

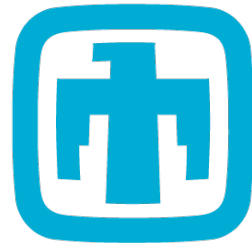
# Bayesian Optimal Experimental Designs for High-dimensional Physics-based Models

*Jim Oreluk, Leonid Sheps, Habib N. Najm*

**joreluk@sandia.gov**

Sandia National Laboratories, Livermore, CA, USA

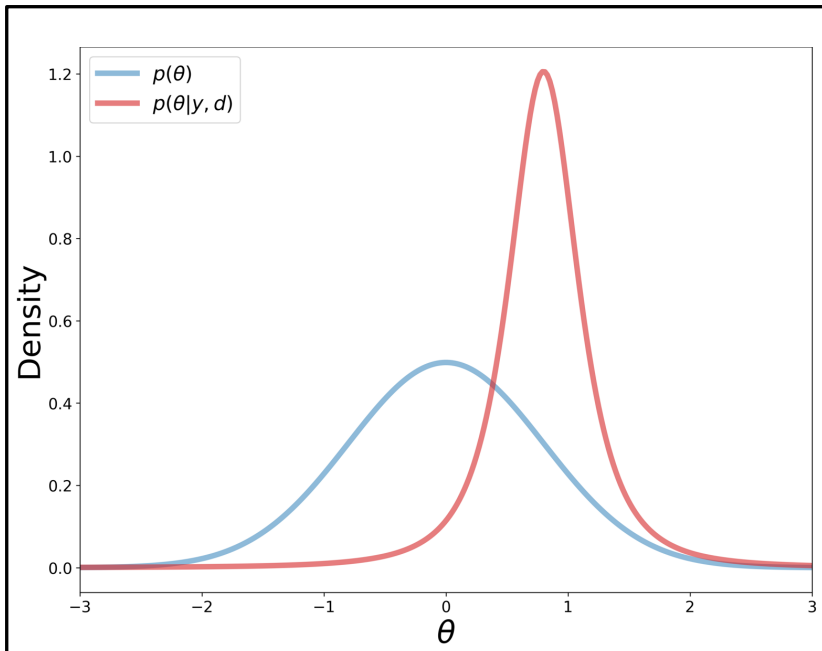
DATAWorks 2023



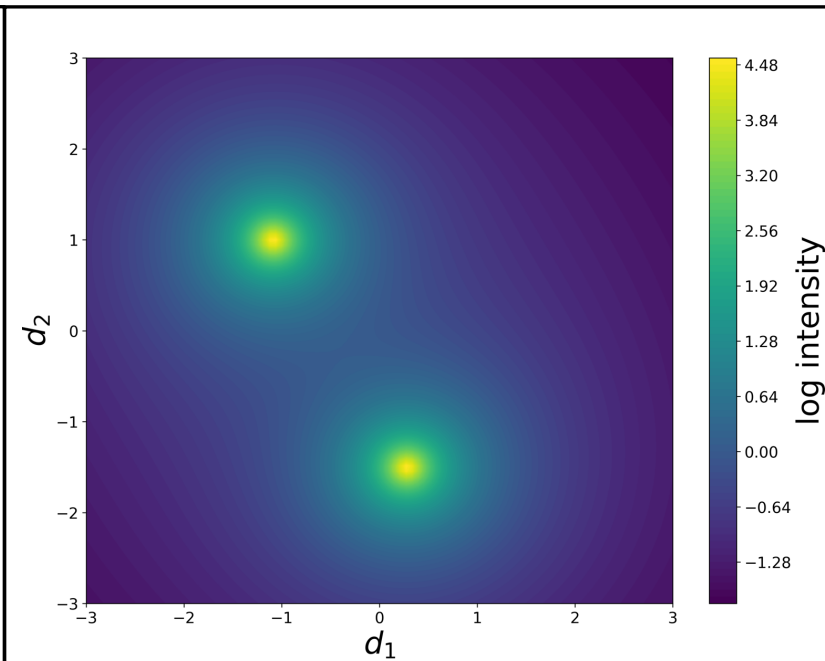
**Sandia  
National  
Laboratories**

# Motivation

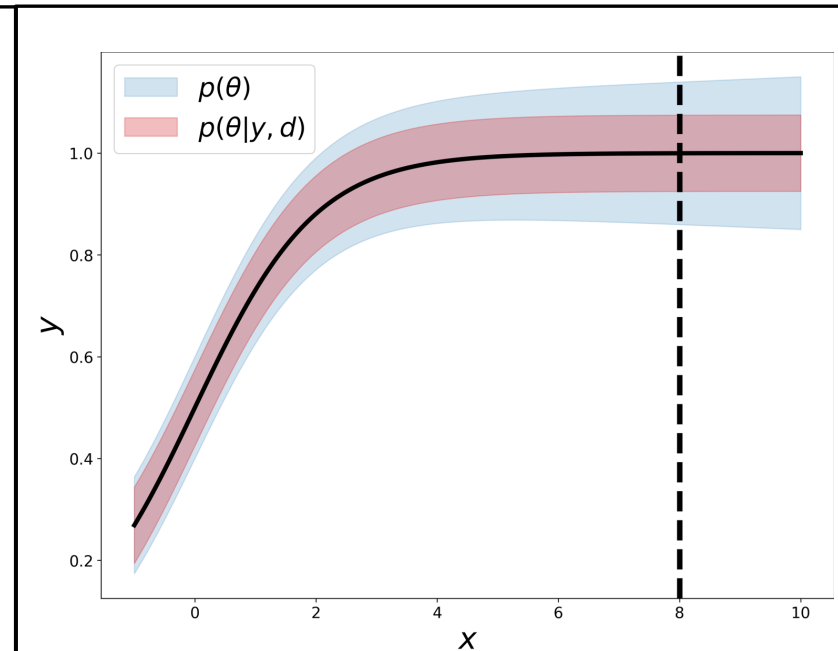
Experimental data helps us improve our understanding of a problem



Learning unknown parameter values



Identifying source locations in a field

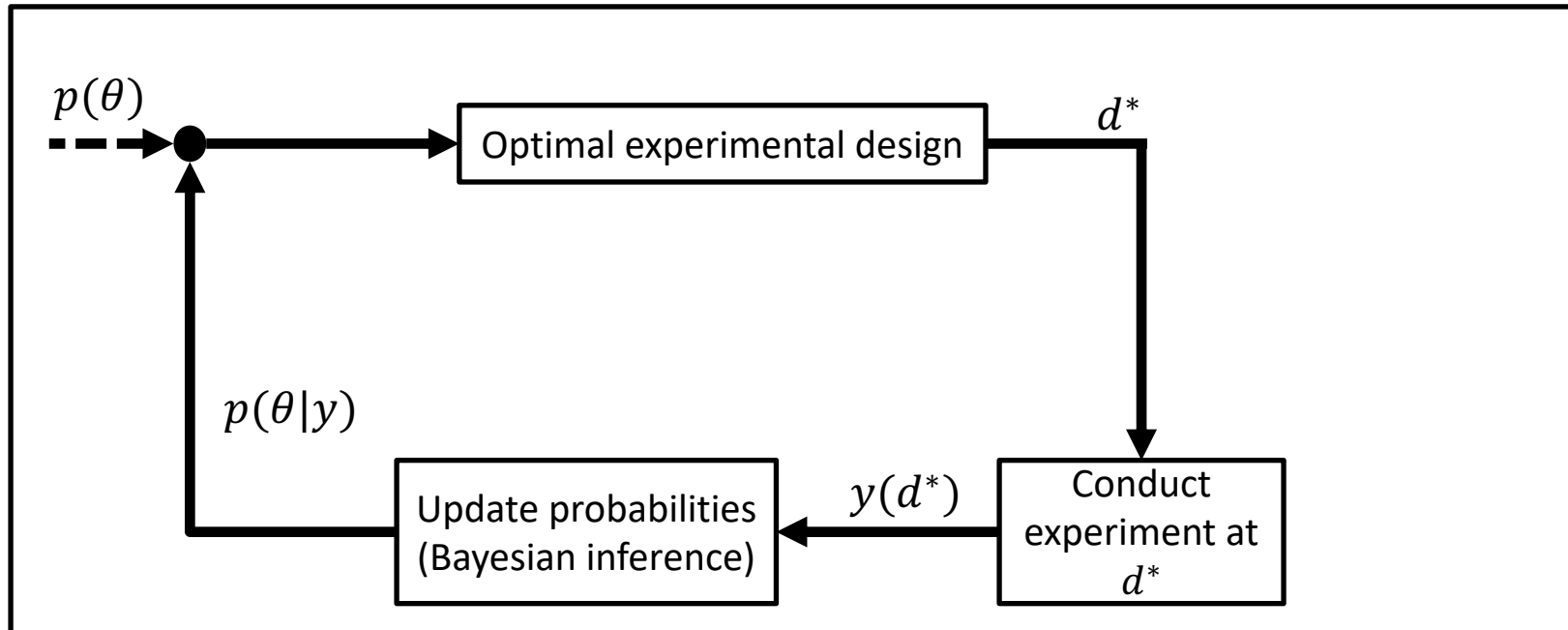


Improving predictivity of a QOI

## Decision making problem

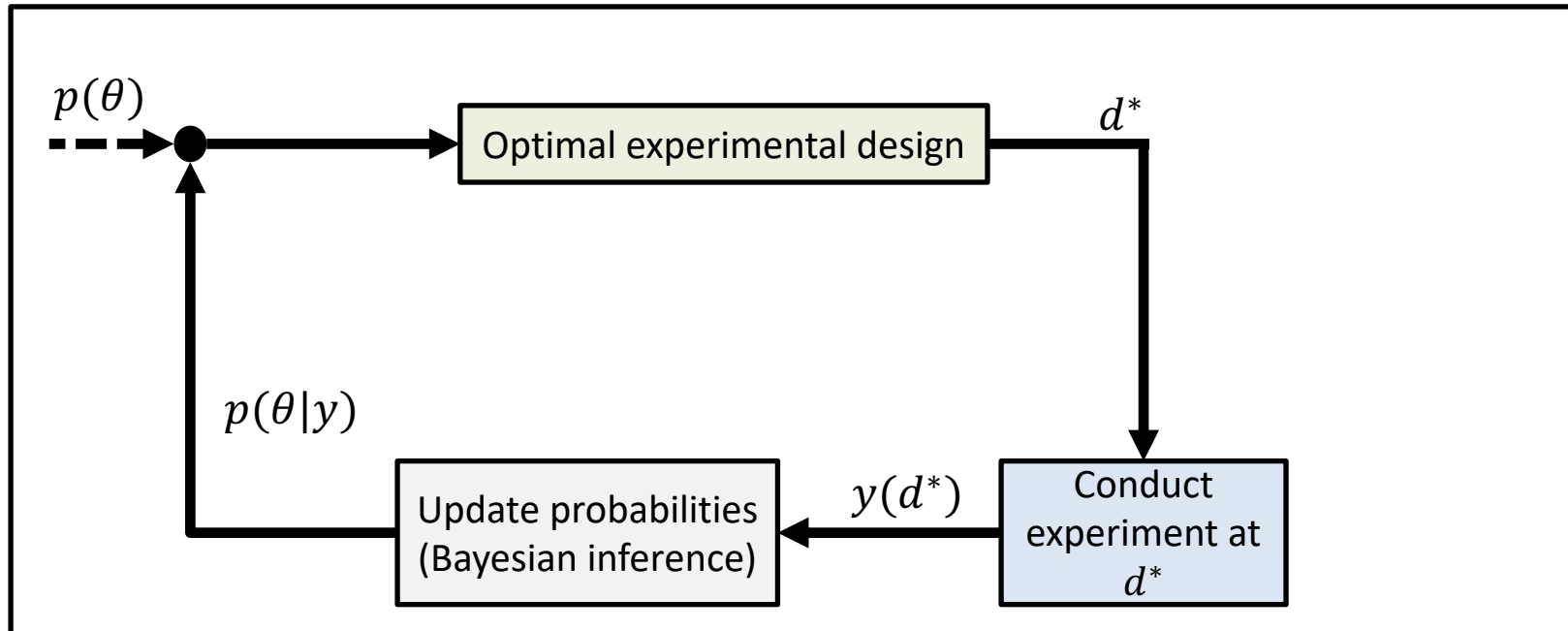
What experiment and/or test should we conduct that is **most** beneficial?

# Optimal experimental design loop



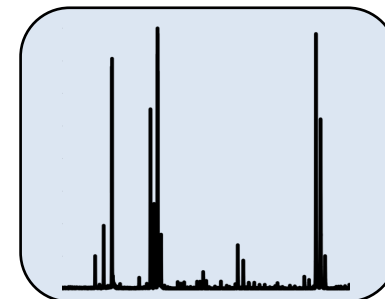
# Optimal experimental design loop

$$d^* = \arg \max_{d \in \mathcal{D}} U(d)$$



Bayes' theorem

$$p(\theta|y) = \frac{p(y|\theta)p(\theta)}{p(y)}$$



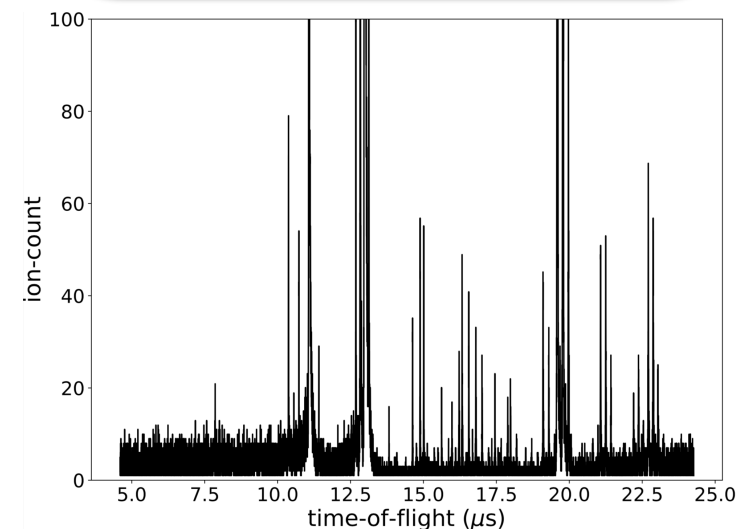
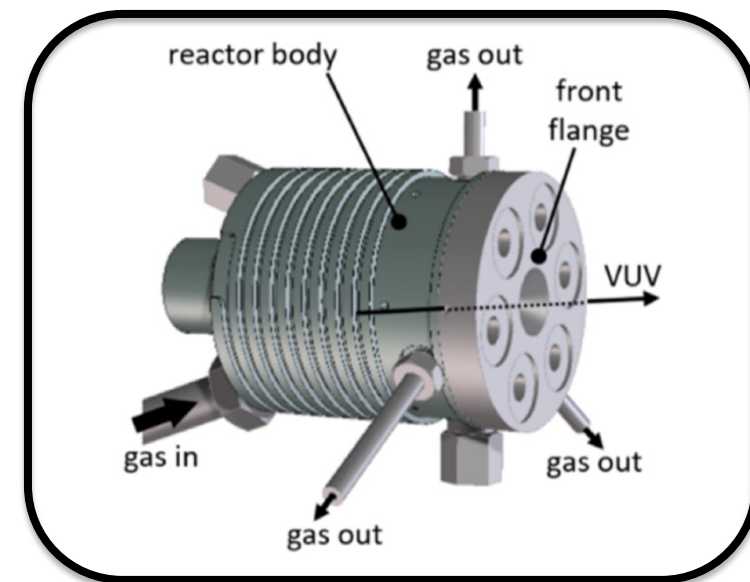
# Experimental overview

**Objective:** study the low-temperature oxidation of propane

- Measure the evolution of highly reactive intermediates and products

## High-pressure photolysis reactor (HPR) experiment

- Premixed mixture flows into a constant pressure reactor
- Photolysis laser fires instantaneously irradiating the gas mixture
  - Chemical precursor breaks down initiating reactions
- Gas mixture exhausts out, sampled by a synchrotron tunable vacuum-ultraviolet (VUV) photoionization mass spectrometer
  - Measurement of time-of-flight mass spectrum



Time-of-flight mass spectrum at a fixed VUV energy (11.3 eV) and kinetic time (60 ms)

L. Sheps, I. Antonov, K. Au. Sensitive mass spectrometer for time-resolved gas-phase chemistry studies at high pressures. *The Journal of Physical Chemistry A* 123.50 (2019) 10804-10814.

# Motivation

- Operation of the real experiment is **costly and laborious**
  - Potential setup time for the apparatus
  - Calibration of the apparatus and instrumentation
- **Limited time** to run experiments
  - Advanced Light Source, Lawrence Berkeley National Laboratory

Can we use a computational model to identify which experimental conditions are expected to be most informative?

# Modeling the HPR experiment

## Data model:

$$y(\mathbf{d}, \mathbf{x}) = \xi(\mathbf{d}, \mathbf{x}) + \epsilon(\mathbf{d}, \mathbf{x})$$

$$y(\mathbf{d}, \mathbf{x}) = f(\boldsymbol{\theta}, \mathbf{d}, \mathbf{x}) + \delta(\mathbf{x}) + \epsilon(\mathbf{d}, \mathbf{x})$$

$$\mathbf{x} = [\tau, t, E], \boldsymbol{\theta} \in \mathbb{R}^{1151}$$

$$\delta(\mathbf{x}) \sim GP(\mu_\delta(\mathbf{x}), \Sigma_\delta(\mathbf{x}, \mathbf{x}')), \epsilon(\mathbf{x}) \sim \mathcal{N}(0, s(\mathbf{x})^2)$$

$\mathbf{d}$ : design conditions

$\boldsymbol{\theta}$ : model parameters

$\mathbf{x}$ : spatial/temporal coordinates

$y(\mathbf{d}, \mathbf{x})$ : ion-count data

$\xi(\mathbf{d}, \mathbf{x})$ : true physical process

$f(\boldsymbol{\theta}, \mathbf{d}, \mathbf{x})$ : physics and instrument model

$\delta(\mathbf{x})$ : model error

$\epsilon(\mathbf{d}, \mathbf{x})$ : observation noise

### Physics model

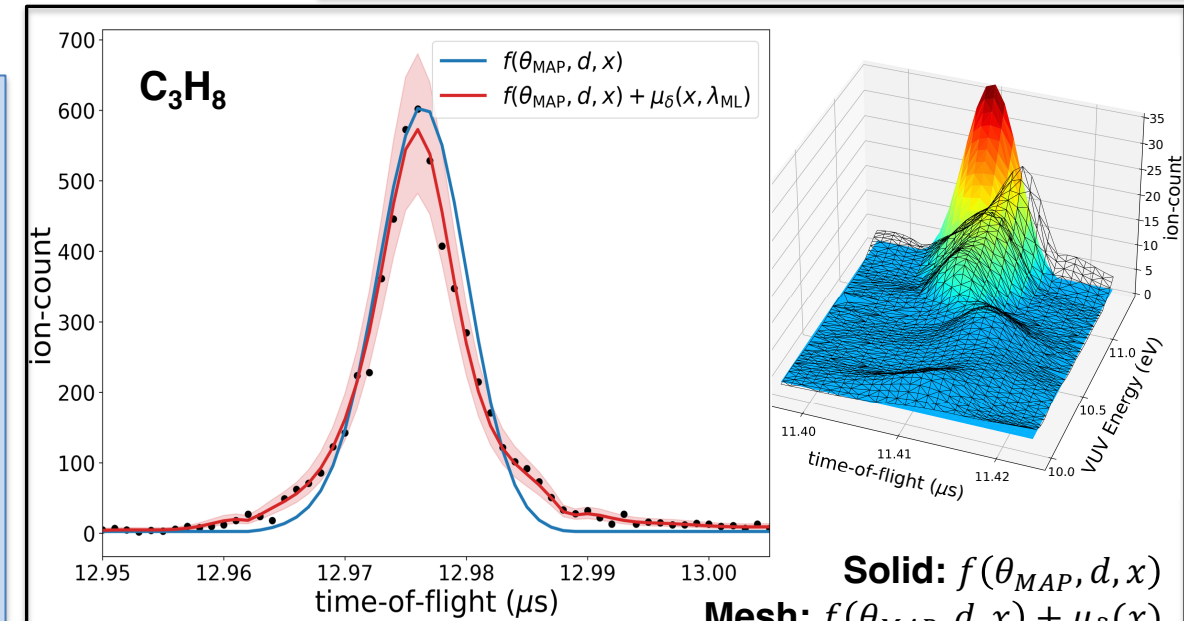
- Batch reactor model
- Photolysis laser and ionization model [1]

### Chemical model

- C0 - C3 chemical mechanism
- 173 species / 1146 reactions [2]

### Instrument model

- Maps concentrations to ion counts
- Peaks idealized as Gaussian distributions



Predictions with model error provides additional fidelity to the physics model while being more inline with the data

[1] Oreluk, et al. *Combustion Theory and Modelling*, 2022.

[2] Miller, et al. *Progress in Energy and Combustion Science*, 2021.

# Bayesian optimal experimental design

## Objective

Find a set of experimental conditions that maximizes the expected utility  $U(d)$

- Utility function models and compares the desirability of outcomes
- Target experiments to learn specific chemical rate constant measurements

$$d^* = \arg \max_{d \in \mathcal{D}} U(d)$$

where,

$$U(d) = \int_{y \in \mathcal{Y}} \int_{\theta \in \Theta} u(y, d, \theta) p(\theta, y | d) d\theta dy$$
$$= \int_{y \in \mathcal{Y}} \int_{\theta \in \Theta} u(y, d, \theta) p(\theta | y, d) p(y | d) d\theta dy$$

### Notation

$d$  : design conditions  
 $\theta$  : model parameters  
 $y$  : data

Expected utility involves a double integral over data and model parameters.

Necessitates a model for the physical system and for the experimental instrument.

Uses **Bayesian posterior on model parameters**, and **Bayesian model evidence on data**.



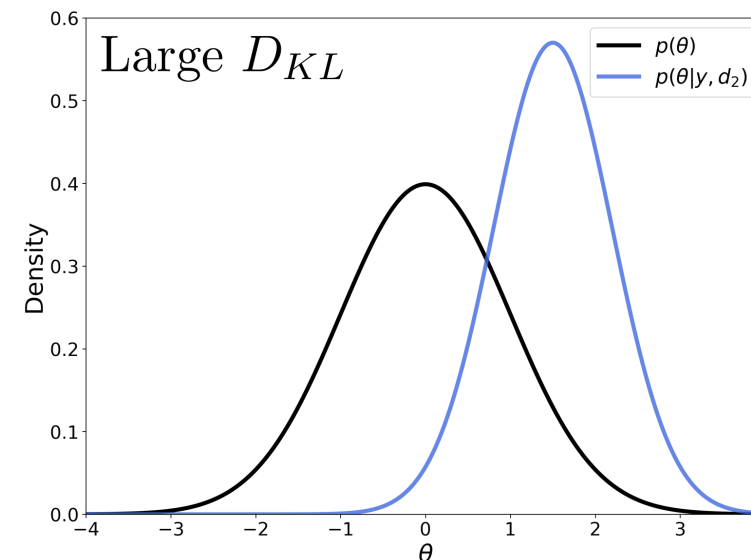
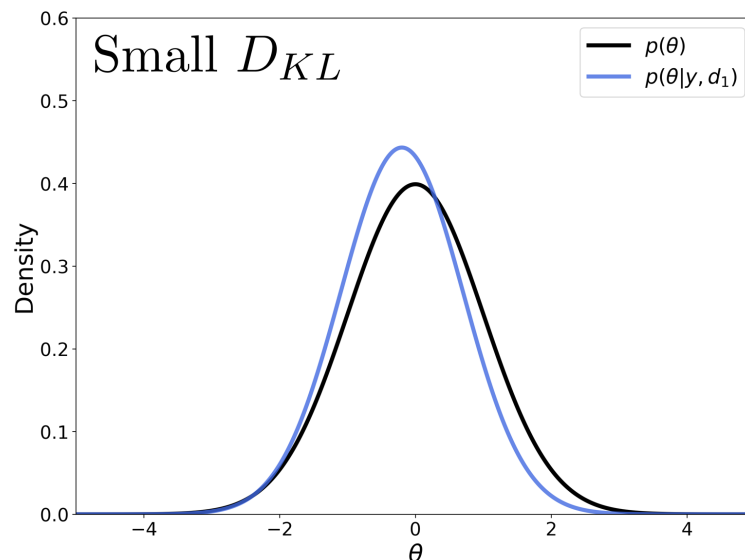
# Selecting a utility function, $u(y, d)$

Utility function should *reflect the goal or purpose* of the experiment

## 1) Improve the estimate of a unknown model parameter

- Information gain of an experiment is closely related to minimizing the parameter uncertainty

$$u(y, d) = D_{KL}(p(\theta|y, d) || p(\theta)) = \int p(\theta|y, d) \log \left[ \frac{p(\theta|y, d)}{p(\theta)} \right] d\theta$$



# Selecting a utility function, $u(y, d)$

Utility function should *reflect the goal or purpose* of the experiment

## 1) Improve the estimate of a unknown model parameter

- Information gain of an experiment is closely related to minimizing the parameter uncertainty

$$u(y, d) = D_{\text{KL}}(p(\theta|y, d) || p(\theta)) = \int p(\theta|y, d) \log \left[ \frac{p(\theta|y, d)}{p(\theta)} \right] d\theta$$

## 2) Improve the predictivity of a quantity of interest, $q$

- Measure divergence between **prior predictive** and **posterior predictive** distributions

$$\begin{aligned} p(q) &= \int_{\theta \in \Theta} p(q|\theta)p(\theta)d\theta & u(y, d) &= D_{\text{KL}}(p(q|y, d) || p(q)) \\ p(q|y, d) &= \int_{\theta \in \Theta} p(q|\theta)p(\theta|y, d)d\theta & &= \int p(q|y, d) \log \left[ \frac{p(q|y, d)}{p(q)} \right] d\theta \end{aligned}$$

# Selecting a utility function, $u(y, d)$

Utility function should *reflect the goal or purpose* of the experiment

## 1) Improve the estimate of a unknown model parameter

- Information gain of an experiment is closely related to minimizing the parameter uncertainty

$$u(y, d) = D_{\text{KL}}(p(\theta|y, d) || p(\theta)) = \int p(\theta|y, d) \log \left[ \frac{p(\theta|y, d)}{p(\theta)} \right] d\theta$$

$$U(d) = \int_{\mathcal{Y}} \left( \int_{\Theta} p(\theta|y, d) \log \left[ \frac{p(\theta|y, d)}{p(\theta)} \right] d\theta \right) p(y|d) dy$$

$$U(d) = \int_{\mathcal{Y}} \int_{\Theta} [\log p(y|\theta, d) - \log p(y|d)] p(y|\theta, d) p(\theta) d\theta dy$$

# Approximating the expected utility

## Numerical approximation:

$$U(d) \approx \frac{1}{N} \sum_{i=0}^N \left[ \log p(y^{(i)} | \theta^{(i)}, d) - \underbrace{\log p(y^{(i)} | d)}_{??} \right] \quad \text{where, } \theta^{(i)} \sim p(\theta) \\ y^{(i)} \sim p(y | \theta^{(i)}, d)$$

## Nested Monte Carlo

$$U(d) \approx \frac{1}{N} \sum_{i=0}^N \left[ \log p(y^{(i)} | \theta^{(i)}, d) - \log \left( \frac{1}{M} \sum_{j=0}^M p(y^{(i)} | \theta^{(j)}, d) \right) \right]$$

$$d^* = \arg \max_{d \in \mathcal{D}} U(d)$$

- *Costly to evaluate  $U(d)$*

-  $\frac{dU}{dd} = ?$

T. Rainforth et al., On nesting Monte Carlo estimators, *International Conference on Machine Learning*. PMLR, 2018.

K.J. Ryan, Estimating expected information gains for experimental designs with application to the random fatigue-limit model, *Journal of Computational and Graphical Statistics* 12 (2003) 585–603.

# Maximizing the expected utility, $U(d)$

$$d^* = \arg \max_{d \in \mathcal{D}} U(d)$$

## Bayesian Optimization

- Construct a GP of the unknown objective function

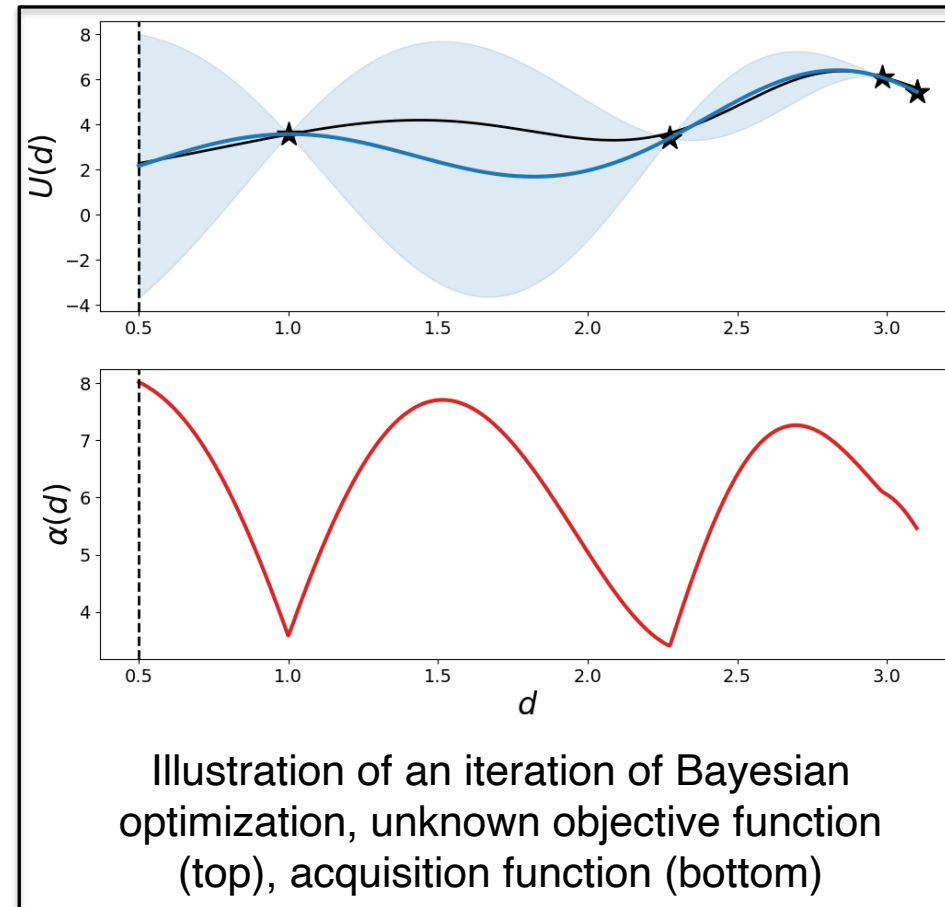
$$U(d) \sim \mathcal{N}(\mu(d), K(d, d'))$$

- Use *acquisition function*  $\alpha(d)$  to select new samples
  - Gaussian Process upper confidence bound (UCB)

$$\alpha_t(d) = \mu_{t-1}(d) + \sqrt{\beta_t} \sigma_{t-1}(d)$$

$$\sigma_{t-1}(d) = \sqrt{K(d, d)}$$

- Exploits regions with a high mean and explores regions of high uncertainty
- Select next optimization iteration at:  $d_t = \arg \max_{d \in \mathcal{D}} \alpha_t(d)$
- Evaluate utility function at  $U(d_t)$



# Challenges

## High-dimensional model output

- High-fidelity physics-based simulations are be expensive to evaluate
- Evaluating  $U(d) \sim \mathcal{O}(NM)$  forward solves

$$U(d) \approx \frac{1}{N} \sum_{i=0}^N \left[ \log p(y^{(i)} | \theta^{(i)}, d) - \log \left( \frac{1}{M} \sum_{j=0}^M p(y^{(i)} | \theta^{(j)}, d) \right) \right]$$

- Constructing a surrogate model addresses the costly run-time

$$f(\boldsymbol{\theta}, \mathbf{d}, x_i) \approx g(\boldsymbol{\theta}, \mathbf{d}, x_i) = \sum_{p=1}^P c_p \Psi_p(\xi)$$

- Total **number of outputs** remains problematic  $y(\mathbf{d}, \mathbf{x}) \in \mathbb{R}^{510,000,000}$

# Challenges

## High-dimensional model output

- High-fidelity physics-based simulations are be expensive to evaluate
- Evaluating  $U(d) \sim \mathcal{O}(NM)$  forward solves

$$U(d) \approx \frac{1}{N} \sum_{i=0}^N \left[ \log p(y^{(i)} | \theta^{(i)}, d) - \log \left( \frac{1}{M} \sum_{j=0}^M p(y^{(i)} | \theta^{(j)}, d) \right) \right]$$

- Constructing a surrogate model addresses the costly run-time

$$f(\boldsymbol{\theta}, \mathbf{d}, x_i) \approx g(\boldsymbol{\theta}, \mathbf{d}, x_i) = \sum_{p=1}^P c_p \Psi_p(\xi)$$

Can we find *low-dimensional* representations of the high-dimensional output?

# Reducing output dimensionality

**Goal:** Map high-dimensional model output to a lower-dimensional space while minimizing loss of information

## Truncated SVD

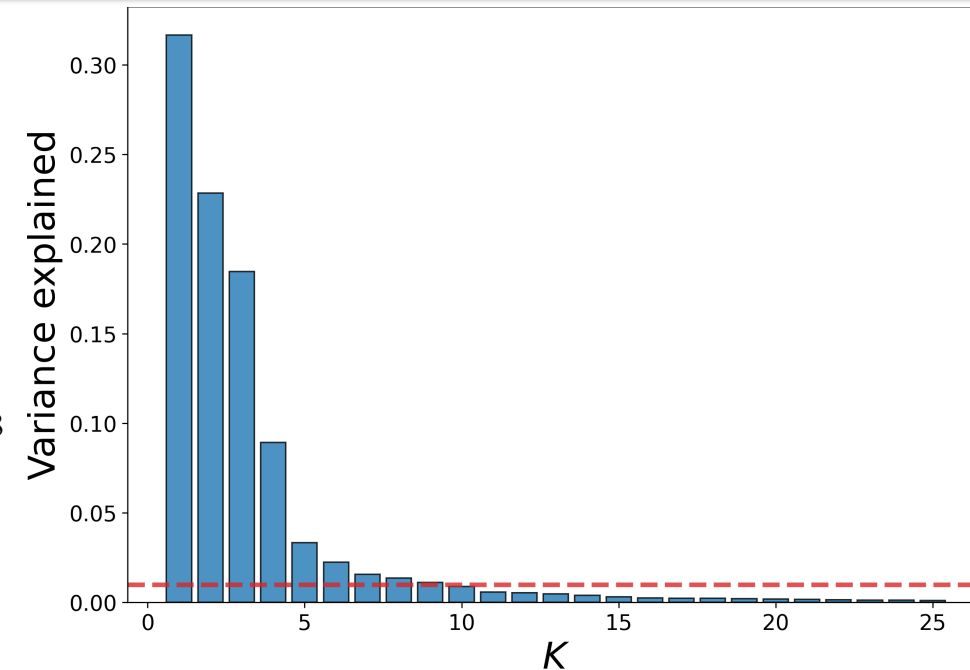
At a fixed design  $d$ ,

- Draw  $n$  samples, where  $\theta^{(i)} \sim p(\theta)$
- Evaluate the model  $f(\theta, d, x) + \delta(\lambda_{\text{ML}}, x) + \epsilon(x)$
- Construct output matrix  $Y = USV^T$ , where  $Y \in \mathbb{R}^{n \times J}$ ,  $J = 5.1 \times 10^8$
- Retain only top  $K$  singular values of  $S$
- Low-rank approximation:  $Y_K = U_K S_K V_K^T$

## Transformation:

$$q(\theta, d, x) = y(\theta, d, x) V_K$$

$$\underbrace{q(\theta, d, x)}_{(1 \times K)} = \underbrace{[f(\theta, d, x) + \delta(\lambda_{\text{ML}}, x) + \epsilon(x)]}_{(1 \times J)} \underbrace{V_K}_{(J \times K)}$$



93% of the total variance is explained by a small number ( $K = 10$ ) components at 620 K, 2 bar,  $n_{\text{C}_3\text{H}_8} = 2\text{e}14$ ,  $n_{\text{O}_2} = 8\text{e}18$ ,  $n_{\text{Cl}} = 1\text{e}13$ .



# Challenges: Target *specific* model parameters

Expected utility measures difference between  $p(\theta)$  and  $p(\theta|y)$ , in other words, difference in *all uncertain parameters*

- Only interested in designs that improve *particular chemical rates*
  - chemical model contains 1146 uncertain chemical rates

## Set of targeted reaction rates

Reaction ID	
	<b>n-R + O2 → QOOH3</b>
1040	O2 + CH3CH2CH2 → CH3CHCH2OOH
	<b>n-R + O2 → QOOH1</b>
1041	O2 + CH3CH2CH2 → CH2CH2CH2OOH
1042	O2 + CH3CH2CH2 → CH2CH2CH2OOH
	<b>n-R + O2 → propene</b>
1043	O2 + CH3CH2CH2 → HO2 + CH3CHCH2
	<b>n-R + O2 → propylene oxide</b>
1044	O2 + CH3CH2CH2 → OH + C-CH2OCH(CH3)
	<b>n-RO2 → products</b>
1045	CH3CH2CH2OO → CH3CHCH2OOH
1046	CH3CH2CH2OO → CH2CH2CH2OOH
1047	CH3CH2CH2OO → CH2CH2CH2OOH
1048	CH3CH2CH2OO → HO2 + CH3CHCH2
1049	CH3CH2CH2OO → OH + C-CH2OCH(CH3)
	<b>QOOH3 → products</b>
1051	CH3CHCH2OOH → HO2 + CH3CHCH2
1052	CH3CHCH2OOH → OH + C-CH2OCH(CH3)
	<b>QOOH1 → propene</b>
1053	CH2CH2CH2OOH → HO2 + CH3CHCH2
	<b>i-R+O2 → products</b>
1056	O2 + CH3CHCH3 → CH3CH(OOH)CH2
1057	O2 + CH3CHCH3 → HO2 + CH3CHCH2
1058	O2 + CH3CHCH3 → OH + C-CH2OCH(CH3)
	<b>i-RO2 → products</b>
1059	CH3CH(OO)CH3 → CH3CH(OOH)CH2
1060	CH3CH(OO)CH3 → HO2 + CH3CHCH2
1061	CH3CH(OO)CH3 → OH + C-CH2OCH(CH3)
	<b>QOOH2 → products</b>
1062	CH3CH(OOH)CH2 → HO2 + CH3CHCH2
1063	CH3CH(OOH)CH2 → OH + C-CH2OCH(CH3)
	<b>O2 + QOOH1 → products</b>
1065	O2 + CH2CH2CH2OOH → OHOCH2CH2CH2OO
	<b>O2 + QOOH3 → products</b>
1069	O2 + CH3CHCH2OOH → CH3CH(OO)CH2OOH
1071	O2 + CH3CHCH2OOH → CH2CH(OOH)CH2OOH
	<b>O2 + QOOH2 → products</b>
1080	O2 + CH3CH(OOH)CH2 → CH3CH(OOH)CH2OO
1081	O2 + CH3CH(OOH)CH2 → CH2CH(OOH)CH2OOH
	<b>3-hydroperoxy-propylperoxyradical → KHP</b>
1088	OHOCH2CH2CH2OO → OH + OCHCH2CH2OOH
	<b>KHP → products</b>
1112	OCHCH2CH2OOH → OCHCH2CH2O + OH

# Challenges: Target *specific* model parameters

Expected utility measures difference between  $p(\theta)$  and  $p(\theta|y)$ , in other words, difference in ***all uncertain parameters***

- Only interested in designs that improve *particular chemical rates*
  - chemical model contains 1146 uncertain chemical rates

**Solution:** Reformulate utility to marginalize over nuisance parameters

Let  $\theta$  be the parameters of interest and  $\eta$  be the nuisance parameters

Nested Monte Carlo estimator becomes:

$$U(d) \approx \frac{1}{N} \sum_{i=1}^N \left[ \log \left( \frac{1}{K} \sum_{k=1}^K p \left( y^{(i)} | \theta^{(i)}, \eta^{(k)}, d \right) \right) - \log \left( \frac{1}{M} \sum_{j=1}^M p \left( y^{(i)} | \theta^{(j)}, \eta^{(j)}, d \right) \right) \right]$$

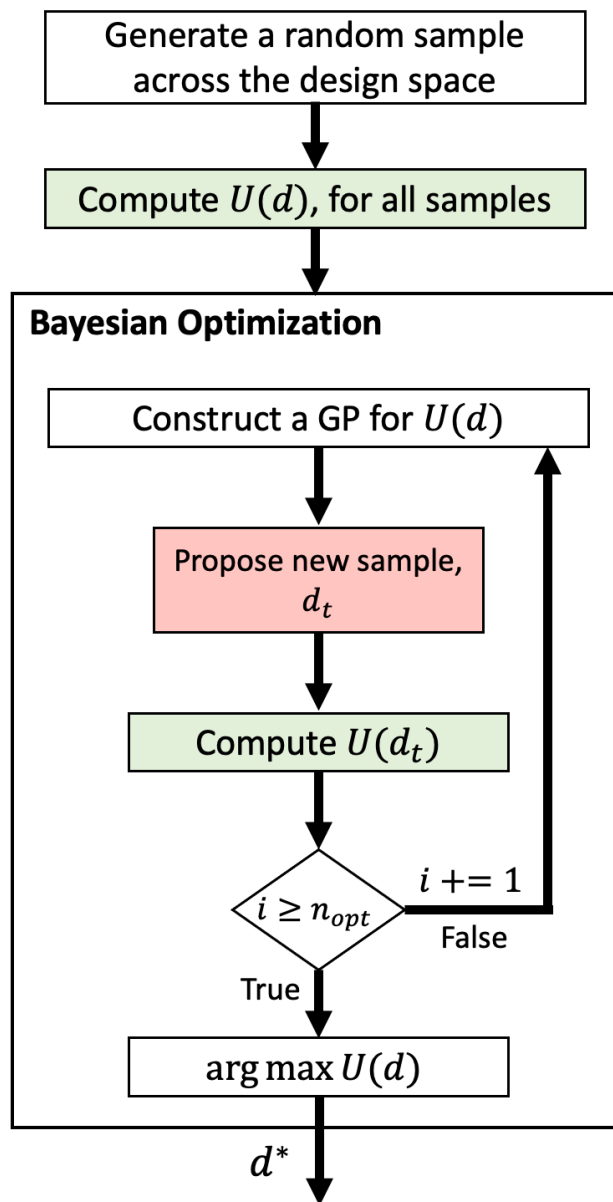
Additional loop in likelihood to marginalize over nuisance parameters

## Set of targeted reaction rates

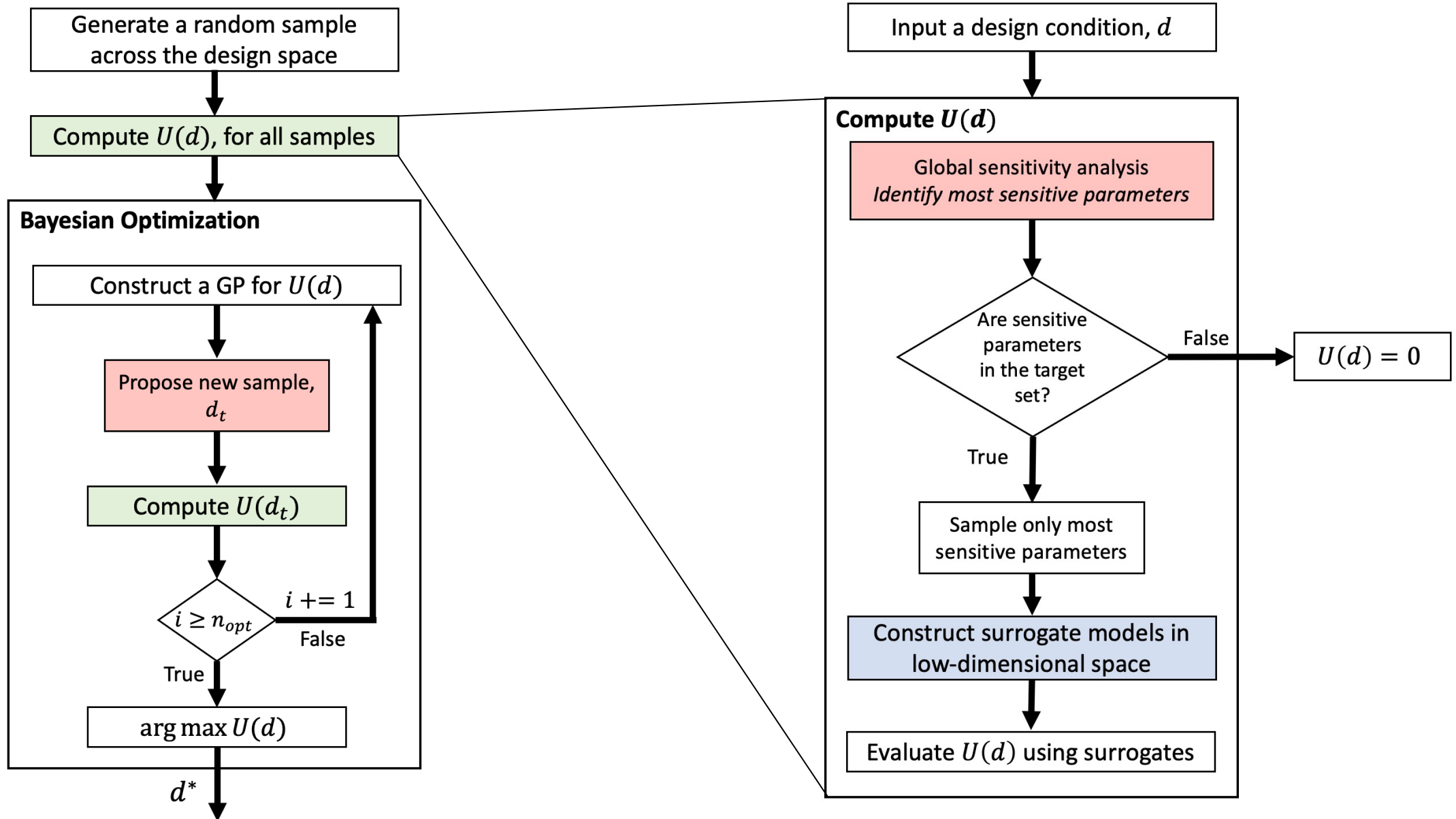
Reaction ID	
	<b>n-R + O2 → QOOH3</b>
1040	O2 + CH3CH2CH2 → CH3CHCH2OOH
	<b>n-R + O2 → QOOH1</b>
1041	O2 + CH3CH2CH2 → CH2CH2CH2OOH
1042	O2 + CH3CH2CH2 → CH2CH2CH2OOH
	<b>n-R + O2 → propene</b>
1043	O2 + CH3CH2CH2 → HO2 + CH3CHCH2
	<b>n-R + O2 → propylene oxide</b>
1044	O2 + CH3CH2CH2 → OH + C-CH2OCH(CH3)
	<b>n-RO2 → products</b>
1045	CH3CH2CH2OO → CH3CHCH2OOH
1046	CH3CH2CH2OO → CH2CH2CH2OOH
1047	CH3CH2CH2OO → CH2CH2CH2OOH
1048	CH3CH2CH2OO → HO2 + CH3CHCH2
1049	CH3CH2CH2OO → OH + C-CH2OCH(CH3)
	<b>QOOH3 → products</b>
1051	CH3CHCH2OOH → HO2 + CH3CHCH2
1052	CH3CHCH2OOH → OH + C-CH2OCH(CH3)
	<b>QOOH1 → propene</b>
1053	CH2CH2CH2OOH → HO2 + CH3CHCH2
	<b>i-R+O2 → products</b>
1056	O2 + CH3CHCH3 → CH3CH(OOH)CH2
1057	O2 + CH3CHCH3 → HO2 + CH3CHCH2
1058	O2 + CH3CHCH3 → OH + C-CH2OCH(CH3)
	<b>i-RO2 → products</b>
1059	CH3CH(OO)CH3 → CH3CH(OOH)CH2
1060	CH3CH(OO)CH3 → HO2 + CH3CHCH2
1061	CH3CH(OO)CH3 → OH + C-CH2OCH(CH3)
	<b>QOOH2 → products</b>
1062	CH3CH(OOH)CH2 → HO2 + CH3CHCH2
1063	CH3CH(OOH)CH2 → OH + C-CH2OCH(CH3)
	<b>O2 + QOOH1 → products</b>
1065	O2 + CH2CH2CH2OOH → OHCH2CH2CH2OO
	<b>O2 + QOOH3 → products</b>
1069	O2 + CH3CHCH2OOH → CH3CH(OO)CH2OOH
1071	O2 + CH3CHCH2OOH → CH2CH(OOH)CH2OOH
	<b>O2 + QOOH2 → products</b>
1080	O2 + CH3CH(OOH)CH2 → CH3CH(OOH)CH2OO
1081	O2 + CH3CH(OOH)CH2 → CH2CH(OOH)CH2OOH
	<b>3-hydroperoxy-propylperoxyradical → KHP</b>
1088	OHCH2CH2CH2OO → OH + OCHCH2CH2OOH
	<b>KHP → products</b>
1112	OCHCH2CH2OOH → OCHCH2CH2O + OH

Feng, Chi, and Youssef M. Marzouk. "A layered multiple importance sampling scheme for focused optimal Bayesian experimental design." *arXiv preprint arXiv:1903.11187* (2019).

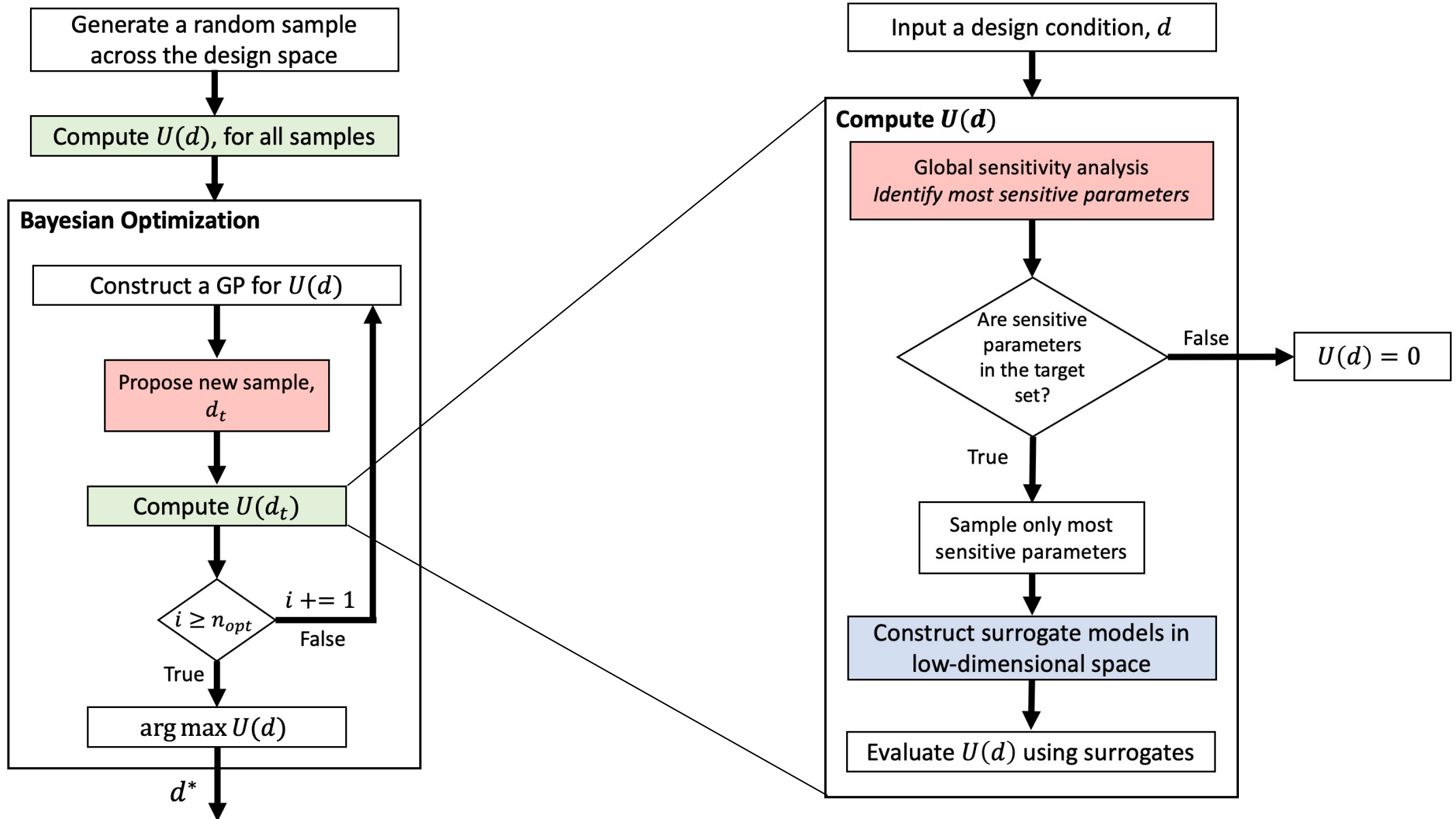
# Targeted OED workflow



# Targeted OED workflow



# Targeted OED workflow



# Results

Original data model,

$$y(\mathbf{d}, \mathbf{x}) = \xi(\mathbf{d}, \mathbf{x}) + \epsilon(\mathbf{d}, \mathbf{x})$$

$$y(\mathbf{d}, \mathbf{x}) = f(\boldsymbol{\theta}, \mathbf{d}, \mathbf{x}) + \delta(\mathbf{x}) + \epsilon(\mathbf{d}, \mathbf{x})$$

$$\mathbf{x} = [\tau, t, E], \boldsymbol{\theta} \in \mathbb{R}^{1151}$$

$$\delta(\mathbf{x}) \sim GP(\mu_\delta(\mathbf{x}), \Sigma_\delta(\mathbf{x}, \mathbf{x}')), \epsilon(\mathbf{x}) \sim \mathcal{N}(0, s(\mathbf{x})^2)$$

Design space:  $\mathbf{d} = [T, p, n_{Cl}, n_{O2}]$

$$\Sigma_\delta = \begin{bmatrix} \Sigma_1 & & & & \\ & \ddots & & & \\ & & \Sigma_D & & \\ & & & \ddots & \\ & & & & \Sigma_{Dt_{e+1}} & & \\ & & & & & \ddots & \\ & & & & & & \Sigma_{D(t_e+1)} \end{bmatrix}$$

## Assumptions:

- **Model error  $\delta(\mathbf{x})$  independent of  $\mathbf{d}$**
- Local GPs constructed independently for each kinetic time
- Considering all 1151 model parameters as uncertain
  - Designs are constructed to target specific sets of the model parameters
- Prior parameter uncertainty is specified by literature

# Results

Original data model,

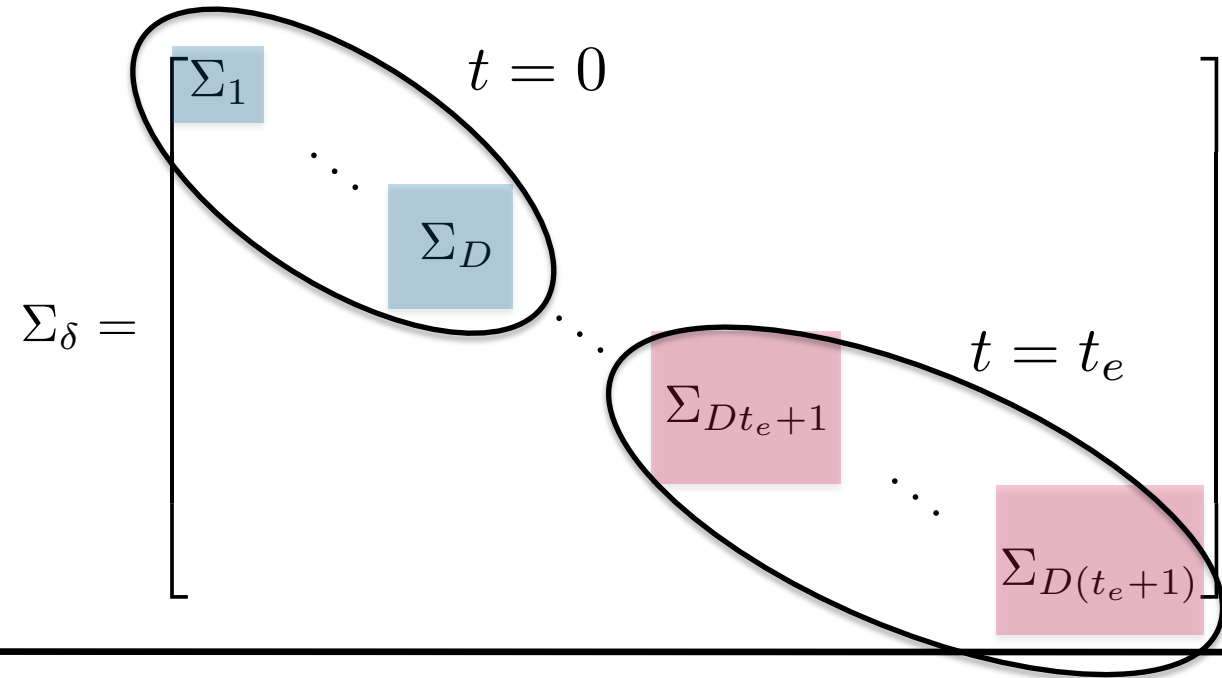
$$y(\mathbf{d}, \mathbf{x}) = \xi(\mathbf{d}, \mathbf{x}) + \epsilon(\mathbf{d}, \mathbf{x})$$

$$y(\mathbf{d}, \mathbf{x}) = f(\boldsymbol{\theta}, \mathbf{d}, \mathbf{x}) + \delta(\mathbf{x}) + \epsilon(\mathbf{d}, \mathbf{x})$$

$$\mathbf{x} = [\tau, t, E], \boldsymbol{\theta} \in \mathbb{R}^{1151}$$

$$\delta(\mathbf{x}) \sim GP(\mu_\delta(\mathbf{x}), \Sigma_\delta(\mathbf{x}, \mathbf{x}')), \epsilon(\mathbf{x}) \sim \mathcal{N}(0, s(\mathbf{x})^2)$$

Design space:  $\mathbf{d} = [T, p, n_{Cl}, n_{O2}]$

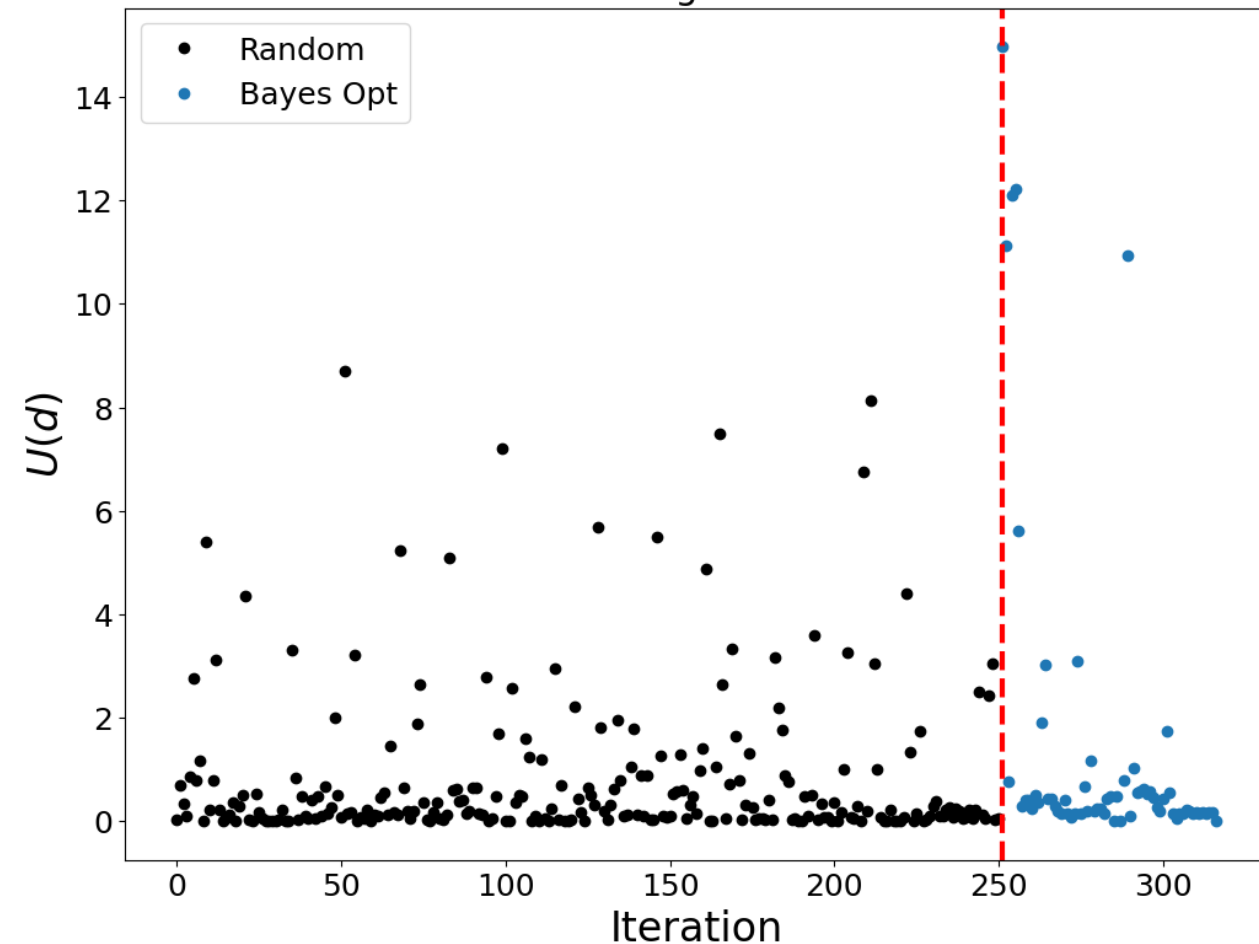


## Assumptions:

- Model error  $\delta(x)$  independent of  $d$
- **Local GPs constructed independently for each kinetic time**
- Considering all 1151 model parameters as uncertain
  - Designs are constructed to target specific sets of the model parameters
- Prior parameter uncertainty is specified by literature

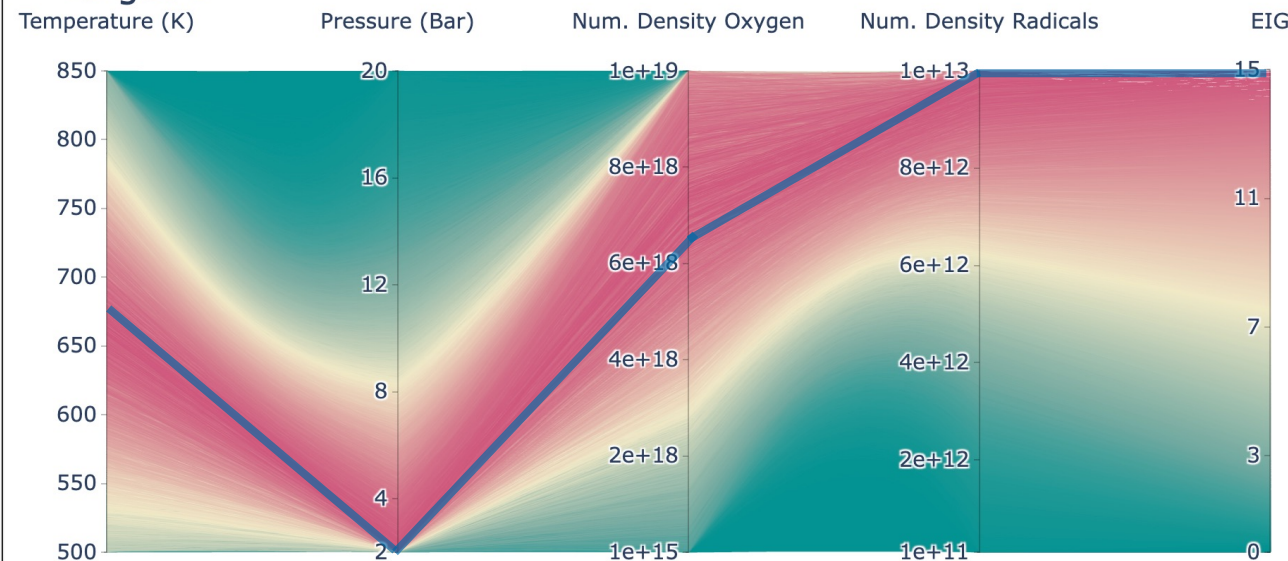
# Results

### Target 2



Optimization iterations for maximizing the expected utility. A number of random samples (black points) are evaluated in parallel prior to Bayesian Optimization (blue points), which quickly improves on the previous maxima.

### Target 2

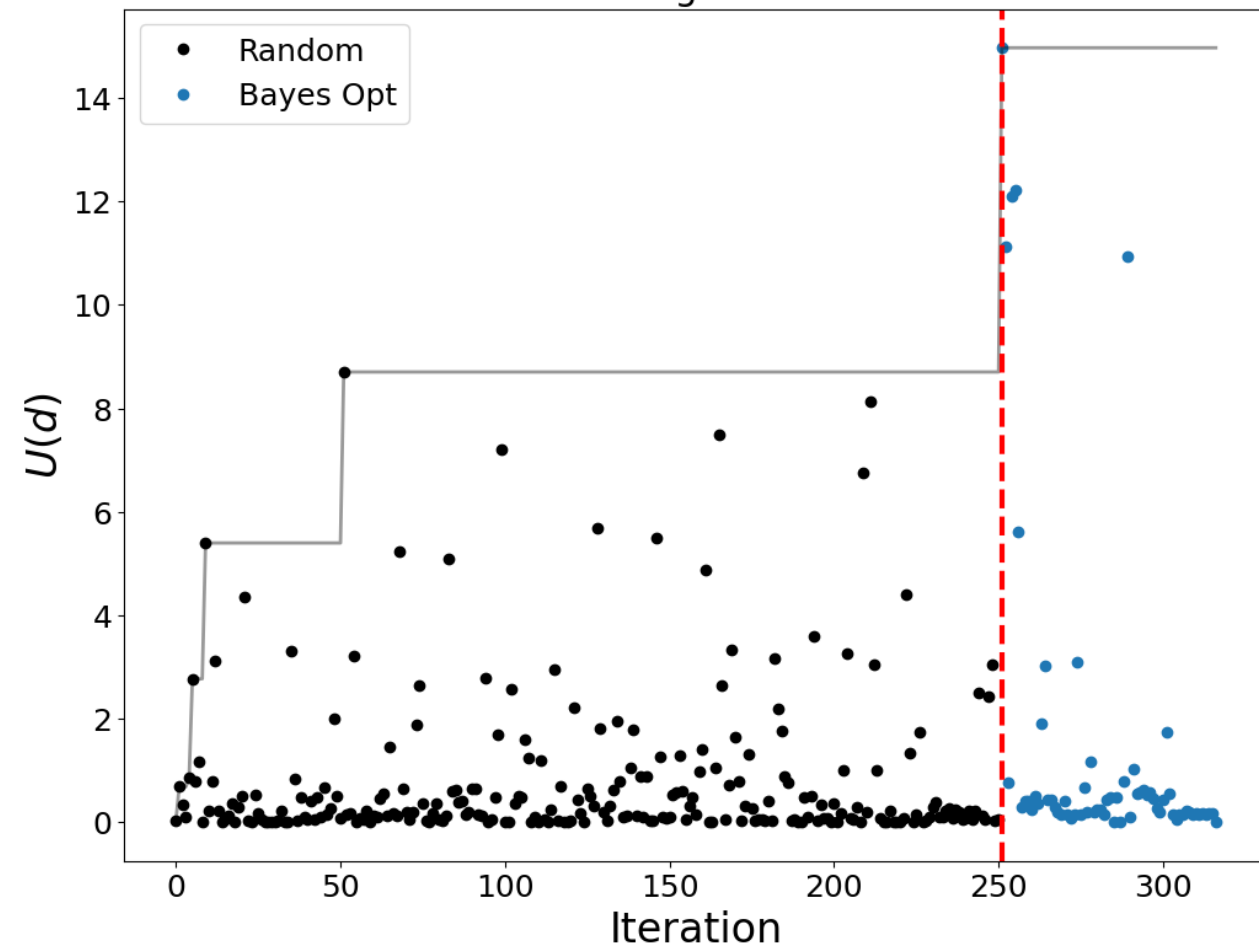


Parallel coordinate plot shows the range of design variables associated with high  $U(d)$  in red. Trajectory of optimal design is highlighted in blue.



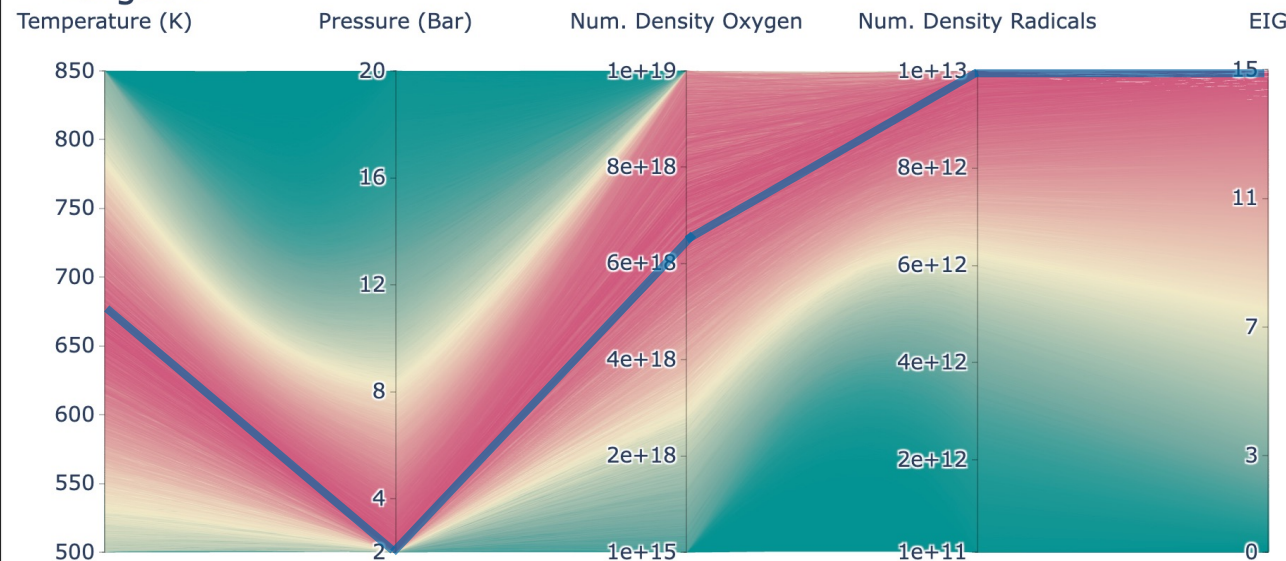
# Results

### Target 2



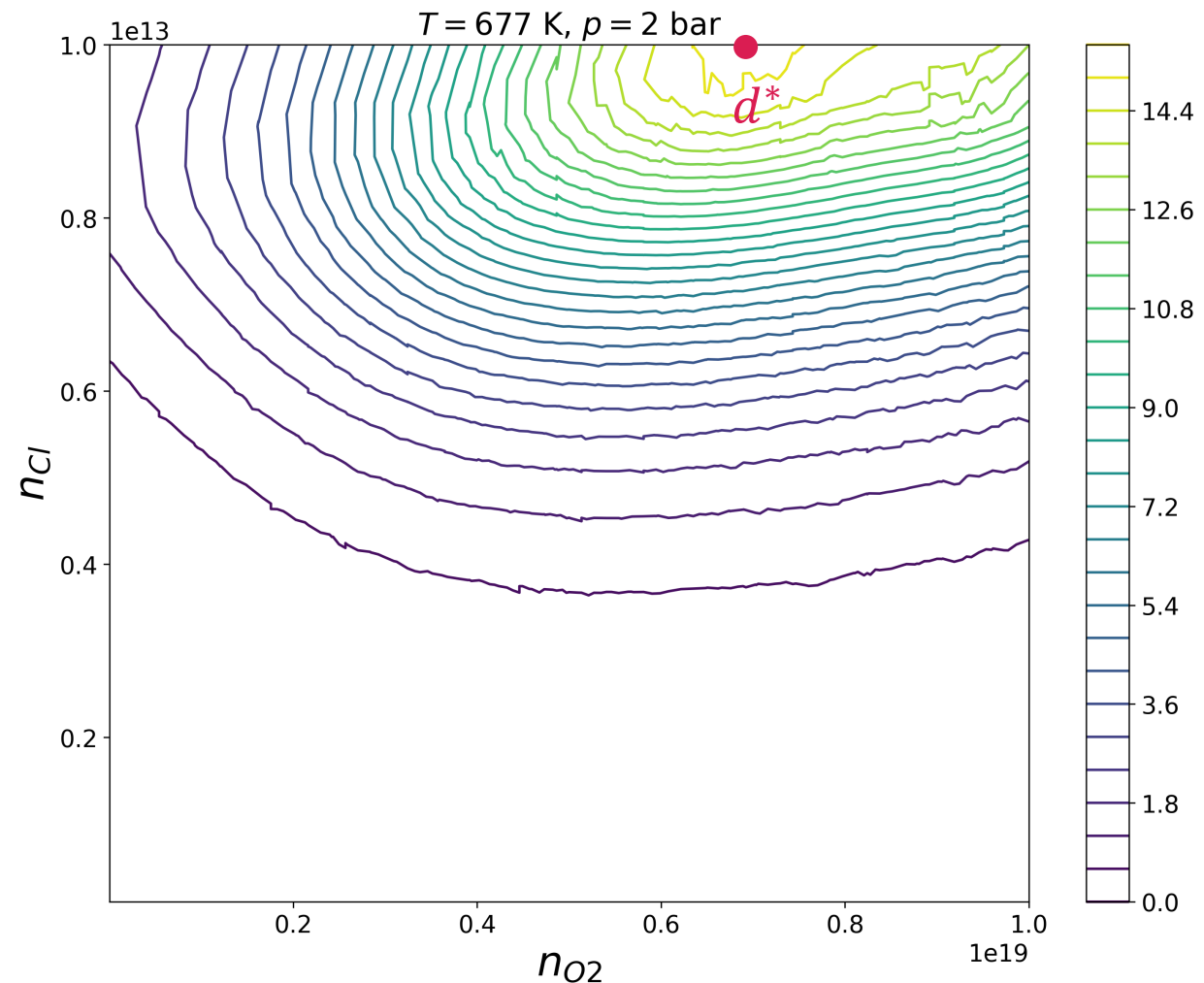
Optimization iterations for maximizing the expected utility. A number of random samples (black points) are evaluated in parallel prior to Bayesian Optimization (blue points), which quickly improves on the previous maxima.

### Target 2

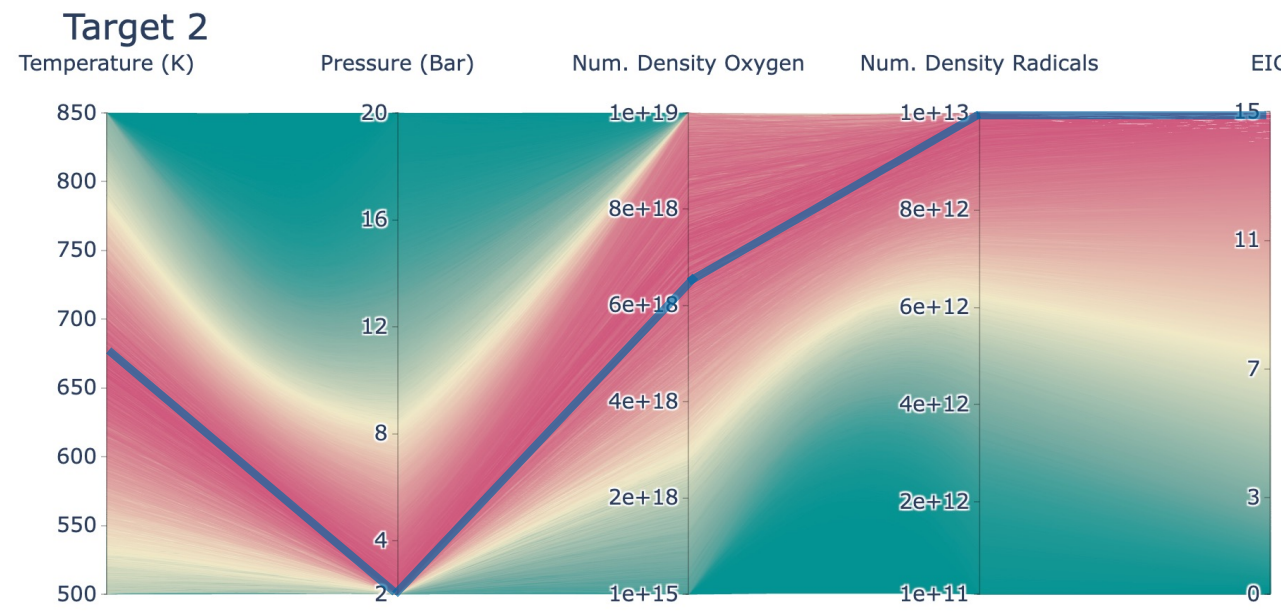


Parallel coordinate plot shows the range of design variables associated with high  $U(d)$  in red. Trajectory of optimal design is highlighted in blue.

# Results



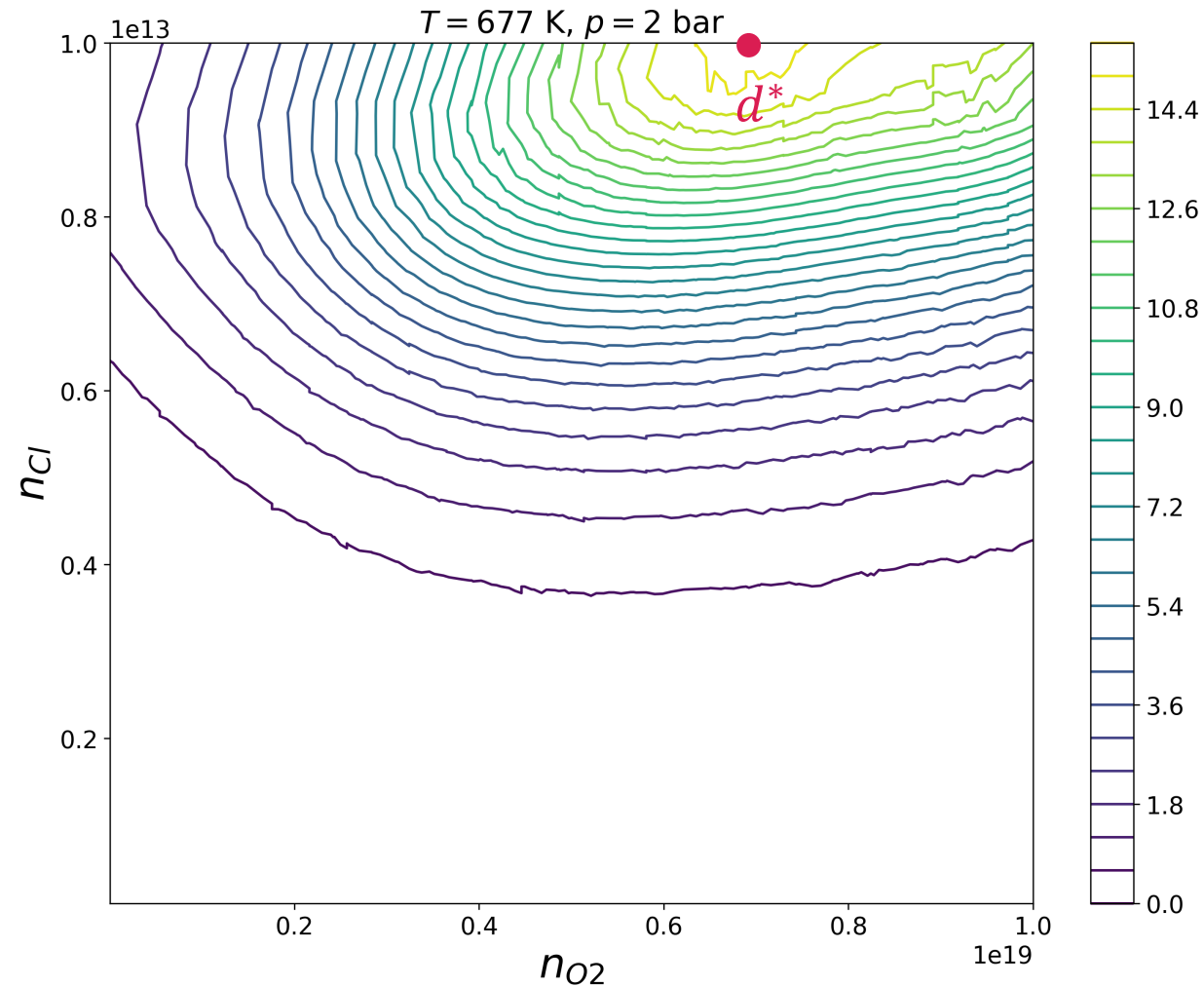
Cross-section of  $U(d)$  for a fixed temperature and pressure.



Parallel coordinate plot shows the range of design variables associated with high  $U(d)$  in red. Trajectory of optimal design is highlighted in blue.

Heatmap enables selection of a suboptimal design which can be desirable if resources are scarce or there are safety concerns

# Results



Cross-section of  $U(d)$  for a fixed temperature and pressure.

Heatmap enables selection of a suboptimal design which can be desirable if resources are scarce or there are safety concerns

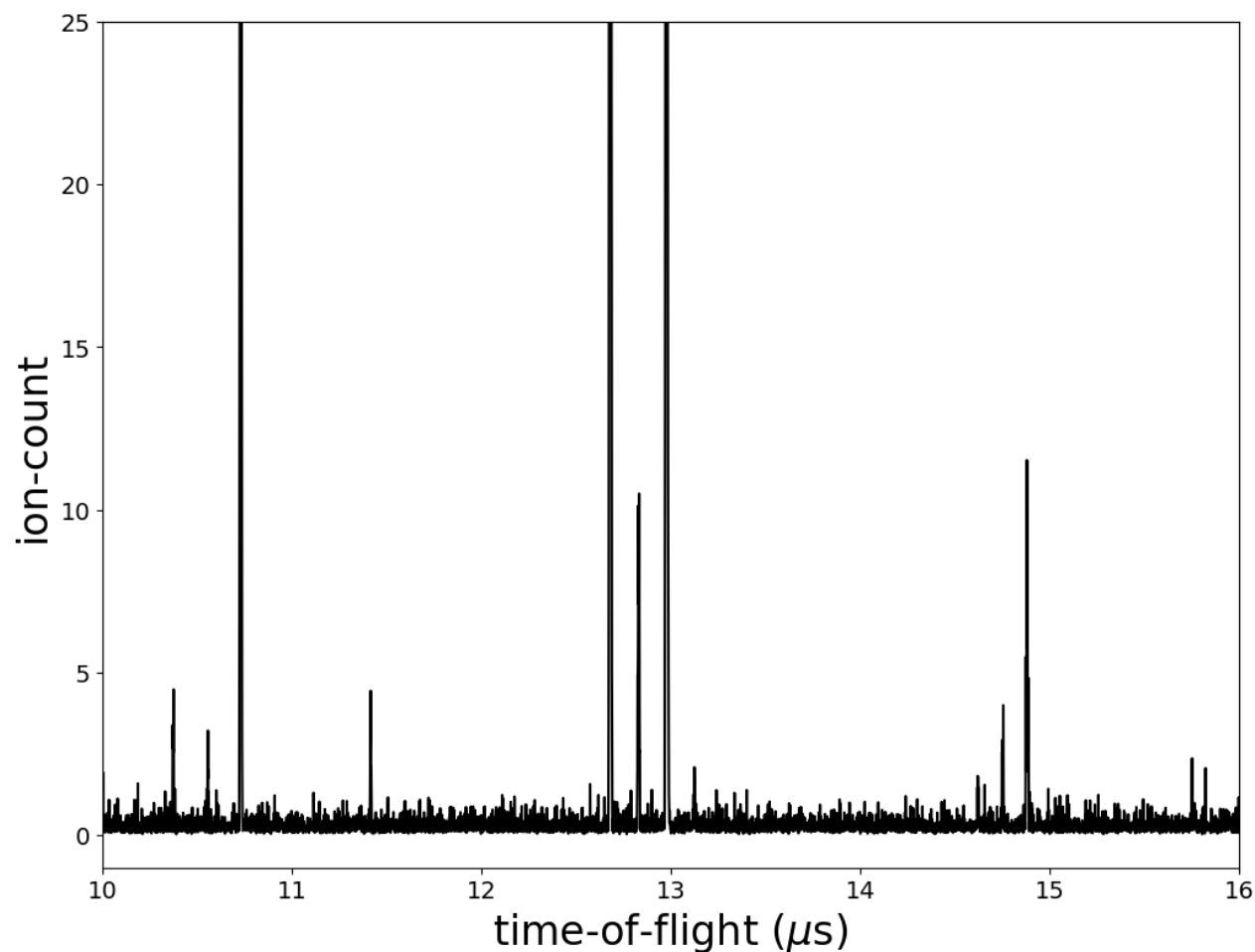
**What does the posterior distribution look like after measuring data at  $d^*$ ?**

- Generate datasets using the probabilistic model at  $d^*$
- Plausible posterior distributions are aggregated using logarithmic pooling:

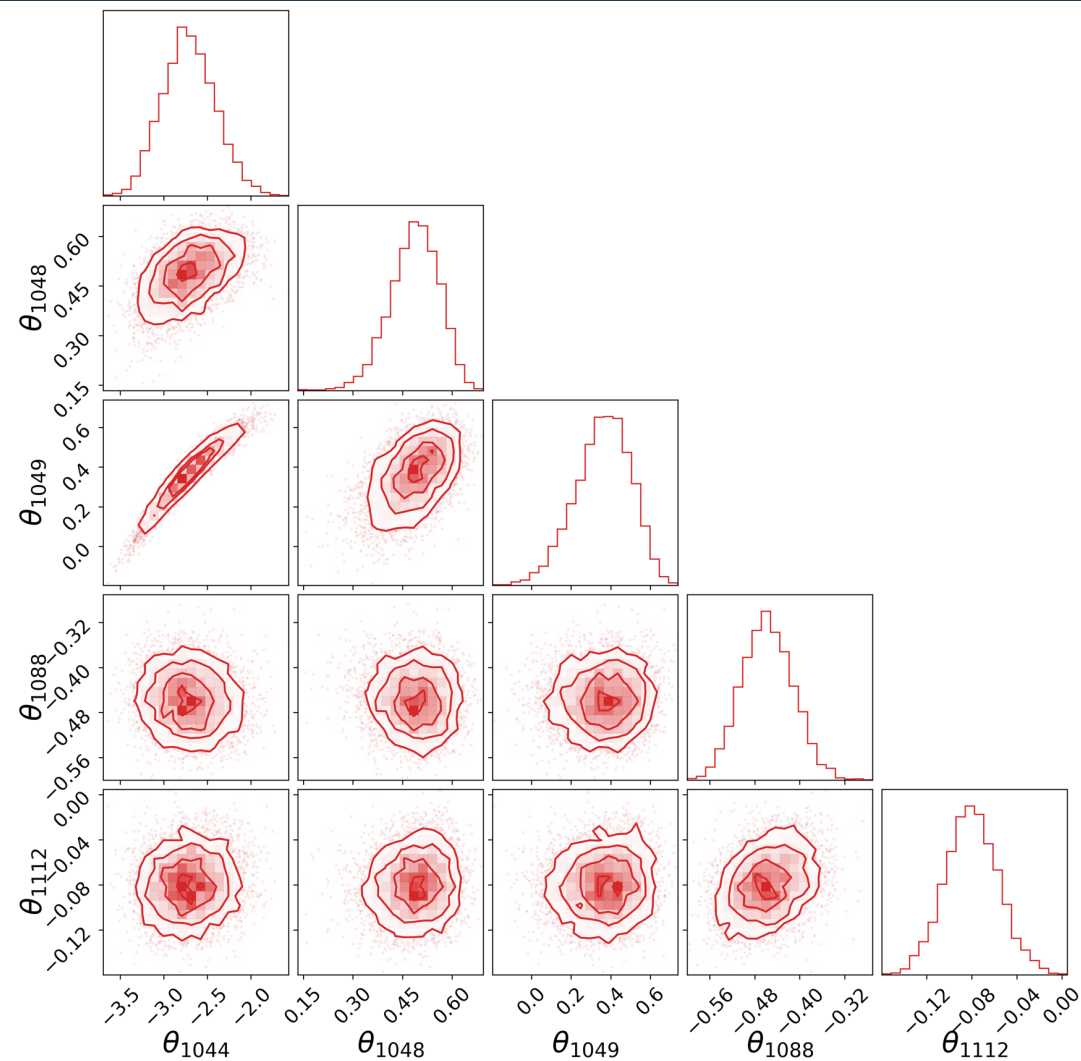
$$p(\theta)_{\log} \approx \left[ \prod_{i=1}^K p(\theta|y_i) \right]^{1/K}$$

$$p(\theta|y_i) \propto p(\theta)p(y_i|\theta)$$

# Results

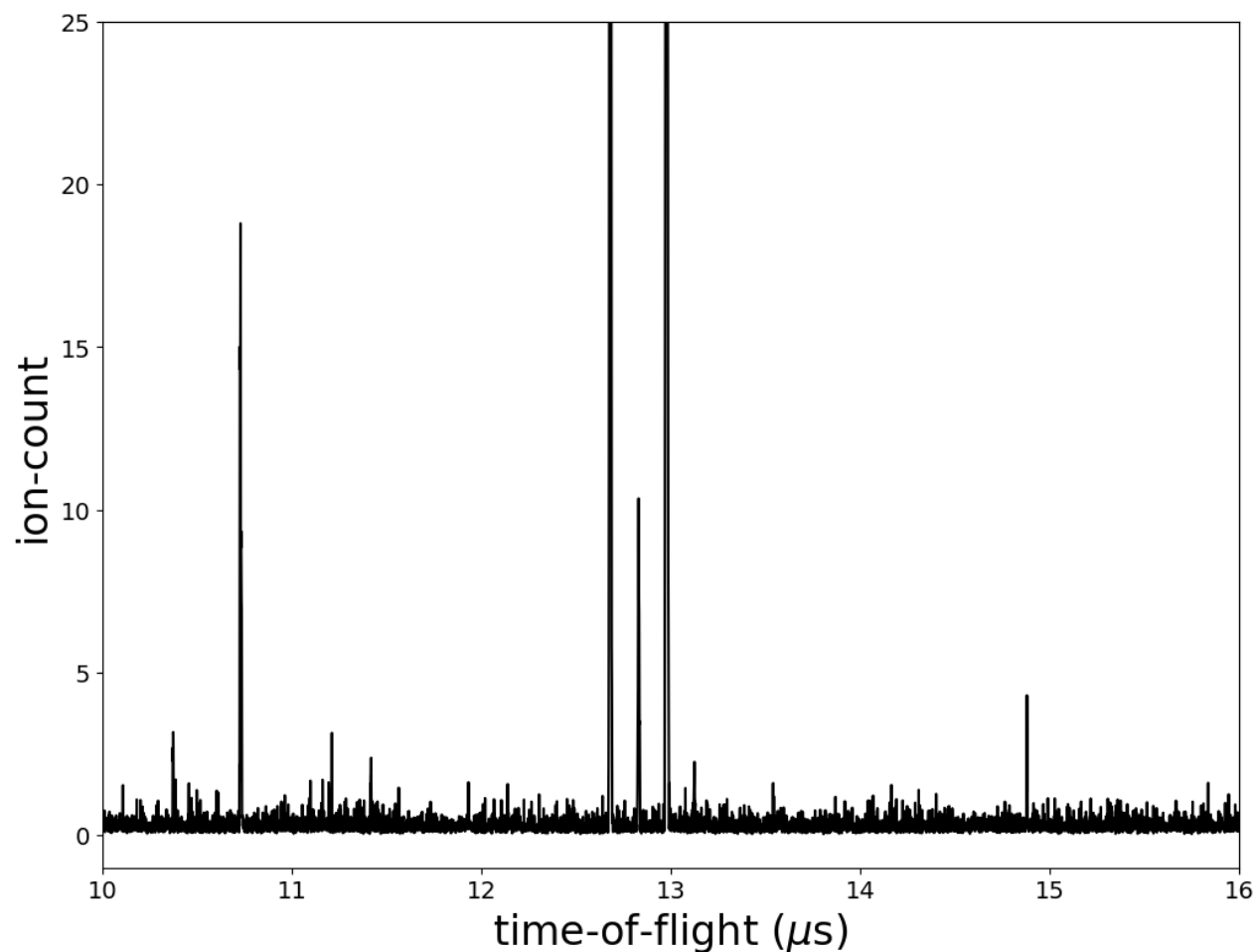


Visualization of the range of plausible datasets generated by the probabilistic model at  $d^*$  both time and VUV energy are fixed at 60ms and 11.3 eV, respectively. New species appear with significant changes to the peak magnitudes.

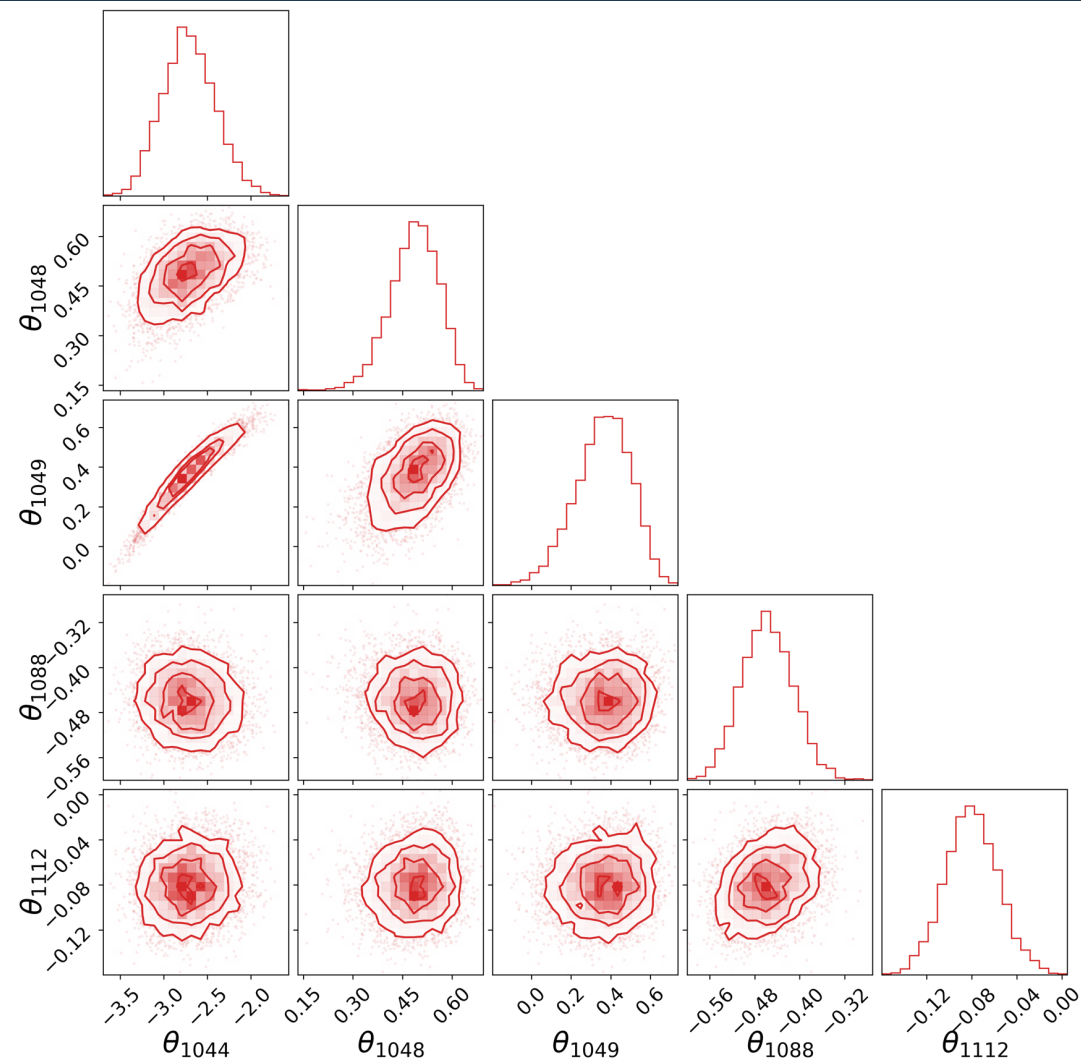


Estimated posterior after log pooling with  $K = 10^3$  plausible datasets generated at  $d^*$ . Prior distributions are a standard normal.

# Results

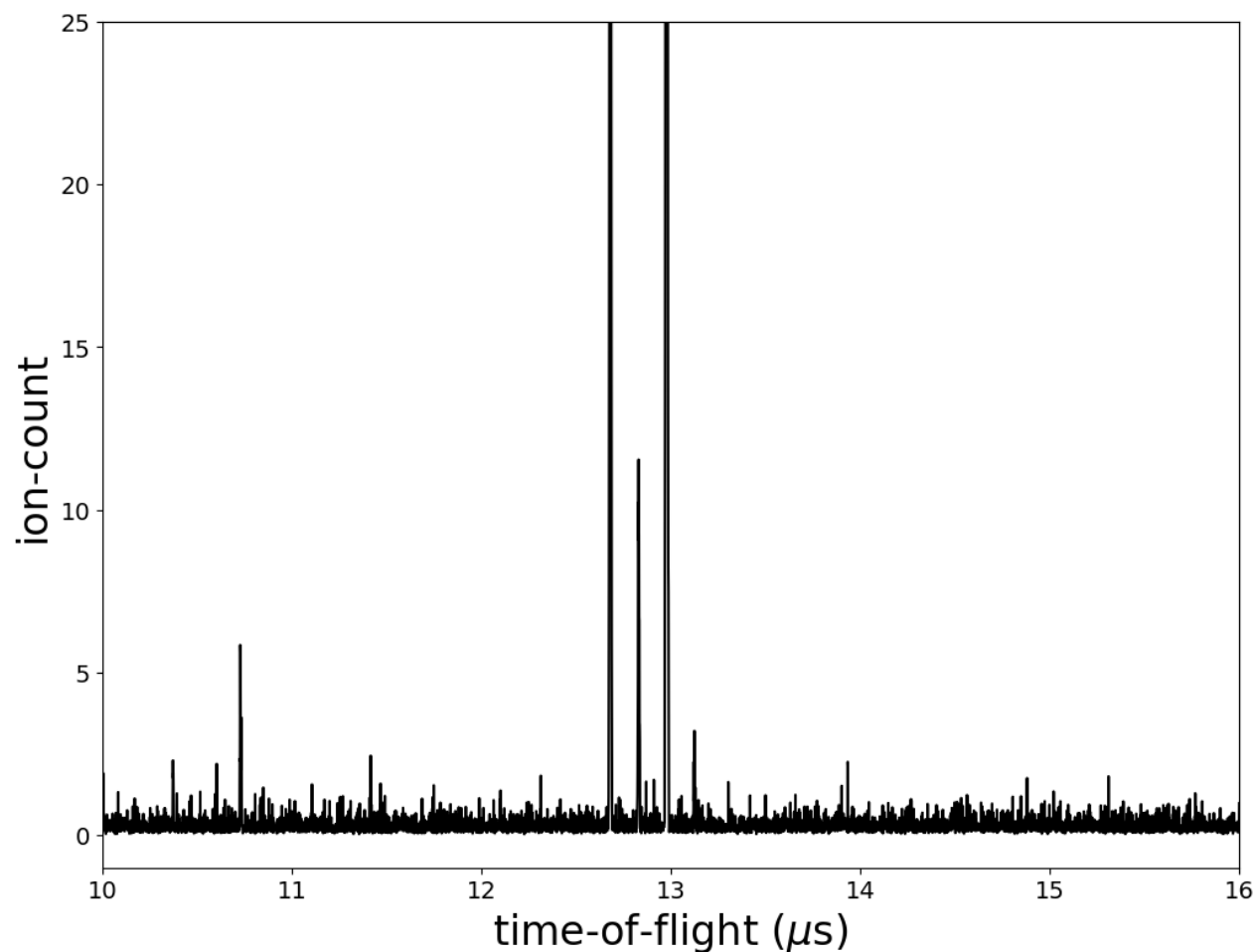


Visualization of the range of plausible datasets generated by the probabilistic model at  $d^*$  both time and VUV energy are fixed at 60ms and 11.3 eV, respectively. New species appear with significant changes to the peak magnitudes.

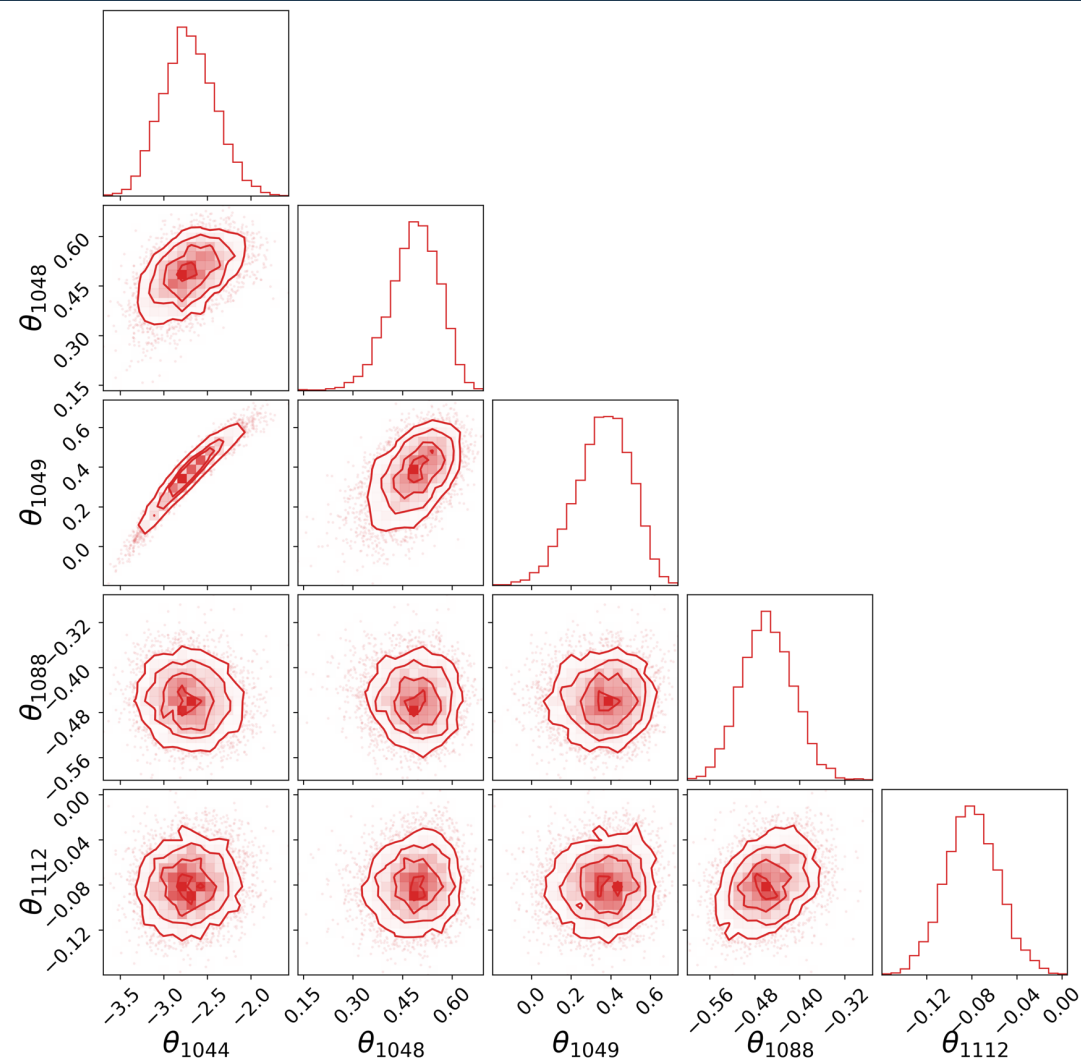


Estimated posterior after log pooling with  $K = 10^3$  plausible datasets generated at  $d^*$ . Prior distributions are a standard normal.

# Results

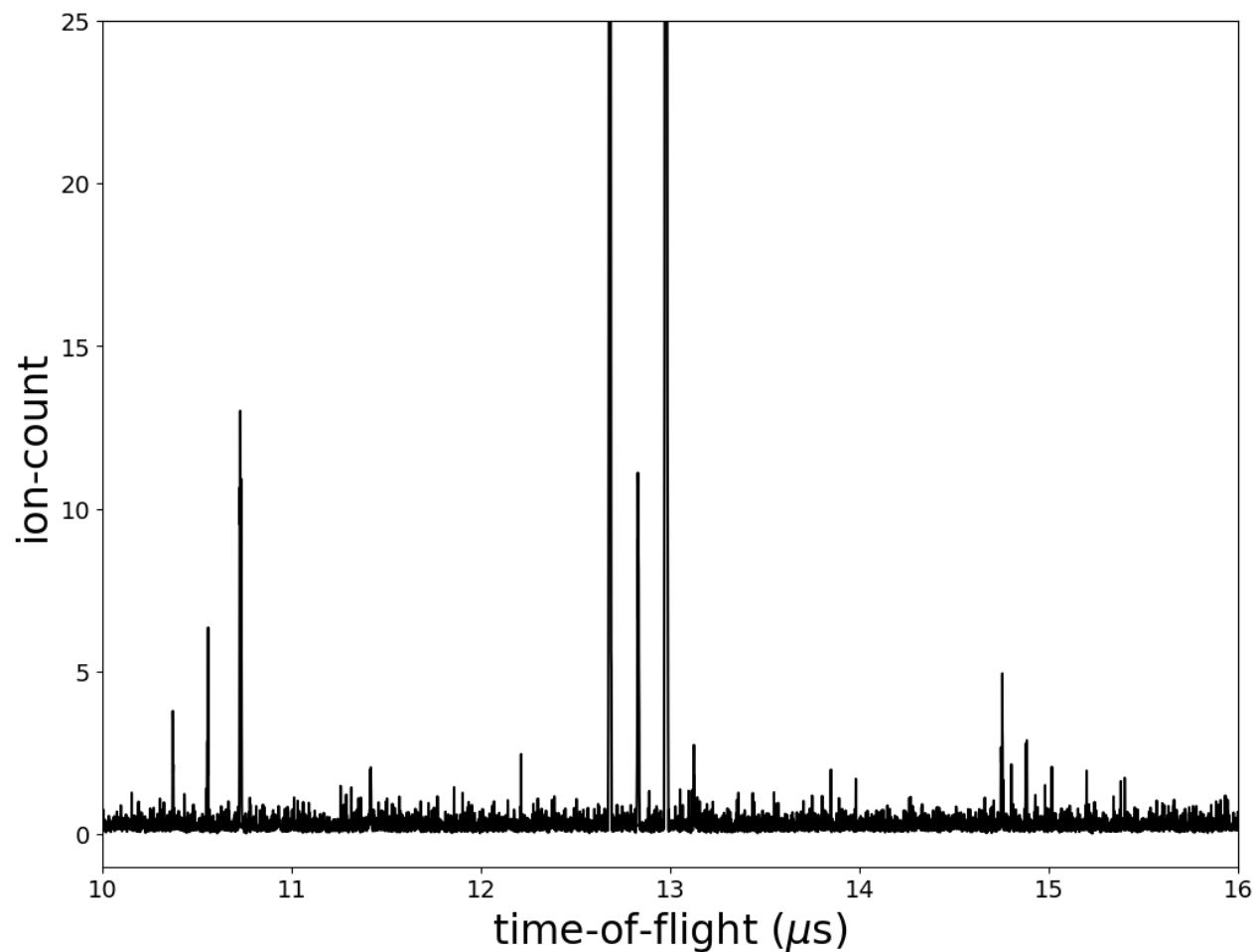


Visualization of the range of plausible datasets generated by the probabilistic model at  $d^*$  both time and VUV energy are fixed at 60ms and 11.3 eV, respectively. New species appear with significant changes to the peak magnitudes.

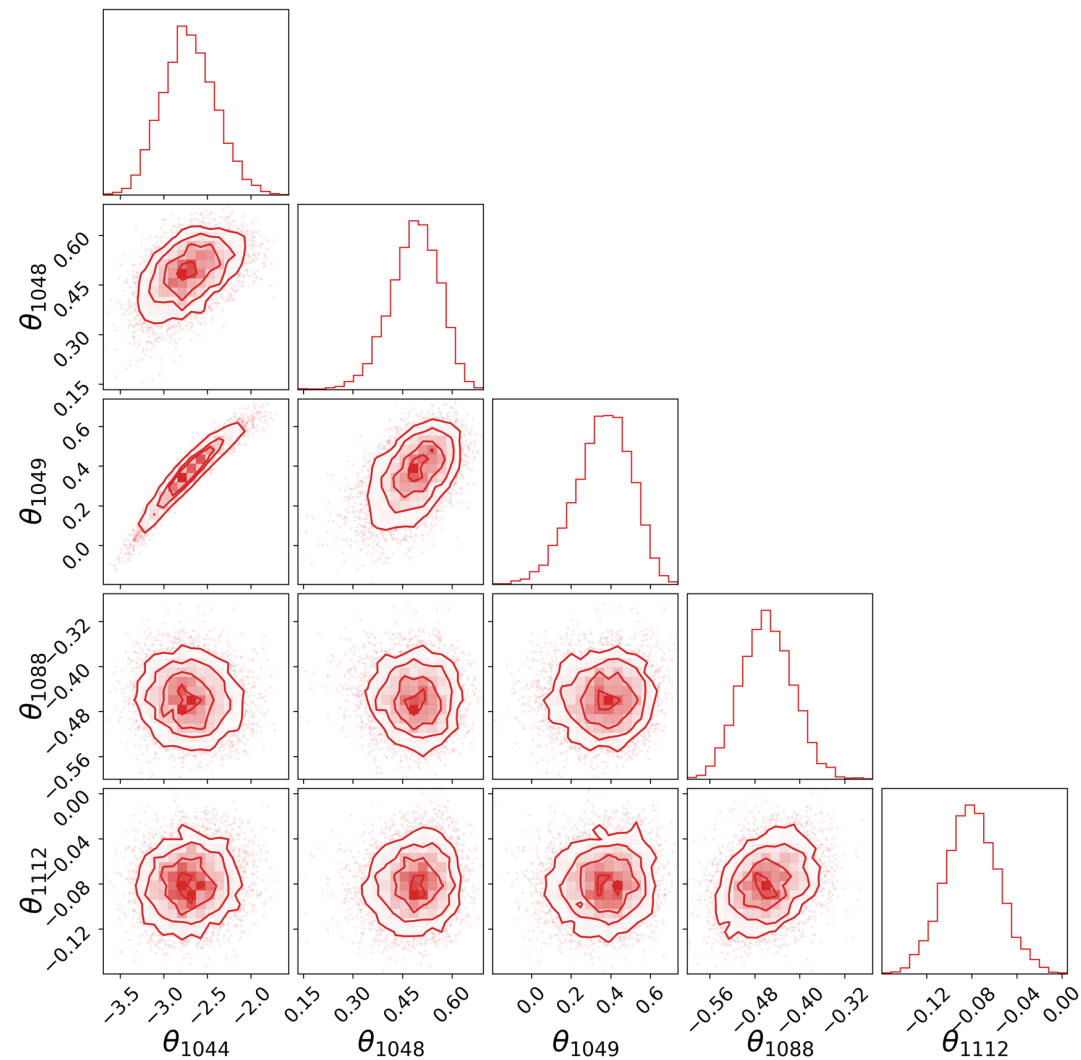


Estimated posterior after log pooling with  $K = 10^3$  plausible datasets generated at  $d^*$ . Prior distributions are a standard normal.

# Results

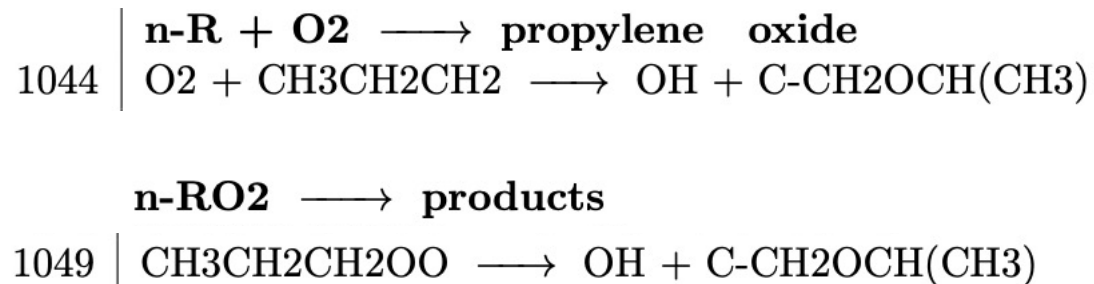


Visualization of the range of plausible datasets generated by the probabilistic model at  $d^*$  both time and VUV energy are fixed at 60ms and 11.3 eV, respectively. New species appear with significant changes to the peak magnitudes.

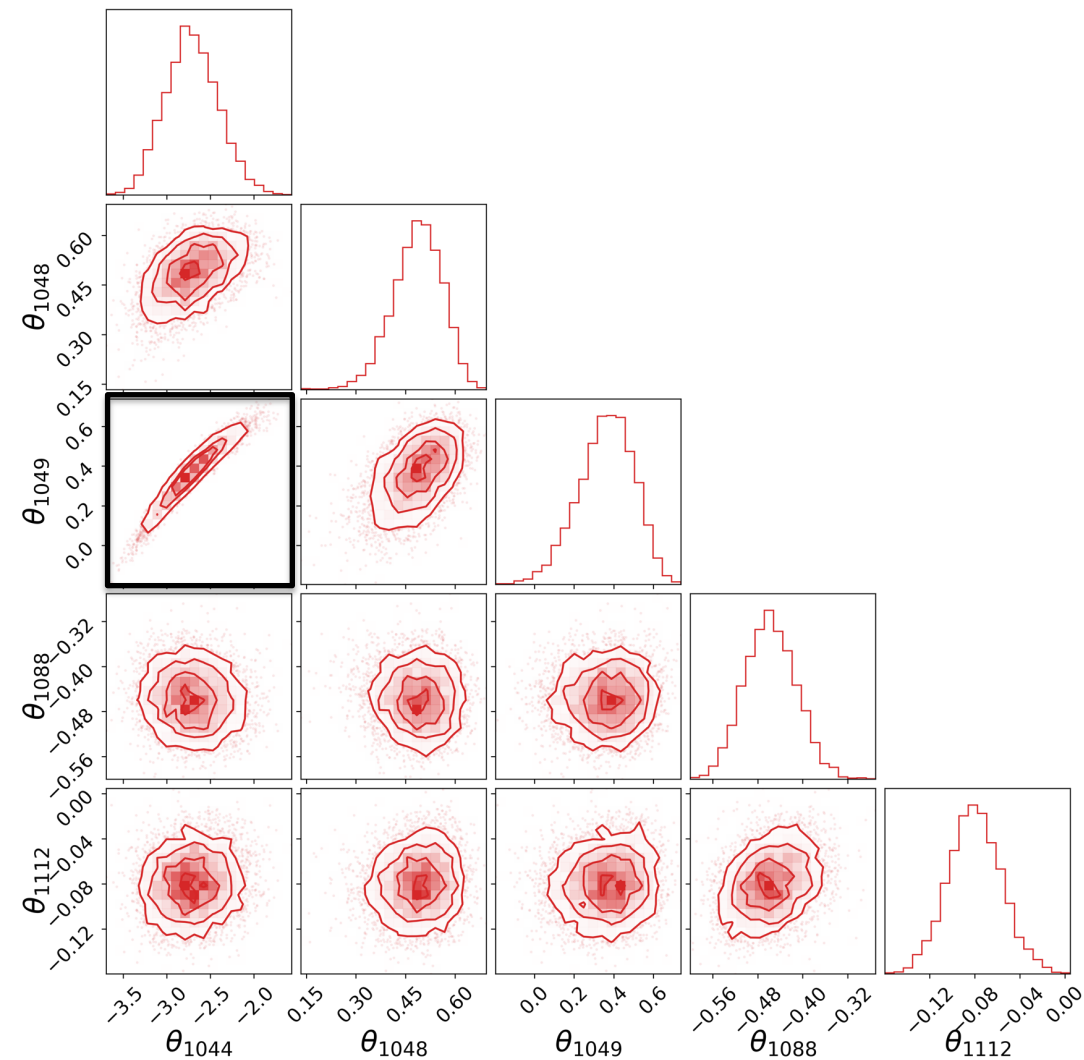


Estimated posterior after log pooling with  $K = 10^3$  plausible datasets generated at  $d^*$ . Prior distributions are a standard normal.

# Results



Posterior distribution provides correlations between two of the sensitive reaction rates, both sharing propylene oxide as a product.



Estimated posterior after log pooling with  $K = 10^3$  plausible datasets generated at  $d^*$ . Prior distributions are a standard normal.



# Conclusion & Future Work

- Demonstrated feasibility and workflow for applying OED for high-dimensional physics-based models
- Bayesian optimization enables efficient exploration of noisy utility functions
- Low-dimensional surrogates of model output provide necessary computational savings for assessing complex models

---

## Future Work

- Demonstrate improvement in the posterior density using measurements at optimal design compared to random designs
- Update initial prior estimates utilizing available measurement data

# Acknowledgements

We would like to thank Oscar Diaz-Ibarra, Kyungjoo Kim, Arun Hegde, Cosmin Safta, and Khachik Sargsyan for helpful conversations about this work.

This work was supported by the US Department of Energy (DOE), Office of Basic Energy Sciences (BES) Division of Chemical Sciences, Geosciences, and Biosciences.

Sandia National Laboratories is a multimission laboratory managed and operated by National Technology and Engineering Solutions of Sandia, LLC., a wholly owned subsidiary of Honeywell International, Inc., for the U.S. Department of Energy's National Nuclear Security Administration under contract DE-NA-0003525.



**Sandia  
National  
Laboratories**

# Additional Slides

# Challenges

## Large number of uncertain model parameters

- Physics and instrument model contain 1151 uncertain model parameters
- Only a *small fraction* have a significant impact on the spectrum at each design

---

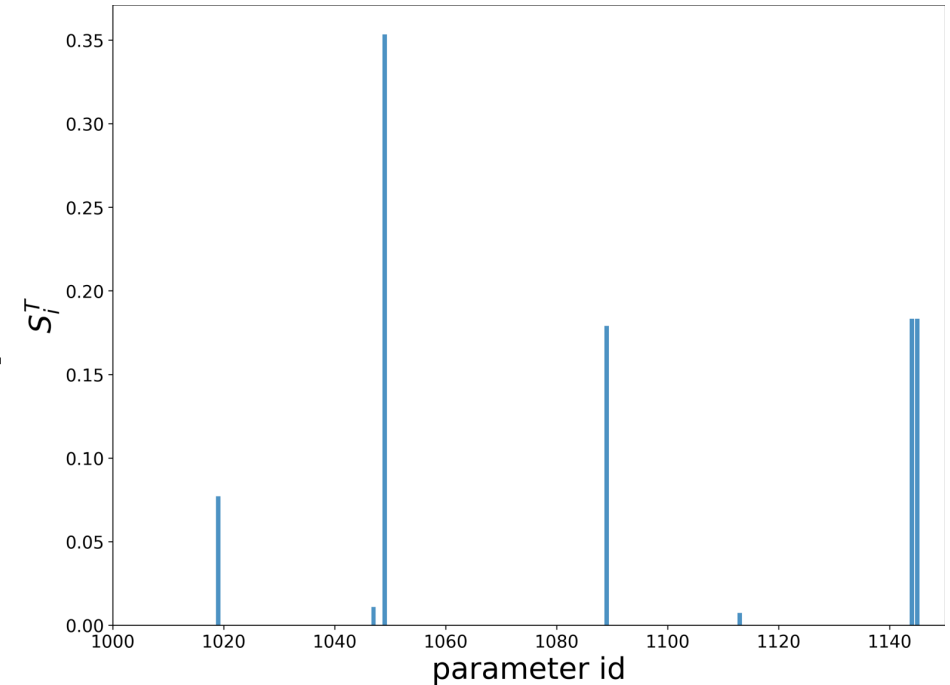
## Variance-based Global Sensitivity Analysis

$$S_i^T = 1 - \frac{\text{Var}_{\theta \sim i} [\mathbb{E}_{\theta_i} (f(\boldsymbol{\theta}) | \theta_{\sim i})]}{\text{Var}_{\boldsymbol{\theta}} [f(\boldsymbol{\theta})]}$$

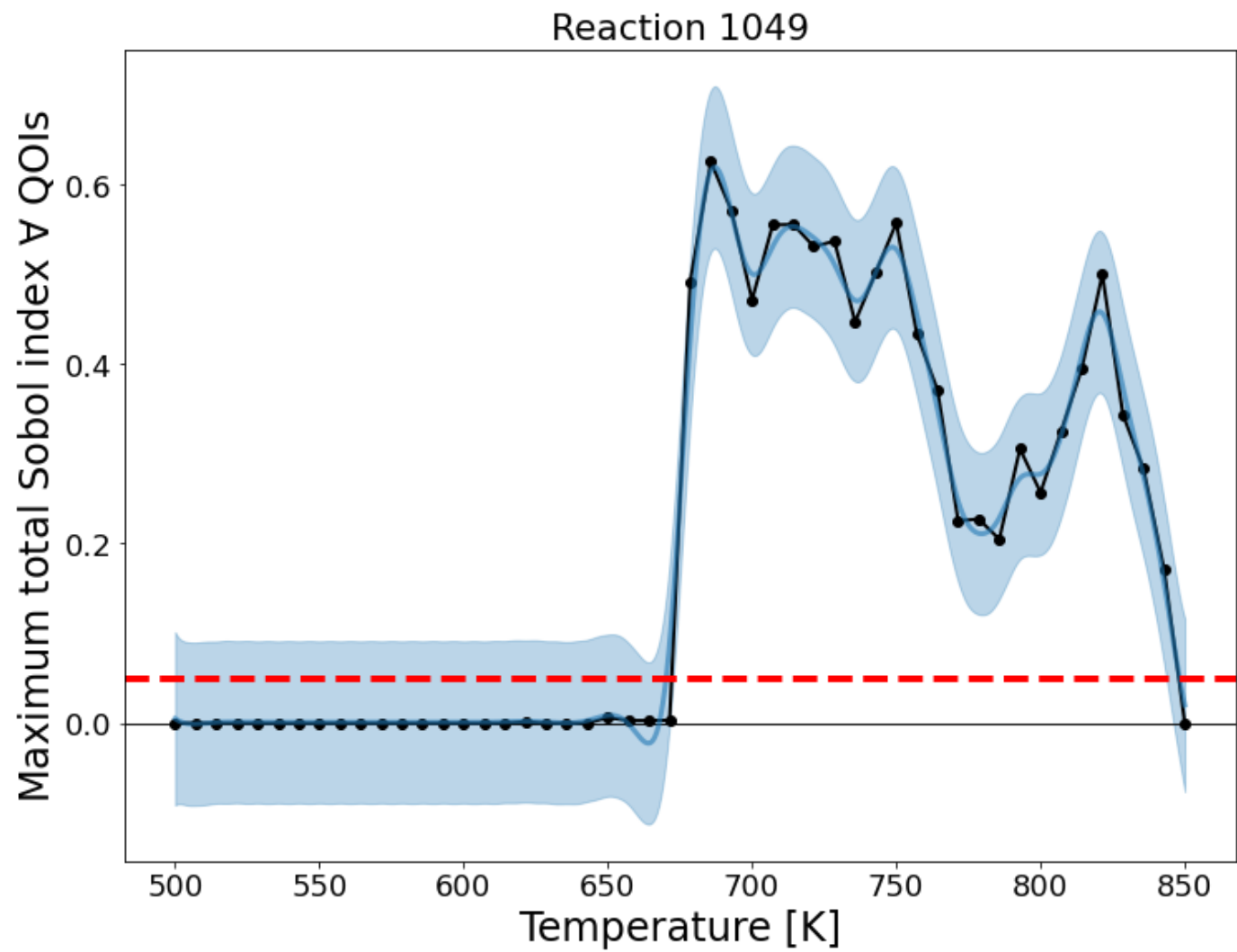
- $\ell_1$  regularization of the polynomial chaos expansion coefficients

$$\mathbf{c}^{CS} = \arg \min_{\mathbf{c}} \|\mathbf{y} - \mathbf{c}\Psi\|_2 + \gamma \|\mathbf{c}\|_1$$

Bruno Sudret. Global sensitivity analysis using polynomial chaos expansions. *Reliability engineering & system safety* 93, no. 7 (2008): 964-979



Only a small number of nonzero total Sobol indices for design conditions 620 K, 2 bar,  $n_{\text{C}_3\text{H}_8} = 2\text{e}14$ ,  $n_{\text{O}_2} = 8\text{e}18$ ,  $n_{\text{Cl}} = 1\text{e}13$ .



# Reducing output dimensionality

**Goal:** Map high-dimensional model output to a lower-dimensional space while minimizing loss of information

$$g : \mathbb{R}^J \mapsto \mathbb{R}^K$$

## Autoencoder

$$\begin{aligned} \text{Encoder network:} \quad & q = g_\psi(y) & q \in \mathbb{R}^K \\ \text{Decoder network:} \quad & \hat{y} = h_\phi(g_\psi(y)) & y \in \mathbb{R}^J \\ & & K \ll J \end{aligned}$$

Minimize reconstruction loss

$$L(\phi, \psi) = \min_{\phi, \psi} \|y - h_\phi(g_\psi(y))\|_2^2$$

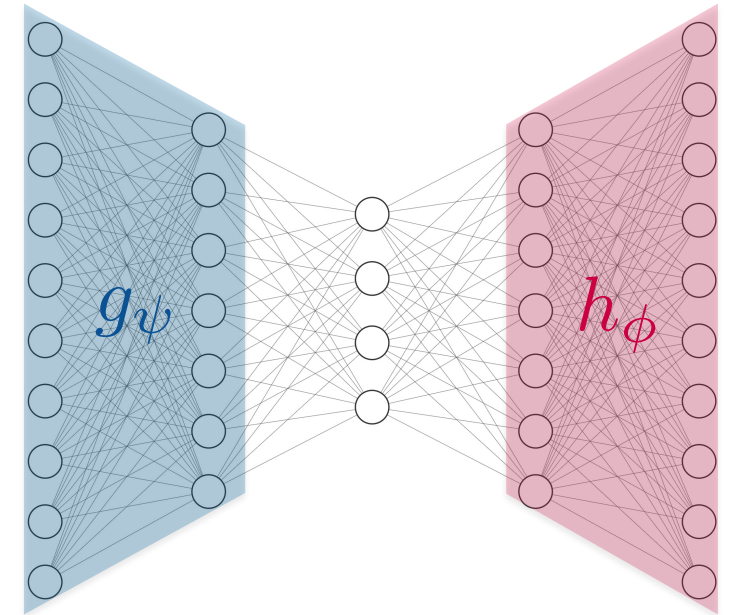
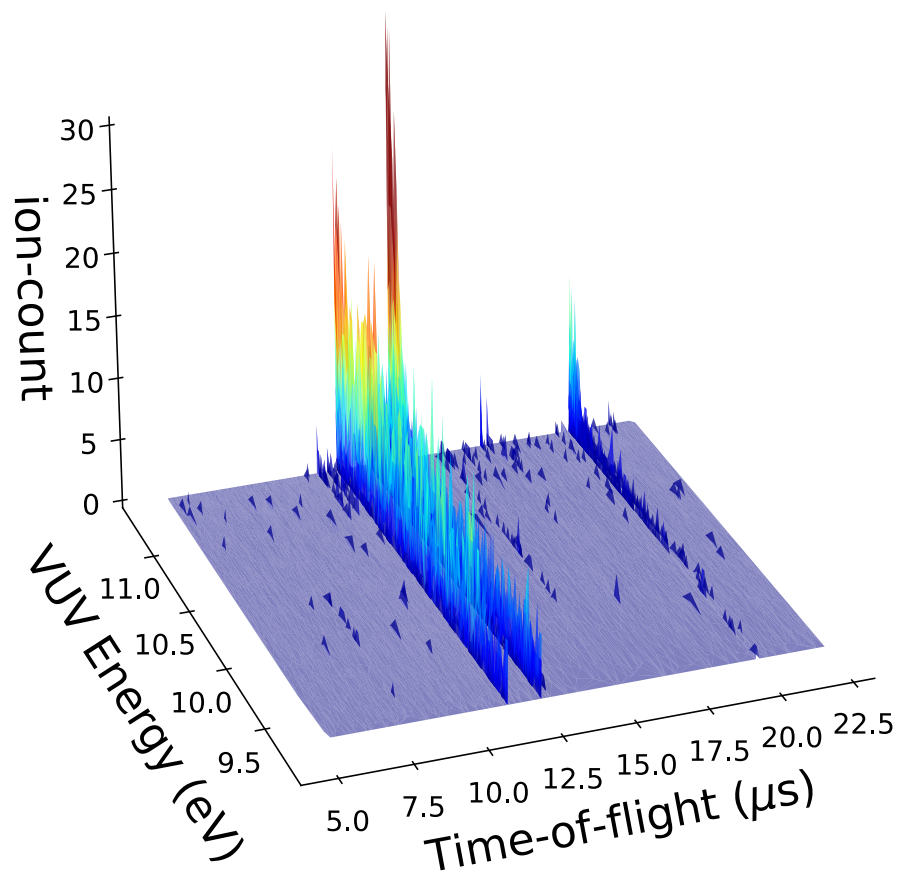
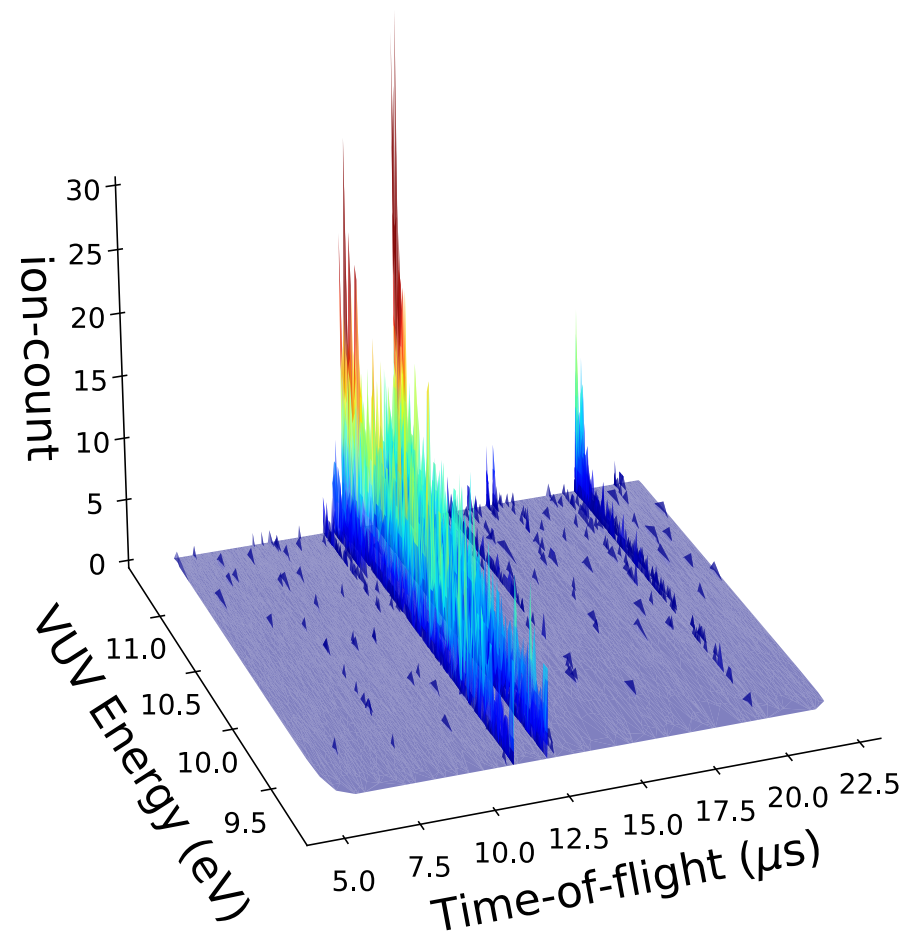


Illustration of an autoencoder architecture

$t = 0$  ms

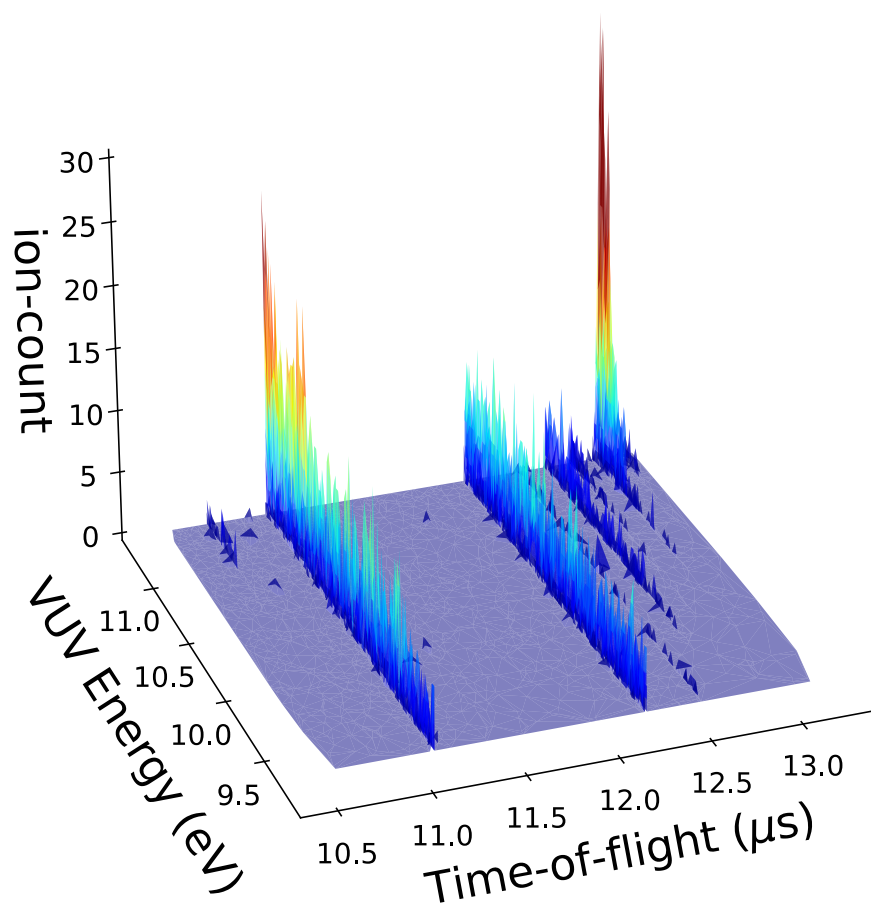


$t = 60$  ms

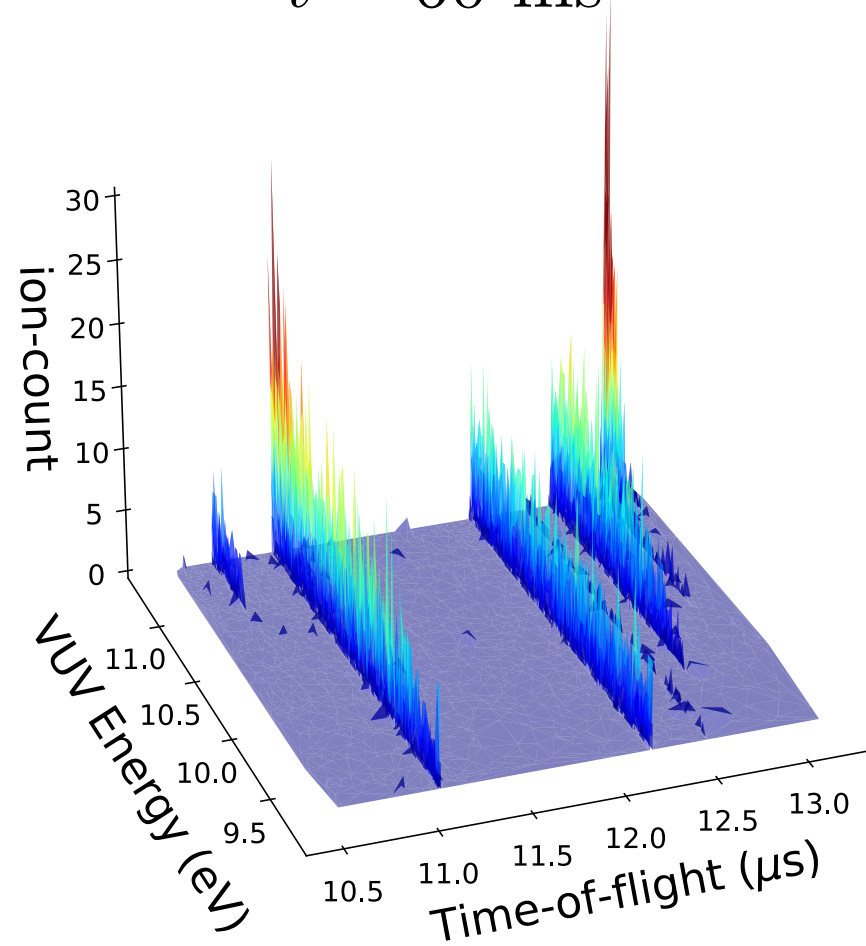


Experiment conducted at  $T = 700$  K,  $p = 10$  bar,  $\chi_{\text{C}_3\text{H}_8} = 9.7 \times 10^{-7}$ ,  $\chi_{\text{O}_2} = 2.9 \times 10^{-2}$ ,  $\chi_{\text{pre}} = 2.2 \times 10^{-4}$  with a helium bath gas

$t = 0$  ms



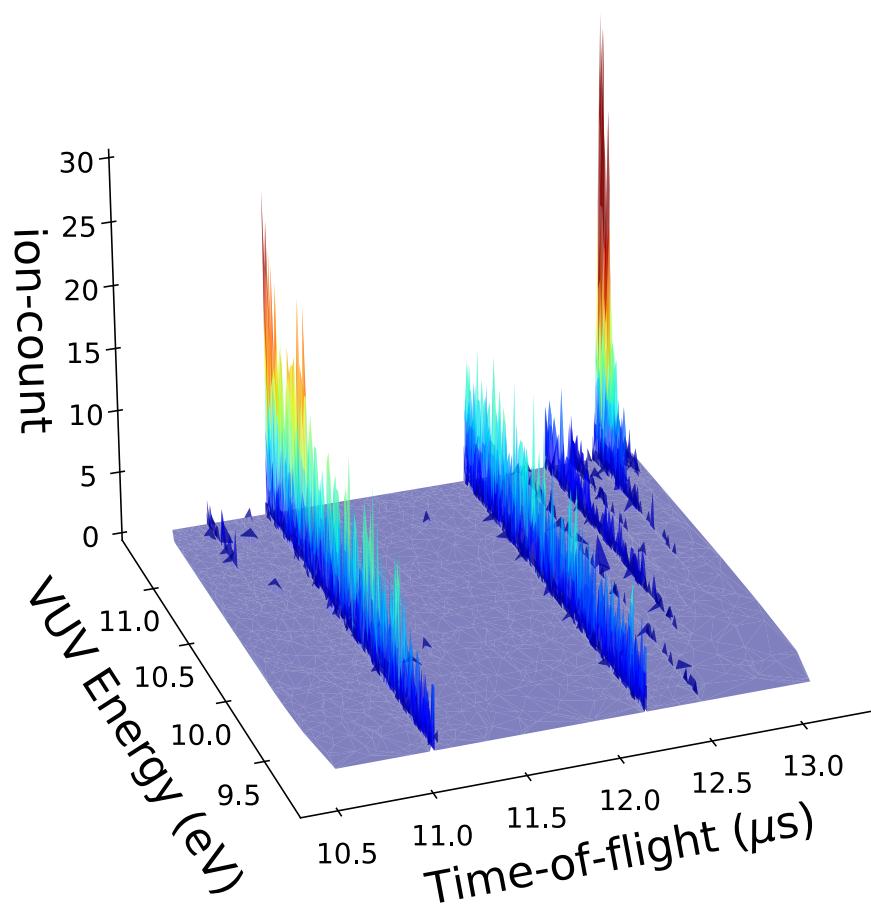
$t = 60$  ms



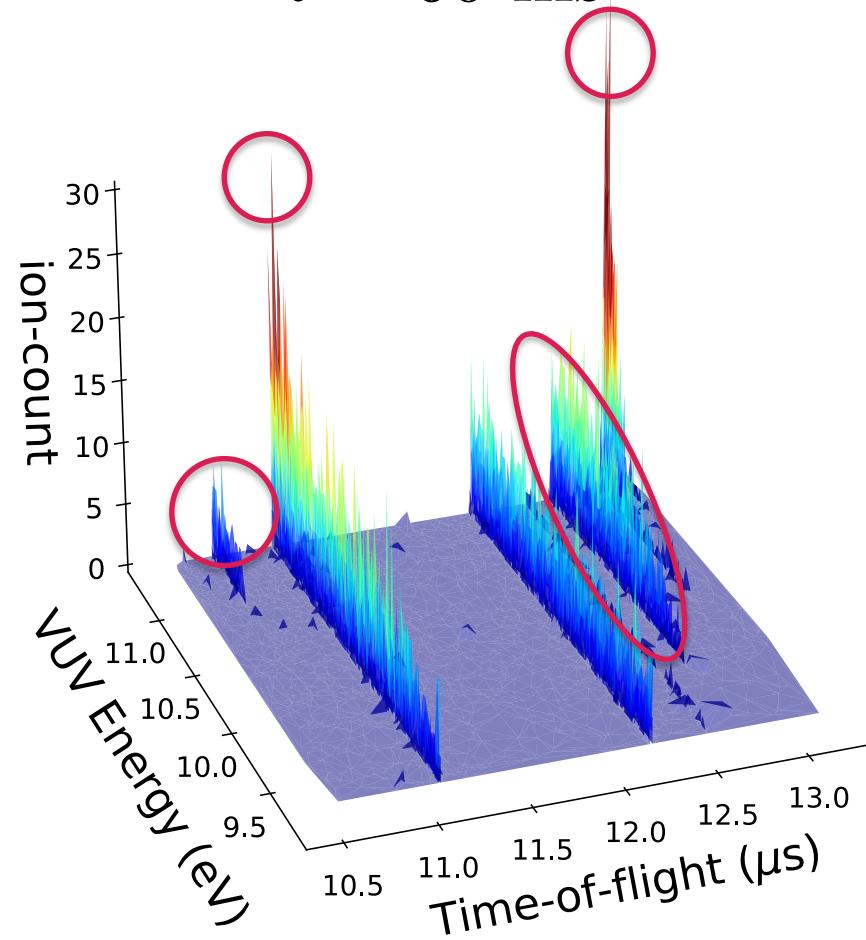
Experiment conducted at  $T = 700$  K,  $p = 10$  bar,  $\chi_{\text{C}_3\text{H}_8} = 9.7 \times 10^{-7}$ ,  $\chi_{\text{O}_2} = 2.9 \times 10^{-2}$ ,  $\chi_{\text{pre}} = 2.2 \times 10^{-4}$  with a helium bath gas



$t = 0$  ms

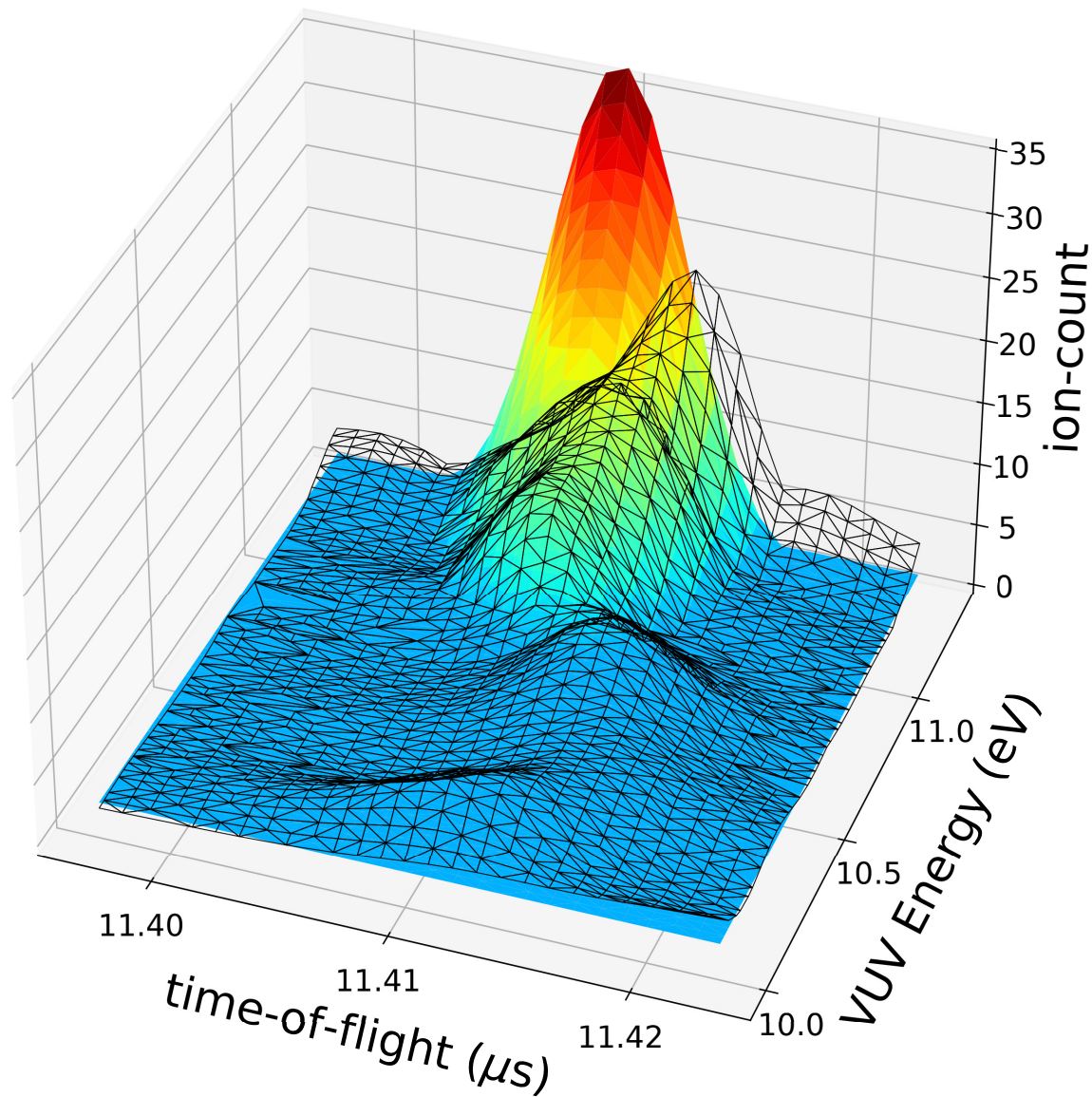


$t = 60$  ms



Experiment conducted at  $T = 700$  K,  $p = 10$  bar,  $\chi_{\text{C}_3\text{H}_8} = 9.7 \times 10^{-7}$ ,  $\chi_{\text{O}_2} = 2.9 \times 10^{-2}$ ,  $\chi_{\text{pre}} = 2.2 \times 10^{-4}$  with a helium bath gas

# Modeling the HPR experiment



**Solid surface** is the prediction of  $f(\theta_{MAP}, d, x)$  for one of the peaks in the time-of-flight spectrum ( $\text{H}_2\text{O}_2$ ).

**Mesh surface** shows the prediction with model error,  $f(\theta_{MAP}, d, x) + \mu_\delta(x)$  which increases the fidelity of the predictive model.

# Numerical approximation

$$u(y, d, \theta) = D_{KL}(p(\theta|y, d)||p(\theta)) = \int_{\Theta} p(\theta|y, d) \log \left[ \frac{p(\theta|y, d)}{p(\theta)} \right] d\theta = u(y, d)$$

$$\begin{aligned} U(d) &= \int_{\mathcal{Y}} \int_{\Theta} u(y, d) p(\theta|y, d) d\theta p(y|d) dy \\ &= \int_{\mathcal{Y}} u(y, d) p(y|d) dy \\ &= \int_{\mathcal{Y}} \left( \int_{\Theta} p(\theta|y, d) \log \left[ \frac{p(\theta|y, d)}{p(\theta)} \right] d\theta \right) p(y|d) dy \end{aligned}$$

Using  $p(\theta|y, d) = p(y|\theta, d)p(\theta)/p(y|d)$ ,

$$\begin{aligned} U(d) &= \int_{\mathcal{Y}} \int_{\Theta} \log \left[ \frac{p(y|\theta, d)}{p(y|d)} \right] p(y|\theta, d) p(\theta) d\theta dy \\ &= \int_{\mathcal{Y}} \int_{\Theta} [\log p(y|\theta, d) - \log p(y|d)] p(y|\theta, d) p(\theta) d\theta dy \end{aligned}$$