

ML for inorganic molecular design: descriptors and similarity in transition metal chemical space

Jon Paul Janet¹ Heather Kulik ¹

¹Department of Chemical Engineering, Massachusetts Institute of Technology



255th ACS National Meeting, New Orleans

03.19.18

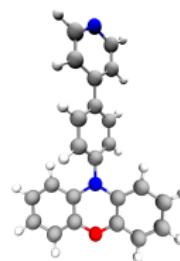
Data-driven molecular design

**Machine learning
is transforming
how we design
new materials...**

Data-driven molecular design

**Machine learning
is transforming
how we design
new materials...**

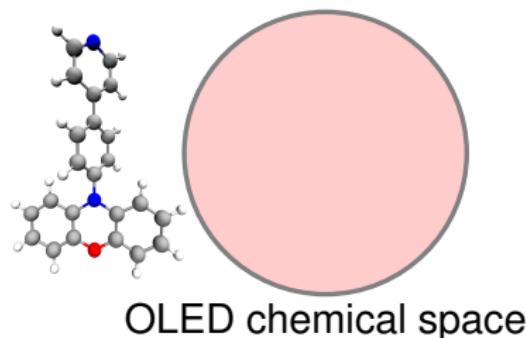
Gomez-Bombarelli, R. et al.. *Nat. Mater.*, 15(10):1120-1127, 2016.



Data-driven molecular design

**Machine learning
is transforming
how we design
new materials...**

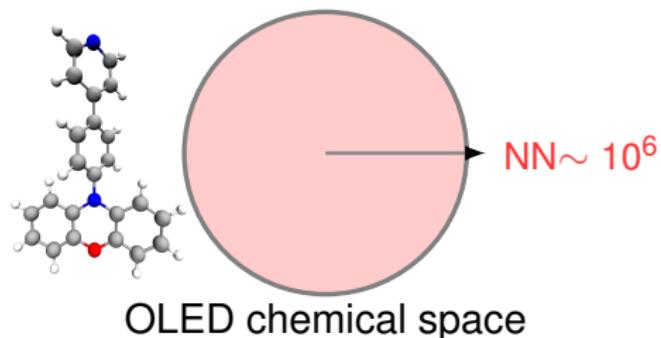
Gomez-Bombarelli, R. et al.. *Nat. Mater.*, 15(10):1120-1127, 2016.



Data-driven molecular design

**Machine learning
is transforming
how we design
new materials...**

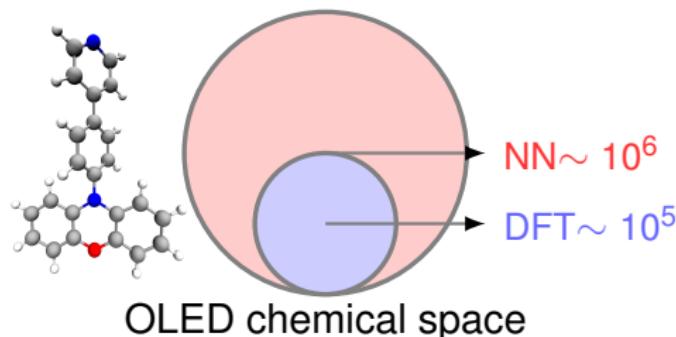
Gomez-Bombarelli, R. et al.. *Nat. Mater.*, 15(10):1120-1127, 2016.



Data-driven molecular design

**Machine learning
is transforming
how we design
new materials...**

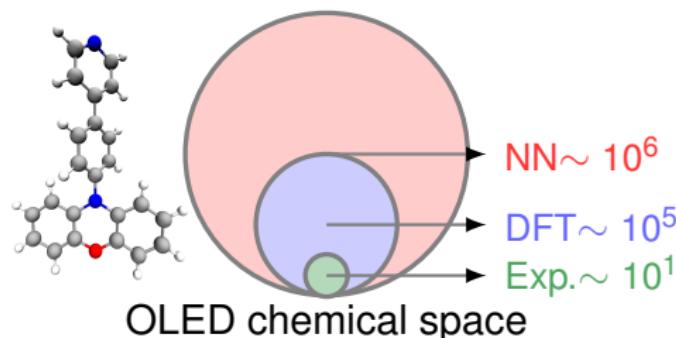
Gomez-Bombarelli, R. et al.. *Nat. Mater.*, 15(10):1120-1127, 2016.



Data-driven molecular design

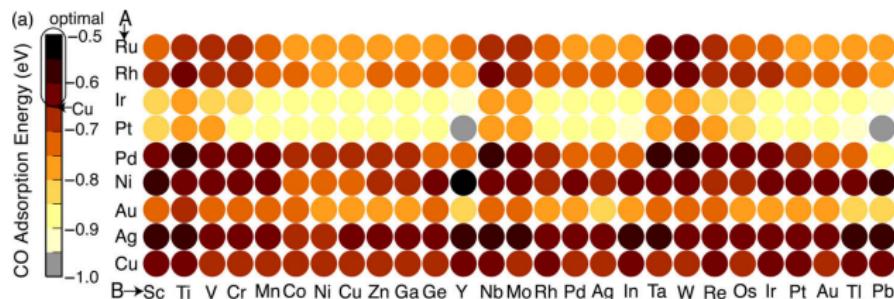
**Machine learning
is transforming
how we design
new materials...**

Gomez-Bombarelli, R. et al.. *Nat. Mater.*, 15(10):1120-1127, 2016.



Data-driven molecular design

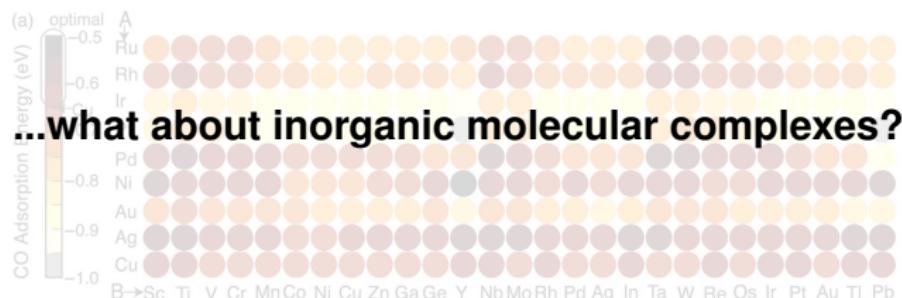
**Machine learning
is transforming
how we design
new materials...**



Ma, X. et al. J. Phys. Chem. Lett., (18):3528-3533, 2015.

Data-driven molecular design

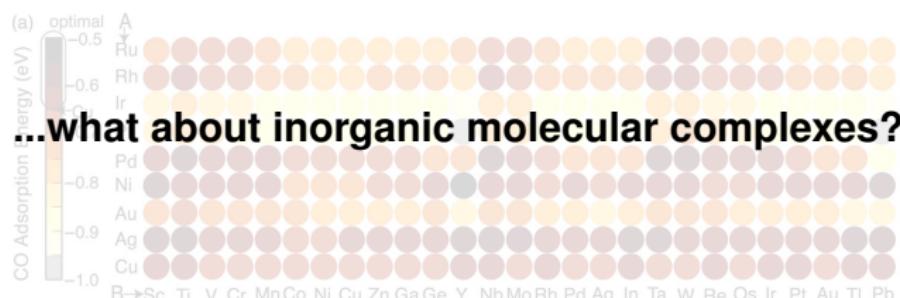
**Machine learning
is transforming
how we design
new materials...**



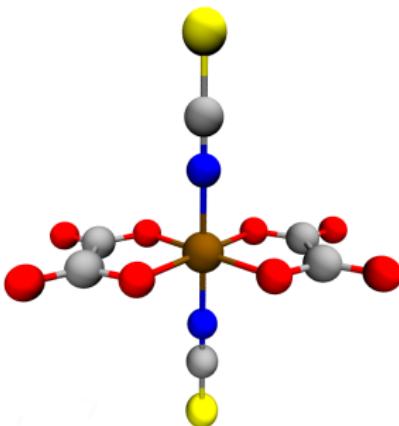
Ma, X. et al. J. Phys. Chem. Lett., (18):3528-3533, 2015.

Data-driven molecular design

**Machine learning
is transforming
how we design
new materials...**

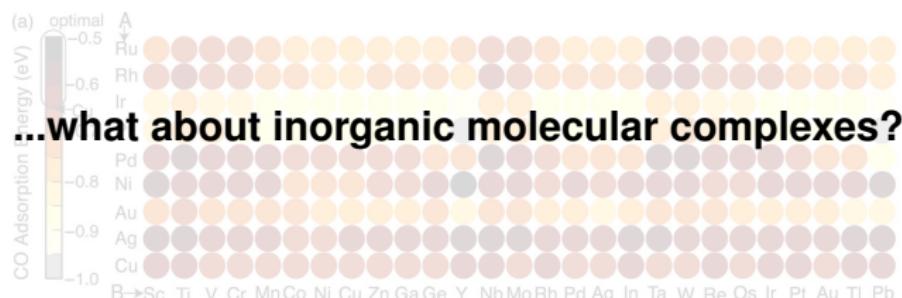


M. X. et al. J. Phys. Chem. Lett., (18):3528-3533, 2015.

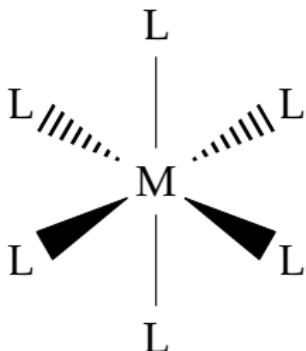


Data-driven molecular design

**Machine learning
is transforming
how we design
new materials...**

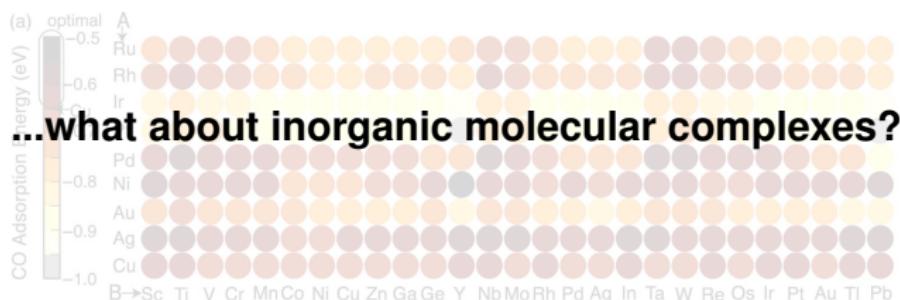


Ma, X. et al. J. Phys. Chem. Lett., (18):3528-3533, 2015.

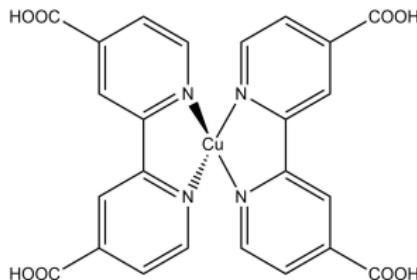
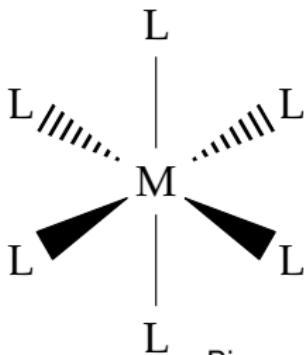


Data-driven molecular design

**Machine learning
is transforming
how we design
new materials...**



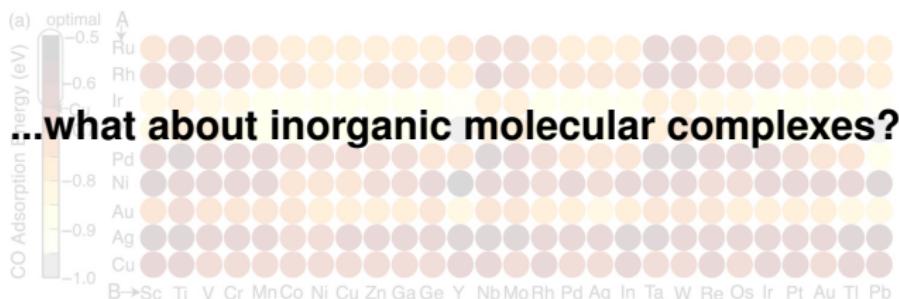
Ma, X. et al. J. Phys. Chem. Lett., (18):3528-3533, 2015.



Bignozzi, C. et al. *Coord. Chem. Rev.*, 257(9), 2013.

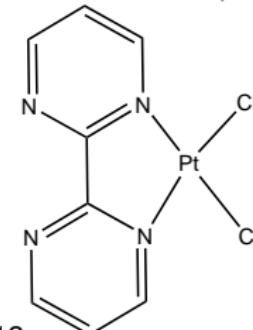
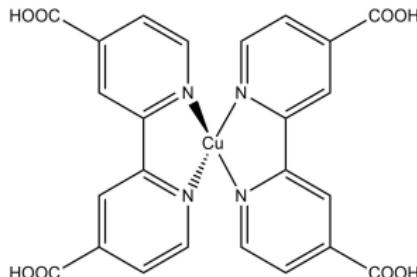
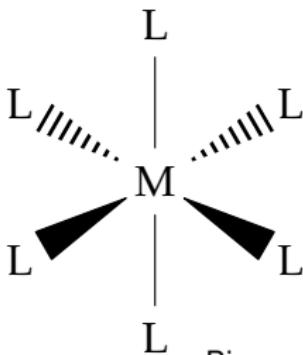
Data-driven molecular design

**Machine learning
is transforming
how we design
new materials...**



Ma, X. et al. J. Phys. Chem. Lett., (18):3528-3533, 2015.

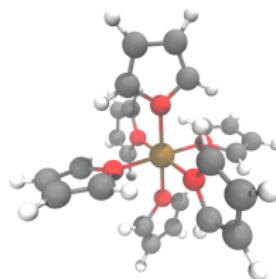
Periana, R. A. et al. *Science*, 280(5363), 1998.



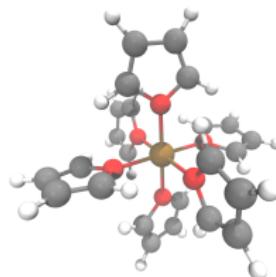
Bignozzi, C. et al. *Coord. Chem. Rev.*, 257(9), 2013.

Transition metal complexes

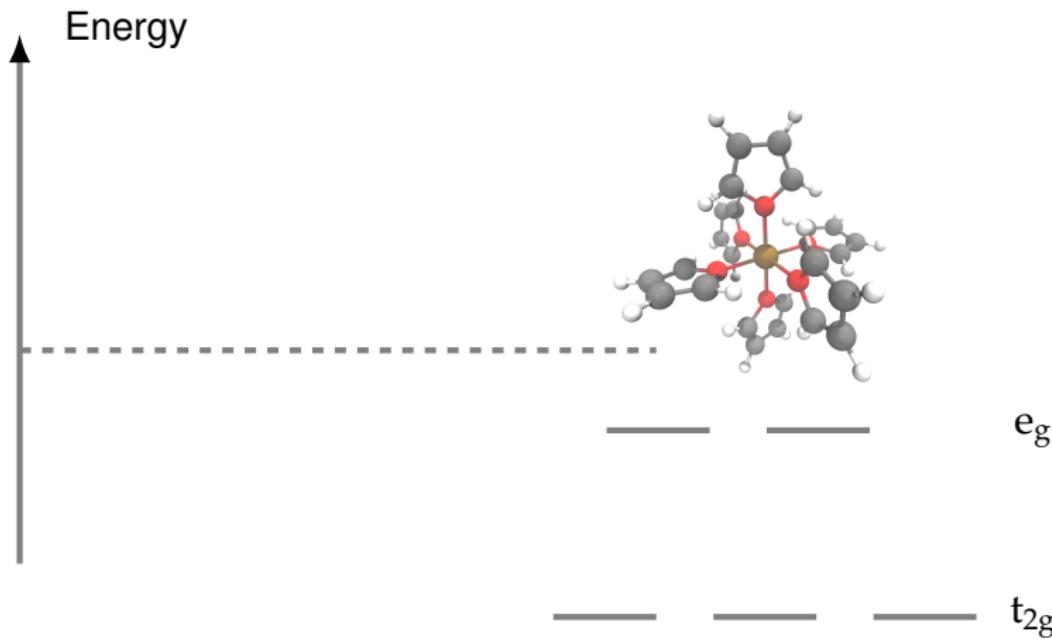
Transition metal complexes



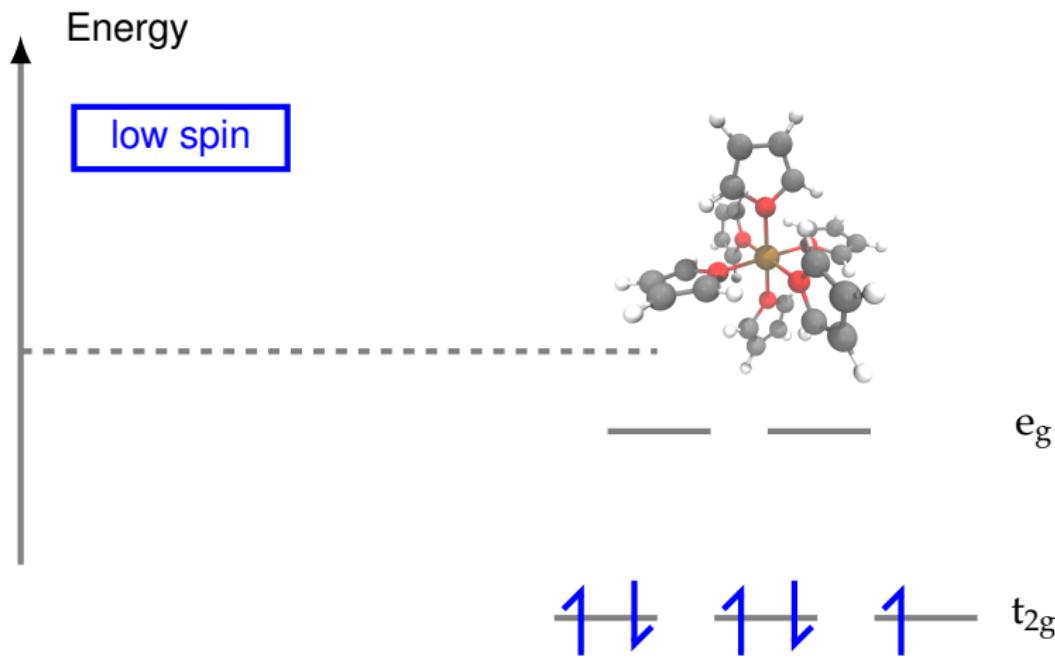
Transition metal complexes

 e_g  t_{2g}

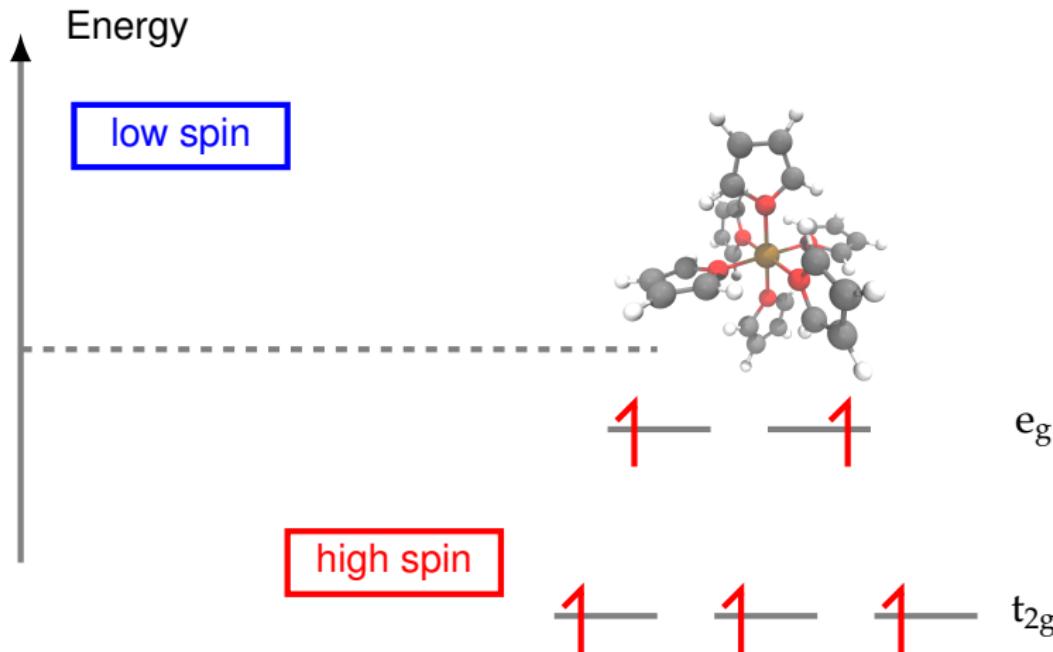
Transition metal complexes



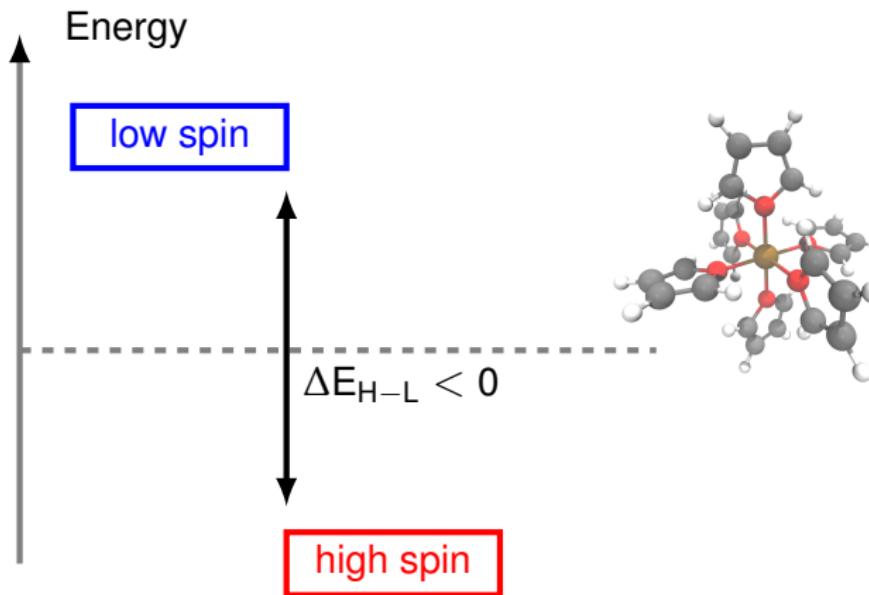
Transition metal complexes



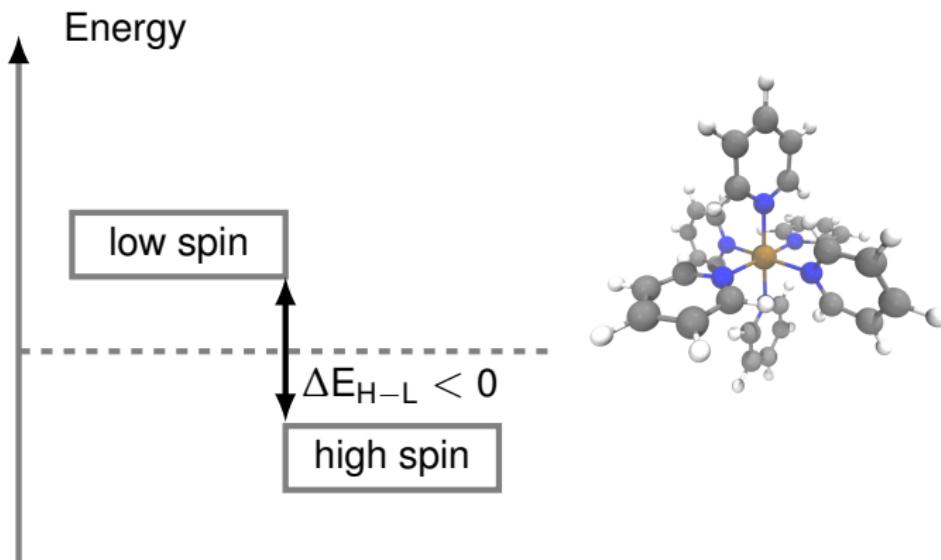
Transition metal complexes



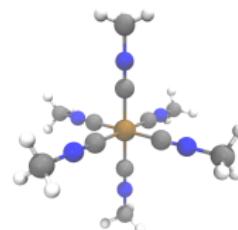
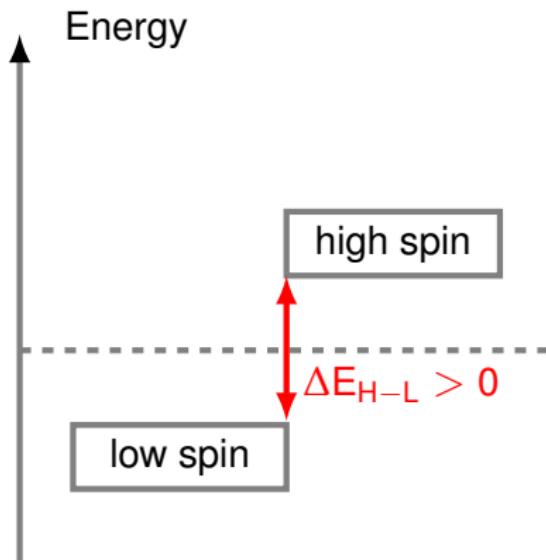
Transition metal complexes



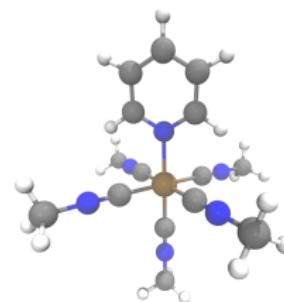
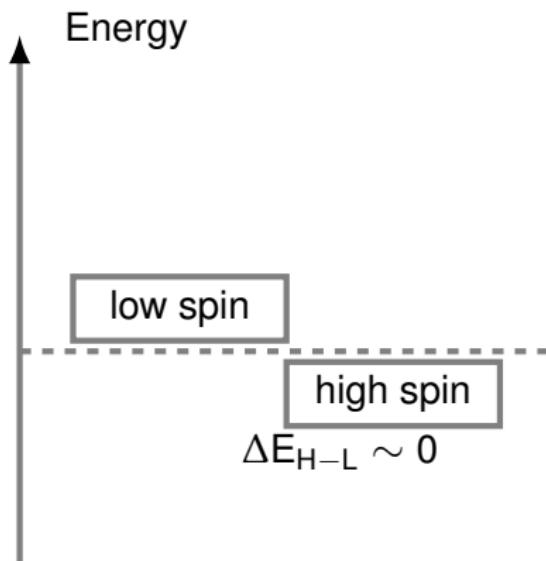
Transition metal complexes



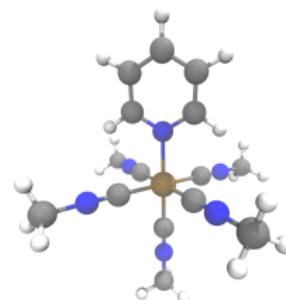
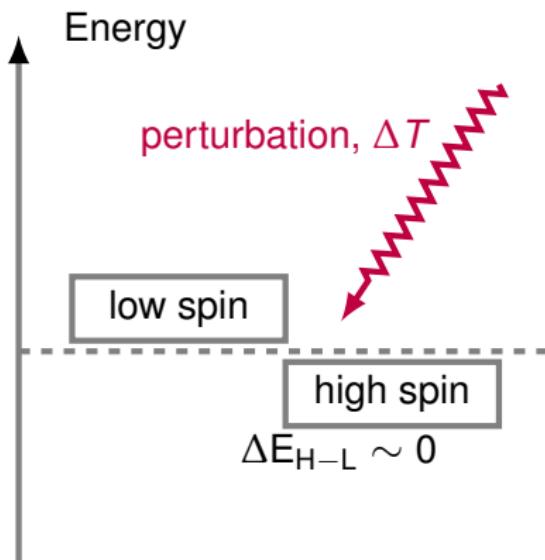
Transition metal complexes



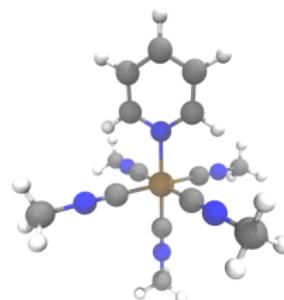
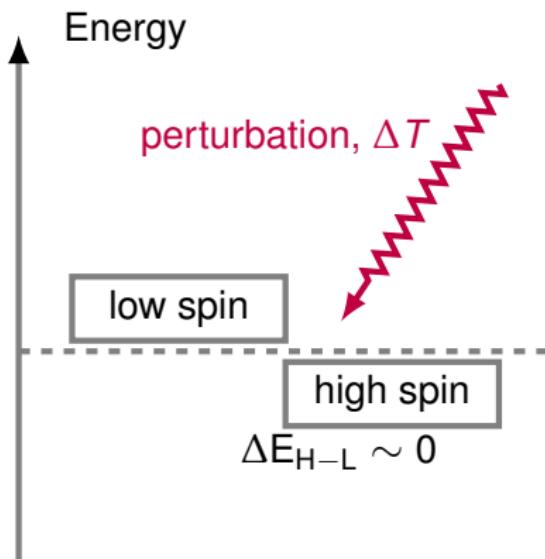
Transition metal complexes



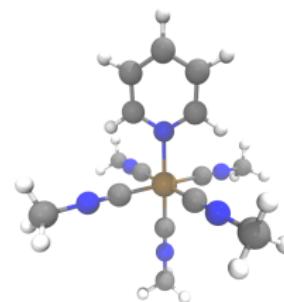
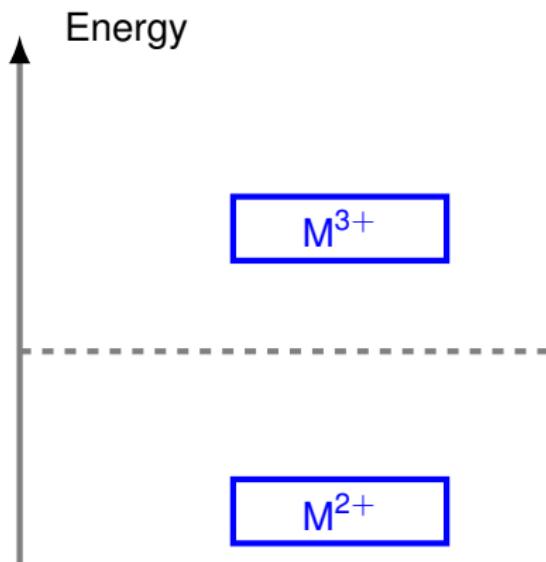
Transition metal complexes



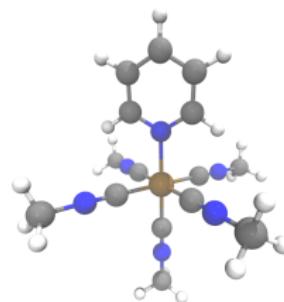
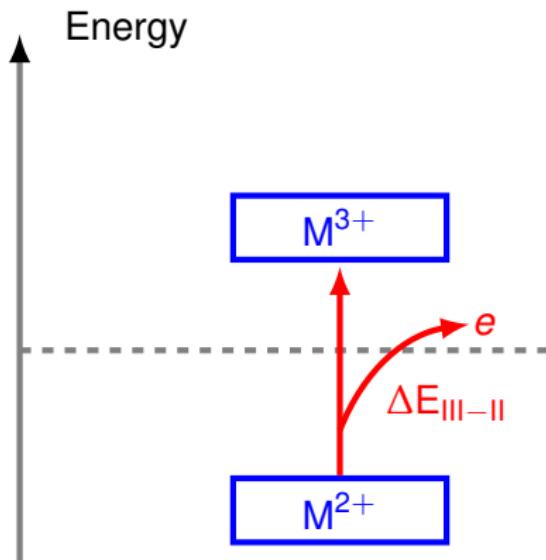
Transition metal complexes



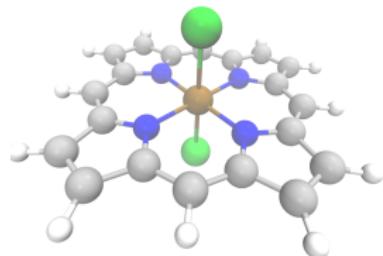
Transition metal complexes



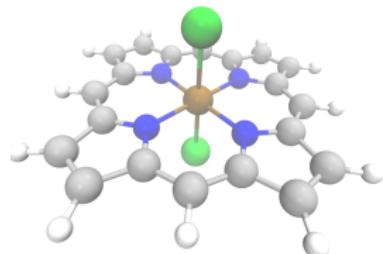
Transition metal complexes



How to estimate properties?

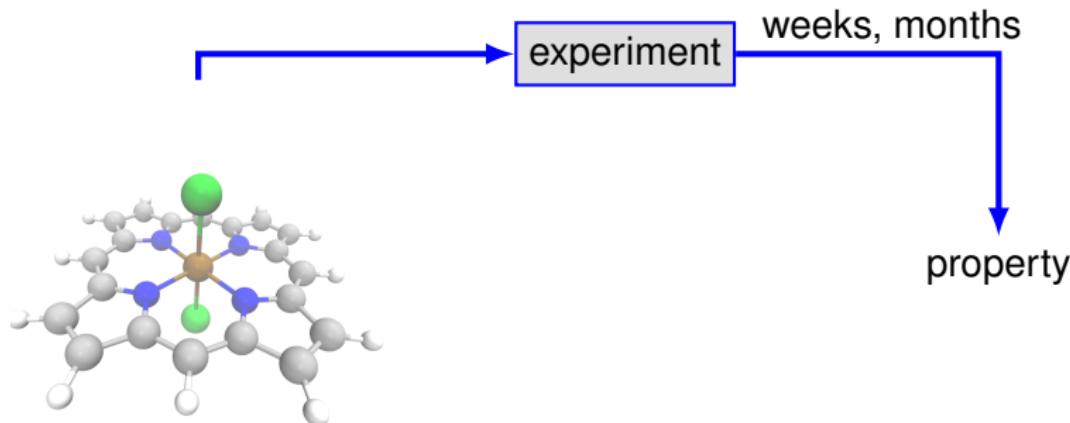


How to estimate properties?

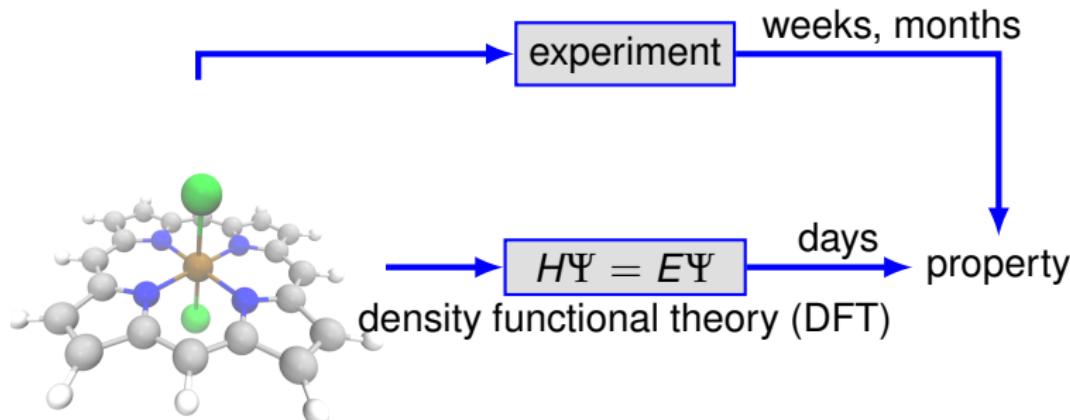


property

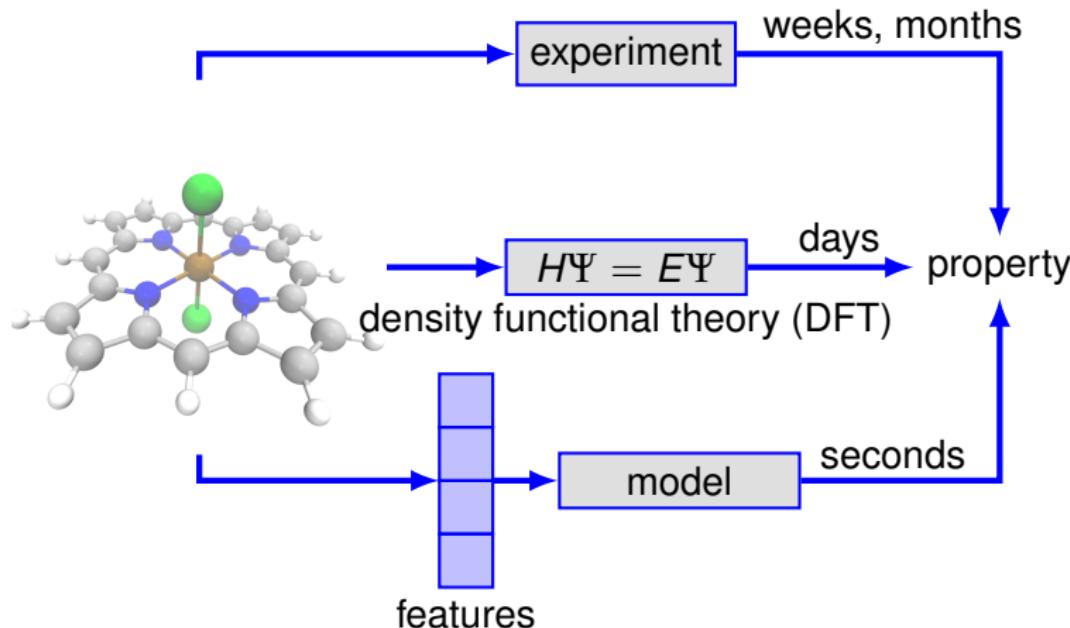
How to estimate properties?



How to estimate properties?

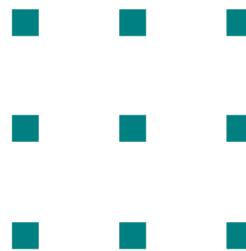


How to estimate properties?



Input space design

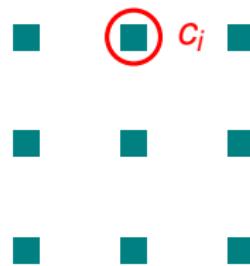
What would be the ideal feature space?



Chemical Space C_f

Input space design

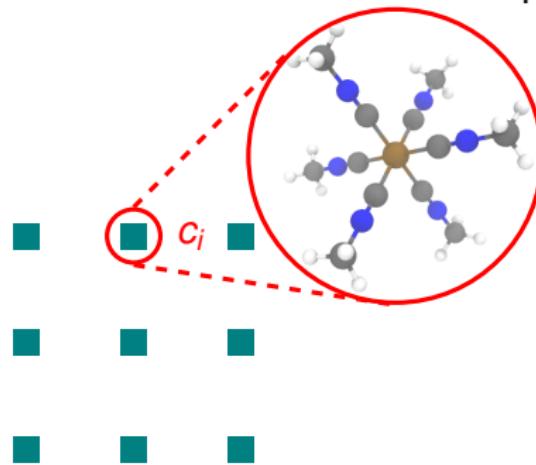
What would be the ideal feature space?



Chemical Space C_f

Input space design

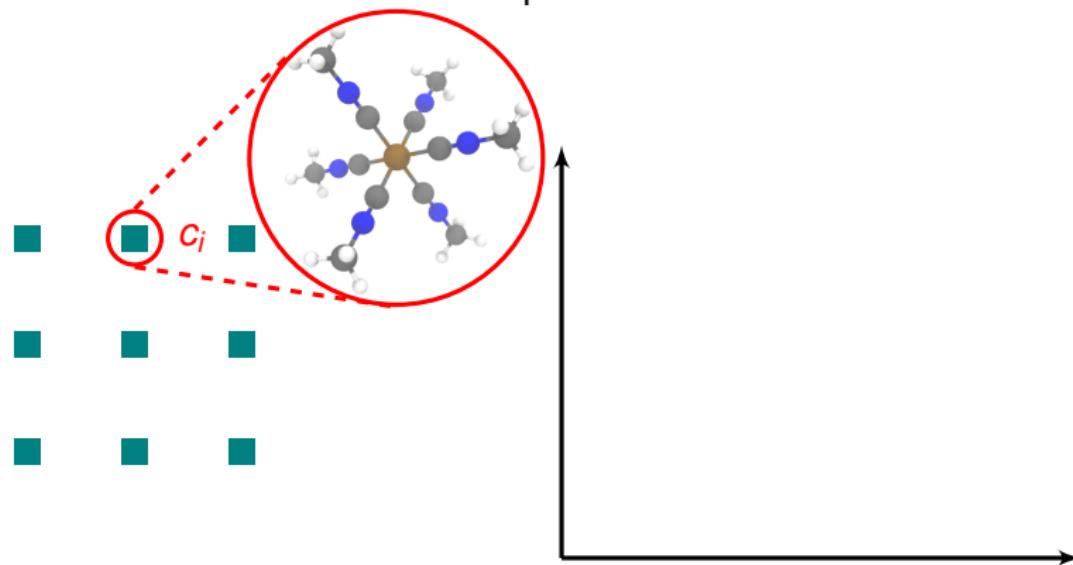
What would be the ideal feature space?



Chemical Space C_f

Input space design

What would be the ideal feature space?

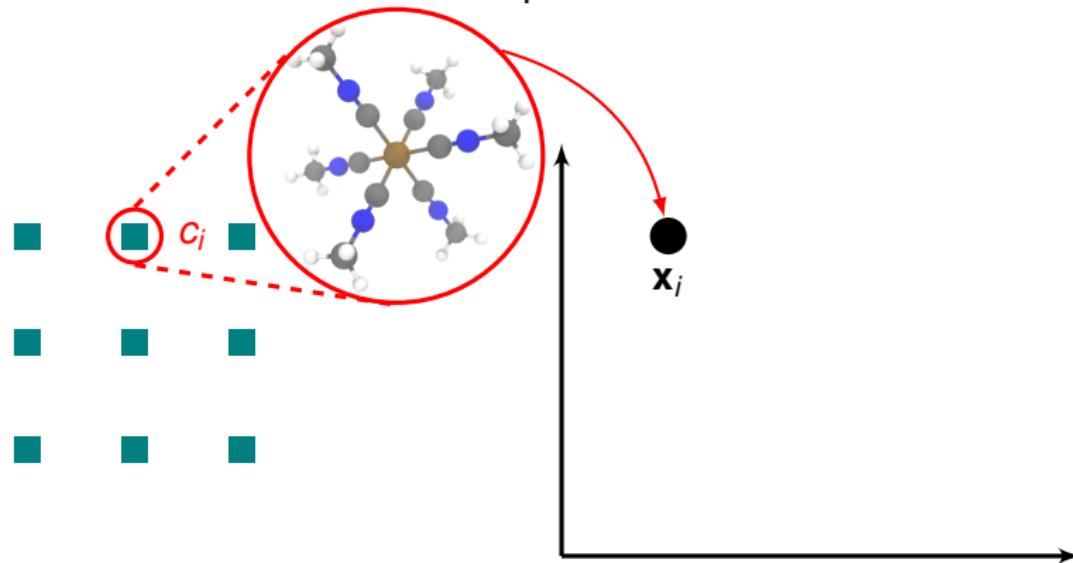


Chemical Space C_f

Descriptor Space $\mathcal{X} \subset \mathbb{R}^d$

Input space design

What would be the ideal feature space?

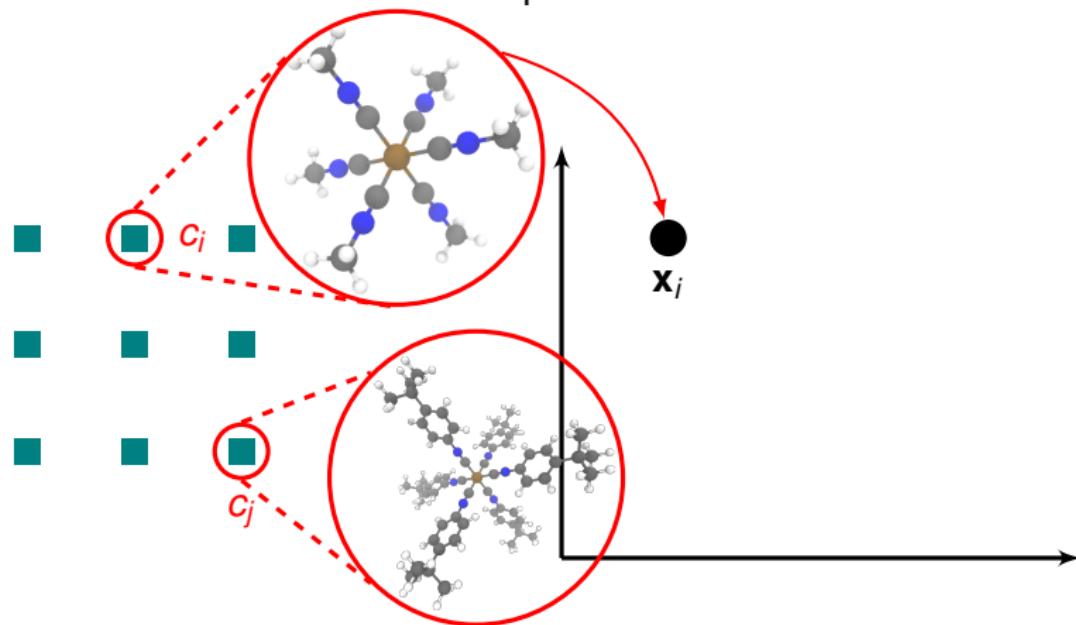


Chemical Space C_f

Descriptor Space $\mathcal{X} \subset \mathbb{R}^d$

Input space design

What would be the ideal feature space?

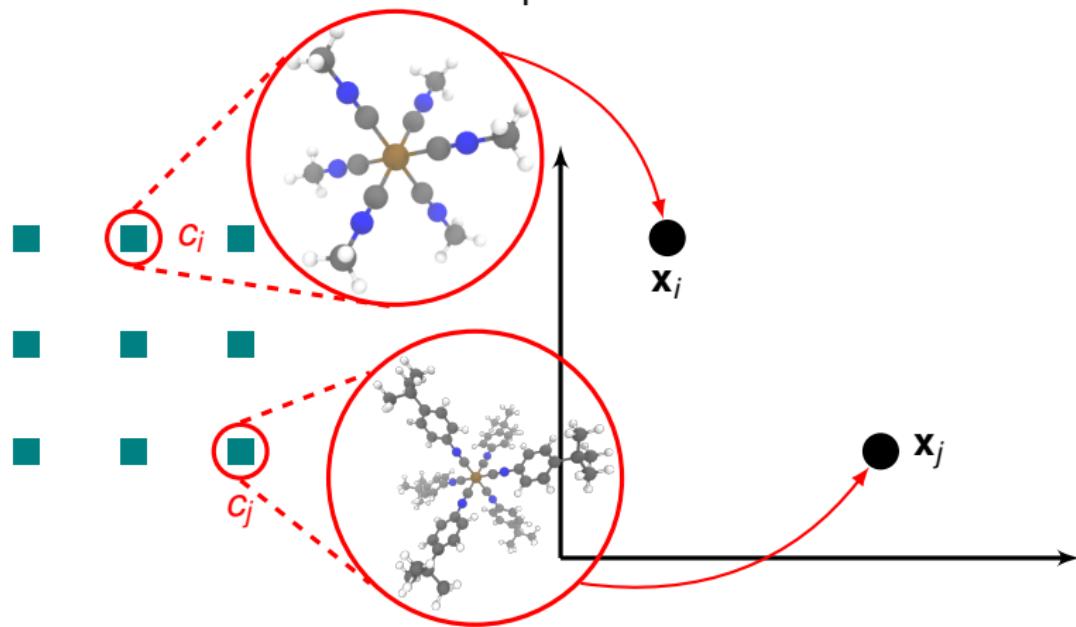


Chemical Space C_f

Descriptor Space $\mathcal{X} \subset \mathbb{R}^d$

Input space design

What would be the ideal feature space?

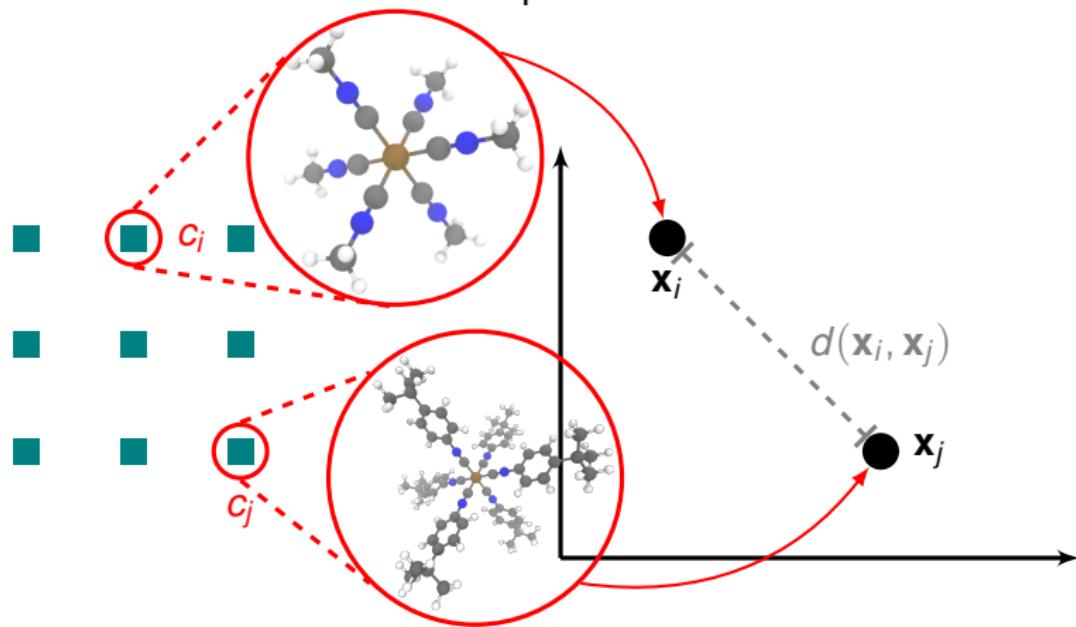


Chemical Space C_f

Descriptor Space $\mathcal{X} \subset \mathbb{R}^d$

Input space design

What would be the ideal feature space?

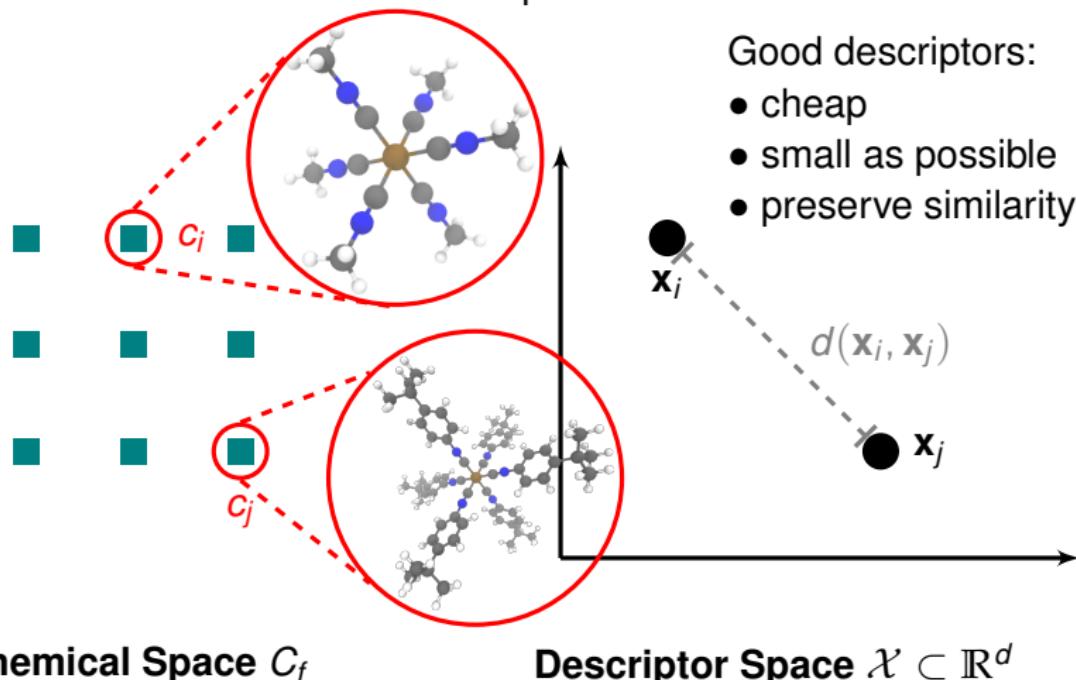


Chemical Space C_f

Descriptor Space $\mathcal{X} \subset \mathbb{R}^d$

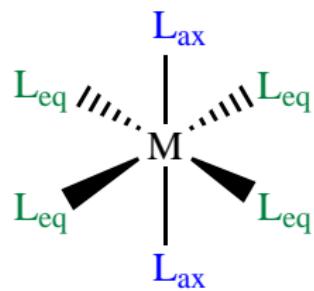
Input space design

What would be the ideal feature space?



Data for spin splitting

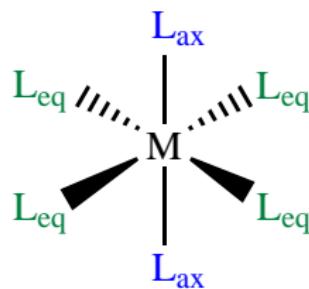
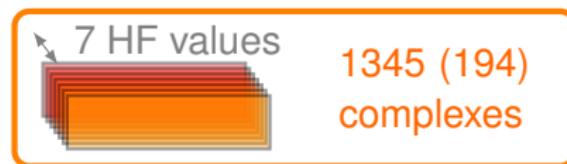
Data for octahedral complexes¹:



¹Janet, J.P., and Kulik, H.J. *Chem. Sci.*, 2017, 8, 5137-5152.

Data for spin splitting

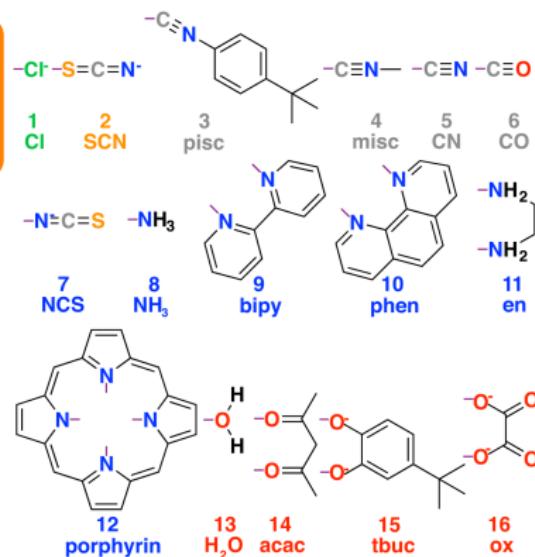
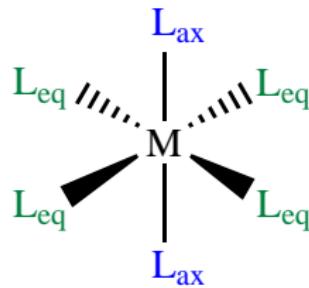
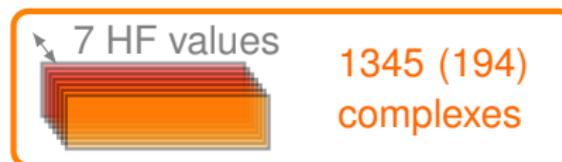
Data for octahedral complexes¹:



¹Janet, J.P., and Kulik, H.J. *Chem. Sci.*, 2017, 8, 5137-5152.

Data for spin splitting

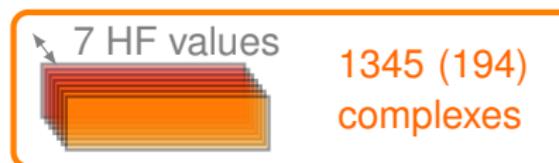
Data for octahedral complexes¹:



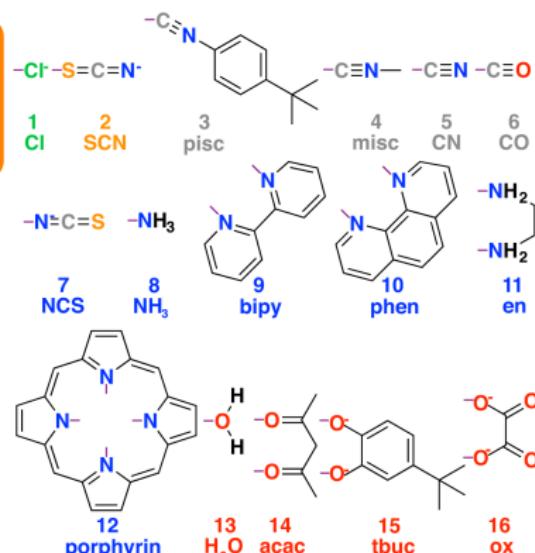
¹Janet, J.P., and Kulik, H.J. *Chem. Sci.*, 2017, 8, 5137-5152.

Data for spin splitting

Data for octahedral complexes¹:



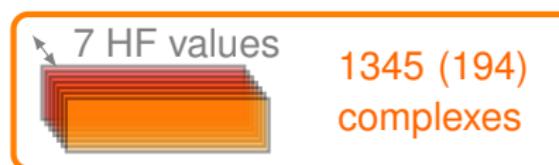
B3LYP-like DFT
HF exchange in 0-30%
gas phase optimizaton
LANL2DZ/6-31G*
high- and low-spin
M(II)/(III)



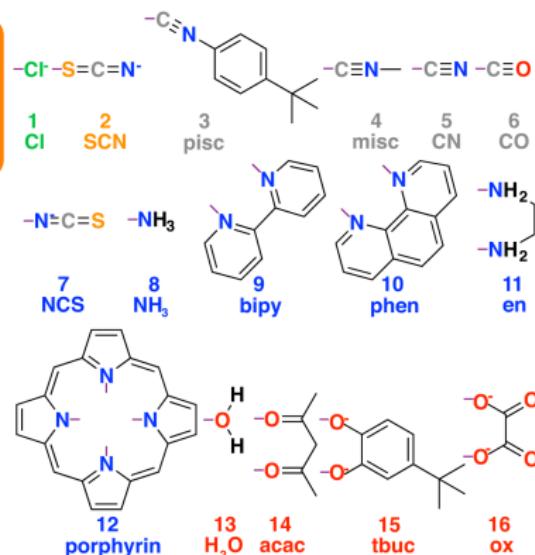
¹Janet, J.P., and Kulik, H.J. *Chem. Sci.*, 2017, 8, 5137-5152.

Data for spin splitting

Data for octahedral complexes¹:



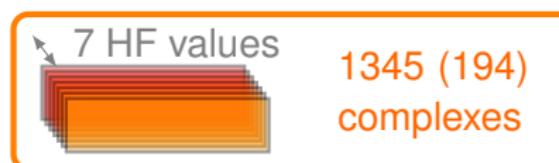
Coulomb matrix eigenspectrum (CM-ES) descriptor & kernel ridge regression (KRR)



¹Janet, J.P., and Kulik, H.J. *Chem. Sci.*, 2017, 8, 5137-5152.

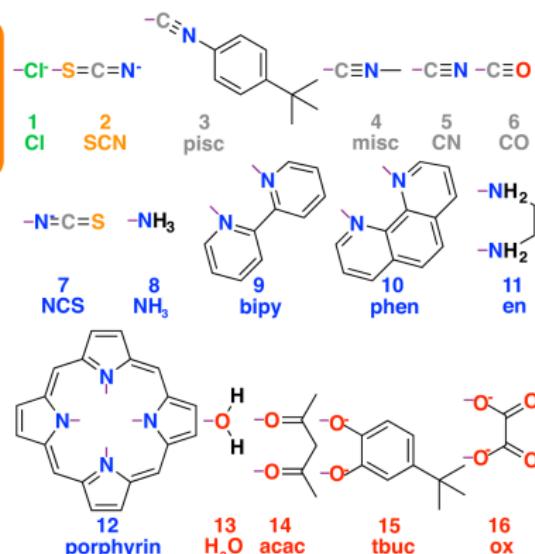
Data for spin splitting

Data for octahedral complexes¹:



Coulomb matrix eigenspectrum (CM-ES) descriptor & kernel ridge regression (KRR)

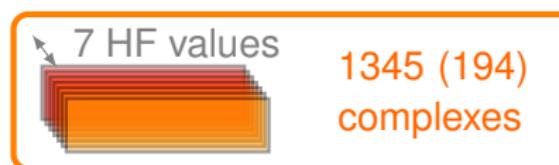
	ΔE_{H-L}	RMSE
CM-ES	19.2	kcal/mol



¹Janet, J.P., and Kulik, H.J. *Chem. Sci.*, 2017, 8, 5137-5152.

Data for spin splitting

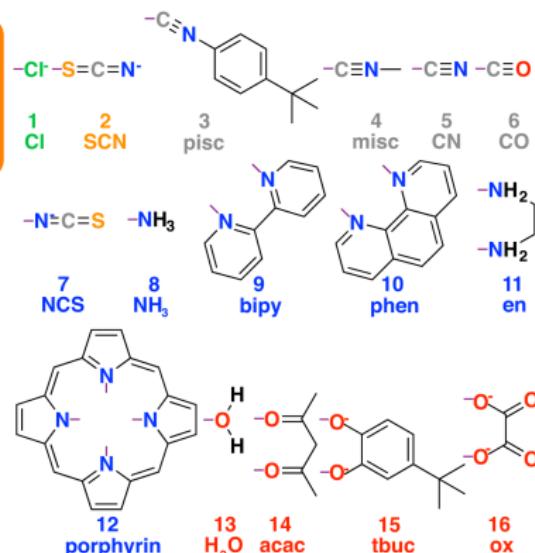
Data for octahedral complexes¹:



Coulomb matrix eigenspectrum (CM-ES) descriptor & kernel ridge regression (KRR)

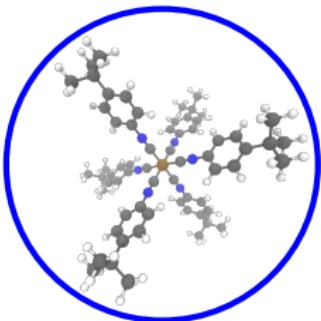
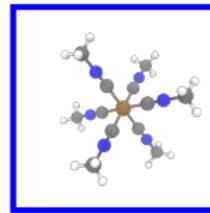
	ΔE_{H-L}	RMSE
CM-ES	19.2	kcal/mol

Why?

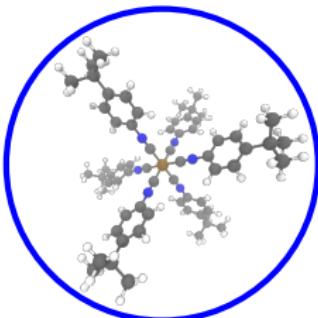


¹Janet, J.P., and Kulik, H.J. *Chem. Sci.*, 2017, 8, 5137-5152.

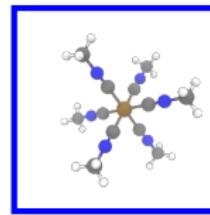
A tale of two complexes

 $\text{Fe}[\text{pisc}]_6^{3+}$  $\text{Fe}[\text{misc}]_6^{3+}$ 

A tale of two complexes

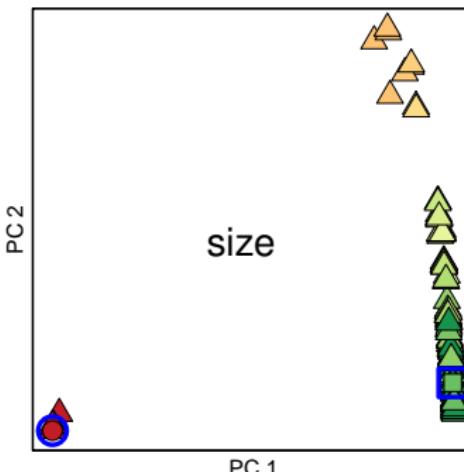
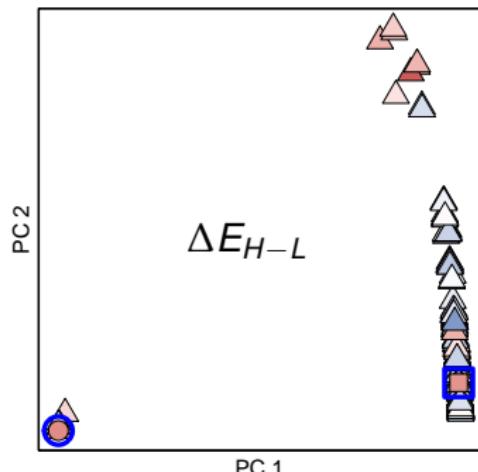
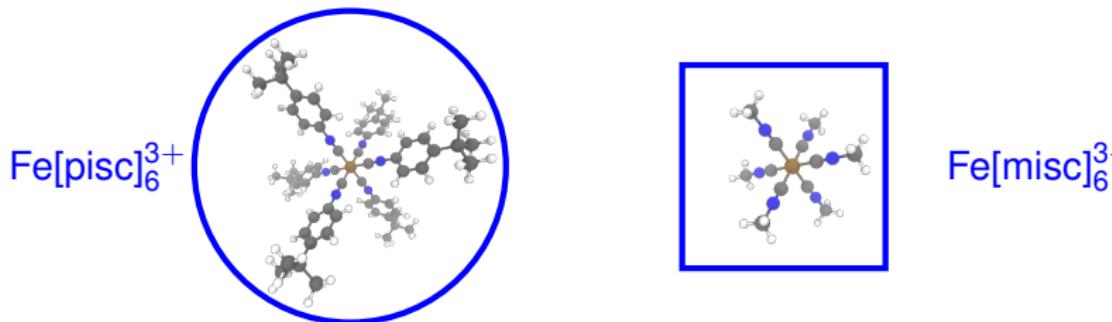
 $\text{Fe}[\text{pisc}]_6^{3+}$ 

$$\Delta E_{\text{H-L}} = 37.7 \text{ kcal/mol}$$

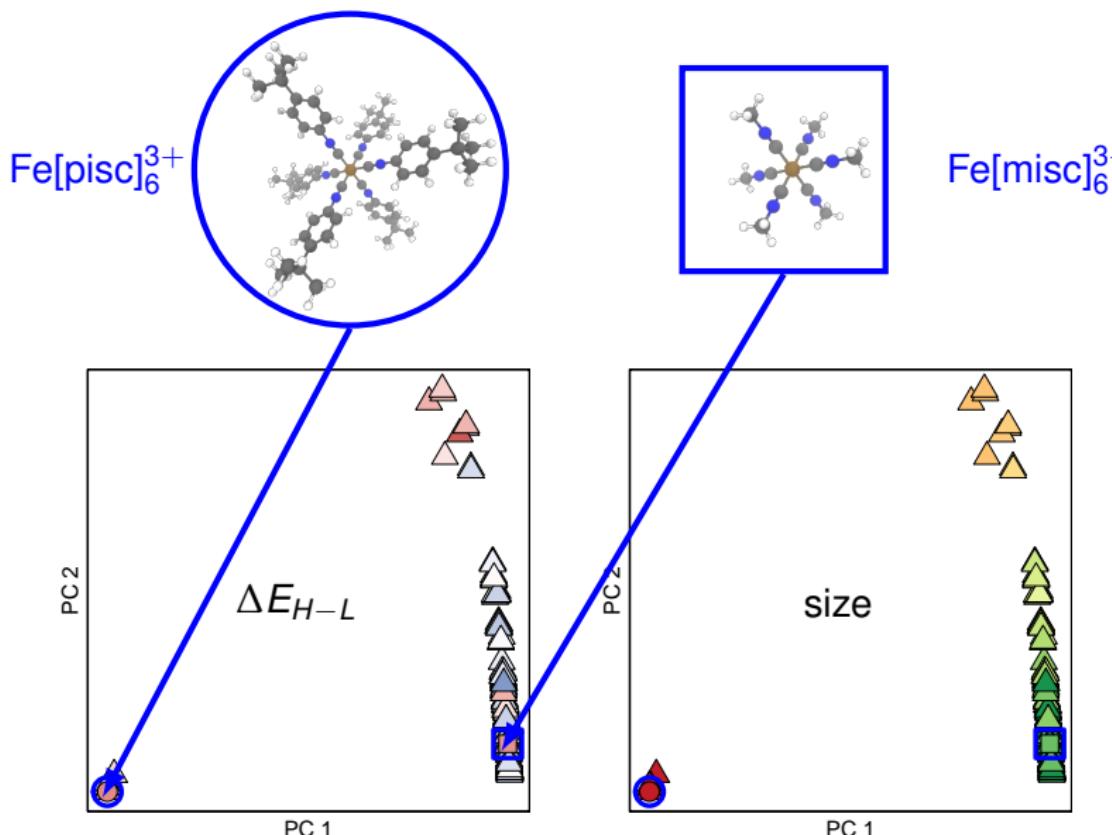
 $\text{Fe}[\text{misc}]_6^{3+}$ 

$$\Delta E_{\text{H-L}} = 40.7 \text{ kcal/mol}$$

A tale of two complexes



A tale of two complexes



MCDL-25

mixed continuous discrete lcoal (MCDL)

MCDL-25

mixed continuous discrete Icoal (MCDL)
metal

properties

identity

Fe(II)

oxidation state

MCDL-25

mixed continuous discrete local (MCDL)

metal local ligand

properties

properties

identity

Fe(II)

oxidation state

$\chi = 3.44$



$\chi = 2.55$

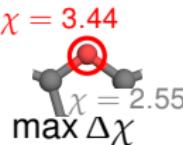
$\max \Delta \chi$

MCDL-25

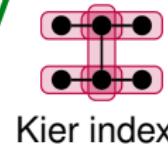
mixed continuous metal properties

identity
Fe(II)
oxidation state

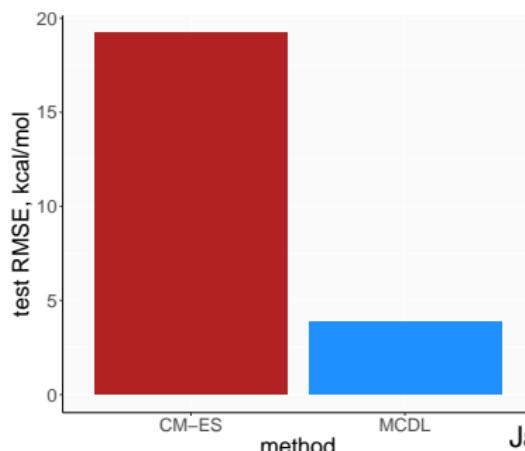
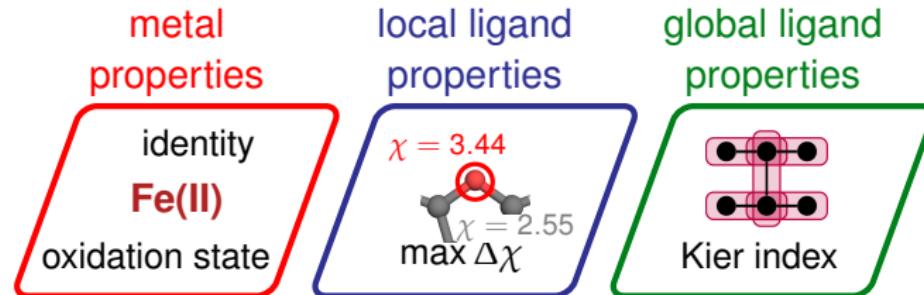
continuous local ligand properties



discrete global ligand properties

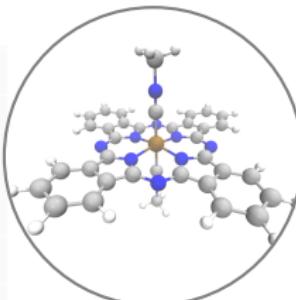
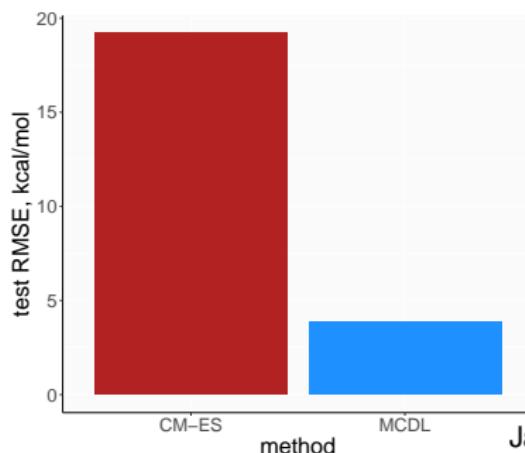
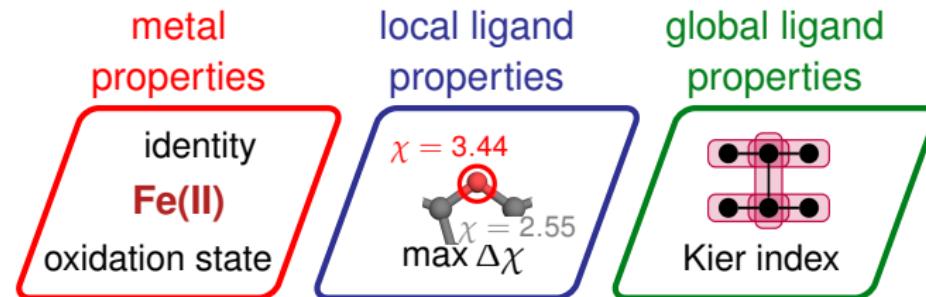


MCDL-25

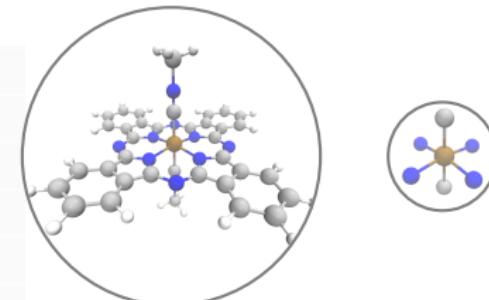
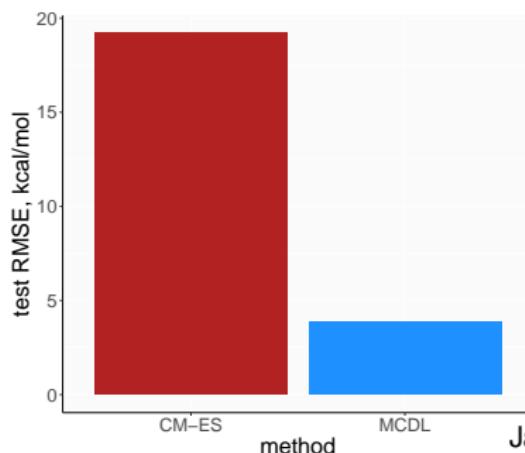
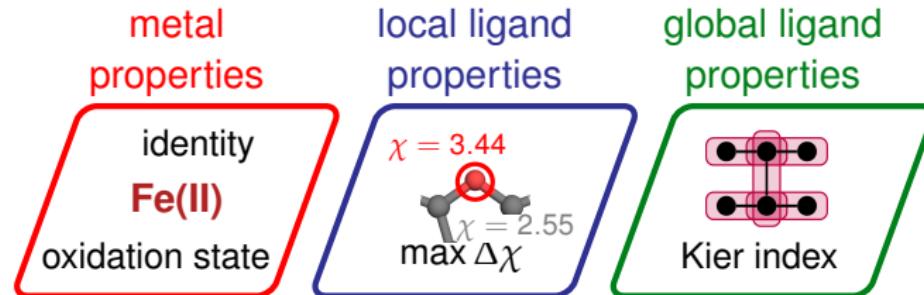


Janet, J.P., and Kulik, H.J. *Chem. Sci.*, 2017, 8, 5137-5152.

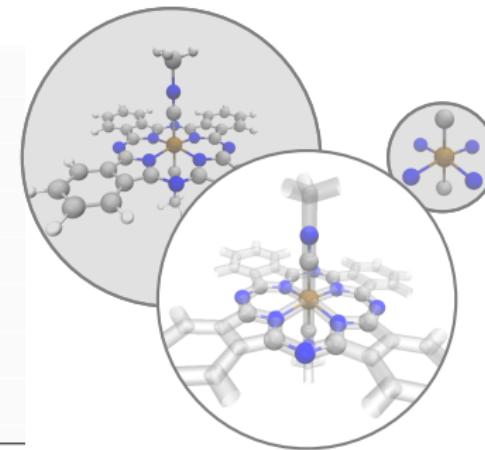
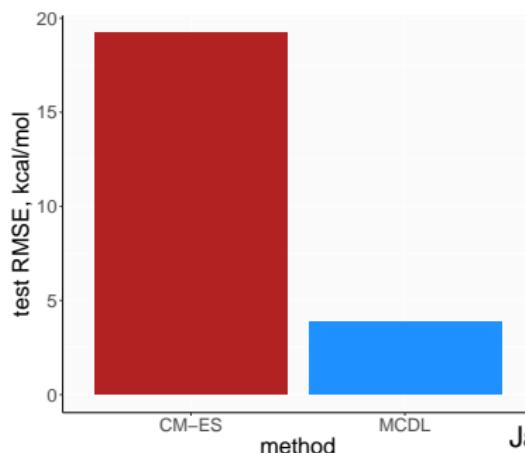
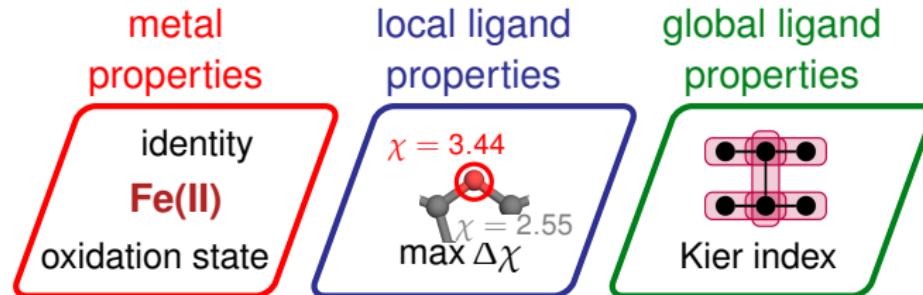
MCDL-25

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

MCDL-25

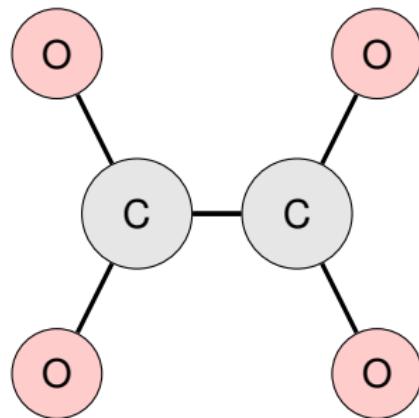


MCDL-25

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

Extensible, continuous descriptors - RACs

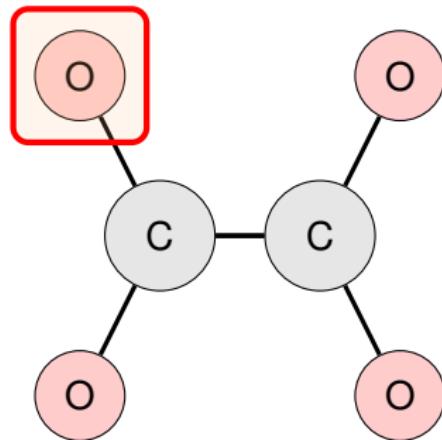
Based on autocorrelations²



²Broto, P., Moreau, G. and Vandycke, C. *Eur. J. Med. Chem.*, 19(1):71-78, 1984.

Extensible, continuous descriptors - RACs

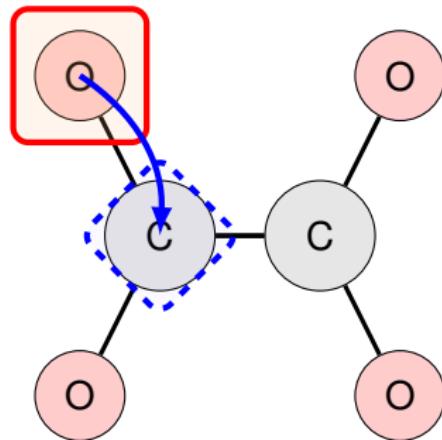
Based on autocorrelations²



²Broto, P., Moreau, G. and Vandycke, C. *Eur. J. Med. Chem.*, 19(1):71-78, 1984.

Extensible, continuous descriptors - RACs

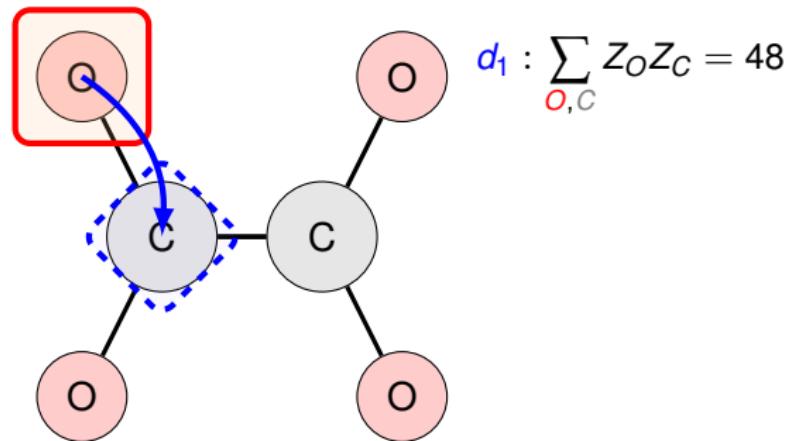
Based on autocorrelations²



²Broto, P., Moreau, G. and Vandycke, C. *Eur. J. Med. Chem.*, 19(1):71-78, 1984.

Extensible, continuous descriptors - RACs

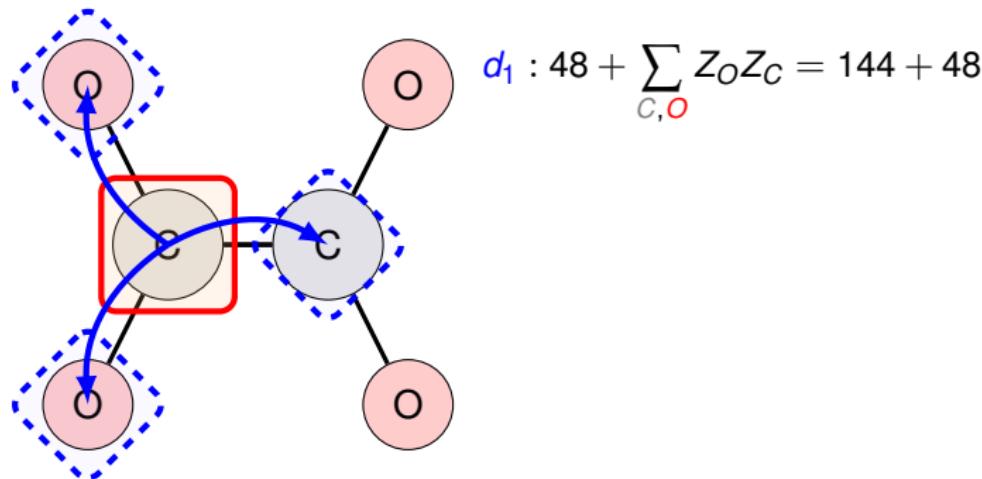
Based on autocorrelations²



²Broto, P., Moreau, G. and Vandycke, C. *Eur. J. Med. Chem.*, 19(1):71-78, 1984.

Extensible, continuous descriptors - RACs

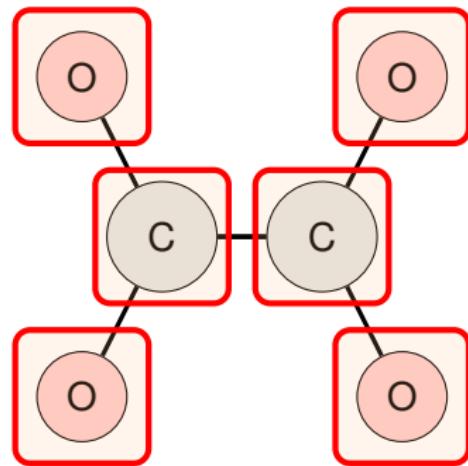
Based on autocorrelations²



²Broto, P., Moreau, G. and Vandycke, C. *Eur. J. Med. Chem.*, 19(1):71-78, 1984.

Extensible, continuous descriptors - RACs

Based on autocorrelations²

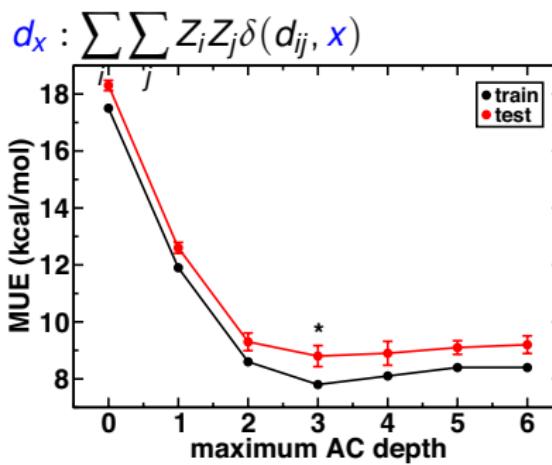
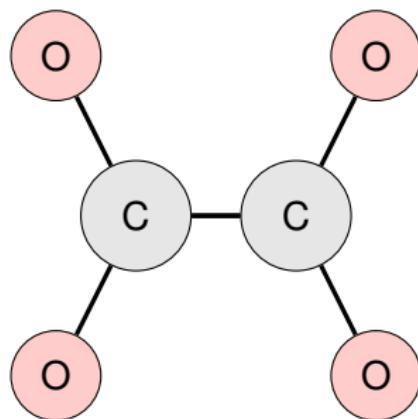


$$d_1 : \sum_i \sum_j Z_i Z_j \delta(d_{i,j}, 1)$$

²Broto, P., Moreau, G. and Vandycke, C. *Eur. J. Med. Chem.*, 19(1):71-78, 1984.

Extensible, continuous descriptors - RACs

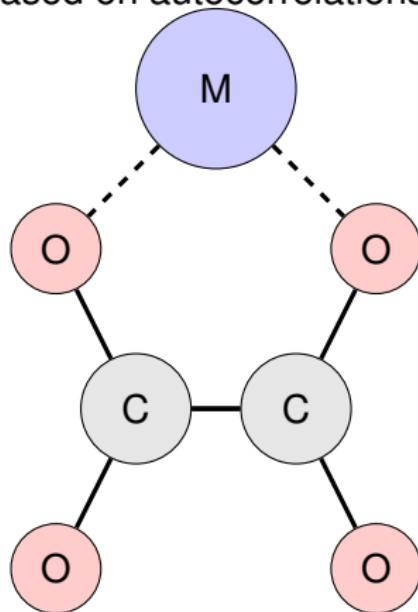
Based on autocorrelations²



²Broto, P., Moreau, G. and Vandycke, C. *Eur. J. Med. Chem.*, 19(1):71-78, 1984.

Extensible, continuous descriptors - RACs

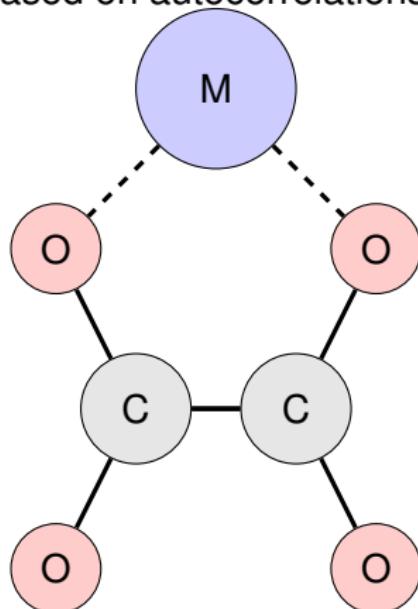
Based on autocorrelations



How to adapt to TM complexes?

Extensible, continuous descriptors - RACs

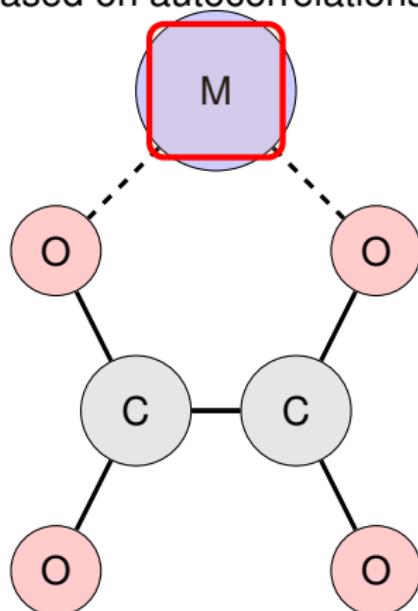
Based on autocorrelations



How to adapt to TM complexes?
restrict the scope to focus on
near-metal atoms

Extensible, continuous descriptors - RACs

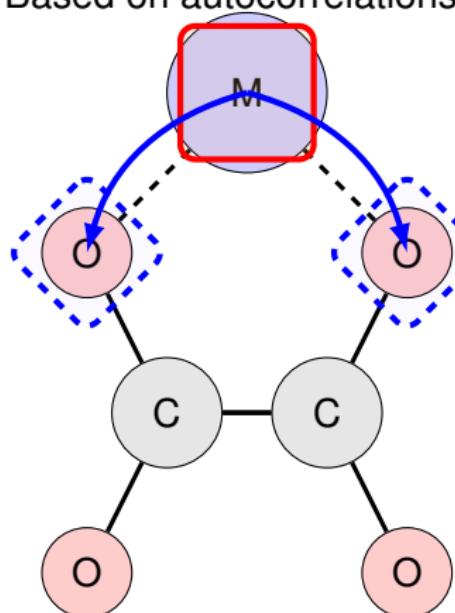
Based on autocorrelations



How to adapt to TM complexes?
restrict the scope to focus on
near-metal atoms

Extensible, continuous descriptors - RACs

Based on autocorrelations

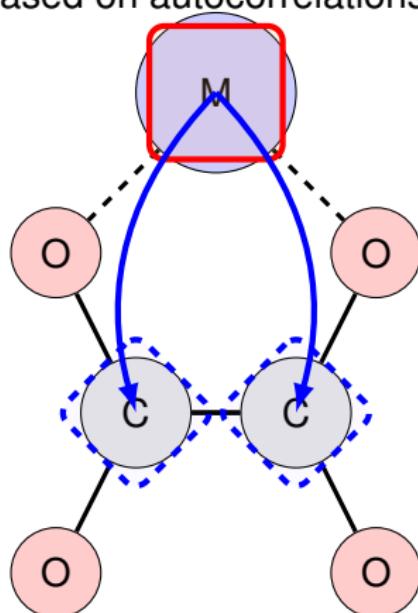


How to adapt to TM complexes?
restrict the scope to focus on
near-metal atoms

$$d_1 : \sum_{M,O} Z_M Z_O$$

Extensible, continuous descriptors - RACs

Based on autocorrelations

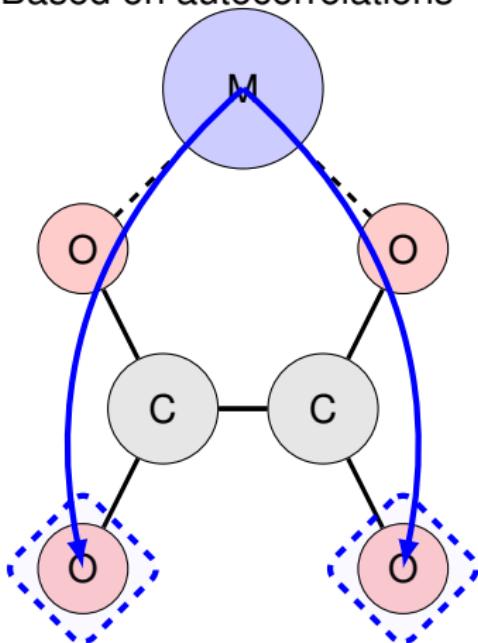


How to adapt to TM complexes?
restrict the scope to focus on
near-metal atoms

$$d_2 : \sum_{M,C} Z_M Z_C$$

Extensible, continuous descriptors - RACs

Based on autocorrelations

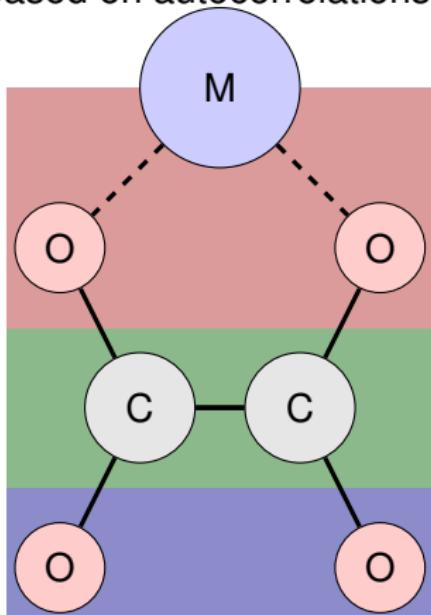


How to adapt to TM complexes?
restrict the scope to focus on
near-metal atoms

$$d_3 : \sum_{M,O} Z_M Z_O$$

Extensible, continuous descriptors - RACs

Based on autocorrelations

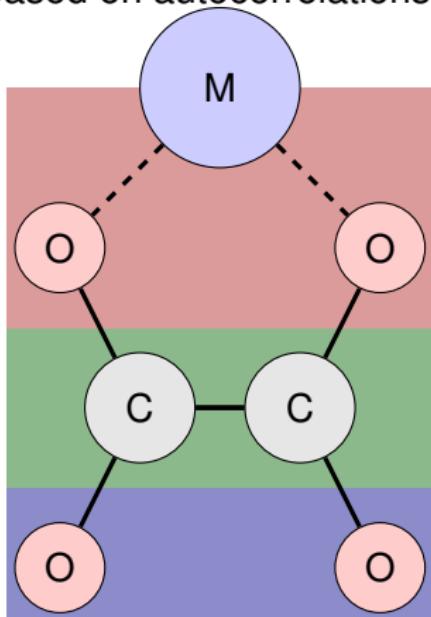


How to adapt to TM complexes?
restrict the scope to focus on
near-metal atoms

$$d_3 : \sum_{M,O} Z_M Z_O$$

Extensible, continuous descriptors - RACs

Based on autocorrelations

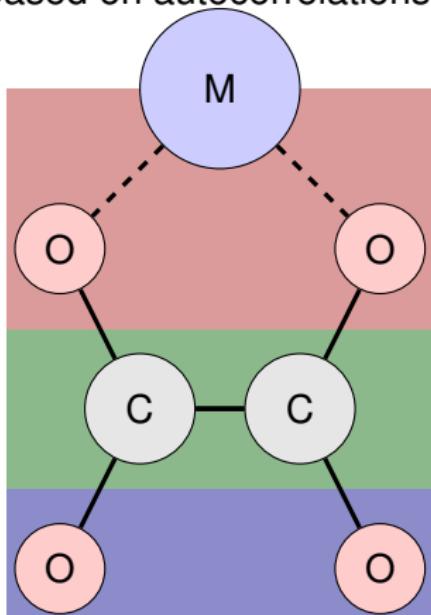


How to adapt to TM complexes?
restrict the scope to focus on
near-metal atoms

$$d_3 : \sum_{M,O} Z_M Z_O (Z_i - Z_j)$$

Extensible, continuous descriptors - RACs

Based on autocorrelations



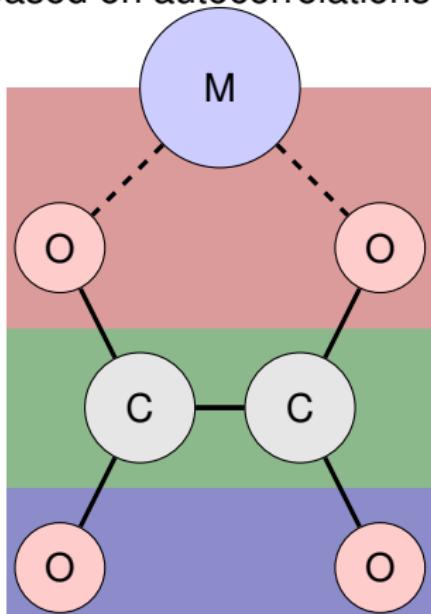
How to adapt to TM complexes?
restrict the scope to focus on
near-metal atoms

$$d_3 : \sum_{M,O} Z_M Z_O (Z_i - Z_j)$$

properties: T, χ, Z, I, S

Extensible, continuous descriptors - RACs

Based on autocorrelations

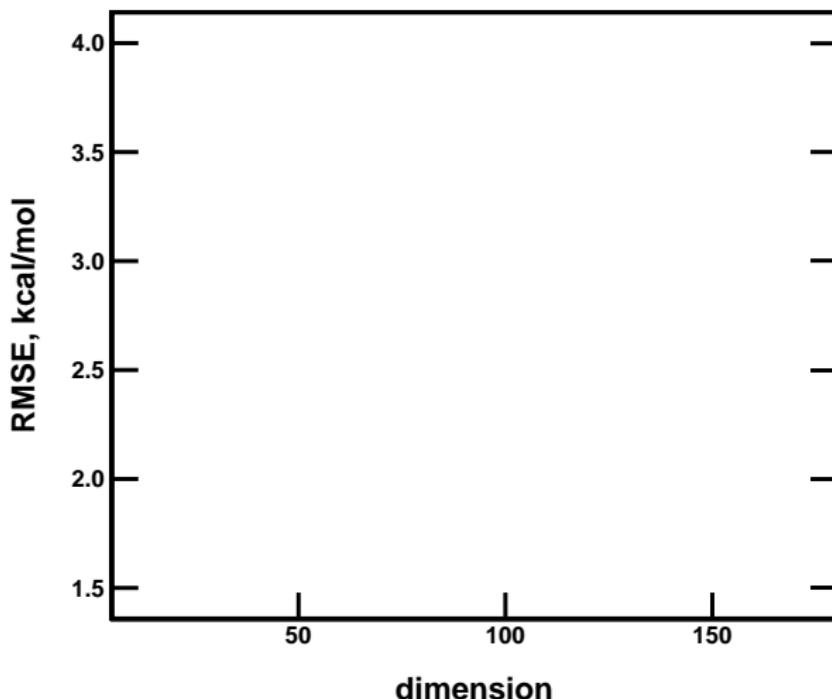


How to adapt to TM complexes?
restrict the scope to focus on
near-metal atoms

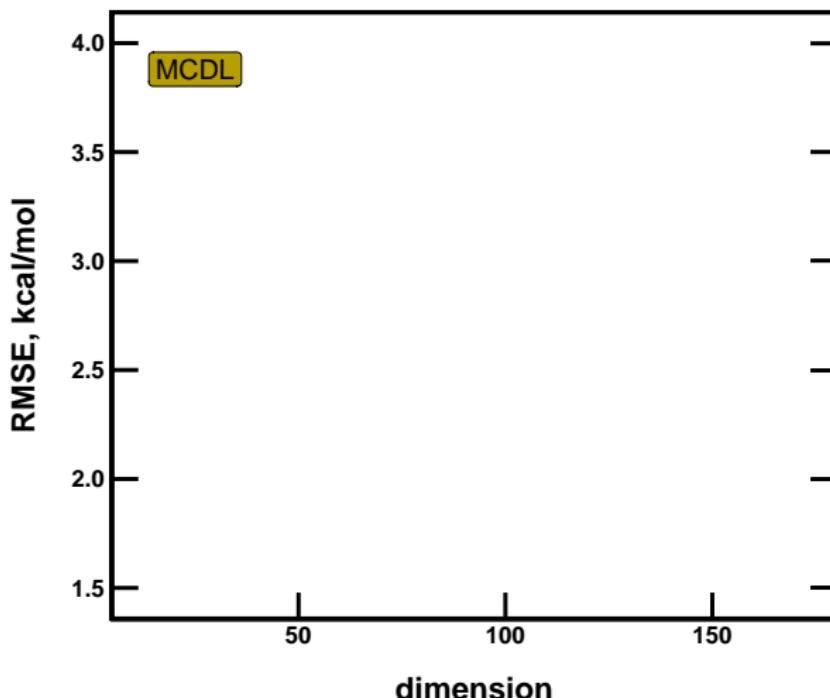
$$d_3 : \sum_{M,O} Z_M Z_O (Z_i - Z_j)$$

~ 160 features in total

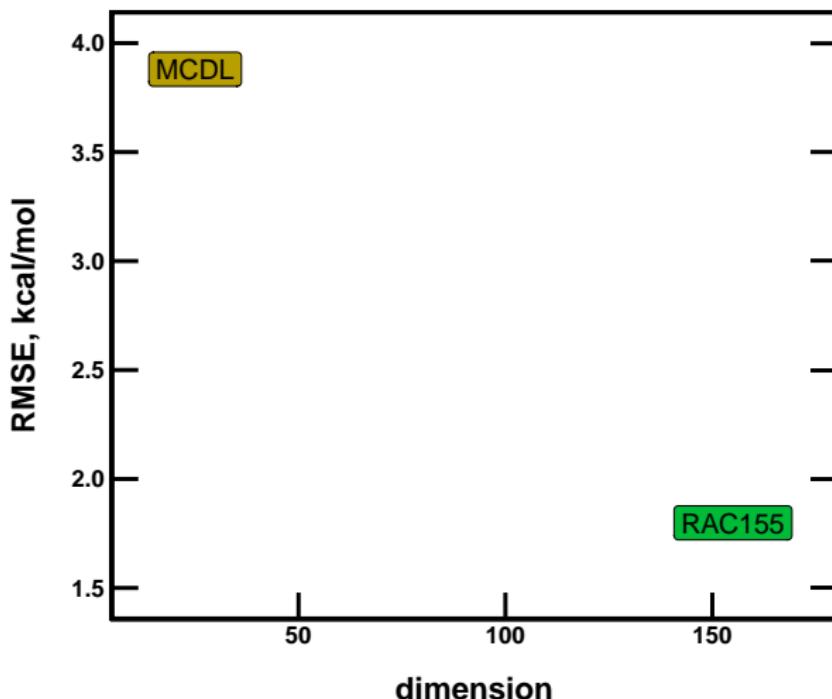
Feature selection



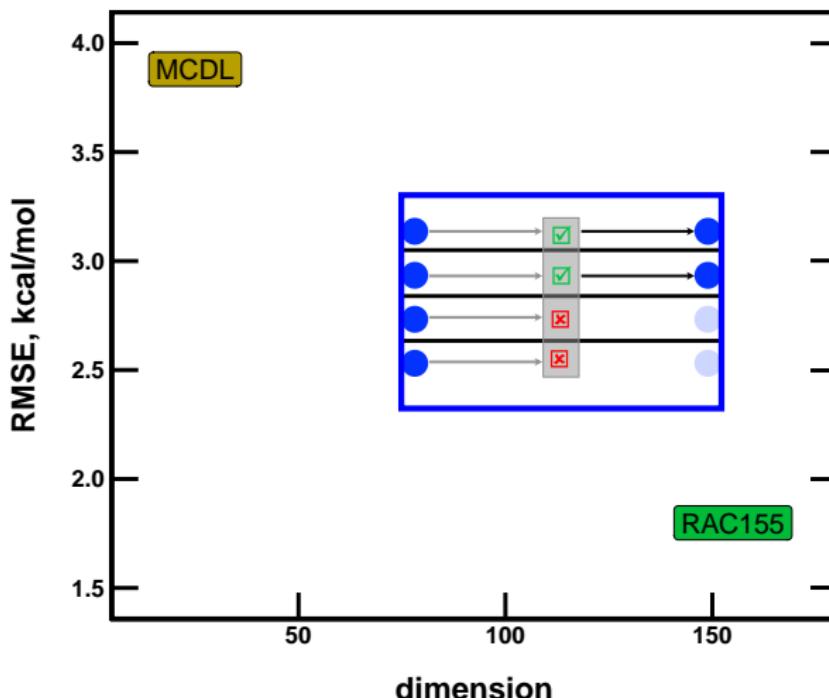
Feature selection



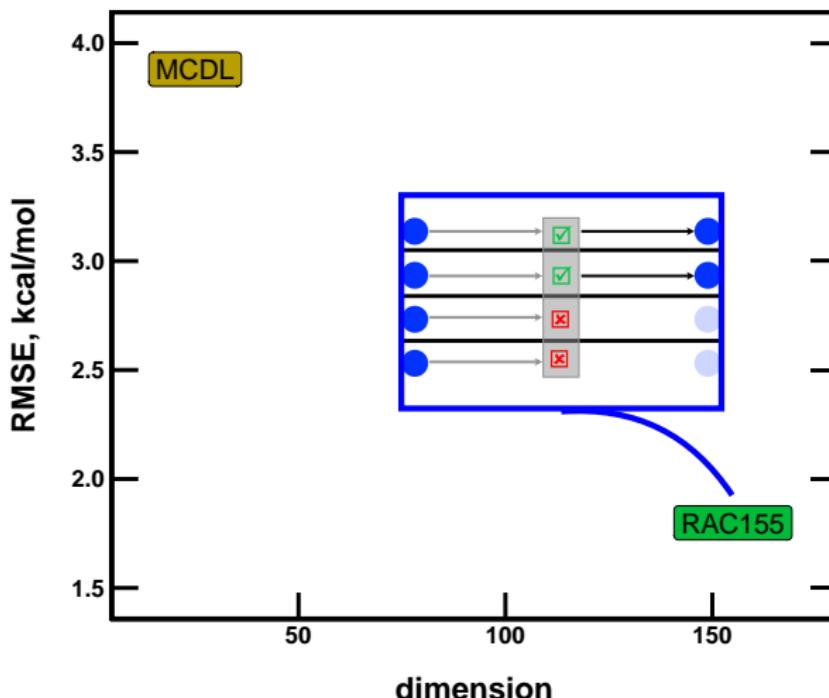
Feature selection



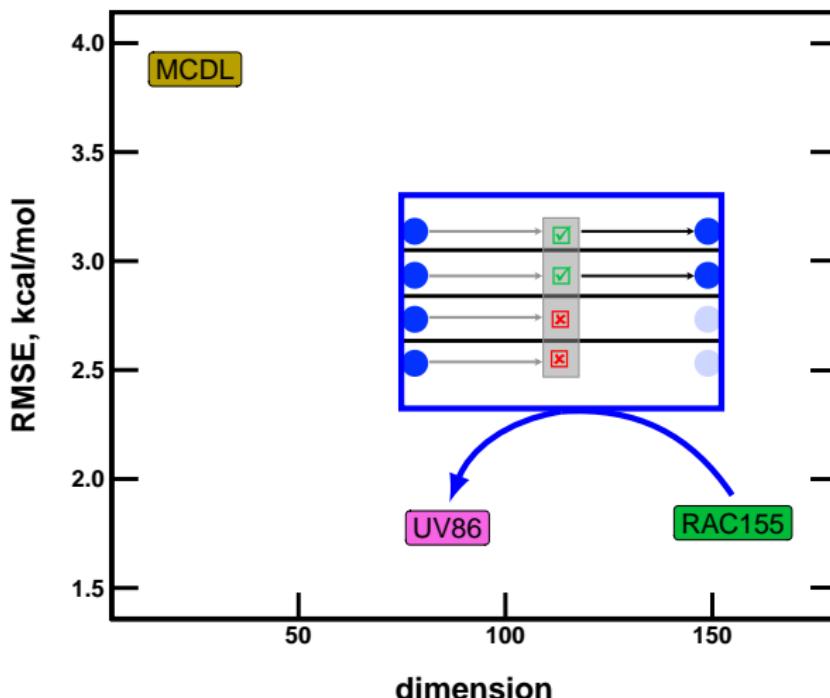
Feature selection



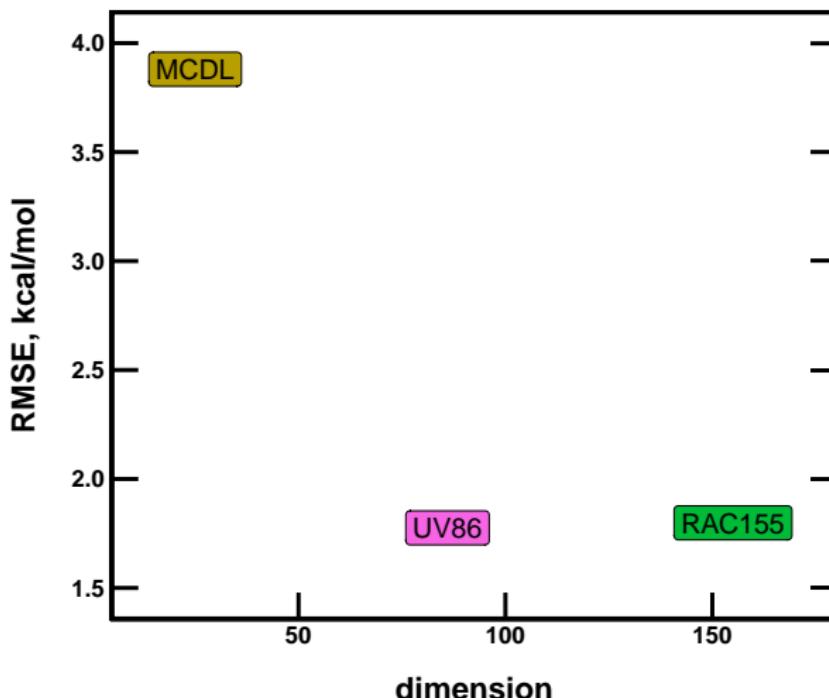
Feature selection



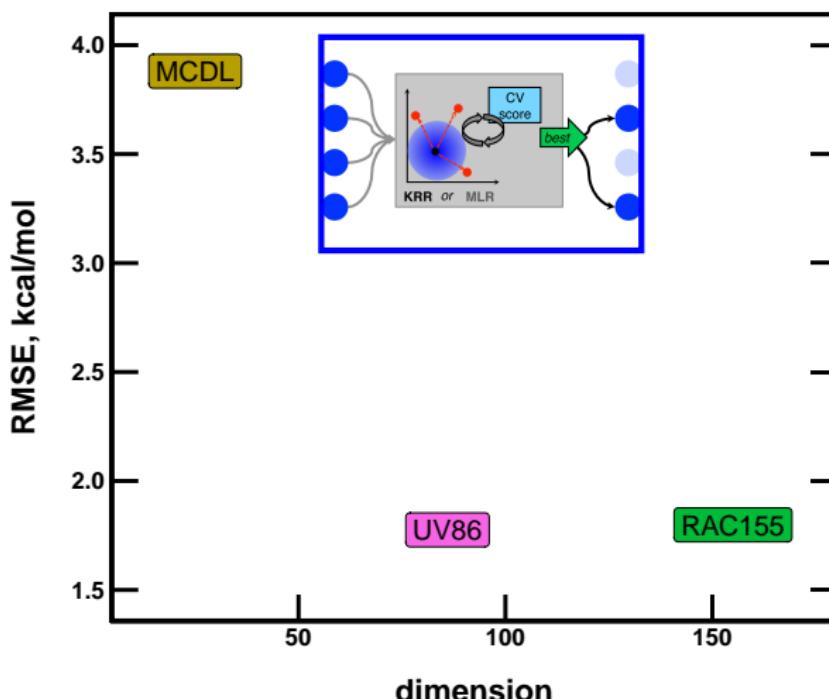
Feature selection



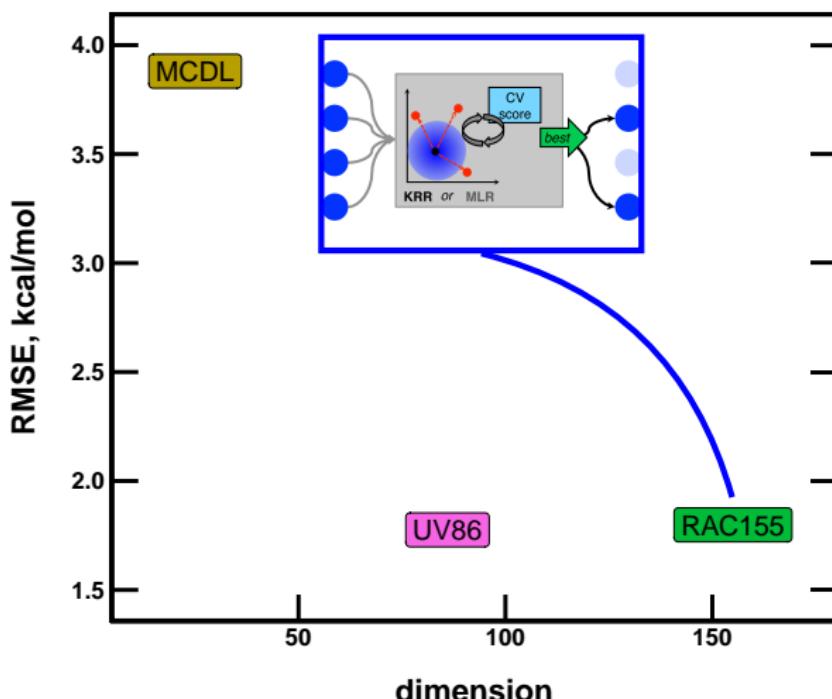
Feature selection



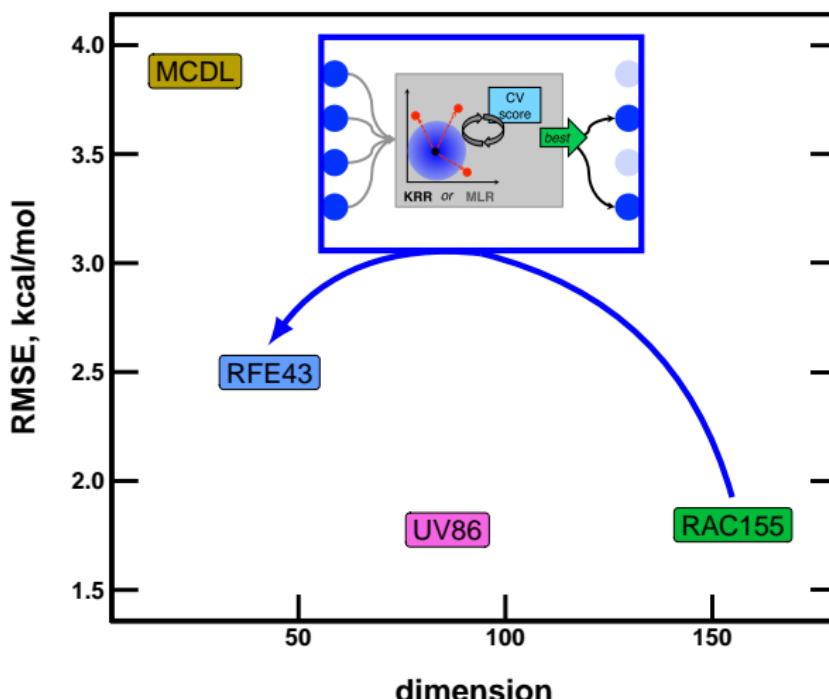
Feature selection



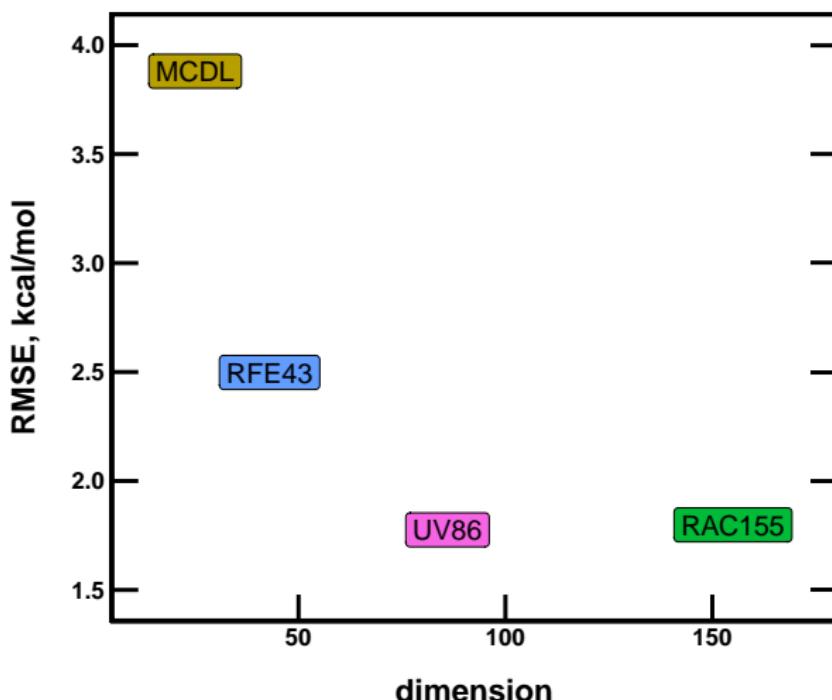
Feature selection



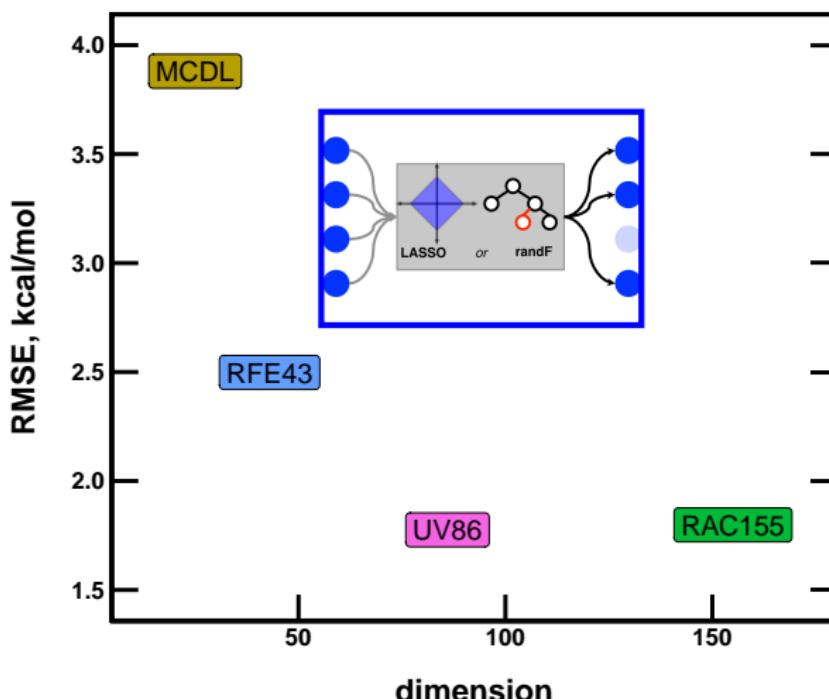
Feature selection



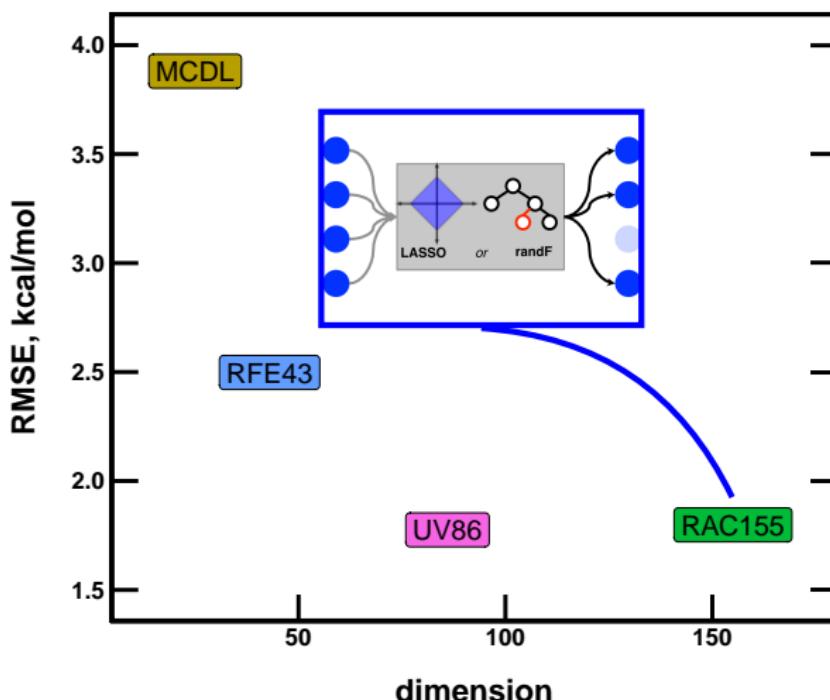
Feature selection



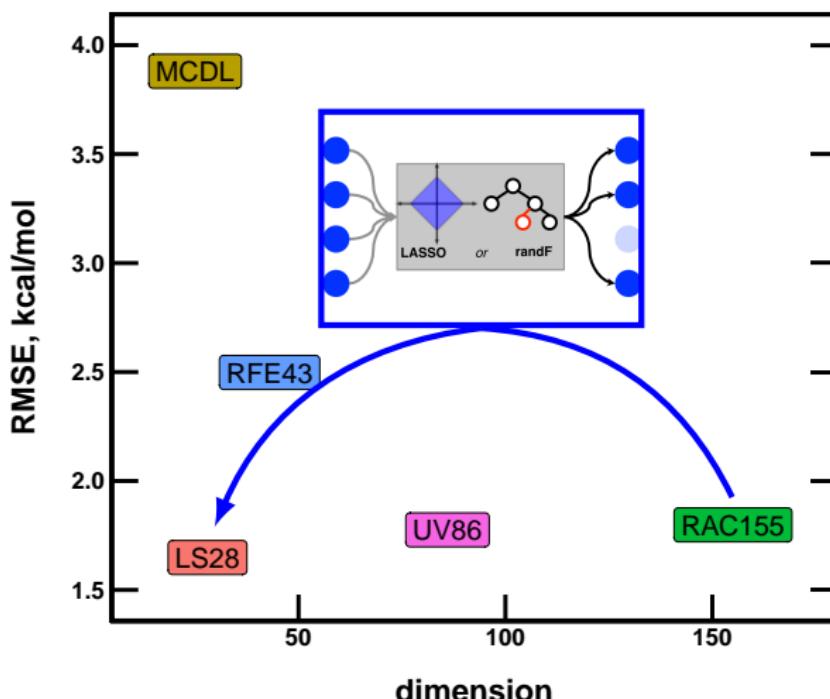
Feature selection



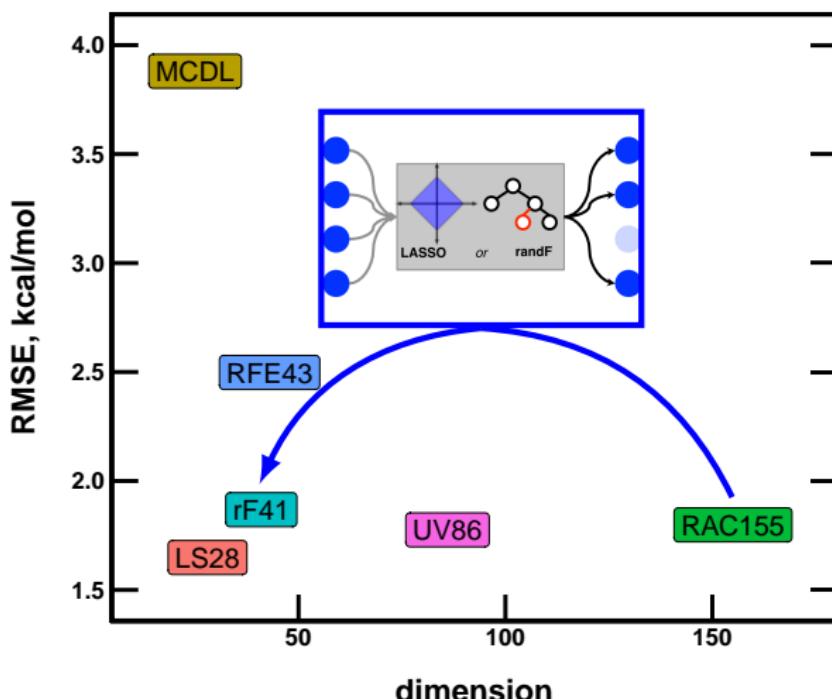
Feature selection



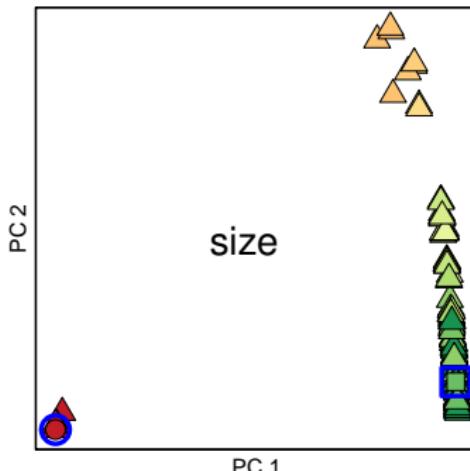
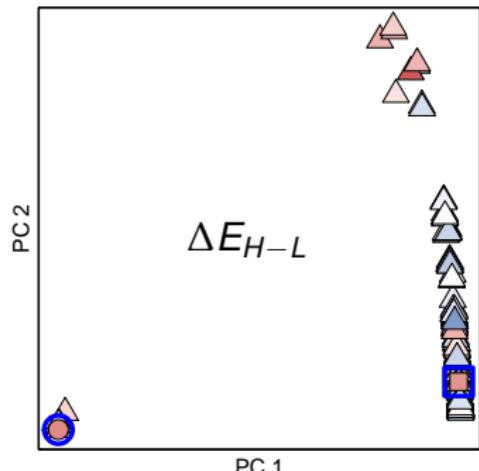
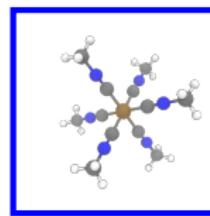
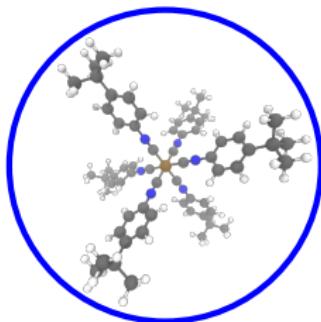
Feature selection



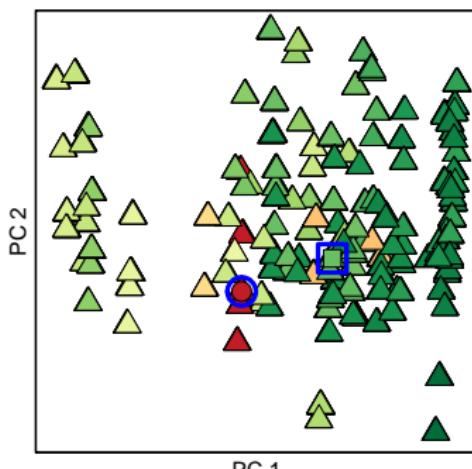
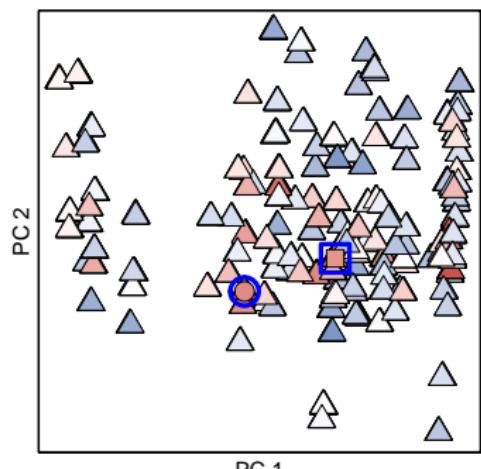
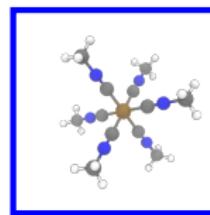
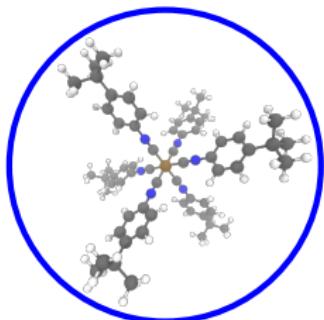
Feature selection



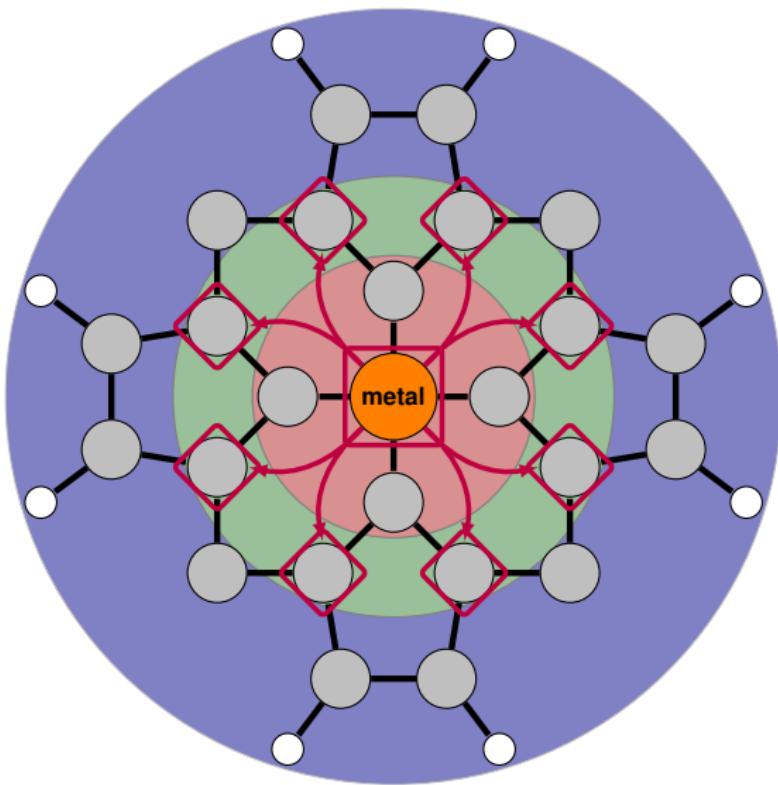
A tale of two complexes, II



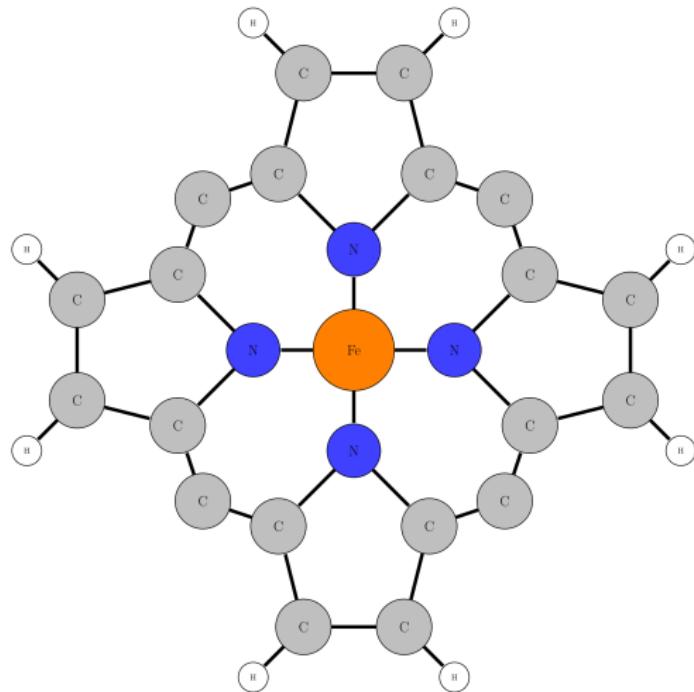
A tale of two complexes, II



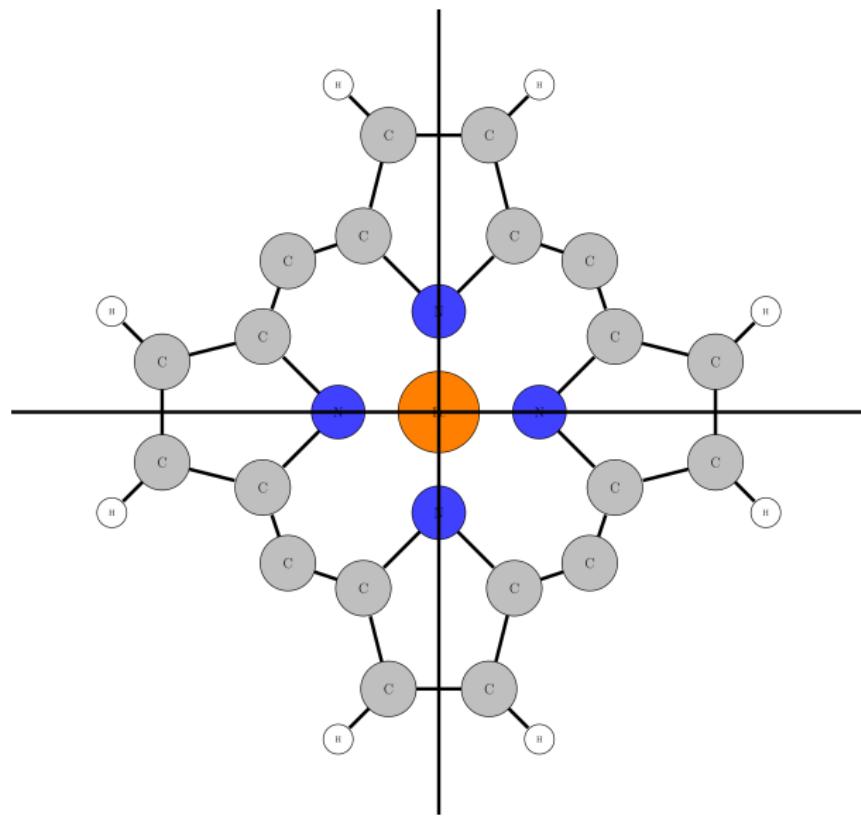
Do features depend on properties?



Do features depend on properties?

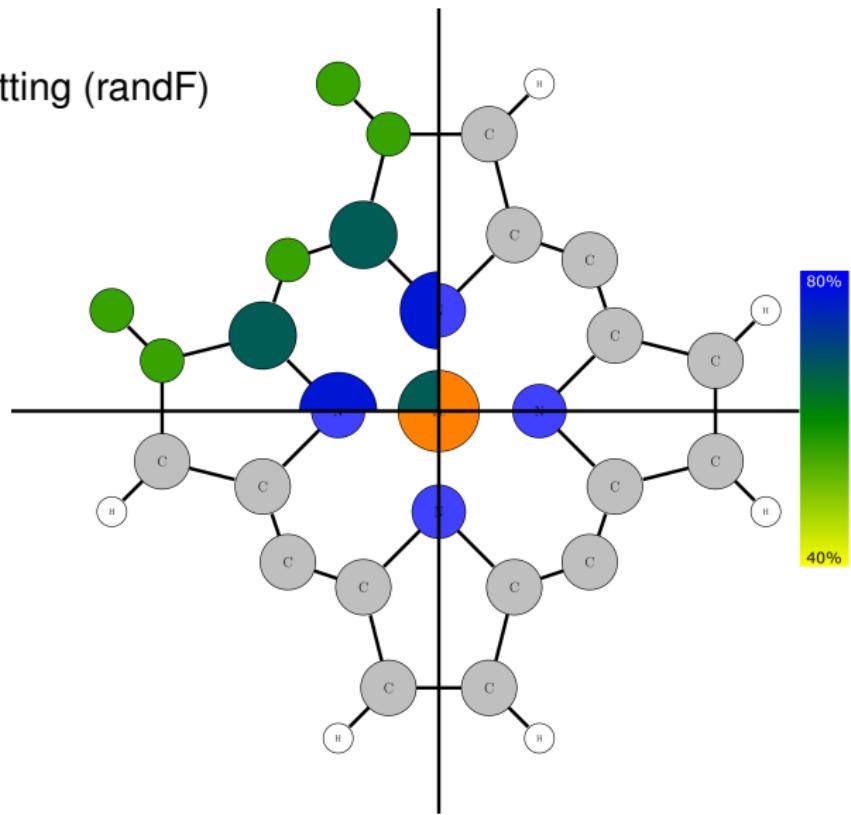


Do features depend on properties?



Do features depend on properties?

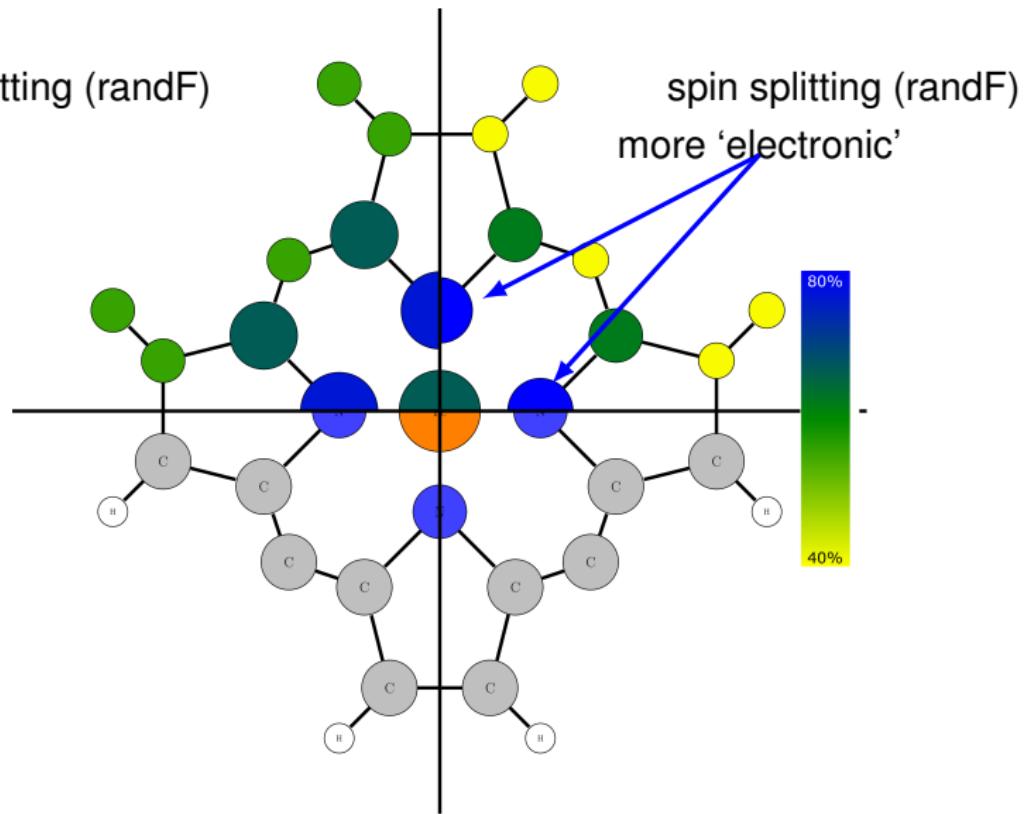
spin splitting (randF)



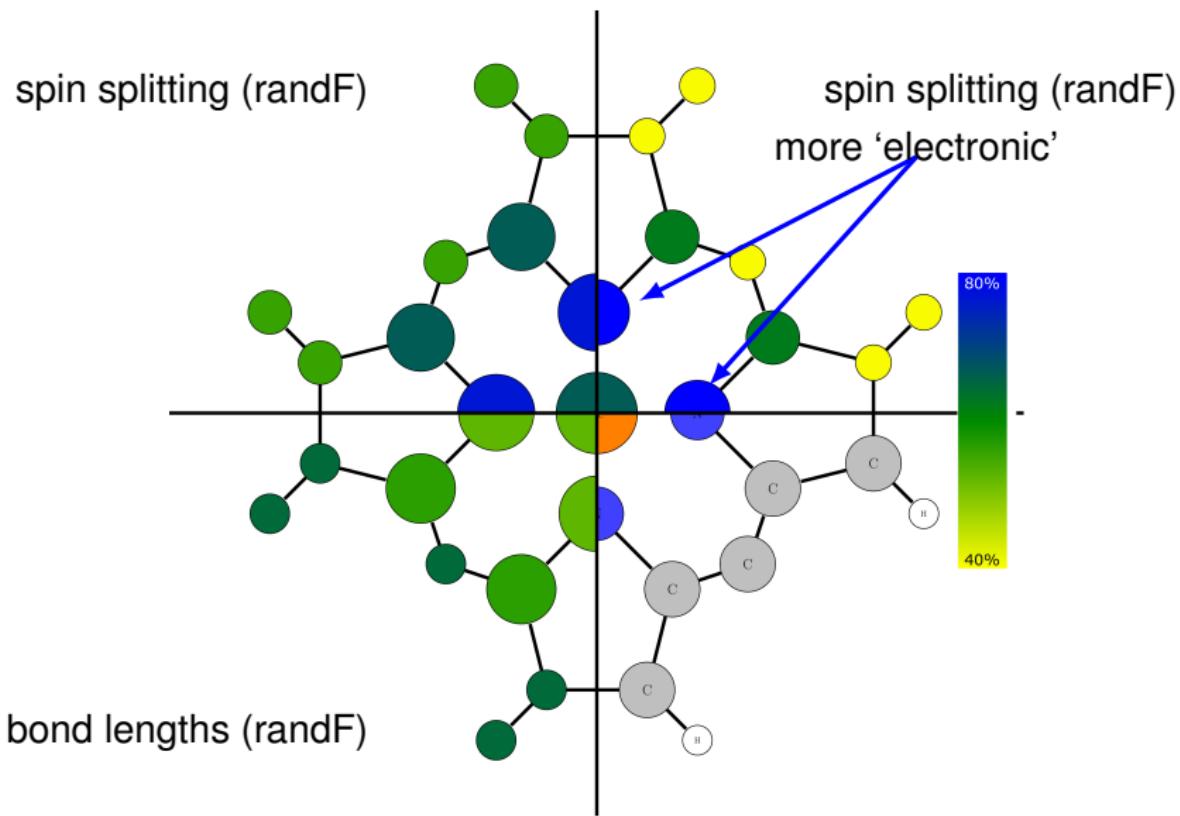
Do features depend on properties?

spin splitting (randF)

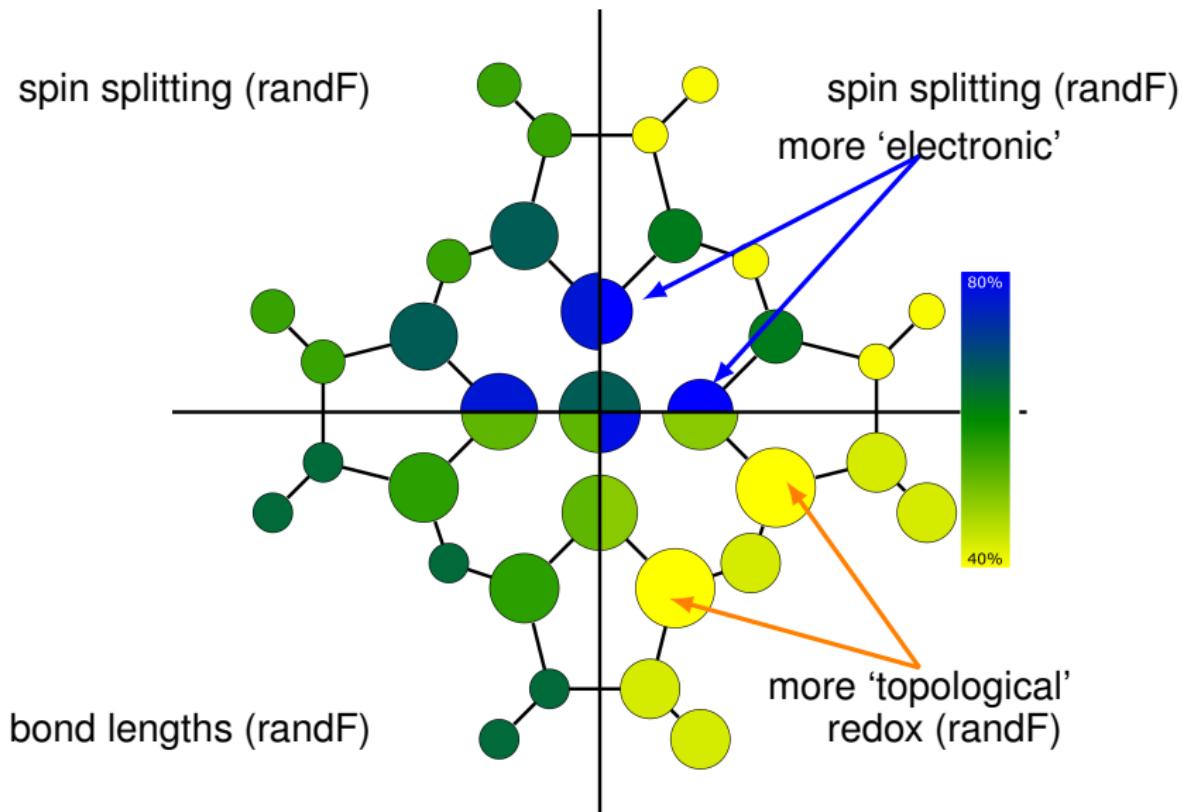
spin splitting (randF)
more 'electronic'



Do features depend on properties?



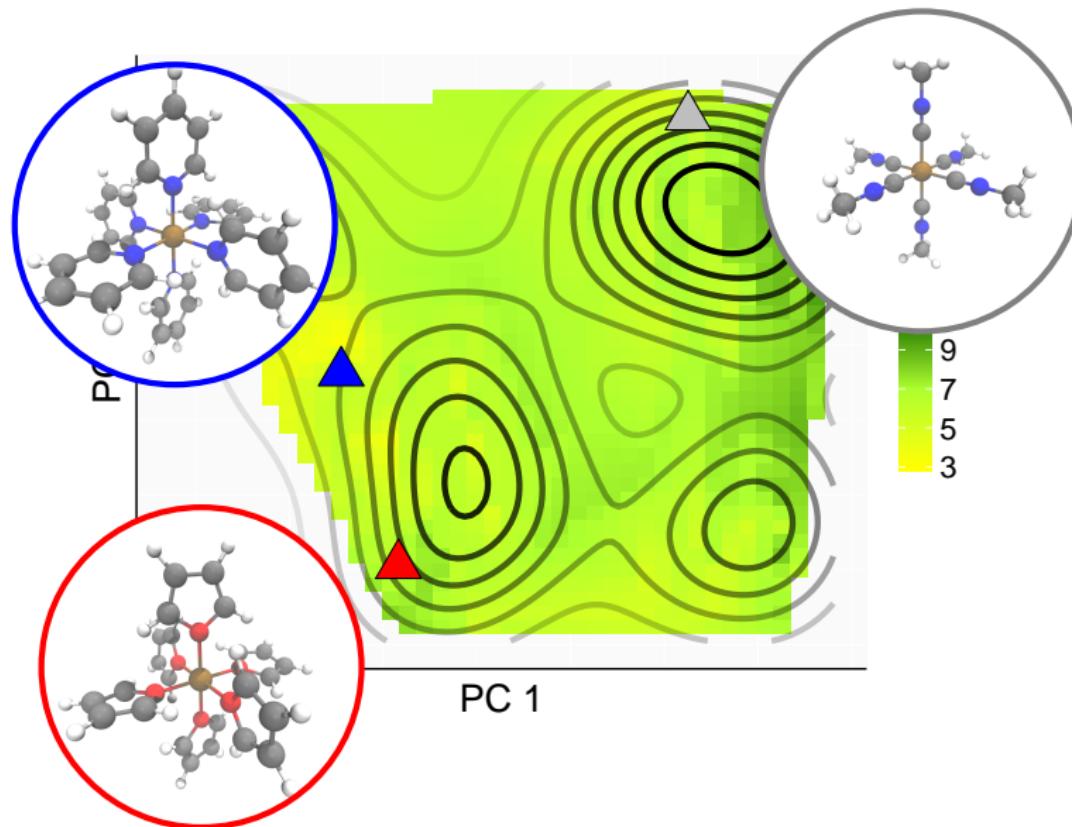
Do features depend on properties?



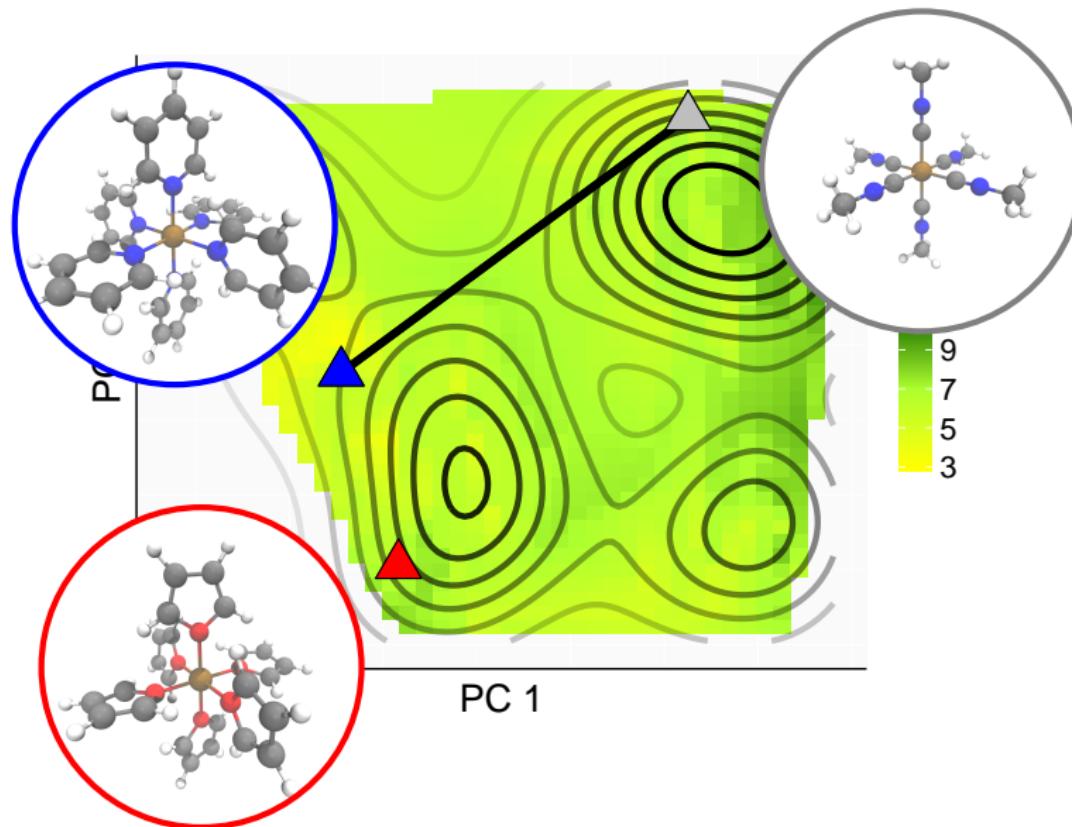
Mapping TM complex space



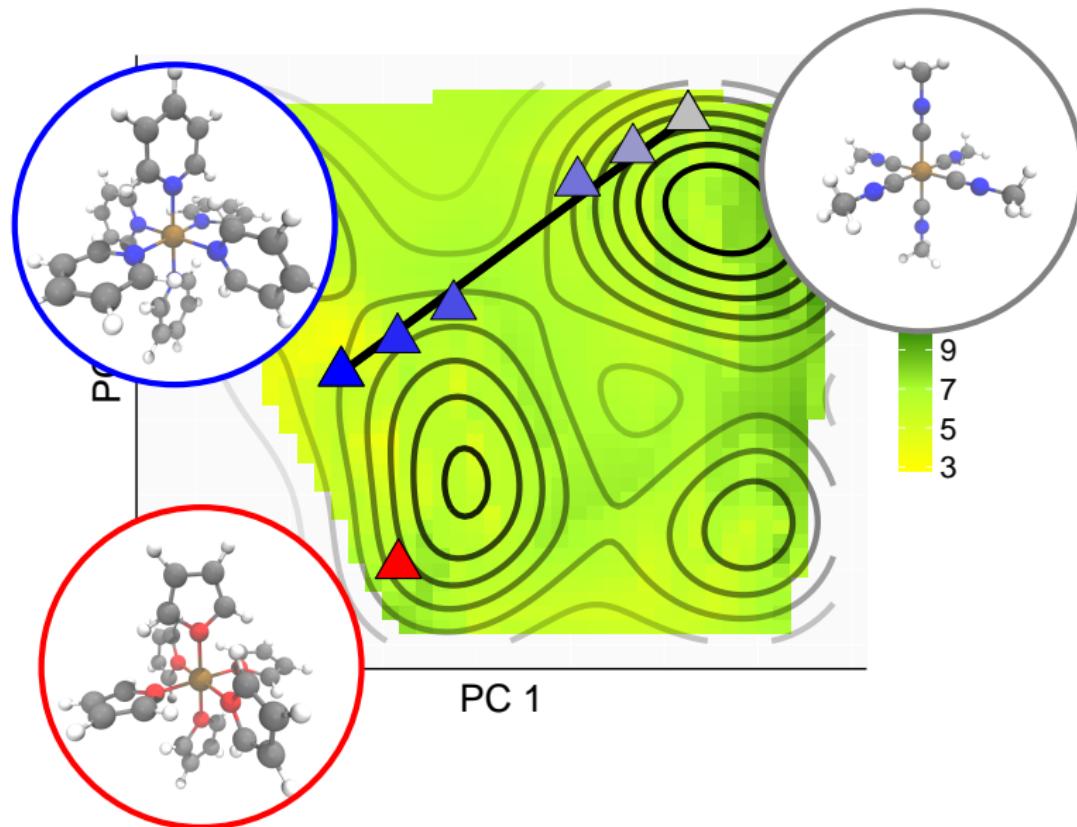
Mapping TM complex space



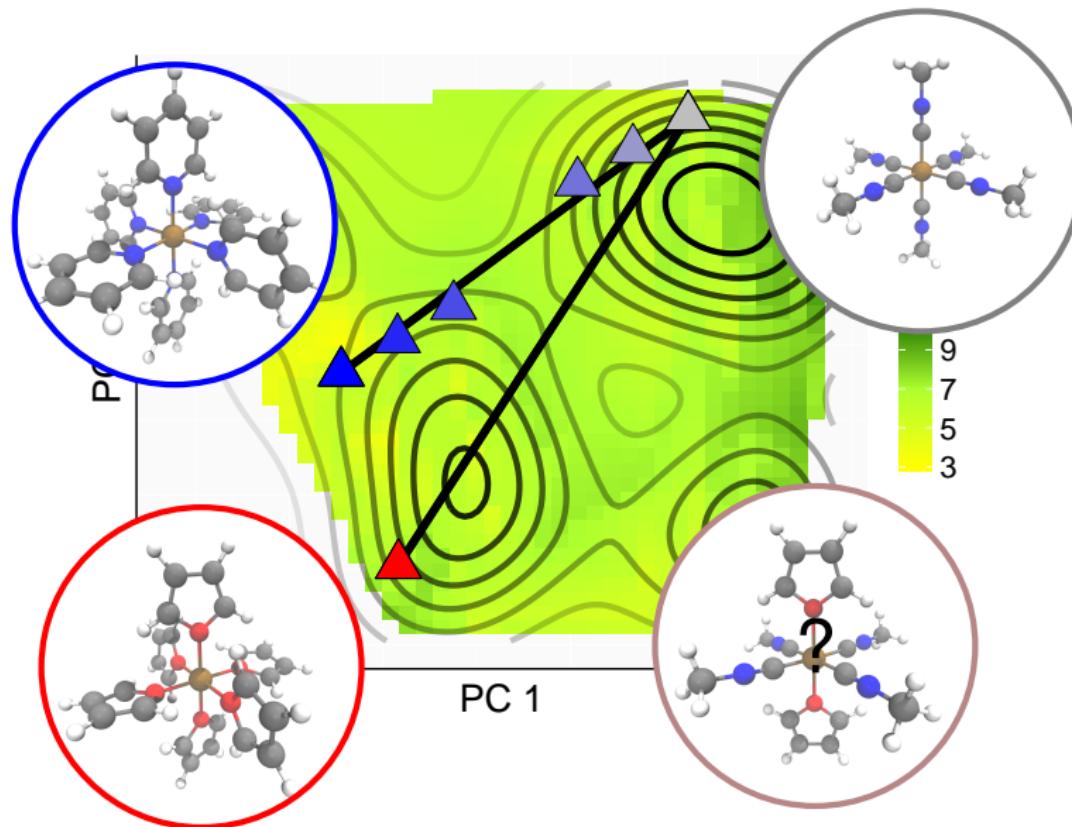
Mapping TM complex space



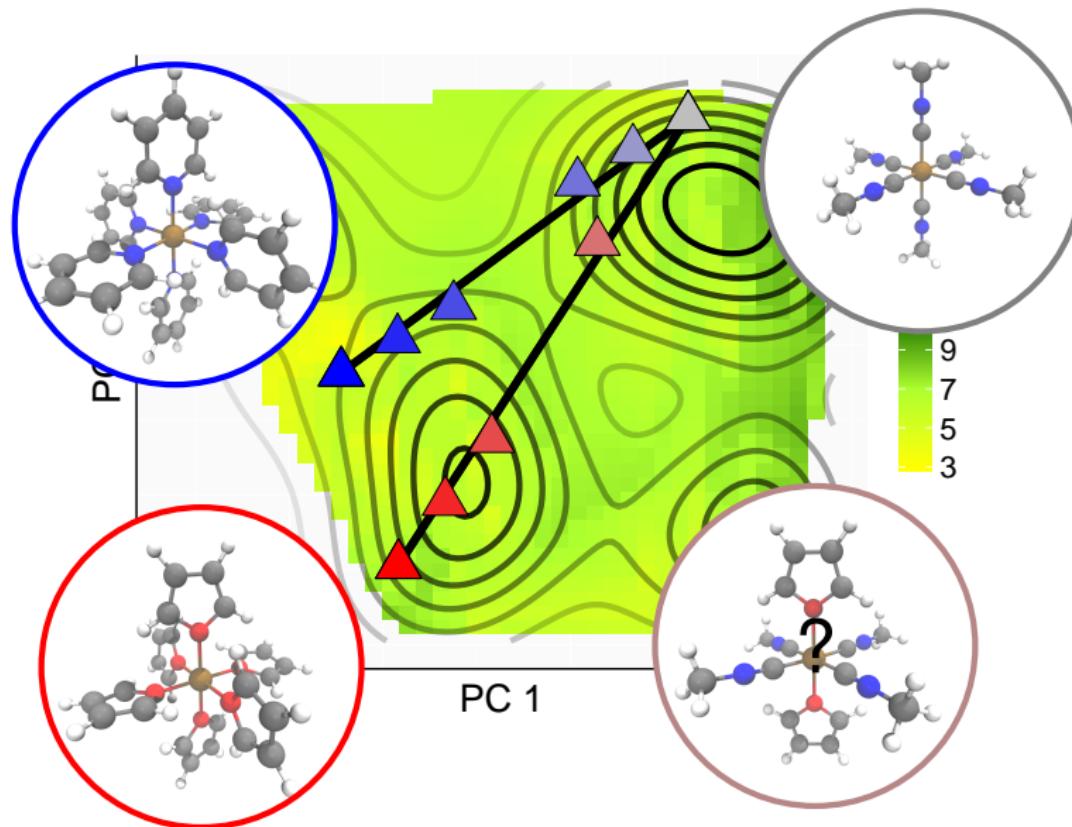
Mapping TM complex space



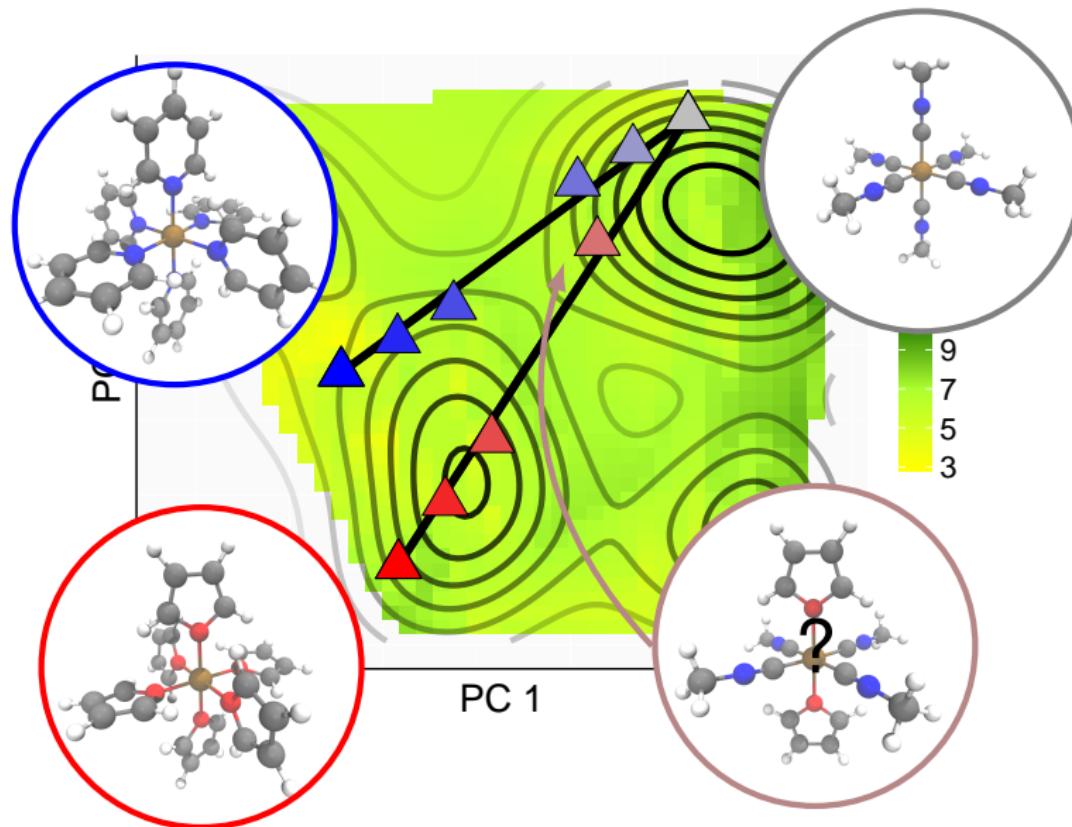
Mapping TM complex space



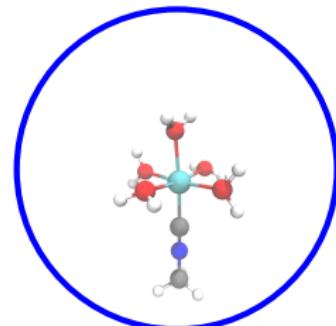
Mapping TM complex space



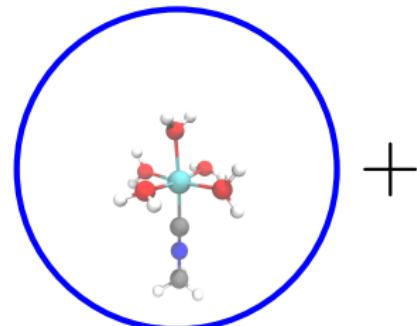
Mapping TM complex space



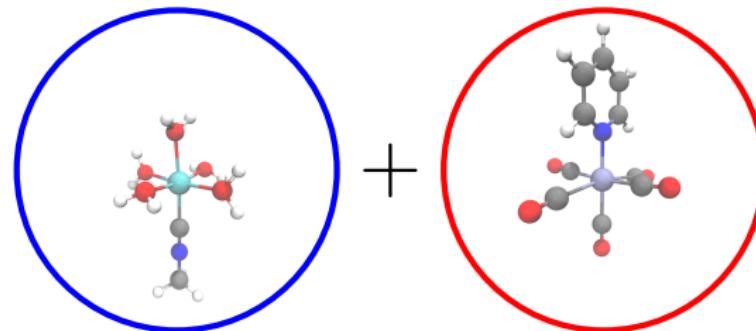
Mapping TM complex space



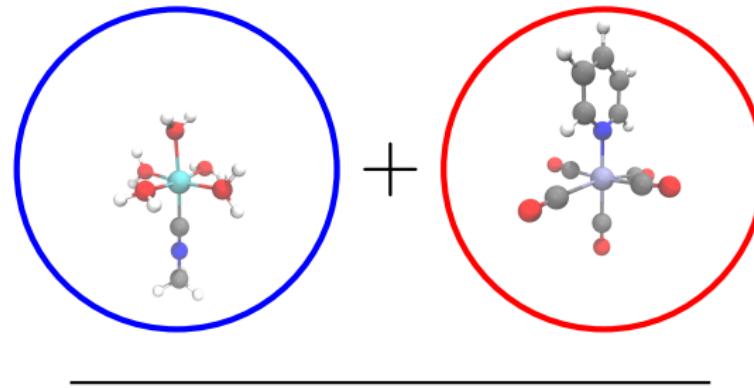
Mapping TM complex space



Mapping TM complex space

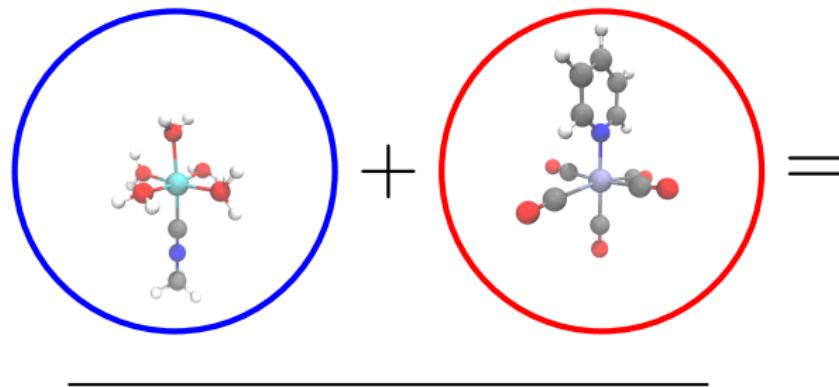


Mapping TM complex space



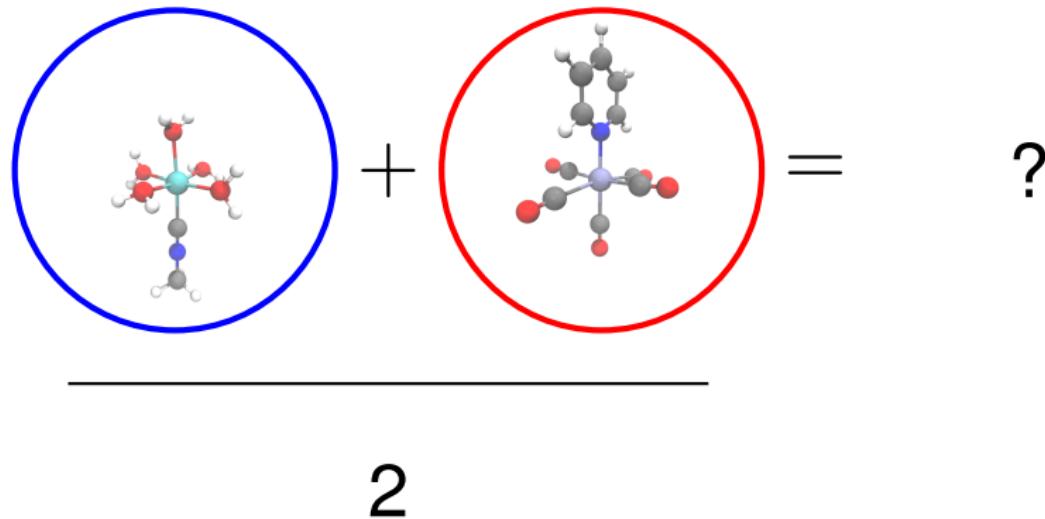
2

Mapping TM complex space

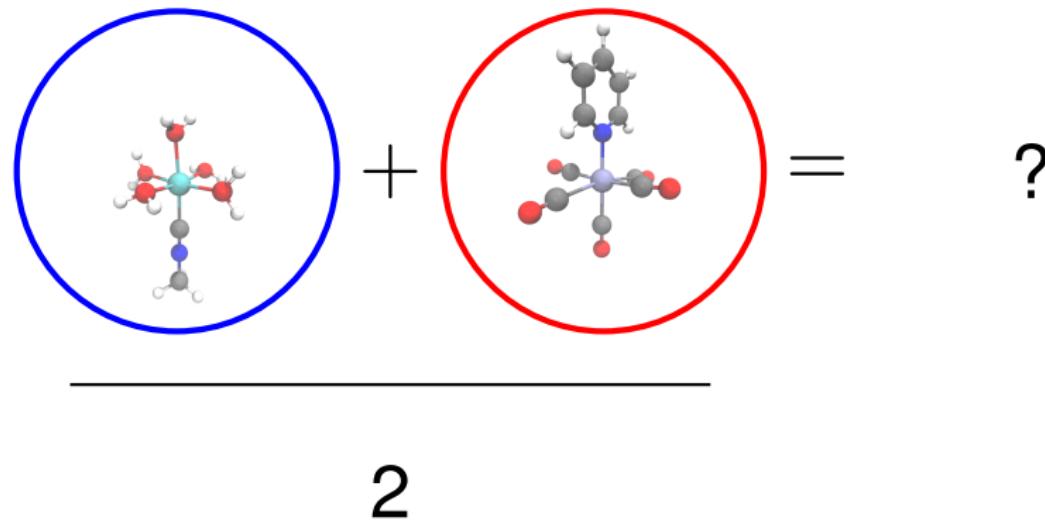


2

Mapping TM complex space



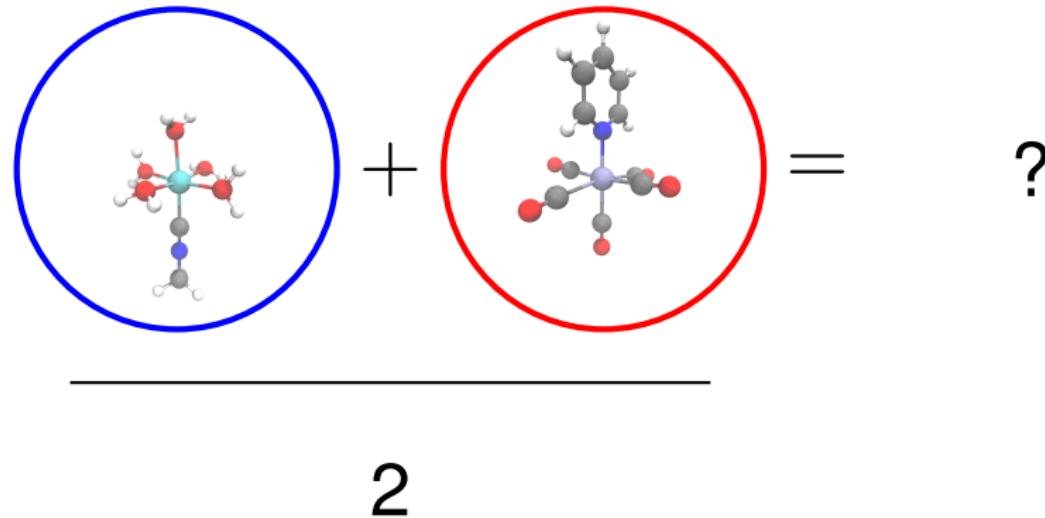
Mapping TM complex space



Cr(II) $[\text{H}_2\text{O}]_5$ [misc]

Co(II) $[\text{CO}]_5$ [pyr]

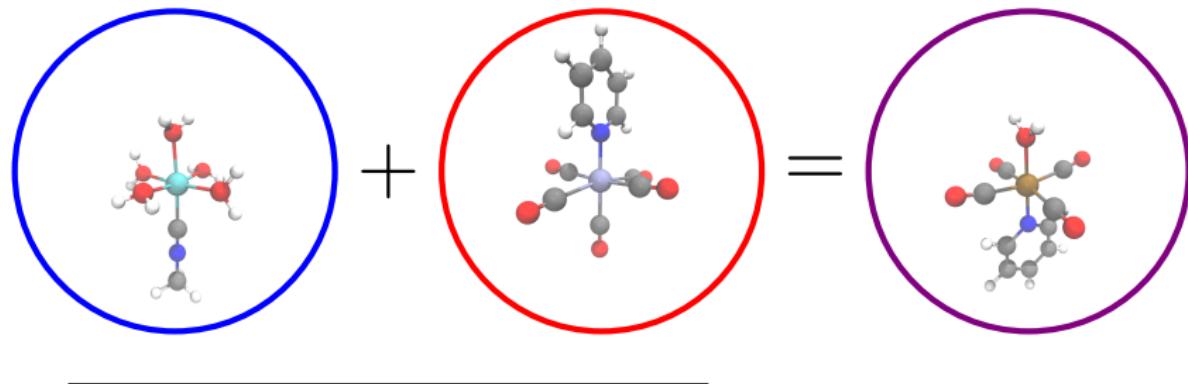
Mapping TM complex space



$\text{Cr(II)} [\text{H}_2\text{O}]_5 [\text{misc}]$
 $\Delta G = 5.3 \text{ eV}$

$\text{Co(II)} [\text{CO}]_5 [\text{pyr}]$
 $\Delta G = 8.1 \text{ eV}$

Mapping TM complex space



2

Cr(II) $[\text{H}_2\text{O}]_5$ [misc]
 $\Delta G = 5.3 \text{ eV}$

Co(II) $[\text{CO}]_5$ [pyr]
 $\Delta G = 8.1 \text{ eV}$

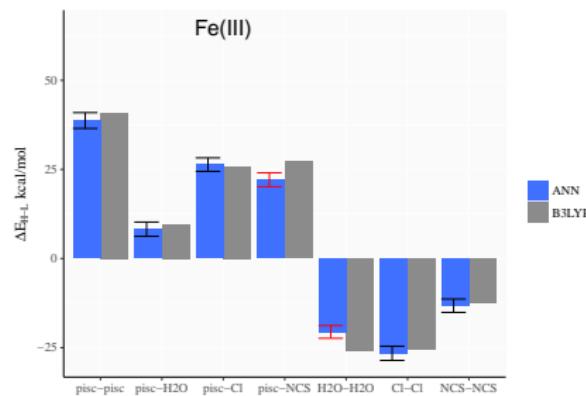
Fe(II) $[\text{CO}]_4$ [pyr][water]
 $\Delta G = 7.8 \text{ eV}$

Model transferability

Test-set performance is not necessarily a good metric for general transferability²:

Model transferability

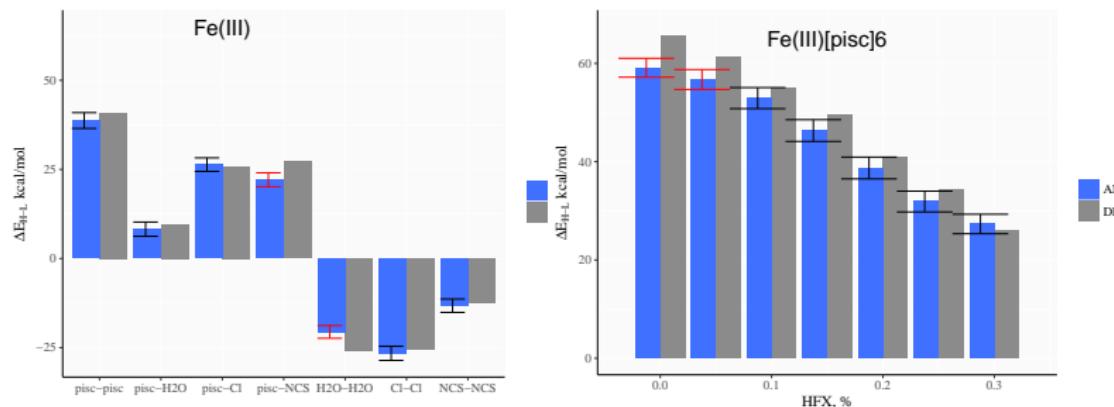
Test-set performance is not necessarily a good metric for general transferability²:



²Janet, J.P., and Kulik, H.J. *Chem. Sci.*, 2017, 8, 5137-5152.

Model transferability

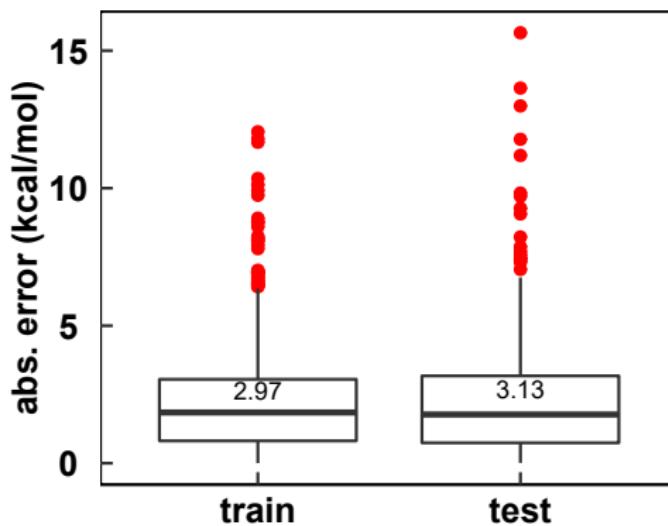
Test-set performance is not necessarily a good metric for general transferability²:



²Janet, J.P., and Kulik, H.J. *Chem. Sci.*, 2017, 8, 5137-5152.

Model transferability

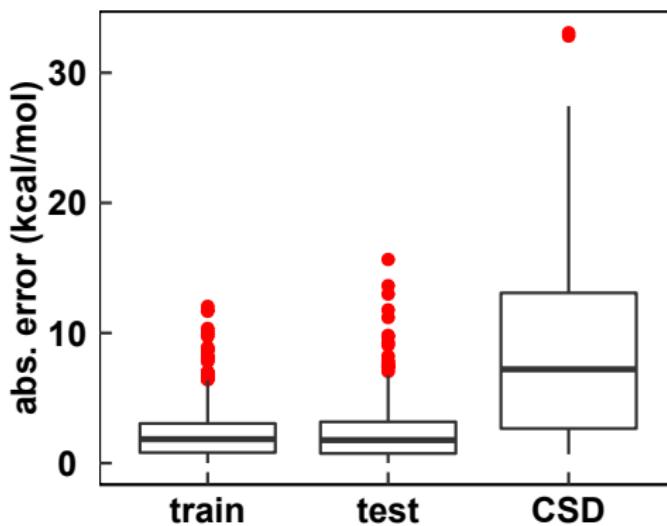
Test-set performance is not necessarily a good metric for general transferability²:



²Janet, J.P., and Kulik, H.J. *Chem. Sci.*, 2017, 8, 5137-5152.

Model transferability

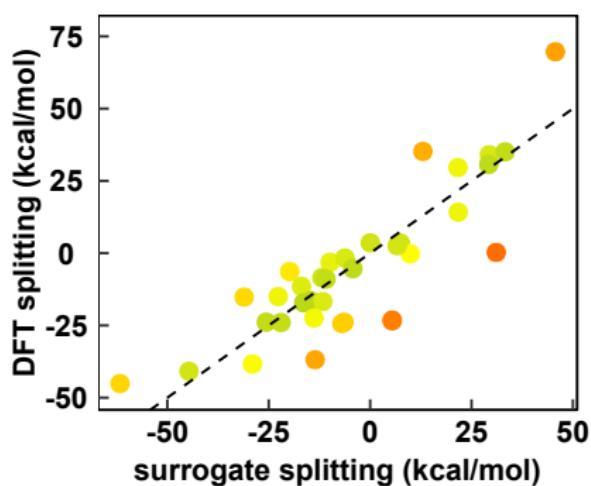
Test-set performance is not necessarily a good metric for general transferability²:



²Janet, J.P., and Kulik, H.J. *Chem. Sci.*, 2017, 8, 5137-5152.

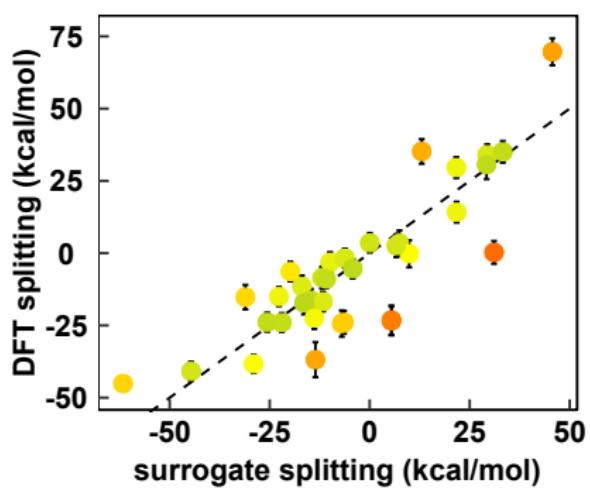
Model transferability

Uncertainty estimates are essential for our surrogate model to explore chemical space:



Model transferability

Uncertainty estimates are essential for our surrogate model to explore chemical space:



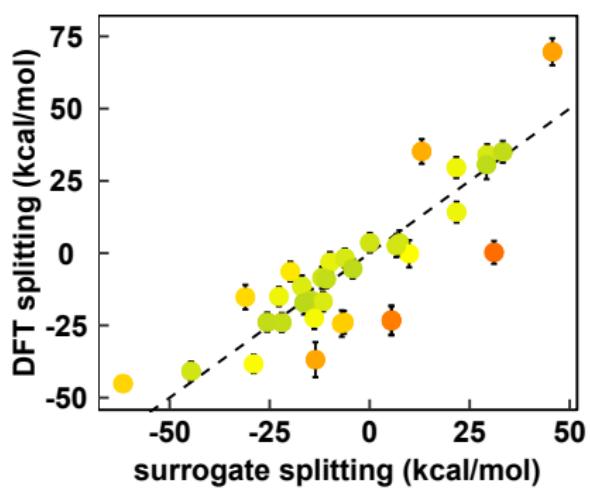
Uncertainty from mc-dropout¹:
ANN model approximates variational inference with GP under some conditions:

$$\text{var}(y^*|x^*) \approx \frac{1}{J} \sum_j \hat{y}_j^T \hat{y}_j + \tau^{-1}$$

Gal, Y. and Ghahramani, Z., 2016. ICMLR 1050-1059

Model transferability

Uncertainty estimates are essential for our surrogate model to explore chemical space:



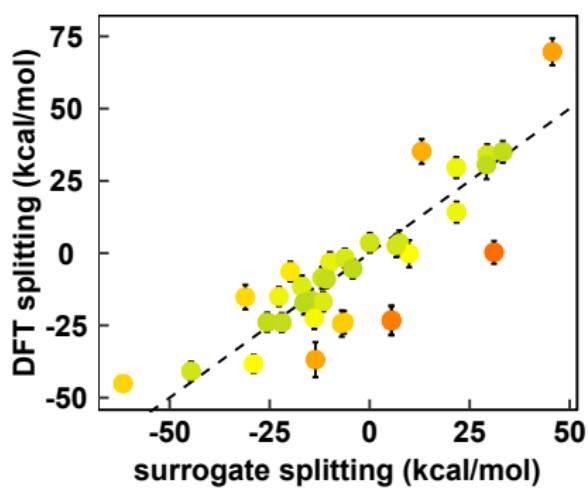
Uncertainty from mc-dropout¹:
ANN model approximates variational inference with GP under some conditions:

$$\text{var}(y^*|x^*) \approx \frac{1}{J} \sum_j \hat{y}_j^T \hat{y}_j + \tau^{-1}$$

Gal, Y. and Ghahramani, Z., 2016. ICMLR 1050-1059

Model transferability

Uncertainty estimates are essential for our surrogate model to explore chemical space:



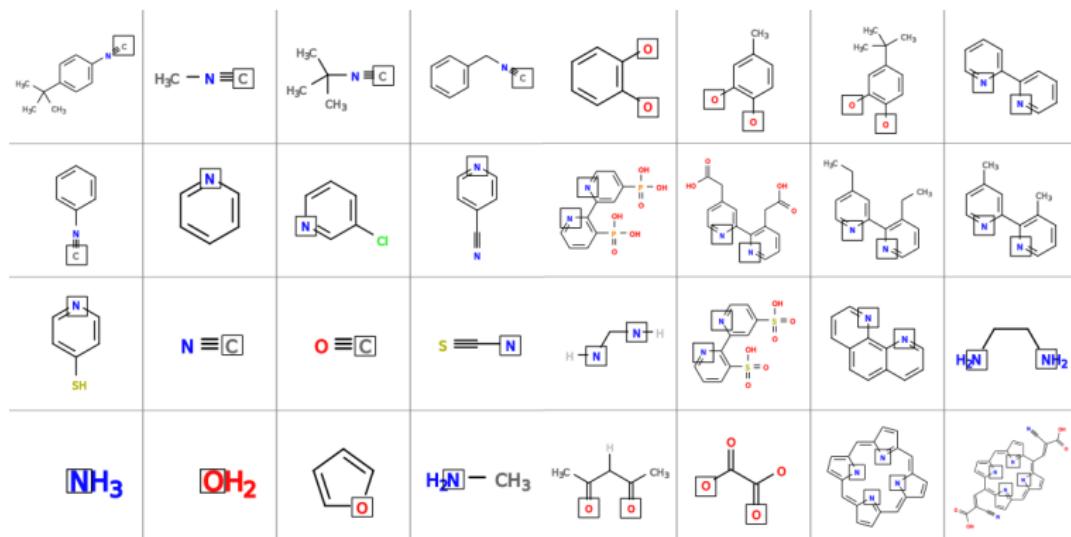
Demonstration

Can we use the ANN model to find new spin-crossover materials,
i.e. $\Delta E_{H-L} = 0$?

³Janet, J.P., Chan, L. and Kulik, H.J. *J. Phys. Chem. Lett.*, 2018, 9, 5, 1064-1071.

Demonstration

Can we use the ANN model to find new spin-crossover materials, i.e. $\Delta E_{H-L} = 0$? Define a space of 32 ligands, 5 metals and with ~ 5600 possible elements with forced axial/equatorial symmetry³:



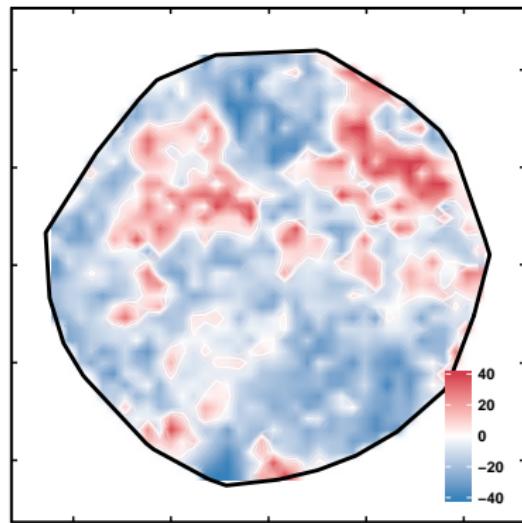
³Janet, J.P., Chan, L. and Kulik, H.J. *J. Phys. Chem. Lett.*, 2018, 9, 5, 1064-1071.

Demonstration

ANN is trained on 14 of these ligands, covers only **2%** of the design space.

Demonstration

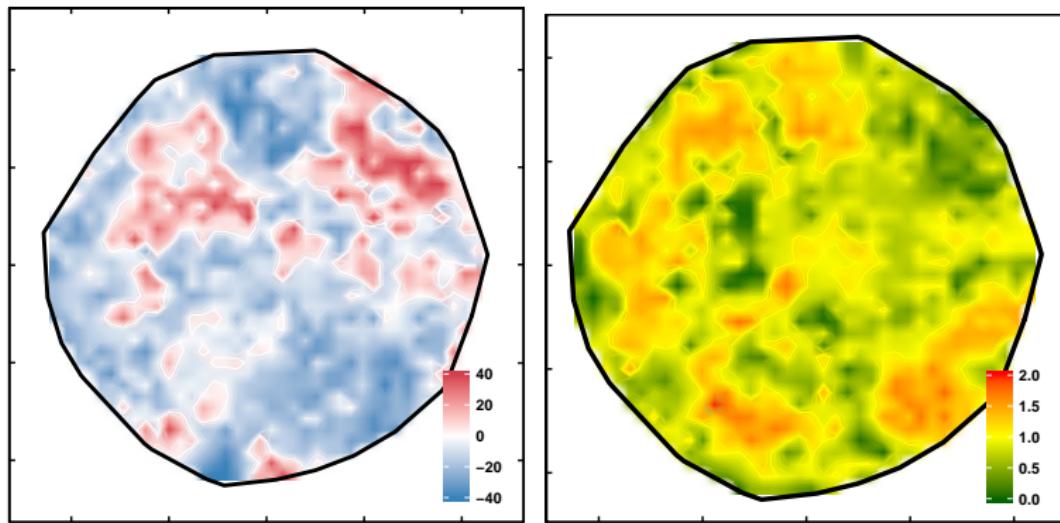
ANN is trained on 14 of these ligands, covers only **2%** of the design space. We can visualize the design space using t-SNE⁴:



⁴Maaten, L., & Hinton, G., 2008. J. Mach. Learn. Res. 2579-2605.

Demonstration

ANN is trained on 14 of these ligands, covers only **2%** of the design space. We can visualize the design space using t-SNE⁴:



⁴Maaten, L., & Hinton, G., 2008. J. Mach. Learn. Res. 2579-2605.

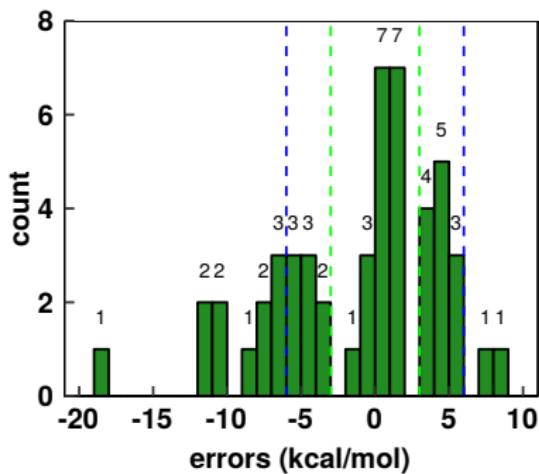
How accurate are we?

Test 51 leads from ANN with DFT⁵:

⁵Janet, J.P., Chan, L. and Kulik, H.J. *J. Phys. Chem. Lett.*, 2018, 9, 5, 1064-1071.

How accurate are we?

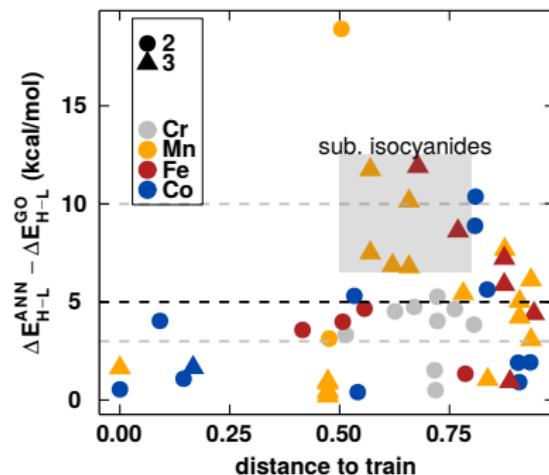
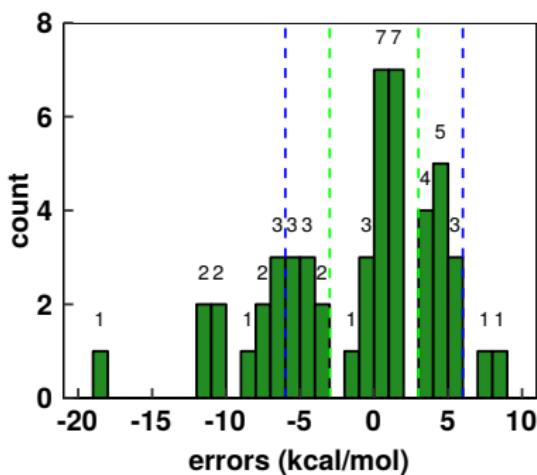
Test 51 leads from ANN with DFT⁵:



⁵Janet, J.P., Chan, L. and Kulik, H.J. *J. Phys. Chem. Lett.*, 2018, 9, 5, 1064-1071.

How accurate are we?

Test 51 leads from ANN with DFT⁵:



⁵Janet, J.P., Chan, L. and Kulik, H.J. *J. Phys. Chem. Lett.*, 2018, 9, 5, 1064-1071.

Conclusions

- choice of molecular representation is important
- different properties depend non-equally on features
- feature-space geometry can provide insight into model reliability
- imbuing ‘chemical intuition’ to descriptor construction can drastically improve learning
- conversely, feature selection can contribute to understanding systems

Acknowledgments

Thanks to the Kulik group and funding partners:

