

# Learning to Optimize as Policy Learning

Yisong Yue

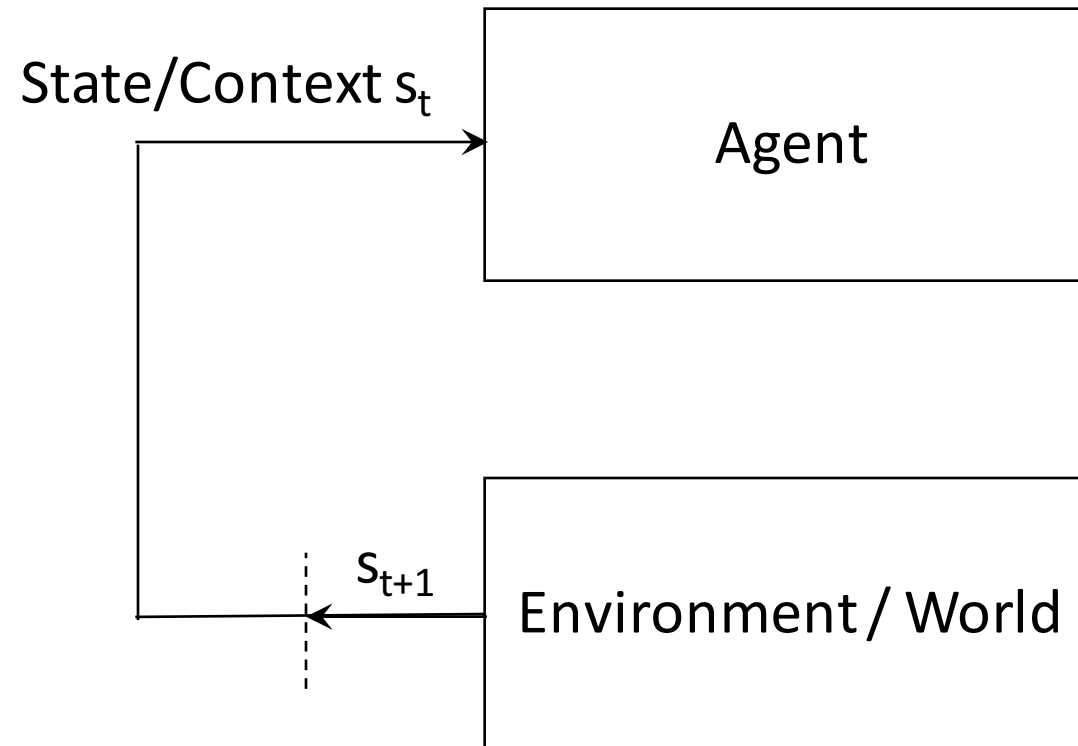
# Policy Learning (Reinforcement & Imitation)

**Goal:** Find “Optimal” Policy

**Imitation Learning:**  
Optimize imitation loss

**Reinforcement Learning:**  
Optimize environmental reward

**Learning-based Approach for  
Sequential Decision Making**



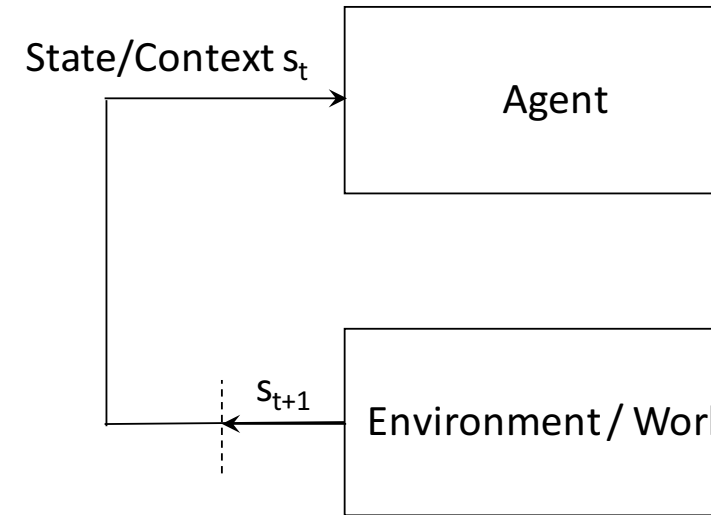
# Basic Formulation

(Typically a Neural Net)

- Policy:  $\pi(s) \rightarrow P(a)$   
State                      Action

- Roll-out:  $\tau = \langle s_0, a_0, s_1, a_1, s_2, \dots \rangle$  (aka trace or trajectory)  
Transition Function:  $P(s' | s, a)$

- Objective:  $\sum_i r(s_i, a_i)$



# Optimization as Sequential Decision Making

- Many Solvers are Sequential
  - Tree-Search
  - Greedy
  - Gradient Descent
- Can view solver as “agent” or “policy”
  - State = intermediate solution
  - Find a state with high reward (solution)
  - **Learn better local decision making**

- Formalize Learning
  - Builds upon mo
- Theoretical Analysis
- Interesting Algorithms



# Example #1: Learning to Search (Discrete)

## Integer Program

$$\max - \sum_{i=1}^5 x_i,$$

subject to:

$$x_1 + x_2 \geq 1,$$

$$x_2 + x_3 \geq 1,$$

$$x_3 + x_4 \geq 1,$$

$$x_3 + x_5 \geq 1,$$

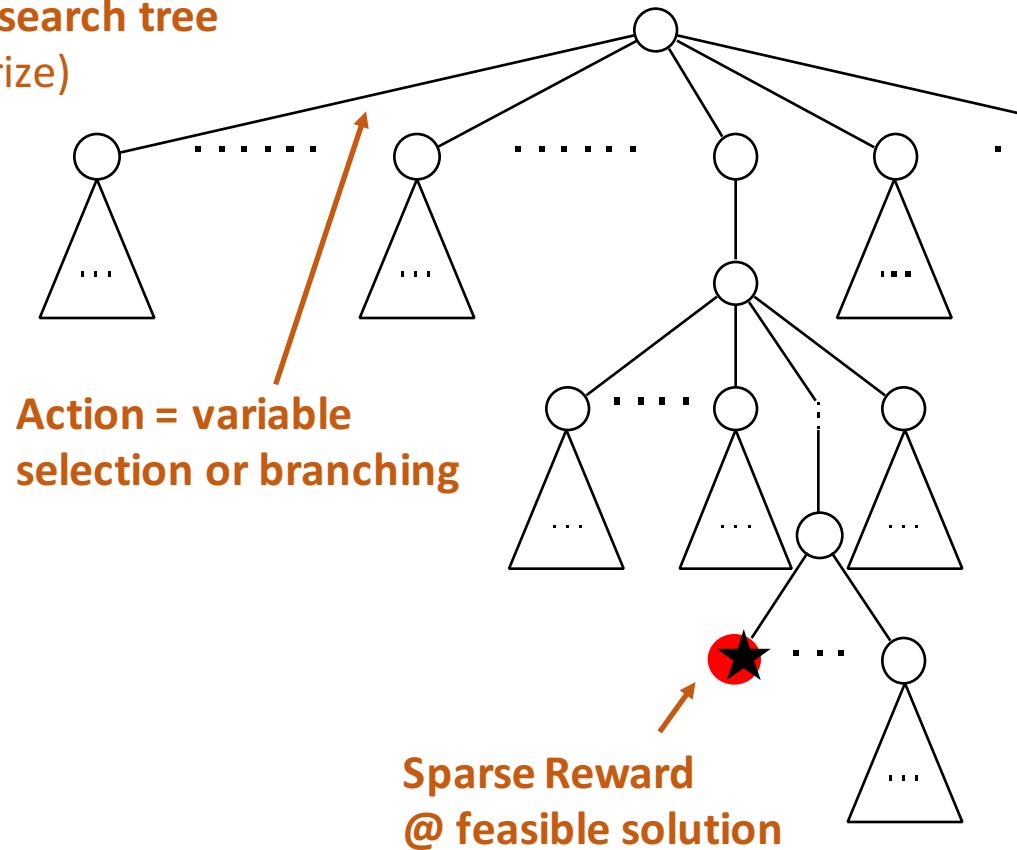
$$x_4 + x_5 \geq 1,$$

$$x_i \in \{0, 1\}, \forall i \in \{1, \dots, 5\}$$

State = partial search tree  
(need to featurize)



## Tree-Search (Branch and Bound)



[He et al., 2014][Khalil et al., 2016] [Song et al., arXiv]

# Example #1: Learning to Search (Discrete)

## Integer Program

$$\max - \sum_{i=1}^5 x_i,$$

subject to:

$$x_1 + x_2 \geq 1,$$

$$x_2 + x_3 \geq 1,$$

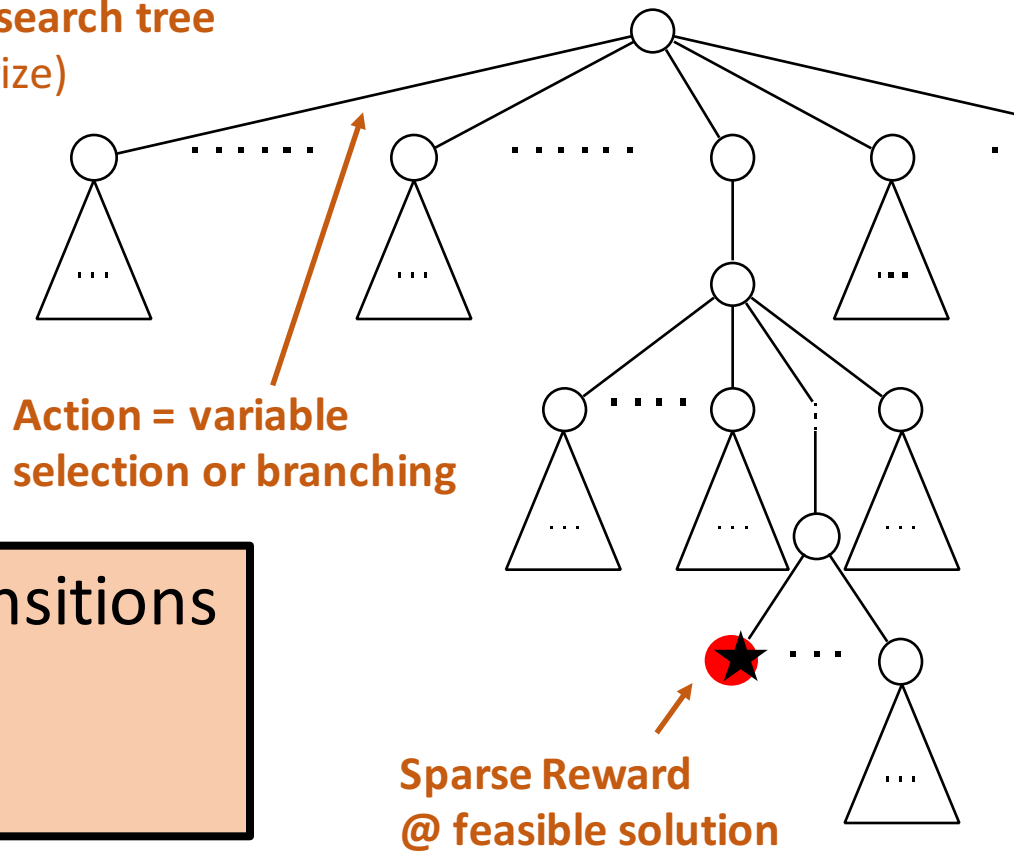
$$x_3 + x_4 \geq 1,$$

- Deterministic State Transitions
- Massive State Space
- Sparse Rewards

State = partial search tree  
(need to featurize)



## Tree-Search (Branch and Bound)



[He et al., 2014][Khalil et al., 2016] [Song et al., arXiv]

# Example #2: Learning Greedy Algorithms (discrete)

## Contextual Submodular Maximization:

- Greedy Sequential Selection:

- $\Psi \leftarrow \Psi \oplus \underset{a}{\operatorname{argmax}} F_x(\Psi \oplus a)$

Not Available at Test Time

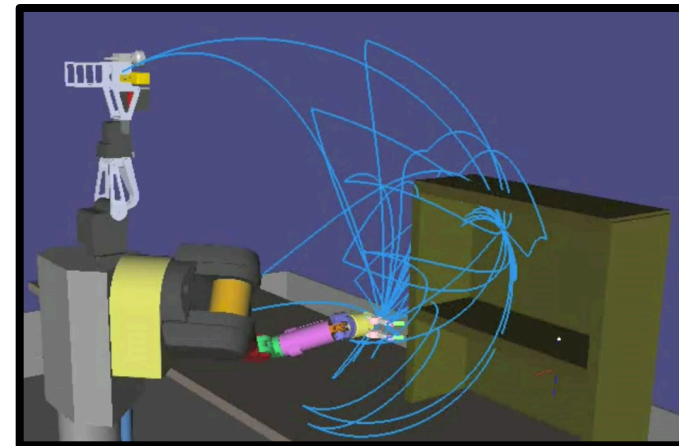
- Train policy to mimic greedy:

- $\pi(s) \rightarrow a$

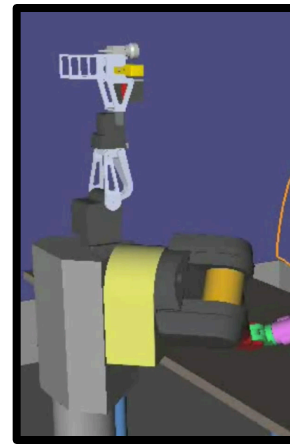
State  $s = (\Psi, x)$

$$\underset{\Psi: |\Psi| \leq B}{\operatorname{argmax}} F_x(\Psi)$$

Context / Environment      Submodular      Selection



Dictionary of Trajectories



Selected

# Example #2: Learning Greedy Algorithms (discret

## Contextual Submodular Maximization:

- Greedy Sequential Selection:

- $\Psi \leftarrow \Psi \oplus \underset{a}{\operatorname{argmax}} F_x(\Psi \oplus a)$

Not Available at Test Time

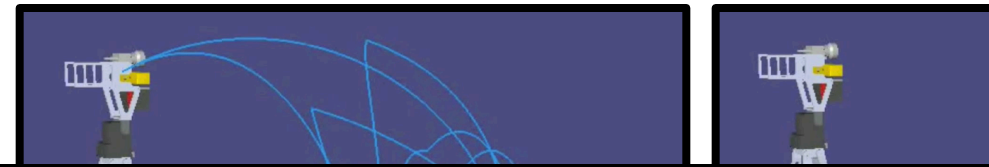
- Train policy to mimic greedy:

- $\pi(s) \rightarrow a$

State  $s = (\Psi, x)$

$$\underset{\Psi: |\Psi| \leq B}{\operatorname{argmax}} F_x(\Psi)$$

Context / Environment      Submodular      Selection

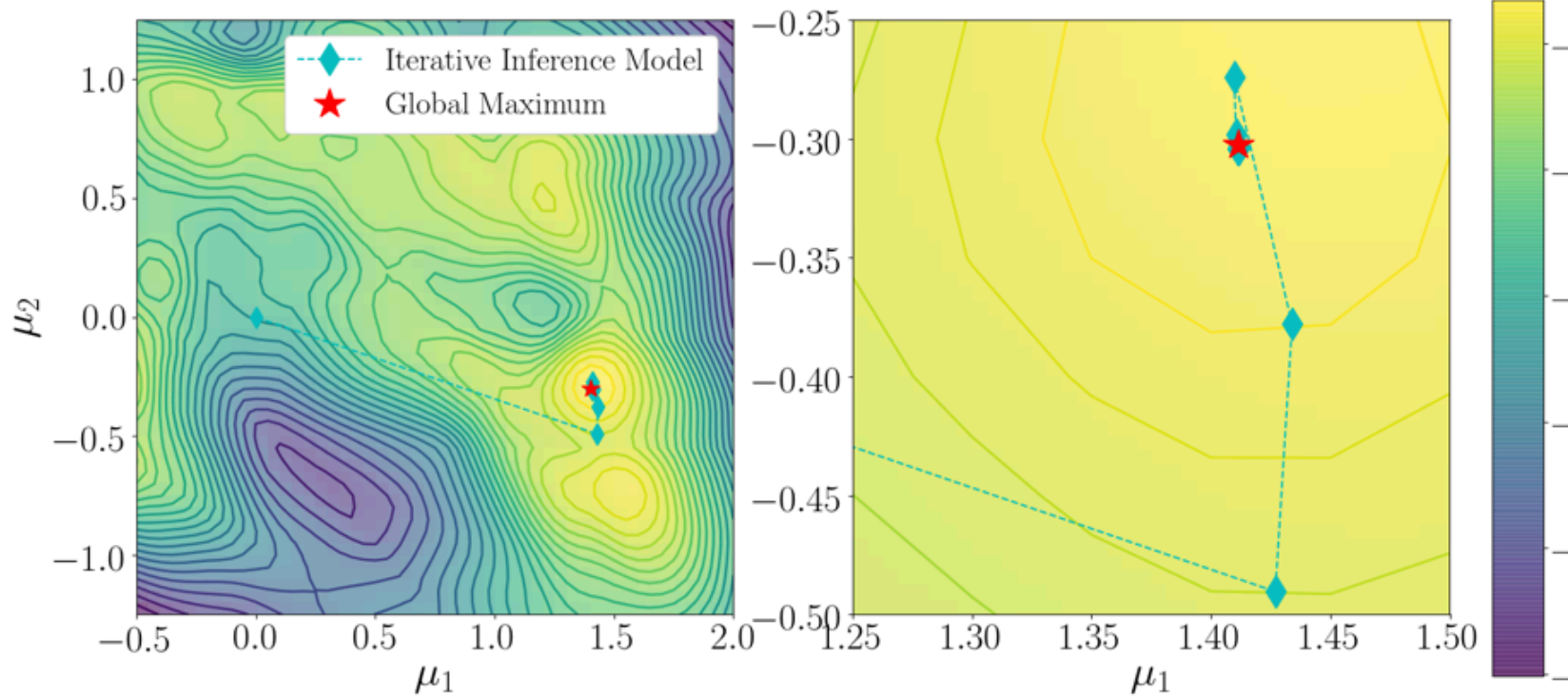


- Deterministic State Transitions
- Massive State Space
- Dense Rewards
- Note: Not Learning Submodular

# Example #3: Iterative Amortized Inference (cont)

## Gradient Descent Style Updates:

- State = description of problem & current parameters
- Action = next point



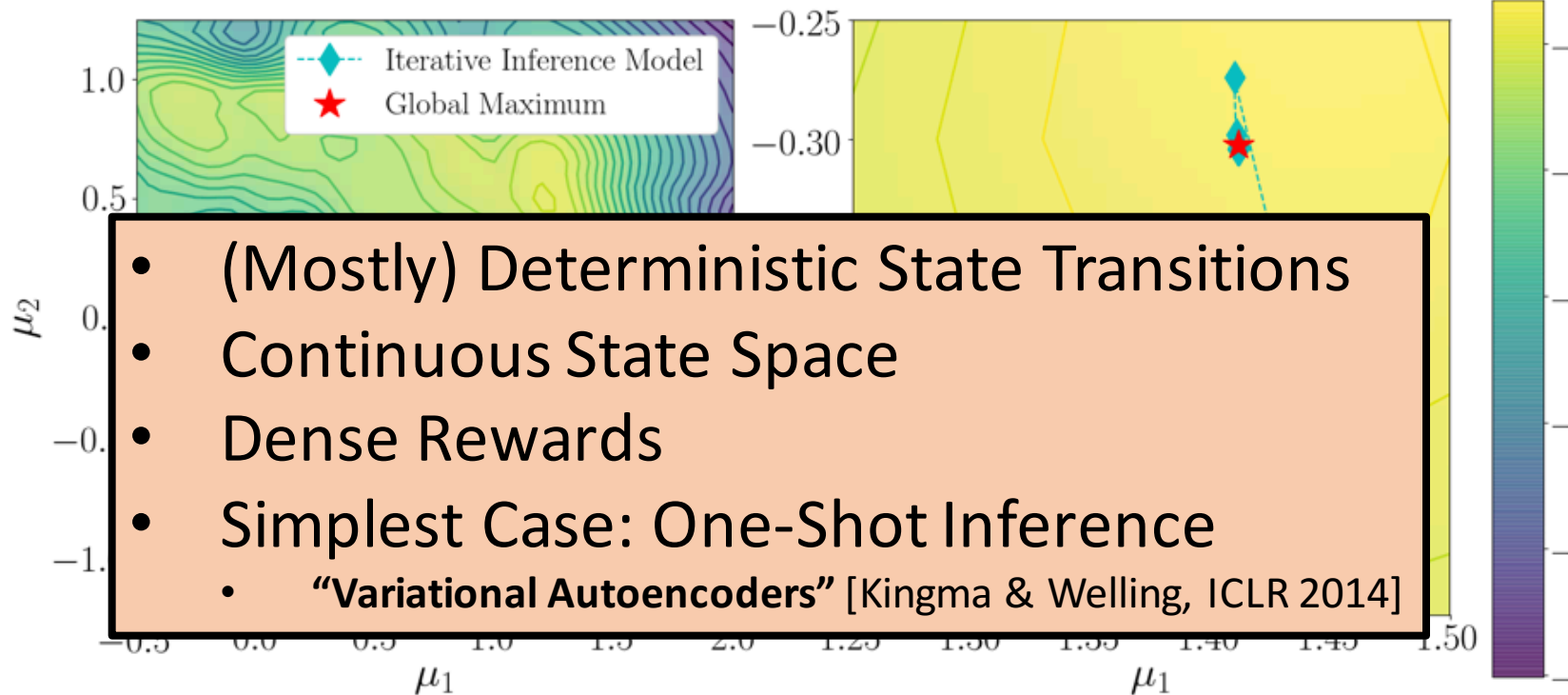
Useful for Accelerating Variational Inference

Iterative Amortized Inference, Joe Marino, Yisong Yue, Stephan Mandt. ICML 2018

# Example #3: Iterative Amortized Inference (cont)

## Gradient Descent Style Updates:

- State = description of problem & current state
- Action = next point



- (Mostly) Deterministic State Transitions
- Continuous State Space
- Dense Rewards
- Simplest Case: One-Shot Inference
  - “Variational Autoencoders” [Kingma & Welling, ICLR 2014]

Useful for Accelerating Variational Inference

# Optimization as Sequential Decision Making

## Learning to Search

- Discrete Optimization (Tree Search), Sparse Rewards
- **Learning to Search via Retrospective Imitation** [arXiv]
- **Co-training for Policy Learning** [UAI 2019]

## Contextual Submodular Maximization

- Discrete Optimization (Greedy), Dense Rewards
- **Learning Policies for Contextual Submodular Prediction** [ICML 2013]

## Learning to Infer

- Continuous Optimization (Gradient-style), Dense Rewards
- **Iterative Amortized Inference** [ICML 2018]
- **A General Method for Amortizing Variational Filtering** [NeurIPS 2018]



Jialin



Stepha



Joe M

# Optimization as Sequential Decision Making

## Learning to Search

- Discrete Optimization (Tree Search), Sparse Rewards
- **Learning to Search via Retrospective Imitation** [arXiv]
- **Co-training for Policy Learning** [UAI 2019]

## Contextual Submodular Maximization

- Discrete Optimization (Greedy), Dense Rewards
- **Learning Policies for Contextual Submodular Prediction** [ICML 2013]

## Learning to Infer

- Continuous Optimization (Gradient-style), Dense Rewards
- **Iterative Amortized Inference** [ICML 2018]
- **A General Method for Amortizing Variational Filtering** [NeurIPS 2018]



Jialin



Stepha

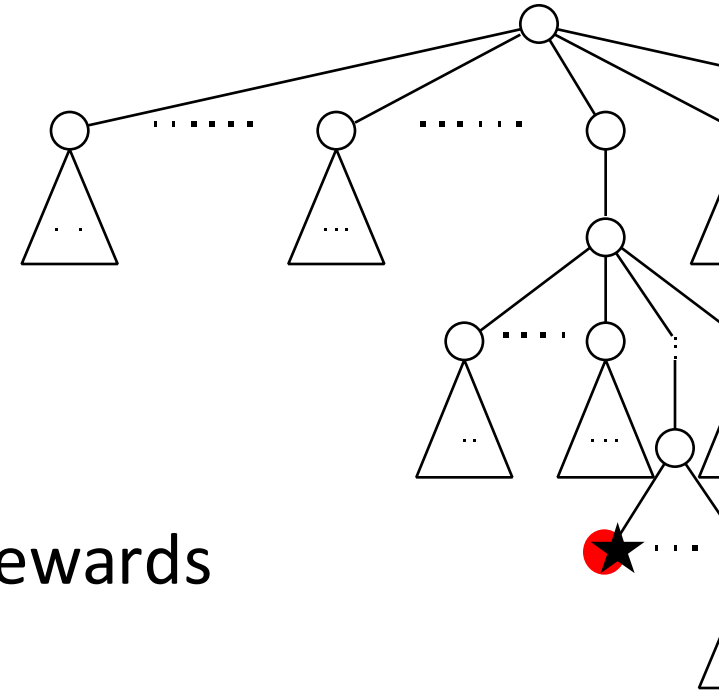


Joe M



# Learning to Optimize for Tree Search

- Idea #1: Treat as Standard RL
- Randomly explore for high rewards
  - **Very hard exploration problem!**
- Issues: massive state space & sparse rewards



# Learning to Optimize for Tree Search

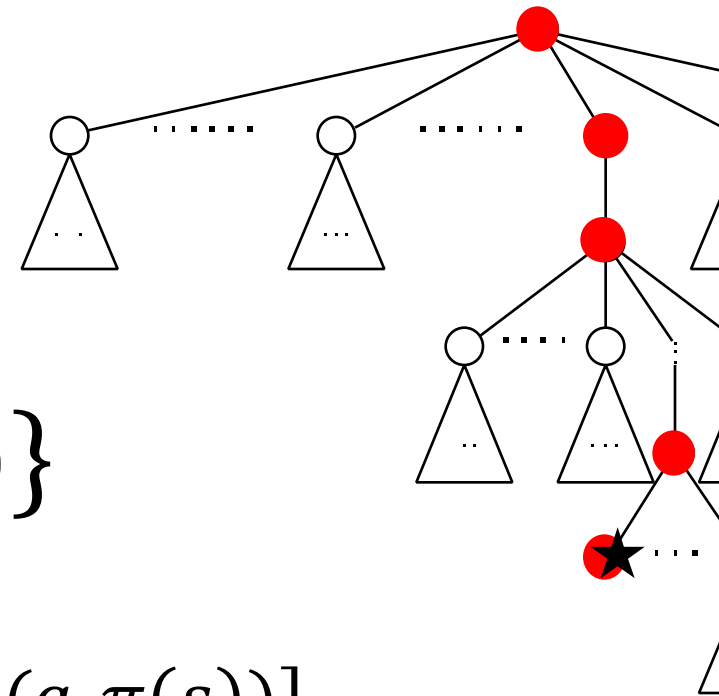
- Idea #2: Treat as Standard IL
- Convert to Supervised Learning
  - Assume access to solved instances

“Demonstration Data”

- Training Data:  $D_0 = \left\{ \left( \begin{array}{c} \text{tree} \\ \text{tree} \end{array}, \begin{array}{c} \text{tree} \\ \text{tree} \end{array} \right) \right\}$

- Basic IL:  $\underset{\pi \in \Pi}{\operatorname{argmin}} L_{D_0}(\pi) \equiv E_{(s,a) \sim D_0} [\ell(a, \pi(s))]$

Behavioral Cloning



# Challenges w/ Imitation Learning

- Issues with Behavioral Cloning
  - Minimize  $L_{D_0}$  ... implications?
  - If  $\pi$  makes a mistake early, subsequent state distribution  $\approx D_0$  ??
  - Some extensions to Interactive IL [He et al., NeurIPS 2014]

**Our Approach is also Interactive IL**

- Demonstrations not Available on Large Problems
  - How to (formally) bootstrap from smaller problems?
  - Bridging the gap between IL & RL

**Our Approach gives one solution**

# Retrospective Imitation



Jialin  
Song

- Given:
  - Family of Distributions of Search problems
    - Family is parameterized by size/difficulty
  - Solved Instances on the Smallest/Easiest Instances
    - “Demonstrations”

Difficulty levels:  $k=1, \dots$

- Goal:
  - Interactive IL approach
  - Can Scale up from Smallest/Easiest Instances
  - Formal Guarantees

Connections to Curriculum  
& Transfer Learning

Learning to Search via Retrospective Imitation, Jialin Song, Ravi Lanka, et al., arXiv

# Retrospective Imitation

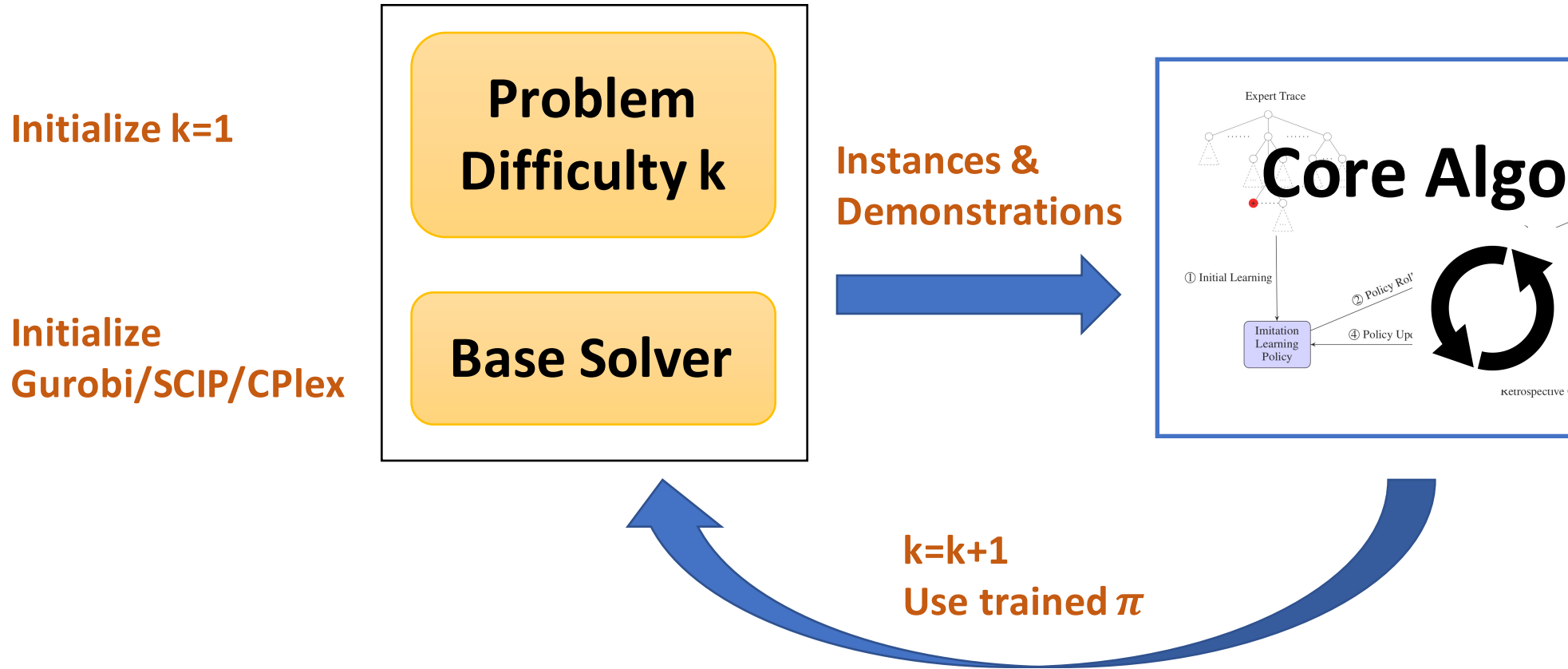
- Two-Stage Algorithm
- Core Algorithm
  - Fixed problem difficulty
  - Reductions to Supervised Learning
- Full Algorithm w/ Scaling Up
  - Uses Core Algorithm as Subroutine

Interactive IL w/ Sparse Environment

**Learning to Search via Retrospective Imitation**, Jialin Song, Ravi Lanka, et al., arXiv

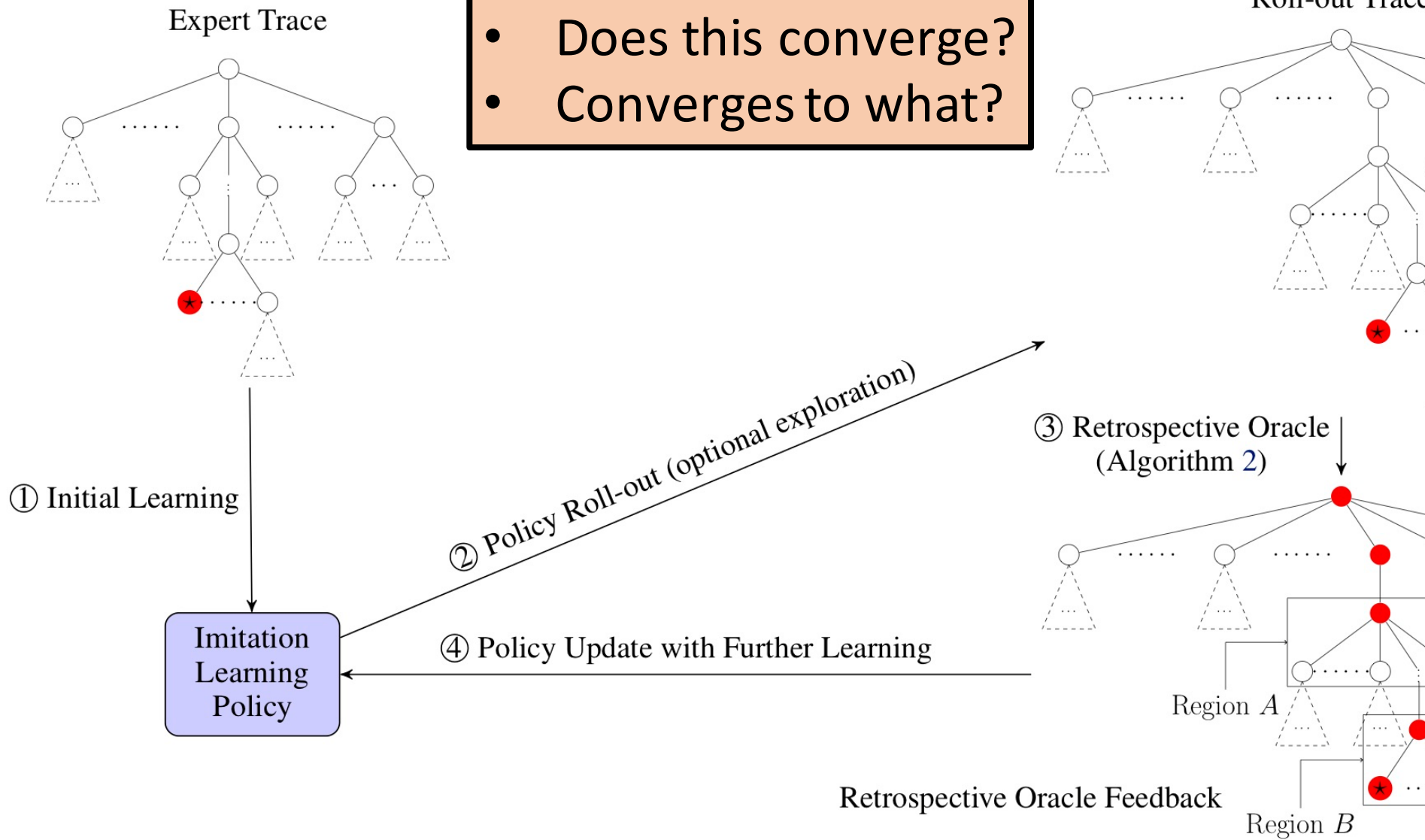


# Retrospective Imitation (Full Algorithm)



# Core Algorithm

- Does this converge?
- Converges to what?



Learning to Search via Retrospective Imitation, Jialin Song, Ravi Lanka, et al., arXiv



# Imitation Learning Tutorial (ICML 2018)

<https://sites.google.com/view/icml2018-imitation-learning/>

**Yisong Yue**



yyue@caltech.edu



@YisongYue



[yisongyue.com](http://yisongyue.com)

**Hoang M. Le**



hmle@caltech.edu

@HoangMinhLe

[hoangle.info](http://hoangle.info)

# Issues w/ Distribution Drift & Imitation S

- Demonstrations from initial Solver:  $D_0 = \left\{ \left( \begin{array}{c} \text{tree} \\ \text{tree} \end{array} \right) \right\}$

“correct” decision in this state

Which input states?

Correct relative to what?

- Supervised learning:  $\operatorname{argmin}_{\pi \in \Pi} L_{D_0}(\pi) \equiv E_{(s,a) \sim D_0} [\ell(a, \pi(s))]$

Oracle call to TensorFlow/PyTorch/etc...

If  $\pi$  achieves low error on  $D_0$ , so what?

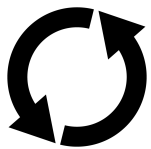
# Interactive Imitation Learning (Core Alg)

- First popularized by [Daume et al., 2009] [Ross et al., 2011]

- Basic idea:

- Train  $\pi_{i-1} = \operatorname{argmin}_{\pi \in \Pi} L_{D_{i-1}}(\pi)$  **Supervised Learning**

**i=i+1**



- Roll-out  $\pi_{i-1}$ , collect traces  $\{\tau\}$  **Run on instances**

- Demonstrator converts  $\{\tau\}$  into per-state feedback:  $\hat{D}_i$  **Depends on**

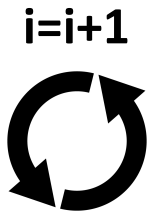
- $D_i = \hat{D}_i \cup D_{i-1}$  **Data aggregation**

**Search-based Structured Prediction**, Daume, Langford, Marcu, Machine Learning Journal 2009

**A Reduction of Imitation Learning and Structured Prediction to No-Regret Online Learning**, Ross, Gordon, Bagrodia

# Interactive Imitation

- First popularized by [Daume et al]
- Basic idea:
  - Train  $\pi_{i-1} = \operatorname{argmin}_{\pi \in \Pi} L_{D_{i-1}}(\pi)$



- Roll-out  $\pi_{i-1}$ , collect traces  $\{\tau\}$

- Demonstrator converts  $\{\tau\}$  into per-state feedback:  $\hat{D}_i$

- $D_i = \hat{D}_i \cup D_{i-1}$

Learns to Correct its Own Mis

Convergence Guarantees:

- $\sum_{i=0}^M L_{D_i}(\pi_i) \rightarrow \min_{\pi \in \Pi} \sum_{i=0}^M L_{D_i}(\pi)$
- Follow-the-Leader argument
- Also studied in [He et al., NeurIPS]

Requires defining “correct”

- Retrospective Oracle

Run on instances

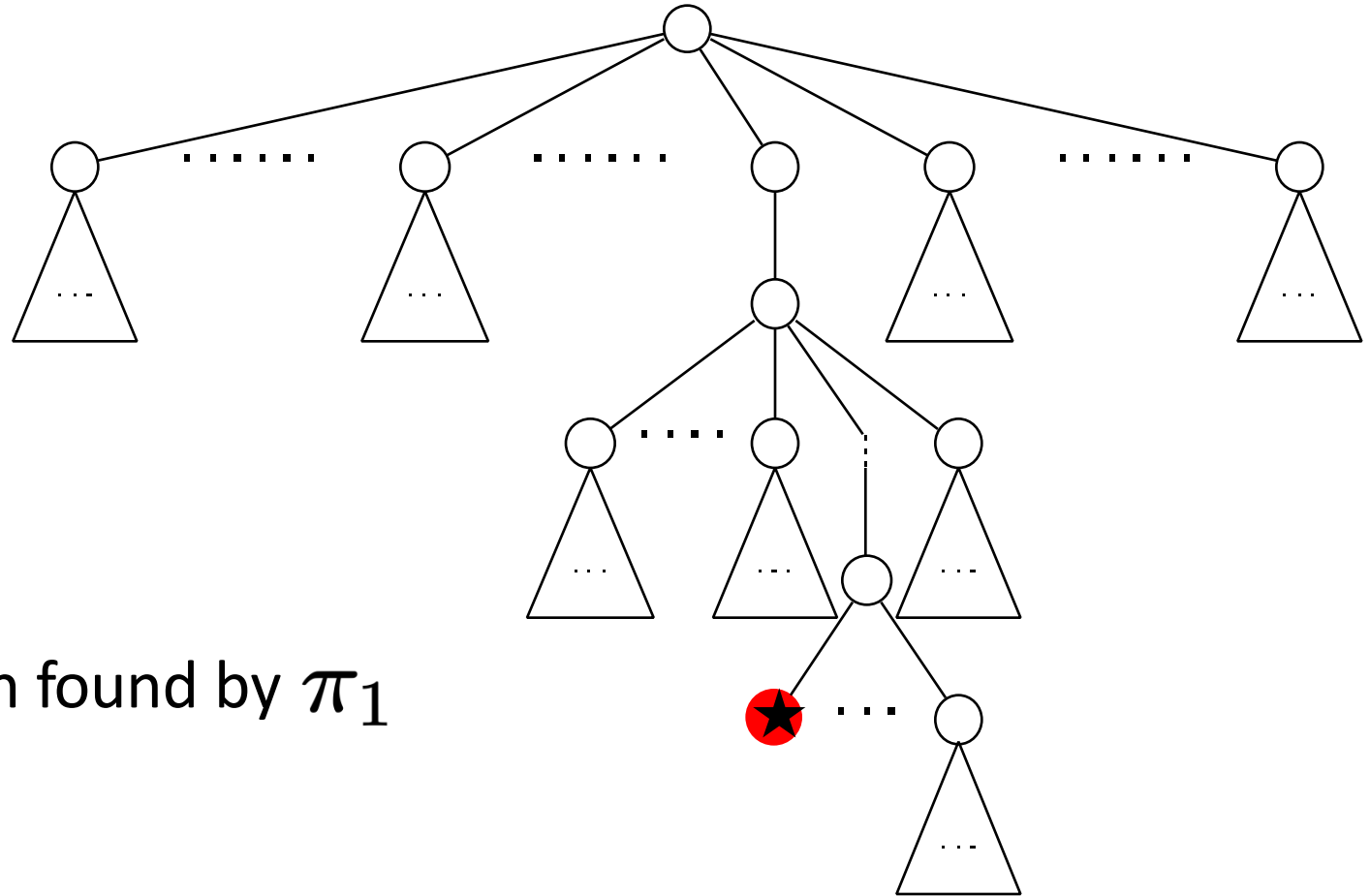
Depends on

Data aggregation

Search-based Structured Prediction, Daume, Langford, Marcu, Machine Learning Journal 2009

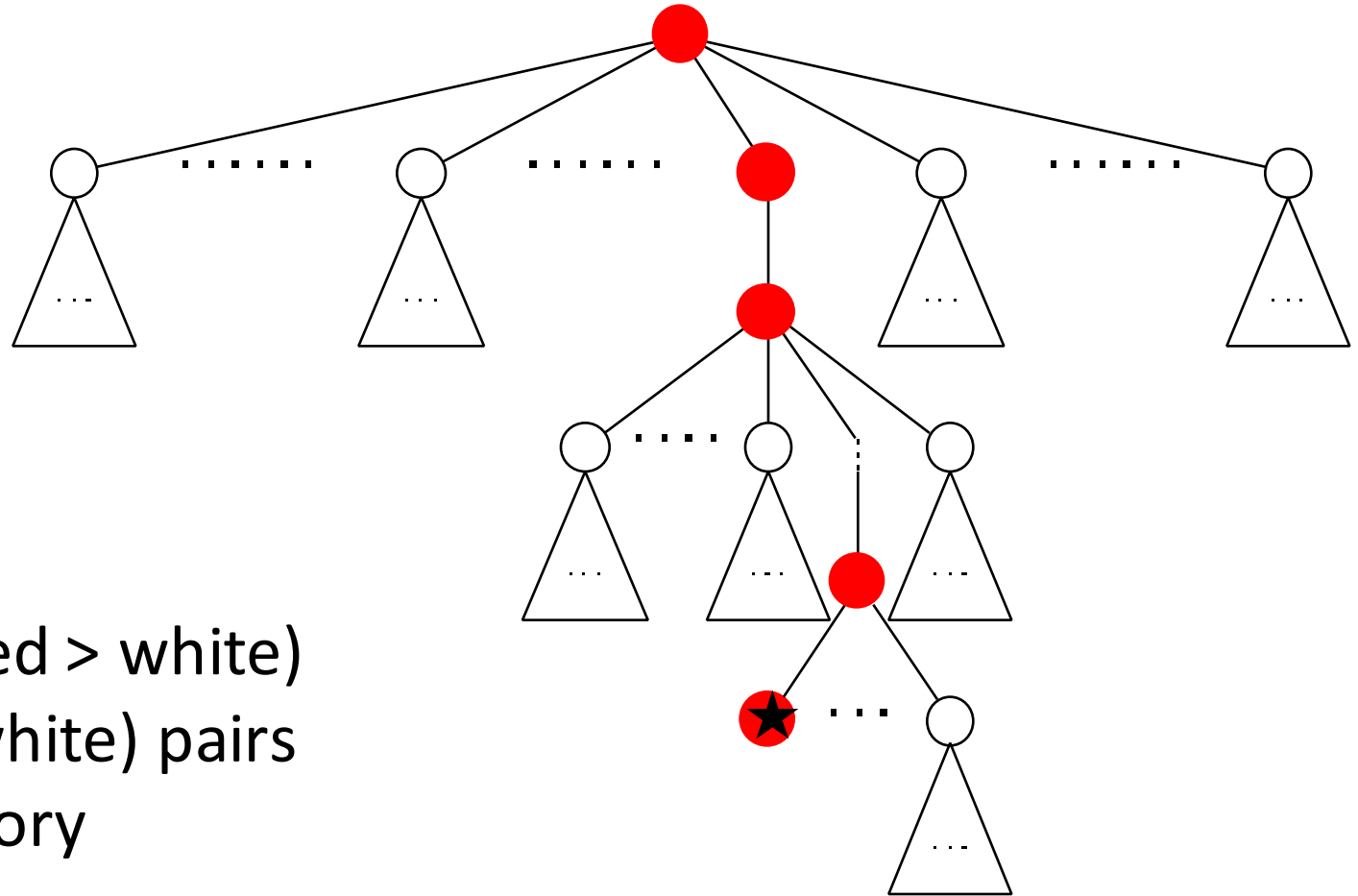
A Reduction of Imitation Learning and Structured Prediction to No-Regret Online Learning, Ross, Gordon, Bagn

# $\pi_1$ Policy Rollout



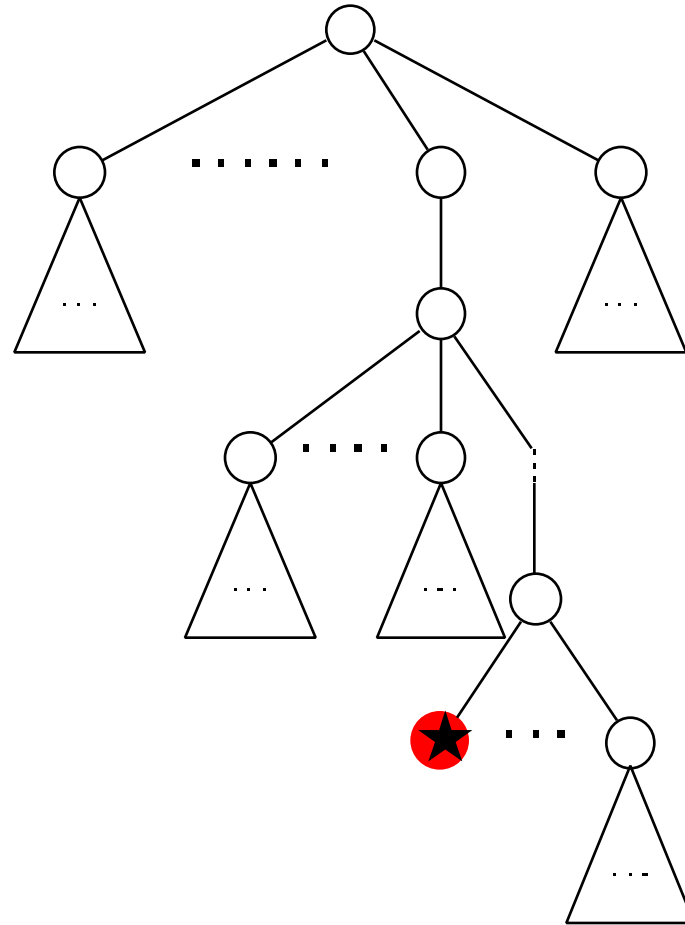
★: best solution found by  $\pi_1$

# Retrospective Oracle Feedback

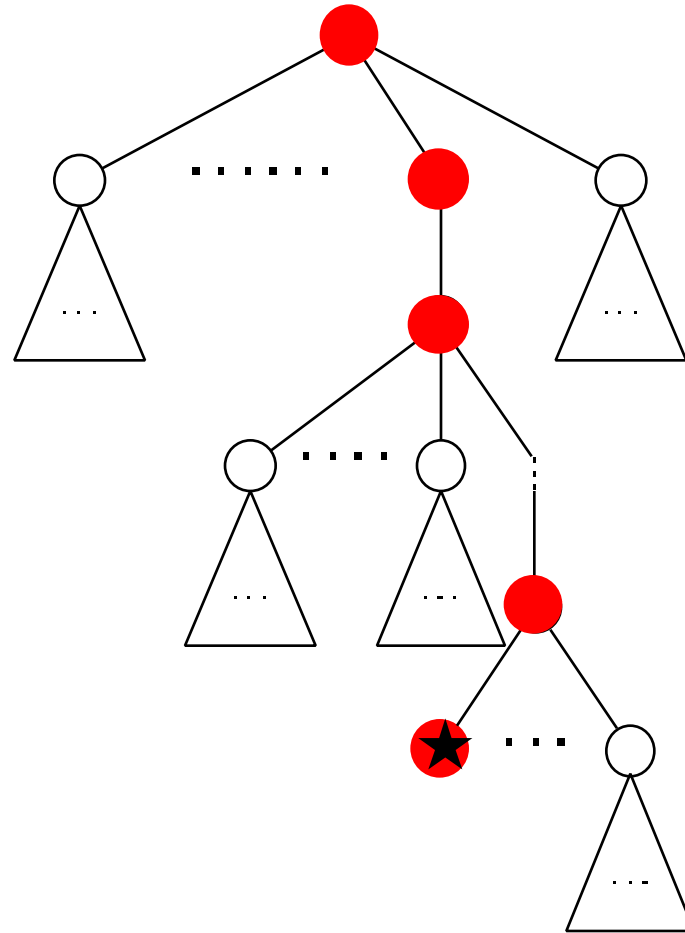


Feedback: (red > white)  
for all (red, white) pairs  
in the trajectory

# $\pi_2$ Policy Rollout



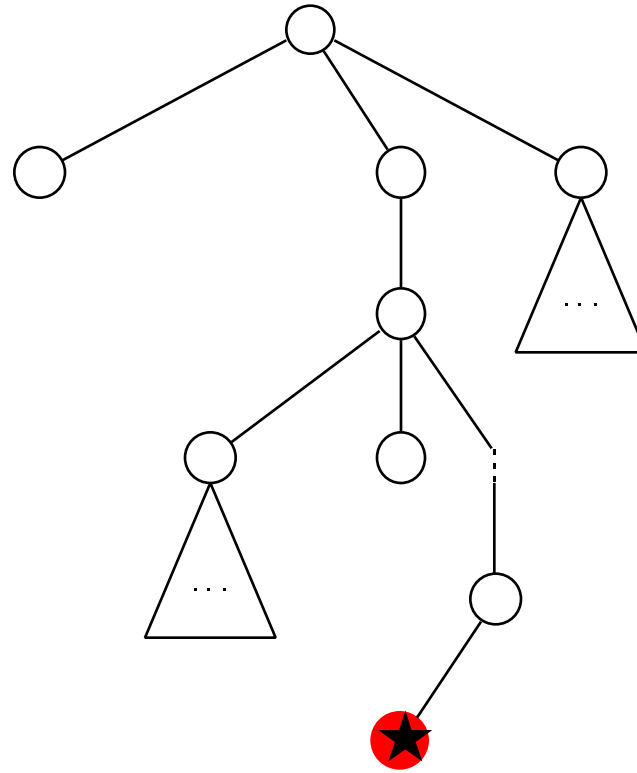
# Retrospective Oracle Feedback



Feedback: (red > white)  
for all (red, white) pairs  
in the trajectory



# $\pi_3$ Policy Rollout



# Core Algorithm Summary

- Sequence of Learning Reductions
- Leverages Retrospective Oracle to Define “Correct”
  - Relies on sparse environmental rewards
- Converges to near-optimal policy in class
  - Offloads computational challenges to Supervised Learning Oracle
- For supervised learning error  $\varepsilon$ :

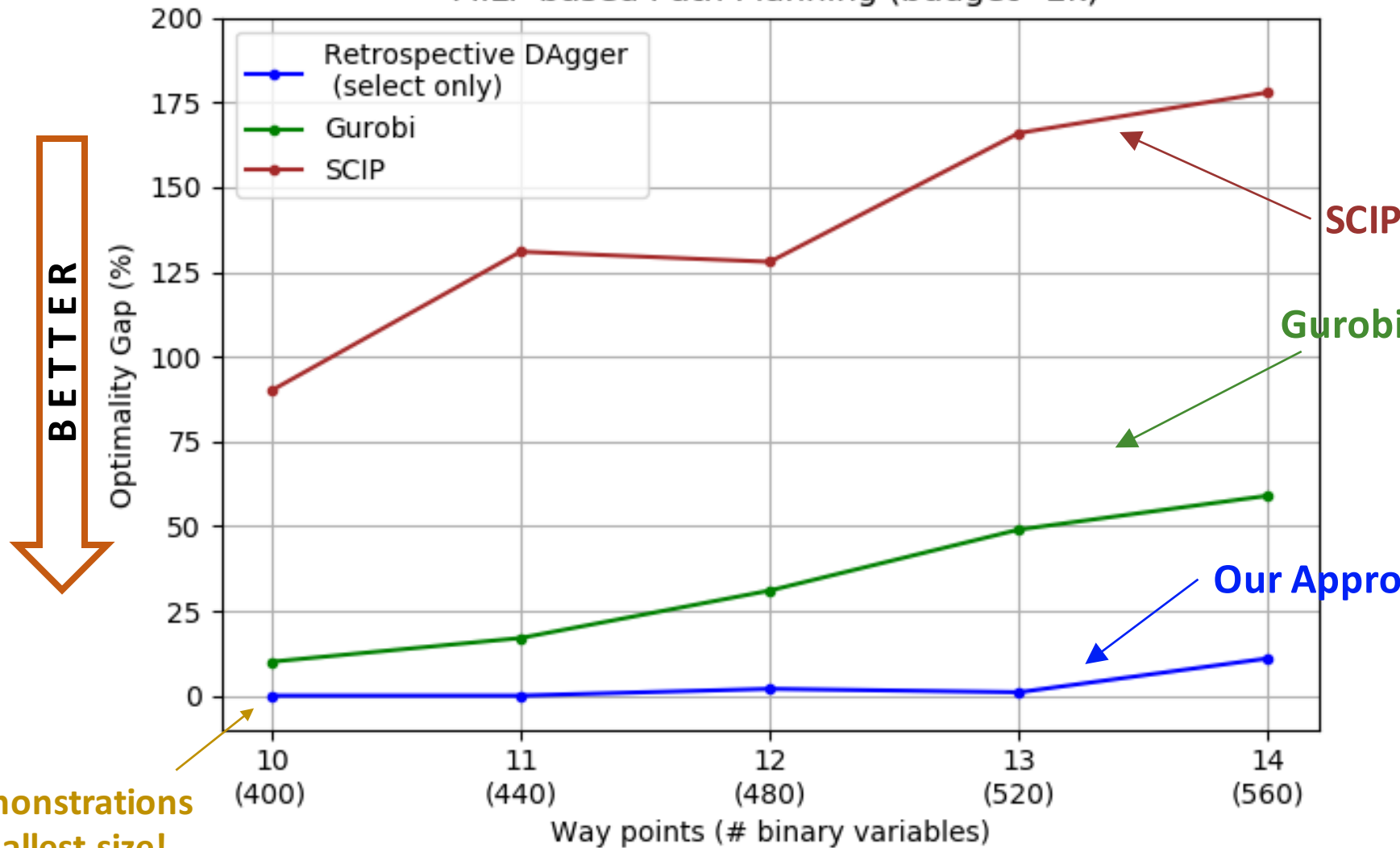
$$\text{Expected Search Length} = \frac{H^*}{1 - 2\varepsilon}$$

Optimal Search Length  
(typically # integers)

# Guarantees for Full Algorithm

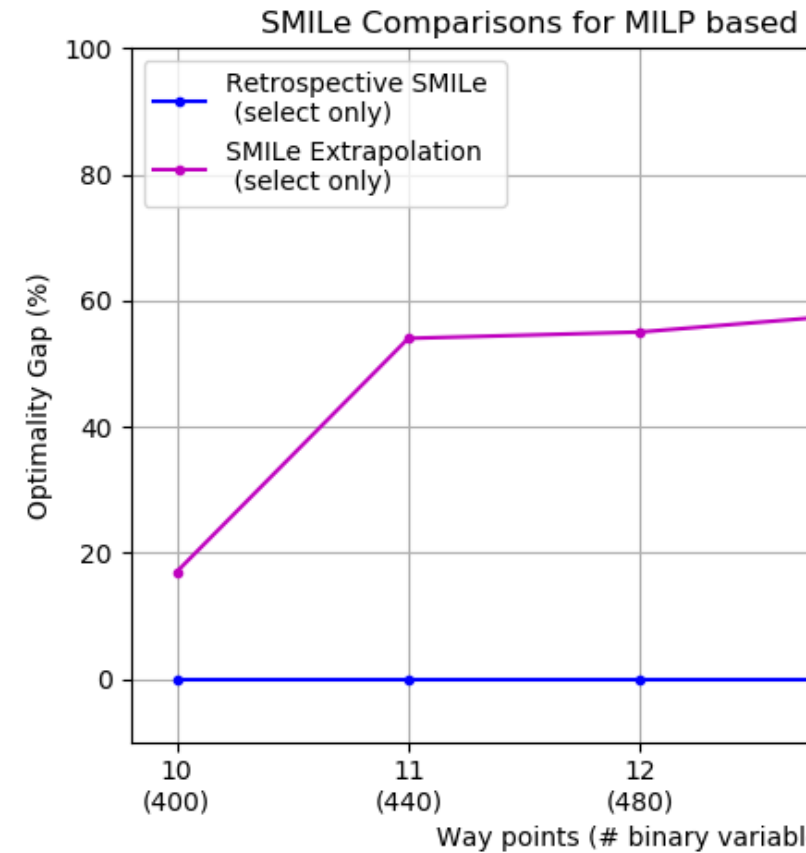
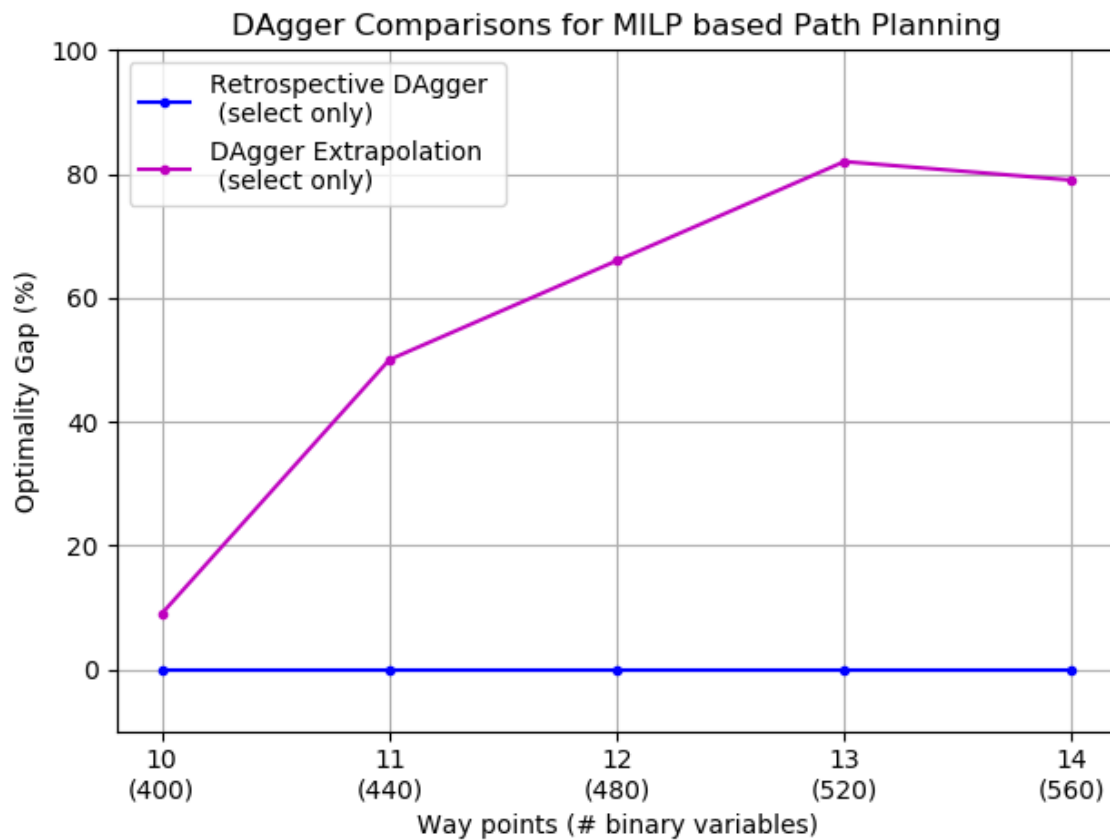
- Run  $\pi^k$  on problems of difficulty  $k+1$ 
  - Initial demonstrations for the harder problem instances
- **Suppose:** we could have run external solver on harder instances  
**Gurobi/SCIP/CPlex/ET**
- **Suppose:** search trace includes feasible solution of external solver
- Then  $\pi^k$  is as good as using original external solver!
  - (might take longer to converge)

# Retrospective DAgger vs Heuristics for MILP based Path Planning (budget=2k)



Learning to Search via Retrospective Imitation, Jialin Song, Ravi Lanka, et al., arXiv

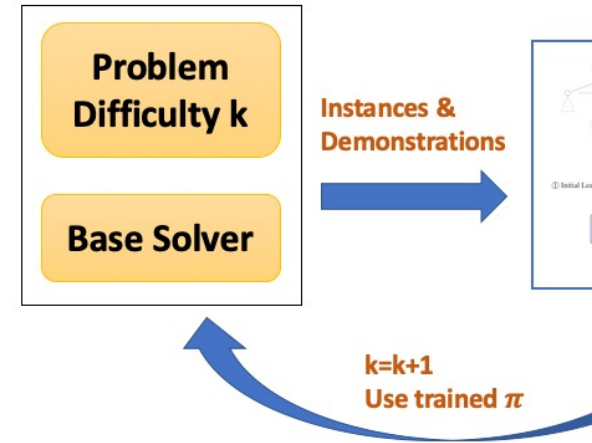
# Comparisons w/ Conventional IL



**Learning to Search via Retrospective Imitation**, Jialin Song, Ravi Lanka, et al., arXiv

# Retrospective Imitation

- Two-Stage Algorithm
  - Leverages Supervised Learning Oracle
- Initial demonstrations on small problems
- Exploits sparse environmental reward
  - “Retrospective Oracle”
- Iteratively scale up to harder problems



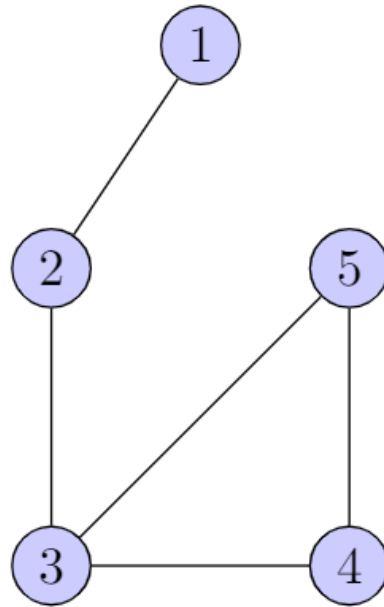
# Co-Training for Policy Learning

(Multiple Views)



Jialin  
Song

## Example: Minimum Vertex Cover



Graph View

[Khalil et al., 2017]

$$\max - \sum_{i=1}^5 x_i,$$

subject to:

$$x_1 + x_2 \geq 1,$$

$$x_2 + x_3 \geq 1,$$

$$x_3 + x_4 \geq 1,$$

$$x_3 + x_5 \geq 1,$$

$$x_4 + x_5 \geq 1,$$

$$x_i \in \{0, 1\}, \forall i \in \{1, \dots, 5\}$$

Integer Program View  
(Branch & Bound View)

[He et al., 2014]

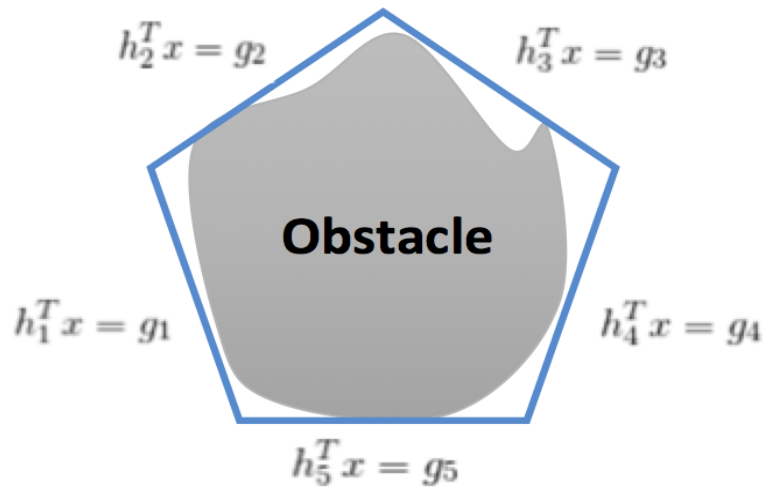
# Co-Training for Policy Learning

(Multiple Views)

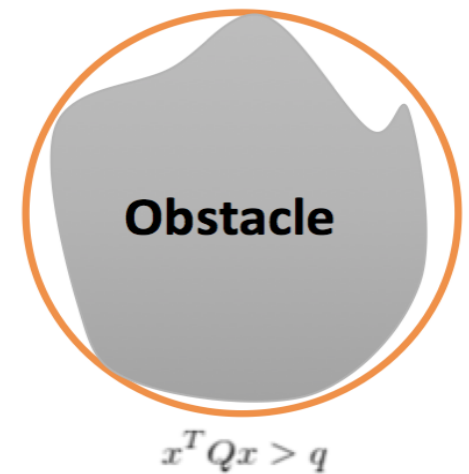


Jialin  
Song

## Example: Different Types of Integer Programs



ILP




QCQP



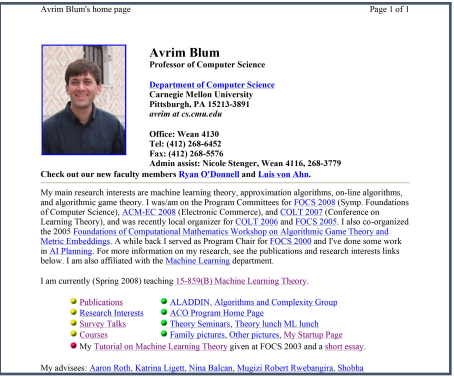
# Co-Training [Blum & Mitchell, 1998]

- Many learning problems have different sources of information
- Webpage Classification: Words vs Hyperlinks

Prof. Avrim Blum      My Advisor



Avrim Blum's home page      Page 1 of 1



**Avrim Blum**  
Professor of Computer Science  
Department of Computer Science  
Carnegie Mellon University  
Pittsburgh, PA 15213-3891  
avrim@cs.cmu.edu

Office: Wean 4130  
Tel: (412) 268-6452  
Fax: (412) 268-5576  
Admin assist: Nicole Stenger, Wean 4116, 268-3779  
Admin assist: Nicole Stenger, Wean 4116, 268-3779

Check out our new faculty members [Ryan D'Donnell](#) and [Lutz von Ahn](#).

My main research interests are machine learning theory, approximation algorithms, on-line algorithms, and algorithmic game theory. I was on the Program Committees for FOCS 2008 (Symp. Foundations of Computer Science), ACM-EC 2008 (Electronic Commerce), and COLT 2007 (Conference on Learning Theory), and was recently local organizer for COLT 2006 and FOCS 2005. I also co-organized the 2005 Foundations of Computational Mathematics Workshop on Algorithmic Game Theory and Metric Embeddings. A while back I served as Program Chair for FOCS 2000 and I've done some work in AI Planning. For more information on my research, see the publications and research interests links below. I am also affiliated with the Machine Learning department.

I am currently (Spring 2008) teaching 15-859(B) Machine Learning Theory.


- Publications
- Research Interests
- Survey Talks
- Courses
- ALADDIN, Algorithms and Complexity Group
- ACO Program Home Page
- Theory Seminars, Theory lunch ML lunch
- Family pictures, Other pictures, My Startup Page

My Tutorial on Machine Learning Theory given at FOCS 2003 and a short essay.

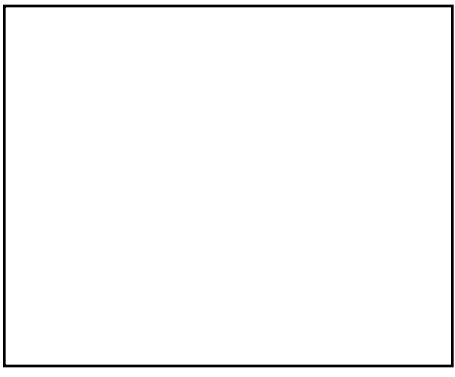
My advisees: Aaron Roth, Katrina Ligeti, Nina Balcan, Magzi Robert Rwehangira, Shobha

x - Link info & Text info

Prof. Avrim Blum      My Advisor




Avrim Blum's home page      Page 1 of 1




x<sub>1</sub> - Link info

Prof. Avrim Blum      My Advisor



Avrim Blum's home page      Page 1 of 1



**Avrim Blum**  
Professor of Computer Science  
Department of Computer Science  
Carnegie Mellon University  
Pittsburgh, PA 15213-3891  
avrim@cs.cmu.edu

Office: Wean 4130  
Tel: (412) 268-6452  
Fax: (412) 268-5576  
Admin assist: Nicole Stenger, Wean 4116, 268-3779  
Admin assist: Nicole Stenger, Wean 4116, 268-3779

Check out our new faculty members [Ryan D'Donnell](#) and [Lutz von Ahn](#).

My main research interests are machine learning theory, approximation algorithms, on-line algorithms, and algorithmic game theory. I was on the Program Committees for FOCS 2008 (Symp. Foundations of Computer Science), ACM-EC 2008 (Electronic Commerce), and COLT 2007 (Conference on Learning Theory), and was recently local organizer for COLT 2006 and FOCS 2005. I also co-organized the 2005 Foundations of Computational Mathematics Workshop on Algorithmic Game Theory and Metric Embeddings. A while back I served as Program Chair for FOCS 2000 and I've done some work in AI Planning. For more information on my research, see the publications and research interests links below. I am also affiliated with the Machine Learning department.

I am currently (Spring 2008) teaching 15-859(B) Machine Learning Theory.

- Publications
- Research Interests
- Survey Talks
- Courses
- ALADDIN, Algorithms and Complexity Group
- ACO Program Home Page
- Theory Seminars, Theory lunch ML lunch
- Family pictures, Other pictures, My Startup Page

My Tutorial on Machine Learning Theory given at FOCS 2003 and a short essay.

My advisees: Aaron Roth, Katrina Ligeti, Nina Balcan, Magzi Robert Rwehangira, Shobha

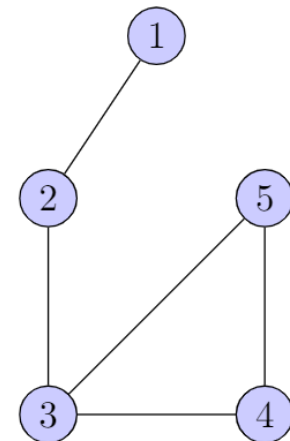
x<sub>2</sub> - Text info

(Taken from Nina Balcan's slides)

# What's Different about Policy Co-Training

- Sequential Decisions vs 1-Shot Decisions
- (Sparse) Environmental Feedback
  - Can collect more “labels”
- Different Action Spaces
  - Graph vs Branch-and-Bound

(Not always applicable)



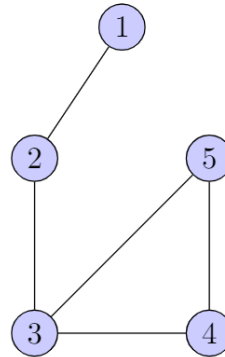
max -  
subject to  
 $x_1 + x_2$   
 $x_2 + x_3$   
 $x_3 + x_4$   
 $x_3 + x_5$   
 $x_4 + x_5$   
 $x_i \in \{0, 1\}$

# Intuition

- [1] “Learning combinatorial optimization algorithms over graphs”
- [2] “Learning to Search in Branch and Bound Algorithms” [He et al.]
- [3] “Learning to Search via Retrospective Imitation” [Song et al.]

MVC Instance

E.g., [1]



E.g., [2,3]

$$\max - \sum_{i=1}^5 x_i,$$

subject to:

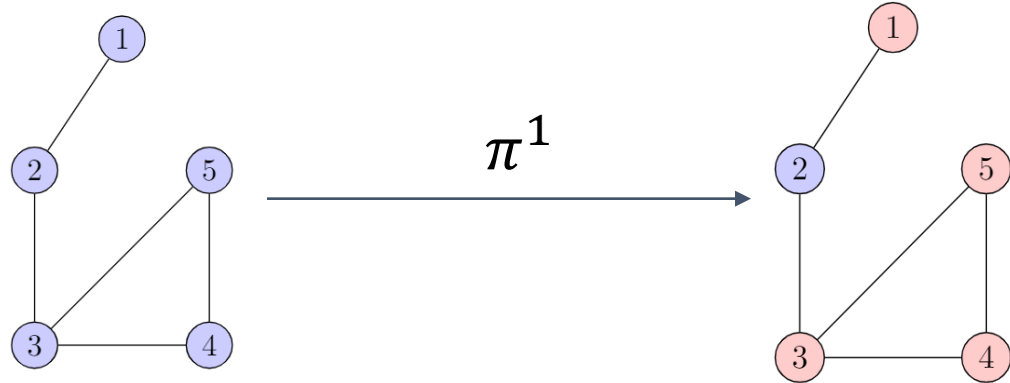
$$x_1 + x_2 \geq 1,$$
$$x_2 + x_3 \geq 1,$$
$$x_3 + x_4 \geq 1,$$
$$x_3 + x_5 \geq 1,$$
$$x_4 + x_5 \geq 1,$$
$$x_i \in \{0, 1\}, \forall i \in \{1, \dots, 5\}$$

# Intuition

- [1] "Learning combinatorial optimization algorithms over graph"
- [2] "Learning to Search in Branch and Bound Algorithms" [He et al.]
- [3] "Learning to Search via Retrospective Imitation" [Song et al.]

MVC Instance

E.g., [1]



E.g., [2,3]

$$\max - \sum_{i=1}^5 x_i,$$

subject to:

$$x_1 + x_2 \geq 1,$$

$$x_2 + x_3 \geq 1,$$

$$x_3 + x_4 \geq 1,$$

$$x_3 + x_5 \geq 1,$$

$$x_4 + x_5 \geq 1,$$

$$x_i \in \{0, 1\}, \forall i \in \{1, \dots, 5\}$$

$\pi^2$

$$x_1 = 0$$

$$x_2 = 1$$

$$x_3 = 1$$

$$x_4 = 1$$

$$x_5 = 0$$

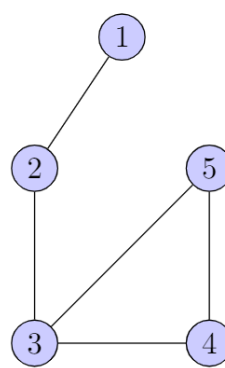
Better

# Intuition

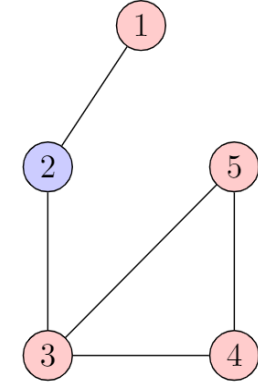
- [1] "Learning combinatorial optimization algorithms over graphs"
- [2] "Learning to Search in Branch and Bound Algorithms" [He et al.]
- [3] "Learning to Search via Retrospective Imitation" [Song et al.]

MVC Instance

E.g., [1]



$\pi^1$



E.g., [2,3]

$$\begin{aligned} & \max - \sum_{i=1}^5 x_i, \\ & \text{subject to:} \\ & x_1 + x_2 \geq 1, \\ & x_2 + x_3 \geq 1, \\ & x_3 + x_4 \geq 1, \\ & x_3 + x_5 \geq 1, \\ & x_4 + x_5 \geq 1, \\ & x_i \in \{0, 1\}, \forall i \in \{1, \dots, 5\} \end{aligned}$$

$\pi^2$

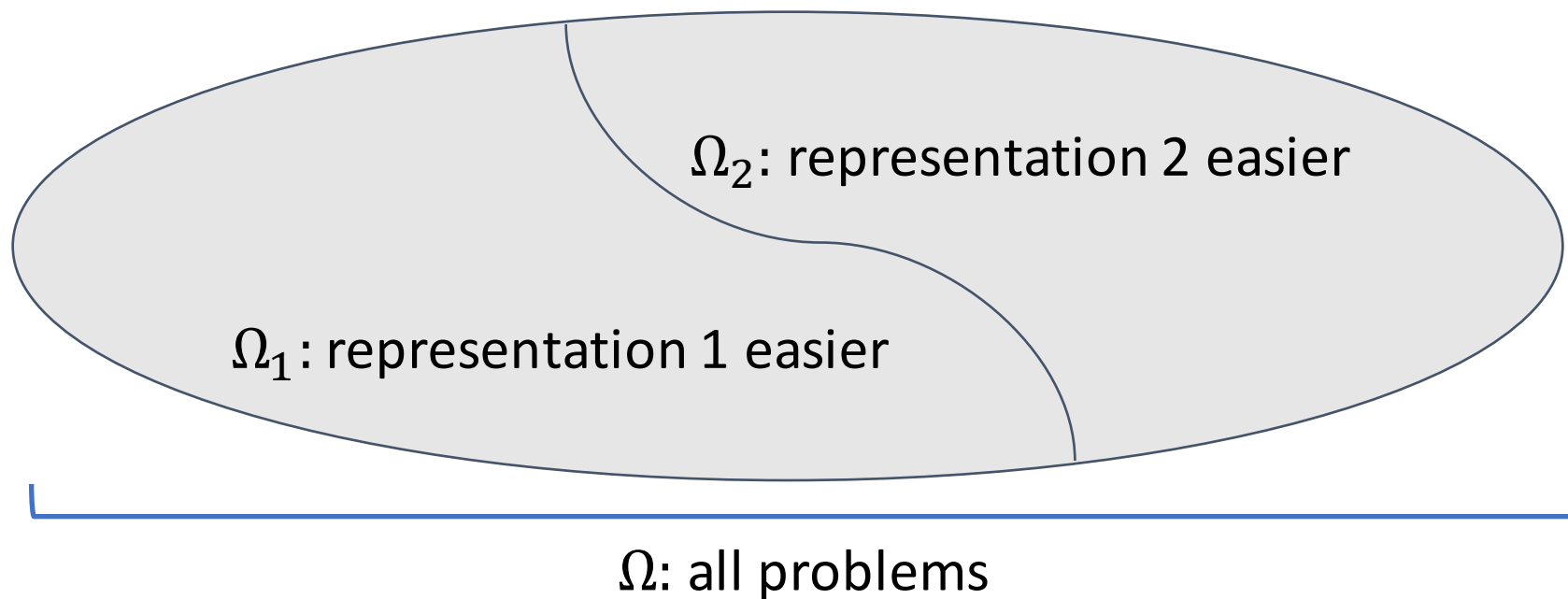
$$\begin{aligned} x_1 &= 0 \\ x_2 &= 1 \\ x_3 &= 1 \\ x_4 &= 1 \\ x_5 &= 0 \end{aligned}$$

Demonstration

Better

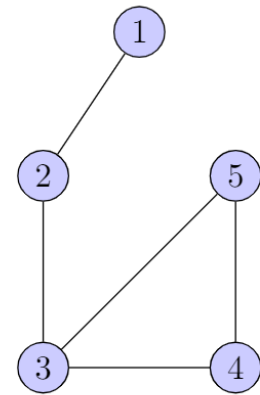
# Theoretical Insight

- Different representations differ in hardness
- Goal: quantify improvement



# (Towards) a Theory of Policy Co-Training

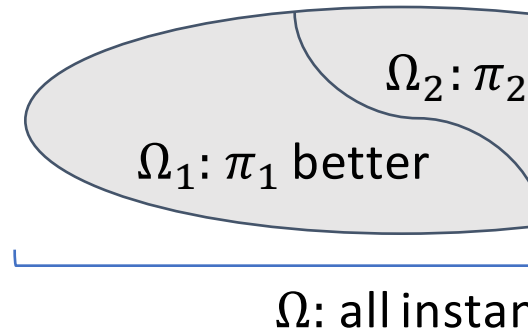
- Two MDP “views”:  $M^1$  &  $M^2$ 
  - $f^{1 \rightarrow 2}(\tau^1) \Rightarrow \tau^2$  (and vice versa)
    - “Trajectory” / “Rollout”
  - Realizing  $\tau^1$  on  $M^1 \Leftrightarrow$  realizing  $\tau^2$  on  $M^2$



max  
subject to  
 $x_1 +$   
 $x_2 +$   
 $x_3 +$   
 $x_3 +$   
 $x_4 +$   
 $x_i \in \{$

- **Question:** when does having two views/policies help?
  - Policy Improvement (next slide)
    - Builds upon [Kang et al., ICML 2018]
  - Optimality Gap for Shared Action Spaces (in paper)
    - Builds upon [DasGupta et al., NeurIPS 2002]

# Policy Improvement Bound



**Standard for Policy Gradient**

**1-step suboptimality of  $\pi^1$  on  $\Omega$**

**JS Divergence of  $\pi^2$  vs  $\pi^1$  on  $\Omega_2$**

**Want**

**KL Divergence of  $\pi^1$  vs  $\pi'^1$  on  $\Omega$**

**1-step suboptimal**

$$J(\pi'^1) \geq J_{\pi^1}(\pi'^1) - \frac{2\gamma(\alpha_{\Omega}^1 \varepsilon_{\Omega}^1 + 4\beta_{\Omega_2}^2 \varepsilon_{\Omega_2}^2)}{(1-\gamma)^2} + \delta_{\Omega_2}^2$$

**Performance of new policy (either RL or IL)**

**Approximation by sampling from  $\pi^1$**

**Discount**

**Performance Gap of  $\pi^2$  over  $\pi^1$  on  $\Omega_2$**   
 $J(\pi^2 | M \sim \Omega_2) - J(\pi^1 | M \sim \Omega_2)$

**Want to Maximize**

Builds upon theoretical results from [Kang et al., ICML 2018]

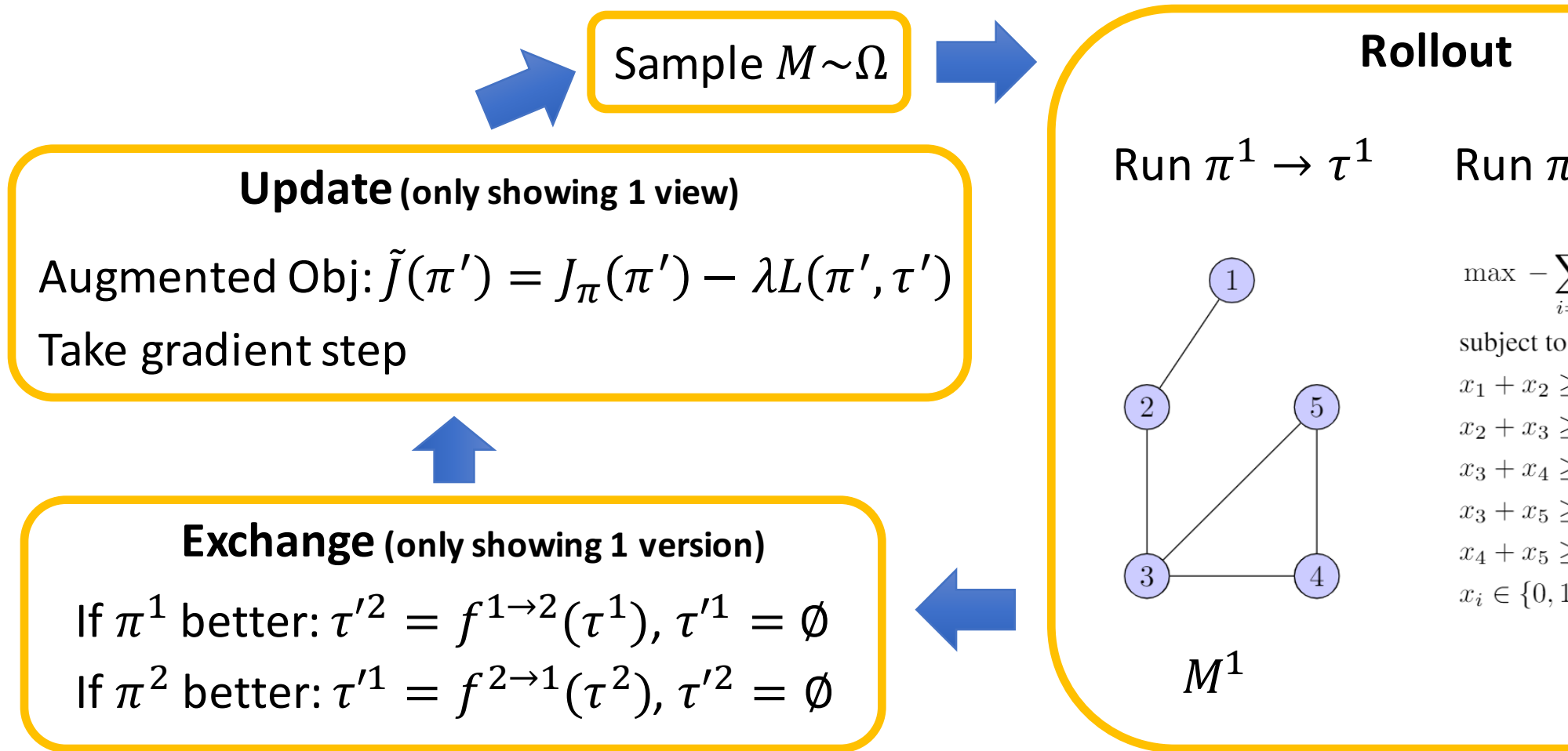


# Policy Improvement Bound (Summary)

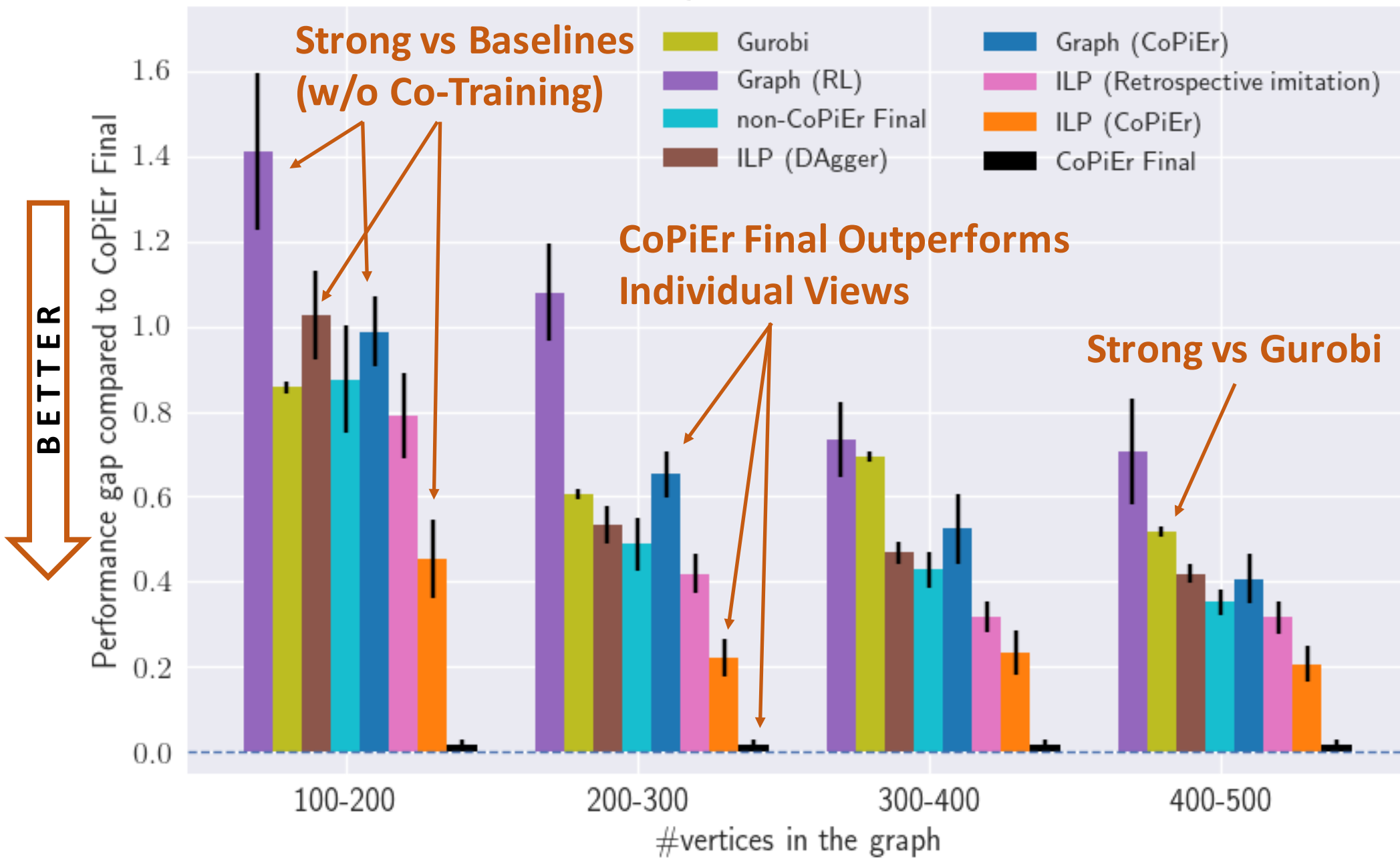
$$J(\pi'^1) \geq J_{\pi^1}(\pi'^1) - \frac{2\gamma(\alpha_{\Omega}^1 \varepsilon_{\Omega}^1 + 4\beta_{\Omega_2}^2 \varepsilon_{\Omega_2}^2)}{(1-\gamma)^2} + \delta_{\Omega}^2$$

- Minimizing  $\beta_{\Omega_2}^2 \rightarrow$  low disagreement between  $\pi^2$  vs  $\pi^1$
- Maximizing  $\delta_{\Omega_2}^2 \rightarrow$  high performance gap  $\pi^2$  over  $\pi^1$  on some

# CoPiEr Algorithm (Co-training for Policy Learning)



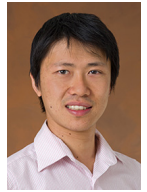
Performance comparison for Minimum Vertex Cover



# Ongoing: Integration with ENav



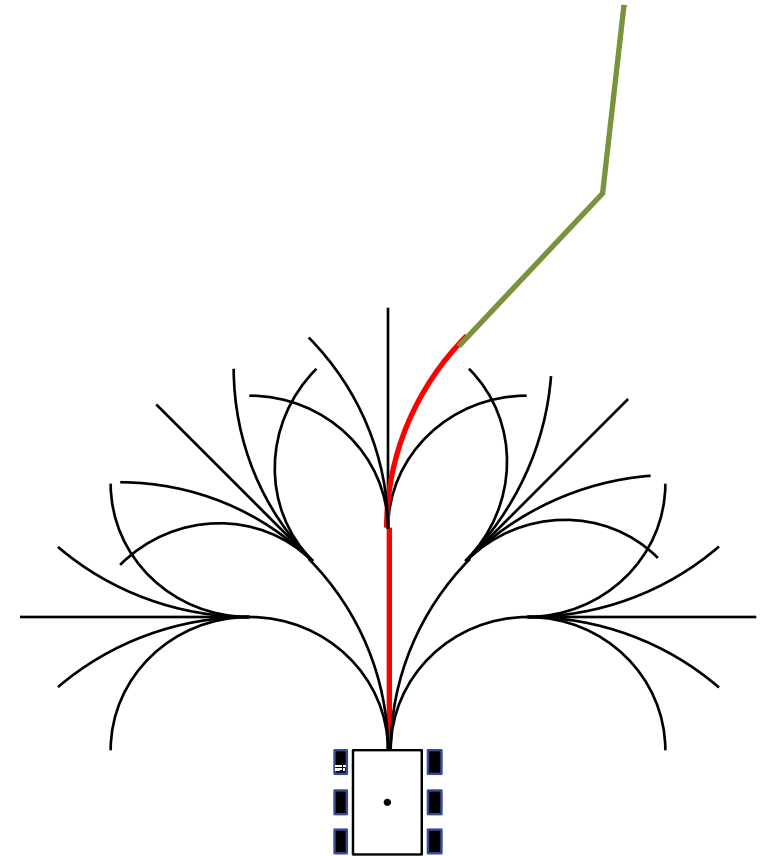
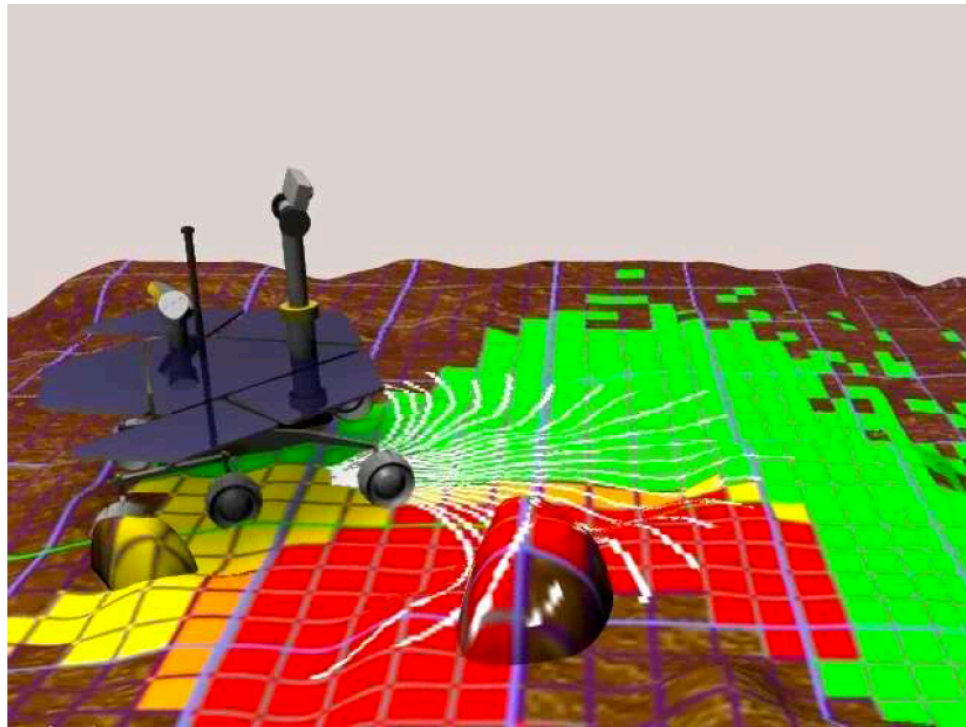
Ravi  
Lanka



Hiro  
Ono



O  
T



# Ongoing: Additive Manufacturing

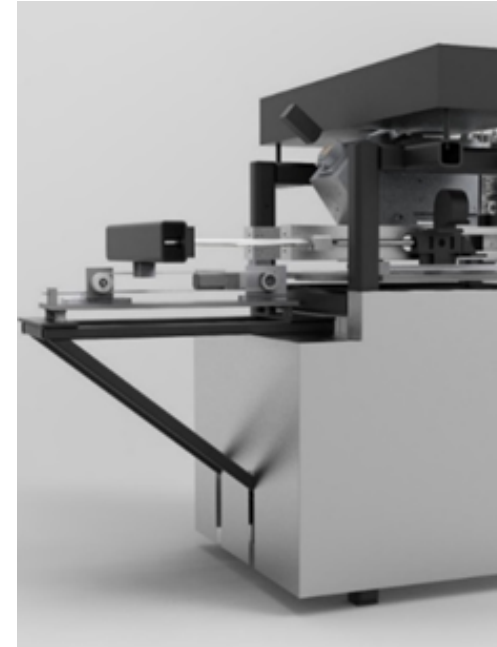
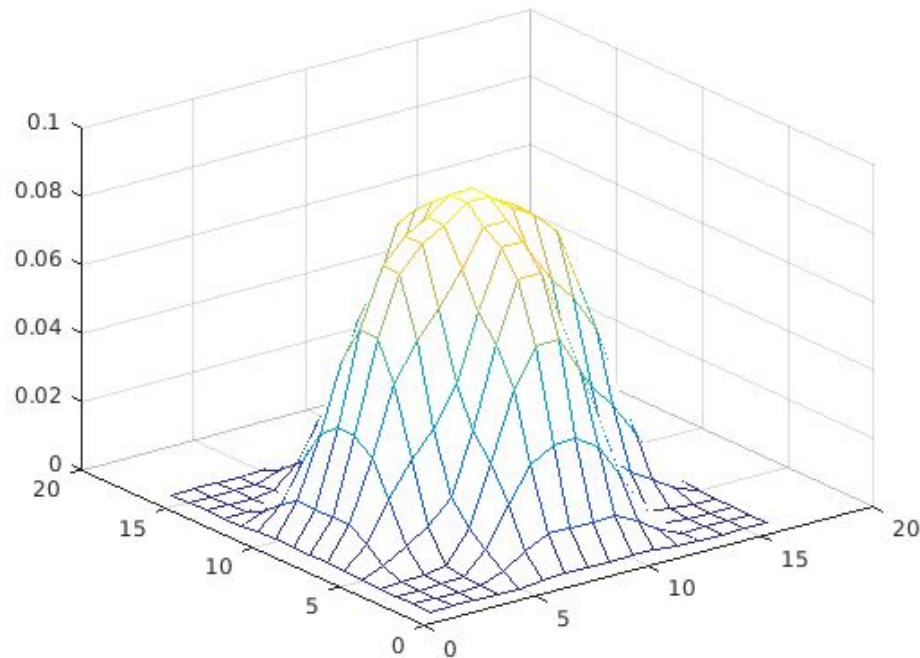


Stephanie  
Ding



Jialin  
Song

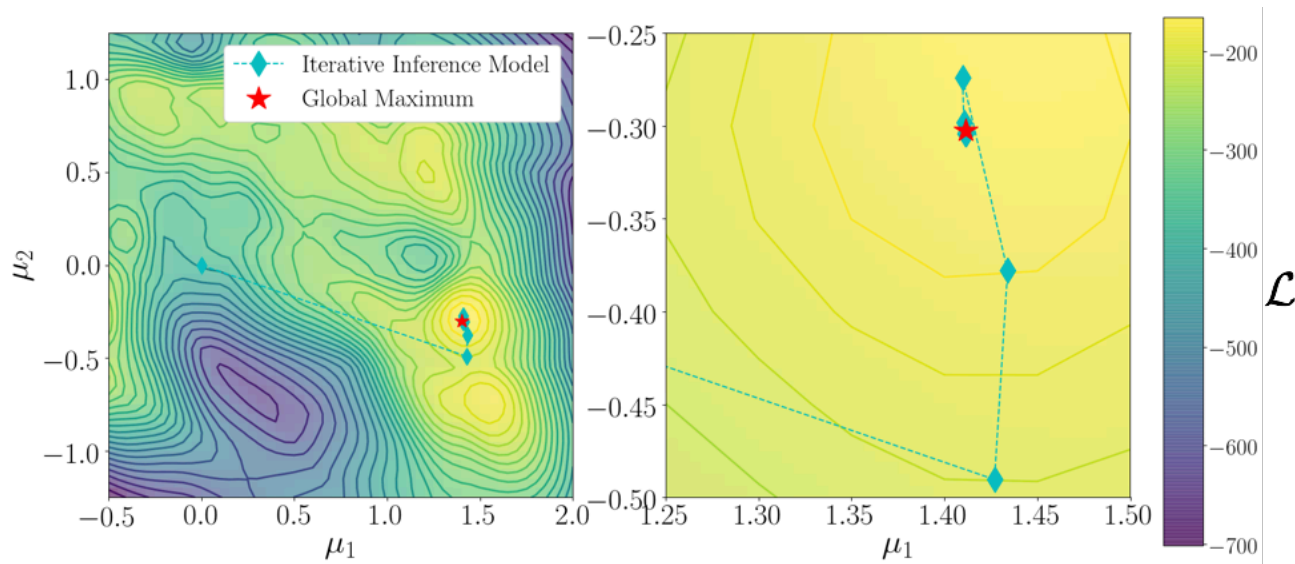
- Planning for 3D Inkjet Droplet Printing



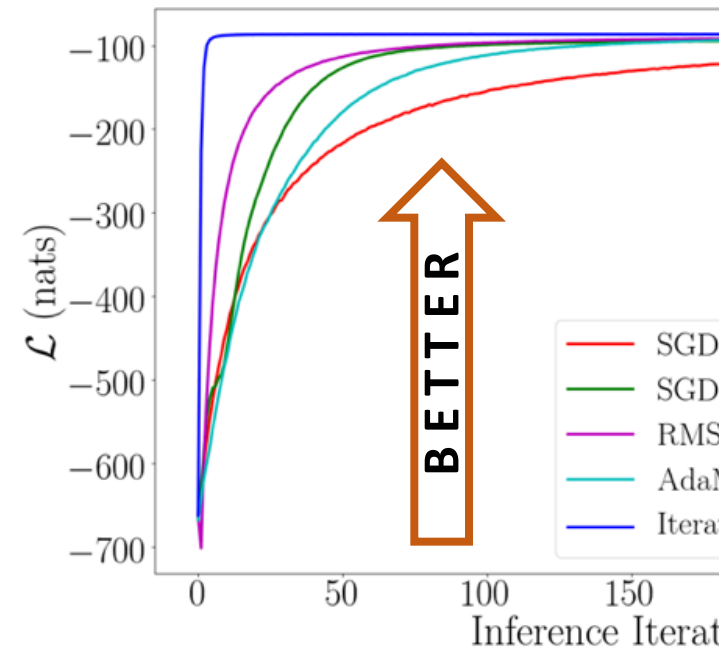
Rensselaer



# Iterative Amortized Inference (for Deep Probabilistic Models)



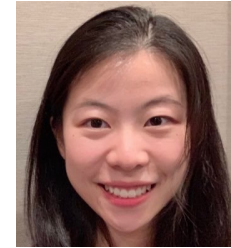
Related to “Learning to Learn” [Andychowicz et al., 2016]



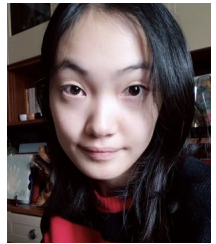
**Iterative Amortized Inference**, Joe Marino et al., ICML 2018

**A General Framework for Amortizing Variational Filtering**, Joe Marino et al, NeurIPS 2018

# Ongoing: Amortized Planning

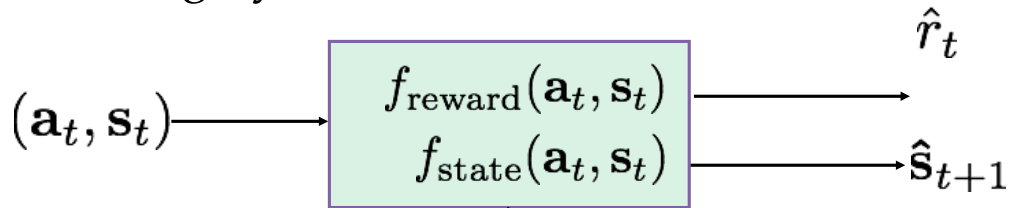


Yujia  
Huang

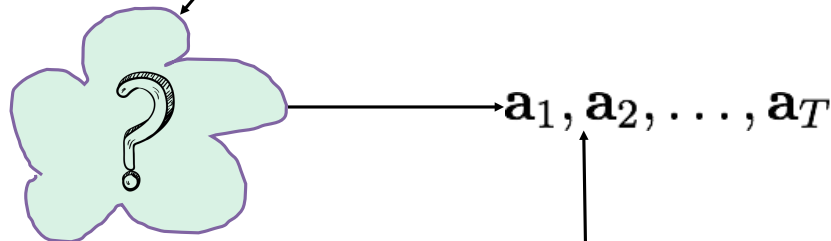


Sophie  
Dai

Learning dynamics:



Planning:



Optimize:

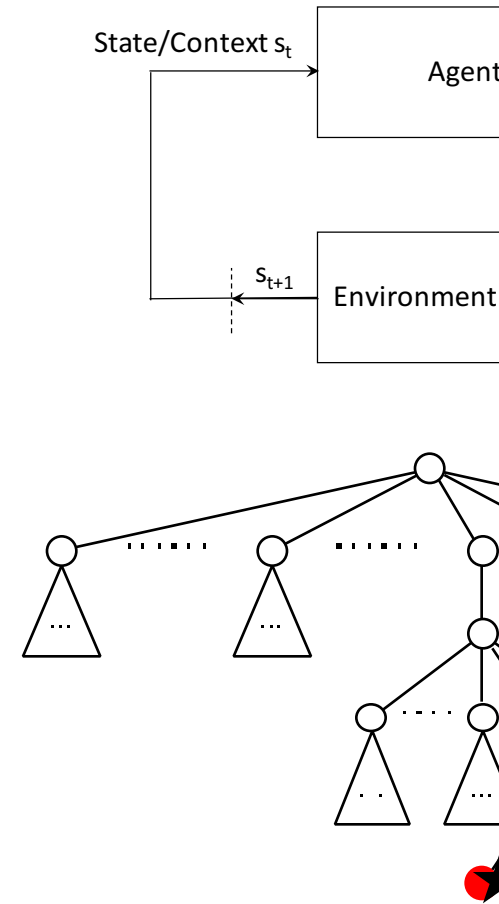
$$\max_{\mathbf{a}_1, \dots, \mathbf{a}_T} \sum_{t=1}^T f_{\text{reward}}(f_{\text{state}}(\hat{\mathbf{s}}_{t-1}, \mathbf{a}_{t-1}), \mathbf{a}_t)$$

Baseline: Gradient-based P

Can use (offline) training to

# Learning to Optimize as Policy Learning

- Optimization as Sequential Decision Making
- Formulate New Learning Problems
  - Builds upon RL/IL
- Interesting Algorithms
  - Theoretical Analysis/Guidance
  - Good Empirical Performance







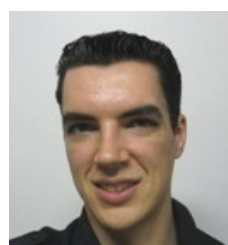
Jialin  
Song



Ravi  
Lanka



Joe  
Marino



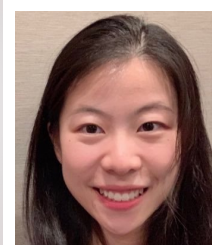
Stephane  
Ross



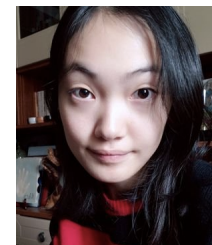
Aadyot  
Bhatnagar



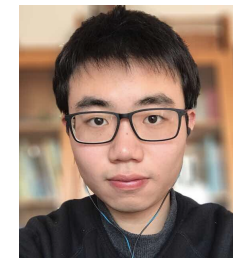
Albert  
Zhao



Yujia  
Huang



Sophie  
Dai



Hao  
Liu



Milan  
Cvitkovic



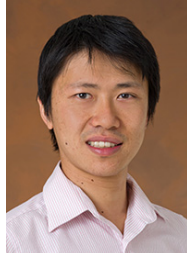
Robin  
Zhou



Debadeepta  
Dey



Stephan  
Mandt



Hiro  
Ono



Drew  
Bagnell



Uduak  
Inyang-Udoh



Sandipan  
Mishra



Olivier  
Toupet

**Learning to Search via Retrospective Imitation**, Jialin Song, Ravi Lanka, et al., arXiv

**Co-Training for Policy Learning**, Jialin Song, Ravi Lanka, et al., UAI 2019

**Learning Policies for Contextual Submodular Optimization**, Stephane Ross et al., ICML

**Iterative Amortized Inference**, Joe Marino et al., ICML 2018

**A General Framework for Amortizing Variational Filtering**, Joe Marino et al, NeurIPS 20

<https://github.com/ravi-lanka-4/CoPiEr>

[https://github.com/joelouismarino/iterative\\_inference](https://github.com/joelouismarino/iterative_inference)