

Incorporating Uncertainty into Reinforcement Learning through Gaussian Processes

Master's Thesis

Markus Kaiser

4. July 2016



SIEMENS

1 The Bicycle Benchmark

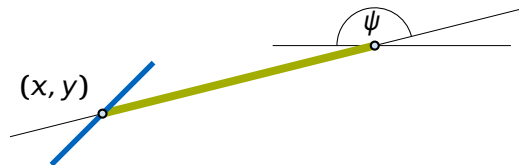
2 Gaussian Processes

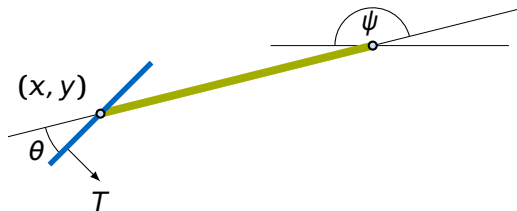
3 Model-Uncertainties in Reinforcement Learning

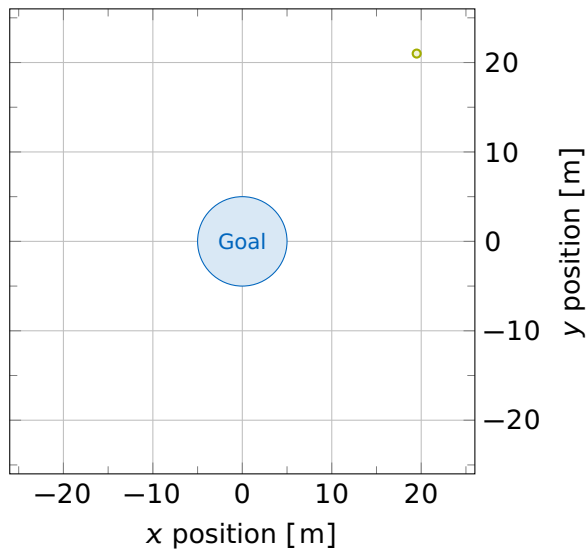
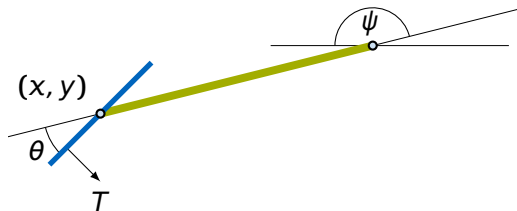
4 Results

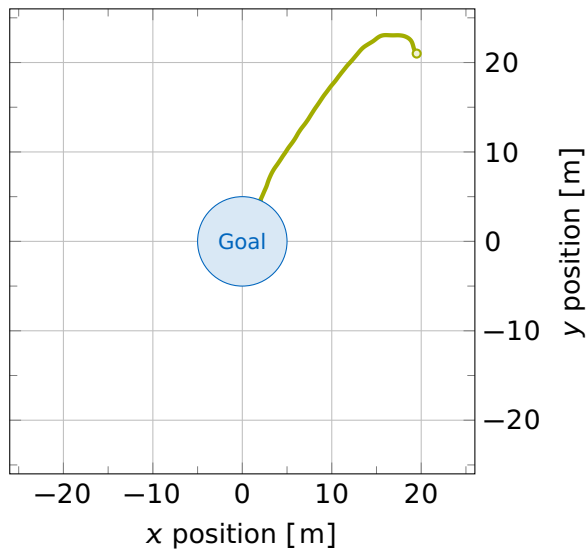
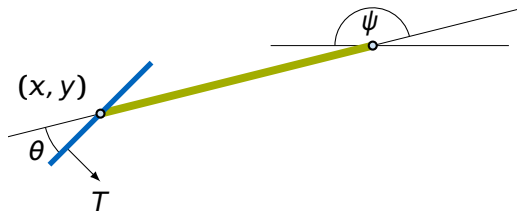


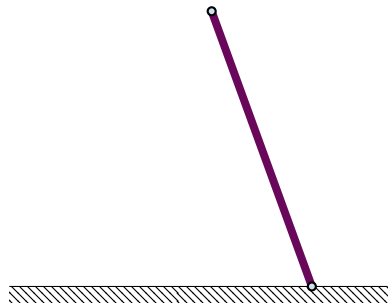


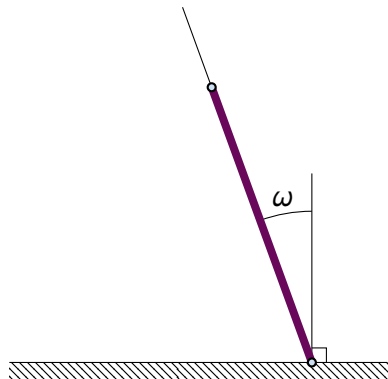


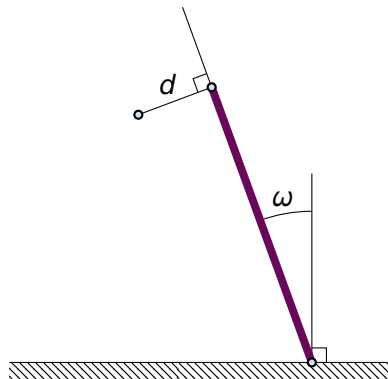


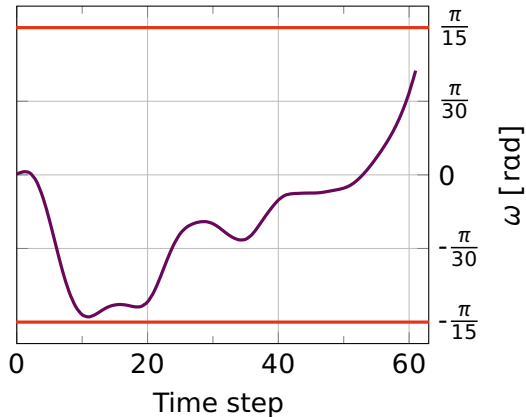
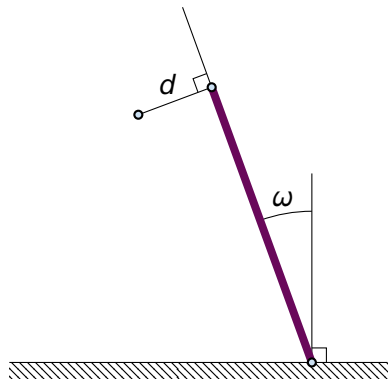












State Variables

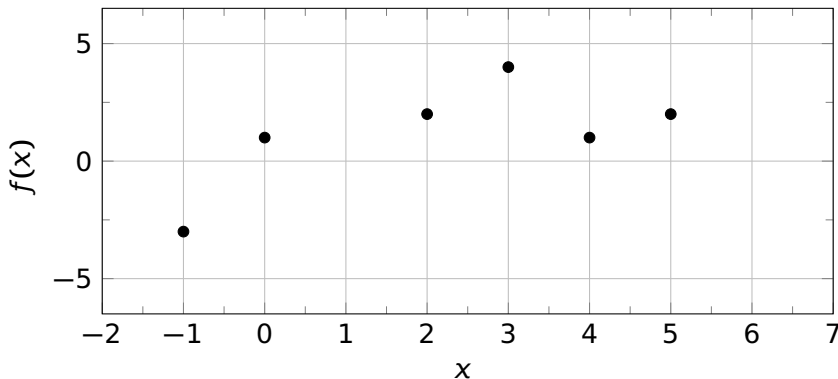
Notation	Description
θ	Steering Angle
$\dot{\theta}$	
ω	Leaning Angle
$\dot{\omega}$	
x	Front tyre position
y	
ψ	Bicycle orientation

Actions

Notation	Description
d	Lean distance
T	Steering force

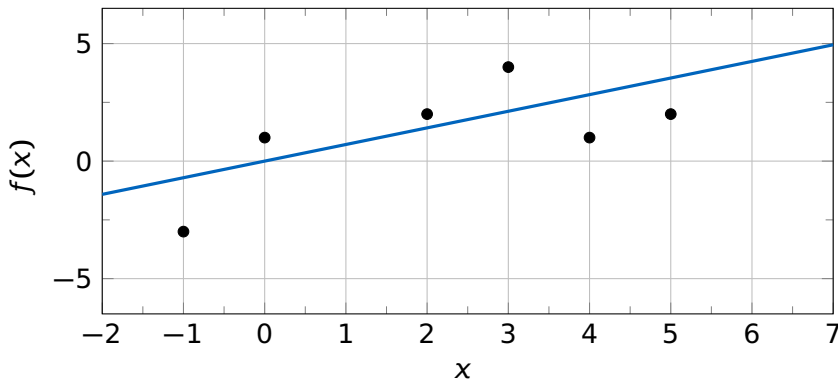
Data

- Observe a **noisy** data set $\{(\mathbf{x}_i, \mathbf{y}_i)\}$
- Assume that $\mathbf{y}_i = f(\mathbf{x}_i) + \epsilon$



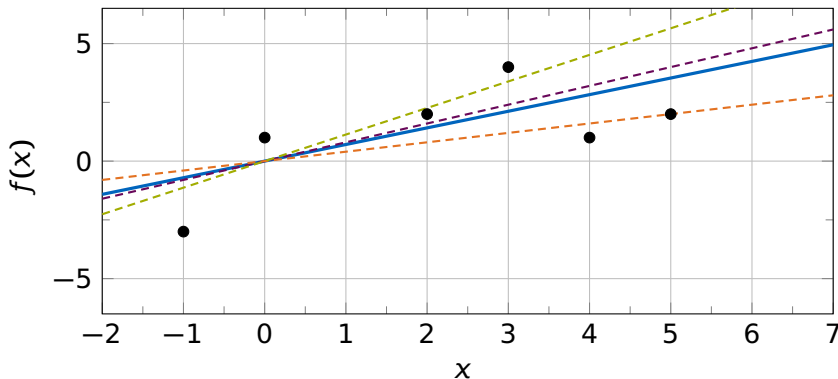
Parametric Models

- Assume **structure** about f , such as $f(\mathbf{x}) = \mathbf{W}\mathbf{x}$
- Find a **single \mathbf{W}** which is optimal w.r.t. some criterion



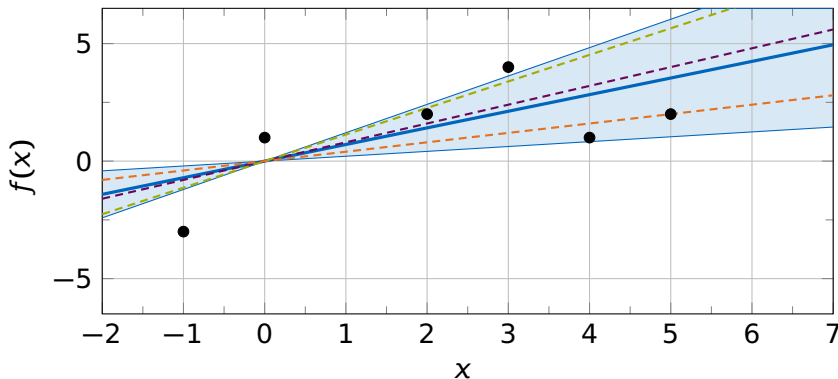
Parametric Models

- Assume **structure** about f , such as $f(\mathbf{x}) = \mathbf{W}\mathbf{x}$
- Find a **single \mathbf{W}** which is optimal w.r.t. some criterion



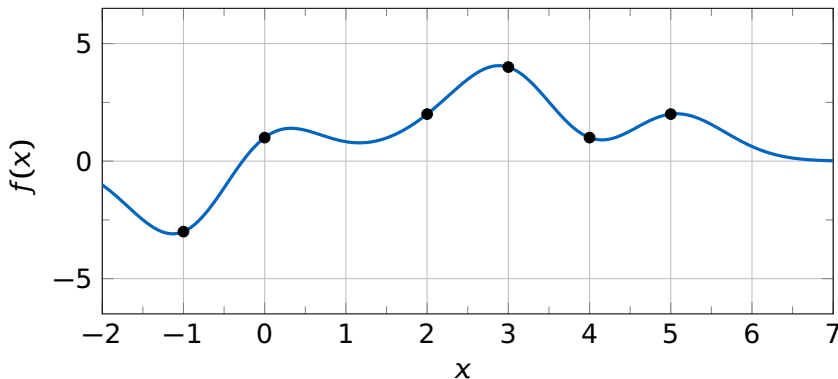
Bayesian Parametric Models

- Assume **structure** about f , such as $f(\mathbf{x}) = \mathbf{W}\mathbf{x}$
- Find a **distribution** $p(\mathbf{W} | \mathbf{X}, \mathbf{y})$ of plausible parameters



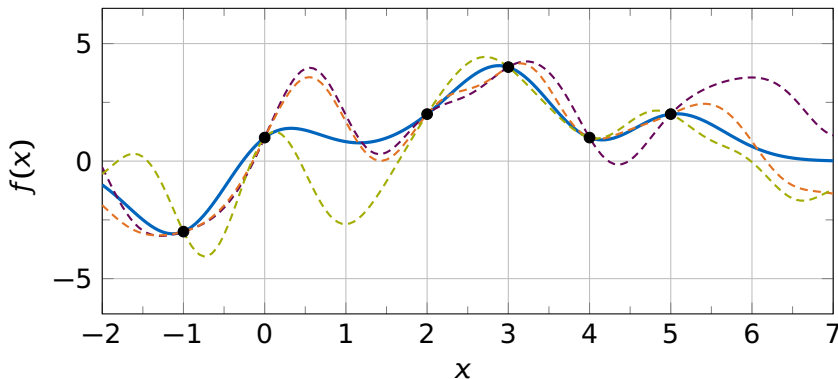
Bayesian Non-Parametric Models

- f can be an arbitrary function
- Find a distribution $p(f | \mathbf{X}, \mathbf{y})$ over functions



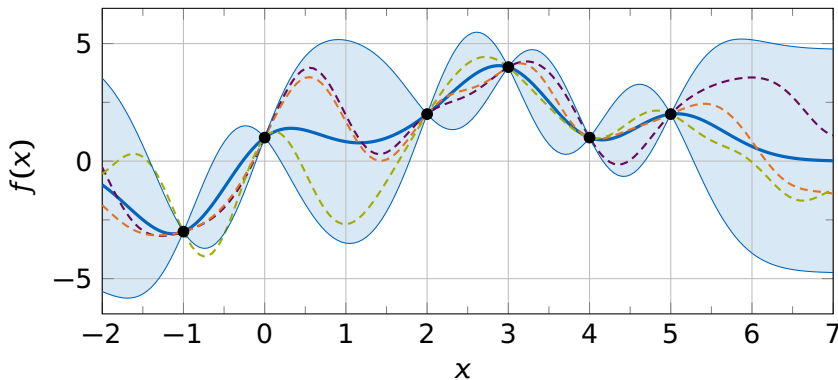
Bayesian Non-Parametric Models

- f can be an arbitrary function
- Find a distribution $p(f | \mathbf{X}, \mathbf{y})$ over functions



Bayesian Non-Parametric Models

- f can be an arbitrary function
- Find a distribution $p(f | \mathbf{X}, \mathbf{y})$ over functions



Definition (Gaussian Process)

A **Gaussian Process (GP)** is a collection of random variables $\{\mathbf{F}_{\mathbf{x}}\}$, any finite subset of which has a joint Gaussian distribution.

- Extension of Gaussians to (infinite) function spaces
- $\mathbf{F}_{\mathbf{x}}$ models the function value $f(\mathbf{x})$

Definition (Gaussian Process)

A **Gaussian Process (GP)** is a collection of random variables $\{\mathbf{F}_{\mathbf{x}}\}$, any finite subset of which has a joint Gaussian distribution.

- Extension of Gaussians to (infinite) function spaces
- $\mathbf{F}_{\mathbf{x}}$ models the function value $f(\mathbf{x})$

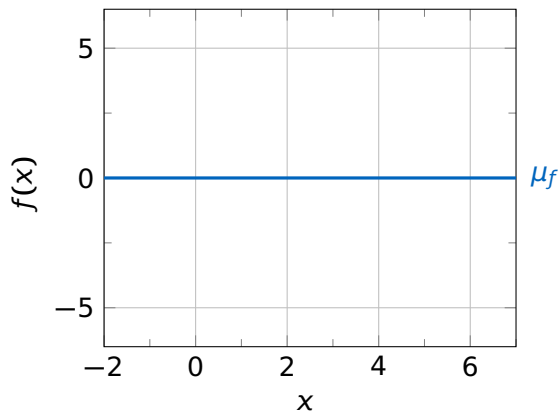
Mean and Kernel Functions

A GP is **completely determined** by two functions.

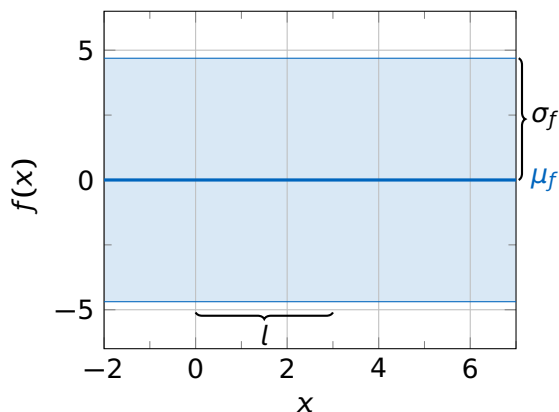
Mean function $\mu_f(\mathbf{x}) = \mathbb{E}[f(\mathbf{x})]$

Kernel function $\mathcal{K}(\mathbf{x}, \mathbf{x}') = \text{cov}[f(\mathbf{x}), f(\mathbf{x}')]]$

- The kernel encodes the **prior assumptions** about the function

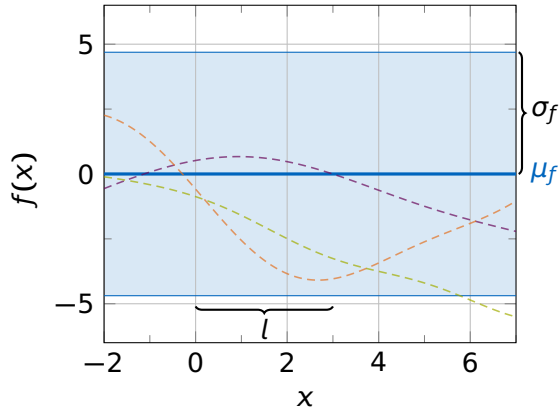


$$\mu_f(x) = 0$$



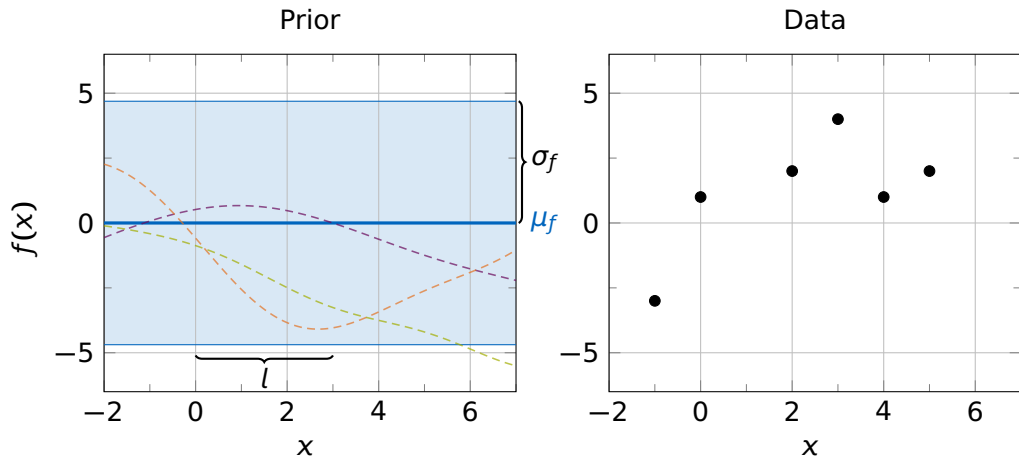
$$\mu_f(x) = 0$$

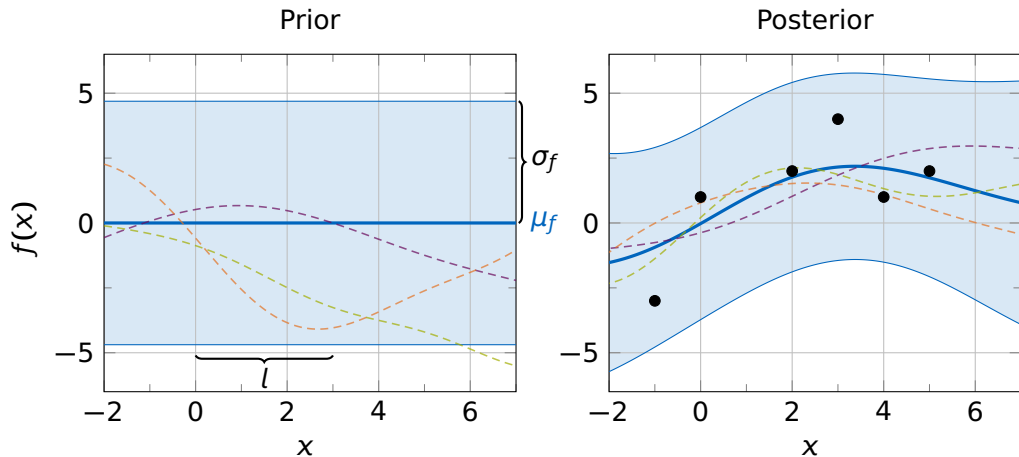
$$\kappa(x, x') = \sigma_f \cdot \exp\left(-\frac{(x - x')^2}{2 \cdot l^2}\right)$$

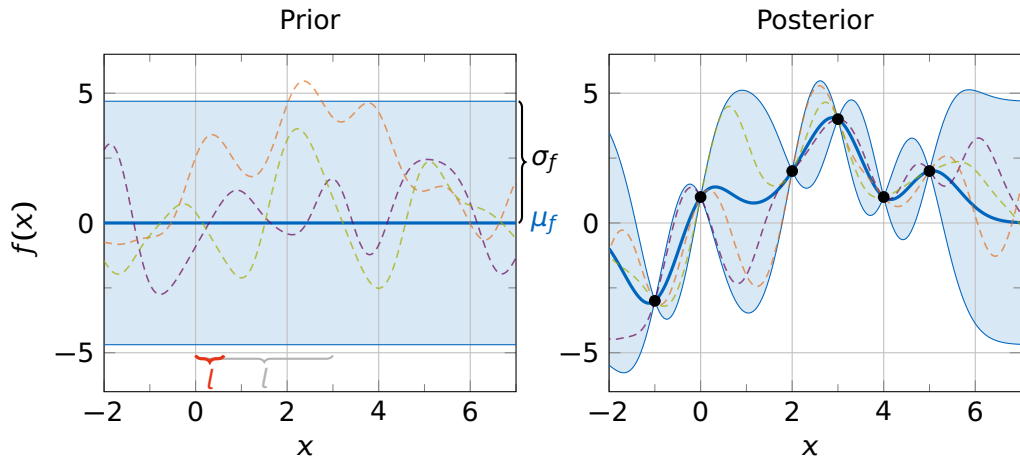


$$\mu_f(x) = 0$$

$$\kappa(x, x') = \sigma_f \cdot \exp\left(-\frac{(x - x')^2}{2 \cdot l^2}\right)$$

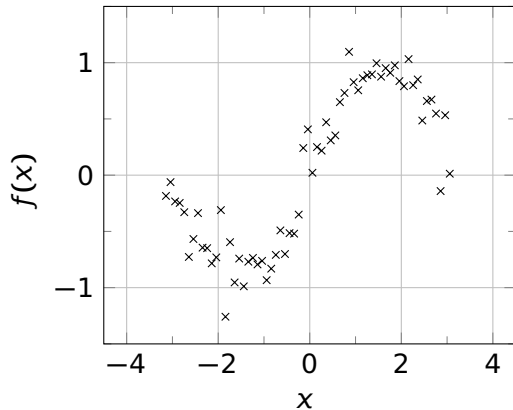






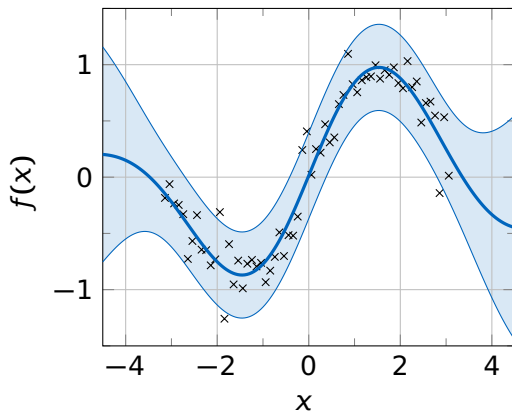
Full Gaussian Processes

- Calculation of the posterior in $\mathcal{O}(N^3)$
- Not feasible for large data sets



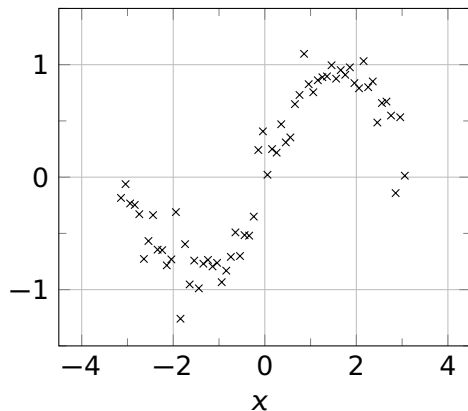
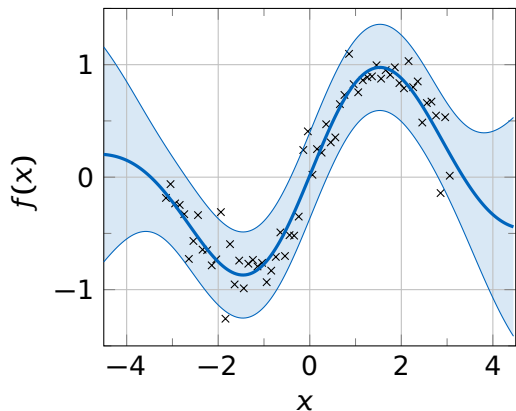
Full Gaussian Processes

- Calculation of the posterior in $\mathcal{O}(N^3)$
- Not feasible for large data sets



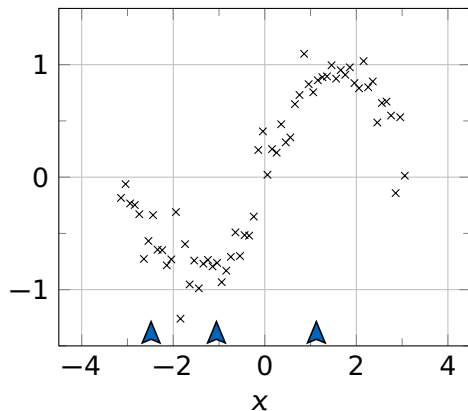
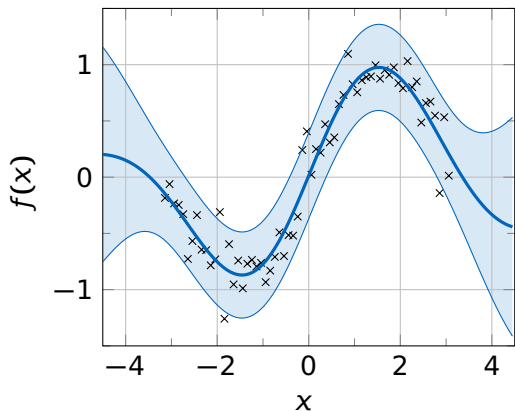
Sparse Gaussian Processes

- Find a set of M Pseudo Inputs
- Calculation of the posterior in $\mathcal{O}(NM^2)$



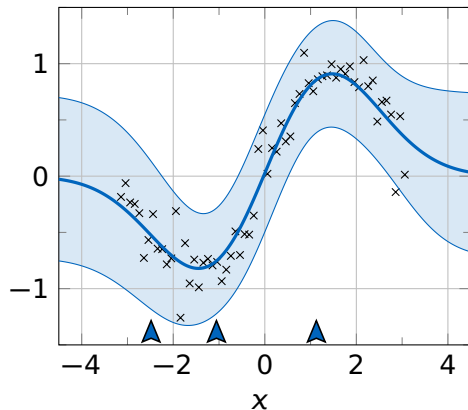
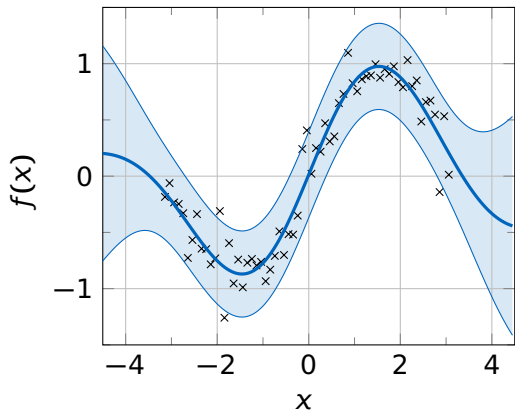
Sparse Gaussian Processes

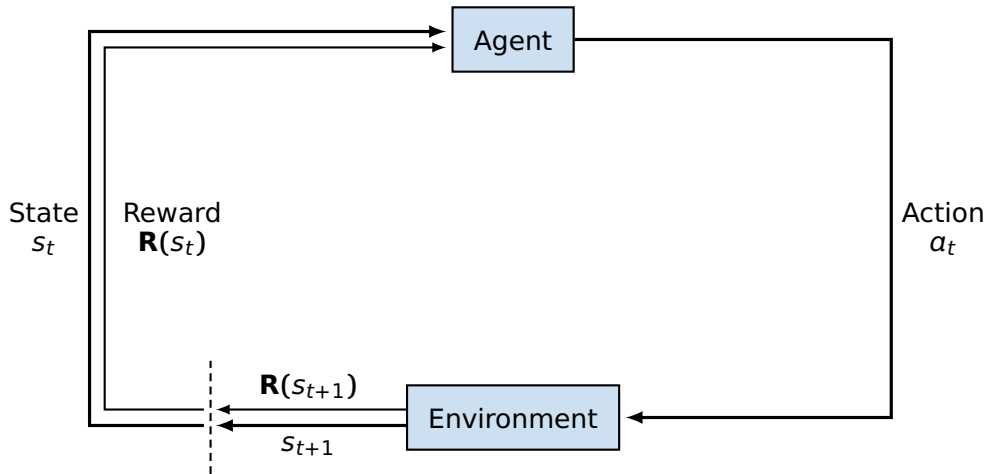
- Find a set of M Pseudo Inputs
- Calculation of the posterior in $\mathcal{O}(NM^2)$

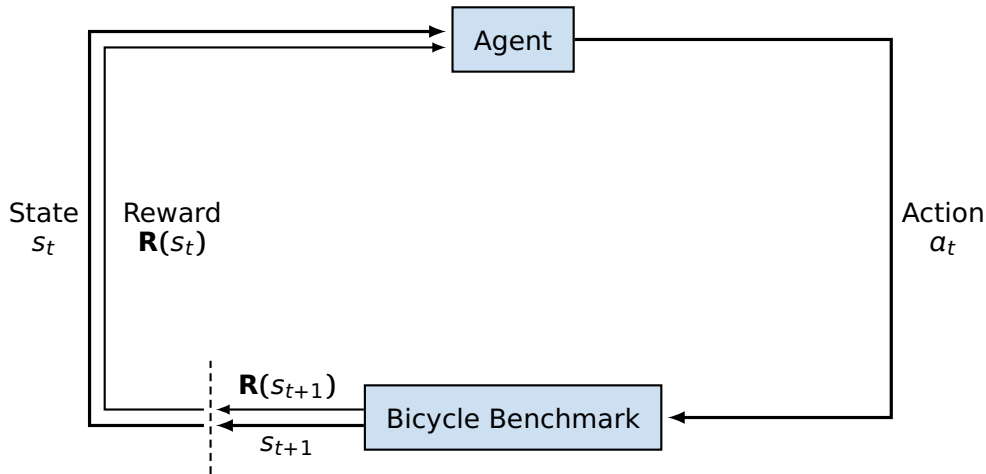


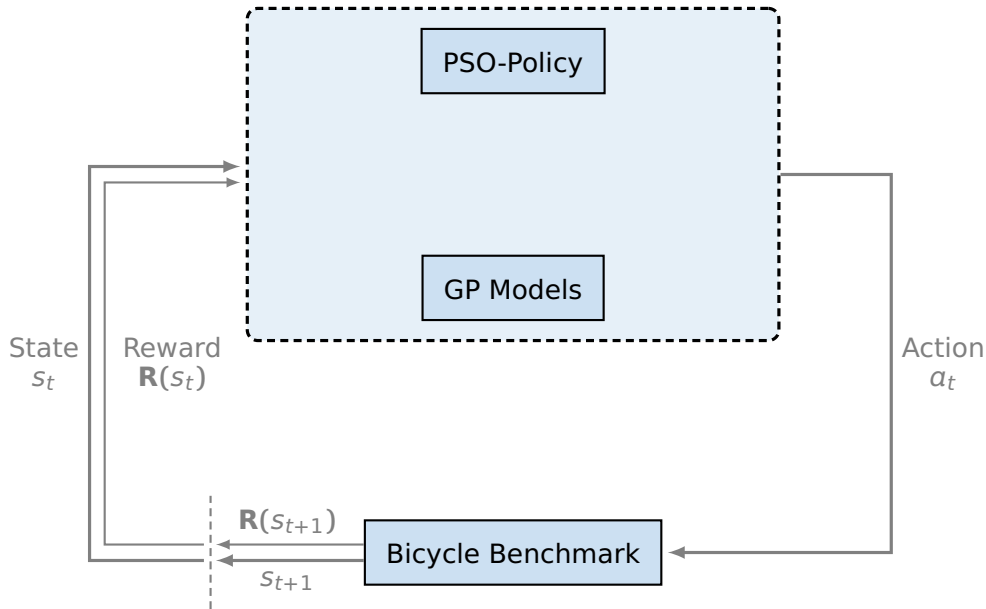
Sparse Gaussian Processes

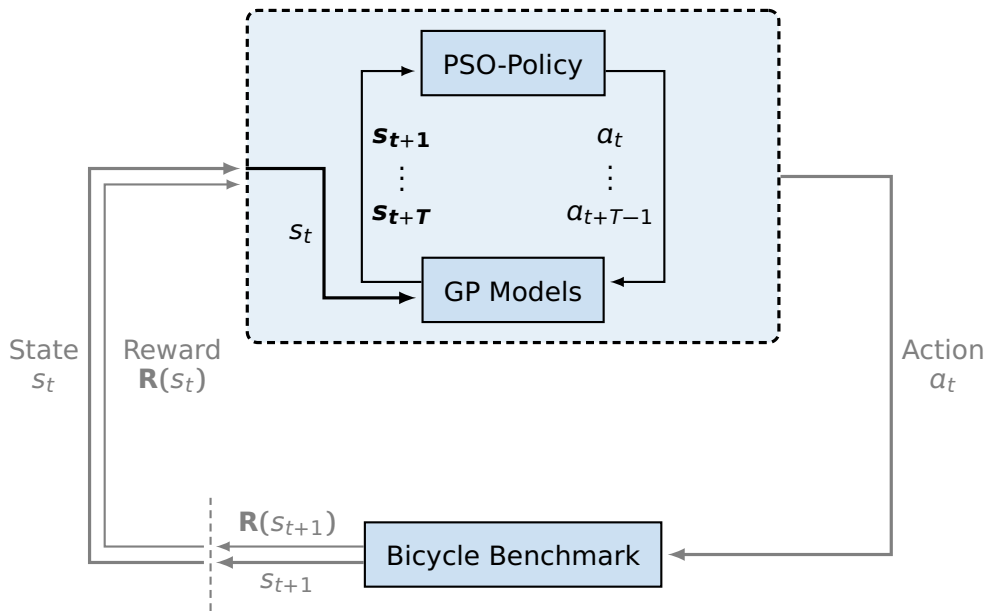
- Find a set of M Pseudo Inputs
- Calculation of the posterior in $\mathcal{O}(NM^2)$

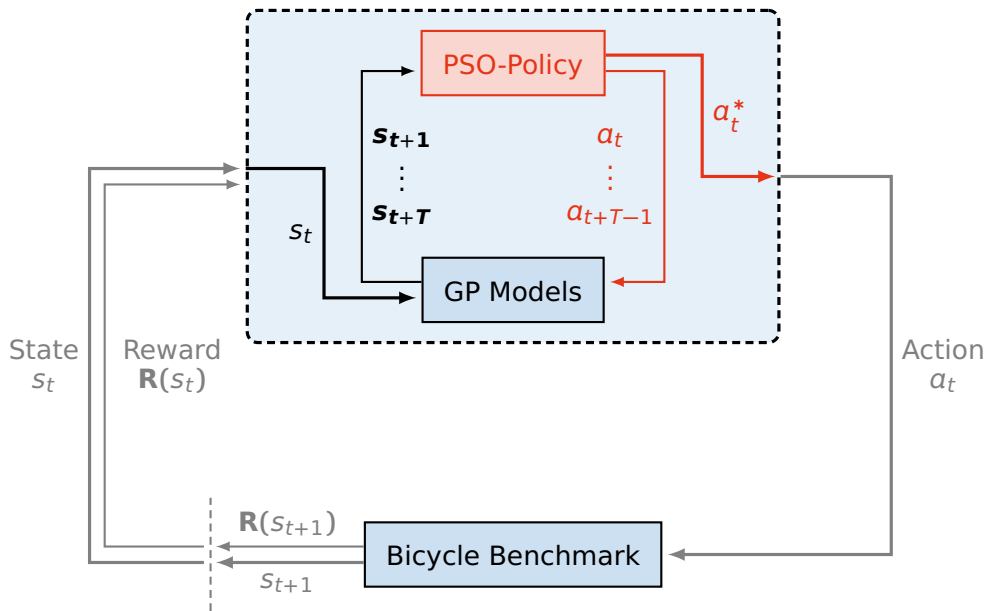












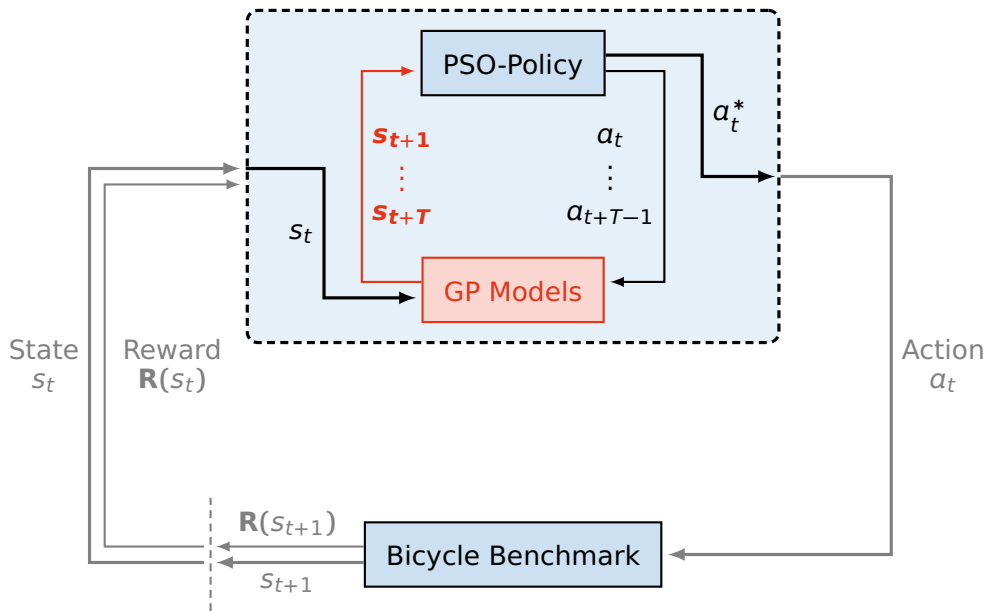
Definition (PSO-Policy)

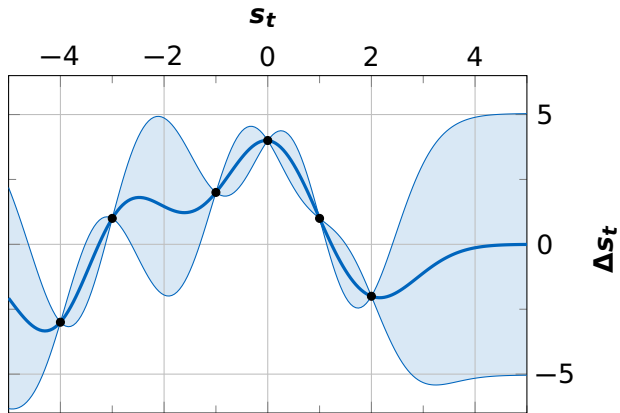
The **Particle Swarm Optimization-Policy (PSO-P)** chooses actions via optimization of the expected accumulated reward.

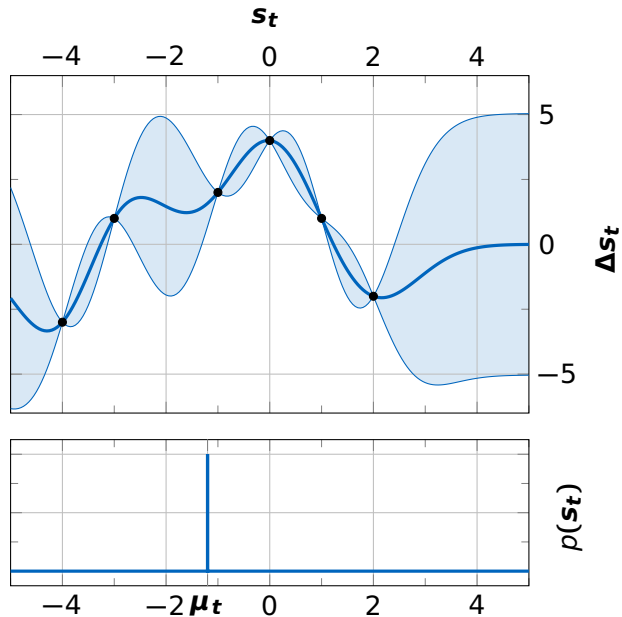
$\pi_{\text{PSO-P}}(\mathbf{s}) := \mathbf{a}_0^*$, where

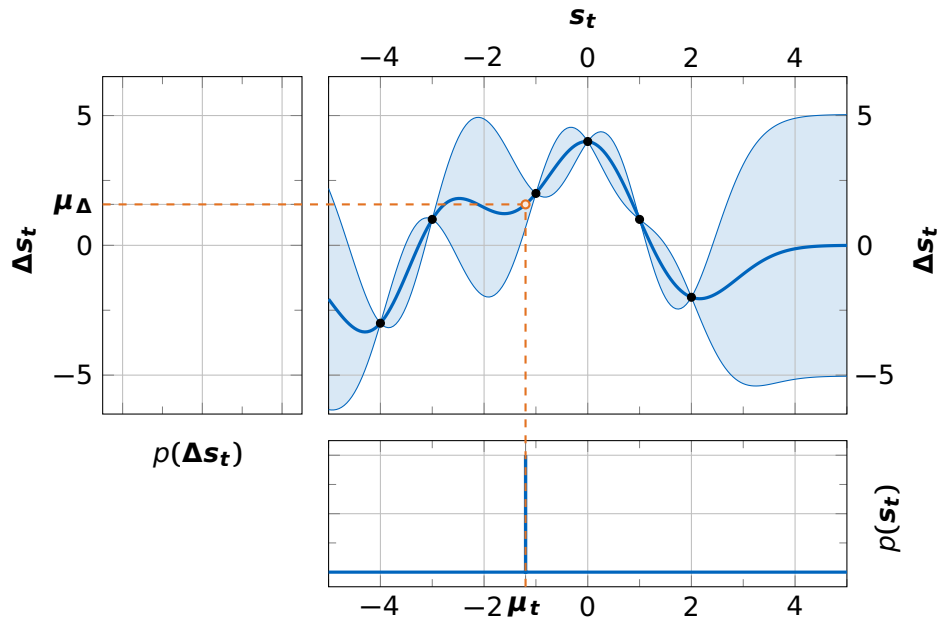
$$\mathbf{a}^* \in \operatorname{argmax}_{\mathbf{a} \in \mathcal{A}^T} \left\{ \mathbb{E} \left[\sum_{t=1}^T \gamma^t \mathbf{R}(\mathbf{s}_t) \mid \mathcal{GP}, \mathbf{s}_0 = \mathbf{s}, \mathbf{a}_0, \dots, \mathbf{a}_{T-1} \right] \right\}$$

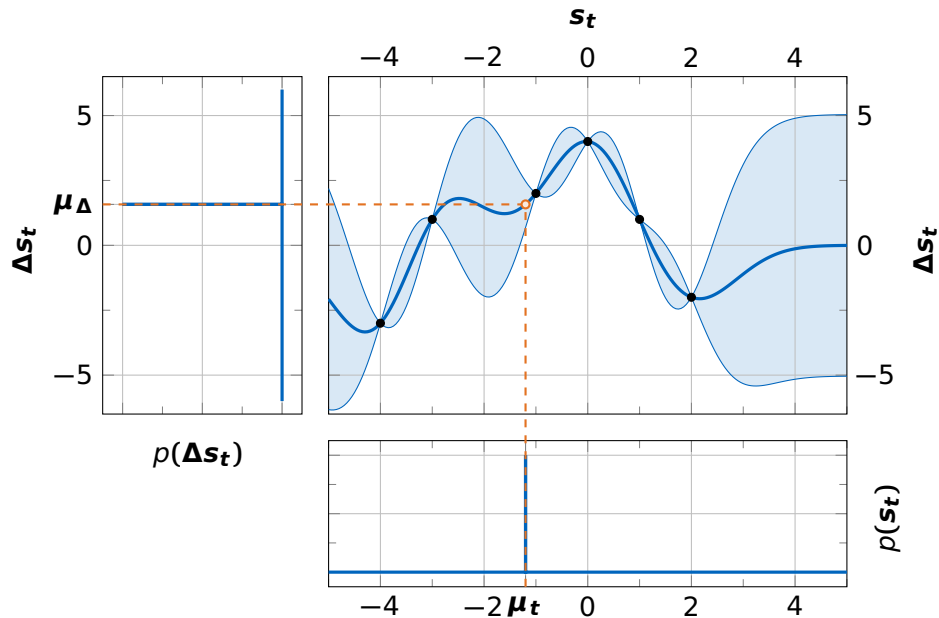
- PSO is a gradient-free heuristic
- Directly exploits the transition models

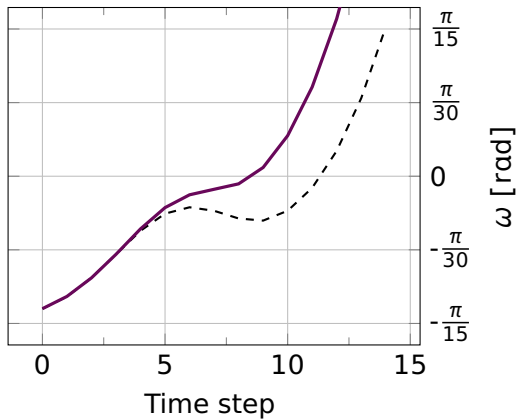
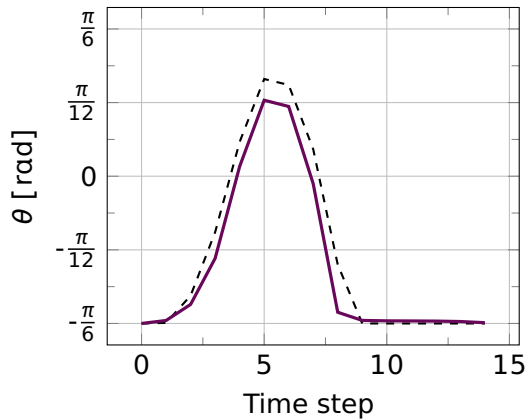


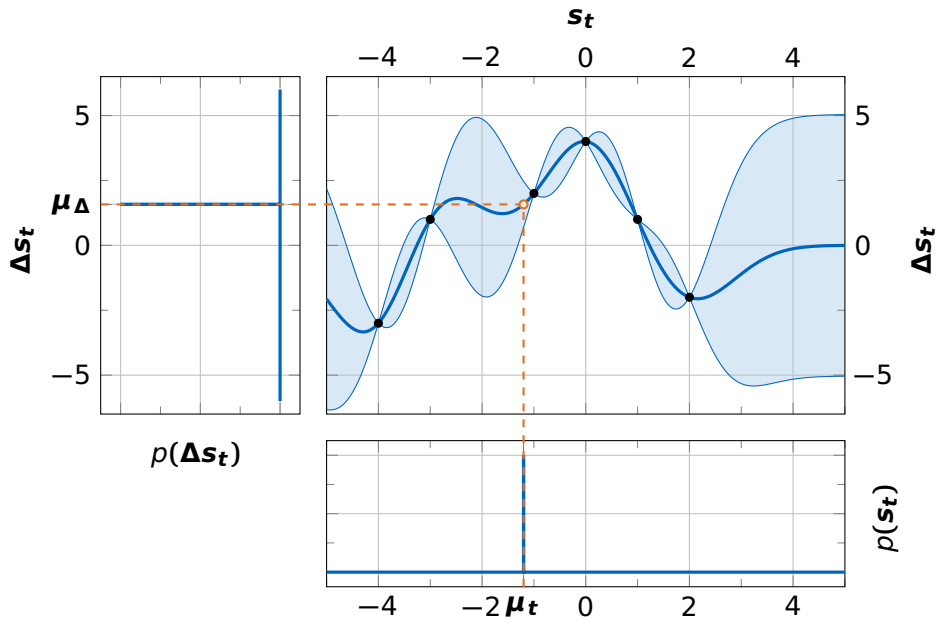


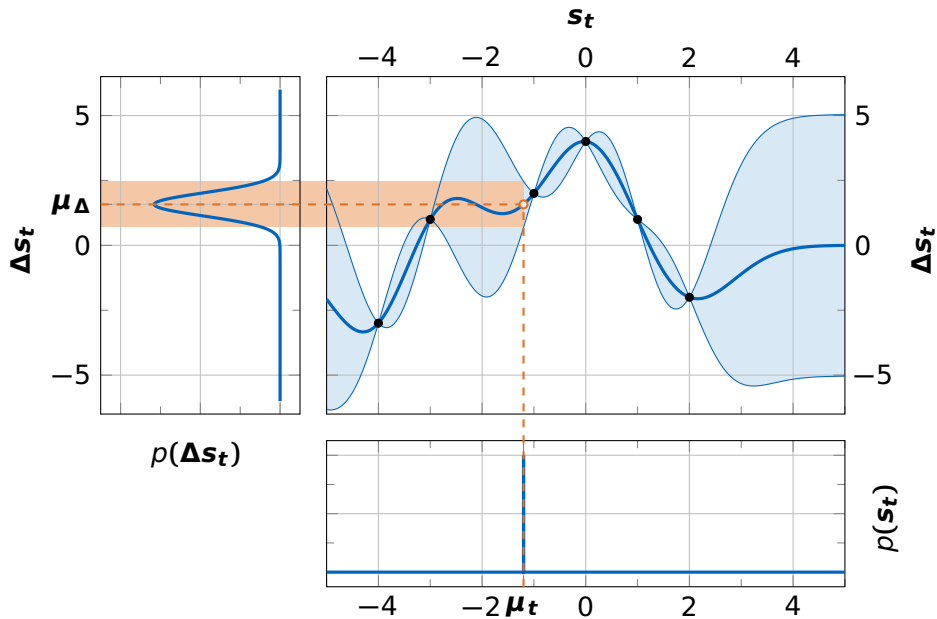


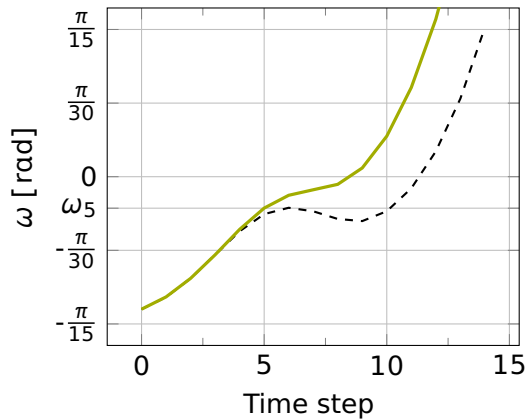


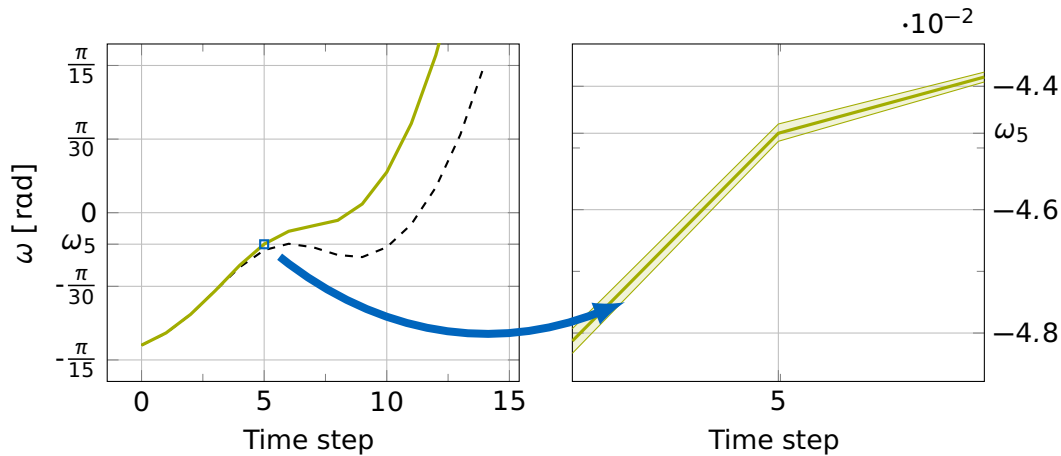


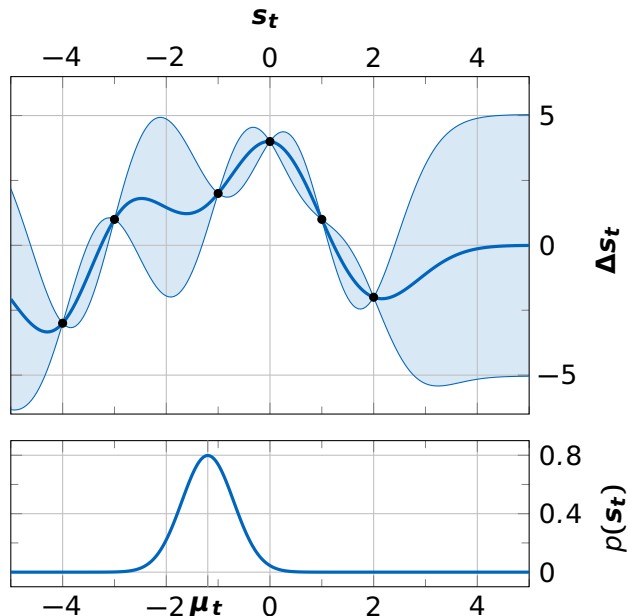


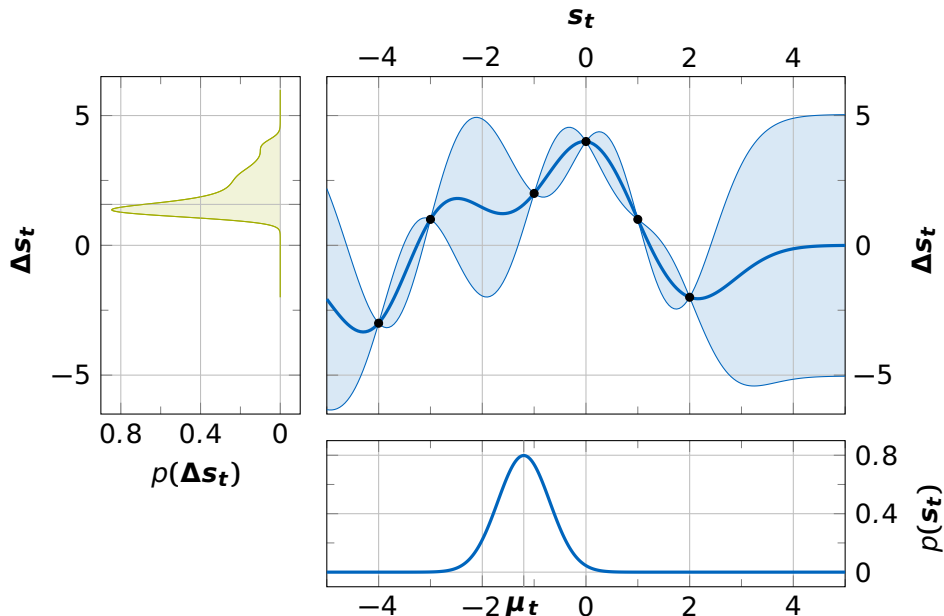


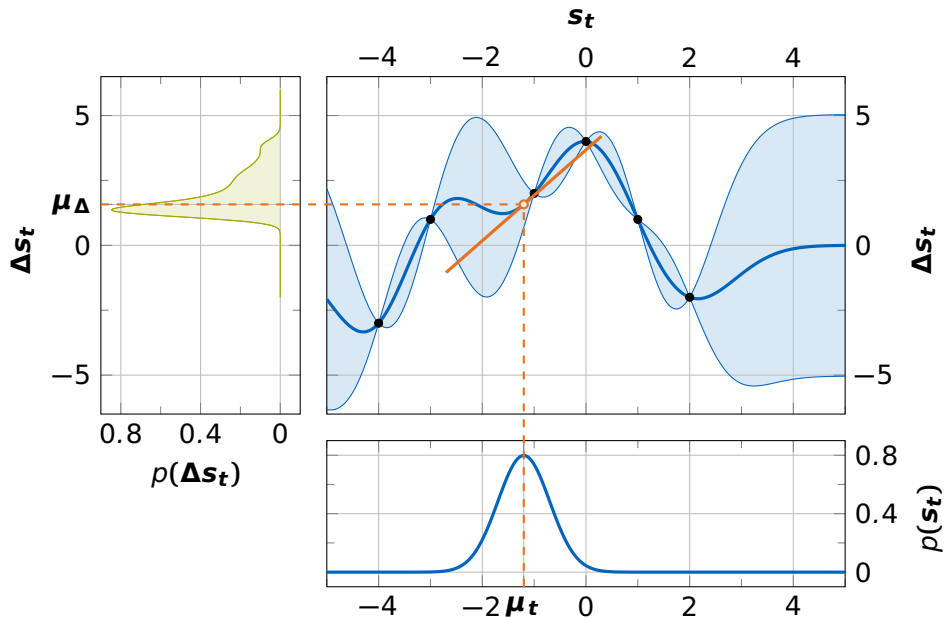


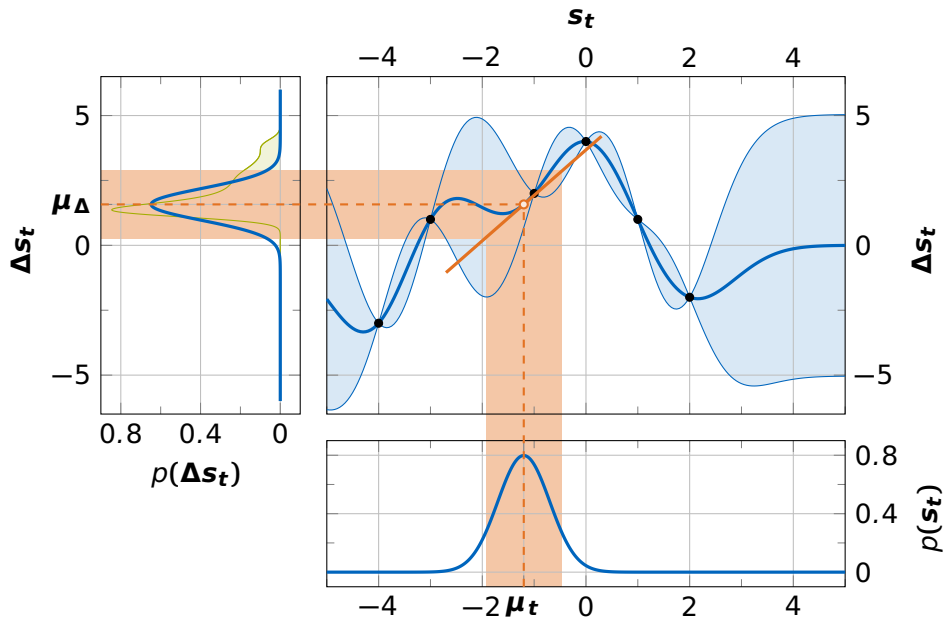


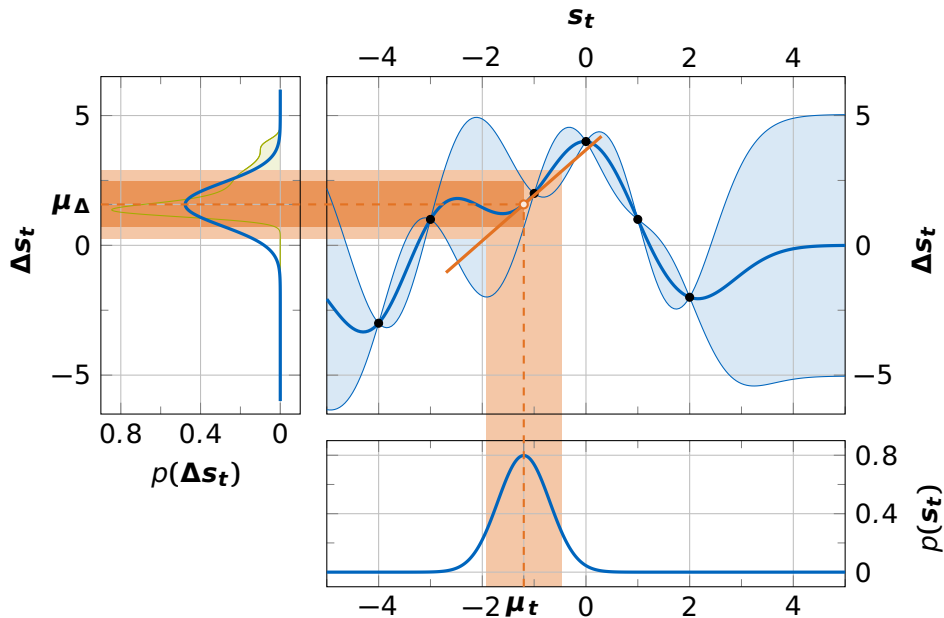


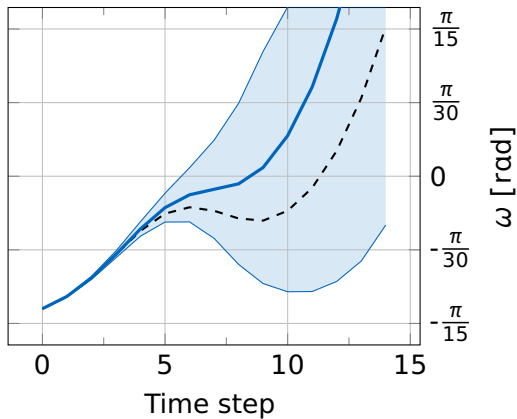
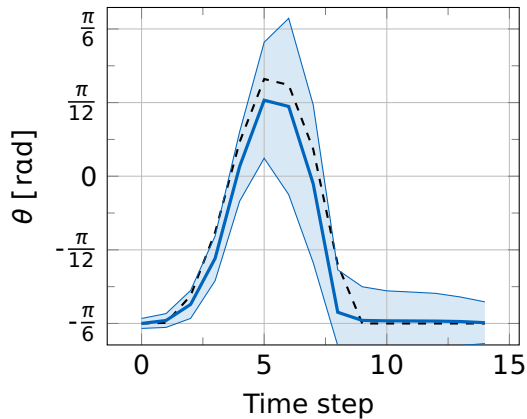


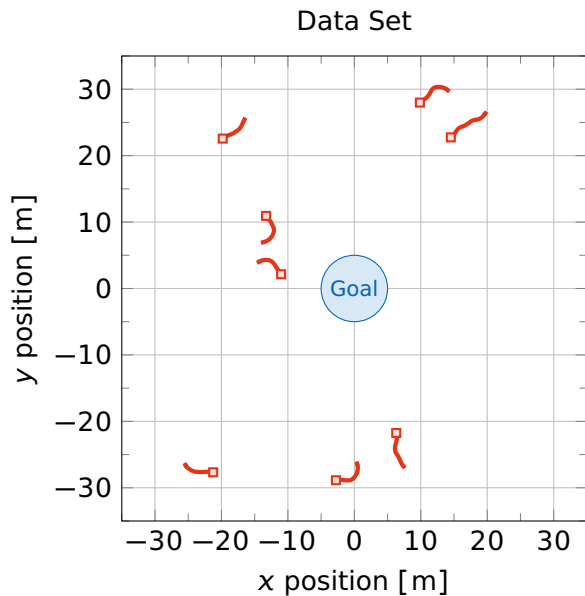




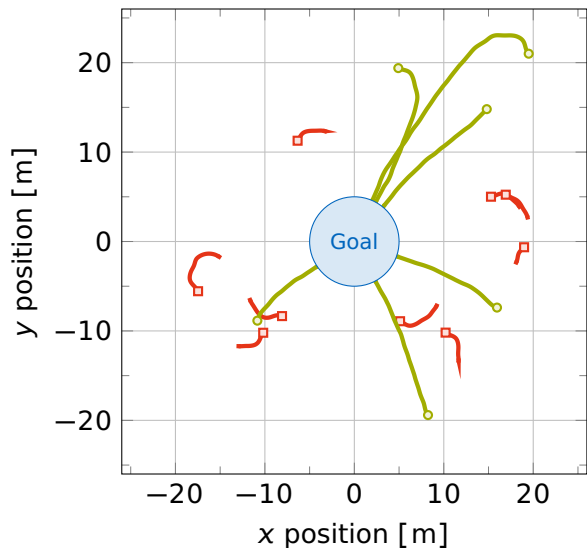




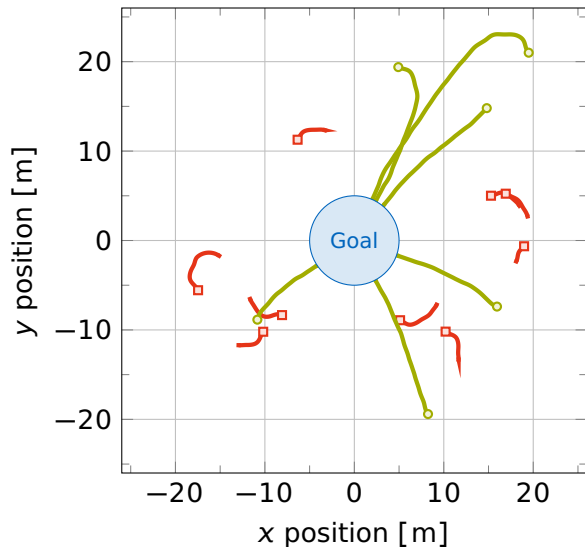




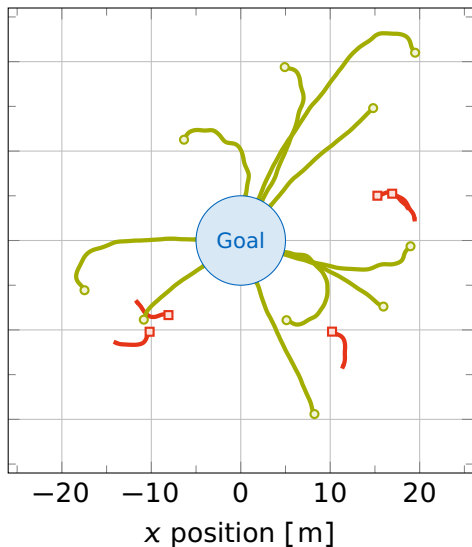
MAP Predictions



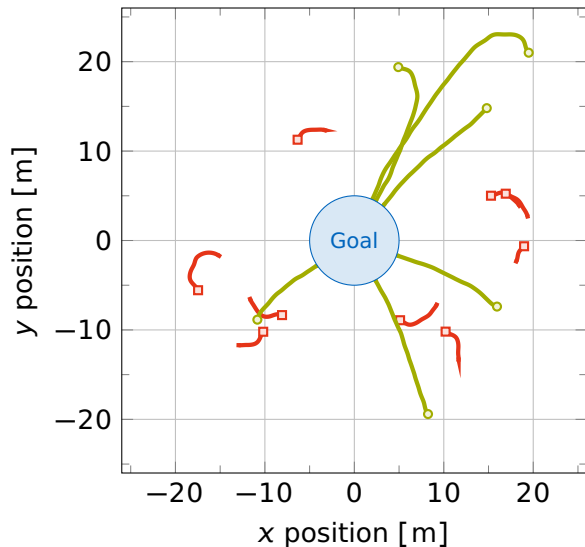
MAP Predictions



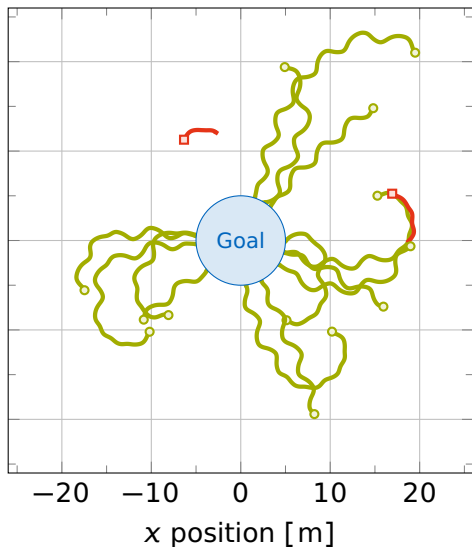
One-Step Uncertainties



MAP Predictions



Multi-Step Uncertainties



Metric	MAP	OS	MS
Trajectories	9660	9660	9660

MAP Deterministic predictions

OS One-Step Uncertainties

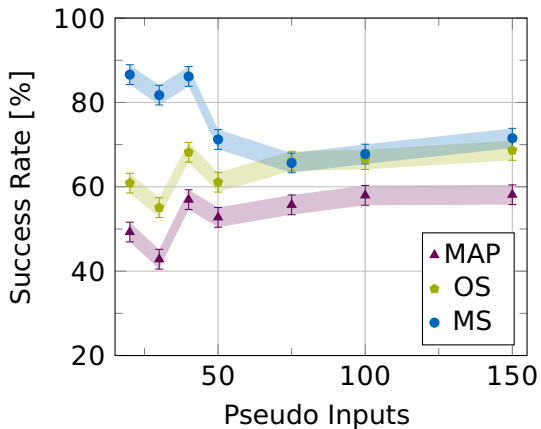
MS Multi-Step Uncertainties

Metric	MAP	OS	MS
Trajectories	9660	9660	9660
Success Rate	53.4 %	63.8 %	75.8 %
Time to Goal			
Mean	59.9	62.0	66.5
Median	60	60	63

MAP Deterministic predictions

OS One-Step Uncertainties

MS Multi-Step Uncertainties



Metric	MAP	OS	MS
Trajectories	9660	9660	9660
Success Rate	53.4 %	63.8 %	75.8 %
Time to Goal			
Mean	59.9	62.0	66.5
Median	60	60	63

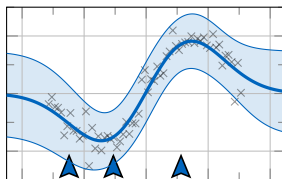
MAP Deterministic predictions

OS One-Step Uncertainties

MS Multi-Step Uncertainties

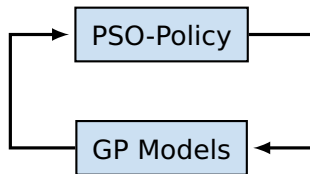
■ Sparse GPs

- Scale to large data sets
- First test at LSY



■ Combination of GP and PSO

- Directly exploit models
- Fully non-parametric

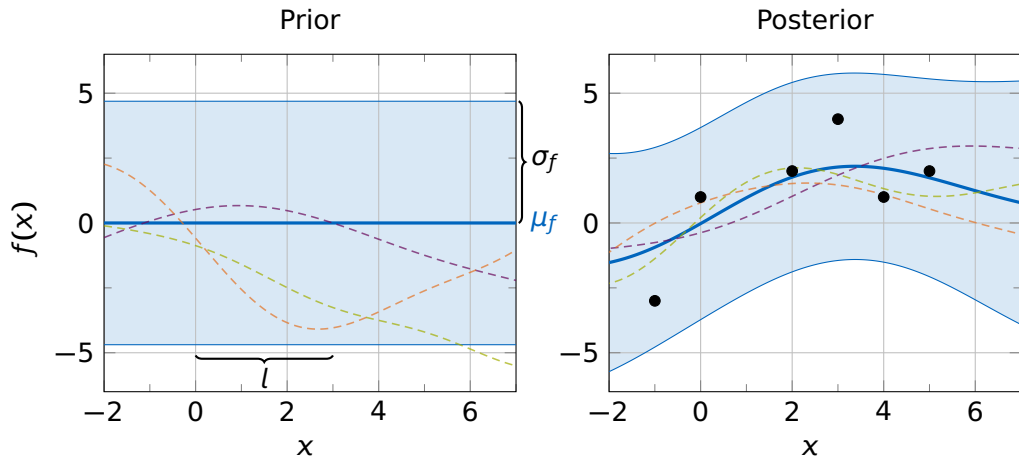


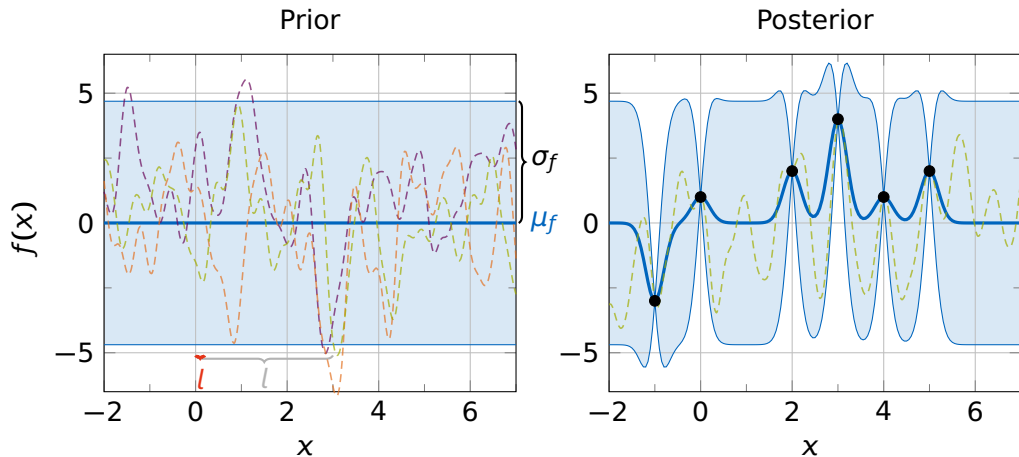
■ Bayesian Models in RL

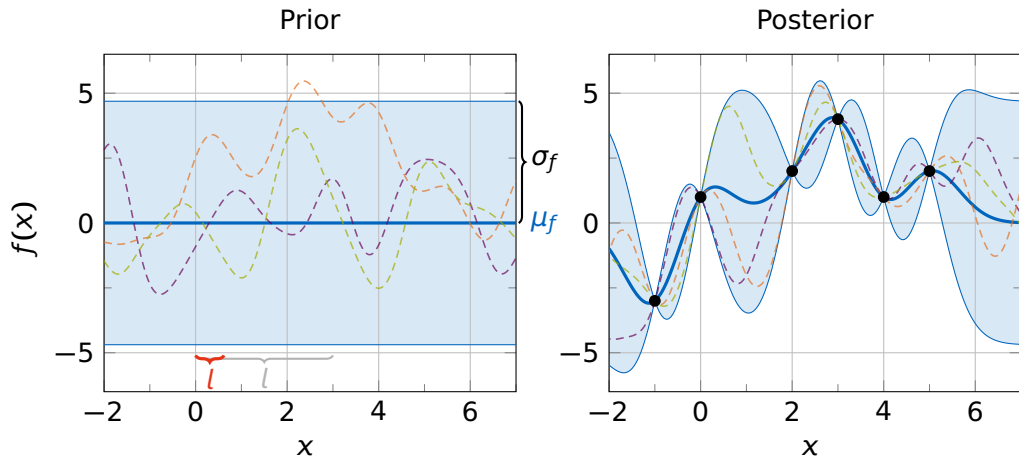
- Reduce model-bias
- Improve performance

Metric	MAP	OS	MS
Success Rate	53.4 %	63.8 %	75.8 %

Additional Material







Definition (Bicycle Reward Function)

The **bicycle reward function** is defined as

$$\mathbf{R}_{\text{bicycle}}(\mathbf{s}) := \begin{cases} 2 & \text{if goal reached} \\ 0 & \text{if fallen down} \\ c \cdot \mathcal{N}(\Delta_{\mathbf{s}}^{\psi} \mid 0, \sigma_{\text{angle}}^2) & \text{otherwise} \end{cases}$$

where $\Delta_{\mathbf{s}}^{\psi}$ denotes the **angle towards the goal**.

