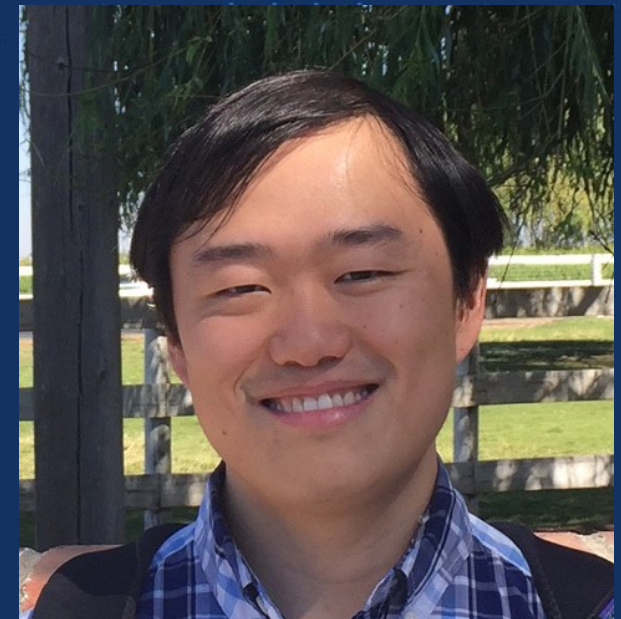


# Finding Good Representation for Search and Exploration in RL

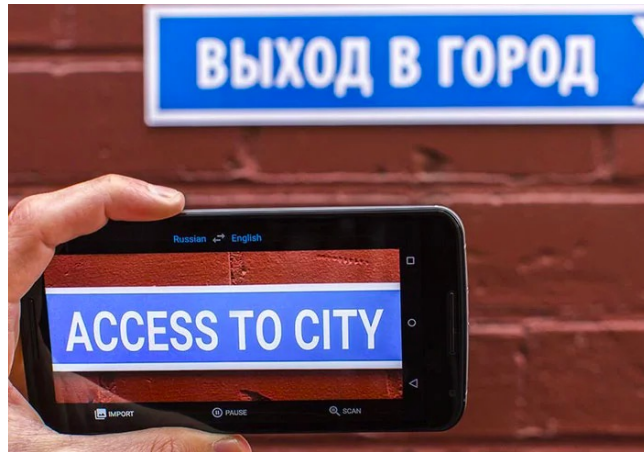
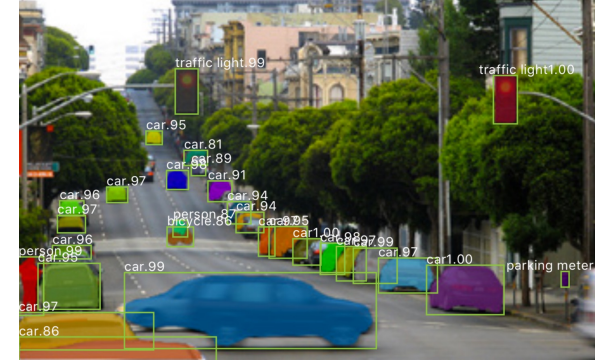
Yuandong Tian

Research Scientist

Facebook AI Research

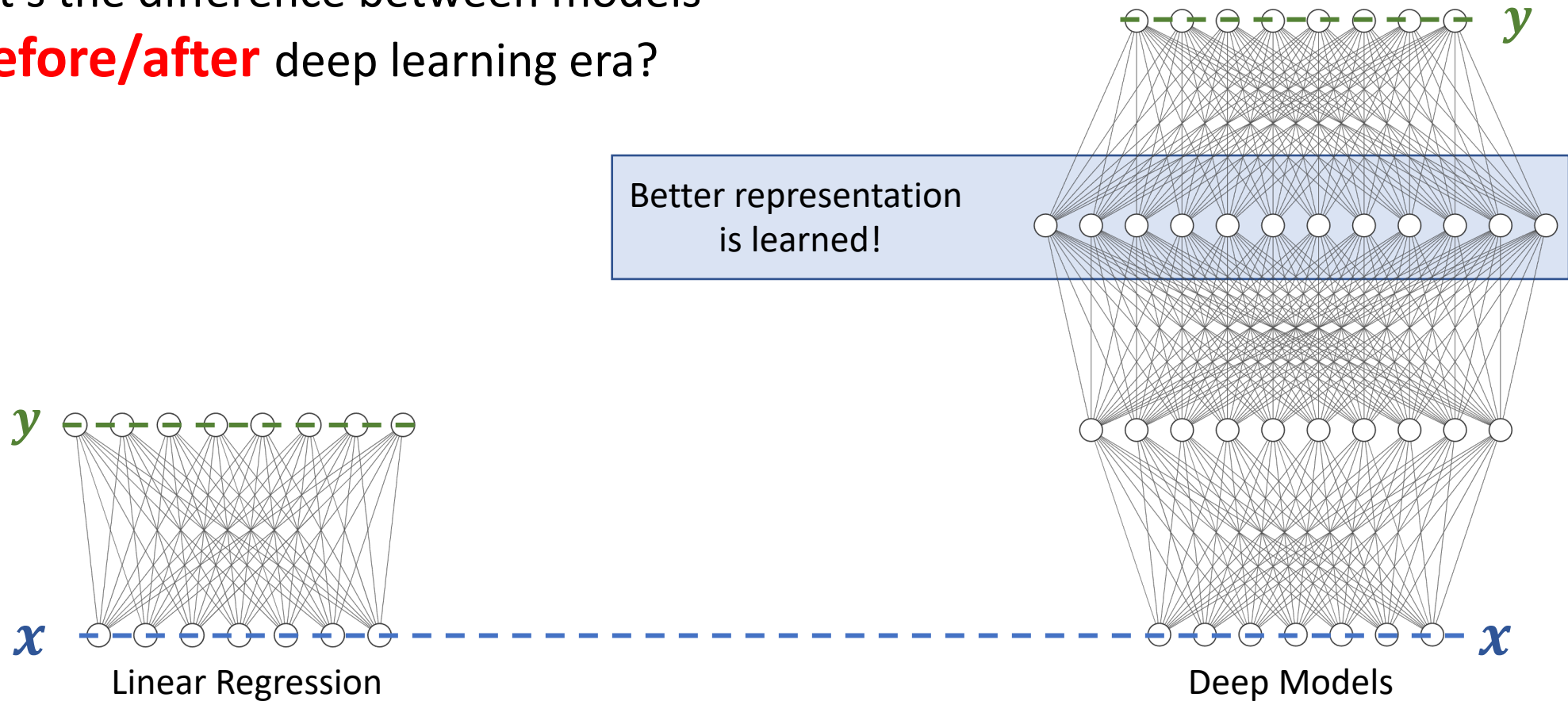


# Great Empirical Success of Deep Models



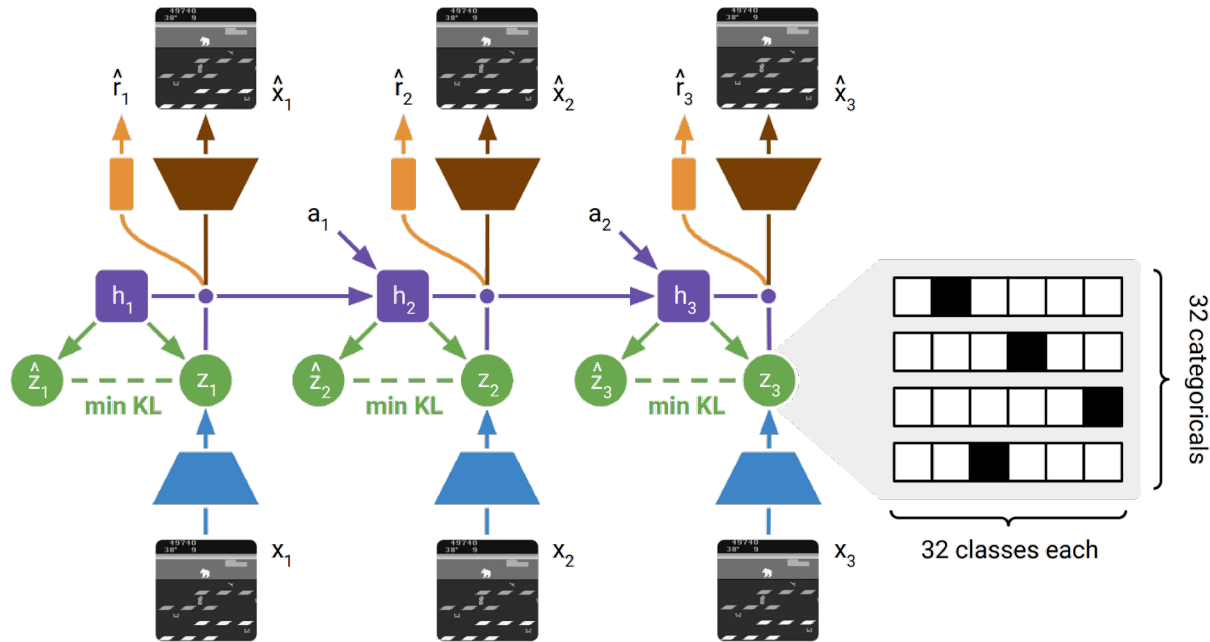
# Representation Learning

What's the difference between models  
**before/after** deep learning era?

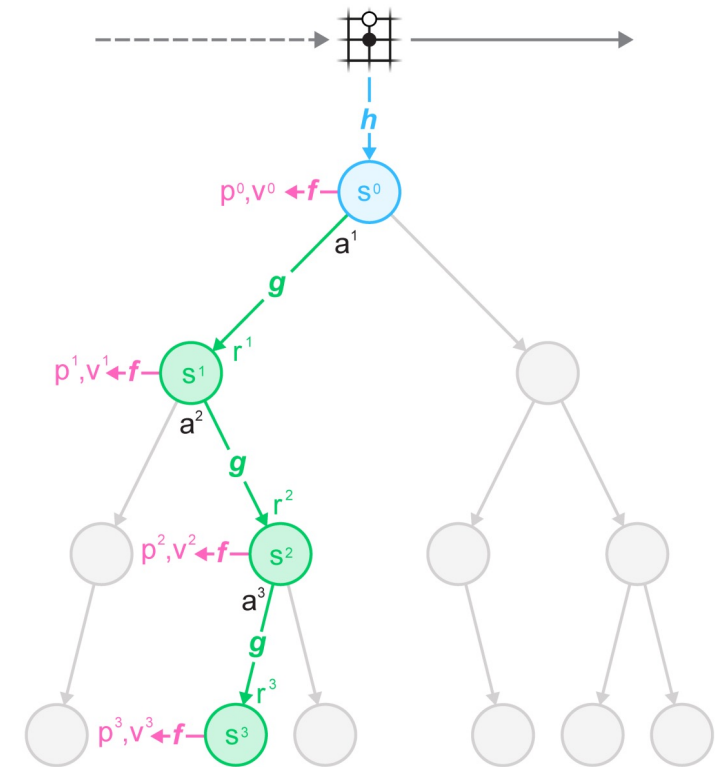


# Representation Learning in RL

- State Representation



[D. Hafner et al, Mastering Atari with Discrete World Models]

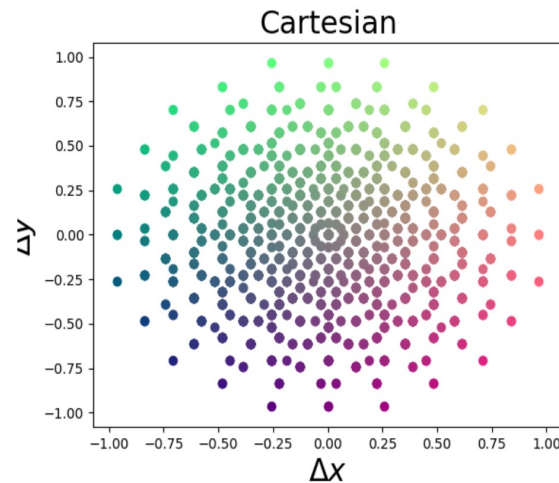
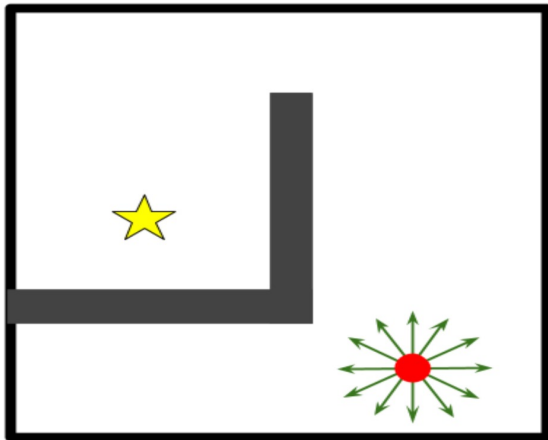


[J. Schrittwieser et al, Mastering Atari, Go, Chess and Shogi by Planning with a Learned Model, Nature]

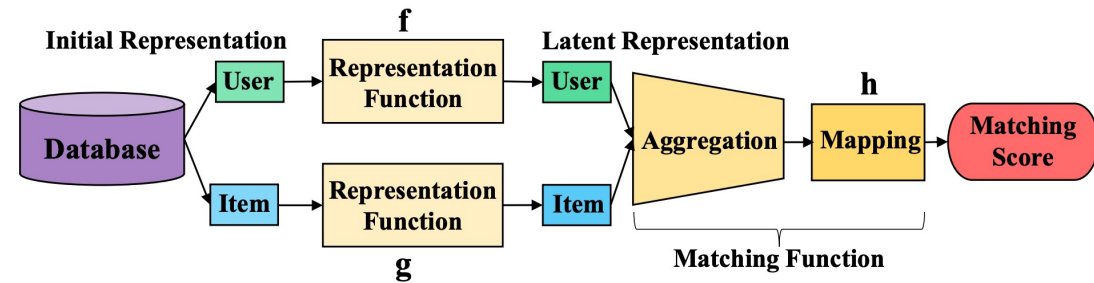
# Representation Learning in RL

- Action Representation

- Action embedding if we have a million of actions to choose from.
- Action embedding to transfer across different tasks.



[Y. Chandak, Learning Action Representations for Reinforcement Learning, ICML 2019]

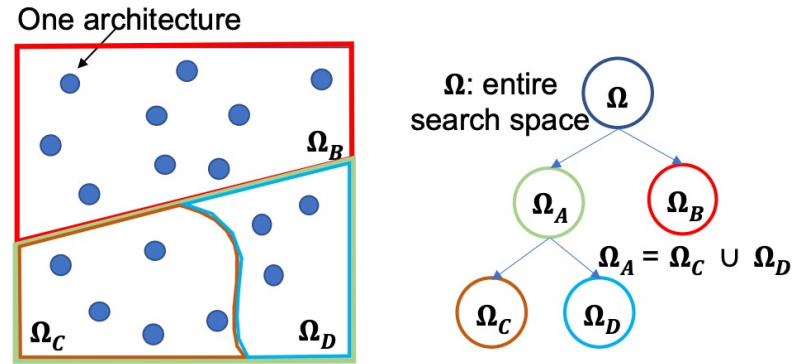


[Z. Deng et al, DeepCF: A Unified Framework of Representation Learning and Matching Function Learning in Recommender System, AAI 2019]

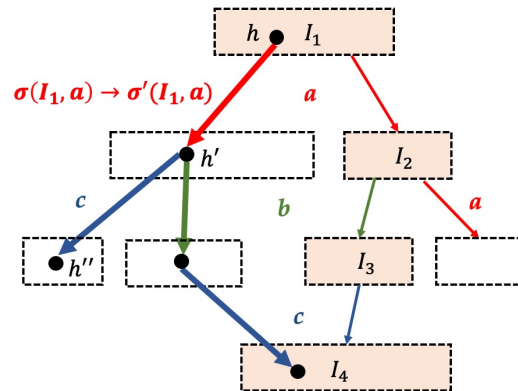
# How about high-level representation?

- Representation of the entire policy?
  - [A. Singh et al, Parrot: Data-Driven Behavioral Priors for Reinforcement Learning, ICLR 2021]
- Representation of the environment (multi-task learning)?

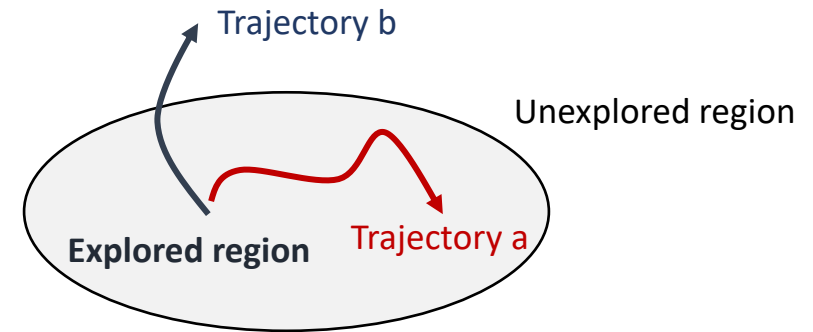
# This Talk



Representation of the Action Space



Representation for Easier Search

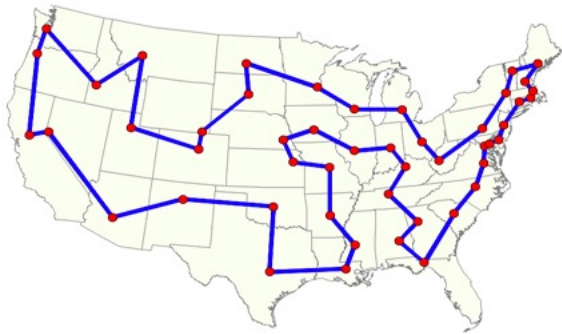


Representation for RL Exploration

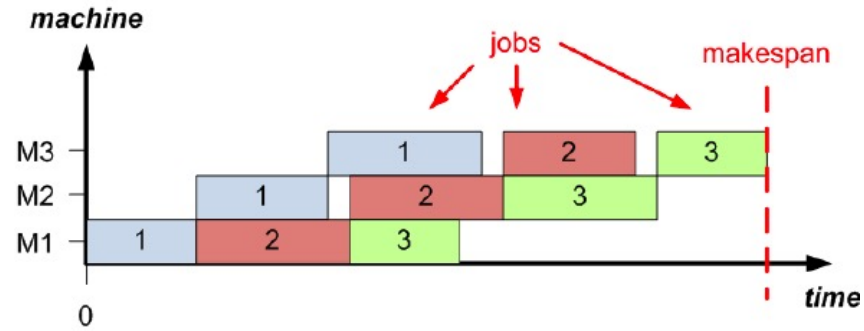
# Representation of the Action Space



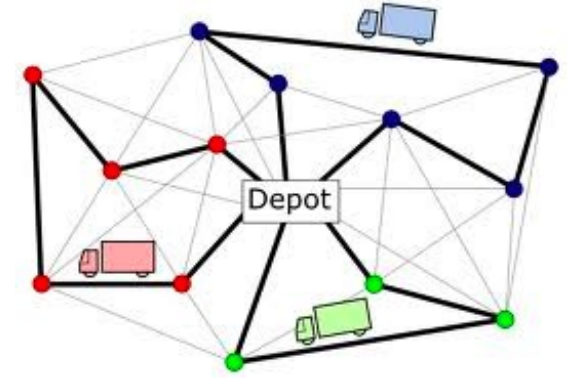
# Optimization Problems



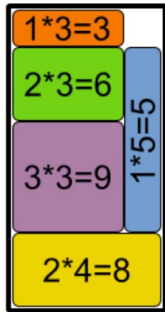
Travel Salesman Problem



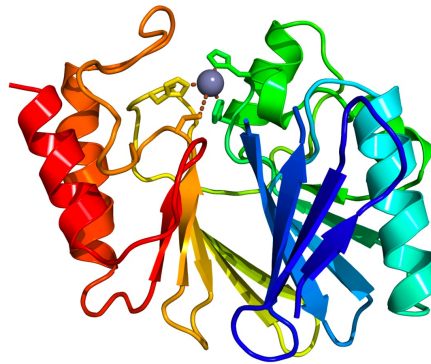
Job Scheduling



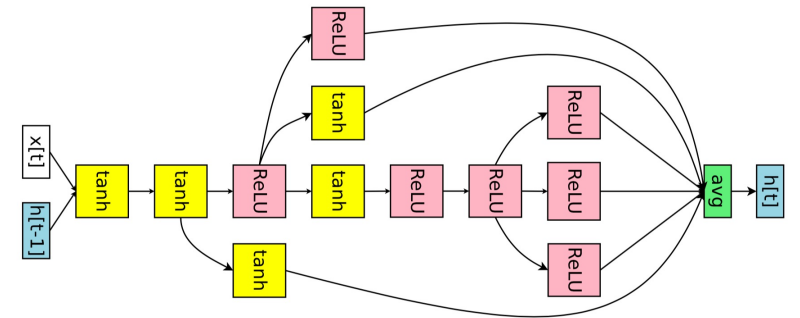
Vehicle Routing



Bin Packing



Protein Folding



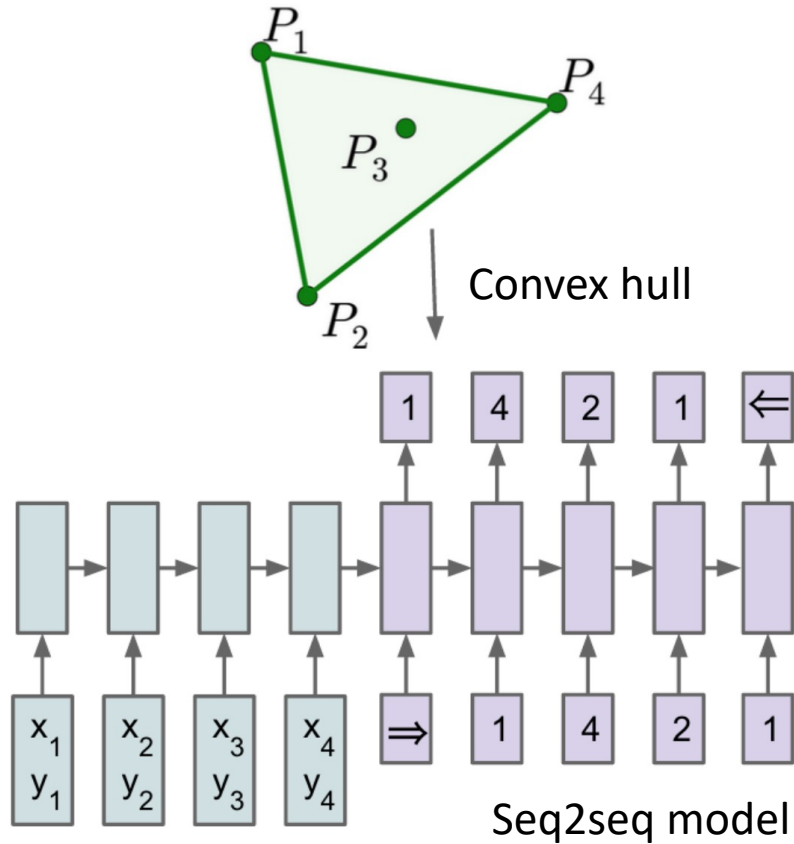
Model-Search

# There exists many MDPs for a single optimization!

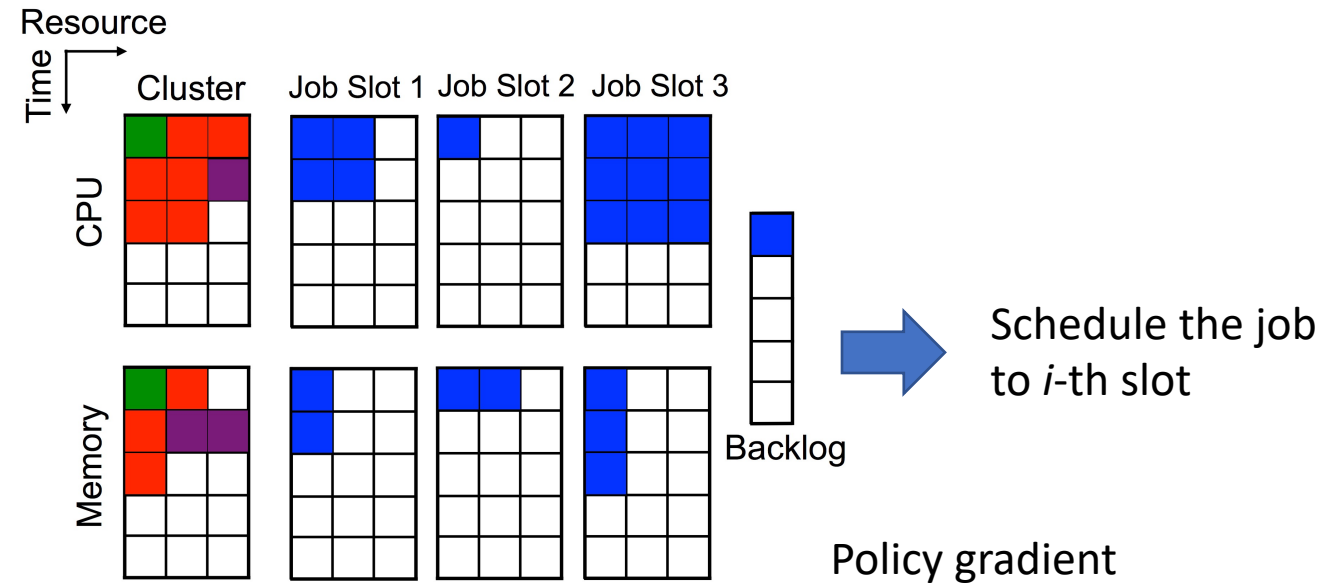
Name	Ways of Parameterization
One-shot Prediction	Spec $\rightarrow$ Solution
Progressive Prediction	Spec $\rightarrow$ SolPart1 $\rightarrow$ SolPart2 $\rightarrow$ SolPart3
Iterative Refinement	Spec $\rightarrow$ Sol1 $\rightarrow$ Sol2 (improved) $\rightarrow$ Sol3 (Better Improved)
Learned Action Space	Spec $\rightarrow$ All solution space $\rightarrow$ Small solution space $\rightarrow$ ...

**Representation Matters!**

# Direct predicting solutions

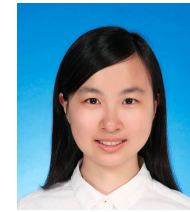


[O. Vinyals. et al, *Pointer Networks*, NIPS 2015]

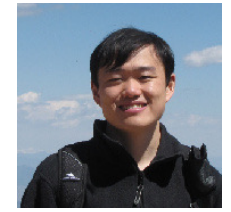


[H. Mao et al, *Resource Management with Deep Reinforcement Learning*, ACM Workshop on Hot Topics in Networks, 2016]

# Local Rewriting Framework

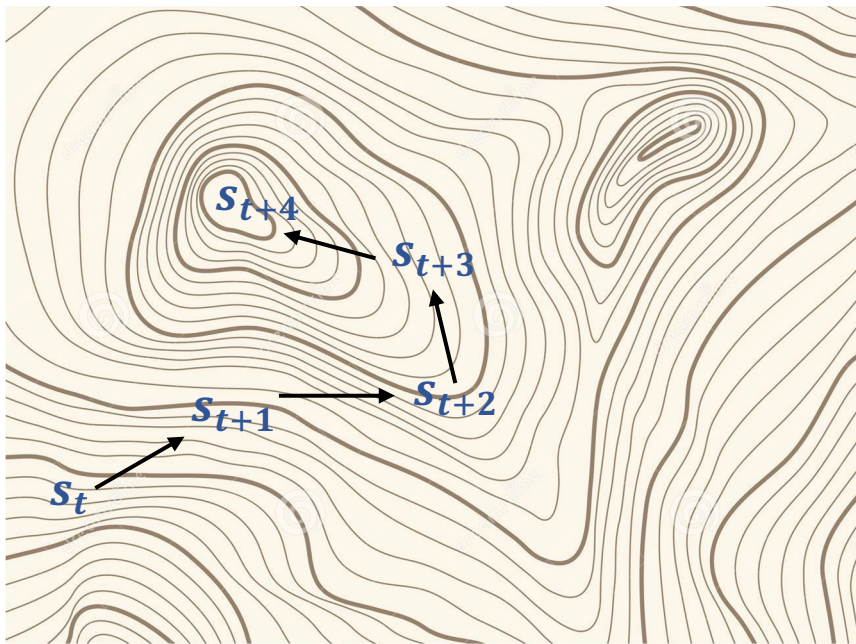


Xinyun Chen

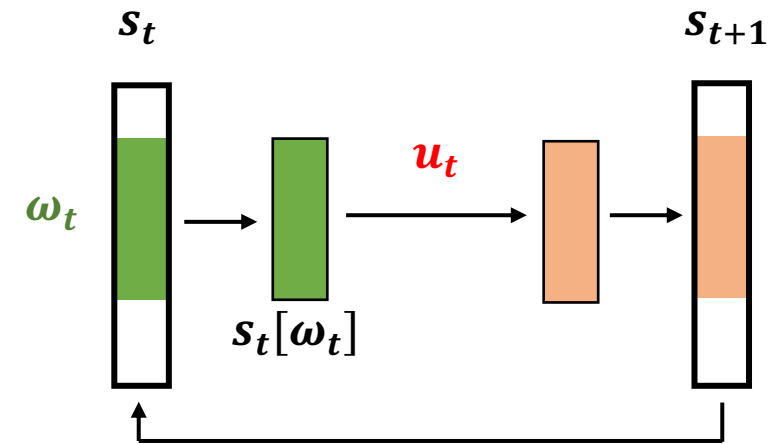


Yuandong Tian

[X. Chen and Y. Tian, *Learning to Perform Local Rewriting for Combinatorial Optimization*, NeurIPS 2019]



Start from a feasible solution and iteratively converges to a good solution



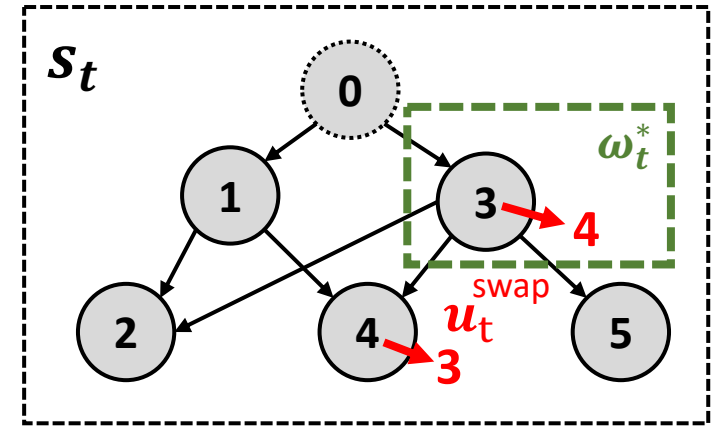
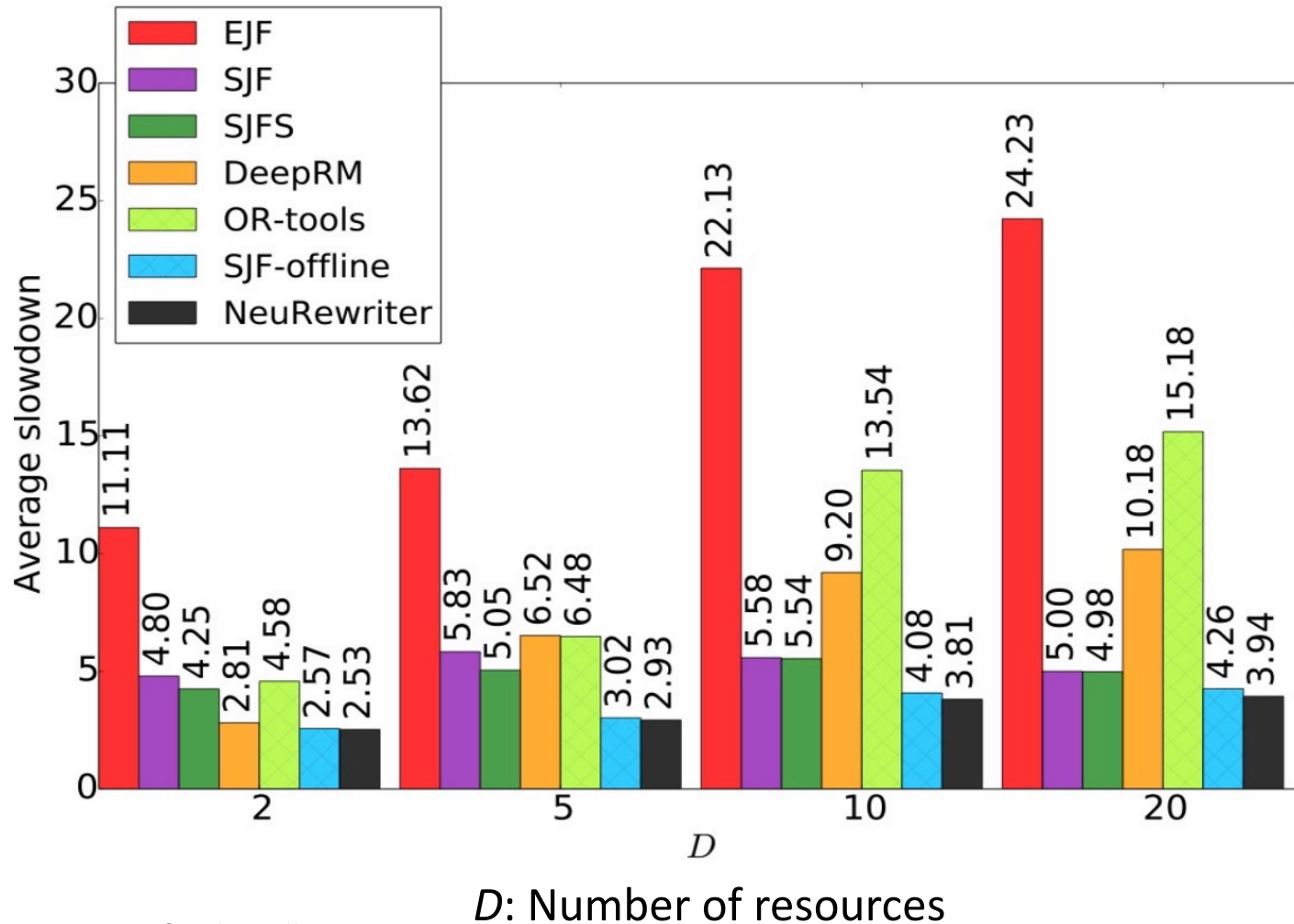
Current State  
(i.e. Solution)

Region-Picker

Rule-Picker

$$\begin{array}{ccccc} s_t & \longrightarrow & \omega_t \sim \pi_\omega(\cdot | s_t) & \longrightarrow & u_t \sim \pi_u(\cdot | s_t[\omega_t]) \\ \uparrow & & & & \uparrow \\ s_{t+1} = f(s_t, \omega_t, u_t) & & & & \end{array}$$

# Online Job Scheduling



## Baselines:

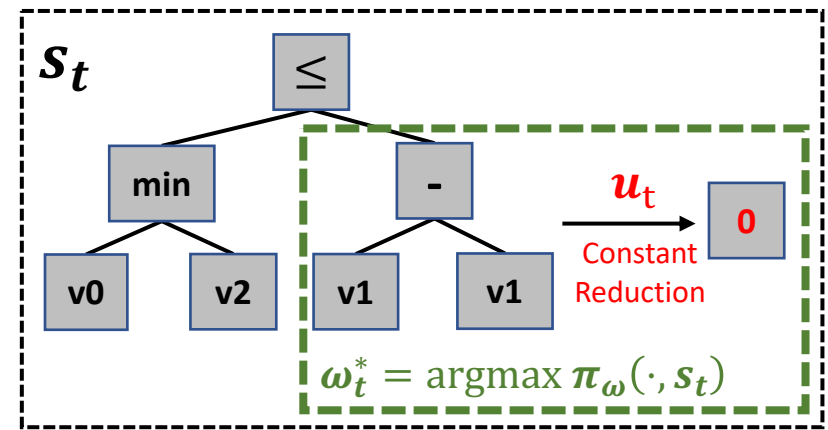
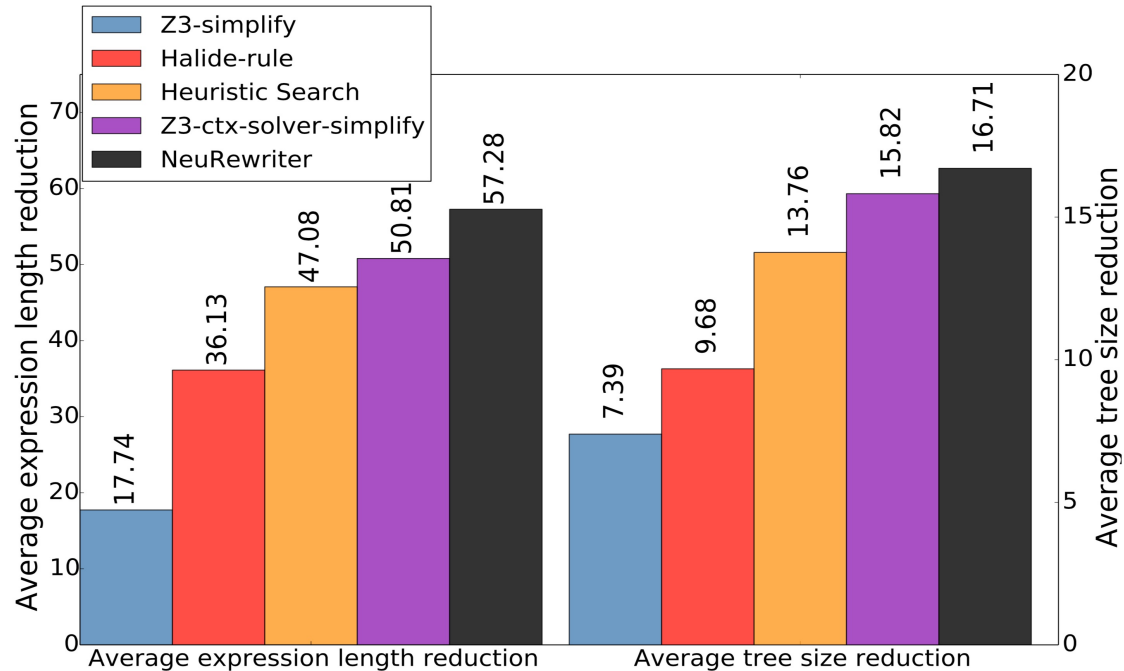
- Earliest Job First (EJF)
- Shortest Job First (SJF)
- Shortest First Search (SJFS)
- DeepRM

## Offline baselines:

- Google OR-tools (OR-tools)
- SJF-offline

	Time (s)
OR-tools	10.0
DeepRM	0.020
NeuRewriter	0.037

# Expression Simplification



Baselines:

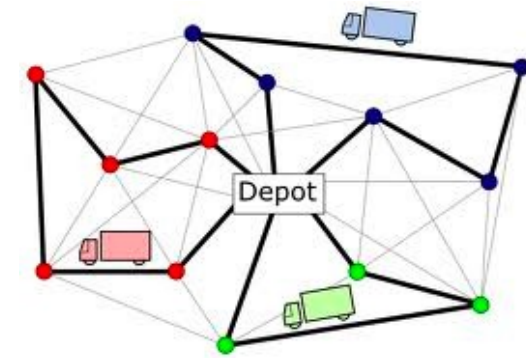
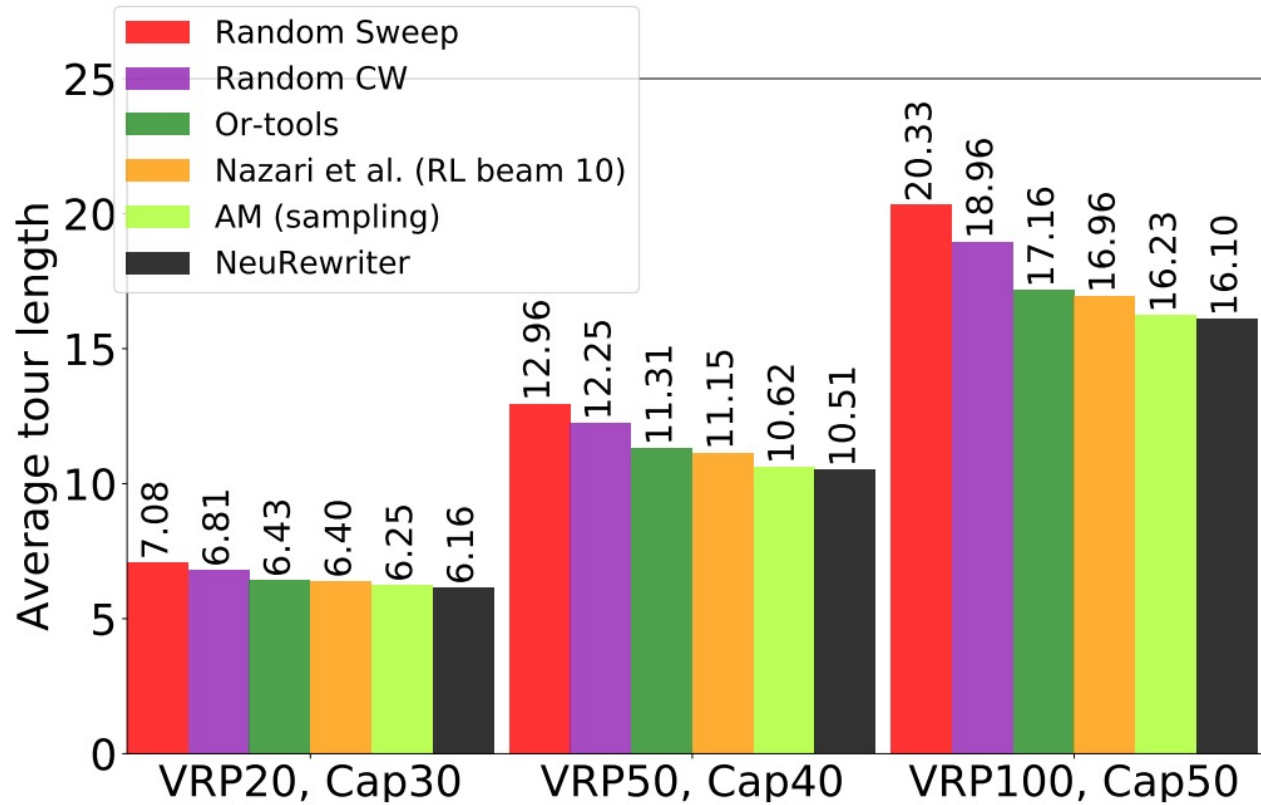
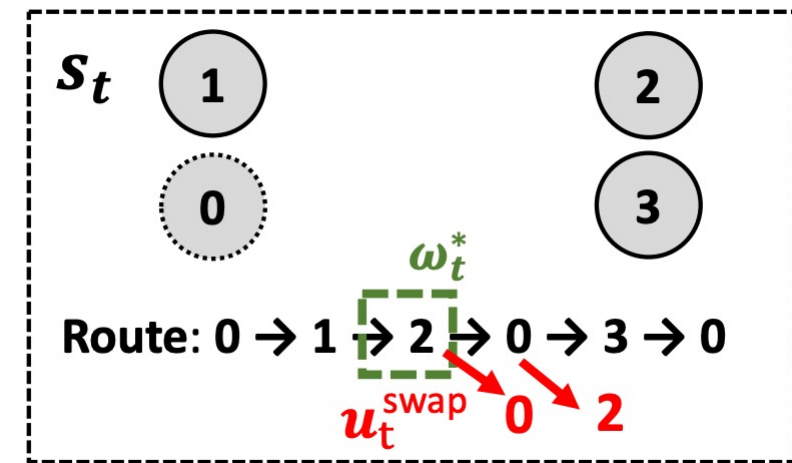
- Z3-simplify
- Z3-ctx-solver-simplify
- Heuristic Search
- Halide rules

	Time (s)
Z3-solver	1.375
NeuRewriter	0.159

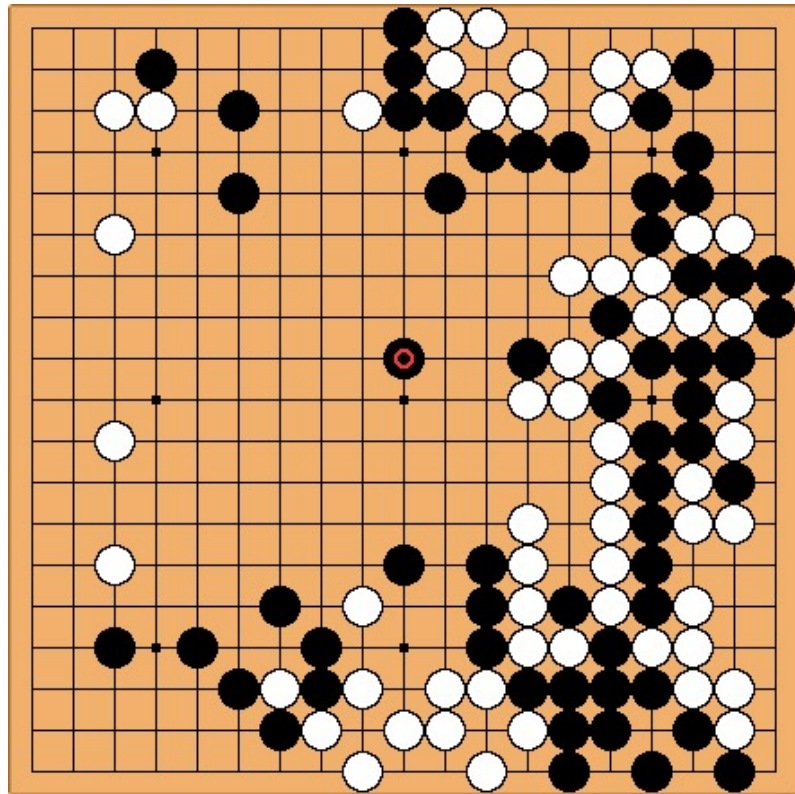
Follow-up work: Getting rid of manually specified rules

[H. Shi et al., Deep Symbolic Superoptimization without Human Knowledge, ICLR 2020]

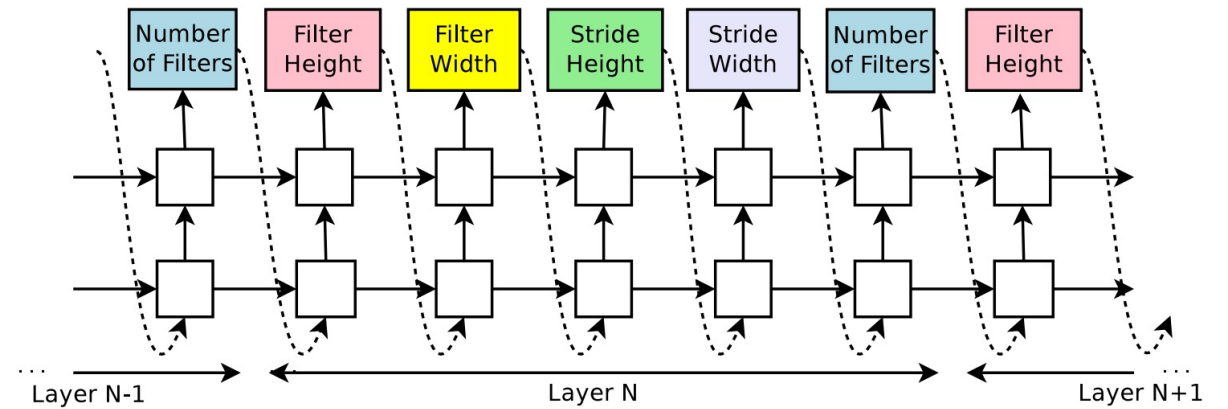
# Capacitated Vehicle Routing



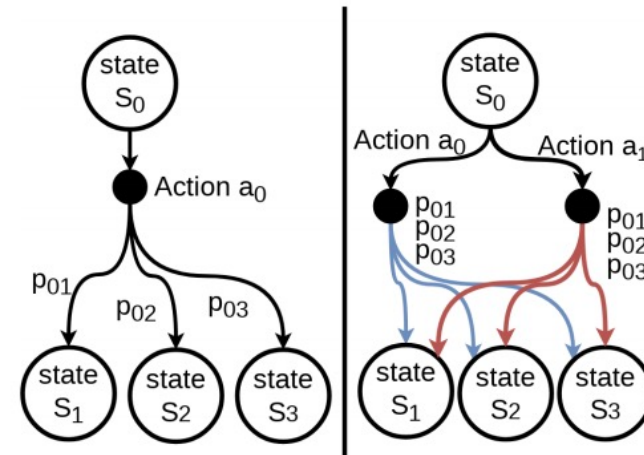
# Predefined Action Space



Fixed action space =  $R^{361}$



[B. Zoph and Q. Le, *Neural Architecture Search with Reinforcement Learning*, 2016]



[G. Malazgirt, *TauRieL: Targeting Traveling Salesman Problem with a deep reinforcement learning inspired architecture*]



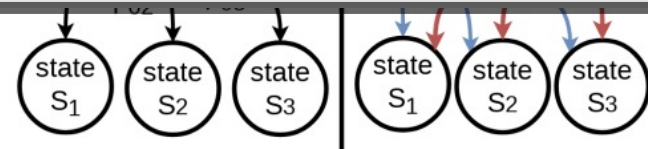
# Predefined Action Space

Number of Filters   Filter Height   Filter Width   Stride Height   Stride Width   Number of Filters   Filter Height

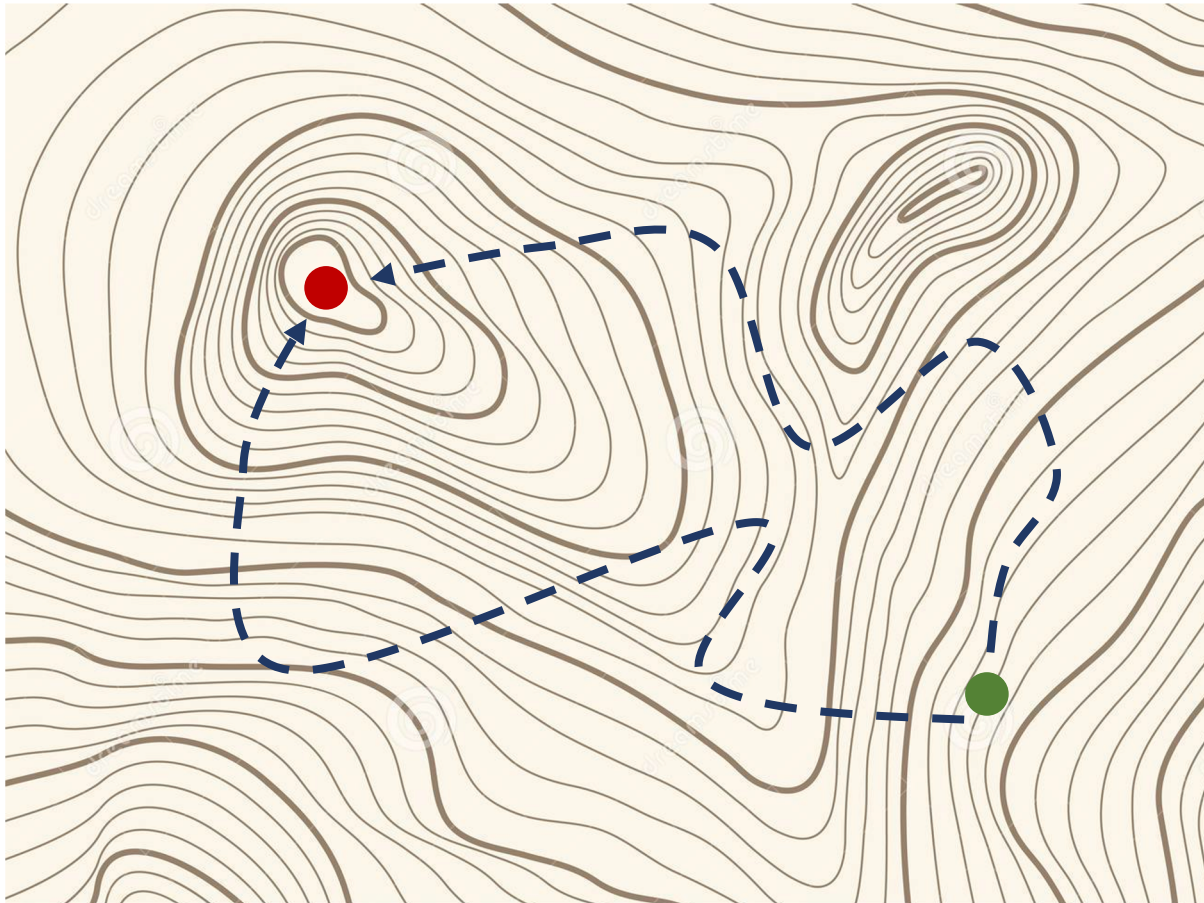
# Why Predefined Action Space?

...  
N+1  
2016]

Fixed action space =  $R^{361}$



# Why Predefined Action Space?



We only care the final solution

We don't care how we get it.

# Different Representation matters

Depth = {1, 2, 3, 4, 5}  
Channels = {32, 64}  
KernelSize = {3x3, 5x5}

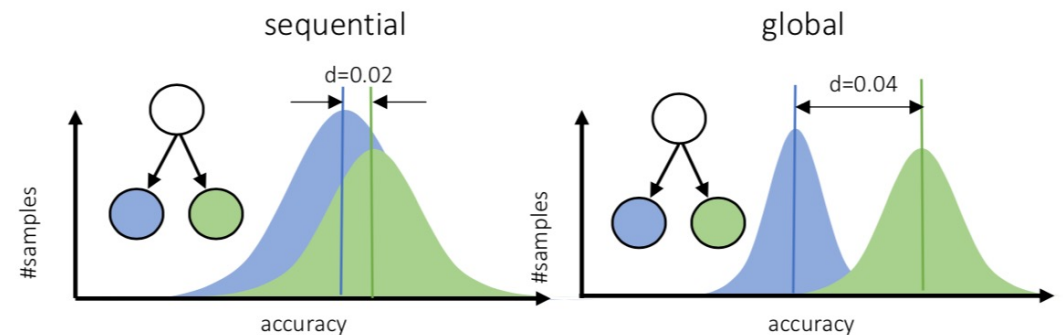
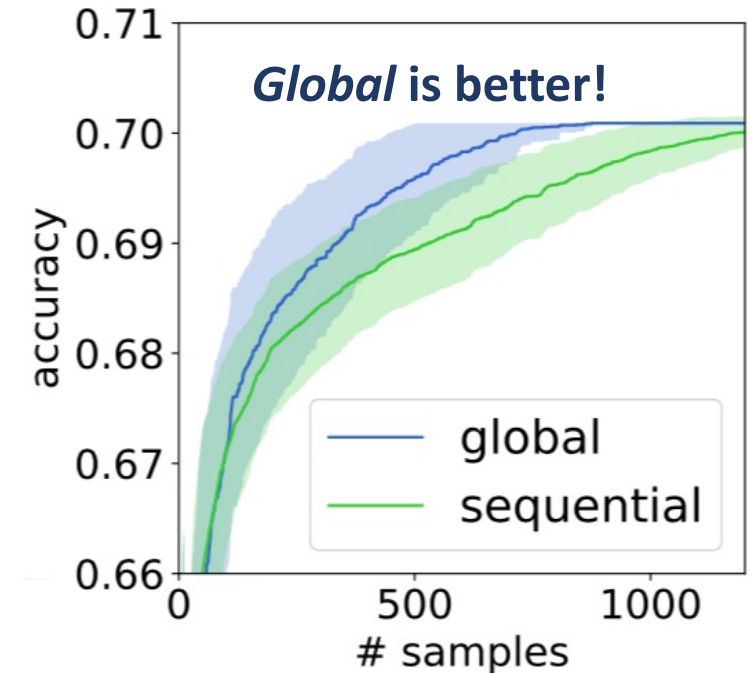
1364 networks.

**Goal:** Find the network  
with the best accuracy using fewest trials.

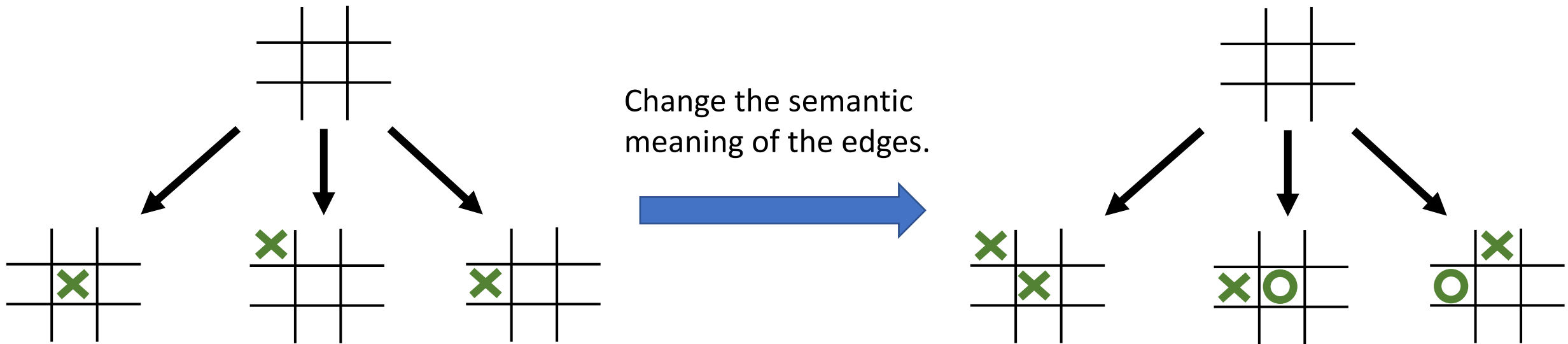
## Representation of action space

*Sequential* = { add a layer, set K, set C }

*Global* = { Set depth, set all K, set all C }



# The Meaning of Learning Action Space



**Not allowed in games, but doable in optimization.**

# Learning Action Space



Linnan Wang



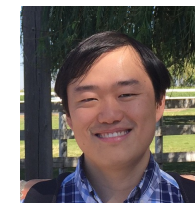
Saining Xie



Teng Li



Rodrigo Fonseca

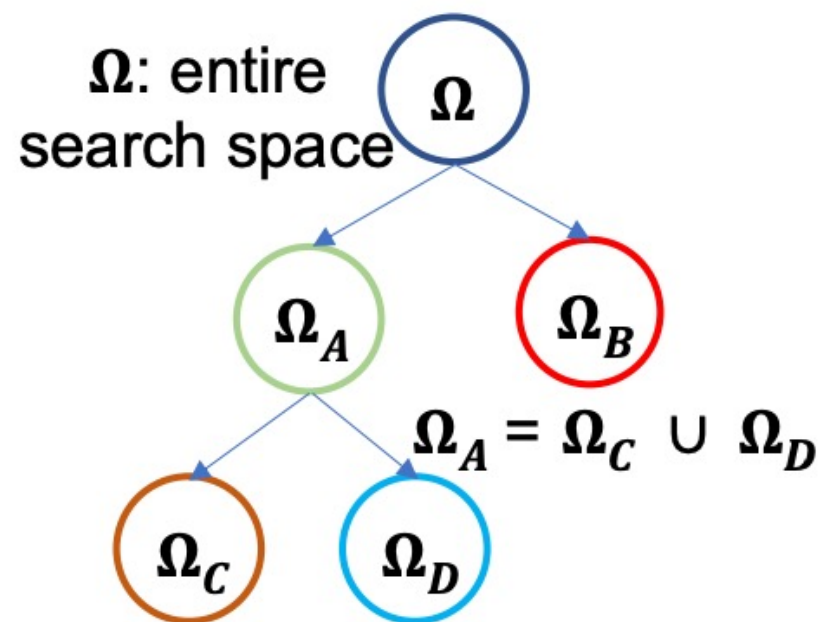
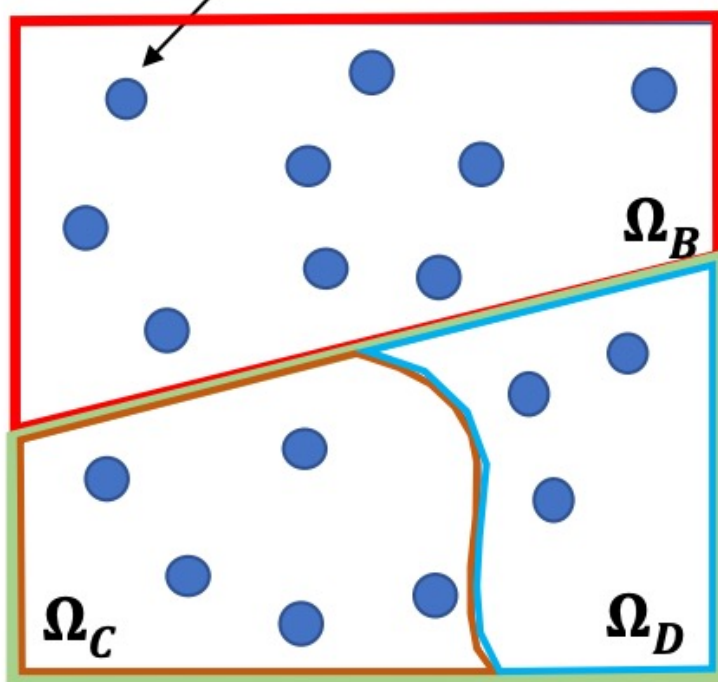


Yuandong Tian

[L. Wang, R. Fonseca, Y. Tian, **Learning Search Space Partition for Black-box Optimization using Monte Carlo Tree Search**, *NeurIPS 2020*]

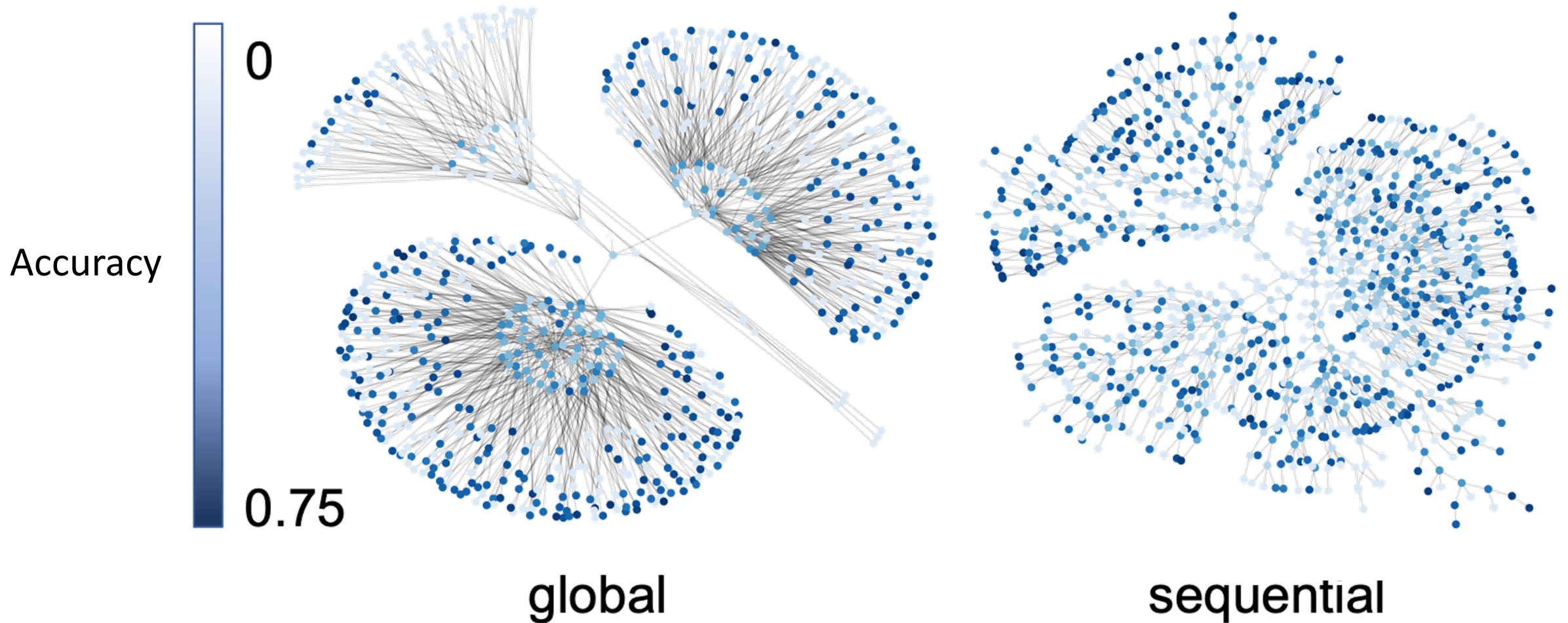
[L. Wang, S. Xie, T. Li, R. Fonseca, Y. Tian, **Sample-Efficient Neural Architecture Search by Learning Action Space**, *TPAMI 2021*]

One architecture

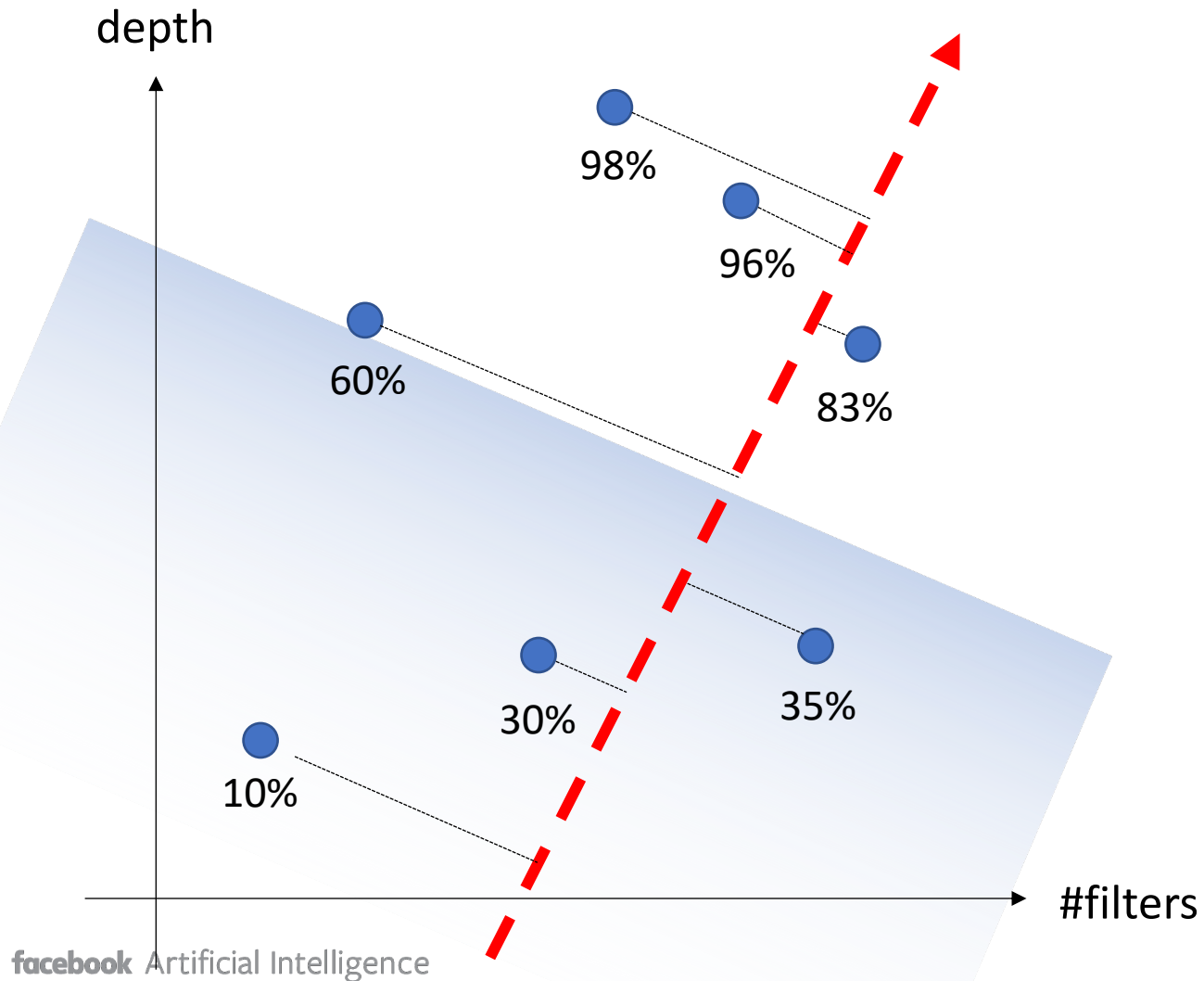


Partition = Action

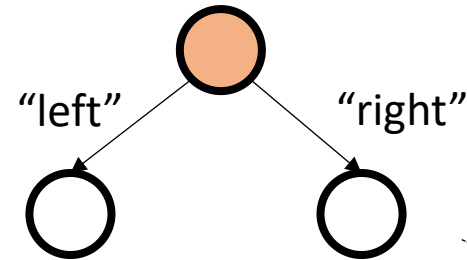
# Different Partition $\rightarrow$ Different Value Distribution



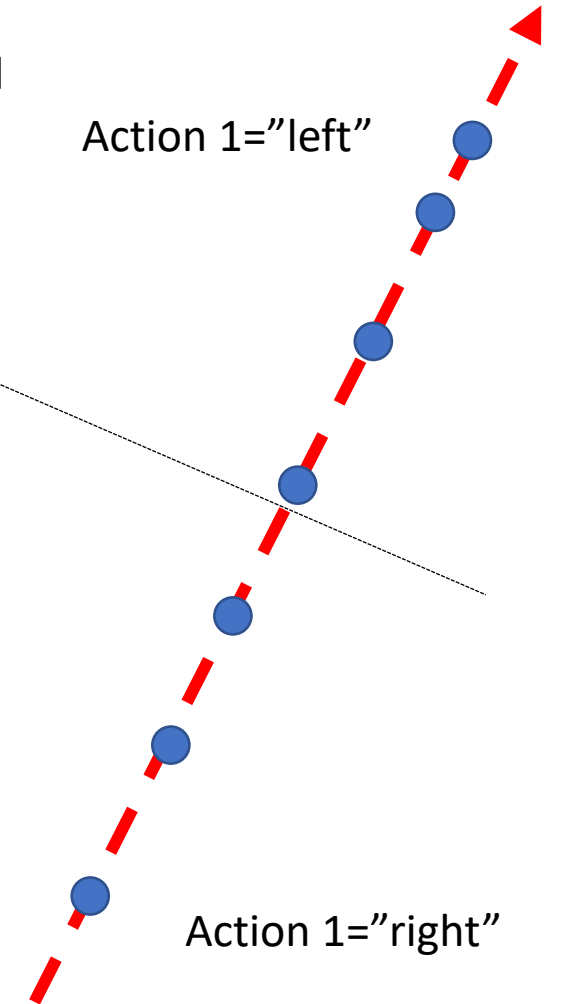
# Learn action space



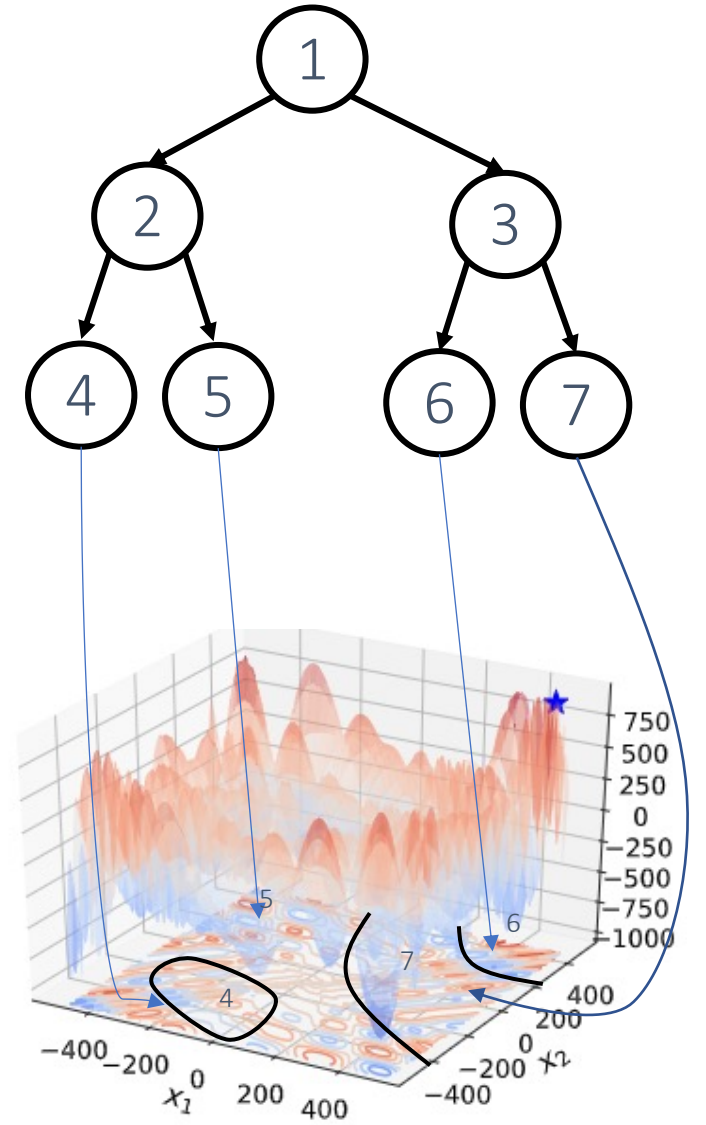
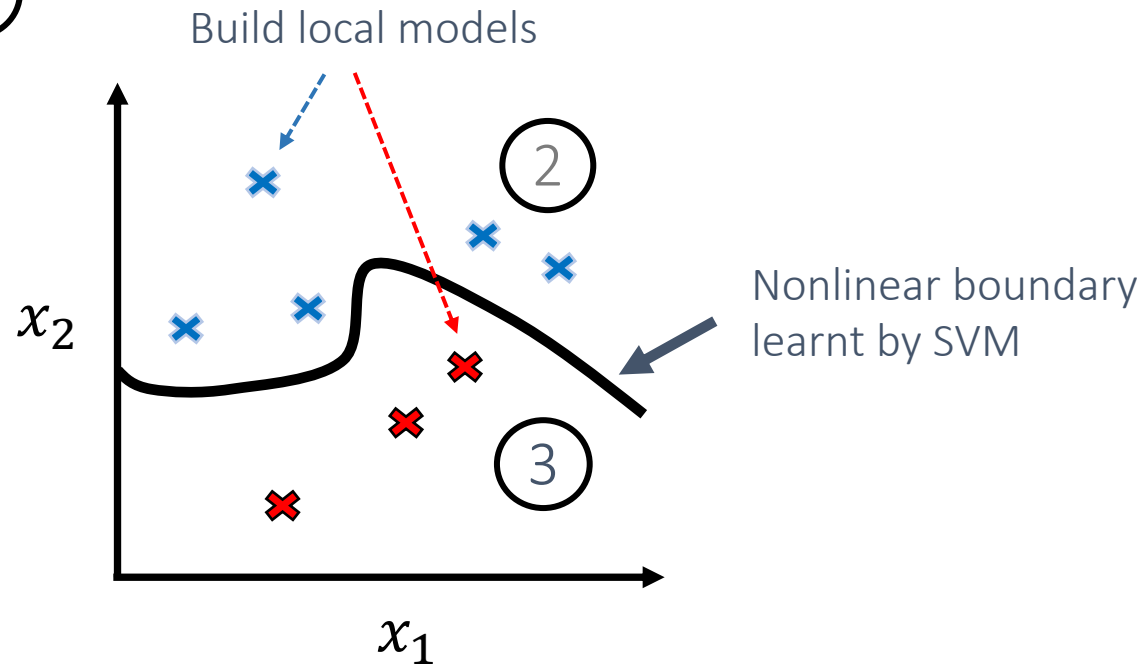
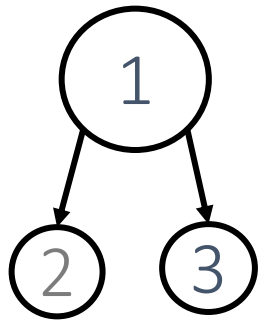
Current node whose action space is learned



Action 1="left"

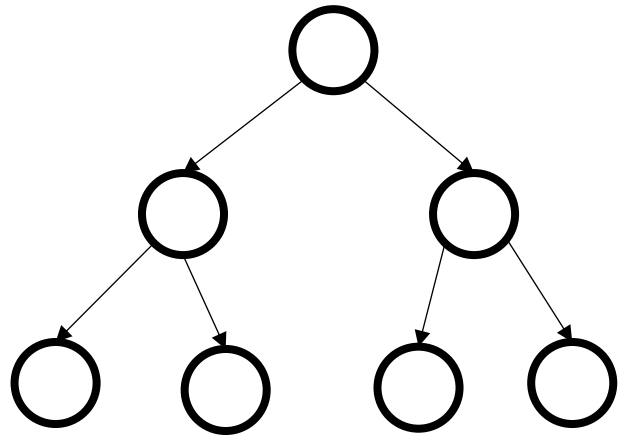


# Nonlinear Partition





# Approach



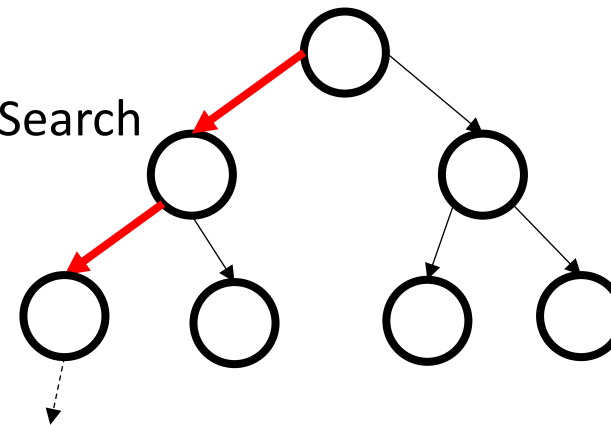
Fixed action branches  
(but not action space)

(a) Train the action space.

	Accuracy
(filter=2, depth=5)	85%
(filter=3, depth=7)	92%

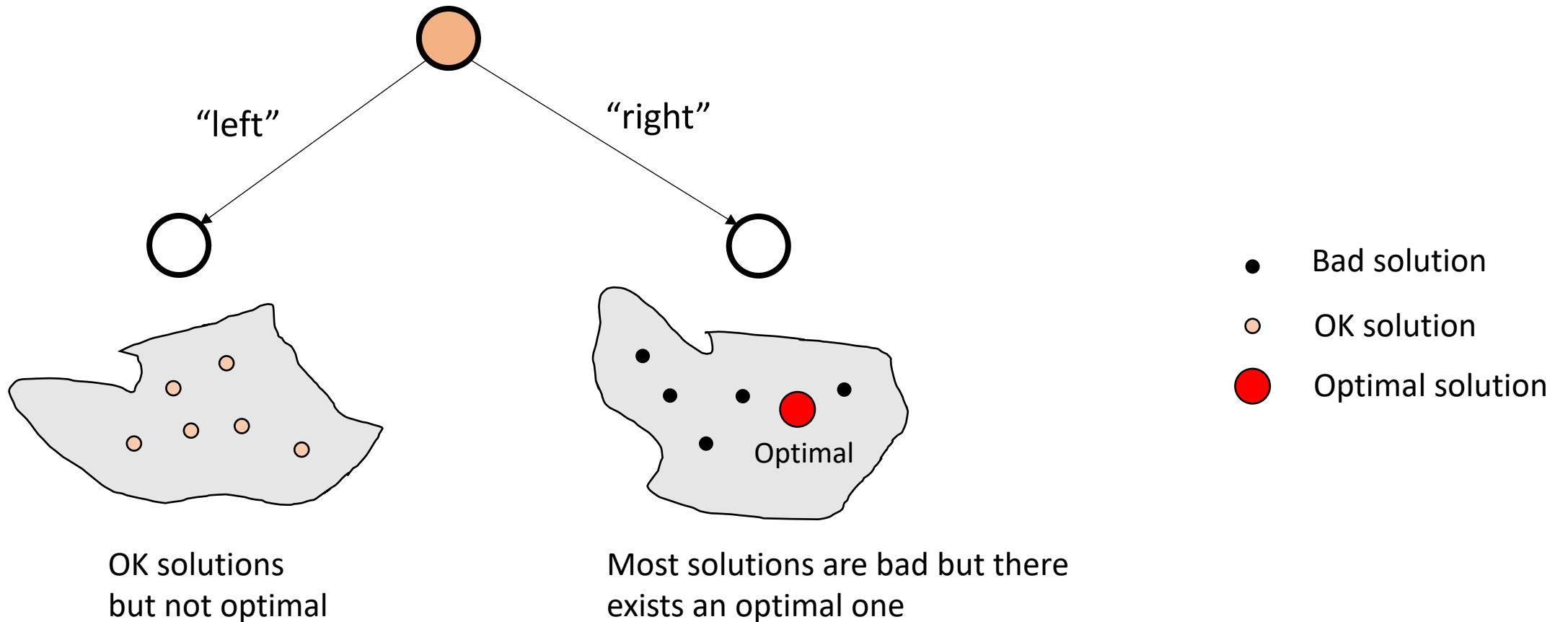
(b) Search using learned action space until a fixed #rollouts are used.

Monte Carlo Tree Search  
(MCTS)



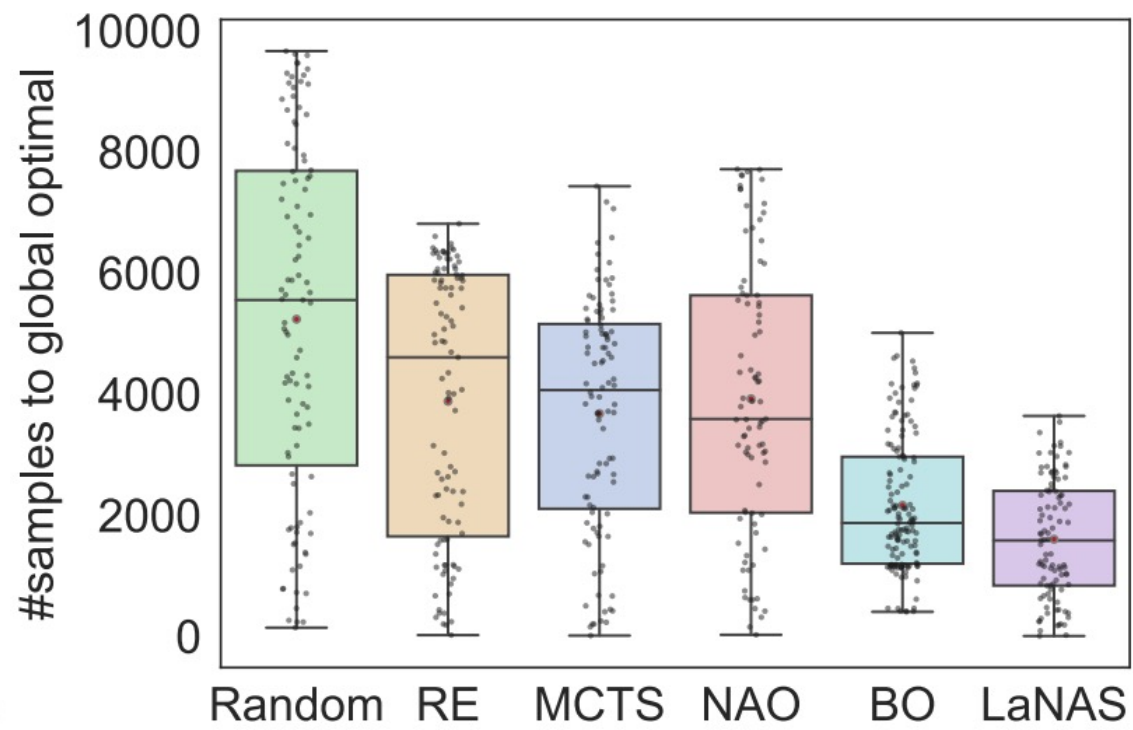
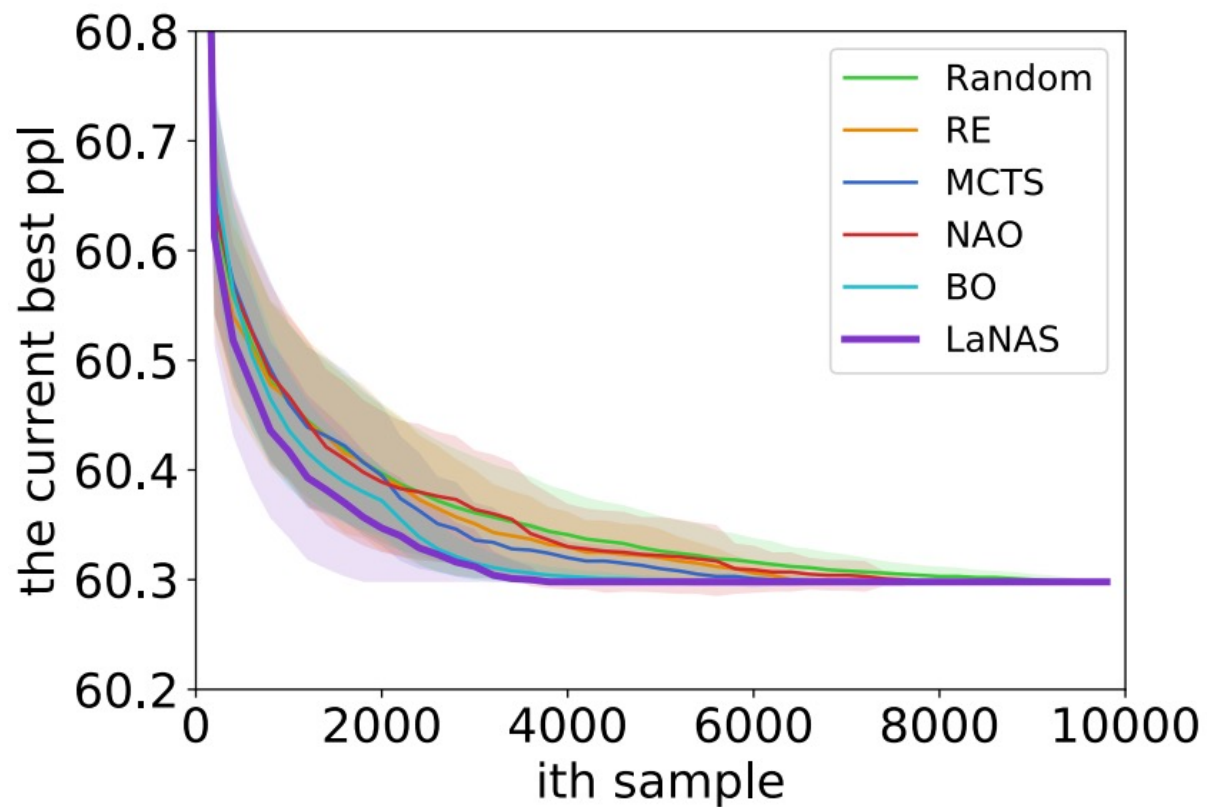
Getting the true quality  $f(x)$  for the solution  $x$

# Why Exploration is Important



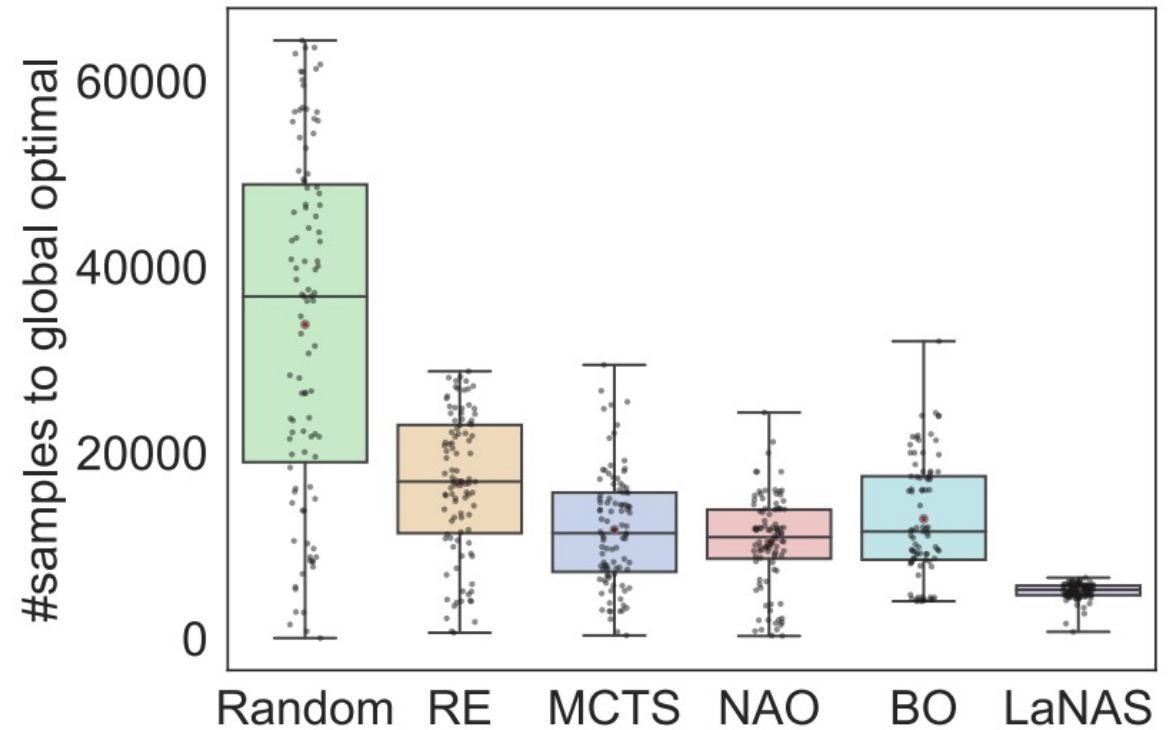
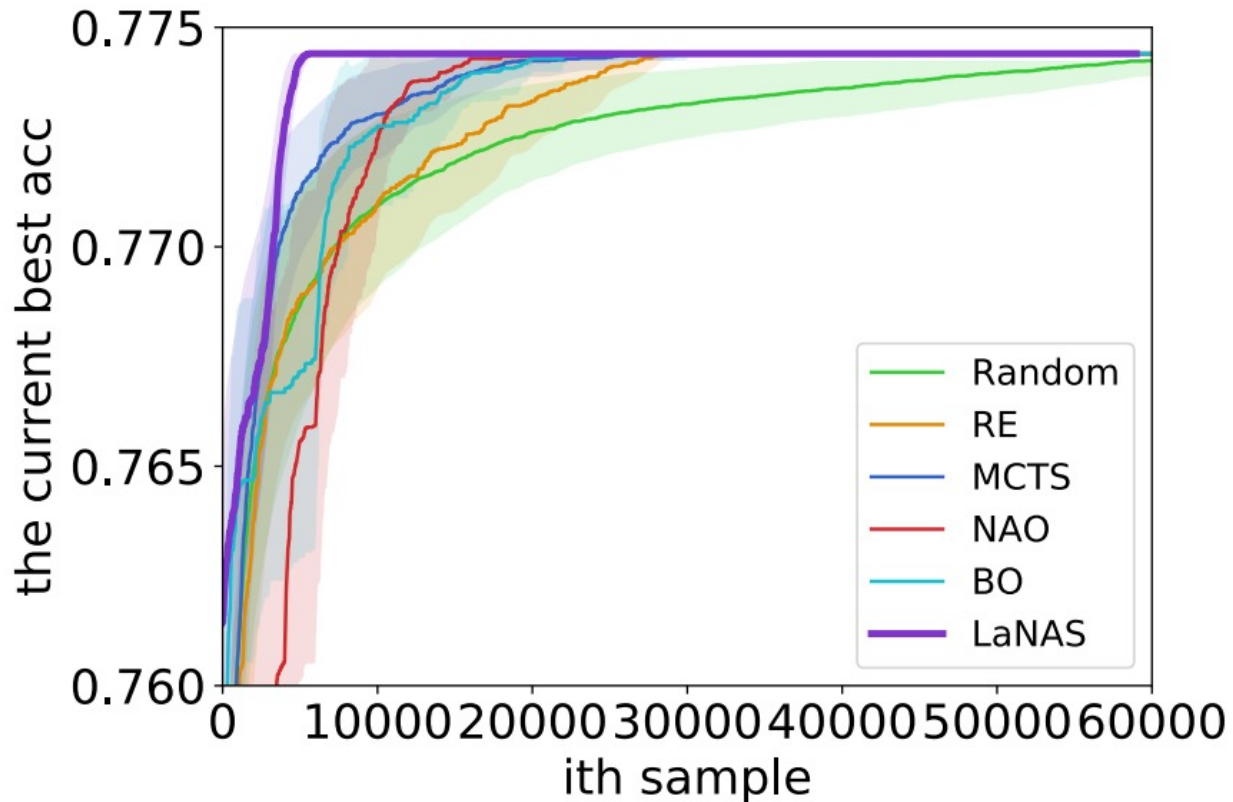
# Performance

Customized dataset: LSTM-10K (PTB)



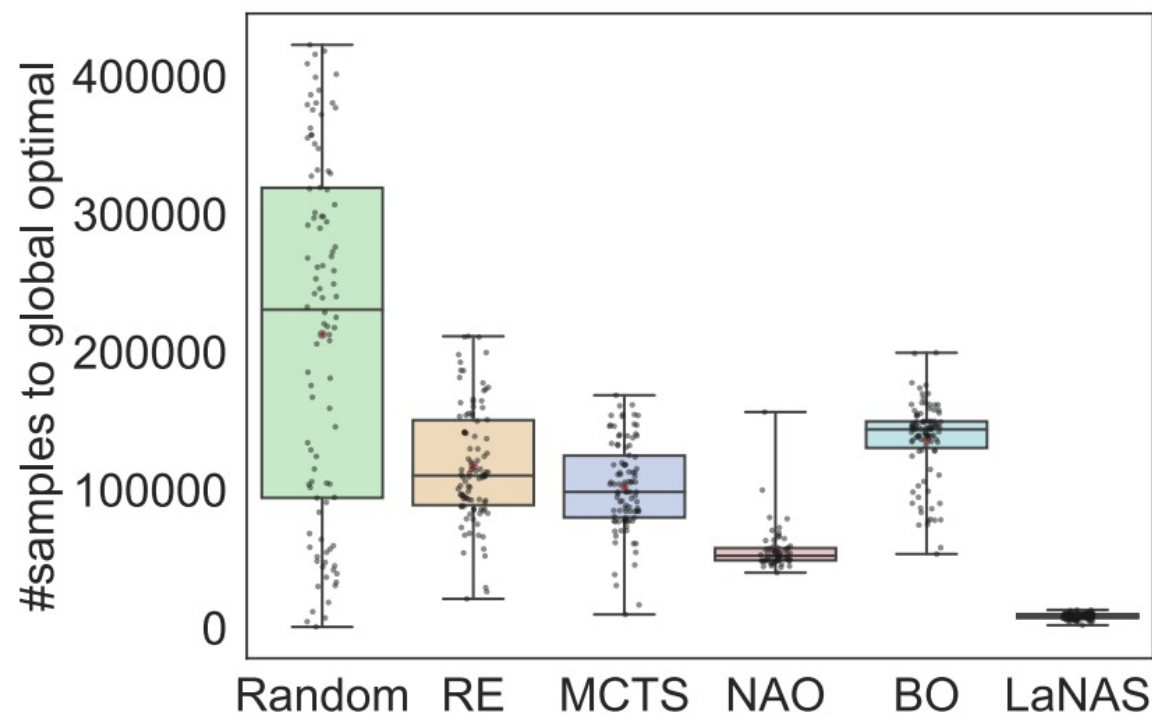
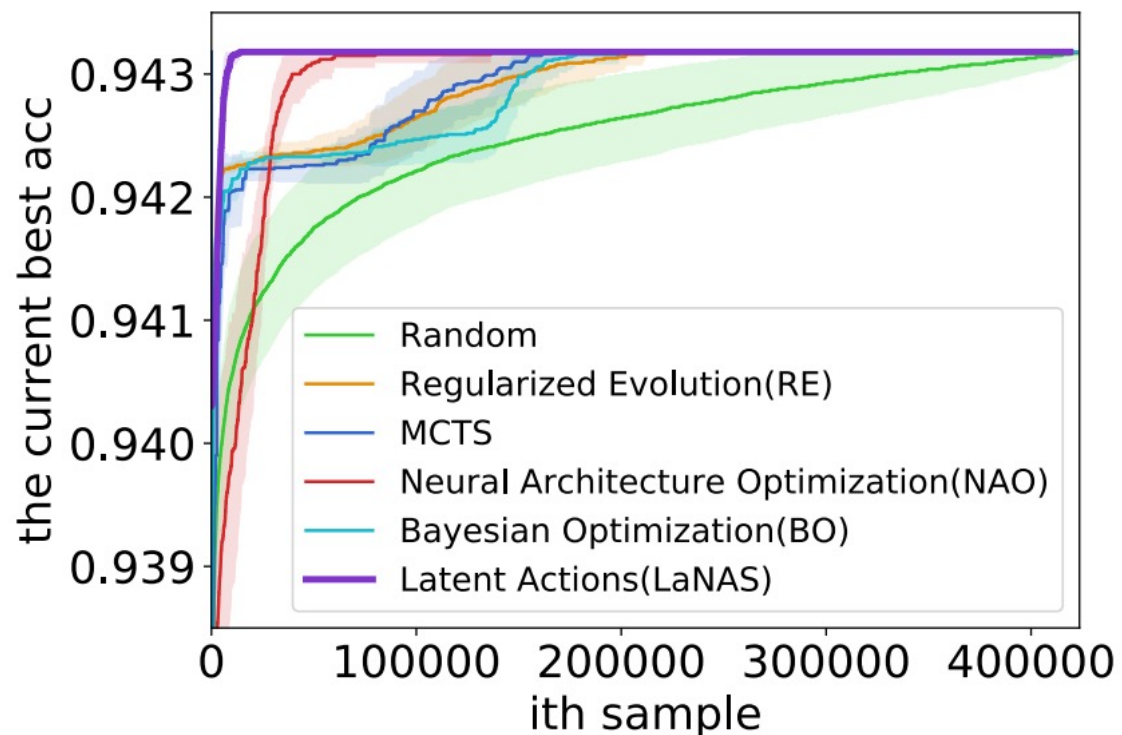
# Performance

Customized dataset: ConvNet-60K (CIFAR-10, VGG style models)



# Performance

NASBench-101 (CIFAR-10, 420k models, NASNet Search Space)



Each curve is repeated 100 times. We randomly pick 2k models to initialize.

# Open Domain

CIFAR-10  
(NASNet style  
architecture)

Model	Using ImageNet	Params	Top1 err	M	GPU days
search based methods					
NASNet-A+c/o [22]	X	3.3 M	2.65	20000	2000
AmoebaNet-B+c/o [10]	X	2.8 M	2.55 $\pm$ 0.05	27000	3150
PNASNet-5 [29]	X	3.2 M	3.41 $\pm$ 0.09	1160	225
NAO+c/o [30]	X	128.0 M	2.11	1000	200
AmoebaNet-B+c/o	X	34.9 M	2.13 $\pm$ 0.04	27000	3150
EfficientNet-B7	✓	64M	1.01		
BiT-M	✓	60M	1.09		
<b>LaNet+c/o</b>	X	3.2 M	<b>1.63</b> $\pm$ 0.05	800	150
<b>LaNet+c/o</b>	X	44.1 M	<b>0.99</b> $\pm$ 0.02	800	150
one-shot NAS based methods					
ENAS+c/o [18]	X	4.6 M	2.89	-	0.45
DARTS+c/o [20]	X	3.3 M	2.76 $\pm$ 0.09	-	1.5
BayesNAS+c/o [31]	X	3.4 M	2.81 $\pm$ 0.04	-	0.2
ASNG-NAS+c/o [32]	X	3.9 M	2.83 $\pm$ 0.14	-	0.11
XNAS+c/o [33]	X	3.7 M	1.81		0.3
<b>oneshot-LaNet+c/o</b>	X	3.6 M	<b>1.68</b> $\pm$ 0.06	-	3
<b>oneshot-LaNet+c/o</b>	X	45.3 M	<b>1.2</b> $\pm$ 0.03	-	3

M: number of samples selected.

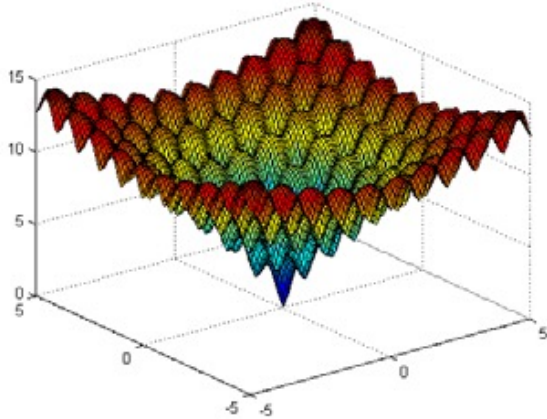
# Open Domain

ImageNet  
(mobile setting  
Flop < 600M)

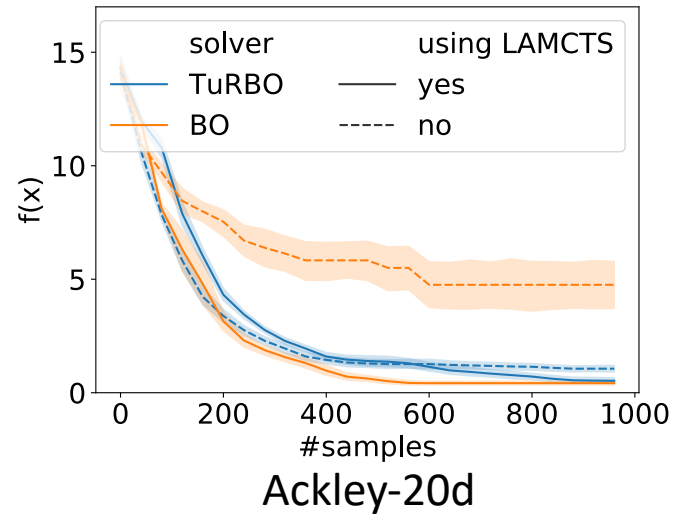
<b>Model</b>	<b>FLOPs</b>	<b>Params</b>	<b>top1 / top5 err</b>
NASNet-A (Zoph et al. (2018))	564M	5.3 M	26.0 / 8.4
NASNet-B (Zoph et al. (2018))	488M	5.3 M	27.2 / 8.7
NASNet-C (Zoph et al. (2018))	558M	4.9 M	27.5 / 9.0
AmoebaNet-A (Real et al. (2018))	555M	5.1 M	25.5 / 8.0
AmoebaNet-B (Real et al. (2018))	555M	5.3 M	26.0 / 8.5
AmoebaNet-C (Real et al. (2018))	570M	6.4 M	<b>24.3 / 7.6</b>
PNASNet-5 (Liu et al. (2018a))	588M	5.1 M	25.8 / 8.1
DARTS (Liu et al. (2018b))	574M	4.7 M	26.7 / 8.7
FBNet-C (Wu et al. (2018))	375M	5.5 M	25.1 / -
RandWire-WS (Xie et al. (2019))	583M	5.6 M	25.3 / 7.8
BayesNAS (Zhou et al. (2019))	-	3.9 M	26.5 / 8.9
LaNet	570M	5.1 M	<b>25.0 / 7.7</b>

# La-MCTS as a meta method

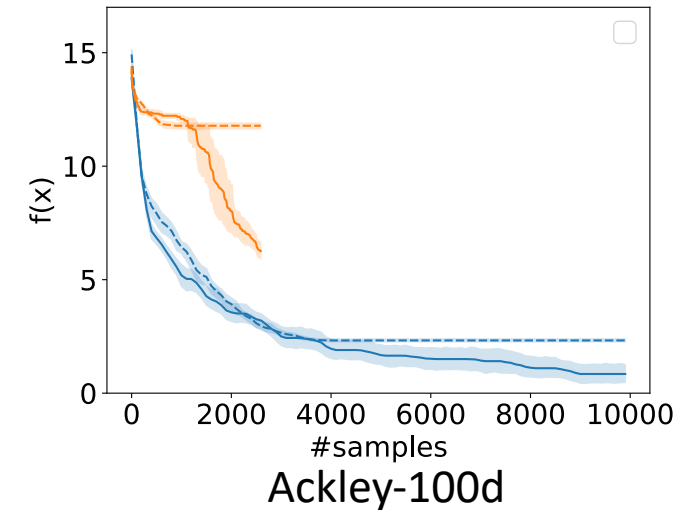
$$x^* = \arg \min_{x \in \Omega} f(x)$$



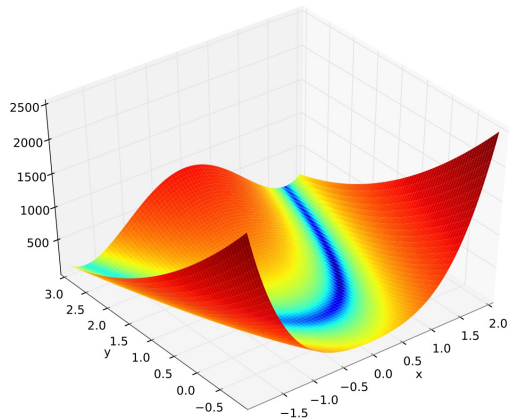
Ackley-2d



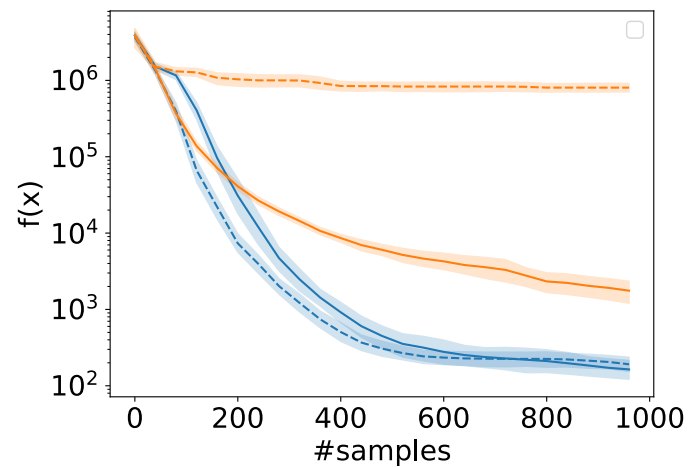
Ackley-20d



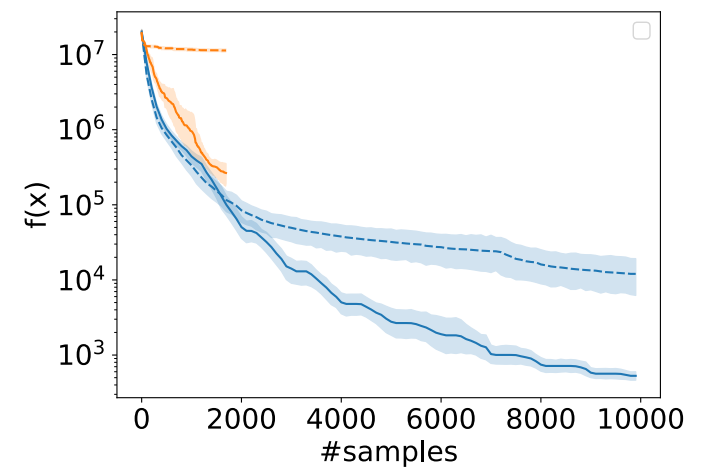
Ackley-100d



Rosenbrock-2d



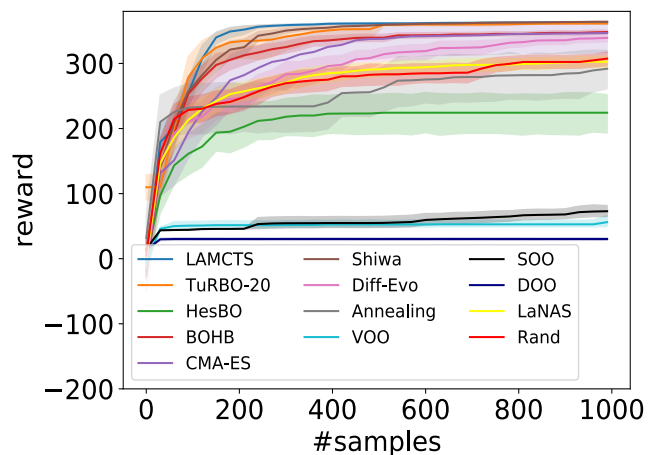
Rosenbrock-20d



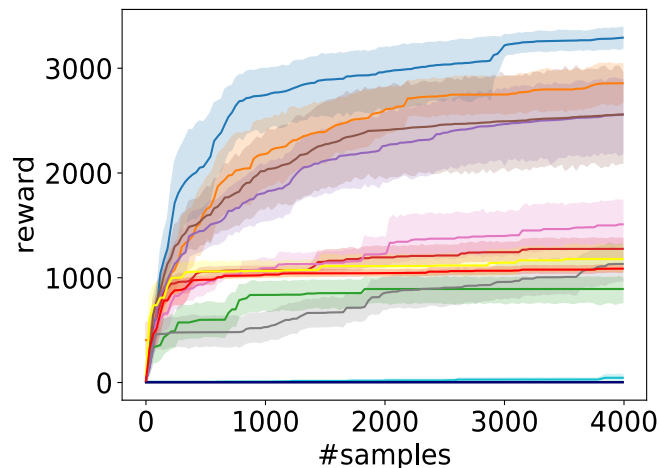
Rosenbrock-100d



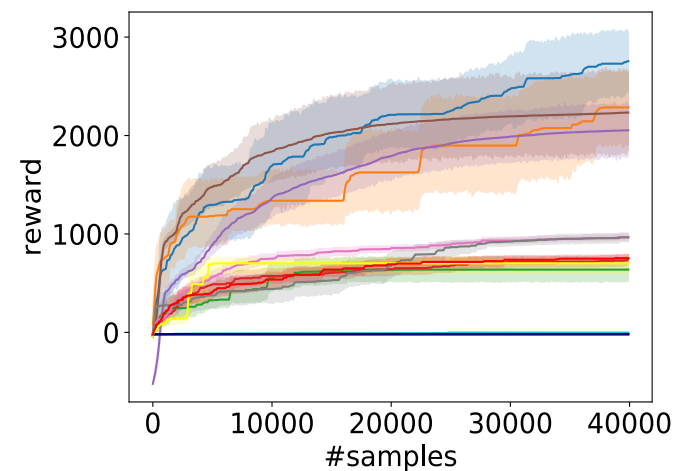
# Optimizing linear policy for Mujoco tasks



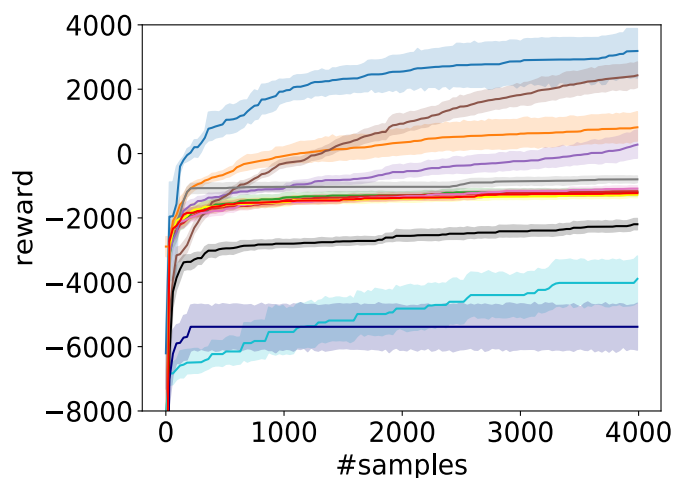
(a) Swimmer, #params = 16



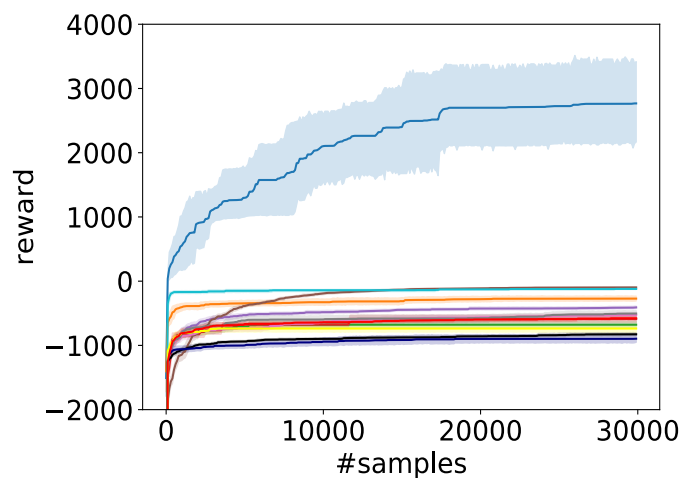
(b) Hopper, #params = 33



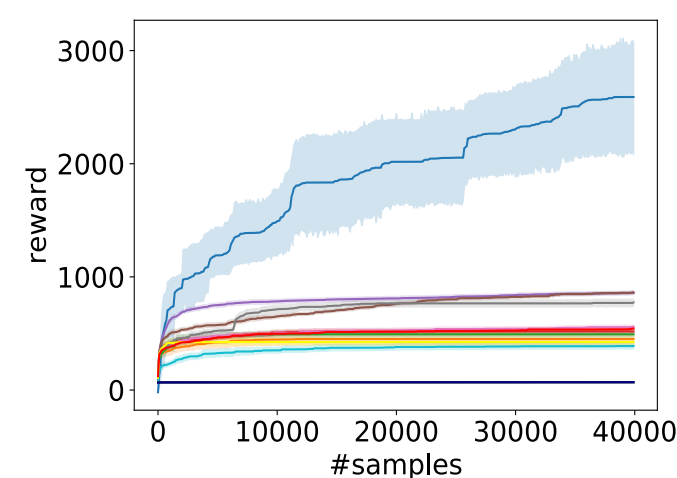
(c) Walker-2d, #params = 102



(d) Half-Cheetah, #params = 102

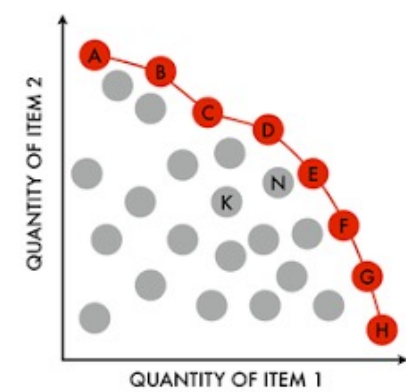


(e) Ant, #params = 888

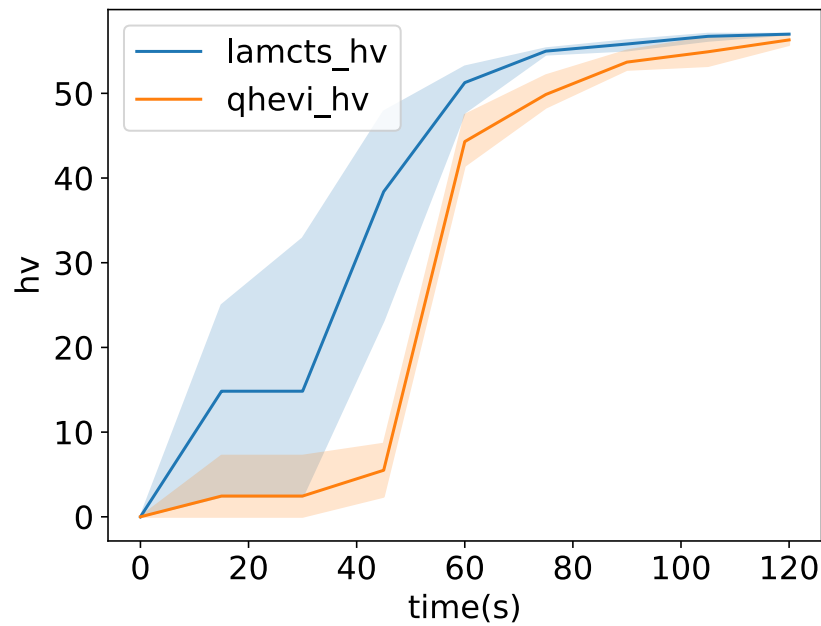


(f) Humanoid, #params = 6392

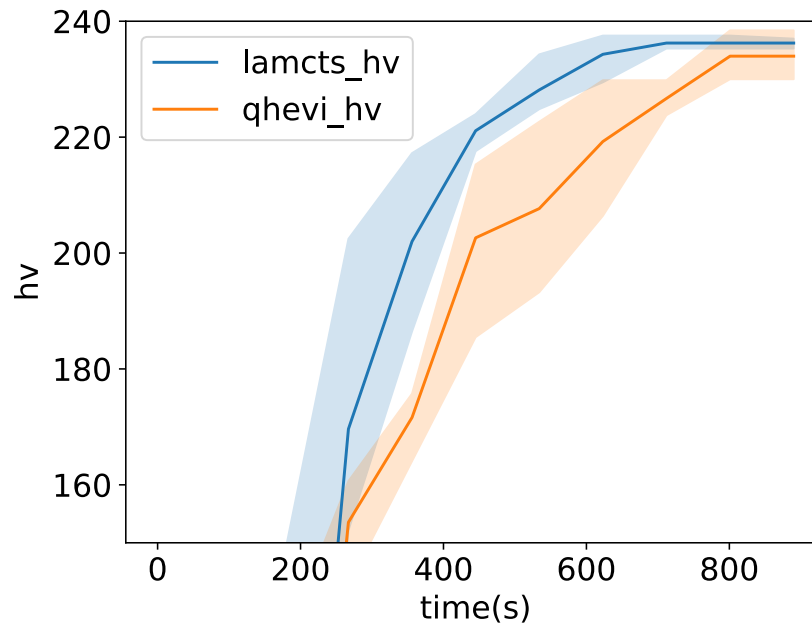
# Multi-Objective Optimization



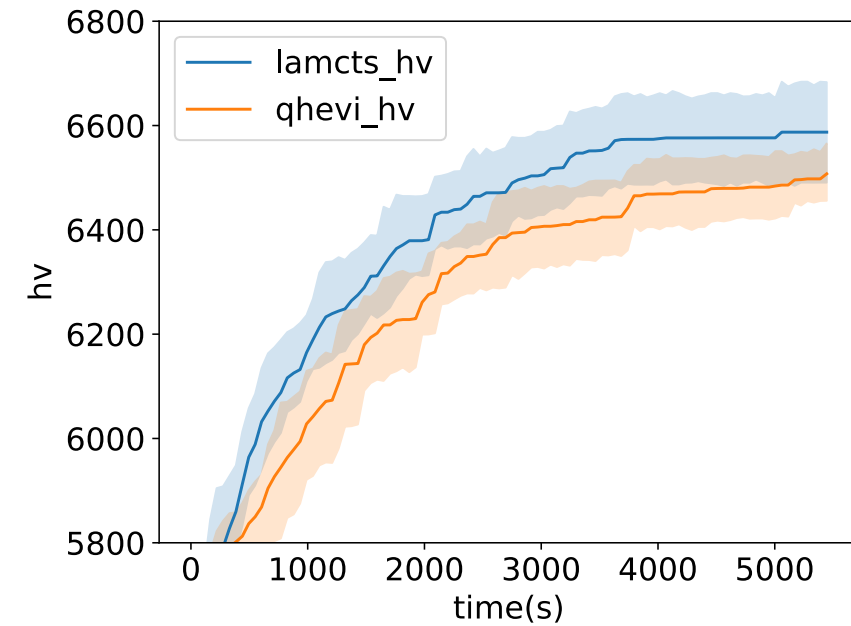
**HV:** Hyper Volume of the Pareto Frontier



Branin-Currin problem  
(2 objective)

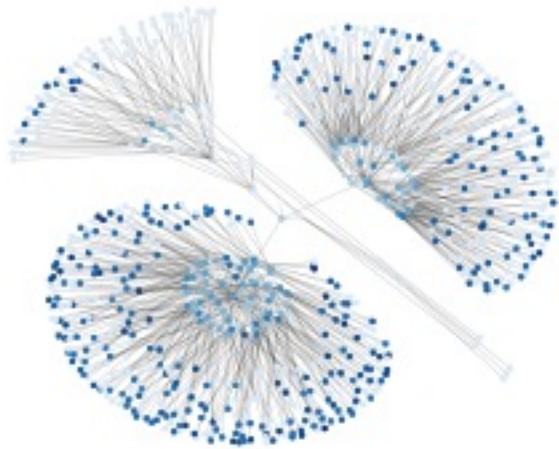


Vehicle Safety  
(3 objective)



Waveguide

Code is public now!



# LA-MCTS

<https://github.com/facebookresearch/LaMCTS>

Both 3<sup>rd</sup> and 8<sup>th</sup> teams in NeurIPS 2020 Black-box optimization competition use our method!



# Representation for Easier Search



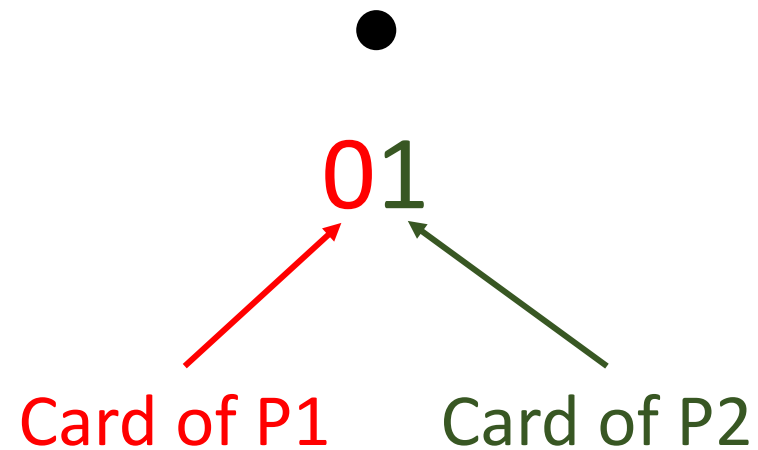
# Imperfect Information Game 101

Deal one *private* card (0/1) to **Player 1**  
Deal one *private* card (0/1) to **Player 2**

4 possibilities:

00    01    10    11

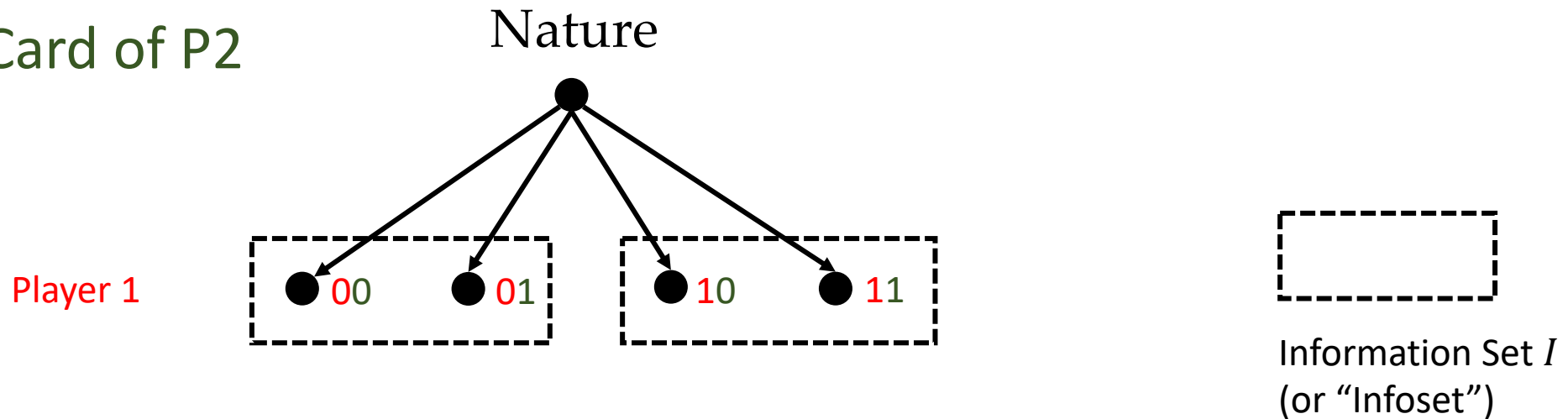
Complete information state  
or “history”  $h$



# Imperfect Information Game 101

01

Card of P1 Card of P2

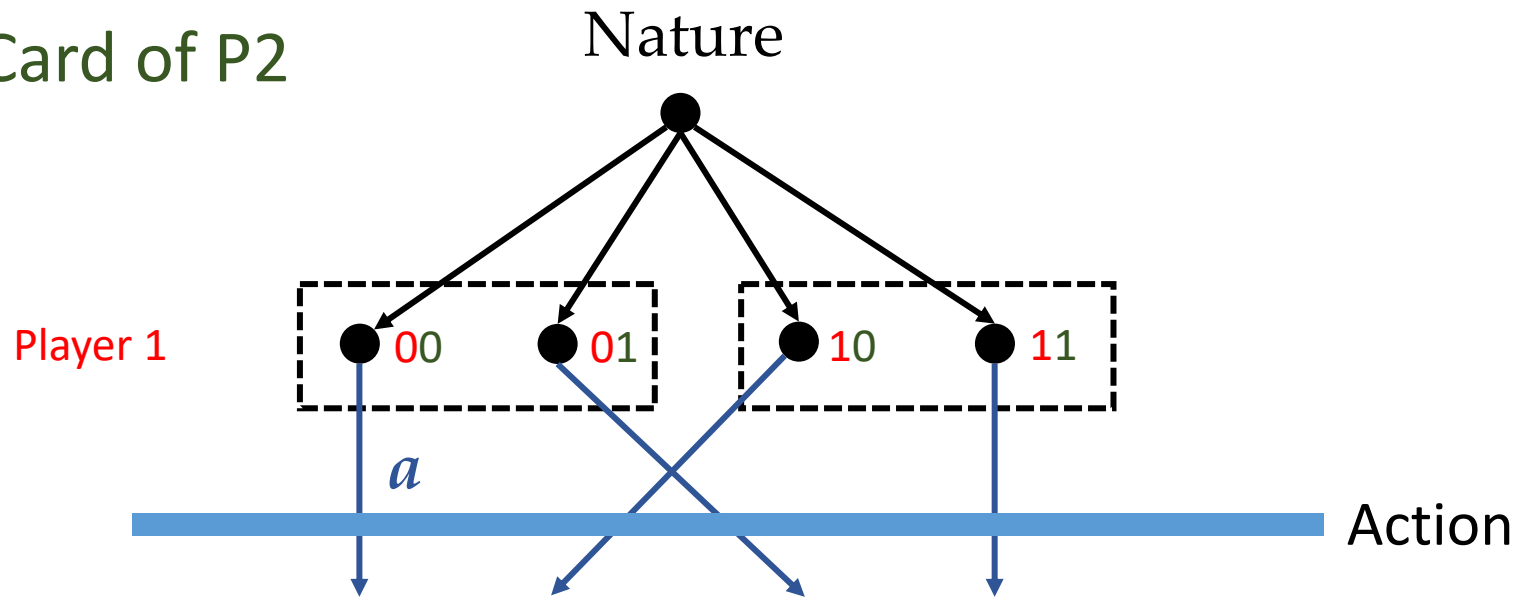


**Complete states within the same info sets have the same policy.  
(since **Player 1** doesn't know the private card of **Player 2**)**

# Imperfect Information Game 101

01

Card of P1 Card of P2

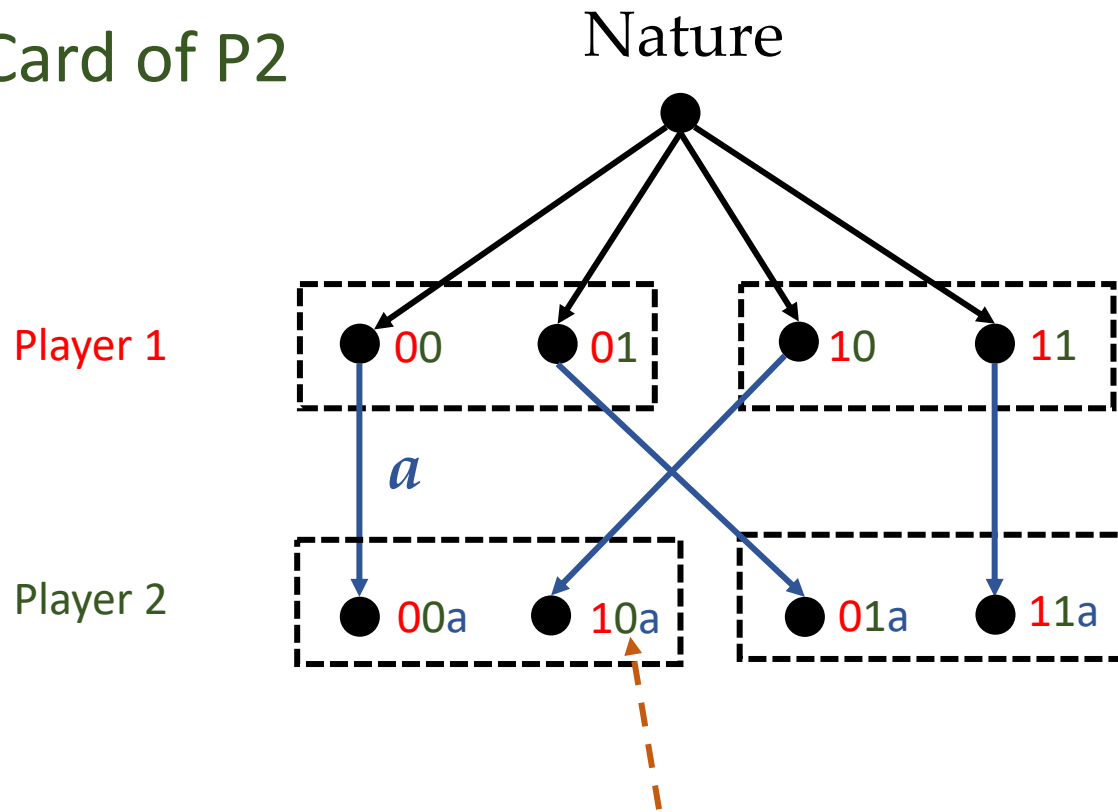




# Imperfect Information Game 101

01

Card of P1 Card of P2

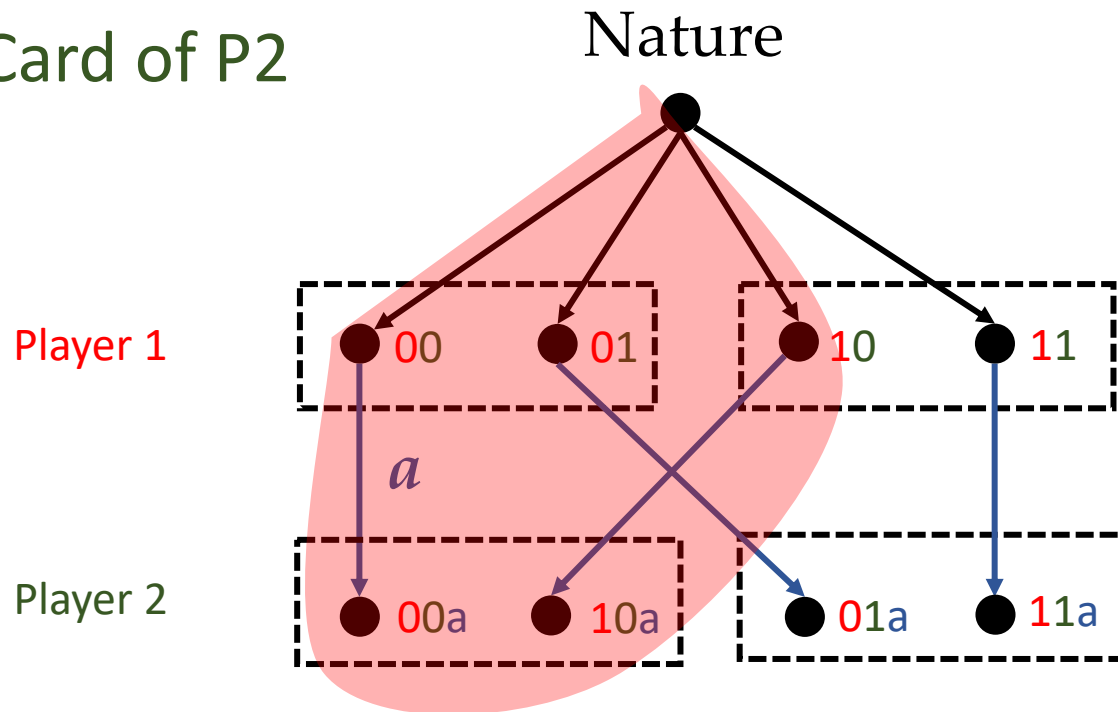


The action is public, so it is in the history of game

# Imperfect Information Game 101

01

Card of P1 Card of P2

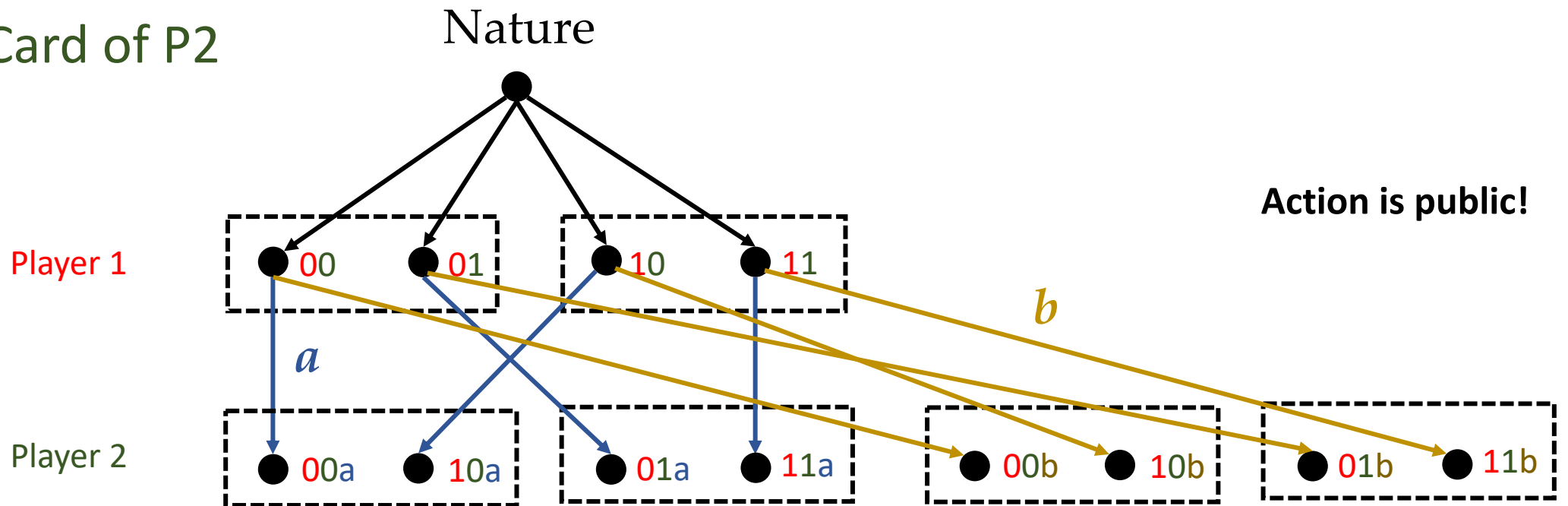


The DAG circle is here.

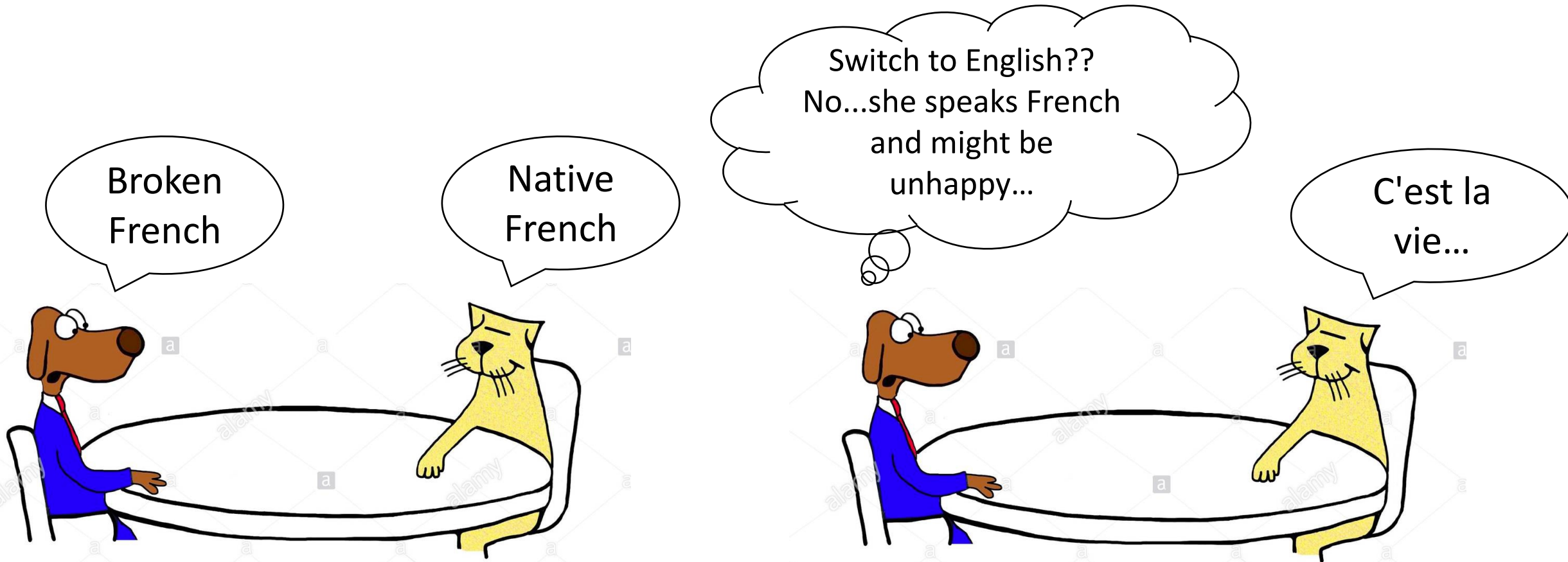
# Imperfect Information Game 101

01

Card of P1 Card of P2

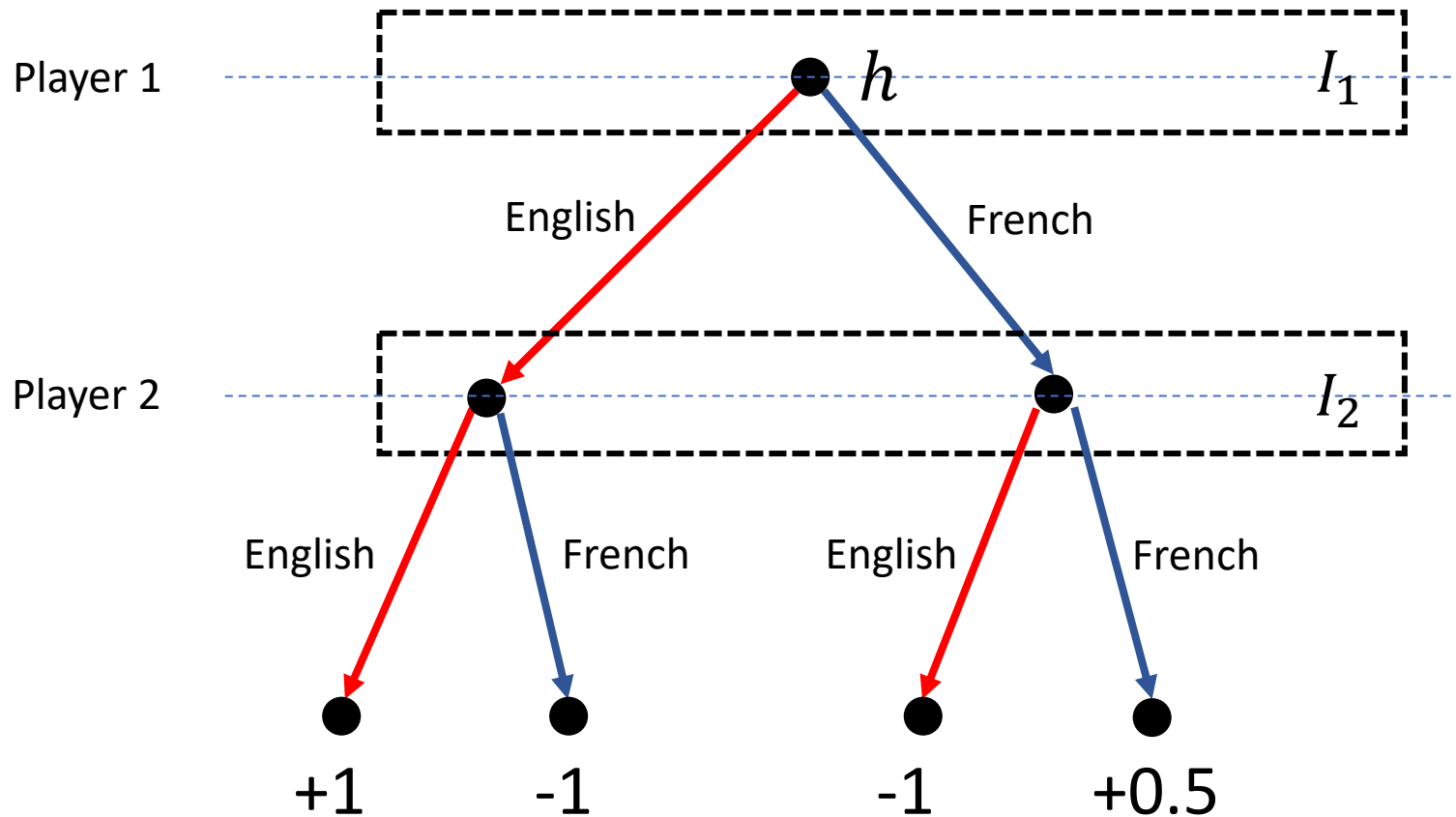


Optimizing one node at a time → Local optimum!



A **unilateral** change of policy doesn't improve co-operative communication  
(many single-agent DRL approach improves by unilateral changes of agent policy)

Optimizing one node at a time  $\rightarrow$  Local optimum!




**Player 2 makes the decision without knowing player 1's action.**

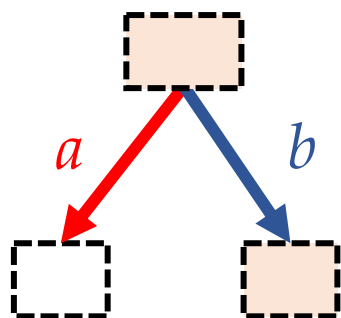
**(French, French):**  
local Nash Equilibrium +0.5

**(English, English):**  
global Nash Equilibrium +1.0

*A joint optimization of policy  $\sigma(I_1)$  and  $\sigma(I_2)$  yields optimal solution*

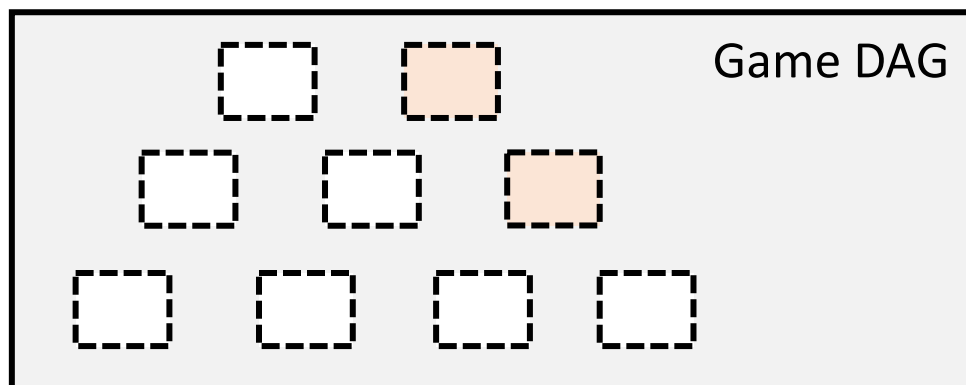
# Naïve Formulation

 *active* infosets  
 $\sigma \rightarrow \sigma'$



Change policy  $\sigma \rightarrow \sigma'$ :

Take action **a** rather than action **b** at infosets  $I$



$h \in Z$  terminal states when the reward is revealed

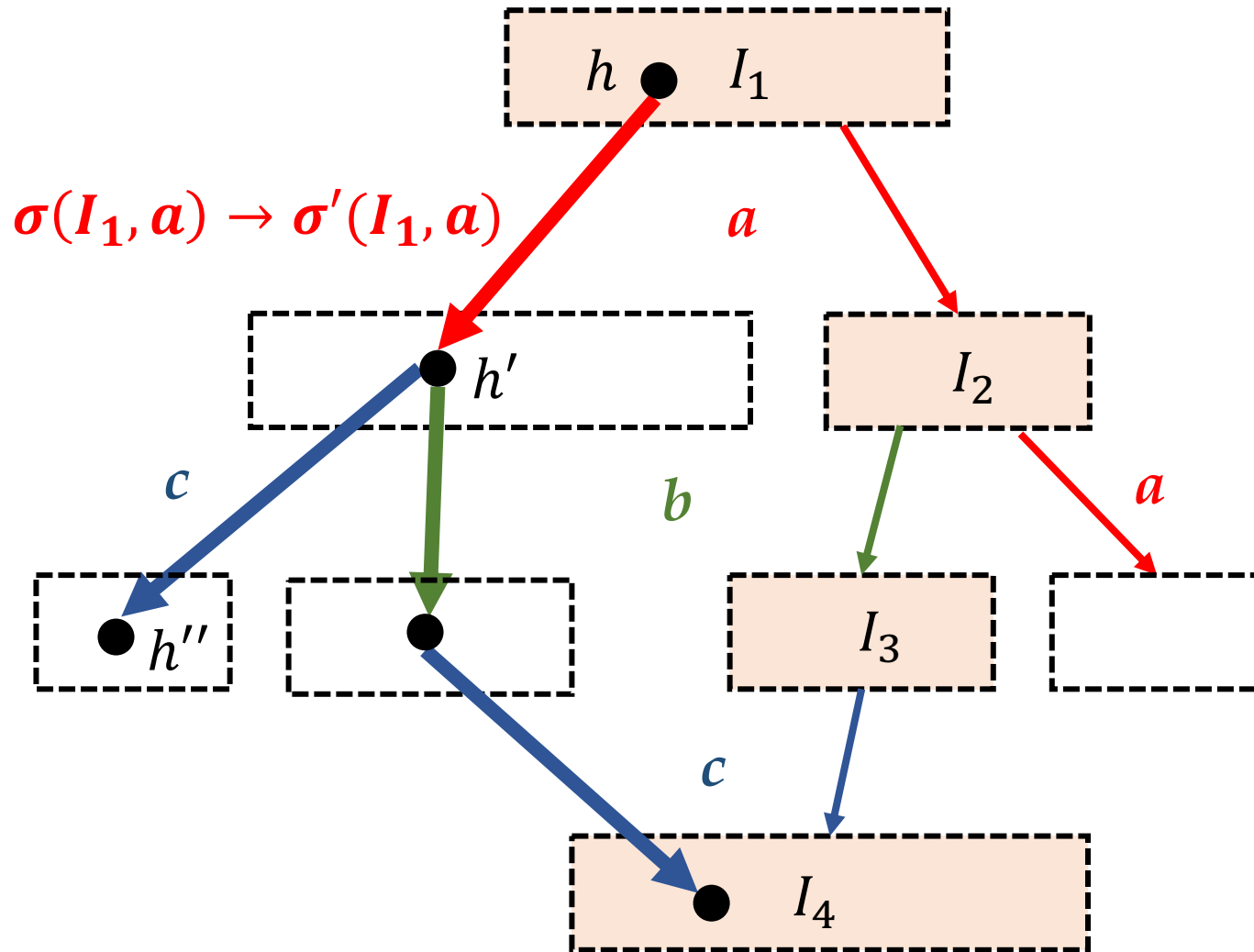
$$\text{Game value } \bar{v}^\sigma = \sum_{h \in Z} \text{Reachability } \pi^\sigma(h) \text{ Value at the terminal } v(h)$$

## Idea:

1. Pick a **subtree** and do *local* improvement
2. Multiple infosets need to be picked for *joint policy search*

But, things are complicated!

# Dependency between policies



A change of  $\sigma(I_1, a)$  affects **all** the reachability of down-stream states and/or infoSets, no matter they are *active* or not.

A trajectory could re-enter into another active set and leave and re-enter again.

The value of an inactive infoSet  $I_3$  will change since the reachability to  $I_3$  changes.

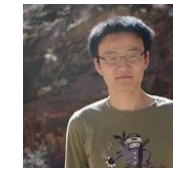
An infoSet might contain both affected states and unaffected states.

**Do we need to consider all infoSets?**

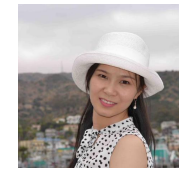
# Policy-change Density



Yuandong Tian



Qucheng Gong



Tina Jiang

[Y. Tian et al, **Joint Policy Search for Multi-agent Collaboration with Imperfect Information**, *NeurIPS 2020*]

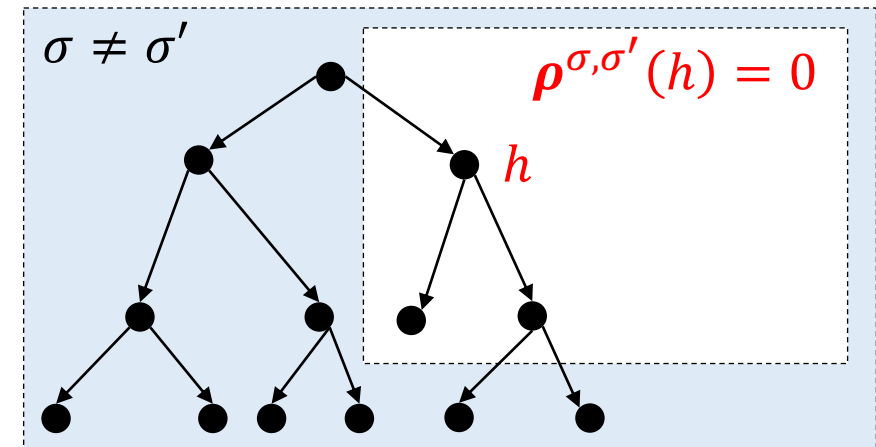
$$\rho^{\sigma, \sigma'}(h) = \pi^{\sigma'}(h) \left[ \sum_{a \in A(I)} \sigma'(I, a) v^{\sigma}(ha) - v^{\sigma}(h) \right]$$

## Two key properties:

(a) Its summation yields overall value changes

$$\bar{v}^{\sigma'} - \bar{v}^{\sigma} = \sum_{h \in Z} \rho^{\sigma, \sigma'}(h)$$

(b) For regions whose policy doesn't change, it vanishes even if policy changes at downstream/upstream states.





# Value Changes w.r.t Localized Policy Change

## Main Theorem (New representation of value change)

$$\underline{\bar{v}^{\sigma'} - \bar{v}^{\sigma}} = \sum_{\substack{I \in \mathcal{I} \\ \text{All active Infosets} \\ \text{where } \sigma' \neq \sigma}} \sum_{h \in I} \underline{\rho^{\sigma, \sigma'}(h)}$$

Overall value changes

Only depend on  $\sigma'$  via reachability  $\pi^{\sigma'}(h)$

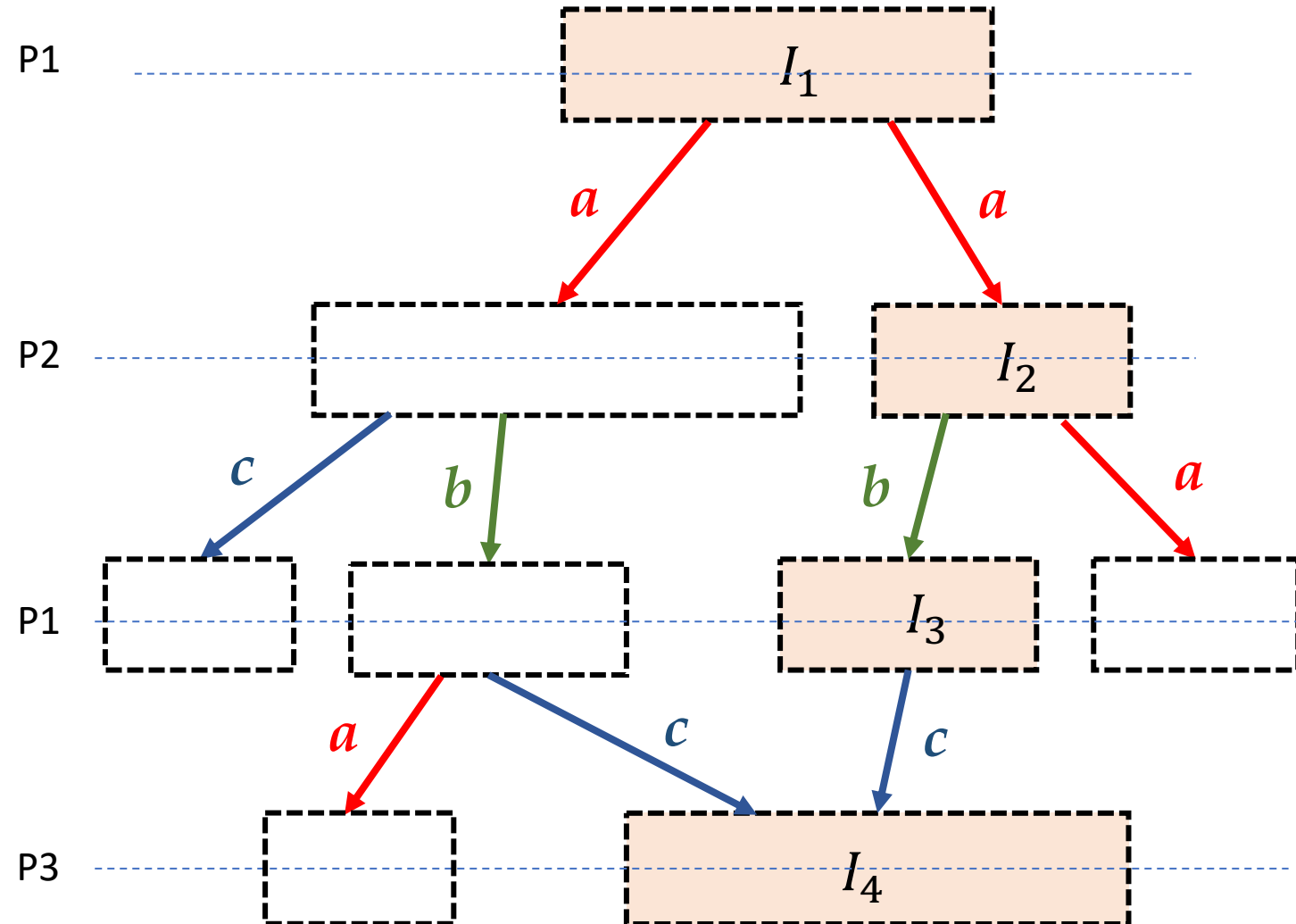
**Inactive Infosets doesn't matter!!**

# JPS (Joint Policy Search)

1. Initial infosets  $I_{\text{cand}} = \{I_1\}$
2. Pick  $I \in I_{\text{cand}}$
3. Pick an action  $a$
4. Set  $\sigma'(I, b) = \delta(a = b)$
5. Compute  $\rho^{\sigma, \sigma'}$
6. Set  $I_{\text{cand}} = \text{Succ}(I, a)$

Repeat until maximal depth  $D$  is reached.

Backtrace  
(depth-first search)



# Performance

	Comm (Def. 1)				Mini-Hanabi [15]	Simple Bidding (Def. 2)			2SuitBridge (Def. 3)		
	$L = 3$	$L = 5$	$L = 6$	$L = 7$		$N = 4$	$N = 8$	$N = 16$	$N = 3$	$N = 4$	$N = 5$
CFR1k [43]	0.89*	0.85	0.85	0.85	9.11*	2.18*	4.96*	10.47	1.01*	1.62*	2.60
<b>CFR1k+JPS</b>	<b>1.00*</b>	<b>1.00*</b>	<b>1.00*</b>	<b>1.00*</b>	<b>9.50*</b>	<b>2.20*</b>	<b>5.00*</b>	<b>10.56*</b>	<b>1.07*</b>	<b>1.71*</b>	<b>2.74*</b>
A2C [26]	0.60*	0.57	0.51	0.02	8.20*	2.19	4.79	9.97	0.66	1.03	1.71
BAD [15]	<b>1.00*</b>	0.88	0.50	0.29	9.47*	<b>2.23*</b>	4.99*	9.81	0.53	0.98	1.31
<b>Best Known</b>	1.00	1.00	1.00	1.00	10	2.25	5.06	10.75	1.13	1.84	2.89
#States	633	34785	270273	2129793	53	241	1985	16129	4081	25576	147421
#Infosets	129	2049	8193	32769	45	61	249	1009	1021	5116	24571

JPS can improve existing policies, and help it jump out of local optima

# Sample-based JPS

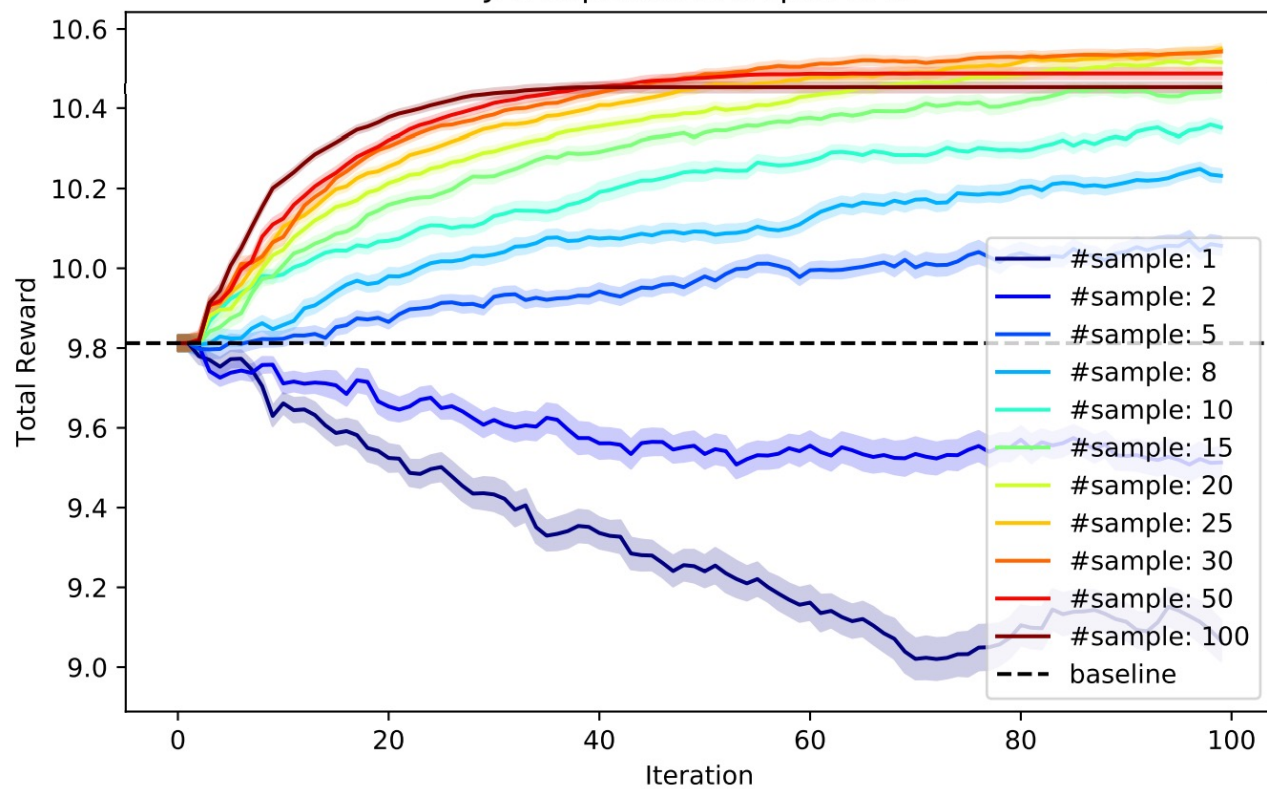
	Initialization	All states	#Samples per info set							
			1	2	5	8	15	20	25	30
Mini-Hanabi [15]	CFR1k [43]	9.50	<b>10.00</b>	9.99	9.95	9.75	9.51	9.51	9.51	9.51
SimpleBidding ( $N = 16$ )	CFR1k [43]	10.56	10.47	10.47	10.49	10.52	10.58	10.60	10.61	<b>10.61</b>
SimpleBidding ( $N = 16$ )	BAD [15]	10.47	9.91	9.95	10.22	10.34	10.50	10.55	<b>10.57</b>	10.55
2-suited Bridge ( $N = 3$ )	BAD [15]	1.12	0.89	1.12	<b>1.13</b>	1.13	1.12	1.12	1.12	1.12
2-suited Bridge ( $N = 4$ )	BAD [15]	1.71	1.23	1.63	1.71	<b>1.71</b>	1.68	1.67	1.68	1.69
2-suited Bridge ( $N = 5$ )	BAD [15]	2.77	2.12	2.51	2.74	<b>2.79</b>	2.79	2.76	2.77	2.78

$$\bar{v}^{\sigma'} - \bar{v}^{\sigma} = \sum_{I \in \mathcal{I}} \sum_{h \in I} \rho^{\sigma, \sigma'}(h)$$

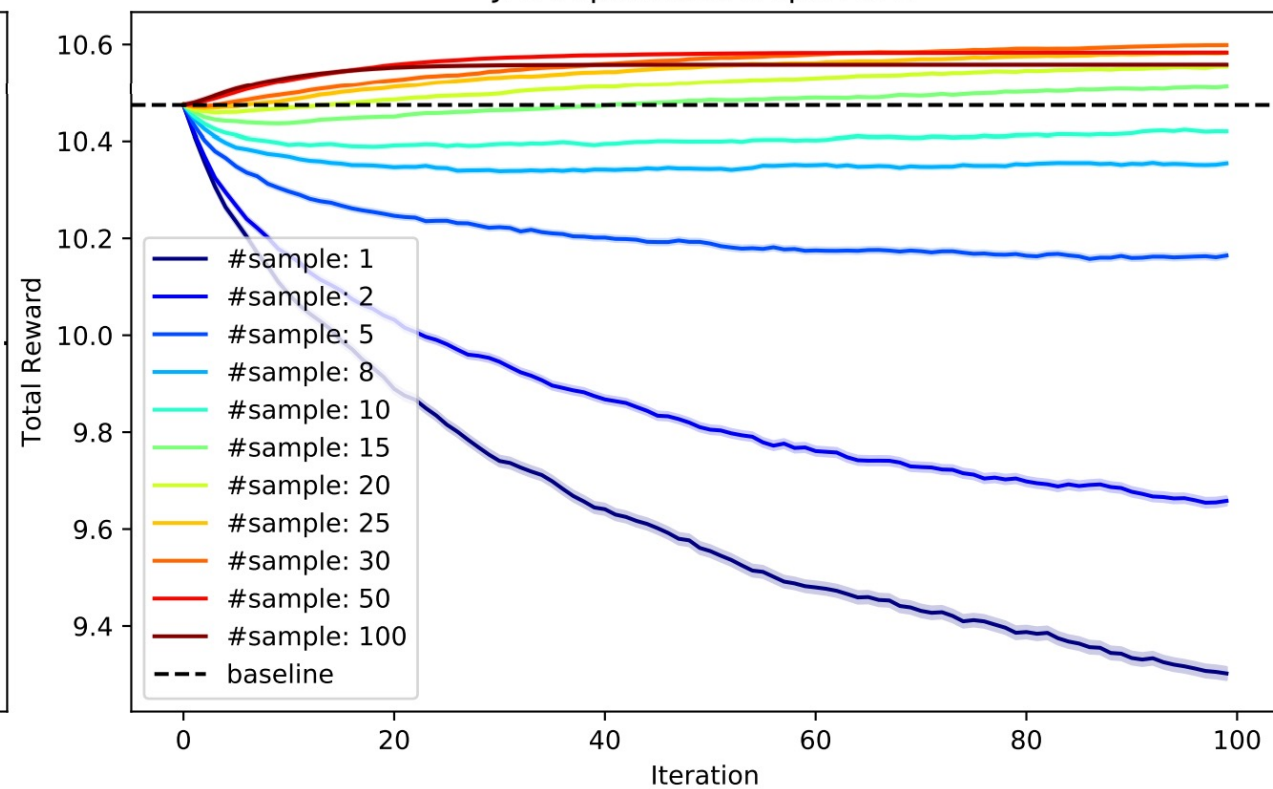
Sample  $h$  in each  $I$

# Sample-based JPS

JPS improvement upon BAD



JPS improvement upon CFR



Simple Biddings ( $N=16$ )

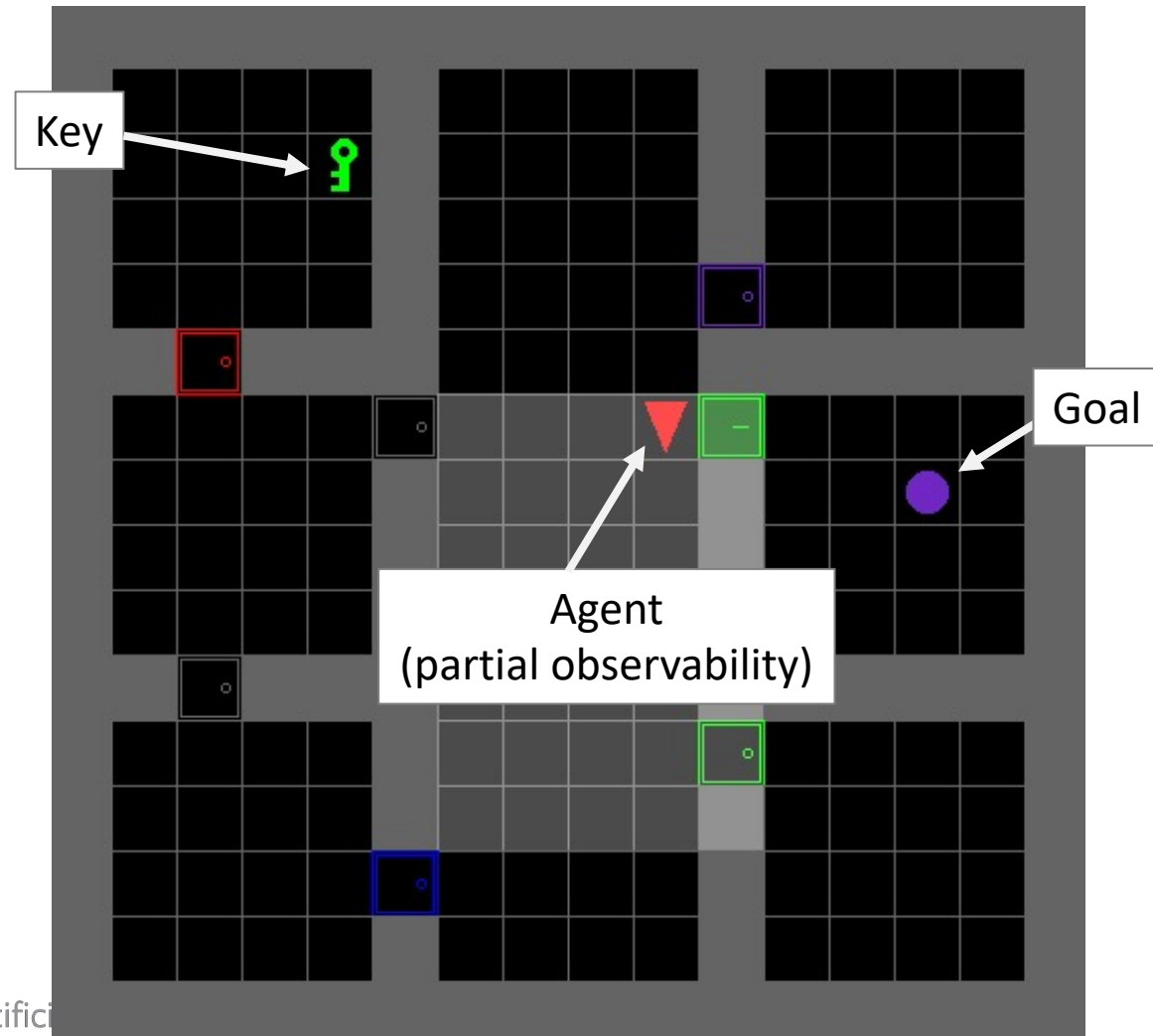
# Contract Bridge Bidding

Methods	Vs. WBridge5 (1000 games) (IMPs/board)
Previous SoTA (Rong et al, 2019)	+ 0.25 (on 64 games)
Our A2C baseline	+ 0.29 ± 0.22
1% JPS (2 days)	+ 0.44 ± 0.20
5% JPS (2 days)	+ 0.37 ± 0.19
1% JPS (14 days)	+ <b>0.63 ± 0.20</b>

**WBridge5:** Champions of computer bridge tournament in 2005, 2007, 2008, 2016-2018

# Representation for RL Exploration

# Exploration in Environment with Sparse Reward

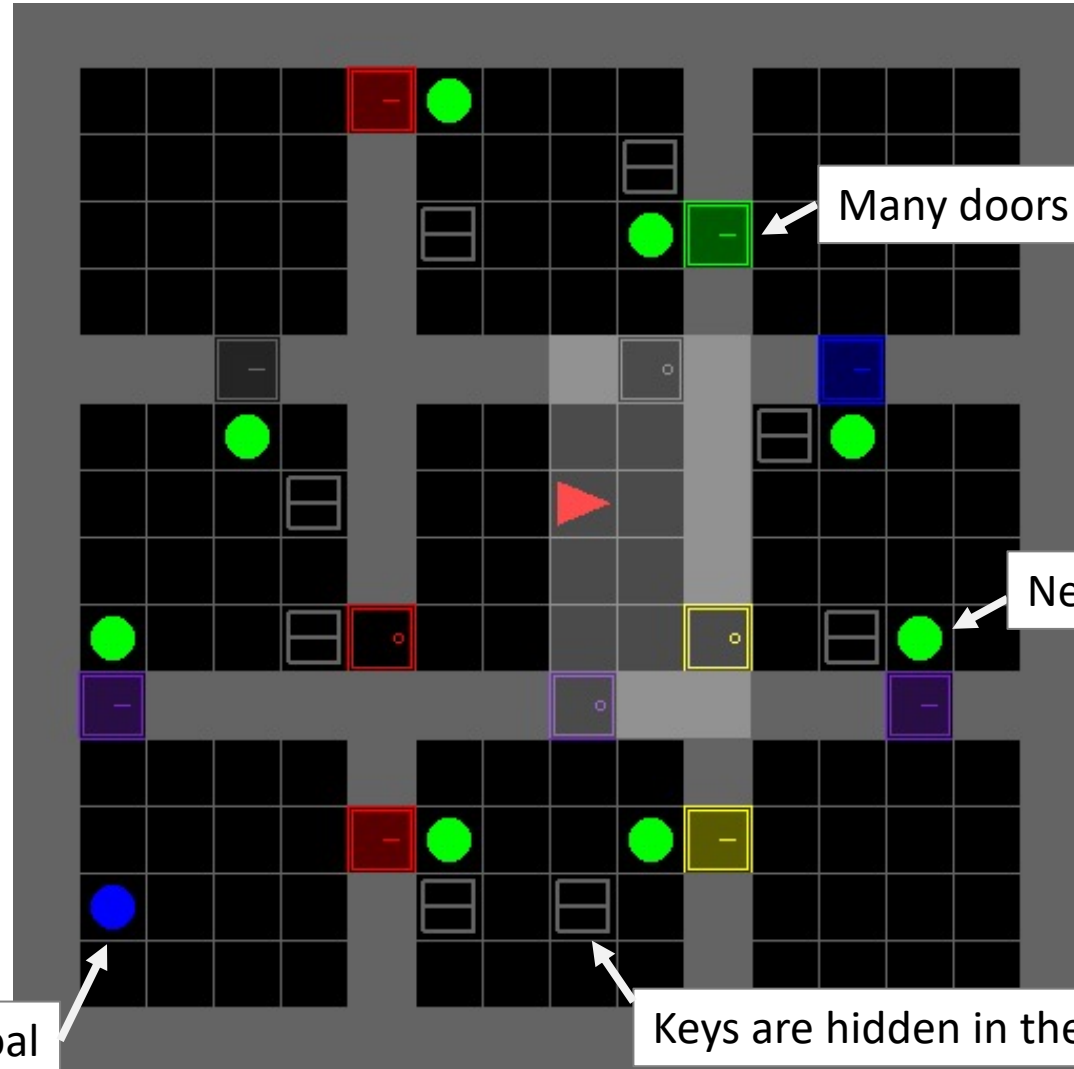


## No external reward

- when agent wanders around.
- when agent picks the key
- when agent opens all doors
- when agent opens the locked door
- ...
- until the agent reaches the goal



# And more complicated situations...



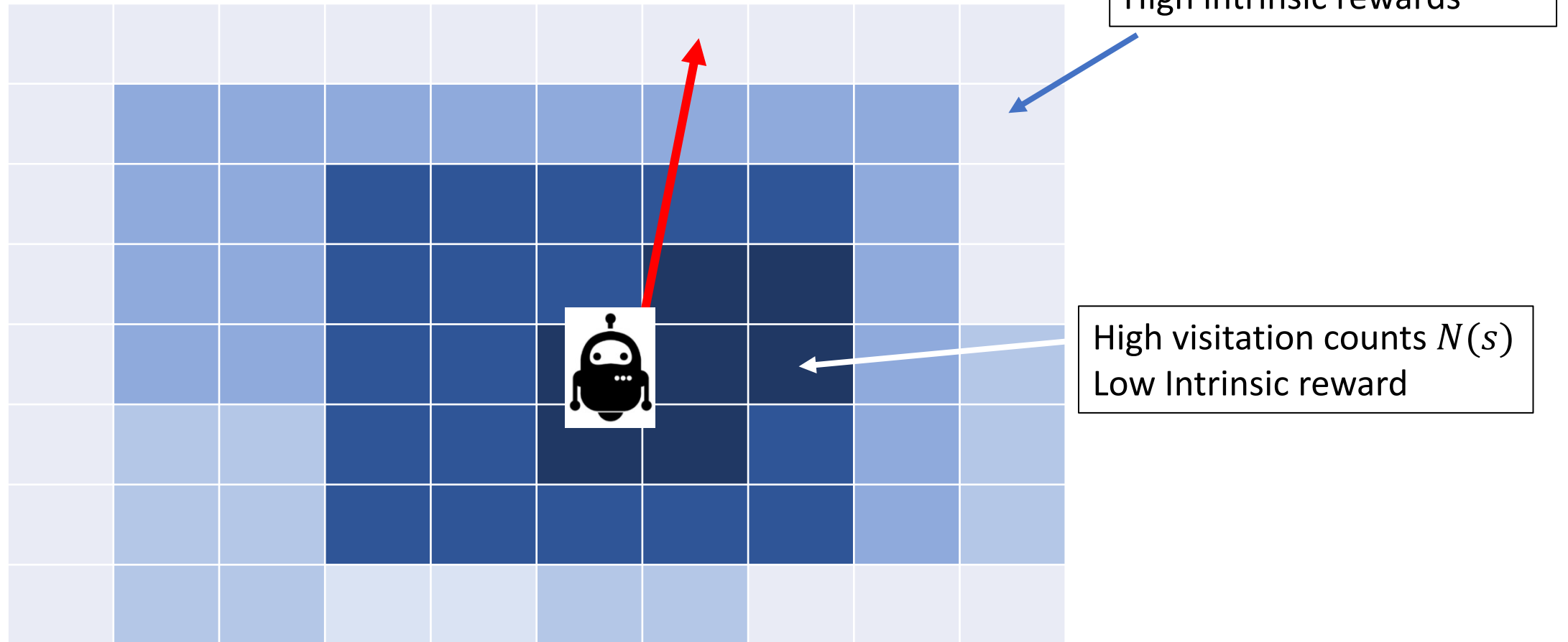
Many doors are locked

Need to move obstacles around

Goal

Keys are hidden in the boxes

# Count-based Exploration

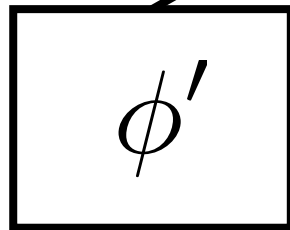


What if we have exponential #states?

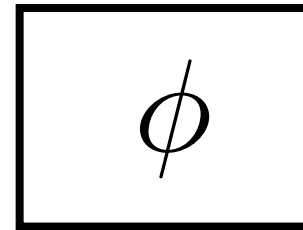
# Random Network Distillations (RND)

$$N(\mathbf{s}) \approx \frac{1}{\|\phi'(\mathbf{s}) - \phi(\mathbf{s})\|}$$

Online Network



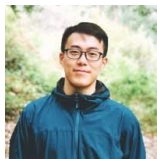
Random fixed target network



$\mathbf{s}$

Familiar state = low prediction error

# BeBold



Tianjun Zhang



Huazhe Xu



Xiaolong Wang



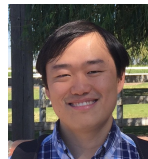
Yi Wu



Kurt Keutzer

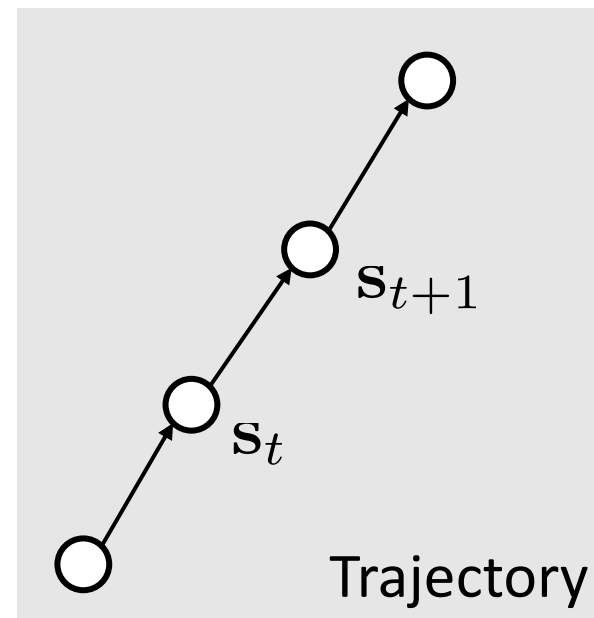


Joseph Gonzalez



Yuandong Tian

BeBold = **B**eyond the **B**oundary of **E**xplored **R**egions



RND for t+1

RND for t

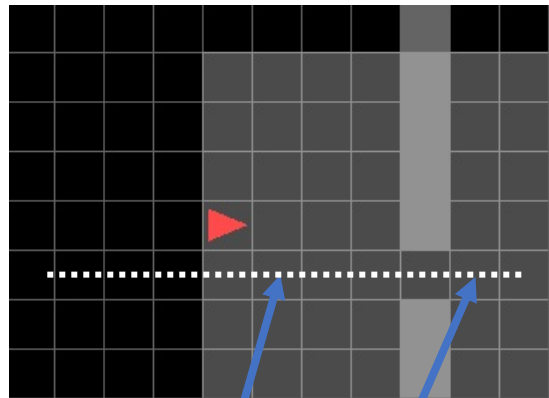
$$r^i(\mathbf{s}_t, \mathbf{a}_t) = \max(\|\phi(\mathbf{o}_{t+1}) - \phi'(\mathbf{o}_{t+1})\|_2 - \|\phi(\mathbf{o}_t) - \phi'(\mathbf{o}_t)\|_2, 0) * \mathbb{1}\{\underline{N_e(\mathbf{o}_{t+1})} = 1\}$$

Observation (rather than full state)

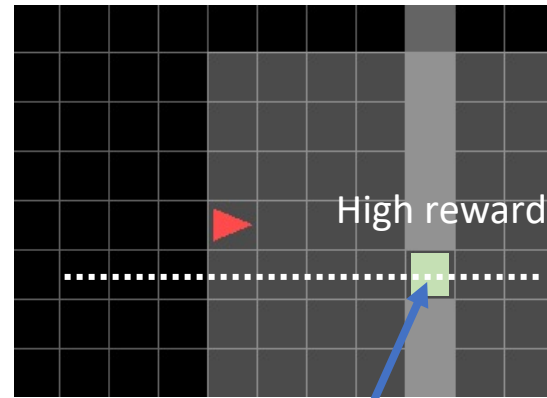
Episodic visitation count  
(Hash Table)

# BeBold

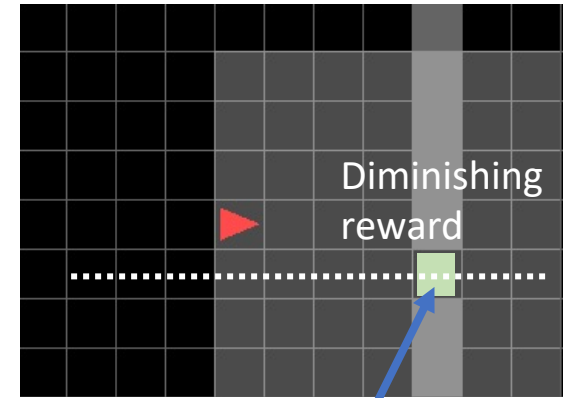
Repeat



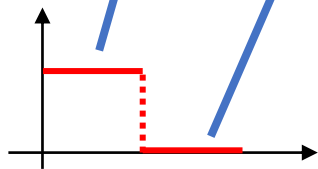
Exploration



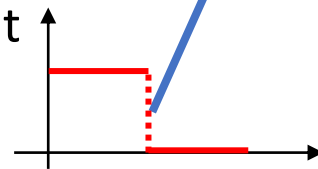
Policy learning



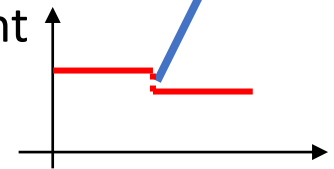
Visitation  
count



Visitation  
count



Visitation  
count



## RND



1. RND assigns high IR (dark green) throughout the environment



2. RND temporarily focuses on the upper right corner (yellow)



3. RND by chance starts exploring the bottom right corner heavily, resulting in the IR at top right higher than bottom right



4. RND re-explores the upper right and forgets the bottom right, gets trapped

## BeBold



1. BeBold assigns high IR (dark red) near the start and low IR for the rest (light red)



2. BeBold pushes every direction to the frontier of exploration uniformly (yellow)



3. BeBold continuously pushes the exploration frontier



4. BeBold reaches the end of exploration

# MiniGrid

	MRN6	MRN7S-8	MRN12-S10	KCS3R3	KCS4R3	KCS5R3	KCS6R3	OM2DI-h	OM2DI-hb	OM1Q	OM2Q	OMFULL
ICM				✓								
RND				✓				✓				
RIDE	✓	✓	✓	✓	✓			✓				
AMIGO				✓								
BeBold	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓

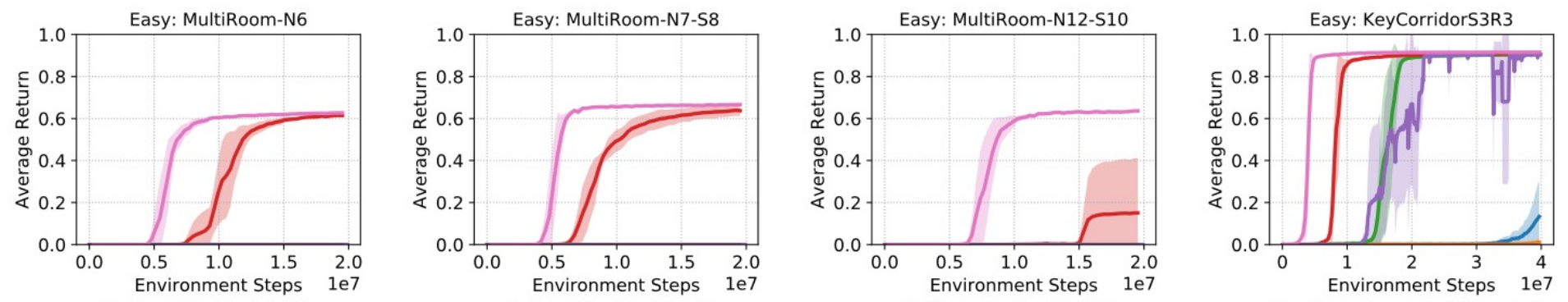
✓ : Solved within 120M steps

\*MR is short for MultiRoom, KC is for KeyCorridor, OM is for ObstructedMaze

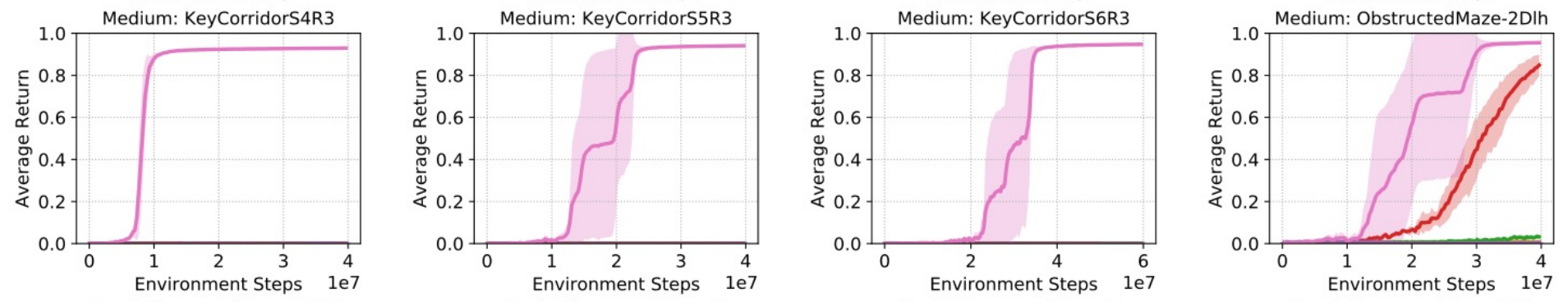
[Chevalier-Boisvert, Maxime, Lucas Willems, and Suman Pal. "Minimalistic gridworld environment for openai gym." GitHub repository (2018)]

IMPALA ICM RND RIDE AMIGO BeBold

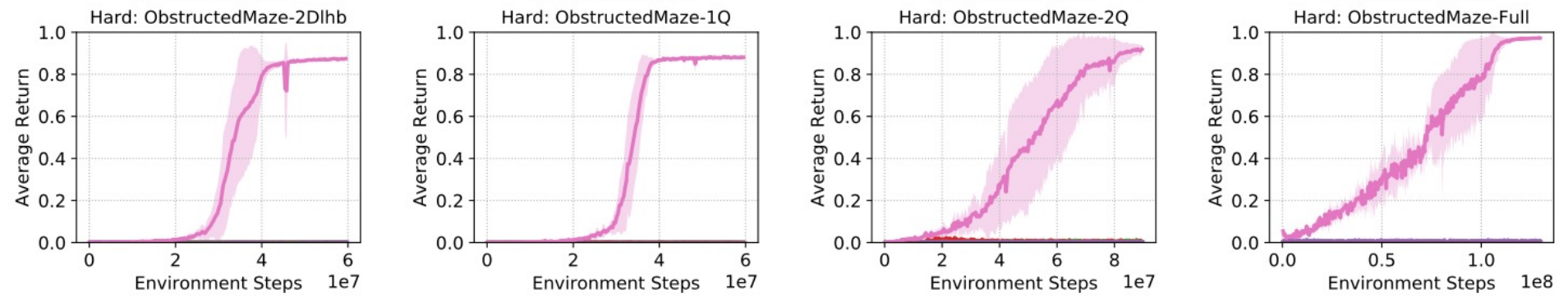
Easy



Medium



Hard

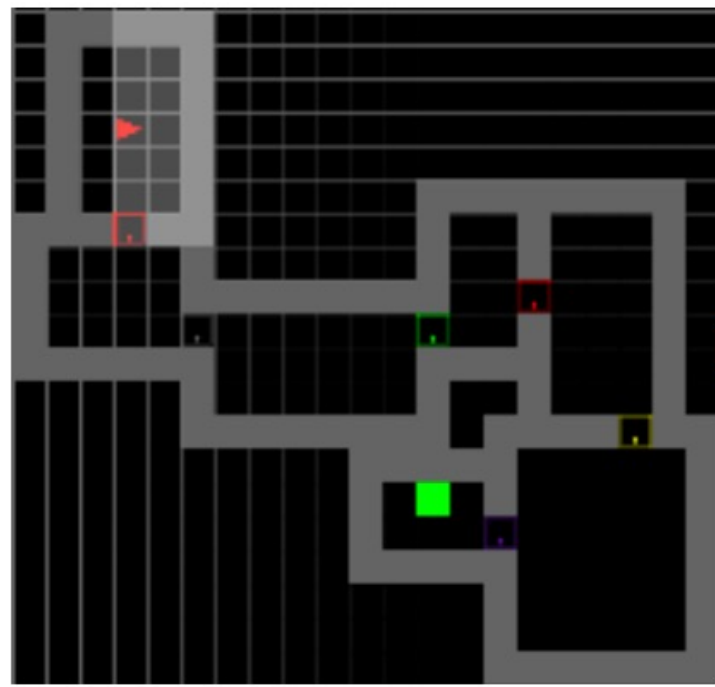


AMIGO: [C. Andres, et al. "Learning with AMIGo: Adversarially Motivated Intrinsic Goals." ICLR 2021]

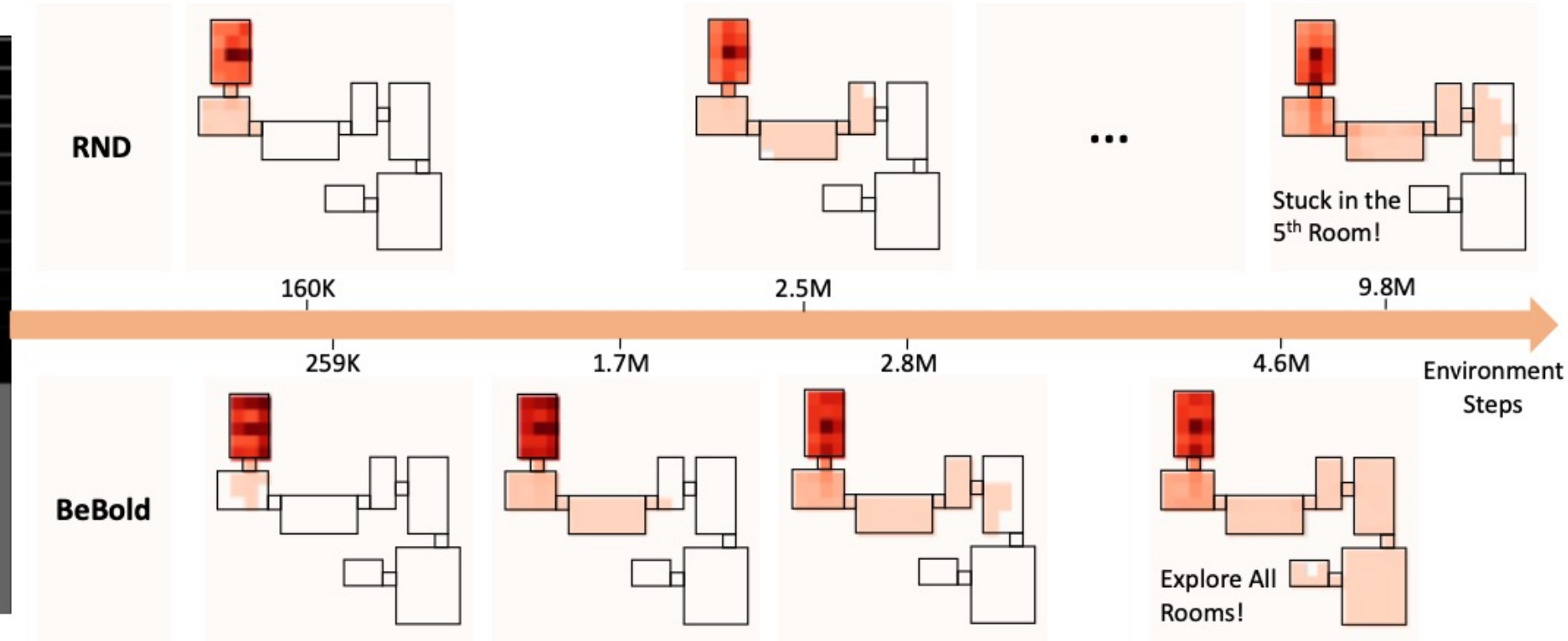
RIDE: [R. Roberta, and Tim Rocktäschel. "RIDE: Rewarding Impact-Driven Exploration for Procedurally-Generated Environments.", ICLR 2020]



# Pure Exploration in MiniGrid



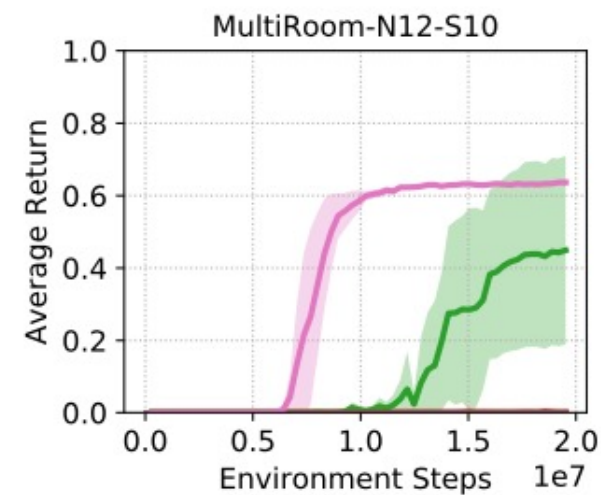
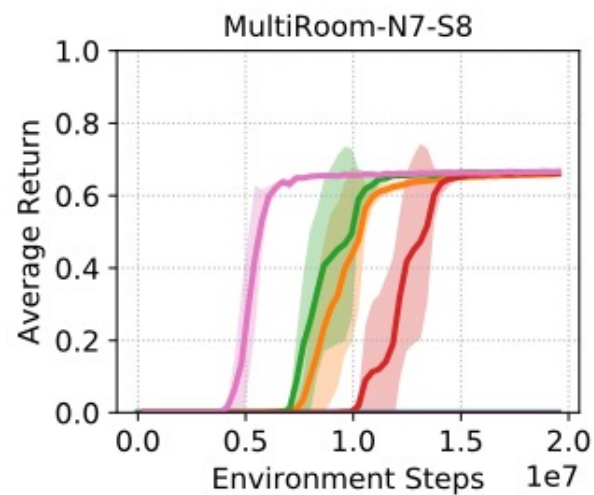
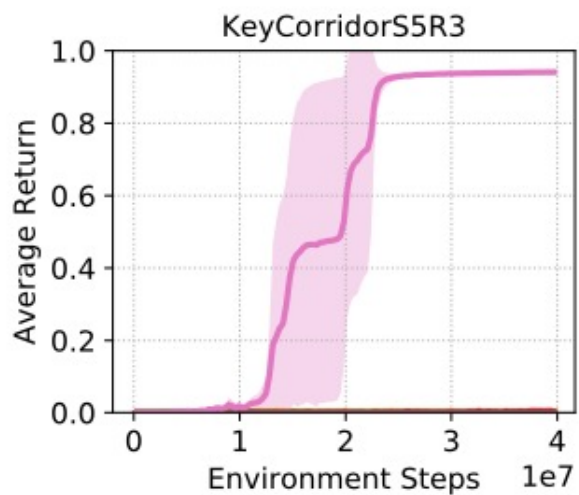
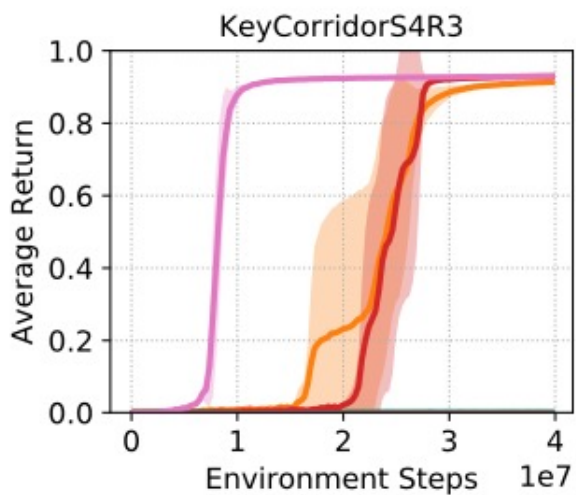
MultiRoomN7S8



# Ablation Study

$$r^i(\mathbf{s}_t, \mathbf{a}_t) = \max(\|\phi(\mathbf{o}_{t+1}) - \phi'(\mathbf{o}_{t+1})\|_2 - \|\phi(\mathbf{o}_t) - \phi'(\mathbf{o}_t)\|_2, 0) * \mathbb{1}\{N_e(\mathbf{o}_{t+1}) = 1\}$$

— RND — RND with ERIR — BeBold w.o. ERIR — BeBold w.o. Clipping — BeBold

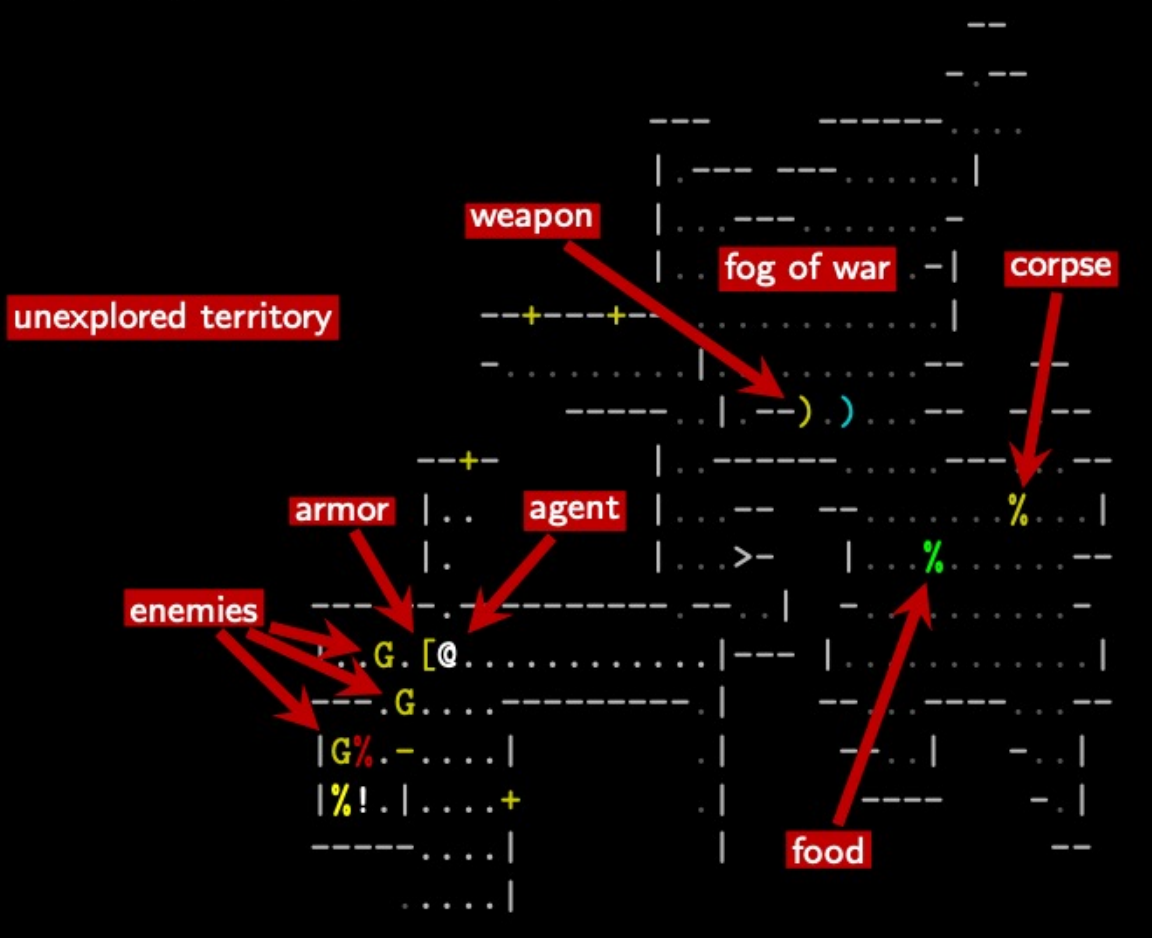


# NetHack

➤ You kill the dwarf! Welcome to experience level 5.--More--

**Legend**

- " -- Amulet
- ) -- Weapon
- [ -- Armor
- ! -- Potion
- ? -- Scroll
- / -- Wand
- = -- Ring
- + -- Spellbook
- \* -- Gem
- ( -- Tool
- 0 -- Boulder
- \$ -- Gold
- % -- Comestible



The image shows a NetHack game map with various annotations. A red box labeled "unexplored territory" points to a region of the map. A red box labeled "weapon" points to a yellow ')' symbol. A red box labeled "fog of war" points to a greyed-out area. A red box labeled "corpse" points to a yellow '%' symbol. A red box labeled "armor" points to a yellow '[' symbol. A red box labeled "agent" points to a yellow '@' symbol. A red box labeled "enemies" points to yellow 'G' symbols. A red box labeled "food" points to a green '%' symbol.

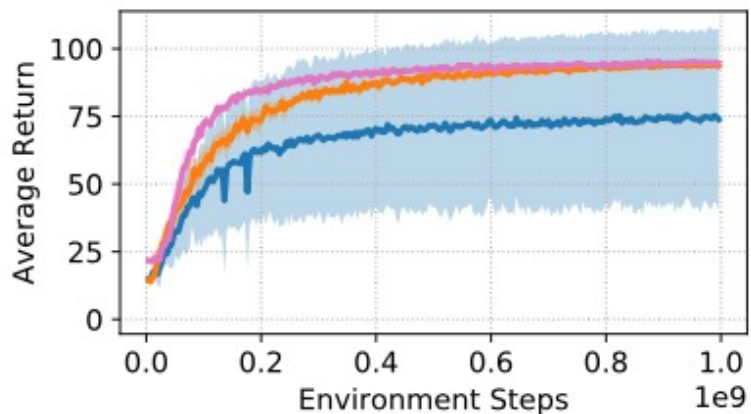
Agent61322 the Novice      St:18/02 Dx:12 Co:12 In:11 Wi:13 Ch:8 Neutral S:  
Dlv1:5 \$:0 HP:37(39) Pw:25(25) AC:5 Xp:5/168 T:768 Hungry

Agent States

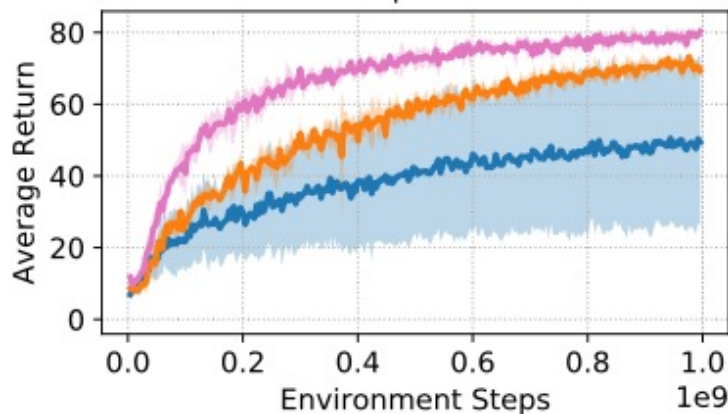
# 6 Tasks in NetHack

— IMPALA — RND — BeBold

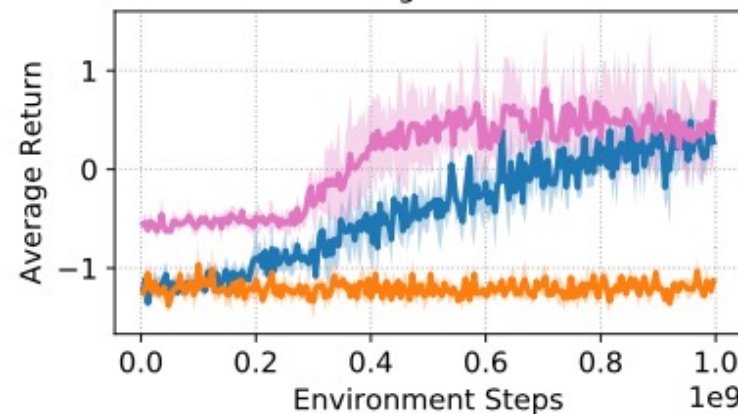
staircase



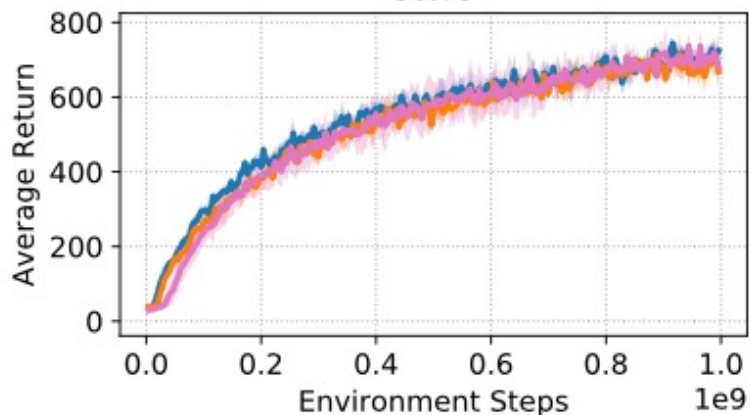
pet



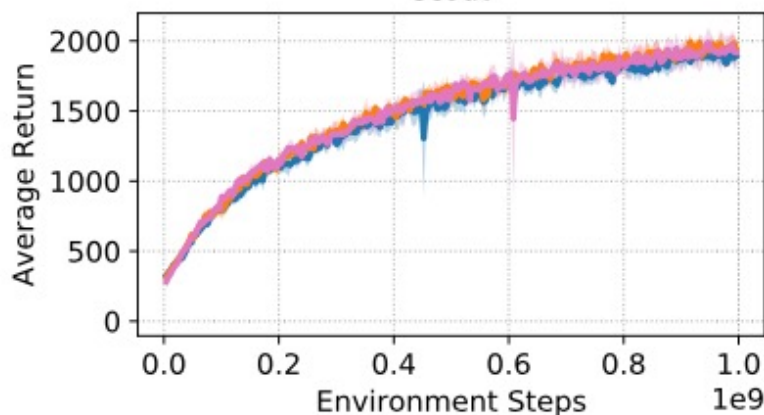
gold



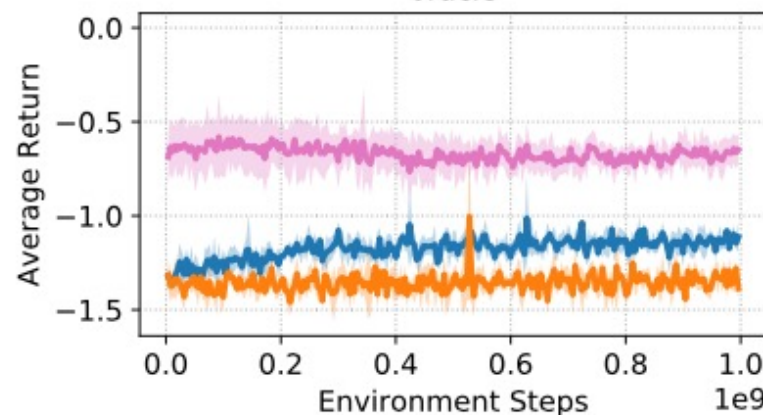
score



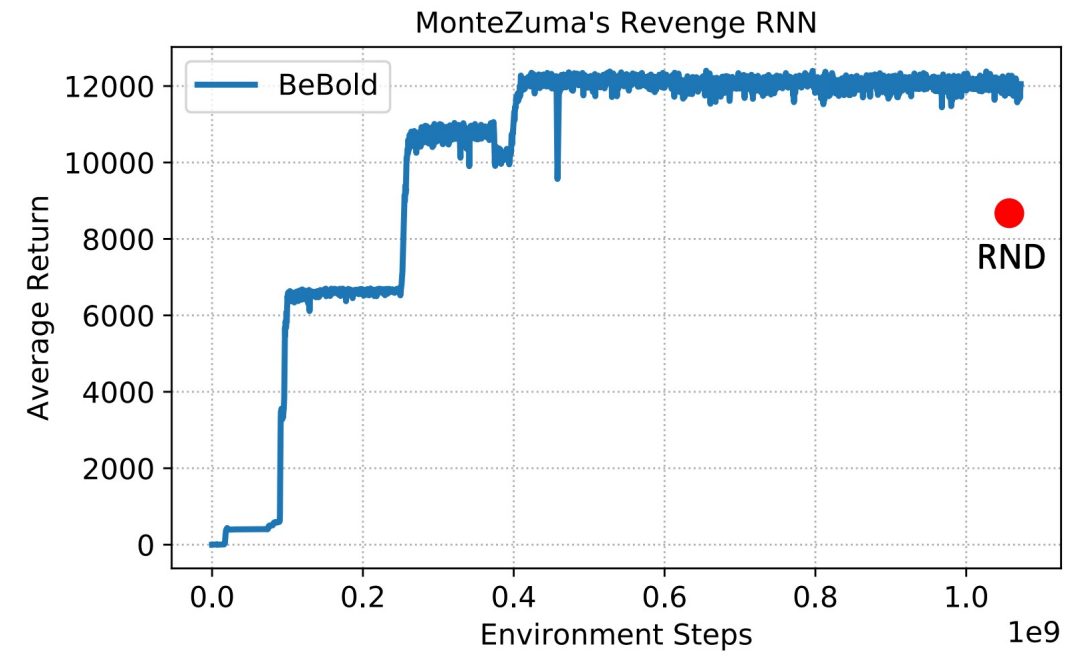
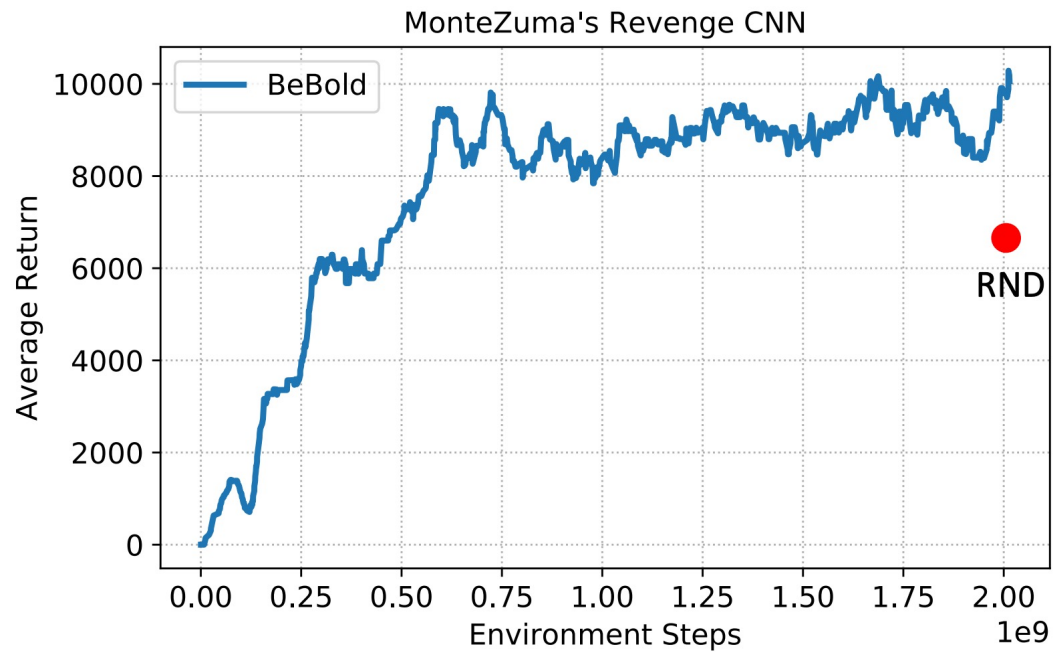
scout



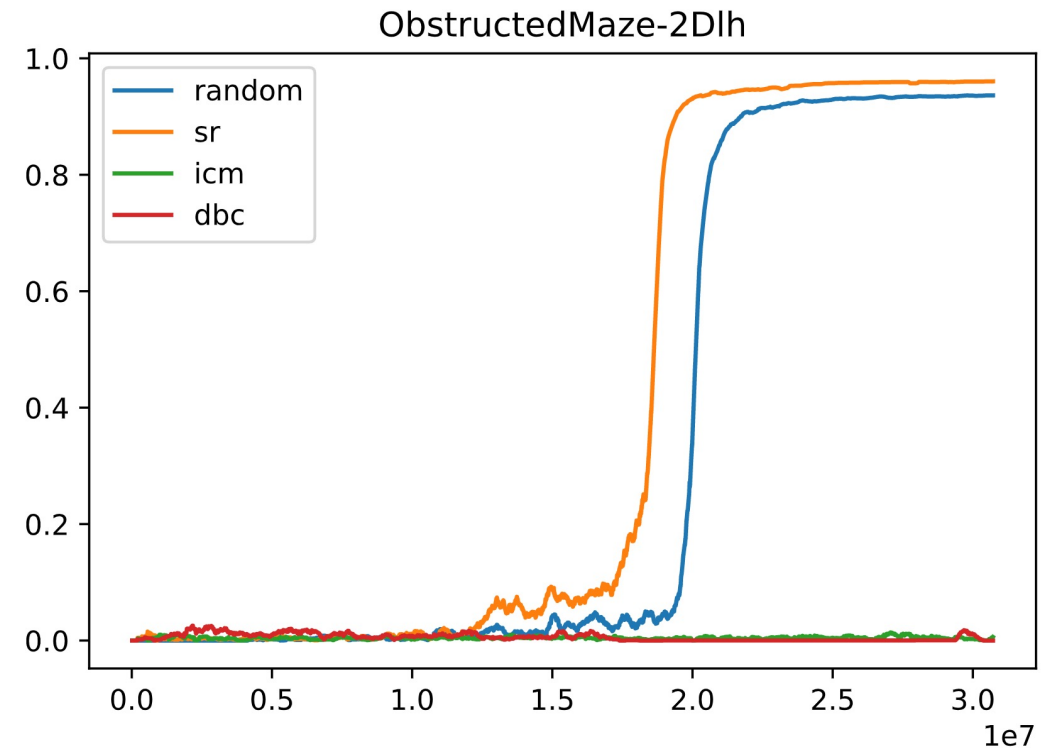
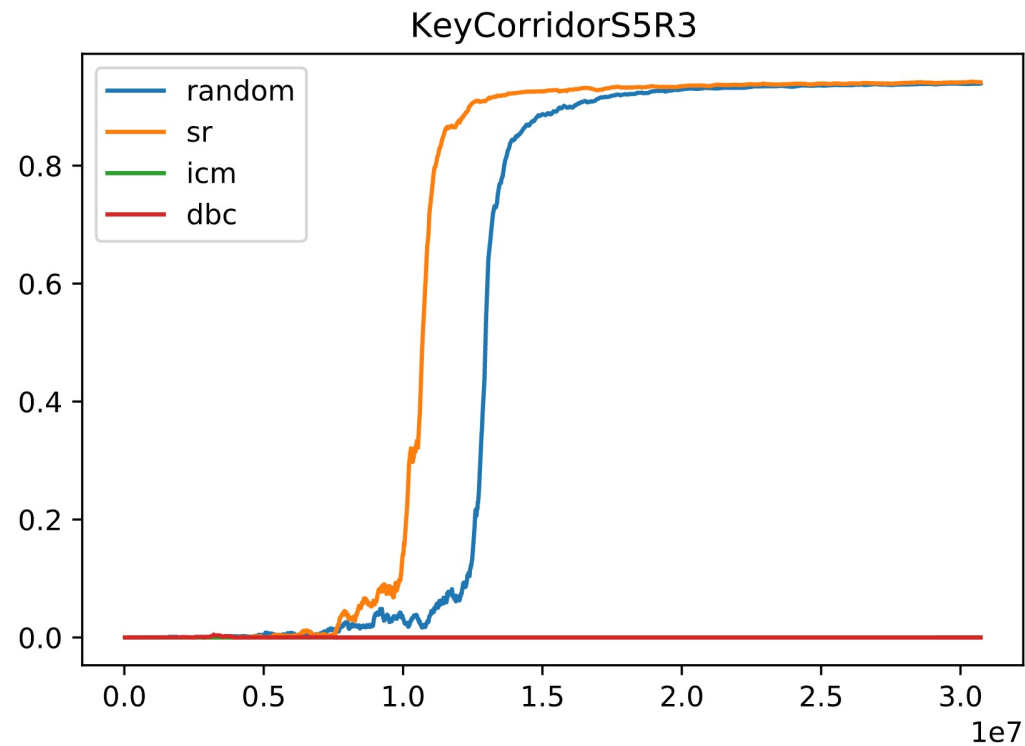
oracle



# MonteZuma's Revenge



# Huge Performance Difference with different Representations



Random = vanilla BeBold  
DBC = deep bisimulation control  
SR = Successor Representation

Thanks!