

Thinking Fast and Slow with Deep Learning and Tree Search

by

Thomas Anthony, Zheng Tian, David Barber
(University College London)



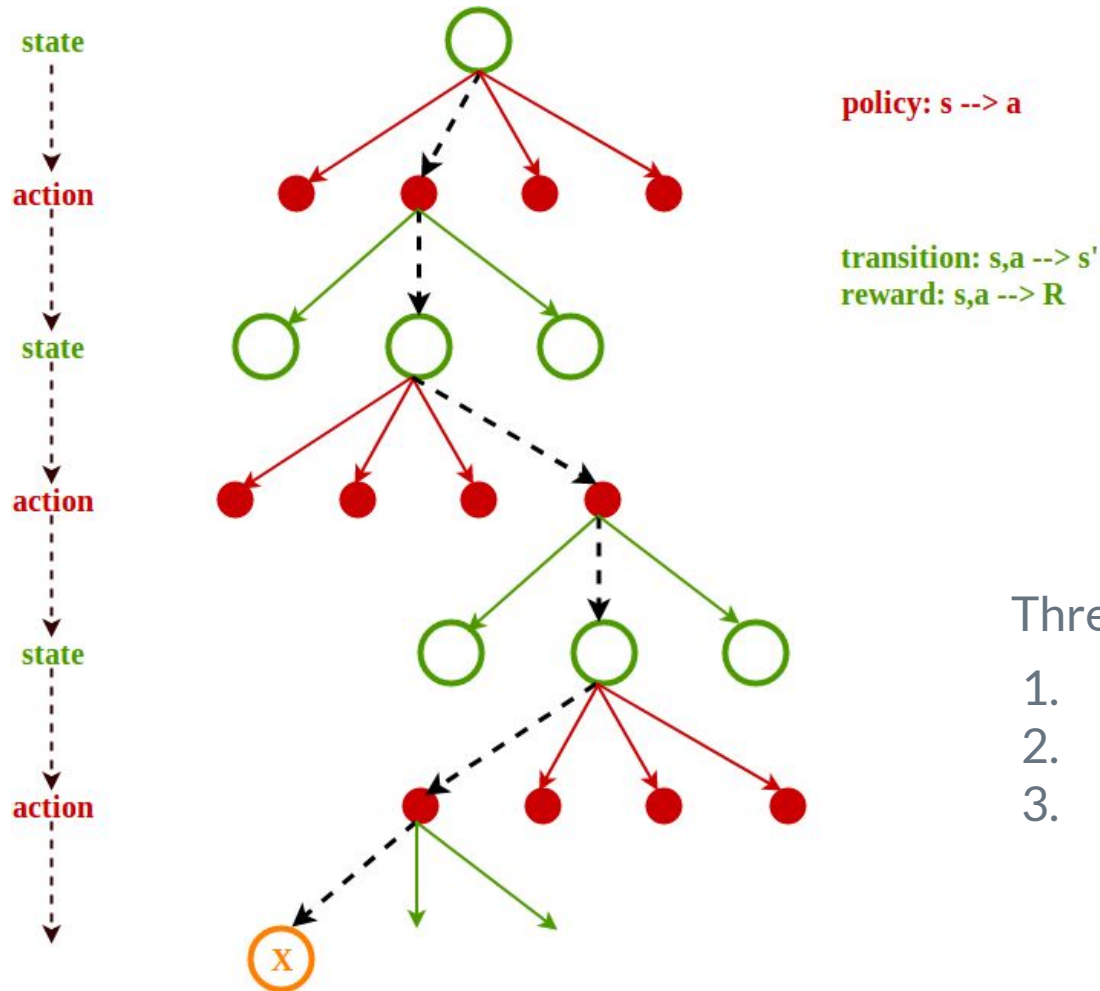
NIPS Debriefing

Thomas Moerland
(TU Delft)

Content

1. Background: Search & Reinforcement Learning
2. Algorithm: Expert Iteration (ExIt) → Iterating Search & RL
3. Results: Hex (game)
4. Discussion: Interpretation

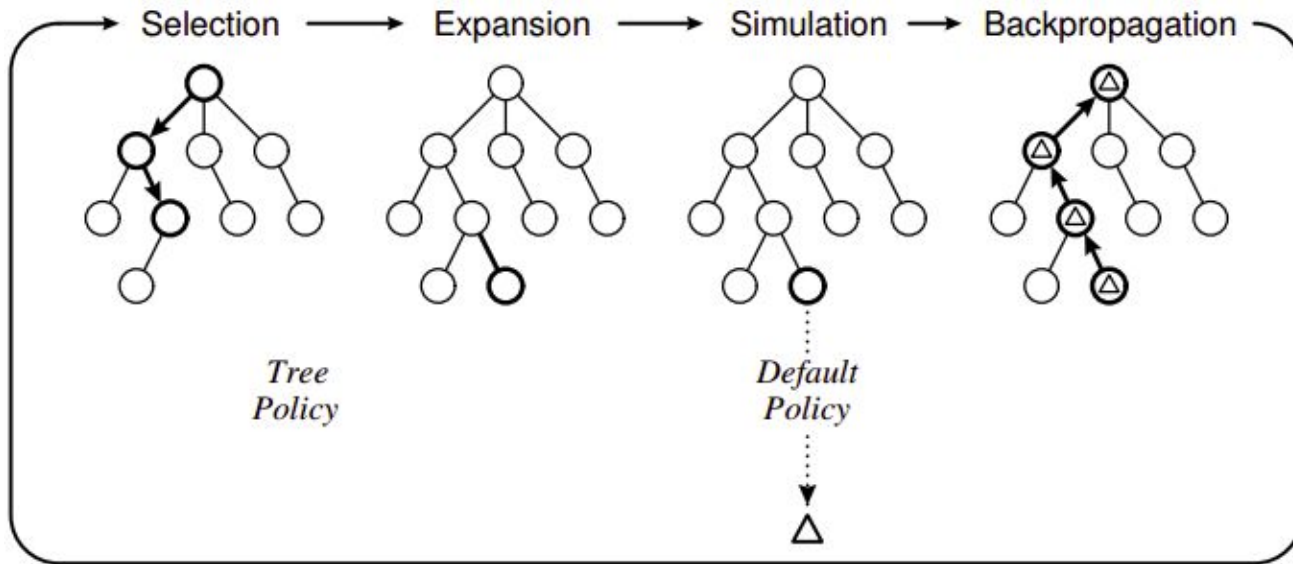
Sequential Decision Making



Three goals:

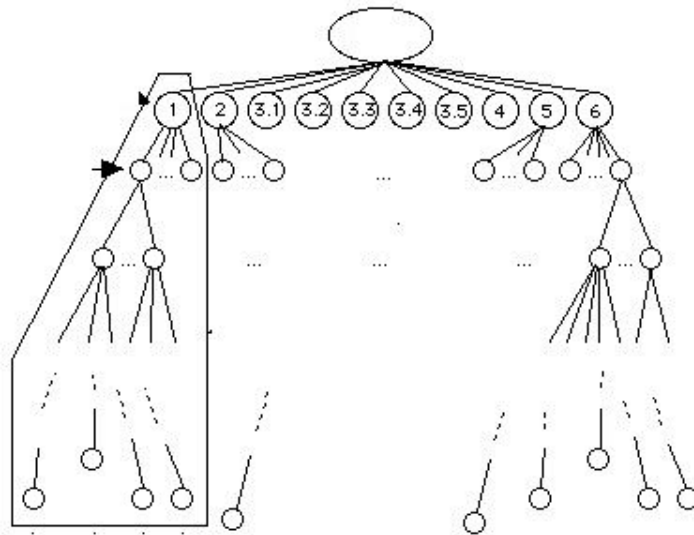
1. Reduce breadth
2. Reduce depth
3. Avoid repeating work

Monte Carlo Tree Search



- + Keep information local
- + Uncertainty for exploration

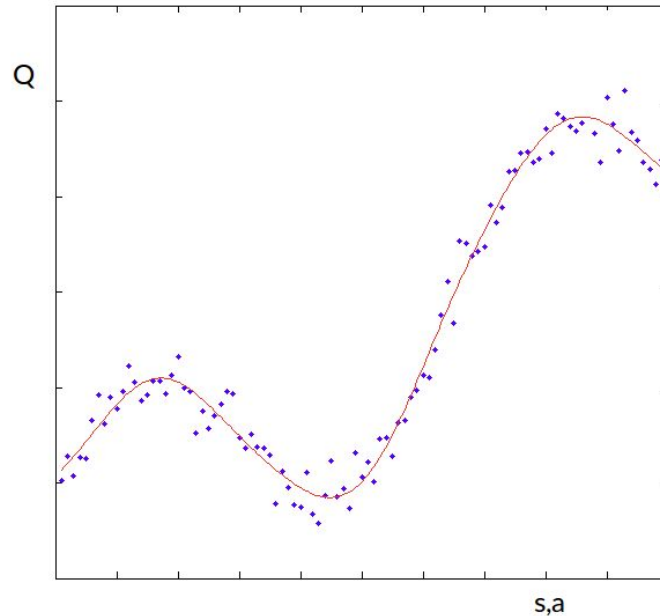
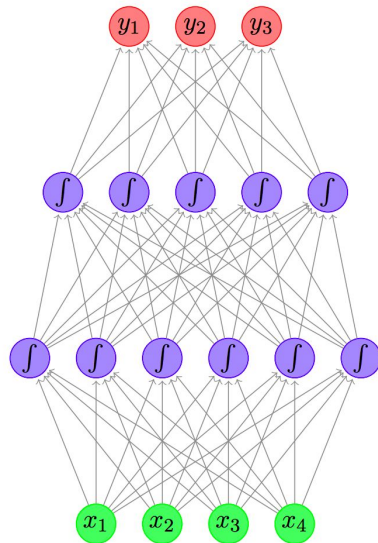
Monte Carlo Tree Search



- + Keep information local
- + Uncertainty for exploration

- No generalization
- Does not scale (memory)

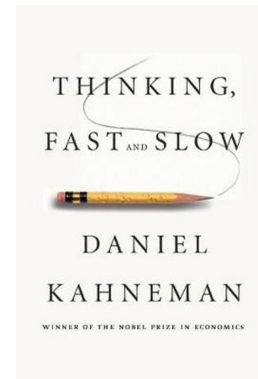
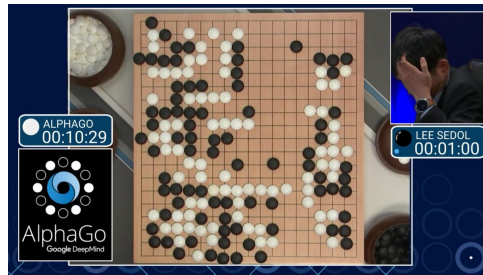
Reinforcement Learning



- + Generalization
- + Bootstrapping

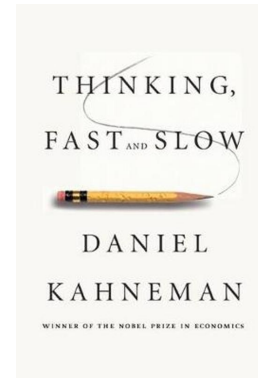
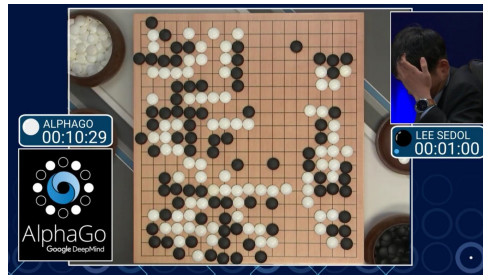
- Unstable / local minima / high variance
- Poor exploration
- (single trace)

Algorithm: Expert Iteration (ExIt)



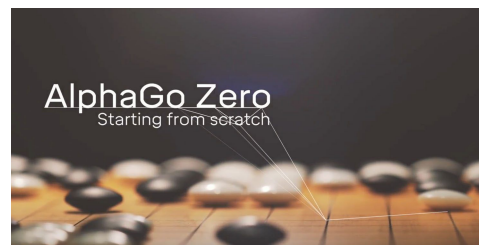
Tree Search	Supervised Learning (RL)
'Thinking slow' <i>Expert</i>	'Thinking fast' <i>Apprentice</i>

Algorithm: Expert Iteration (Exit)

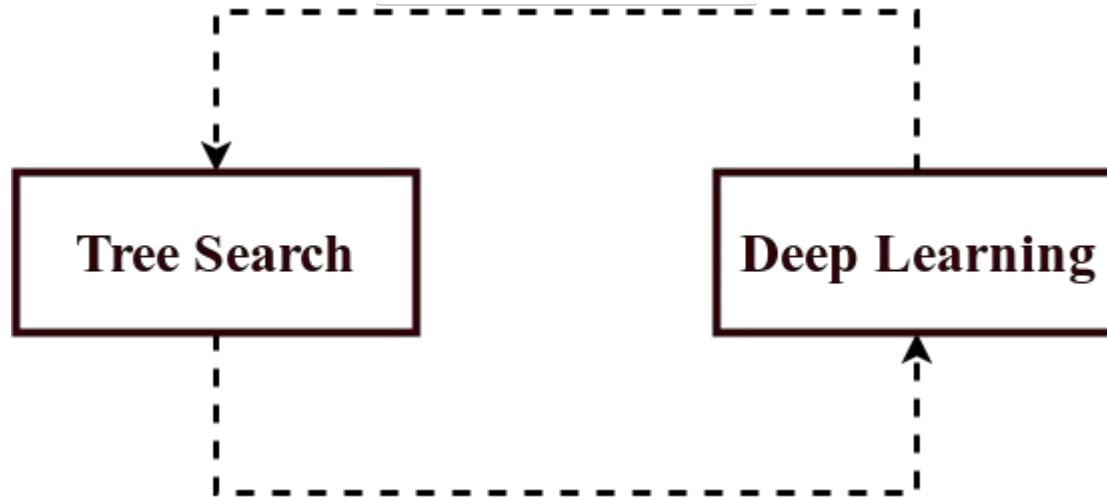


Tree Search	Supervised Learning (RL)
'Thinking slow' <i>Expert</i>	'Thinking fast' <i>Apprentice</i>

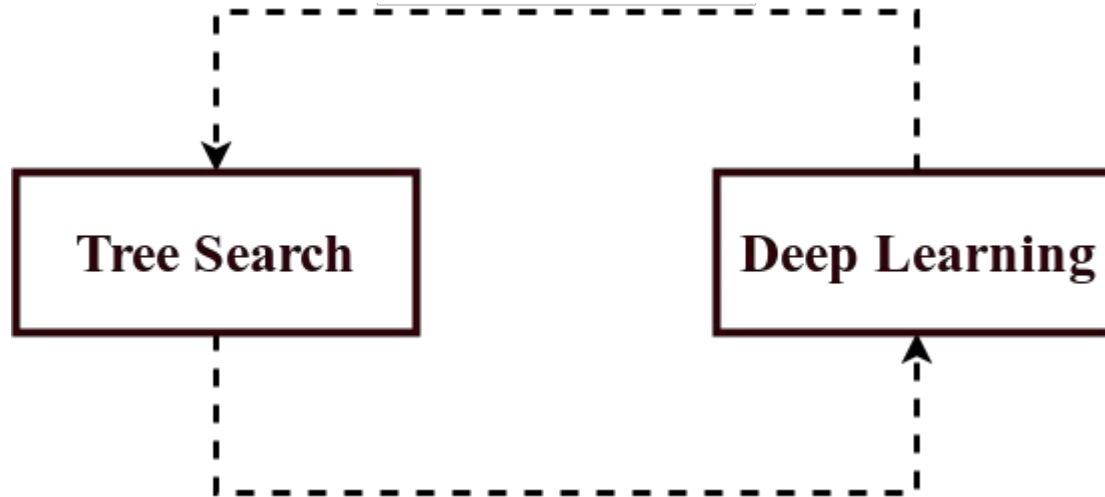
=



Iterating (local) search & generalization

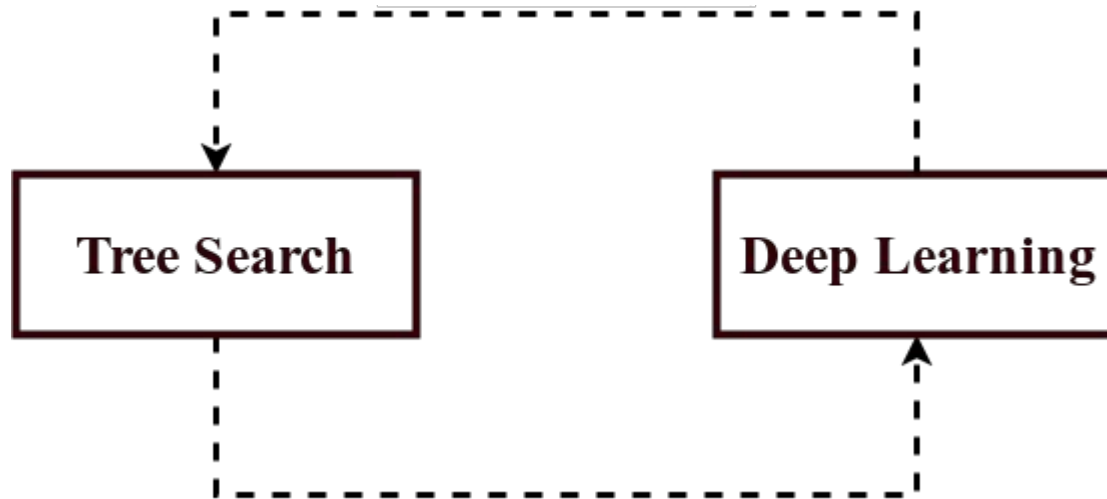


Iterating (local) search & generalization



What to store?

Iterating (local) search & generalization



What to store?

I. Reduce breadth

Policy

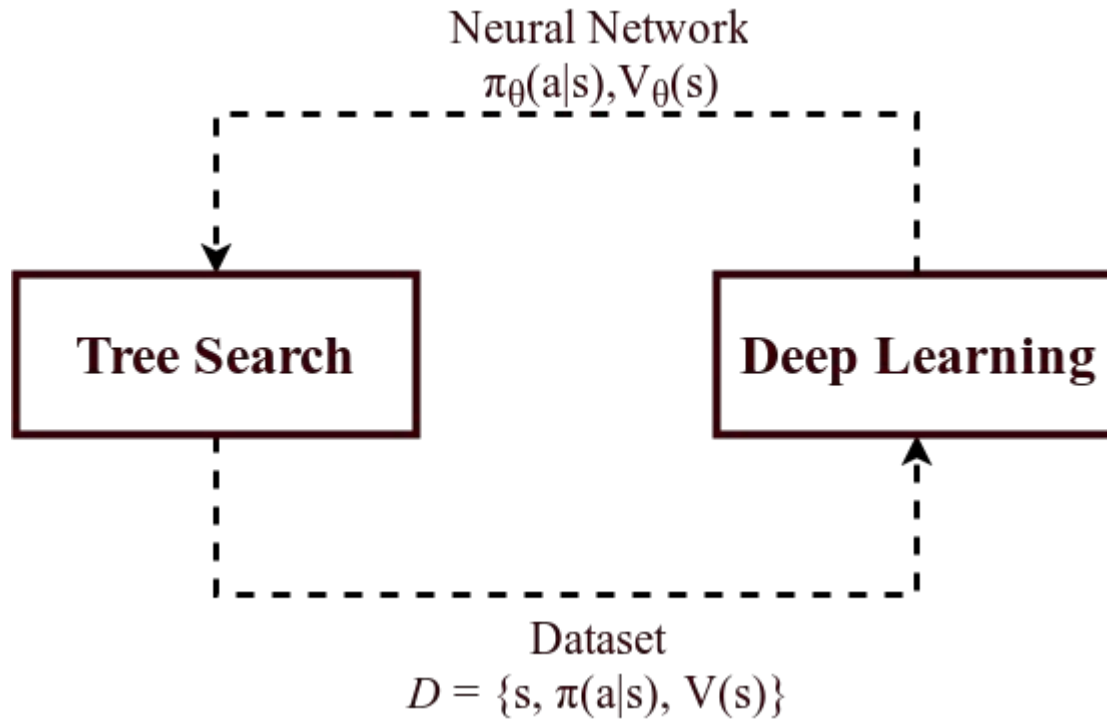
$s \rightarrow \pi(a)$

II. Reduce depth

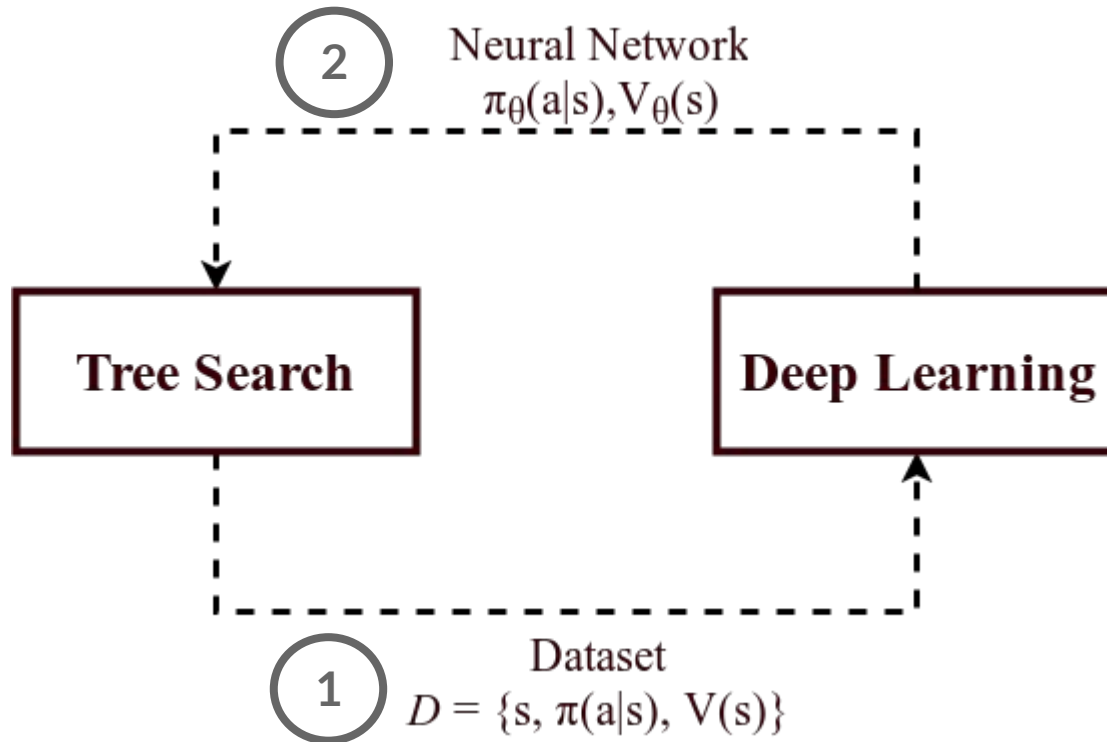
Value

$s \rightarrow V$

Iterating (local) search & generalization

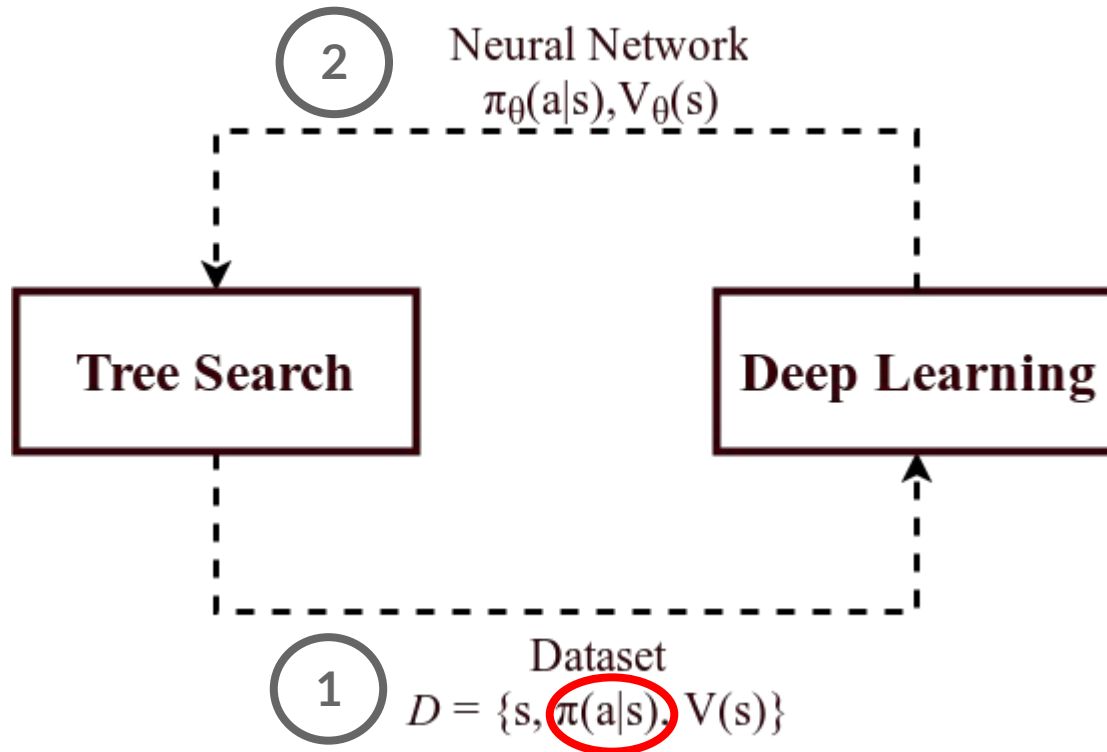


Iterating (local) search & generalization



1. How to process the tree search information?
2. How to use the network in the tree search?

Iterating (local) search & generalization



1. How to process the tree search information?
2. How to use the network in the tree search?

How to process the tree search result?

Policy: MCTS picks action with highest $n(s,a)$ at the root.

How to process the tree search result?

Policy: MCTS picks action with highest $n(s,a)$ at the root.

1. Chosen-action target (CAT):

$$a^* = \operatorname{argmax}_a (n(s, a))$$

$$\mathcal{L}_{\text{CAT}} = -\log[\pi(a^*|s)]$$

How to process the tree search result?

Policy: MCTS picks action with highest $n(s,a)$ at the root.

1. Chosen-action target (CAT):

$$a^* = \operatorname{argmax}_a (n(s, a))$$

$$\mathcal{L}_{\text{CAT}} = -\log[\pi(a^*|s)]$$

2. Tree-policy Targets (TPT):

$$\pi^*(a|s) = \frac{n(s, a)}{n(s)}$$

$$\mathcal{L}_{\text{TPT}} = -\sum_a \frac{n(s, a)}{n(s)} \log[\pi(a|s)]$$

How to process the tree search result?

Policy: MCTS picks action with highest $n(s,a)$ at the root.

1. Chosen-action target (CAT):

$$a^* = \operatorname{argmax}_a (n(s, a))$$

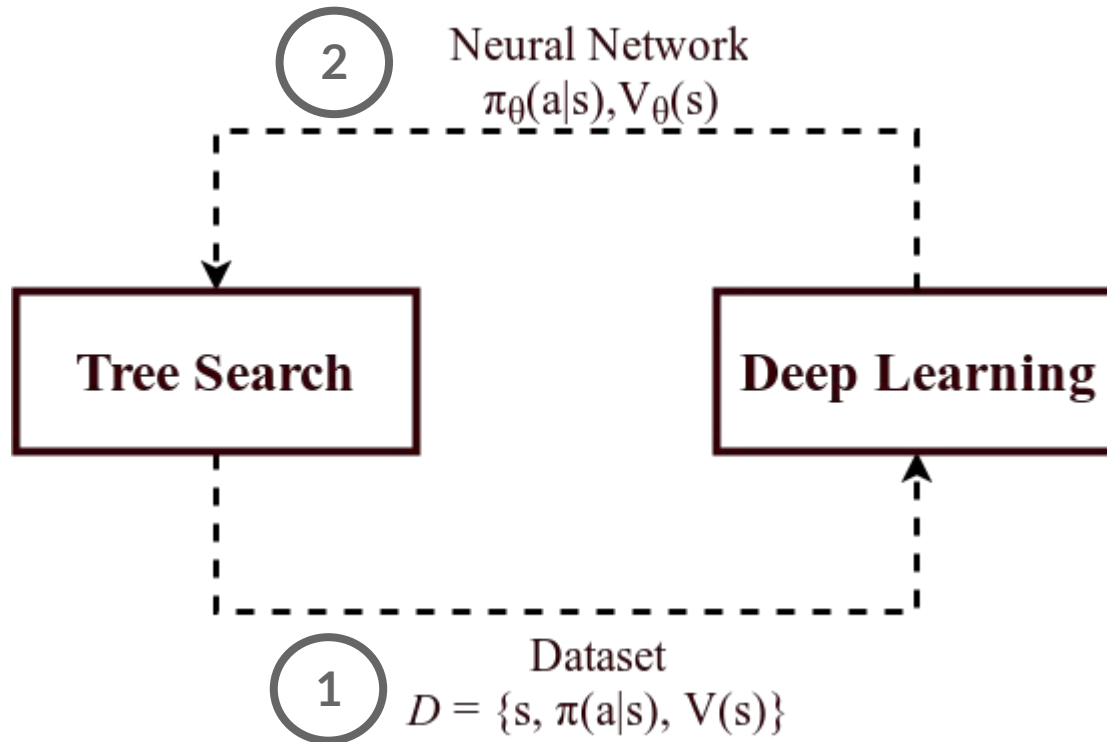
$$\mathcal{L}_{\text{CAT}} = -\log[\pi(a^*|s)]$$

2. Tree-policy Targets (TPT):

$$\pi^*(a|s) = \frac{n(s, a)}{\sum_a n(s, a)}$$

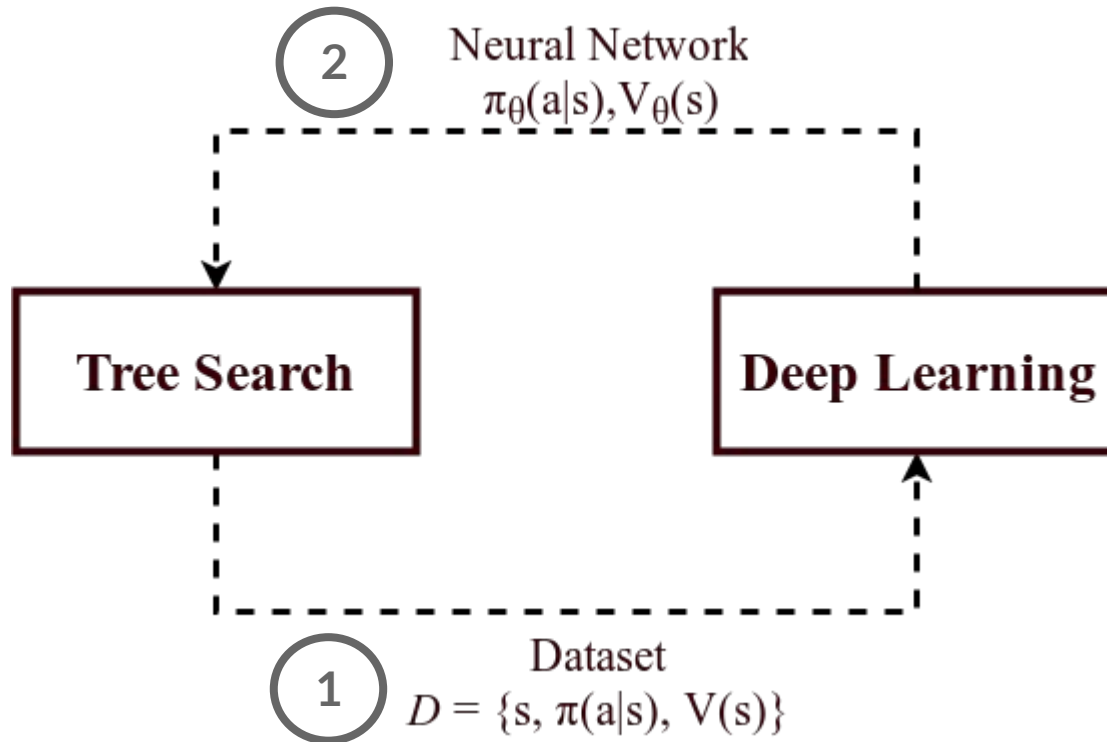
$$\mathcal{L}_{\text{TPT}} = -\sum_a \frac{n(s, a)}{n(s)} \log[\pi(a|s)]$$

Iterating (local) search & generalization



1. How to process the tree search information?
2. How to use the network in the tree search?

Iterating (local) search & generalization



1. How to process the tree search information?
2. How to use the network in the tree search?

How to warm start the tree search?

Policy (reduce breadth):

$$\text{UCT}(s, a) = \frac{r(s, a)}{n(s, a)} + c_b \sqrt{\frac{\log n(s)}{n(s, a)}}$$

How to warm start the tree search?

Policy:

This paper (ExIt)

$$\text{UCT}(s, a) = \frac{r(s, a)}{n(s, a)} + c_b \sqrt{\frac{\log n(s)}{n(s, a)}} + w_a \frac{\hat{\pi}(a|s)}{n(s, a) + 1}$$

How to warm start the tree search?

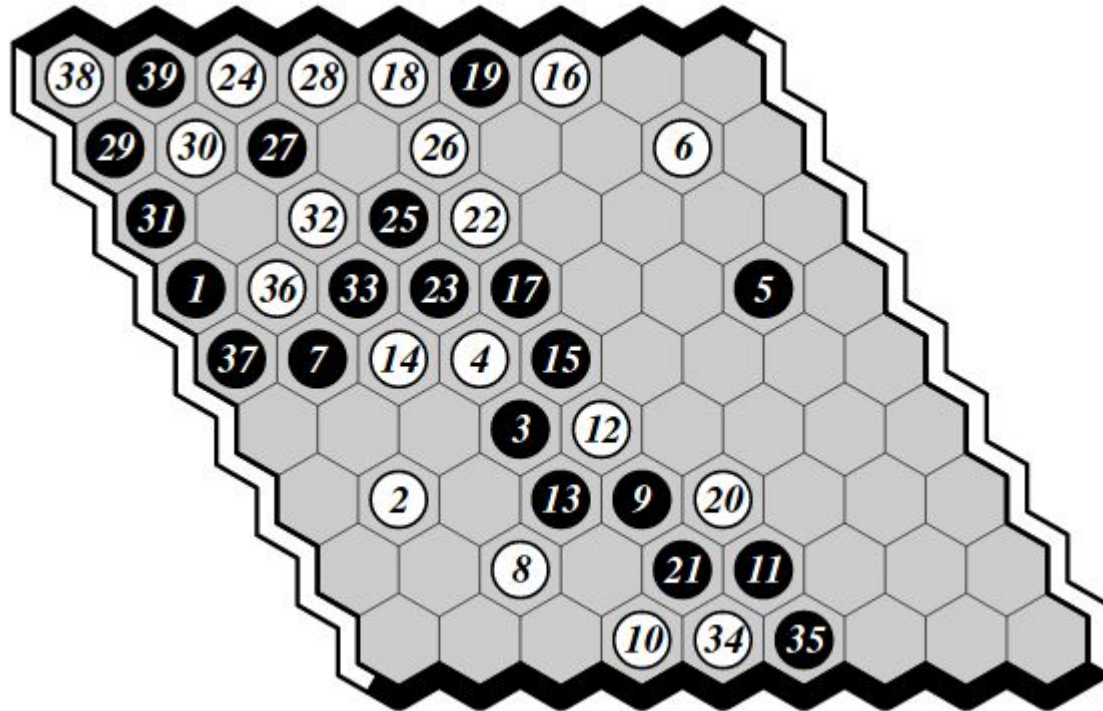
Policy:

AlphaGo Zero

$$\text{UCT}(s, a) = \frac{r(s, a)}{n(s, a)} + w_a \frac{\hat{\pi}(a|s)}{n(s, a) + 1}$$

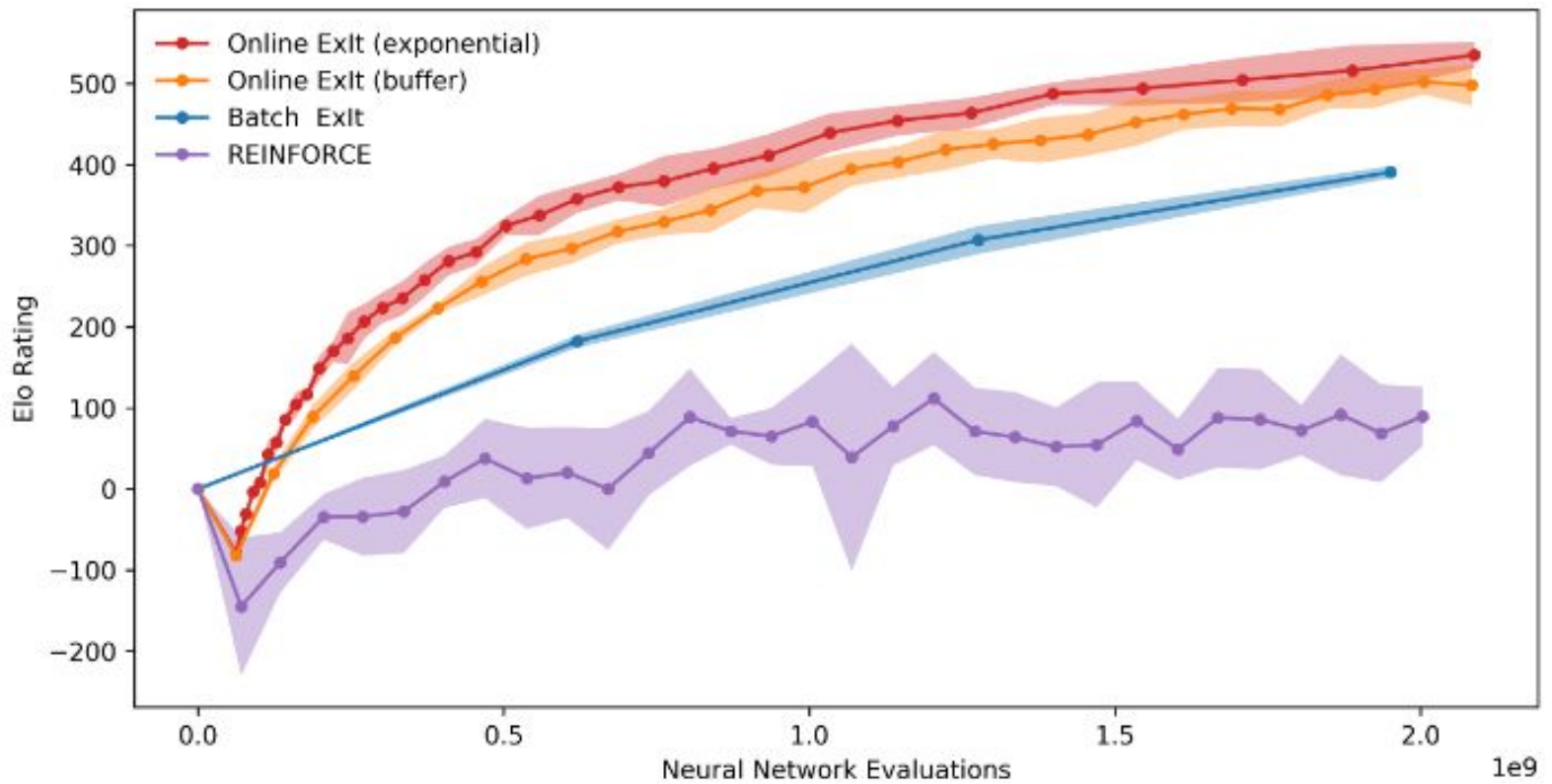
Results

Game = Hex



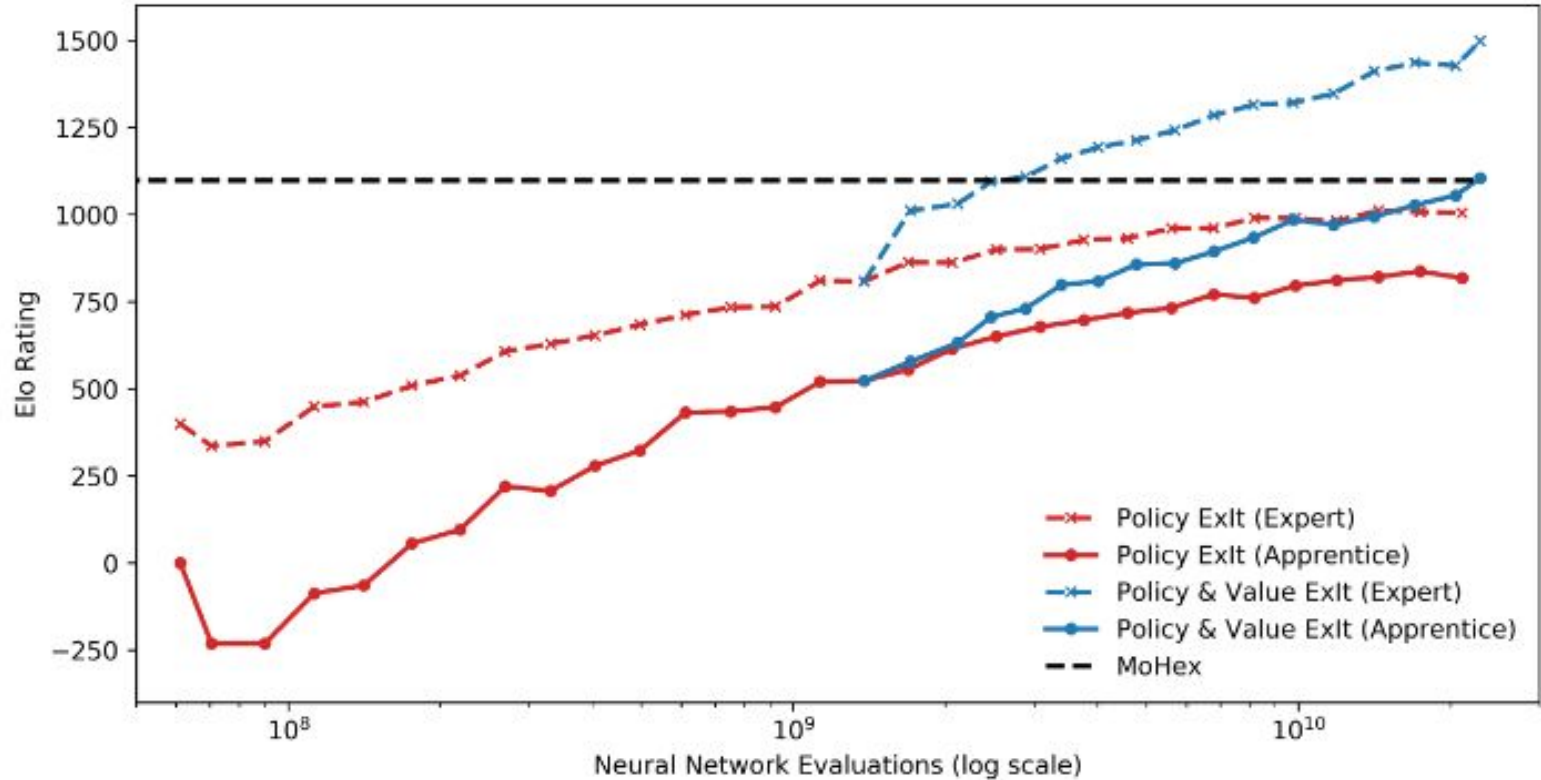
Results

Comparison to vanilla RL (policy network only)



Results

Added benefit of value network



Results

Battle versus MoHex

(search-based, game-specific pruning and end-game solving)

EXIT Setting	Time/move	EXIT win rate	MOHEX Setting	Solver	Time/move
10^4 iterations	$\sim 0.3s$	75.3%	10^4 iterations	No	$\sim 0.2s$
10^4 iterations	$\sim 0.3s$	59.3%	10^5 iterations	No	$\sim 2s$
10^4 iterations	$\sim 0.3s$	55.6%	4s/move	Yes	4s

Results

Battle versus MoHex

(search-based, game-specific pruning and end-game solving)

EXIT Setting	Time/move	EXIT win rate	MOHEX Setting	Solver	Time/move
10^4 iterations	$\sim 0.3s$	75.3%	10^4 iterations	No	$\sim 0.2s$
10^4 iterations	$\sim 0.3s$	59.3%	10^5 iterations	No	$\sim 2s$
10^4 iterations	$\sim 0.3s$	55.6%	4s/move	Yes	4s

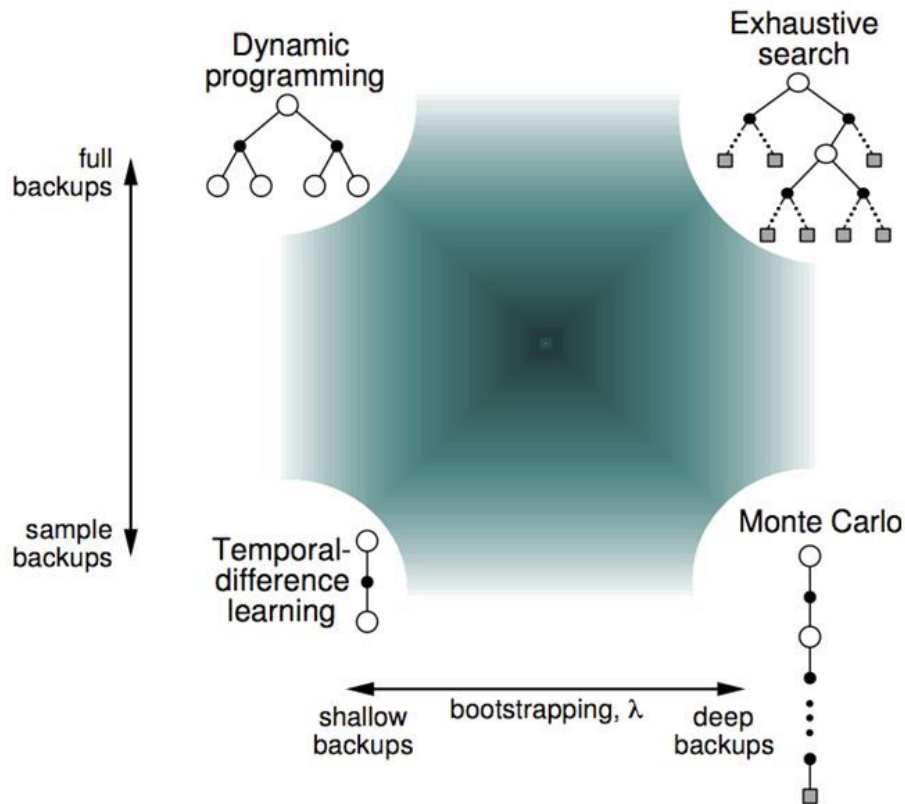
Maybe not as impressive as AlphaGo Zero

But: Training time ~ 100.000 times lower than AGO

(AGO = $\sim 1.6e12$ traces)

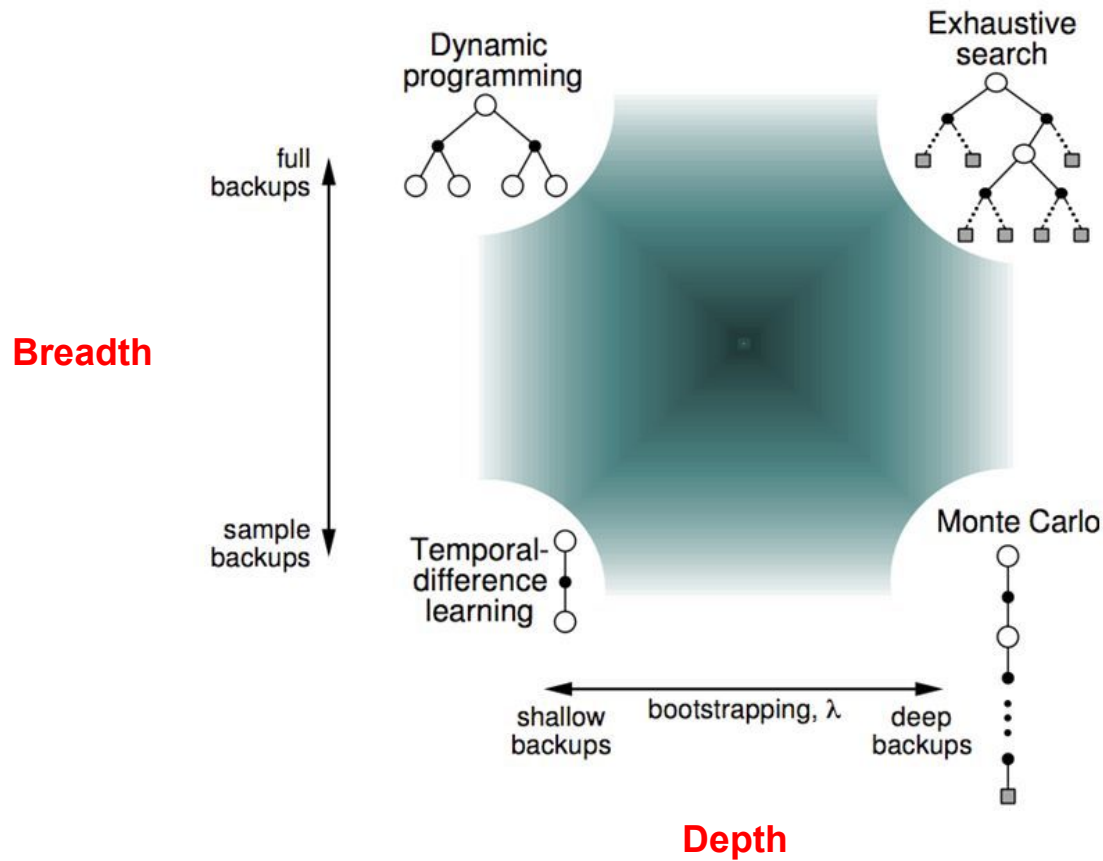
Discussion: Search and RL

Unified View



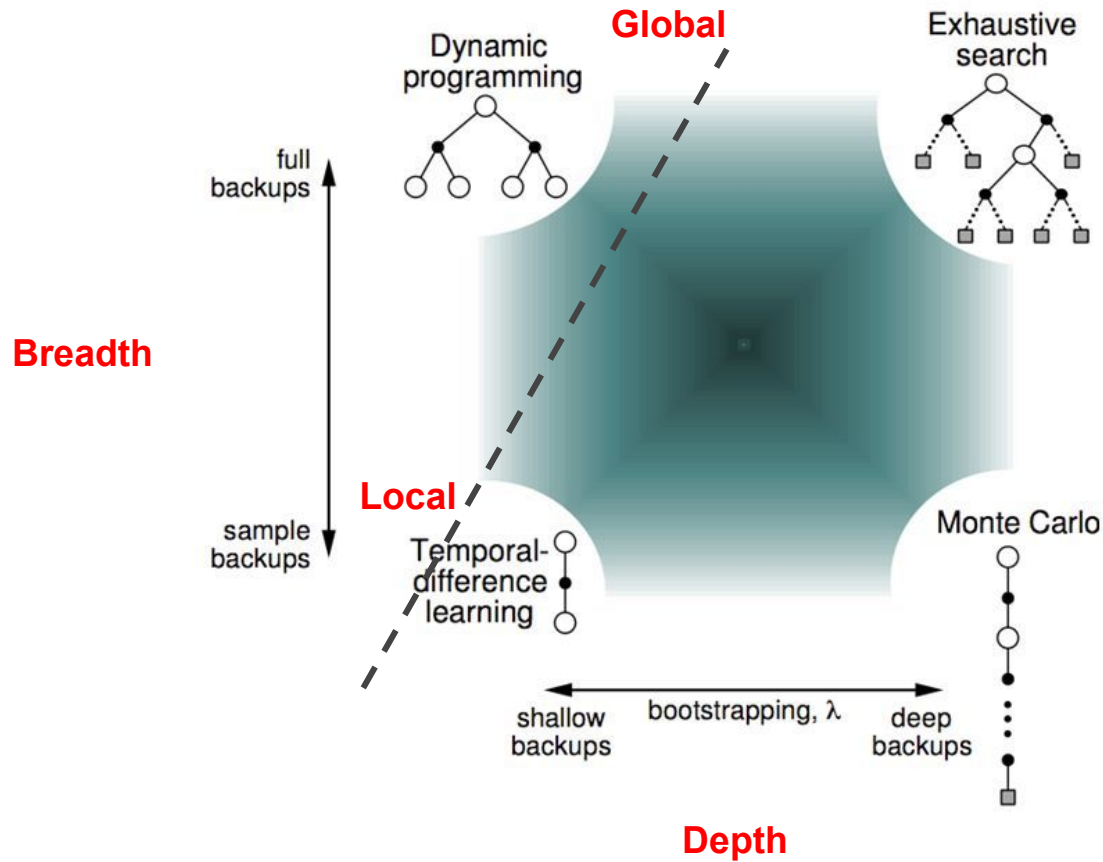
Discussion: Search and RL

Unified View



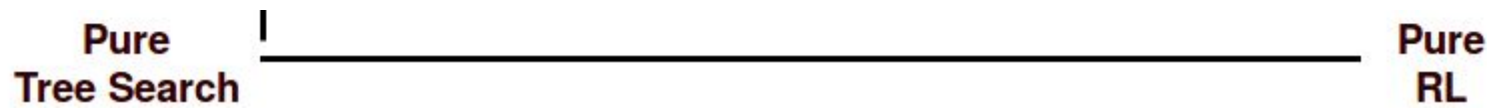
Discussion: Search and RL

Unified View



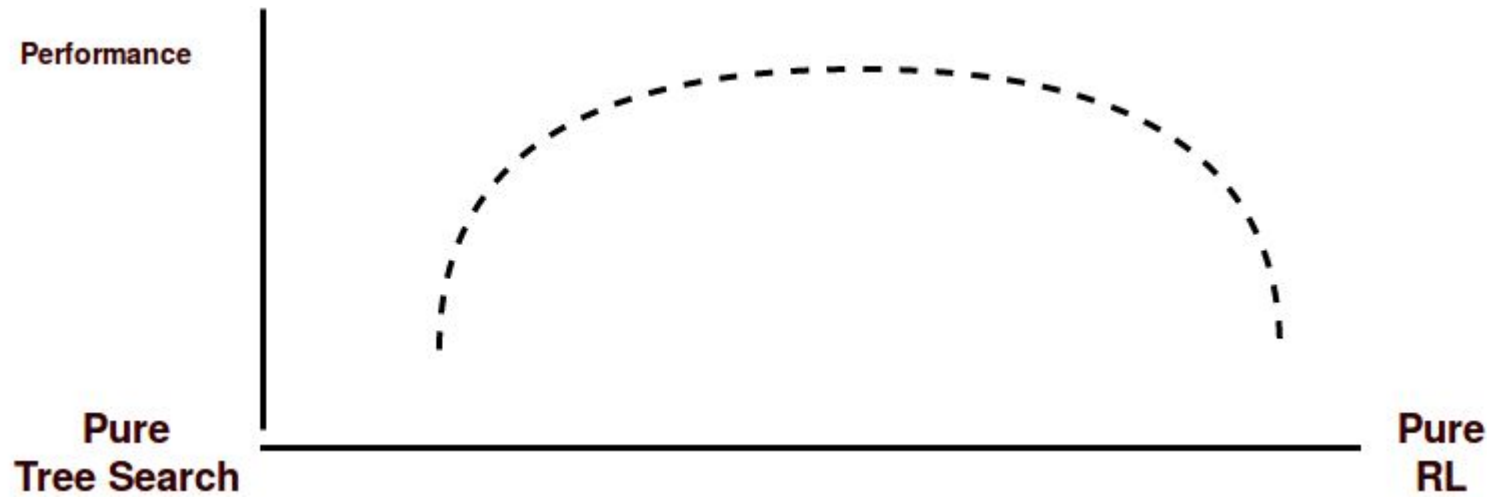
Discussion

Balancing local differentiation and global generalization



Discussion

Balancing local differentiation and global generalization



Discussion

- I empirically observed the problem of too early generalization in RL.

Discussion

- I empirically observed the problem of too early generalization in RL.

Open Questions:

- *Tree search*
- *What to store?*
- *How to steer search from a NN?*
- *Balance search/function approximation*

Thanks!