

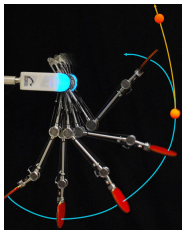
Useful Models for Robot Learning

Marc Deisenroth

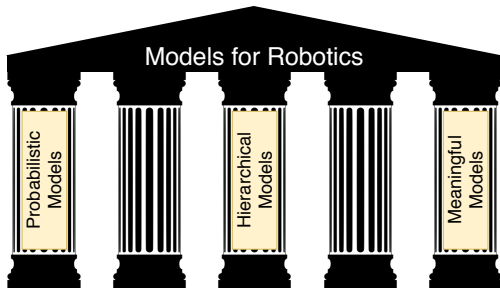
Centre for Artificial Intelligence
Department of Computer Science
University College London

ELLIS Workshop

March 5, 2020



- Automatic adaption in robotics ► **Learning**
- Practical constraint: **data efficiency**
- Models are useful for data-efficient learning in robotics



1 Probabilistic models

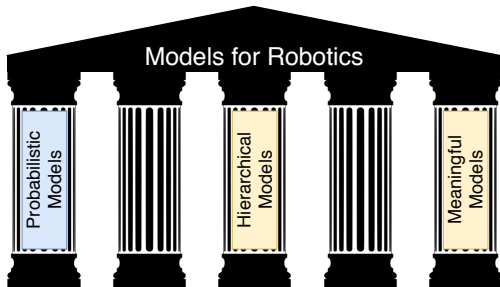
- ▶ Fast reinforcement learning

2 Hierarchical models

- ▶ Infer task similarities within a meta-learning framework

3 Physically meaningful models

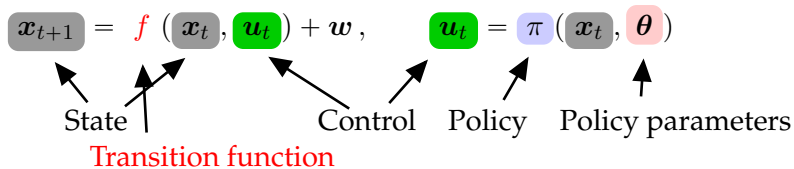
- ▶ Encode real-world constraints into learning

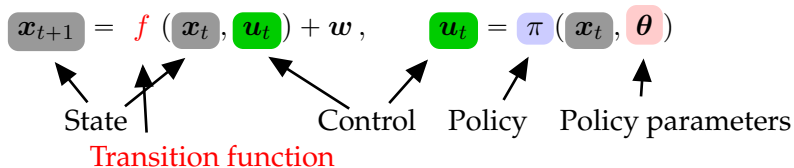


Carl Rasmussen



Dieter Fox





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 $J(\theta) = \sum_t \mathbb{E}[c(\mathbf{x}_t)|\theta]$

Objective

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

PILCO Framework: High-Level Steps

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

Objective

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

PILCO Framework: High-Level Steps

- 1 Probabilistic model for transition function
 - ▶ System identification
- 2 Compute long-term state evolution $p(\mathbf{x}_1|\theta), \dots, p(\mathbf{x}_T|\theta)$

Objective

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

PILCO Framework: High-Level Steps

- 1 Probabilistic model for transition function
 - ▶▶ System identification
- 2 Compute long-term state evolution $p(\mathbf{x}_1|\theta), \dots, p(\mathbf{x}_T|\theta)$
- 3 Policy improvement

Objective

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

PILCO Framework: High-Level Steps

- 1 Probabilistic model for transition function
 - ▶▶ System identification
- 2 Compute long-term state evolution $p(\mathbf{x}_1|\theta), \dots, p(\mathbf{x}_T|\theta)$
- 3 Policy improvement
- 4 Apply controller

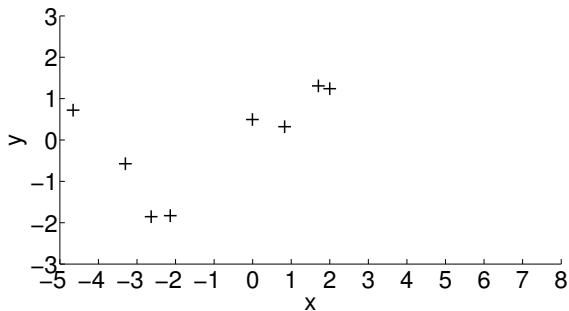
Objective

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

PILCO Framework: High-Level Steps

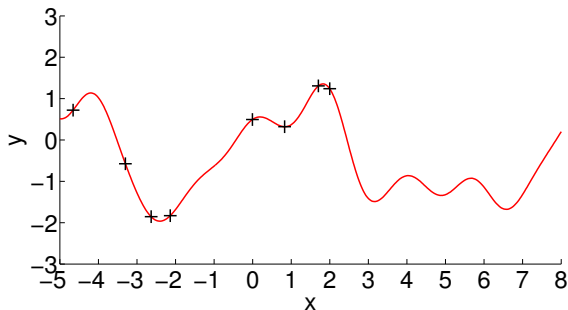
- 1 Probabilistic model for transition function f**
▶ System identification
- 2 Compute long-term predictions $p(\mathbf{x}_1|\theta), \dots, p(\mathbf{x}_T|\theta)$**
- 3 Policy improvement**
- 4 Apply controller**

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



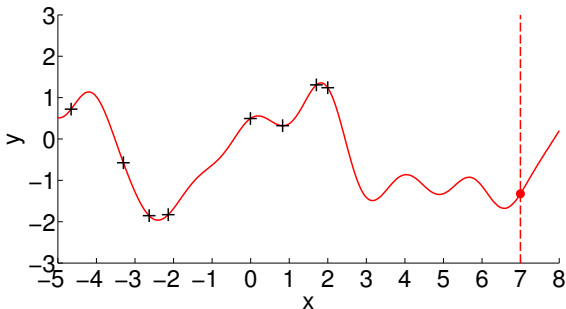
Observed function values

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



Plausible model

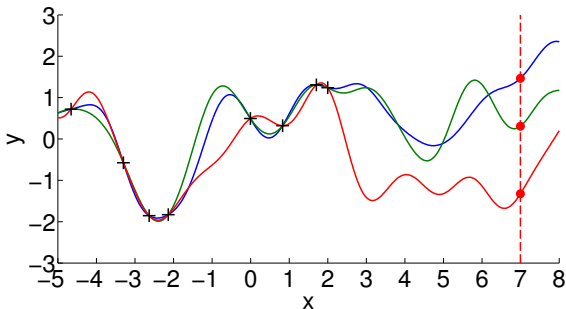
Model learning problem: Find a function $f : x \mapsto f(x) = y$



Plausible model

Predictions? Decision Making?

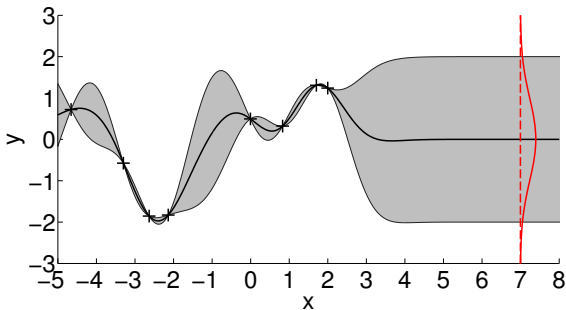
Model learning problem: Find a function $f : x \mapsto f(x) = y$



More plausible models

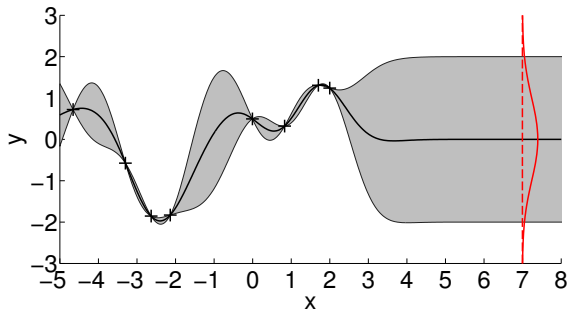
Predictions? Decision Making? Model Errors!

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

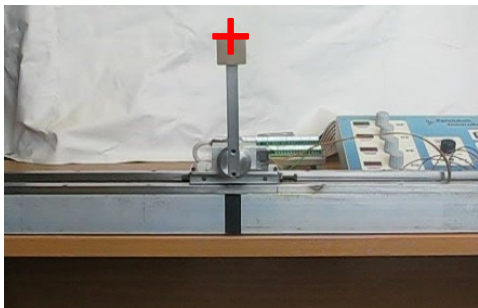
- ▶ Express **uncertainty** about the underlying function to be **robust to model errors**
- ▶ **Gaussian process** for model learning (Rasmussen & Williams, 2006)

Objective

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

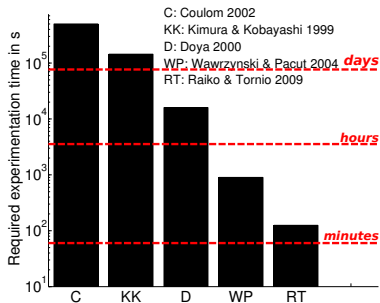
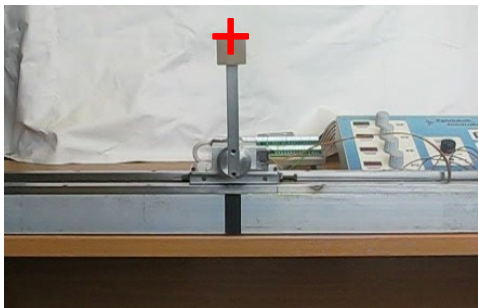
PILCO Framework: High-Level Steps

- 1 Probabilistic model for transition function f
 - ▶▶ System identification
- 2 Compute long-term predictions $p(\mathbf{x}_1|\theta), \dots, p(\mathbf{x}_T|\theta)$
- 3 Policy optimization via gradient descent
- 4 Apply controller

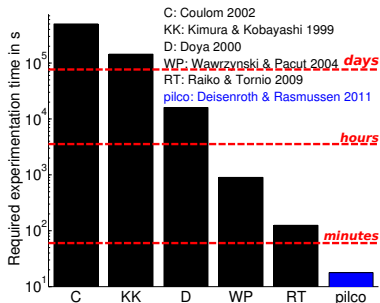
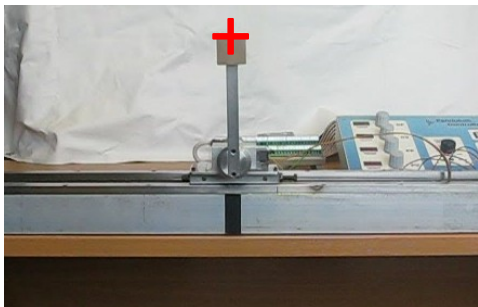


- Swing up and balance a freely swinging pendulum on a cart
- No knowledge about nonlinear dynamics ►► Learn from scratch
- Cost function $c(\mathbf{x}) = 1 - \exp(-\|\mathbf{x} - \mathbf{x}_{\text{target}}\|^2)$

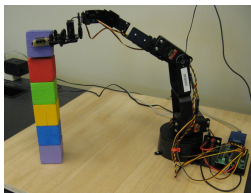
- Code: <https://github.com/ICL-SML/pilco-matlab>



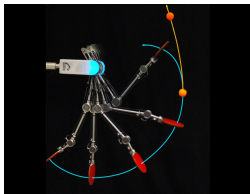
- Swing up and balance a freely swinging pendulum on a cart
- No knowledge about nonlinear dynamics ► Learn from scratch
- Cost function $c(\mathbf{x}) = 1 - \exp(-\|\mathbf{x} - \mathbf{x}_{\text{target}}\|^2)$
- Code: <https://github.com/ICL-SML/pilco-matlab>



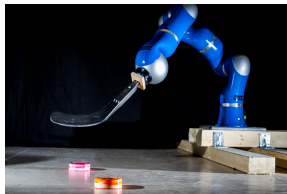
- Swing up and balance a freely swinging pendulum on a cart
- No knowledge about nonlinear dynamics ►► Learn from scratch
- Cost function $c(\mathbf{x}) = 1 - \exp(-\|\mathbf{x} - \mathbf{x}_{\text{target}}\|^2)$
- **Unprecedented learning speed** compared to state-of-the-art
- Code: <https://github.com/ICL-SML/pilco-matlab>



with D Fox



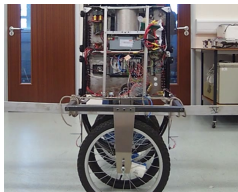
with P Englert, A Paraschos, J Peters



with A Kupcsik, J Peters, G Neumann



B Bischoff (Bosch), ESANN 2013



A McHutchon (U Cambridge)



B Bischoff (Bosch), ECML 2013

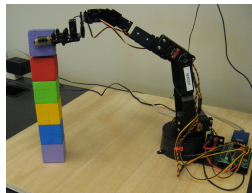
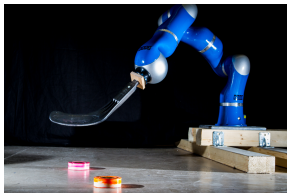
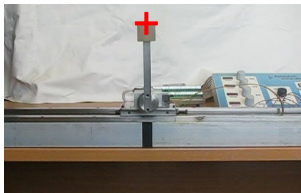
►► Application to a wide range of robotic systems

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

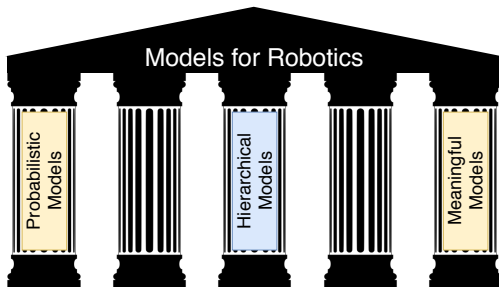
Englert et al. (ICRA, 2013): *Model-based Imitation Learning by Probabilistic Trajectory Matching*

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*



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



Steindór Sæmundsson



Katja Hofmann



Meta Learning (Schmidhuber 1987)

Generalize knowledge from known tasks to new (related) tasks



Meta Learning (Schmidhuber 1987)

Generalize knowledge from known tasks to new (related) tasks

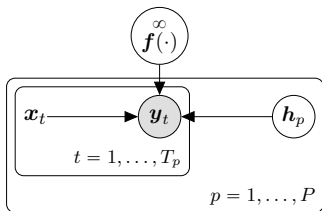
- Different robot configurations (link lengths, weights, ...)
- Re-use experience gathered so far generalize learning to new dynamics that are similar
 - ▶ Accelerated learning



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

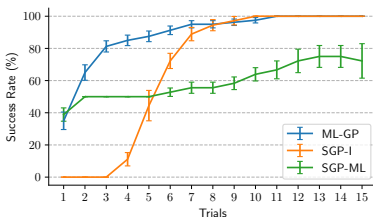


- **Separate** global and task-specific properties
- Shared global parameters describe general dynamics
- Describe task-specific (local) properties with latent variable
- Online variational inference of local properties



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

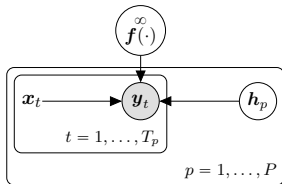
- GP captures **global properties** of the dynamics
- Latent variable \mathbf{h}_p encodes **local properties**
 - ▶▶ Variational inference to find a posterior on latent task



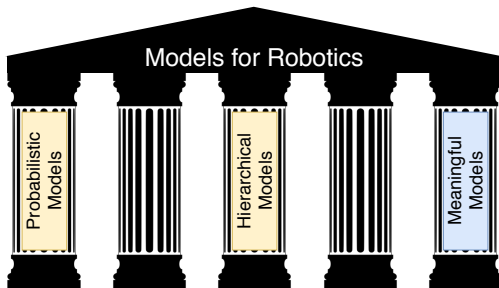
- Train on 6 tasks with different configurations (length/mass)
- Few-shot generalization on 4 unseen configurations
- Success: solve all 10 (6 training + 4 test) tasks
- **Meta learning: blue**
- **Independent (GP-MPC): orange**
- **Aggregated experience model (no latents): green**

▶▶ **Meta RL generalizes well to unseen tasks**

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



- 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



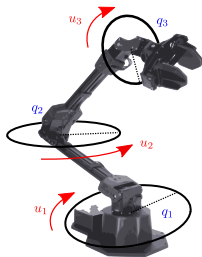
Steindór Sæmundsson



Alexander Terenin



Katja Hofmann



Equations of motion

$$u = \frac{d}{dt} \left(\frac{\partial L}{\partial \dot{q}} \right) - \frac{\partial L}{\partial q}$$

- Goal: **Data efficiency** and **interpretability**
- **Inductive biases** to account for physical/mechanical properties (e.g., conservation laws, configuration constraints)
 - ▶ Learn dynamical systems that are “meaningful”

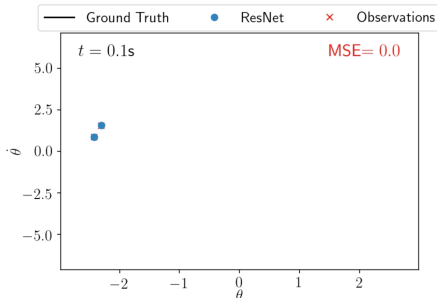
Approach:

- **Euler discretization** of continuous-time dynamical system

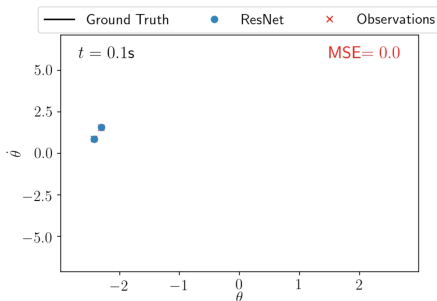
$$\mathbf{x}(T)|\mathbf{x}_0 = \int_{t=0}^T f_{\theta}(\mathbf{x}(t))dt \approx \mathbf{x}_0 + h \sum_{t=0}^{T-1} f_{\theta}(\mathbf{x}_t, t)$$

- **Deep residual network**
(E, 2017; Haber & Ruthotto, 2017; Chen et al., 2018)

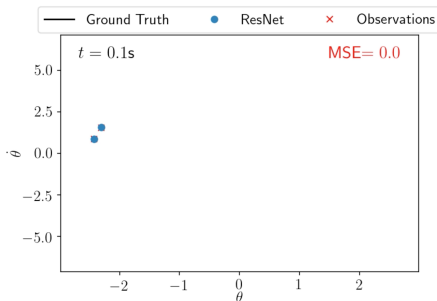
- ODE: $\ddot{\theta} = -\frac{g}{l} \sin \theta$
- Observation: $\mathbf{y} = [\theta, \dot{\theta}]^\top$
- Training data: 15 seconds (150 data points)



- ODE: $\ddot{\theta} = -\frac{g}{l} \sin \theta$
- Observation: $\mathbf{y} = [\theta, \dot{\theta}]^\top + \boldsymbol{\epsilon}$, $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, 0.33^2 \mathbf{I})$
- Training data: 15 seconds



- ODE: $\ddot{\theta} = -\frac{g}{l} \sin \theta$
- Observation: $\mathbf{y} = [\theta, \dot{\theta}]^\top + \boldsymbol{\epsilon}$, $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, 0.33^2 \mathbf{I})$
- Training data: 15 seconds



- Low prediction quality
- Does not obey physics
- ResNet does not conserve energy

- **Lagrangian:** Encodes “type” of physics, symmetries.

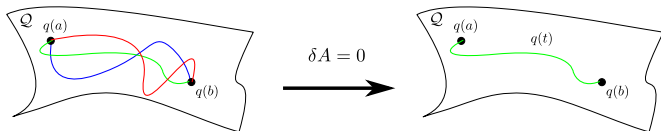
$$L(\mathbf{q}(t), \dot{\mathbf{q}}(t))$$

- **Lagrangian:** Encodes “type” of physics, symmetries.

$$L(\mathbf{q}(t), \dot{\mathbf{q}}(t))$$

- **Hamilton's Principle:**

$$A = \int_a^b L(\mathbf{q}(t), \dot{\mathbf{q}}(t)) dt, \quad \frac{\delta A}{\delta \mathbf{q}(t)} = \mathbf{0}$$

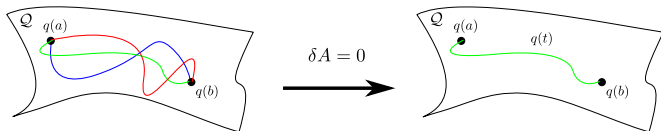


- **Lagrangian:** Encodes “type” of physics, symmetries.

$$L(\mathbf{q}(t), \dot{\mathbf{q}}(t))$$

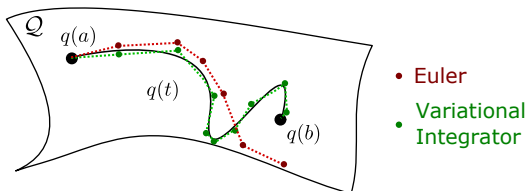
- **Hamilton's Principle:**

$$A = \int_a^b L(\mathbf{q}(t), \dot{\mathbf{q}}(t)) dt, \quad \frac{\delta A}{\delta \mathbf{q}(t)} = \mathbf{0}$$

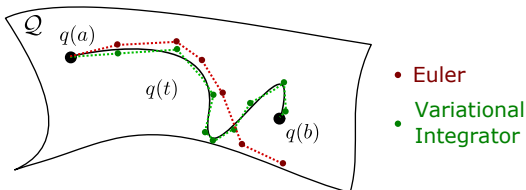


First idea:

- Learn Lagrangian L instead of dynamics
- Encode physical properties via L (e.g., Lutter et al., 2019; Greydanus et al., 2019)



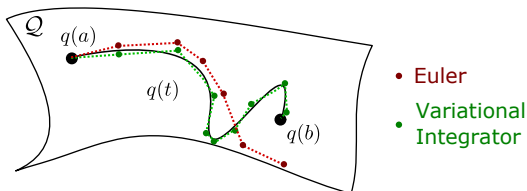
Second idea: Discretize in a way that preserves the physics



Second idea: Discretize in a way that preserves the physics

■ Conservative, separable Newtonian system:

$$L_{\theta}(\mathbf{q}, \dot{\mathbf{q}}) = T_{\theta}(\dot{\mathbf{q}}) - U_{\theta}(\mathbf{q}) = \underbrace{\frac{1}{2} \dot{\mathbf{q}}^{\top} \mathbf{M}_{\theta} \dot{\mathbf{q}}}_{\text{kinetic}} - \underbrace{U_{\theta}(\mathbf{q})}_{\text{potential}}$$

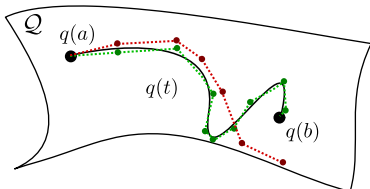


Second idea: Discretize in a way that preserves the physics

- Conservative, separable Newtonian system:

$$L_{\theta}(\mathbf{q}, \dot{\mathbf{q}}) = \underbrace{T_{\theta}(\dot{\mathbf{q}})}_{\text{kinetic}} - \underbrace{U_{\theta}(\mathbf{q})}_{\text{potential}}$$

- Discretize action integral A



- Euler
- Variational Integrator

Second idea: Discretize in a way that preserves the physics

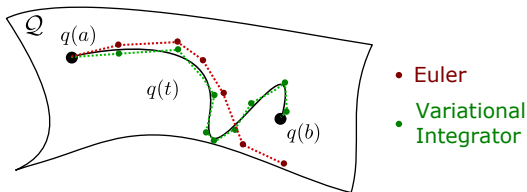
- Conservative, separable Newtonian system:

$$L_\theta(\mathbf{q}, \dot{\mathbf{q}}) = T_\theta(\dot{\mathbf{q}}) - U_\theta(\mathbf{q}) = \underbrace{\frac{1}{2} \dot{\mathbf{q}}^\top \mathbf{M}_\theta \dot{\mathbf{q}}}_{\text{kinetic}} - \underbrace{U_\theta(\mathbf{q})}_{\text{potential}}$$

- Discretize action integral A
- Explicit variational integrator

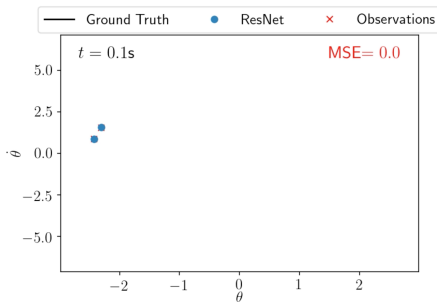
$$\mathbf{x}_{t+1} = f_\theta(\mathbf{x}_1, t, h), \quad \mathbf{x}_t := [\mathbf{q}_t, \mathbf{q}_{t-1}]$$

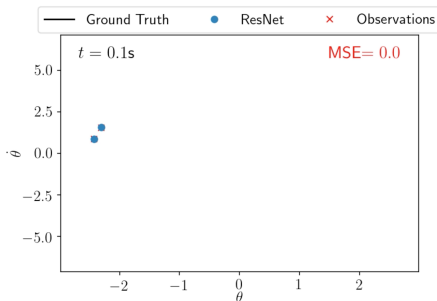
with initial condition \mathbf{x}_1



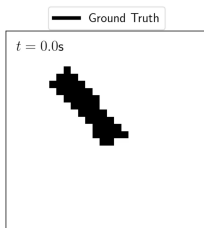
- Physical properties (e.g., conservation laws) automatically enforced
- Flexibility retained to model U_θ (e.g., with a neural network)
- Notions of kinetic and potential energy
 - ▶▶ Increased interpretability

Example: Pendulum with Noisy Observations



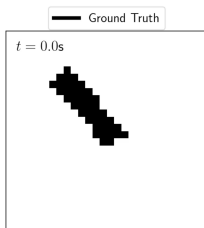


- Good predictive performance
- Obeys physics
- Conserves energy



Setting:

- Observations: 28×28 pixel images
- Training data: 60 images (6 seconds of pendulum movement)

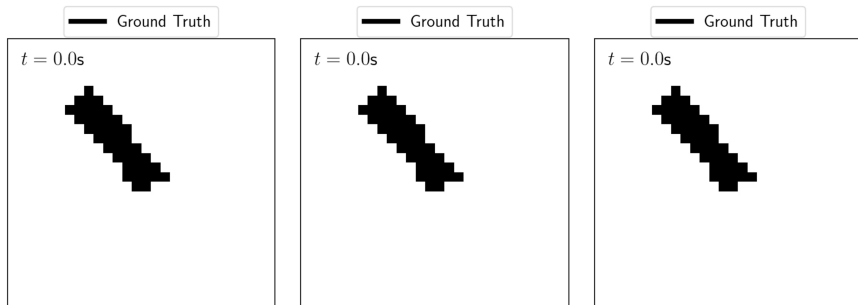


Setting:

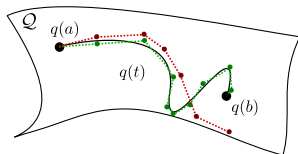
- Observations: 28×28 pixel images
- Training data: 60 images (6 seconds of pendulum movement)

Approach:

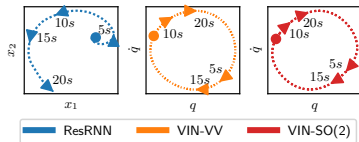
- Variational auto-encoder to embed pixels in low-dimensional space
- VIN within low-dimensional space



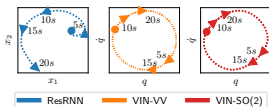
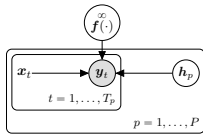
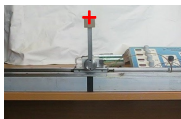
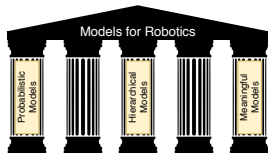
- Residual RNN
- VIN
- VIN on SO(2)
- Code: <https://tinyurl.com/yx3yhhvo>



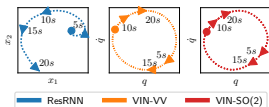
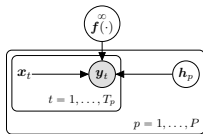
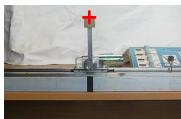
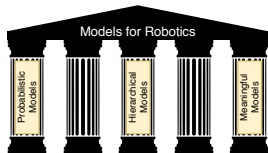
- Euler
- Variational Integrator



- Encode physics constraints when learning predictive models
- Variational integrator instead of Euler discretization
- Can be combined with VAE to learn predictive models from image observations
- Data efficient and interpretable



- **Data efficiency** is a practical challenge for autonomous robots
- Three useful models for data-efficient learning in robotics
 - 1 **Probabilistic models** for fast reinforcement learning
 - 2 **Hierarchical models** for learning task similarities within a meta-learning framework
 - 3 **Physically meaningful models** to encode real-world constraints into learning

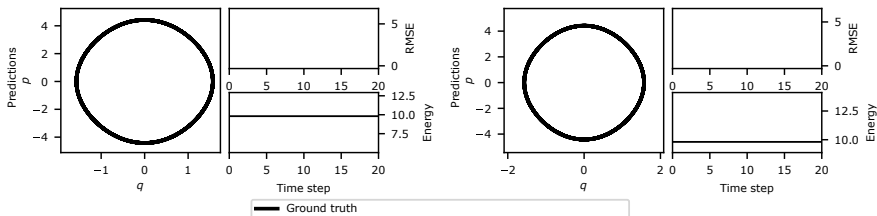


- **Data efficiency** is a practical challenge for autonomous robots
- Three useful models for data-efficient learning in robotics
 - 1 **Probabilistic models** for fast reinforcement learning
 - 2 **Hierarchical models** for learning task similarities within a meta-learning framework
 - 3 **Physically meaningful models** to encode real-world constraints into learning

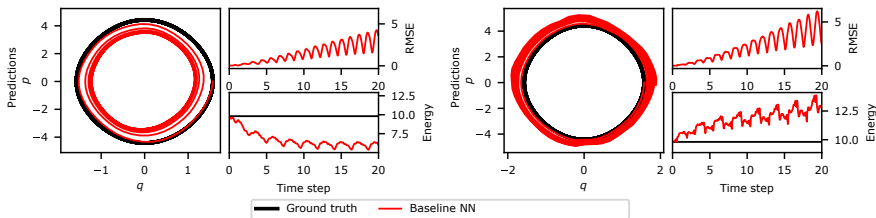
Thank you for your attention

- [1] D. P. Bertsekas. *Dynamic Programming and Optimal Control*, volume 1 of *Optimization and Computation Series*. Athena Scientific, Belmont, MA, USA, 3rd edition, 2005.
- [2] D. P. Bertsekas. *Dynamic Programming and Optimal Control*, volume 2 of *Optimization and Computation Series*. Athena Scientific, Belmont, MA, USA, 3rd edition, 2007.
- [3] 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.
- [4] Y. Chen, A. Huang, Z. Wang, I. Antonoglou, J. Schrittwieser, D. Silver, and N. de Freitas. Bayesian Optimization in AlphaGo. *arXiv:1812.06855*, 2018.
- [5] M. P. Deisenroth, P. Englert, J. Peters, and D. Fox. Multi-Task Policy Search for Robotics. In *Proceedings of the 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*, 2011.
- [9] W. E. A Proposal on Machine Learning via Dynamical Systems. *Communications in Mathematics and Statistics*, 5(1):1–11, 3 2017.
- [10] 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.
- [11] P. Englert, A. Paraschos, J. Peters, and M. P. Deisenroth. Probabilistic Model-based Imitation Learning. *Adaptive Behavior*, 21:388–403, 2013.
- [12] S. Greydanus, M. Dzamba, and J. Yosinski. Hamiltonian Neural Networks. In *Advances in Neural Information Processing Systems*, 2019.

- [13] E. Haber and L. Ruthotto. Stable Architectures for Deep Neural Networks. *Inverse Problems*, 34(1):014004, 2017.
- [14] 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.
- [15] M. Lutter, C. Ritter, and J. Peters. Deep Lagrangian Networks: Using Physics as Model Prior for Deep Learning. In *Proceedings of the International Conference on Learning Representations*, 2019.
- [16] C. E. Rasmussen and C. K. I. Williams. *Gaussian Processes for Machine Learning*. The MIT Press, 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] J. Schmidhuber. Evolutionary Principles in Self-Referential Learning. Master's thesis, 1987.
- [19] S. Sæmundsson, A. Terenin, K. Hofmann, and M. P. Deisenroth. Variational Integrator Networks for Physically Structured Embeddings. In *arXiv:1910.09349*, 2019.
- [20] S. Sæmundsson, A. Terenin, K. Hofmann, and M. P. Deisenroth. Variational Integrator Networks for Physically Structured Embeddings. In *Proceedings of the International Conference on Artificial Intelligence and Statistics*, 2020.

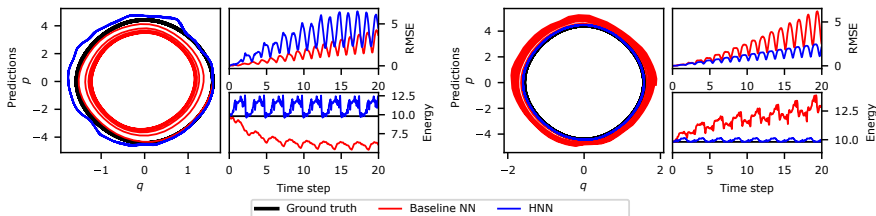


Pendulum System. **Left:** 150 observations; **Right:** 750 observations.



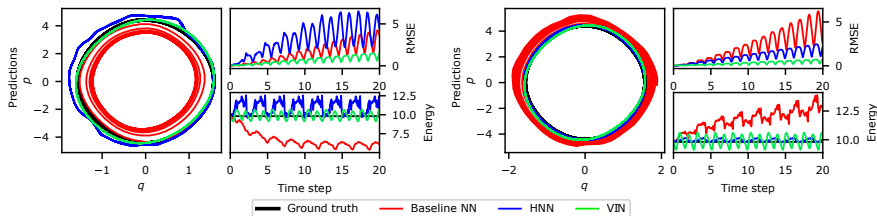
Pendulum System. **Left:** 150 observations; **Right:** 750 observations.

- **Baseline neural network:** Dissipates/adds energy for low and moderate data



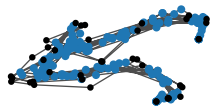
Pendulum System. **Left:** 150 observations; **Right:** 750 observations.

- **Baseline neural network:** Dissipates/adds energy for low and moderate data
- **Hamiltonian neural network** (Greydanus et al., 2019): Overfits in low-data regime

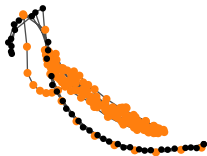


Pendulum System. **Left:** 150 observations; **Right:** 750 observations.

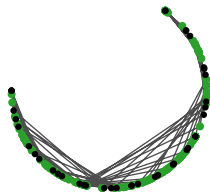
- **Baseline neural network:** Dissipates/adds energy for low and moderate data
- **Hamiltonian neural network** (Greydanus et al., 2019): Overfits in low-data regime
- **Variational integrator network:** Conserves energy and generalizes better in both regimes



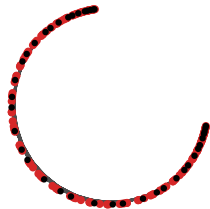
(a) VAE



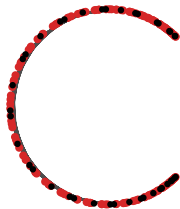
(b) Dynamic VAE



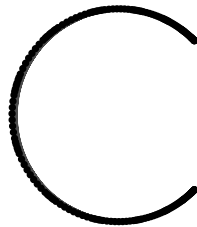
(c) Lie Group VAE



(d) VIN- $SO(2)$



(e) VIN- $SO(2)$ with fixed M



(f) Ground Truth