

Differentiable Planners

Saurabh Gupta

Solving a RL Problem

Structured
Policies

Better Reward Signals

Sim2Real

Better Optimization

Convert into a
Supervised Training
Problem

Solve a Related but
Supervision-rich Problem

Build Models and Plan
with Them

Model-free RL
with sparse
rewards

Known reward,
known model.
Model-based RL



Structured Policies

Borrow policy structures from classical pipelines, and inject learning.

```
graph LR; A[Observations] --> B[State Estimation]; B --> C[Planning]; C --> D[Control];
```

Observations

State
Estimation

Planning

Control

Typical Robotics Pipeline

Observations



State
Estimation



Planning



Control

Manipulation



Observed Images



6DOF Pose



Grasp Motion
Planning

Observations



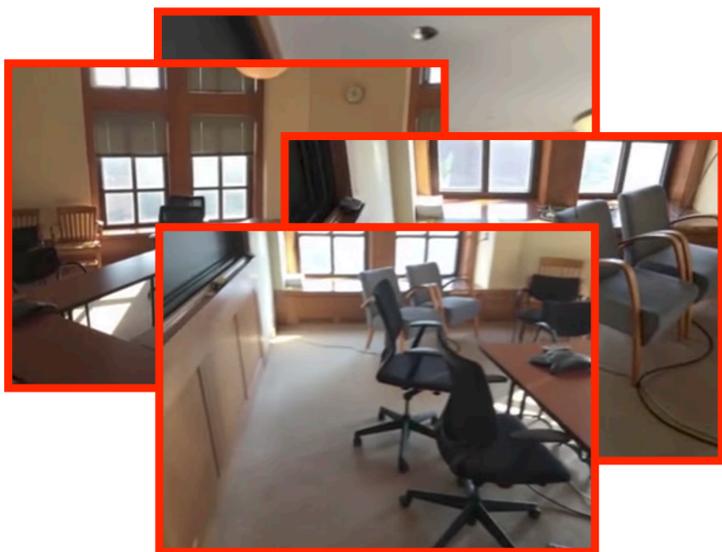
State Estimation



Planning



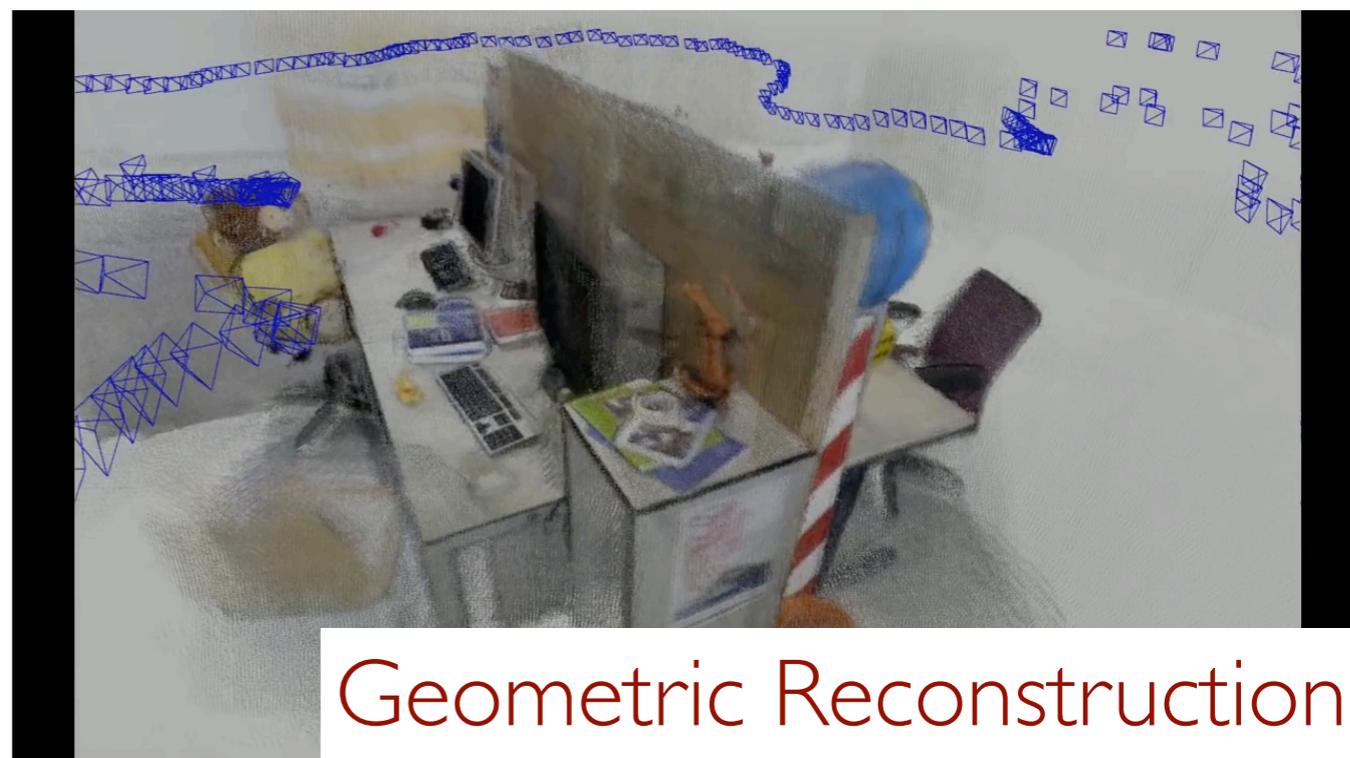
Control



Observed Images

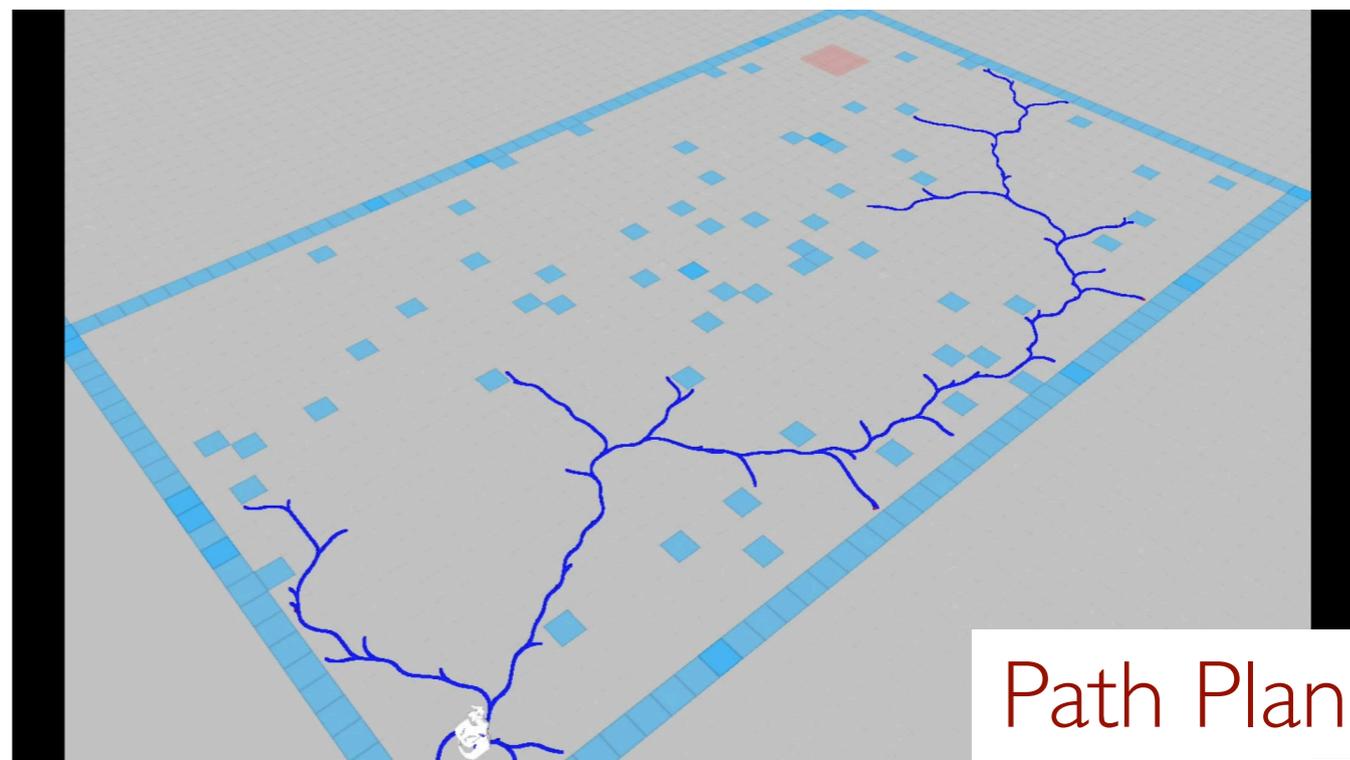
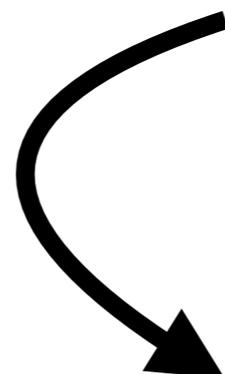


Mapping



Geometric Reconstruction

Planning



Path Plan

Hartley and Zisserman. 2000. Multiple View Geometry in Computer Vision

Thrun, Burgard, Fox. 2005. Probabilistic Robotics

Canny. 1988. The complexity of robot motion planning.

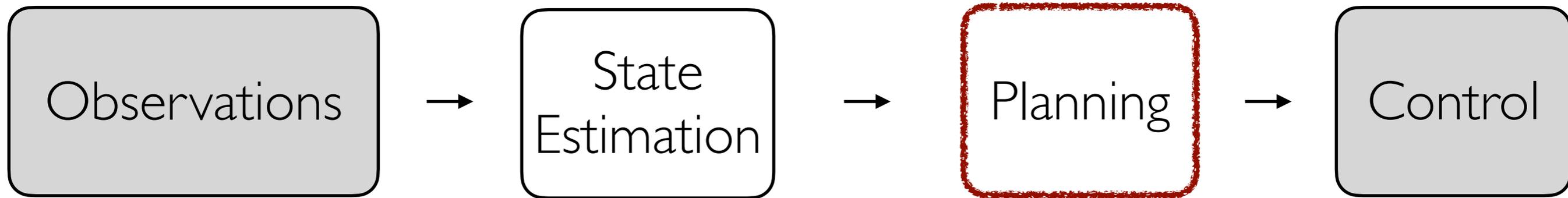
Kavraki et al. RAI 1996. Probabilistic roadmaps for path planning in high-dimensional configuration spaces.

Lavalle and Kuffner. 2000. Rapidly-exploring random trees: Progress and prospects.

Video Credits: Mur-Artal et al., Palmieri et al.

Structured Policies

Borrow policy structures from classical pipelines, and inject learning.

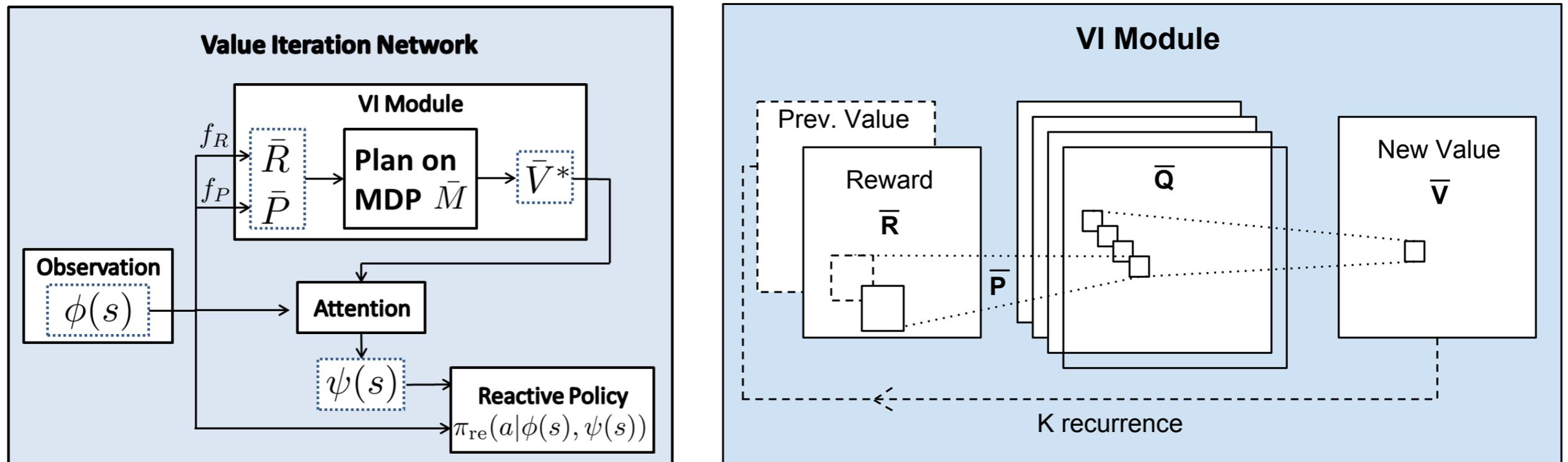


If we can differentiate through the path planning step, we can train state estimators, directly for outputting good control.

- Analytically differentiate through a planner
 - eg: Differentiating through value iteration
 - eg: Differentiating through MPC (this paper)
- Train a planner (and differentiate through it)
 - eg: Value iteration networks
 - eg: Spatial Planning using Transformers

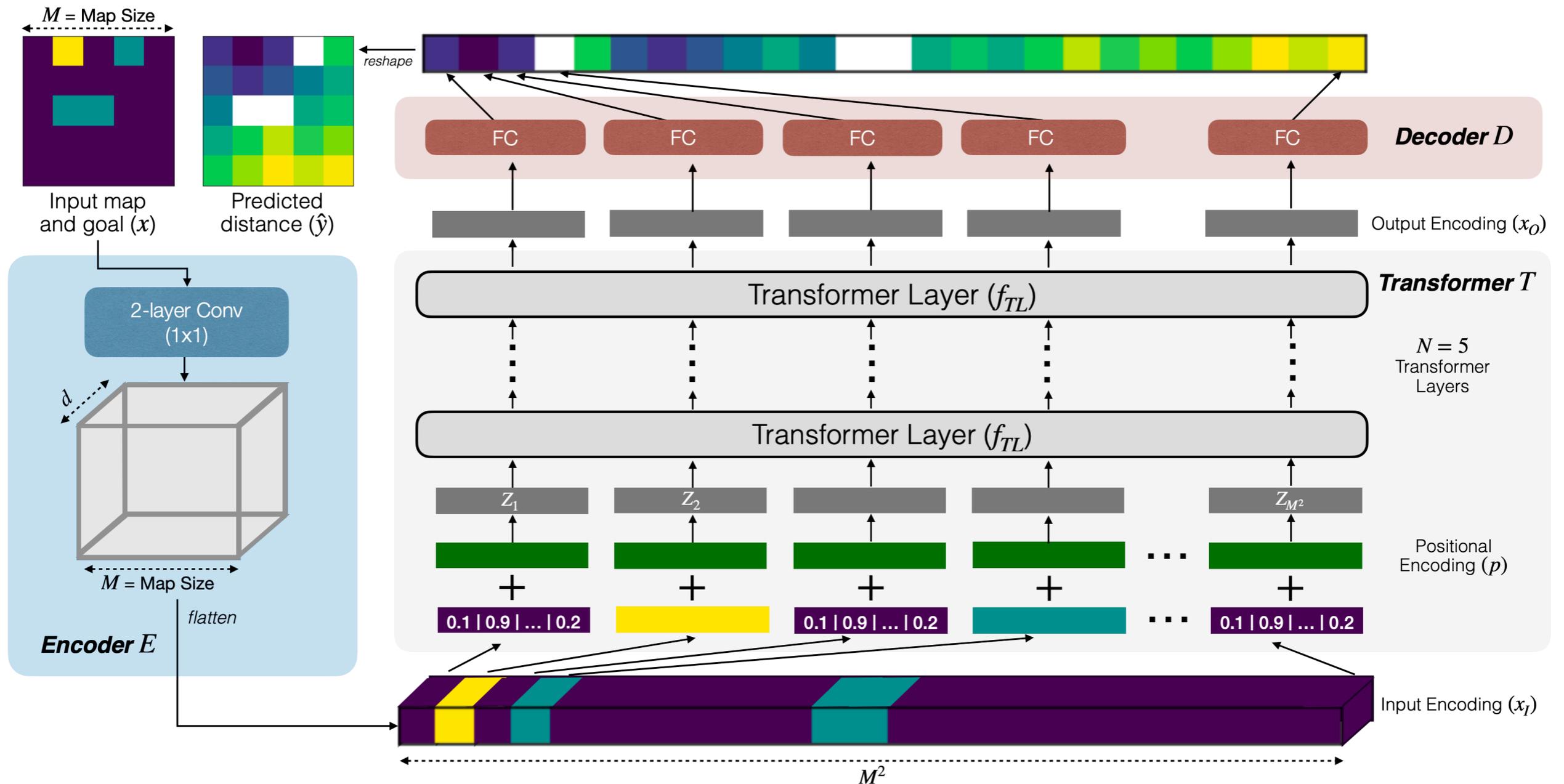
Value iteration

Value iteration networks



express value iteration computation as a convolutional network
(with max-pooling across channels)

Spatial Planning using Transformers



long-range value prediction via transformers

Differentiable MPC for End-to-end Planning and Control

Brandon Amos¹ Ivan Dario Jimenez Rodriguez² Jacob Sacks²

Byron Boots² J. Zico Kolter^{1,3}

¹Carnegie Mellon University

²Georgia Tech

³Bosch Center for AI

- $\tau_{1:T}^* = \operatorname{argmin}_{\tau_{1:T}} \sum_{t=1}^T \frac{1}{2} \tau_t^\top C_t \tau_t + c_t^\top \tau_t$ subject to $x_1 = x_{\text{init}}, x_{t+1} = F_t \tau_t + f_t$.

- $\theta = \{C, c, F, f\}, \tau_{1:T}^* = \{x_t, u_t\}_{1:T}$

- Paper is about, how to compute gradient: $\frac{\partial \ell}{\partial \theta} = \frac{\partial \ell}{\partial \tau_{1:T}^*} \frac{\partial \tau_{1:T}^*}{\partial \theta}$.

Builds upon OptNet

$$z_{i+1} = \operatorname{argmin}_z \frac{1}{2} z^T Q(z_i) z + q(z_i)^T z$$

subject to $A(z_i)z = b(z_i)$

$$G(z_i)z \leq h(z_i)$$

-
- Optimization layer:
 - inputs z_i describe a constrained optimization problem
 - output z_{i+1} is the minimizer of the optimization problem

Builds upon OptNet

$$z_{i+1} = \underset{z}{\operatorname{argmin}} \frac{1}{2} z^T Q(z_i) z + q(z_i)^T z$$

$$\text{subject to } A(z_i) z = b(z_i)$$

$$G(z_i) z \leq h(z_i)$$

-

- Compute Lagrangian, with Lagrange multipliers, $\lambda \geq 0$, ν :

$$L(z, \nu, \lambda) = \frac{1}{2} z^T Q z + q^T z + \nu^T (A z - b) + \lambda^T (G z - h)$$

- For a convex optimization problem, at optima, KKT conditions, complementary slackness must hold:

$$Q z^* + q + A^T \nu^* + G^T \lambda^* = 0$$

$$A z^* - b = 0$$

$$D(\lambda^*)(G z^* - h) = 0,$$

-

Builds upon OptNet

- Get gradients via implicit differentiation:

$$\begin{aligned}dQz^* + Qdz + dq + dA^T \nu^* + \\ A^T d\nu + dG^T \lambda^* + G^T d\lambda = 0 \\ dAz^* + Adz - db = 0\end{aligned}$$

- $D(Gz^* - h)d\lambda + D(\lambda^*)(dGz^* + Gdz - dh) = 0$

$$\begin{bmatrix} Q & G^T & A^T \\ D(\lambda^*)G & D(Gz^* - h) & 0 \\ A & 0 & 0 \end{bmatrix} \begin{bmatrix} dz \\ d\lambda \\ d\nu \end{bmatrix} =$$
$$- \begin{bmatrix} dQz^* + dq + dG^T \lambda^* + dA^T \nu^* \\ D(\lambda^*)dGz^* - D(\lambda^*)dh \\ dAz^* - db \end{bmatrix}.$$

- Rearrange:

$$\begin{bmatrix} Q & G^T & A^T \\ D(\lambda^*)G & D(Gz^* - h) & 0 \\ A & 0 & 0 \end{bmatrix} \begin{bmatrix} dz \\ d\lambda \\ d\nu \end{bmatrix} = - \begin{bmatrix} dQz^* + dq + dG^T \lambda^* + dA^T \nu^* \\ D(\lambda^*)dGz^* - D(\lambda^*)dh \\ dAz^* - db \end{bmatrix}.$$

- Rearrange:
- Can compute partial derivatives: $\partial z / \partial b$, by solving system of equations
- Better still, can directly compute $\frac{\partial l}{\partial b} = \frac{\partial l}{\partial z} \frac{\partial z}{\partial b}$

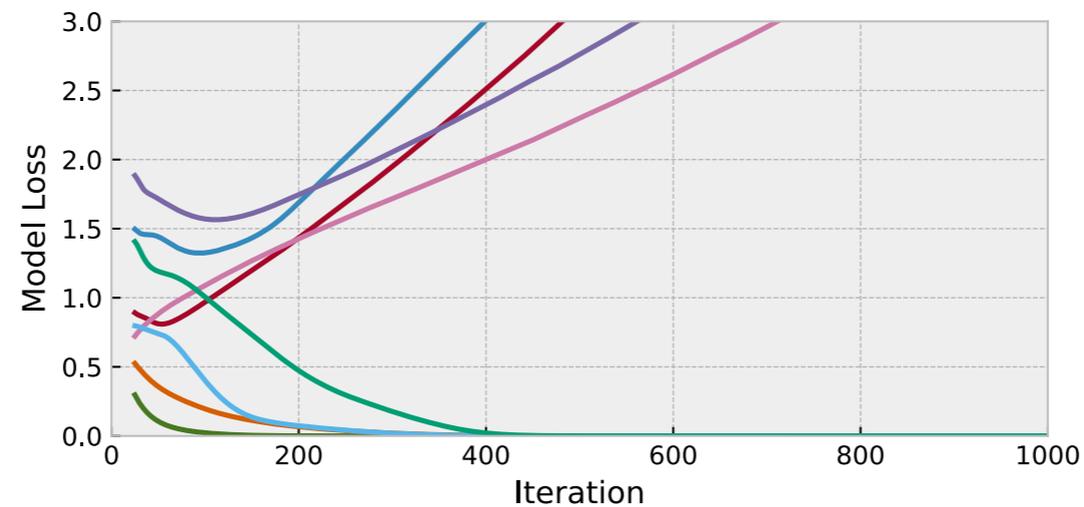
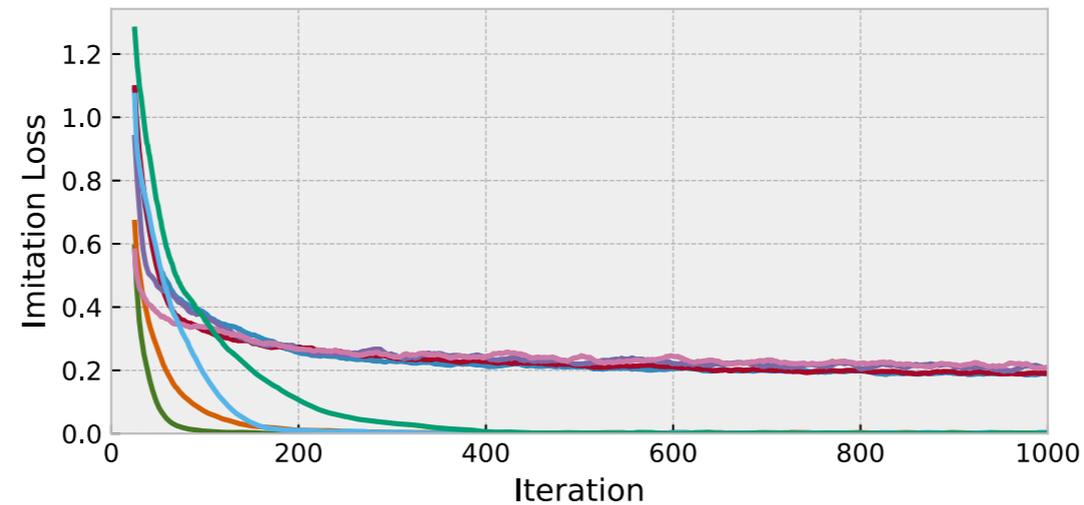
$$\begin{bmatrix} d_z \\ d_\lambda \\ d_\nu \end{bmatrix} = - \begin{bmatrix} Q & G^T D(\lambda^*) & A^T \\ G & D(Gz^* - h) & 0 \\ A & 0 & 0 \end{bmatrix}^{-1} \begin{bmatrix} \left(\frac{\partial l}{\partial z^*}\right)^T \\ 0 \\ 0 \end{bmatrix}$$

$$\nabla_Q \ell = \frac{1}{2}(d_z z^T + z d_z^T) \quad \nabla_q \ell = d_z$$

$$\nabla_A \ell = d_\nu z^T + \nu d_z^T \quad \nabla_b \ell = -d_\nu$$

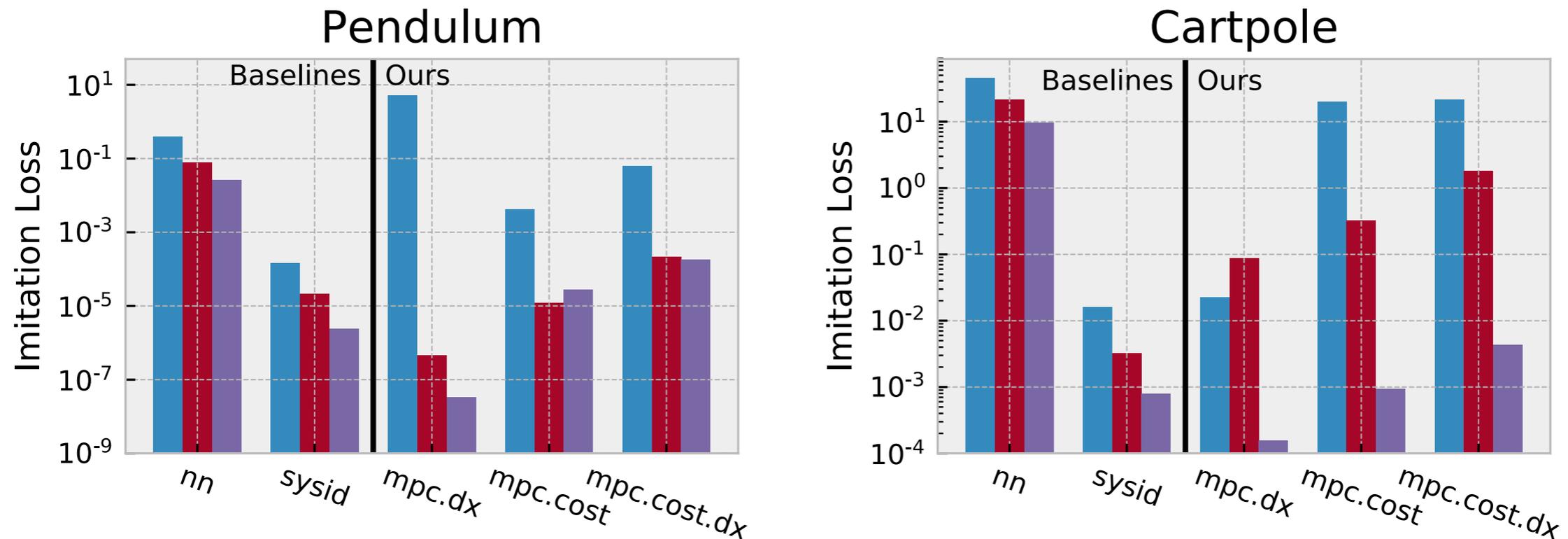
$$\nabla_G \ell = D(\lambda^*)(d_\lambda z^T + \lambda d_z^T) \quad \nabla_h \ell = -D(\lambda^*)d_\lambda$$

Experiments: Learning LQR Models



- Learning dynamics function (A and B matrices) for LQR.

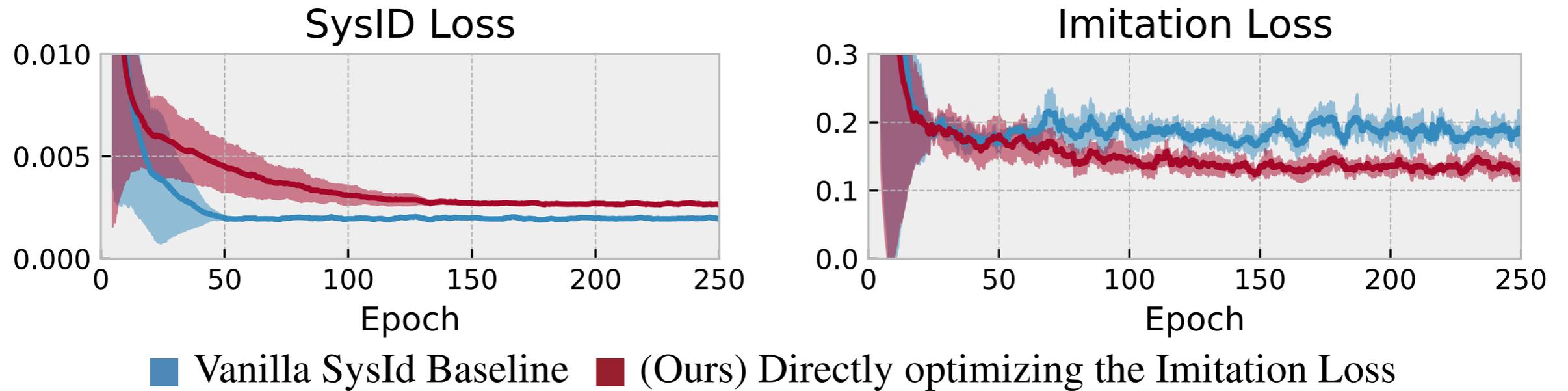
Experiments: Learning in non-convex problems



■ #Train: 10 ■ #Train: 50 ■ #Train: 100

- nn: LSTM to predict next action
- sysid: known cost, estimates parameters for 1 step dynamics functions
- mpc.dx: known cost, estimates parameters for 1 step dynamics functions
- mpc.cost: known dynamics, estimates cost parameters
- mpc.cost.dx: estimates both cost and dynamics function parameters

Experiments: SysID with a non-realizable expert



- Figure 5: Convergence results in the non-realizable Pendulum task.
- sysid: known cost, estimates parameters for 1 step dynamics functions
- mpc.dx: known cost, estimates parameters for 1 step dynamics functions
- SysID does better at estimating parameters, but worse at imitating

Thank you