# Conditioning by adaptive sampling for robust design

**David Brookes**

Biophysics Graduate Group

University California, Berkeley
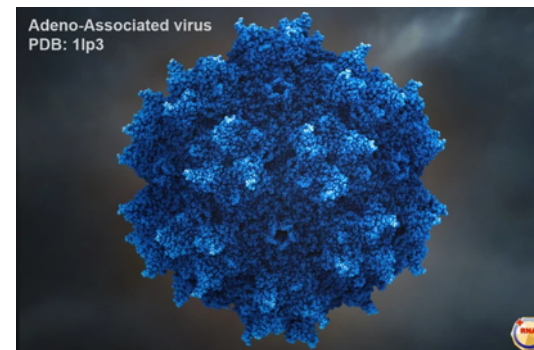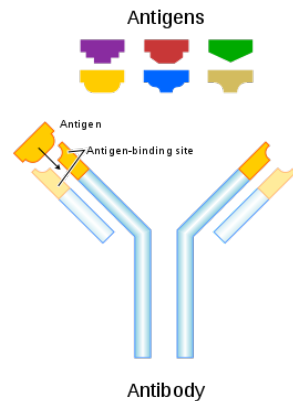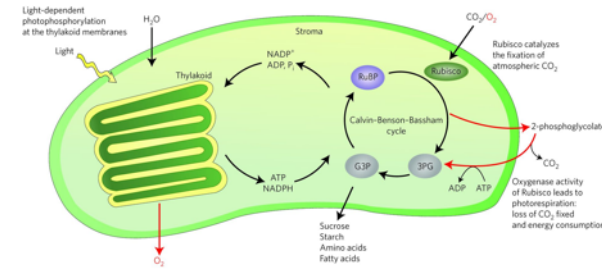
Jennifer Listgarten

EECS and Center for Computational Biology

University California, Berkeley

# Motivating problem: design protein sequences

- Proteins are made up of sequences of amino acids (20 possibilities)
- Huge variety of proteins whose function we would like to improve





Antigens

Antigen

Antigen-binding site

Antibody

Adeno-Associated virus
PDB: 1lp3

# Motivating problem: design protein sequences

- Proteins are made up of sequences of amino acids (20 possibilities)
- Huge variety of proteins whose function we would like to improve



Proteins that fluoresce





Antigens

Antigen
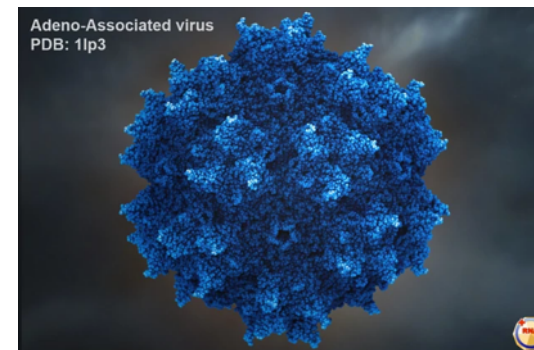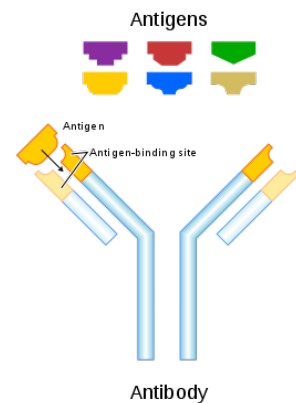
Antigen-binding site

Antibody

# Motivating problem: design protein sequences

- Proteins are made up of sequences of amino acids (20 possibilities)
- Huge variety of proteins whose function we would like to improve



Proteins that fluoresce
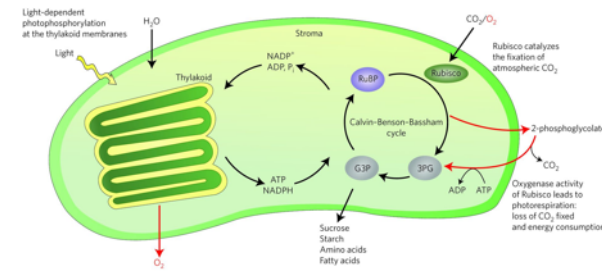




... that act as drugs

# Motivating problem: design protein sequences

- Proteins are made up of sequences of amino acids (20 possibilities)
- Huge variety of proteins whose function we would like to improve



Proteins that fluoresce



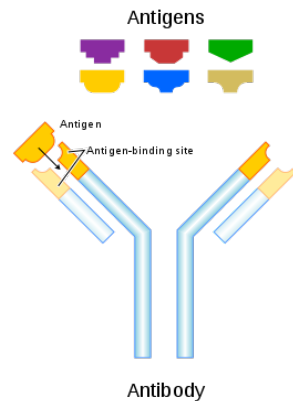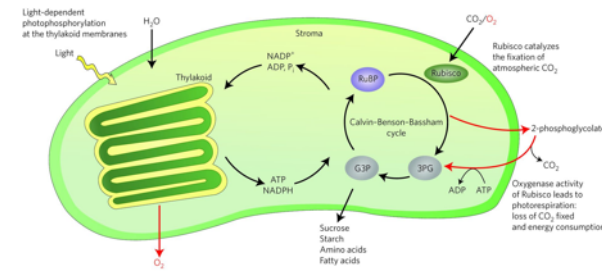… that fixate carbon in the atmosphere



… that act as drugs

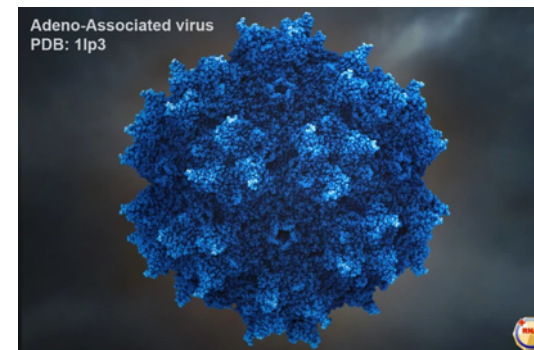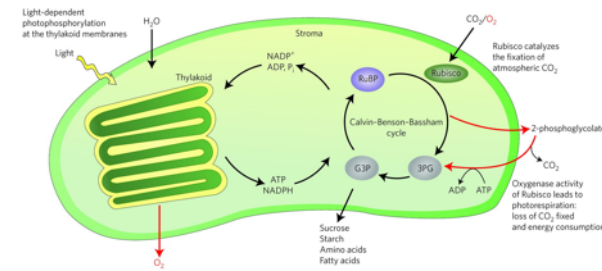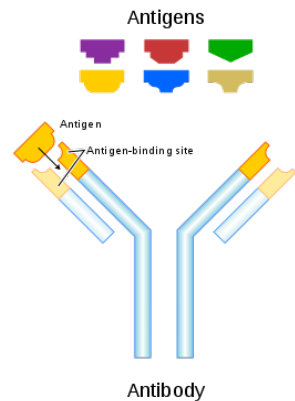# Motivating problem: design protein sequences

- Proteins are made up of sequences of amino acids (20 possibilities)
- Huge variety of proteins whose function we would like to improve



Proteins that fluoresce



... that fixate carbon in the atmosphere



.... that act as drugs



... that deliver gene-editing tools to tissues

# How to map sequence to function?

A law of molecular biology:

Sequence           Structure           Function



ex: fluorescence

Hughes A, Mort M, Carlisle F, *et al* B04 Alternative Splicing In Htt *Journal of Neurology, Neurosurgery & Psychiatry* 2014;**85**:A10.
http://www.rcsb.org/structure/6FWW

# Bypassing the structure relationships

A law of molecular biology:

Sequence         Structure         Function



High throughput experiments (& ML)

Hughes A, Mort M, Carlisle F, *et al* B04 Alternative Splicing In Htt *Journal of Neurology, Neurosurgery & Psychiatry* 2014;**85**:A10.
http://www.rcsb.org/structure/6FWW

# Can we solve the inverse problem?

A law of molecular biology:



Sequence        Structure        Function

Design problem: Given a model, find sequences with desired function

Hughes A, Mort M, Carlisle F, *et al* B04 Alternative Splicing In Htt *Journal of Neurology, Neurosurgery & Psychiatry* 2014;**85**:A10.
http://www.rcsb.org/structure/6FWW

# Why is protein design difficult?

- Huge, rugged search space

$\implies$ size scales as $20^L$

# Why is protein design difficult?

- Huge, rugged search space

$\implies$ size scales as $20^L$

- Discrete search space (no gradients)

# Why is protein design difficult?

- Huge, rugged search space

$\implies$ size scales as $20^L$

- Discrete search space (no gradients)
- Uncertainty in predictor



https://livingthing.danmackinlay.name/gaussian_processes.html69

# Possible solution: model-based optimization (MBO)

Idea: replace the standard (hard) objective

$$\max_{\boldsymbol{x} \in \mathcal{X}} f(\boldsymbol{x})$$

e.g. the space of sequences

# Possible solution: model-based optimization (MBO)

Idea: replace the standard (hard) objective with a potentially easier one

$$\max_{\boldsymbol{x} \in \mathcal{X}} f(\boldsymbol{x}) \qquad \Longrightarrow \qquad \max_{\boldsymbol{\theta} \in \mathbb{R}^d} \mathbb{E}_{p(\boldsymbol{x}|\boldsymbol{\theta})}[f(\boldsymbol{x})]$$

the space of sequences

model over sequence space

# Possible solution: model-based optimization (MBO)

Idea: replace the standard (hard) objective with a potentially easier one

$$\max_{\boldsymbol{x} \in \mathcal{X}} f(\boldsymbol{x}) \quad \Longrightarrow \quad \max_{\boldsymbol{\theta} \in \mathbb{R}^d} \mathbb{E}_{p(\boldsymbol{x}|\boldsymbol{\theta})}[f(\boldsymbol{x})]$$

Solution approach is to iterate:

1. Sample from "search model" $p(x|\theta)$

2. Evaluate samples on $f(x)$

3. Adjust $\theta$ so the model favors samples with large function evals

# Possible solution: model-based optimization (MBO)

Idea: replace the standard (hard) objective with a potentially easier one

$$\max_{\boldsymbol{x} \in \mathcal{X}} f(\boldsymbol{x}) \qquad \Longrightarrow \qquad \max_{\boldsymbol{\theta} \in \mathbb{R}^d} \mathbb{E}_{p(\boldsymbol{x}|\boldsymbol{\theta})}[f(\boldsymbol{x})]$$

Solution approach is to iterate:

1.  Sample from "search model" $p(x|\theta)$
2.  Evaluate samples on $f(x)$
3.  Adjust $\theta$ so the model favors sequences with large function evals

✓ Model can sample broad areas of sequence space
✓ Does not require gradients of $f$
✓ Can incorporate uncertainty

# First attempt at MBO for protein design:
## Design by Adaptive Sampling (DbAS)

Our aim is solve the MBO objective:

$$\underset{\boldsymbol{\theta}}{\operatorname{argmax}} \log \mathbb{E}_{p(\mathbf{x}|\boldsymbol{\theta})} \left[ P(S|\mathbf{x}) \right]$$

# First attempt at MBO for protein design: Design by Adaptive Sampling (DbAS)

Our aim is solve the MBO objective:

$$\underset{\boldsymbol{\theta}}{\operatorname{argmax}} \log \mathbb{E}_{p(\mathbf{x}|\boldsymbol{\theta})} \left[ P(S|\mathbf{x}) \right]$$

where

- $p(x|\theta)$ is the search model (VAE, HMM...)

# First attempt at MBO for protein design: Design by Adaptive Sampling (DbAS)

Our aim is solve the MBO objective:

$$\underset{\boldsymbol{\theta}}{\arg\max} \log \mathbb{E}_{p(\mathbf{x}|\boldsymbol{\theta})}\left[P(S|\mathbf{x})\right]$$

where

- $p(x|\theta)$ is the search model (VAE, HMM…)
- $S$ is desired set of property values

    → *e.g.* fluorescence $> \alpha$

# First attempt at MBO for protein design:
## Design by Adaptive Sampling (DbAS)

Our aim is solve the MBO objective:

$$\underset{\boldsymbol{\theta}}{\operatorname{argmax}} \log \mathbb{E}_{p(\mathbf{x}|\boldsymbol{\theta})} \left[ P(S|\mathbf{x}) \right]$$

where

- $p(x|\theta)$ is the search model (VAE, HMM…)

- $S$ is desired set of property values

  → *e.g.* fluorescence $> \alpha$

- $P(S|x)$ is a stochastic predictive model ("oracle") that maps sequences to property

# Design by Adaptive Sampling (cont.)

**Two issues:**
1. $\theta$ **is in the expectation distribution.**

$$\underset{\boldsymbol{\theta}}{\operatorname{argmax}} \log \mathbb{E}_{p(\mathbf{x}|\boldsymbol{\theta})} \left[ P(S|\mathbf{x}) \right]$$

# Design by Adaptive Sampling (cont.)

*maximize a lower bound*

**Two issues:**

1. ~~$\theta$ is in the expectation distribution.~~

$$\underset{\boldsymbol{\theta}}{\arg\max} \log \mathbb{E}_{p(\mathbf{x}|\boldsymbol{\theta})}\left[P(S|\mathbf{x})\right],$$

$\geq$

$$\underset{\boldsymbol{\theta}}{\arg\max} \mathbb{E}_{p(\mathbf{x}|\boldsymbol{\theta}^{(t)})}\left[P(S|\mathbf{x}) \log p(\mathbf{x}|\boldsymbol{\theta})\right]$$

# Design by Adaptive Sampling (cont.)

**Two issues:**

1. ~~$\theta$ is in the expectation distribution.~~

2. MC estimates for rare events.

*maximize a lower bound*

$$\underset{\theta}{\mathrm{argmax}}\ \log \mathbb{E}_{p(\mathbf{x}|\theta)}\left[P(S|\mathbf{x})\right],$$

$\geq$

$$\underset{\theta}{\mathrm{argmax}}\ \mathbb{E}_{p(\mathbf{x}|\theta^{(t)})}\left[P(S|\mathbf{x})\log p(\mathbf{x}|\theta)\right]$$

# Design by Adaptive Sampling (cont.)

**Two issues:**

1.  ~~$\theta$ is in the expectation distribution.~~

2.  ~~MC estimates for rare events.~~

*maximize a lower bound*

$$\operatorname*{argmax}_{\theta} \log \mathbb{E}_{p(\mathbf{x}|\theta)} \left[ P(S|\mathbf{x}) \right],$$

$$\downarrow \quad \geq$$

$$\operatorname*{argmax}_{\theta} \mathbb{E}_{p(\mathbf{x}|\theta^{(t)})} \left[ P(S|\mathbf{x}) \log p(\mathbf{x}|\theta) \right]$$

*anneal a sequence of relaxations:*
$S^t \rightarrow S$, *where* $S^t \supset S^{t+1}$

# Design by Adaptive Sampling (cont.)

*maximize a lower bound*

**Two issues:**

1. ~~$\theta$ is in the expectation distribution.~~

2. ~~MC estimates for rare events.~~

$$\underset{\theta}{\arg\max}\ \log \mathbb{E}_{p(\mathbf{x}|\theta)}\left[P(S|\mathbf{x})\right],$$

$\Downarrow$ $\geq$

$$\underset{\theta}{\arg\max}\ \mathbb{E}_{p(\mathbf{x}|\theta^{(t)})}\left[P(S|\mathbf{x})\log p(\mathbf{x}|\theta)\right]$$

$\Downarrow$ Anneal and MC

$$\theta^{(t+1)} = \underset{\theta}{\arg\max}\ \sum_{i=1}^{M} P(S^{(t)}|\mathbf{x}_i^{(t)})\log p(\mathbf{x}_i^{(t)}|\theta)$$

# Design by Adaptive Sampling (cont.)

*maximize a lower bound*

**Two issues:**

1. $\theta$ is in th~~e~~ ~~distribut~~

2. ~~MC estimates for rare~~ ~~events.~~

$$\arg\max_{\theta} \log \mathbb{E}_{p(\mathbf{x}|\theta)}[P(S|\mathbf{x})],$$

Assumes oracle is unbiased and has good uncertainty estimates

$$\max_{\theta} \mathbb{E}_{p(\mathbf{x}|\theta^{(t)})}\left[P(S|\mathbf{x})\log p(\mathbf{x}|\theta)\right]$$

↓ Anneal and MC

$$\boldsymbol{\theta}^{(t+1)} = \arg\max_{\boldsymbol{\theta}} \sum_{i=1}^{M} P(S^{(t)}|\mathbf{x}_i^{(t)}) \log p(\mathbf{x}_i^{(t)}|\boldsymbol{\theta})$$

# How pathological oracles lead you astray

# How pathological oracles lead you astray
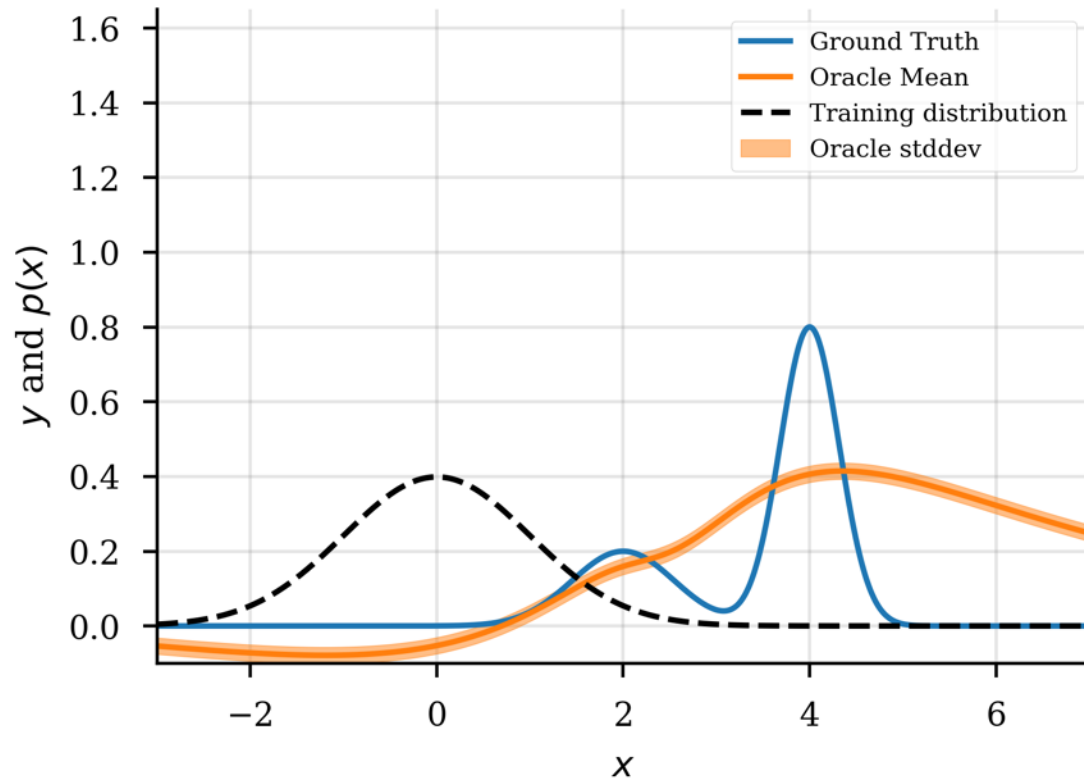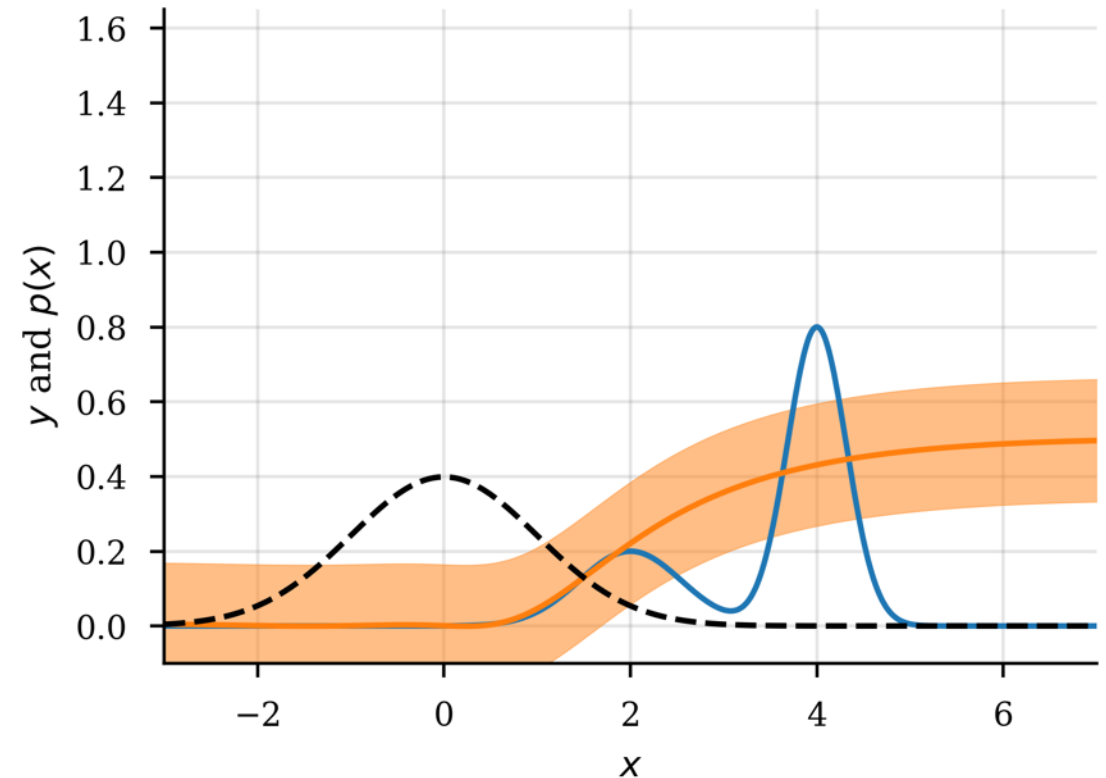
Acceptable



Many training examples

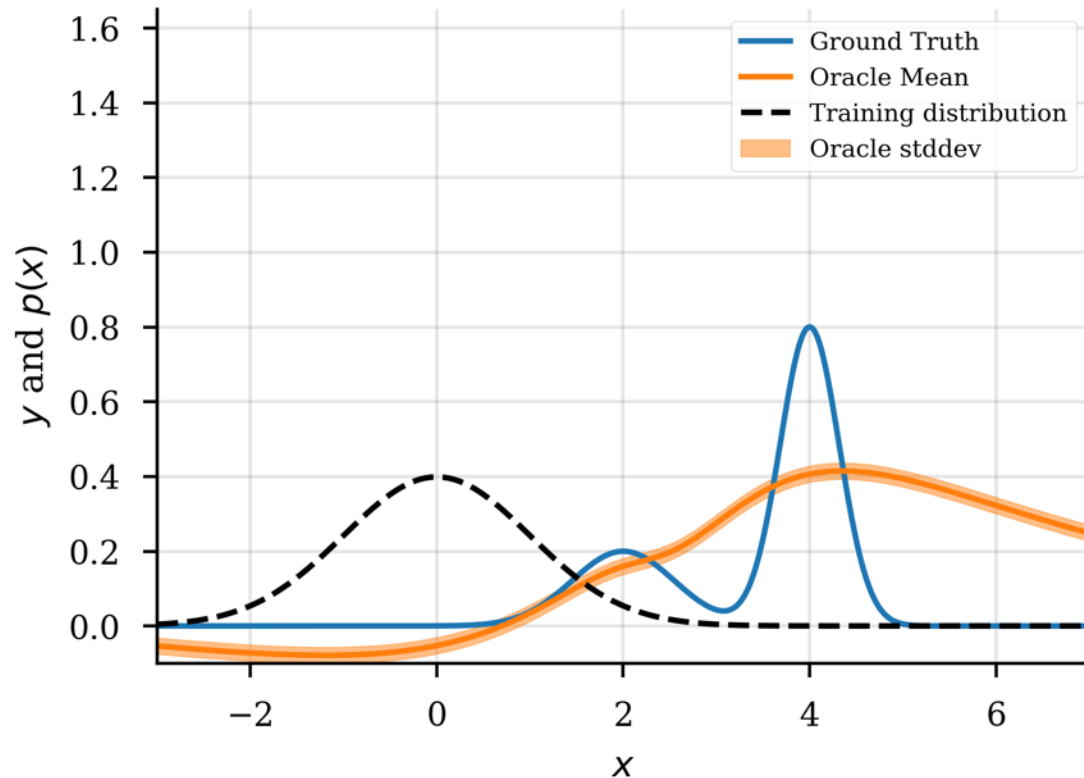# How pathological oracles lead you astray



Acceptable
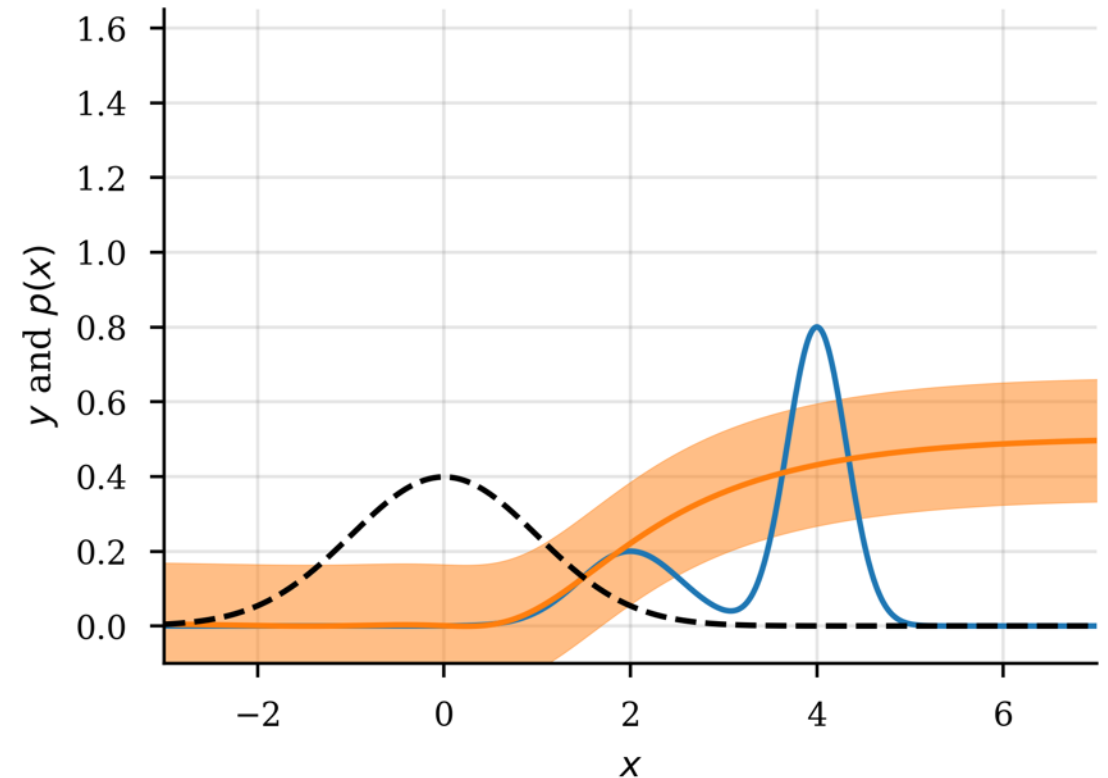
Pathological

Many training examples

Fewer training examples

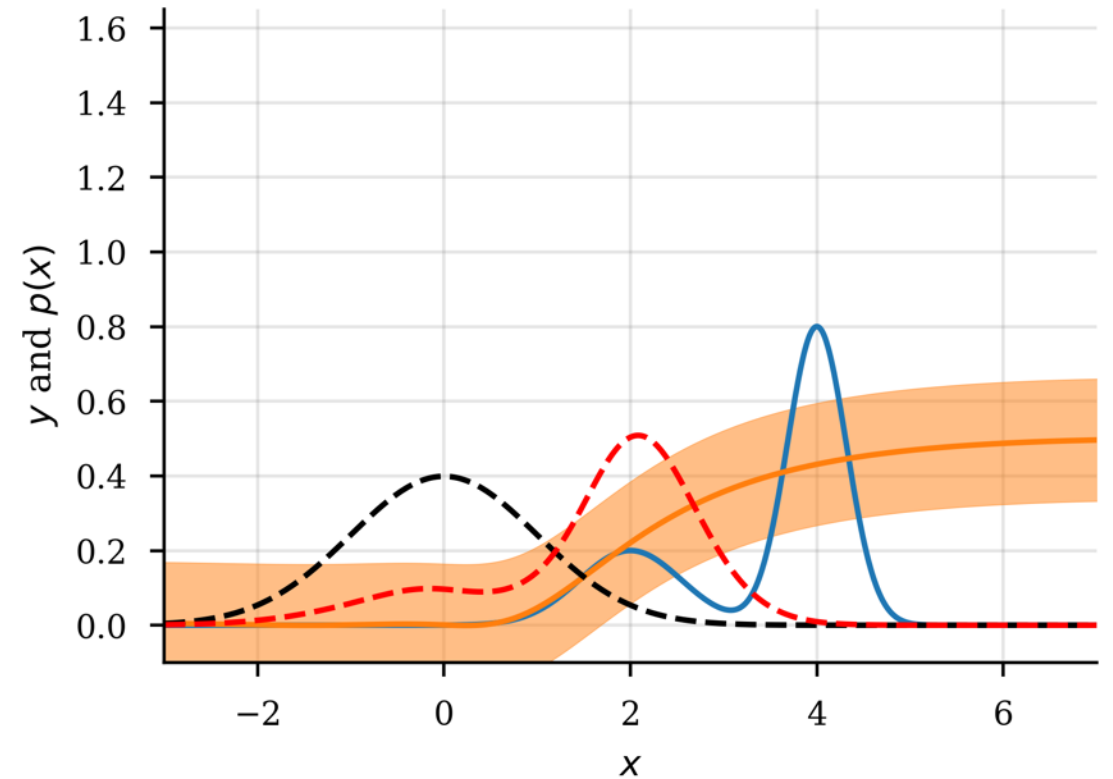# How pathological oracles lead you astray

Acceptable

Pathological



Idea: estimate training distribution of x *conditioned* on high values of oracle
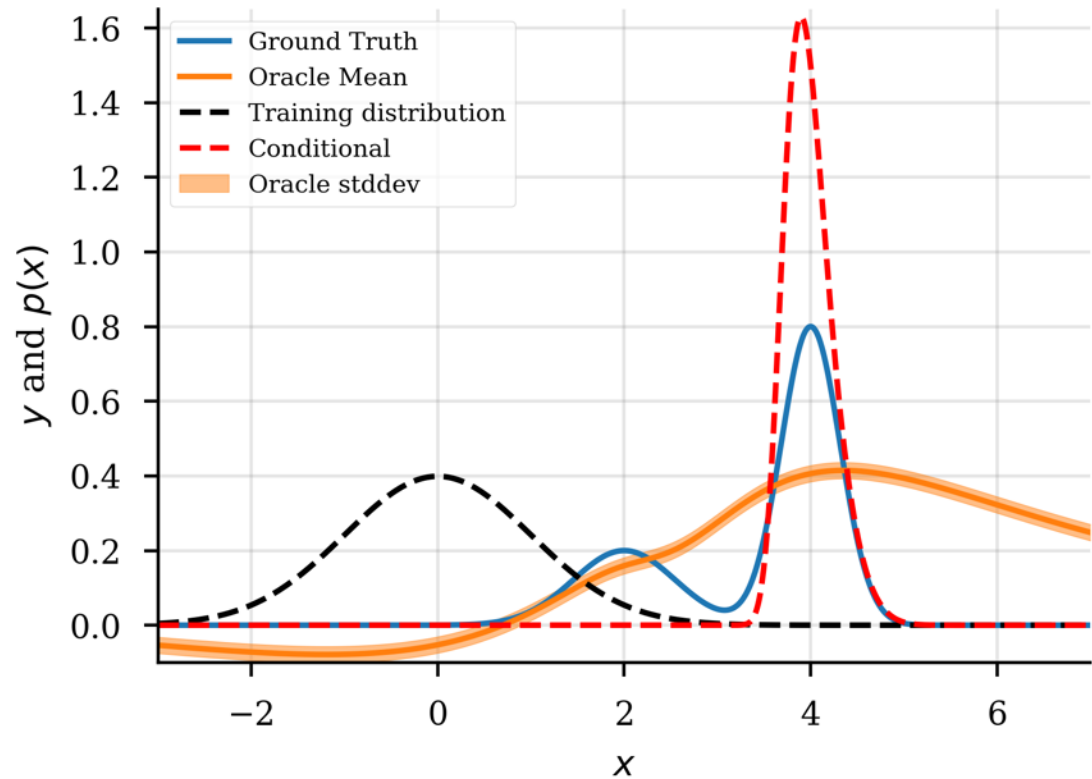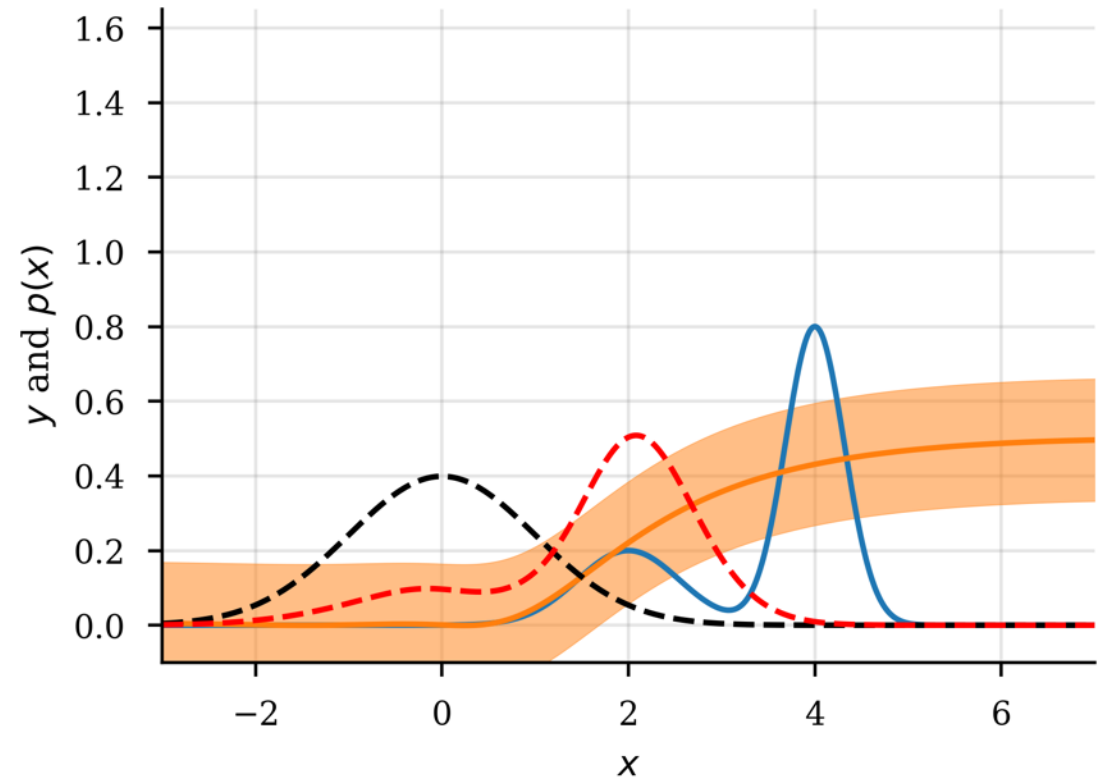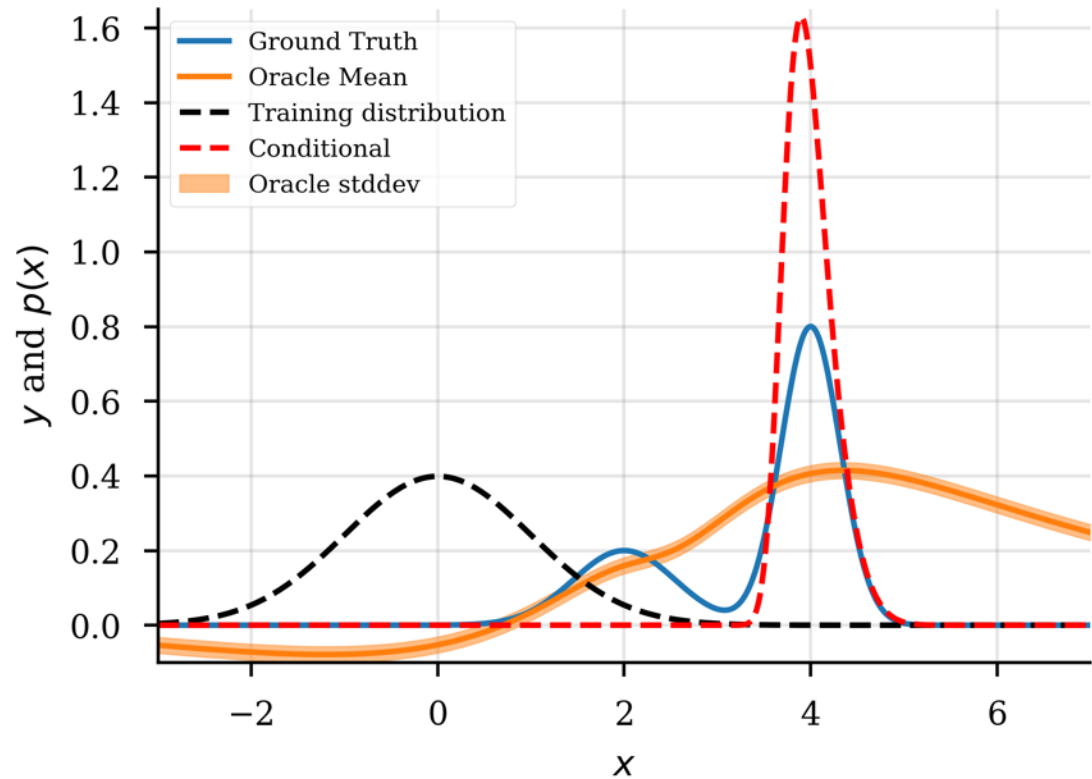
# Fixing pathological oracles w/ conditioning



Idea: estimate training distribution of x *conditioned* on high values of oracle

# Fixing pathological oracles w/ conditioning



Idea: estimate training distribution of x *conditioned* on high values of oracle

Don't have access to training distribution, but can build a model $p(\boldsymbol{x}|\boldsymbol{\theta}^{(0)})$ to approximate it

# Conditioning by Adaptive Sampling (CbAS)

**Previous formulation:**

$$\underset{\boldsymbol{\theta}}{\mathrm{argmax}} \log \mathbb{E}_{p(\mathbf{x}|\boldsymbol{\theta})} [P(S|\mathbf{x})]$$

$\downarrow \quad \geq$

$$\underset{\boldsymbol{\theta}}{\mathrm{argmax}} \mathbb{E}_{p(\mathbf{x}|\boldsymbol{\theta}^{(t)})} [P(S|\mathbf{x}) \log p(\mathbf{x}|\boldsymbol{\theta})]$$

$\downarrow$ Anneal and MC

$$\boldsymbol{\theta}^{(t+1)} = \underset{\boldsymbol{\theta}}{\mathrm{argmax}} \sum_{i=1}^{M} P(S^{(t)}|\mathbf{x}_i^{(t)}) \log p(\mathbf{x}_i^{(t)}|\boldsymbol{\theta})$$

**New formulation:**

$$\underset{\boldsymbol{\theta}}{\mathrm{argmin}} \, D_{KL} \left( p(\mathbf{x}|S, \boldsymbol{\theta}^{(0)}) || p(\mathbf{x}|\boldsymbol{\theta}) \right)$$

$p(x|\boldsymbol{\theta}^{(0)})$ models the training distribution

# Conditioning by Adaptive Sampling (CbAS)

**Previous formulation:**

$$\underset{\boldsymbol{\theta}}{\operatorname{argmax}} \log \mathbb{E}_{p(\mathbf{x}|\boldsymbol{\theta})} \left[ P(S|\mathbf{x}) \right]$$

$\downarrow \quad \geq$

$$\underset{\boldsymbol{\theta}}{\operatorname{argmax}} \mathbb{E}_{p(\mathbf{x}|\boldsymbol{\theta}^{(t)})} \left[ P(S|\mathbf{x}) \log p(\mathbf{x}|\boldsymbol{\theta}) \right]$$

$\downarrow$ Anneal and MC

$$\boldsymbol{\theta}^{(t+1)} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \sum_{i=1}^{M} P(S^{(t)}|\mathbf{x}_i^{(t)}) \log p(\mathbf{x}_i^{(t)}|\boldsymbol{\theta})$$

**New formulation:**

$$\underset{\boldsymbol{\theta}}{\operatorname{argmin}} D_{KL} \left( p(\mathbf{x}|S, \boldsymbol{\theta}^{(0)}) || p(\mathbf{x}|\boldsymbol{\theta}) \right)$$

$\downarrow \quad =$

$$\underset{\boldsymbol{\theta}}{\operatorname{argmax}} \mathbb{E}_{p(\mathbf{x}|\boldsymbol{\theta}^{(0)})} \left[ P(S|\mathbf{x}) \log p(\mathbf{x}|\boldsymbol{\theta}) \right]$$

# Conditioning by Adaptive Sampling (CbAS)

**Previous formulation:**

$$\underset{\boldsymbol{\theta}}{\arg\max} \log \mathbb{E}_{p(\mathbf{x}|\boldsymbol{\theta})}\left[P(S|\mathbf{x})\right]$$

$\Big\downarrow$    $\geq$

$$\underset{\boldsymbol{\theta}}{\arg\max} \mathbb{E}_{p(\mathbf{x}|\boldsymbol{\theta}^{(t)})}\left[P(S|\mathbf{x})\log p(\mathbf{x}|\boldsymbol{\theta})\right]$$

$\Big\downarrow$    Anneal and MC

$$\boldsymbol{\theta}^{(t+1)} = \underset{\boldsymbol{\theta}}{\arg\max} \sum_{i=1}^{M} P(S^{(t)}|\mathbf{x}_i^{(t)})\log p(\mathbf{x}_i^{(t)}|\boldsymbol{\theta})$$

**New formulation:**

$$\underset{\boldsymbol{\theta}}{\arg\min} D_{KL}\left(p(\mathbf{x}|S,\boldsymbol{\theta}^{(0)}) || p(\mathbf{x}|\boldsymbol{\theta})\right)$$

$\Big\downarrow$    $=$

$$\underset{\boldsymbol{\theta}}{\arg\max} \mathbb{E}_{p(\mathbf{x}|\boldsymbol{\theta}^{(0)})}\left[P(S|\mathbf{x})\log p(\mathbf{x}|\boldsymbol{\theta})\right]$$

Can't anneal when sampling dist. doesn't change!

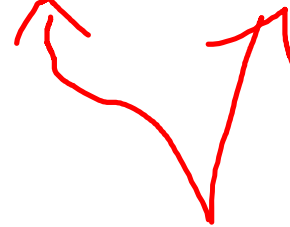# Conditioning by Adaptive Sampling (CbAS)

Previous formulation:

New formulation:

$$\operatorname*{argmin}_{\boldsymbol{\theta}} D_{KL}\left(p(\mathbf{x}|S, \boldsymbol{\theta}^{(0)})||p(\mathbf{x}|\boldsymbol{\theta})\right)$$

$\geq$

$=$

$$\operatorname*{argmax}_{\boldsymbol{\theta}} \mathbb{E}_{p(\mathbf{x}|\boldsymbol{\theta}^{(t)})}\left[P(S|\mathbf{x})\log p(\mathbf{x}|\boldsymbol{\theta})\right]$$

$$\operatorname*{argmax}_{\boldsymbol{\theta}} \mathbb{E}_{p(\mathbf{x}|\boldsymbol{\theta}^{(0)})}\left[P(S|\mathbf{x})\log p(\mathbf{x}|\boldsymbol{\theta})\right]$$

Anneal and MC

$=$

$$\operatorname*{argmax}_{\boldsymbol{\theta}} \mathbb{E}_{p(\mathbf{x}|\boldsymbol{\theta}^{(t)})}\left[\frac{p(\mathbf{x}|\boldsymbol{\theta}^{(0)})}{p(\mathbf{x}|\boldsymbol{\theta}^{(t)})}P(S|\mathbf{x})\log p(\mathbf{x}|\boldsymbol{\theta})\right]$$

$$\boldsymbol{\theta}^{(t+1)} = \operatorname*{argmax}_{\boldsymbol{\theta}} \sum_{i=1}^{M} P(S^{(t)}|\mathbf{x}_i^{(t)})\log p(\mathbf{x}_i^{(t)}|\boldsymbol{\theta})$$

Importance sampling proposal dist.

# Conditioning by Adaptive Sampling (CbAS)

Previous formulation:

$$\operatorname*{argmax}_{\boldsymbol{\theta}} \log \mathbb{E}_{p(\mathbf{x}|\boldsymbol{\theta})} \left[ P(S|\mathbf{x}) \right]$$

$\geq$

$$\operatorname*{argmax}_{\boldsymbol{\theta}} \mathbb{E}_{p(\mathbf{x}|\boldsymbol{\theta}^{(t)})} \left[ P(S|\mathbf{x}) \log p(\mathbf{x}|\boldsymbol{\theta}) \right]$$

Anneal and MC

$$\boldsymbol{\theta}^{(t+1)} = \operatorname*{argmax}_{\boldsymbol{\theta}} \sum_{i=1}^{M} P(S^{(t)}|\mathbf{x}_i^{(t)}) \log p(\mathbf{x}_i^{(t)}|\boldsymbol{\theta})$$

New formulation:

$$\operatorname*{argmin}_{\boldsymbol{\theta}} D_{KL} \left( p(\mathbf{x}|S, \boldsymbol{\theta}^{(0)}) || p(\mathbf{x}|\boldsymbol{\theta}) \right)$$

$=$

$$\operatorname*{argmax}_{\boldsymbol{\theta}} \mathbb{E}_{p(\mathbf{x}|\boldsymbol{\theta}^{(0)})} \left[ P(S|\mathbf{x}) \log p(\mathbf{x}|\boldsymbol{\theta}) \right]$$

$=$

$$\operatorname*{argmax}_{\boldsymbol{\theta}} \mathbb{E}_{p(\mathbf{x}|\boldsymbol{\theta}^{(t)})} \left[ \frac{p(\mathbf{x}|\boldsymbol{\theta}^{(0)})}{p(\mathbf{x}|\boldsymbol{\theta}^{(t)})} P(S|\mathbf{x}) \log p(\mathbf{x}|\boldsymbol{\theta}) \right]$$
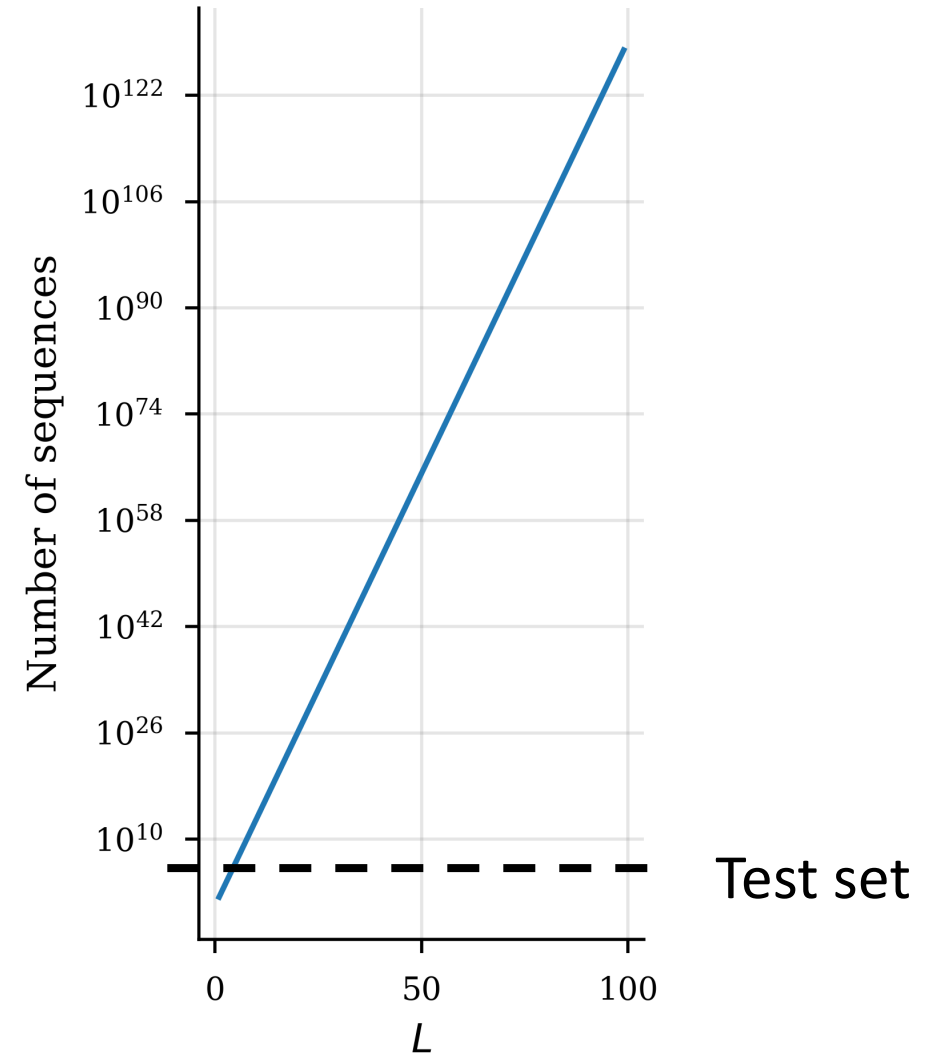
Anneal and MC

$$\operatorname*{argmax}_{\boldsymbol{\theta}} \sum_{i=1}^{M} \frac{p(\mathbf{x}_i^{(t)}|\boldsymbol{\theta}^{(0)})}{p(\mathbf{x}_i^{(t)}|\boldsymbol{\theta}^{(t)})} P(S^{(t)}|\mathbf{x}_i^{(t)}) \log p(\mathbf{x}_i^{(t)}|\boldsymbol{\theta})$$

# Testing is fundamentally different

- We don't trust our oracle and generally can't query the ground truth
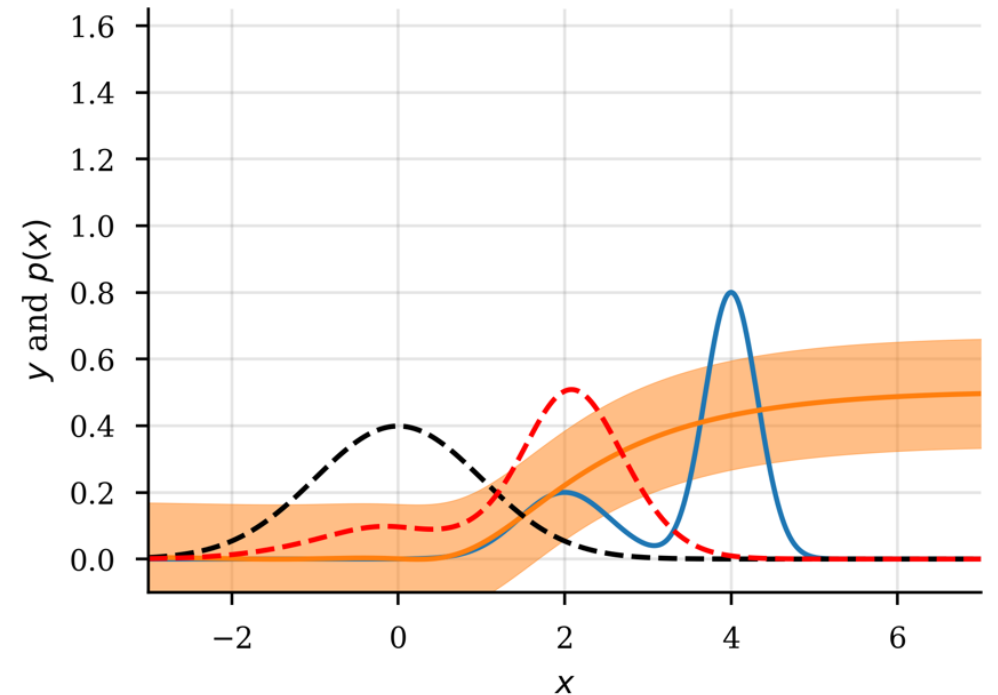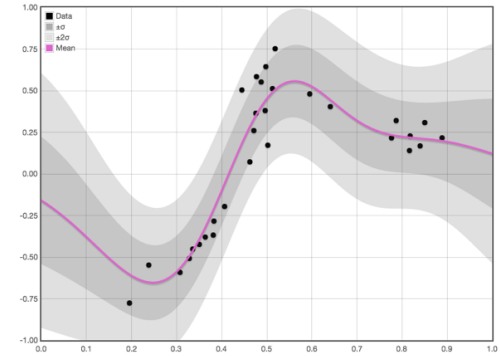
# Testing is fundamentally different

- We don't trust our oracle and generally can't query the ground truth

- We can't hold-out a test set of good sequences
  - Near-zero chance of any of these sequences being found by the method

# Testing is fundamentally different

- We don't trust our oracle and generally can't query the ground truth

- We can't hold-out a test set of good sequences
  - Near-zero chance of any of these sequences being found by the method

- We can't use some canonical test function as the oracle
  - In our problem it is untrustworthy

# Testing strategy

- Simulate a ground truth based on real data

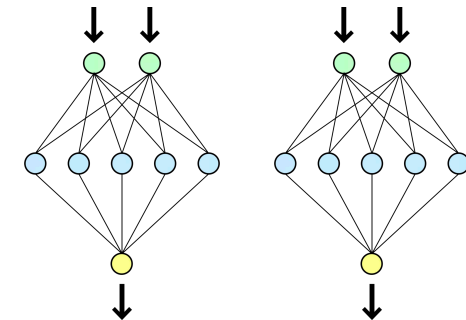  → "Ground truth" is a GP mean function



Ground truth GP

# Testing strategy

- Simulate a ground truth based on real data
  - → "Ground truth" is a GP mean function
- Ground truth vales values are sampled from the GP for given sequences
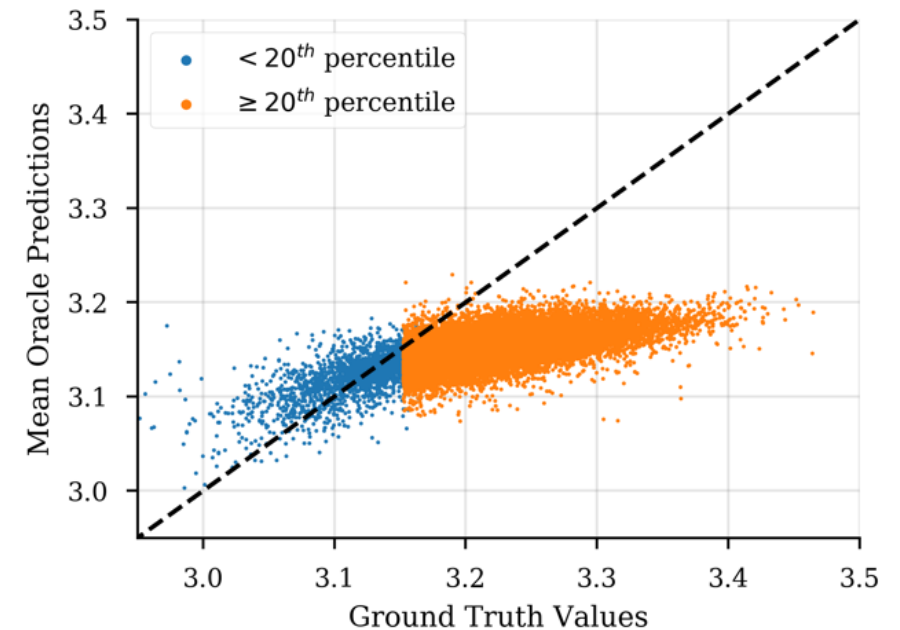- Use these input-output pairs to train oracles.
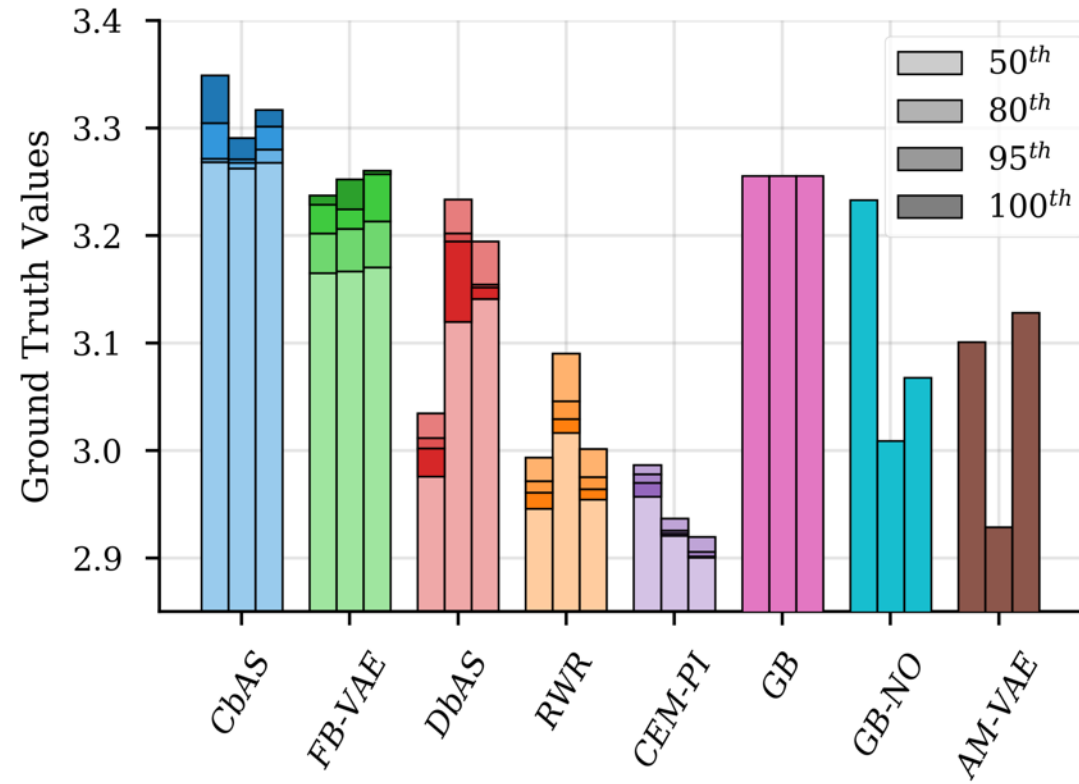


Ground truth GP

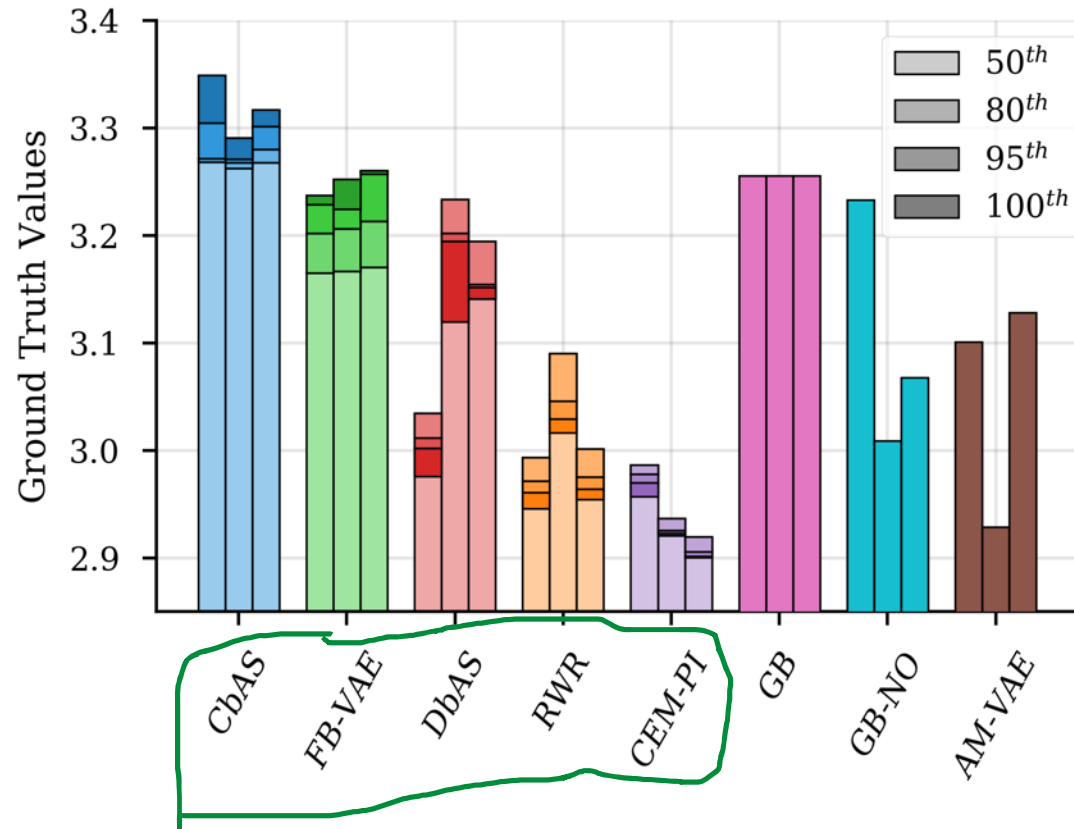Training data

Oracles

# Testing strategy

- Simulate a ground truth based on real data

    → "Ground truth" is a GP mean function

- Ground truth vales values are sampled from the GP for given sequences

- Use these input-output pairs to train oracles

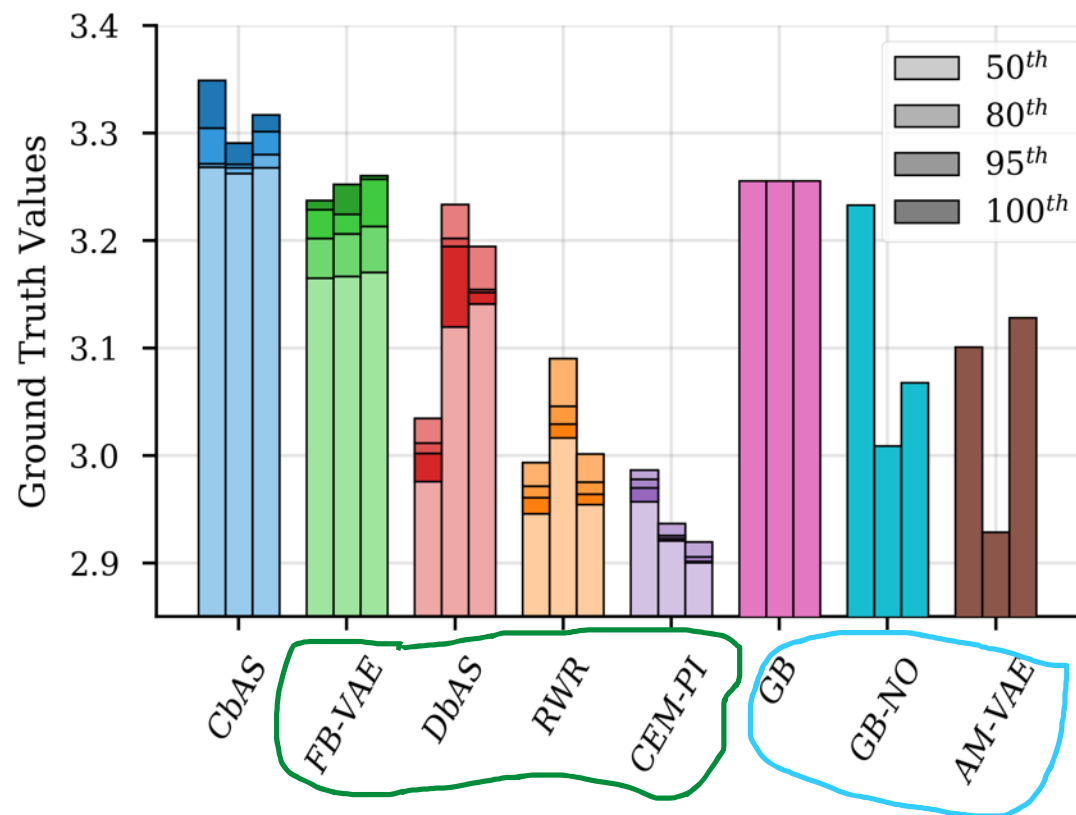- Coerce training set so these oracles exhibit pathologies

# Results

# Results

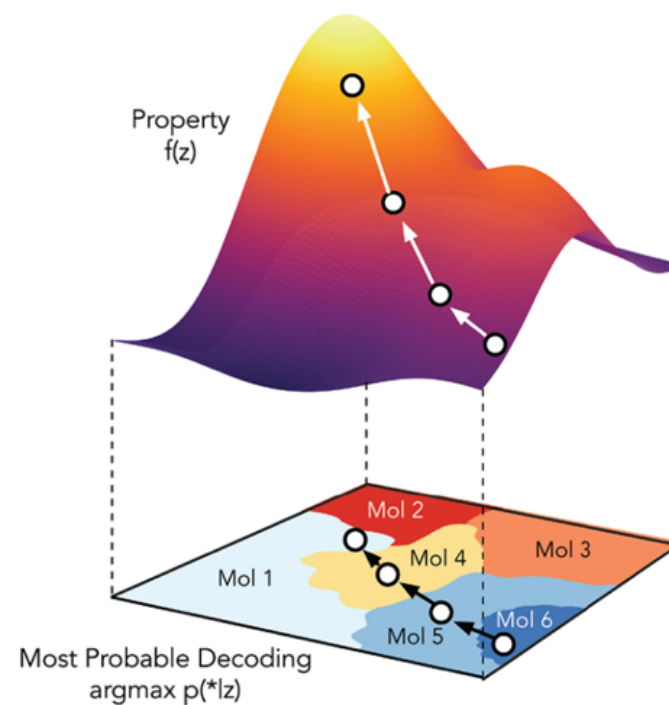Use weighted ML updates with weights:

- *CbAS*: $\dfrac{p(x|\theta^{(0)})}{p(x|\theta^{(t)})} P(S^{(t)}|x)$

- *DbAS*: $P(S^{(t)}|x)$

- *RWR*: $e^{\alpha\, f(x)}$

- CEM-PI: $\mathbb{1}_{\{PI(x)>\gamma^{(t)}\}}(x)$

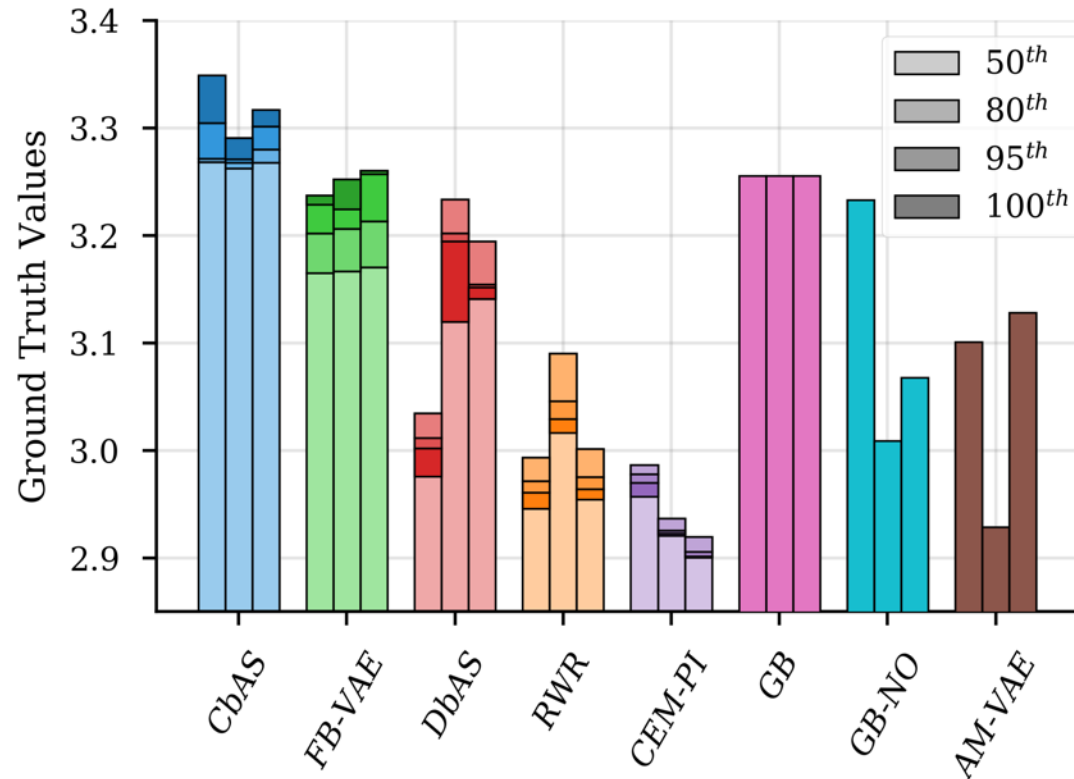- FB-VAE: $\mathbb{1}_{\{f(x)>\gamma^{(t)}\}}(x)$ w/ additional considerations

# Results



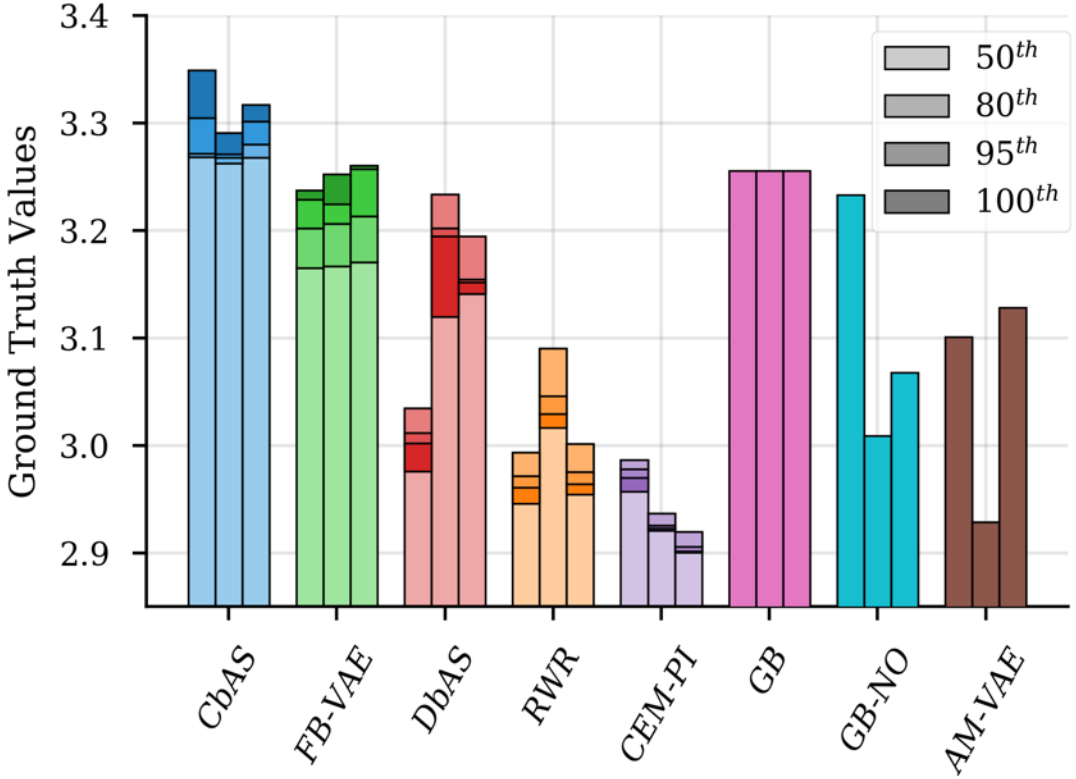Model-based optimizations

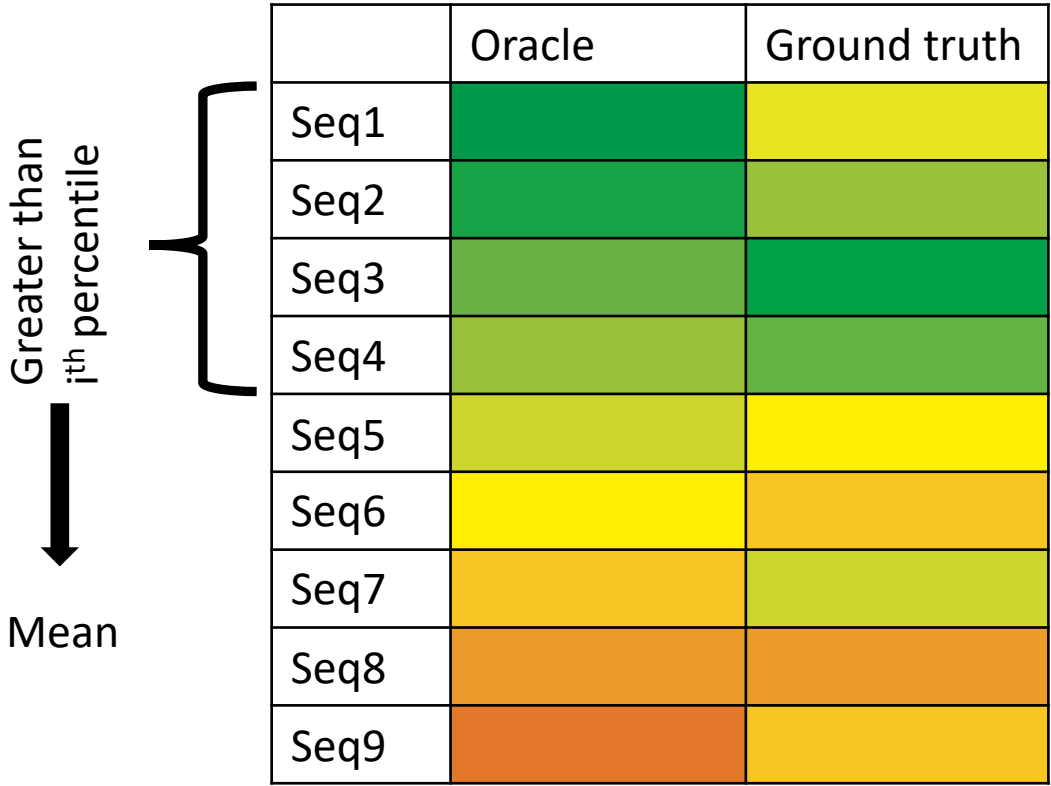Gradient descent on latent spaces

# Results



What does each bar represent?

| | Oracle | Ground truth |
|---|---|---|
| Seq1 | | |
| Seq2 | | |
| Seq3 | | |
| Seq4 | | |
| Seq5 | | |
| Seq6 | | |
| Seq7 | | |
| Seq8 | | |
| Seq9 | | |

# Results



What does each bar represent?

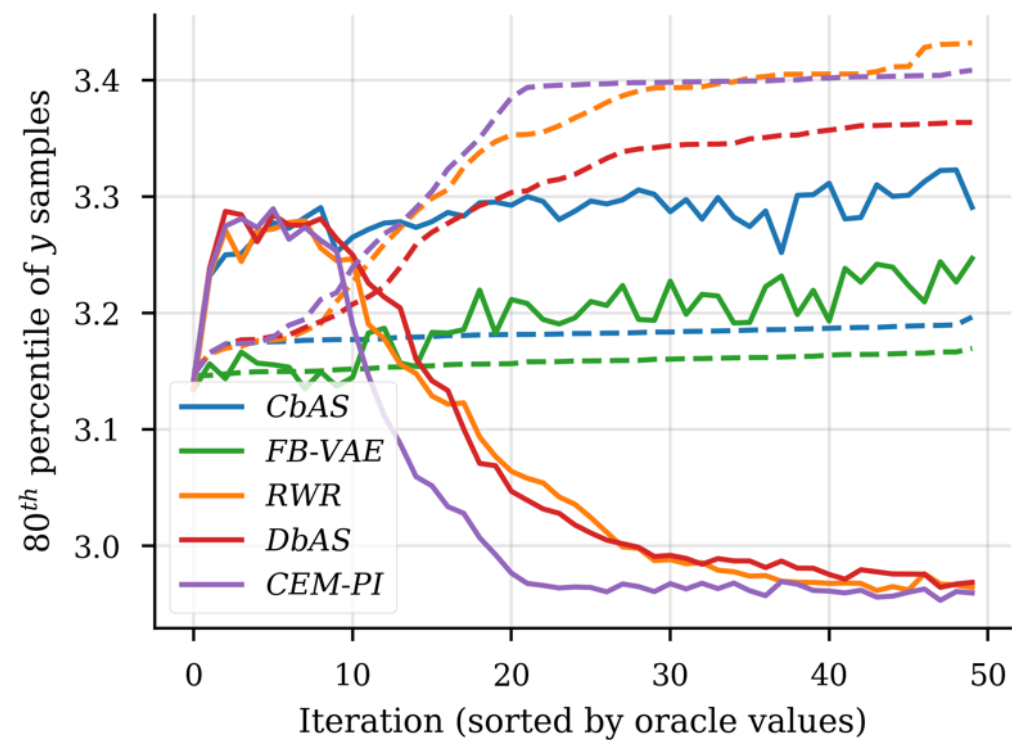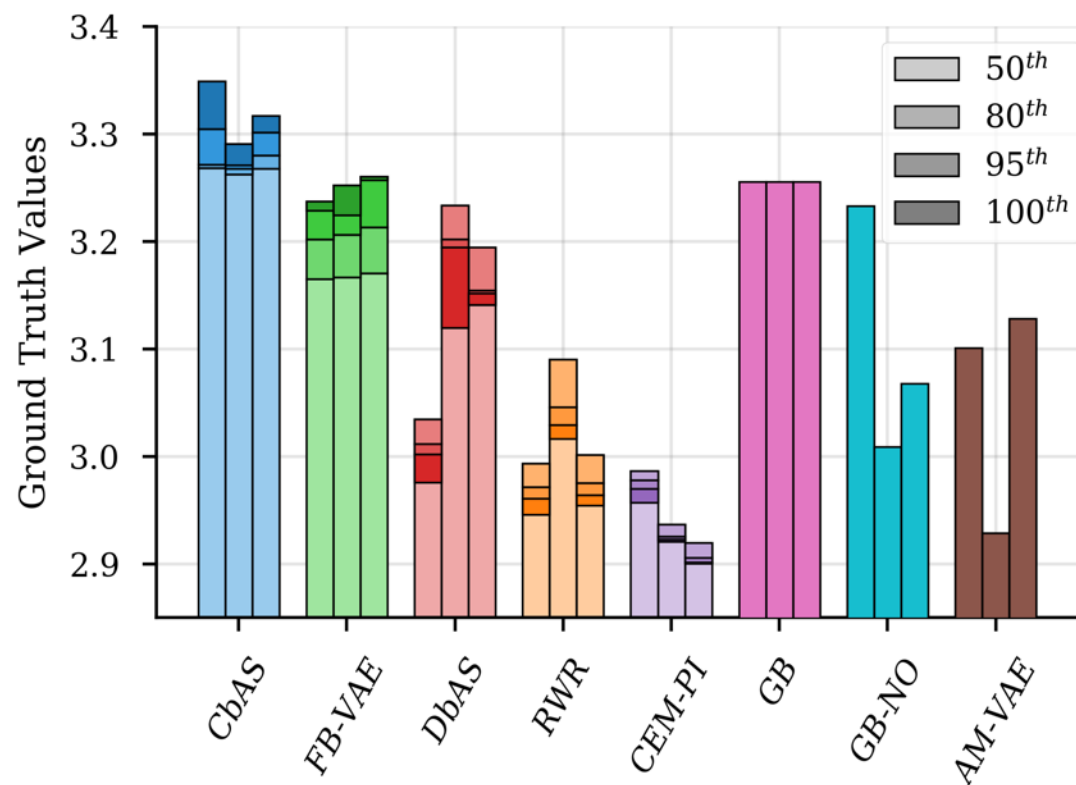| | Oracle | Ground truth |
|------|--------|--------------|
| Seq1 | | |
| Seq2 | | |
| Seq3 | | |
| Seq4 | | |
| Seq5 | | |
| Seq6 | | |
| Seq7 | | |
| Seq8 | | |
| Seq9 | | |

Greater than i$^{th}$ percentile

Mean

# Results

# Wrap-up

- Introduced a new model-based optimization method that is robust to pathological oracles

- Specifically targeted for discrete design problems

- Ongoing work to move beyond proof-of-principle:
  - Collaboration with wet-lab to perform end-to-end validation

# Thanks!