# ProBO: Versatile Bayesian Optimization Using Any Probabilistic Programming Language

## W. Neiswanger et al. 2019

Paper review by Wanru Zhao

# Bayesian Optimization (BO)

Consider a 'well behaved' function $f : \mathcal{X} \to \mathbb{R}$ where $\mathcal{X} \subseteq \mathbb{R}^D$ is a bounded domain.

$$x_M = \arg\min_{x \in \mathcal{X}} f(x).$$

- Expensive to evaluate

- Black box

- Derivative-free

- (Maybe) noisy

# Bayesian Optimization (BO)

**Flow**

- Methodology to perform global optimisation of multimodal black-box functions.

- 1. Choose some prior measure over the space of possible objectives $f$.

- 2. Combine prior and the likelihood to get a posterior measure over the objective given some observations.

- 3. Use the posterior to decide where to take the next evaluation according to some acquisition function.

- 4. Augment the data.

- Iterate between 2 and 4 until the evaluation budget is over.
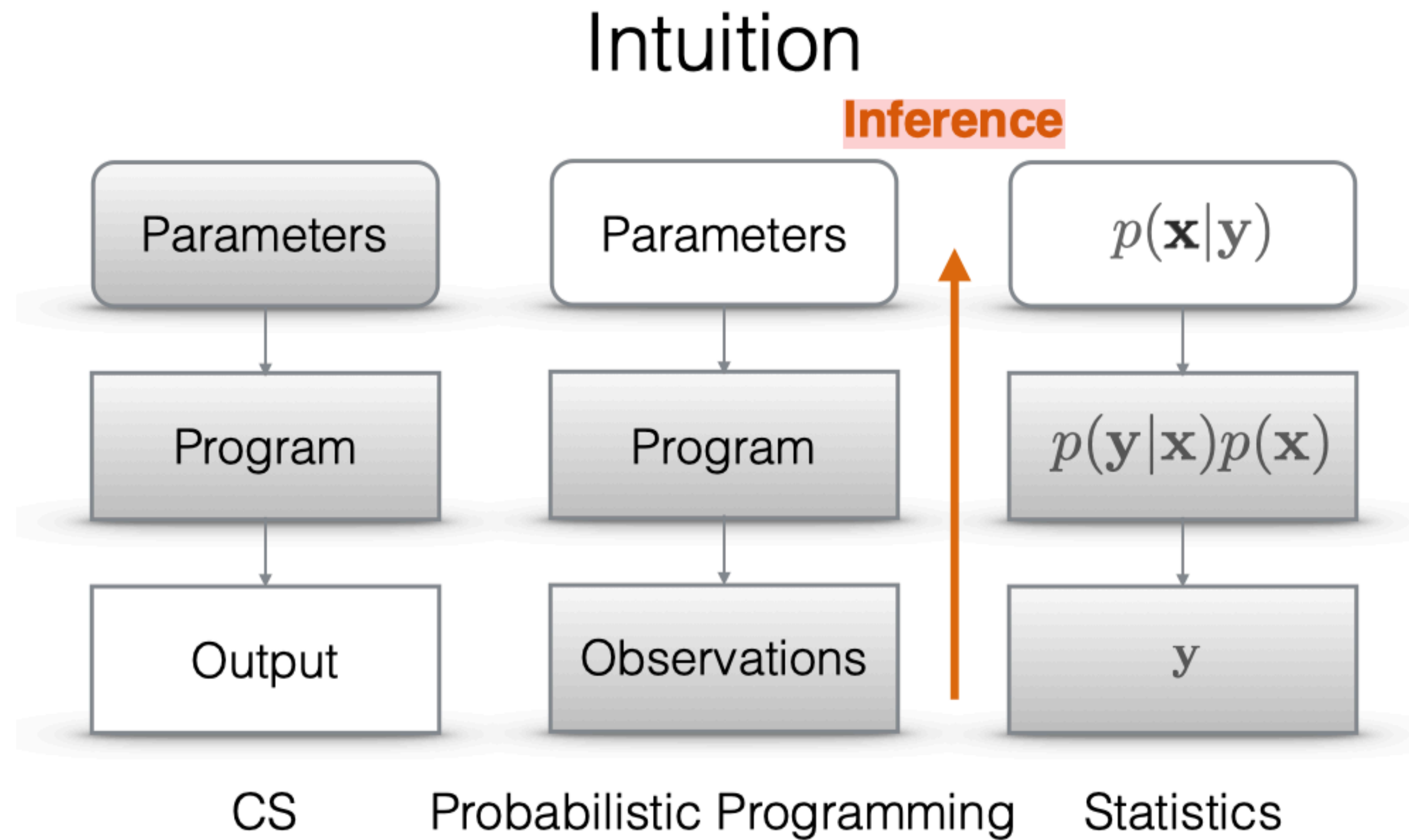
# Bayesian Optimization (BO)

## Surrogate Model

- Gaussian process

- Random Forrest

- t-Student processes

- Neural Networks

# Bayesian Optimization (BO)

## Acquisition Function

- expected improvement (EI)

- probability of improvement (PI)

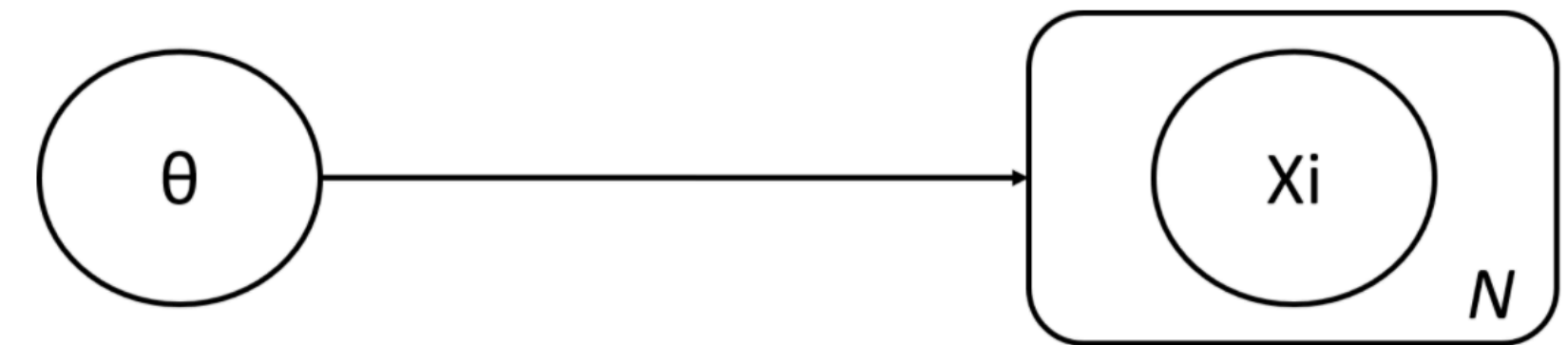- GP upper confidence bound (UCB)

- Thompson sampling (TS)

# Probabilistic Programming Languages (PPLs)

# Probabilistic Programming Languages (PPLs)

## Example: a biased coin toss

- Calculate the bias of a coin:

- Bernoulli distribution with latent variable $\theta$

- $P(x_i = 1 \mid \theta) = \theta$ and $P(x_i = 0 \mid \theta) = 1 - \theta$

- Infer $\theta$ based on previous results of coin toss - $P(\theta \mid x_1, x_2, \ldots, x_N)$



```
# Model
theta = Uniform(0.0, 1.0)
x = Bernoulli(probs=theta, sample_shape=10)
Data 5 data = np.array([0, 1, 0, 0, 0, 0, 0, 0, 0, 1])
Inference
qtheta = Empirical( 8 tf.Variable(tf.ones(1000) * 0.5))
inference = ed.HMC({theta: qtheta},
data={x: data})
inference.run()
Results 13 mean, stddev = ed.get_session().run( [qtheta.mean(),qtheta.stddev()])
print("Posterior mean:", mean)
print("Posterior stddev:", stddev)
```

# Probabilistic Programming Languages (PPLs)

**Probabilistic Programming (PP) in Julia: New Inference Algorithms**

*Day of a biologist who wants to use Gaussian Mixtures in his research*

# Probabilistic Programming Languages (PPLs)

**Recent popular PPL examples**

- Often built upon existing languages

- PyMC3/PyMC4 (Python)

- Stan (C++, Python, R)

- Turing.jl (Julia)

- WebPPL (JavaScript)

- Edward (Tensorflow)

- Pyro (PyTorch)

# ProBO

## a BO system for PPL models

- Computes and optimizes acquisition functions via operations that can be implemented in a broad variety of PPLs

- Goal: allow a custom model written in an arbitrary PPL to be plugged in and immediately used in BO.

# Related Work

- BOPP

  - describes a BO method for marginal maximum a posteriori (MMAP) estimates of latent variables

  - uses BO (with GP models) to help estimate latent variables in a given PPL

- BOAT

  - provides a custom PPL involving composed GP models with parametric mean functions for use in BO

  - uses exact inference & expected improvement

# ProBO

## Abstraction for Probabilistic Programs

- Three core PPL operations:

  - `inf(D)` - returns `post` (PPL dependent)

  - `post(s)` - returns a sample from the posterior distribution

  - `gen(x, z, s)` - returns sample from generative distribution

---

**Algorithm 1** ProBO($\mathcal{D}_0$, `inf`, `gen`)

---

1: **for** $n = 1, \ldots, N$ **do**
2: $\quad$ post $\leftarrow$ inf($\mathcal{D}_{n-1}$) $\qquad\qquad\qquad\qquad\quad$ ▷ Run inference algorithm to compute post
3: $\quad$ $x_n \leftarrow \text{argmin}_{x \in \mathcal{X}}\, a(x, \text{post}, \text{gen})$ $\qquad$ ▷ Optimize acquisition using post and gen
4: $\quad$ $y_n \sim s(x_n)$ $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ ▷ Observe system at $x_n$
5: $\quad$ $\mathcal{D}_n \leftarrow \mathcal{D}_{n-1} \cup (x_n, y_n)$ $\qquad\qquad\qquad$ ▷ Add new observations to dataset
6: Return $\mathcal{D}_N$.

---

# ProBO

## PPL Acquisition Functions

- Expected Improvement (EI)

- Probability of Improvement (PI)

- Upper Confident Bound (UCB)

- Thompson Sampling (TS)

---

**Algorithm 2** $a_{\mathrm{EI}}(x, \texttt{post}, \texttt{gen})$      ▷ EI

1: **for** $m = 1, \ldots, M$ **do**
2:     $z_m \leftarrow \texttt{post}(s_m)$
3:     $y_m \leftarrow \texttt{gen}(x, z_m, s_m)$
4: $f_{\min} \leftarrow \min_{y \in \mathcal{D}} f(y)$
5: Return $\sum_{m=1}^{M} \mathbb{1}\left[f(y_m) \leq f_{\min}\right] (f_{\min} - f(y_m))$

---

**Algorithm 3** $a_{\mathrm{PI}}(x, \texttt{post}, \texttt{gen})$      ▷ PI

1: **for** $m = 1, \ldots, M$ **do**
2:     $z_m \leftarrow \texttt{post}(s_m)$
3:     $y_m \leftarrow \texttt{gen}(x, z_m, s_m)$
4: $f_{\min} \leftarrow \min_{y \in \mathcal{D}} f(y)$
5: Return $\sum_{m=1}^{M} \mathbb{1}\left[f(y_m) \leq f_{\min}\right]$

---

**Algorithm 4** $a_{\mathrm{UCB}}(x, \texttt{post}, \texttt{gen})$      ▷ UCB

1: **for** $m = 1, \ldots, M$ **do**
2:     $z_m \leftarrow \texttt{post}(s_m)$
3:     $y_m \leftarrow \texttt{gen}(x, z_m, s_m)$
4: Return $\widehat{\mathrm{LCB}}\left(f(y_m)_{m=1}^{M}\right)$    ▷ See text for details

---

**Algorithm 5** $a_{\mathrm{TS}}(x, \texttt{post}, \texttt{gen})$      ▷ TS

1: $z \leftarrow \texttt{post}(s_1)$
2: **for** $m = 1, \ldots, M$ **do**
3:     $y_m \leftarrow \texttt{gen}(x, z, s_m)$
4: Return $\sum_{m=1}^{M} f(y_m)$

# ProBO
## Computational Considerations

- `inf()` cost dependent on PPLs inference algorithm

  - e.g. MCMC algorithms - $O(n)$ per iteration

- `inf()` only executed **once per query**

- Acquisition optimisation executed 100s times per query

  - `post()` & `gen()` cheaply implemented - $O(1)$

# ProBO

## Acquisition function optimisation

- `post()` & `gen()` not analytically differentiable

- Authors explored zeroth-order optimisation of $a_{MF}$

  - `post()` & `gen()` called $M_f$ times

  - Any zeroth-order optimisation algorithm can be used

---

**Algorithm 6** $a_{\mathrm{MF}}(x, \mathbf{post}, \mathbf{gen})$

1: $a_{\min} \leftarrow$ Min value of $a$ seen so far
2: $\ell = -\infty$, $f = 1$
3: **while** $\ell \leq a_{\min}$ **do**
4: $\quad \ell \leftarrow$ LCB-bootstrap $(\mathbf{post}, \mathbf{gen}, M_f)$
5: $\quad f \leftarrow f + 1$
6: Return $a(x, \mathbf{post}, \mathbf{gen})$ using $M = M_f$

---

**Algorithm 7** LCB-bootstrap$(\mathbf{post}, \mathbf{gen}, M_f)$

1: $y_{1:M_f} \leftarrow$ Call **post** and **gen** $M_f$ times
2: **for** $j = 1, \ldots, B$ **do**
3: $\quad \tilde{y}_{1:M_f} \leftarrow$ Resample$(y_{1:M_f})$
4: $\quad a_j \leftarrow \lambda(\tilde{y}_{1:M_f})$ $\quad \triangleright$ See text for details
5: Return LCB $(a_{1:B})$

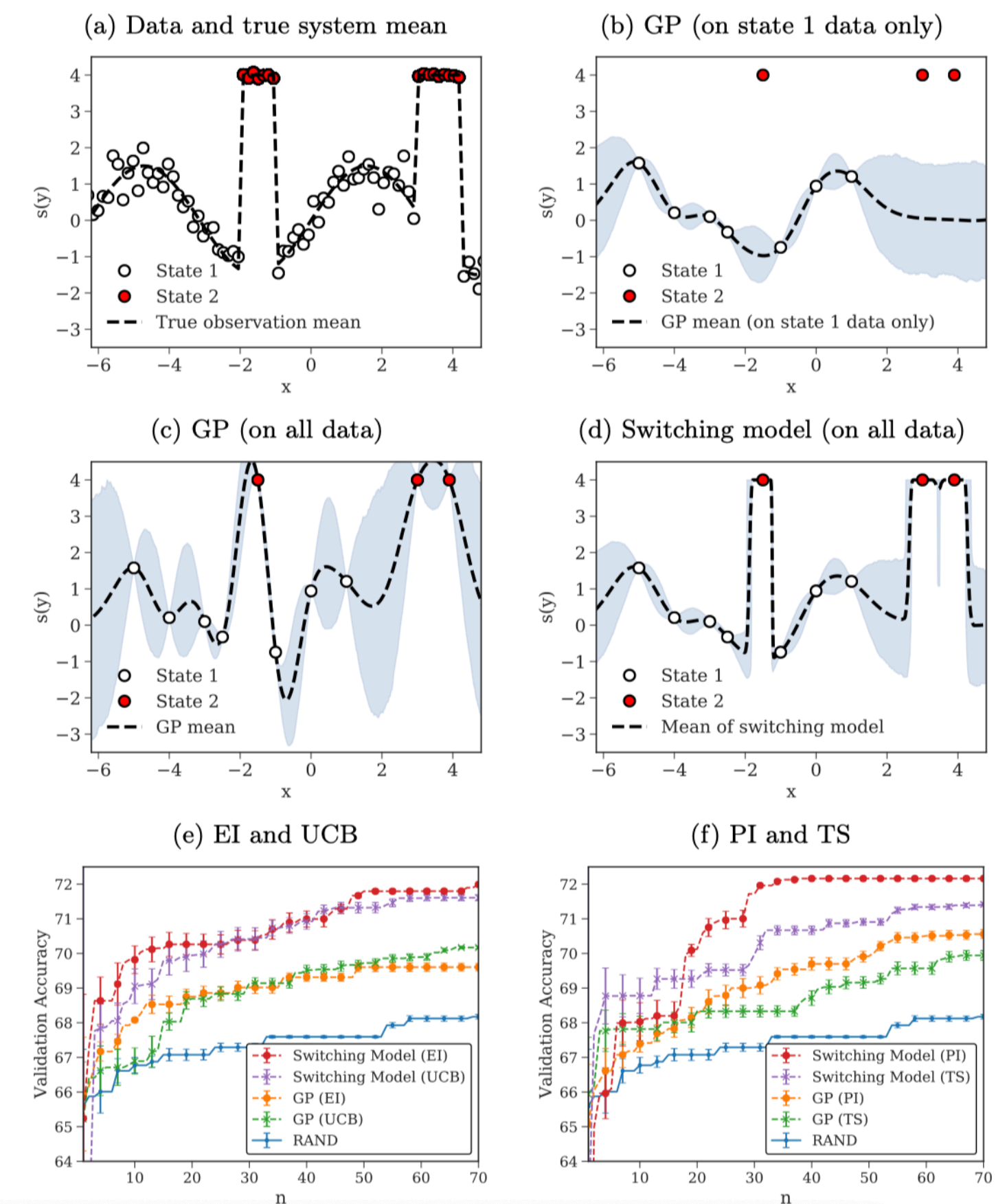# Examples and Experiments

**Experiment Setting**

- PPL Implementations:

  - Stan + No U-Turn Sampler (a form of Hamiltonian Monte Carlo)

  - Edward + black box variational inference

- GP comparisons:

  - George

  - GPy

# Examples and Experiments

## BO with State Observations

- Switching Model: ProBO using a dynamic value of $M_f$

- Task: neural network architecture and hyperparameter search

  - multi-layer perceptron (MLP) neural networks

  - x = (number of layers, layer width, learning rate,  batch size)

  - Pima Indians Diabetes Dataset

- ProBO+switching model **vs** BO+GP



*BO with State Observations*

(a) Data and true system mean
(b) GP (on state 1 data only)
(c) GP (on all data)
(d) Switching model (on all data)
(e) EI and UCB
(f) PI and TS

# Examples and Experiments
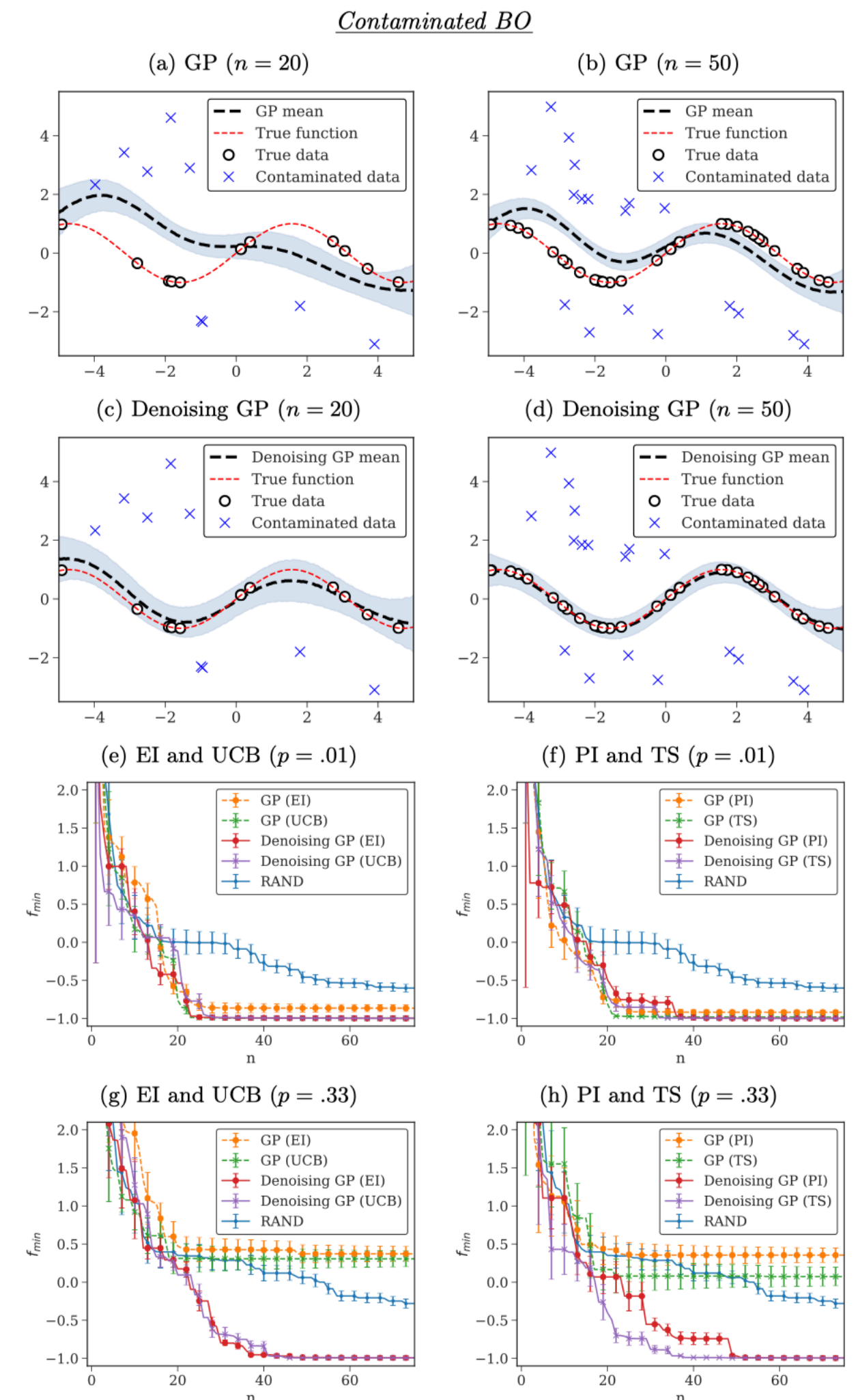## Robust Models for Contaminated BO

- not have access to state observations

- Task: synthetic optimization task

  - system model $M_s$

  - contamination model $M_c$

  - denoising model $y \sim w_s M_s(\cdot \mid z_s; x) + w_c M_c(\cdot \mid z_c; x)$

- ProBO+denoising model **vs** BO+GP



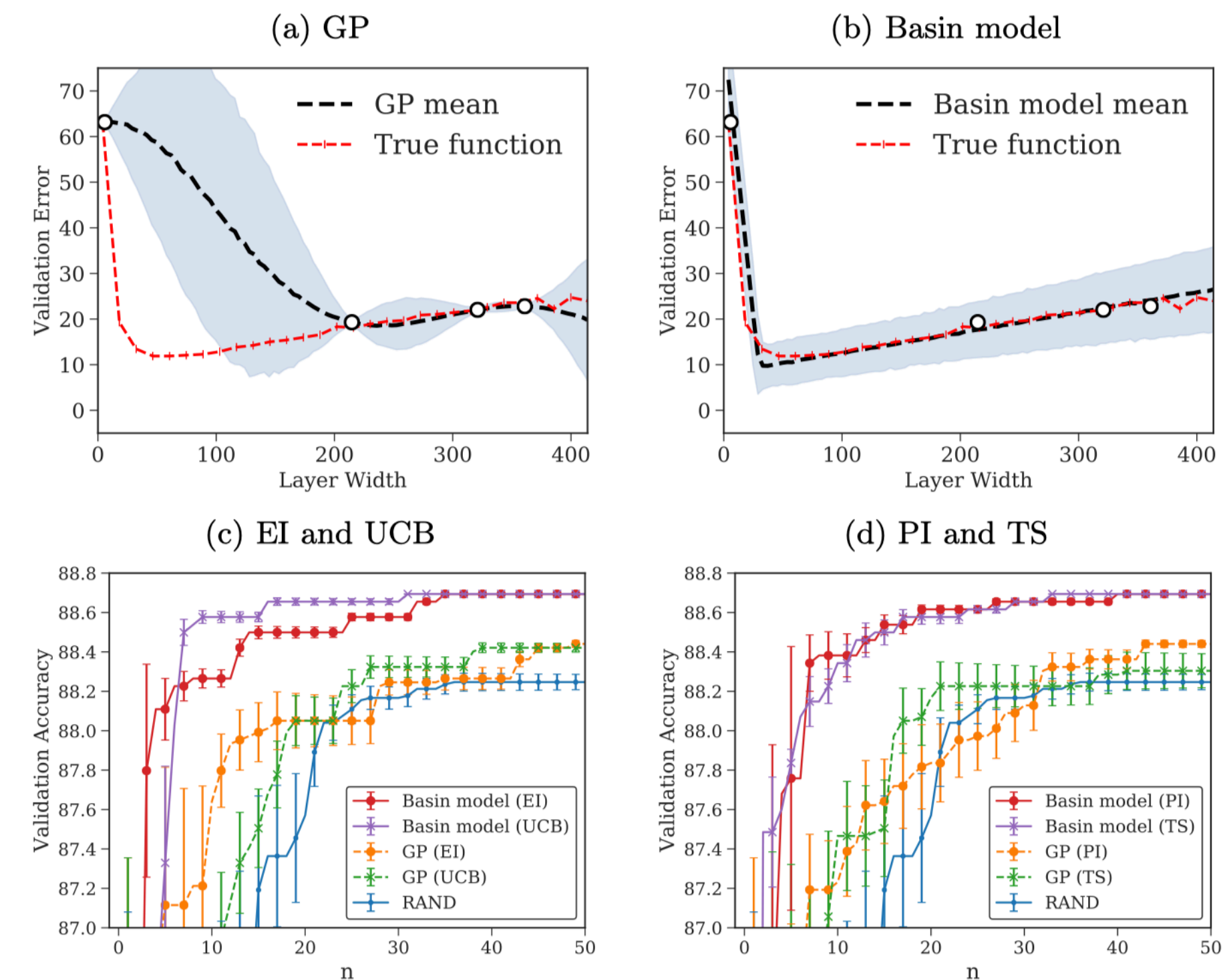*Contaminated BO*

# Examples and Experiments
## BO with Prior Structure on the Objective Function

- basin model: have prior knowledge about properties of the objective function

- Task: tuning model complexity

  - number of units of hidden layers in 4-layer MLP

  - Wisconsin Breast Cancer Diagnosis dataset
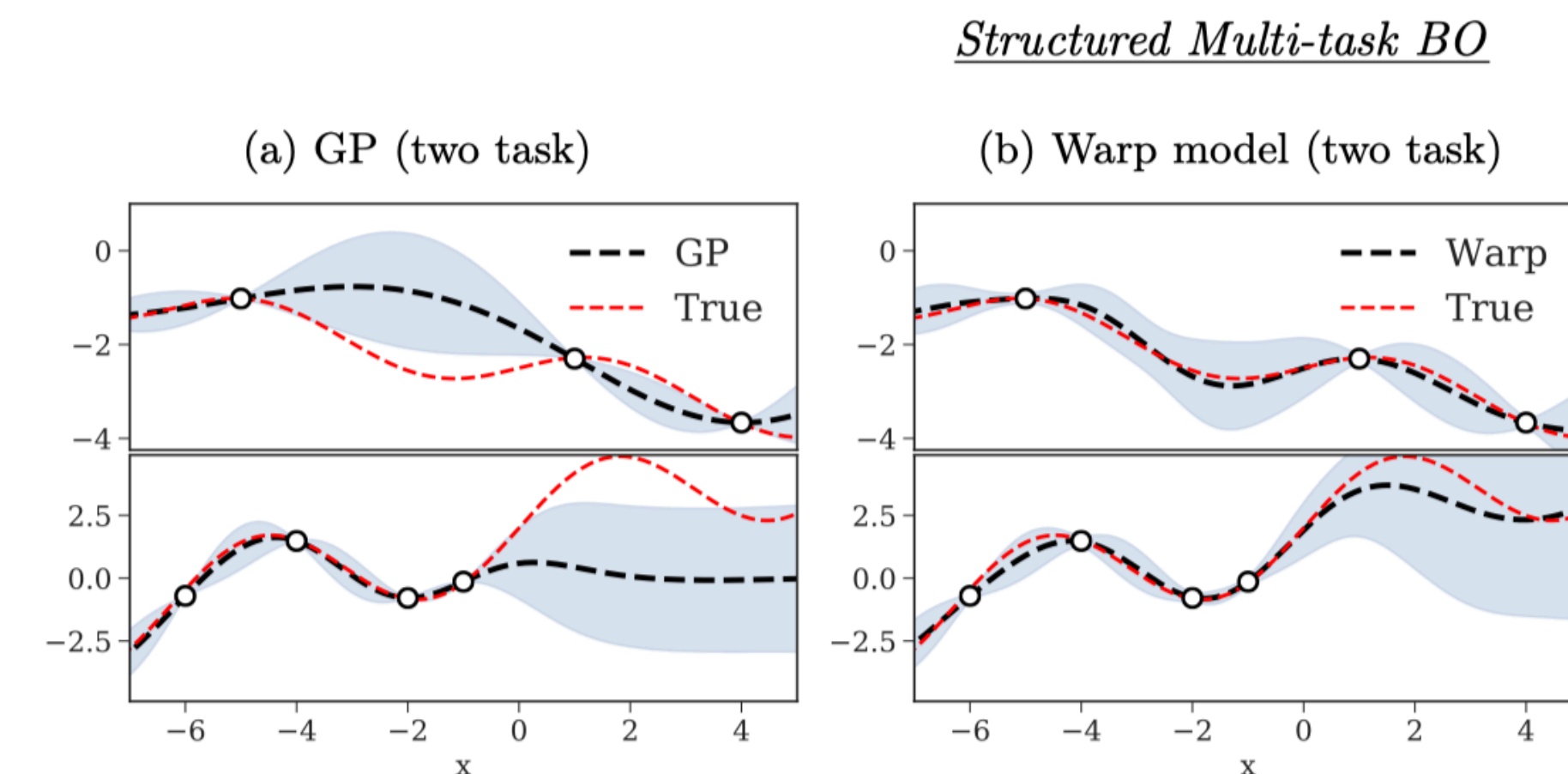
- ProBO+basin model **vs** BO+GP



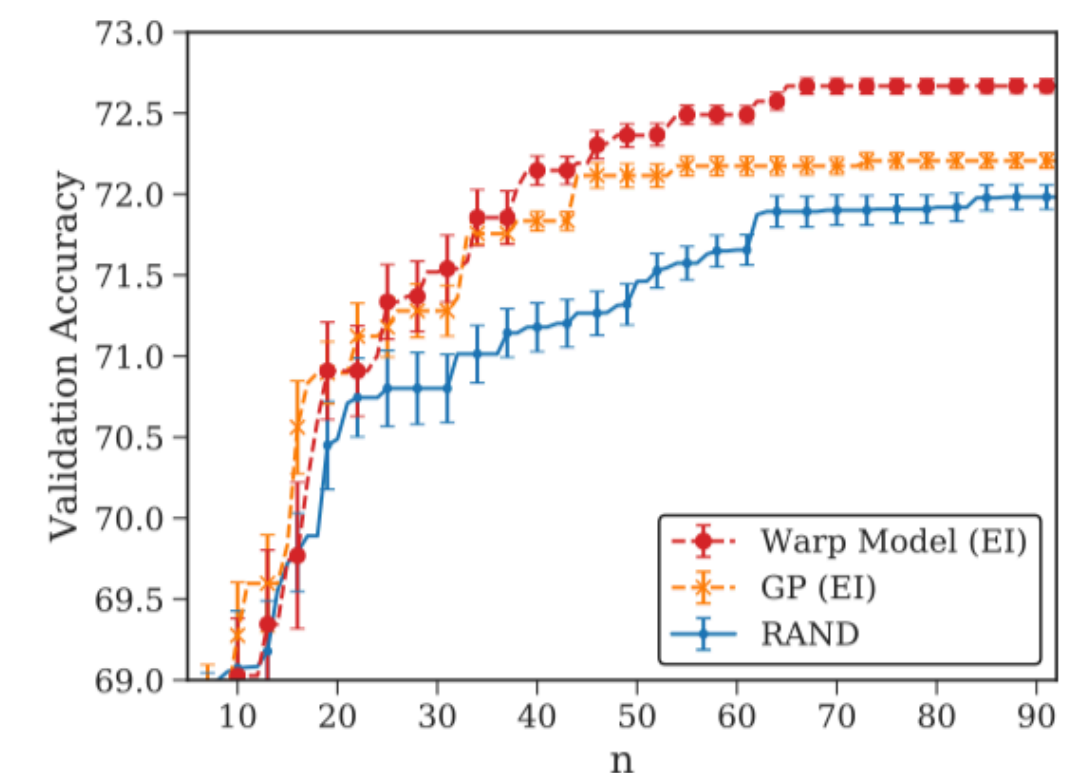*BO with Prior Structure on the Objective Function*

# Examples and Experiments
## Structured Models for Multi-task and Contextual BO, and Model Ensembles

- optimize multiple systems jointly, where there is some known relation between the systems

- a finite set of systems (multi-task BO)

- systems are each indexed by a context vector c ∈ R d (contextual BO)

- warp model: incorporate prior structure about the relationship among these systems, warps a latent model based on context/task-specific parameters

- parametric model: model with a specific trend, shape, or specialty for a subset of the data

- posterior predictive densities of multiple PPL models, using only our three PPL operations

*Structured Multi-task BO*

(a) GP (two task)

(b) Warp model (two task)

- - - GP
- - - True

- - - Warp
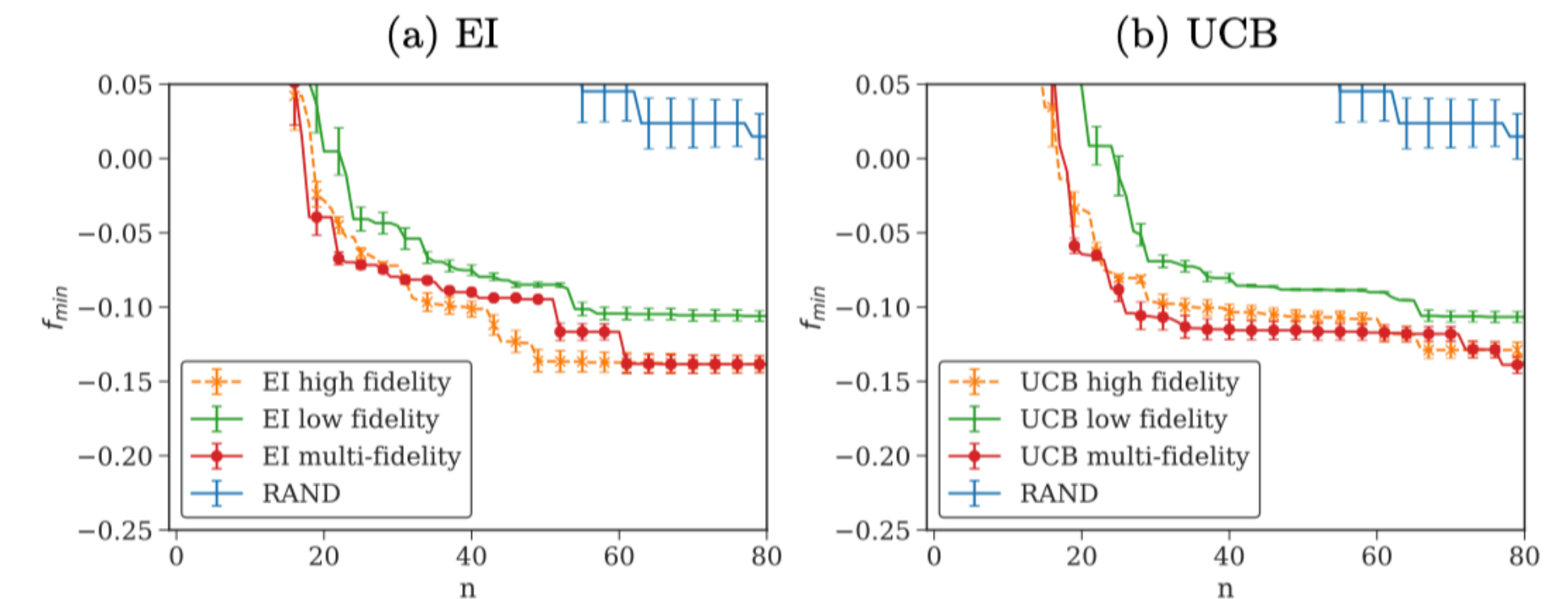- - - True

(c) EI

Warp Model (EI)
GP (EI)
RAND

# Examples and Experiments

## Multi-fidelity Acquisition Optimization

- two-fidelity setting:

  - high-fidelity a (M = 1000)

  - low-fidelity a (M = 10)

  - multi-fidelity $a_{\mathrm{MF}}$

- 3x better performance than high-fidelity in terms of calls to gen()

*Multi-fidelity Acquisition Functions*



(a) EI     (b) UCB

(c) Calls to **gen**

| PPL acquisition method $a(x)$ | Avg. number $\mathrm{gen}/a(x)$ |
|---|---|
| EI high-fidelity | 1000 |
| EI multi-fidelity | 347.89 |
| EI low-fidelity | 10 |
| UCB high-fidelity | 1000 |
| UCB multi-fidelity | 324.65 |
| UCB low-fidelity | 10 |

# Opinion of the paper

## Key Takeaway

- Present ProBO, a system for versatile Bayesian optimization using models from any PPL

- Use PPLs to implement new models for optimization settings that are difficult for standard BO methods and models

# Following Work

- BoTorch: A Framework for Efficient Monte-Carlo Bayesian Optimization. NIPS 2020.

  - a library for Bayesian Optimization built on PyTorch.

  - it benefit from gradient-based optimization provided by differentiable programming, as well as algebraic methods designed to exploit GPU acceleration.

- BANANAS: Bayesian Optimization with Neural Architectures for Neural Architecture Search

  - BO + neural predictor framework as a high-performance framework for NAS.

  - use the ProBO implementation.

# Opinion of the paper
## Criticism

- The writing

  - the structure is clear and easy to follow

  - related work is not sufficient enough

- The release of ProBO is not open-sourced at all

  - Make it difficult to understand the implementation details

  - Limit the spread of the author's idea

  - No further maintenance or extension

# References

- A Gentle Introduction to Probabilistic Programming Languages

- Bayesian Methods for Hackers: Probabilistic Programming and Bayesian Inference

- An Introduction to Probabilistic Programming

# Thanks for listening!

# Q&A