

Symmetries and the Intrinsic Dimensionality of Chemical Space

Guido Falk von Rudorff, University of Kassel

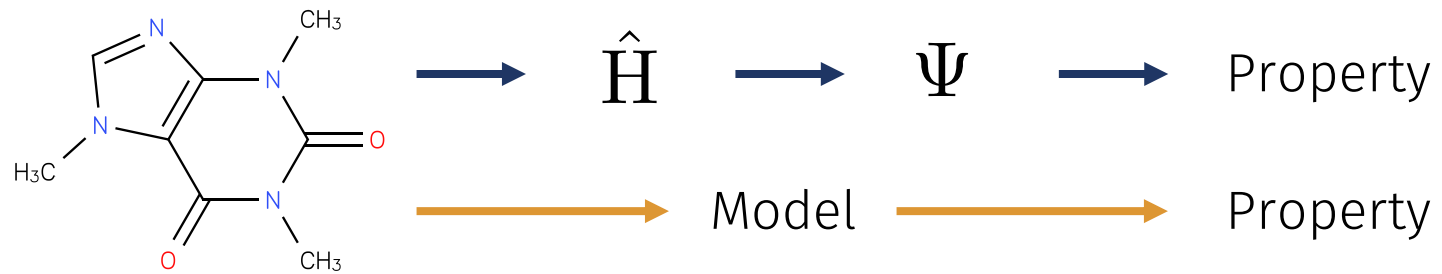
 vonrudorff@uni-kassel.de

 nablachem.org/talks

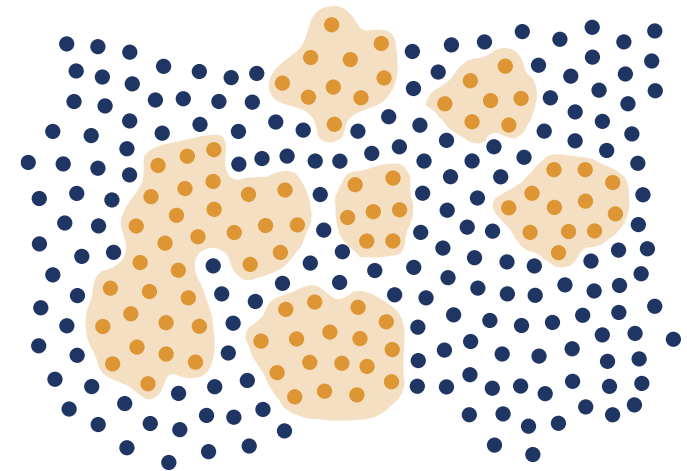
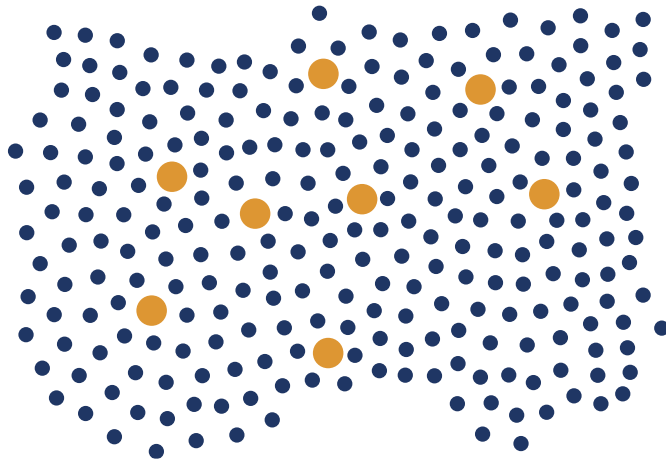
 [ferchault](https://github.com/ferchault)

 [@ferchault](https://twitter.com/ferchault)

Properties of a single system



Speed does not matter:
even enumeration is impossible.

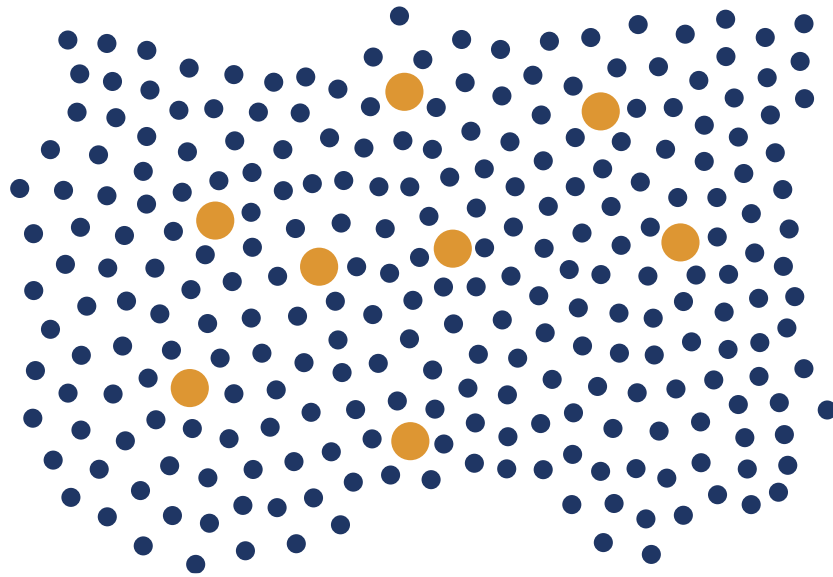


Chemical Space

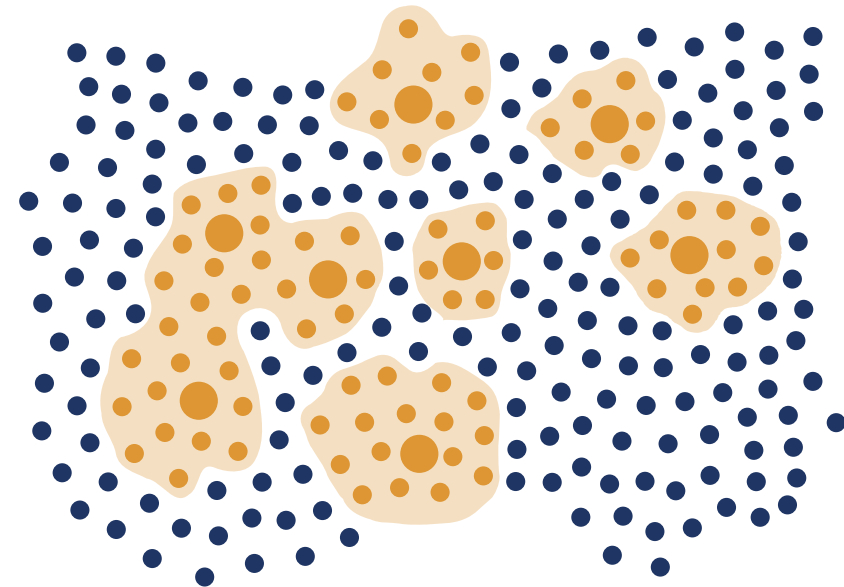
All thermally accessible configurations
for some constraints on Z_i .

$$\hat{H} = \hat{H}(\underbrace{Z_i}_{4N}, \underbrace{\mathbf{R}_i}_{1D, \text{ close to } \sum_i Z_i}, \underbrace{N_e}_{1D}, \sigma)$$

Without Perturbation



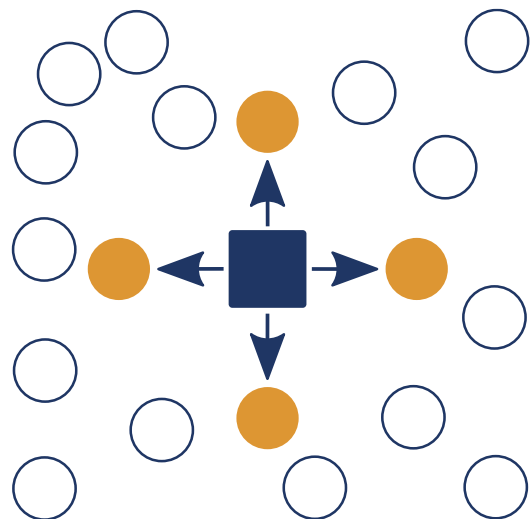
With Perturbation



Systems/Molecules

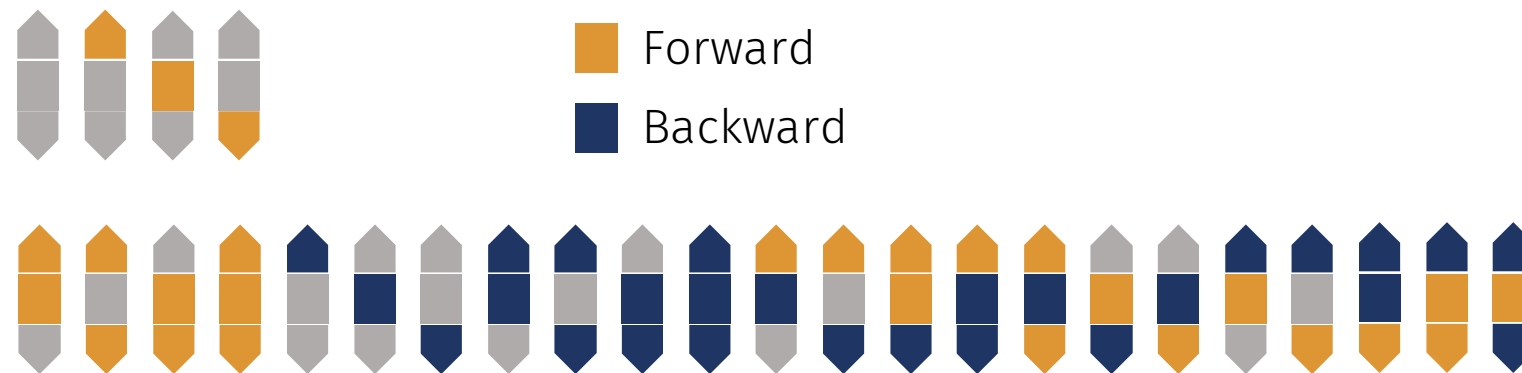
- Any
- Known
- Approximated

Quantum Alchemy



Taylor/Padé approximant

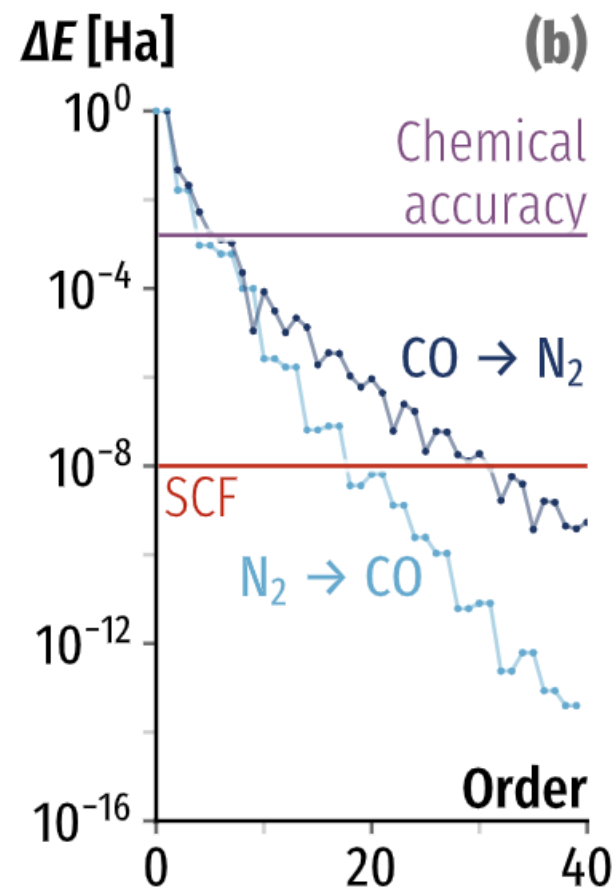
- Energy function of
 - Geometry Forces, Vibrations
 - Nuclear charges Alchemical changes
- Idea: obtain dominant leading derivatives, predict many systems



$$Q(\mathbf{x}) \simeq \sum_{|\alpha| \leq k} \frac{\partial^{|\alpha|} Q(\mathbf{a})}{\partial x_1^{\alpha_1} \cdots \partial x_n^{\alpha_n}} \frac{(\mathbf{x} - \mathbf{a})^\alpha}{\alpha!}$$

Differentiable / Analytic / Converge Quickly

- ✓ Total Energy [1,2]
- ✓ Electron density [1,2]
- ✓ Orbital eigenvalues [2]
- ✓ Dipole moments [2]
- ✓ Non-covalent interactions [1]
- ✓ Binding energies [1,2]
- ✓ Deprotonation energies [3]
- ✓ Ionisation Energy [4]
- ✓ Electron Affinity [4]
- ✓ Photoelectron circular dichroism [5]



1 | GFvR, OA von Lilienfeld, *Phys. Rev. Res.*, 2020. 2 | GFvR, *J. Chem. Phys.*, 2021. 3 | GFvR, OA von Lilienfeld, *Phys. Chem. Chem. Phys.*, 2020. 4 | E Eikey, A Maldonado, C Griego, GFvR, J Keith, *J. Chem. Phys.*, 2022. 5 | GFvR, A Artemyev, B Lagutin, P Demekhin, *J. Chem. Phys.*, accepted.

Angular emission (Photoelectron circular dichroism)

- Expensive to calculate
- Highly coupled degrees of freedom: multidimensional expansion

$$\frac{d\sigma^\pm}{d\Omega} = \frac{\sigma}{4\pi} \left[1 \pm \beta_1 P_1(\cos \theta) - \frac{1}{2} \beta_2 P_2(\cos \theta) \right]$$

dichroic
parameter

anisotropy
parameter

$$\frac{\partial \beta_i}{\partial R_x}$$

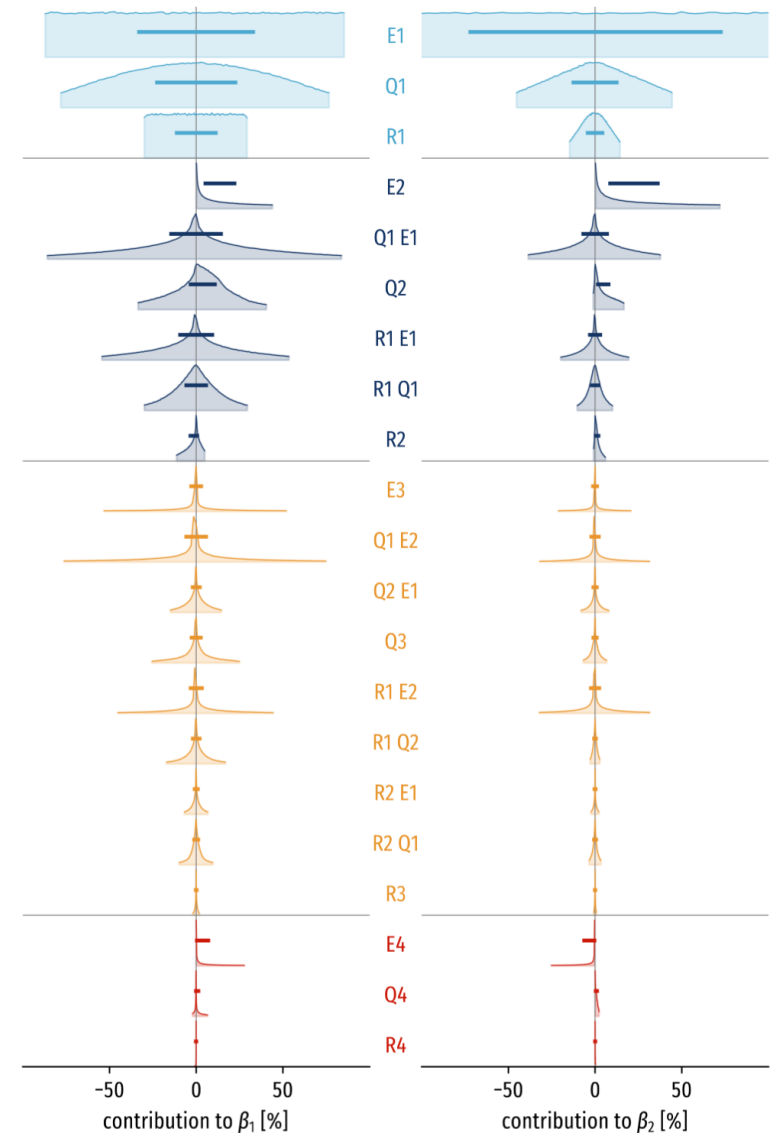
∈ R1

$$\frac{\partial^2 \beta_i}{\partial Q_y^2}$$

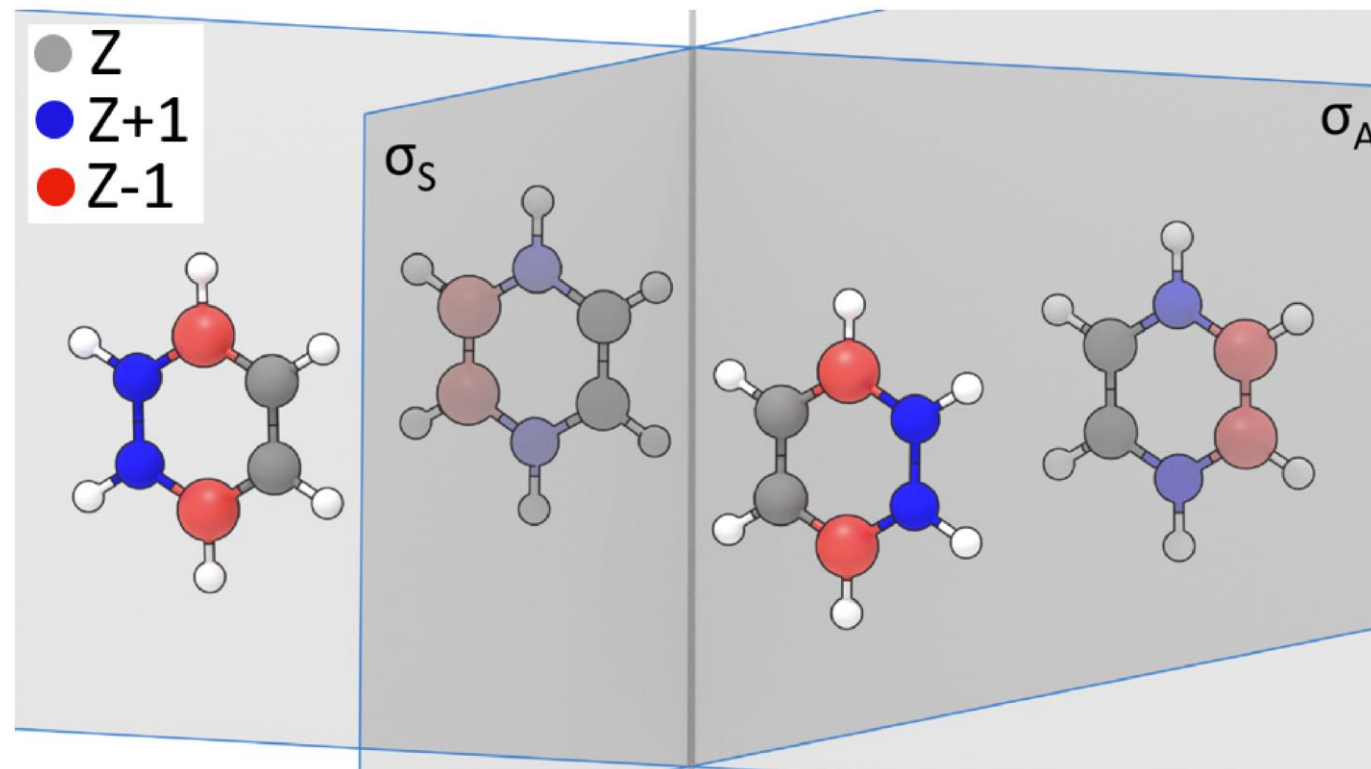
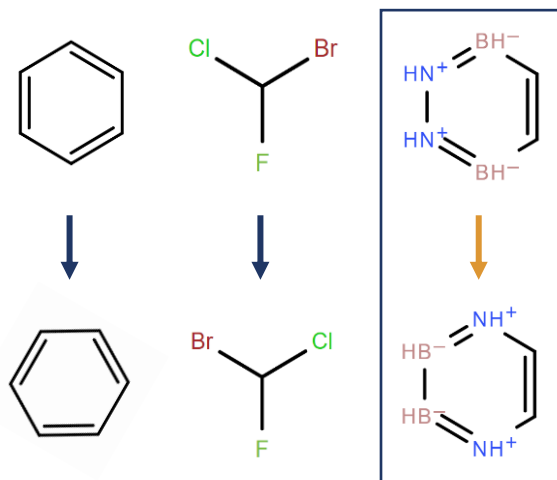
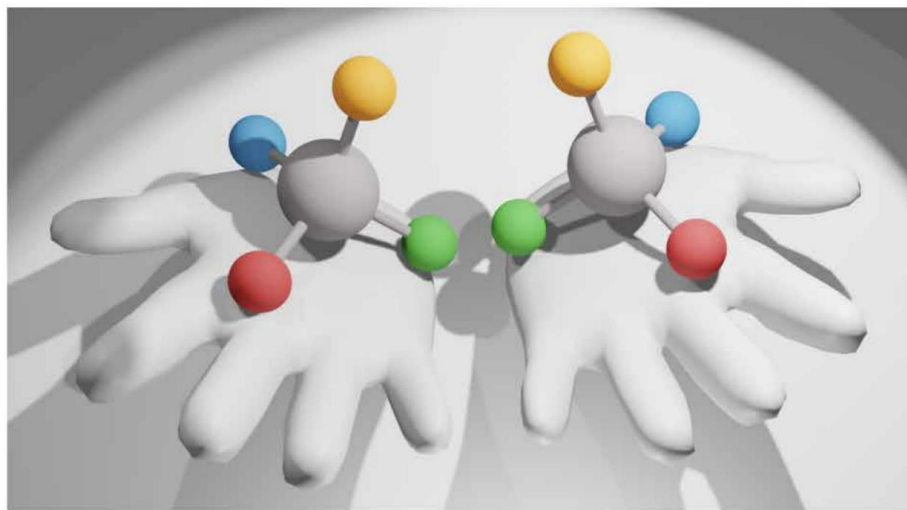
∈ Q2

$$\frac{\partial^2 \beta_i}{\partial R_x \partial Q_y}$$

∈ R1 Q1

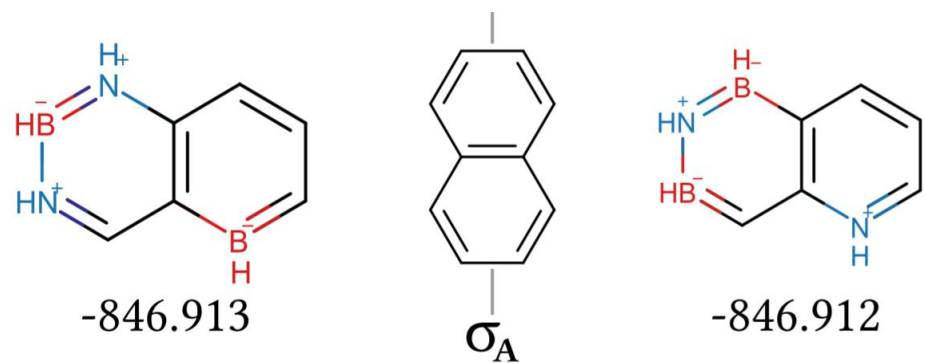


Alchemical Enantiomers

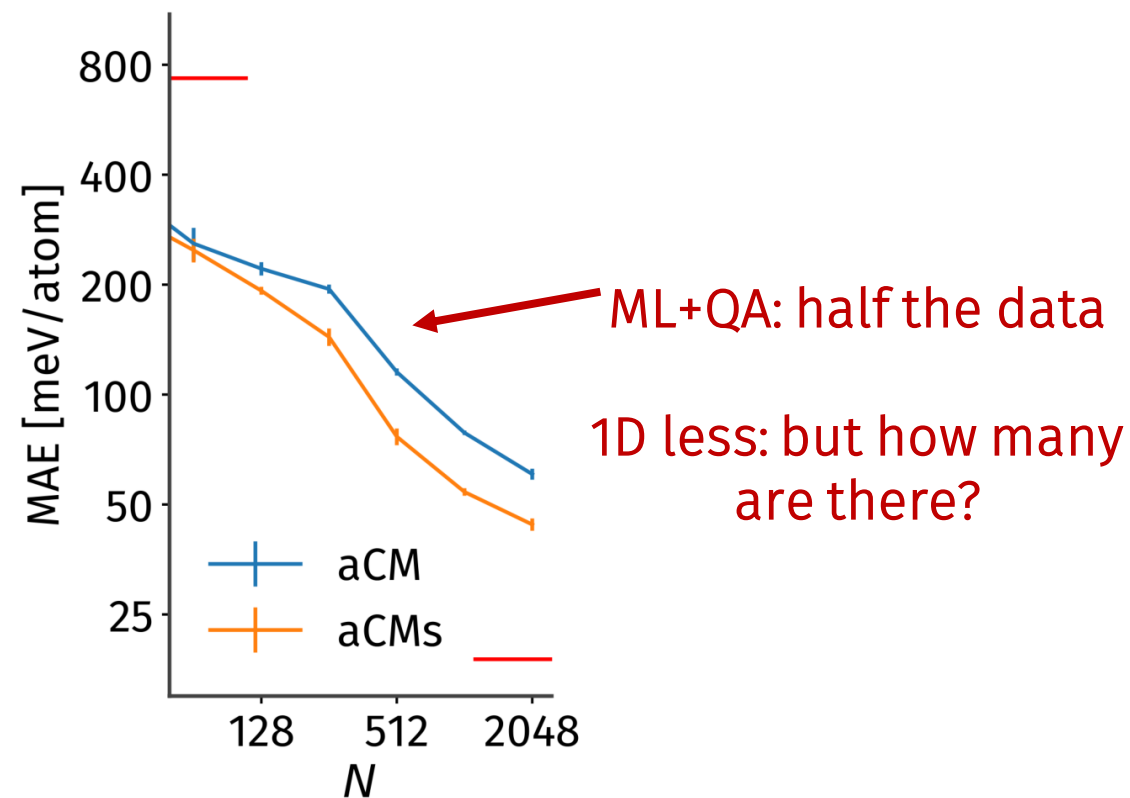


Fundamentally new symmetry

Electronic energy only



Speed up machine learning



Definition Intrinsic Dimension

Minimal number of degrees of freedom to describe a property.

≠ intrinsic dimension of a point cloud!

$$f(x, y) = x + y$$

Example: Dimers

- Energy: 3 dimensions
- Net charge: 1 dimension

Example: Atom

- 1 dimension
- 1 dimension

$$E(\mathbf{R}_I, \mathbf{Z}_I)$$

Distance Net charge Z_1+Z_2
Asymmetry Z_1-Z_2



Ali Banjafar

How to measure?

Dynamics in $4N$ embedding using quantum alchemy, keeping the property Q constant

$$U(\Delta\mathbf{R}, \Delta\mathbf{Z}) = k_1 [Q(\Delta\mathbf{R}, \Delta\mathbf{Z}) - Q_{\text{ref}}]^2 + k_2 \|\Delta\mathbf{R} \oplus \Delta\mathbf{Z}\|$$

Trajectories form a level set / hyperplane of the null space.

(Local) intrinsic dimension D_Q

Apply point cloud estimators to level set: ID of null-space $D_{\text{NS},Q}$

$$D_Q = 4N - D_{\text{NS},Q}$$

Limitations

- Sampling (long enough and bias-free)
- Point cloud ID estimators

Small molecules (up to six atoms)

- Equilibrium and non-equilibrium ID identical
- Scaling linearly, prefactor $\ll 4$

Allow for data-driven fundamental statements
“Most molecules do X”, “High X means low Y”

Transferability

More reliable understanding of trends

Lower data bias

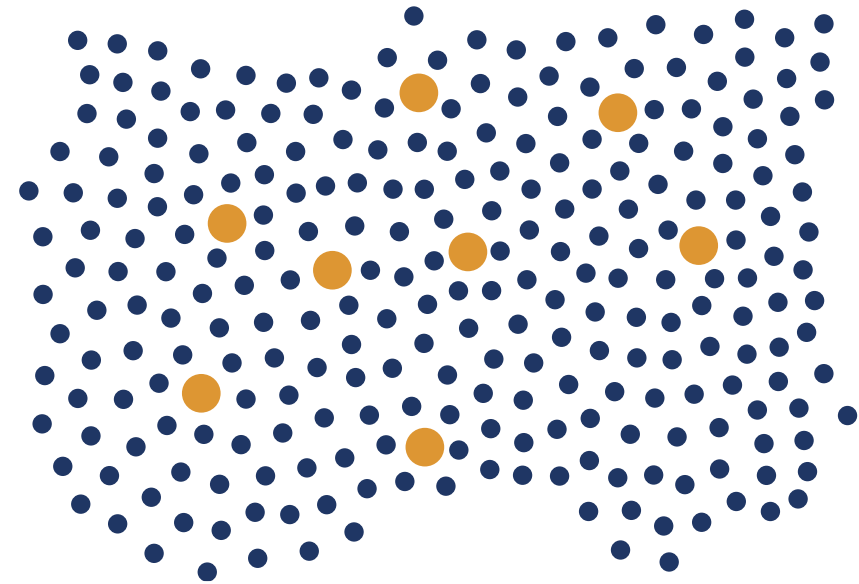
More realistic generalisation error

Formal statements

Often require uniform sampling

Problems

- Total number unknown
- Distribution unknown

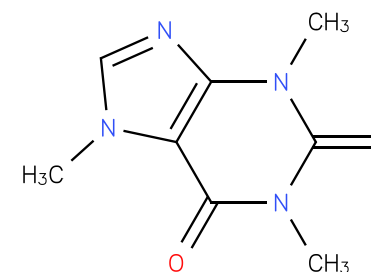
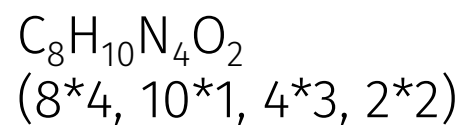


Goal

Sample all molecules (with given constraints) with known probabilities.

Sampling

- Choose **weighted** random sum formula
- Choose **weighted** random degree sequence
- Choose **weighted** random molecular graph



Requirements

- Find all sum formulas and degree sequences
- Sample loop-free multigraphs with given degree sequences uniformly
- Find **weights**

Solved

Solved, Seconds^[1]

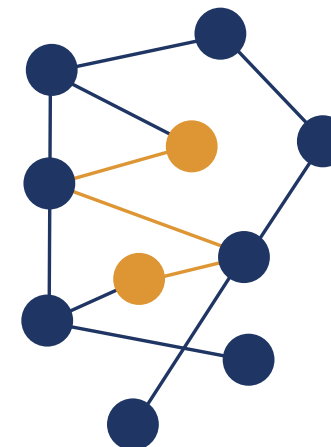
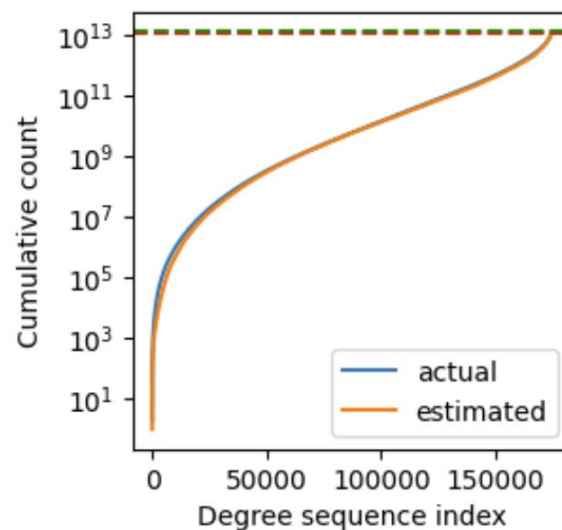
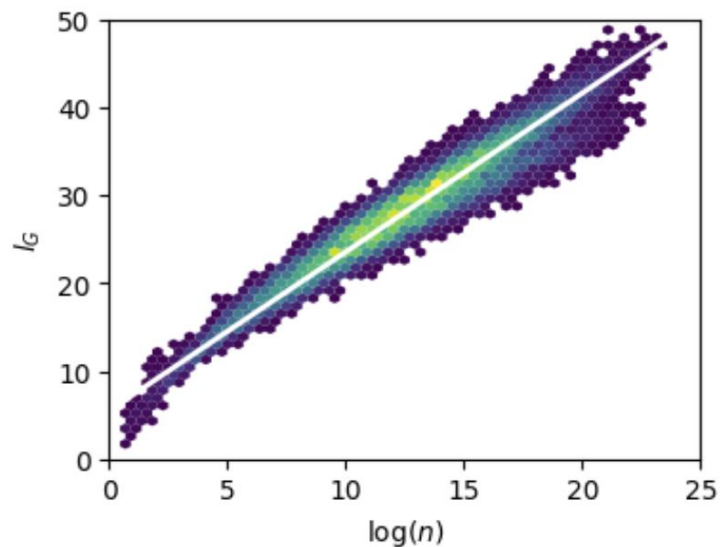
Goal




Estimate number n of loop-free multigraphs with given degree sequence.

Average Path Length l_G

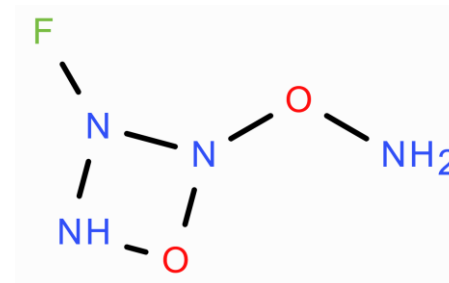
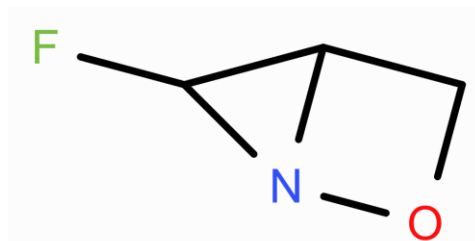
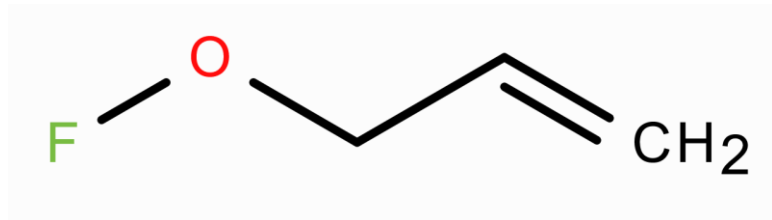
Sample from random molecule pairs

$$l_G \sim \log n$$



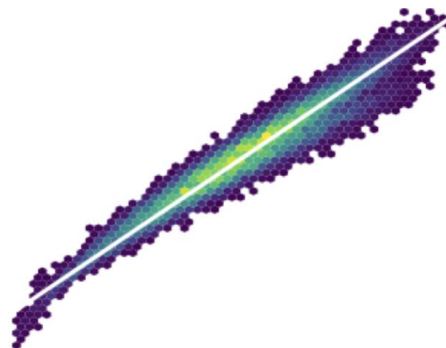
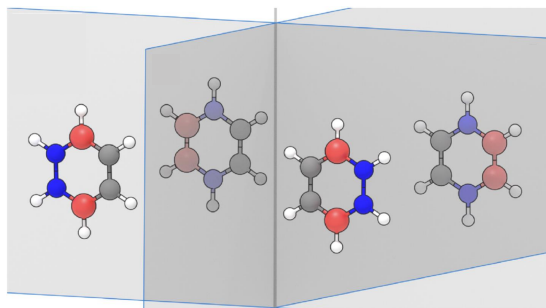
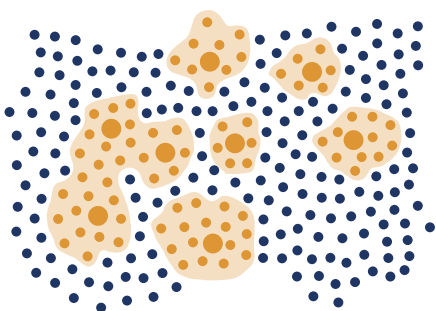
-  Molecule
-  Identical up to one bond
-  Minimum path

10 atoms, CHONF, at least 3 hydrogens and one fluorine



Available conditions in this approach

- Number of atoms
- Atoms per element (and combinations)
- Valences
- Via rejection sampling:
 - Bonds (count and bond orders)
 - Ring presence / membership
 - Substructures
 - ...



Symmetries | Reducing (“folding”) search space

Dimensionality | Limit data efficiency

Random Sampling | Allows for collective statements

Thanks

Anton Artemyev

Ali Banjafar

Emily Eikey

Philipp Demekhin

Chasz Griego

John Keith

Boris Lagutin

Anatole von Lilienfeld

Alex Maldonado

Sana Qureshi