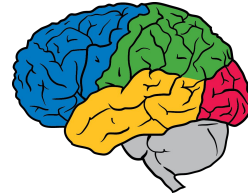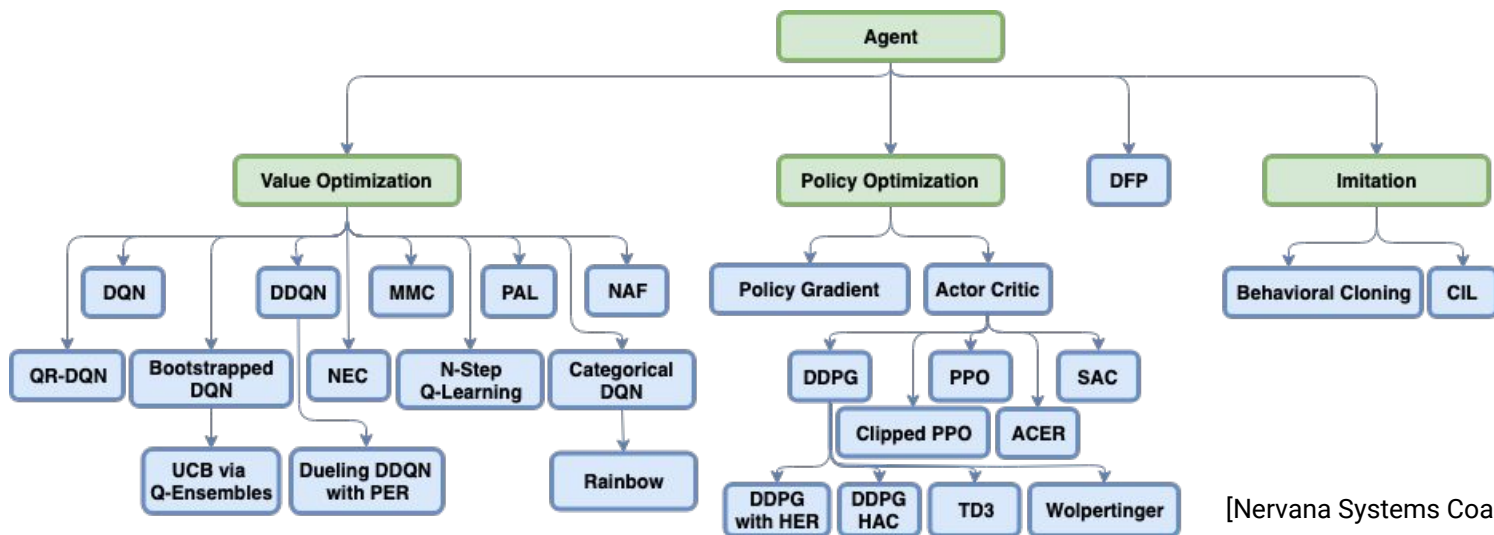# Evolving Reinforcement Learning Algorithms

JD Co-Reyes, Yingjie Miao, Daiyi Peng, Esteban Real, Sergey Levine, Quoc V. Le, Honglak Lee, Aleksandra Faust
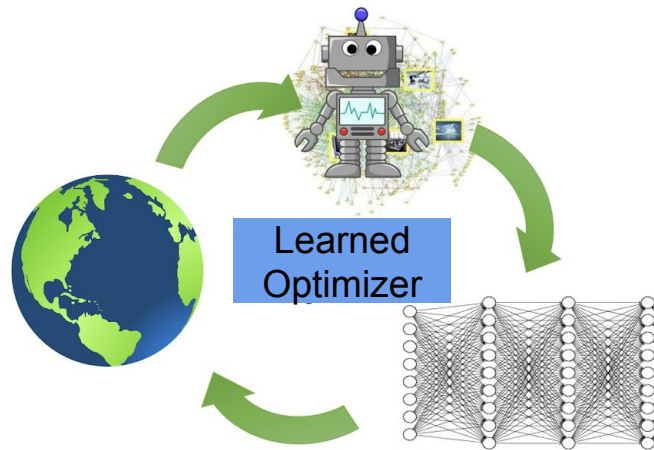
# Wide Choice of RL Algorithms



[Nervana Systems Coach Diagram]

Desire: General purpose RL algorithms without manual effort.
Problem: Can we meta-learn RL algorithms that generalize well on unseen tasks?

Google Research

# RL Algorithm as a Learned Optimizer



reinforcement learning

Learned Optimizer

- Learning procedure which takes in MDP and transforms experience into optimal behavior

- Can we meta-learn the optimizer?

  ○ Improved performance

  ○ Generalize to unseen environments

  ○ Interpretable

  ○ Scale with data and compute

Google Research

# Example: Simple Modifications to Existing Algorithms

$$\delta^2 = (Q(s_t, a_t) - (r_t + \gamma * \max_a Q(s_{t+1}, a)))^2$$

[1] Kumar, A., Zhou, A., Tucker, G., & Levine, S. (2020). Conservative Q-Learning for Offline Reinforcement Learning. *ArXiv, abs/2006.04779*.

Google Research

# Example: Simple Modifications to Existing Algorithms

CQL: adds scaled log softmax policy to TD error

$$\delta^2 + \beta \log \sum_{a} \exp\left(Q(s_t, a)\right) - Q(s_t, a_t)$$

[1] Kumar, A., Zhou, A., Tucker, G., & Levine, S. (2020). Conservative Q-Learning for Offline Reinforcement Learning. *ArXiv, abs/2006.04779*.

Google Research

# Example: Simple Modifications to Existing Algorithms

CQL: adds scaled log softmax policy to TD error

$$\delta^2 + \beta \log \sum_a \exp\left(Q(s_t, a)\right) - Q(s_t, a_t)$$
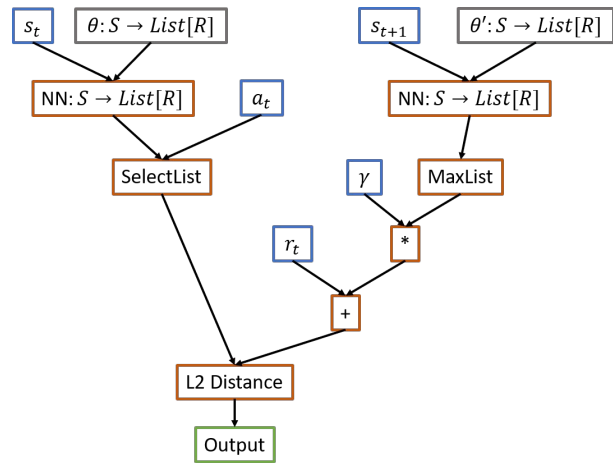
M-DQN: adds scaled log policy to reward

$$\hat{q}_{\text{m-dqn}}(r_t, s_{t+1}) = r_t + \alpha\tau \ln \pi_{\bar{\theta}}(a_t|s_t) + \gamma \sum_{a' \in \mathcal{A}} \pi_{\bar{\theta}}(a'|s_{t+1})\left(q_{\bar{\theta}}(s_{t+1}, a') - \tau \ln \pi_{\bar{\theta}}(a'|s_{t+1})\right)$$

[1] Kumar, A., Zhou, A., Tucker, G., & Levine, S. (2020). Conservative Q-Learning for Offline Reinforcement Learning. *ArXiv, abs/2006.04779*.
[2] Vieillard, N., Pietquin, O., & Geist, M. (2020). Munchausen Reinforcement Learning. *ArXiv, abs/2007.14430*.
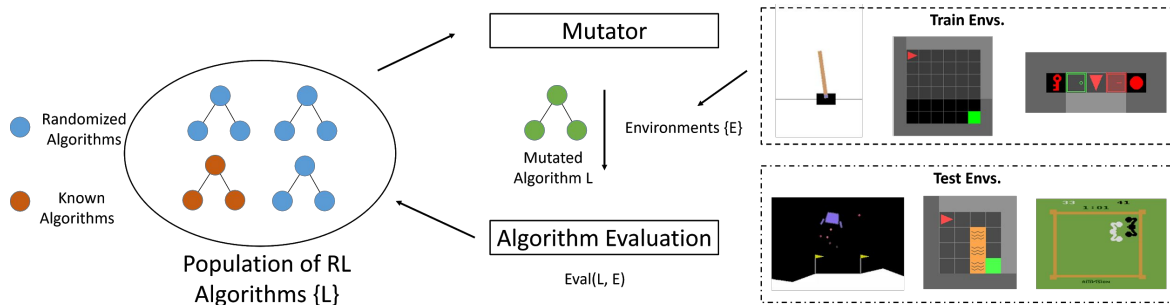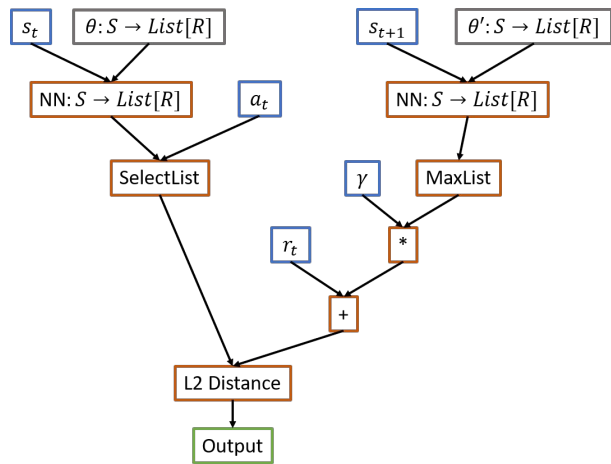
Google Research

# Evolving RL Algorithms

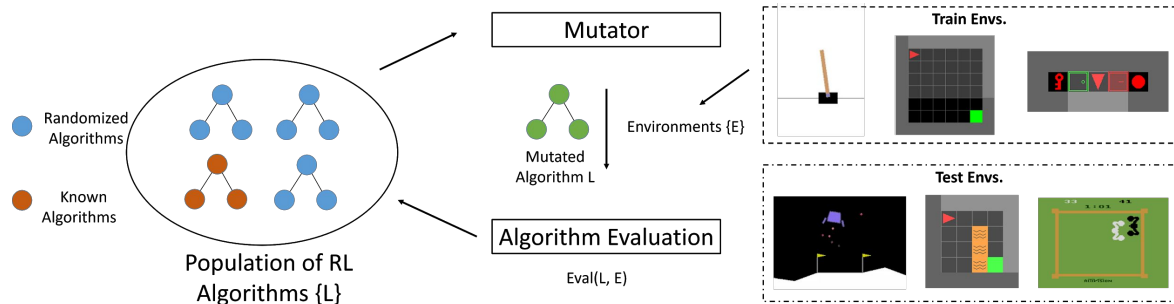- **Insight:** RL algorithm as a computational graph

Google Research

# Evolving RL Algorithms

- **Insight:** RL algorithm as a computational graph
- **Method:** Evolve population of graphs by mutating, training, and evaluating RL agents

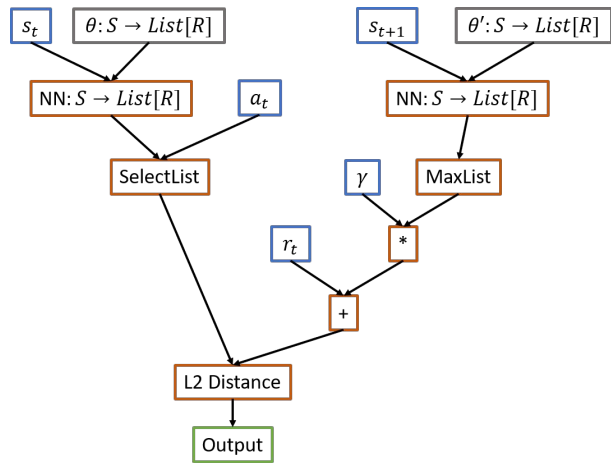# Evolving RL Algorithms

- **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

Google Research

# 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
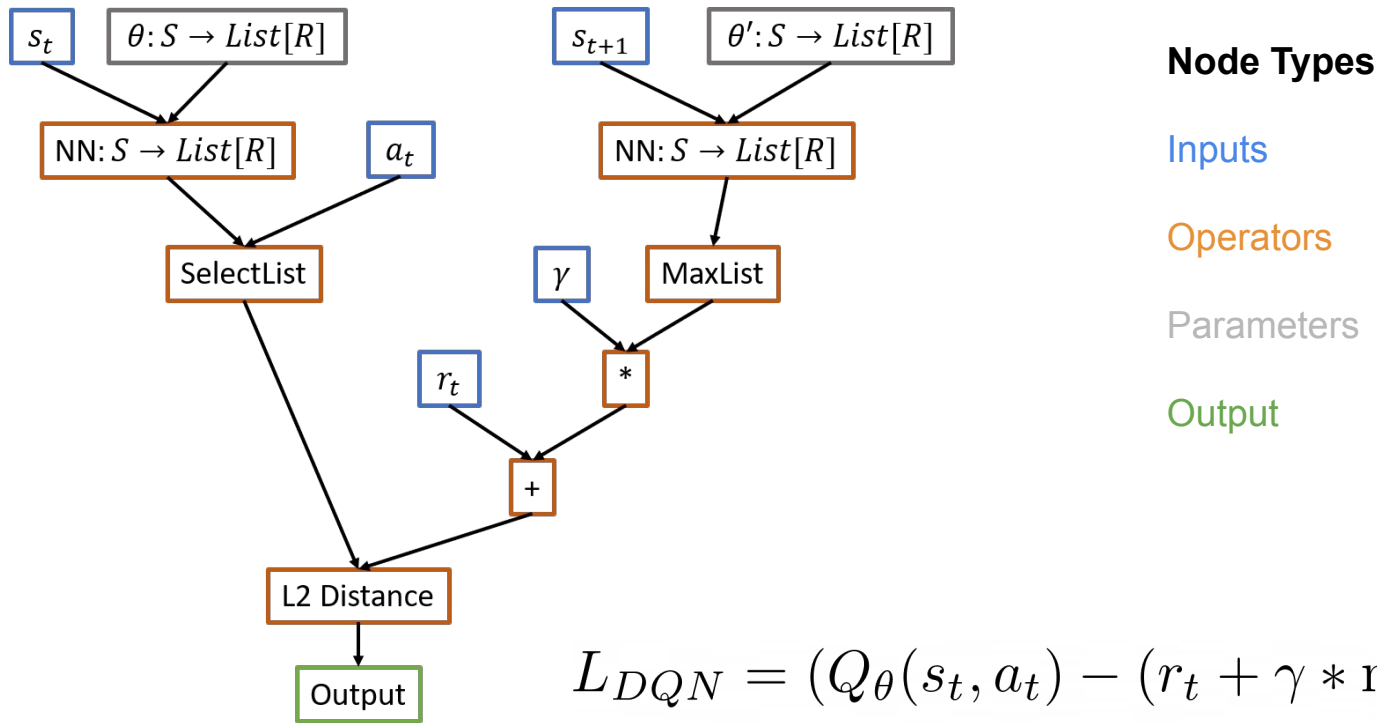  - Not interpretable
  - **Exploration**: Alet et al. 2020
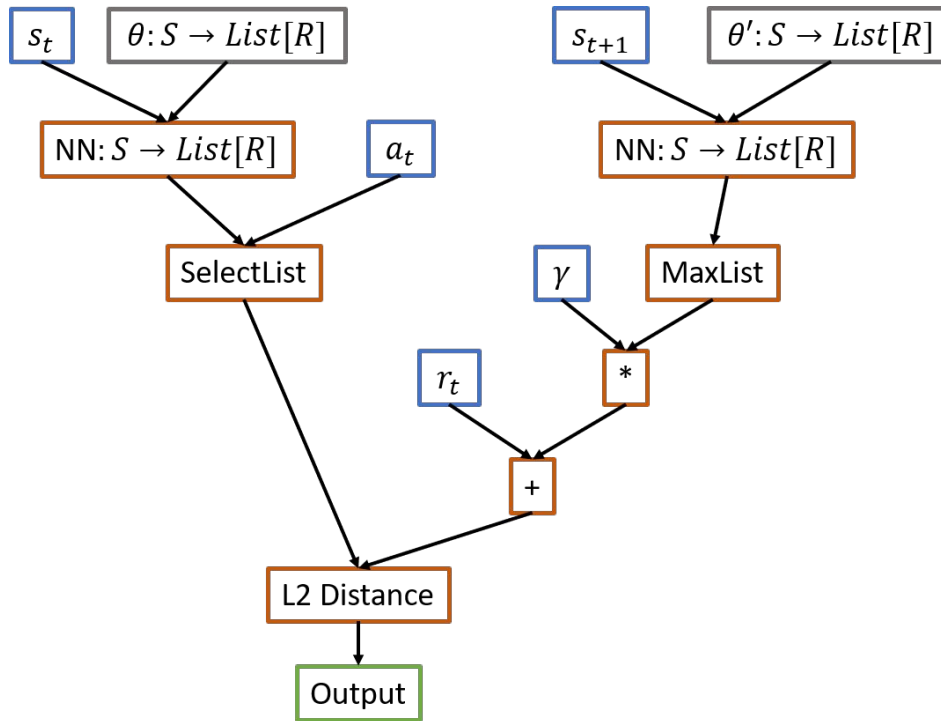
# Algorithm Representation

**Expressive**                    **Interpretable**                    **Generalizable**

# RL Algorithm as a Computational Graph



**Node Types**

Inputs

Operators

Parameters

Output

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

Google Research

# RL Algorithm as a Computational Graph



Typing allows for domain agnostic programs and type checking

Google Research

# RL Algorithm as a Computational Graph



| Operation | Input Types | Output Type |
|---|---|---|
| Add | $\mathbb{X}, \mathbb{X}$ | $\mathbb{X}$ |
| Subtract | $\mathbb{X}, \mathbb{X}$ | $\mathbb{X}$ |
| Max | $\mathbb{X}, \mathbb{X}$ | $\mathbb{X}$ |
| Min | $\mathbb{X}, \mathbb{X}$ | $\mathbb{X}$ |
| DotProduct | $\mathbb{X}, \mathbb{X}$ | $\mathbb{R}$ |
| Div | $\mathbb{X}, \mathbb{X}$ | $\mathbb{X}$ |
| L2Distance | $\mathbb{X}, \mathbb{X}$ | $\mathbb{R}$ |
| MaxList | $List[\mathbb{R}]$ | $\mathbb{R}$ |
| MinList | $List[\mathbb{R}]$ | $\mathbb{R}$ |
| ArgMaxList | $List[\mathbb{R}]$ | $\mathbb{Z}$ |
| SelectList | $List[\mathbb{X}], \mathbb{Z}$ | $\mathbb{X}$ |
| MeanList | $List[\mathbb{X}]$ | $\mathbb{X}$ |
| VarianceList | $List[\mathbb{X}]$ | $\mathbb{X}$ |
| Log | $\mathbb{X}$ | $\mathbb{X}$ |
| Exp | $\mathbb{X}$ | $\mathbb{X}$ |
| Abs | $\mathbb{X}$ | $\mathbb{X}$ |
| (C)NN:$\mathbb{S} \to List[\mathbb{R}]$ | $\mathbb{S}$ | $List[\mathbb{R}]$ |
| (C)NN:$\mathbb{S} \to \mathbb{R}$ | $\mathbb{S}$ | $\mathbb{R}$ |
| (C)NN:$\mathbb{S} \to \mathbb{V}$ | $\mathbb{V}$ | $\mathbb{V}$ |
| Softmax | $List[\mathbb{R}]$ | $\mathbb{P}$ |
| KLDiv | $\mathbb{P}, \mathbb{P}$ | $\mathbb{R}$ |
| Entropy | $\mathbb{P}$ | $\mathbb{R}$ |
| Constant | | 1, 0.5, 0.2, 0.1, 0.01 |
| MultiplyTenth | $\mathbb{X}$ | $\mathbb{X}$ |
| Normal(0, 1) | | $\mathbb{R}$ |
| Uniform(0, 1) | | $\mathbb{R}$ |

Google Research

# Outer loop Optimization

How to scale with
compute?

# Evolving RL Algorithms



Evolving Reinforcement Learning Algorithms, Co-Reyes, Miao, Peng, Real, Levine, Le, Lee, Faust, ICLR 2021 https://sites.google.com/view/evolvingrl

Google Research

# Meta-Learn RL Algorithms



Randomized Algorithms

Known Algorithms

Population of RL Algorithms {L}

Mutator

Mutated Algorithm L

Environments {E}

Algorithm Evaluation

Eval(L, E)

**Algorithm 1** Algorithm Evaluation, $\mathrm{Eval}(L, \mathcal{E})$

1: **Input:** RL Algorithm $L$, Environment $\mathcal{E}$, training episodes $M$
2: **Initialize:** Q-value parameters $\theta$, target parameters $\theta'$ empty replay buffer $\mathcal{D}$
3: **for** $i = 1$ **to** $M$ **do**
4:      **for** $t = 0$ **to** $T$ **do**
5:          With probability $\epsilon$, select a random action $a_t$,
6:          otherwise select $a_t = \arg\max_a Q(s_t, a)$
7:          Step environment $s_{t+1}, r_t \sim \mathcal{E}(a_t, s_t)$
8:          $\mathcal{D} \leftarrow \mathcal{D} \cup \{s_t, a_t, r_t, s_{t+1}\}$
9:          Update parameters $\theta \leftarrow \theta - \nabla_\theta L(s_t, a_t, r_t, s_{t+1}, \theta, \gamma)$
10:          Update target $\theta' \leftarrow \theta$
11:      **end for**
12:      Compute episode return $R_m = \sum_{t=0}^T r_t$
13: **end for**
14: **Output:**
15:      Normalized training performance $\frac{1}{M} \sum_{m=1}^M \frac{R_m - R_{min}}{R_{max} - R_{min}}$

- Learn loss function for DQN style update procedure
- Score each algorithm with normalized training performance

Google Research

# Meta-Learn RL Algorithms



Mutator

Environments {E}

Mutated Algorithm L

Randomized Algorithms

Known Algorithms

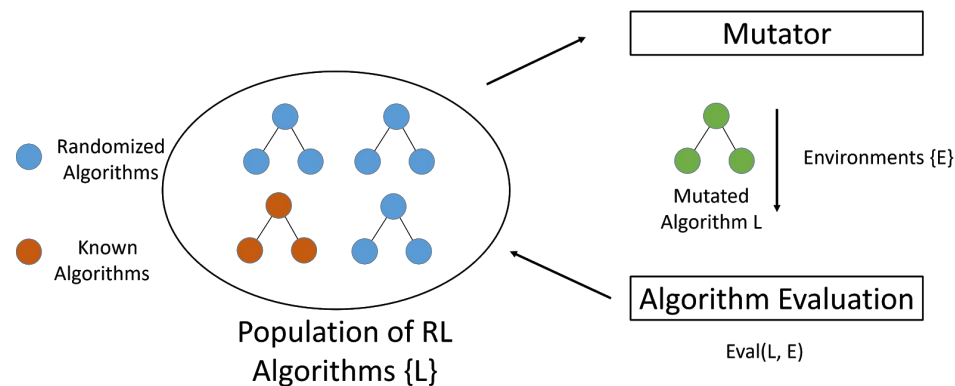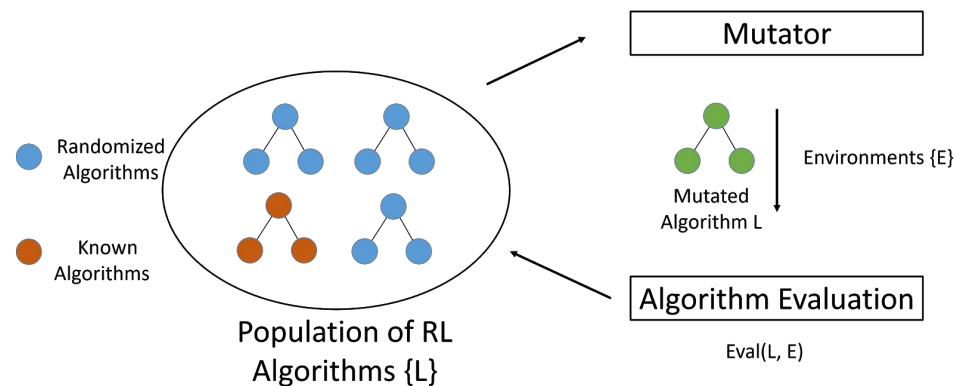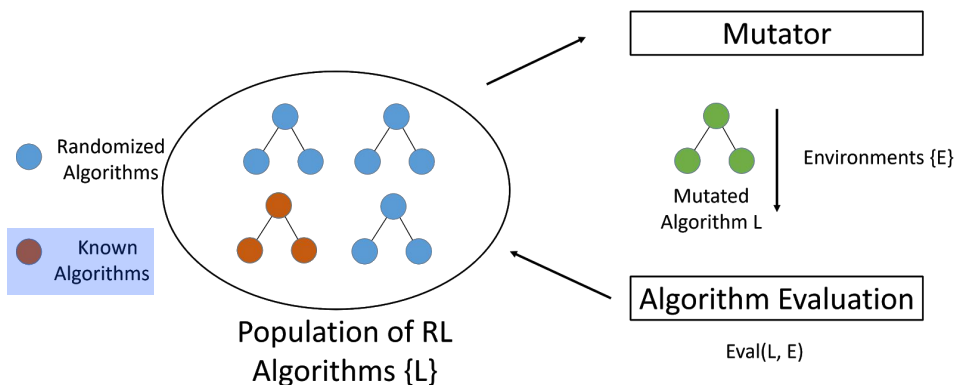Population of RL Algorithms {L}

Algorithm Evaluation

Eval(L, E)

**Algorithm 2** Evolving RL Algorithms

1: **Input:** Training environments $\{\mathcal{E}\}$, hurdle environment $\mathcal{E}_h$, hurdle threshold $\alpha$, optional existing algorithm $A$
2: **Initialize:** Population $P$ of RL algorithms $\{L\}$, history $H$, randomized inputs $I$. If bootstrapping, initialize $P$ with $A$.
3: Score each L in $P$ with $H[L].score \leftarrow \sum_{\mathcal{E}} \text{Eval}(L, \mathcal{E})$
4: **for** $c = 0$ **to C do**
5:     Sample tournament $T \sim Uniform(P)$
6:     Parent algorithm $L \leftarrow$ highest score algorithm in $T$
7:     Child algorithm $L' \leftarrow \text{Mutate}(L)$
8:     $H[L'].hash \leftarrow \text{Hash}(L'(I))$
9:     **if** $H[L'].hash$ was new **and** $\text{Eval}(L', \mathcal{E}_h) > \alpha$ **then**
10:         $H[L'].score \leftarrow \sum_{\mathcal{E}} \text{Eval}(L', \mathcal{E})$
11:     **end if**
12:     Add $L'$ to population $P$
13:     Remove oldest $L$ from population $P$
14: **end for**
15: **Output:** Algorithm L with highest score

- Regularized Evolution for outer loop optimization

Google Research

# Optimizations



**Mutator**

Mutated Algorithm L

Environments {E}

**Algorithm Evaluation**

Eval(L, E)

Randomized Algorithms

Known Algorithms

Population of RL Algorithms {L}

---

**Algorithm 2** Evolving RL Algorithms

1: **Input:** Training environments $\{\mathcal{E}\}$, hurdle environment $\mathcal{E}_h$, hurdle threshold $\alpha$, optional existing algorithm $A$.
2: **Initialize:** Population $P$ of RL algorithms $\{L\}$, history $H$, randomized inputs $I$. If bootstrapping, initialize $P$ with $A$.
3: Score each L in $P$ with $H[L].score \leftarrow \sum_{\mathcal{E}} \text{Eval}(L, \mathcal{E})$
4: **for** $c = 0$ **to C do**
5:     Sample tournament $T \sim Uniform(P)$
6:     Parent algorithm $L \leftarrow$ highest score algorithm in $T$
7:     Child algorithm $L' \leftarrow$ Mutate(L)
8:     $H[L'].hash \leftarrow \text{Hash}(L'(I))$
9:     **if** $H[L'].hash$ was new **and** $\text{Eval}(L', \mathcal{E}_h) > \alpha$ **then**
10:         $H[L'].score \leftarrow \sum_{\mathcal{E}} \text{Eval}(L', \mathcal{E})$
11:     **end if**
12:     Add $L'$ to population $P$
13:     Remove oldest $L$ from population
14: **end for**
15: **Output:** Algorithm L with highest score
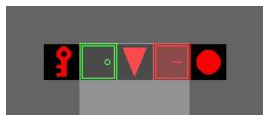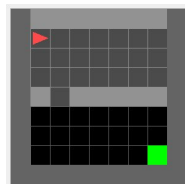
- Don't reevaluate functionally equivalent or duplicate programs
- Saves 70% of computation

Google Research

# Optimizations



**Algorithm 2** Evolving RL Algorithms

1: **Input:** Training environments $\{\mathcal{E}\}$, hurdle environment $\mathcal{E}_h$, hurdle threshold $\alpha$, optional existing algorithm $A$
2: **Initialize:** Population $P$ of RL algorithms $\{L\}$, history $H$, randomized inputs $I$. If bootstrapping, initialize $P$ with $A$.
3: Score each L in $P$ with $H[L].score \leftarrow \sum_{\mathcal{E}} \text{Eval}(L, \mathcal{E})$
4: **for** $c = 0$ **to C do**
5:     Sample tournament $T \sim Uniform(P)$
6:     Parent algorithm $L \leftarrow$ highest score algorithm in $T$
7:     Child algorithm $L' \leftarrow \text{Mutate}(L)$
8:     $H[L'].hash \leftarrow \text{Hash}(L'(I))$
9:     **if** $H[L'].hash$ was new **and** $\text{Eval}(L', \mathcal{E}_h) > \alpha$ **then**
10:         $H[L'].score \leftarrow \sum_{\mathcal{E}} \text{Eval}(L', \mathcal{E})$
11:     **end if**
12:     Add $L'$ to population $P$
13:     Remove oldest $L$ from population
14: **end for**
15: **Output:** Algorithm L with highest score

- Stop early if performance on hurdle environment is bad
- Saves additional 30% of computation

Google Research

# Bootstrap from existing algorithms



**Algorithm 2** Evolving RL Algorithms

1: **Input:** Training environments $\{\mathcal{E}\}$, hurdle environment $\mathcal{E}_h$, hurdle threshold $\alpha$, optional existing algorithm $A$
2: **Initialize:** Population $P$ of RL algorithms $\{L\}$, history $H$, randomized inputs $I$. If bootstrapping, initialize $P$ with $A$.
3: Score each L in $P$ with $H[L].score \leftarrow \sum_{\mathcal{E}} \text{Eval}(L, \mathcal{E})$
4: **for** $c = 0$ **to C do**
5:     Sample tournament $T \sim Uniform(P)$
6:     Parent algorithm $L \leftarrow$ highest score algorithm in $T$
7:     Child algorithm $L' \leftarrow \text{Mutate}(L)$
8:     $H[L'].hash \leftarrow \text{Hash}(L'(I))$
9:     **if** $H[L'].hash$ was new **and** $\text{Eval}(L', \mathcal{E}_h) > \alpha$ **then**
10:         $H[L'].score \leftarrow \sum_{\mathcal{E}} \text{Eval}(L', \mathcal{E})$
11:     **end if**
12:     Add $L'$ to population $P$
13:     Remove oldest $L$ from population
14: **end for**
15: **Output:** Algorithm L with highest score

- Can initialize population with existing algorithms

Google Research

# Environments



Train

Test

Hurdle Env.

- Want training environments that are computationally cheap but diverse
- Test environments include completely different state and action sizes (including image observations)

Google Research

# Results

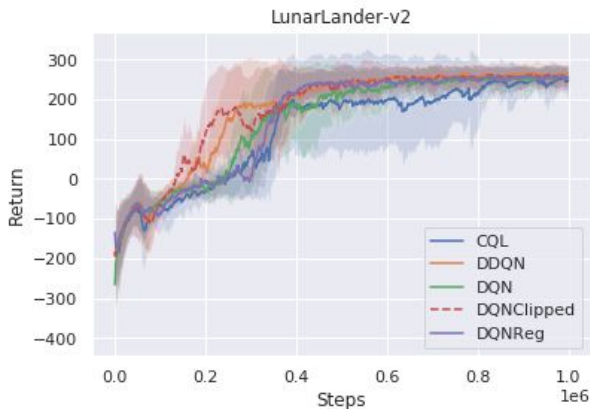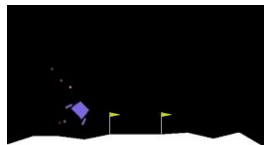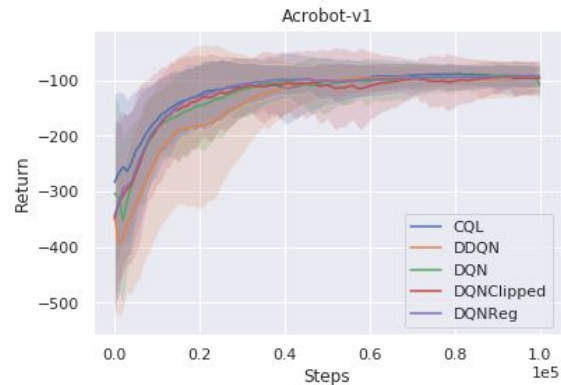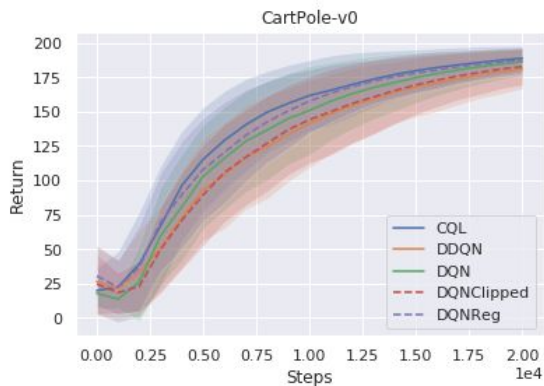# Learned Algorithm 1: DQN_Clipped as Constrained Optimization

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

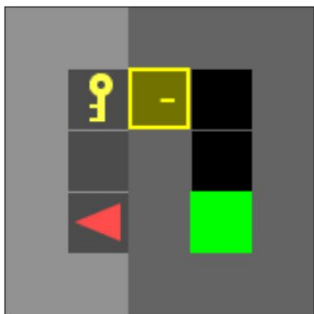$$L_{\text{DQNClipped}} = \max\left[Q(s_t, a_t), \delta^2 + Y_t\right] + \max\left[Q(s_t, a_t) - Y_t, \gamma(\max_a Q_{targ}(s_t, a))^2\right]$$

Google Research

# Learned Algorithm 1: DQN_Clipped as Constrained Optimization

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

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

Case 2: $Q(s_t, a_t) - Y_t > \delta^2$

- Minimize Q

Case 3: $Q(s_t, a_t) - Y_t \leq \delta^2$

- Minimize normal TD error



MiniGrid-DoorKey-5x5-v0

case1
case2
case3
case4

Google Research

# Learned Algorithm 1: DQN_Clipped as Constrained Optimization

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

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

Case 2: $Q(s_t, a_t) - Y_t > \delta^2$

- Minimize Q

Case 3: $Q(s_t, a_t) - Y_t \leq \delta^2$

- Minimize normal TD error



MiniGrid-DoorKey-5x5-v0

Update Rule Proportion

case1
case2
case3
case4

Steps

Google Research

# Learned Algorithm 2: DQN_Reg as Soft Constraint

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

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

Google Research

# Learned Algorithm 2: DQN_Reg as Soft Constraint

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

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

$$L_{CQL} = \beta \log \sum_a \exp\left(Q(s_t, a)\right) - Q(s_t, a_t) + \delta^2$$

Google Research

# Learned Algorithm 2: DQN_Reg as Soft Constraint

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

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

$$L_{CQL} = \beta \log \sum_a \exp\left(Q(s_t, a)\right) - Q(s_t, a_t) + \delta^2$$

- DQNReg as version of entropy regularization that penalizes Q-values on dataset to prevent overfitting

Google Research

# Generalize to Unseen Environments

Google Research

# DQNReg Outperforms on Sparse Reward Train Envs.



use the key to open the door and then get to the goal

get to the green goal square

pick up the green ball

# DQNReg Generalizes to Sparse Reward Test Envs.



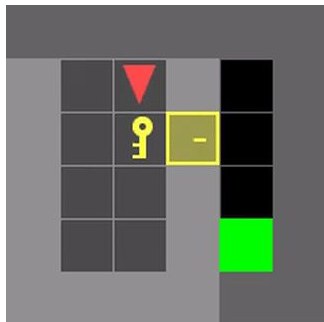use the key to open the door and then get to the goal



open the door


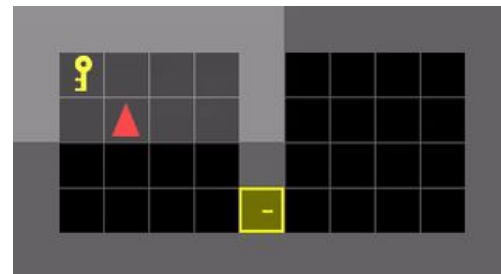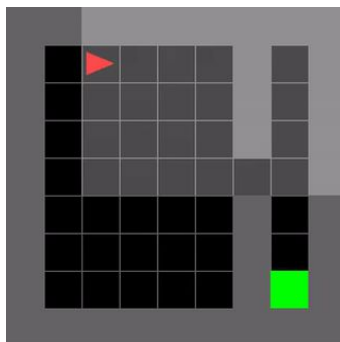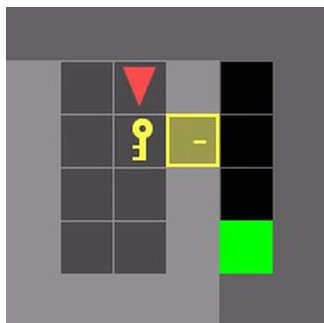
avoid the lava and get to the green goal square



MiniGrid-DoorKey-6x6-v0



MiniGrid-Unlock-v0



MiniGrid-LavaGapS5-v0

Google Research

# Generalize to Unseen Environments

DQN

DQNReg

Google Research

# Atari Performance



| Env | DQN | DDQN | PPO | DQNReg |
|---|---|---|---|---|
| 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** |

*Baselines taken from reported numbers.*

Learned algorithm (DQNReg) generalizes to Atari games when meta-training was on non-image based environments. Not tuned to Atari games.
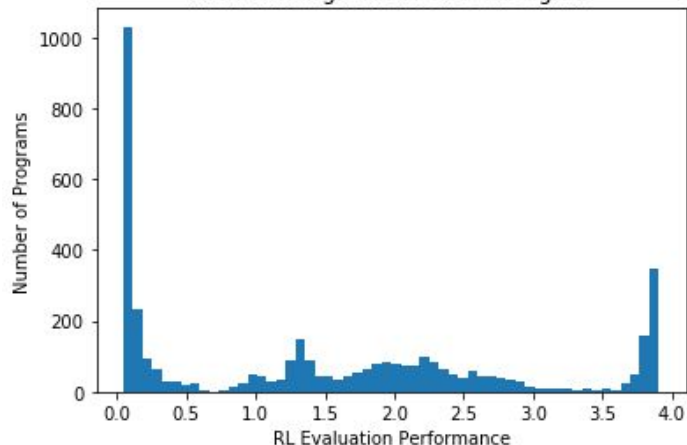
Google Research

# DQNReg



MiniGrid-DoorKey-5x5-v0

DQN overestimates Q values while learned algorithms DQNClipped and DQNReg overcome this issue and underestimate Q values

Google Research

# Learning Convergence



Meta-Training Performance Histogram

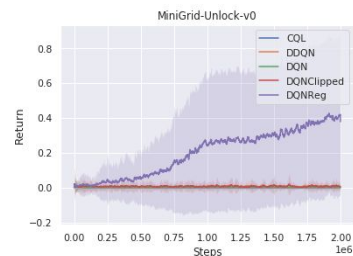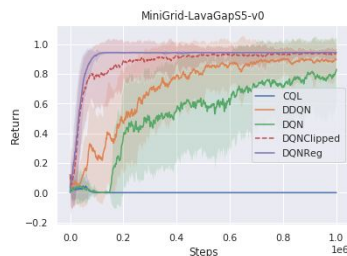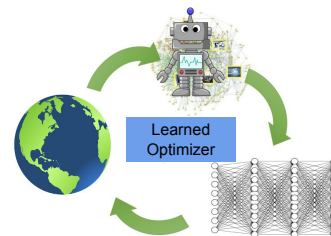| Raw Equation | Simplified Equation | Score | Rank |
|---|---|---|---|
| $\delta^2 + 0.1 * Q(s_t, a_t) + r_t - (\gamma * Q_{targ} - 0.1 * Q(s_t, a_t))$ | $\delta^2 + 0.2 * Q(s_t, a_t)$ | 3.903 | 2 |
| $\delta^2 + 0.1 * Q(s_t, a_t) - \gamma + Q_{targ}$ | $\delta^2 + 0.1 * Q(s_t, a_t)$ | 3.902 | 3 |
| $\delta^2 + ((r_t + \gamma * Q_{targ} + Q(s_t, a_t)) * (\gamma - \max(\gamma, 0.1 * Q(s_t, a_t))$ $- \gamma * Q_{targ} - 0.1 * Q(s_t, a_t))$ | NA | 3.846 | 11146 |
| $\delta^2 + (\delta^2 + 0.1 * Q(s_t, a_t))^2$ | NA | 3.65 | 12146 |
| $\delta^2 + Q(s_t, a_t)$ | $\delta^2 + Q(s_t, a_t)$ | 2.8 | 12446 |
| $\delta^2$ | $\delta^2$ | 2.28 | 13246 |

Top performing algorithms have similar structure

With different training environments or initialization, could find other families of models with better performance

Google Research

# Conclusion

- RL algorithm as a computational graph

- Evolve new RL algorithms

- Learned algorithms generalize to unseen environments
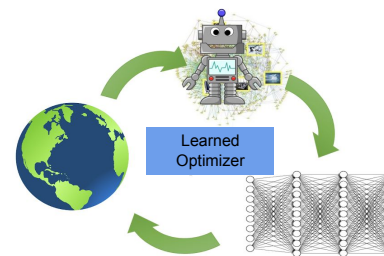
reinforcement learning

Learned Optimizer

MiniGrid-LavaGapS5-v0

MiniGrid-Unlock-v0

| Env | DQN | DDQN | PPO | DQNReg |
|-----|-----|------|-----|--------|
| 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** |

Google Research

# Discussion



reinforcement learning

Learned Optimizer

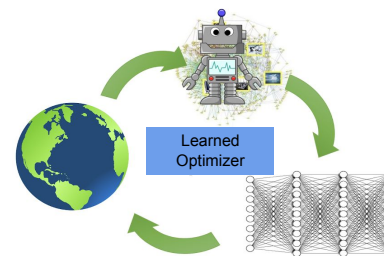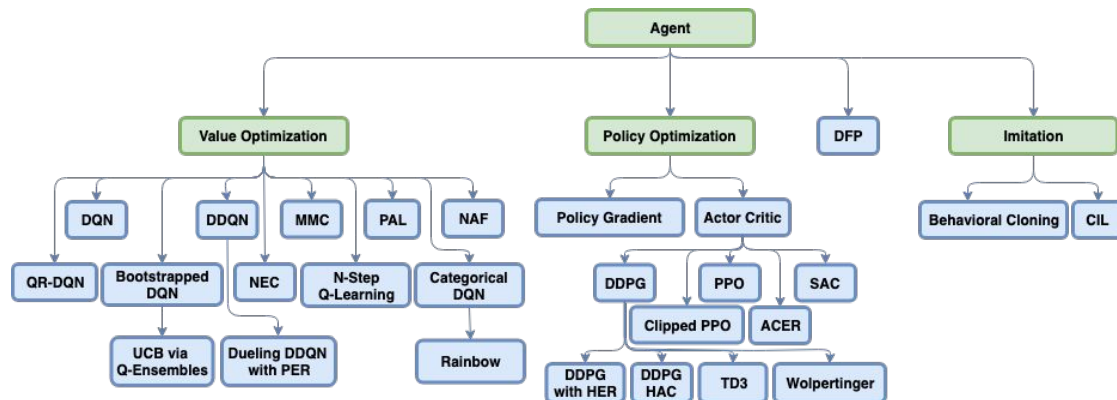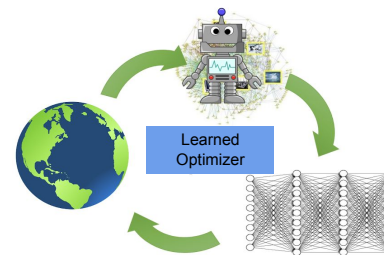- Incorporate learned modifications into existing algorithms

Google Research

# Discussion

reinforcement learning



- Incorporate learned modifications into existing algorithms

- Machine assisted algorithm development
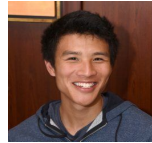
Google Research

# Discussion



reinforcement learning

- Incorporate learned modifications into existing algorithms

- Machine assisted algorithm development



- Extend to other families: actor critic, offline RL

Google Research

# Thank you to collaborators!



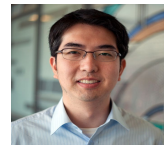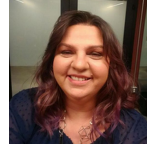JD Co-Reyes    Yingjie Miao    Daiyi Peng    Esteban Real

Sergey Levine    Quoc Le    Honglak Lee    Aleksandra Faust

# Questions?

BAIR
BERKELEY ARTIFICIAL INTELLIGENCE RESEARCH

Google Research