

Design in Compound Space With Machine Learning and Quantum Alchemy

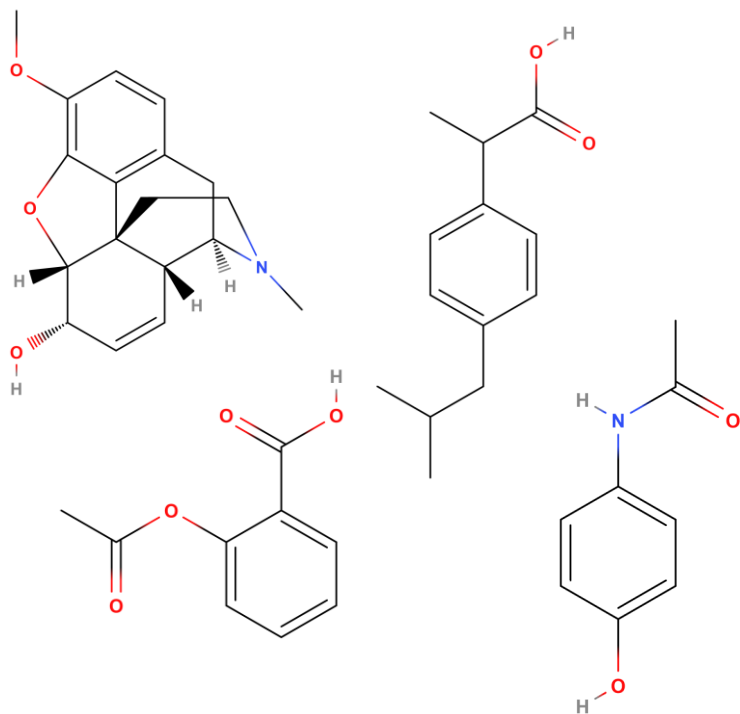
Guido Falk von Rudorff, University of Vienna

 ferchault

 @ferchault

 guido.vonrudorff.de

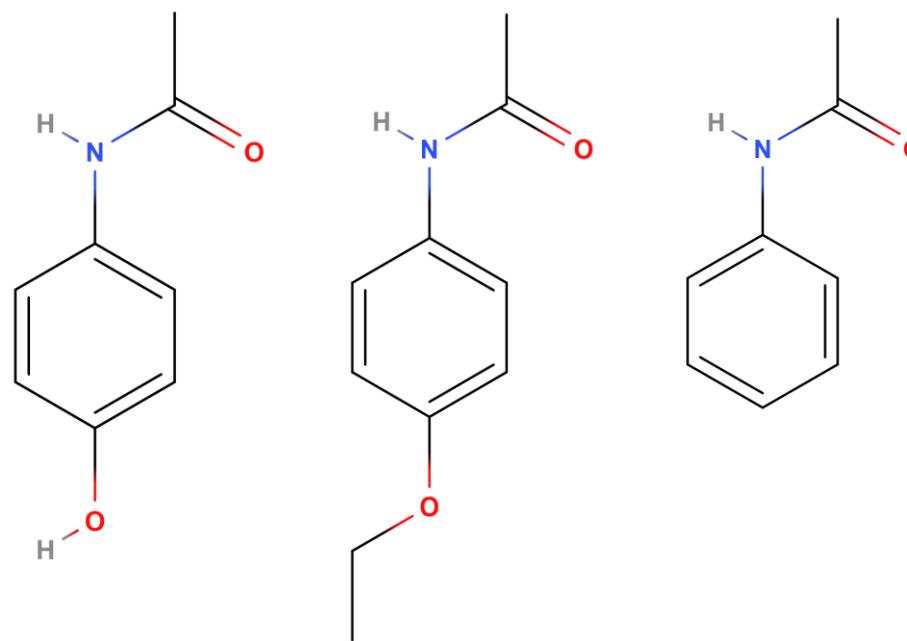
Design: sample by guided trial-and-error.



Global Search Problem

Which class of compounds?

Drug-like: 10^{60}

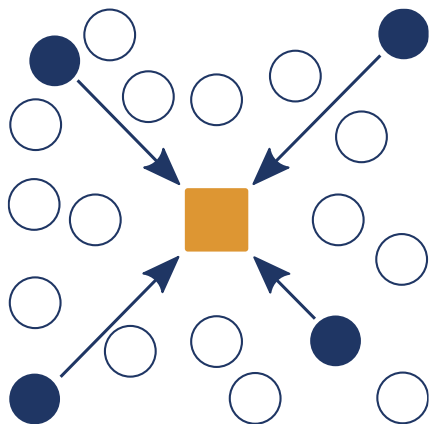


Local Search Problem

Which particular species within that class?

BN-doped 8x8 graphene: 10^{50}

Machine Learning



Foundations | Statistical modelling

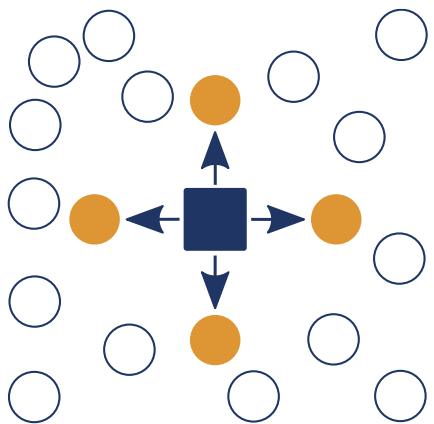
Accuracy | Systematically improvable through data and training

Specialty | Universal, scale-bridging, data-driven approach

Limitation | Requires training data, no black box

Chemistry: exciting new domains become accessible!

Quantum Alchemy



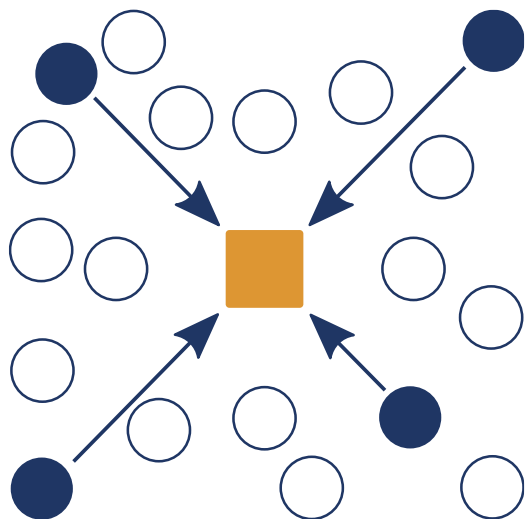
Foundations | Perturbation theory

Accuracy | Systematically improvable through higher orders terms

Specialty | Combinatorial scaling with chemical diversity

Limitation | Finite range in chemical space

Machine Learning



Kernel-Ridge-Regression

- Efficient in the low-data regime (around 1k points)

- Ingredients

- Representation

 \mathbf{M}

- Similarity measure

 $k(\mathbf{M}_i, \mathbf{M}_j)$

- Observed properties

 \mathbf{y}

- Training

- Pairwise similarities

 \mathbf{K}

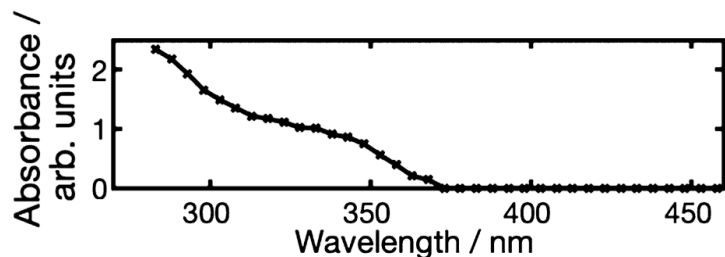
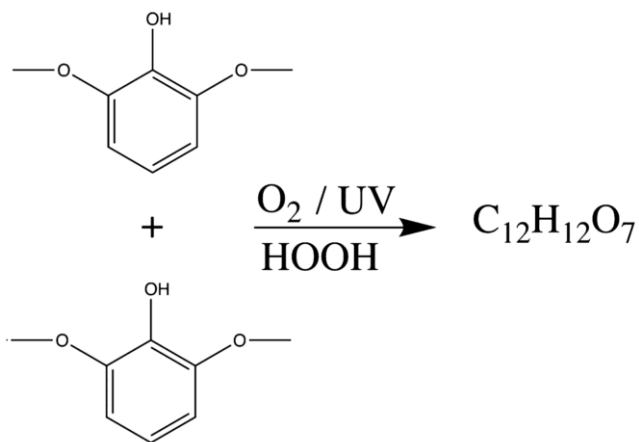
- Model coefficients

 $\boldsymbol{\alpha} = (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{y}$

- Predictions

- Compare to training $\tilde{q}(\mathbf{M}) = \sum_i \alpha_i k(\mathbf{M}, \mathbf{M}_i)$

Experiment

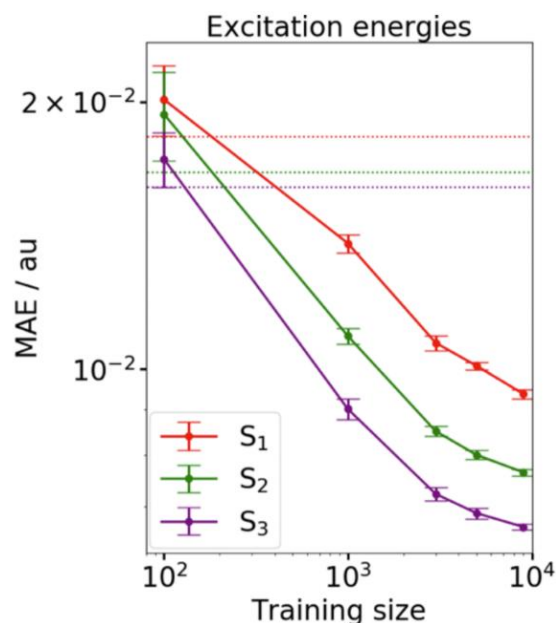


Search space

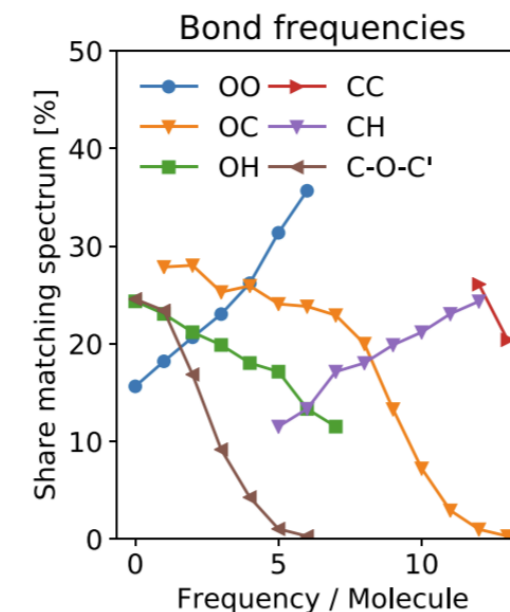
Molecular graphs: 264 M
Stable molecules: 123 M

Model

Excitation energies
Oscillator strengths

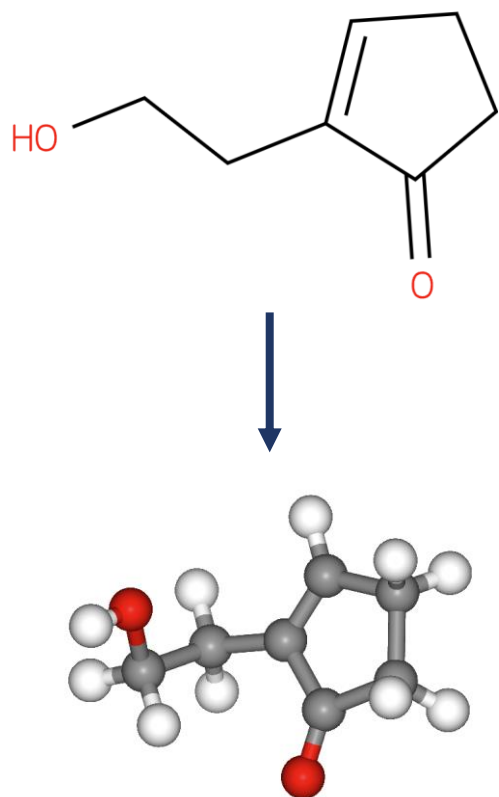


Identified features

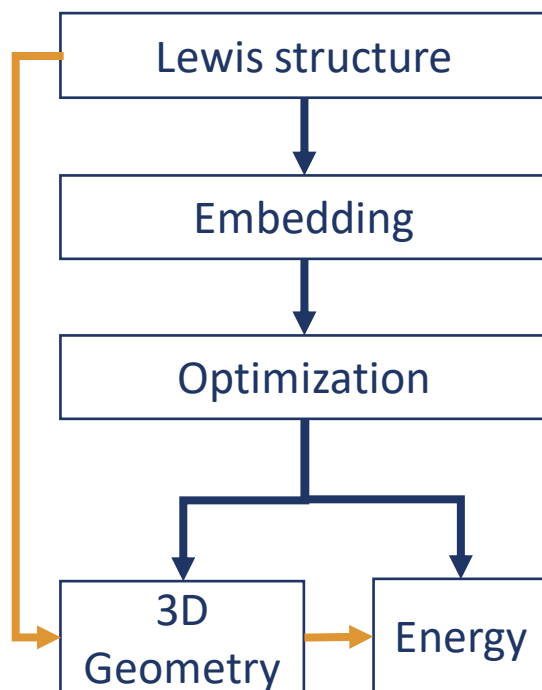


Guide experiment

How many molecules are left and which feature to measure next?



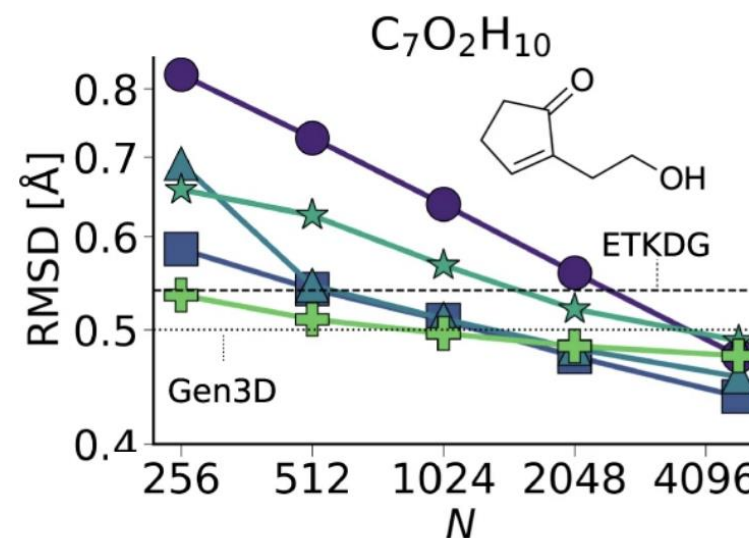
Traditional / G2S

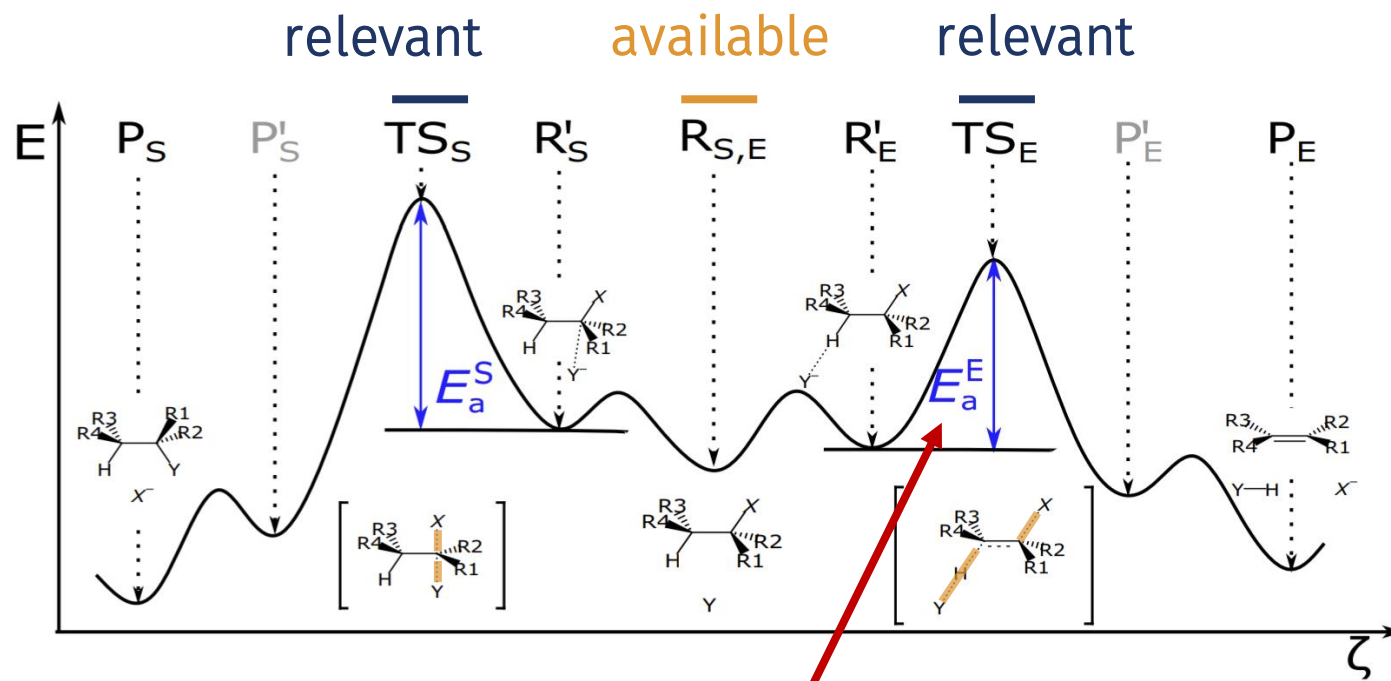


ML: 100.000x faster

G2S

- Closer to DFT than common methods
 - Small molecules
- Applicable to complex chemical spaces
 - Transition state geometries
 - Carbenes





ML: 10.000x faster

Competing reactions: E2, S_N2

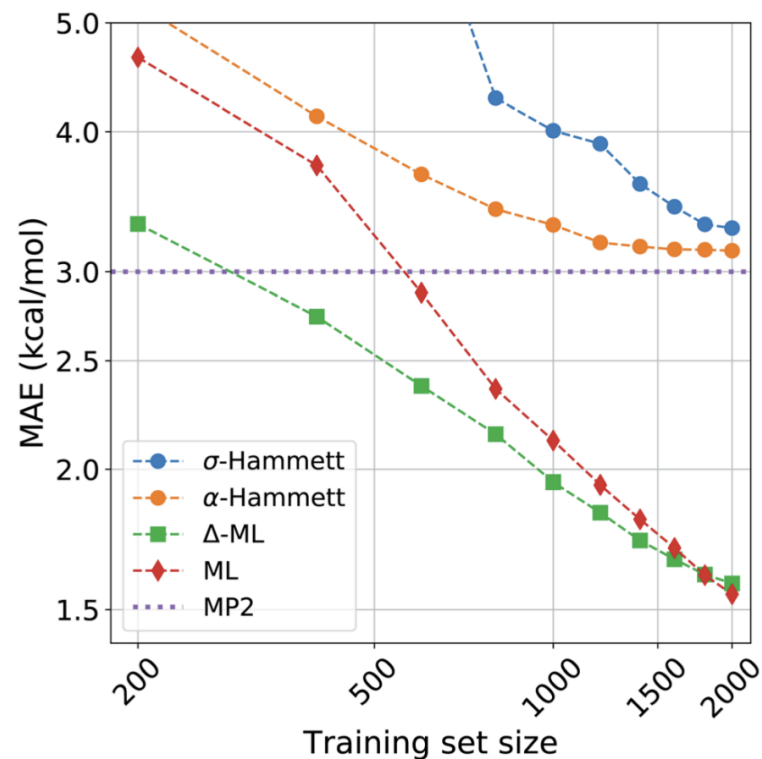
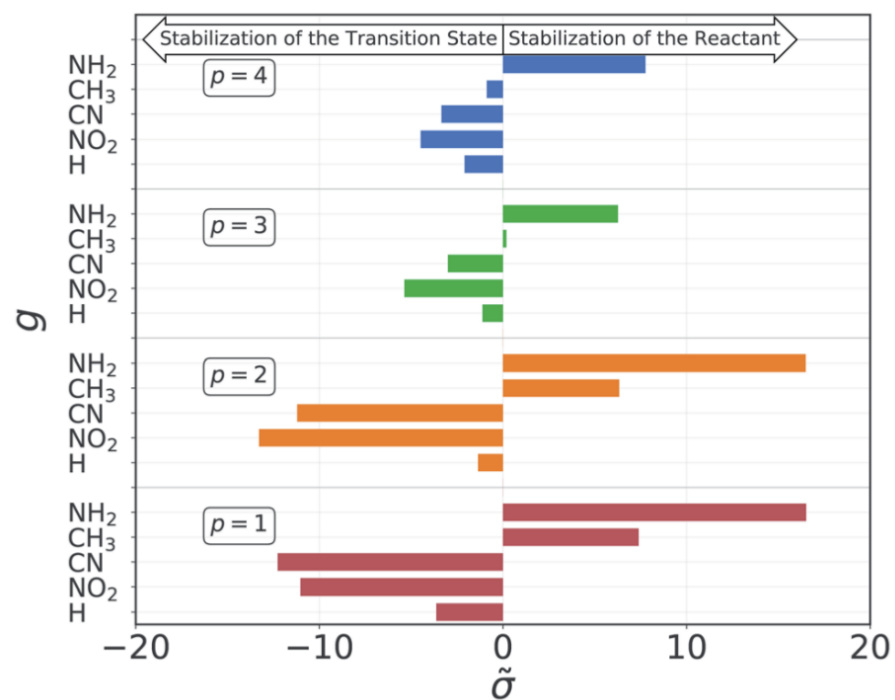
- 4.5k reactions in one new dataset
- Learning activation energies from reactants only reaching 2.5 kcal/mol with 800 data points
- Learning geometries of transition states
 - direct
 - 0.05 Angstrom for distances
 - G2S
 - 0.45 Angstrom heavy-atom RMSD

Hammond's postulate

- valid for S_N2 , not for E2
- E2 preferred in 75% of the cases: less sensitive to reactant complex geometry

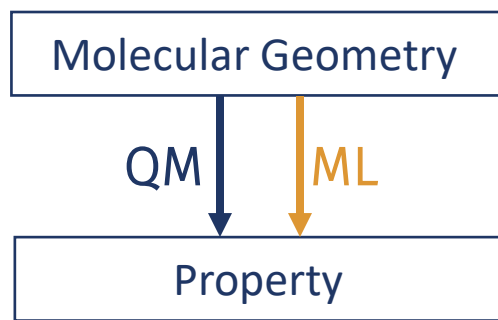
GFvR, S.N. Heinen, M. Bragato, O. A. von Lilienfeld, *Mach. Learn.: Sci. Technol.* 2020.

S.N. Heinen, GFvR, O. A. von Lilienfeld, *J. Chem. Phys.*, 2021.



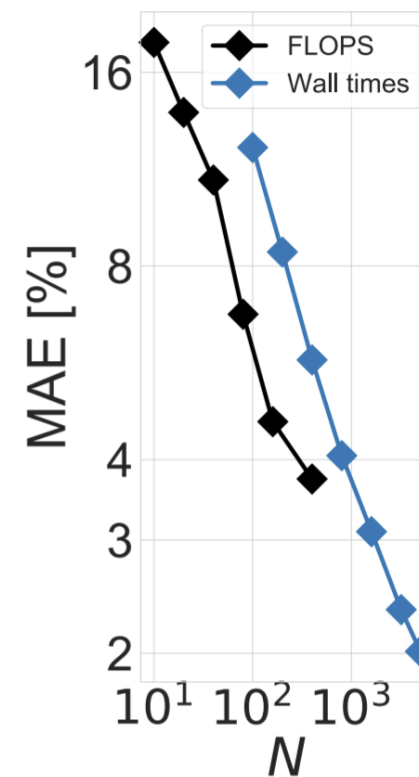
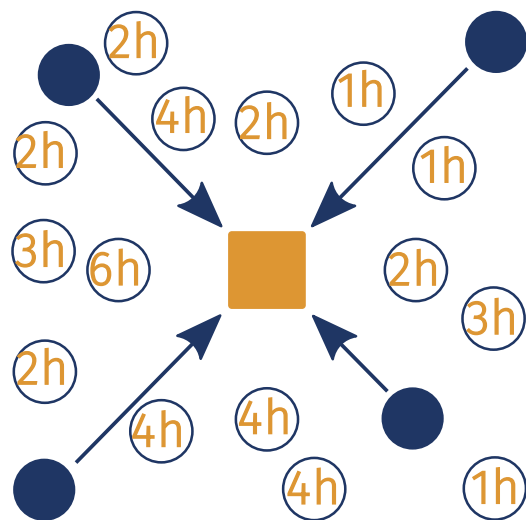
Hammett equation
 Electron withdrawing/-
 donating order of
 substituents
 recoverable


- Physics-based model
- Use Hammett as baseline



Computational effort as molecular property

- Improves models
- Accuracy depends on problem
 - Single points: 2%
 - Transition state search: 25%
 - Geometry optimisations: 40%



 chemspacelab/Enhanced-Hammett

 qmlcode/qml

 ferchault/mlscheduling

 ...

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

leruli.com

Leruli

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

Search

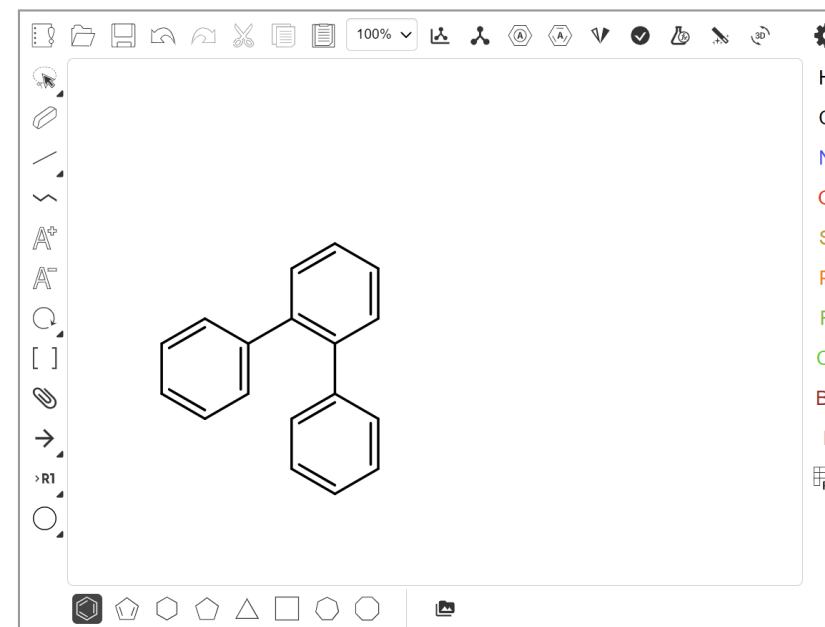
Examples: [C₈O₂H₁₈](#) [Resveratrol](#) [C1COCCO1](#)



Draw



Upload



Leruli

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

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

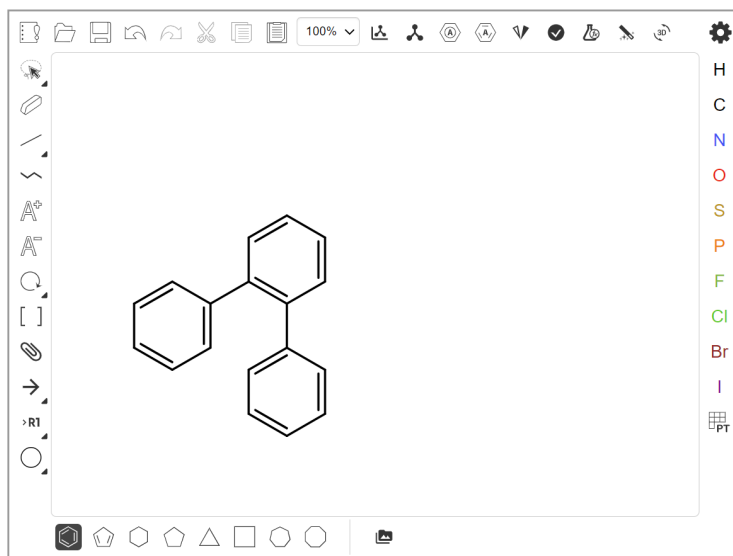
Examples: C8O2H18 Resveratrol C1COCCO1



Draw



Upload



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

Share your results by copying the URL or clicking

Lewis Structure

Reference : Indigo

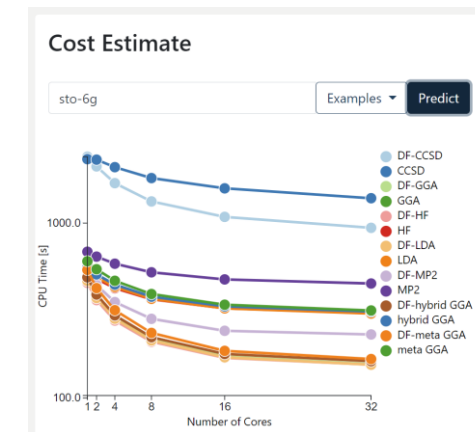
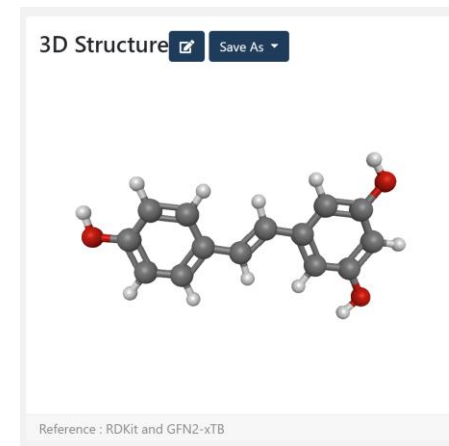
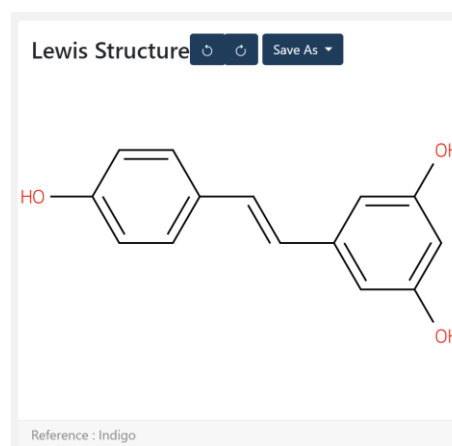
Drug Likeness	Score
Lipinski	5/5
Ghose	4/4
Weber	2/2
Rheos	7/7
Rule of 3	3/5
Drug-Like	5/6

Lipinski	Score
Molecular Weight <= 500	228.25 g/Mol
logP <= 5.0	3.3
HBD <= 5	3
HBA <= 5	3
Rotatable Bonds <= 5	5

Ghose	Score
160 <= Molecular Weight <= 480	228.25 g/Mol
-0.4 <= logP <= 5.6	3.3
20 <= Atoms <= 70	29
40 <= Molecular Refractivity <= 130	66.58

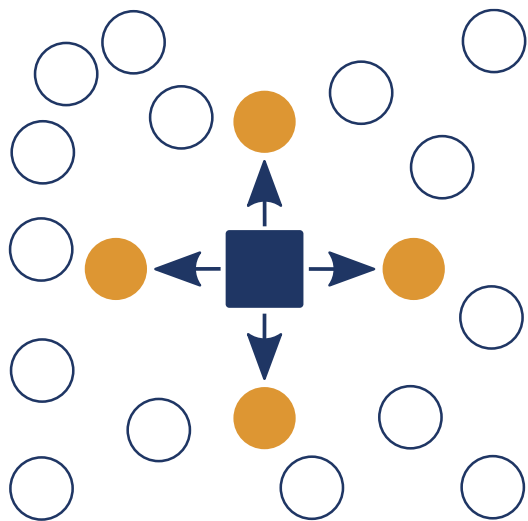
Weber	Score
Topo. Polar Surface Area <= 140	60.69
Rotatable Bonds <= 10	5

REOS	Score
200 <= Molecular Weight <= 500	228.25 g/Mol
-5.0 <= logP <= 5.0	3.3
0 <= HBD <= 5	3
0 <= HBA <= 10	3
0 <= Rotatable Bonds <= 8	5
-2 <= Formal Charge <= 2	0
15 <= Heavy Atoms <= 50	17



Want to include your model? Let me know

Quantum Alchemy

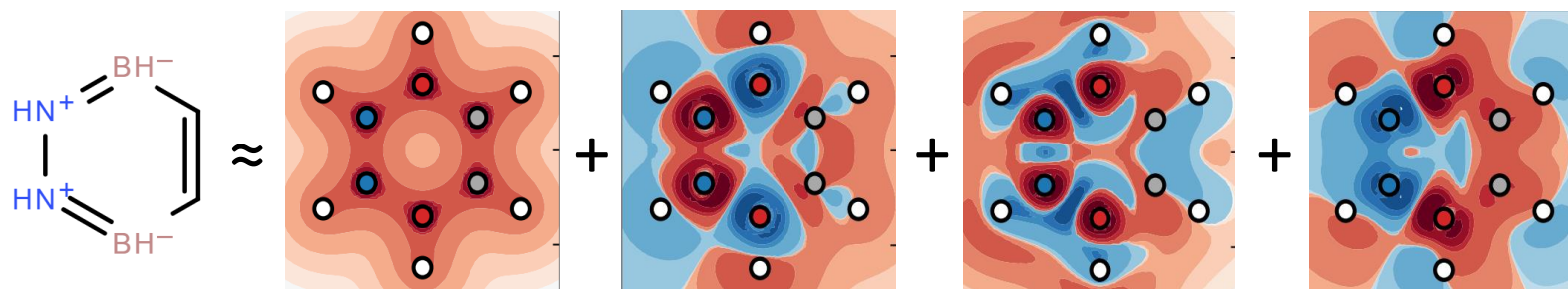


Taylor expansion

- Energy function of
 - Geometry Forces, Vibrations
 - Nuclear charges Alchemical changes
- Idea: obtain dominant leading derivatives, predict many systems

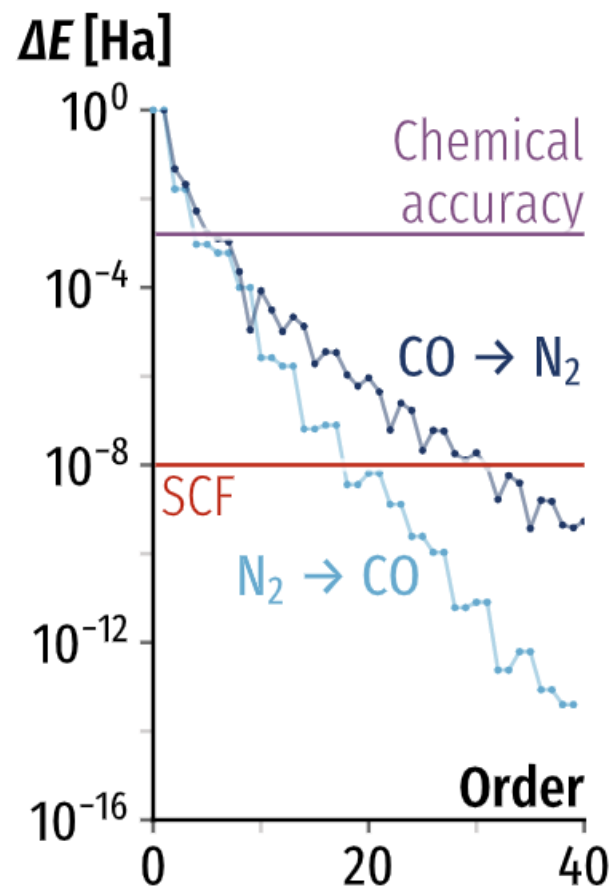
$$\hat{H}(\lambda) \equiv \lambda \hat{H}_t + (1 - \lambda) \hat{H}_r \quad \lambda \in [0, 1]$$

$$E_t = E_r + \sum_{n=1}^{\infty} \frac{1}{n!} \left. \frac{\partial^n E(\lambda)}{\partial \lambda^n} \right|_{\lambda=0}$$



E. B. Wilson, *J. Chem. Phys.* 1962.

GFvR, O. A. von Lilienfeld, *Phys. Rev. Res.*, 2020.

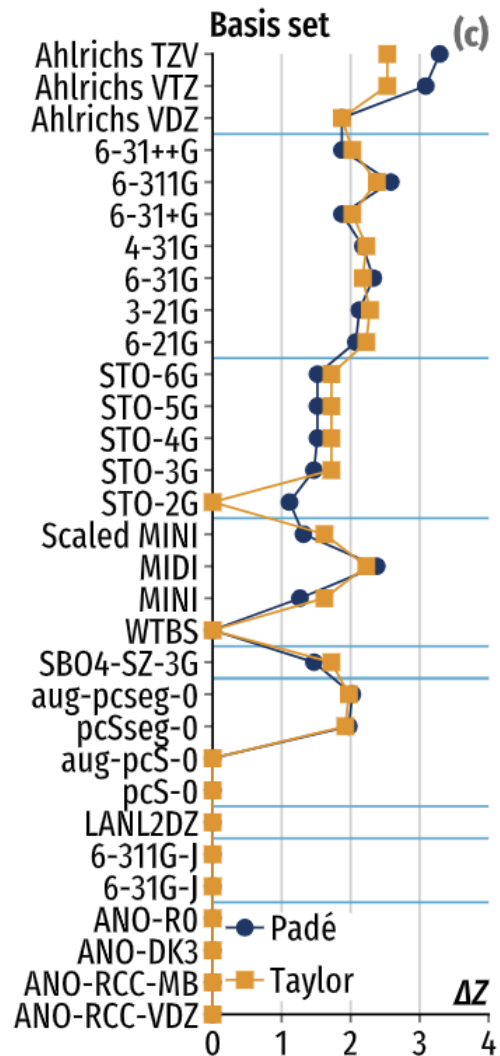
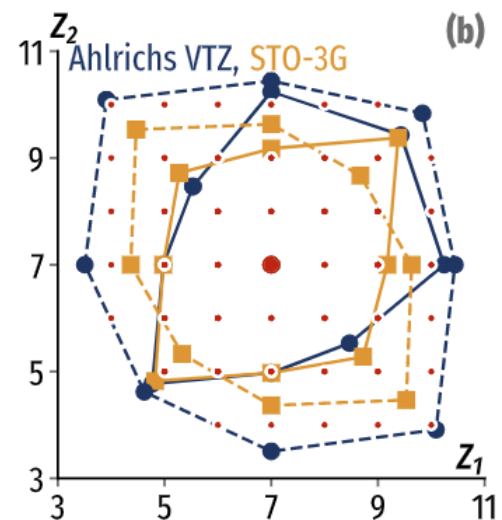
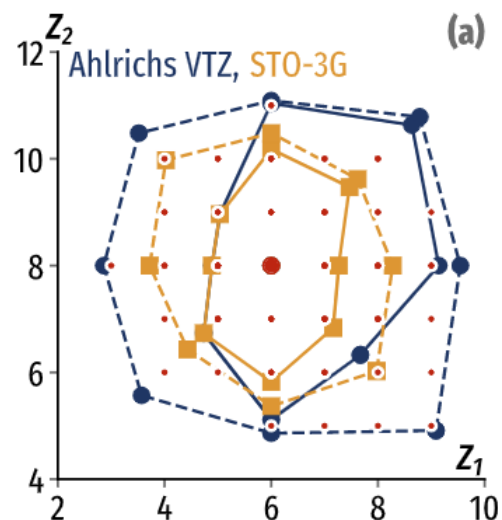


Taylor expansion

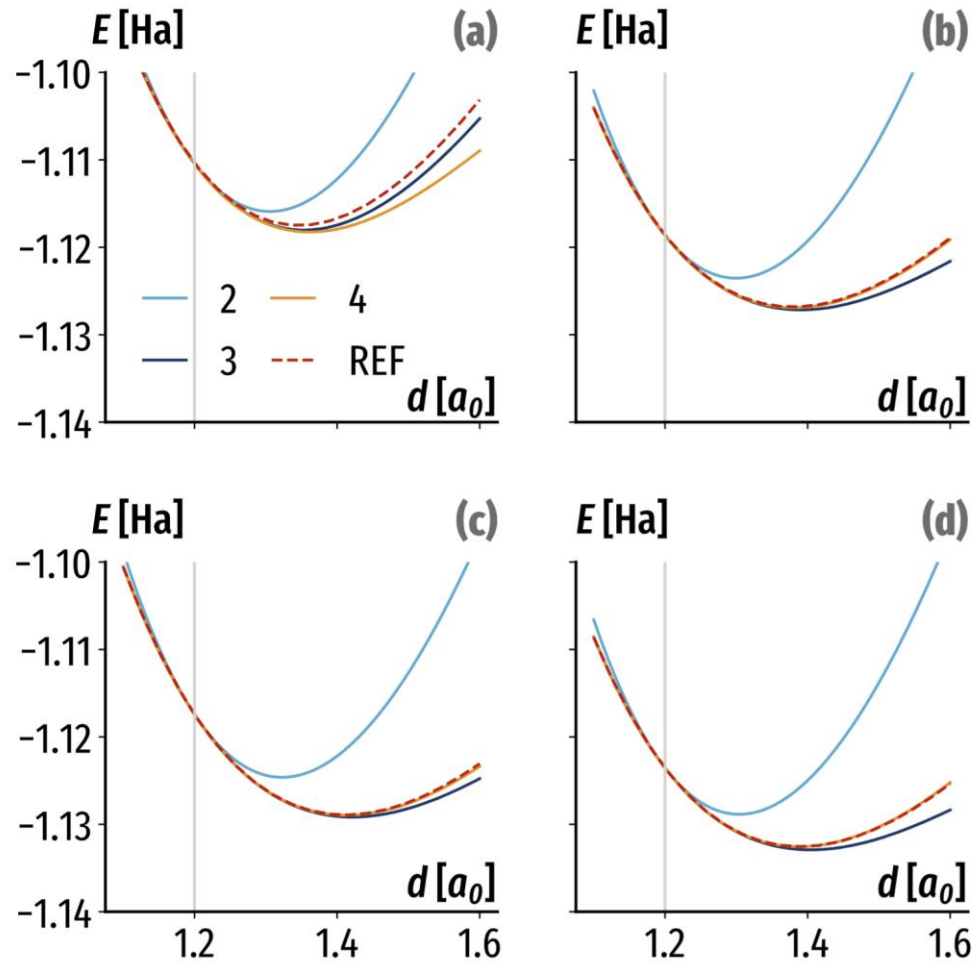
- First terms accurate enough
 - Truncate early
- Scales with chemical space

Paradigm shift

Few highly accurate calculations instead of many intermediate ones



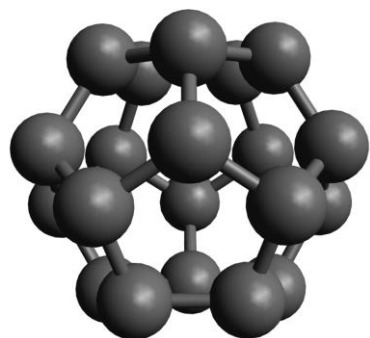
Padé approximant
Larger convergence radius



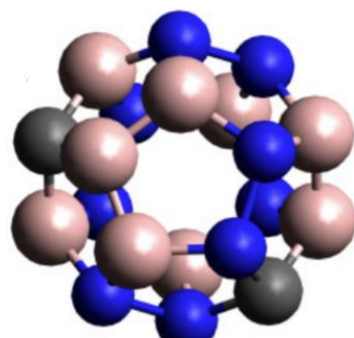
Mixed derivatives
A path to relaxation

Scaling with chemical space

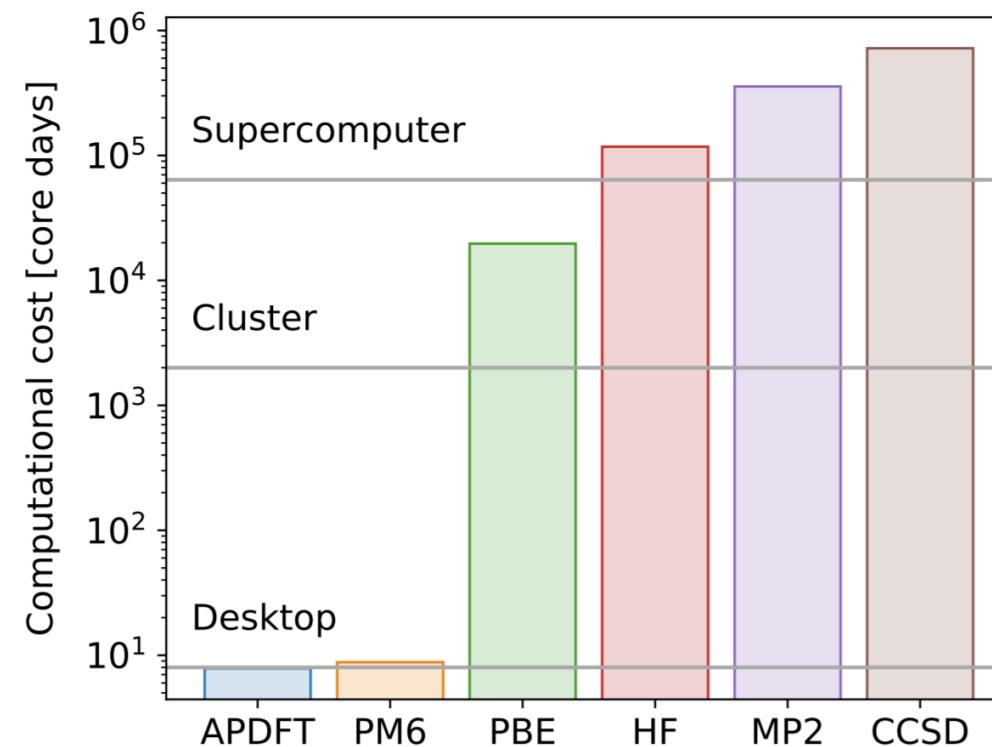
- 1 derivative for second order
- 5 derivatives for third order



C_{20}



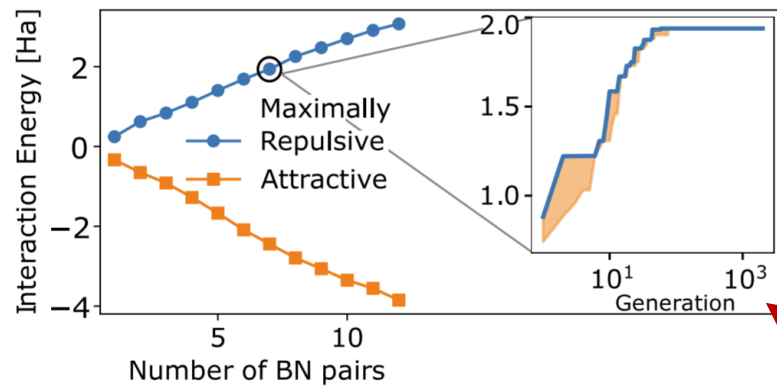
$3.1 \cdot 10^6$
targets



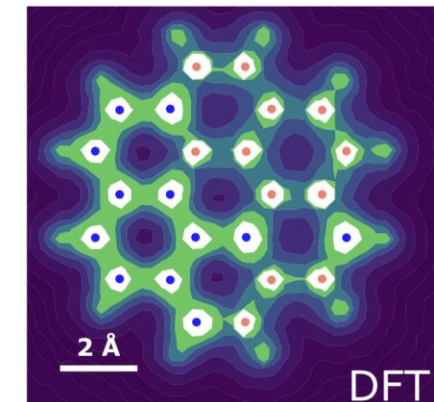
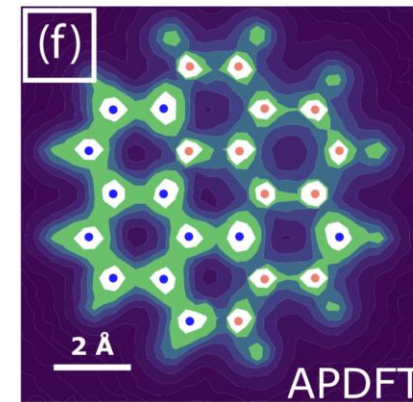
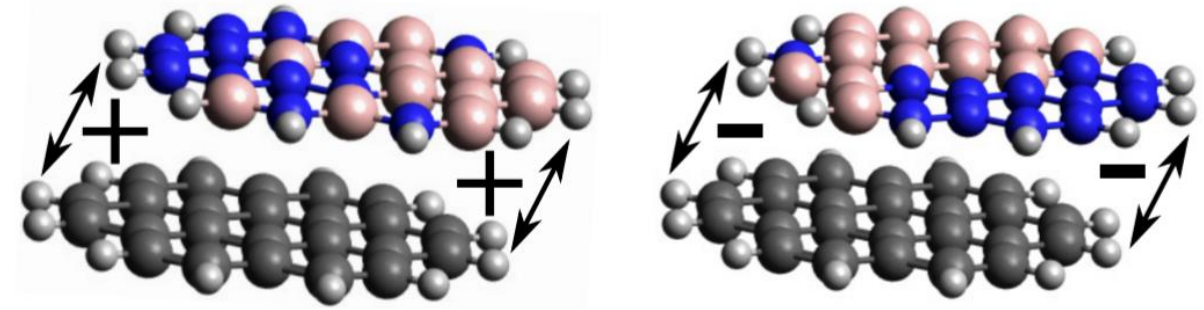
QA: 80.000x faster

BN-doped coronene dimer

- Identify most/least attractive doping pattern
- Design example!



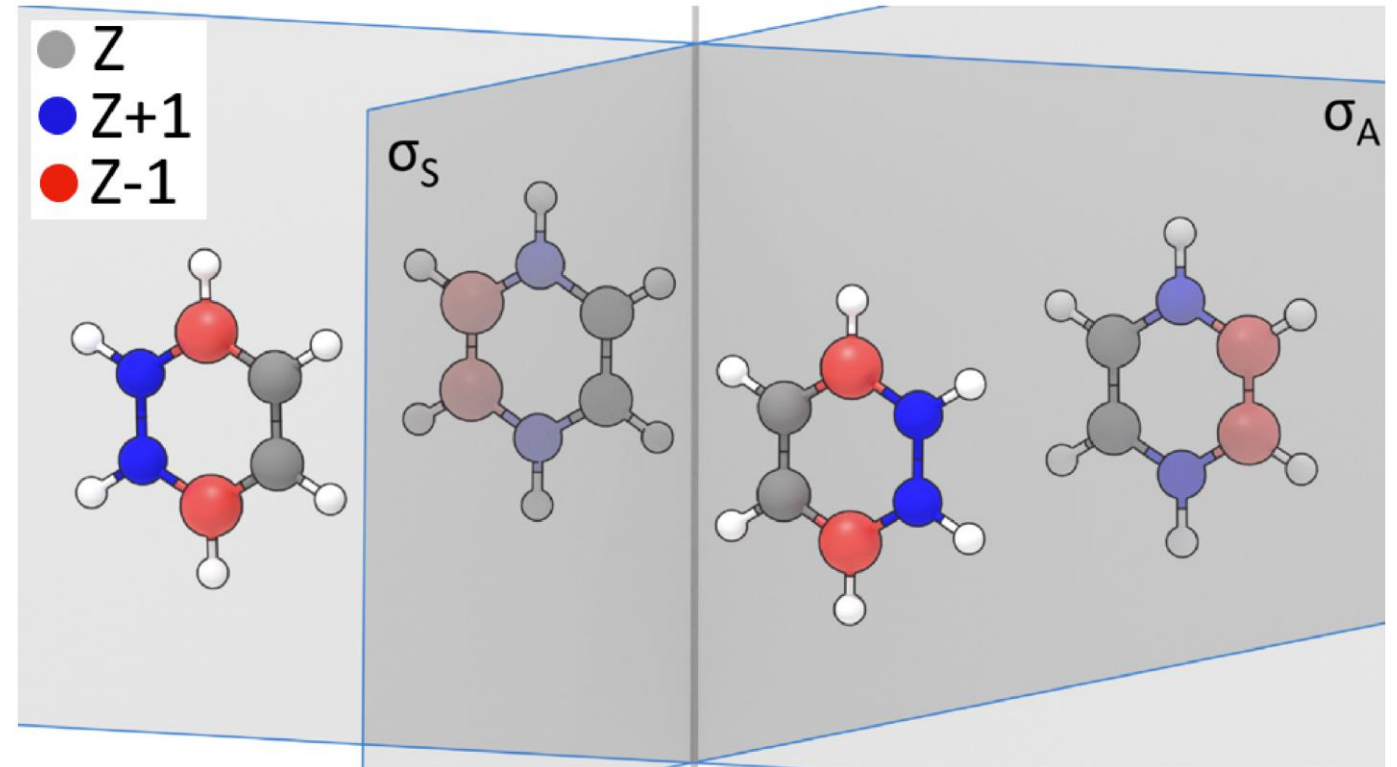
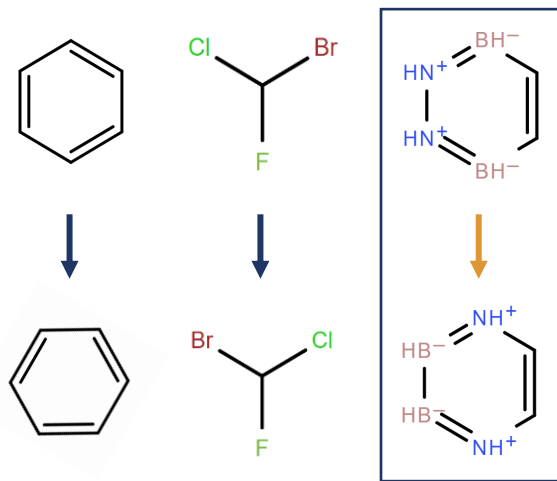
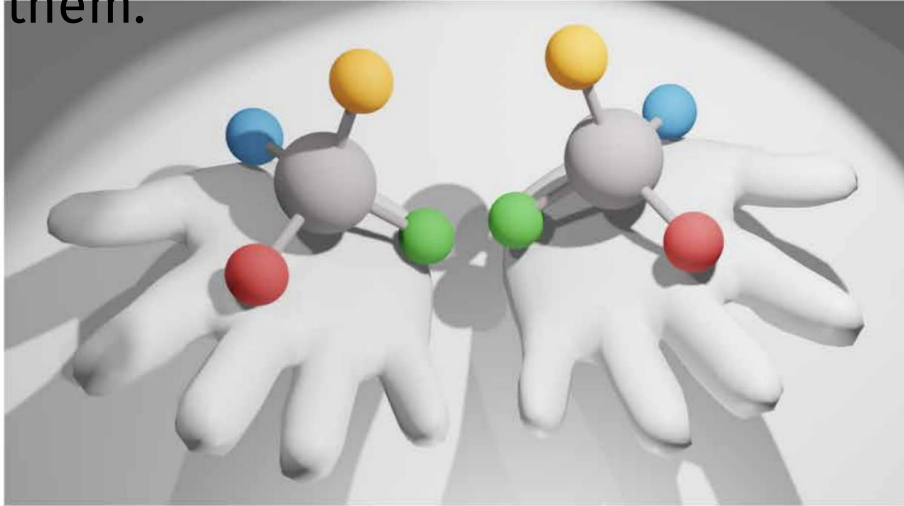
QA: 20.000x faster



$2.8 \cdot 10^{10}$ targets

Alchemical enantiomers

Quasi-degeneracy for systems if this symmetry applies to them.



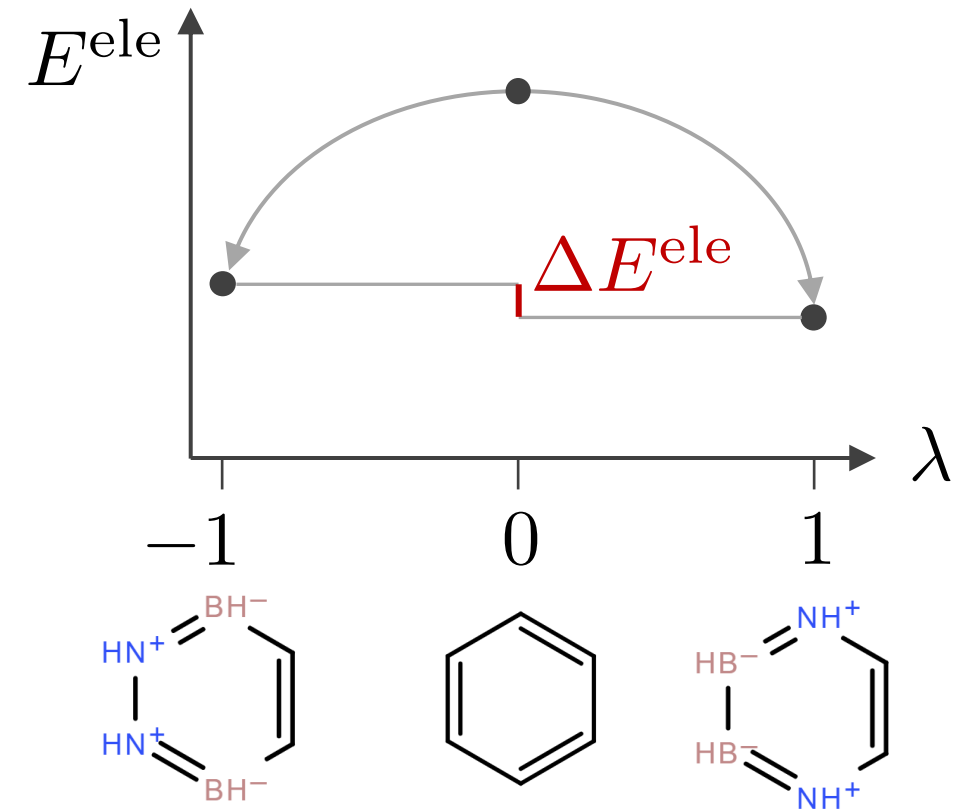
$$\Delta E_{ij}^{\text{ele}} = E_0 - E_0 + \int_{\Omega} d\mathbf{r} \sum_{n=0}^{\infty} \frac{\Delta v_i}{(n+1)!} \left[\frac{\partial^n \rho}{\partial \lambda_i^n} + \frac{\partial^n \rho}{\partial \lambda_j^n} \right]$$

$$\Delta E_{(0)}^{\text{ele}} = E_0 - E_0 = 0$$

$$\Delta E_{(1)}^{\text{ele}} = 2 \int_{\Omega} \Delta v \rho = \int_{\Omega} e \cdot o = 0$$

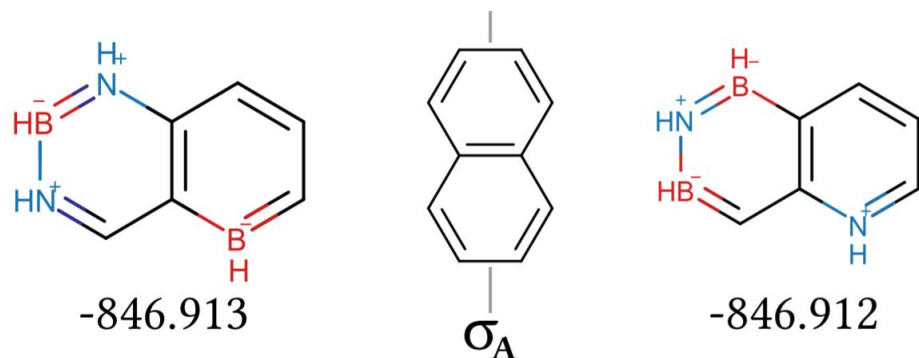
$$\Delta E_{(2)}^{\text{ele}} = \int_{\Omega} \Delta v \left[\frac{\partial \rho}{\partial \lambda_i} + \frac{\partial \rho}{\partial \lambda_j} \right]$$

$$= \int_{\Omega} \Delta v \left[\sum_I \frac{\partial \rho}{\partial Z_I} \frac{\partial Z_I}{\partial \lambda_i} + \frac{\partial \rho}{\partial Z_I} \frac{\partial Z_I}{\partial \lambda_j} \right] = 0$$



Fundamentally new symmetry

Electronic energy only



Bond energy rules

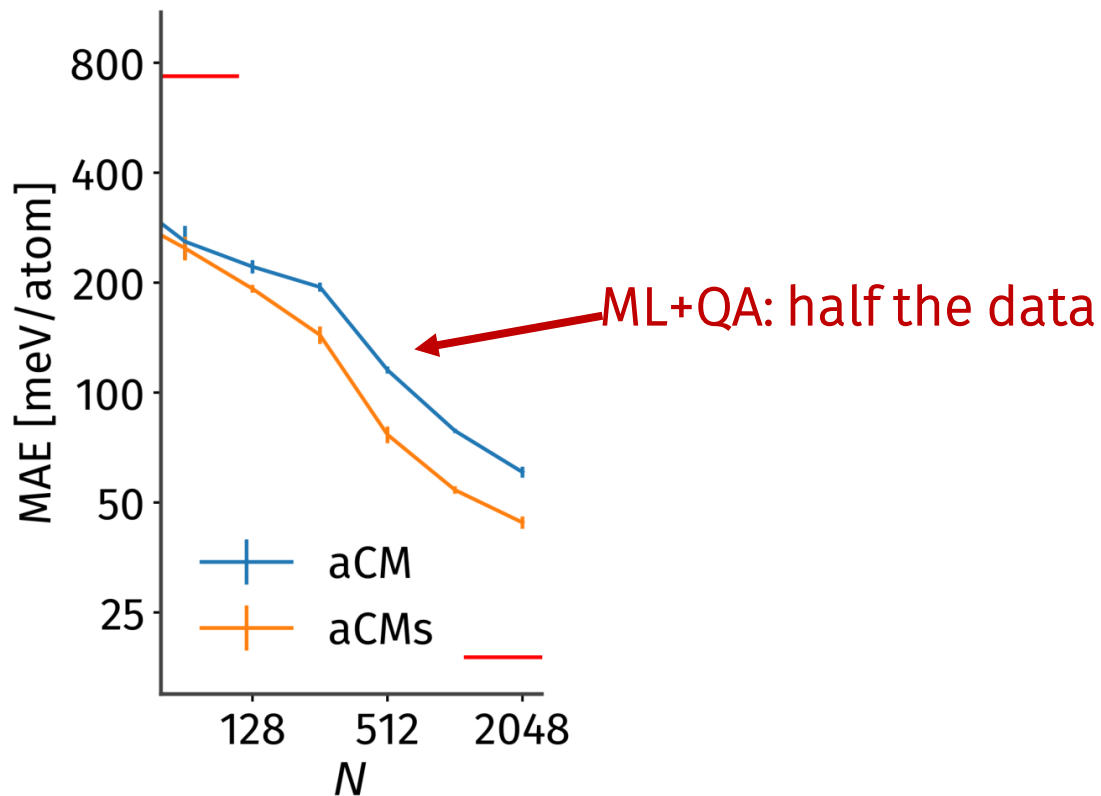
Consecutive Elements

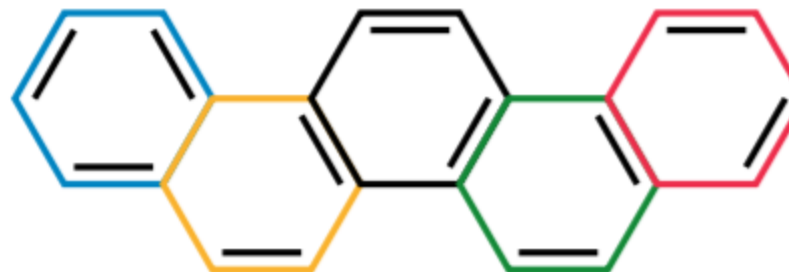
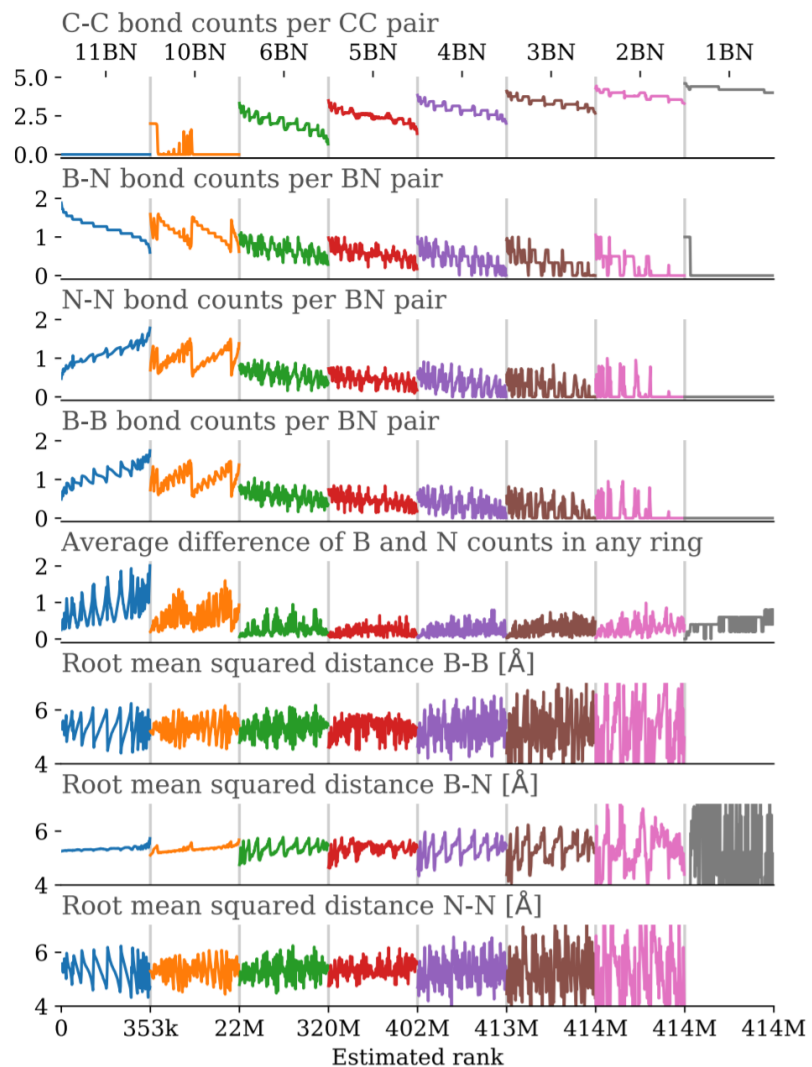
Q R S

B C N

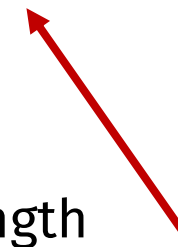
$$E_{QR} \simeq E_{SR} + 0.5(E_{QQ} - E_{SS})$$

Speed up machine learning





x 414 M



Design rules in order of decreasing strength

- Add BN pairs
- Maximize CC bonds
- Substitute sites shared between rings
- Maximize BN bonds
- Avoid N substitutions on rings sharing a larger amount of bonds with other rings
- Balance BN substitutions in each ring

QA: Millions at once!

Not a single QM calculation required!

Machine Learning | Global search, information transfer

Quantum Alchemy | Local gradients, structure in chemical space, replaces brute

force
Reaction Barriers | S.N. Heinen, GFvR, O. A. von Lilienfeld, *J. Chem. Phys.*, 2021.

Geometry Learning | D. Lemm, GFvR, O. A. von Lilienfeld, *Nat. Commun.* 2021.

Quantum Alchemy | GFvR, *J. Chem. Phys.* 2021.

Symmetry | GFvR, O. A. von Lilienfeld, *Science Adv.* 2021.



Prof. von Lilienfeld



Dominik Lemm



Marco Bragato



Stefan Heinen



Dr. Max Schwilk