

CSCI 699 - ProbGen

# Probabilistic and Generative Models

Willie Neiswanger

# **Lecture 12 - Predictive UQ, Active Learning, and Bayesian Optimization**

# Today

# Today

**Lecture:** Using predictive UQ models for active learning and sequential decision making, including Bayesian optimization and optimal experimental design.

# Today

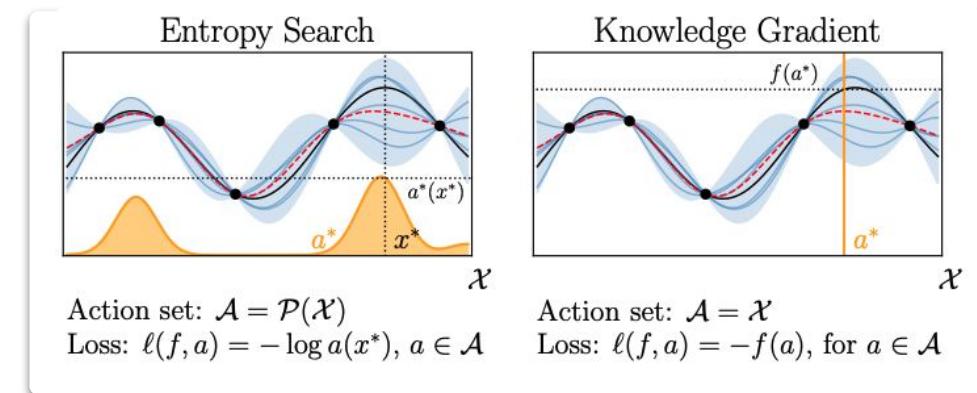
**Lecture:** Using predictive UQ models for active learning and sequential decision making, including Bayesian optimization and optimal experimental design.

- Finish: deep uncertainty quantification.

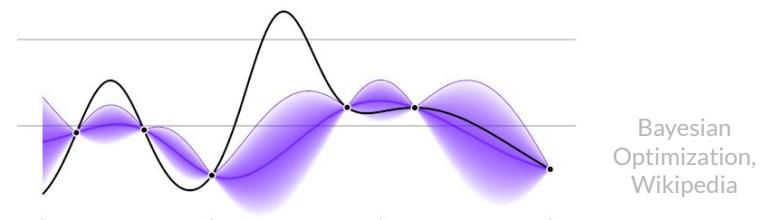
# Today

**Lecture:** Using predictive UQ models for active learning and sequential decision making, including Bayesian optimization and optimal experimental design.

- Finish: deep uncertainty quantification.
- Decision making under uncertainty.
- Active learning.
- Bayesian optimization.
- E.g., UCB, PI, EI, KG, ES, etc.



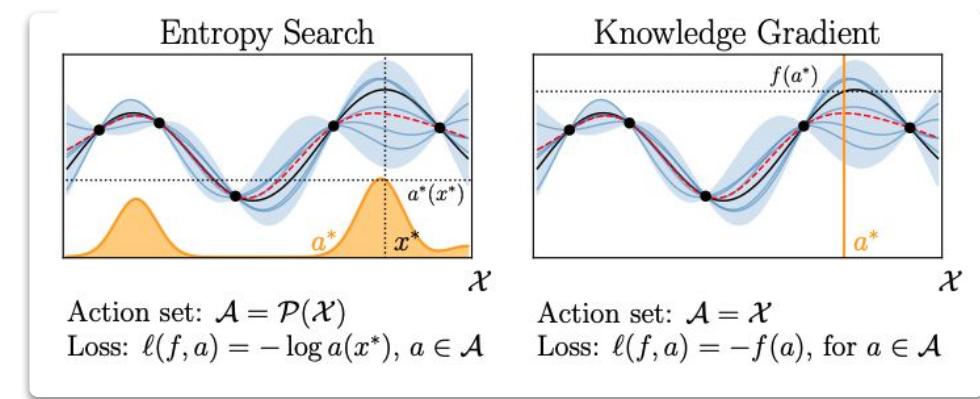
"Generalizing Bayesian Optimization with Decision-theoretic Entropies", 2022



# Today

**Lecture:** Using predictive UQ models for active learning and sequential decision making, including Bayesian optimization and optimal experimental design.

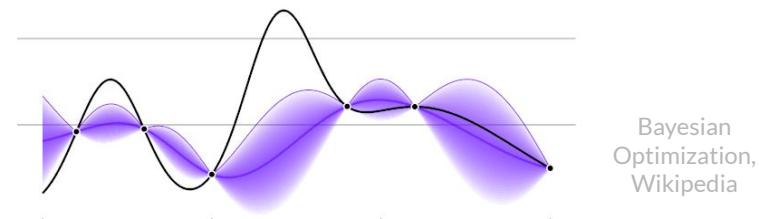
- Finish: deep uncertainty quantification.
- Decision making under uncertainty.
- Active learning.
- Bayesian optimization.
- E.g., UCB, PI, EI, KG, ES, etc.



"Generalizing Bayesian Optimization with Decision-theoretic Entropies", 2022

## After:

- Paper presentations from students.



# Today

After: Next batch of four paper presentations.

**TextGrad: Automatic “Differentiation” via Text**

Mert Yuksekoglu<sup>1\*</sup>, Federico Bianchi<sup>1\*</sup>, Joseph Boen<sup>2</sup>, Sheng Liu<sup>2</sup>, Zhi Hu<sup>2</sup>, Carlos Guestrin<sup>3</sup>, James Zou<sup>1,2,3</sup>

<sup>1</sup>DEPARTMENT OF COMPUTER SCIENCE, STANFORD UNIVERSITY  
<sup>2</sup>DEPARTMENT OF BIOMEDICAL DATA SCIENCE, STANFORD UNIVERSITY  
<sup>3</sup>CHAN ZUCKERBERG BIOHUB

CORRESPONDENCE: MERTY@STANFORD.EDU AND JAMESZ@STANFORD.EDU

[REPOSITORY AND TUTORIALS](#)

**Abstract**

AI is undergoing a paradigm shift, with breakthroughs achieved by systems orchestrating multiple large language models (LLMs) and other complex compounds. As a result, developing principled and automated optimization methods for compound AI systems is one of the most important research challenges in neural networks. In just a few years, we have gone from days until backpropagation and automatic differentiation transformed the field by making optimization turn-key. Inspired by this, we introduce **TEXTGRAD**, a powerful framework performing automatic “differentiation” via text. **TEXTGRAD** backpropagates textual feedback provided by LLMs to improve individual components of a compound AI system. In our framework, LLMs provide rich, general, natural language summaries of the model’s internal state, in the form of graphs, natural language code snippets to manipulate structures. **TEXTGRAD** follows PyTorch’s syntax and abstraction and is flexible and easy-to-use. It works out-of-the-box for a variety of tasks, where the users only provide the objective function without tuning components or prompts of the framework. We showcase **TEXTGRAD**’s effectiveness and generality across a diverse range of applications, including image editing and molecule optimization, and demonstrate its use in agent planning. Without modifying the framework, **TEXTGRAD** improves the zero-shot accuracy of GPT-4o in Google-Prompt Question Answering from 51% to 55%, yields 20% relative performance gain in optimizing LeetCode-Hard coding problem solutions, improves prompts for reasoning, designs new druglike small molecules with desirable *in silico* binding, and designs radiation oncology treatment plans with high specificity. **TEXTGRAD** lays a foundation to accelerate the development of the next generation of AI systems.

**1 Introduction**

There is an emerging paradigm shift in how AI systems are built, owing to the breakthroughs of Large Language Models (LLMs) [1–6]. The new generation of AI applications are increasingly compound systems involving multiple LLMs, which will likely be used in conjunction with other AI tools such as a simulator, or web search. For instance, a system of LLMs communicating with symbolic solvers can solve olympiad-level math problems [7]; a system of LLMs using search engines and code interpreter tools performs comparably to human competitive programmers [8] and are solving real-world

\*Corresponding authors: merty@stanford.edu, yuanzehuan@bytedance.com, † project lead

arXiv:2406.07496v1 [cs.CL] 11 Jun 2024

**Efficient Online Bayesian Inference for Neural Bandits**

Gerardo Duran-Martin<sup>1</sup>, Aleyna Kara<sup>2</sup>, and Kevin Murphy<sup>3</sup>

<sup>1</sup>Queen Mary University, UK  
<sup>2</sup>Bogaziçi University, Turkey  
<sup>3</sup>Google Research, USA

December 2, 2021

**Abstract**

In this paper we present a new algorithm for online (sequential) inference in Bayesian neural networks, and show its superiority for solving constrained multi-armed bandit problems. The key idea is to combine the Kullback-Leibler (which measures the likelihood function) with a learned (or random) low-dimensional affine subspace for the parameters; the use of a subspace enables us to scale our algorithm to problems with  $\sim 1M$  parameters. While most other neural bandit methods need to store the entire history of the problem “online”, our approach only stores the last  $\sim 100$  observations, thus saving constant memory. This is possible because we represent uncertainty about all the parameters in the model, not just the final linear layer. We show good results on the “Deep Bayesian Bandit Showdown” benchmark, as well as MNIST and a recommender system.

**1 Introduction**

Contextual bandit problems (see e.g., [LS19; SH19]) are a special case of reinforcement learning, in which the state (context) at each time step is chosen independently, rather than being dependent on the past history of states and actions. Despite this limitation, contextual bandits are widely used in real-world applications of machine learning systems [Li+10; Guo+20], advertising [MoM+19; Das+21], healthcare [Gos+17; AKR22], etc. The goal is to find the best action  $a_t$  by observing previous actions  $a_i$  in response to each input context or state  $s_i$ . To do this, the decision making agent must learn a reward model  $E[y|s_i, a_i, \theta] = f(s_i, a_i; \theta)$ , where  $\theta$  are the unknown model parameters. Under supervised learning, the agent does not get to see the “correct” output, but instead only gets feedback on whether the choice it made was good or bad (in terms of the reward). If the agent  $a$  could pick the best action  $a^*$  using  $a^* = \arg\max_a f(s_i, a; \theta)$ , however, since  $\theta$  is unknown, the agent must “explore”, so it can gather information about the reward function, before it can “exploit” its model.

In the bandit literature, the two most common solutions to solving the exploration-exploitation dilemma are based on the Upper Confidence Bound (UCB) method (see e.g., [Li+10; KCI12]) and the Thompson Sampling (TS) method (see e.g., [AG13; Li+18]). The key difference between UCB and TS is that while computing the posterior  $p(\theta|D_{t-1})$  is an online fashion, where  $D_{t-1} = \{(s_i, a_i, p_i) : i = 1 : t\}$  is all the data seen so far. This can be done in closed form for linear-Gaussian models, but for nonlinear models, such as deep neural networks (DNNs), it is computationally infeasible.

In this paper, we propose to use the extended Kalman filter to recursively approximate the parameter posterior  $p(\theta|D_t)$  using constant time and memory (i.e., independent of  $T$ ). The main novelty of our approach is that we show how to scale the EKF to large neural networks by leveraging recent results that show that deep neural networks often have very few “degrees of freedom” (see e.g., [Li+18; Im+19; Lar+21]). Thus we can compute a low-dimensional subspace and perform Bayesian filtering in the subspace rather than the original parameter space. We therefore call our method “Bayesian subspace bandits”.

1

Corresponding authors: wangw@pku.edu.cn, yuanzehuan@bytedance.com, † project lead

Preprint. Under review.

**Visual Autoregressive Modeling: Scalable Image Generation via Next-Scale Prediction**

Keyu Tian<sup>1,2</sup>, Yi Jiang<sup>1,2</sup>, Zabuan Yuan<sup>2,3</sup>, Bingyue Peng<sup>2</sup>, Liwei Wang<sup>1,4</sup>

<sup>1</sup>Peking University  
<sup>2</sup>Bytedance Inc.  
<sup>3</sup>yuanzehuan@bytedance.com, bingyue.peng@bytedance.com, wanglw@pku.edu.cn

Try and explore our online demo at: <https://var.vision>  
Codes and models: <https://github.com/FoundationVision/VAR>



Figure 1: Generated samples from Visual Auto-Regressive (VAR) transformer trained on ImageNet. We show 512x512 samples (top), 256x256 samples (middle), and zero-shot image editing results (bottom).

**Abstract**

We present Visual AutoRegressive (VAR), a new generation paradigm that models the visual world via autoregressive “next-token prediction”, diverging from the standard raster-scan “text-to-image” prediction. This simple, intuitive methodology allows autoregressive (AR) models to learn visually conditioned distributions and can generalize well: VAR, for the first time, matches GPT-style AR models’ success in generating images and in image generation. On ImageNet 256x256 benchmark, VAR significantly improve AR baseline by improving Fréchet Inception Distance (FID) from 18.65 to 1.73, inception score (IS) from 10.80 to 18.00, and 20+ failure cases. It is also empirically verified that VAR outperforms the Diffusion Transformer (DT) in multiple dimensions including image quality, inference speed, data efficiency, and scalability. Scaling up VAR models exhibits clear power-law scaling laws similar to those of DTs. VAR also exhibits strong zero-shot generalization ability and cold evidence. VAR further showcases zero-shot generalization ability in downstream tasks including image-in-painting, out-painting, and editing. These results suggest VAR has great empirical potential in improving the performance of LLMs: Scaling Laws and zero-shot generalization. We have released all models and code to promote the exploration of AR/VAR models for visual generation and unified learning.

\*Denote equal contribution  
More results on our project page: [timothysbrooks.com/instruct-pix2pix](https://timothysbrooks.com/instruct-pix2pix)

**InstructPix2Pix: Learning to Follow Image Editing Instructions**

Tim Brooks\* Aleksander Holynski\* Alexei A. Efros  
University of California, Berkeley



Figure 2: Given an image and an instruction for how to edit it, our model performs the appropriate edit. Our model does not require full descriptions for the input or output image, and edits images in the forward pass without per-example inversion or fine-tuning.

**1. Introduction**

We present a method for teaching a generative model to follow human-written instructions for image editing. Since training data for this task is difficult to acquire at scale, we propose an approach for generating a paired dataset that can be used to train a diffusion model. We build on existing models: a large language model GPT-3 [17] and a text-to-image model (Stable Diffusion) [53]. These two models capture complementary knowledge about language and images that can be combined to create paired training data for a task spanning multiple domains.

Using our generated paired data, we train a conditional diffusion model that, given an input image and a text instruction for how to edit it, generates the edited image. Our model directly performs the edits in the forward pass, and does not require any additional edits in the backward pass, or example finetuning. Despite being trained entirely on synthetic examples (i.e., both generated written instructions and generated

1

# Today

After: Next batch of four paper presentations.

**TextGrad: Automatic “Differentiation” via Text**

Mert Yuksekoglu<sup>1\*</sup>, Federico Bianchi<sup>1\*</sup>, Joseph Boen<sup>2</sup>, Sheng Liu<sup>3</sup>, Zhi Hu<sup>3</sup>, Carlos Guestrin<sup>1</sup>, James Zou<sup>1,2,3</sup>

<sup>1</sup>DEPARTMENT OF COMPUTER SCIENCE, STANFORD UNIVERSITY  
<sup>2</sup>DEPARTMENT OF BIOLOGICAL DATA SCIENCE, STANFORD UNIVERSITY  
<sup>3</sup>CHAN ZUCKERBERG BIHUB

CORRESPONDENCE: MERTY@STANFORD.EDU AND JAMESZ@STANFORD.EDU

[REPOSITORY AND TUTORIALS](#)

**Abstract**

AI is undergoing a paradigm shift, with breakthroughs achieved by systems orchestrating multiple large language models (LLMs) and other complex compounds. As a result, developing principled and automated optimization methods for compound AI systems is one of the most important research challenges in the field. In this paper, we show that in just a few days and backpropagation and automatic differentiation transformed the field by making optimization turn-key. Inspired by this, we introduce **TEXTGRAD**, a powerful framework performing automatic “differentiation” via text. **TEXTGRAD** backpropagates textual feedback provided by LLMs to improve individual components of a compound AI system. In our framework, LLMs provide rich, general, natural language summaries of the state-of-the-art in decision graph, neural network, code snippets to manipulate structures. **TEXTGRAD** follows PyTorch’s syntax and abstraction and is flexible and easy-to-use. It works out-of-the-box for a variety of tasks, where the users only provide the objective function without tuning components or prompts of the framework. We showcase **TEXTGRAD**’s effectiveness and generality across a diverse range of applications, including reinforcement learning and molecule optimization, zero-shot image planning, without modifying the framework. **TEXTGRAD** improves the zero-shot accuracy of GPT-4o in Google-Prompt Question Answering from 51% to 55%, yields 20% relative performance gain in optimizing LeetCode-Hard coding problem solutions, improves prompts for reasoning, designs new druglike small molecules with desirable *in silico* binding, and designs radiation oncology treatment plans with high specificity. **TEXTGRAD** lays a foundation to accelerate the development of the next generation of AI systems.

**1 Introduction**

There is an emerging paradigm shift in how AI systems are built, owing to the breakthroughs of Large Language Models (LLMs) [1–6]. The new generation of AI applications are increasingly compound systems involving multiple LLMs, which will be used in conjunction with other AI tools such as a simulator, or web search. For instance, a system of LLMs communicating with symbolic solvers can solve olympiad-level math problems [7]; a system of LLMs using search engines and code interpreter tools performs comparably to human competitive programmers [8] and are solving real-world

<sup>\*</sup>Co-first authors.

arXiv:2406.07496v1 [cs.CL] 11 Jun 2024

1

**Efficient Online Bayesian Inference for Neural Bandits**

Gerardo Duran-Martin<sup>1</sup>, Aleyna Kara<sup>2</sup>, and Kevin Murphy<sup>3</sup>

<sup>1</sup>Queen Mary University, UK  
<sup>2</sup>Bogaziçi University, Turkey  
<sup>3</sup>Google Research, USA

December 2, 2021

**Abstract**

In this paper we present a new algorithm for online (sequential) inference in Bayesian neural networks, and show its superiority for solving constrained multi-armed bandit problems. Our approach (which we term the likelihood function) uses a subspace to (learned or random) low-dimensional affine subspaces for the parameters; the use of a subspace enables us to scale our algorithm to problems with  $\sim 1M$  parameters. While most other neural bandit methods need to store the entire history of the problem “in memory or forget”, our approach only stores the most recent memory. This is possible because we represent uncertainty about all the parameters in the model, not just the final linear layer. We show good results on the “Deep Bayesian Bandit Showdown” benchmark, as well as MNIST and a recommender system.

**1 Introduction**

Contextual bandit problems (see e.g., [LS19; SH19]) are a special case of reinforcement learning, in which the state (context) at each time step is chosen independently, rather than being dependent on the past history of states and actions. Despite this limitation, contextual bandits are widely used in real-world applications such as recommendation systems [Li+10; Guo+20], advertising [Mo+17; Das+21], healthcare [Gos+17; AKR22], etc. The goal is to find the best action  $a_t$  by observing previous actions  $a_i$  in response to each input context or state  $s_i$ . To do this, the decision making agent must learn a reward model  $E[y|s_i, a_i, \theta] = f(s_i, a_i; \theta)$ , where  $\theta$  are the unknown model parameters. Under supervised learning, the agent does not get to see the “correct” output, but instead only gets feedback on whether the choice it made was good or bad (in terms of the reward). If the agent  $a$  would pick the wrong action  $a^*$  using  $a^* = \arg\max_a f(s_i, a; \theta)$ , however, since  $\theta$  is unknown, the agent must “explore”, so it can gather information about the reward function, before it can “exploit” its model.

In the bandit literature, the two most common solutions to solving the exploration-exploitation dilemma are based on the Upper Confidence Bound (UCB) method (see e.g., [Li+10; KCI12]) and the Thompson Sampling (TS) method (see e.g., [AG13; Li+18]). The key idea behind UCB and TS is to directly compute the posterior  $p(a|D_{t-1})$  in an online fashion, where  $D_{t-1} = \{(s_i, a_i, r_i)\}_{i=1}^{t-1}$  is all the data seen so far. This can be done in closed form for linear-Gaussian models, but for nonlinear models, such as deep neural networks (DNNs), it is computationally infeasible.

In this paper, we propose to use LLMs to extend the extended Kalman filter to recursively approximate the parameter posterior  $p(\theta|D_t)$  using constant time and memory (i.e., independent of  $T$ ). The main novelty of our approach is that we show how to scale the EKF to large neural networks by leveraging recent results that show that deep neural networks often have very few “degrees of freedom” (see e.g., [Li+18; Izm+19; Lar+21]). Thus we can compute a low-dimensional subspace and perform Bayesian filtering in the subspace rather than the original parameter space. We therefore call our method “Bayesian subspace bandits”.

Corresponding authors: wangw@pku.edu.cn, yuanzhan@bytedance.com, † project lead

Printed: Under review.

arXiv:2112.00195v1 [cs.LG] 1 Dec 2021

1

**Visual Autoregressive Modeling: Scalable Image Generation via Next-Scale Prediction**

Keyu Tian<sup>1,2</sup>, Yi Jiang<sup>3,4</sup>, Zabuan Yuan<sup>2,5</sup>, Bingyue Peng<sup>2</sup>, Liwei Wang<sup>1,4</sup>

<sup>1</sup>Peking University  
<sup>2</sup>Bytedance Inc.  
<sup>3</sup>keytian@stu.pku.edu.cn, jiangyi\_enjoy@bytedance.com, yuanzhanhua@bytedance.com, bingyue\_peng@bytedance.com, wangliw@pku.edu.cn

Try and explore our online demo at: <https://var.vision>  
Codes and models: <https://github.com/FoundationVision/VAR>



Figure 1: Generated samples from Visual Auto-Regressive (VAR) transformer trained on ImageNet. We show 512x512 samples (top), 256x256 samples (middle), and zero-shot image editing results (bottom).

**Abstract**

We present Visual AutoRegressive (VAR), a new generation paradigm that models the visual world via autoregressive “next-token prediction” or “next-resolution prediction”, diverging from the standard raster-scan “text-to-image” prediction. This simple, intuitive methodology allows autoregressive (AR) transformers to learn from raw images and can generalize well: VAR, for the first time, matches GPT-style AR models’ sample efficiency and a state-of-the-art image model (Stable Diffusion) – to generate a large dataset of image editing examples. Our conditional diffusion model, InstructPix2Pix, is trained on our generated data, and generalizes to real images and user-written instructions at inference time. The proposed model does not require per-example fine-tuning or inversion, our model edits images quickly, in a matter of seconds. We show compelling editing results for a diverse collection of input images and written instructions.

<sup>1</sup>Donors equal contribution  
More results on our project page: [timothysbrooks.com/instruct-pix2pix](https://timothysbrooks.com/instruct-pix2pix)

arXiv:2404.02905v2 [cs.CV] 10 Jun 2024

1

**InstructPix2Pix: Learning to Follow Image Editing Instructions**

Tim Brooks\* Aleksander Holynski\* Alexei A. Efros  
University of California, Berkeley



Figure 1: Given an image and an instruction for how to edit it, our model performs the appropriate edit. Our model does not require full descriptions for the input or output image, and edits images in the forward pass without per-example inversion or fine-tuning.

**1. Introduction**

We present a method for teaching a generative model to follow human-written instructions for image editing. Since training data for this task is difficult to acquire at scale, we propose an approach for generating a paired dataset that tells the model what to do: our model follows these instructions to edit the image. To obtain training data for this problem, we combine the knowledge of image generation models (e.g., DALL-E [Ch+21] and a state-of-the-art image model (Stable Diffusion) [SD23]). These two models capture complementary knowledge about language and images that can be combined to create paired training data for a task spanning multiple domains.

Using our generated paired data, we train a conditional diffusion model that, given an input image and a text instruction for how to edit it, generates the edited image. Our model directly performs the requested edit in the forward pass, and does not require any additional edit-in-the-backpass or description of the input/output images, or example fine-tuning. Despite being trained entirely on synthetic examples (i.e., both generated written instructions and generated

arXiv:2211.09890v2 [cs.CV] 18 Jan 2023

1

# Today

After: Next batch of four paper presentations.

**TextGrad: Automatic “Differentiation” via Text**

Mert Yuksekoglu<sup>1\*</sup>, Federico Bianchi<sup>1\*</sup>, Joseph Boen<sup>2</sup>, Sheng Liu<sup>2</sup>, Zhi Hu<sup>2</sup>, Carlos Guestrin<sup>3</sup>, James Zou<sup>1,2,3</sup>

<sup>1</sup>DEPARTMENT OF COMPUTER SCIENCE, STANFORD UNIVERSITY  
<sup>2</sup>DEPARTMENT OF BIOMEDICAL DATA SCIENCE, STANFORD UNIVERSITY  
<sup>3</sup>CHAN ZUCKERBERG BIOHUB

CORRESPONDENCE: MERTY@STANFORD.EDU AND JAMESZ@STANFORD.EDU

[REPOSITORY AND TUTORIALS](#)

**Abstract**

AI is undergoing a paradigm shift, with breakthroughs achieved by systems orchestrating multiple large language models (LLMs) and other complex compounds. As a result, developing principled and automated optimization methods for compound AI systems is one of the most important research challenges in the field. In this paper, we show that in just a few days and with backpropagation and automatic differentiation transformed the field by making optimization turn-key. Inspired by this, we introduce **TEXTGRAD**, a powerful framework performing automatic “differentiation” via text. **TEXTGRAD** backpropagates textual feedback provided by LLMs to improve individual components of a compound AI system. In our framework, LLMs provide rich, general, natural language summaries of the problem in question, graphically, naturally, and code respects its underlying structures. **TEXTGRAD** follows PyTorch’s syntax and abstraction and is flexible and easy-to-use. It works out-of-the-box for a variety of tasks, where the users only provide the objective function without tuning components or prompts of the framework. We showcase **TEXTGRAD**’s effectiveness and generality across a diverse range of applications, including reinforcement learning and molecule optimization, and demonstrate its use in planning. Without modifying the framework, **TEXTGRAD** improves the zero-shot accuracy of GPT-4o in Google-Proof Question Answering from 51% to 55%, yields 20% relative performance gain in optimizing LeetCode-Hard coding problem solutions, improves prompts for reasoning, designs new druglike small molecules with desirable *in silico* binding, and designs radiation oncology treatment plans with high specificity. **TEXTGRAD** lays a foundation to accelerate the development of the next generation of AI systems.

**1 Introduction**

There is an emerging paradigm shift in how AI systems are built, owing to the breakthroughs of Large Language Models (LLMs) [1–6]. The new generation of AI applications are increasingly compound systems involving multiple LLMs, which can be combined in various ways to build complex AI systems such as a simulator, or web search. For instance, a system of LLMs communicating with symbolic solvers can solve olympiad-level math problems [7]; a system of LLMs using search engines and code interpreter tools performs comparably to human competitive programmers [8] and are solving real-world

\*Co-first authors.

arXiv:2406.07496v1 [cs.CL] 11 Jun 2024

**Efficient Online Bayesian Inference for Neural Bandits**

Gerardo Duran-Martin<sup>1</sup>, Aleyna Kara<sup>2</sup>, and Kevin Murphy<sup>3</sup>

<sup>1</sup>Queen Mary University, UK  
<sup>2</sup>Bogaziçi University, Turkey  
<sup>3</sup>Google Research, USA

December 2, 2021

**Abstract**

In this paper we present a new algorithm for online (sequential) inference in Bayesian neural networks, and show its superiority for solving contextual bandit problems. The key idea is to combine the Kullback-Leibler (which measures the likelihood function) with a learned (or random) low-dimensional affine subspace for the parameters; the use of a subspace enables us to scale our algorithm to problems with  $\sim 1M$  parameters. While most other neural bandit methods need to store the entire history of the problem in “context memory”, our approach only needs to store a small amount of memory. This is possible because we represent uncertainty about all the parameters in the model, not just the final linear layer. We show good results on the “Deep Bayesian Bandit Showdown” benchmark, as well as MNIST and a recommender system.

**1 Introduction**

Contextual bandit problems (see e.g., [LS19; SH19]) are a special case of reinforcement learning, in which the state (context) at each time step is chosen independently, rather than being dependent on the past history of states and actions. Despite this limitation, contextual bandits are widely used in real-world applications such as recommendation systems [Li+10; Guo+20], advertising [Mo+17; Das+21], healthcare [Gos+17; AKR22], etc. The goal is to learn a policy that maximizes the reward by choosing  $a_t$  in response to each input context or state  $s_t$ . To do this, the decision making agent must learn a reward model  $E[y_t|s_t, a_t, \theta] = f(s_t, a_t; \theta)$ , where  $\theta$  are the unknown model parameters. Unlike supervised learning, the agent does not get to see the “correct” output, but instead only gets feedback on whether the choice it made was good or bad (in terms of the reward). If the agent  $a_t$  would pick the wrong action using  $a_t = \arg\max_a f(s_t, a; \theta)$ , however, since  $\theta$  is unknown, the agent must “explore”, so it can gather information about the reward function, before it can “exploit” its model.

In the bandit literature, the two most common solutions to solving the exploration-exploitation dilemma are based on the Upper Confidence Bound (UCB) method (see e.g., [Li+10; KCI12]) and the Thompson Sampling (TS) method (see e.g., [AG13; Li+18]). The key difference between UCB and TS is that while UCB is primarily computing the posterior  $p(a|D_{t-1})$  in an online fashion, where  $D_{t-1} = \{(s_i, a_i, r_i)\}_{i=1}^{t-1}$  is all the data seen so far. This can be done in closed form for linear-Gaussian models, but for nonlinear models, such as deep neural networks (DNNs), it is computationally infeasible.

In this paper, we propose to use the extended Kalman filter to recursively approximate the parameter posterior  $p(\theta|D_t)$  under constant time and memory (i.e., independent of  $T$ ). The main novelty of our approach is that we show how to scale the EKF to large neural networks by leveraging recent results that show that deep neural networks often have very few “degrees of freedom” (see e.g., [Li+18; Im+19; Lar+21]). Thus we can compute a low-dimensional subspace and perform Bayesian filtering in the subspace rather than the original parameter space. We therefore call our method “Bayesian subspace bandits”.

1

Corresponding authors: wangw@pku.edu.cn, yuanzhan@bytedance.com, † project lead

Preprint. Under review.

arXiv:2112.00195v1 [cs.LG] 1 Dec 2021

**Visual Autoregressive Modeling: Scalable Image Generation via Next-Scale Prediction**

Keyu Tian<sup>1,2</sup>, Yi Jiang<sup>1,2</sup>, Zabuan Yuan<sup>2,3</sup>, Bingyue Peng<sup>2</sup>, Liwei Wang<sup>1,4</sup>

<sup>1</sup>Peking University  
<sup>2</sup>Bytedance Inc.  
<sup>3</sup>yuanzhanhua@bytedance.com, bingyue.peng@bytedance.com, wangliw@pku.edu.cn

Try and explore our online demo at: <https://var.vision>  
Codes and models: <https://github.com/FoundationVision/VAR>



Figure 1: Generated samples from Visual Auto-Regressive (VAR) transformer trained on ImageNet. We show 512x512 samples (top), 256x256 samples (middle), and zero-shot image editing results (bottom).

**Abstract**

We present Visual AutoRegressive (VAR), a new generation paradigm that models the visual world with a sequence of autoregressive “next-token prediction”, diverging from the standard raster-scan “next-tokens prediction”. This simple, intuitive methodology allows autoregressive (AR) transformers to learn visually-aware representations and can generalize well: VAR, for the first time, matches GPT-style AR models’ superior performance in image generation. On ImageNet 256x256 benchmark, VAR significantly improve AR baseline by improving Fréchet Inception Distance (FID) from 18.65 to 1.73, inception score (IS) from 18.03 to 20.8, and a 20x speedup. It is also empirically verified that VAR outperforms the Diffusion Transformer (DT) in multiple dimensions including image quality, inference speed, data efficiency, and scalability. Scaling up VAR models exhibits clear power-law scaling laws similar to those of DT. VAR also exhibits strong generalization ability with limited visual evidence. VAR further showcases zero-shot generalization ability in downstream tasks including image-in-painting, out-painting, and editing. These results suggest VAR has great empirical potential in improving the performance of LLMs. Scientific Laws and research generalization. We have released all models and code to promote the exploration of AR/VAE models for visual generation and unified learning.

\*Denote equal contribution.  
More results on our project page: [timothysbrooks.com/instruct-pix2pix](https://timothysbrooks.com/instruct-pix2pix)

arXiv:2404.02905v2 [cs.CV] 10 Jun 2024

**InstructPix2Pix: Learning to Follow Image Editing Instructions**

Tim Brooks\*, Aleksander Holynski\*, Alexei A. Efros  
University of California, Berkeley

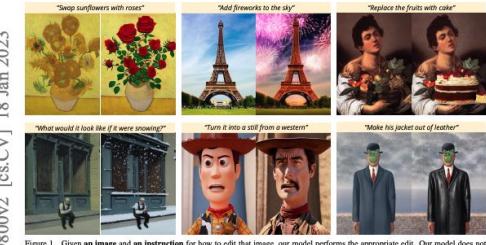


Figure 1: Given an image and an instruction for how to edit it, our model performs the appropriate edit. Our model does not require full descriptions for the input or output image, and edits images in the forward pass without per-example inversion or fine-tuning.

**1. Introduction**

We present a method for teaching a generative model to follow human-written instructions for image editing. Since training data for this task is difficult to acquire at scale, we propose an approach for generating a paired dataset that tells the model what to do, our model follows these instructions to edit the image. To obtain training data for this problem, we combine the knowledge of image generation models (e.g., DALL-E [Ch+21] and a state-of-the-art image model (Stable Diffusion) [SD23]). These two models capture complementary knowledge about language and images that can be combined to create paired training data for a task spanning both domains.

Using our generated paired data, we train a conditional diffusion model that, given an input image and a text instruction for how to edit it, generates the edited image. Our model directly performs the requested edit in the forward pass, and does not require any additional computation or fine-tuning. Despite being trained entirely on synthetic examples (i.e., both generated written instructions and generated

arXiv:2211.09890v2 [cs.CV] 18 Jan 2023

# Today

After: Next batch of four paper presentations.

**TextGrad: Automatic “Differentiation” via Text**

Mert Yuksekoglu<sup>1\*</sup>, Federico Bianchi<sup>1\*</sup>, Joseph Boen<sup>2</sup>, Sheng Liu<sup>3</sup>, Zhi Hu<sup>3</sup>, Carlos Guestrin<sup>1</sup>, James Zou<sup>1,2,3</sup>

<sup>1</sup>DEPARTMENT OF COMPUTER SCIENCE, STANFORD UNIVERSITY  
<sup>2</sup>DEPARTMENT OF BIOMEDICAL DATA SCIENCE, STANFORD UNIVERSITY  
<sup>3</sup>CHAN ZUCKERBERG BIOHUB

CORRESPONDENCE: MERTY@STANFORD.EDU AND JAMESZ@STANFORD.EDU

[REPOSITORY AND TUTORIALS](#)

**Abstract**

AI is undergoing a paradigm shift, with breakthroughs achieved by systems orchestrating multiple large language models (LLMs) and other complex compounds. As a result, developing principled and automated optimization methods for compound AI systems is one of the most important research challenges in neural networks. In just a few years, we have gone from days until backpropagation and automatic differentiation transformed the field by making optimization turn-key. Inspired by this, we introduce **TEXTGRAD**, a powerful framework performing automatic “differentiation” via text. **TEXTGRAD** backpropagates textual feedback provided by LLMs to improve individual components of a compound AI system. In our framework, LLMs provide rich, general, natural language summaries of the internal state of the system, graphically, naturally, and code respects its modular structures. **TEXTGRAD** follows PyTorch’s syntax and abstraction and is flexible and easy-to-use. It works out-of-the-box for a variety of tasks, where the users only provide the objective function without tuning components or prompts of the framework. We showcase **TEXTGRAD**’s effectiveness and generality across a diverse range of applications, including image editing and molecule optimization, and demonstrate its use in agent planning. Without modifying the framework, **TEXTGRAD** improves the zero-shot accuracy of GPT-4o in Google-Poof Question Answering from 51% to 55%, yields 20% relative performance gain in optimizing LeetCode-Hard coding problem solutions, improves prompts for reasoning, designs new druglike small molecules with desirable *in silico* binding, and designs radiation oncology treatment plans with high specificity. **TEXTGRAD** lays a foundation to accelerate the development of the next generation of AI systems.

**1 Introduction**

Contextual bandit problems (see e.g., [LS19; SH19]) are a special case of reinforcement learning, in which the state (context) at each time step is chosen independently, rather than being dependent on the past history of states and actions. Despite this limitation, contextual bandits are widely used in real-world applications such as recommendation systems [Li+10; Guo+20], advertising [Mo+17; Du+21], healthcare [Gao+17; AKR21], etc. The goal is to find the best action  $a_t$  by observing previous actions  $a_i$  in response to each input context or state  $s_i$ . To do this, the decision making agent must learn a reward model  $E[y|s_i, a_i, \theta] = f(s_i, a_i; \theta)$ , where  $\theta$  are the unknown model parameters. Under supervised learning, the agent does not get to see the “correct” output, but instead only gets feedback on whether the choice it made was good or bad (in terms of the reward). If the agent  $a$  would pick the wrong action  $a^*$  using  $a^* = \arg\max_a f(s_i, a; \theta)$ , however, since  $\theta$  is unknown, the agent must “explore”, so it can gather information about the reward function, before it can “exploit” its model.

In the bandit literature, the two most common solutions to solving the exploration-exploitation dilemma are based on the Upper Confidence Bound (UCB) method (see e.g., [Li+10; KCI12]) and the Thompson Sampling (TS) method (see e.g., [AG13; Li+18]). The key idea behind UCB and TS is efficiently computing the posterior  $p(\theta|D_{t-1})$  in an online fashion, where  $D_{t-1} = \{(s_i, a_i, y_i) : i = 1 : t\}$  is all the data seen so far. This can be done in closed form for linear-Gaussian models, but for nonlinear models, such as deep neural networks (DNNs), it is computationally infeasible.

In this paper, we propose to use **TEXTGRAD**, the extended Kalman filter to recursively approximate the parameter posterior  $p(\theta|D_t)$  using constant time and memory (i.e., independent of  $T$ ). The main novelty of our approach is that we show how to scale the EKF to large neural networks by leveraging recent results that show that deep neural networks often have very few “degrees of freedom” (see e.g., [Li+18; Izm+19; Lar+21]). Thus we can compute a low-dimensional subspace and perform Bayesian filtering in the subspace rather than the original parameter space. We therefore call our method “Bayesian subspace bandits”.

\*Co-first authors.

arXiv:2406.07496v1 [cs.CL] 11 Jun 2024

**Efficient Online Bayesian Inference for Neural Bandits**

Gerardo Duran-Martin<sup>1</sup>, Aleyna Kara<sup>2</sup>, and Kevin Murphy<sup>3</sup>

<sup>1</sup>Queen Mary University, UK  
<sup>2</sup>Bogaziçi University, Turkey  
<sup>3</sup>Google Research, USA

December 2, 2021

**Abstract**

In this paper we present a new algorithm for online (sequential) inference in Bayesian neural networks, and show its superiority for solving contextual bandit problems. The basic idea is to combine the standard Kalman filter (which is linear in the likelihood function) with a learned (learned or random) low-dimensional affine subspace for the parameters; the use of a subspace enables us to scale our algorithm to problems with  $\sim 1M$  parameters. While most other neural bandit methods need to store the entire history of the problem (“contextual forgetting”), our approach only stores the most recent memory. This is possible because we represent uncertainty about all the parameters in the model, not just the final linear layer. We show good results on the “Deep Bayesian Bandit Showdown” benchmark, as well as MNIST and a recommender system.

**1 Introduction**

Contextual bandit problems (see e.g., [LS19; SH19]) are a special case of reinforcement learning, in which the state (context) at each time step is chosen independently, rather than being dependent on the past history of states and actions. Despite this limitation, contextual bandits are widely used in real-world applications such as recommendation systems [Li+10; Guo+20], advertising [Mo+17; Du+21], healthcare [Gao+17; AKR21], etc. The goal is to find the best action  $a_t$  by observing previous actions  $a_i$  in response to each input context or state  $s_i$ . To do this, the decision making agent must learn a reward model  $E[y|s_i, a_i, \theta] = f(s_i, a_i; \theta)$ , where  $\theta$  are the unknown model parameters. Under supervised learning, the agent does not get to see the “correct” output, but instead only gets feedback on whether the choice it made was good or bad (in terms of the reward). If the agent  $a$  would pick the wrong action  $a^*$  using  $a^* = \arg\max_a f(s_i, a; \theta)$ , however, since  $\theta$  is unknown, the agent must “explore”, so it can gather information about the reward function, before it can “exploit” its model.

In the bandit literature, the two most common solutions to solving the exploration-exploitation dilemma are based on the Upper Confidence Bound (UCB) method (see e.g., [Li+10; KCI12]) and the Thompson Sampling (TS) method (see e.g., [AG13; Li+18]). The key idea behind UCB and TS is efficiently computing the posterior  $p(\theta|D_{t-1})$  in an online fashion, where  $D_{t-1} = \{(s_i, a_i, y_i) : i = 1 : t\}$  is all the data seen so far. This can be done in closed form for linear-Gaussian models, but for nonlinear models, such as deep neural networks (DNNs), it is computationally infeasible.

In this paper, we propose to use **TEXTGRAD**, the extended Kalman filter to recursively approximate the parameter posterior  $p(\theta|D_t)$  using constant time and memory (i.e., independent of  $T$ ). The main novelty of our approach is that we show how to scale the EKF to large neural networks by leveraging recent results that show that deep neural networks often have very few “degrees of freedom” (see e.g., [Li+18; Izm+19; Lar+21]). Thus we can compute a low-dimensional subspace and perform Bayesian filtering in the subspace rather than the original parameter space. We therefore call our method “Bayesian subspace bandits”.

Corresponding authors: wangli@pku.edu.cn, yuanzehua@bytedance.com, † project lead

Preprint. Under review.

arXiv:2112.00195v1 [cs.LG] 1 Dec 2021

**Visual Autoregressive Modeling: Scalable Image Generation via Next-Scale Prediction**

Keyu Tian<sup>1,2</sup>, Yi Jiang<sup>2,†</sup>, Zabuan Yuan<sup>2,§</sup>, Bingyue Peng<sup>2</sup>, Liwei Wang<sup>1,¶</sup>

<sup>1</sup>Peking University  
<sup>2</sup>Bytedance Inc.  
keytianetu@pku.edu.cn, jiangyi\_enjoy@bytedance.com, yuanzehua@bytedance.com, bingyue peng@bytedance.com, wangli@pku.edu.cn

Try and explore our online demo at: <https://var.vision>  
Codes and models: <https://github.com/FoundationVision/VAR>



Figure 1: Generated samples from Visual Auto-Regressive (VAR) transformer trained on ImageNet. We show 512x512 samples (top), 256x256 samples (middle), and zero-shot image editing results (bottom).

**Abstract**

We present Visual AutoRegressive (VAR), a new generation paradigm that models the visual world via autoregressive modeling. Unlike existing methods for “text-to-image prediction” or “next-resolution prediction”, diverging from the standard raster-scan “next-token prediction”. This simple, intuitive methodology allows autoregressive (AR) transformers to learn visually-aware representations and can generalize well: VAR, for the first time, matches GPT-style AR models significantly in terms of image generation. On ImageNet 256x256 benchmark, VAR significantly improve AR baseline by improving Fréchet inception distance (FID) from 18.65 to 1.73, inception score (IS) from 6.80 to 8.03, and 20+ failure cases. In addition, it is also empirically verified that VAR outperforms the Diffusion Transformer (DT) in multiple dimensions including image quality, inference speed, data efficiency, and scalability. Scaling up VAR models exhibits clear power-law scaling laws similar to those of DT. VAR also exhibits clear linear scaling laws with respect to added evidence. VAR further showcases zero-shot generalization ability in downstream tasks including image-in-painting, out-painting, and editing. These results suggest VAR has great empirical potential in improving the performance of LLMs: Scaling Laws and zero-shot generalization. We have released all models and code to promote the exploration of AR/VAE models for visual generation and unified learning.

†Denotes equal contribution.  
More results on our project page: [timothysbrooks.com/instruct-pix2pix](https://timothysbrooks.com/instruct-pix2pix)

arXiv:2404.02905v2 [cs.CV] 10 Jun 2024

**InstructPix2Pix: Learning to Follow Image Editing Instructions**

Tim Brooks\*, Aleksander Holynski\*, Alexei A. Efros  
University of California, Berkeley

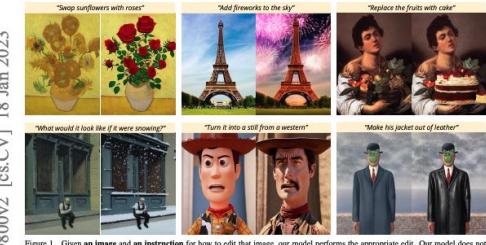


Figure 2: Given an image and an instruction for how to edit it, our model performs the appropriate edit. Our model does not require full descriptions for the input or output image, and edits images in the forward pass without per-example inversion or fine-tuning.

**1. Introduction**

We propose a method for editing images from human instructions, given an input image and a written instruction that tells the model what to do. Our model follows these instructions to edit the image. To obtain training data for this problem, we combine the knowledge of image generation models (such as DALL-E 2 [Chen+22]) and a text-image model (such as GPT-3 [Chowdhery+22]) to generate a large dataset of image editing examples. Our conditional diffusion model, InstructPix2Pix, is trained on our generated data, and generalizes to real images and user-written instructions at inference time. The key idea is that our model follows the instruction, does not require per-example fine-tuning or inversion, our model edits images quickly, in a matter of seconds. We show compelling editing results for a diverse collection of input images and written instructions.

1 Co-first authors.

Using our generated paired data, we train a conditional diffusion model that, given an input image and a text instruction for how to edit it, generates the edited image. Our model directly performs the edit in the forward pass, and does not require any additional edits, per-example inversion, or per-example fine-tuning. Despite being trained entirely on synthetic examples (i.e., both generated written instructions and generated

arXiv:2211.09890v2 [cs.CV] 18 Jan 2023

# Today

After: Next batch of four paper presentations.

**TextGrad: Automatic “Differentiation” via Text**

Mert Yuksekoglu<sup>1\*</sup>, Federico Bianchi<sup>1\*</sup>, Joseph Boen<sup>2</sup>, Sheng Liu<sup>3</sup>, Zhi Hu<sup>3</sup>, Carlos Guestrin<sup>1</sup>, James Zou<sup>1,2,3</sup>

<sup>1</sup>DEPARTMENT OF COMPUTER SCIENCE, STANFORD UNIVERSITY  
<sup>2</sup>DEPARTMENT OF BIOMEDICAL DATA SCIENCE, STANFORD UNIVERSITY  
<sup>3</sup>CHAN ZUCKERBERG BIOHUB

CORRESPONDENCE: MERTY@STANFORD.EDU AND JAMESZ@STANFORD.EDU

[REPOSITORY AND TUTORIALS](#)

**Abstract**

AI is undergoing a paradigm shift, with breakthroughs achieved by systems orchestrating multiple large language models (LLMs) and other complex compounds. As a result, developing principled and automated optimization methods for compound AI systems is one of the most important research challenges in neural networks. In just a few years, we have gone from days until backpropagation and automatic differentiation transformed the field by making optimization turn-key. Inspired by this, we introduce **TEXTGRAD**, a powerful framework performing automatic “differentiation” via text. **TEXTGRAD** backpropagates textual feedback provided by LLMs to improve individual components of a compound AI system. In our framework, LLMs provide rich, general language support for differentiable optimization in the form of graph, text, and code concepts to model complex structures. **TEXTGRAD** follows PyTorch’s syntax and abstraction and is flexible and easy-to-use. It works out-of-the-box for a variety of tasks, where the users only provide the objective function without tuning components or prompts of the framework. We showcase **TEXTGRAD**’s effectiveness and generality across a diverse range of applications, including editing and codebase optimization, zero-shot reinforcement planning, without modifying the framework. **TEXTGRAD** improves the zero-shot accuracy of GPT-4o in Google-Poof Question Answering from 51% to 55%, yields 20% relative performance gain in optimizing LeetCode-Hard coding problem solutions, improves prompts for reasoning, designs new druglike small molecules with desirable *in silico* binding, and designs radiation oncology treatment plans with high specificity. **TEXTGRAD** lays a foundation to accelerate the development of the next generation of AI systems.

**1 Introduction**

There is an emerging paradigm shift in how AI systems are built, owing to the breakthroughs of Large Language Models (LLMs) [1–6]. The new generation of AI applications are increasingly compound systems involving multiple LLMs, which will likely be used in conjunction with other AI tools such as a simulator, or web search. For instance, a system of LLMs communicating with symbolic solvers can solve olympiad-level math problems [7]; a system of LLMs using search engines and code interpreter tools performs comparably to human competitive programmers [8] and are solving real-world

\*Corresponding authors.

arXiv:2406.07496v1 [cs.CL] 11 Jun 2024

**Efficient Online Bayesian Inference for Neural Bandits**

Gerardo Duran-Martin<sup>1</sup>, Aleyna Kara<sup>2</sup>, and Kevin Murphy<sup>3</sup>

<sup>1</sup>Queen Mary University, UK  
<sup>2</sup>Bogaziçi University, Turkey  
<sup>3</sup>Google Research, USA

December 2, 2021

**Abstract**

In this paper we present a new algorithm for online (sequential) inference in Bayesian neural networks, and show its utility for solving contextual bandit problems. The key idea is to combine the Kalman filter (which provides the likelihood function) with a learned (or random) low-dimensional affine subspace for the parameters; the use of a subspace enables us to scale our algorithm to problems with  $\sim 1M$  parameters. While most other neural bandit methods need to store the entire history of the problem (“contextual forgetting”), our approach only stores a constant amount of memory. This is possible because we represent uncertainty about all the parameters in the model, not just the final linear layer. We show good results on the “Deep Bayesian Bandit Showdown” benchmark, as well as MNIST and a recommender system.

**1 Introduction**

Contextual bandit problems (see e.g., [LS19; SH19]) are a special case of reinforcement learning, in which the state (context) at each time step is chosen independently, rather than being dependent on the past history of states and actions. Despite this limitation, contextual bandits are widely used in real-world applications of machine learning systems [Li+10; Guo+20], advertising [MoM+19; Das+21], healthcare [Gos+17; AKR22], etc. The goal is to learn a policy that maximizes the reward by choosing  $a_t$  in response to each input context or state  $s_t$ . To do this, the decision making agent must learn a reward model  $E[y_t|s_t, a_t, \theta] = f(s_t, a_t; \theta)$ , where  $\theta$  are the unknown model parameters. Under supervised learning, the agent does not get to see the “correct” output, but instead only gets feedback on whether the choice it made was good or bad (in terms of the reward). If the agent  $a_t$  would pick the wrong action using  $a_t = \arg\max_a f(s_t, a; \theta)$ , however, since  $\theta$  is unknown, the agent must “explore”, so it can gather information about the reward function, before it can “exploit” its model.

In the bandit literature, the two most common solutions to solving the exploration-exploitation dilemma are based on the ε-greedy (EG) [19] method (see e.g., [Li+10; KCI12]) and the Thompson Sampling (TS) method (see e.g., [AG13; Li+18]). The key difference between EG and TS is that EG sequentially computes the posterior  $p(\theta|D_{t-1})$  in an online fashion, where  $D_{t-1} = \{(s_i, a_i, r_i) : i = 1 : t\}$  is all the data seen so far. This can be done in closed form for linear-Gaussian models, but for nonlinear models, such as deep neural networks (DNNs), it is computationally infeasible.

In this paper, we propose to use the extended Kalman filter to recursively approximate the parameter posterior  $p(\theta|D_t)$  over constant time and memory (i.e., independent of  $T$ ). The main novelty of our approach is that we show how to scale the EKF to large neural networks by leveraging recent results that show that deep neural networks often have very few “degrees of freedom” (see e.g., [Li+18; Im+19; Lar+21]). Thus we can compute a low-dimensional subspace and perform Bayesian filtering in the subspace rather than the original parameter space. We therefore call our method “Bayesian subspace bandits”.

1

Corresponding authors: wangw@pku.edu.cn, yuanzehuan@bytedance.com, † project lead

Preprint. Under review.

**Visual Autoregressive Modeling: Scalable Image Generation via Next-Scale Prediction**

Keyu Tian<sup>1,2</sup>, Yi Jiang<sup>1,2</sup>, Zabuan Yuan<sup>2,3\*</sup>, Bingyue Peng<sup>2</sup>, Liwei Wang<sup>1,4</sup>

<sup>1</sup>Peking University  
<sup>2</sup>Bytedance Inc.  
<sup>3</sup>yuanzehuan@bytedance.com, bingyue.peng@bytedance.com, wangliw@pku.edu.cn

Try and explore our online demo at: <https://var.vision>  
Codes and models: <https://github.com/FoundationVision/VAR>



Figure 1: Generated samples from Visual Auto-Regressive (VAR) transformer trained on ImageNet. We show 512x512 samples (top), 256x256 samples (middle), and zero-shot image editing results (bottom).

**Abstract**

We present Visual AutoRegressive (VAR), a new generation paradigm that models the visual world via autoregressive prediction, allowing the model to “edit” images via “next-token prediction”. This simple, intuitive methodology allows autoregressive (AR) transformers to learn image generation and can generalize well: VAR, for the first time, matches GPT-style AR models’ success in generating images on ImageNet 256x256 benchmark. VAR significantly improve FID score by improving Fréchet Inception Distance (FID) from 18.65 to 1.73, inception score (IS) from 18.03 to 20.8, and aperceptual score (AP) from 18.03 to 20.8. It is also empirically verified that VAR outperforms the Diffusion Transformer (DT) in multiple dimensions including image quality, inference speed, data efficiency, and scalability. Scaling up VAR models exhibits clear power-law scaling laws similar to those of DTs. VAR also exhibits strong zero-shot generalization ability and cold evidence. VAR further showcases zero-shot generalization ability in downstream tasks including image-in-painting, out-painting, and editing. These results suggest VAR has great empirical potential in improving the performance of LLMs. See the Laws and research generalization. We have released all models and code to promote the exploration of AR/VAE models for visual generation and unified learning.

\*Denotes equal contribution.  
More results on our project page: [timothysbrooks.com/instruct-pix2pix](https://timothysbrooks.com/instruct-pix2pix)

**InstructPix2Pix: Learning to Follow Image Editing Instructions**

Tim Brooks\*, Aleksander Holynski\*, Alexei A. Efros  
University of California, Berkeley



Figure 2: Given an image and an instruction for how to edit it, our model performs the appropriate edit. Our model does not require full descriptions for the input or output image, and edits images in the forward pass without per-example inversion or fine-tuning.

**1. Introduction**

We present a method for editing images from human instructions, given an input image and a written instruction that tells the model what to do. Our model follows these instructions to edit the image. To obtain training data for this problem, we combine the knowledge of image transformation models (ITMs) [Berman+19] and a text-image model (Stable Diffusion) – to generate a large dataset of image editing examples. Our conditional diffusion model, InstructPix2Pix, is trained on our generated data, and generalizes to real images and user-written instructions at inference time. Our model is trained in the forward pass and does not require per-example finetuning or inversion, our model edits images quickly, in a matter of seconds. We show compelling editing results for a diverse collection of input images and written instructions.

\*Denotes equal contribution.  
More results on our project page: [timothysbrooks.com/instruct-pix2pix](https://timothysbrooks.com/instruct-pix2pix)

InstructPix2Pix: model to perform edits on images

# Quick Aside

## Midway Reports and Grading:

- We (Oliver and I) finished reviewing through and grading all midway reports.
  - Everyone did a nice job :D!
- We have a few light comments on each report – which could be useful for final report.
  - We will likely send these via our group message threads on slack.
- And will also upload grades on Brightspace.

# Quick Aside

## Reminder – Final Course Project Timeline

- Final Presentations (30 mins long)
  - **April 25** - Groups 1-5.
  - **May 02** - Groups 6-9.

# Quick Aside

## Reminder – Final Course Project Timeline

- Final Presentations (30 mins long)
  - **April 25** - Groups 1-5.
  - **May 02** - Groups 6-9.
- Final Report (8-10 pages, extended version of midway report)
  - (LaTeX template will be provided soon).
  - **Due May 9, EOD**

# **Review: Predictive UQ and Gaussian Processes**

# Review – Gaussian Processes

## Review – Gaussian Processes

Intuitively, a GP is a probabilistic model, which can be viewed as a distribution over functions.

## Review – Gaussian Processes

Intuitively, a GP is a probabilistic model, which can be viewed as a distribution over functions.

But it uses a (multivariate) **Gaussian distribution** to model these functions !

## Review – Gaussian Processes

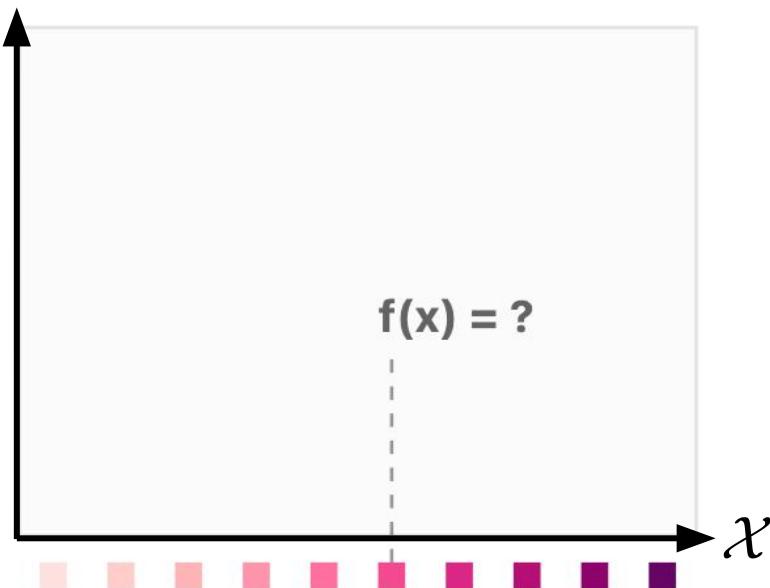
Intuitively, a GP is a probabilistic model, which can be viewed as a distribution over functions.

But it uses a (multivariate) **Gaussian distribution** to model these functions !

Let's visualize this...

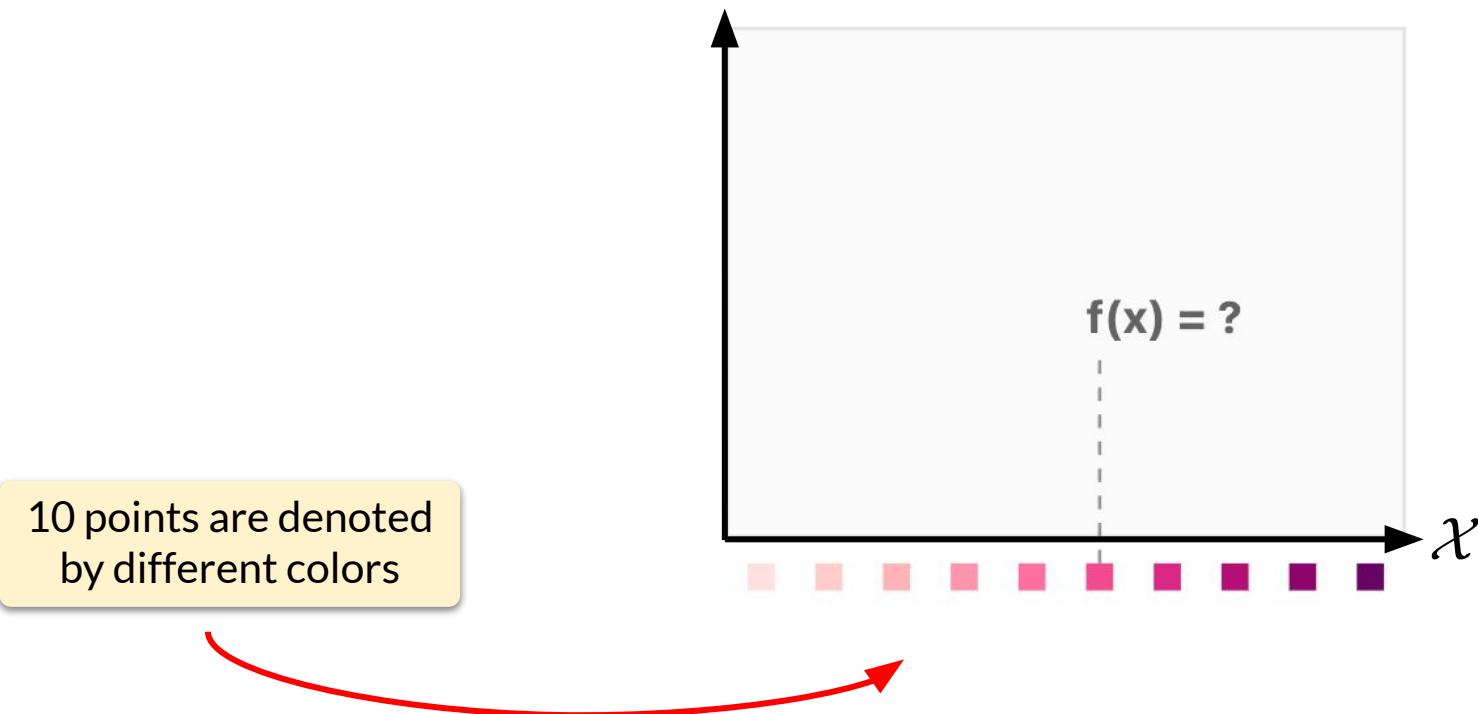
# Review – Gaussian Processes

Suppose we pick a set of  $d=10$  points on the  $x$ -axis, and want to define a **function** (or set of functions) over these points:



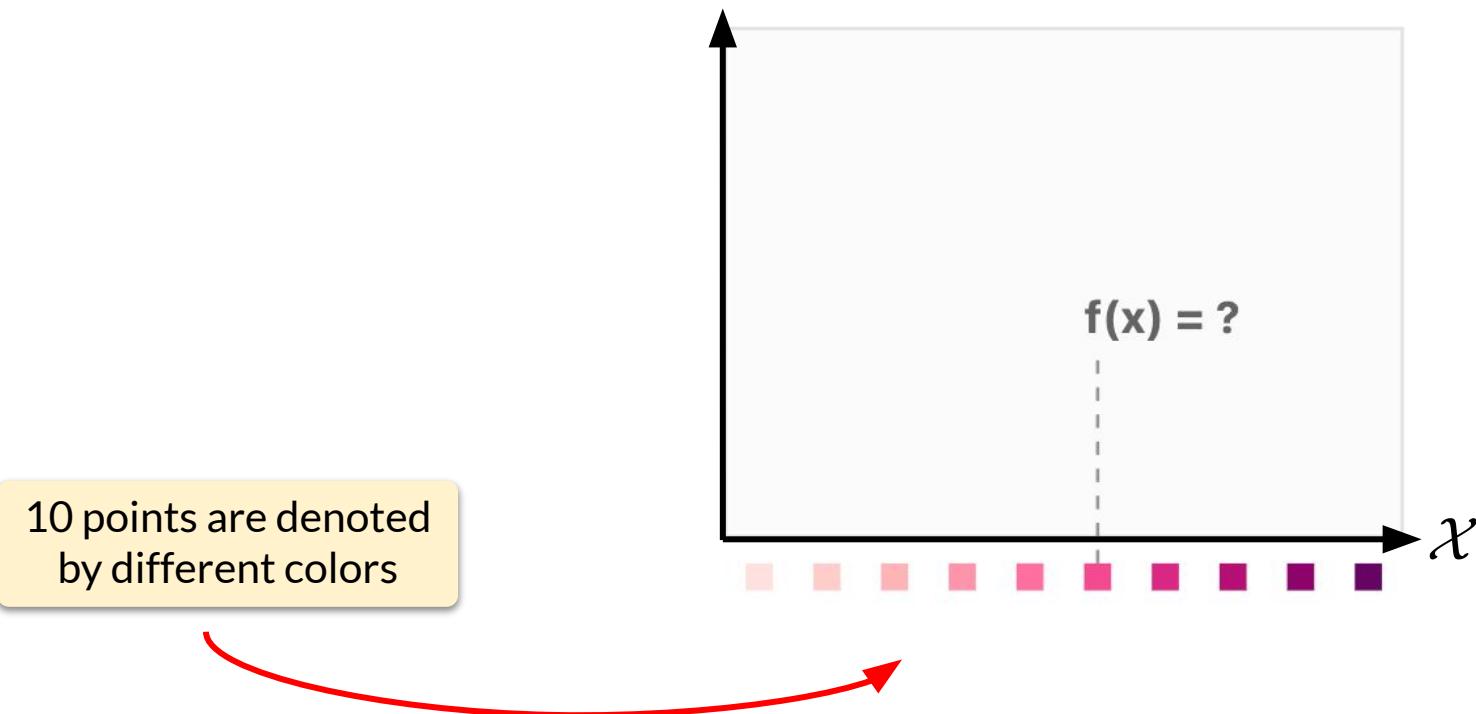
# Review – Gaussian Processes

Suppose we pick a set of  $d=10$  points on the  $x$ -axis, and want to define a **function** (or set of functions) over these points:



# Review – Gaussian Processes

Suppose we pick a set of  $d=10$  points on the  $x$ -axis, and want to define a **function** (or set of functions) over these points:



We can define a *distribution of functions*, defined on these 10 points, in the following way...

## Review – Gaussian Processes

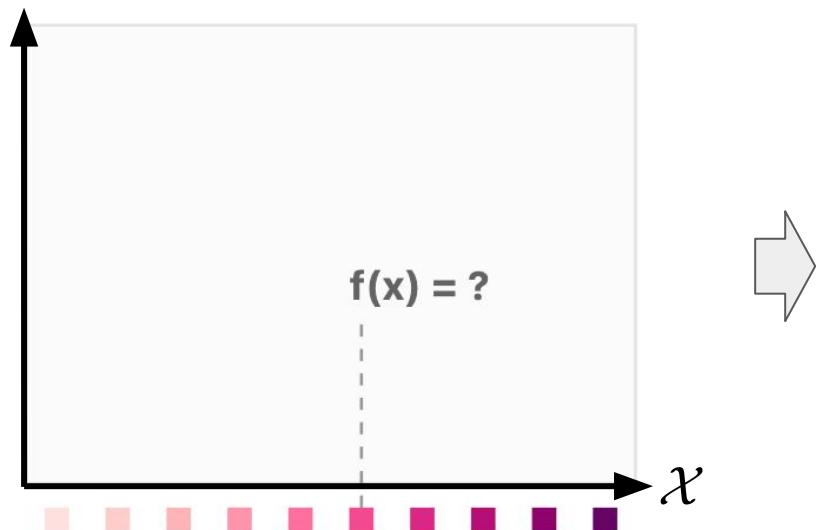
We define a MVN of  **$d$  dimensions**, with a mean and covariance matrix parameter.

Mean parameter  $\mu \in \mathbb{R}^d$  – assume we set this to **0** for now.

Covariance matrix  $\Sigma \in \mathbb{R}^{d \times d}$  – define using a *kernel function* based on distance between input points (next slide).

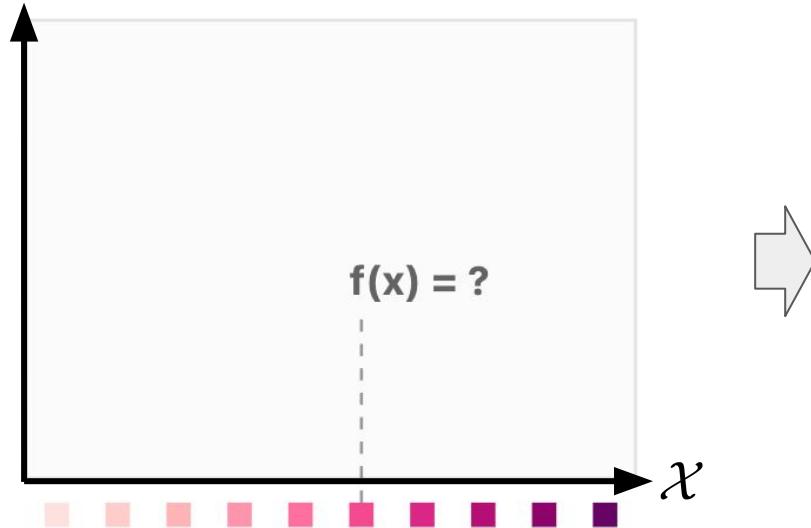
# Review – Gaussian Processes

For these 10 points:

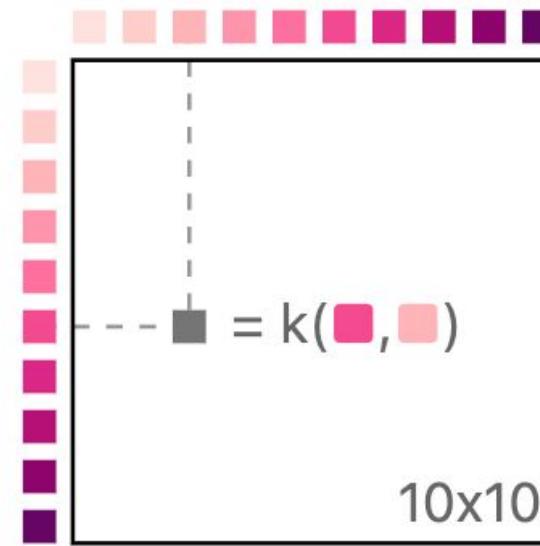


# Review – Gaussian Processes

For these 10 points:



Covariance matrix:



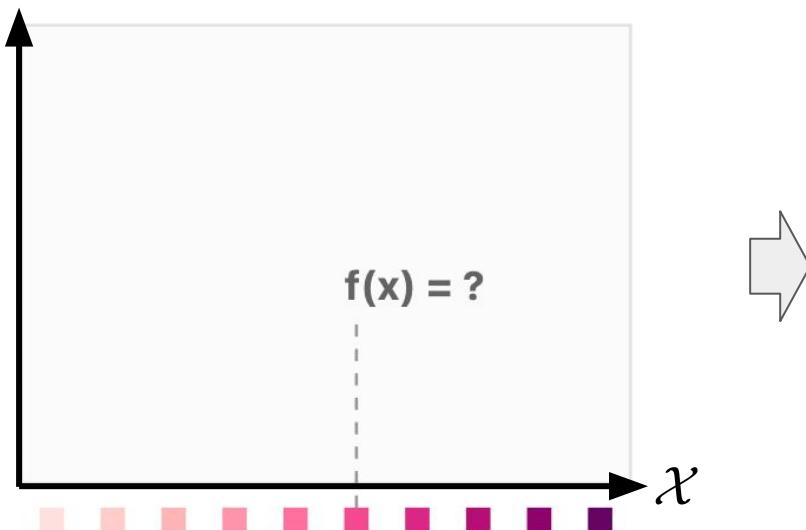
Define covariance between any pair of inputs using kernel function  $k(\cdot, \cdot)$ .

One example: *RBF kernel*,

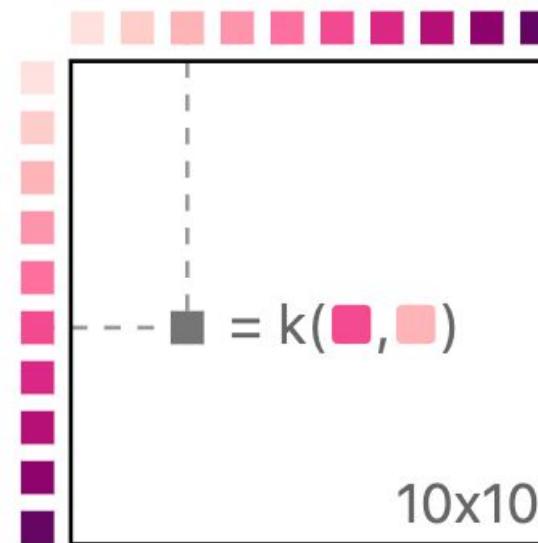
$$k(x_1, x_2) = \sigma^2 \left( -\frac{\|x_1 - x_2\|^2}{2\ell^2} \right)$$

# Review – Gaussian Processes

For these 10 points:



Covariance matrix:



Define covariance between any pair of inputs using kernel function  $k(\cdot, \cdot)$ .

One example: *RBF kernel*,

$$k(x_1, x_2) = \sigma^2 \left( -\frac{\|x_1 - x_2\|^2}{2\ell^2} \right)$$

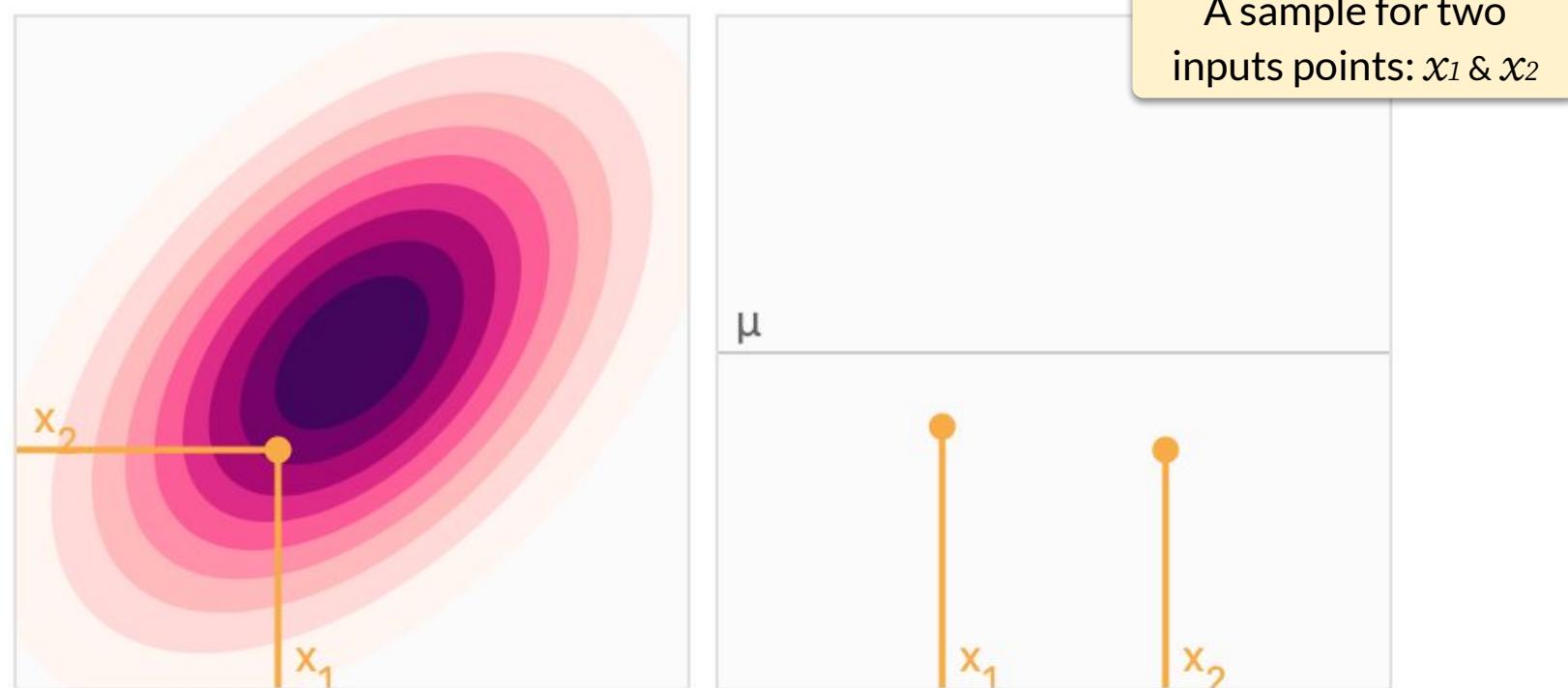
Highest when  $x_1 = x_2$ , lowest when  $x_2$  very different from  $x_2$ .

## Review – Gaussian Processes

This yields a MVN, where samples give function values for the  $d$  input points:

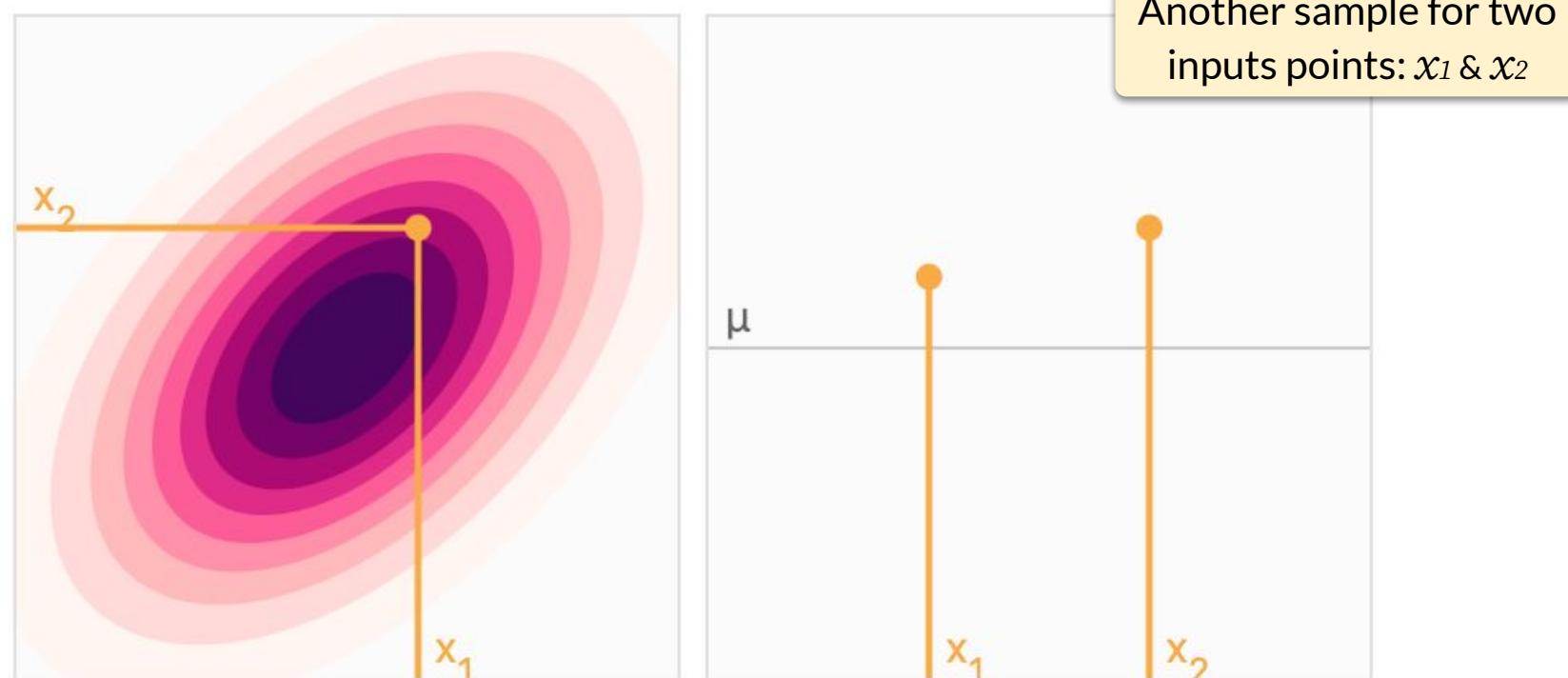
# Review – Gaussian Processes

This yields a MVN, where samples give function values for the  $d$  input points:



# Review – Gaussian Processes

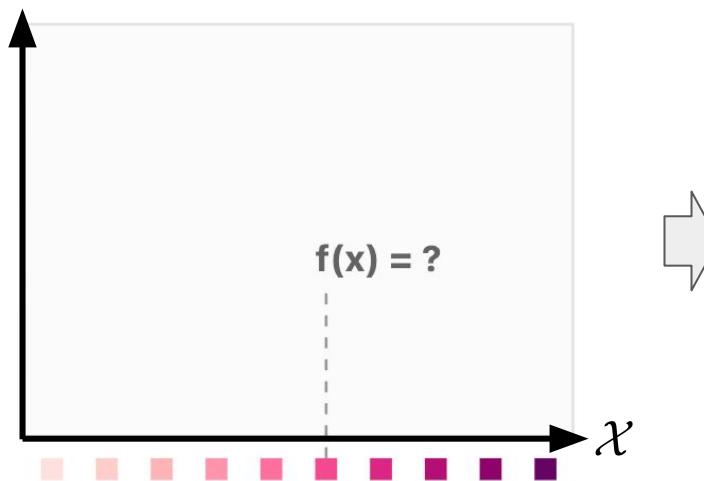
This yields a MVN, where samples give function values for the  $d$  input points:



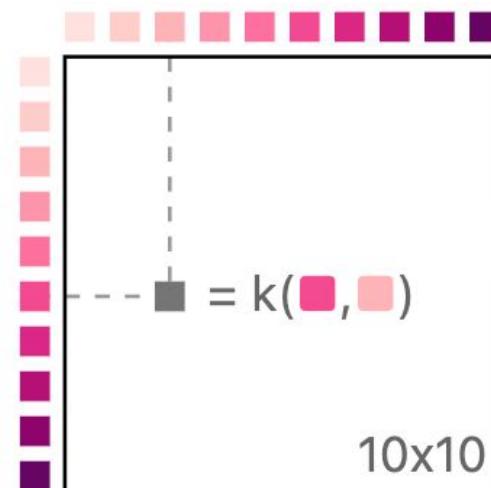
# Review – Gaussian Processes

This yields a MVN, where samples give function values for the  $d$  input points:

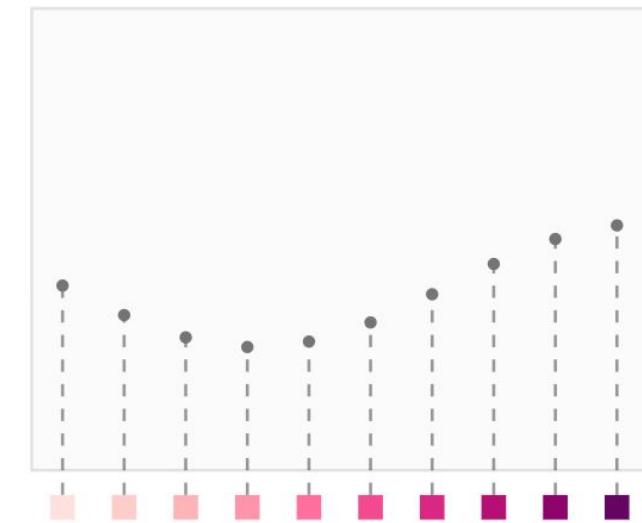
For these 10 points:



Covariance matrix:



Sampled function values



## Review – Gaussian Processes

We define a **Gaussian Process (GP)** as a distribution over functions, parameterized by a *mean function*  $m(\cdot)$  and *covariance kernel function*  $k(\cdot, \cdot)$ , written

$$f \sim \text{GP}(m(\cdot), k(\cdot, \cdot))$$

## Review – Gaussian Processes

We define a **Gaussian Process (GP)** as a distribution over functions, parameterized by a *mean function*  $m(\cdot)$  and *covariance kernel function*  $k(\cdot, \cdot)$ , written

$$f \sim \text{GP}(m(\cdot), k(\cdot, \cdot))$$

Where, for any finite collection of  $d$  input points, we can:

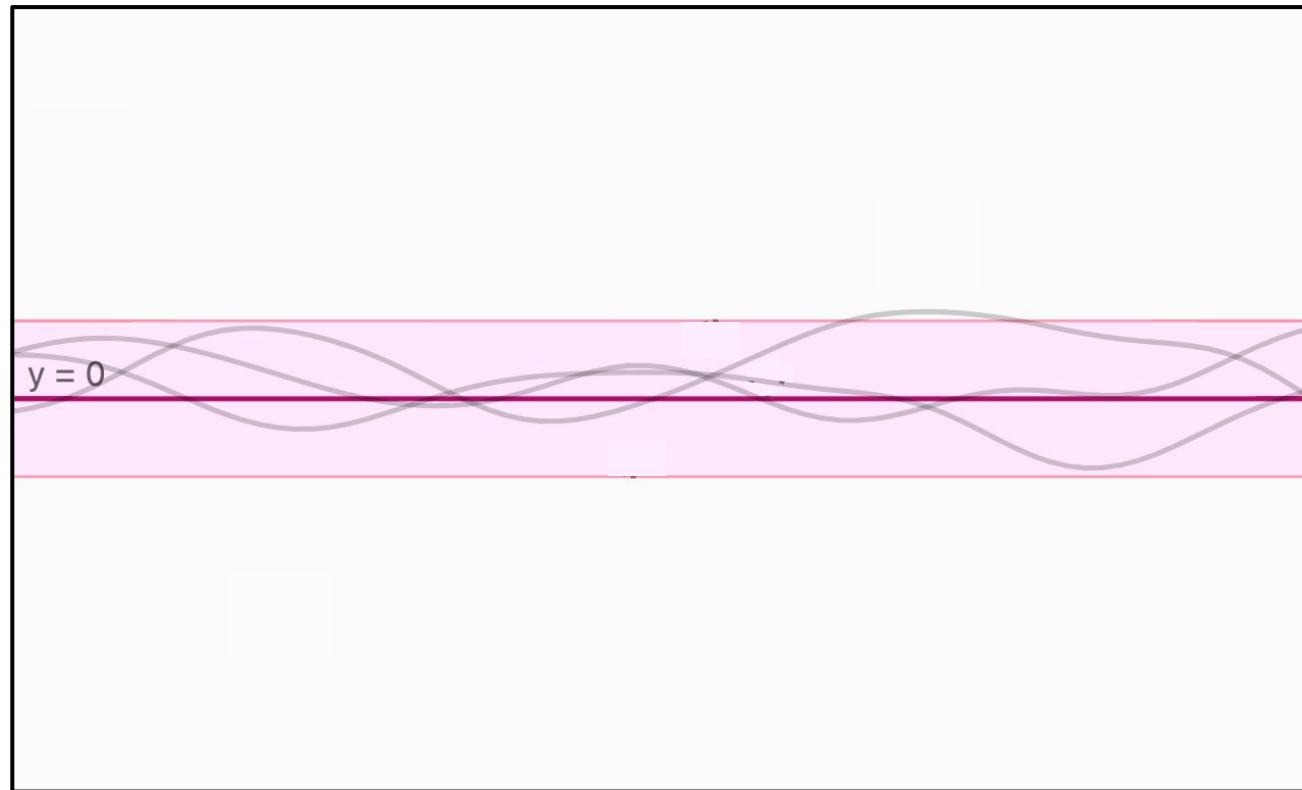
- Define a MVN on these  $d$  points, using the mean/kernel functions to define a mean parameter and covariance matrix.
- And this MVN yields a distribution over function values.

# Review – Gaussian Processes

What does this distribution look like?

# Review – Gaussian Processes

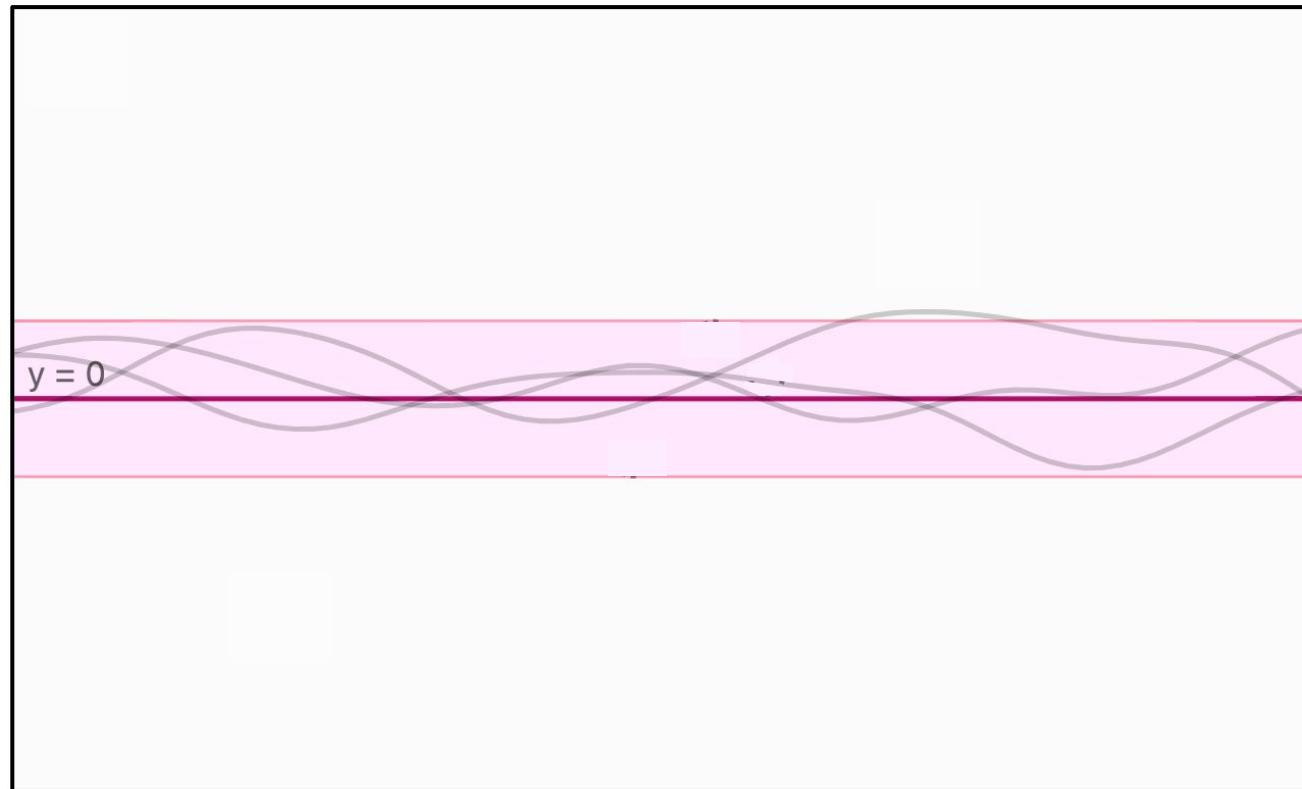
What does this distribution look like? → often drawn like this:



# Review – Gaussian Processes

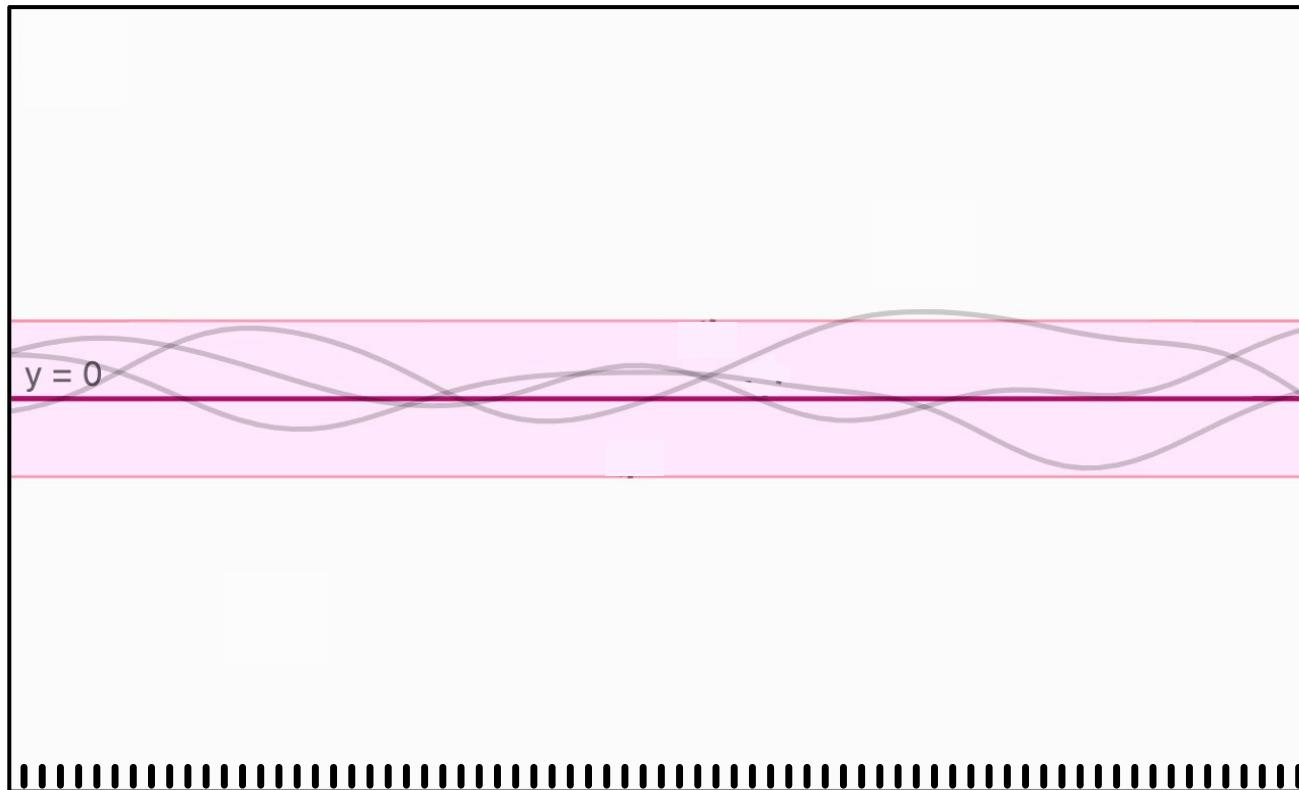
Called a *Sausage Plot* because it looks like a sausage

What does this distribution look like? → often drawn like this:



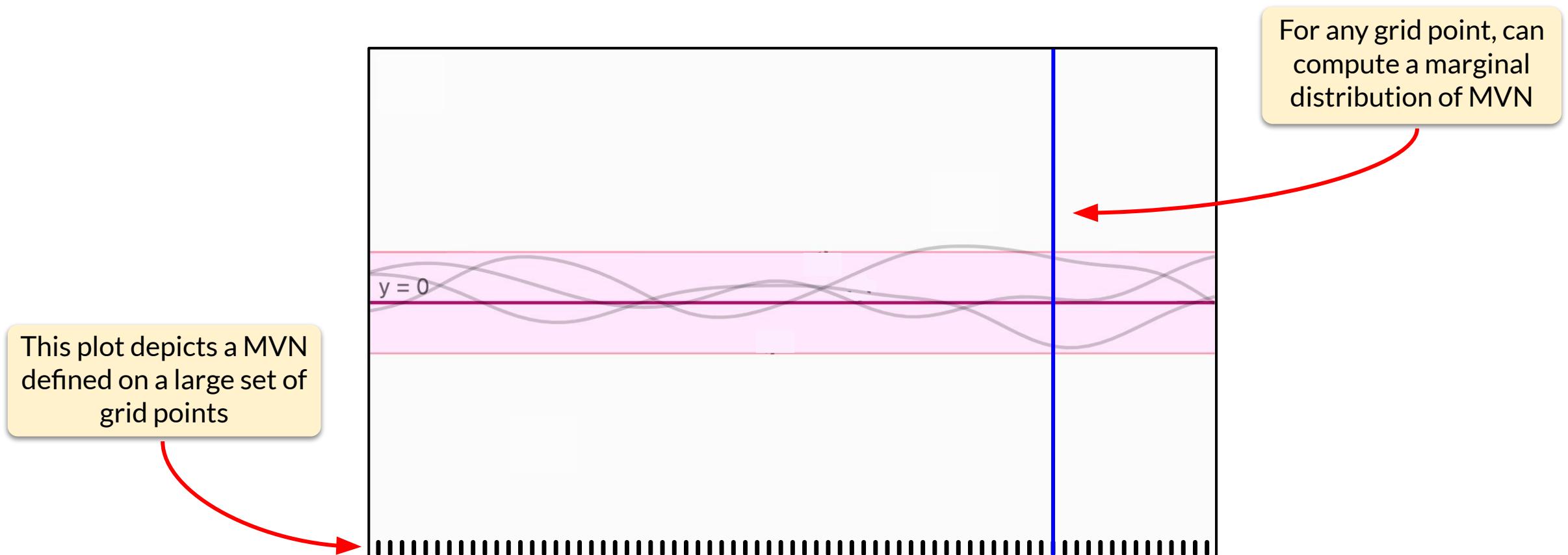
# Review – Gaussian Processes

What does this distribution look like? → often drawn like this:



# Review – Gaussian Processes

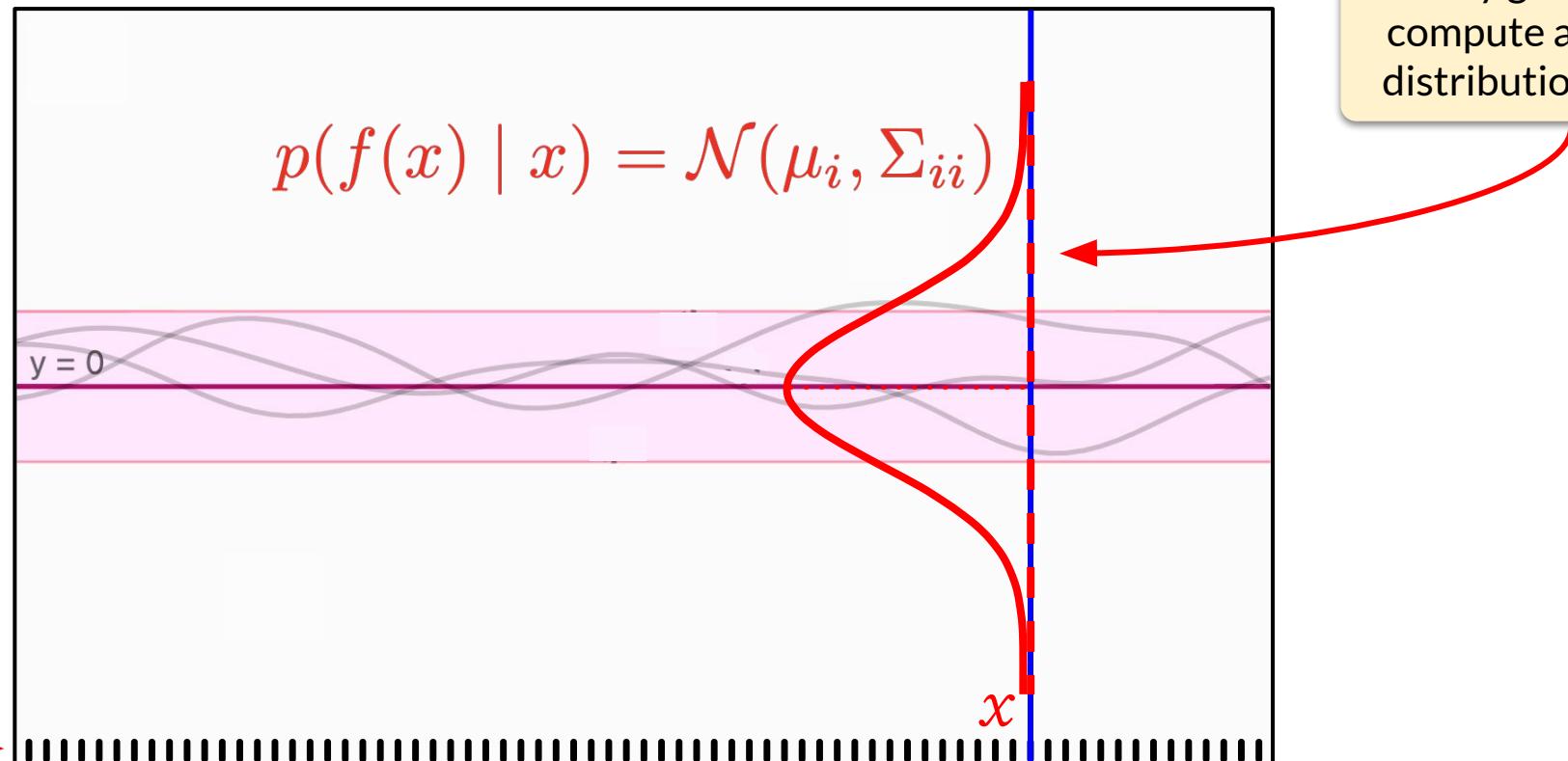
What does this distribution look like? → often drawn like this:



# Review – Gaussian Processes

What does this distribution look like? → often drawn like this:

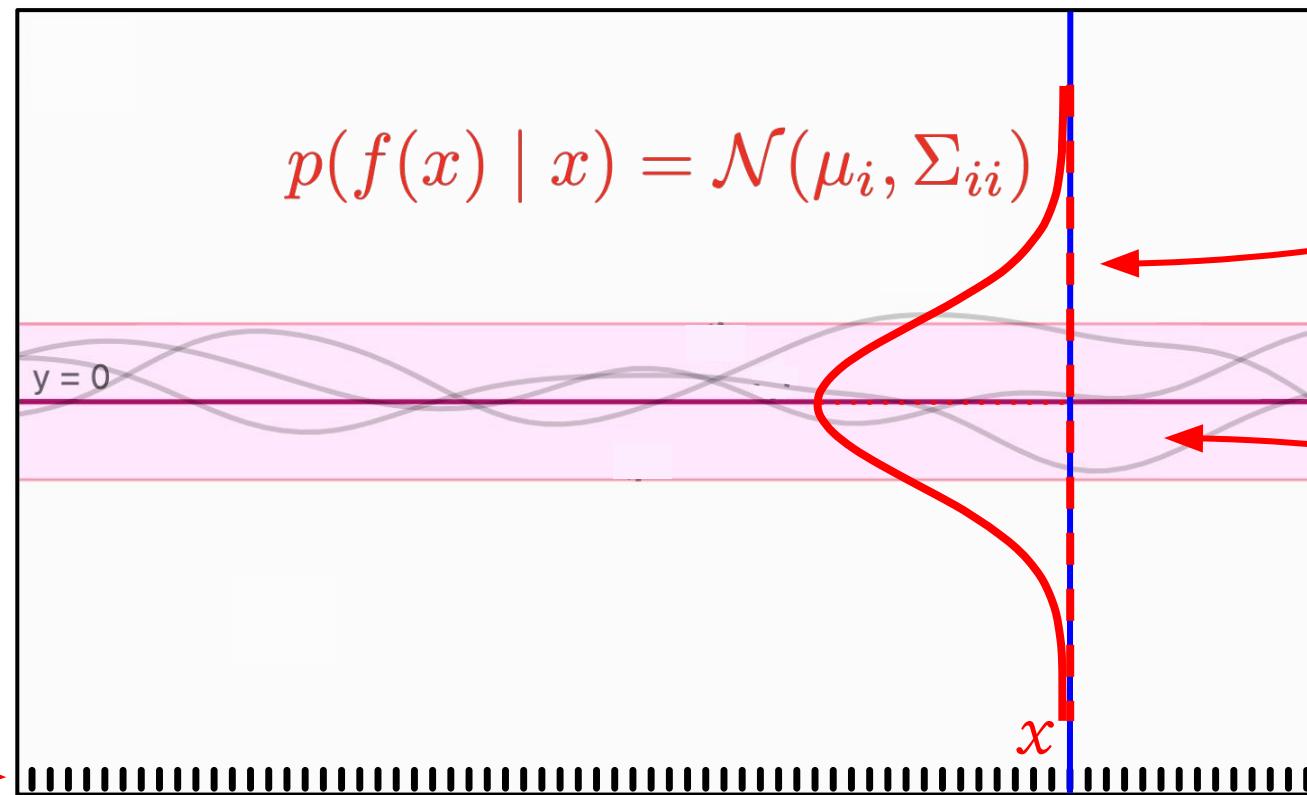
This plot depicts a MVN defined on a large set of grid points



# Review – Gaussian Processes

What does this distribution look like? → often drawn like this:

This plot depicts a MVN defined on a large set of grid points



For any grid point, can compute a marginal distribution of MVN

Samples from the MVN are shown in grey (appear as functions!)

# Review – Gaussian Process Posterior

How do we use a GP as a predictive uncertainty model?

## Review – Gaussian Process Posterior

How do we use a GP as a predictive uncertainty model?

⇒ Use it as a prior distribution and do Bayesian inference!

# Review – Gaussian Process Posterior

How do we use a GP as a predictive uncertainty model?

⇒ Use it as a prior distribution and do Bayesian inference!

In particular, we will:

- Use a GP to define a prior over functions.
- Observe a set of  $(x, y)$  pairs.
- Compute a posterior GP over functions.

# Review – Gaussian Process Posterior

How do we use a GP as a predictive uncertainty model?

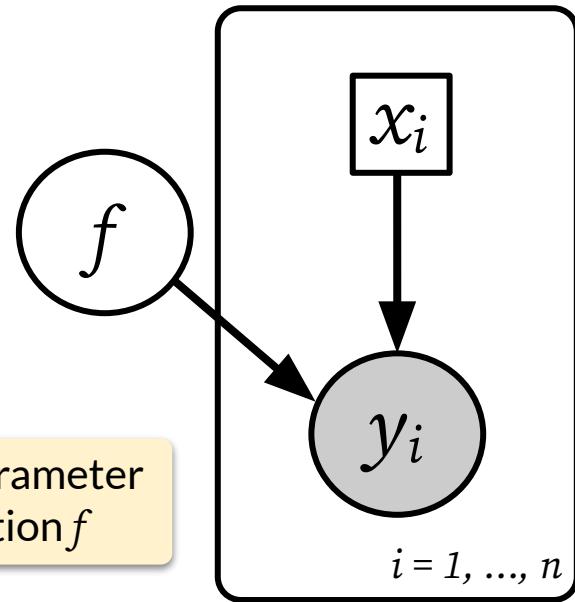
⇒ Use it as a prior distribution and do Bayesian inference!

In particular, we will:

- Use a GP to define a prior over functions.
- Observe a set of  $(x, y)$  pairs.
- Compute a posterior GP over functions.

In GPs, the parameter  
is the function  $f$

Prior on Parameter



# Review – Gaussian Process Posterior

Computing the posterior:

# Review – Gaussian Process Posterior

Computing the posterior:

- We have a GP prior (*essentially a multivariate normal*):  $f \sim \text{GP}(m(\cdot), k(\cdot, \cdot))$

# Review – Gaussian Process Posterior

Computing the posterior:

- We have a GP prior (*essentially a multivariate normal*):  $f \sim \text{GP}(m(\cdot), k(\cdot, \cdot))$
- We have a Gaussian likelihood for observations:  $y \sim p(y \mid f ; x) = \mathcal{N}(f, \sigma^2)$

# Review – Gaussian Process Posterior

Computing the posterior:

- We have a GP prior (*essentially a multivariate normal*):  $f \sim \text{GP}(m(\cdot), k(\cdot, \cdot))$
- We have a Gaussian likelihood for observations:  $y \sim p(y \mid f ; x) = \mathcal{N}(f, \sigma^2)$
- $\Rightarrow$  Posterior is (*essentially*) also a multivariate normal distribution.

# Review – Gaussian Process Posterior

Computing the posterior:

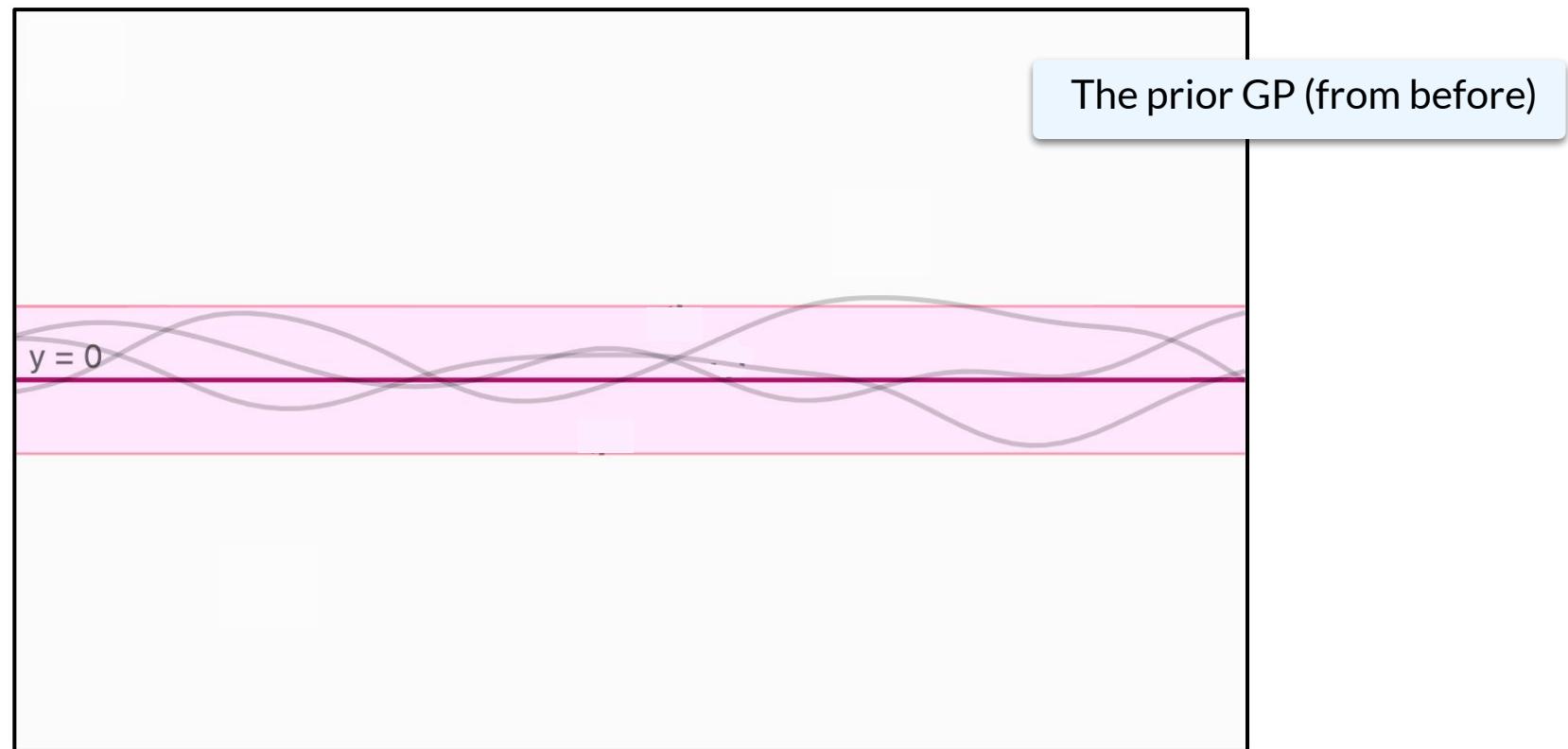
- We have a GP prior (*essentially a multivariate normal*):  $f \sim \text{GP}(m(\cdot), k(\cdot, \cdot))$
- We have a Gaussian likelihood for observations:  $y \sim p(y \mid f ; x) = \mathcal{N}(f, \sigma^2)$
- $\Rightarrow$  Posterior is (*essentially*) also a multivariate normal distribution.

Posterior is technically a **posterior Gaussian process**, so it's another GP with a mean function and covariance kernel parameters

$\Rightarrow$  which can be computed in closed form.

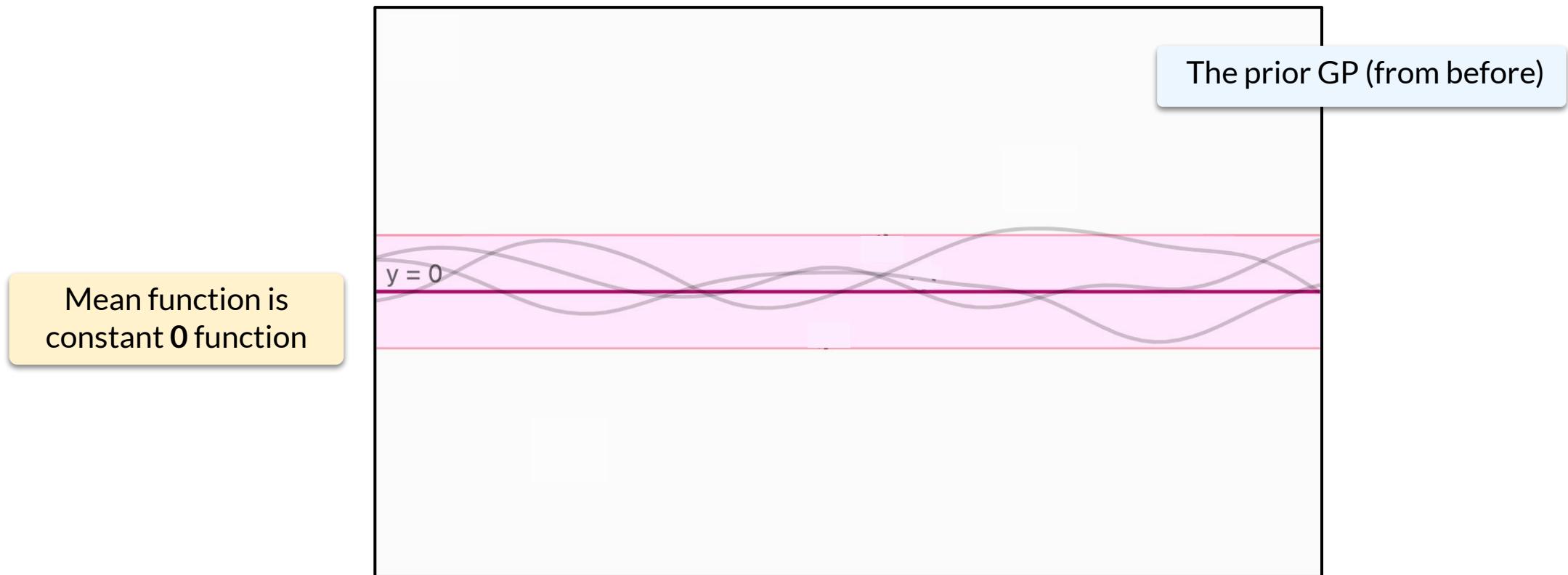
# Review – Gaussian Process Posterior

Visualizing this:



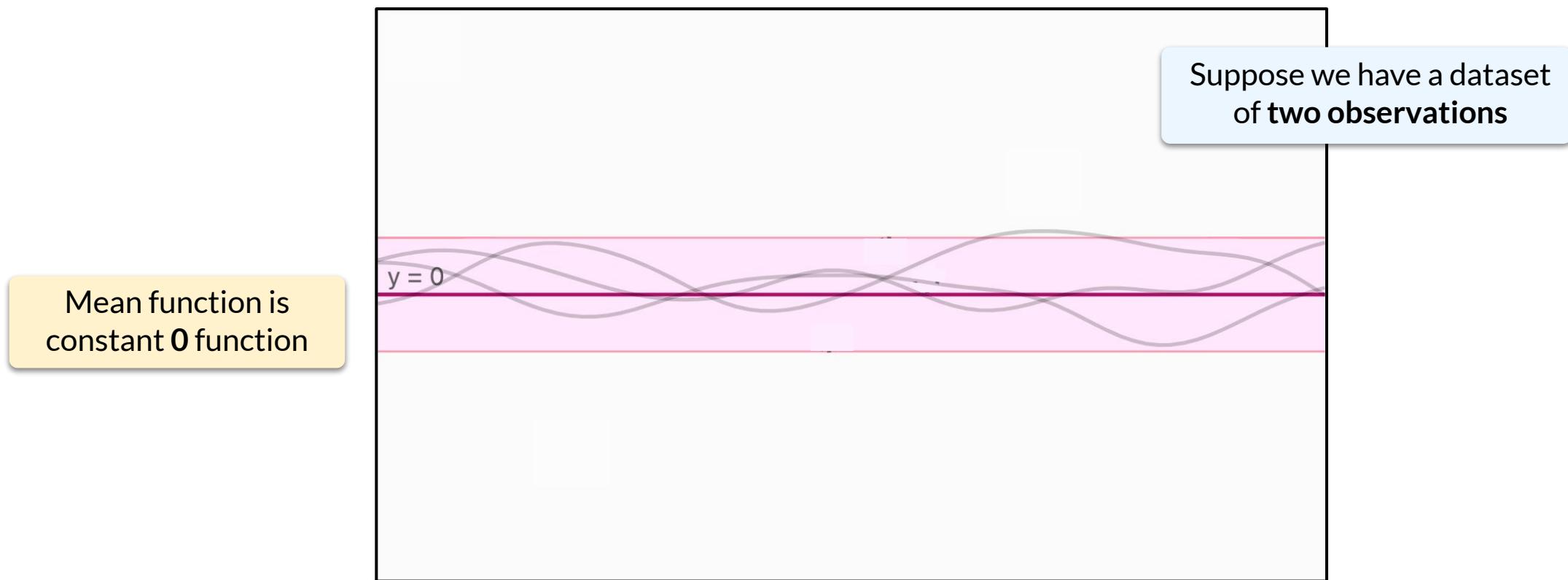
# Review – Gaussian Process Posterior

Visualizing this:



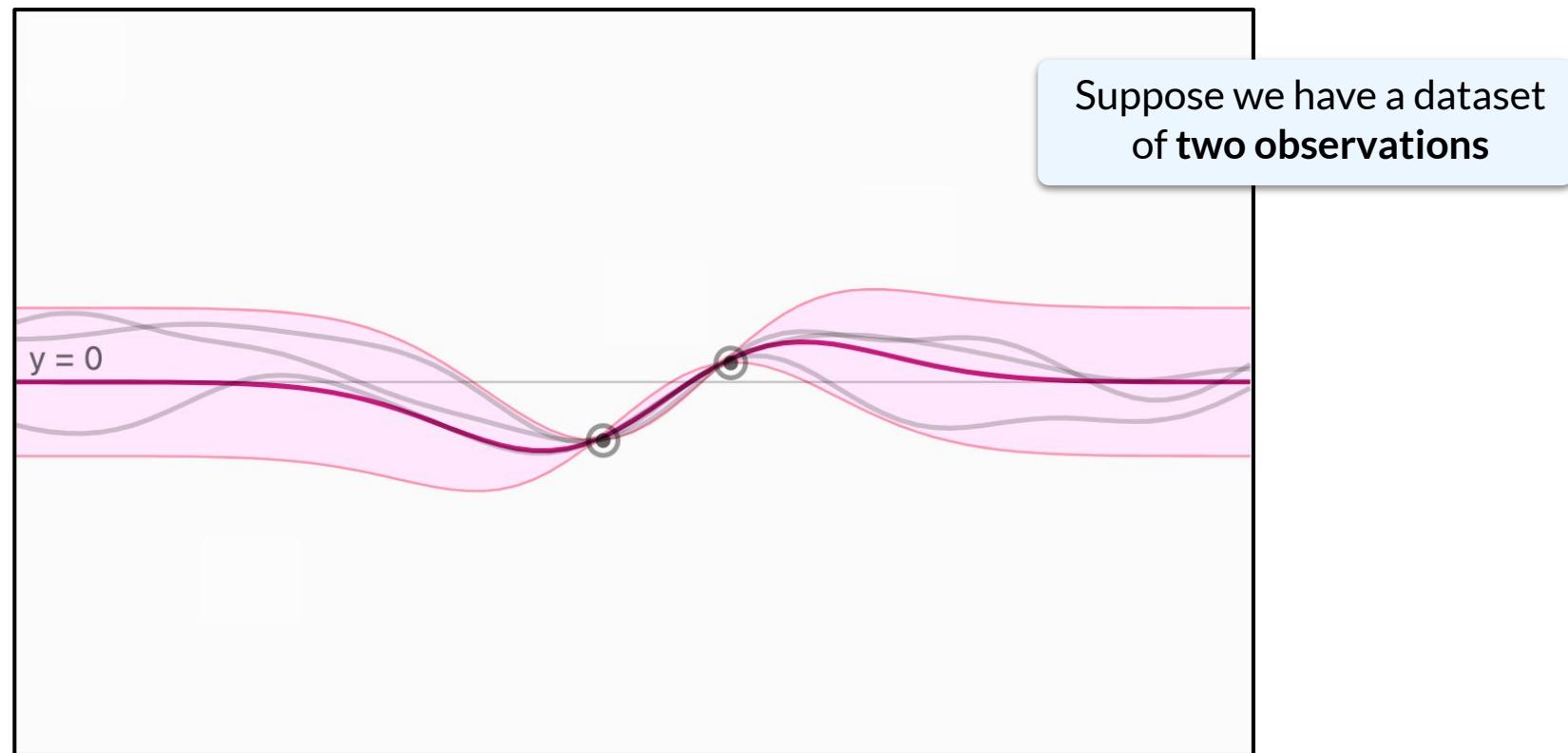
# Review – Gaussian Process Posterior

Visualizing this:



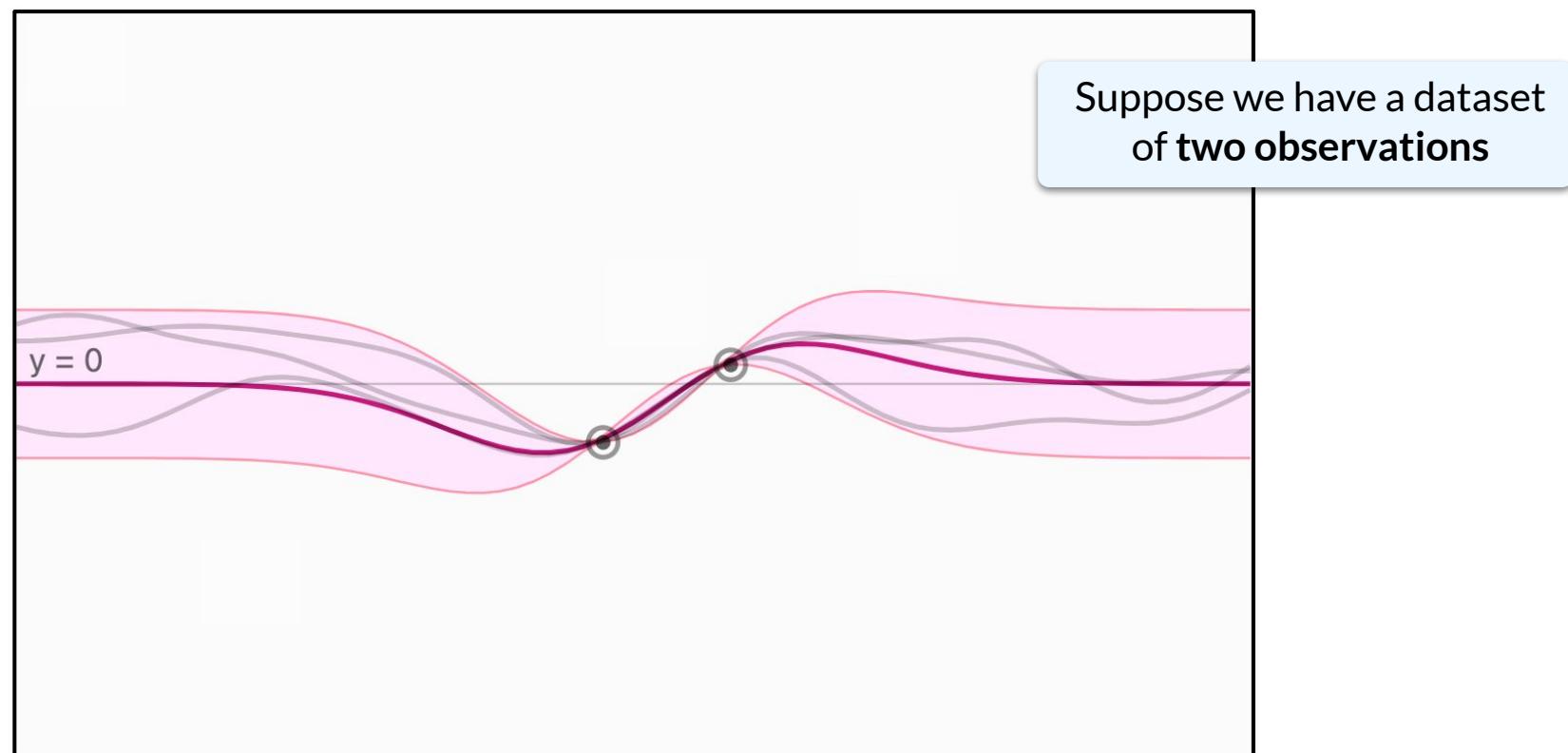
# Review – Gaussian Process Posterior

Visualizing this:



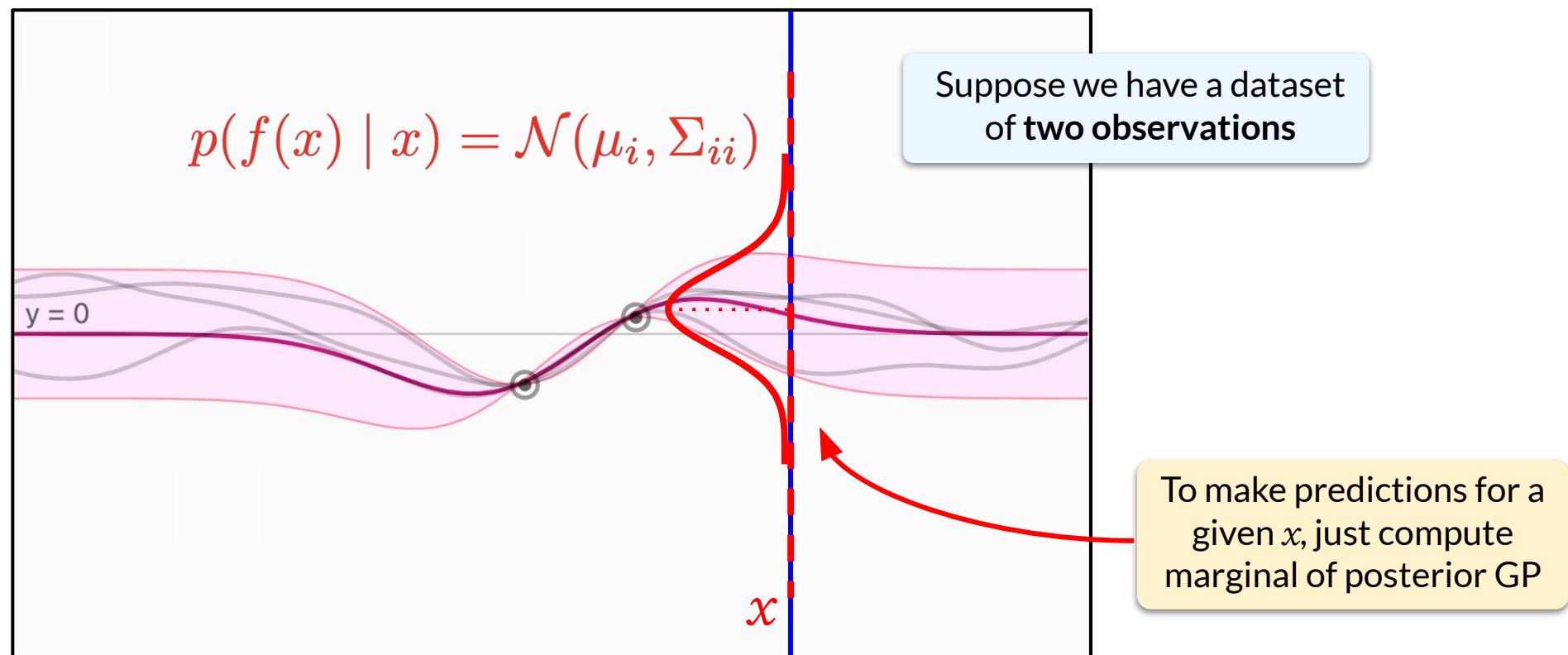
# Review – Gaussian Process Posterior

Visualizing this:



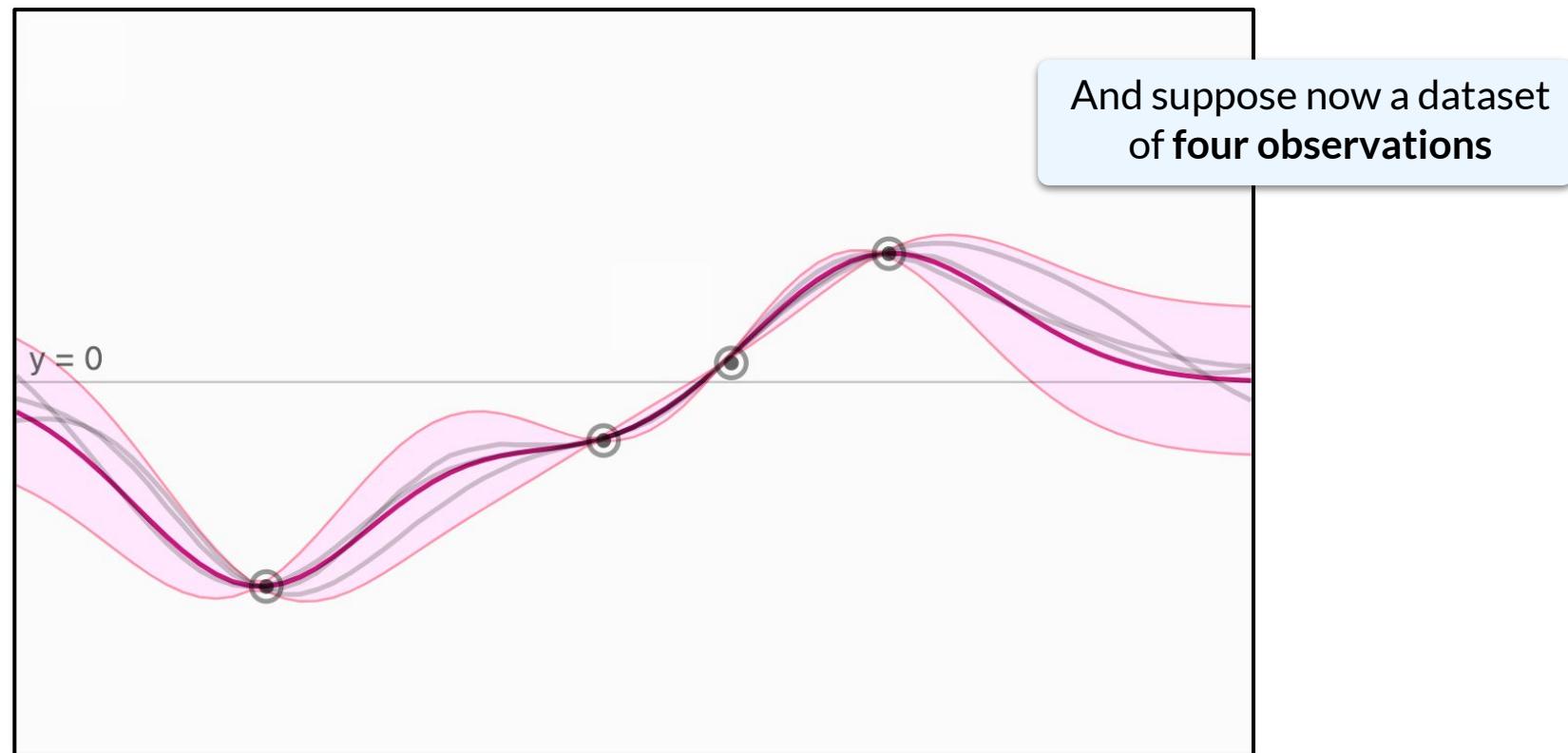
# Review – Gaussian Process Posterior

Visualizing this:



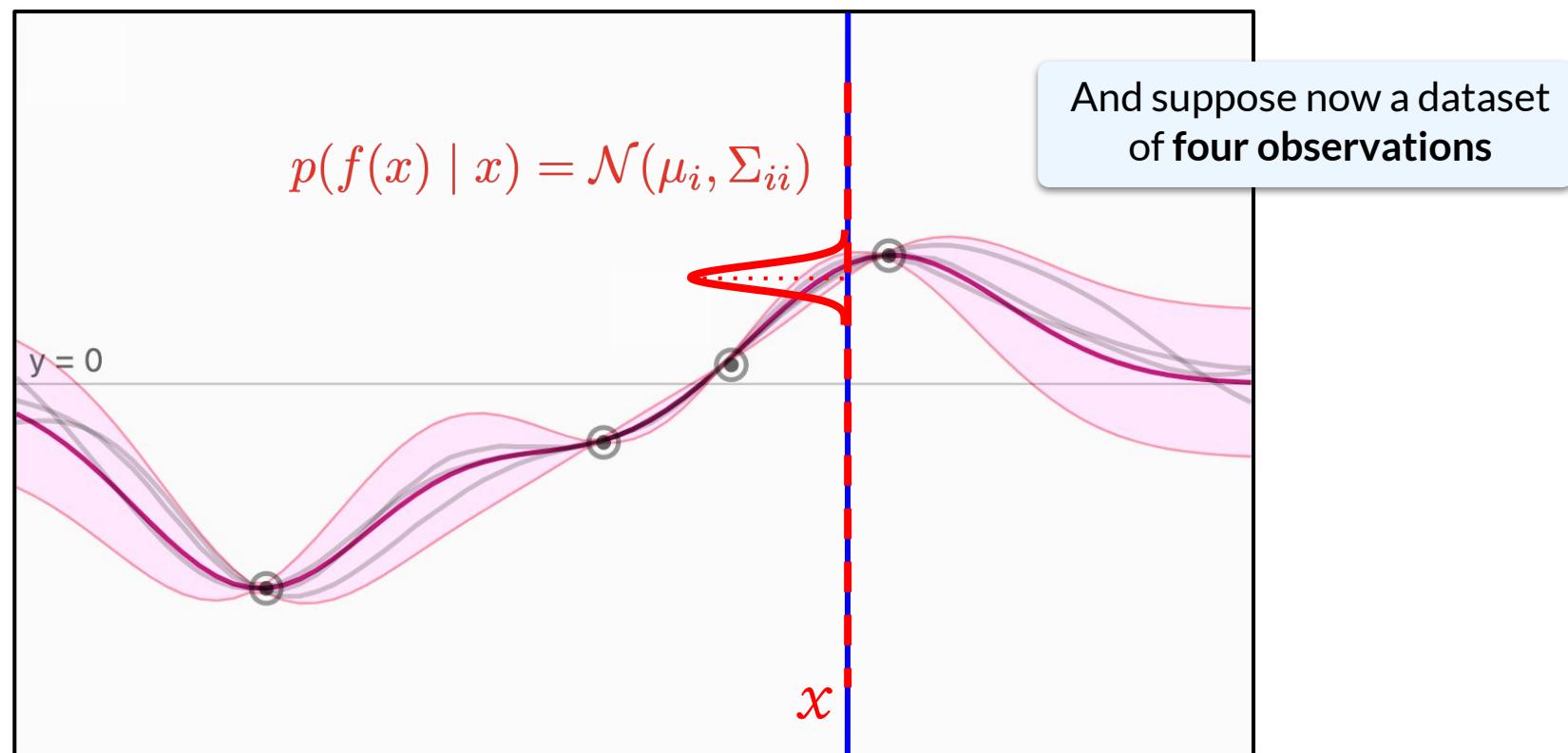
# Review – Gaussian Process Posterior

Visualizing this:



# Review – Gaussian Process Posterior

Visualizing this:



# Gaussian Processes – Summary

GP Summary:

# Gaussian Processes – Summary

GP Summary:

- GP is parameterized (defined) by a *mean function* and *covariance kernel function*.

# Gaussian Processes – Summary

## GP Summary:

- GP is parameterized (defined) by a *mean function* and *covariance kernel function*.
- For any set of input points  $x$ , these parameters yield a MVN distribution

# Gaussian Processes – Summary

## GP Summary:

- GP is parameterized (defined) by a *mean function* and *covariance kernel function*.
- For any set of input points  $x$ , these parameters yield a MVN distribution
- And this MVN can be viewed as a **distribution over function values** (for the input points).

# Gaussian Processes – Summary

GP Summary:

- GP is parameterized (defined) by a *mean function* and *covariance kernel function*.
- For any set of input points  $x$ , these parameters yield a MVN distribution
- And this MVN can be viewed as a **distribution over function values** (for the input points).
- We can then do Bayesian inference:

# Gaussian Processes – Summary

## GP Summary:

- GP is parameterized (defined) by a *mean function* and *covariance kernel function*.
- For any set of input points  $x$ , these parameters yield a MVN distribution
- And this MVN can be viewed as a **distribution over function values** (for the input points).
- We can then do Bayesian inference:
  - Use GP as a prior distribution.
  - Define a Gaussian likelihood for observations.
  - Then compute a GP posterior, given a set of data.

# Gaussian Processes – Summary

## GP Summary:

- GP is parameterized (defined) by a *mean function* and *covariance kernel function*.
- For any set of input points  $x$ , these parameters yield a MVN distribution
- And this MVN can be viewed as a **distribution over function values** (for the input points).
- We can then do Bayesian inference:
  - Use GP as a prior distribution.
  - Define a Gaussian likelihood for observations.
  - Then compute a GP posterior, given a set of data.
- ⇒ This yields a predictive UQ model!

# **Deep Uncertainty Models**

# Deep Uncertainty Models

We can also define predictive UQ models using neural networks.

# Deep Uncertainty Models

We can also define predictive UQ models using neural networks.

I will go through a few examples of these.

# Deep Uncertainty Models

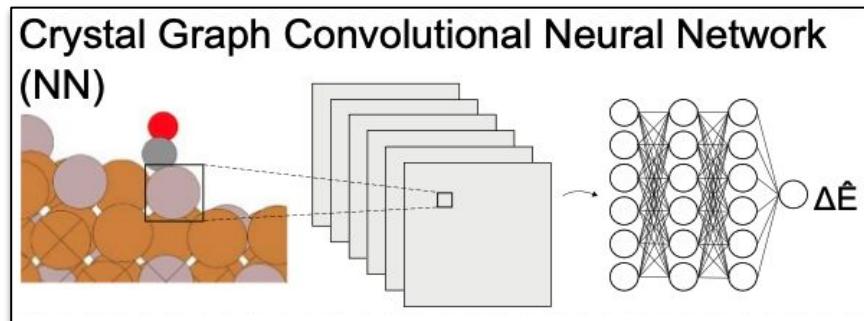
As a running example, we'll look at methods for incorporating predictive uncertainty into the following neural network:

# Deep Uncertainty Models – Example: Computational Catalyst Design

As a running example, we'll look at methods for incorporating predictive uncertainty into the following neural network:

# Deep Uncertainty Models – Example: Computational Catalyst Design

As a running example, we'll look at methods for incorporating predictive uncertainty into the following neural network:



Xie, Tian, and Jeffrey C. Grossman. "Crystal graph convolutional neural networks for an accurate and interpretable prediction of material properties." *Physical review letters*, 2018.

CGCNN Model (*a graph convolutional neural network model*)

- *Inputs:* three-dimensional atomic structure (a graph).
- *Outputs:* DFT-calculated site adsorption energies  $\Delta E$ .  
( $\Rightarrow$  Regression).

"Methods for comparing uncertainty quantifications for material property predictions", \*Tran, \*Neiswanger, et al., MLST. 2020

"Computational catalyst discovery: Active classification through myopic multiscale sampling", \*Tran, \*Neiswanger, et al., J. Chem. Phys. 2021

## Deep Uncertainty Models – Example: Computational Catalyst Design

Suppose we want to *incorporate uncertainty into this neural network.*

(I.e., use neural network for both point predictions and uncertainties).

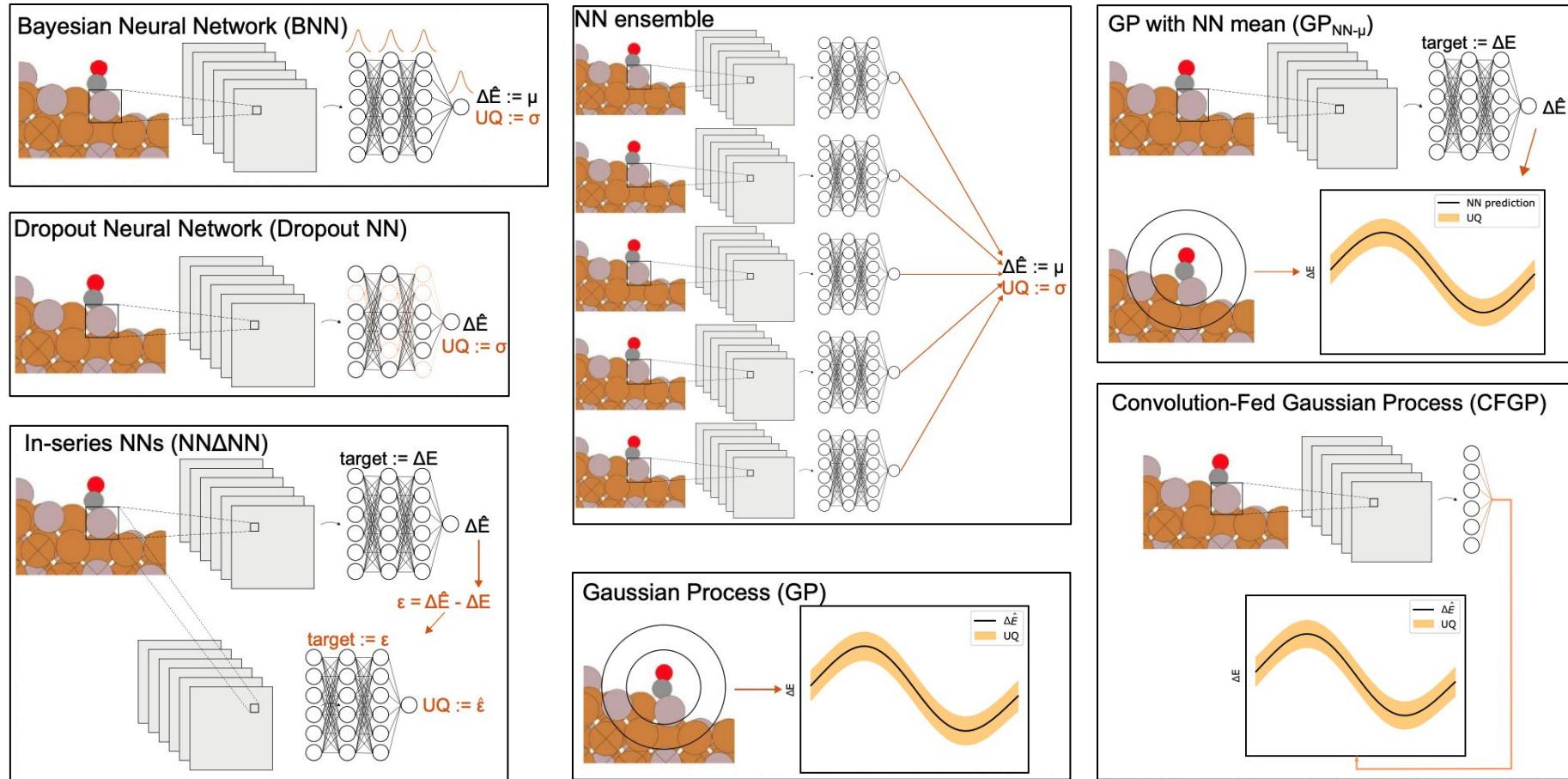
## Deep Uncertainty Models – Example: Computational Catalyst Design

Suppose we want to *incorporate uncertainty into this neural network.*

(I.e., use neural network for both point predictions and uncertainties).

⇒ there are many ways this can be done, and it's an active area of research!

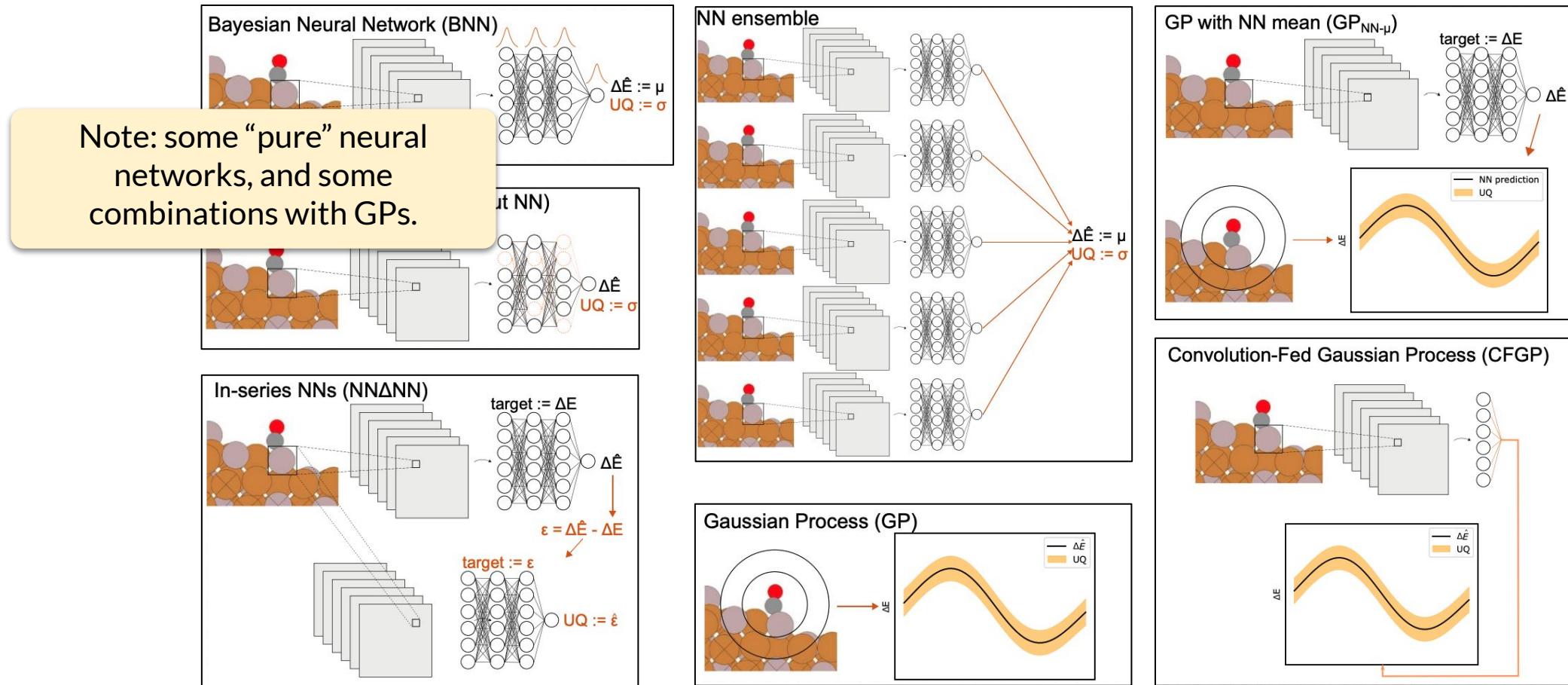
# Deep Uncertainty Models – Example: Computational Catalyst Design



"Methods for comparing uncertainty quantifications for material property predictions", \*Tran, \*Neiswanger, et al., MLST. 2020

"Computational catalyst discovery: Active classification through myopic multiscale sampling", \*Tran, \*Neiswanger, et al., J. Chem. Phys. 2021

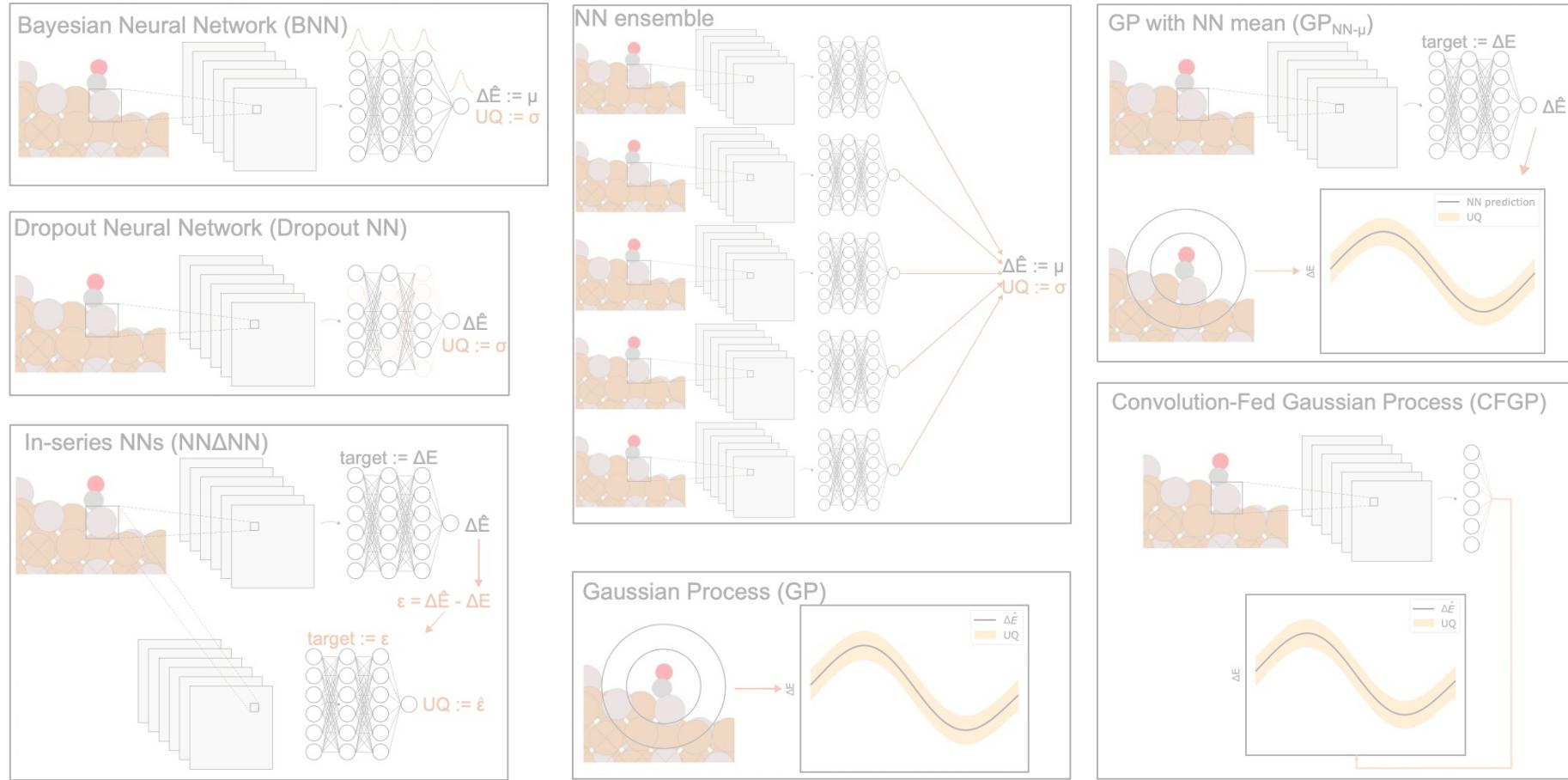
# Deep Uncertainty Models – Example: Computational Catalyst Design



“Methods for comparing uncertainty quantifications for material property predictions”, \*Tran, \*Neiswanger, et al., MLST. 2020

“Computational catalyst discovery: Active classification through myopic multiscale sampling”, \*Tran, \*Neiswanger, et al., J. Chem. Phys. 2021

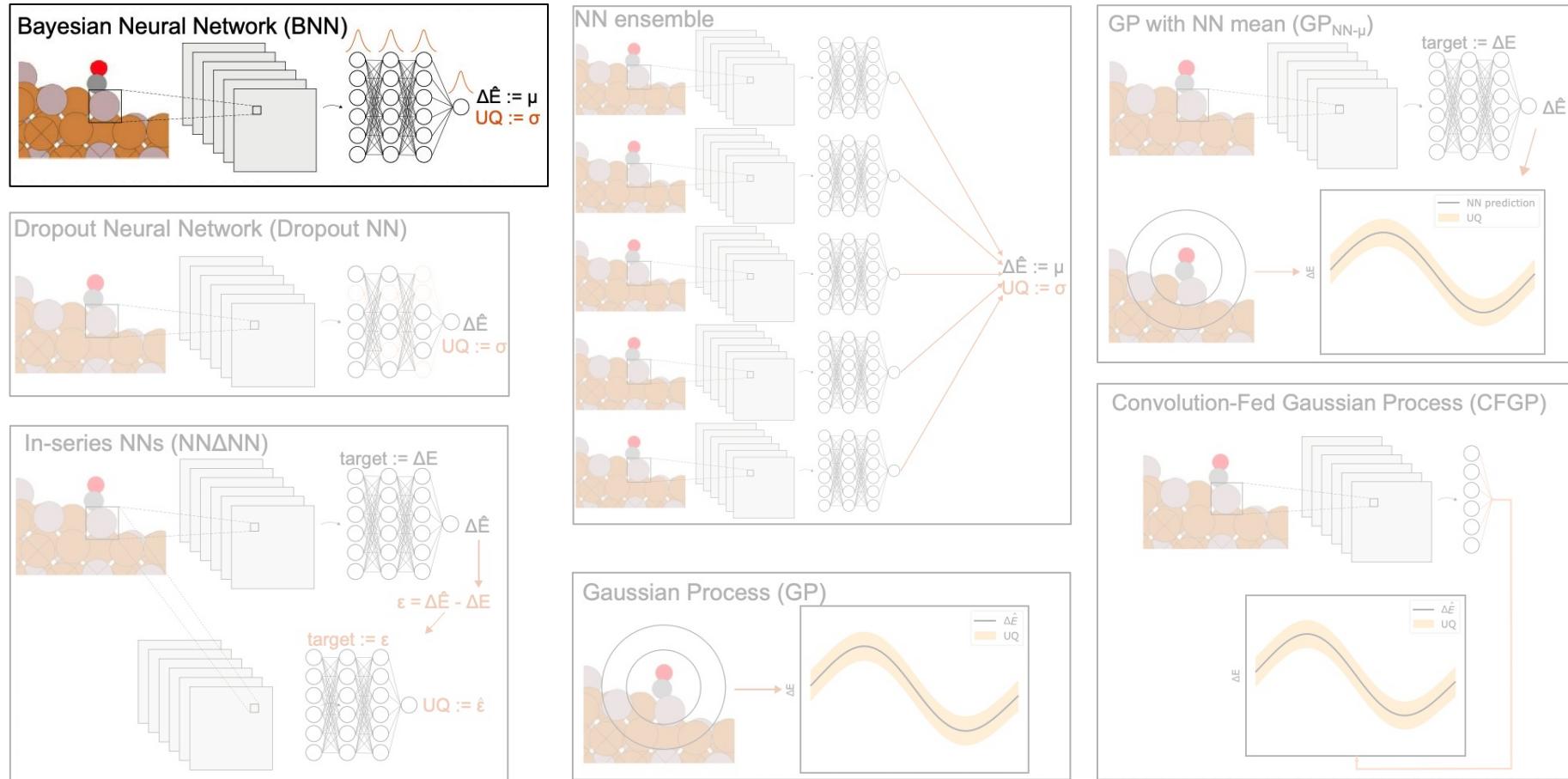
# Deep Uncertainty Models – Example: Computational Catalyst Design



"Methods for comparing uncertainty quantifications for material property predictions", \*Tran, \*Neiswanger, et al., MLST. 2020

"Computational catalyst discovery: Active classification through myopic multiscale sampling", \*Tran, \*Neiswanger, et al., J. Chem. Phys. 2021

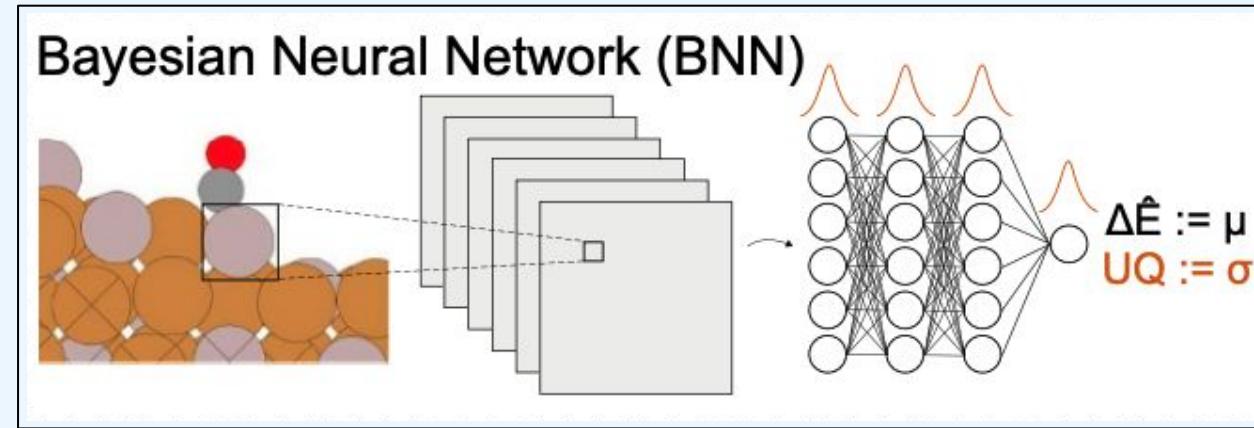
# Deep Uncertainty Models – Example: Computational Catalyst Design



"Methods for comparing uncertainty quantifications for material property predictions", \*Tran, \*Neiswanger, et al., MLST. 2020

"Computational catalyst discovery: Active classification through myopic multiscale sampling", \*Tran, \*Neiswanger, et al., J. Chem. Phys. 2021

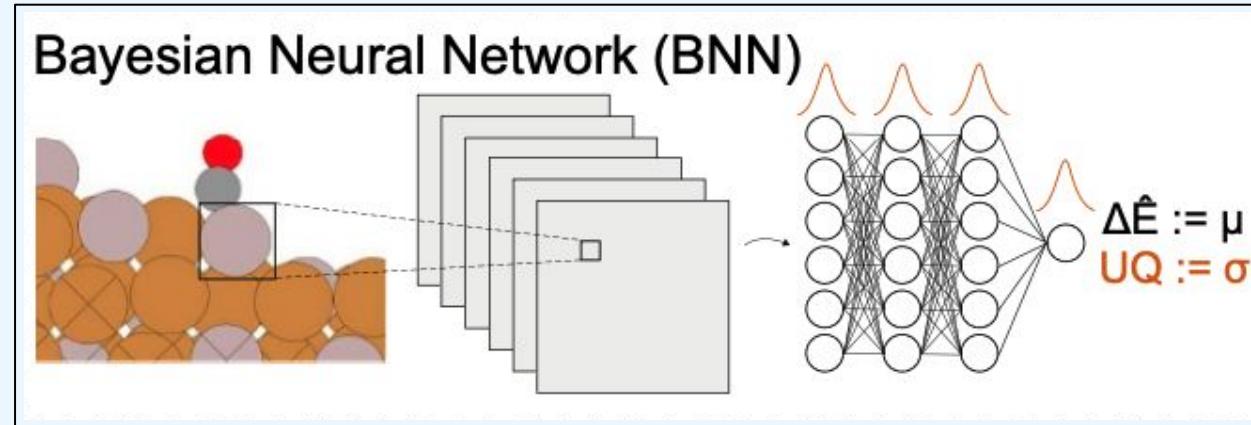
# Deep Uncertainty Models – Example: Computational Catalyst Design



"Methods for comparing uncertainty quantifications for material property predictions", \*Tran, \*Neiswanger, et al., MLST. 2020

"Computational catalyst discovery: Active classification through myopic multiscale sampling", \*Tran, \*Neiswanger, et al., J. Chem. Phys. 2021

# Deep Uncertainty Models – Example: Computational Catalyst Design

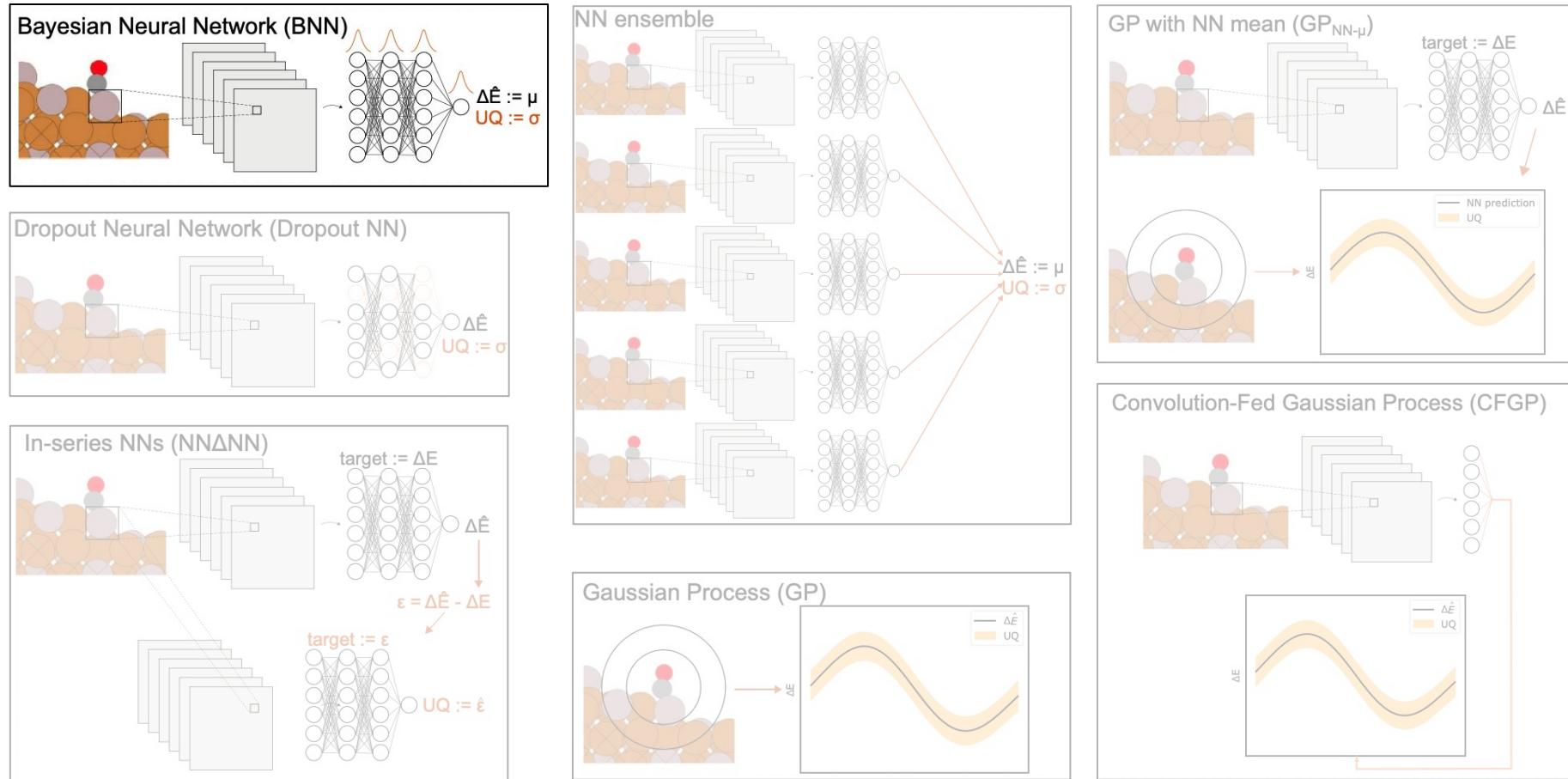


- Define a prior on neural network weights.
- In training, compute posterior over neural network weights (rather than a point estimate).
- Often using some type of MCMC (e.g., LMC) or variational inference.

"Methods for comparing uncertainty quantifications for material property predictions", \*Tran, \*Neiswanger, et al., MLST. 2020

"Computational catalyst discovery: Active classification through myopic multiscale sampling", \*Tran, \*Neiswanger, et al., J. Chem. Phys. 2021

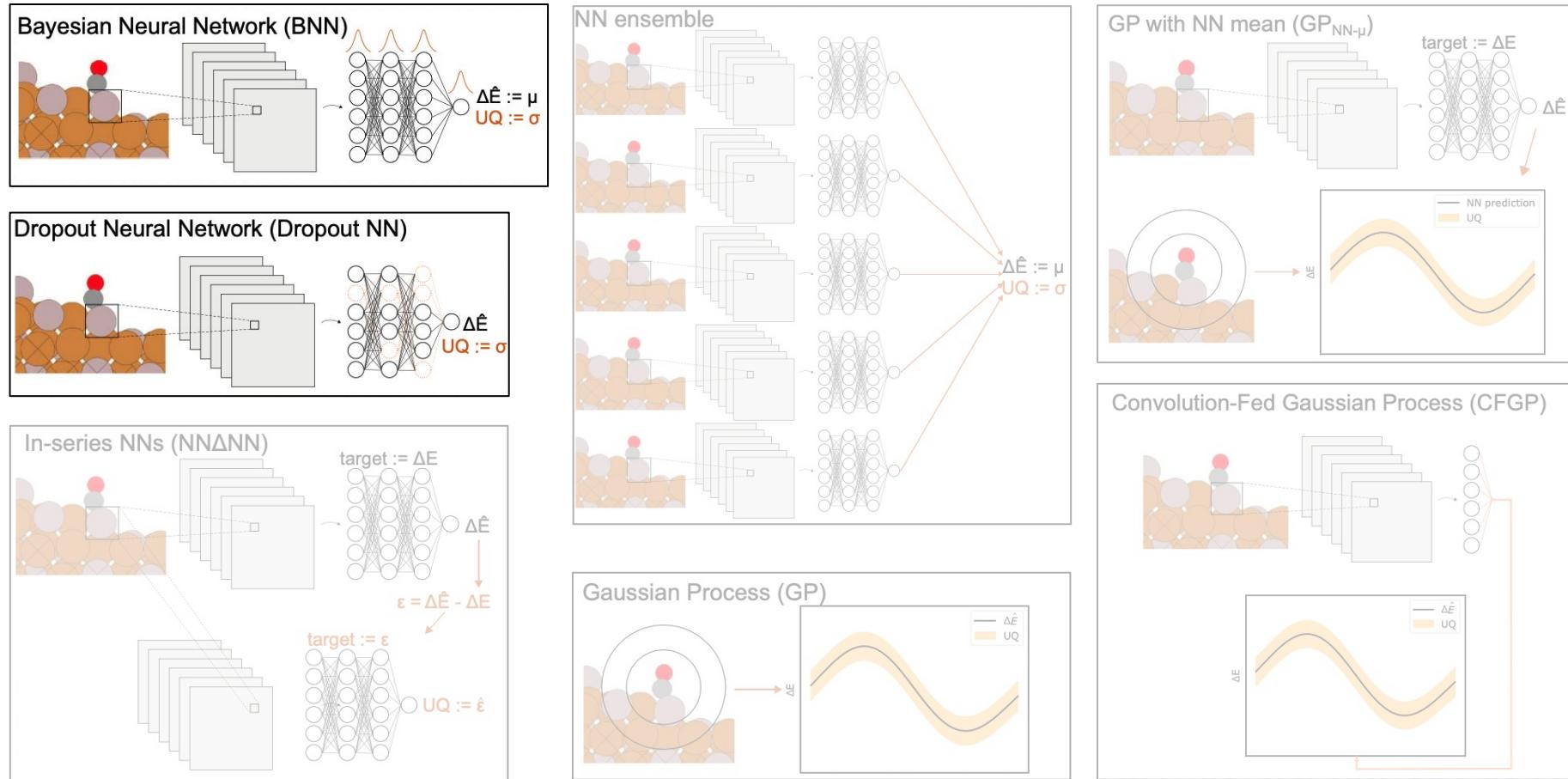
# Deep Uncertainty Models – Example: Computational Catalyst Design



"Methods for comparing uncertainty quantifications for material property predictions", \*Tran, \*Neiswanger, et al., MLST. 2020

"Computational catalyst discovery: Active classification through myopic multiscale sampling", \*Tran, \*Neiswanger, et al., J. Chem. Phys. 2021

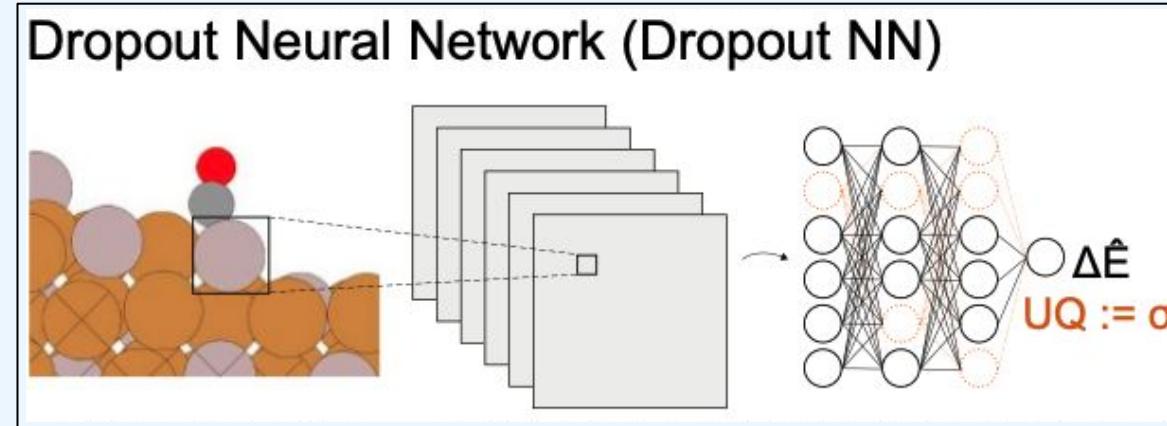
# Deep Uncertainty Models – Example: Computational Catalyst Design



"Methods for comparing uncertainty quantifications for material property predictions", \*Tran, \*Neiswanger, et al., MLST. 2020

"Computational catalyst discovery: Active classification through myopic multiscale sampling", \*Tran, \*Neiswanger, et al., J. Chem. Phys. 2021

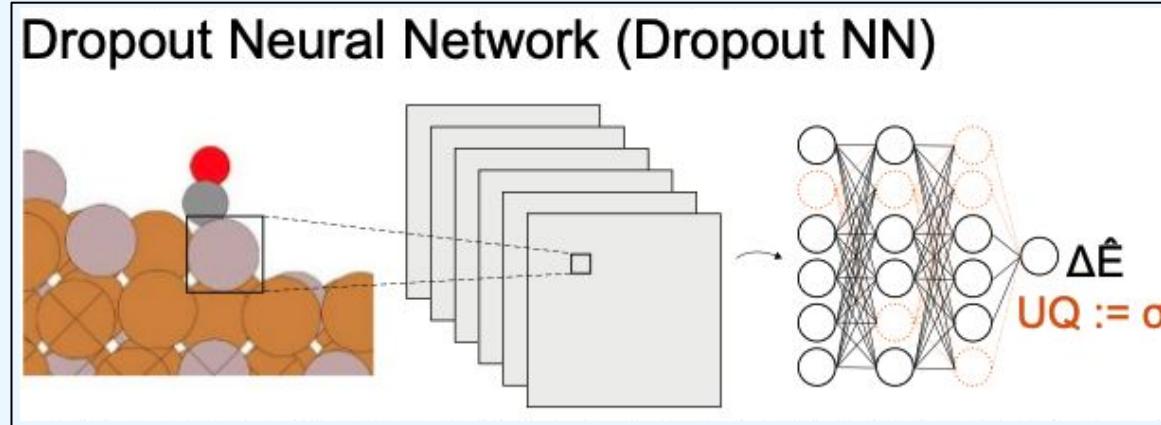
# Deep Uncertainty Models – Example: Computational Catalyst Design



"Methods for comparing uncertainty quantifications for material property predictions", \*Tran, \*Neiswanger, et al., MLST. 2020

"Computational catalyst discovery: Active classification through myopic multiscale sampling", \*Tran, \*Neiswanger, et al., J. Chem. Phys. 2021

# Deep Uncertainty Models – Example: Computational Catalyst Design

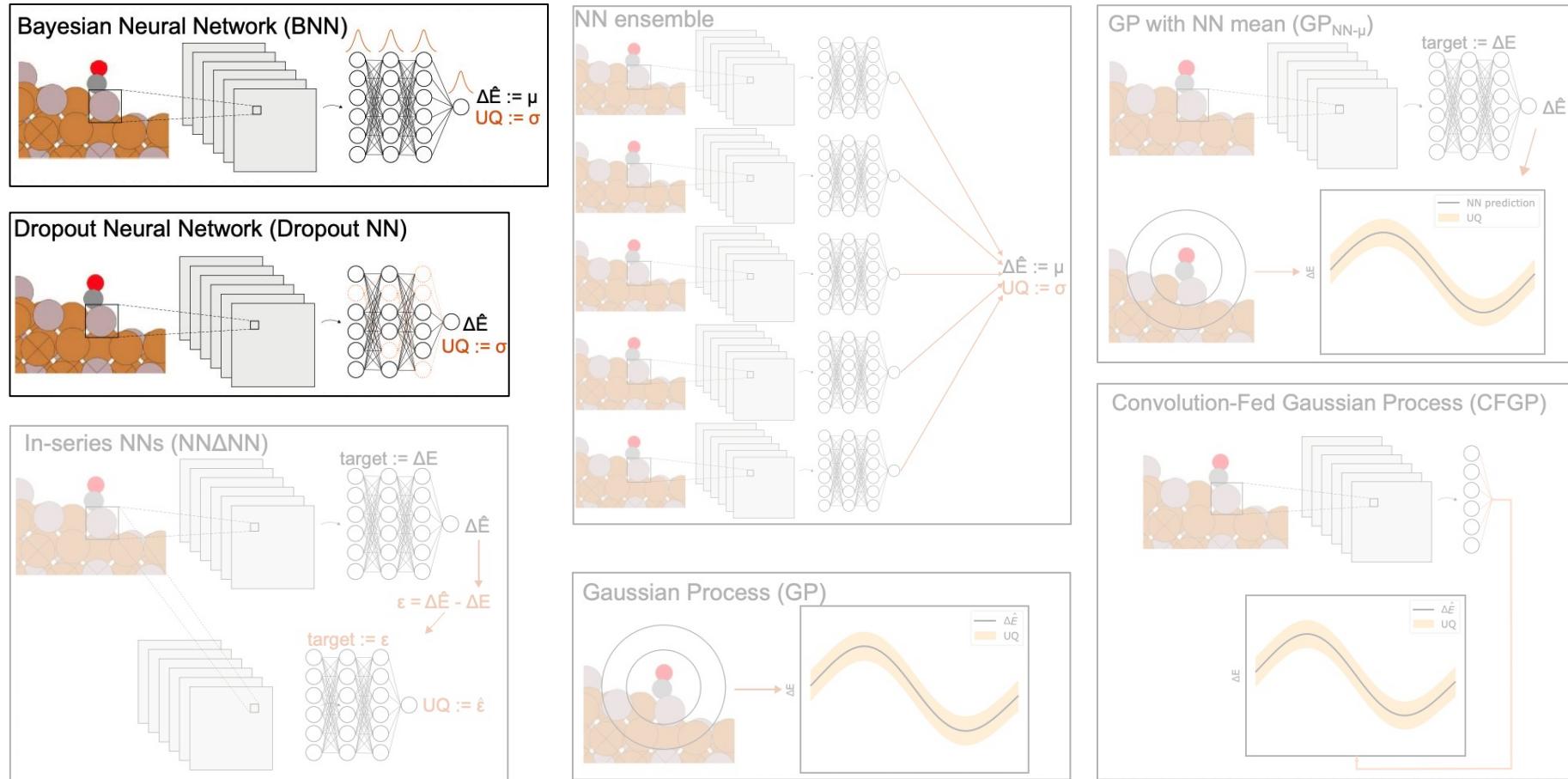


- During train/inference-time, perform *Monte Carlo Dropout*.
- (*I.e.*, set a random fraction of the neuron's outputs to zero during forward pass of the model).
- ⇒ A set of forward passes yields a set of samples from an output distribution.

"Methods for comparing uncertainty quantifications for material property predictions", \*Tran, \*Neiswanger, et al., MLST. 2020

"Computational catalyst discovery: Active classification through myopic multiscale sampling", \*Tran, \*Neiswanger, et al., J. Chem. Phys. 2021

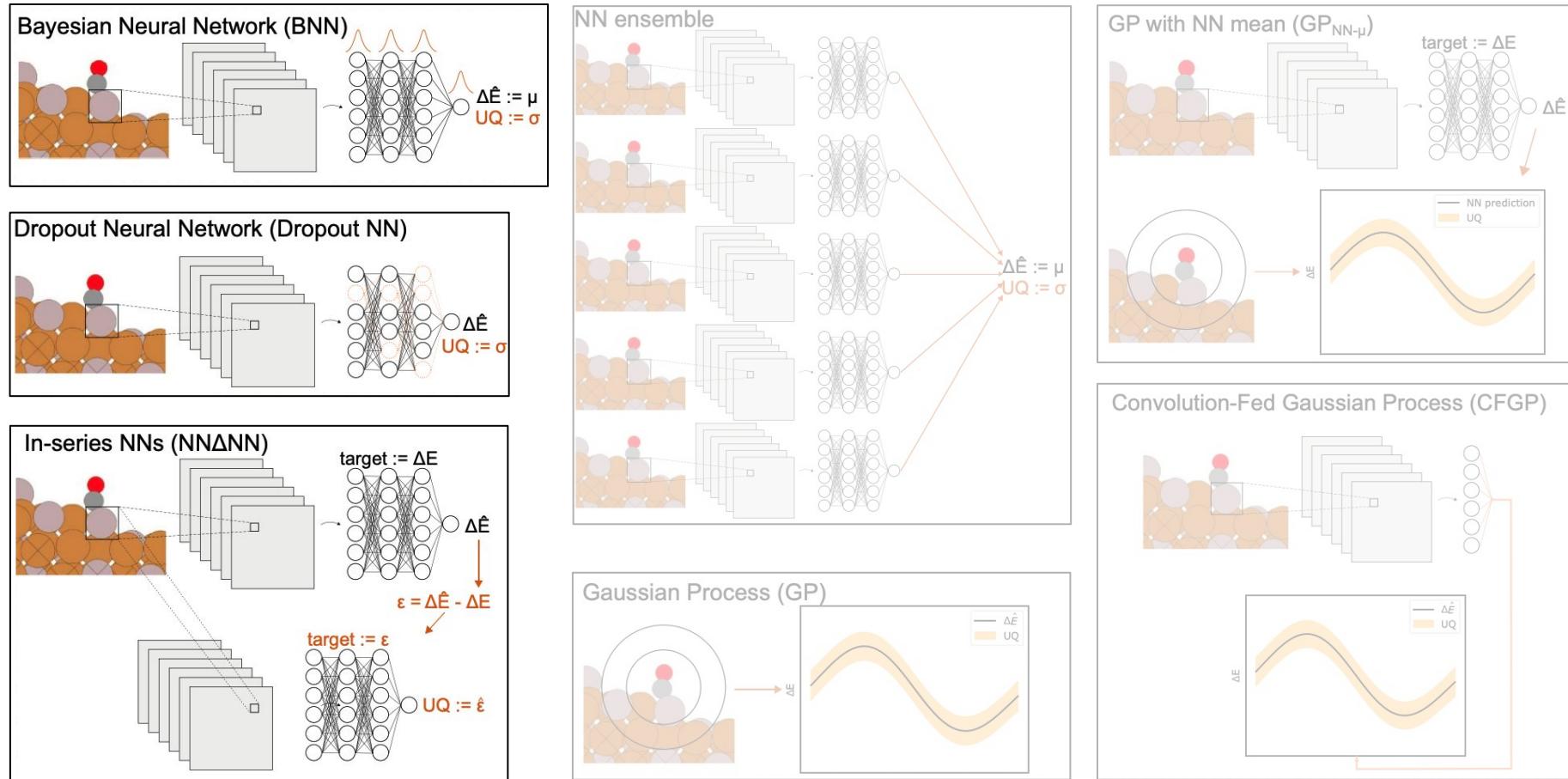
# Deep Uncertainty Models – Example: Computational Catalyst Design



"Methods for comparing uncertainty quantifications for material property predictions", \*Tran, \*Neiswanger, et al., MLST. 2020

"Computational catalyst discovery: Active classification through myopic multiscale sampling", \*Tran, \*Neiswanger, et al., J. Chem. Phys. 2021

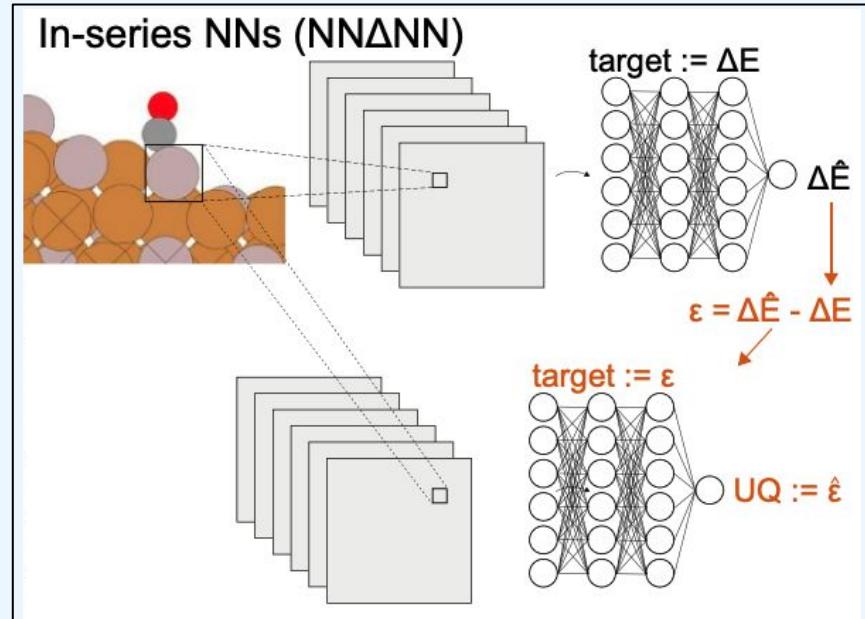
# Deep Uncertainty Models – Example: Computational Catalyst Design



"Methods for comparing uncertainty quantifications for material property predictions", \*Tran, \*Neiswanger, et al., MLST. 2020

"Computational catalyst discovery: Active classification through myopic multiscale sampling", \*Tran, \*Neiswanger, et al., J. Chem. Phys. 2021

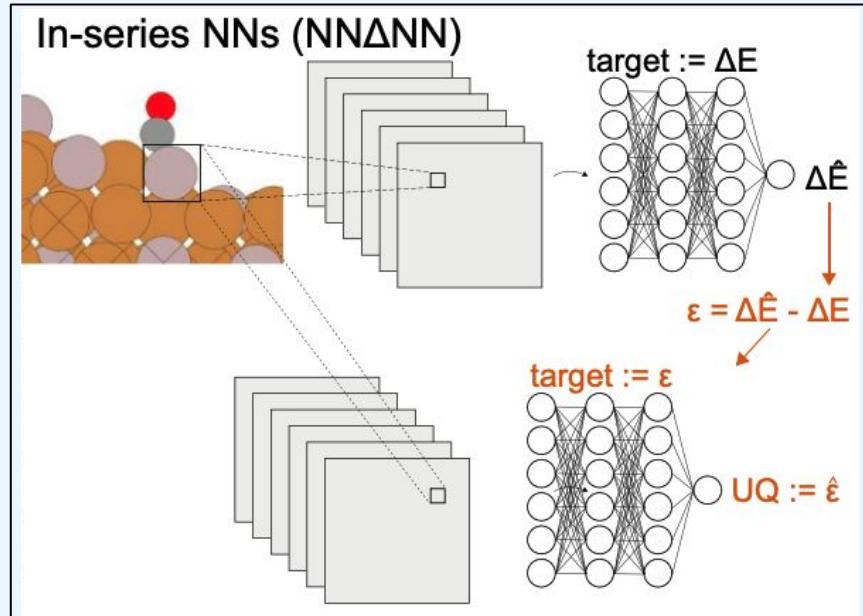
# Deep Uncertainty Models – Example: Computational Catalyst Design



"Methods for comparing uncertainty quantifications for material property predictions", \*Tran, \*Neiswanger, et al., MLST. 2020

"Computational catalyst discovery: Active classification through myopic multiscale sampling", \*Tran, \*Neiswanger, et al., J. Chem. Phys. 2021

# Deep Uncertainty Models – Example: Computational Catalyst Design

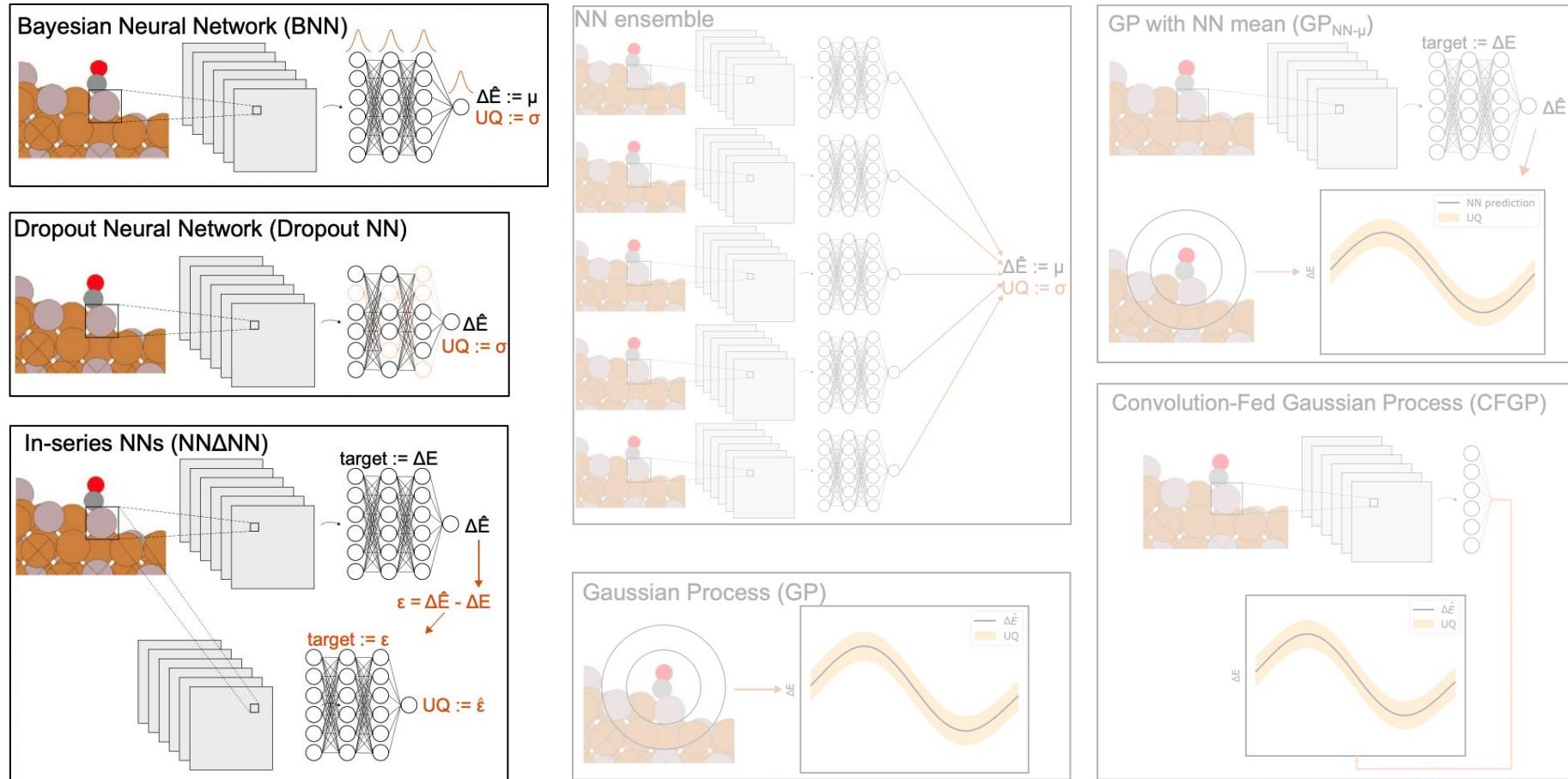


- First train a neural network for point prediction.
- Then compute residuals on a held-out set.
- Then use another neural network to try and predict residuals.
- ⇒ Second network predicts error, which can be used as UQ model.

"Methods for comparing uncertainty quantifications for material property predictions", \*Tran, \*Neiswanger, et al., MLST. 2020

"Computational catalyst discovery: Active classification through myopic multiscale sampling", \*Tran, \*Neiswanger, et al., J. Chem. Phys. 2021

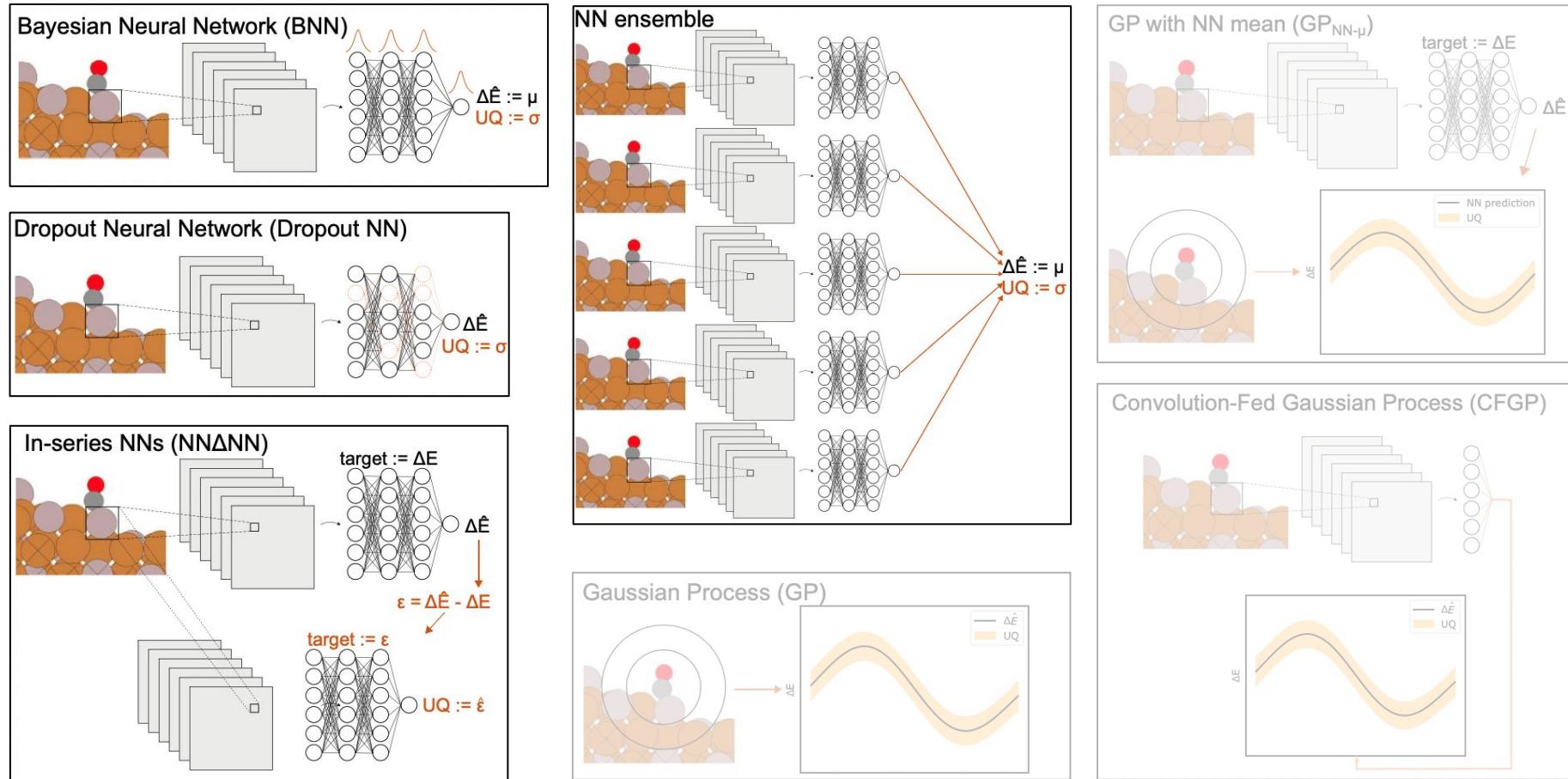
# Deep Uncertainty Models – Example: Computational Catalyst Design



"Methods for comparing uncertainty quantifications for material property predictions", \*Tran, \*Neiswanger, et al., MLST. 2020

"Computational catalyst discovery: Active classification through myopic multiscale sampling", \*Tran, \*Neiswanger, et al., J. Chem. Phys. 2021

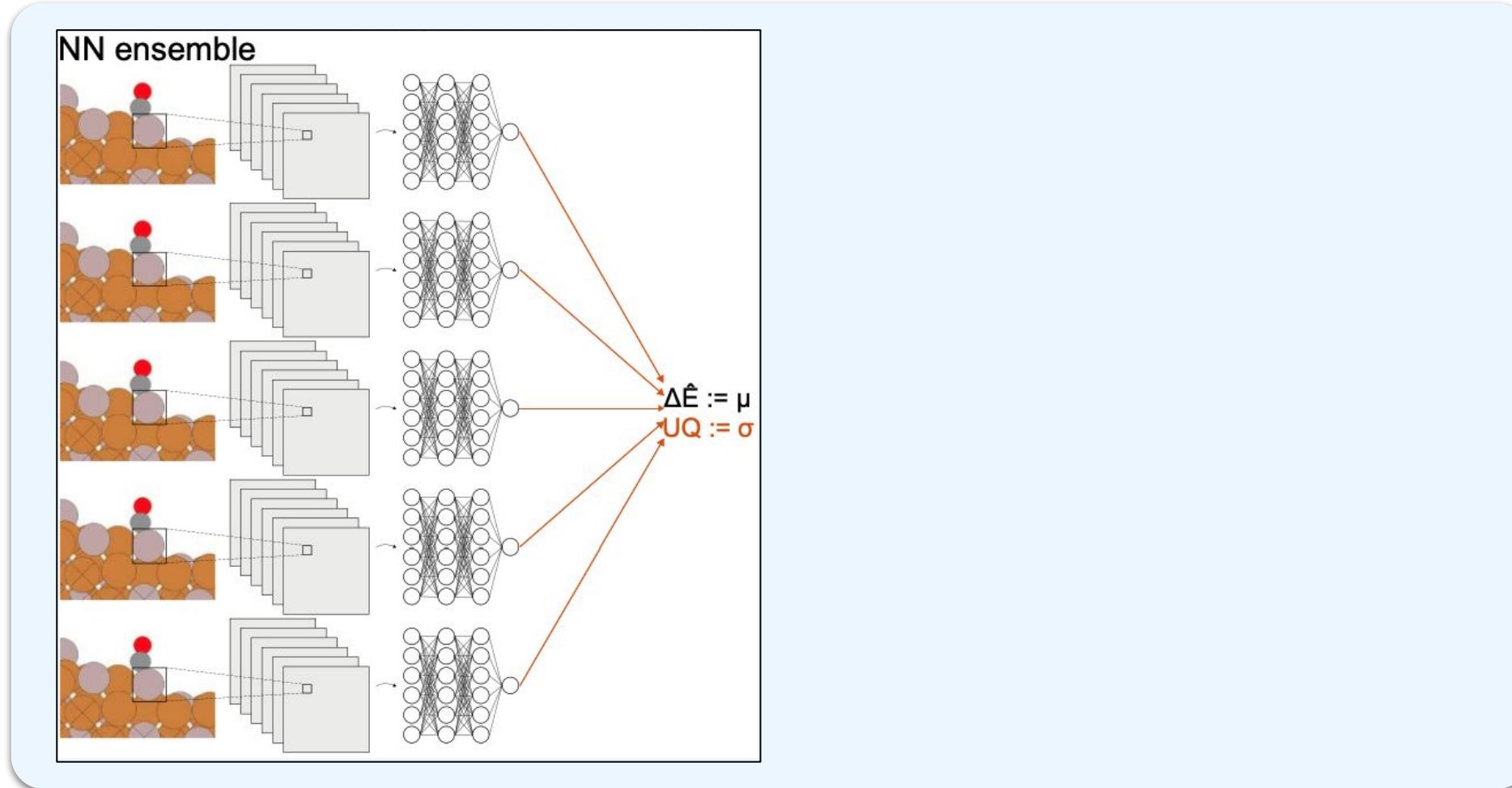
# Deep Uncertainty Models – Example: Computational Catalyst Design



"Methods for comparing uncertainty quantifications for material property predictions", \*Tran, \*Neiswanger, et al., MLST. 2020

"Computational catalyst discovery: Active classification through myopic multiscale sampling", \*Tran, \*Neiswanger, et al., J. Chem. Phys. 2021

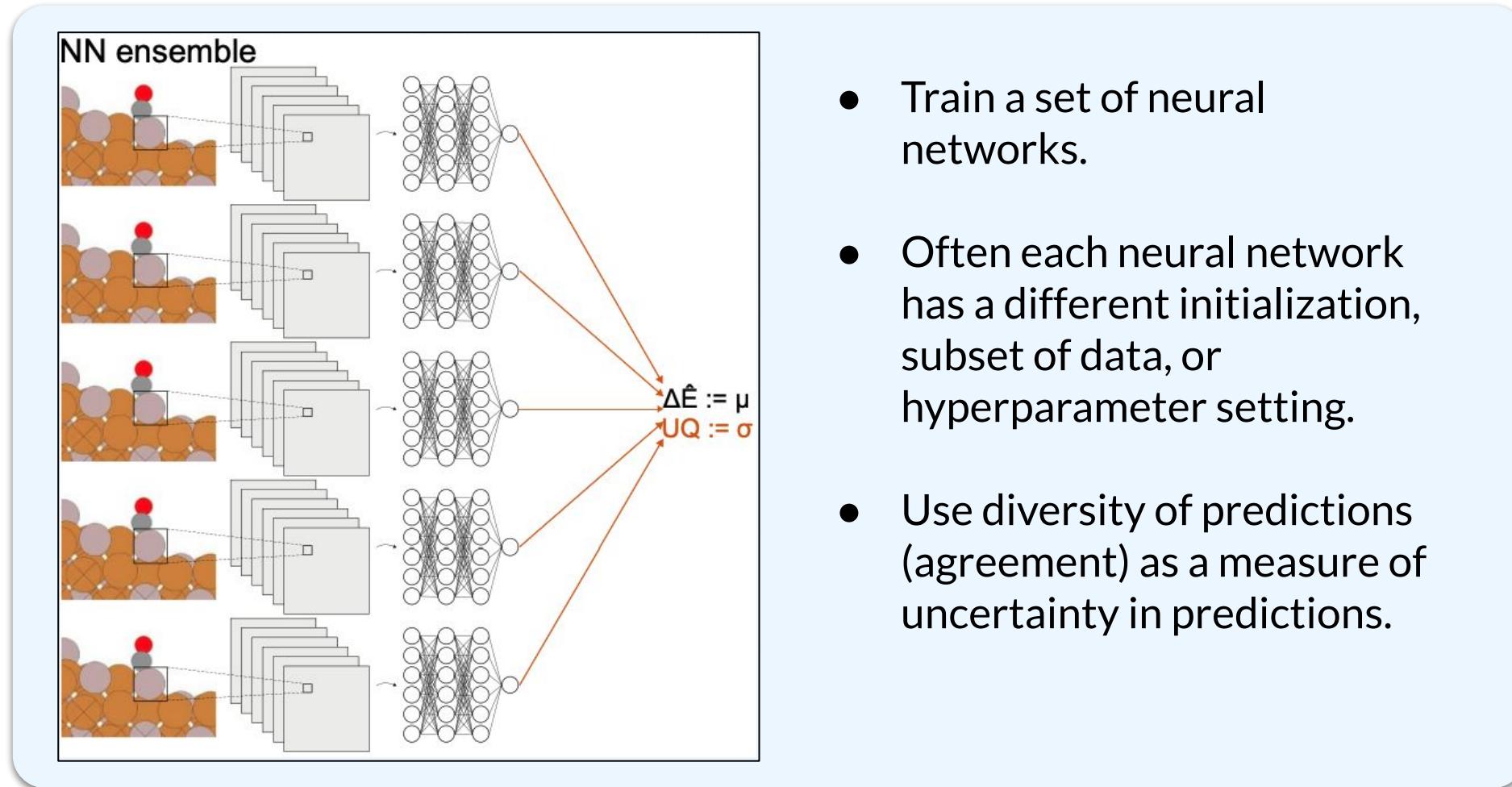
# Deep Uncertainty Models – Example: Computational Catalyst Design



"Methods for comparing uncertainty quantifications for material property predictions", \*Tran, \*Neiswanger, et al., MLST. 2020

"Computational catalyst discovery: Active classification through myopic multiscale sampling", \*Tran, \*Neiswanger, et al., J. Chem. Phys. 2021

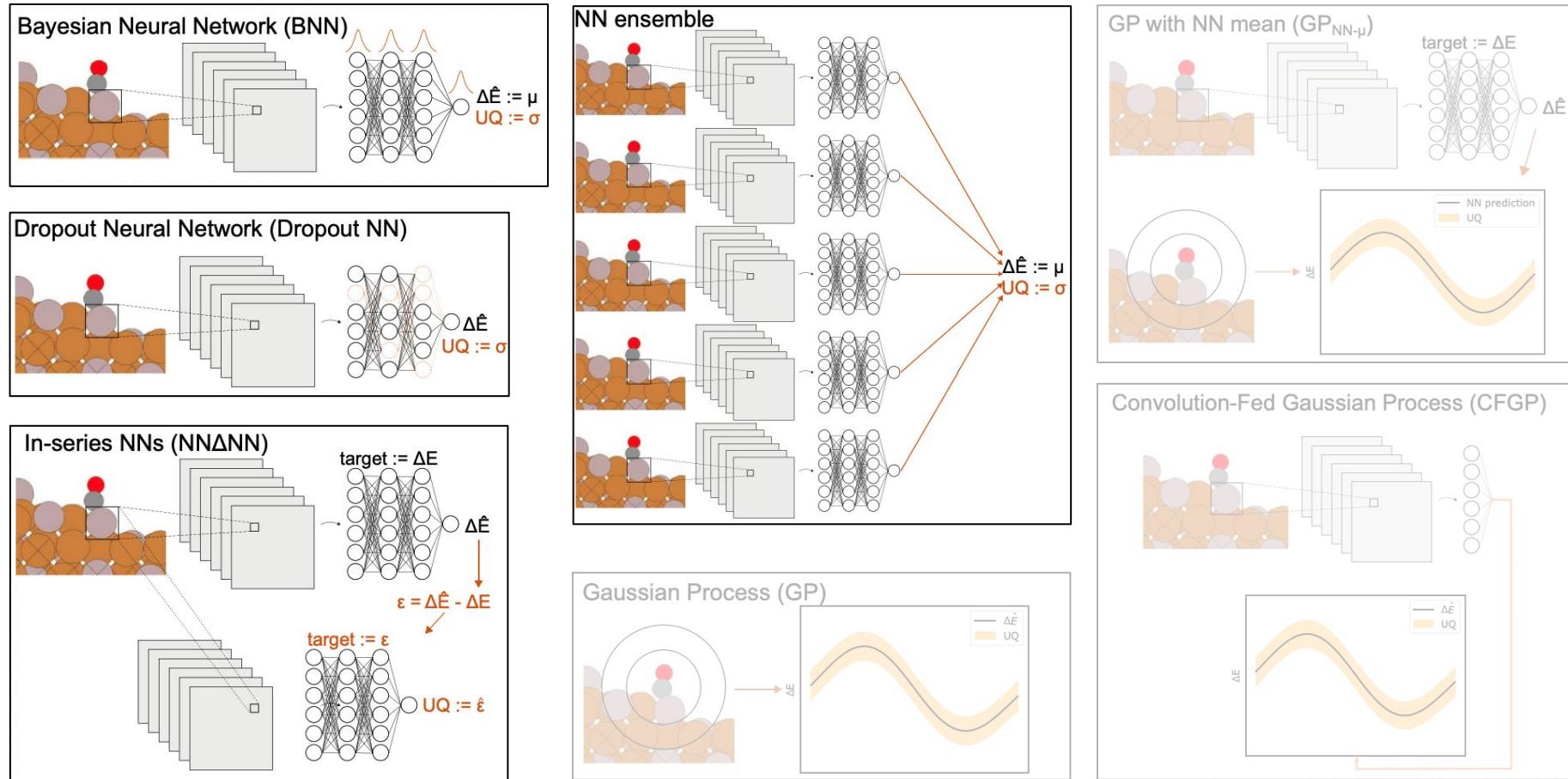
# Deep Uncertainty Models – Example: Computational Catalyst Design



"Methods for comparing uncertainty quantifications for material property predictions", \*Tran, \*Neiswanger, et al., MLST. 2020

"Computational catalyst discovery: Active classification through myopic multiscale sampling", \*Tran, \*Neiswanger, et al., J. Chem. Phys. 2021

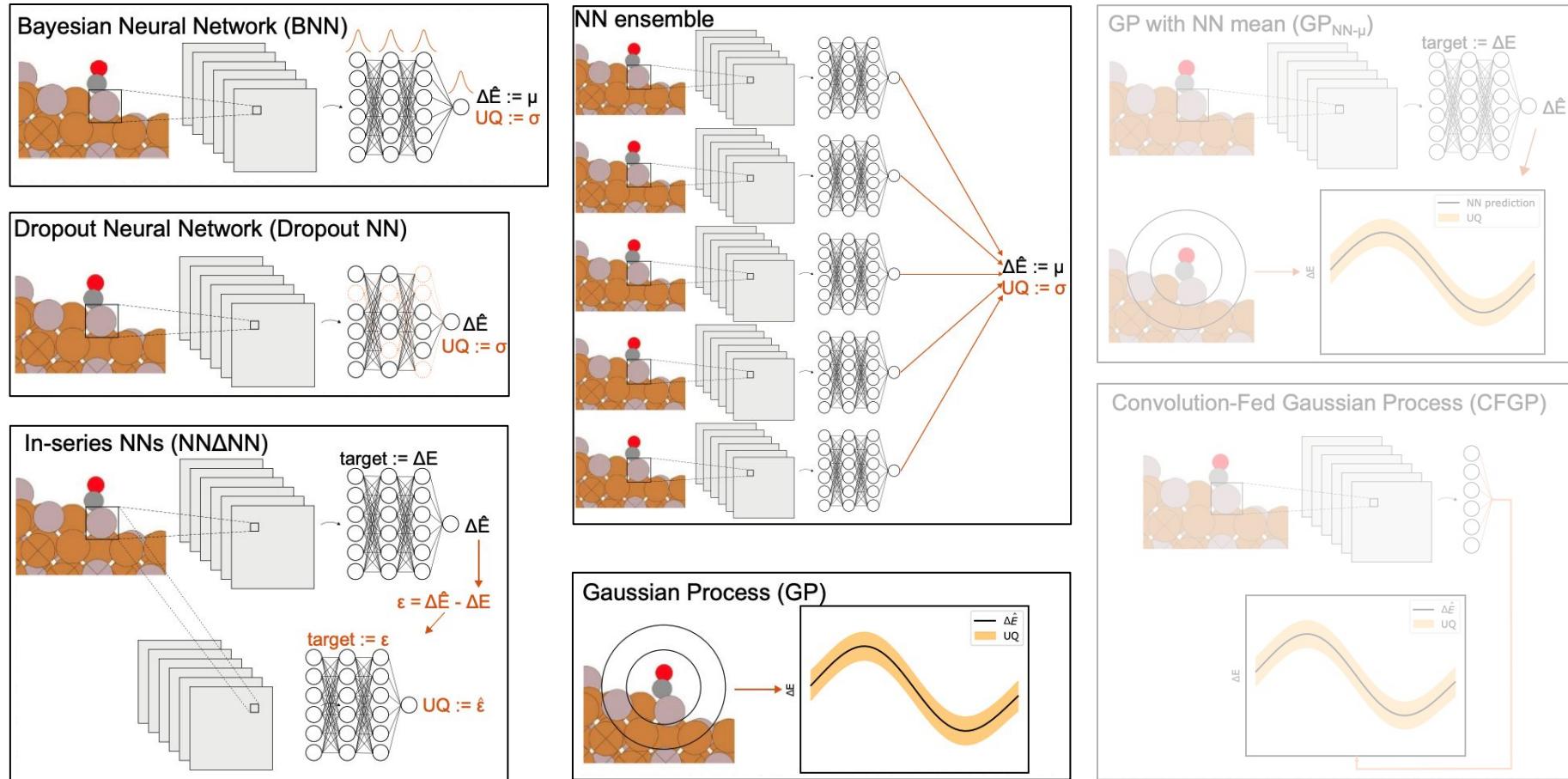
# Deep Uncertainty Models – Example: Computational Catalyst Design



"Methods for comparing uncertainty quantifications for material property predictions", \*Tran, \*Neiswanger, et al., MLST. 2020

"Computational catalyst discovery: Active classification through myopic multiscale sampling", \*Tran, \*Neiswanger, et al., J. Chem. Phys. 2021

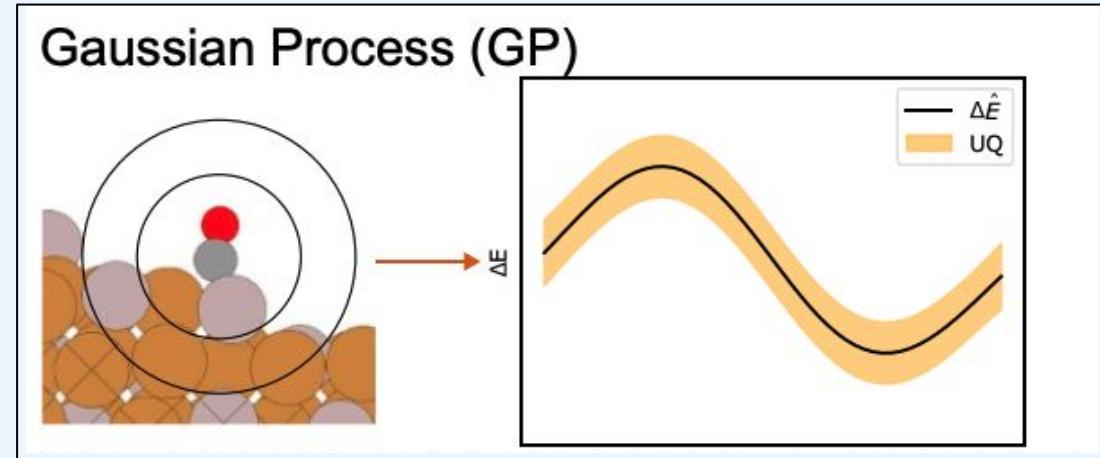
# Deep Uncertainty Models – Example: Computational Catalyst Design



"Methods for comparing uncertainty quantifications for material property predictions", \*Tran, \*Neiswanger, et al., MLST. 2020

"Computational catalyst discovery: Active classification through myopic multiscale sampling", \*Tran, \*Neiswanger, et al., J. Chem. Phys. 2021

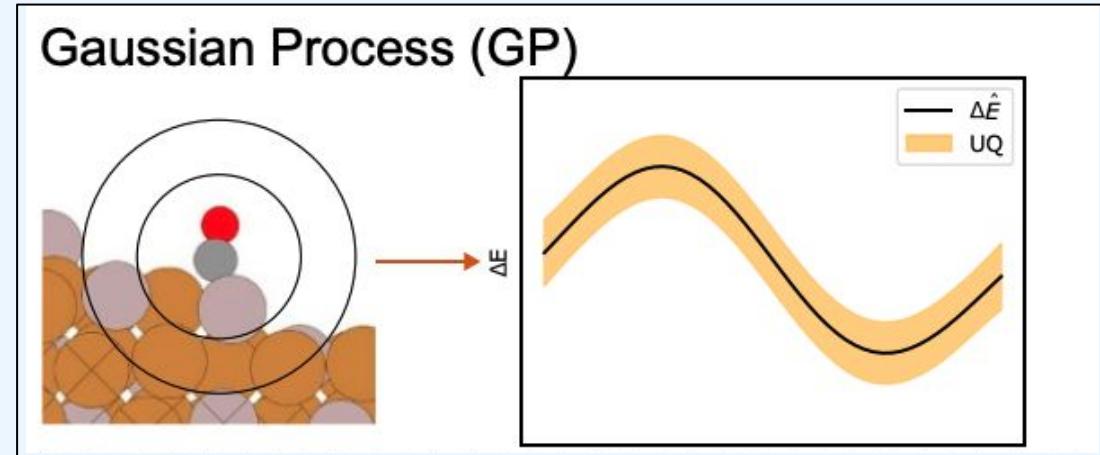
# Deep Uncertainty Models – Example: Computational Catalyst Design



"Methods for comparing uncertainty quantifications for material property predictions", \*Tran, \*Neiswanger, et al., MLST. 2020

"Computational catalyst discovery: Active classification through myopic multiscale sampling", \*Tran, \*Neiswanger, et al., J. Chem. Phys. 2021

# Deep Uncertainty Models – Example: Computational Catalyst Design

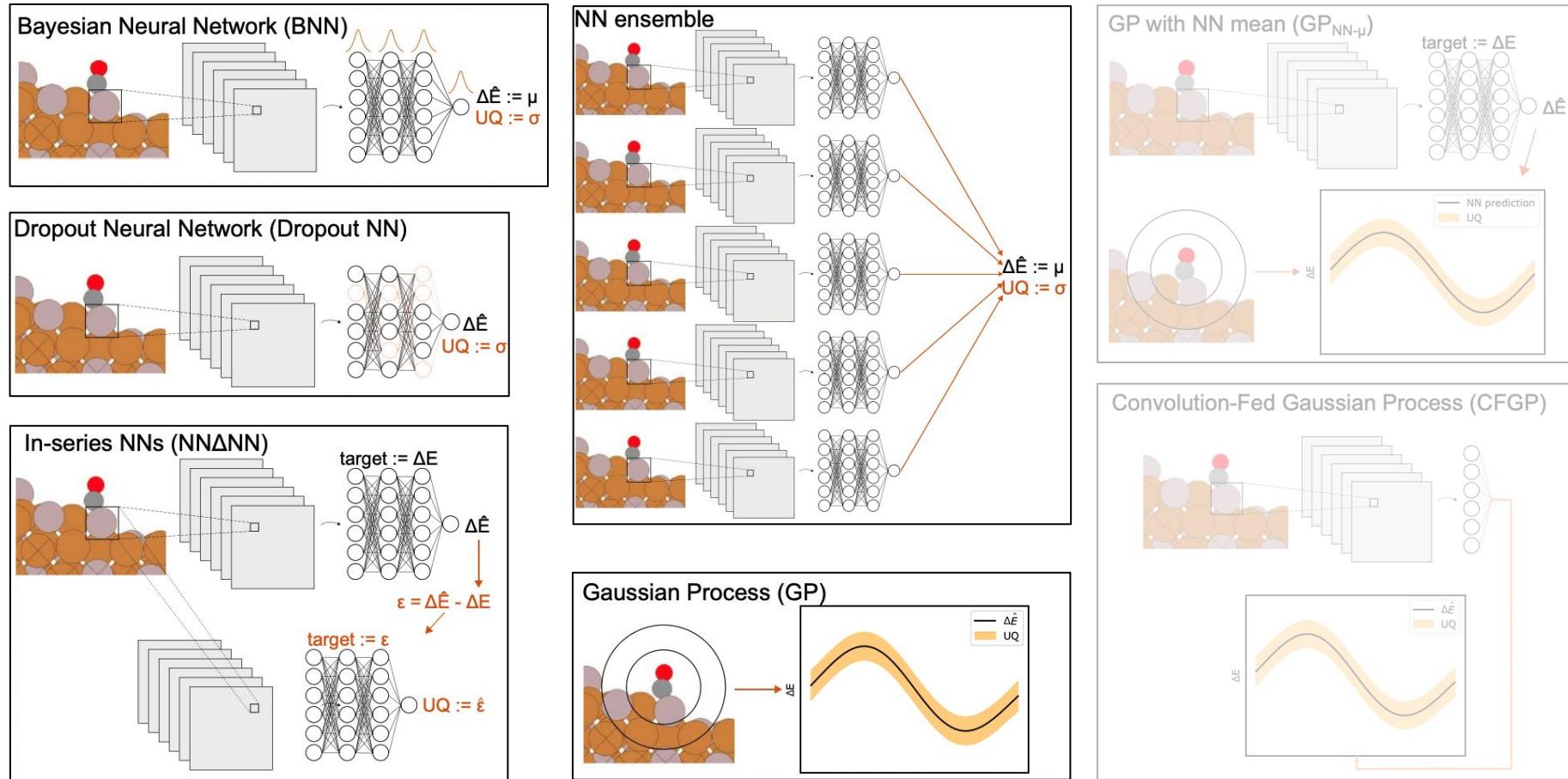


- Define some appropriate mean and kernel function on graph inputs.
- This yields a valid GP, which can make predictions on new inputs (new graphs)
- *Difficulty:* hard to hand-design good kernel function, & doesn't use NN.

“Methods for comparing uncertainty quantifications for material property predictions”, \*Tran, \*Neiswanger, et al., MLST. 2020

“Computational catalyst discovery: Active classification through myopic multiscale sampling”, \*Tran, \*Neiswanger, et al., J. Chem. Phys. 2021

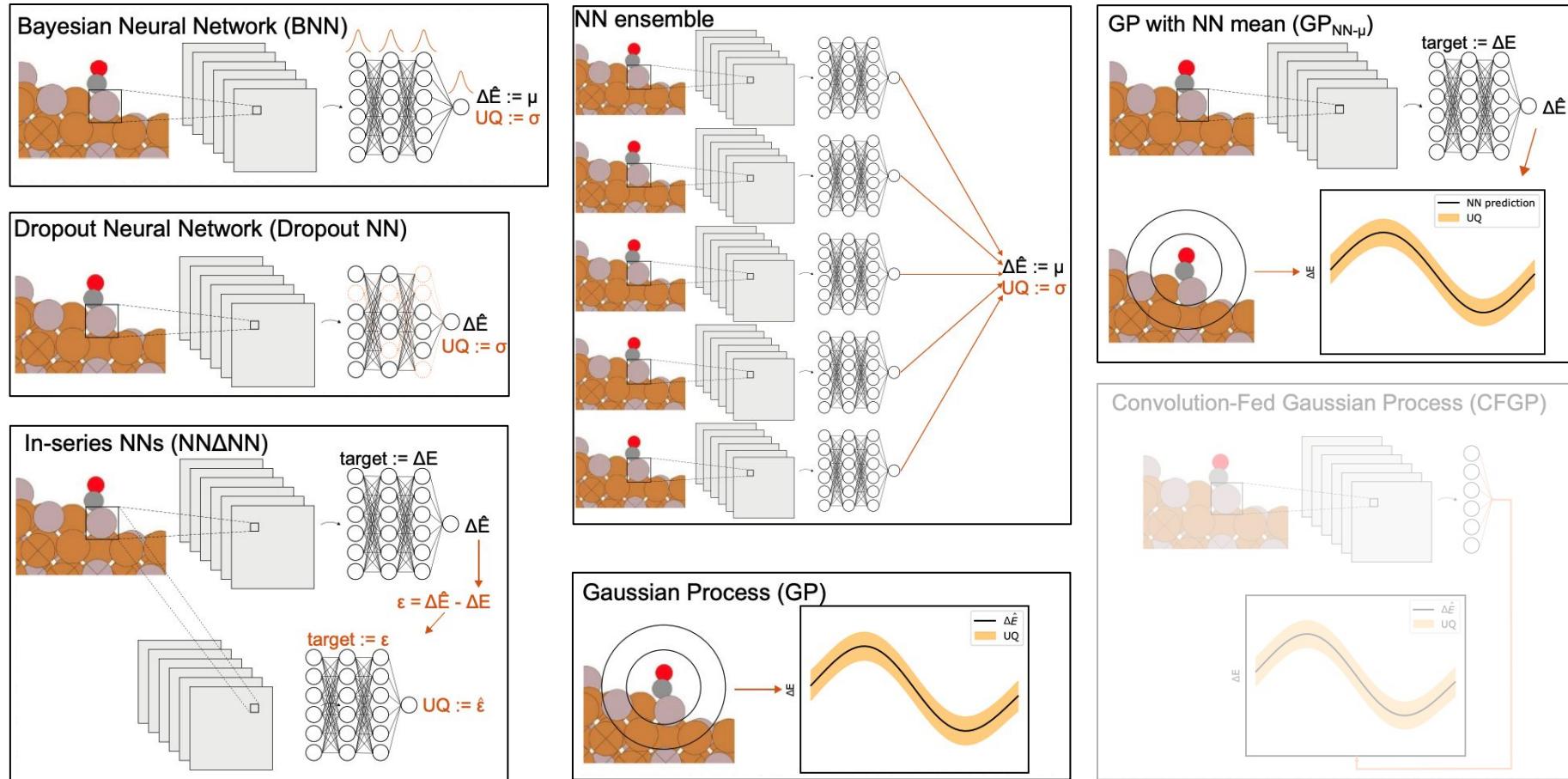
# Deep Uncertainty Models – Example: Computational Catalyst Design



"Methods for comparing uncertainty quantifications for material property predictions", \*Tran, \*Neiswanger, et al., MLST. 2020

"Computational catalyst discovery: Active classification through myopic multiscale sampling", \*Tran, \*Neiswanger, et al., J. Chem. Phys. 2021

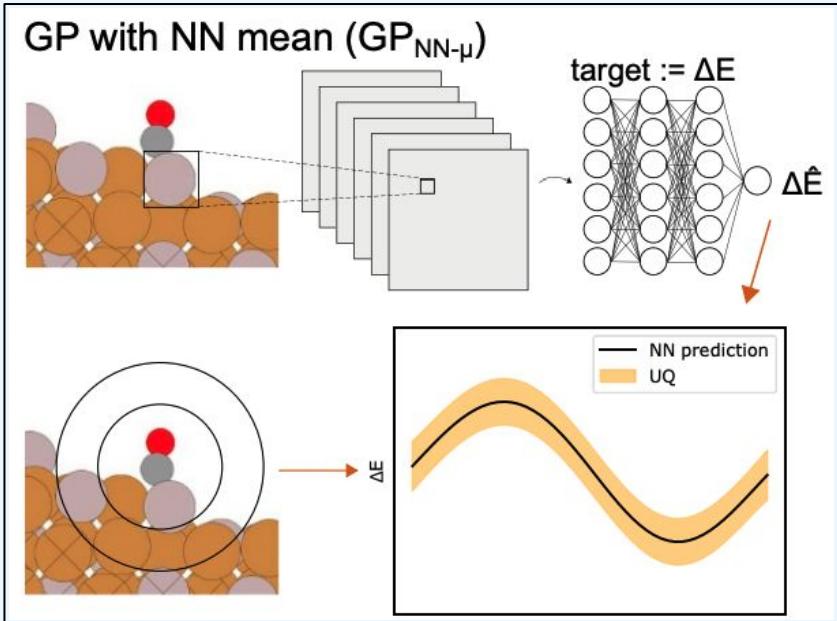
# Deep Uncertainty Models – Example: Computational Catalyst Design



"Methods for comparing uncertainty quantifications for material property predictions", \*Tran, \*Neiswanger, et al., MLST. 2020

"Computational catalyst discovery: Active classification through myopic multiscale sampling", \*Tran, \*Neiswanger, et al., J. Chem. Phys. 2021

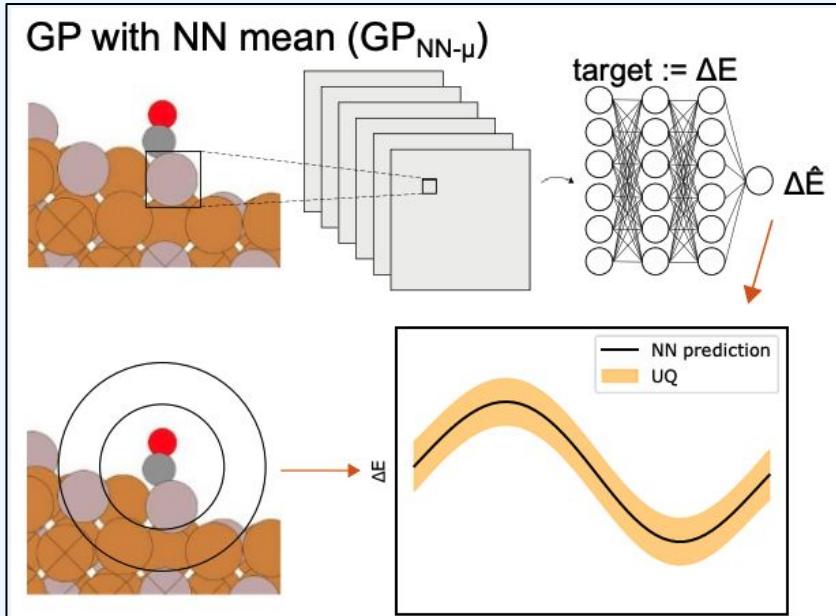
# Deep Uncertainty Models – Example: Computational Catalyst Design



"Methods for comparing uncertainty quantifications for material property predictions", \*Tran, \*Neiswanger, et al., MLST. 2020

"Computational catalyst discovery: Active classification through myopic multiscale sampling", \*Tran, \*Neiswanger, et al., J. Chem. Phys. 2021

# Deep Uncertainty Models – Example: Computational Catalyst Design

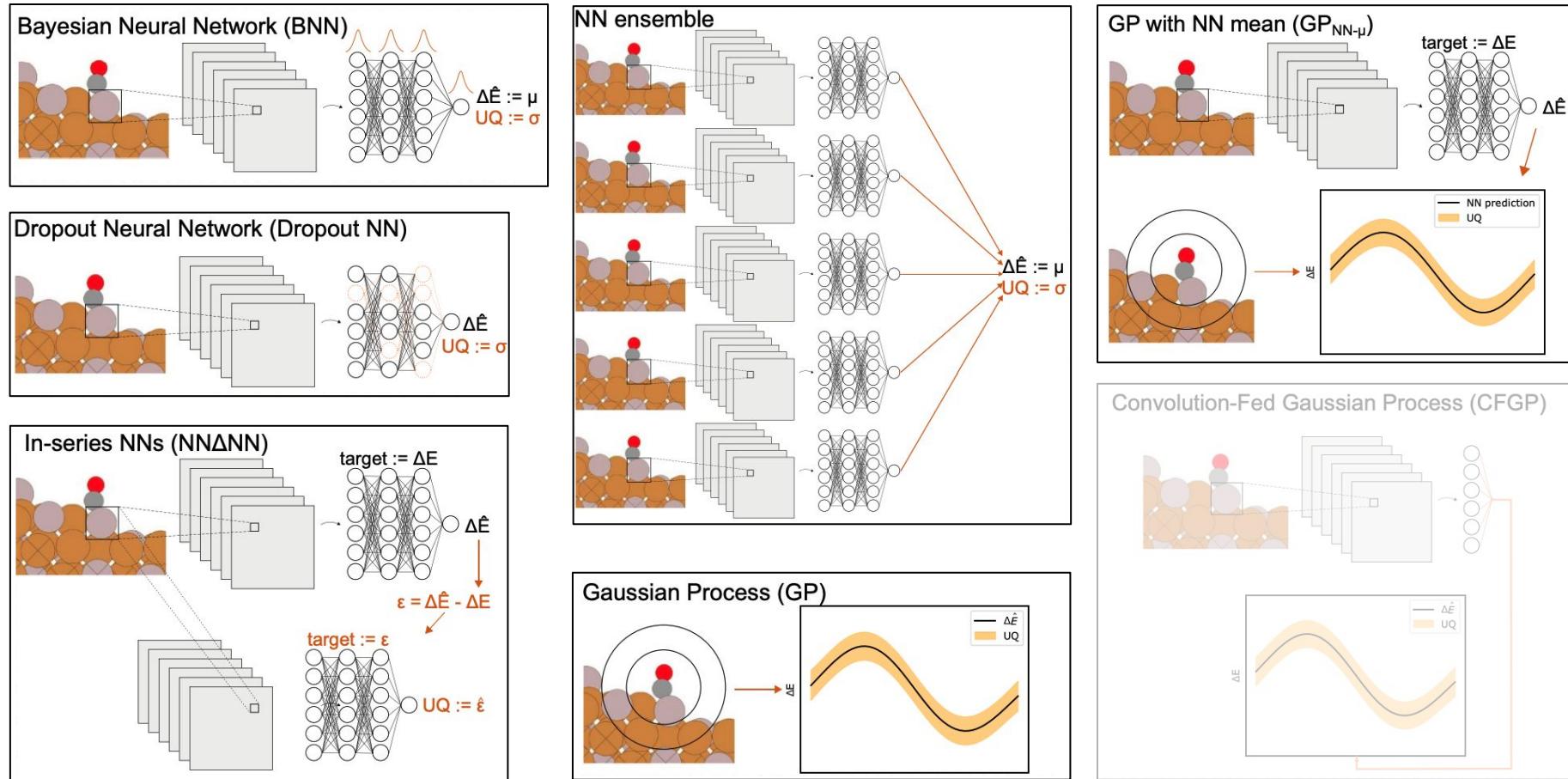


- In GPs, the prior mean fn is usually set to be **0**.
- But it could be anything...
- Here, we can set the prior mean function to be the output of a neural network model!
- Downside: still need to define a kernel function on graph inputs.

"Methods for comparing uncertainty quantifications for material property predictions", \*Tran, \*Neiswanger, et al., MLST. 2020

"Computational catalyst discovery: Active classification through myopic multiscale sampling", \*Tran, \*Neiswanger, et al., J. Chem. Phys. 2021

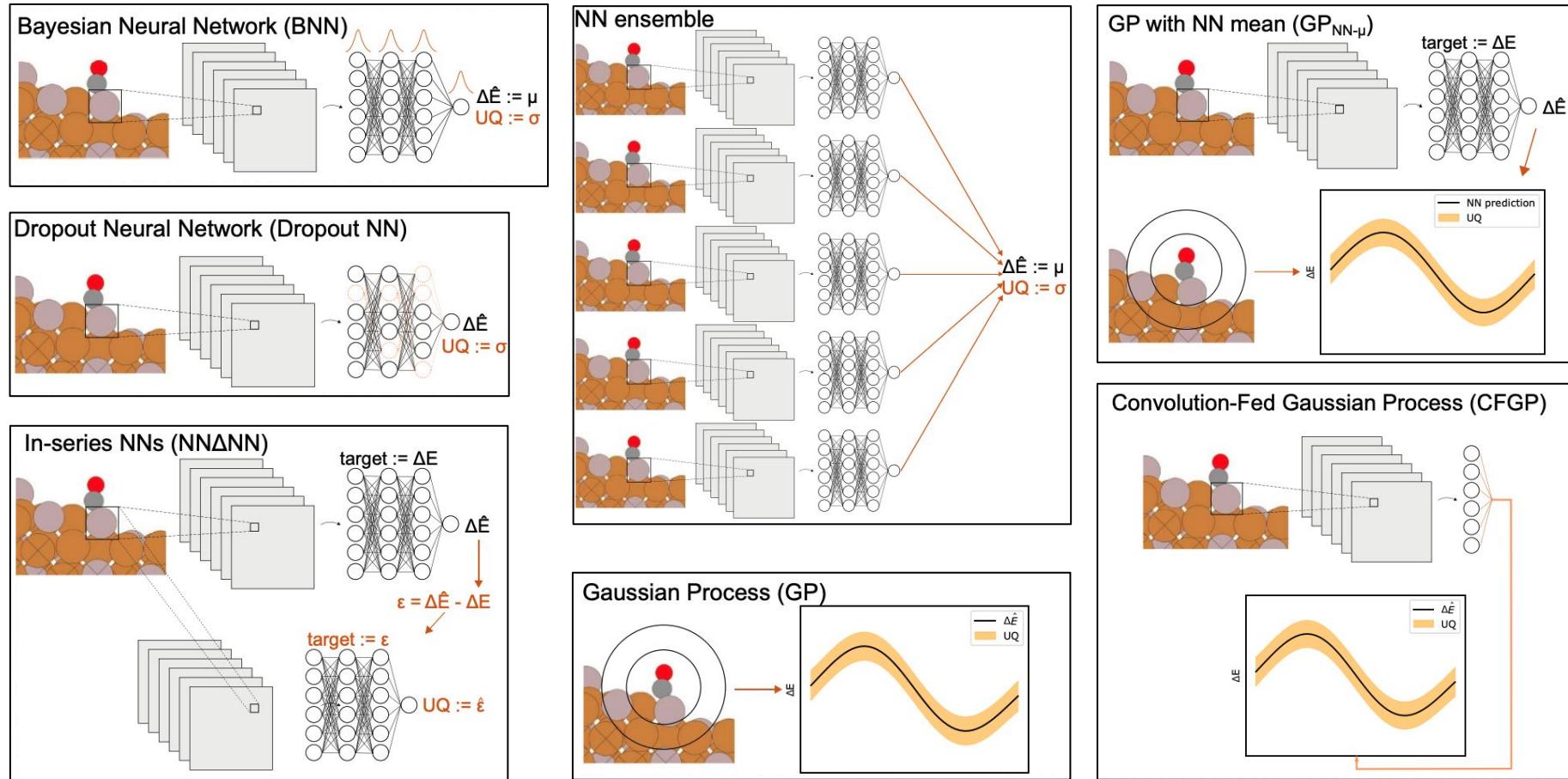
# Deep Uncertainty Models – Example: Computational Catalyst Design



"Methods for comparing uncertainty quantifications for material property predictions", \*Tran, \*Neiswanger, et al., MLST. 2020

"Computational catalyst discovery: Active classification through myopic multiscale sampling", \*Tran, \*Neiswanger, et al., J. Chem. Phys. 2021

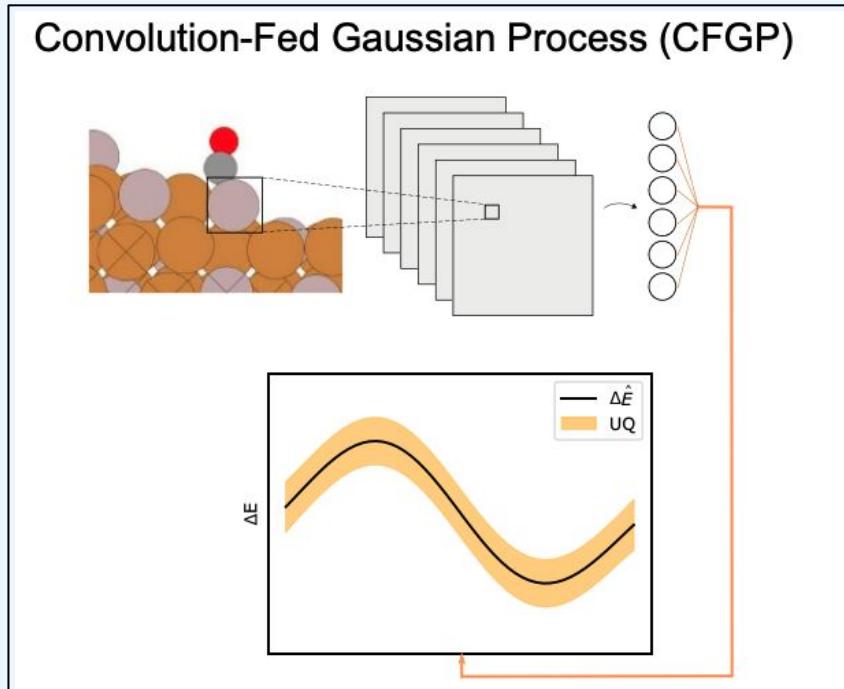
# Deep Uncertainty Models – Example: Computational Catalyst Design



"Methods for comparing uncertainty quantifications for material property predictions", \*Tran, \*Neiswanger, et al., MLST. 2020

"Computational catalyst discovery: Active classification through myopic multiscale sampling", \*Tran, \*Neiswanger, et al., J. Chem. Phys. 2021

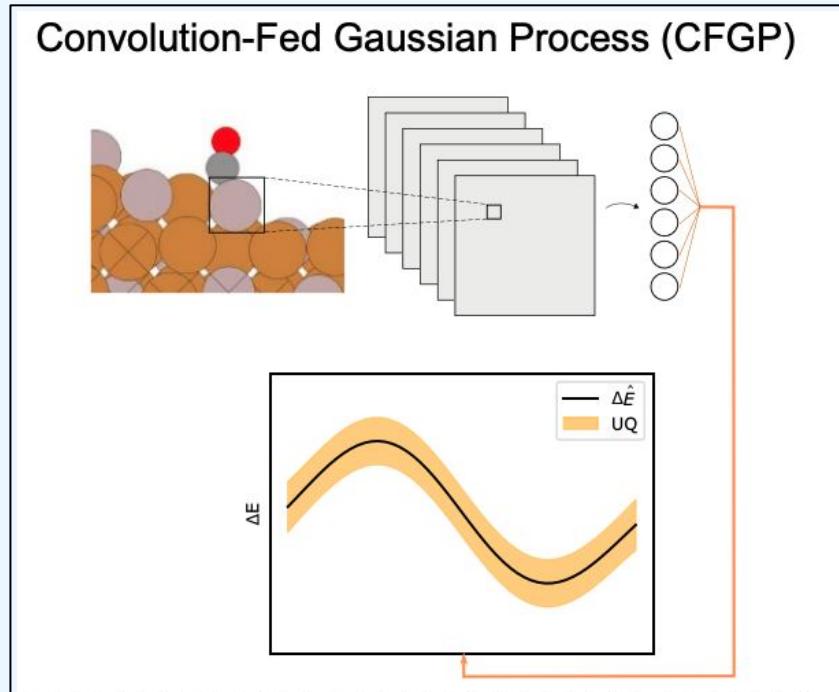
# Deep Uncertainty Models – Example: Computational Catalyst Design



"Methods for comparing uncertainty quantifications for material property predictions", \*Tran, \*Neiswanger, et al., MLST. 2020

"Computational catalyst discovery: Active classification through myopic multiscale sampling", \*Tran, \*Neiswanger, et al., J. Chem. Phys. 2021

# Deep Uncertainty Models – Example: Computational Catalyst Design

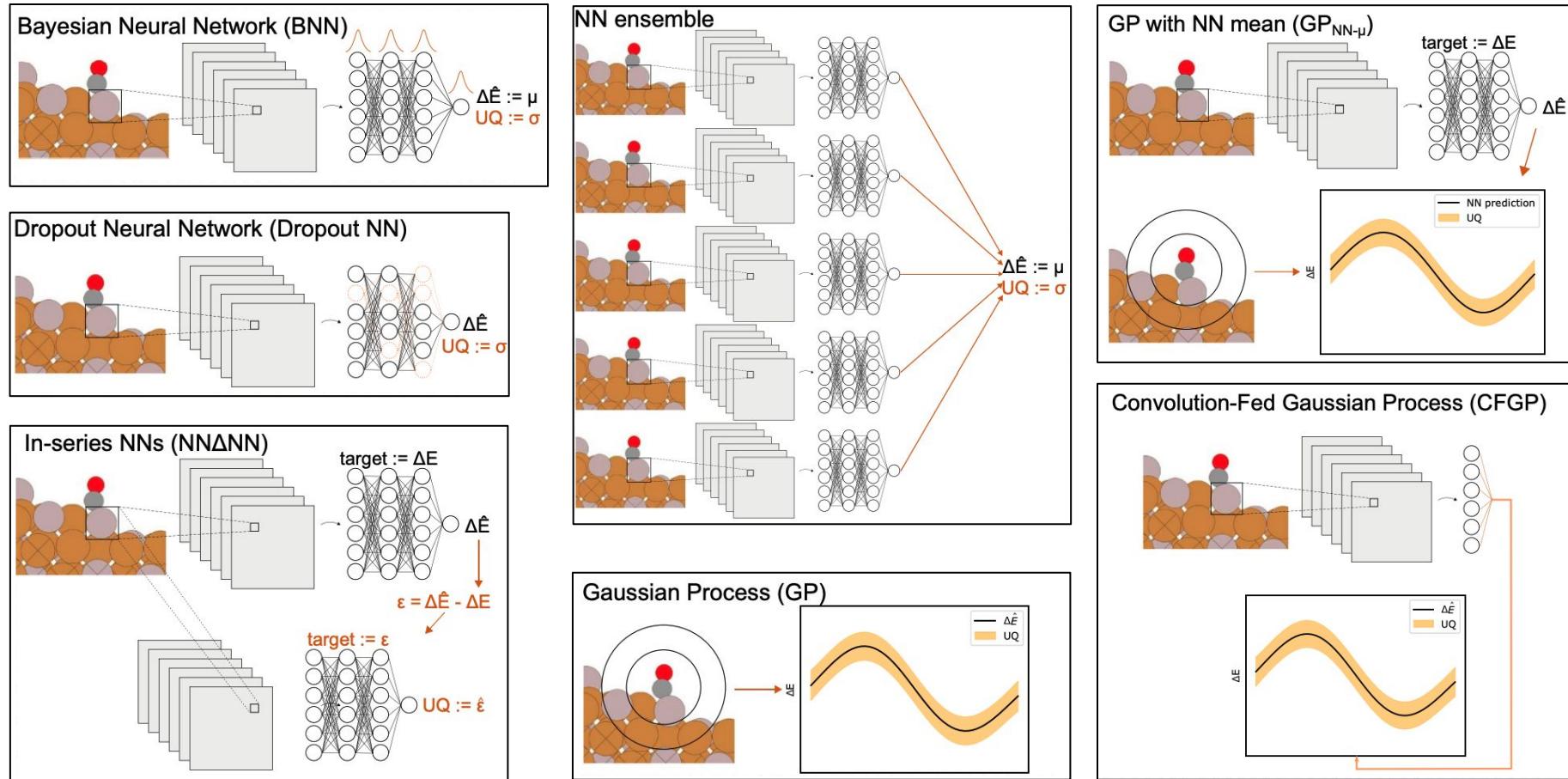


- Here, we use the neural network to learn good representations (e.g., for a GP).
- For any input, first make a prediction w/ neural network.
- Use output at an intermediate layer as new transformed input (into a continuous space).
- Then define a Gaussian process on this new representation of inputs.

"Methods for comparing uncertainty quantifications for material property predictions", \*Tran, \*Neiswanger, et al., MLST. 2020

"Computational catalyst discovery: Active classification through myopic multiscale sampling", \*Tran, \*Neiswanger, et al., J. Chem. Phys. 2021

# Deep Uncertainty Models – Example: Computational Catalyst Design



"Methods for comparing uncertainty quantifications for material property predictions", \*Tran, \*Neiswanger, et al., MLST. 2020

"Computational catalyst discovery: Active classification through myopic multiscale sampling", \*Tran, \*Neiswanger, et al., J. Chem. Phys. 2021

# **Bayesian Optimization**

# Expensive Black-box Functions

# Expensive Black-box Functions

Many real-world systems can be modeled as an **expensive black-box function  $f$** .



# Expensive Black-box Functions

Many real-world systems can be modeled as an **expensive black-box function  $f$** .

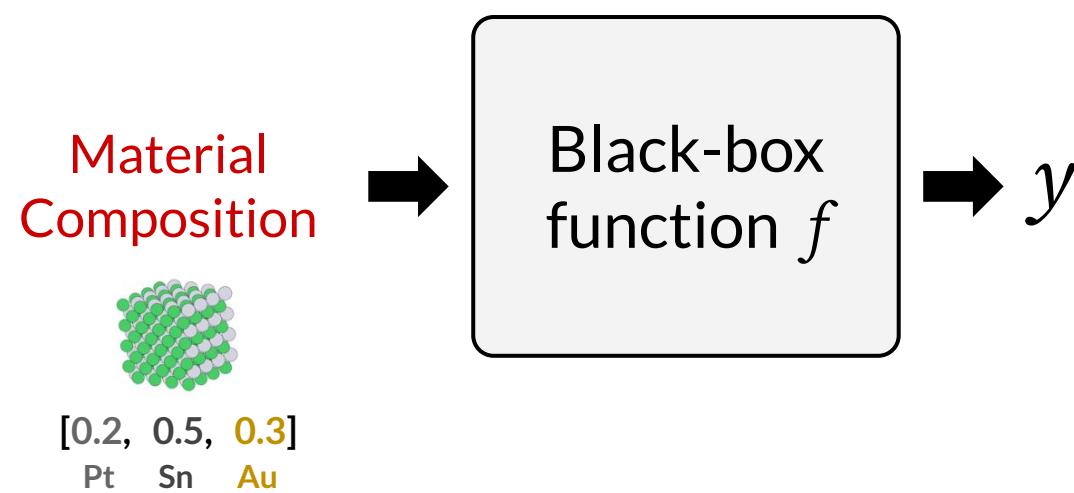
MATERIALS  
SCIENCE



# Expensive Black-box Functions

Many real-world systems can be modeled as an **expensive black-box function**  $f$ .

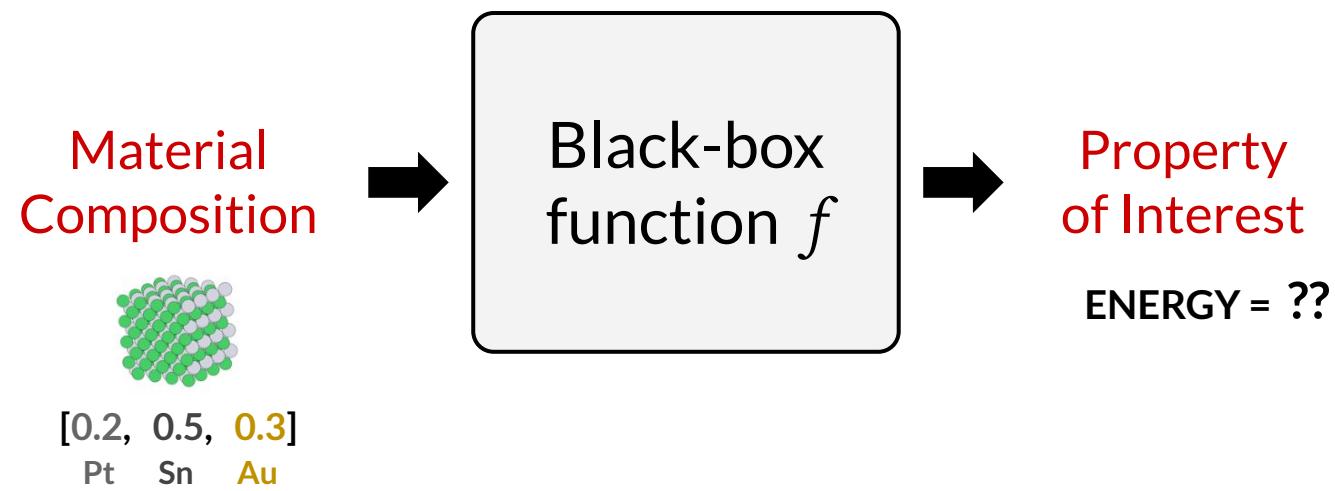
MATERIALS  
SCIENCE



# Expensive Black-box Functions

Many real-world systems can be modeled as an **expensive black-box function  $f$** .

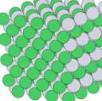
MATERIALS  
SCIENCE

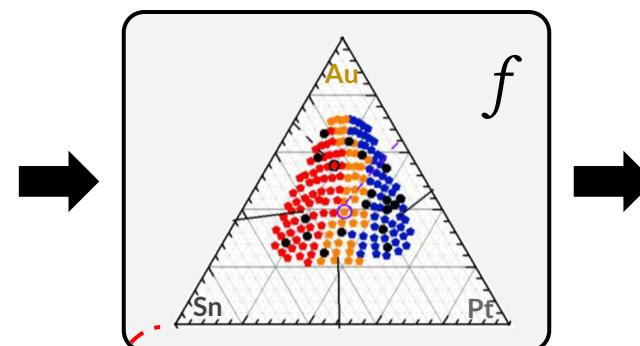


# Expensive Black-box Functions

Many real-world systems can be modeled as an **expensive black-box function**  $f$ .

MATERIALS  
SCIENCE

Material  
Composition  
  
[0.2, 0.5, 0.3]  
Pt Sn Au



Measured  
Property  
ENERGY = 5.5

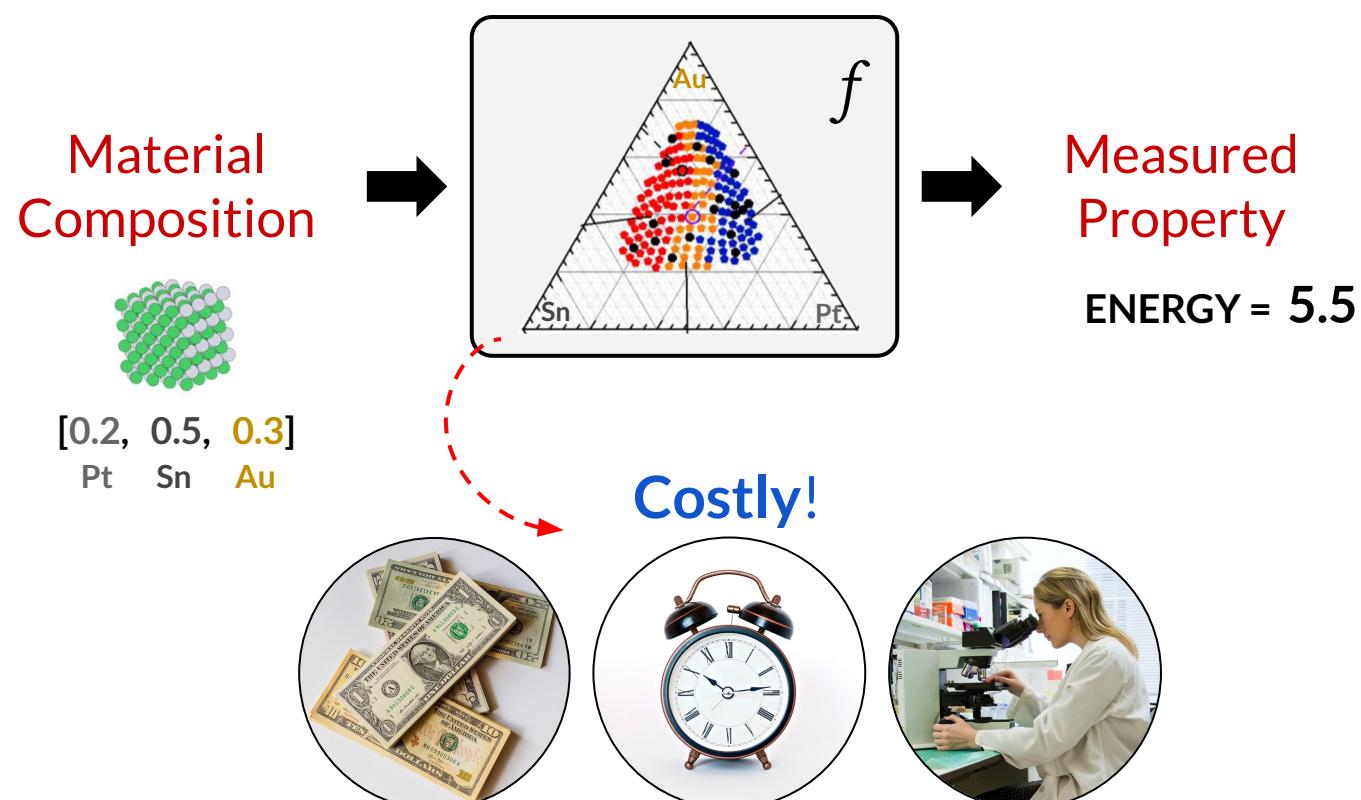
QUERY:  
Synthesize  
and Measure  
Property



# Expensive Black-box Functions

Many real-world systems can be modeled as an **expensive black-box function**  $f$ .

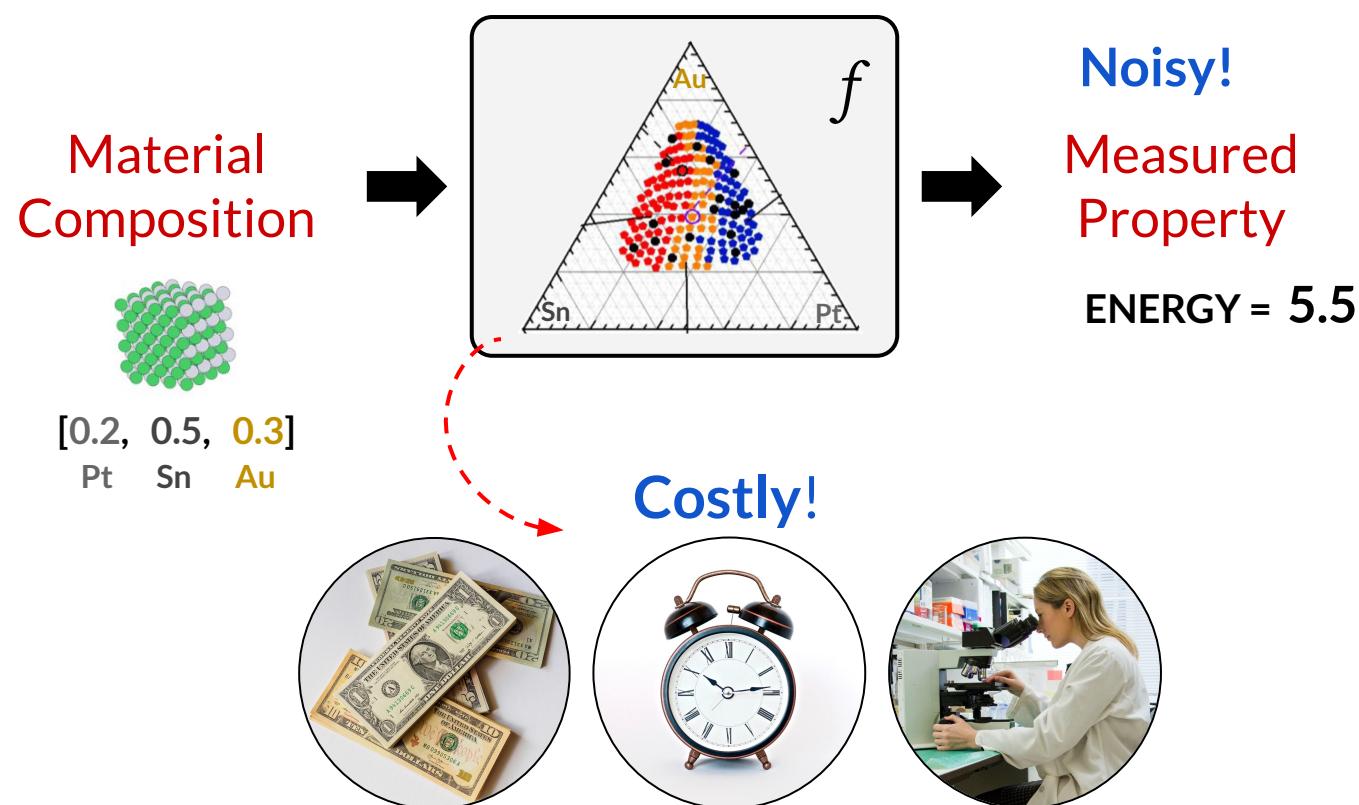
MATERIALS  
SCIENCE



# Expensive Black-box Functions

Many real-world systems can be modeled as an **expensive black-box function**  $f$ .

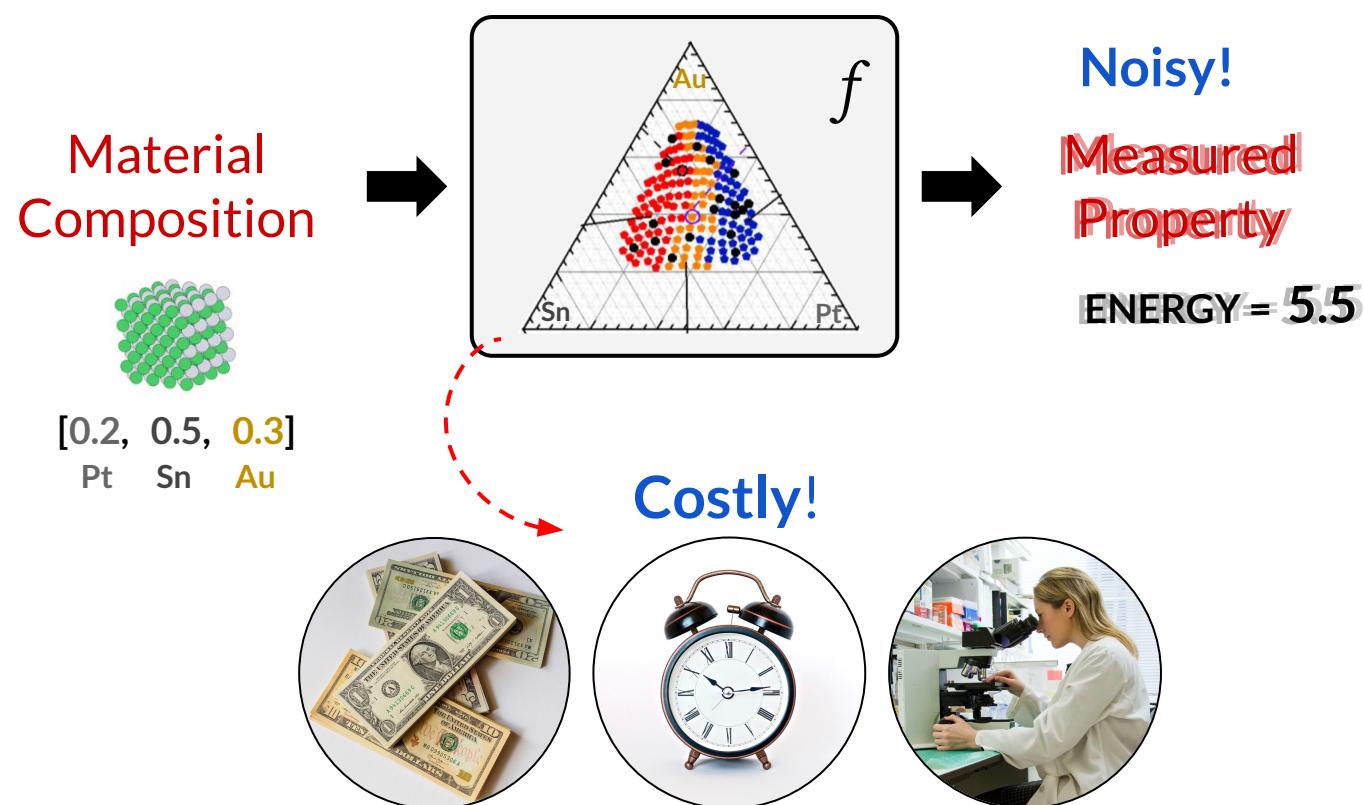
MATERIALS  
SCIENCE



# Expensive Black-box Functions

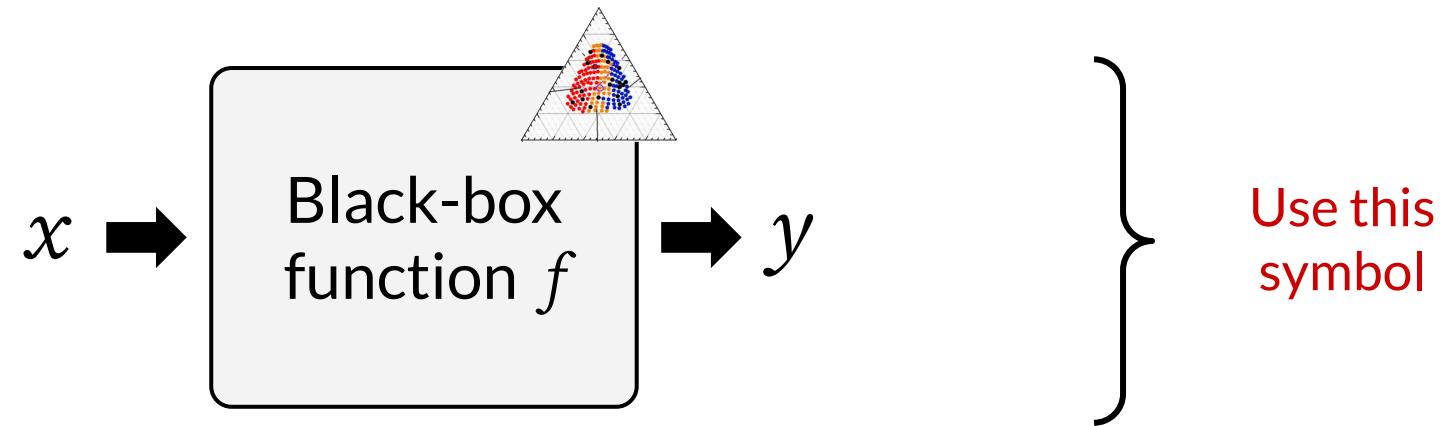
Many real-world systems can be modeled as an **expensive black-box function**  $f$ .

MATERIALS  
SCIENCE



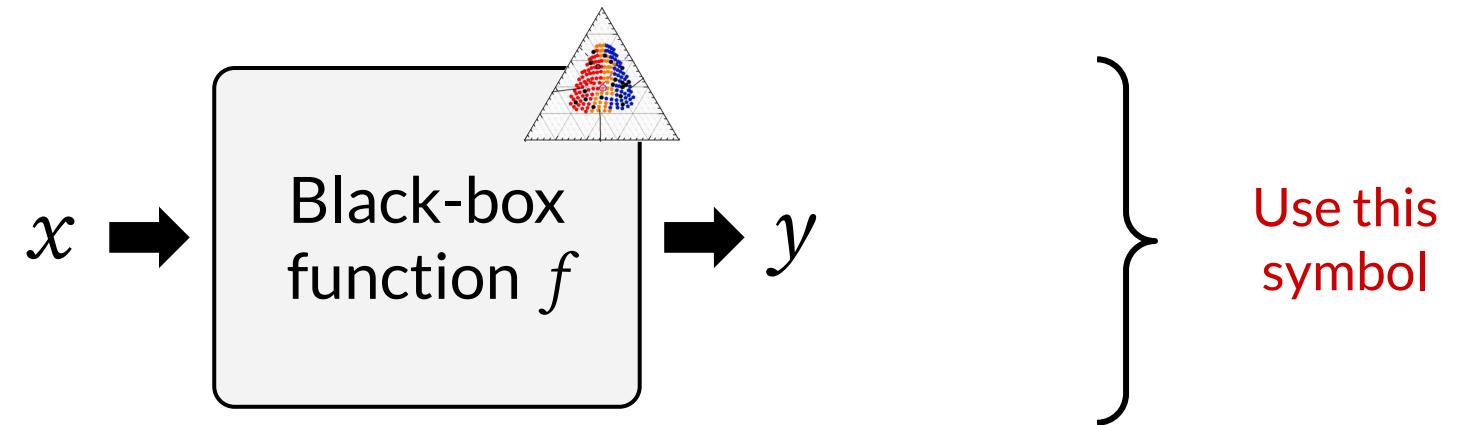
# Expensive Black-box Functions

Many real-world systems can be modeled as an **expensive black-box function**  $f$ .



# Expensive Black-box Functions

Many real-world systems can be modeled as an **expensive black-box function**  $f$ .

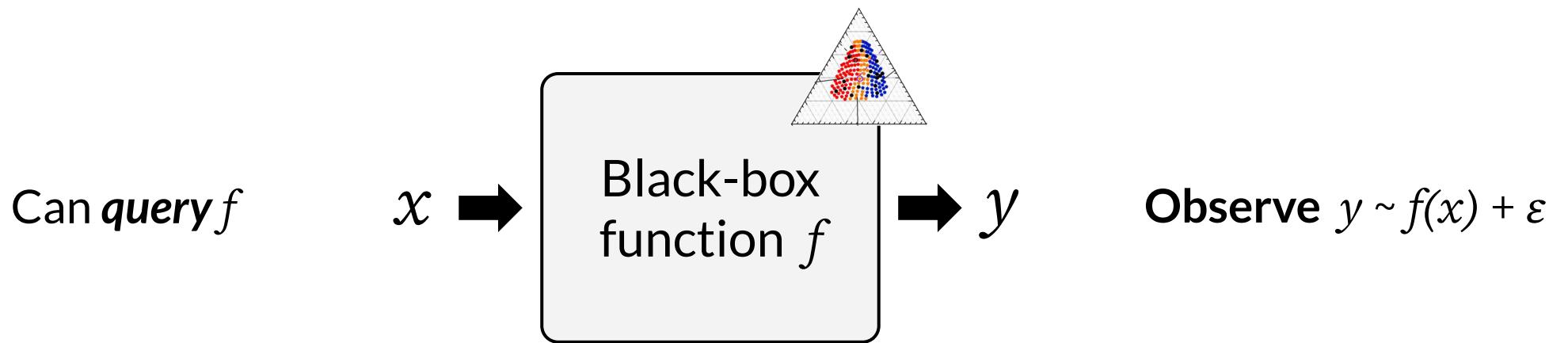


Leads to task of **black-box global optimization**.

# **Black-Box Global Optimization**

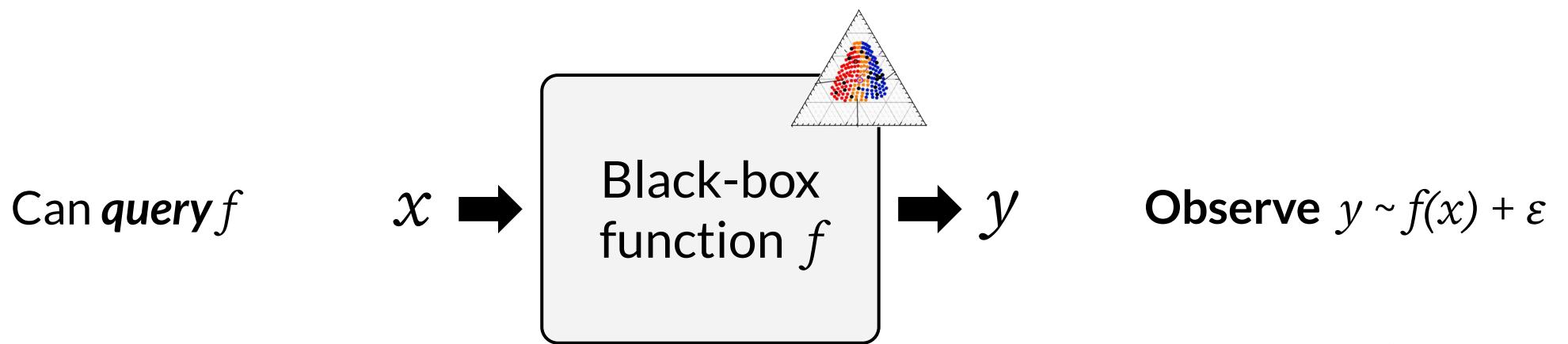
# Black-Box Global Optimization

Suppose we have an **expensive black-box function**  $f : \mathcal{X} \rightarrow \mathbb{R}$ .



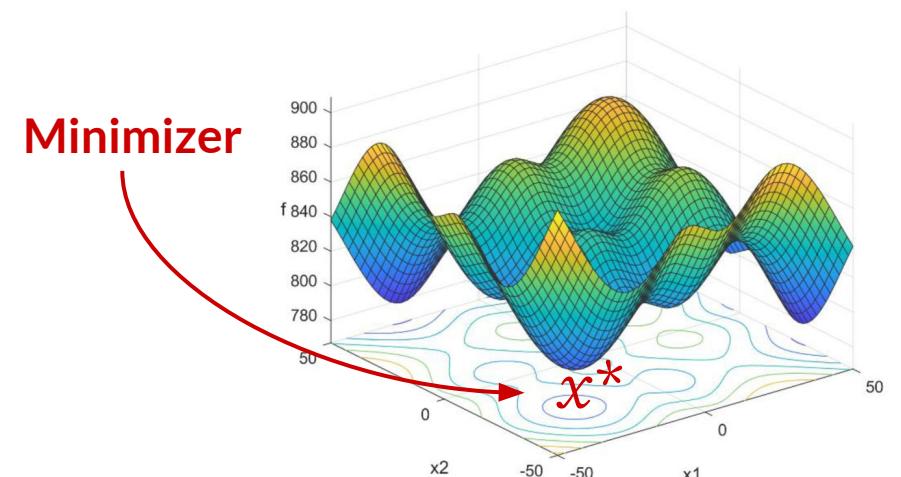
# Black-Box Global Optimization

Suppose we have an **expensive black-box function**  $f : \mathcal{X} \rightarrow \mathbb{R}$ .



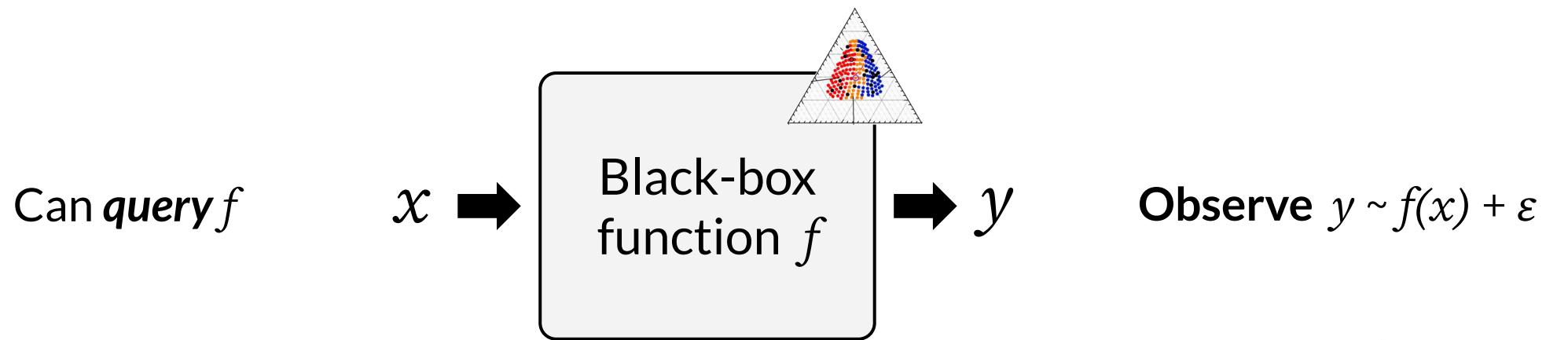
**Goal:** estimate the location of a global optima.

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



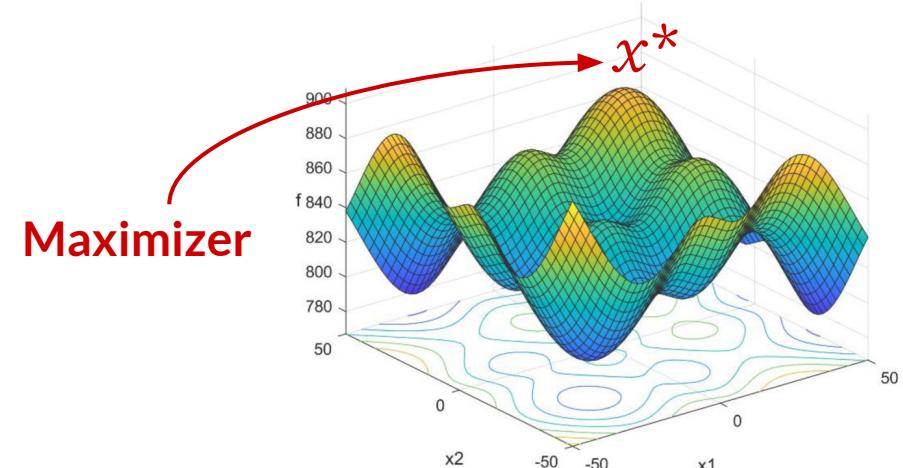
# Black-Box Global Optimization

Suppose we have an **expensive black-box function**  $f : \mathcal{X} \rightarrow \mathbb{R}$ .



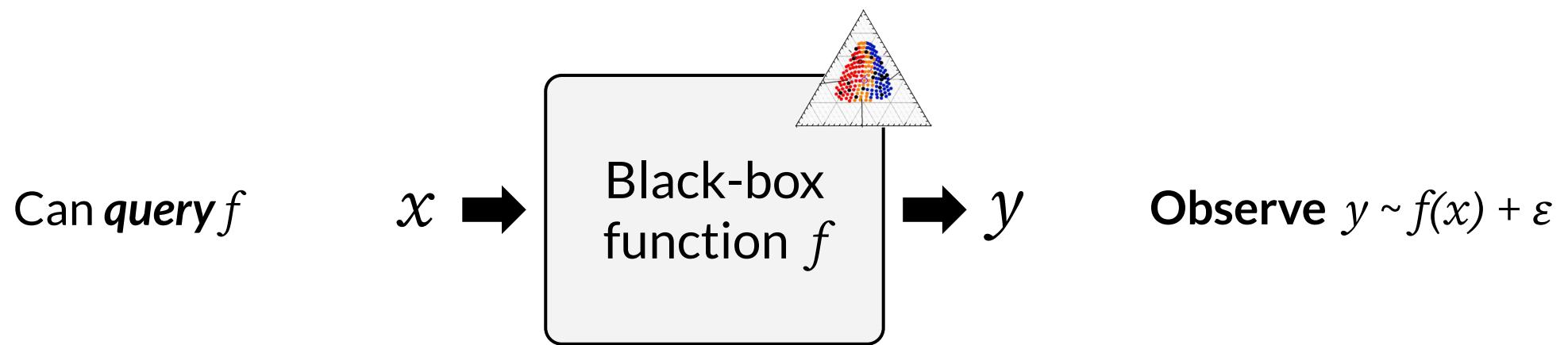
**Goal:** estimate the location of a global optima.

$$x^* = \arg \max_{x \in \mathcal{X}} f(x)$$



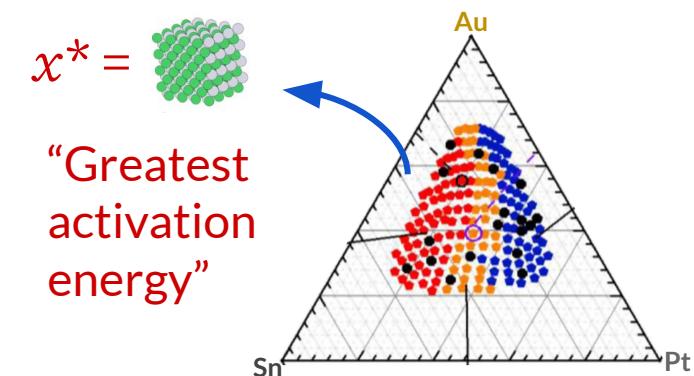
# Black-Box Global Optimization

Suppose we have an **expensive black-box function**  $f : \mathcal{X} \rightarrow \mathbb{R}$ .



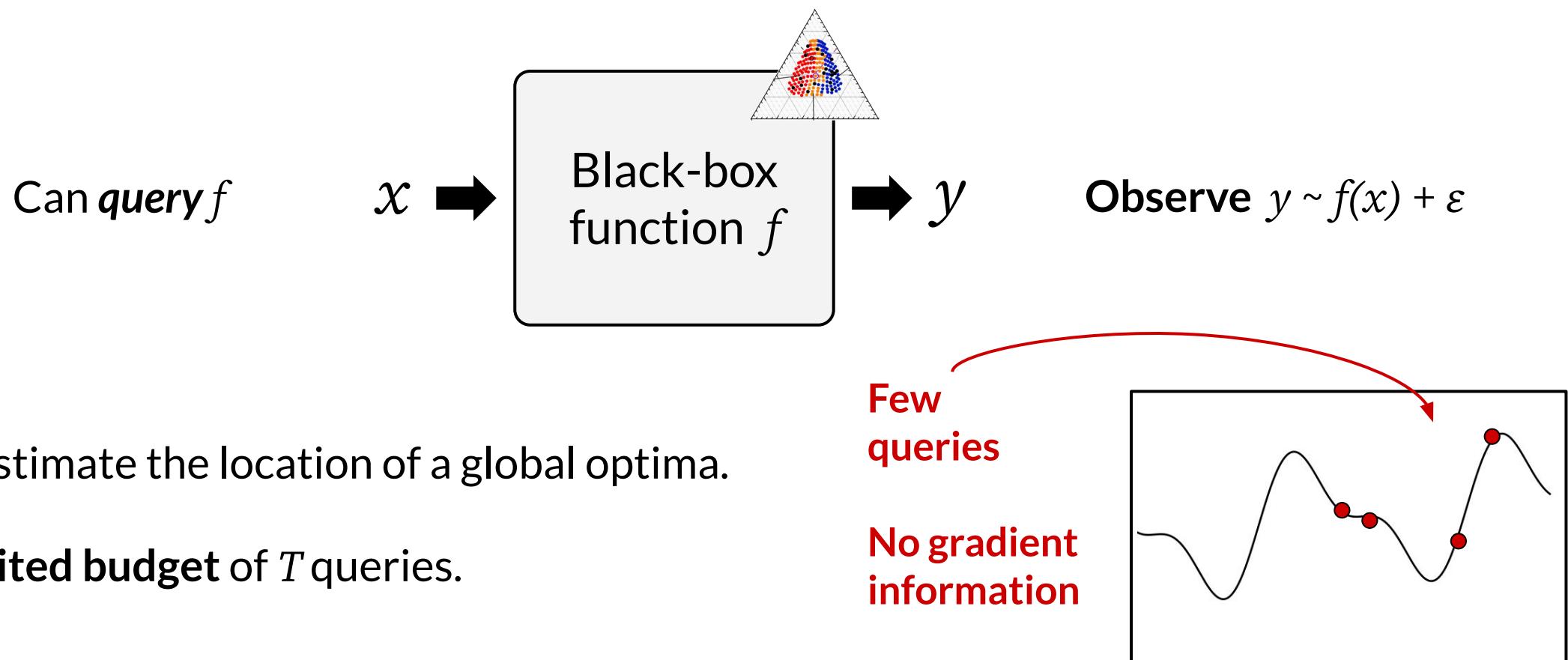
**Goal:** estimate the location of a global optima.

$$x^* = \arg \max_{x \in \mathcal{X}} f(x)$$

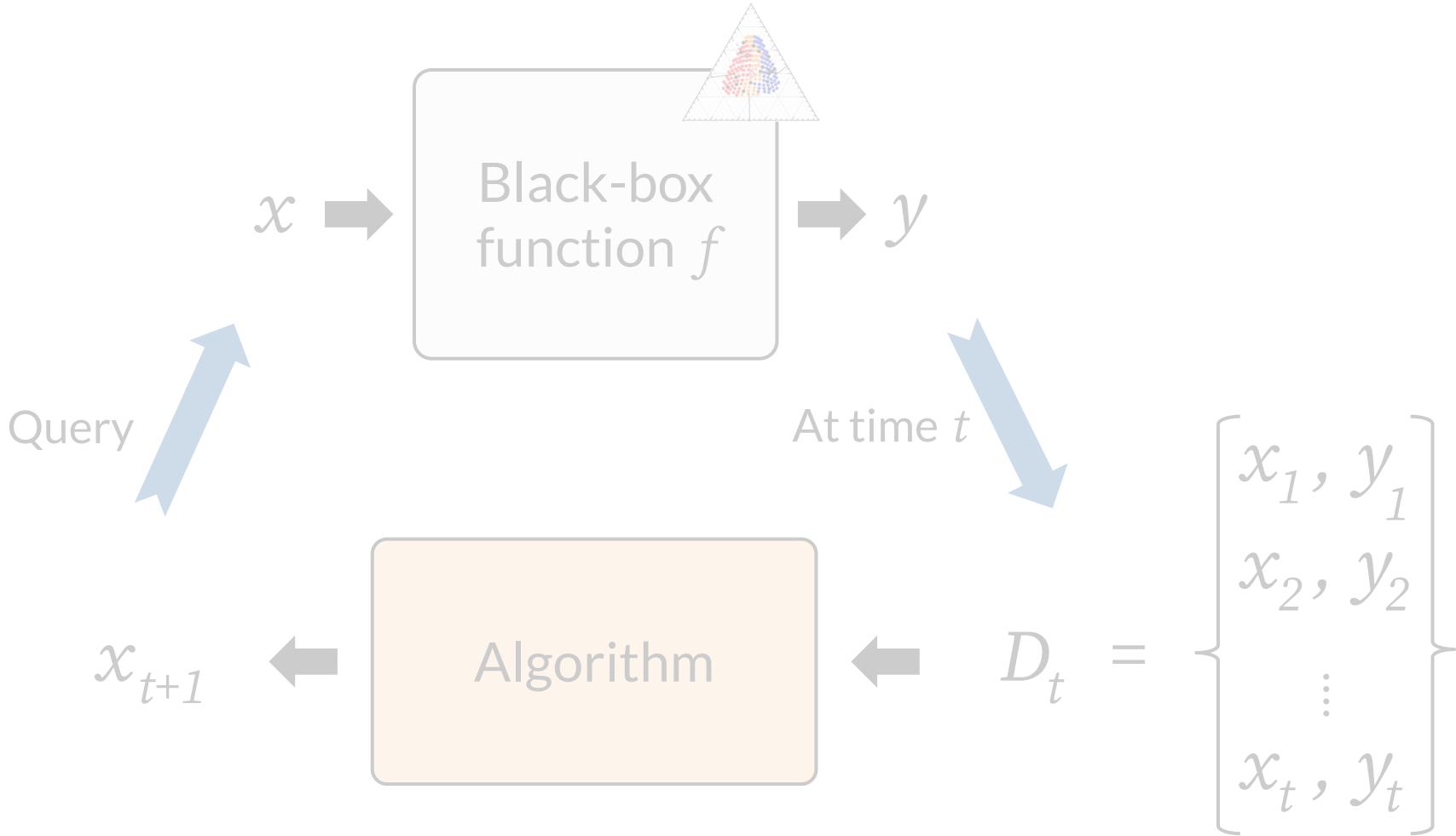


# Black-Box Global Optimization

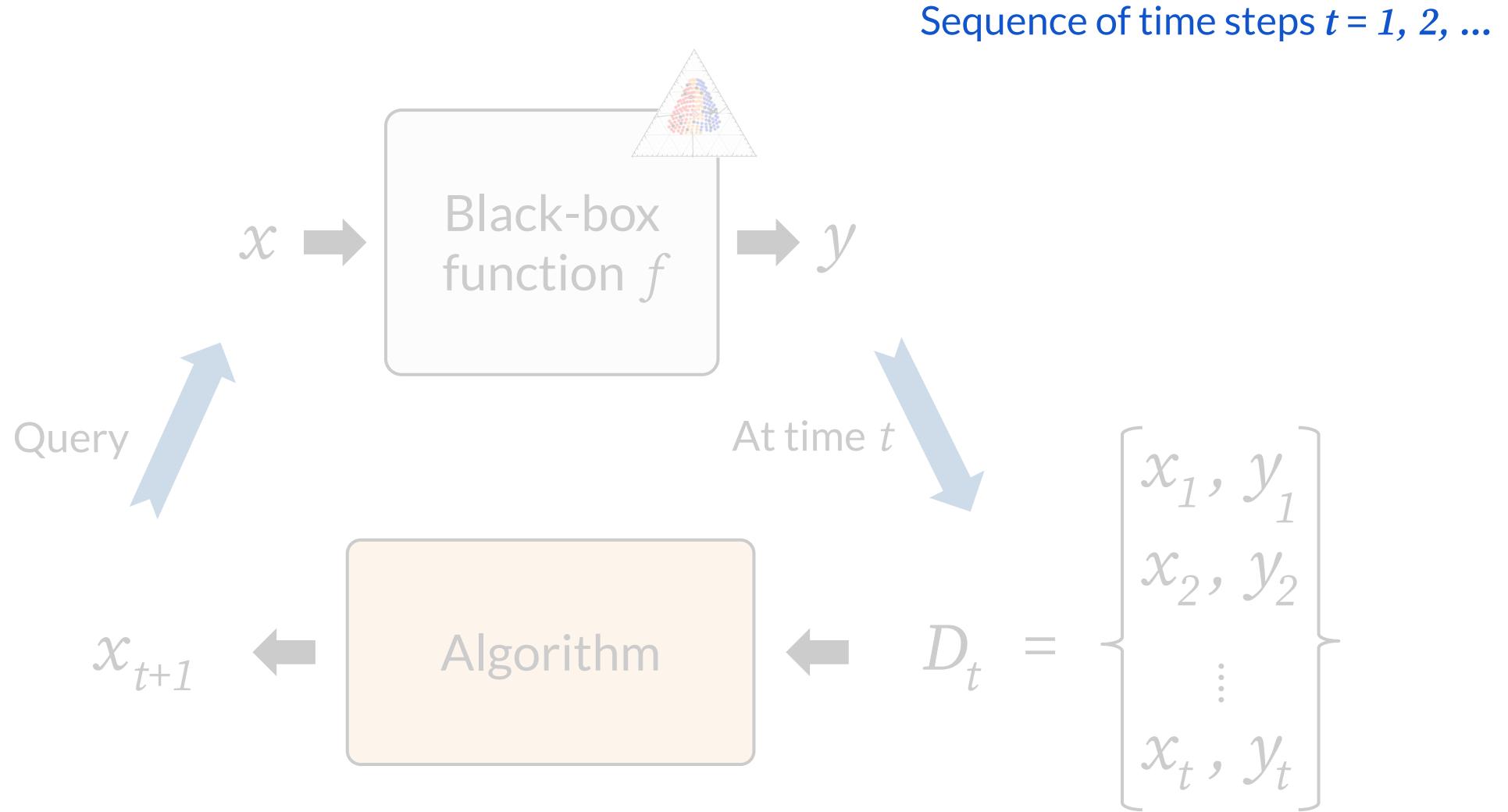
Suppose we have an **expensive black-box function**  $f : \mathcal{X} \rightarrow \mathbb{R}$ .



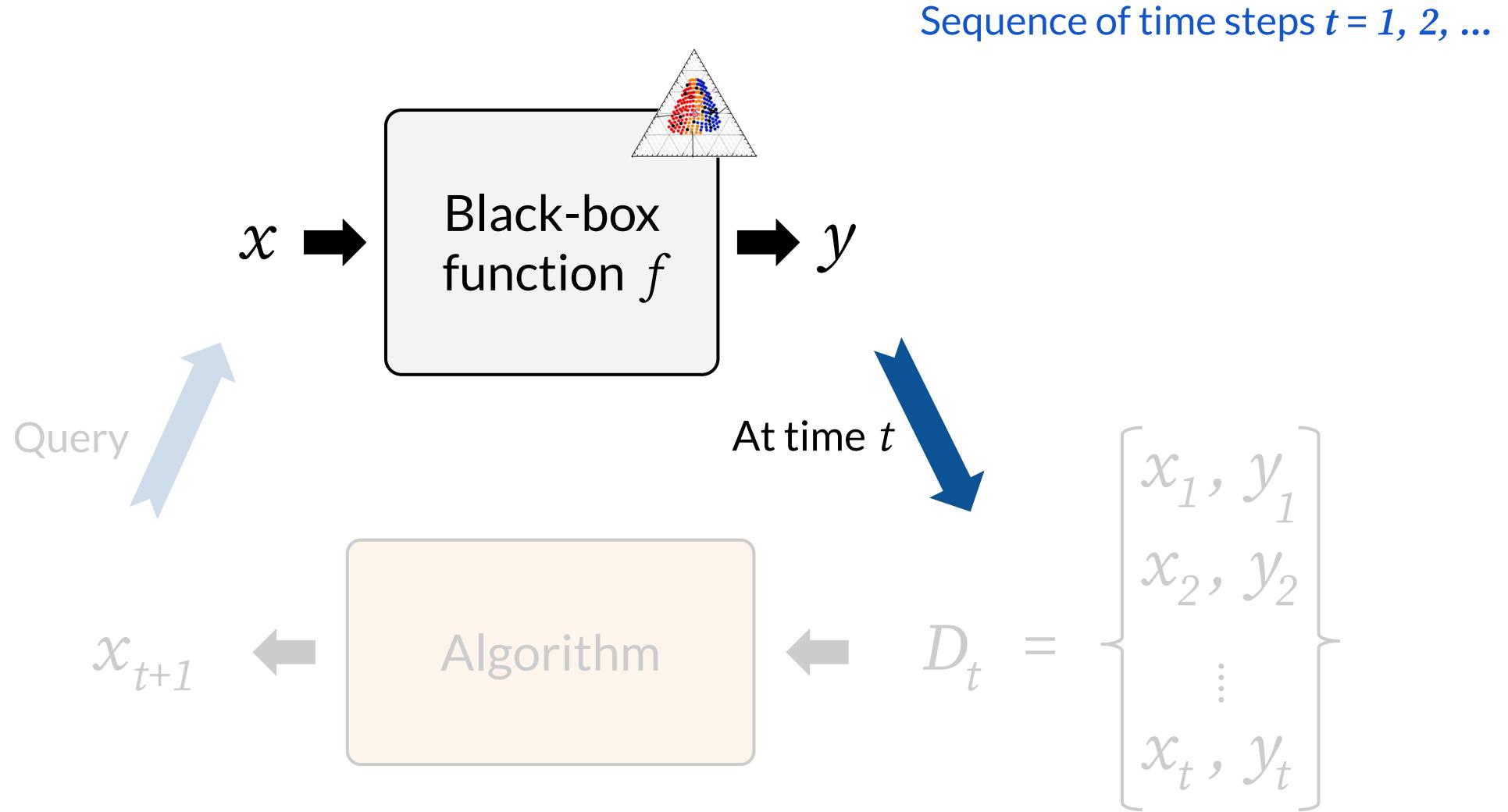
# Black-Box Global Optimization – Usual Setup



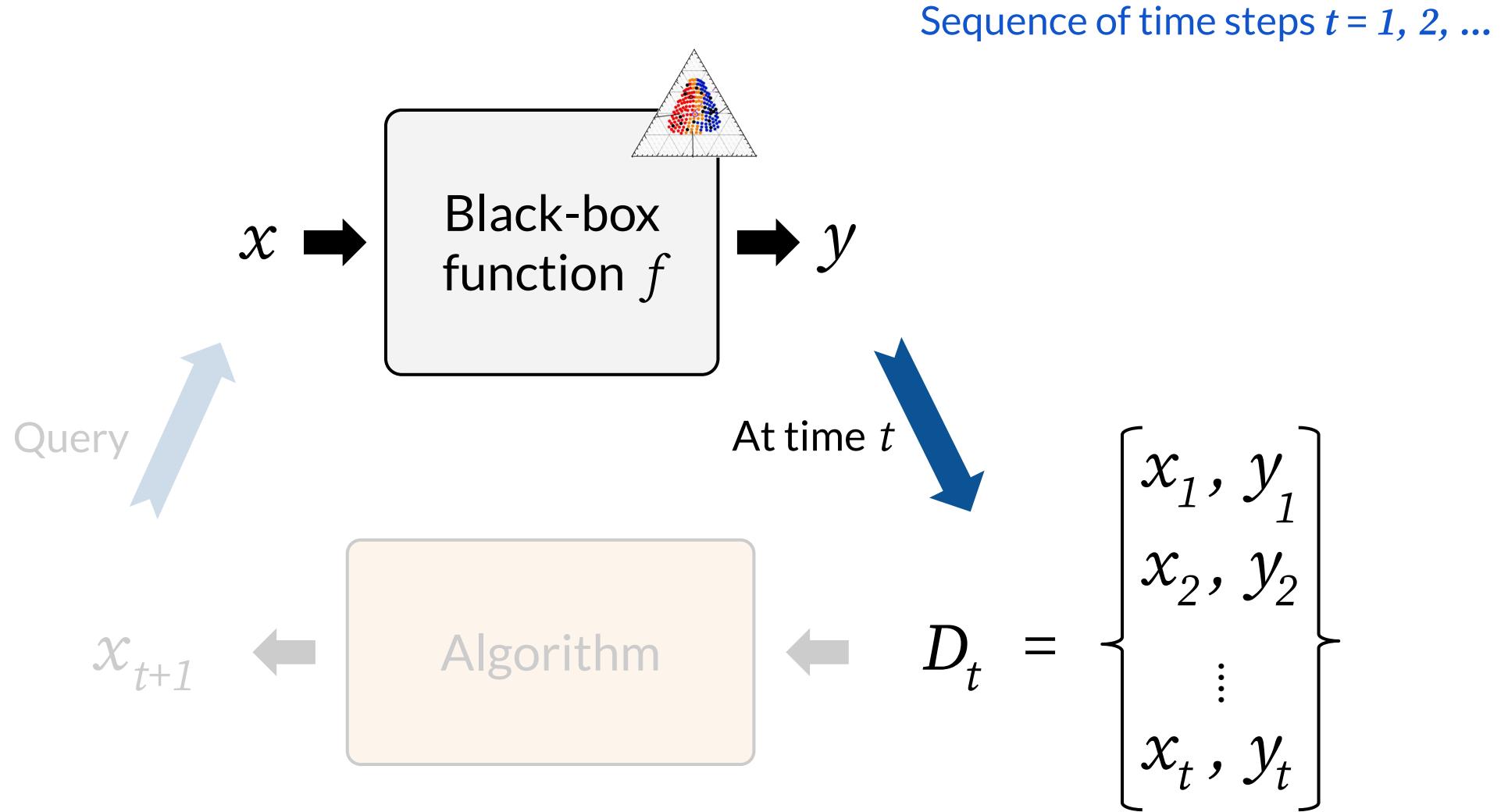
# Black-Box Global Optimization – Usual Setup



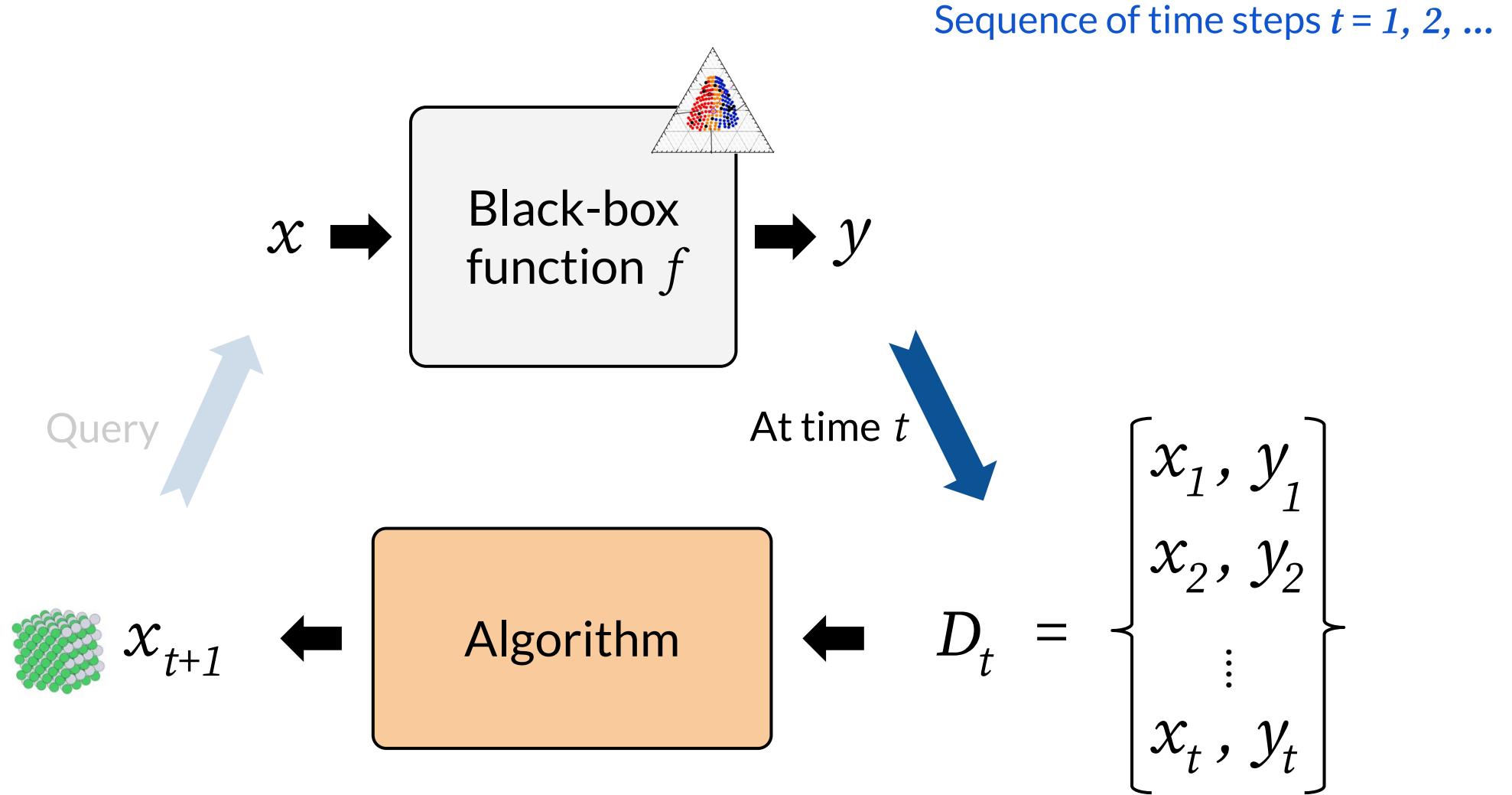
# Black-Box Global Optimization – Usual Setup



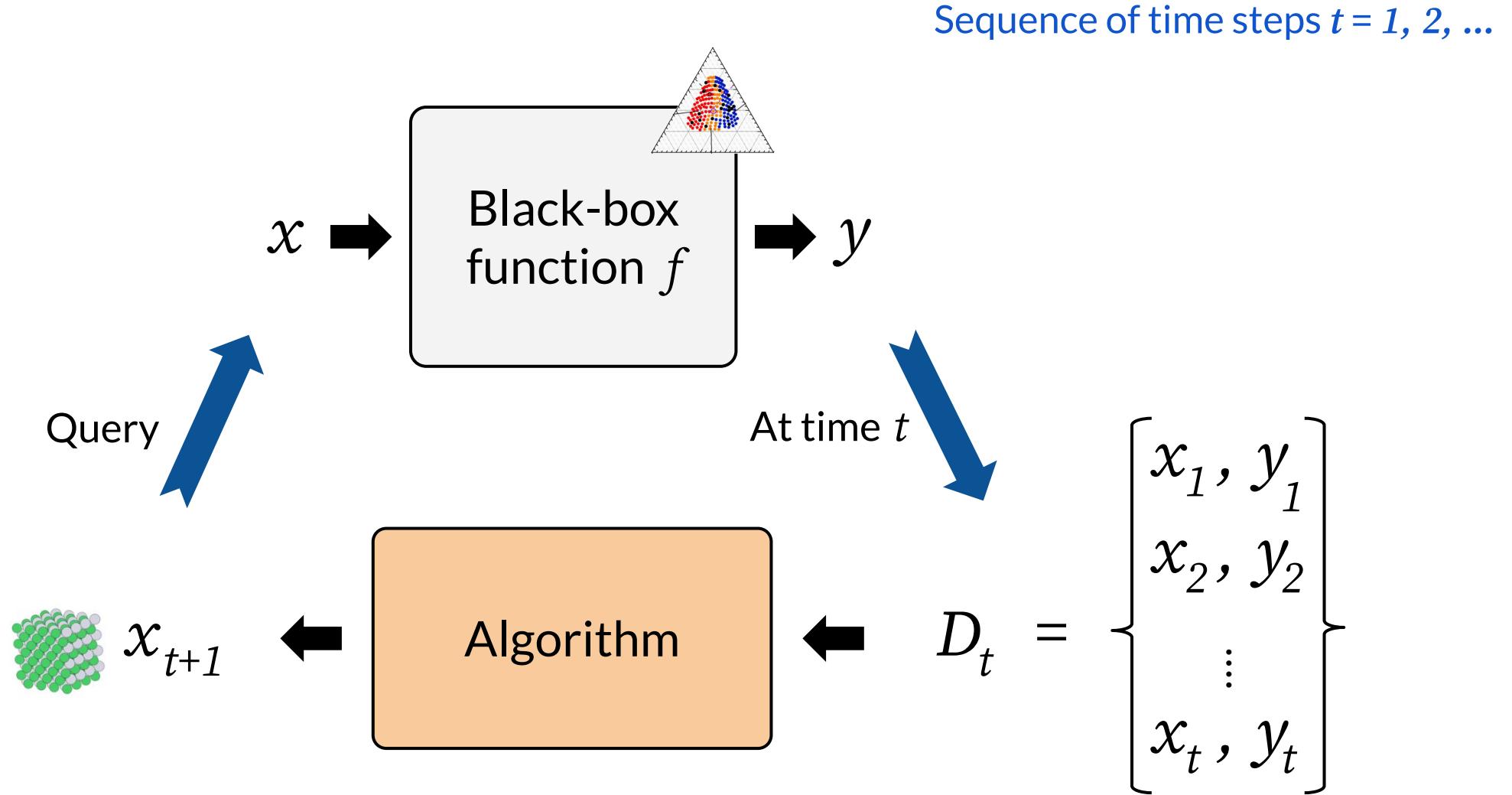
# Black-Box Global Optimization – Usual Setup



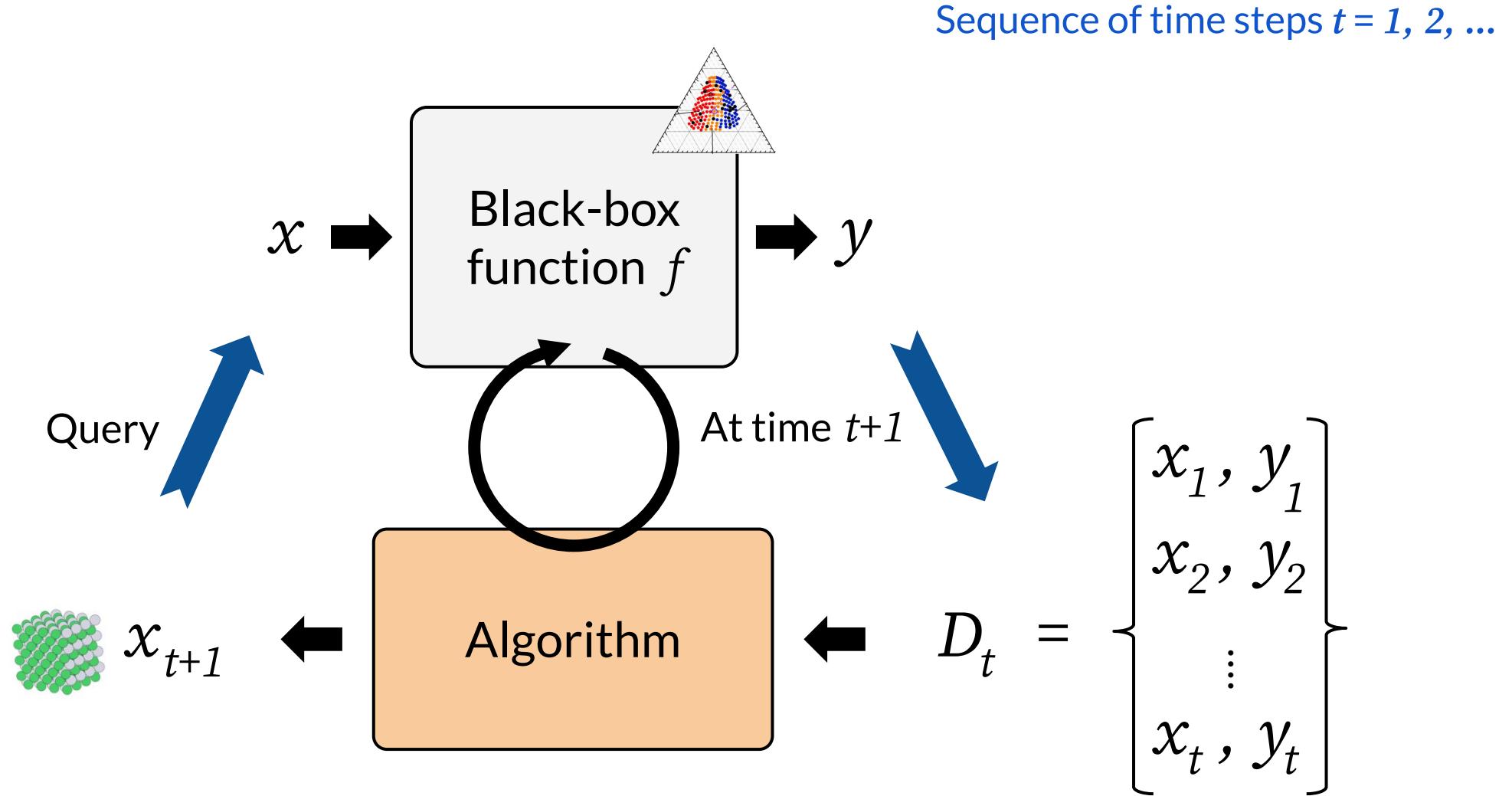
# Black-Box Global Optimization – Usual Setup



# Black-Box Global Optimization – Usual Setup



# Black-Box Global Optimization – Usual Setup



# Bayesian Optimization – Motivation

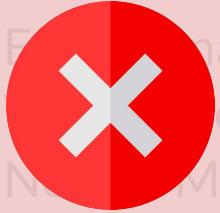
Derivative-free optimization:

- Evolutionary or genetic algorithms [FOW, 1966]
- Trust region methods [LM, 1944-1963]
- Nelder-Mead / simplex methods [NM, 1965]
- Simulated annealing [P, 1970; KGV, 1983]

# Bayesian Optimization – Motivation

Derivative-free optimization:

- Evolutionary computation methods [FOW, 1966]
- Nonlinear simplex methods [LM, 1944-1963]
- Simulated annealing [P, 1970; KGV, 1983]



**Can require  
1000s of queries!**

Need:

- Better sample-efficiency.
- Adaptive to history (data-driven).
- (Ideally) ability to incorporate any prior knowledge/beliefs about function.

# Bayesian Optimization – Motivation

Derivative-free optimization:

- Evolutionary computation methods [FOW, 1966]
- Randomization methods [LM, 1944-1963]
- Nelder-Mead/Simplex method [NMS, 1965]
- Simulated annealing [P, 1970; KGV, 1983]



**Can require  
1000s of queries!**

Need:

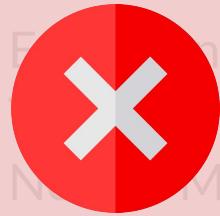
- Better sample-efficiency.
- Adaptive to history (data-driven).
- (Ideally) ability to incorporate any prior knowledge/beliefs about function.

What is Bayesian Opt? Uses a (probabilistic) predictive model to aid in optimization.

# Bayesian Optimization – Motivation

Derivative-free optimization:

- Evolutionary computation methods [FOW, 1966]
- Nelder-Mead/Simplex method [NED, 1965]
- Simulated annealing [P, 1970; KGV, 1983]

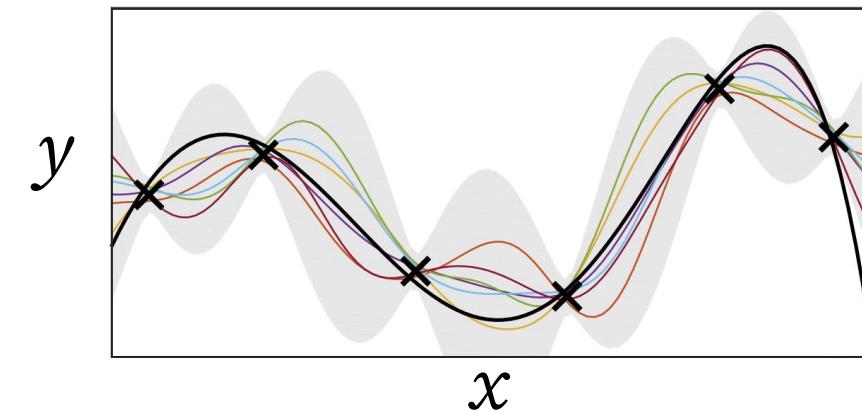
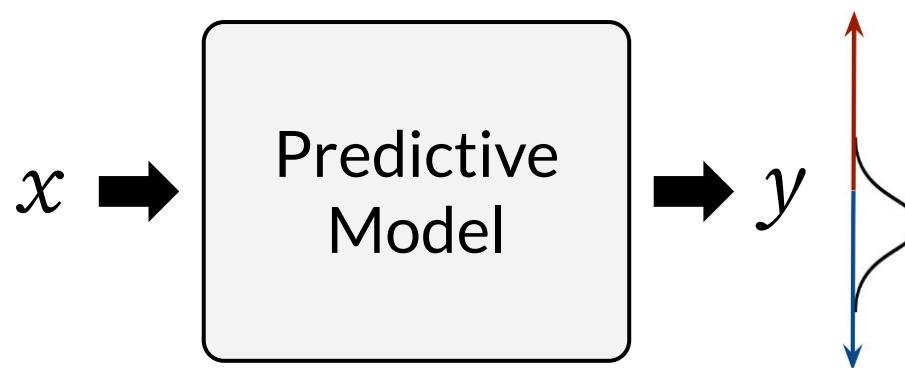


**Can require  
1000s of queries!**

Need:

- Better sample-efficiency.
- Adaptive to history (data-driven).
- (Ideally) ability to incorporate any prior knowledge/beliefs about function.

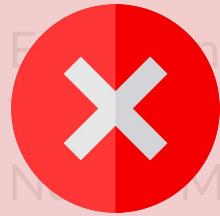
What is Bayesian Opt? Uses a (probabilistic) predictive model to aid in optimization.



# Bayesian Optimization – Motivation

Derivative-free optimization:

- Evolutionary computation methods [FOW, 1966]
- Nelder-Mead/Simplex method [NED, 1965]
- Simulated annealing [P, 1970; KGV, 1983]

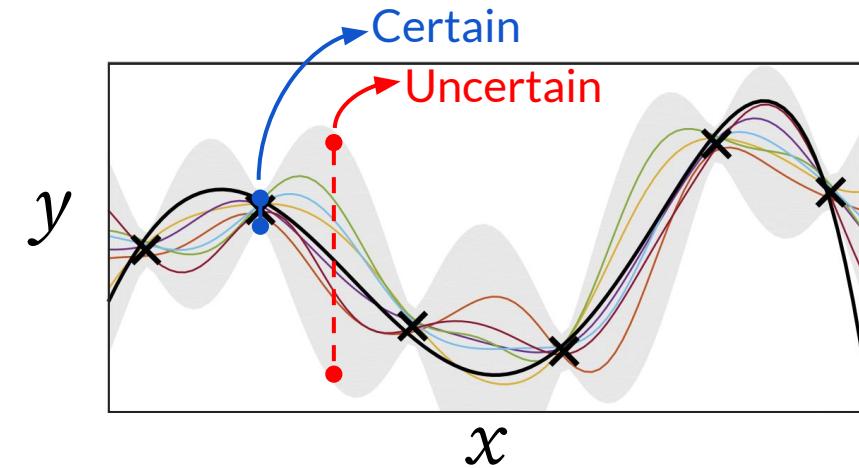
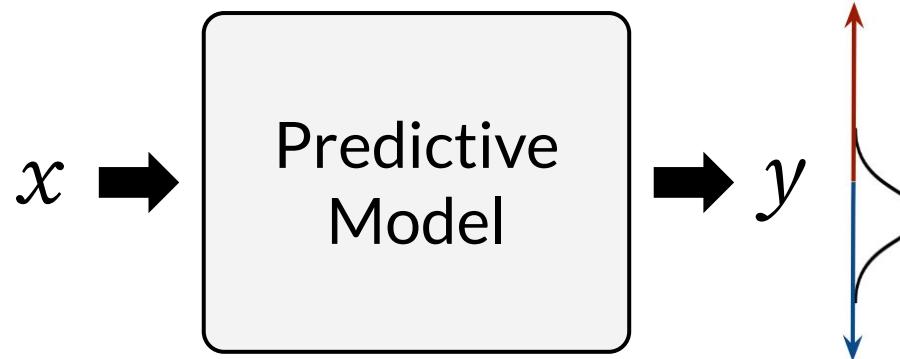


**Can require  
1000s of queries!**

Need:

- Better sample-efficiency.
- Adaptive to history (data-driven).
- (Ideally) ability to incorporate any prior knowledge/beliefs about function.

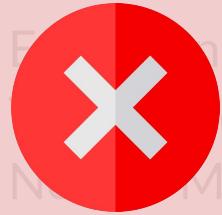
What is Bayesian Opt? Uses a (probabilistic) predictive model to aid in optimization.



# Bayesian Optimization – Motivation

Derivative-free optimization:

- Evolutionary computation methods [FOW, 1966]
- Nelder-Mead/Simplex method [NED, 1965]
- Simulated annealing [P, 1970; KGV, 1983]



**Can require  
1000s of queries!**

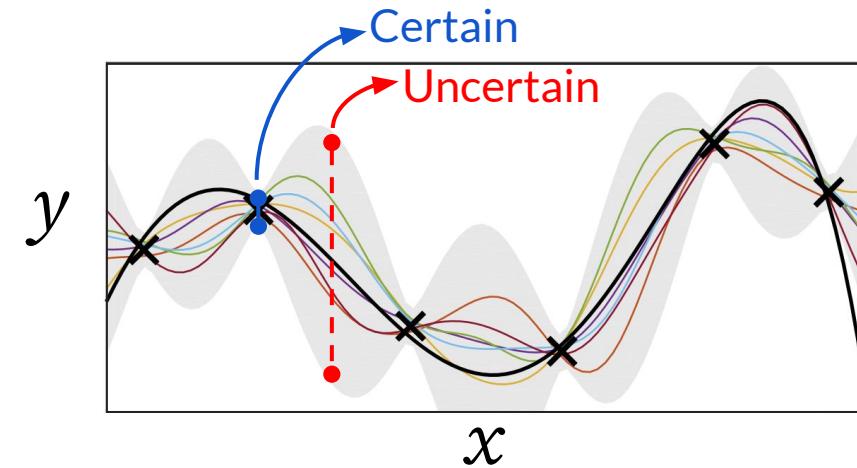
Need:

- Better sample-efficiency.
- Adaptive to history (data-driven).
- (Ideally) ability to incorporate any prior knowledge/beliefs about function.

What is Bayesian Opt? Uses a (probabilistic) predictive model to aid in optimization.



Model used for *intelligent* decision making  $\Rightarrow$  sample-efficient opt



# Bayesian Optimization – Algorithm Overview

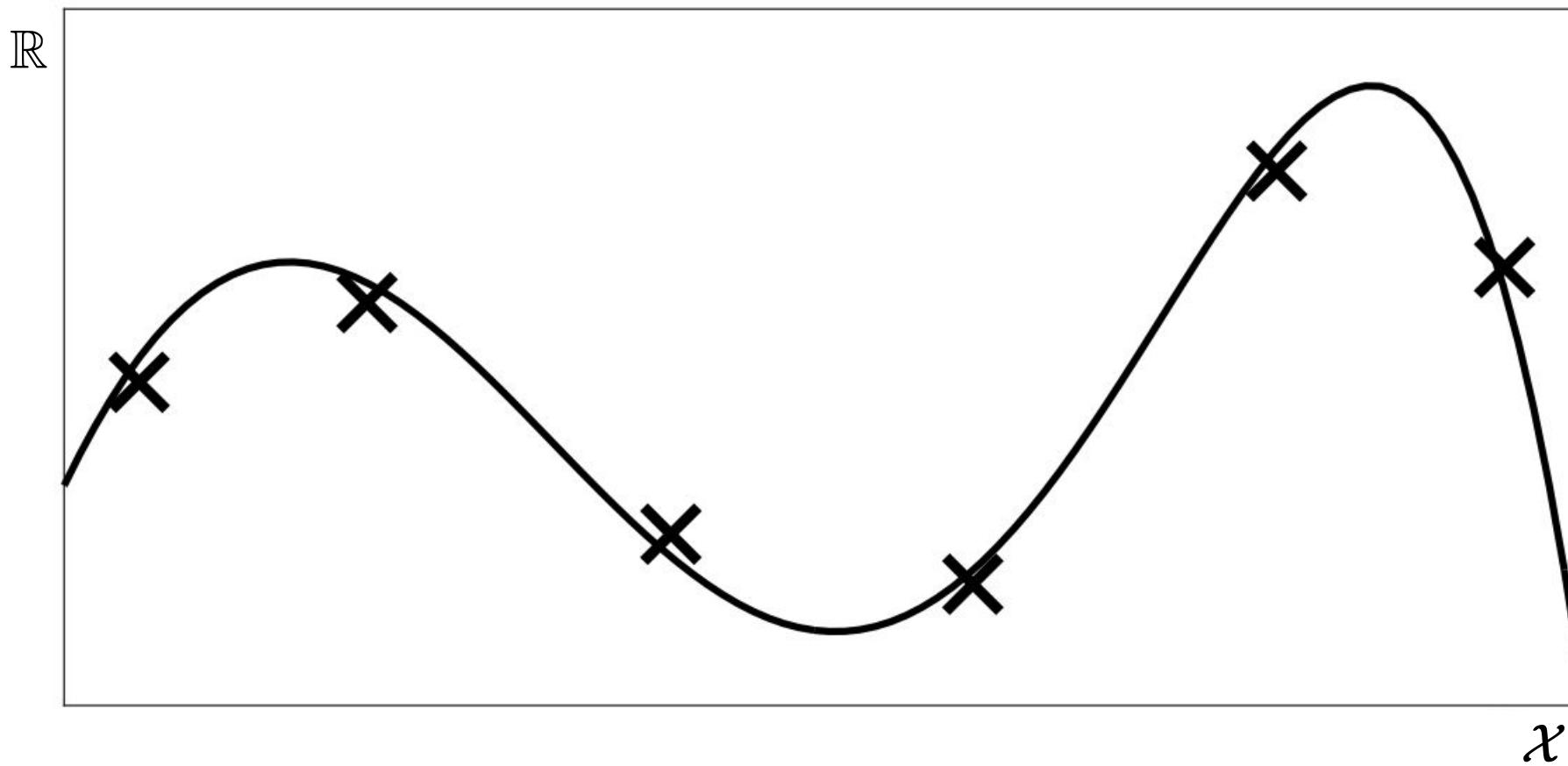
At each iteration of BO, three steps:

# Bayesian Optimization – Algorithm Overview

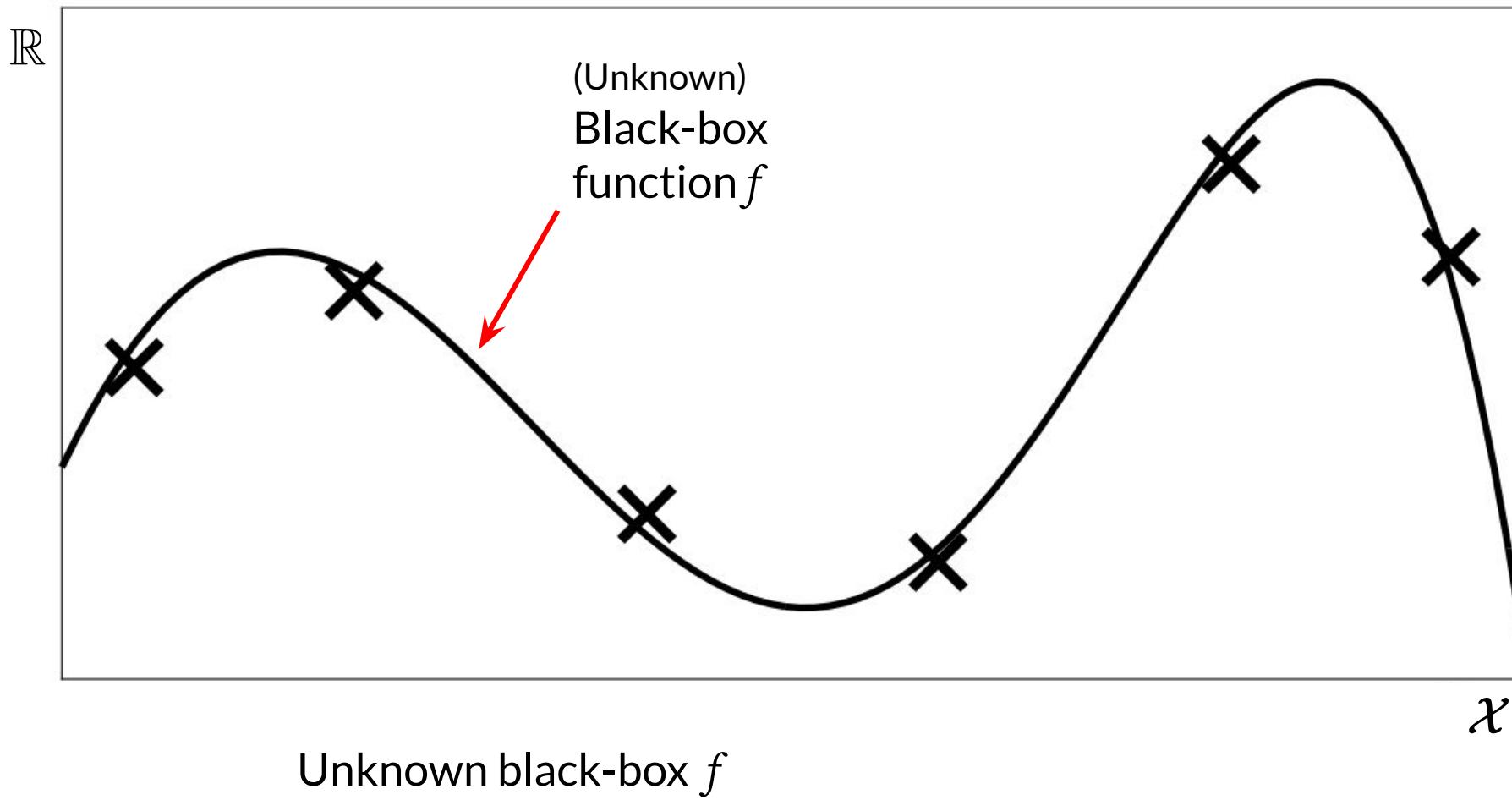
At each iteration of BO, three steps:

1. Fit model to dataset.
2. Optimize an *acquisition function* to select next point to query.
3. Query black-box function.

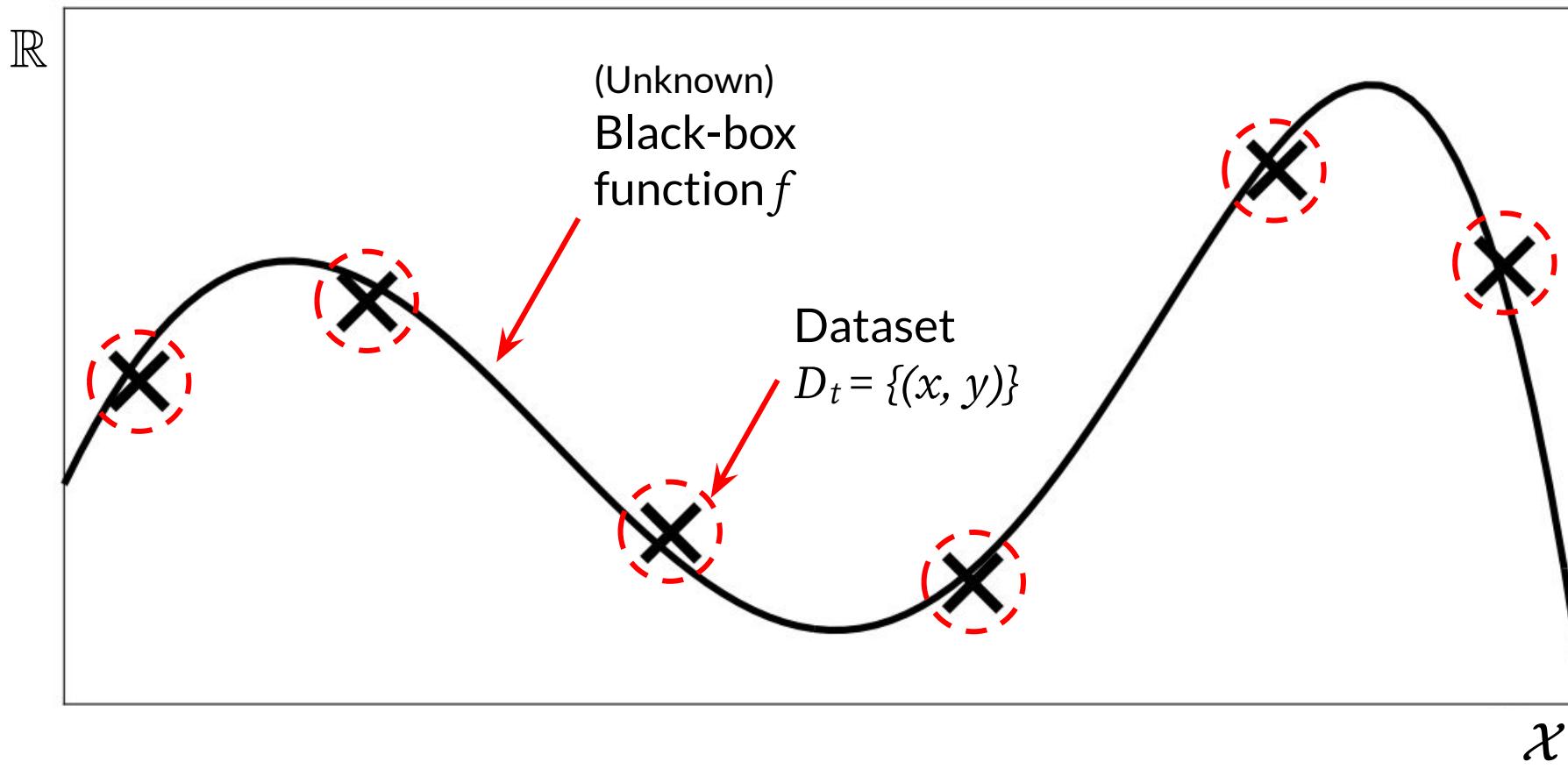
# Bayesian Optimization – Visualization



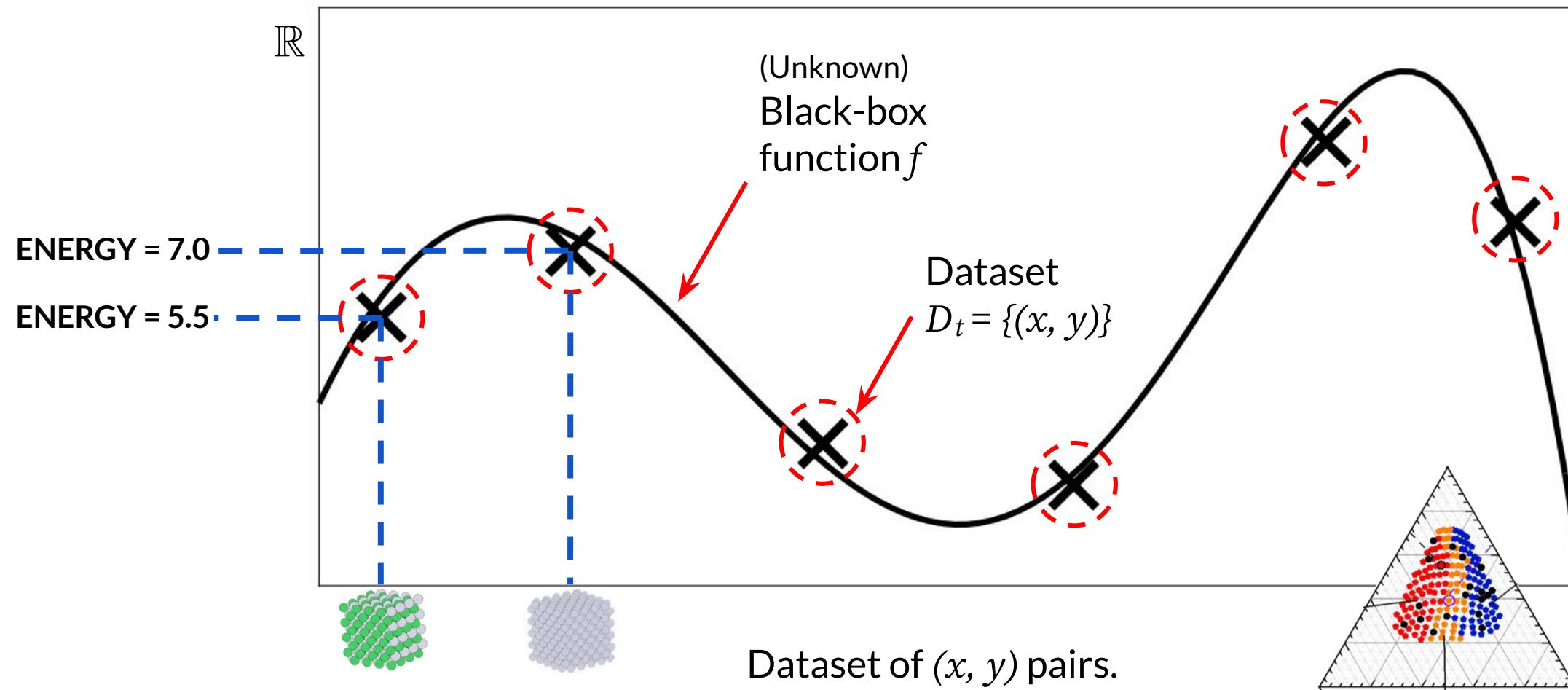
# Bayesian Optimization – Visualization



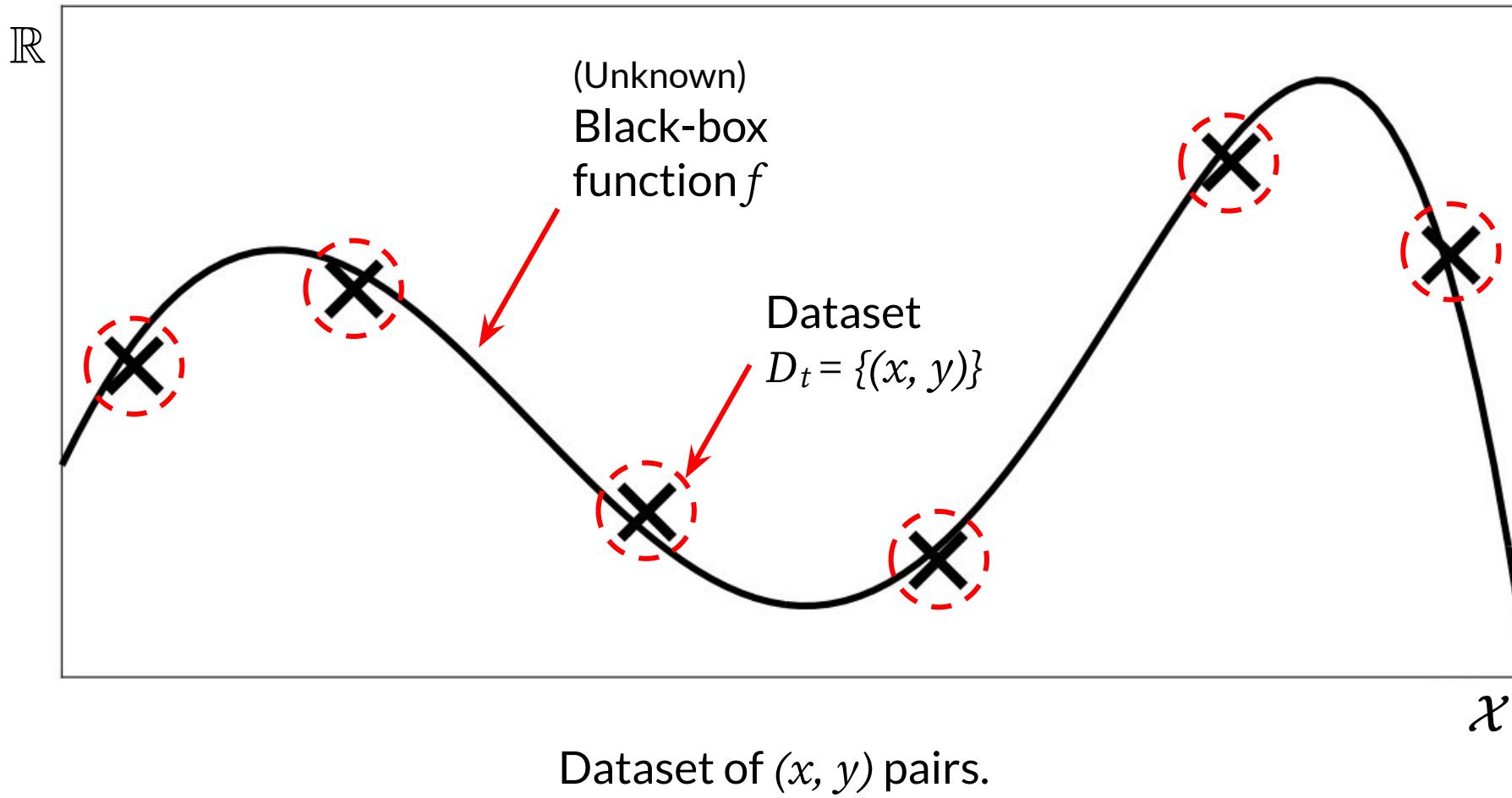
# Bayesian Optimization – Visualization



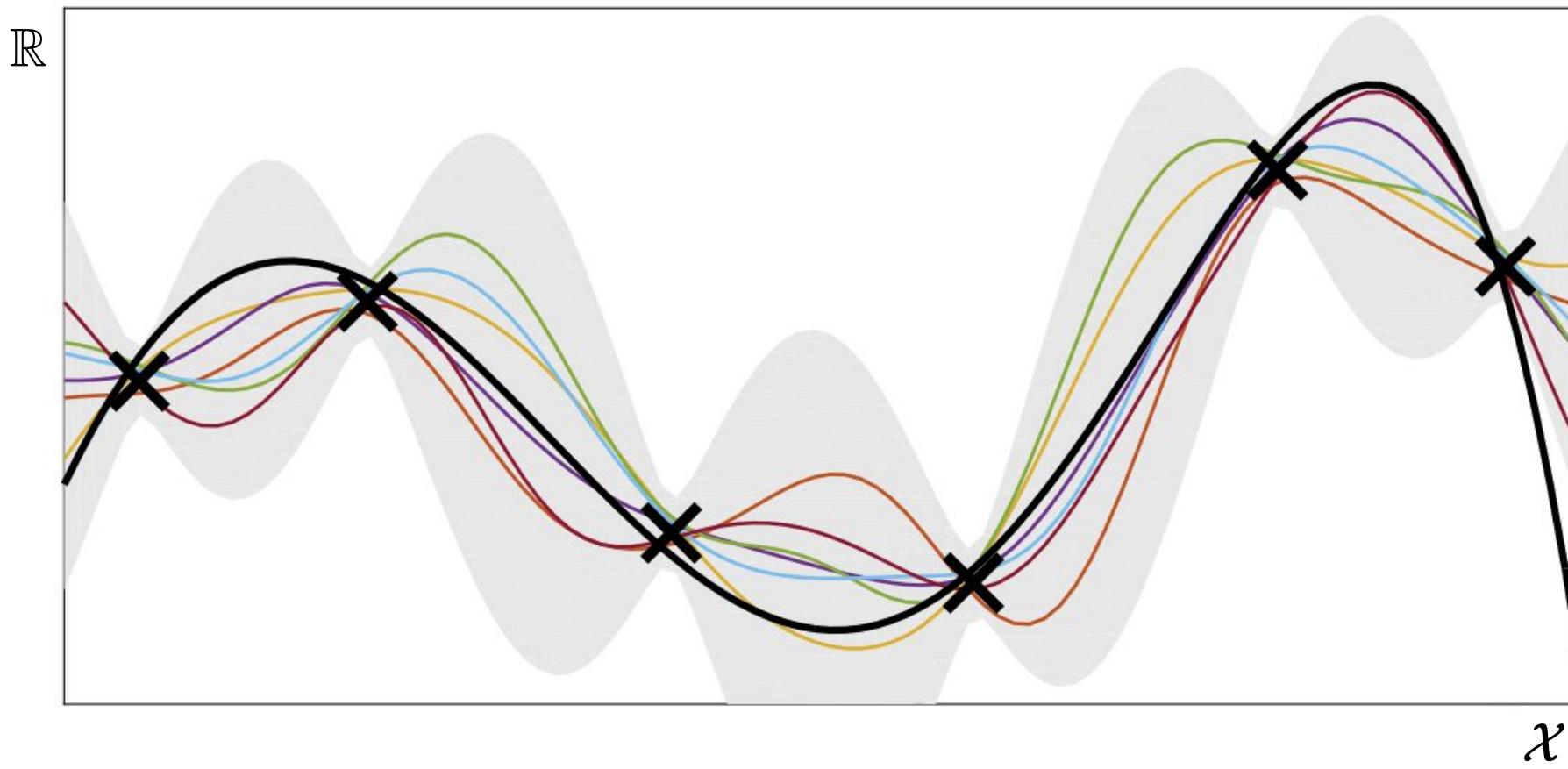
# Bayesian Optimization – Visualization



# Bayesian Optimization – Visualization

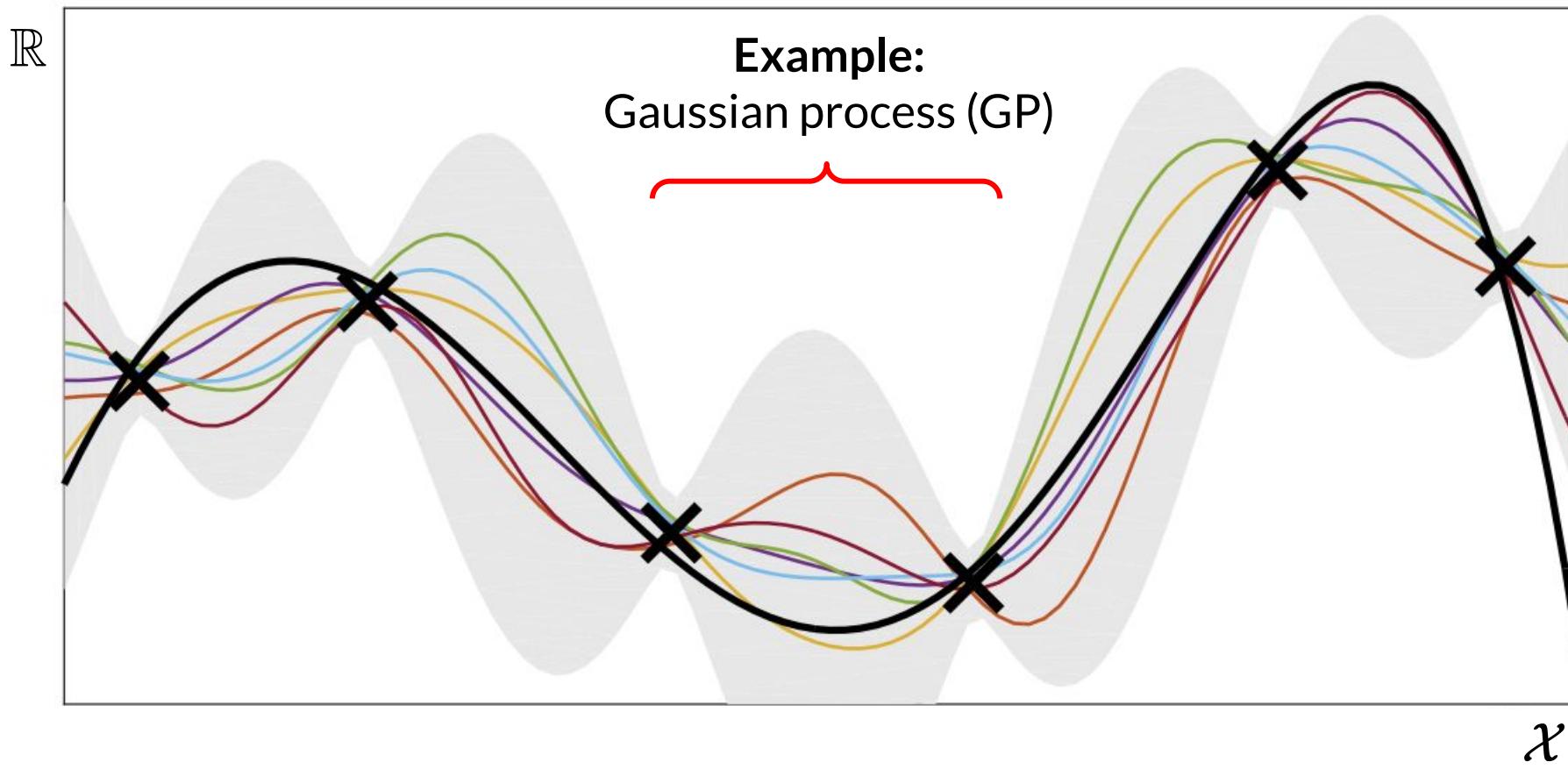


# Bayesian Optimization – Visualization



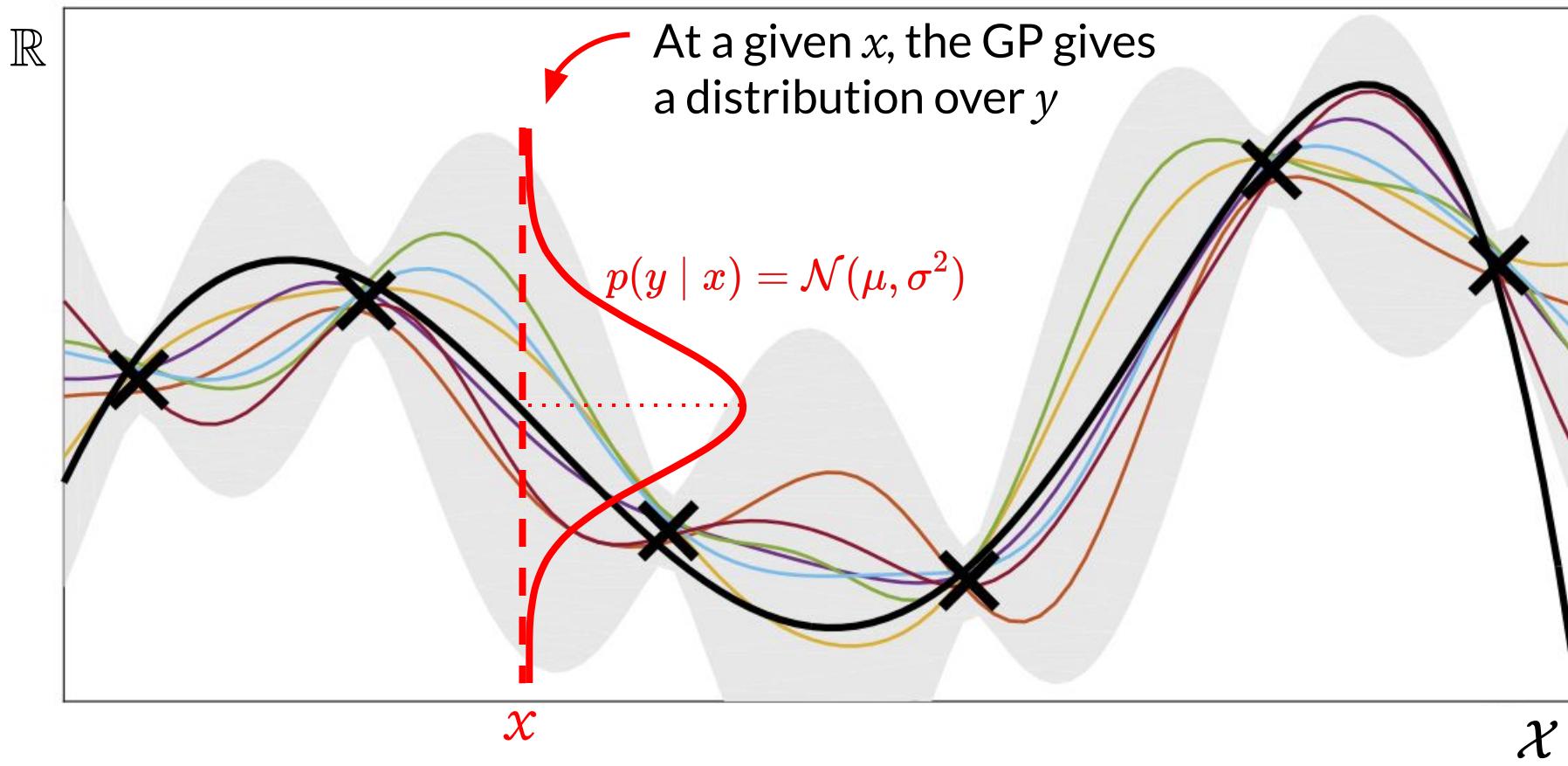
Given this dataset, can fit probabilistic model.

# Bayesian Optimization – Visualization

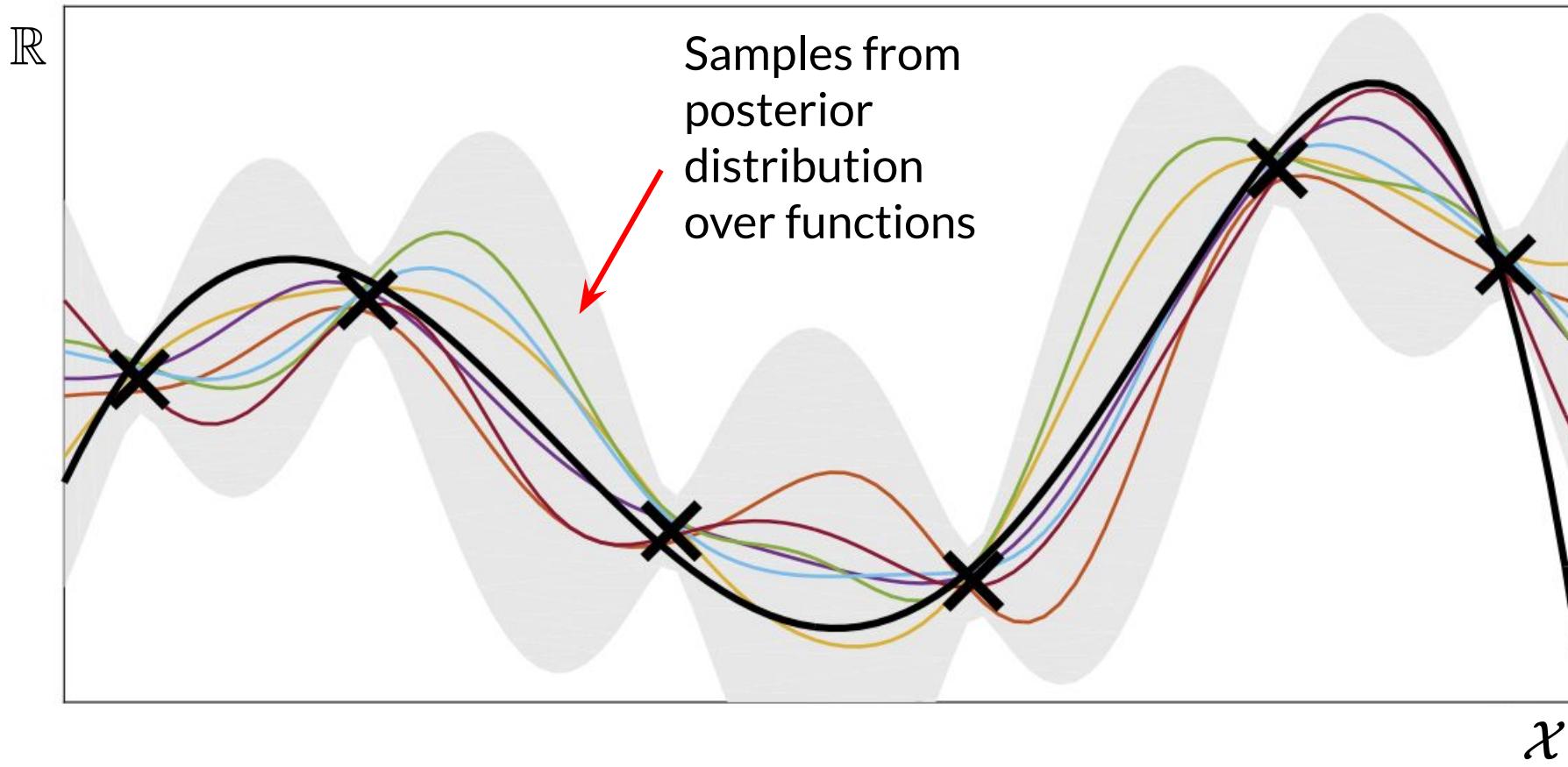


A popular choice for this model is a **Gaussian Process**.

# Bayesian Optimization – Visualization

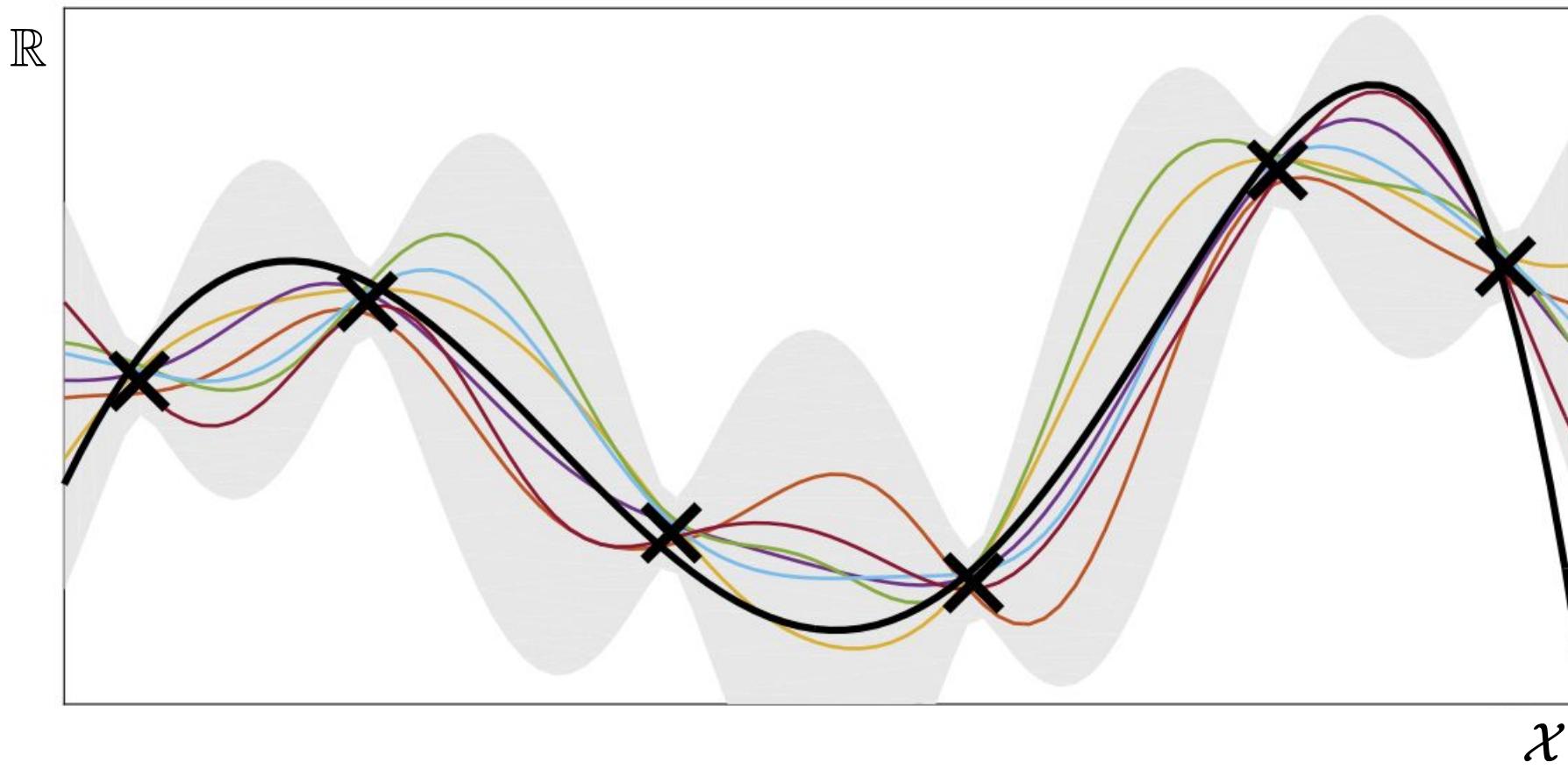


# Bayesian Optimization – Visualization



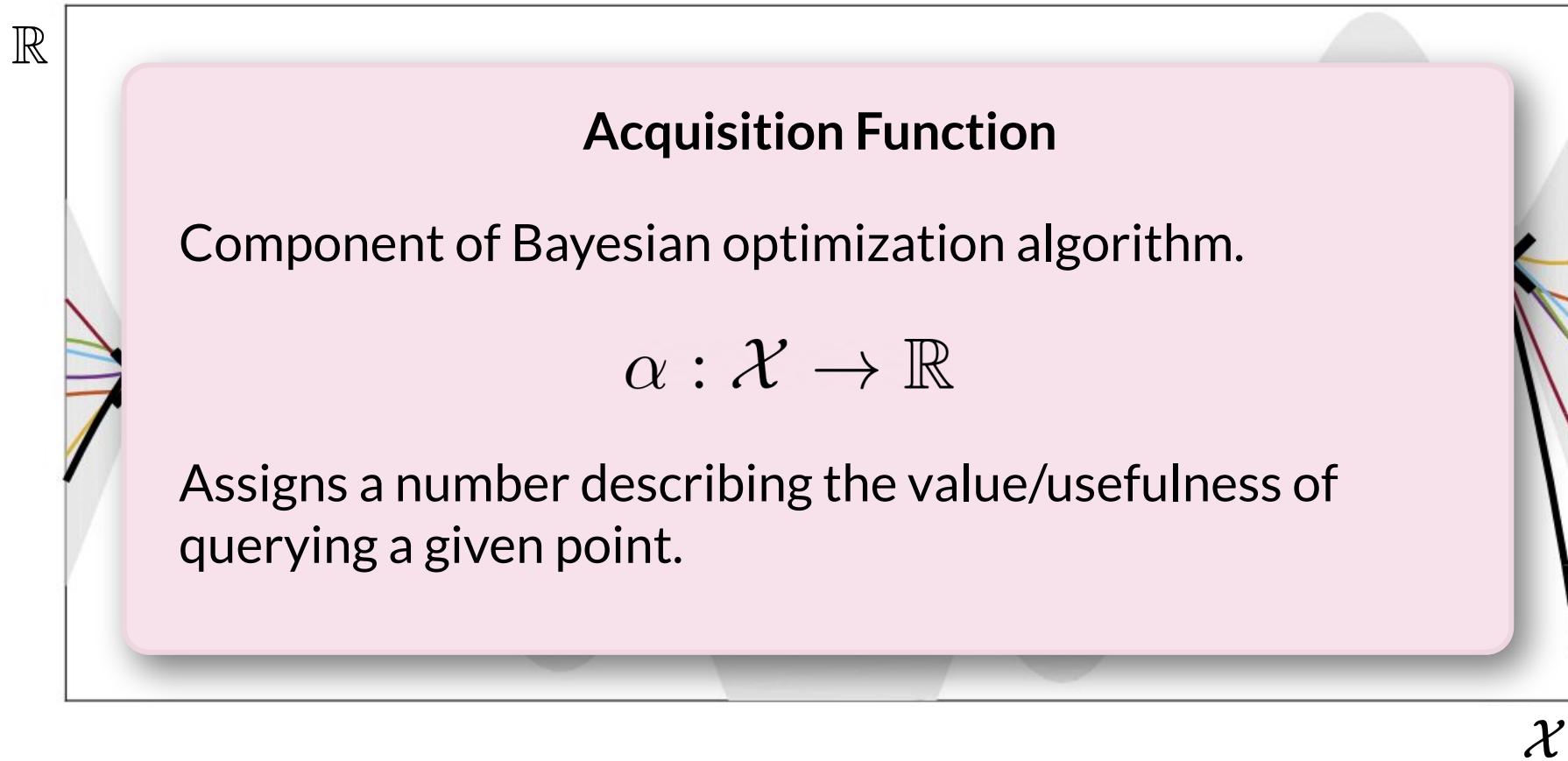
A popular choice for this model is a **Gaussian Process**.

# Bayesian Optimization – Visualization

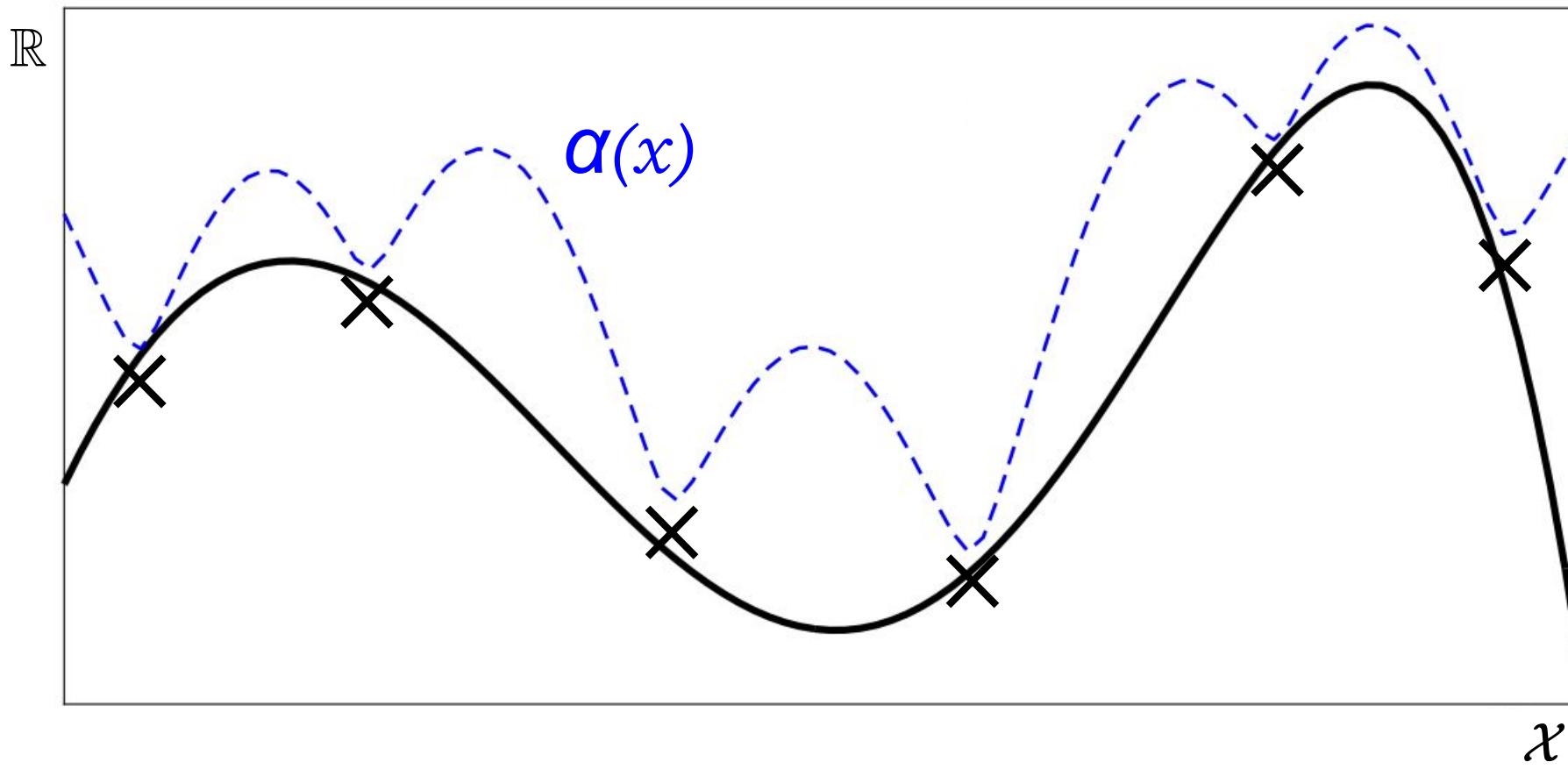


Define acquisition function using model.

# Bayesian Optimization – Visualization

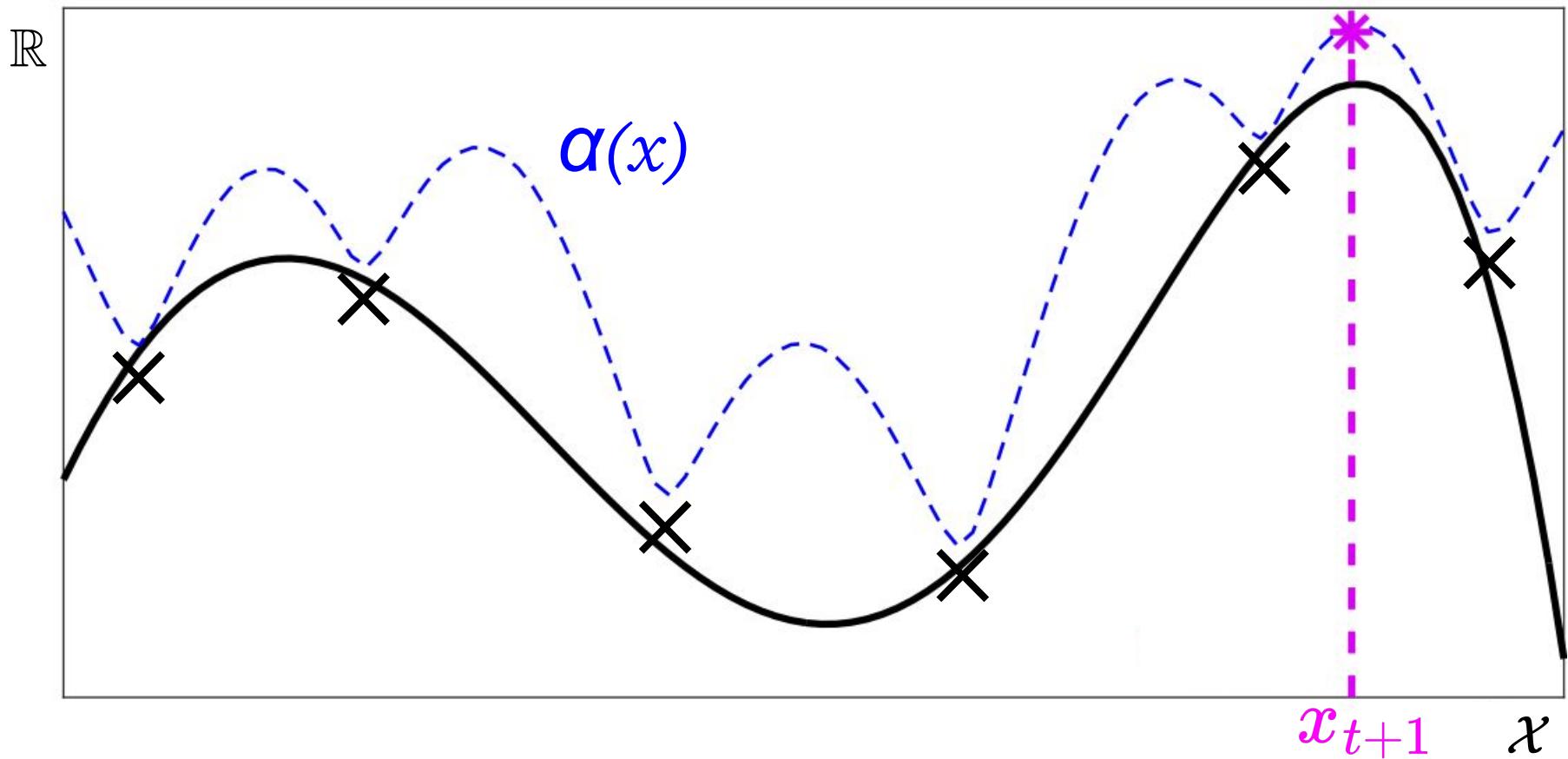


# Bayesian Optimization – Visualization



Define acquisition function using model.

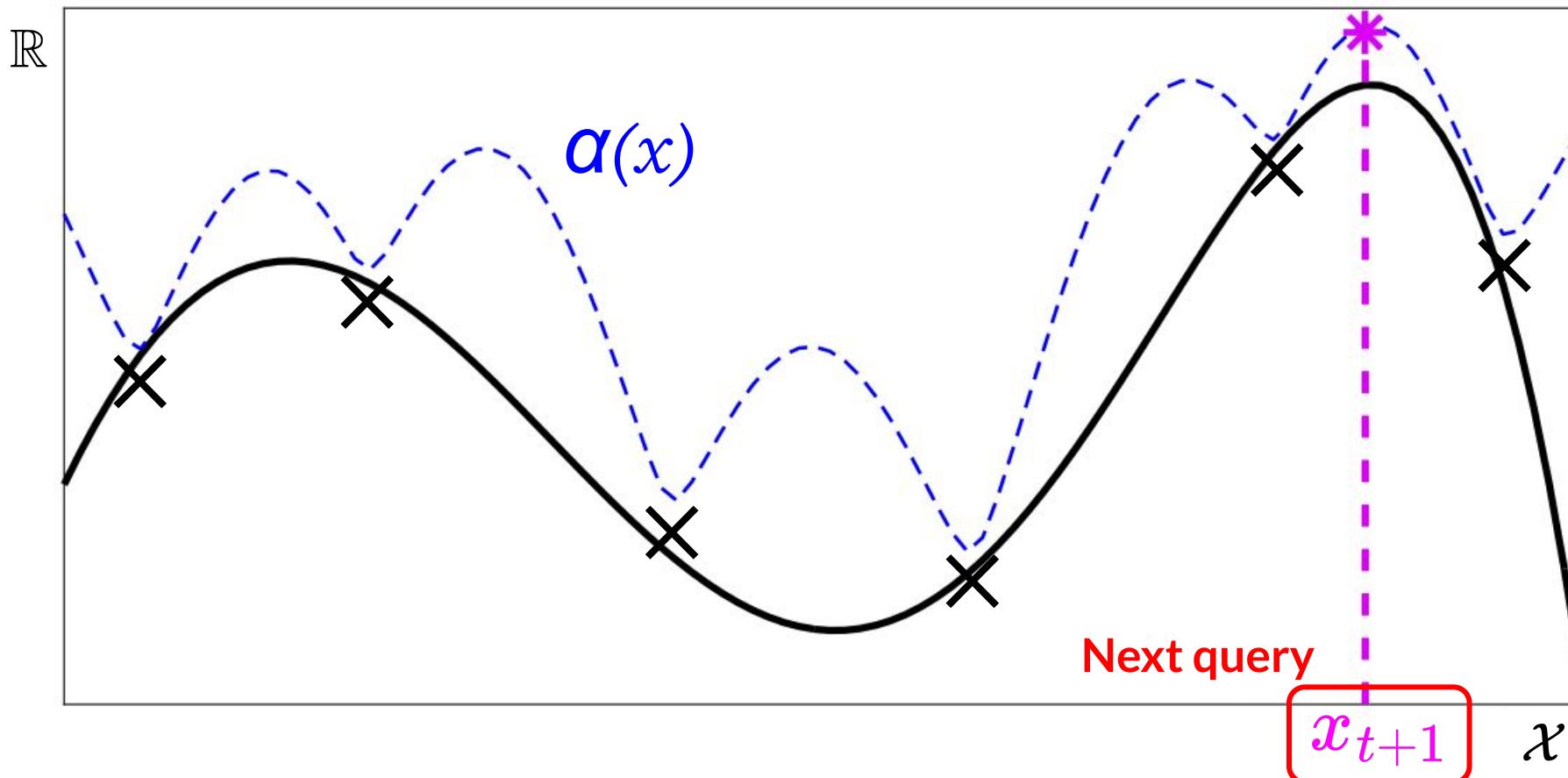
# Bayesian Optimization – Visualization



Optimize acquisition function  $\Rightarrow$  yields next point to query.

$$x_{t+1} = \arg \max_{x \in \mathcal{X}} \alpha(x)$$

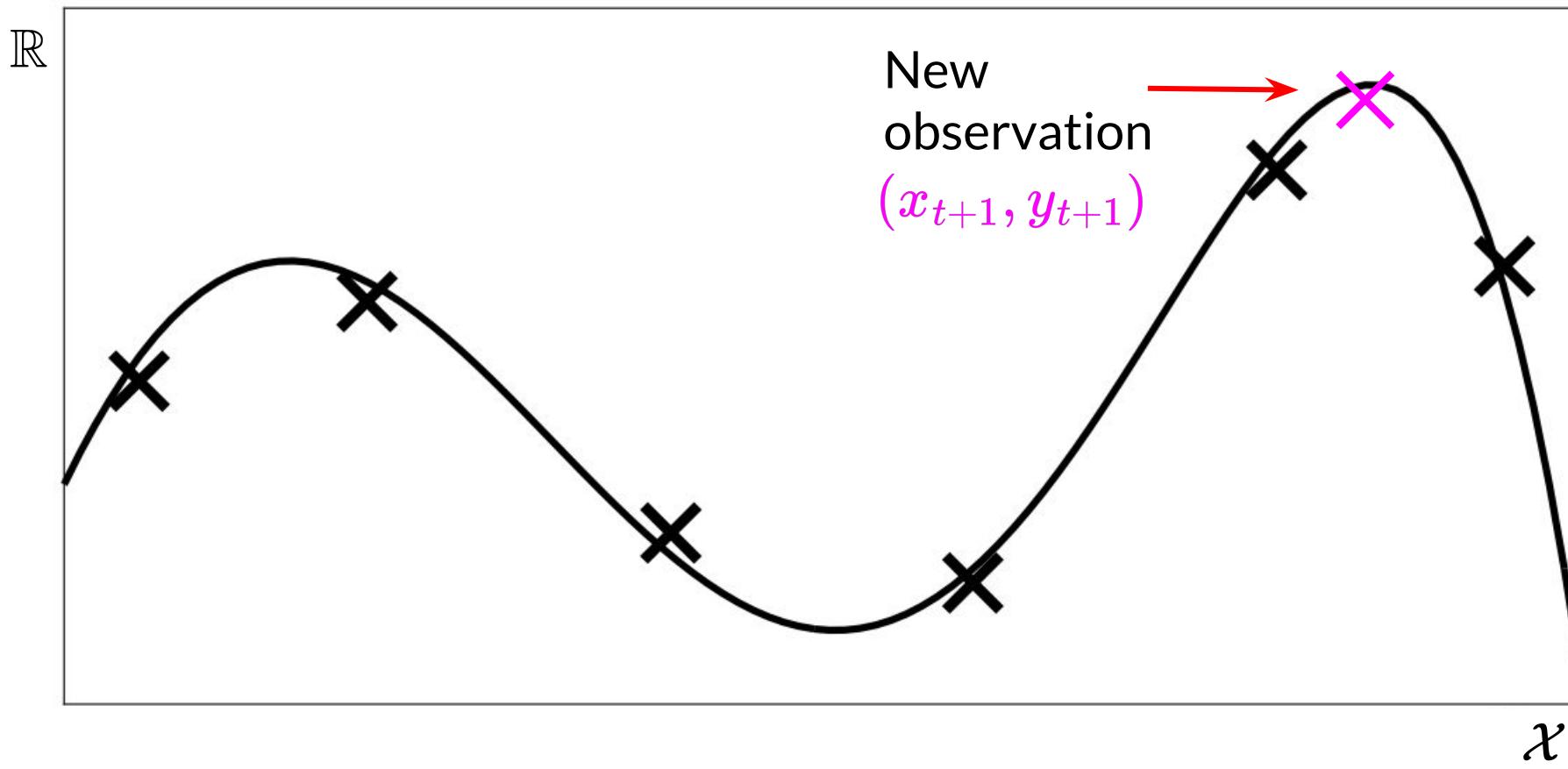
# Bayesian Optimization – Visualization



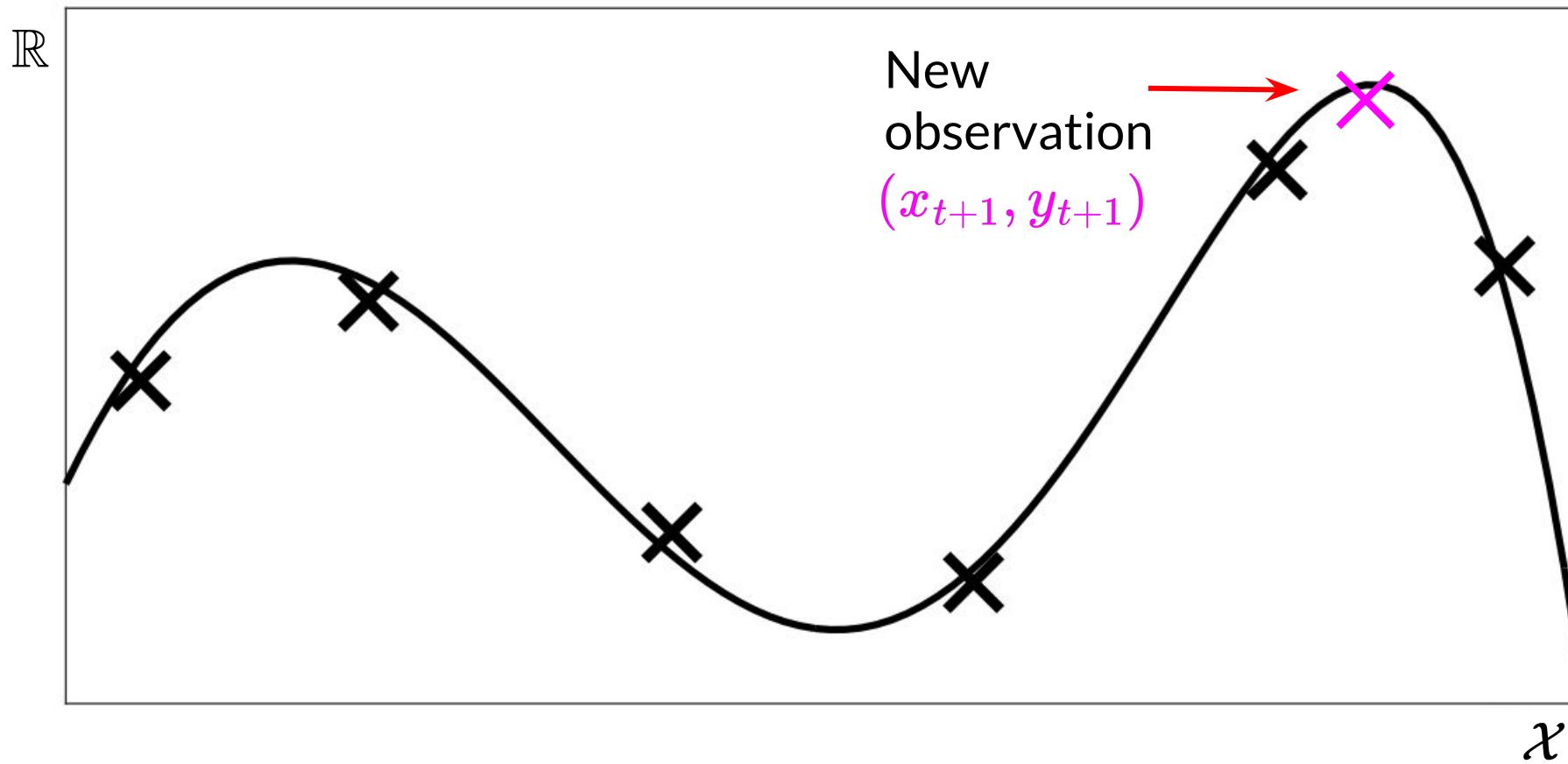
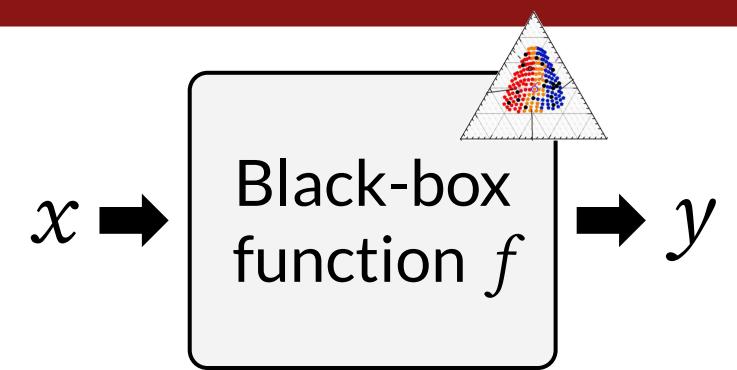
Optimize acquisition function  $\Rightarrow$  yields next point to query.

$$x_{t+1} = \arg \max_{x \in \mathcal{X}} \alpha(x)$$

# Bayesian Optimization – Visualization



# Bayesian Optimization – Visualization



Query black-box  $f$  at  $x$ , observe  $y$ , and update dataset.

# **Acquisition Function – Key Component of BO**

## Acquisition Function – Key Component of BO

Acquisition functions are, effectively, (one of) the main things that change between different BO procedures.

⇒ they dictate the BO algorithm, and choose the point queried at each iteration.

# Acquisition Function – Key Component of BO

Acquisition functions are, effectively, (one of) the main things that change between different BO procedures.

⇒ they dictate the BO algorithm, and choose the point queried at each iteration.

Goal when designing acquisition functions:

- Want to be “optimal” in some sense; to choose the point that is *most informative*.
- And/or: most-likely to accomplish goals of algorithm, e.g., carry out optimization.

## Acquisition Function – Key Component of BO

Next we will go through a few popular acquisition functions (for Bayes opt) in more detail:

## Acquisition Function – Key Component of BO

Next we will go through a few popular acquisition functions (for Bayes opt) in more detail:

**UCB** (*upper confidence bound*) – which we saw visualized before.

## Acquisition Function – Key Component of BO

Next we will go through a few popular acquisition functions (for Bayes opt) in more detail:

**UCB** (*upper confidence bound*) – which we saw visualized before.

**PI** (*probability of improvement*) – one of the simplest acquisition functions.

# Acquisition Function – Key Component of BO

Next we will go through a few popular acquisition functions (for Bayes opt) in more detail:

**UCB** (*upper confidence bound*) – which we saw visualized before.

**PI** (*probability of improvement*) – one of the simplest acquisition functions.

**EI** (*expected improvement*) – perhaps the most popular of the acquisition functions.

# Acquisition Function – Key Component of BO

Next we will go through a few popular acquisition functions (for Bayes opt) in more detail:

**UCB** (*upper confidence bound*) – which we saw visualized before.

**PI** (*probability of improvement*) – one of the simplest acquisition functions.

**EI** (*expected improvement*) – perhaps the most popular of the acquisition functions.

**KG** (*knowledge gradient*) – bit more complex, but performs well and is well-justified.

# Acquisition Function – Key Component of BO

Next we will go through a few popular acquisition functions (for Bayes opt) in more detail:

**UCB** (*upper confidence bound*) – which we saw visualized before.

**PI** (*probability of improvement*) – one of the simplest acquisition functions.

**EI** (*expected improvement*) – perhaps the most popular of the acquisition functions.

**KG** (*knowledge gradient*) – bit more complex, but performs well and is well-justified.

**ES** (*entropy search*) – rooted in information theory and *Bayesian optimal experimental design (BOED)*.

# Acquisition Function – Key Component of BO

Next we will go through a few popular acquisition functions (for Bayes opt) in more detail:

**UCB** (*upper confidence bound*) – which we saw visualized before.

**PI** (*probability of improvement*) – one of the simplest acquisition functions.

**EI** (*expected improvement*) – perhaps the most popular of the acquisition functions.

**KG** (*knowledge gradient*) – bit more complex, but performs well and is well-justified.

**ES** (*entropy search*) – rooted in information theory and *Bayesian optimal experimental design (BOED)*.

**But first, a quick discussion on terminology!**

# **Terminology – Active Learning, Bayes Opt, Experimental Design, etc.**

## Terminology – Active Learning, Bayes Opt, Experimental Design, etc.

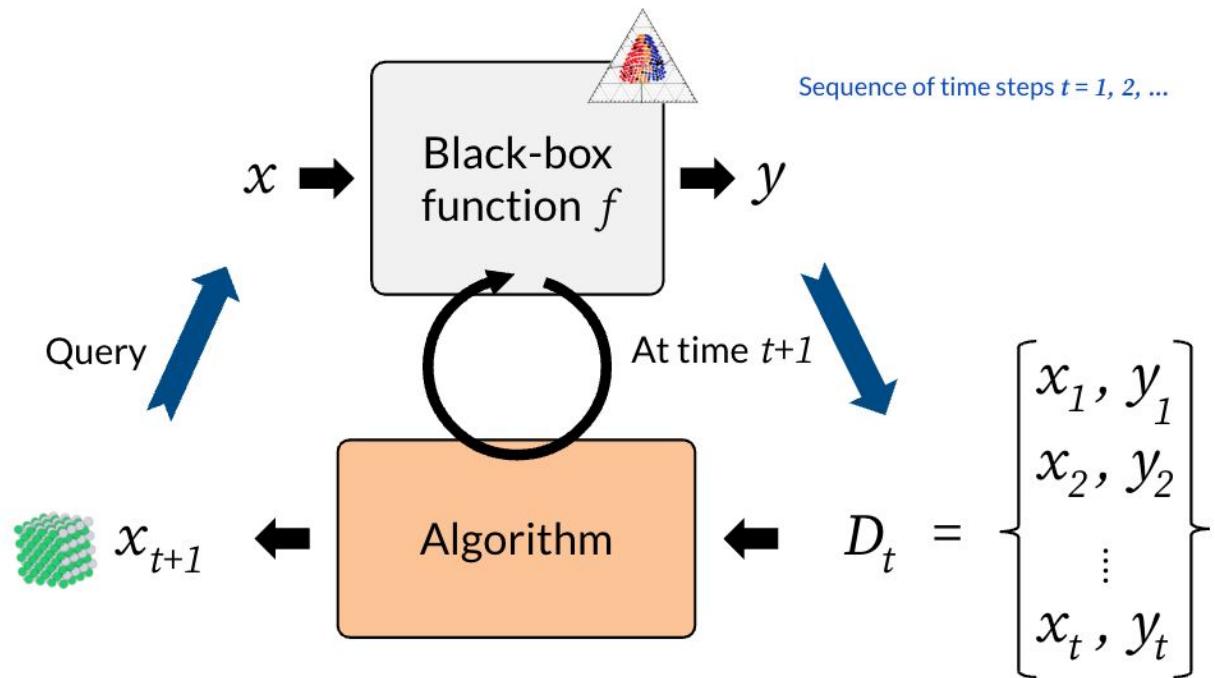
These terms all refer to related procedures, but there are slight differences between them...

**Similarities:** these all refer to sequential decision making (data acquisition) procedures.

# Terminology – Active Learning, Bayes Opt, Experimental Design, etc.

These terms all refer to related procedures, but there are slight differences between them...

**Similarities:** these all refer to sequential decision making (data acquisition) procedures.



# **Terminology – Active Learning, Bayes Opt, Experimental Design, etc.**

These terms all refer to related procedures, but there are slight differences between them...

**Differences:**

# Terminology – Active Learning, Bayes Opt, Experimental Design, etc.

These terms all refer to related procedures, but there are slight differences between them...

## Differences:

### Active Learning

Originally from CS/ML.

*Classically:* iteratively choosing points to label from a set.

*More recently:* any learning method where you choose a data point to observe at each iteration .

# Terminology – Active Learning, Bayes Opt, Experimental Design, etc.

These terms all refer to related procedures, but there are slight differences between them...

## Differences:

### Active Learning

Originally from CS/ML.

*Classically:* iteratively choosing points to label from a set.

*More recently:* any learning method where you choose a data point to observe at each iteration .

### Experimental Design

Originally from statistics.

*Classically:* goal is to efficiently estimate a model parameter (often by reducing posterior entropy).

*Can view parameter estimation as one strategy for active learning/optimization.*

# Terminology – Active Learning, Bayes Opt, Experimental Design, etc.

These terms all refer to related procedures, but there are slight differences between them...

## Differences:

### Active Learning

Originally from CS/ML.

*Classically:* iteratively choosing points to label from a set.

*More recently:* any learning method where you choose a data point to observe at each iteration .

### Experimental Design

Originally from statistics.

*Classically:* goal is to efficiently estimate a model parameter (often by reducing posterior entropy).

*Can view parameter estimation as one strategy for active learning/optimization.*

### Bayesian Optimization

Originally from operations research (among other fields).

*Can view as:* a special case of active learning, aiming to carry out the task of optimization.

(Which could potentially use an experimental design strategy, e.g., in *entropy search*).

# Terminology – Active Learning, Bayes Opt, Experimental Design, etc.

A bit more on experimental design:

## Terminology – Active Learning, Bayes Opt, Experimental Design, etc.

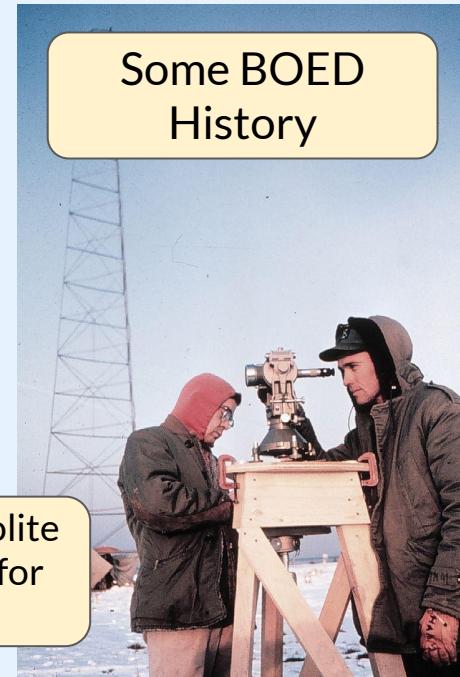
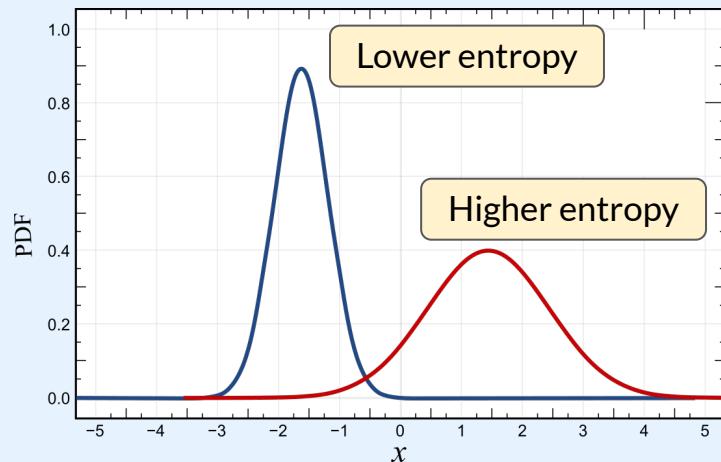
A bit more on experimental design: often referred to as *information-based experimental design* (if it involves reducing posterior entropy), or *Bayesian optimal experimental design (BOED)*.

# Terminology – Active Learning, Bayes Opt, Experimental Design, etc.

A bit more on experimental design: often referred to as *information-based experimental design* (if it involves reducing posterior entropy), or *Bayesian optimal experimental design (BOED)*.

**BOED:** have model with an (unknown) parameter of interest.

- Choose experiments that most reduce uncertainty about parameter.
- **Uncertainty:** **entropy** of posterior distribution over parameter.



Reducing theodolite measurements for surveying.

## Terminology – Active Learning, Bayes Opt, Experimental Design, etc.

A bit more on experimental design: often referred to as *information-based experimental design* (if it involves reducing posterior entropy), or *Bayesian optimal experimental design (BOED)*.

Relatedly – there is *information-based Bayesian optimization*:

**Info-based BO:** entropy search (ES)<sup>1</sup>, predictive ES<sup>2</sup>, max-value ES<sup>3</sup>.

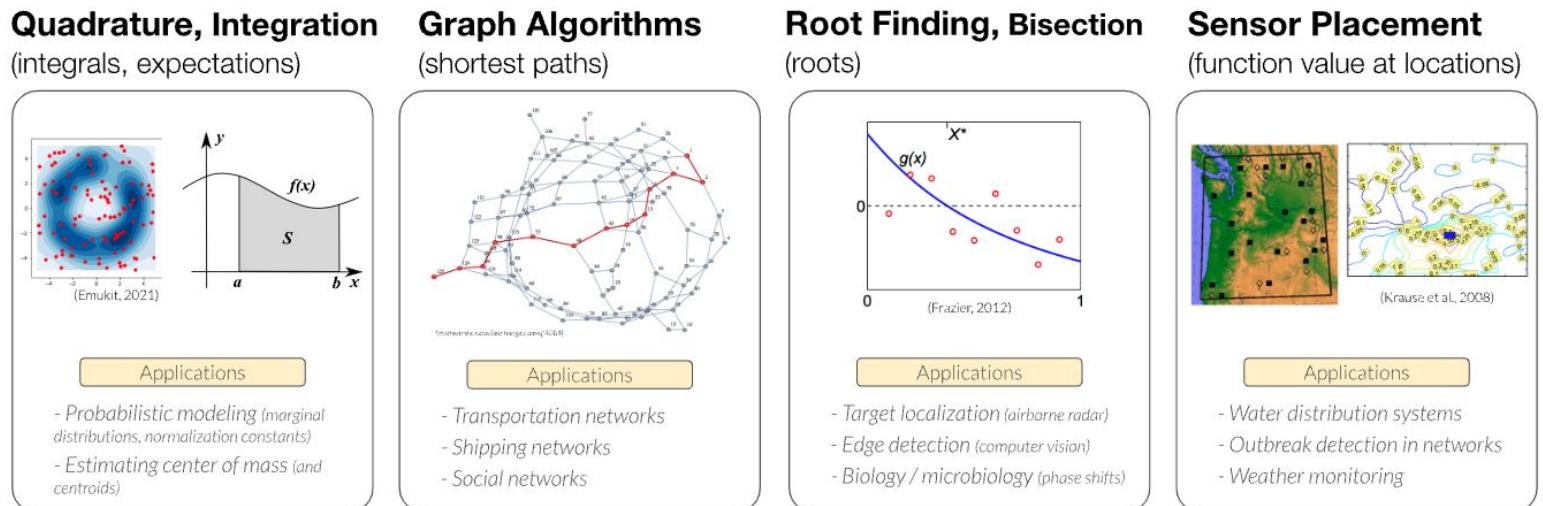
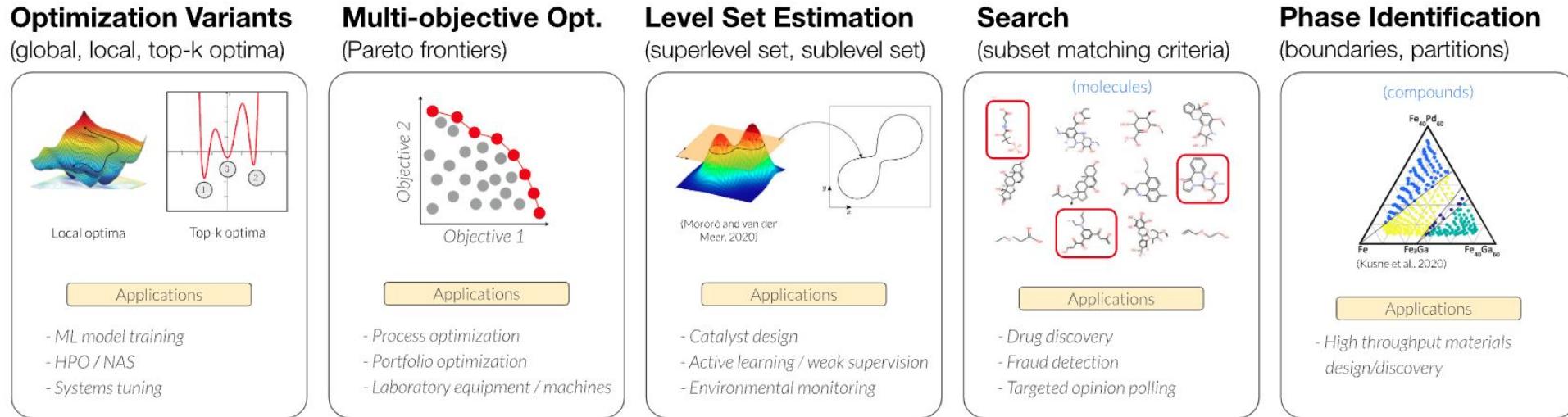
- Rooted in **Bayesian optimal experimental design (BOED)**.

Which we will talk about a little later...

## Terminology – Active Learning, Bayes Opt, Experimental Design, etc.

*Side note:* there are many other potential special cases of active learning beyond global optimization, as in BO...

# Terminology – Active Learning, Bayes Opt, Experimental Design, etc.



# Acquisition Functions – Back to Acquisition Functions

# Acquisition Functions – Back to Acquisition Functions

## Acquisition Function

Component of Bayesian optimization algorithm.

$$\alpha : \mathcal{X} \rightarrow \mathbb{R}$$

Assigns a number describing the value/usefulness of querying a given point.

# Acquisition Functions – Upper Confidence Bound (UCB)

## Acquisition Functions – Upper Confidence Bound (UCB)

The UCB acquisition function is often written:

$$\alpha_t(x) = \mu(x) + \beta_t \sigma(x)$$

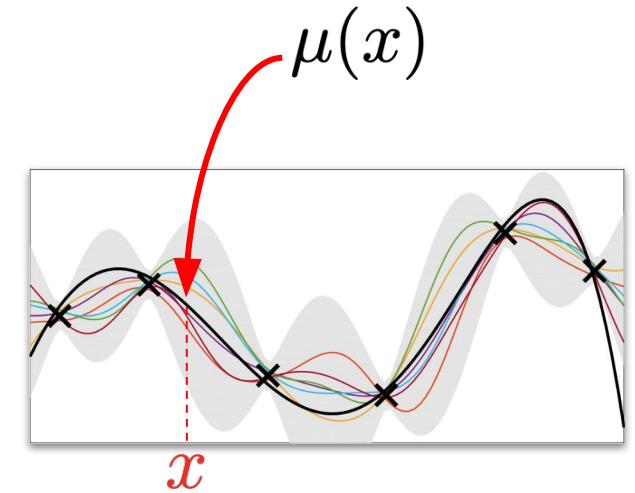
# Acquisition Functions – Upper Confidence Bound (UCB)

The UCB acquisition function is often written:

$$\alpha_t(x) = \mu(x) + \beta_t \sigma(x)$$

where:

$\mu(x)$  - Mean of *marginal posterior* at input  $x$ .



# Acquisition Functions – Upper Confidence Bound (UCB)

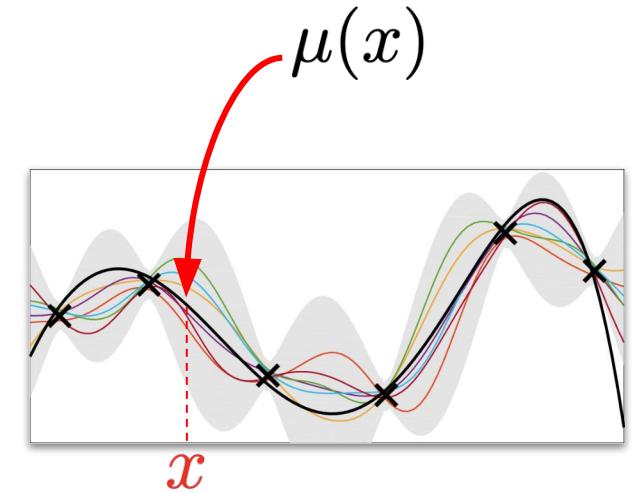
The UCB acquisition function is often written:

$$\alpha_t(x) = \mu(x) + \beta_t \sigma(x)$$

where:

$$p(f(x) | \mathcal{D}_t)$$

$\mu(x)$  - Mean of *marginal posterior* at input  $x$ .



# Acquisition Functions – Upper Confidence Bound (UCB)

The UCB acquisition function is often written:

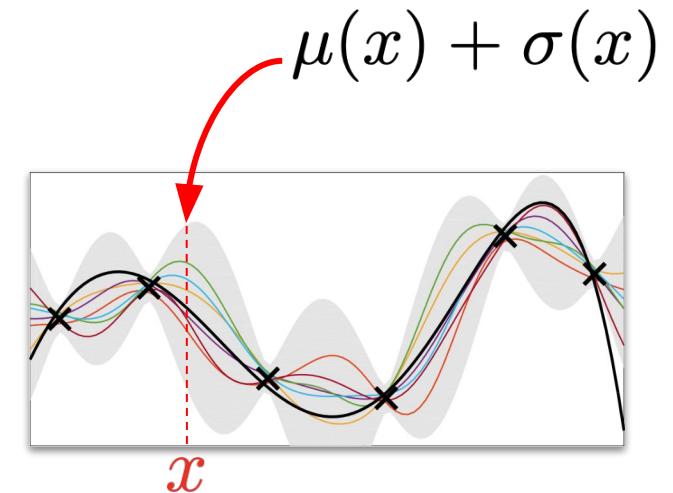
$$\alpha_t(x) = \mu(x) + \beta_t \sigma(x)$$

where:

$$p(f(x) | \mathcal{D}_t)$$

$\mu(x)$  - Mean of *marginal posterior* at input  $x$ .

$\sigma(x)$  - Standard deviation of *marginal posterior* at input  $x$ .



# Acquisition Functions – Upper Confidence Bound (UCB)

The UCB acquisition function is often written:

$$\alpha_t(x) = \mu(x) + \beta_t \sigma(x)$$

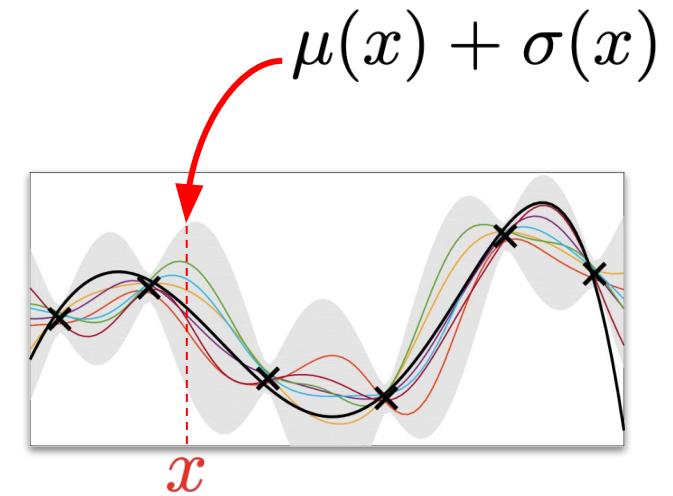
where:

$$p(f(x) | \mathcal{D}_t)$$

$\mu(x)$  - Mean of *marginal posterior* at input  $x$ .

$\sigma(x)$  - Standard deviation of *marginal posterior* at input  $x$ .

$\beta_t$  - Hyperparameter controlling tradeoff between exploration and exploitation.

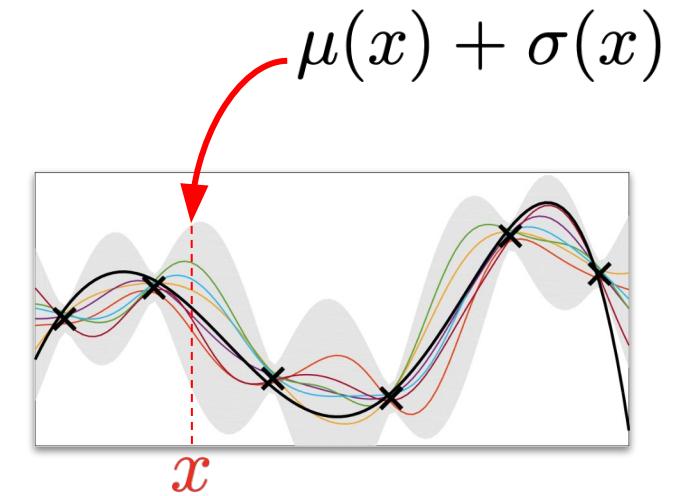


# Acquisition Functions – Upper Confidence Bound (UCB)

The UCB acquisition function is often written:

$$\alpha_t(x) = \mu(x) + \beta_t \sigma(x)$$

Intuition:



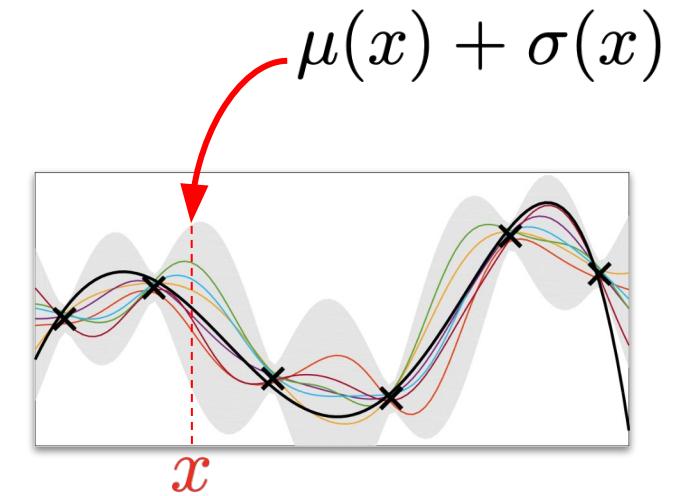
# Acquisition Functions – Upper Confidence Bound (UCB)

The UCB acquisition function is often written:

$$\alpha_t(x) = \mu(x) + \beta_t \sigma(x)$$

Intuition:

Maximizing variance  $\rightarrow$  queries *most unknown* part of function  
 $\Rightarrow$  most “informative”.



# Acquisition Functions – Upper Confidence Bound (UCB)

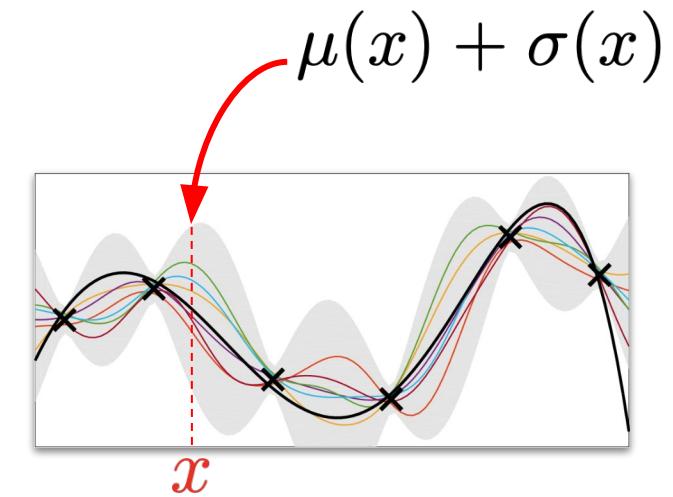
The UCB acquisition function is often written:

$$\alpha_t(x) = \mu(x) + \beta_t \sigma(x)$$

Intuition:

Maximizing variance → queries *most unknown* part of function  
⇒ most “informative”.

Maximizing mean → queries *best guess of highest value*, based on current knowledge of function.



# Acquisition Functions – Upper Confidence Bound (UCB)

The UCB acquisition function is often written:

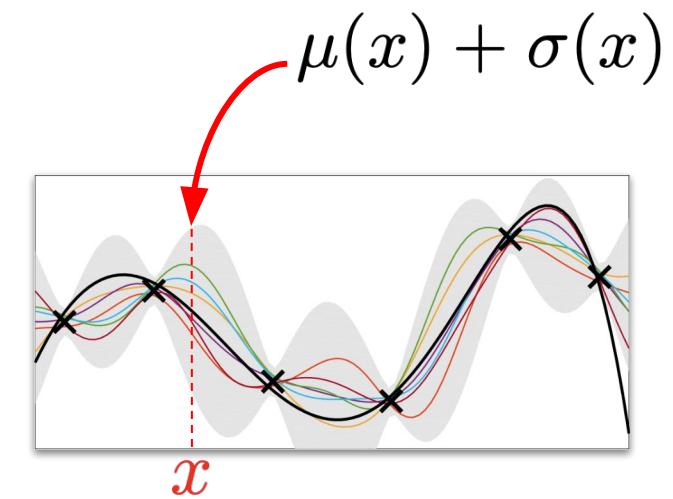
$$\alpha_t(x) = \mu(x) + \beta_t \sigma(x)$$

Intuition:

Maximizing variance → queries *most unknown* part of function  
⇒ most “informative”.

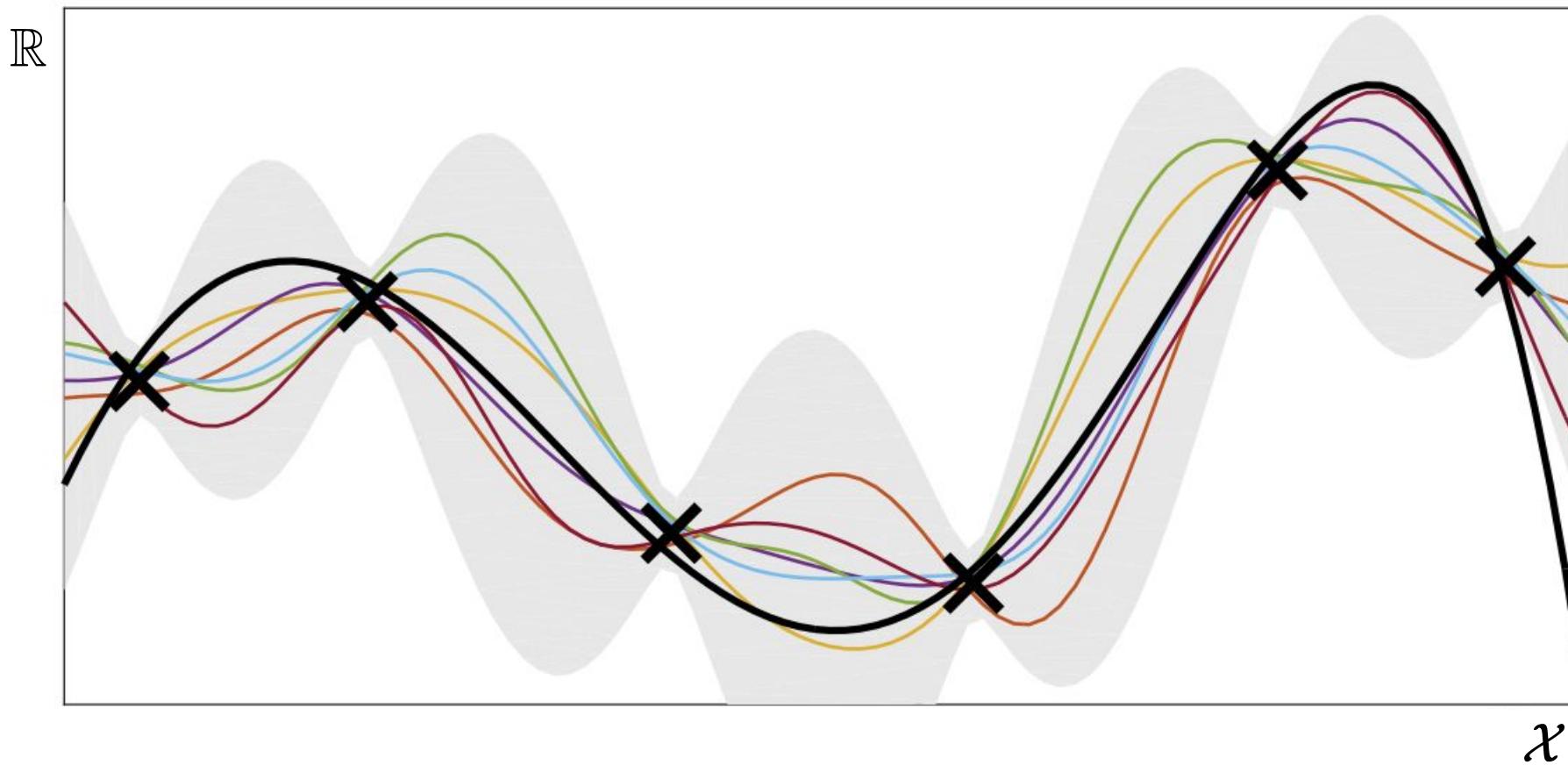
Maximizing mean → queries *best guess of highest value*, based on current knowledge of function.

We want to trade off between these two quantities.



# Acquisition Functions – Upper Confidence Bound (UCB)

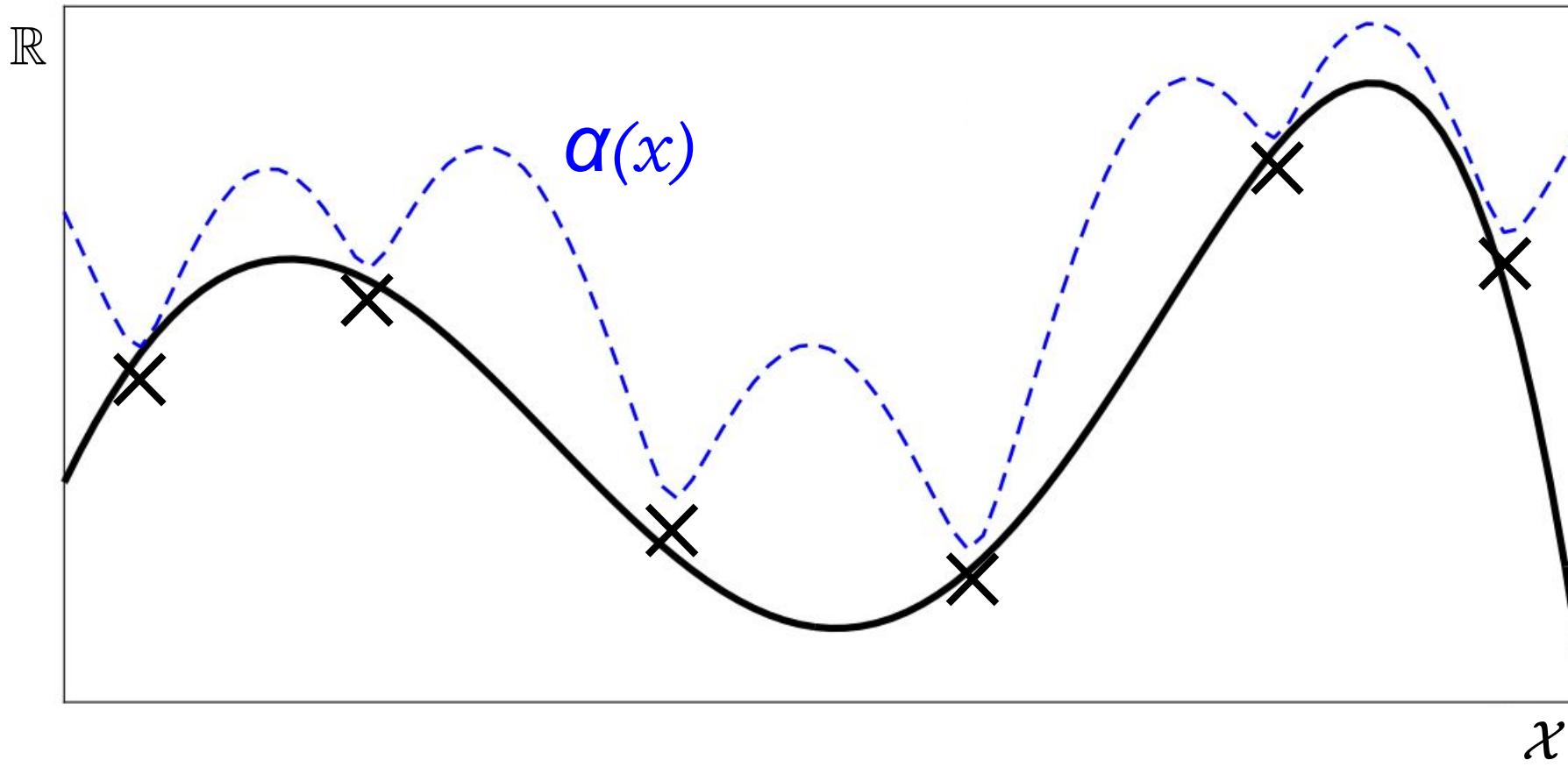
Visualizing UCB



Given our probabilistic model...

# Acquisition Functions – Upper Confidence Bound (UCB)

Visualizing UCB



... define the **UCB** acquisition function.

## Acquisition Functions – Experimental Design (EIG about $f$ )

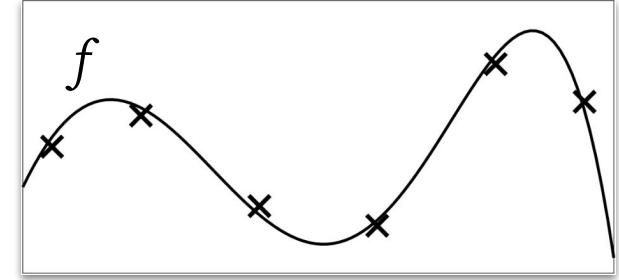
Quick aside on experimental design, to give some perspective on UCB:

## Acquisition Functions – Experimental Design (EIG about $f$ )

Quick aside on experimental design, to give some perspective on UCB:

Suppose you don't want to do optimization, but instead want to learn the full function landscape (i.e., learn "all of  $f$ ").

⇒ technically speaking, this is not *optimization* anymore.

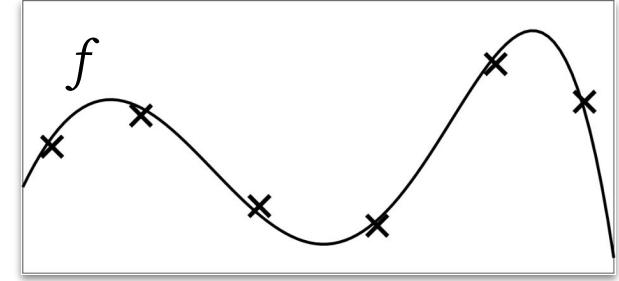


## Acquisition Functions – Experimental Design (EIG about $f$ )

Quick aside on experimental design, to give some perspective on UCB:

Suppose you don't want to do optimization, but instead want to learn the full function landscape (i.e., learn "all of  $f$ ").

⇒ technically speaking, this is not *optimization* anymore.



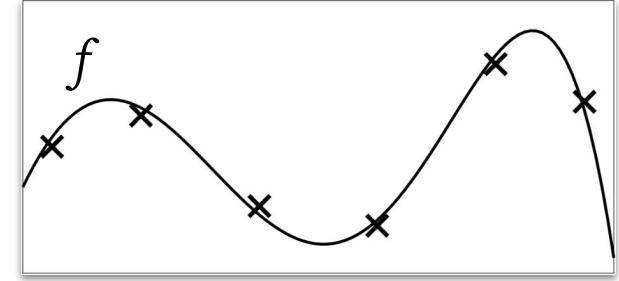
An *experimental design approach* would be to query the point that **maximally reduces the posterior entropy about  $f$** .

# Acquisition Functions – Experimental Design (EIG about $f$ )

Quick aside on experimental design, to give some perspective on UCB:

Suppose you don't want to do optimization, but instead want to learn the full function landscape (i.e., learn "all of  $f$ ").

⇒ technically speaking, this is not *optimization* anymore.



An *experimental design approach* would be to query the point that **maximally reduces the posterior entropy about  $f$** .

(This is equivalently also commonly called: *expected information gain*, or *EIG, about  $f$* .)

## Acquisition Functions – Experimental Design (EIG about $f$ )

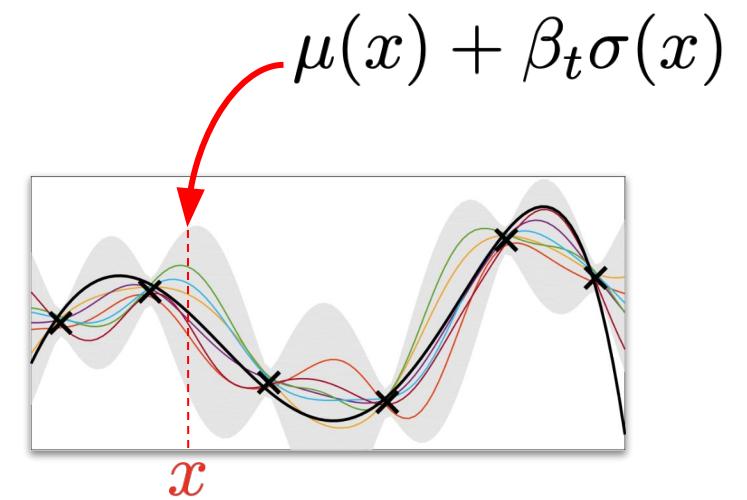
You can prove (at least for GP models) that maximally reducing posterior entropy is equivalent to querying point with maximum standard deviation (of posterior marginal).

## Acquisition Functions – Experimental Design (EIG about $f$ )

You can prove (at least for GP models) that maximally reducing posterior entropy is equivalent to querying point with maximum standard deviation (of posterior marginal).

So in UCB  $\Rightarrow$

The second term (standard deviation term) is in fact balancing exploration (EIG about  $f$ ) with exploitation (best expected input  $x$  ).



# Acquisition Functions – Probability of Improvement (PI)

## **Acquisition Functions – Probability of Improvement (PI)**

Next is one of the simplest acquisition functions (though not used a lot in practice).

## Acquisition Functions – Probability of Improvement (PI)

Next is one of the simplest acquisition functions (though not used a lot in practice).

The probability of improvement (PI) acquisition function is written:

## Acquisition Functions – Probability of Improvement (PI)

Next is one of the simplest acquisition functions (though not used a lot in practice).

The probability of improvement (PI) acquisition function is written:

$$\alpha_t(x) = p(f(x) > f(x^*) \mid \mathcal{D}_t)$$

## Acquisition Functions – Probability of Improvement (PI)

Next is one of the simplest acquisition functions (though not used a lot in practice).

The probability of improvement (PI) acquisition function is written:

$$\alpha_t(x) = p(f(x) > f(x^*) \mid \mathcal{D}_t)$$

where  $f(x^*)$  is the function value of the best point observed so far in  $\mathcal{D}_t$ .

*I.e., the posterior probability that the  $f(x)$  is greater than the best point observed so far.*

## Acquisition Functions – Probability of Improvement (PI)

Next is one of the simplest acquisition functions (though not used a lot in practice).

The probability of improvement (PI) acquisition function is written:

$$\alpha_t(x) = p(f(x) > f(x^*) \mid \mathcal{D}_t)$$

where  $f(x^*)$  is the function value of the best point observed so far in  $\mathcal{D}_t$ .

*I.e., the posterior probability that the  $f(x)$  is greater than the best point observed so far.*

You can prove this is equal to ...

## Acquisition Functions – Probability of Improvement (PI)

Next is one of the simplest acquisition functions (though not used a lot in practice).

The probability of improvement (PI) acquisition function is written:

$$\begin{aligned}\alpha_t(x) &= p(f(x) > f(x^*) \mid \mathcal{D}_t) \\ &= \Phi\left(\frac{\mu(x) - f(x^*)}{\sigma(x)}\right)\end{aligned}$$

## Acquisition Functions – Probability of Improvement (PI)

Next is one of the simplest acquisition functions (though not used a lot in practice).

The probability of improvement (PI) acquisition function is written:

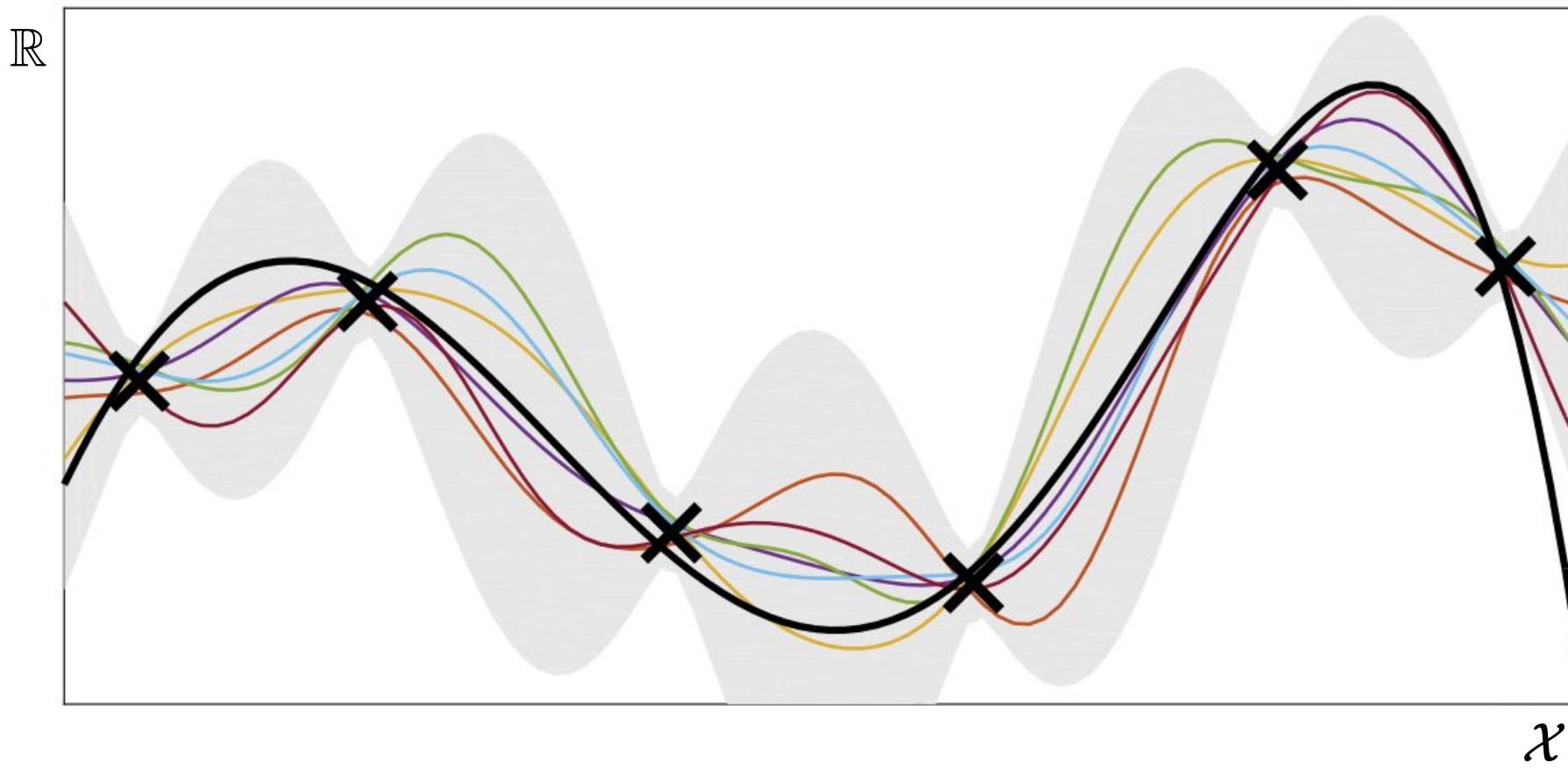
$$\begin{aligned}\alpha_t(x) &= p(f(x) > f(x^*) \mid \mathcal{D}_t) \\ &= \Phi\left(\frac{\mu(x) - f(x^*)}{\sigma(x)}\right)\end{aligned}$$

where  $\Phi(\cdot)$  denotes the CDF of the standard normal distribution.

(And, as before,  $\mu(x)$  and  $\sigma(x)$  are the mean and SD of marginal posterior at  $x$ .)

# Acquisition Functions – Probability of Improvement (PI)

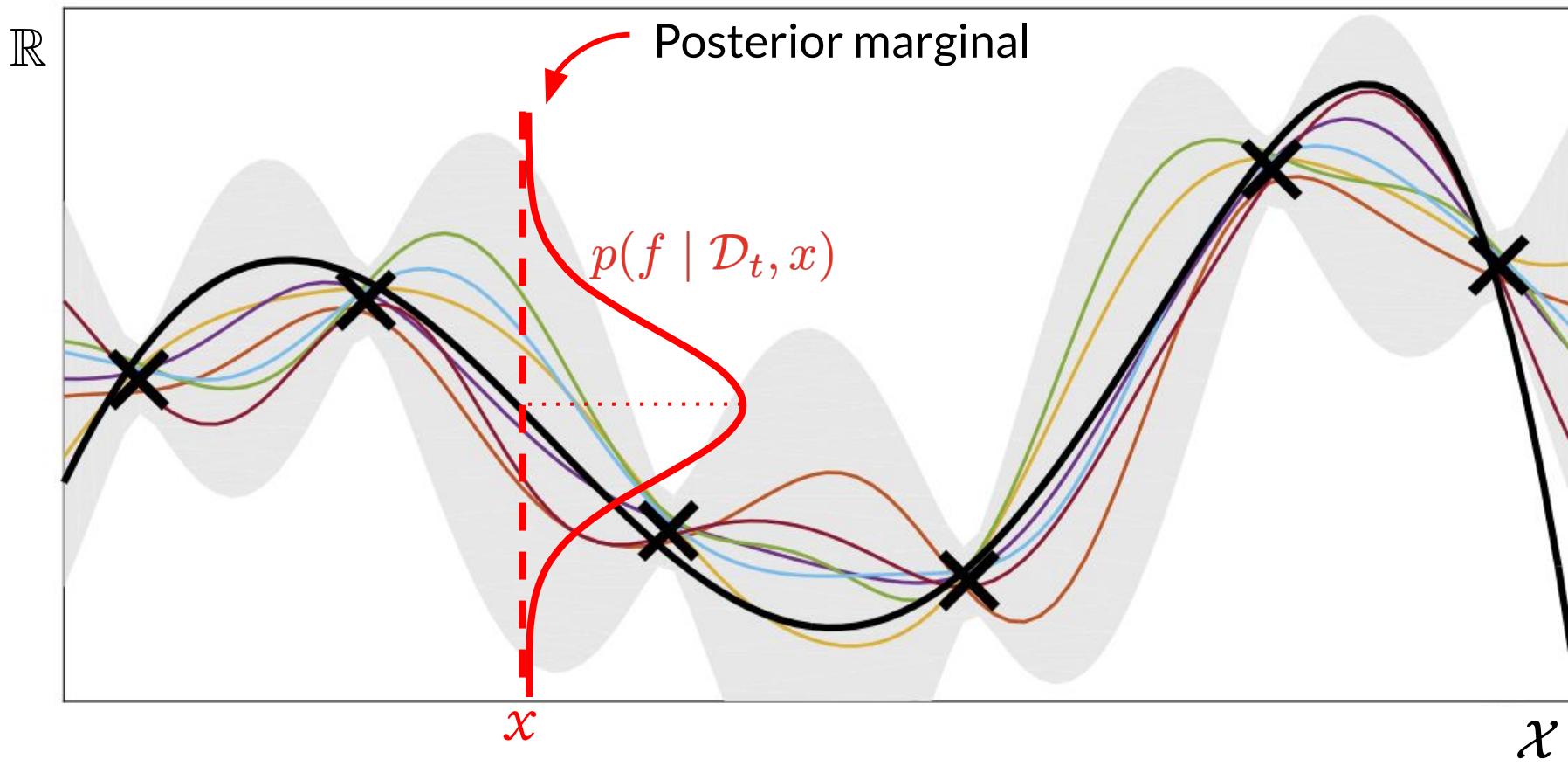
Visualizing PI



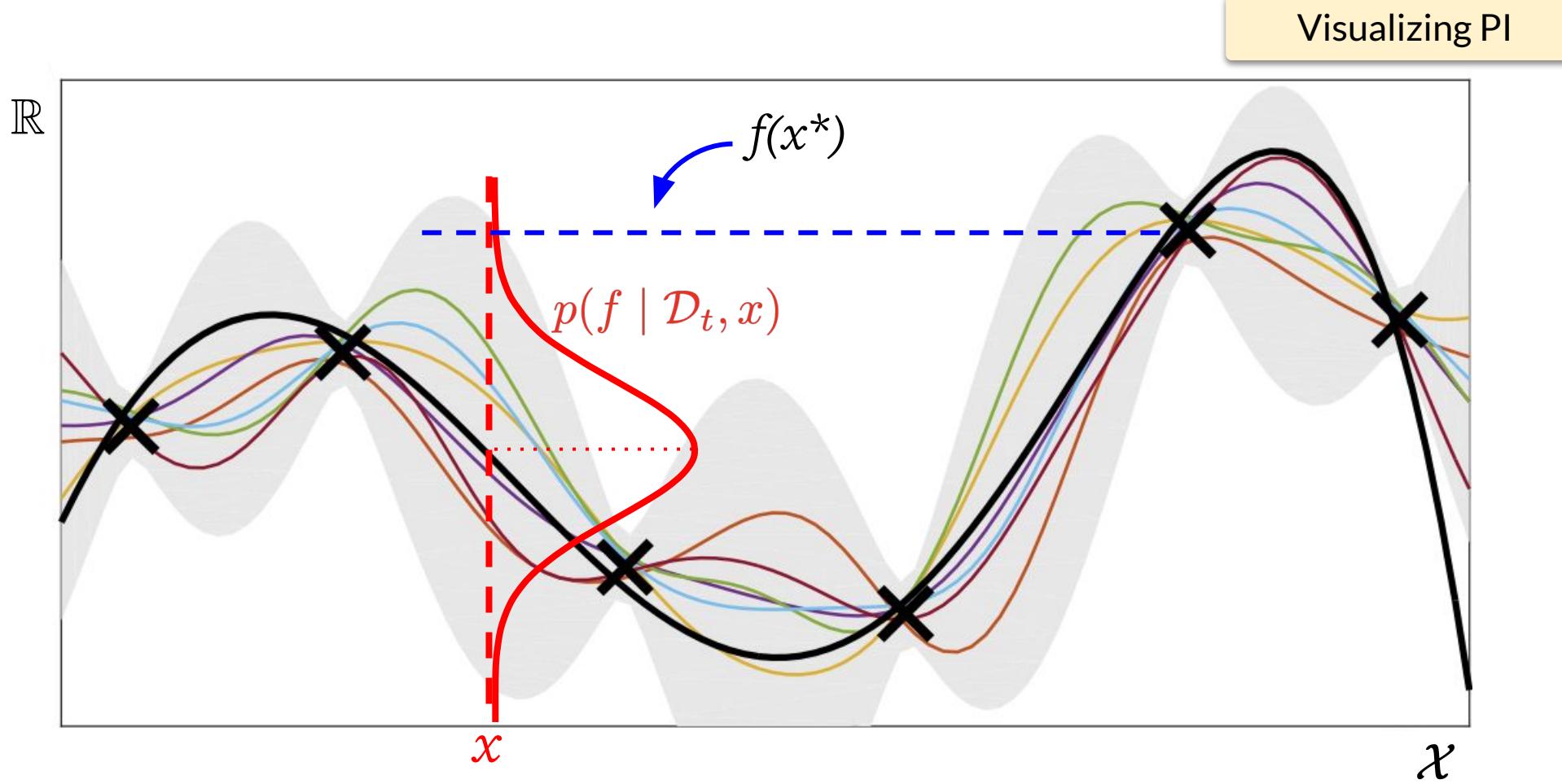
Given our probabilistic model...

# Acquisition Functions – Probability of Improvement (PI)

Visualizing PI

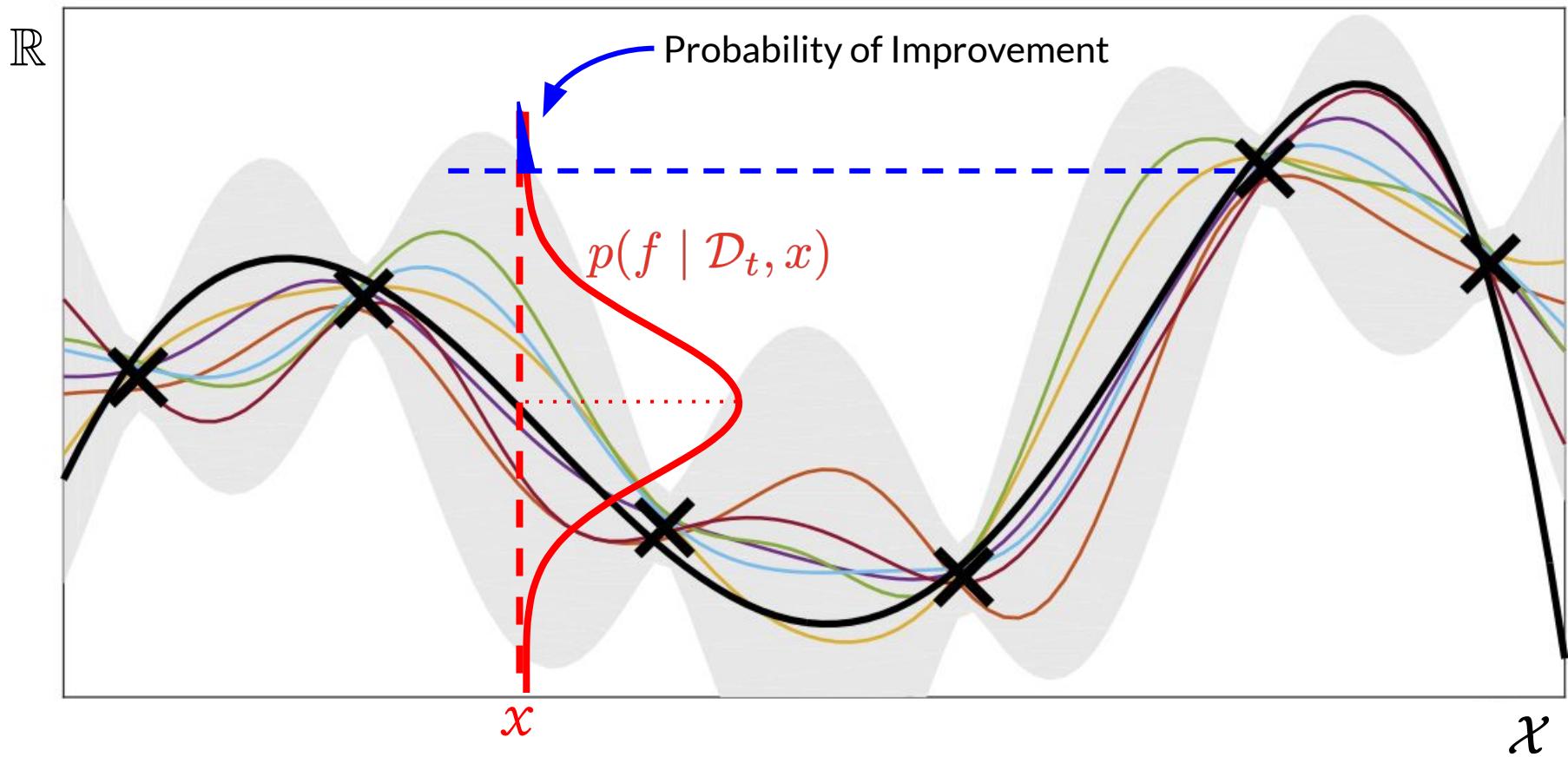


# Acquisition Functions – Probability of Improvement (PI)



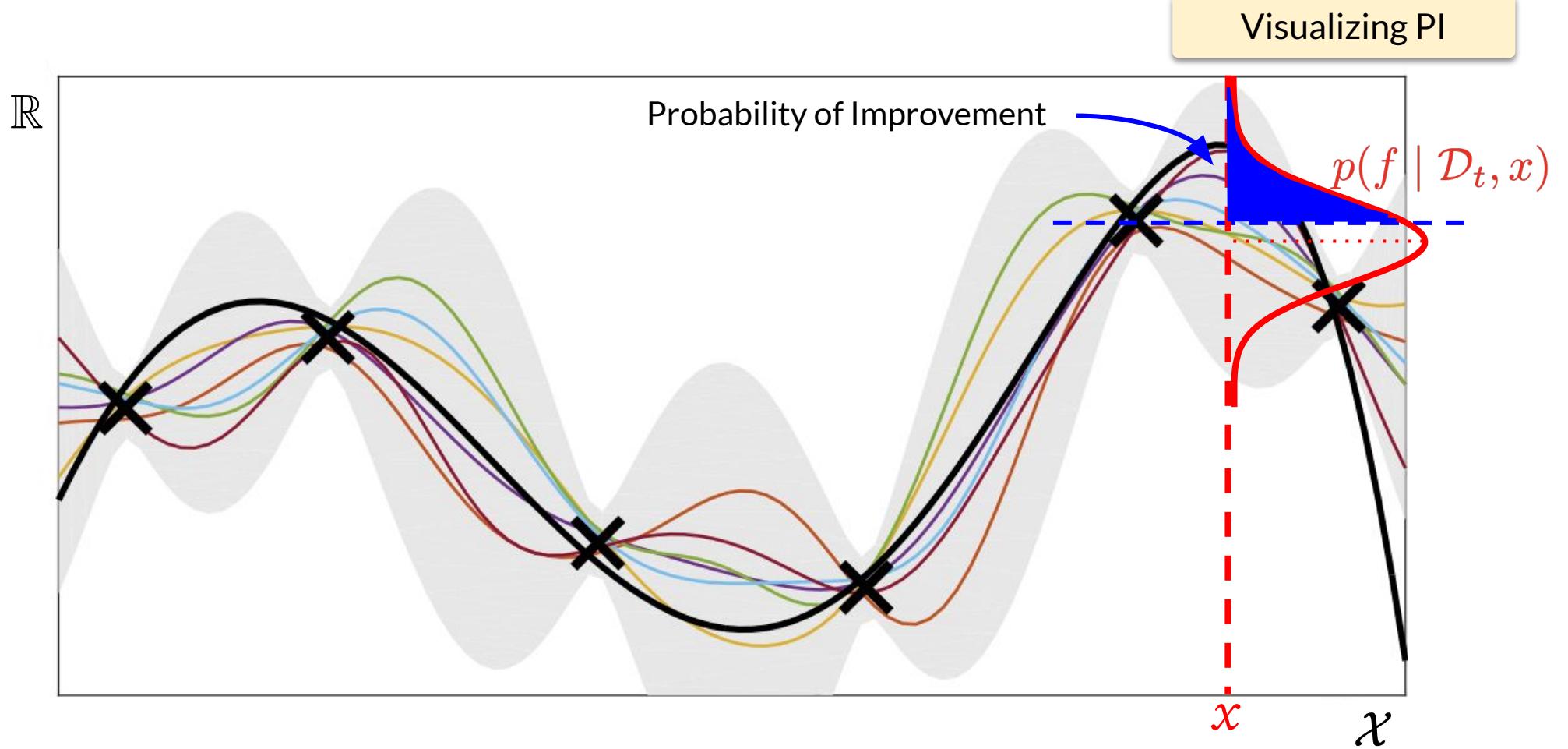
# Acquisition Functions – Probability of Improvement (PI)

Visualizing PI



... define the PI acquisition function.

# Acquisition Functions – Probability of Improvement (PI)



## **Acquisition Functions – Probability of Improvement (PI)**

However, note that:

## Acquisition Functions – Probability of Improvement (PI)

However, note that:

PI maximizes the **probability** that  $f(x)$  is greater than then best point so far.

Does **\*not\*** try to maximize the **amount** that  $f(x)$  is greater than the best point so far.

## Acquisition Functions – Probability of Improvement (PI)

However, note that:

PI maximizes the **probability** that  $f(x)$  is greater than then best point so far.

Does **\*not\*** try to maximize the **amount** that  $f(x)$  is greater than the best point so far.

(i.e., it might choose an  $x$  that is only *slightly better* than best point so far, if it is more confident about it – rather than a point that is *significantly greater in expectation*, but it is less confident about it).

## Acquisition Functions – Probability of Improvement (PI)

However, note that:

PI maximizes the **probability** that  $f(x)$  is greater than then best point so far.

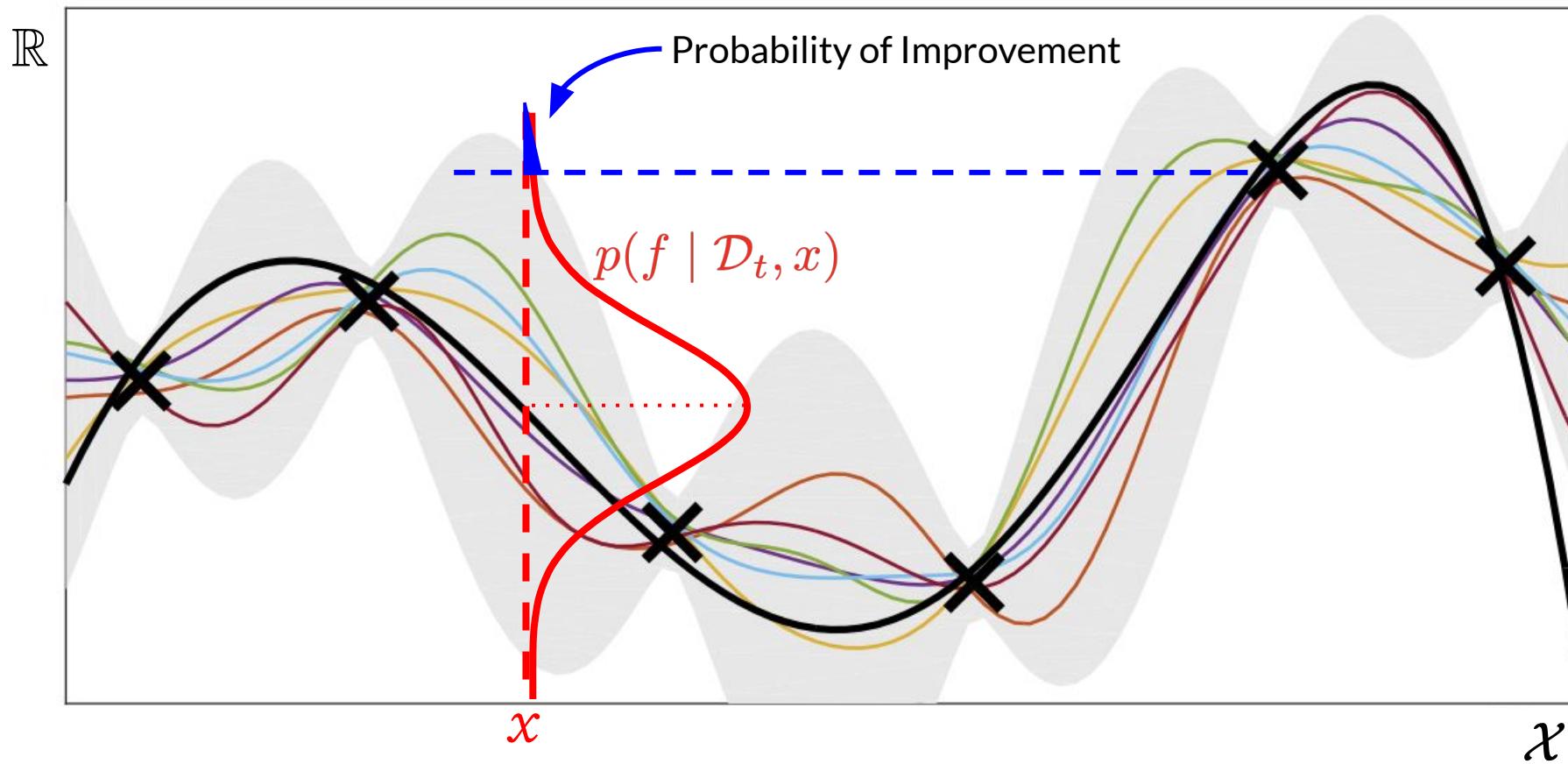
Does \*not\* try to maximize the **amount** that  $f(x)$  is greater than the best point so far.

(i.e., it might choose an  $x$  that is only *slightly better* than best point so far, if it is more confident about it – rather than a point that is *significantly greater in expectation*, but it is less confident about it).

Visualizing this...

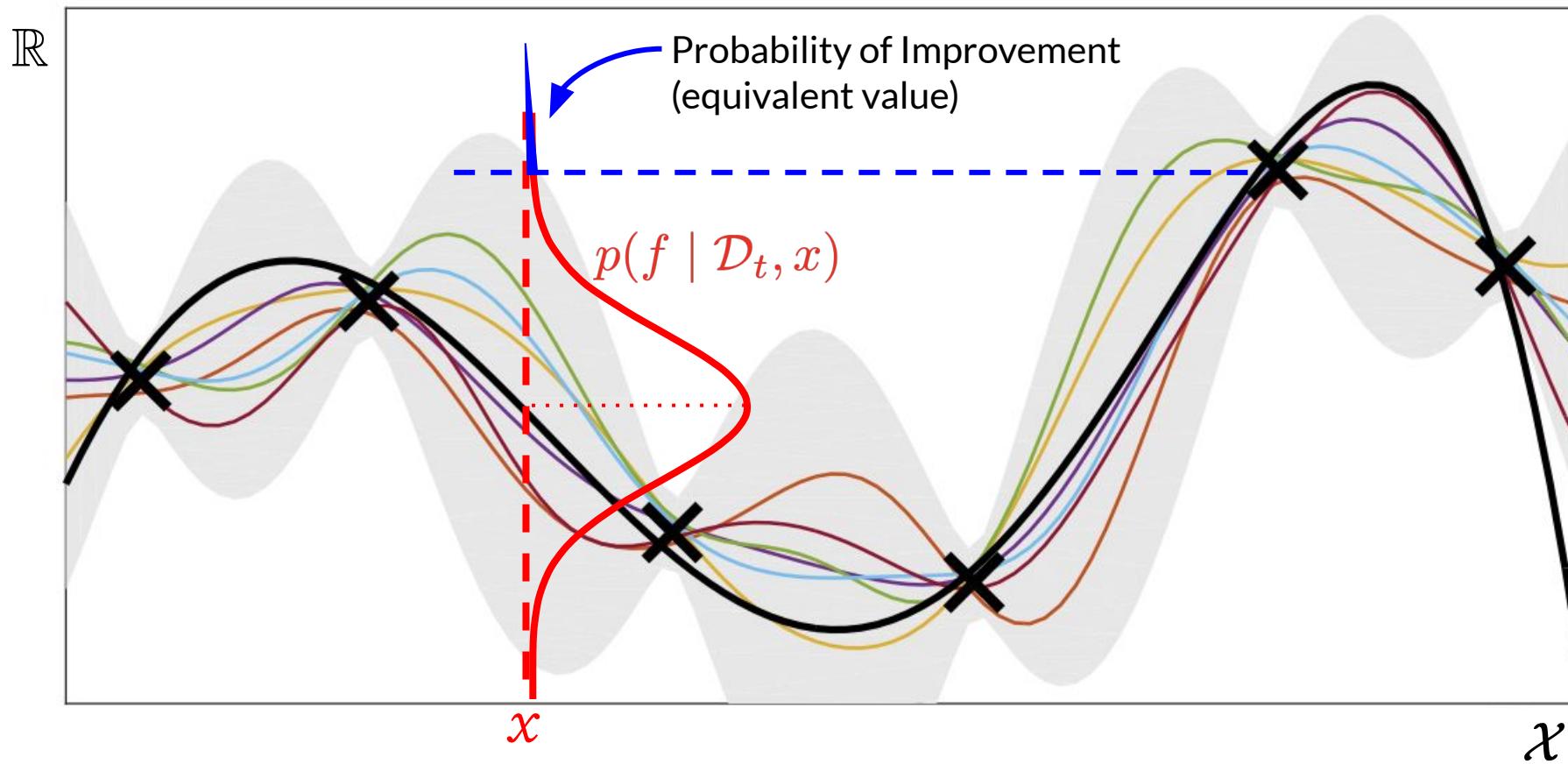
# Acquisition Functions – Probability of Improvement (PI)

Visualizing PI

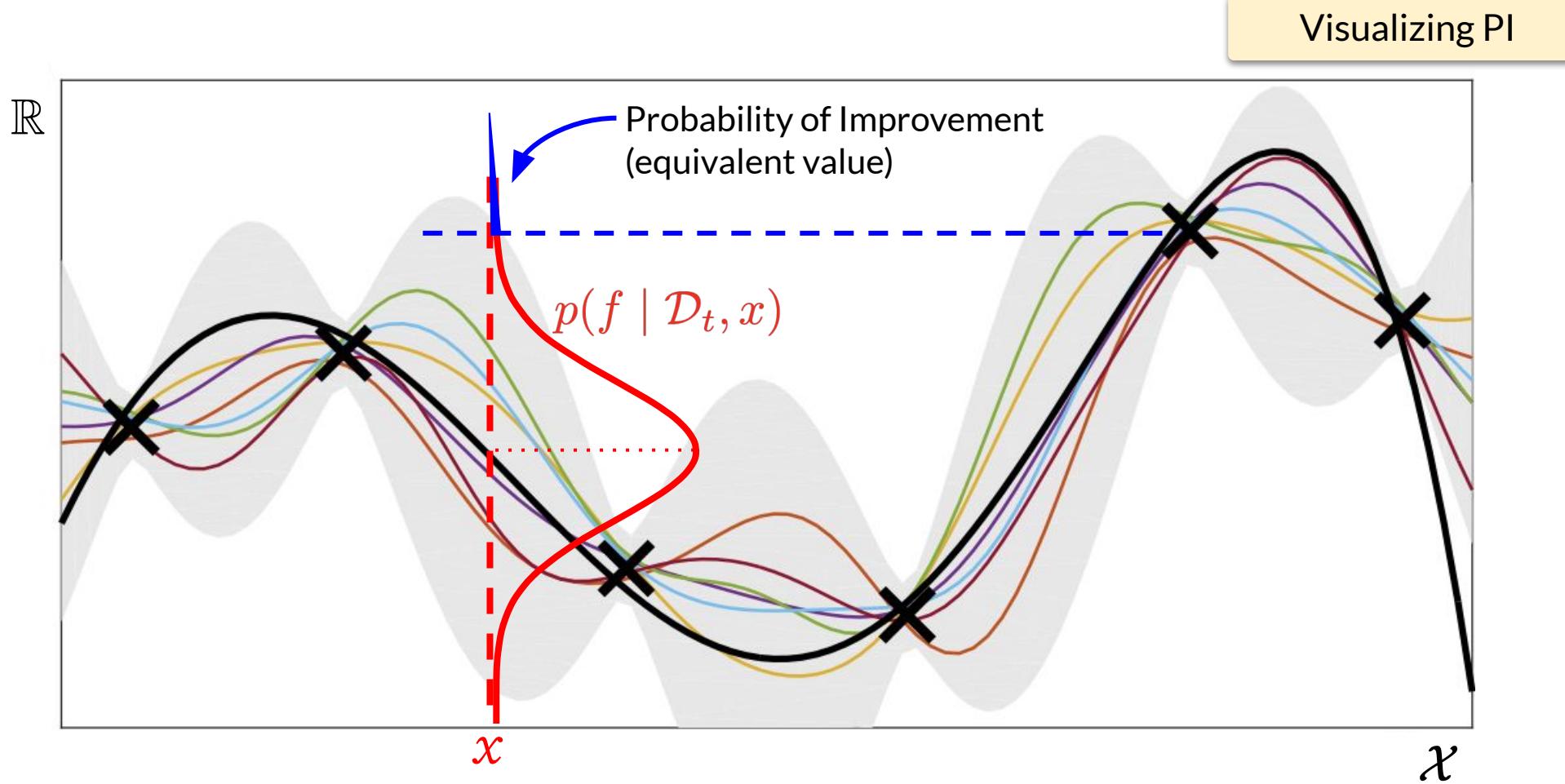


# Acquisition Functions – Probability of Improvement (PI)

Visualizing PI



# Acquisition Functions – Probability of Improvement (PI)



In contrast, we could try to maximize the *expected amount of improvement* over the best point so far...

# **Acquisition Functions – Expected Improvement (EI)**

## Acquisition Functions – Expected Improvement (EI)

This leads us to one of the most famous and widely used BO acquisition functions (except possibly for UCB).

## Acquisition Functions – Expected Improvement (EI)

This leads us to one of the most famous and widely used BO acquisition functions (except possibly for UCB).

The **expected improvement (EI)** acquisition function is written:

## Acquisition Functions – Expected Improvement (EI)

This leads us to one of the most famous and widely used BO acquisition functions (except possibly for UCB).

The **expected improvement (EI)** acquisition function is written:

$$\alpha_t(x) = \mathbb{E}_{p(f|\mathcal{D}_t)} [\max\{0, f(x) - f(x^*)\}]$$

## Acquisition Functions – Expected Improvement (EI)

This leads us to one of the most famous and widely used BO acquisition functions (except possibly for UCB).

The **expected improvement (EI)** acquisition function is written:

$$\alpha_t(x) = \mathbb{E}_{p(f|\mathcal{D}_t)} [\max\{0, f(x) - f(x^*)\}]$$

where  $f(x^*)$  is the function value of the best point observed so far in  $\mathcal{D}_t$ .

## Acquisition Functions – Expected Improvement (EI)

This leads us to one of the most famous and widely used BO acquisition functions (except possibly for UCB).

The **expected improvement (EI)** acquisition function is written:

$$\alpha_t(x) = \mathbb{E}_{p(f|\mathcal{D}_t)} [\max\{0, f(x) - f(x^*)\}]$$

where  $f(x^*)$  is the function value of the best point observed so far in  $\mathcal{D}_t$ .

*I.e., the expected amount that  $f(x)$  is greater than the best point observed so far (while ignoring the magnitude if  $f(x)$  is less than the best point so far).*

## Acquisition Functions – Expected Improvement (EI)

This leads us to one of the most famous and widely used BO acquisition functions (except possibly for UCB).

The **expected improvement (EI)** acquisition function is written:

$$\alpha_t(x) = \mathbb{E}_{p(f|\mathcal{D}_t)} [\max\{0, f(x) - f(x^*)\}]$$

where  $f(x^*)$  is the function value of the best point observed so far in  $\mathcal{D}_t$ .

*I.e., the expected amount that  $f(x)$  is greater than the best point observed so far (while ignoring the magnitude if  $f(x)$  is less than the best point so far).*

You can prove this is equal to ...

## Acquisition Functions – Expected Improvement (EI)

This leads us to one of the most famous and widely used BO acquisition functions (except possibly for UCB).

The **expected improvement (EI)** acquisition function is written:

$$\begin{aligned}\alpha_t(x) &= \mathbb{E}_{p(f|\mathcal{D}_t)} [\max\{0, f(x) - f(x^*)\}] \\ &= (\mu(x) - f(x^*))\Phi\left(\frac{\mu(x) - f(x^*)}{\sigma(x)}\right) + \sigma(x)\phi\left(\frac{\mu(x) - f(x^*)}{\sigma(x)}\right)\end{aligned}$$

## Acquisition Functions – Expected Improvement (EI)

This leads us to one of the most famous and widely used BO acquisition functions (except possibly for UCB).

The **expected improvement (EI)** acquisition function is written:

$$\begin{aligned}\alpha_t(x) &= \mathbb{E}_{p(f|\mathcal{D}_t)} [\max\{0, f(x) - f(x^*)\}] \\ &= (\mu(x) - f(x^*))\Phi\left(\frac{\mu(x) - f(x^*)}{\sigma(x)}\right) + \sigma(x)\phi\left(\frac{\mu(x) - f(x^*)}{\sigma(x)}\right)\end{aligned}$$

where  $\Phi(\cdot)$  denotes the CDF of the standard normal distribution and  $\phi(\cdot)$  denotes the PDF of the standard normal distribution.

(And, as before,  $\mu(x)$  and  $\sigma(x)$  are the mean and SD of marginal posterior at  $x$ .)

# Acquisition Functions – Expected Improvement (EI)

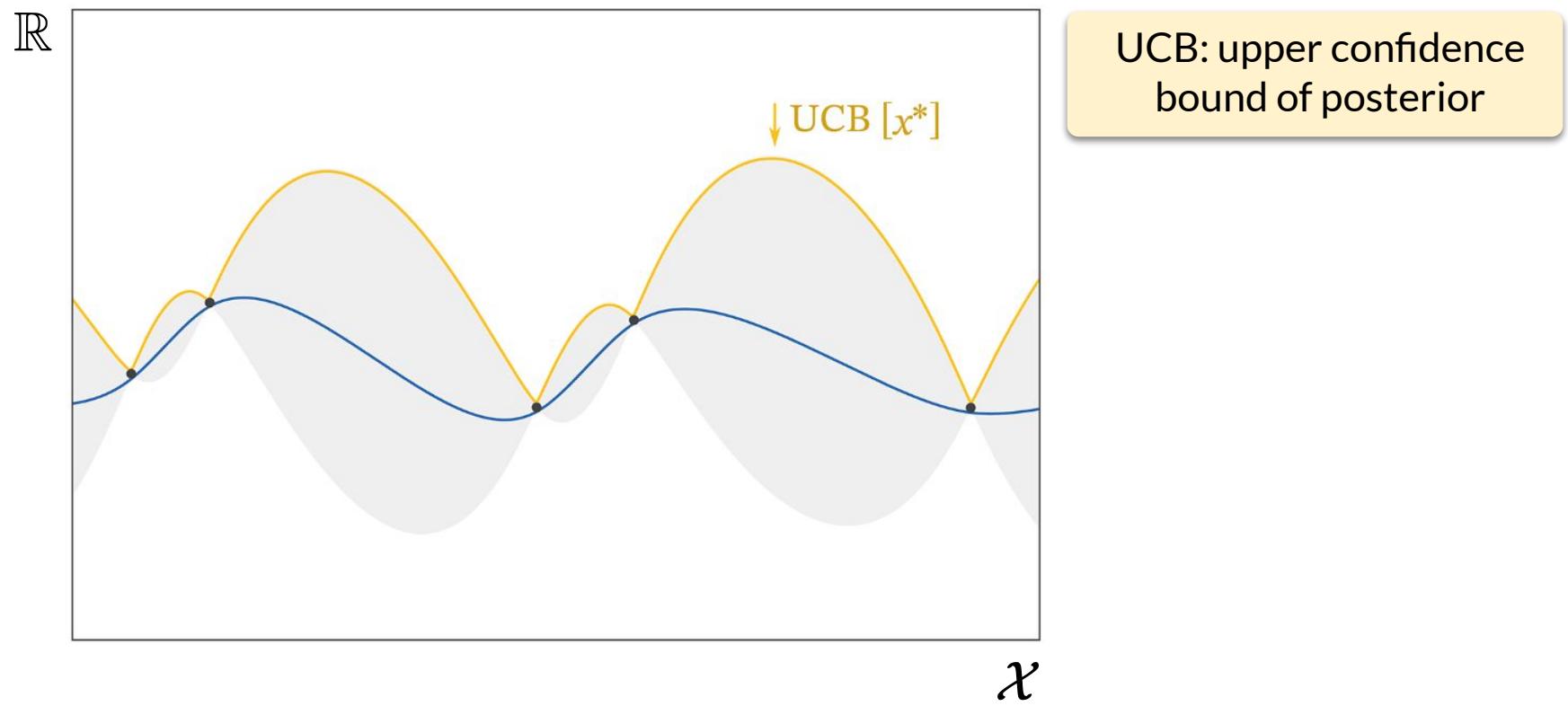
Nice visual comparing UCB vs. PI vs. EI:

Source: M. O. Ahmed, S. Prince, "Bayesian optimization"

# Acquisition Functions – Expected Improvement (EI)

Nice visual comparing UCB vs. PI vs. EI:

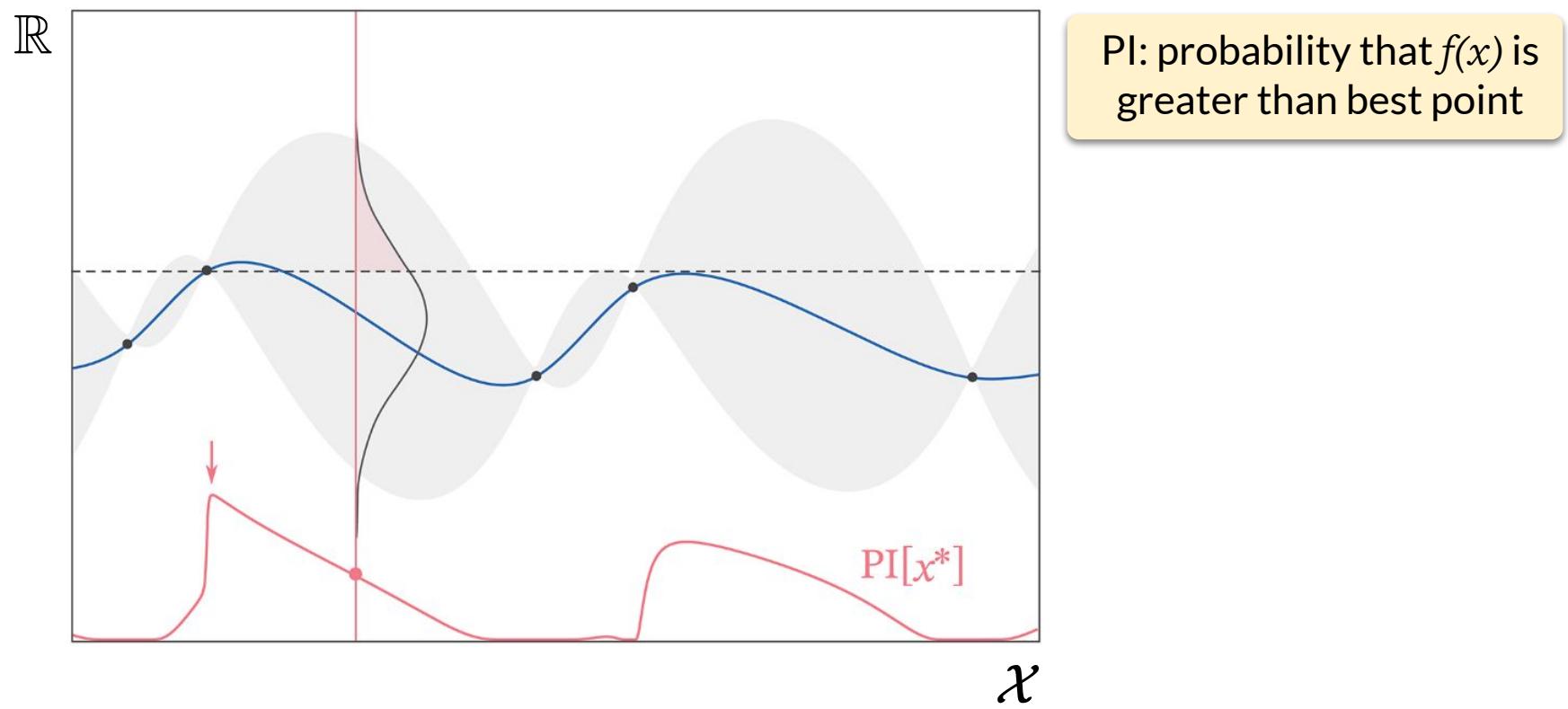
Source: M. O. Ahmed, S. Prince, "Bayesian optimization"



# Acquisition Functions – Expected Improvement (EI)

Nice visual comparing UCB vs. PI vs. EI:

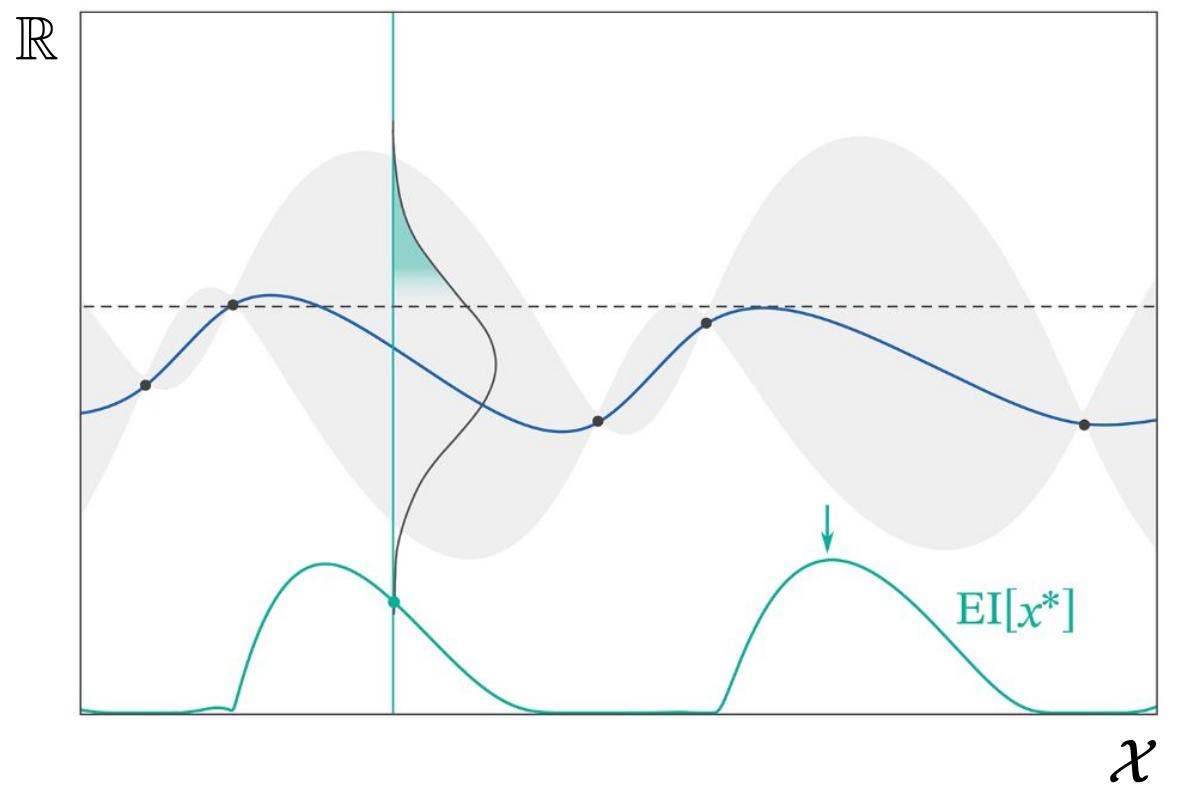
Source: M. O. Ahmed, S. Prince, "Bayesian optimization"



# Acquisition Functions – Expected Improvement (EI)

Nice visual comparing UCB vs. PI vs. EI:

Source: M. O. Ahmed, S. Prince, "Bayesian optimization"



EI: *expected amount  $f(x)$  is greater than best point (and ignoring amount that it's less than best point)*

# Acquisition Functions – Entropy Search (ES)

## Acquisition Functions – Entropy Search (ES)

Another strategy is the “experimental design” strategy, but *applied to BO setting*.

## Acquisition Functions – Entropy Search (ES)

Another strategy is the “experimental design” strategy, but *applied to BO setting*.

Specifically, this means:

Reducing posterior uncertainty (entropy) over quantity of interest – in this case, the location  $x^*$  of the function maximizer – rather than, e.g., the full landscape of  $f$ .

## Acquisition Functions – Entropy Search (ES)

Another strategy is the “experimental design” strategy, but *applied to BO setting*.

Specifically, this means:

Reducing posterior uncertainty (entropy) over quantity of interest – in this case, the location  $x^*$  of the function maximizer – rather than, e.g., the full landscape of  $f$ .

As mentioned, this is an *info-based* BO acquisition function called **entropy search**.

## Acquisition Functions – Entropy Search (ES)

Another strategy is the “experimental design” strategy, but *applied to BO setting*.

Specifically, this means:

Reducing posterior uncertainty (entropy) over quantity of interest – in this case, the location  $x^*$  of the function maximizer – rather than, e.g., the full landscape of  $f$ .

As mentioned, this is an *info-based* BO acquisition function called **entropy search**.

Why would you want to use this acquisition function?

⇒ Intuitively: if you get rewarded for *how good your guess of  $x^*$  is* – rather than, the value of the maximizer,  $f(x^*)$  – ES should be more optimal.

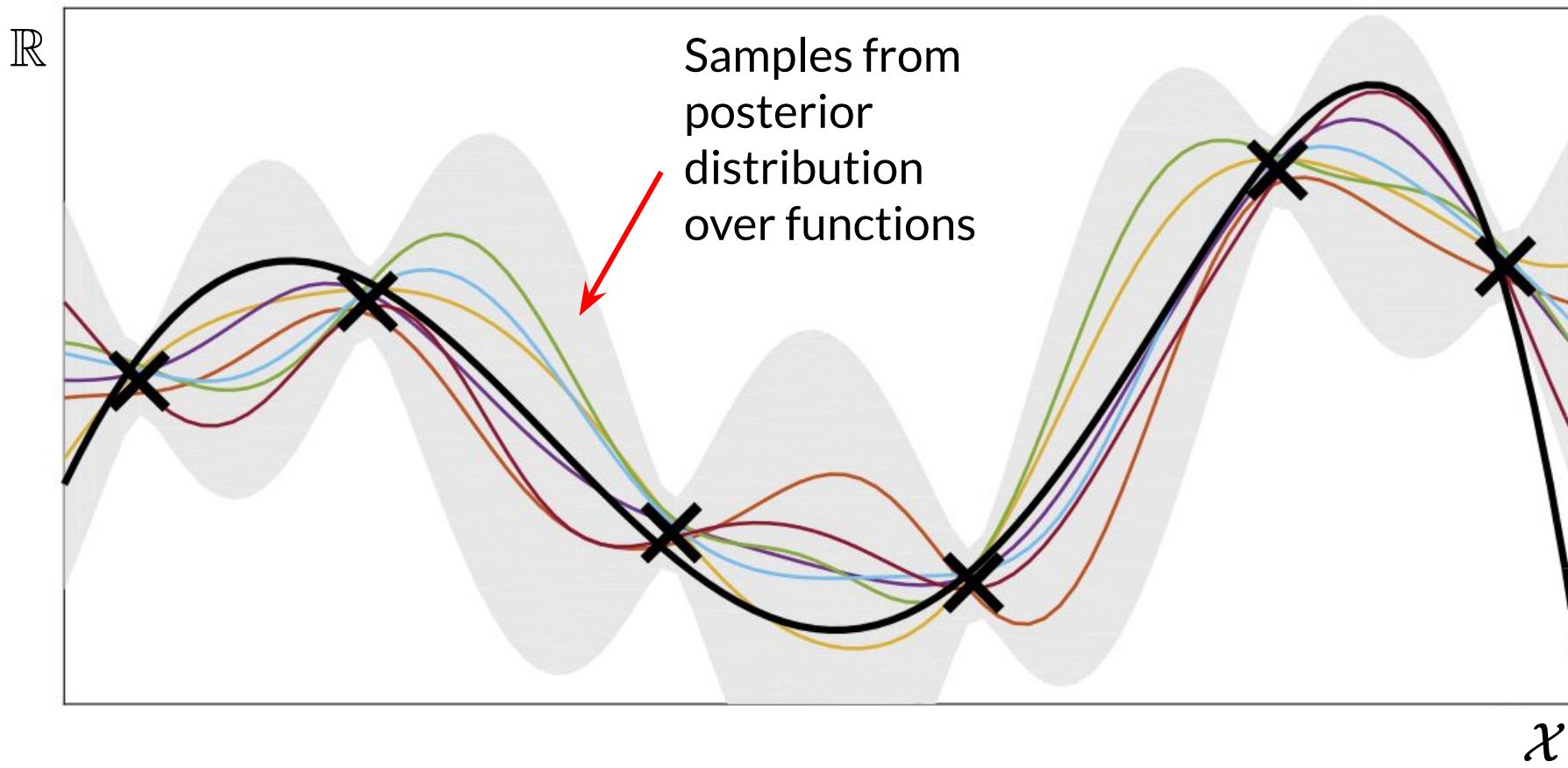
# Acquisition Functions – Entropy Search

Visualizing ES

To explain this acquisition function, let's visualize...

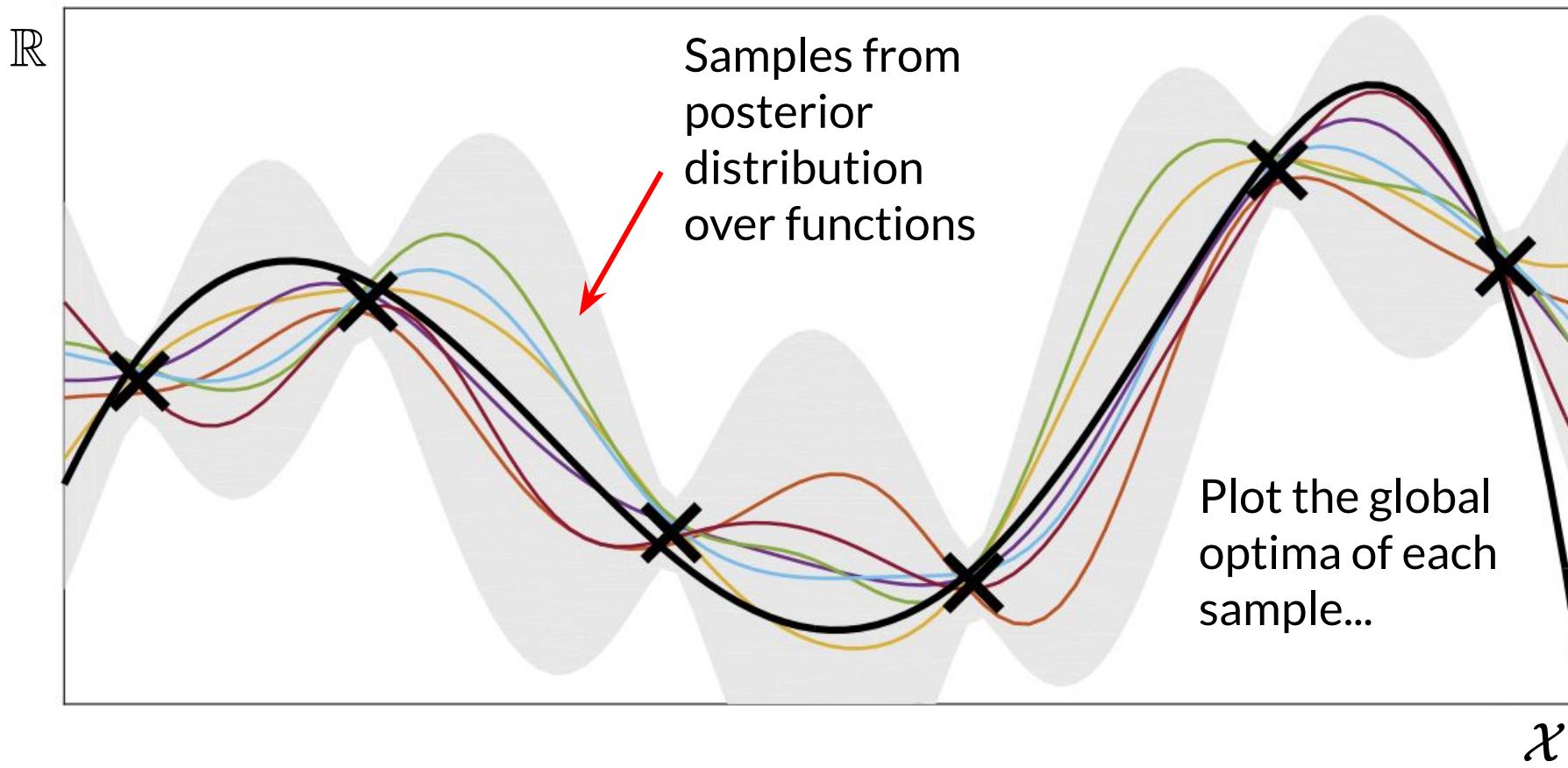
# Acquisition Functions – Entropy Search

Visualizing ES



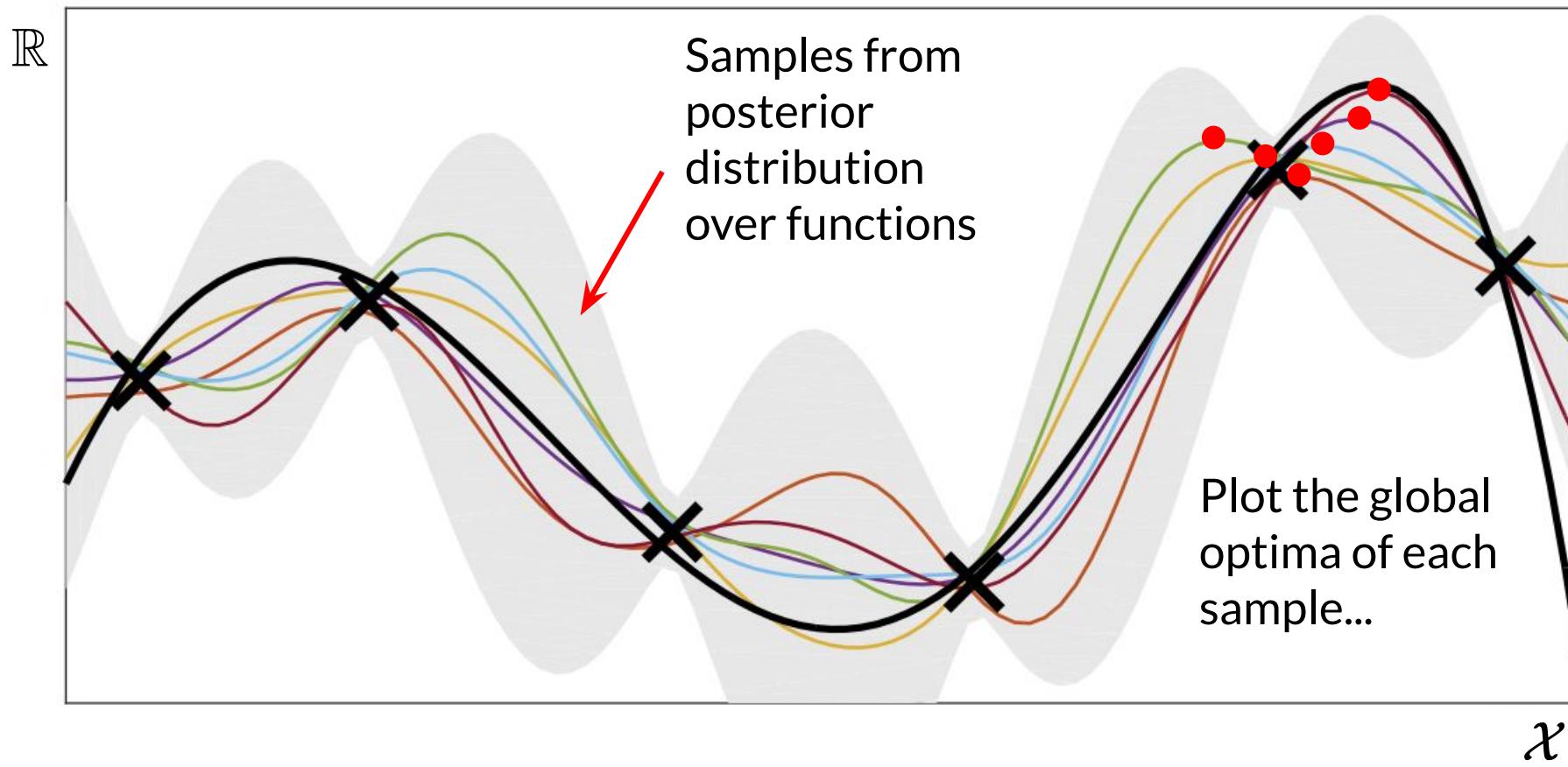
# Acquisition Functions – Entropy Search

Visualizing ES



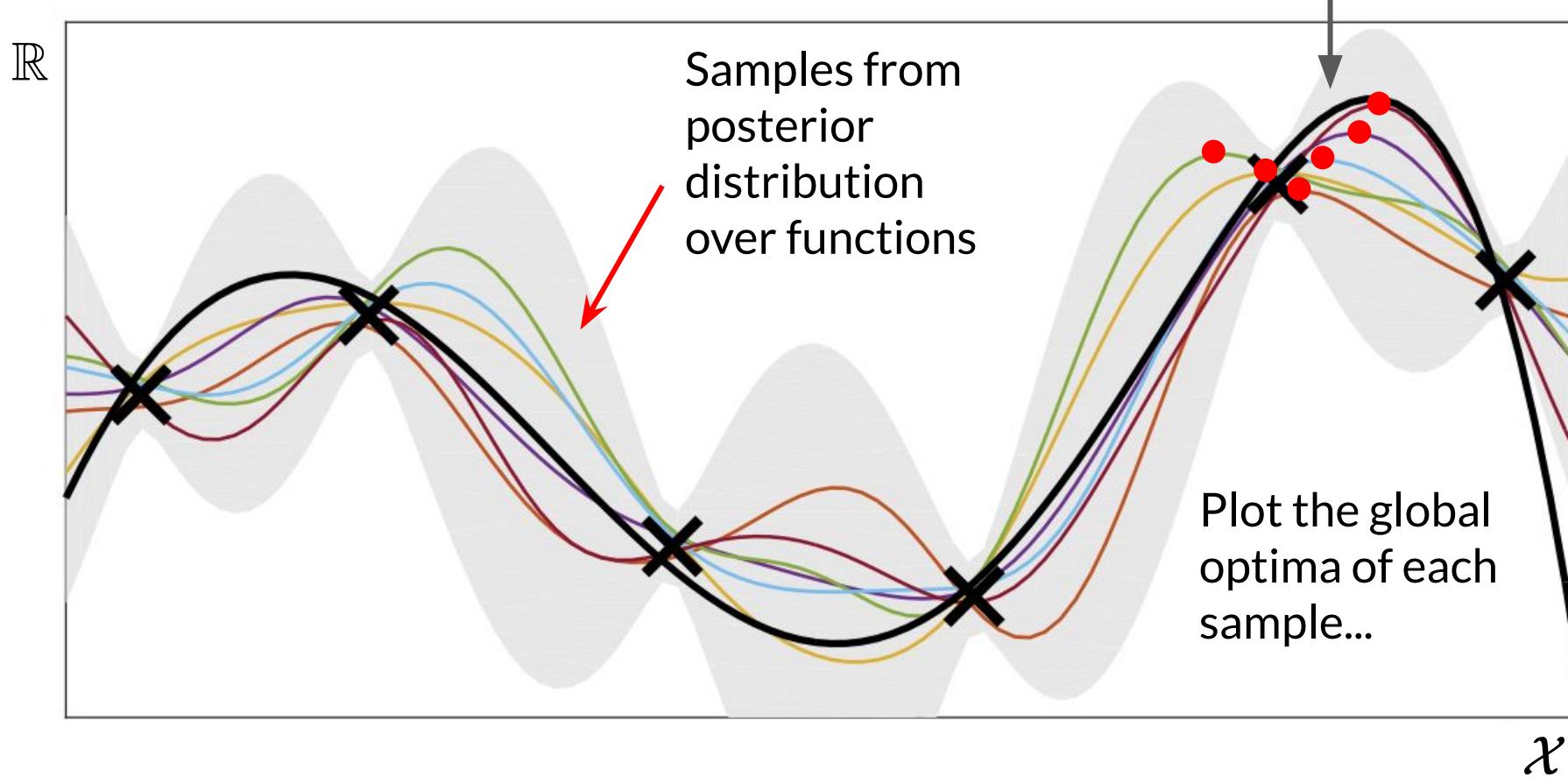
# Acquisition Functions – Entropy Search

Visualizing ES



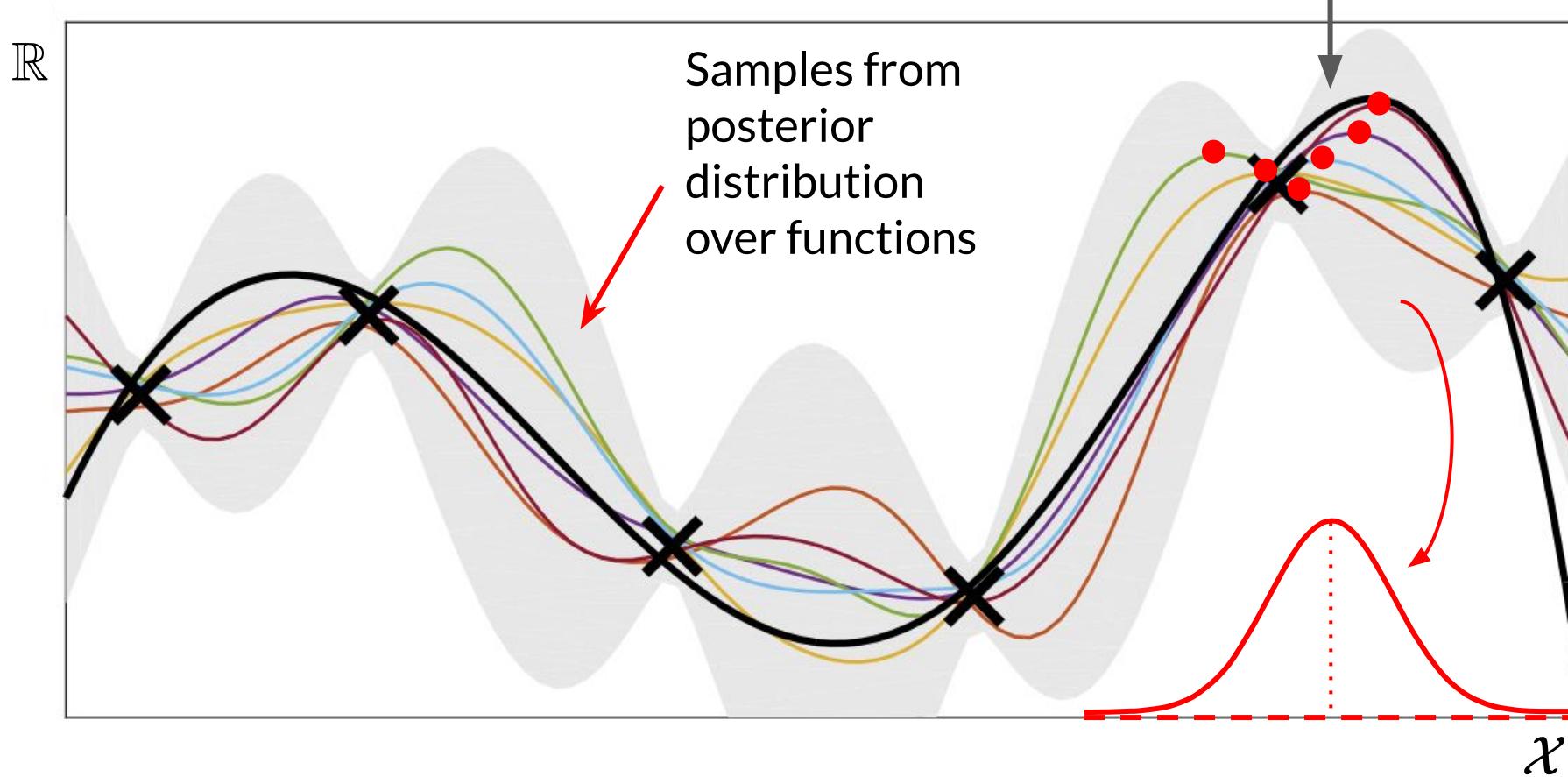
# Acquisition Functions – Entropy Search

There is a posterior distribution over global optima induced by probabilistic model:  $p(x^* | \mathcal{D}_t)$



# Acquisition Functions – Entropy Search

There is a posterior distribution over global optima induced by probabilistic model:  $p(x^* | \mathcal{D}_t)$



# Acquisition Functions – Entropy Search

This leads us to the acquisition function:

*e.g. used in entropy search (ES), predictive entropy search (PES)*

# Acquisition Functions – Entropy Search

This leads us to the acquisition function:

e.g. used in entropy search (ES), predictive entropy search (PES)

$$\alpha_t(x) = H[p(x^* \mid \mathcal{D}_t)] - \mathbb{E}_{p(y_x \mid \mathcal{D}_t)}[H[p(x^* \mid \mathcal{D}_t \cup \{(x, y_x)\})]]$$

entropy of posterior distribution over  
global optima  $x^*$

minus

expected entropy of posterior distribution  
over global optima  $x^*$ ,  
... if we were to make a query at  $x$

# Acquisition Functions – Entropy Search

This leads us to the acquisition function:

e.g. used in entropy search (ES), predictive entropy search (PES)

$$\alpha_t(x) = H[p(x^* | \mathcal{D}_t)] - \mathbb{E}_{p(y_x | \mathcal{D}_t)}[H[p(x^* | \mathcal{D}_t \cup \{(x, y_x)\})]]$$

entropy of posterior distribution over  
global optima  $x^*$

minus

expected entropy of posterior distribution  
over global optima  $x^*$ ,  
... if we were to make a query at  $x$

# Acquisition Functions – Entropy Search

This leads us to the acquisition function:

e.g. used in entropy search (ES), predictive entropy search (PES)

$$\alpha_t(x) = H[p(x^* \mid \mathcal{D}_t)] - \mathbb{E}_{p(y_x \mid \mathcal{D}_t)}[H[p(x^* \mid \mathcal{D}_t \cup \{(x, y_x)\})]]$$

entropy of posterior distribution over global optima  $x^*$

minus

expected entropy of posterior distribution over global optima  $x^*$ ,  
... if we were to make a query at  $x$

```
graph LR; A["\alpha_t(x) = H[p(x^* | \mathcal{D}_t)] - \mathbb{E}_{p(y_x | \mathcal{D}_t)}[H[p(x^* | \mathcal{D}_t \cup \{(x, y_x)\})]]"] --> B["entropy of posterior distribution over global optima x*"]; A --> C["minus"]; A --> D["expected entropy of posterior distribution over global optima x*, ... if we were to make a query at x"]
```

# Acquisition Functions – Entropy Search

This leads us to the acquisition function:

e.g. used in entropy search (ES), predictive entropy search (PES)

$$\alpha_t(x) = H[p(x^* \mid \mathcal{D}_t)] - \mathbb{E}_{p(y_x \mid \mathcal{D}_t)}[H[p(x^* \mid \mathcal{D}_t \cup \{(x, y_x)\})]]$$

entropy of posterior distribution over global optima  $x^*$

minus

expected entropy of posterior distribution over global optima  $x^*$ ,  
... if we were to make a query at  $x$

# Acquisition Functions – Entropy Search

This leads us to the acquisition function:

e.g. used in entropy search (ES), predictive entropy search (PES)

$$\alpha_t(x) = H[p(x^* \mid \mathcal{D}_t)] - \mathbb{E}_{p(y_x \mid \mathcal{D}_t)}[H[p(x^* \mid \mathcal{D}_t \cup \{(x, y_x)\})]]$$

entropy of posterior distribution over  
global optima  $x^*$

minus

expected entropy of posterior distribution  
over global optima  $x^*$ ,  
... if we were to make a query at  $x$

# Acquisition Functions – Entropy Search

This leads us to the acquisition function:

e.g. used in entropy search (ES), predictive entropy search (PES)

$$\alpha_t(x) = H[p(x^* \mid \mathcal{D}_t)] - \mathbb{E}_{p(y_x \mid \mathcal{D}_t)}[H[p(x^* \mid \mathcal{D}_t \cup \{(x, y_x)\})]]$$

entropy of posterior distribution over global optima  $x^*$

minus

expected entropy of posterior distribution over global optima  $x^*$ ,  
... if we were to make a query at  $x$

**Expected information gain (EIG)** – expected decrease in entropy if we were to query  $f$  at  $x$ .

# **Acquisition Functions – Additional**

## **Acquisition Functions – Additional**

A few acquisition functions we didn't get to today:

## Acquisition Functions – Additional

A few acquisition functions we didn't get to today:

- **Knowledge gradient** – somewhat complex, but one of the *most-optimal* (both theoretically and in practice) BO acquisition functions.

## Acquisition Functions – Additional

A few acquisition functions we didn't get to today:

- **Knowledge gradient** – somewhat complex, but one of the *most-optimal* (both theoretically and in practice) BO acquisition functions.
- **Thompson sampling** – a stochastic acquisition function (*i.e.*, draw one posterior, and use it within an acquisition function)

# Acquisition Functions – Additional

A few acquisition functions we didn't get to today:

- **Knowledge gradient** – somewhat complex, but one of the *most-optimal* (both theoretically and in practice) BO acquisition functions.
- **Thompson sampling** – a stochastic acquisition function (*i.e.*, draw one posterior, and use it within an acquisition function)
- **Non-BO acq fns.** – acquisition functions for non-optimization tasks, such as level sets, quadrature, phase boundaries, and more-general tasks.

# **Next Time**

## **Next Time**

Finish up our discussion of active learning methods!

