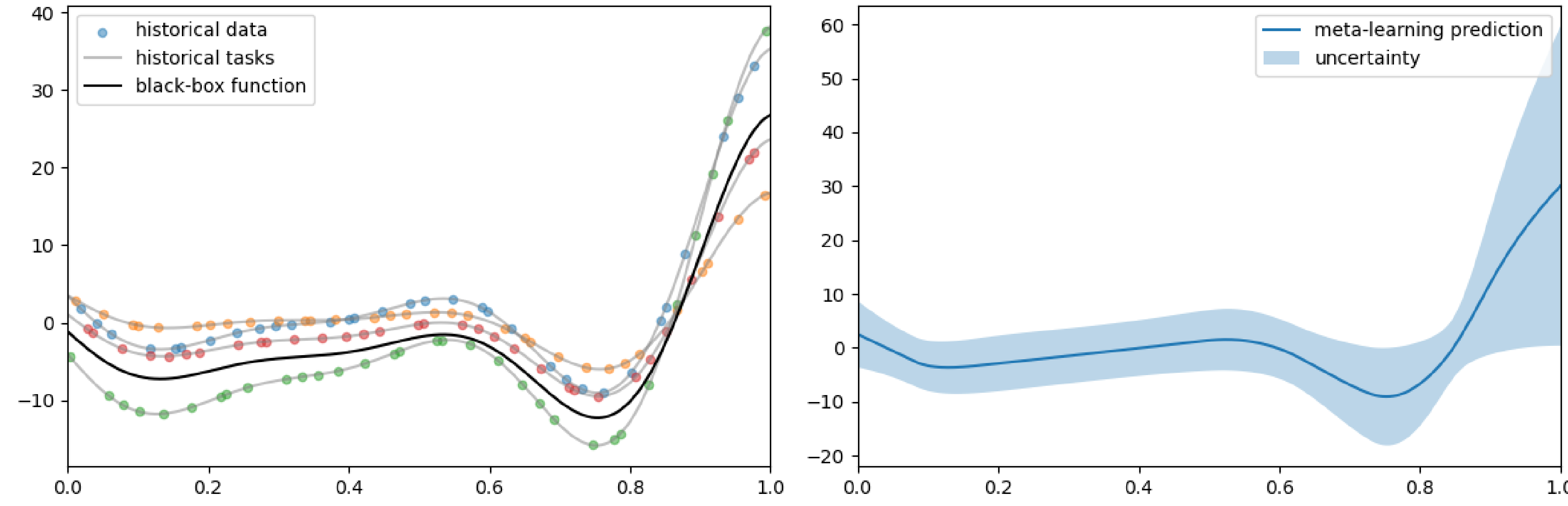


Motivation and Setup

Bayesian optimization (BO) aims to optimize an expensive black-box function:

$$\mathbf{x}^* = \arg \min_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x}) \quad (1)$$

Meta-learning BO leverages information from past optimization experiences to accelerate the current optimization process.

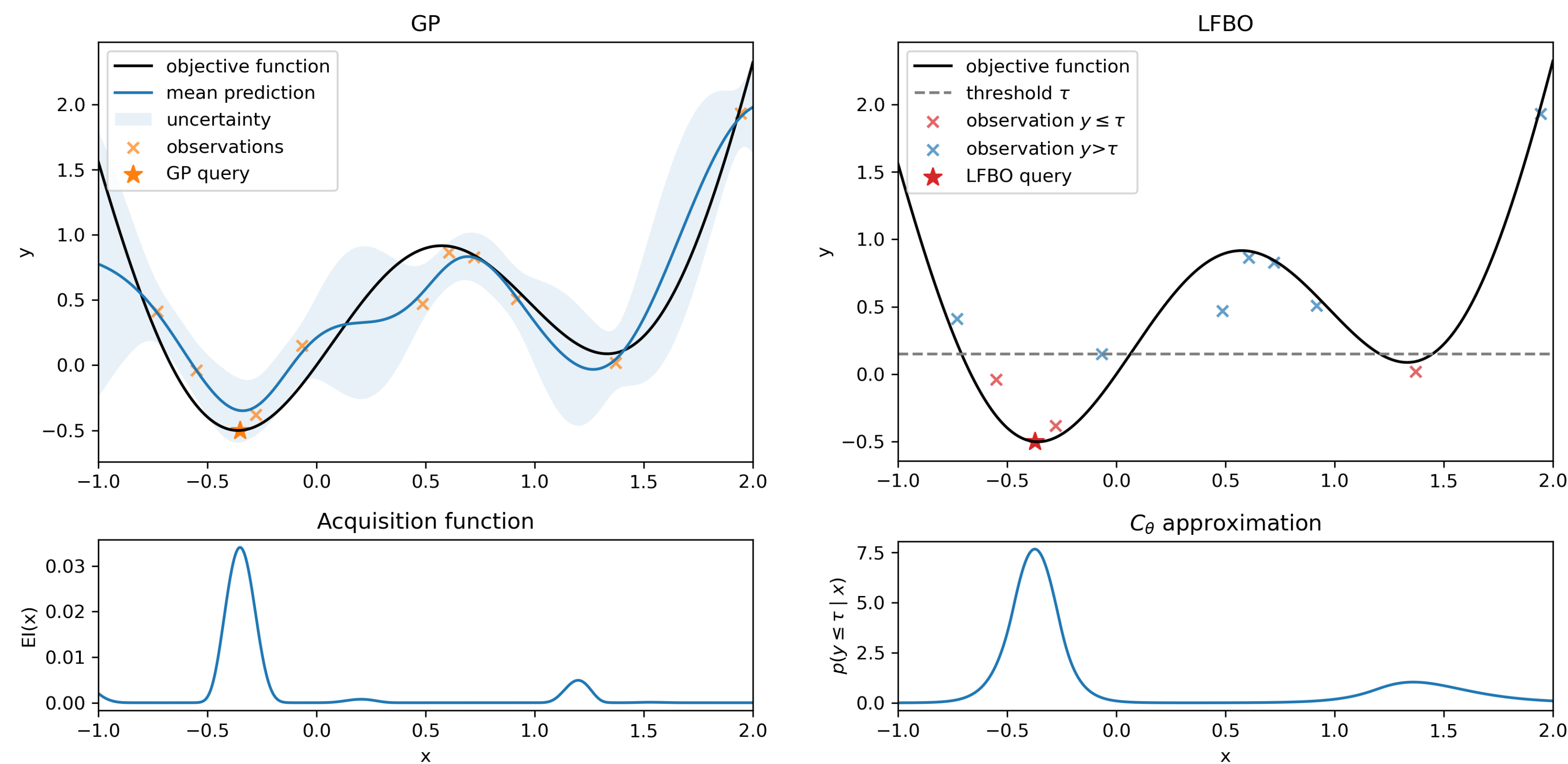


Contribution

- **Directly meta-learn the acquisition function** (optimization strategy) with robustness to heterogeneous scales and noises across tasks.
- Use a **probabilistic task adaptation** to account for task uncertainty.
- Propose a novel adaptation procedure with **residual prediction to ensure robust adaptation**.

Background: Likelihood-free Bayesian optimization (LFBO)

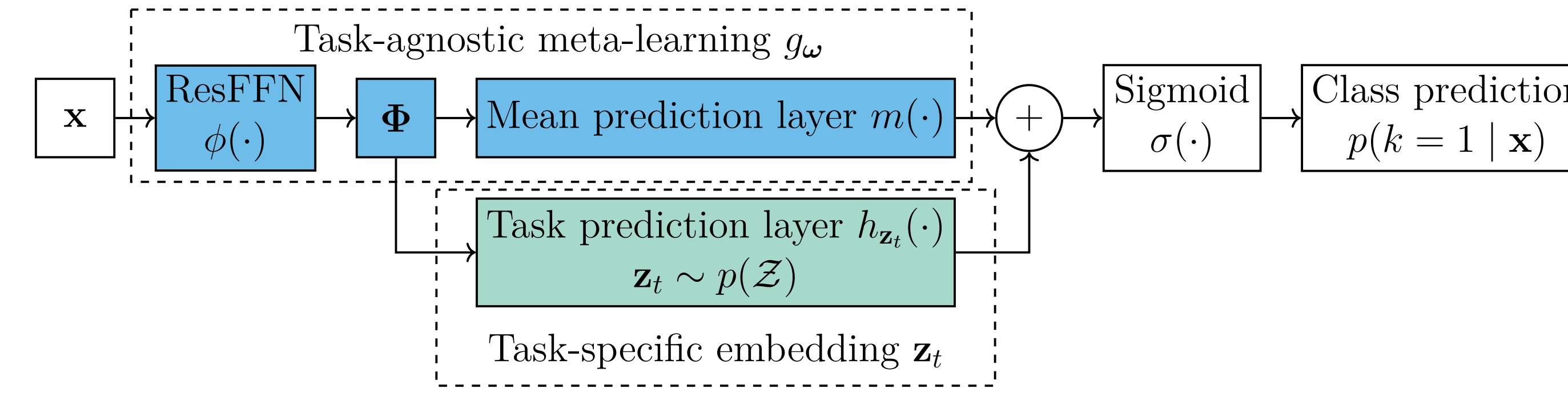
LFBO directly approximate the acquisition function by casting the approximation to a classification problem. The trained classifier C_θ directly output the approximation.



$$\mathcal{L}^{\text{LFBO}}(\theta; \mathcal{D}_N, \tau) = -\mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}_N} [\max(\tau - y, 0) \ln C_\theta(\mathbf{x}) + \ln(1 - C_\theta(\mathbf{x}))]. \quad (2)$$

MALIBO

Meta-learning acquisition function



The meta-learning optimize the LFBO loss across tasks:

$$\mathcal{L}^{\text{meta}}(\omega, \{\mathbf{z}_t\}_{t=1}^T) = \frac{1}{T} \sum_{t=1}^T \mathcal{L}^{\text{LFBO}}(\omega, \mathbf{z}_t; \mathcal{D}^t, \tau) + \lambda \mathcal{R}(\{\mathbf{z}_t\}_{t=1}^T; p(\mathcal{Z})) \quad (3)$$

Directly optimizing $\mathcal{L}^{\text{LFBO}}$ leads to unreliable task adaptation. We regularize \mathcal{Z} to follow a Gaussian prior with \mathcal{R} for Bayesian task adaptation.

Probabilistic task adaptation

We need the posterior of task embedding $p(\mathbf{z} | \mathcal{D}_N)$ for prediction:

$$C(\mathbf{x}; \omega, \mathcal{D}_N) = \int p(k=1 | \omega, \mathbf{z}) p(\mathbf{z} | \mathcal{D}_N) d\mathbf{z} \quad (4)$$

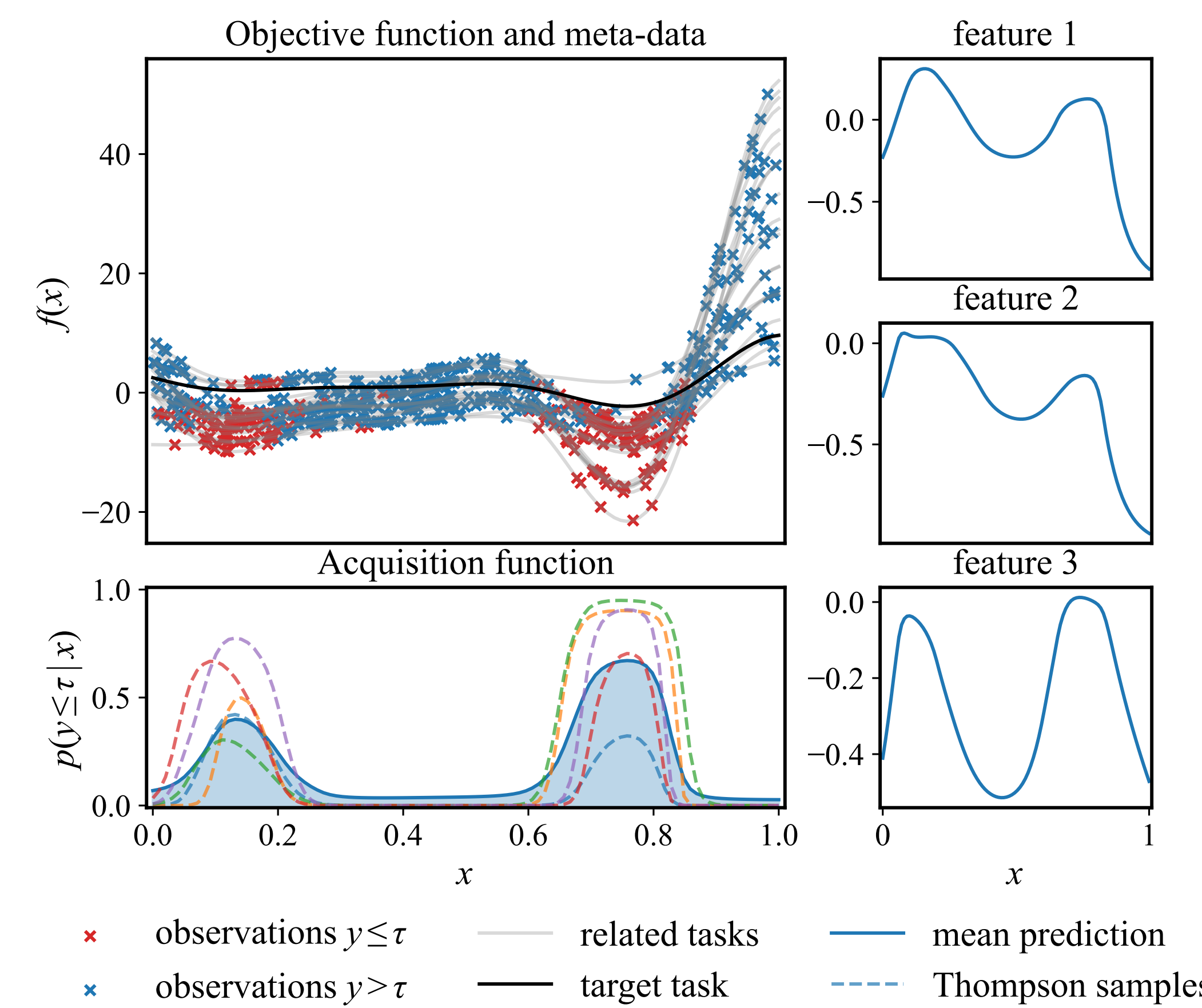
Obtaining the posterior is intractable due to $\sigma(\cdot)$ and the marginal. We use Laplace approximation to obtain $p(\mathbf{z} | \mathcal{D}_N) \simeq q(\mathbf{z}) = \mathcal{N}(\mathbf{z} | \mathbf{z}_{\text{MAP}}, \Sigma_N)$:

$$\mathbf{z}_{\text{MAP}} = \arg \min_{\mathbf{z}} \frac{1}{2} \mathbf{z}^T \mathbf{z} - \sum_{n=1}^N (k_n(\tau - y) \ln \hat{k}_n + \ln(1 - \hat{k}_n)) \quad (5)$$

$$\Sigma_N^{-1} = \Sigma_0^{-1} + \sum_{n=1}^N (k_n(\tau - y) + 1) \hat{k}_n (1 - \hat{k}_n) \Phi_n \Phi_n^T, \quad (6)$$

The output is a sample of the approximated posterior.

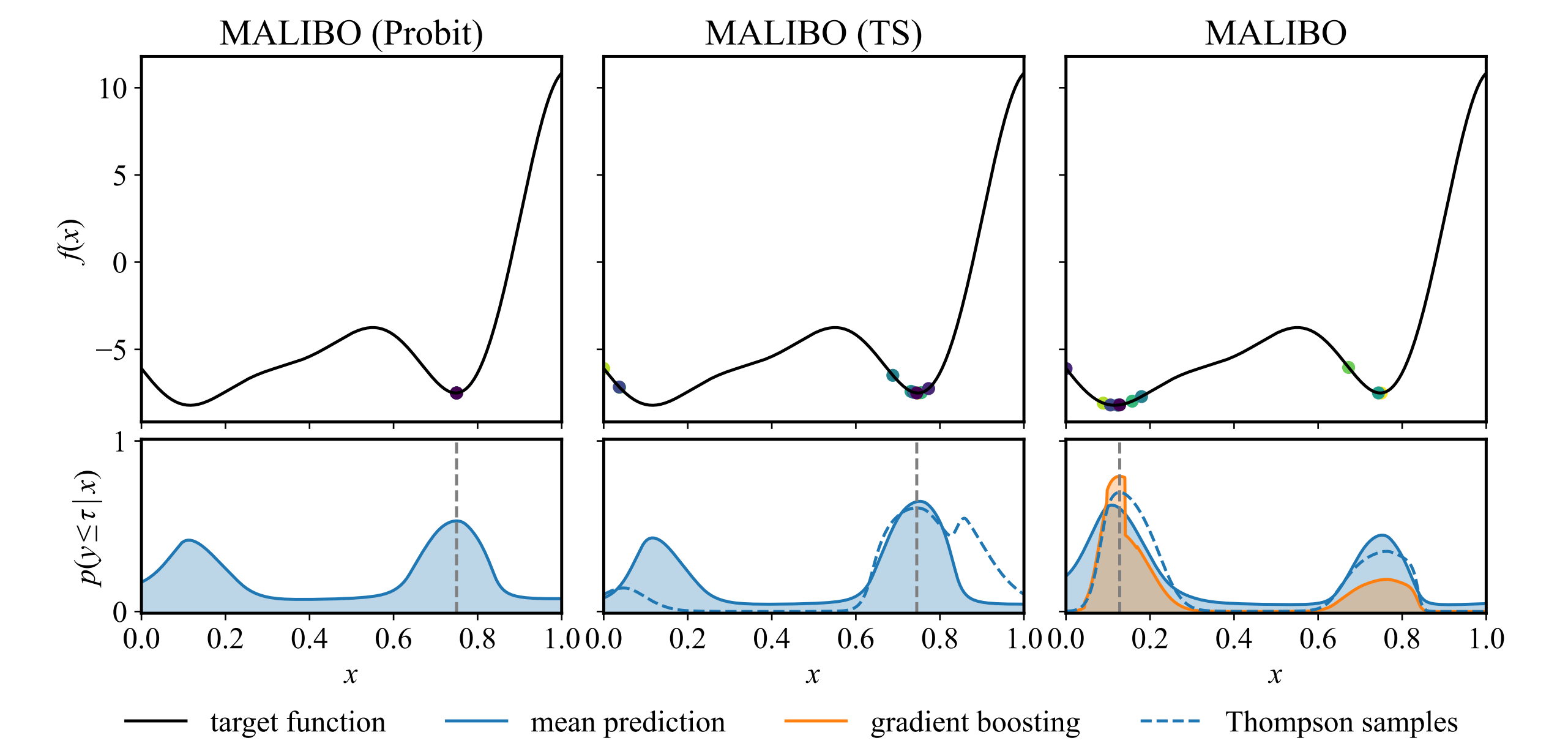
$$\hat{C}(\mathbf{x}) = \sigma(m(\phi(\mathbf{x})) + h_{\hat{\mathbf{z}}}(\phi(\mathbf{x}))), \quad \hat{\mathbf{z}} \sim q(\mathbf{z}) \quad (7)$$



Robust adaptation via residual prediction

To overcome the distribution shift between meta-data and the non-i.i.d. optimization data, we use a Gradient Boosting tree (GB) to model the residual between the model prediction and the black-box function.

$$C_{\text{GB}}(\mathbf{x}) = \sigma \left(m(\phi(\mathbf{x})) + h_{\hat{\mathbf{z}}}(\phi(\mathbf{x})) + \sum_{i=1}^M r_i(\mathbf{x}) \right) \quad (8)$$



Experiments

