

Predictive Entropy Search for Efficient Global Optimization of Black-box Functions

José Miguel Hernández-Lobato, Matthew W. Hoffman, Zoubin Ghahramani

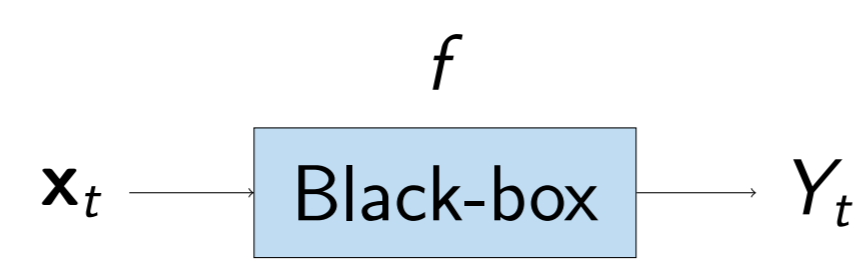
University of Cambridge



UNIVERSITY OF CAMBRIDGE

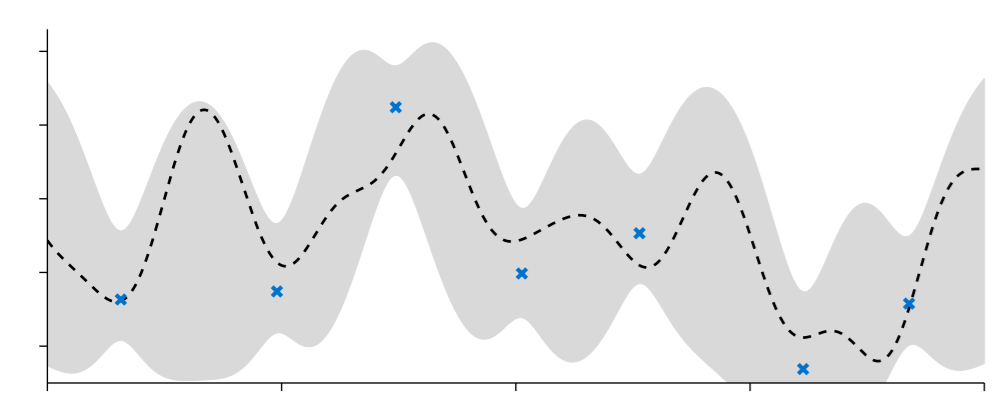
Bayesian black-box optimization

We are interested in solving black-box optimization problems of the form $\mathbf{x}^* = \arg \max_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x})$ by sequentially **querying** points \mathbf{x}_t and observing Y_t .



Here **black-box** means:

- we may only be able to observe the function value, i.e. no gradients
- our observations may be corrupted by noise, i.e. we observe Y_t , but $f(\mathbf{x}_t) = \mathbb{E}[Y_t | \mathbf{x}_t]$



Given some function to optimize (the dashed line) we can make observations Y_t and construct a **Bayesian posterior** over the (unknown!) function.

This posterior can then be used to guide our search, i.e. the selection of inputs \mathbf{x}_t .

A framework for Bayesian optimization

The following pseudocode outlines an algorithmic framework for Bayesian optimization:

- | | | |
|--|-----------------------|---|
| 1: for $t = 1, \dots, T$ do | | |
| 2: select point $\mathbf{x}_t = \arg \max_{\mathbf{x}} \alpha_{t-1}(\mathbf{x})$ | <i>Exploration</i> | select the next point to query by maximizing some acquisition function |
| 3: observe $y_t \sim p(\cdot \mathbf{x}_t)$ | | |
| 4: update posterior $p(f \mathcal{D}_t)$ | <i>Prediction</i> | posterior model for prediction; generally use a GP |
| 5: end for | | |
| 6: return $\tilde{\mathbf{x}}_T$ | <i>Recommendation</i> | return final recommendation; can maximize the posterior mean |

The **acquisition function** used to explore the function should in some sense try to gain as much information about the optimizer location as possible.

Predictive Entropy Search

A common active learning approach is to select points which maximize the expected reduction in posterior entropy about some predictor. Applying this to optimization as in [2, 3], let \mathbf{x}_* be the unknown optimizer and write its information gain:

$$\alpha_t(\mathbf{x}) = H[\mathbf{x}_* | \mathcal{D}_t] - \mathbb{E}_{P(y | \mathcal{D}_t, \mathbf{x})} [H[\mathbf{x}_* | \mathcal{D}_t \cup \{(\mathbf{x}, y)\}]] \quad (\text{ES})$$

- But:
- the distribution $P(\mathbf{x}_* | \dots)$ has no closed form expression; and
 - this expensive (and rough) approximation must be done for every (\mathbf{x}, y)

Note, however, that the information gain is symmetric. Changing the order of these arguments we can write the acquisition as

$$\alpha_t(\mathbf{x}) = H[y | \mathcal{D}_t, \mathbf{x}] - \mathbb{E}_{P(\mathbf{x}_* | \mathcal{D}_t)} [H[y | \mathcal{D}_t, \mathbf{x}, \mathbf{x}_*]] \quad (\text{PES})$$

which we call **Predictive Entropy Search**.

- This requires:
- sampling from the distribution over maximizers \mathbf{x}_* , and
 - computing the predictive entropy conditioned on this maximizer.

Sampling optima

To sample an optimum we need only sample $f \sim p(\cdot | \mathcal{D}_t)$ and return $\arg \max_{\mathbf{x}} f(\mathbf{x})$.

- if our set of possibly query points is discrete this is **Thompson sampling**;
- however, f is an infinite dimensional object!

Instead we will approximately sample a $f(\cdot) = \phi(\cdot)^T \theta$ where ϕ consist of **random features** and θ are sampled from the resulting approximate posterior.

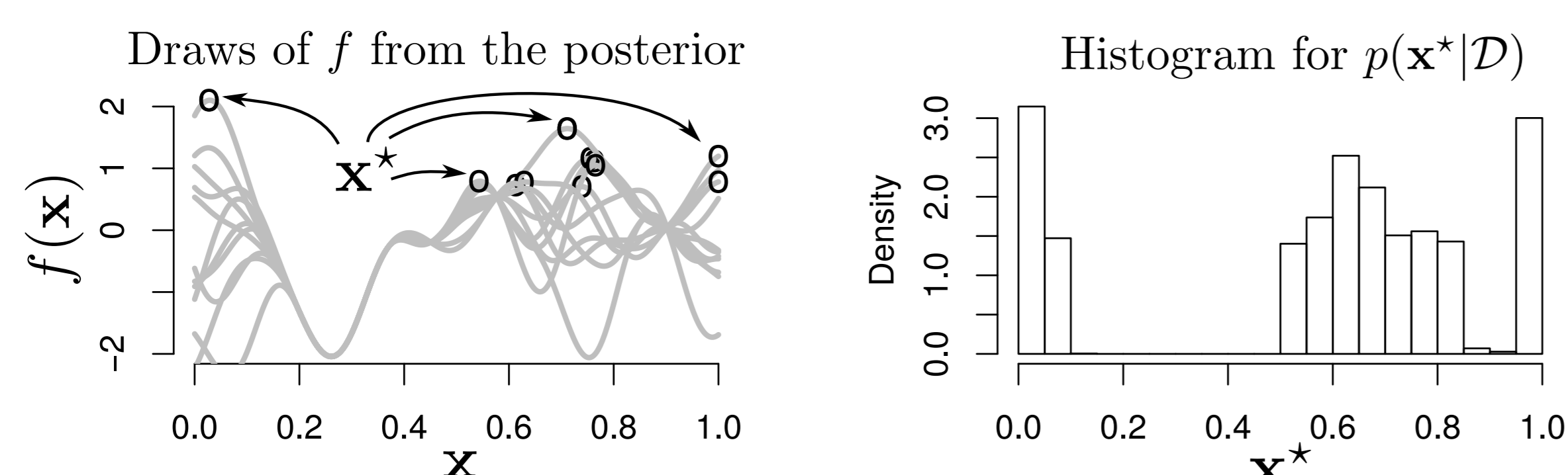
Via Bochner's theorem [1] a shift-invariant kernel can be written as the Fourier transform of its spectral density. Treating this as a probability distribution $p(\mathbf{w})$ we can write:

$$k(\mathbf{x}, \mathbf{x}') = \alpha \mathbb{E}_{p(\mathbf{w})} [e^{-i\mathbf{w}^T(\mathbf{x}-\mathbf{x}')}] = 2\alpha \mathbb{E}_{p(\mathbf{w}, b)} [\cos(\mathbf{w}^T \mathbf{x} + b) \cos(\mathbf{w}^T \mathbf{x}' + b)] \quad (1)$$

where b is uniformly distributed in $[0, 2\pi]$. Letting

$$\phi(\mathbf{x}) = \sqrt{2\alpha/m} \cos(\mathbf{W}\mathbf{x} + \mathbf{b})$$

be the m -dimensional feature map we can approximate the GP posterior with a simple linear-Gaussian model; can sample θ directly.



Approximating the entropy

To construct $\alpha_t(\mathbf{x})$ we can approximate the fact that \mathbf{x}_* is a maximum with the following constraints:

- $\nabla f(\mathbf{x}_*) = 0$;
- $\text{diag}[\nabla^2 f(\mathbf{x}_*)] < 0$;
- $f(\mathbf{x}_*) > f(\mathbf{x})$.
- $\text{upper}[\nabla^2 f(\mathbf{x}_*)] = 0$;
- $f(\mathbf{x}_*) > \max_t f(\mathbf{x}_t)$.

The first two constraints can be incorporated exactly, producing $\mathcal{N}(\mathbf{z} | \mathbf{m}_0, \mathbf{V}_0)$ where \mathbf{z} are the latent values for the maximizer and Hessian. The second set of constraints can be approximated as

$$\mathcal{N}(\mathbf{z} | \mathbf{m}_0, \mathbf{V}_0) \times \underbrace{\Phi_{\sigma^2}(f(\mathbf{x}_*) - y_{\max})}_{\text{maximizer constraint}} \times \underbrace{\prod_{i=1}^d \mathbb{I}([\nabla^2 f(\mathbf{x}_*)]_{ii} \leq 0)}_{\text{Hessian constraints}}$$

Using **Expectation Propagation** we can approximate this density with a single multivariate Gaussian, leading to a Gaussian distribution over the latent $f(\mathbf{x}_*)$. This need only be computed once per iteration.

Finally, given any \mathbf{x} we can compute the joint $p(f(\mathbf{x}), f(\mathbf{x}_*))$. An additional factor can be incorporated requiring $f(\mathbf{x}) < f(\mathbf{x}_*)$ and approximated using a single step of EP.

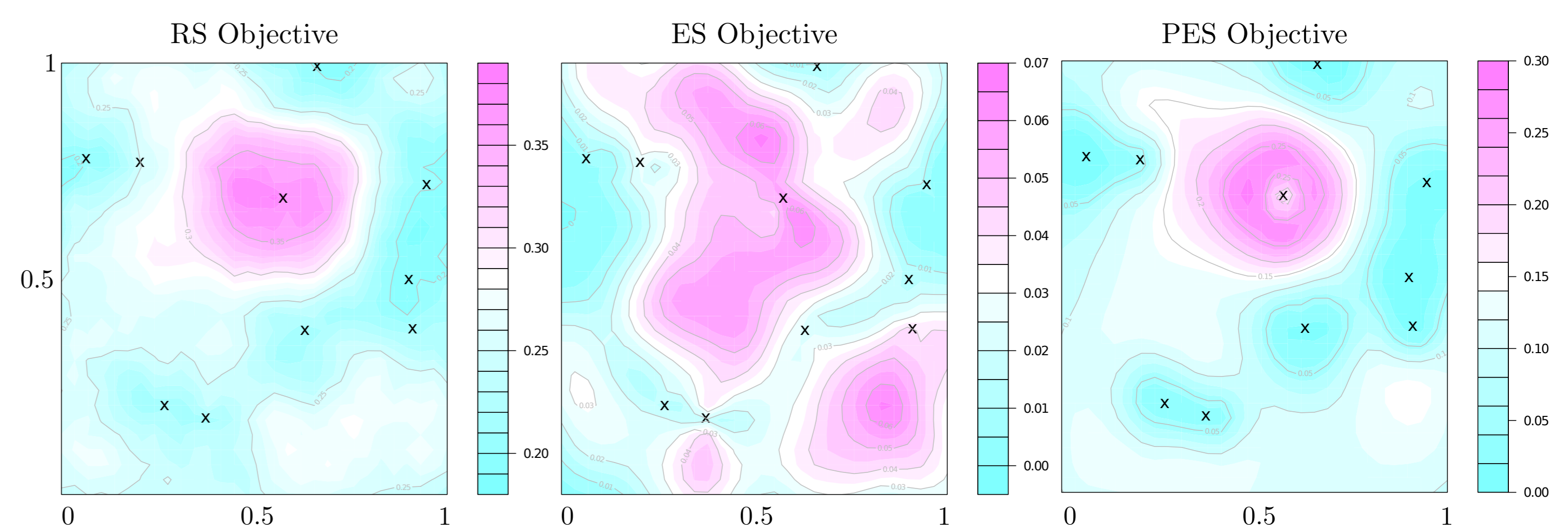
Taking entropies of these Gaussian distributions the acquisition function is:

$$\alpha_n(\mathbf{x}) = 0.5 \log[v_n(\mathbf{x}) + \sigma^2] - 0.5 \log[v_n(\mathbf{x} | \mathbf{x}_*) + \sigma^2]$$

where $v_n(\mathbf{x})$ and $v_n(\mathbf{x} | \mathbf{x}_*)$ are the unconditioned and conditioned variances.

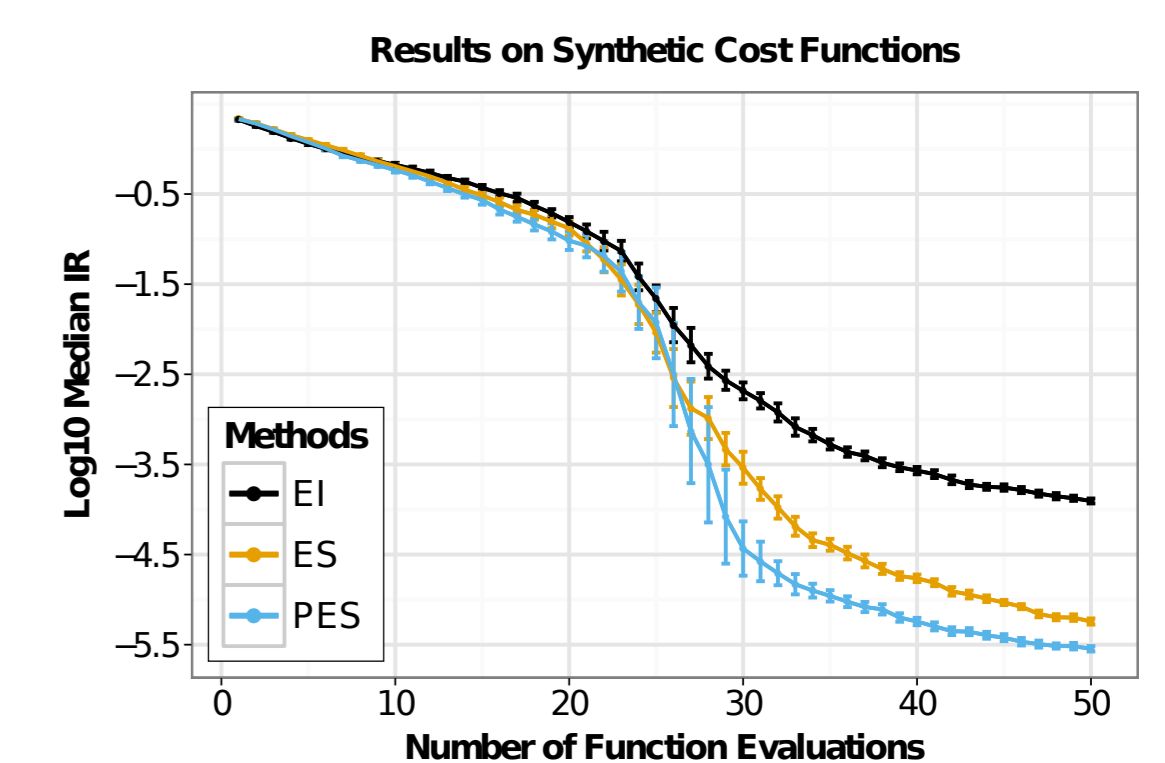
Accuracy of the PES acquisition

The following compares a fine-grained random sampling (RS) scheme to compute the ground truth objective with ES and PES. We see PES provides a much better approximation.

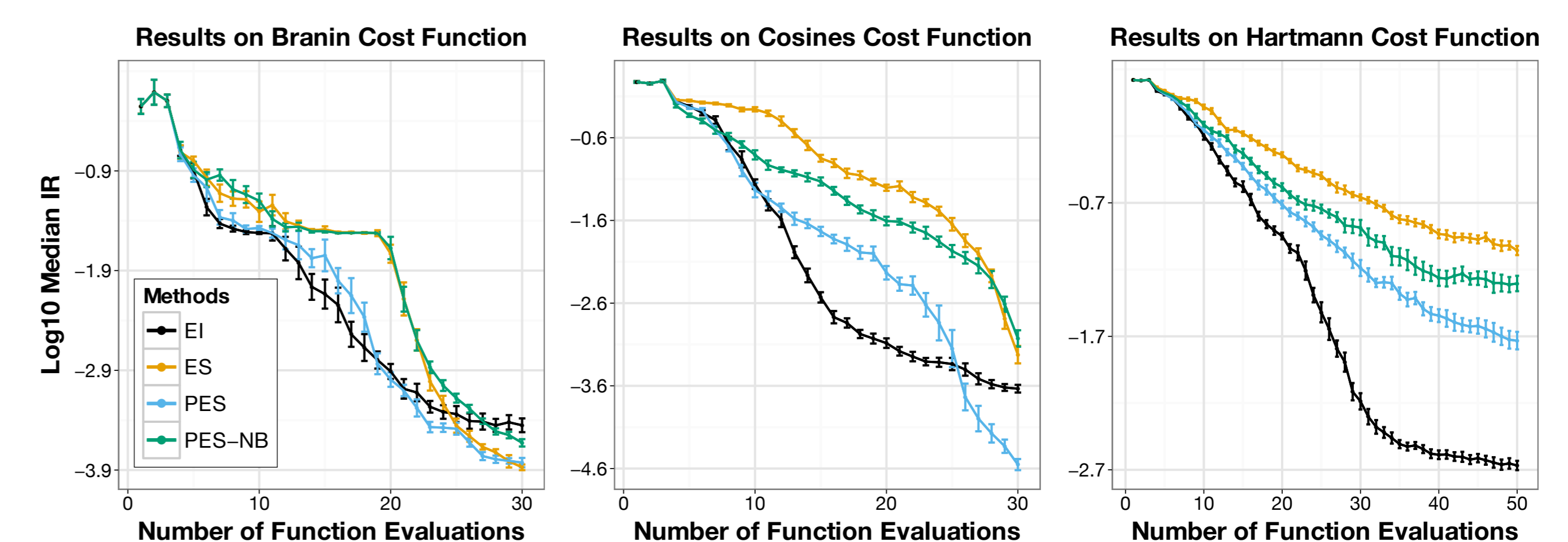


Performance of PES

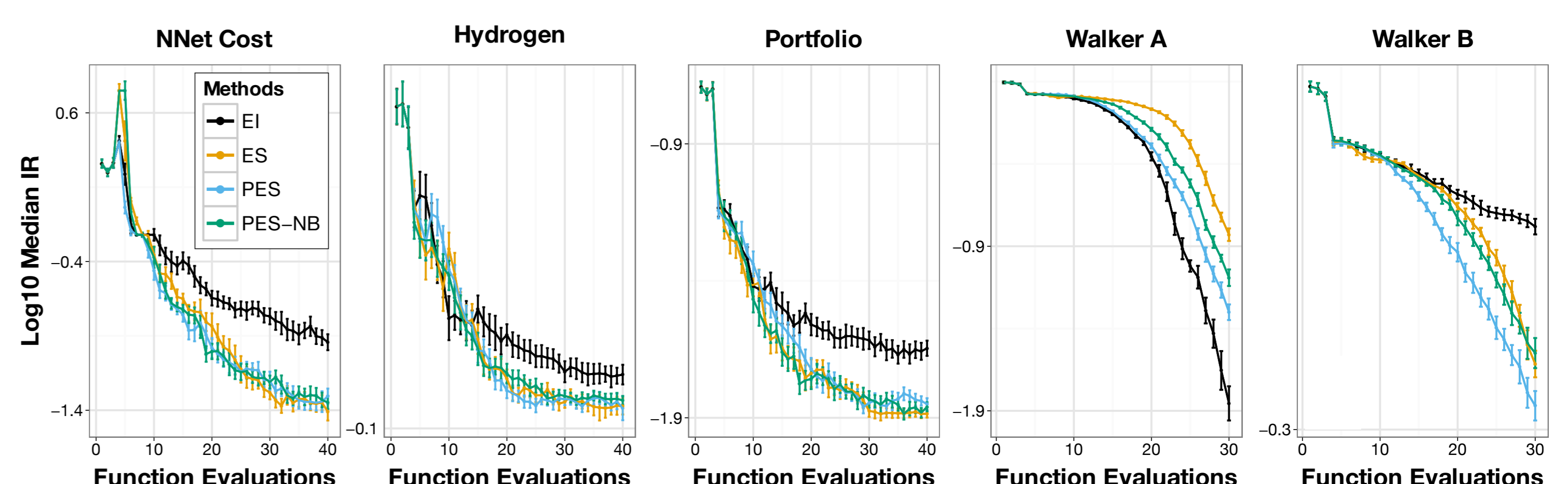
The *Right* plot shows the performance on a function randomly sampled from the given GP prior.



The *Bottom* plots compare the performance on various global optimization test problems.



We further consider several real-world experiments. (NNet) the hyperparameters of a neural network; (Hydrogen) the hydrogen production of a particular bacteria; (Portfolio) the Sharpe ratio of 1-year simulated returns; and the speed of a bipedal robot under (Walker A) noiseless and (Walker B) noisy observations.



Note also: the Walker-A experiment had less noise, and hence here EI showed better performance, due to its more exploitative behavior.

[1] S. Bochner. Lectures on Fourier integrals. Princeton University Press, 1959.

[2] P. Hennig and C. J. Schuler. Entropy search for information-efficient global optimization. *Journal of Machine Learning Research*, 13, 2012.

[3] J. Villemonteix, E. Vazquez, and E. Walter. An informational approach to the global optimization of expensive-to-evaluate functions. *Journal of Global Optimization*, 44(4):509-534, 2009.