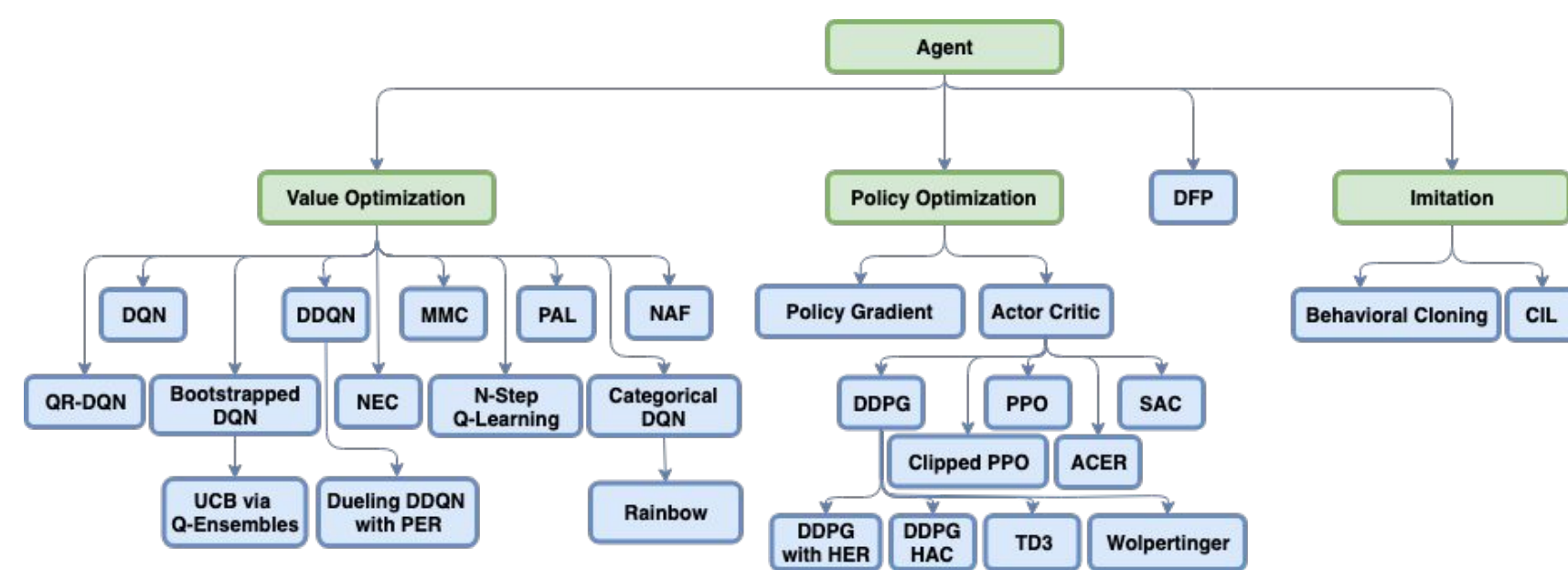
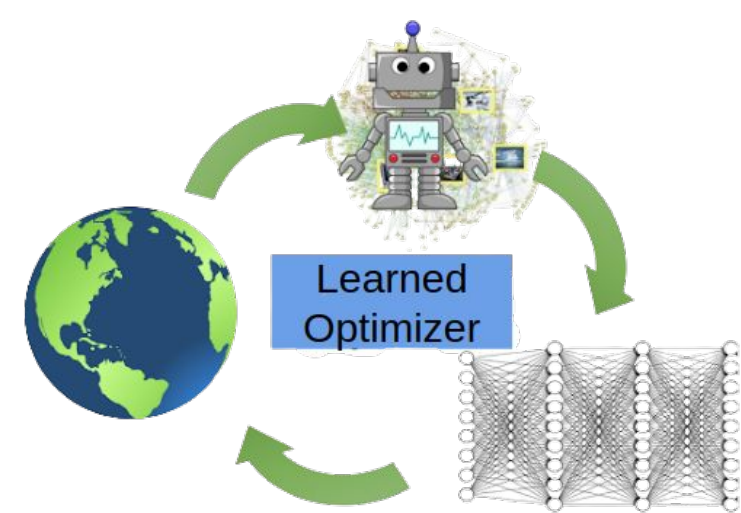


Motivation



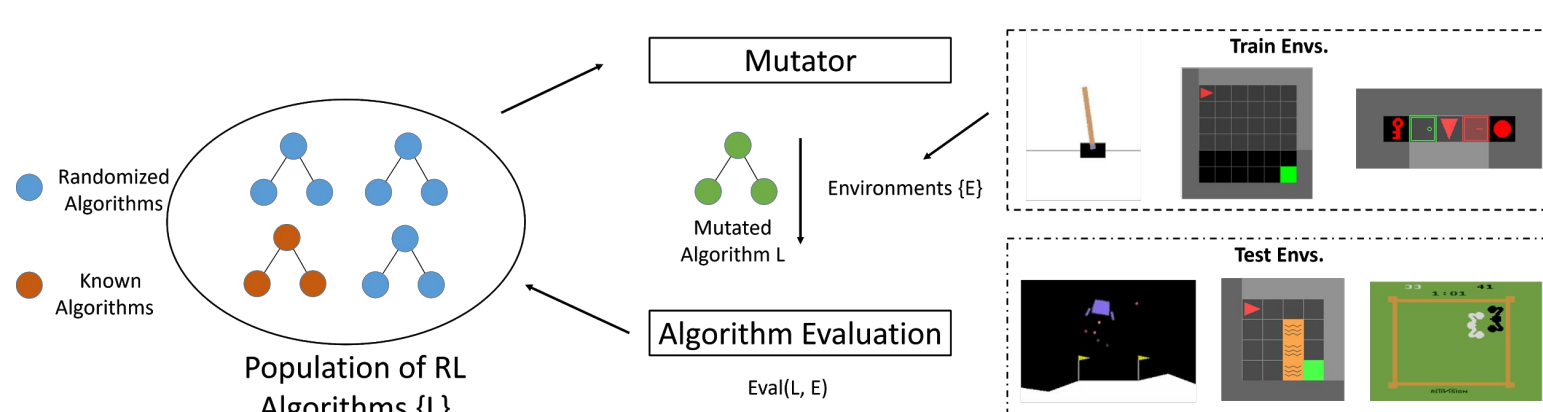
- Desire general purpose RL algorithms without manual effort.
- Problem: Meta-learn RL algorithms that generalize.



- RL algorithm as a learned optimizer
- Meta-learn the optimizer
 - Improved performance
 - Generalizable
 - Interpretable
 - Scale with compute

Overview

- **Insight:** RL algorithm as a computational graph
- **Method:** Evolve population of graphs by mutating, training, and evaluating RL agents
- **Result:** Learn new algorithms which generalize to unseen environments

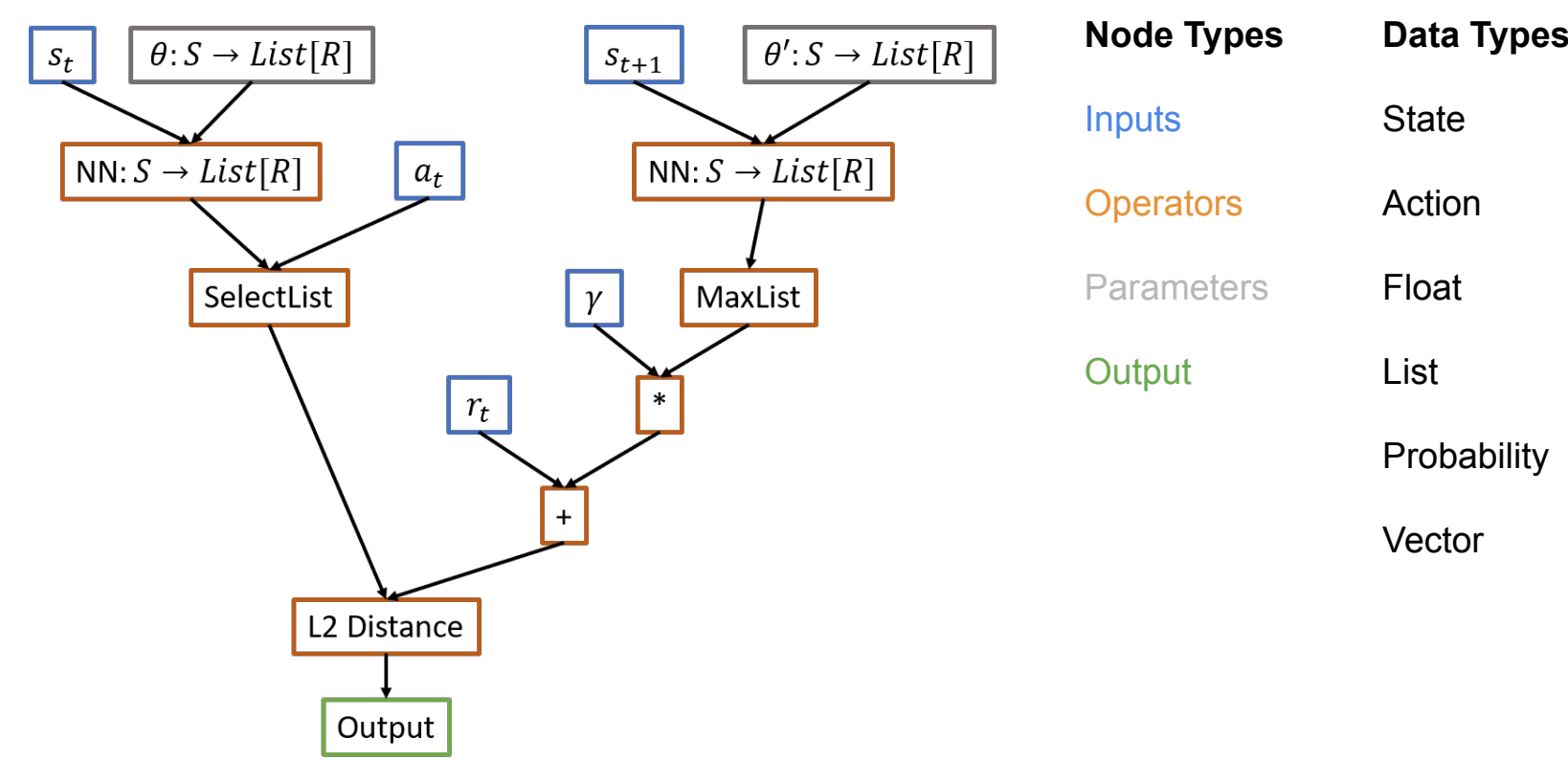


Prior Work

- **Genetic Programming**
 - Holland 1975, Koza 1993, Schmidhuber 1987
 - **AutoML:** Zoph & Le 2016, Hutter 2018, Real et al. 2020
 - Mostly applied to SL
- **Meta-learning in RL**
 - **Adaptation:** Finn & Levine 2018
 - **RNNs:** Duan et al. 2016, Wang et al. 2017
 - Not domain agnostic
- **Learning RL Algorithms**
 - **Metagradients:** Kirsch et al. 2020, Oh et al. 2020
 - Not interpretable
 - **Exploration:** Alet et al. 2020

RL Algorithm as a Computational Graph

Computational graph computes loss function for agent to optimize.



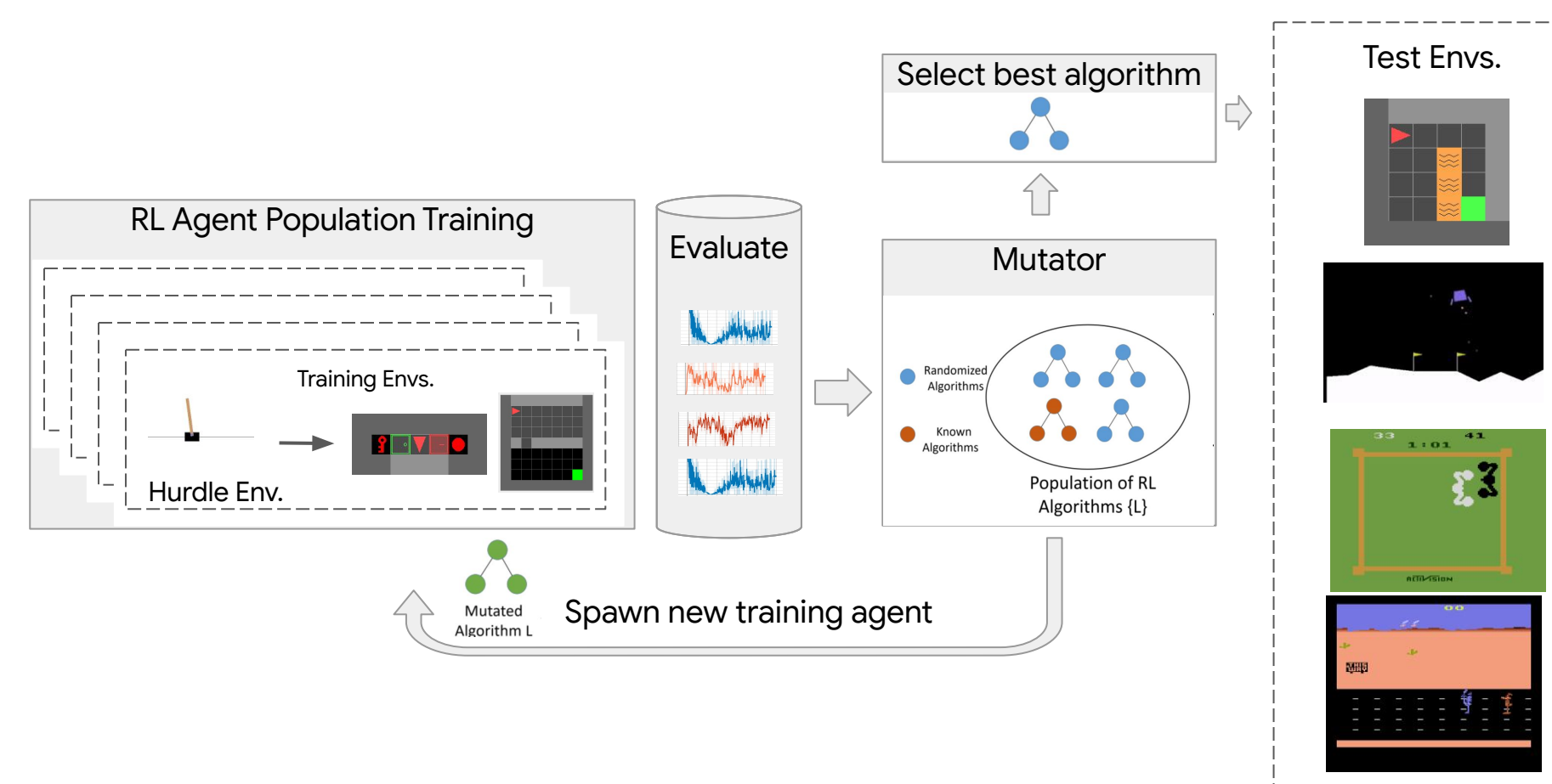
$$L_{DQN} = (Q_{\theta}(s_t, a_t) - (r_t + \gamma * \max_a Q_{\theta'}(s_{t+1}, a)))^2$$

Example graph for DQN loss function which computes Bellman squared error using two Q value networks.

- Directed acyclic graph computes DQN style loss function for value-based agent
- Node types include neural network operators to support more complex architectures
- Data types allows for type checking and ruling out invalid graphs
- Functional equivalence checker skips graphs that are functionally the same
- Representation is expressive, interpretable, and generalizable

Operation	Input Types	Output Type
Add	X, X	X
Subtract	X, X	X
Max	X, X	X
Min	X, X	X
DotProduct	X, X	R
Div	X, X	X
L2Distance	X, X	R
MaxList	List[R]	R
MinList	List[R]	R
ArgMaxList	List[X], Z	Z
SelectList	List[X], Z	X
MeanList	List[X]	X
VarianceList	List[X]	X
Log	X	X
Exp	X	X
Abs	X	X
(CNN: S → List[R])	S	List[R]
(CNN: S → R)	S	R
(CNN: S → V)	S	V
Softmax	List[R]	F
KLDiv	F, F	R
Entropy	F	R
Constant		1, 0.5, 0.2, 0.1, 0.01
MultiplyTenth	X	X
Normal(0, 1)		R
Uniform(0, 1)		R

Evolve Population of RL Algorithms



- Initialize population of ~ 300 RL algorithms with randomized computation graphs
- Evaluate performance by training over set of diverse but cheap envs.
- Mutate most promising algorithms to spawn new agents for evaluation
- Regularized Evolution removes oldest algorithms from population
- Hurdle env. stops bad algorithms early
- Can bootstrap from existing algorithms
- Evolutionary method can scale with compute

Learned Algorithms

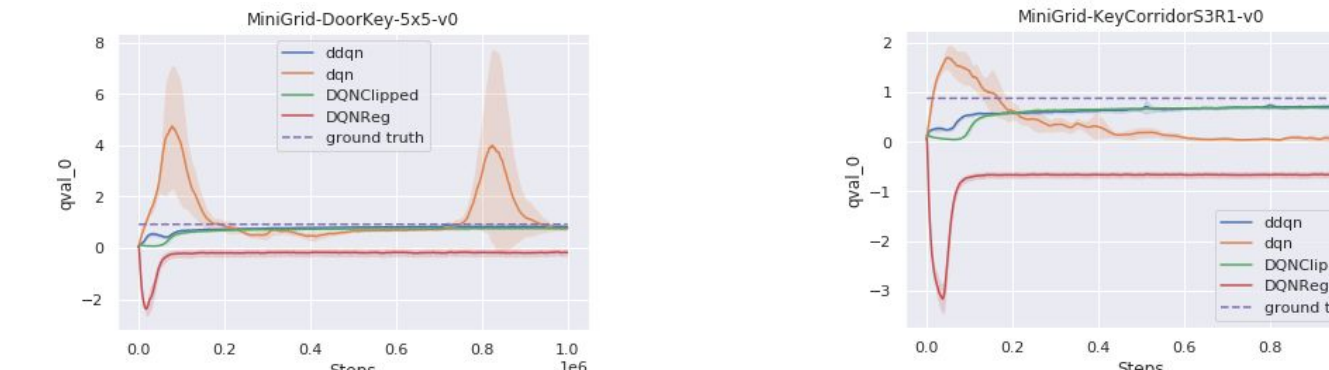
$$Y_t = r_t + \gamma * \max_a Q_{target}(s_t, a), \text{ and } \delta = Q(s_t, a_t) - Y_t.$$

$$L_{DQNClipped} = \max [Q(s_t, a_t), \delta^2 + Y_t] + \max [Q(s_t, a_t) - Y_t, \gamma(\max_a Q_{target}(s_t, a))^2]$$

- DQNClipped as constrained optimization

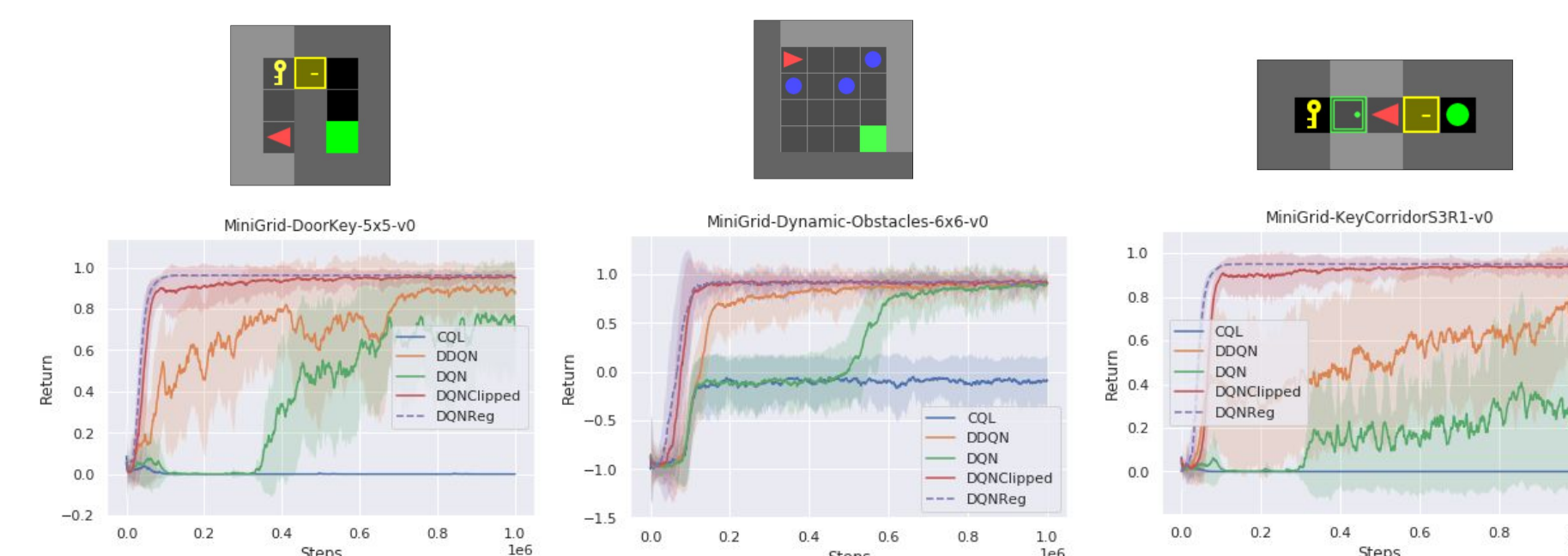
$$L_{DQNRReg} = 0.1 * Q(s_t, a_t) + \delta^2$$

- DQNRReg as entropy regularization

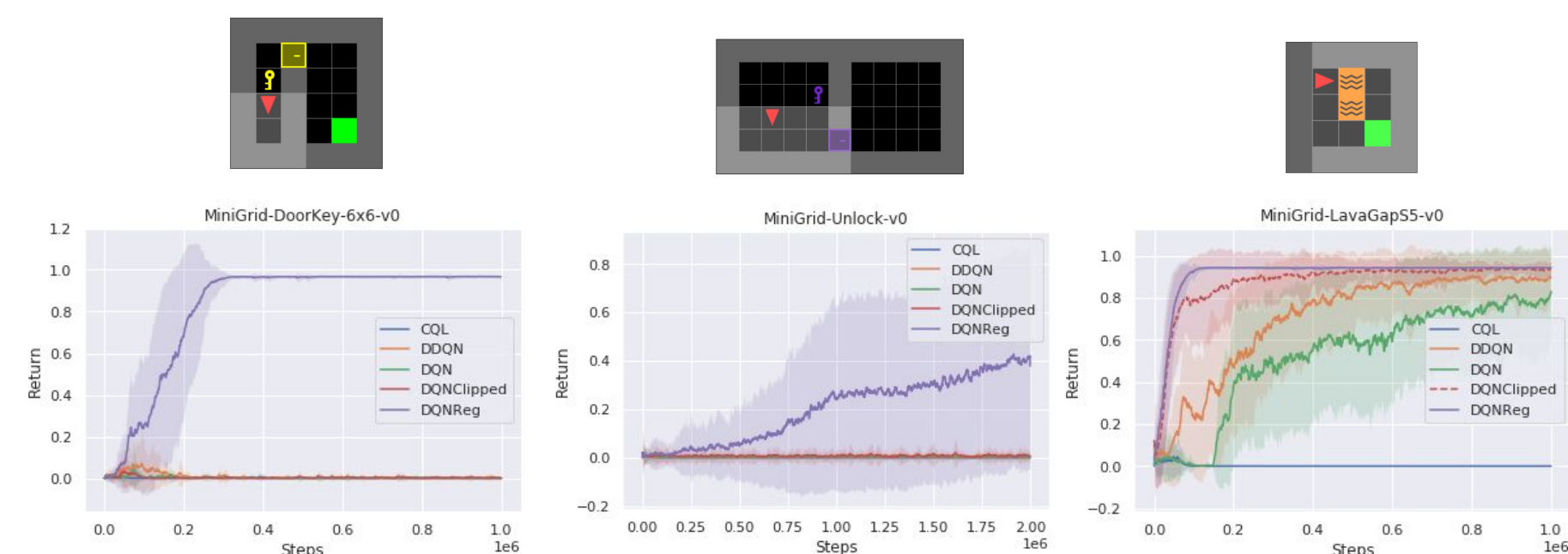


- Learned algorithms prevent value overestimation in different ways

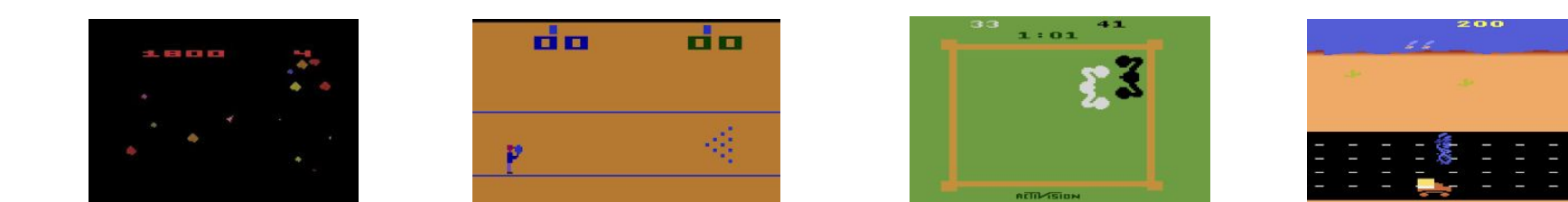
Results



- Outperform baselines on train envs.



- Generalize to unseen environments.



Env	DQN	DDQN	PPO	DQNRReg
Asteroid	1364.5	734.7	2097.5	2390.4
Bowling	50.4	68.1	40.1	80.5
Boxing	88.0	91.6	94.6	100.0
RoadRunner	39544.0	44127.0	35466.0	65516.0

- Benefits on Atari even though training envs. were non-image based

Future Work

- Extensions to actor critic, offline RL, representation learning
- Analyze and incorporate learned algorithms into existing ones
- Machine assisted algorithm development

