

# Machine-learning assisted workflows for inorganic molecular discovery

Jon Paul Janet<sup>1</sup>    Chenru Duan<sup>2</sup>    Aditya Nandy<sup>2</sup>  
Heather Kulik<sup>1</sup>

<sup>1</sup>Department of Chemical Engineering, Massachusetts Institute of Technology

<sup>2</sup>Department of Chemistry, Massachusetts Institute of Technology



Machine Learning and Informatics for Chemistry and Materials

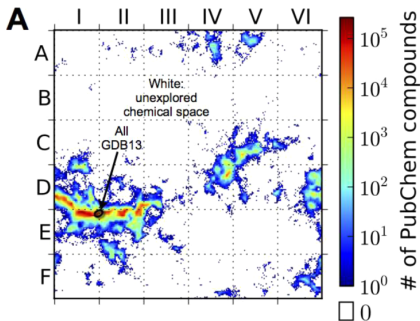
10.02.19

# Motivation: chemical discovery

## How can we design new materials using computers?

The space of possible chemistries is incredibly vast, with  $\mathcal{O}(10^{60})$  small organic molecules.

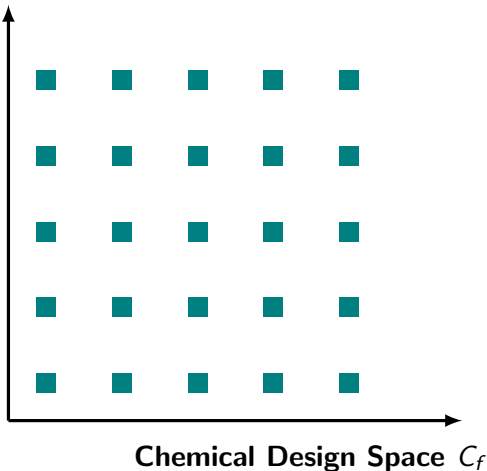
All potentially undiscovered medicines, catalysts and materials are somewhere, out in this huge space.



Virshup *et al.*, *J. Am. Chem. Soc.*, 135(19): 7296–7303, 2013.

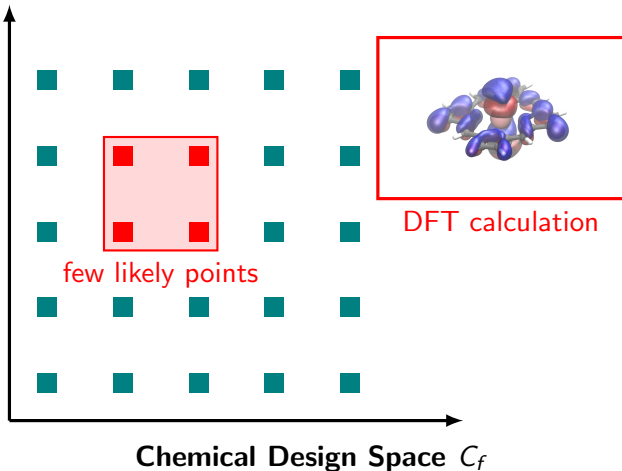
# Motivation: chemical discovery

How can we design new materials using computers?



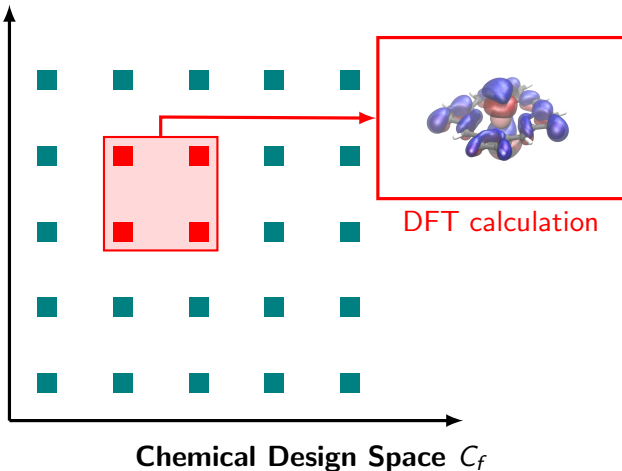
# Motivation: chemical discovery

How can we design new materials using computers?



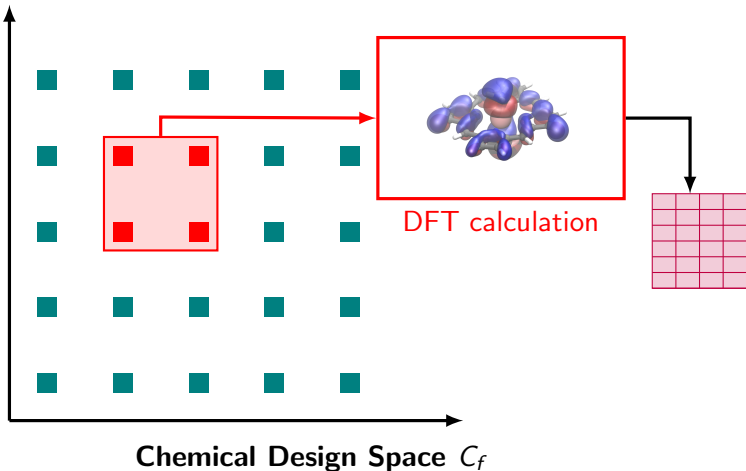
# Motivation: chemical discovery

How can we design new materials using computers?



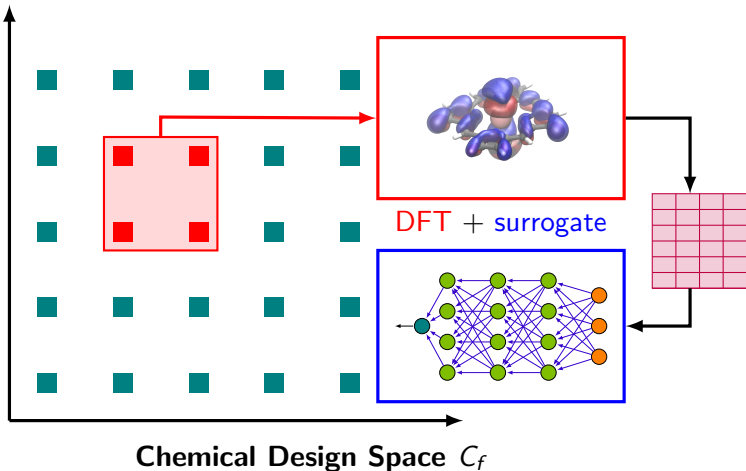
# Motivation: chemical discovery

How can we design new materials using computers?



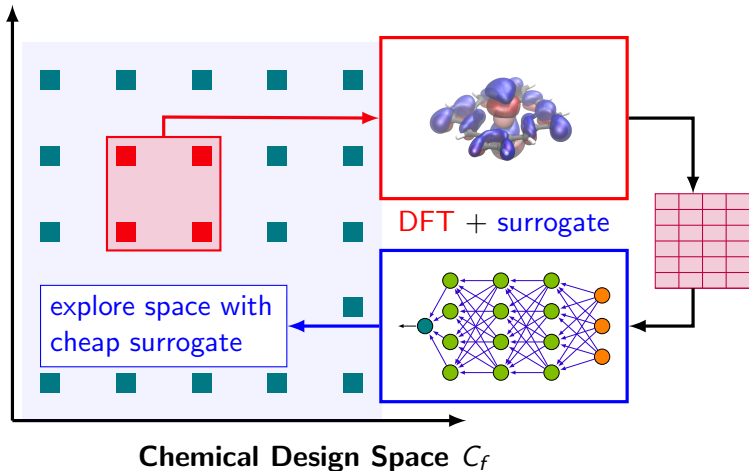
# Motivation: chemical discovery

How can we design new materials using computers?



# Motivation: chemical discovery

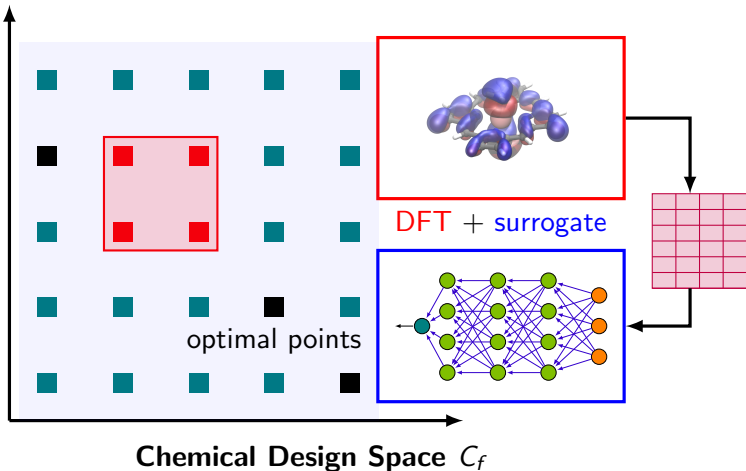
How can we design new materials using computers?





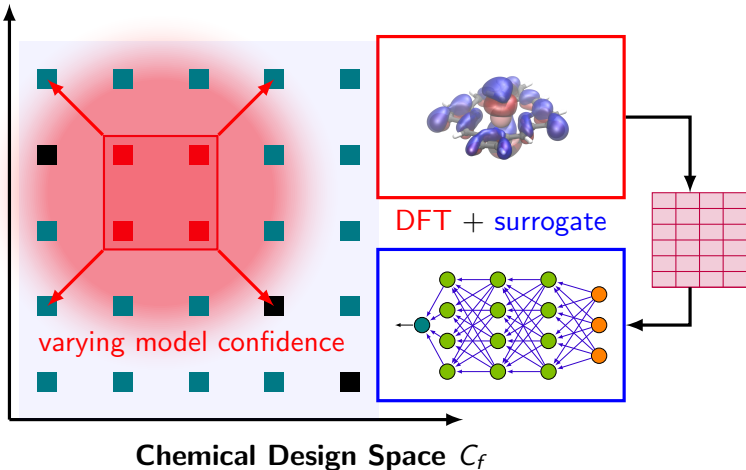
# Motivation: chemical discovery

How can we design new materials using computers?



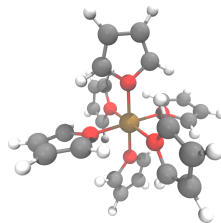
# Motivation: chemical discovery

How can we design new materials using computers?

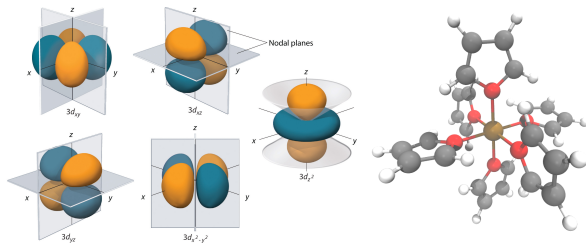


# Transition metal complexes

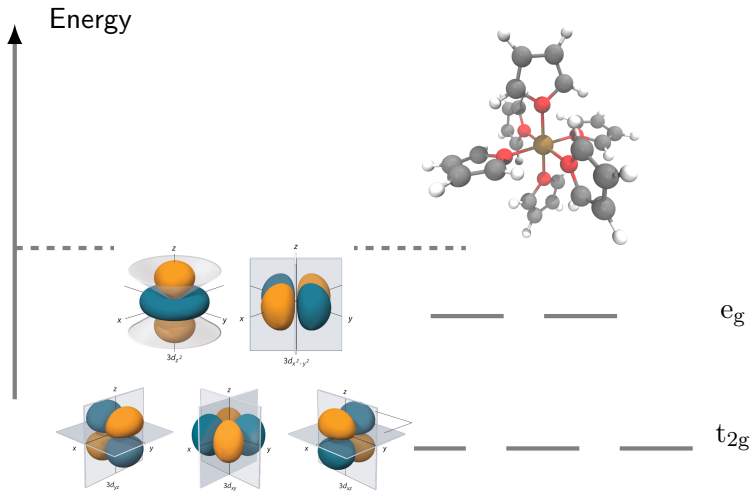
# Transition metal complexes



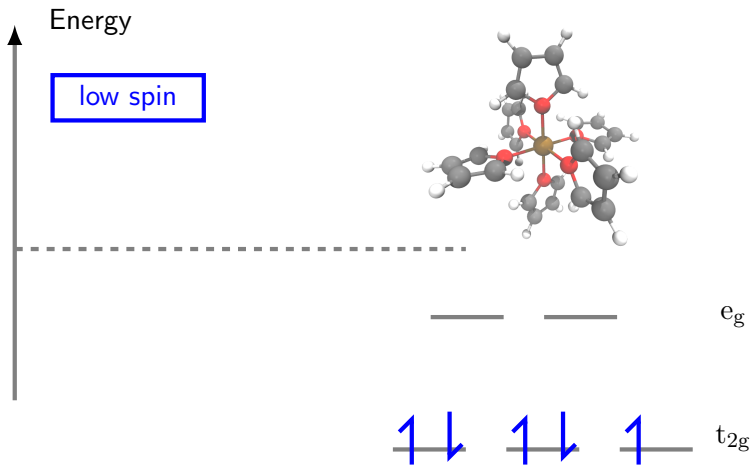
# Transition metal complexes



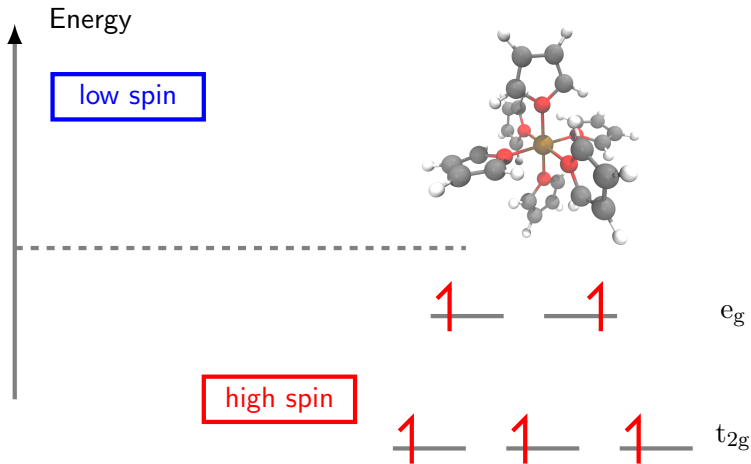
# Transition metal complexes



# Transition metal complexes

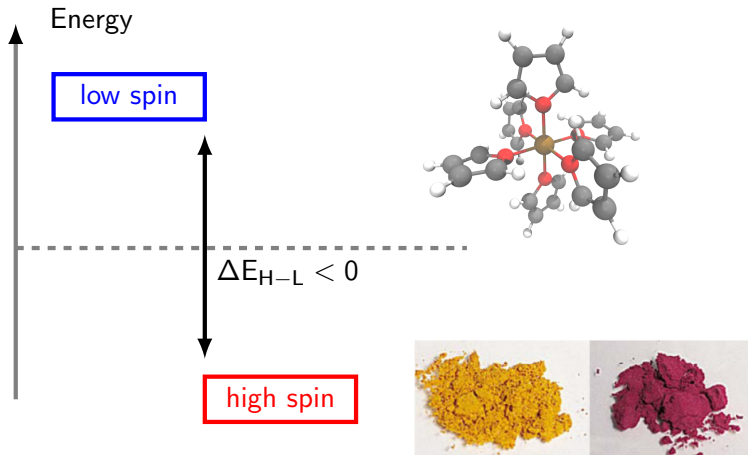


# Transition metal complexes

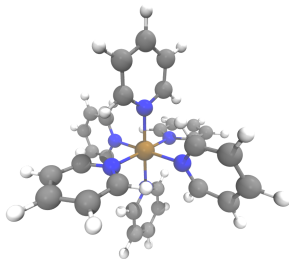
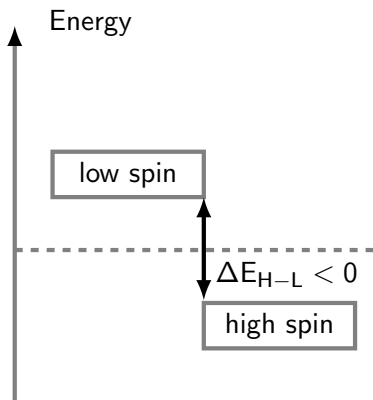




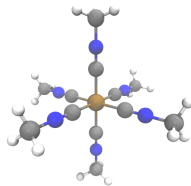
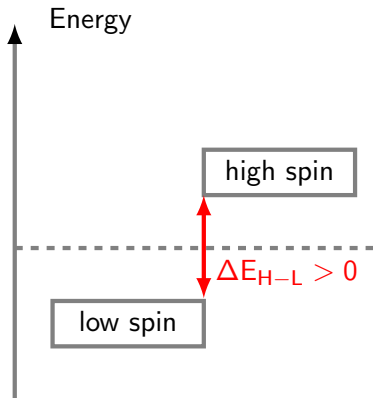
# Transition metal complexes



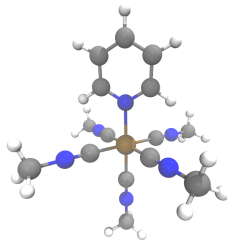
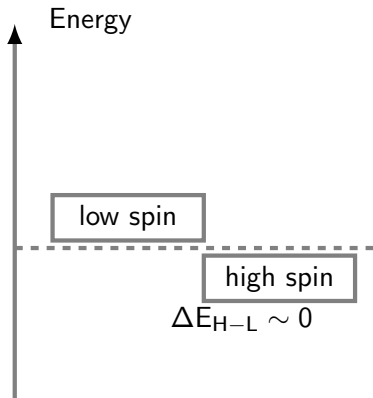
# Transition metal complexes



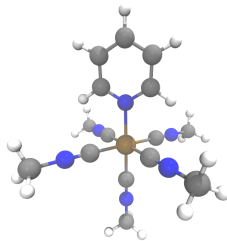
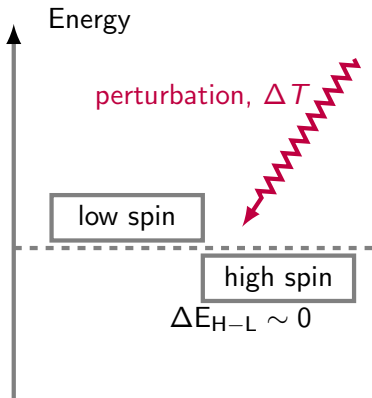
# Transition metal complexes



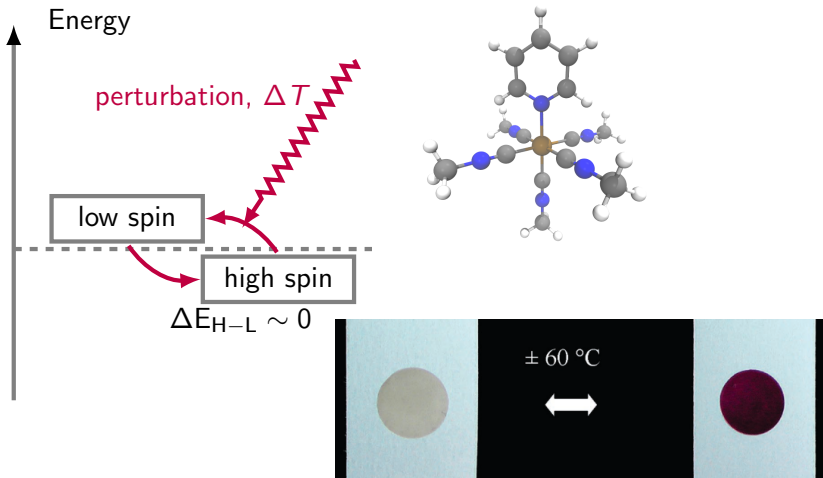
# Transition metal complexes



# Transition metal complexes

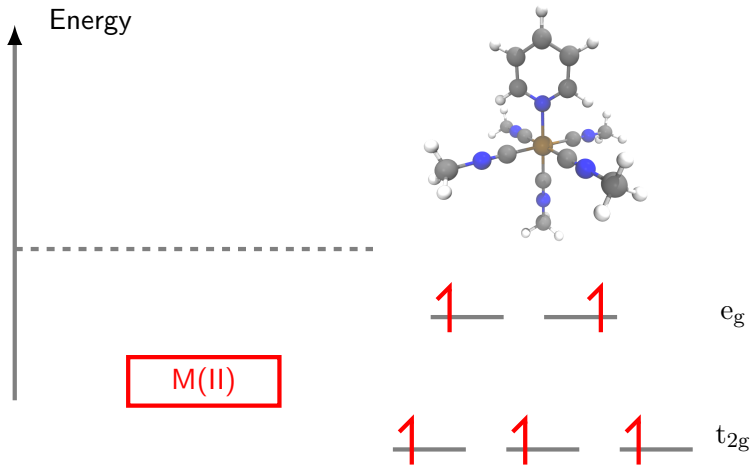


# Transition metal complexes

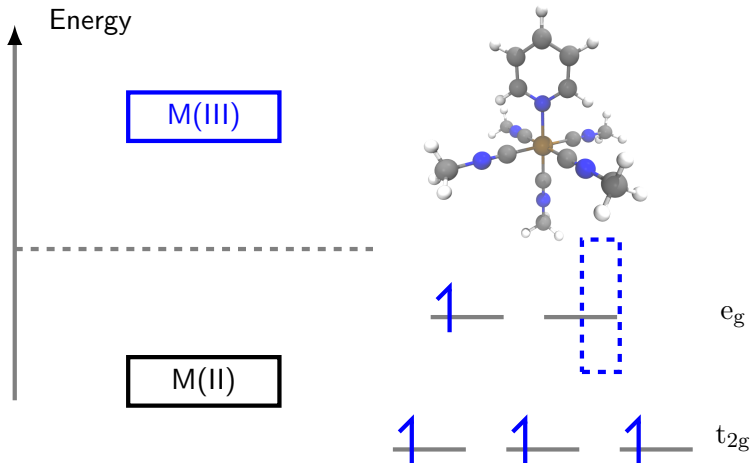


Seredyuk, M et al., *Chem. Mater.*, 18(10):2513–2519, 2006.

# Transition metal complexes

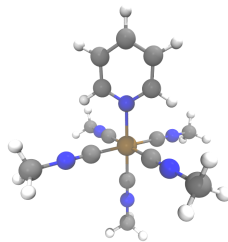
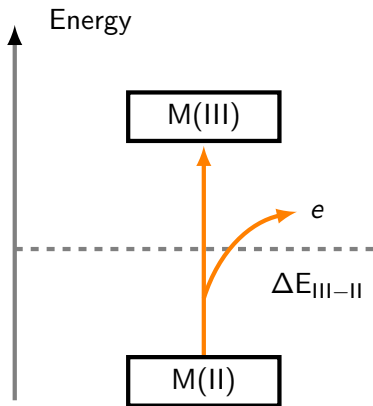


# Transition metal complexes





# Transition metal complexes

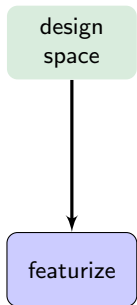


# Algorithmic chemical discovery

design  
space

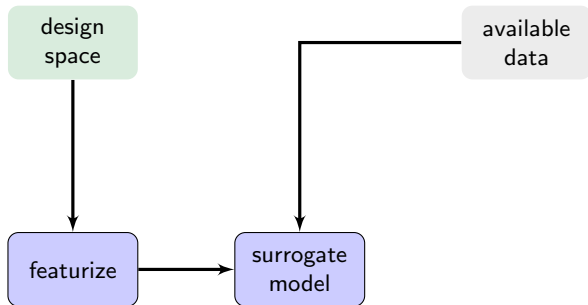
candidate  
materials

# Algorithmic chemical discovery



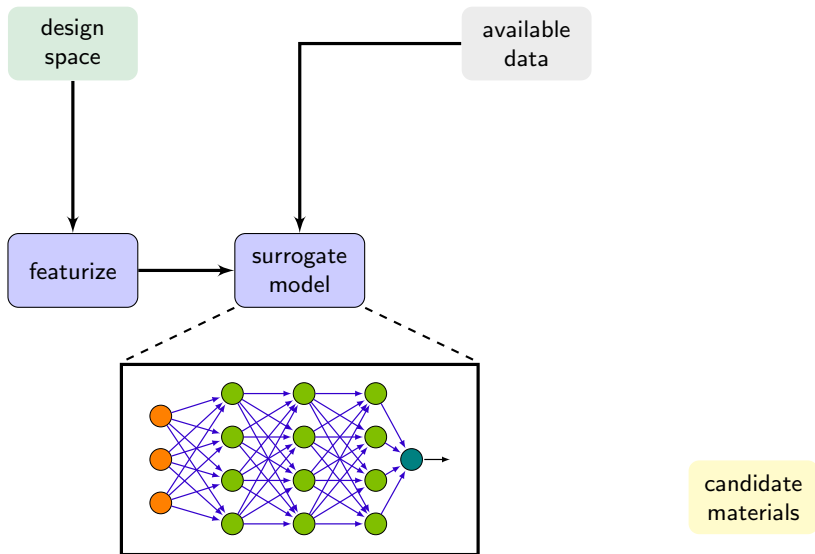
candidate materials

# Algorithmic chemical discovery

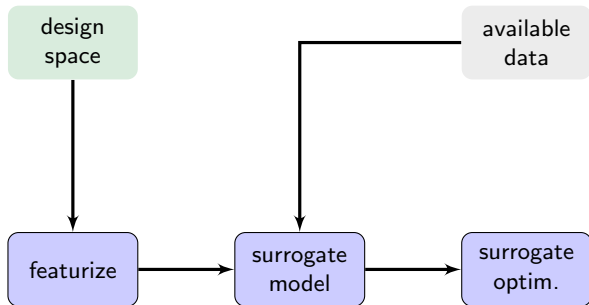


candidate materials

# Algorithmic chemical discovery

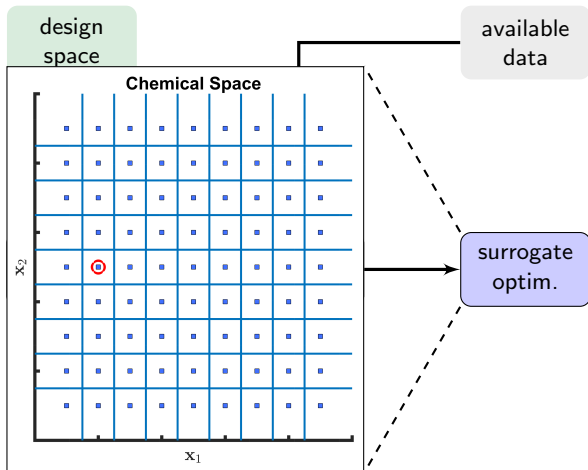


# Algorithmic chemical discovery



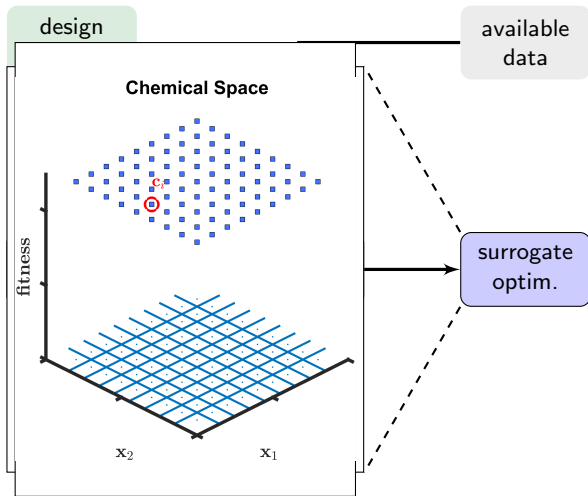
candidate  
materials

# Algorithmic chemical discovery



candidate materials

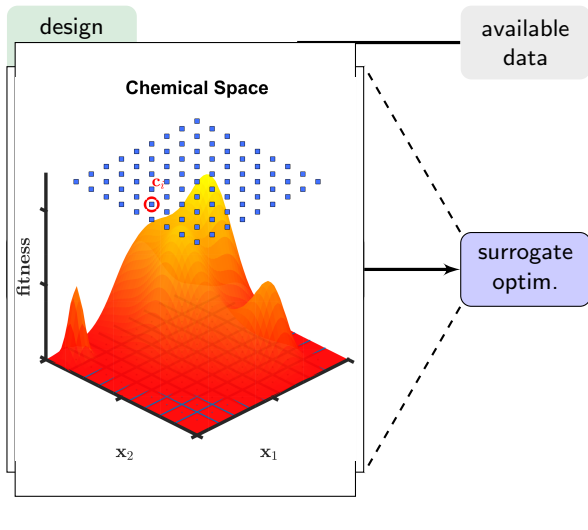
# Algorithmic chemical discovery



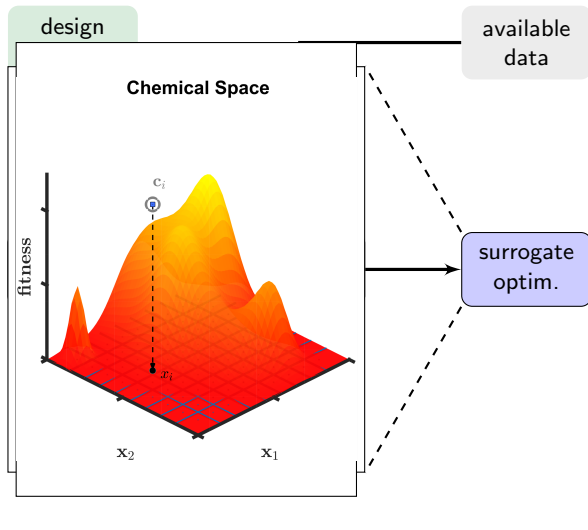
candidate  
materials



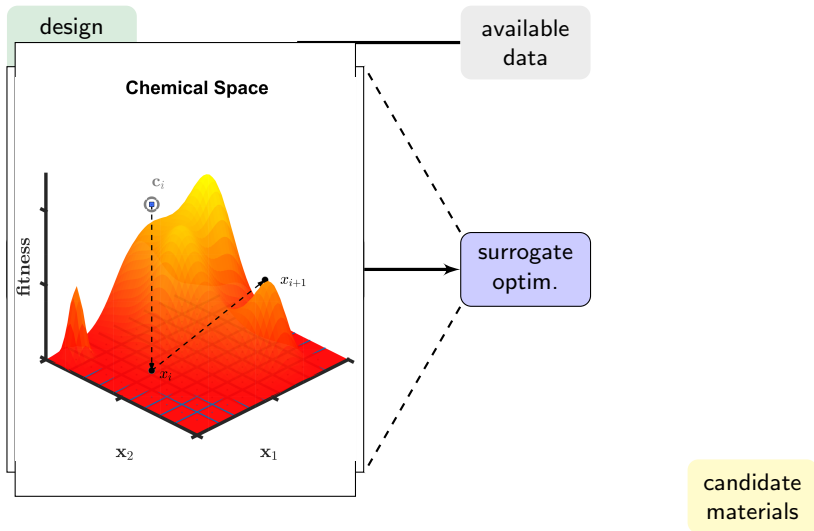
# Algorithmic chemical discovery



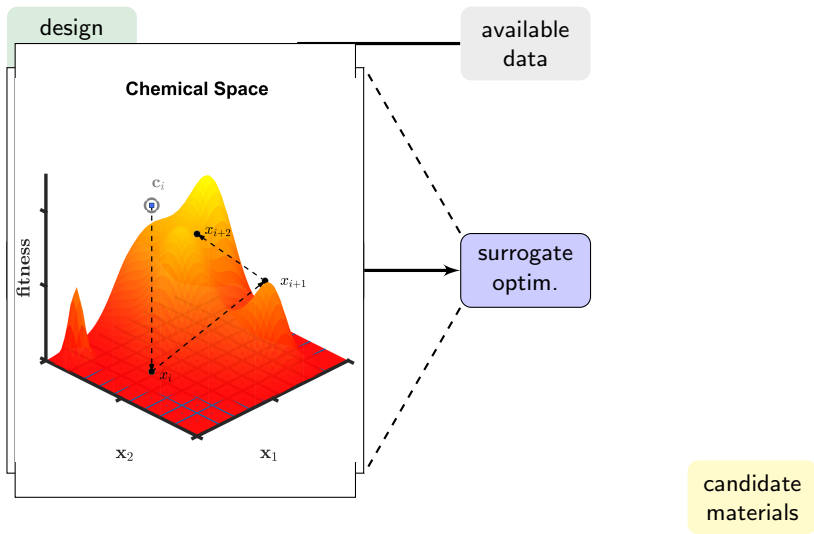
# Algorithmic chemical discovery



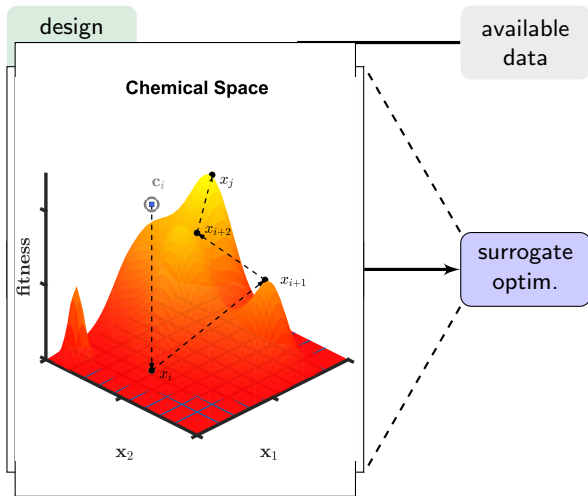
# Algorithmic chemical discovery



# Algorithmic chemical discovery

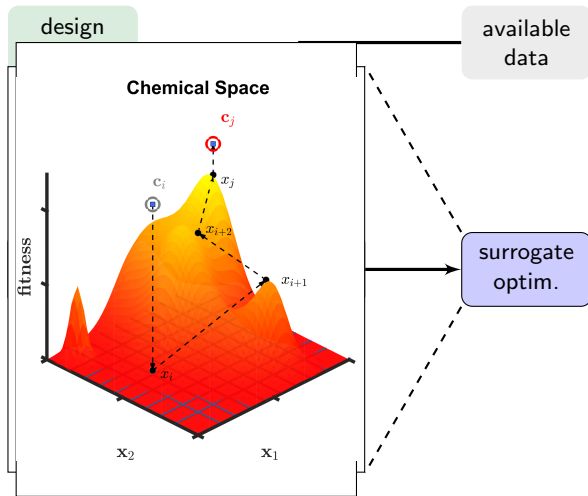


# Algorithmic chemical discovery



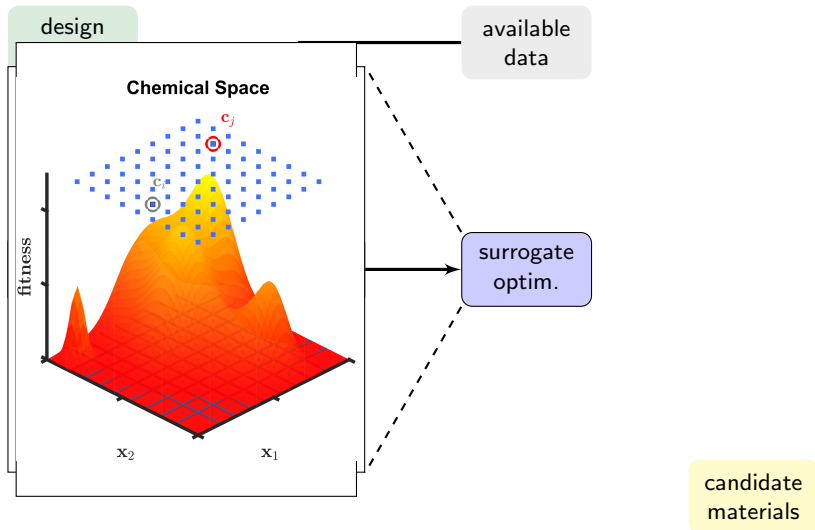
candidate materials

# Algorithmic chemical discovery

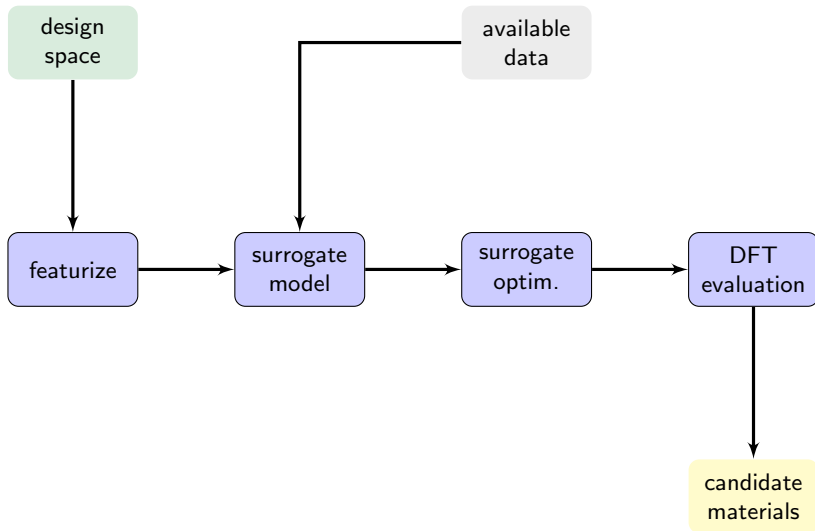


candidate  
materials

# Algorithmic chemical discovery

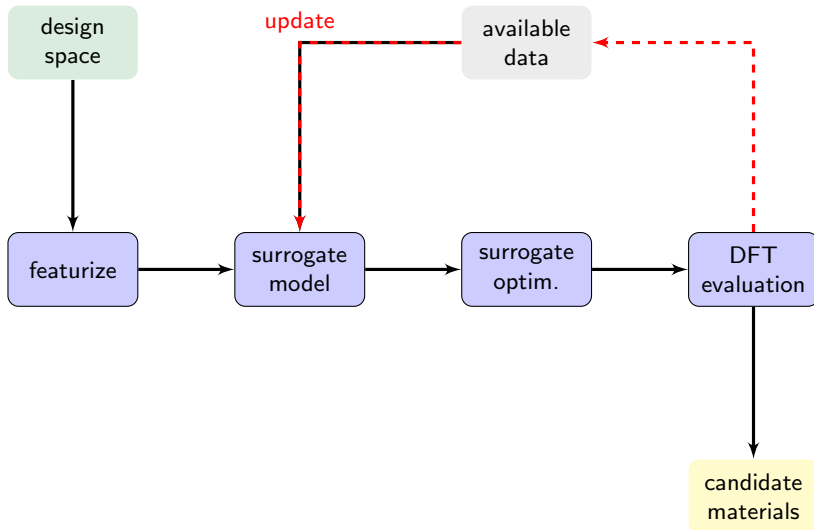


# Algorithmic chemical discovery





# Algorithmic chemical discovery

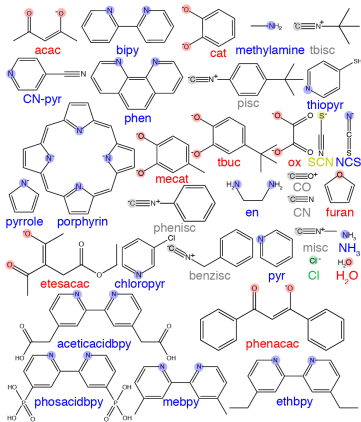


# First-principles calculations

train on  $\sim 100$ – $2000$  DFT calculations:

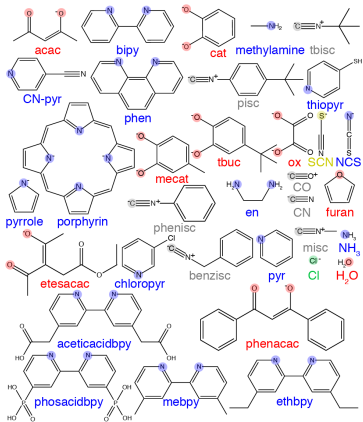
# First-principles calculations

train on  $\sim 100$ – $2000$  DFT calculations:



# First-principles calculations

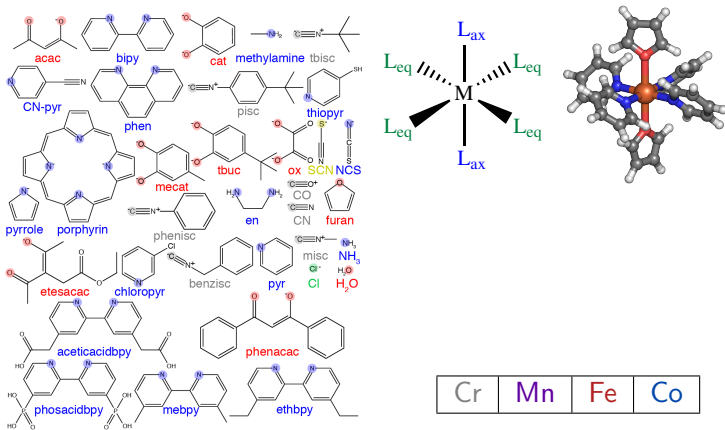
train on  $\sim 100\text{--}2000$  DFT calculations:



Cr	Mn	Fe	Co
----	----	----	----

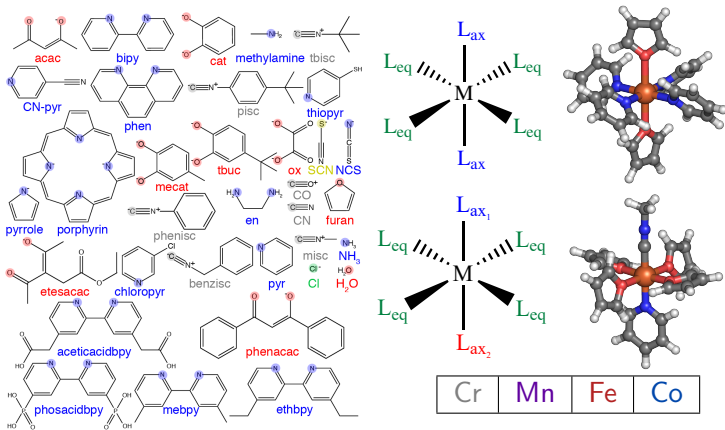
# First-principles calculations

train on ~ 100–2000 DFT calculations:



# First-principles calculations

train on  $\sim 100\text{--}2000$  DFT calculations:

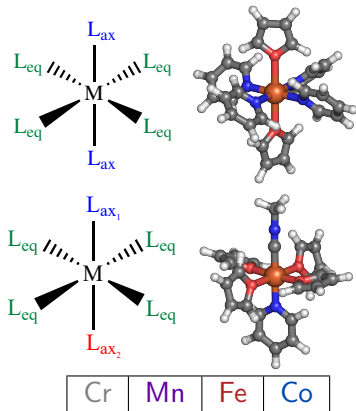


# First-principles calculations

train on  $\sim 100$ – $2000$  DFT calculations:

Details:

- B3LYP-like DFT
- gas phase optimization
- LANL2DZ/6-31G\*
- COSMO solvents
- high- and low-spin  
M(II)/(III)



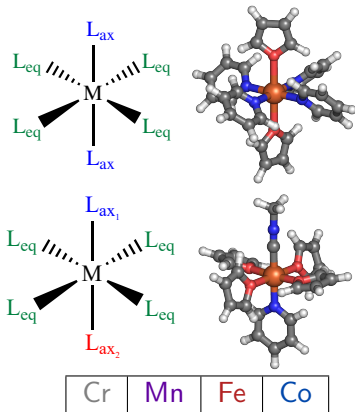
# First-principles calculations

train on  $\sim 100$ – $2000$  DFT calculations:

Details:

- B3LYP-like DFT
- gas phase optimization
- LANL2DZ/6-31G\*
- COSMO solvents
- high- and low-spin  
M(II)/(III)

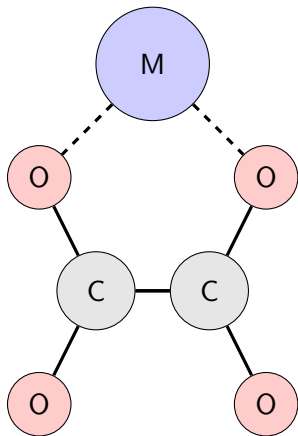
↘ HF exchange varied 0–30%





# Featurization with RACs

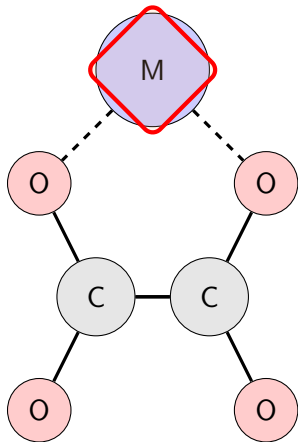
Graph-based features designed for TM complexes



Janet, J.P., and Kulik, H.J., *J. Phys. Chem. A*, 121(46):8939–8954, 2017.

# Featurization with RACs

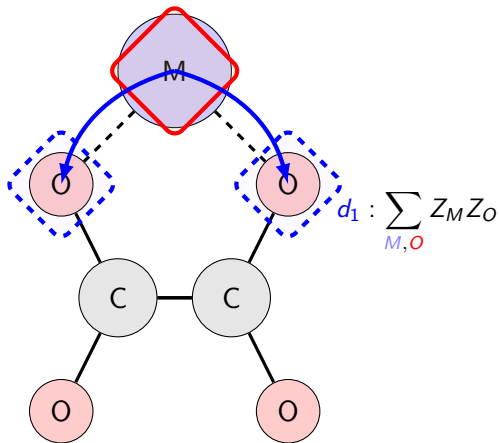
Graph-based features designed for TM complexes



Janet, J.P., and Kulik, H.J., *J. Phys. Chem. A*, 121(46):8939–8954, 2017.

# Featurization with RACs

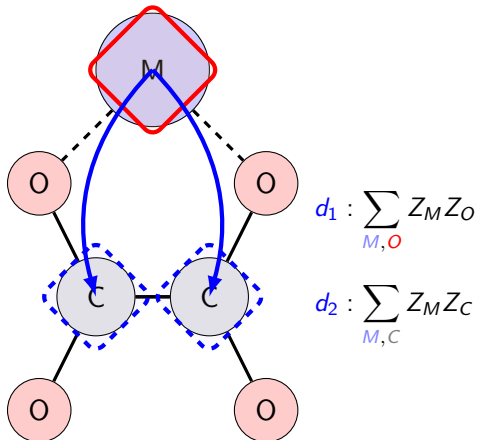
Graph-based features designed for TM complexes



Janet, J.P., and Kulik, H.J., *J. Phys. Chem. A*, 121(46):8939–8954, 2017.

# Featurization with RACs

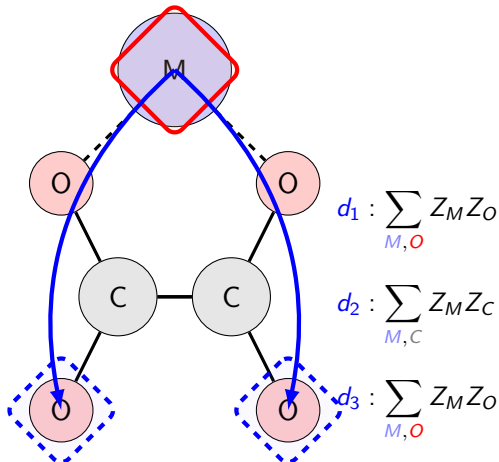
Graph-based features designed for TM complexes



Janet, J.P., and Kulik, H.J., *J. Phys. Chem. A*, 121(46):8939–8954, 2017.

# Featurization with RACs

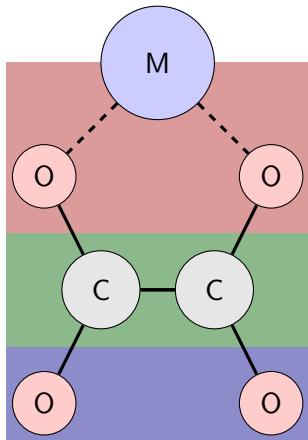
Graph-based features designed for TM complexes



Janet, J.P., and Kulik, H.J., *J. Phys. Chem. A*, 121(46):8939–8954, 2017.

# Featurization with RACs

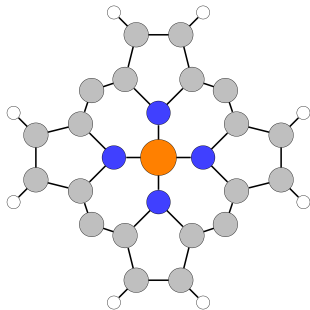
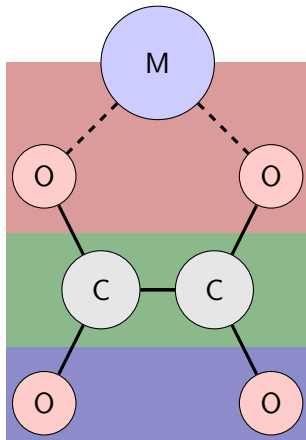
Graph-based features designed for TM complexes



Janet, J.P., and Kulik, H.J., *J. Phys. Chem. A*, 121(46):8939–8954, 2017.

# Featurization with RACs

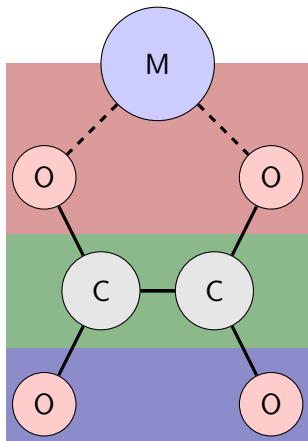
Graph-based features designed for TM complexes



Janet, J.P., and Kulik, H.J., *J. Phys. Chem. A*, 121(46):8939–8954, 2017.

# Featurization with RACs

Graph-based features designed for TM complexes



spin splitting  $\Delta E_{H-L}$

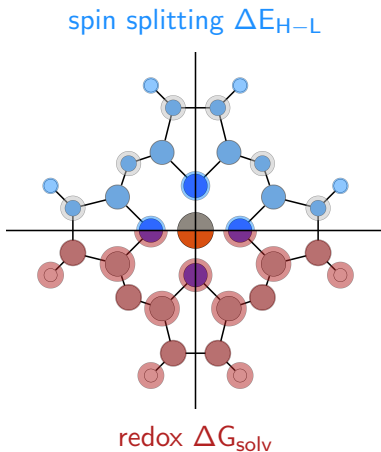
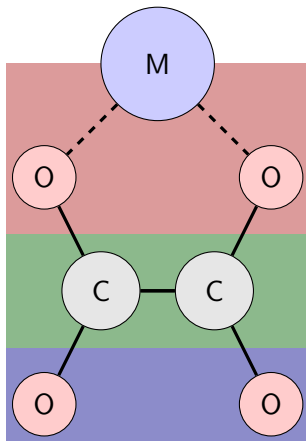
redox  $\Delta G_{\text{solv}}$

Janet, J.P., and Kulik, H.J., *J. Phys. Chem. A*, 121(46):8939–8954, 2017.



# Featurization with RACs

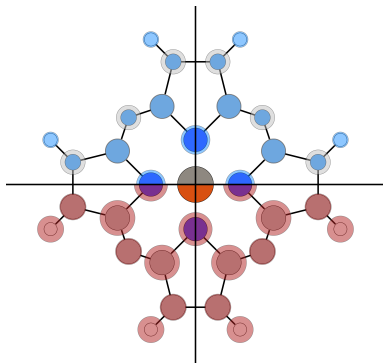
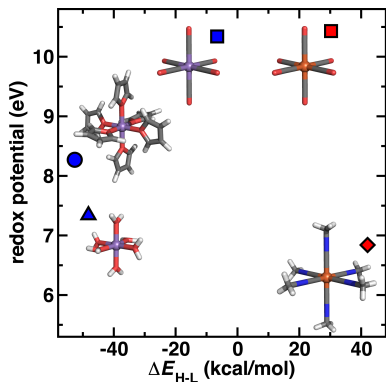
Graph-based features designed for TM complexes



Janet, J.P., and Kulik, H.J., *J. Phys. Chem. A*, 121(46):8939–8954, 2017.

# Featurization with RACs

Graph-based features designed for TM complexes



Janet, J.P., and Kulik, H.J., *J. Phys. Chem. A*, 121(46):8939–8954, 2017.

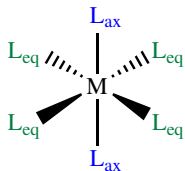
Janet, J.P. et al., *Inorg. Chem.*, 58(16):10592–10606, 2019.

# Property inference with ANNs

Estimate properties using small artificial neural networks (ANNs)

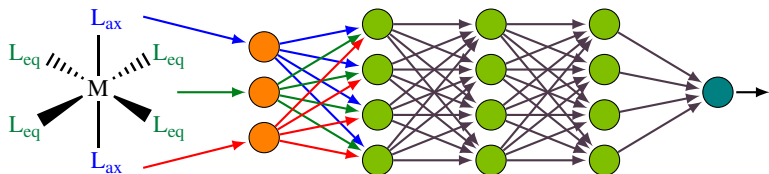
# Property inference with ANNs

Estimate properties using small artificial neural networks (ANNs)



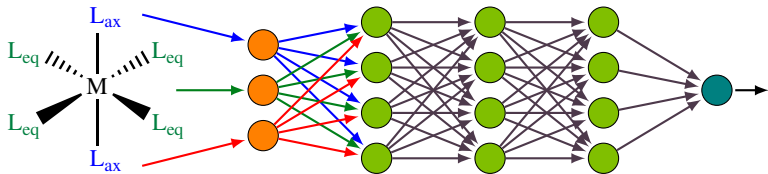
# Property inference with ANNs

Estimate properties using small artificial neural networks (ANNs)

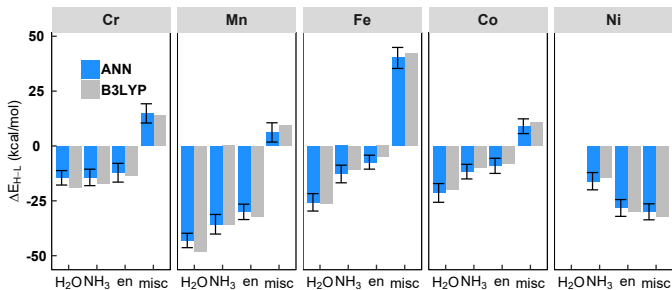


# Property inference with ANNs

Estimate properties using small artificial neural networks (ANNs)

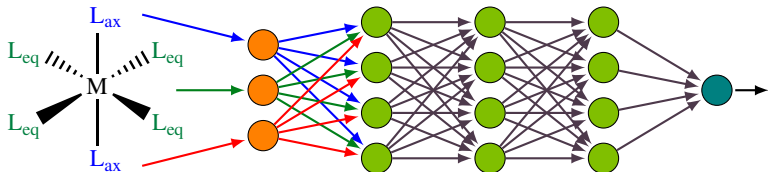


spin splitting energies

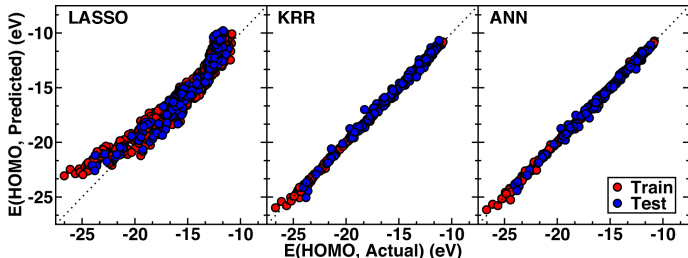


# Property inference with ANNs

Estimate properties using small artificial neural networks (ANNs)



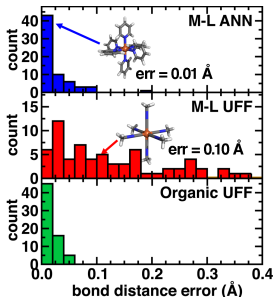
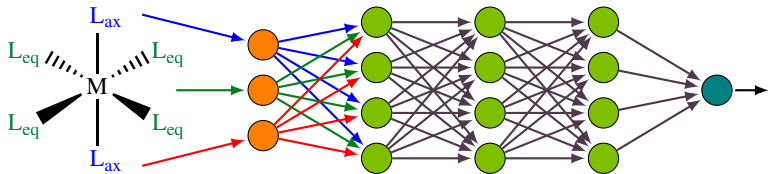
frontier orbital properties



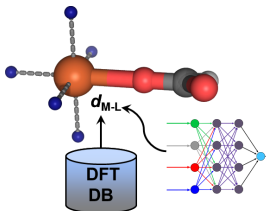
ANN: 500 × 500 ReLU nodes, fully connected

# Property inference with ANNs

Estimate properties using small artificial neural networks (ANNs)



DFT equilibrium bond lengths



Janet, J.P. et al., *Ind. Eng. Chem. Res.*, 56(17):4898–4910, 2017.

Janet, J.P. et al., *Inorg. Chem.*, 58(16):10592–10606, 2019. 300 × 200 × 200 tanh nodes, fully connected



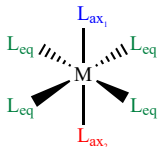
# Beyond prediction: live job management

In high-throughput DFT screening, job failure is a frequent issue:

# Beyond prediction: live job management

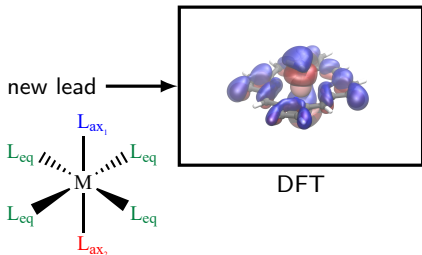
In high-throughput DFT screening, job failure is a frequent issue:

new lead



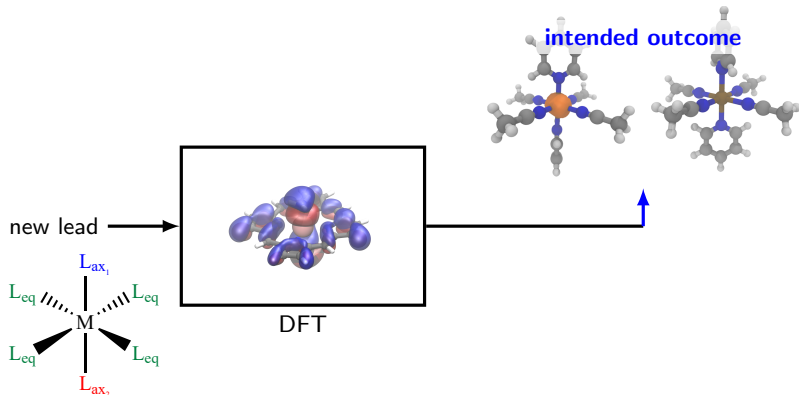
# Beyond prediction: live job management

In high-throughput DFT screening, job failure is a frequent issue:



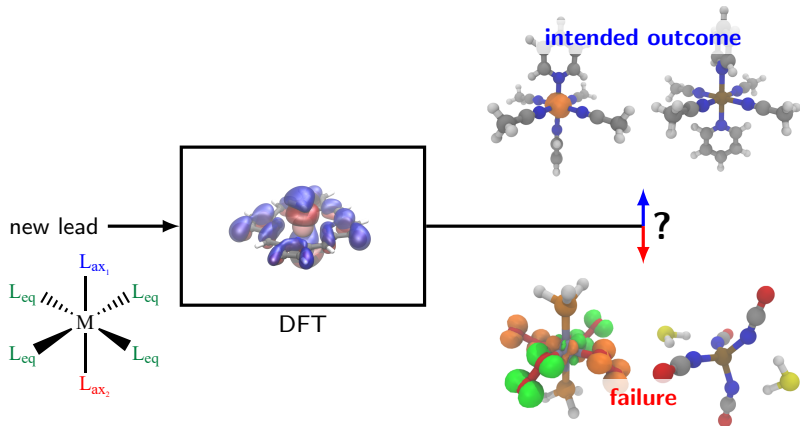
# Beyond prediction: live job management

In high-throughput DFT screening, job failure is a frequent issue:

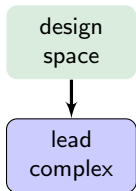


# Beyond prediction: live job management

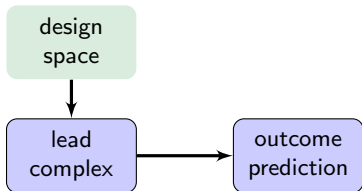
In high-throughput DFT screening, job failure is a frequent issue:



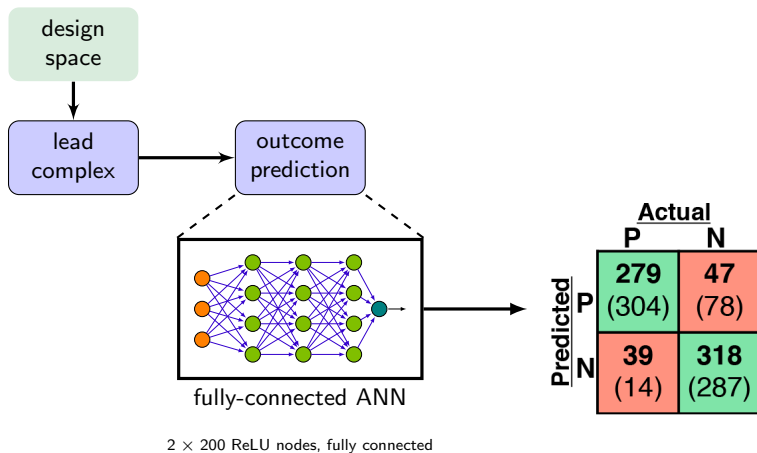
# Beyond prediction: live job management



# Beyond prediction: live job management

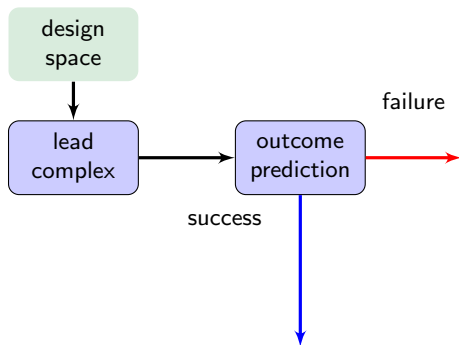


# Beyond prediction: live job management

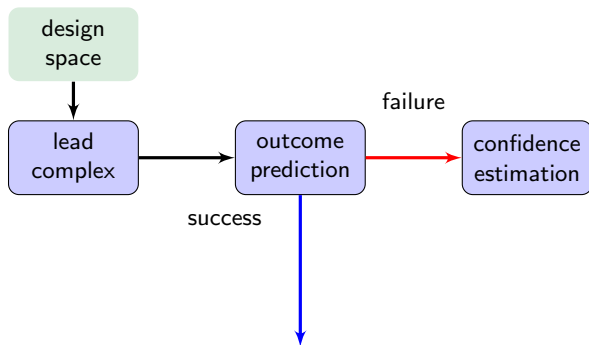




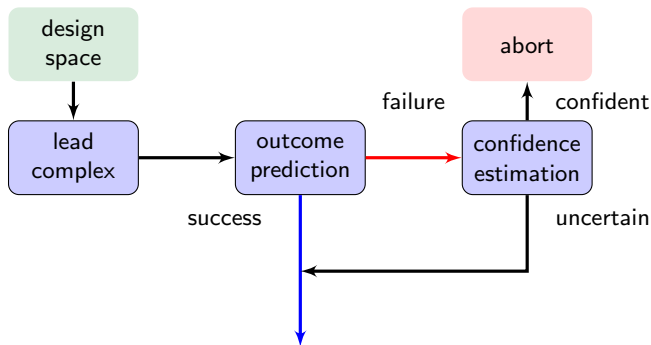
# Beyond prediction: live job management



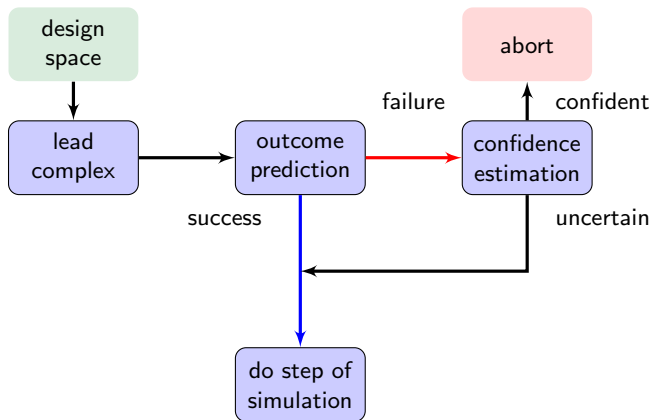
# Beyond prediction: live job management



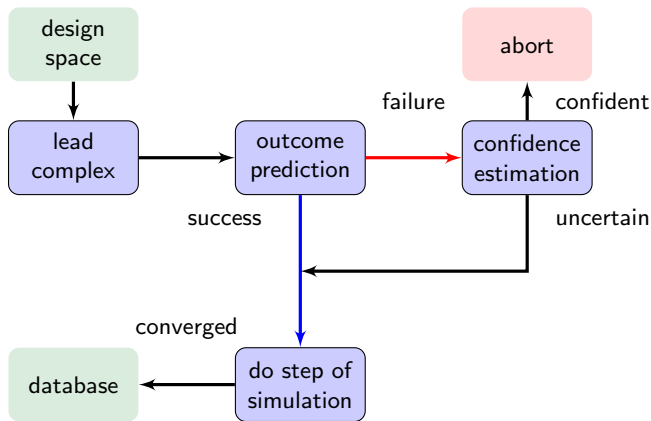
# Beyond prediction: live job management



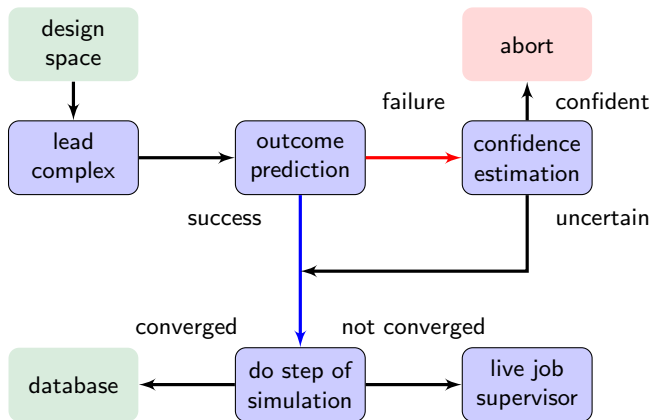
# Beyond prediction: live job management



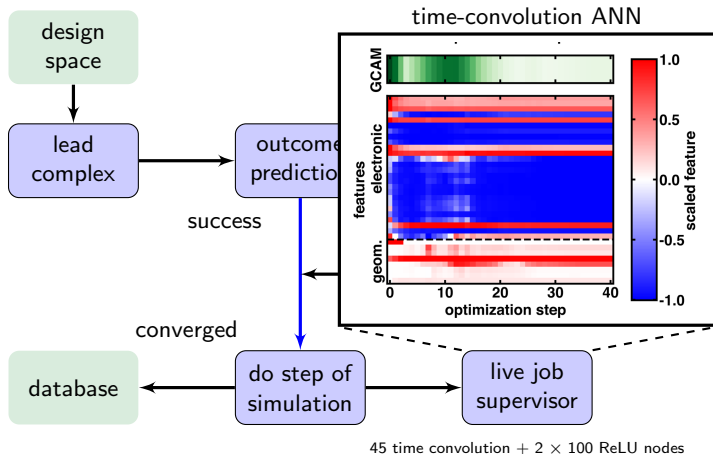
# Beyond prediction: live job management



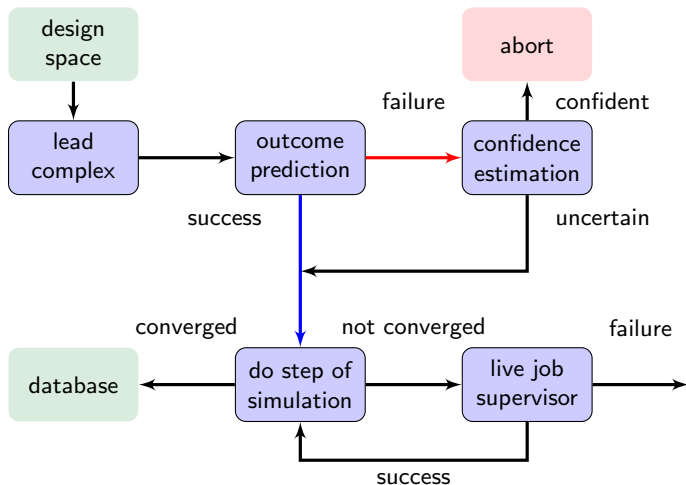
# Beyond prediction: live job management



# Beyond prediction: live job management

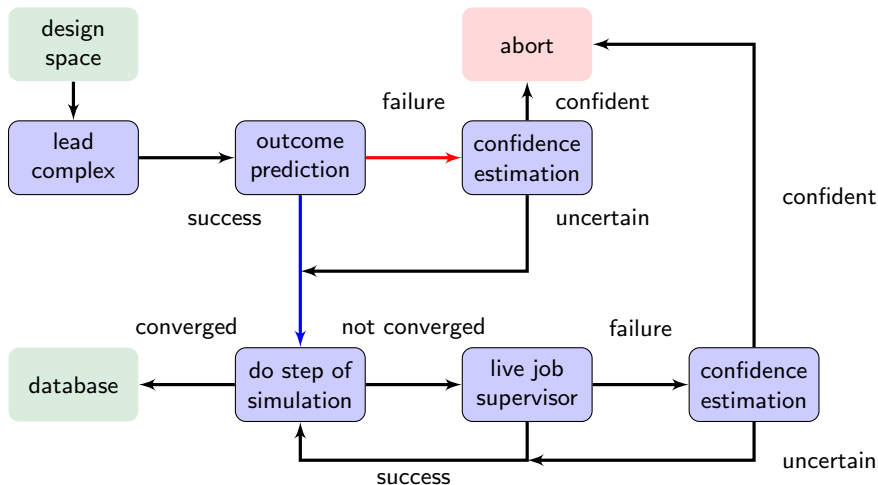


# Beyond prediction: live job management

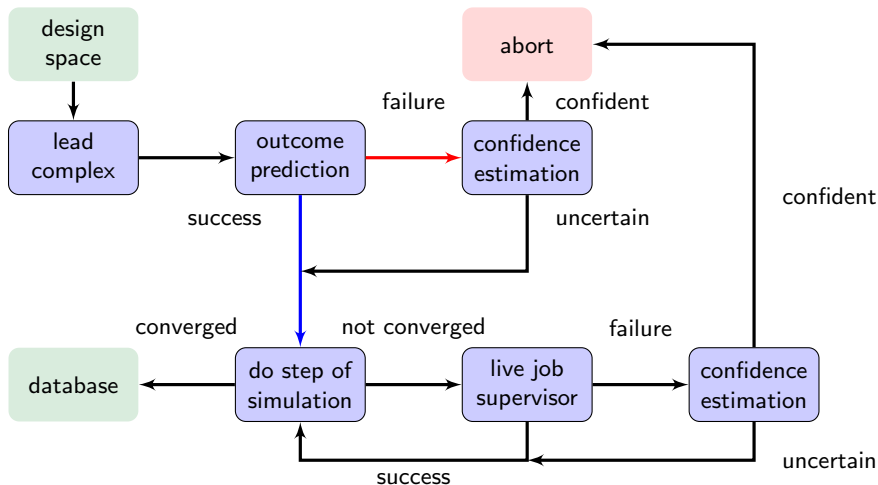




# Beyond prediction: live job management



# Beyond prediction: live job management



This leads to about **40% time savings** and can abort almost all failures.

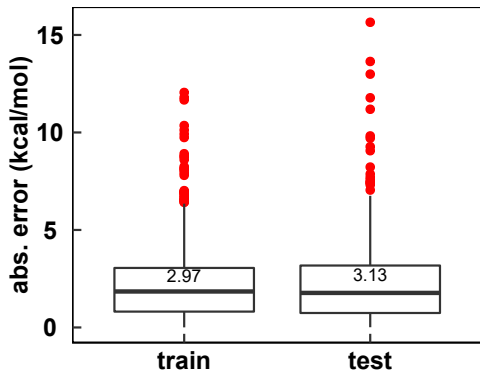
Duan, C., Janet, J.P. et al., *J. Chem. Theory. Comp.*, 15(4):2331–2345, 2019.

# Model transferability

Test-set performance is not necessarily a good metric for general transferability:

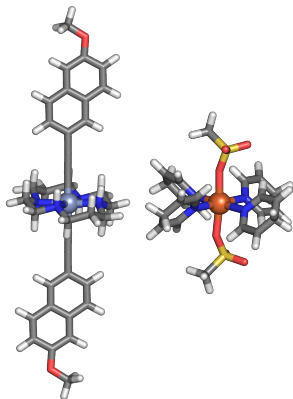
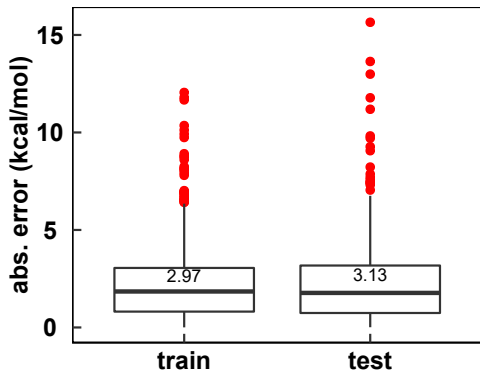
# Model transferability

Test-set performance is not necessarily a good metric for general transferability:



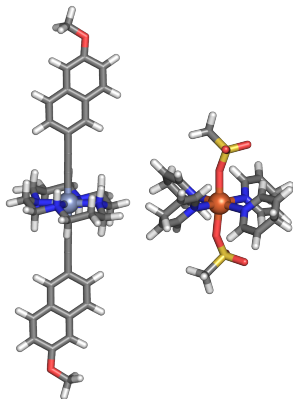
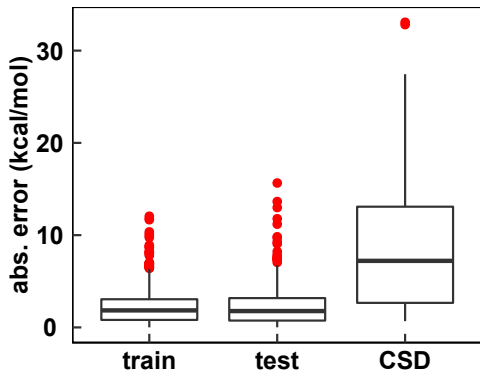
# Model transferability

Test-set performance is not necessarily a good metric for general transferability:



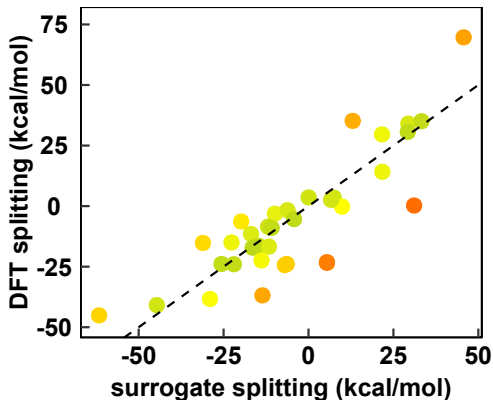
# Model transferability

Test-set performance is not necessarily a good metric for general transferability:

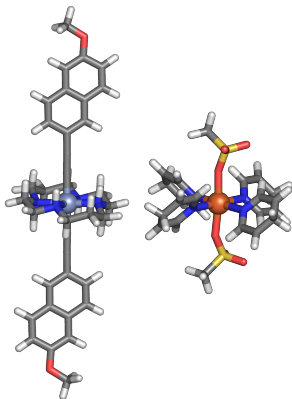


# Model transferability

Test-set performance is not necessarily a good metric for general transferability:



Janet, J.P., and Kulik, H.J., *Chem. Sci.*, 8:5137–5152, 2017.



50 × 50 tanh nodes fully connected ANN

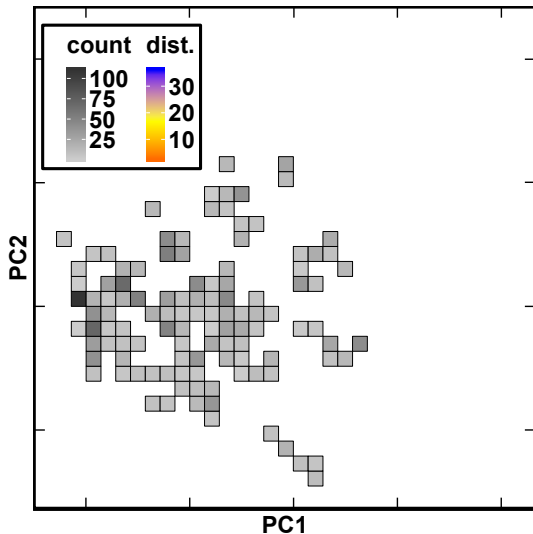
# How far can we extrapolate?

'Out-of-distribution' test:  
spin-splitting energies of  
116 structures from the  
CSD, from training-like to  
very different.



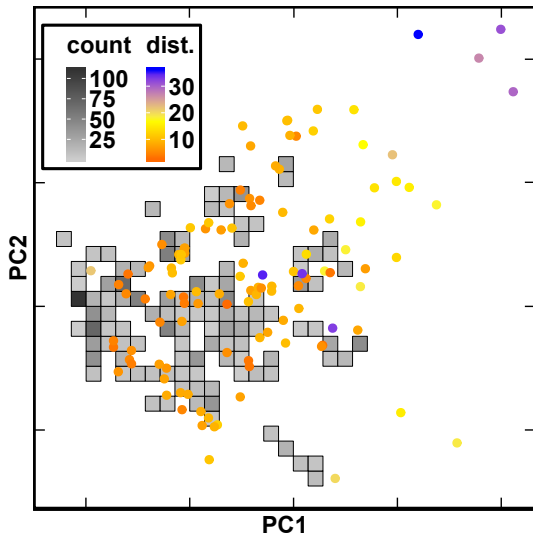
## How far can we extrapolate?

'Out-of-distribution' test:  
spin-splitting energies of  
116 structures from the  
CSD, from training-like to  
very different.



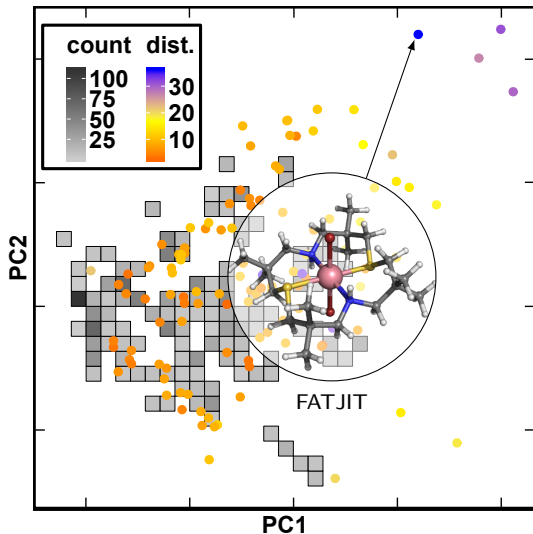
## How far can we extrapolate?

'Out-of-distribution' test:  
spin-splitting energies of  
116 structures from the  
CSD, from training-like to  
very different.



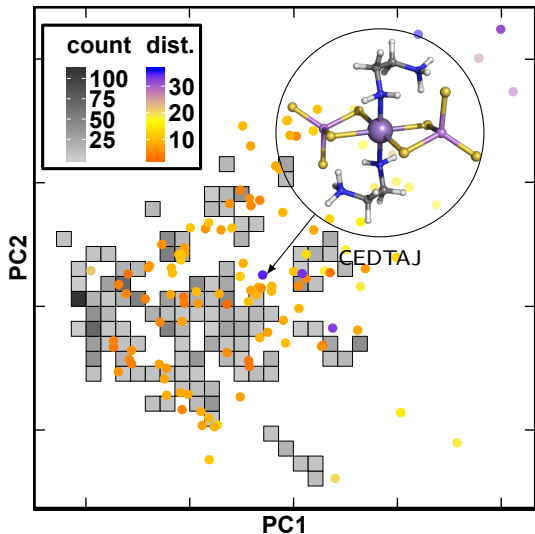
## How far can we extrapolate?

'Out-of-distribution' test:  
spin-splitting energies of  
116 structures from the  
CSD, from training-like to  
very different.



## How far can we extrapolate?

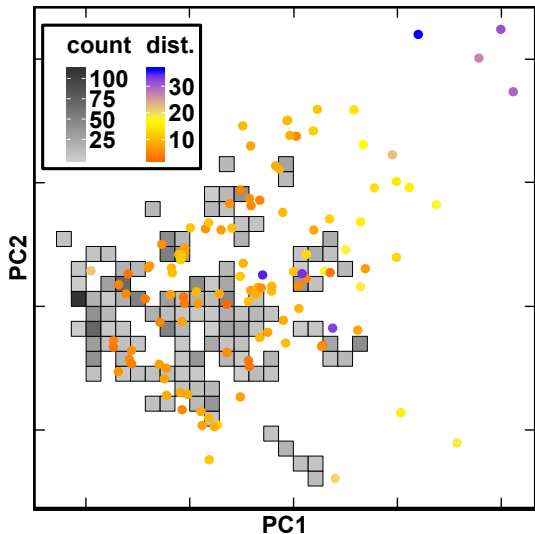
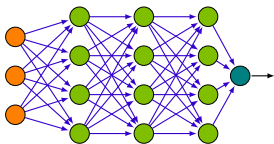
'Out-of-distribution' test:  
spin-splitting energies of  
116 structures from the  
CSD, from training-like to  
very different.



## How far can we extrapolate?

'Out-of-distribution' test:  
spin-splitting energies of  
116 structures from the  
CSD, from training-like to  
very different.

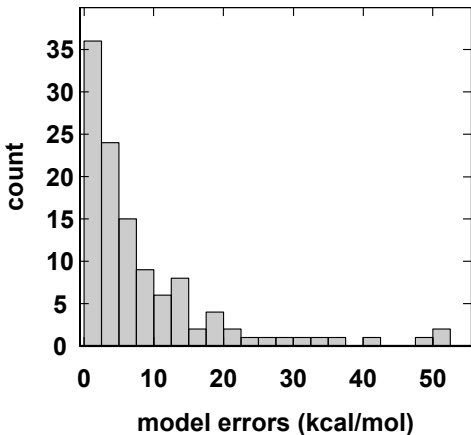
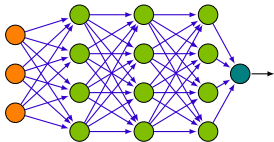
Train 3-layer fully con-  
nected ANN on 1900 DFT  
results on simple ligands:



# How far can we extrapolate?

'Out-of-distribution' test:  
spin-splitting energies of  
116 structures from the  
CSD, from training-like to  
very different.

Train 3-layer fully con-  
nected ANN on 1900 DFT  
results on simple ligands:



# Latent distance similarity

# Latent distance similarity

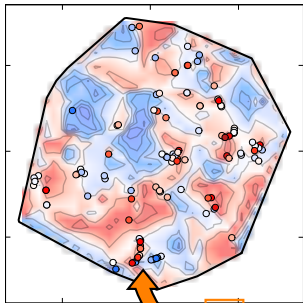
input molecule





# Latent distance similarity

feature space

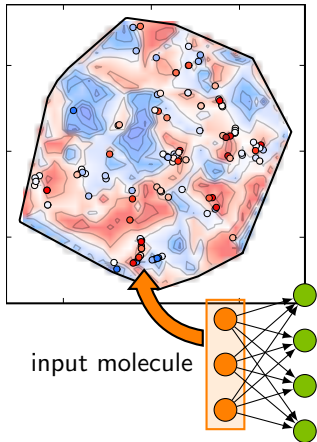


input molecule



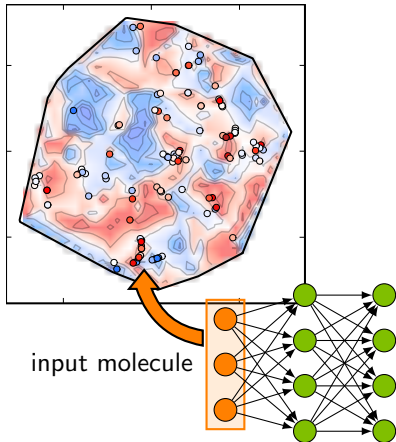
# Latent distance similarity

feature space



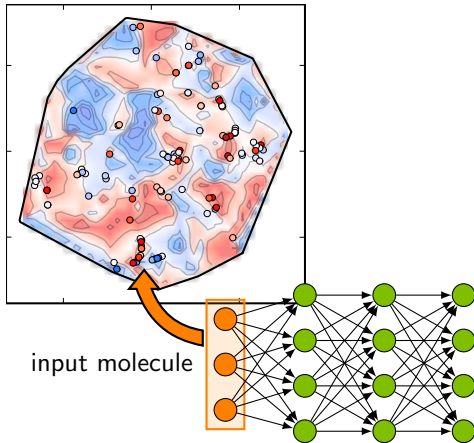
# Latent distance similarity

feature space



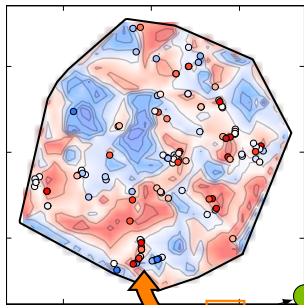
# Latent distance similarity

feature space

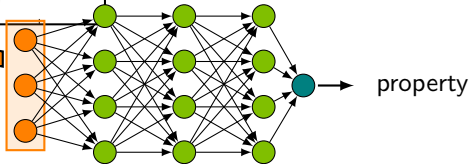


# Latent distance similarity

feature space

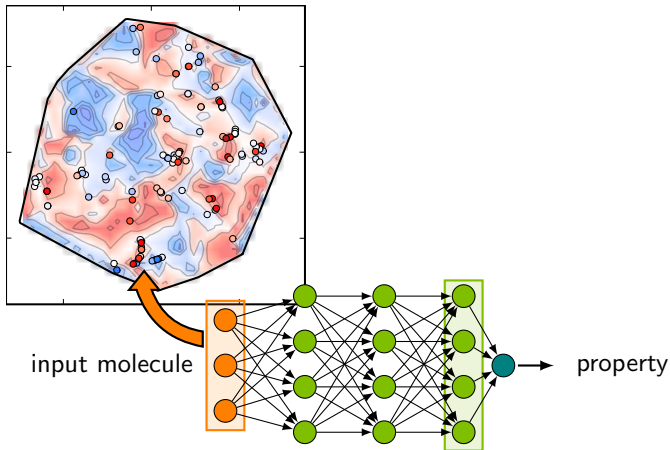


input molecule



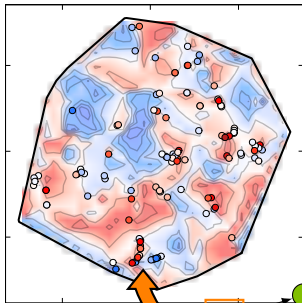
# Latent distance similarity

feature space

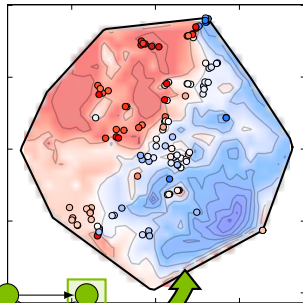


# Latent distance similarity

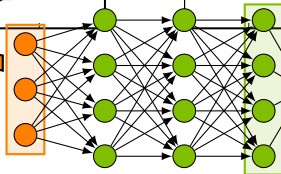
feature space



latent space



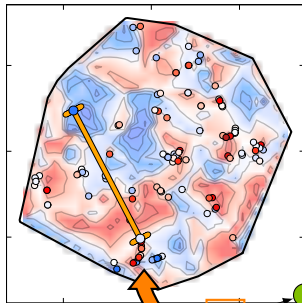
input molecule



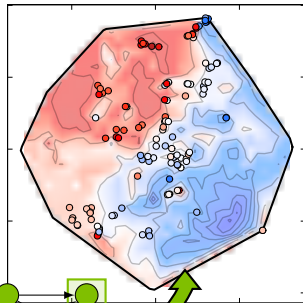
property

# Latent distance similarity

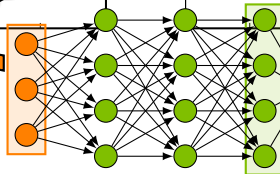
feature space geometry



latent space



input molecule

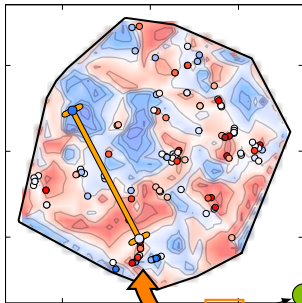


property

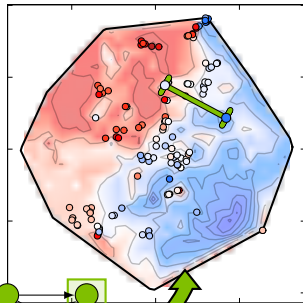


# Latent distance similarity

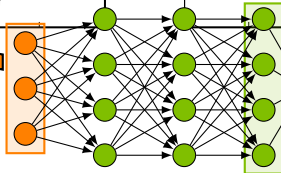
feature space geometry



latent space geometry



input molecule



property

## From distance to energy

Propose a simple conditionally-Gaussian model for predicting error distribution with latent distance,  $d$ :

$$\varepsilon(d) \sim \mathcal{N}(0, \sigma_1^2 + d\sigma_2^2)$$

## From distance to energy

Propose a simple conditionally-Gaussian model for predicting error distribution with latent distance,  $d$ :

$$\varepsilon(d) \sim \mathcal{N}(0, \sigma_1^2 + d\sigma_2^2)$$

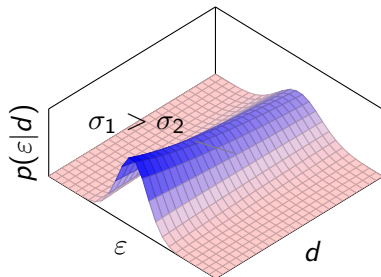
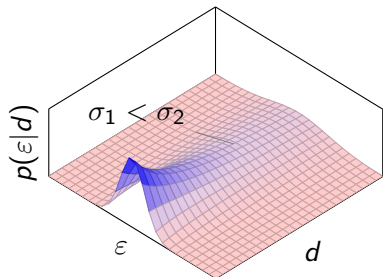
Estimate  $\sigma_1$ ,  $\sigma_2$  by max log likelihood, stably estimated even using few out-of-sample points

## From distance to energy

Propose a simple conditionally-Gaussian model for predicting error distribution with latent distance,  $d$ :

$$\varepsilon(d) \sim \mathcal{N}(0, \sigma_1^2 + d\sigma_2^2)$$

Estimate  $\sigma_1, \sigma_2$  by max log likelihood, stably estimated even using few out-of-sample points



## Other UQ metrics

1) Data-sampling ensembles:

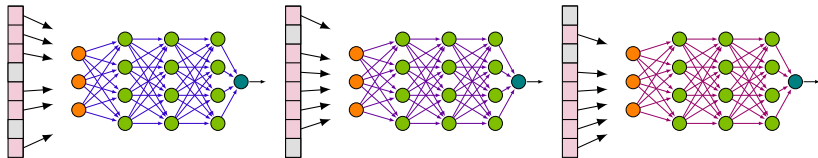
# Other UQ metrics

## 1) Data-sampling ensembles:



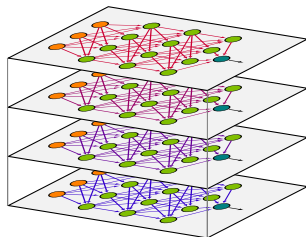
# Other UQ metrics

## 1) Data-sampling ensembles:



# Other UQ metrics

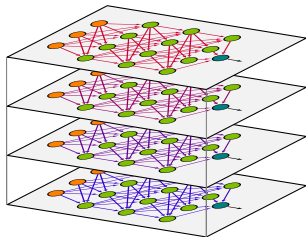
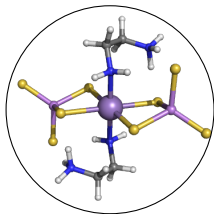
## 1) Data-sampling ensembles:





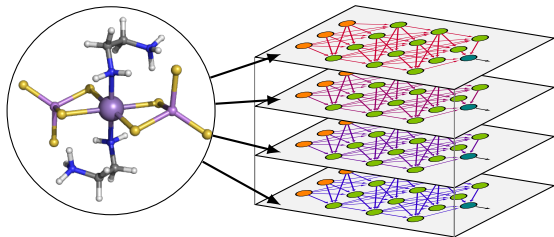
# Other UQ metrics

## 1) Data-sampling ensembles:



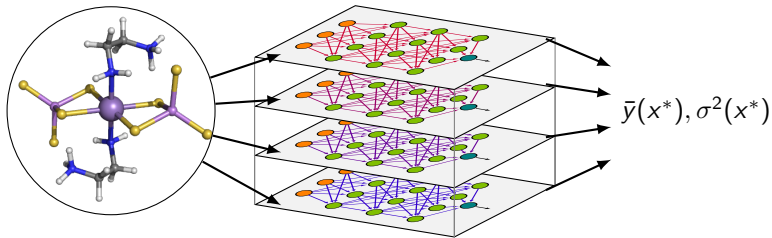
# Other UQ metrics

## 1) Data-sampling ensembles:



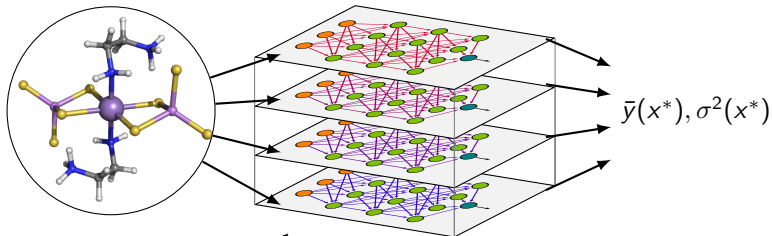
# Other UQ metrics

## 1) Data-sampling ensembles:

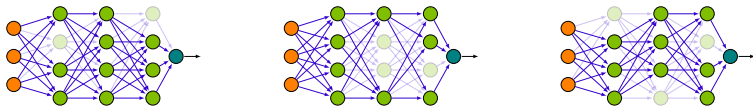


## Other UQ metrics

### 1) Data-sampling ensembles:



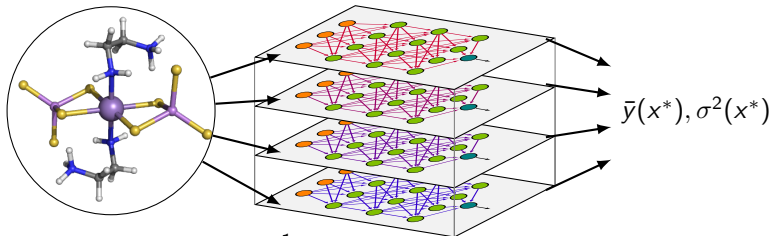
### 2) Monte Carlo dropout<sup>1</sup>:



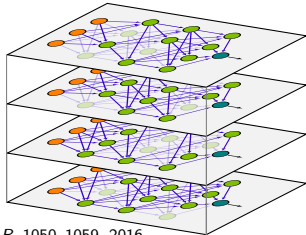
<sup>1</sup>:Gal, Y. and Ghahramani, Z., *ICMLR*, 1050–1059, 2016.

# Other UQ metrics

## 1) Data-sampling ensembles:



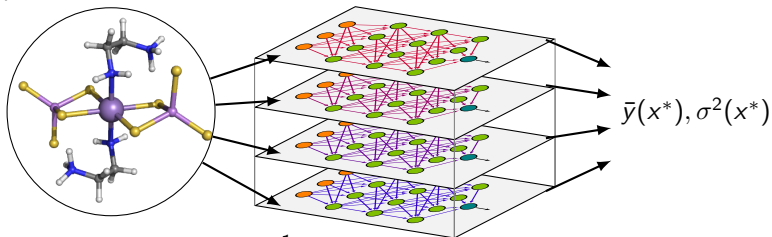
## 2) Monte Carlo dropout<sup>1</sup>:



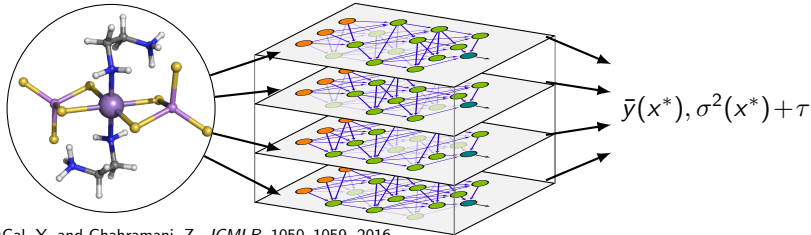
<sup>1</sup>:Gal, Y. and Ghahramani, Z., *ICMLR*, 1050–1059, 2016.

# Other UQ metrics

## 1) Data-sampling ensembles:

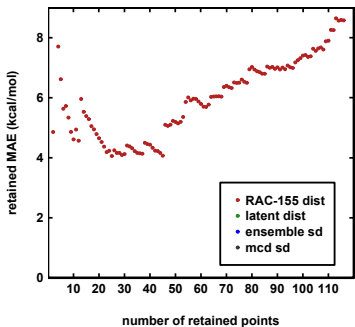
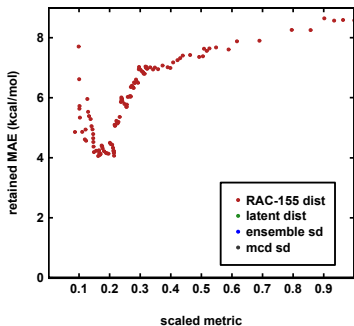


## 2) Monte Carlo dropout<sup>1</sup>:



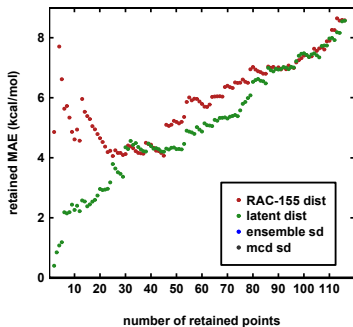
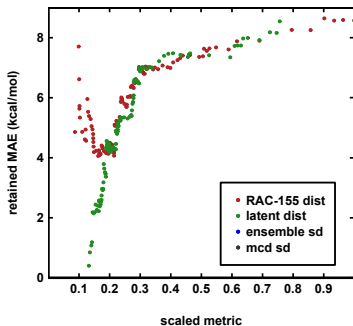
<sup>1</sup>:Gal, Y. and Ghahramani, Z., *ICMLR*, 1050–1059, 2016.

# How do these distributions compare?



# How do these distributions compare?

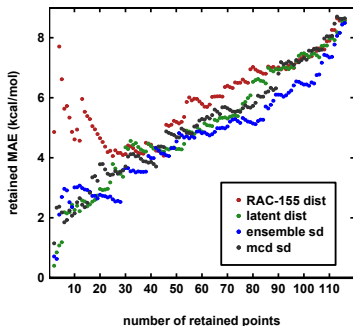
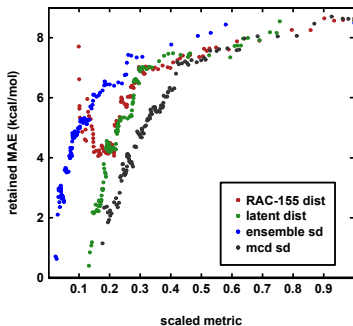
latent distances are superior to feature space distances





# How do these distributions compare?

latent distances are superior to feature space distances



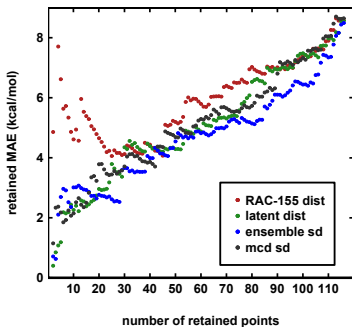
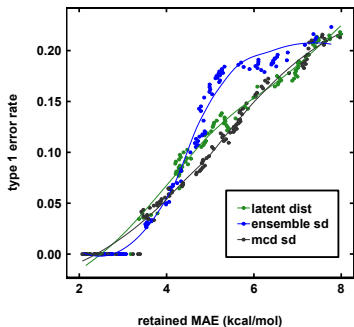
comparable with ensembles and mc dropout

# How do these distributions compare?

Comparison in energy units:

$$\varepsilon(d) \sim \mathcal{N}(0, \sigma_1^2 + d\sigma_2^2)$$

latent distances are superior to feature space distances



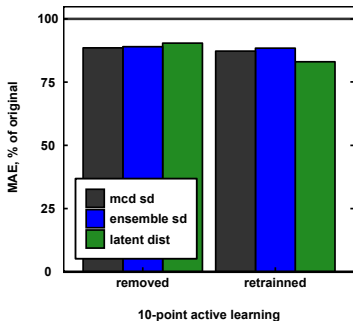
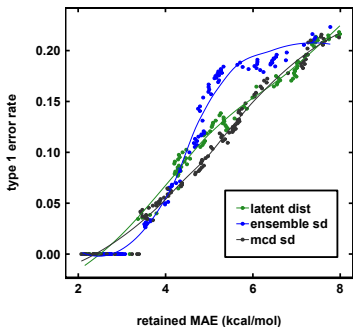
comparable with ensembles and mc dropout

# How do these distributions compare?

Comparison in energy units:

$$\varepsilon(d) \sim \mathcal{N}(0, \sigma_1^2 + d\sigma_2^2)$$

latent distances are superior to feature space distances

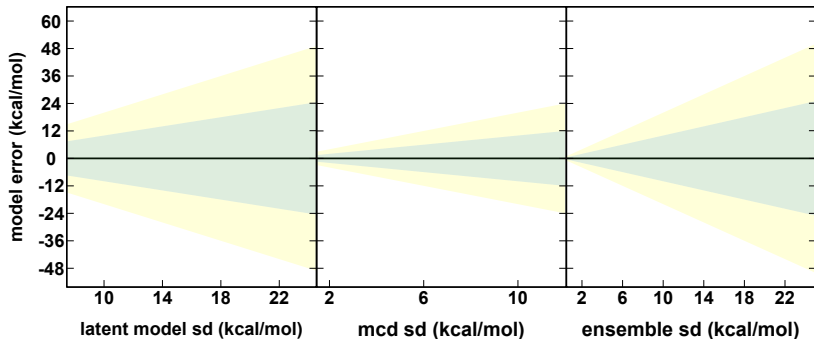


comparable with ensembles and mc dropout

# How do these distributions compare?

Comparison in energy units:

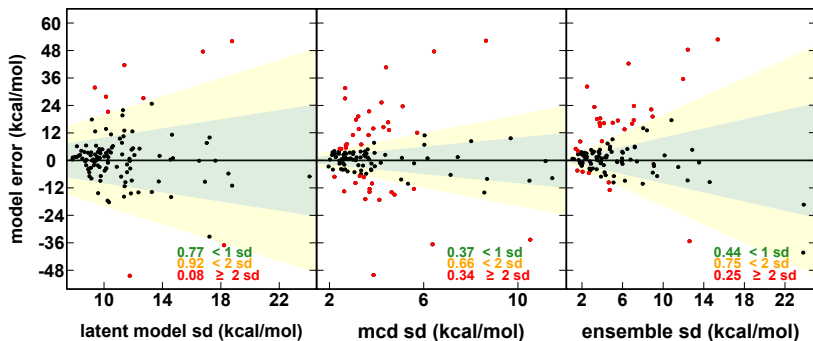
$$\varepsilon(d) \sim \mathcal{N}(0, \sigma_1^2 + d\sigma_2^2)$$



# How do these distributions compare?

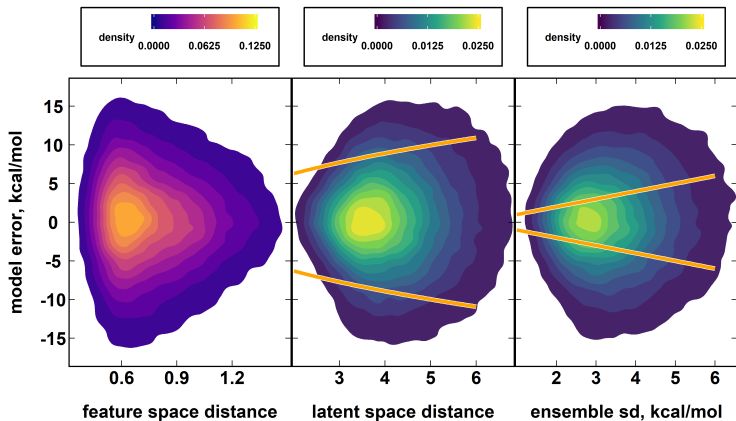
Comparison in energy units:

$$\varepsilon(d) \sim \mathcal{N}(0, \sigma_1^2 + d\sigma_2^2)$$



# How do these distributions compare?

Comparison in energy units: QM9 atomization energy data

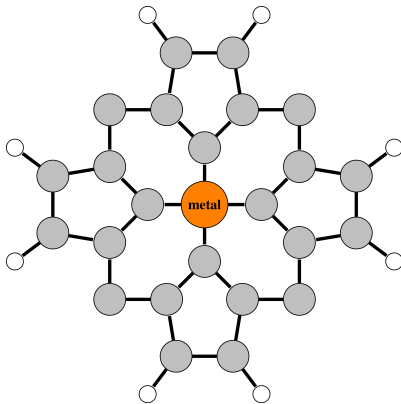


# Targeted data acquisition

Lack of data is a persistent issue. We can exploit knowledge of feature importance to generate data:

# Targeted data acquisition

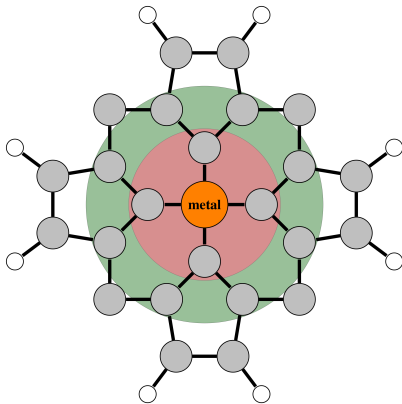
Lack of data is a persistent issue. We can exploit knowledge of feature importance to generate data:





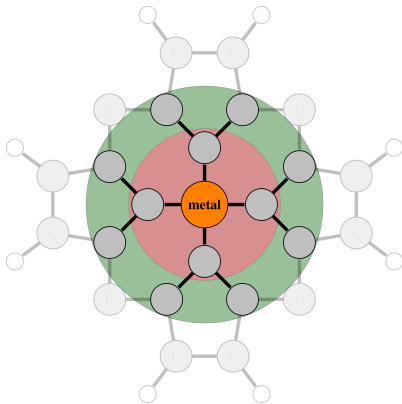
# Targeted data acquisition

Lack of data is a persistent issue. We can exploit knowledge of feature importance to generate data:



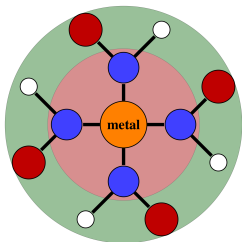
# Targeted data acquisition

Lack of data is a persistent issue. We can exploit knowledge of feature importance to generate data:



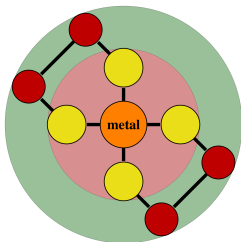
# Targeted data acquisition

Lack of data is a persistent issue. We can exploit knowledge of feature importance to generate data:



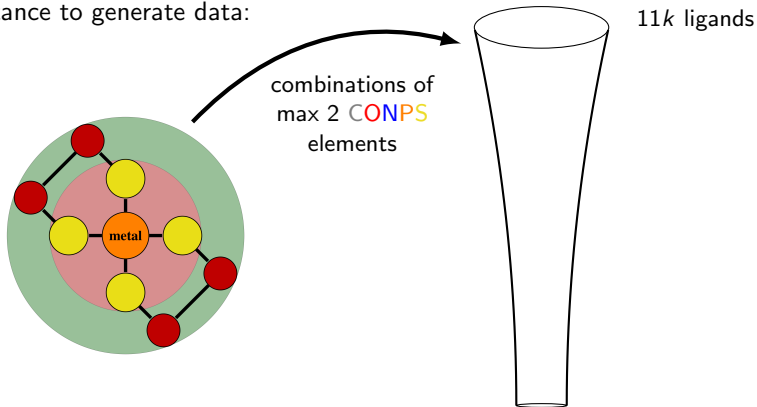
# Targeted data acquisition

Lack of data is a persistent issue. We can exploit knowledge of feature importance to generate data:



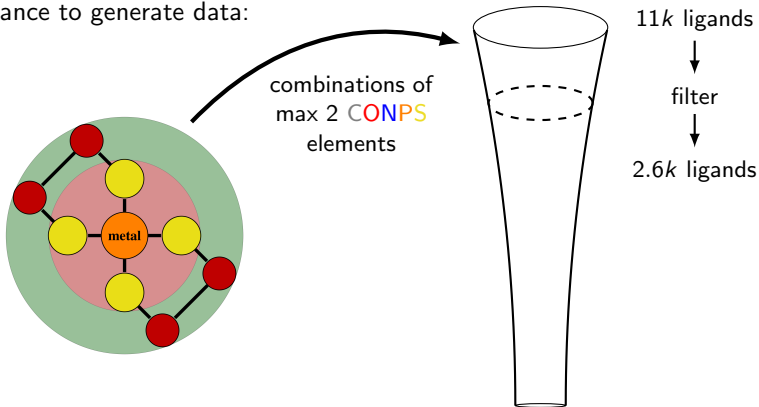
# Targeted data acquisition

Lack of data is a persistent issue. We can exploit knowledge of feature importance to generate data:



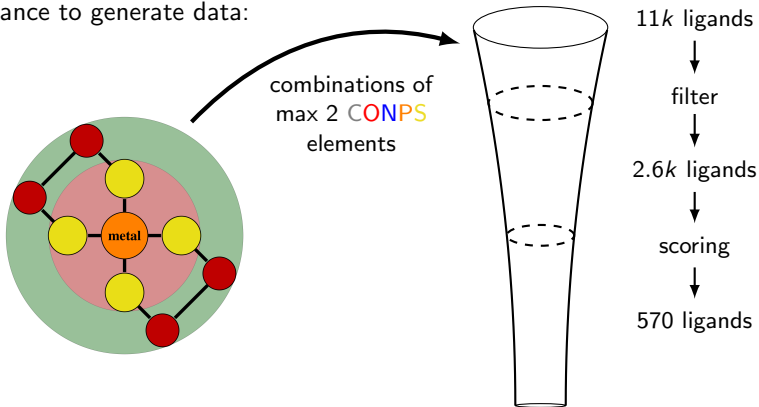
# Targeted data acquisition

Lack of data is a persistent issue. We can exploit knowledge of feature importance to generate data:



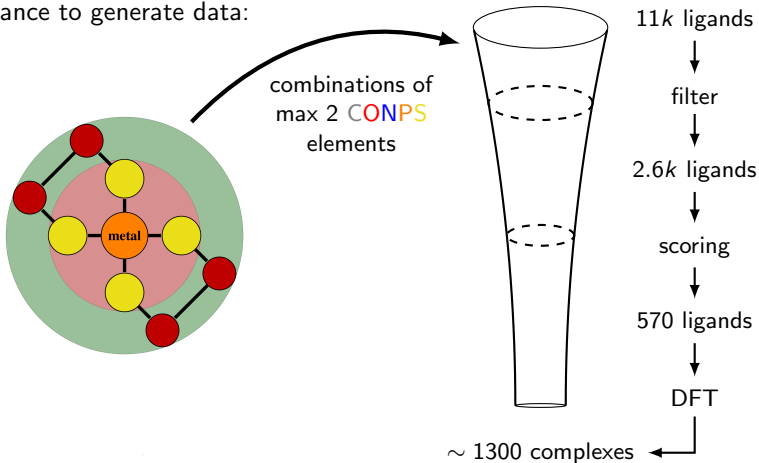
# Targeted data acquisition

Lack of data is a persistent issue. We can exploit knowledge of feature importance to generate data:



# Targeted data acquisition

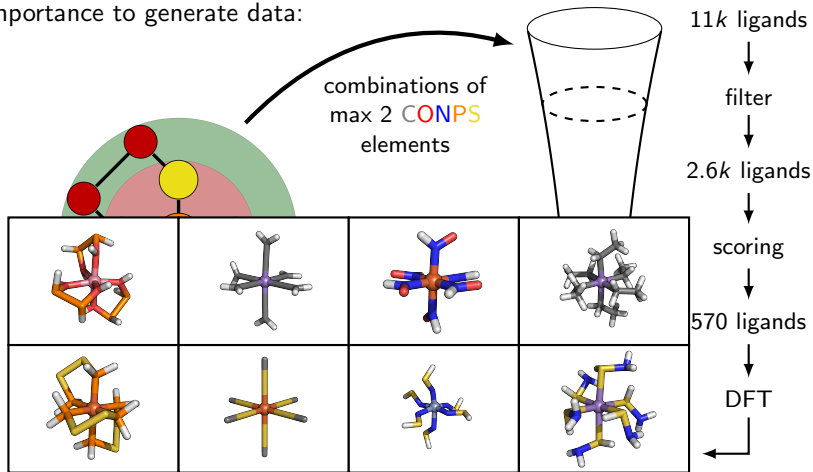
Lack of data is a persistent issue. We can exploit knowledge of feature importance to generate data:





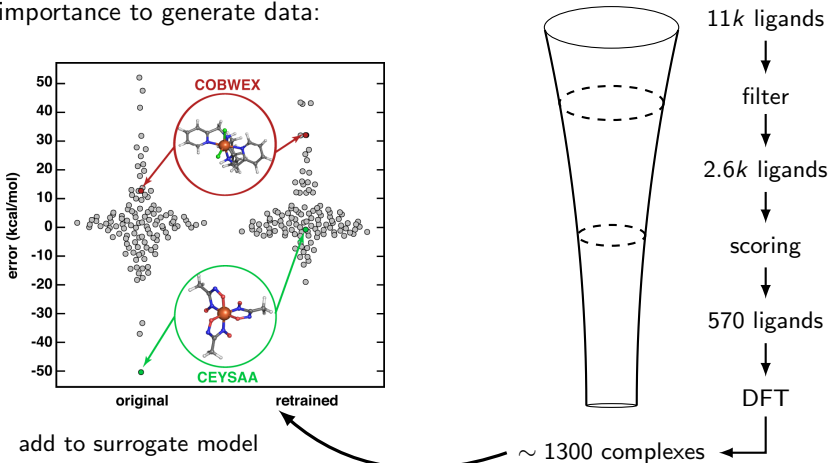
# Targeted data acquisition

Lack of data is a persistent issue. We can exploit knowledge of feature importance to generate data:



# Targeted data acquisition

Lack of data is a persistent issue. We can exploit knowledge of feature importance to generate data:

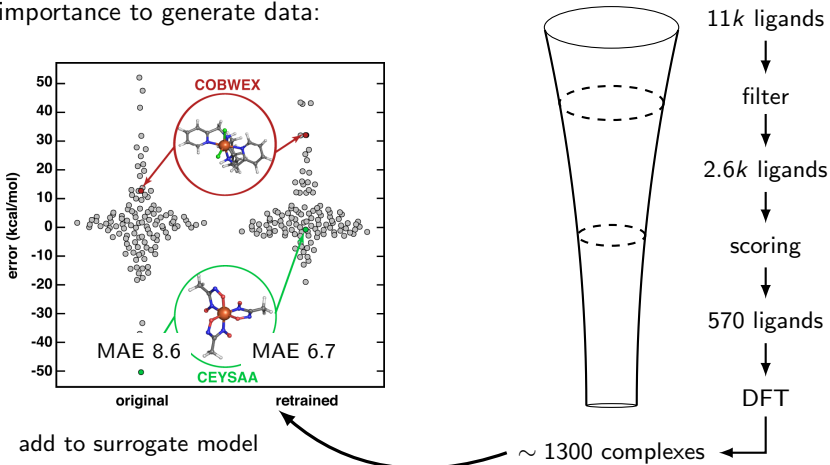


$3 \times 100$  tanh nodes, fully connected ANN

Gugler, S., Janet, J.P., and Kulik, H.J., *Mol. Sys. Des. Eng.*, Advance Article, 2019.

# Targeted data acquisition

Lack of data is a persistent issue. We can exploit knowledge of feature importance to generate data:



$3 \times 100$  tanh nodes, fully connected ANN

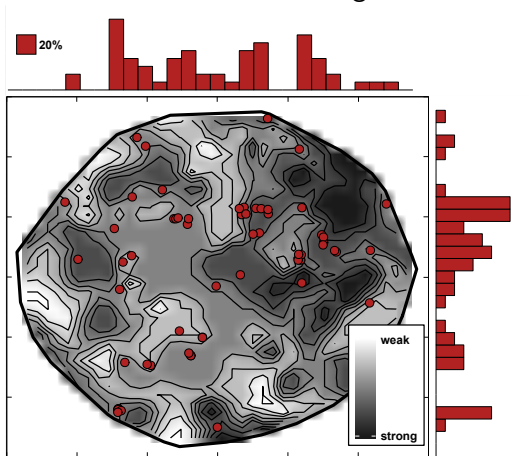
Gugler, S., Janet, J.P., and Kulik, H.J., *Mol. Sys. Des. Eng.*, Advance Article, 2019.

# Hedging against DFT uncertainty

Because we have trained our models with varying exact exchange, we can tune functionals for design:

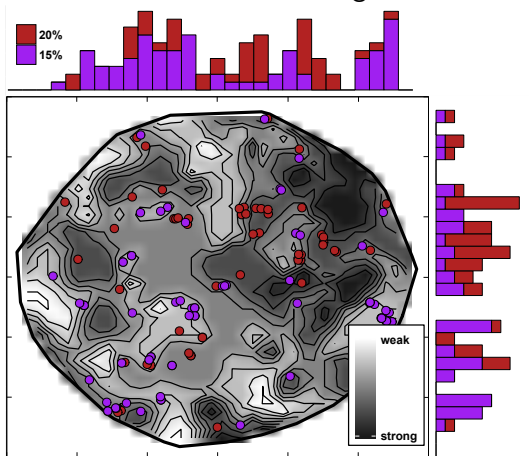
# Hedging against DFT uncertainty

Because we have trained our models with varying exact exchange, we can tune functionals for design:



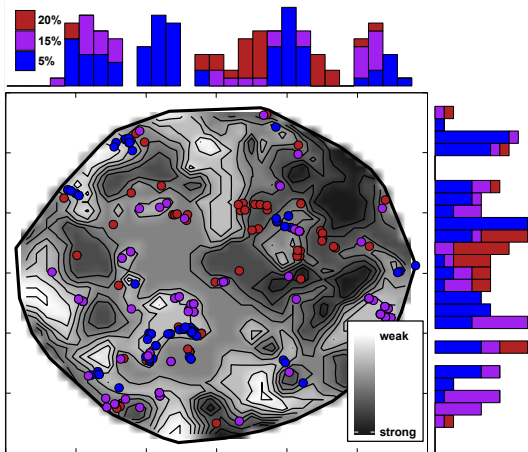
# Hedging against DFT uncertainty

Because we have trained our models with varying exact exchange, we can tune functionals for design:



# Hedging against DFT uncertainty

Because we have trained our models with varying exact exchange, we can tune functionals for design:



## Using models for discovery

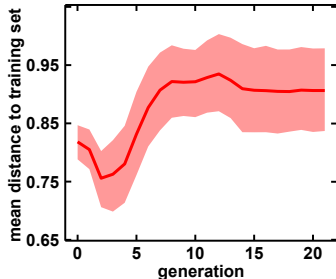
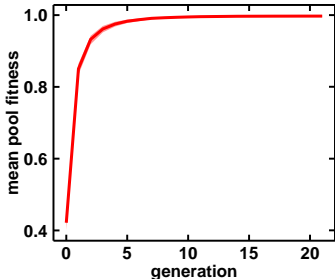
We utilize evolutionary algorithms to conduct design, guided by uncertainty metrics



# Using models for discovery

We utilize evolutionary algorithms to conduct design, guided by uncertainty metrics

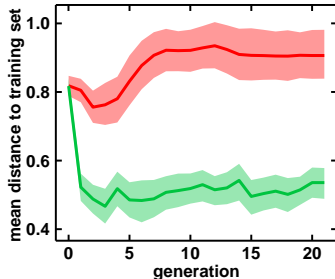
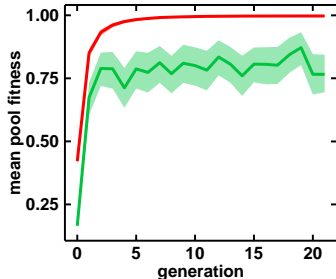
$$F_{s+d}(x) := \exp \left[ - \left( \frac{\Delta E_{H-L}(x)}{\rho \Delta E_{H-L}} \right)^2 \right]$$



# Using models for discovery

We utilize evolutionary algorithms to conduct design, guided by uncertainty metrics

$$F_{s+d}(x) := \exp \left[ - \left( \frac{\Delta E_{H-L}(x)}{\rho \Delta E_{H-L}} \right)^2 \right] \exp \left[ - \left( \frac{d(x)}{\rho_d} \right)^2 \right]$$



## Using models for discovery

We utilize evolutionary algorithms to conduct design, guided by uncertainty metrics

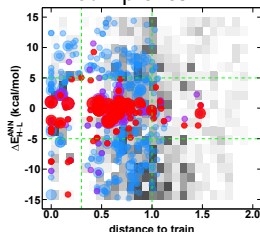
$$F_{s+d}(x) := \exp \left[ - \left( \frac{\Delta E_{H-L}(x)}{p \Delta E_{H-L}} \right)^2 \right] \exp \left[ - \left( \frac{d(x)}{p_d} \right)^2 \right]$$

# Using models for discovery

We utilize evolutionary algorithms to conduct design, guided by uncertainty metrics

$$F_{s+d}(x) := \exp \left[ - \left( \frac{\Delta E_{H-L}(x)}{p\Delta E_{H-L}} \right)^2 \right] \exp \left[ - \left( \frac{d(x)}{p_d} \right)^2 \right]$$

spin crossover  
complexes



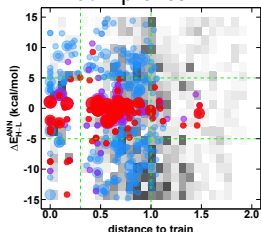
Janet, J.P., Chan, L. and Kulik, H.J., *J. Phys. Chem. Lett.*, 9(5):1064–1071, 2018.

# Using models for discovery

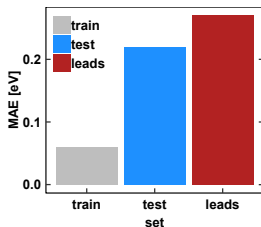
We utilize evolutionary algorithms to conduct design, guided by uncertainty metrics

$$F_{s+d}(x) := \exp \left[ - \left( \frac{\Delta E_{H-L}(x)}{p\Delta E_{H-L}} \right)^2 \right] \exp \left[ - \left( \frac{d(x)}{p_d} \right)^2 \right]$$

spin crossover  
complexes



frontier orbital  
properties



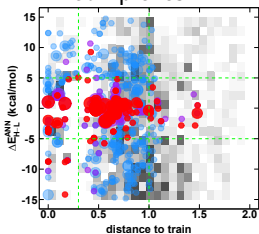
Janet, J.P., Chan, L. and Kulik, H.J., *J. Phys. Chem. Lett.*, 9(5):1064–1071, 2018.  
Nandy, A. et al., *Ind. Eng. Chem. Res.*, 57(42):13973–13986, 2018.

# Using models for discovery

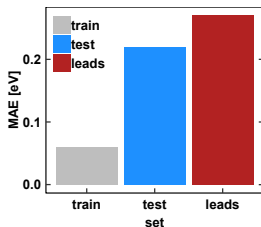
We utilize evolutionary algorithms to conduct design, guided by uncertainty metrics

$$F_{s+d}(x) := \exp \left[ - \left( \frac{\Delta E_{H-L}(x)}{p \Delta E_{H-L}} \right)^2 \right] \exp \left[ - \left( \frac{d(x)}{p_d} \right)^2 \right]$$

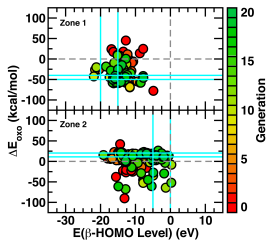
spin crossover  
complexes



frontier orbital  
properties



unusual catalytic  
reaction energies



Janet, J.P., Chan, L. and Kulik, H.J., *J. Phys. Chem. Lett.*, 9(5):1064–1071, 2018.

Nandy, A. et al., *Ind. Eng. Chem. Res.*, 57(42):13973–13986, 2018.

Nandy, A. et al., *ACS Catal.*, 9(9):8243–8255, 2019.

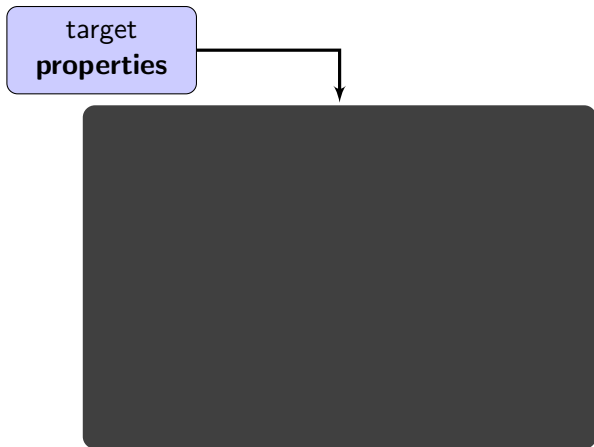
# More than just ML

# More than just ML

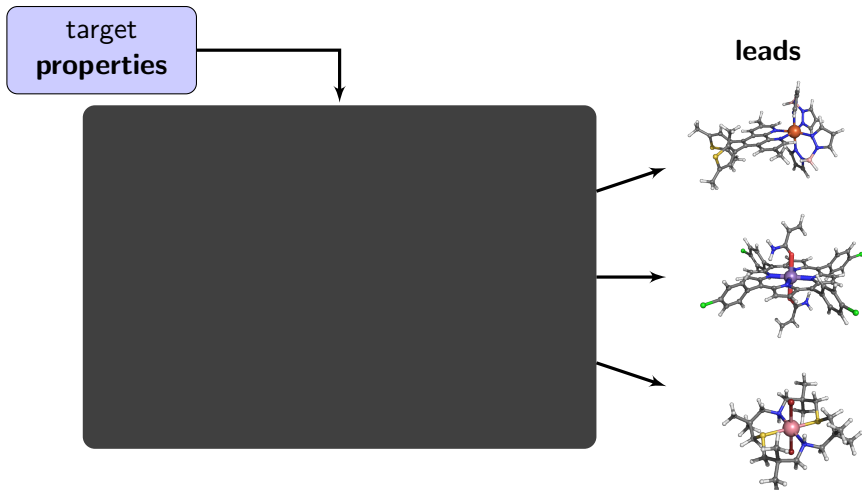
target  
**properties**



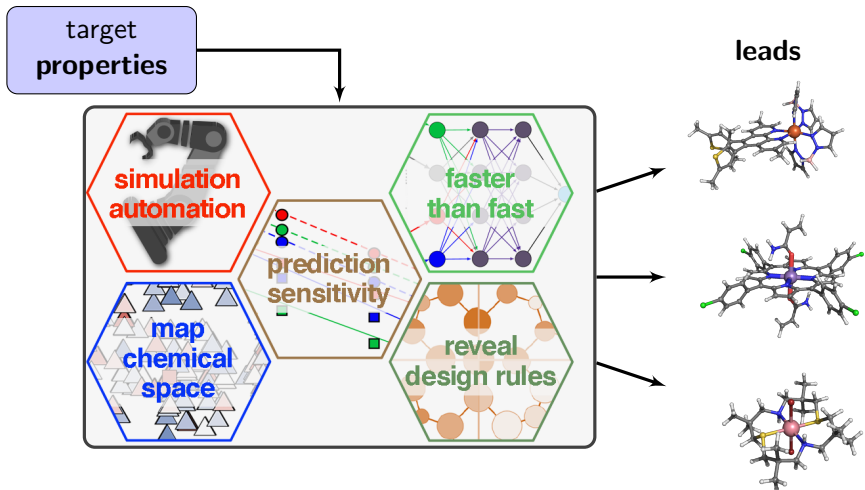
# More than just ML



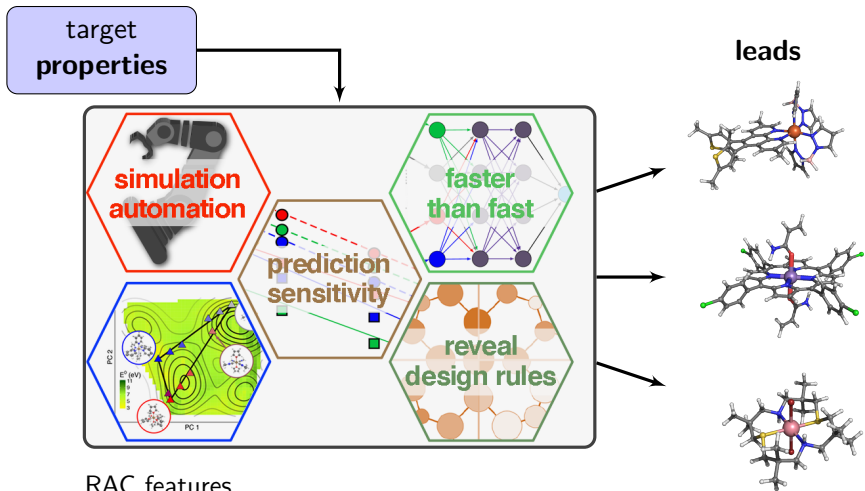
# More than just ML



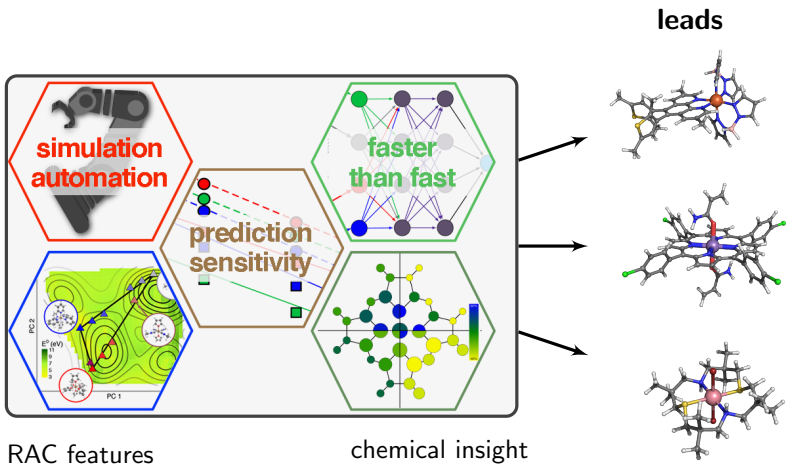
# More than just ML



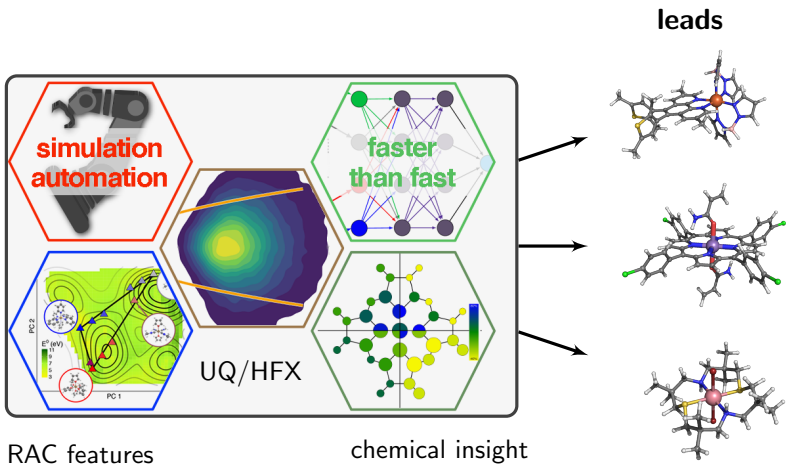
# More than just ML



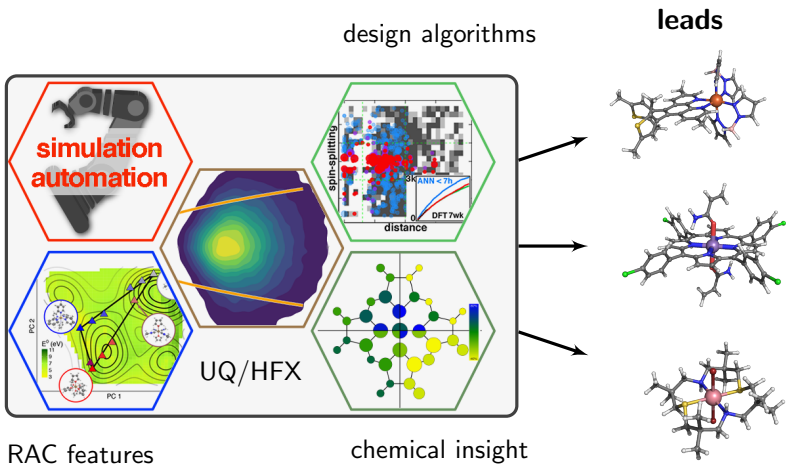
# More than just ML



# More than just ML



# More than just ML



# More than just ML

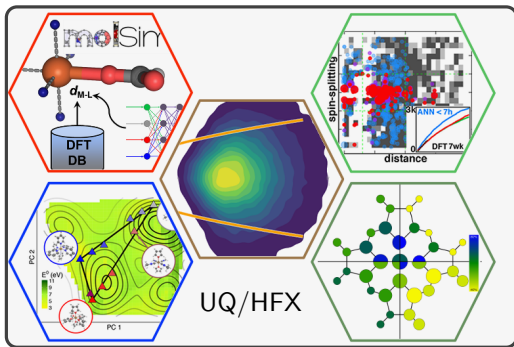


[github.com/hjkgrp](https://github.com/hjkgrp)

mAD/molSimplify

design algorithms

leads



RAC features

chemical insight



# More than just ML

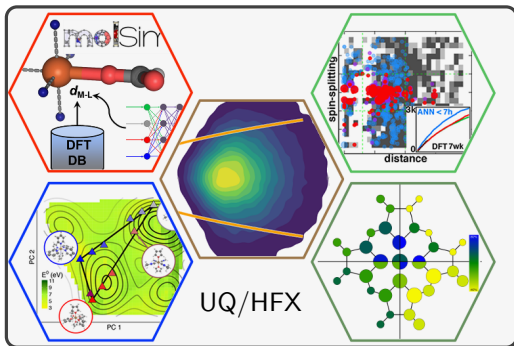


[github.com/hjkgrp](https://github.com/hjkgrp)

mAD/molSimplify

design algorithms

leads



RAC features

chemical insight

An end-to-end approach for rational TM complex design

## Case study: redox couples

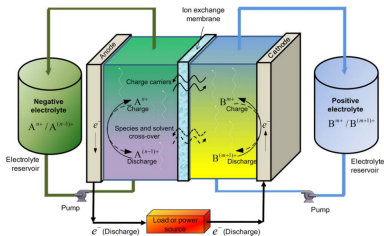
Consider a multiobjective optimization problem:

Redox flow batteries (RFBs)  
are a promising option for  
scalable energy storage:

# Case study: redox couples

Consider a multiobjective optimization problem:

Redox flow batteries (RFBs)  
are a promising option for  
scalable energy storage:



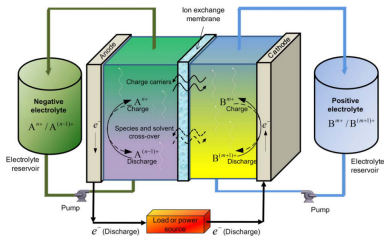
Perry, M.L. and Adam, Z., *J. Electrochem. Soc.*, 163(1):A5064–A5067, 2018.

## Case study: redox couples

Consider a multiobjective optimization problem:

Redox flow batteries (RFBs) are a promising option for scalable energy storage:

Transition metal complexes make attractive redox couples for RFBs



Perry, M.L. and Adam, Z., *J. Electrochem. Soc.*, 163(1):A5064–A5067, 2018.

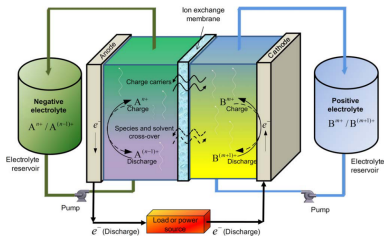
## Case study: redox couples

Consider a multiobjective optimization problem:

Redox flow batteries (RFBs) are a promising option for scalable energy storage:

Transition metal complexes make attractive redox couples for RFBs

- good ion stability (compared to organics)



Perry, M.L. and Adam, Z., *J. Electrochem. Soc.*, 163(1):A5064–A5067, 2018.

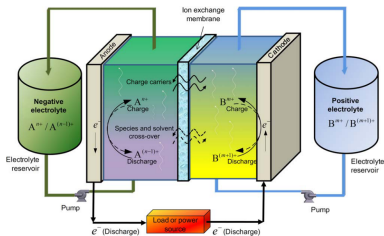
## Case study: redox couples

Consider a multiobjective optimization problem:

Redox flow batteries (RFBs) are a promising option for scalable energy storage:

Transition metal complexes make attractive redox couples for RFBs

- good ion stability (compared to organics)
- good range of redox potentials available



Perry, M.L. and Adam, Z., *J. Electrochem. Soc.*, 163(1):A5064–A5067, 2018.

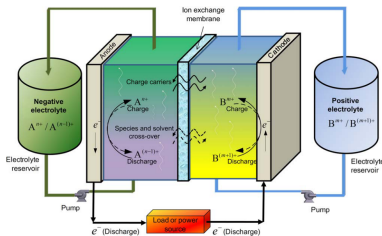
## Case study: redox couples

Consider a multiobjective optimization problem:

Redox flow batteries (RFBs) are a promising option for scalable energy storage:

Transition metal complexes make attractive redox couples for RFBs

- good ion stability (compared to organics)
- good range of redox potentials available
- **solubility is an issue!**

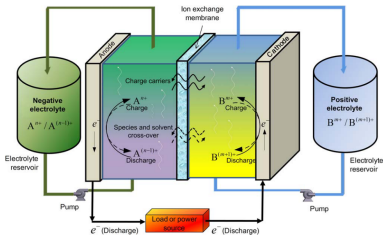


Perry, M.L. and Adam, Z., *J. Electrochem. Soc.*, 163(1):A5064–A5067, 2018.

## Case study: redox couples

Consider a multiobjective optimization problem:

Redox flow batteries (RFBs) are a promising option for scalable energy storage:



Transition metal complexes make attractive redox couples for RFBs

- good ion stability (compared to organics)
- good range of redox potentials available
- **solubility is an issue!**

$$E_{\text{cell}} = 0.5 \times V_{\text{cell}} \times C \times n \times F$$

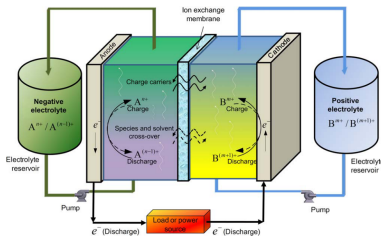
Perry, M.L. and Adam, Z., *J. Electrochem. Soc.*, 163(1):A5064–A5067, 2018.



## Case study: redox couples

Consider a multiobjective optimization problem:

Redox flow batteries (RFBs) are a promising option for scalable energy storage:



Perry, M.L. and Adam, Z., *J. Electrochem. Soc.*, 163(1):A5064–A5067, 2018.

Transition metal complexes make attractive redox couples for RFBs

- good ion stability (compared to organics)
- good range of redox potentials available
- **solubility is an issue!**

$$E_{\text{cell}} = 0.5 \times V_{\text{cell}} \times C \times n \times F$$

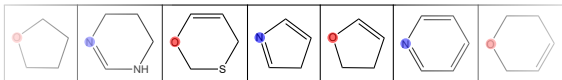
We need complexes that have high redox potential **and** good solubility

# Design space construction

# Design space construction

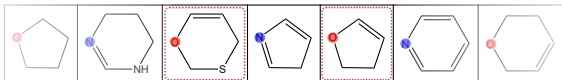


# Design space construction



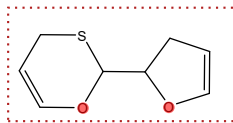
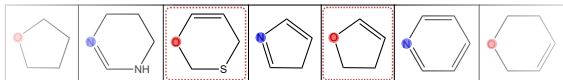
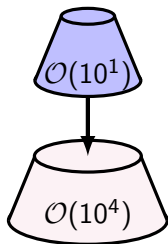
40 heterocycles

# Design space construction



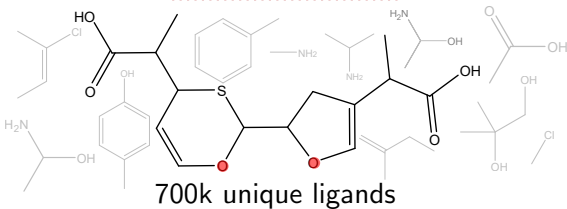
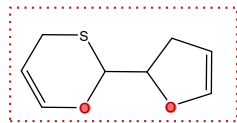
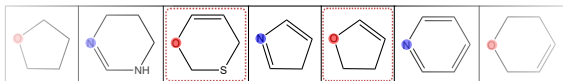
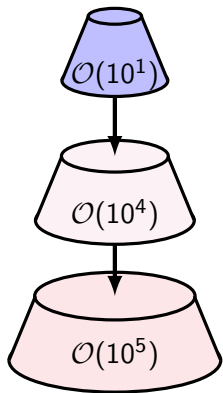
40 heterocycles

# Design space construction

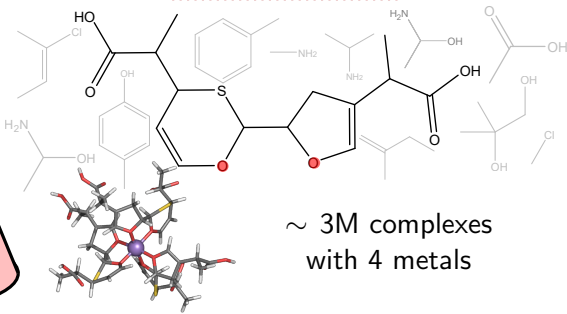
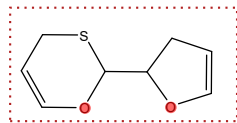
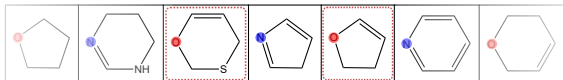
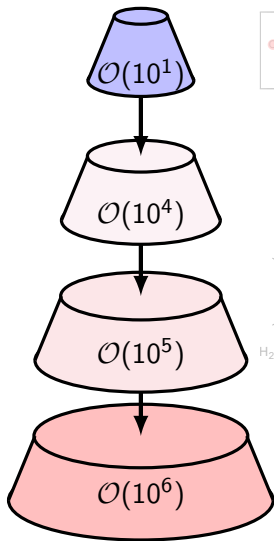


800 base ligands

# Design space construction



# Design space construction





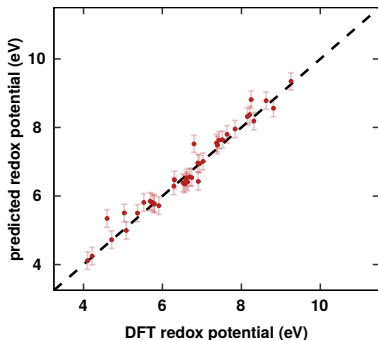
# Multiobjective framework

We can predict quantities of interest for our RFBs:

# Multiobjective framework

We can predict quantities of interest for our RFBs:

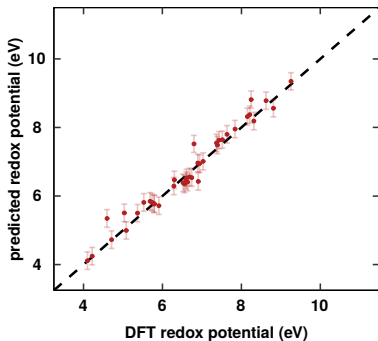
$$\Delta G_{\text{solv}} \approx \Delta E_{\text{III-II}} + \Delta \Delta G_s$$



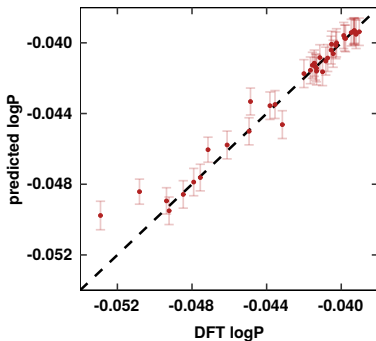
# Multiobjective framework

We can predict quantities of interest for our RFBs:

$$\Delta G_{\text{soln}} \approx \Delta E_{\text{III-II}} + \Delta \Delta G_s$$



$$\log P \approx \log \frac{\Delta G_{s,\text{II,octanol}}}{\Delta G_{s,\text{II,water}}}$$

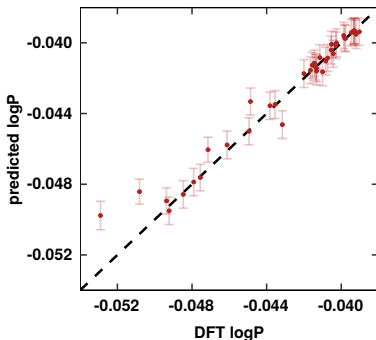
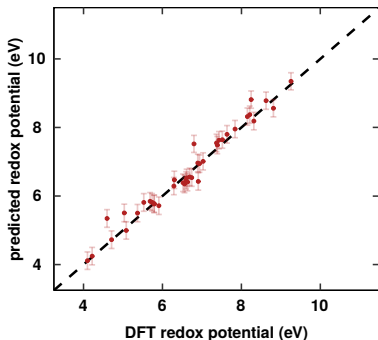


# Multiobjective framework

We can predict quantities of interest for our RFBs:

$$\Delta G_{\text{solv}} \approx \Delta E_{\text{III-II}} + \Delta \Delta G_s$$

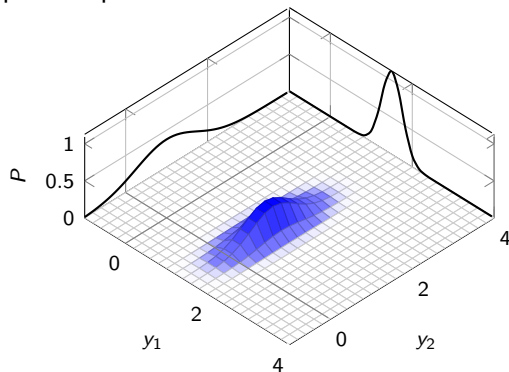
$$\log P \approx \log \frac{\Delta G_{s,\text{II},\text{octanol}}}{\Delta G_{s,\text{II},\text{water}}}$$



$$\begin{bmatrix} \Delta G_{\text{solv}} \\ \log P \end{bmatrix} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \end{bmatrix} \sim \mathcal{N} \left( \begin{bmatrix} \hat{\mu}_1 \\ \hat{\mu}_2 \end{bmatrix}, \begin{bmatrix} \hat{\sigma}_1^2 & 0 \\ 0 & \hat{\sigma}_2^2 \end{bmatrix} \right)$$

# Multiobjective framework

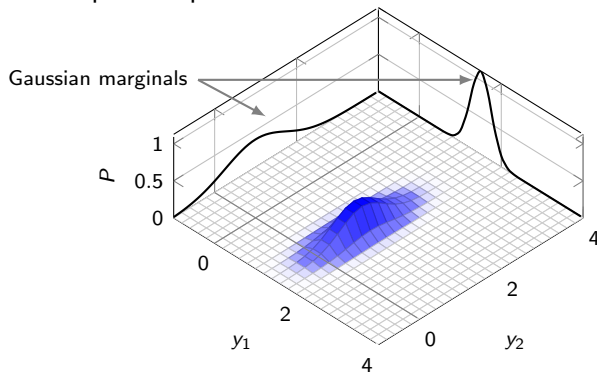
We can predict quantities of interest for our RFBs:



$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \end{bmatrix} \sim \mathcal{N} \left( \begin{bmatrix} \hat{\mu}_1 \\ \hat{\mu}_2 \end{bmatrix}, \begin{bmatrix} \hat{\sigma}_1^2 & 0 \\ 0 & \hat{\sigma}_2^2 \end{bmatrix} \right)$$

# Multiobjective framework

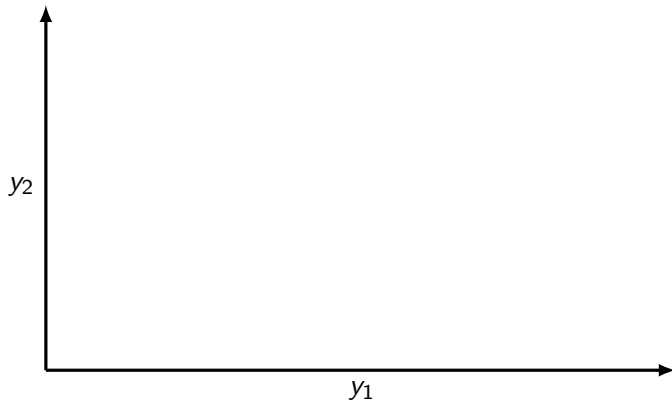
We can predict quantities of interest for our RFBs:



$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \end{bmatrix} \sim \mathcal{N} \left( \begin{bmatrix} \hat{\mu}_1 \\ \hat{\mu}_2 \end{bmatrix}, \begin{bmatrix} \hat{\sigma}_1^2 & 0 \\ 0 & \hat{\sigma}_2^2 \end{bmatrix} \right)$$

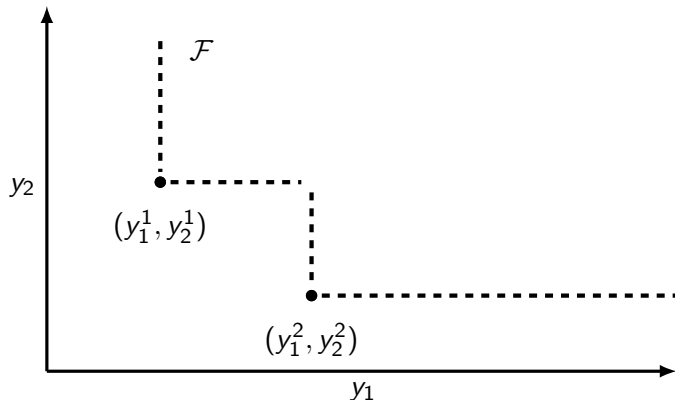
## 2D EGO Illustration

We will use a multiobjective expected improvement framework:



## 2D EGO Illustration

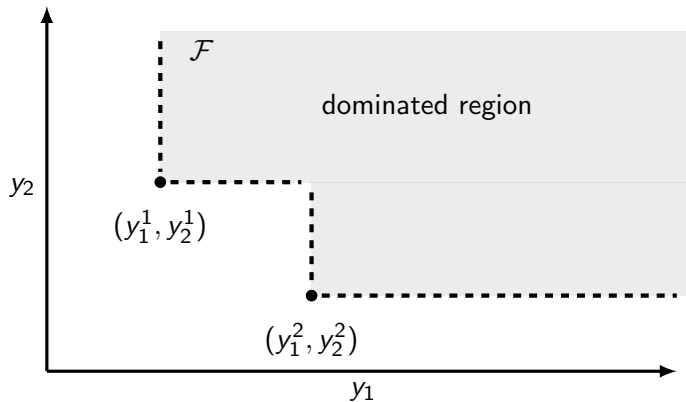
We will use a multiobjective expected improvement framework:





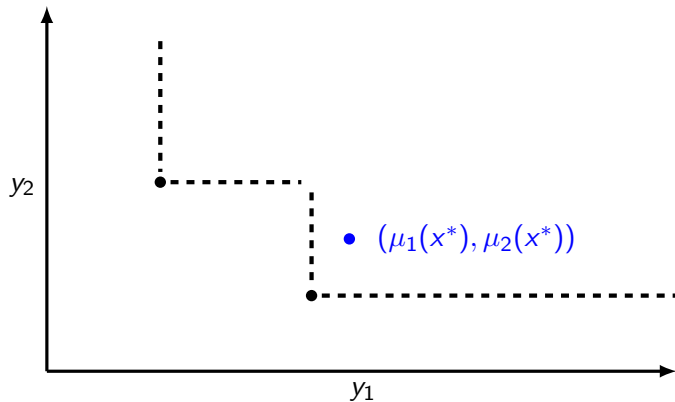
## 2D EGO Illustration

We will use a multiobjective expected improvement framework:



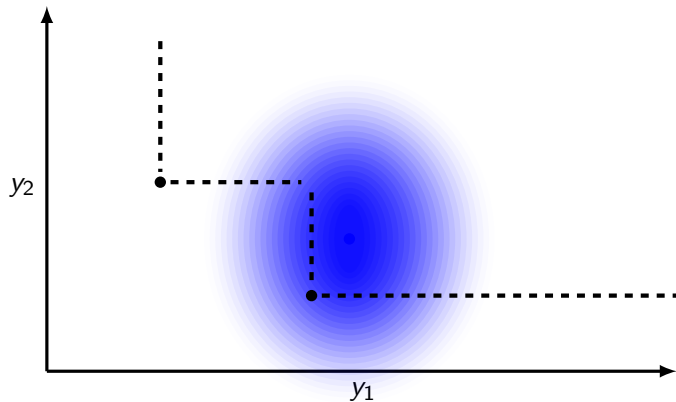
## 2D EGO Illustration

We will use a multiobjective expected improvement framework:



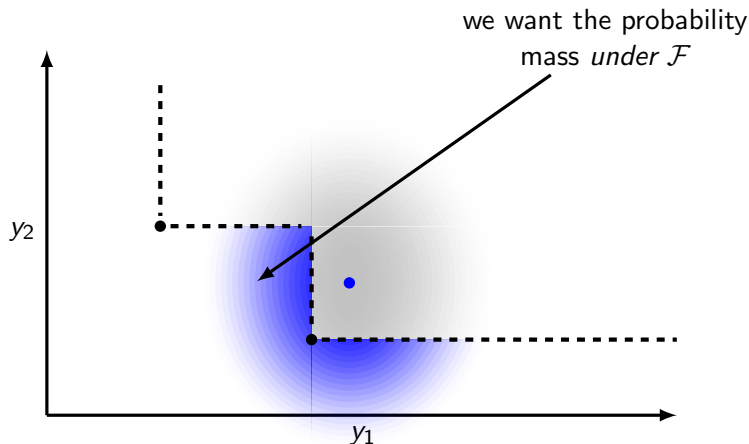
## 2D EGO Illustration

We will use a multiobjective expected improvement framework:



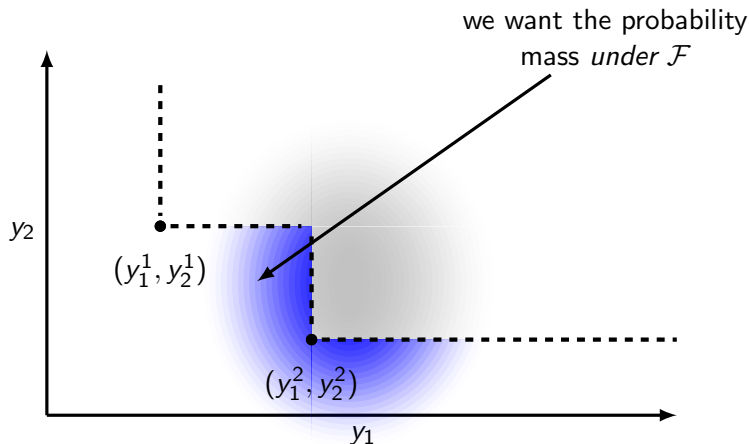
## 2D EGO Illustration

We will use a multiobjective expected improvement framework:



## 2D EGO Illustration

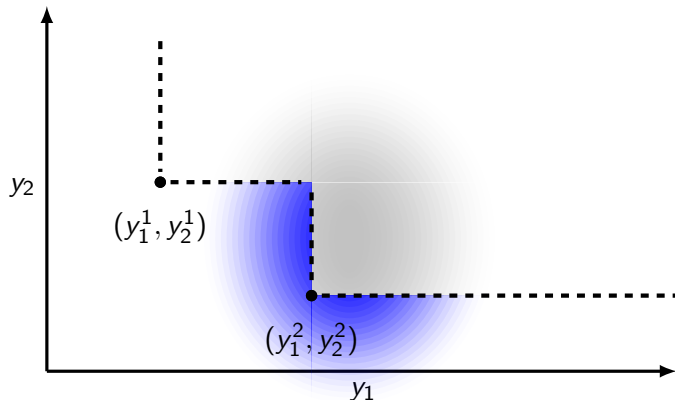
We will use a multiobjective expected improvement framework:



## 2D EGO Illustration

We will use a multiobjective expected improvement framework:

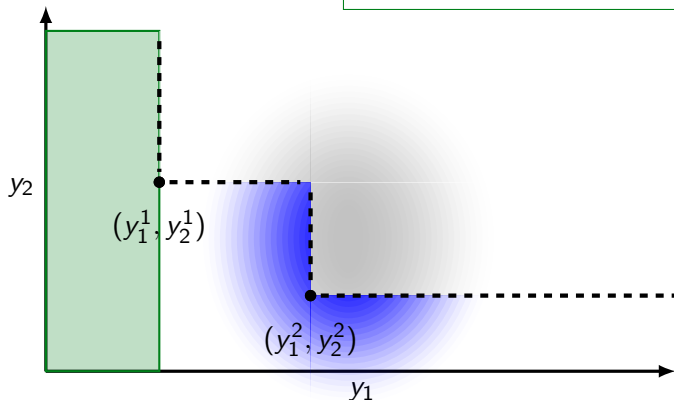
$$P(I) =$$



## 2D EGO Illustration

We will use a multiobjective expected improvement framework:

$$P(I) = \int_{-\infty}^{y_1^1} \int_{-\infty}^{\infty} p(\hat{y}_1(x^*), \hat{y}_2(x^*)) dy_1 dy_2$$

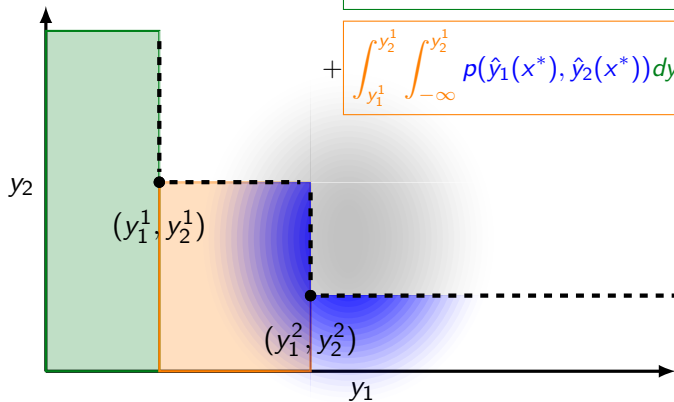


## 2D EGO Illustration

We will use a multiobjective expected improvement framework:

$$P(I) = \int_{-\infty}^{y_1^1} \int_{-\infty}^{\infty} p(\hat{y}_1(x^*), \hat{y}_2(x^*)) dy_1 dy_2$$

$$+ \int_{y_1^1}^{y_2^1} \int_{-\infty}^{y_2^1} p(\hat{y}_1(x^*), \hat{y}_2(x^*)) dy_1 dy_2$$





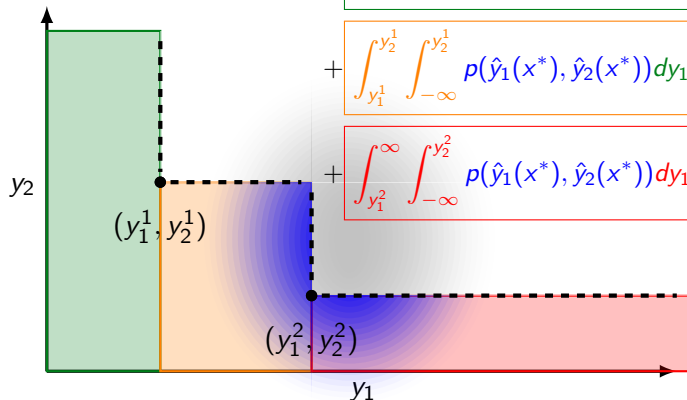
## 2D EGO Illustration

We will use a multiobjective expected improvement framework:

$$P(I) = \int_{-\infty}^{y_1^1} \int_{-\infty}^{\infty} p(\hat{y}_1(x^*), \hat{y}_2(x^*)) dy_1 dy_2$$

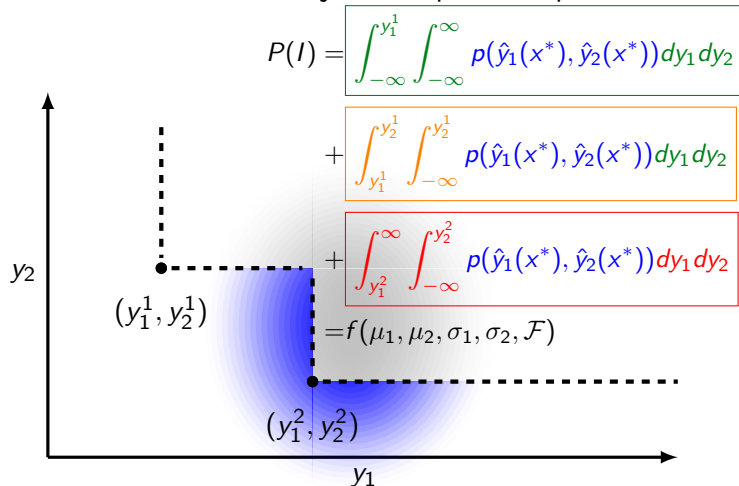
$$+ \int_{y_1^1}^{y_2^1} \int_{-\infty}^{y_2^1} p(\hat{y}_1(x^*), \hat{y}_2(x^*)) dy_1 dy_2$$

$$+ \int_{y_1^2}^{\infty} \int_{-\infty}^{y_2^2} p(\hat{y}_1(x^*), \hat{y}_2(x^*)) dy_1 dy_2$$



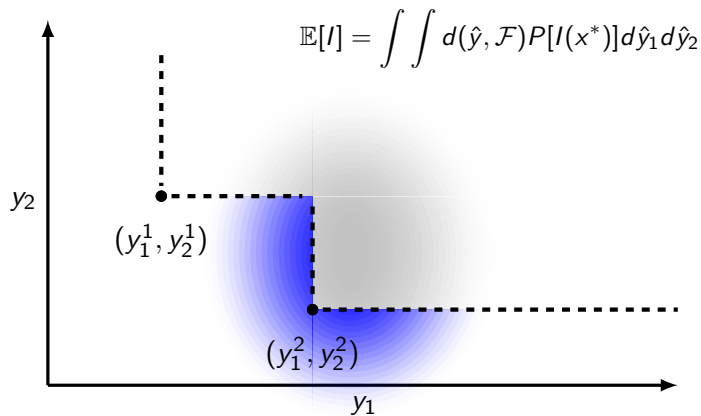
## 2D EGO Illustration

We will use a multiobjective expected improvement framework:



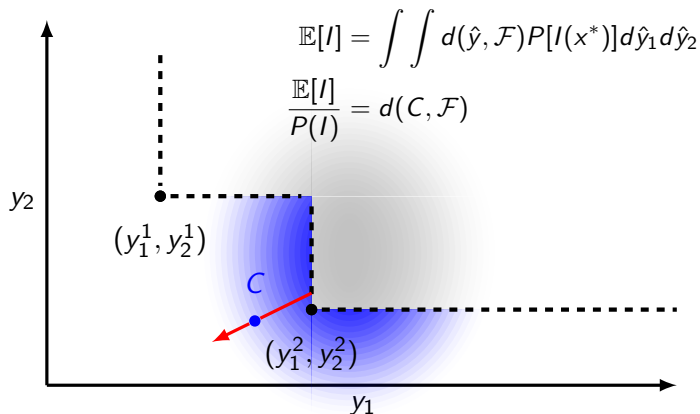
## 2D EGO Illustration

We will use a multiobjective expected improvement framework:



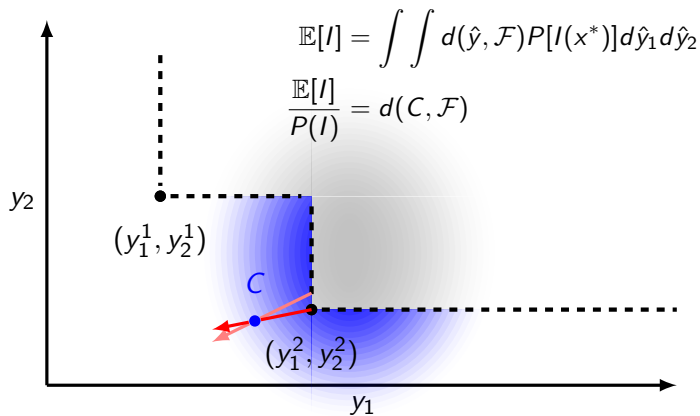
## 2D EGO Illustration

We will use a multiobjective expected improvement framework:



## 2D EGO Illustration

We will use a multiobjective expected improvement framework:

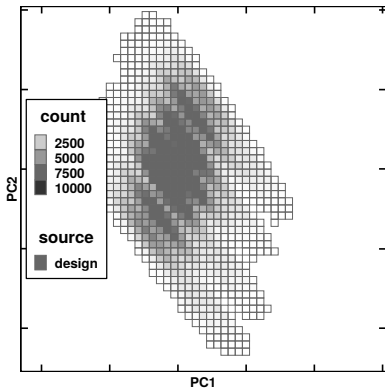


# Design space and existing data

We observe poor overlap between existing data and design space:

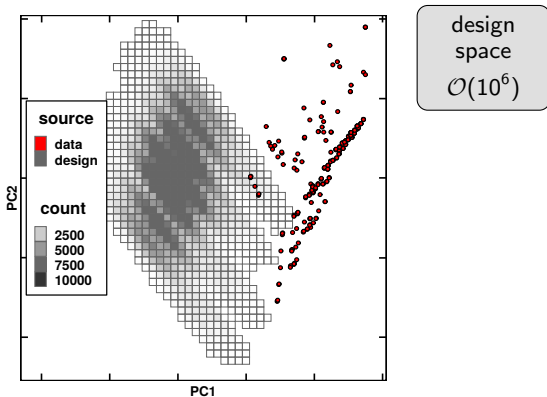
# Design space and existing data

We observe poor overlap between existing data and design space:



# Design space and existing data

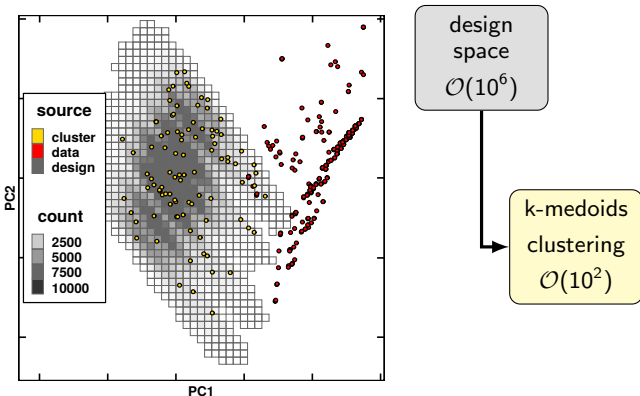
We observe poor overlap between existing data and design space:





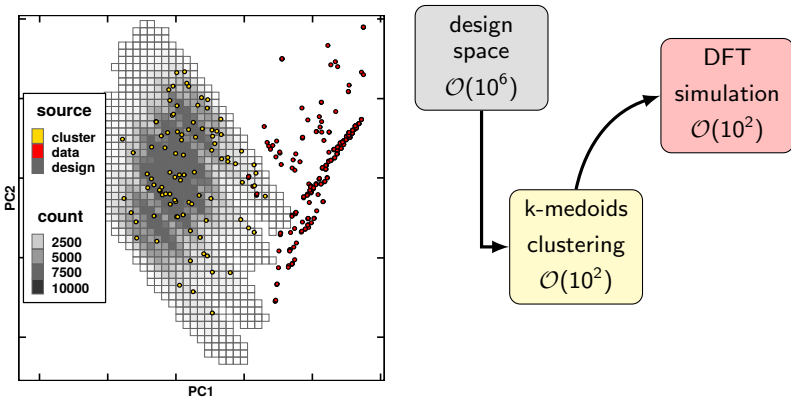
# Design space and existing data

We observe poor overlap between existing data and design space:



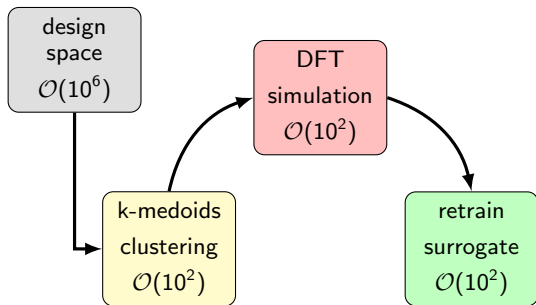
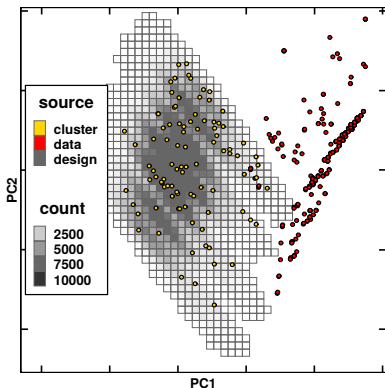
# Design space and existing data

We observe poor overlap between existing data and design space:



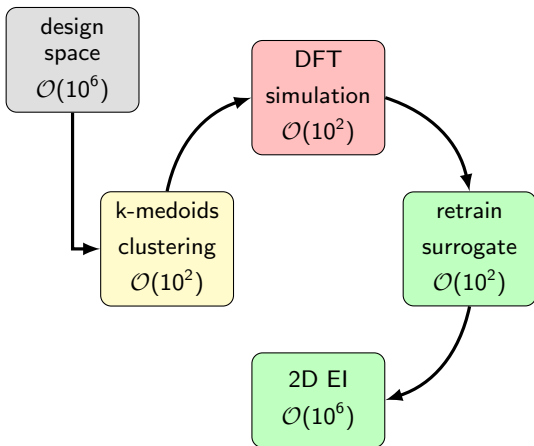
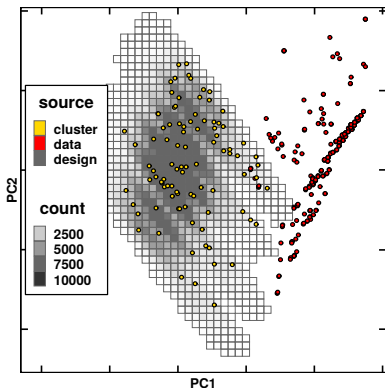
# Design space and existing data

We observe poor overlap between existing data and design space:



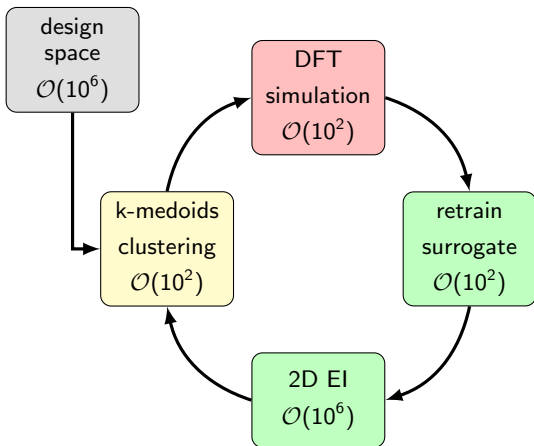
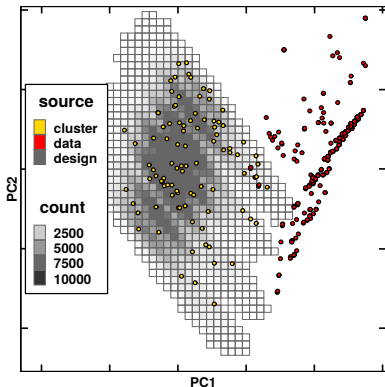
# Design space and existing data

We observe poor overlap between existing data and design space:



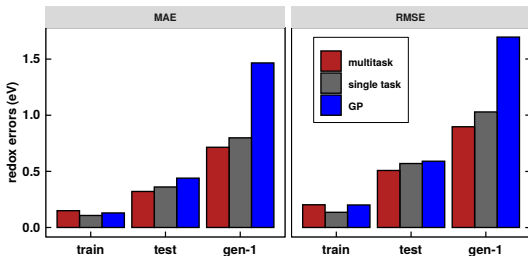
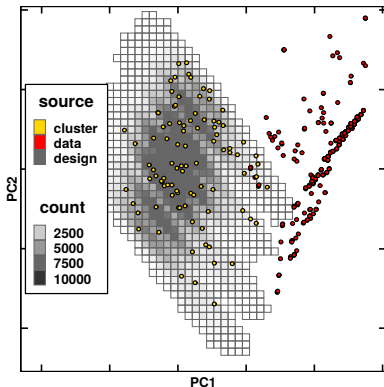
# Design space and existing data

We observe poor overlap between existing data and design space:



# Design space and existing data

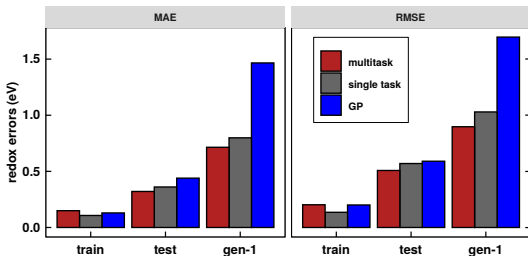
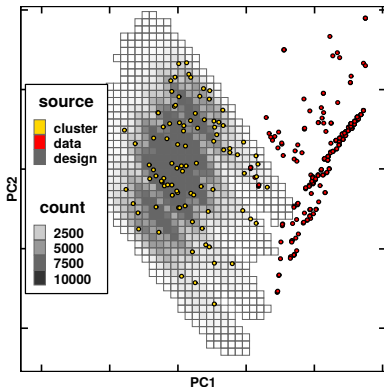
We observe poor overlap between existing data and design space:



single task redox:  $3 \times 100$  tanh nodes, fully connected  
single task log P:  $2 \times 50$  ReLU nodes, skip + residual connections  
multitask:  $2 \times 100$  tanh nodes, fully connected

# Design space and existing data

We observe poor overlap between existing data and design space:



Therefore, we will use multitask ANN to drive the design process

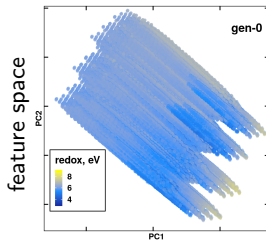
single task redox:  $3 \times 100$  tanh nodes, fully connected

single task log P:  $2 \times 50$  ReLU nodes, skip + residual connections

multitask:  $2 \times 100$  tanh nodes, fully connected

# Active learning

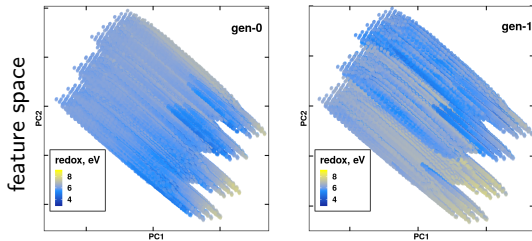
We can monitor the evolution of the model as data is acquired





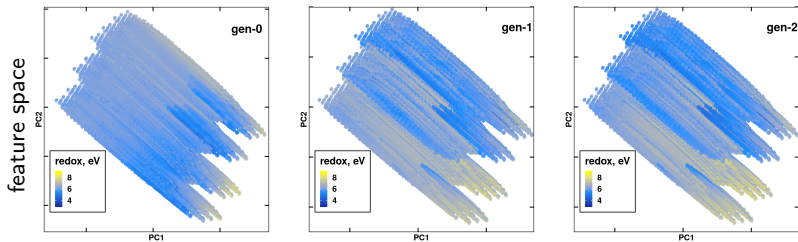
# Active learning

We can monitor the evolution of the model as data is acquired



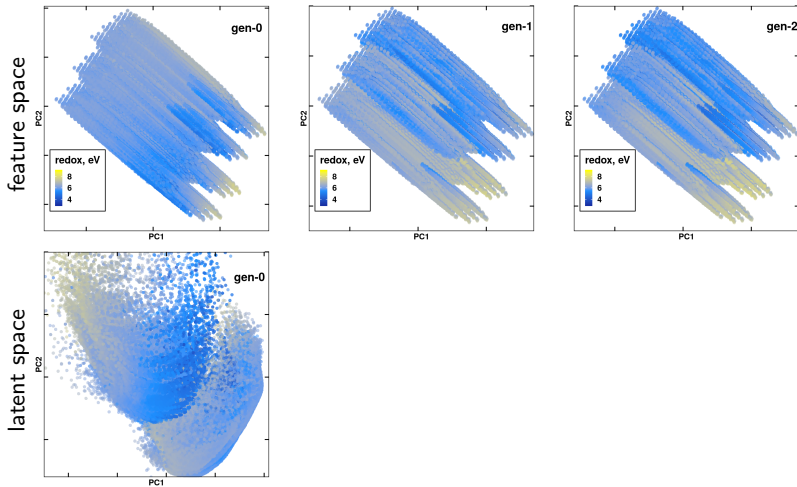
# Active learning

We can monitor the evolution of the model as data is acquired



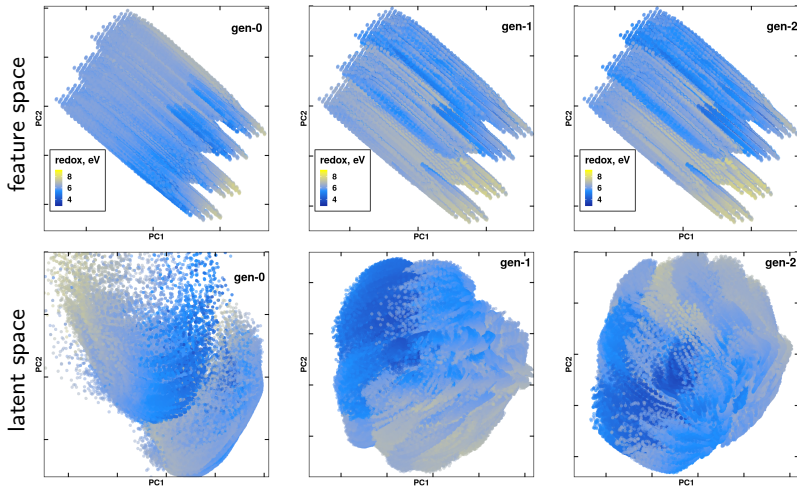
# Active learning

We can monitor the evolution of the model as data is acquired



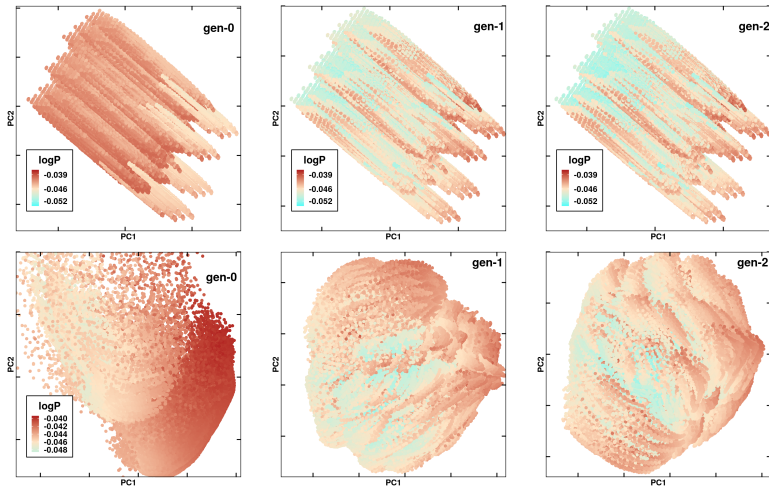
# Active learning

We can monitor the evolution of the model as data is acquired



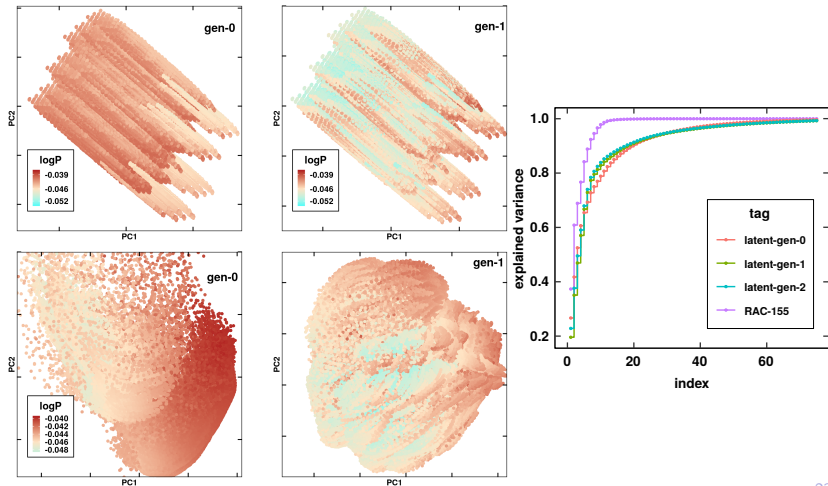
# Active learning

We can monitor the evolution of the model as data is acquired

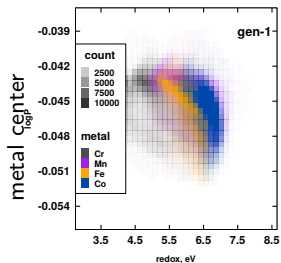


# Active learning

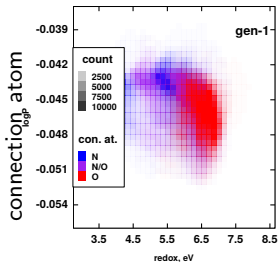
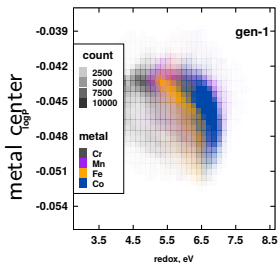
We can monitor the evolution of the model as data is acquired



# Mapping output space

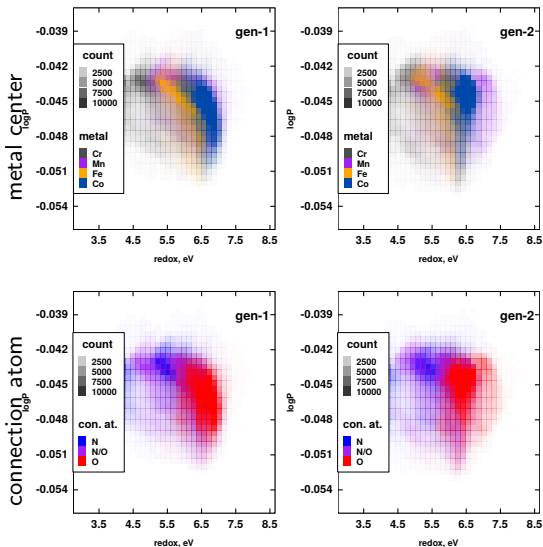


# Mapping output space

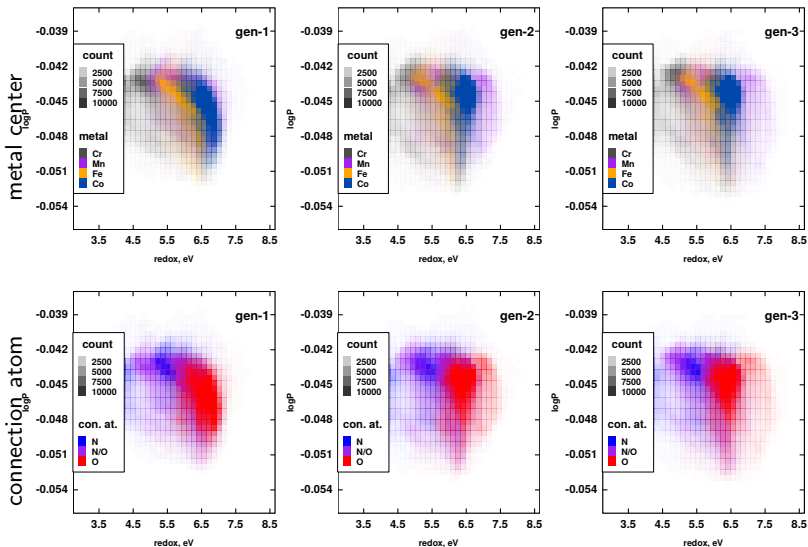




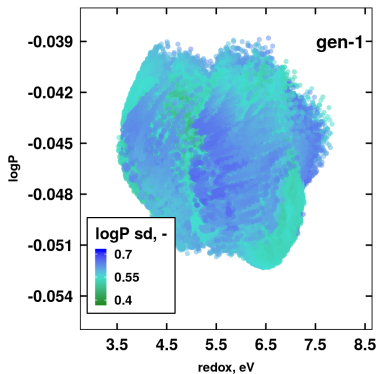
# Mapping output space



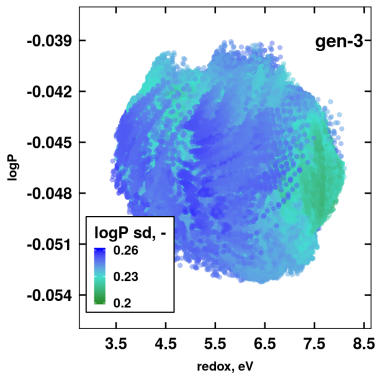
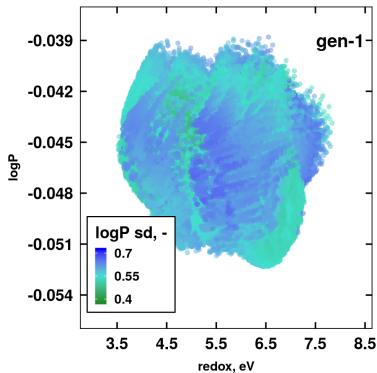
# Mapping output space



# Mapping output space



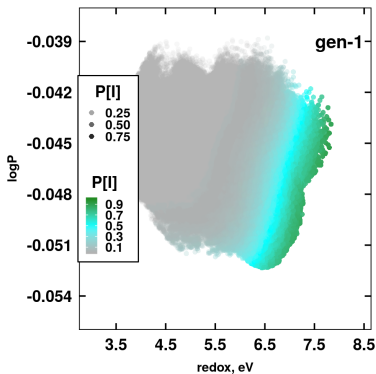
# Mapping output space



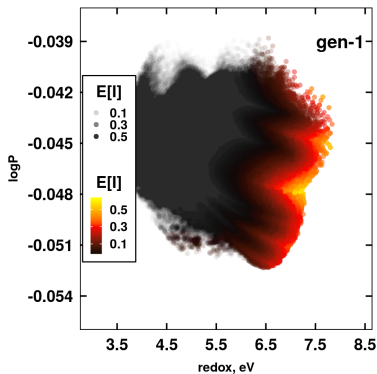
model confidence is localized in target region!

# EGO results

probability of improvement

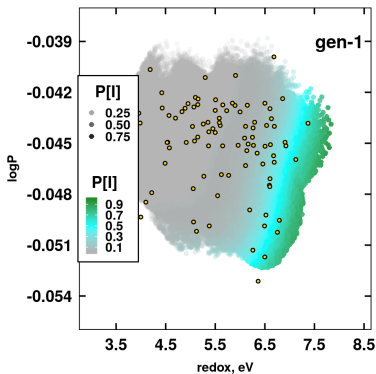


expected improvement

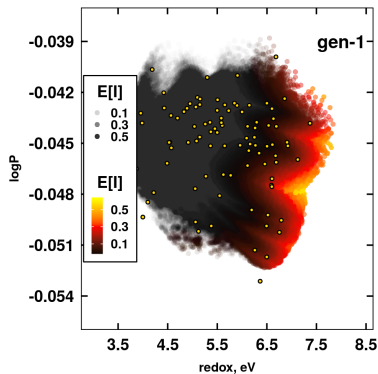


# EGO results

probability of improvement

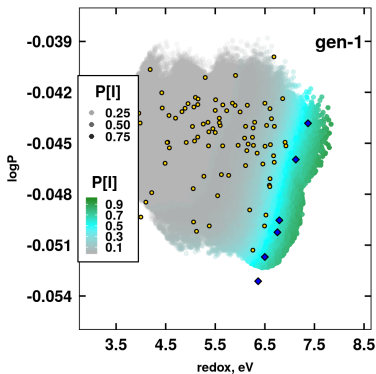


expected improvement

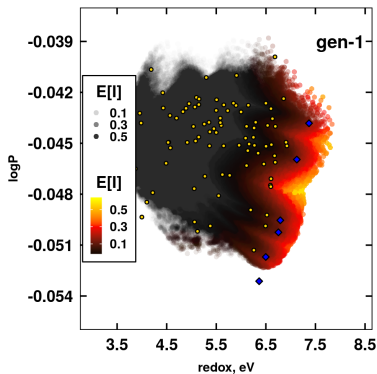


# EGO results

probability of improvement

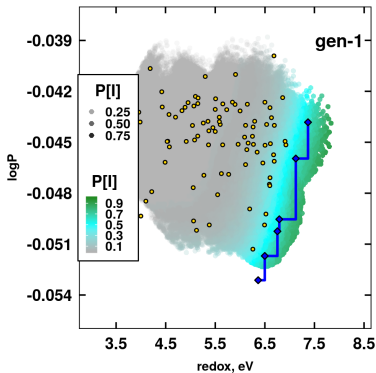


expected improvement

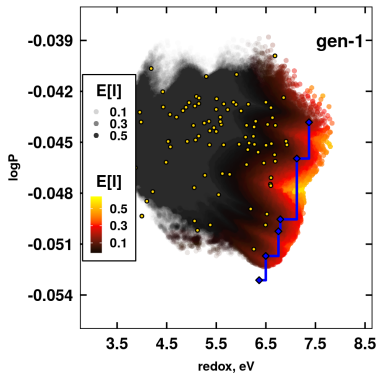


# EGO results

probability of improvement



expected improvement



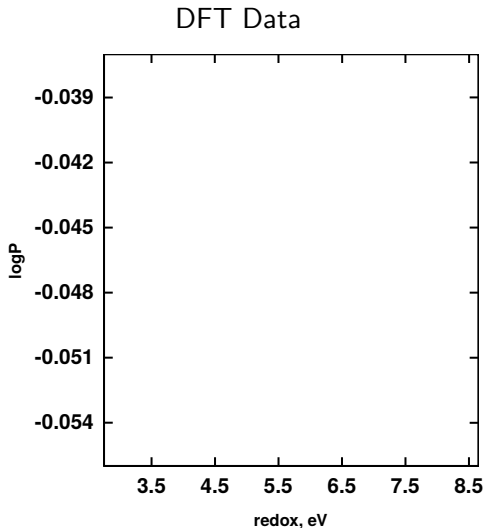


# EGO results

probability of improvement

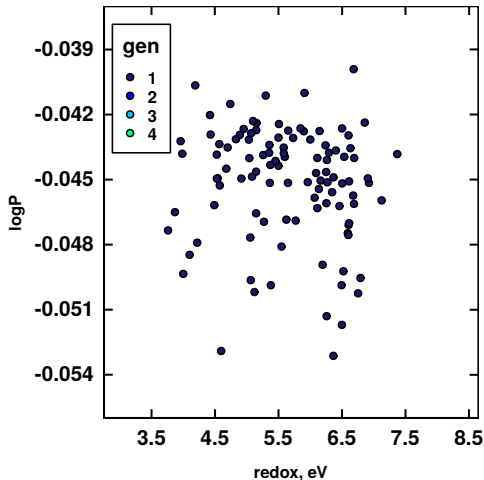
expected improvement

# DFT results



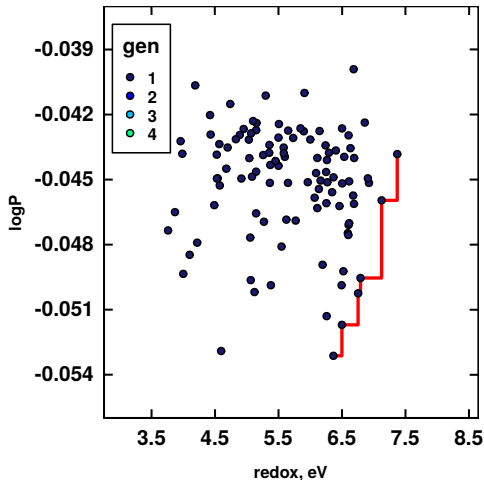
# DFT results

k-medoids points (gen 1: 107 complexes)



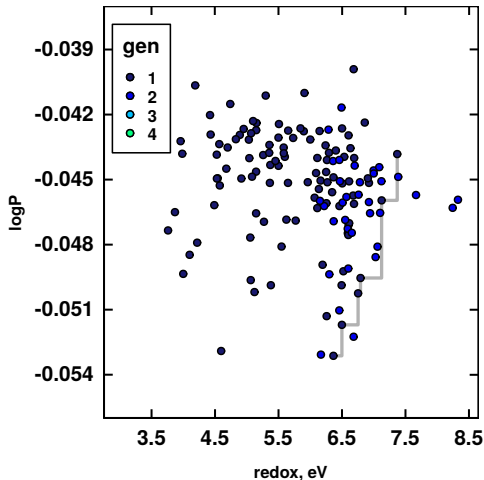
# DFT results

pareto front (gen 1)



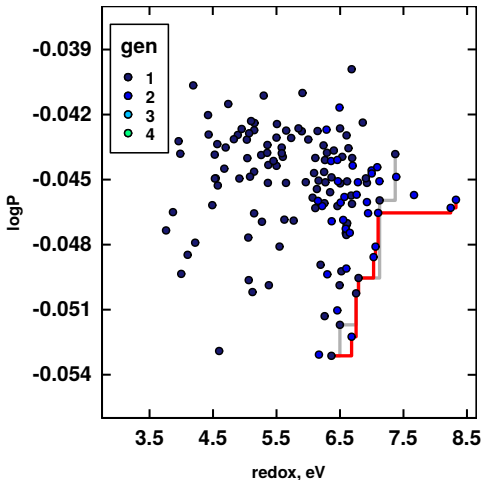
# DFT results

El points (gen 2: 34 complexes)



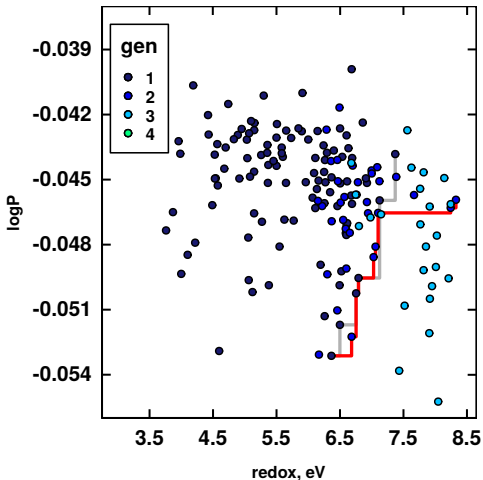
# DFT results

pareto front (gen 2)



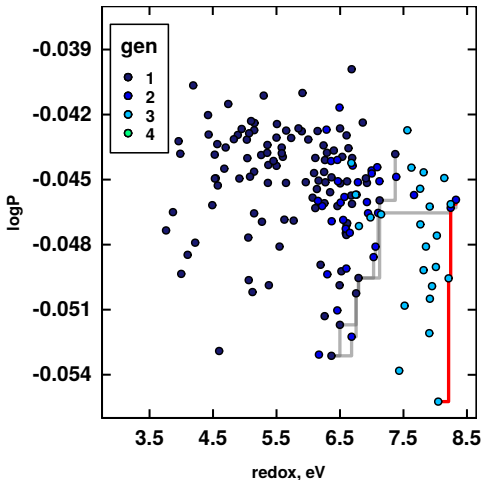
# DFT results

El points (gen 3: 24 complexes)



# DFT results

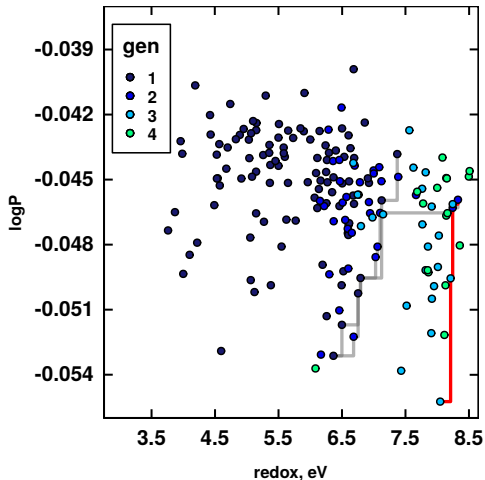
pareto front (gen 3)





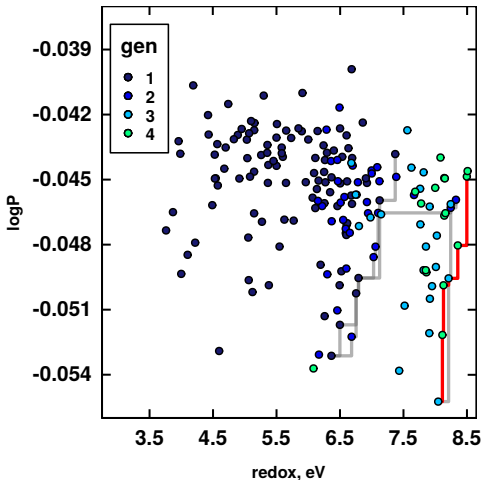
# DFT results

El points (gen 4: 15 complexes)

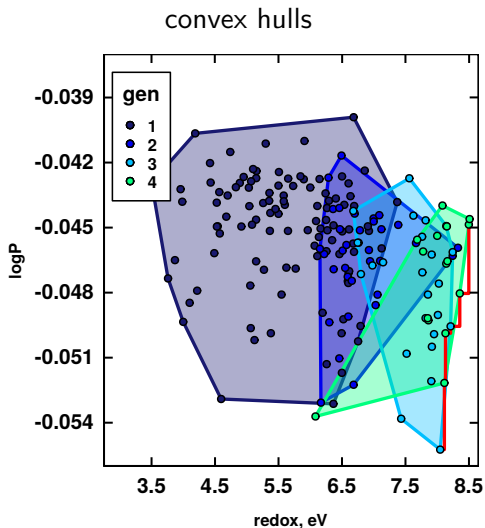


# DFT results

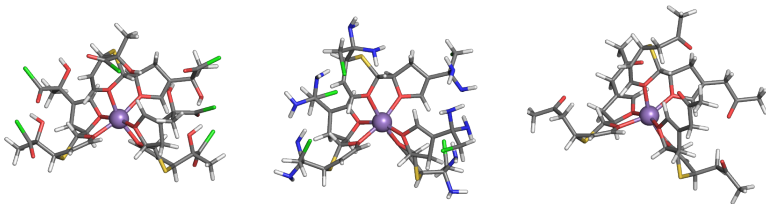
pareto front (gen 4)



# DFT results



# DFT results

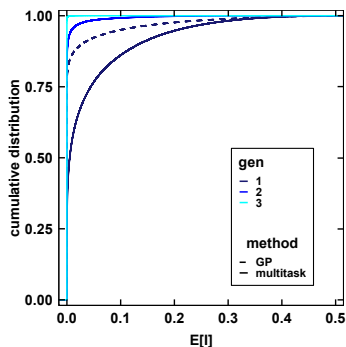


Janet, J.P., et al., *in preparation*.

# Case study conclusions

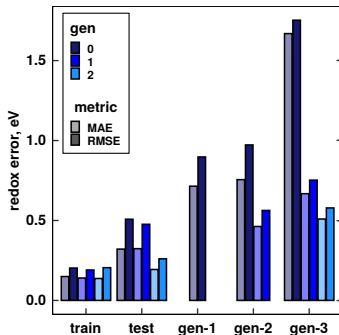
## Case study conclusions

- EI framework provides high resolution in the region of interest (c.f. maximum uncertainty), converges quickly



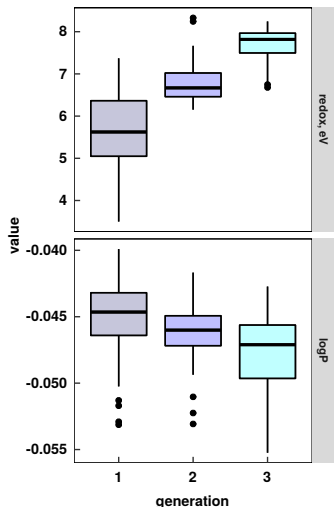
## Case study conclusions

- EI framework provides high resolution in the region of interest (c.f. maximum uncertainty), converges quickly
- We are able to identify fruitful regions from large chemical spaces based on few DFT evaluations



## Case study conclusions

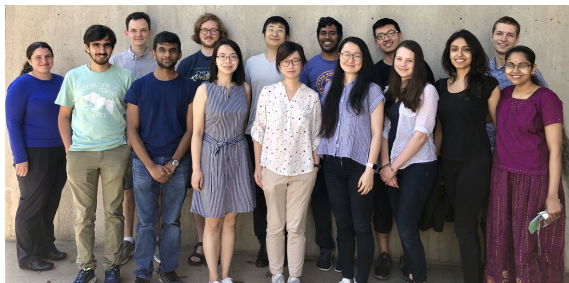
- EI framework provides high resolution in the region of interest (c.f. maximum uncertainty), converges quickly
- We are able to identify fruitful regions from large chemical spaces based on few DFT evaluations
- Multiobjective DFT optimization guided by data-driven method efficiency generates lead complexes





# Acknowledgments

Thanks to the Kulik group and funding partners:



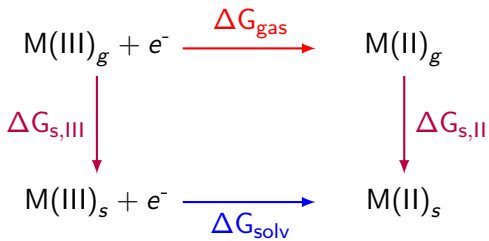
# Thermodynamic cycle



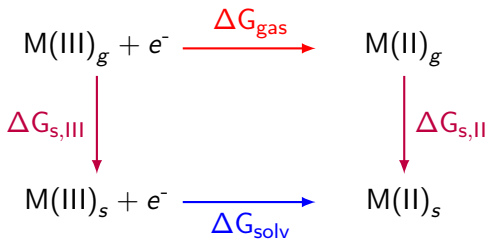
# Thermodynamic cycle



# Thermodynamic cycle

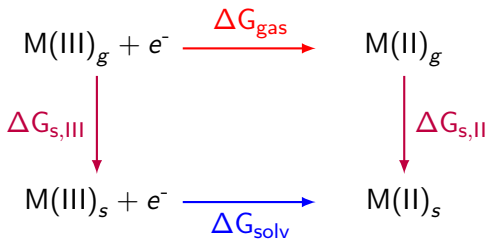


# Thermodynamic cycle



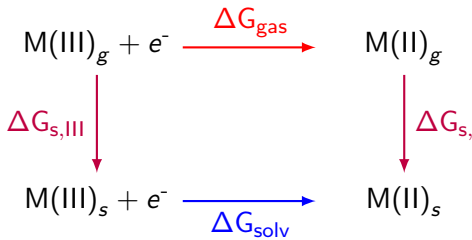
$$\Delta G_{\text{solv}} = \Delta G_{\text{gas}} + \Delta\Delta G_{\text{s}}$$

# Thermodynamic cycle

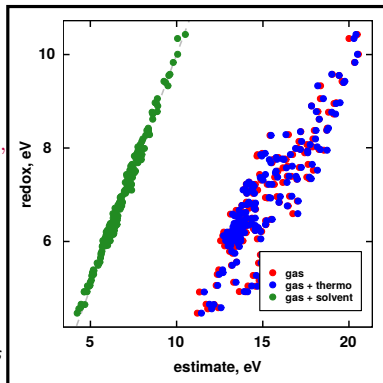


$$\begin{aligned} \Delta G_{\text{solv}} &= \Delta G_{\text{gas}} + \Delta\Delta G_{\text{s}} \\ &\approx \Delta E_{\text{gas,III-II}} + \Delta\Delta G_{\text{s}} \end{aligned}$$

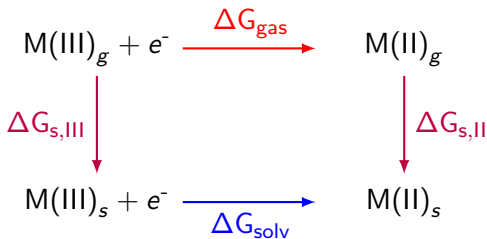
# Thermodynamic cycle



$$\begin{aligned} \Delta G_{\text{solv}} &= \Delta G_{\text{gas}} + \Delta\Delta G_{\text{s}} \\ &\approx \Delta E_{\text{gas,III-II}} + \Delta\Delta G_{\text{s}} \end{aligned}$$



# Thermodynamic cycle



$$\begin{aligned} \Delta G_{\text{solv}} &= \Delta G_{\text{gas}} + \Delta\Delta G_{\text{s}} \\ &\approx \Delta E_{\text{gas,III-II}} + \Delta\Delta G_{\text{s}} \end{aligned}$$

$$\log P \approx \log \frac{\Delta G_{\text{s,II,octanol}}}{\Delta G_{\text{s,II,water}}}$$

- Spin state taken as the M(II) ground state
- Need 3 geometry optimizations at minimum