

Data-Efficient Reinforcement Learning for Autonomous Robots

Marc Deisenroth

Centre for Artificial Intelligence
Department of Computer Science
University College London

m.deisenroth@ucl.ac.uk



 @mpd37

IBIS 2019, Nagoya

November 20, 2019



- **Vision:** Autonomous robots support humans in everyday activities ➤ Fast learning and automatic adaptation



- **Vision:** Autonomous robots support humans in everyday activities ➤ Fast learning and automatic adaptation
- Currently: Data-hungry learning or human guidance

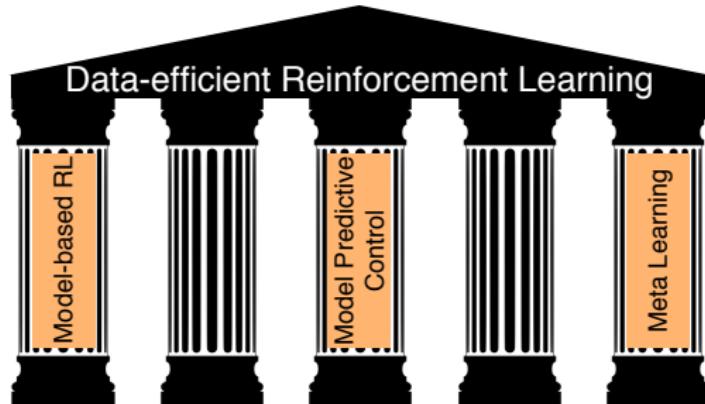


- **Vision:** Autonomous robots support humans in everyday activities ➤ Fast learning and automatic adaptation
- Currently: Data-hungry learning or human guidance

Fully **autonomous learning and decision making with little data** in real-life situations

Data-Efficient Reinforcement Learning

Ability to learn and make decisions in complex domains without requiring large quantities of data



1 Model-based RL

- ▶ Data-efficient decision making

2 Model predictive RL

- ▶ Speed up learning further by online planning

3 Meta learning

- ▶ Generalization of knowledge to new situations



Reinforcement Learning

$$x_{t+1} = f(x_t, u_t) + w, \quad u_t = \pi(x_t, \theta)$$

Diagram illustrating the Reinforcement Learning update rule:

- State**: Represented by x_t (grey box).
- Control**: Represented by u_t (green box).
- Policy**: Represented by π (purple box).
- Policy parameters**: Represented by θ (pink box).

Arrows indicate dependencies:
- x_{t+1} depends on x_t (State) and u_t (Control).
- u_t depends on x_t (State) and θ (Policy parameters).
- π depends on x_t (State) and θ (Policy parameters).
- A red arrow labeled "Transition function" points from x_t to x_{t+1} .

$$\mathbf{x}_{t+1} = f(\mathbf{x}_t, \mathbf{u}_t) + \mathbf{w}, \quad \mathbf{u}_t = \pi(\mathbf{x}_t, \boldsymbol{\theta})$$

↑
State ↑
Control ↑
Policy ↑
Policy parameters
Transition function

Objective (Controller Learning)

Find policy parameters $\boldsymbol{\theta}^*$ that minimize the expected long-term cost

$$J(\boldsymbol{\theta}) = \sum_{t=1}^T \mathbb{E}[c(\mathbf{x}_t) | \boldsymbol{\theta}], \quad p(\mathbf{x}_0) = \mathcal{N}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0).$$

Instantaneous cost $c(\mathbf{x}_t)$, e.g., $\|\mathbf{x}_t - \mathbf{x}_{\text{target}}\|^2$

- ▶ Typical objective in **optimal control** and **reinforcement learning** (Bertsekas, 2005; Sutton & Barto, 1998)

Objective

Minimize expected long-term cost $J(\theta) = \sum_t \mathbb{E}[c(x_t)|\theta]$

PILCO Framework: High-Level Steps

- 1 Probabilistic model for transition function f
 - ▶ System identification

Objective

Minimize expected long-term cost $J(\theta) = \sum_t \mathbb{E}[c(x_t)|\theta]$

PILCO Framework: High-Level Steps

- 1 Probabilistic model for transition function f
 - ▶ System identification
- 2 Compute long-term predictions $p(x_1|\theta), \dots, p(x_T|\theta)$

Objective

Minimize expected long-term cost $J(\theta) = \sum_t \mathbb{E}[c(x_t)|\theta]$

PILCO Framework: High-Level Steps

- 1 Probabilistic model for transition function f
 - ▶ System identification
- 2 Compute long-term predictions $p(x_1|\theta), \dots, p(x_T|\theta)$
- 3 Policy improvement

Objective

Minimize expected long-term cost $J(\theta) = \sum_t \mathbb{E}[c(x_t)|\theta]$

PILCO Framework: High-Level Steps

- 1 Probabilistic model for transition function f
 - ▶ System identification
- 2 Compute long-term predictions $p(x_1|\theta), \dots, p(x_T|\theta)$
- 3 Policy improvement
- 4 Apply controller

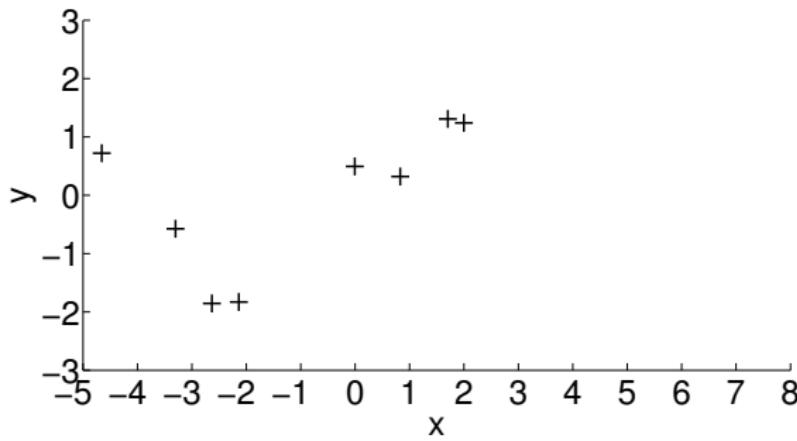
Objective

Minimize expected long-term cost $J(\theta) = \sum_t \mathbb{E}[c(x_t)|\theta]$

PILCO Framework: High-Level Steps

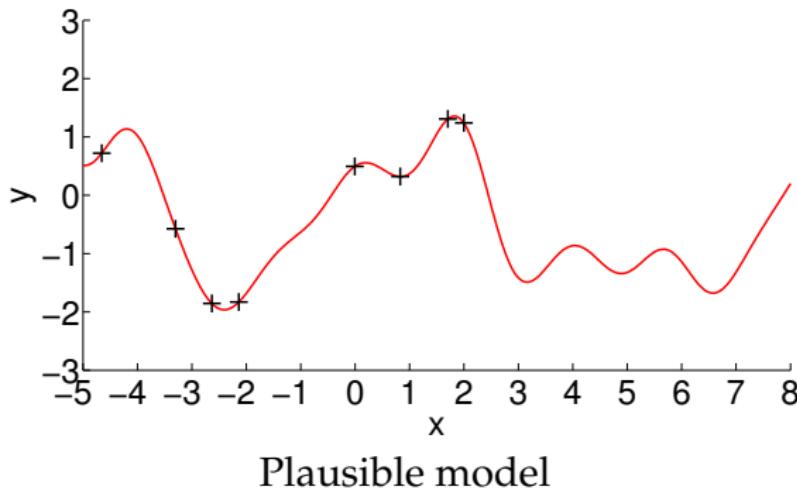
- 1 **Probabilistic model for transition function f**
 - **System identification**
- 2 Compute long-term predictions $p(x_1|\theta), \dots, p(x_T|\theta)$
- 3 Policy improvement
- 4 Apply controller

Model learning problem: Find a function $f : x \mapsto f(x) = y$

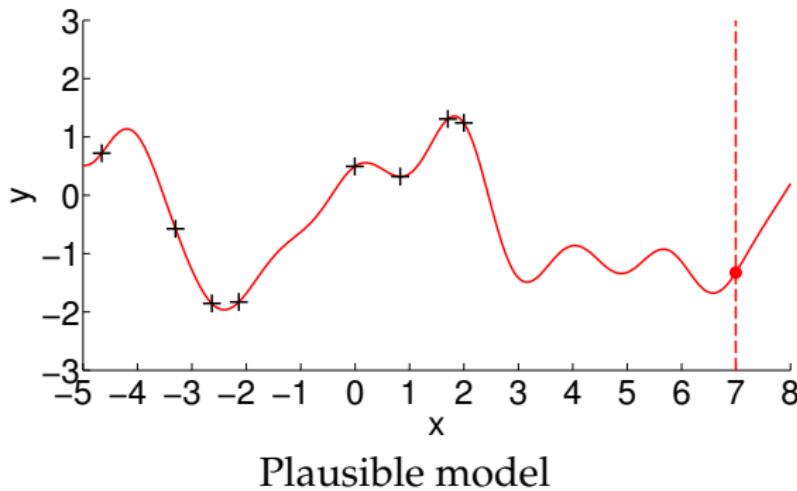


Observed function values

Model learning problem: Find a function $f : x \mapsto f(x) = y$

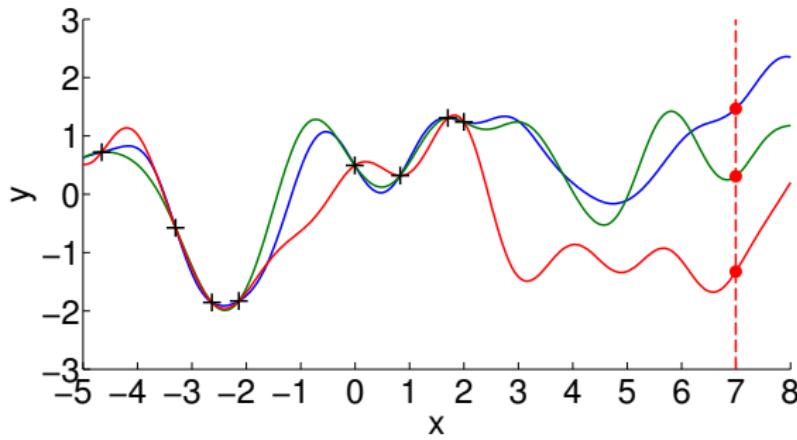


Model learning problem: Find a function $f : x \mapsto f(x) = y$



Predictions? Decision Making?

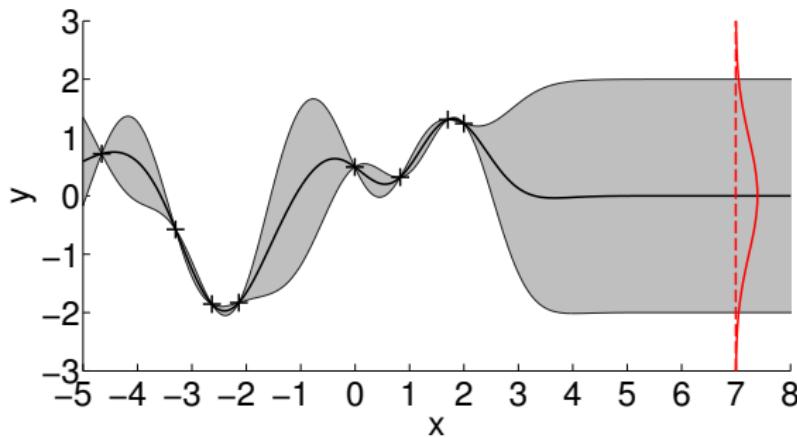
Model learning problem: Find a function $f : x \mapsto f(x) = y$



More plausible models

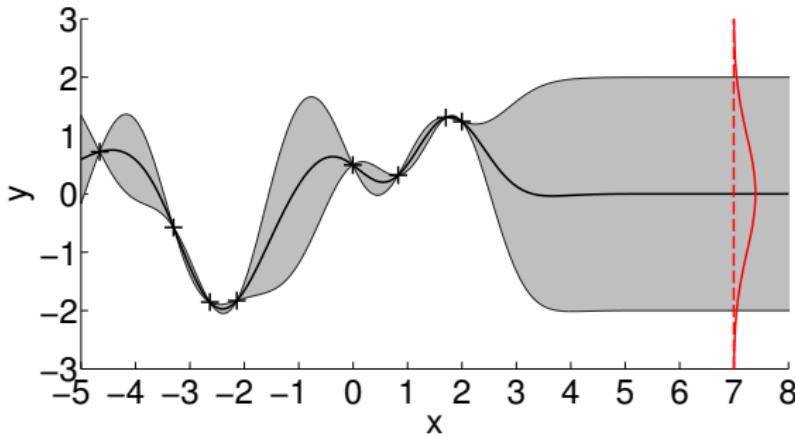
Predictions? Decision Making? Model Errors!

Model learning problem: Find a function $f : x \mapsto f(x) = y$



Distribution over plausible functions

Model learning problem: Find a function $f : x \mapsto f(x) = y$



Distribution over plausible functions

- ▶ Express **uncertainty** about the underlying function to be **robust to model errors**
- ▶ **Gaussian process** for model learning
(Rasmussen & Williams, 2006)

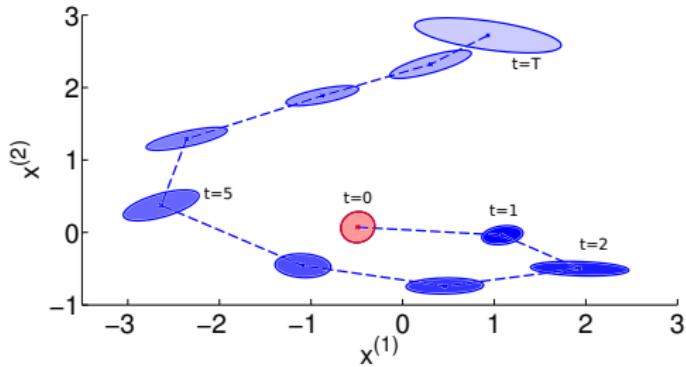
Objective

Minimize expected long-term cost $J(\theta) = \sum_t \mathbb{E}[c(x_t)|\theta]$

PILCO Framework: High-Level Steps

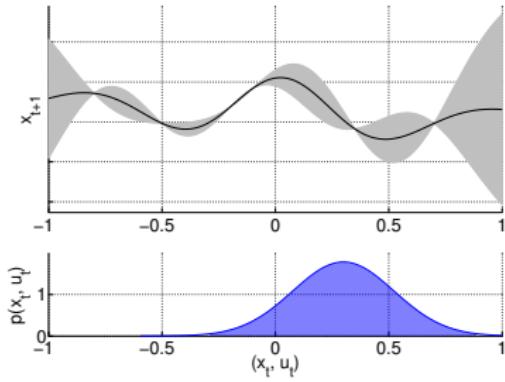
- 1 Probabilistic model for transition function f
 - ▶ System identification
- 2 **Compute long-term predictions** $p(x_1|\theta), \dots, p(x_T|\theta)$
- 3 Policy improvement
- 4 Apply controller

Long-Term Predictions



- Iteratively compute $p(x_1|\theta), \dots, p(x_T|\theta)$

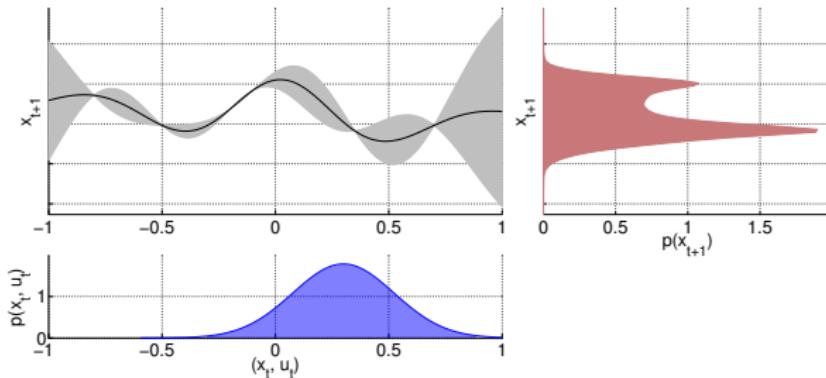
Long-Term Predictions



- Iteratively compute $p(\mathbf{x}_1|\boldsymbol{\theta}), \dots, p(\mathbf{x}_T|\boldsymbol{\theta})$

$$\underbrace{p(\mathbf{x}_{t+1}|\mathbf{x}_t, \mathbf{u}_t)}_{\text{GP prediction}} \quad \underbrace{p(\mathbf{x}_t, \mathbf{u}_t|\boldsymbol{\theta})}_{\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})}$$

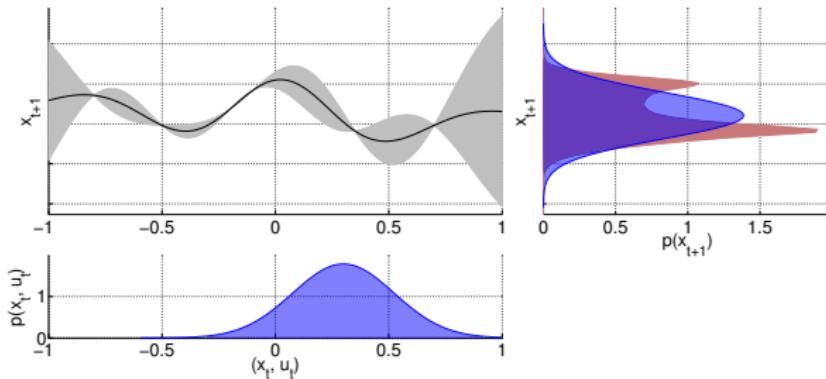
Long-Term Predictions



- Iteratively compute $p(\mathbf{x}_1|\boldsymbol{\theta}), \dots, p(\mathbf{x}_T|\boldsymbol{\theta})$

$$p(\mathbf{x}_{t+1}|\boldsymbol{\theta}) = \iiint \underbrace{p(\mathbf{x}_{t+1}|\mathbf{x}_t, \mathbf{u}_t)}_{\text{GP prediction}} \underbrace{p(\mathbf{x}_t, \mathbf{u}_t|\boldsymbol{\theta})}_{\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})} d\mathbf{f} d\mathbf{x}_t d\mathbf{u}_t$$

Long-Term Predictions



- Iteratively compute $p(\mathbf{x}_1|\boldsymbol{\theta}), \dots, p(\mathbf{x}_T|\boldsymbol{\theta})$

$$p(\mathbf{x}_{t+1}|\boldsymbol{\theta}) = \iiint \underbrace{p(\mathbf{x}_{t+1}|\mathbf{x}_t, \mathbf{u}_t)}_{\text{GP prediction}} \underbrace{p(\mathbf{x}_t, \mathbf{u}_t|\boldsymbol{\theta})}_{\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})} df d\mathbf{x}_t d\mathbf{u}_t$$

► GP moment matching

(Girard et al., 2002; Quiñonero-Candela et al., 2003)

Deisenroth et al. (IEEE-TPAMI, 2015): *Gaussian Processes for Data-Efficient Learning in Robotics and Control*

Objective

Minimize expected long-term cost $J(\theta) = \sum_t \mathbb{E}[c(x_t)|\theta]$

PILCO Framework: High-Level Steps

- 1 Probabilistic model for transition function f
 - ▶ System identification
- 2 Compute long-term predictions $p(x_1|\theta), \dots, p(x_T|\theta)$
- 3 **Policy improvement**
 - Compute expected long-term cost $J(\theta)$
 - Find parameters θ that minimize $J(\theta)$
- 4 Apply controller

Objective

Minimize expected long-term cost $J(\theta) = \sum_t \mathbb{E}[c(x_t)|\theta]$

- (Gaussian) state distributions $p(x_1|\theta), \dots, p(x_T|\theta)$

Objective

Minimize expected long-term cost $J(\boldsymbol{\theta}) = \sum_t \mathbb{E}[c(\mathbf{x}_t) | \boldsymbol{\theta}]$

- (Gaussian) state distributions $p(\mathbf{x}_1 | \boldsymbol{\theta}), \dots, p(\mathbf{x}_T | \boldsymbol{\theta})$
- Compute

$$\mathbb{E}[c(\mathbf{x}_t) | \boldsymbol{\theta}] = \int c(\mathbf{x}_t) \mathcal{N}(\mathbf{x}_t | \boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t) d\mathbf{x}_t, \quad t = 1, \dots, T,$$

and sum them up to obtain $J(\boldsymbol{\theta})$

Objective

Minimize expected long-term cost $J(\boldsymbol{\theta}) = \sum_t \mathbb{E}[c(\mathbf{x}_t) | \boldsymbol{\theta}]$

- (Gaussian) state distributions $p(\mathbf{x}_1 | \boldsymbol{\theta}), \dots, p(\mathbf{x}_T | \boldsymbol{\theta})$
- Compute

$$\mathbb{E}[c(\mathbf{x}_t) | \boldsymbol{\theta}] = \int c(\mathbf{x}_t) \mathcal{N}(\mathbf{x}_t | \boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t) d\mathbf{x}_t, \quad t = 1, \dots, T,$$

and sum them up to obtain $J(\boldsymbol{\theta})$

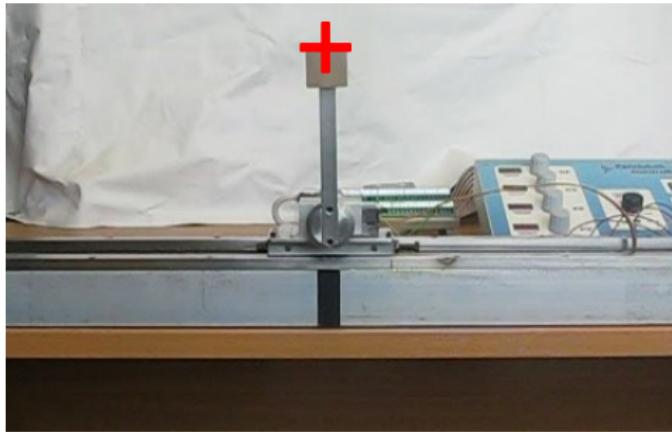
- Analytically compute gradient $dJ(\boldsymbol{\theta})/d\boldsymbol{\theta}$
- Standard gradient-based optimizer (e.g., BFGS) to find $\boldsymbol{\theta}^*$

Objective

Minimize expected long-term cost $J(\theta) = \sum_t \mathbb{E}[c(x_t)|\theta]$

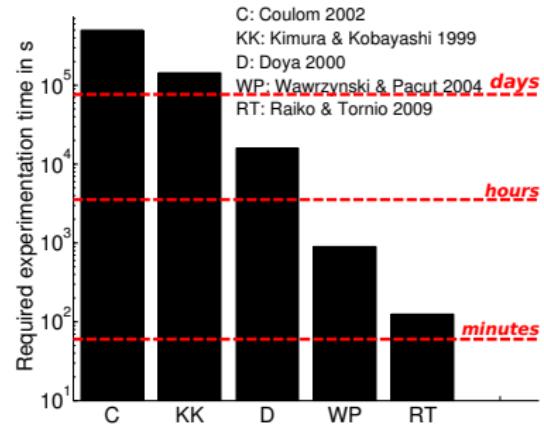
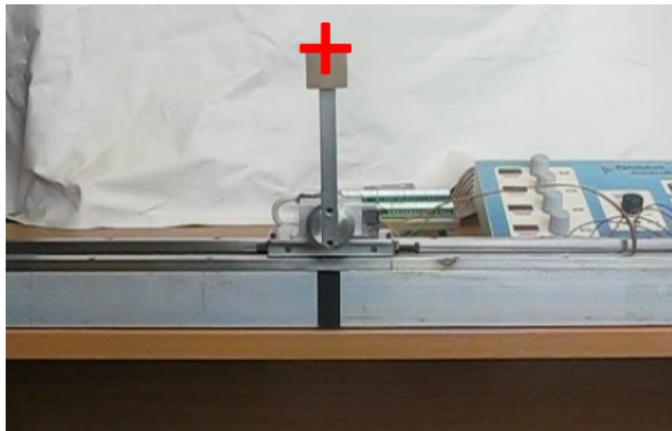
PILCO Framework: High-Level Steps

- 1 Probabilistic model for transition function f
 - ▶ System identification
- 2 Compute long-term predictions $p(x_1|\theta), \dots, p(x_T|\theta)$
- 3 Policy improvement
- 4 **Apply controller**



- Swing up and balance a freely swinging pendulum on a cart
- No knowledge about nonlinear dynamics ➤ Learn from scratch
- Cost function $c(\mathbf{x}) = 1 - \exp(-\|\mathbf{x} - \mathbf{x}_{\text{target}}\|^2)$
- Code: <https://github.com/ICL-SML/pilco-matlab>

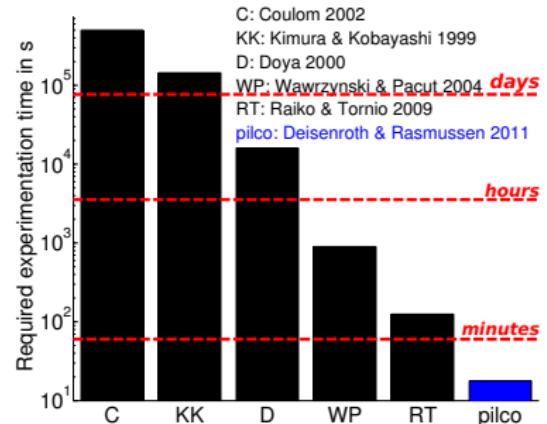
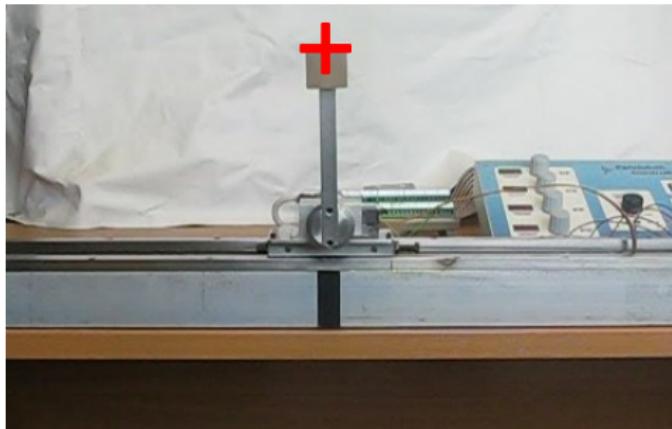
Standard Benchmark: Cart-Pole Swing-up



- Swing up and balance a freely swinging pendulum on a cart
- No knowledge about nonlinear dynamics ➤ Learn from scratch
- Cost function $c(\mathbf{x}) = 1 - \exp(-\|\mathbf{x} - \mathbf{x}_{\text{target}}\|^2)$

- Code: <https://github.com/ICL-SML/pilco-matlab>

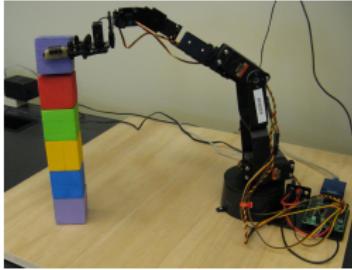
Standard Benchmark: Cart-Pole Swing-up



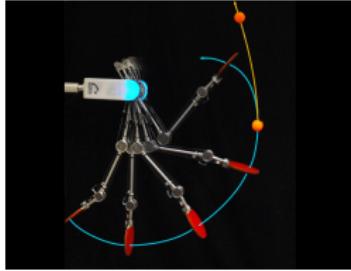
- Swing up and balance a freely swinging pendulum on a cart
- No knowledge about nonlinear dynamics ➤ Learn from scratch
- Cost function $c(\mathbf{x}) = 1 - \exp(-\|\mathbf{x} - \mathbf{x}_{\text{target}}\|^2)$
- **Unprecedented learning speed** compared to state-of-the-art
- Code: <https://github.com/ICL-SML/pilco-matlab>

Deisenroth & Rasmussen (ICML, 2011): PILCO: A Model-based and Data-efficient Approach to Policy Search

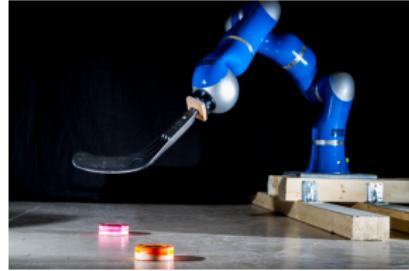
Wide Applicability



with D Fox



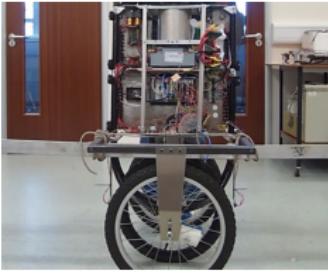
with P Englert, A Paraschos, J Peters



with A Kupcsik, J Peters, G Neumann



B Bischoff (Bosch), ESANN 2013



A McHutchon (U Cambridge)



B Bischoff (Bosch), ECML 2013

► Application to a wide range of robotic systems

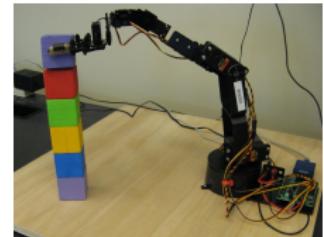
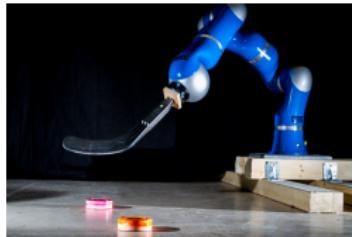
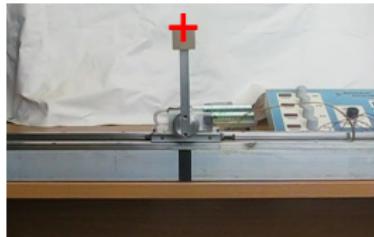
Deisenroth et al. (RSS, 2011): *Learning to Control a Low-Cost Manipulator using Data-efficient Reinforcement Learning*

Englert et al. (ICRA, 2013): *Model-based Imitation Learning by Probabilistic Trajectory Matching*

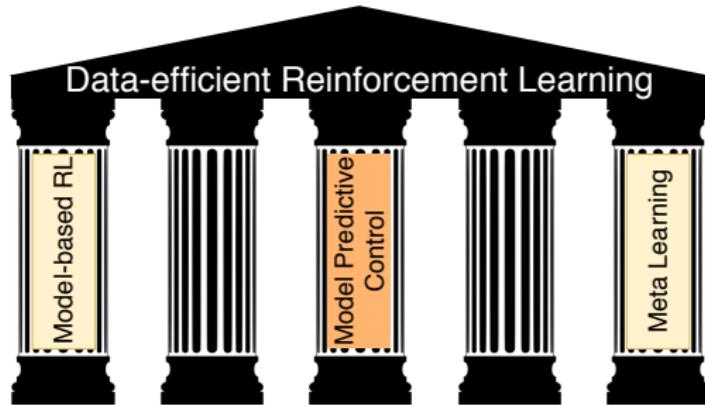
Deisenroth et al. (ICRA, 2014): *Multi-Task Policy Search for Robotics*

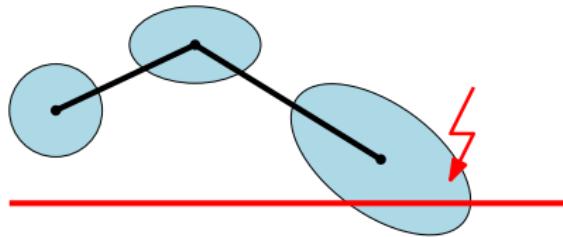
Kupcsik et al. (AIJ, 2017): *Model-based Contextual Policy Search for Data-Efficient Generalization of Robot Skills*

Summary (1)

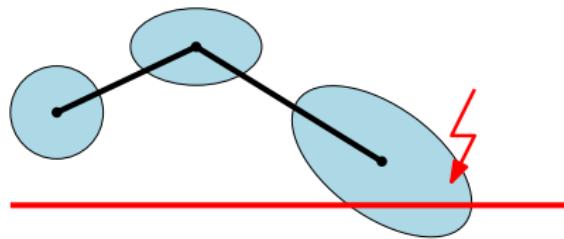


- In robotics, **data-efficient** learning is critical
- Probabilistic, model-based RL approach
 - Reduce model bias
 - Unprecedented learning speed
 - Wide applicability





- Deal with real-world **safety constraints** (states/controls)
- Use probabilistic model to predict whether state constraints are violated (e.g., Sui et al., 2015; Berkenkamp et al., 2017)
- Adjust policy if necessary (during policy learning)

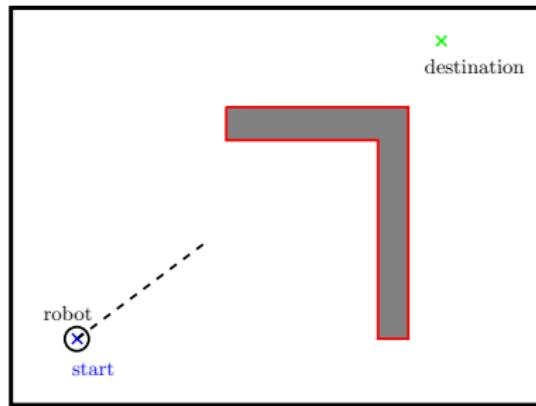


- Deal with real-world **safety constraints** (states/controls)
- Use probabilistic model to predict whether state constraints are violated (e.g., Sui et al., 2015; Berkenkamp et al., 2017)
- Adjust policy if necessary (during policy learning)
- ▶ Safe exploration within an MPC-based RL setting
- ▶ Optimize control signals u_t directly (no policy parameters)

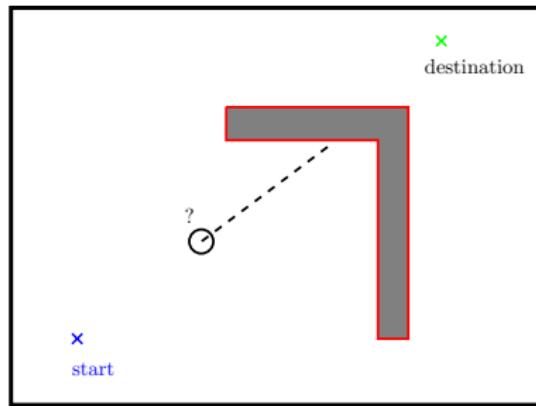
- Idea: Optimize control signals directly (instead of policy parameters)
- Few parameters to optimize ➤ Low-dimensional search space
- Open-loop control
 - No chance of success (with minor model inaccuracies)

- Idea: Optimize control signals directly (instead of policy parameters)
- Few parameters to optimize ➤ Low-dimensional search space
- Open-loop control
 - No chance of success (with minor model inaccuracies)
- Model Predictive Control (MPC) turns this into a closed-loop control approach

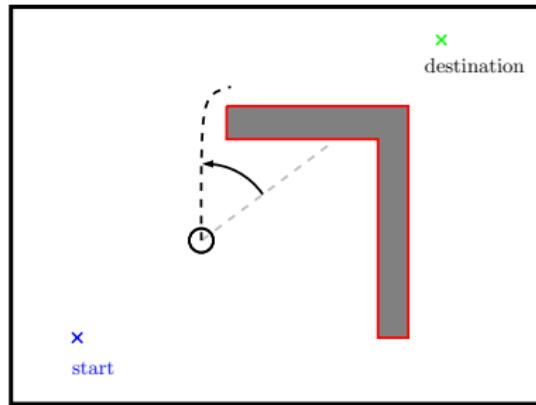
- Idea: Optimize control signals directly (instead of policy parameters)
- Few parameters to optimize ➤ Low-dimensional search space
- Open-loop control
 - No chance of success (with minor model inaccuracies)
- Model Predictive Control (MPC) turns this into a closed-loop control approach
- Positive side-effect: Increase robustness to model errors (online approach) ➤ Increase data efficiency



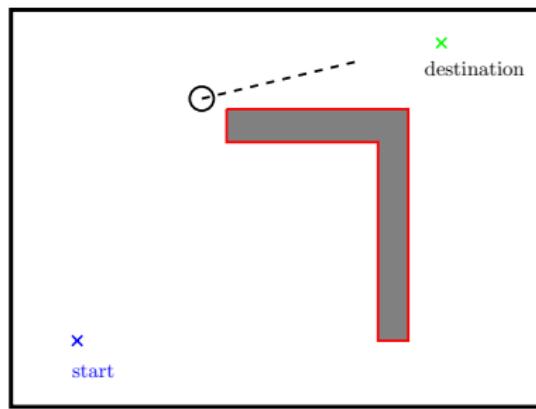
- Given a state x_t , plan (open loop) over a **short horizon** of length H to get an open-loop control sequence $u_{t+0}^*, \dots, u_{t+H-1}^*$
- After transitioning into a new state x_{t+1} , re-plan (as previously): Get $u_{t+1+0}^*, \dots, u_{t+1+H-1}^*$ ► **closed-loop/feedback control**



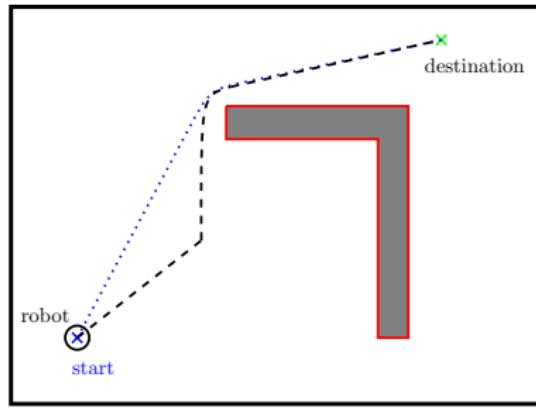
- Given a state x_t , plan (open loop) over a **short horizon** of length H to get an open-loop control sequence $u_{t+0}^*, \dots, u_{t+H-1}^*$
- After transitioning into a new state x_{t+1} , re-plan (as previously): Get $u_{t+1+0}^*, \dots, u_{t+1+H-1}^*$ ► **closed-loop/feedback control**



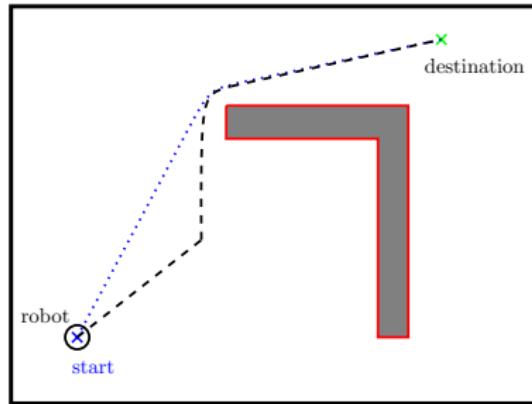
- Given a state x_t , plan (open loop) over a **short horizon** of length H to get an open-loop control sequence $u_{t+0}^*, \dots, u_{t+H-1}^*$
- After transitioning into a new state x_{t+1} , re-plan (as previously): Get $u_{t+1+0}^*, \dots, u_{t+1+H-1}^*$ ► **closed-loop/feedback control**



- Given a state x_t , plan (open loop) over a **short horizon** of length H to get an open-loop control sequence $u_{t+0}^*, \dots, u_{t+H-1}^*$
- After transitioning into a new state x_{t+1} , re-plan (as previously): Get $u_{t+1+0}^*, \dots, u_{t+1+H-1}^*$ ► **closed-loop/feedback control**



- Given a state x_t , plan (open loop) over a **short horizon** of length H to get an open-loop control sequence $u_{t+0}^*, \dots, u_{t+H-1}^*$
- After transitioning into a new state x_{t+1} , re-plan (as previously): Get $u_{t+1+0}^*, \dots, u_{t+1+H-1}^*$ ► **closed-loop/feedback control**



- Given a state x_t , plan (open loop) over a **short horizon** of length H to get an open-loop control sequence $u_{t+0}^*, \dots, u_{t+H-1}^*$
- After transitioning into a new state x_{t+1} , re-plan (as previously): Get $u_{t+1+0}^*, \dots, u_{t+1+H-1}^*$ ► **closed-loop/feedback control**
- Use this within a trial-and-error model-based RL setting

- Learn GP model for transition dynamics
- Repeat (while executing the policy):
 - 1 In current state x_t , determine optimal control sequence u_0^*, \dots, u_{H-1}^*
 - 2 Apply first control u_0^* in state x_t
 - 3 Transition to next state x_{t+1}
 - 4 Update GP transition model with observed transition $((x_t, u_0^*), x_{t+1})$

- Uncertainty propagation is deterministic (GP moment matching)
 - Re-formulate system dynamics:

$$\begin{aligned} \mathbf{z}_{t+1} &= f_{MM}(\mathbf{z}_t, \mathbf{u}_t) \\ \mathbf{z}_t &= \{\boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t\} \quad \text{► Collects moments} \end{aligned}$$

- Uncertainty propagation is deterministic (GP moment matching)
 - Re-formulate system dynamics:

$$\mathbf{z}_{t+1} = f_{MM}(\mathbf{z}_t, \mathbf{u}_t)$$

$\mathbf{z}_t = \{\boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t\}$ ► Collects moments

- Deterministic system function that propagates moments
- Lipschitz continuity (under mild assumptions) implies that we can apply Pontryagin's minimum principle
 - Principled treatment of constraints on controls

Theoretical Results

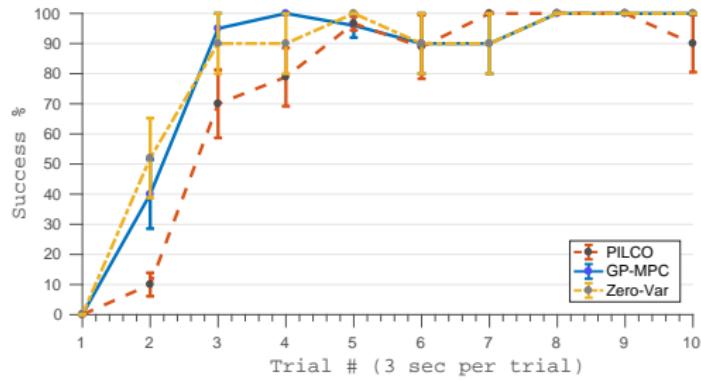
- Uncertainty propagation is deterministic (GP moment matching)
 - Re-formulate system dynamics:

$$\mathbf{z}_{t+1} = f_{MM}(\mathbf{z}_t, \mathbf{u}_t)$$

$$\mathbf{z}_t = \{\boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t\} \quad \text{► Collects moments}$$

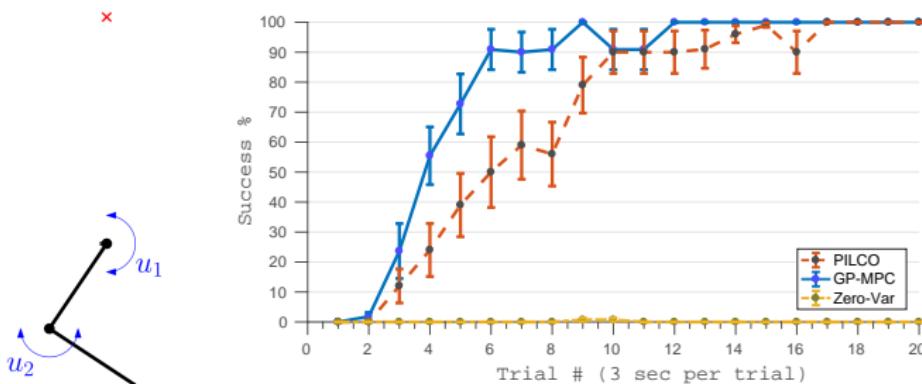
- Deterministic system function that propagates moments
- Lipschitz continuity (under mild assumptions) implies that we can apply Pontryagin's minimum principle
 - Principled treatment of constraints on controls
- Use predictive uncertainty to check violation of state constraints

Learning Speed (Cart Pole)



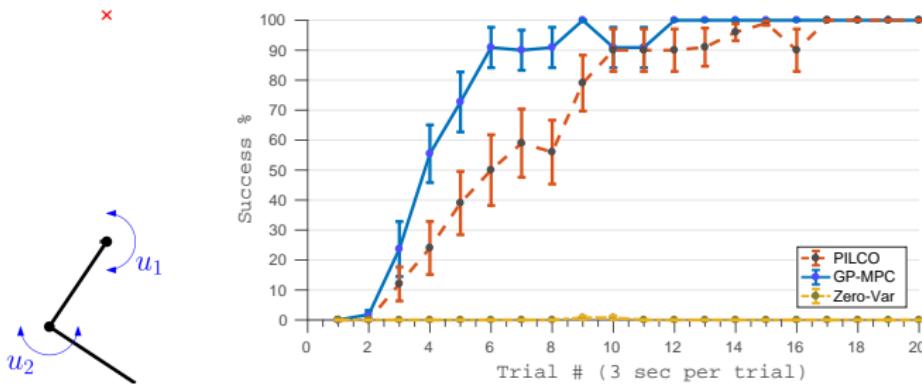
- Zero-Var: Only use the mean of the GP, discard variances for long-term predictions
- MPC: Increased data efficiency (40% less experience required than PILCO)
 - MPC more robust to model inaccuracies than a parametrized feedback controller

Learning Speed (Double Pendulum)



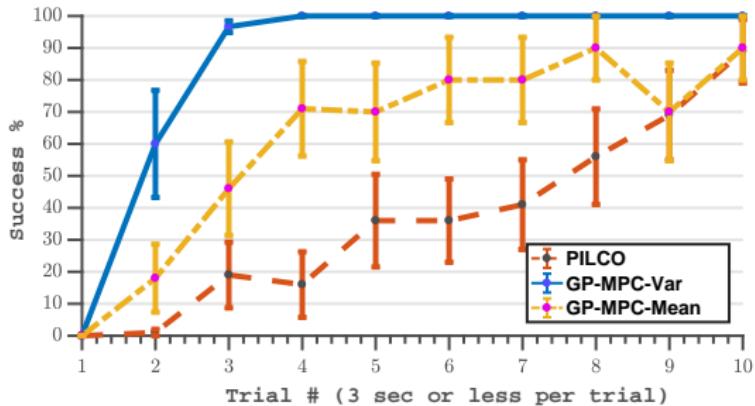
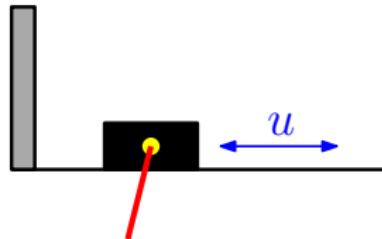
- GP-MPC maintains the same improvement in data efficiency
- Zero-Var fails:
 - Gets stuck in local optimum near start state
 - Insufficient exploration due to lack of uncertainty propagation

Learning Speed (Double Pendulum)



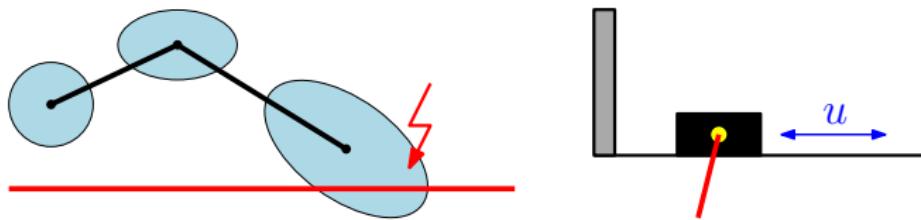
- GP-MPC maintains the same improvement in data efficiency
- Zero-Var fails:
 - Gets stuck in local optimum near start state
 - Insufficient exploration due to lack of uncertainty propagation
- Although MPC is fairly robust to model inaccuracies we cannot get away without uncertainty propagation

Safety Constraints (Cart Pole)



PILCO	16/100	constraint violations
GP-MPC-Mean	21/100	constraint violations
GP-MPC-Var	3/100	constraint violations

► Propagating model uncertainty important for safety



- Probabilistic prediction models for safe exploration
- Uncertainty propagation can be used to reduce violation of safety constraints
- MPC framework increases robustness to model errors
 - ▶ Increased data efficiency





Meta Learning

Generalize knowledge from known tasks to new (related) tasks



Meta Learning

Generalize knowledge from known tasks to new (related) tasks

- Different robot configurations (link lengths, weights, ...)
- Re-use experience gathered so far generalize learning to new dynamics that are similar
 - ▶ Accelerated learning

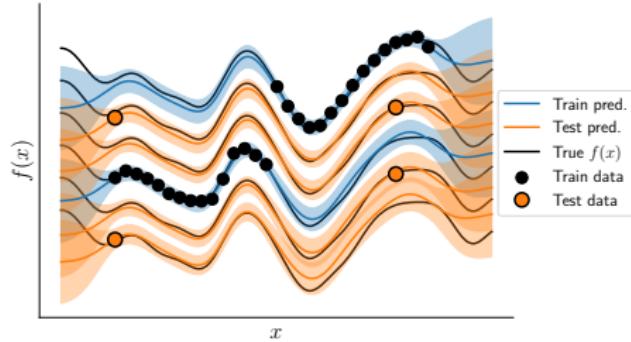
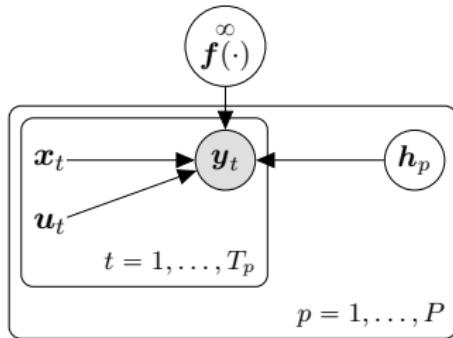


- Separate global and task-specific properties
- Shared global parameters describe general dynamics
- Describe task-specific (local) configurations with latent variable



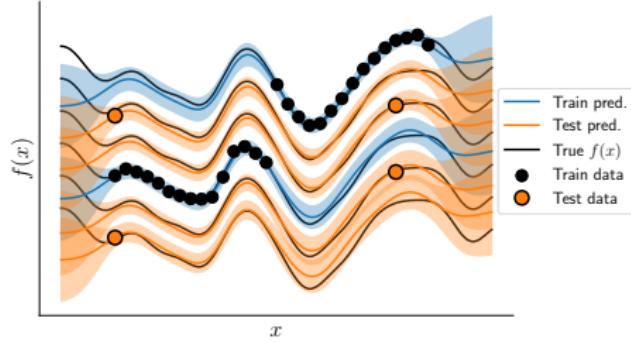
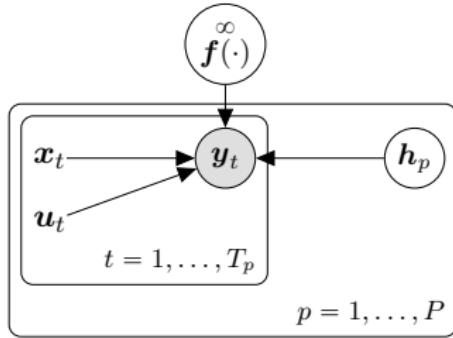
- Separate global and task-specific properties
- Shared global parameters describe general dynamics
- Describe task-specific (local) configurations with latent variable
- Online variational inference of (unseen) configurations
- Few-shot model-based RL

Meta Model Learning with Latent Variables



$$\mathbf{y}_t = f(\mathbf{x}_t, \mathbf{u}_t, \mathbf{h}_p)$$

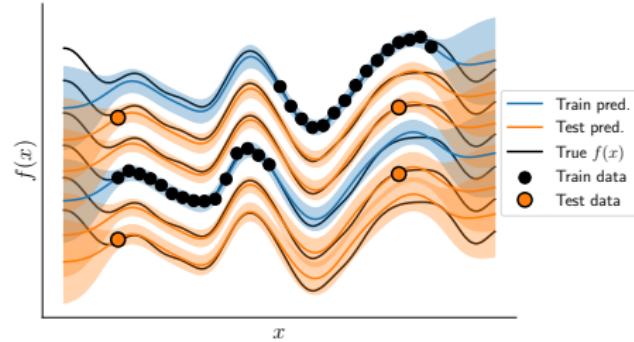
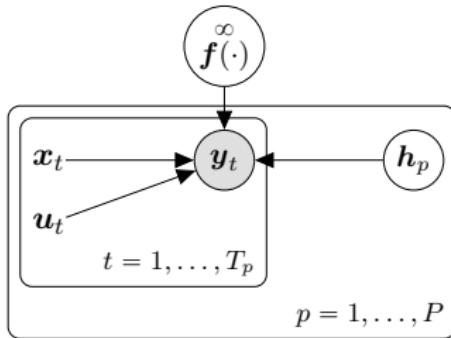
Meta Model Learning with Latent Variables



$$\mathbf{y}_t = \mathbf{f}(\mathbf{x}_t, \mathbf{u}_t, \mathbf{h}_p)$$

- GP captures global properties of the dynamics

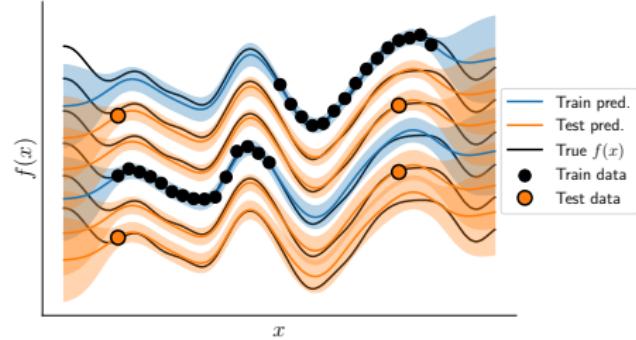
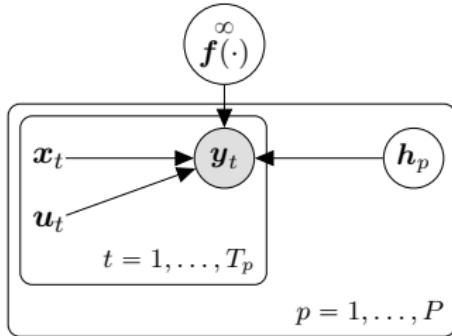
Meta Model Learning with Latent Variables



$$y_t = f(x_t, u_t, h_p)$$

- GP captures global properties of the dynamics
- Latent variable h_p describes local configuration
 - ▶ Variational inference to find a posterior on latent configuration

Meta Model Learning with Latent Variables

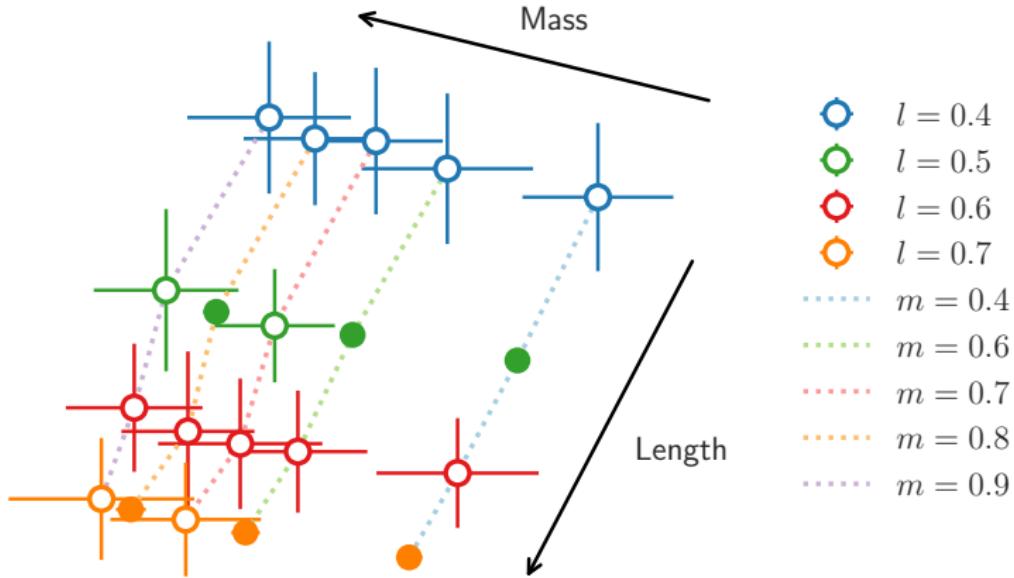


$$y_t = f(x_t, u_t, h_p)$$

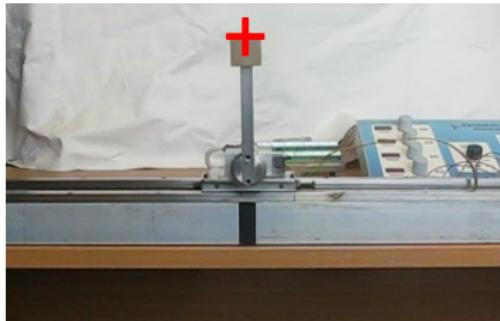
- GP captures global properties of the dynamics
- Latent variable h_p describes local configuration
 - ▶ Variational inference to find a posterior on latent configuration
- Fast online inference of new configurations (no model re-training required)

Sæmundsson et al. (UAI, 2018): *Meta Reinforcement Learning with Latent Variable Gaussian Processes*

Latent Embeddings



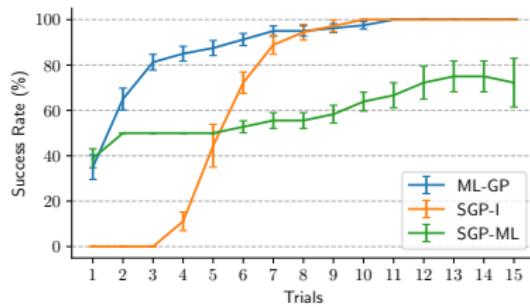
- Latent variable h encodes length l and mass m of the cart pole
- 6 training tasks, 14 held-out test tasks



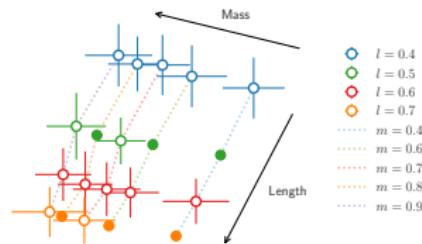
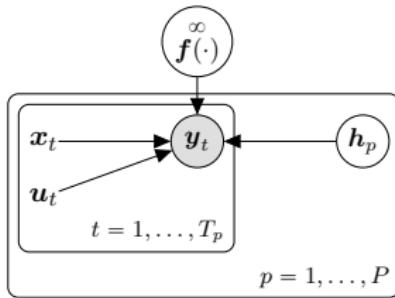
- Pre-trained on 6 training configurations until solved

Model	Training (s)	Description
Independent	16.1 ± 0.4	Independent GP-MPC
Aggregated	23.7 ± 1.4	Aggregated experience (no latents)
Meta learning	15.1 ± 0.5	Aggregated experience (with latents)

► Meta learning can help speeding up RL

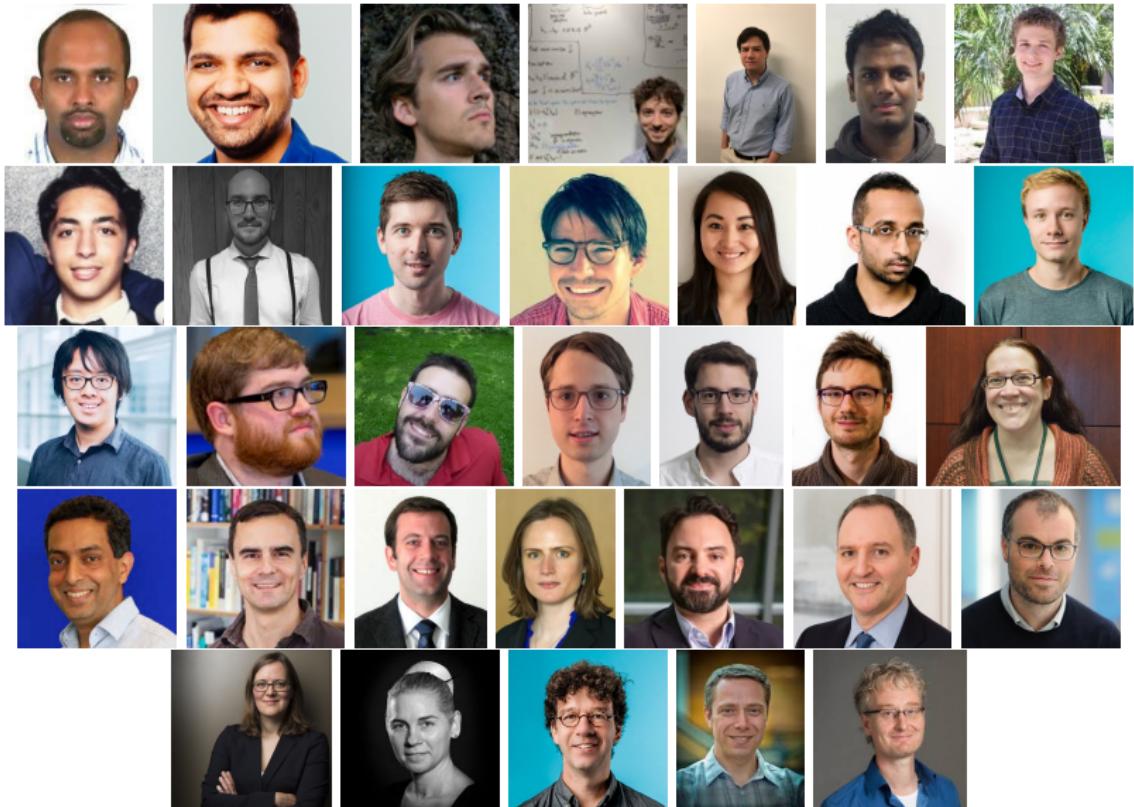


- Few-shot generalization on 4 unseen configurations
 - Success: solve all 10 (6 training + 4 test) tasks
 - Meta learning: blue
 - Independent (GP-MPC): orange
 - Aggregated experience model (no latents): green
- **Meta RL generalizes well to unseen tasks**

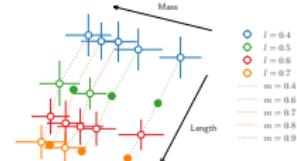
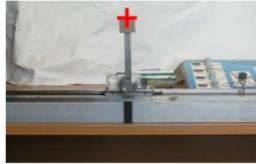


- Generalize knowledge from known situations to unseen ones
 - ▶ **Few-shot learning**
- Latent variable can be used to **infer task similarities**
- Significant speed-up in model learning and model-based RL

Team and Collaborators

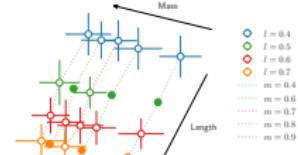
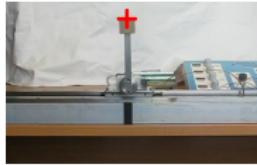


Wrap-Up



- **Data efficiency** is a practical challenge for autonomous robots
- Three pillars of data-efficient reinforcement learning for autonomous robots
 - 1 **Model-based reinforcement learning** with learned probabilistic models for fast learning from scratch
 - 2 **Model predictive control** with learned dynamics models accelerate learning and allow for safe exploration
 - 3 **Meta learning** using latent variables to generalize knowledge to new situations
- **Key to success:** Probabilistic modeling and Bayesian inference

Wrap-Up



- **Data efficiency** is a practical challenge for autonomous robots
- Three pillars of data-efficient reinforcement learning for autonomous robots
 - 1 **Model-based reinforcement learning** with learned probabilistic models for fast learning from scratch
 - 2 **Model predictive control** with learned dynamics models accelerate learning and allow for safe exploration
 - 3 **Meta learning** using latent variables to generalize knowledge to new situations
- **Key to success:** Probabilistic modeling and Bayesian inference

Thank you for your attention

ありがとうございました

References I

- [1] F. Berkenkamp, M. Turchetta, A. P. Schoellig, and A. Krause. Safe Model-based Reinforcement Learning with Stability Guarantees. In *Advances in Neural Information Processing Systems*, 2017.
- [2] D. P. Bertsekas. *Dynamic Programming and Optimal Control*, volume 1 of *Optimization and Computation Series*. Athena Scientific, Belmont, MA, USA, 3rd edition, 2005.
- [3] D. P. Bertsekas. *Dynamic Programming and Optimal Control*, volume 2 of *Optimization and Computation Series*. Athena Scientific, Belmont, MA, USA, 3rd edition, 2007.
- [4] B. Bischoff, D. Nguyen-Tuong, T. Koller, H. Markert, and A. Knoll. Learning Throttle Valve Control Using Policy Search. In *Proceedings of the European Conference on Machine Learning and Knowledge Discovery in Databases*, 2013.
- [5] M. P. Deisenroth, P. Englert, J. Peters, and D. Fox. Multi-Task Policy Search for Robotics. In *Proceedings of the IEEE International Conference on Robotics and Automation*, 2014.
- [6] M. P. Deisenroth, D. Fox, and C. E. Rasmussen. Gaussian Processes for Data-Efficient Learning in Robotics and Control. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37(2):408–423, 2015.
- [7] M. P. Deisenroth and C. E. Rasmussen. PILCO: A Model-Based and Data-Efficient Approach to Policy Search. In *Proceedings of the International Conference on Machine Learning*, 2011.
- [8] M. P. Deisenroth, C. E. Rasmussen, and D. Fox. Learning to Control a Low-Cost Manipulator using Data-Efficient Reinforcement Learning. In *Proceedings of Robotics: Science and Systems*, Los Angeles, CA, USA, June 2011.
- [9] P. Englert, A. Paraschos, J. Peters, and M. P. Deisenroth. Model-based Imitation Learning by Probabilistic Trajectory Matching. In *Proceedings of the IEEE International Conference on Robotics and Automation*, 2013.
- [10] P. Englert, A. Paraschos, J. Peters, and M. P. Deisenroth. Probabilistic Model-based Imitation Learning. *Adaptive Behavior*, 21:388–403, 2013.
- [11] A. Girard, C. E. Rasmussen, and R. Murray-Smith. Gaussian Process Priors with Uncertain Inputs: Multiple-Step Ahead Prediction. Technical Report TR-2002-119, University of Glasgow, 2002.
- [12] S. Kamthe and M. P. Deisenroth. Data-Efficient Reinforcement Learning with Probabilistic Model Predictive Control. In *Proceedings of the International Conference on Artificial Intelligence and Statistics*, 2018.

References II

- [13] A. Kupcsik, M. P. Deisenroth, J. Peters, L. A. Poha, P. Vadakkepata, and G. Neumann. Model-based Contextual Policy Search for Data-Efficient Generalization of Robot Skills. *Artificial Intelligence*, 2017.
- [14] T. X. Nghiem and C. N. Jones. Data-driven Demand Response Modeling and Control of Buildings with Gaussian Processes. In *Proceedings of the American Control Conference*, 2017.
- [15] J. Quiñonero-Candela, A. Girard, J. Larsen, and C. E. Rasmussen. Propagation of Uncertainty in Bayesian Kernel Models—Application to Multiple-Step Ahead Forecasting. In *IEEE International Conference on Acoustics, Speech and Signal Processing*, volume 2, pages 701–704, Apr. 2003.
- [16] C. E. Rasmussen and C. K. I. Williams. *Gaussian Processes for Machine Learning*. Adaptive Computation and Machine Learning. The MIT Press, Cambridge, MA, USA, 2006.
- [17] S. Sæmundsson, K. Hofmann, and M. P. Deisenroth. Meta Reinforcement Learning with Latent Variable Gaussian Processes. In *Proceedings of the Conference on Uncertainty in Artificial Intelligence*, 2018.
- [18] Y. Sui, A. Gotovos, J. W. Burdick, and A. Krause. Safe Exploration for Optimization with Gaussian Processes. In *Proceedings of the International Conference on Machine Learning*, 2015.
- [19] M. K. Titsias. Variational Learning of Inducing Variables in Sparse Gaussian Processes. In *Proceedings of the International Conference on Artificial Intelligence and Statistics*, 2009.

■ Controller:

$$\tilde{\pi}(\mathbf{x}, \boldsymbol{\theta}) = \sum_{k=1}^K w_k \exp \left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^\top \boldsymbol{\Lambda} (\mathbf{x} - \boldsymbol{\mu}_k) \right)$$
$$u = \pi(\mathbf{x}, \boldsymbol{\theta}) = u_{\max} \sigma(\tilde{\pi}(\mathbf{x}, \boldsymbol{\theta})) \in [-u_{\max}, u_{\max}] ,$$

where σ is a squashing function.

■ Controller:

$$\tilde{\pi}(\mathbf{x}, \boldsymbol{\theta}) = \sum_{k=1}^K w_k \exp \left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^\top \boldsymbol{\Lambda} (\mathbf{x} - \boldsymbol{\mu}_k) \right)$$
$$u = \pi(\mathbf{x}, \boldsymbol{\theta}) = u_{\max} \sigma(\tilde{\pi}(\mathbf{x}, \boldsymbol{\theta})) \in [-u_{\max}, u_{\max}] ,$$

where σ is a squashing function.

■ Parameters:

$$\boldsymbol{\theta} := \{w_k, \boldsymbol{\mu}_k, \boldsymbol{\Lambda}\}$$

- Controller:

$$\tilde{\pi}(\mathbf{x}, \boldsymbol{\theta}) = \sum_{k=1}^K w_k \exp \left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^\top \boldsymbol{\Lambda} (\mathbf{x} - \boldsymbol{\mu}_k) \right)$$
$$u = \pi(\mathbf{x}, \boldsymbol{\theta}) = u_{\max} \sigma(\tilde{\pi}(\mathbf{x}, \boldsymbol{\theta})) \in [-u_{\max}, u_{\max}] ,$$

where σ is a squashing function.

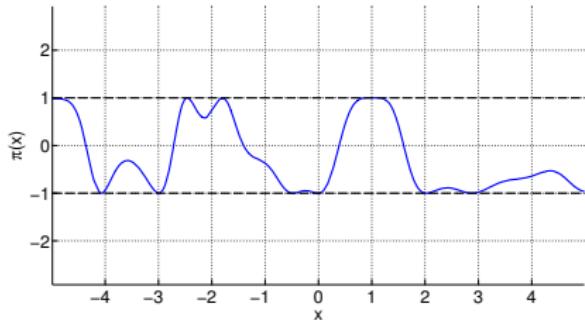
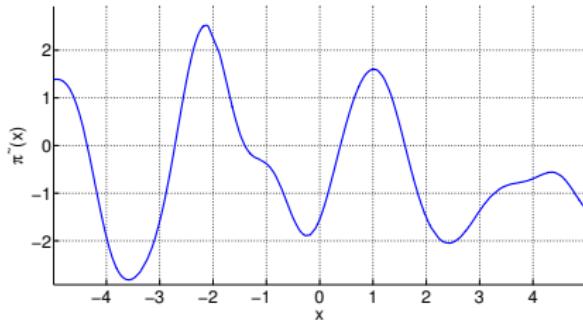
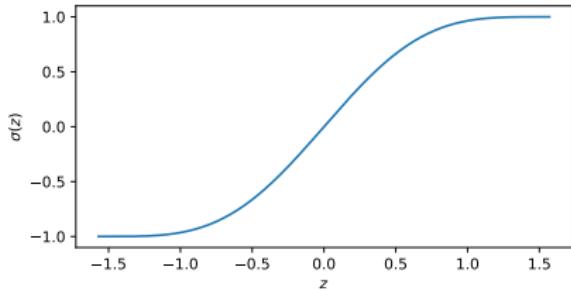
- Parameters:

$$\boldsymbol{\theta} := \{w_k, \boldsymbol{\mu}_k, \boldsymbol{\Lambda}\}$$

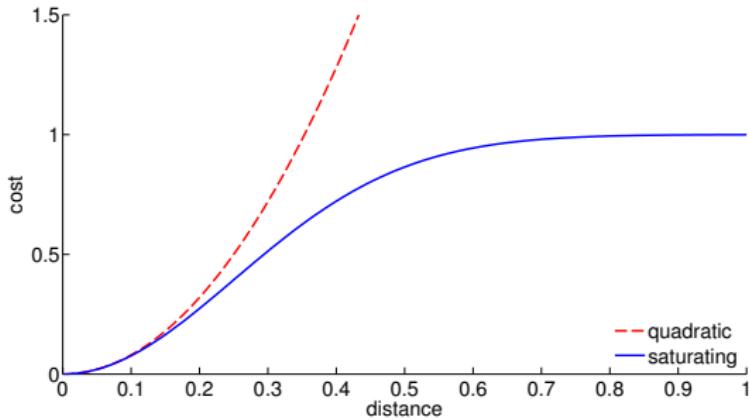
- Squashing function:

$$\sigma(z) = \frac{9}{8} \sin(z) + \frac{1}{8} \sin(3z)$$

Squashing Function

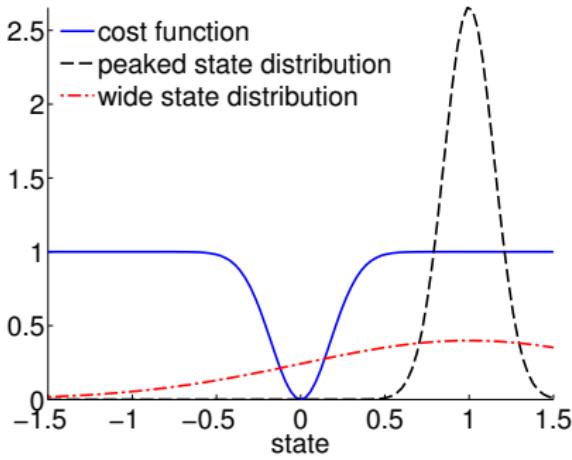


- **Quadratic cost** $c(\mathbf{x}) = (\mathbf{x} - \mathbf{x}_{\text{target}})^{\top} \mathbf{W} (\mathbf{x} - \mathbf{x}_{\text{target}})$
- **Saturating cost** $c(\mathbf{x}) = 1 - \exp \left(- (\mathbf{x} - \mathbf{x}_{\text{target}})^{\top} \mathbf{W} (\mathbf{x} - \mathbf{x}_{\text{target}}) \right)$



- Quadratic cost pays a lot of attention to states “far away”
 - ▶ Bad idea for exploration

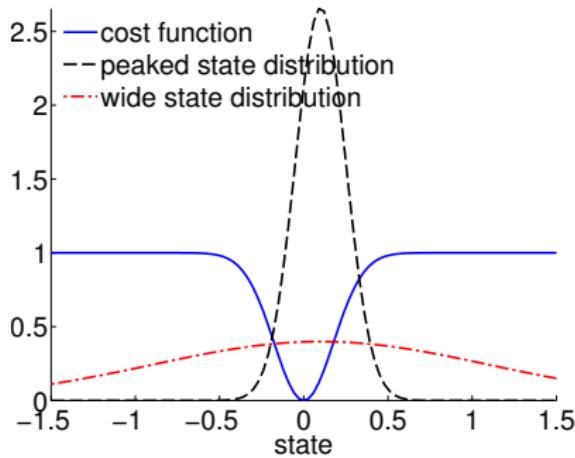
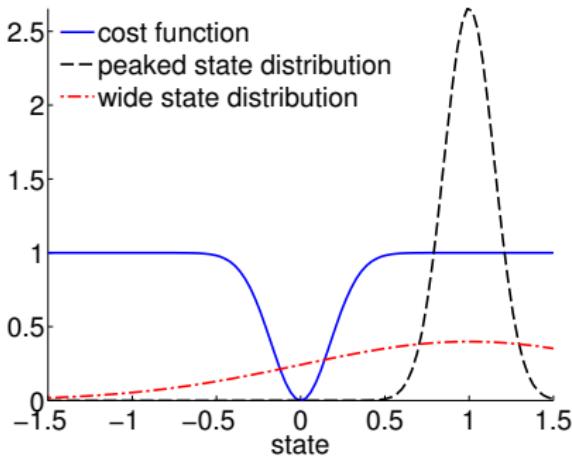
Task: Minimize $\mathbb{E}[c(\mathbf{x}_t)]$



- In the **early stages of learning**, state predictions are expected to be far away from the target

Natural Exploration with the Saturating Cost

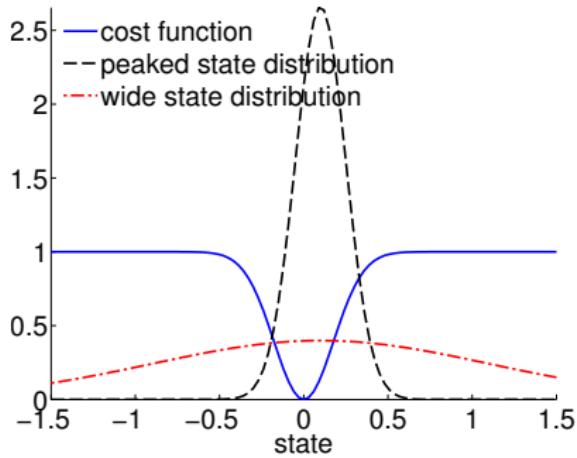
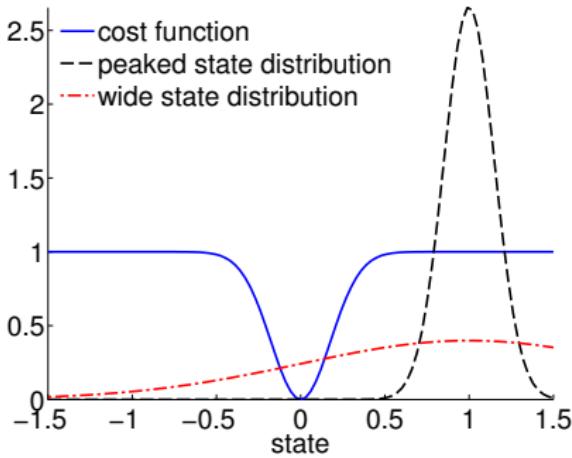
Task: Minimize $\mathbb{E}[c(x_t)]$



- In the **early stages of learning**, state predictions are expected to be far away from the target ► **Exploration** favored

Natural Exploration with the Saturating Cost

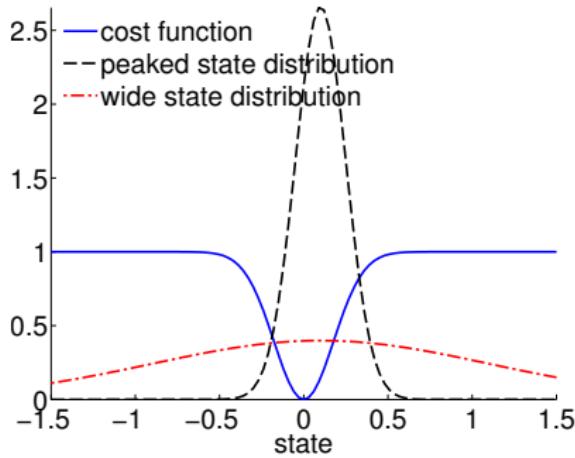
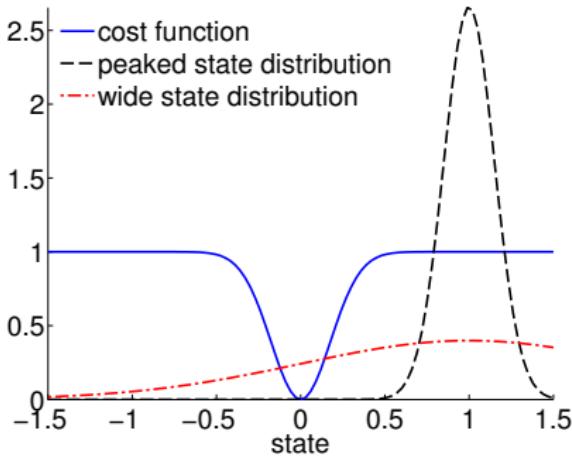
Task: Minimize $\mathbb{E}[c(x_t)]$



- In the **early stages of learning**, state predictions are expected to be far away from the target ► **Exploration** favored
- In the **final stages of learning**, state predictions are expected to be close to the target

Natural Exploration with the Saturating Cost

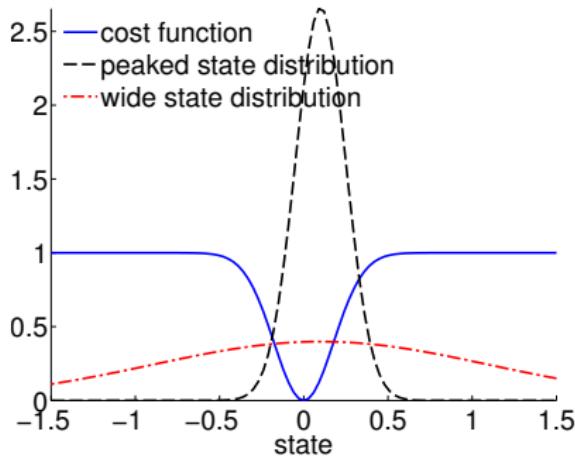
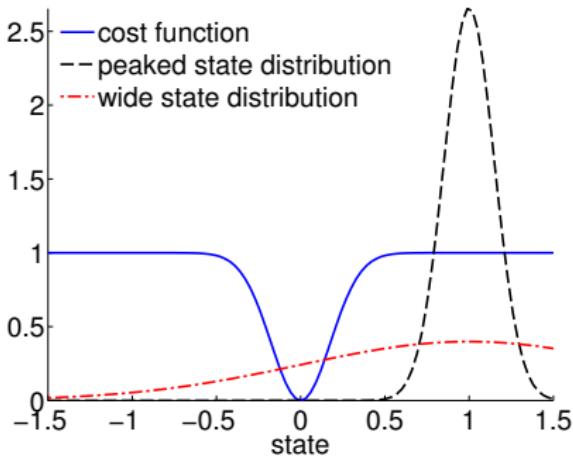
Task: Minimize $\mathbb{E}[c(x_t)]$



- In the **early stages of learning**, state predictions are expected to be far away from the target ► **Exploration** favored
- In the **final stages of learning**, state predictions are expected to be close to the target ► **Exploitation** favored

Natural Exploration with the Saturating Cost

Task: Minimize $\mathbb{E}[c(x_t)]$



- In the **early stages of learning**, state predictions are expected to be far away from the target ➤ **Exploration** favored
 - In the **final stages of learning**, state predictions are expected to be close to the target ➤ **Exploitation** favored
- Bayesian treatment: **Natural exploration/exploitation trade-off**

$$f \sim GP(0, k), \quad \text{Training data: } \mathbf{X}, \mathbf{y}$$
$$\mathbf{x}_* \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

- Compute $\mathbb{E}[f(\mathbf{x}_*)]$

$$f \sim GP(0, k), \quad \text{Training data: } \mathbf{X}, \mathbf{y}$$
$$\mathbf{x}_* \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

- Compute $\mathbb{E}[f(\mathbf{x}_*)]$

$$\mathbb{E}_{f, \mathbf{x}_*}[f(\mathbf{x}_*)] = \mathbb{E}_{\mathbf{x}} \left[\mathbb{E}_f[f(\mathbf{x}_*) | \mathbf{x}_*] \right] = \mathbb{E}_{\mathbf{x}_*} [m_f(\mathbf{x}_*)]$$

$$f \sim GP(0, k), \quad \text{Training data: } \mathbf{X}, \mathbf{y}$$
$$\mathbf{x}_* \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

- Compute $\mathbb{E}[f(\mathbf{x}_*)]$

$$\begin{aligned}\mathbb{E}_{f, \mathbf{x}_*}[f(\mathbf{x}_*)] &= \mathbb{E}_{\mathbf{x}} \left[\mathbb{E}_f[f(\mathbf{x}_*) | \mathbf{x}_*] \right] = \mathbb{E}_{\mathbf{x}_*} [m_f(\mathbf{x}_*)] \\ &= \mathbb{E}_{\mathbf{x}_*} \left[k(\mathbf{x}_*, \mathbf{X})(\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{y} \right]\end{aligned}$$

$$f \sim GP(0, k), \quad \text{Training data: } \mathbf{X}, \mathbf{y}$$

$$\mathbf{x}_* \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

- Compute $\mathbb{E}[f(\mathbf{x}_*)]$

$$\begin{aligned}\mathbb{E}_{f, \mathbf{x}_*}[f(\mathbf{x}_*)] &= \mathbb{E}_{\mathbf{x}} \left[\mathbb{E}_f[f(\mathbf{x}_*) | \mathbf{x}_*] \right] = \mathbb{E}_{\mathbf{x}_*} [m_f(\mathbf{x}_*)] \\ &= \mathbb{E}_{\mathbf{x}_*} \left[k(\mathbf{x}_*, \mathbf{X})(\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{y} \right] \\ &= \boldsymbol{\beta}^\top \int k(\mathbf{X}, \mathbf{x}_*) \mathcal{N}(\mathbf{x}_* | \boldsymbol{\mu}, \boldsymbol{\Sigma}) d\mathbf{x}_* \\ \boldsymbol{\beta} := (\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{y} \quad &\blacktriangleright \text{independent of } \mathbf{x}_*\end{aligned}$$

$$f \sim GP(0, k), \quad \text{Training data: } \mathbf{X}, \mathbf{y}$$

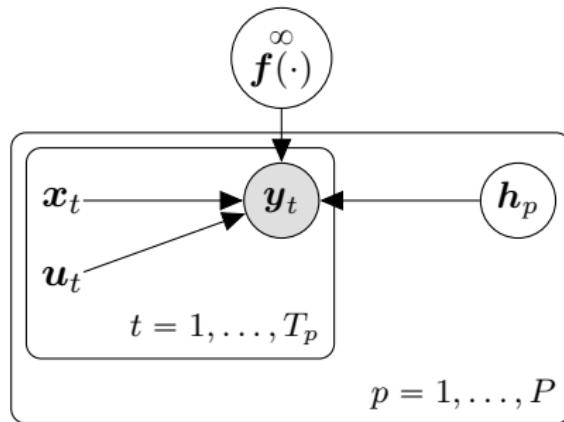
$$\mathbf{x}_* \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

- Compute $\mathbb{E}[f(\mathbf{x}_*)]$

$$\begin{aligned}\mathbb{E}_{f, \mathbf{x}_*}[f(\mathbf{x}_*)] &= \mathbb{E}_{\mathbf{x}} \left[\mathbb{E}_f[f(\mathbf{x}_*) | \mathbf{x}_*] \right] = \mathbb{E}_{\mathbf{x}_*} [m_f(\mathbf{x}_*)] \\ &= \mathbb{E}_{\mathbf{x}_*} \left[k(\mathbf{x}_*, \mathbf{X})(\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{y} \right] \\ &= \boldsymbol{\beta}^\top \int k(\mathbf{X}, \mathbf{x}_*) \mathcal{N}(\mathbf{x}_* | \boldsymbol{\mu}, \boldsymbol{\Sigma}) d\mathbf{x}_* \\ \boldsymbol{\beta} &:= (\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{y} \quad \text{► independent of } \mathbf{x}_*\end{aligned}$$

- If k is a Gaussian (squared exponential) kernel, this integral can be solved analytically
- Variance of $f(\mathbf{x}_*)$ can be computed similarly

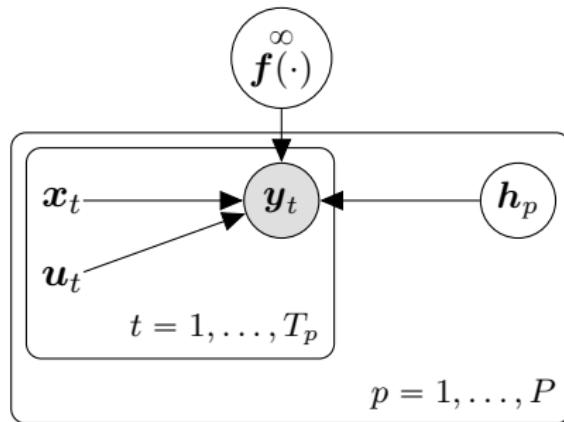
Meta Learning Model



$$\mathbf{f}(\cdot) \sim GP$$

$$p(\mathbf{H}) = \prod_p p(\mathbf{h}_p), \quad p(\mathbf{h}_p) = \mathcal{N}(\mathbf{0}, \mathbf{I})$$

Meta Learning Model

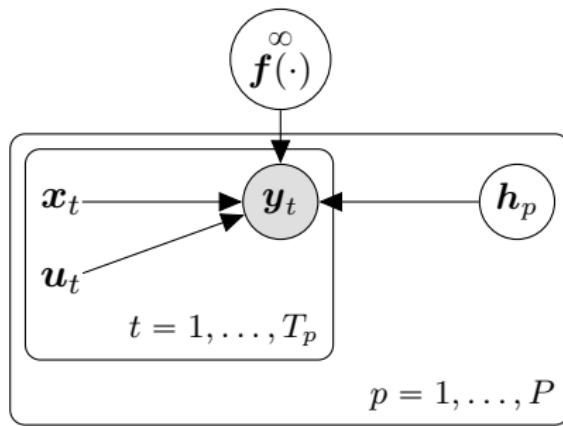


$$\mathbf{f}(\cdot) \sim GP$$

$$p(\mathbf{H}) = \prod_p p(\mathbf{h}_p), \quad p(\mathbf{h}_p) = \mathcal{N}(\mathbf{0}, \mathbf{I})$$

$$p(\mathbf{Y}, \mathbf{H}, \mathbf{f}(\cdot) | \mathbf{X}, \mathbf{U}) = \prod_{p=1}^P p(\mathbf{h}_p) \prod_{t=1}^{T_p} p(\mathbf{y}_t | \mathbf{x}_t, \mathbf{u}_t, \mathbf{h}_p, \mathbf{f}(\cdot)) p(\mathbf{f}(\cdot))$$

$$\mathbf{y}_t = \mathbf{x}_{t+1} - \mathbf{x}_t$$



Mean-field variational family:

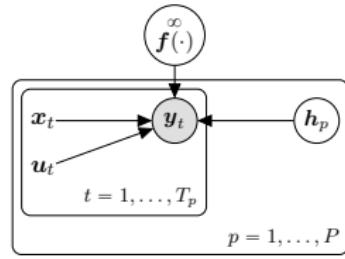
$$q(\mathbf{f}(\cdot), \mathbf{H}) = q(\mathbf{f}(\cdot))q(\mathbf{H})$$

$$q(\mathbf{H}) = \prod_{p=1}^P \mathcal{N}(\mathbf{h}_p | \mathbf{n}_p, \mathbf{T}_p),$$

$$q(\mathbf{f}(\cdot)) = \int p(\mathbf{f}(\cdot) | \mathbf{f}_Z) q(\mathbf{f}_Z) d\mathbf{f}_Z \quad \blacktriangleright \text{SV-GP (Titsias, 2009)}$$

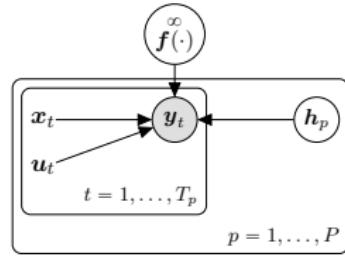
Evidence Lower Bound

$$ELBO = \mathbb{E}_{q(\mathbf{f}(\cdot), \mathbf{H})} \left[\log \frac{p(\mathbf{Y}, \mathbf{H}, \mathbf{f}(\cdot) | \mathbf{X}, \mathbf{U})}{q(\mathbf{f}(\cdot), \mathbf{H})} \right]$$



Evidence Lower Bound

$$\begin{aligned} ELBO &= \mathbb{E}_{q(\mathbf{f}(\cdot), \mathbf{H})} \left[\log \frac{p(\mathbf{Y}, \mathbf{H}, \mathbf{f}(\cdot) | \mathbf{X}, \mathbf{U})}{q(\mathbf{f}(\cdot), \mathbf{H})} \right] \\ &= \sum_{p=1}^P \sum_{t=1}^{T_p} \mathbb{E}_{q(\mathbf{f}_t | \mathbf{x}_t, \mathbf{u}_t, \mathbf{h}_p) q(\mathbf{h}_p)} \left[\log p(\mathbf{y}_t | \mathbf{f}_t) \right] \\ &\quad - \text{KL}(q(\mathbf{H}) || p(\mathbf{H})) - \text{KL}(q(\mathbf{f}(\cdot)) || p(\mathbf{f}(\cdot))) \end{aligned}$$



Evidence Lower Bound

$$\begin{aligned}
 ELBO &= \mathbb{E}_{q(\mathbf{f}(\cdot), \mathbf{H})} \left[\log \frac{p(\mathbf{Y}, \mathbf{H}, \mathbf{f}(\cdot) | \mathbf{X}, \mathbf{U})}{q(\mathbf{f}(\cdot), \mathbf{H})} \right] \\
 &= \sum_{p=1}^P \sum_{t=1}^{T_p} \mathbb{E}_{q(\mathbf{f}_t | \mathbf{x}_t, \mathbf{u}_t, \mathbf{h}_p) q(\mathbf{h}_p)} \left[\log p(\mathbf{y}_t | \mathbf{f}_t) \right] \\
 &\quad - \text{KL}(q(\mathbf{H}) || p(\mathbf{H})) - \text{KL}(q(\mathbf{f}(\cdot)) || p(\mathbf{f}(\cdot))) \\
 &\quad \underbrace{\qquad\qquad\qquad}_{\text{Monte Carlo estimate}} \\
 &= \sum_{p=1}^P \sum_{t=1}^{T_p} \overbrace{\mathbb{E}_{q(\mathbf{f}_t | \mathbf{x}_t, \mathbf{u}_t, \mathbf{h}_p) q(\mathbf{h}_p)} \left[\log p(\mathbf{y}_t | \mathbf{f}_t) \right]}^{\text{closed-form solution}} \\
 &\quad \underbrace{- \text{KL}(q(\mathbf{H}) || p(\mathbf{H})) - \text{KL}(q(\mathbf{F}_Z) || p(\mathbf{F}_Z))}_{\text{closed-form solution}}
 \end{aligned}$$

