

Motivation

- Safety of deployed machine learning models is highly important in many applications.
- Data is costly to obtain in small sample settings such as in engineering or medical applications.
- Identify valid subdomains of input space with a model error smaller than some required tolerance.

a) Background & Problem

- Validate a model $f_M: \mathbb{X} \rightarrow \mathbb{Y}$ over $\mathbb{X} \subset \mathbb{R}^d$
- Expensive observations $Y_{\mathbf{x}} = f_E(\mathbf{x}) + \epsilon$ subject to homoscedastic Gaussian noise ϵ
- **Validation Metric.**

$$P(-\xi < f_D(\mathbf{x}) < \xi),$$

with tolerance $\xi \in \mathbb{R}_{>0}$ and model discrepancy $f_D(\mathbf{x}) := f_M(\mathbf{x}) - Y_{\mathbf{x}}$.

- Reformulation: $P(g(\mathbf{x}) > 0)$ with *limit state function* $g(\mathbf{x}) := \xi - |f_D|$

Definitions

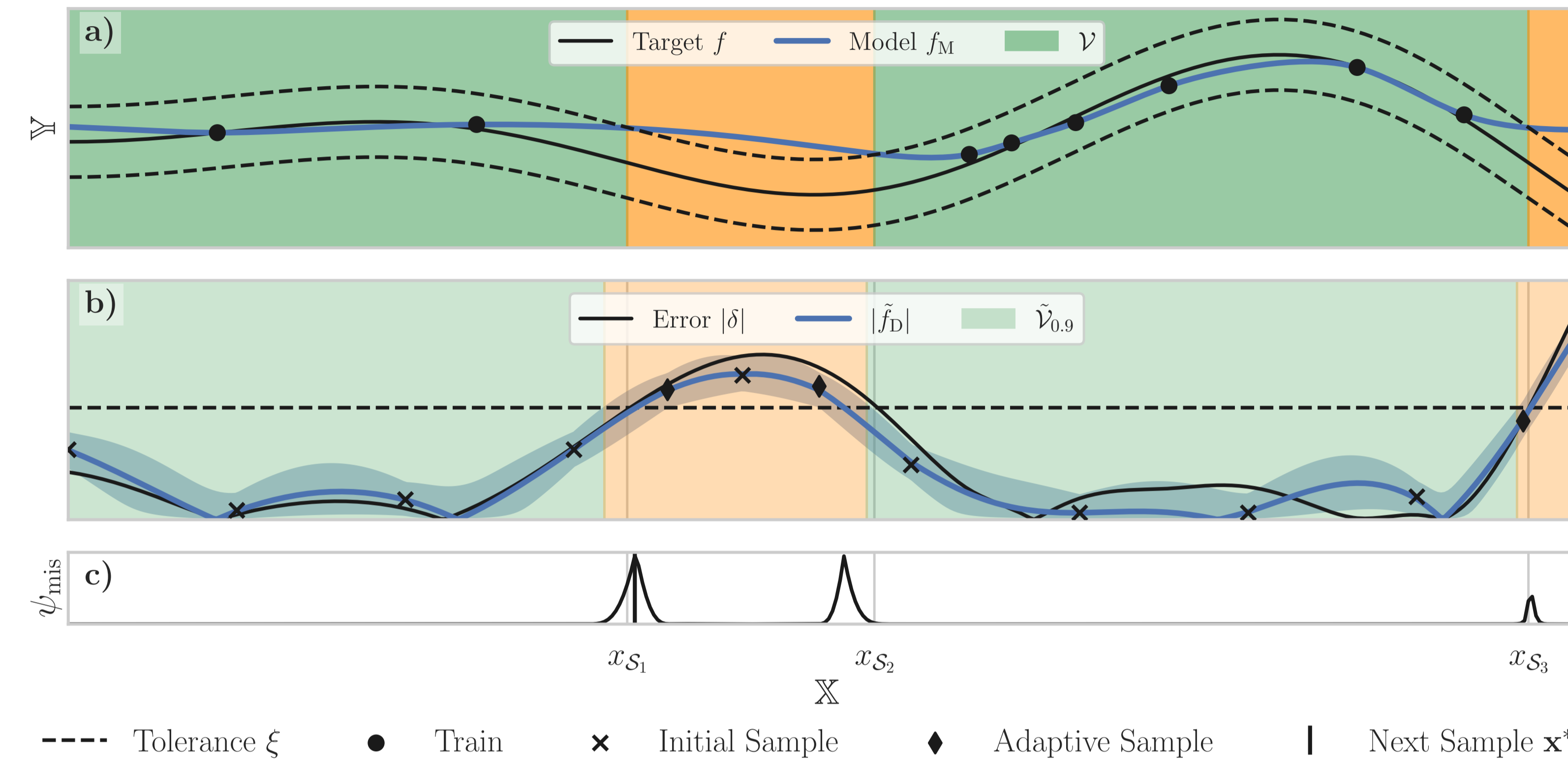
Local Validity. A model f_M is locally valid at \mathbf{x} , given a tolerance level ξ , if $\xi - |\delta(\mathbf{x})| \geq 0$. Then, the valid region of f_M is

$$\mathcal{V} = \{\mathbf{x} \in \mathbb{X}: \xi - |\delta(\mathbf{x})| \geq 0\},$$

with noiseless discrepancy $\delta(\mathbf{x}) = f_M(\mathbf{x}) - f_E(\mathbf{x})$.

Limit State. The limit state of f_M is given by

$$\mathcal{S} = \{\mathbf{x} \in \mathbb{X}: \xi - |\delta(\mathbf{x})| = 0\}.$$



b) Gaussian Process (GP) Error Model

Use a Gaussian process (GP) to learn the limit state

$$\hat{g} = \xi - |\tilde{f}_D| \quad \tilde{f}_D \sim \mathcal{GP}(\mu, k).$$

The prediction is a folded Gaussian posterior, available in closed-form.

c) Learning the Limit State with MC-Prob.

Bayesian Active Learning. A new query \mathbf{x}^* for evaluation is obtained as

$$\mathbf{x}^* = \arg \max_{\mathbf{x} \in \mathcal{C}} \psi_{\text{mis}}(\mathbf{x}),$$

with candidates \mathcal{C} .

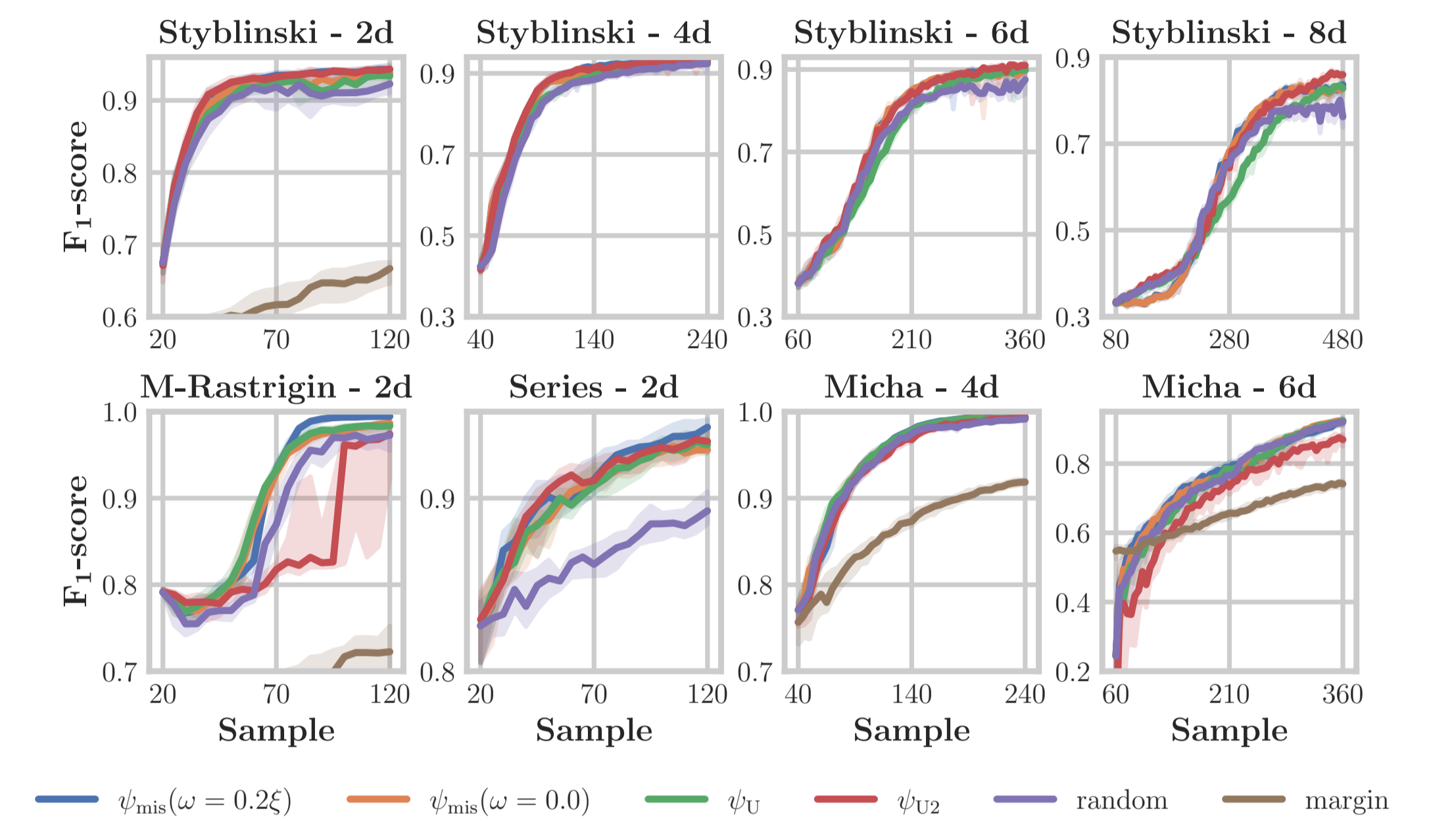
Acquisition Function. We use the misclassification probability (MC-Prob.) as acquisition function

$$\psi_{\text{mis}}(\mathbf{x}; \omega) = \begin{cases} P(\hat{G}_{\mathbf{x}} \leq -\omega), & \text{for } |\mu_{y|\mathcal{D}}(\mathbf{x})| \leq \xi \\ 1 - P(\hat{G}_{\mathbf{x}} \leq \omega), & \text{for } |\mu_{y|\mathcal{D}}(\mathbf{x})| > \xi, \end{cases}$$

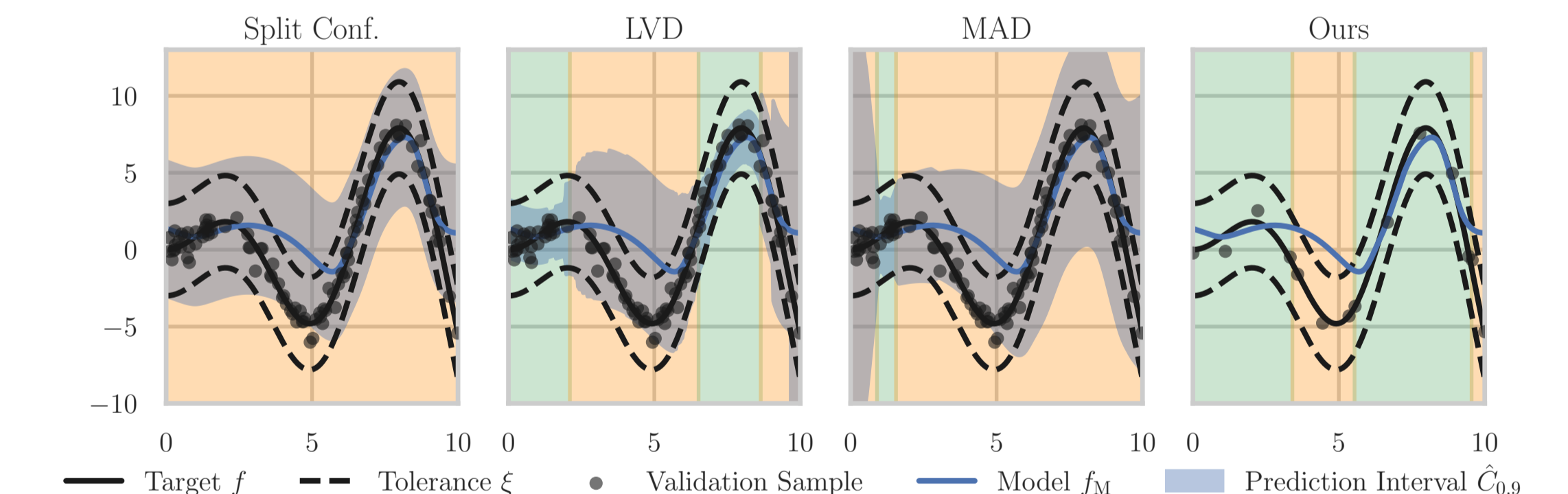
with hyperparameter $\omega \in \mathbb{R}_+$ to control the exploration-exploitation trade-off.

Experiments

Benchmark results



Comparison with conformal prediction



Conclusion

- Novel formulation for local validation, inspired by active learning reliability
- Misclassification probability (MC-Prob) based on epistemic uncertainty is used
- Higher sample efficiency and thus accuracy in limited sample settings compared to previous work

Paper
+ Code

