



University of  
BRISTOL



SIEMENS  
*Ingenuity for life*

# Interpretable Dynamics Models

For Data-Efficient Reinforcement Learning

---

Markus Kaiser, Clemens Otte, Thomas A. Runkler, Carl Henrik Ek

[markus.kaiser@siemens.com](mailto:markus.kaiser@siemens.com)

April 24, 2019

Siemens AG, Technical University of Munich, University of Bristol

# Wet-Chicken Benchmark<sup>1</sup>



---

<sup>1</sup>Tresp 1994; Hans and Udluft 2009.

# Wet-Chicken Benchmark<sup>1</sup>

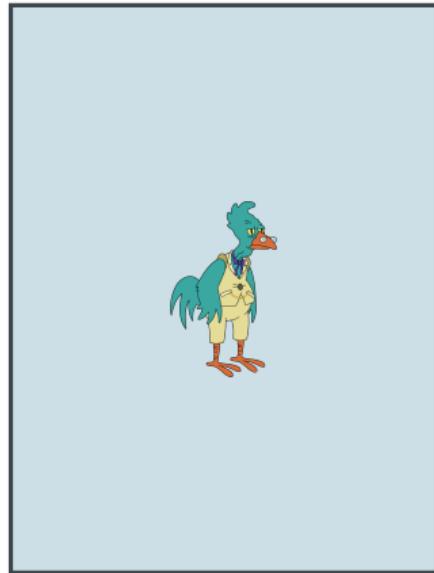


---

<sup>1</sup>Tresp 1994; Hans and Udluft 2009.

# Wet-Chicken Benchmark<sup>1</sup>

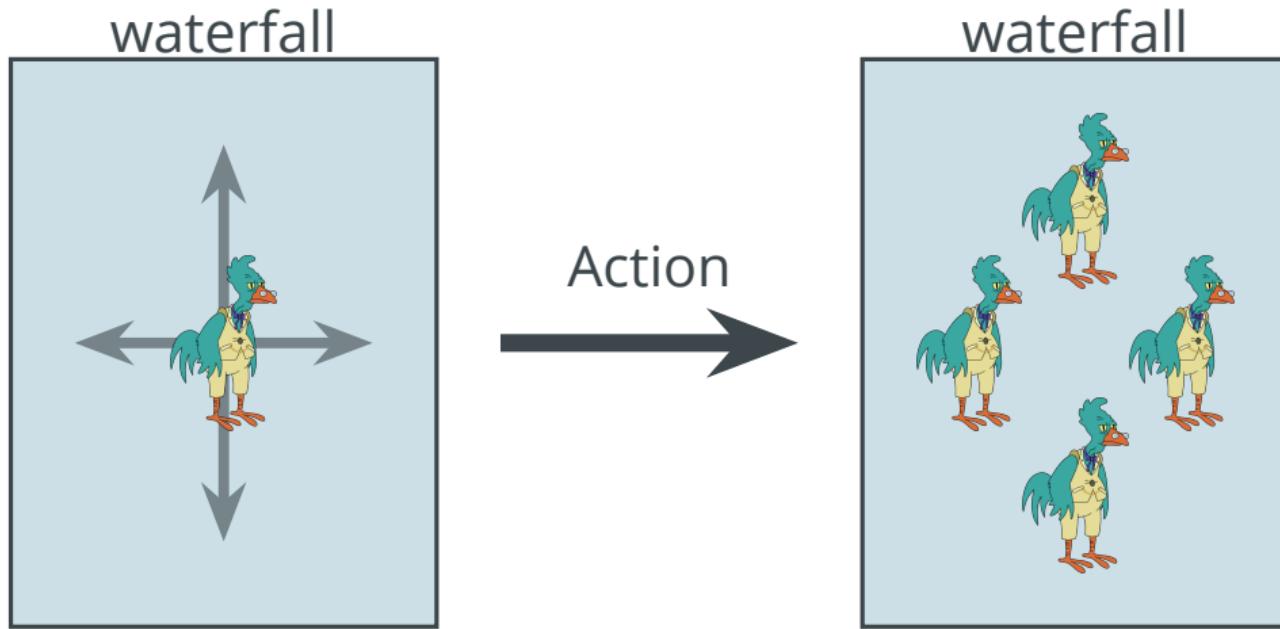
waterfall



---

<sup>1</sup>Tresp 1994; Hans and Udluft 2009.

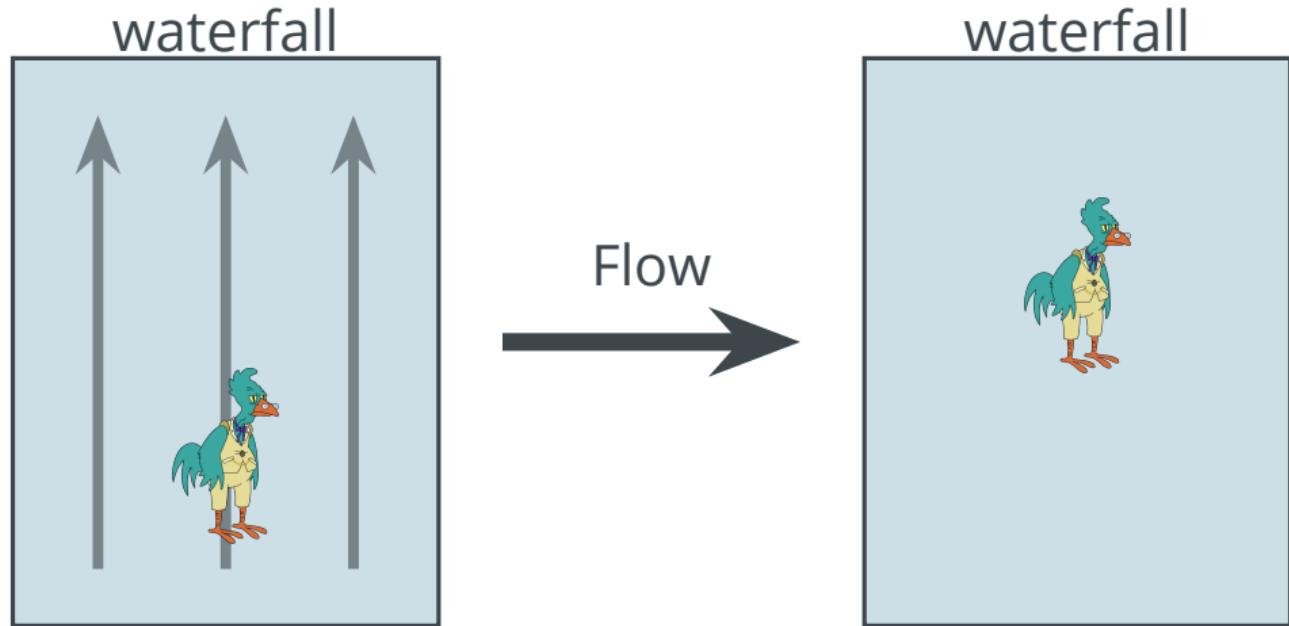
## Wet-Chicken Benchmark<sup>1</sup>



---

<sup>1</sup>Tresp 1994; Hans and Udluft 2009.

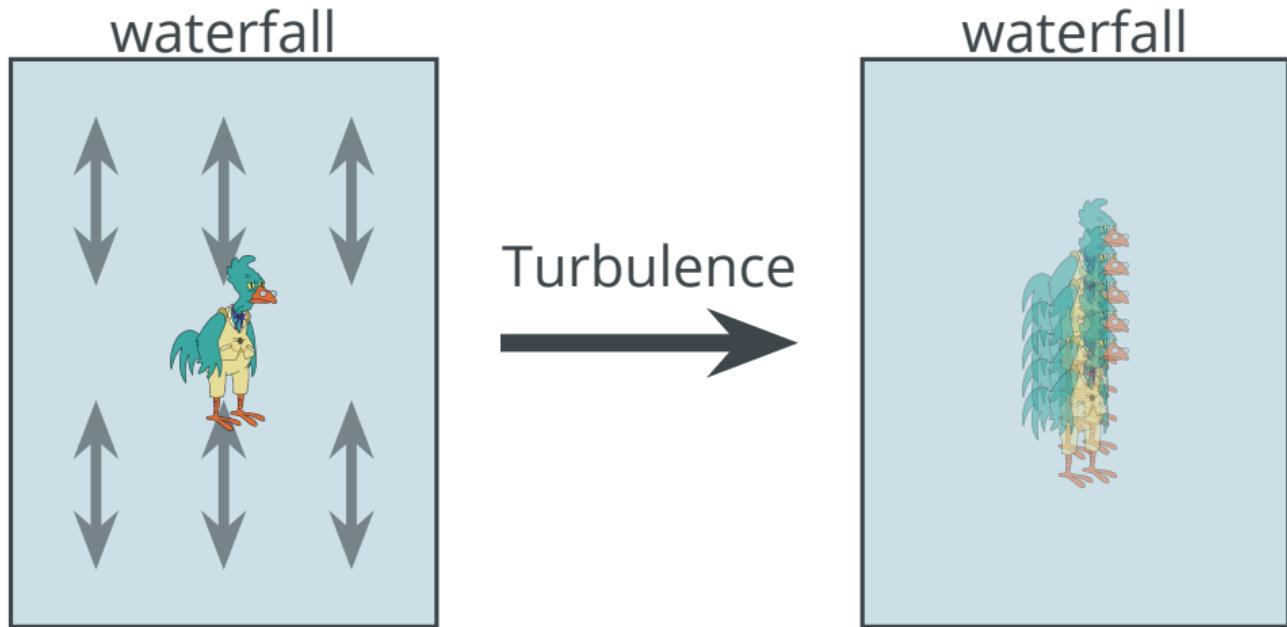
## Wet-Chicken Benchmark<sup>1</sup>



---

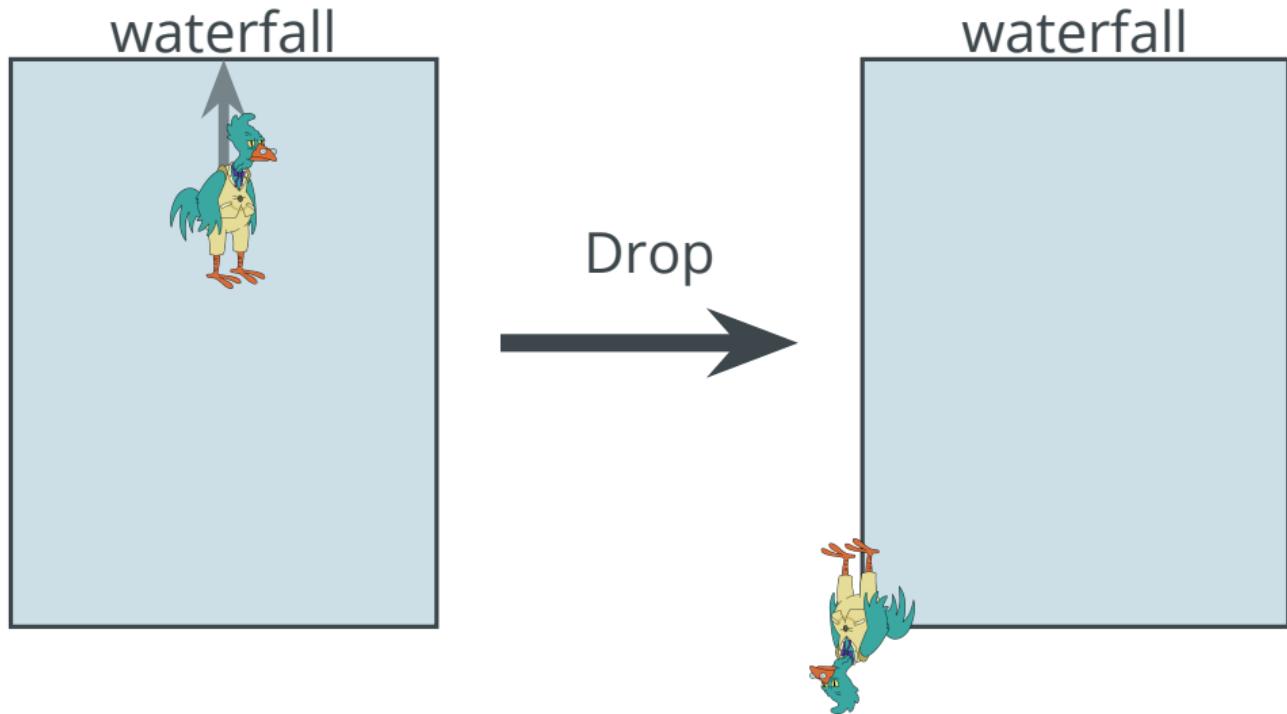
<sup>1</sup>Tresp 1994; Hans and Udluft 2009.

## Wet-Chicken Benchmark<sup>1</sup>



<sup>1</sup>Tresp 1994; Hans and Udluft 2009.

## Wet-Chicken Benchmark<sup>1</sup>

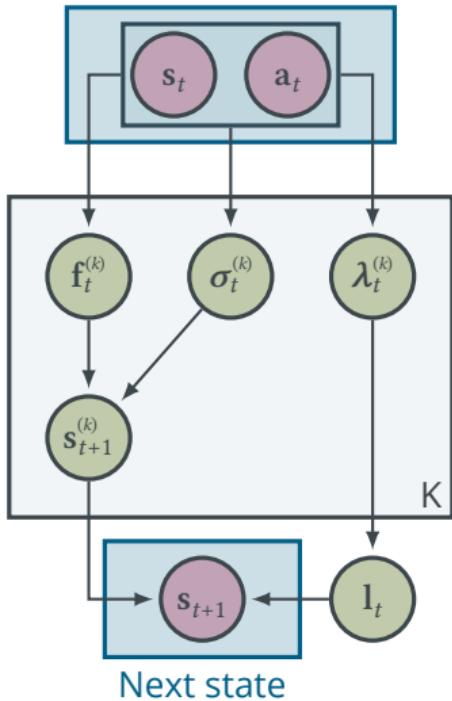


---

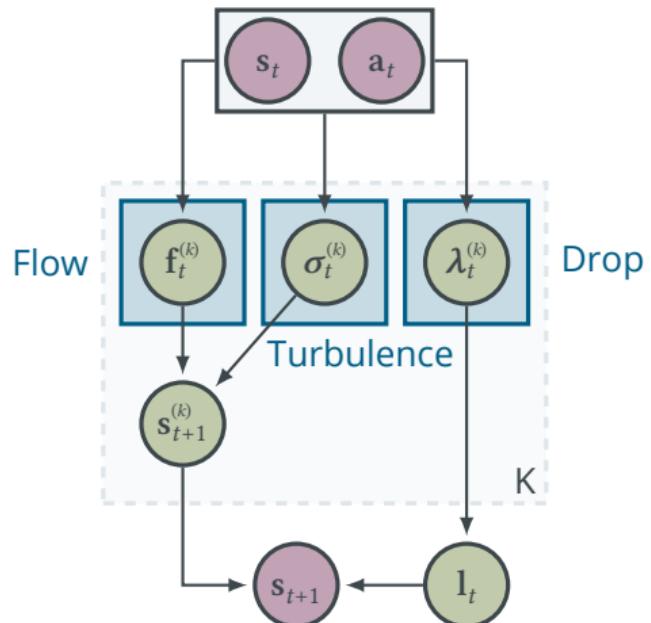
<sup>1</sup>Tresp 1994; Hans and Udluft 2009.

## Dynamics: Graphical Model

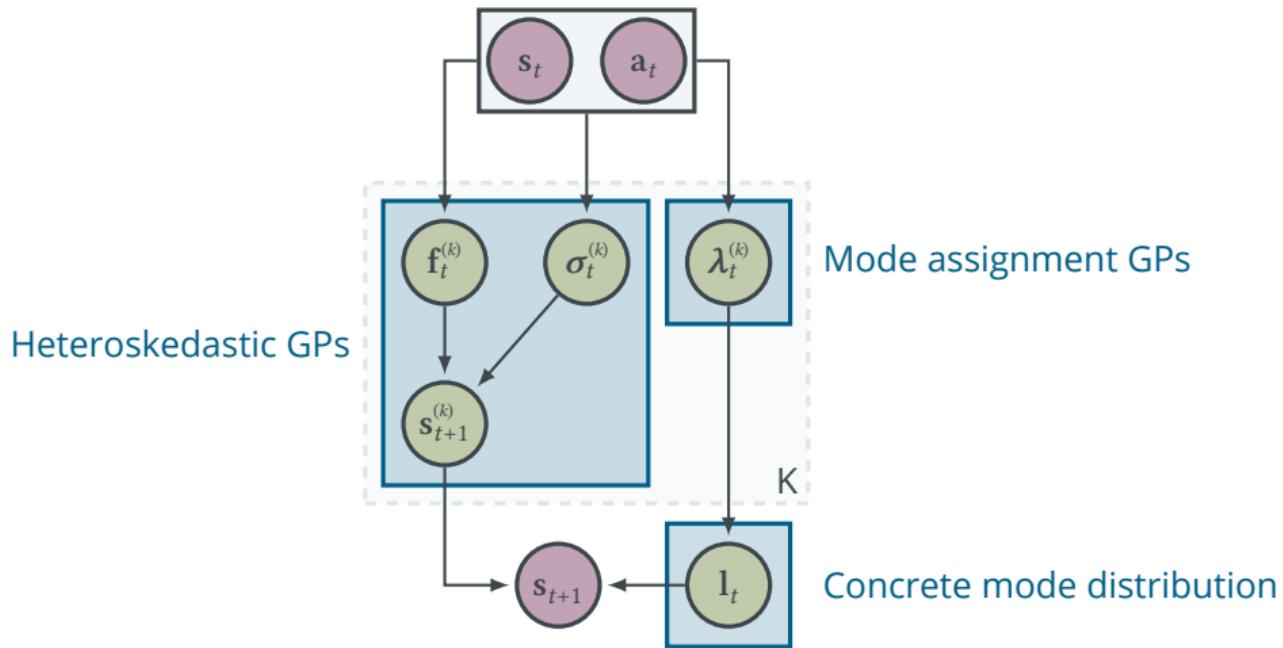
Current state and action



## Dynamics: Graphical Model



## Dynamics: Graphical Model



## Dynamics: Posterior

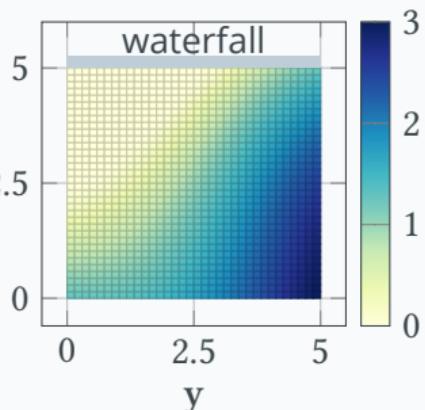
$$\begin{aligned} p(\Delta \mathbf{x}_{t+1}) &= p(\Delta \mathbf{x}_{t+1} | \text{drop}) \cdot p(\text{drop}) \\ &\quad + p(\Delta \mathbf{x}_{t+1} | \text{no drop}) \cdot p(\text{no drop}) \end{aligned}$$

## Dynamics: Posterior

$$\begin{aligned} p(\Delta \mathbf{x}_{t+1}) &= p(\Delta \mathbf{x}_{t+1} | \text{drop}) \cdot p(\text{drop}) \\ &\quad + p(\Delta \mathbf{x}_{t+1} | \text{no drop}) \cdot p(\text{no drop}) \end{aligned}$$

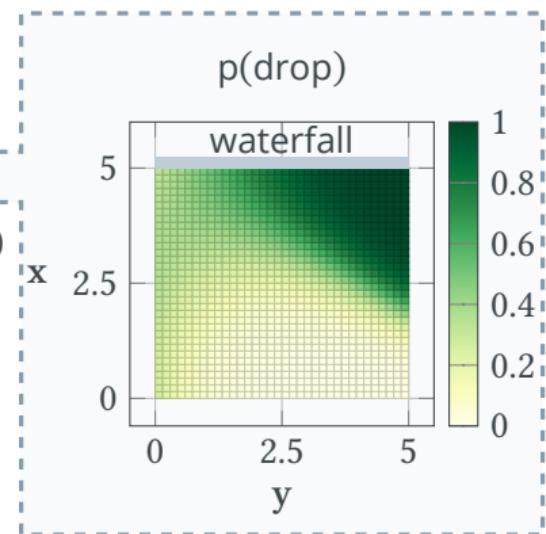
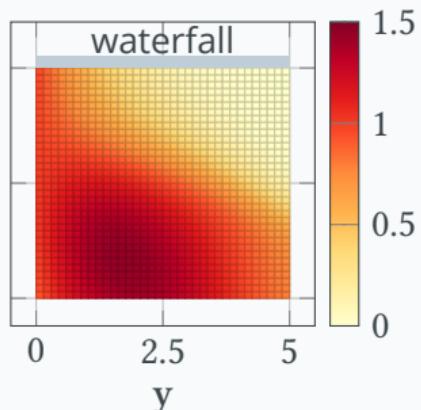
Flow

$$\mathbb{E}[\Delta \mathbf{x}_{t+1} | \text{no drop}]$$

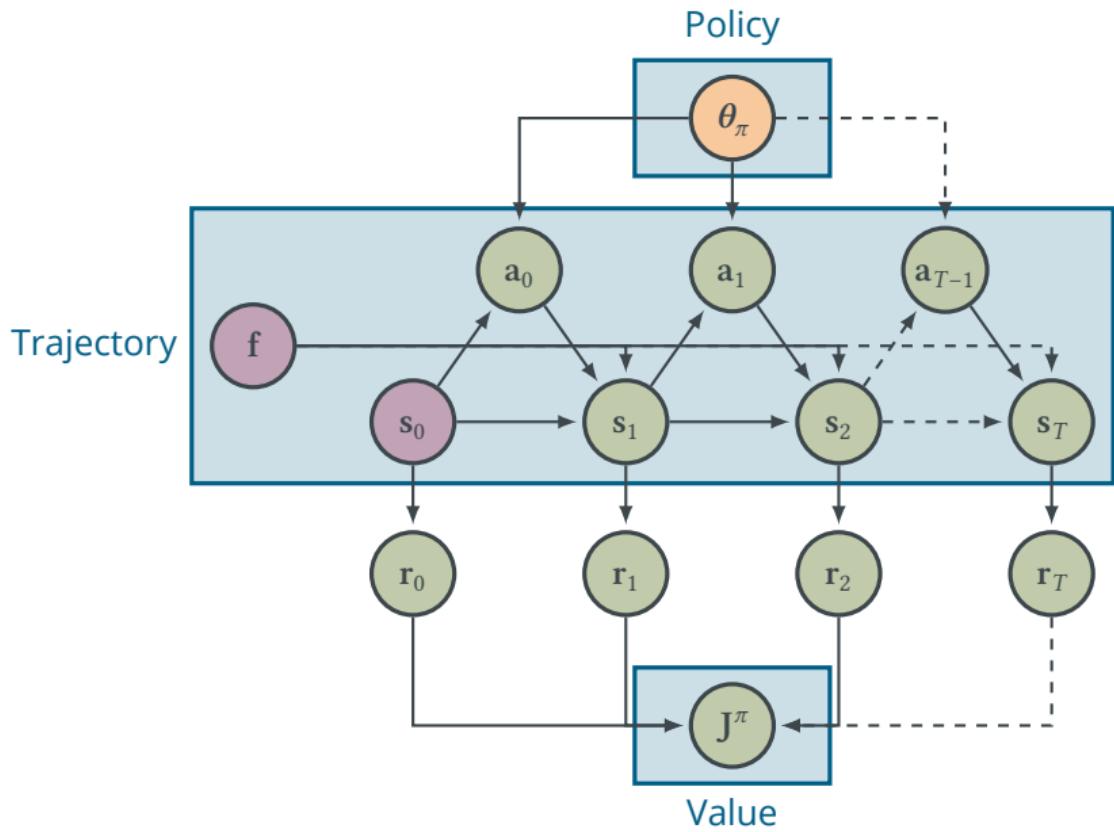


Turbulence

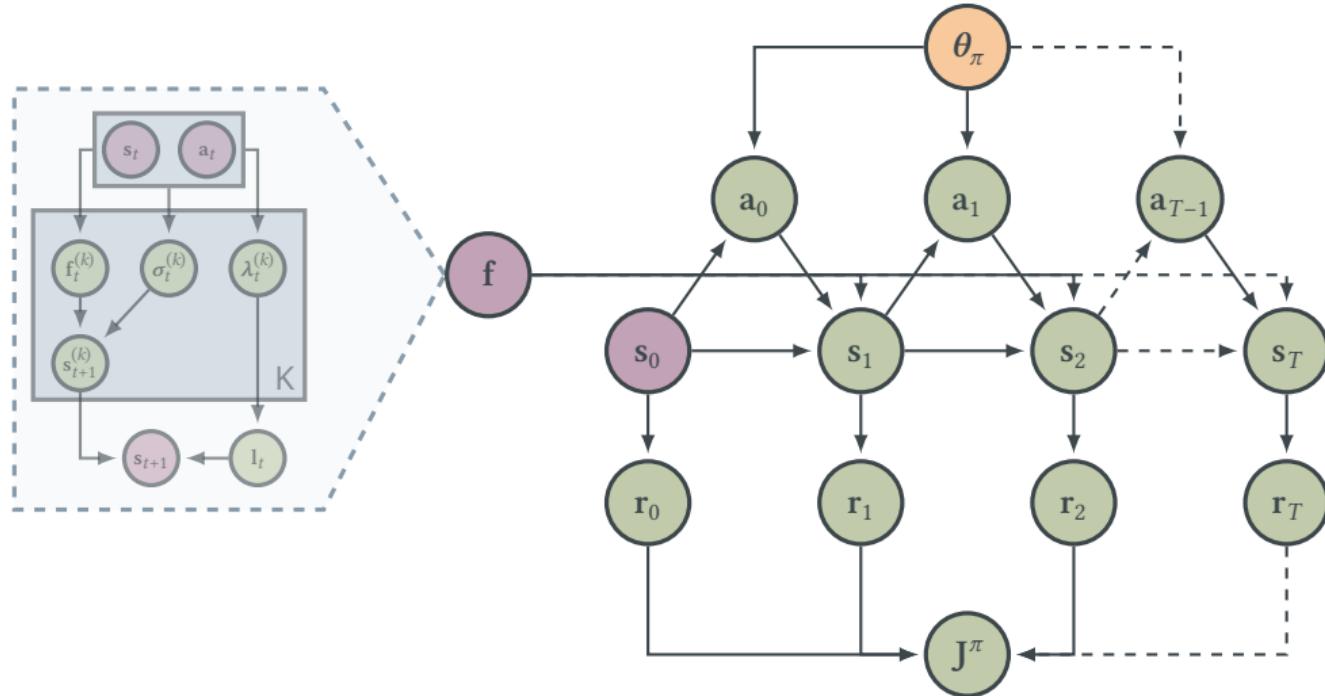
$$\sigma[\Delta \mathbf{x}_{t+1} | \text{no drop}]$$



# Policy: Graphical Model



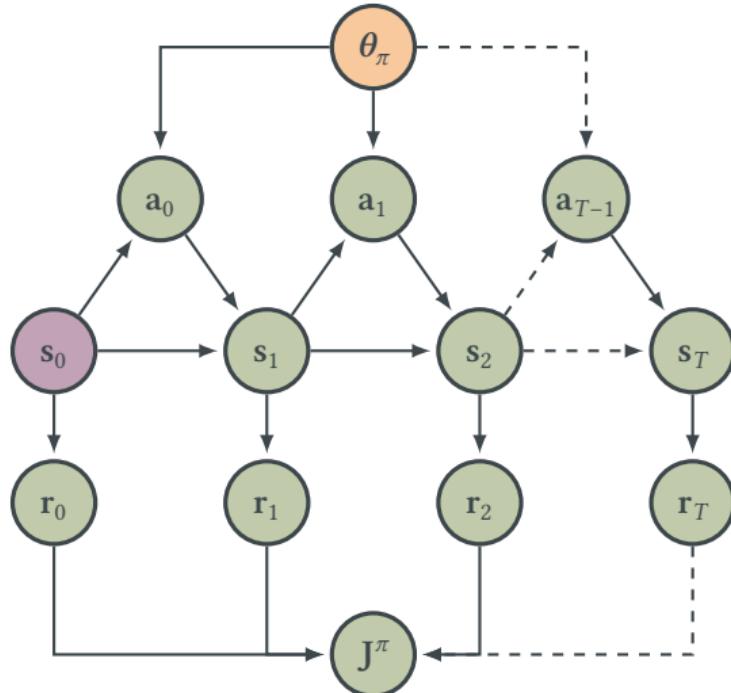
## Policy: Graphical Model



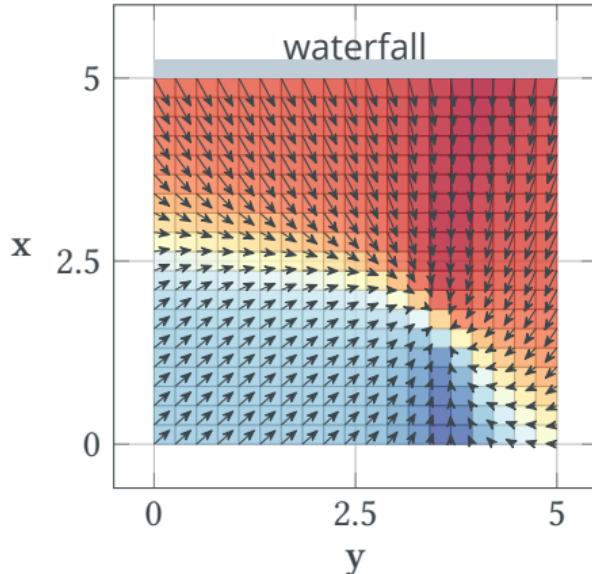
## Policy: Graphical Model

$$\begin{aligned}\mathbb{E}[J^\pi(\theta_\pi)] &= \sum_{t=0}^T \gamma^t \mathbb{E}_{p(s_t|\theta_\pi)}[r_t] \\ &\approx \frac{1}{P} \sum_{p=1}^P \sum_{t=0}^T \gamma^t r_t^{(p)}\end{aligned}$$

$$\nabla J^\pi(\theta_\pi) \approx \frac{1}{P} \sum_{p=1}^P \sum_{t=0}^T \gamma^t \nabla_{\theta_\pi} r_t^{(p)}$$



## Policy: Posterior



N	NFQ <sup>2</sup>	GP <sup>3</sup>	Ours
100	0.66(16)	<b>1.41(1)</b>	1.18(9)
250	1.71(7)	1.54(1)	<b>2.33(1)</b>
500	1.60(10)	1.56(1)	<b>2.25(1)</b>
1000	1.99(6)	2.13(1)	<b>2.32(1)</b>
2500	2.26(2)	1.91(1)	<b>2.28(1)</b>
5000	<b>2.33(1)</b>	1.91(1)	2.28(1)

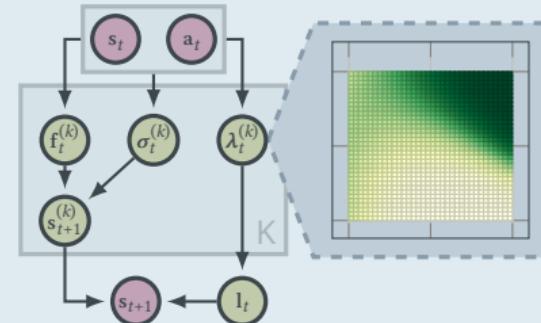
<sup>2</sup>Riedmiller 2005

<sup>3</sup>Deisenroth and Rasmussen 2011

# Industrial Reinforcement Learning

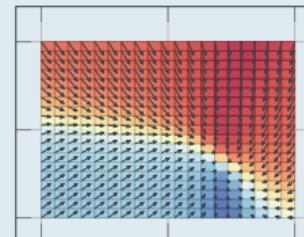
## Interaction with domain experts

- Incomplete system knowledge
- Hierarchical priors
- Interpretable sub-models



## Trustworthy decision making

- Uncertainty due to incomplete data
- Stochastic systems
- Robust and efficient inference



## References

---

-  Deisenroth, Marc and Carl E. Rasmussen (2011). "PILCO: A Model-Based and Data-Efficient Approach to Policy Search". In: *Proceedings of the 28th International Conference on Machine Learning (ICML-11)*, pp. 465–472.
-  Hans, Alexander and Steffen Udluft (2009). "Efficient Uncertainty Propagation for Reinforcement Learning with Limited Data". In: *International Conference on Artificial Neural Networks*. Springer, pp. 70–79.
-  Kaiser, Markus et al. (Oct. 16, 2018). "Data Association with Gaussian Processes". In: arXiv: 1810.07158 [cs, stat].
-  Riedmiller, Martin (2005). "Neural Fitted Q Iteration - First Experiences with a Data Efficient Neural Reinforcement Learning Method". In: *European Conference on Machine Learning*. Springer, pp. 317–328.
-  Tresp, Volker (1994). "The Wet Game of Chicken". In: *Siemens AG, CT IC 4, Technical Report*.

## Variational Bound

$$\begin{aligned}\mathcal{L}_{\text{DAGP}} &= \mathbb{E}_{q(F, \lambda, U)} \left[ \log \frac{p(S', L, F, \lambda, U | S)}{q(F, \lambda, U)} \right] \\ &= \sum_{n=1}^N \mathbb{E}_{q(f_n)} [\log p(s'_n | f_n, l_n)] + \sum_{n=1}^N \mathbb{E}_{q(\lambda_n)} [\log p(l_n | \lambda_n)] \\ &\quad - \sum_{k=1}^K \text{KL}(q(u^{(k)}) \| p(u^{(k)} | Z^{(k)})) - \sum_{k=1}^K \text{KL}(q(u_\lambda^{(k)}) \| p(u_\lambda^{(k)} | Z_\lambda^{(k)}))\end{aligned}$$

## Predictive Posterior

$$\begin{aligned} q(s_{t+1} | s_t) &= \int \sum_{k=1}^K q(l_t^{(k)} | s_t) q(s_{t+1}^{(k)} | s_t) dl_t^{(1)} \dots dl_t^{(K)} \\ &\approx \sum_{k=1}^K \hat{l}_t^{(k)} \hat{s}_{t+1}^{(k)} \end{aligned}$$