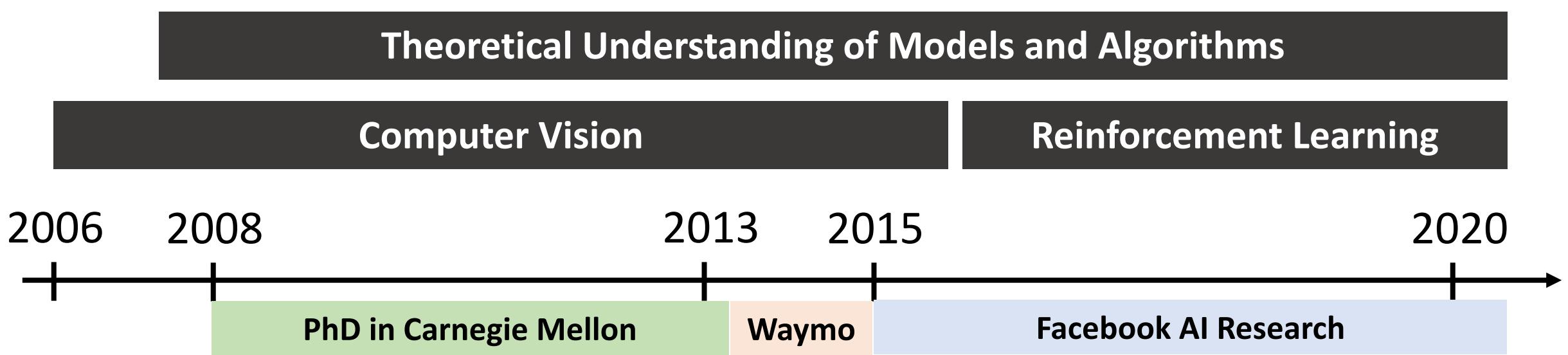
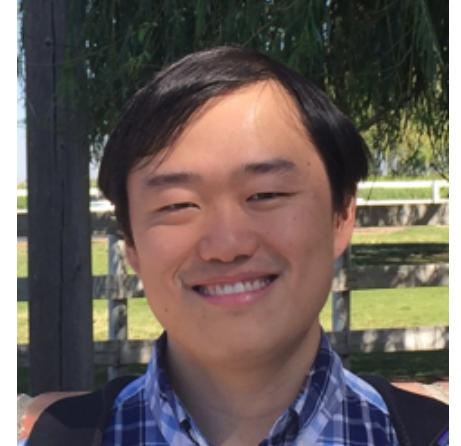


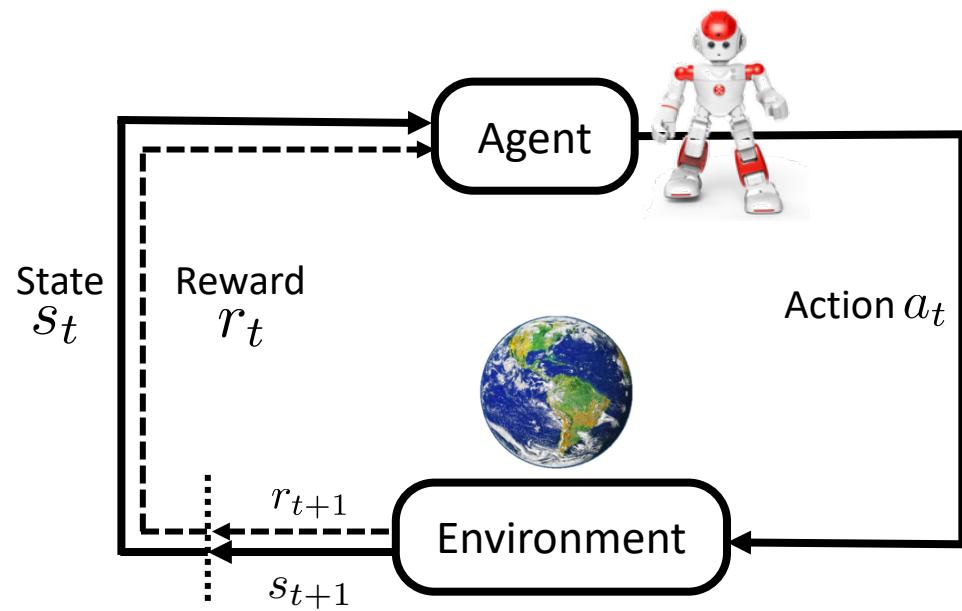
Rethinking Challenging Problems in Reinforcement Learning

Yuandong Tian
Research Scientist
Facebook AI Research

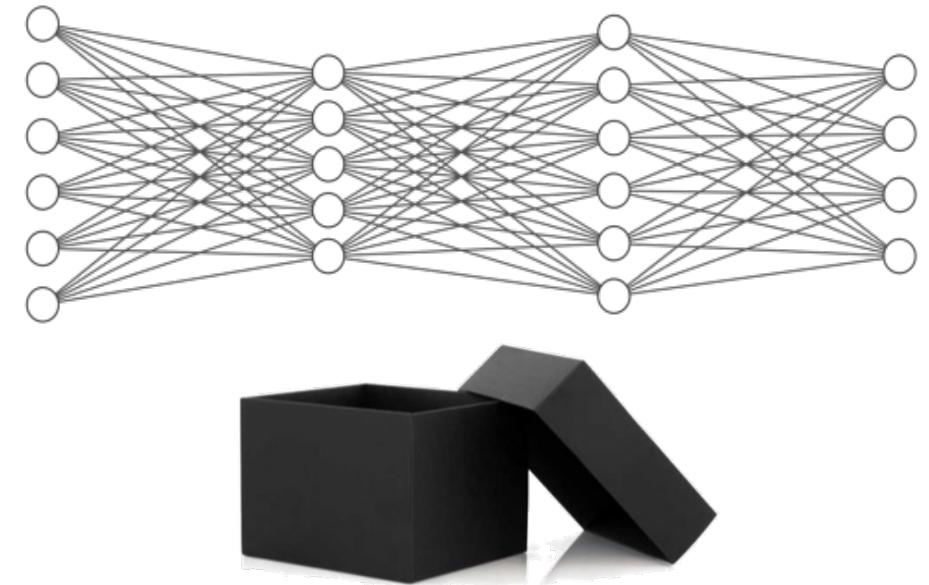
Career Path



Research Directions



Reinforcement Learning



Understanding of Deep Models

Challenging Problems in RL

Compounding
Errors

Exploration

Demonstration

Instability

Safety

Sample Complexity

Agent Collaboration

Sim2Real

State Decomposition

Off-policy Data

Convergence

Skill Acquisition

Hierarchical RL

Challenging Problems in RL

Compounding
Errors

Exploration

Demonstration

Instability

Safety

Sample Complexity

Agent Collaboration

Sim2Real

State Decomposition

Off-policy Data

Convergence

Hierarchical RL

Skill Acquisition

Joint Policy Search for Multi-agent Collaboration with Imperfect Information



Yuandong Tian



Qucheng Gong

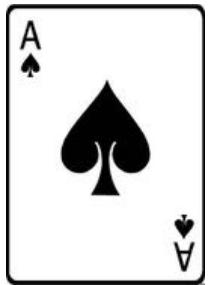


Tina Jiang

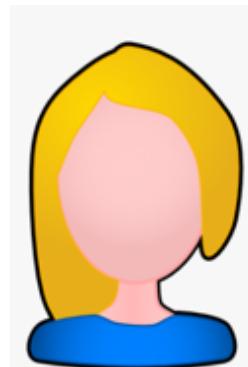
Facebook AI Research

An Illustrative Example

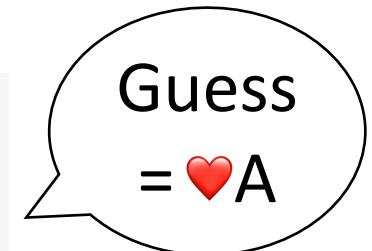
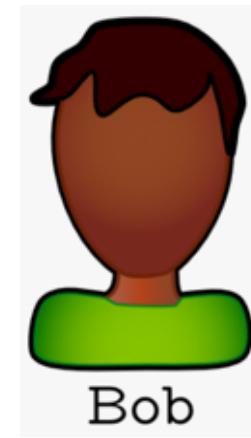
Private Card



or



Public Signal
1 or 2 or 3



One possible solution (6 symmetric solutions):

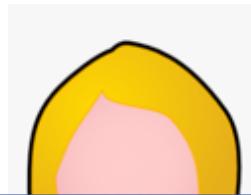
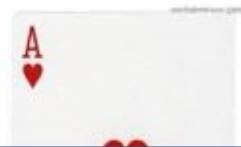
Private card	Alice's Action	Bob's Action
♥ A	1	Guess ♥ A
♠ A	3	Guess ♠ A
--	2	--

Not used

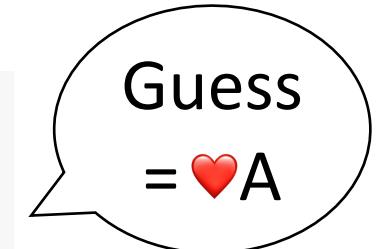
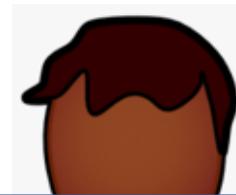
What if Alice and Bob never use signal 2,
but sending signal 2 come with additional rewards?

An Illustrative Example

Private Card



Public Signal
1 or 2 or 3



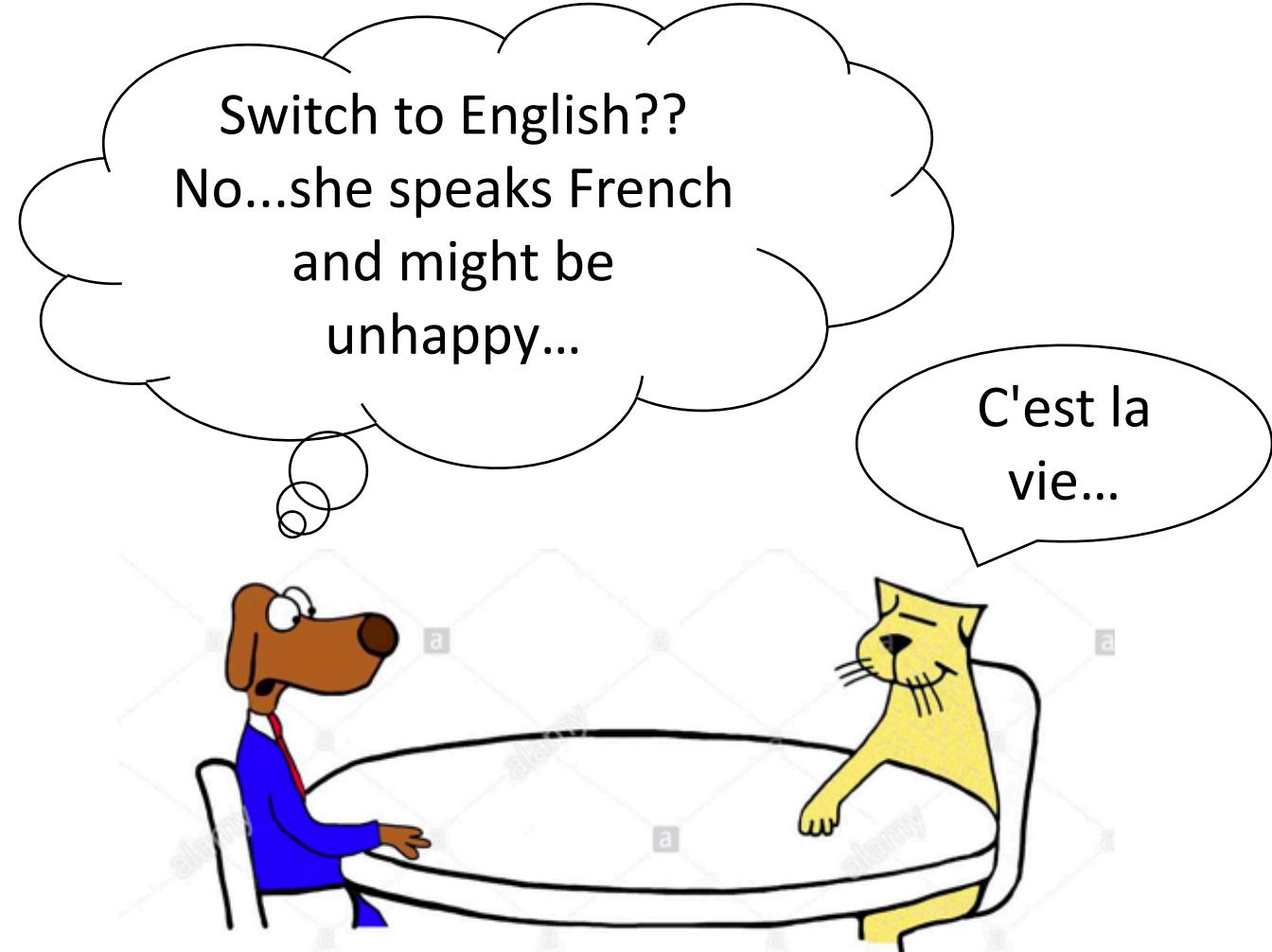
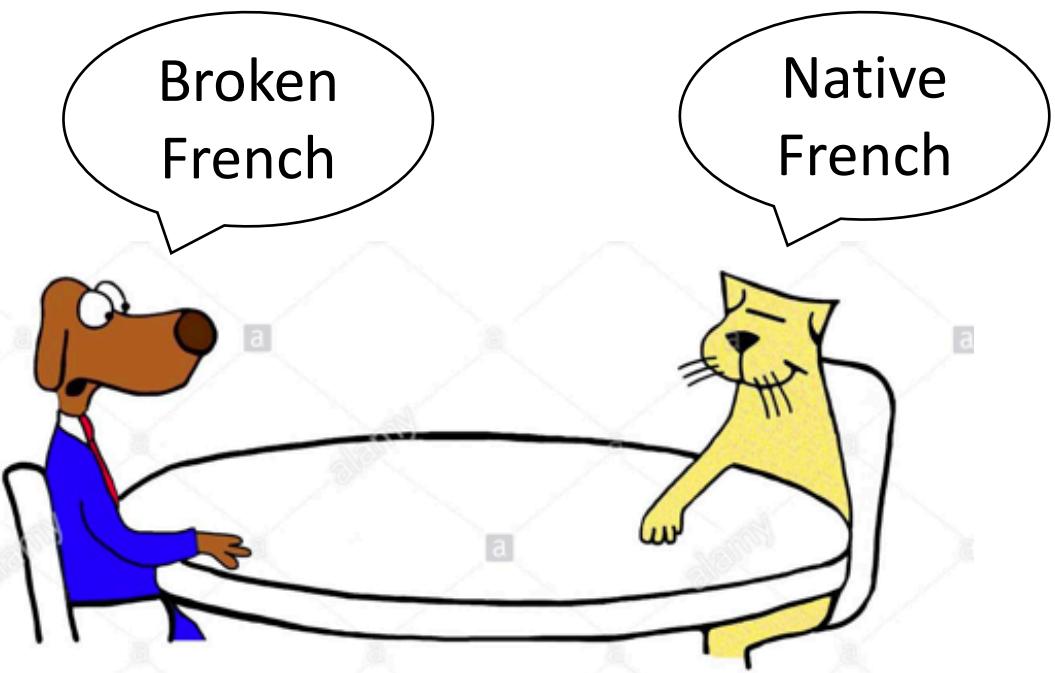
For pure multi-agent collaborative games, A unilateral optimization of policy doesn't improve overall value.

▼ A	1	Guess ▼ A
♠ A	3	Guess ♠ A
--	2	--

Not used

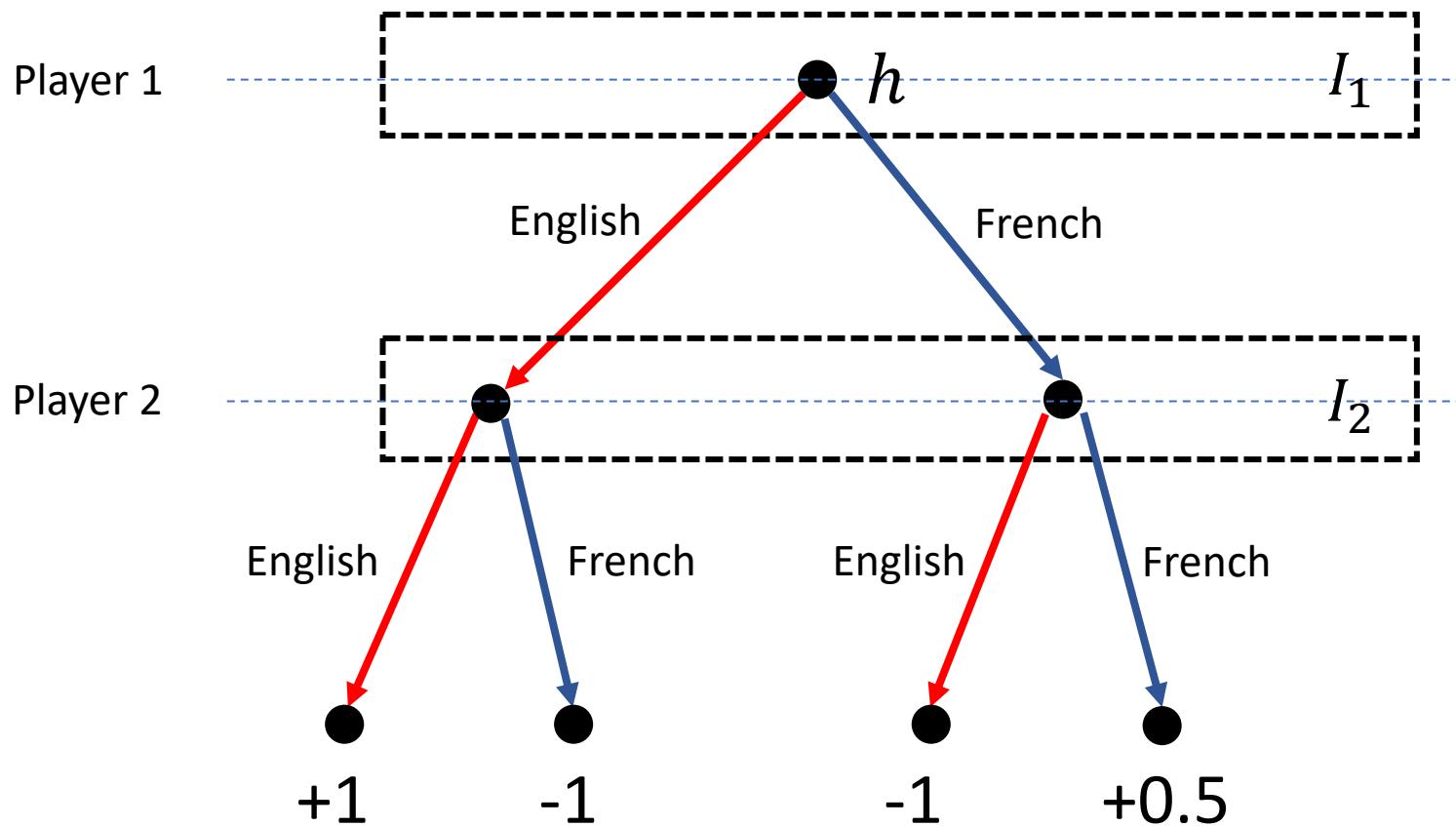
but sending signal 2 come with additional rewards?

Another example



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

Communication Game



InfoSet



Complete state (h)

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

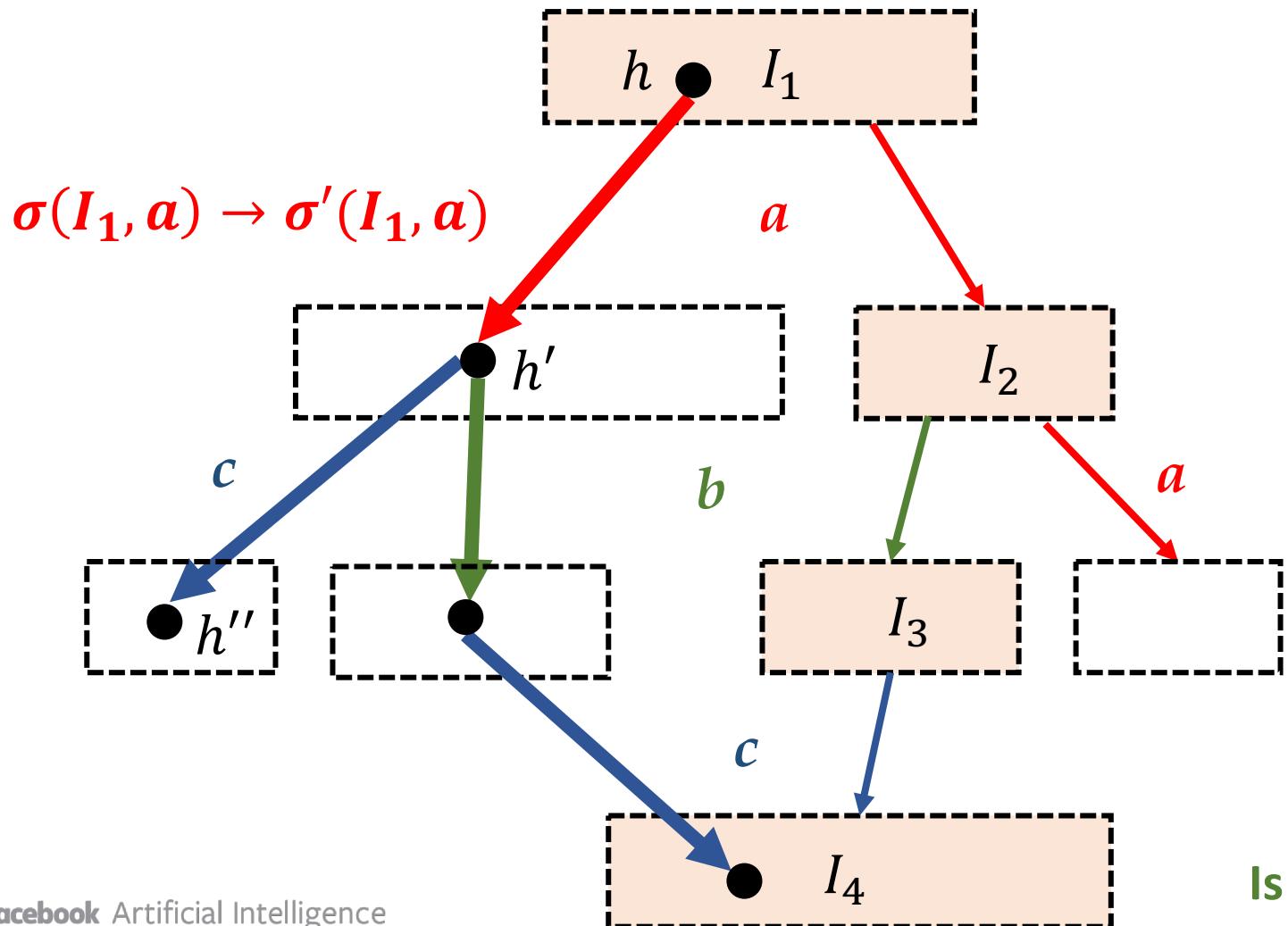
(French, French):
local Nash Equilibrium +0.5

(English, English):
global Nash Equilibrium +1.0

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

Dependency between policies

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



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.

Is there a good way to track value changes?

Policy-change Density

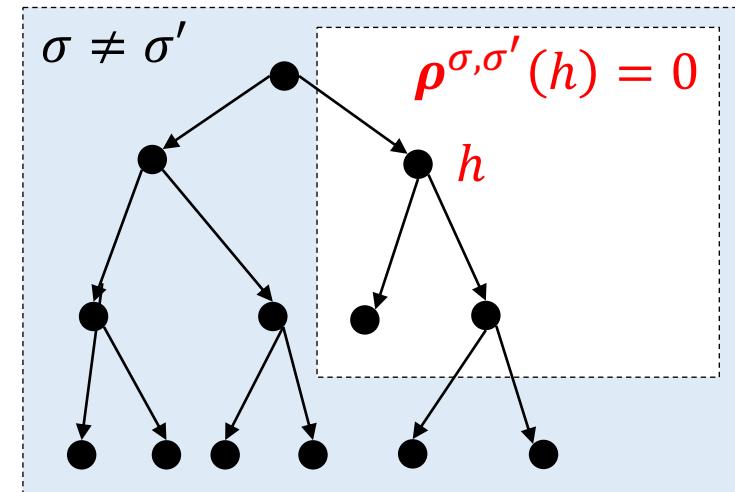
Density $\rho^{\sigma, \sigma'}(h) = \pi^{\sigma'}(h) \left[\sum_{a \in A(I)} \sigma'(I, a) v^\sigma(ha) - v^\sigma(h) \right]$

Two key properties:

- (a) Its summation yields overall value changes

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

- (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

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

Overall value changes
due to policy change

All active Infosets
 $(\sigma' \neq \sigma)$

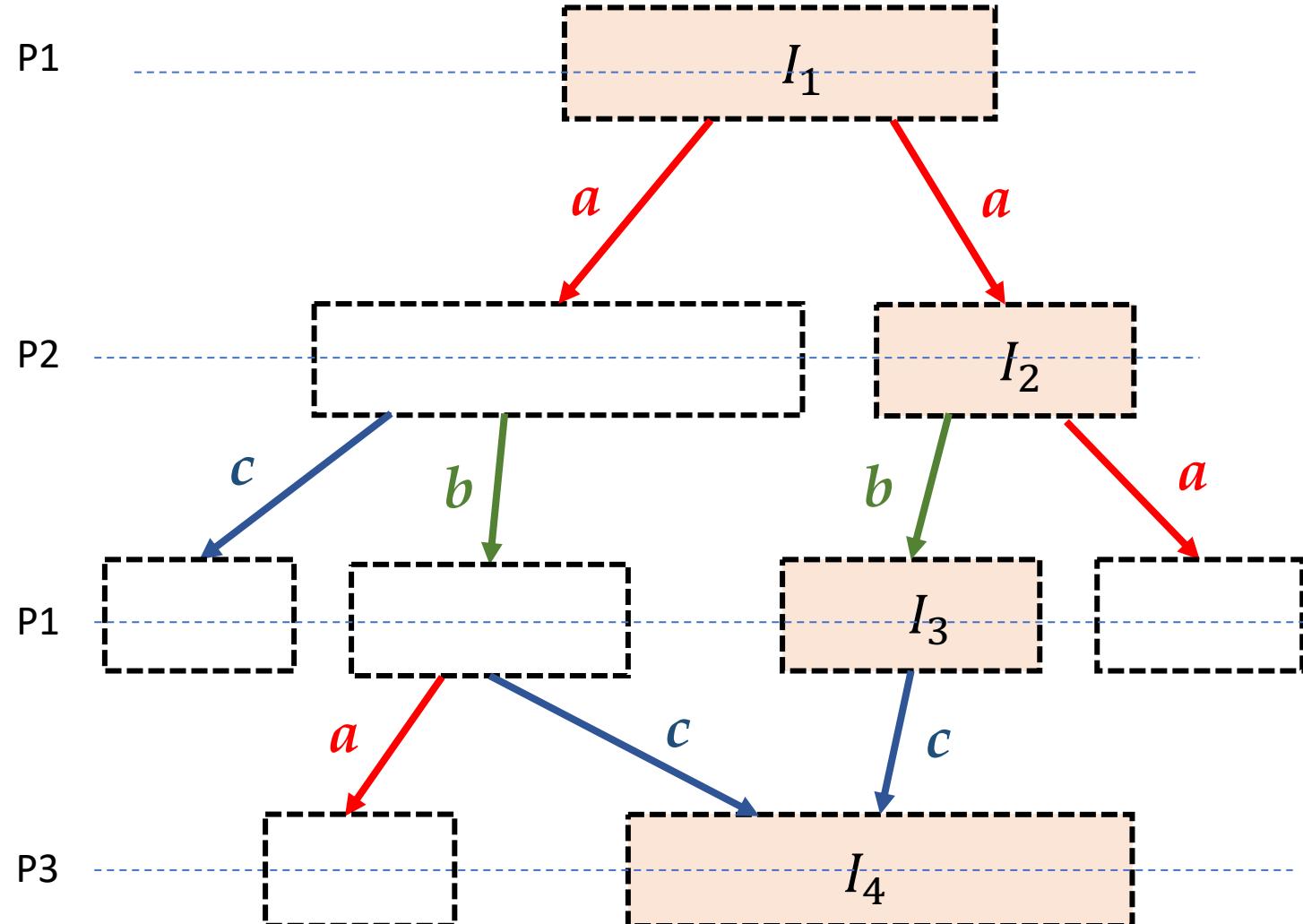
Inactive Infosets doesn't matter!!

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)					[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
CFR1k+JPS	1.00*	1.00*	1.00*	1.00*		9.50*	2.20*	5.00*	10.56*	1.07*	1.71*	2.74*
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]	1.00*	0.88	0.50	0.29		9.47*	2.23*	4.99*	9.81	0.53	0.98	1.31
Best Known	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

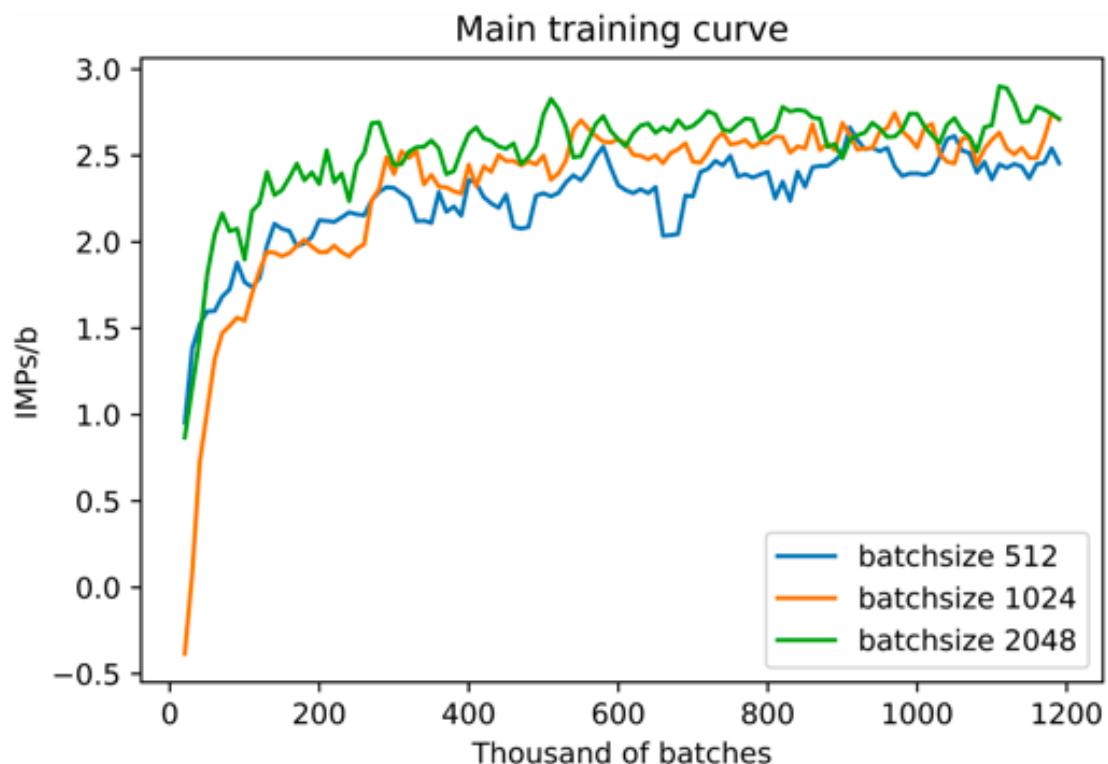
JPS can improve existing policies, and help it jump out of local optima

Contract Bridge Bidding

	N	West	North	East	South
W	♠A9743 ♥K8763 ♦A6 ♣7	2♠ ¹ Pass	2NT ² 4♣ ³ Pass	Pass	1♠ 3♣ 4NT ⁴ 7♠
E	♠Q82 ♥104 ♦QJ85432 ♣J	Pass	5♣ ⁵ Pass	Pass	
S	♠KJ1065 ♥A ♦K7 ♣A6543	Pass	Pass	Pass	

(1) Hearts and a minor. (2) Spade support, forcing to game. (3) Short clubs. (4) Keycard Blackwood. (5) Two key cards and the queen of spades, treating his fifth card as the equivalent of the queen.

- 100 years of history
- Imperfect Information
- Collaborative + Competitive
- Large State Space (5.4×10^{28})



A2C Self-play

Double-Dummy Evaluation against SoTA software

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)	+ 0.63 ± 0.20

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

BeBold: Exploration Beyond the Boundary of Explored Regions

Tianjun Zhang^{1,4}



Huazhe Xu^{1,4}



Xiaolong Wang^{1,2}



Yi Wu³



Kurt Keutzer¹



Joseph E. Gonzalez¹

Yuandong Tian⁴



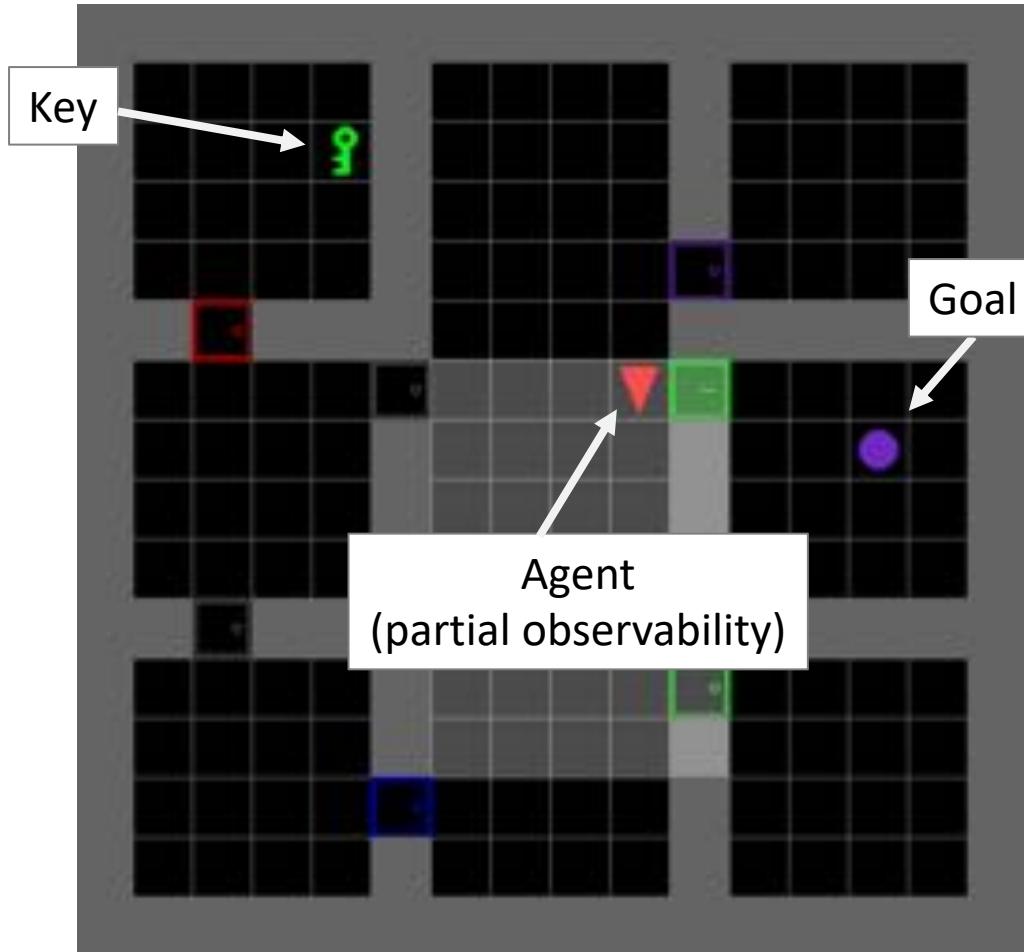
¹UC Berkeley

²UCSD

³Tsinghua University

⁴FaceBook AI Research

Environment with Sparse Reward



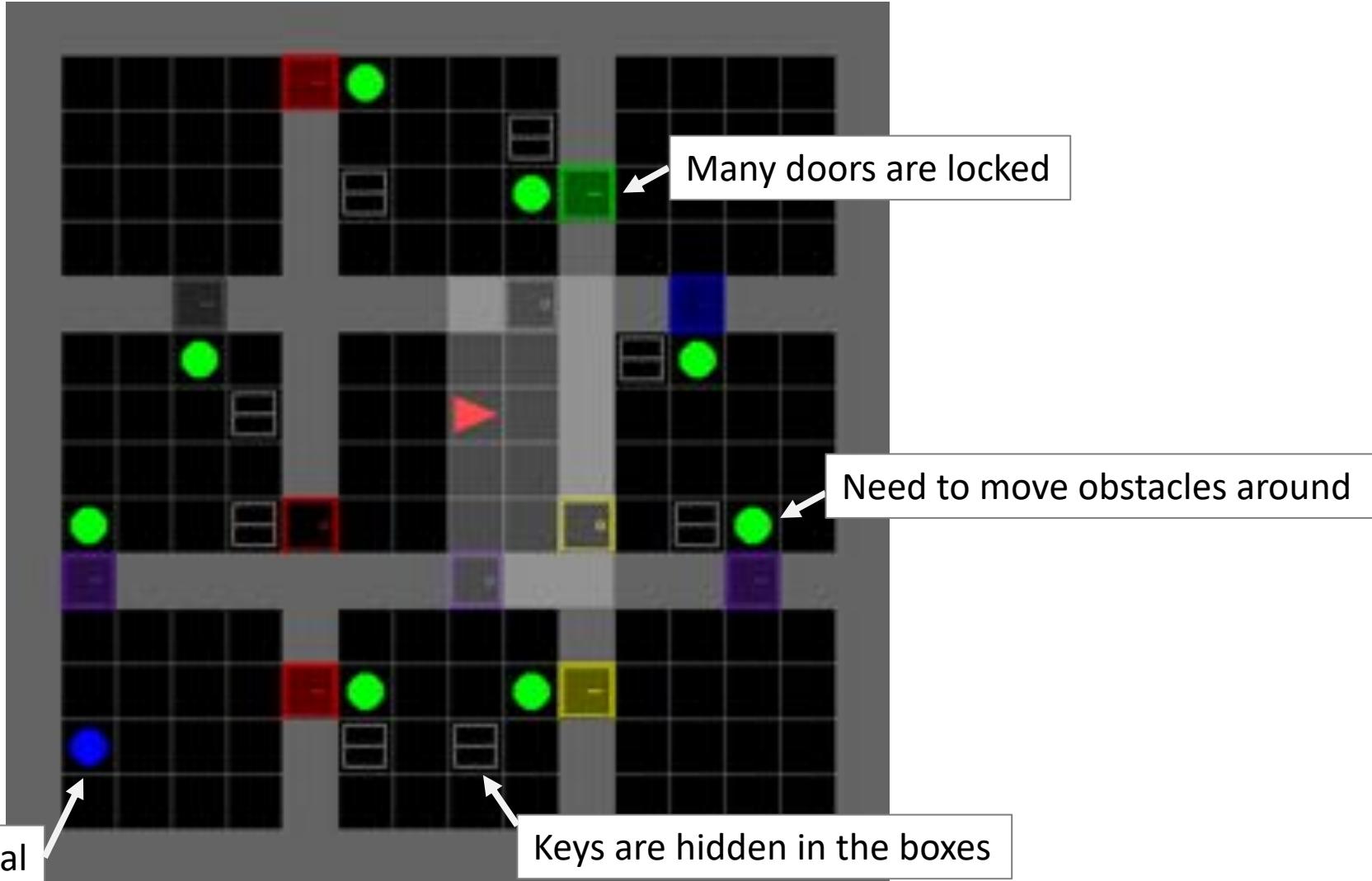
No external reward

when agent wonders around.
when agent picks the key
when agent opens all doors
when agent opens the locked door

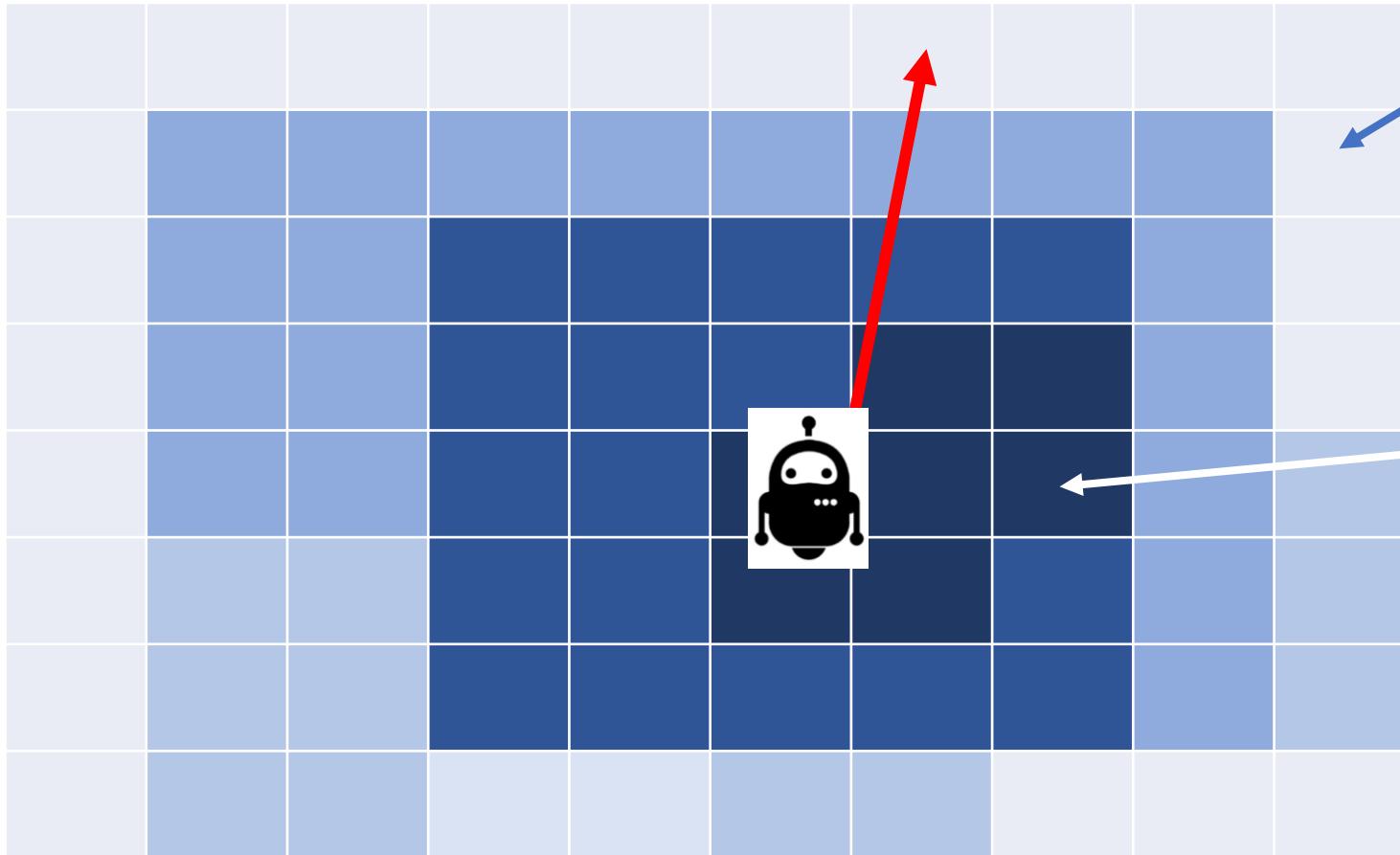
...

until the agent reaches the goal

And more complicated situations...



Count-based Exploration



Low visitation counts $N(s)$
High intrinsic rewards

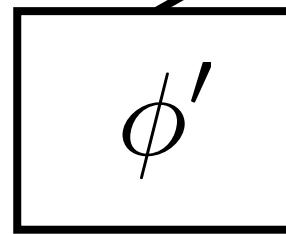
High visitation counts $N(s)$
Low Intrinsic reward

What if we have exponential #states?

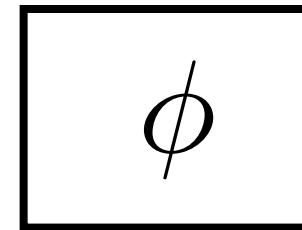
Random Network Distillations (RND)

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

Student Network
(learning from teacher)



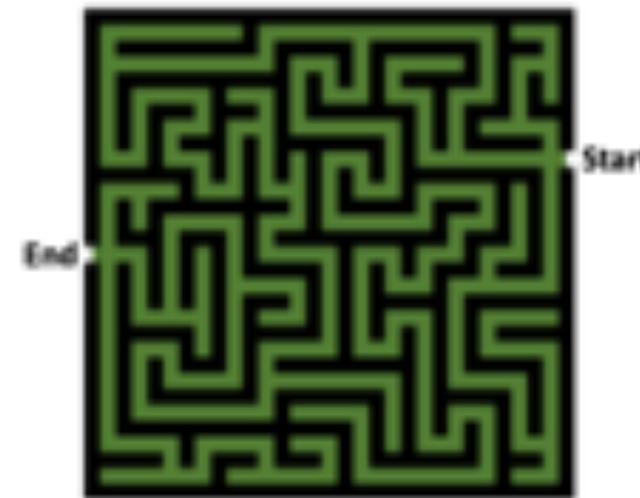
Random fixed
teacher network



s

Familiar state = low prediction error

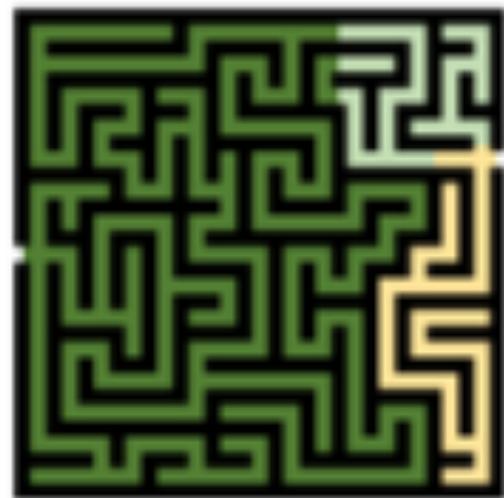
Issues in RND



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



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

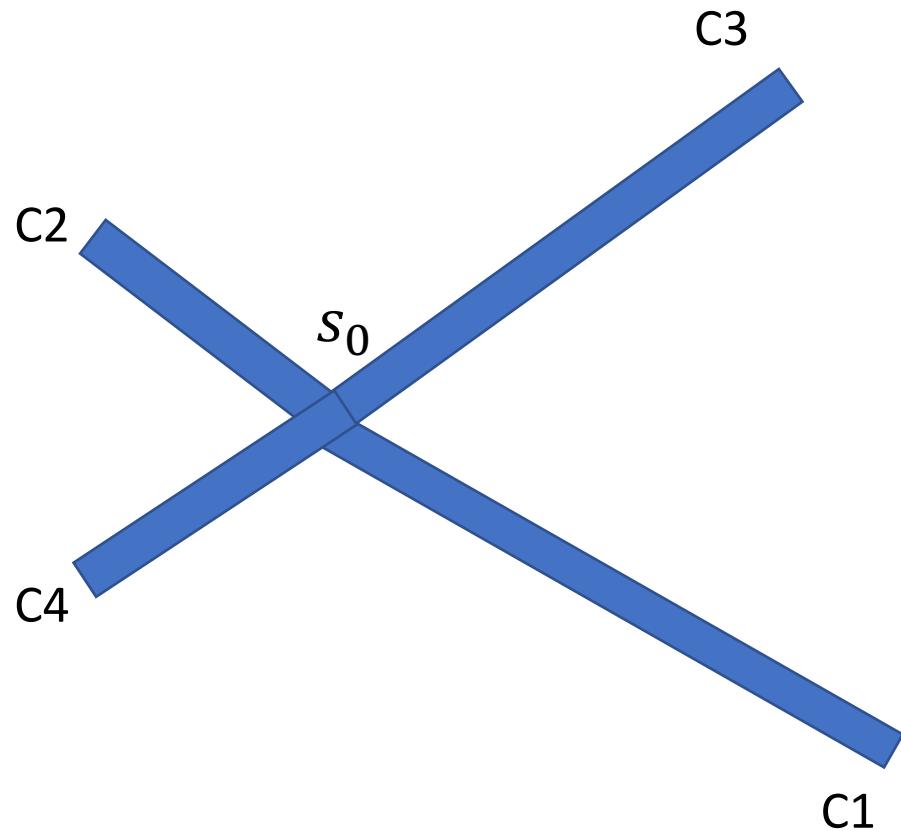


3. RND by chance starts
exploring the bottom right
corner heavily, resulting in
the IR at top right higher
than bottom right



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

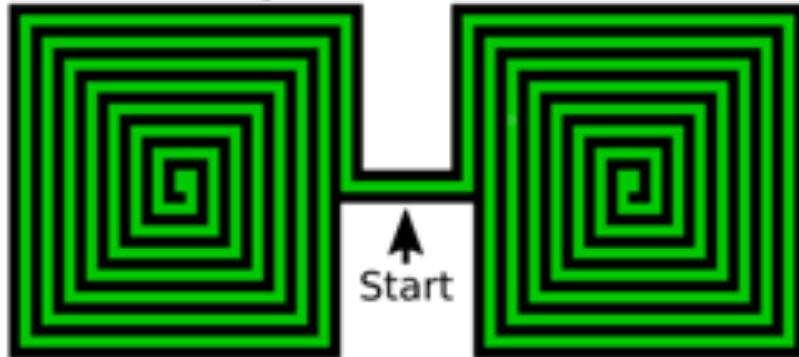
Multi-Corridor Problem



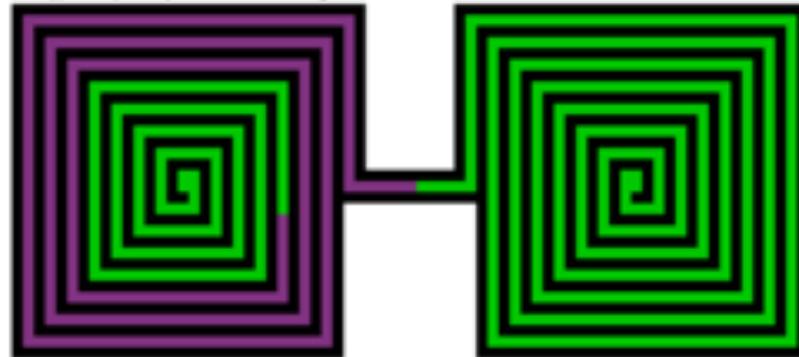
	C1	C2	C3	C4	Entropy
Length	40	10	30	10	-
Count-Based	68K	2K	3K	36K	1.23
BeBold Tabular	27K	51K	13K	17K	1.80
RND	26K	38K	14K	32K	1.92
BeBold	27K	24K	33K	27K	1.99

Detachment in RND

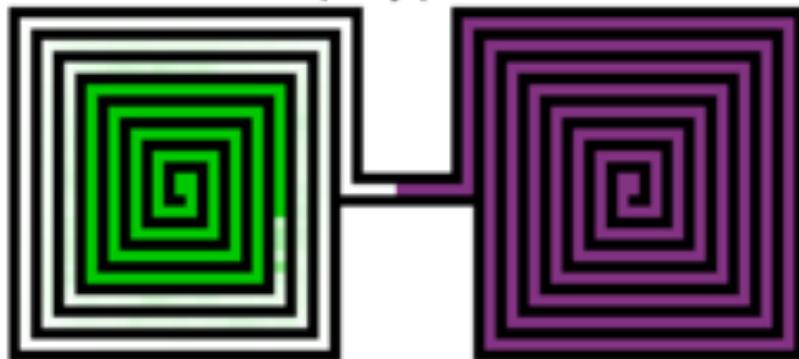
1. Intrinsic reward (green) is distributed throughout the environment



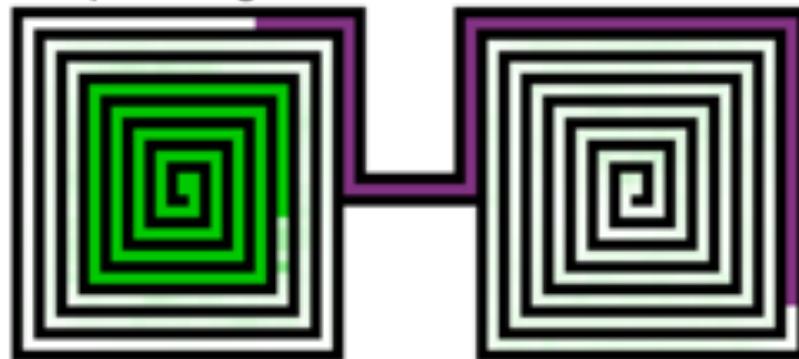
2. An IM algorithm might start by exploring (purple) a nearby area with intrinsic reward



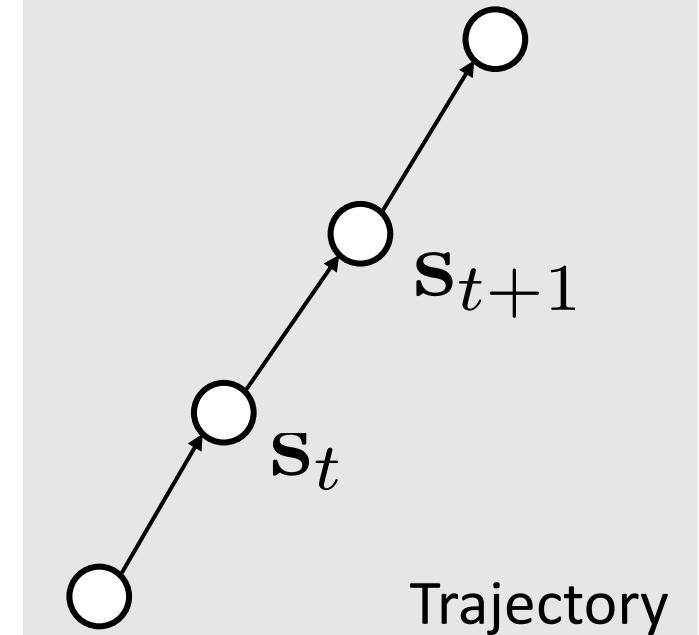
3. By chance, it may explore another equally profitable area



4. Exploration fails to rediscover promising areas it has detached from



BeBold



$$\underline{r^i(\mathbf{s}_t, \mathbf{a}_t)} = \max \left(\frac{1}{\underline{N(\mathbf{s}_{t+1})}} - \frac{1}{\underline{N(\mathbf{s}_t)}}, 0 \right) * \mathbb{1}\{\underline{N_e(\mathbf{s}_{t+1})} = 1\}$$

Annotations for the equation:

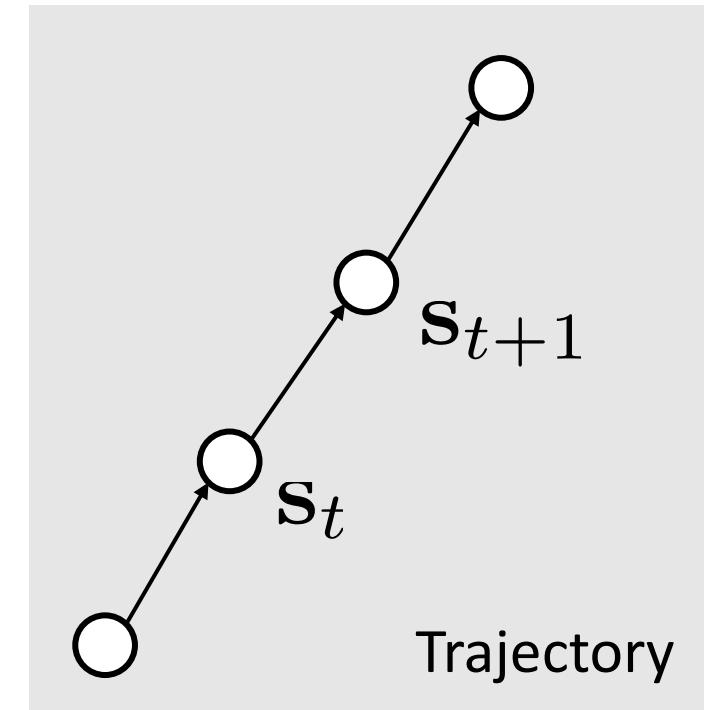
- Intrinsic Reward**: Points to the underlined term $r^i(\mathbf{s}_t, \mathbf{a}_t)$.
- Inverse of visitation counts**: Points to the terms $\frac{1}{N(\mathbf{s}_{t+1})}$ and $\frac{1}{N(\mathbf{s}_t)}$.
- Episodic visitation count**: Points to the underlined term $N_e(\mathbf{s}_{t+1})$.

BeBold (Beyond the Boundary of Explored Regions)

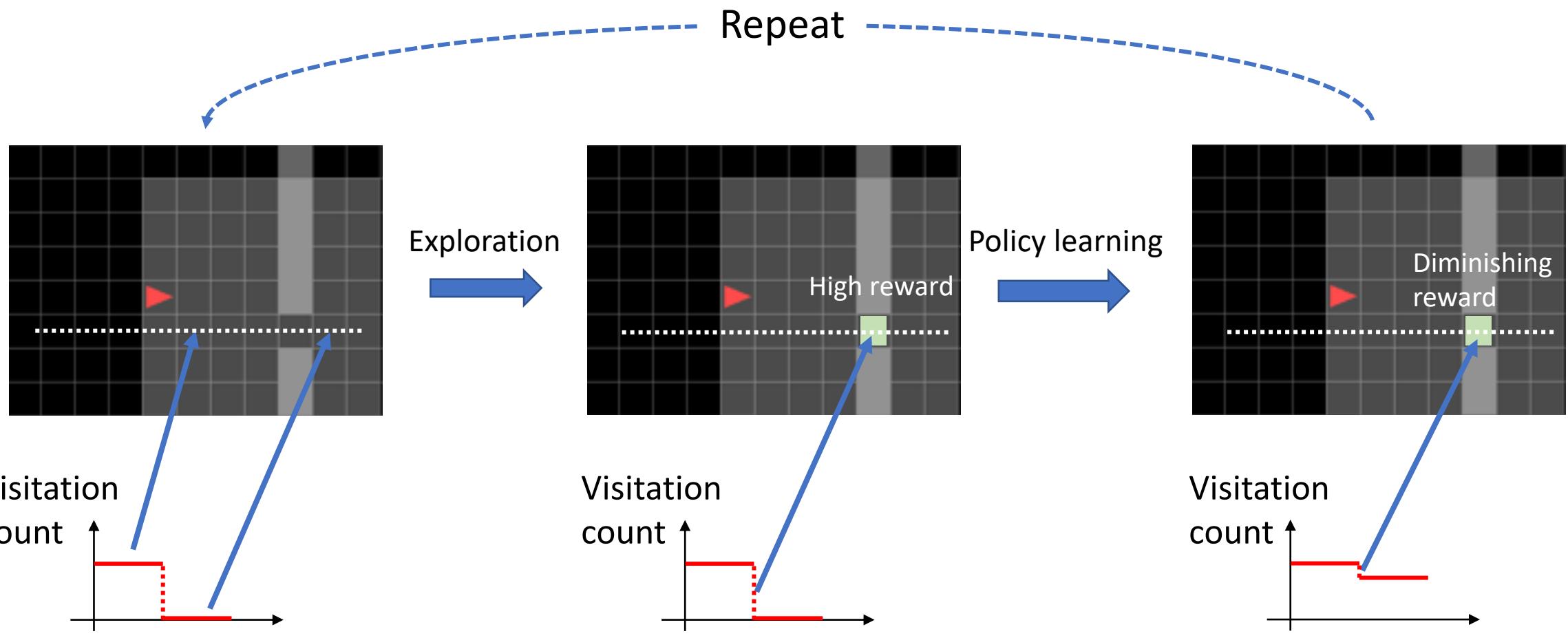
$$\underline{r^i(\mathbf{s}_t, \mathbf{a}_t)} = \max \left(\frac{1}{N(\mathbf{s}_{t+1})} - \frac{1}{N(\mathbf{s}_t)}, 0 \right)$$

Intrinsic Reward

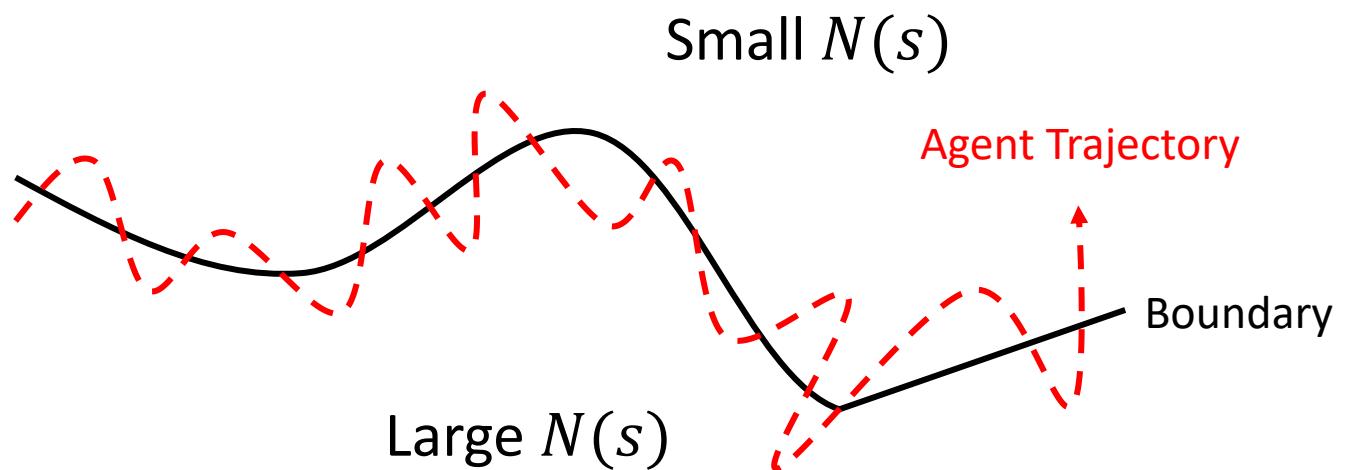
Inverse of visitation counts



BeBold (Beyond the Boundary of Explored Regions)



BeBold



$$\underline{r^i(\mathbf{s}_t, \mathbf{a}_t)} = \max \left(\frac{1}{N(\mathbf{s}_{t+1})} - \frac{1}{N(\mathbf{s}_t)}, 0 \right) * \mathbb{1}\{\underline{N_e(\mathbf{s}_{t+1})} = 1\}$$

Intrinsic Reward

Episodic
visitation count

BeBold (final form)

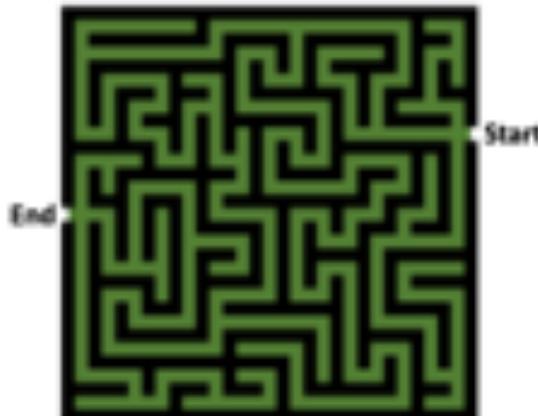
$$r^i(\mathbf{s}_t, \mathbf{a}_t) = \max(||\phi(\mathbf{o}_{t+1}) - \phi'(\mathbf{o}_{t+1})||_2 - ||\phi(\mathbf{o}_t) - \phi'(\mathbf{o}_t)||_2, 0) * \mathbb{1}\{\underline{N_e(\mathbf{o}_{t+1}) = 1}\}$$

RND for t+1 **RND for t**

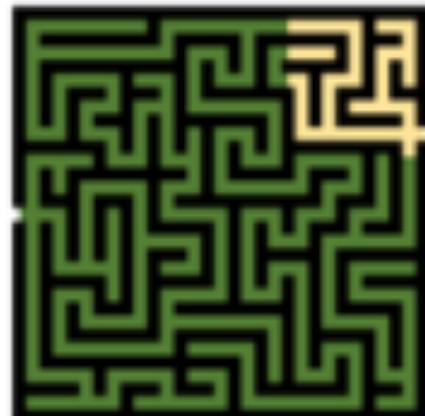
Observation (rather than full state)

Hash Table

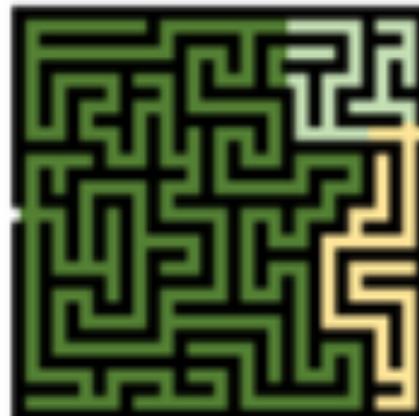
RND



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

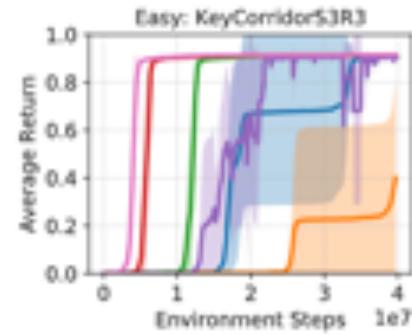
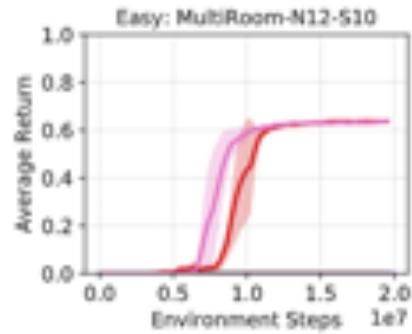
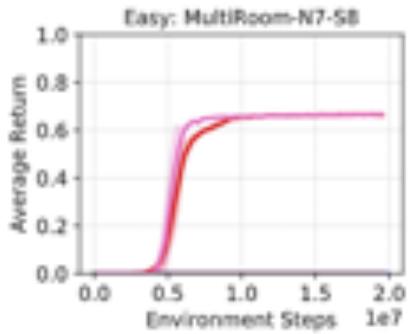
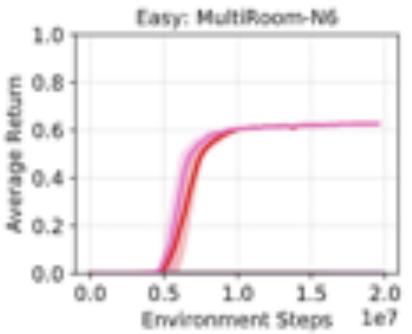
*MR is short for MultiRoom, KC is for KeyCorridor, OM is for ObstructedMaze

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

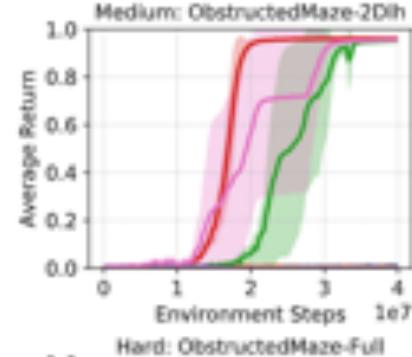
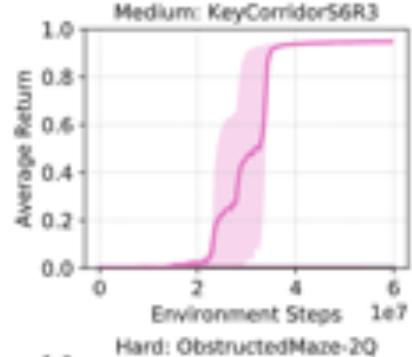
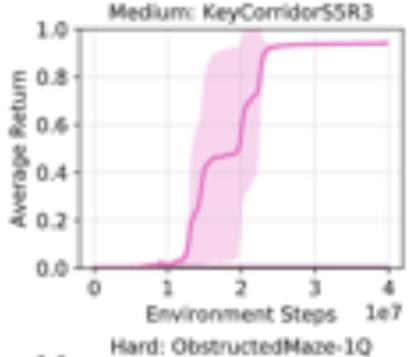
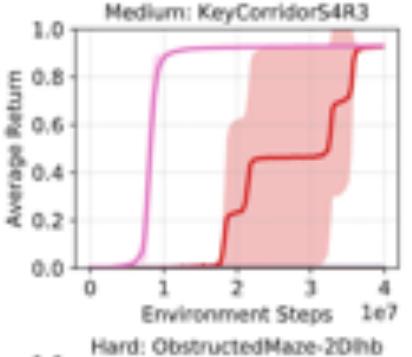
MiniGri

d

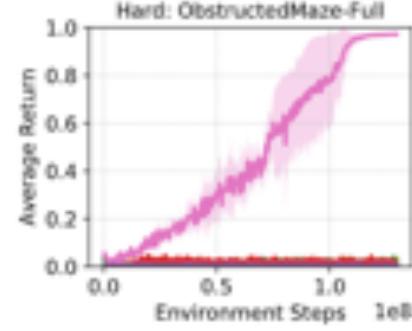
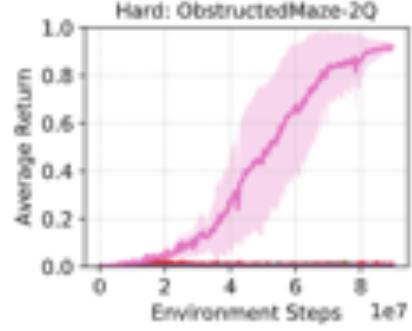
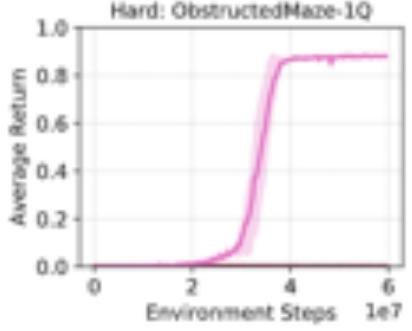
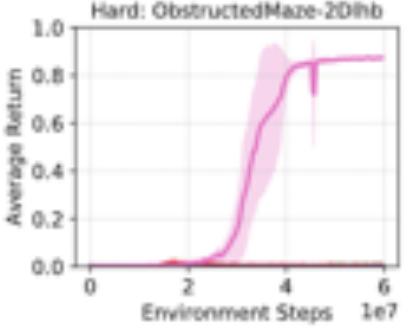
Easy



Medium



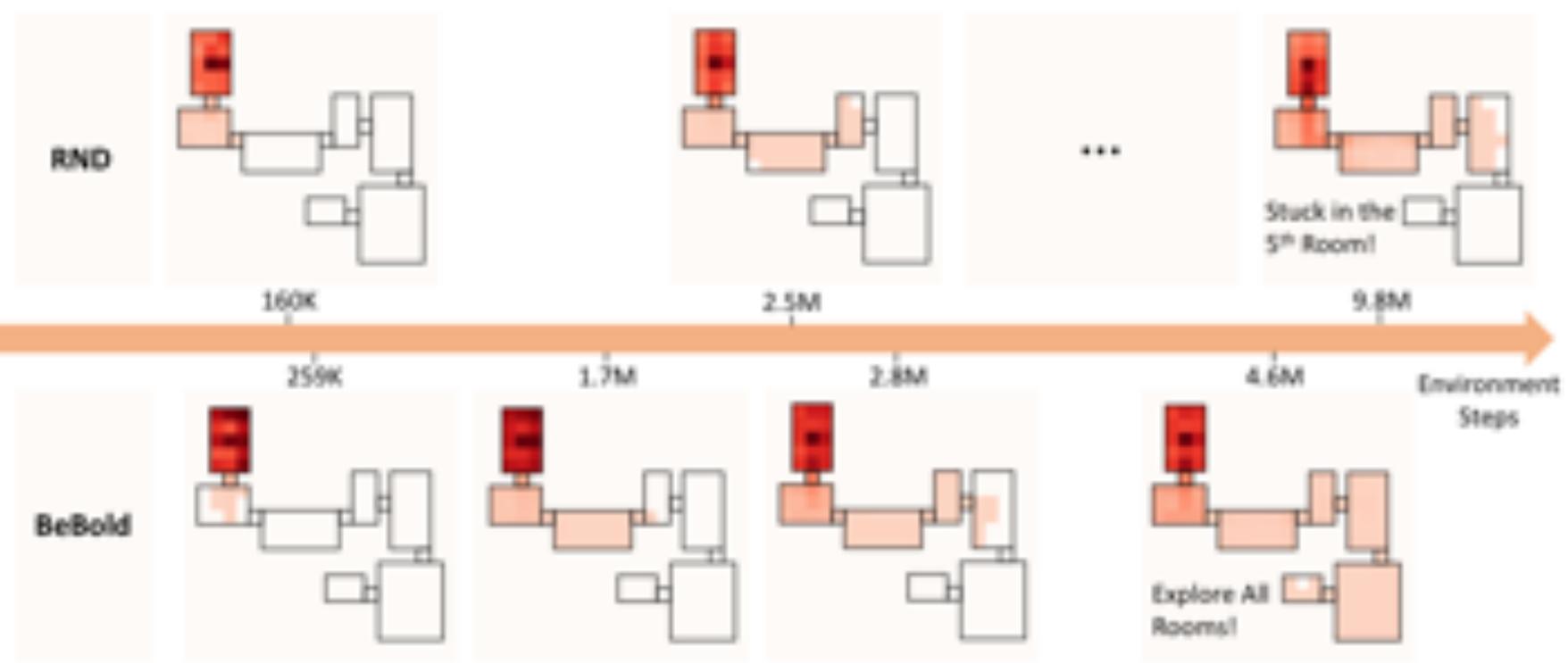
Hard



AMIGO: [Campero, Andres, et al. "Learning with AMIGO: Adversarially Motivated Intrinsic Goals." arXiv preprint arXiv:2006.12122 (2020)]

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

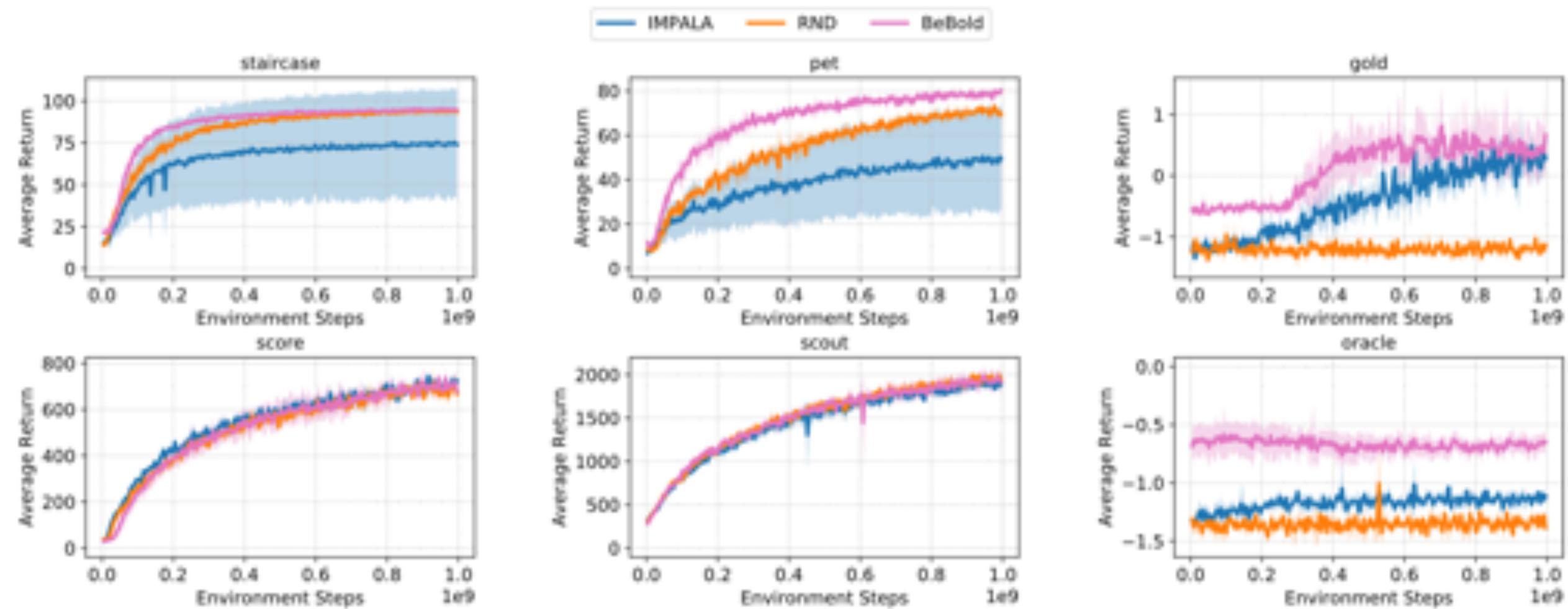
Pure Exploration



NetHack



6 Tasks in NetHack



Learning Action Space in Monte Carlo Tree Search



Linnan Wang¹



Saining Xie²



Teng Li²



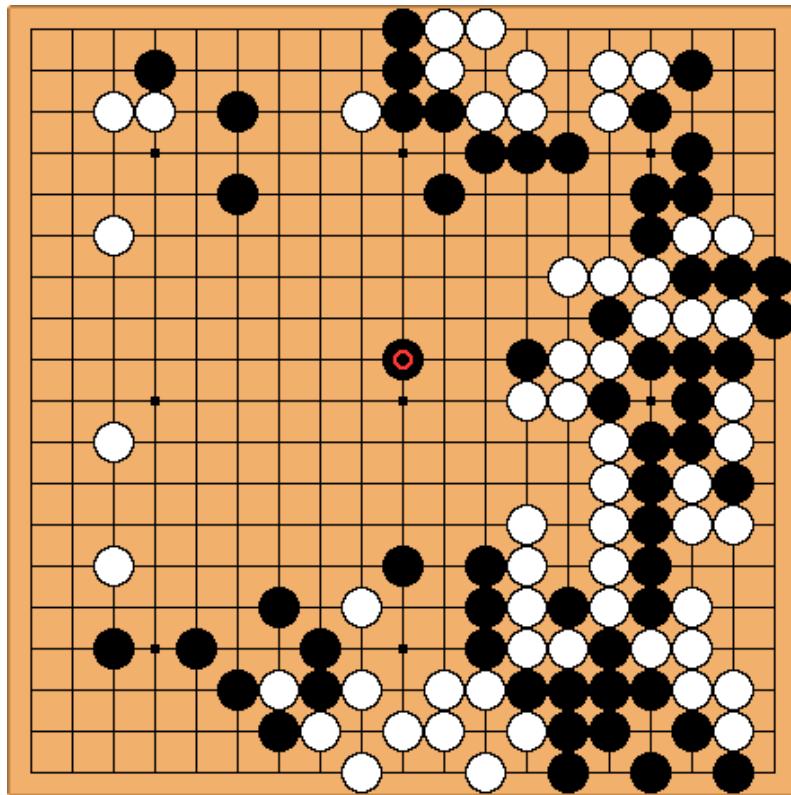
Rodrigo Fonseca¹



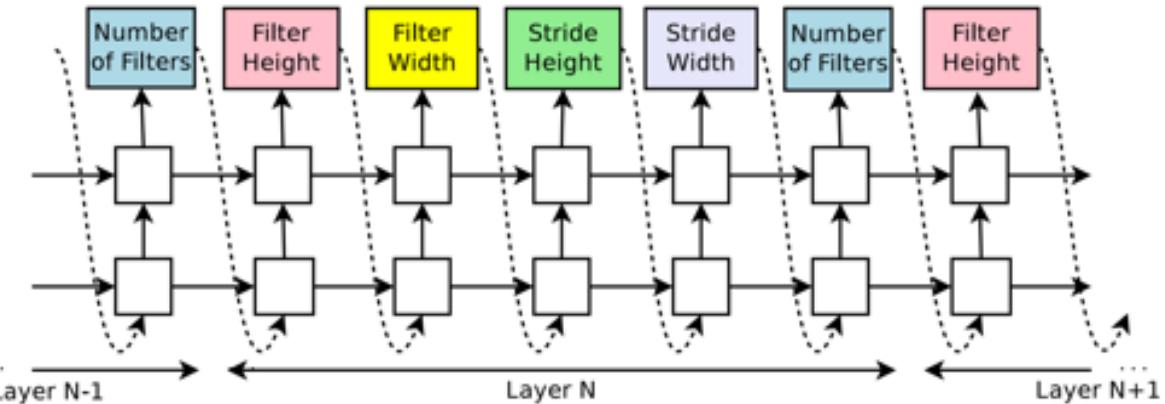
Yuandong Tian²

¹Brown University, ²Facebook AI Research

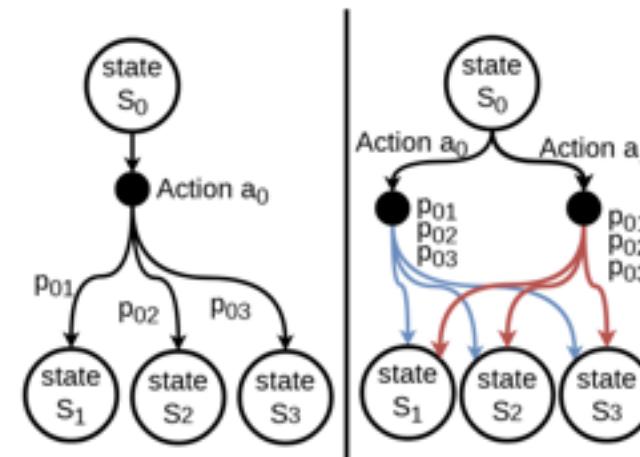
Predefined Action Space



Fixed action space = R^{361}

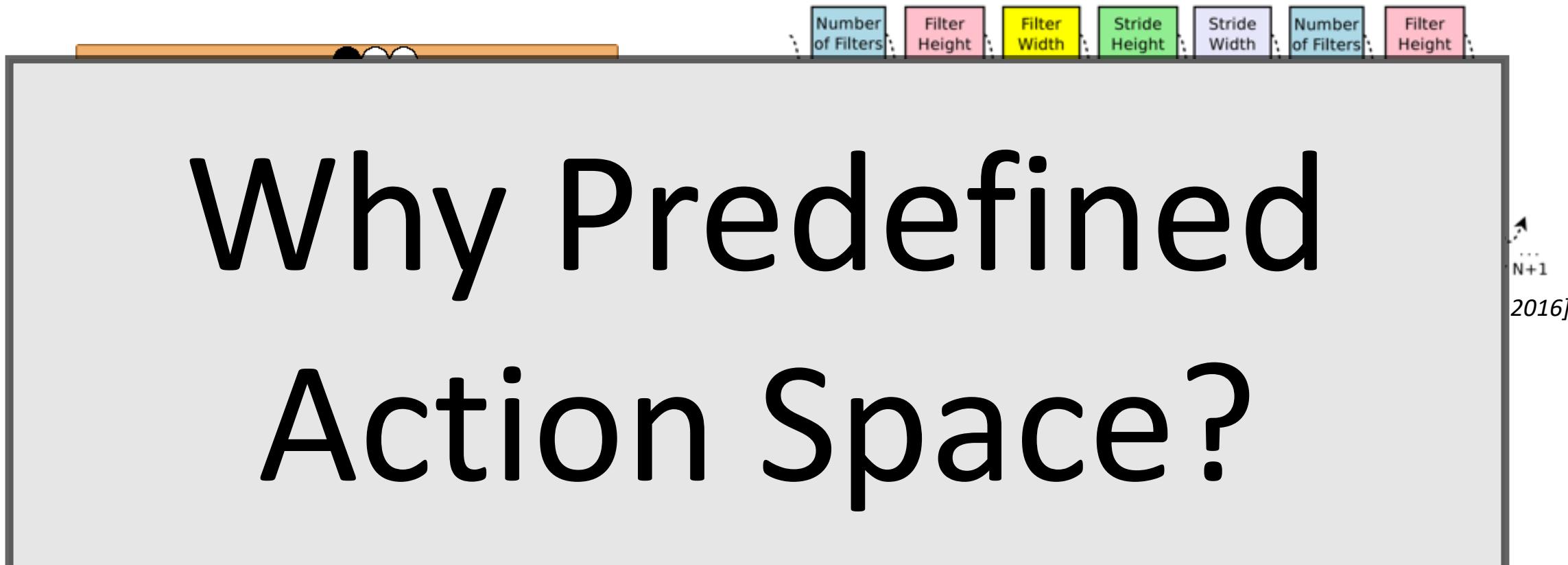


[B. Zoph and Q. Le, *Neural Architecture Search with Reinforcement Learning*, 2016]

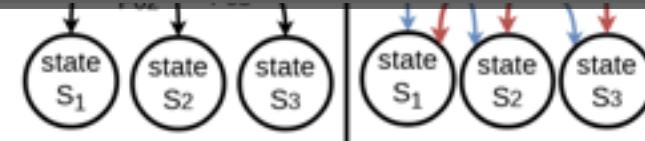


[G. Malazgirt, *TauRieL: Targeting Traveling Salesman Problem with a deep reinforcement learning inspired architecture*]

Predefined Action Space

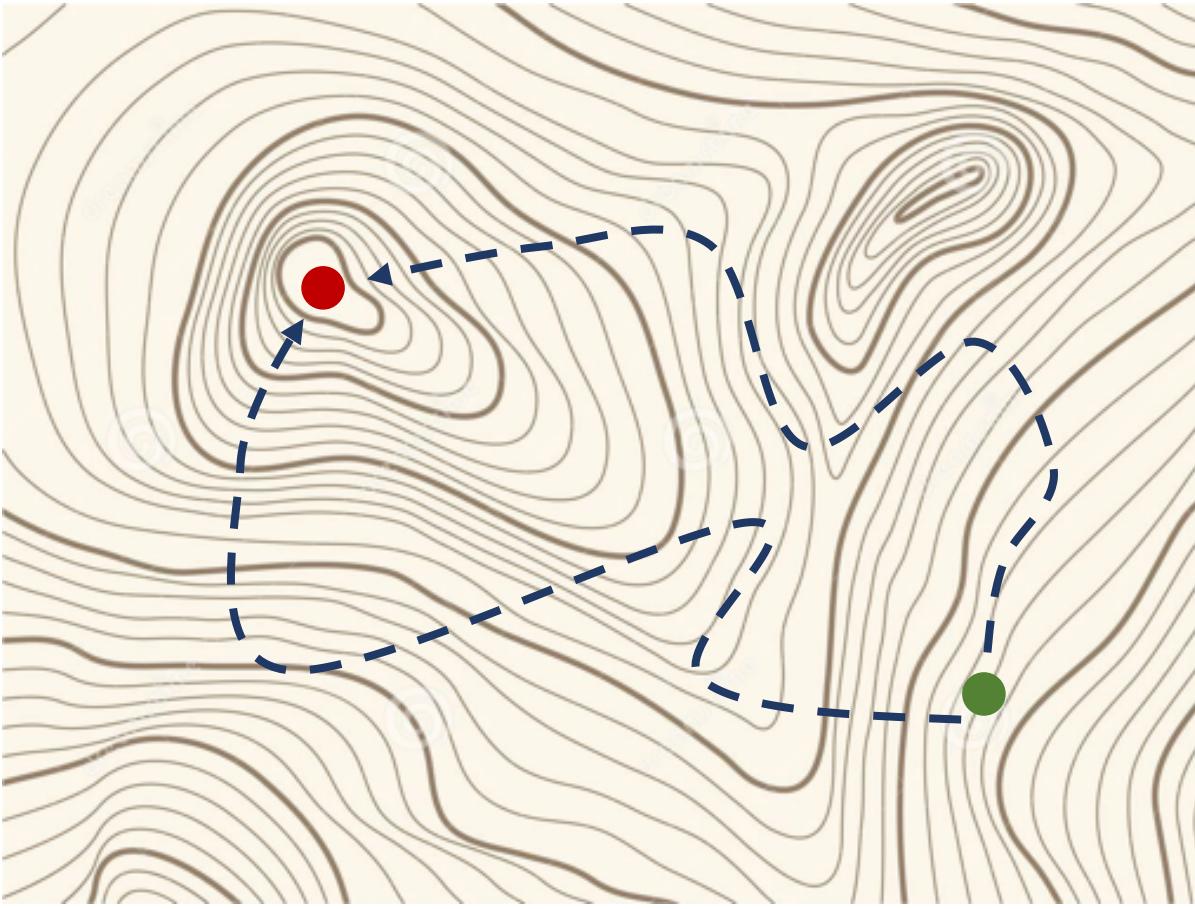


Fixed action space = R^{361}



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

Why Predefined Action Space?



We only care the final solution

We don't care how we get it.

Motivating Examples

Depth = {1, 2, 3, 4, 5}

Channels = {32, 64}

KernelSize = {3x3, 5x5}

1364 networks.

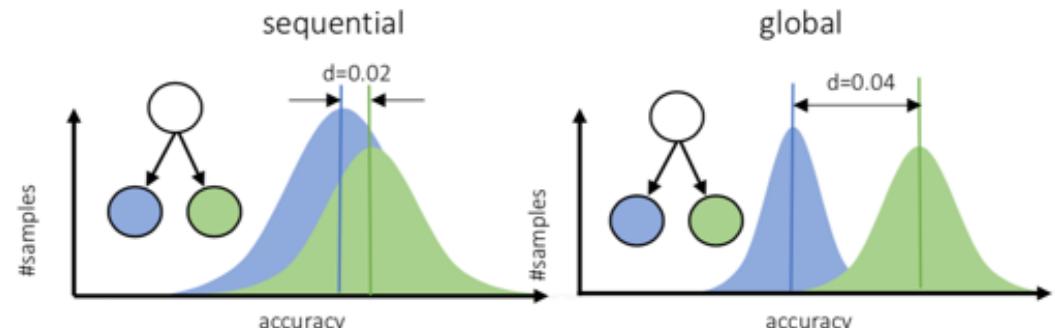
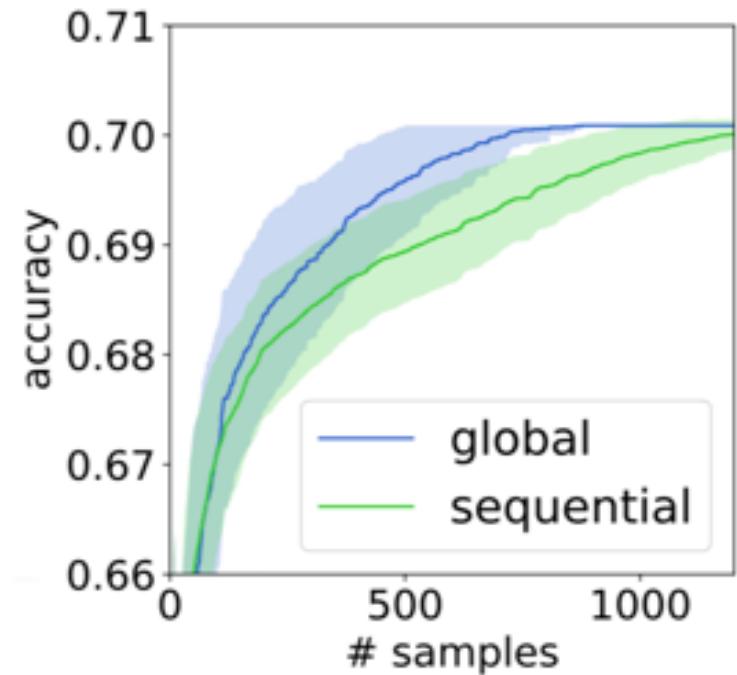
Goal: Find the network
with the best accuracy using fewest trials.

Action space

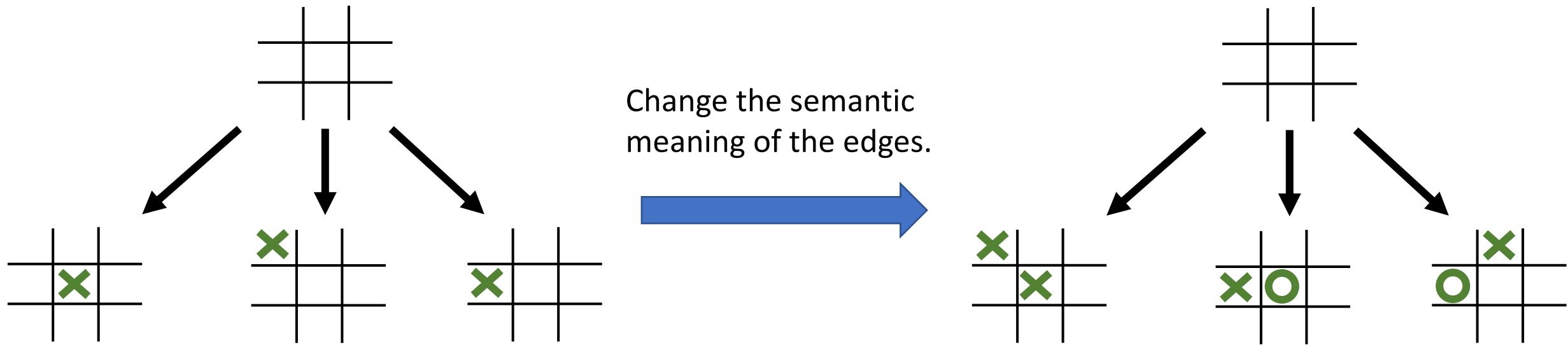
Sequential = { add a layer, set K, set C }

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

Global is better!

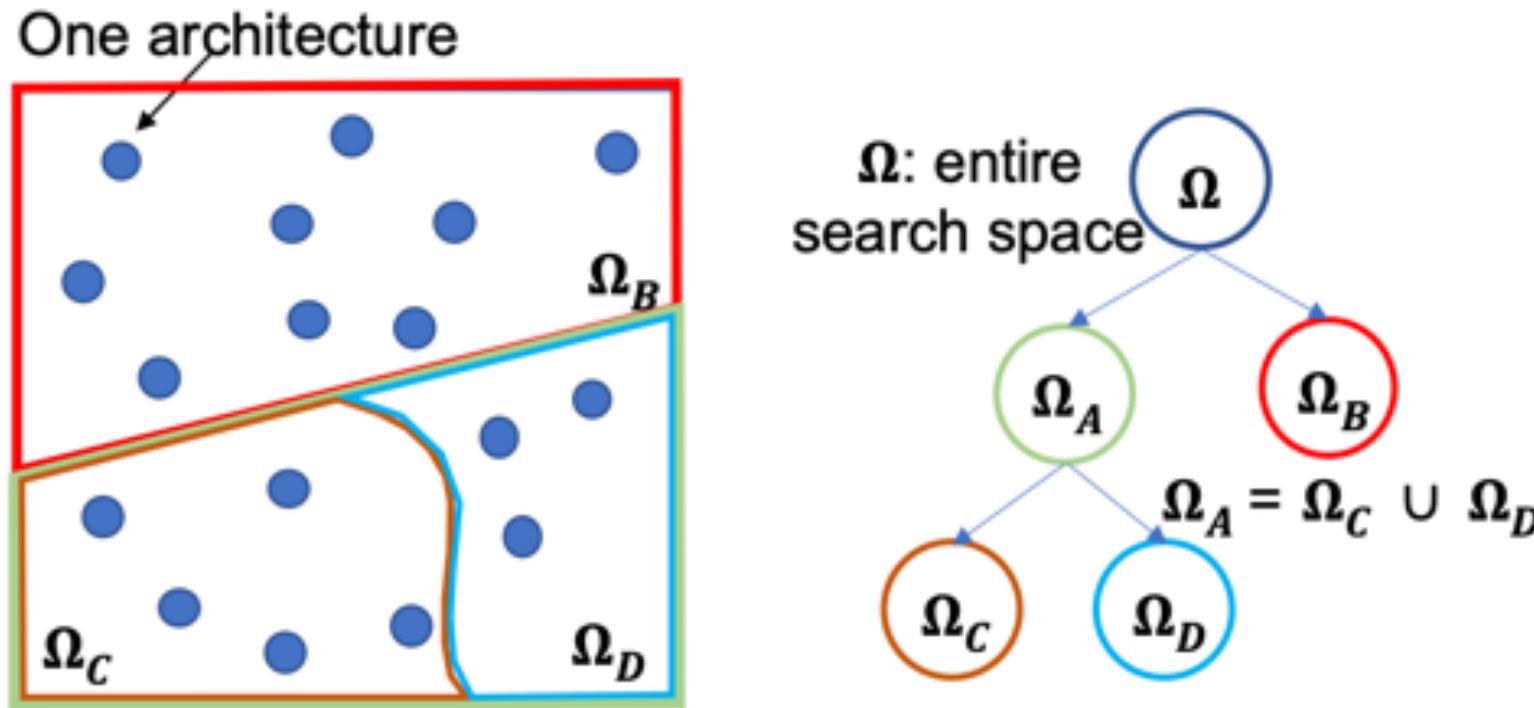


The Meaning of Learning Action Space



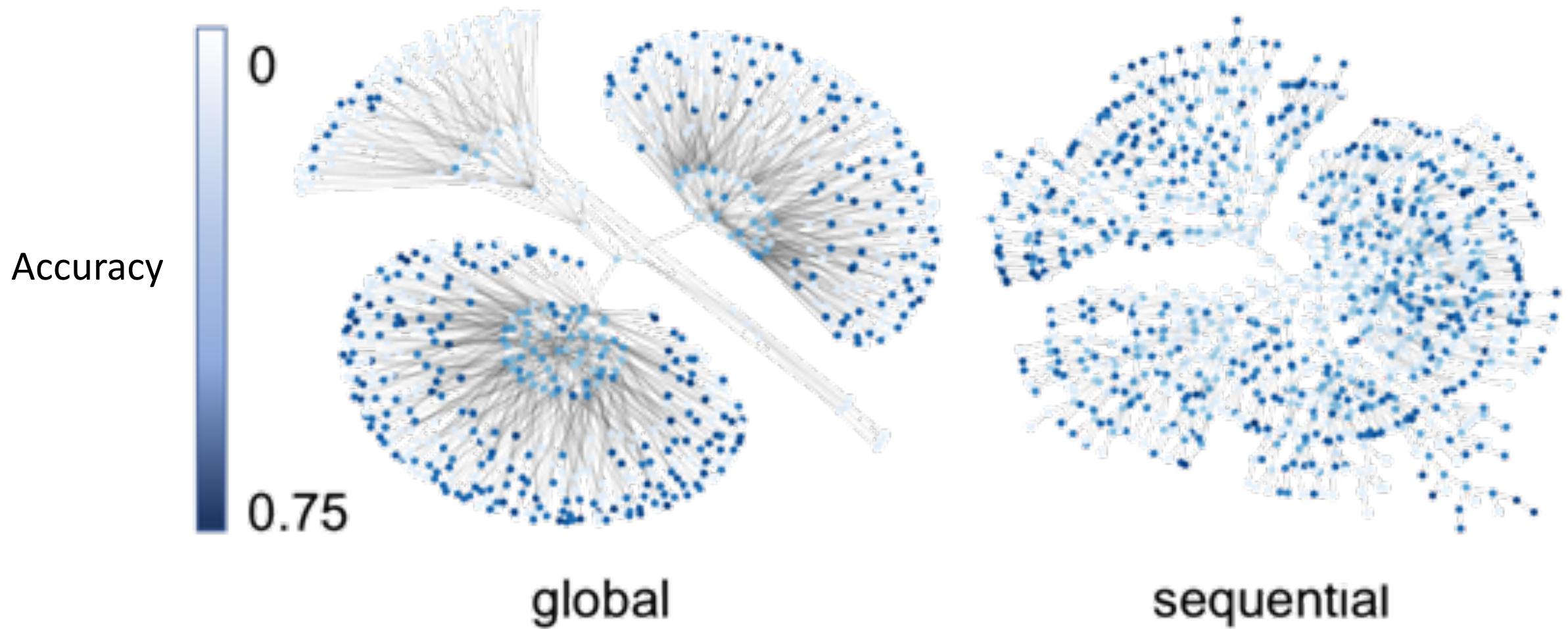
Not allowed in games, but doable in optimization.

How to learn the action space?

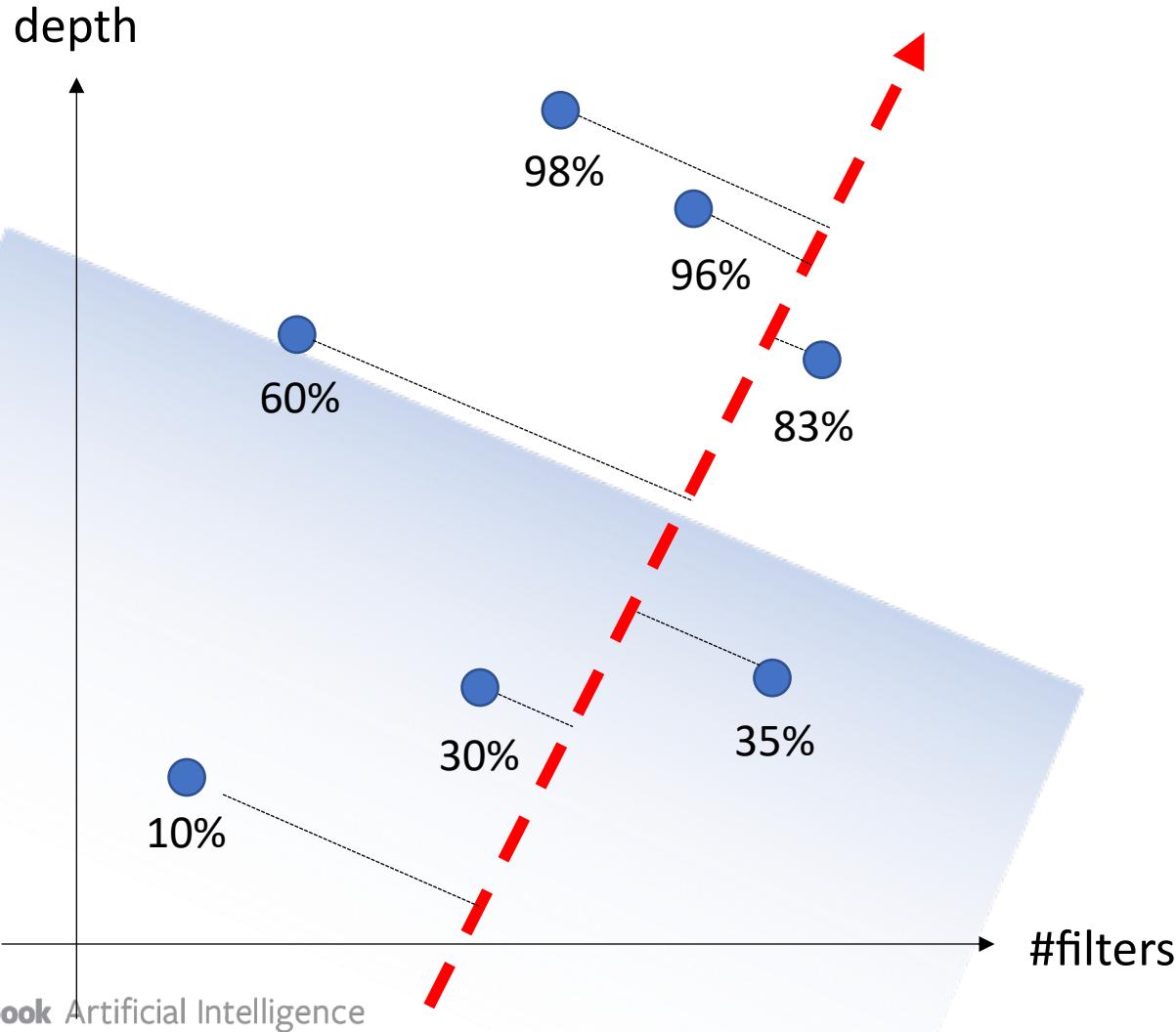


Partition = Action

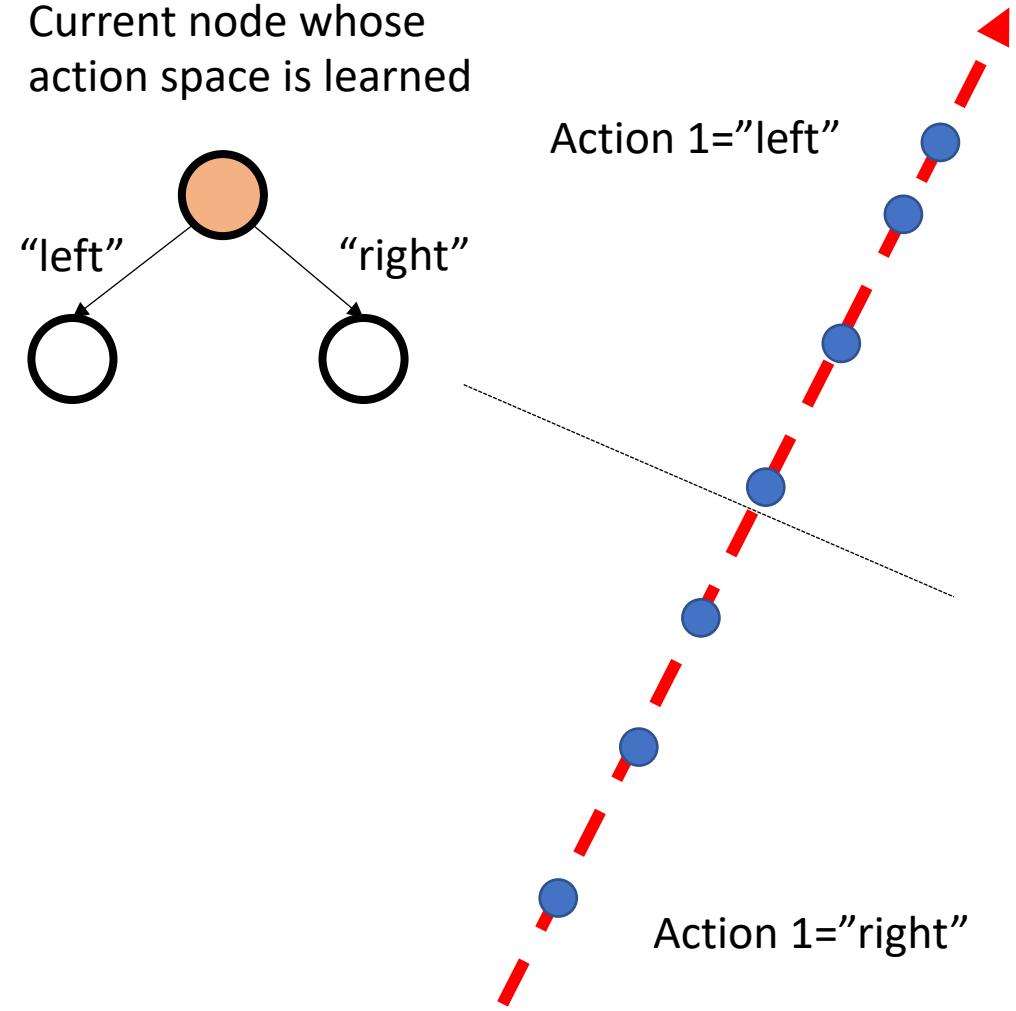
Different Partition → Different Value Distribution



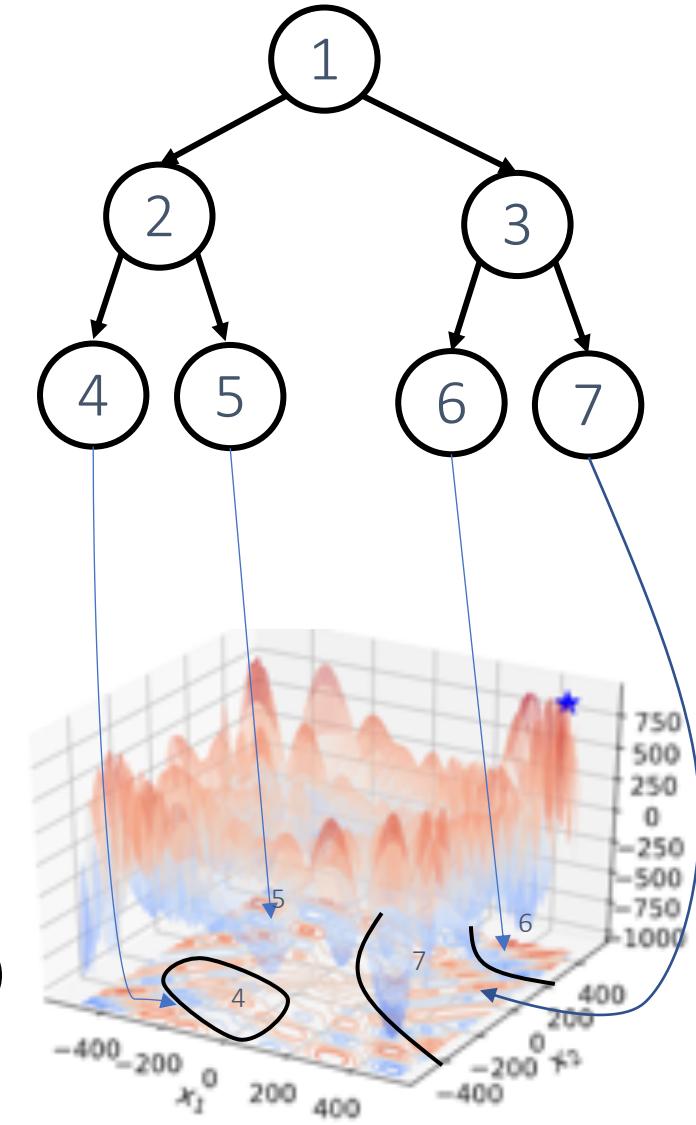
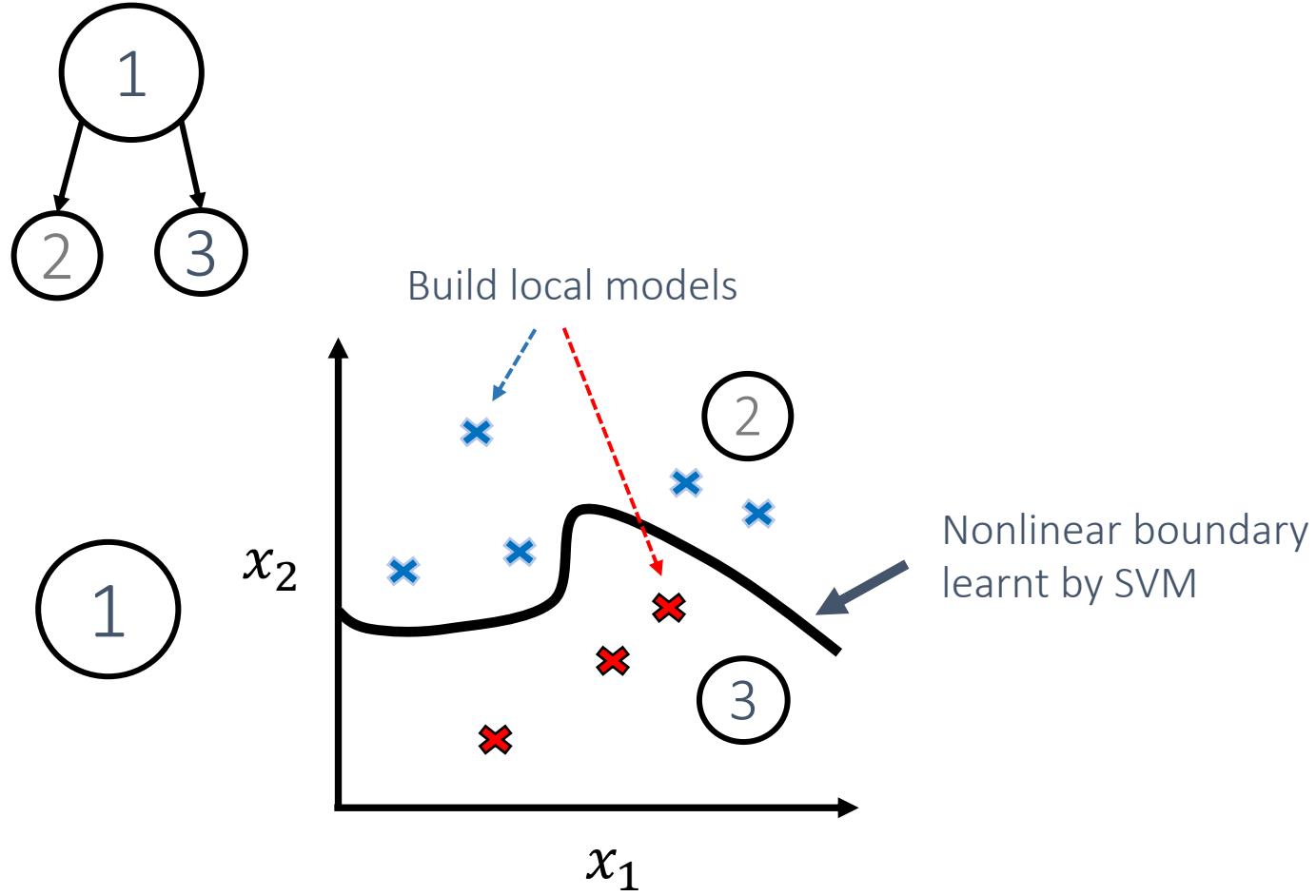
Learn action space



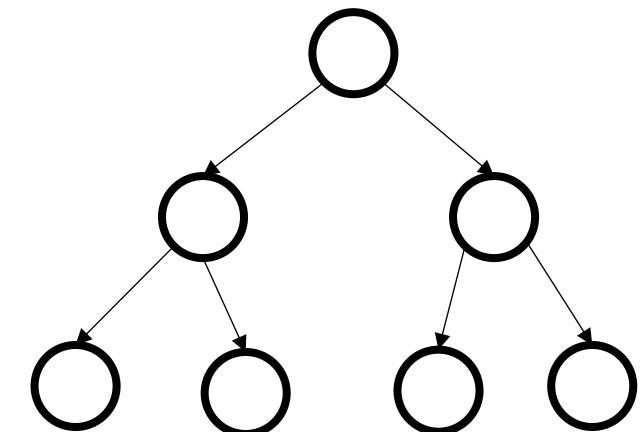
Current node whose
action space is learned



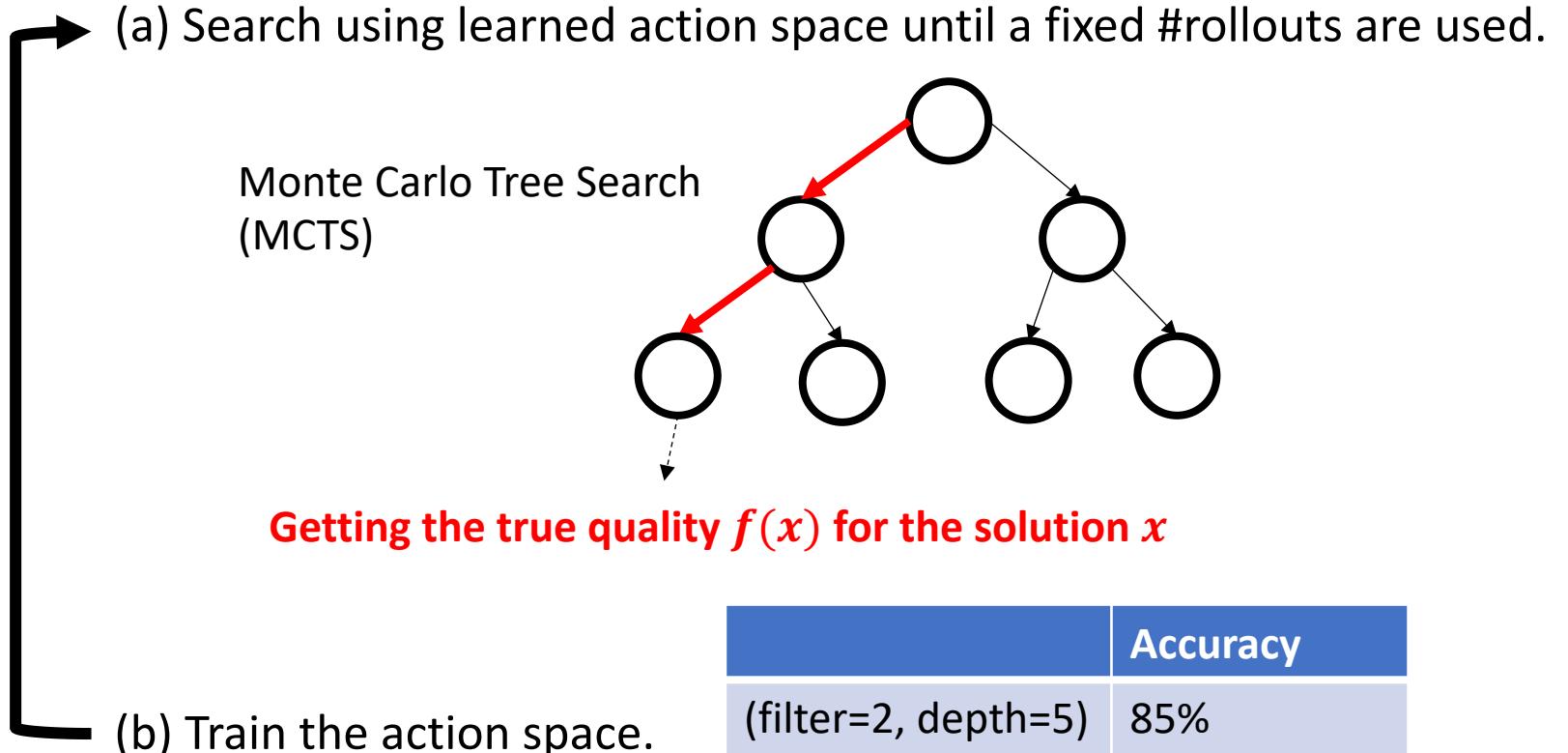
Nonlinear Partition



Approach



Fixed action branches
(but not action space)



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



Linnan Wang¹



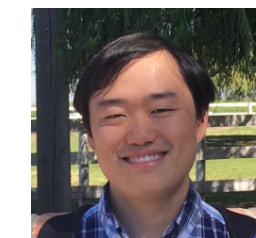
Saining Xie²



Teng Li²



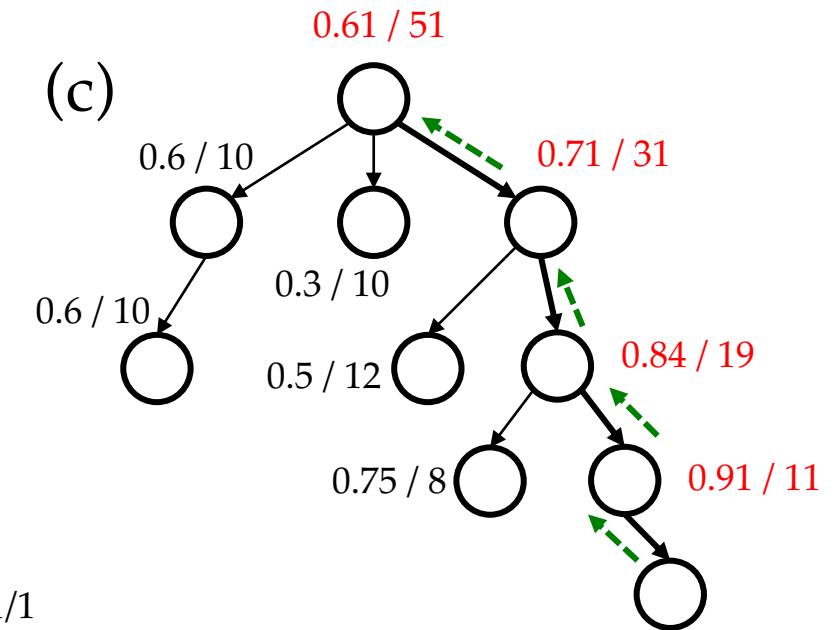
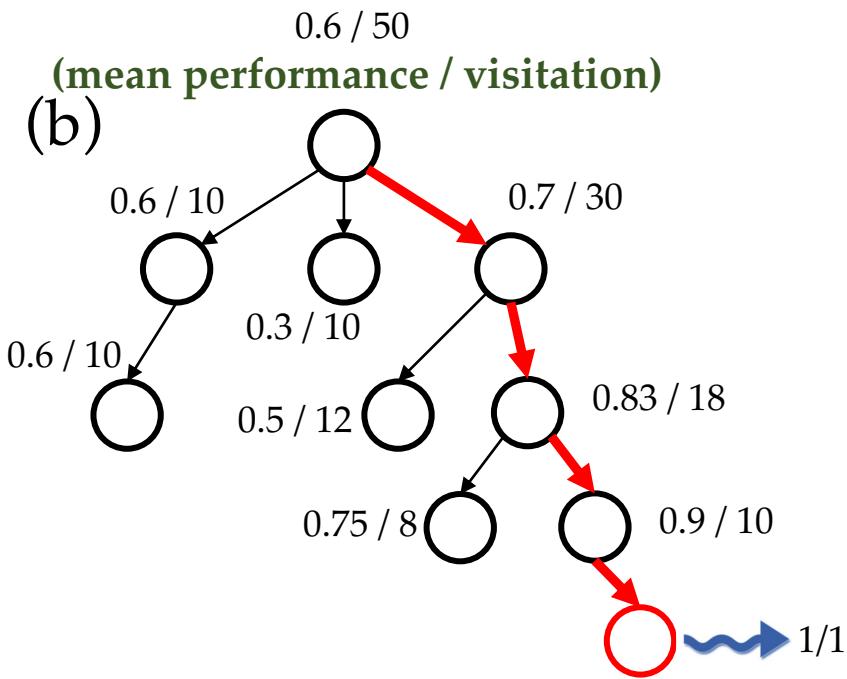
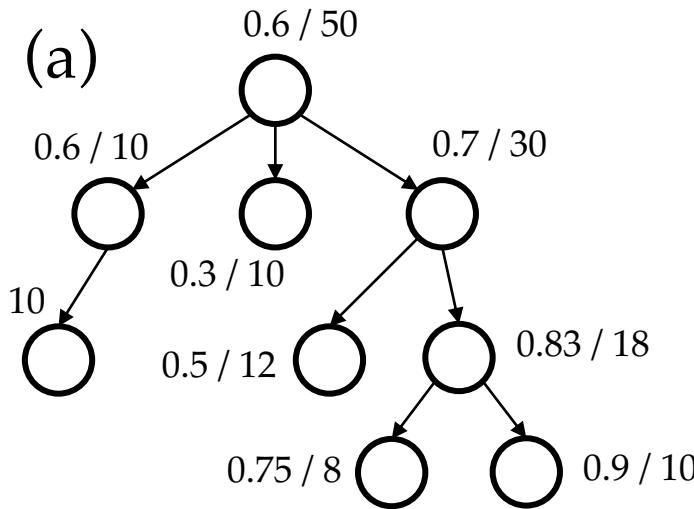
Rodrigo Fonseca¹



Yuandong Tian²

Monte Carlo Tree Search

Search towards the good nodes while keeping exploration in mind

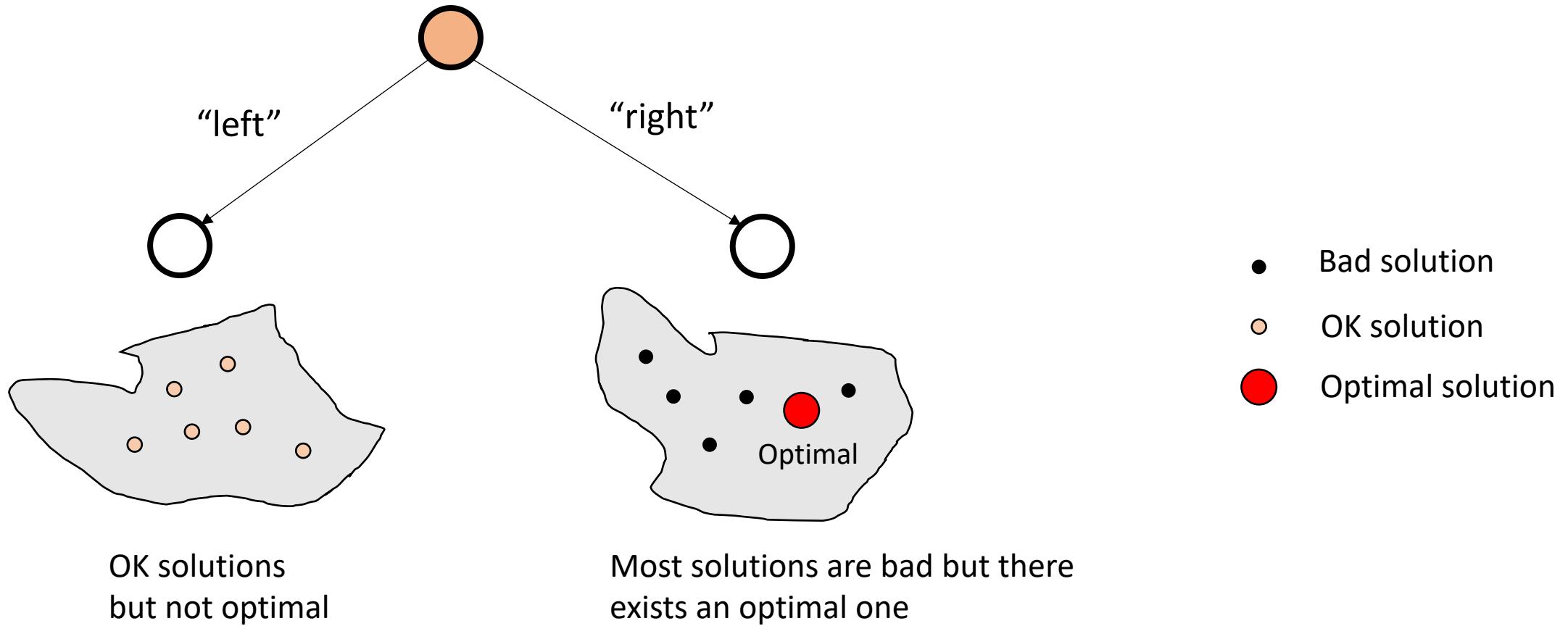


$$a_t = \arg \max_a Q(s_t, a) + u(s_t, a)$$

Exploration

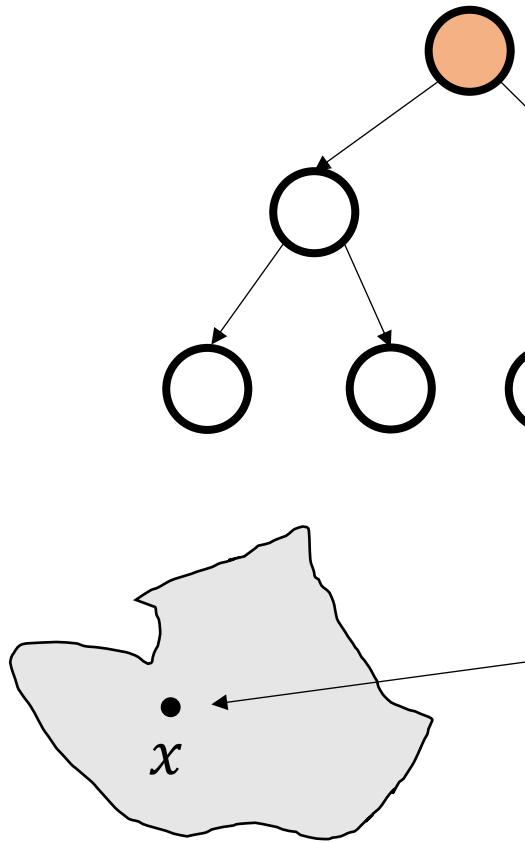
$$u(s, a) = c_{\text{puct}} P(s, a) \frac{\sqrt{\sum_b N(s, b)}}{1 + N(s, a)}$$

Why Exploration is Important



Sample in a Leaf

How?

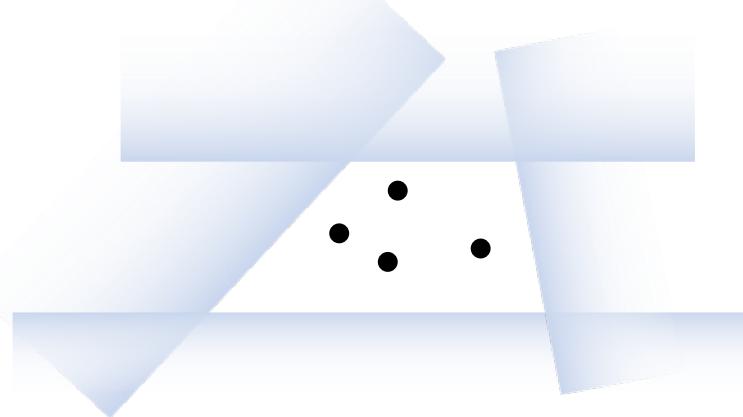


Random sample a point
Get the true function value

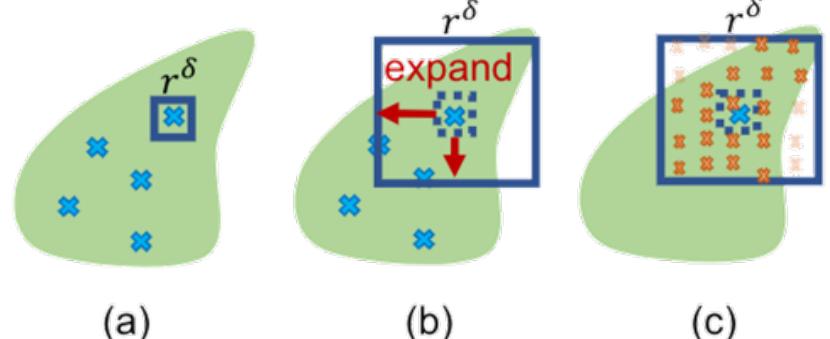
$$x \rightarrow f(x)$$

(can be expensive)

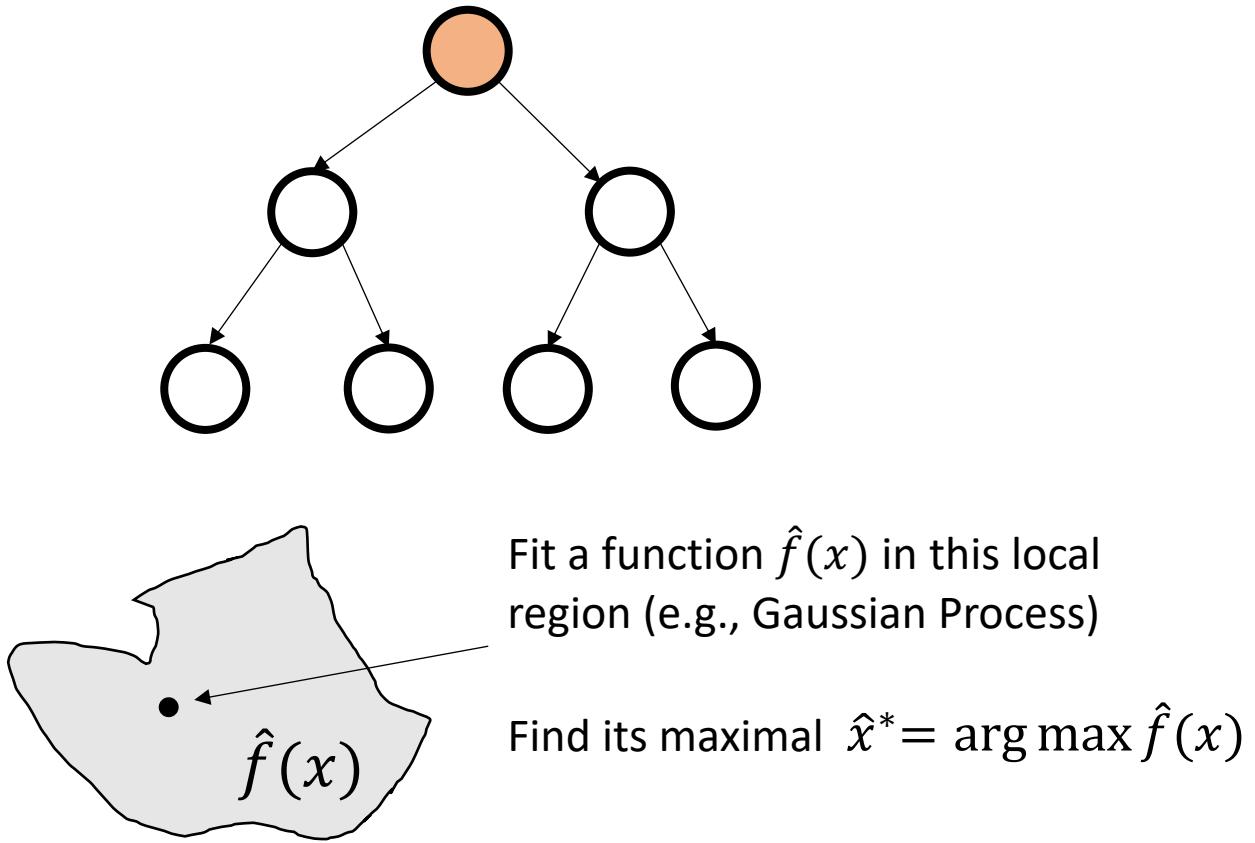
Linear case (sample within a polytope)



Nonlinear case (reject sampling)



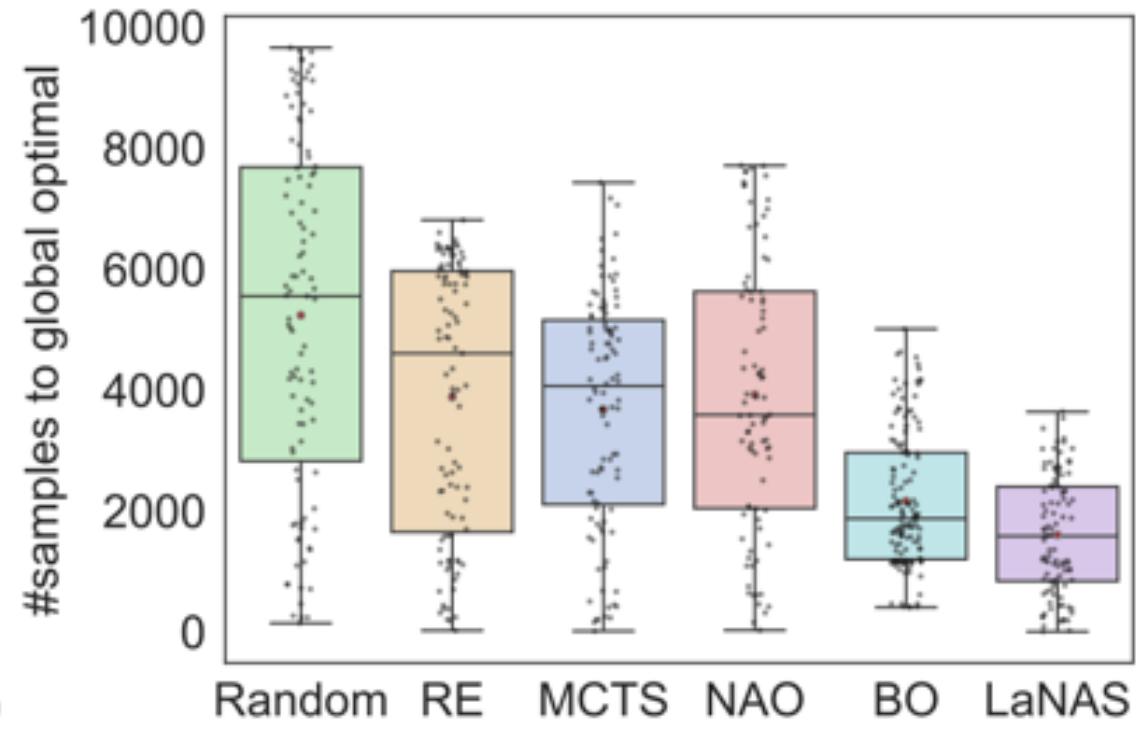
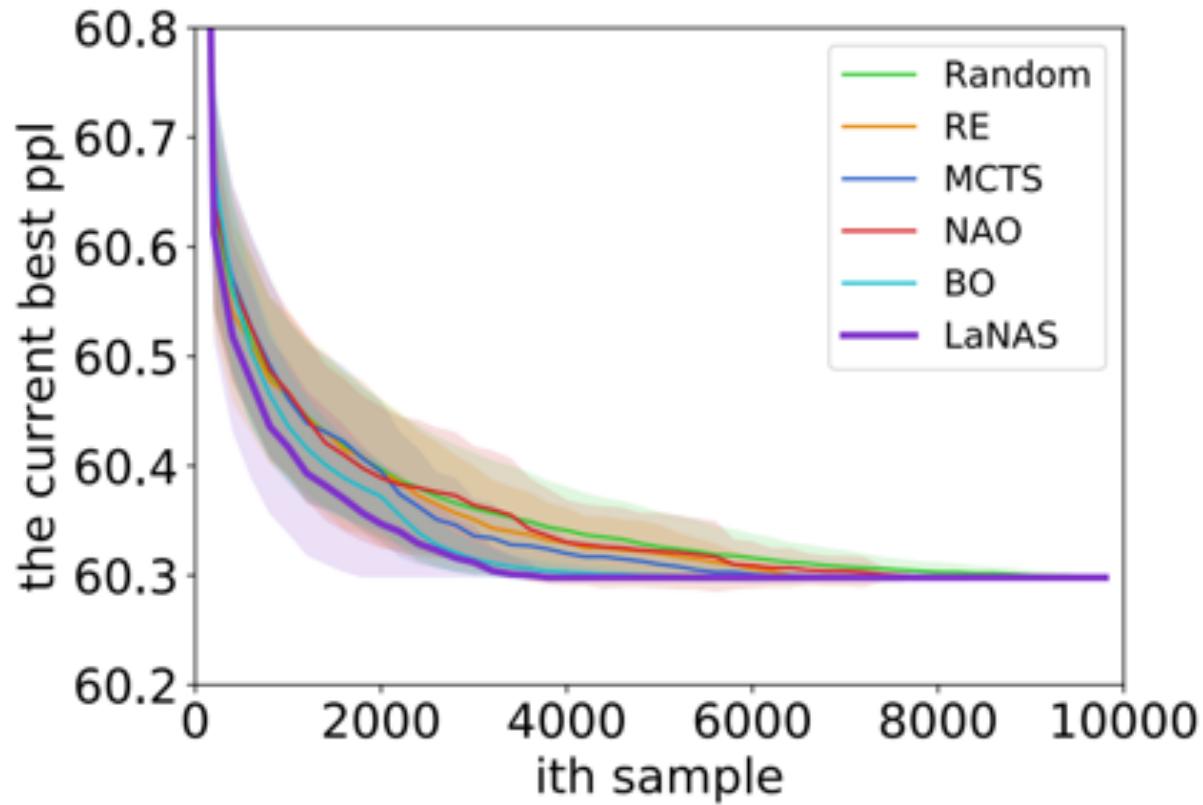
Sample in a Leaf



- Only need to fit a **local** GP
(much lower sample complexity)
- Function value is more **uniform** within a leaf
(easy for fitting)
- The fitting is needed only when #samples is small
Low computational cost
- Accurate** estimation of the leaf value
(uniform sample gives mean rather than max)

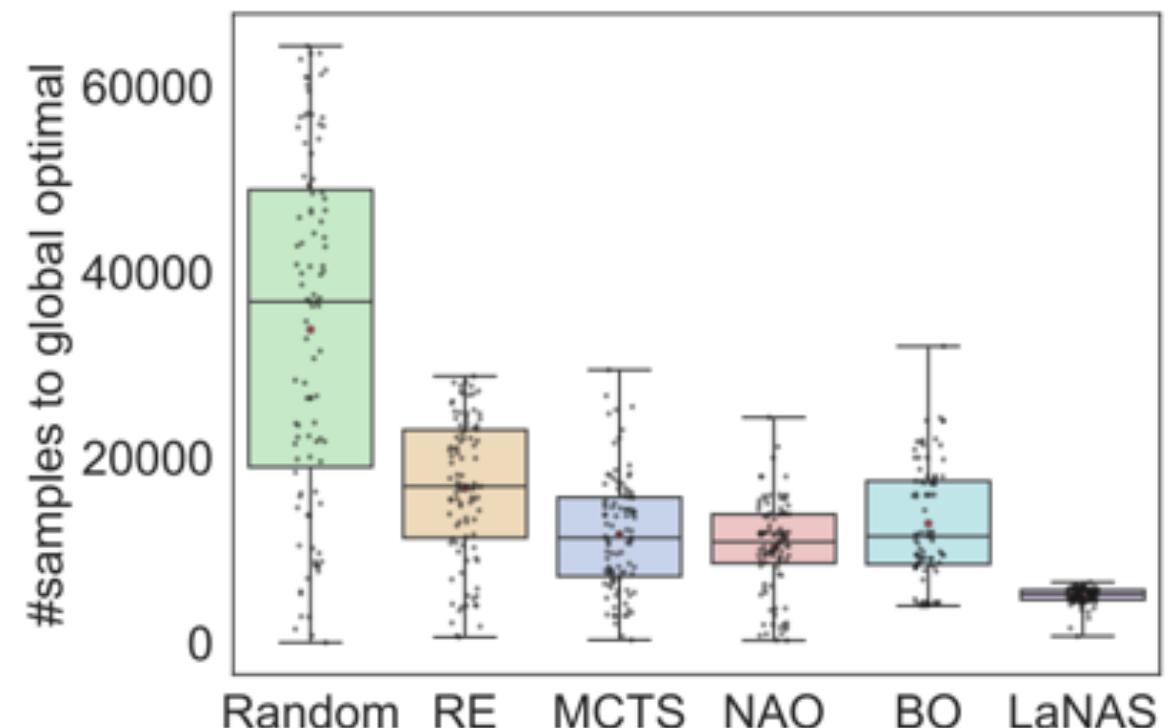
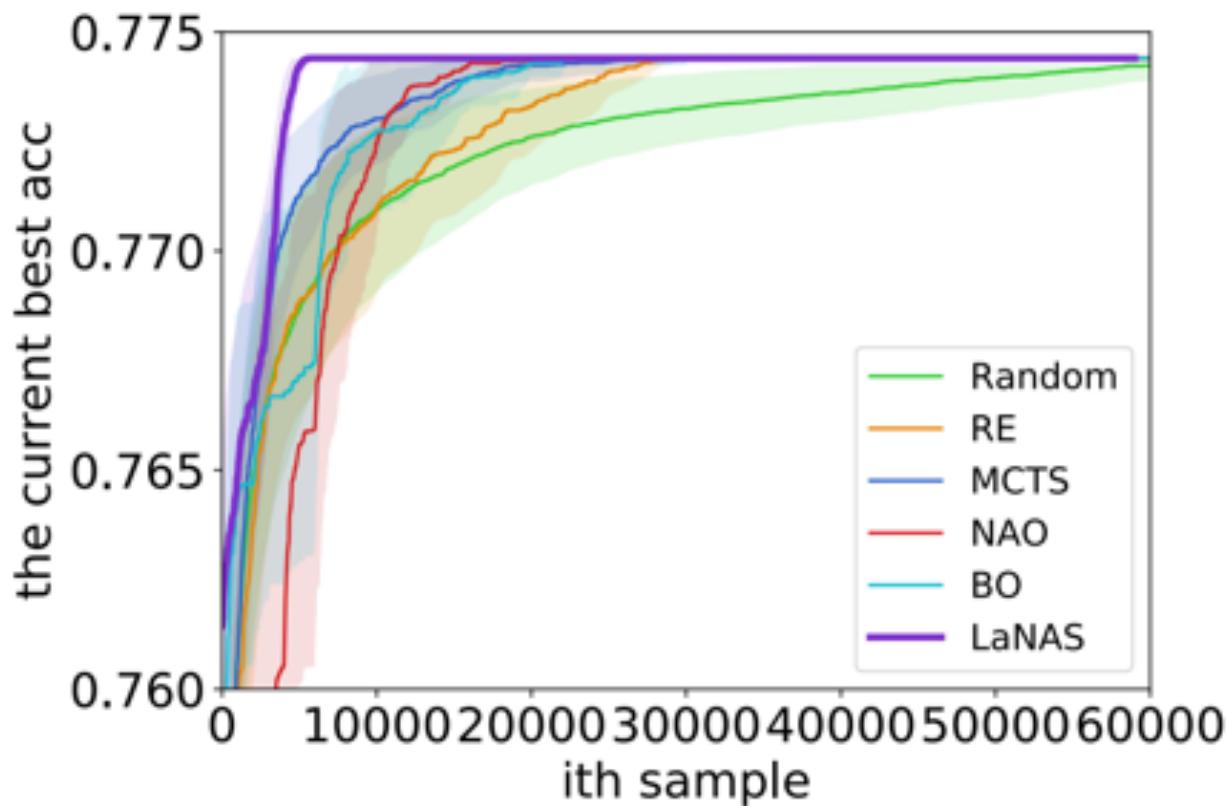
Performance

Customized dataset: LSTM-10K (PTB)



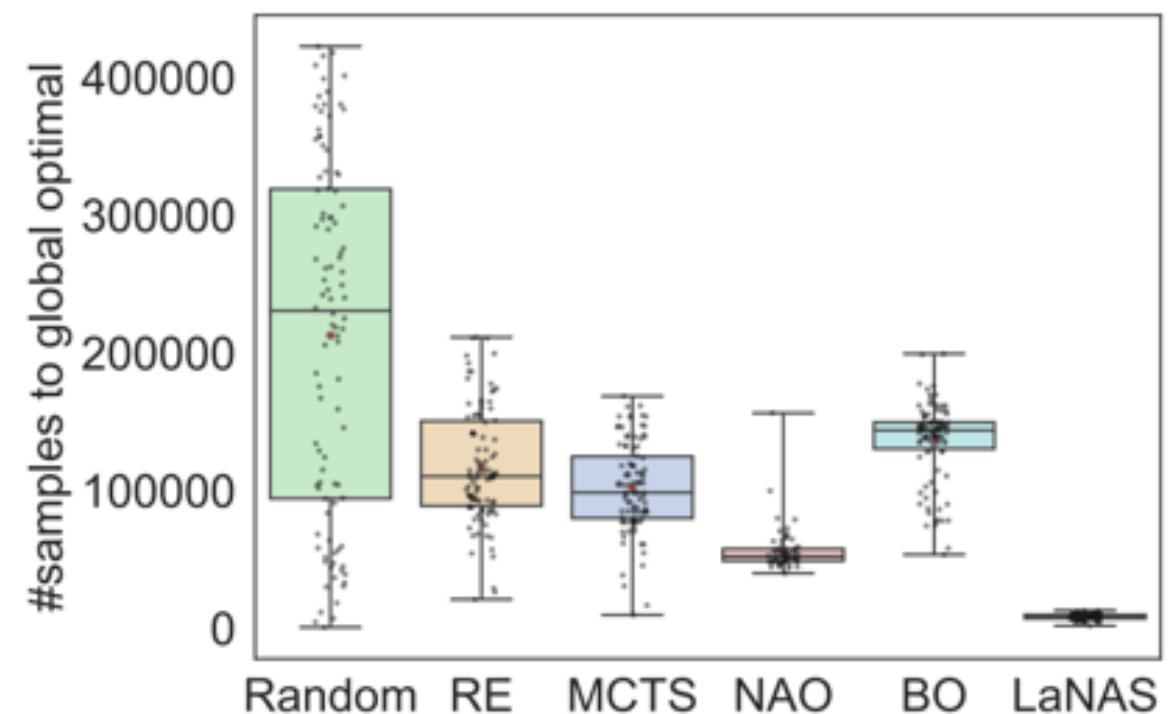
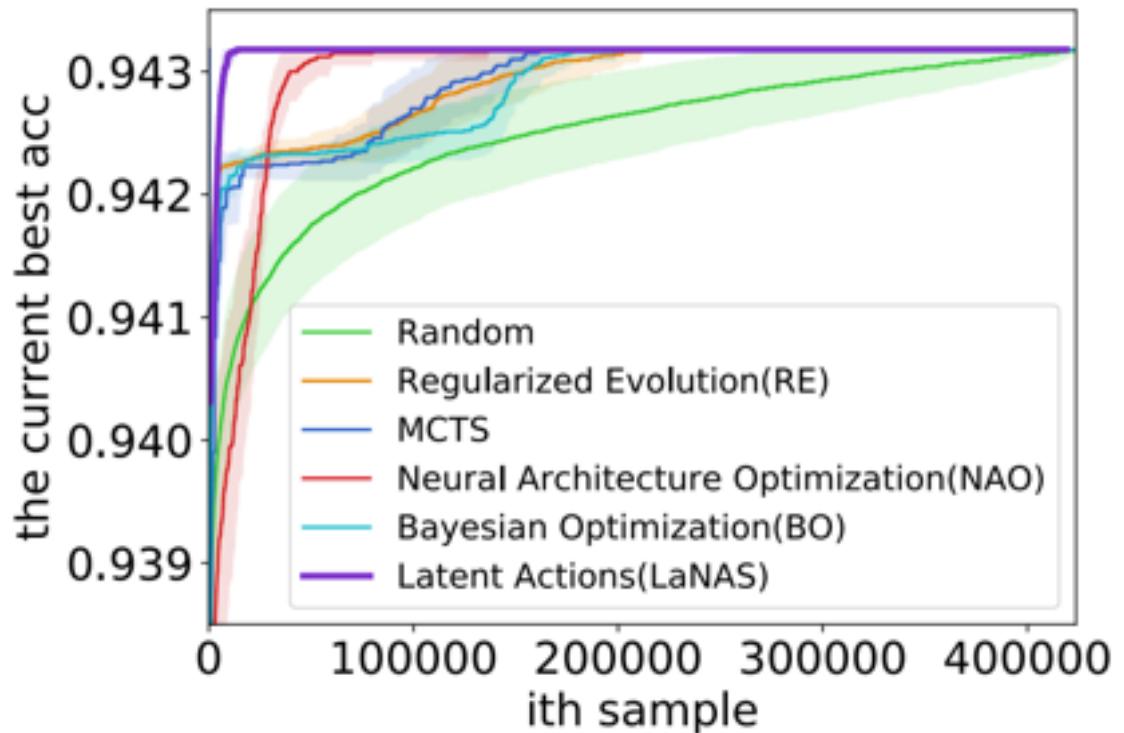
Performance

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 ± 0.05	27000	3150
PNASNet-5 [29]	X	3.2 M	3.41 ± 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 ± 0.04	27000	3150
EfficientNet-B7	✓	64M	1.01		
BiT-M	✓	60M	1.09		
LaNet+c/o	X	3.2 M	1.63 ± 0.05	800	150
LaNet+c/o	X	44.1 M	0.99 ± 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 ± 0.09	-	1.5
BayesNAS+c/o [31]	X	3.4 M	2.81 ± 0.04	-	0.2
ASNG-NAS+c/o [32]	X	3.9 M	2.83 ± 0.14	-	0.11
XNAS+c/o [33]	X	3.7 M	1.81		0.3
oneshot-LaNet+c/o	X	3.6 M	1.68 ± 0.06	-	3
oneshot-LaNet+c/o	X	45.3 M	1.2 ± 0.03	-	3

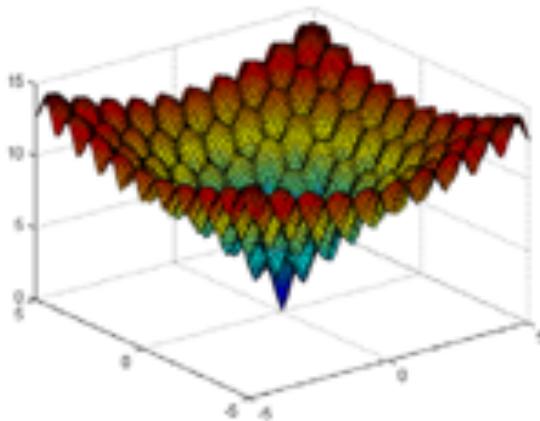
M: number of samples selected.

Open Domain

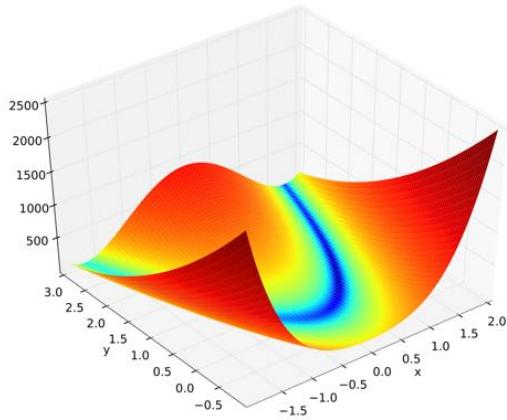
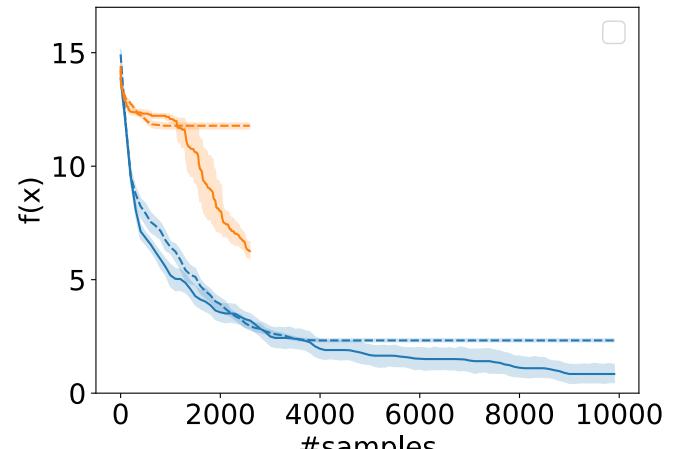
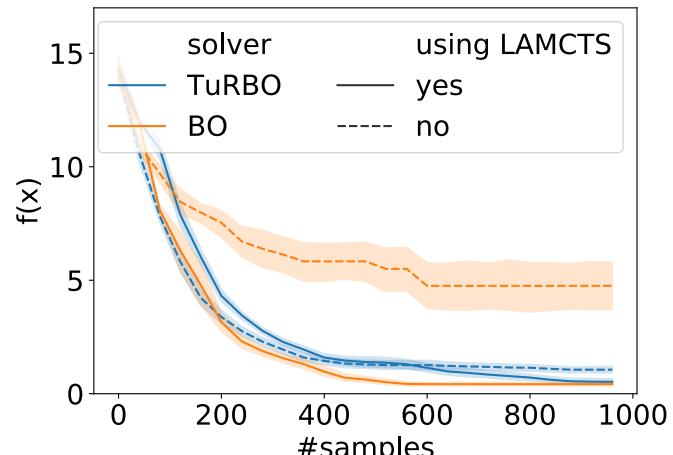
ImageNet
(mobile setting
Flop < 600M)

Model	FLOPs	Params	top1 / top5 err
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	24.3 / 7.6
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	25.0 / 7.7

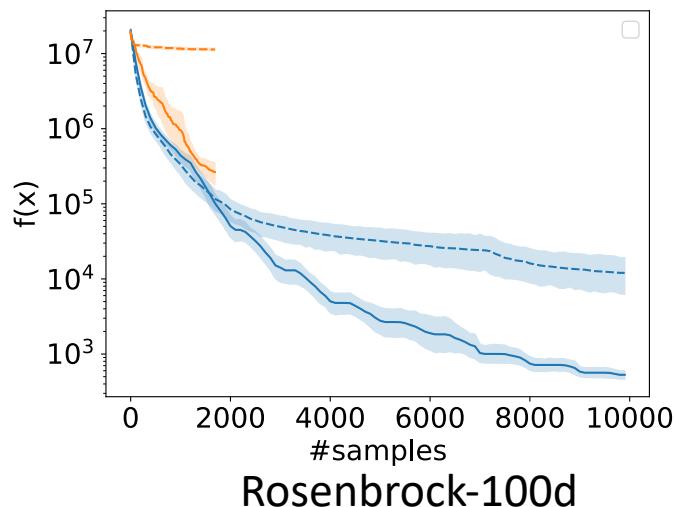
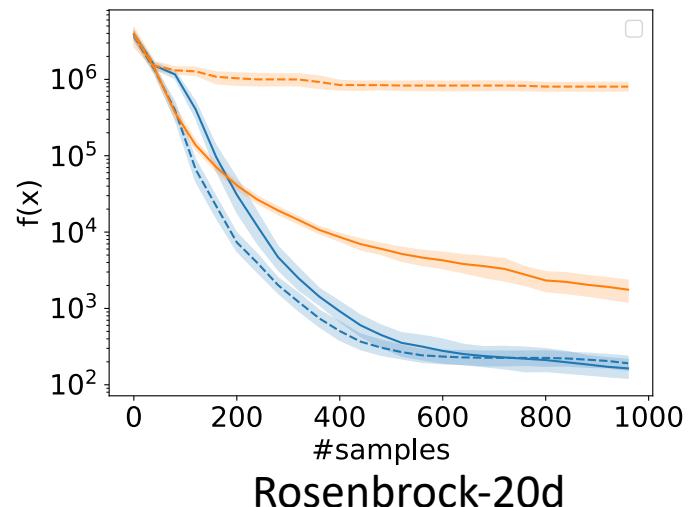
La-MCTS as a meta method $x^* = \arg \min_{x \in \Omega} f(x)$



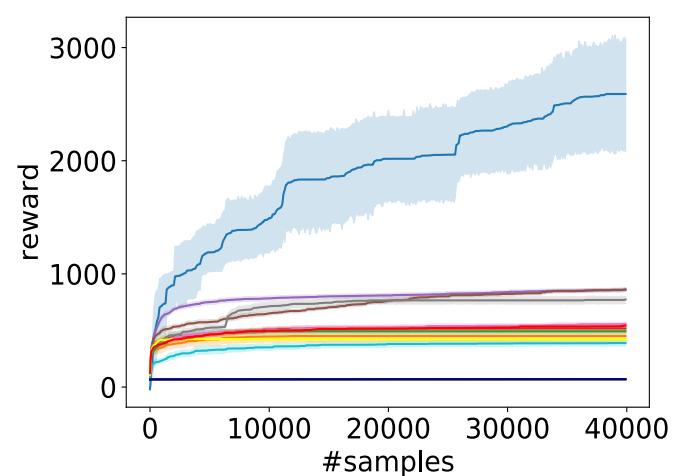
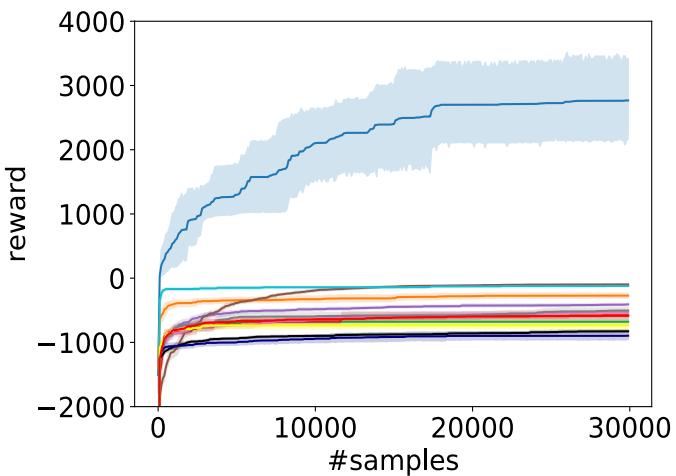
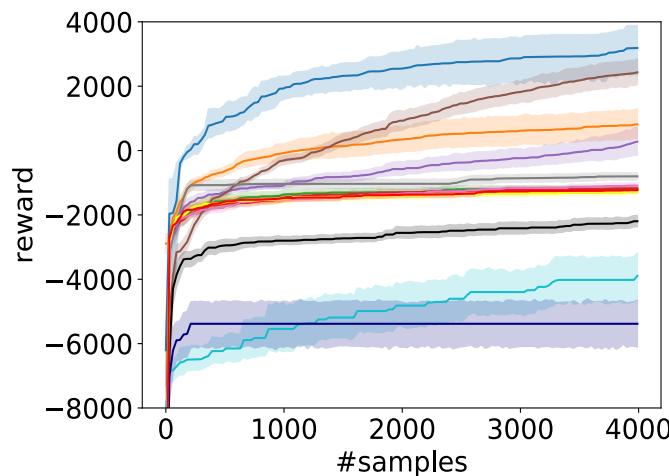
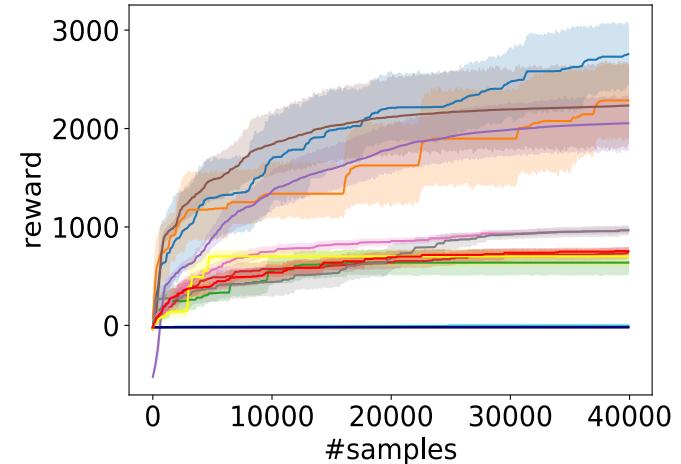
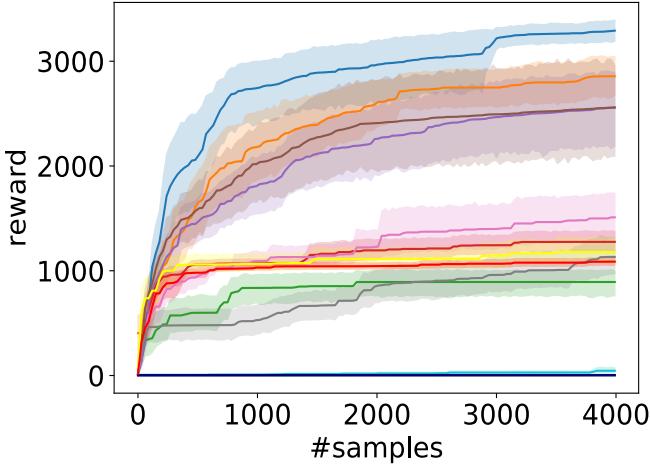
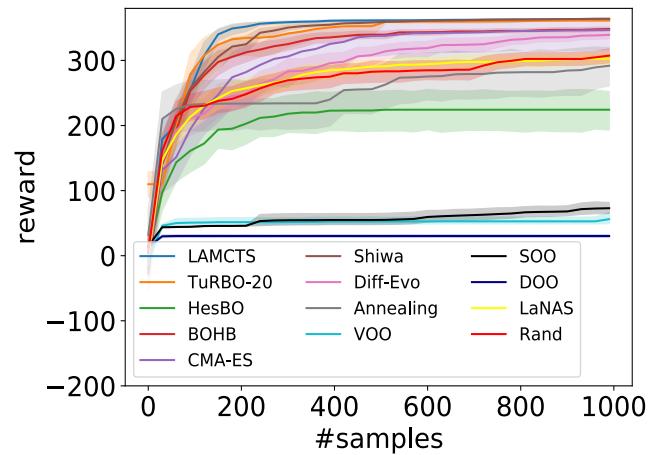
Ackley-2d



Rosenbrock-2d



Optimizing linear policy for Mujoco tasks



Limitations

Task	Reward Threshold	The average episodes (#samples) to reach the threshold				
		LA-MCTS	ARS V2-t [54]	NG-lin [55]	NG-rbf [55]	TRPO-nn [54]
Swimmer-v2	325	132	427	1450	1550	N/A
Hopper-v2	3120	2897	1973	13920	8640	10000
HalfCheetah-v2	3430	3877	1707	11250	6000	4250
Walker2d-v2	4390	N/A($r_{best} = 3314$)	24000	36840	25680	14250
Ant-v2	3580	N/A($r_{best} = 2791$)	20800	39240	30000	73500
Humanoid-v2	6000	N/A($r_{best} = 3384$)	142600	130000	130000	unknown

N/A stands for not reaching reward threshold.

r_{best} stands for the best reward achieved by LA-MCTS under the budget in Fig. 3.

Too many explorations might hurt in Mujoco tasks.

Code is public now!



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





Thanks!