

# CONTINUAL LEARNING VIA SEQUENTIAL FUNCTION-SPACE VARIATIONAL INFERENCE



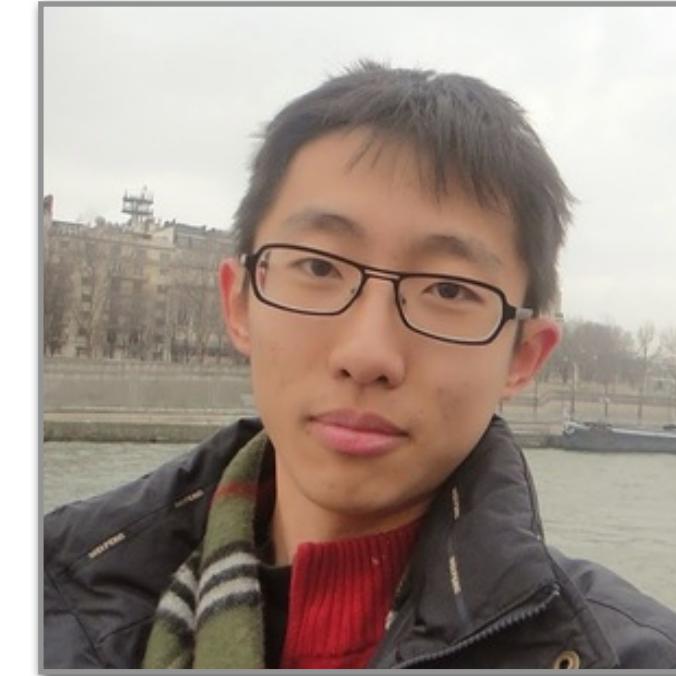
**TIM G. J. RUDNER**

@timrudner



**FREDDIE BICKFORD SMITH**

@fbickfordsmith



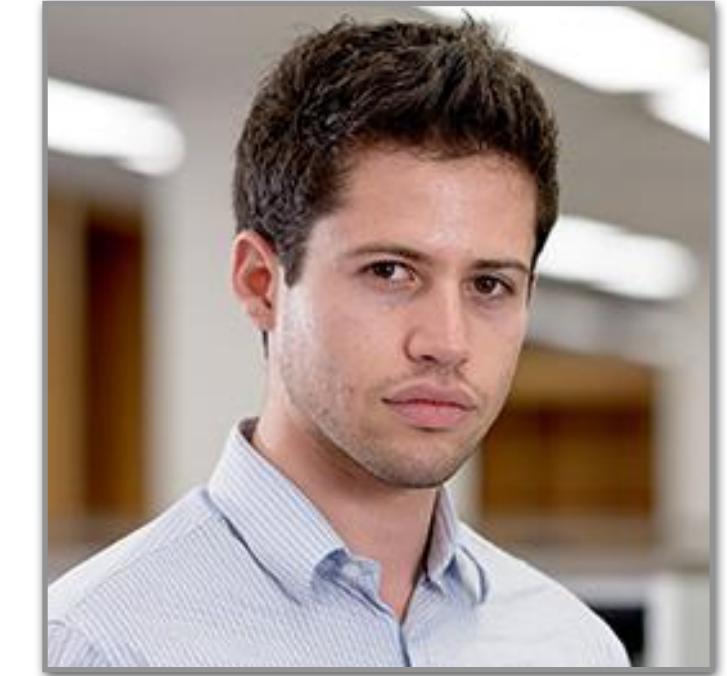
**QIXUAN FENG**

@qixuan\_feng



**YEE WHYE TEH**

@yeewhye



**YARIN GAL**

@yaringal



INTERNATIONAL CONFERENCE ON MACHINE LEARNING 2022

Correspondence to

[tim.rudner@cs.ox.ac.uk](mailto:tim.rudner@cs.ox.ac.uk)

**Paper:** <http://timrudner.com/sfsvi>

**Code:** <http://timrudner.com/sfsvi-code>



# CONTINUAL LEARNING

## Goal:

- **Adapt** to new tasks **without forgetting** previously **attained abilities**.

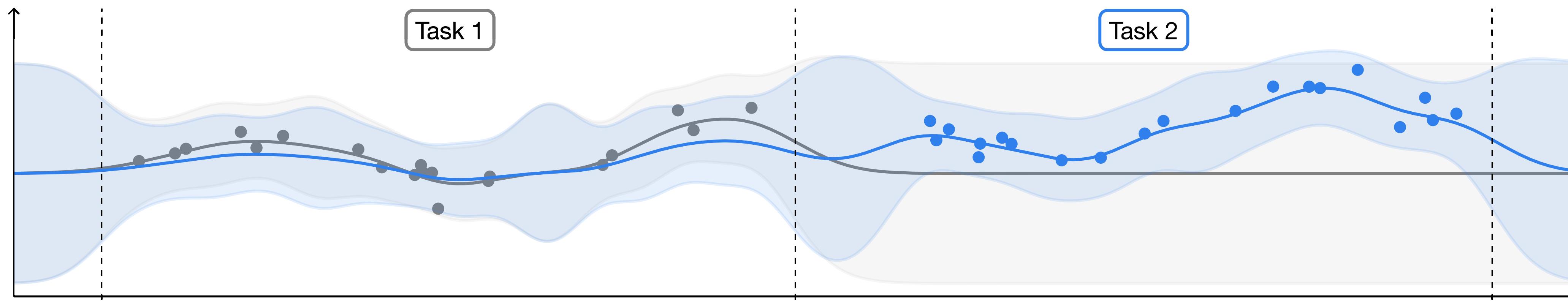
## In practice:

- Learn to make accurate *predictions* for inputs associated with new tasks without “un-learning” to make accurate predictions on inputs associated with previous tasks.

# CONTINUAL LEARNING VIA FUNCTION-SPACE VARIATIONAL INFERENCE

## Key ideas

1. Learning new tasks by selecting “good” functions from a distribution over functions
2. Incorporating prior knowledge via **prior distributions over functions**



# SEQUENTIAL BAYESIAN INFERENCE OVER STOCHASTIC FUNCTIONS

## Sequential Bayesian Inference

- ▶ Posterior distribution over stochastic functions:

$$p(f|\mathcal{D}_1, \dots, \mathcal{D}_t) \propto p(\mathcal{D}_t|f)p(f|\mathcal{D}_1, \dots, \mathcal{D}_{t-1})$$

- ▶ Sequential variational problem:

$$\iff \min_{q_t(f) \in \mathcal{Q}_f} \{\mathbb{D}_{\text{KL}}(q_t(f) \| p_t(f \mid \mathcal{D}_1, \dots, \mathcal{D}_t))\}$$

$$\max_{q_t(f) \in \mathcal{Q}_f} \left\{ \mathbb{E}_{q_t(f)} [\log p(\mathbf{y}_t \mid f(\mathbf{X}_t))] - \mathbb{D}_{\text{KL}}(q_t(f) \| p_t(f \mid \mathcal{D}_1, \dots, \mathcal{D}_{t-1})) \right\}$$

# SEQUENTIAL BAYESIAN INFERENCE OVER STOCHASTIC FUNCTIONS

## Sequential Bayesian Inference

- ▶ Posterior distribution over stochastic functions:

$$p(f|\mathcal{D}_1, \dots, \mathcal{D}_t) \propto p(\mathcal{D}_t|f)p(f|\mathcal{D}_1, \dots, \mathcal{D}_{t-1})$$

- ▶ Sequential variational problem:

$$\iff \min_{q_t(f) \in \mathcal{Q}_f} \{\mathbb{D}_{\text{KL}}(q_t(f) \| p_t(f \mid \mathcal{D}_1, \dots, \mathcal{D}_t))\}$$

$$\max_{q_t(f) \in \mathcal{Q}_f} \left\{ \mathbb{E}_{q_t(f)} [\log p(\mathbf{y}_t \mid f(\mathbf{X}_t))] - \mathbb{D}_{\text{KL}}(q_t(f) \| q_{t-1}(f)) \right\}$$

# BACKGROUND: FUNCTION-SPACE VARIATIONAL INFERENCE

## Function-Space Variational Inference (Rudner et al. [2021])

- **Distributions over functions** induced by distributions over parameters:

$$p(\boldsymbol{\theta}) \Rightarrow p(f(\cdot; \boldsymbol{\theta})) \quad q(\boldsymbol{\theta}) \Rightarrow q(f(\cdot; \boldsymbol{\theta}))$$

- **Variational problem:**

$$\iff$$

$$\min_{q(\boldsymbol{\theta}) \in \mathcal{Q}_{\Theta}} \mathbb{D}_{\text{KL}}(q(f(\cdot; \boldsymbol{\theta})) \| p(f(\cdot; \boldsymbol{\theta}) | \mathcal{D}))$$

$$\max_{q(\boldsymbol{\theta}) \in \mathcal{Q}_{\Theta}} \left\{ \mathbb{E}_{q(f(\mathbf{X}_{\mathcal{D}}; \boldsymbol{\theta}))} [\log p(\mathbf{y} | f(\mathbf{X}_{\mathcal{D}}; \boldsymbol{\theta}))] - \mathbb{D}_{\text{KL}}(q(f(\cdot; \boldsymbol{\theta})) \| p(f(\cdot; \boldsymbol{\theta}))) \right\}$$

# BACKGROUND: FUNCTION-SPACE VARIATIONAL INFERENCE

## Function-Space Variational Inference (Rudner et al. [2021])

- **Distributions over functions** induced by distributions over parameters:

$$p(\boldsymbol{\theta}) \Rightarrow p(f(\cdot; \boldsymbol{\theta})) \quad q(\boldsymbol{\theta}) \Rightarrow q(f(\cdot; \boldsymbol{\theta}))$$

- **Variational problem:**

$$\iff \min_{q(\boldsymbol{\theta}) \in \mathcal{Q}_{\Theta}} \mathbb{D}_{\text{KL}}(q(f(\cdot; \boldsymbol{\theta})) \| p(f(\cdot; \boldsymbol{\theta}) | \mathcal{D}))$$

$$\max_{q(\boldsymbol{\theta}) \in \mathcal{Q}_{\Theta}} \left\{ \mathbb{E}_{q(f(\mathbf{X}_{\mathcal{D}}; \boldsymbol{\theta}))} [\log p(\mathbf{y} | f(\mathbf{X}_{\mathcal{D}}; \boldsymbol{\theta}))] - \sup_{\mathbf{X} \in \mathcal{X}_{\mathbb{N}}} \mathbb{D}_{\text{KL}}(q(f(\mathbf{X}; \boldsymbol{\theta})) \| p(f(\mathbf{X}; \boldsymbol{\theta}))) \right\}$$

(Sun et al. [2019])

# BACKGROUND: FUNCTION-SPACE VARIATIONAL INFERENCE

## Function-Space Variational Inference (Rudner et al. [2021])

- ▶ **KL Approximation:**
  - ▶ Linearize mapping:  $f(\cdot; \Theta) \approx \tilde{f}(\cdot; \Theta) \doteq f(\cdot; \mathbf{m}) + \mathcal{J}_{\mathbf{m}}(\cdot)(\Theta - \mathbf{m})$

with 
$$\mathcal{J}_{\mathbf{m}}(\cdot) \doteq \left. \frac{\partial f(\cdot; \Theta)}{\partial \Theta} \right|_{\Theta=\mathbf{m}}$$

# BACKGROUND: FUNCTION-SPACE VARIATIONAL INFERENCE

## Function-Space Variational Inference (Rudner et al. [2021])

- ▶ **KL Approximation:**
  - ▶ Distributions over functions under linearized mapping

$$q(f(\cdot; \theta)) \approx \tilde{q}(\tilde{f}(\cdot; \theta)) \quad p(f(\cdot; \theta)) \approx \tilde{p}(\tilde{f}(\cdot; \theta))$$

# BACKGROUND: FUNCTION-SPACE VARIATIONAL INFERENCE

## Function-Space Variational Inference (Rudner et al. [2021])

- ▶ **KL Approximation:**
  - ▶ Approximate KL divergence:

$$\sup_{\mathbf{X} \in \mathcal{X}_{\mathbb{N}}} \mathbb{D}_{\text{KL}}(q(f(\mathbf{X}; \boldsymbol{\theta})) \| p(f(\mathbf{X}; \boldsymbol{\theta}))) \approx \sup_{\mathbf{X} \in \mathcal{X}_{\mathbb{N}}} \mathbb{D}_{\text{KL}}(\tilde{q}(\tilde{f}(\mathbf{X}; \boldsymbol{\theta})) \| \tilde{p}(\tilde{f}(\mathbf{X}; \boldsymbol{\theta})))$$

# BACKGROUND: FUNCTION-SPACE VARIATIONAL INFERENCE

## Function-Space Variational Inference (Rudner et al. [2021])

- ▶ **KL Approximation:**

- ▶ Approximate KL divergence:

$$\sup_{\mathbf{X} \in \mathcal{X}_{\mathbb{N}}} \mathbb{D}_{\text{KL}}(q(f(\mathbf{X}; \boldsymbol{\theta})) \| p(f(\mathbf{X}; \boldsymbol{\theta}))) \approx \sup_{\mathbf{X} \in \mathcal{X}_{\mathbb{N}}} \mathbb{D}_{\text{KL}}(\tilde{q}(\tilde{f}(\mathbf{X}; \boldsymbol{\theta})) \| \tilde{p}(\tilde{f}(\mathbf{X}; \boldsymbol{\theta})))$$

- ▶ Supremum Estimation:

$$\sup_{\mathbf{X} \in \mathcal{X}_{\mathbb{N}}} \mathbb{D}_{\text{KL}}(\tilde{q}(\tilde{f}(\mathbf{X}; \boldsymbol{\theta})) \| \tilde{p}(\tilde{f}(\mathbf{X}; \boldsymbol{\theta}))) \approx \max_{\mathbf{X} \in \mathcal{X}_{\mathcal{C}}^S} \mathbb{D}_{\text{KL}}(\tilde{q}(\tilde{f}(\mathbf{X}; \boldsymbol{\theta})) \| \tilde{p}(\tilde{f}(\mathbf{X}; \boldsymbol{\theta})))$$

where  $\mathcal{X}_{\mathcal{C}}^S = \{\mathbf{X}^{(i)}\}_{i=1}^S$  with  $\mathbf{X}^{(i)} = \{\mathbf{x}^{(j)}\}_{j=1}^K$  sampled from  $p_{\mathcal{X}_{\mathcal{C}}}$

# BACKGROUND: FUNCTION-SPACE VARIATIONAL INFERENCE

## Function-Space Variational Inference (Rudner et al. [2021])

- ▶ **Approximate variational problem:**

$$\max_{q(\boldsymbol{\theta}) \in \mathcal{Q}_{\Theta}} \left\{ \mathbb{E}_{q(f(\mathbf{X}_{\mathcal{D}}; \boldsymbol{\theta}))} [\log p(\mathbf{y} \mid f(\mathbf{X}_{\mathcal{D}}; \boldsymbol{\theta}))] - \max_{\mathbf{X} \in \mathcal{X}_{\mathcal{C}}^S} \mathbb{D}_{\text{KL}}(\tilde{q}(\tilde{f}(\mathbf{X}; \boldsymbol{\theta})) \parallel \tilde{p}(\tilde{f}(\mathbf{X}; \boldsymbol{\theta}))) \right\}$$

where  $\mathcal{X}_{\mathcal{C}}^S = \{\mathbf{X}^{(i)}\}_{i=1}^S$  with  $\mathbf{X}^{(i)} = \{\mathbf{x}^{(j)}\}_{j=1}^K$  sampled from  $p_{\mathcal{X}_{\mathcal{C}}}$

# SEQUENTIAL FUNCTION-SPACE VARIATIONAL INFERENCE

## Empirical prior and variational distributions over functions

- ▶ Define prior and variational distributions over parameters:

$$p_t(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\mu}_{t-1}, \boldsymbol{\Sigma}_{t-1})$$

$$q_t(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t)$$

- ▶ Tractable prior and variational distributions over functions:

$$\tilde{p}_t(\tilde{f}(\mathbf{X}; \boldsymbol{\theta})) = \tilde{q}_{t-1}(\tilde{f}(\mathbf{X}; \boldsymbol{\theta})) = \mathcal{N}(f(\mathbf{X}; \boldsymbol{\mu}_{t-1}), \mathcal{J}_{\boldsymbol{\mu}_{t-1}}(\mathbf{X}) \boldsymbol{\Sigma}_{t-1} \mathcal{J}_{\boldsymbol{\mu}_{t-1}}(\mathbf{X}')^\top)$$

$$\tilde{q}_t(\tilde{f}(\mathbf{X}; \boldsymbol{\theta})) = \mathcal{N}(f(\mathbf{X}; \boldsymbol{\mu}_t), \mathcal{J}_{\boldsymbol{\mu}_t}(\mathbf{X}) \boldsymbol{\Sigma}_t \mathcal{J}_{\boldsymbol{\mu}_t}(\mathbf{X}')^\top)$$

# SEQUENTIAL FUNCTION-SPACE VARIATIONAL INFERENCE

## Sequential Function-Space Variational Inference (S-FSVI)

- Approximate function-space variational objective under linearization:

$$\mathbb{E}_{q(f(\mathbf{X}_{\mathcal{D}_t}; \boldsymbol{\theta}))} [\log p(\mathbf{y}_{\mathcal{D}_t} \mid f(\mathbf{X}_{\mathcal{D}_t}; \boldsymbol{\theta}))] - \max_{\mathbf{X} \in \mathcal{X}_{\mathcal{C}}^S} \mathbb{D}_{\text{KL}}(\tilde{q}_t(\tilde{f}(\mathbf{X}; \boldsymbol{\theta})) \parallel \tilde{q}_{t-1}(\tilde{f}(\mathbf{X}; \boldsymbol{\theta})))$$

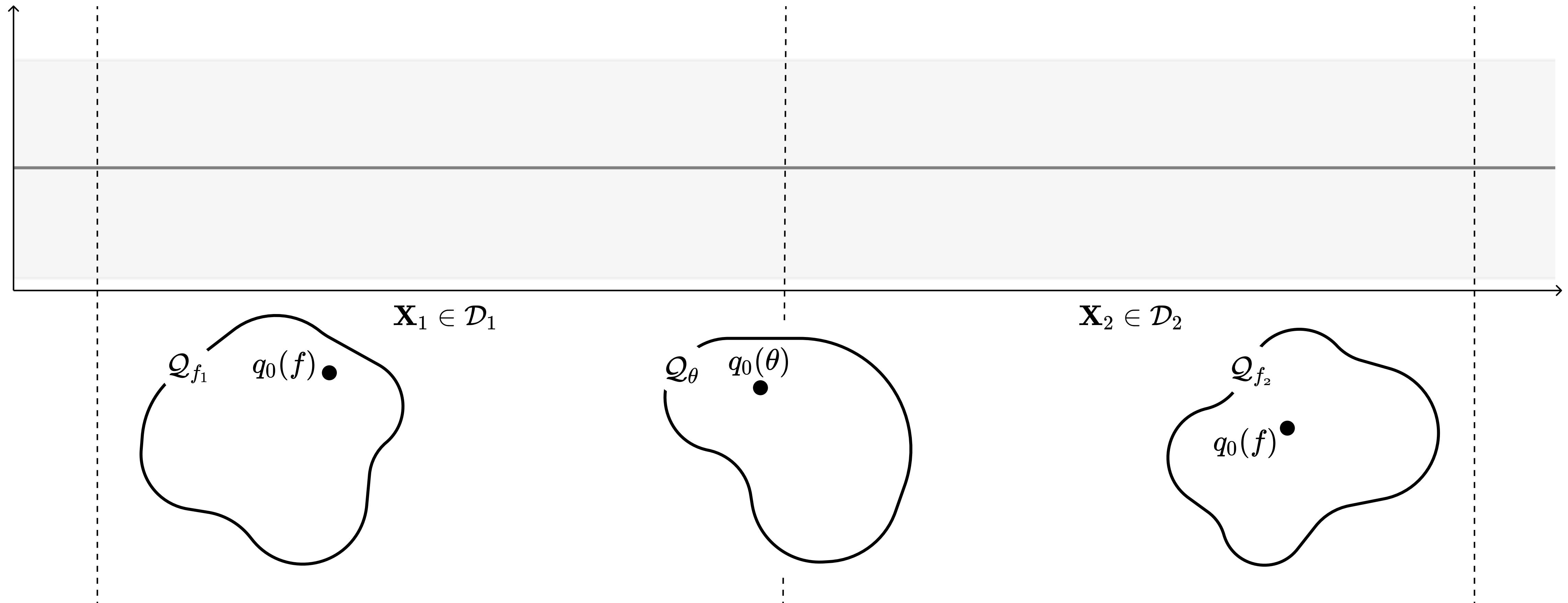
where  $\mathcal{X}_{\mathcal{C}}^S = \{\mathbf{X}^{(i)}\}_{i=1}^S$  with  $\mathbf{X}^{(i)} = \{\mathbf{x}^{(j)}\}_{j=1}^K$  sampled from  $p_{\mathcal{X}_{\mathcal{C}}}$

- Crucially,  $p_{\mathcal{X}_{\mathcal{C}}}$  is defined over a set of **task-specific coresets**

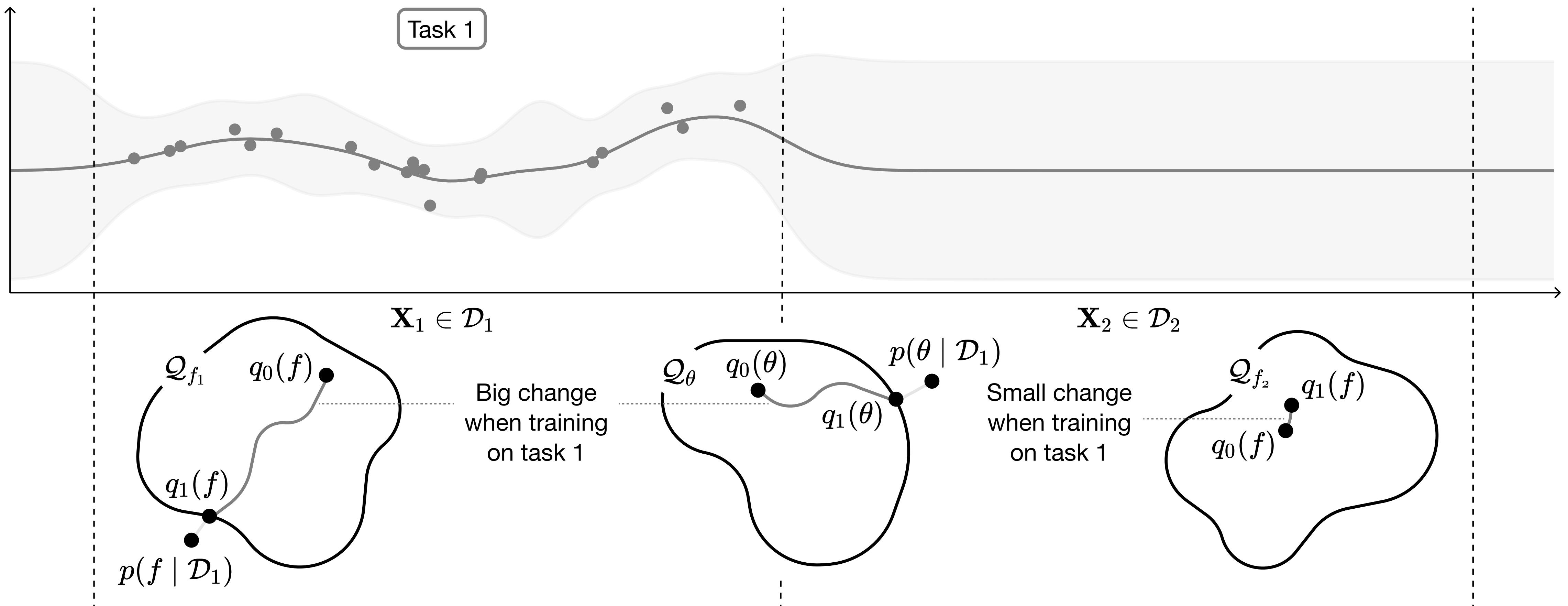
$$\mathcal{X}_{1:t} \subseteq \mathcal{P}\left(\{\mathbf{x}_t^{(n)}\}_{n=1}^{N_t} \cup \bigcup_{t'=1}^{t-1} \{\mathbf{x}_{t'}^{(n)}\}_{n=1}^{S_{t'}}\right)$$

- Explicitly discourages divergence from prior distribution over functions

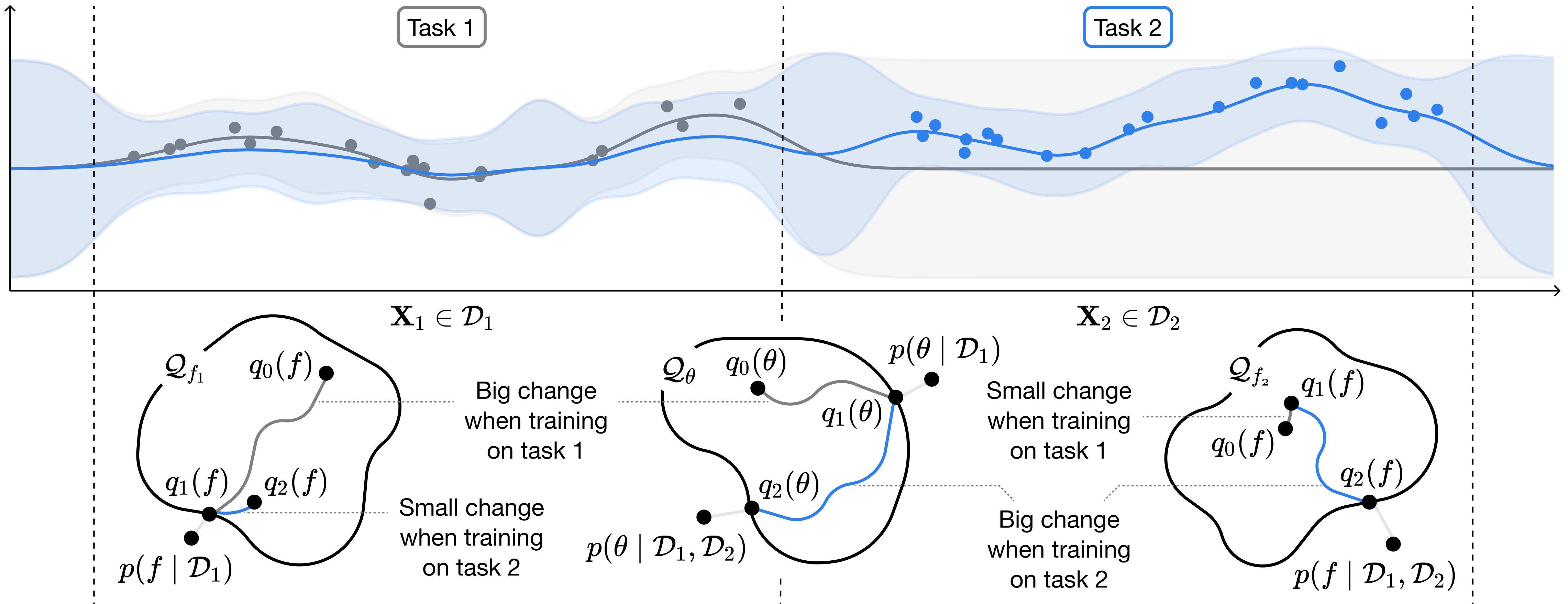
# SEQUENTIAL FUNCTION-SPACE VARIATIONAL INFERENCE



# CONTINUAL LEARNING VIA FUNCTION-SPACE VARIATIONAL INFERENCE



# CONTINUAL LEARNING VIA FUNCTION-SPACE VARIATIONAL INFERENCE

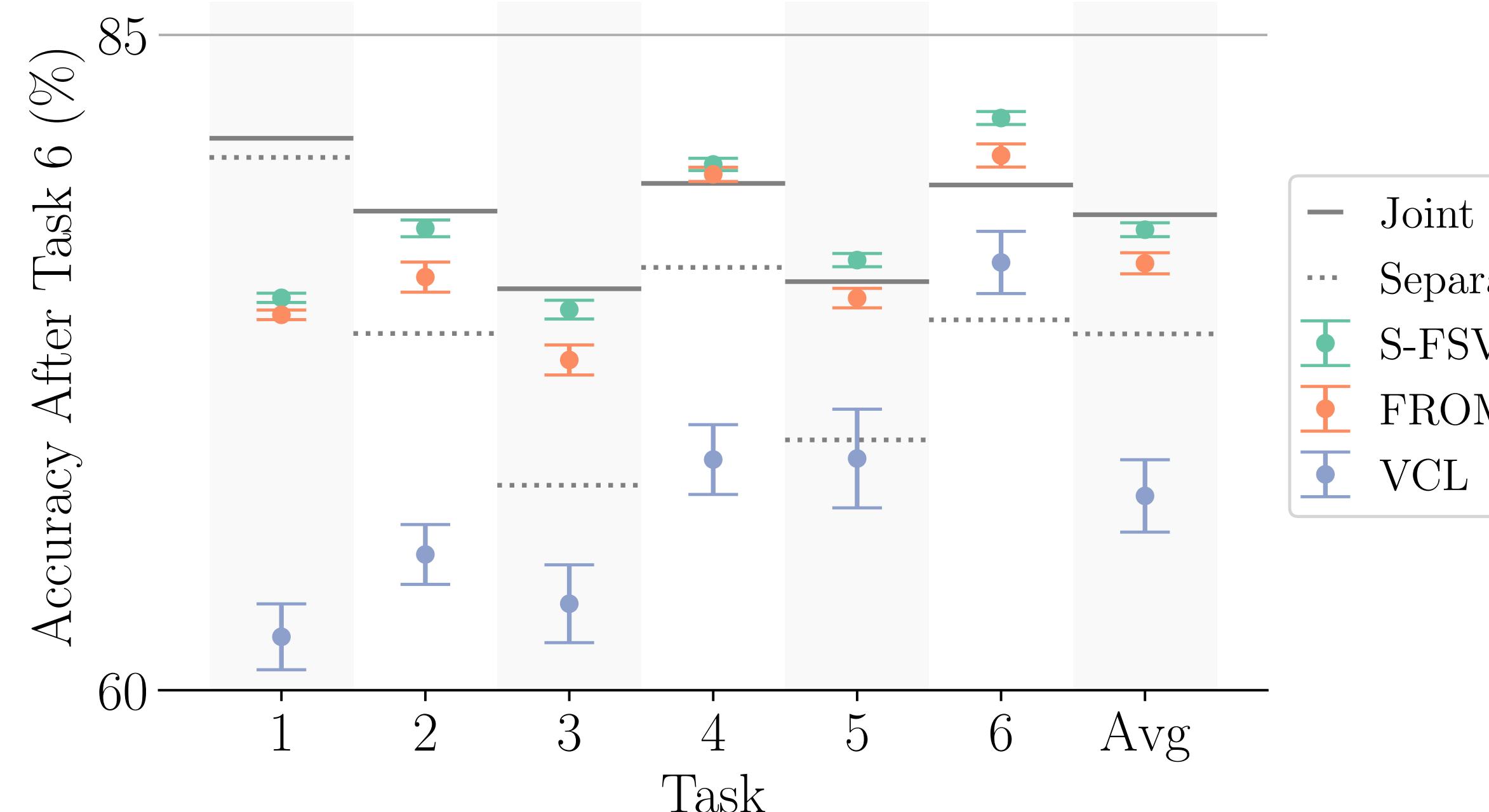


# PREDICTIVE PERFORMANCE: MNIST

Method	S-MNIST (MH)	S-FMNIST (MH)	P-MNIST (SH)	S-MNIST (SH)
EWC (Kirkpatrick et al., 2017)	63.10%	—	84.00%	—
SI (Zenke et al., 2017)	98.90%	—	86.00%	—
VCL (Nguyen et al., 2018) <sup>1</sup>	98.40%	98.60% $\pm$ 0.04	93.00%	32.11% $\pm$ 1.16
VCL (no coresnet)	97.00%	89.60% $\pm$ 1.75	87.50% $\pm$ 0.61	17.74% $\pm$ 1.20
FRCL (Titsias et al., 2020) <sup>3</sup>	97.80% $\pm$ 0.22	97.28% $\pm$ 0.17	94.30% $\pm$ 0.06	—
FROMP (Pan et al., 2020)	99.00% $\pm$ 0.04	99.00% $\pm$ 0.03	94.90% $\pm$ 0.04	35.29% $\pm$ 0.52
VAR-GP (Kapoor et al., 2021)	—	—	97.20% $\pm$ 0.08	90.57% $\pm$ 1.06
S-FSVI (ours) <sup>2</sup>	99.54% $\pm$ 0.04	99.19% $\pm$ 0.02	95.76% $\pm$ 0.02	92.87% $\pm$ 0.14
<b>S-FSVI Ablation Study:</b>				
S-FSVI (larger networks) <sup>4</sup>	99.76% $\pm$ 0.00	99.16% $\pm$ 0.03	97.50% $\pm$ 0.01	93.38% $\pm$ 0.10
S-FSVI (no coresnet) <sup>5</sup>	99.62% $\pm$ 0.02	99.54% $\pm$ 0.01	84.06% $\pm$ 0.46	20.15% $\pm$ 0.52
S-FSVI (minimal coresnet) <sup>6</sup>	—	—	89.59% $\pm$ 0.30	51.44% $\pm$ 1.22

- ▶ Multi-head settings: virtually solved (with and without coresets)
- ▶ Single-head settings: still require coresets

# PREDICTIVE PERFORMANCE: SPLIT CIFAR & SEQUENTIAL OMNIGLOT



Split CIFAR

Method	Test Accuracy
Learning Without Forgetting <sup>1</sup>	62.06% $\pm$ 2.0
EWC	67.32% $\pm$ 4.7
Online EWC <sup>2</sup>	69.99% $\pm$ 3.2
Progress & Compress <sup>3</sup>	70.32% $\pm$ 3.3
FRCL <sup>4</sup>	81.47% $\pm$ 1.6
<b>S-FSVI (ours)<sup>5</sup></b>	<b>83.29%<math>\pm</math>1.2</b>

Sequential Omniglot

# SUMMARY

## Sequential Function-Space Variational Inference

- ▶ Is based on **Bayesian inference over stochastic functions**
- ▶ **Outperforms related** objective- and replay-based **methods**
- ▶ Can be scaled to
  - ▶ **large network architectures**
  - ▶ **long task sequences**

# THANK YOU!

**Correspondence to**

tim.rudner@cs.ox.ac.uk

**Paper:** <http://timrudner.com/sfsvi>

**Code:** <http://timrudner.com/sfsvi-code>



**TIM G. J. RUDNER**  
@timrudner



**FREDDIE BICKFORD SMITH**  
@fbickfordsmith



**QIXUAN FENG**  
@qixuan\_feng



**YEE WHYE TEH**  
@yeewhye



**YARIN GAL**  
@yaringal