# Design in Compound Space with Machine Learning and Quantum Alchemy

Guido Falk von Rudorff, University of Kassel
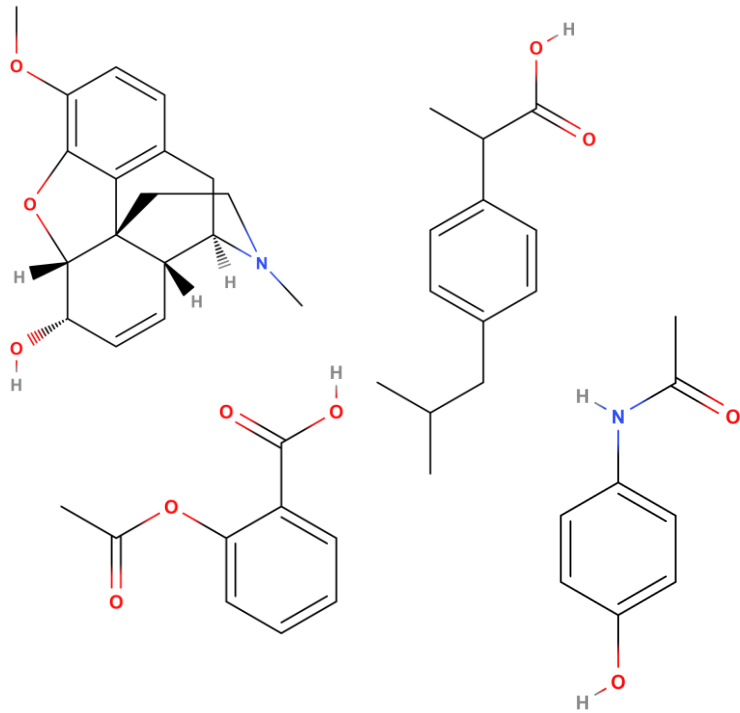
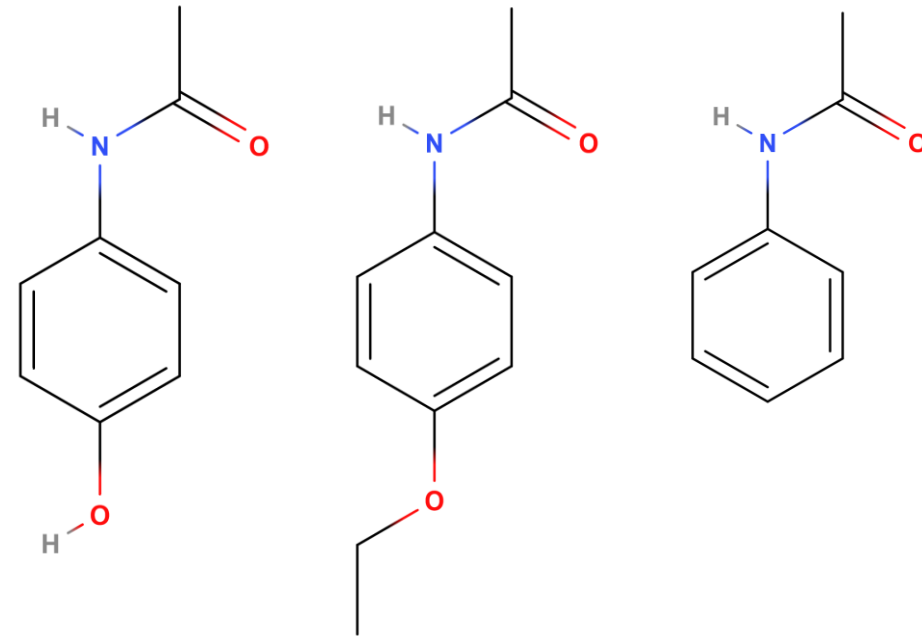ferchault          @ferchault          nablachem.org
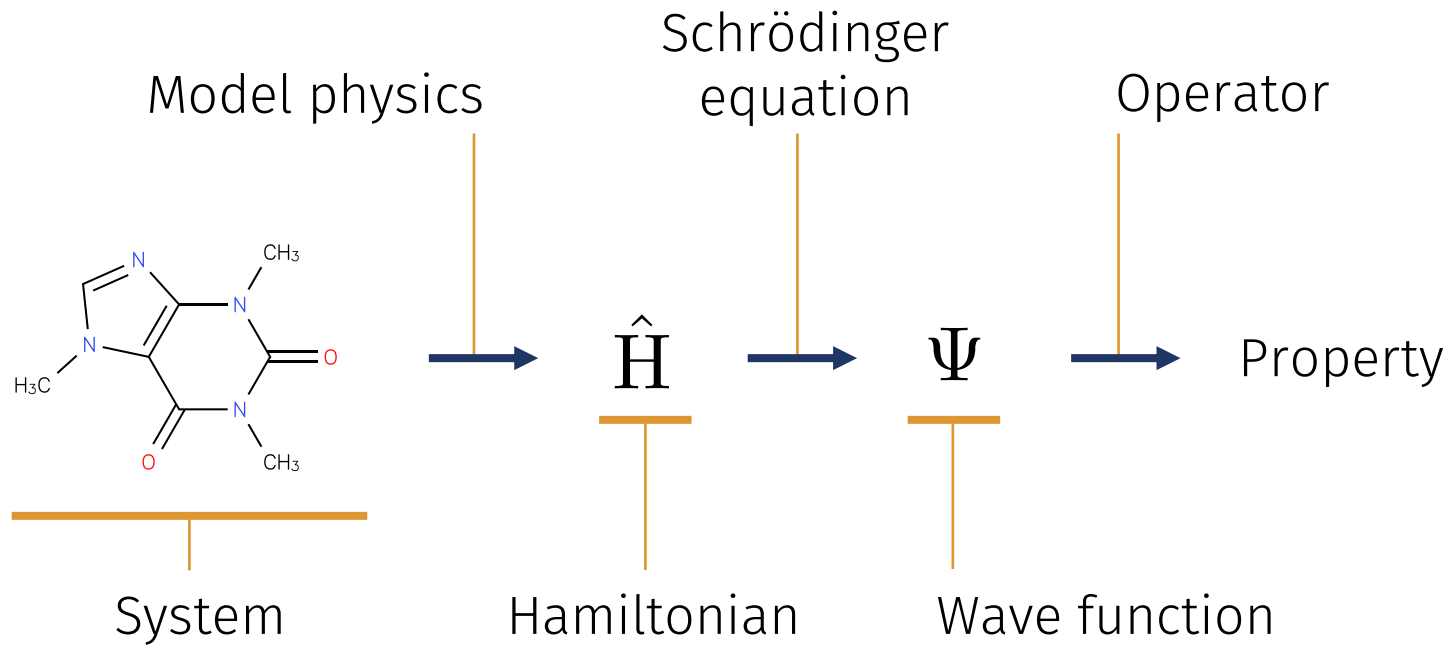
**Global Search Problem**
Which class of compounds?

**Local Search Problem**
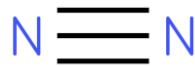Which particular species within that class?

Model physics

Schrödinger
equation

Operator

Ĥ

Ψ

Property

System

Hamiltonian

Wave function

$$\Psi = \Psi(\mathbf{r}_1, \mathbf{r}_2, \mathbf{r}_3, \ldots, \mathbf{r}_n)$$

Methane

$CH_4$

$N_2$

$N \equiv N$

Ethanol

$H_3C$     $OH$
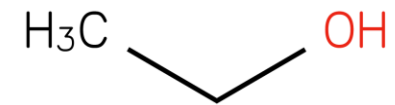
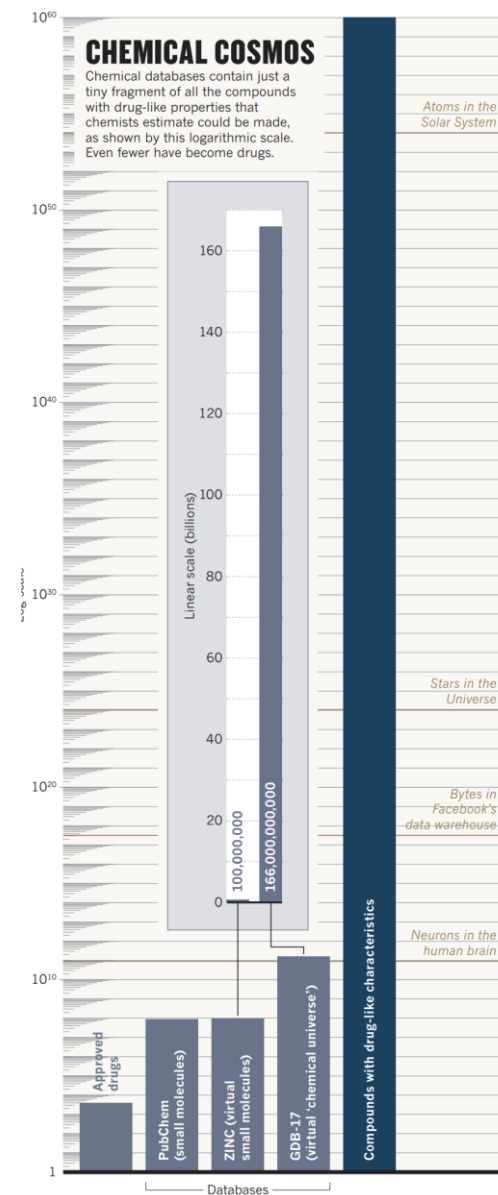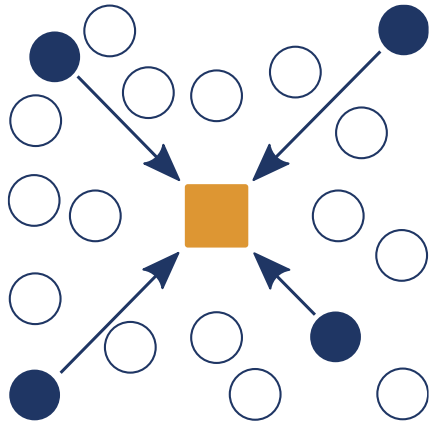Solved by approximations in computational chemistry?

## Commercial databases
- 164 million molecules
- 15k added daily

## Scale
- One person: 1 million compounds/second
- 10 billion people on earth
- $10^{26}$ universe ages to go through



A. Mullard, *Nature*, 2017.
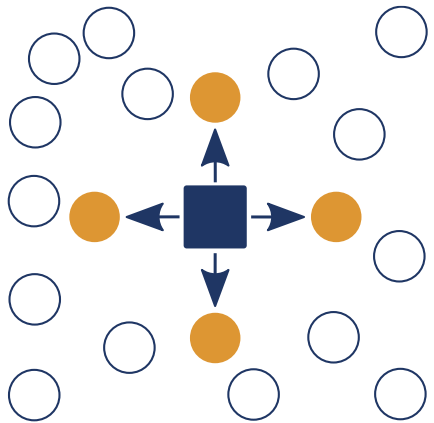
## Machine Learning
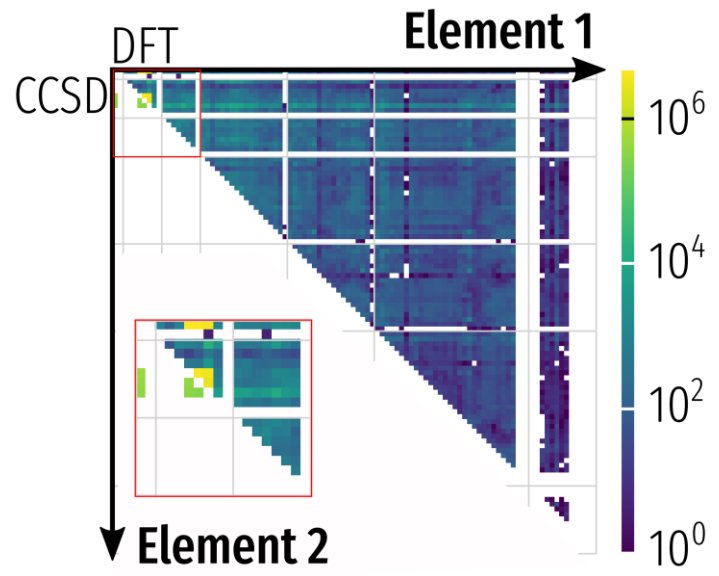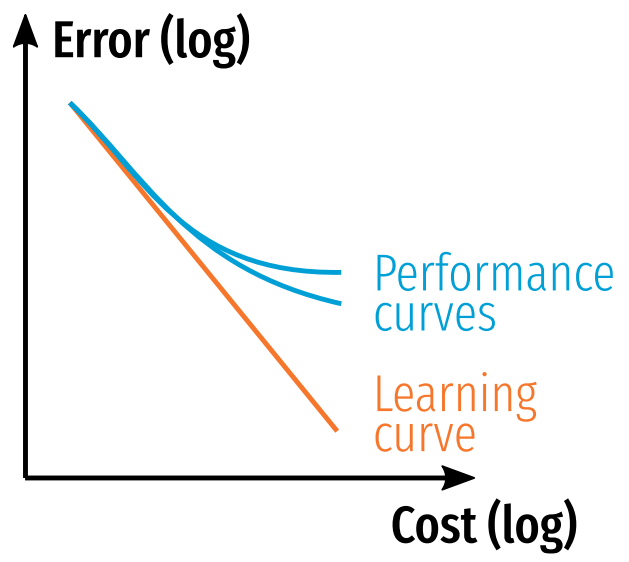


**Foundations |** Statistical modelling

**Accuracy |** Systematically improvable through data and training

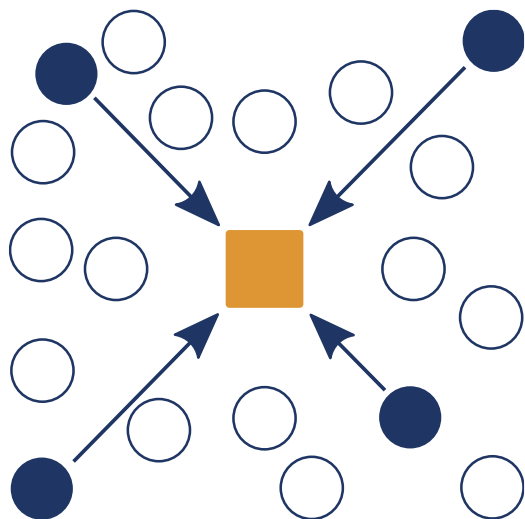**Specialty |** Universal, scale-bridging, data-driven approach

**Limitation |** Requires training data, no black box

## Quantum Alchemy

Error (log) vs. Cost (log): Performance curves, Learning curve



DFT, CCSD — Element 1, Element 2; color scale $10^0$ to $10^6$

## Machine Learning



## Kernel-Ridge-Regression

- Efficient in the low-data regime (around 1k points)
- Ingredients
  - Representation $\quad \mathbf{M}$
  - Similarity measure $\quad k(\mathbf{M}_i, \mathbf{M}_j)$
  - Observed properties $\quad \mathbf{y}$
- Training
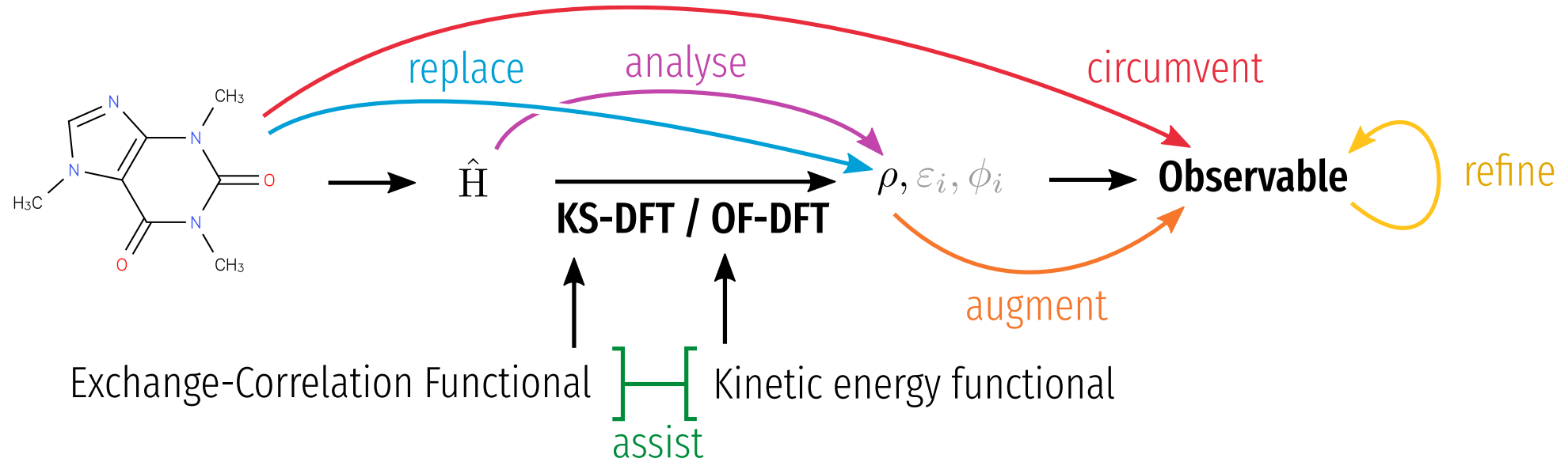  - Pairwise similarities $\quad \mathbf{K}$
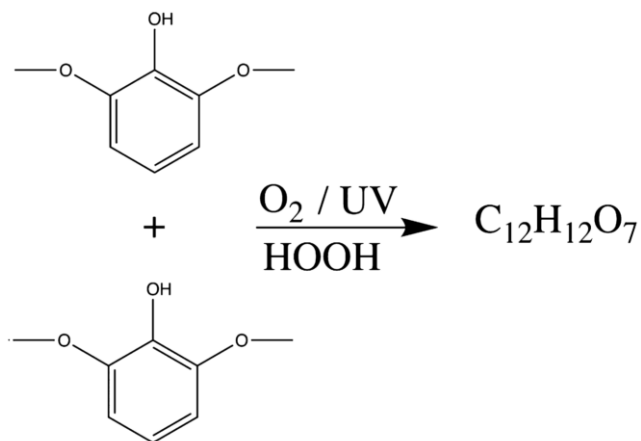  - Model coefficients $\quad \boldsymbol{\alpha} = (\mathbf{K} + \lambda\mathbf{I})^{-1}\mathbf{y}$
- Predictions
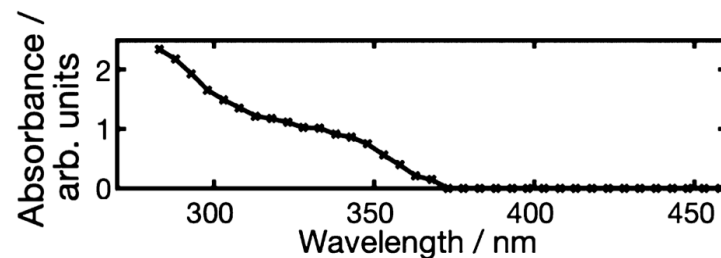  - Compare to training $\quad \tilde{q}(\mathbf{M}) = \sum_i \alpha_i k(\mathbf{M}, \mathbf{M}_i)$

V. N. Vapnik: The Nature of Statistical Learning Theory, *Springer* 2000.
B. Schölkopf, A. J. Smola: Learning with Kernels, *MIT Press*, 2001.

replace

analyse

circumvent

$\hat{H}$      $\rho, \varepsilon_i, \phi_i$      **Observable**

refine

**KS-DFT / OF-DFT**

augment

Exchange-Correlation Functional    Kinetic energy functional

assist

## Experiment





$$\underset{}{\text{OH}} + \underset{}{\text{OH}} \xrightarrow[\text{HOOH}]{O_2 / UV} C_{12}H_{12}O_7$$

## Identified features



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

## Search space
Molecular graphs:        264 M
Stable molecules:        123 M

Traditional / G2S

Lewis structure

Embedding

Optimization

3D Geometry

Energy
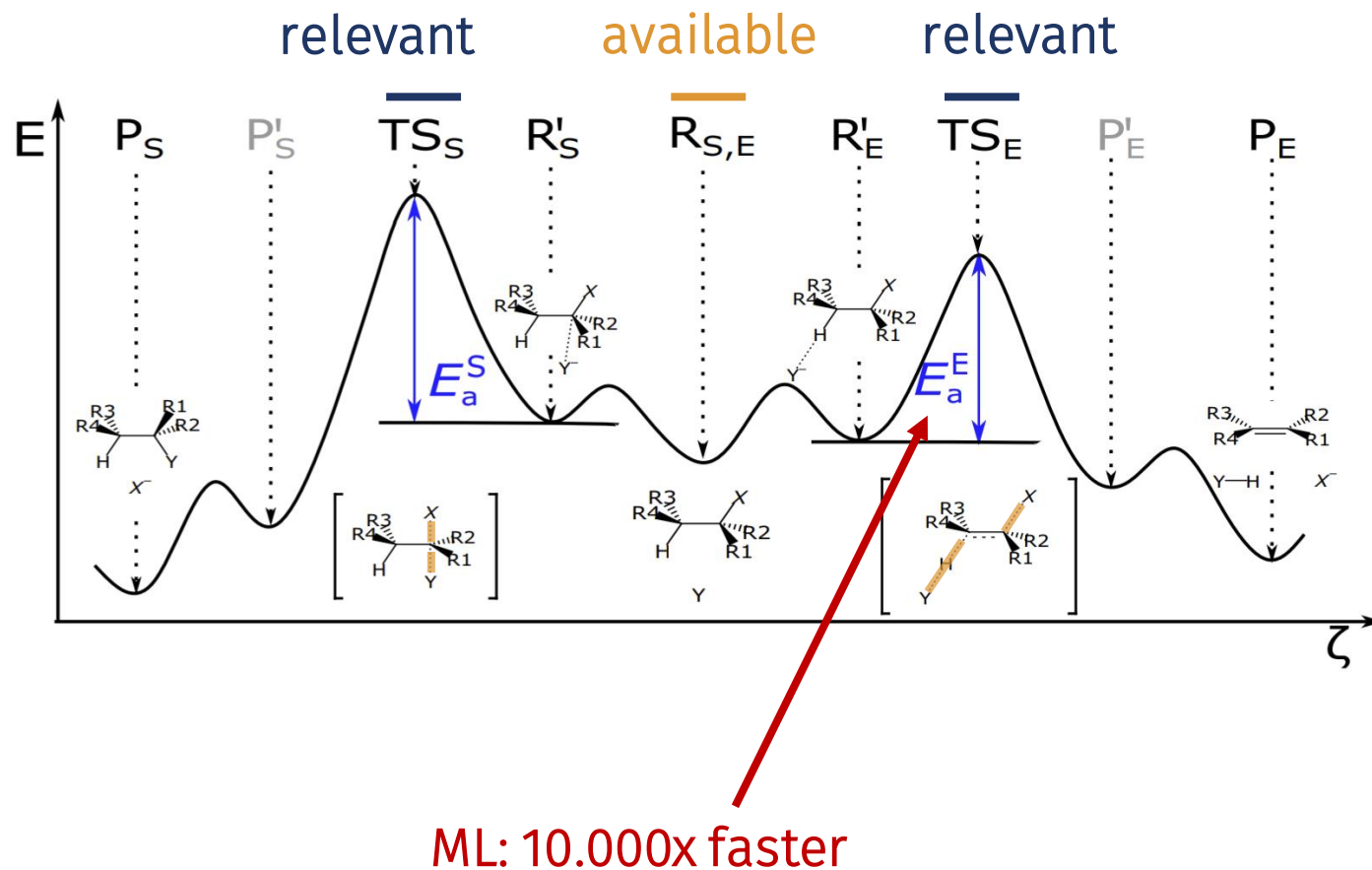
ML: 100.000x faster
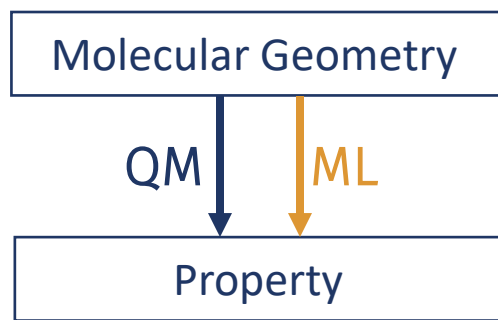
G2S
- Closer to DFT than common methods
    - Small molecules

- Applicable to complex chemical spaces
    - Transition state geometries
    - Carbenes
    - Elpasolite crystals

relevant     available     relevant

$E_a^S$     $E_a^E$
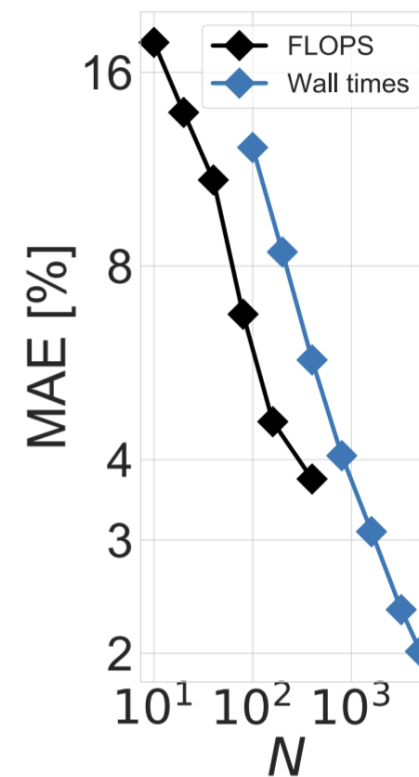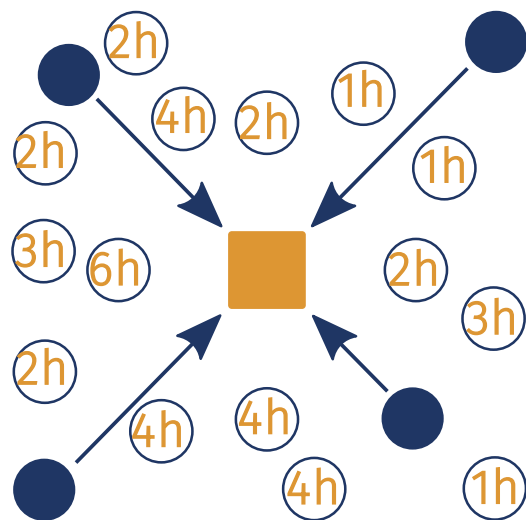
ML: 10.000x faster

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

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**.

## Molecular Geometry
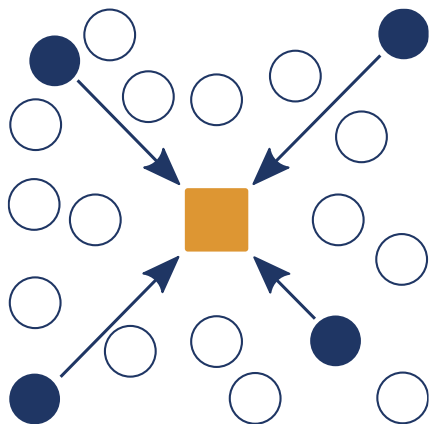
QM | ML

Property



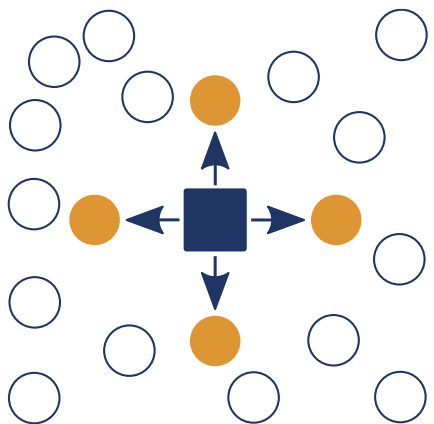## Computational effort as molecular property
- Improves models
- Accuracy depends on problem
    - Single points:              2%
    - Transition state search:   25%
    - Geometry optimisations: 40%



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

## Machine Learning



## Quantum Alchemy



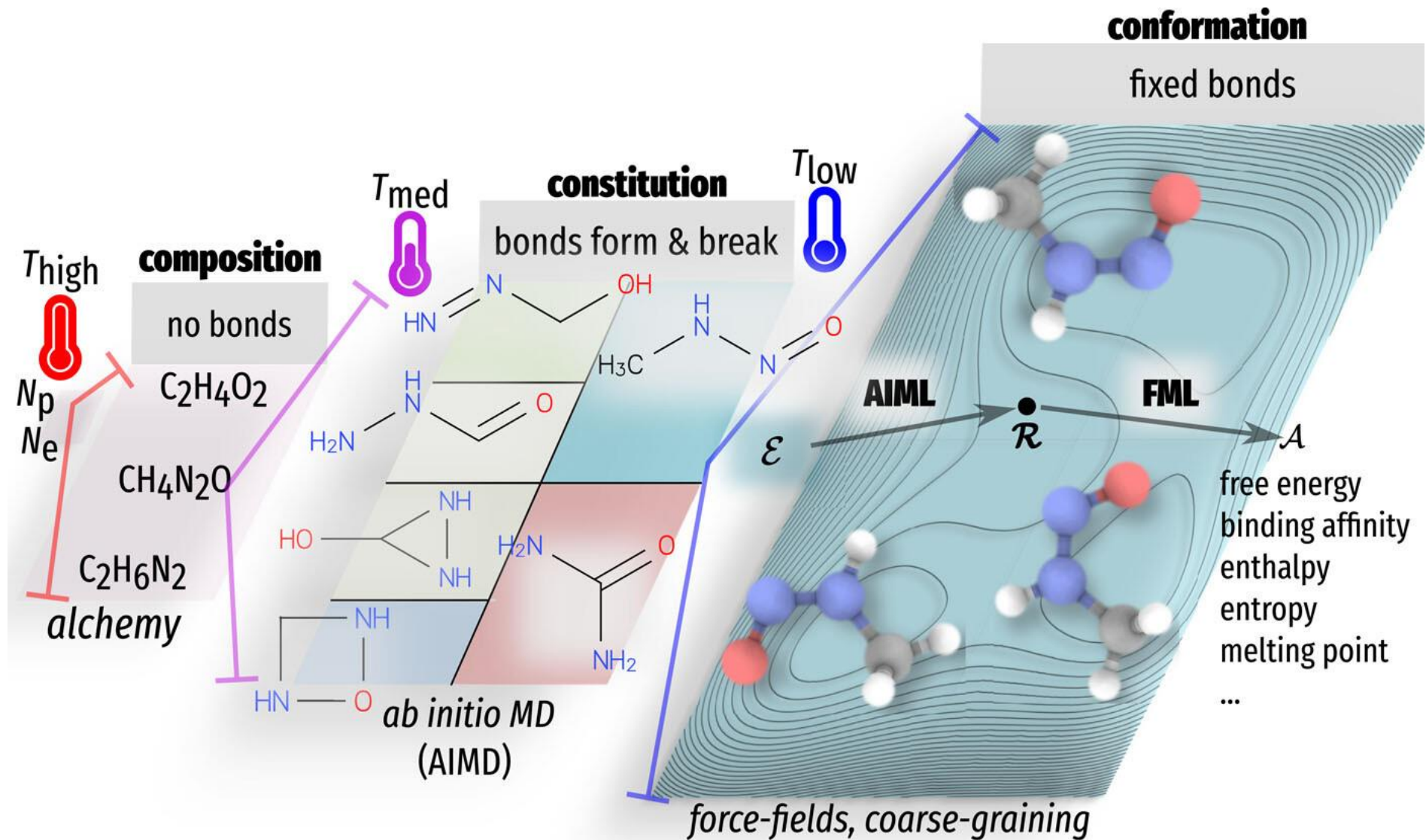Foundations | Perturbation theory

Accuracy | Systematically improvable through higher orders terms

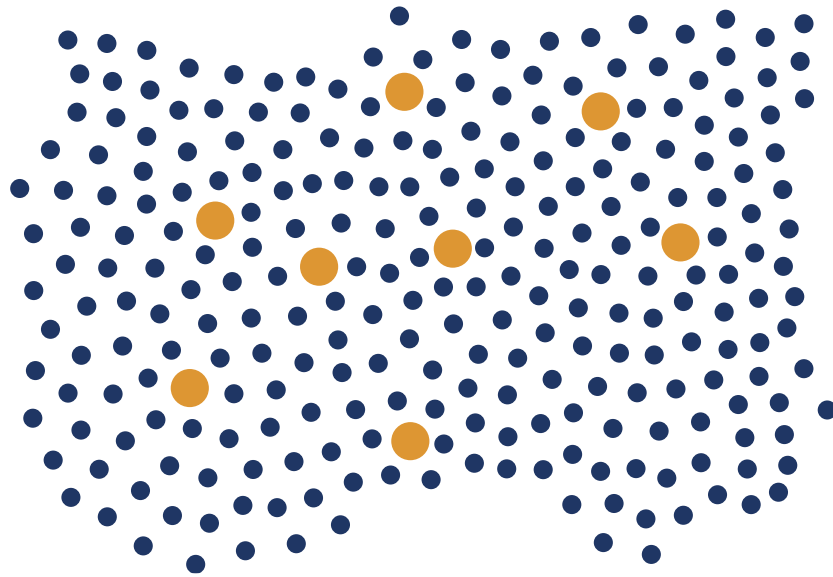Specialty | Combinatorial scaling with chemical diversity
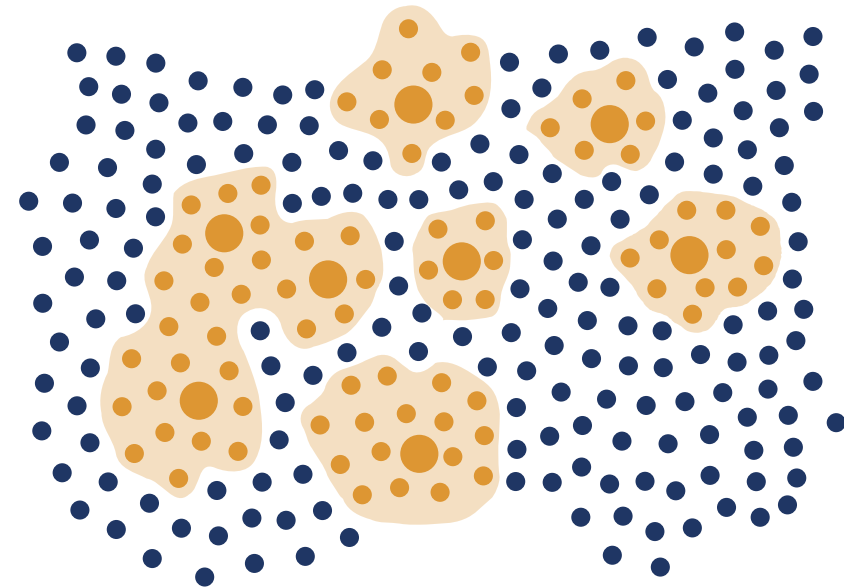
Limitation | Finite range in chemical space

J. Weinreich, D. Lemm, **GFvR**, O.A. von Lilienfeld, *J. Chem. Phys.*, **2022**.

Without Perturbation

With Perturbation

Systems/Molecules
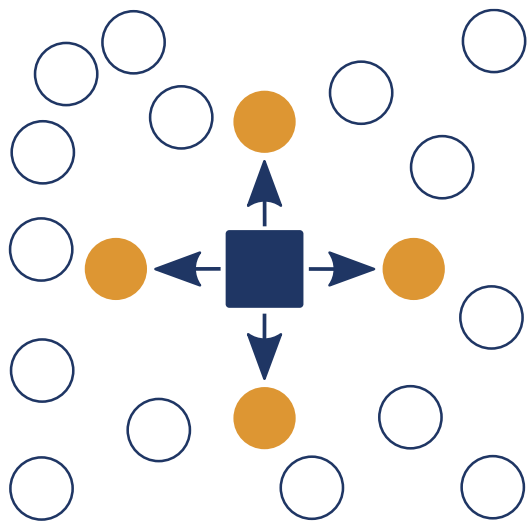- Any
- Known
- Approximated

**Perspective shift**
Few highly accurate calculations
instead of many intermediate ones

$$\hat{\mathrm{H}} = \hat{\mathrm{H}}(Z_i, \mathbf{R}_i, N_e, \sigma)$$
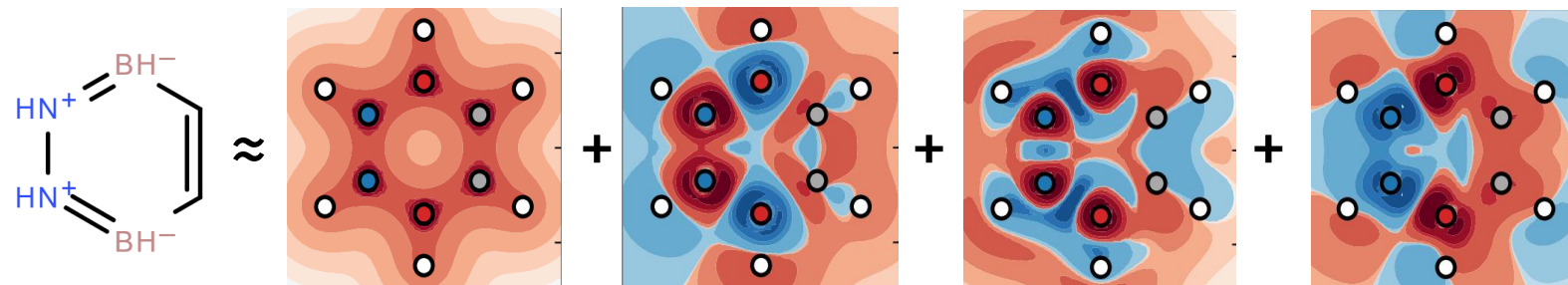
1D

4N    1D, close to $\sum_i Z_i$

## Quantum Alchemy



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



E. B. Wilson, *J. Chem. Phys.* 1962.
**GFvR**, O. A. von Lilienfeld, *Phys. Rev. Res.*, 2020.

Interpolate between molecular isoelectronic Hamiltonians

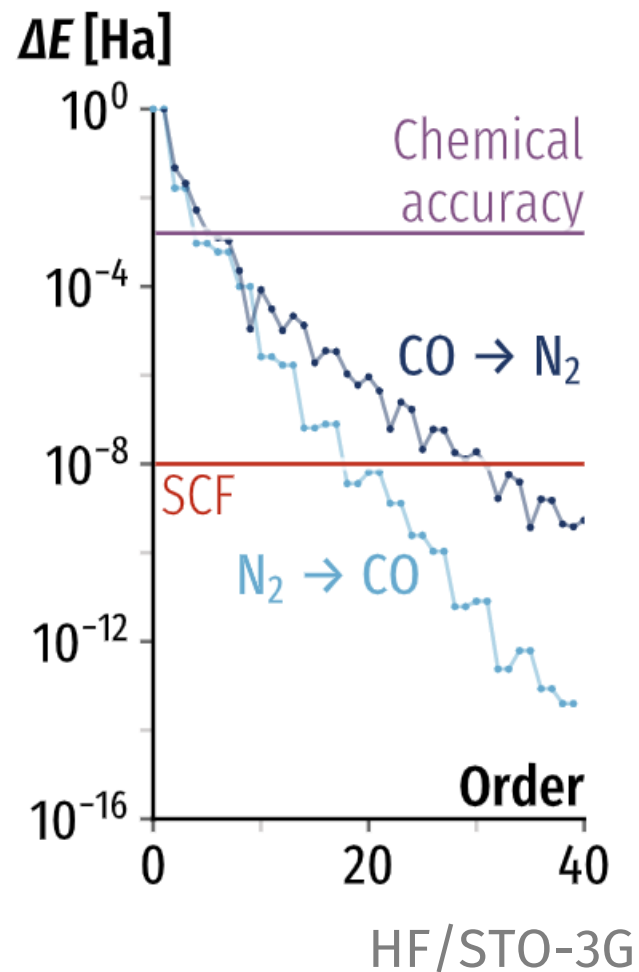$$\hat{H}(\lambda) \equiv \lambda \hat{H}_{\mathrm{t}} + (1 - \lambda)\hat{H}_{\mathrm{r}}$$

Taylor expansion around reference molecule

$$E_{\mathrm{t}} = E_{\mathrm{r}} + \Delta E^{\mathrm{NN}} + \int_{\Omega} d\mathbf{r} \sum_{n=0}^{\infty} \frac{1}{(n+1)!} \Delta v \frac{\partial^n \rho_{\lambda}(\mathbf{r})}{\partial \lambda^n}\bigg|_{\lambda=0}$$

- Gives consistent energies, densities, forces, …
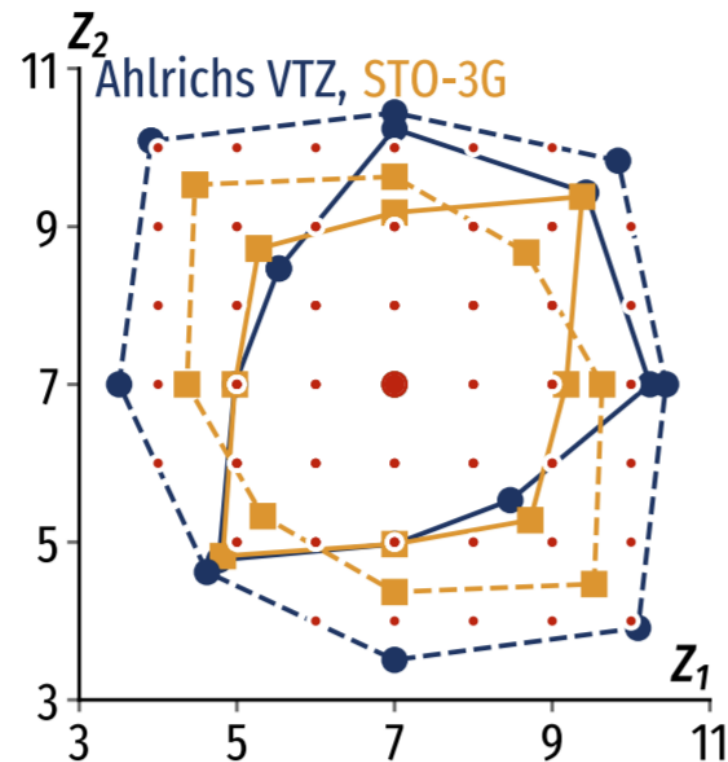- Uses the same derivatives for all predictions

AT THIS POINT, YOU'RE PROBABLY THINKING, "I LOVE THIS EQUATION AND WISH IT WOULD NEVER END!"
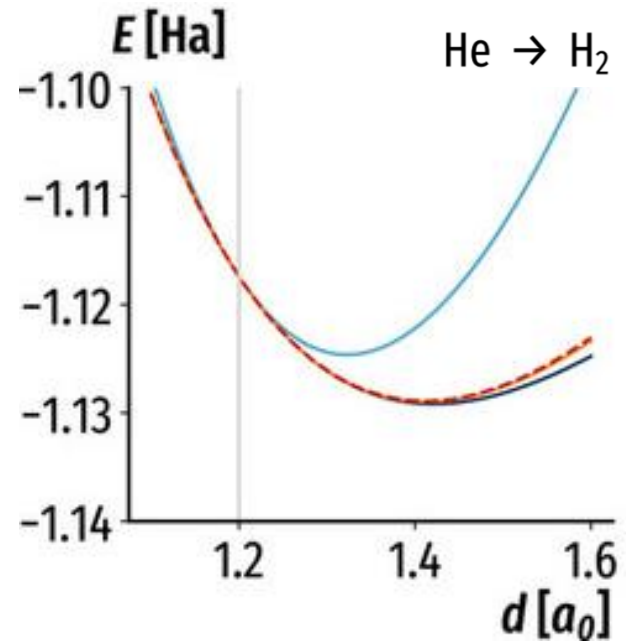
WELL, GOOD NEWS!

HF/STO-3G

## Taylor expansion
- First terms accurate enough
  - Truncate early
- Converges to the right value
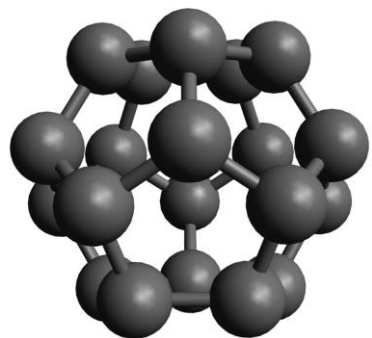- Large convergence radius
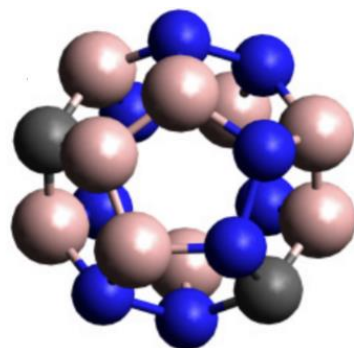- Scales with chemical space

**Taylor expansion**
- Large changes still converge (more slowly)
- Geometric response can be recovered
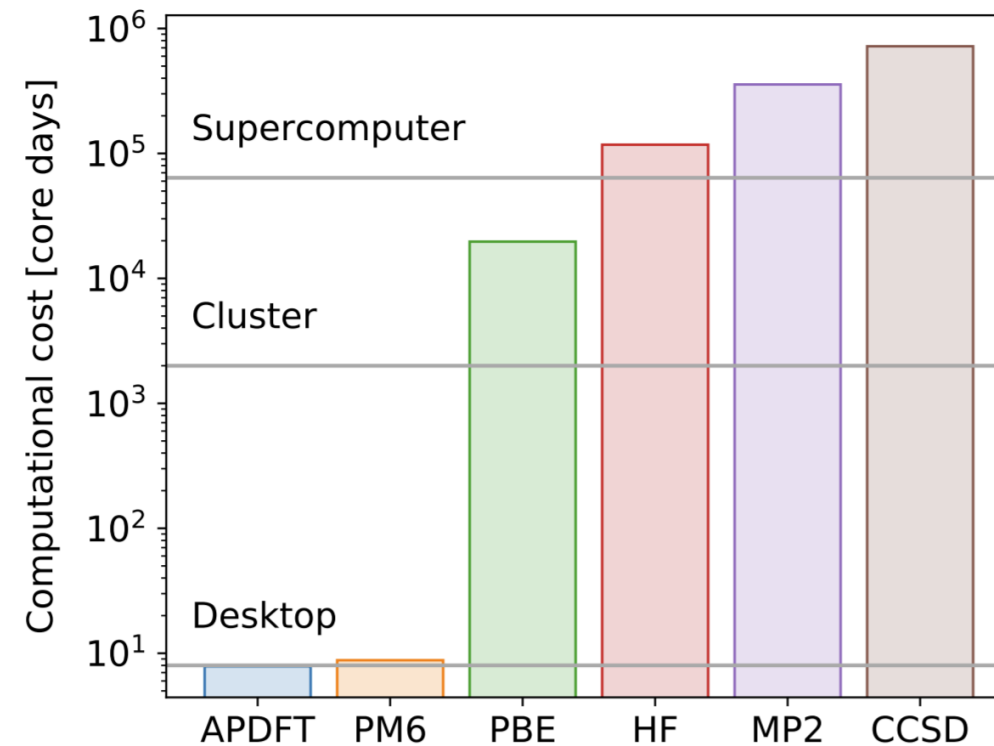
## Scaling with chemical space

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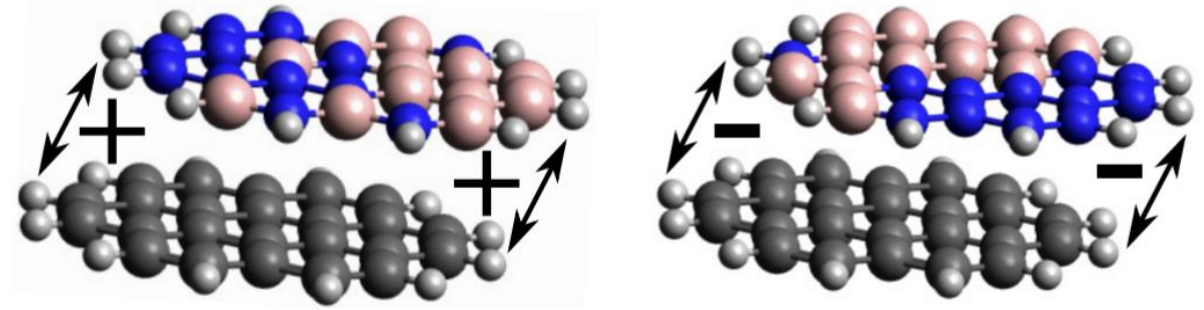- 5 derivatives for third order



$C_{20}$

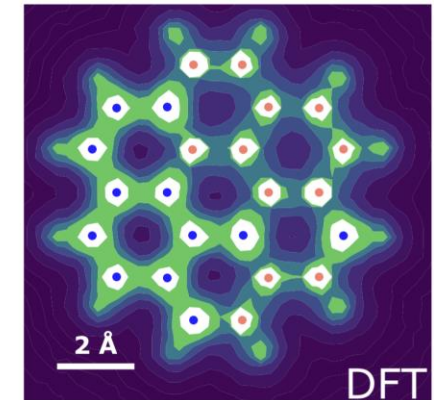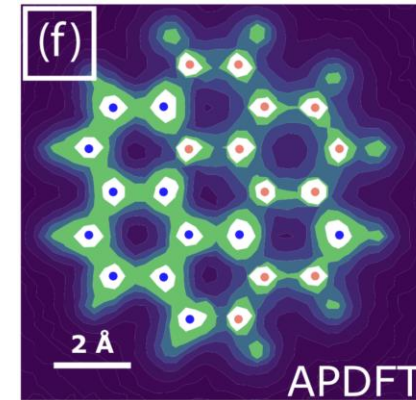$3.1 \cdot 10^6$ targets

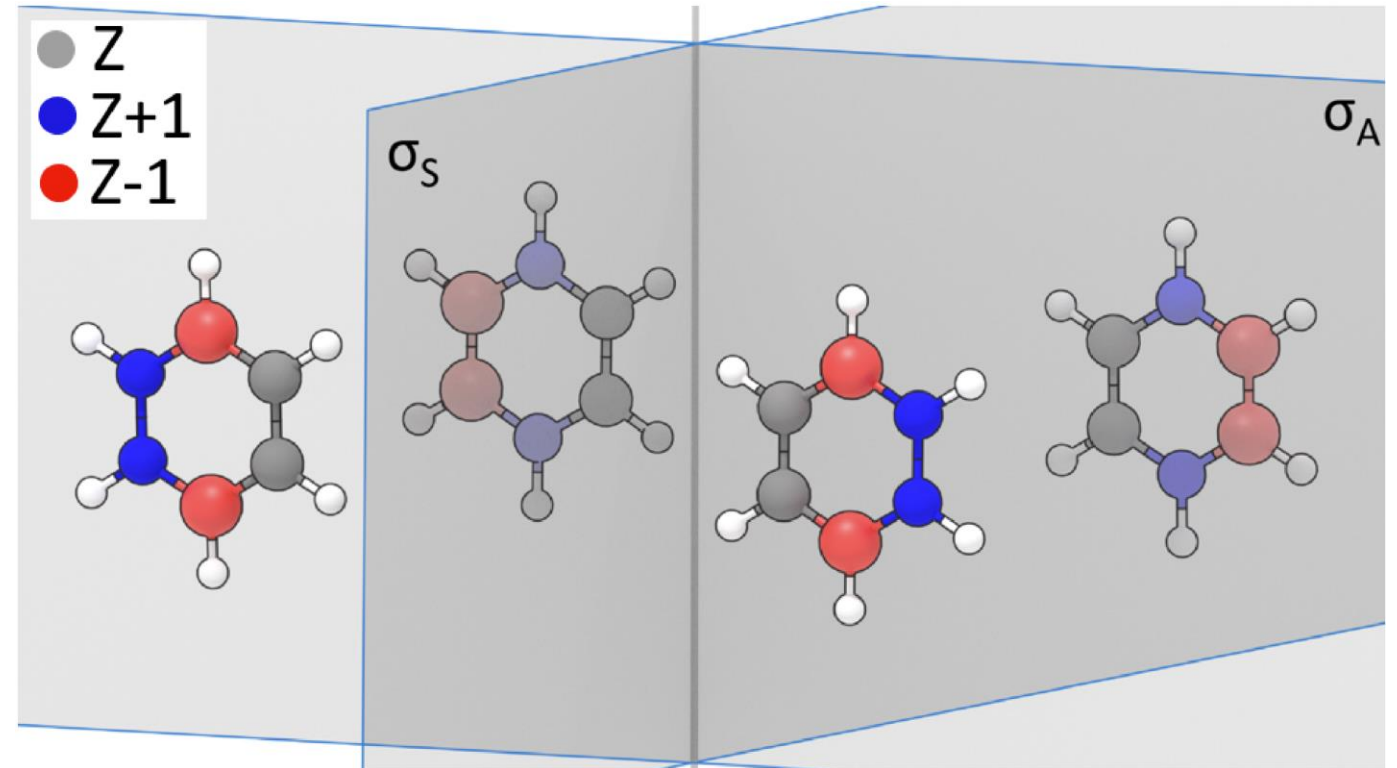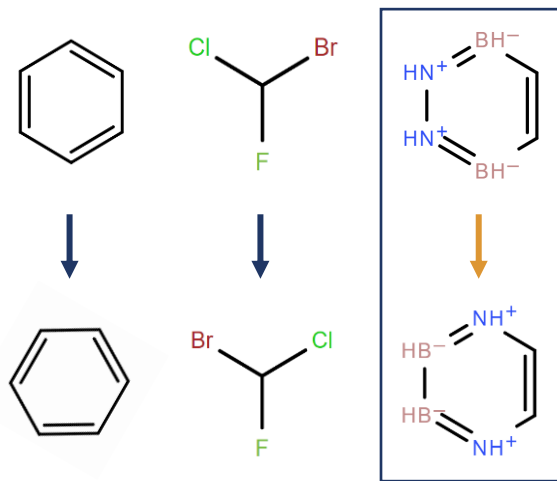**QA: 80.000x faster**

**BN-doped coronene dimer**
- Identify most/least attractive doping pattern
- Design case



QA: 20.000x faster

$2.8 \cdot 10^{10}$ targets

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

## Fundamentally new symmetry
Electronic energy only



-846.913     $\sigma_A$     -846.912

## 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



ML+QA: half the data

MAE [meV/atom]

800

400

200

100

50

25

— aCM
— aCMs

128    512   2048
*N*

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

x 414 M

QA: Millions at once!

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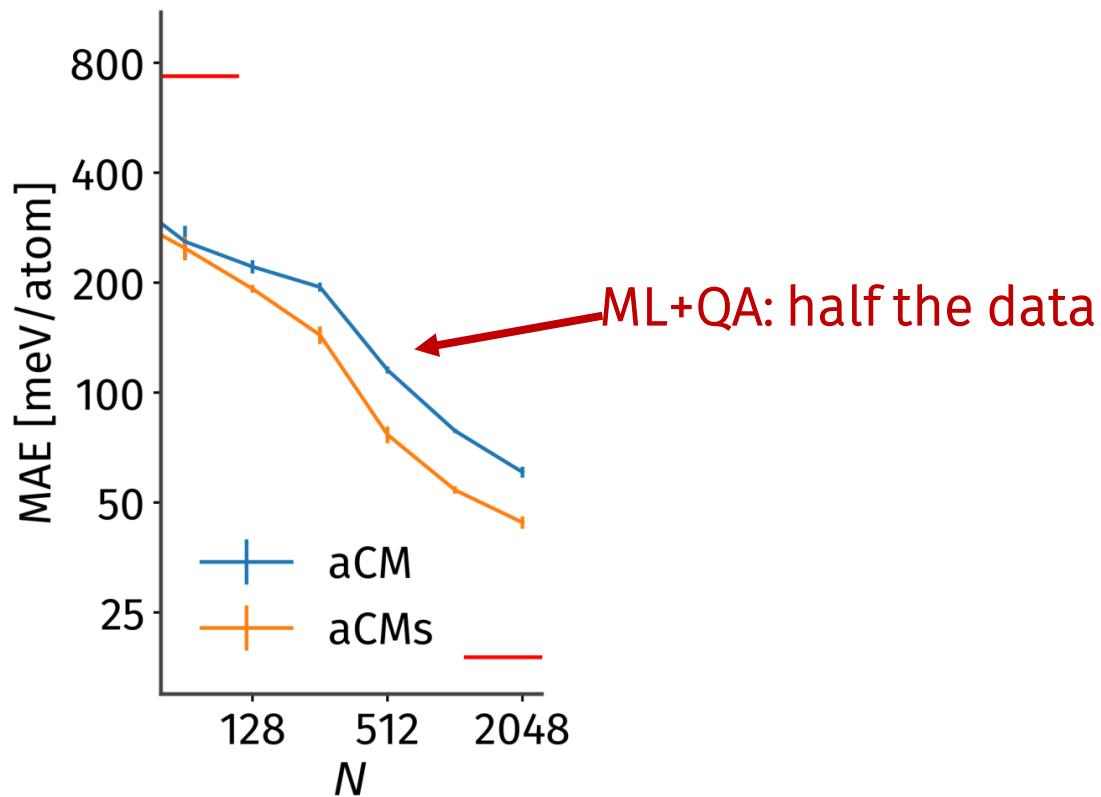- Add BN pairs
- Maximize CC bonds
- Substitute sites shared between rings
- Maximize BN bonds
- Avoid N substitutions on rings sharing a larger amount of bonds with other rings
- Balance BN substitutions in each ring

Not a single QM calculation required!

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

**Efficient |** Re-use knowledge, no one-by-one

**Symmetries |** Reducing ("folding") search space

**Constraints? |** Exclude regions of interest

**Differentiable Chemistry? |** Arbitrary derivatives

**Representations? |** Better data efficiency

**Ensembles? |** Derivatives on dynamic observables

**Thanks**
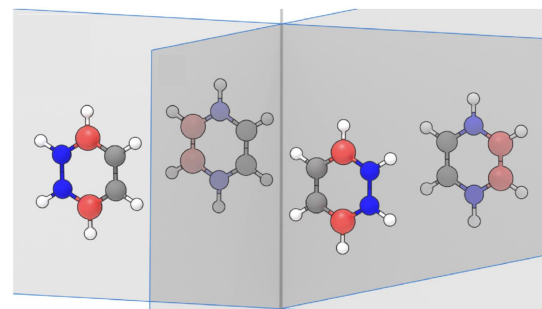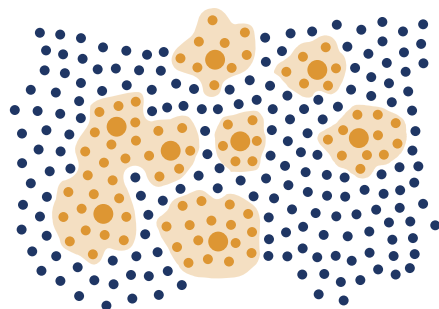Marco Bragato
Giorgio Domenichini
Emily Eikey
Stefan Heinen
Chasz Griego
Konstantin Karandashev
John Keith
Mario Krenn
Simon Krug
Dominik Lemm
Anatole von Lilienfeld
Alex Maldonado
Michael Sahre
Max Schwilk
Enrico Tapavicza
Jan Weinreich

ferchault    @ferchault    nablachem.org

Interpolate between molecular isoelectronic Hamiltonians

$$\hat{H}(\lambda) \equiv \lambda \hat{H}_{\mathrm{t}} + (1 - \lambda)\hat{H}_{\mathrm{r}} \qquad \lambda \in [0, 1]$$

Taylor expansion around reference molecule

$$E_{\mathrm{t}} = \sum_{n=0}^{\infty} \frac{1}{n!} \frac{\partial^n}{\partial \lambda^n} \left\langle \psi_\lambda \left| \hat{H}(\lambda) \right| \psi_\lambda \right\rangle \bigg|_{\lambda=0} = E_{\mathrm{r}} + \sum_{n=1}^{\infty} \frac{1}{n!} \frac{\partial^n E(\lambda)}{\partial \lambda^n} \bigg|_{\lambda=0}$$

Hellmann-Feynman theorem

$$\partial_\lambda E = \left\langle \psi_\lambda \left| \hat{H}_{\mathrm{t}} - \hat{H}_{\mathrm{r}} \right| \psi_\lambda \right\rangle = \Delta E^{\mathrm{NN}} + \int_\Omega d\mathbf{r} \underbrace{(v_{\mathrm{t}}(\mathbf{r}) - v_{\mathrm{r}}(\mathbf{r}))}_{\equiv \Delta v} \rho_\lambda(\mathbf{r})$$

O. A. von Lilienfeld, *J. Chem. Phys.* **2009**.

Alchemical Perturbation Density Functional Theory (APDFT)

$$E_t = E_r + \Delta E^{\mathrm{NN}} + \int_\Omega d\mathbf{r} \sum_{n=0}^{\infty} \frac{1}{(n+1)!} \Delta v \frac{\partial^n \rho_\lambda(\mathbf{r})}{\partial \lambda^n}\bigg|_{\lambda=0}$$

$$\rho_t = \sum_{n=0}^{\infty} \frac{1}{n!} \frac{\partial^n \rho}{\partial \lambda^n}\bigg|_{\lambda=0}$$

- Gives consistent energies, densities, forces, …
- Uses the same derivatives for all predictions

ferchault/APDFT