

Bayesian Optimization

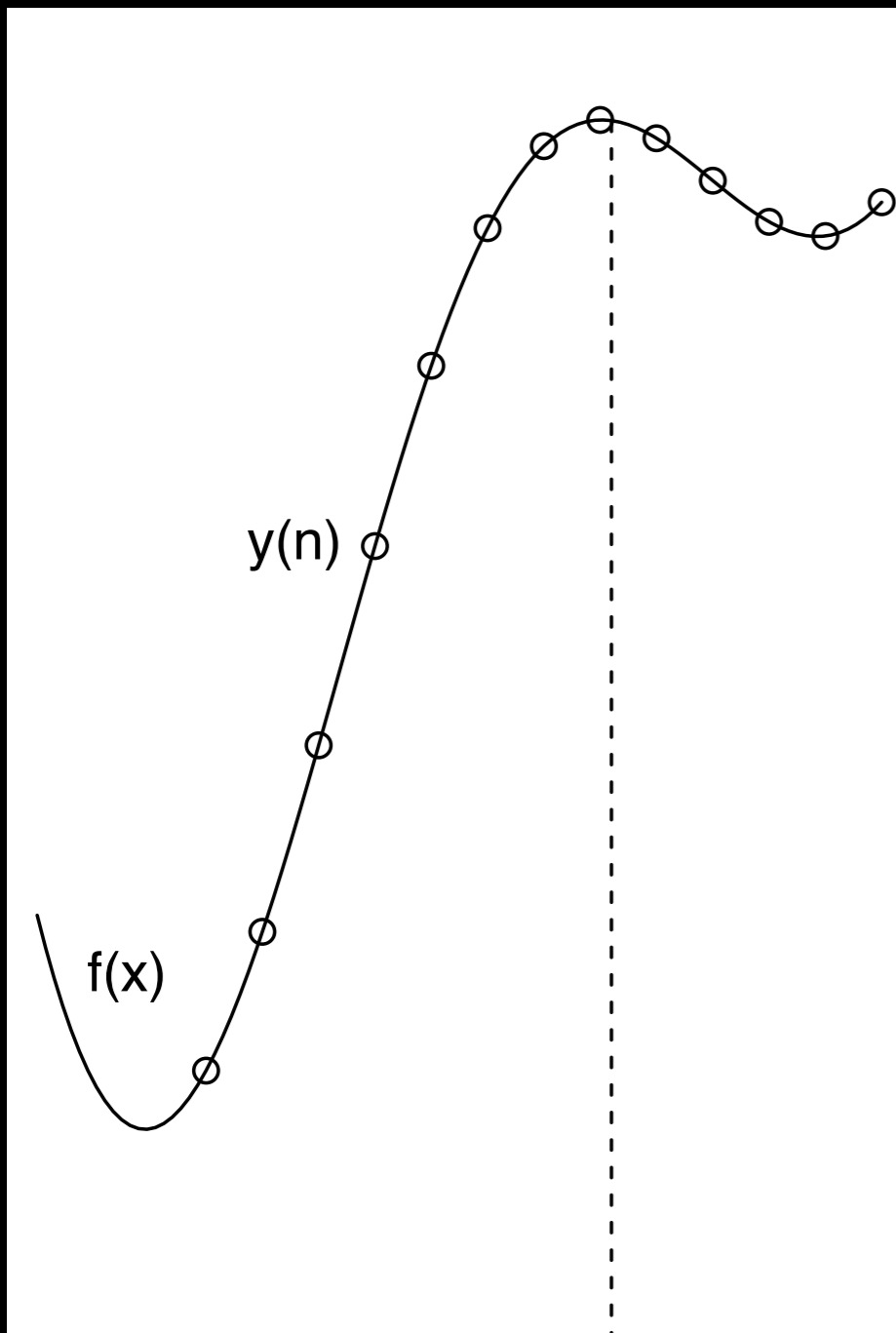
Peter Frazier

Associate Professor, Operations Research & Information Engineering, Cornell
Staff Data Scientist, Uber

For details & links to code, see

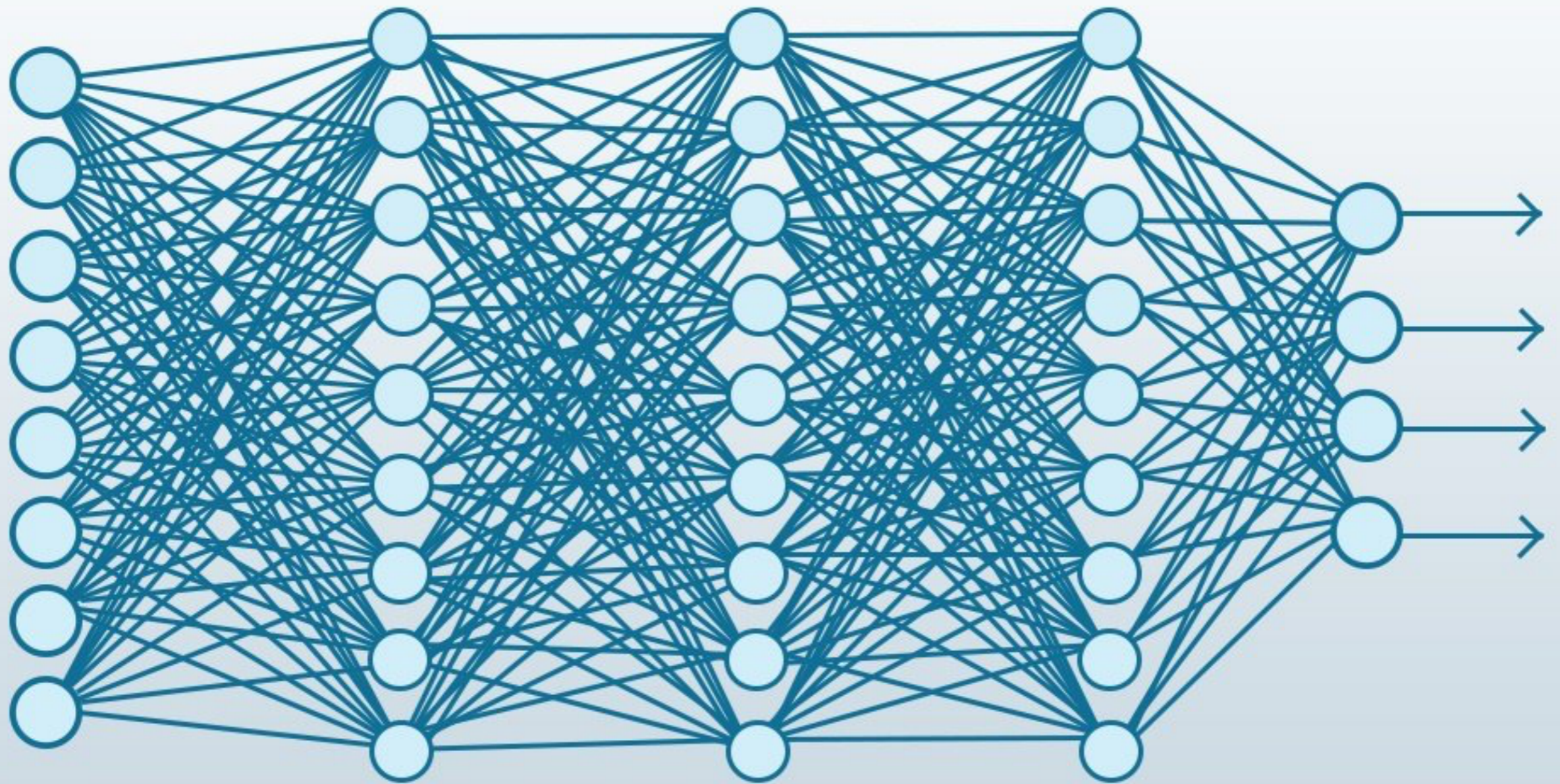
P. Frazier, "A Tutorial on Bayesian Optimization" <https://arxiv.org/abs/1807.02811>

This is the standard problem in Bayesian Optimization



- We'd like to optimize $F : \mathbb{R}^d \rightarrow \mathbb{R}$, where $d < 20$.
- F 's feasible set A is simple, e.g., box constraints.
- F is continuous but lacks special structure, e.g., concavity, that would make it easy to optimize.
- F is derivative-free: evaluations do not give gradient information.
- F is expensive to evaluate: the # of times we can evaluate it is severely limited.
- F may be noisy. If noise is present, we'll assume it is independent and normally distributed, with common but unknown variance.

Optimization of expensive functions arises when fitting machine learning models



Optimization of expensive functions arises when tuning algorithms via backtesting



Optimization of expensive functions arises when tuning websites with A/B testing

yelp Find coffee Near San Francisco, CA Sign Up Log In

Home About Me Write a Review Find Friends Messages Talk Events

coffee San Francisco, CA Showing 1-10 of 6567

Browse Category: Coffee & Tea Show Filters

1. Blue Bottle Coffee Hayes Valley
★★★★★ 1558 reviews
\$\$ · Coffee & Tea
315 Linden St
San Francisco, CA 94102
(510) 653-3394

This Blue Bottle location is so cute and tiny. Way tinier than their other locations--it almost looks like a little pop up shop in a garage. Good thing it still brews their super yummy **coffee**...

2. Philz Coffee 748 Van Ness Ave
★★★★★ 1216 reviews
\$\$ · Coffee & Tea
San Francisco, CA 94102
(415) 292-7660

The hype about Philz is real. Gingersnap iced **coffee**, where have you been my whole life? Not too **coffee**-y and not too gingersnap-y. And the girl working was the one who suggested it when we...

3. Blue Bottle Coffee Co SoMa
★★★★★ 1524 reviews
\$\$ · Coffee & Tea
66 Mint St
San Francisco, CA 94103
(510) 653-3394

Excellent iced **coffee**- location is tucked away but is the best **coffee** we've found this close to Mosso.

Mo' Map Redo search when map moved

Google Maps Map Data Terms of Use Report a map error

Optimization of expensive functions arises when tuning transportation markets

U B E R

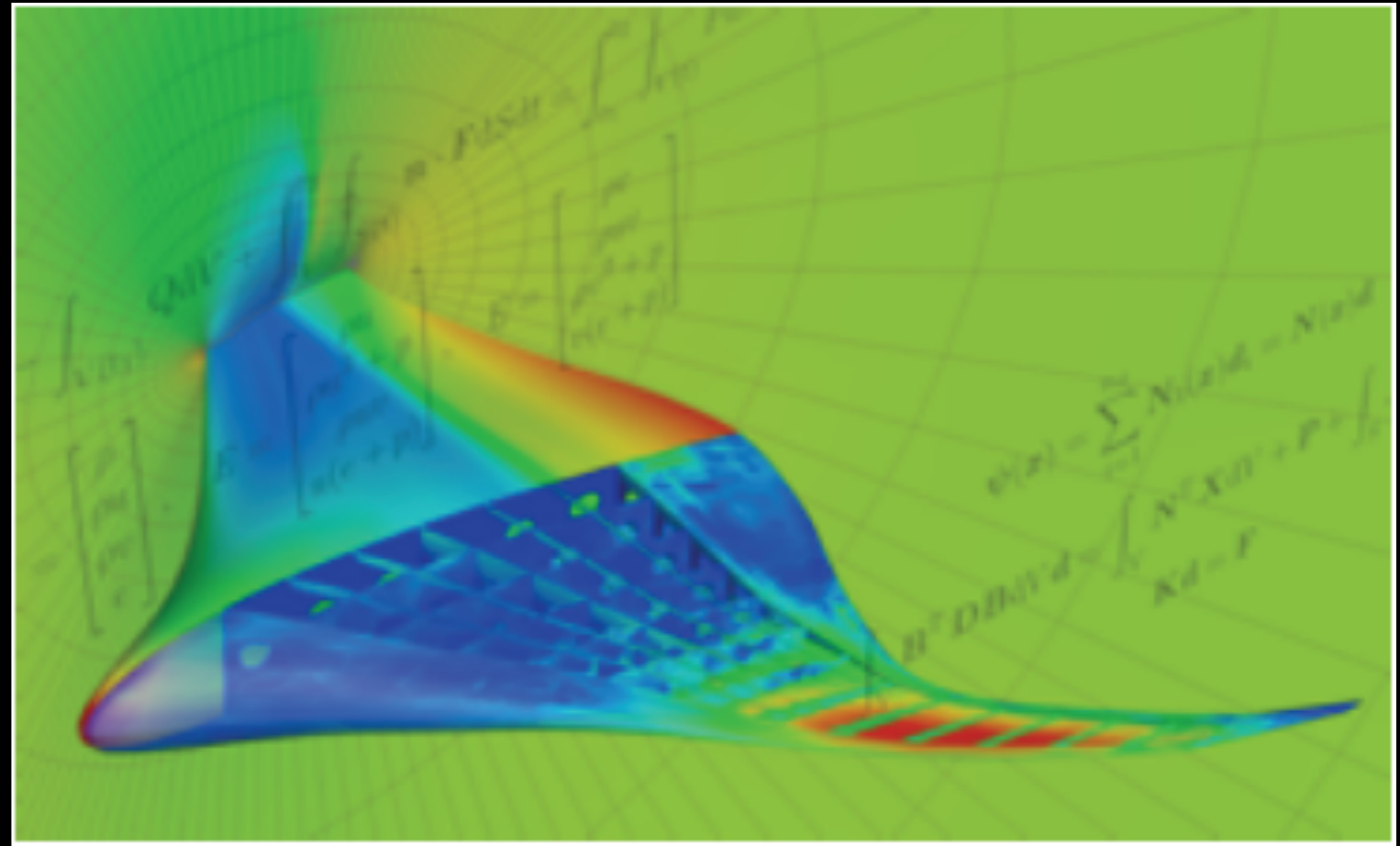
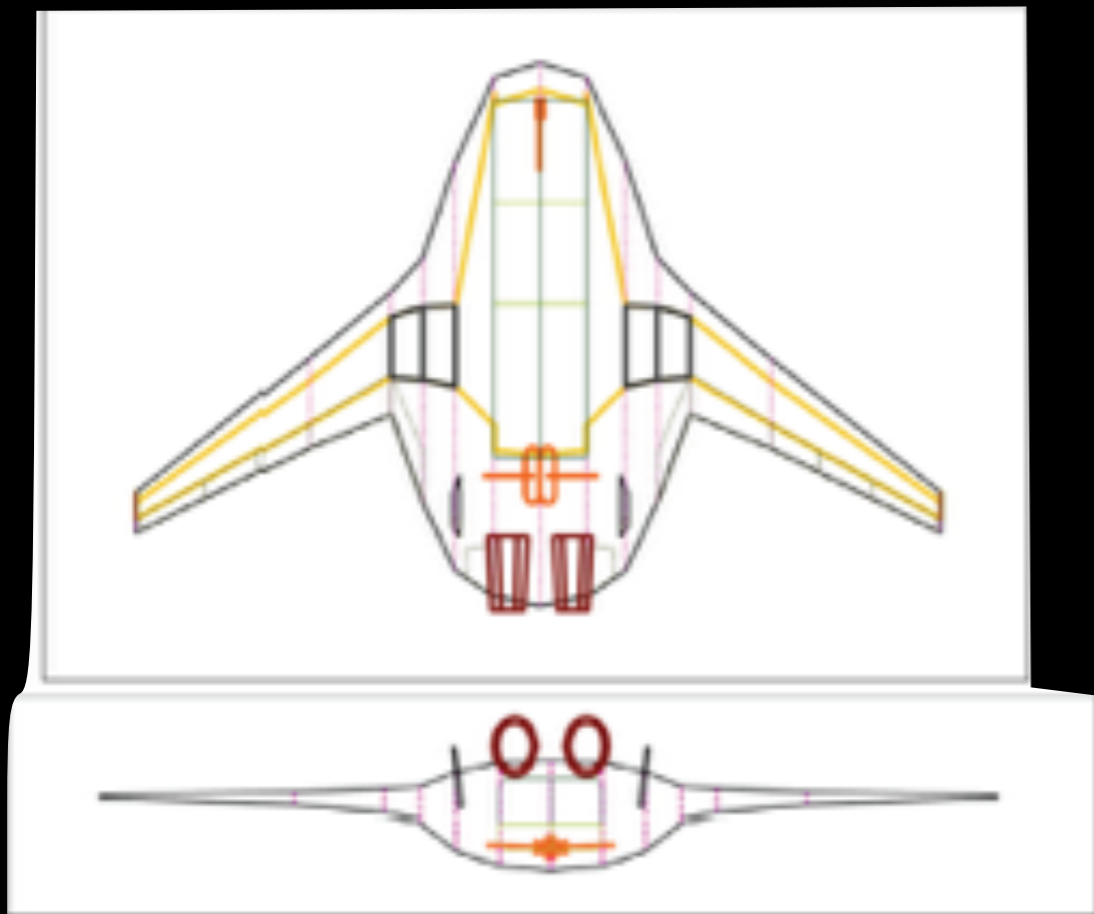
 uberPOOL

SHARE YOUR RIDE, SPLIT THE COST

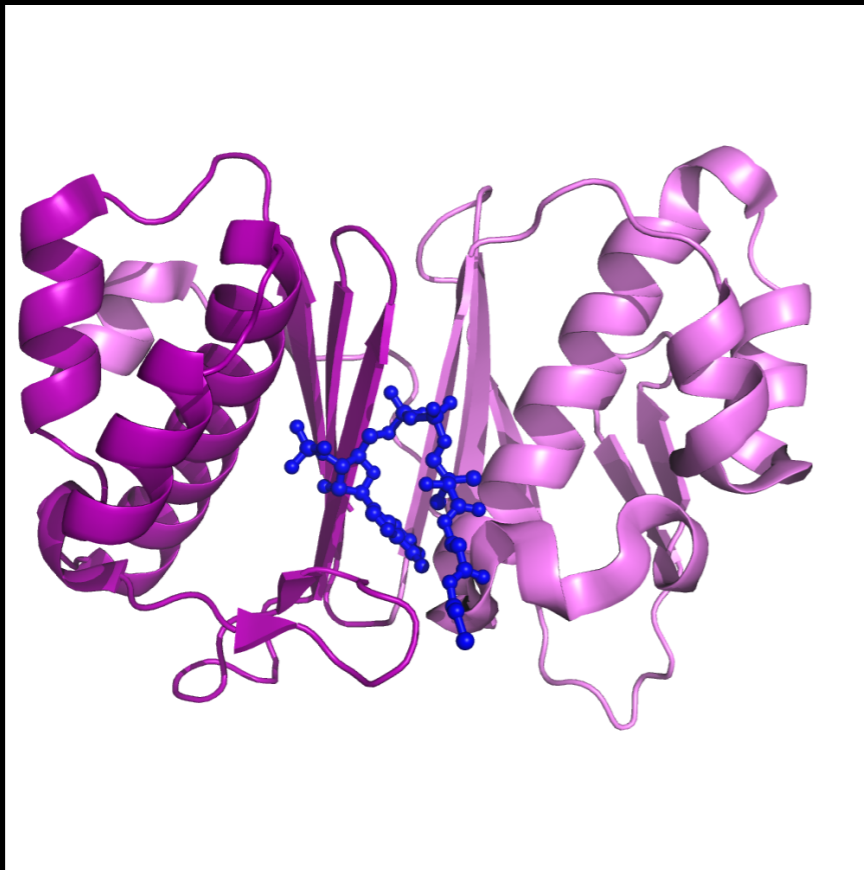
uberPOOL matches you with another rider heading in the same direction. It adds only a few minutes, and you both save big. Trips are up to 50% less than uberX. From home to work to play, uberPOOL gets you there for way, way less.

[SIGN UP FOR UBER](#)

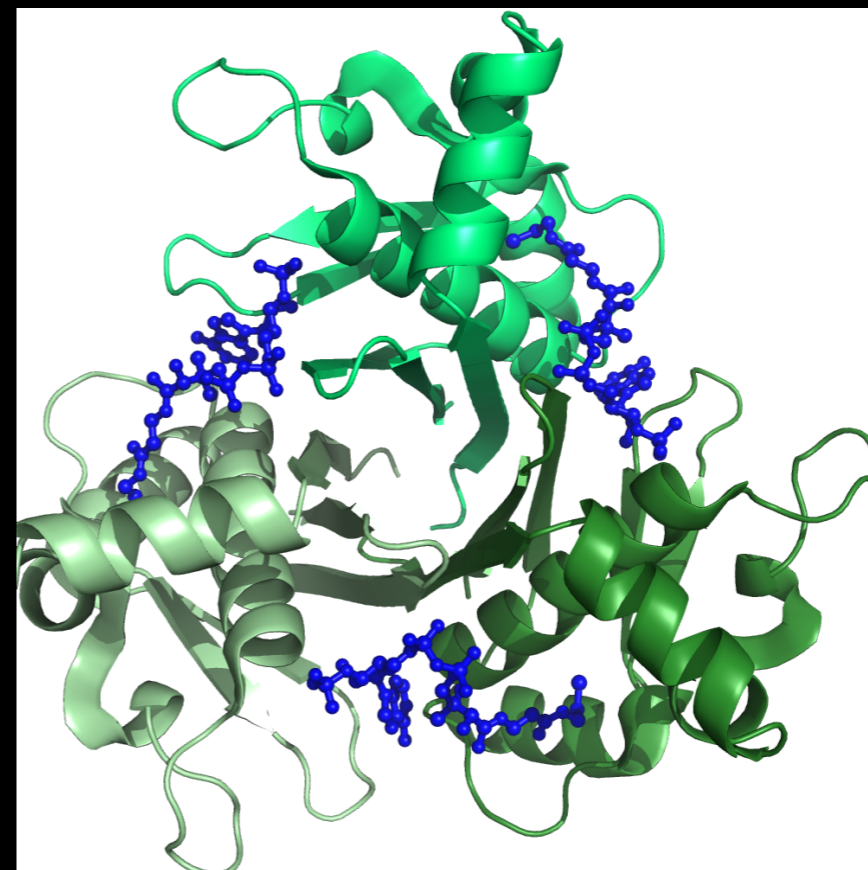
Optimization of expensive functions arises when optimizing physics-based models



Optimization of expensive functions arises in drug and materials discovery



Sfp
(a protein-modifying enzyme)



AcpS
(another protein-modifying enzyme)

Bayesian Optimization is one way to optimize expensive functions

Assume a Bayesian prior on F
(usually a Gaussian process prior)

while (budget is not exhausted) {

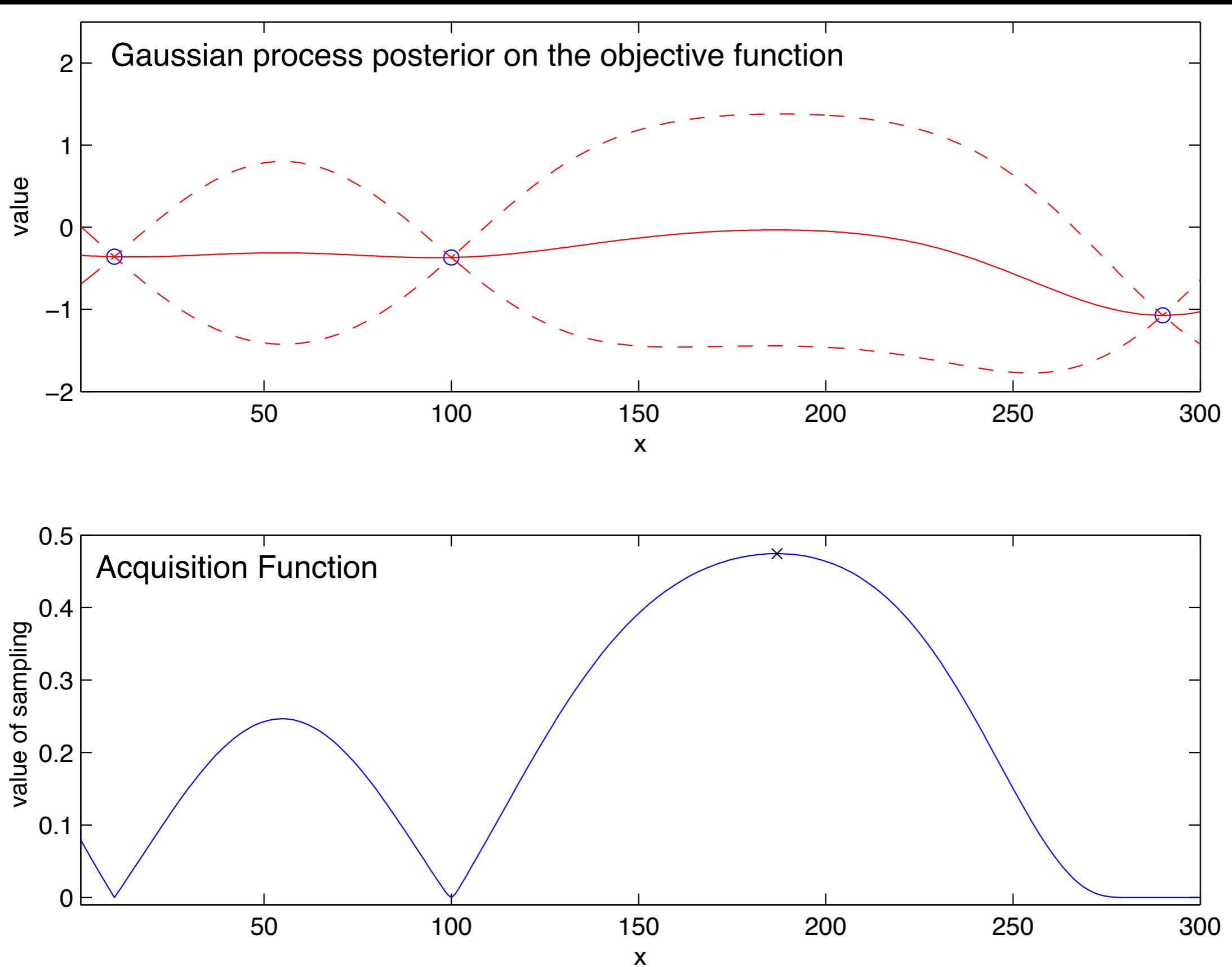
 Find x that maximizes $\text{acquisition}(x, \text{posterior})$

 Sample x & observe $F(x)$

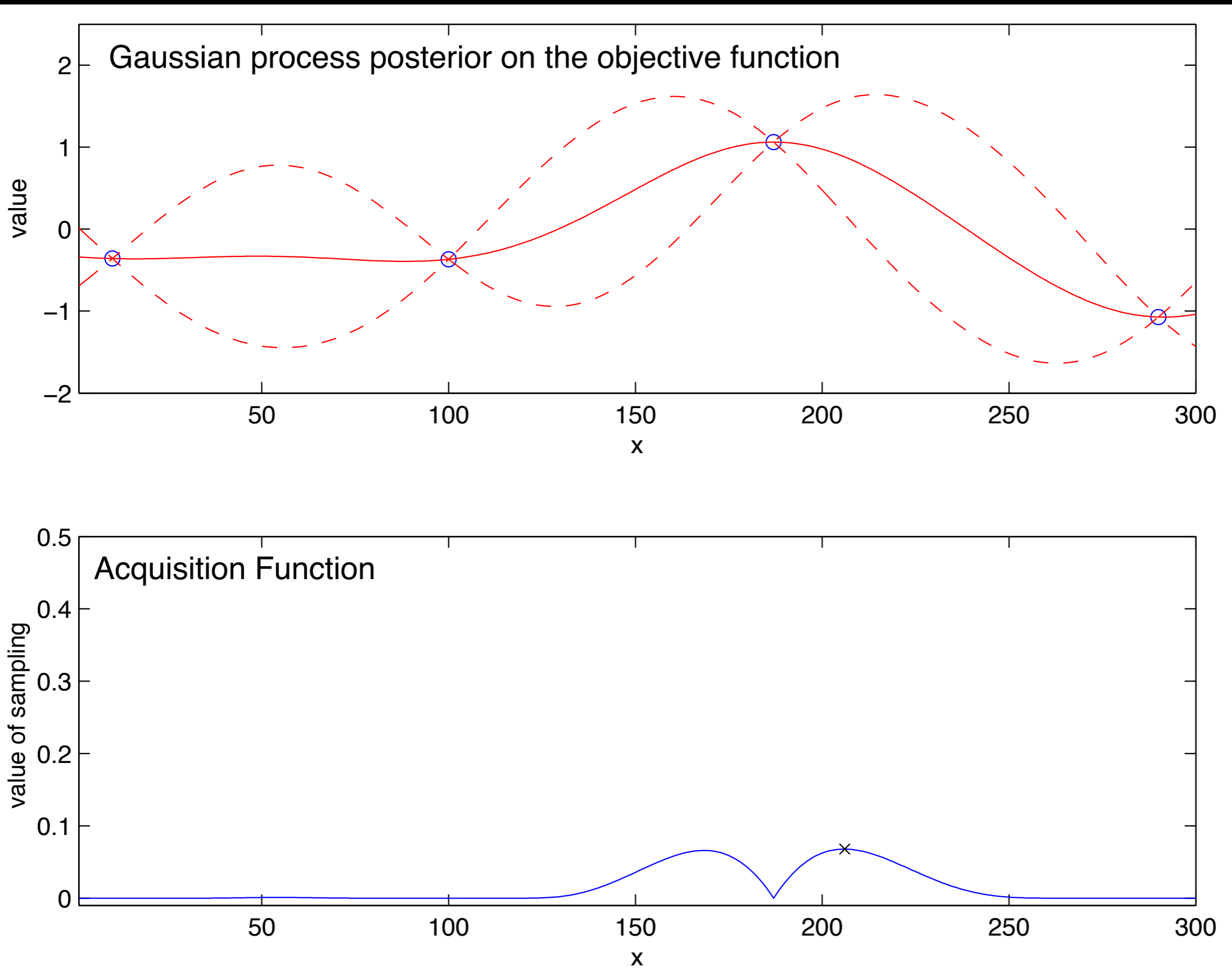
 Update the posterior distribution on F

}

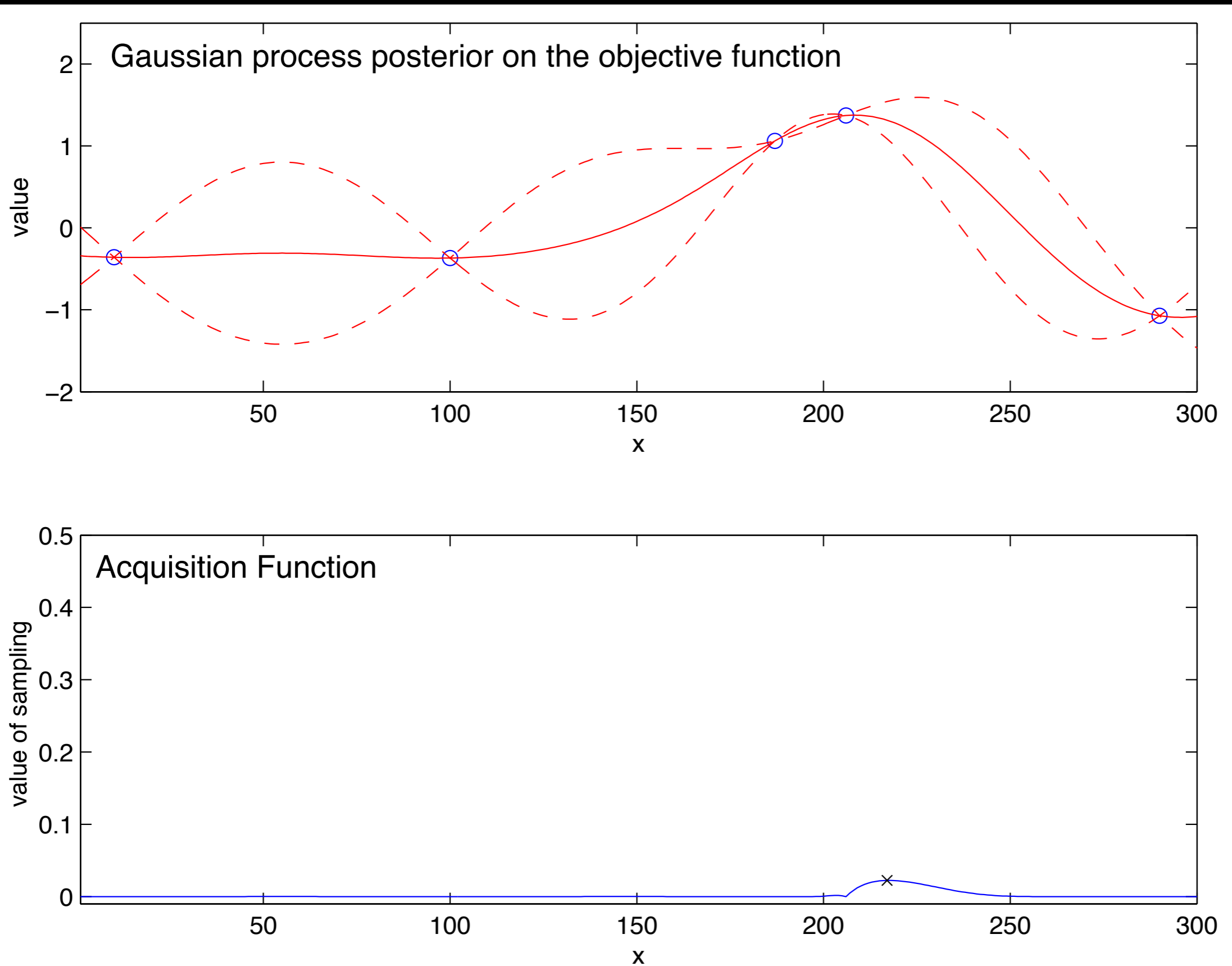
Here's Bayesian optimization, optimizing a 1-dim objective



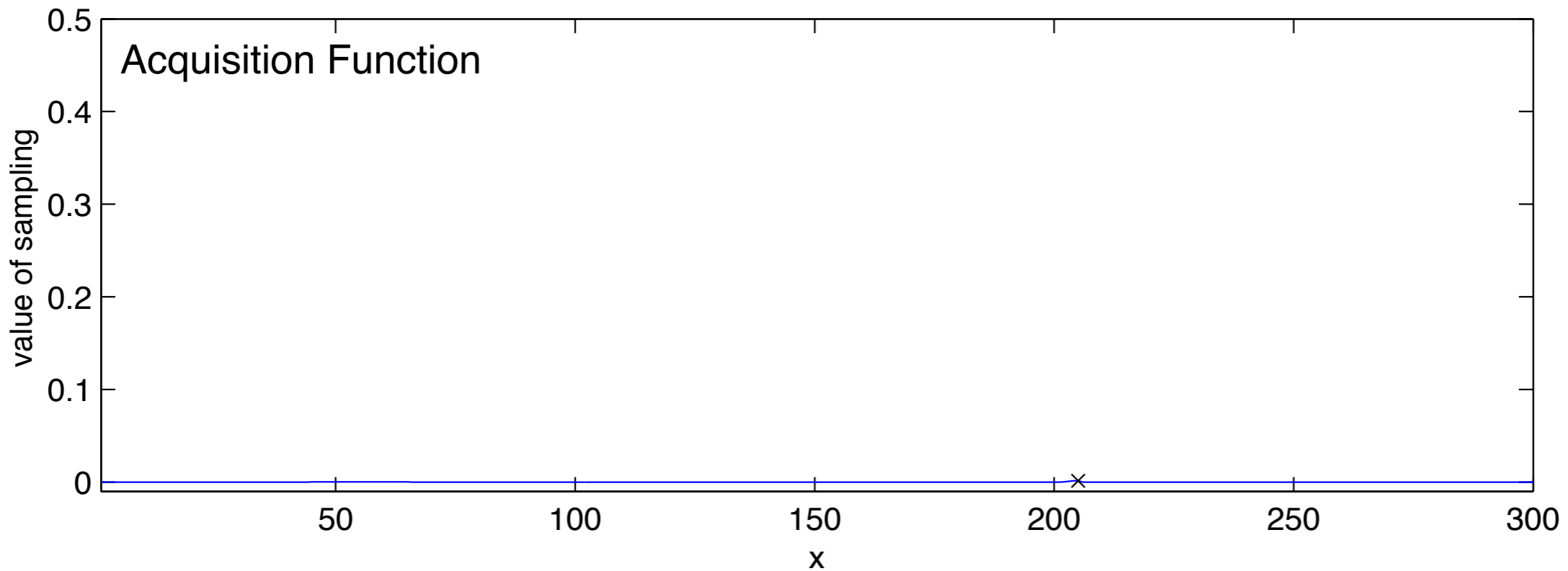
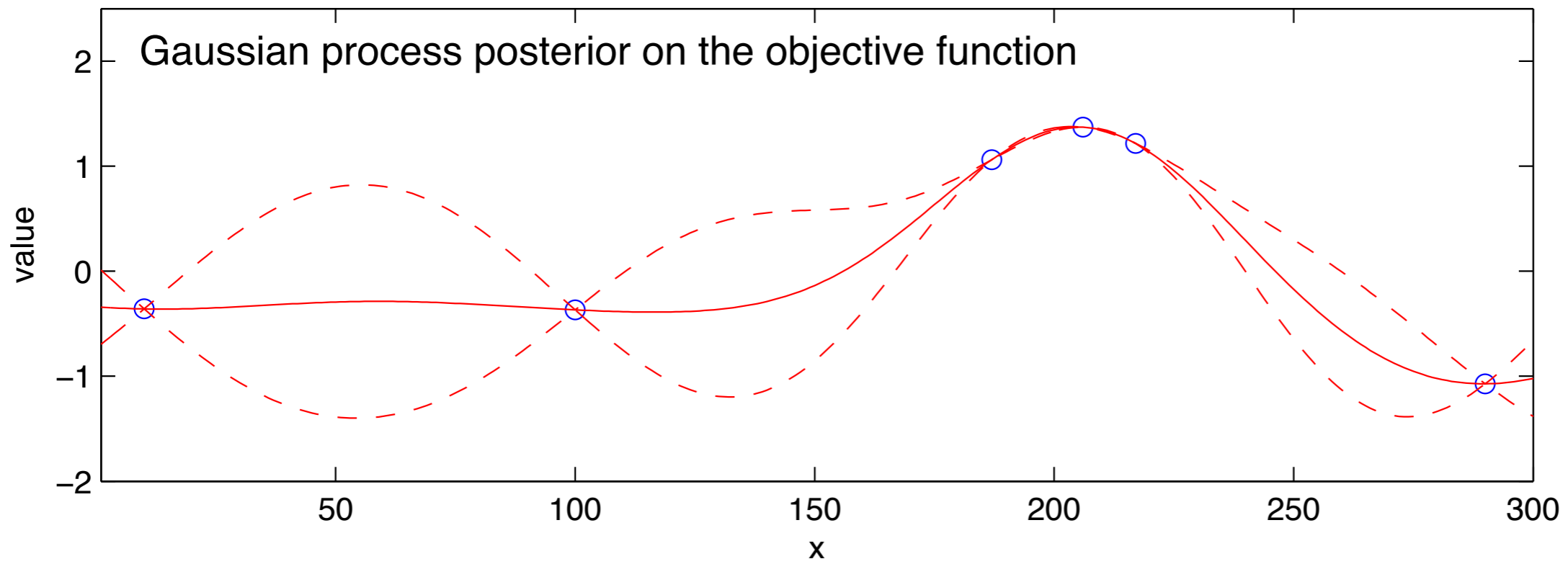
Here's Bayesian optimization, optimizing a 1-dim objective



Here's Bayesian optimization, optimizing a 1-dim objective



Here's Bayesian optimization, optimizing a 1-dim objective



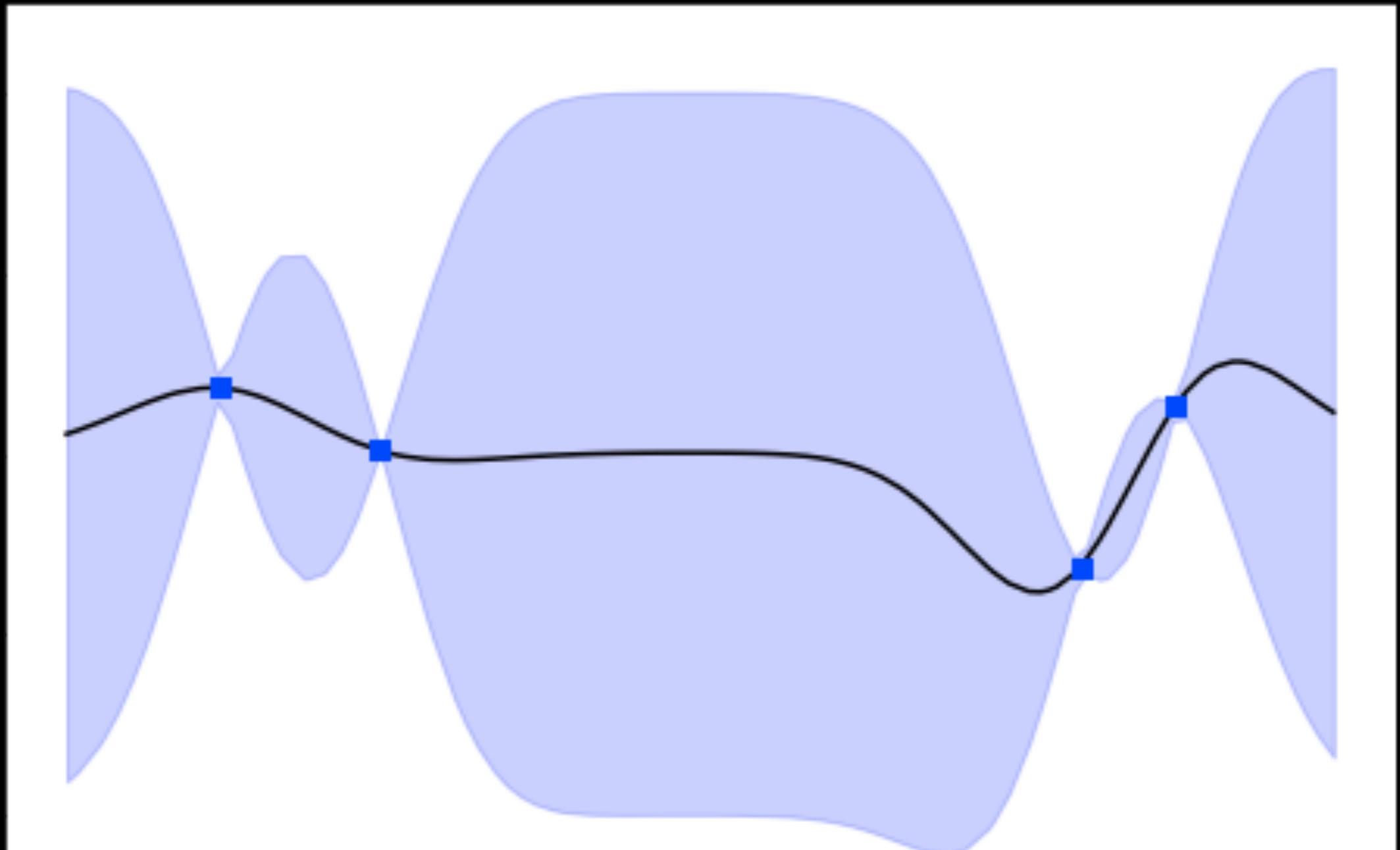
Agenda

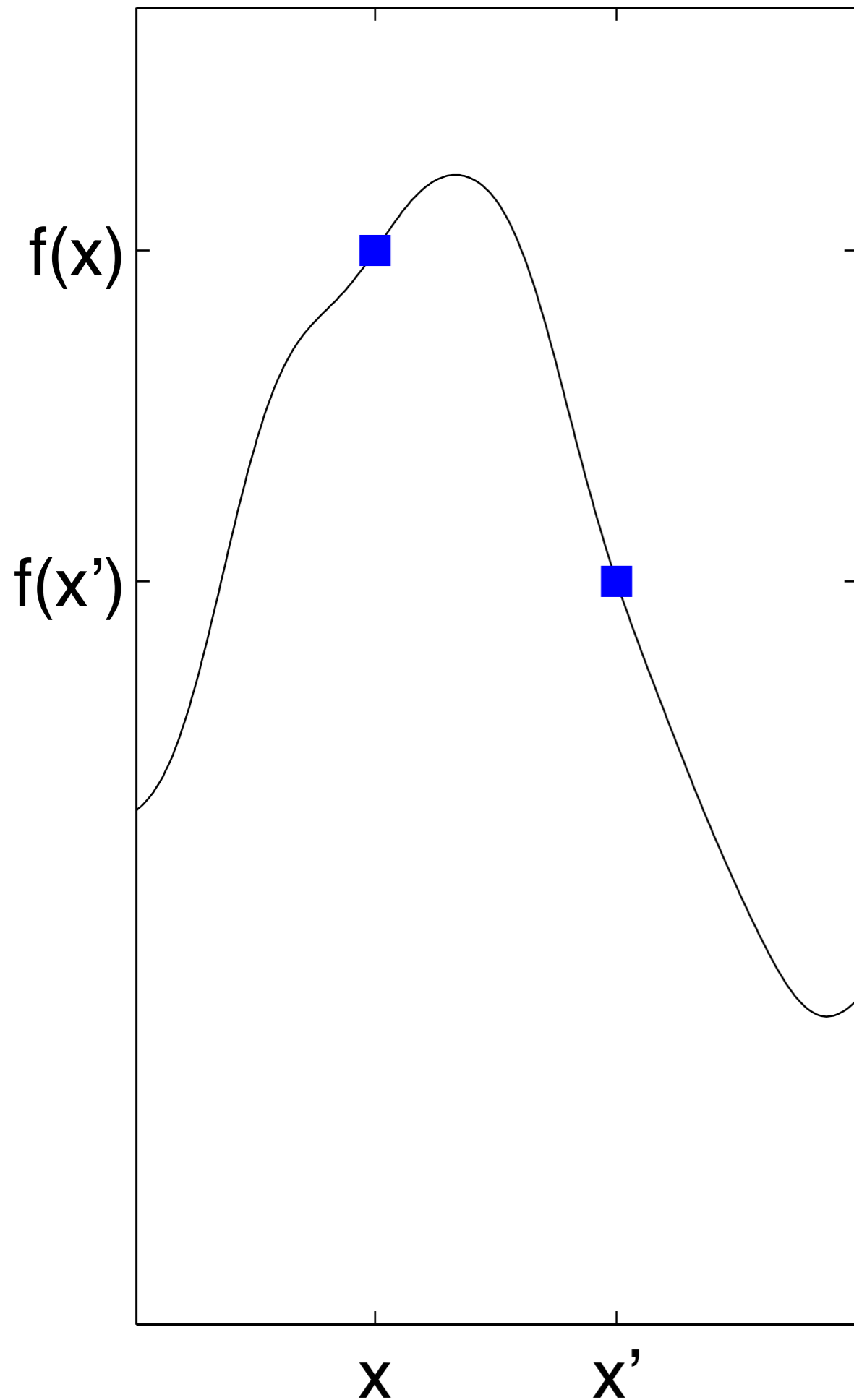
- Gaussian process regression
- Expected improvement
- Beyond the standard problem

Agenda

- **Gaussian process regression**
- Expected improvement (EI)
- Parallel BayesOpt (using EI)
- Parallel BayesOpt with gradients and/or noise (using KG)
- Code & Research directions

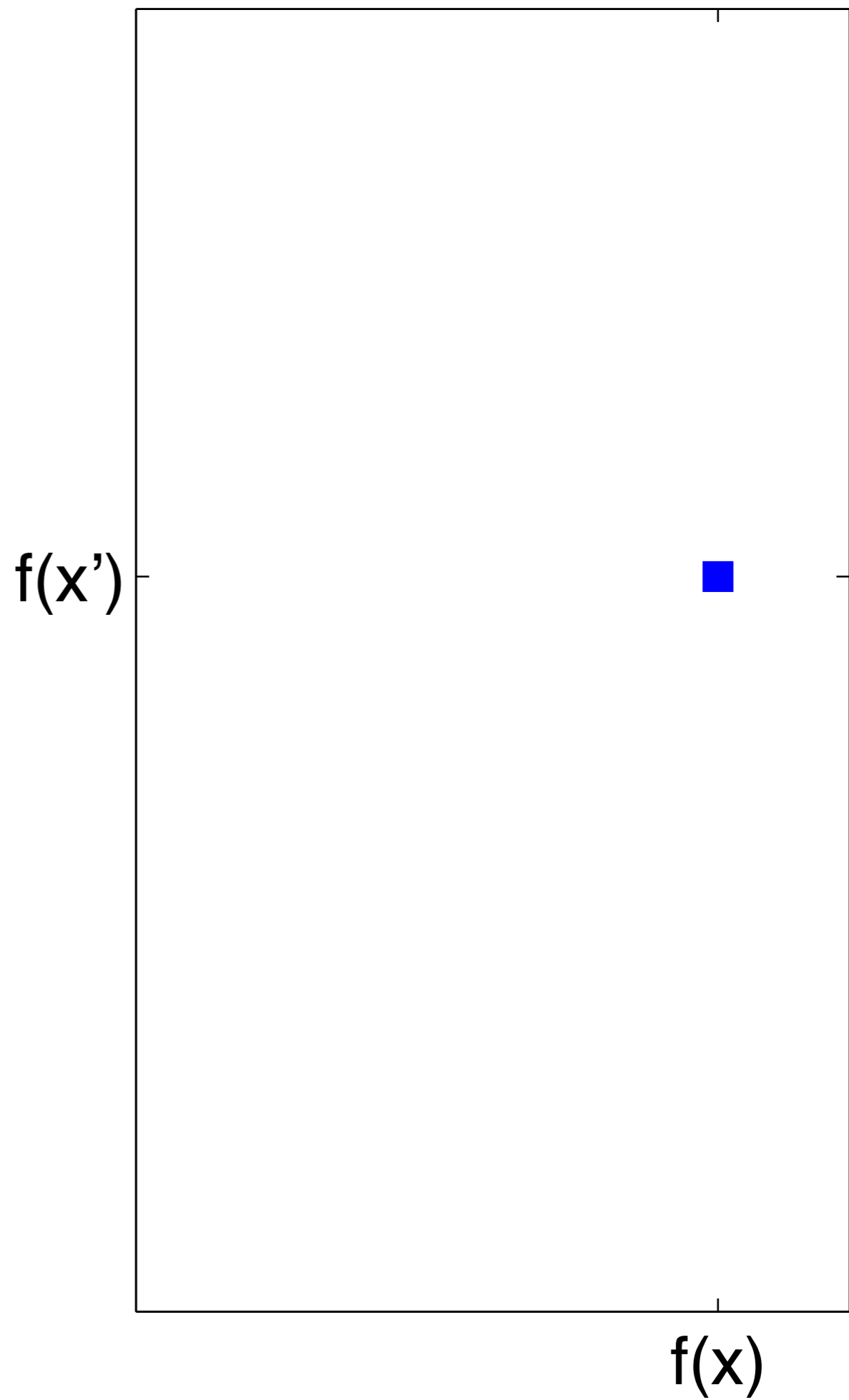
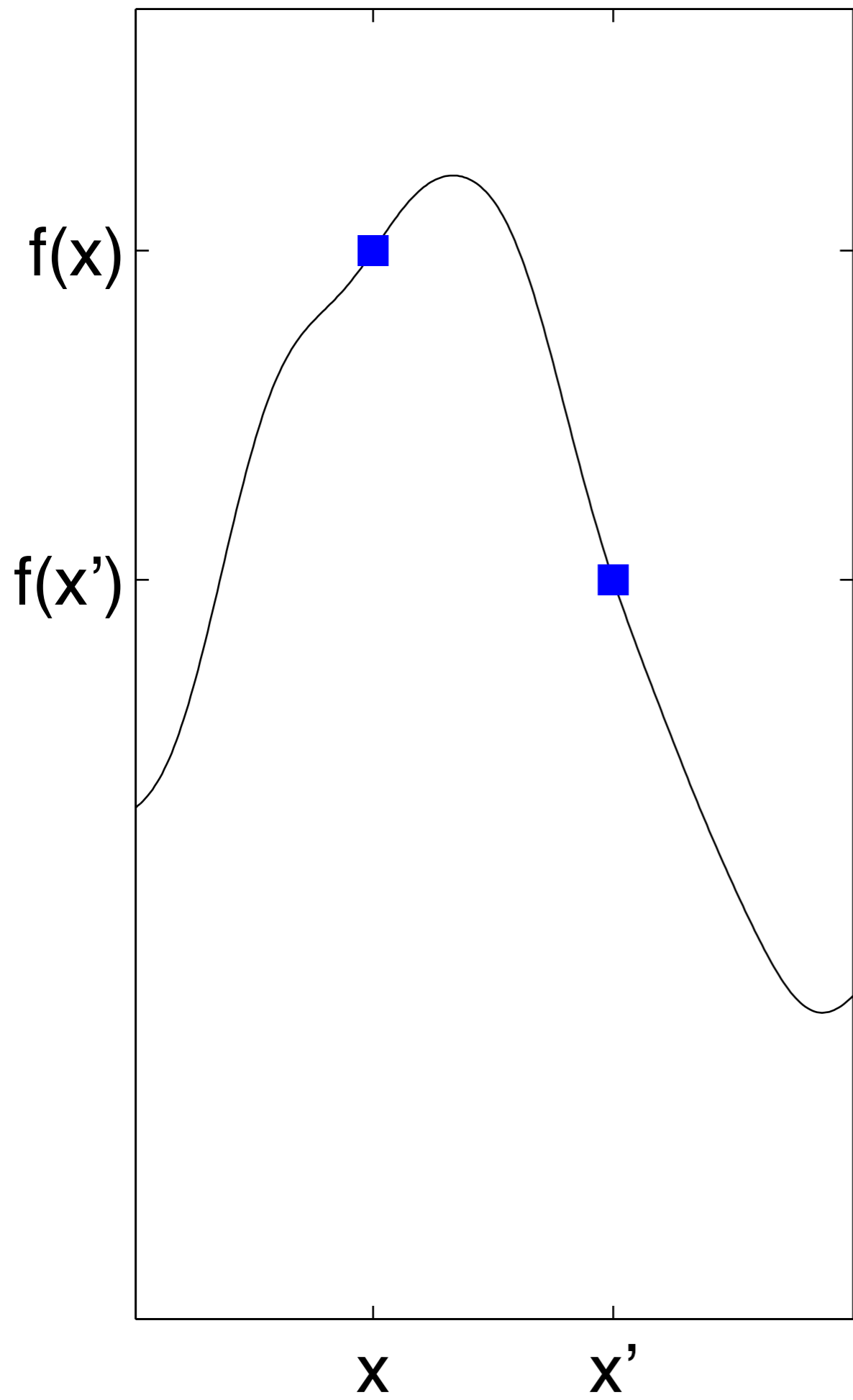
Bayesian optimization usually uses
Gaussian process regression

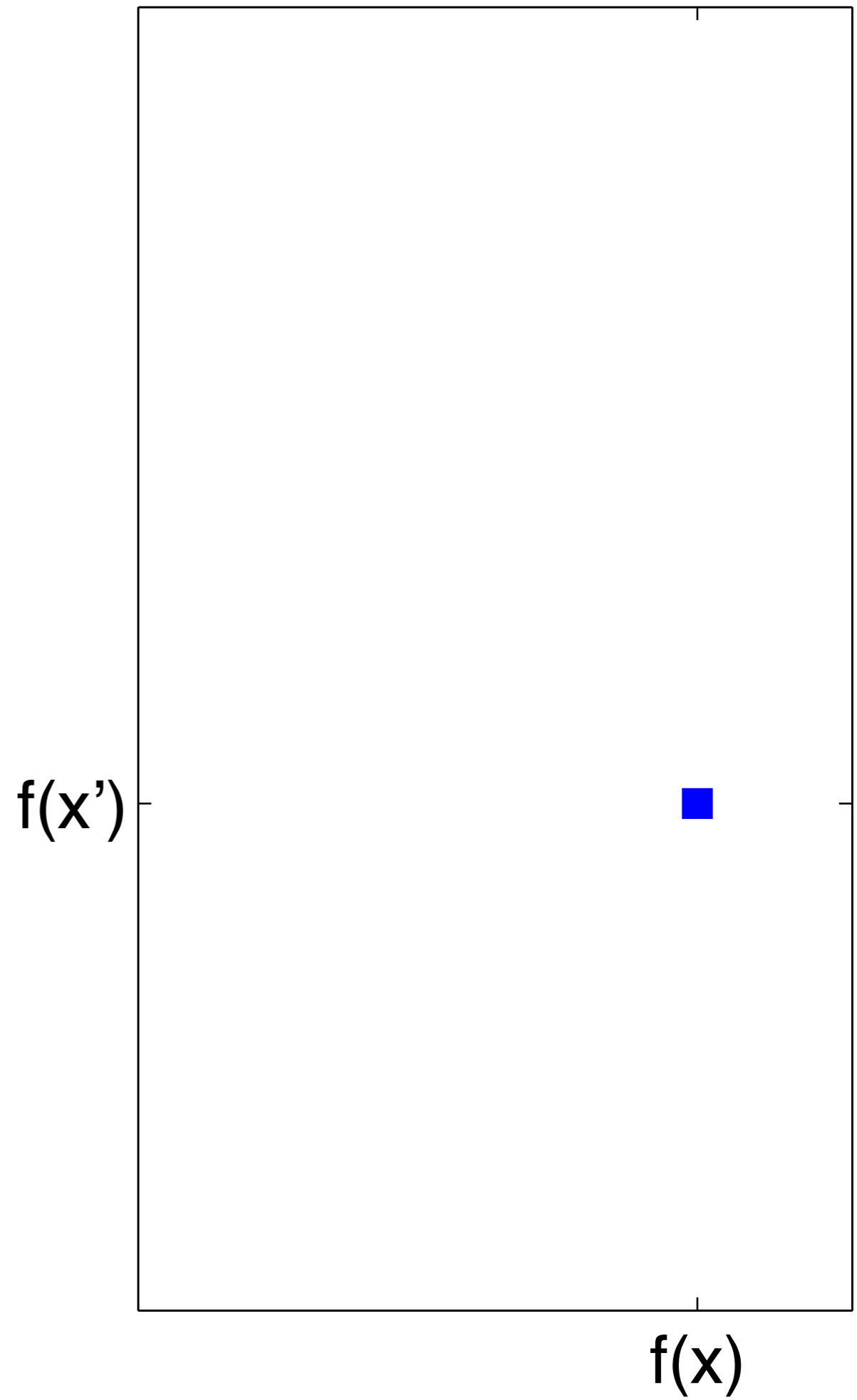
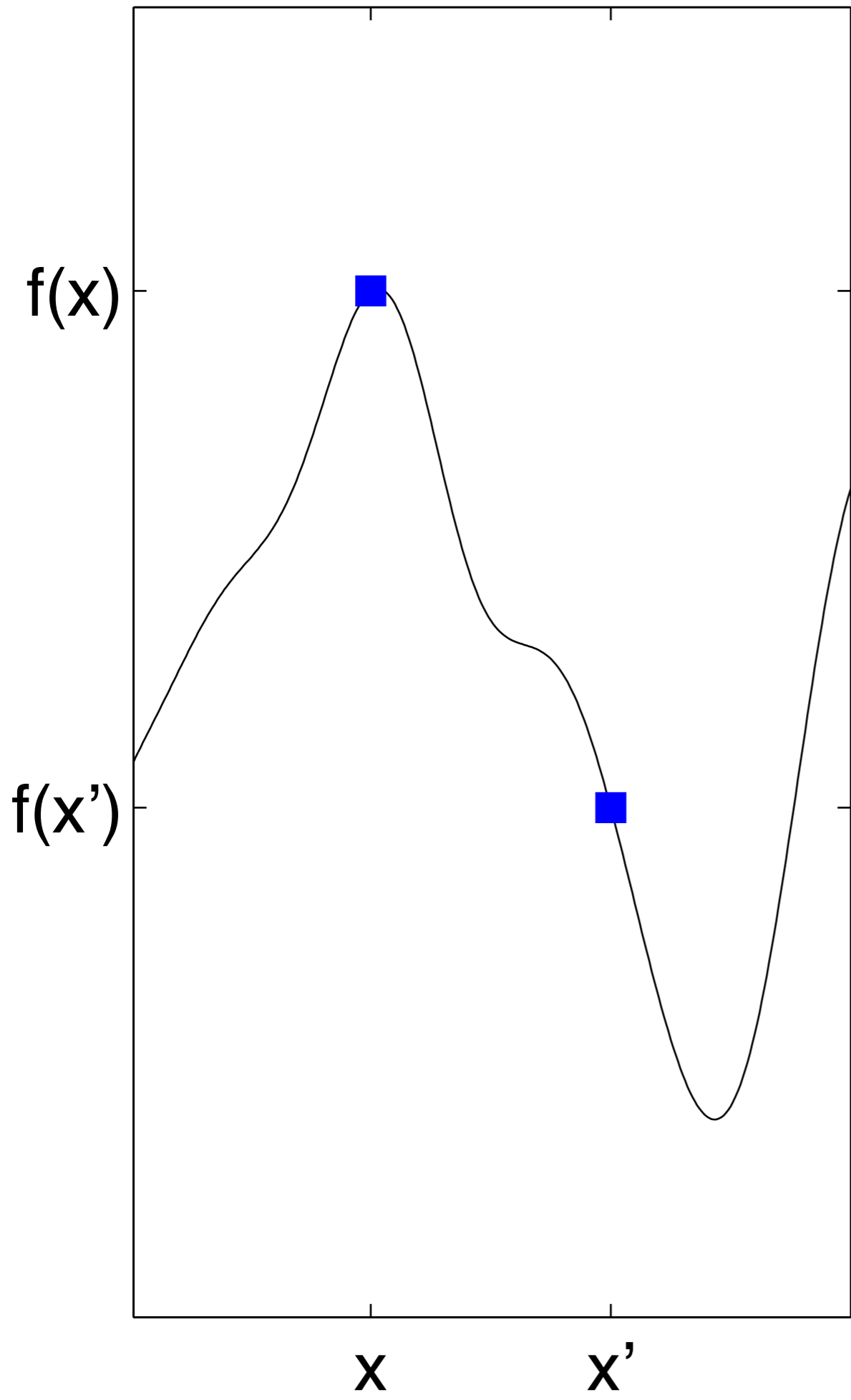


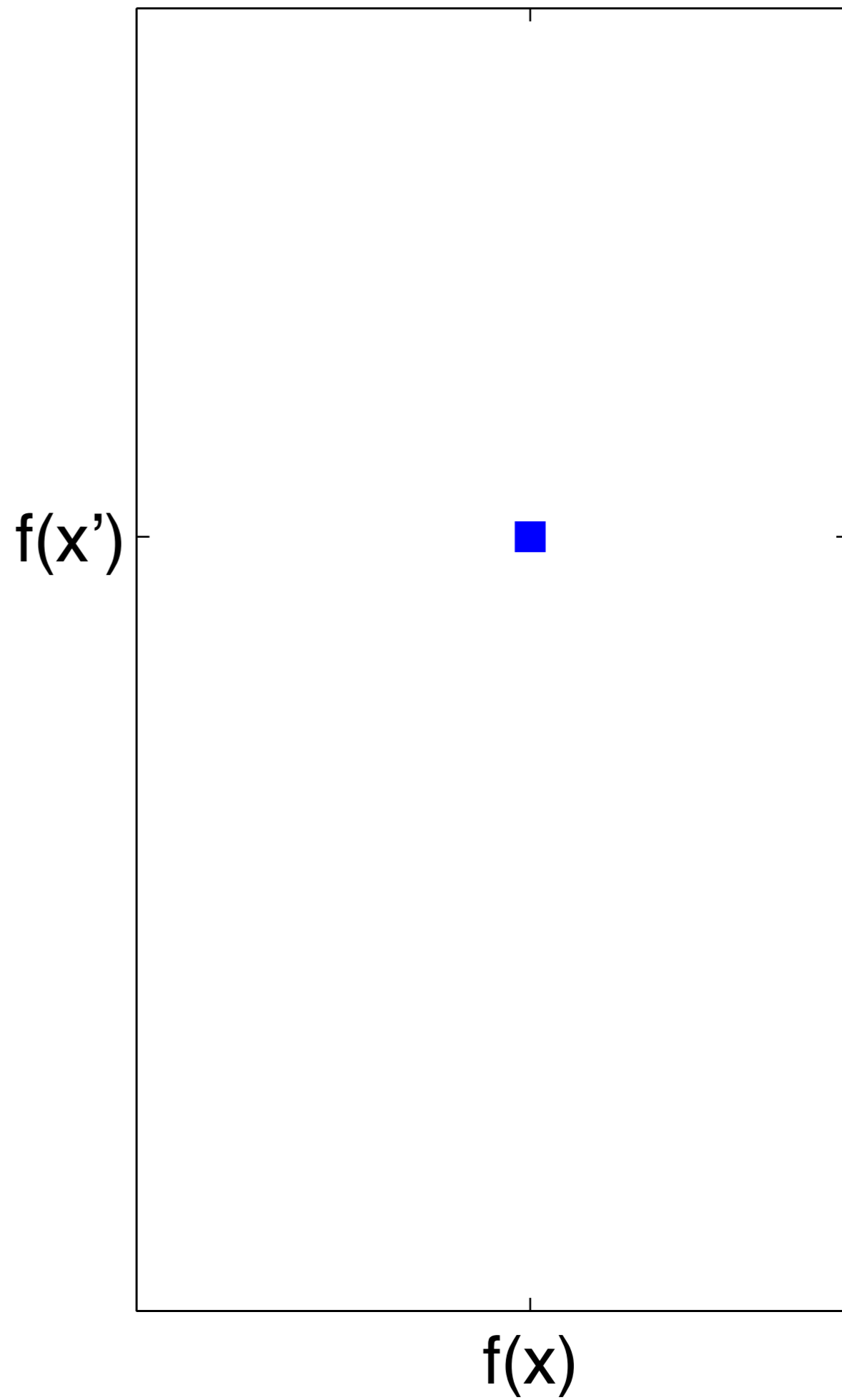
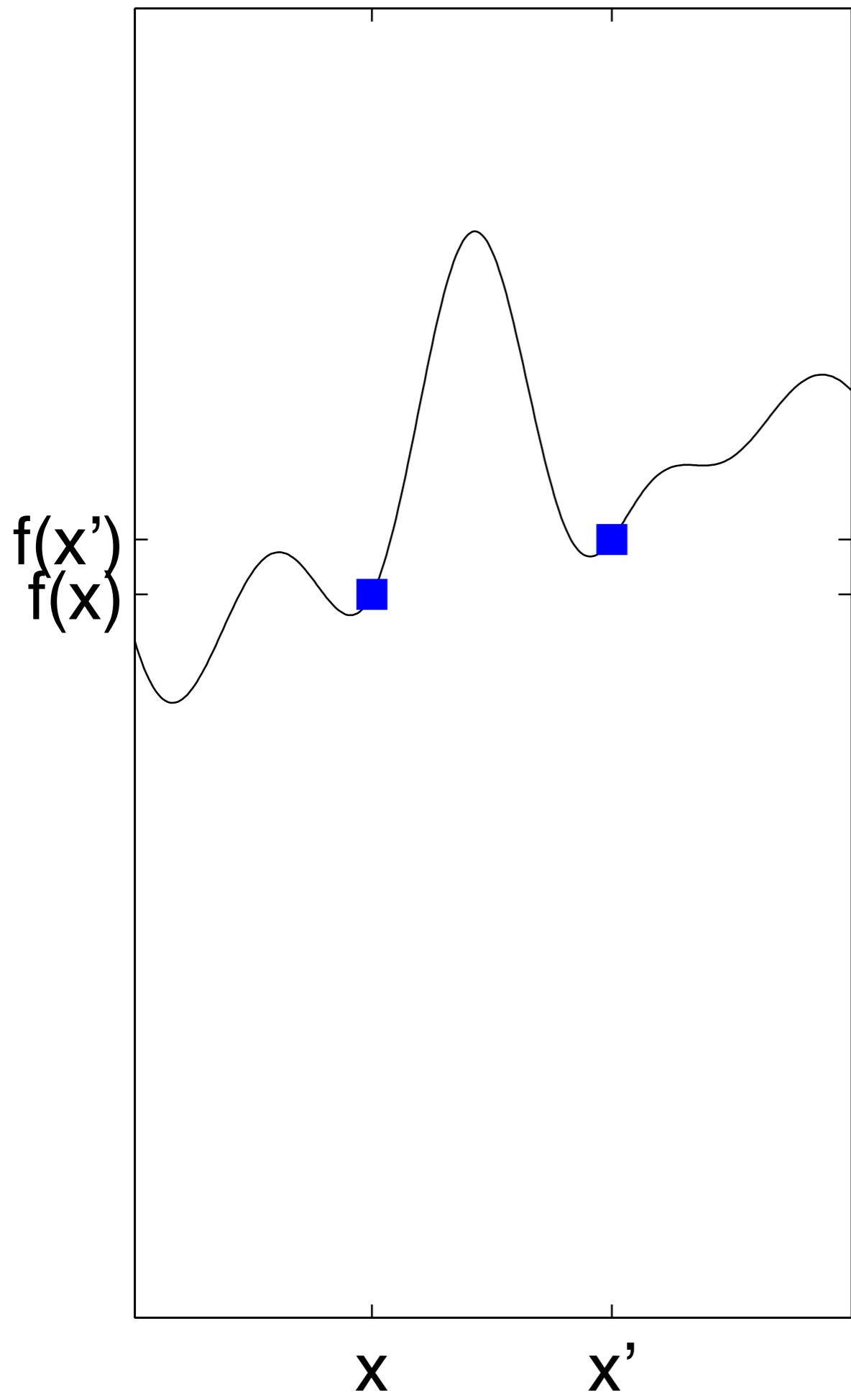


Let's start simply

- Fix 2 points x and x'
- Consider the vector $[f(x), f(x')]$



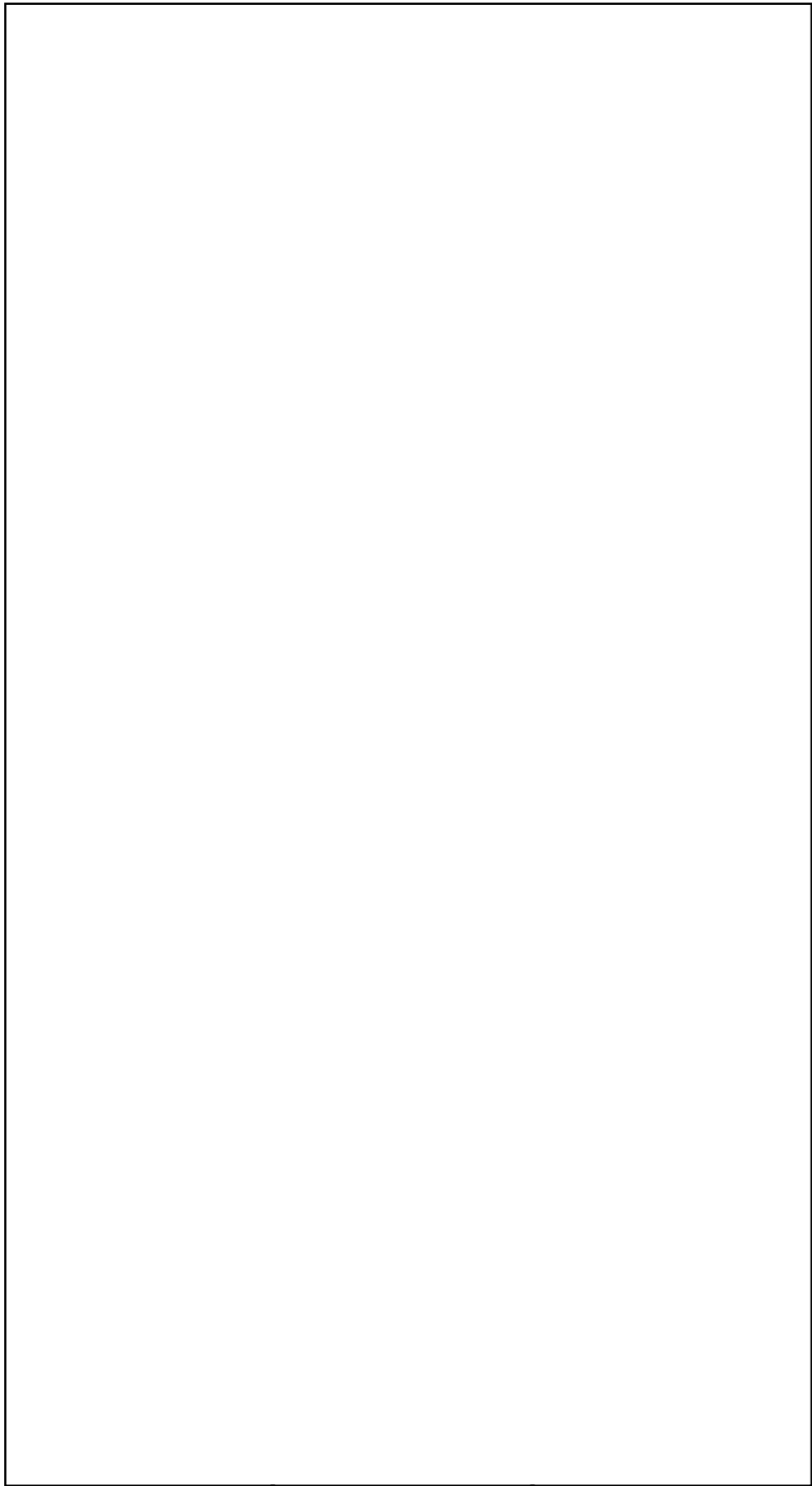




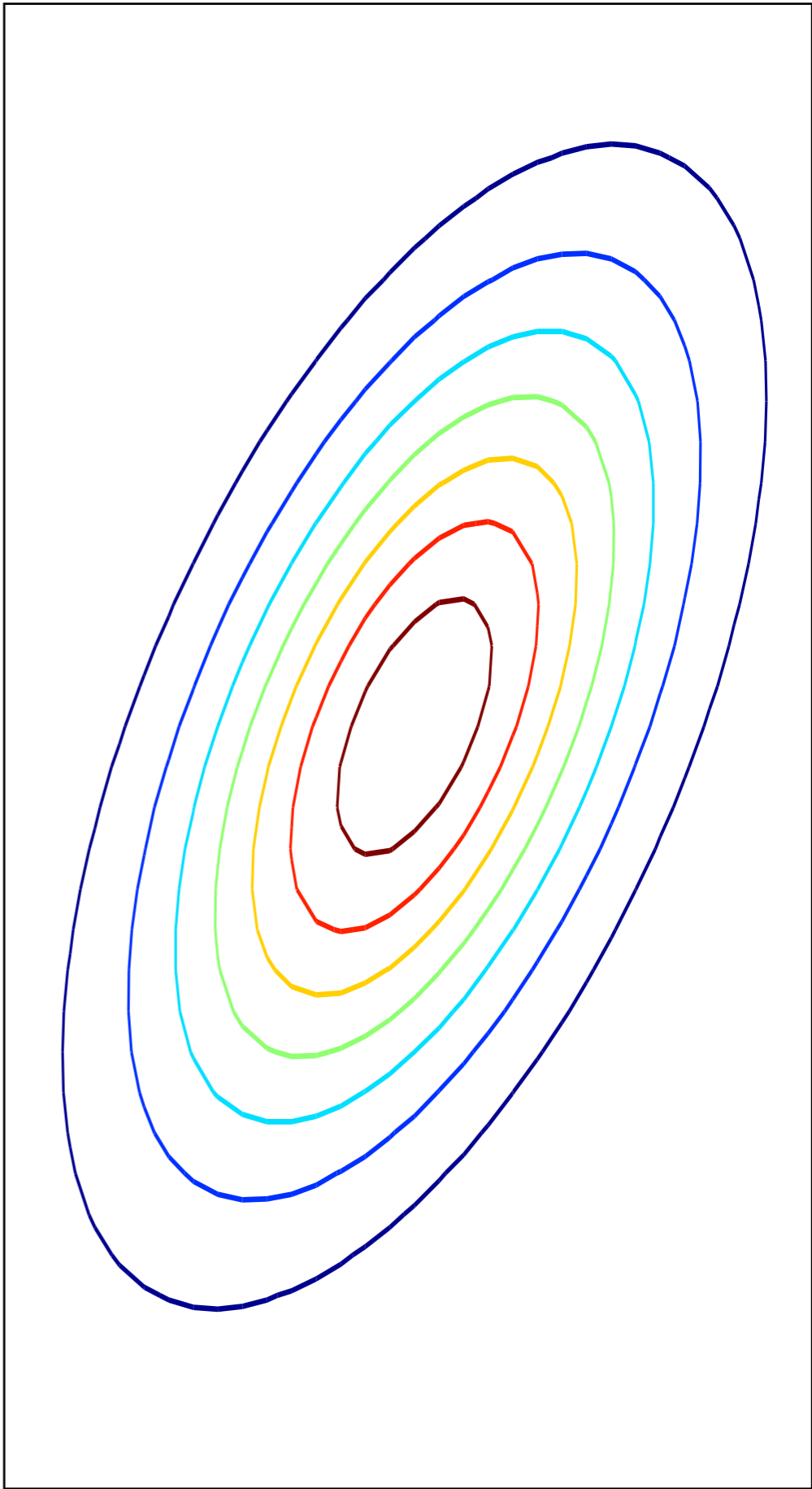
Let's place a multivariate normal prior on $[f(x), f(x')]$

$$\begin{bmatrix} f(x) \\ f(x') \end{bmatrix} \sim N \left(\begin{bmatrix} \mu_0(x) \\ \mu_0(x') \end{bmatrix}, \begin{bmatrix} \Sigma_0(x, x) & \Sigma_0(x, x') \\ \Sigma_0(x', x) & \Sigma_0(x', x') \end{bmatrix} \right)$$

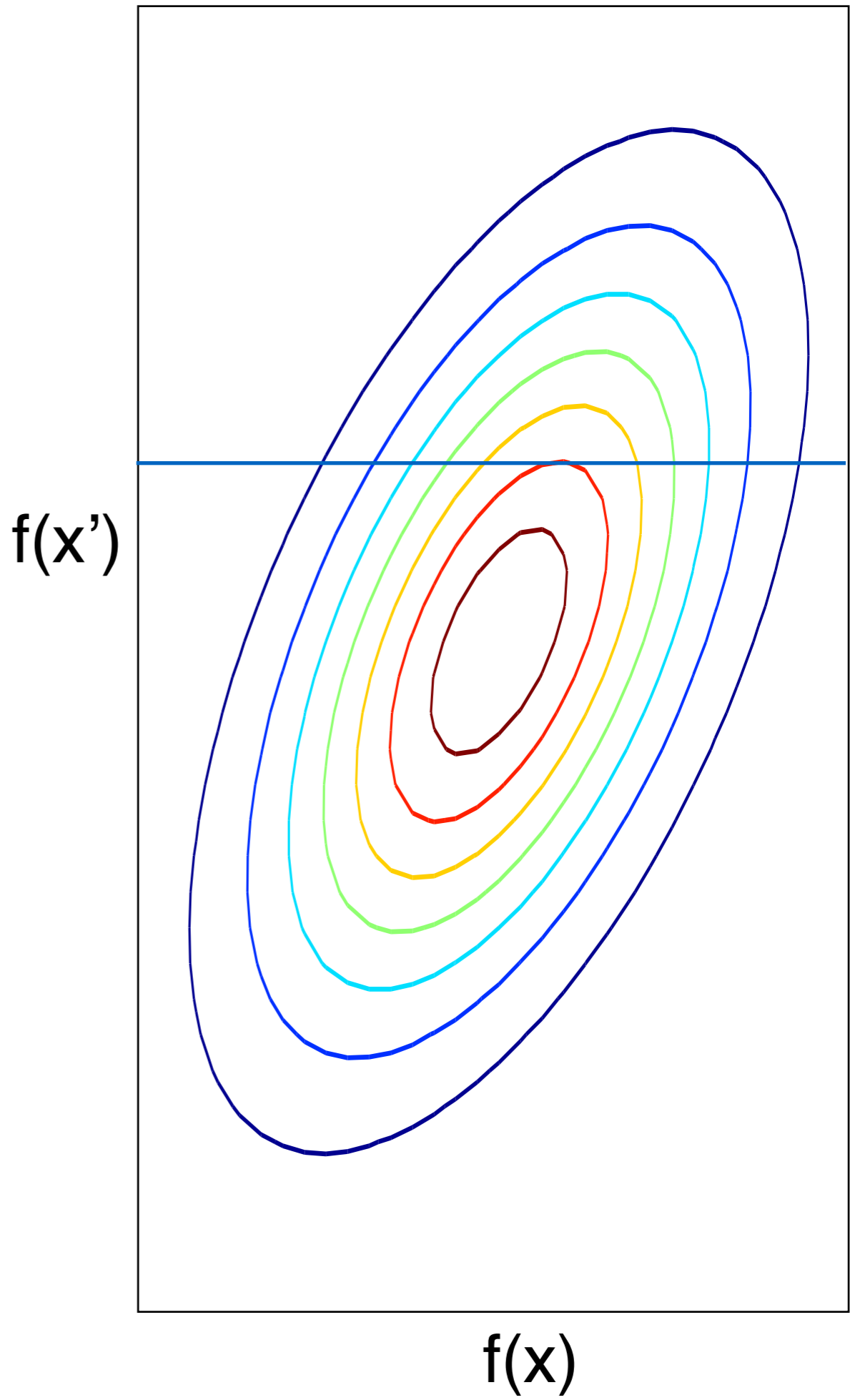
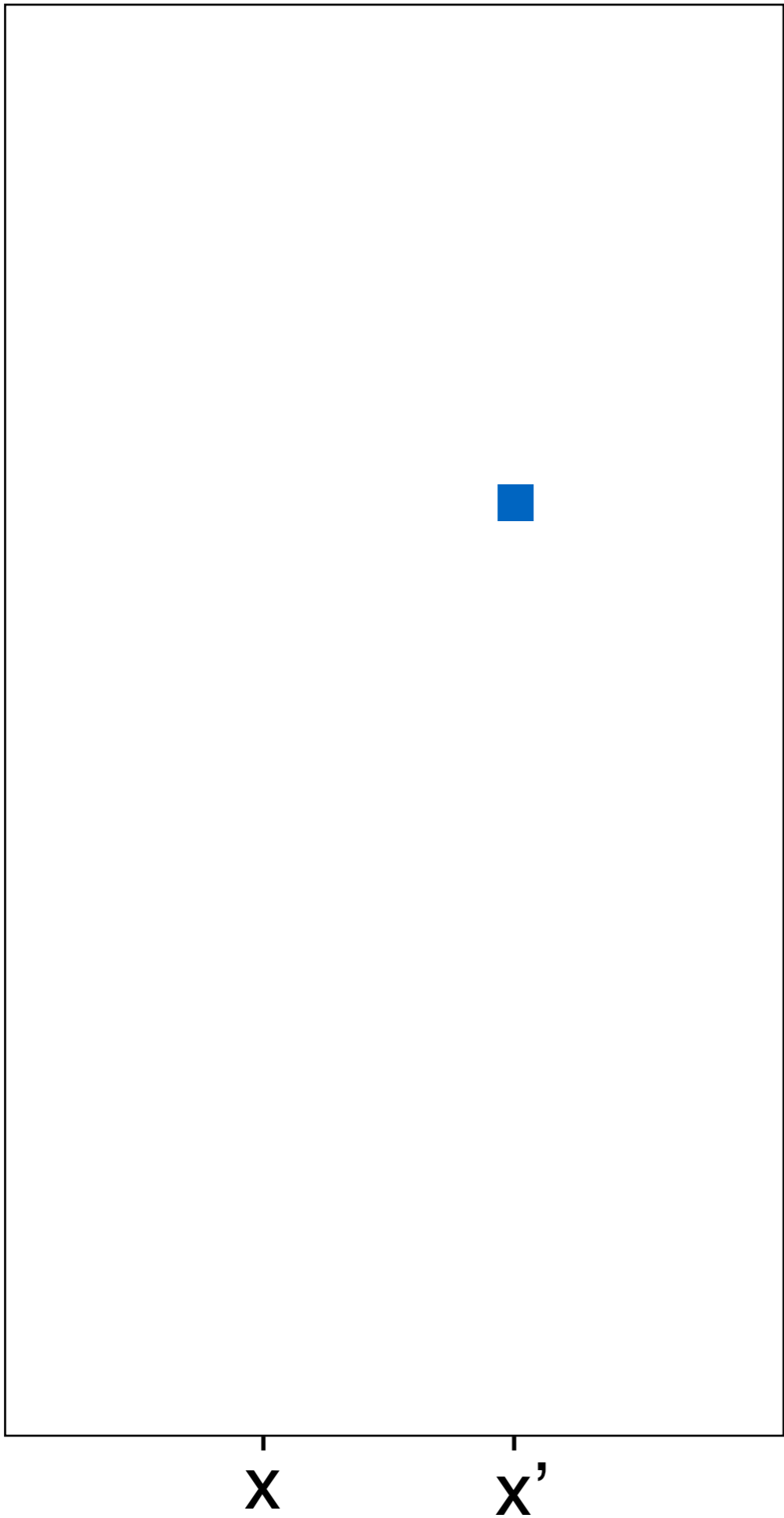
- Σ_0 is a function that decreases with $\|x-x'\|$
- μ is another function, often set to a constant

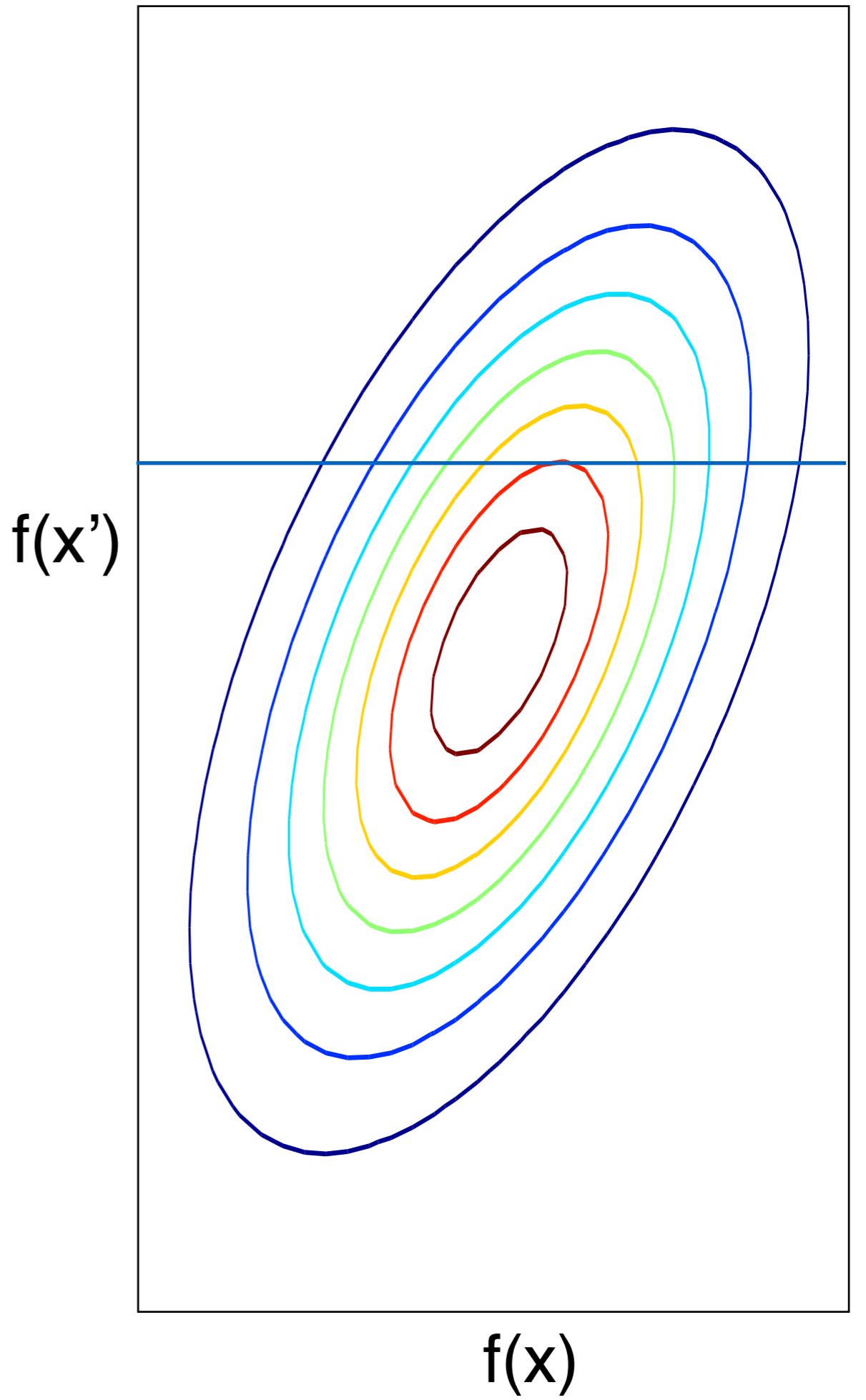
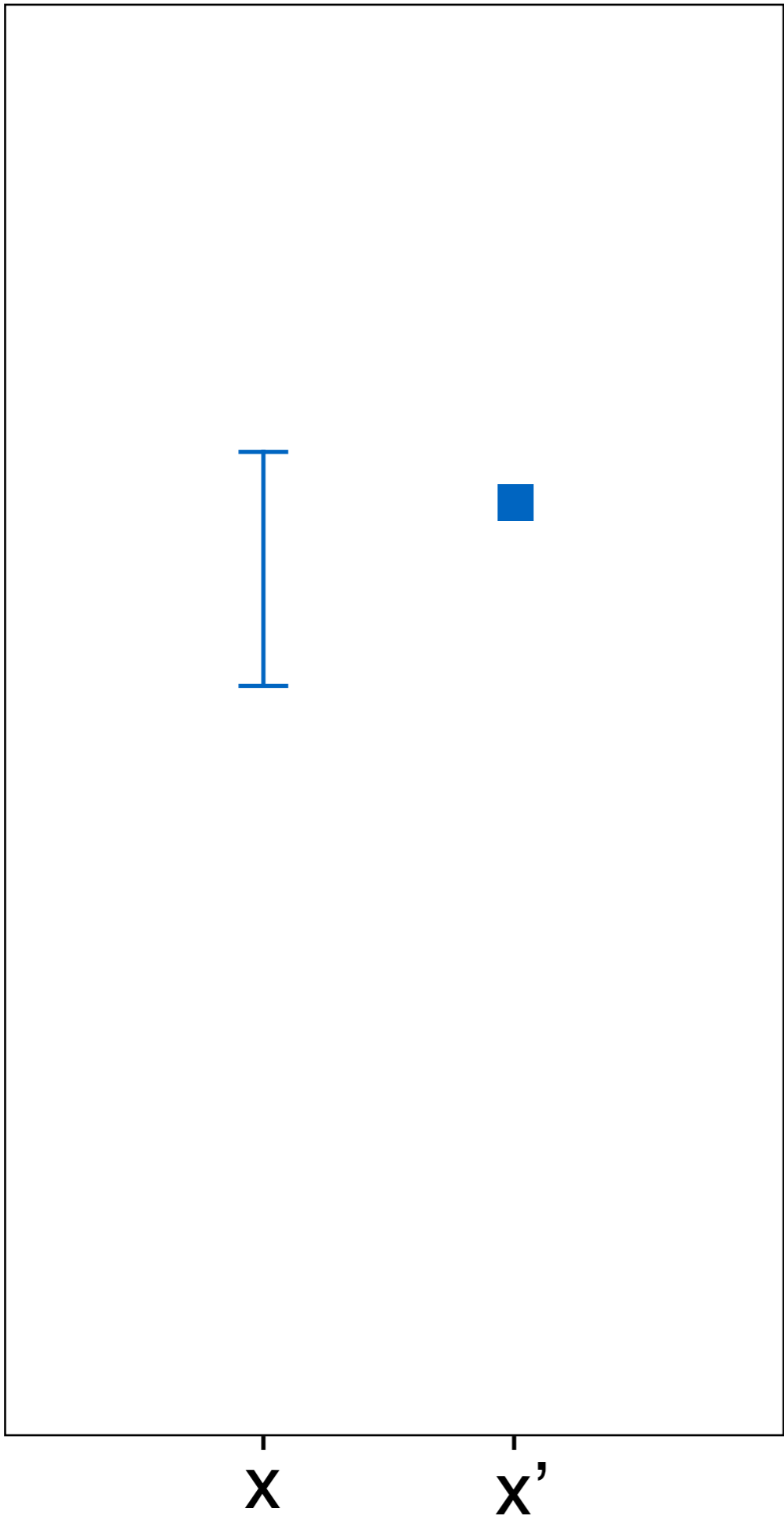


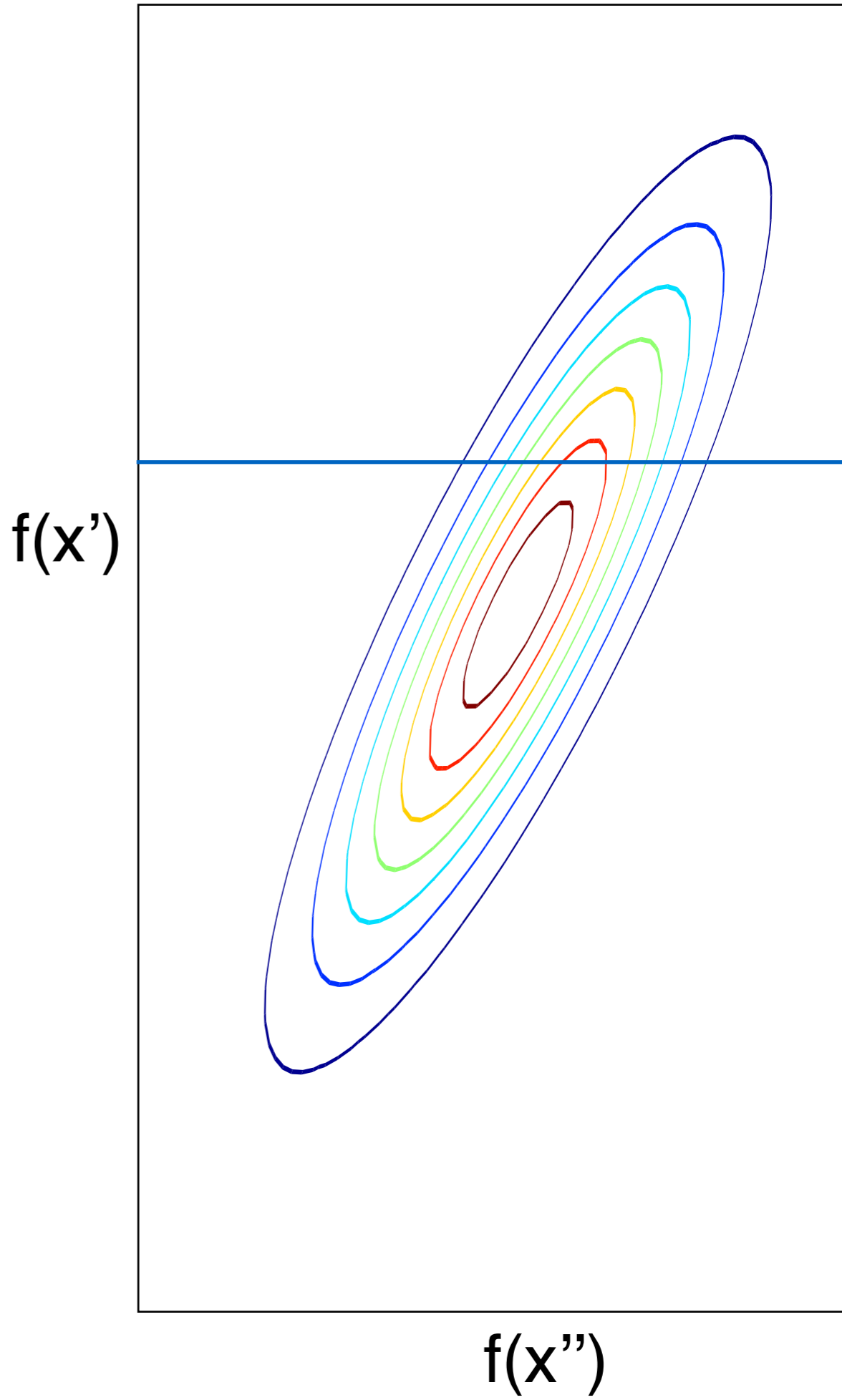
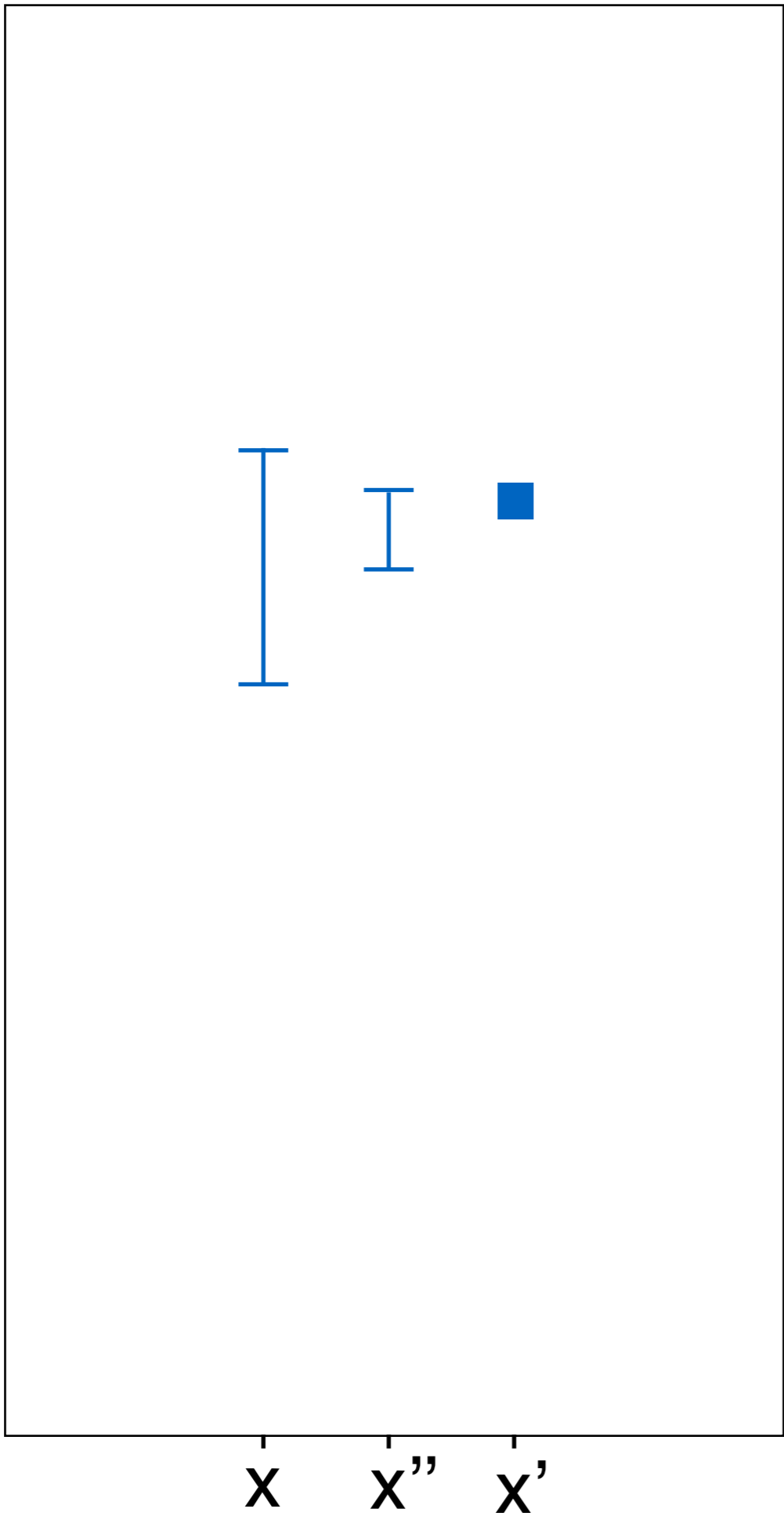
$f(x')$



$f(x)$







Gaussian Process Regression

- A prior on a function f is a **Gaussian process prior** if, for any collection of points x_1, \dots, x_k ,

$$\begin{bmatrix} f(x_1) \\ \vdots \\ f(x_k) \end{bmatrix} \sim N \left(\begin{bmatrix} \mu_0(x_1) \\ \vdots \\ \mu_0(x_k) \end{bmatrix}, \begin{bmatrix} \Sigma_0(x_1, x_1) & \dots & \Sigma_0(x_1, x_k) \\ \vdots & \ddots & \vdots \\ \Sigma_0(x_k, x_1) & \dots & \Sigma_0(x_k, x_k) \end{bmatrix} \right)$$

- we call Σ_0 the “kernel” or “covariance function”
- we call μ_0 the “mean function”

We can compute the posterior analytically

- The posterior on $f(x')$ given $f(x_1), \dots, f(x_n)$ is normal with mean $\mu_n(x')$ and variance $\sigma_n^2(x')$:
 - $\mu_n(x') = \Sigma_0(x', X_{1:n}) \Sigma_0(X_{1:n}, X_{1:n})^{-1} f(X_{1:n})$
 - $\sigma_n^2(x') = \Sigma_0(x', x') - \Sigma_0(x', X_{1:n}) \Sigma_0(X_{1:n}, X_{1:n})^{-1} \Sigma_0(X_{1:n}, x')$
 - Formula assumes $\mu_0=0$.
There is a similar formula for general μ_0 .

How should we choose the kernel & mean functions?

- 1. Designate a parametric family.
- 2. Take an initial stage of samples from points sampled uniformly at random
- 3. Choose hyperparameters using either:
 - 2a. Maximum likelihood estimation
 - 2b. Maximum *a posteriori* estimation, with a prior on the hyperparameters
 - 2c. Sample the hyperparameters from their posterior

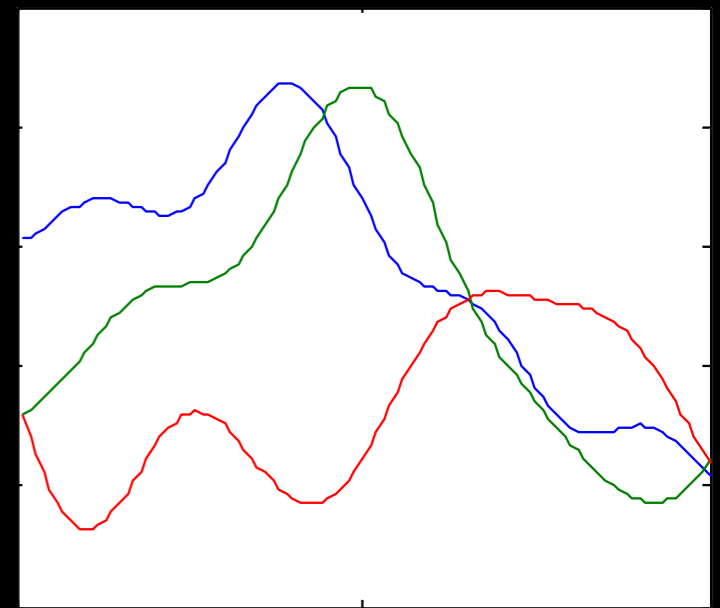
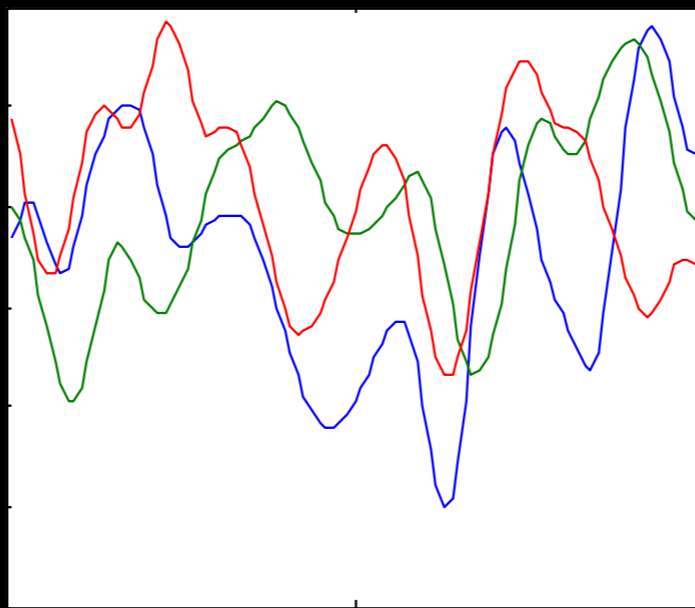
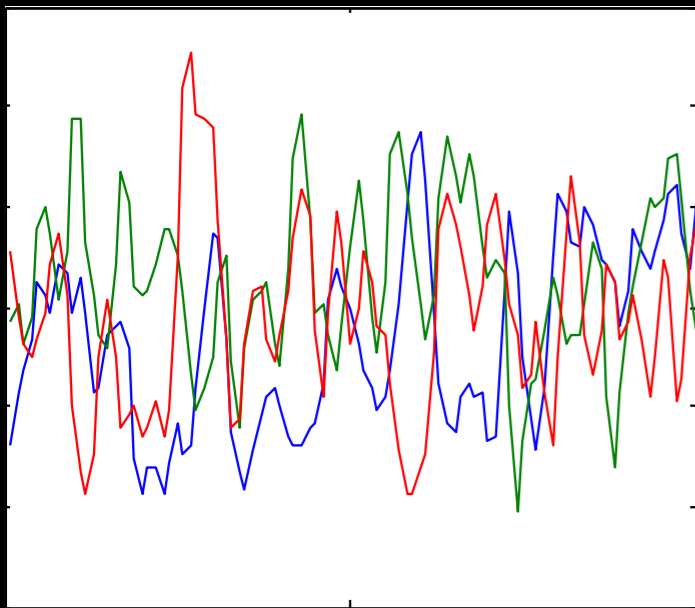
How should we choose the mean function?

- Parametric Family:
 - $\mu_0(x) = \beta$
 - $\mu_0(x) = \beta_0 + \sum_j \beta_j \Psi_j(x)$
- Prior:
 - Flat over \mathbb{R} (non-informative)
 - Uniform
 - Normal

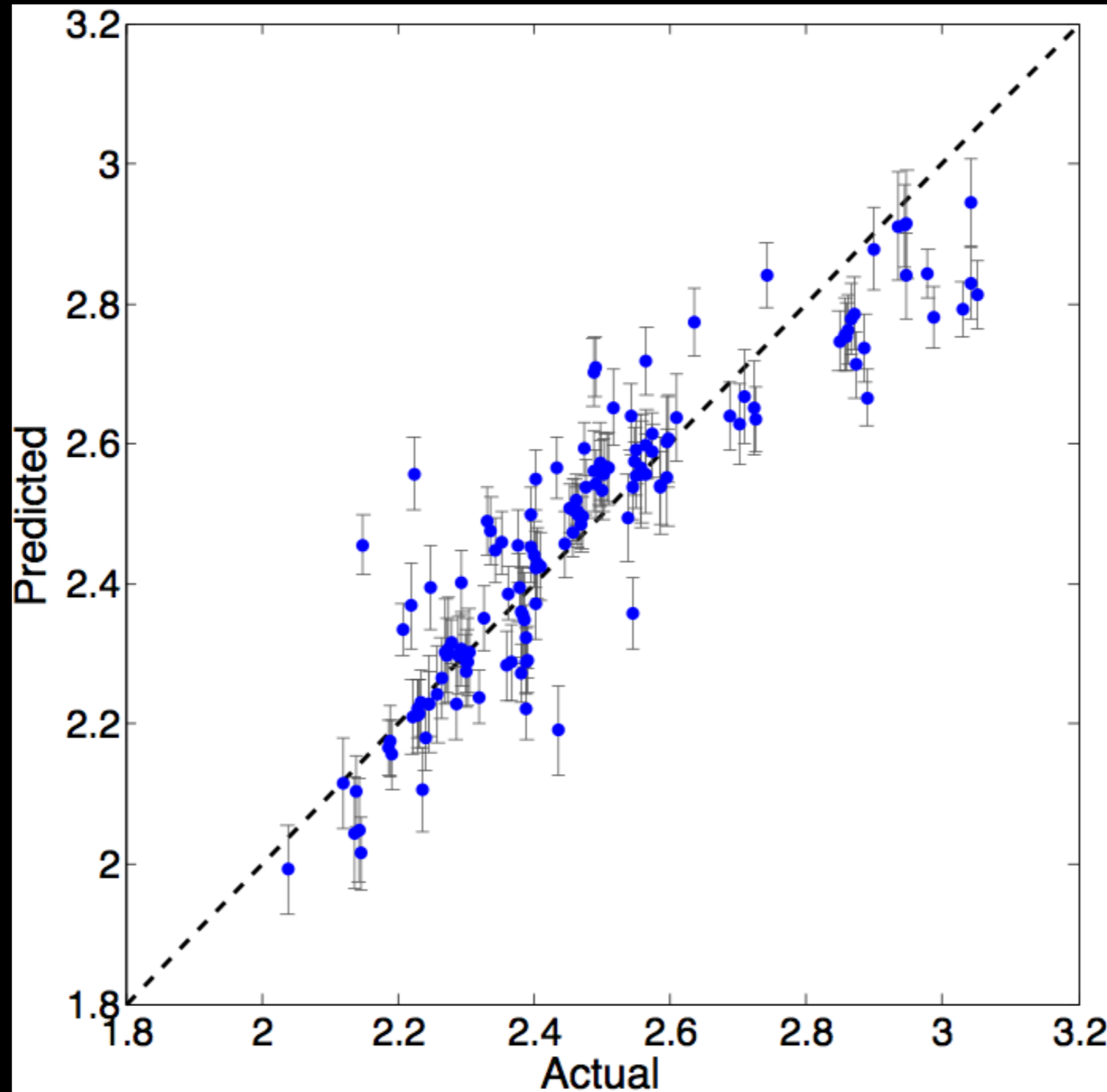
How should we choose the kernel?

- Parametric Family:
 - $\Sigma_0(x, x') = \alpha_0 \exp(-\|x-x'\|^\rho)$
 - $\Sigma_0(x, x') = \alpha_0 2^{1-\nu} z^\nu K_\nu(z) / \Gamma(\nu)$
where $z = (2\nu)^{1/2} \|x-x'\| / \rho$,
 $K_\nu(\cdot)$ is the modified Bessel function
 - The above with x_d rescaled by α_d
- Prior: flat, uniform, or normal

How should we choose the kernel?



Leave one-out cross-validation is worth doing



Noise can be incorporated

$$\mu_n(x') = \Sigma_0(x', x_{1:n}) [\Sigma_0(x_{1:n}, x_{1:n}) + \lambda^2 I]^{-1} f(x_{1:n})$$

$$\sigma_n^2(x') = \Sigma_0(x', x') - \Sigma_0(x', x_{1:n}) [\Sigma_0(x_{1:n}, x_{1:n}) + \lambda^2 I]^{-1} \Sigma_0(x_{1:n}, x')$$

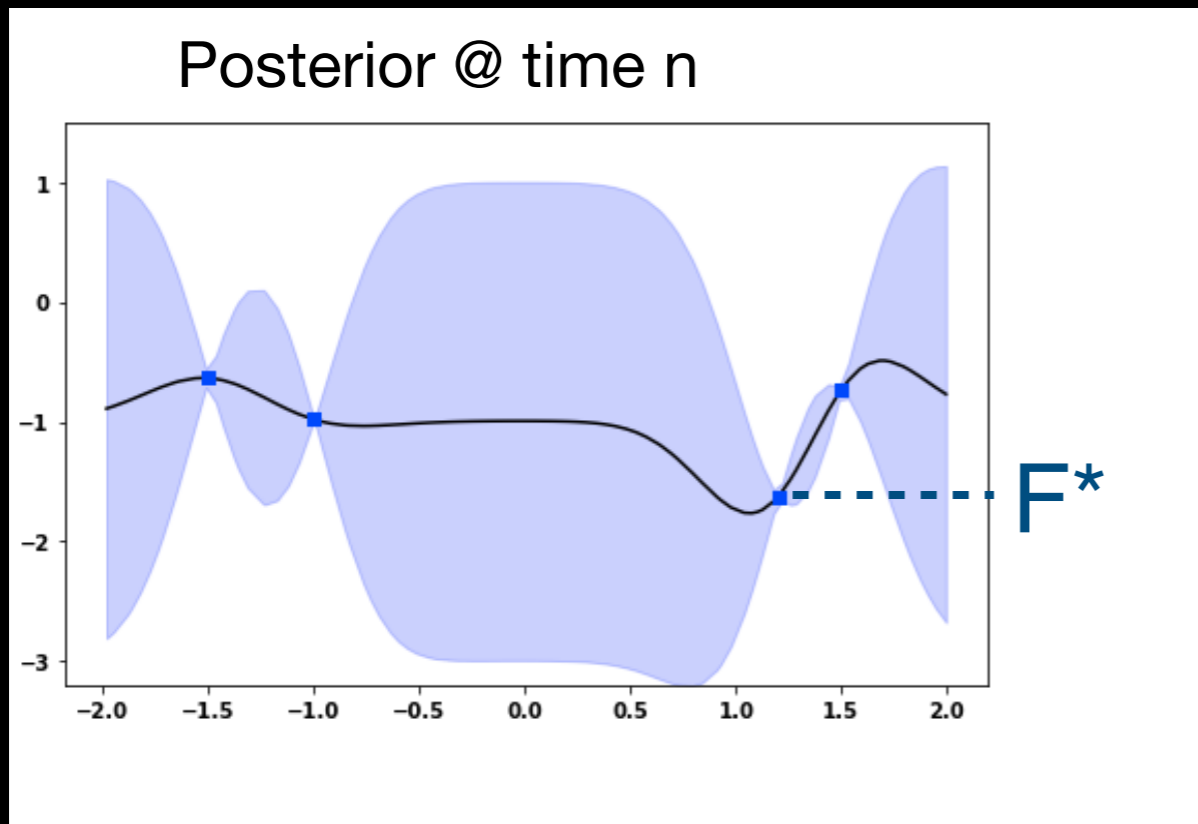
For heteroscedastic noise:

Kersting et al. "Most likely heteroscedastic gaussian process regression" ICML 2007

Agenda

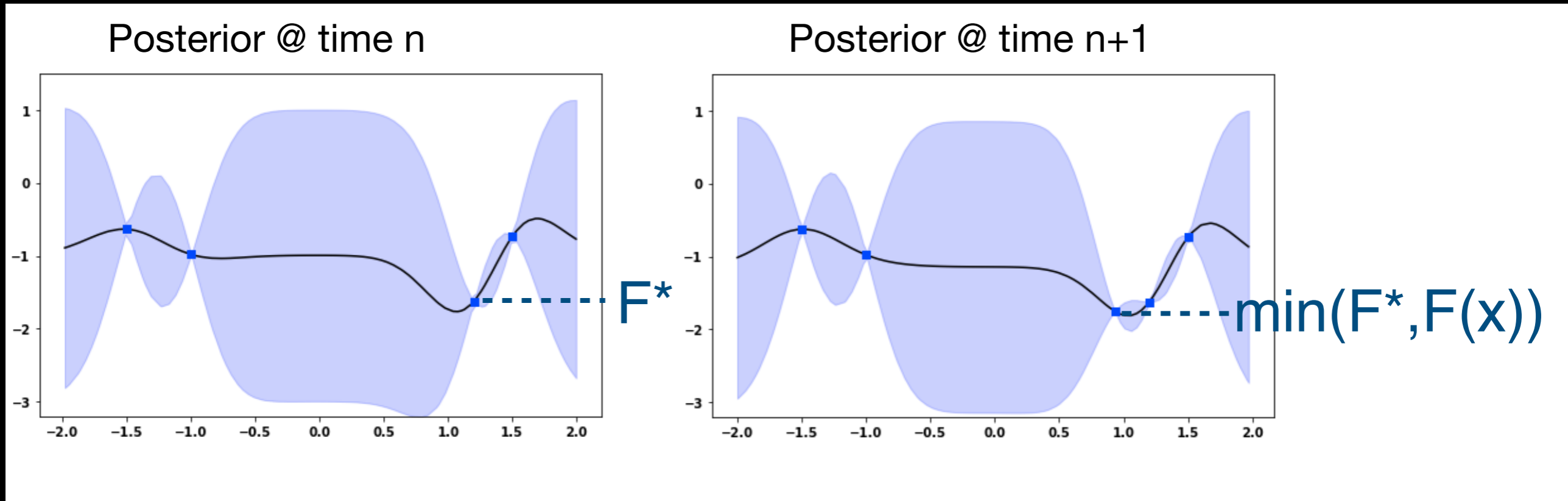
- Gaussian process regression
- **Expected improvement (EI)**
- Parallel BayesOpt (using EI)
- Parallel BayesOpt with gradients and/or noise (using KG)
- Code & Research directions

This is the **Expected Improvement (EI)** acquisition function
[Mockus 1989; Jones, Schonlau & Welch 1998]



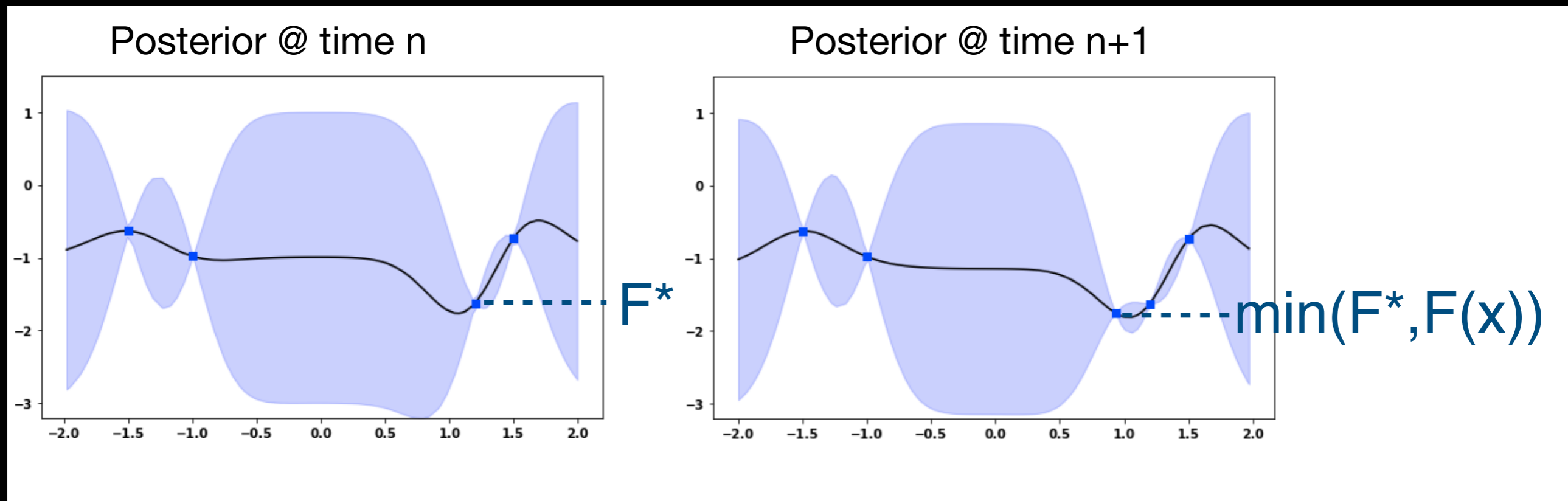
- Loss if we stop now:
 F^*

This is the **Expected Improvement (EI)** acquisition function [Mockus 1989; Jones, Schonlau & Welch 1998]



- Loss if we stop now:
 F^*
- Loss if we stop after sampling $F(x)$:
 $\min(F^*, F(x))$

This is the **Expected Improvement (EI)** acquisition function [Mockus 1989; Jones, Schonlau & Welch 1998]



- Loss if we stop now:
 F^*
- Loss if we stop after sampling $F(x)$:
 $\min(F^*, F(x))$
- Reduction in loss due to sampling:
 $E_n[F^* - \min(F^*, F(x))] = E_n[(F^* - F(x))^+] = EI(x)$

Expected improvement is Bayes-optimal
(in the noise-free standard BO problem)
under some assumptions:

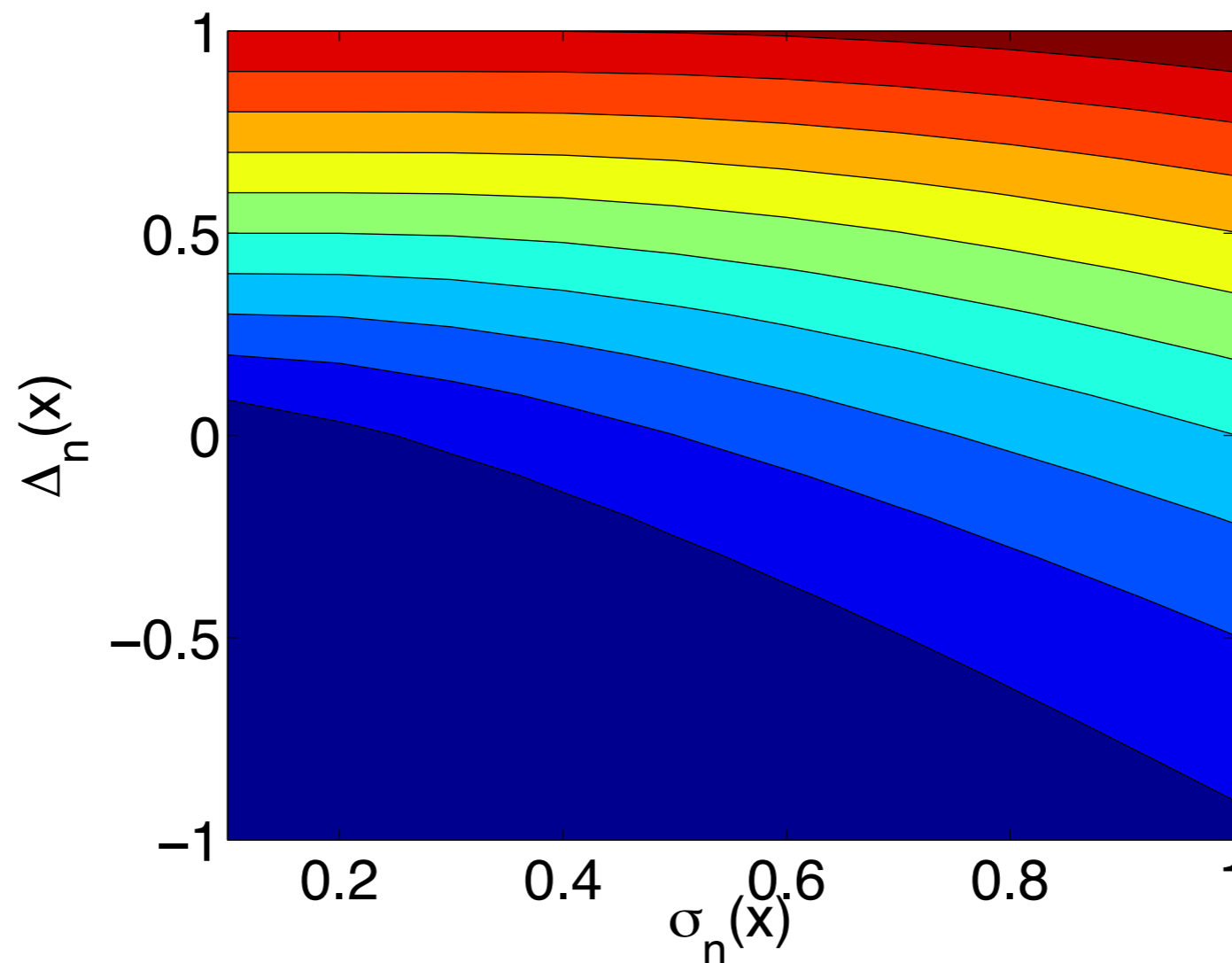
- this is our last evaluation
- we are risk neutral
- we are only willing to select a previously evaluated point as a final solution

You can compute expected improvement in closed form

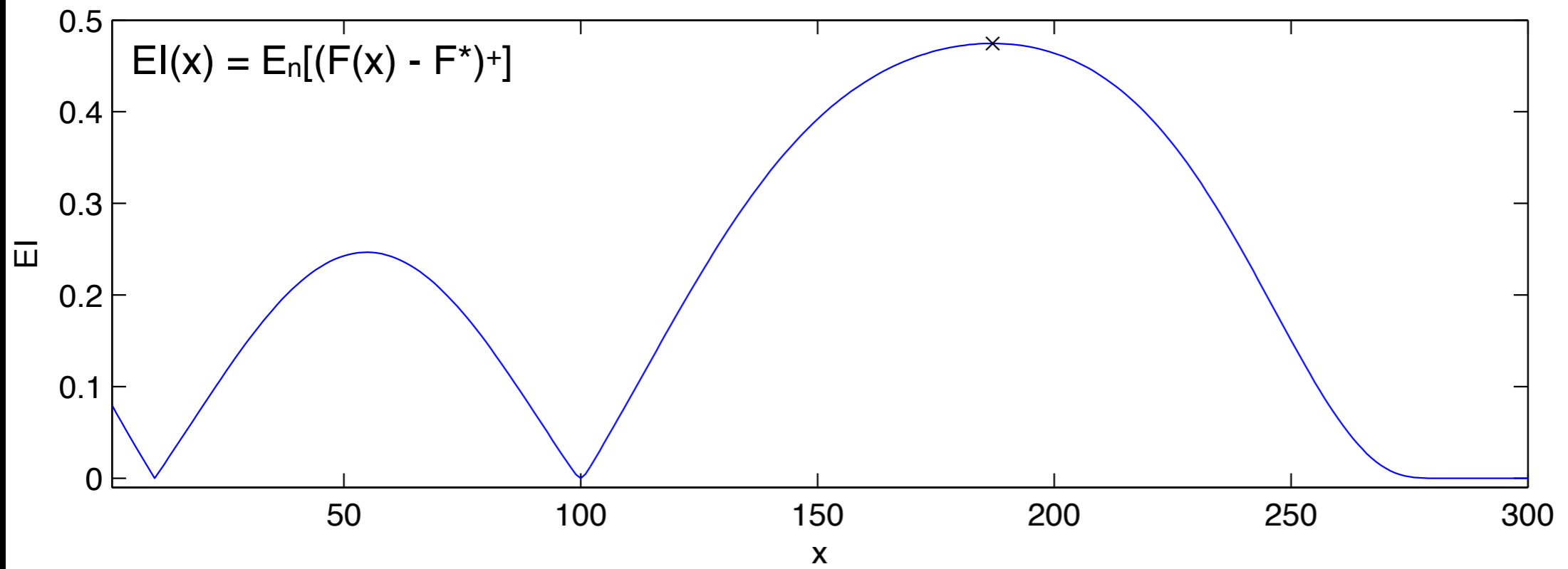
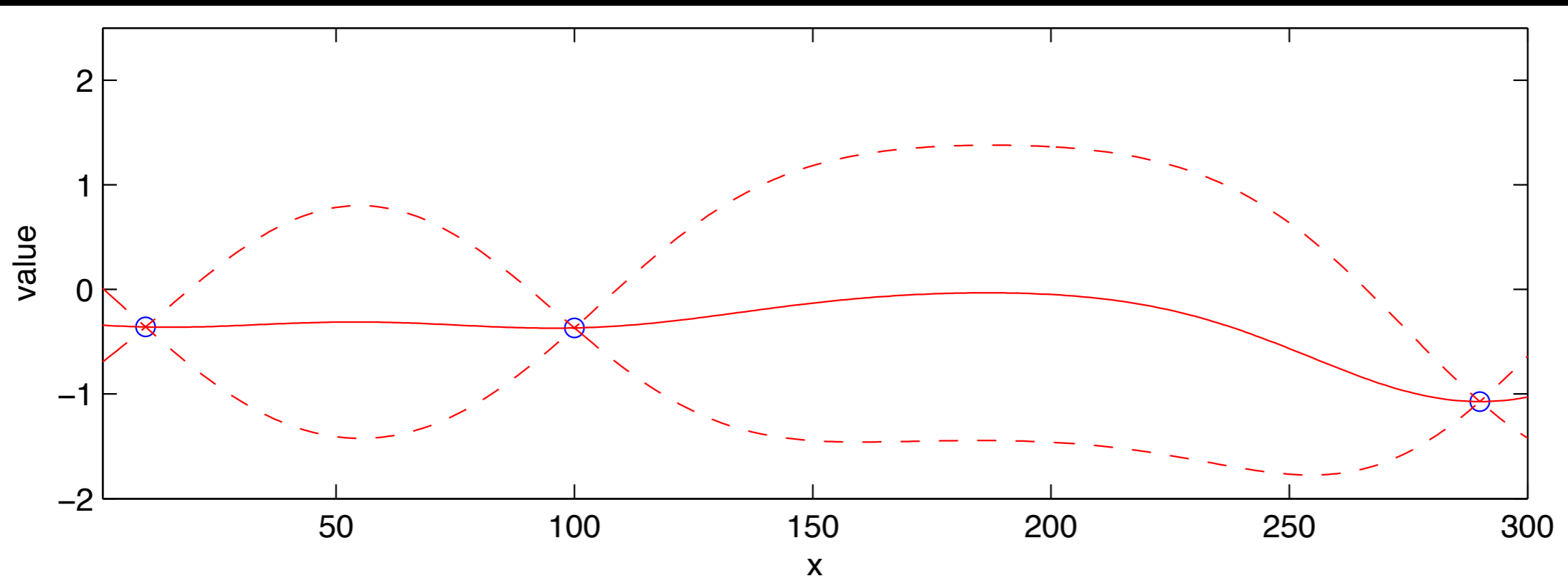
$$\mathbf{EI}_n(\mathbf{x}) = [\Delta_n(\mathbf{x})]^+ + \sigma_n(\mathbf{x})\varphi\left(\frac{\Delta_n(\mathbf{x})}{\sigma_n(\mathbf{x})}\right) - |\Delta_n(\mathbf{x})|\Phi\left(-\frac{|\Delta_n(\mathbf{x})|}{\sigma_n(\mathbf{x})}\right)$$

where $\Delta_n(\mathbf{x}) = F_n^* - \mu_n(\mathbf{x})$

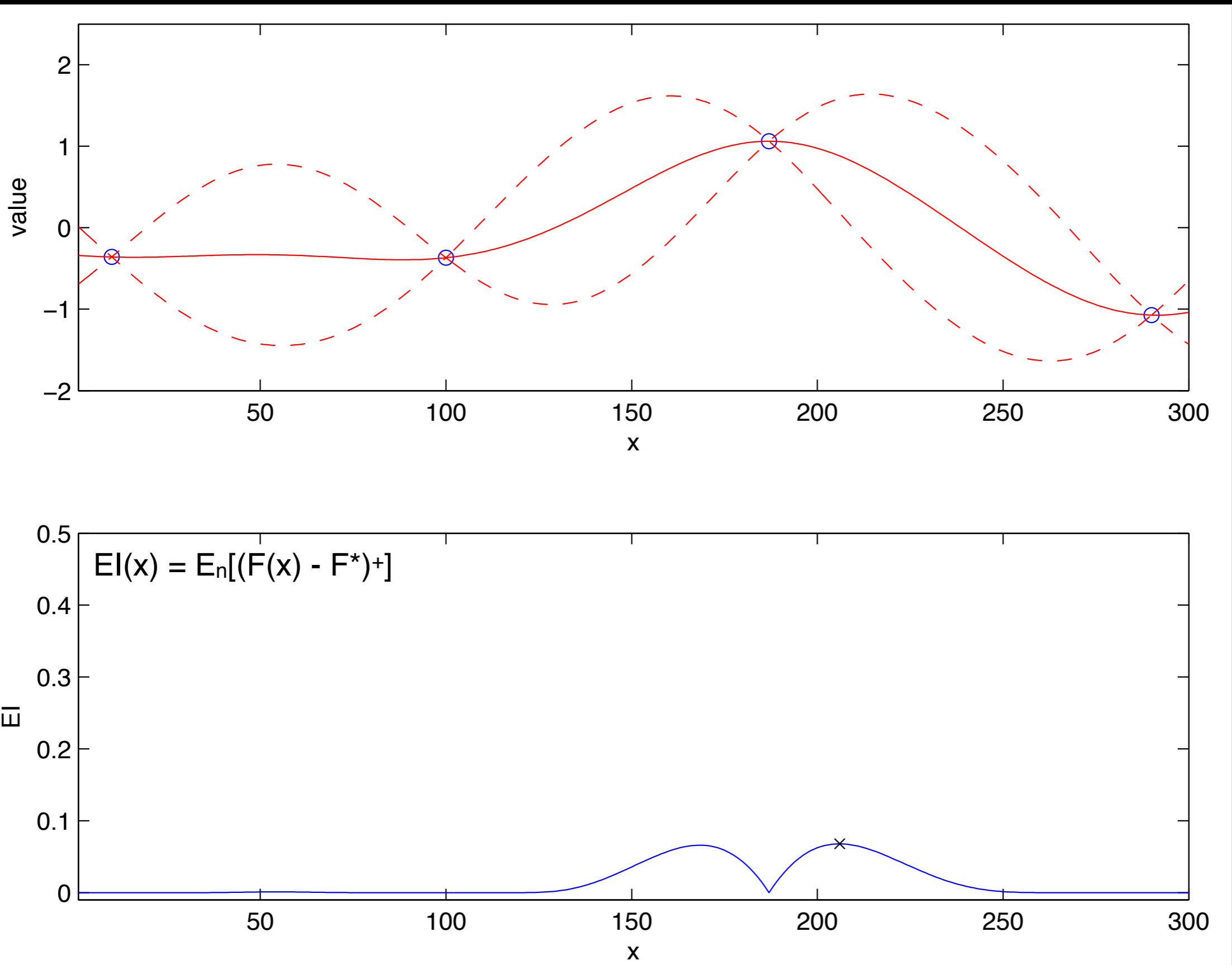
Expected improvement trades exploration ($\sigma_n(x)$) vs. exploitation ($\Delta_n(x)$)



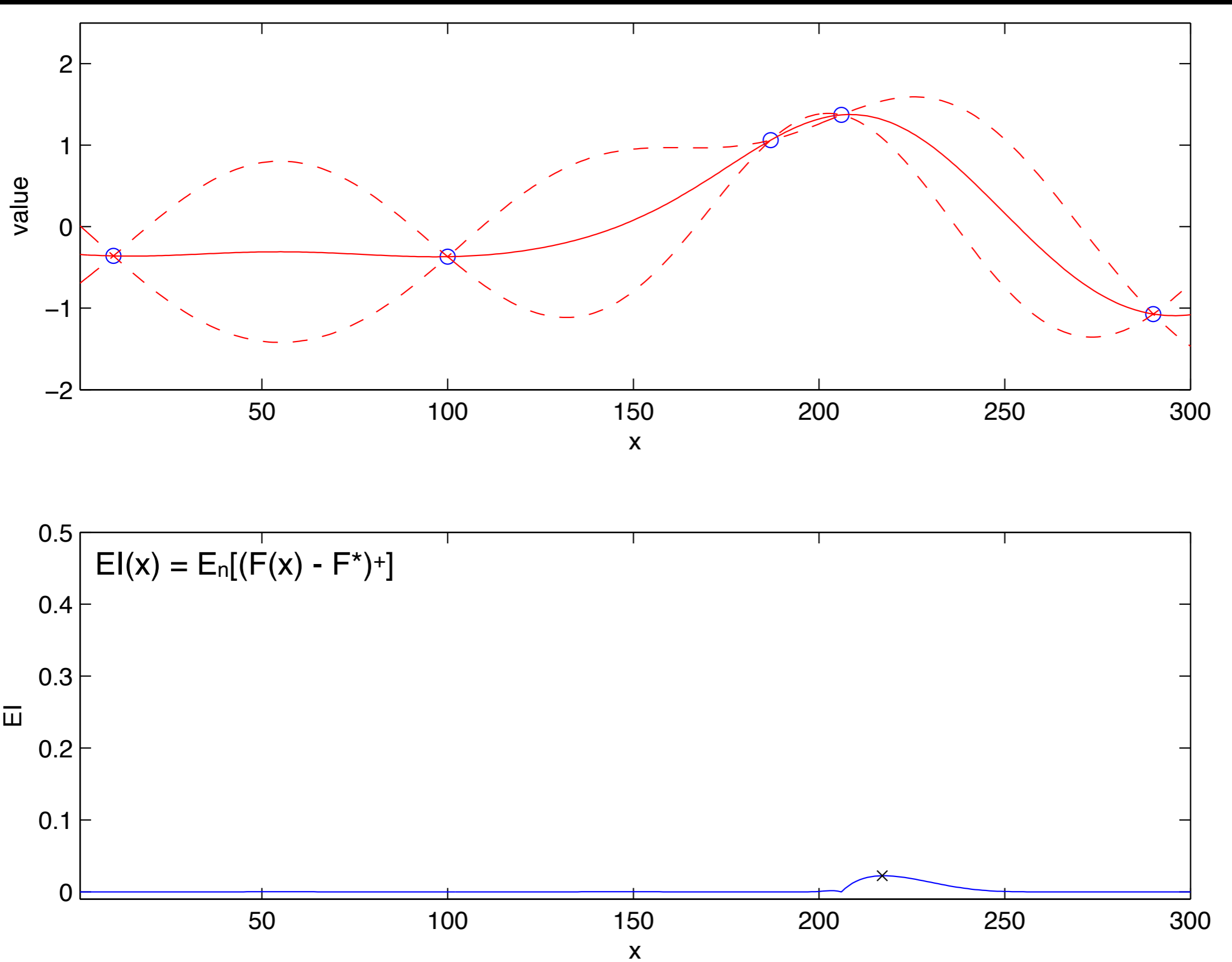
Here is expected improvement maximizing a 1-dim objective



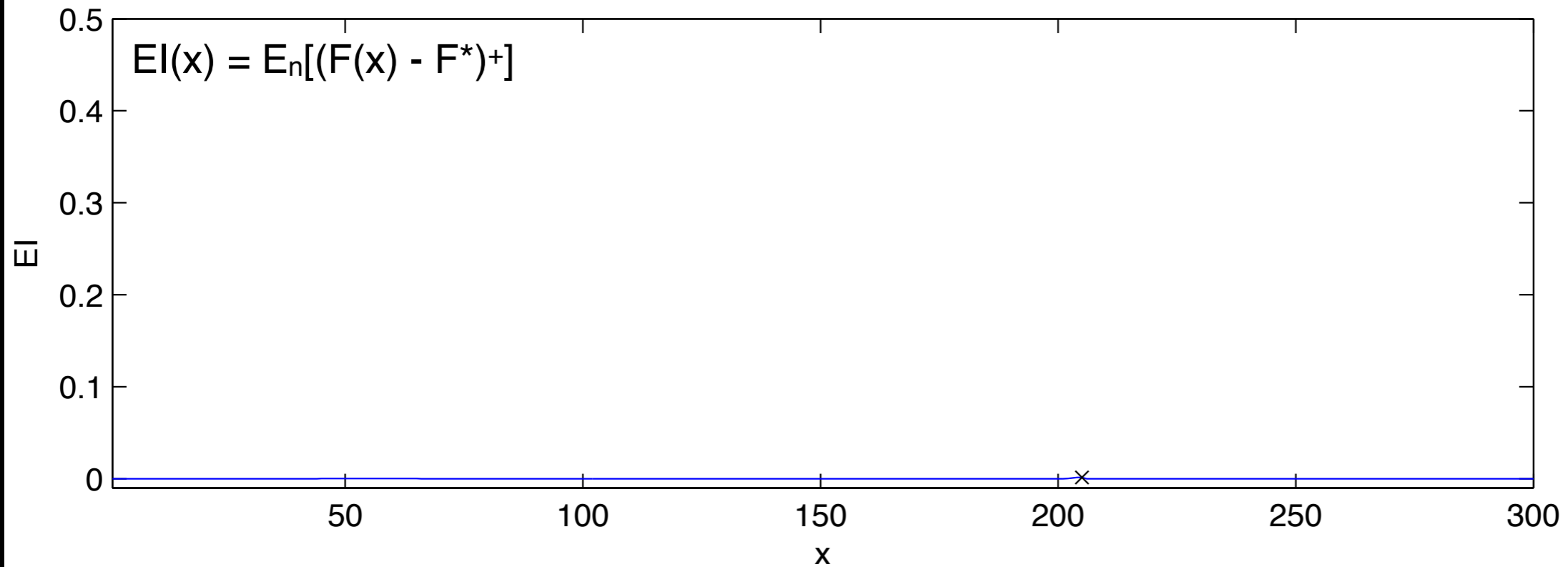
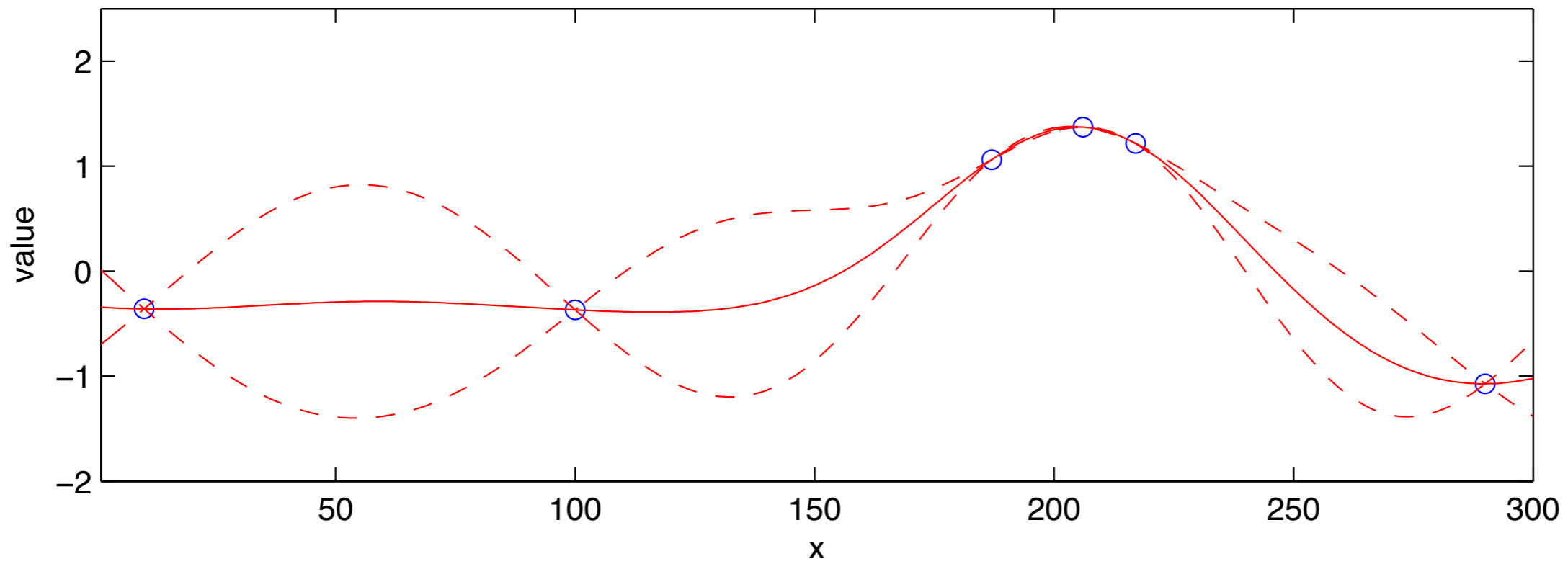
Here is expected improvement maximizing a 1-dim objective



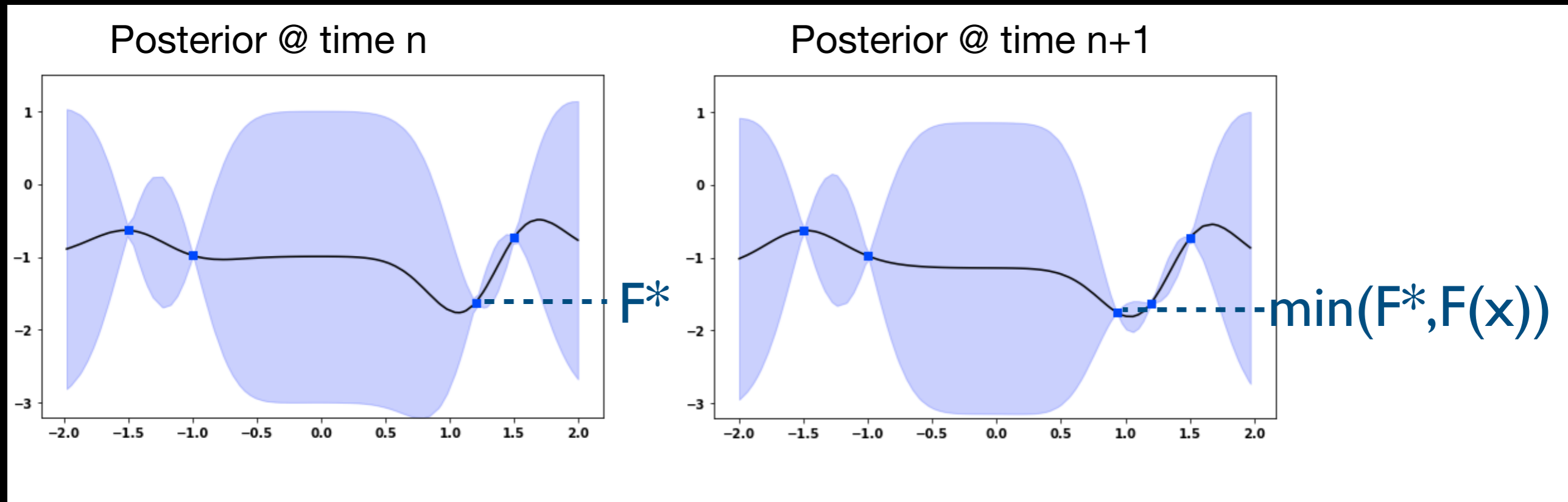
Here is expected improvement maximizing a 1-dim objective



Here is expected improvement maximizing a 1-dim objective



Expected Improvement is one-step Bases-Optimal under assumptions



Expected improvement is Bayes-Optimal
(in the noise-free standard BayesOpt problem) if:

- this is our last evaluation
- we are risk neutral
- we are only willing to select a previously evaluated point as a final solution

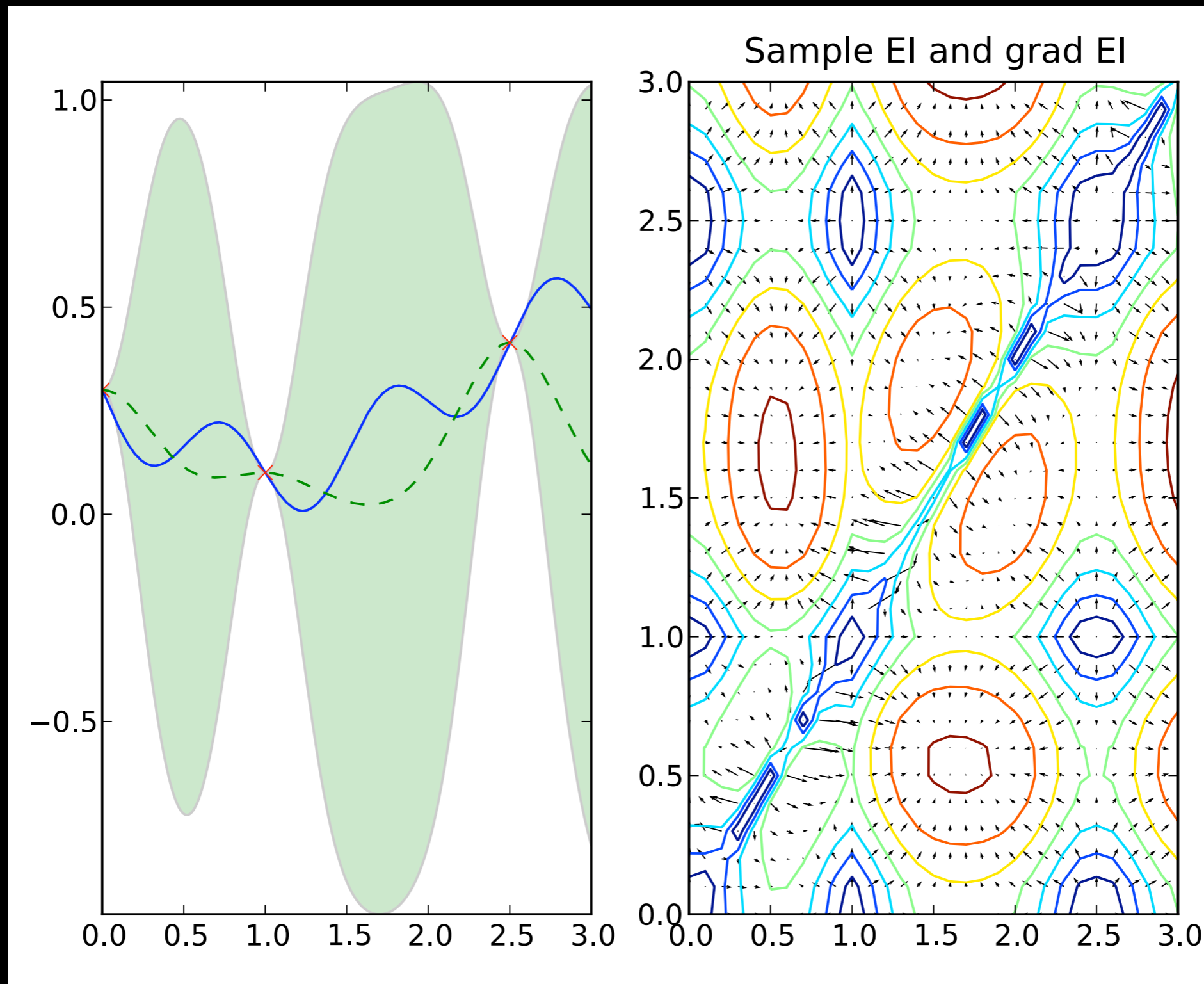
Agenda

- Gaussian process regression
- Expected improvement (EI)
- **Parallel BayesOpt (using EI)**
- Parallel BayesOpt with gradients and/or noise (using KG)
- Code & Research directions

We can parallelize EI

$$EI(x_{1:q}) = E_n[(F^* - \min\{F(x_1), \dots, F(x_q)\})^+]$$

We can parallelize EI



$$EI(x_{1:q}) = E_n[(F^* - \min\{F(x_1), \dots, F(x_q)\})^+]$$

Here's how to maximize parallel EI

1. Construct an unbiased estimator of the gradient of $EI(x_{1:q})$ using **infinitesimal perturbation analysis (IPA)**.
2. Use **multistart stochastic gradient ascent** to approximately maximize $EI(x_{1:q})$

Here's how we estimate ∇EI

- $Y=[f(x_1),\dots,f(x_q)]'$ is multivariate normal
- Let $m = E[Y]$ and $C = \text{Chol}(\text{Cov}[Y])$
- $Y=m+CZ$, where Z is a vector of independent standard normals
- $EI(x_{1:q}) = E[h(Y)]$ where $h(Y) = (F^* - \min\{Y\})^+$
- If our problem is well-behaved,
we can switch derivative and expectation:
$$\nabla EI(x_{1:q}) = E[\nabla h(m+CZ)]$$

Here's how we estimate ∇EI

1. Simulate a vector Z of independent standard normals
2. Calculate $m = E[Y]$, $C = \text{Chol}(\text{Cov}[Y])$
3. Our estimator of $\nabla EI(x_1, \dots, x_q)$ is $\nabla h(m + CZ)$
where $h(Y) = (F^* - \min\{Y\})^+$

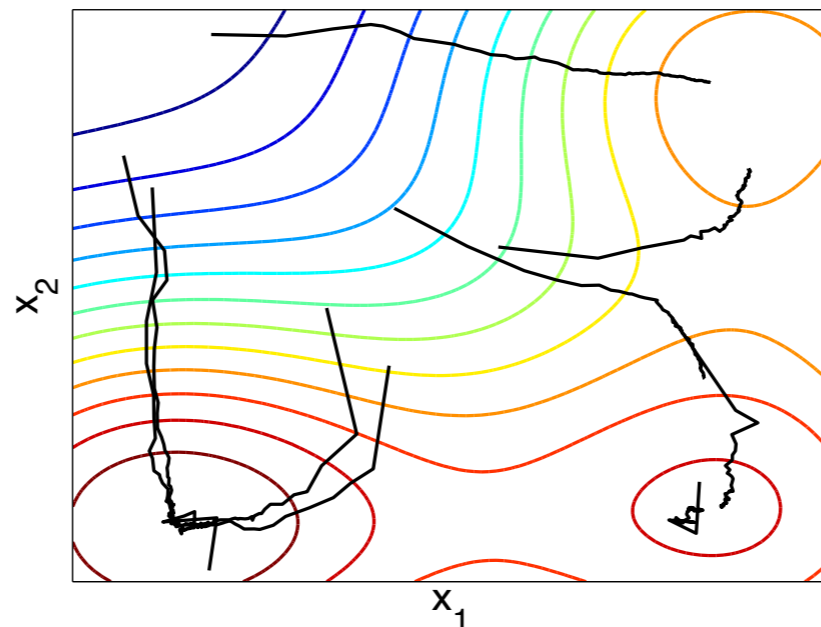
We use this estimator of ∇E in multistart stochastic gradient ascent

- 1 Select several starting points, uniformly at random.
- 2 From each starting point, iterate using the stochastic gradient method until convergence.

$$(\vec{x}_1, \dots, \vec{x}_q) \leftarrow (\vec{x}_1, \dots, \vec{x}_q) + \alpha_n g(\vec{x}_1, \dots, \vec{x}_q, \omega),$$

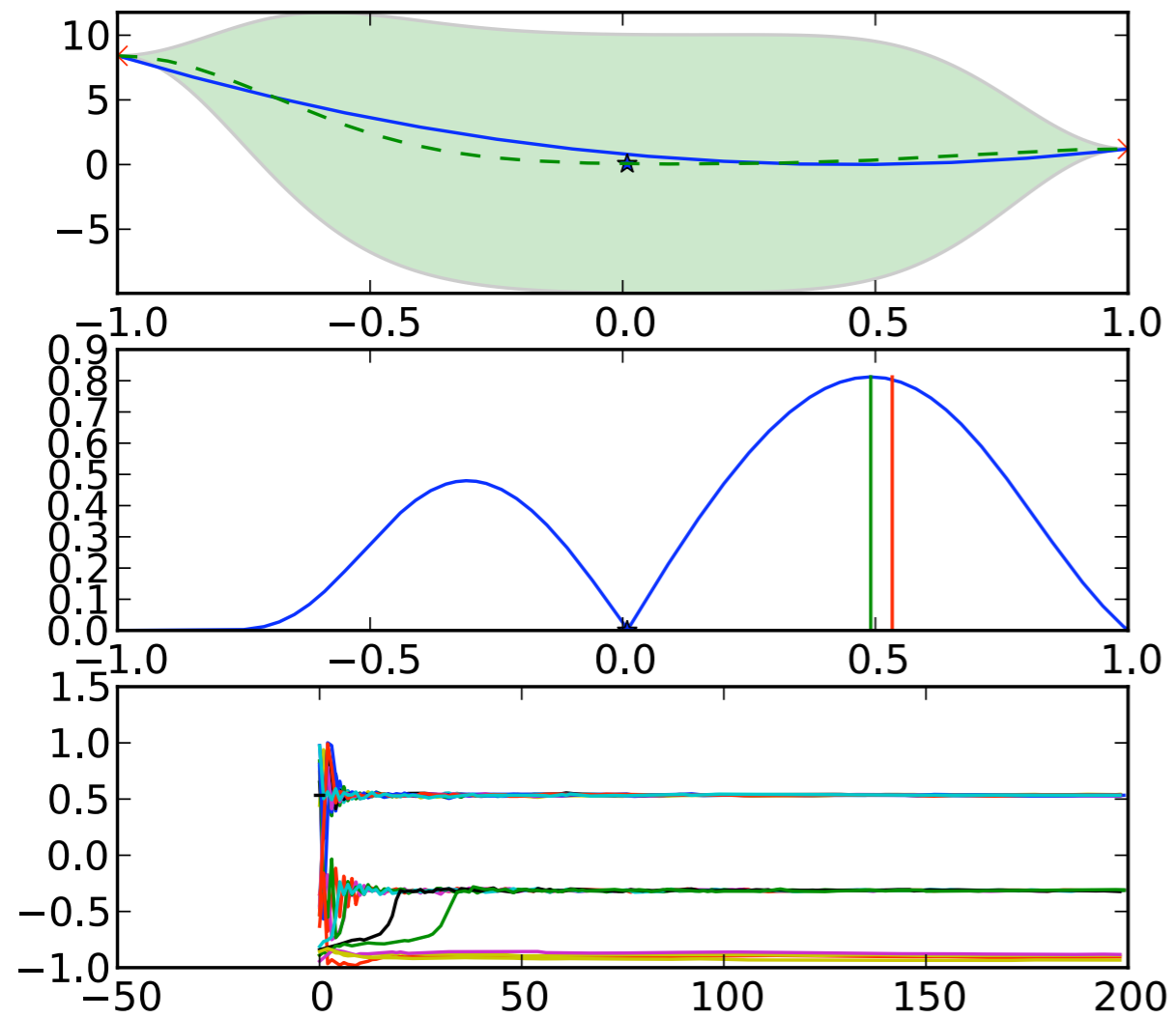
where (α_n) is a stepsize sequence.

- 3 For each starting point, average the iterates to get an estimated stationary point. (Polyak-Ruppert averaging)
- 4 Select the estimated stationary point with the best estimated value as the solution.

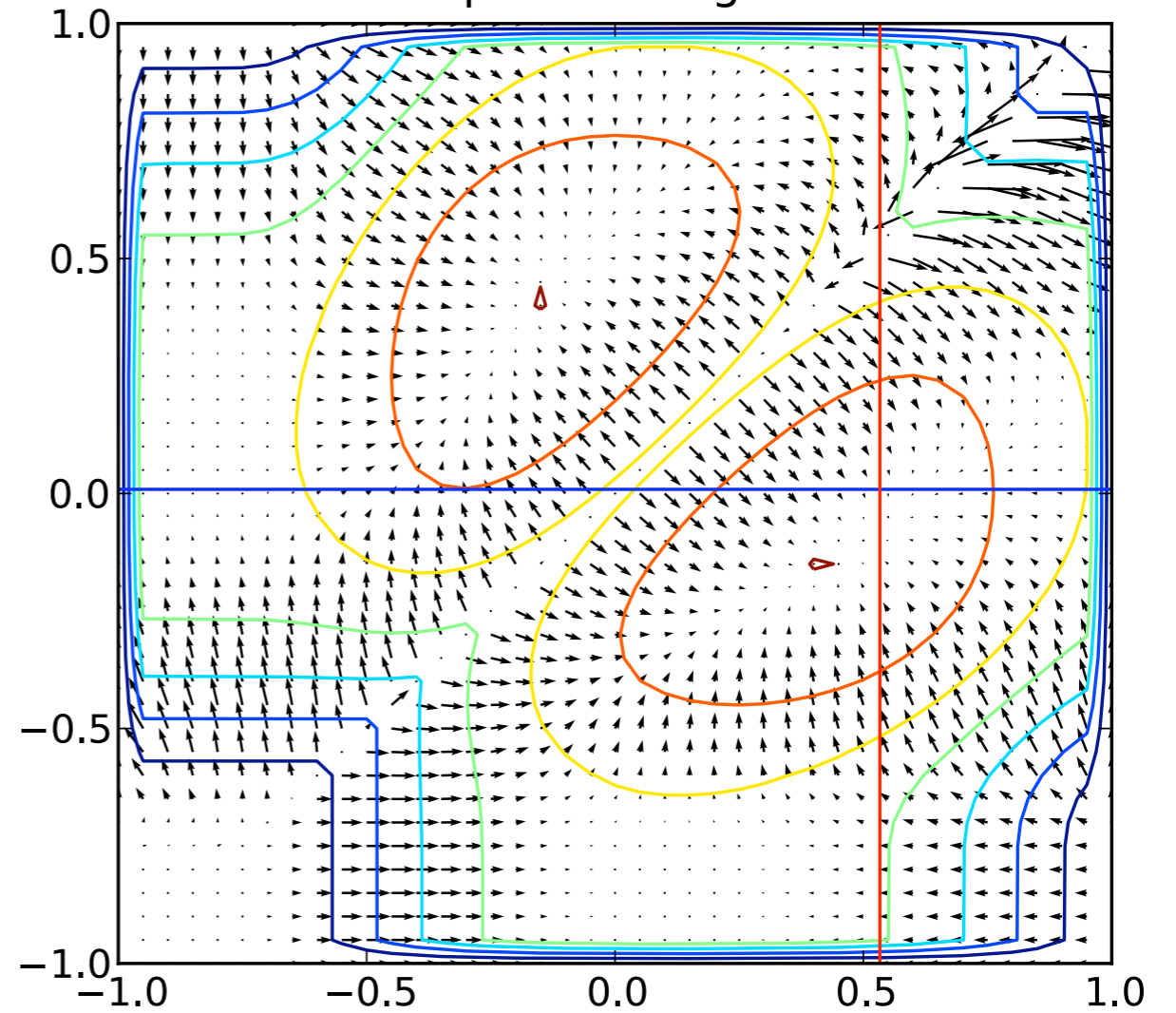


Animation

GPP of points sampled

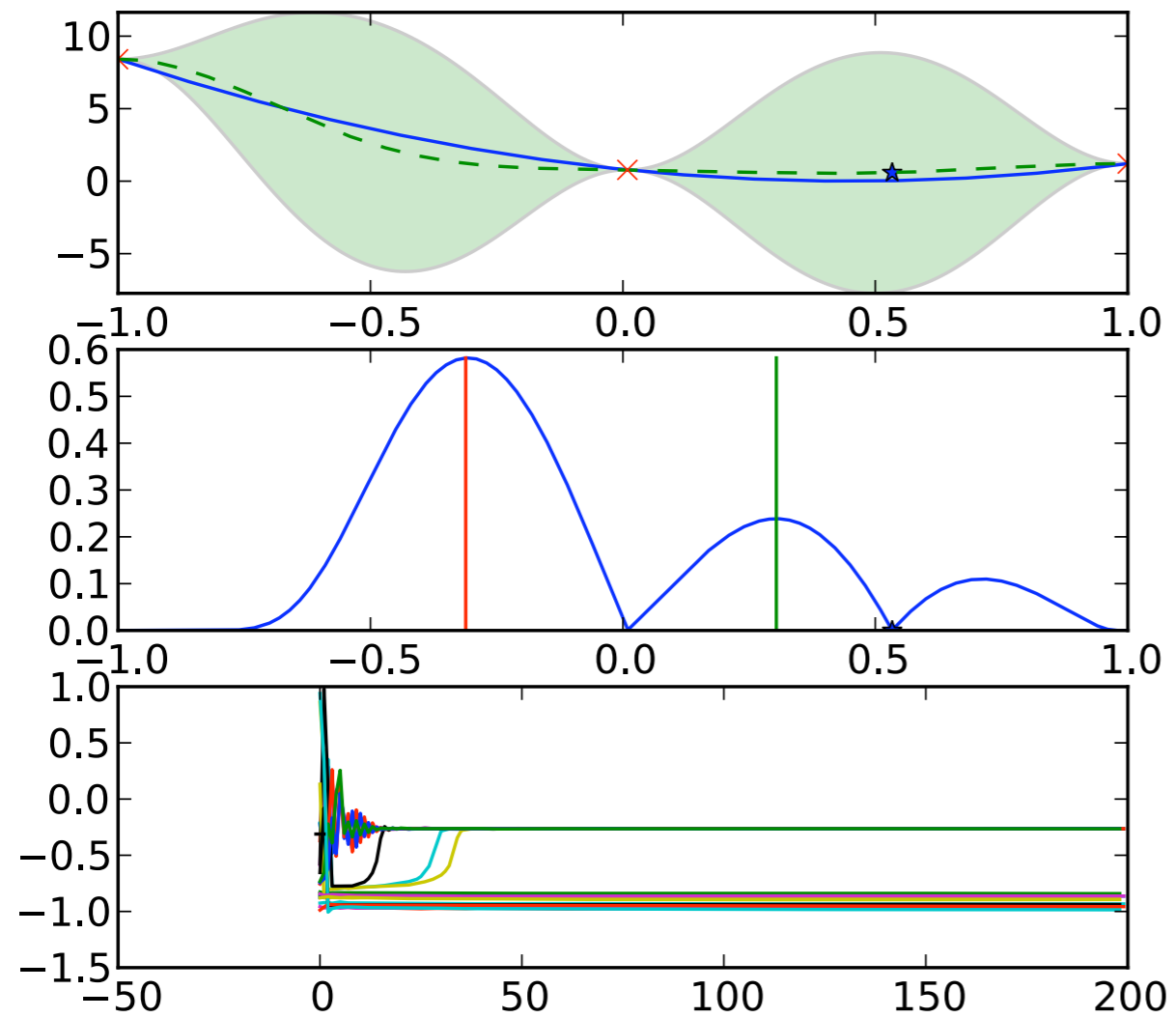


Sample EI and grad EI

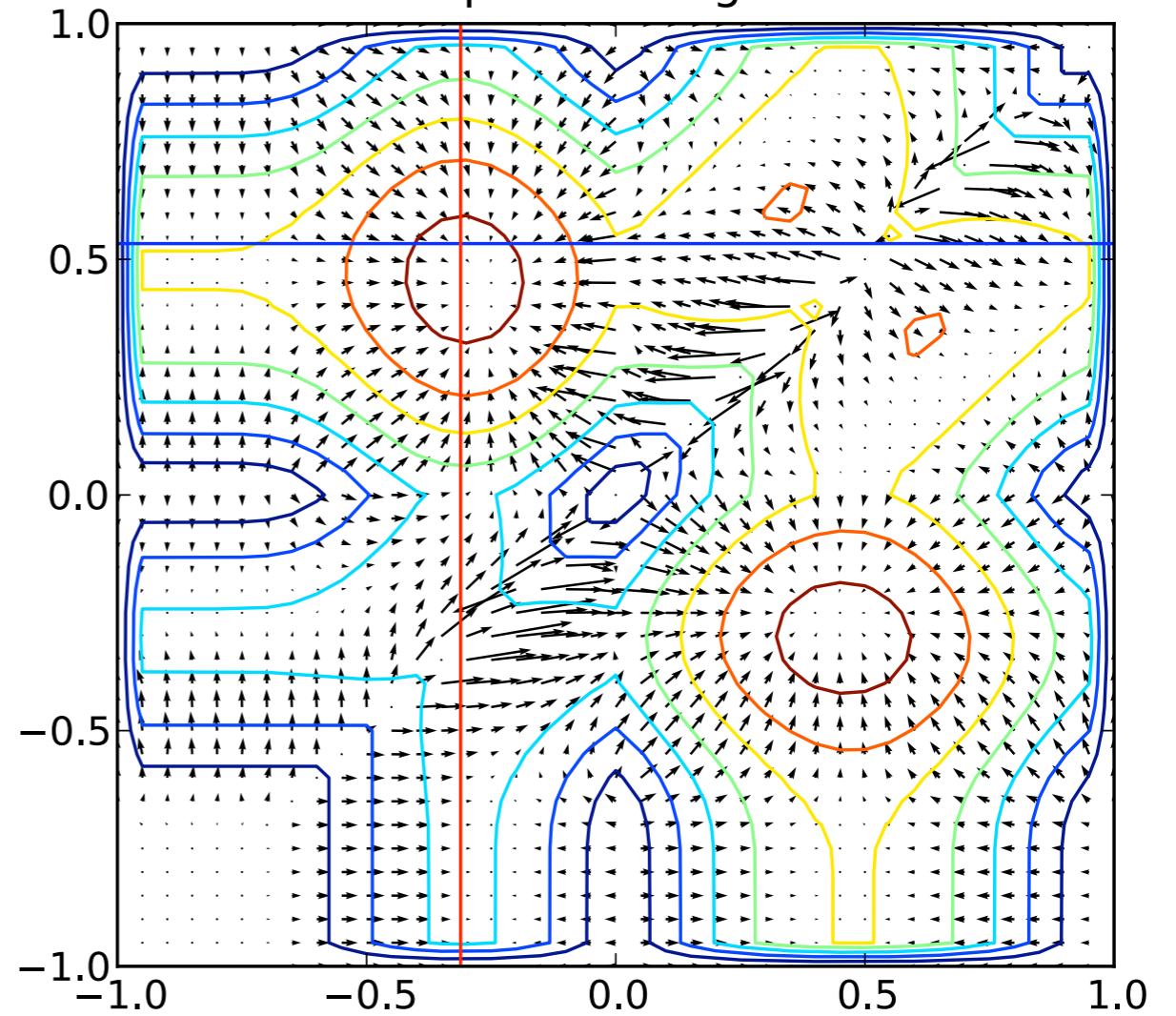


Animation

GPP of points sampled

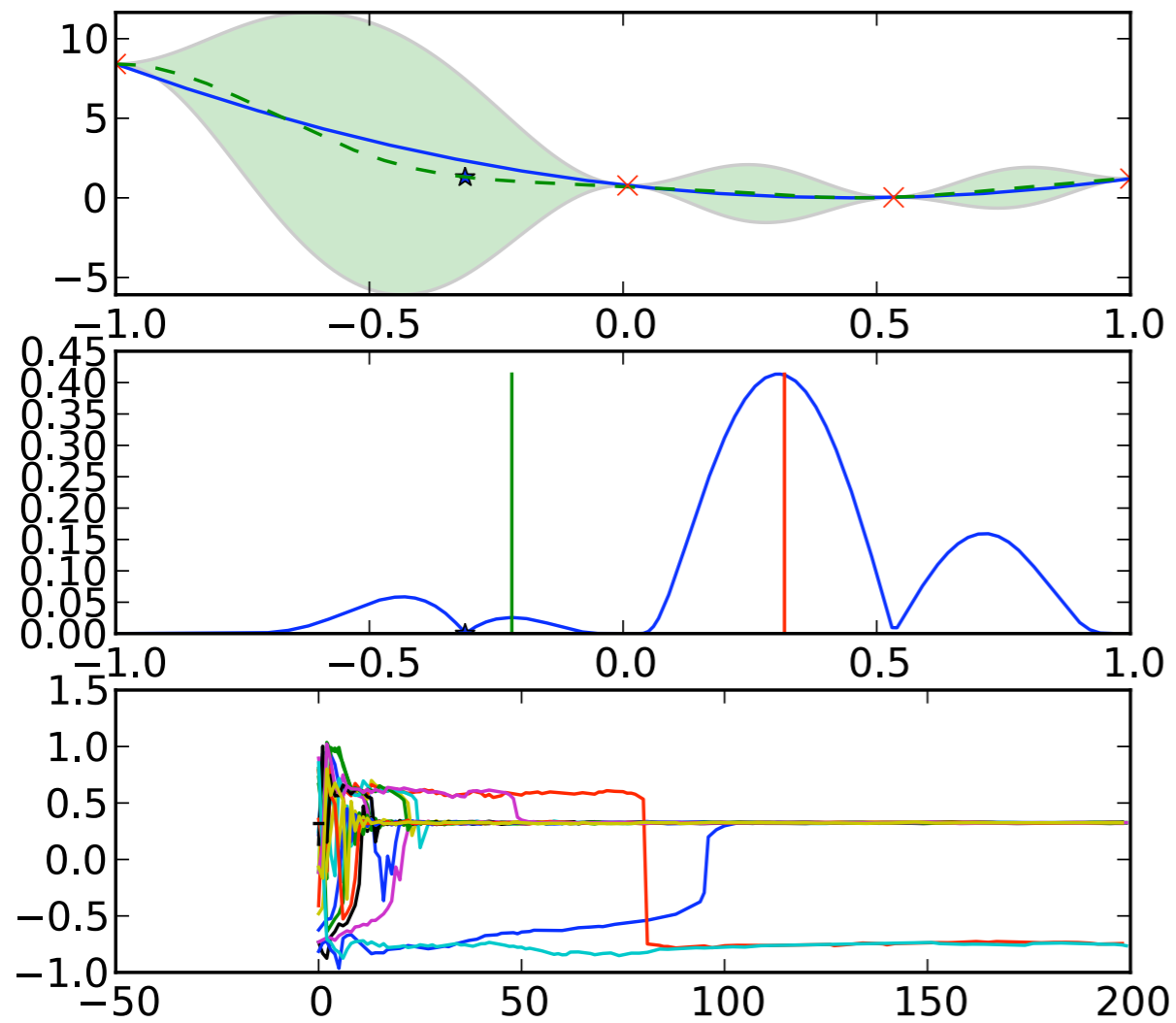


Sample EI and grad EI

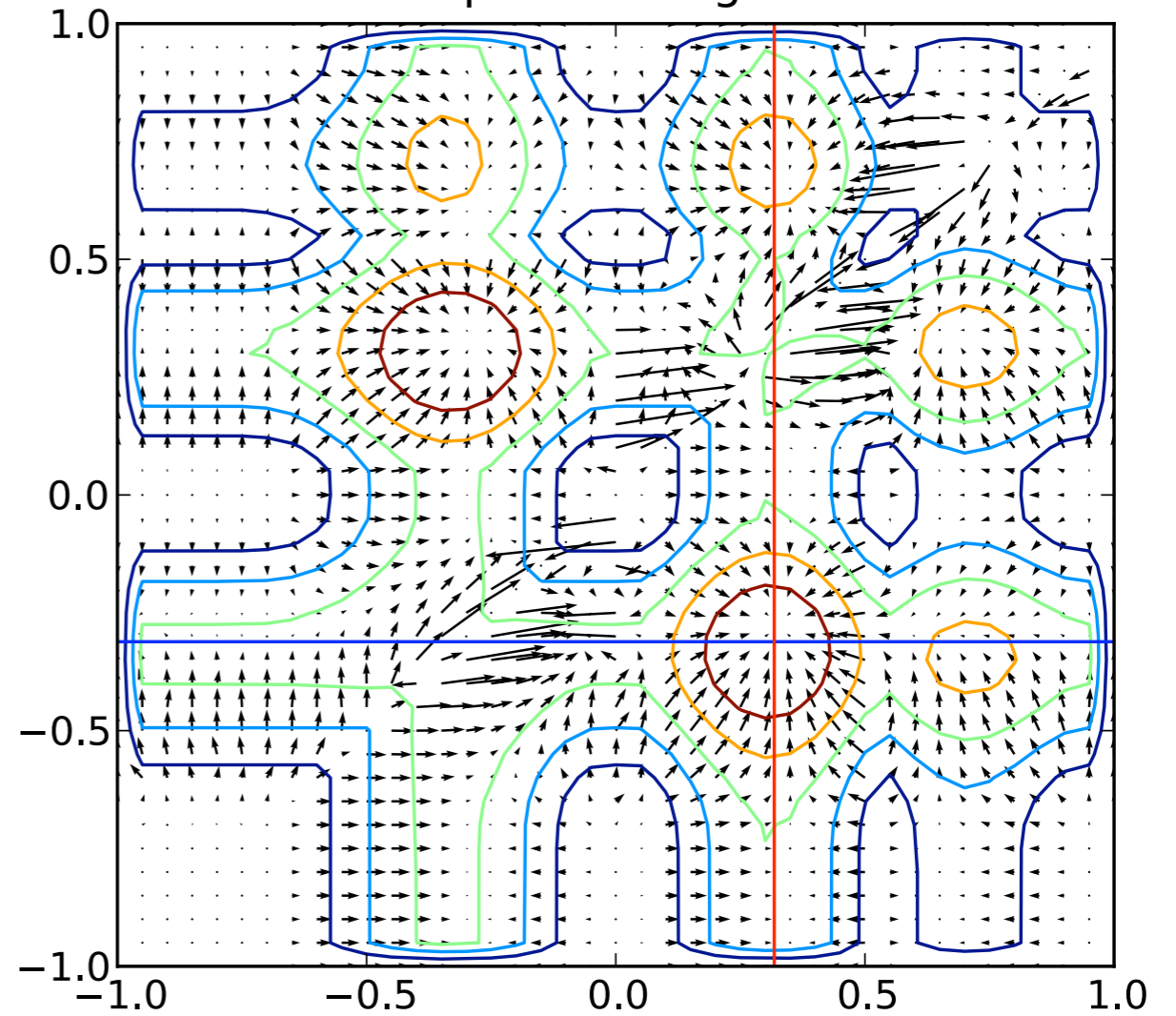


Animation

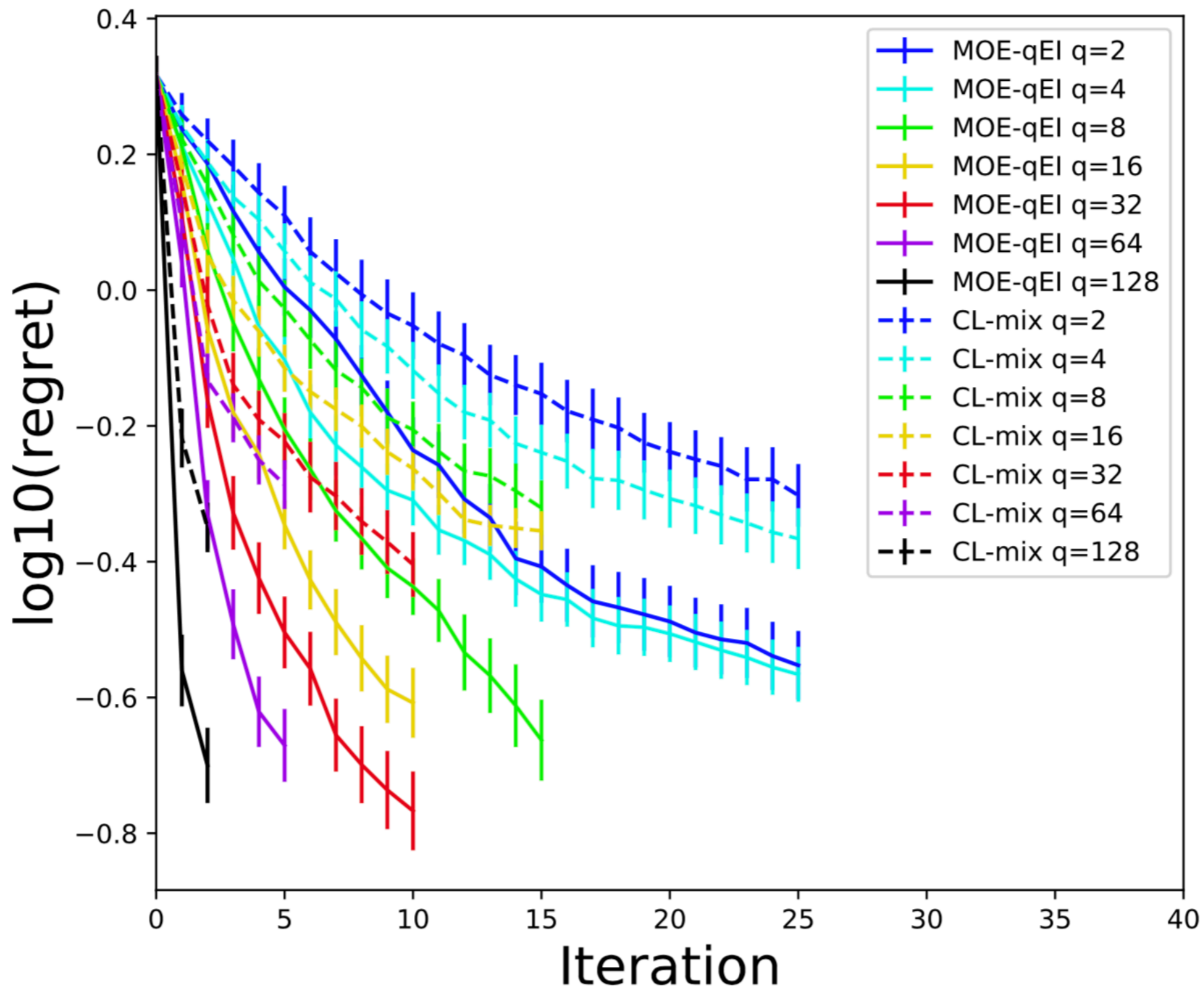
GPP of points sampled



Sample EI and grad EI



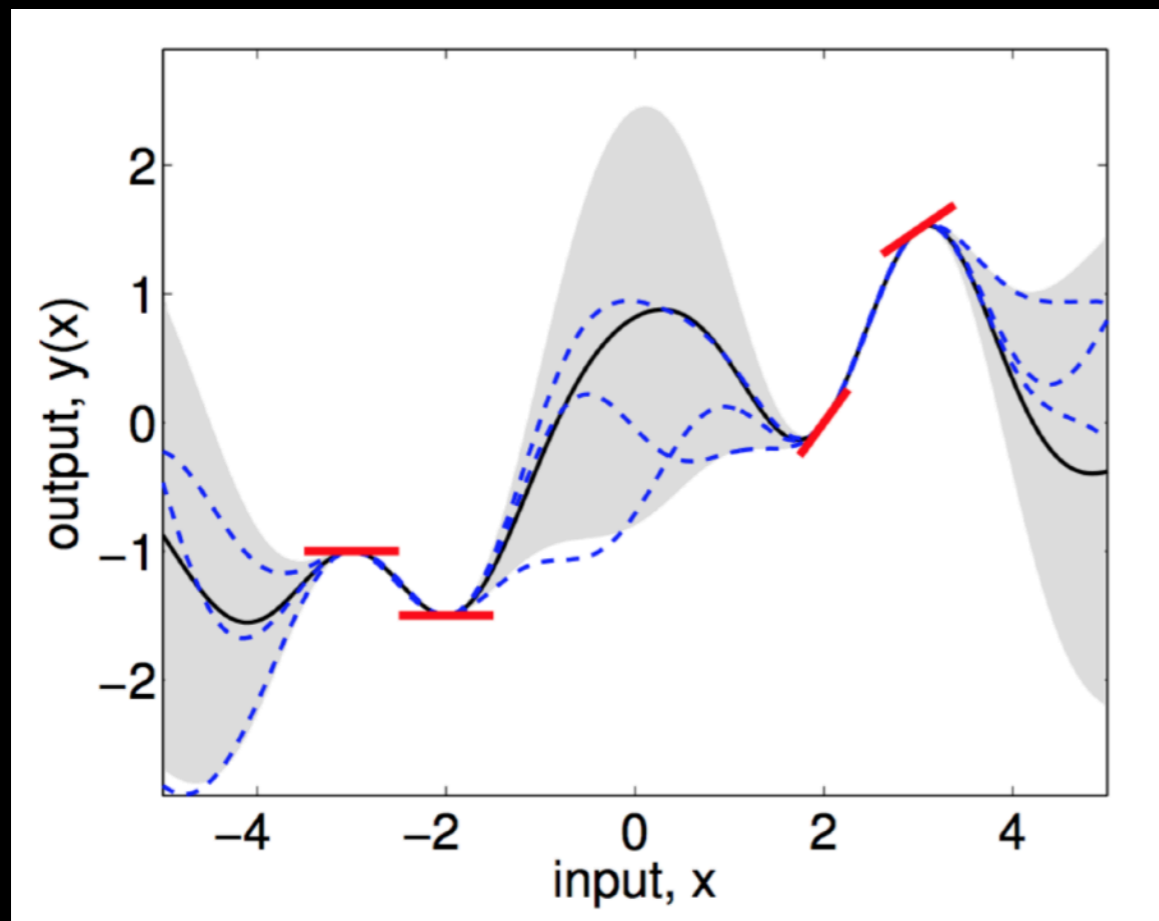
Here's a demonstration on a 6-dimensional BayesOpt problem with up to 128 parallel evaluations



Agenda

- Gaussian process regression
- Expected improvement (EI)
- Parallel BayesOpt (using EI)
- **Parallel BayesOpt with gradients/noise (using KG)**
- Code & Research directions

Bayesian (global) optimization with gradients



Evaluating at x provides

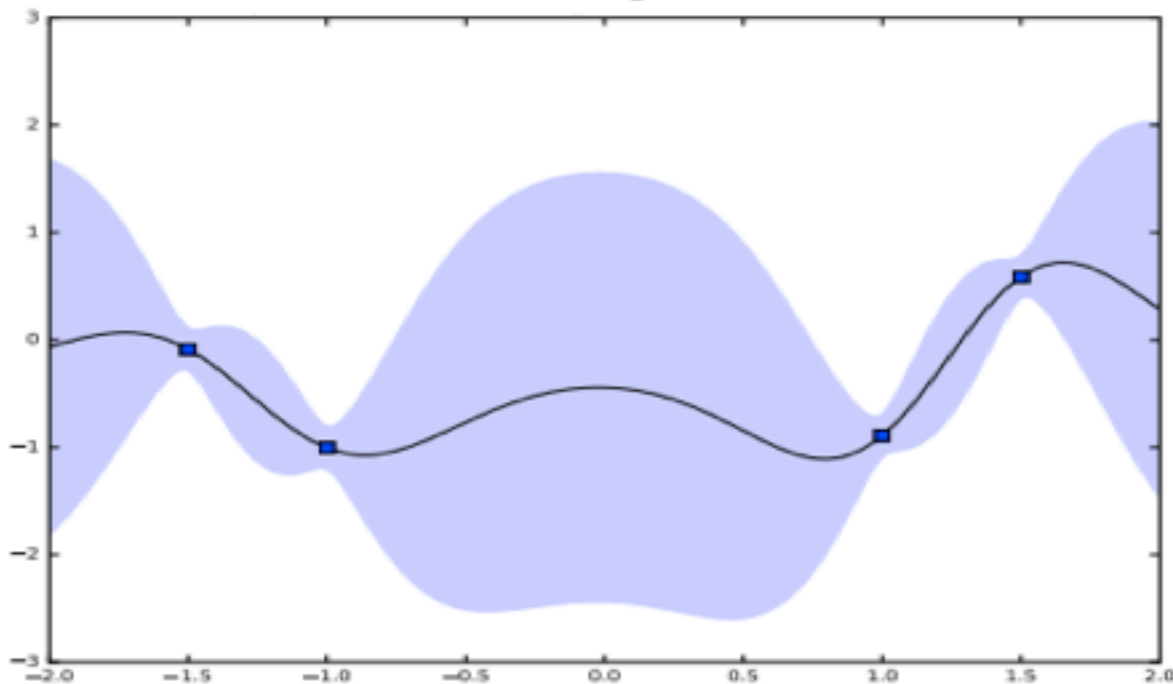
- $F(x)$, optionally with noise
- optionally incomplete/ noisy/ biased observations of $\nabla F(x)$

Evaluations can be sequential
or in **batches**

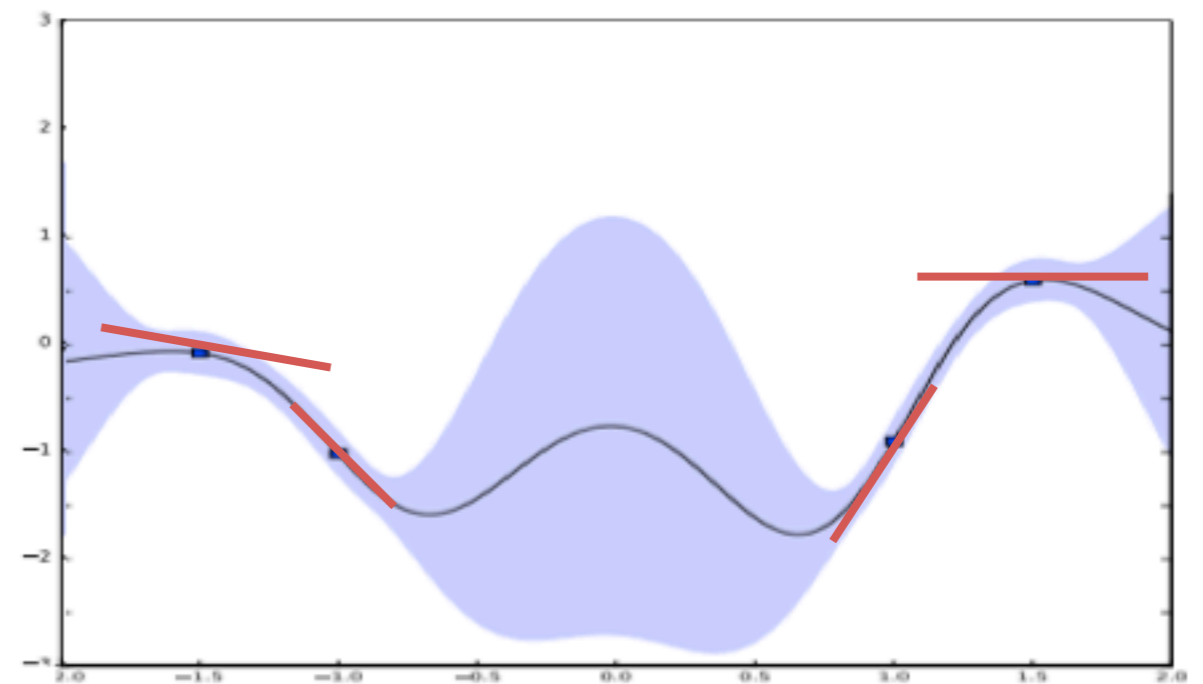
GP regression handles gradients naturally

[Rasmussen & Williams 2006]

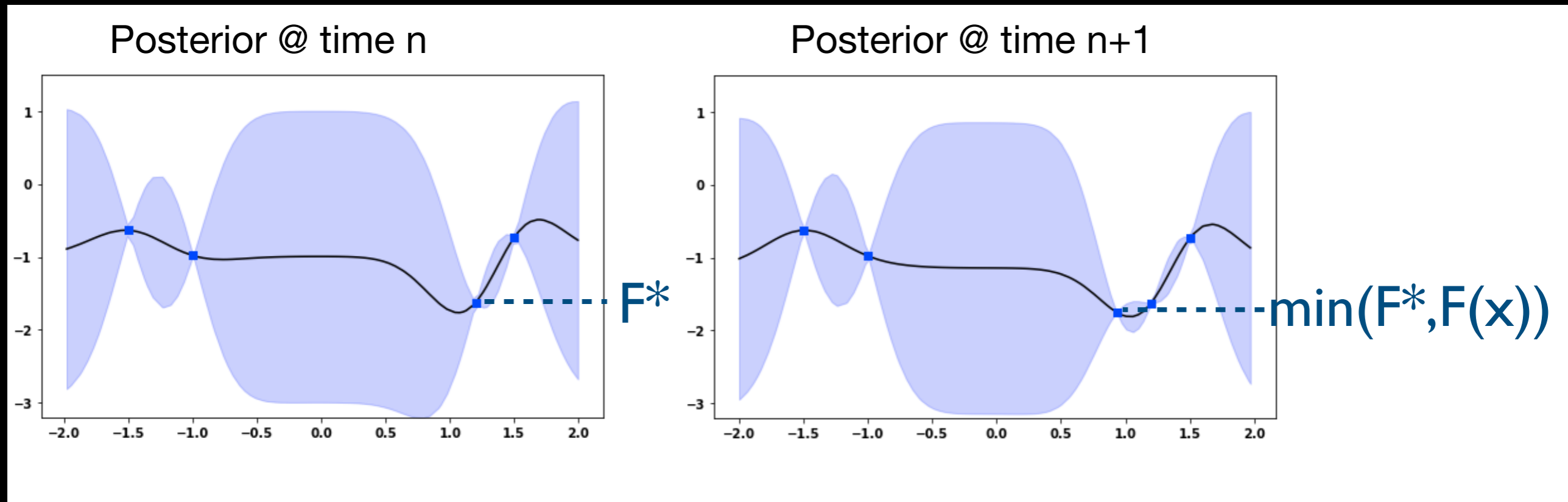
Posterior **without** gradient observations



Posterior **with** gradient observations

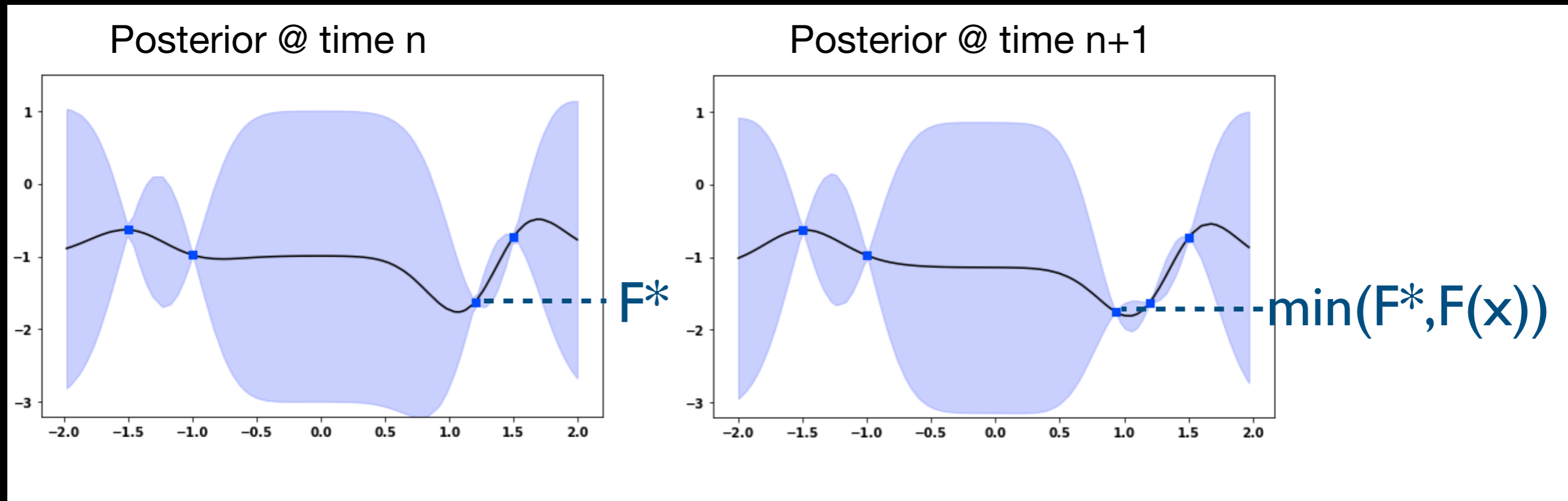


Here is EI in a standard BO problem

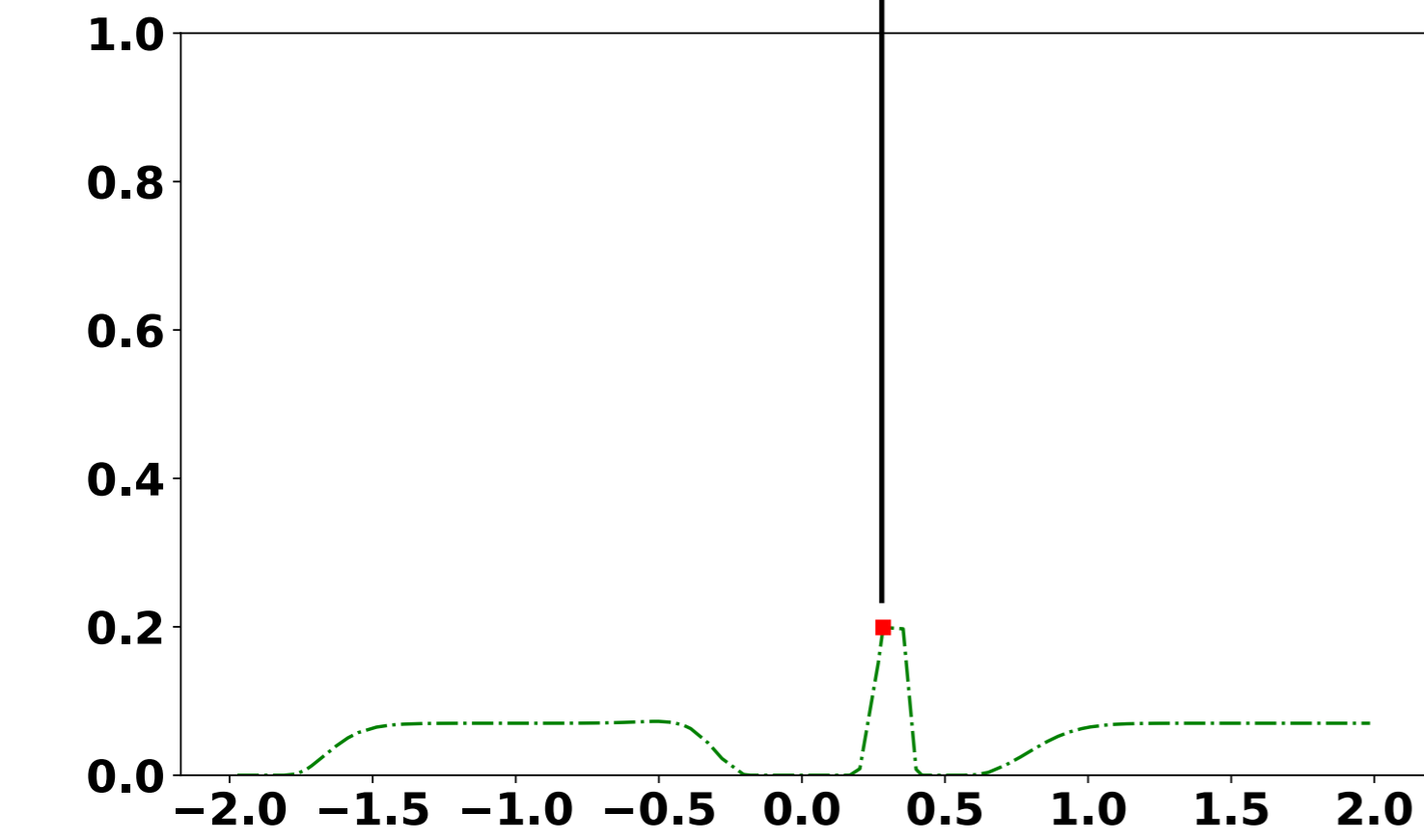
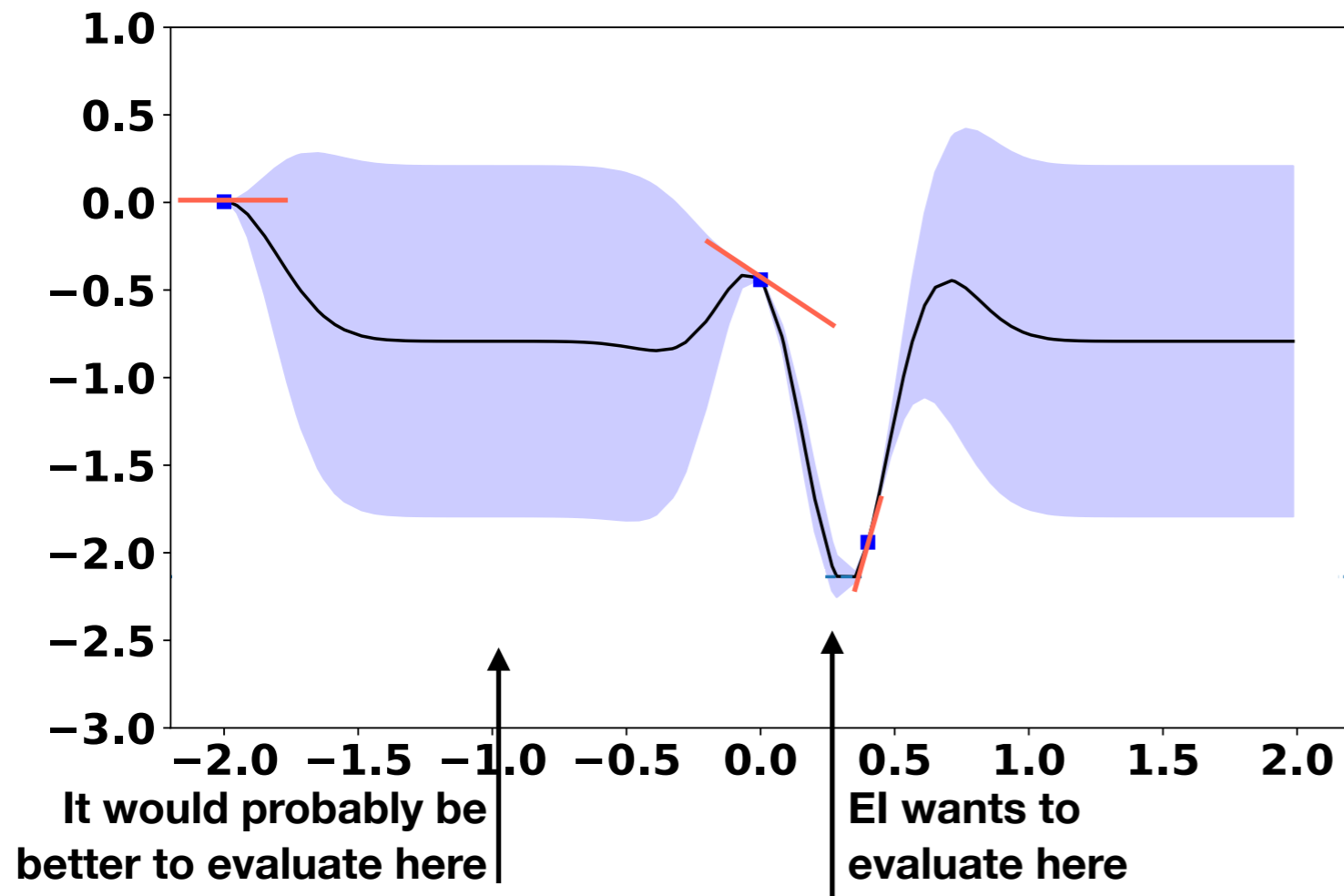


- Loss if we stop now:
 F^*
- Loss if we stop after sampling $F(x)$:
 $\min(F^*, F(x))$
- Reduction in loss due to sampling:
 $E_n[F^* - \min(F^*, F(x))] = E_n[(F^* - F(x))^+] = EI(x)$

El places **no value** on gradient information



- Loss if we stop now:
 F^*
- Loss if we stop after sampling $F(x)$ & its gradient:
 $\min(F^*, F(x))$
- Reduction in loss due to sampling:
 $E_n[F^* - \min(F^*, F(x))] = E_n[(F^* - F(x))^+] = \text{El}(x)$



El can make
poor decisions

Recall: Expected improvement is Bayes-Optimal
(in the noise-free standard BO problem) if:

- this is our last evaluation
- we are risk neutral
- we are only willing to select a previously evaluated point as a final solution

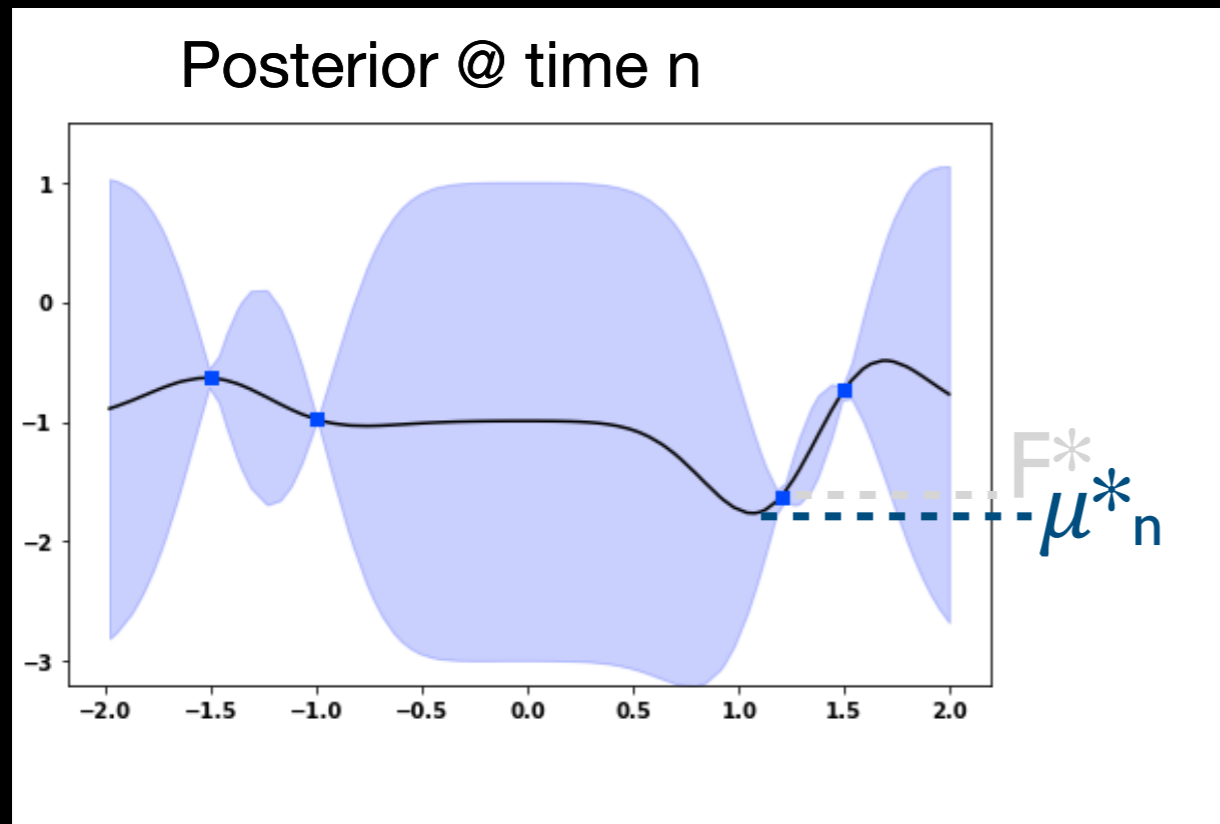
Let's **drop** one of these assumptions

- this is our last evaluation
- we are risk neutral
- ~~we are only willing to select a previously evaluated point as a final solution~~

This gives the Knowledge-Gradient acquisition function

- this is our last evaluation
- we are risk neutral
- ~~we are only willing to select a previously evaluated point as a final solution~~

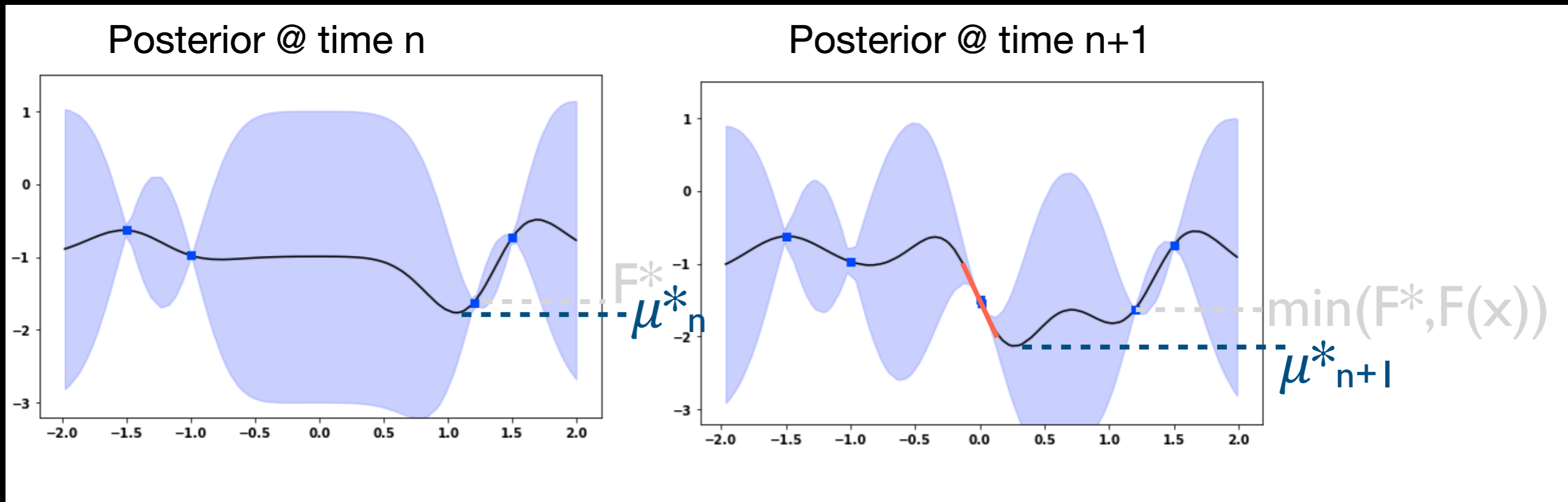
The KG acquisition function, $\text{KG}(\mathbf{x})$



- Loss if we stop now:
 $\mu_n^* = \min_{\mathbf{x}} \mu_n(\mathbf{x})$

$\mu_n(\mathbf{x}) := E_n[F(\mathbf{x})]$ is the expected value of our objective under the posterior @ time n

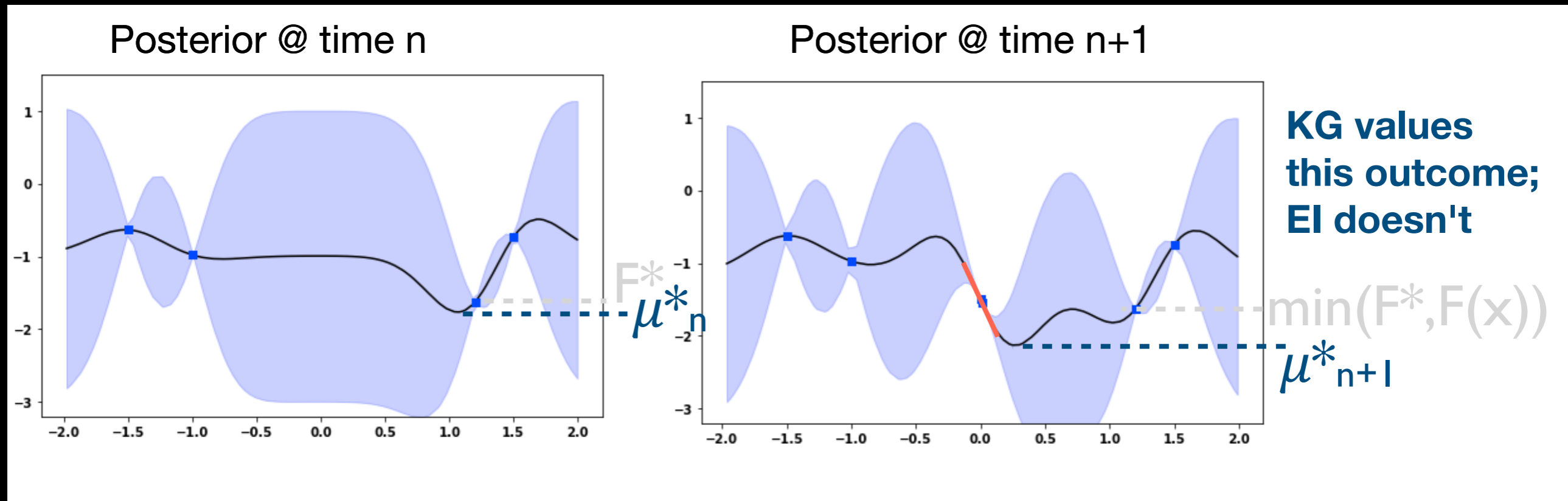
The KG acquisition function, $\text{KG}(\mathbf{x})$



- Loss if we stop now:
 $\mu_n^* = \min_{\mathbf{x}} \mu_n(\mathbf{x})$
- Loss if we stop after sampling $F(\mathbf{x})$ and its gradient:
 $\mu_{n+1}^* = \min_{\mathbf{x}} \mu_{n+1}(\mathbf{x})$

$\mu_n(\mathbf{x}) := E_n[F(\mathbf{x})]$ is the expected value of our objective under the posterior @ time n

The KG acquisition function, $KG(x)$



- Loss if we stop now:

$$\mu_n^* = \min_x \mu_n(x)$$

- Loss if we stop after sampling $F(x)$ and its gradient:

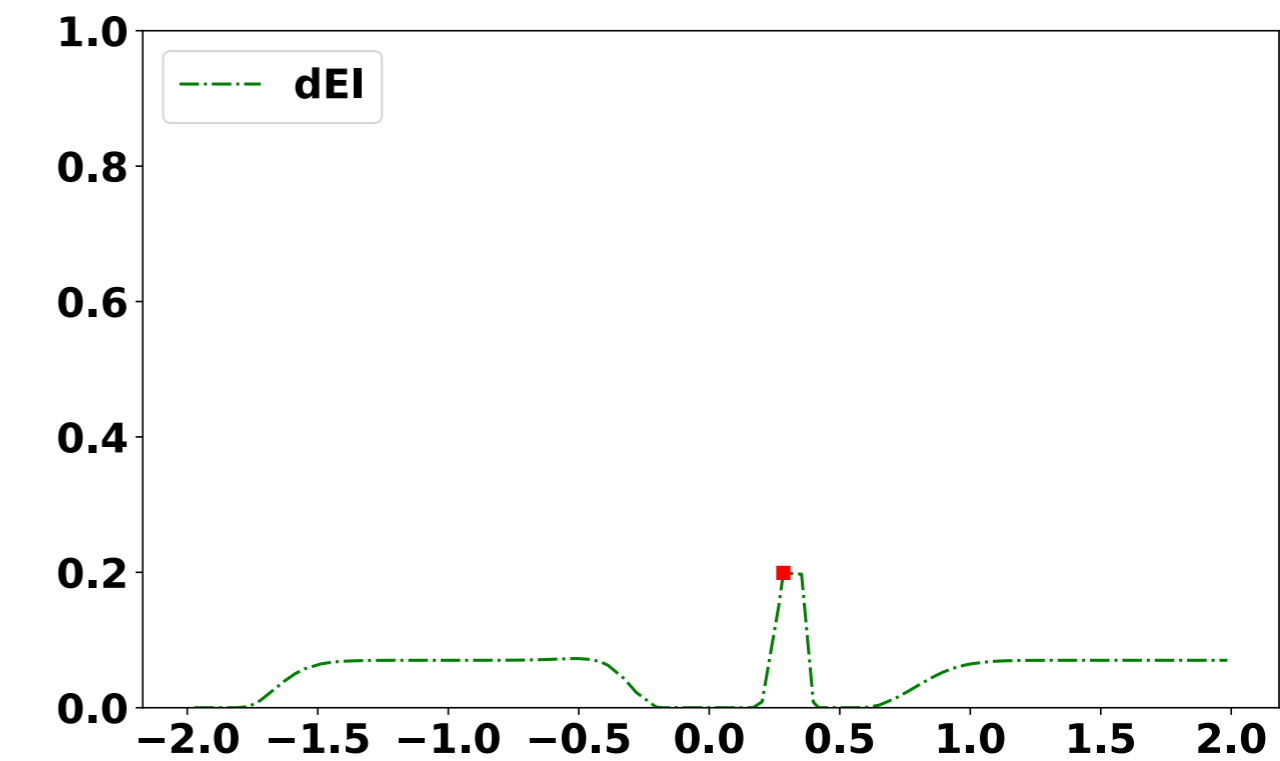
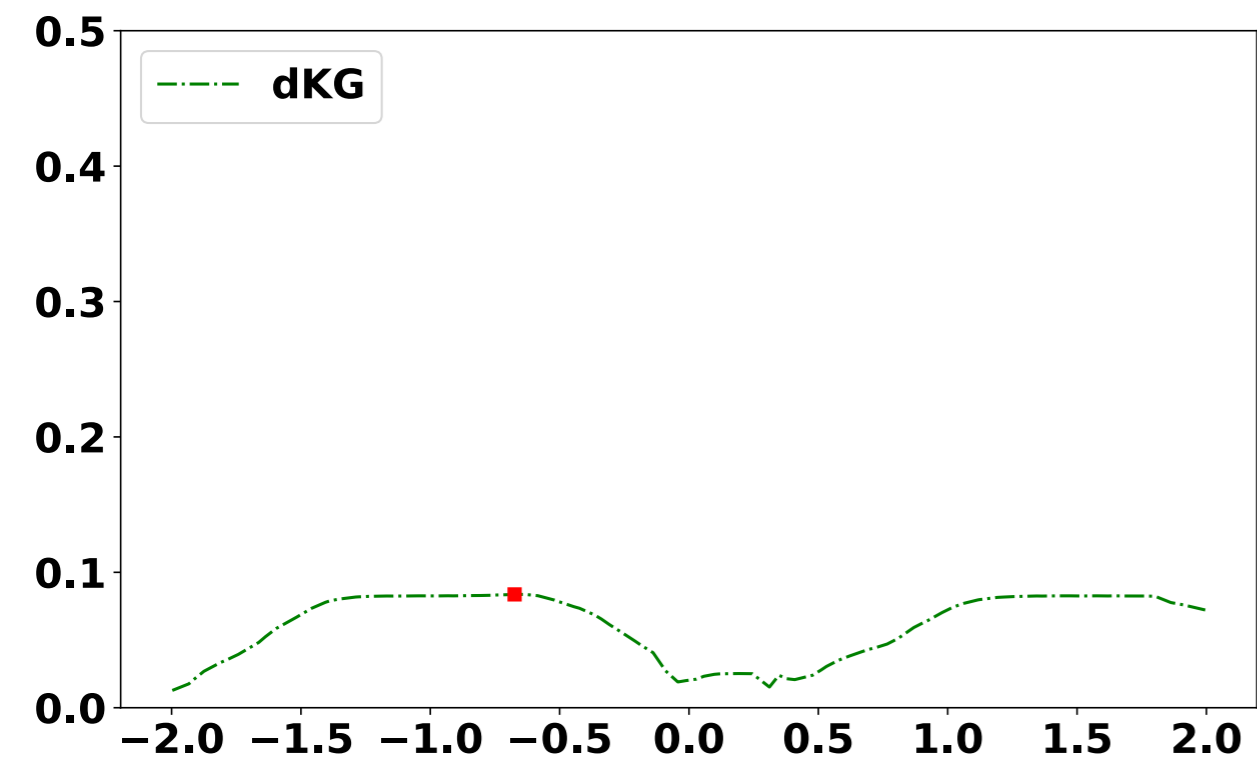
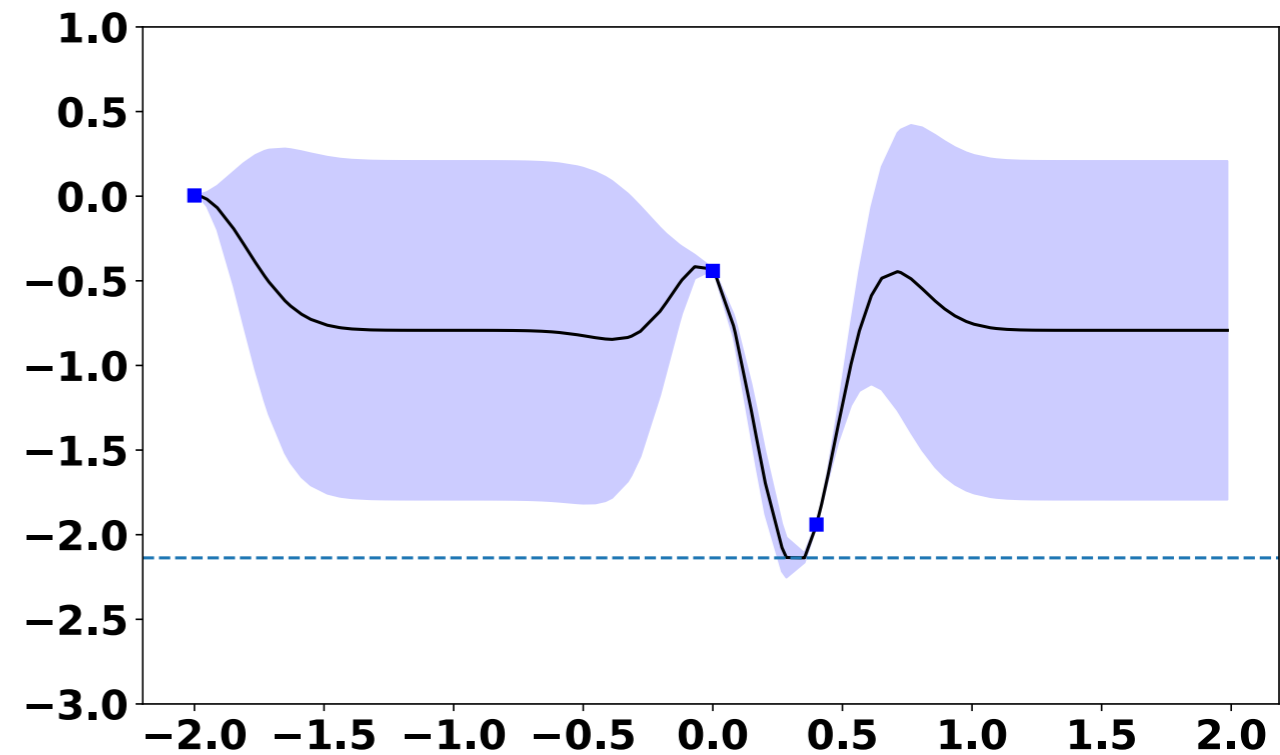
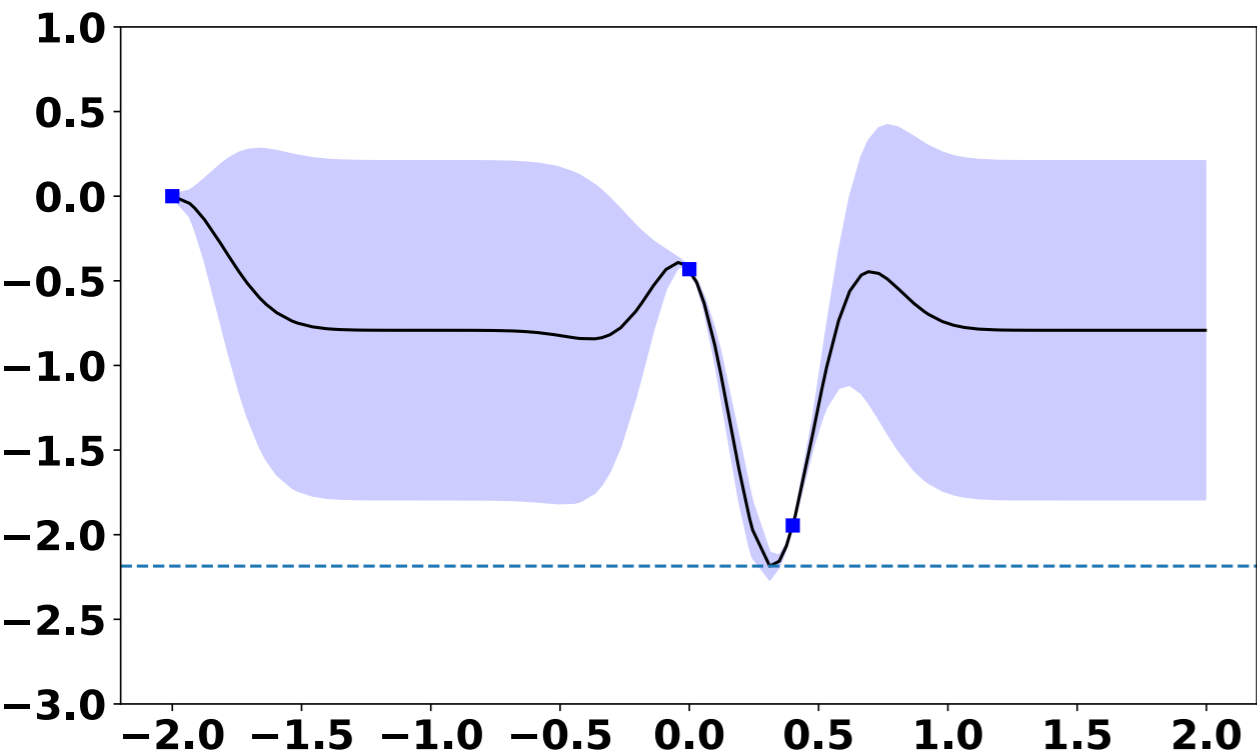
$$\mu_{n+1}^* = \min_x \mu_{n+1}(x)$$

- Reduction in loss due to sampling:

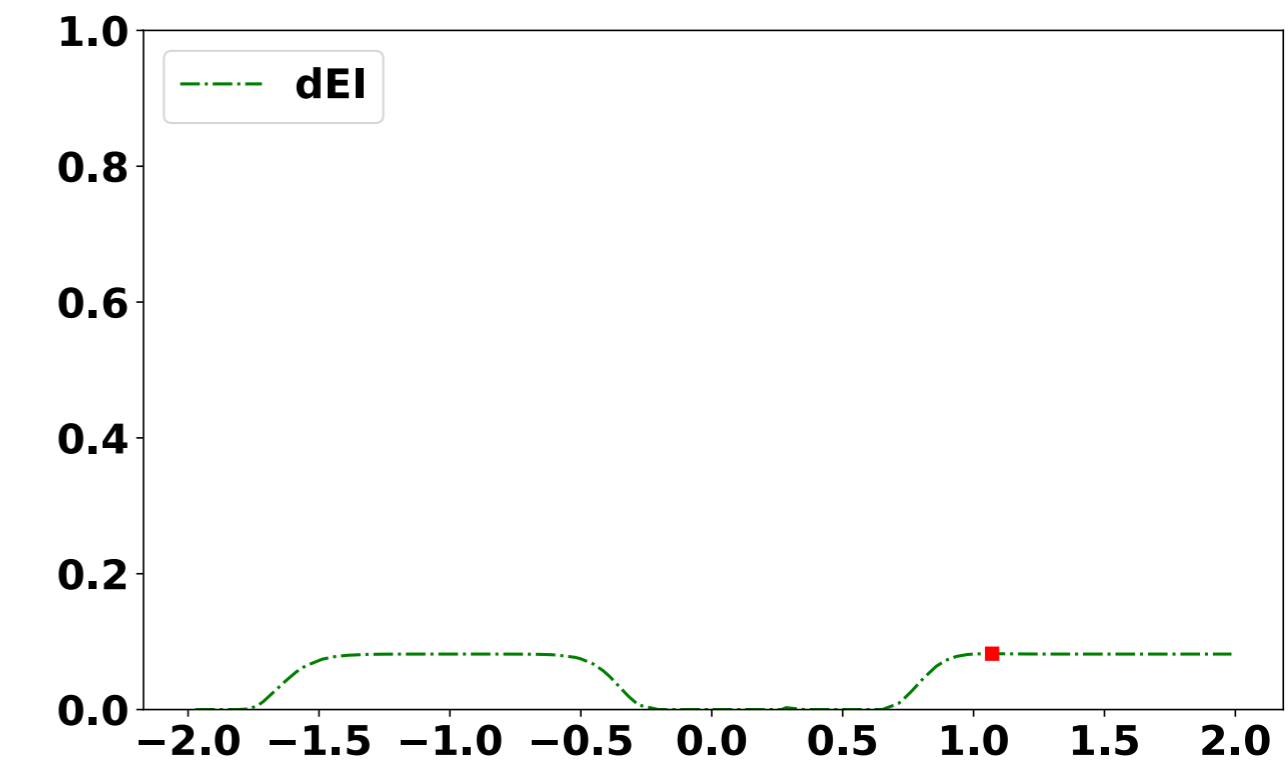
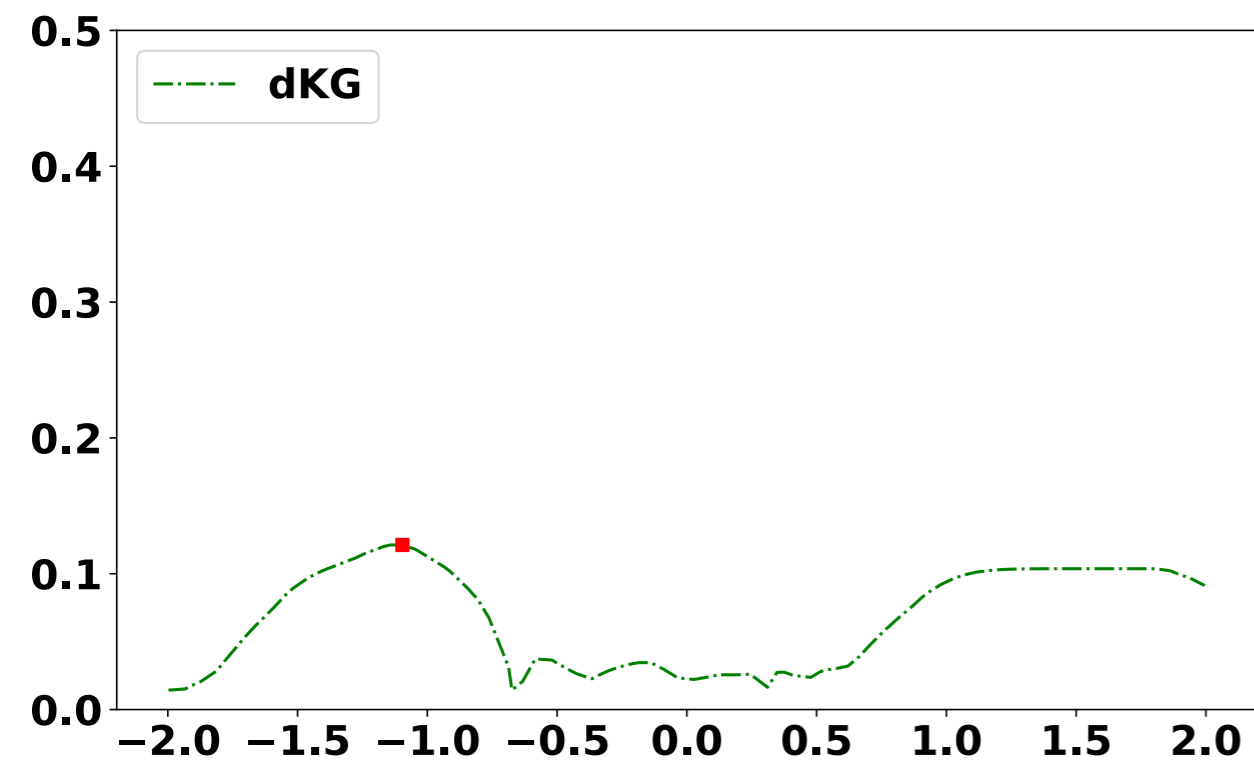
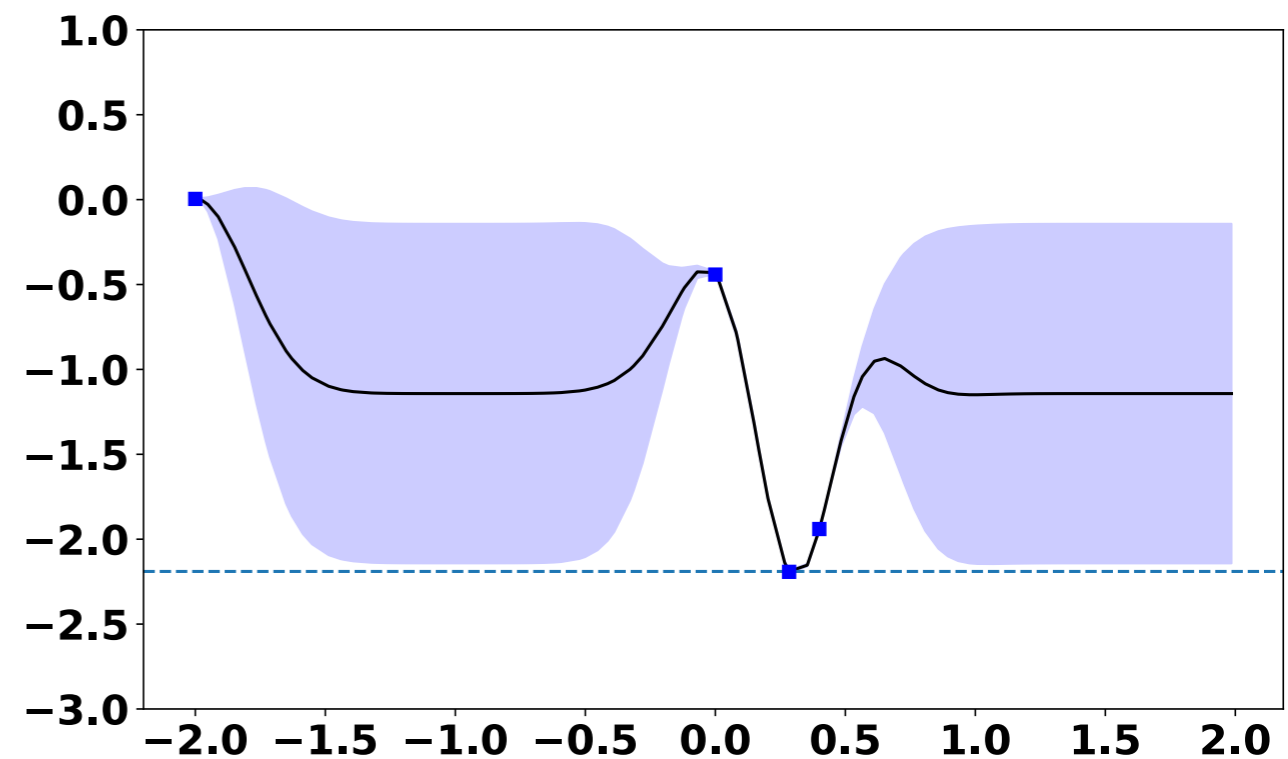
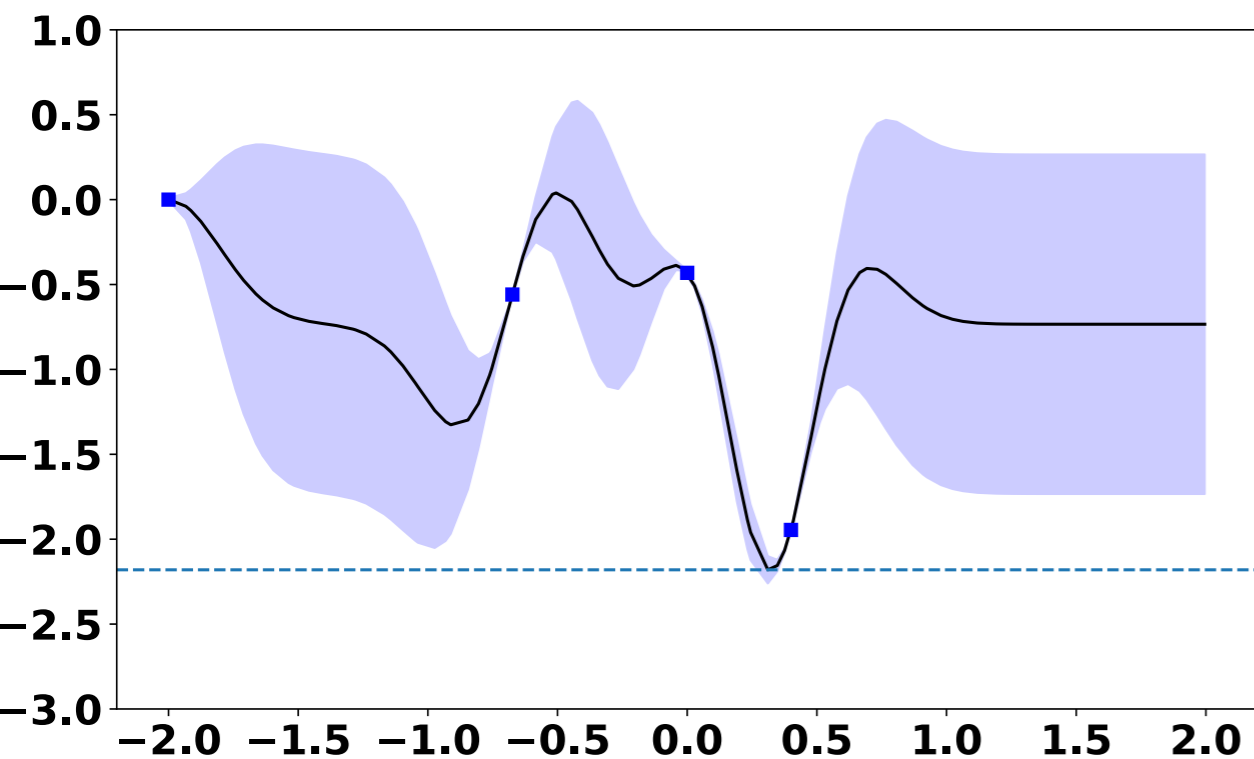
$$KG(x) = E_n[\mu_n^* - \mu_{n+1}^* \mid \text{sample } x]$$

$\mu_n(x) := E_n[F(x)]$ is the expected value of our objective under the posterior @ time n

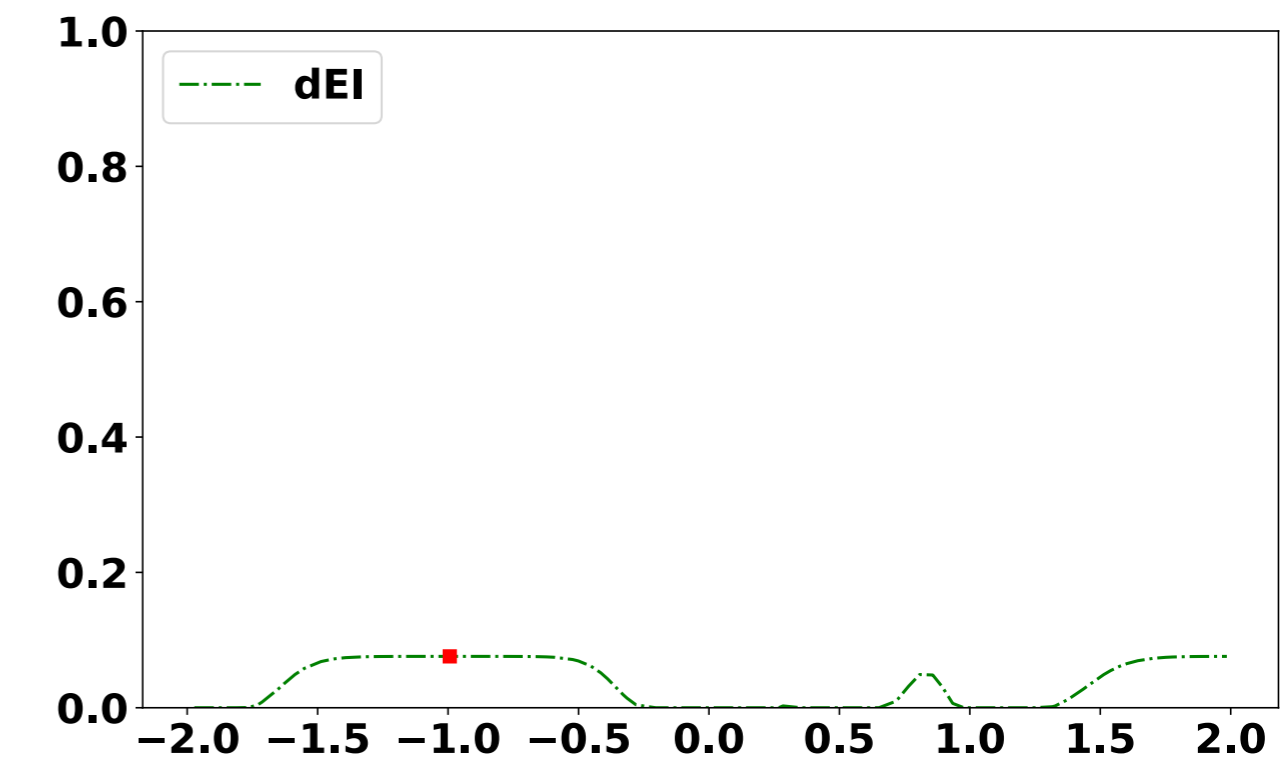
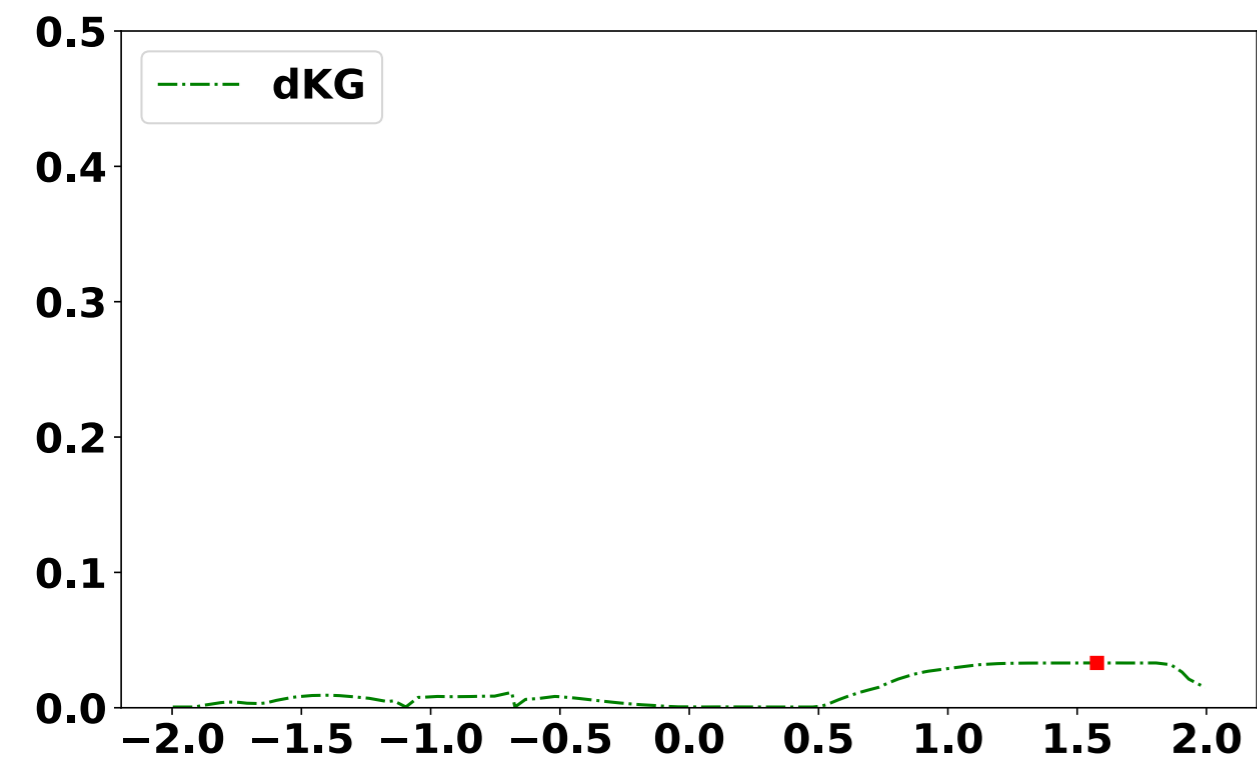
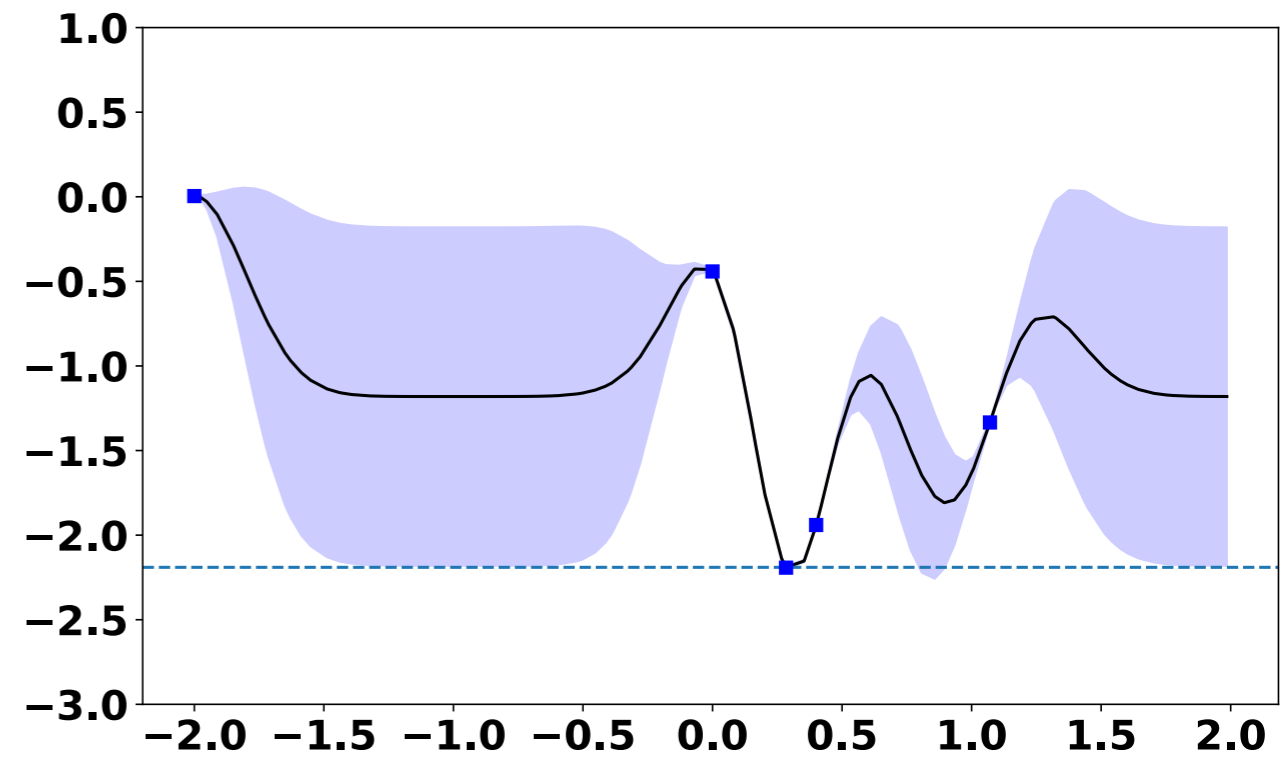
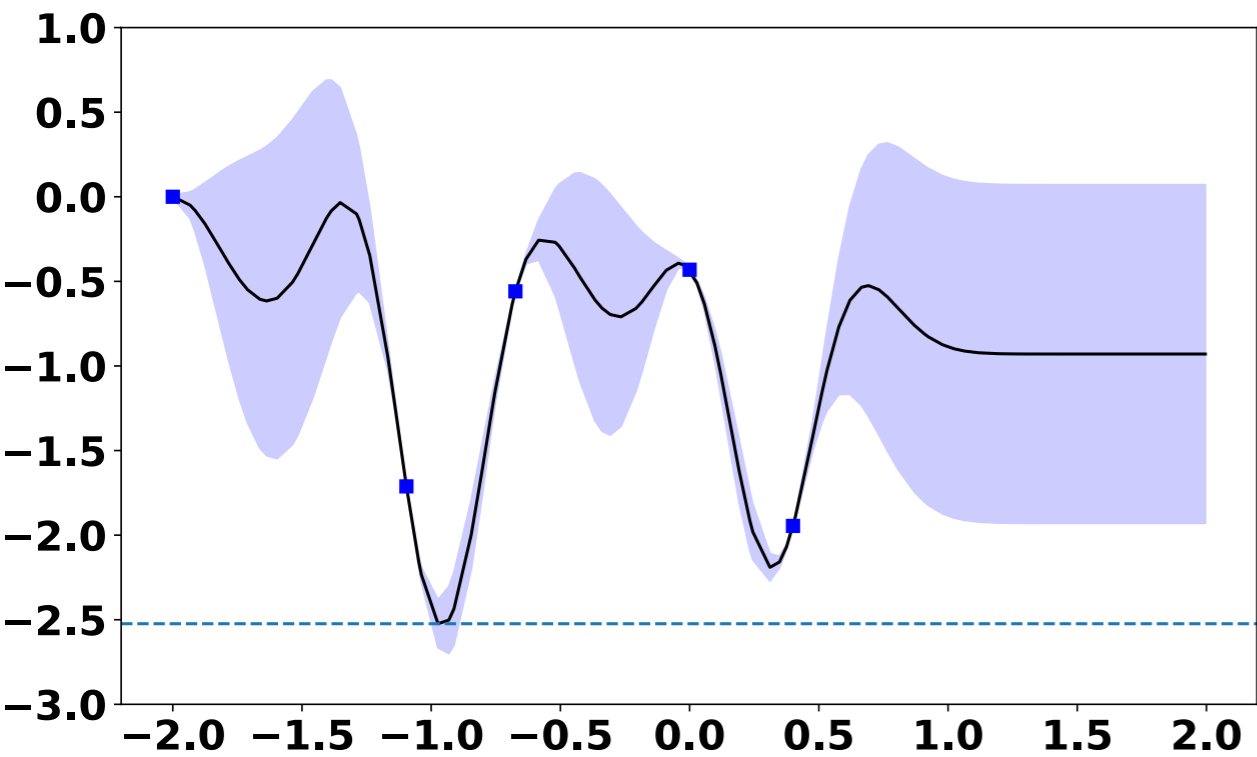
KG explores more effectively than EI when we have gradients



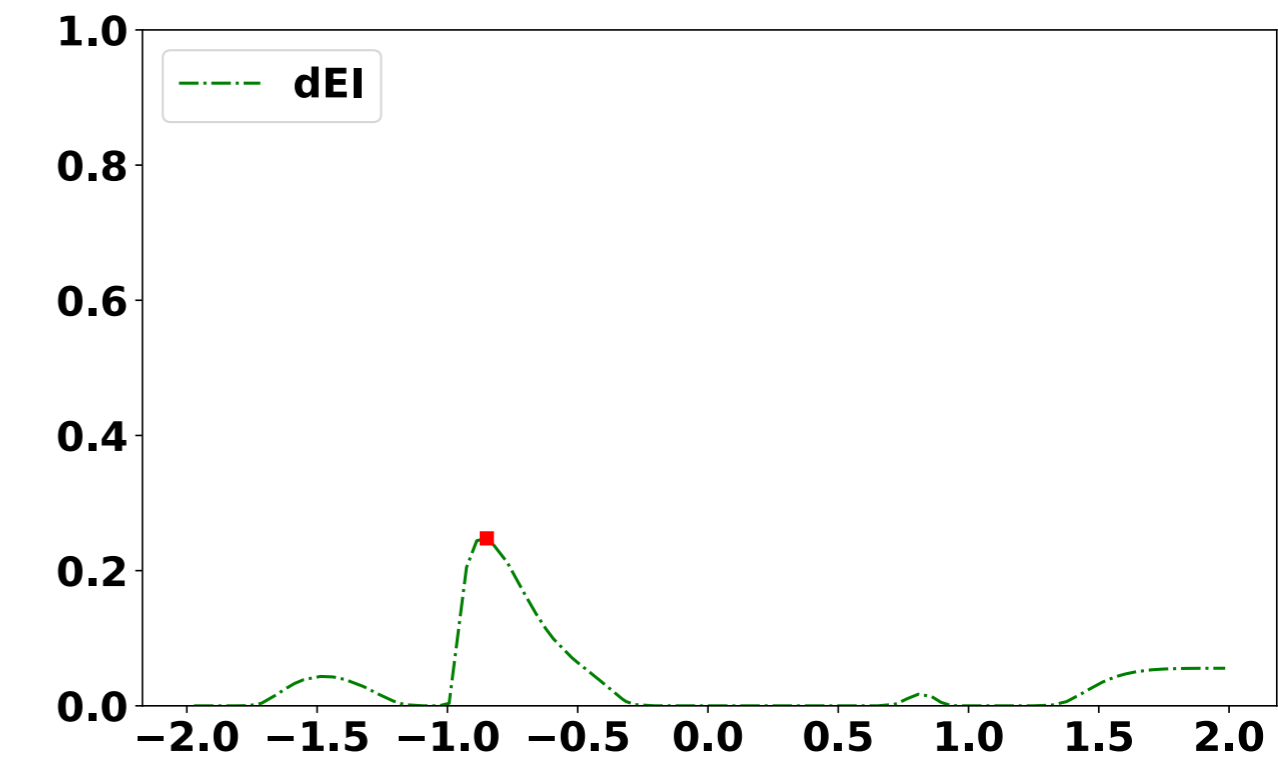
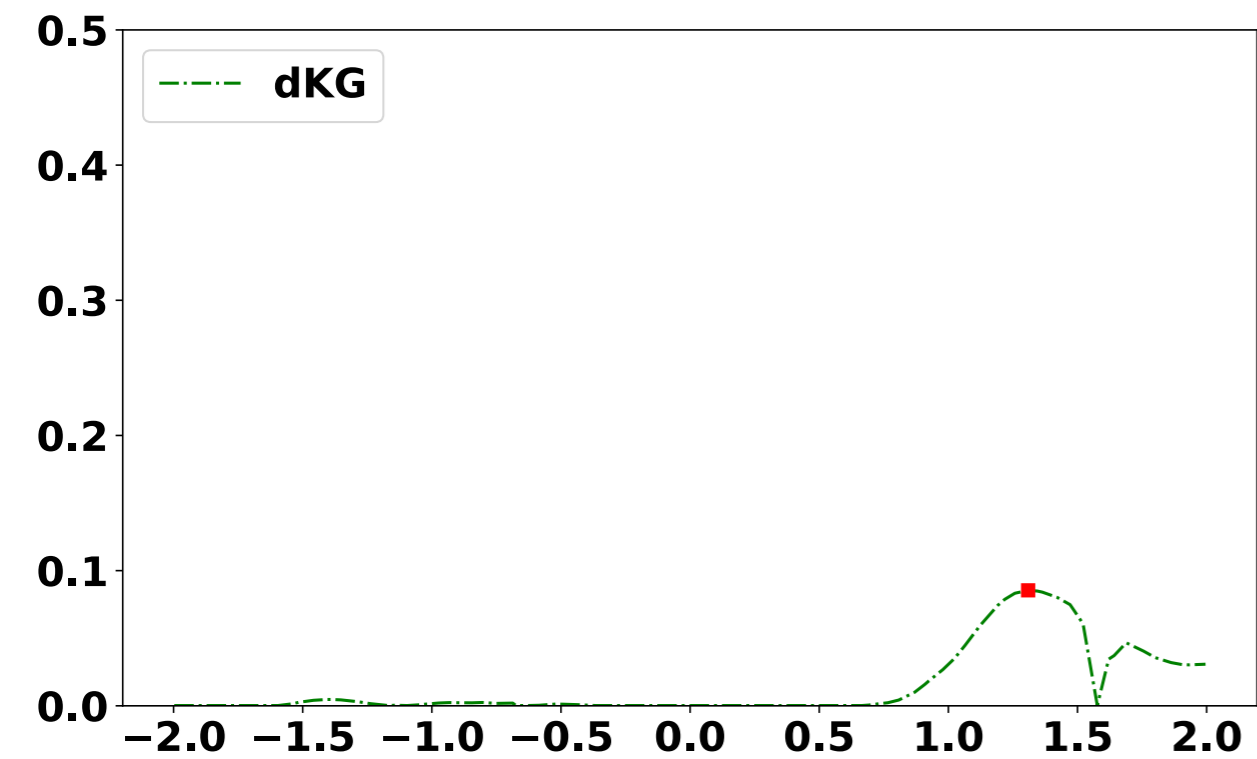
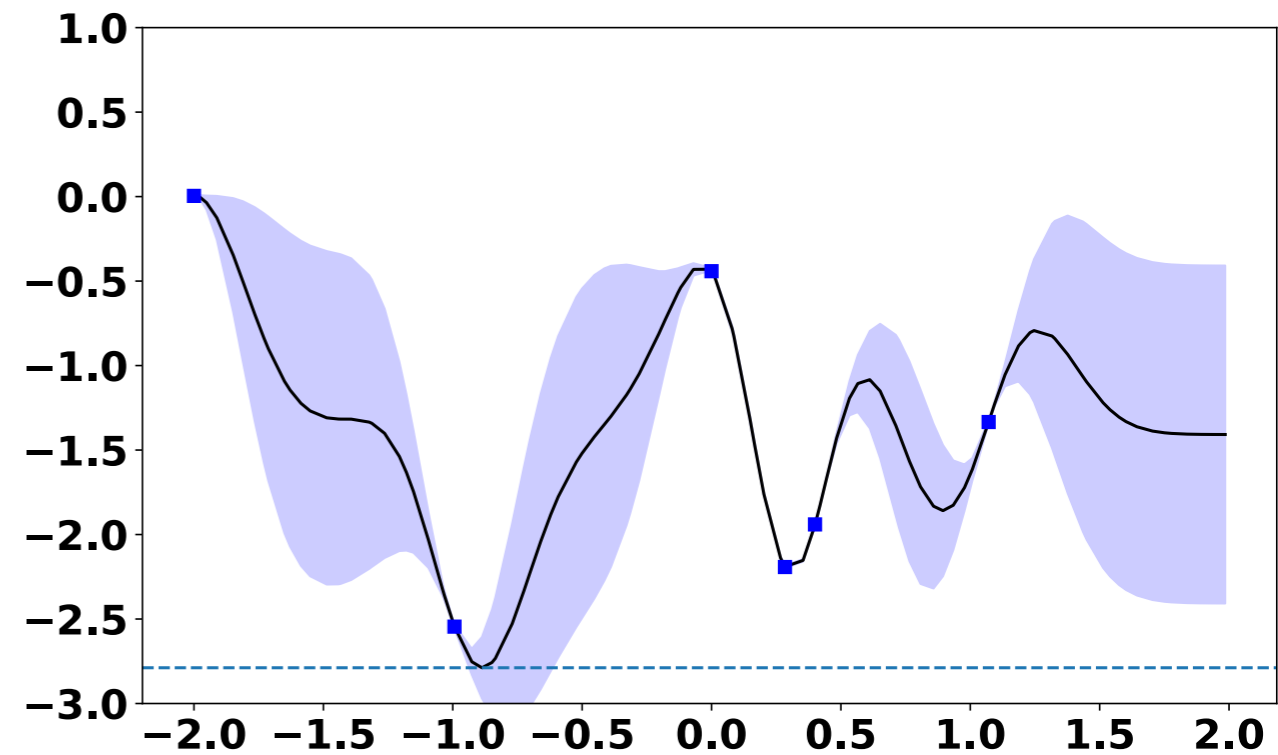
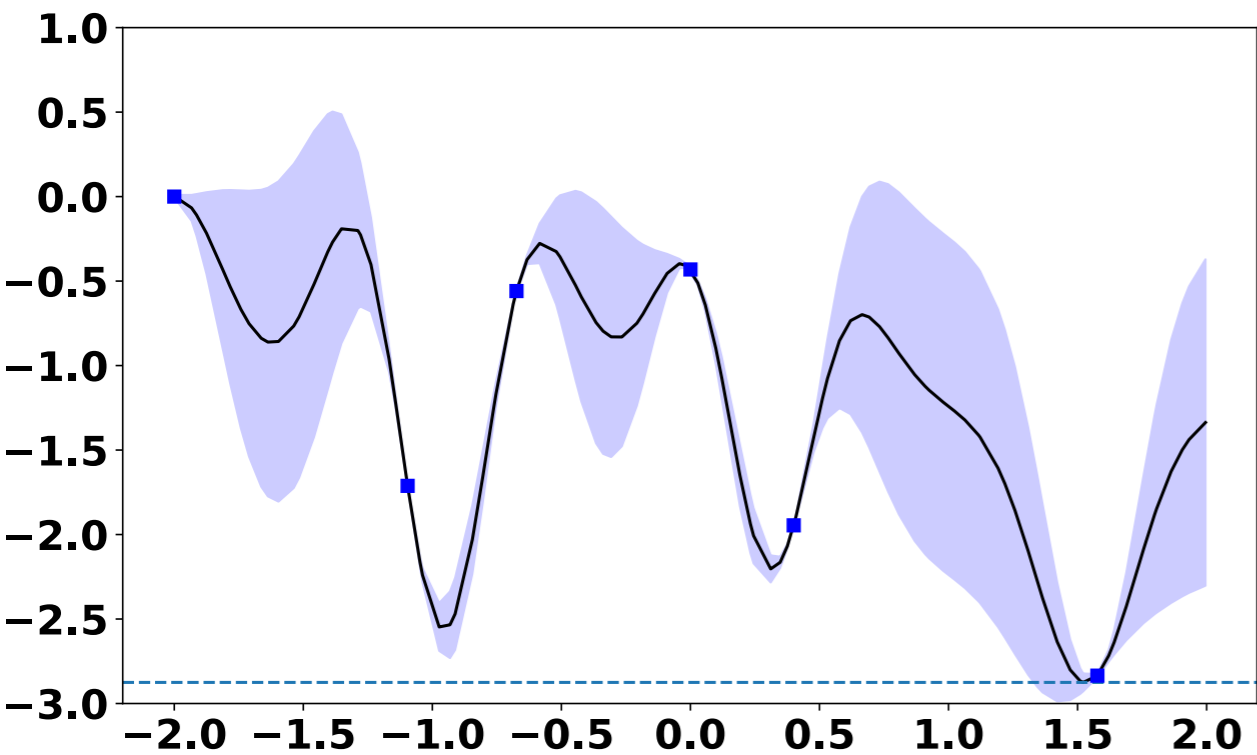
KG explores more effectively than EI when we have gradients



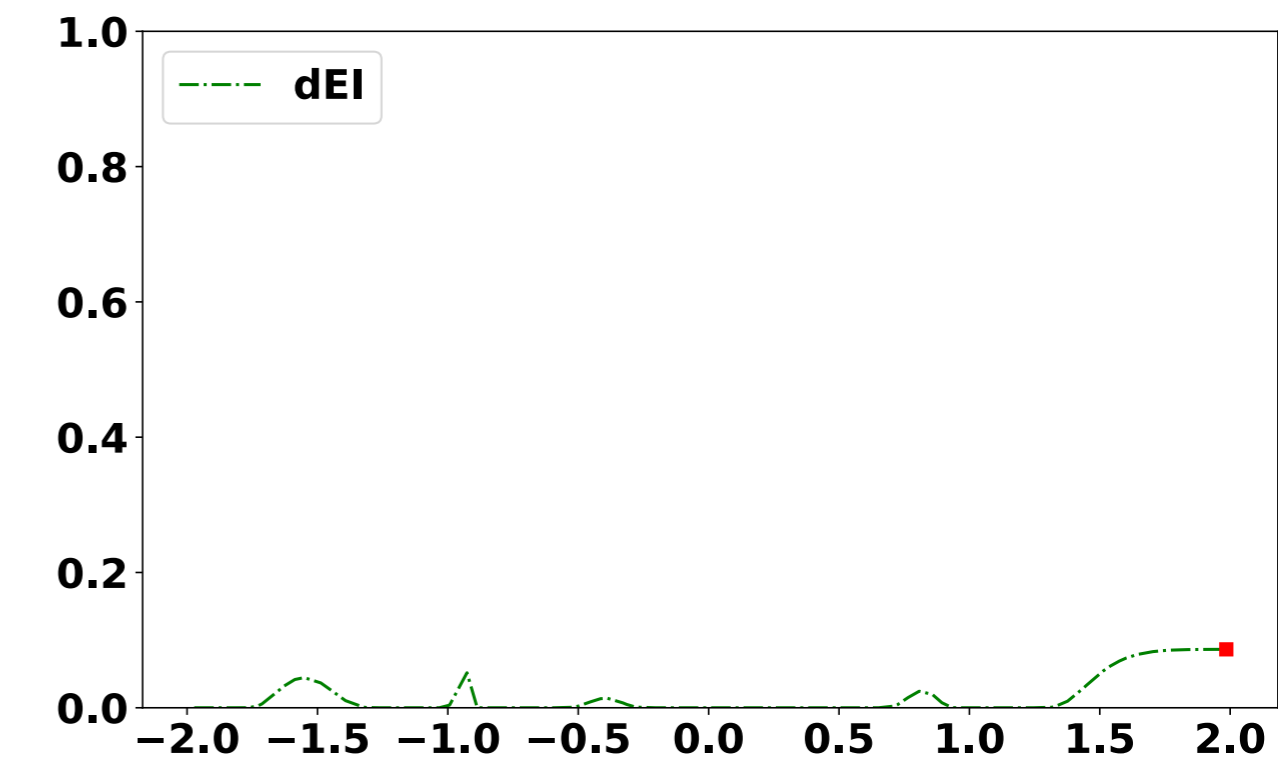
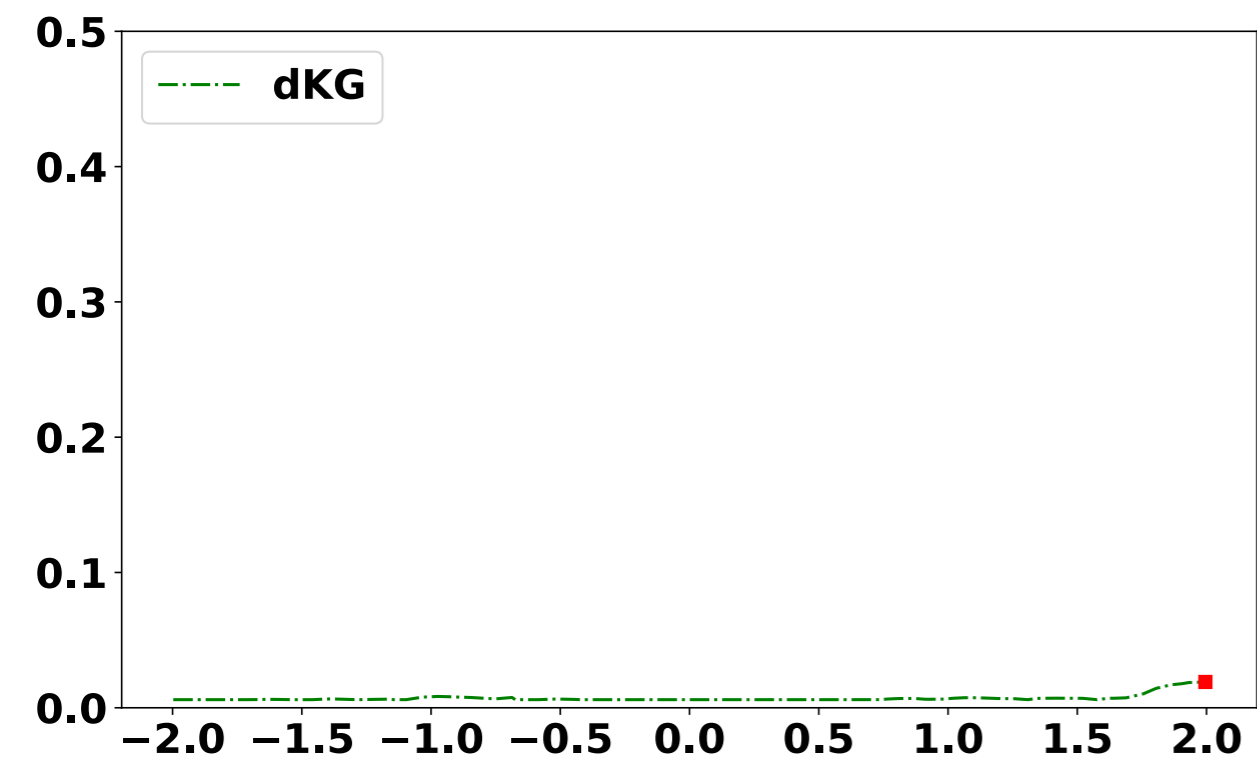
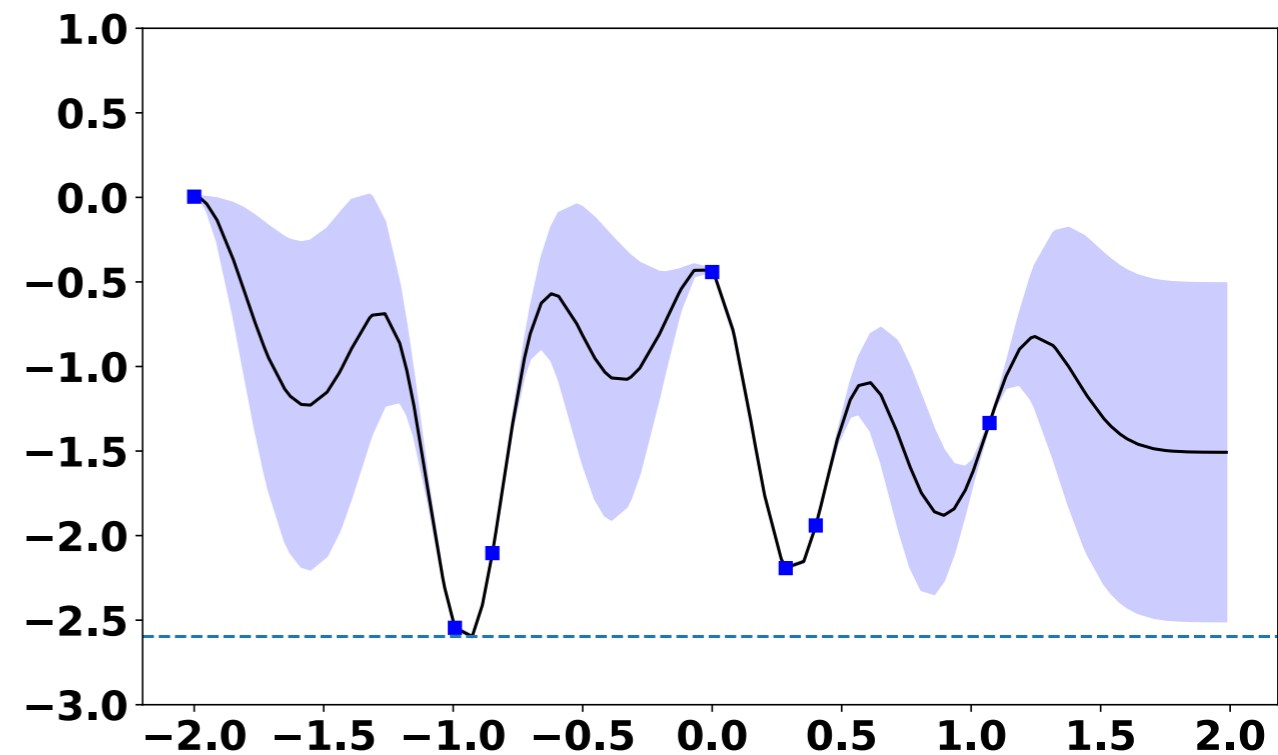
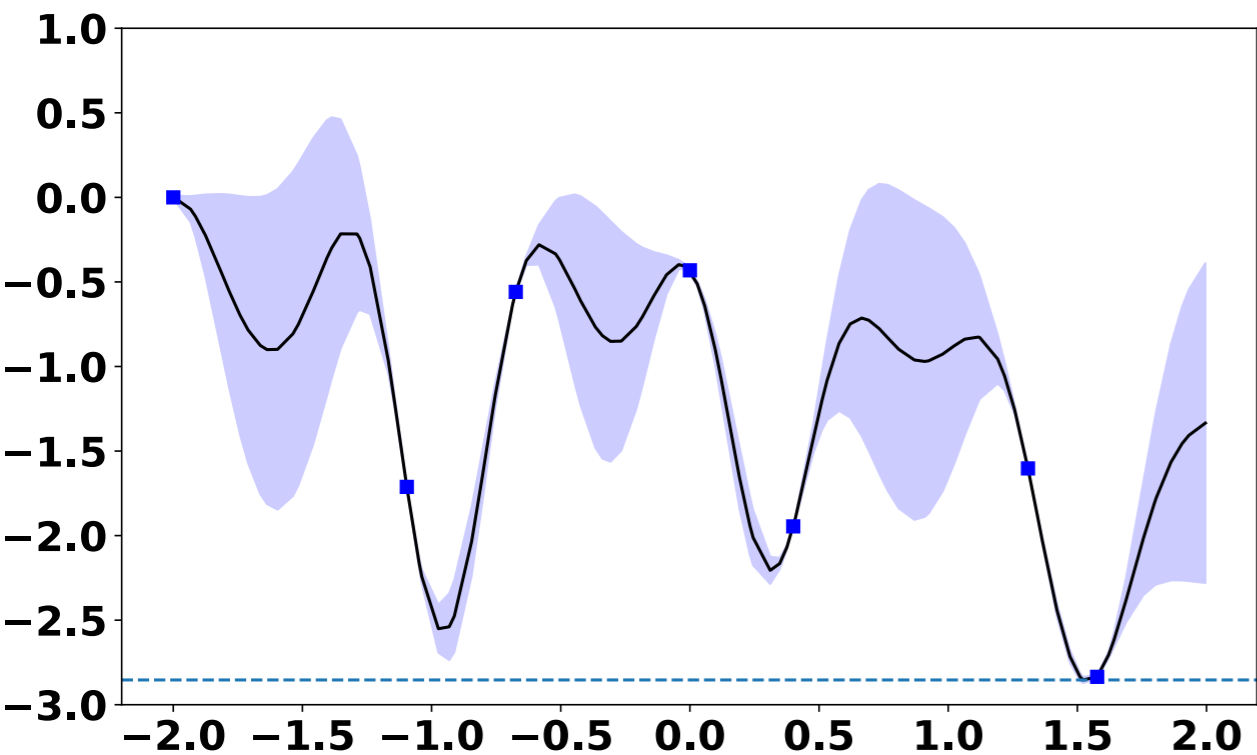
KG explores more effectively than EI when we have gradients



KG explores more effectively than EI when we have gradients



KG explores more effectively than EI when we have gradients





**How can we
optimize KG
efficiently?**

Here's how to **compute** KG

For i in 1:replications

- Simulate $Y(\mathbf{x}_{1:q}) = [F(\mathbf{x}_1), \nabla F(\mathbf{x}_1), \dots, F(\mathbf{x}_q), \nabla F(\mathbf{x}_q)]$
from the posterior on F
(plus noise, if observations are noisy)
- Calculate $\mu_{n+1}^* = \min_{\mathbf{x}'} \mu_{n+1}(\mathbf{x}')$ using $Y(\mathbf{x}_{1:q})$
- Calculate $\mu_n^* - \mu_{n+1}^*$

$\text{KG}(\mathbf{x}_{1:q})$ is the average of the simulated $\mu_n^* - \mu_{n+1}^*$

Our approach for optimizing parallel EI
also works for optimizing KG

1. Estimate $\nabla \text{KG}(\mathbf{x}_{1:q})$
using infinitesimal perturbation analysis (IPA)
& the envelope theorem
2. Use multistart stochastic gradient ascent to
maximize $\text{KG}(\mathbf{x}_{1:q})$

Here's how to **estimate** $\nabla KG(x_{1:q})$

To estimate $\nabla KG(x_{1:q})$, use:

infinitesimal perturbation analysis (IPA) (Ho & Cao 2012)

$$\begin{aligned}\nabla KG(\mathbf{x}) &= \nabla E_n[\mu^*_n - \mu^*_{n+1}(Y(\mathbf{x}), \mathbf{x})] \\ &= E_n[\nabla \mu^*_n - \mu^*_{n+1}(Y(\mathbf{x}), \mathbf{x})]\end{aligned}$$

& the envelope theorem (Milgrom & Segal 2002)

$$\begin{aligned}\nabla \mu^*_{n+1}(Y(\mathbf{x}), \mathbf{x}) &= \nabla \min_{\mathbf{x}'} \mu_{n+1}(\mathbf{x}'; Y(\mathbf{x}), \mathbf{x}) \\ &= \nabla \min_{\mathbf{x}'} \mu_{n+1}(\mathbf{x}'; m(\mathbf{x}) + C(\mathbf{x})Z, \mathbf{x}) \\ &= \nabla \mu_{n+1}(\mathbf{x}^*; m(\mathbf{x}) + C(\mathbf{x})Z, \mathbf{x})\end{aligned}$$

we write \mathbf{x} instead of $\mathbf{x}_{1:q}$ to reduce clutter

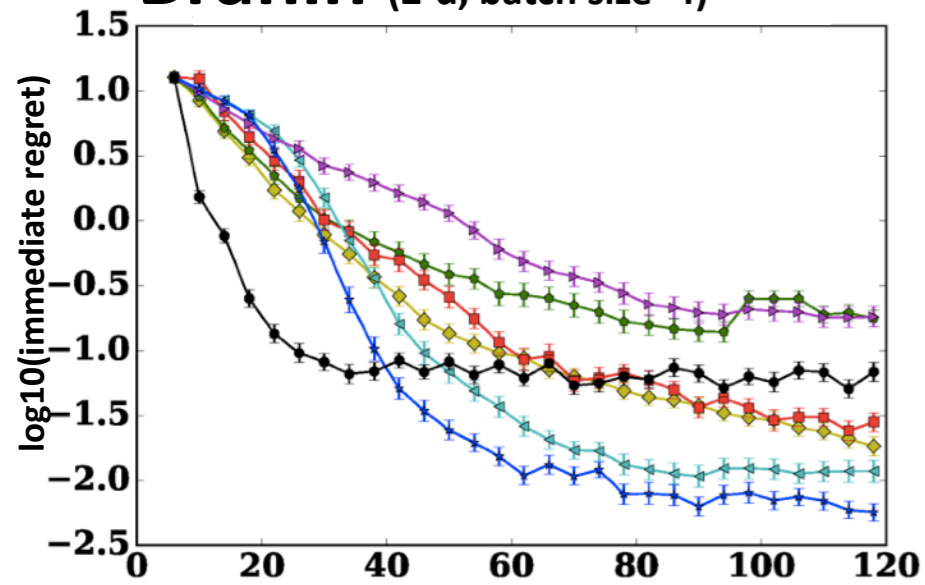
Here's how to **estimate** $\nabla\text{KG}(x_{1:q})$

1. Simulate a vector Z of independent standard normals
2. Calculate $m = E[Y]$, $C = \text{Chol}(\text{Cov}[Y])$,
where $Y = [F(x_1), \nabla F(x_1), \dots, F(x_q), \nabla F(x_q)]$
3. Solve $\min_{x'} \mu_{n+1}(x'; m(x) + C(x)Z, x)$ using L-BFGS-B.
Let x^* be the solution.
4. Our estimator of $\nabla\text{KG}(x)$ is $\nabla\mu_{n+1}(x^*; m(x) + C(x)Z, x)$,
calculating the gradient holding x^* fixed

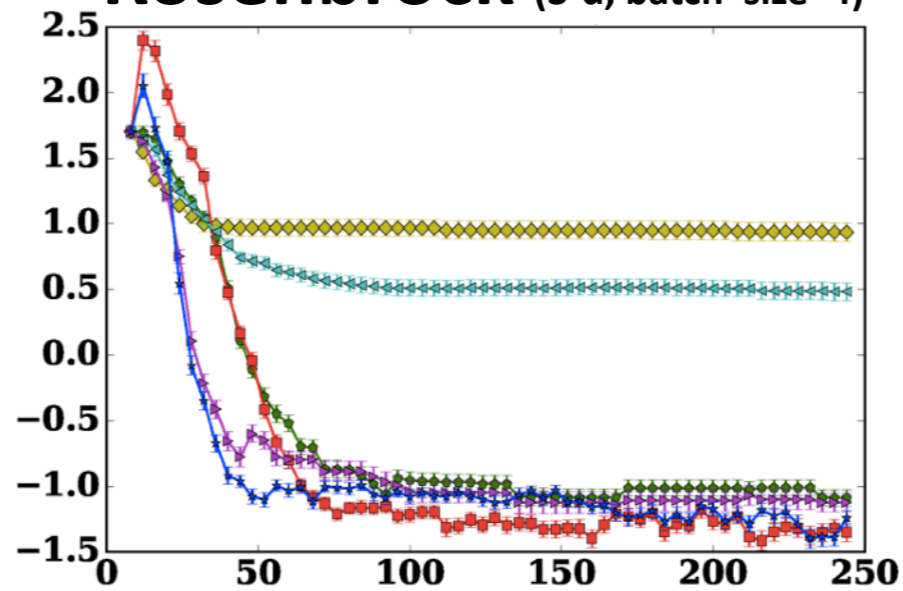
we write x instead of $x_{1:q}$ to reduce clutter

KG outperforms EI substantially, if you observe gradients

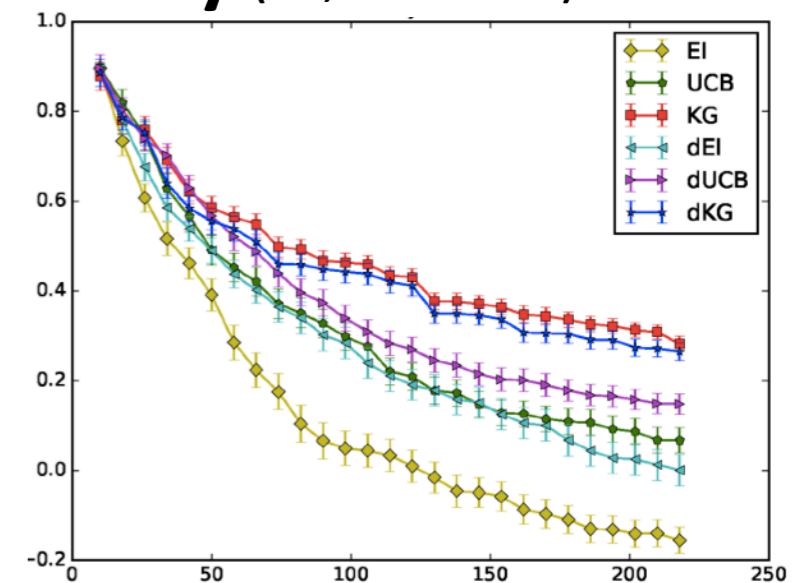
Branin (2-d, batch size=4)



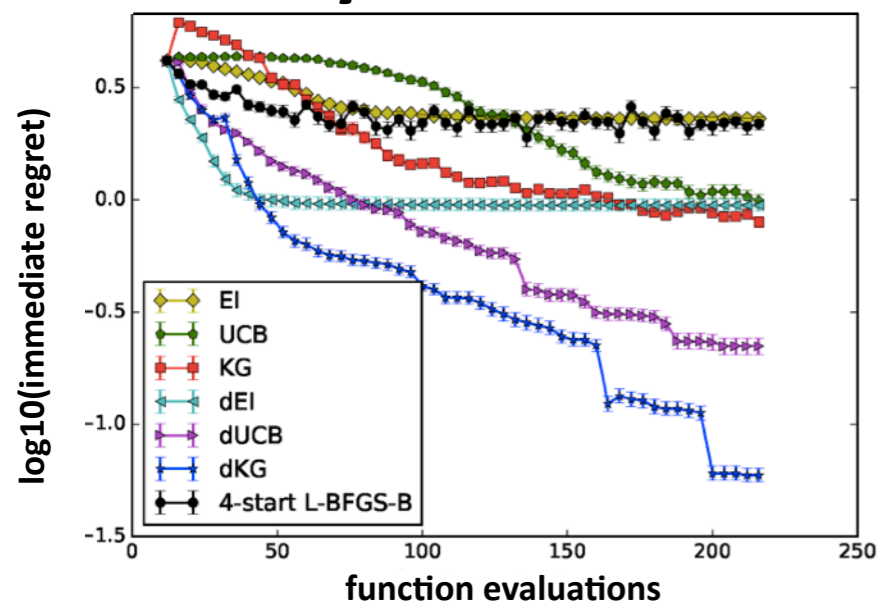
Rosenbrock (3-d, batch size=4)



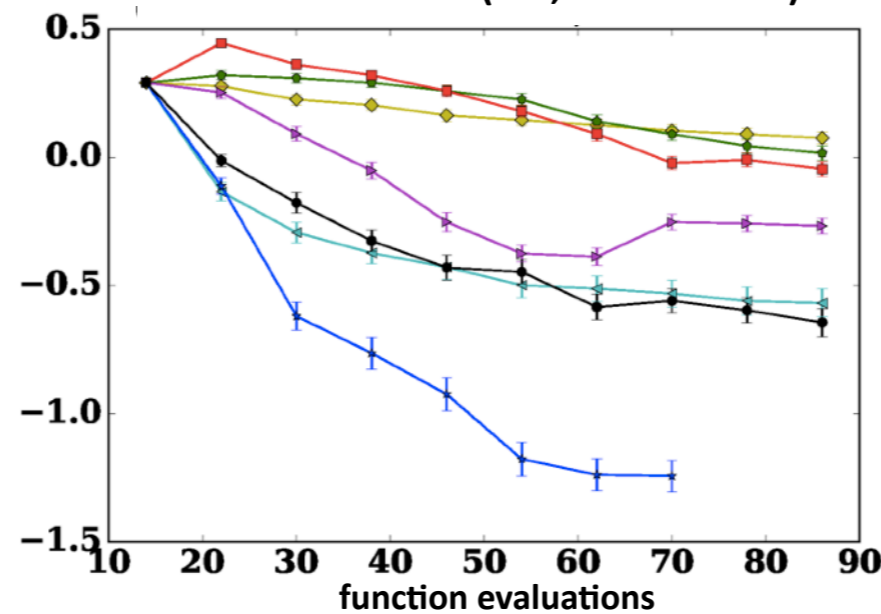
Levy (4-d, batch size=8)



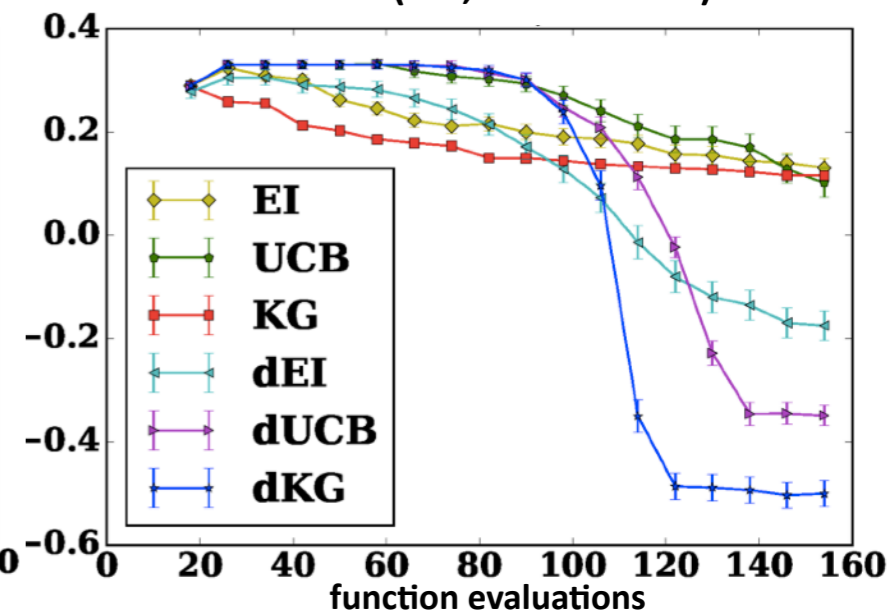
Ackley (5-d, batch size=4)



Hartmann (6-d, batch size=8)

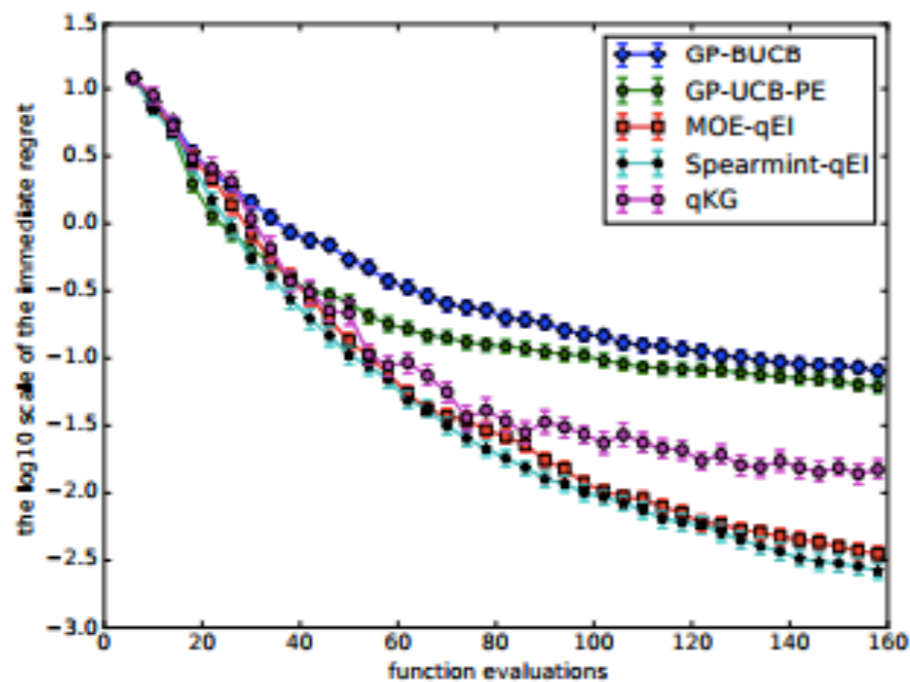


Cosine (8-d, batch size=4)

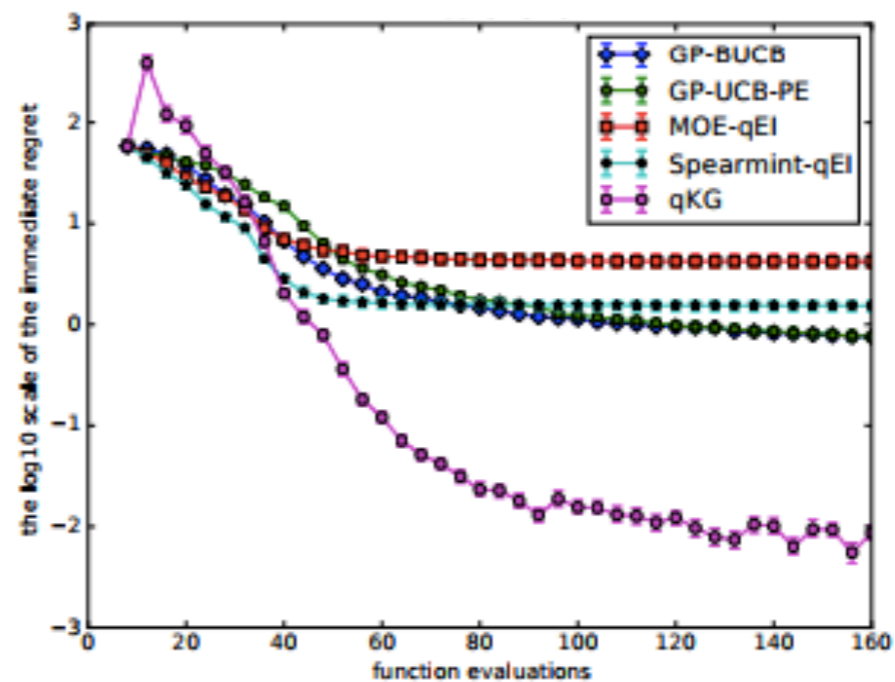


KG outperforms EI substantially
if you **don't observe gradients**,
but have noise or are in $\text{dim} > 1$

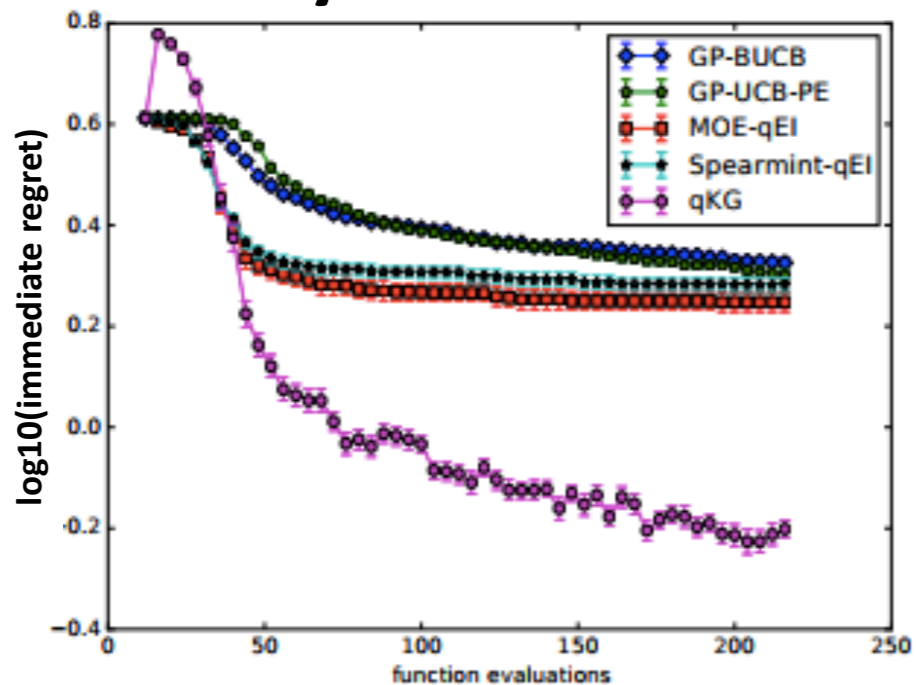
Branin (2-d, batch size=4, no noise)



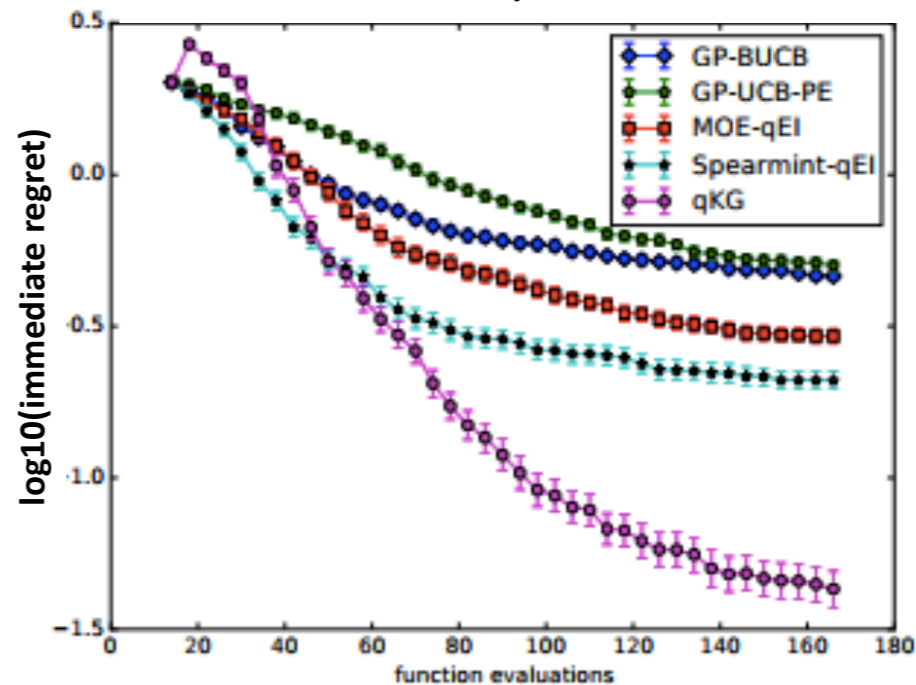
Rosenbrock (3-d, batch size=4, no noise)



Ackley (5-d, batch size=4, no noise)

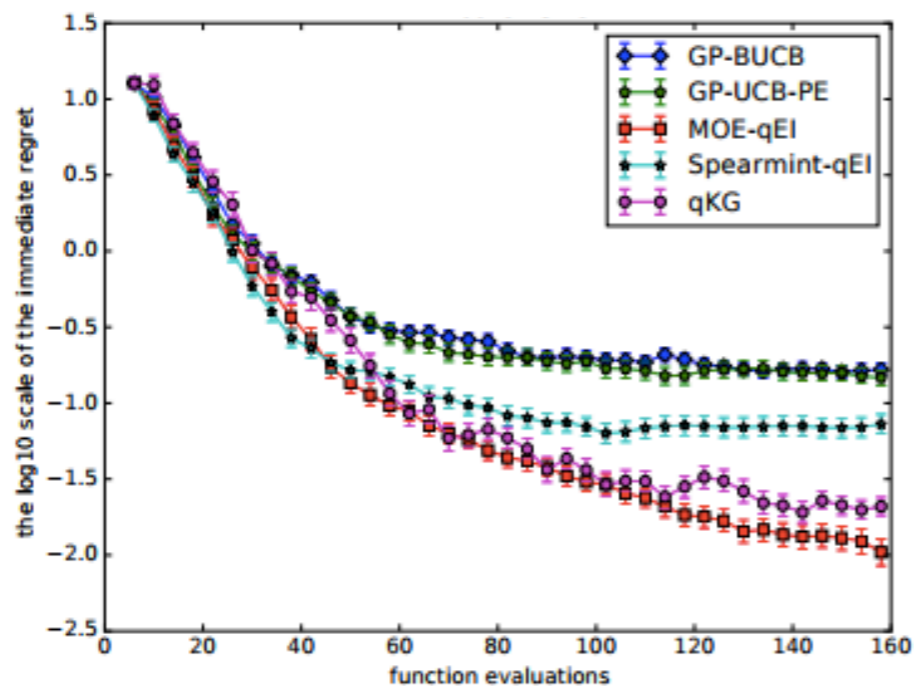


Hartmann (6-d, batch size=4, no noise)

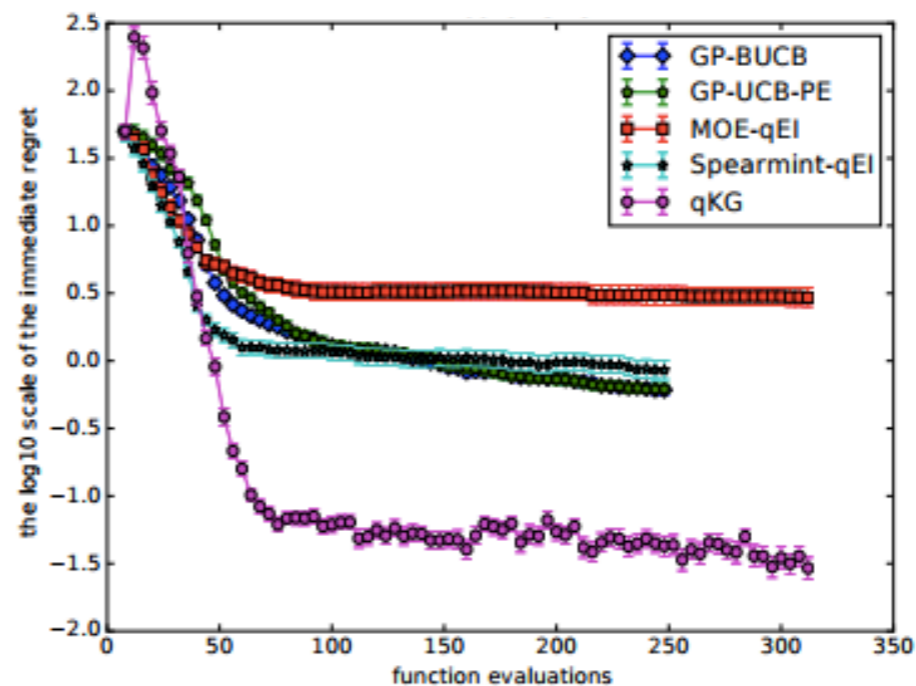


KG outperforms EI substantially
if you **don't observe gradients**,
but have noise or are in $\text{dim} > 1$

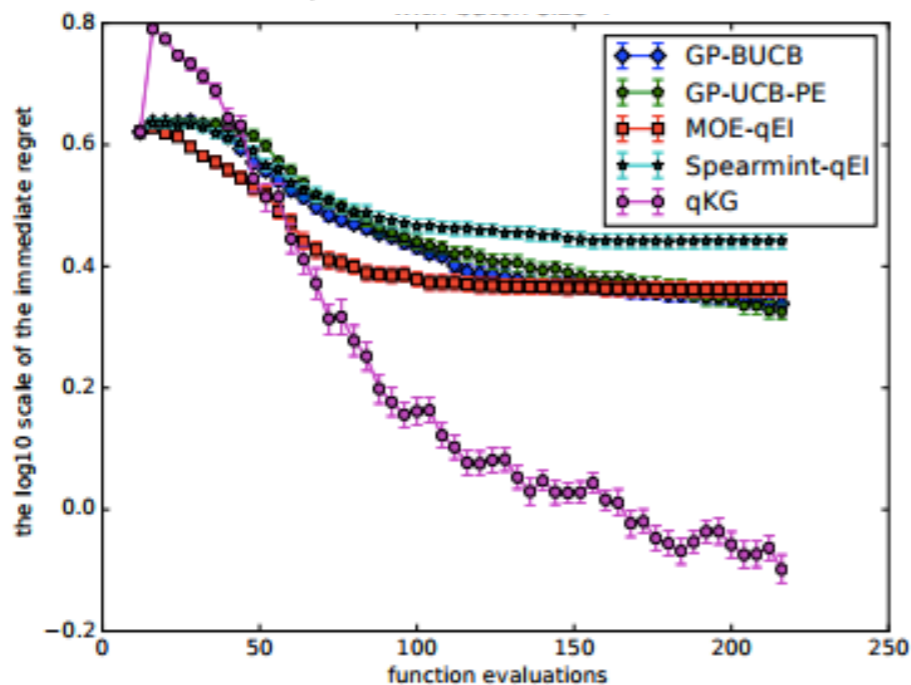
Branin (2-d, batch size=4, noisy)



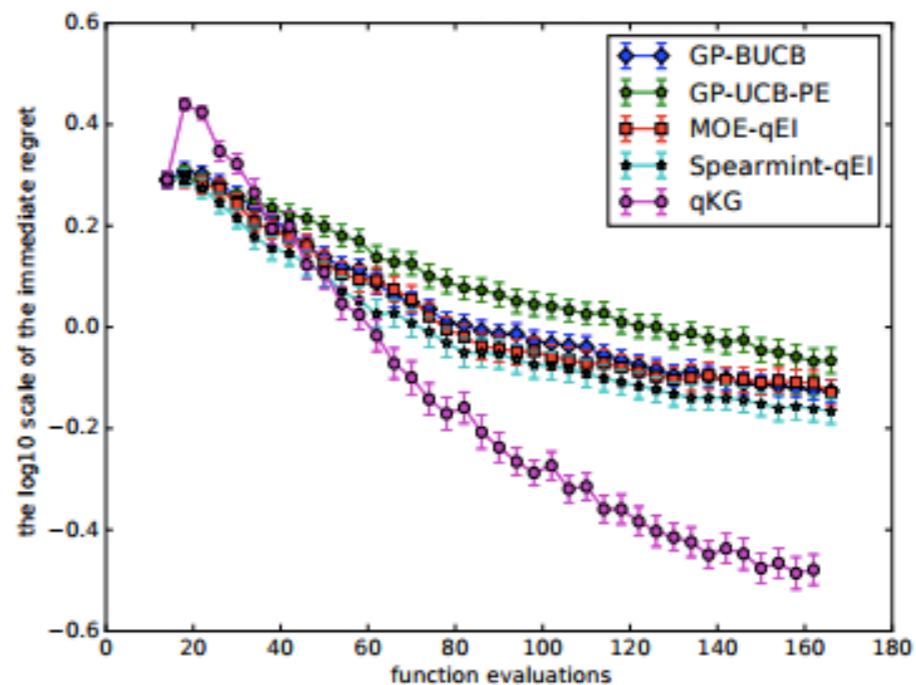
Rosenbrock (3-d, batch size=4, noisy)



Ackley (5-d, batch size=4, noisy)



Hartmann (5-d, batch size=4, noisy)



Agenda

- Gaussian process regression
- Expected improvement (EI)
- Parallel BayesOpt (using EI)
- Parallel BayesOpt with gradients and/or noise (using KG)
- **Code & Research directions**

Here are some open-source codes for Bayesian optimization

- Cornell Metrics Optimization Engine (MOE):
<https://github.com/wujian16/Cornell-MOE>
based on Yelp's MOE codebase
<https://github.com/Yelp/MOE>
- GPyOpt:
<https://github.com/SheffieldML/GPyOpt>
- Spearmint:
<https://github.com/JasperSnoek/spearmint>

There are also commercial services

Secure | https://sigopt.com



See Demo

Pricing

FAQ

Log In

Sign Up

Improve ML models 100x faster

SigOpt's API tunes your model's parameters through *state-of-the-art* Bayesian optimization.

- Exponentially faster and more accurate than grid search. Faster, more stable, and easier to use than open source solutions.
- Extracts additional revenue and performance left on the table by conventional tuning.



Optimizing in-production models for

hotwire™

IQT.
IN · Q · TEL

WorkFusion

MIT

What didn't we cover?

- Statistical models beyond Gaussian processes
 - Warping
 - Student-t processes
 - Neural networks
 - Models for high-dimensional problems
- Other acquisition functions:
 - Entropy Search and predictive entropy search
 - Thompson sampling
 - Probability of improvement
- Convergence results & regret guarantees
- Problems with exotic structure
 - Multi-fidelity & multi-information source optimization
 - Optimization with trace observations
 - Multi-task optimization

Research Directions

- Applications
- Why do 1-step algorithms work so well?
Meaningful theoretical guarantees
- Going beyond Gaussian processes
- High-dimensional BayesOpt
- Problems with exotic structure

Thanks! Any Questions?

- Cornell Metrics Optimization Engine (MOE):
<https://github.com/wujian16/Cornell-MOE>
- "A Tutorial on Bayesian Optimization"
<https://arxiv.org/abs/1807.02811>

Backup

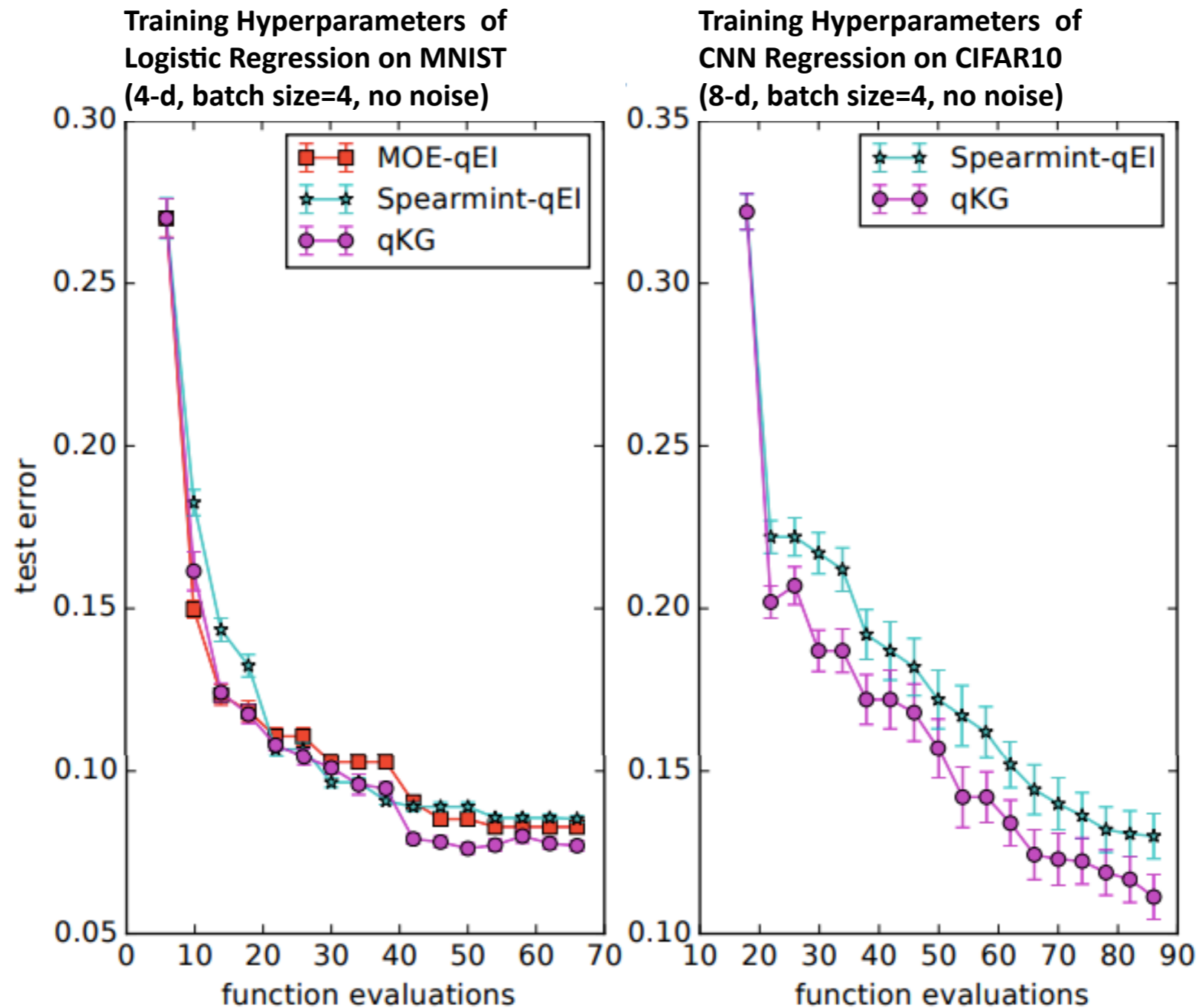
KG converges to a globally optimal solution over a discrete domain

Theorem: When the domain A is discrete & finite, the KG algorithm is consistent, i.e.,

$$\lim_{N \rightarrow \infty} f(x^*(\text{dKG}, N)) = \max_{x \in A} f(x)$$

almost surely under the prior, where $x^*(\text{KG}, N)$ is the solution computed by KG after N batches of samples.

KG provides substantial value over EI when we **don't have gradients** when there is noise, or we are in > 1 dim.



GP regression handles gradients naturally

[Rasmussen & Williams 2006]

1. $f \sim \text{GP}(\mu, K)$
2. $x \rightarrow [f(x), \nabla f(x)]$ is a vector-valued function.
3. Our GP prior on f implies this GP prior on $f, \nabla f$:

$$\begin{bmatrix} f \\ \nabla f \end{bmatrix} \sim \text{GP} \left(\begin{bmatrix} \mu \\ \nabla \mu \end{bmatrix}, \begin{bmatrix} K & \nabla K \\ \nabla K & \nabla^2 K \end{bmatrix} \right)$$

(Assuming some regularity conditions on K)

4. Do standard GP regression on $f, \nabla f$
Calculate the posterior on the objective f as needed

Bayesian optimization is really **flexible**

Elicit a prior distribution on the function f
while (budget is not exhausted) {

Find the **information source** whose value
of information is the largest.

Query that **information source**

Update the posterior distribution.

}

We can use Bayesian optimization
when we have multiple models,
with different accuracies on different subsystems



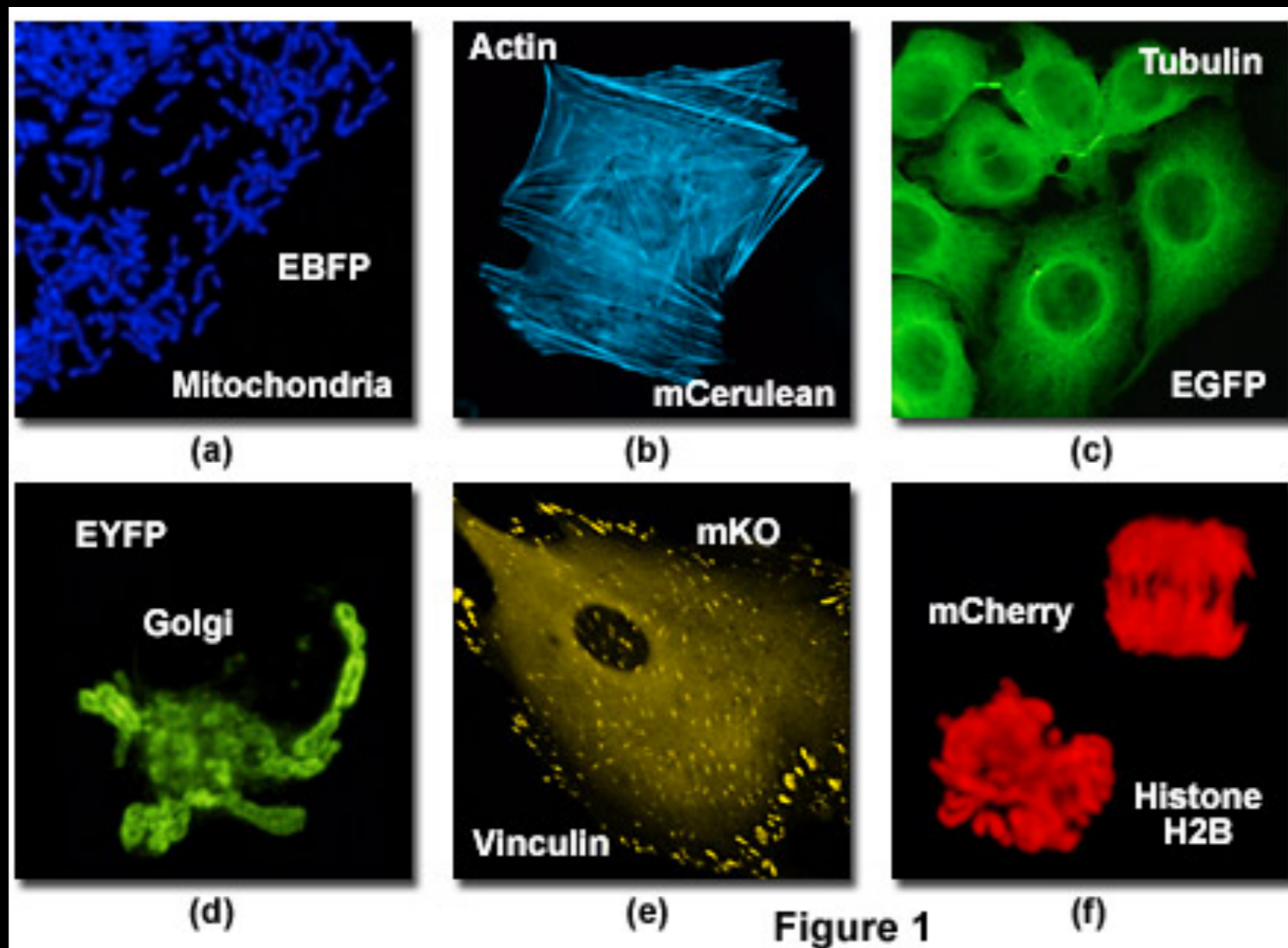
Matthias Poloczek



Jialei Wang



We can use Bayesian optimization when searching over combinatorial spaces



Examples of fluorescent protein tags (not from our collaboration)



Jialei Wang



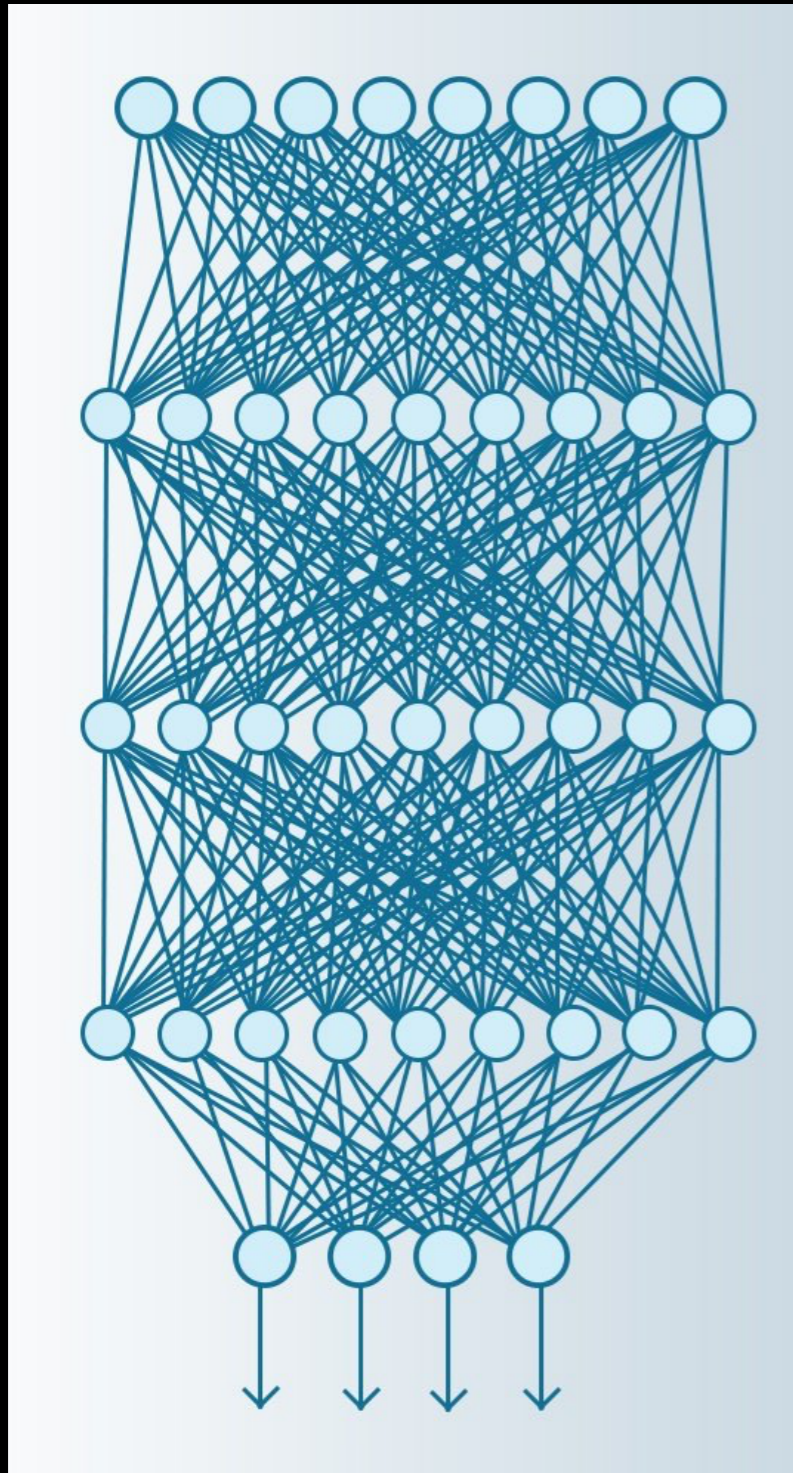
Michael
Burkart



Nathan
Gianneschi

in submission

We can use Bayesian optimization when optimizing integrals & sums



$$\min_{x \in A \subset \mathbb{R}^d} \int F(x, w) p(dw)$$
$$\max_{x \in A \subset \mathbb{R}^d} \frac{1}{n} \sum F(x, w)$$



Saul Toscano-Palmerin