# Learning Reaction Barriers And Transition State Geometries

Guido Falk von Rudorff, University of Vienna

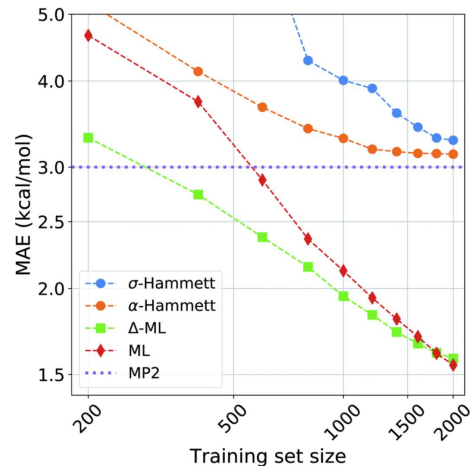ferchault     @ferchault     guido.vonrudorff.de

- Reactions: complicated landscape
- Not only expensive but also hard problem
- Even if the reaction mechanism is known:
  - Find reactant (R)/product (P) complexes
  - Find transition state (TS) geometries
  - Describe energy near TS
  - Low level of automation available
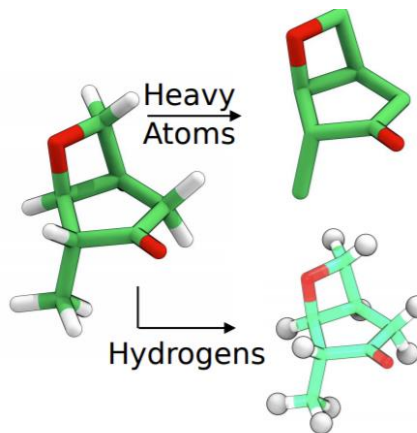
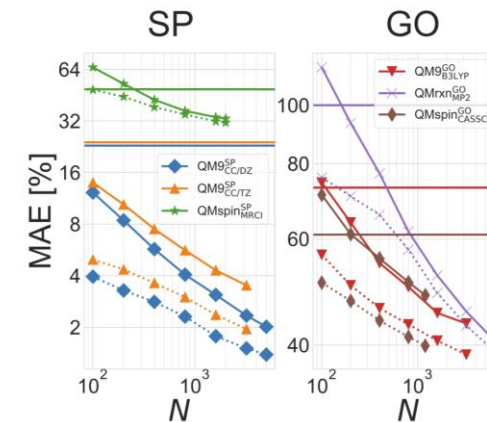- Machine learning: accelerate & less supervised

Energies with KRR

Detrending with Hammett's equation

Geometries with Graph2Structure

Estimate computational cost

qmlcode/qml

chemspacelab/ Enhanced-Hammett

qmlcode/qml

ferchault/mlscheduling

## Idea

- Molecular representation for each molecule $i$
  - CM, BoB, FCHL, SLATM, ...
- Distance metric
  - Typically L1 or L2 norm
- Kernel function
  - Laplacian, Gaussian

$$\mathbf{M}_i$$

$$d_{ij} \equiv d(\mathbf{M}_i, \mathbf{M}_j)$$
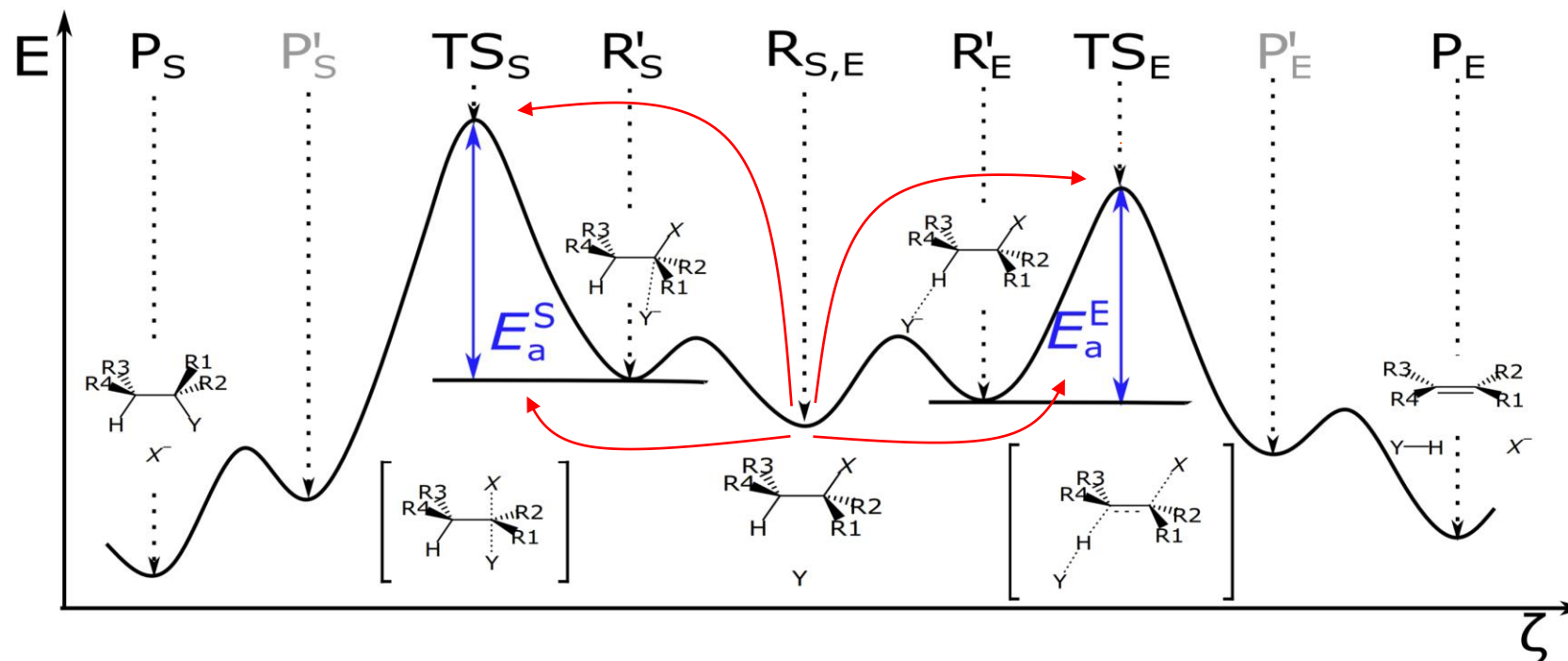
$$k_{ij} \equiv k(d_{ij})$$

## Procedure

- Get $i$ data points with scalar property (label) $\qquad \{q_i\}$
  - E.g. atomisation energy
- Calculate all representations $\qquad\qquad\qquad \{\mathbf{M}_i\}$
  - typically ~1k
- Find distance and kernel matrices $\qquad\qquad \mathbf{D}, \mathbf{K}$
  - Symmetric
- Train model for predictions $\qquad\qquad\qquad \{\tilde{q}_i\}$
- Find best hyperparameters (cross-validation)

$$\arg\min_{\alpha} \sum_i (q_i - \tilde{q}_i)^2 + \lambda \sum_{ij} \alpha_i \alpha_j k_{ij}$$

$$\Rightarrow \alpha = (\mathbf{K} + \lambda \mathbf{I})^{-1} \qquad\qquad \tilde{q}(\mathbf{M}) = \sum_i \alpha_i k(\mathbf{M}, \mathbf{M}_i)$$
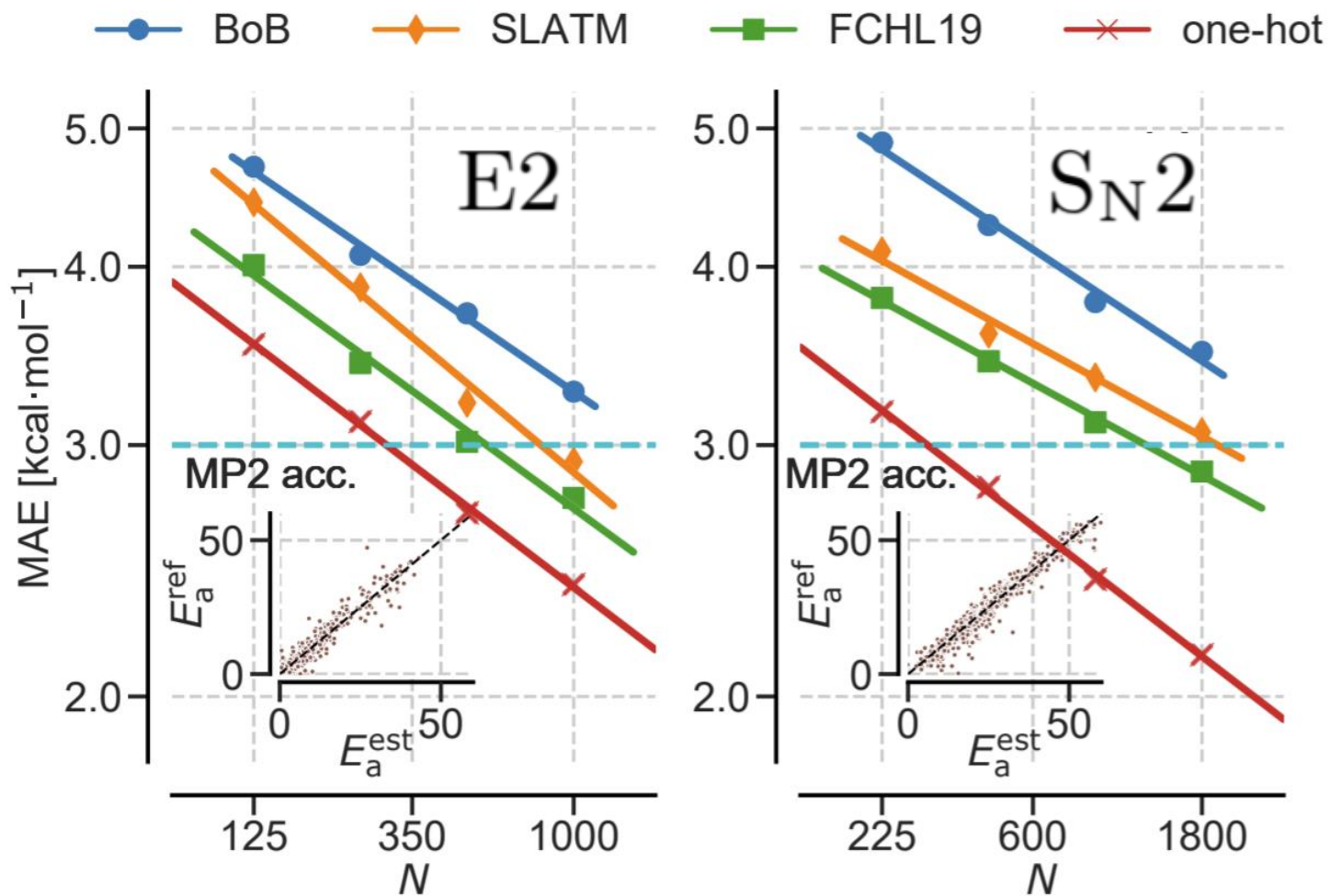
- Activation energies $E_a$
- Transition state geometries
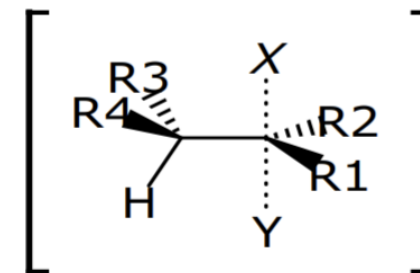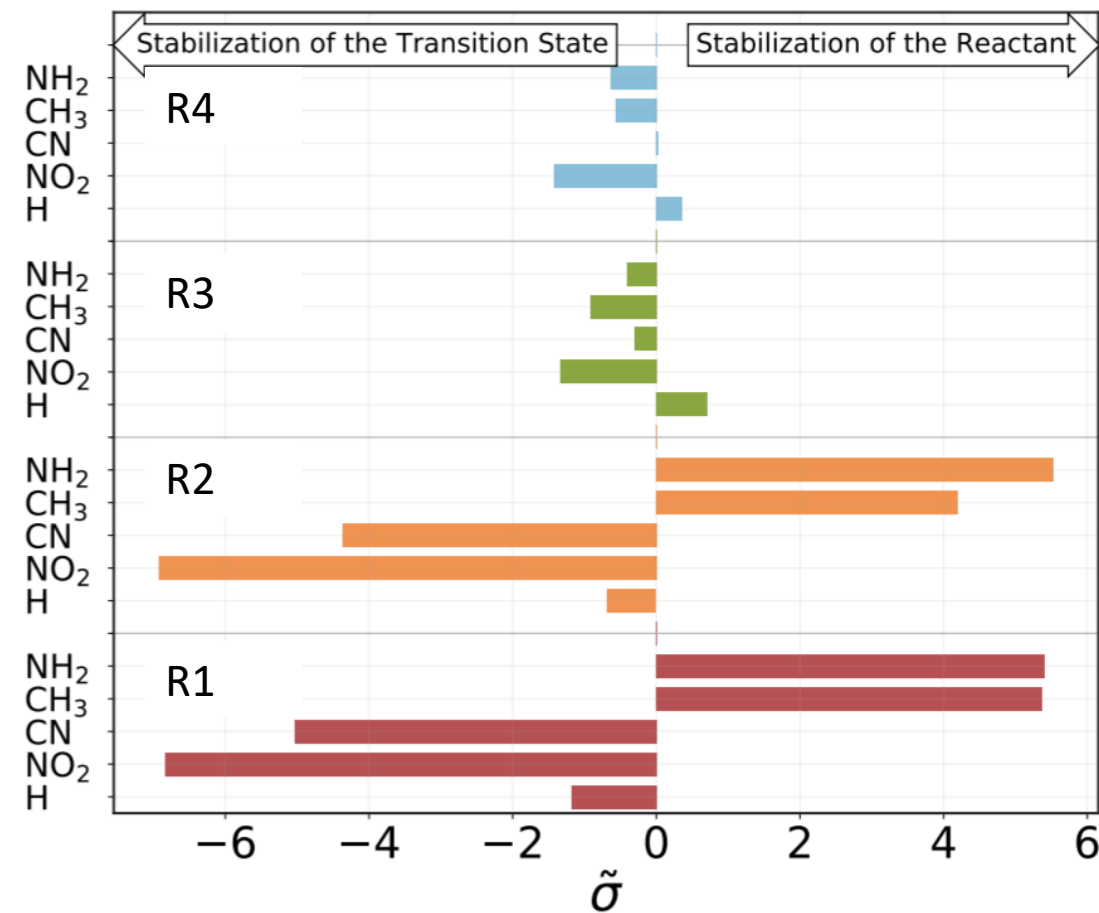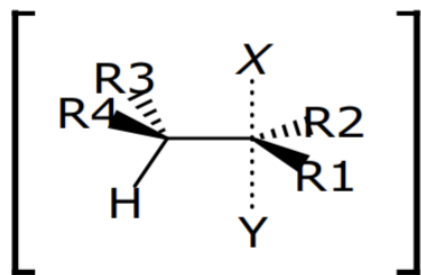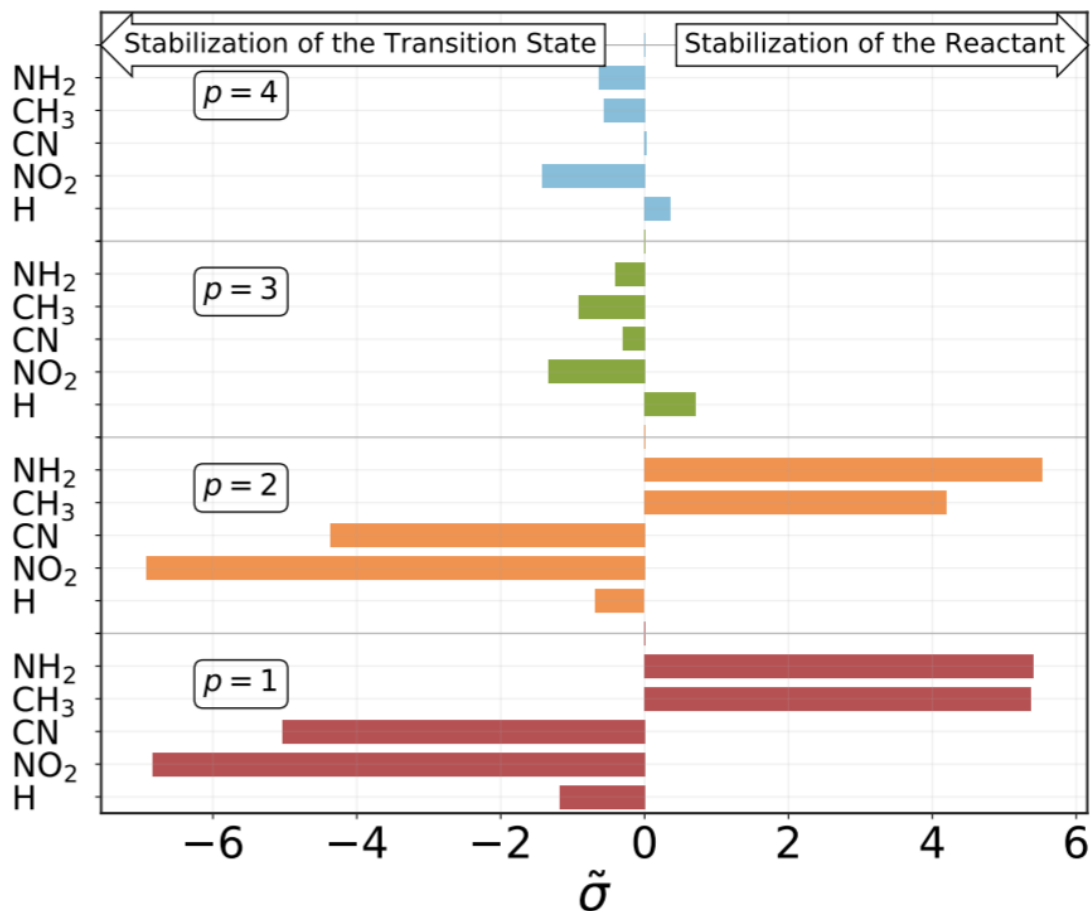- Dataset of 4.5k transition states, 143k reactant geometries, part MP2, part DF-LCCSD

- Geometry-based representations on lowest conformer
  - BoB
  - SLATM
  - FCHL19
- Graph-based representations
  - One-hot

S. Heinen, GFvR, O. A. von Lilienfeld, *arXiv* 2009.13429.

- Often in chemistry: trends obscure relevant detail
  - Electron density dominated by individual atoms
  - Energies dominated by elemental composition
  - Bond energies dominated by element pairs
  - …



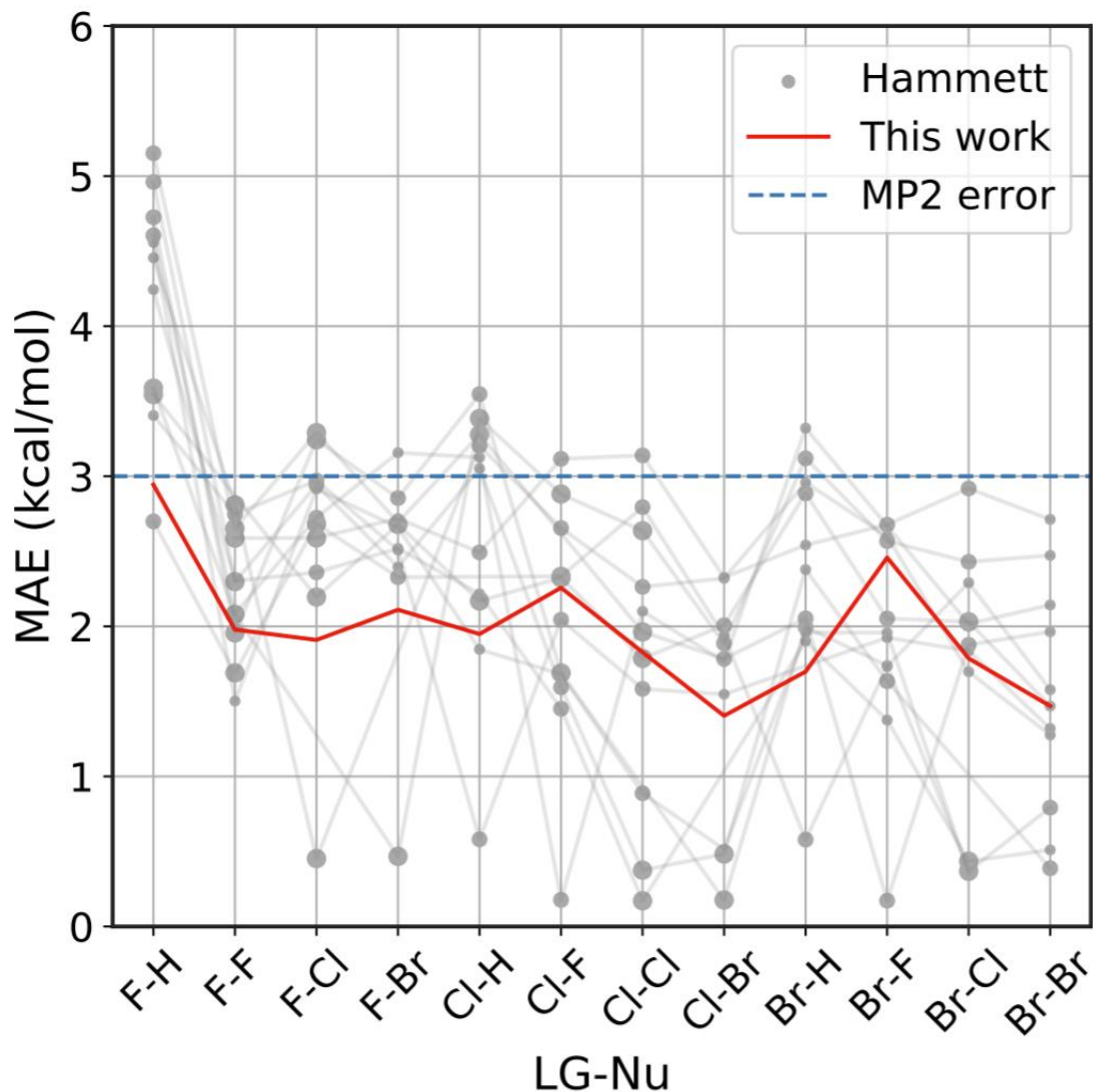M. Bragato, GFVR, O. A. von Lilienfeld, *Chem. Sci. 2020 (arXiv 2004.14946).*

Hammett's equation (1935):

$$\log\left(\frac{K}{K_0}\right) \simeq \rho\sigma$$

Can be used to remove linear trends in the data

1. Find two aspects (e.g. solute/solvent) that are orthogonal and approximately balanced in the data set
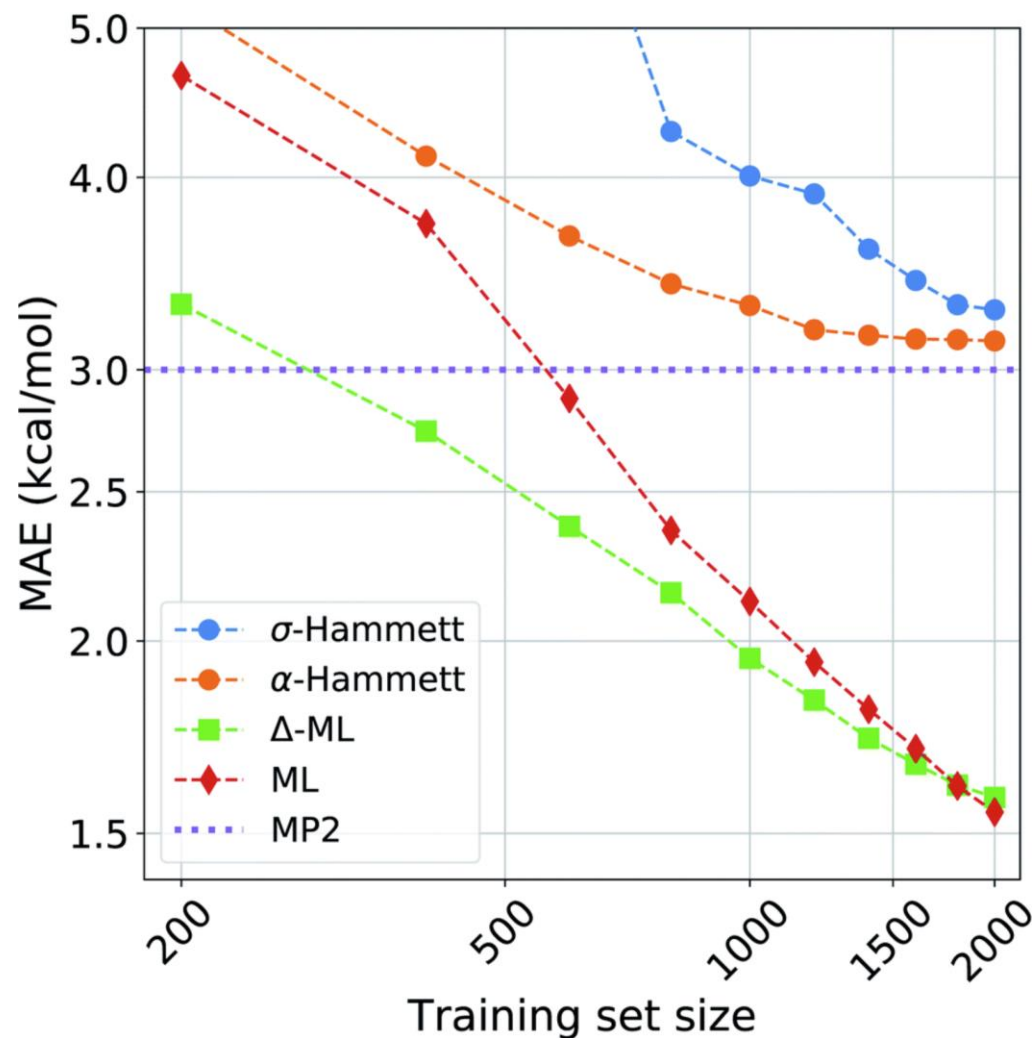2. Fit rho, sigma in a robust manner

Hammett's equation (1935):

$$\log\left(\frac{K}{K_0}\right) \simeq \rho\sigma$$

Can be used to remove linear trends in the data

1. Find two aspects (e.g. solute/solvent) that are orthogonal and approximately balanced in the data set
2. Fit rho, sigma in a robust manner

M. Bragato, GFVR, O. A. von Lilienfeld, *Chem. Sci. 2020 (arXiv 2004.14946).*
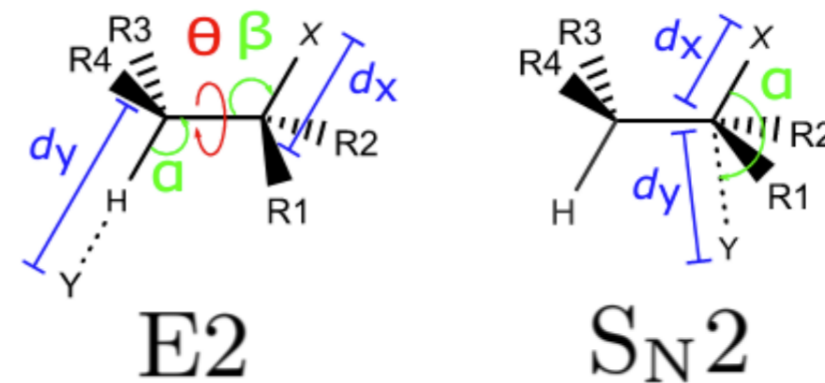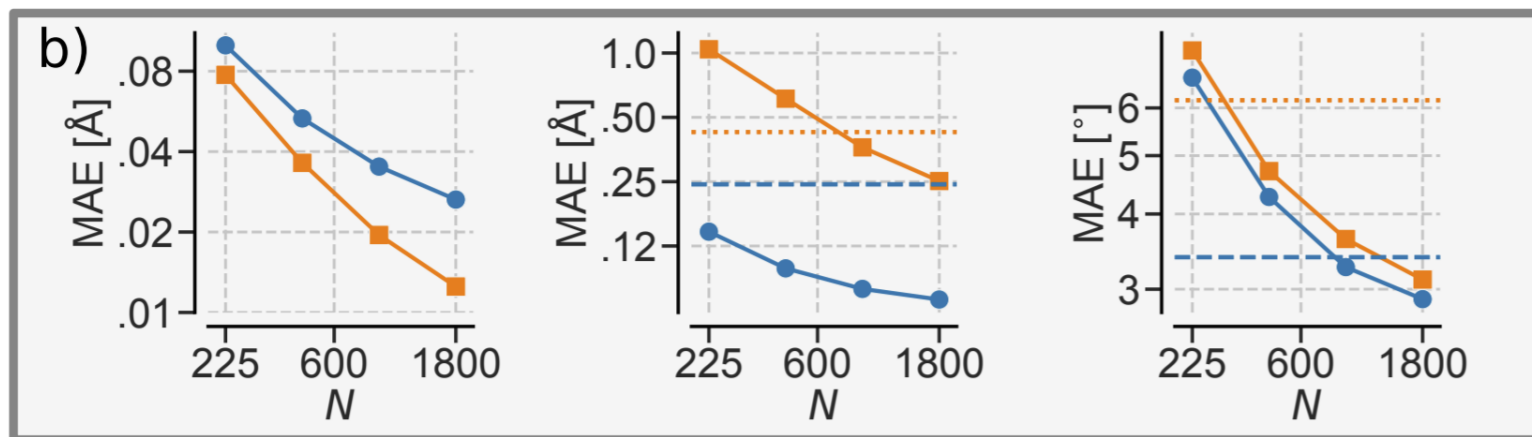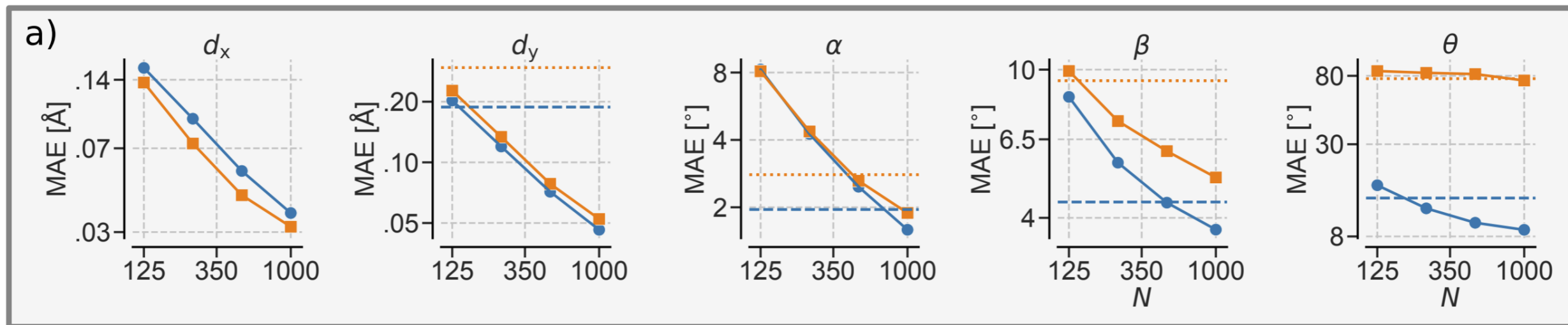
- Hammett nearly reaches MP2 accuracy
- Residuals are easier to learn
- Preprocessing of datasets most helpful for small training sets
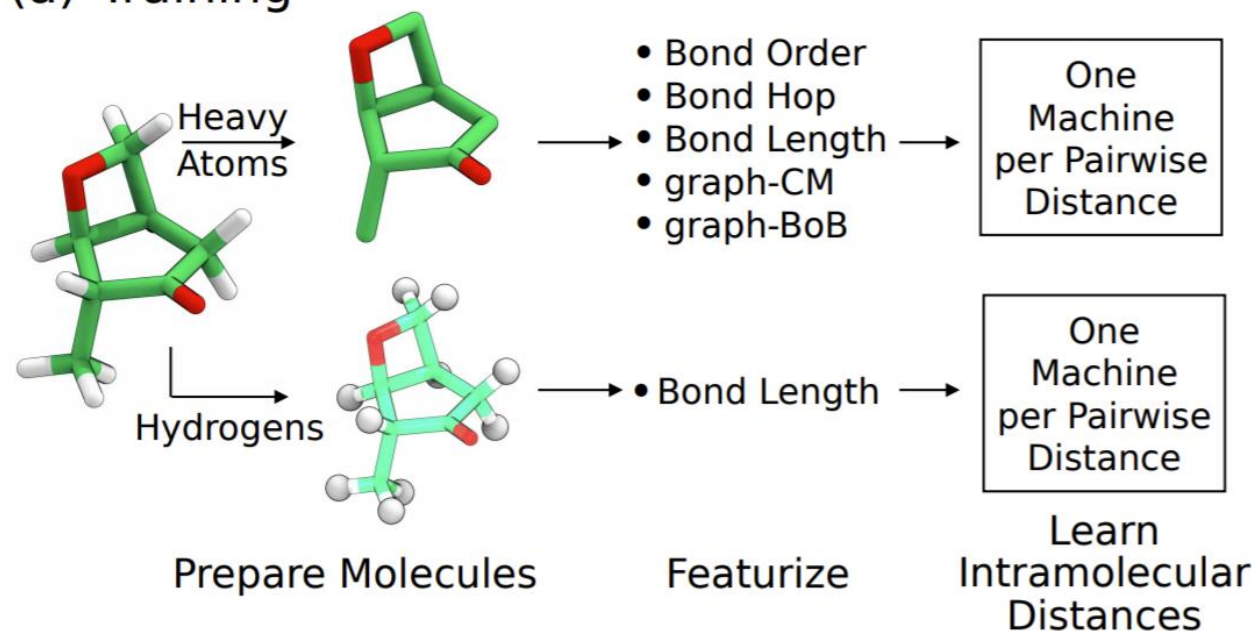
chemspacelab/Enhanced-Hammett

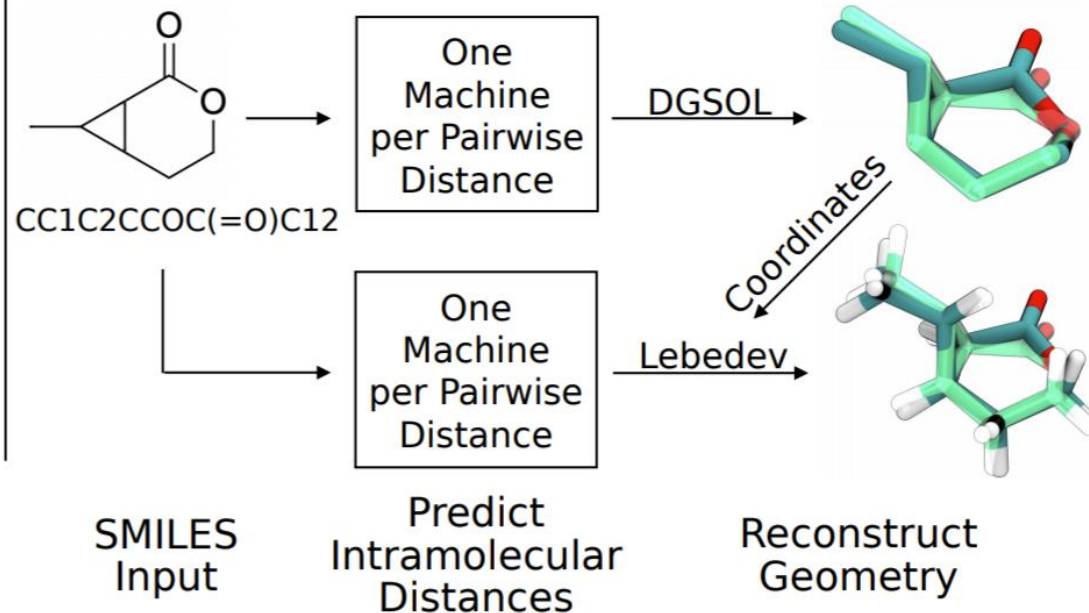M. Bragato, GFVR, O. A. von Lilienfeld, *Chem. Sci. 2020 (arXiv 2004.14946).*

(a) Training

Prediction

Prepare Molecules → Featurize → Learn Intramolecular Distances

Heavy Atoms

Hydrogens

- Bond Order
- Bond Hop
- Bond Length
- graph-CM
- graph-BoB

One Machine per Pairwise Distance

- Bond Length

One Machine per Pairwise Distance

SMILES Input

CC1C2CCOC(=O)C12

One Machine per Pairwise Distance → DGSOL

One Machine per Pairwise Distance → Lebedev

Predict Intramolecular Distances

Coordinates

Reconstruct Geometry

(a)

- Learns standard chemistry, but also carbenes, transition state geometries
- More accurate w.r.t. to QM calculations than state-of-the-art embedding methods (which only do standard chemistry)
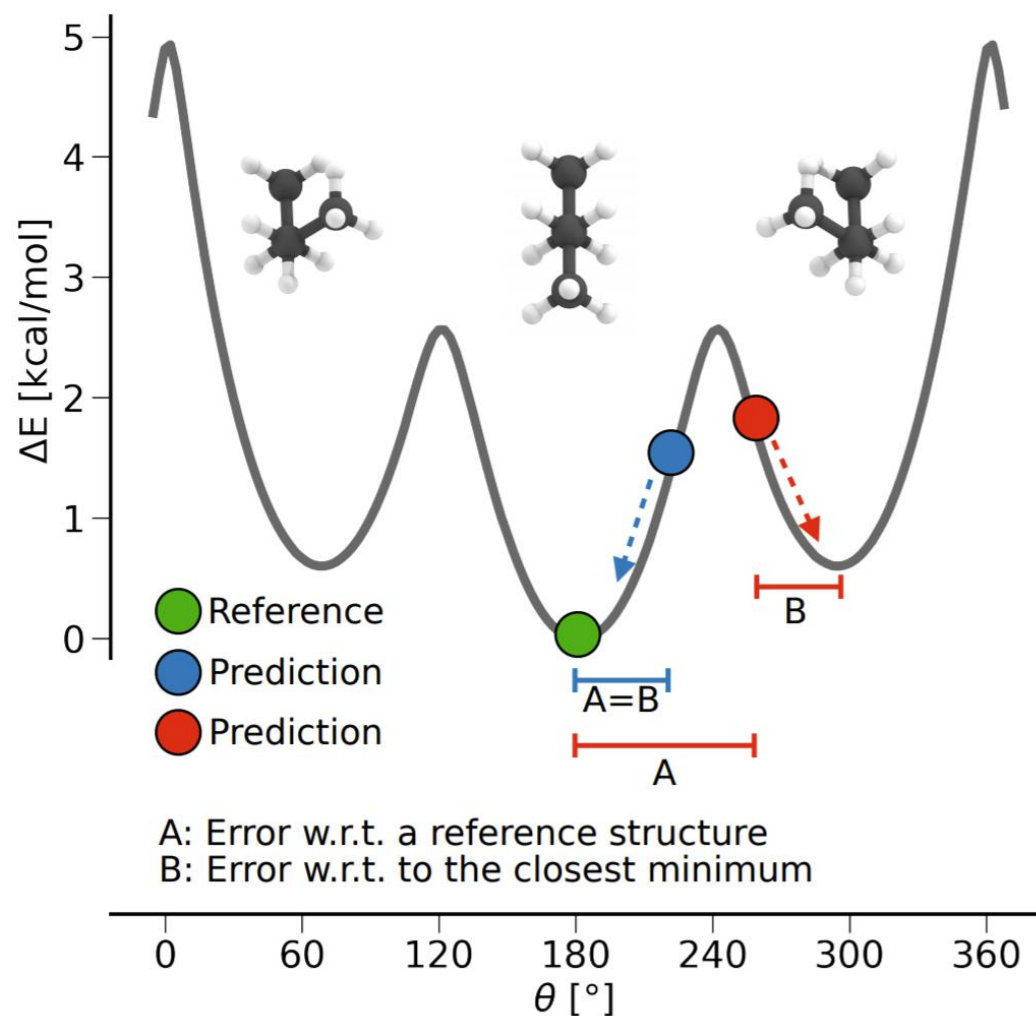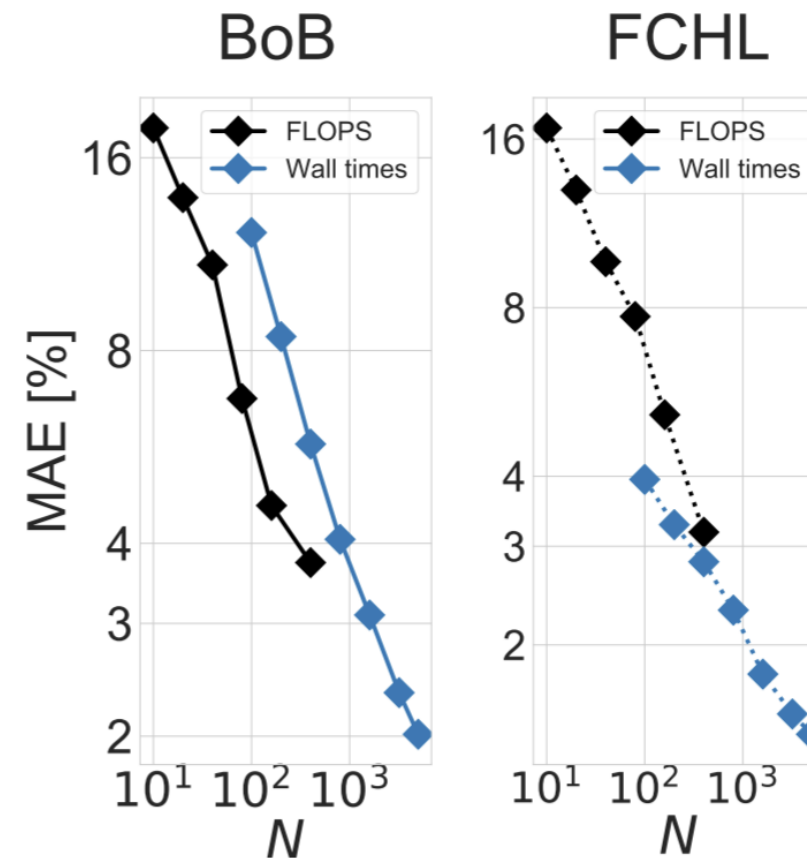- Can produce initial guesses for e.g. transition state searches

- Most efficient for cases with wide minima
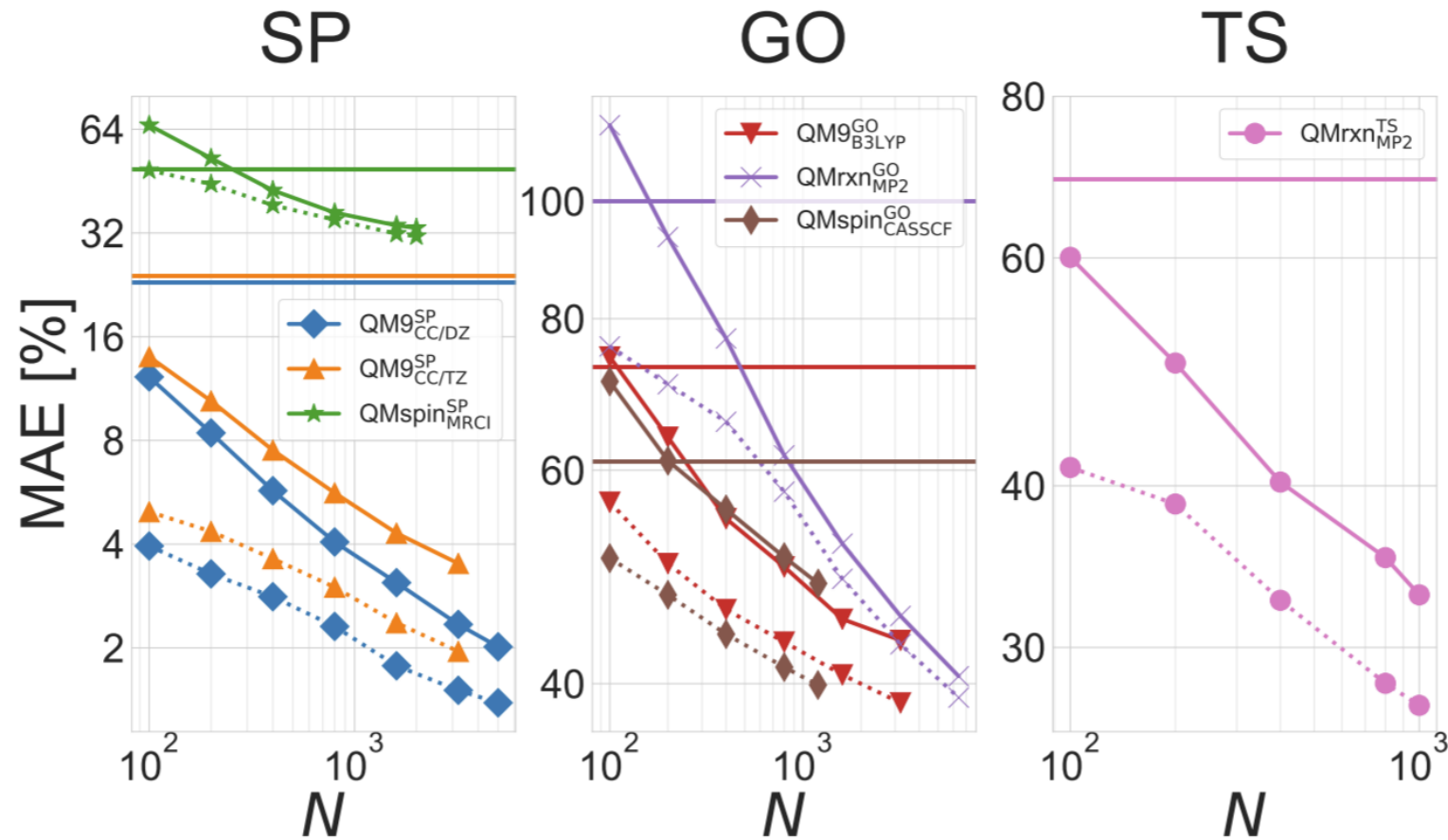- Multiple small minima e.g. stereoisomers account for most of the error

qmlcode/qml

- Training data expensive
- Not all training points equally expensive
  - Geometry optimisations may take longer
  - SCF might not converge or converge more slow
- Treating this as "molecular property"


- Best case: controlled environment, single points (QM9)



S. Heinen, M. Schwilk, GFvR, O. A. von Lilienfeld, *Mach. Learn.: Sci. Technol.* 2020 (*arXiv* 1908.06714).

- Realistic case: I/O, noise: different machines



ferchault/mlscheduling

S. Heinen, M. Schwilk, GFvR, O. A. von Lilienfeld, *Mach. Learn.: Sci. Technol.* 2020 (*arXiv* 1908.06714).

chemspacelab/Enhanced-Hammett

qmlcode/qml

ferchault/mlscheduling

...

- Dependencies might break
- Might be an old model
- Tedious/risky to get started
- Therefore:

leruli.com

BETA

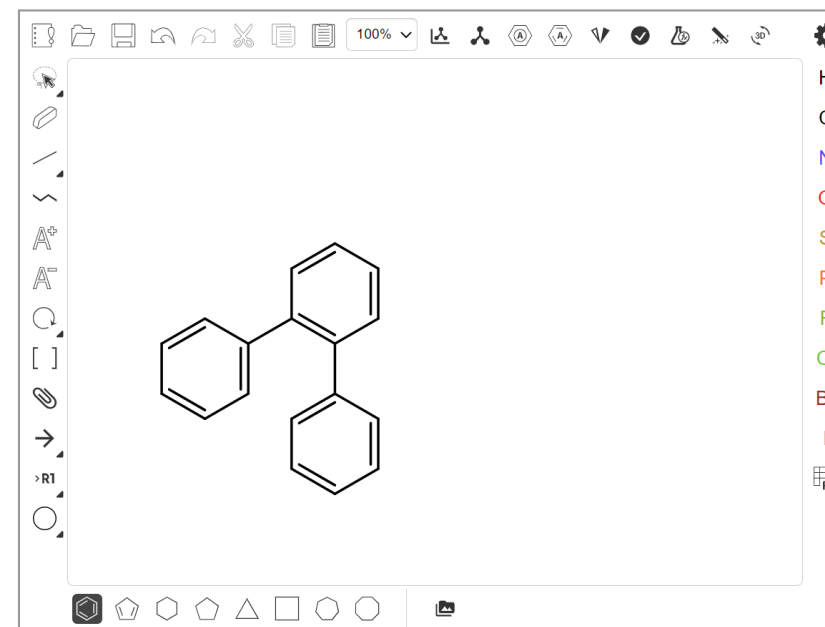Search Sum Formula, Compound Name, SMILES, SMARTS, SELFIES, InchI

C1=CC=CC=C1C1=C(C2=CC=CC=C2)C=CC=C1        Search

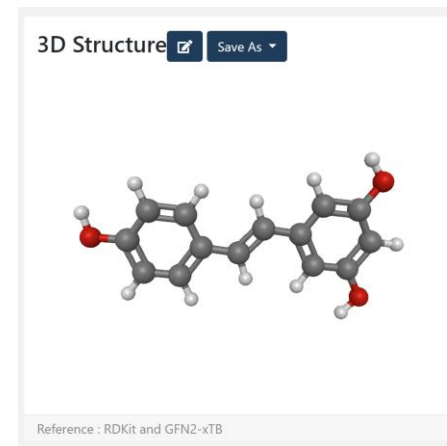Examples:   $C_8O_2H_{18}$   Resveratrol   C1COCCO1

Draw       Upload

100% ∨

H
C
N
O
S
P
F
Cl
Br
I
PT

# Model availability

Want to include your model? Let me know

Estimate computational cost



Detrending with Hammett's equation



Energies with KRR



Geometries with Graph2Structure

ferchault/mlscheduling

chemspacelab/ Enhanced-Hammett

qmlcode/qml

qmlcode/qml

# Acknowledgements

Prof. von Lilienfeld    Dominik Lemm    Marco Bragato    Stefan Heinen    Dr. Max Schwilk

ferchault          @ferchault          guido.vonrudorff.de