

Approximate Inference, The Right Way

John Cunningham



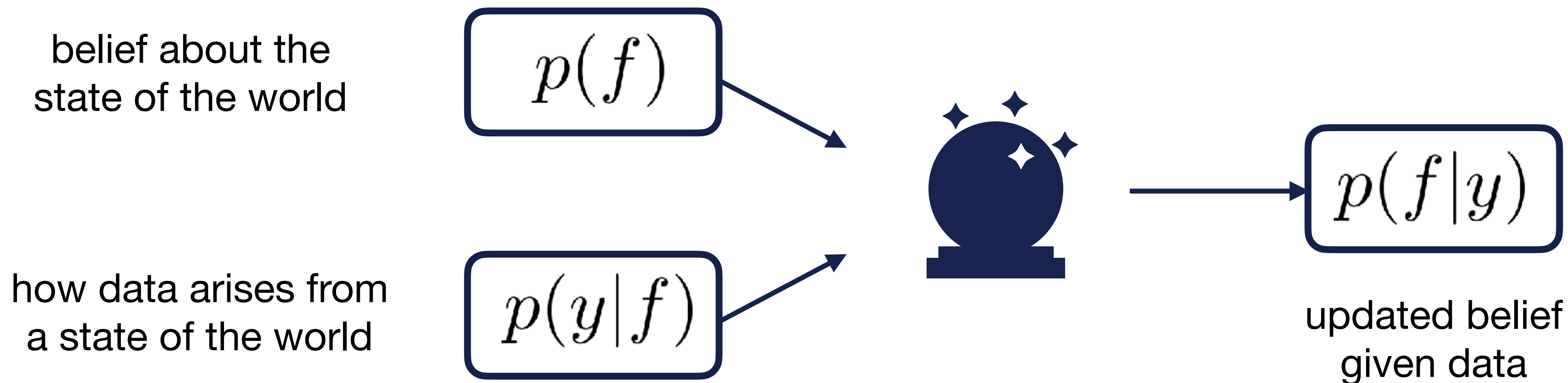
PN School — 27 March 2023

Outline

- ▶ The Promise of Probabilistic Machine Learning
- ▶ Gaussian Process Introduction
- ▶ Scaling Gaussian Processes, and Implications
- ▶ Approximate Gaussian Process Inference, The Right Way
- ▶ iterGP as Probabilistic Numerics
- ▶ Broader Implications

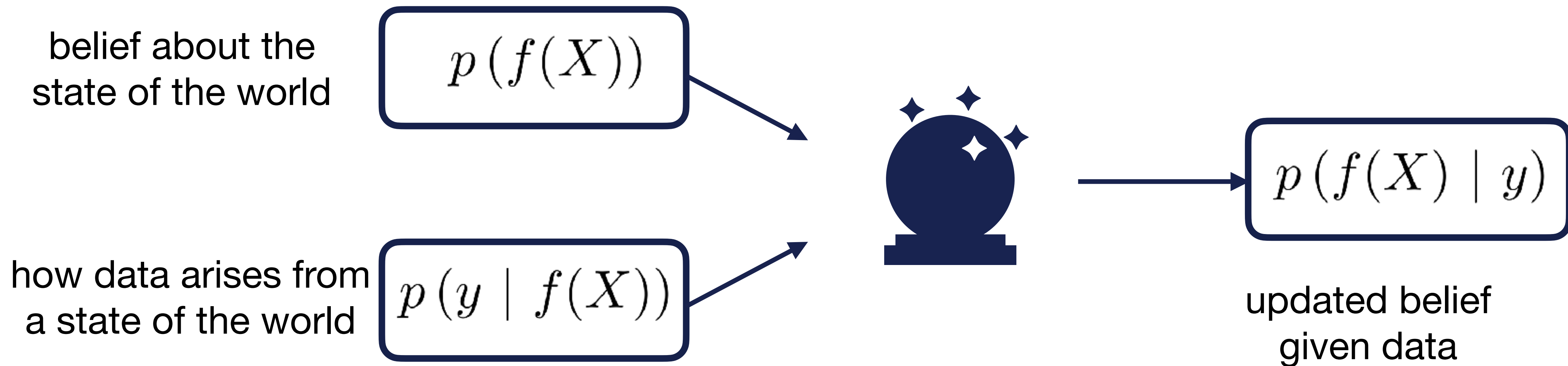
Probabilistic Machine Learning

- Fundamental premise: *treat **all** unknown quantities as random variables.* Bayes Rule does the rest



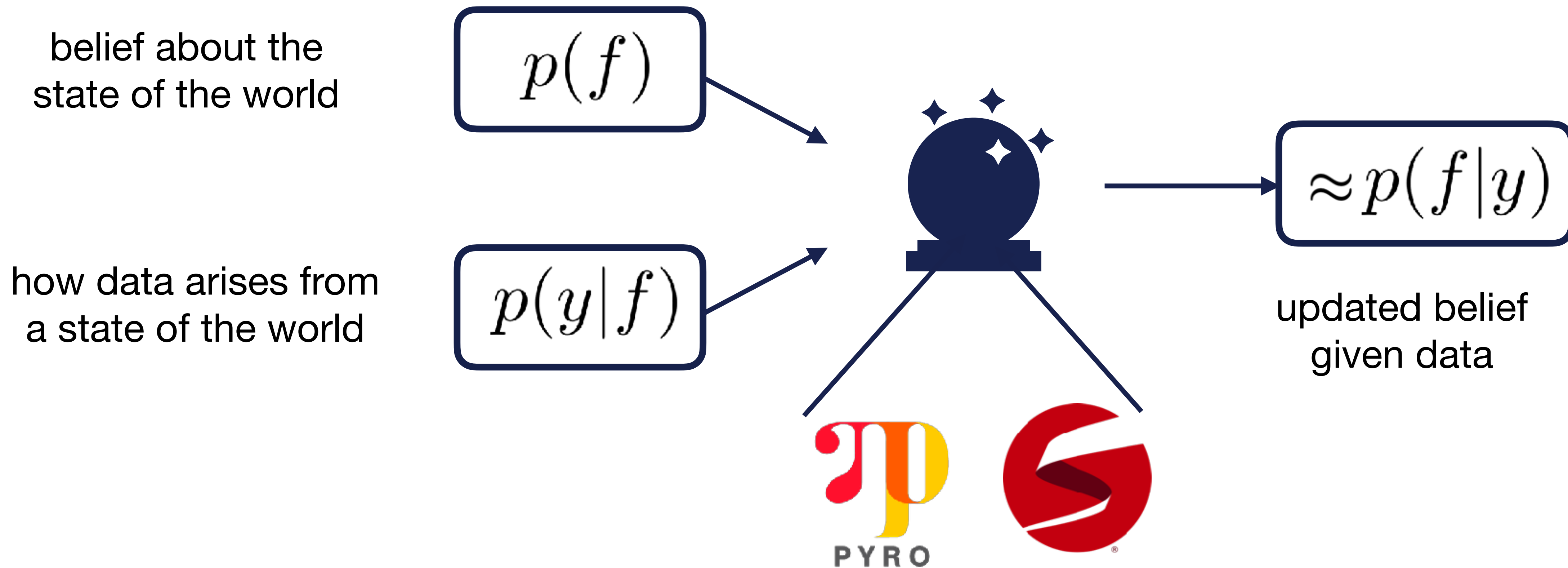
Probabilistic Machine Learning

- Fundamental premise: *treat **all** unknown quantities as random variables.* Bayes Rule does the rest



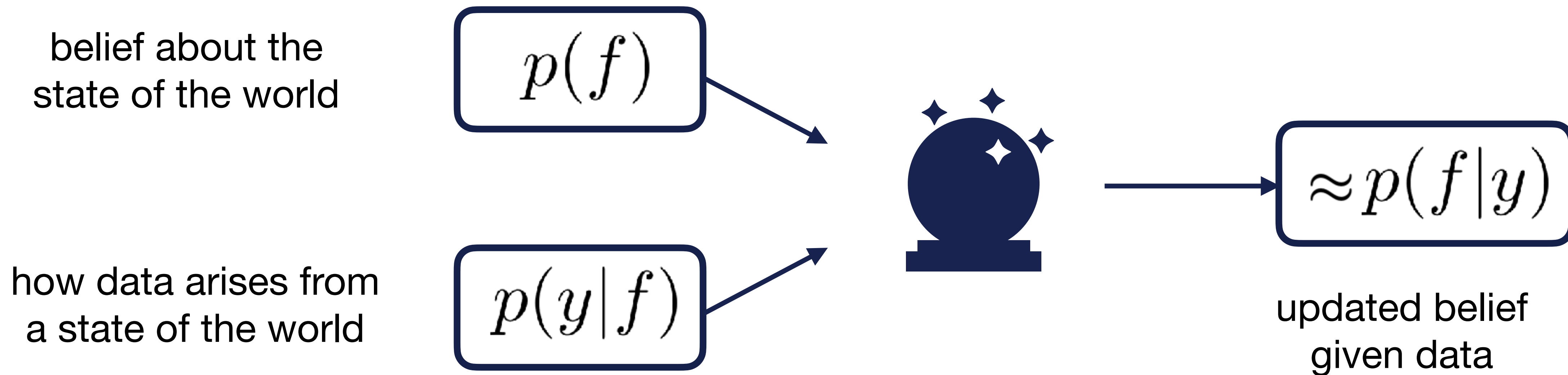
Approximate Inference

- ▶ But this is hard! As soon as you learn the magic of Bayes Rule, you learn that we need to approximately solve it.



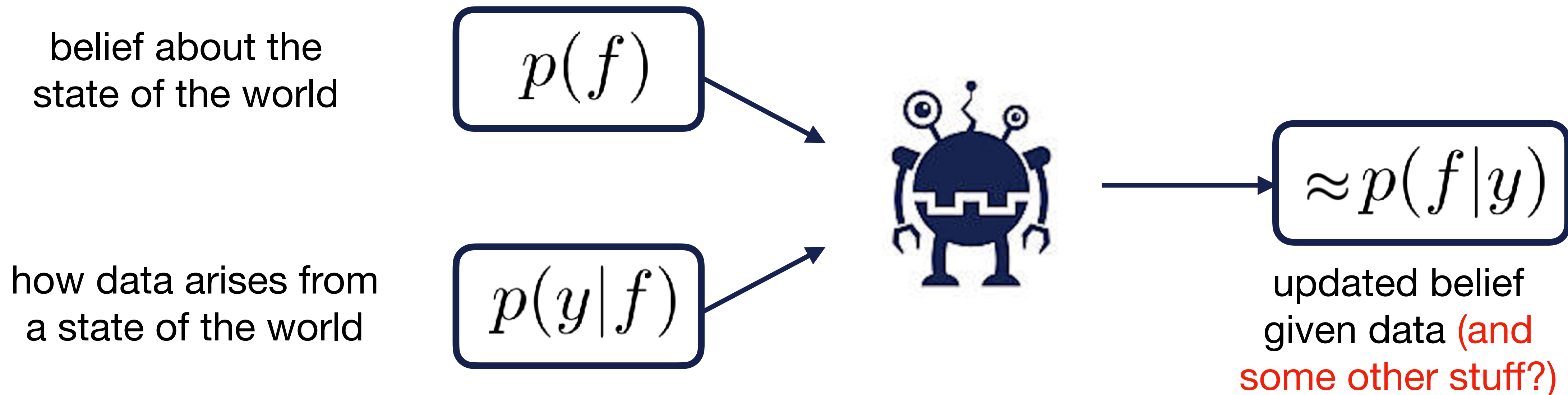
Approximate Inference

- ▶ But this is hard! As soon as you learn the magic of Bayes Rule, you learn that we need to approximately solve it.



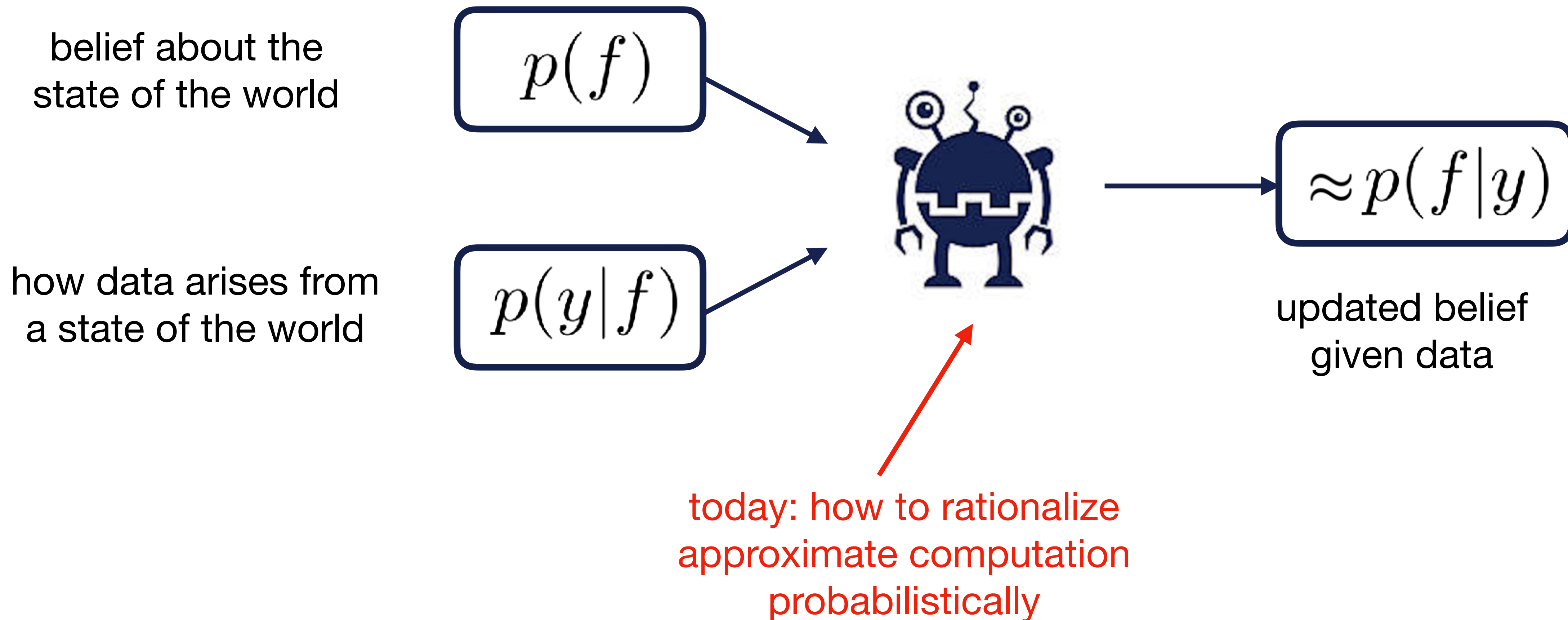
Approximate Inference

- ▶ The crystal ball is actually a computing machine, making assumptions and choices of its own, but we have not accounted for them! All our assumptions, all our unknowns... we were supposed to reason about them probabilistically.



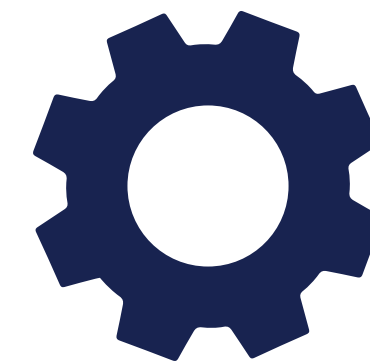
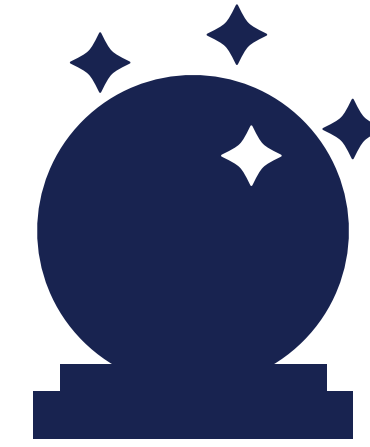
Approximate Inference

- ▶ This statement may seem inherent in the definition of “approximate”, but it is not! The entire point of this talk is to pay off this claim.



Outline

- ▶ The Promise of Probabilistic Machine Learning
- ▶ **Gaussian Process Introduction**
- ▶ Scaling Gaussian Processes, and Implications
- ▶ Approximate Gaussian Process Inference, The Right Way
- ▶ iterGP as Probabilistic Numerics
- ▶ Broader Implications



Gaussian Process Prior

- ▶ Setup:

- ▶ Learn a function $h : \mathcal{X} \rightarrow \mathcal{Y}$

- ▶ Training inputs $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n) \in \mathbb{R}^{n \times d}$

- ▶ Training outputs $\mathbf{y} = (y_1, \dots, y_n)^\top \in \mathbb{R}^n$

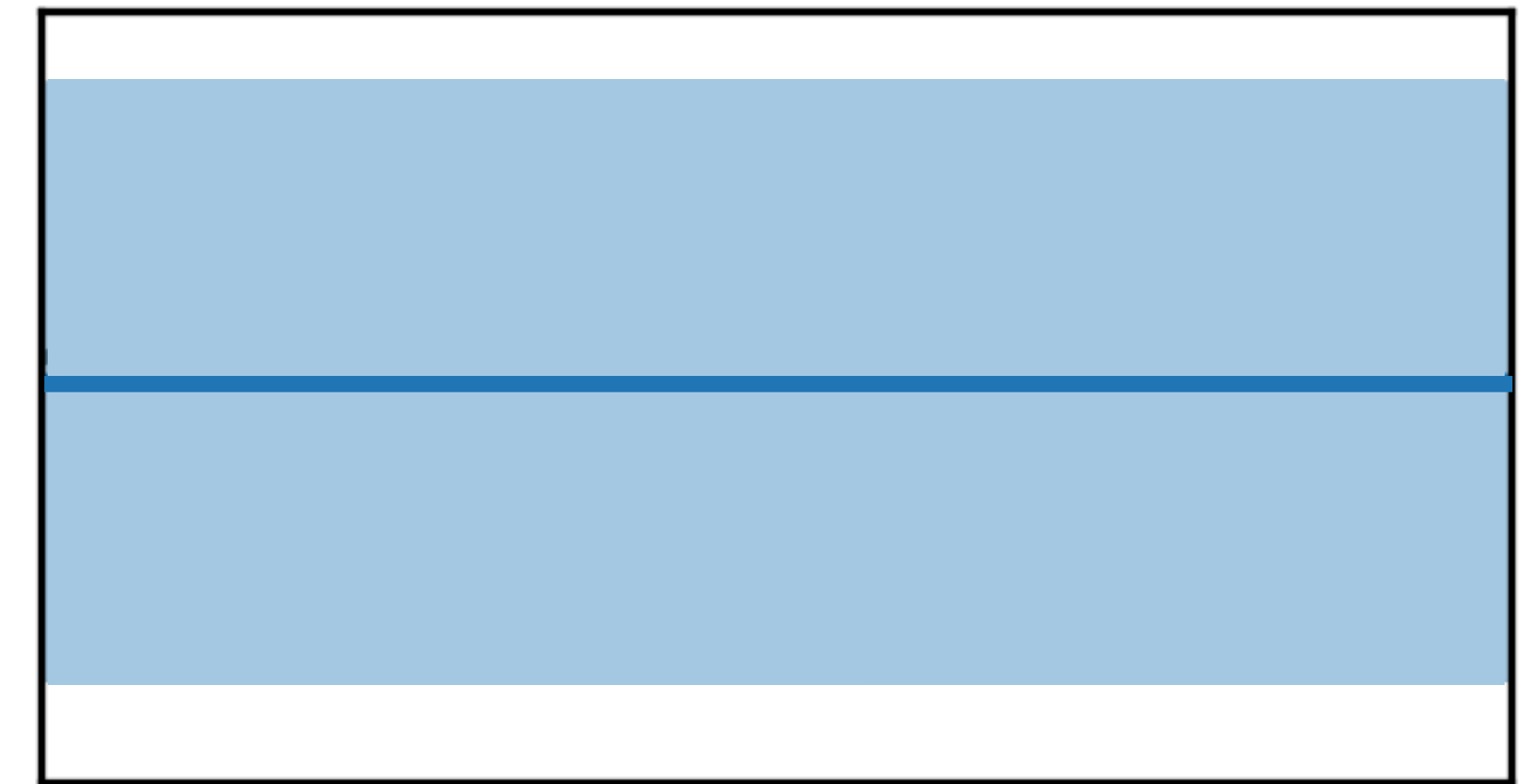
- ▶ Gaussian Process: $f \sim \mathcal{GP}(\mu, k)$

- ▶ Mean function $\mu : \mathbb{R}^d \rightarrow \mathbb{R}$

- ▶ Covariance kernel $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$

- ▶ Joint Gaussian for all \mathbf{X} : $\mathbf{f} = (f(\mathbf{x}_1), \dots, f(\mathbf{x}_n))^\top \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{K})$

$$f \sim \mathcal{GP}(\mu, k)$$



— Latent function

● Data

— Posterior mean

■ Posterior uncertainty

$$\begin{aligned} \mu_j &= \mu(\mathbf{x}_j) \\ \mathbf{K}_{ij} &= k(\mathbf{x}_i, \mathbf{x}_j) \end{aligned}$$

Gaussian Process Posterior

► With likelihood: $\mathbf{y} \mid \mathbf{f} \sim \mathcal{N}(\mathbf{f}, \sigma^2 \mathbf{I})$

► And test inputs: \mathbf{X}_\diamond

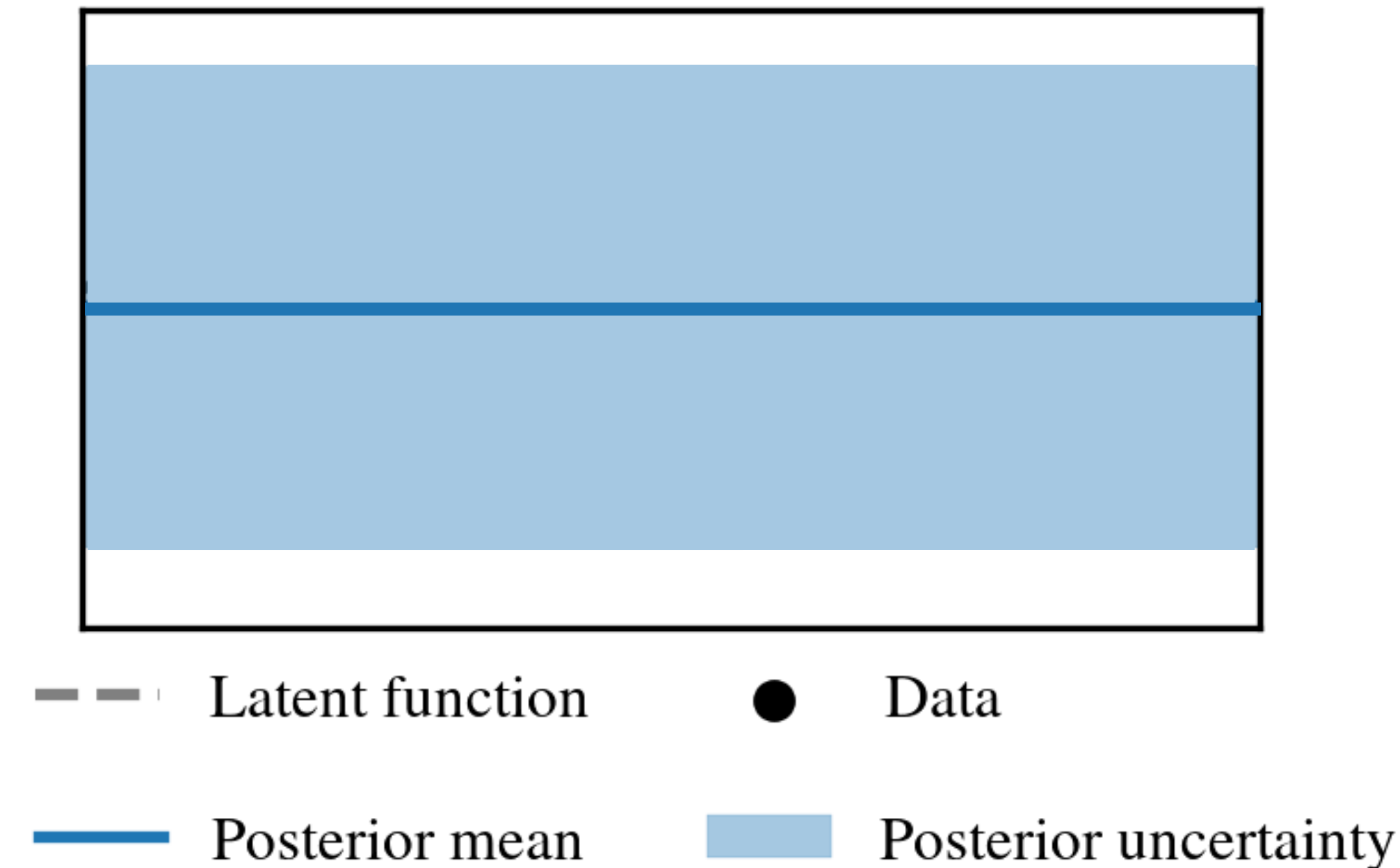
► Induces the posterior:

$$\mathbf{f}_\diamond \sim \mathcal{N}(\mu_*(\mathbf{X}_\diamond), k_*(\mathbf{X}_\diamond, \mathbf{X}_\diamond))$$

► where: $\mu_*(\cdot) = \mu(\cdot) + k(\cdot, \mathbf{X}) \hat{\mathbf{K}}^{-1}(\mathbf{y} - \boldsymbol{\mu})$

$$k_*(\cdot, \cdot) = k(\cdot, \cdot) - k(\cdot, \mathbf{X}) \hat{\mathbf{K}}^{-1} k(\mathbf{X}, \cdot)$$

$$\hat{\mathbf{K}} := \mathbf{K} + \sigma^2 \mathbf{I} \in \mathbb{R}^{n \times n}$$



Gaussian Process Posterior

► With likelihood: $\mathbf{y} \mid \mathbf{f} \sim \mathcal{N}(\mathbf{f}, \sigma^2 \mathbf{I})$

► And test inputs: \mathbf{X}_\diamond

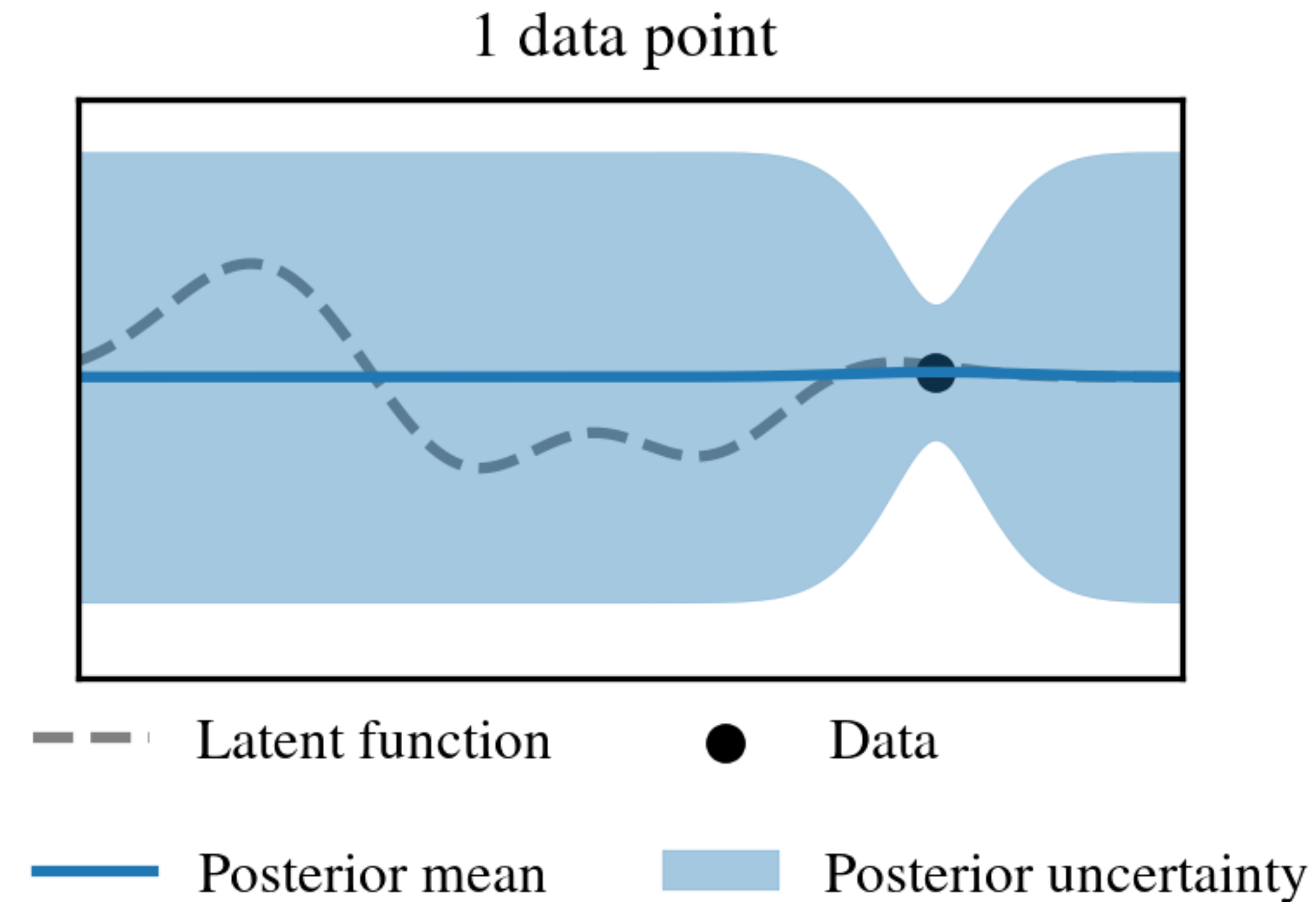
► Induces the posterior:

$$\mathbf{f}_\diamond \sim \mathcal{N}(\mu_*(\mathbf{X}_\diamond), k_*(\mathbf{X}_\diamond, \mathbf{X}_\diamond))$$

► where: $\mu_*(\cdot) = \mu(\cdot) + k(\cdot, \mathbf{X}) \hat{\mathbf{K}}^{-1}(\mathbf{y} - \boldsymbol{\mu})$

$$k_*(\cdot, \cdot) = k(\cdot, \cdot) - k(\cdot, \mathbf{X}) \hat{\mathbf{K}}^{-1} k(\mathbf{X}, \cdot)$$

$$\hat{\mathbf{K}} := \mathbf{K} + \sigma^2 \mathbf{I} \in \mathbb{R}^{n \times n}$$



Gaussian Process Posterior

► With likelihood: $\mathbf{y} \mid \mathbf{f} \sim \mathcal{N}(\mathbf{f}, \sigma^2 \mathbf{I})$

► And test inputs: \mathbf{X}_\diamond

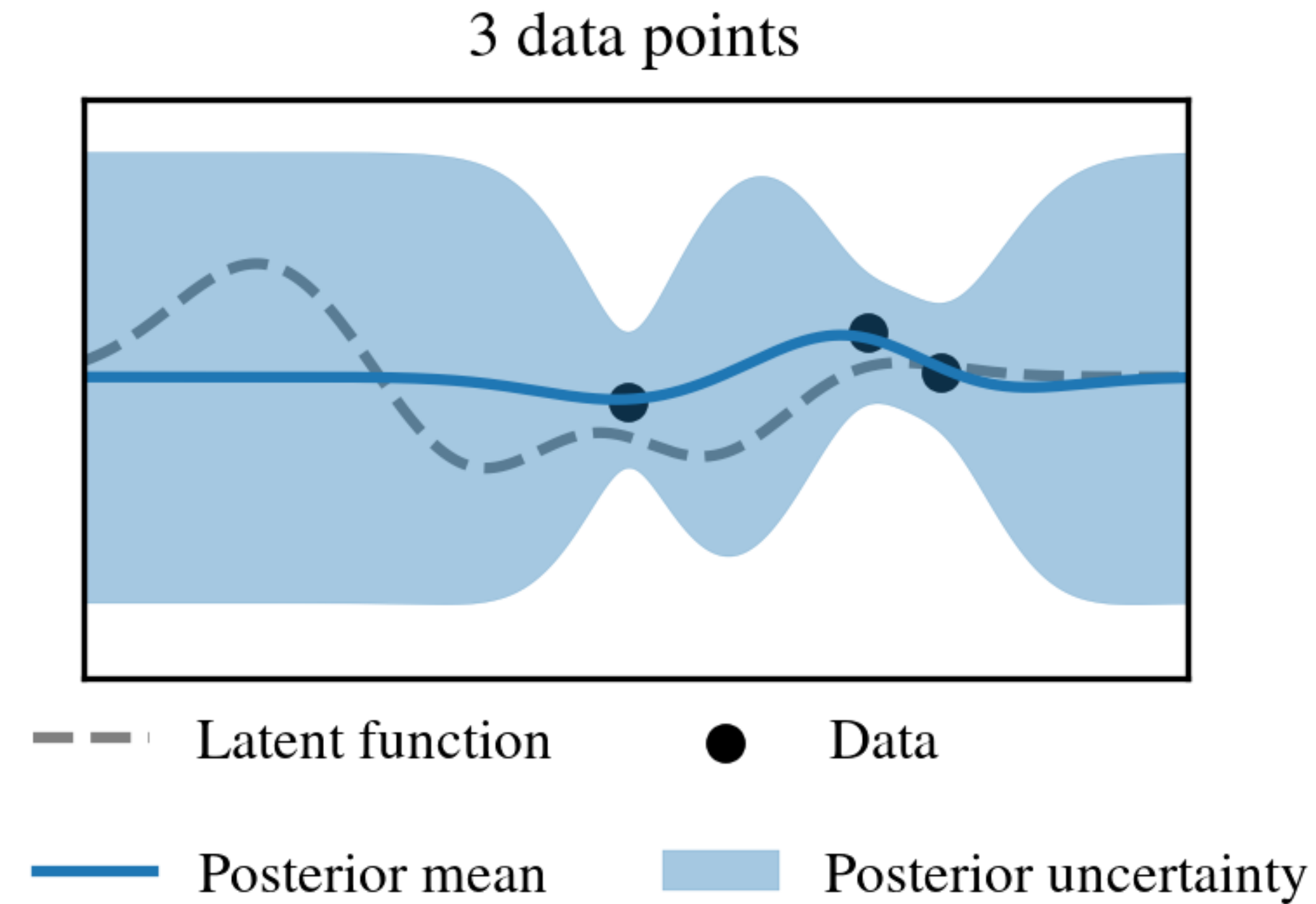
► Induces the posterior:

$$\mathbf{f}_\diamond \sim \mathcal{N}(\mu_*(\mathbf{X}_\diamond), k_*(\mathbf{X}_\diamond, \mathbf{X}_\diamond))$$

► where: $\mu_*(\cdot) = \mu(\cdot) + k(\cdot, \mathbf{X}) \hat{\mathbf{K}}^{-1}(\mathbf{y} - \boldsymbol{\mu})$

$$k_*(\cdot, \cdot) = k(\cdot, \cdot) - k(\cdot, \mathbf{X}) \hat{\mathbf{K}}^{-1} k(\mathbf{X}, \cdot)$$

$$\hat{\mathbf{K}} := \mathbf{K} + \sigma^2 \mathbf{I} \in \mathbb{R}^{n \times n}$$



Gaussian Process Posterior

► With likelihood: $\mathbf{y} \mid \mathbf{f} \sim \mathcal{N}(\mathbf{f}, \sigma^2 \mathbf{I})$

► And test inputs: \mathbf{X}_\diamond

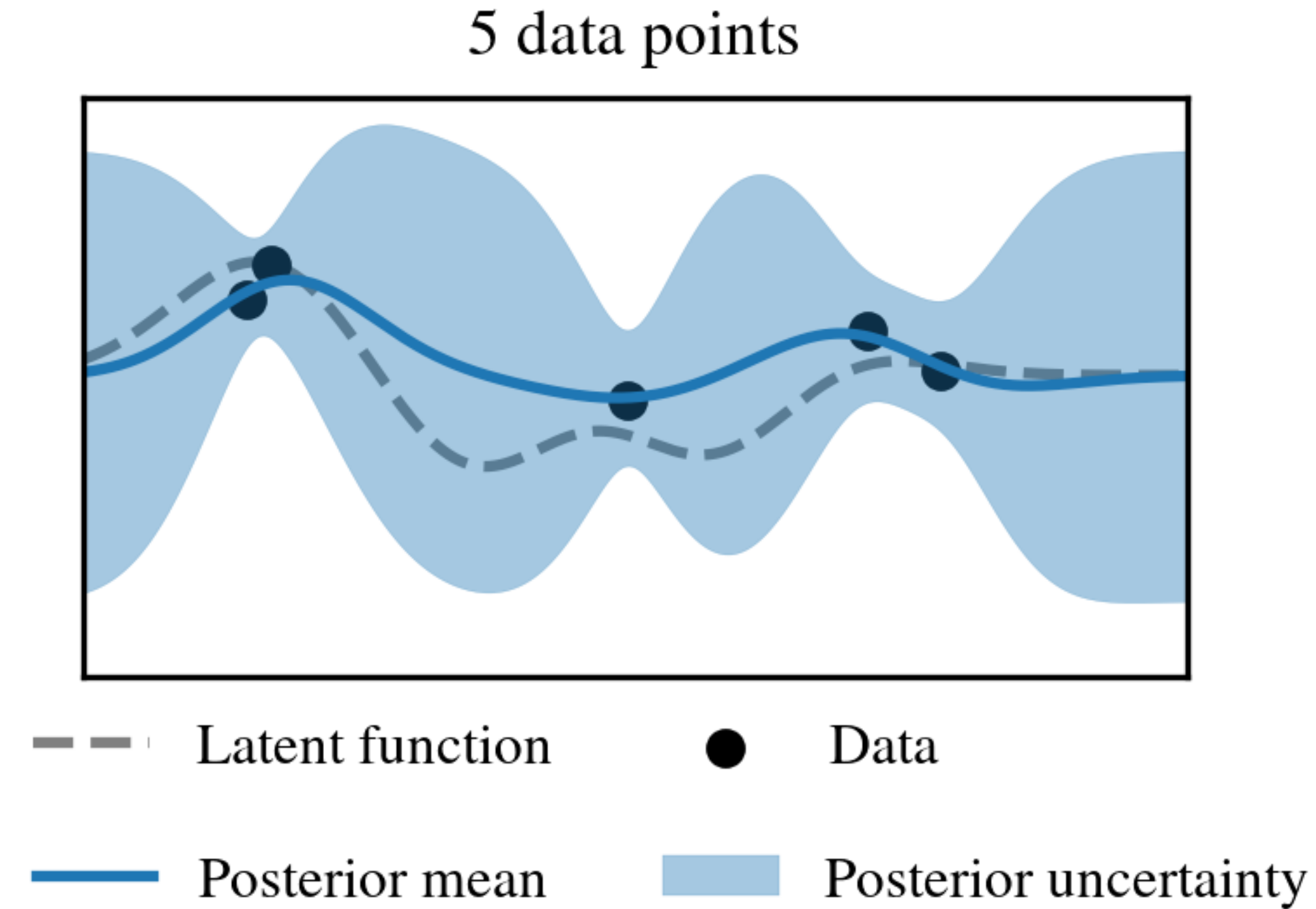
► Induces the posterior:

$$\mathbf{f}_\diamond \sim \mathcal{N}(\mu_*(\mathbf{X}_\diamond), k_*(\mathbf{X}_\diamond, \mathbf{X}_\diamond))$$

► where: $\mu_*(\cdot) = \mu(\cdot) + k(\cdot, \mathbf{X}) \hat{\mathbf{K}}^{-1}(\mathbf{y} - \boldsymbol{\mu})$

$$k_*(\cdot, \cdot) = k(\cdot, \cdot) - k(\cdot, \mathbf{X}) \hat{\mathbf{K}}^{-1} k(\mathbf{X}, \cdot)$$

$$\hat{\mathbf{K}} := \mathbf{K} + \sigma^2 \mathbf{I} \in \mathbb{R}^{n \times n}$$



Gaussian Process Posterior

► With likelihood: $\mathbf{y} \mid \mathbf{f} \sim \mathcal{N}(\mathbf{f}, \sigma^2 \mathbf{I})$

► And test inputs: \mathbf{X}_\diamond

► Induces the posterior:

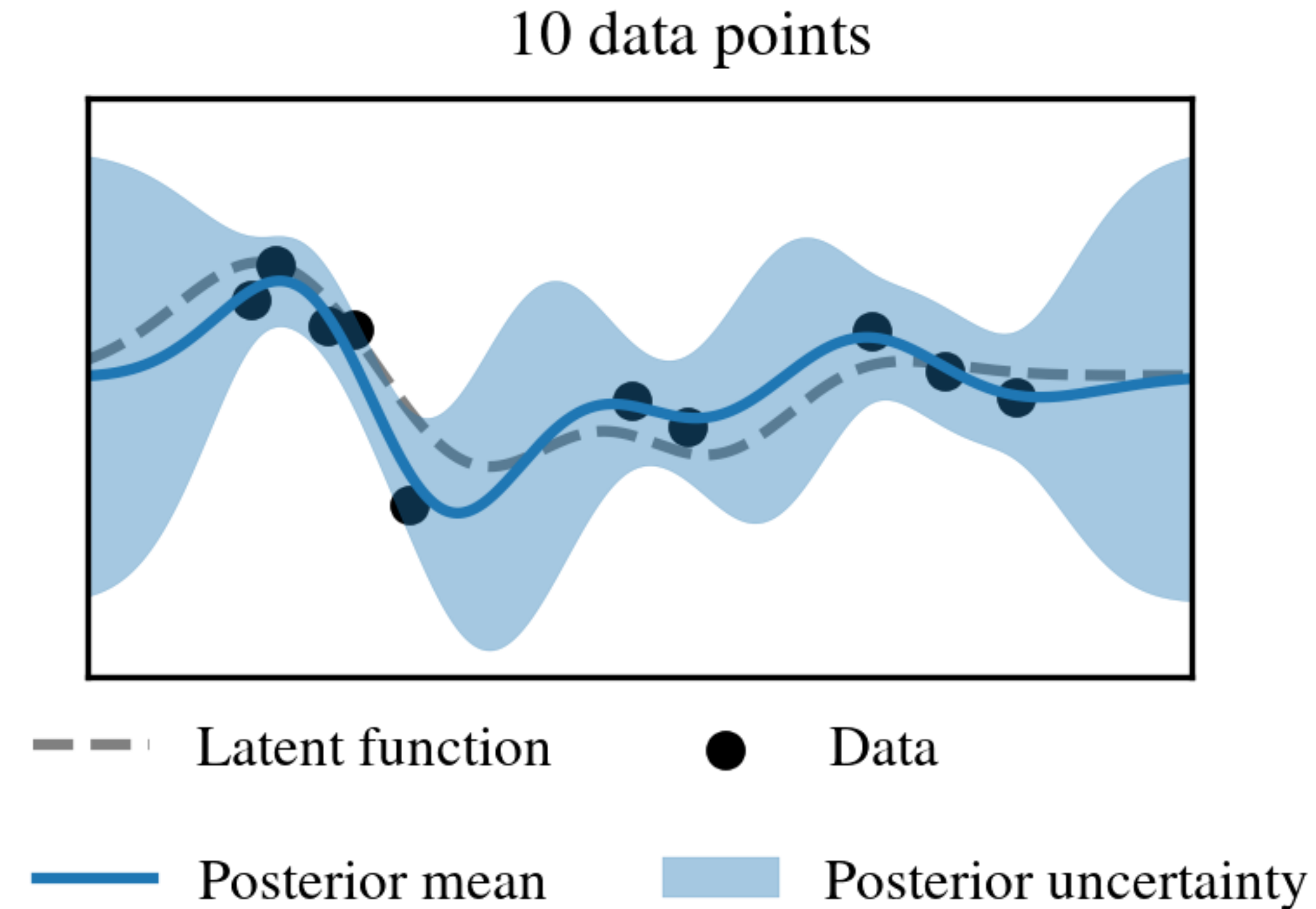
$$\mathbf{f}_\diamond \sim \mathcal{N}(\mu_*(\mathbf{X}_\diamond), k_*(\mathbf{X}_\diamond, \mathbf{X}_\diamond))$$

► where: $\mu_*(\cdot) = \mu(\cdot) + k(\cdot, \mathbf{X}) \hat{\mathbf{K}}^{-1}(\mathbf{y} - \boldsymbol{\mu})$

$$k_*(\cdot, \cdot) = k(\cdot, \cdot) - k(\cdot, \mathbf{X}) \hat{\mathbf{K}}^{-1} k(\mathbf{X}, \cdot)$$

$$\hat{\mathbf{K}} := \mathbf{K} + \sigma^2 \mathbf{I} \in \mathbb{R}^{n \times n}$$

► We will have considerable work to deal with this linear solve



GP Weight Space View

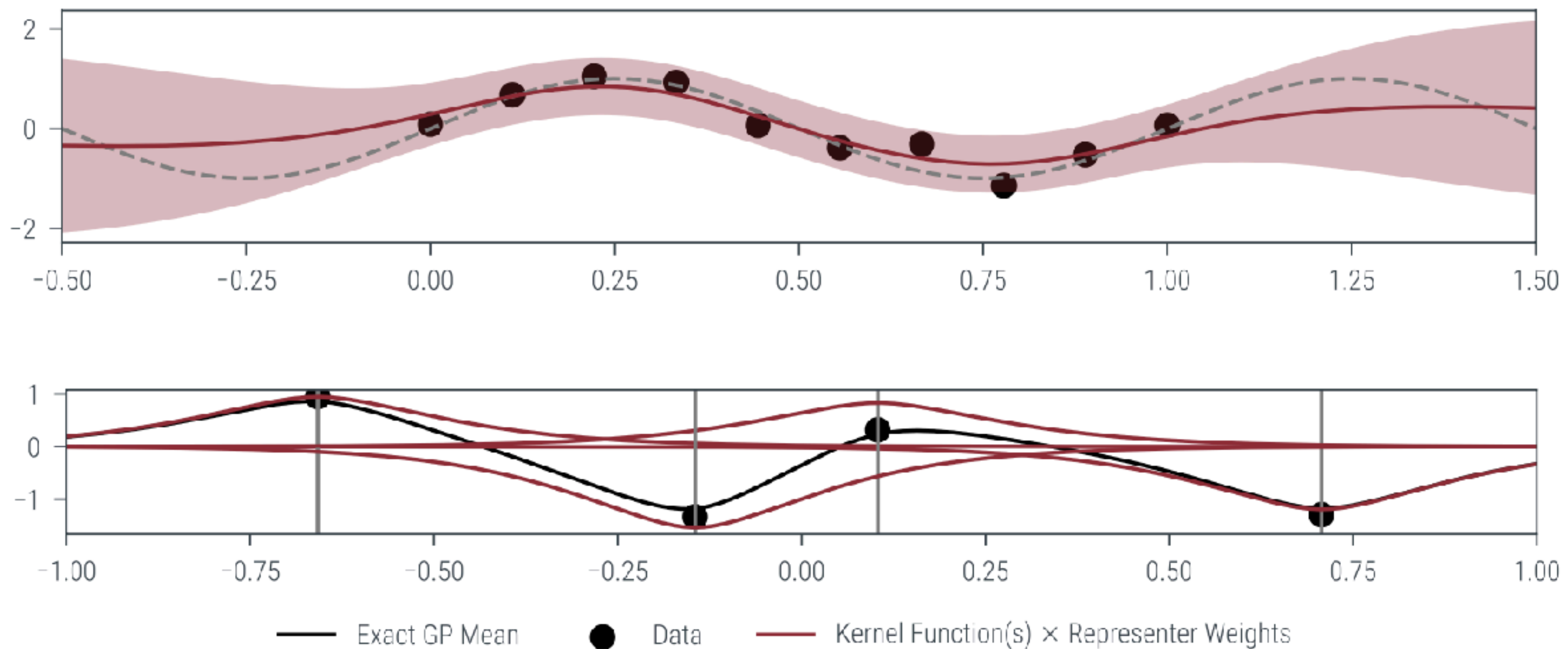
- ▶ Let us also consider the *representer weights* $\mathbf{v}_* = \hat{\mathbf{K}}^{-1}(\mathbf{y} - \boldsymbol{\mu})$

$$\mu_*(\cdot) = \mu(\cdot) + k(\cdot, \mathbf{X}) \overbrace{\hat{\mathbf{K}}^{-1}(\mathbf{y} - \boldsymbol{\mu})}^{\mathbf{v}_*}$$

- ▶ View as kernel basis function regression: place a kernel shaped bump at each input, and the representer weights are the regression coefficients that map that onto the GP posterior mean.

- ▶ We are going to spend time reasoning over these weights

- ▶ (If we think of linear solves probabilistically, then a posterior on our representer weights will provide exactly what we need...)



For Completeness: Learning GP Hyperparameters

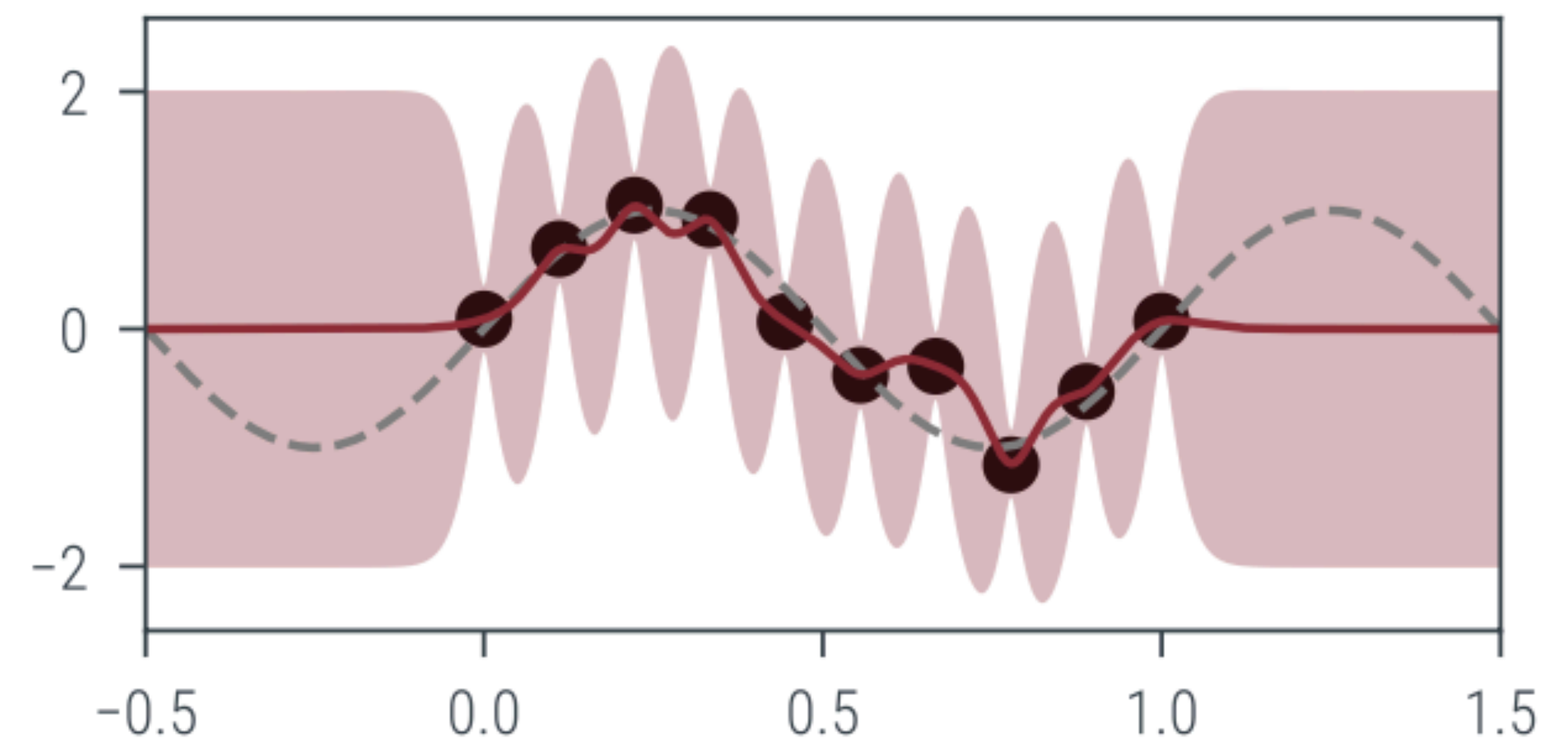
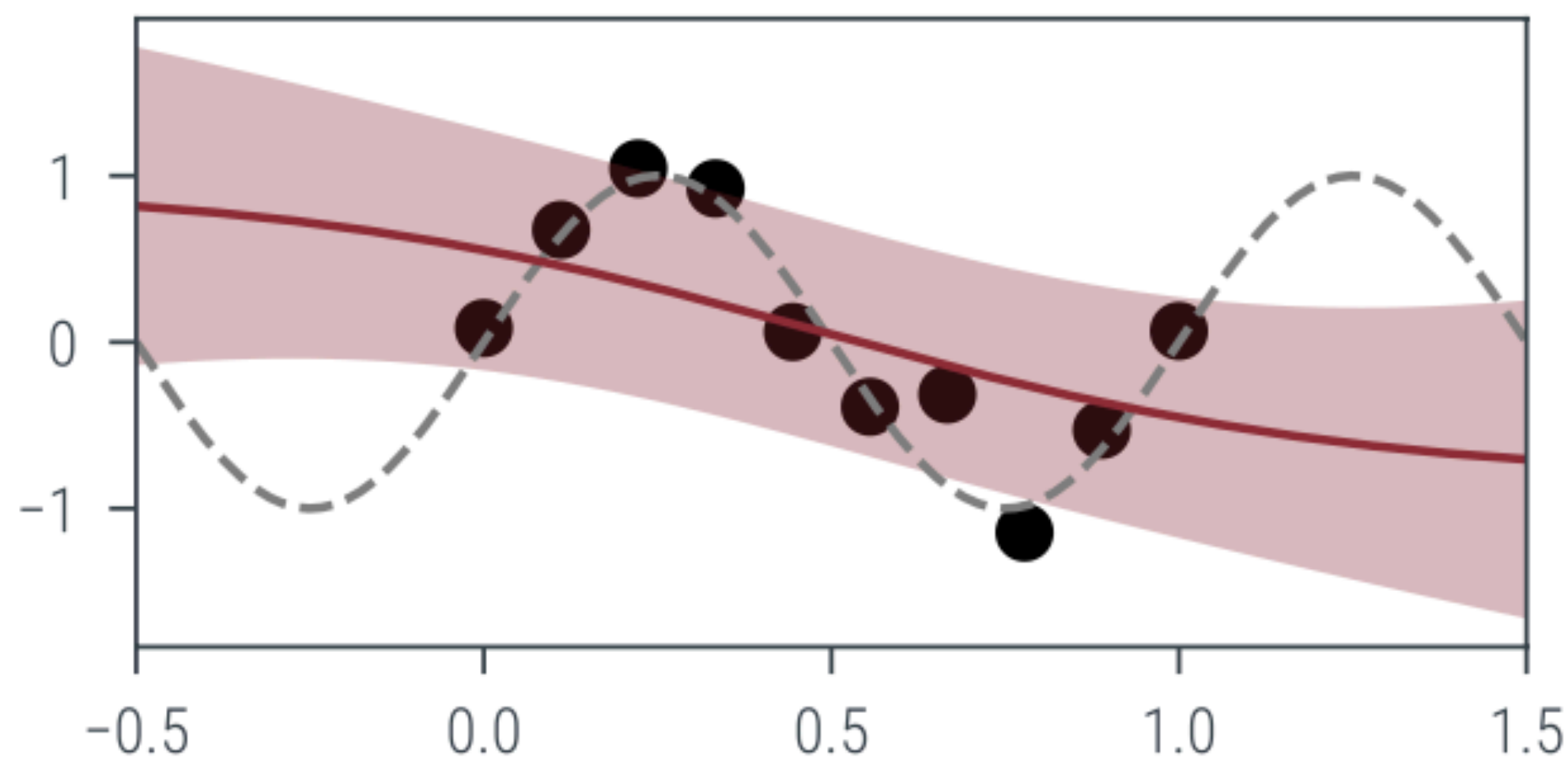
- ▶ Of course, a huge part of GP is learning the hyperparams: model selection or training
- ▶ There's a ton of great work on how to do this, and how to scale it too.

[Gardner et al 2018]
(and many others)

$$\arg \max_{\theta} \mathcal{L}(\theta) = \arg \max_{\theta} -\frac{1}{2}(\mathbf{y} - \boldsymbol{\mu})^{\top} \hat{\mathbf{K}}^{-1}(\mathbf{y} - \boldsymbol{\mu}) - \frac{1}{2} \log |\hat{\mathbf{K}}|$$

“data fit”

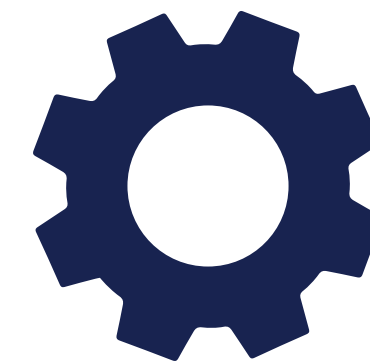
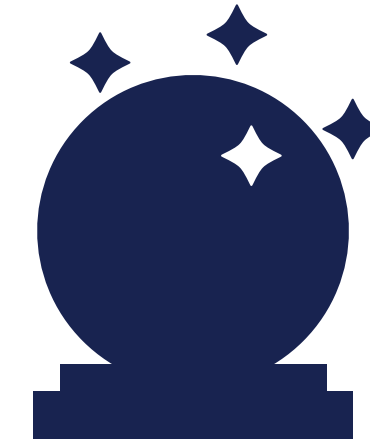
“model complexity”



[Bosch et al 2022]

Outline

- ▶ The Promise of Probabilistic Machine Learning
- ▶ Gaussian Process Introduction
- ▶ **Scaling Gaussian Processes, and Implications**
- ▶ Approximate Gaussian Process Inference, The Right Way
- ▶ iterGP as Probabilistic Numerics
- ▶ Broader Implications



Enter Approximate Inference

- ▶ The core GP object (as far as computation is concerned): $\hat{\mathbf{K}} := \mathbf{K} + \sigma^2 \mathbf{I} \in \mathbb{R}^{n \times n}$
- ▶ It appears over and over:

$$\mu_*(\cdot) = \mu(\cdot) + k(\cdot, \mathbf{X}) \hat{\mathbf{K}}^{-1} (\mathbf{y} - \boldsymbol{\mu})$$

$$k_*(\cdot, \cdot) = k(\cdot, \cdot) - k(\cdot, \mathbf{X}) \hat{\mathbf{K}}^{-1} k(\mathbf{X}, \cdot)$$

$$\arg \max_{\boldsymbol{\theta}} -\frac{1}{2} (\mathbf{y} - \boldsymbol{\mu})^\top \hat{\mathbf{K}}^{-1} (\mathbf{y} - \boldsymbol{\mu}) - \frac{1}{2} \log |\hat{\mathbf{K}}|$$

- ▶ Generally these computations are all cubic in time (and quadratic in storage)
 - ▶ Much literature thus to deal with this cubic scaling
 - ▶ (And note that in cases of special structure, notably $d=1,2,3$, this cost can be much reduced)

[Karvonen and Sarkka 2016 IEEE]

[Loper et al 2021 JMLR]

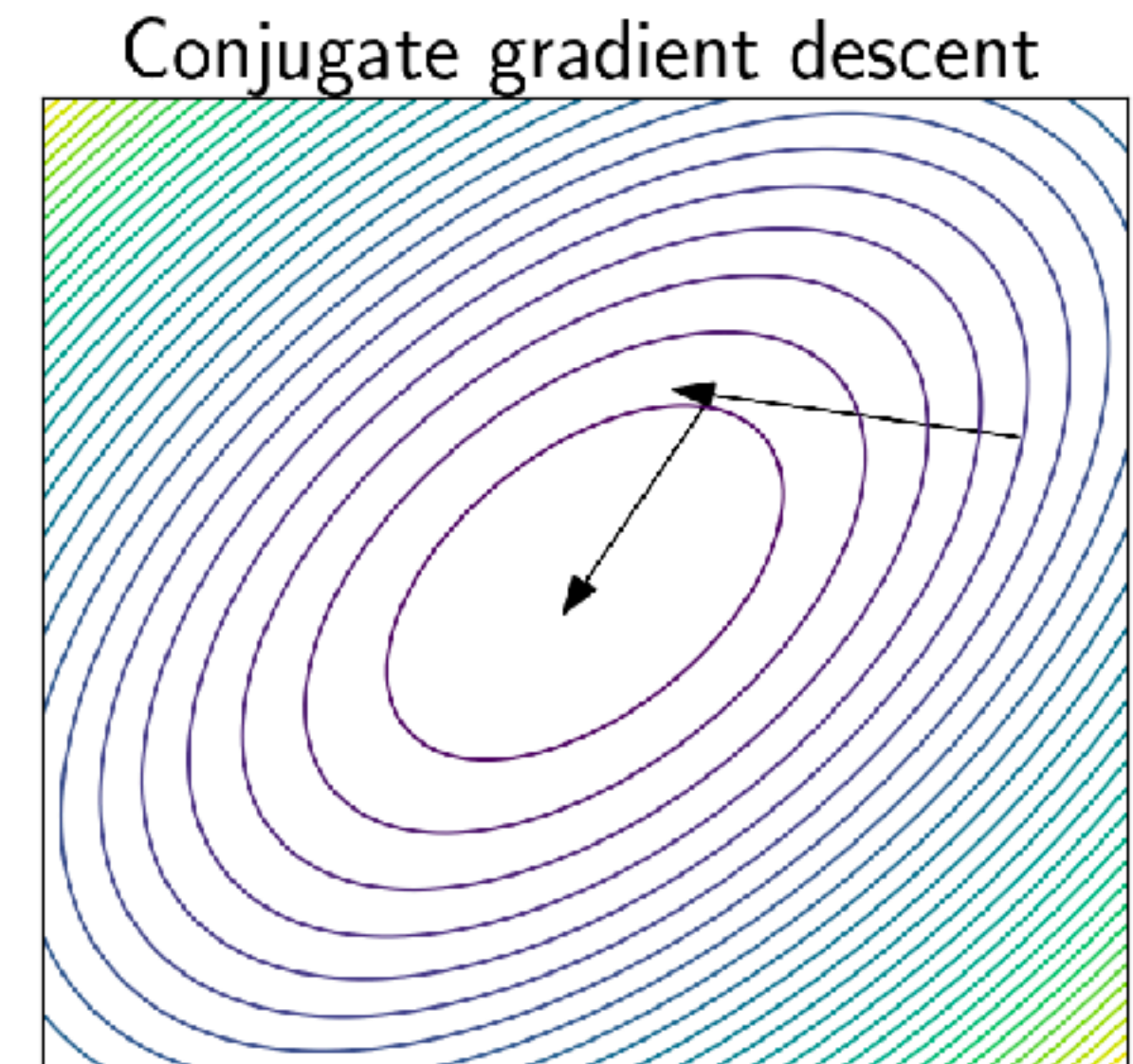
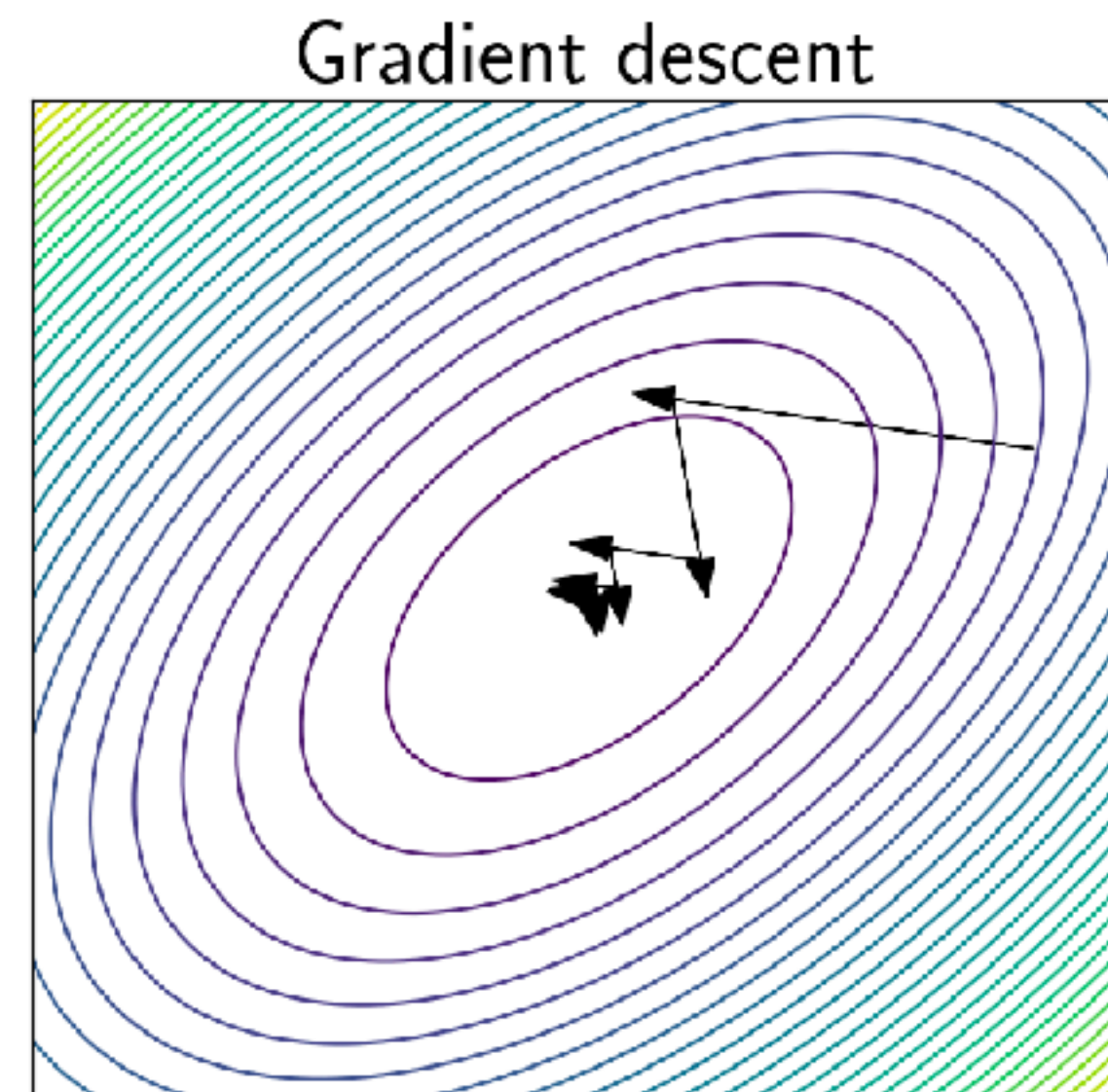
[Greengard et al 2022 arXiv]

Fast Inference 1: Conjugate Gradients

- ▶ View linear solves as an optimization:

$$\min_x \frac{1}{2} \mathbf{x}^\top \hat{\mathbf{K}} \mathbf{x} + \mathbf{x}^\top (\mathbf{y} - \boldsymbol{\mu})$$

- ▶ Gradient descent is slow, so instead take *conjugate* gradient steps, namely steps in $\hat{\mathbf{K}}$ norm.
- ▶ Guaranteed to converge in n steps, but in practice converges to high precision much faster
- ▶ As with many (all) other fast linear solvers, it operates only with forward multiplies of form $\hat{\mathbf{K}} \mathbf{x}$



[Andy Jones]

Fast Inference 1: Conjugate Gradients

- ▶ View linear solves as an optimization:

$$\min_x \frac{1}{2} \mathbf{x}^\top \hat{\mathbf{K}} \mathbf{x} + \mathbf{x}^\top (\mathbf{y} - \boldsymbol{\mu})$$

- ▶ Gradient descent is slow, so instead take *conjugate* gradient steps, namely steps in $\hat{\mathbf{K}}$ norm (via a clever recursion).
- ▶ Guaranteed to converge in n steps, but in practice converges to high precision much faster
- ▶ As with many (all) other fast linear solvers, it operates only with forward multiplies of form $\hat{\mathbf{K}} \mathbf{x}$

Algorithm S3: Preconditioned Conjugate Gradient Method [38]

Input: kernel matrix $\hat{\mathbf{K}}$, labels \mathbf{y} , prior mean $\boldsymbol{\mu}$, preconditioner $\hat{\mathbf{P}}$

Output: representer weights $\mathbf{v}_i \approx \hat{\mathbf{K}}^{-1}(\mathbf{y} - \boldsymbol{\mu})$

```
1 procedure CG( $\hat{\mathbf{K}}, \mathbf{y} - \boldsymbol{\mu}, \hat{\mathbf{P}}$ )
2    $\mathbf{v}_0 \leftarrow \mathbf{0}$ 
3    $\mathbf{s}_0 \leftarrow \mathbf{0}$ 
4   while  $\|\mathbf{r}_i\|_2 > \max(\delta_{\text{rtol}} \|\mathbf{y}\|_2, \delta_{\text{atol}})$  and  $i < i_{\text{max}}$  do
5      $\mathbf{r}_{i-1} \leftarrow (\mathbf{y} - \boldsymbol{\mu}) - \hat{\mathbf{K}} \mathbf{v}_{i-1}$ 
6      $\mathbf{s}_i \leftarrow \hat{\mathbf{P}}^{-1} \mathbf{r}_{i-1} - \frac{(\hat{\mathbf{P}}^{-1} \mathbf{r}_{i-1})^\top \hat{\mathbf{K}} \mathbf{s}_{i-1}}{\mathbf{s}_{i-1}^\top \hat{\mathbf{K}} \mathbf{s}_{i-1}} \mathbf{s}_{i-1}$ 
7      $\mathbf{v}_i \leftarrow \mathbf{v}_{i-1} + \frac{(\hat{\mathbf{P}}^{-1} \mathbf{r}_{i-1})^\top \mathbf{r}_{i-1}}{\mathbf{s}_i^\top \hat{\mathbf{K}} \mathbf{s}_i} \mathbf{s}_i$ 
8   return  $\mathbf{v}$ 
```

[Wenger et al 2022 NeurIPS]

[Cutajar et al 2016 ICML]

[Cunningham et al 2007 ICML]

Fast Inference 1: Conjugate Gradients

- (For completeness) CG is well used for log likelihood computation also:

Algorithm 1: log-Marginal Likelihood

Input: \mathbf{y} (labels), $\hat{\mathbf{K}}$ (kernel matrix), $\hat{\mathbf{P}}$ (preconditioner), ℓ (# of random STE vectors), m (# of CG iterations)

```

1 procedure LOGMARGLIKELIHOOD( $\mathbf{y}, \hat{\mathbf{K}}, \hat{\mathbf{P}}, \ell, m$ )
2    $\mathbf{u} \leftarrow \text{CG}(\hat{\mathbf{K}}, \mathbf{y}, \hat{\mathbf{P}}, m)$   $\triangleright \approx \hat{\mathbf{K}}^{-1} \mathbf{y}$ 
3    $\tau_{\hat{\mathbf{P}}}^{\log} \leftarrow \log \det(\hat{\mathbf{P}})$ 
4   for  $i = 1, \dots, \ell$  do
5      $\mathbf{z}_i \leftarrow \tilde{\mathbf{z}}_i / \|\tilde{\mathbf{z}}_i\|_2$  for rand. vector  $\tilde{\mathbf{z}}_i$ 
6      $\mathbf{T} \leftarrow \text{CG}(\hat{\mathbf{K}}, \mathbf{z}_i, \hat{\mathbf{P}}, m)$   $\triangleright$  equiv. to LANCZOS
7      $[\mathbf{W}, \boldsymbol{\lambda}] \leftarrow \text{EIGENDECOMP}(\mathbf{T})$   $\triangleright \mathbf{T}$  tridiagonal
8      $\omega_j \leftarrow (\mathbf{e}_1^\top \mathbf{w}_j)^2$  for  $j = 0, \dots, m$   $\triangleright$  quad. weights
9      $\gamma_i \leftarrow \sum_{j=0}^m \omega_j \log(\lambda_j)$   $\triangleright \approx \mathbf{z}_i^\top \boldsymbol{\Delta}_{\log} \mathbf{z}_i$ 
10   $\tau_*^{\log} \leftarrow \tau_{\hat{\mathbf{P}}}^{\log} + \frac{n}{\ell} \sum_{i=1}^{\ell} \gamma_i$   $\triangleright \approx \log \det(\hat{\mathbf{K}})$ 
11  return  $-\frac{1}{2}(\mathbf{y}^\top \mathbf{u} + \tau_*^{\log} + n \log(2\pi))$   $\triangleright \approx \mathcal{L}(\boldsymbol{\theta})$ 
```

Algorithm 2: Derivative of the log-Marginal Likelihood

Input: \mathbf{y} (labels), $\hat{\mathbf{K}}$ (kernel matrix), $\hat{\mathbf{P}}$ (preconditioner), ℓ (# of random STE vectors), m (# of CG iterations), $\frac{\partial \hat{\mathbf{K}}}{\partial \boldsymbol{\theta}} / \frac{\partial \hat{\mathbf{P}}}{\partial \boldsymbol{\theta}}$ (functions for computing kernel / preconditioner derivatives)

```

1 procedure DERIVATIVE( $\mathbf{y}, \hat{\mathbf{K}}, \frac{\partial \hat{\mathbf{K}}}{\partial \boldsymbol{\theta}}, \hat{\mathbf{P}}, \frac{\partial \hat{\mathbf{P}}}{\partial \boldsymbol{\theta}}, \ell, m$ )
2    $\mathbf{u} \leftarrow \text{CG}(\hat{\mathbf{K}}, \mathbf{y}, \hat{\mathbf{P}}, m)$   $\triangleright \approx \hat{\mathbf{K}}^{-1} \mathbf{y}$ 
3    $\tau_{\hat{\mathbf{P}}}^{\text{inv}\partial} \leftarrow \text{tr}(\hat{\mathbf{P}}^{-1} \frac{\partial \hat{\mathbf{P}}}{\partial \boldsymbol{\theta}})$ 
4   for  $i = 1, \dots, \ell$  do
5      $\mathbf{z}_i \leftarrow \tilde{\mathbf{z}}_i / \|\tilde{\mathbf{z}}_i\|_2$  for rand. vector  $\tilde{\mathbf{z}}_i$ 
6      $\mathbf{w}_i \leftarrow \text{CG}(\hat{\mathbf{K}}, \frac{\partial \hat{\mathbf{K}}}{\partial \boldsymbol{\theta}} \mathbf{z}_i, \hat{\mathbf{P}}, m)$   $\triangleright \approx \hat{\mathbf{K}}^{-1} \frac{\partial \hat{\mathbf{K}}}{\partial \boldsymbol{\theta}} \mathbf{z}_i$ 
7      $\tilde{\mathbf{w}}_i \leftarrow \hat{\mathbf{P}}^{-1} \frac{\partial \hat{\mathbf{P}}}{\partial \boldsymbol{\theta}} \mathbf{z}_i$ 
8      $\gamma_i \leftarrow \mathbf{z}_i^\top (\mathbf{w}_i - \tilde{\mathbf{w}}_i)$   $\triangleright \approx \mathbf{z}_i^\top \boldsymbol{\Delta}_{\text{inv}\partial} \mathbf{z}_i$ 
9      $\tau_*^{\text{inv}\partial} \leftarrow \tau_{\hat{\mathbf{P}}}^{\text{inv}\partial} + \frac{n}{\ell} \sum_{i=1}^{\ell} \gamma_i$   $\triangleright \approx \text{tr}(\mathbf{K}^{-1} \frac{\partial \mathbf{K}}{\partial \boldsymbol{\theta}})$ 
10  return  $\frac{1}{2}(\mathbf{u}^\top \frac{\partial \mathbf{K}}{\partial \boldsymbol{\theta}} \mathbf{u} - \tau_*^{\text{inv}\partial})$   $\triangleright \approx \frac{\partial}{\partial \boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta})$ 
```

Fast Inference 2: Inducing Point Methods

- ▶ *Inducing points* have been used under many names and variations
 - ▶ SoR, Nystrom, SVGP, DTC,...
 - ▶ These methods differ in their training and posterior covariance assumptions
 - ▶ Detailing all of them is a few lectures in its own right; here I will lay out the essentials

- ▶ Choose a set of locations $Z \in \mathbb{R}^{m \times d}$

- ▶ Posterior then uses the approximation $\mathbf{K}_{XX} \approx \mathbf{K}_{XZ} \mathbf{K}_{ZZ}^{-1} \mathbf{K}_{ZX}$
 - ▶ Well-chosen inducing points make should be a good low rank approximation

- ▶ While methods differ in posterior covariance, they share a posterior mean [Wild et al 2023]

$$\mu(\cdot) = k(\cdot, Z) \mathbf{K}_{ZZ}^{-1} \mathbf{K}_{ZX} \mathbf{K}_{XZ} (\mathbf{K}_{ZX} (\mathbf{K}_{XZ} \mathbf{K}_{ZZ}^{-1} \mathbf{K}_{ZX} + \sigma^2 I) \mathbf{K}_{XZ})^{-1} \mathbf{K}_{ZX} (\mathbf{y} - \boldsymbol{\mu})$$

- ▶ TL;DR: all inversions take place in the inducing point space:

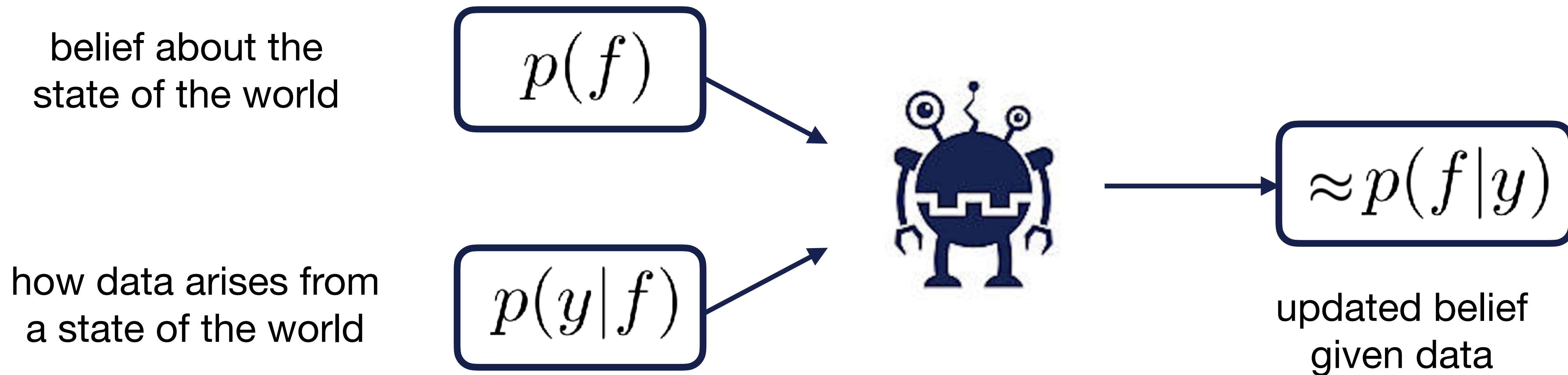
- ▶ Meaningful cost savings
- ▶ Tendency to be overconfident, especially where inducing points are placed.

[Seeger and Williams 2001 NIPS]

[Wu et al 2022 ICML]
[Wu et al 2021 AISTATS]

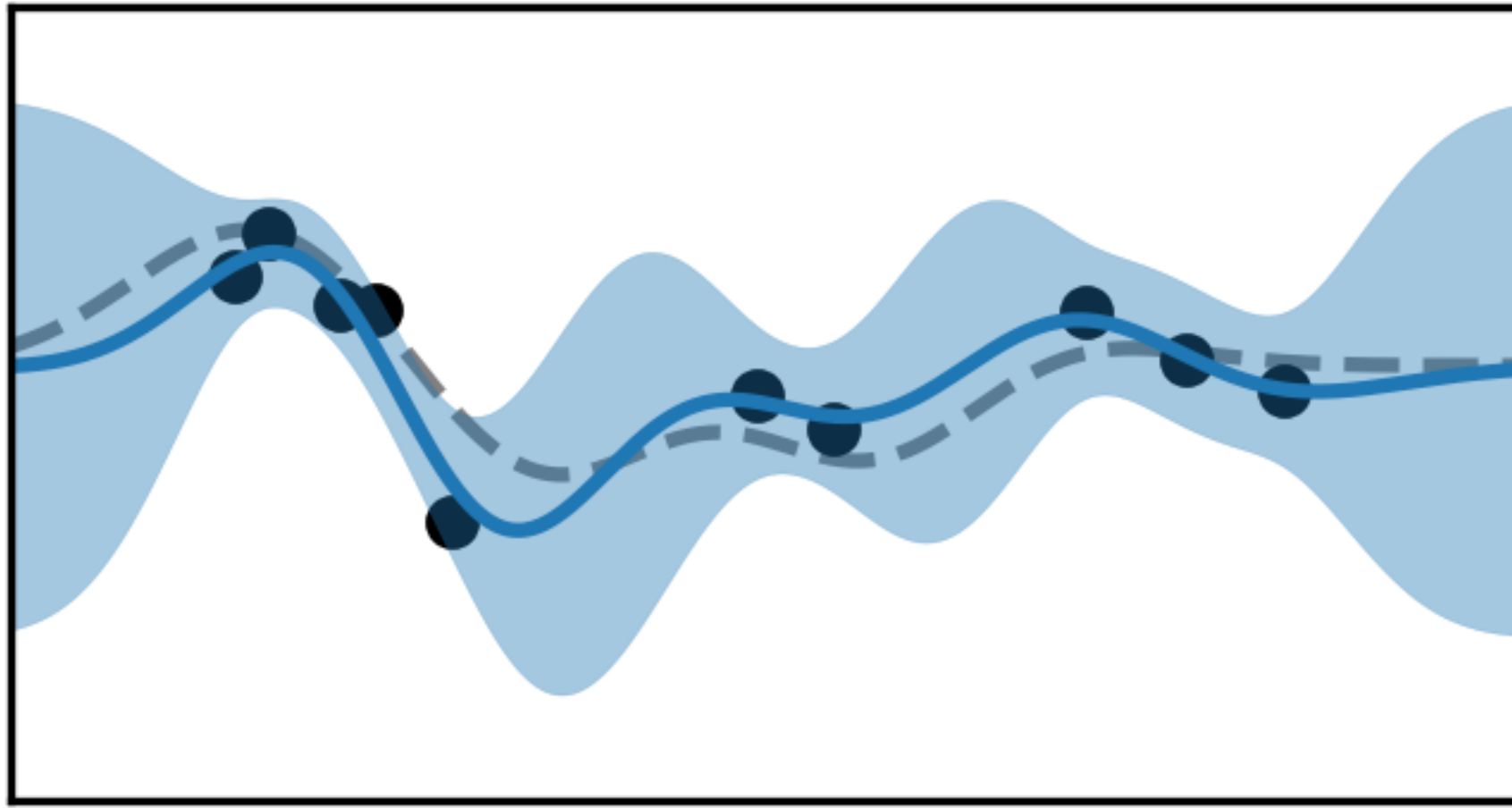
Approximate Inference

- ▶ The crystal ball is actually a computing machine, making assumptions and choices of its own, but we have not accounted for them. All our assumptions, all our unknowns ... we were supposed to reason about them probabilistically!



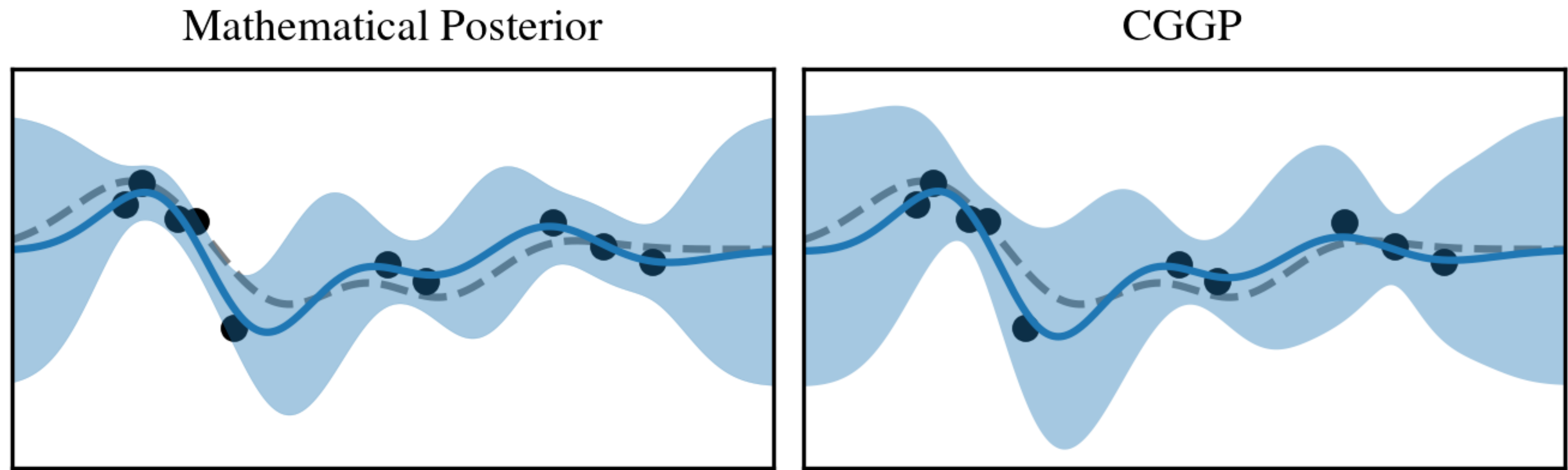
Approximate Inference in Practice

Mathematical Posterior



— — Latent function ● Data — Posterior mean Posterior uncertainty

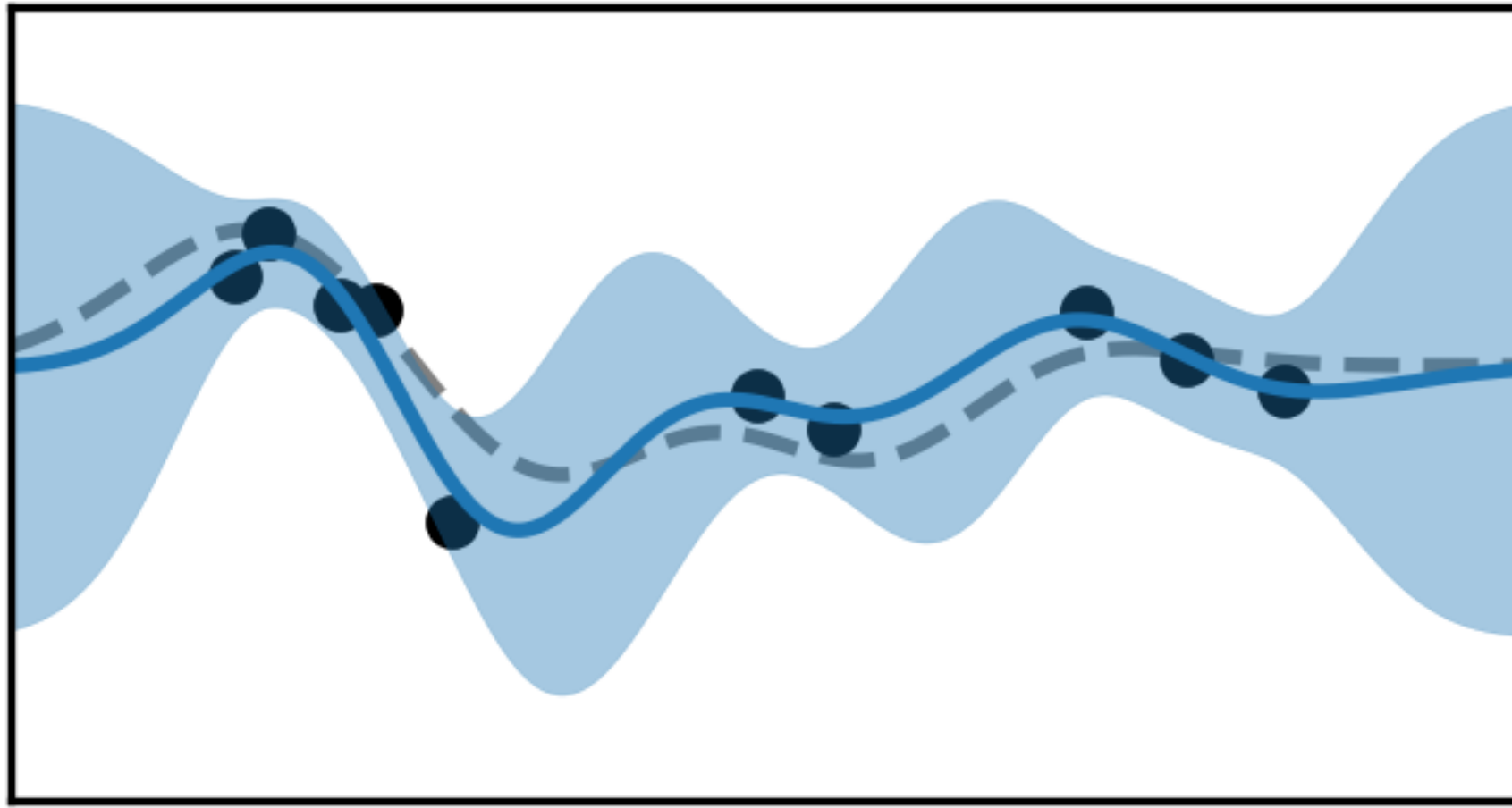
Approximate Inference in Practice



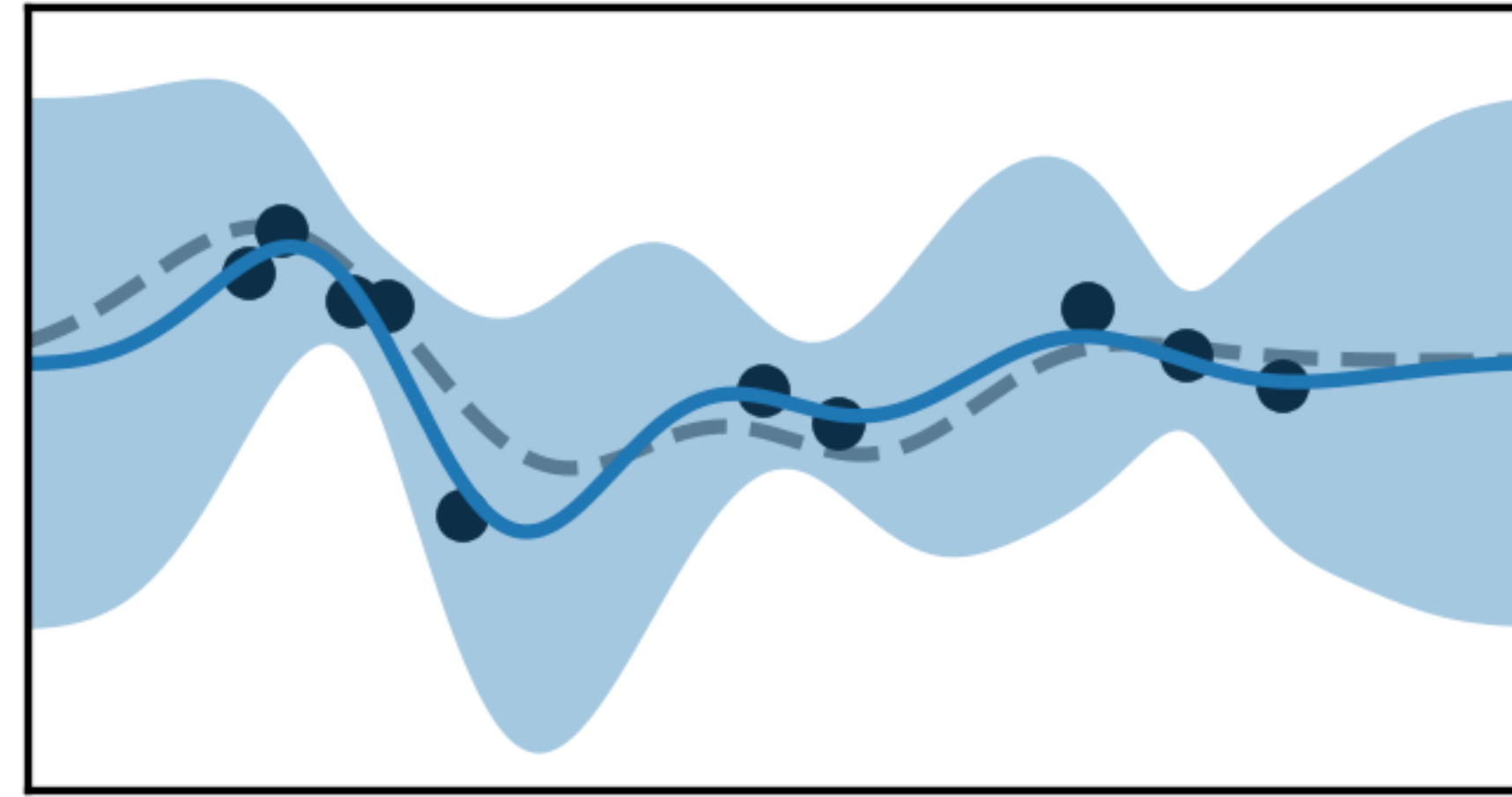
— — Latent function ● Data — Posterior mean Posterior uncertainty

Approximate Inference in Practice

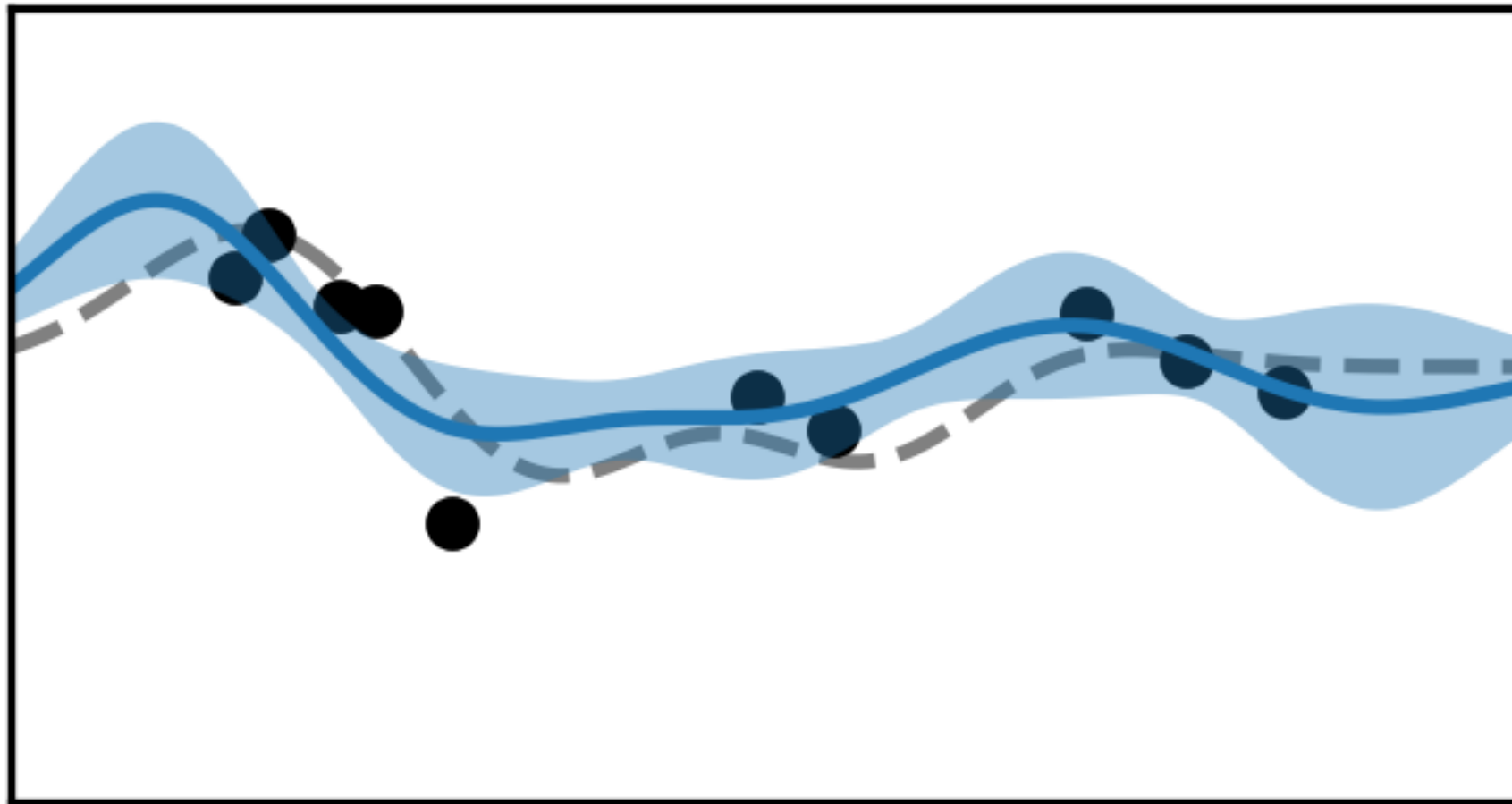
Mathematical Posterior



CGGP



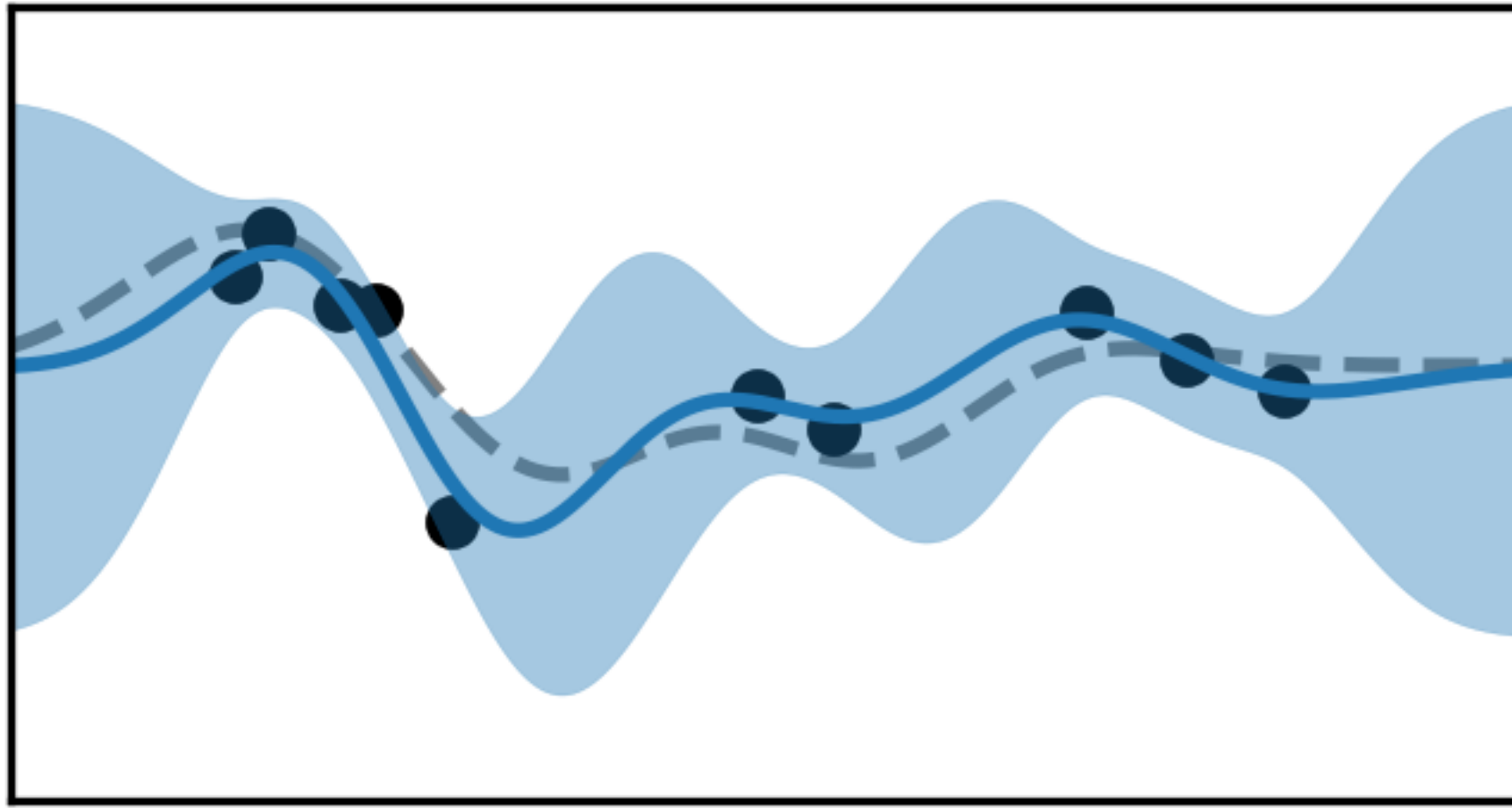
Nyström (SoR)



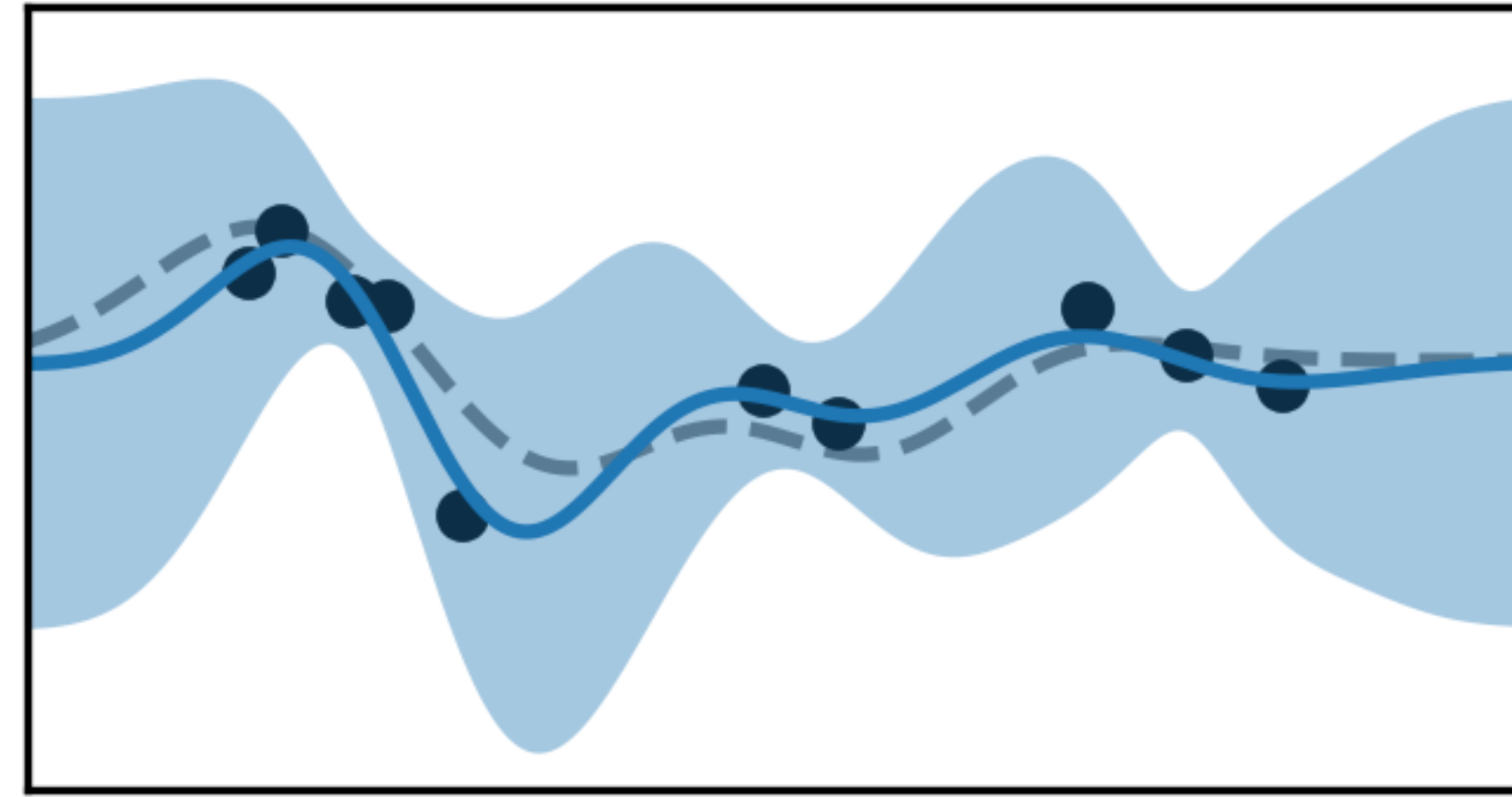
— — Latent function ● Data — Posterior mean Posterior uncertainty

Approximate Inference in Practice

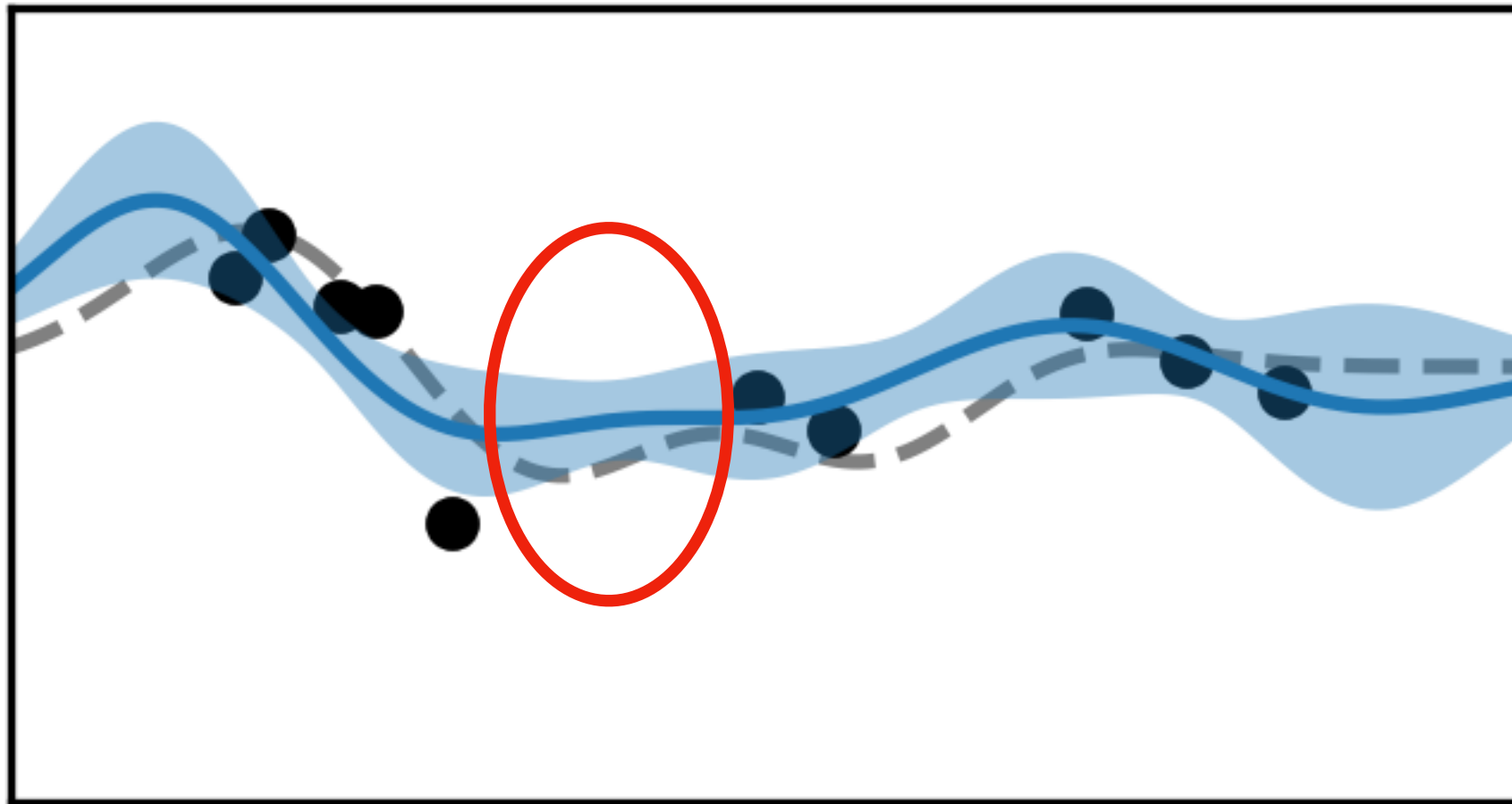
Mathematical Posterior



CGGP



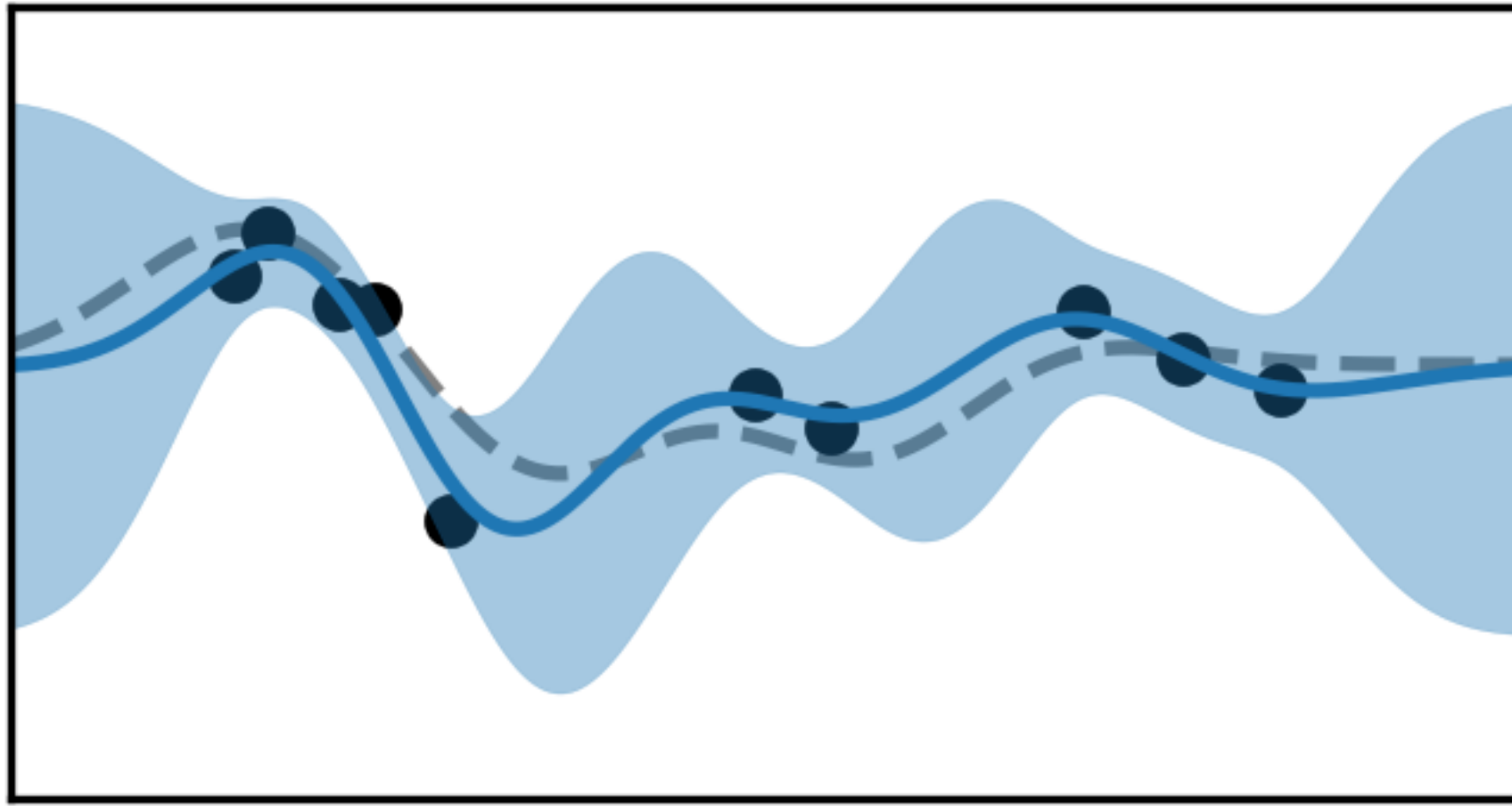
Nyström (SoR)



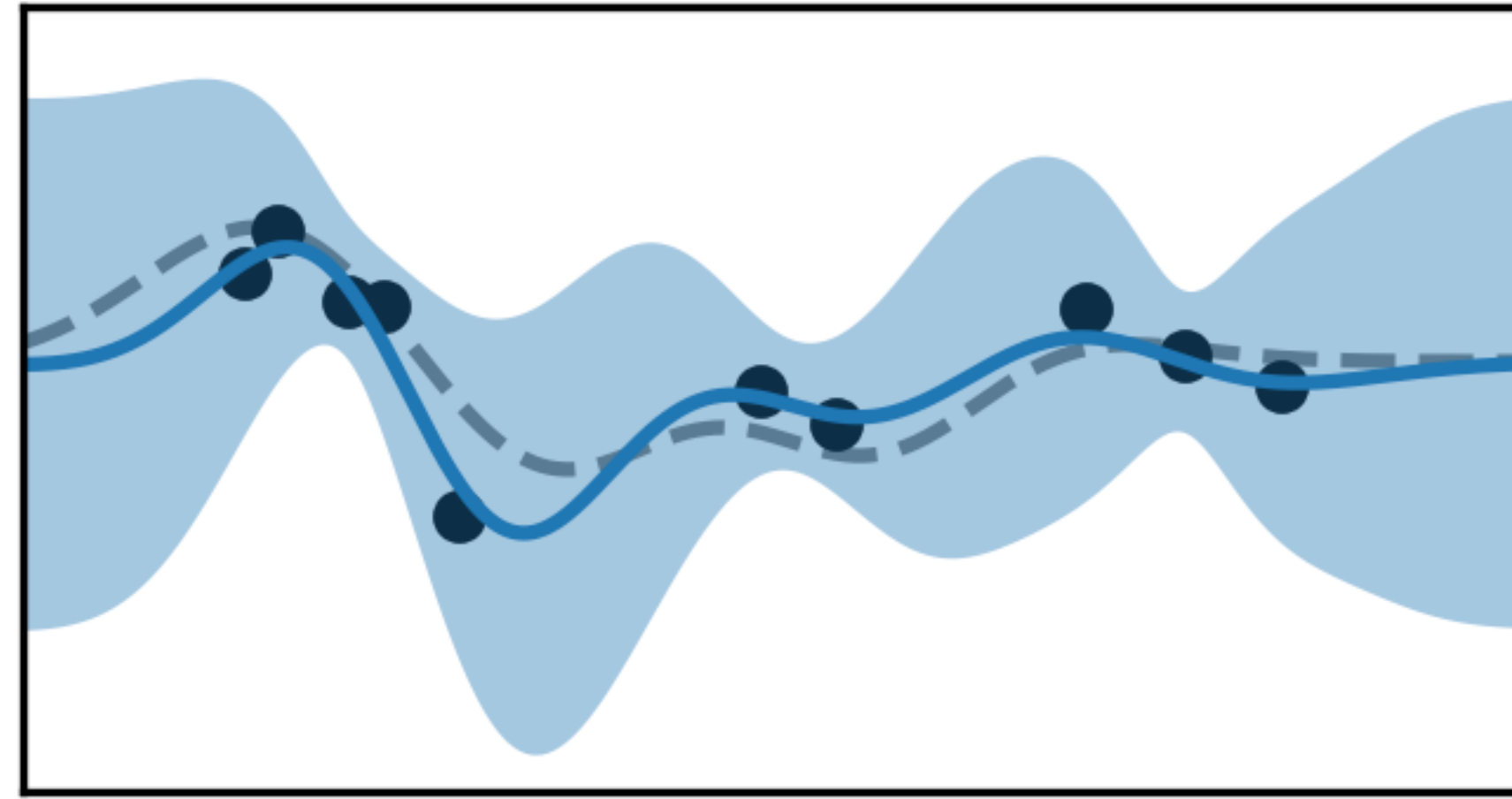
--- Latent function ● Data — Posterior mean Posterior uncertainty

Approximate Inference in Practice

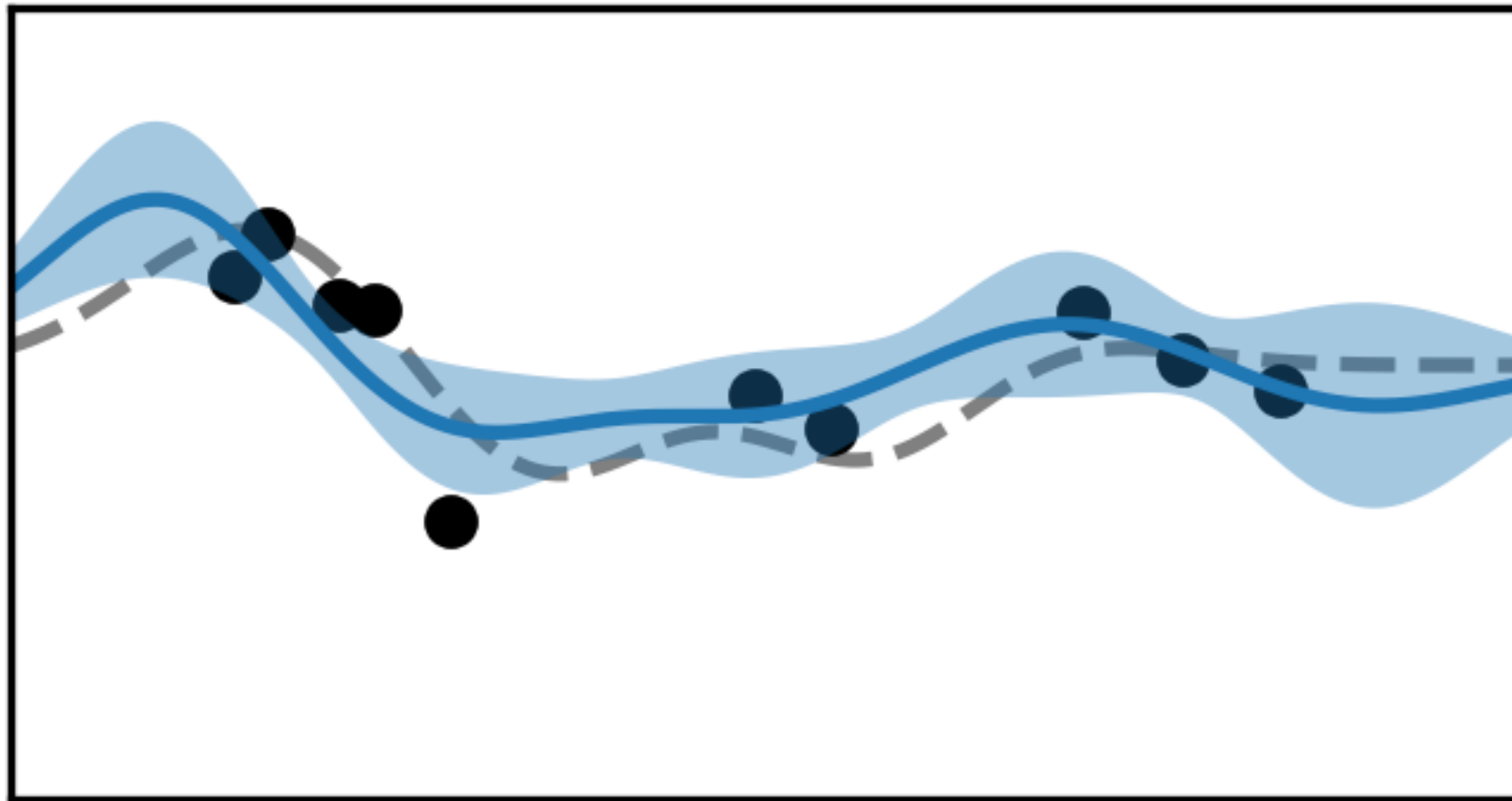
Mathematical Posterior



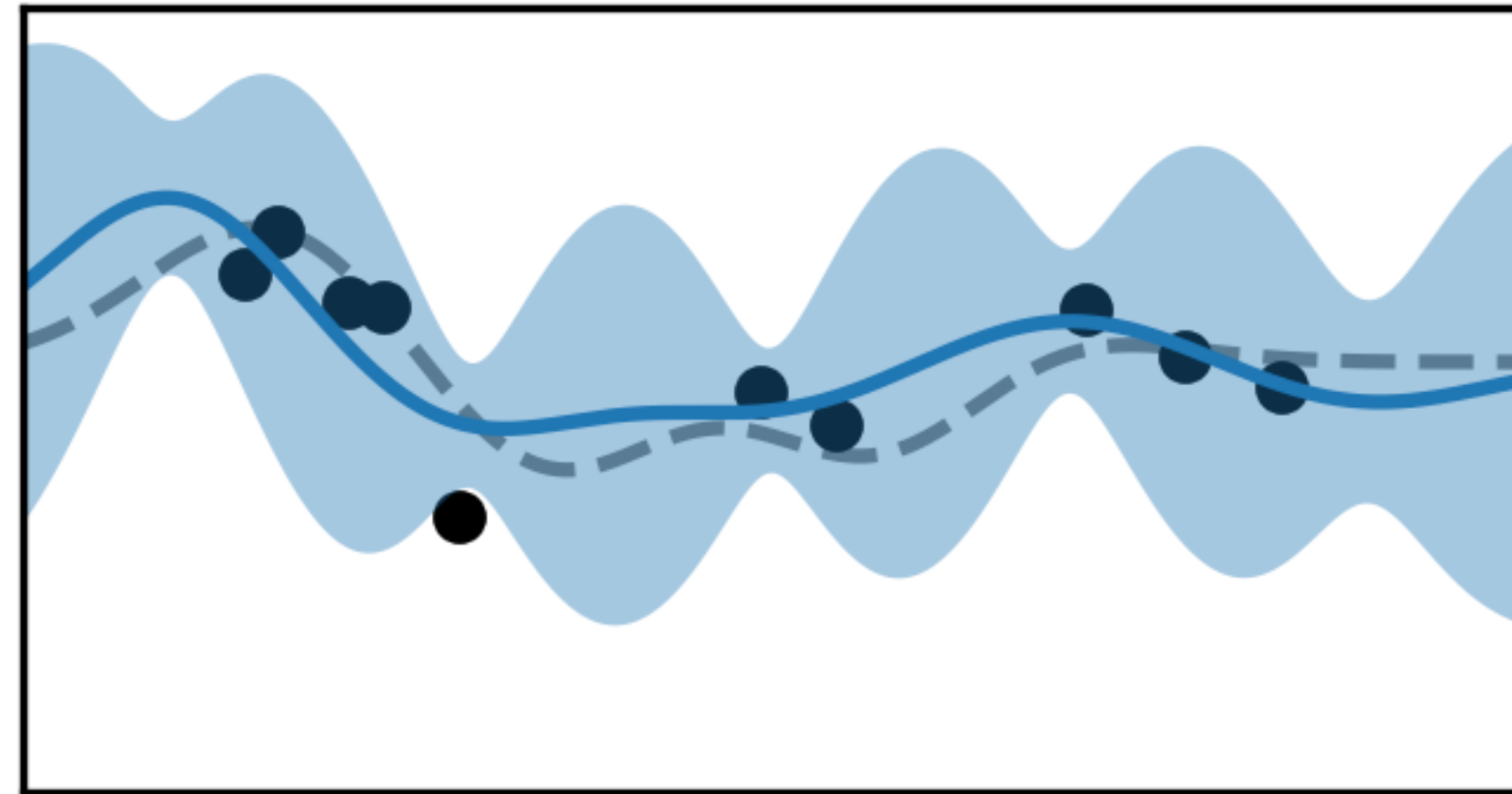
CGGP



Nystrom (SoR)



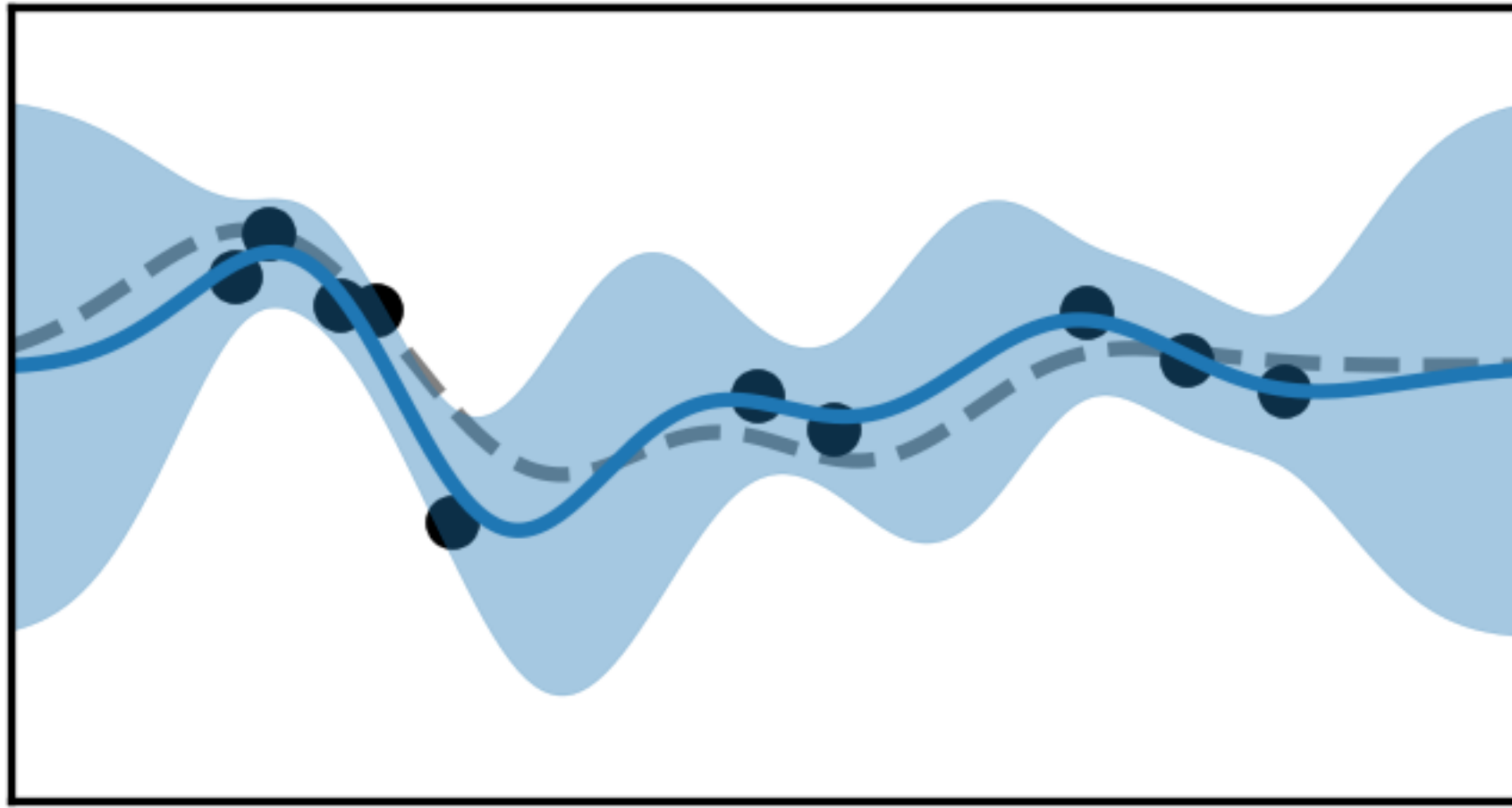
SVGP



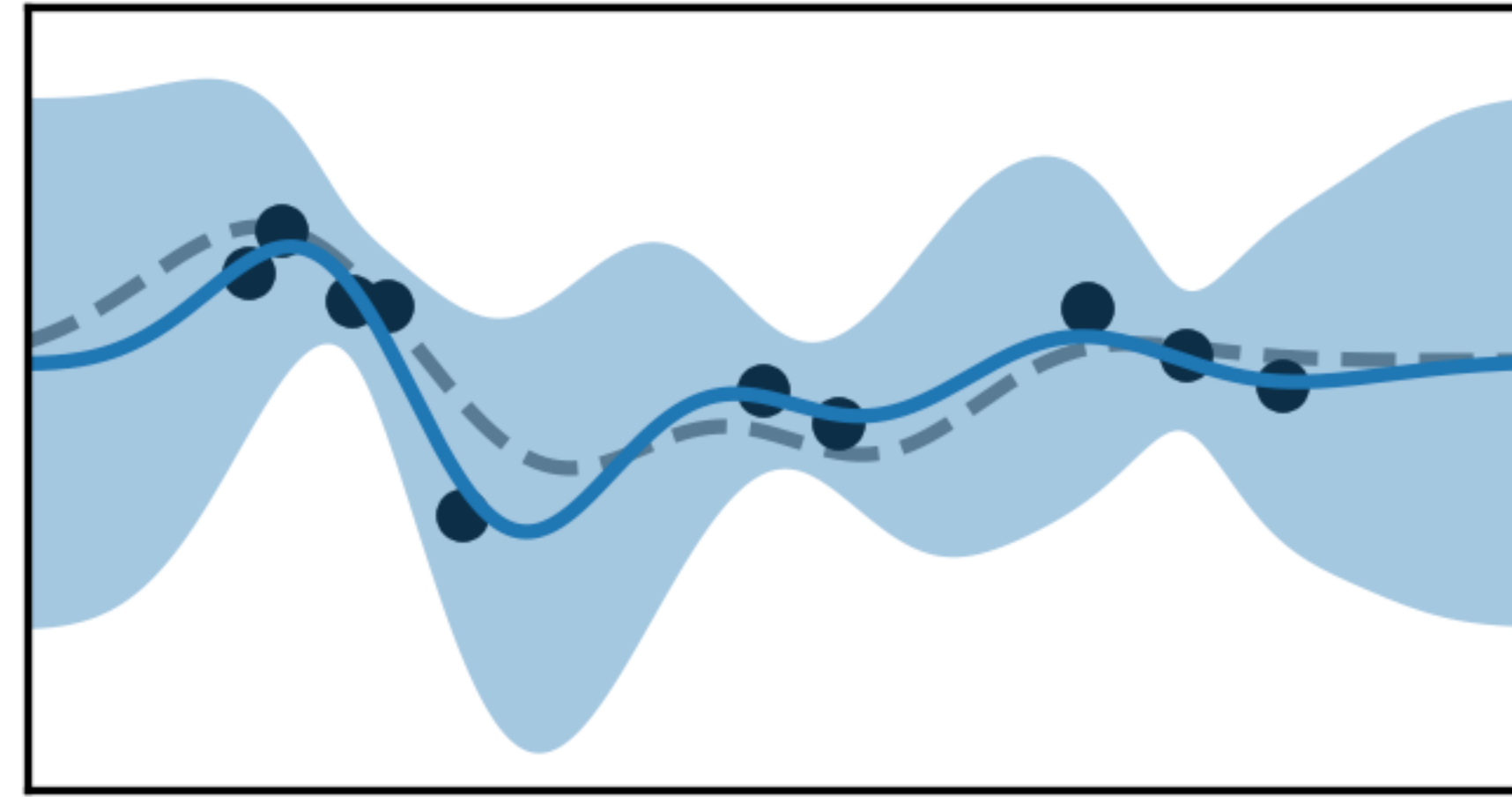
--- Latent function ● Data — Posterior mean Posterior uncertainty

Approximate Inference in Practice

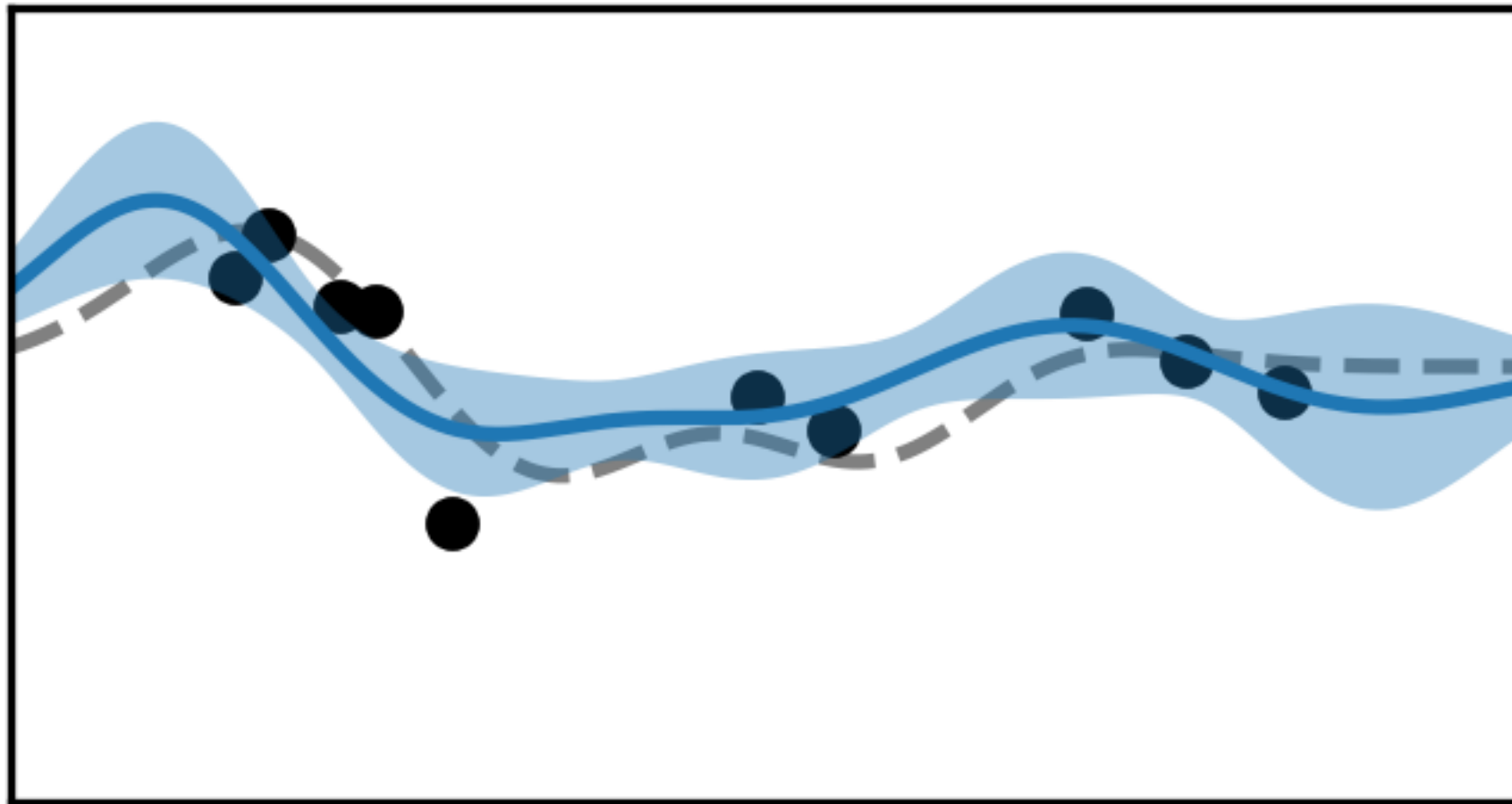
Mathematical Posterior



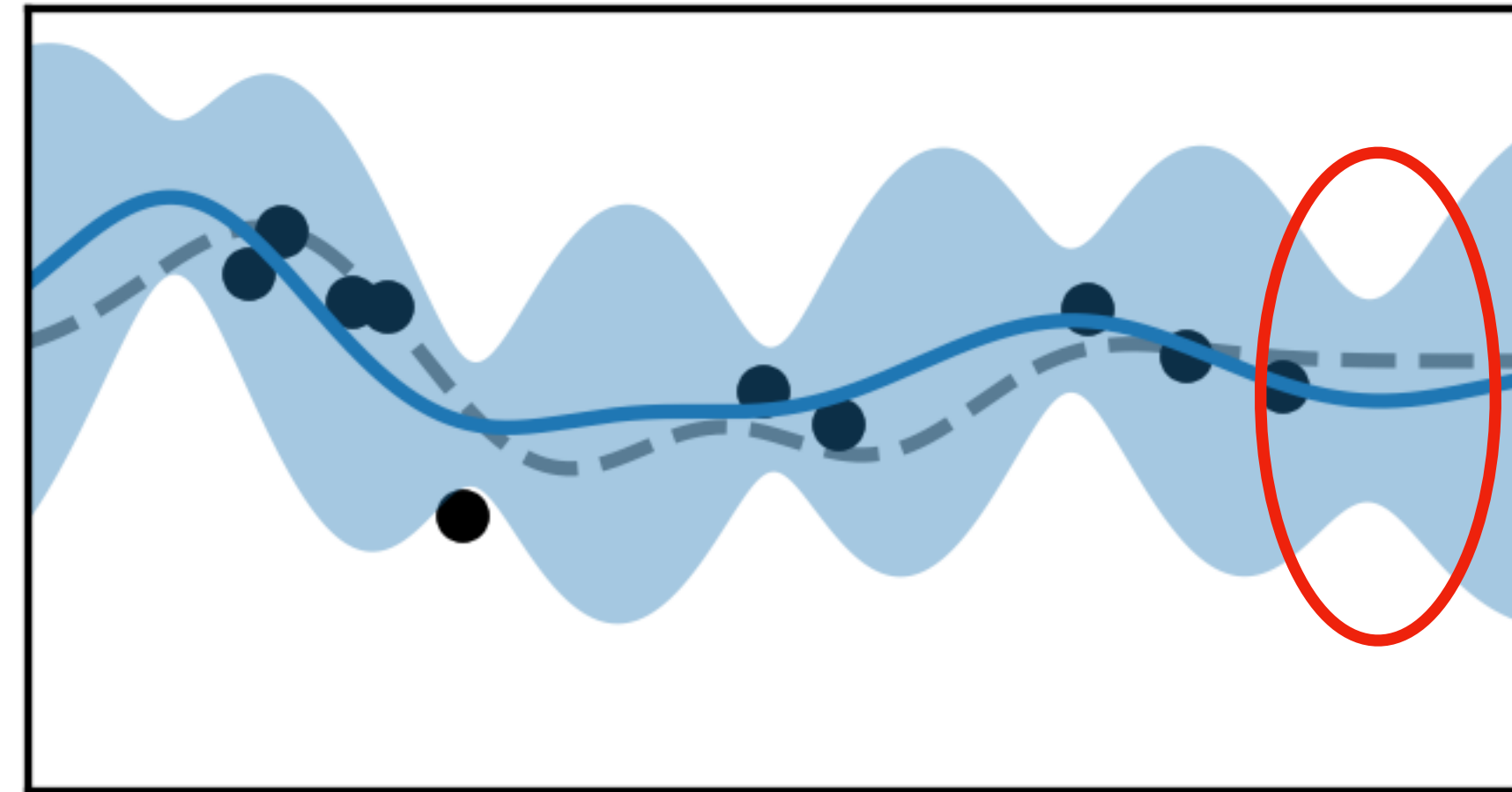
CGGP



Nyström (SoR)



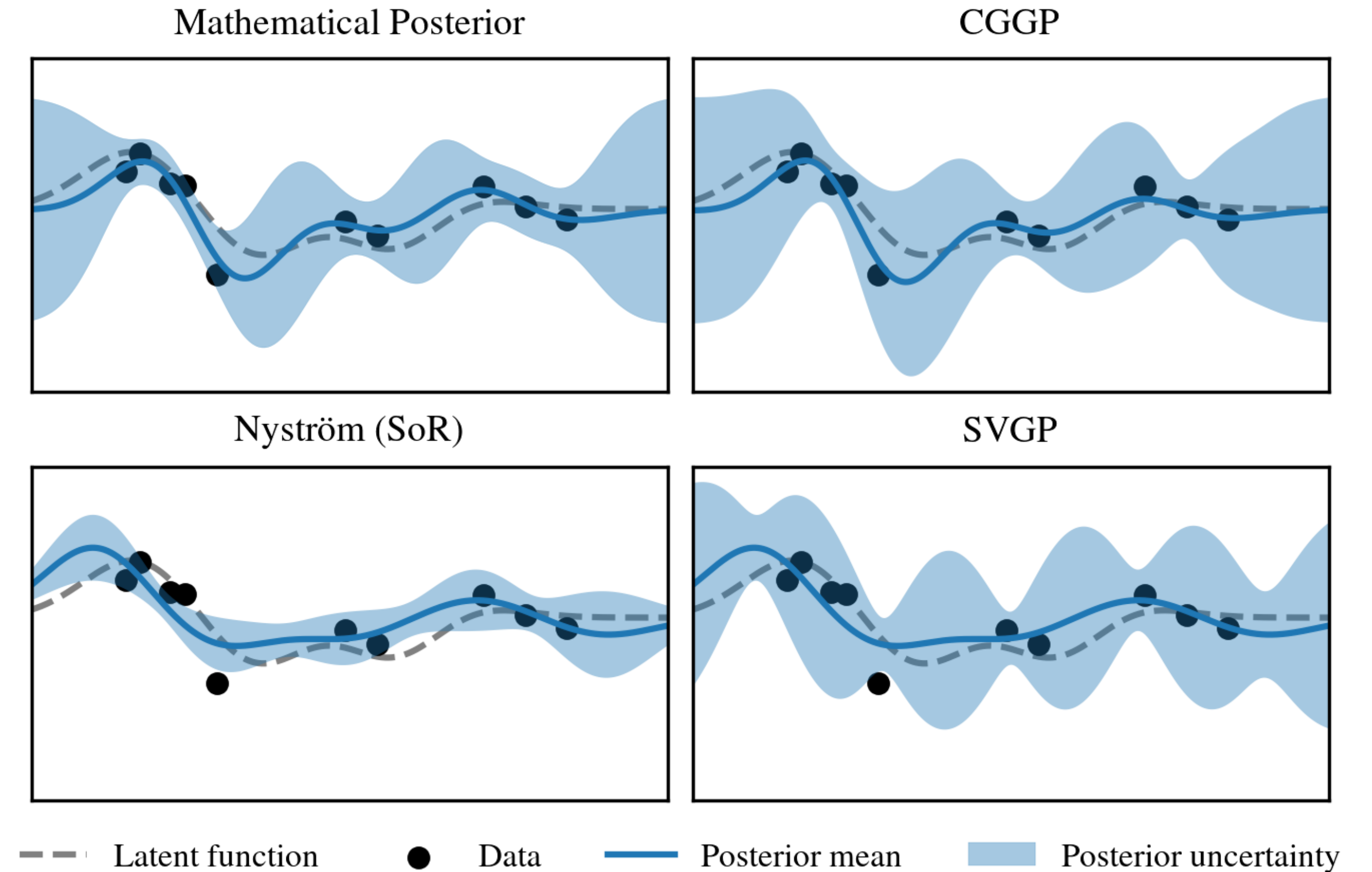
SVGP



--- Latent function ● Data — Posterior mean Posterior uncertainty

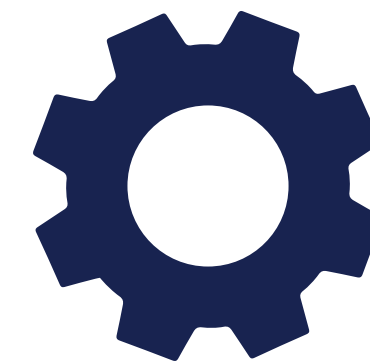
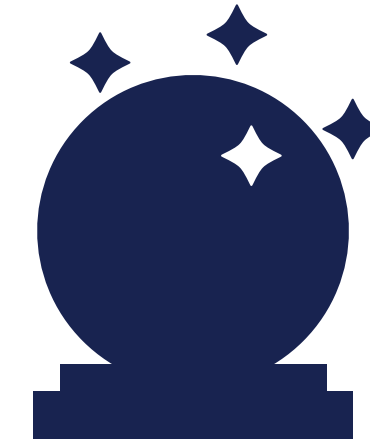
Approximate Inference in Practice

- ▶ The point: the approximation method is making strong (and different) statements about what you know
- ▶ But this effect is ignored...



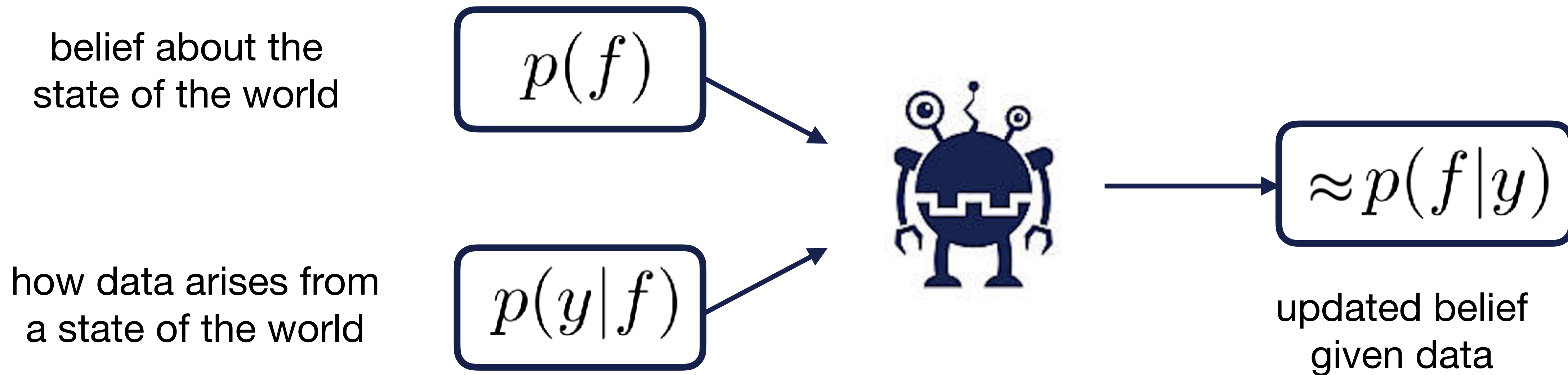
Outline

- ▶ The Promise of Probabilistic Machine Learning
- ▶ Gaussian Process Introduction
- ▶ Scaling Gaussian Processes, and Implications
- ▶ **Approximate Gaussian Process Inference, The Right Way**
- ▶ iterGP as Probabilistic Numerics
- ▶ Broader Implications



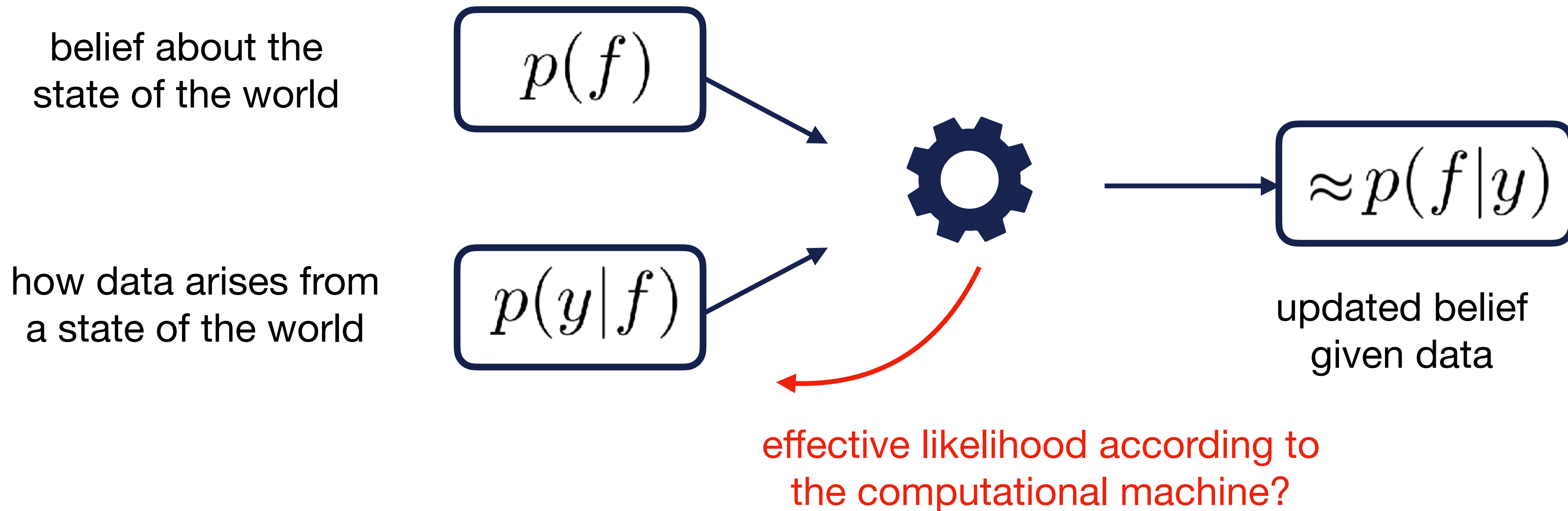
Approximate Inference as Exact Inference

- ▶ Apparently approximate inference is exact inference, but under some different model



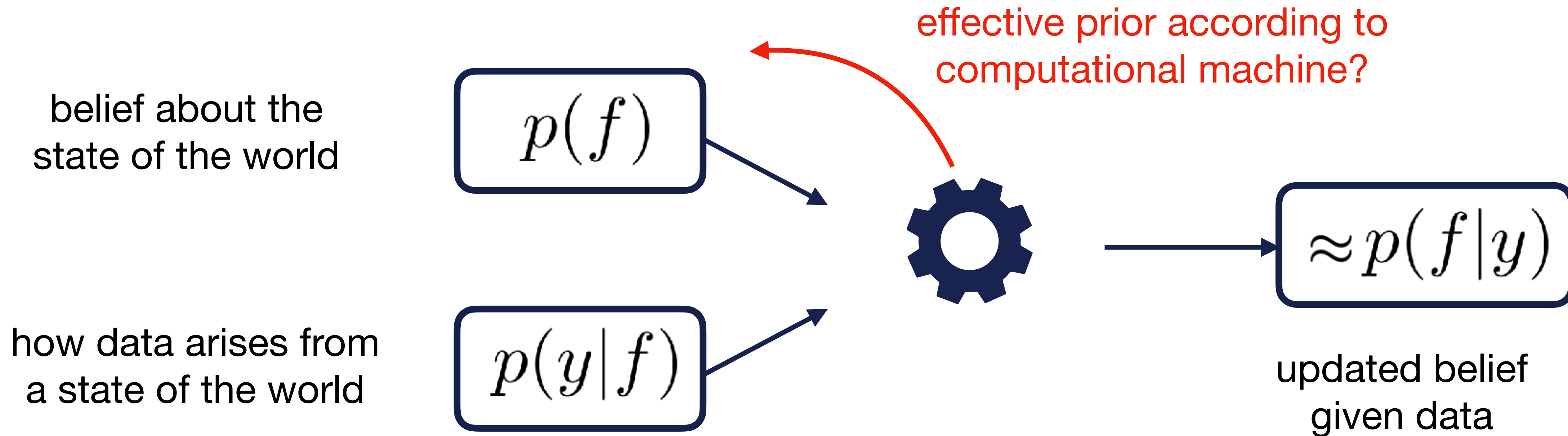
Approximate Inference as Exact Inference

- ▶ Apparently approximate inference is exact inference, but under some different model



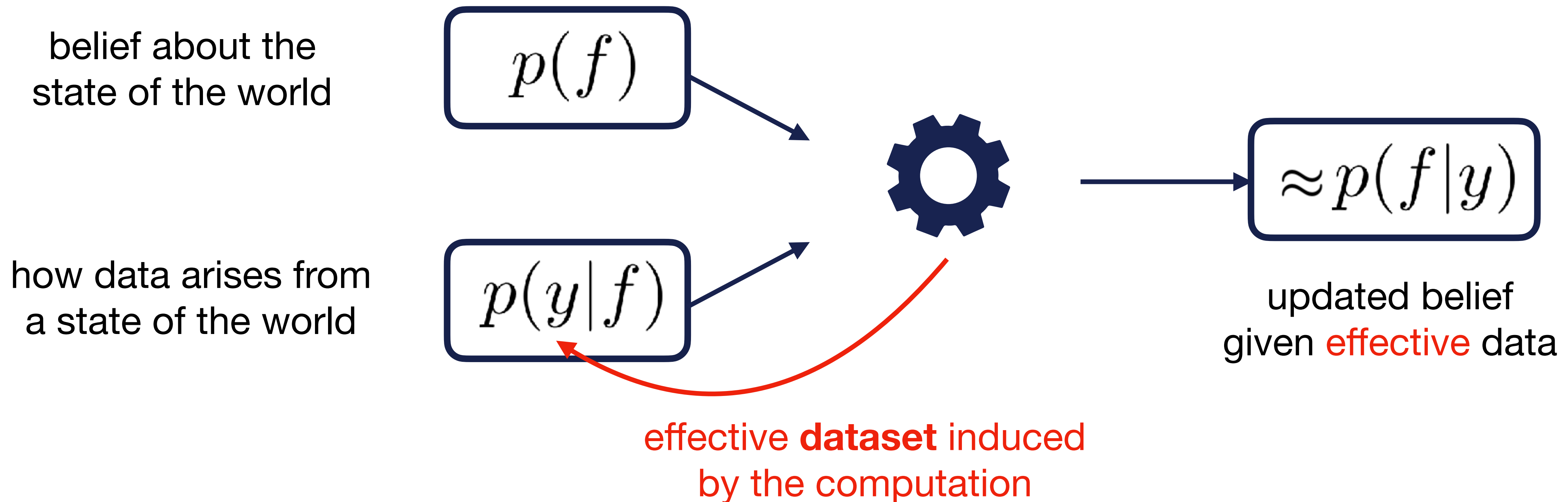
Approximate Inference as Exact Inference

- ▶ Apparently approximate inference is exact inference, but under some different model



Approximate Inference as Exact Inference

- ▶ Apparently approximate inference is exact inference, but under some different **data**



Effective Dataset

- ▶ Iterative numerical methods (for GP) take linear combinations of data
- ▶ Let us define the set of actions taken by a given (approximate) solver as

$$\mathbf{S}_i = [\mathbf{s}_1 \ \mathbf{s}_2 \ \dots \ \mathbf{s}_i] \in \mathbb{R}^{n \times i}$$

- ▶ This should feel plausible; consider the following actions:

| Actions \mathbf{s}_i | Classic Analog |
|---|------------------------------------|
| \mathbf{e}_i | (partial) Cholesky |
| $\text{ev}_i(\hat{\mathbf{K}})$ | (partial) EVD / SVD |
| $\mathbf{s}_i^{\text{PCG}}$ or $\hat{\mathbf{P}}^{-1} \mathbf{r}_i$ | (preconditioned) CG |
| $k(\mathbf{X}, \mathbf{z}_i)$ | \approx Nyström (SoR, DTC), SVGP |

- ▶ (We will make this rigorous shortly, for now we just establish the connection)

Gaussian Process Posterior

► With likelihood: $\mathbf{y} \mid \mathbf{f} \sim \mathcal{N}(\mathbf{f}, \sigma^2 \mathbf{I})$

► And test inputs: \mathbf{X}_\diamond

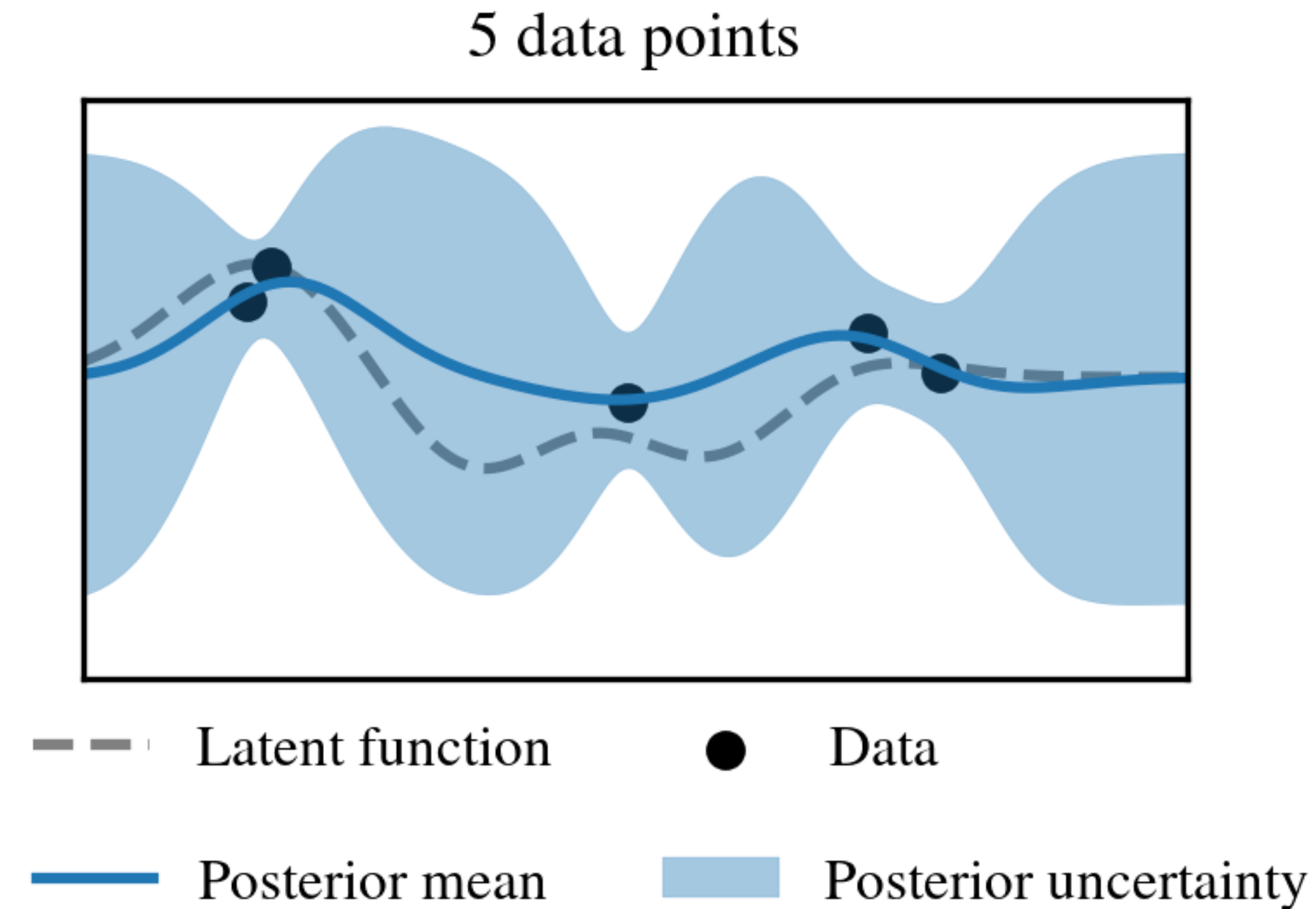
► Induces the posterior:

$$\mathbf{f}_\diamond \sim \mathcal{N}(\mu_*(\mathbf{X}_\diamond), k_*(\mathbf{X}_\diamond, \mathbf{X}_\diamond))$$

► where: $\mu_*(\cdot) = \mu(\cdot) + k(\cdot, \mathbf{X}) \hat{\mathbf{K}}^{-1}(\mathbf{y} - \boldsymbol{\mu})$

$$k_*(\cdot, \cdot) = k(\cdot, \cdot) - k(\cdot, \mathbf{X}) \hat{\mathbf{K}}^{-1} k(\mathbf{X}, \cdot)$$

$$\hat{\mathbf{K}} := \mathbf{K} + \sigma^2 \mathbf{I} \in \mathbb{R}^{n \times n}$$



Gaussian Process Posterior

► With likelihood: $\mathbf{S}_i^\top \mathbf{y} \mid \mathbf{f} \sim \mathcal{N}(\mathbf{S}_i^\top \mathbf{f}, \sigma^2 \mathbf{S}_i^\top \mathbf{S}_i)$

► And test inputs: \mathbf{X}_\diamond

► Induces the posterior:

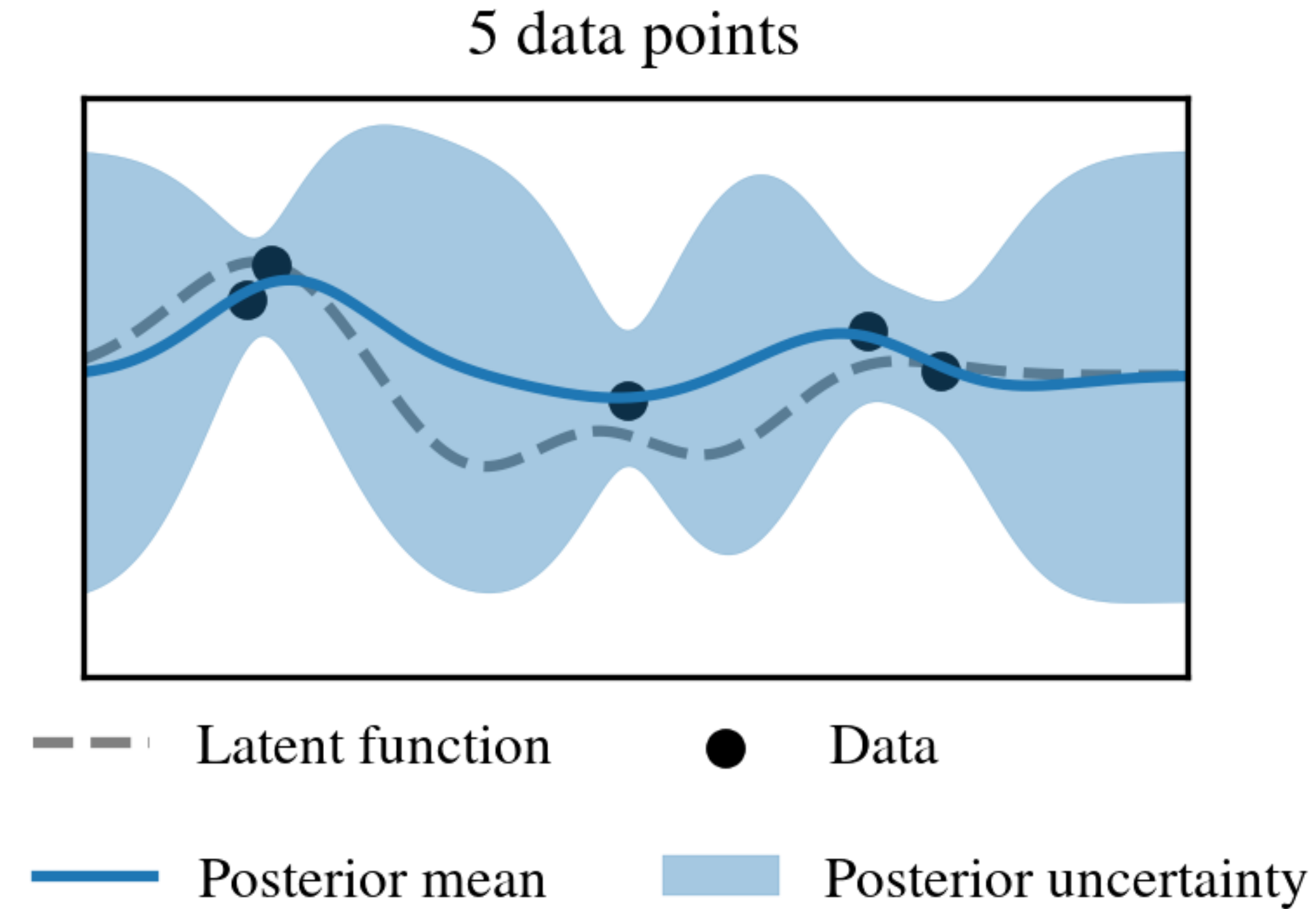
$$\mathbf{f}_\diamond \sim \mathcal{N}(\mathbf{f}_\diamond; \mu_i(\mathbf{X}_\diamond), k_i(\mathbf{X}_\diamond, \mathbf{X}_\diamond))$$

► where: $\mu_i(\cdot) = \mu(\cdot) + k(\cdot, \mathbf{X})\mathbf{C}_i(\mathbf{y} - \mu)$

$$k_i(\cdot, \cdot) = k(\cdot, \cdot) - k(\cdot, \mathbf{X})\mathbf{C}_i k(\mathbf{X}, \cdot)$$

$$\mathbf{C}_i = \mathbf{S}_i(\mathbf{S}_i^\top \hat{\mathbf{K}} \mathbf{S}_i)^{-1} \mathbf{S}_i^\top$$

$$\hat{\mathbf{K}} := \mathbf{K} + \sigma^2 \mathbf{I} \in \mathbb{R}^{n \times n}$$



Combined Uncertainty

- ▶ *Combined Uncertainty* admits a clean decomposition:

$$\underbrace{k(\cdot, \cdot) - k(\cdot, \mathbf{X})\mathbf{C}_i k(\mathbf{X}, \cdot)}_{\text{combined uncertainty } \textcolor{purple}{\bullet}} = \underbrace{k(\cdot, \cdot) - k(\cdot, \mathbf{X})\hat{\mathbf{K}}^{-1} k(\mathbf{X}, \cdot)}_{\text{mathematical uncertainty } \textcolor{blue}{\bullet}} + \underbrace{k(\cdot, \mathbf{X})\boldsymbol{\Sigma}_i k(\mathbf{X}, \cdot)}_{\text{computational uncertainty } \textcolor{green}{\bullet}}$$

$$\mathbf{C}_i = \mathbf{S}_i (\mathbf{S}_i^\top \hat{\mathbf{K}} \mathbf{S}_i)^{-1} \mathbf{S}_i^\top$$

$$\hat{\mathbf{K}} := \mathbf{K} + \sigma^2 \mathbf{I}$$

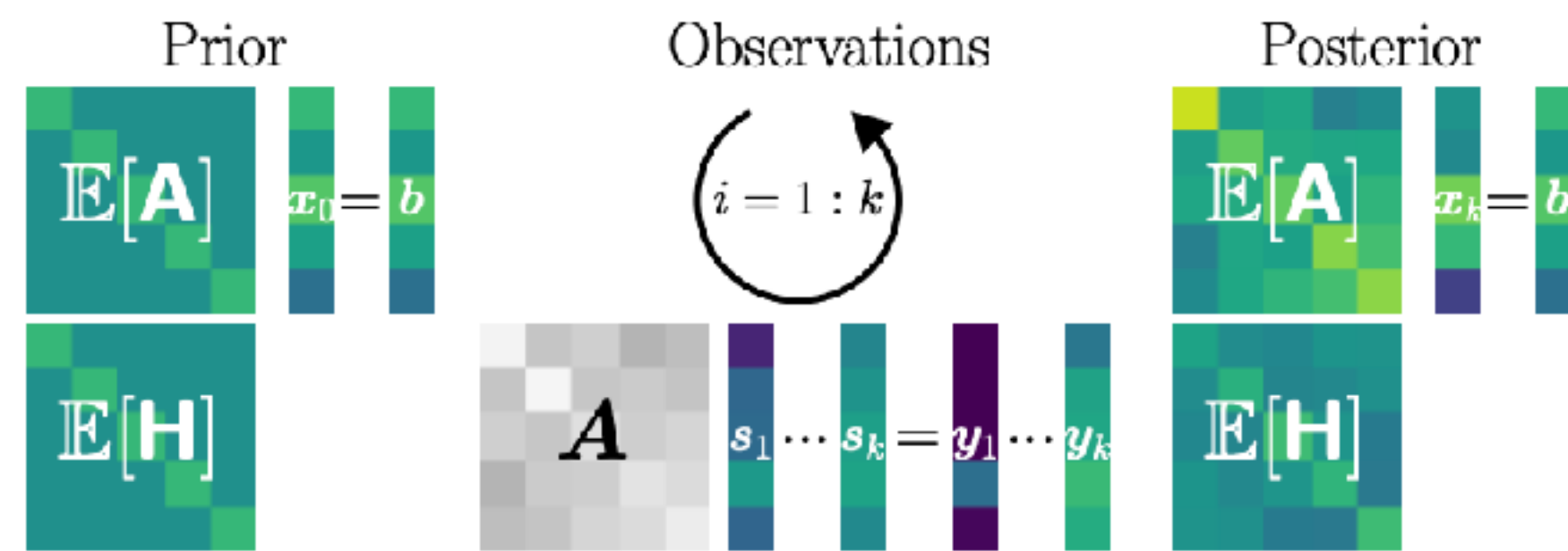
Computational
uncertainty turns out
to be exactly the
uncertainty on the
representer weights

- ▶ Very satisfying outcome:
 - ▶ Our combined (effective) uncertainty is a result of both finite data and finite compute
 - ▶ We call *mathematical uncertainty* the exact posterior with the given data
 - ▶ Computational uncertainty is a direct use of Probabilistic Numerics!

Interlude: so is this a PN story or a GP story?

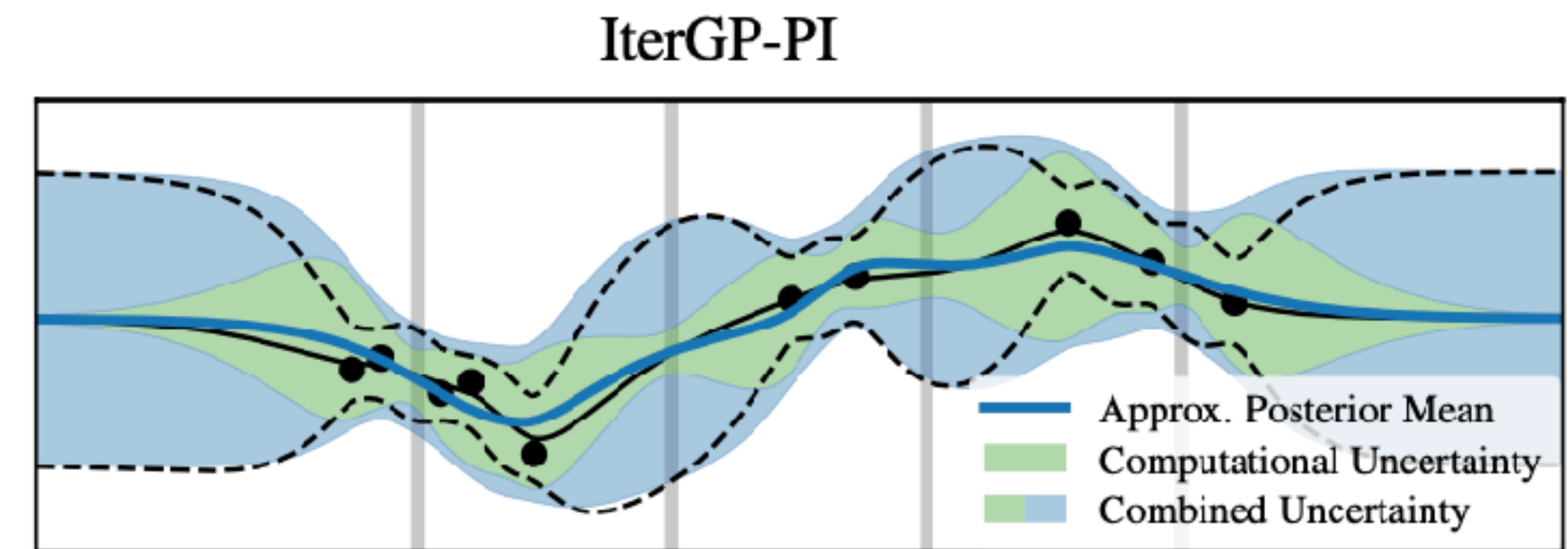
- Both. What I am telling you today builds on huge amounts GP approximate inference work and a significant amount of PN work, two most notably:

PN: linear solvers can be seen as probabilistic inference methods



[Wenger and Hennig NeurIPS 2020]

GP: inference is all linear solves, from which arises the effective dataset



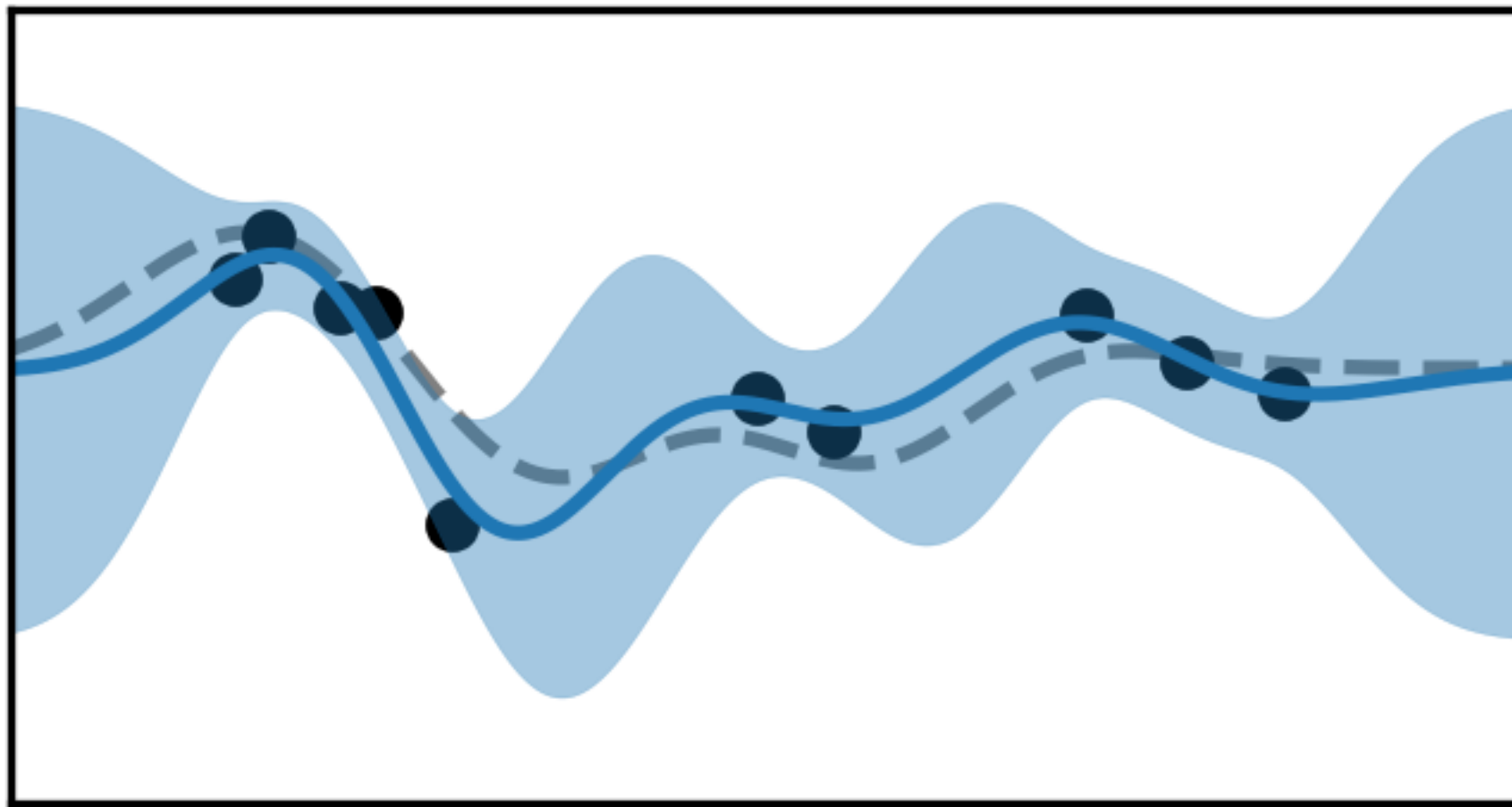
[Wenger et al NeurIPS 2022]

- IterGP* is the family of methods producing combined GP uncertainty from iterative solvers

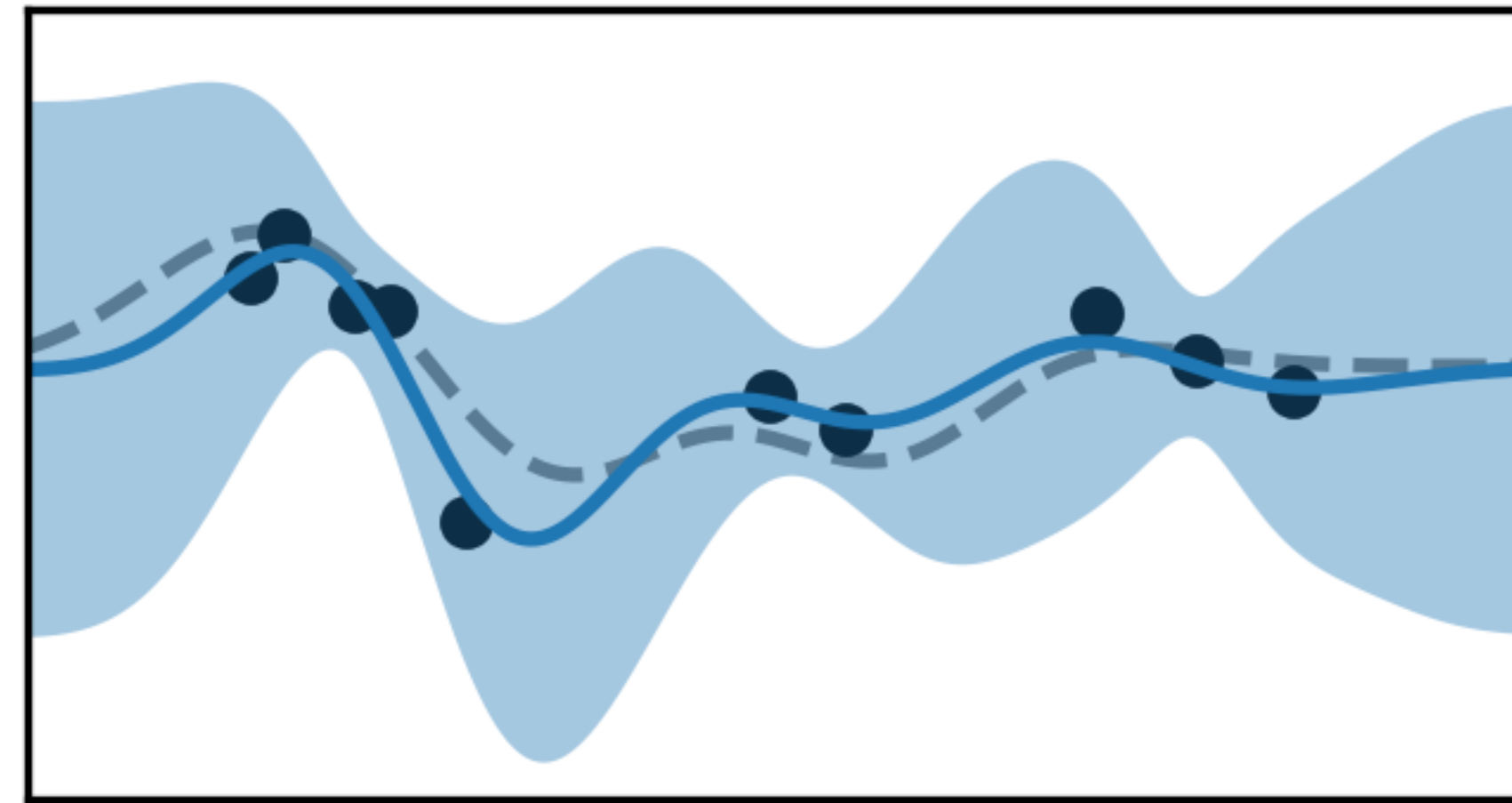
Old Way: Approximate Inference in Practice



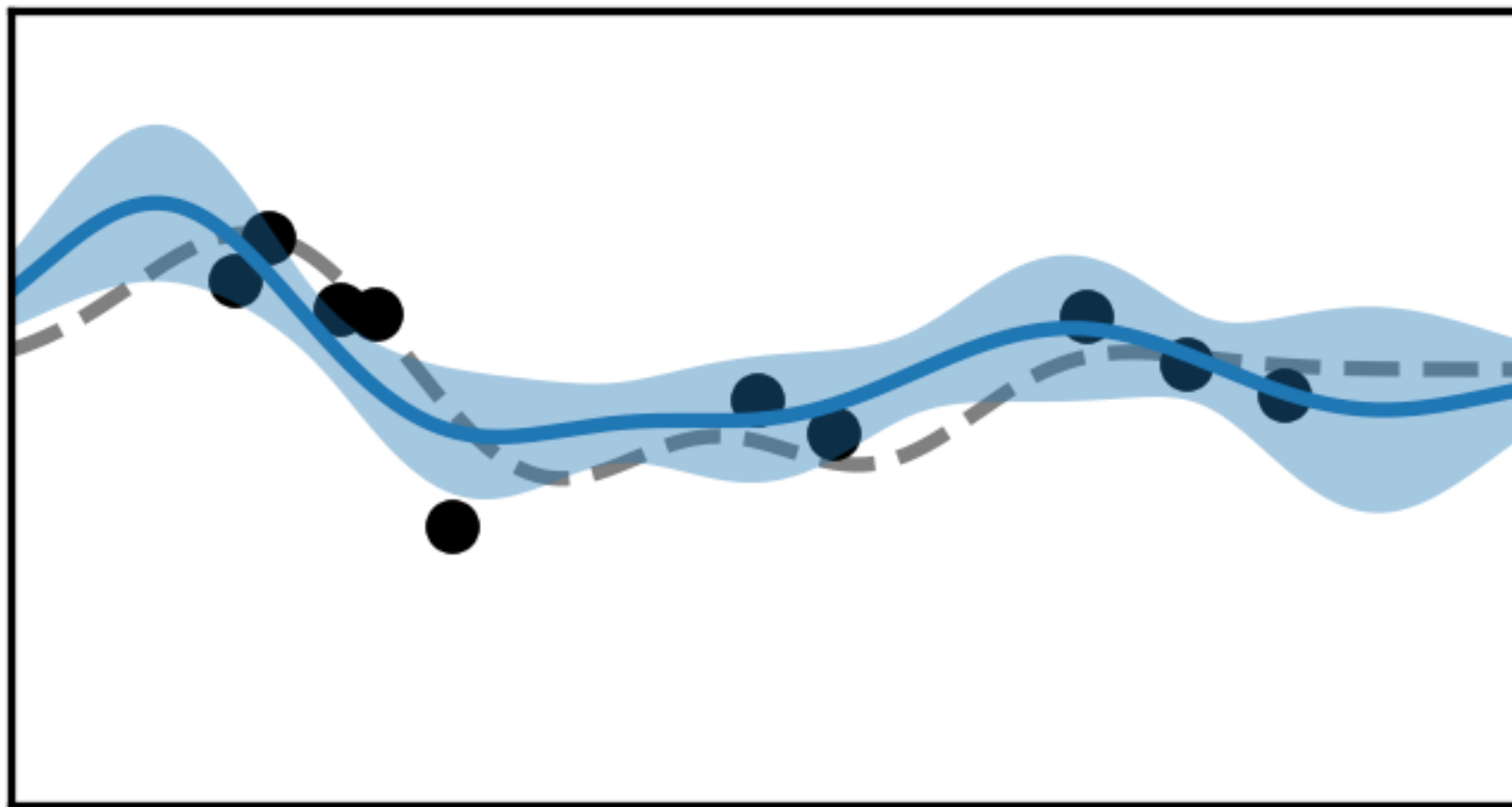
Mathematical Posterior



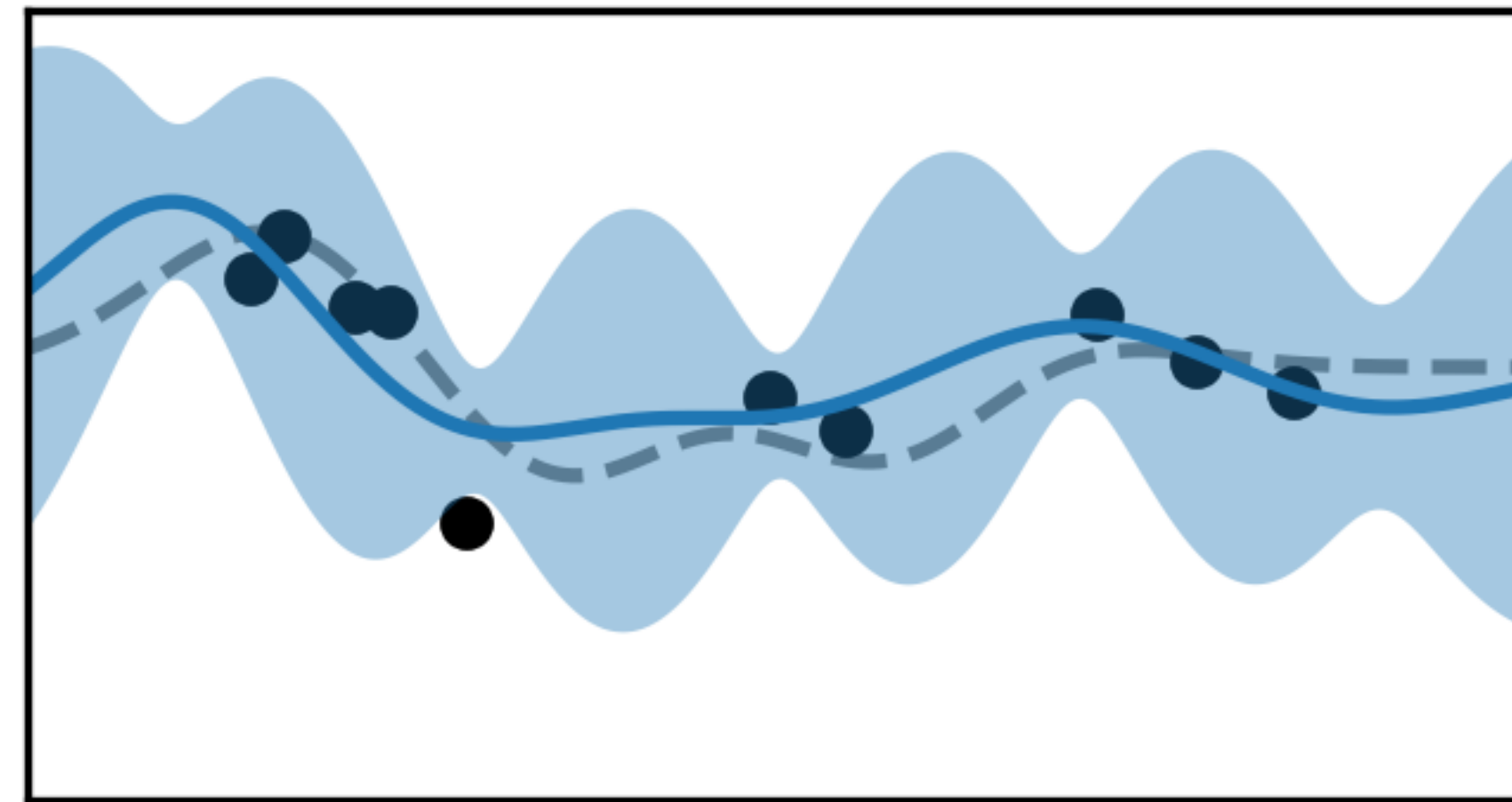
CGGP



Nyström (SoR)

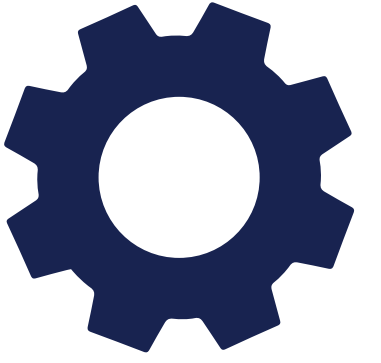
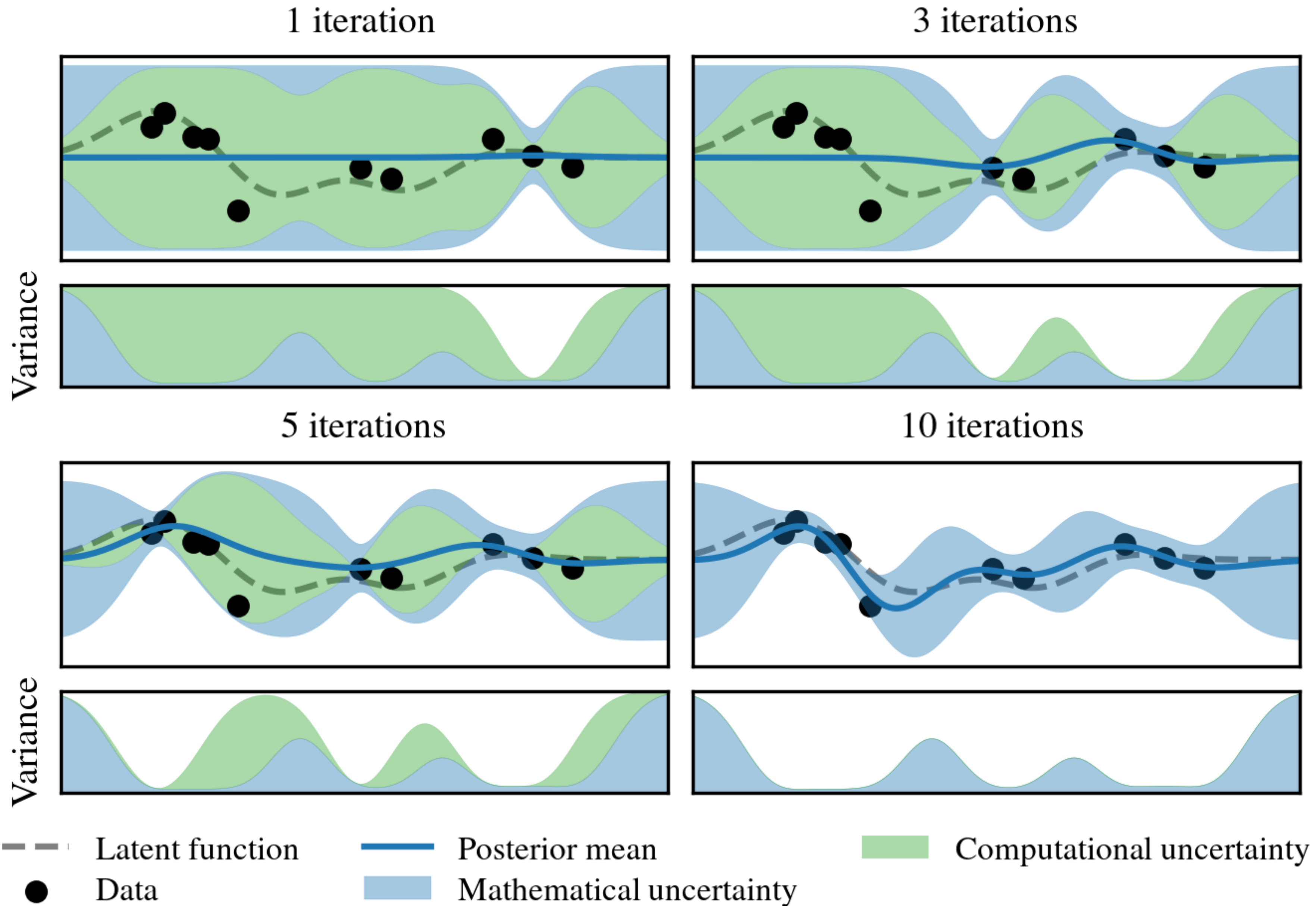


SVGP



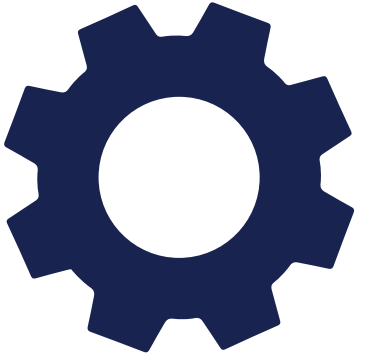
— — Latent function ● Data — Posterior mean Posterior uncertainty

New Way: Sequentially Updating Data Points as IterGP



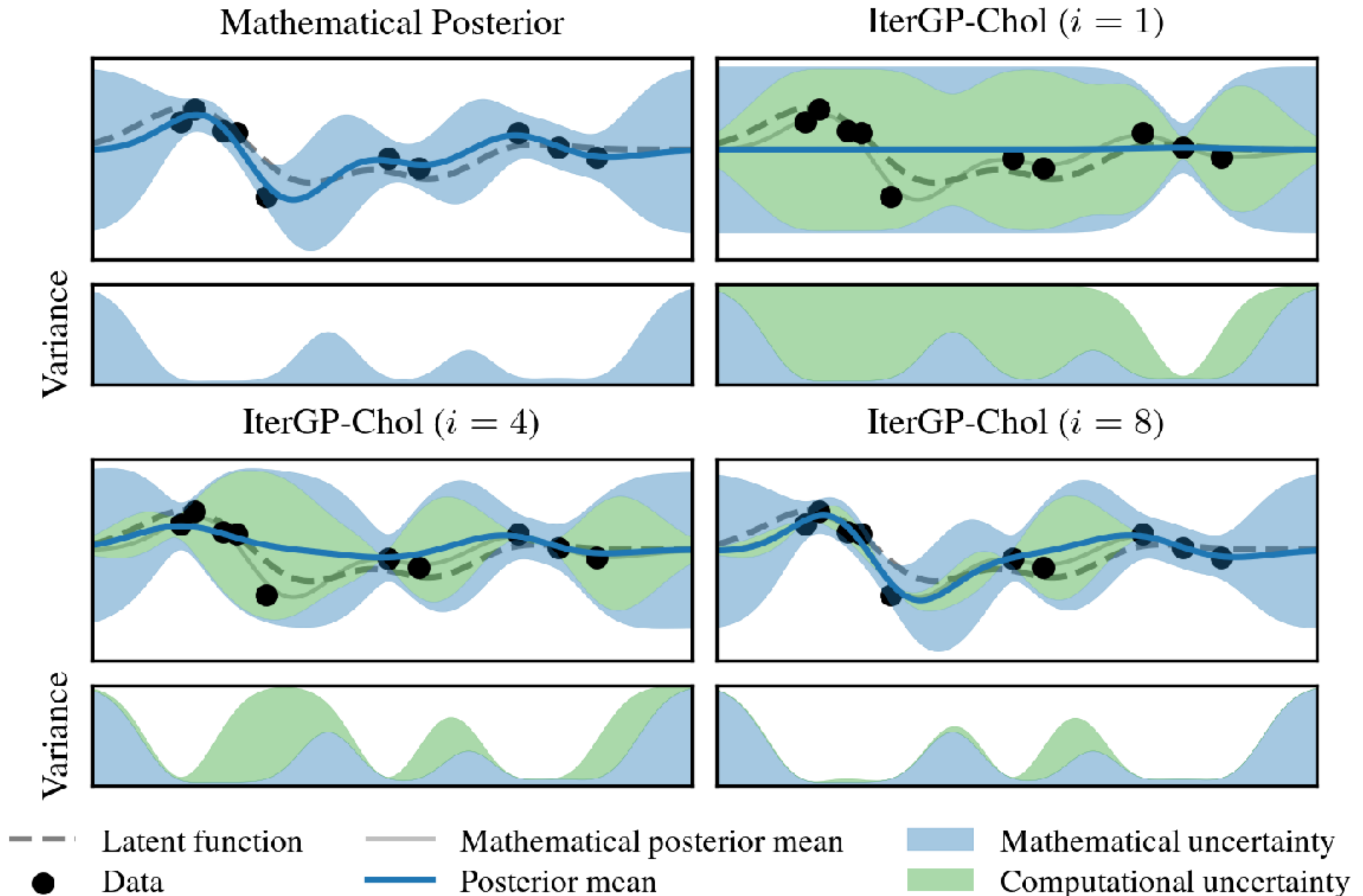
New Way: IterGP-Cholesky

| Method | Actions \mathbf{s}_i | Classic Analog |
|-------------|------------------------|--------------------|
| IterGP-Chol | \mathbf{e}_i | (partial) Cholesky |

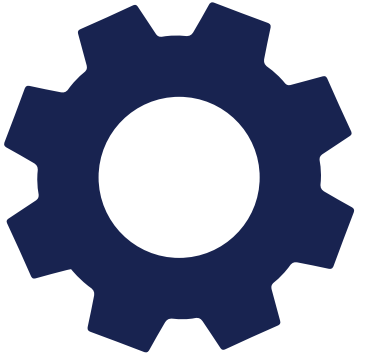


► Recall:

$$\mathbf{S}_i^\top \mathbf{y} \mid \mathbf{f} \sim \mathcal{N}(\mathbf{S}_i^\top \mathbf{f}, \sigma^2 \mathbf{S}_i^\top \mathbf{S}_i)$$



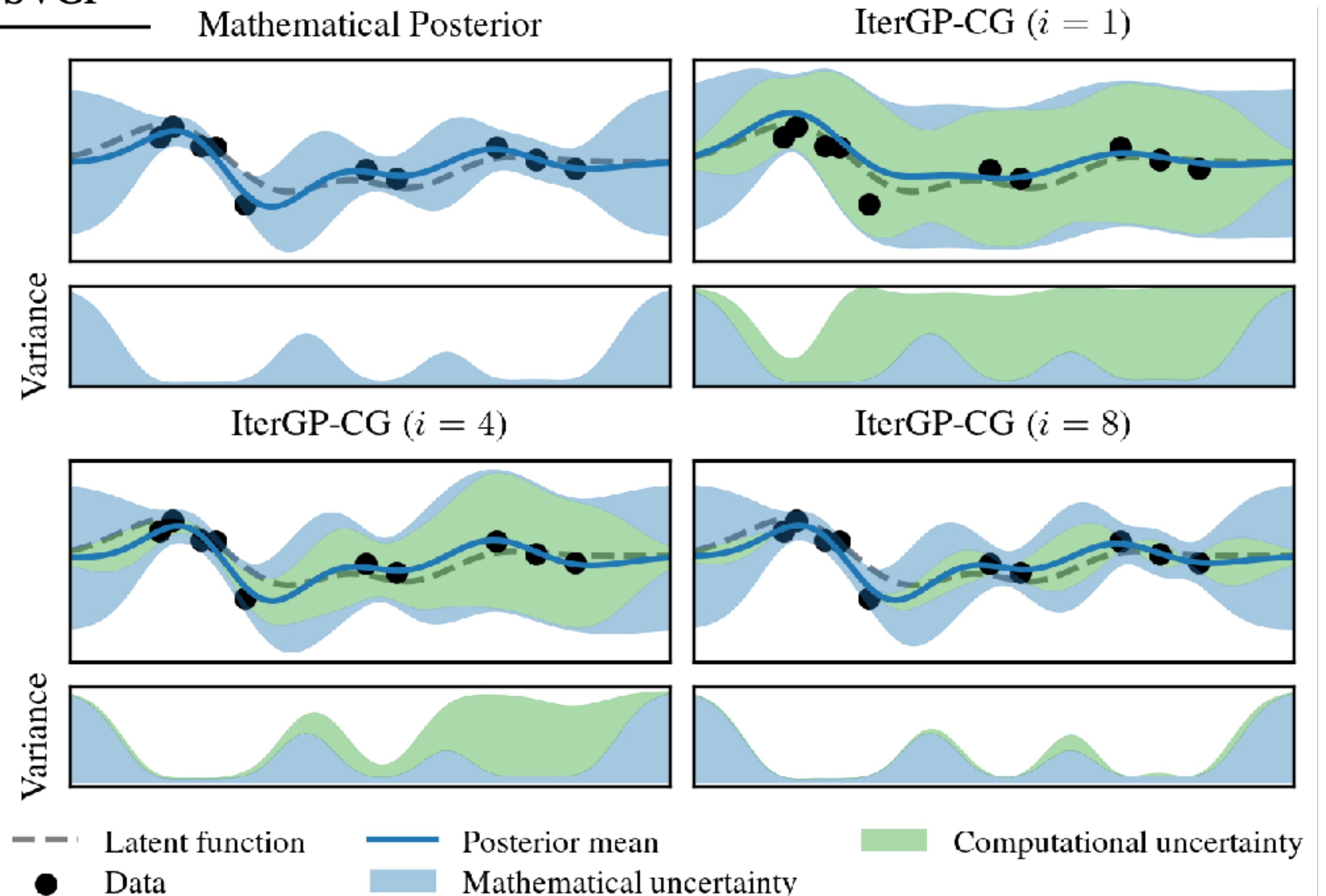
New Way: IterGP-CG



| Method | Actions \mathbf{s}_i | Classic Analog |
|-------------|---|------------------------------------|
| IterGP-Chol | \mathbf{e}_i | (partial) Cholesky |
| IterGP-PBR | $\text{ev}_i(\hat{\mathbf{K}})$ | (partial) EVD / SVD |
| IterGP-CG | $\mathbf{s}_i^{\text{PCG}}$ or $\hat{\mathbf{P}}^{-1} \mathbf{r}_i$ | (preconditioned) CG |
| IterGP-PI | $k(\mathbf{X}, \mathbf{z}_i)$ | \approx Nyström (SoR, DTC), SVGP |

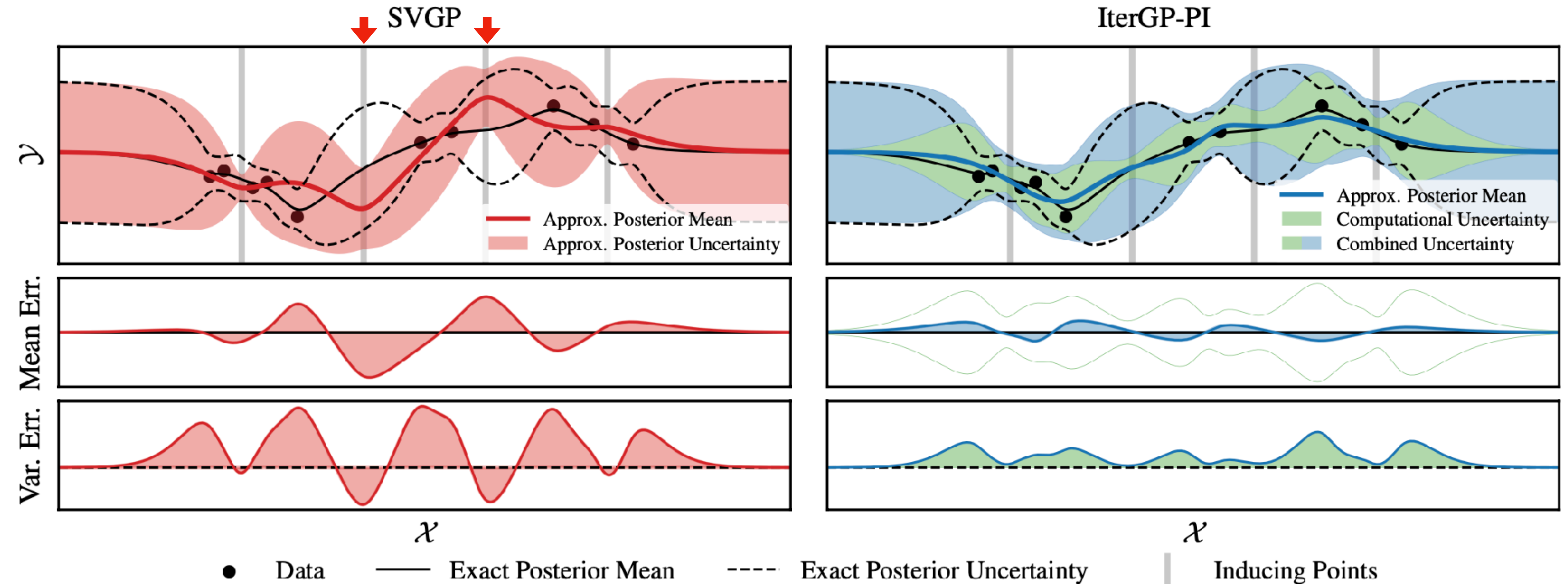
► Notice:

- Computational uncertainty adds to mathematical uncertainty... of course it should
- Untouched data (eg $i=1$ at right) has high combined uncertainty \rightarrow prior



Compare: iterGP-Pseudoinput vs SVGP

- ▶ Combined uncertainty reliably overestimates the truth \rightarrow this is desirable (and correct)!
- ▶ SVGP is overconfident at its inducing points, combined uncertainty corrects this.
- ▶ Overconfidence can produce large mean errors (↓)



Compare: iterGP-Pseudoinput vs SVGP

- ▶ Looking at posterior means is instructive:
 - ▶ Same as before, but for comparison let

$$q(\cdot, \cdot) = k(\cdot, \mathbf{Z}) \mathbf{K}_{\mathbf{Z}\mathbf{Z}}^{-1} k(\mathbf{Z}, \cdot)$$

- ▶ Then SVGP has posterior mean:

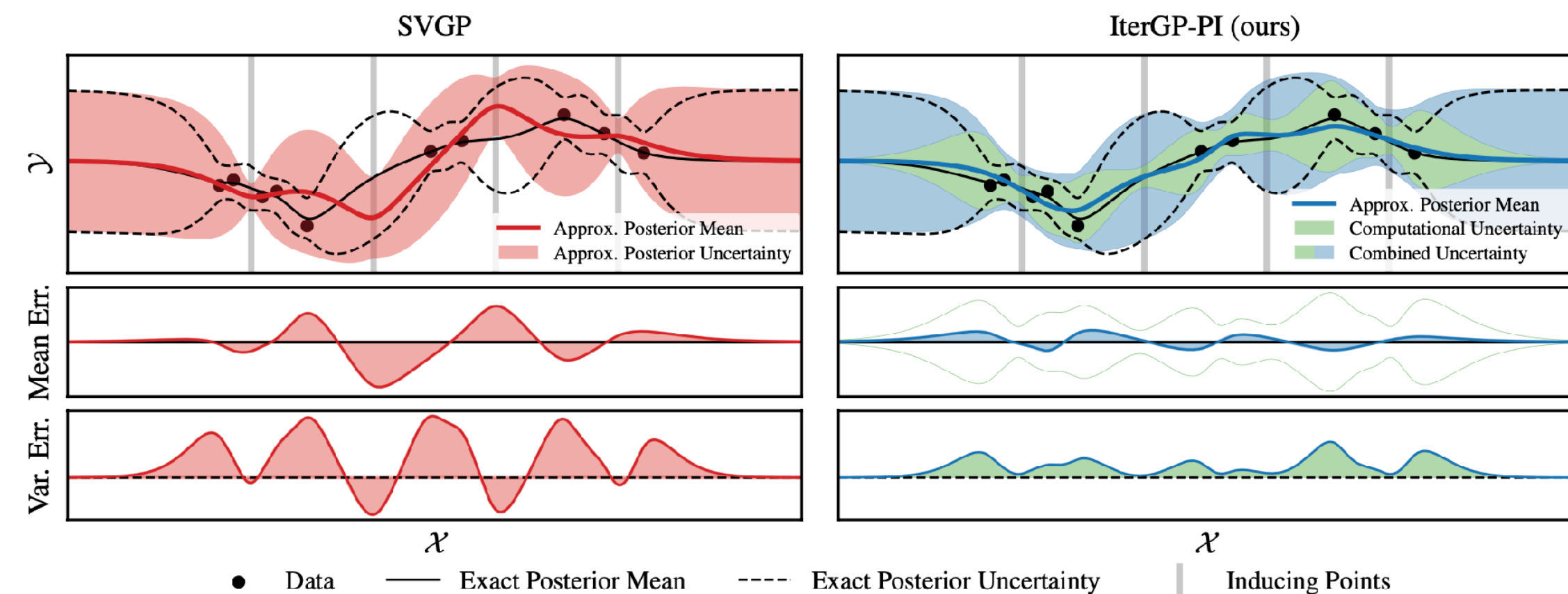
$$\mu_{\text{SVGP}}(\cdot) = q(\cdot, \mathbf{X}) \mathbf{K}_{\mathbf{X}\mathbf{Z}} (\mathbf{K}_{\mathbf{Z}\mathbf{X}} (q(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I}) \mathbf{K}_{\mathbf{X}\mathbf{Z}})^{-1} \mathbf{K}_{\mathbf{Z}\mathbf{X}} (\mathbf{y} - \boldsymbol{\mu})$$

- ▶ And IterGP-PI has posterior mean (for actions. $k(\mathbf{X}, \mathbf{z}_i)$, recall $\mathbf{C}_i = \mathbf{S}_i (\mathbf{S}_i^\top \hat{\mathbf{K}} \mathbf{S}_i)^{-1} \mathbf{S}_i^\top$)

$$\mu_i(\cdot) = k(\cdot, \mathbf{X}) \mathbf{K}_{\mathbf{X}\mathbf{Z}} (\mathbf{K}_{\mathbf{Z}\mathbf{X}} (k(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I}) \mathbf{K}_{\mathbf{X}\mathbf{Z}})^{-1} \mathbf{K}_{\mathbf{Z}\mathbf{X}} (\mathbf{y} - \boldsymbol{\mu})$$

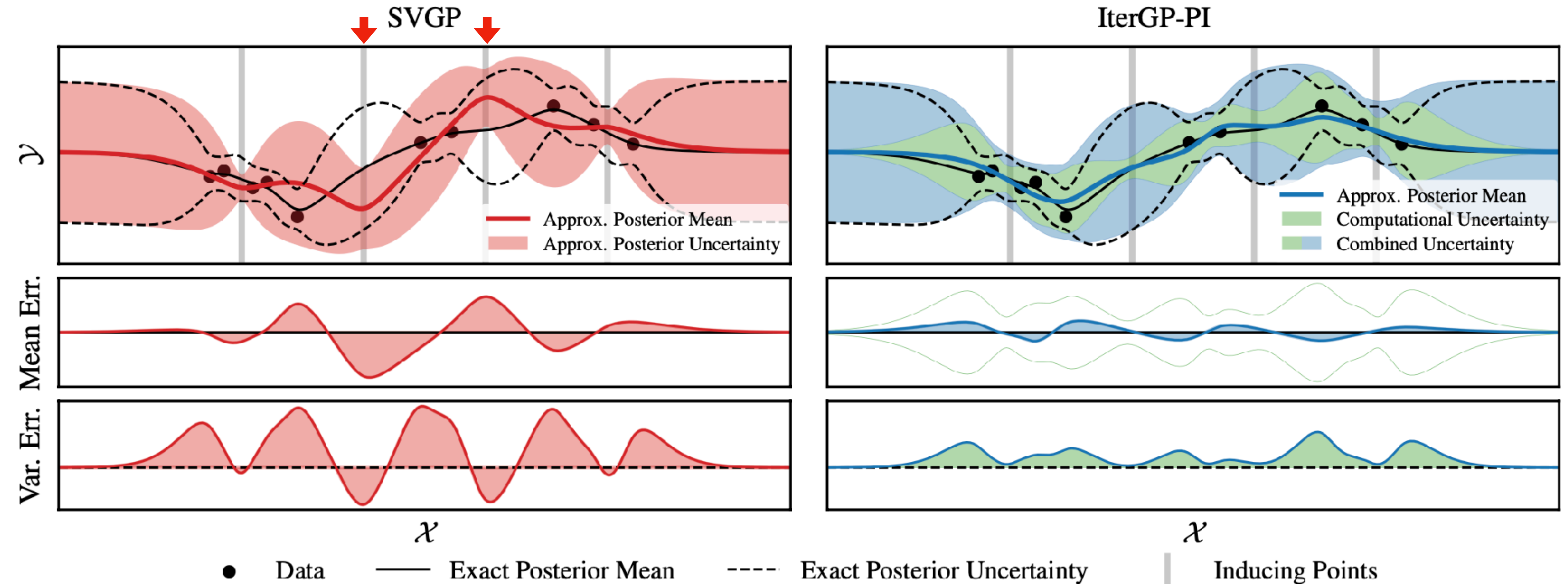
- ▶ Speedup in SVGP comes at cost of overconfidence, since $q(\mathbf{X}, \mathbf{X}) \preceq k(\mathbf{X}, \mathbf{X})$

- ▶ SVGP has too strong belief in representer weights, leading to potentially large errors even with a good variational fit. Computational uncertainty reduces that confidence.



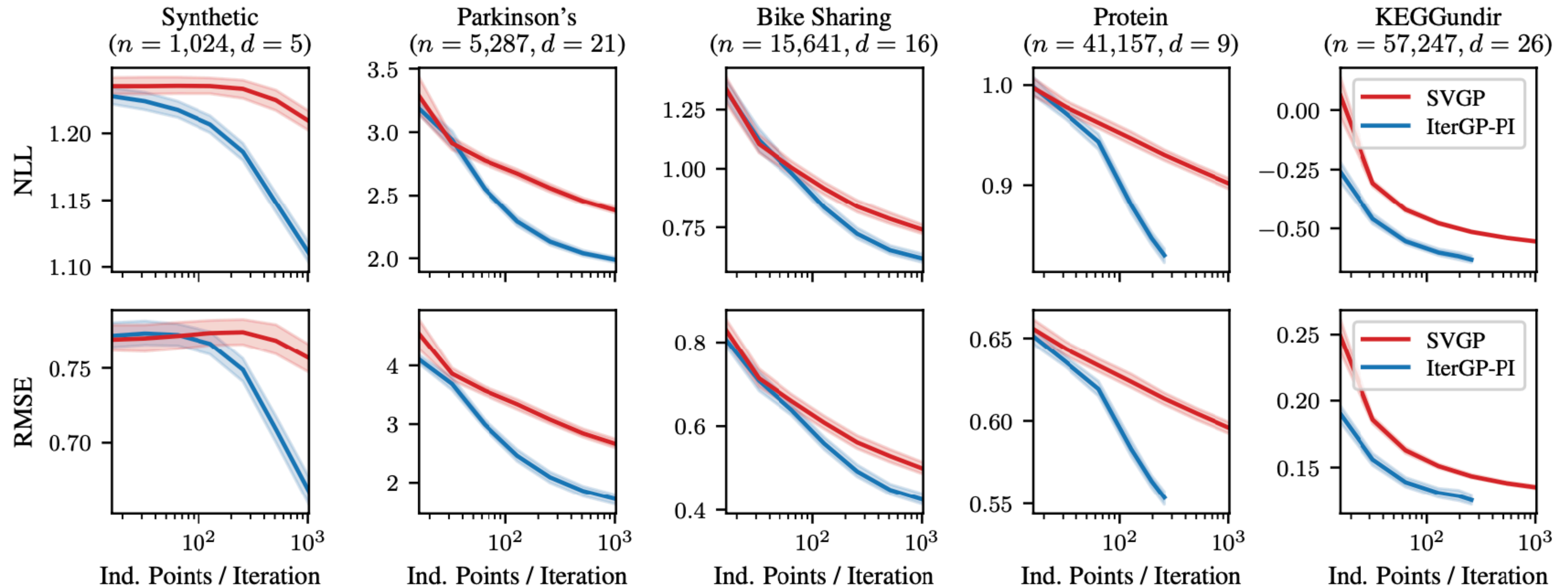
Compare: iterGP-Pseudoinput vs SVGP

- ▶ Combined uncertainty reliably overestimates the truth \rightarrow this is desirable (and correct)!
- ▶ SVGP is overconfident at its inducing points, combined uncertainty corrects this.
- ▶ Overconfidence can produce large mean errors (↓)



Computational Uncertainty Matters

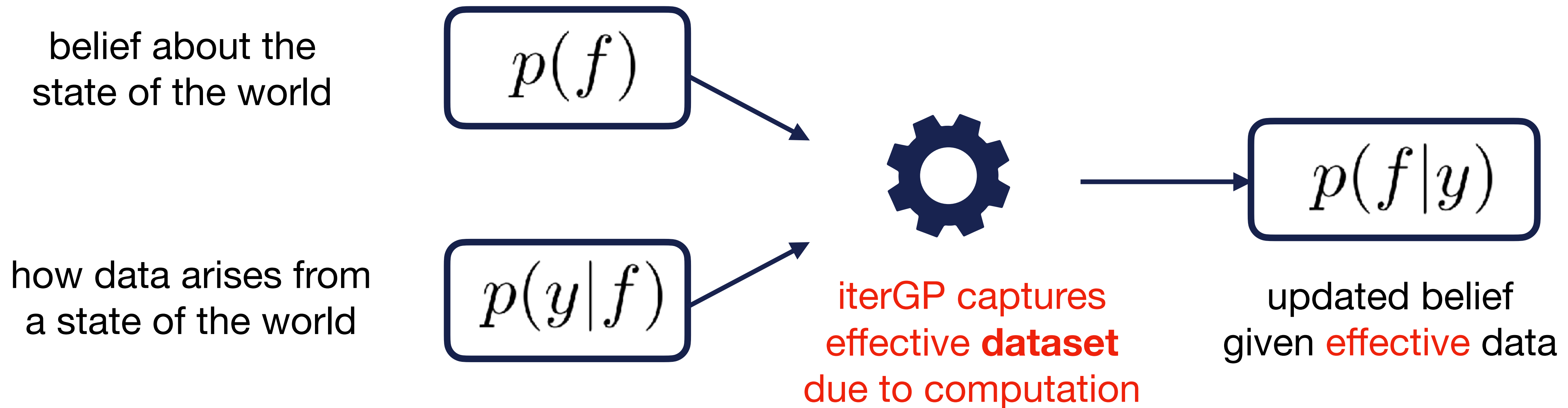
- ▶ IterGP achieves better generalization with a smaller number of inducing points (vs SVGP)



- ▶ Note however that this is a statement about use of inducing points, not performance per flop

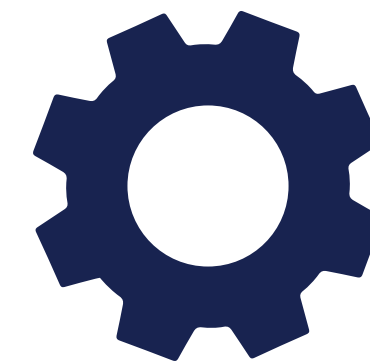
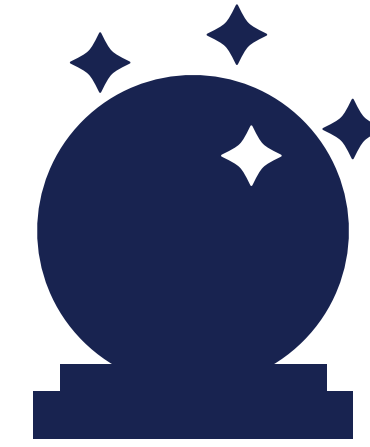
Approximate Inference as Exact Inference (on Different Data)

- ▶ Result: if you update your GP via matrix-vector multiplication, then the combined uncertainty of the IterGP algorithm is precisely the correct object to capture your belief.



Outline

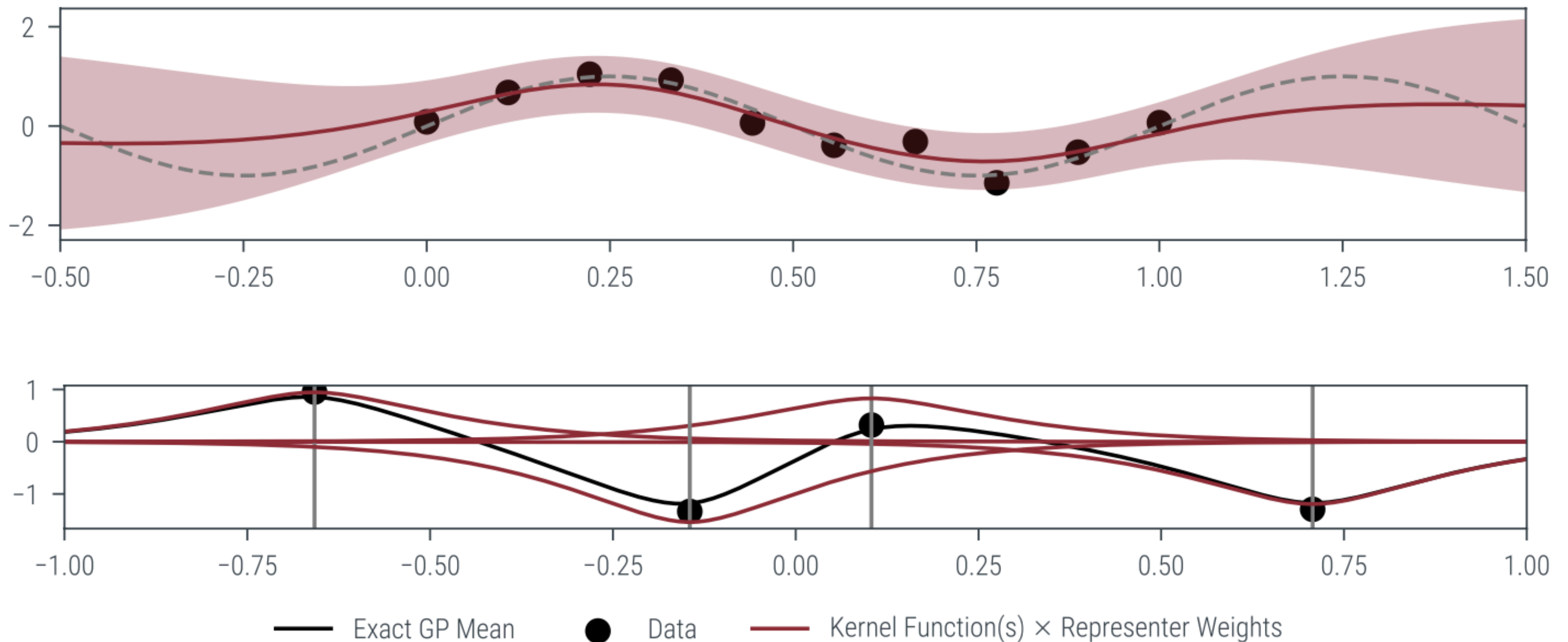
- ▶ The Promise of Probabilistic Machine Learning
- ▶ Gaussian Process Introduction
- ▶ Scaling Gaussian Processes, and Implications
- ▶ Approximate Gaussian Process Inference, The Right Way
- ▶ **iterGP as Probabilistic Numerics**
- ▶ Broader Implications



Return to GP Weight Space View

- Let us again consider the *representer weights* $\mathbf{v}_* = \hat{\mathbf{K}}^{-1}(\mathbf{y} - \boldsymbol{\mu})$

$$\mu_*(\cdot) = \mu(\cdot) + k(\cdot, \mathbf{X}) \overbrace{\hat{\mathbf{K}}^{-1}(\mathbf{y} - \boldsymbol{\mu})}^{\mathbf{v}_*}$$



Learning Representer Weights

- ▶ Consider the GP latent, conditioned on the representer weights

$$p(\mathbf{f}_\diamond \mid \mathbf{v}_*) = \mathcal{N}(\mu(\mathbf{X}_\diamond) + k(\mathbf{X}_\diamond, \mathbf{X})\mathbf{v}_*, k_*(\mathbf{X}_\diamond, \mathbf{X}_\diamond))$$

- ▶ First: suppose representer weights are known (linear solve is done!), we recover **posterior**:

$$\mathbf{f}_\diamond \sim \mathcal{N}(\mu_*(\mathbf{X}_\diamond), k_*(\mathbf{X}_\diamond, \mathbf{X}_\diamond))$$

$$\mu_*(\cdot) = \mu(\cdot) + k(\cdot, \mathbf{X}) \overbrace{\hat{\mathbf{K}}^{-1}(\mathbf{y} - \boldsymbol{\mu})}^{\mathbf{v}_*} \quad k_*(\cdot, \cdot) = k(\cdot, \cdot) - k(\cdot, \mathbf{X})\hat{\mathbf{K}}^{-1}k(\mathbf{X}, \cdot)$$

- ▶ Second: suppose representer weights are unknown (no solve yet!), we recover **prior**:

$$p(\mathbf{v}_*) = \mathcal{N}(\mathbf{v}_*; \mathbf{0}, \hat{\mathbf{K}}^{-1}) \quad \int p(\mathbf{f}_\diamond \mid \mathbf{v}_*)p(\mathbf{v}_*) d\mathbf{v}_*$$

$$\mathbf{f}_\diamond \sim \mathcal{N}(\mu(\mathbf{X}_\diamond), k(\mathbf{X}_\diamond, \mathbf{X}_\diamond))$$

- ▶ This should already feel promising...

Learning Representer Weights (the details)

► Now the numerical method amounts to iteratively updating belief on representer weights:

► Assume an existing belief: $p(\mathbf{v}_*) = \mathcal{N}(\mathbf{v}_*; \mathbf{v}_{i-1}, \Sigma_{i-1})$

► The error/residual of that belief is: $\mathbf{r}_{i-1} = (\mathbf{y} - \mu) - \hat{\mathbf{K}} \mathbf{v}_{i-1}$

► ...this makes sense; the true representer weights are given by $\overbrace{\hat{\mathbf{K}}^{-1}(\mathbf{y} - \mu)}^{\mathbf{v}_*}$.

► Return of actions: $\alpha_i := \mathbf{s}_i^\top \mathbf{r}_{i-1} = \mathbf{s}_i^\top ((\mathbf{y} - \mu) - \hat{\mathbf{K}} \mathbf{v}_{i-1}) = \mathbf{s}_i^\top \hat{\mathbf{K}} (\mathbf{v}_* - \mathbf{v}_{i-1})$

► Conditioning on this projected residual, we result in: $p(\mathbf{v}_*) = \mathcal{N}(\mathbf{v}_* \mid \mathbf{v}_i, \Sigma_i)$

$$\mathbf{v}_i = \mathbf{v}_{i-1} + \underbrace{\Sigma_{i-1} \hat{\mathbf{K}} \mathbf{s}_i}_{=: \mathbf{d}_i} \underbrace{(\mathbf{s}_i^\top \hat{\mathbf{K}} \Sigma_{i-1} \hat{\mathbf{K}} \mathbf{s}_i)^{-1}}_{=: \eta_i} \underbrace{\mathbf{s}_i^\top \hat{\mathbf{K}} (\mathbf{v}_* - \mathbf{v}_{i-1})}_{=: \alpha_i} = \mathbf{C}_i (\mathbf{y} - \mu)$$

$$\Sigma_i = \Sigma_{i-1} - \underbrace{\Sigma_{i-1} \hat{\mathbf{K}} \mathbf{s}_i}_{=: \mathbf{d}_i} \underbrace{(\mathbf{s}_i^\top \hat{\mathbf{K}} \Sigma_{i-1} \hat{\mathbf{K}} \mathbf{s}_i)^{-1}}_{=: \eta_i} \underbrace{\mathbf{s}_i^\top \hat{\mathbf{K}} \Sigma_{i-1}}_{=: \mathbf{d}_i^\top} = \hat{\mathbf{K}}^{-1} - \mathbf{C}_i.$$

$$\mathbf{C}_i := \sum_{j=1}^i \frac{1}{\eta_j} \mathbf{d}_j \mathbf{d}_j^\top = \mathbf{S}_i (\mathbf{S}_i^\top \hat{\mathbf{K}} \mathbf{S}_i)^{-1} \mathbf{S}_i^\top$$

Learning Representer Weights (the details)

► NB this all is captured quite cleanly in an iterative (probabilistic) numerical method

► Accompanying theorems in the paper add considerable strength to these claims

► Matrix vector multiplies imply a computational cost of $\mathcal{O}(n^2 i)$

► Storage is linear as the full covariance matrix needs not be represented in memory

Algorithm 1: A Class of Computation-Aware Iterative Methods for GP Approximation

Input: prior mean function μ , prior covariance function / kernel k , training inputs \mathbf{X} , labels \mathbf{y}

Output: (combined) GP posterior $\mathcal{GP}(\mu_i, k_i)$

```

1  procedure ITERGP( $\mu, k, \mathbf{X}, \mathbf{y}$ )
2     $(\mu_0, k_0) \leftarrow (\mu, k)$                                 ▷ Initialize mean and covariance function with prior.
3     $\boldsymbol{\mu} \leftarrow \mu(\mathbf{X})$                                 ▷ Prior predictive mean.
4     $\hat{\mathbf{K}} \leftarrow k(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I}$                     ▷ Prior predictive kernel matrix.
5    while not STOPPINGCRITERION() do                        ▷ Stopping criterion.
6       $\mathbf{s}_i \leftarrow \text{POLICY}()$                                 ▷ Select action via policy (see Table 1 for examples).
7       $\mathbf{r}_{i-1} \leftarrow (\mathbf{y} - \boldsymbol{\mu}) - \hat{\mathbf{K}} \mathbf{v}_{i-1}$                 ▷ Predictive residual.
8       $\alpha_i \leftarrow \mathbf{s}_i^\top \mathbf{r}_{i-1}$                             ▷ Observation via information operator.
9       $\mathbf{d}_i \leftarrow \boldsymbol{\Sigma}_{i-1} \hat{\mathbf{K}} \mathbf{s}_i = (\mathbf{I} - \mathbf{C}_{i-1} \hat{\mathbf{K}}) \mathbf{s}_i$     ▷ Search direction.
10      $\eta_i \leftarrow \mathbf{s}_i^\top \hat{\mathbf{K}} \boldsymbol{\Sigma}_{i-1} \hat{\mathbf{K}} \mathbf{s}_i = \mathbf{s}_i^\top \hat{\mathbf{K}} \mathbf{d}_i$     ▷ Normalization constant.
11      $\mathbf{C}_i \leftarrow \mathbf{C}_{i-1} + \frac{1}{\eta_i} \mathbf{d}_i \mathbf{d}_i^\top$                 ▷ Precision matrix approximation  $\mathbf{C}_i \approx \hat{\mathbf{K}}^{-1}$ .
12      $\mathbf{Q}_i \leftarrow \mathbf{Q}_{i-1} + \frac{1}{\eta_i} \hat{\mathbf{K}} \mathbf{d}_i \mathbf{d}_i^\top \hat{\mathbf{K}}$         ▷ Kernel matrix approximation  $\mathbf{Q}_i \approx \hat{\mathbf{K}}$ .
13      $\mathbf{v}_i \leftarrow \mathbf{v}_{i-1} + \frac{\alpha_i}{\eta_i} \mathbf{d}_i$                     ▷ Representer weights estimate.
14      $\boldsymbol{\Sigma}_i \leftarrow \boldsymbol{\Sigma}_0 - \mathbf{C}_i$                         ▷ Computational representer weights uncertainty.
15      $p(\mathbf{v}_*) \leftarrow \mathcal{N}(\mathbf{v}_*; \mathbf{v}_i, \boldsymbol{\Sigma}_i)$                 ▷ Belief about representer weights.
16      $\mu_i(\cdot) \leftarrow \mu(\cdot) + k(\cdot, \mathbf{X}) \mathbf{v}_i$                 ▷ Approximate posterior mean function.
17      $k_i(\cdot, \cdot) \leftarrow k(\cdot, \cdot) - k(\cdot, \mathbf{X}) \mathbf{C}_i k(\mathbf{X}, \cdot)$     ▷ Combined uncertainty.
18     return  $\mathcal{GP}(\mu_i, k_i)$ 

```

Greyed out quantities are *not* needed to compute the combined posterior and are only included for clarity of exposition.

[Wenger et al NeurIPS 2022]

Learning Representer Weights

► Notice what has happened:

- Our belief on representer weights captures all computation $p(\mathbf{v}_*) = \mathcal{N}(\mathbf{v}_* \mid \mathbf{v}_i, \Sigma_i)$
- We have the conditional on the latent $p(\mathbf{f}_\diamond \mid \mathbf{v}_*) = \mathcal{N}(\mu(\mathbf{X}_\diamond) + k(\mathbf{X}_\diamond, \mathbf{X})\mathbf{v}_*, k_*(\mathbf{X}_\diamond, \mathbf{X}_\diamond))$
- And voila! $p(\mathbf{f}_\diamond) = \int p(\mathbf{f}_\diamond \mid \mathbf{v}_*)p(\mathbf{v}_*)d\mathbf{v}_* = \mathcal{N}(\mathbf{f}_\diamond; \mu_i(\mathbf{X}_\diamond), k_i(\mathbf{X}_\diamond, \mathbf{X}_\diamond))$

$$\mu_i(\cdot) = \mu(\cdot) + k(\cdot, \mathbf{X})\mathbf{C}_i(\mathbf{y} - \mu)$$

$$k_i(\cdot, \cdot) = \underbrace{k(\cdot, \cdot) - k(\cdot, \mathbf{X})\hat{\mathbf{K}}^{-1}k(\mathbf{X}, \cdot)}_{\text{mathematical uncertainty } \textcolor{blue}{\bullet}} + \underbrace{k(\cdot, \mathbf{X})\Sigma_i k(\mathbf{X}, \cdot)}_{\text{computational uncertainty } \textcolor{green}{\bullet}} = \underbrace{k(\cdot, \cdot) - k(\cdot, \mathbf{X})\mathbf{C}_i k(\mathbf{X}, \cdot)}_{\text{combined uncertainty } \textcolor{purple}{\bullet}}$$

- This is of course exactly the form we ended up with earlier $\mathbf{S}_i^\top \mathbf{y} \mid \mathbf{f} \sim \mathcal{N}(\mathbf{S}_i^\top \mathbf{f}, \sigma^2 \mathbf{S}_i^\top \mathbf{S}_i)$
 - ...paying off the claim that “Effective Dataset” == “Representer Weight Belief Update”

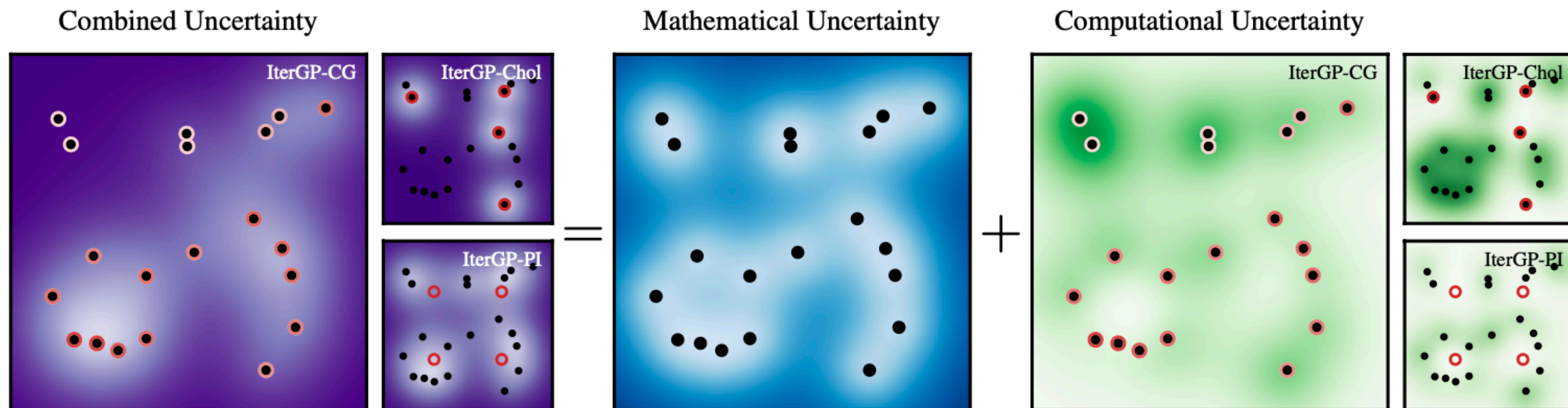
Learning Representer Weights

$$\underbrace{k(\cdot, \cdot) - k(\cdot, \mathbf{X})\mathbf{C}_i k(\mathbf{X}, \cdot)}_{\text{combined uncertainty } \textcolor{purple}{\bullet}} = \underbrace{k(\cdot, \cdot) - k(\cdot, \mathbf{X})\hat{\mathbf{K}}^{-1}k(\mathbf{X}, \cdot)}_{\text{mathematical uncertainty } \textcolor{blue}{\bullet}} + \underbrace{k(\cdot, \mathbf{X})\boldsymbol{\Sigma}_i k(\mathbf{X}, \cdot)}_{\text{computational uncertainty } \textcolor{green}{\bullet}}$$

$$\mathbf{C}_i = \mathbf{S}_i(\mathbf{S}_i^\top \hat{\mathbf{K}} \mathbf{S}_i)^{-1} \mathbf{S}_i^\top$$

$$\hat{\mathbf{K}} := \mathbf{K} + \sigma^2 \mathbf{I}$$

$$p(\mathbf{v}_*) = \mathcal{N}(\mathbf{v}_* \mid \mathbf{v}_i, \boldsymbol{\Sigma}_i)$$



Learning Representer Weights

$$\underbrace{k(\cdot, \cdot) - k(\cdot, \mathbf{X})\mathbf{C}_i k(\mathbf{X}, \cdot)}_{\text{combined uncertainty } \textcolor{purple}{\bullet}} = \underbrace{k(\cdot, \cdot) - k(\cdot, \mathbf{X})\hat{\mathbf{K}}^{-1} k(\mathbf{X}, \cdot)}_{\text{mathematical uncertainty } \textcolor{blue}{\bullet}} + \underbrace{k(\cdot, \mathbf{X})\boldsymbol{\Sigma}_i k(\mathbf{X}, \cdot)}_{\text{computational uncertainty } \textcolor{green}{\bullet}}$$

$$\mathbf{C}_i = \mathbf{S}_i(\mathbf{S}_i^\top \hat{\mathbf{K}} \mathbf{S}_i)^{-1} \mathbf{S}_i^\top$$

$$\hat{\mathbf{K}} := \mathbf{K} + \sigma^2 \mathbf{I}$$

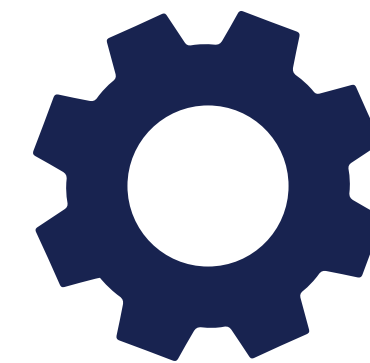
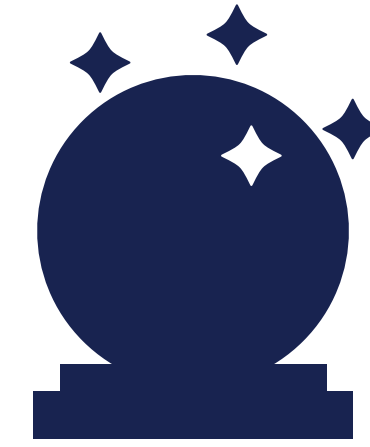
$$p(\mathbf{v}_*) = \mathcal{N}(\mathbf{v}_* \mid \mathbf{v}_i, \boldsymbol{\Sigma}_i)$$

► Parting thoughts:

- Combined uncertainty is tractable and tells us exactly what data we actually consumed
- Mathematical uncertainty is revealed to be conceptual (sure it's there, but only at cubic cost)
- Computational uncertainty then is exactly the uncertainty on the representer weights, which carry all computational updates.
- Here then “Approximate GP Inference the Right Way” and “PN treatment of representer weights” are shown to be **identical**.

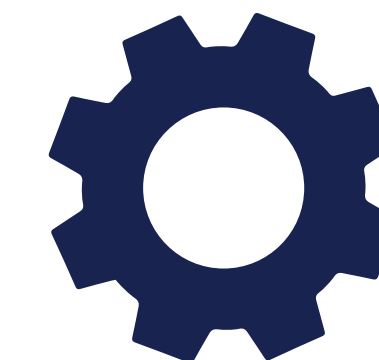
Outline

- ▶ The Promise of Probabilistic Machine Learning
- ▶ Gaussian Process Introduction
- ▶ Scaling Gaussian Processes, and Implications
- ▶ Approximate Gaussian Process Inference, The Right Way
- ▶ iterGP as Probabilistic Numerics
- ▶ **Broader Implications**



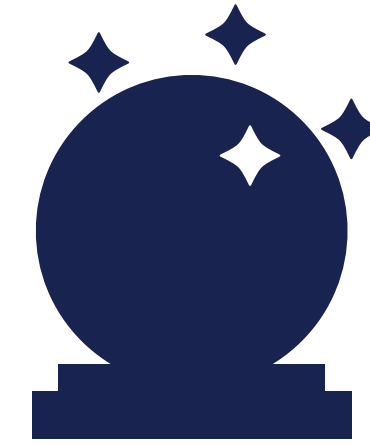
Takeaways

- ▶ Approximate inference **should** be taken into account in the Bayesian framework
 - ▶ iterGP does so for Gaussian Processes
 - ▶ Absent that, computational uncertainty is untracked, and may dominate
- ▶ **Data is as data does:**
 - ▶ In a very concrete sense, iterGP shows that data only “exists” to the extent that it is ingested by the compute/solver (note: this is next-level PN thinking)
 - ▶ Note then for example the practical ease then of iterGP online (actions on new data)
- ▶ **Challenge:** What other inference settings admit (tractable) inference of combined uncertainty?
- ▶ Active learning / BO foreshadow: here is a means to precisely trade off the cost of collecting another data point vs running more compute on existing data.
 - ▶ Combined uncertainty is **exact** regardless of how much compute you do (but of course limited by how much you do)



Thanks

- ▶ Questions?
- ▶ Special thanks to Jonathan, who has led iterGP (and who will be leading the tutorial this afternoon)



- ▶ Jonathan Wenger
- ▶ Geoff Pleiss
- ▶ Luhuan Wu
- ▶ Jacob Gardner
- ▶ Marvin Pförtner
- ▶ Philipp Hennig
- ▶ Dan Biderman
- ▶ Andres Potapczynski
- ▶ Kelly Buchanan
- ▶ Taiga Abe
- ▶ Others...



GATSBY



THE MCKNIGHT FOUNDATION

