# Imitation-Projected Programmatic Reinforcement Learning

**Abhinav Verma\***
*Rice University*

**Hoang M. Le\***
*Caltech*

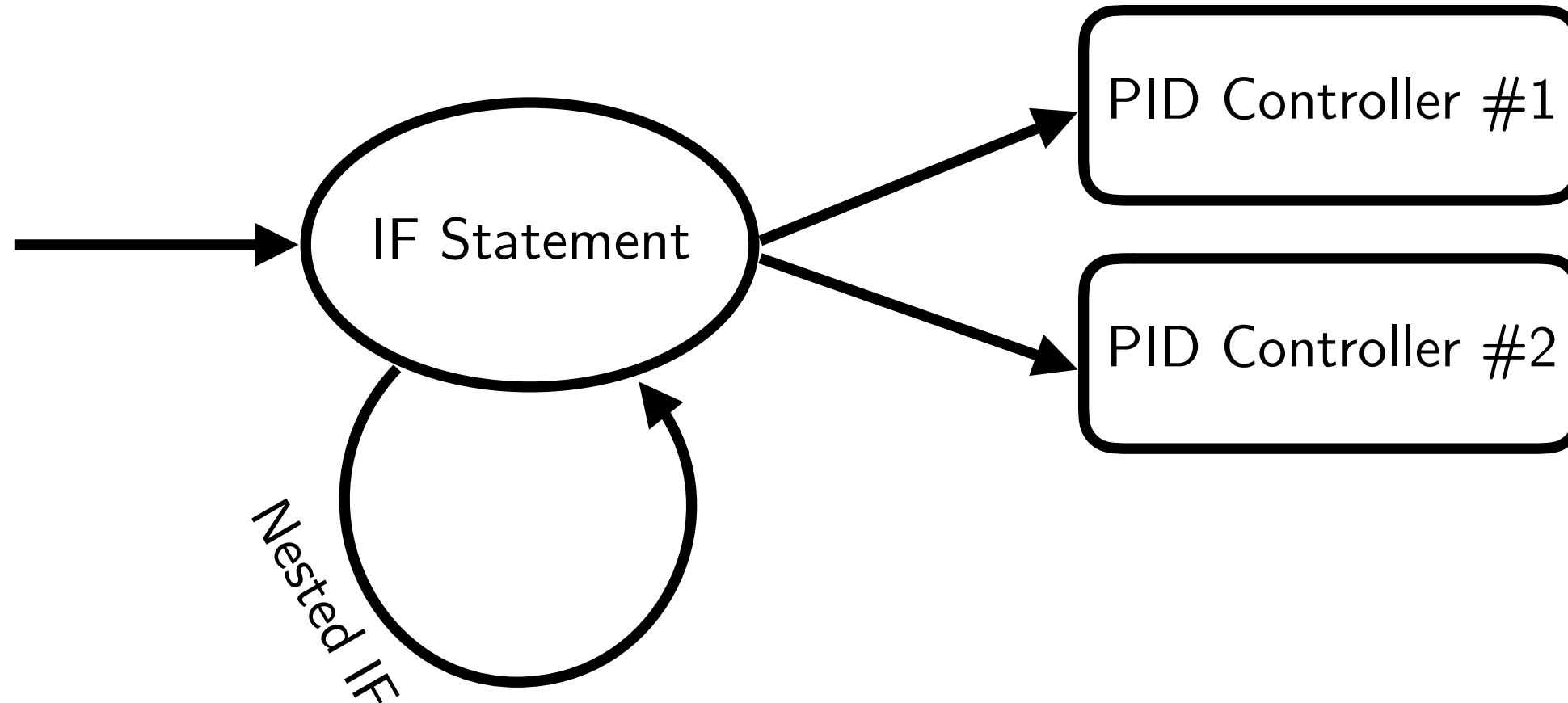**Yisong Yue**
*Caltech*

**Swarat Chaudhuri**
*Rice University*

*(* equal contribution)*

# Why learning programmatic policies

- Natural way to encode structure into policy class

- Benefits: policy-based guarantee, symbolic verification of programs (& interpretable)

- Example: domain-specific language over a family of parameterized simple controllers

# Programmatic reinforcement learning

- The program space $\Pi$

    - language (arithmetic, boolean, relational) over simple policies

- **Goal:** find the best program

$$\pi^* = \mathrm{argmin}_{\pi \in \Pi} J(\pi)$$

- Learning programmatic policies is challneging:
  - Program Synthesis: highly structured nature of policy space
  - Deep RL: programmatic policy representation incompatible with deep policy search

- **<u>Approach Overview:</u>**

    Building program structure into policy search via "lift-and-project"

# Making use of hybrid representation

- Neural policy class $F$ : deep RL representation
  - flexible, but unstable and does not satisfy desired property

- Programmatic policy class $\Pi$
  - less flexible, but certifiable

- Hybrid representation (functional regularization)

$$H \equiv \Pi \oplus F$$

$$h \equiv \pi + \lambda f \quad \text{defined as} \quad h(x) = \pi(x) + \lambda f(x)$$

# Imitation-**pro**jected **p**rogrammatic **re**inforcement **l**earning (PROPEL)
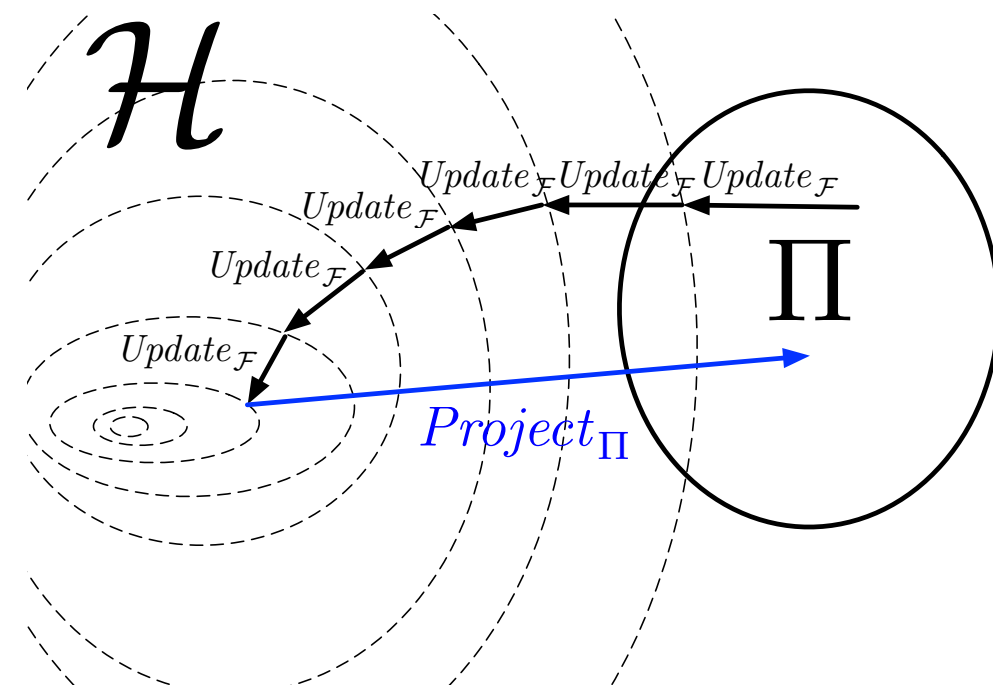
hybrid class:     $H \equiv \Pi \oplus F$

each iteration:   $h_t \leftarrow \text{UPDATE}_F(\pi_{t-1})$

$\pi_t \leftarrow \text{PROJECT}_\Pi(h_t)$

UPDATE:     $f \leftarrow f - \eta\lambda \nabla_F J(\pi + \lambda f)$

$h \leftarrow \pi + \lambda f$

PROJECT:     imitation learning

# Analysis: approximate (functional) mirror descent perspective

hybrid class: $\quad H \equiv \Pi \oplus F$

each iteration: $\quad h_t \leftarrow \text{UPDATE}_F(\pi_{t-1}) \approx \text{UPDATE}_H(\pi_{t-1})$
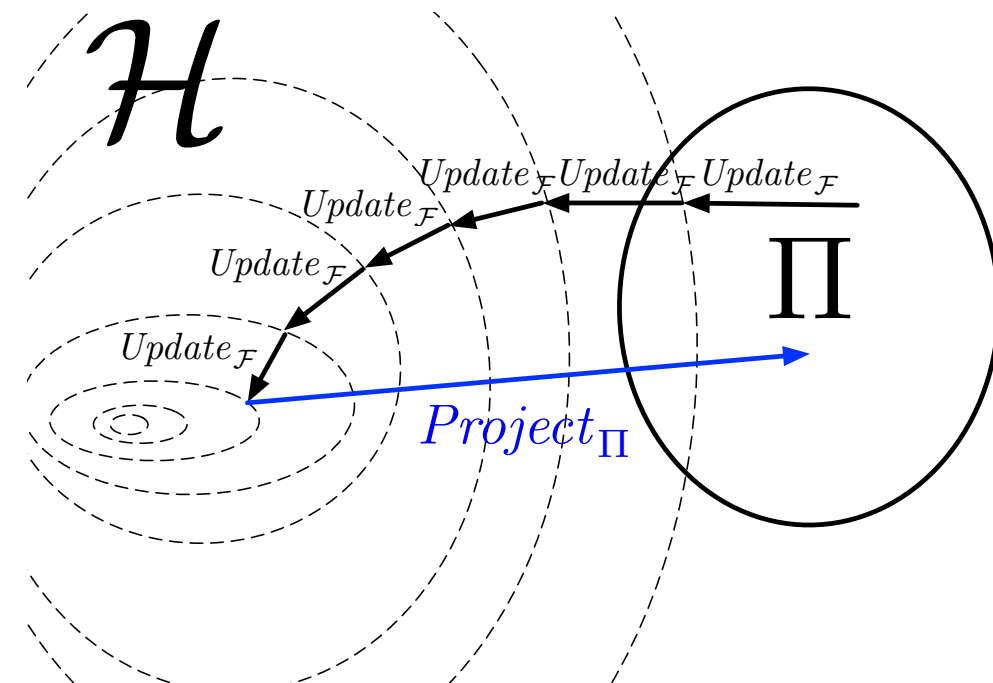
$$\pi_t \leftarrow \text{PROJECT}_\Pi(h_t) \approx \text{argmin}_{\pi \in \Pi} ||\pi - h_t||^2$$

UPDATE: $\quad f \leftarrow f - \eta \lambda \, \nabla_F J(\pi + \lambda f)$

$$h \leftarrow \pi + \lambda f$$

$$\text{UPDATE}_H(\pi_{t-1}) = \pi_{t-1} - \nabla_H J(\pi_{t-1})$$

$\text{PROJECT}_\Pi$ : finite sample error

# Experiment in TORCS

Program space: decision tree with simple, parameterized PID controllers at leaf nodes



"if the car is aligned with the axis of the track…"

$$\textbf{if } (\textbf{obs}_{\texttt{TrackPos}}(0) < 0.001 \textbf{ and } \textbf{obs}_{\texttt{TrackPos}}(0) > -0.001)$$
$$\textbf{then } PID_{\textbf{rpm}}(0.44, 4.92, 0.89, 49.79)$$
$$\textbf{else } PID_{\textbf{rpm}}(0.40, 4.92, 0.89, 49.79)$$

"then accelerate, otherwise slow down"

# Experiment in TORCS - Results

Performance: Lap time in seconds / Crash-ratio  (over 25 seeds, lower is better)

| | G-TRACK | E-ROAD | AALBORG | RUUDSKOGEN | ALPINE-2 |
|---|---|---|---|---|---|
| LENGTH | 3186M | 3260M | 2588M | 3274M | 3774M |
| PRIOR | 312.92 / 0.0 | 322.59 / 0.0 | 244.19 / 0.0 | 340.29 / 0.0 | 402.89 / 0.0 |
| DDPG | 78.82 / 0.24 | 89.71 / 0.28 | 101.06 / 0.4 | CR / 0.68 | CR / 0.92 |
| NDPS | 108.25 / 0.24 | 126.80 / 0.28 | 163.25 / 0.4 | CR / 0.68 | CR / 0.92 |
| VIPER | 83.60 / 0.24 | 87.53 / 0.28 | 110.57 / 0.4 | CR / 0.68 | CR / 0.92 |
| PROPELPROG | 93.67 / 0.4 | 119.17 / 0.4 | 147.28 / 0.12 | 124.58 / 0.16 | 256.59 / 0.16 |
| PROPELTREE | 78.33 / 0.4 | 79.39 / 0.4 | 109.83 / 0.16 | 118.80 / 0.24 | 236.01 / 0.36 |

Generalization: PROPEL completed 10/20 unseen tracks, DDPG completed 2/20

| | G-TRACK | E-ROAD | AALBORG | RUUDSKOGEN | ALPINE-2 |
|---|---|---|---|---|---|
| G-TRACK | - | 124 / CR | CR / CR | CR / CR | CR / CR |
| E-ROAD | 102 / 92 | - | CR / CR | CR / CR | CR / CR |
| AALBORG | 201 / 91 | 228 / CR | - | 217 / CR | CR / CR |
| RUUDSKOGEN | 131 / CR | 135 / CR | CR / CR | - | CR / CR |
| ALPINE-2 | 222 / CR | 231 / CR | 184 / CR | CR / CR | - |

*Code available at:*  *https://bitbucket.org/averma8053/propel*