

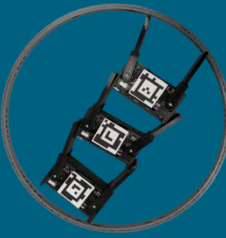
Robot learning on the edge

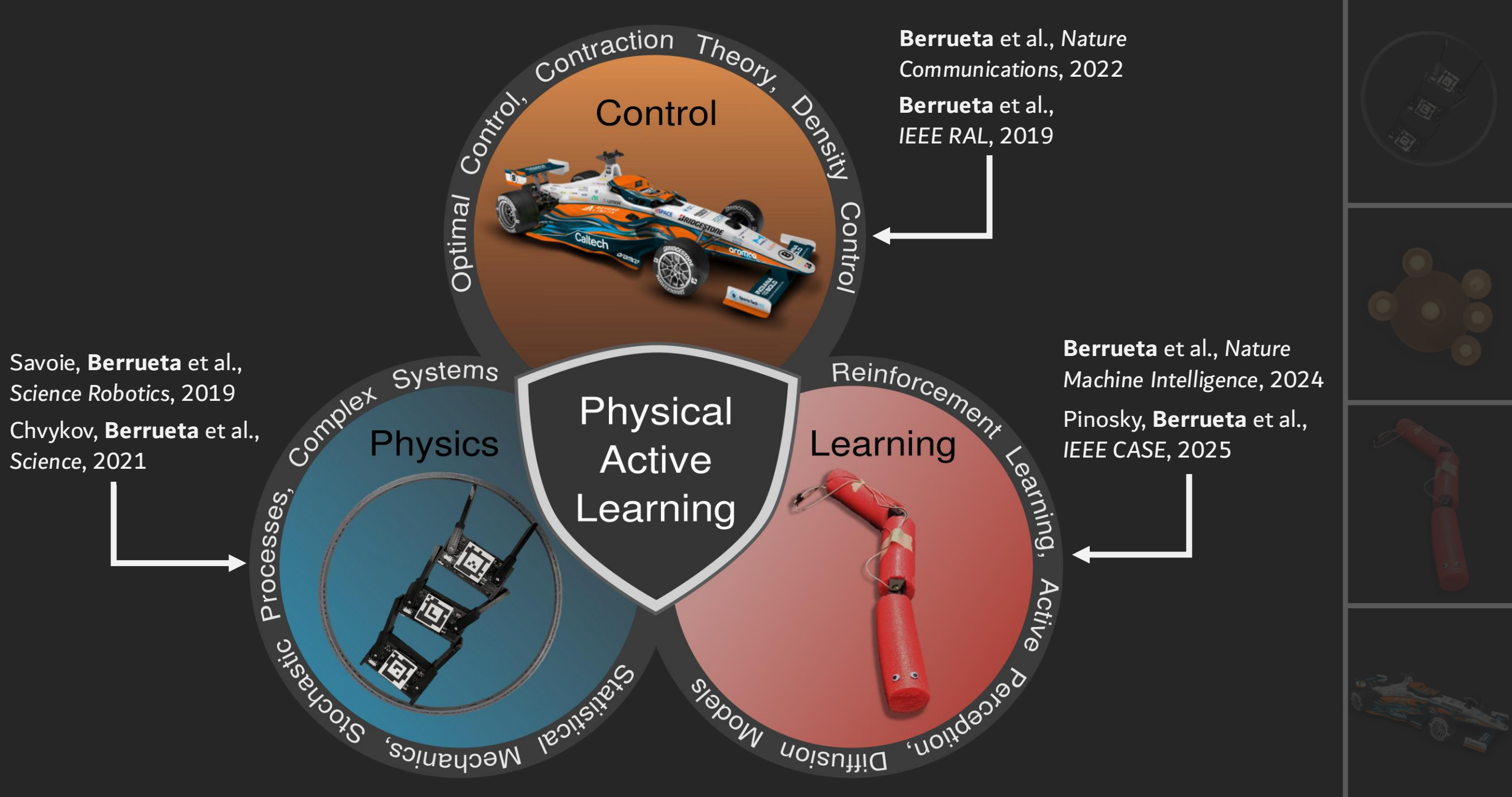
Online learning in hardware

Thomas A. Berrueta

Postdoctoral Scholar, Computing + Mathematical Sciences
ELLIIT Robot Learning Symposium (2025/11/19)

Caltech







CAST

LAP 1 89.556

LAP 2 100.547

LAP 3 110.962

LAP 4 111.436

LAP 5 121.915

LAP 6 122.360

LAP 7 133.129

LAP 8 133.514

LAP 9 144.130

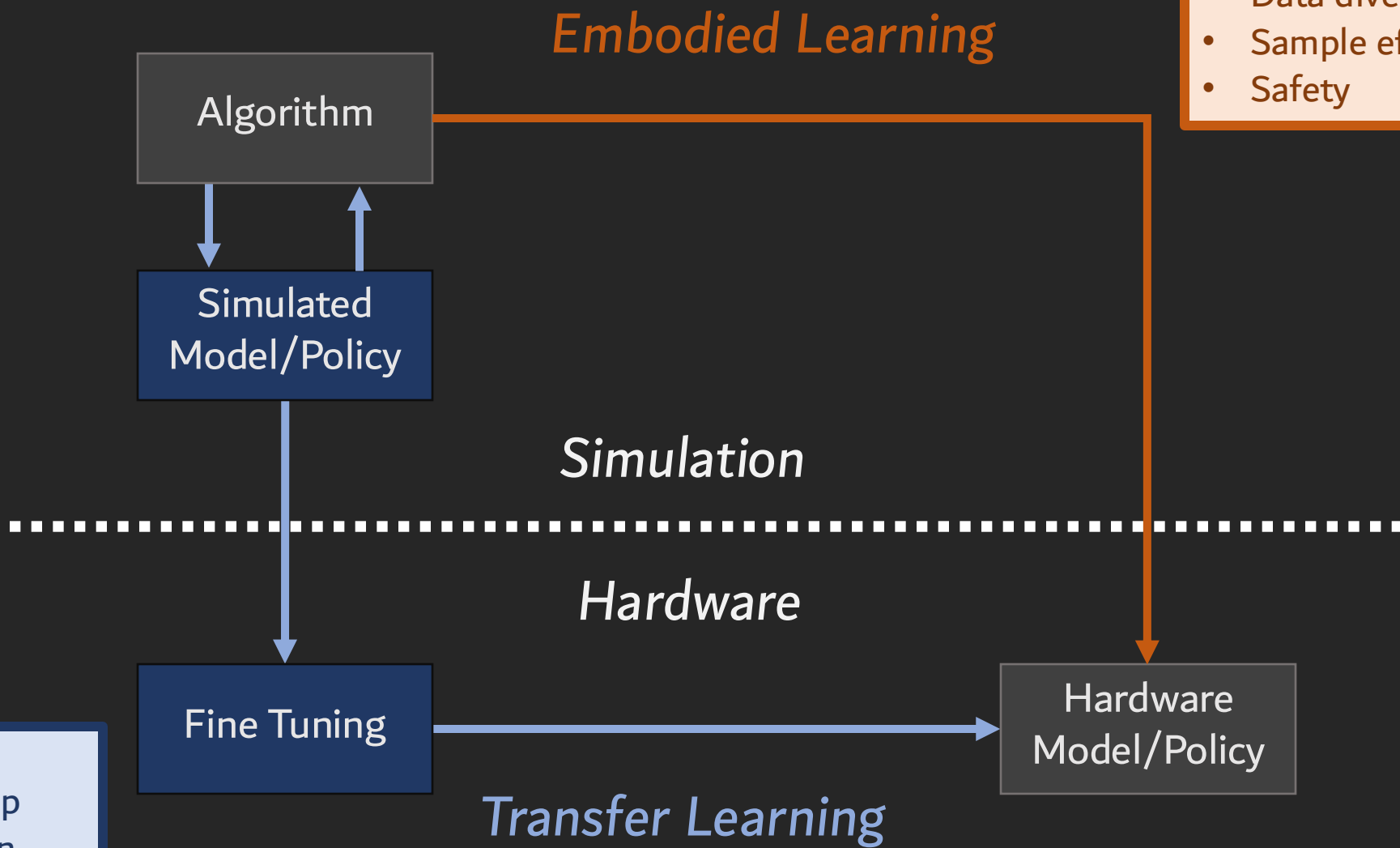


CAST Racer

THROTTLE

BRAKE

SPEED 146.4 mph GEAR 5 RPM 5704

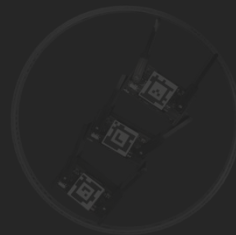


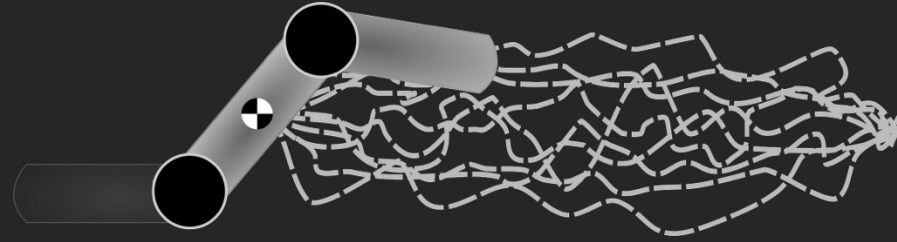
Challenges:

- Data diversity
- Sample efficiency
- Safety

Challenges:

- Reality gap
- Perception
- Hardware

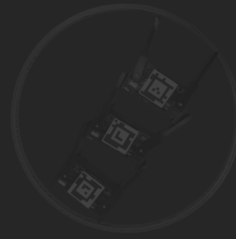


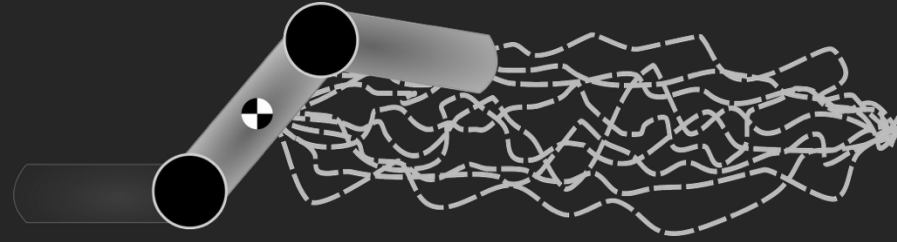


Path decorrelation for efficient online learning



Operational safety through modular design

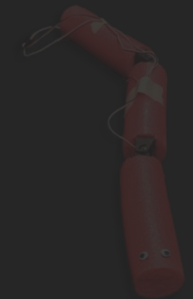
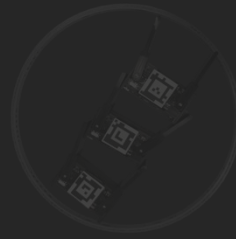




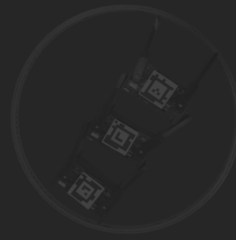
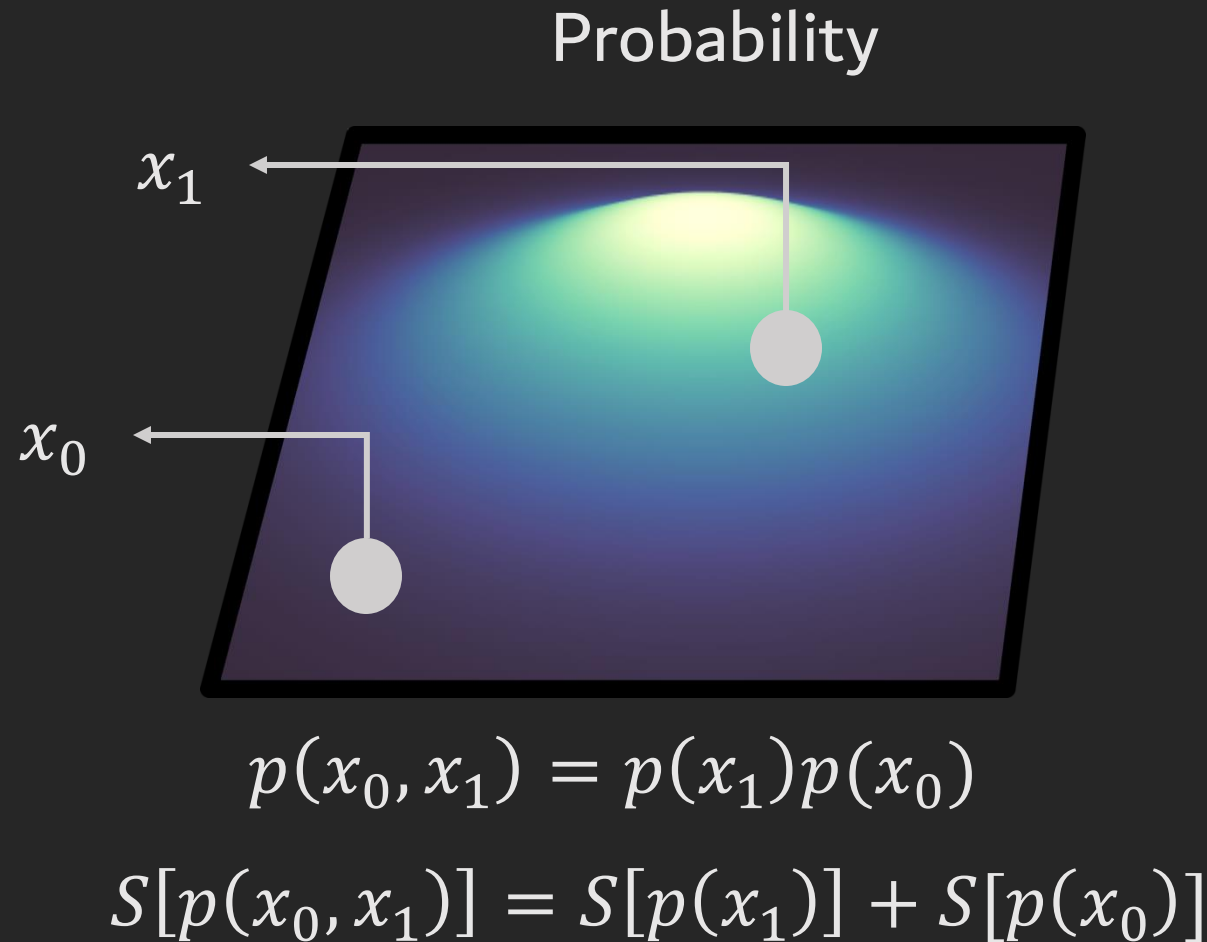
Path decorrelation for efficient online learning



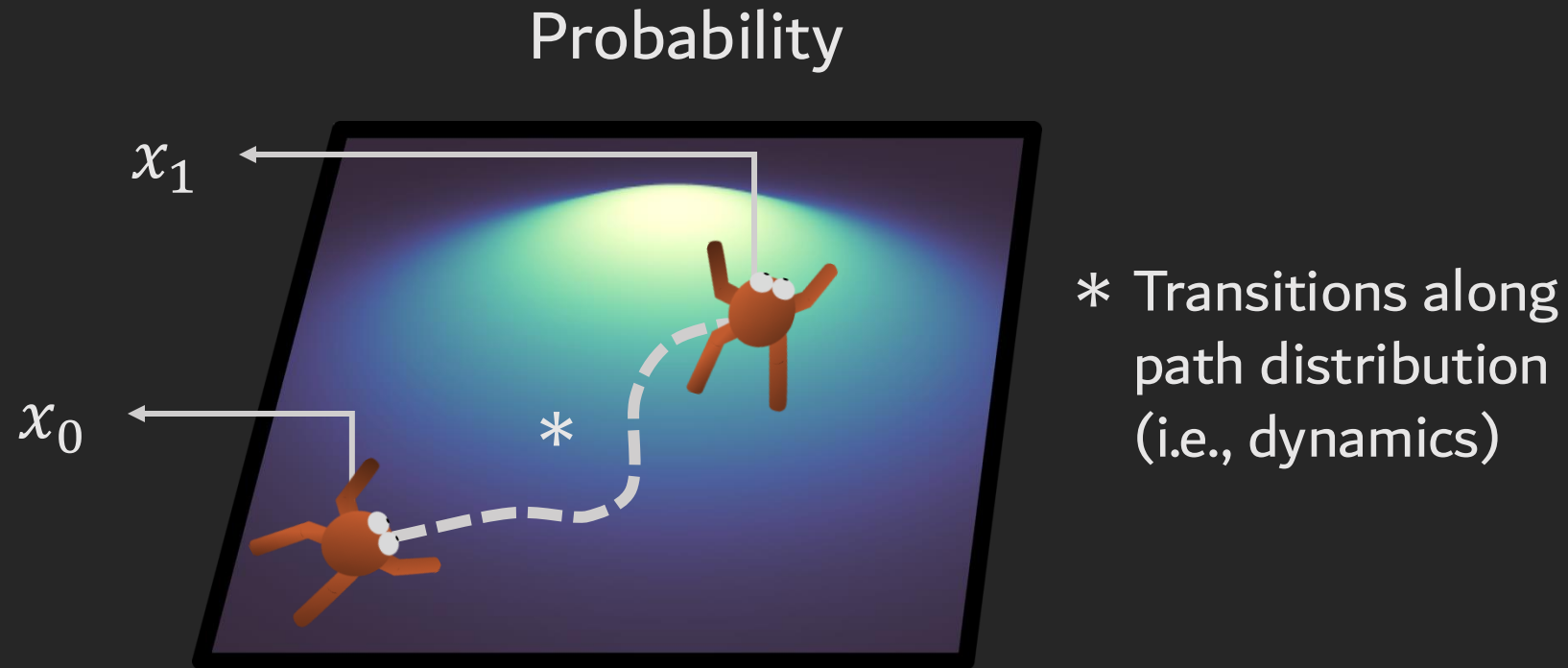
Operational safety through modular design



Physics gets in the way of data diversity



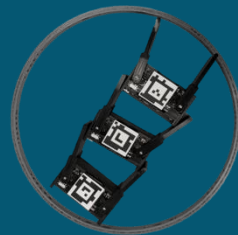
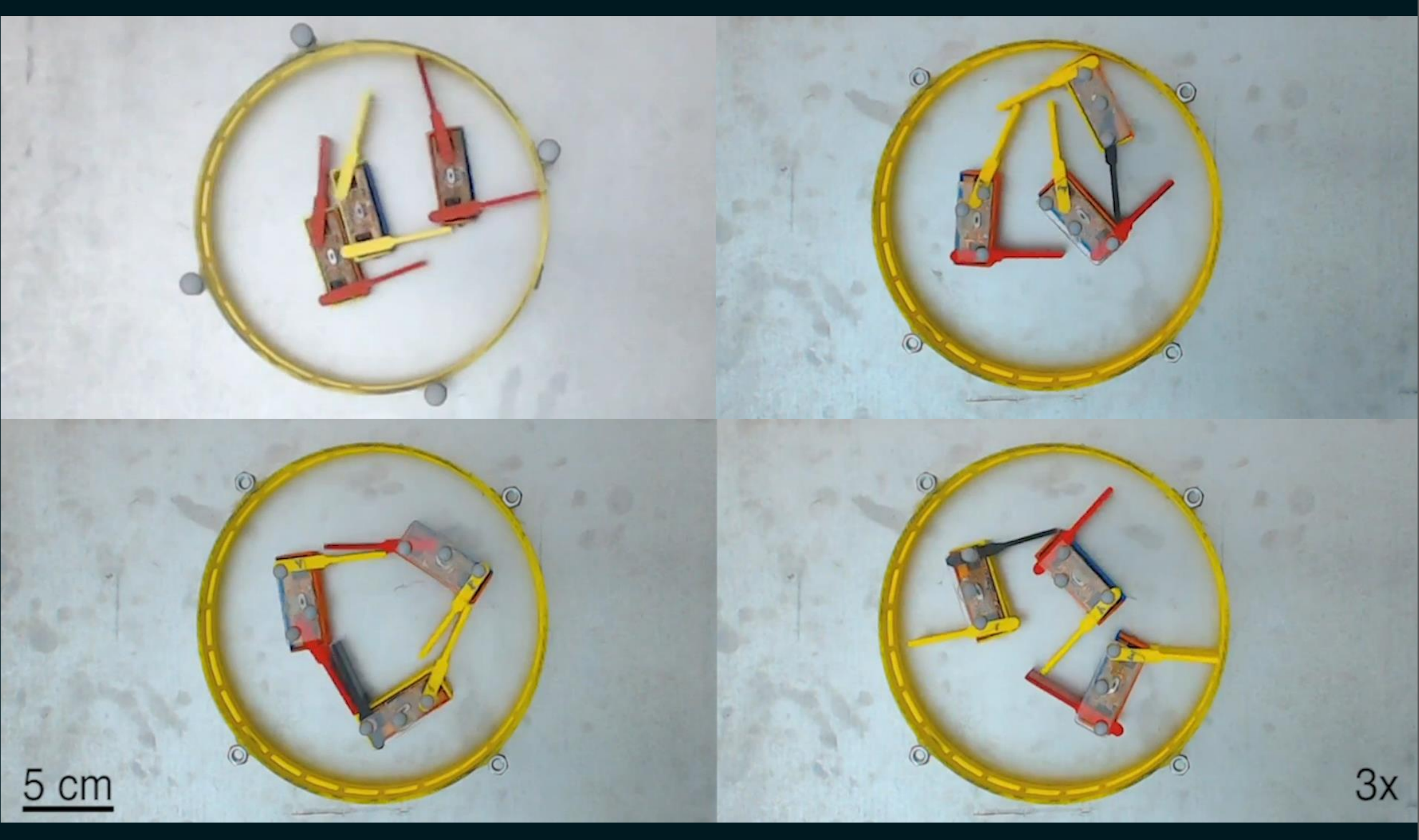
Physics gets in the way of data diversity

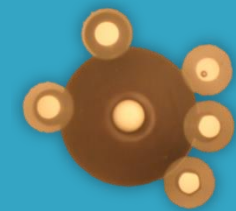
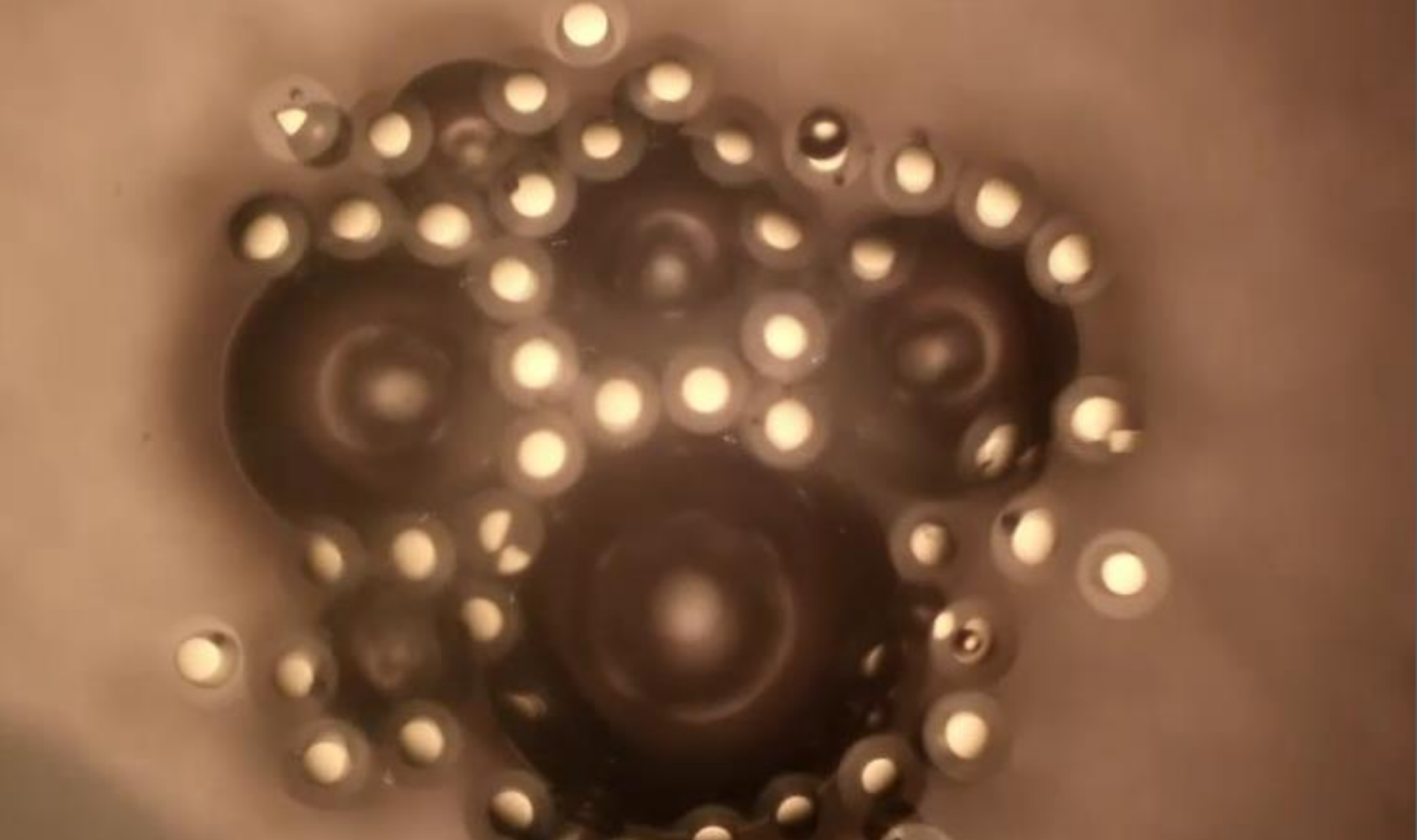


$$p(x_0, x_1) = p(x_1|x_0)p(x_0)$$

$$S[p(x_0, x_1)] \leq S[p(x_1)] + S[p(x_0)]$$

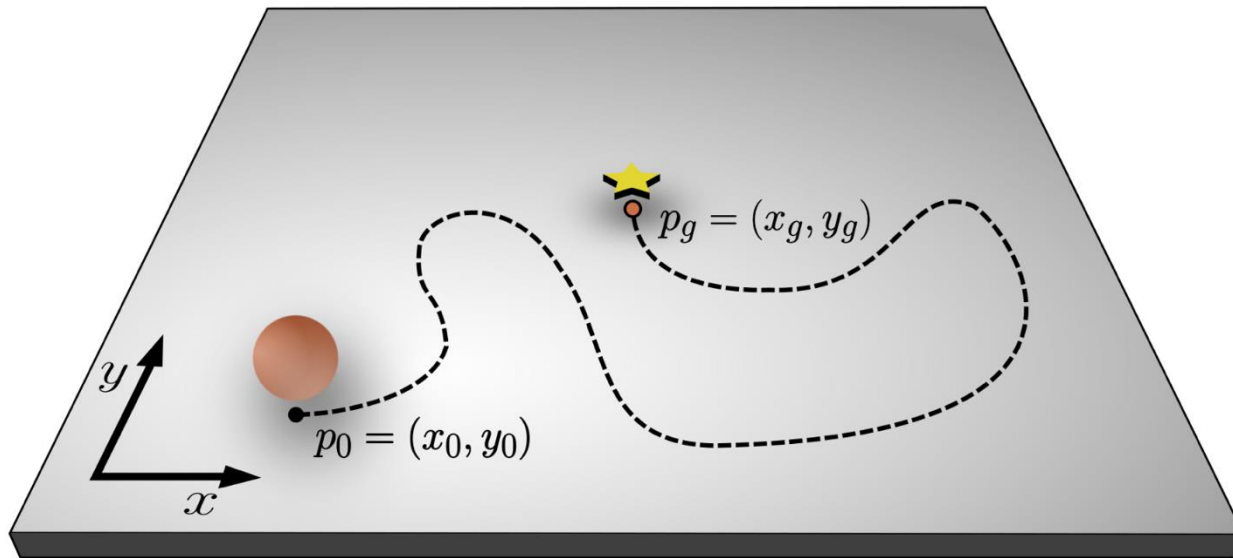






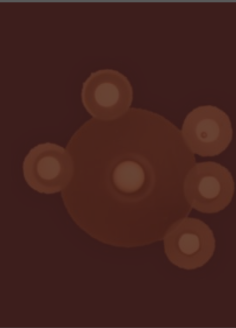
Decorrelating sample paths

Point Mass

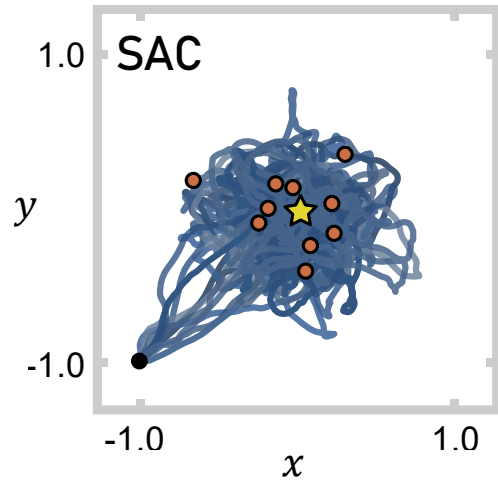
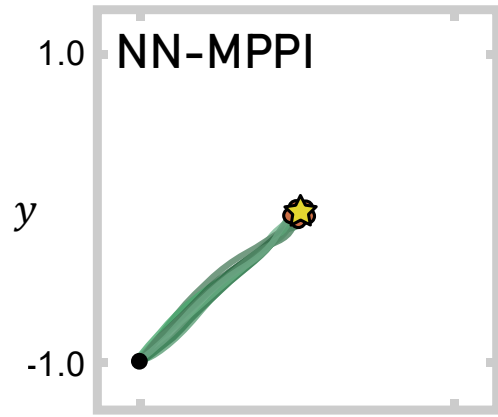


Dynamics: $\vec{x}_{t+1} = A\vec{x}_t + B\vec{u}_t$

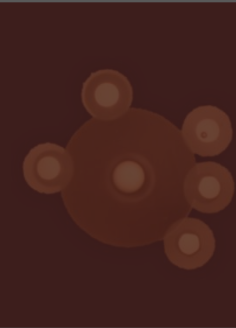
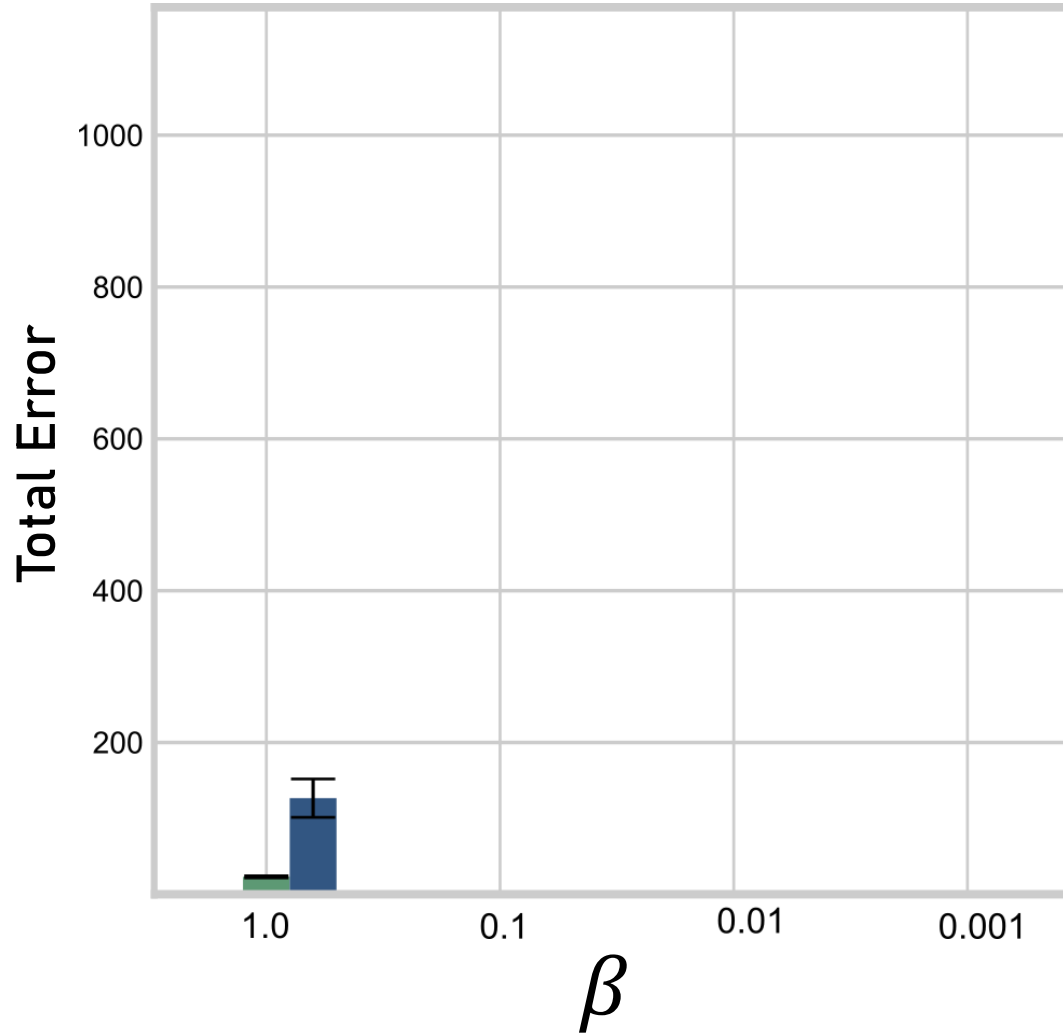
$$A = \begin{bmatrix} 1 & 0 & \beta & 0 \\ 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \quad B = \begin{bmatrix} 0 & 0 \\ 0 & 0 \\ 1 & 0 \\ 0 & 1 \end{bmatrix}$$



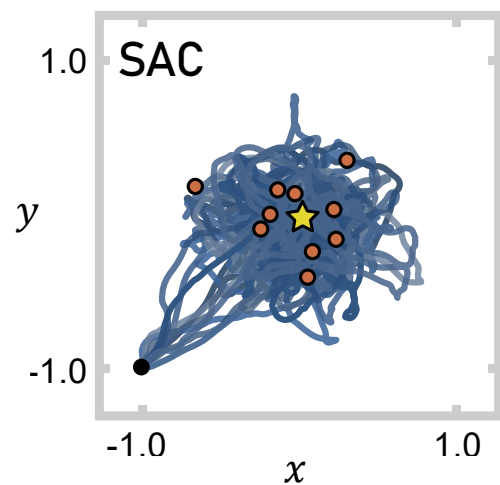
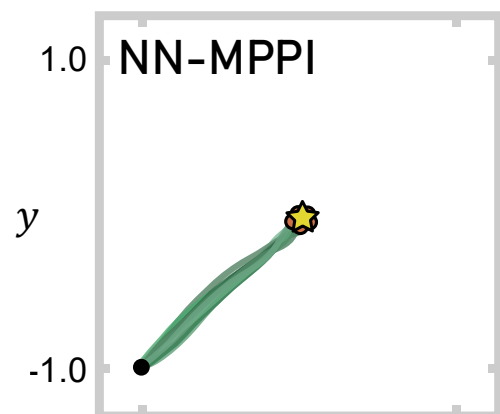
Decorrelating sample paths



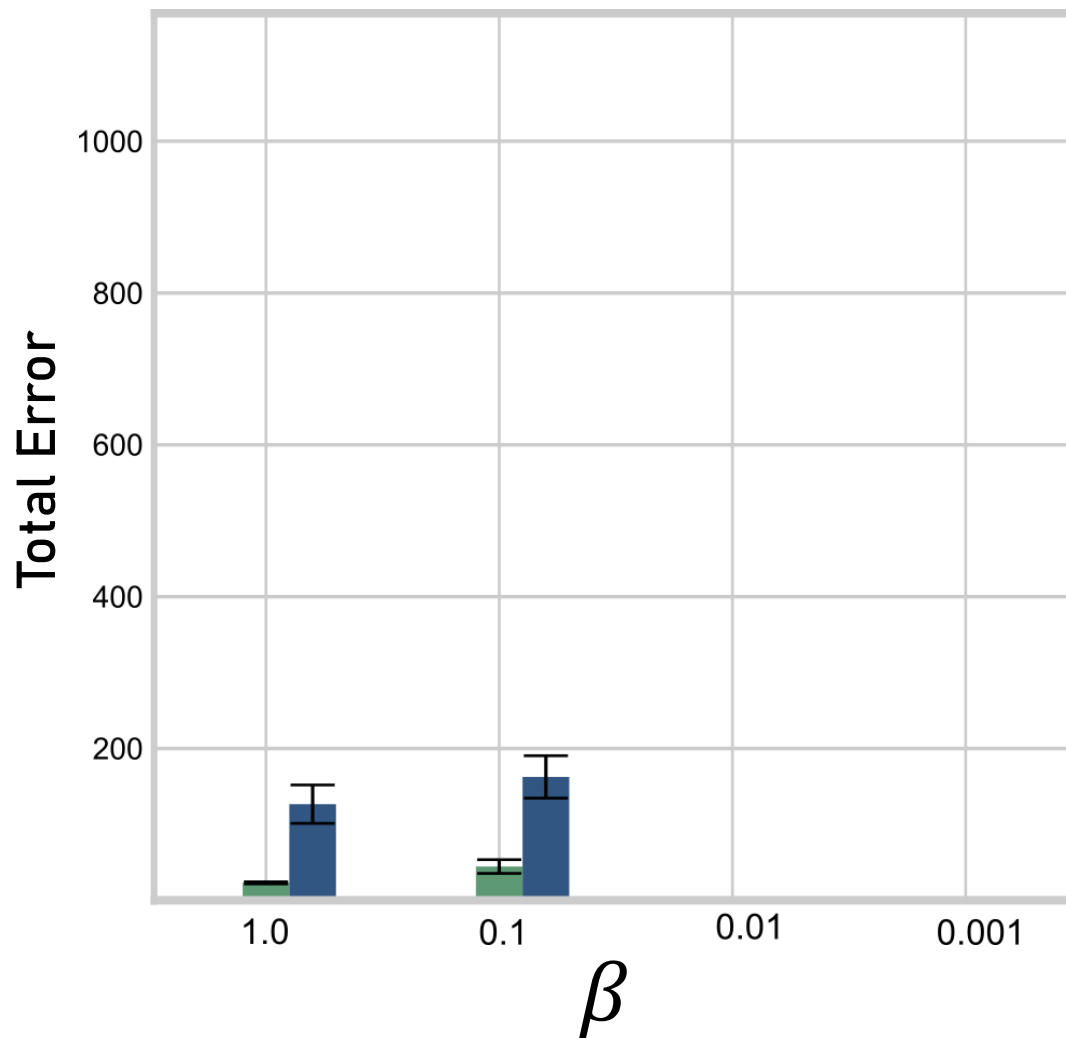
$$\beta = 1.0$$



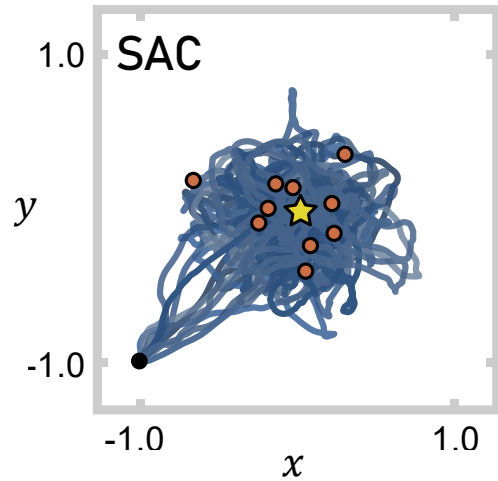
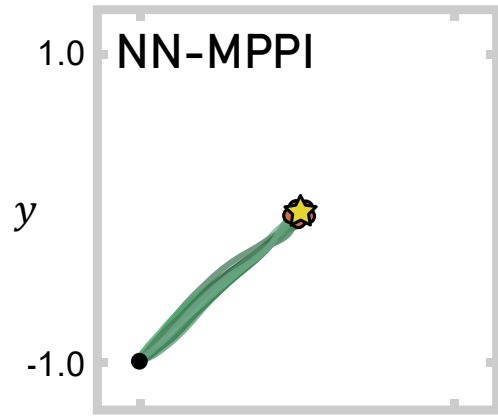
Decorrelating sample paths



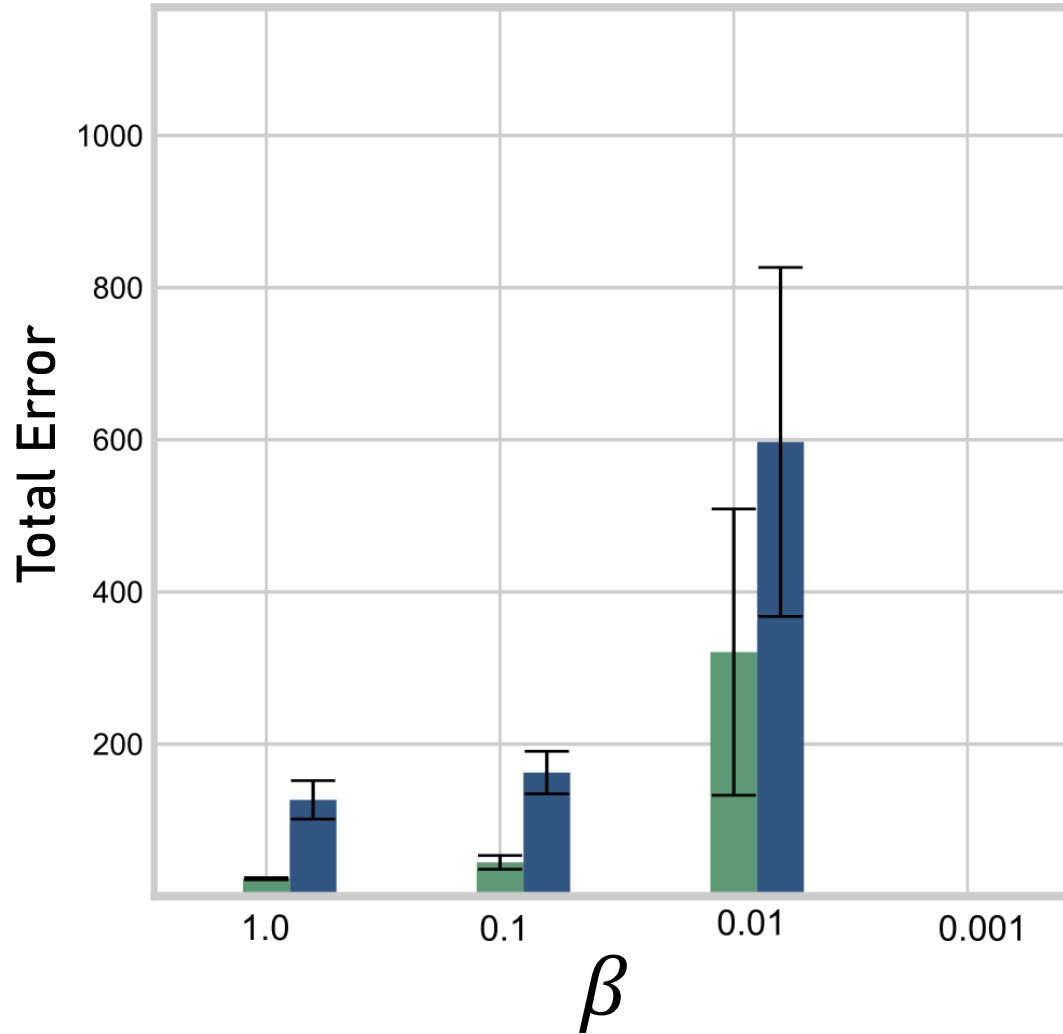
$$\beta = 1.0$$



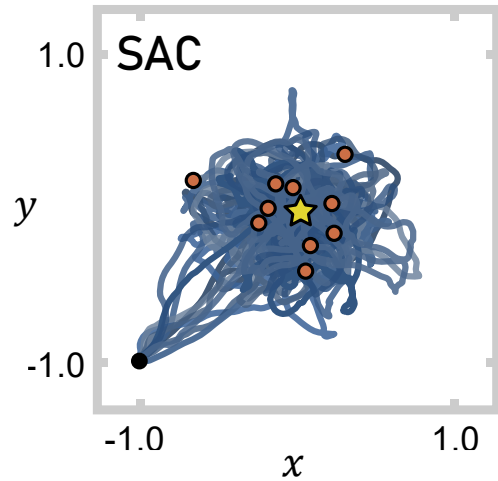
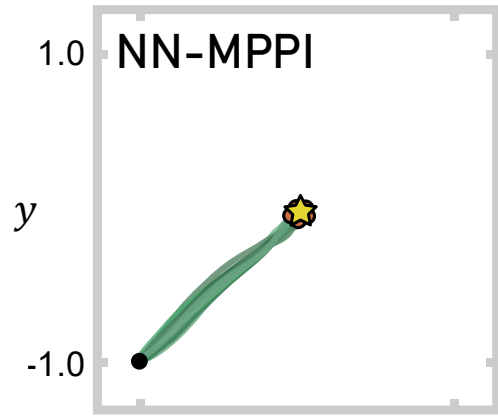
Decorrelating sample paths



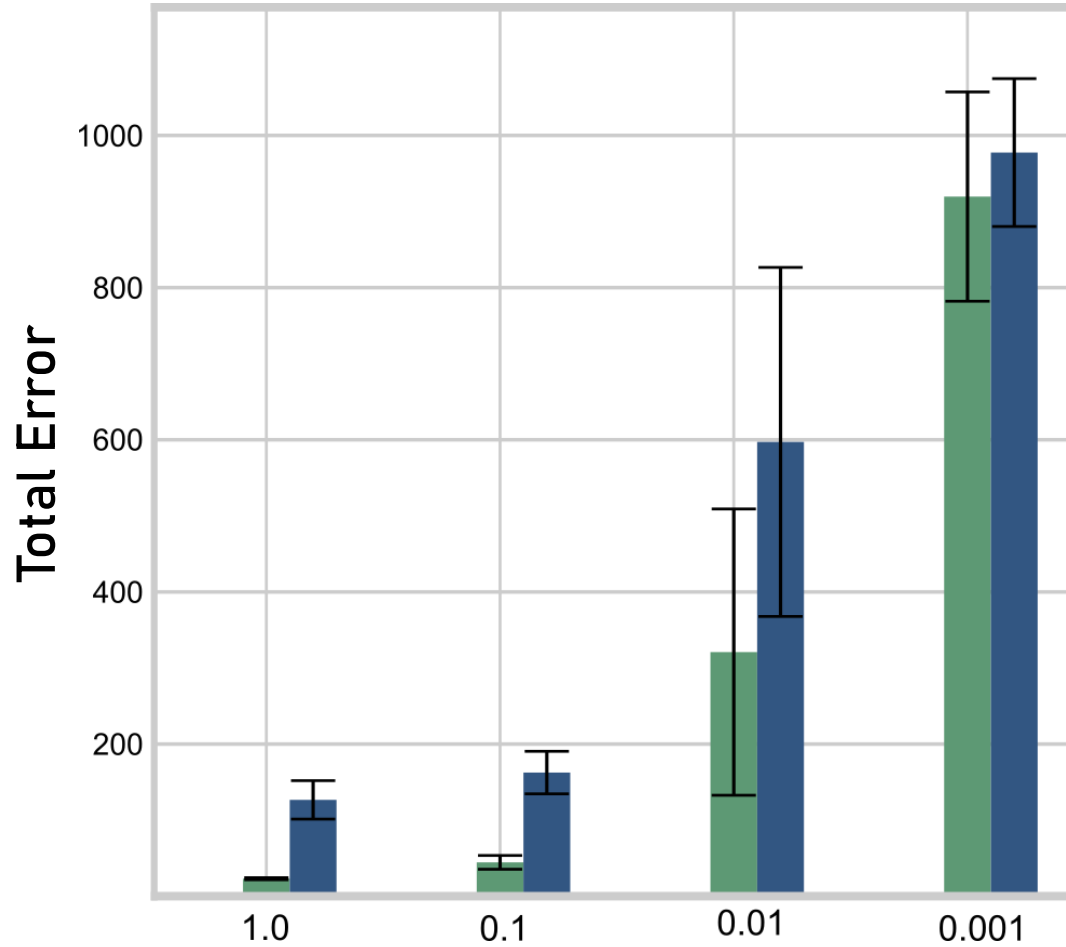
$$\beta = 1.0$$



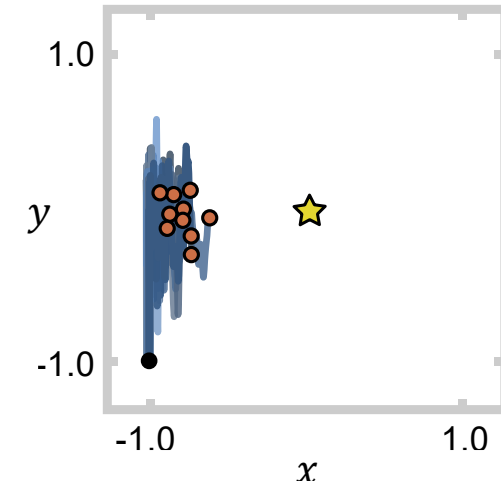
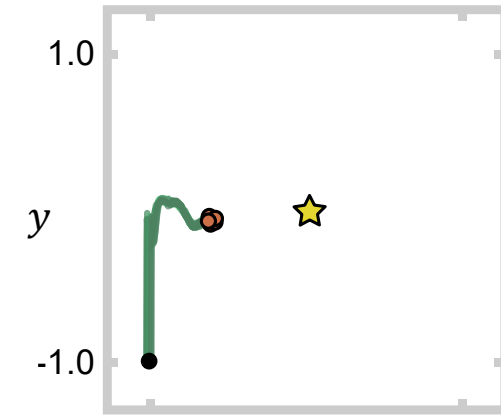
Decorrelating sample paths



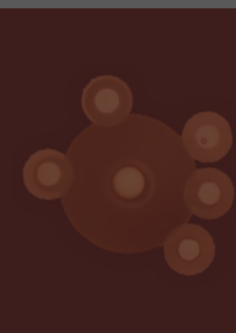
$\beta = 1.0$



β



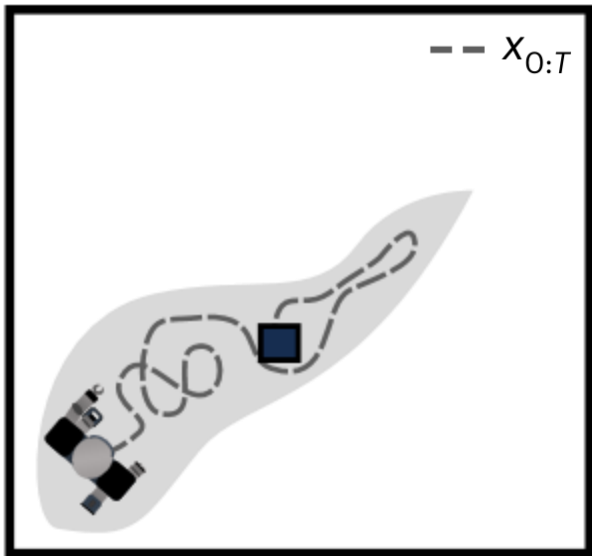
$\beta = 0.001$



Decorrelating sample paths

- Entropy maximization as a means of sample path decorrelation.
- However, this is intractable for general unknown dynamics.

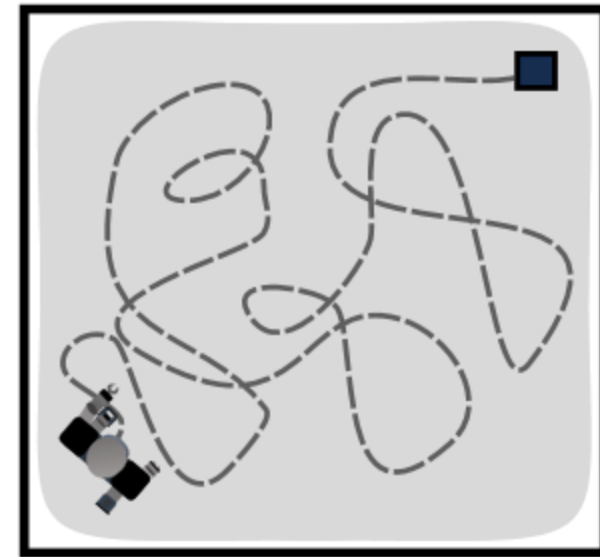
$$P[x_{0:T}] = p(x_0) \prod_{t=0}^T p(x_{t+1}|x_t)$$



$$\operatorname{argmax}_{P[x_{0:T}]} S[P[x_{0:T}]]$$



$$P_{\max}[x_{0:T}]$$



Decorrelating sample paths

$$\begin{aligned} \operatorname{argmax}_{P[x_{0:T}]} & - \int P[x_{0:T}] \log P[x_{0:T}] \mathcal{D}x_{0:T} - \lambda_0 \left(\int P[x_{0:T}] \mathcal{D}x_{0:T} - 1 \right) \\ & - \int \operatorname{Tr} \left(\Lambda(x^*)^\top (\langle \Delta x_{0:T}^2 \rangle_{x^*} - \mathbf{C}[x^*]) \right) dx^* \end{aligned}$$

where:

$$\begin{aligned} \langle \Delta x_{0:T}^2 \rangle_{x^*} &= \int P[x_{0:T}] \left[\sum_{i=0}^T (x_{i+1} - x_i)^\top (x_{i+1} - x_i) \delta(x_i - x^*) \right] \mathcal{D}x_{0:T} \\ \mathbf{C}[x^*] &= \sum_{\tau=t_i}^{t_i+\Delta t} K_{XX}(t_i, \tau), \quad (w/x_{t_i} = x^*) \end{aligned}$$



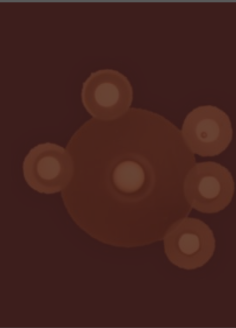
Decorrelating sample paths

- For systems with continuous sample paths, we can prove that:

$$\max_P S[P[x_{0:T}]] \leq \sum_{t=0}^T \boxed{\frac{1}{2} \log \det \mathbf{C}[x_t]} \propto \text{log-Volume of locally reachable states}$$

which is a concave and easily computable quantity.

- Optimizing this expression leads to diffusive exploration, because its optimum describes the sample paths of a class of diffusion processes.



Maximum diffusion reinforcement learning

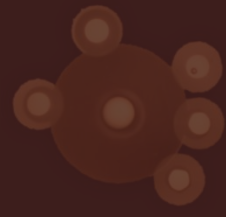
- To illustrate these results, we developed an RL pipeline based on our derivations:

$$\operatorname{argmax}_{\pi} E_{p,\pi} \left[\sum_{t=0}^T \hat{r}(x_t, u_t) \right]$$

- The rewards are augmented with a term that decorrelates sample paths:

$$\hat{r}(x_t, u_t) = r(x_t, u_t) + \frac{\alpha}{2} \log \det \mathbf{C}[x_t]$$

- Agents that optimize MaxDiff objectives are ergodic and asymptotically inherit robustness and online learning guarantees.



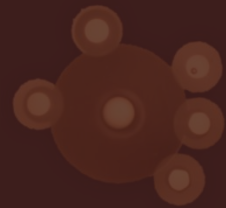
Maximum diffusion reinforcement learning

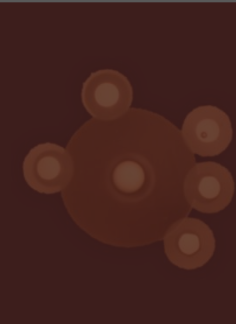
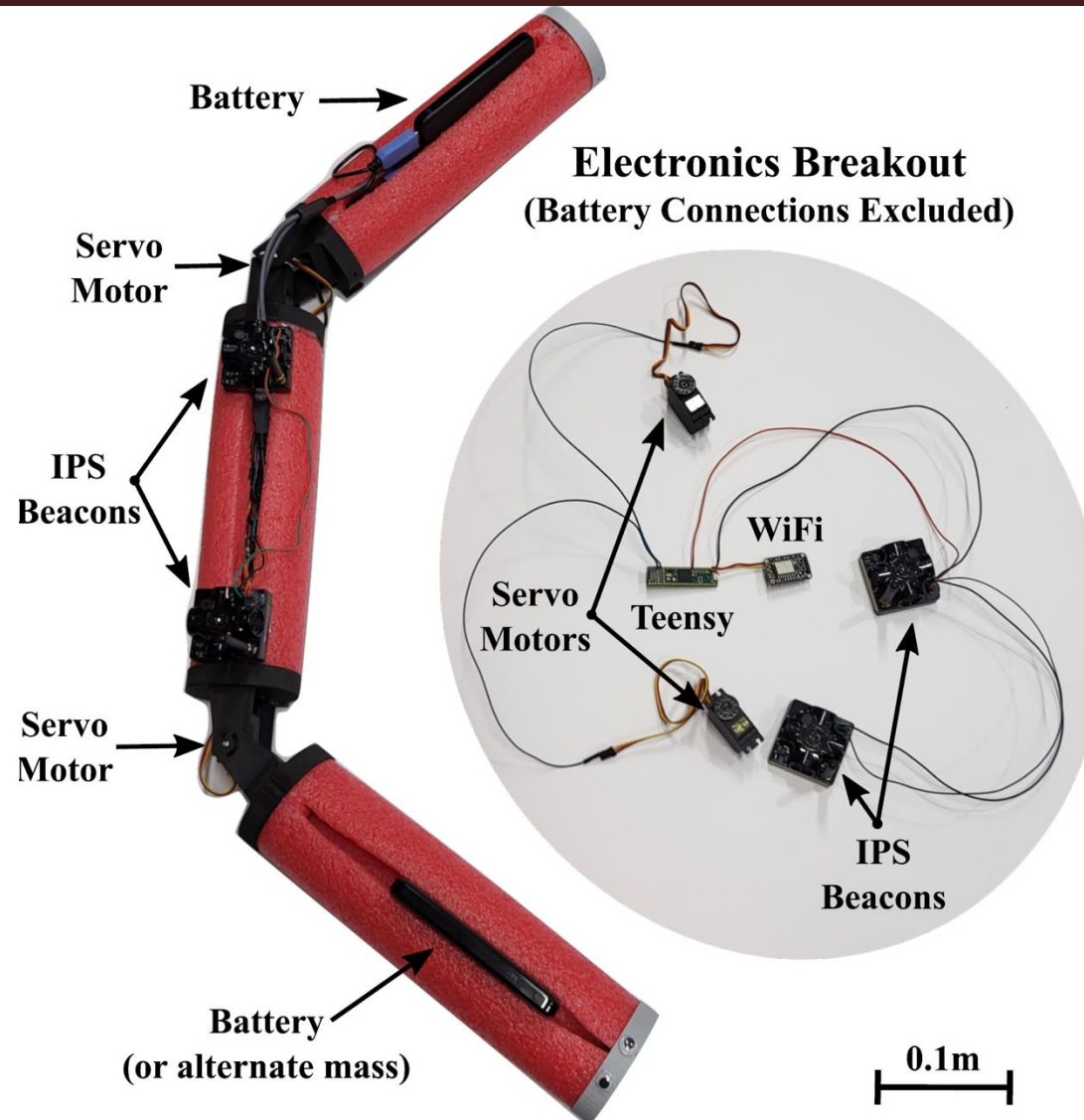
Theorem 2. (*MaxDiff RL Agents are Reliable*) If there exists a PAC-MDP algorithm with policy π^{\max} , then the Markov chain induced by π^{\max} is ergodic and π^{\max} will be ϵ -optimal regardless of initialization.

Theorem 3. (*MaxDiff RL Agents can Learn in Single-Shot*) If there exists a PAC-MDP algorithm with policy π^{\max} , then the Markov chain induced by π^{\max} is ergodic and any realization of π^{\max} will asymptotically achieve the same ϵ -optimality as an ensemble.

- To be PAC-MDP is ϵ -optimal $(1 - \delta)$ -percent of the time:

$$\Pr(\mathcal{V}_{\pi^*}(x_0) - \mathcal{V}_{\pi^{\max}}(x_0) \leq \epsilon) \geq 1 - \delta$$





Seed 1



Seed 2



Seed 3



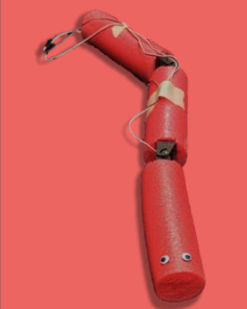
Seed 4



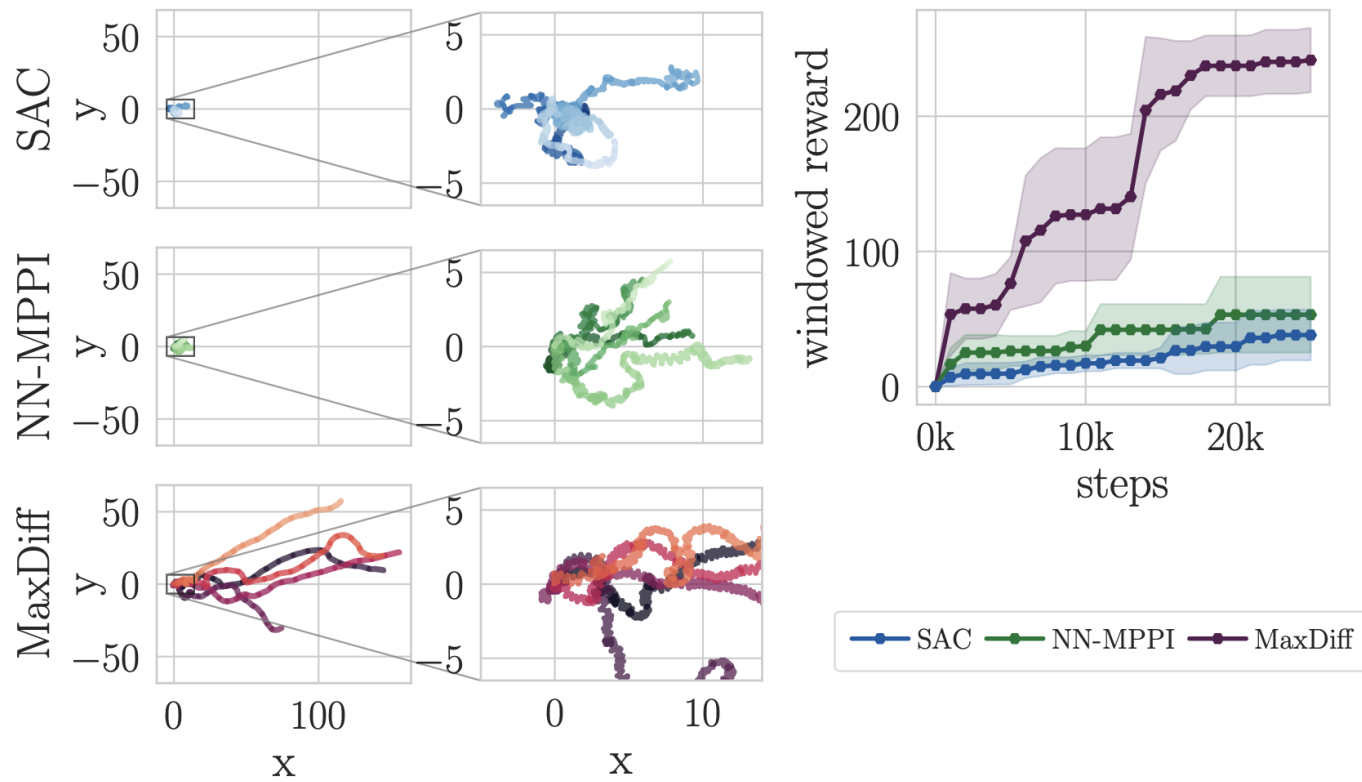
Seed 5



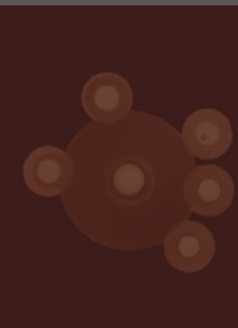
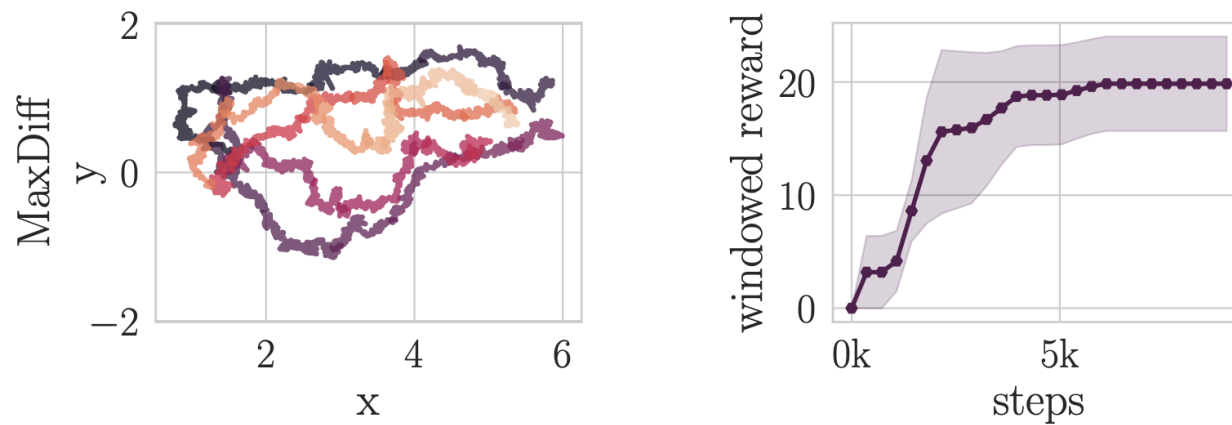
All videos are at 10x speed



Single-shot Simulation



Single-shot Hardware

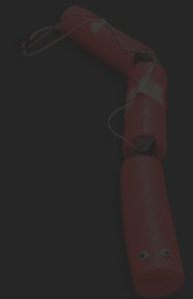
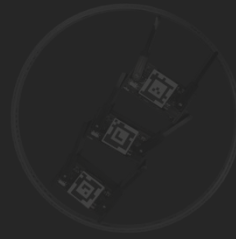




Path decorrelation for efficient online learning



Operational safety through modular design



Best Lap: 1 2:37.130

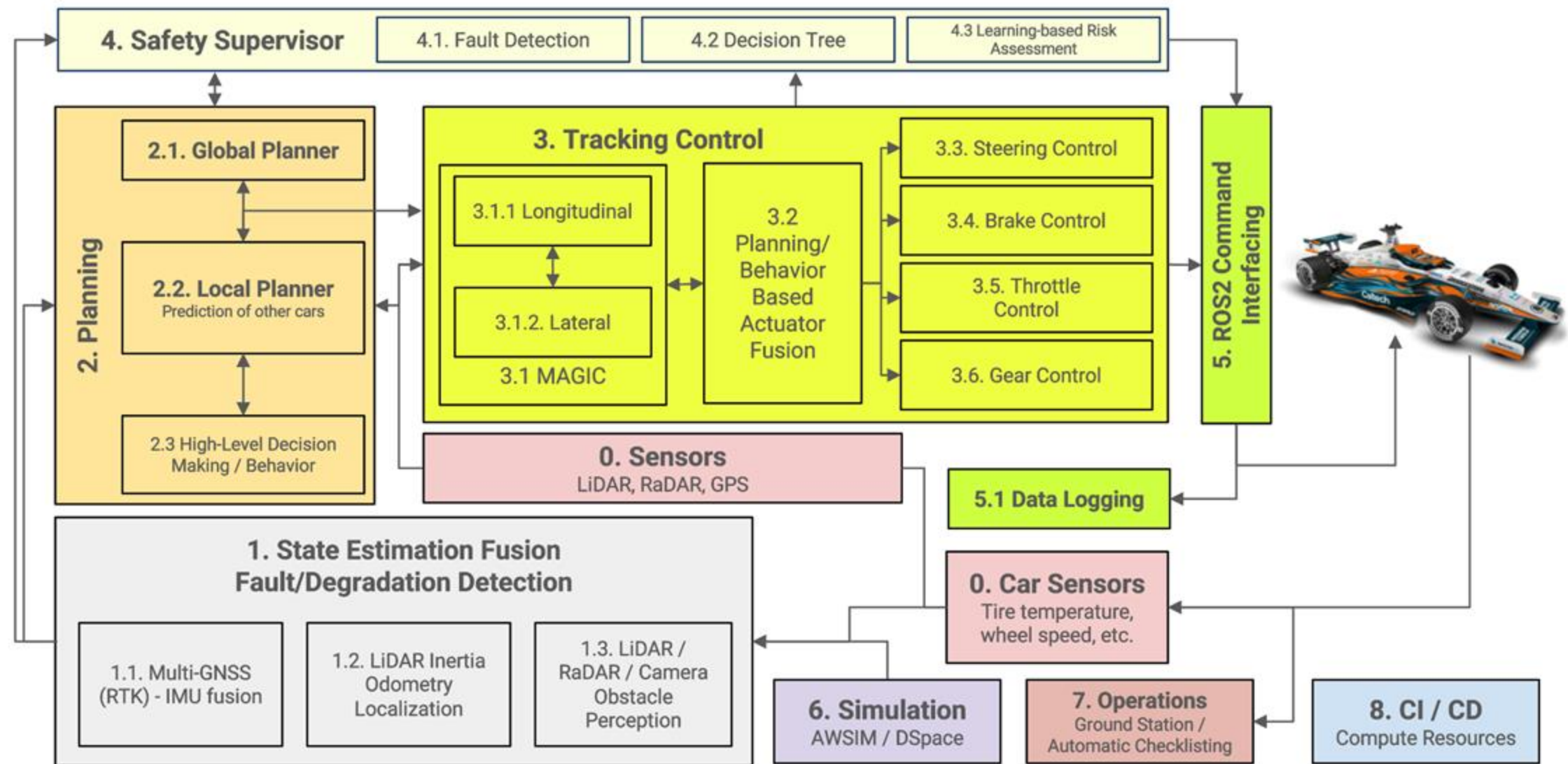
LAP 2

1:34.023

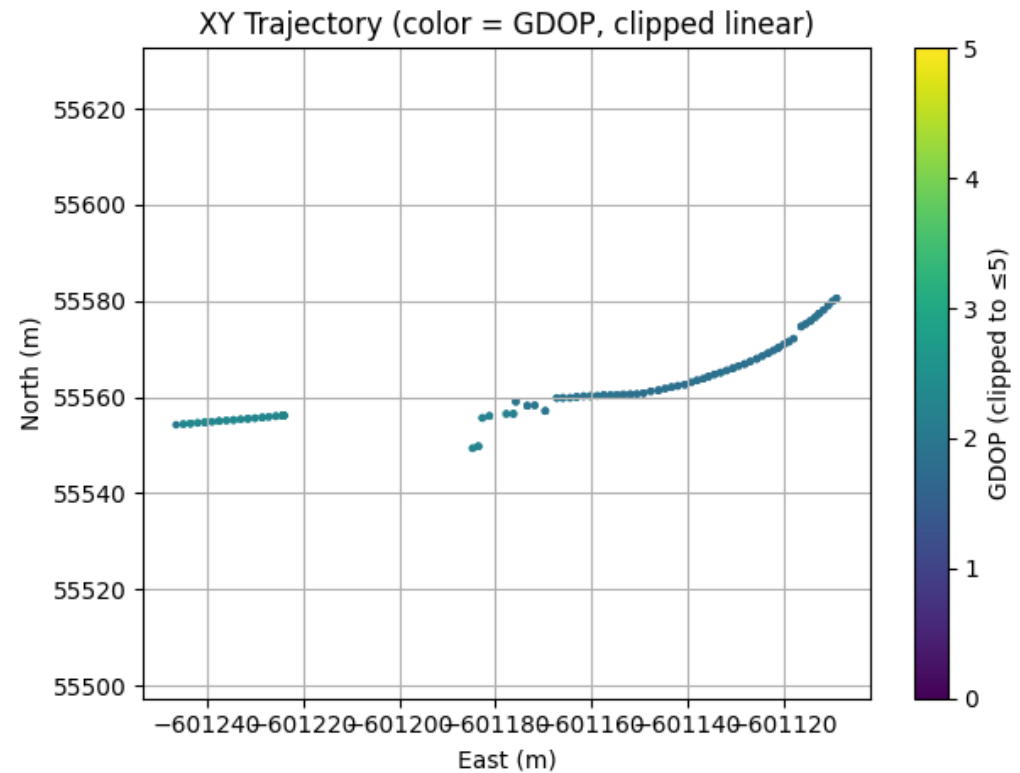


1 2 3



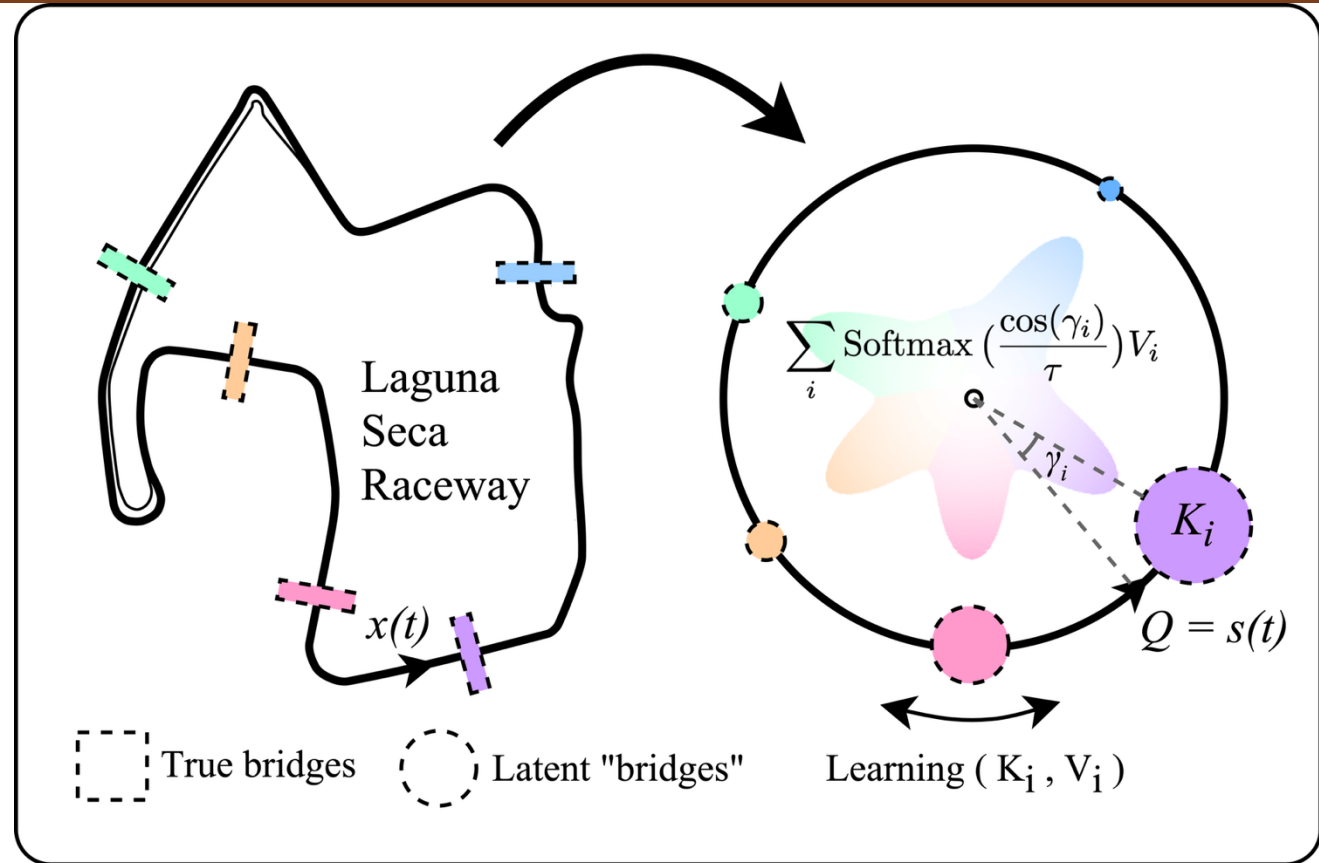
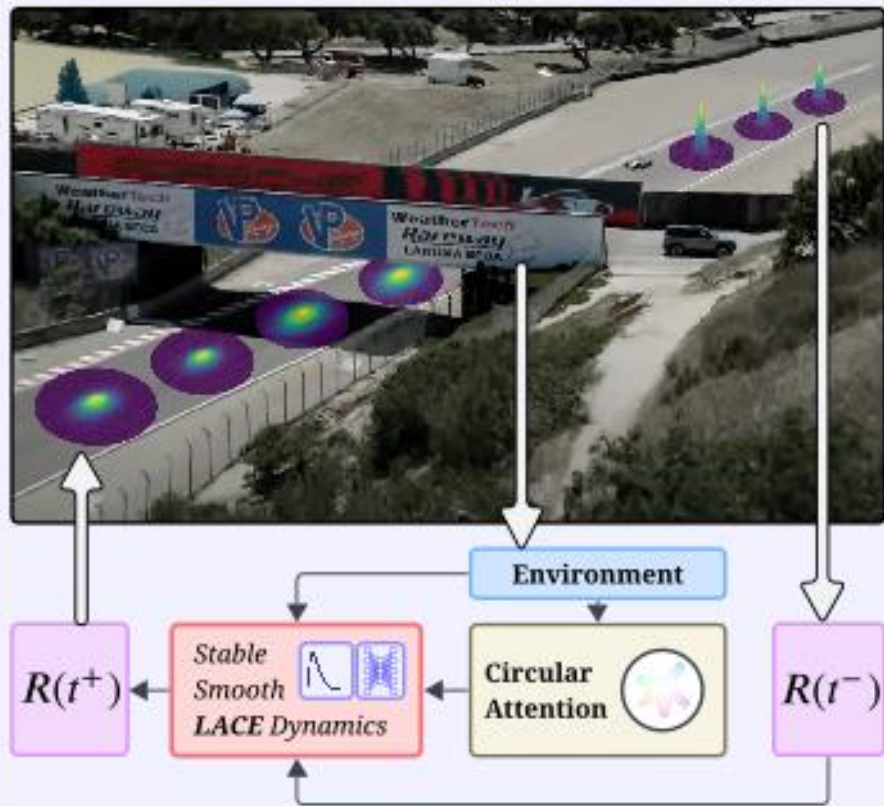


Modeling the dynamics of uncertainty

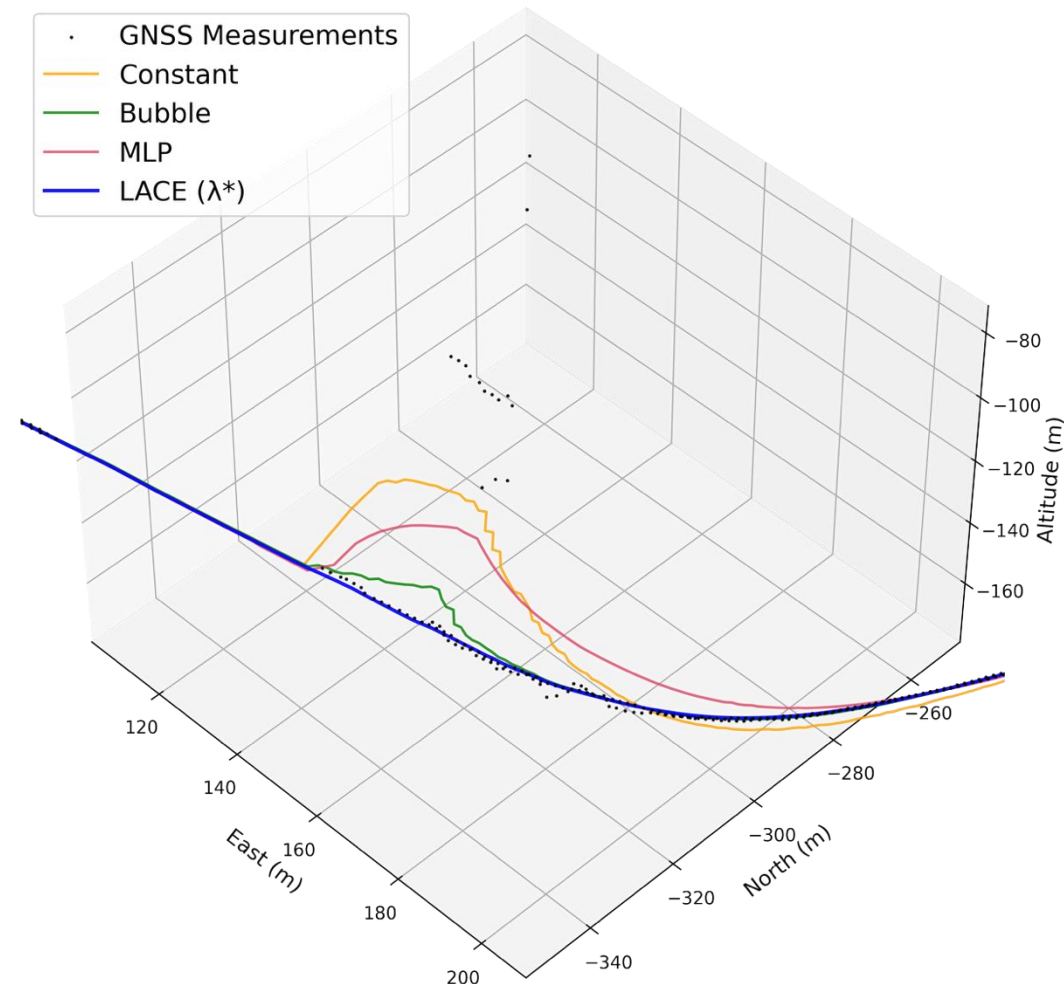
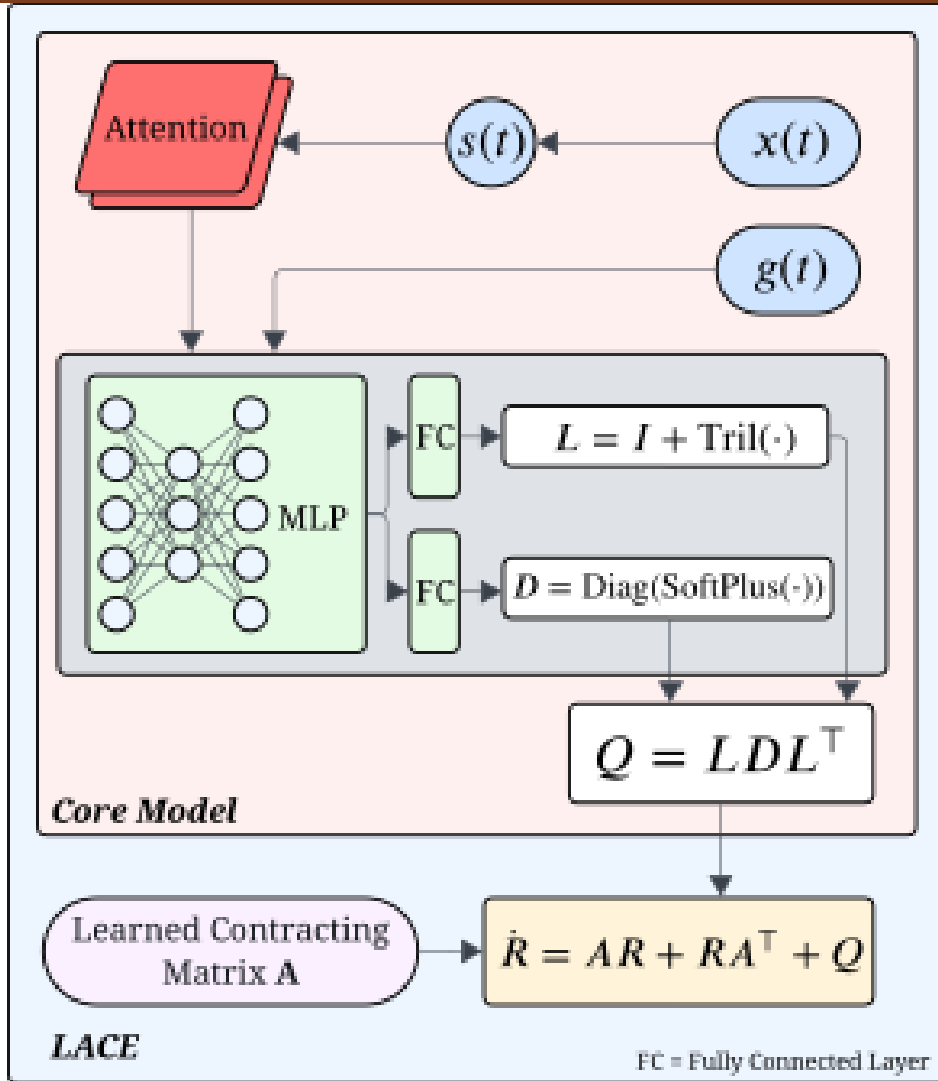




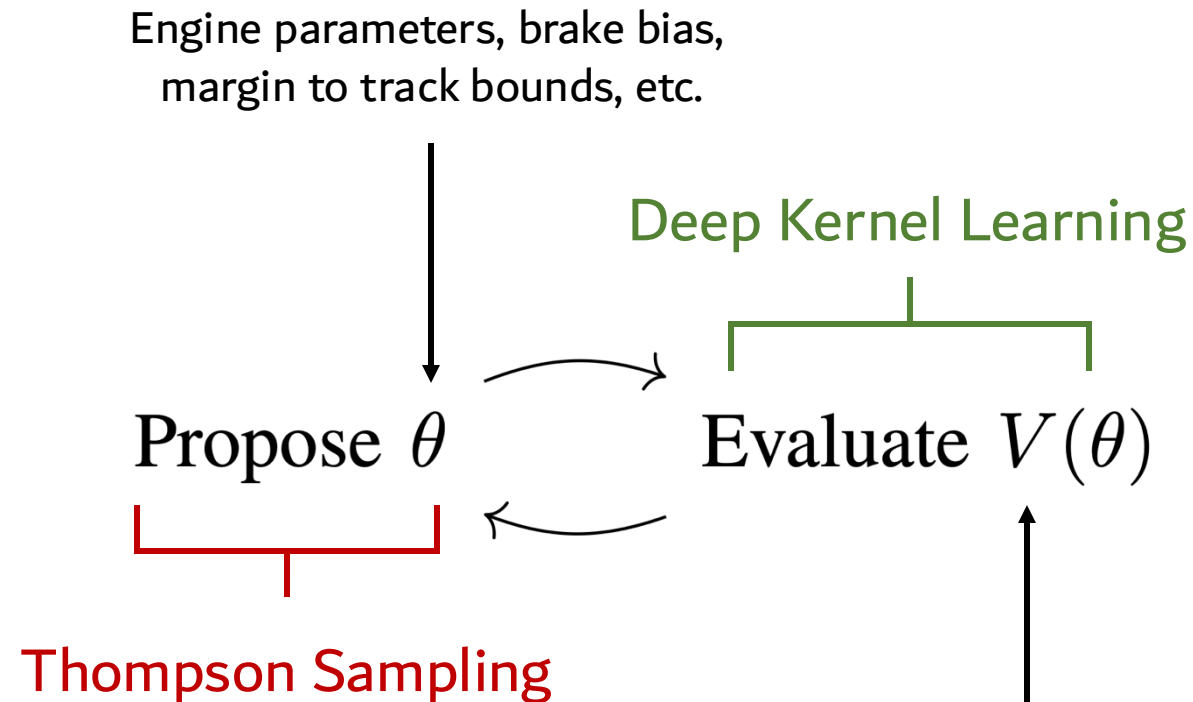
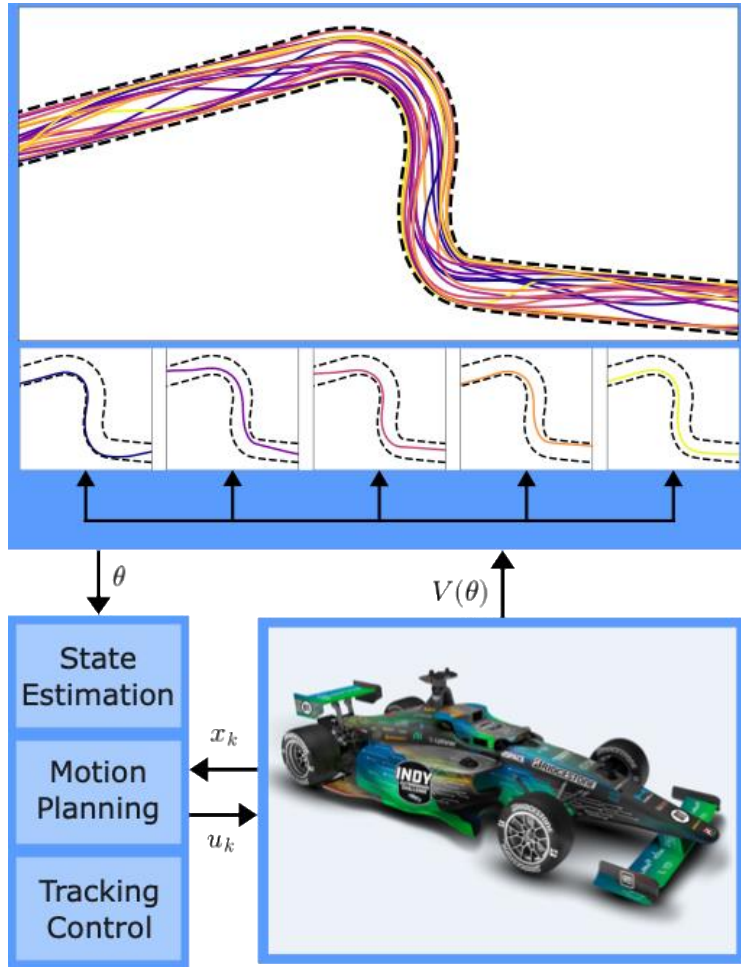
Modeling the dynamics of uncertainty



Modeling the dynamics of uncertainty



Online learning of closed-loop racelines

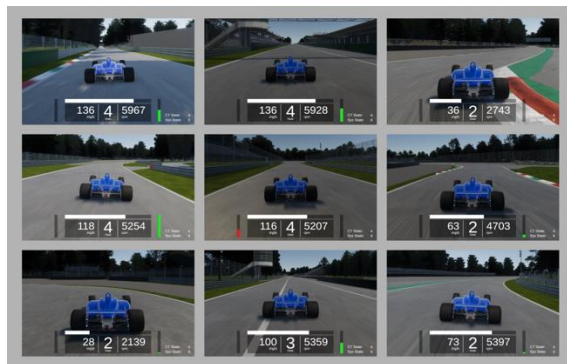


Online learning of closed-loop racelines

Experiments



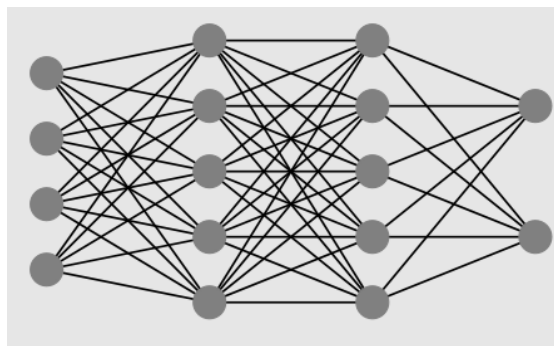
Simulations



Expert data



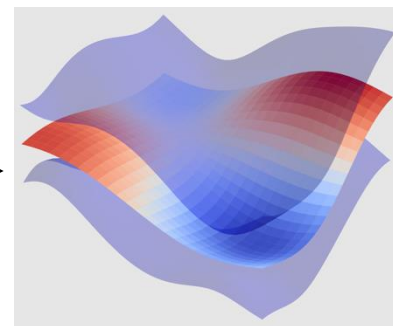
Representation learning



γ

Deep GP with compound kernel

$f_{\gamma}(\theta)$



$$\hat{V}(\theta') \sim N(\mu(\theta'), k^*(f_{\gamma}(\theta), f_{\gamma}(\theta'))))$$

$$k^* = k(f_{\gamma}(\theta), f_{\gamma}(\theta')) + k(\theta, \theta')$$

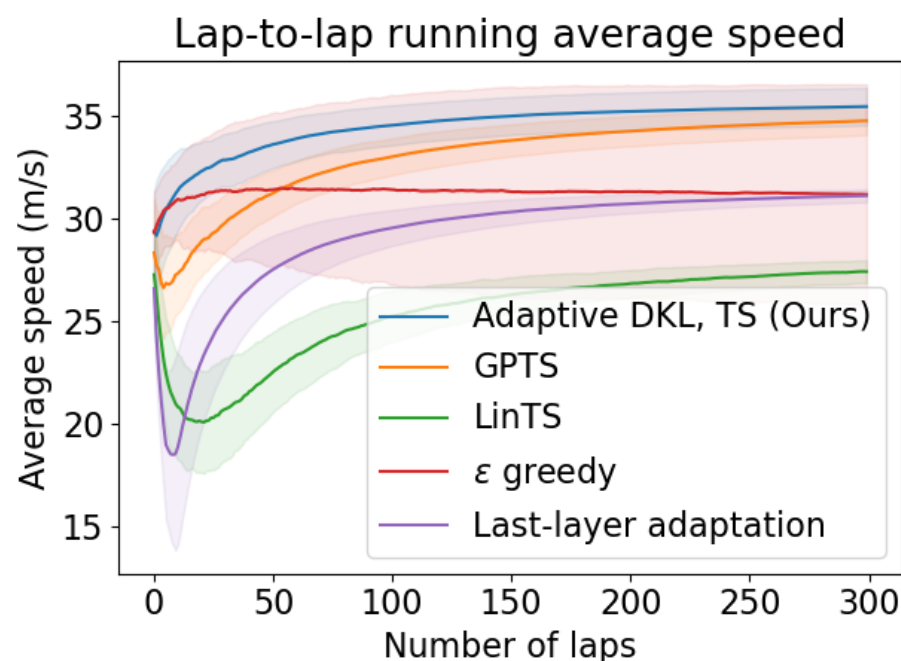
Offline

Online



Online learning of closed-loop racelines

- Our approach can leverage strong priors to learn faster.
- Representation learning is key to finding better lap times.
- Since the MPC is the ultimate determinant of safety, this work is directly transferable to hardware.



Embodied Learning

Challenges:

- Data diversity
- Sample efficiency
- Safety

Algorithm

Simulated
Model/Policy

Simulation

Hardware

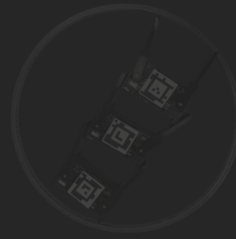
Fine Tuning

Hardware
Model/Policy

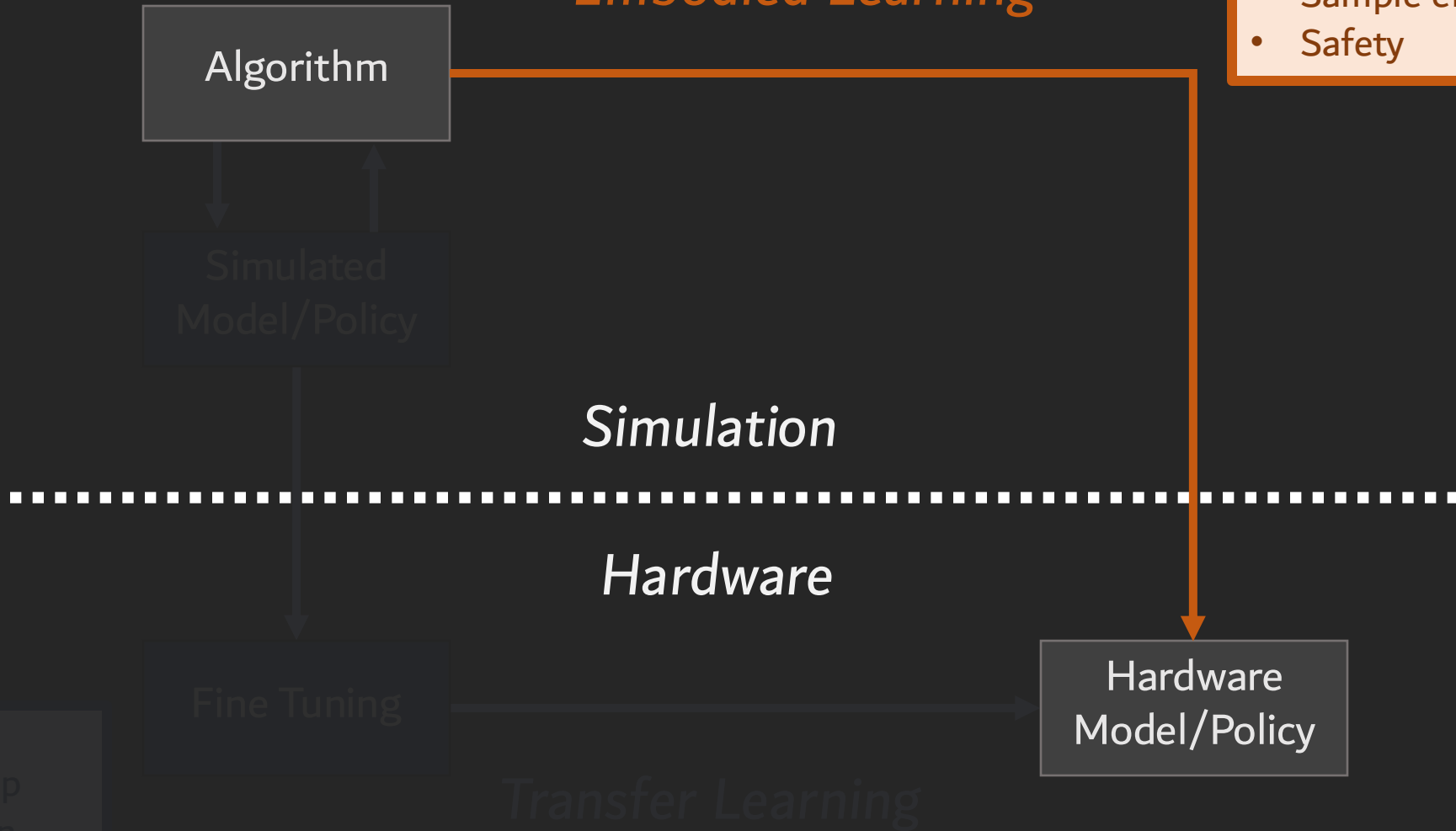
Transfer Learning

Challenges:

- Reality gap
- Perception
- Hardware



Embodied Learning

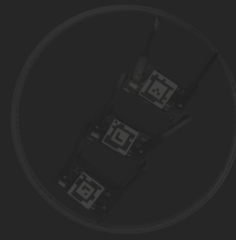


Challenges:

- Data diversity
- Sample efficiency
- Safety

Challenges:

- Reality gap
- Perception
- Hardware



Questions?

