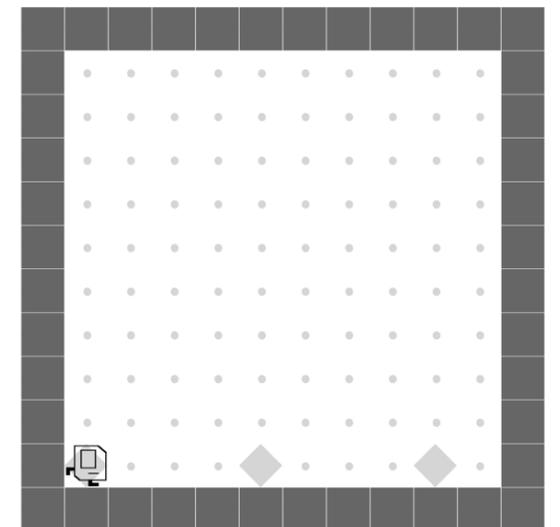


LEAPS

TopOff

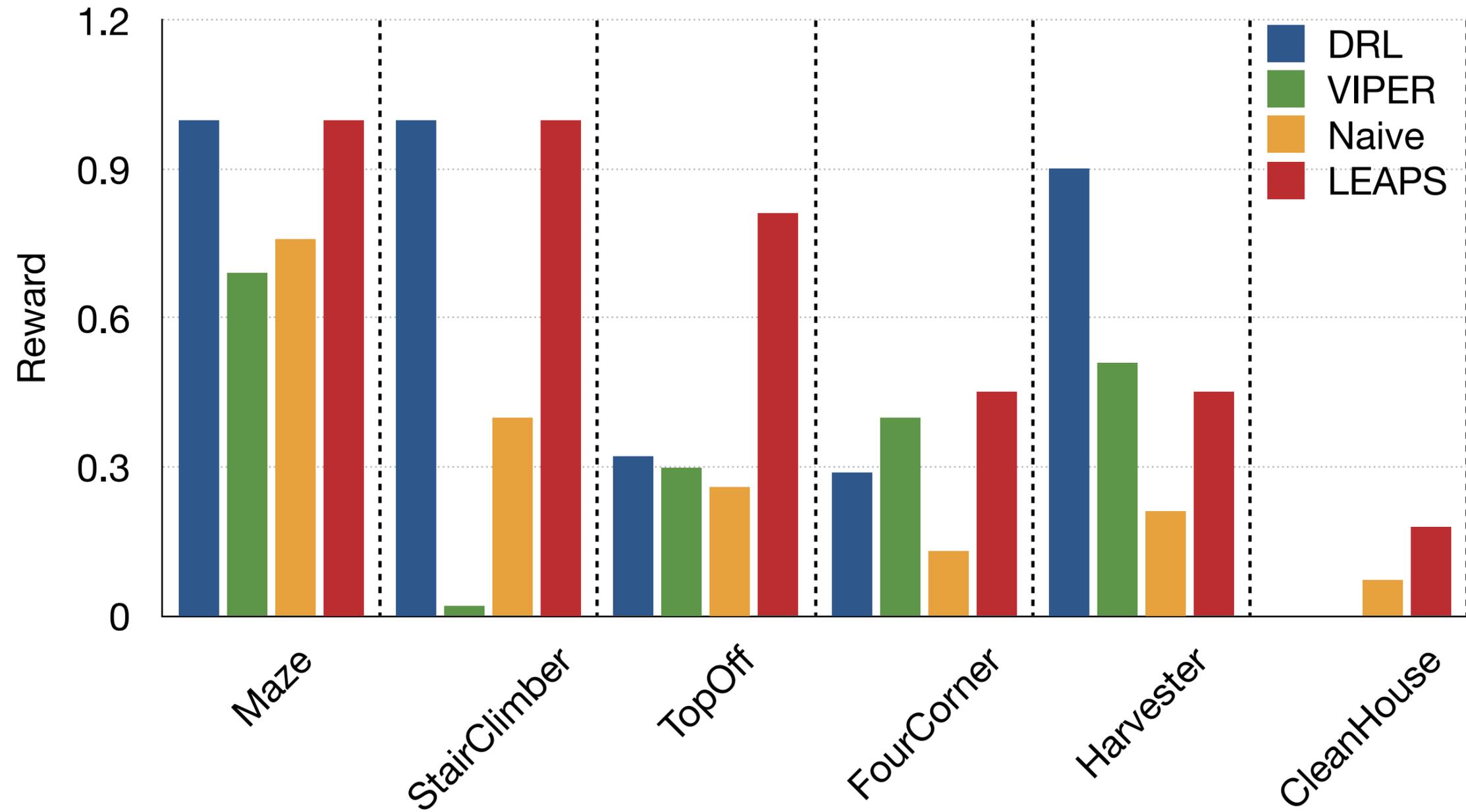


DRL

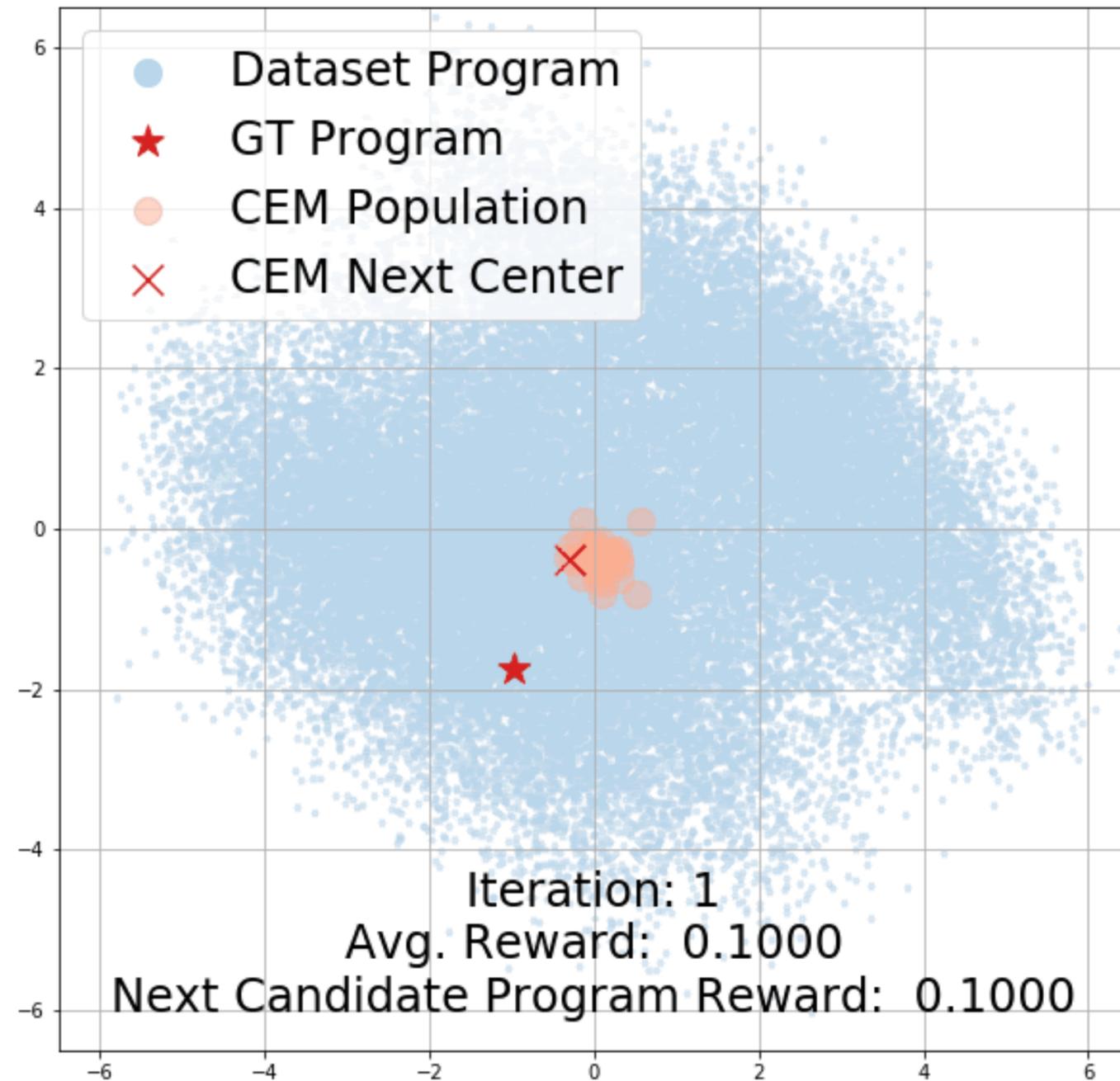


LEAPS

Results



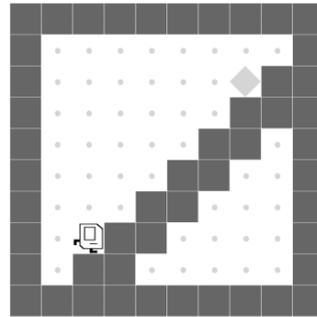
Results - CEM trajectory Visualization



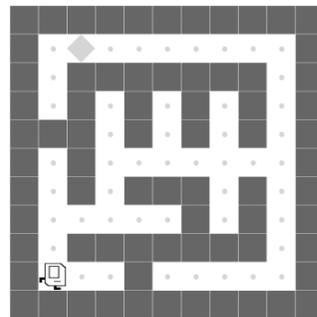
Zero-shot Generalization

Learning

StairClimber



Maze

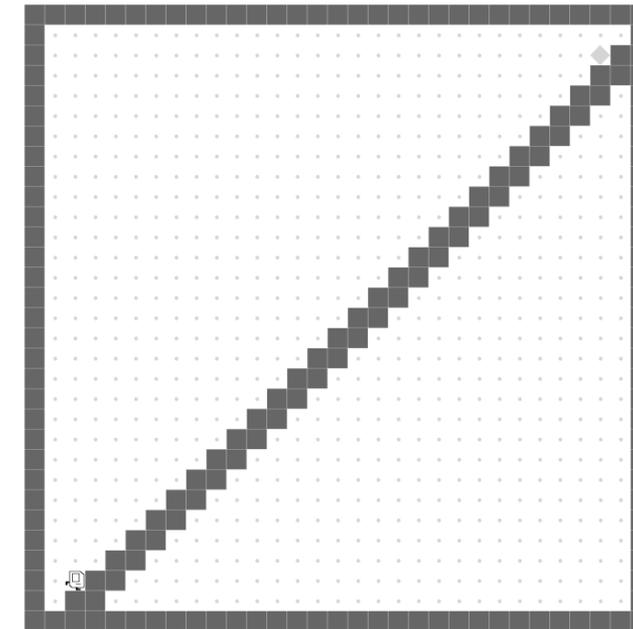
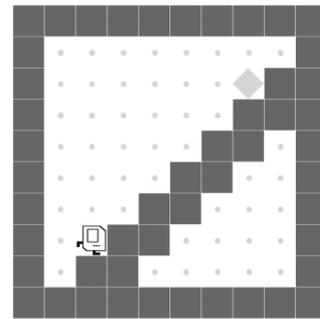


Zero-shot Generalization

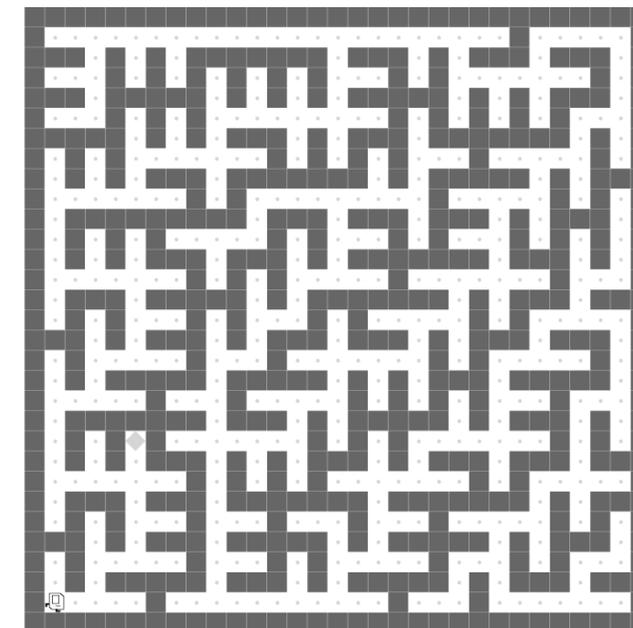
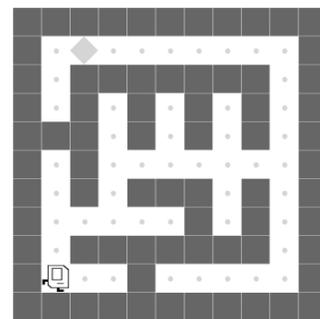
Learning

Evaluation on 100 x 100

StairClimber



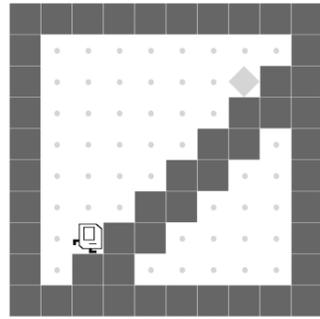
Maze



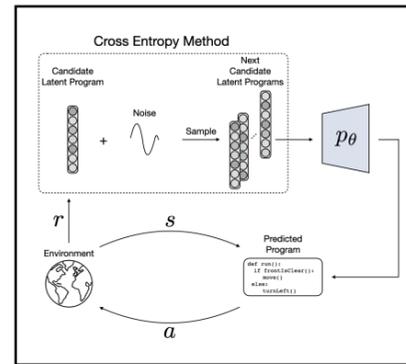
LEAPS Zero-shot Generalization

Learning

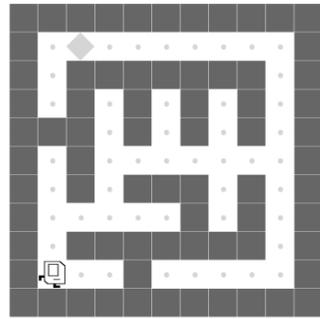
StairClimber



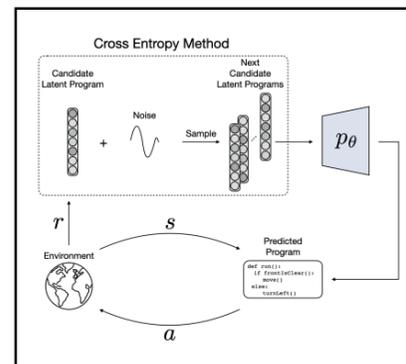
CEM search



Maze



CEM search



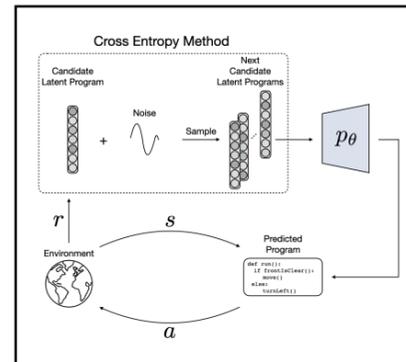
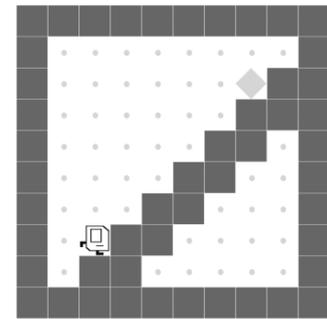
LEAPS Zero-shot Generalization

Learning

LEAPS policy

StairClimber

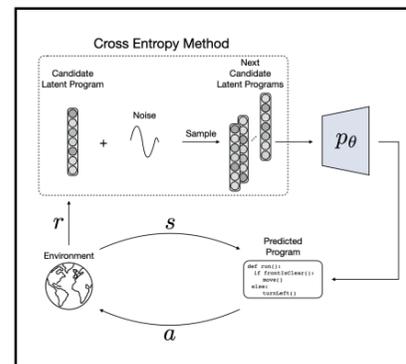
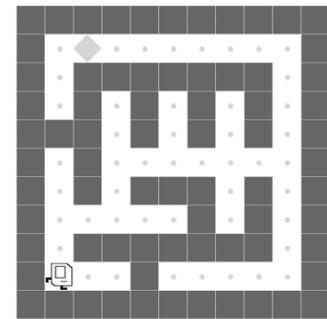
CEM search



```
DEF run()  
  WHILE noMarkersPresent()  
    turnRight  
    move  
  WHILE rightIsClear()  
    turnLeft
```

Maze

CEM search



```
DEF run()  
  IF frontIsClear()  
    turnLeft  
  WHILE noMarkersPresent()  
    turnRight  
    move
```

LEAPS Zero-shot Generalization

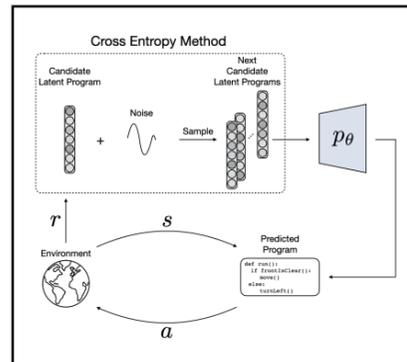
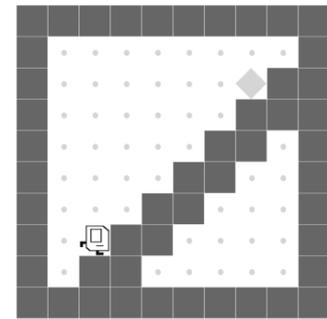
Learning

LEAPS policy

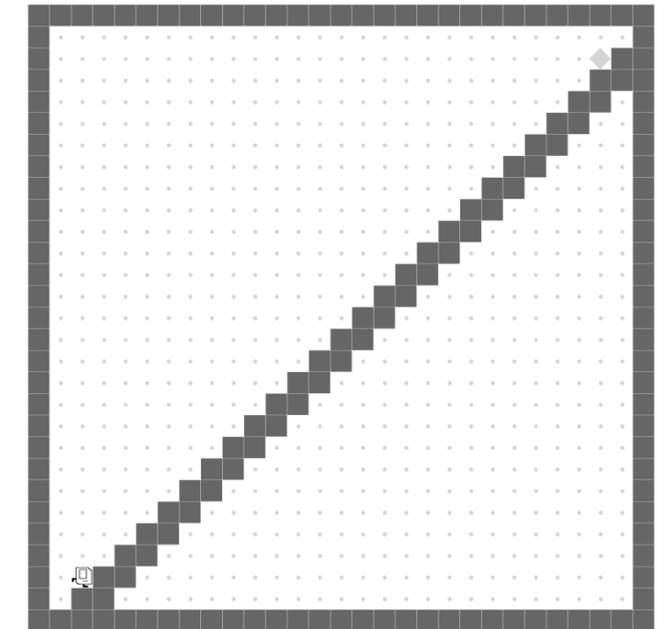
Evaluation on 100 x 100

StairClimber

CEM search

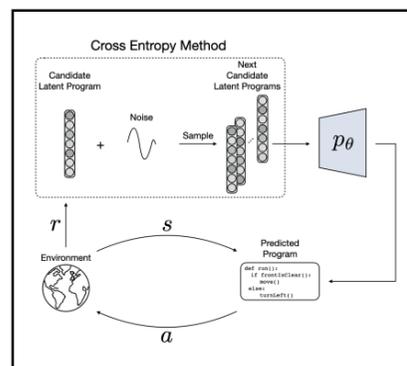
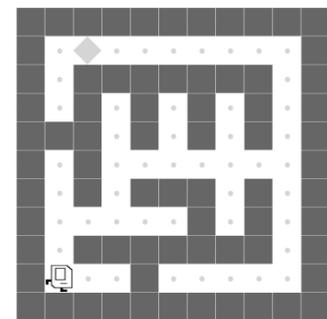


```
DEF run()  
  WHILE noMarkersPresent()  
    turnRight  
    move  
  WHILE rightIsClear()  
    turnLeft
```

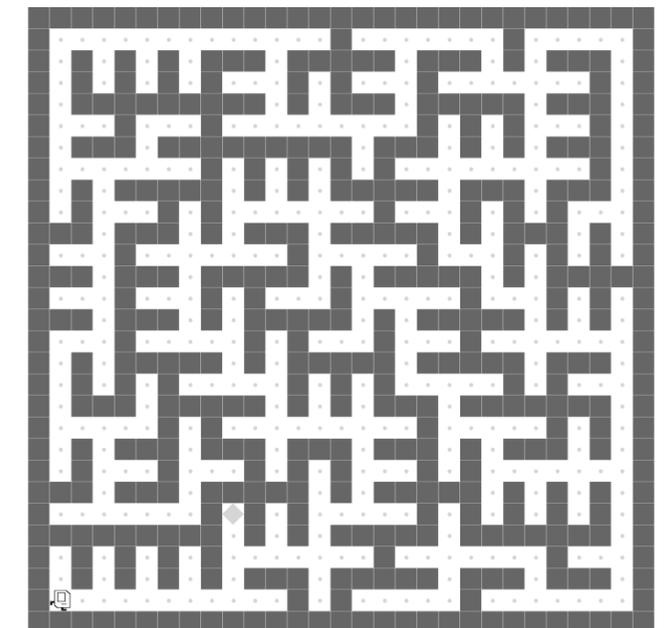


Maze

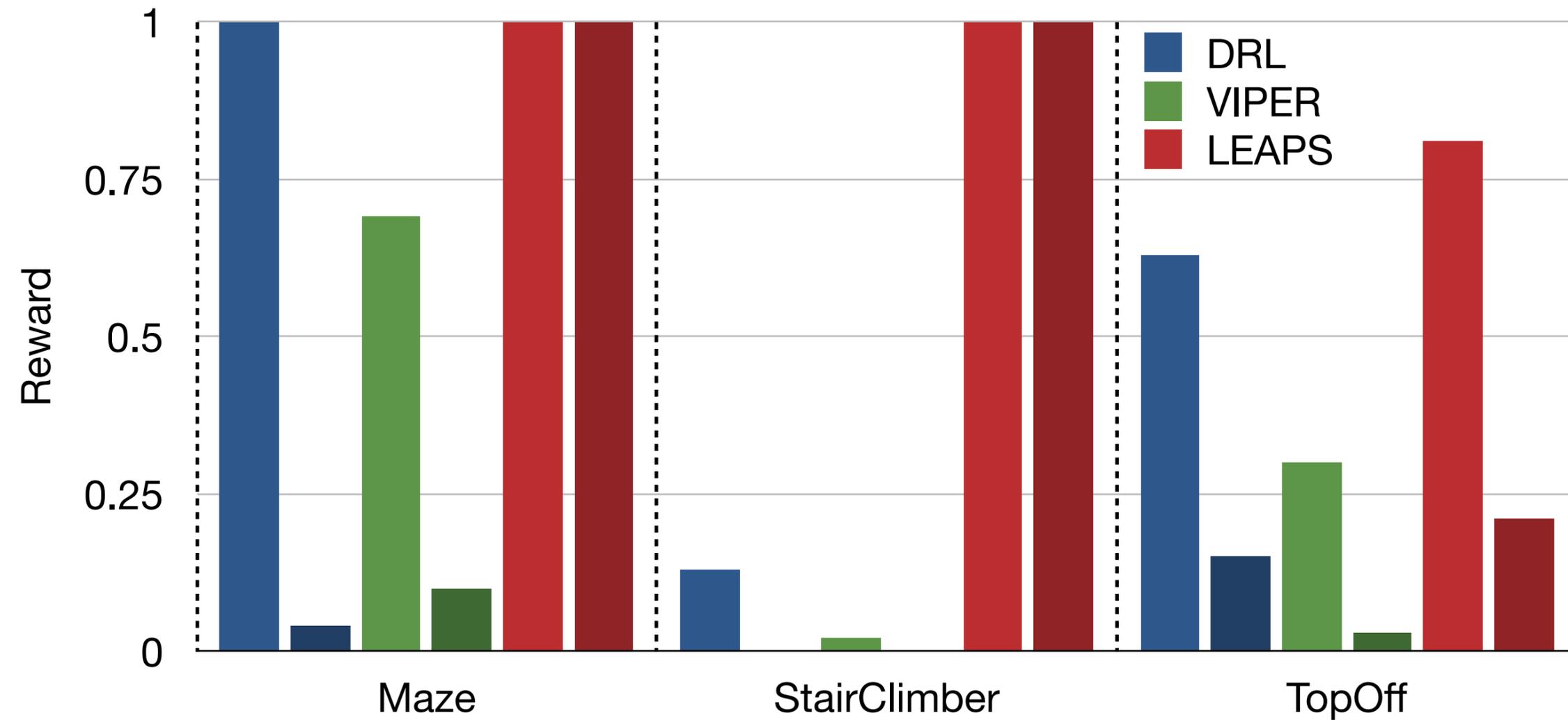
CEM search



```
DEF run()  
  IF frontIsClear()  
    turnLeft  
  WHILE noMarkersPresent()  
    turnRight  
    move
```

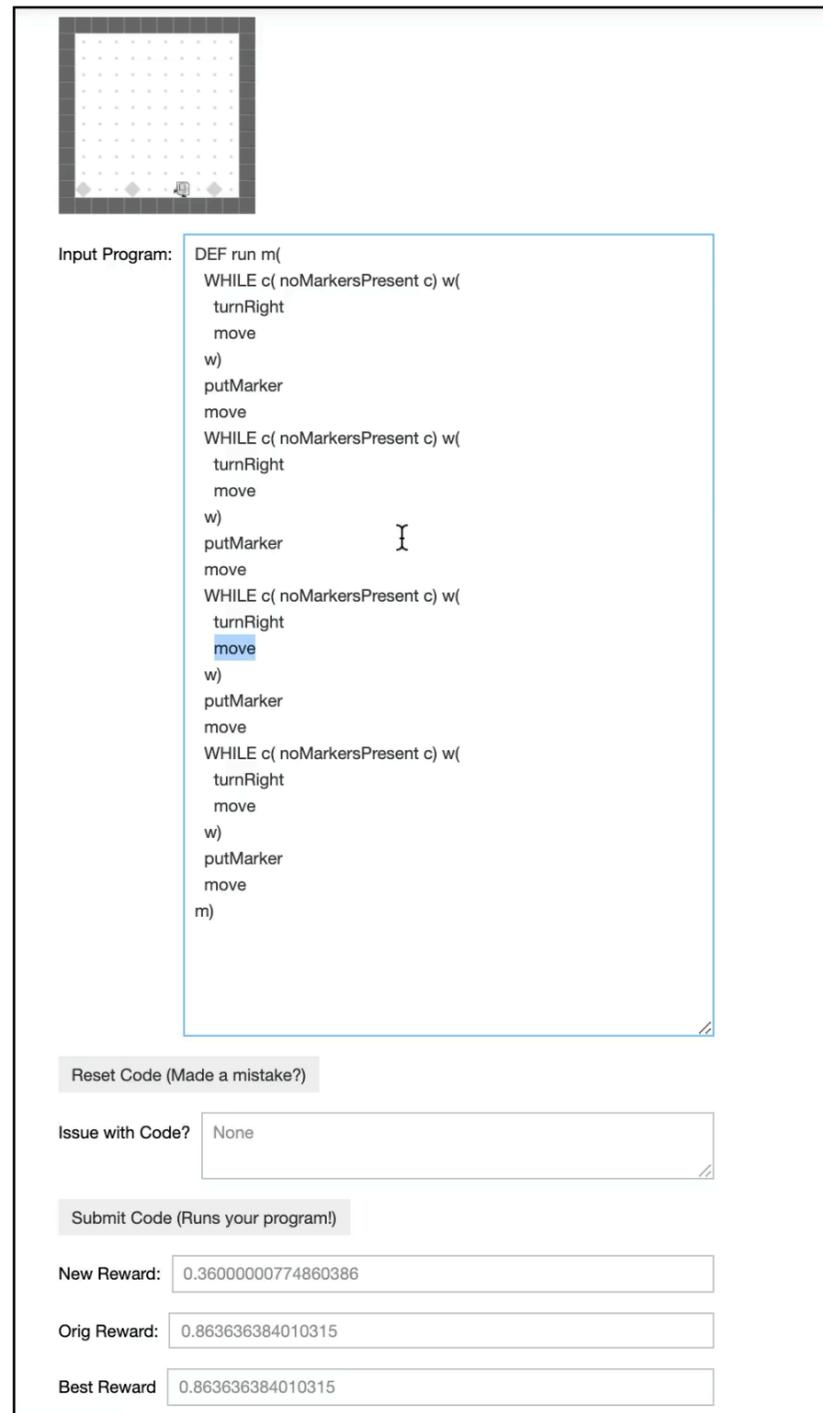


Results - Zero-shot Generalization



Interpretability

Human Debugging Interface



Input Program:

```
DEF run m(  
  WHILE c( noMarkersPresent c) w(  
    turnRight  
    move  
  w)  
  putMarker  
  move  
  WHILE c( noMarkersPresent c) w(  
    turnRight  
    move  
  w)  
  putMarker  
  move  
  WHILE c( noMarkersPresent c) w(  
    turnRight  
    move  
  w)  
  putMarker  
  move  
  WHILE c( noMarkersPresent c) w(  
    turnRight  
    move  
  w)  
  putMarker  
  move  
  m)
```

Reset Code (Made a mistake?)

Issue with Code?

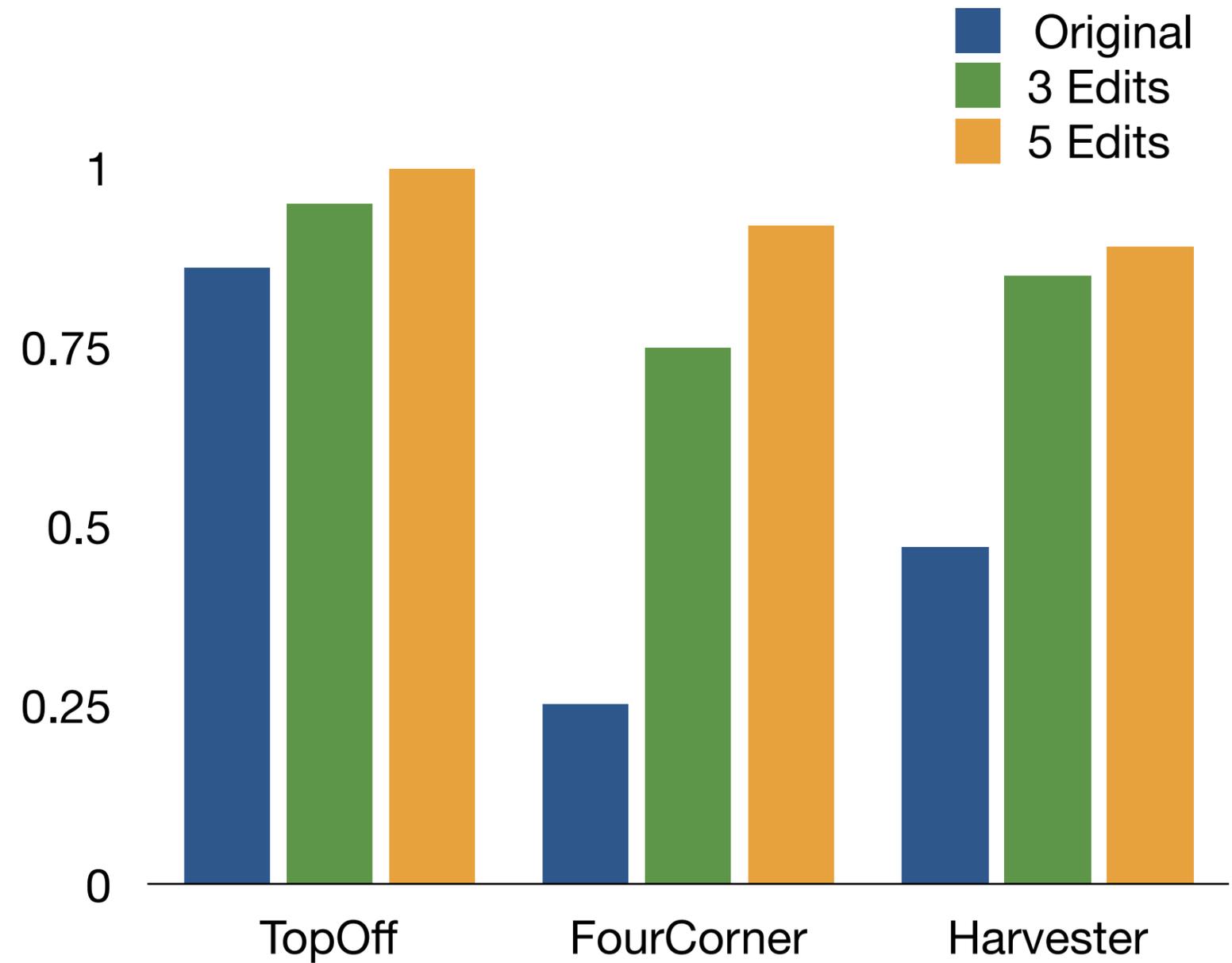
Submit Code (Runs your program!)

New Reward:

Orig Reward:

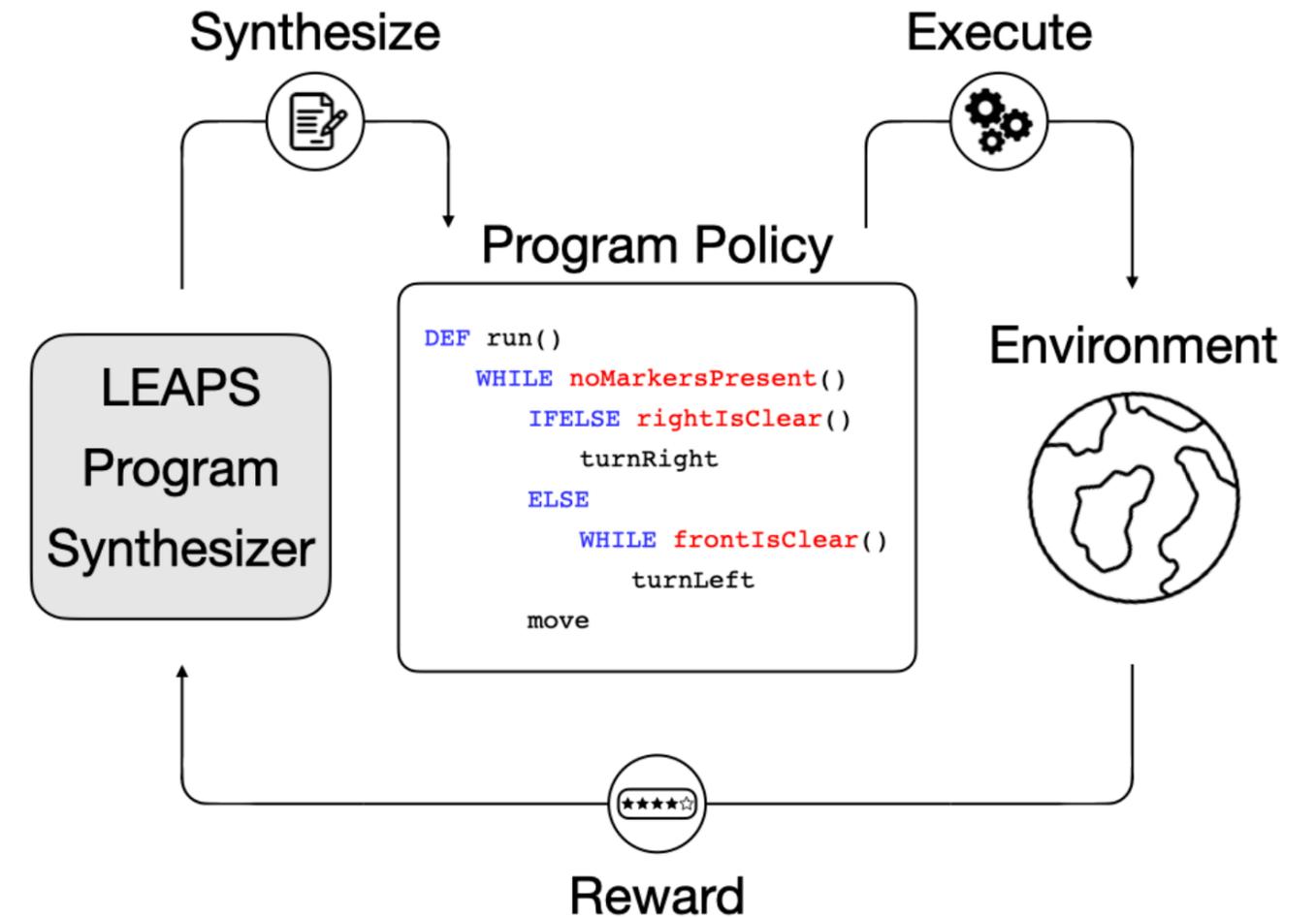
Best Reward

Performance Improvement



Takeaways

- We learn to synthesize a program as a policy purely from reward
- We first learn a program embedding space and then search for a task-solving program
- Our synthesized programs achieve good performance, and are more generalizable and interpretable



Learning to Synthesize Programs as Interpretable and Generalizable Policies

Dweep Trivedi*, Jesse Zhang*, Shao-Hua Sun*, Joseph J Lim

For more details...

Paper and code: clvrai.com/leaps

