



UNIVERSITY OF
CAMBRIDGE

Lancaster
University



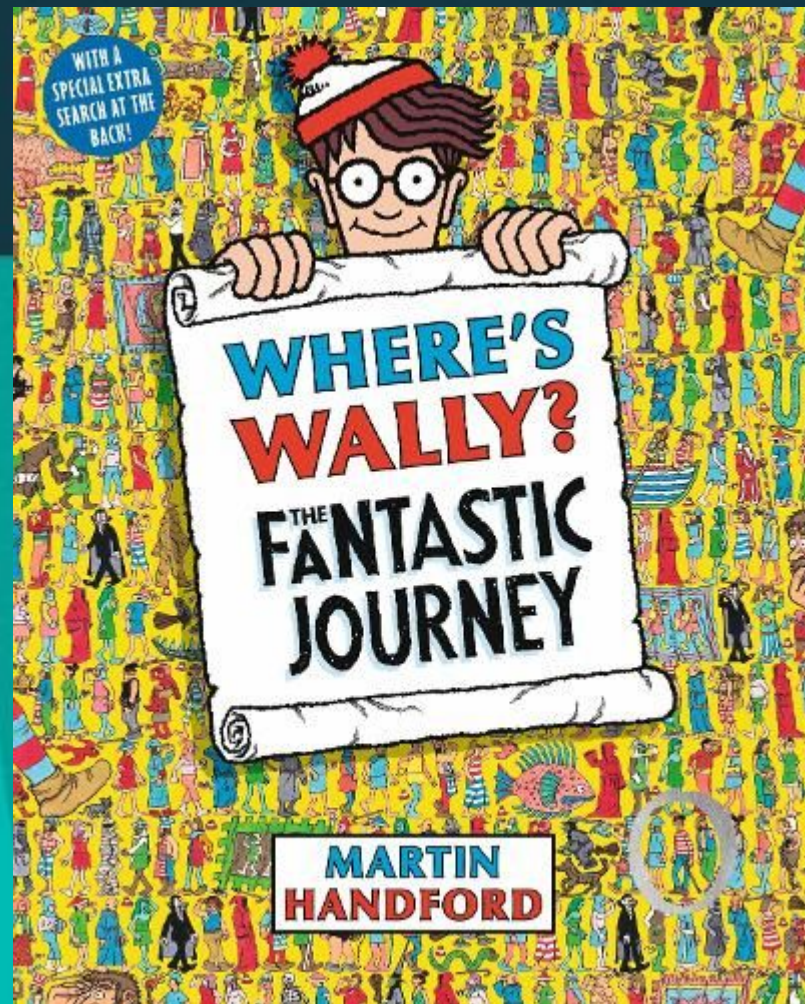
To Bayesian Optimisation and Beyond

Gaussian Processes as Decision Makers

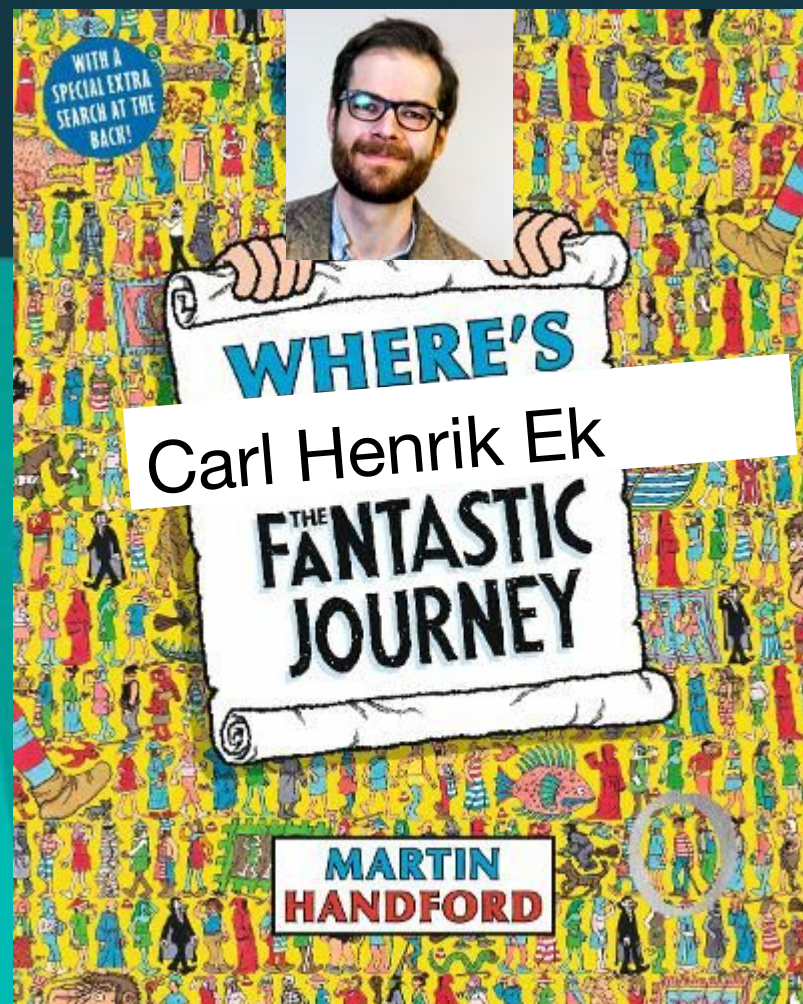
Henry Moss

Bayesian Search

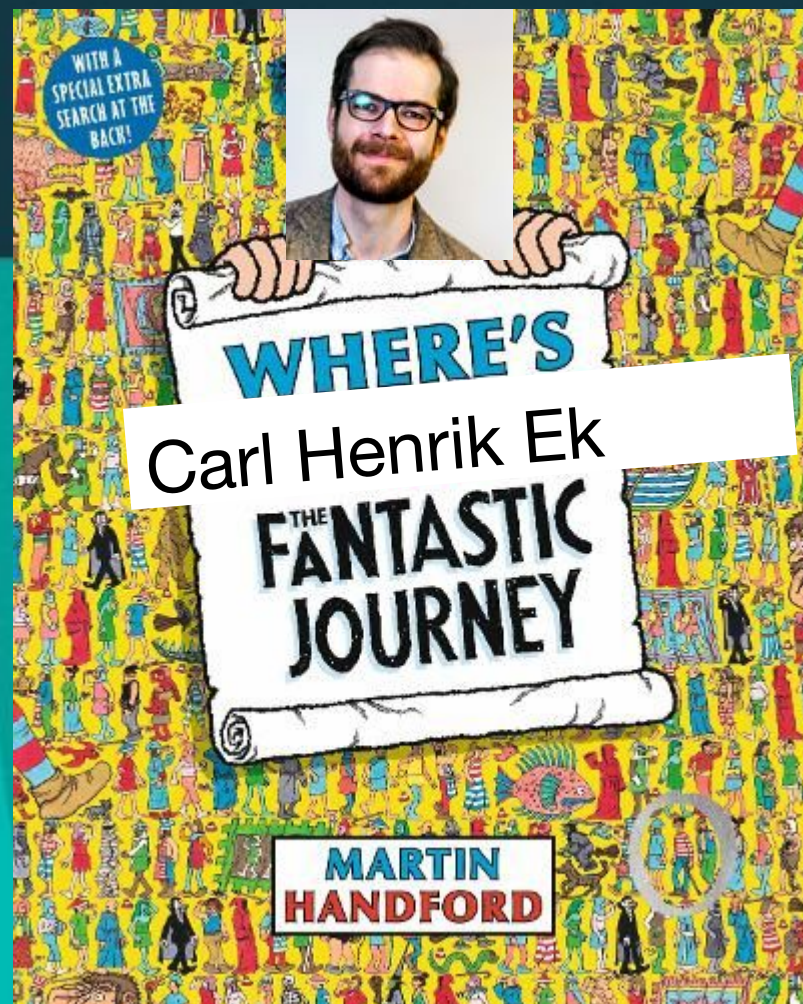
Bayesian Search



Bayesian Search

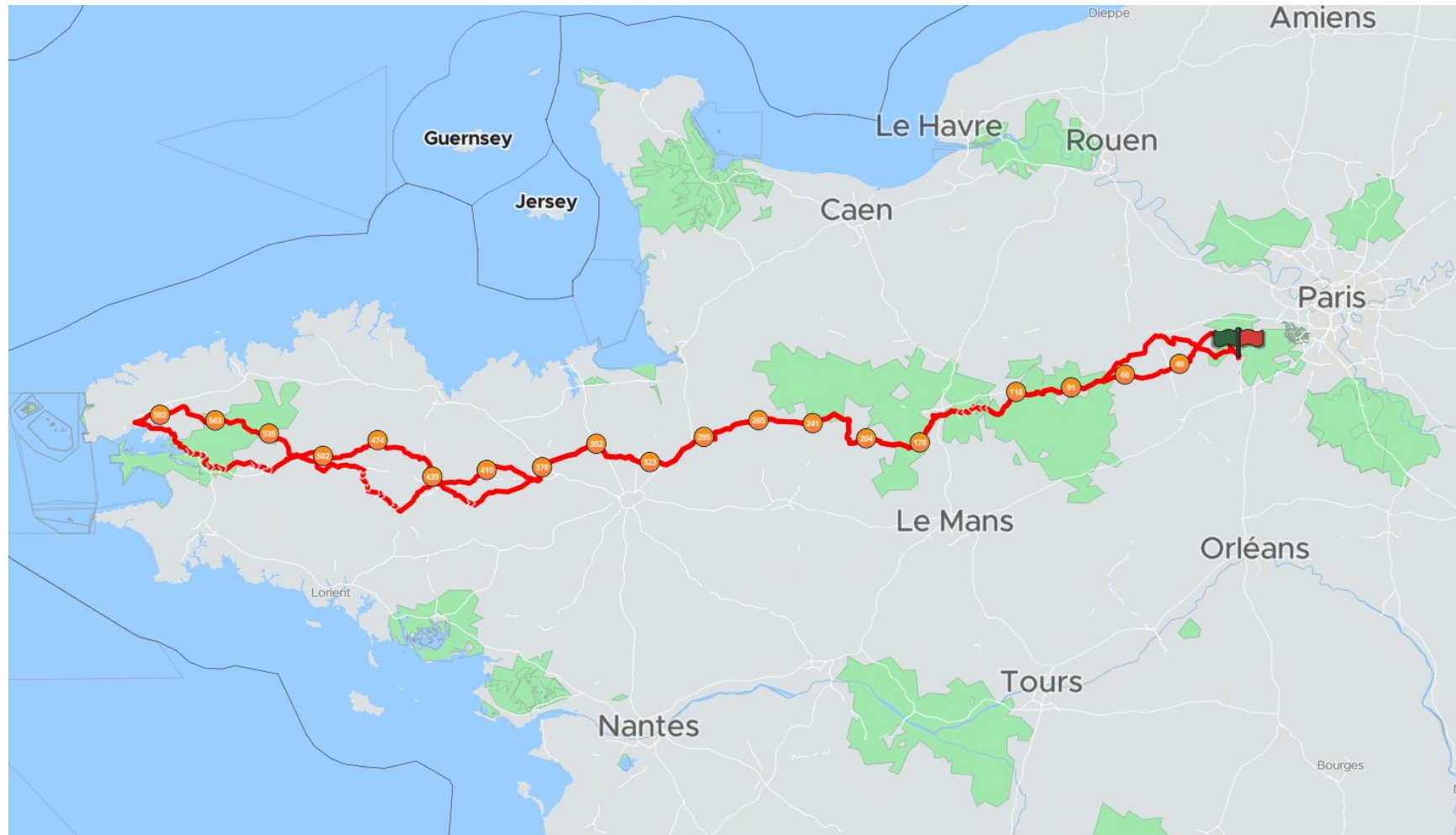


Bayesian Search

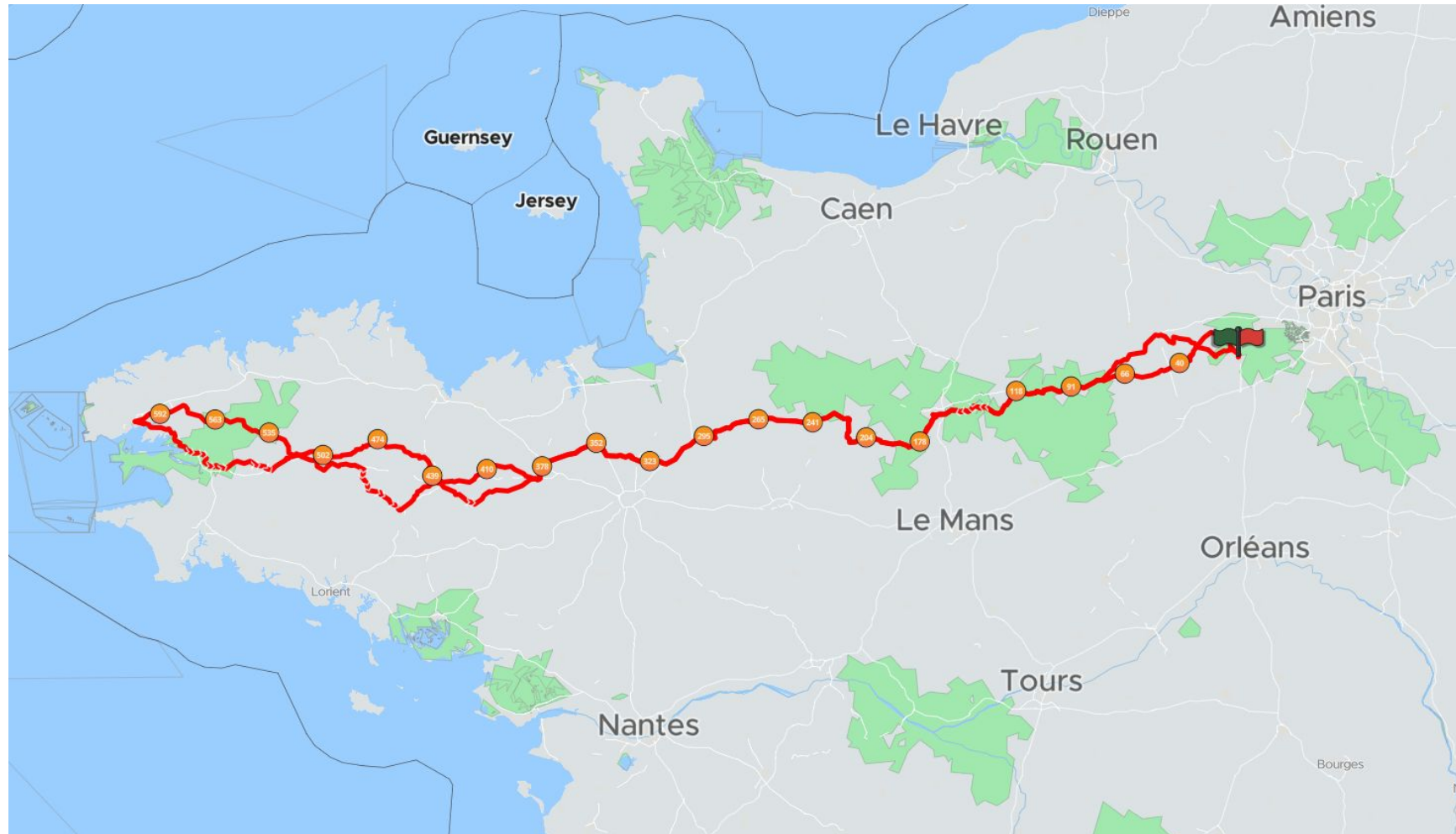




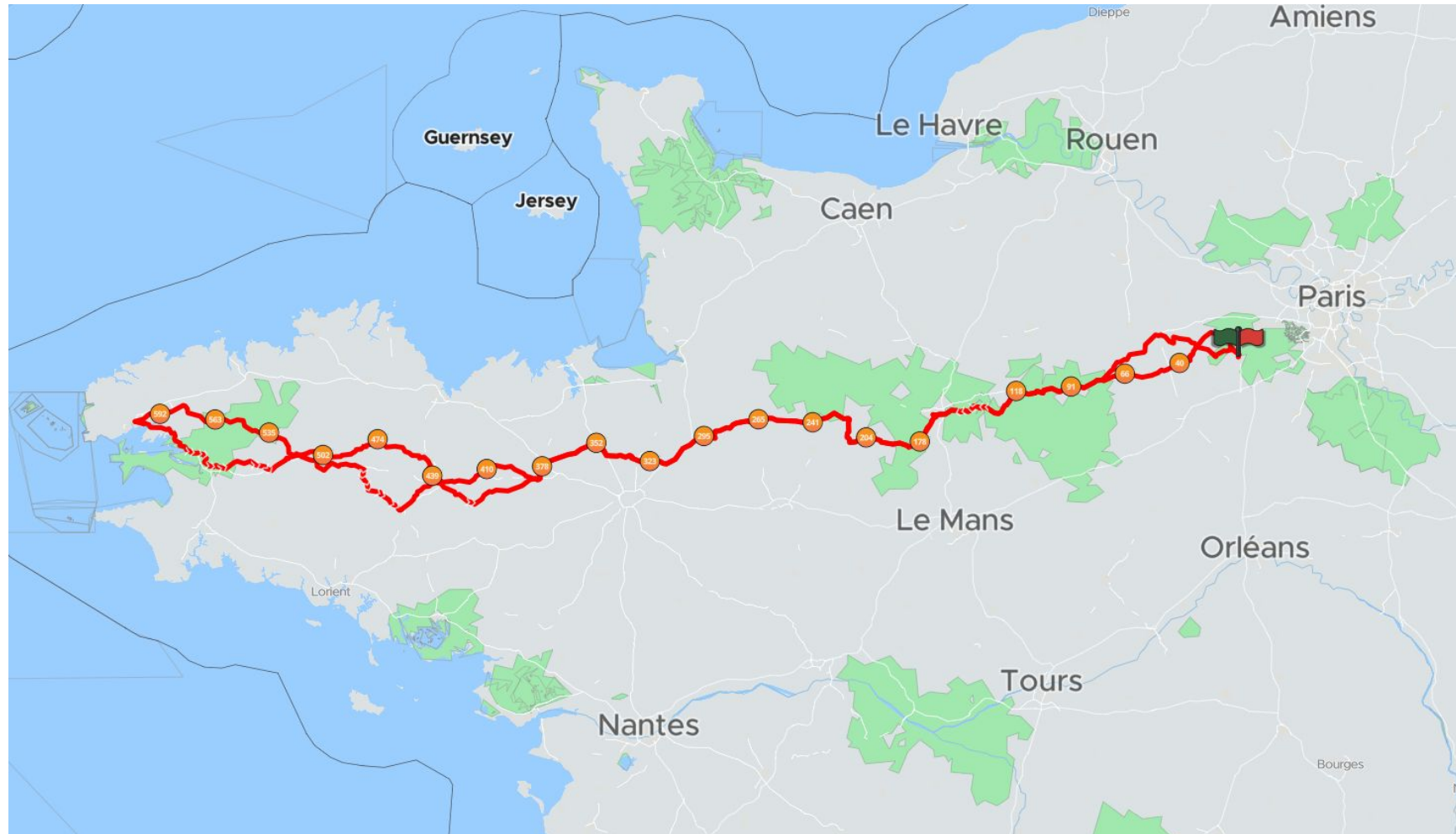
Where is Carl Henrik?



Where is Carl Henrik?



Where is Carl Henrik?





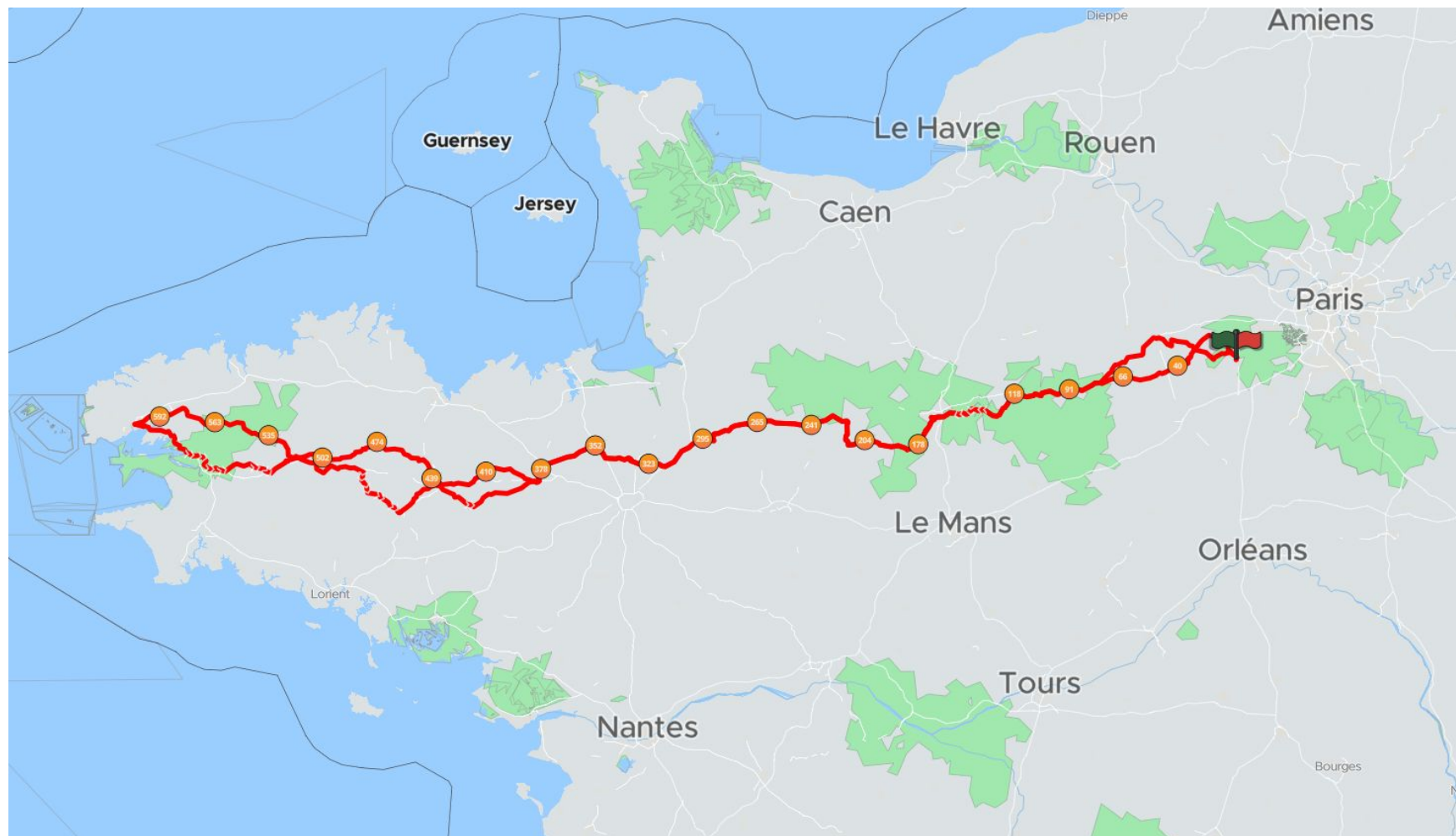
Where is Carl Henrik?





Where is Carl Henrik?

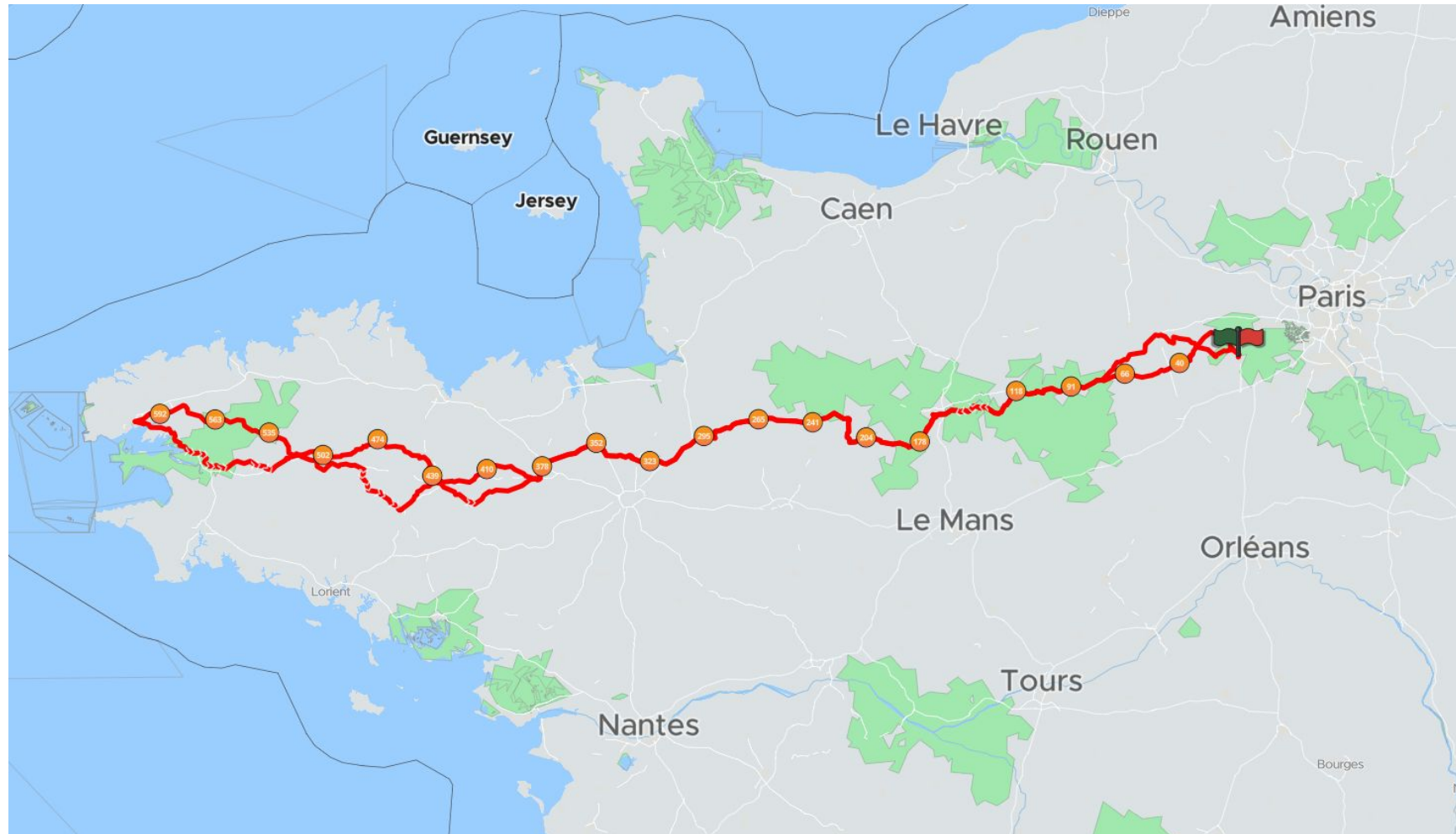
At 3:30 AM?



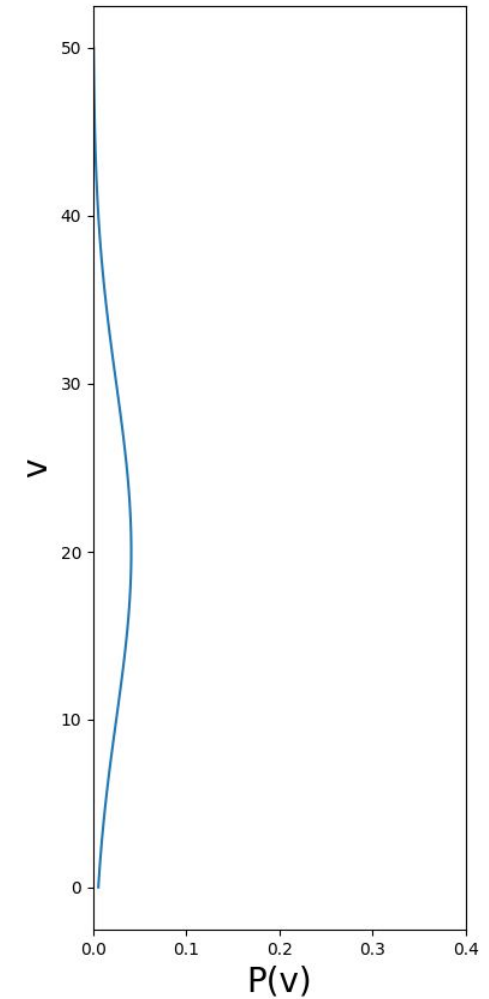
$$d = v \times t$$

Where is Carl Henrik?

At 3:30 AM?

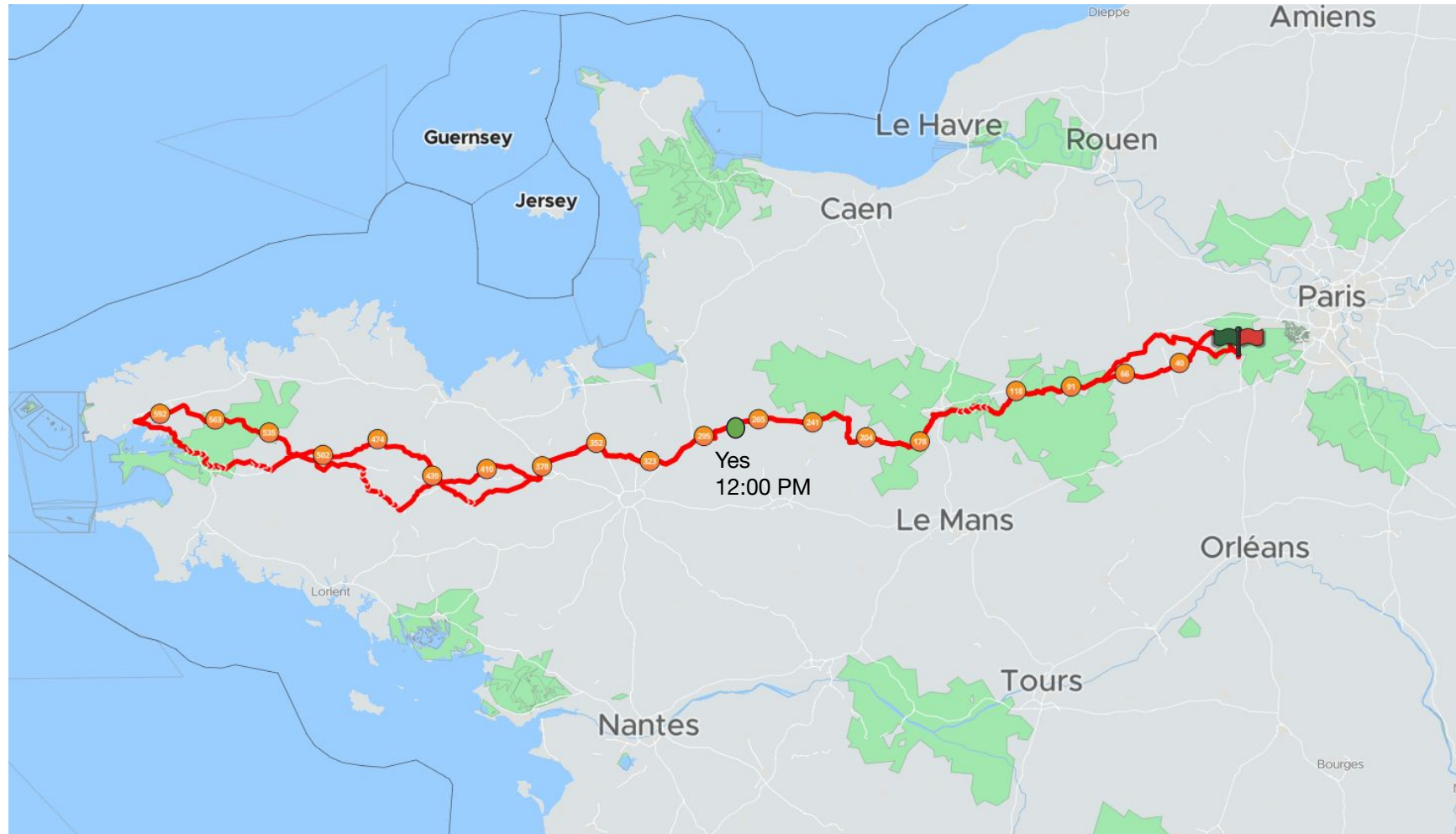


$$d = v \times t$$

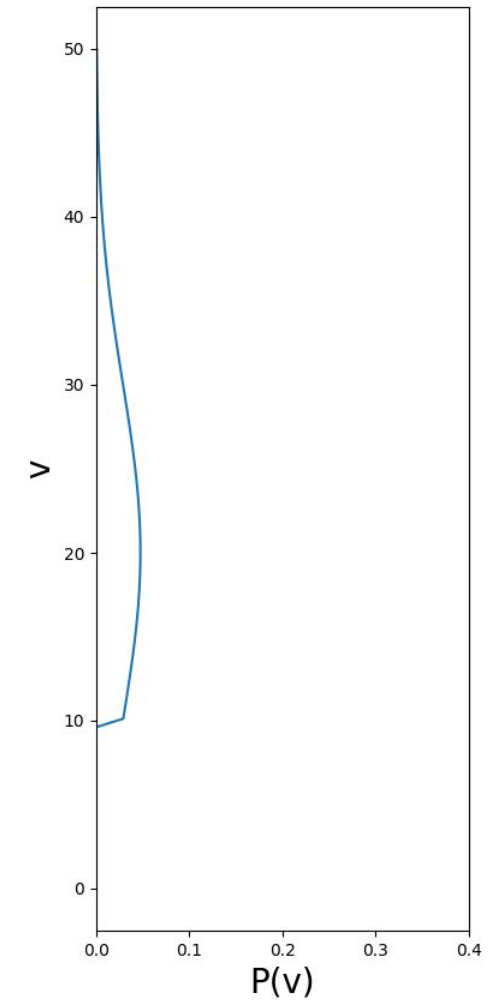


Where is Carl Henrik?

At 3:30 AM?

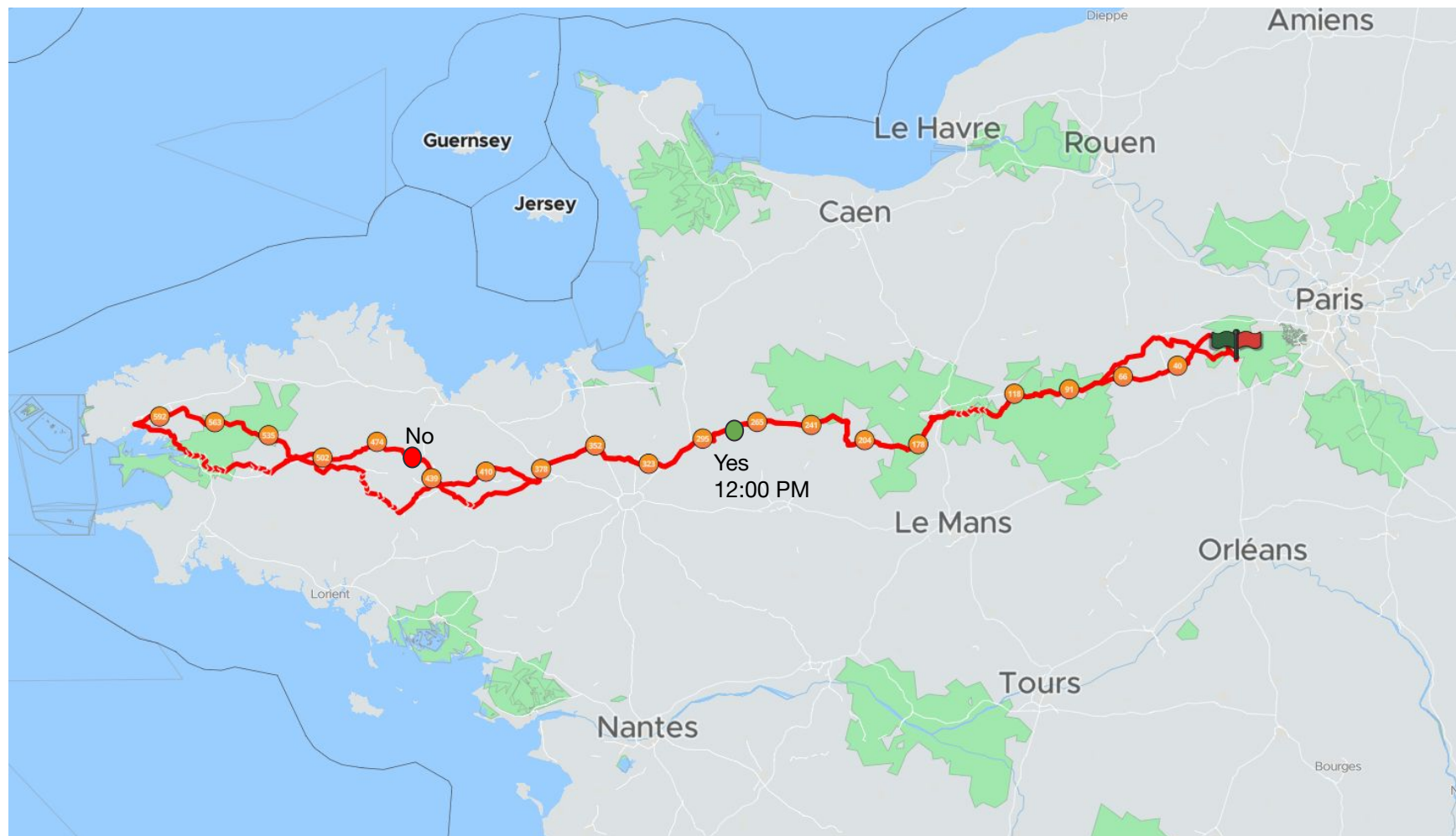


$$d = v \times t$$

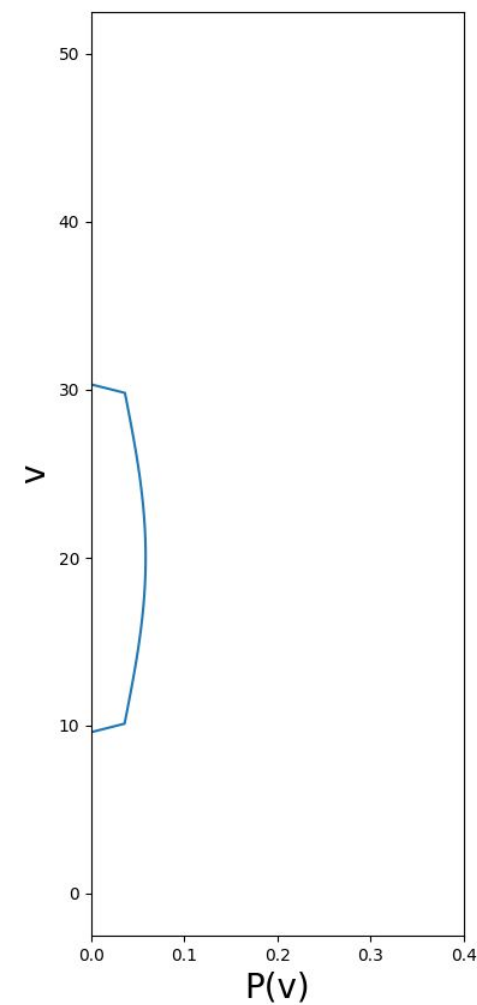


Where is Carl Henrik?

At 3:30 AM?

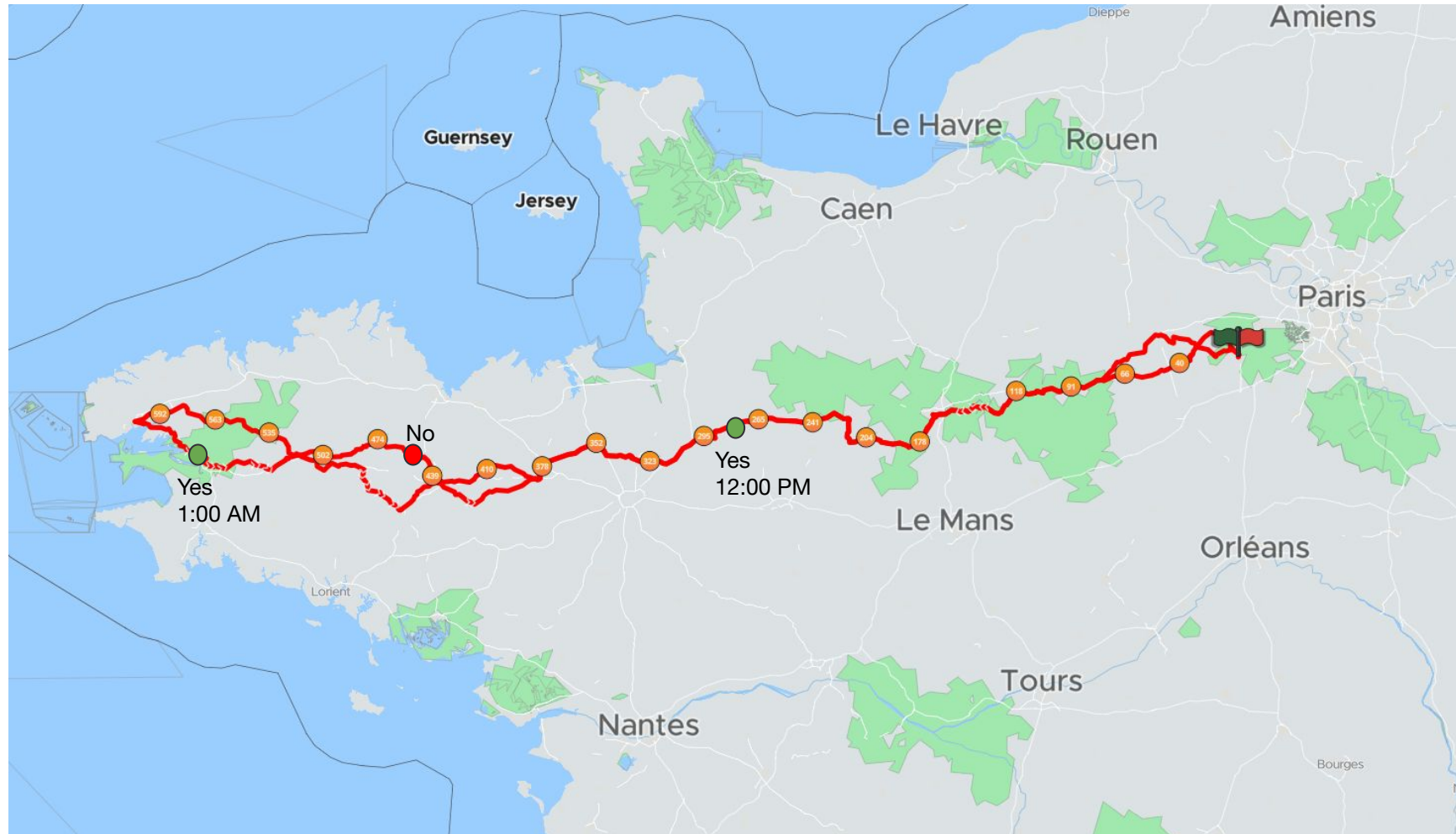


$$d = v \times t$$

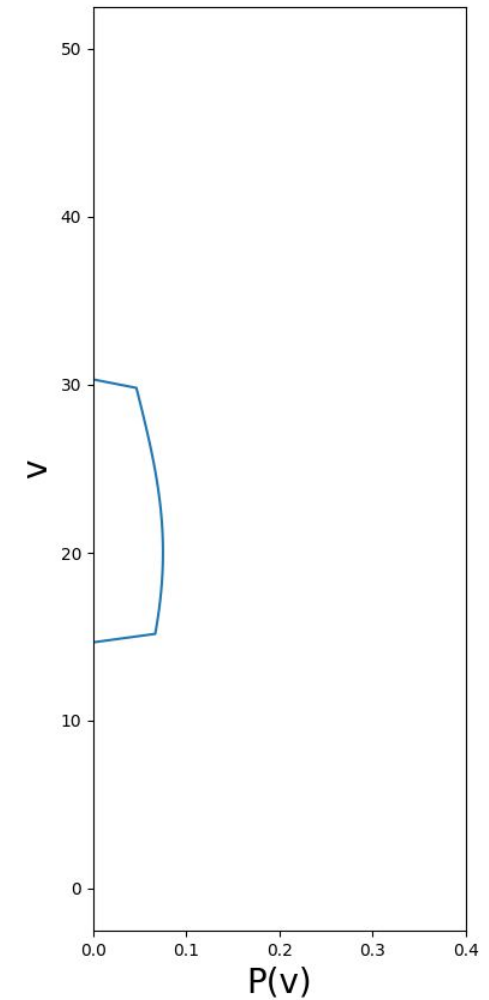


Where is Carl Henrik?

At 3:30 AM?

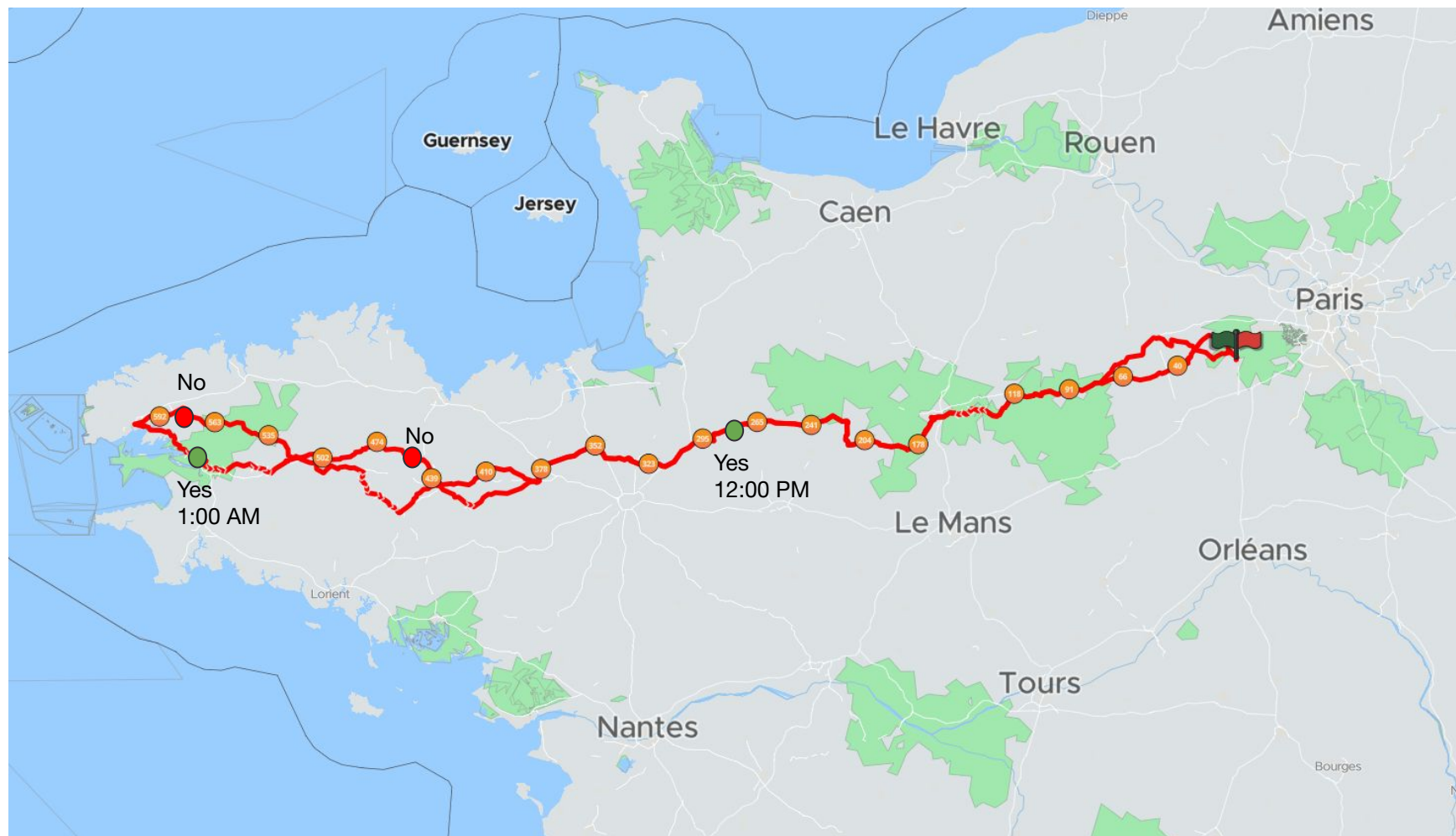


$$d = v \times t$$

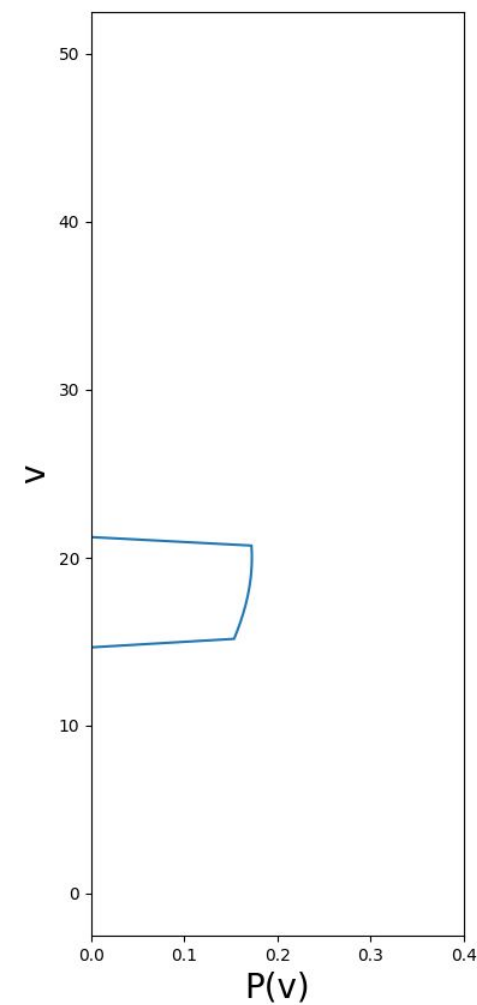


Where is Carl Henrik?

At 3:30 AM?

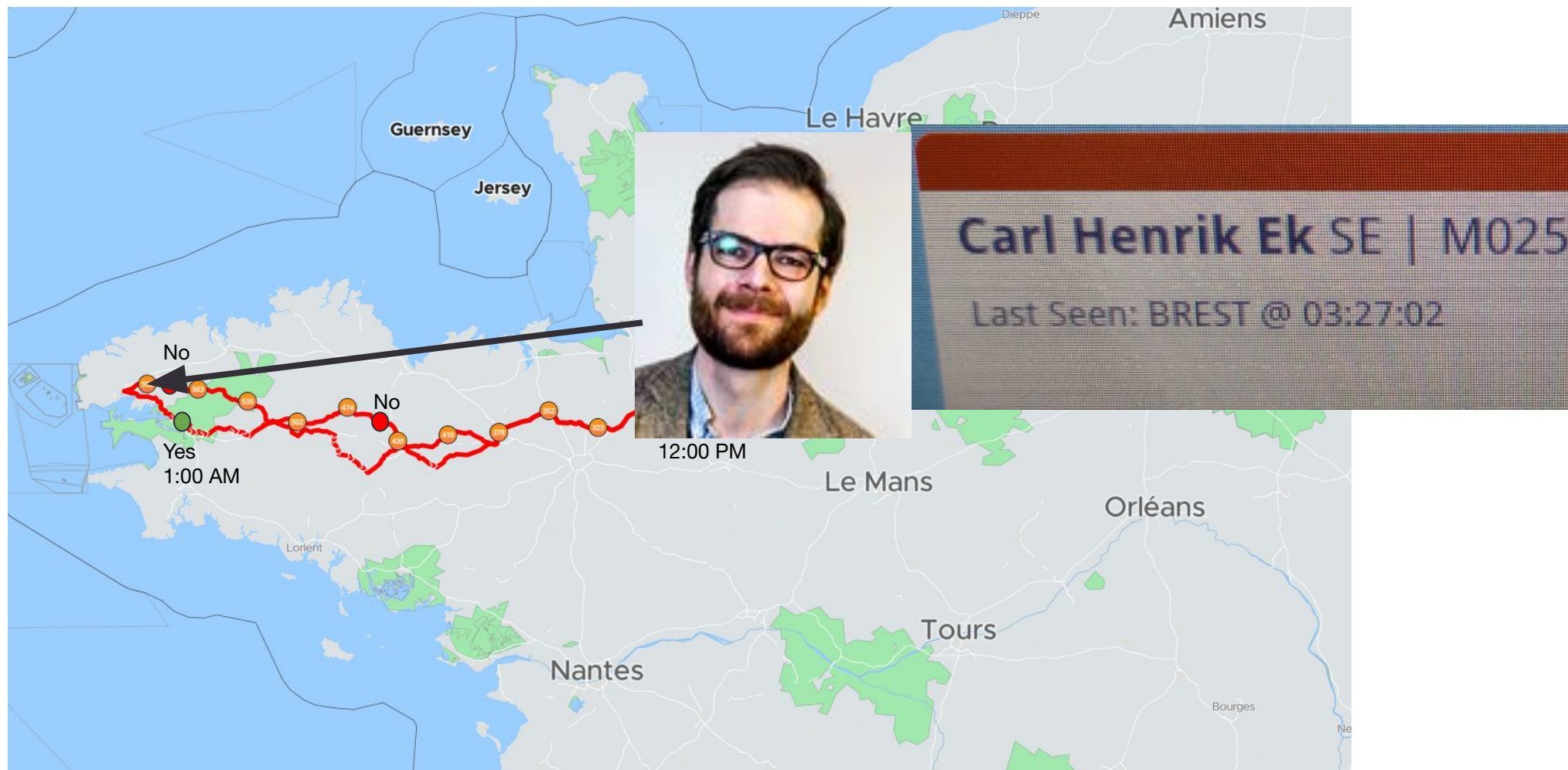


$$d = v \times t$$



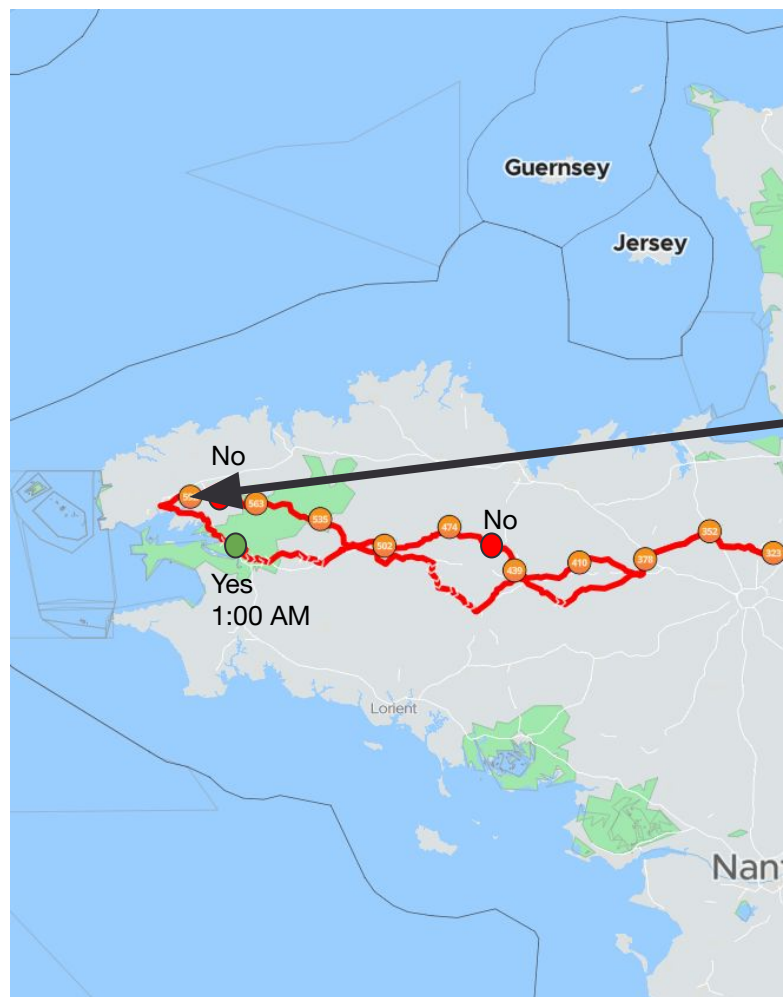
Where is Carl Henrik?

At 3:30 AM?



Where is Carl Henrik?

At 3:30 AM?



INTERPOL WANTED



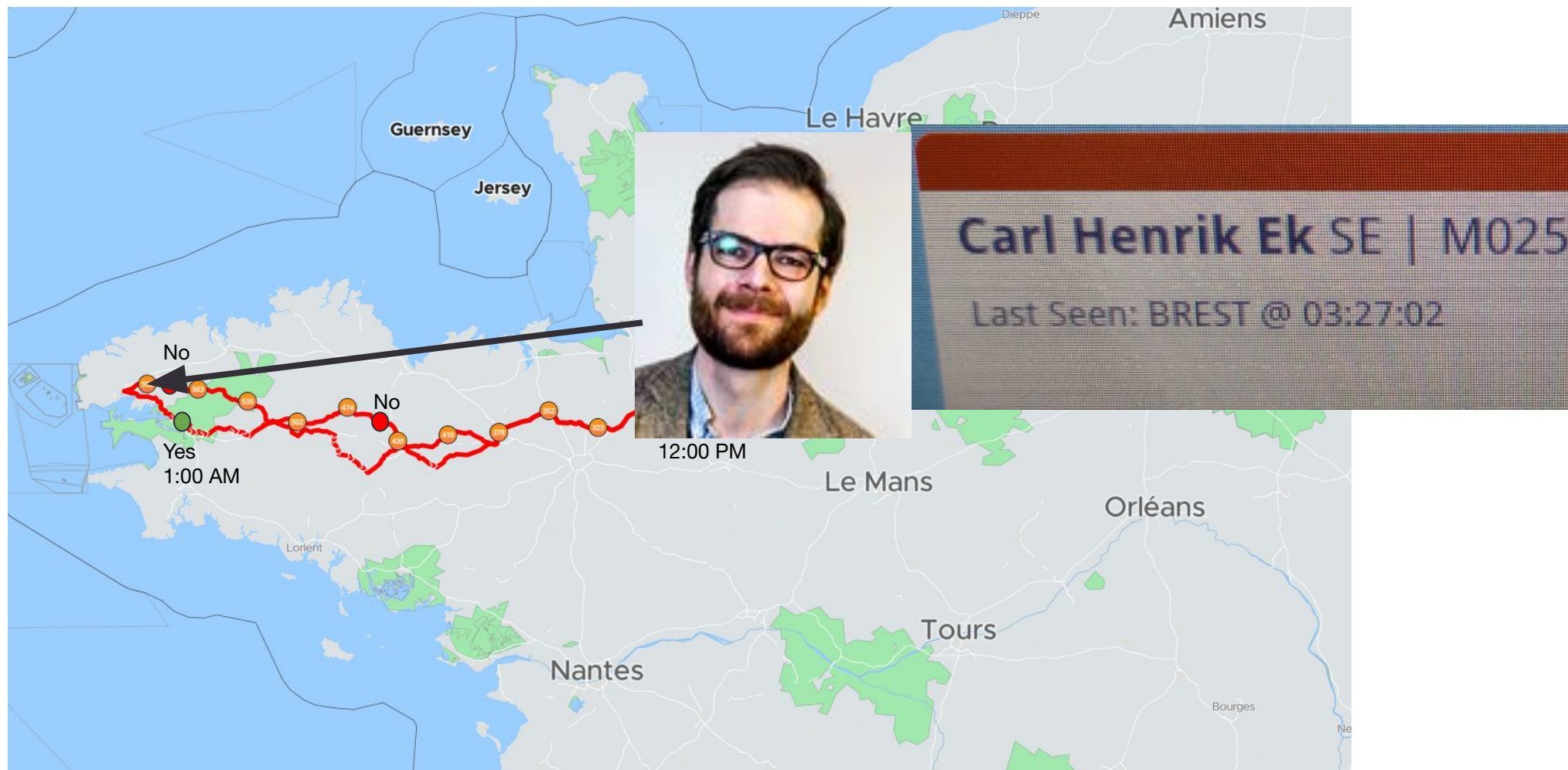
Carl Henrik Ek SE | M025

Last Seen: BREST @ 03:27:02

Your tip could be the missing piece in the puzzle.
If you have any information, contact your local police or go to
www.interpol.int ► Wanted persons

Where is Carl Henrik?

At 3:30 AM?



But can we do better than **random**???

What is Active Learning?

Bayesian search for learning functions

Sequential data collection

Let's make use of uncertainty estimates to make better models



Sequential data collection

Let's make use of uncertainty estimates to make better models

Collect initial data

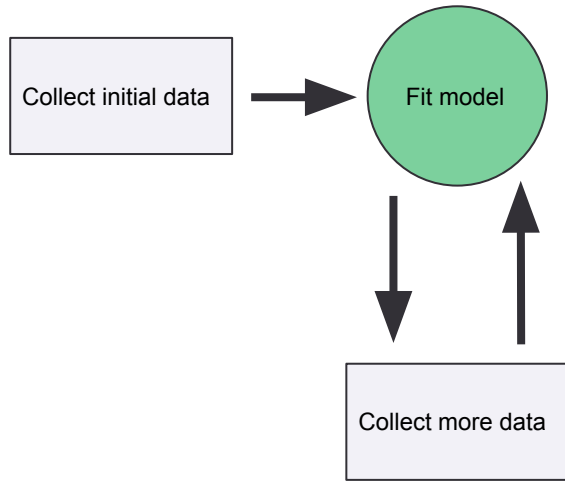
Sequential data collection

Let's make use of uncertainty estimates to make better models



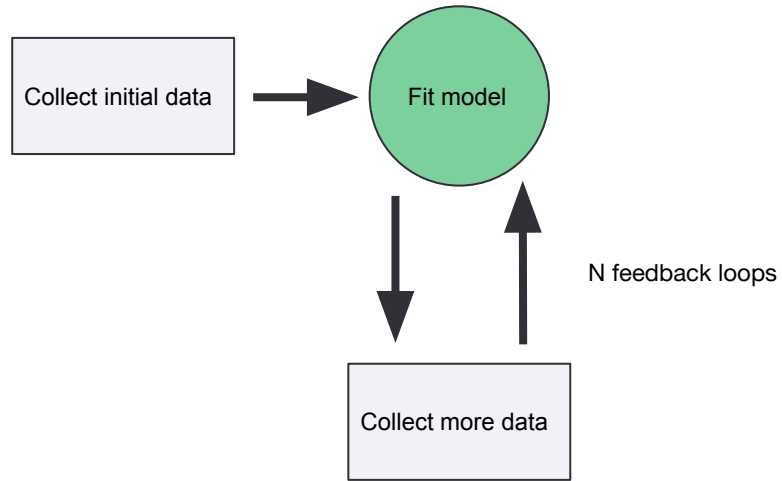
Sequential data collection

Let's make use of uncertainty estimates to make better models



Sequential data collection

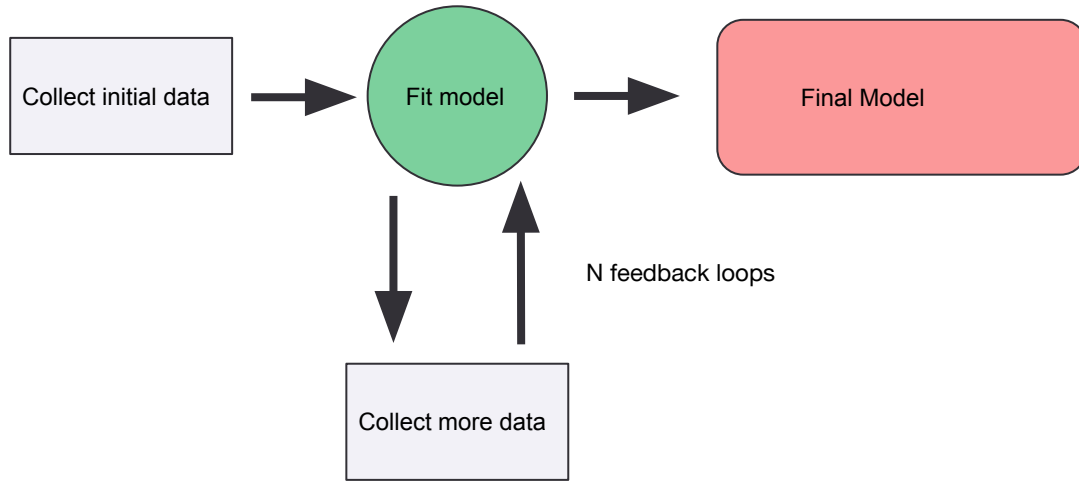
Let's make use of uncertainty estimates to make better models





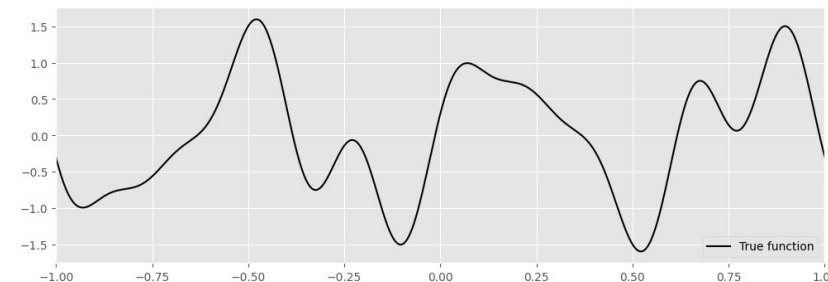
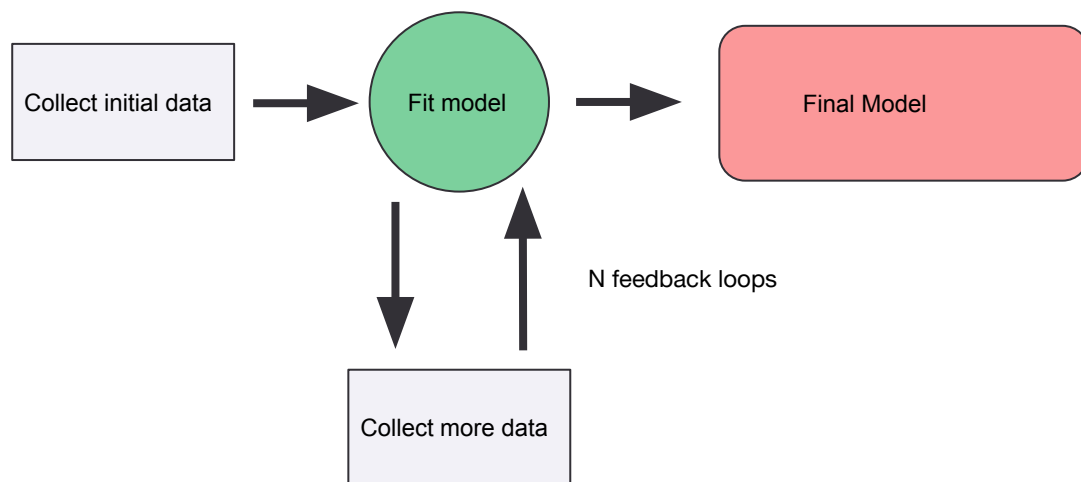
Sequential data collection

Let's make use of uncertainty estimates to make better models



Sequential data collection

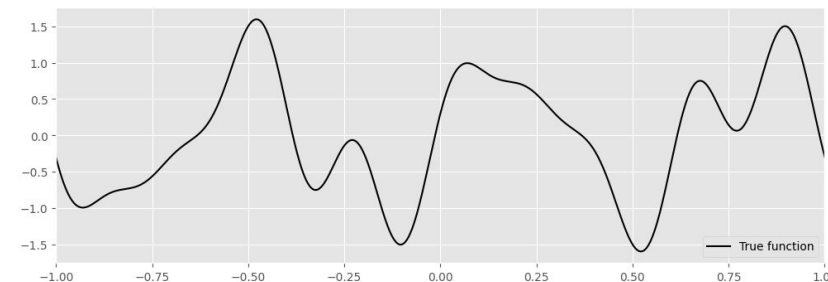
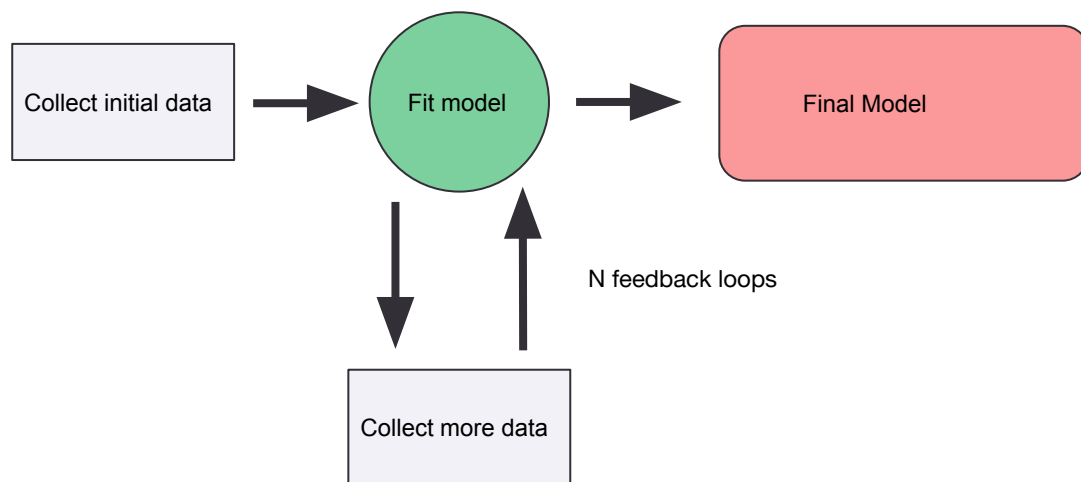
Let's make use of uncertainty estimates to make better models



0

Sequential data collection

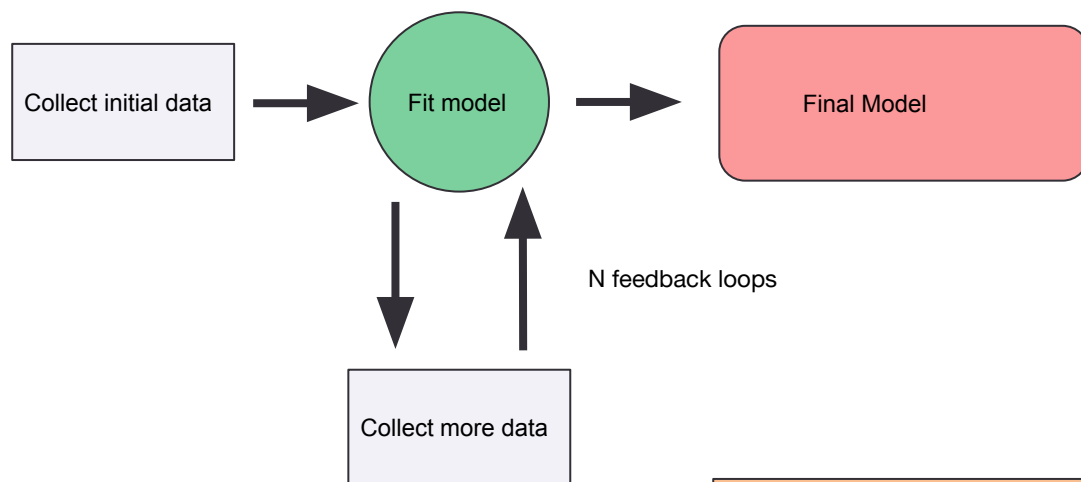
Let's make use of uncertainty estimates to make better models



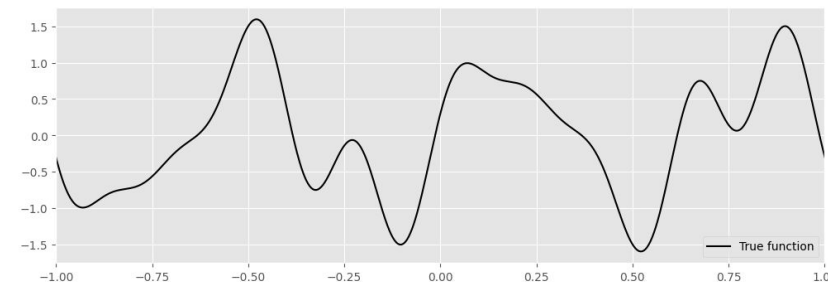
0

Sequential data collection

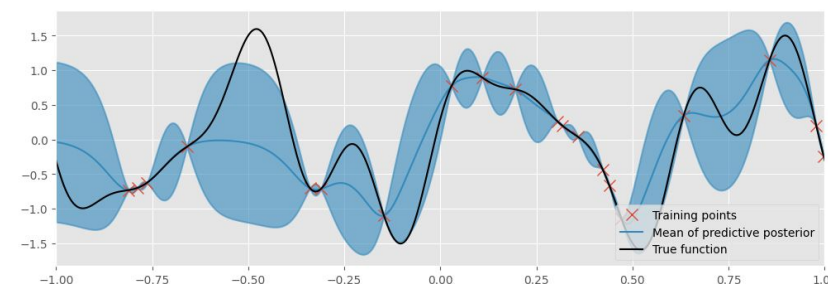
Let's make use of uncertainty estimates to make better models



These damn sausage
plots!!!



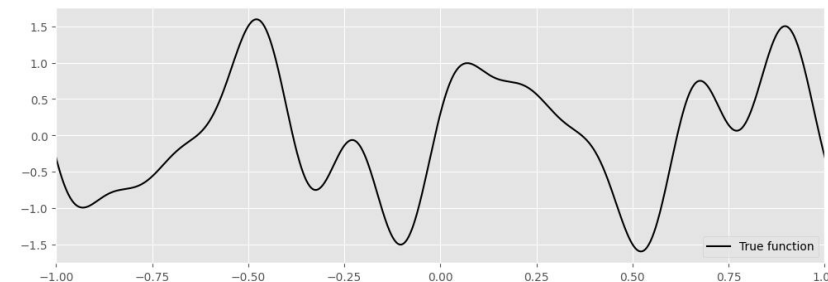
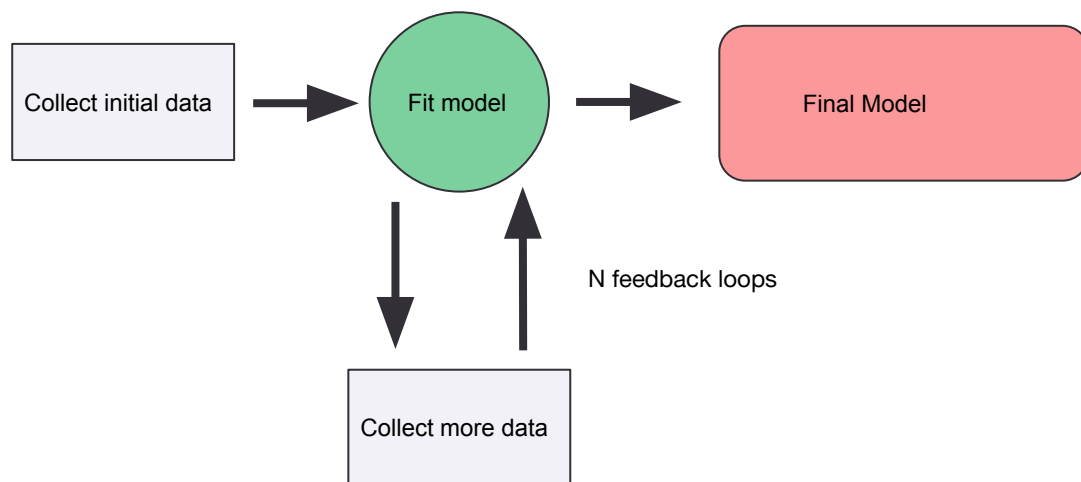
0



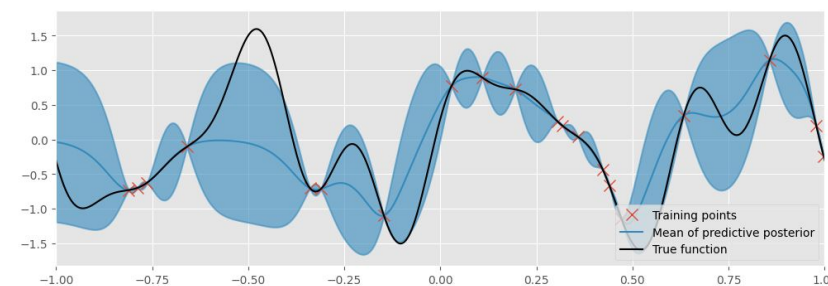
10

Sequential data collection

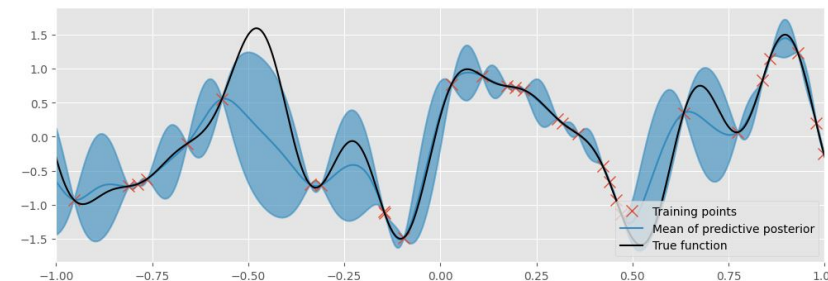
Let's make use of uncertainty estimates to make better models



0



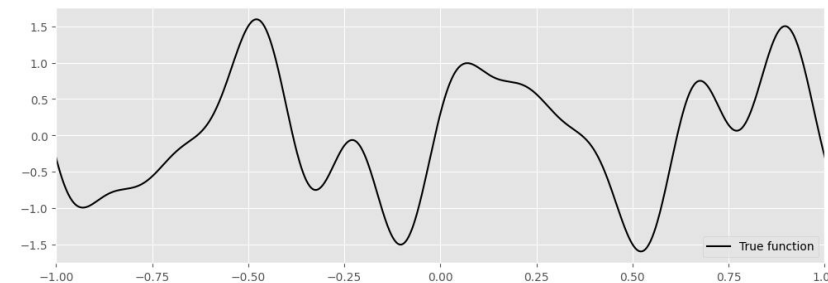
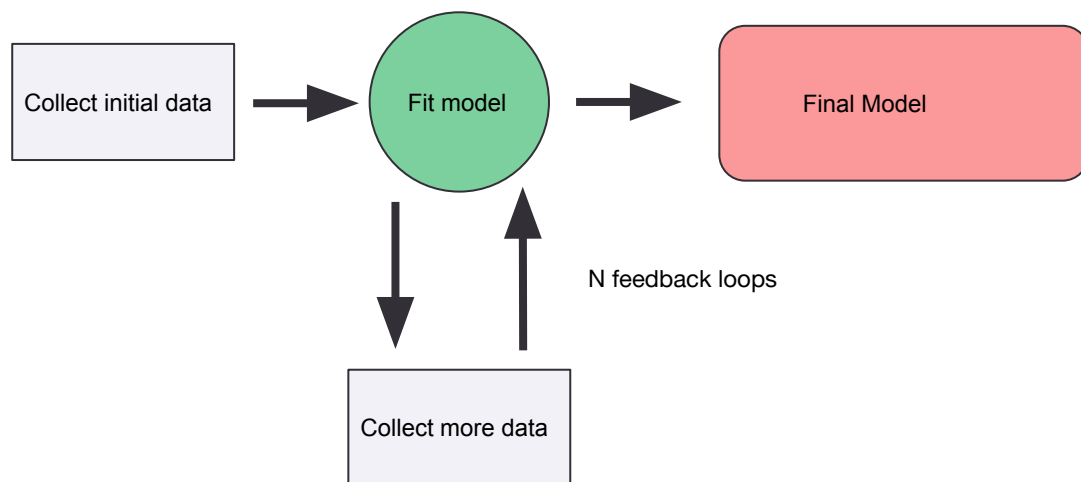
10



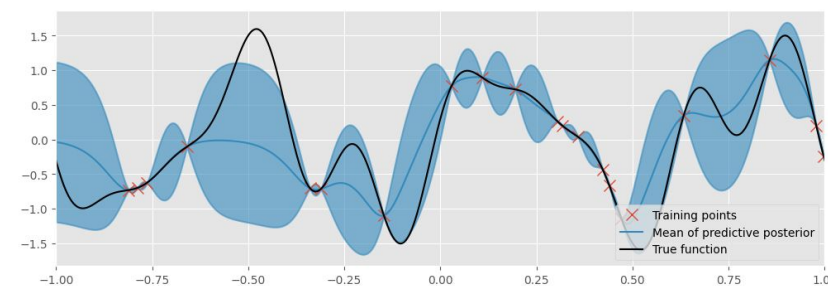
20

Sequential data collection

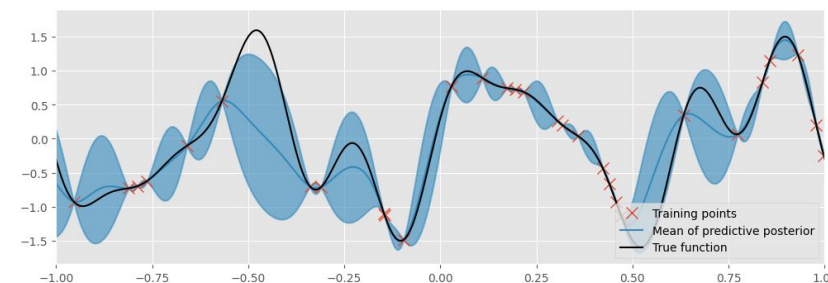
Let's make use of uncertainty estimates to make better models



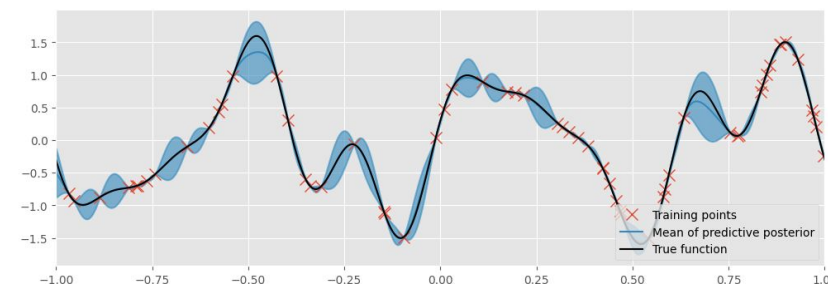
0



10



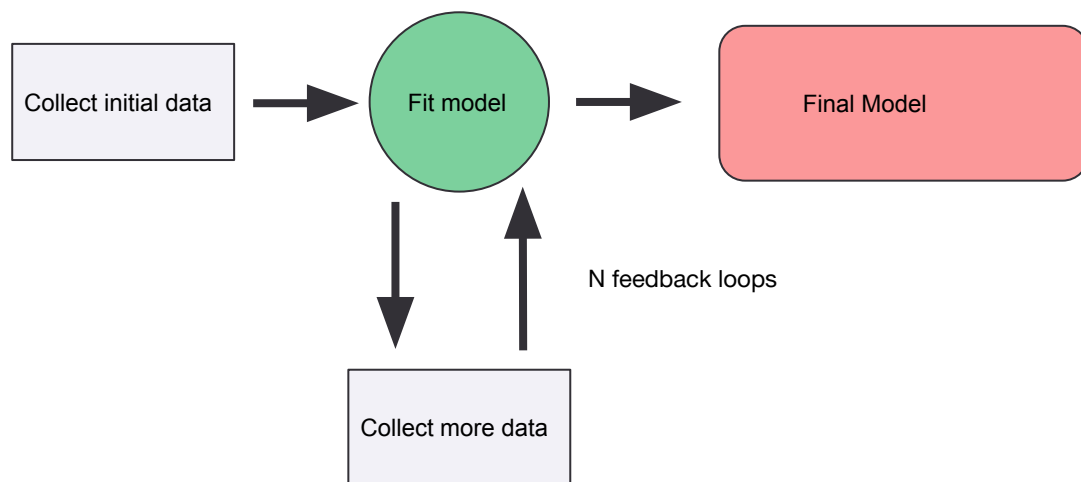
20



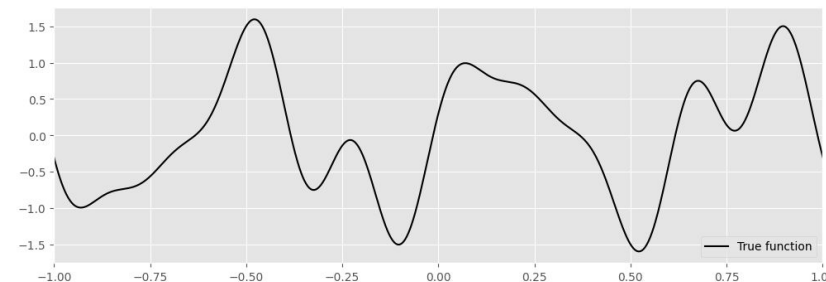
30

Sequential data collection

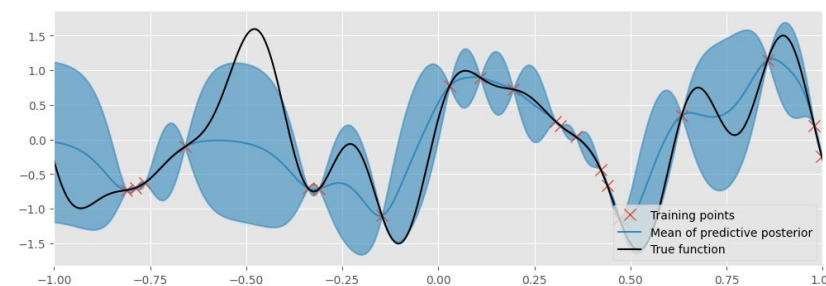
Let's make use of uncertainty estimates to make better models



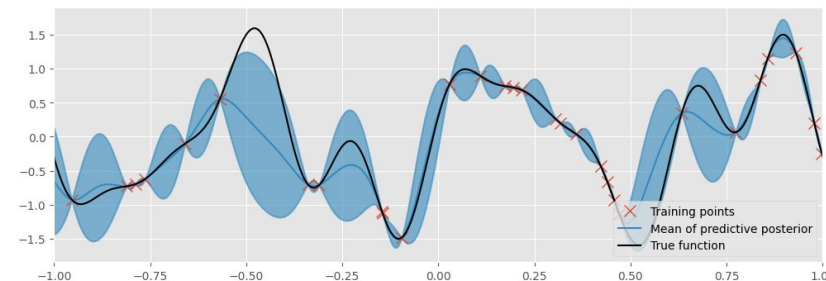
But can we do better than **random**???



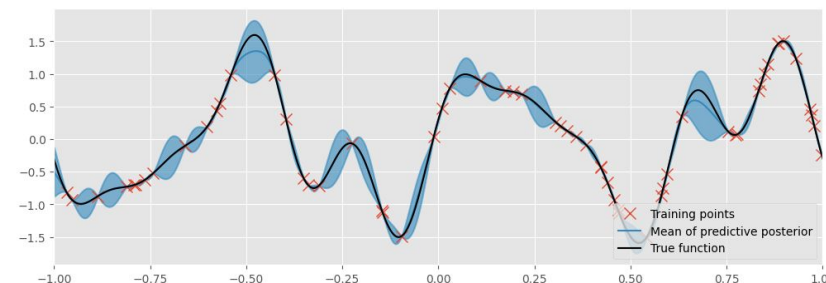
0



10



20



30

Active learning

Sequentially collecting more data to improve your model for the task at hand

Active learning

Sequentially collecting more data to improve your model for the task at hand

- I care about **regression** —> collect data to improve global model accuracy

Active learning

Sequentially collecting more data to improve your model for the task at hand

- I care about **regression** —> collect data to improve global model accuracy
- I care about the **maximum** value of my process —> collect data in promising regions (Bayesian Optimisation)

Active learning

Sequentially collecting more data to improve your model for the task at hand

- I care about **regression** —> collect data to improve global model accuracy
- I care about the **maximum** value of my process —> collect data in promising regions (Bayesian Optimisation)
- I'm interested in **multiple objectives** -> populate the Pareto front (Multi-objective Bayesian Optimisation)

Active learning

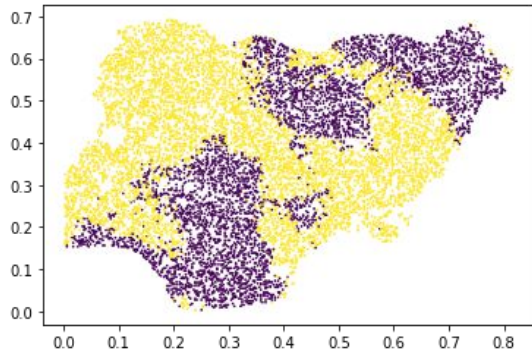
Sequentially collecting more data to improve your model for the task at hand

- I care about **regression** —> collect data to improve global model accuracy
- I care about the **maximum** value of my process —> collect data in promising regions (Bayesian Optimisation)
- I'm interested in **multiple objectives** -> populate the Pareto front (Multi-objective Bayesian Optimisation)
- I care about predicting a **threshold** -> choose data close to threshold (level-set design)

Active learning

Sequentially collecting more data to improve your model for the task at hand

- I care about **regression** —> collect data to improve global model accuracy
- I care about the **maximum** value of my process —> collect data in promising regions (Bayesian Optimisation)
- I'm interested in **multiple objectives** -> populate the Pareto front (Multi-objective Bayesian Optimisation)
- I care about predicting a **threshold** -> choose data close to threshold (level-set design)

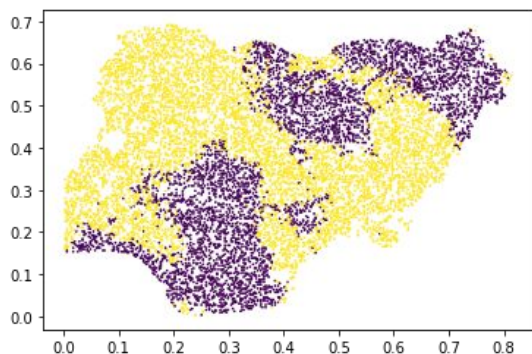


Malaria incidence
in Nigeria

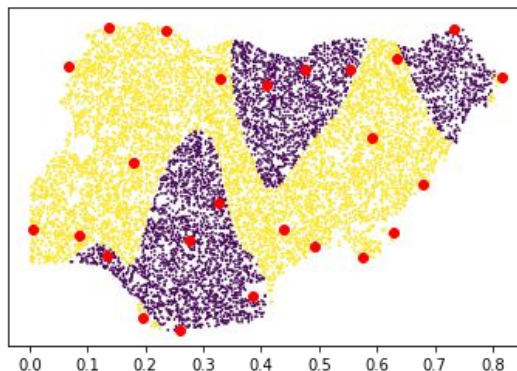
Active learning

Sequentially collecting more data to improve your model for the task at hand

- I care about **regression** —> collect data to improve global model accuracy
- I care about the **maximum** value of my process —> collect data in promising regions (Bayesian Optimisation)
- I'm interested in **multiple objectives** -> populate the Pareto front (Multi-objective Bayesian Optimisation)
- I care about predicting a **threshold** -> choose data close to threshold (level-set design)



Malaria incidence
in Nigeria

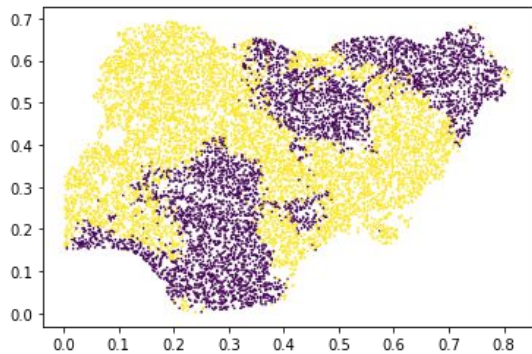


Model on Random
data

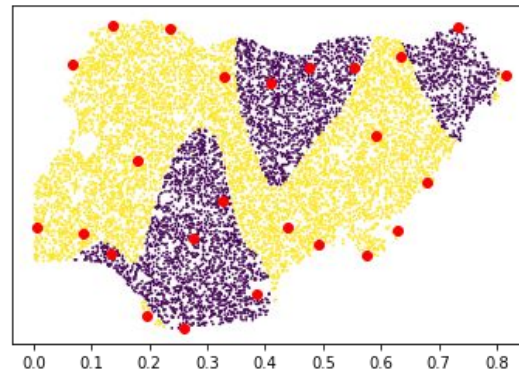
Active learning

Sequentially collecting more data to improve your model for the task at hand

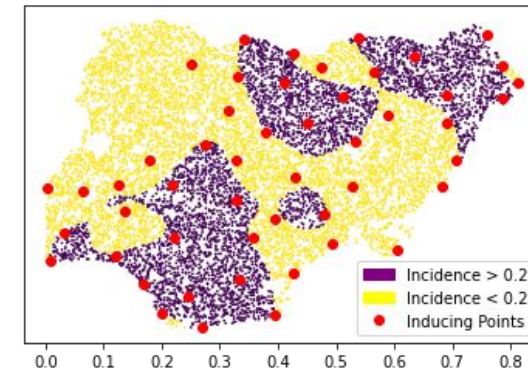
- I care about **regression** —> collect data to improve global model accuracy
- I care about the **maximum** value of my process —> collect data in promising regions (Bayesian Optimisation)
- I'm interested in **multiple objectives** -> populate the Pareto front (Multi-objective Bayesian Optimisation)
- I care about predicting a **threshold** -> choose data close to threshold (level-set design)



Malaria incidence
in Nigeria



Model on Random
data



Model from data
chosen by Active
learning

So, Bayesian Optimisation?

i.e. Active learning for optimisation

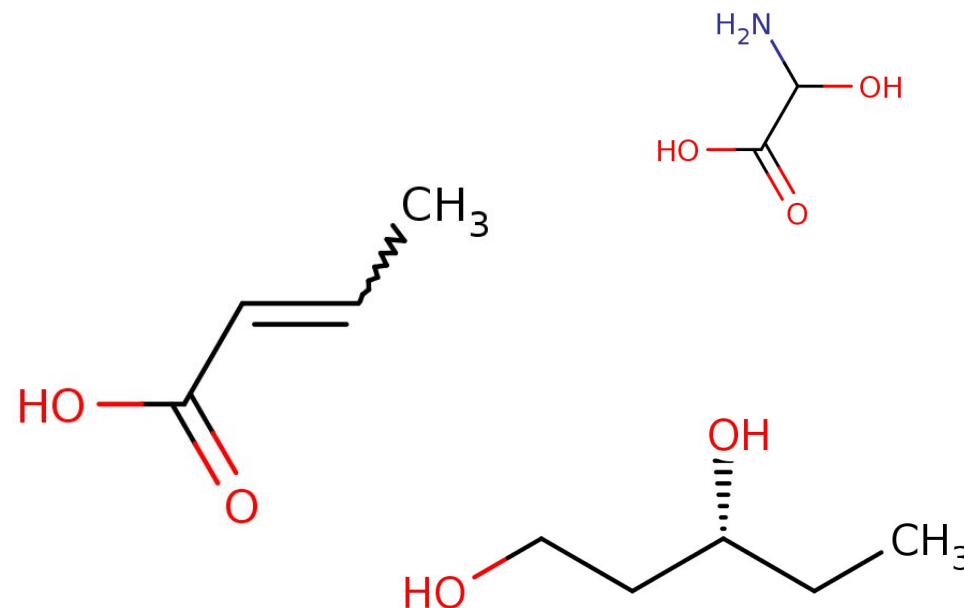
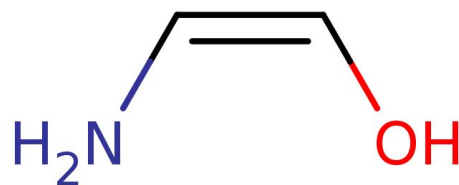
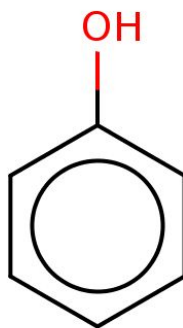
A molecular design pipeline

Efficiently explore molecule space

A molecular design pipeline

Efficiently explore molecule space

- **Large** library of candidates

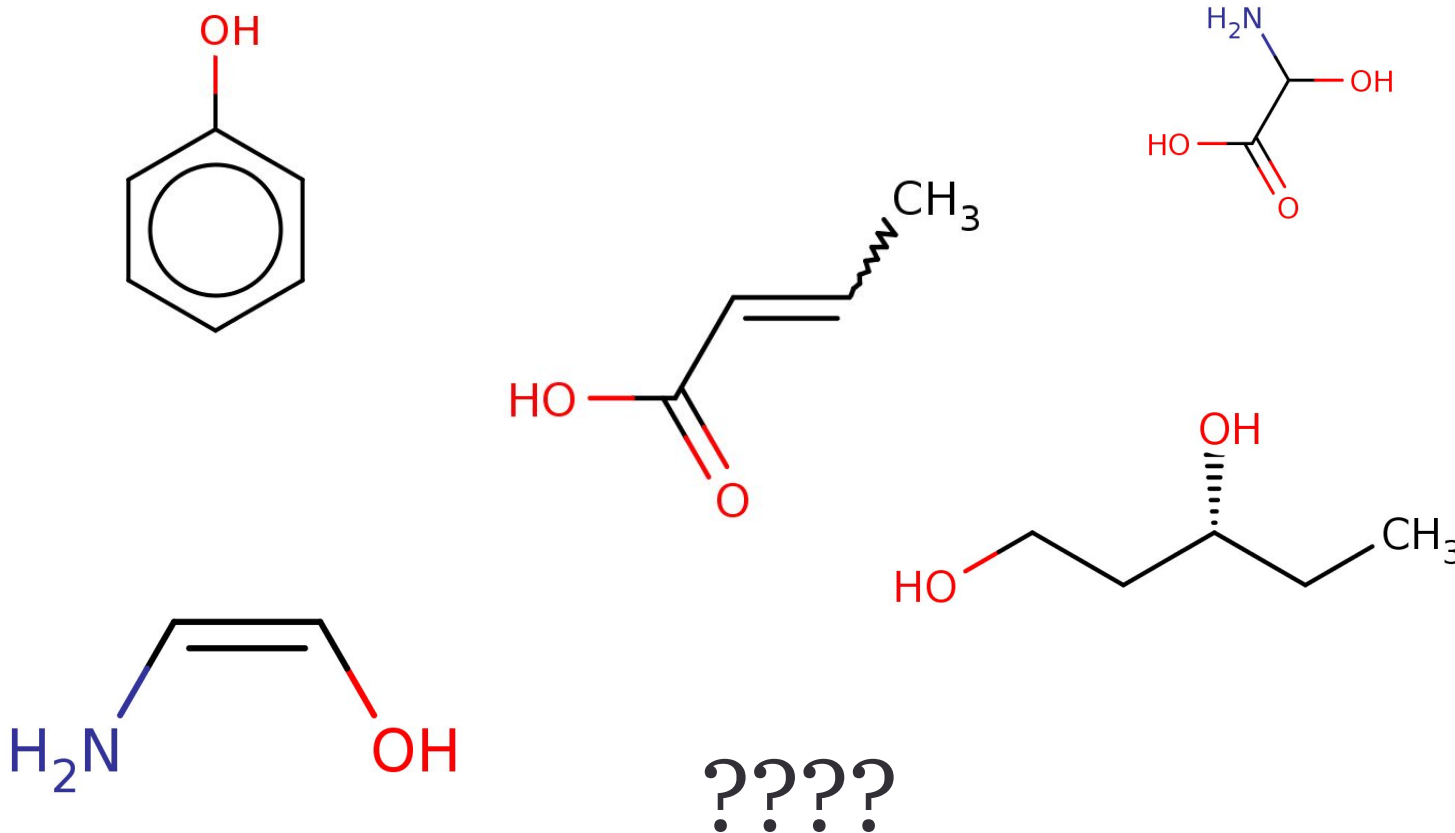


?????

A molecular design pipeline

Efficiently explore molecule space

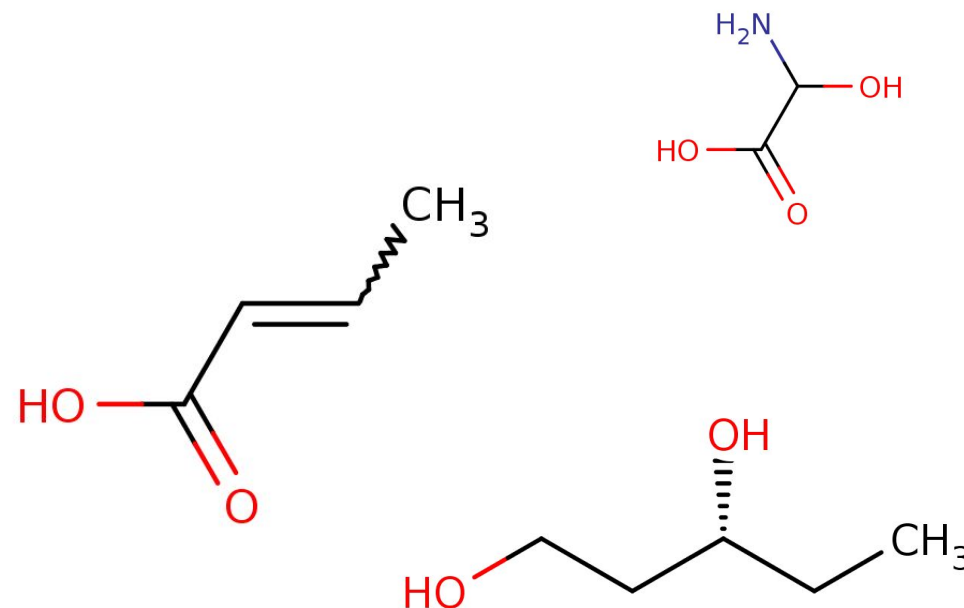
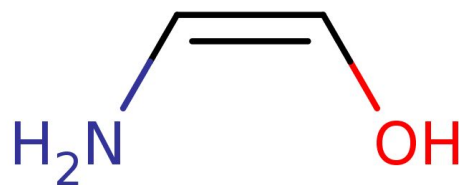
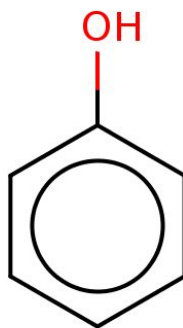
- **Large** library of candidates
- **Expensive** experiments (<10)



A molecular design pipeline

Efficiently explore molecule space

- **Large** library of candidates
- **Expensive** experiments (<10) (**IN A LAB !!!**)

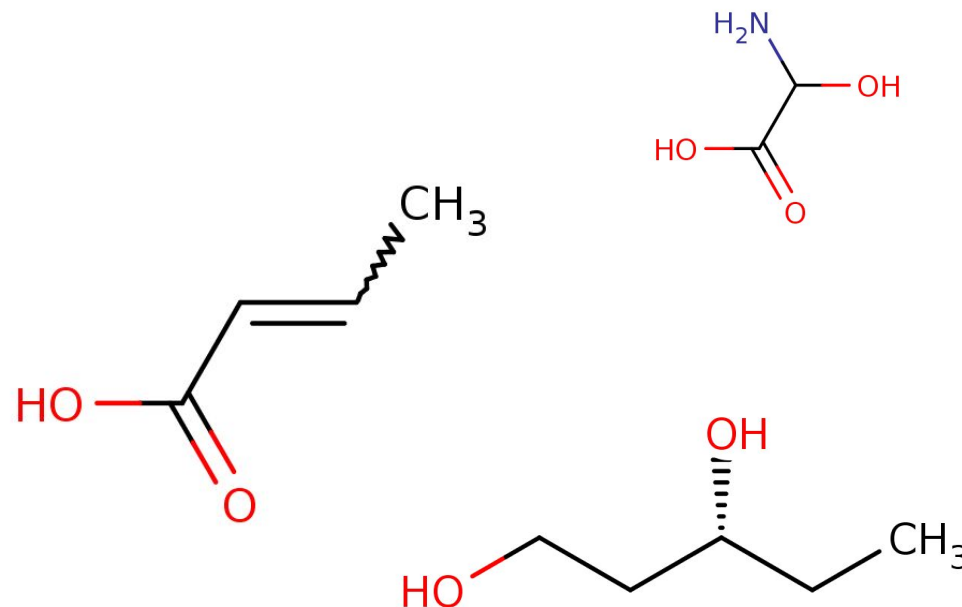
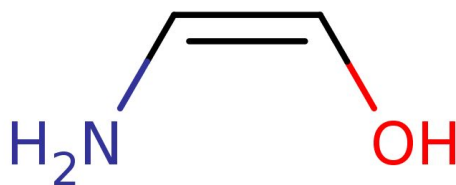
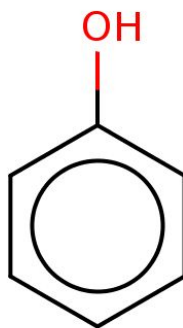


?????

A molecular design pipeline

Efficiently explore molecule space

- **Large** library of candidates
- **Expensive** experiments (<10)
- High degree of **parallelism**

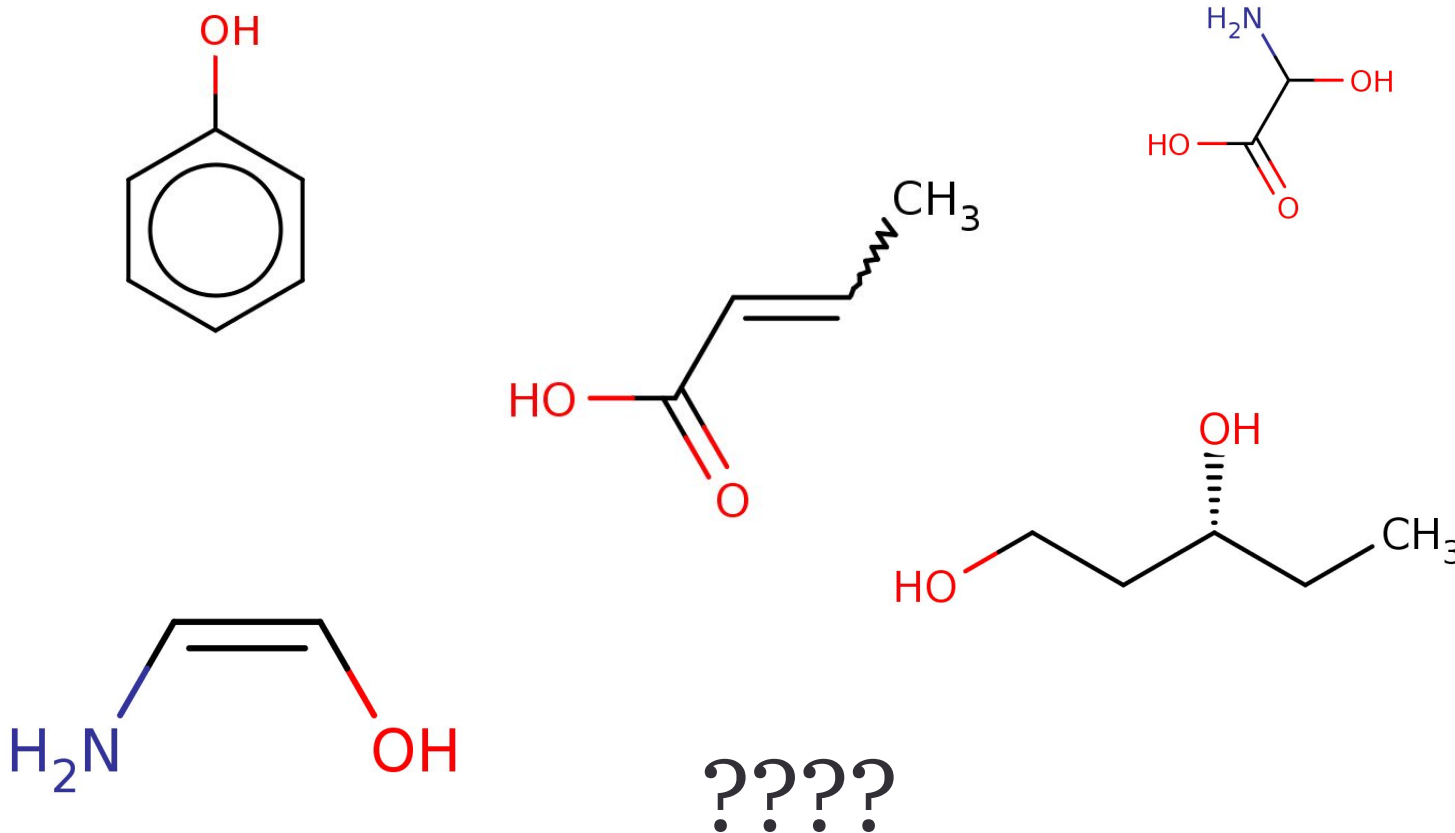


?????

A molecular design pipeline

Efficiently explore molecule space

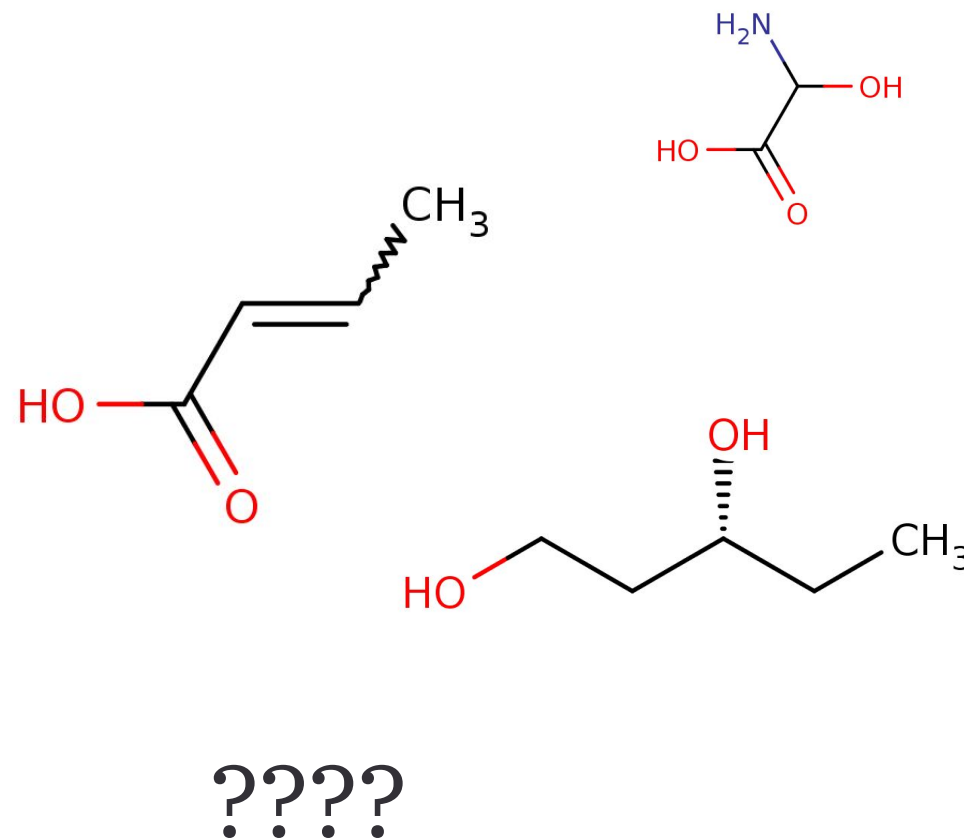
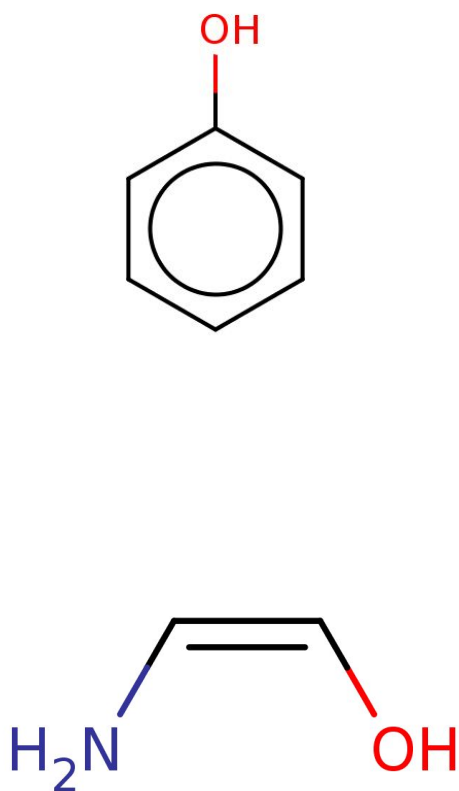
- **Large** library of candidates
- **Expensive** experiments (<10)
- High degree of **parallelism**
- Want molecules with high **affinity**



A molecular design pipeline

Efficiently explore molecule space

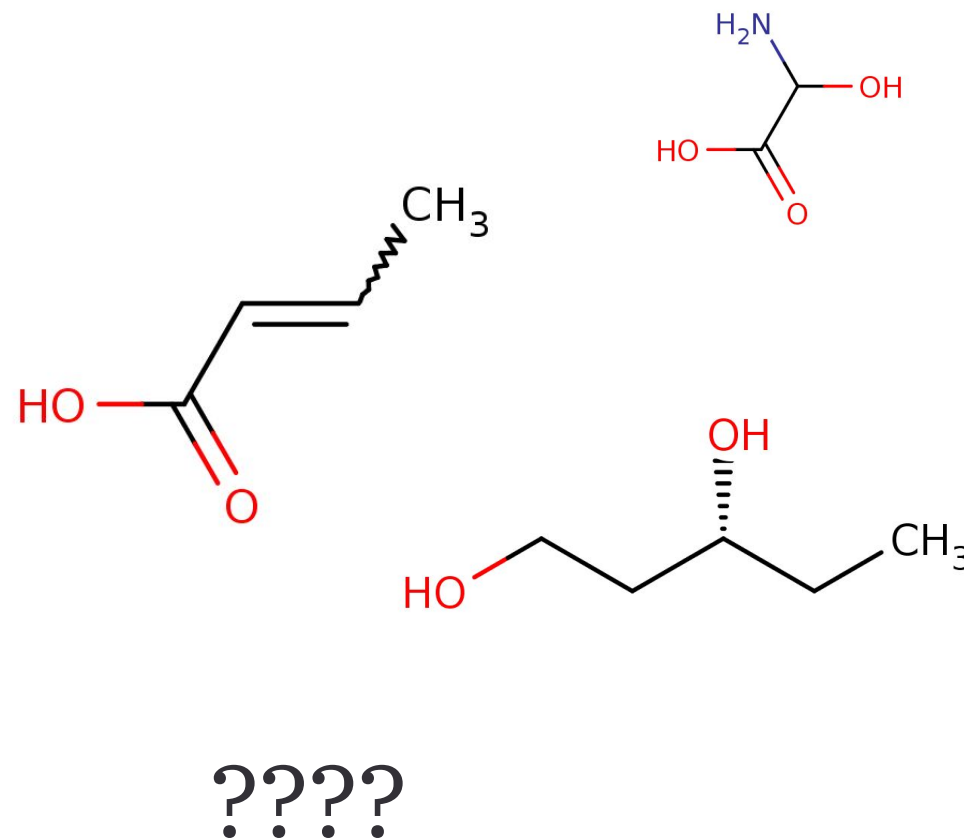
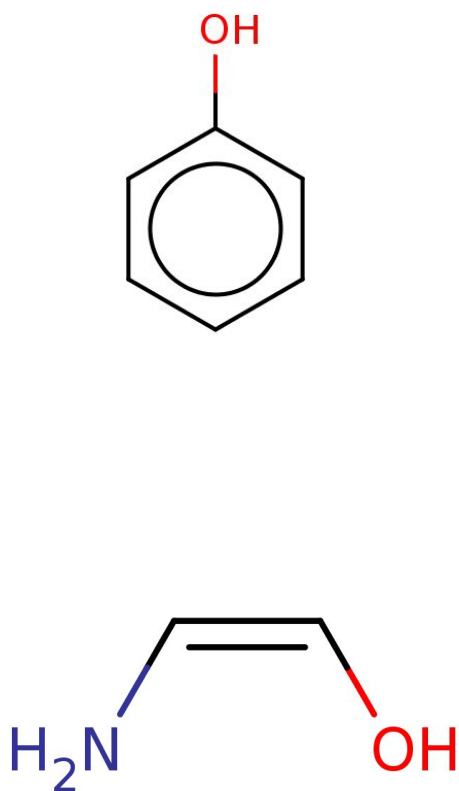
- **Large** library of candidates
- **Expensive** experiments (<10)
- High degree of **parallelism**
- Want molecules with high **affinity**
 - Also easy to make



A molecular design pipeline

Efficiently explore molecule space

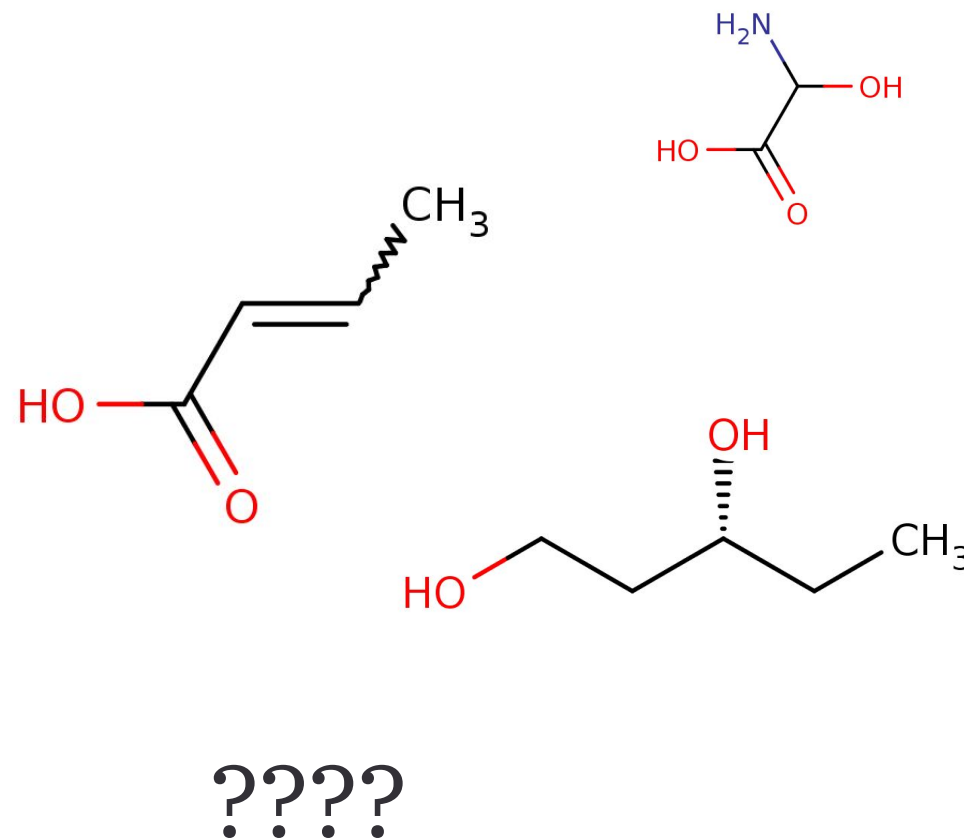
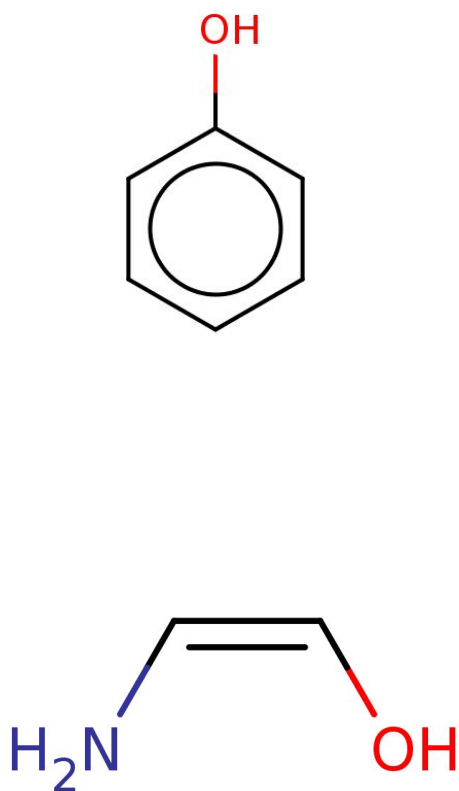
- **Large** library of candidates
- **Expensive** experiments (<10)
- High degree of **parallelism**
- Want molecules with high **affinity**
 - Also easy to make
 - Don't stick to themselves



A molecular design pipeline

Efficiently explore molecule space

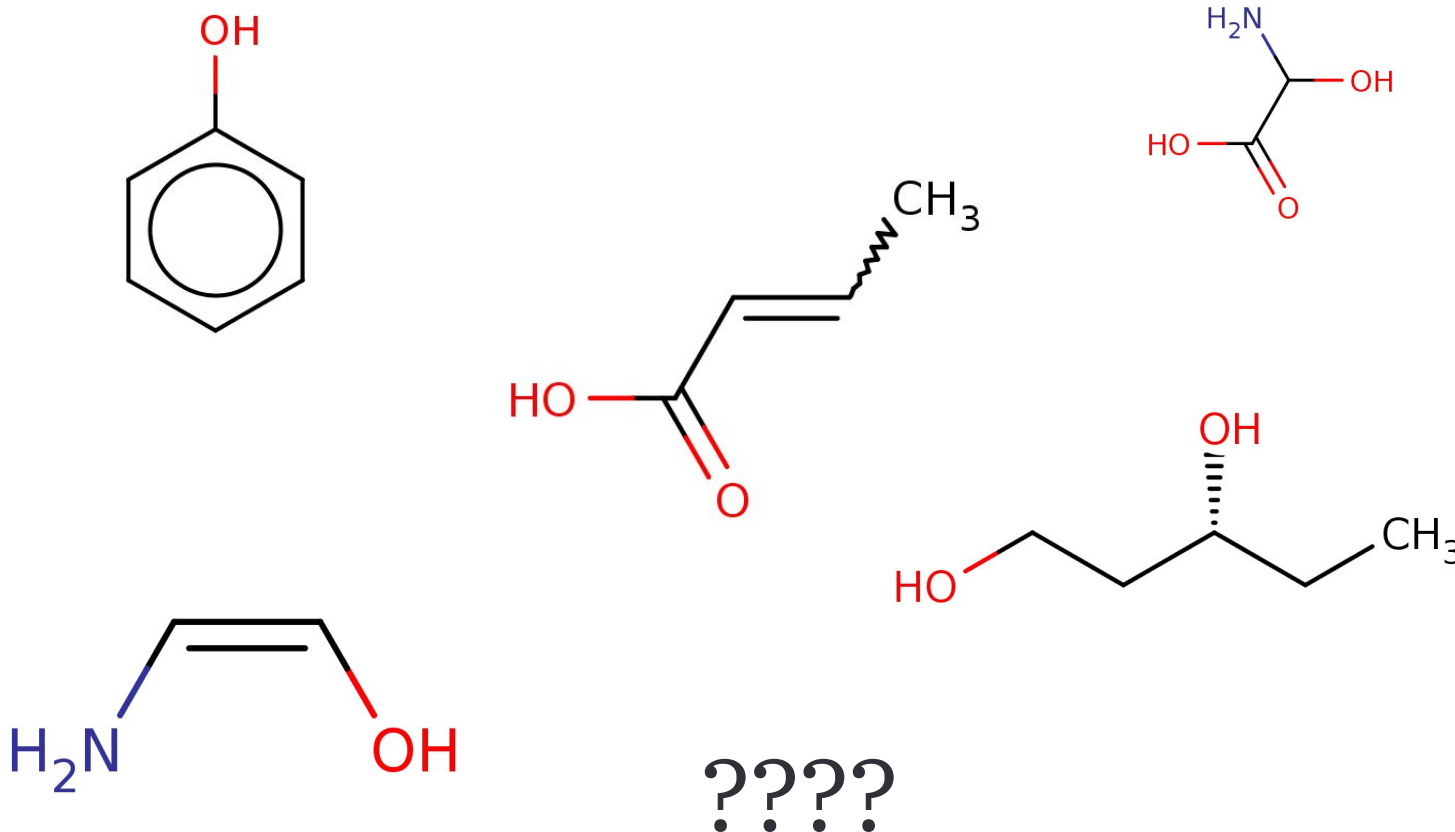
- **Large** library of candidates
- **Expensive** experiments (<10)
- High degree of **parallelism**
- Want molecules with high **affinity**
 - Also easy to make
 - Don't stick to themselves
 - Stable



A molecular design pipeline

Efficiently explore molecule space

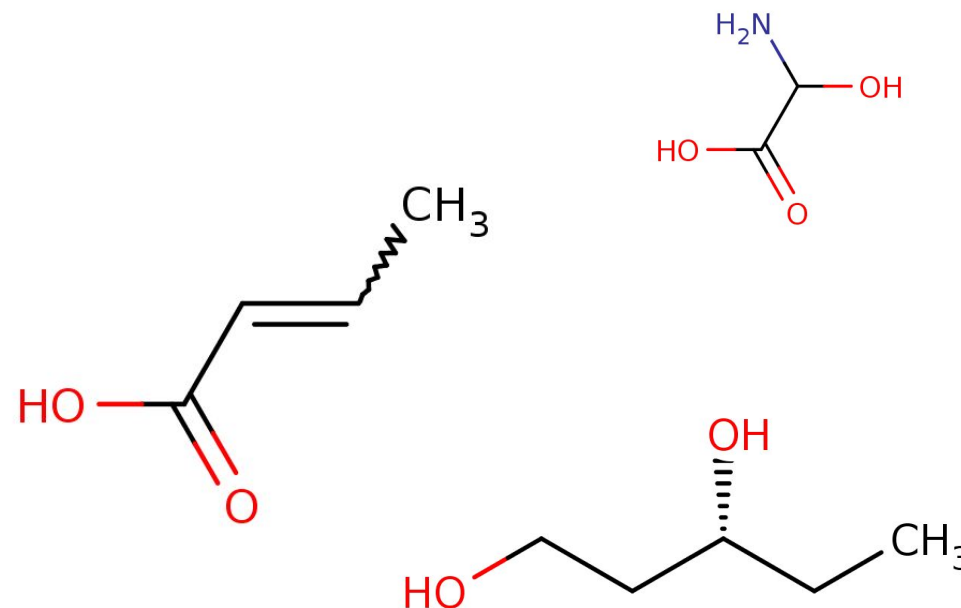
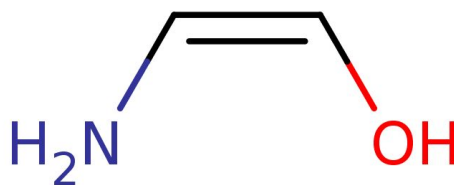
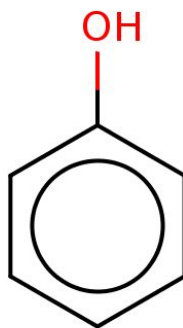
- **Large** library of candidates
- **Expensive** experiments (<10)
- High degree of **parallelism**
- Want molecules with high **affinity**
 - Also easy to make
 - Don't stick to themselves
 - Stable
 - In a new area of "patent space"



A molecular design pipeline

Efficiently explore molecule space

- **Large** library of candidates
- **Expensive** experiments (<10)
- High degree of **parallelism**
- Want molecules with high **affinity**
 - Also easy to make
 - Don't stick to themselves
 - Stable
 - In a new area of "patent space"

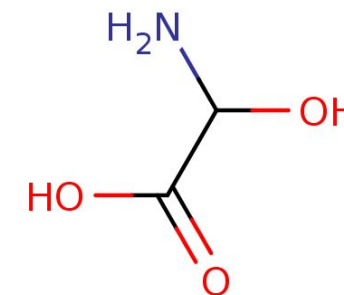
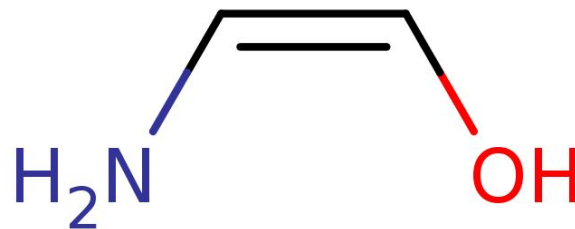
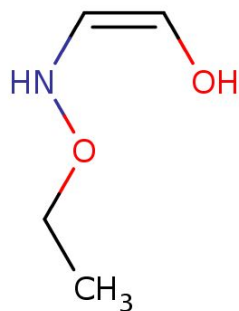
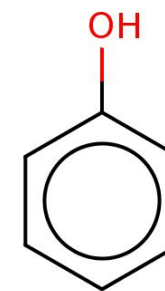
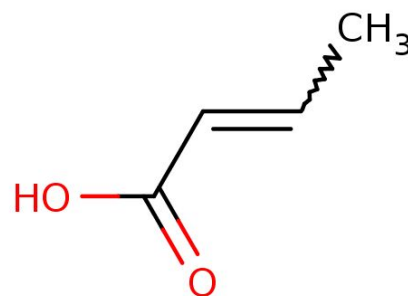
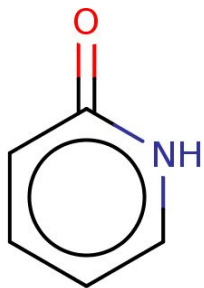


?????

Any ideas?

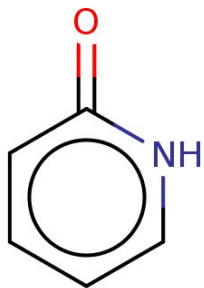
A Simpler Example

Can evaluate **at most** 4

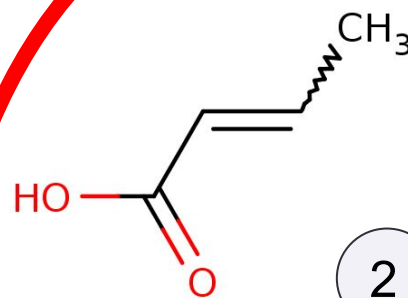
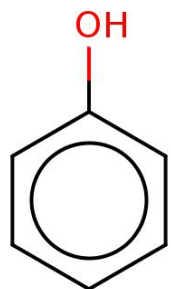


A Simpler Example (grouped)

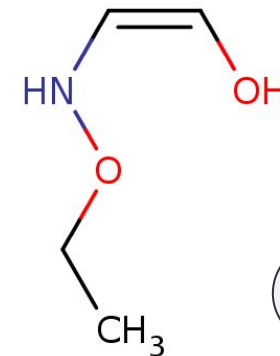
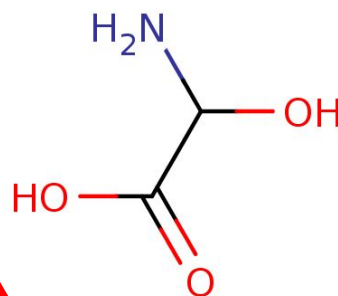
Can evaluate **at most** 4



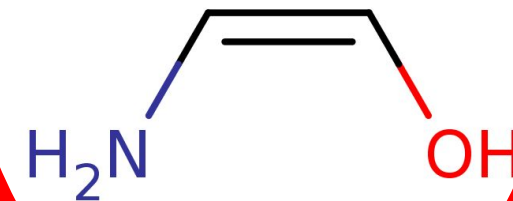
1



2

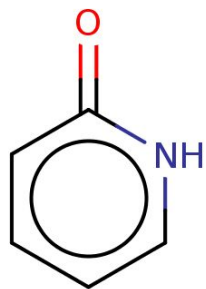


3

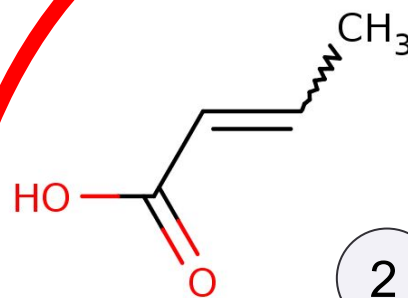
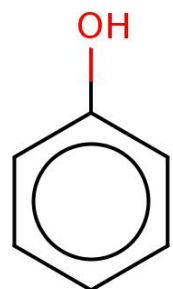


A Simpler Example (grouped)

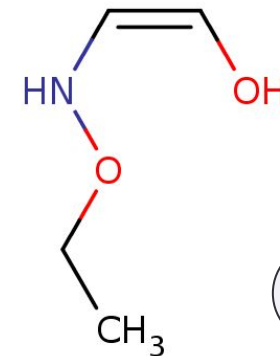
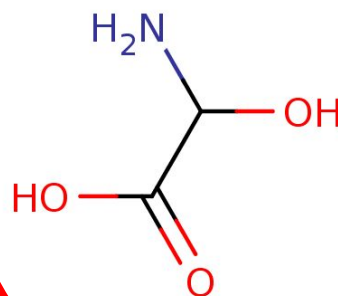
Can evaluate **at most** 4



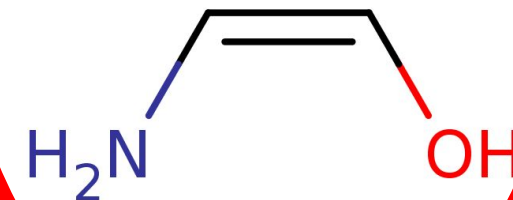
1



2



3

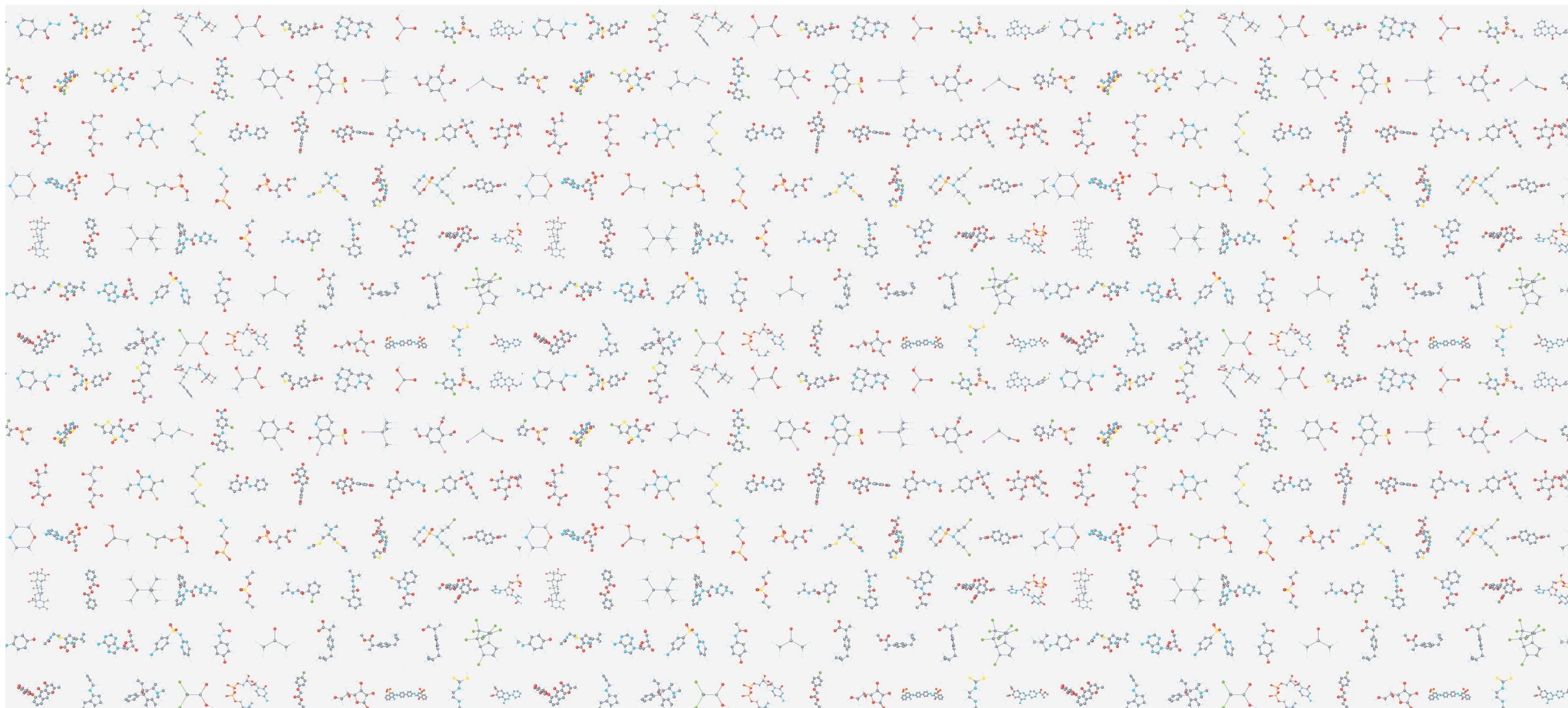


Explore v.s. exploit?



What about at scale?


week





What about at scale?

week



Use a GP!

An Aside: GPs for Molecules

Structured Input Spaces

$$y_i = f(\text{molecule}_i) + \epsilon_i$$

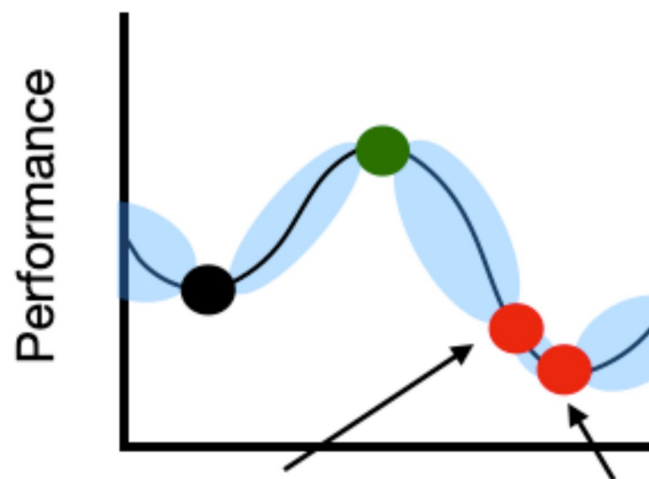
$$D_N = \{(\text{molecule}_i, y_i)\}_i^N$$

An Aside: GPs for Molecules

Structured Input Spaces

$$y_i = f(\text{molecule}_i) + \epsilon_i$$

$$D_N = \{(\text{molecule}_i, y_i)\}_i^N$$



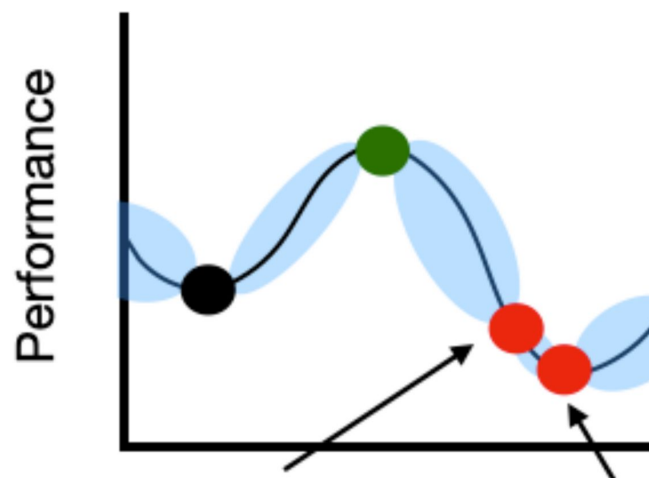
What do we require to
define a GP?

An Aside: GPs for Molecules

Structured Input Spaces

$$y_i = f(\text{molecule}_i) + \epsilon_i$$

$$D_N = \{(\text{molecule}_i, y_i)\}_i^N$$



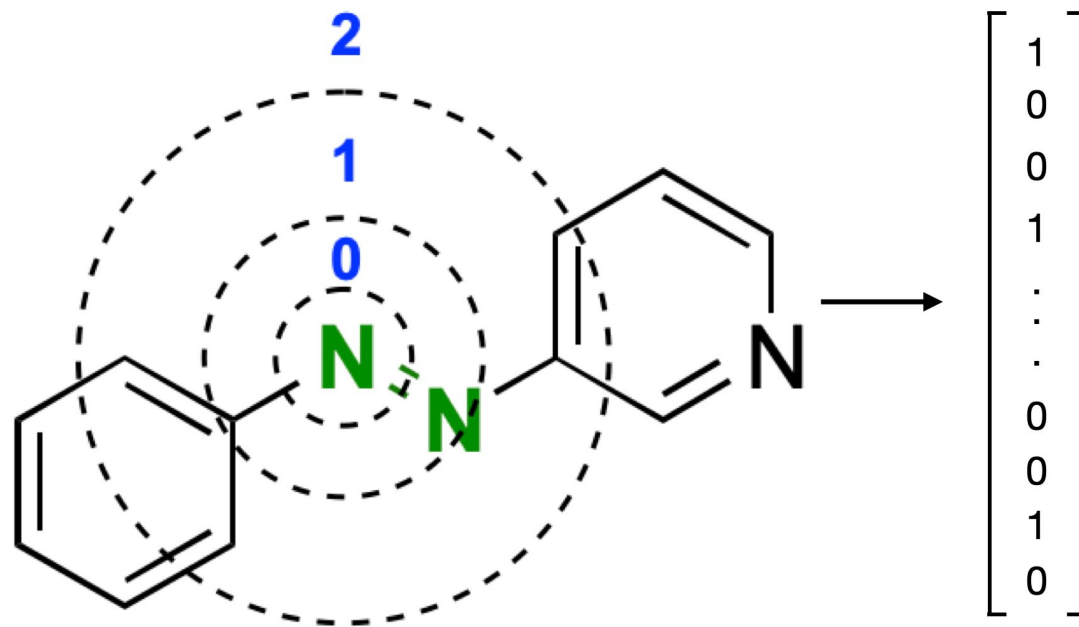
What do we require to define a GP?

$$k(\text{molecule}_i, \text{molecule}_j) = ?$$

An Aside: GPs for Molecules

Fingerprint Kernels

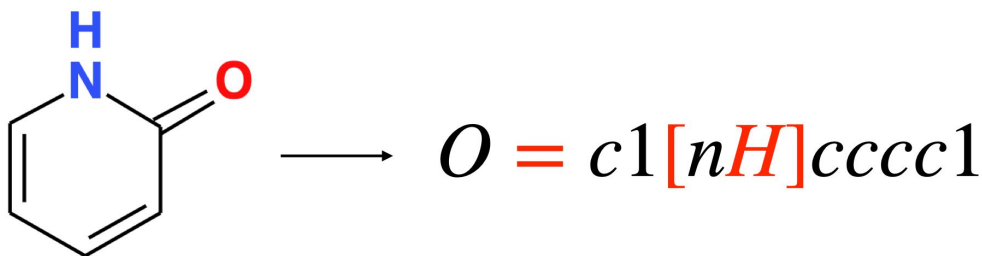
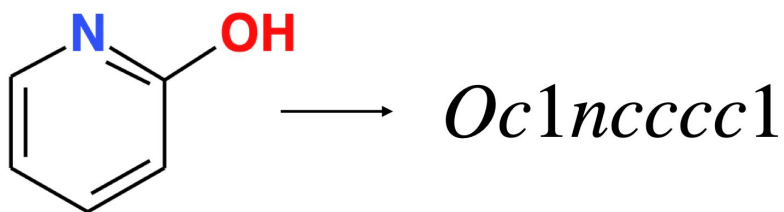
$$k(\text{molecule}_i, \text{molecule}_j) = k_{\text{linear}}(\Phi(\text{molecule}_i), \Phi(\text{molecule}_j))$$



An Aside: GPs for Molecules

String kernels between SMILES strings

$$k(\text{mol}_i, \text{mol}_j) = k(\text{str}(\text{mol}_i), \text{str}(\text{mol}_j))$$



Automatically choosing next molecules

Using GP posteriors and utility functions


Automatically choosing next molecules

Using GP posteriors and utility functions

- $U_f(\text{molecule})$: what is the utility of evaluating  (if it will return f)

Automatically choosing next molecules

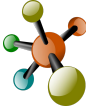
Using GP posteriors and utility functions

- $U_f(\text{molecule})$: what is the utility of evaluating  (if it will return f)
- f^\star Is best so far



Automatically choosing next molecules


Using GP posteriors and utility functions

- $U_f(\text{molecule})$: what is the utility of evaluating  (if it will return f)
- f^\star Is best so far
- Has there been an improvement? $U_f(\text{molecule}) = \mathbb{1}_{(f > f^\star)}$



Automatically choosing next molecules

Using GP posteriors and utility functions

- $U_f(\text{molecule})$: what is the utility of evaluating  (if it will return f)
- f^\star Is best so far
- Has there been an improvement? $U_f(\text{molecule}) = \mathbb{1}_{(f > f^\star)}$
- How big was the improvement? $U_f(\text{molecule}) = \max(f - f^\star, 0)$

Automatically choosing next molecules

Using GP posteriors and utility functions

- $\alpha(\text{molecule}) = \mathbb{E}_f[U_f(\text{molecule})]$: what utility is predicted by my model of f



Automatically choosing next molecules

Using GP posteriors and utility functions

- $\alpha(\text{molecule}) = \mathbb{E}_f[U_f(\text{molecule})]$: what utility is predicted by my model of f

- What the probability of improvement? $\alpha_{\text{PI}}(\text{molecule}) = \mathbb{E}_f[\mathbb{1}_{(f > f^*)}]$

Automatically choosing next molecules

Using GP posteriors and utility functions

- $\alpha(\text{molecule}) = \mathbb{E}_f[U_f(\text{molecule})]$: what utility is predicted by my model of f
 - What the probability of improvement? $\alpha_{\text{PI}}(\text{molecule}) = \mathbb{E}_f[\mathbb{1}_{(f > f^*)}]$
 - How much improvement do we expect? $\alpha_{\text{EI}}(\text{molecule}) = \mathbb{E}_f[\max(f - f^*, 0)]$

Automatically choosing next molecules

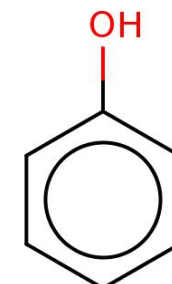
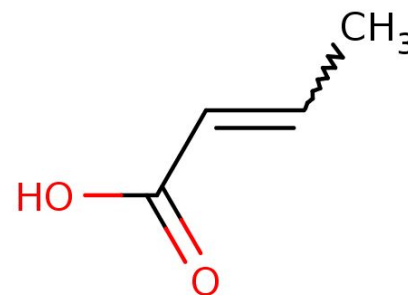
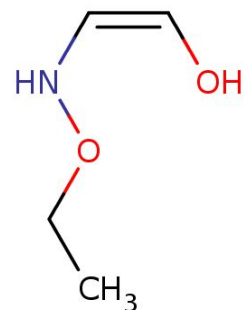
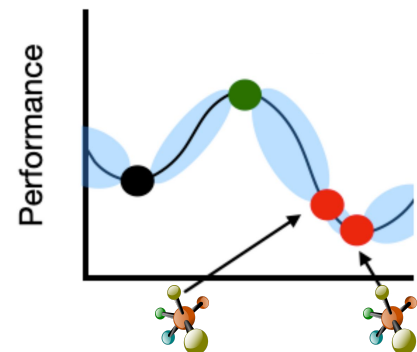
Using GP posteriors and utility functions

- $\alpha(\text{molecule}) = \mathbb{E}_f[U_f(\text{molecule})]$: what utility is predicted by my model of f
 - What the probability of improvement? $\alpha_{\text{PI}}(\text{molecule}) = \mathbb{E}_f[\mathbb{1}_{(f > f^*)}]$
 - How much improvement do we expect? $\alpha_{\text{EI}}(\text{molecule}) = \mathbb{E}_f[\max(f - f^*, 0)]$

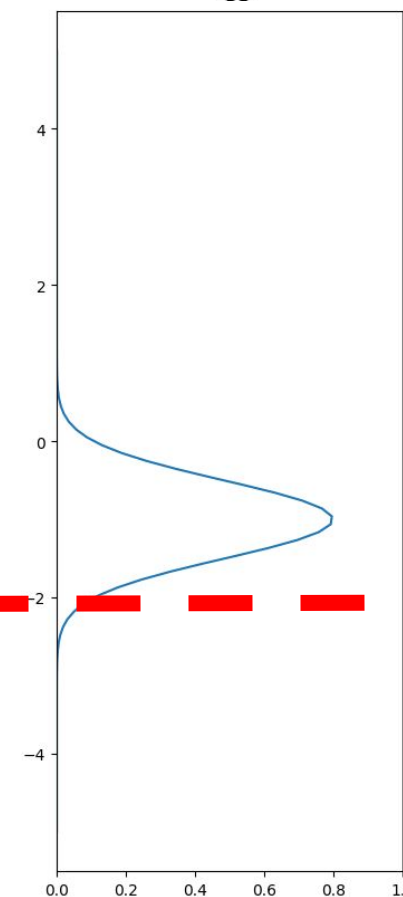
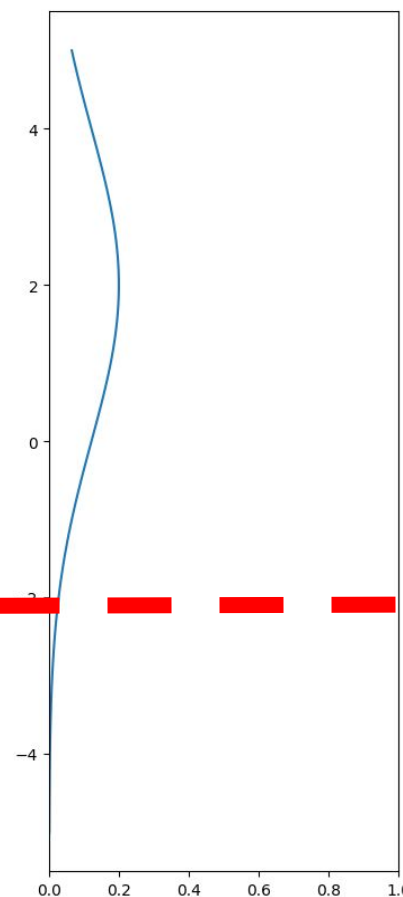
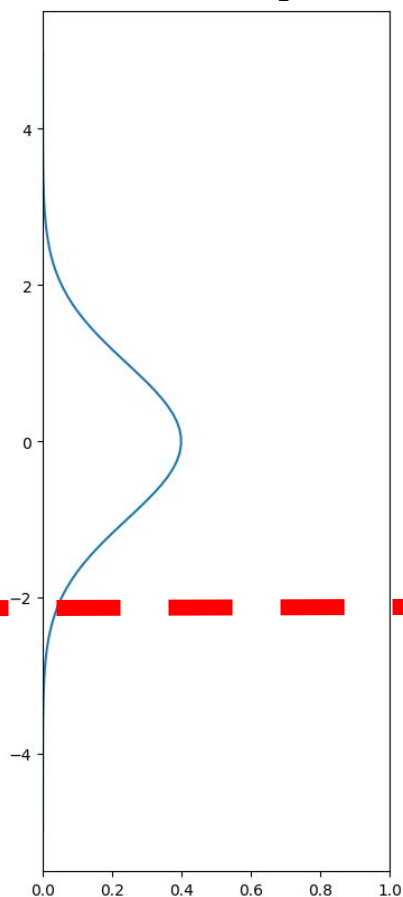
$$f \sim \mathcal{N}(\mu, \sigma^2)$$

Automatically choosing next molecules

Using GP posteriors

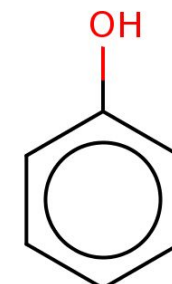
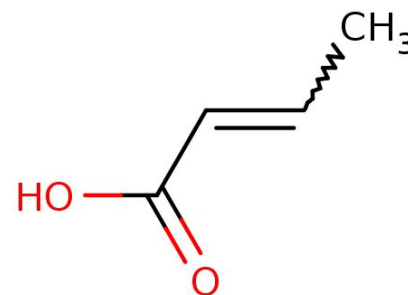
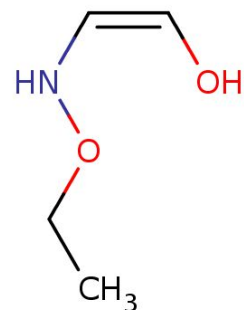
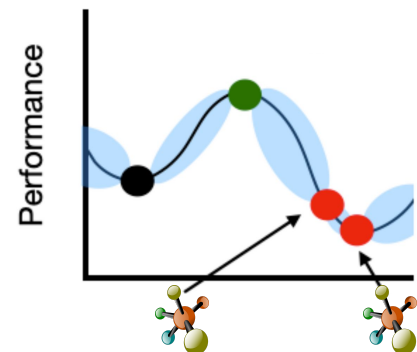


f^*

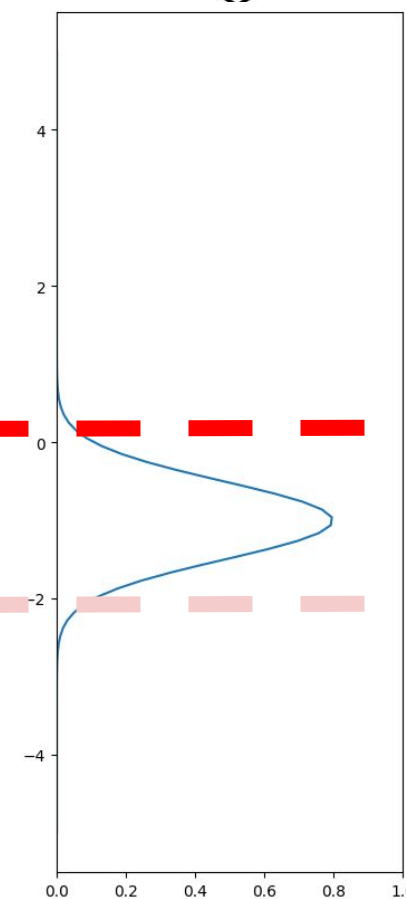
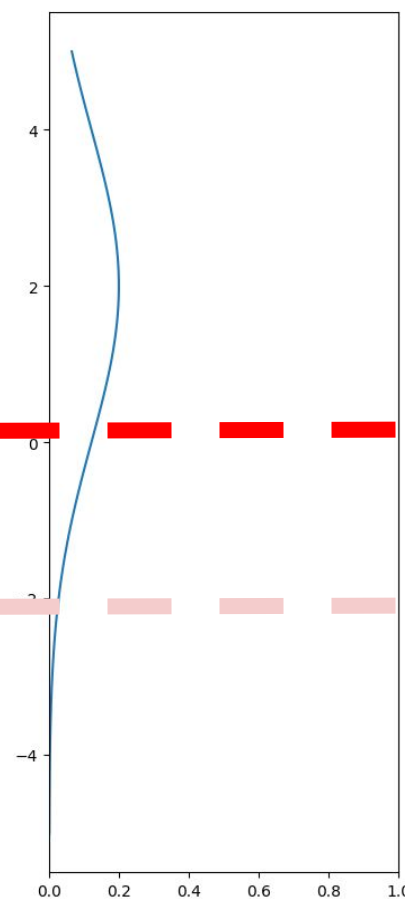
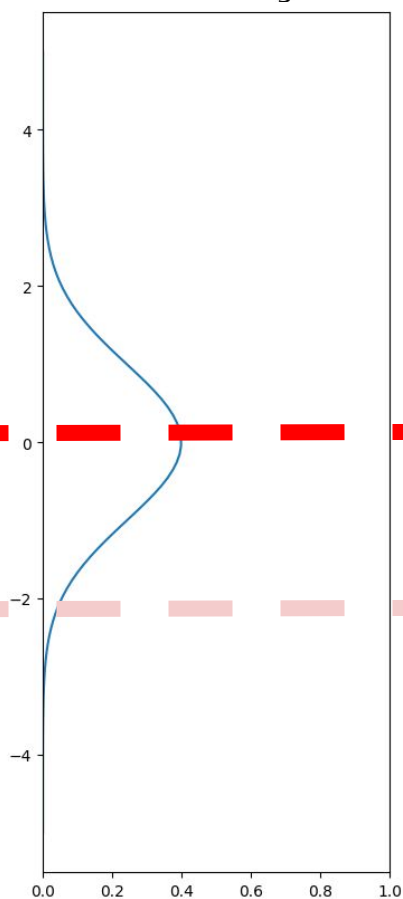


Automatically choosing next molecules

Using GP posteriors



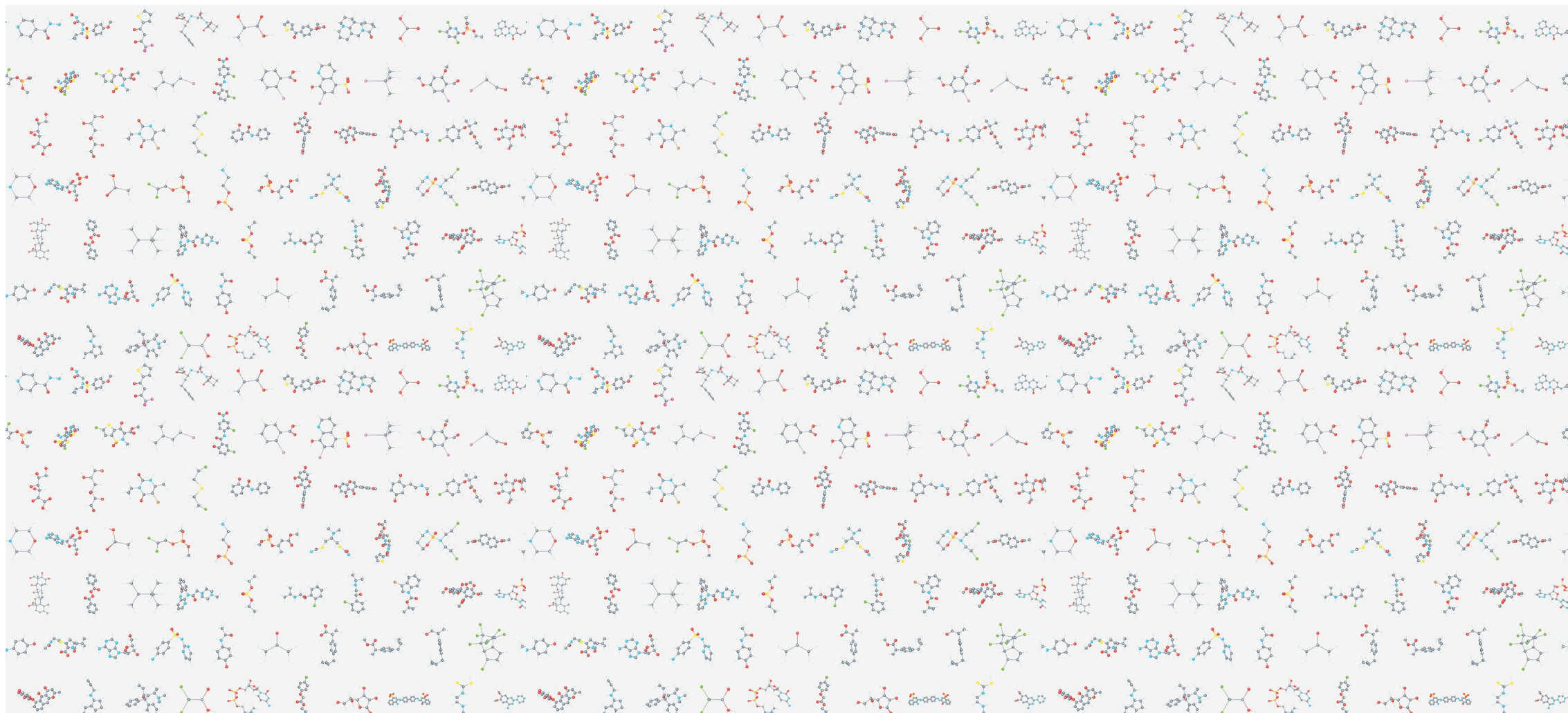
f^*





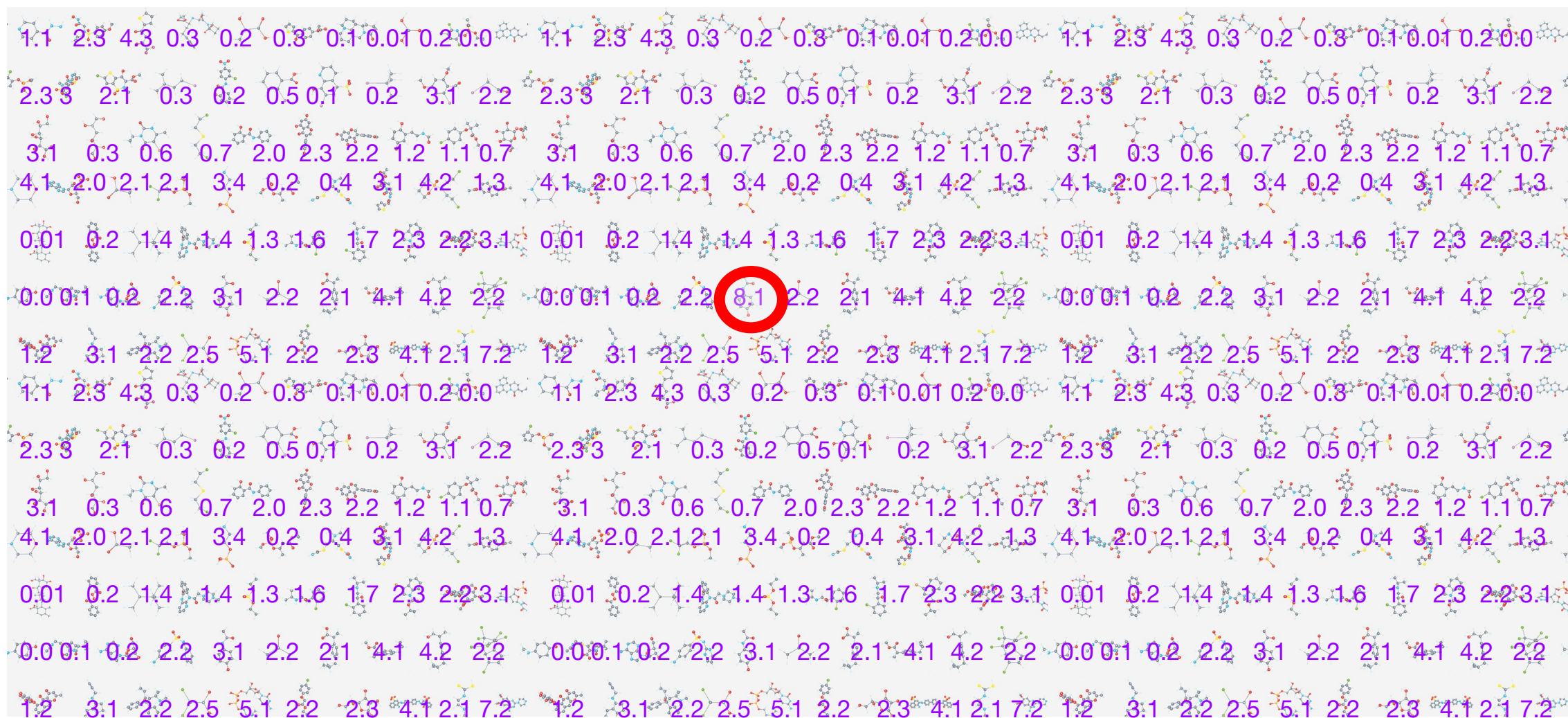
Automatically choosing next molecules

Calc acquisition function and pick best



Automatically choosing next molecules

Calc acquisition function and pick **best**



Automatically choosing next molecules

Full Bayesian optimisation loop

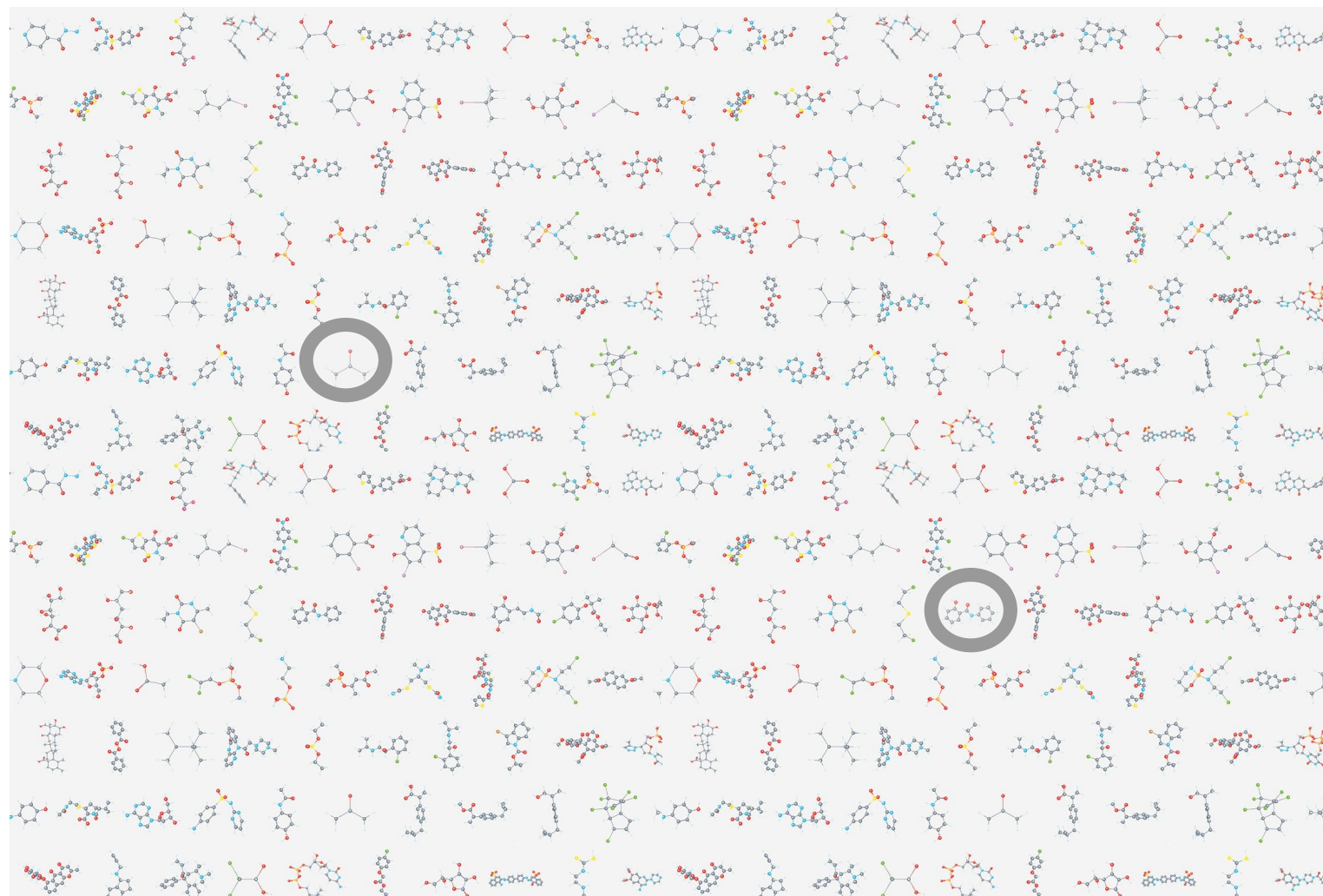
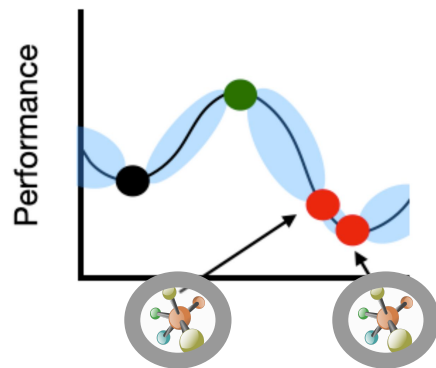
1. Evaluate 2 random molecules



Automatically choosing next molecules

Full Bayesian optimisation loop

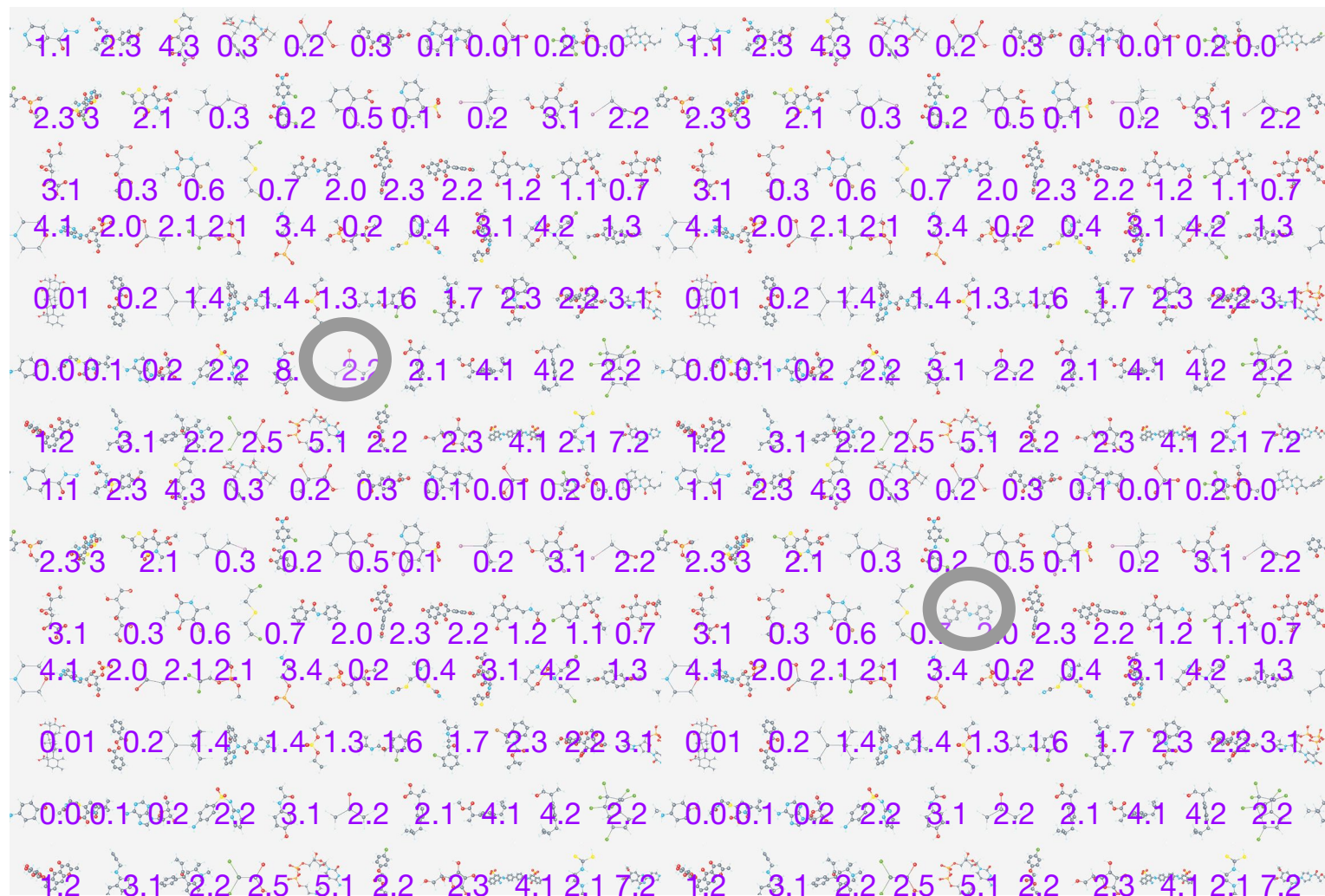
1. Evaluate 2 random molecules
2. Fit GP model to measurements



Automatically choosing next molecules

Full Bayesian optimisation loop

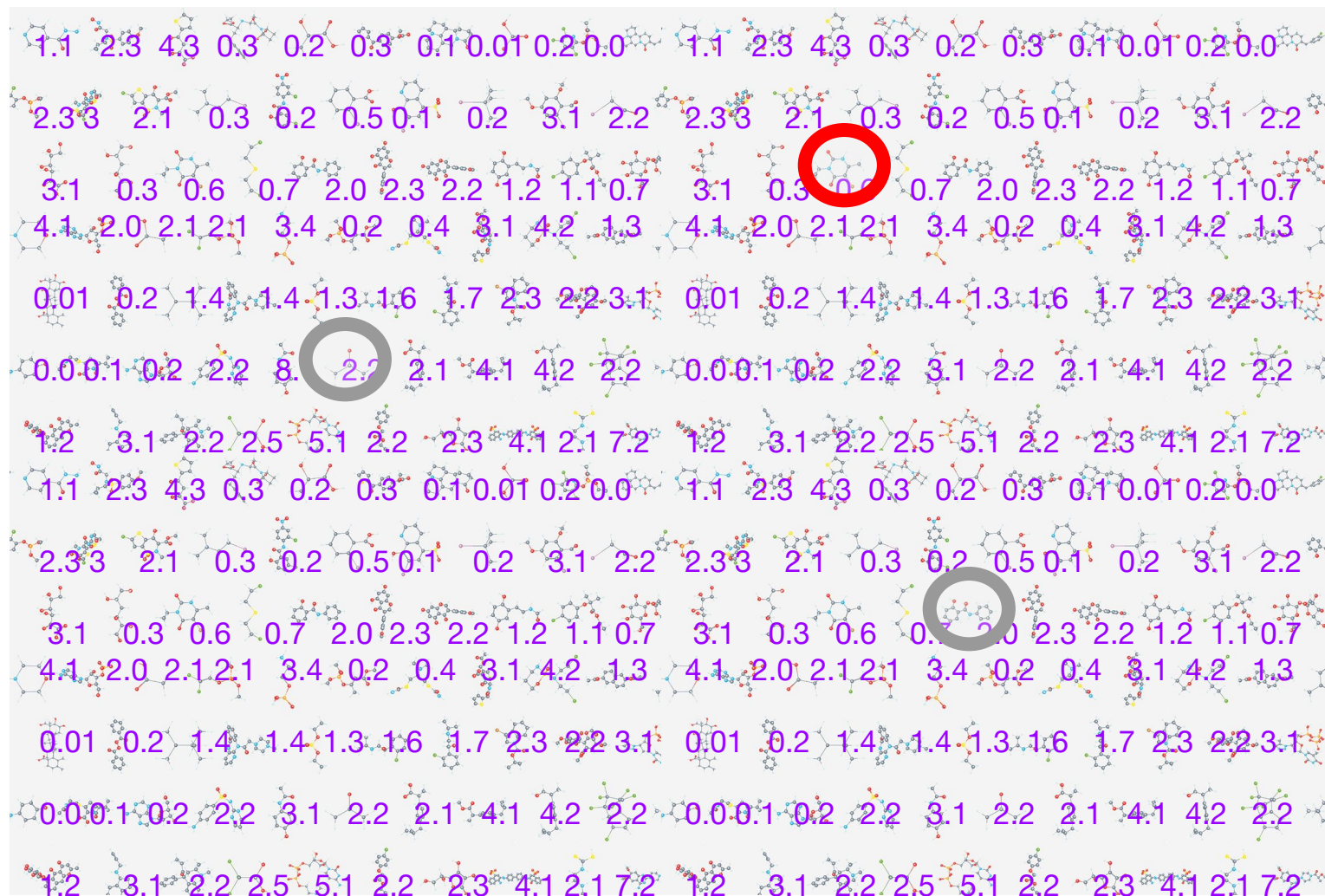
1. Evaluate 2 random molecules
2. Fit GP model to measurements
3. Calc acquisition function



Automatically choosing next molecules

Full Bayesian optimisation loop

1. Evaluate 2 random molecules
2. Fit GP model to measurements
3. Calc acquisition function
4. Choose **new molecule**



Automatically choosing next molecules

Full Bayesian optimisation loop

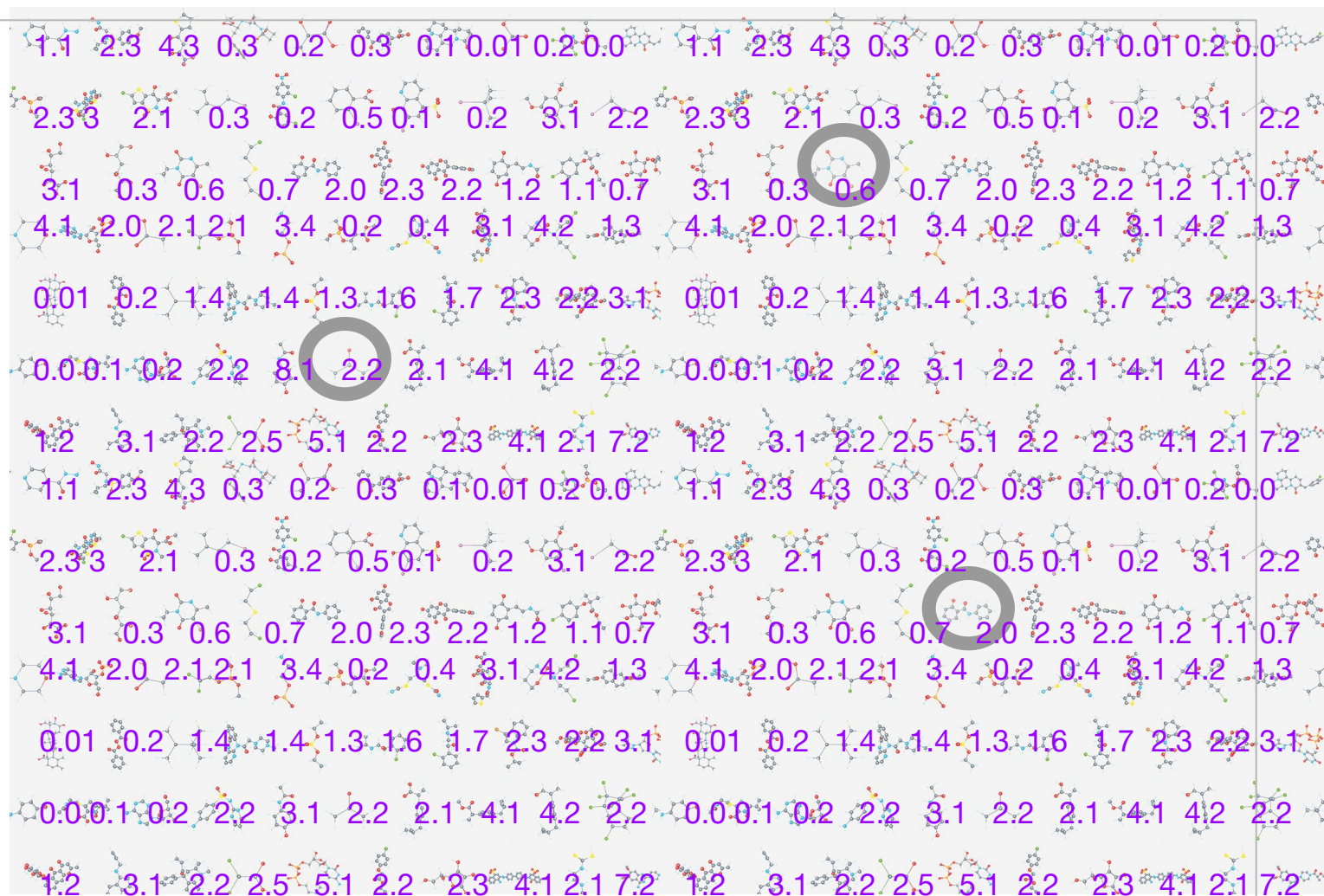
1. Evaluate 2 random molecules
2. Fit GP model to measurements
3. Calc acquisition function
4. Choose new molecule
5. Go to step 2.



Automatically choosing next molecules

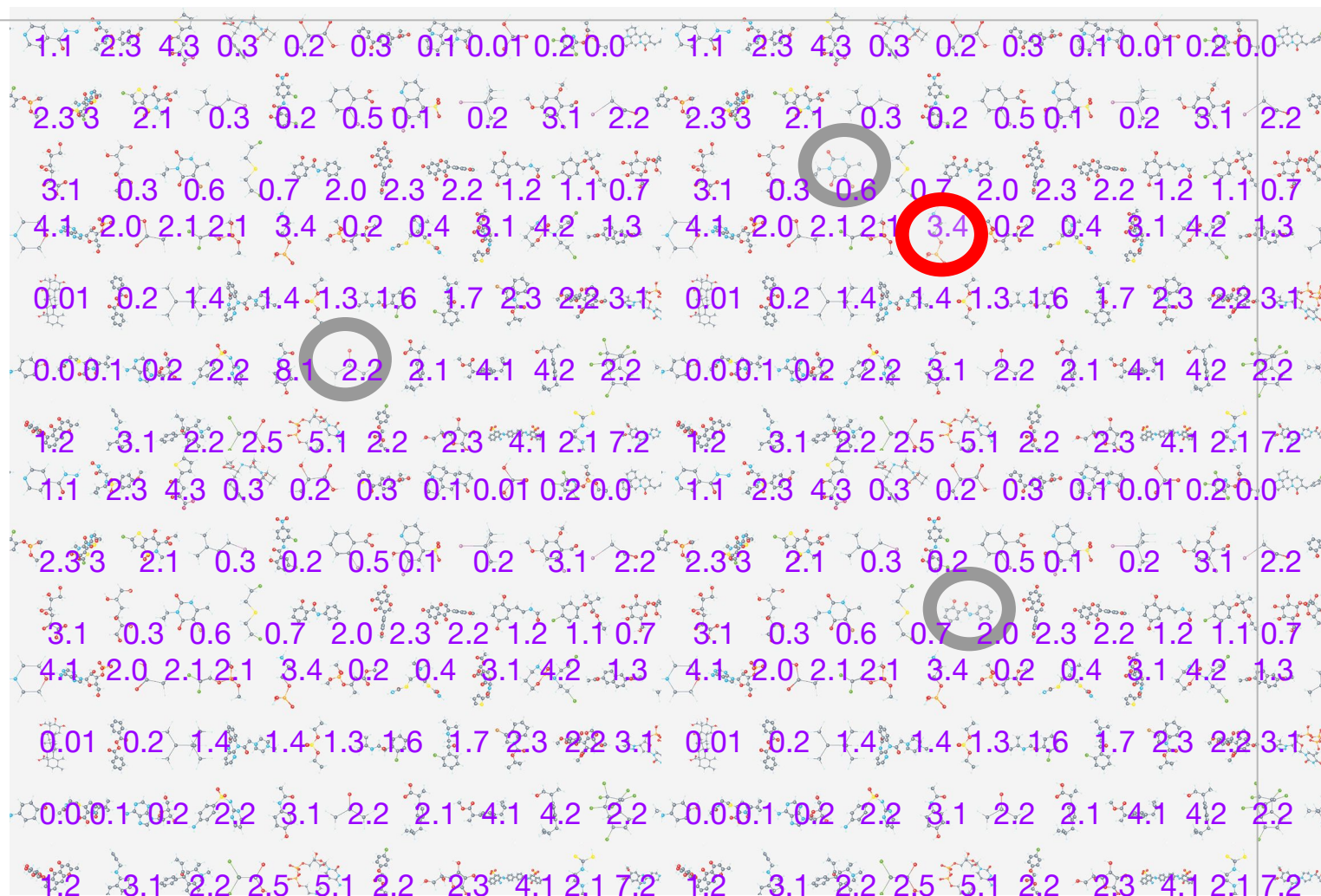
Full Bayesian optimisation loop

1. Evaluate 2 random molecules
2. Fit GP model to measurements
3. Calc new acquisition function
4. Choose new molecule
5. Go to step 2.



Full Bayesian optimisation loop

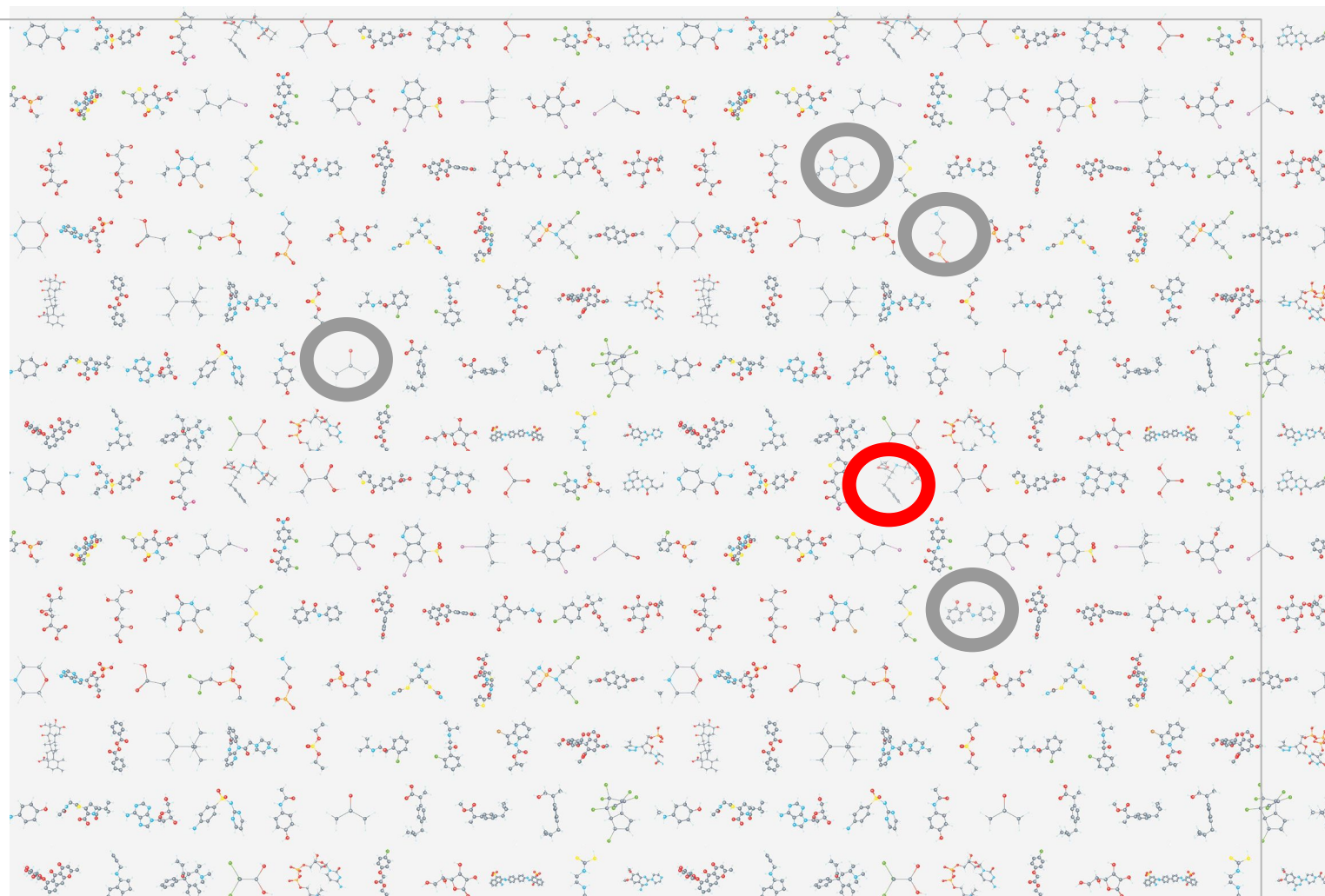
1. Evaluate 2 random molecules
2. Fit GP model to measurements
3. Calc new acquisition function
4. Choose new molecule
5. Go to step 2.



Automatically choosing next molecules

Full Bayesian optimisation loop

1. Evaluate 2 random molecules
2. Fit GP model to measurements
3. Calc new acquisition function
4. Choose new molecule
5. Go to step 2.



Automatically choosing next molecules

Full Bayesian optimisation loop

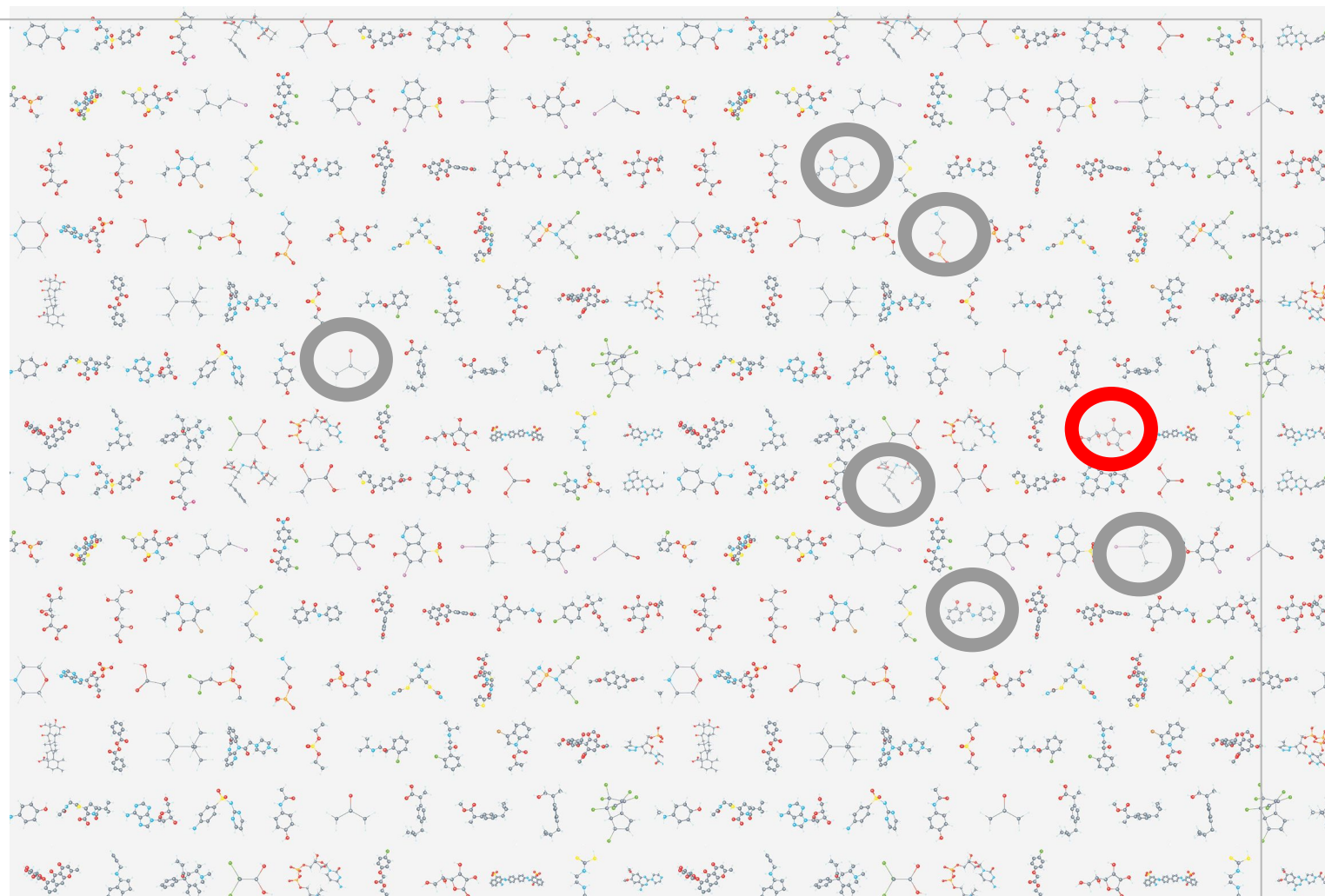
1. Evaluate 2 random molecules
2. Fit GP model to measurements
3. Calc new acquisition function
4. Choose new molecule
5. Go to step 2.



Automatically choosing next molecules

Full Bayesian optimisation loop

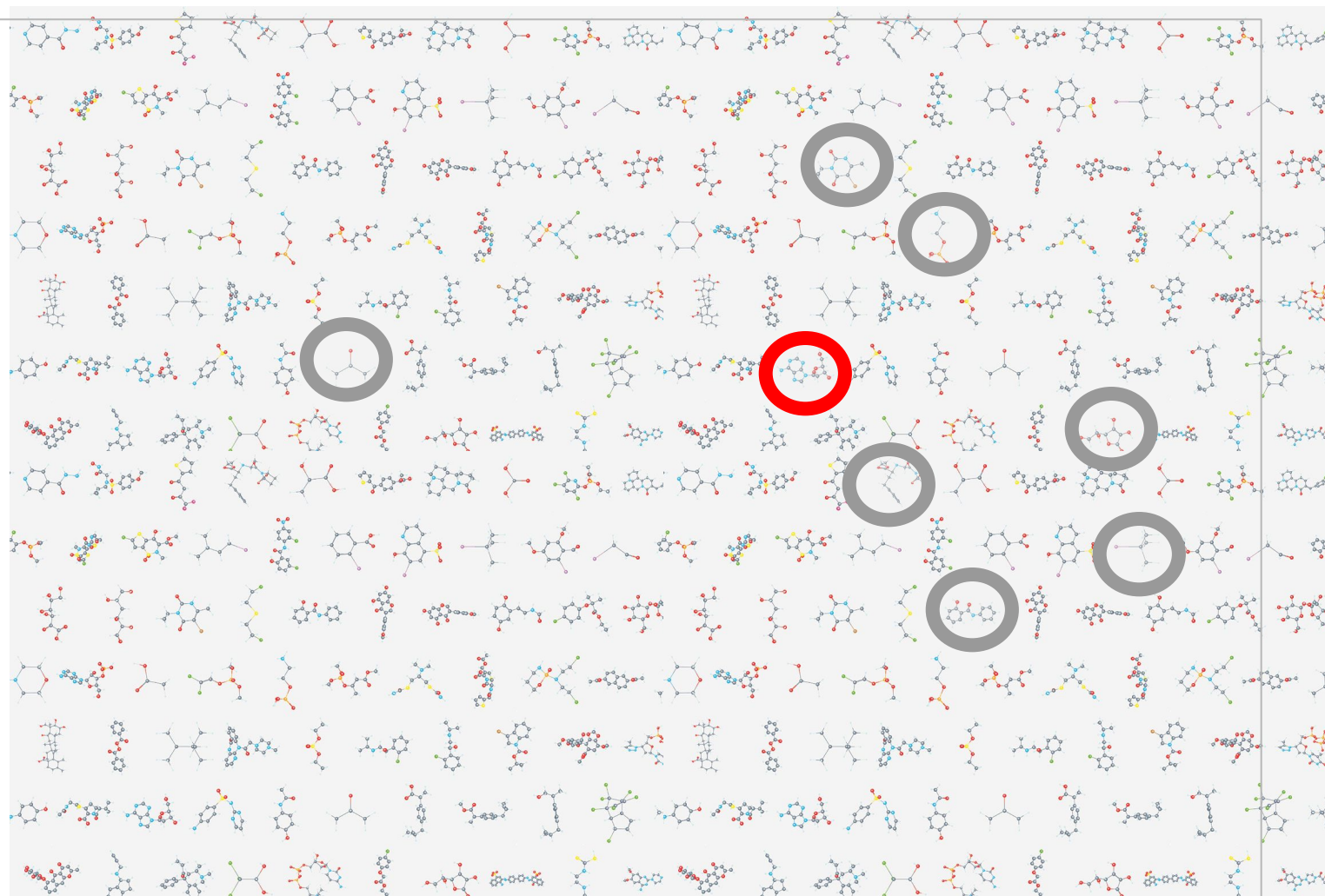
1. Evaluate 2 random molecules
2. Fit GP model to measurements
3. Calc new acquisition function
4. Choose new molecule
5. Go to step 2.



Automatically choosing next molecules

Full Bayesian optimisation loop

1. Evaluate 2 random molecules
2. Fit GP model to measurements
3. Calc new acquisition function
4. Choose new molecule
5. Go to step 2.

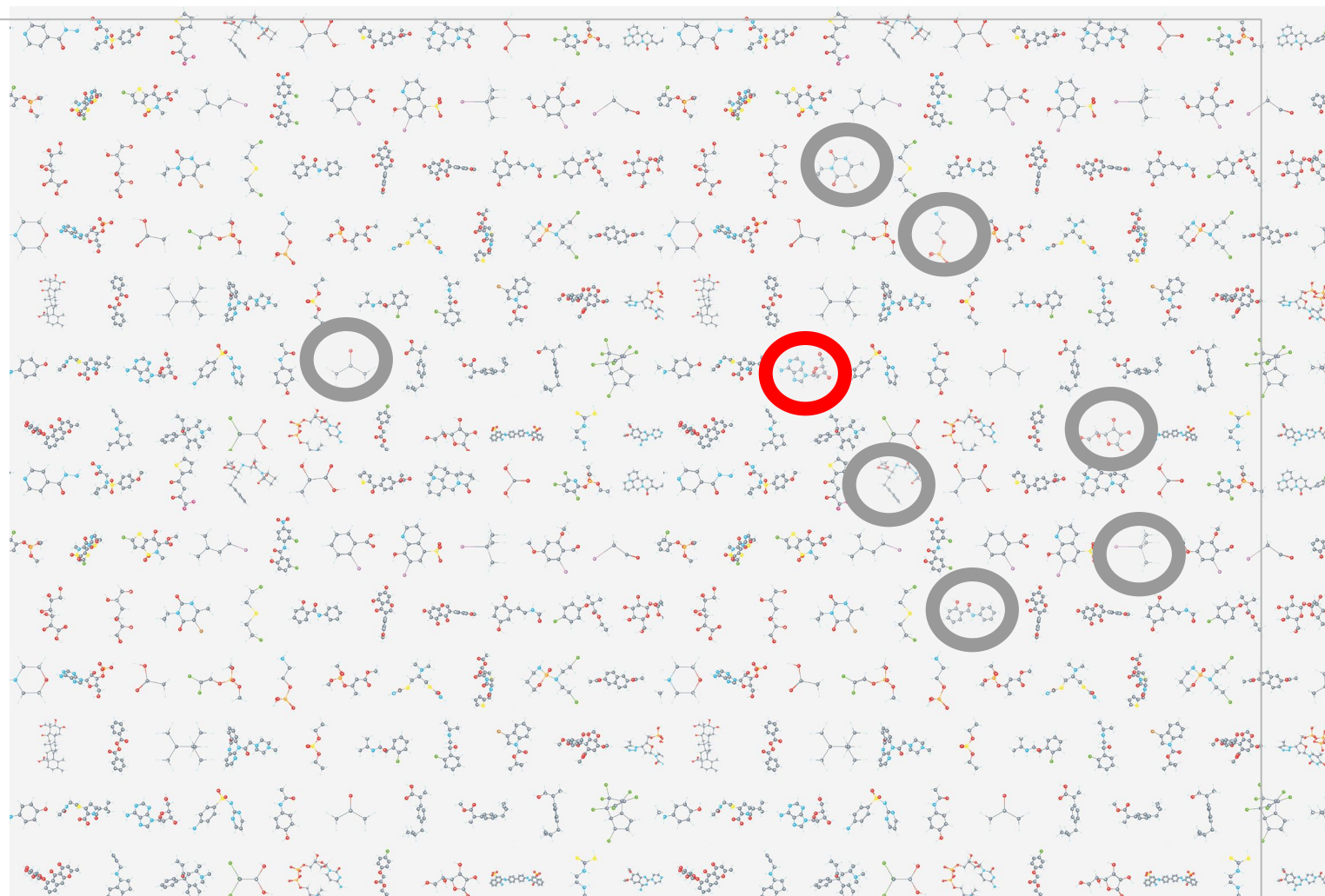


Automatically choosing next molecules

Full Bayesian optimisation loop

1. Evaluate 2 random molecules
2. Fit GP model to measurements
3. Calc new acquisition function
4. Choose new molecule
5. Go to step 2.

And so on





UNIVERSITY OF
CAMBRIDGE

Lancaster
University



What about standard optimisation problems?

i.e. infinite candidates

BO Demo

Let's find the maximum of a 1D function:

BO Demo

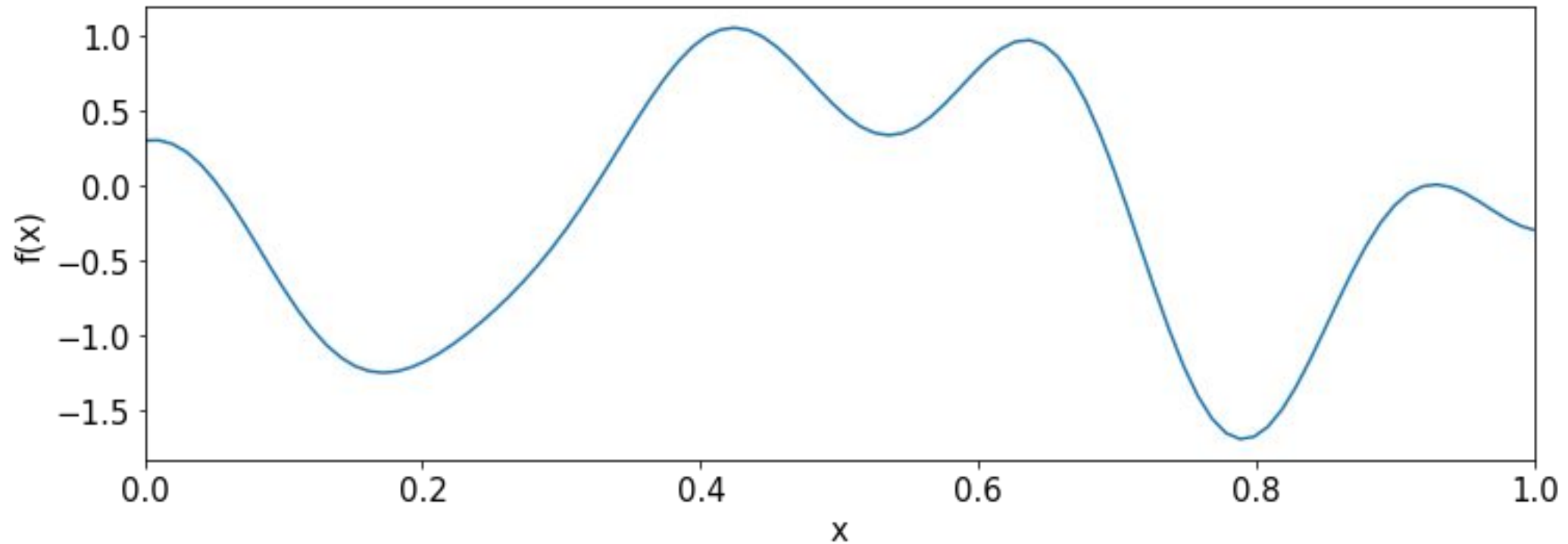
Let's find the maximum of a 1D function:

Using as **few** function evaluations as possible!

BO Demo

Let's find the maximum of a 1D function:

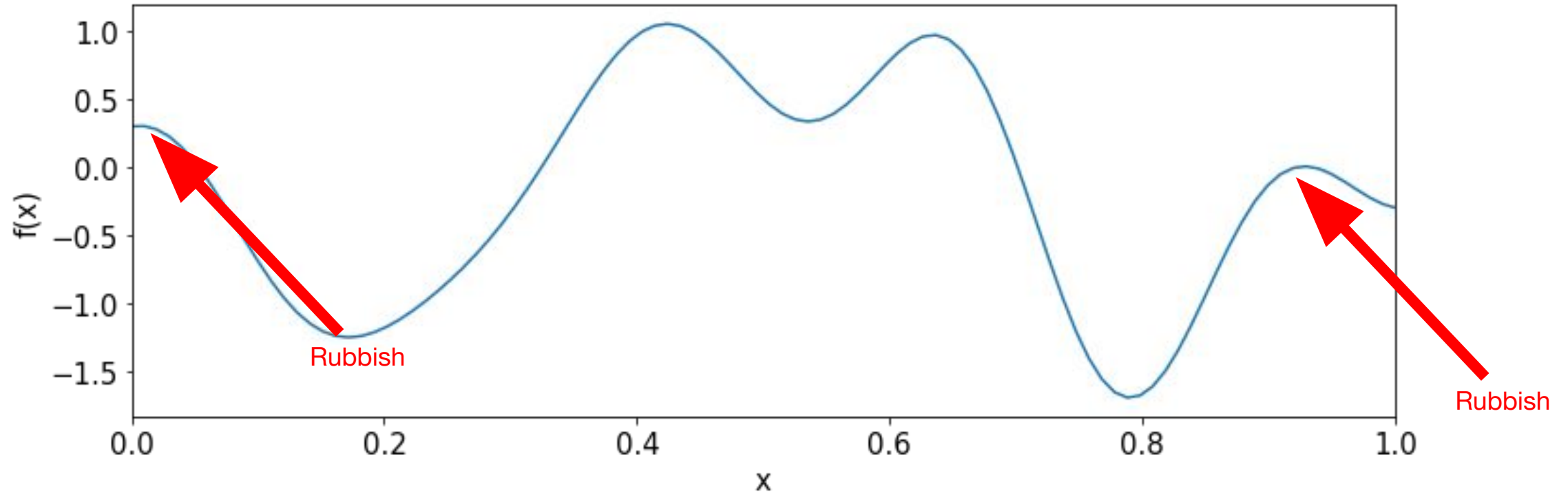
Using as **few** function evaluations as possible!



BO Demo

Let's find the maximum of a 1D function:

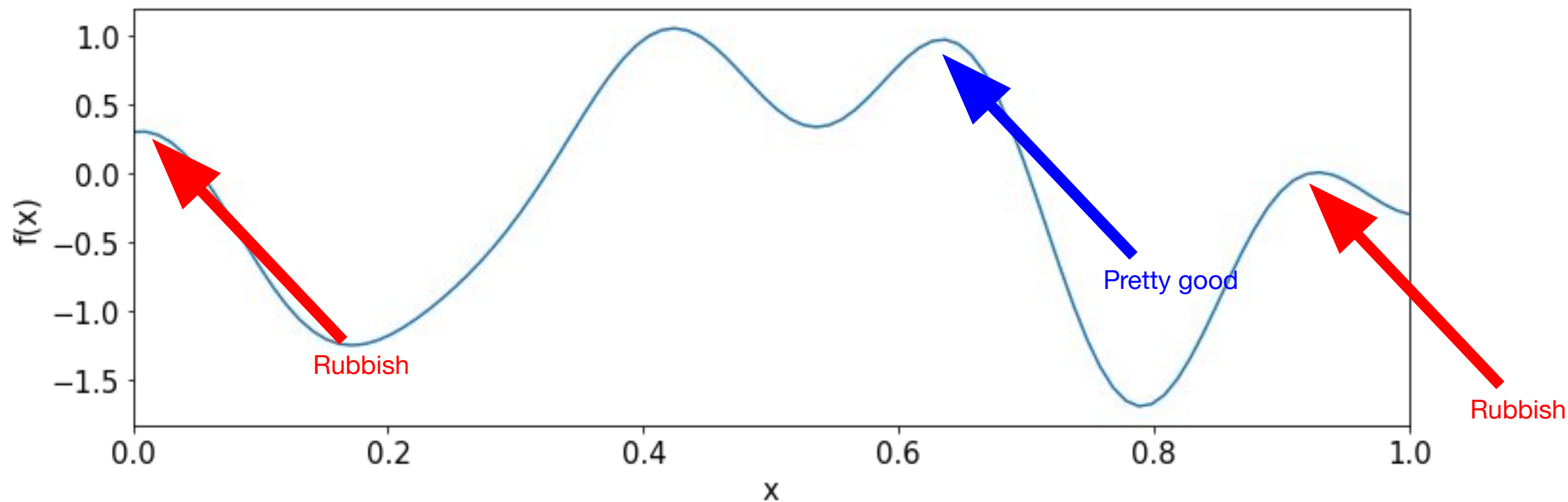
Using as **few** function evaluations as possible!



BO Demo

Let's find the maximum of a 1D function:

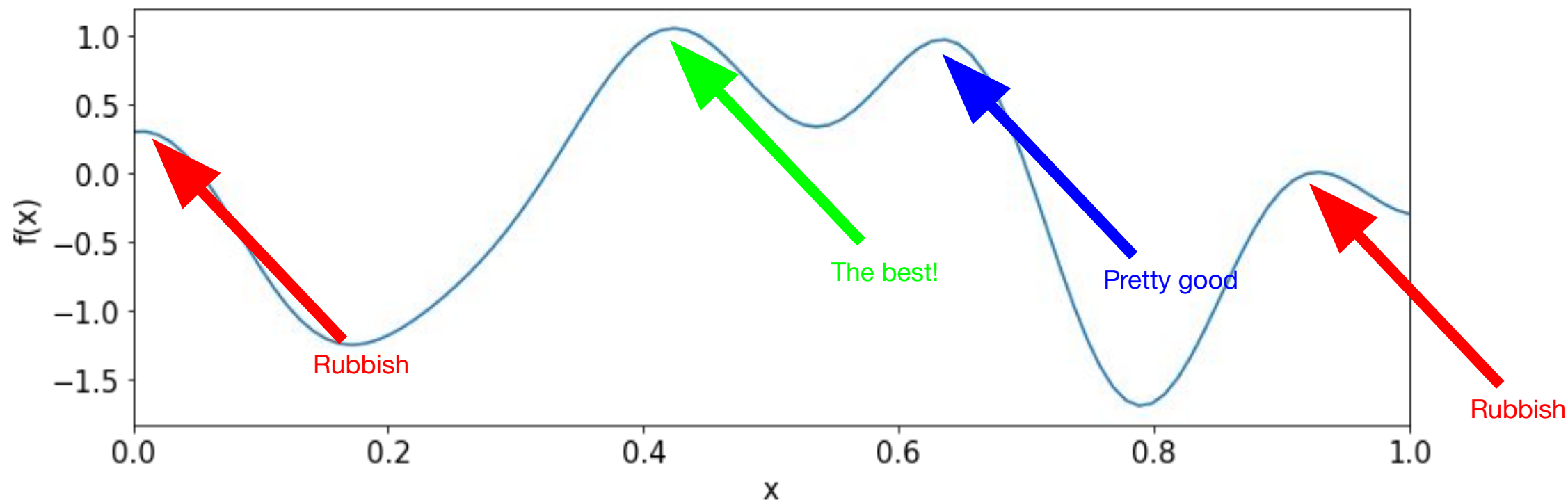
Using as **few** function evaluations as possible!



BO Demo

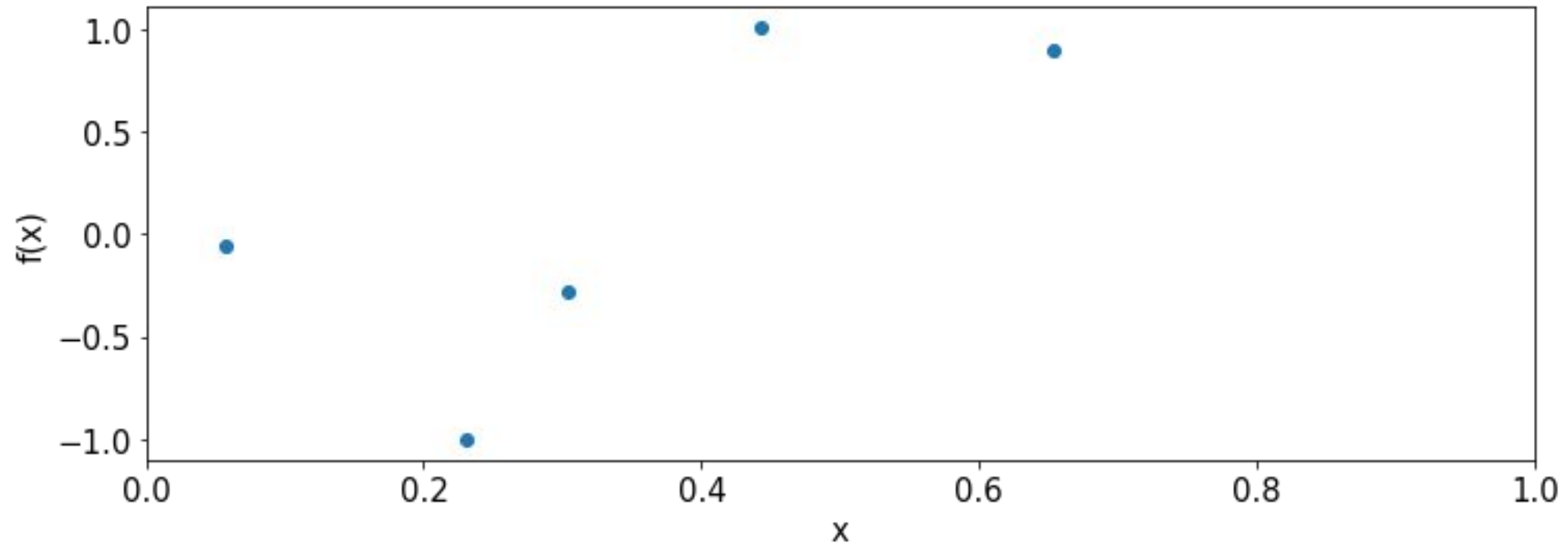
Let's find the maximum of a 1D function:

Using as **few** function evaluations as possible!



BO Demo

Suppose we make 5 evaluations

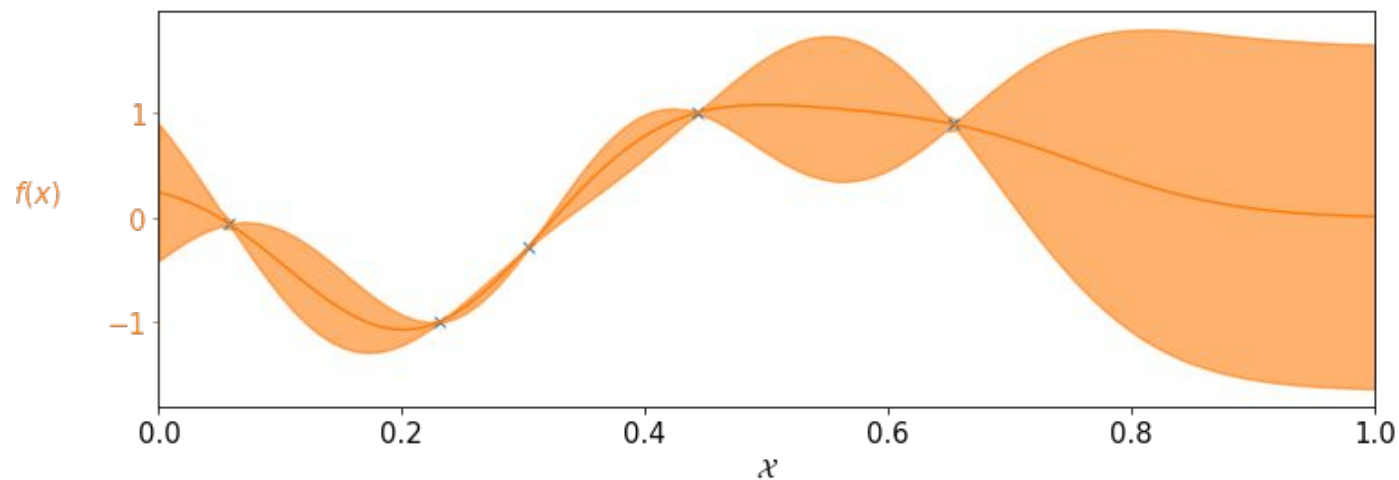
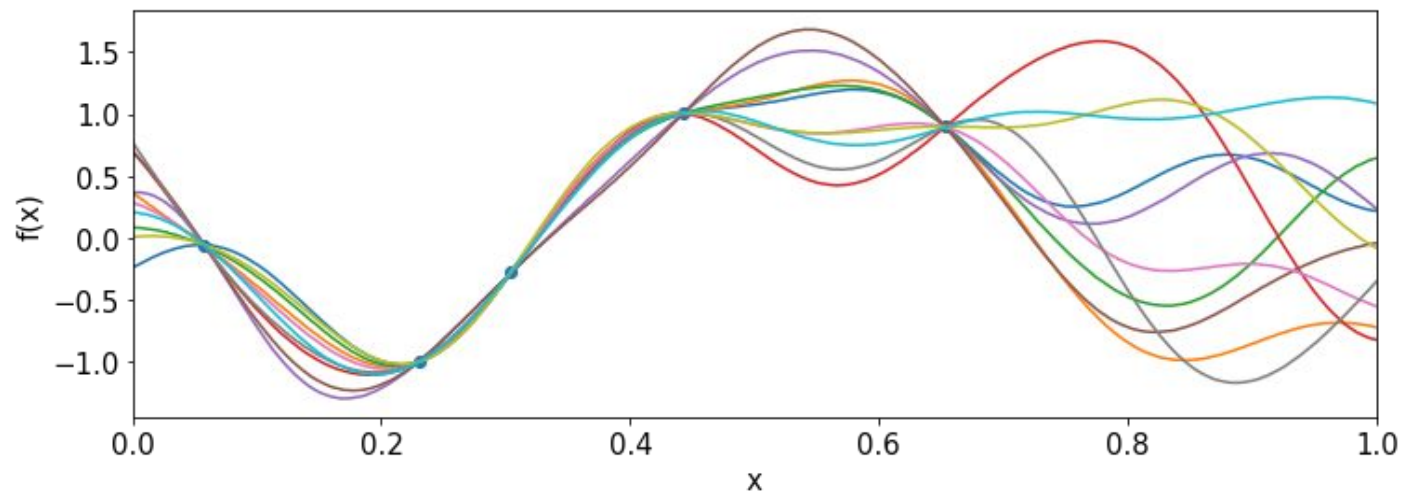


Where should we next evaluate? Explore/Exploit?



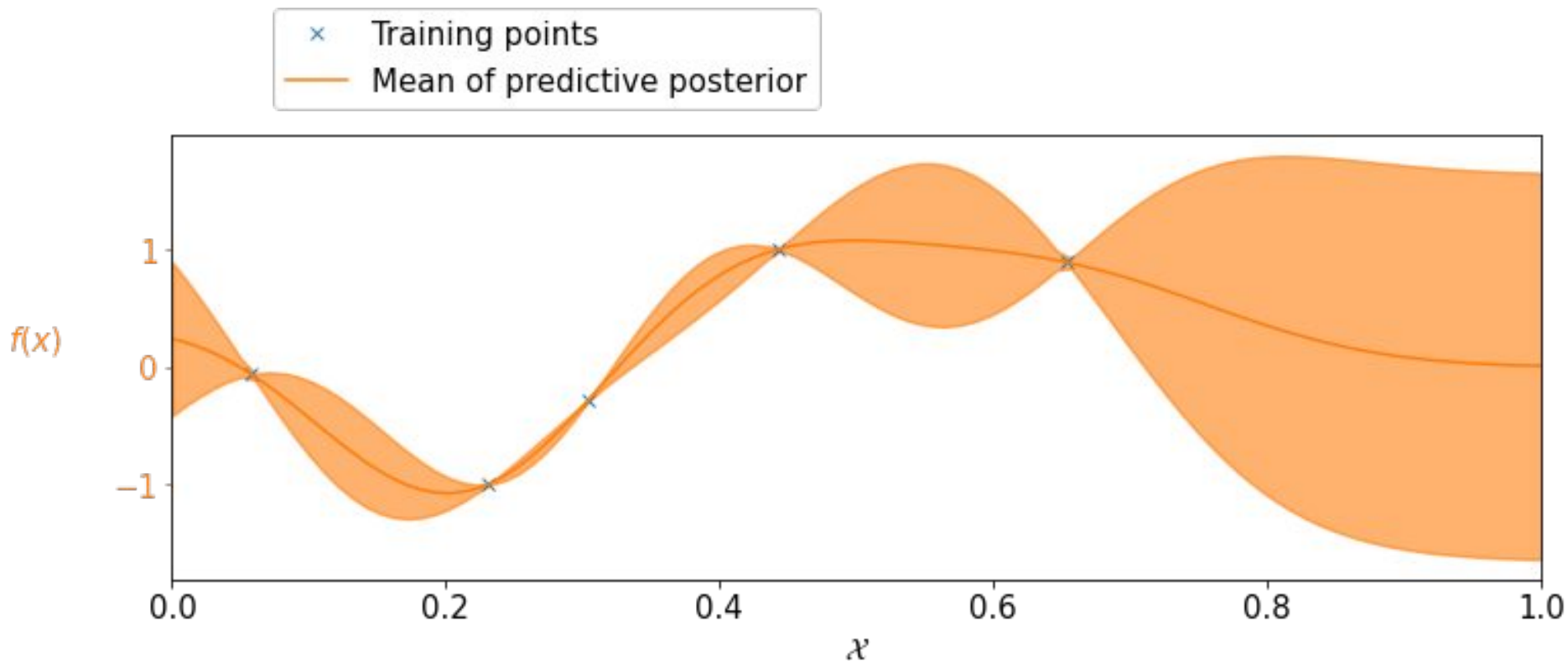
How to automate BO: step 1

Use a statistical model like a Gaussian process



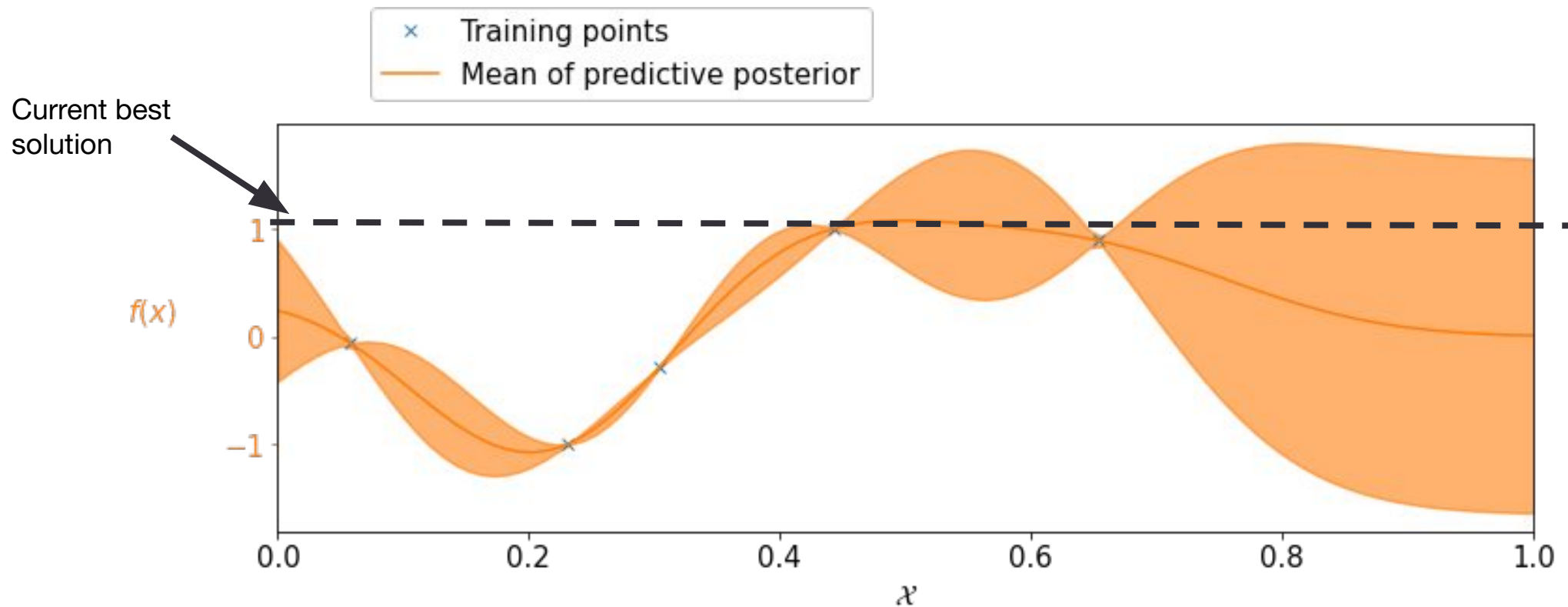
How to automate BO: step 2

Automated decision making via an acquisition function like expected improvement



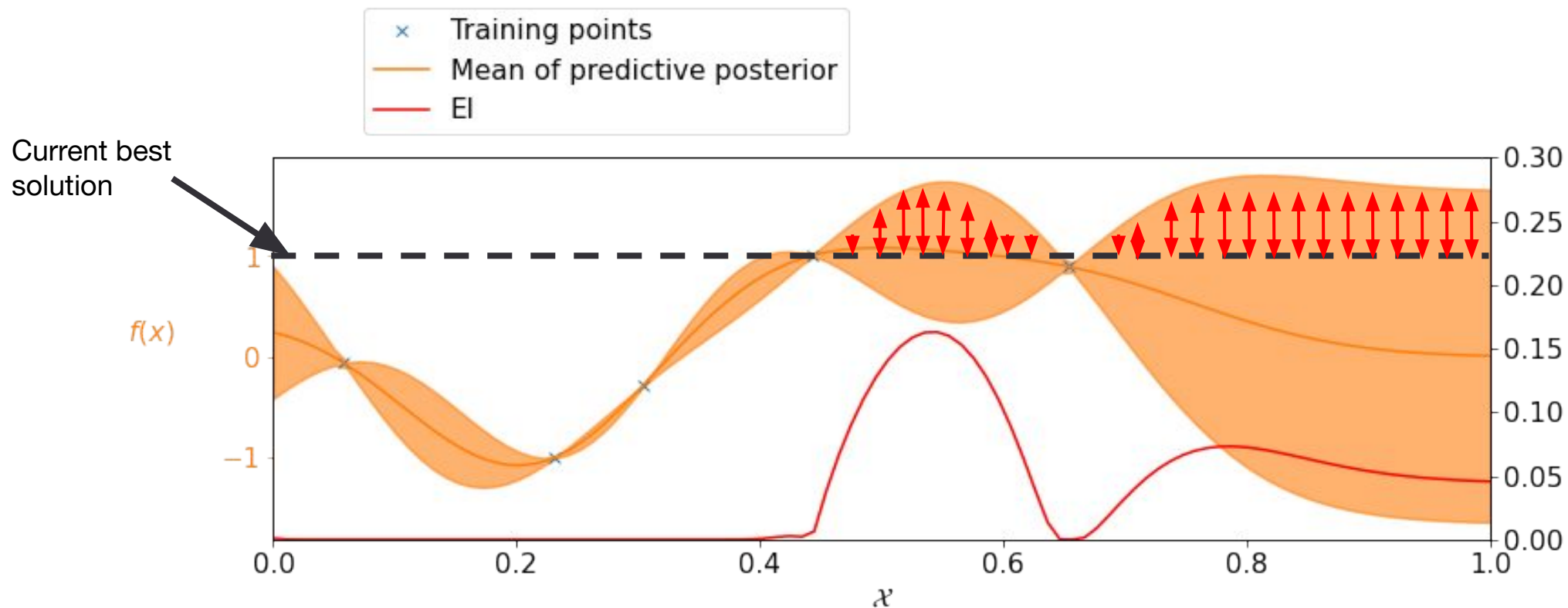
How to automate BO: step 2

Automated decision making via an acquisition function like expected improvement



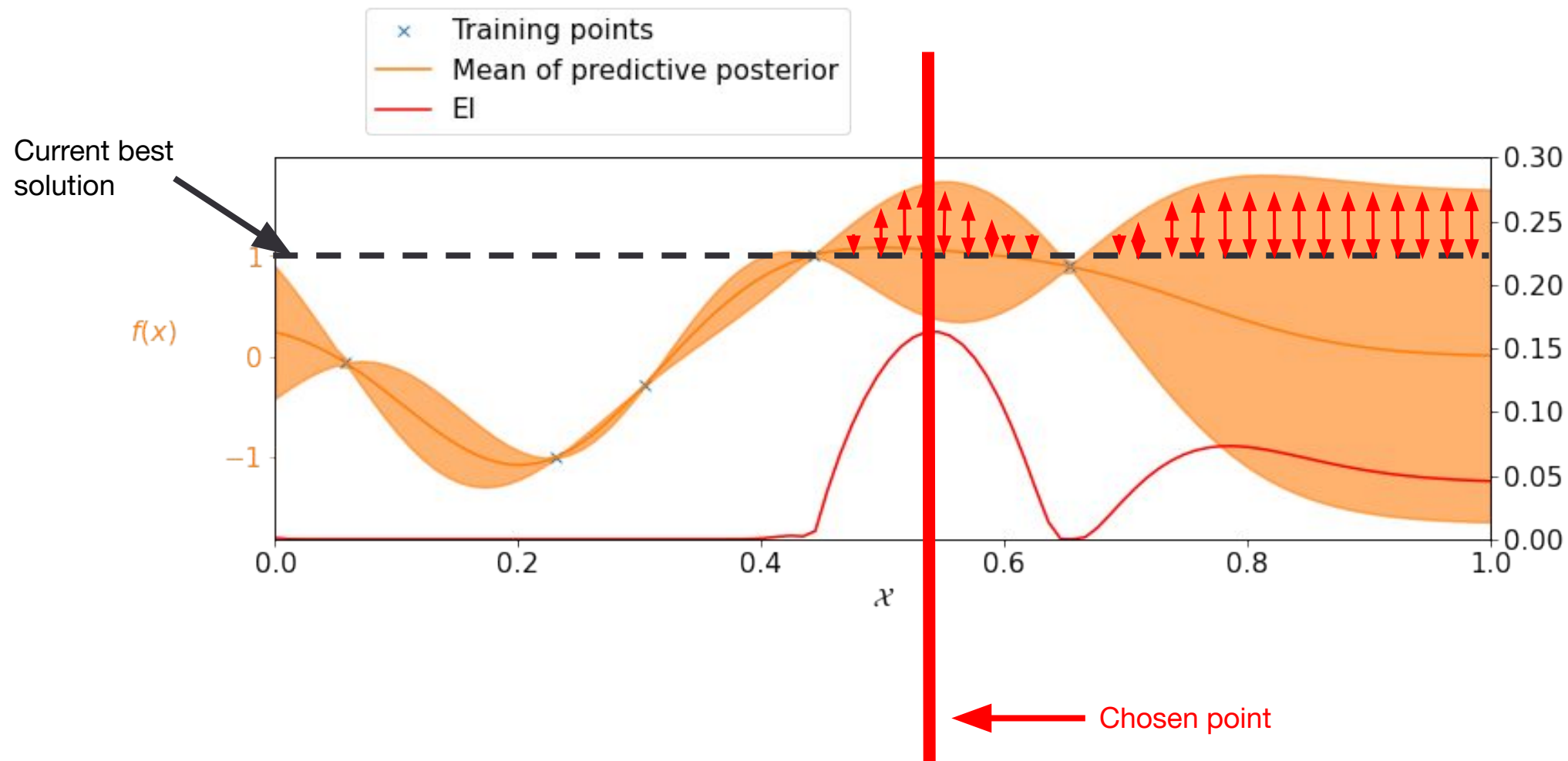
How to automate BO: step 2

Automated decision making via an acquisition function like expected improvement



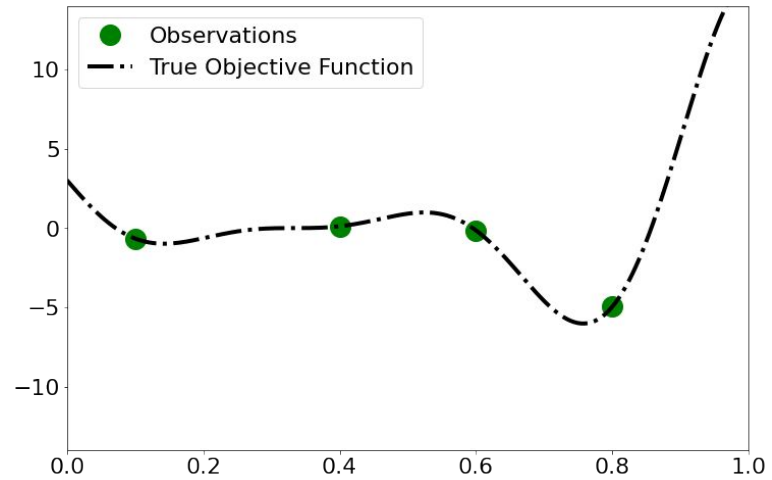
How to automate BO: step 2

Automated decision making via an acquisition function like expected improvement



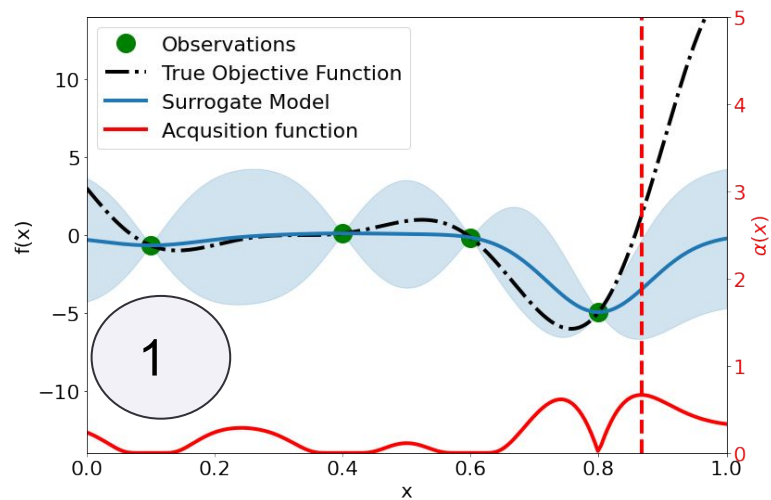
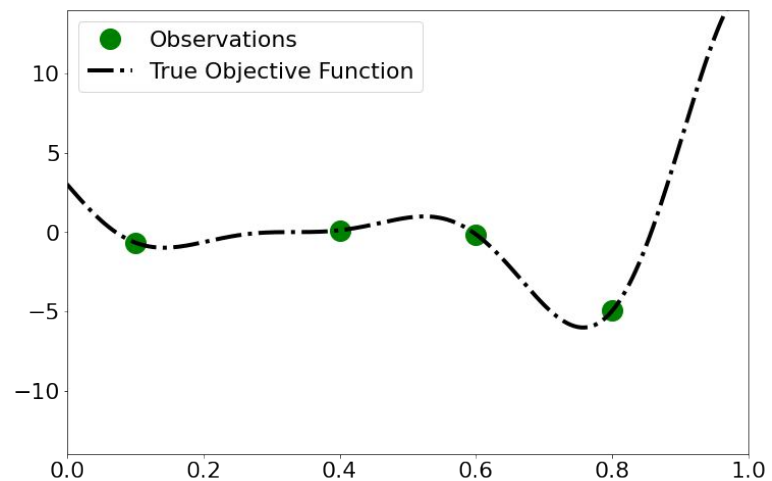
Expected Improvement

Demo BO loop



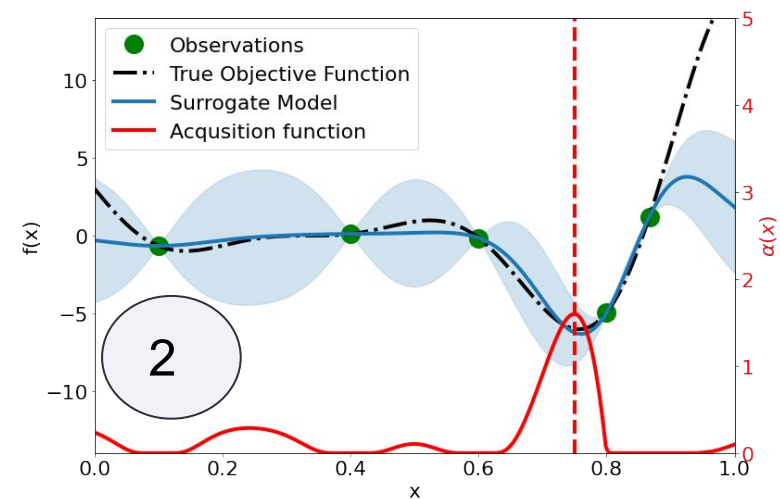
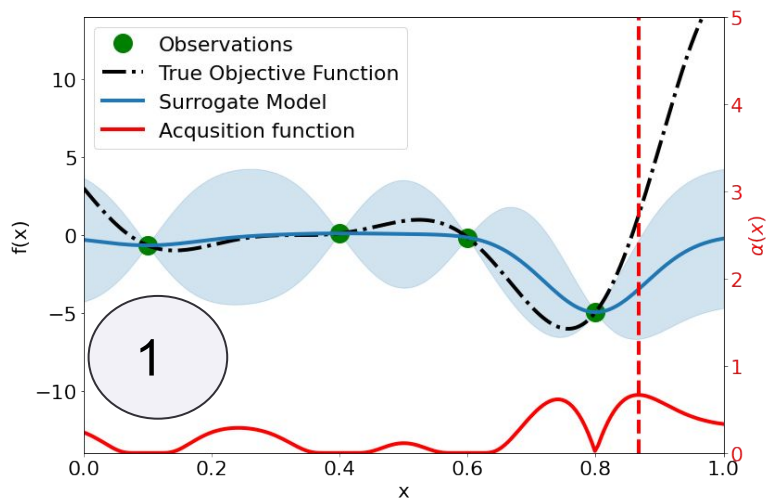
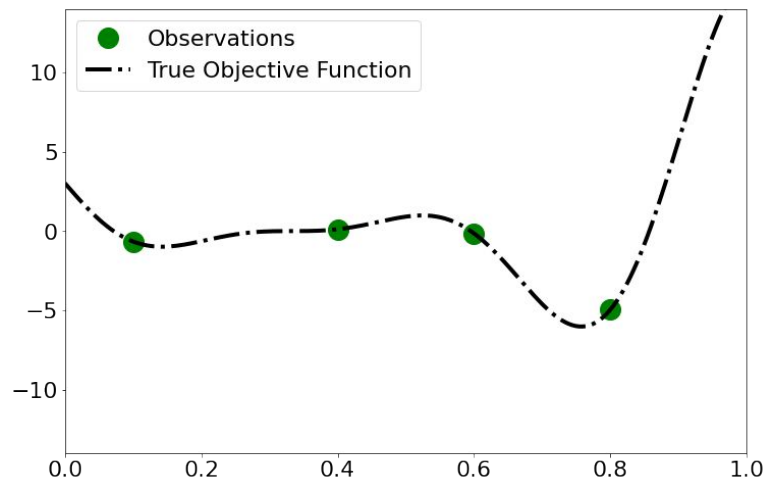
Expected Improvement

Demo BO loop



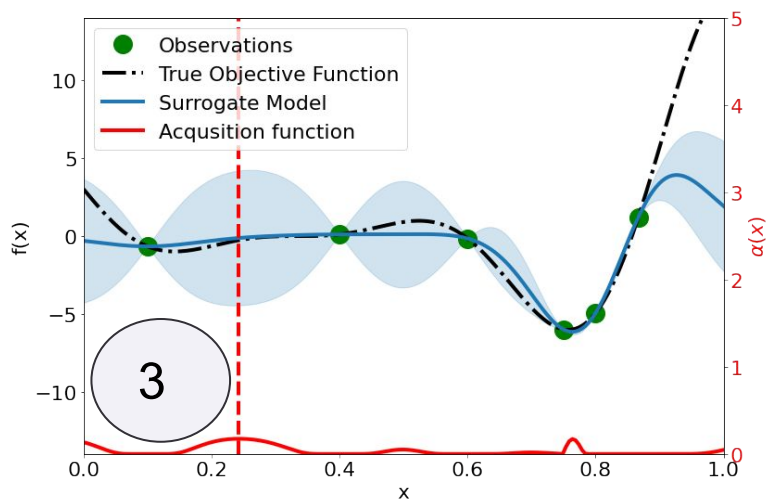
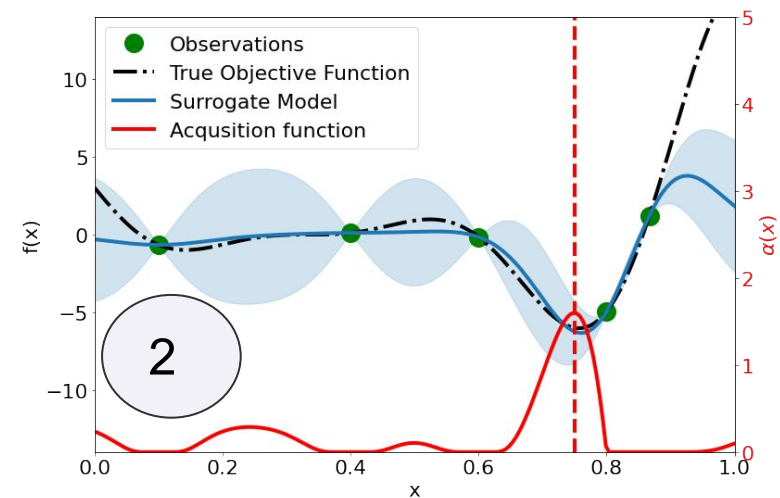
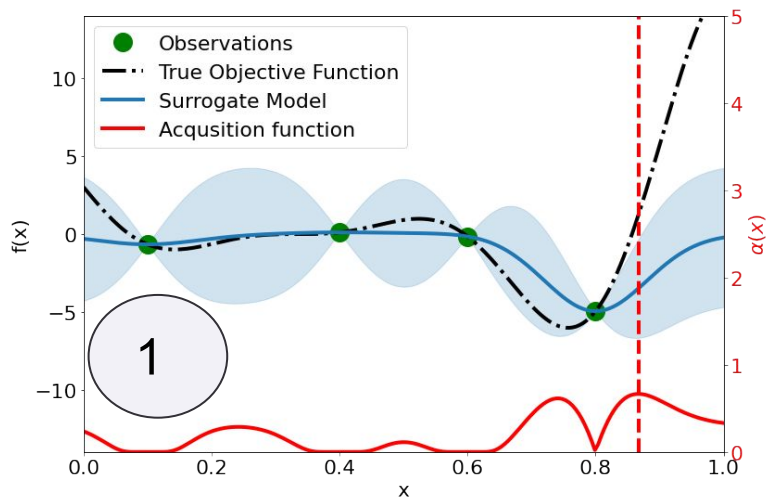
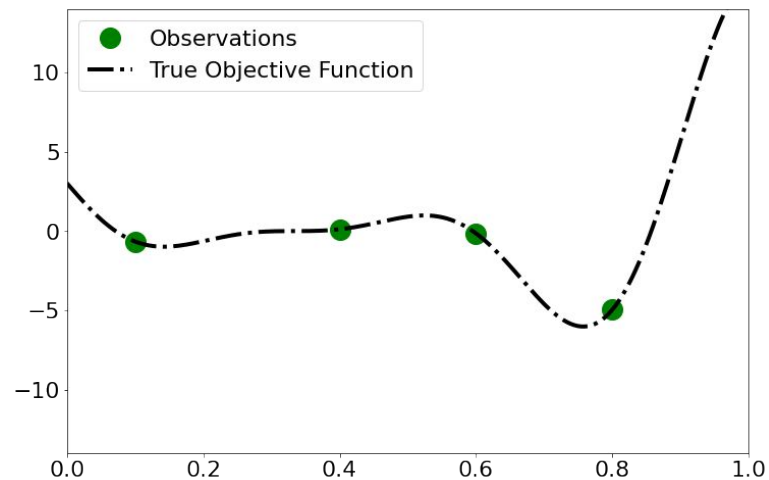
Expected Improvement

Demo BO loop



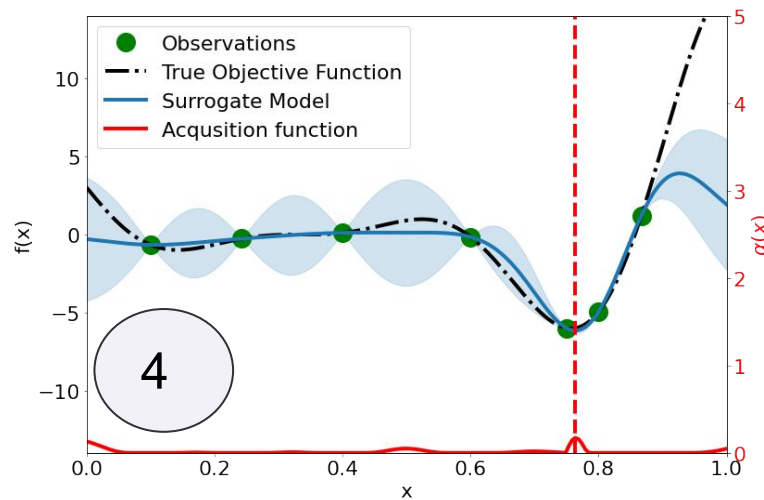
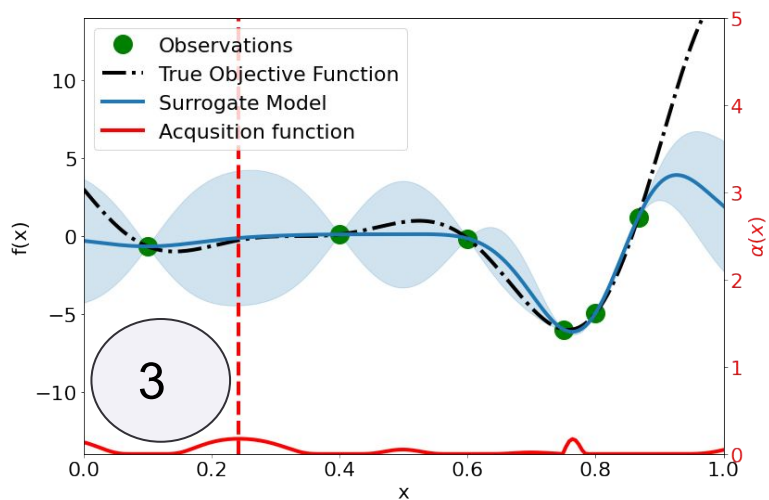
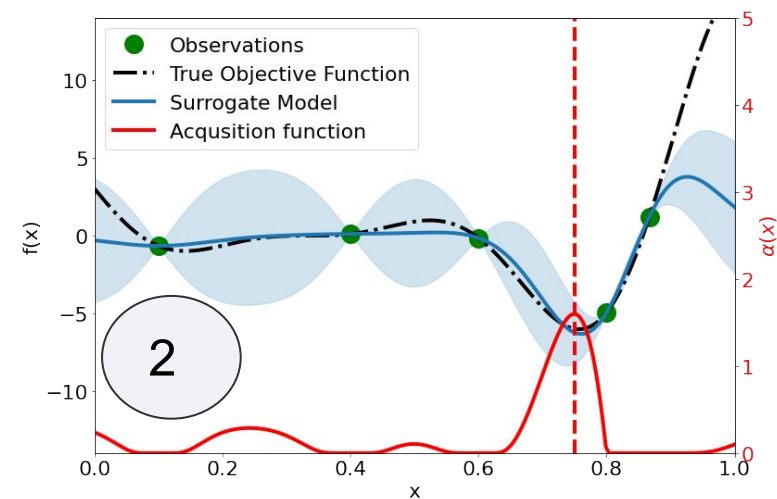
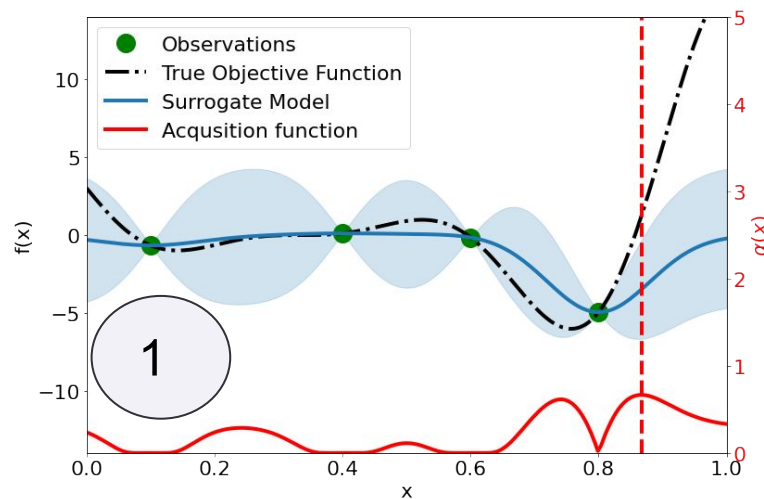
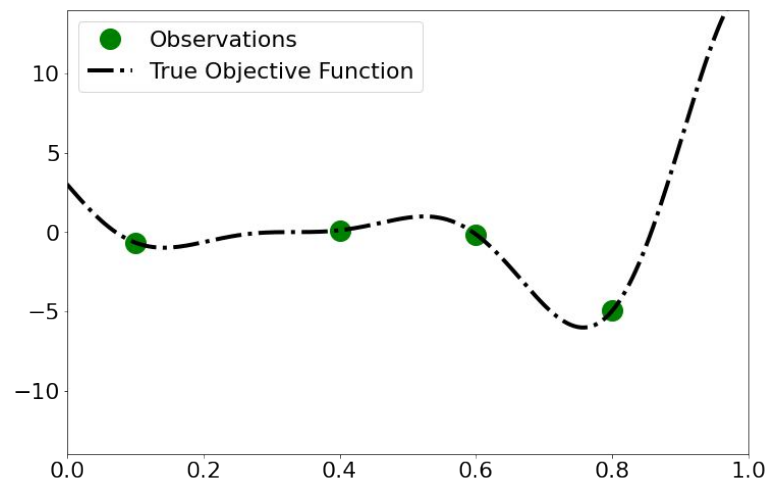
Expected Improvement

Demo BO loop



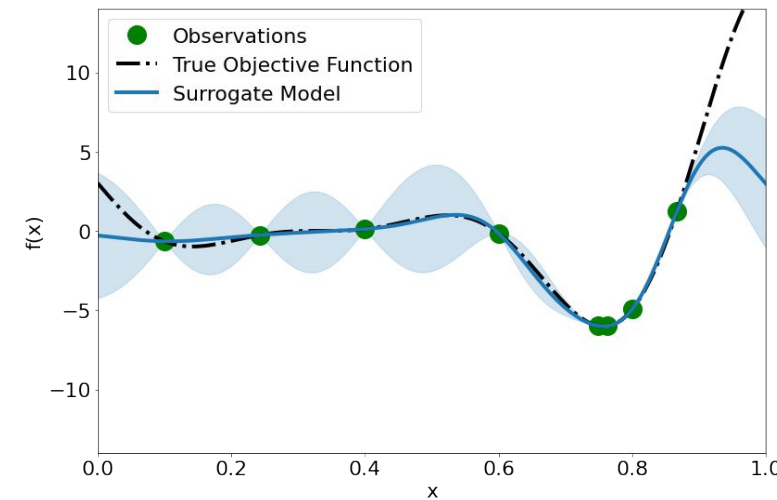
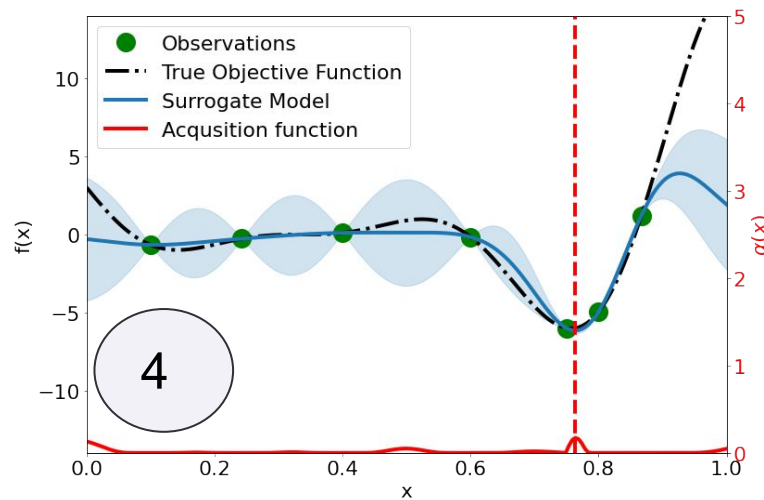
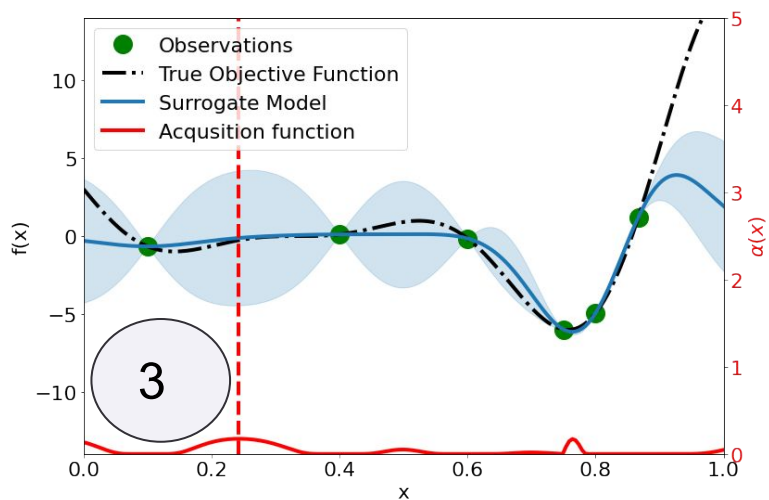
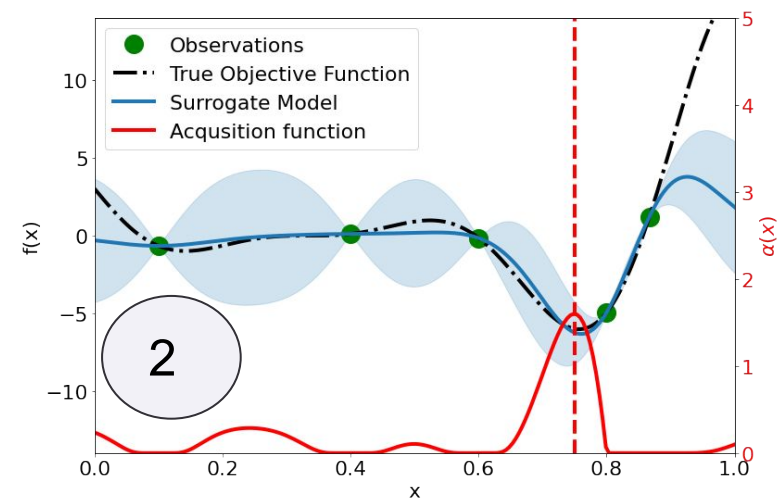
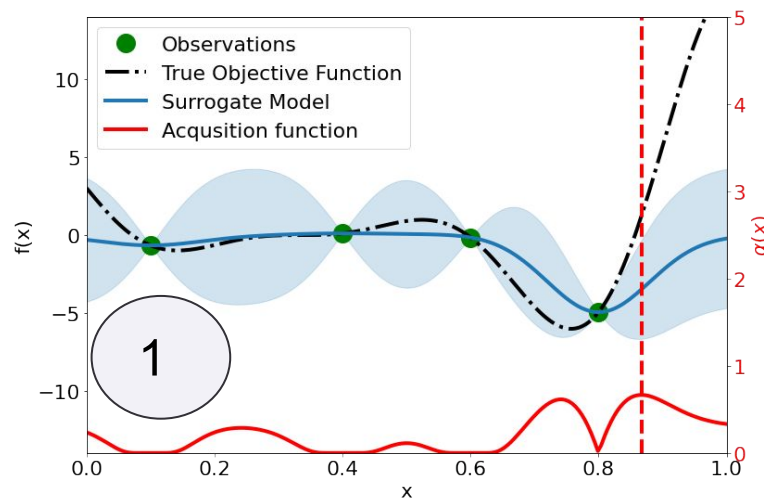
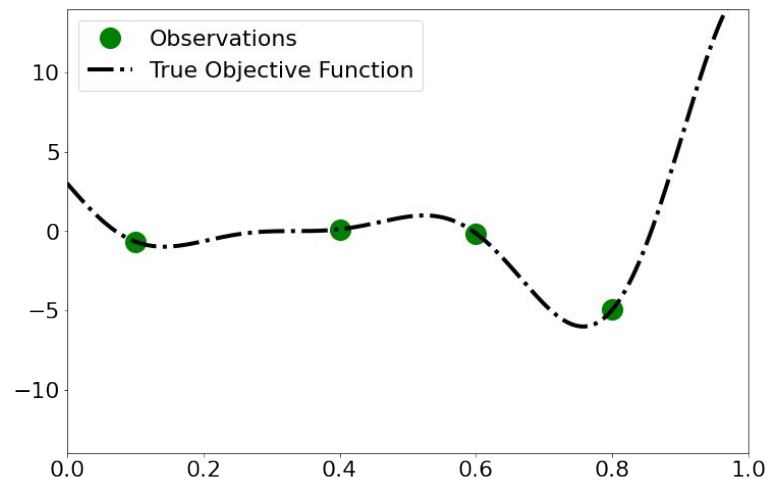
Expected Improvement

Demo BO loop



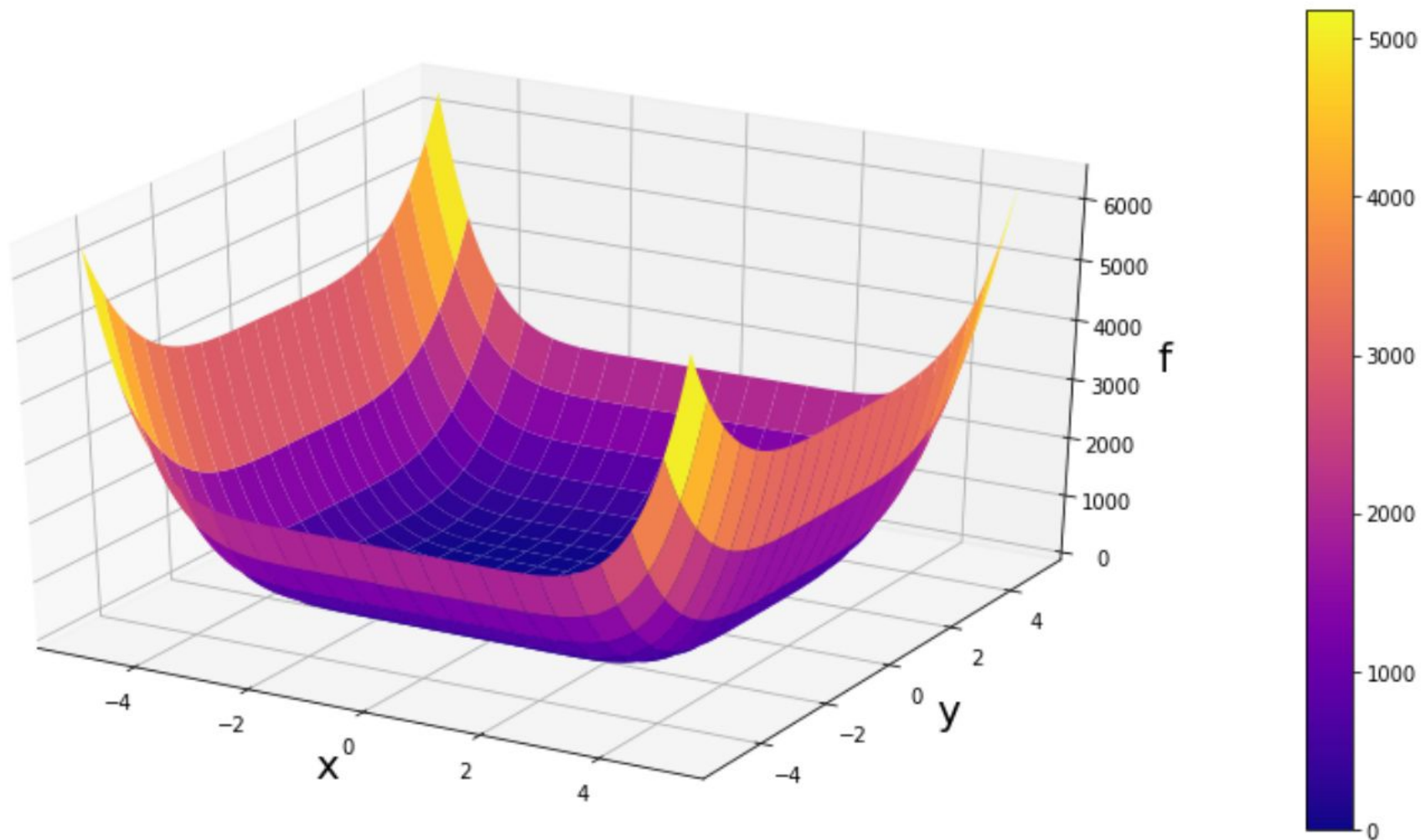
Expected Improvement

Demo BO loop



BO Demo 2

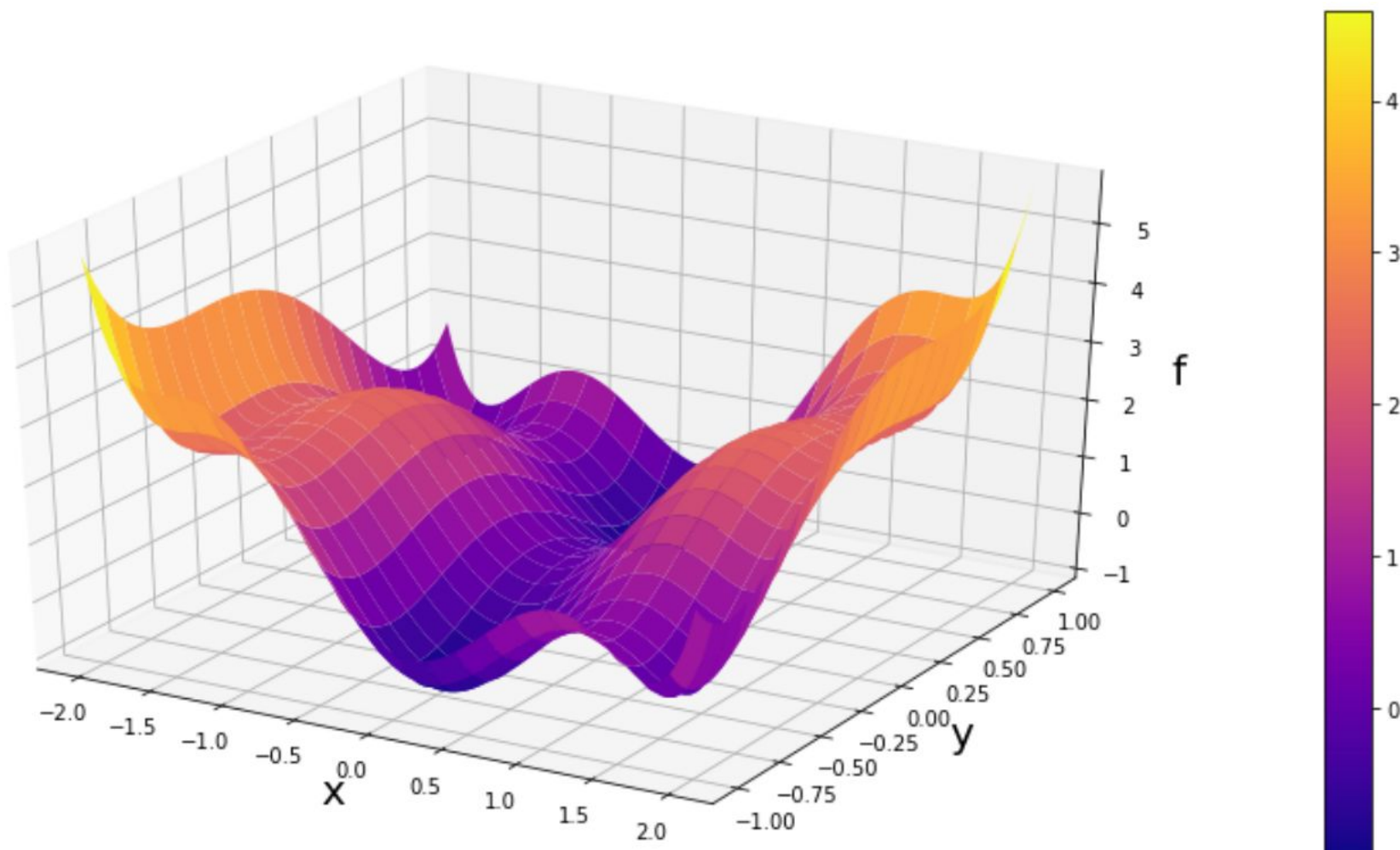
Let minimize the 6 Hump Camel function



Looks like we **can** use a local optimizer!

BO Demo 2

Zoom in: Perhaps not quite as easy?

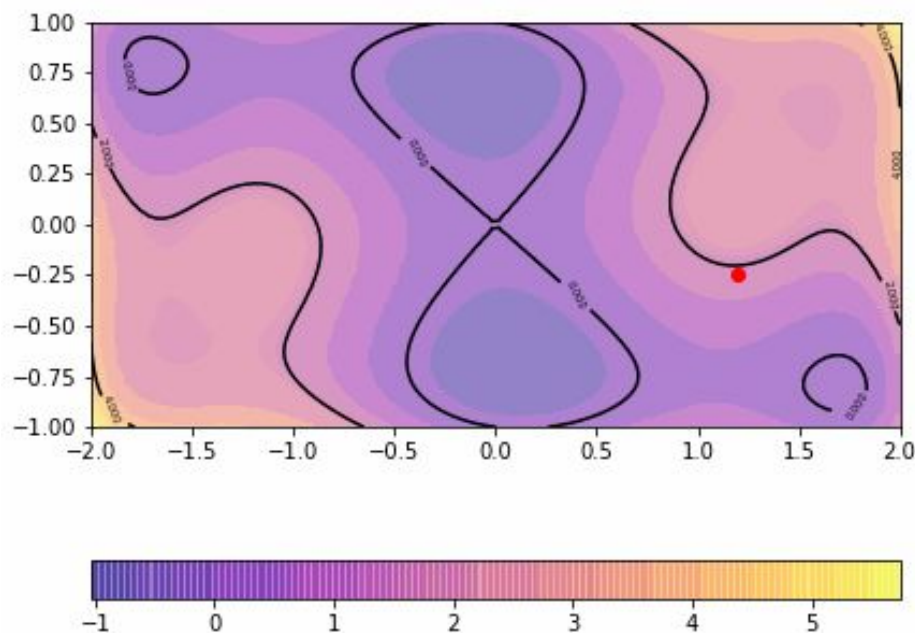


Looks like we **cannot** use a local optimizer!

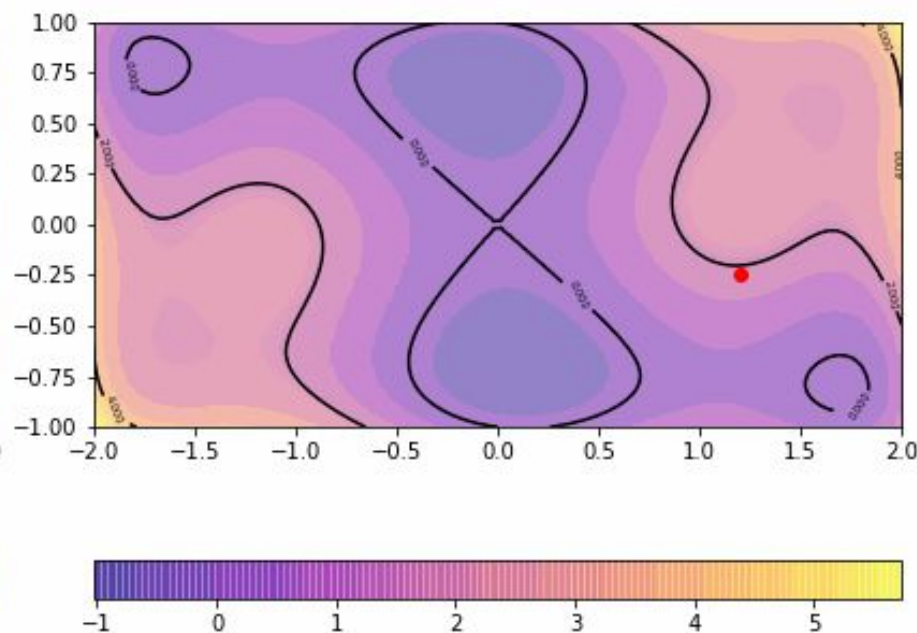
BO Demo 2

Bayesian optimization is a global optimizer

Bayesian optimization (global)

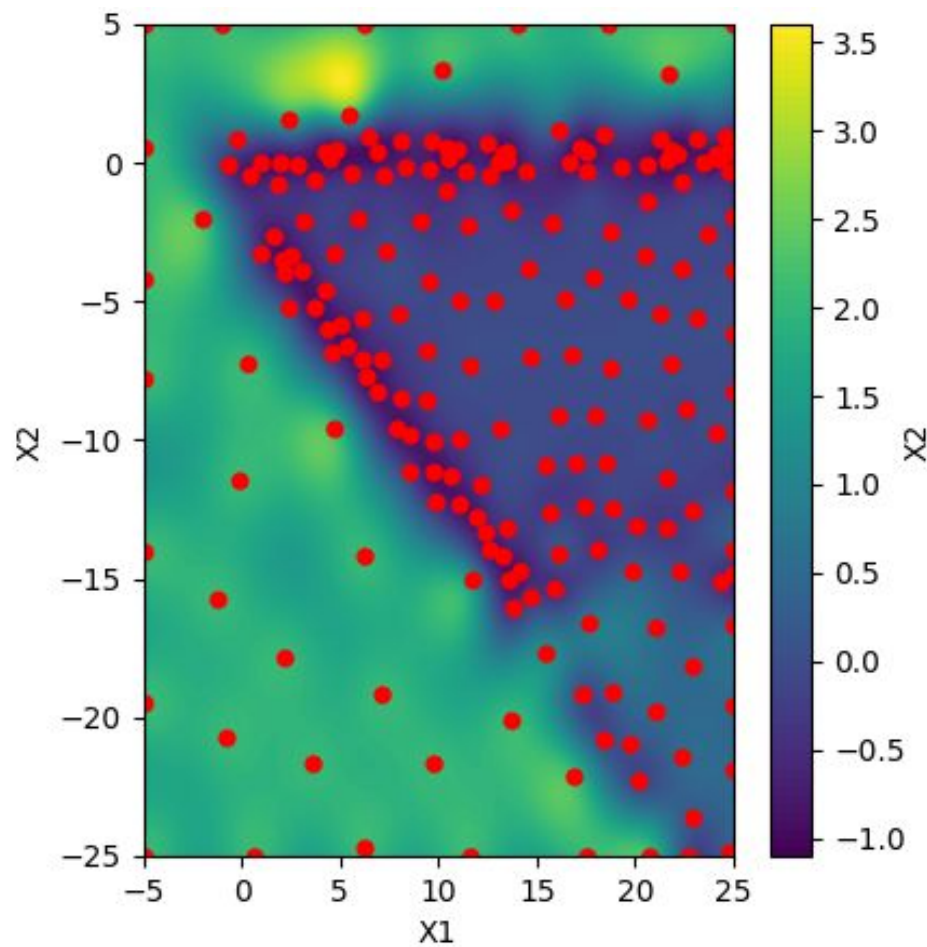


Gradient descent (local)



BO Demo 3

Efficient coverage of the search space



So why do we care about Bayesian Optimization?

So why do we care about Bayesian Optimization?

- BO performs **global** optimization (good for multi-modal functions)

So why do we care about Bayesian Optimization?

- BO performs **global** optimization (good for multi-modal functions)
- BO can optimize under a **limited evaluation budget** (great for problems with high evaluation costs)


So why do we care about Bayesian Optimization?

- BO performs **global** optimization (good for multi-modal functions)
- BO can optimize under a **limited evaluation budget** (great for problems with high evaluation costs)
 - Simulating performance of a car engine (mins)
 - Training a large ML model (hours)
 - Synthesising a new molecule (weeks)
 - Testing performance of a wind turbine in real world (months)




Increasing cost

So why do we care about Bayesian Optimization?

- BO performs **global** optimization (good for multi-modal functions)
 - BO can optimize under a **limited evaluation budget** (great for problems with high evaluation costs)
 - Simulating performance of a car engine (mins)
 - Training a large ML model (hours)
 - Synthesising a new molecule (weeks)
 - Testing performance of a wind turbine in real world (months)
- Increasing cost
- 
- We do not need gradients or noiseless observations (i.e. **black-box** optimization)

So why do we care about Bayesian Optimization?

- BO performs **global** optimization (good for multi-modal functions)
 - BO can optimize under a **limited evaluation budget** (great for problems with high evaluation costs)
 - Simulating performance of a car engine (mins)
 - Training a large ML model (hours)
 - Synthesising a new molecule (weeks)
 - Testing performance of a wind turbine in real world (months)
- Increasing cost
- 
- We do not need gradients or noiseless observations (i.e. **black-box** optimization)

BO: clever modelling rather than brute force!

Cool things that you can do with BO

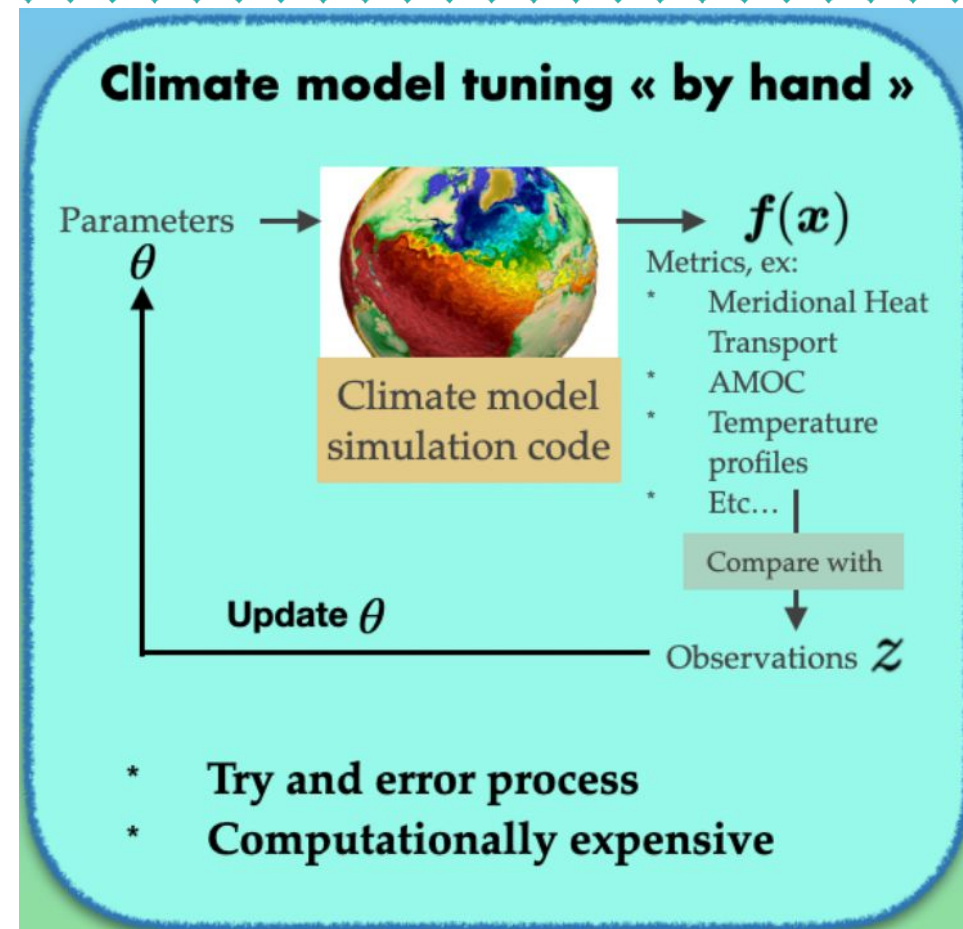
- Fine-tune the performance of AlphaGO (<https://arxiv.org/abs/1812.06855>)
- Allow Amazon Alexa learn how to speak with new voices (<https://arxiv.org/abs/2002.01953>)
- Efficiently find new molecules / genes (<https://arxiv.org/abs/2010.00979>)
- Fine-tune electric car engines
- Optimize large climate models

A great new reference for BO: **<https://bayesoptbook.com/>**

So, Climate model
calibration?

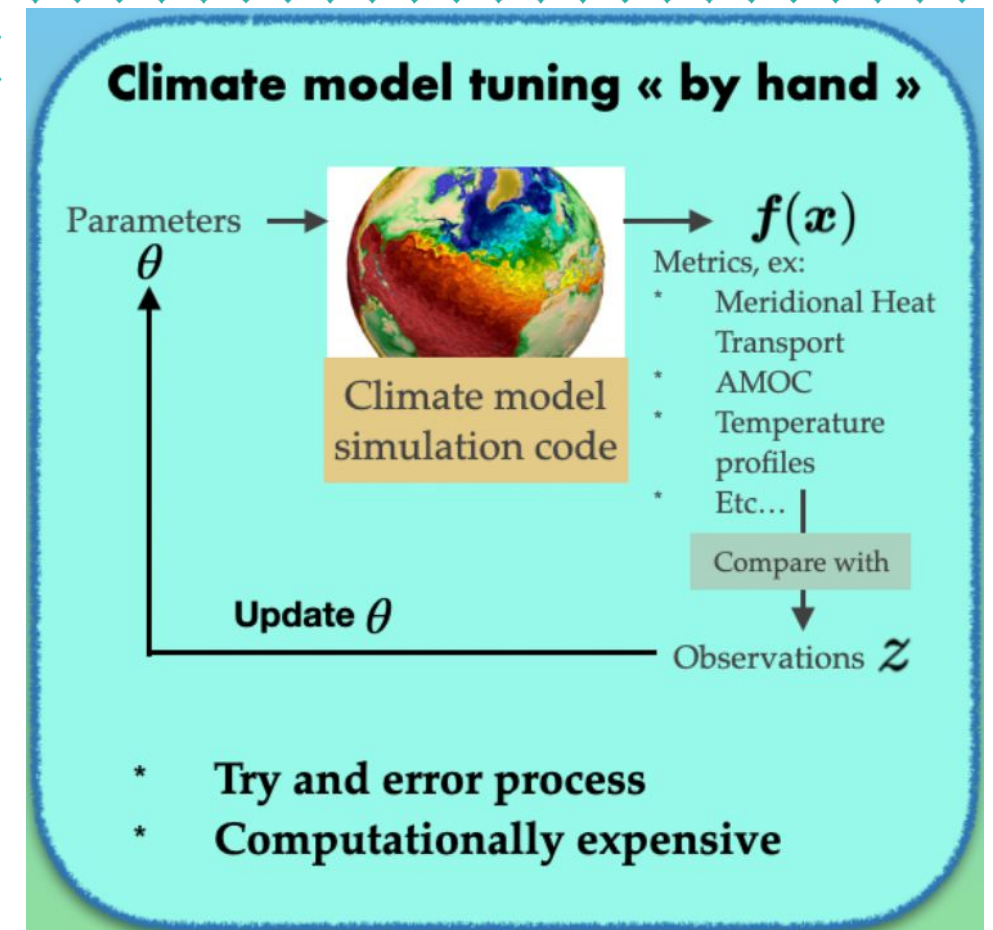
Climate model calibration

Identifying reasonable values for model parameters



Climate model calibration

Identifying reasonable values for model parameters

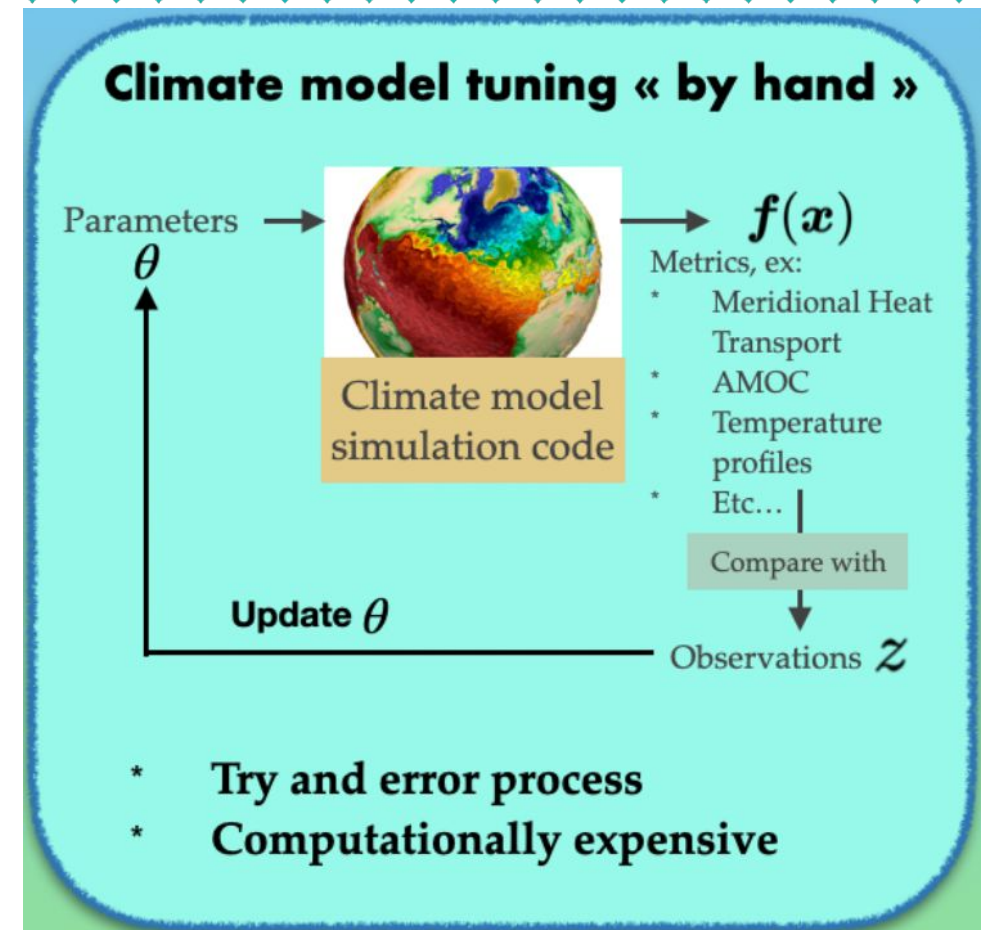


Lguensat et al. 2022.

- Need to find parameters that give high plausibility to historical data —————> a **function maximisation** problem

Climate model calibration

Identifying reasonable values for model parameters

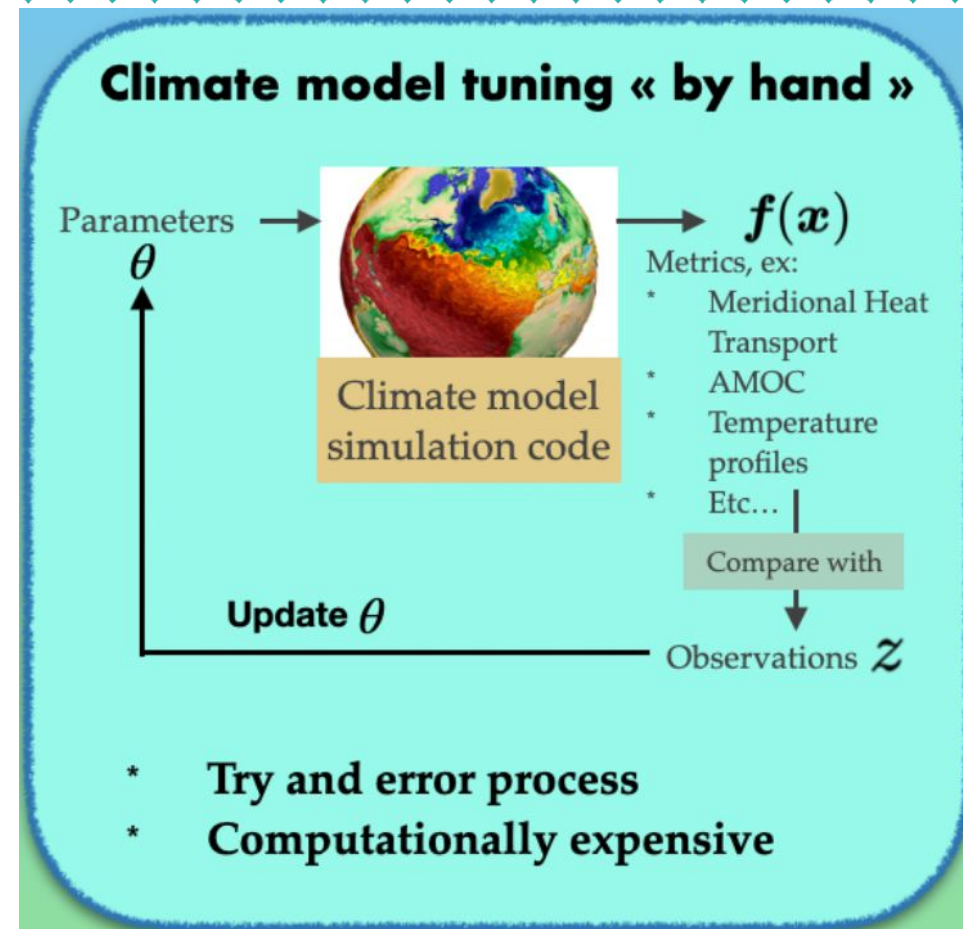


Lguensat et al. 2022.

- Need to find parameters that give high plausibility to historical data —————> a **function maximisation** problem
- Climate models are expensive —————> can only afford a **limited number of evaluations** (no grid!)

Climate model calibration

Identifying reasonable values for model parameters



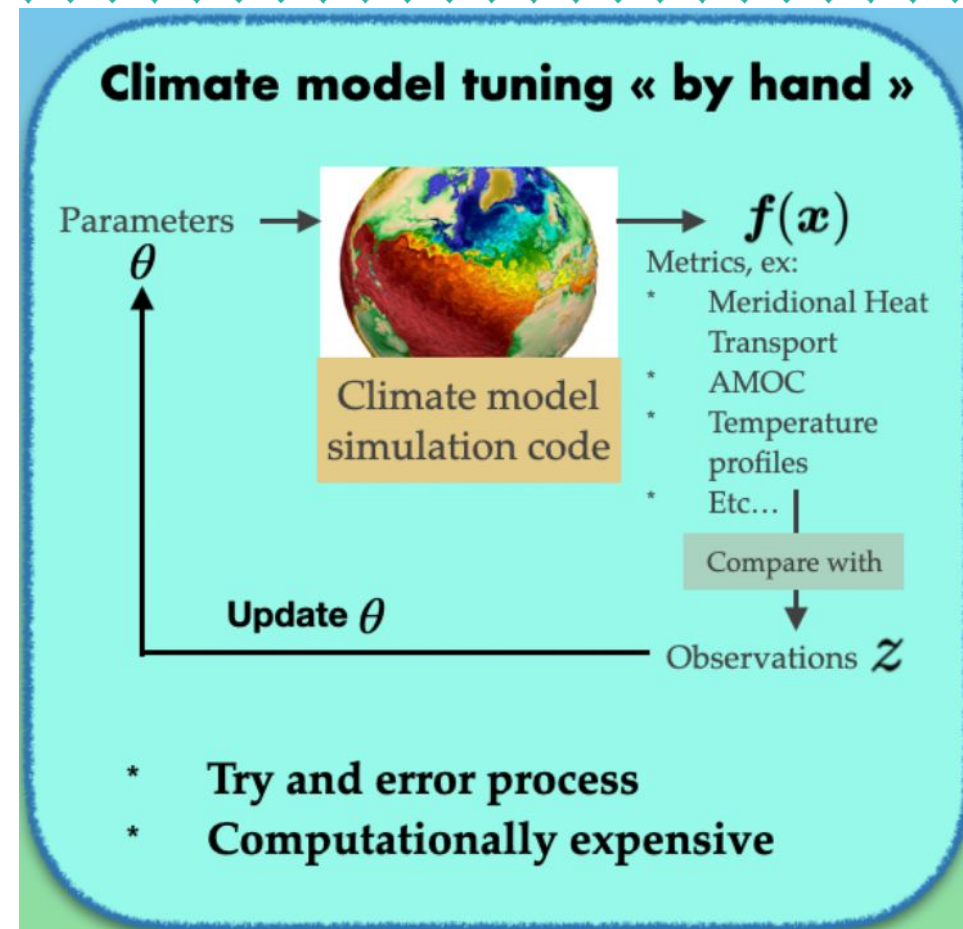
Lguensat et al. 2022.

- Need to find parameters that give high plausibility to historical data —————> a **function maximisation** problem
- Climate models are expensive —————> can only afford a **limited number of evaluations** (no grid!)
- We do not have gradients (easily) and limited prior knowledge —————> a **black-box** objective function

Climate model calibration

Identifying reasonable values for model parameters

So we have a resource-constrained black-box function optimisation!



Lguensat et al. 2022.

- Need to find parameters that give high plausibility to historical data —————> a **function maximisation** problem
- Climate models are expensive —————> can only afford a **limited number of evaluations** (no grid!)
- We do not have gradients (easily) and limited prior knowledge —————> a **black-box** objective function

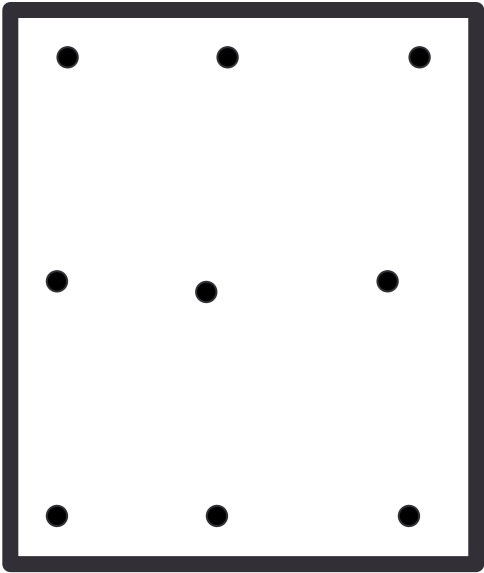
Climate model calibration by iteratively refocusing

sequentially whittle down the plausible region



Climate model calibration by iteratively refocusing

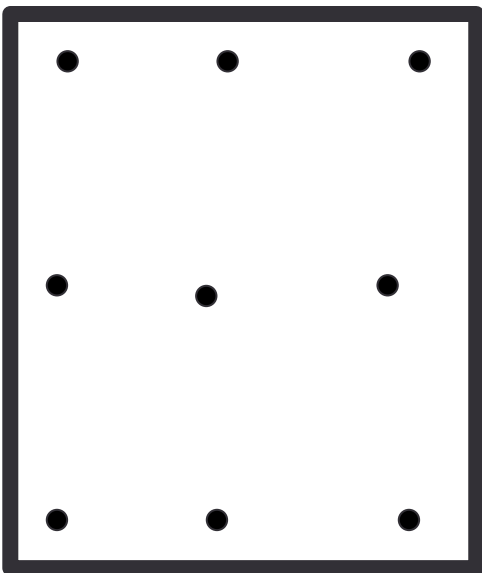
sequentially whittle down the plausible region



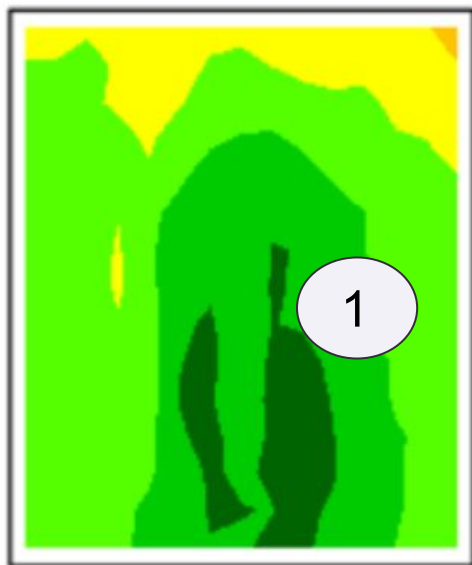
Initial Design

Climate model calibration by iteratively refocusing

sequentially whittle down the plausible region



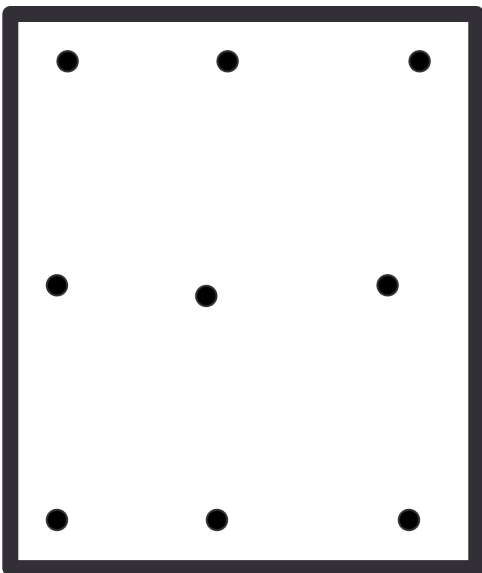
Initial Design



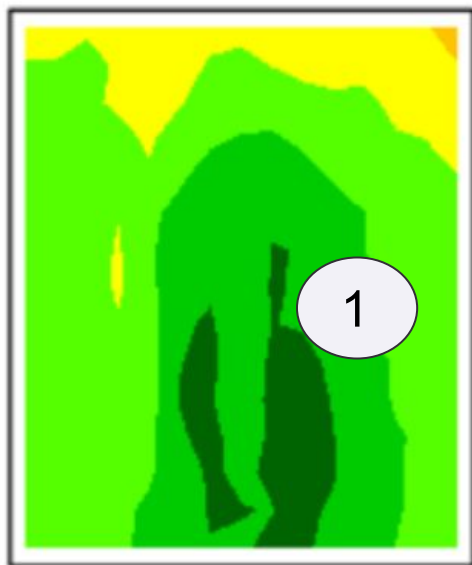
Predicted
implausibility

Climate model calibration by iteratively refocusing

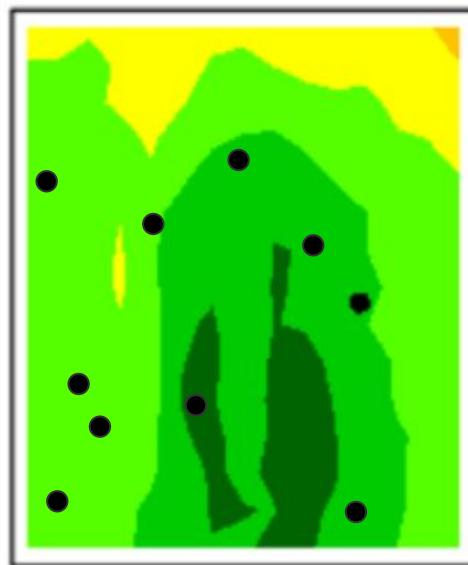
sequentially whittle down the plausible region



Initial Design



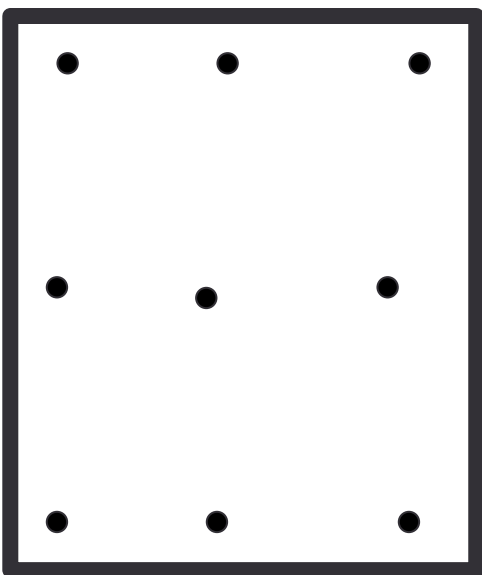
Predicted
implausibility



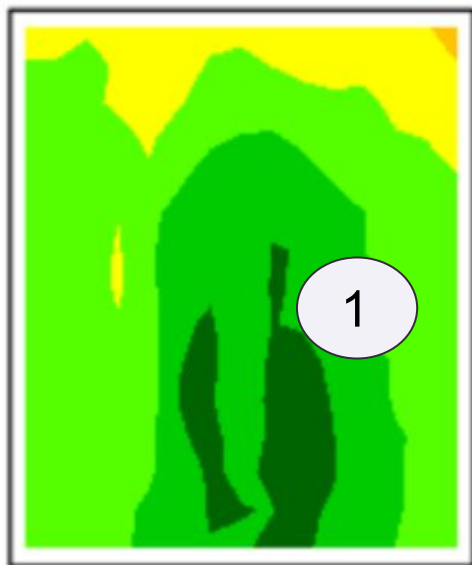
1st set of
evaluations

Climate model calibration by iteratively refocusing

sequentially whittle down the plausible region



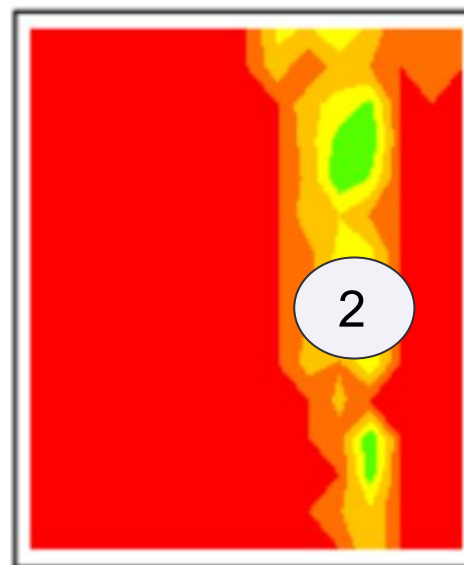
Initial Design



Predicted
implausibility



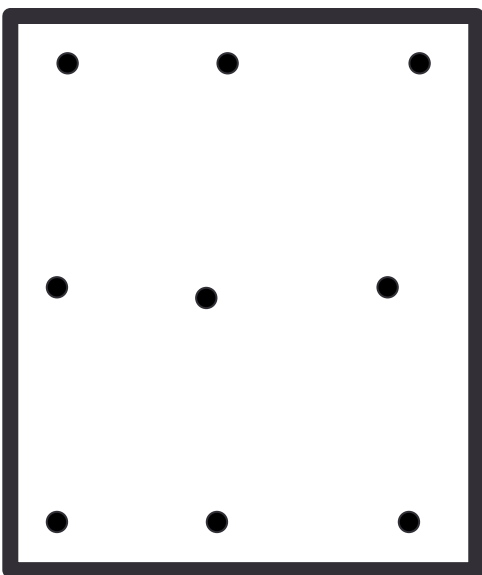
1st set of
evaluations



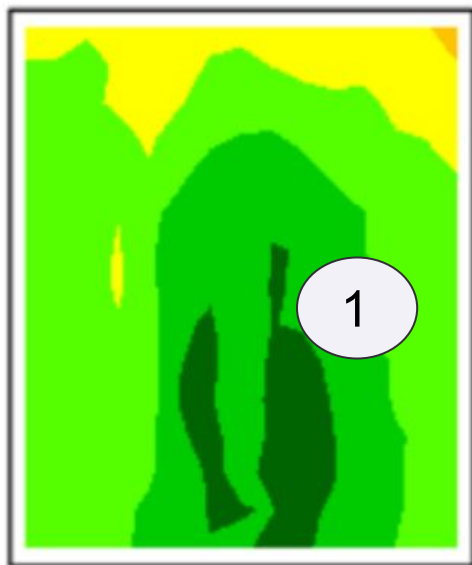
Predicted
implausibility

Climate model calibration by iteratively refocusing

sequentially whittle down the plausible region



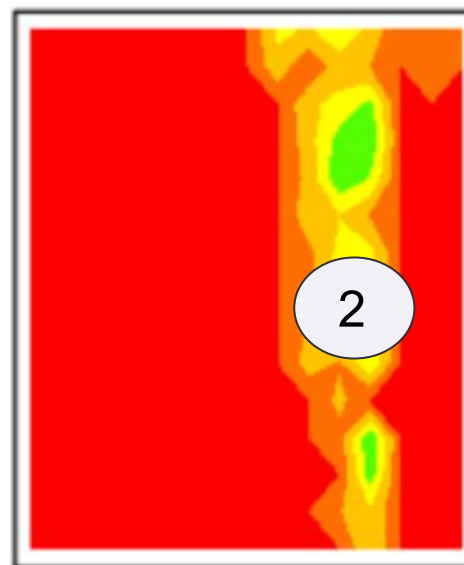
Initial Design



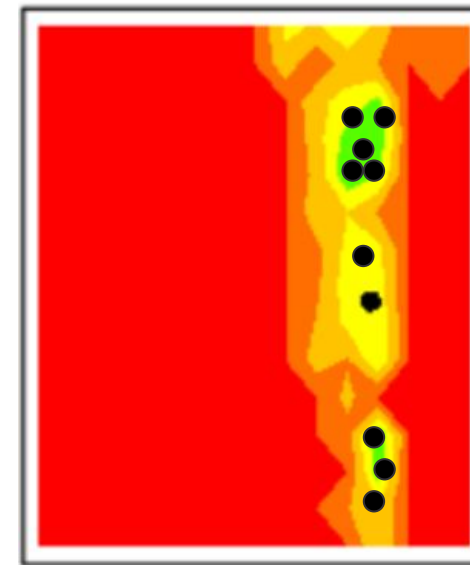
Predicted
implausibility



1st set of
evaluations



Predicted
implausibility



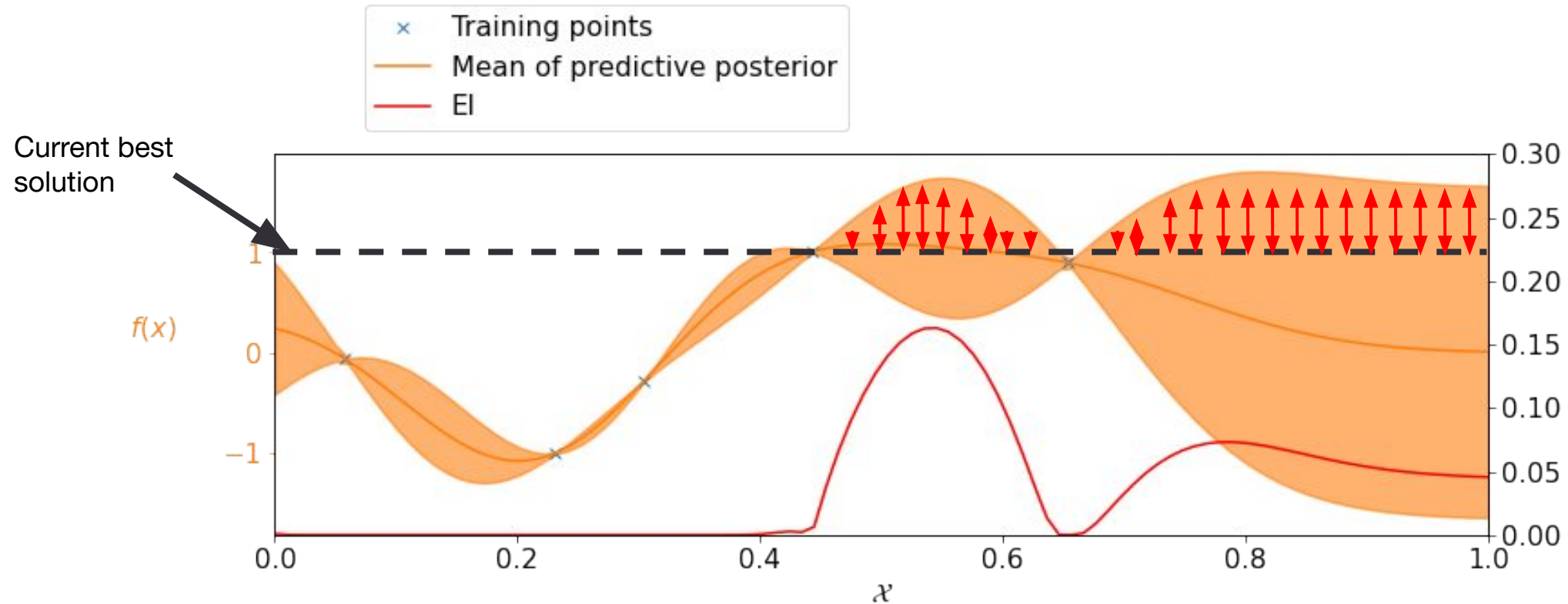
2nd set of
evaluations

Back to molecular design

Large batches

Automatically choosing batches of points

Using GP posteriors and utility functions



How to pick **3** points ?



Automatically choosing batches of molecules

Using GP posteriors and utility functions

- $\alpha_{\text{EI}}(\text{molecule}) = \mathbb{E}_f[\max(f - f^*, 0)] \quad f \sim \mathcal{N}(\mu, \sigma^2)$



Automatically choosing batches of molecules

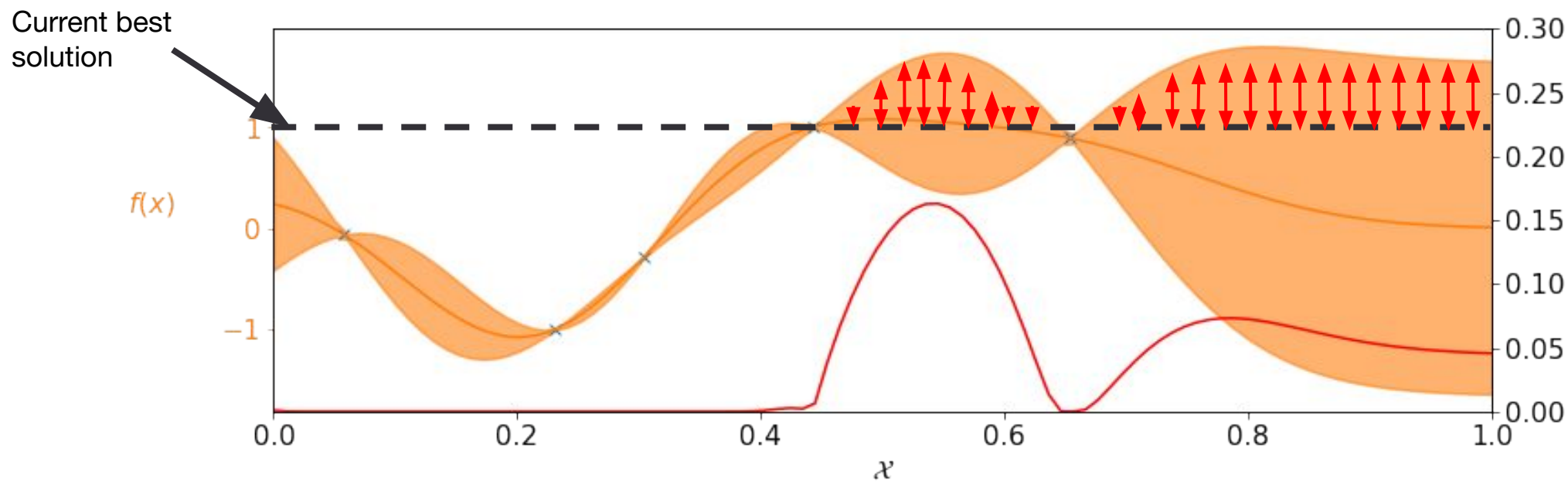
Using GP posteriors and utility functions

- $\alpha_{\text{EI}}(\text{molecule}) = \mathbb{E}_f[\max(f - f^*, 0)]$
- $\alpha_{\text{EI}}(\{\text{molecule}_i, \text{molecule}_j\}) = ???$

Automatically choosing batches of molecules

Using GP posteriors and utility functions

- $\alpha_{\text{EI}}(\text{molecule}) = \mathbb{E}_f[\max(f - f^*, 0)]$
- $\alpha_{\text{EI}}(\{\text{molecule}_i, \text{molecule}_j\}) = \mathbb{E}_{f_i, f_j}[\max(f_i - f^*, f_j - f^*, 0)]$





Automatically choosing batches of molecules

Using GP posteriors and utility functions

- $\alpha_{\text{EI}}(\text{molecule}_i) = \mathbb{E}_f[\max(f - f^*, 0)]$
- $\alpha_{\text{EI}}(\{\text{molecule}_i, \text{molecule}_j\}) = \mathbb{E}_{f_i, f_j}[\max(f_i - f^*, f_j - f^*, 0)]$

$$\begin{pmatrix} f_i \\ f_j \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} \mu_i \\ \mu_j \end{pmatrix}, \begin{pmatrix} \Sigma_{i,i} & \Sigma_{i,j} \\ \Sigma_{j,i} & \Sigma_{j,j} \end{pmatrix} \right)$$



Automatically choosing batches of molecules

Using GP posteriors and utility functions

- $\alpha_{\text{EI}}(\text{molecule}) = \mathbb{E}_f[\max(f - f^*, 0)]$
- $\alpha_{\text{EI}}(\{\text{molecule}_i, \text{molecule}_j\}) = \mathbb{E}_{f_i, f_j}[\max(f_i - f^*, f_j - f^*, 0)]$
- $\alpha_{\text{EI}}(\{\text{molecule}_1, \dots, \text{molecule}_B\}) = ???$

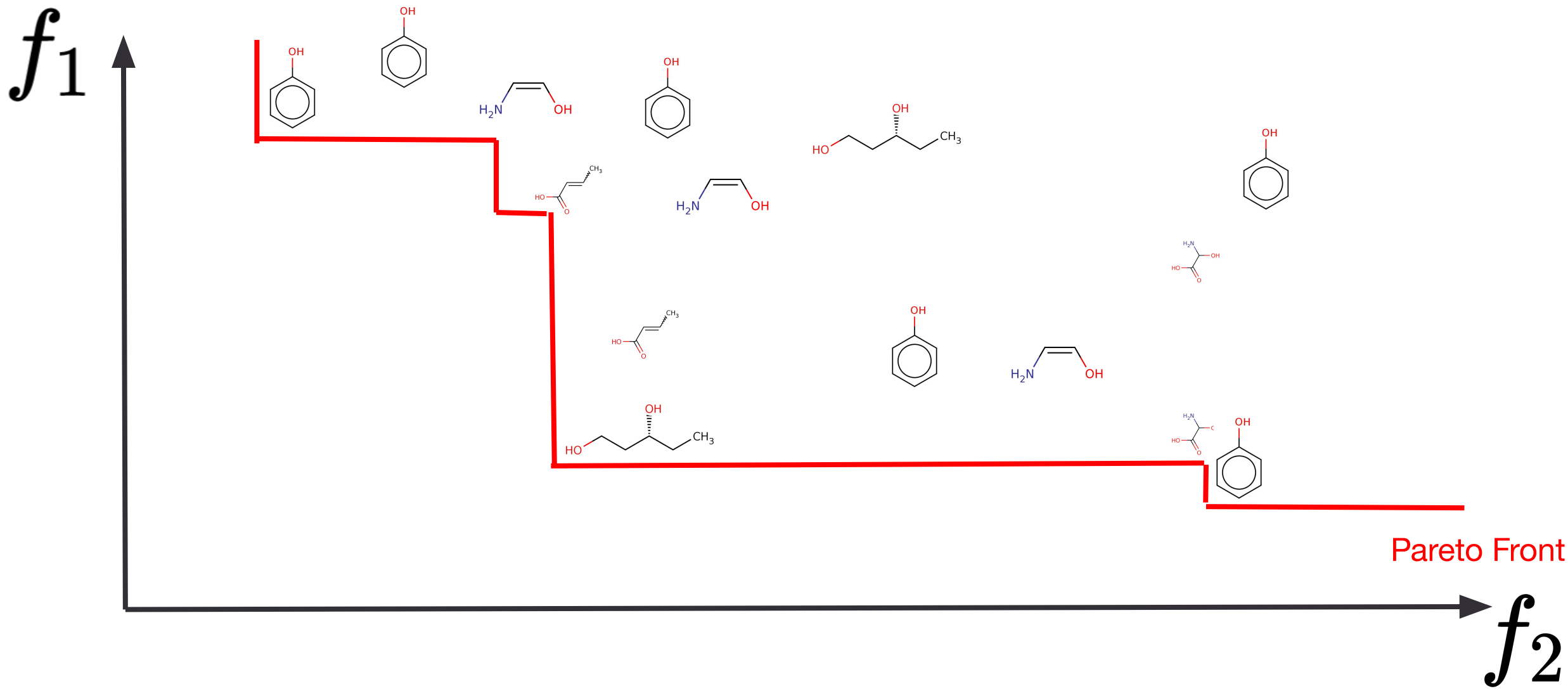
Back to molecular design

Multiple objectives


The figure displays a 2D plot with axes labeled f_1 (vertical) and f_2 (horizontal). The plot contains 15 chemical structures, which are scattered across the space. The structures include various organic molecules, such as phenols, alkenes, and alcohols, represented by their skeletal structures. The distribution shows some clustering, with a group of structures in the upper-left region and another group in the lower-right region.

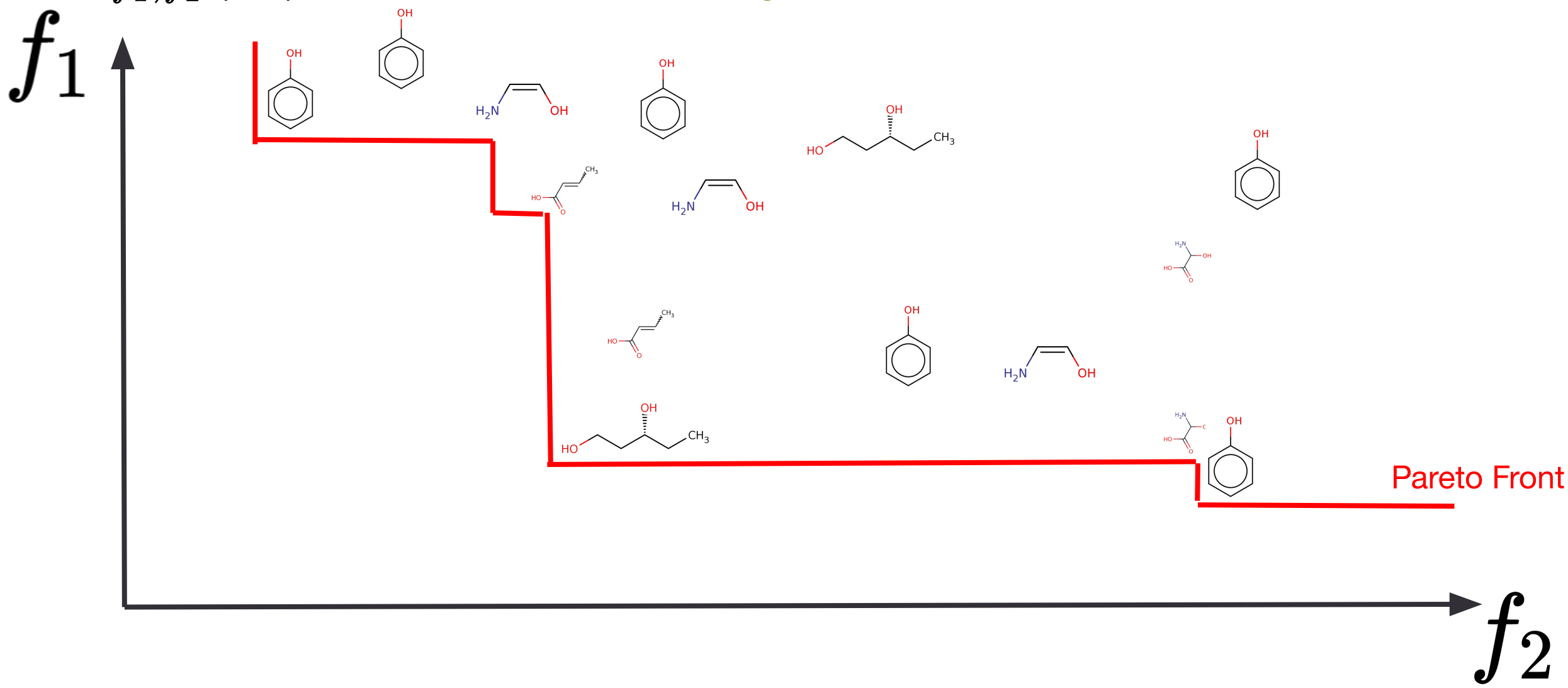
Multi-objective Optimisation

>1 competing objectives




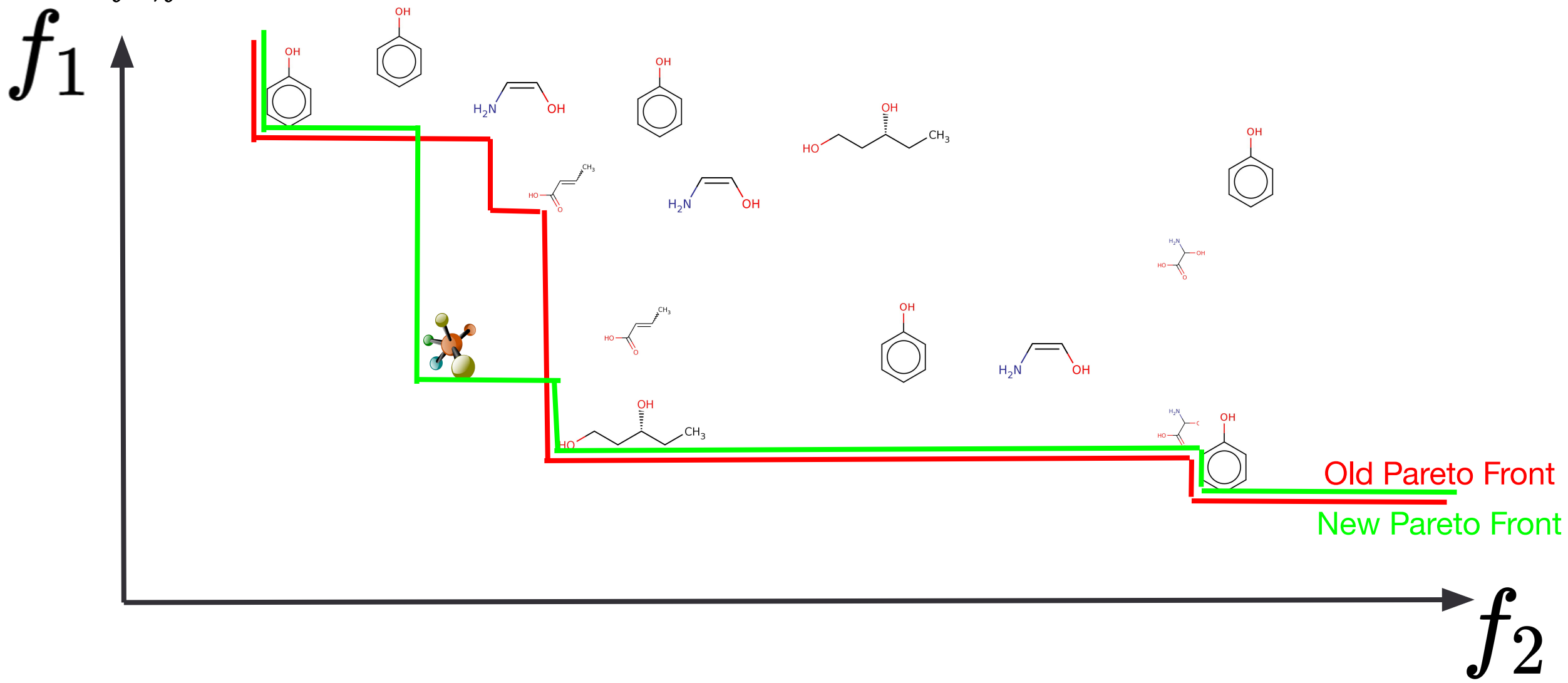
Multi-objective Optimisation

$U_{f_1, f_2}(\text{molecule})$: what is the utility of evaluating  if it will return (f_1, f_2)




Multi-objective Optimisation

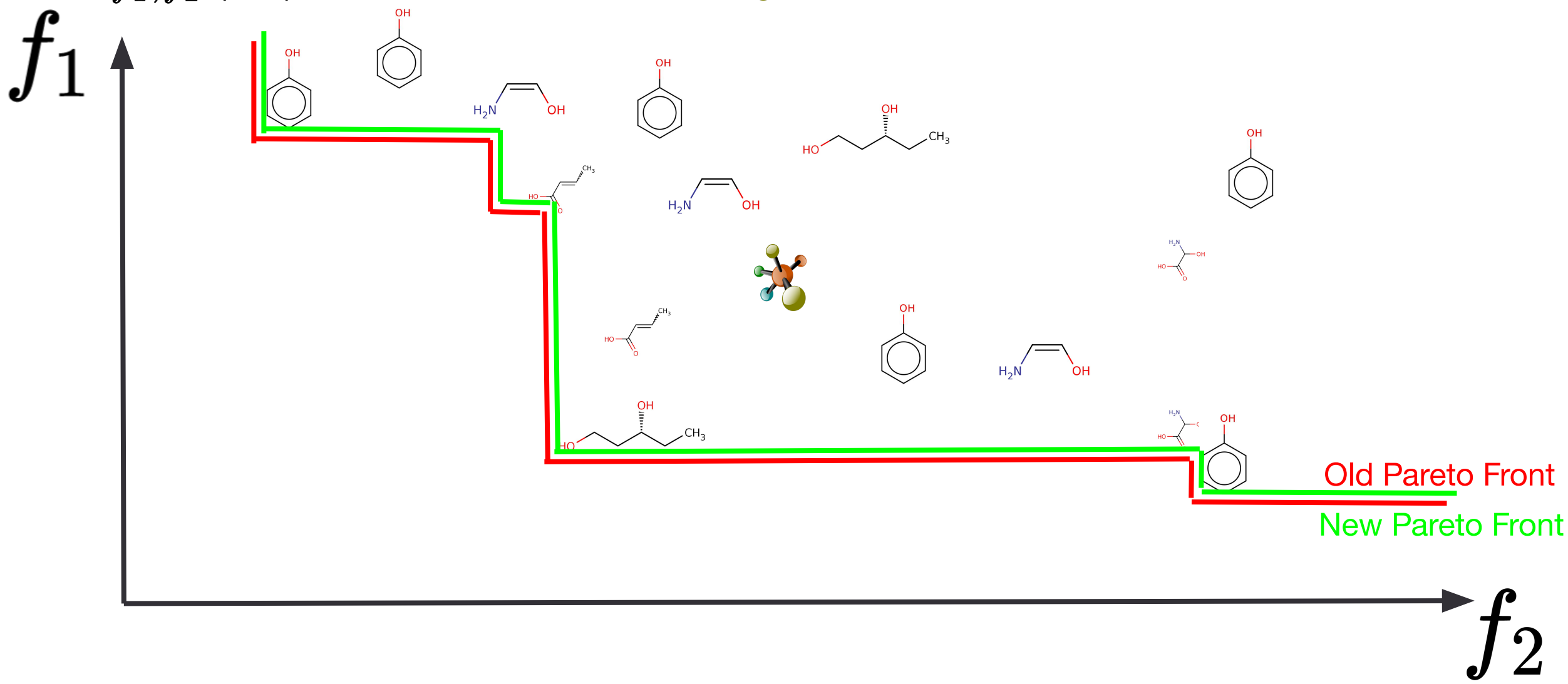
$U_{f_1, f_2}(\text{molecule})$: what is the utility of evaluating  if it will return (f_1, f_2)




The figure displays a multi-objective optimization result. The axes are f_1 (vertical) and f_2 (horizontal). The 'Old Pareto Front' is shown in red, and the 'New Pareto Front' is shown in green. The green frontier represents a set of solutions that are generally superior to the old ones, as they achieve lower f_1 values for the same f_2 values. Chemical structures are used as data points, with some specifically placed on the frontier lines to illustrate the chemical space covered by the optimization process.

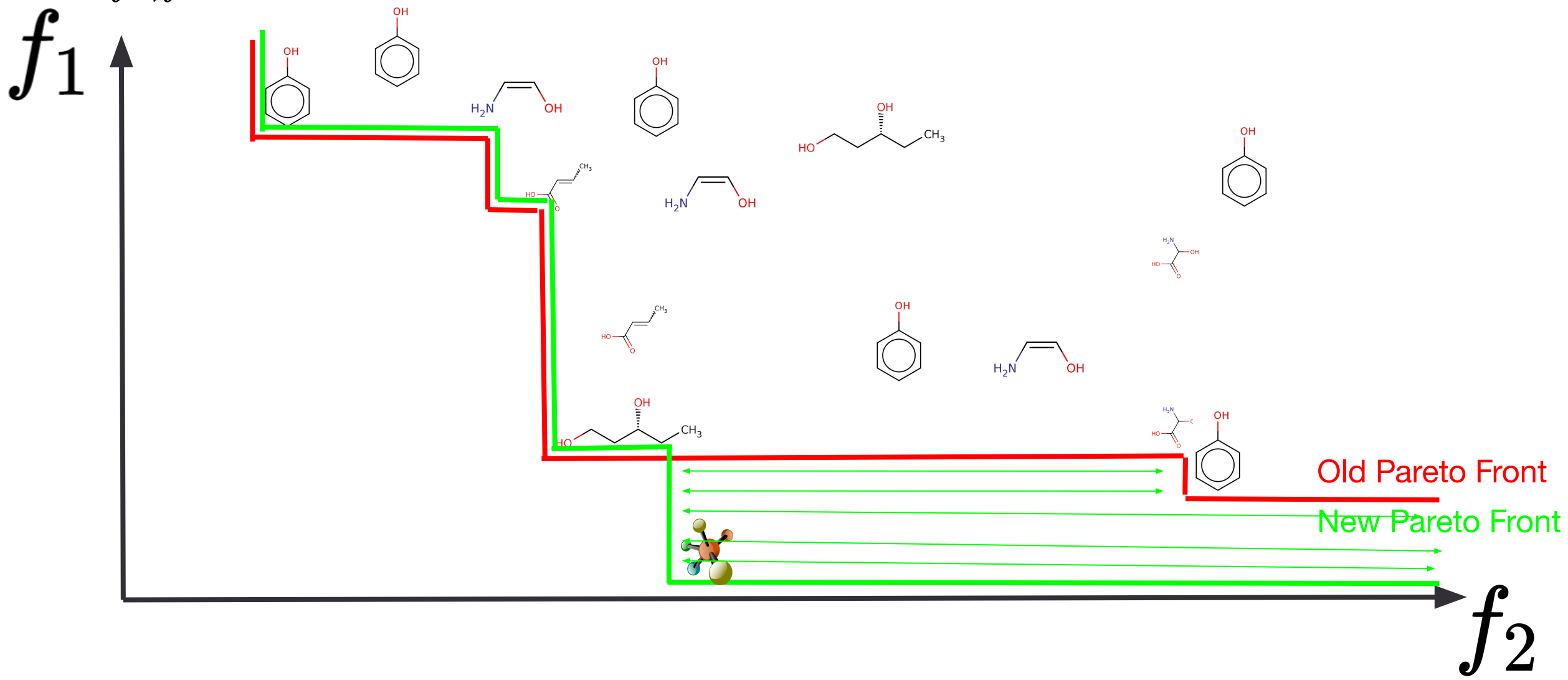
Multi-objective Optimisation

$U_{f_1, f_2}(\text{molecule})$: what is the utility of evaluating  if it will return (f_1, f_2)



Multi-objective Optimisation

$U_{f_1, f_2}(\text{molecule})$: what is the utility of evaluating  if it will return (f_1, f_2)



Multi-objective Optimisation


$U_{f_1, f_2}(\text{🧬})$: what is the utility of evaluating 🧬 if it will return (f_1, f_2)

- Use expected hyper-volume improvement $\alpha_{\text{EHVI}}(\text{🧬}) = \mathbb{E}_{f_1, f_2}(U_{f_1, f_2}(\text{🧬}))$

$$f_1 \sim \mathcal{N}(\mu_1, \sigma_1^2)$$

$$f_2 \sim \mathcal{N}(\mu_2, \sigma_2^2)$$

Multi-objective Optimisation

$U_{f_1, f_2}(\text{molecule})$: what is the utility of evaluating  if it will return (f_1, f_2)

- Use expected hyper-volume improvement $\alpha_{\text{EHVI}}(\text{molecule}) = \mathbb{E}_{f_1, f_2}(U_{f_1, f_2}(\text{molecule}))$

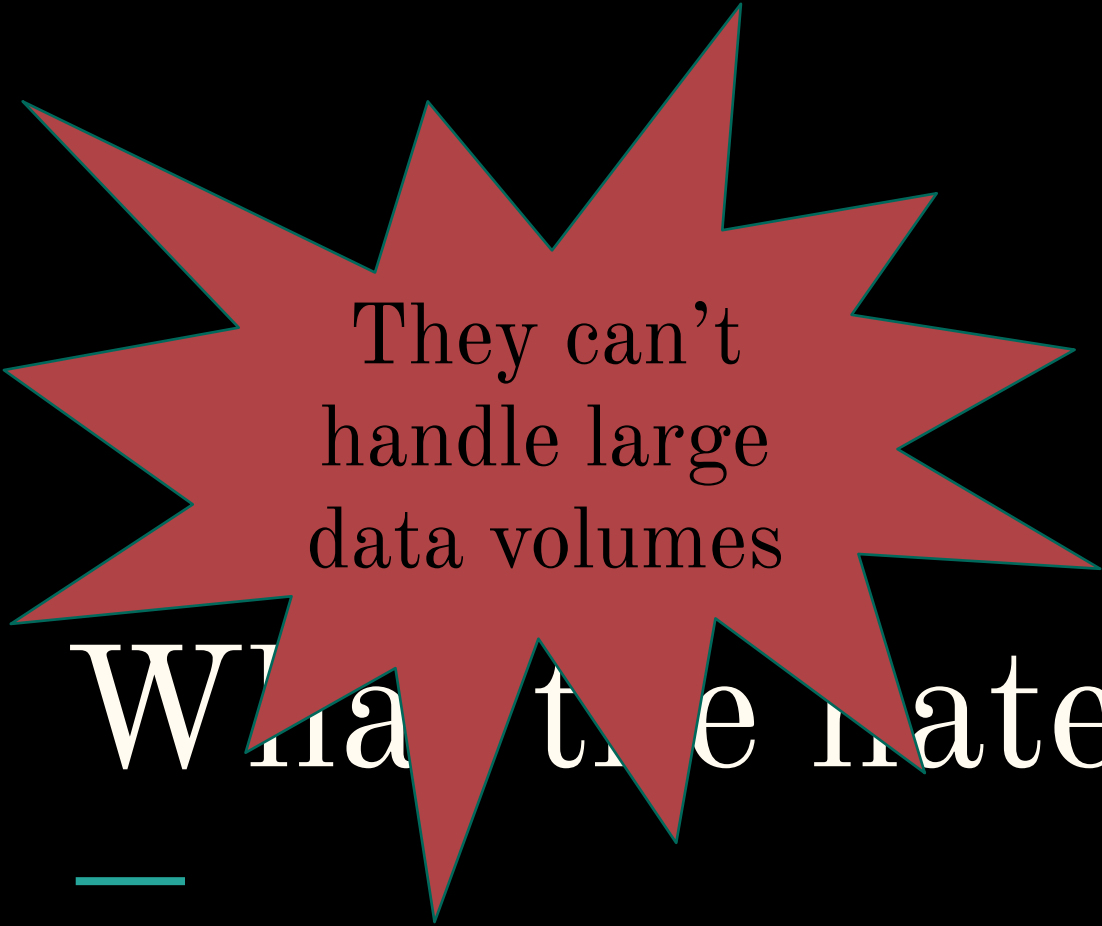
$$f_1 \sim \mathcal{N}(\mu_1, \sigma_1^2)$$

$$f_2 \sim \mathcal{N}(\mu_2, \sigma_2^2)$$

$$\alpha_{\text{EHVI}}(\{\text{molecule}_i, \text{molecule}_j\}) = ???$$


What the haters say

—



They can't
handle large
data volumes


What the natters say



They can't
handle large
data volumes

Only for
Gaussian
data.....

maters say



They can't
handle large
data volumes

They can't
handle
high-dimension
al data

Only for
Gaussian
data.....

maters say

They can't
handle large
data volumes

They can't
handle
high-dimension
al data

Only for
Gaussian
data.....

They are not
interpretable

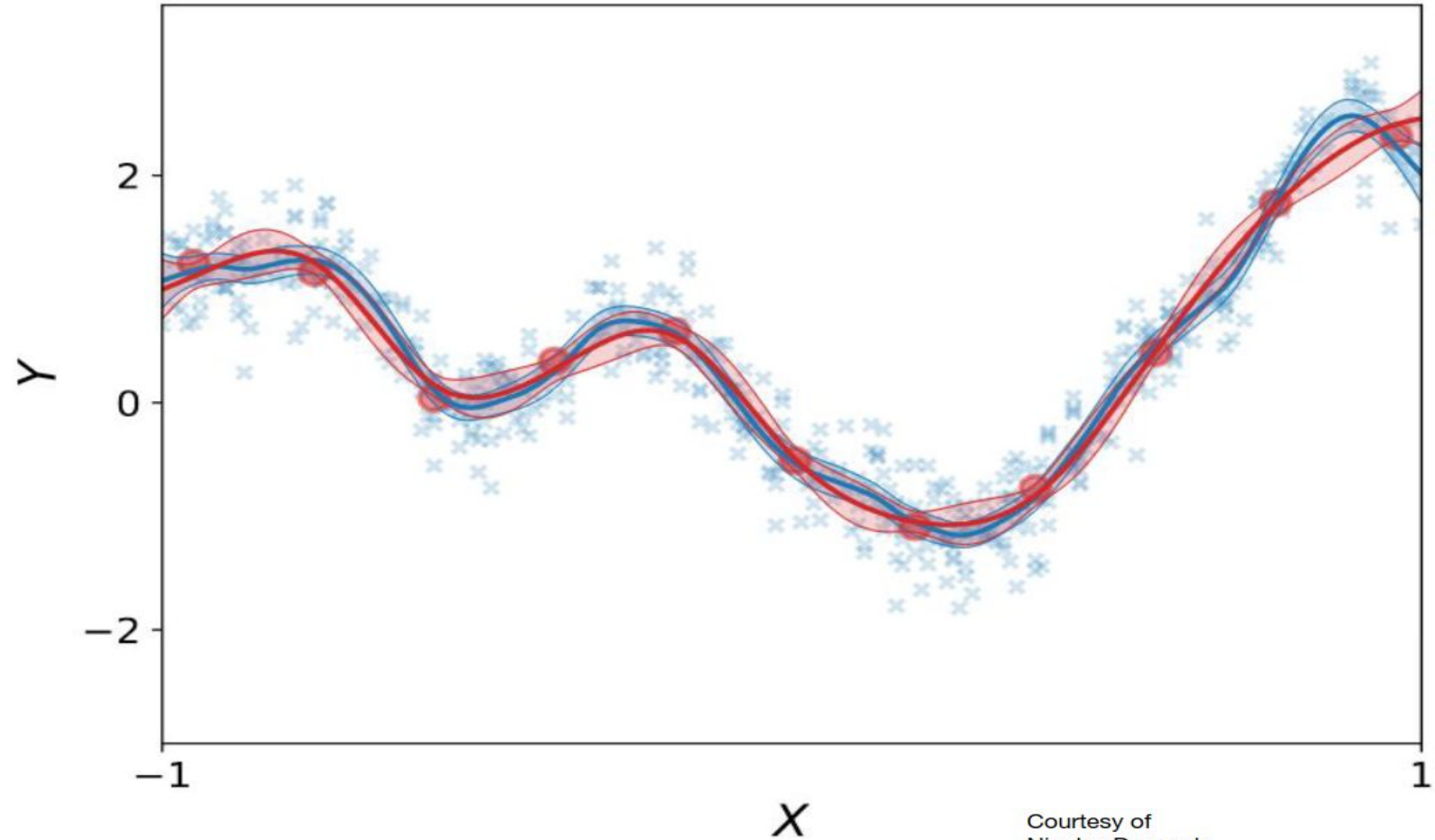
Matrix

...

GPs for big data?

- Use Sparse variational GP
- Replace with $M \ll N$

representative points

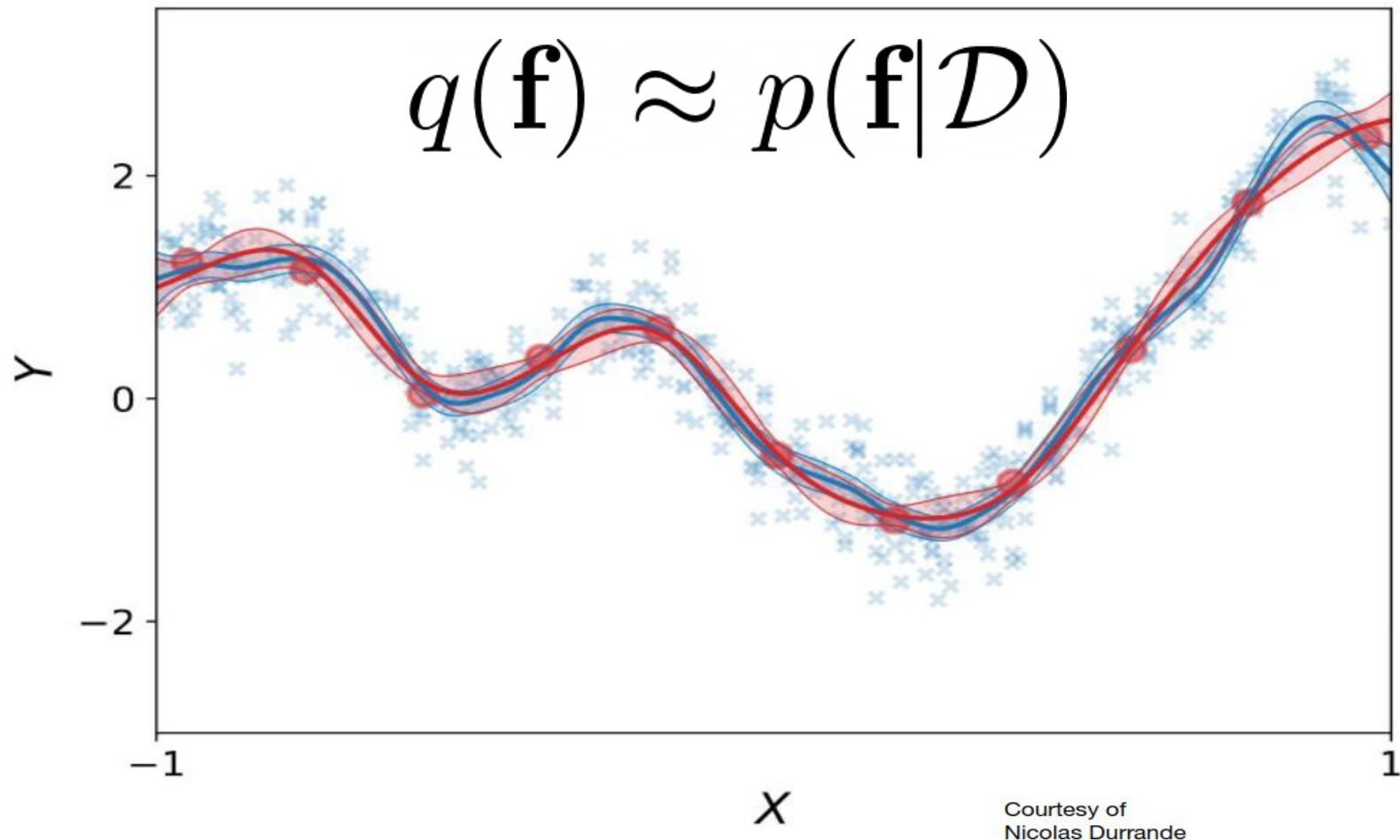


Courtesy of
Nicolas Durrande

GPs for big data?

- Use Sparse variational GP
- Replace with $M \ll N$

representative points

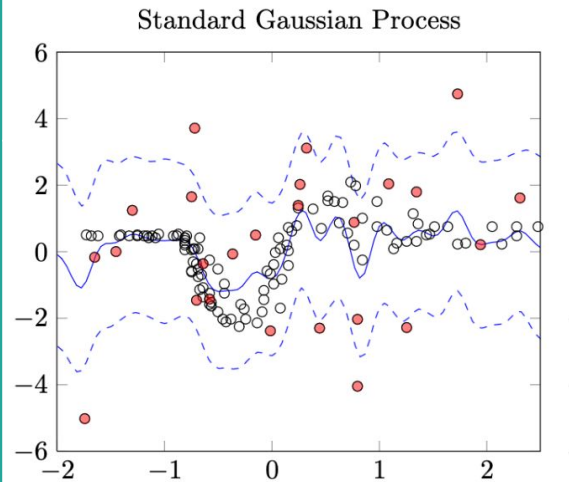


$$\begin{aligned}\text{ELBO}(q(\mathbf{f})) &= \int q(\mathbf{f}) \log p(\mathbf{y}|\mathbf{f}) d\mathbf{f} - \mathcal{KL}(q(\mathbf{f}), p(\mathbf{f})) \\ &= \sum_{i=1}^N \int q(f_i) \log p(y_i|f_i) d\mathbf{f} - \mathcal{KL}(q(\mathbf{f}), p(\mathbf{f}))\end{aligned}$$

$$y_i \sim \mathcal{N}(f(\mathbf{x}_i), \sigma^2)$$

SVGPs for non-Gaussian data?

(Hensman et al. 2015, Saul et al. 2016)

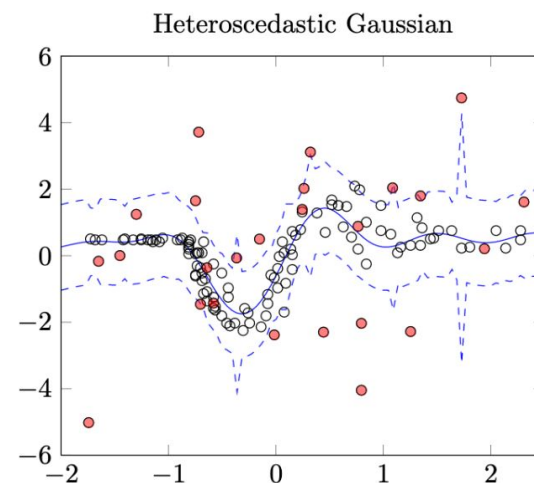
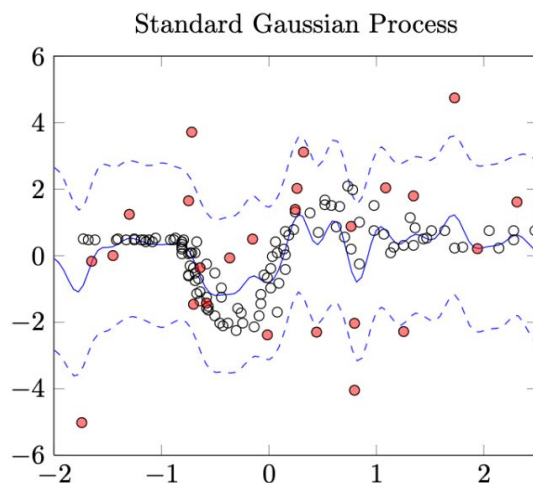


SVGPs for non-Gaussian data?

(Hensman et al. 2015, Saul et al. 2016)

~~$$y_i \sim \mathcal{N}(f(\mathbf{x}_i), \sigma^2)$$~~

$$y_i \sim \mathcal{N}(f_0(\mathbf{x}_i), e^{f_1(\mathbf{x}_i)})$$



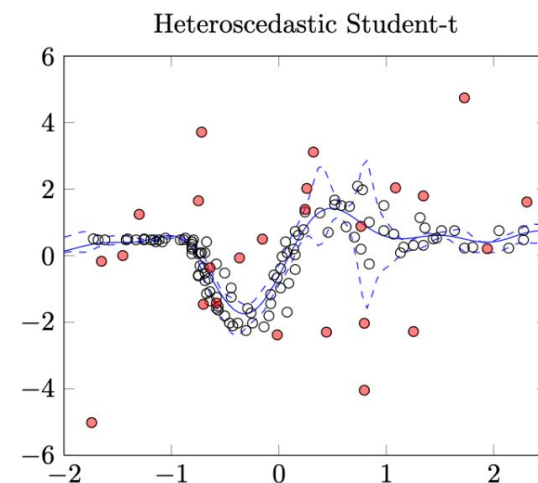
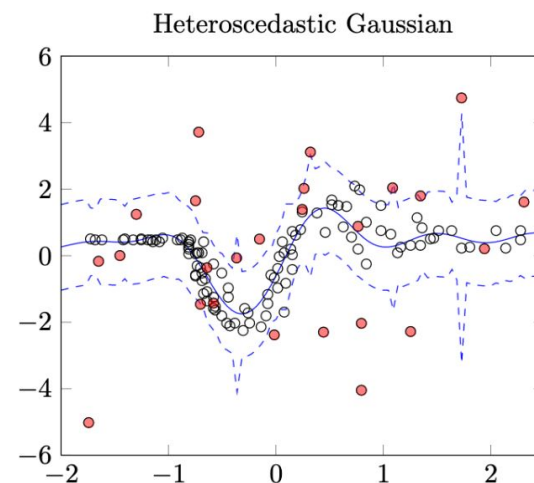
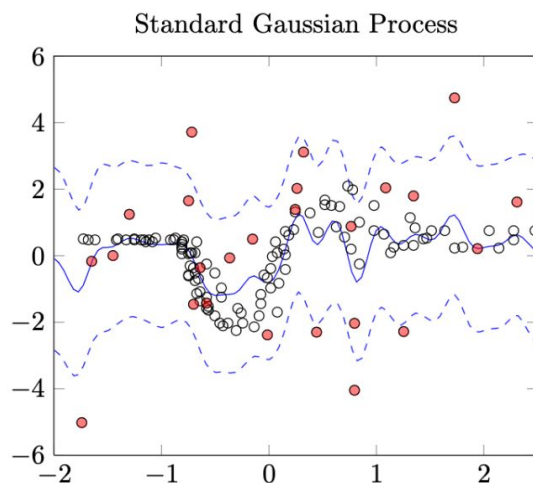
SVGPs for non-Gaussian data?

(Hensman et al. 2015, Saul et al. 2016)

~~$$y_i \sim \mathcal{N}(f(\mathbf{x}_i), \sigma^2)$$~~

~~$$y_i \sim \mathcal{N}(f_0(\mathbf{x}_i), e^{f_1(\mathbf{x}_i)})$$~~

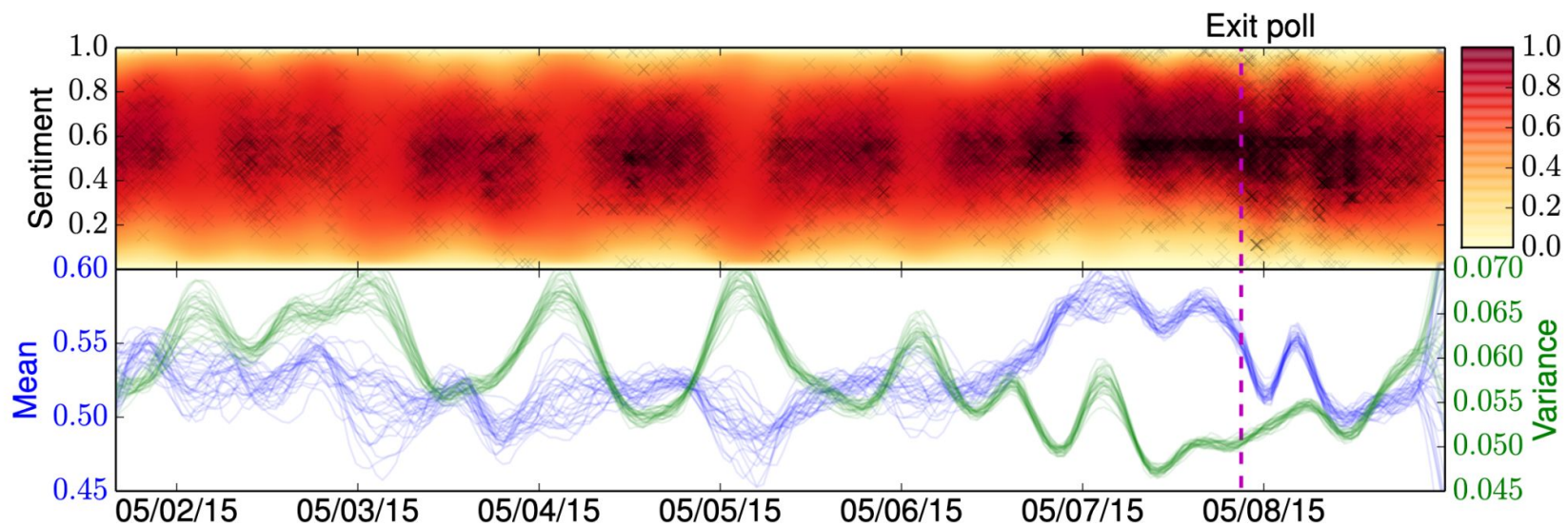
$$y_i \sim St(f_0(\mathbf{x}_i), e^{f_1(\mathbf{x}_i)}, \nu)$$



SVGPs for non-Gaussian data?

(Hensman et al. 2015, Saul et al. 2016)

$$y_i \sim \mathcal{B}(\alpha = f_0(\mathbf{x}_i), \beta = e^{f_1(\mathbf{x}_i)})$$



GPs for
high-dim
data?



Beware the curse of
dimensionality

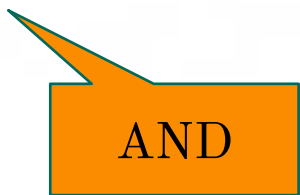


GPs for high-dim data?

- GPs are great in high-dim
- RBF kernels are not.....
- $l_i \propto \sqrt{D}$

$$k(\mathbf{x}, \mathbf{y}) = e^{-\frac{||\mathbf{x} - \mathbf{y}||^2}{2l^2}}$$

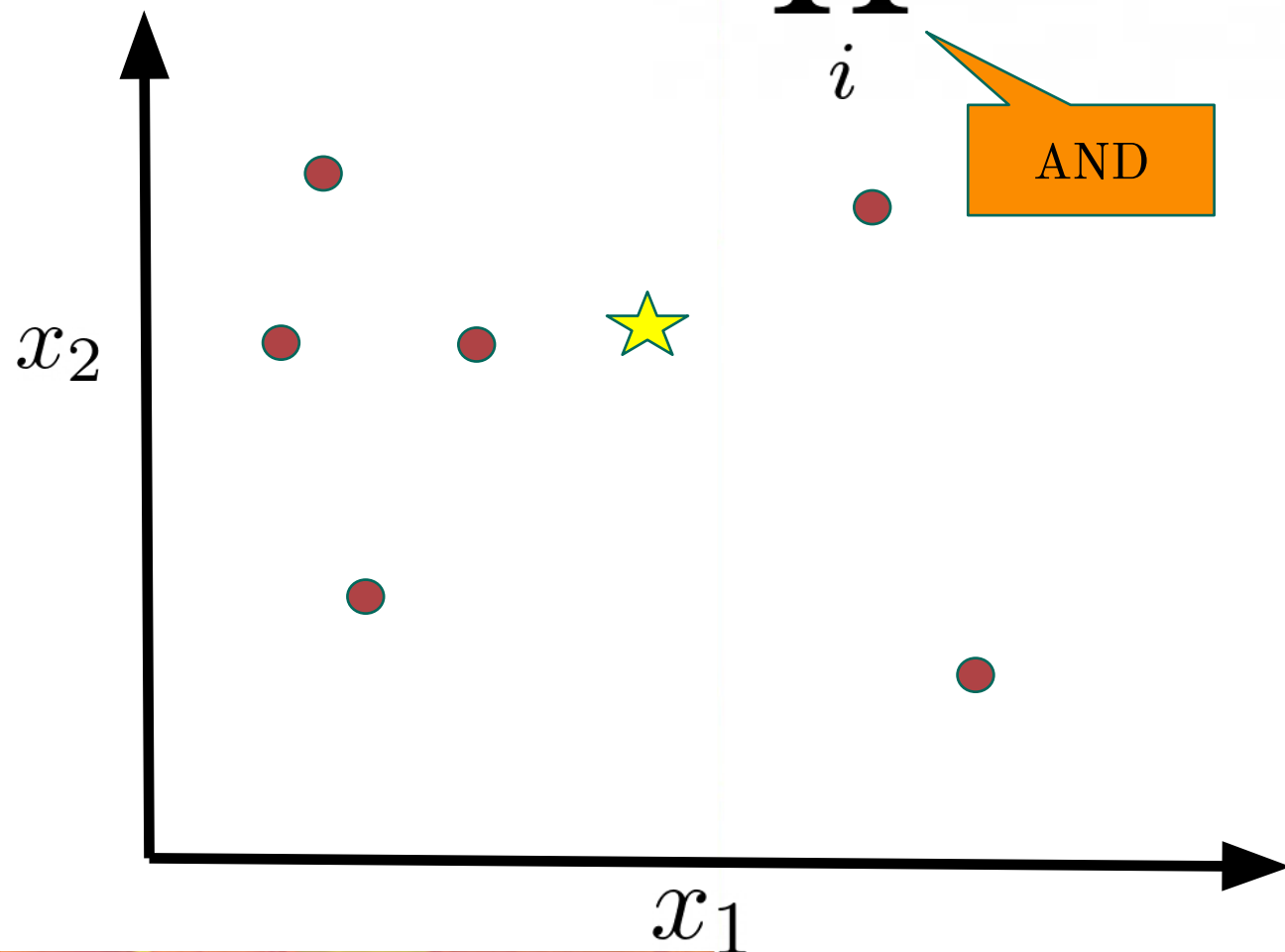
GPs for
high-dim
data?

$$k(\mathbf{x}, \mathbf{y}) = e^{-\frac{||\mathbf{x} - \mathbf{y}||^2}{2l^2}}$$
$$= \prod_i^d k_i(x_i, y_i)$$


AND

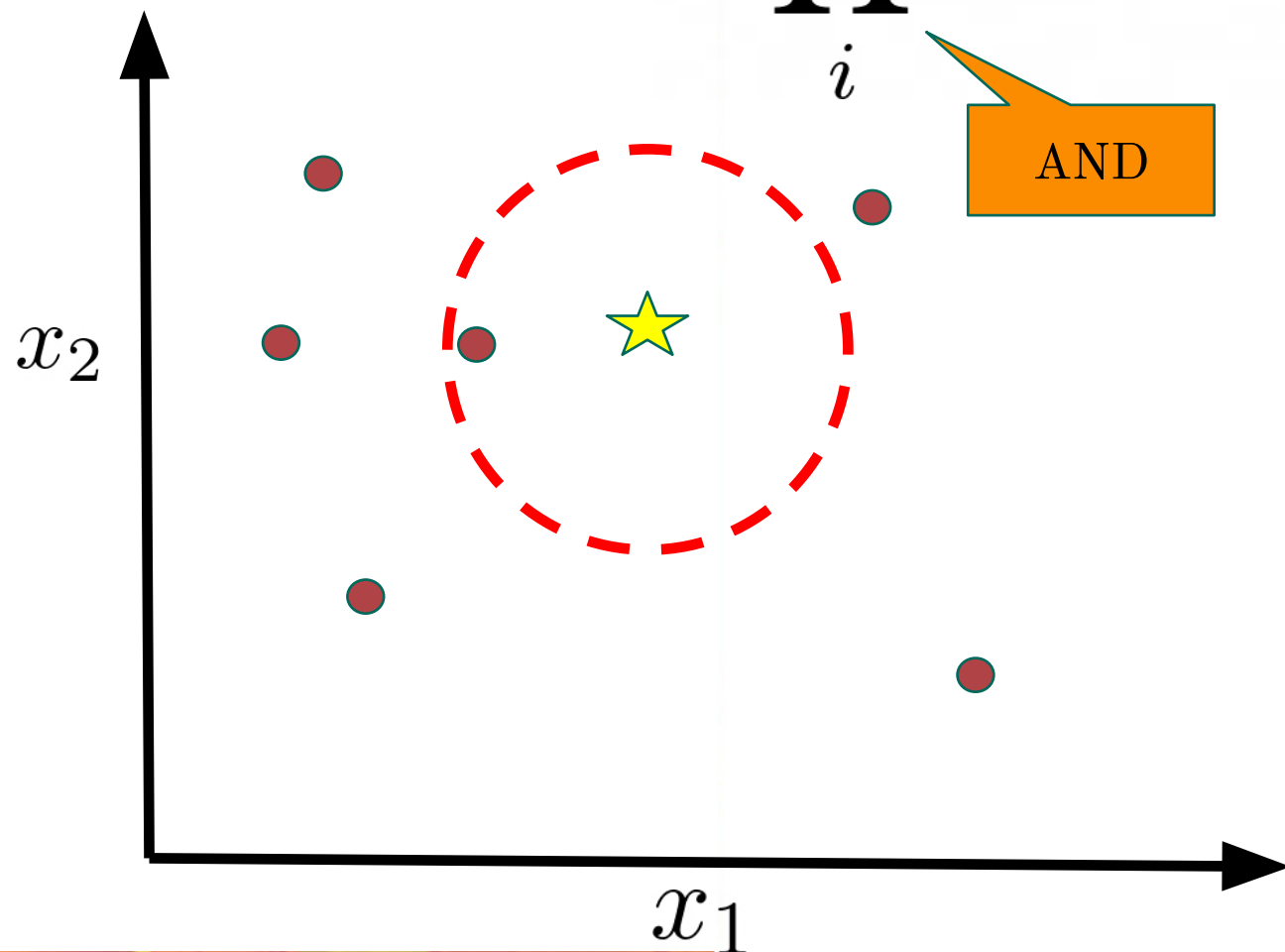
GPs for
high-dim
data?

$$k(\mathbf{x}, \mathbf{y}) = e^{-\frac{||\mathbf{x} - \mathbf{y}||^2}{2l^2}}$$
$$= \prod_i^d k_i(x_i, y_i)$$

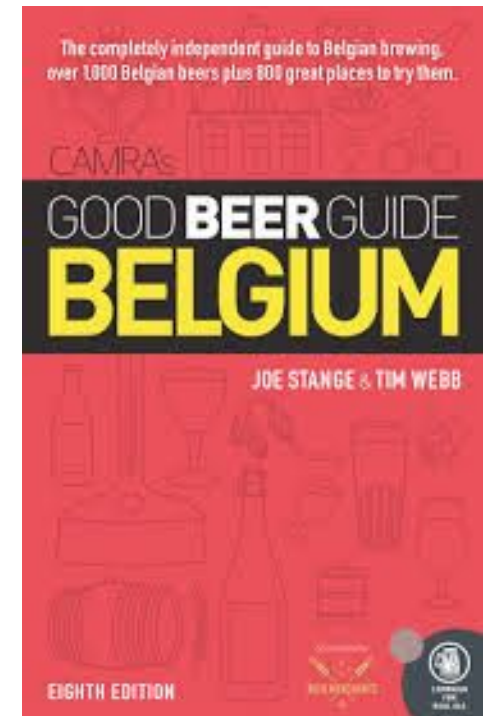


GPs for
high-dim
data?

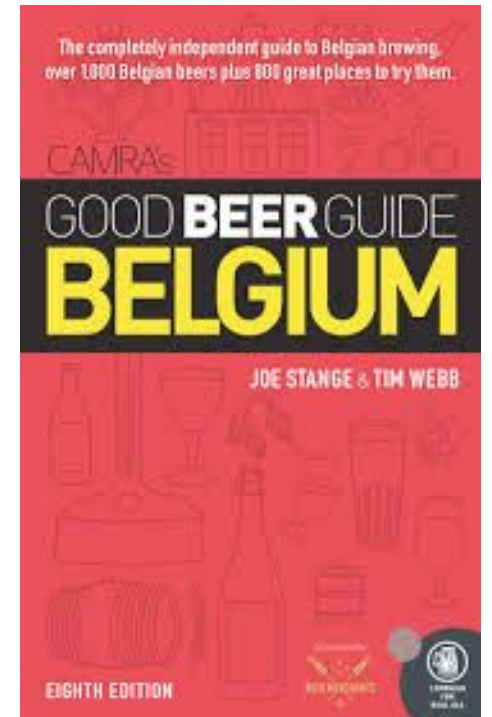
$$k(\mathbf{x}, \mathbf{y}) = e^{-\frac{||\mathbf{x} - \mathbf{y}||^2}{2l^2}}$$
$$= \prod_i^d k_i(x_i, y_i)$$



GPs for high-dim data?

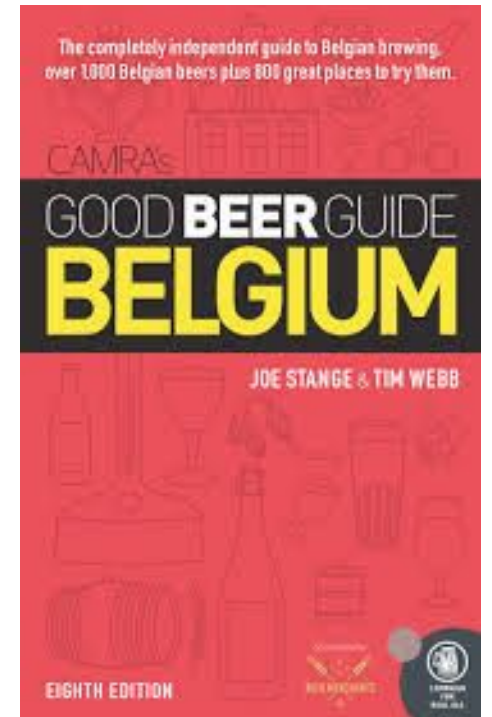


GPs for high-dim data?




GPs for high-dim data?

- Type of fermentation (wild yeast?)
- Ingredients (orange peel?, coriander??????)
- Strength
- Brewed by a monk?
- Barrel-aged?





GPs for
high-dim
data?

$$k(\mathbf{x}, \mathbf{y}) = e^{-\frac{||\mathbf{x} - \mathbf{y}||^2}{2l^2}}$$
$$= \prod_i^d k_i(x_i, y_i)$$


AND

GPs for high-dim data?

$$k(\mathbf{x}, \mathbf{y}) = e^{-\frac{\|\mathbf{x} - \mathbf{y}\|^2}{2l^2}}$$
$$= \prod_i^d k_i(x_i, y_i)$$


$$k_1(\mathbf{x}, \mathbf{y}) = \sum_i^d k_i(x_i, y_i)$$


$$k_2(\mathbf{x}, \mathbf{y}) = \sum_{i < j}^d k_i(x_i, y_i) k_j(x_j, y_j)$$

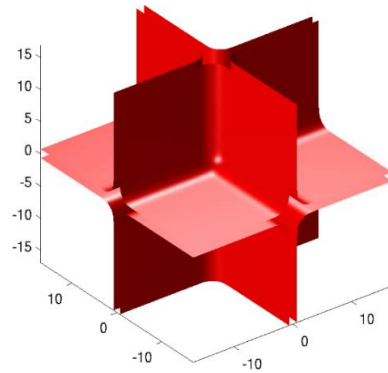
Additive Gaussian Processes

$$k(x, y) = \overset{0}{k_0} + \sum \overset{1}{k_i(x_i, y_i)} + \sum_{i < j} \overset{2}{k_i(x_i, y_i)k_j(x_j, y_j)}$$

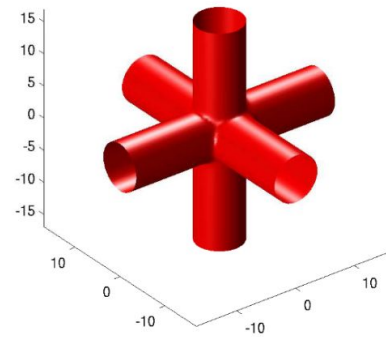
Additive Gaussian Processes

$$k(x, y) = \overset{0}{k_0} + \overset{1}{\sum k_i(x_i, y_i)} + \overset{2}{\sum_{i < j} k_i(x_i, y_i)k_j(x_j, y_j)}$$

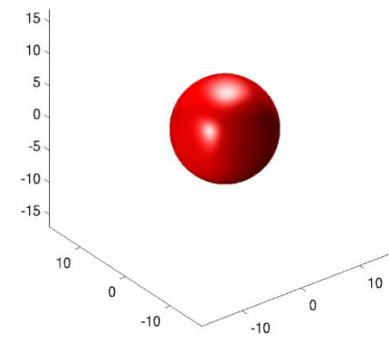
(Duvenaud et al 2011)



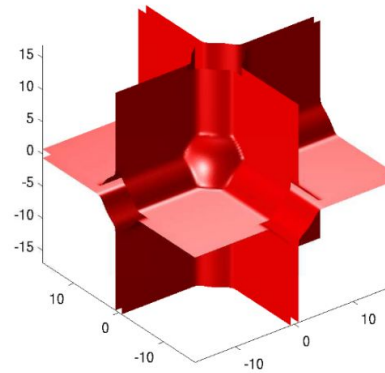
1st order interactions
 $k_1 + k_2 + k_3$



2nd order interactions
 $k_1k_2 + k_2k_3 + k_1k_3$



3rd order interactions
 $k_1k_2k_3$
(Squared-exp kernel)



All interactions
(Additive kernel)

Additive Gaussian Processes

$$\begin{aligned} k(x, y) &= \overset{0}{k_0} + \overset{1}{\sum k_i(x_i, y_i)} + \overset{2}{\sum_{i < j} k_i(x_i, y_i) k_j(x_j, y_j)} \\ &\quad \Updownarrow \\ f(\mathbf{x}) &= f_0 + \sum f_i(x_i) + \sum_{i < j} f_{ij}(x_i, x_j) \end{aligned}$$

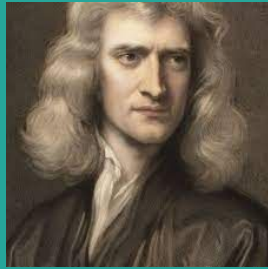
Ginsbourger et al. (2016)

Additive Gaussian Processes

$$\begin{aligned} k(x, y) &= \overset{0}{k_0} + \overset{1}{\sum} k_i(x_i, y_i) + \overset{2}{\sum_{i < j}} k_i(x_i, y_i) k_j(x_j, y_j) \\ &\quad \Updownarrow \\ f(\mathbf{x}) &= f_0 + \sum f_i(x_i) + \sum_{i < j} f_{ij}(x_i, x_j) \end{aligned}$$

- Standard RBF $\rightarrow O(d(N^2 + NM))$
- d additive RBF $\rightarrow O(2^d(N^2 + NM))$

Additive Gaussian Processes



- Newton Girard (Duvenaud et al 2011)

$$\begin{aligned} k(x, y) &= \overset{0}{k_0} + \sum \overset{1}{k_i(x_i, y_i)} + \sum \overset{2}{k_i(x_i, y_i)k_j(x_j, y_j)} \\ &\quad \Updownarrow \\ f(\mathbf{x}) &= f_0 + \sum f_i(x_i) + \sum f_{ij}(x_i, x_j) \end{aligned}$$

- Standard RBF $\rightarrow O(d(N^2 + NM))$
- d additive RBF $\rightarrow O(2^d(N^2 + NM))$
- d additive BBF (NG) $\rightarrow O(d^2(N^2 + NM))$

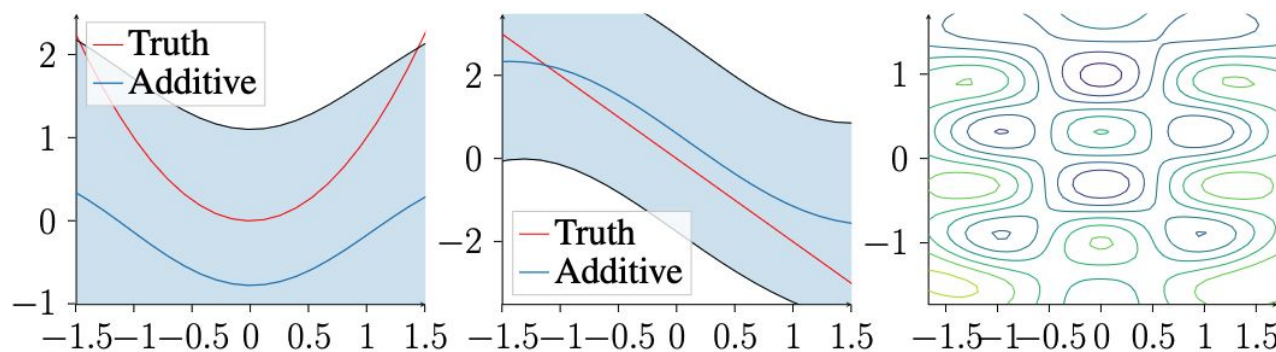
Additive Gaussian Processes



$$\begin{aligned}
 k(x, y) &= \overset{0}{k_0} + \overset{1}{\sum k_i(x_i, y_i)} + \overset{2}{\sum_{i < j} k_i(x_i, y_i) k_j(x_j, y_j)} \\
 &\quad \Updownarrow \\
 f(\mathbf{x}) &= f_0 + \sum f_i(x_i) + \sum_{i < j} f_{ij}(x_i, x_j)
 \end{aligned}$$

Ginsbourger et al. (2016)

$$f(x_1, x_2) = x_1^2 - 2x_2 + \cos(3x_1)\sin(5x_2)$$



$$E[f_i(x_i) | \mathcal{D}] = \overset{(a) f_1}{k_i(x_i, X)} \overset{(b) f_2}{K(X, X)^{-1}} \overset{(c) \text{Interaction}}{\mathbf{y}}$$

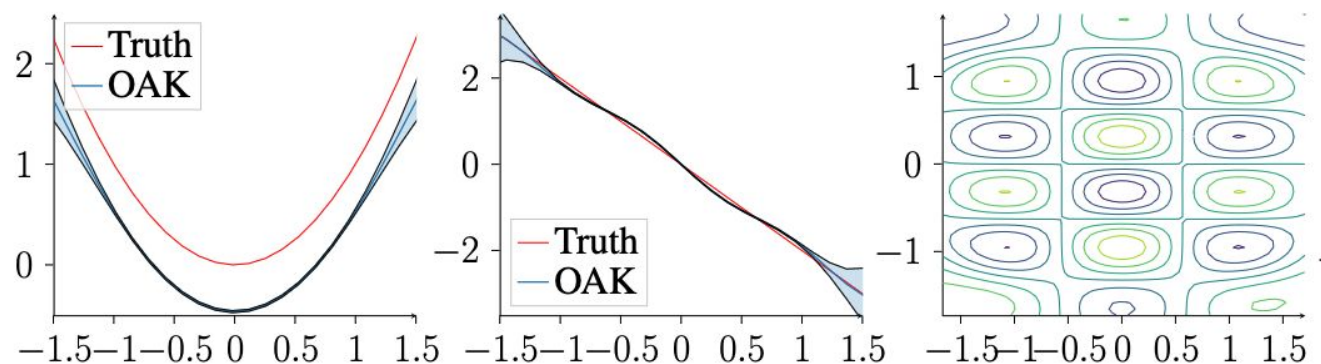
Additive Gaussian Processes

- **Orthogonalise** (Durrande et al 2012)

$$f(x_1, x_2) = (f_1(x_1) + \delta) + (f_2(x_2) - \delta)$$

$$k(x, y) = \overset{0}{k_0} + \overset{1}{\sum k_i(x_i, y_i)} + \overset{2}{\sum_{i < j} k_i(x_i, y_i) k_j(x_j, y_j)}$$
$$\Updownarrow$$
$$f(\mathbf{x}) = f_0 + \sum f_i(x_i) + \sum_{i < j} f_{ij}(x_i, x_j)$$

$$f(x_1, x_2) = x_1^2 - 2x_2 + \cos(3x_1)\sin(5x_2)$$



Additive Gaussian Processes

- **Orthogonalise** (Durrande et al 2012)

$$f(x_1, x_2) = (f_1(x_1) + \delta) + (f_2(x_2) - \delta)$$

- By conditioning

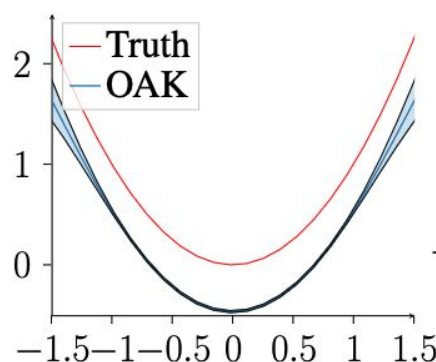
$$f_i(x_i) \Big| \int f_i(x_i) p(x_i) dx_i = 0$$

$$k(x, y) = \overset{0}{k_0} + \overset{1}{\sum k_i(x_i, y_i)} + \overset{2}{\sum_{i < j} k_i(x_i, y_i) k_j(x_j, y_j)}$$

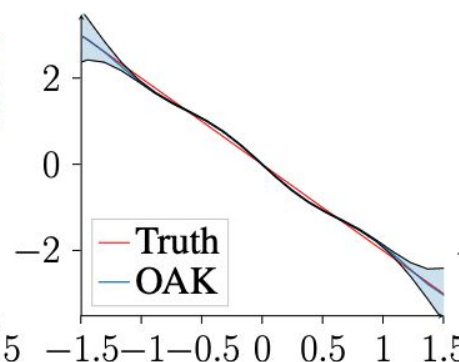
$$\Updownarrow$$

$$f(\mathbf{x}) = f_0 + \sum f_i(x_i) + \sum_{i < j} f_{ij}(x_i, x_j)$$

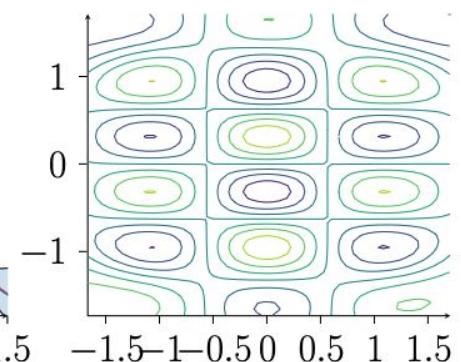
$$f(x_1, x_2) = x_1^2 - 2x_2 + \cos(3x_1)\sin(5x_2)$$



(f) f_1



(g) f_2



(h) Interaction

Additive Gaussian Processes

- **Orthogonalise** (Durrande et al 2012)

$$f(x_1, x_2) = (f_1(x_1) + \delta) + (f_2(x_2) - \delta)$$

- By conditioning

$$f_i(x_i) \Big| \int f_i(x_i) p(x_i) dx_i = 0$$

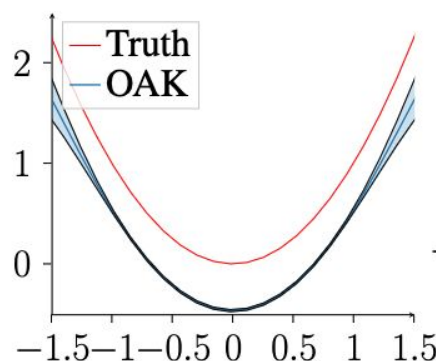
$k(x, y)$

This model is quite interpretable.....

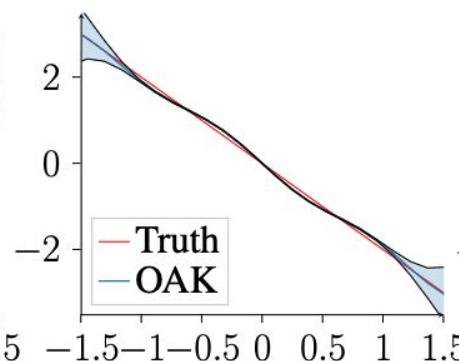
$k(x, y)$

$$f(\mathbf{x}) = f_0 + \sum f_i(x_i) + \sum_{i < j} f_{ij}(x_i, x_j)$$

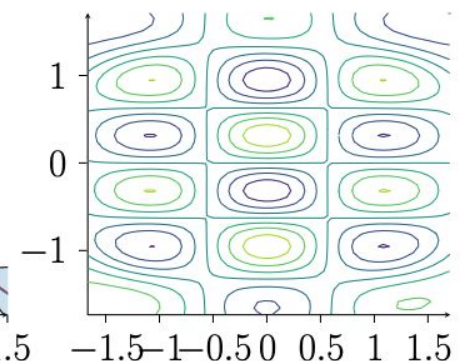
$$f(x_1, x_2) = x_1^2 - 2x_2 + \cos(3x_1)\sin(5x_2)$$



(f) f_1



(g) f_2



(h) Interaction

Additive Gaussian Processes

- **Orthogonalise** (Durrande et al 2012)

$$f(x_1, x_2) = (f_1(x_1) + \delta) + (f_2(x_2) - \delta)$$

- By conditioning

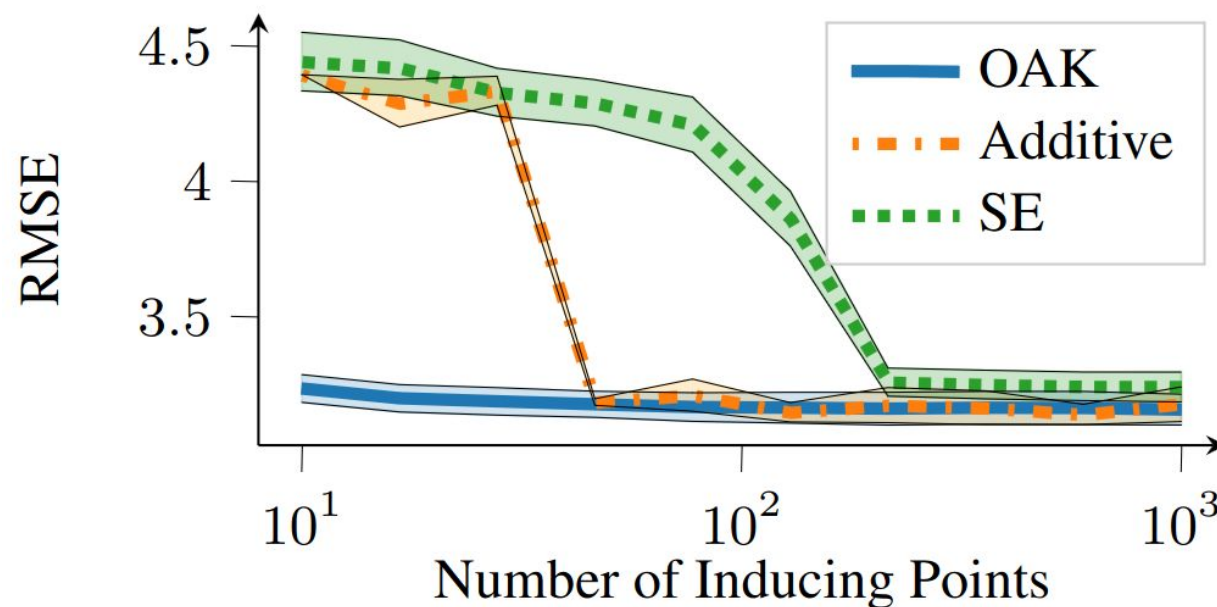
$$f_i(x_i) \Big| \int f_i(x_i) p(x_i) dx_i = 0$$

$k(x, y)$

This model is quite interpretable.....

$k(x_i, y_j)$

$$f(\mathbf{x}) = f_0 + \sum f_i(x_i) + \sum_{i < j} f_{ij}(x_i, x_j)$$



Thanks for listening



**UNIVERSITY OF
CAMBRIDGE**

**Lancaster
University**

