# Computational Alchemy

Guido Falk von Rudorff, University of Vienna

ferchault    @ferchault    guido.vonrudorff.de

**Guido Falk von Rudorff**

| | | |
|---|---|---|
| BSc, Msc | Freie Universität Berlin | Atom clusters, force fields |
| PhD | London | Electronic structure calculations |
| PostDoc | Basel, Vienna | Alchemy and machine learning |

What about you?
- Background
- Why interest in computational alchemy?
- Aims for this lecture
- Experience with classical or electronic structure calculations?
- How comfortable with Python?

Each Thursday:
- 10:15-11:00		Lecture
- 11:15-12:00		Lecture
- 12:15-13:00		Exercise (starting next week)

Exercises:
- 10 sets, two exercises each
- 40 points in total, 20 points needed

Setting:
- Interactive and informal: less lecture, more conversation

Joseph Wright, 1771

- Why do we care?
- Classical alchemy: free energy calculations
- Quantum alchemy
- High-level introduction to quantum mechanical calculations
- High-level introduction to numerical derivatives and automatic differentiation
- Alchemical Perturbation Density Functional Theory
- Fundamental consequences:
    - Atoms in molecules
    - Alchemical Enantiomers
    - Alchemical normal modes
- Applications
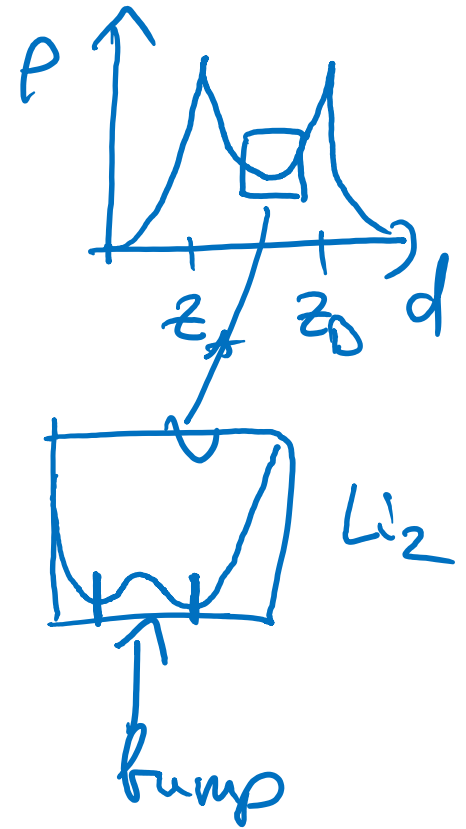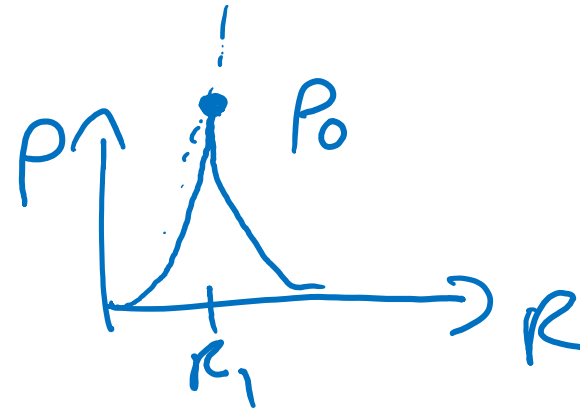    - Electrostatic potentials
    - Non-covalent interactions

$$E \approx -Z^{7/3}$$

- New tool: interpolation in chemical space
  - Machine learning
- Numerical differentiation/integration
  - Computer Algebra Systems
- Automatic differentiation
  - Machine learning
  - Optimization techniques
- Better understanding of electronic structure calculations

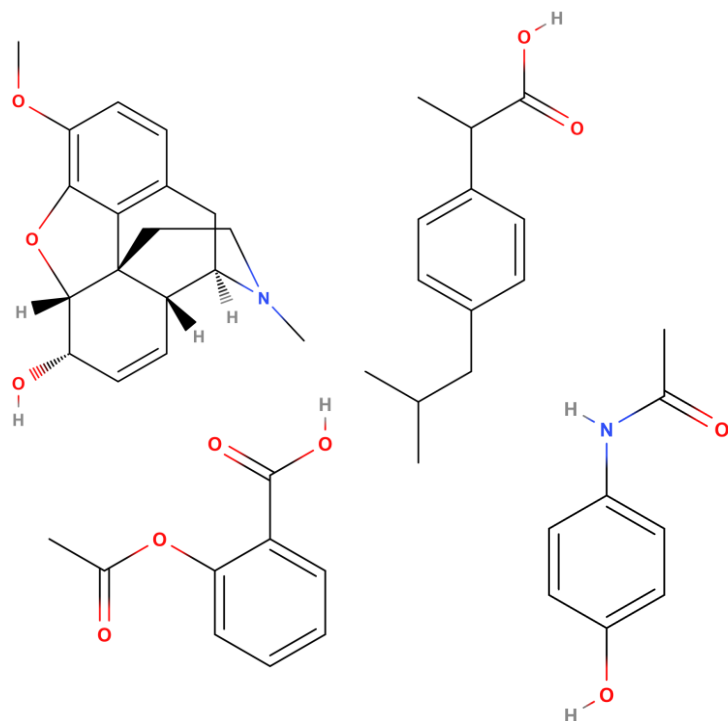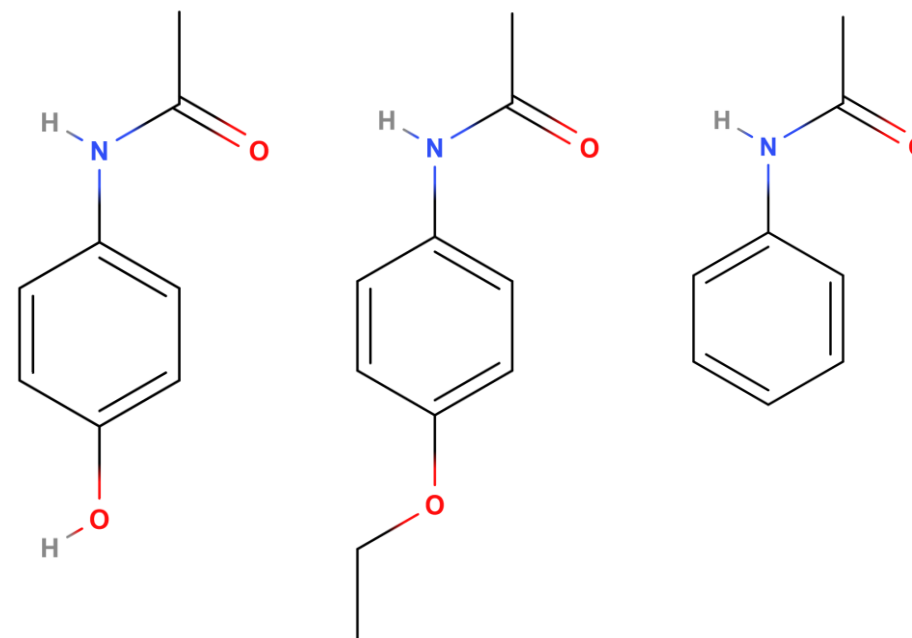- $V$ representability $\quad p \rightsquigarrow V$

- Kato's cusp

$p$ ↑

$P_0$

$R_1$ → $R$

$p$ ↑

$z$ $z_0$ $d$

$Li_2$

pump

# Motivation

Materials / compound design efforts face a vast search space
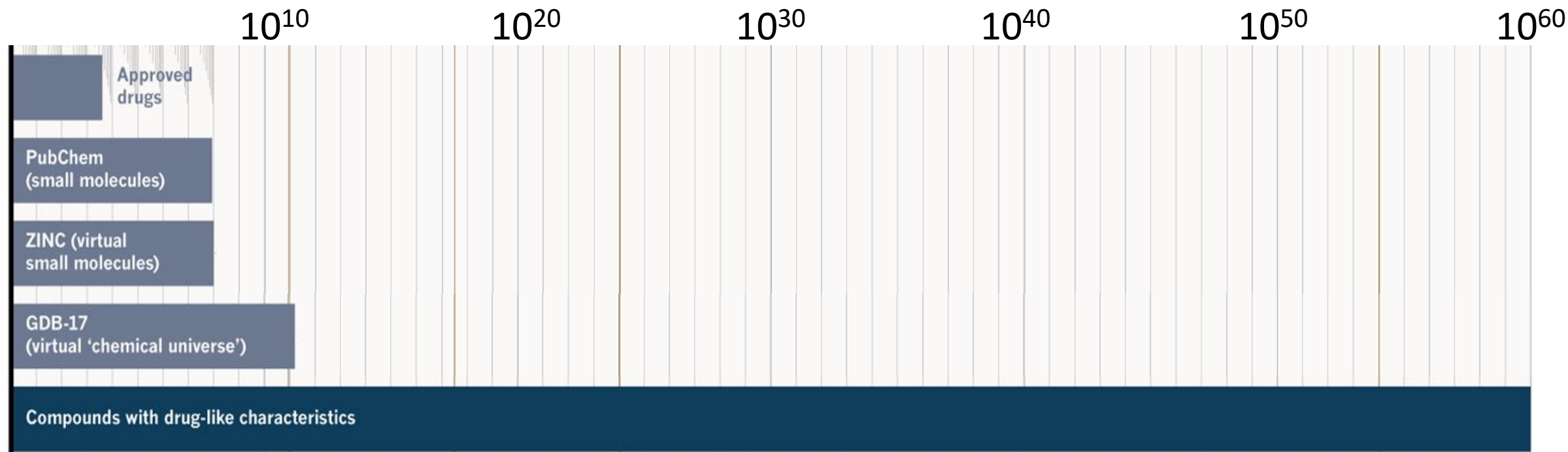


**Global Search Problem**
*Which class of compounds?*

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

$10^{10}$    $10^{20}$    $10^{30}$    $10^{40}$    $10^{50}$    $10^{60}$

Approved drugs

PubChem (small molecules)

ZINC (virtual small molecules)

GDB-17 (virtual 'chemical universe')

Compounds with drug-like characteristics

**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

Mullard, *Nature,* 2017.

Necessity
- Only way to cover problem size
- Still open to systematic evaluation
- Often used as prefiltering step

- Complicated chemistry
- Tricky / error-prone reference calculations

Convenience
- Can be done more accurately
- Uneconomical/cumbersome reference method
- Often used as direct but optional substitute

- Standard energy calculations of well-behaved systems
- Semi-emipirical level sufficient

Restriction
- Subspaces: conformers, constitutional isomers
- Configurations: minima, transition states

Interpolation
- Set of molecules as fixed reference
- Define interpolant
- Small data sets: e.g. KRR
- Large data sets: e.g. NN

Expansion
- Molecules get perturbed
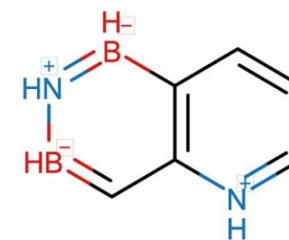- Quantum Alchemy e.g. APDFT
- Taylor expansion w.r.t. elements

$$E_{\mathrm{ML}}$$

$$Z_I, \mathbf{R}_I, N_e, \sigma$$

- Computational alchemy
  True functional value is always defined / obtainable

- Mathematical interpolation
  True value is a model construct and depends on the model employed (s)

Smoothly connecting discrete states

- Classical     thermodynamic integration
  interpolate parameters representing elements

- Quantum     no interpolant needed in many directions (Schrödinger's equation)
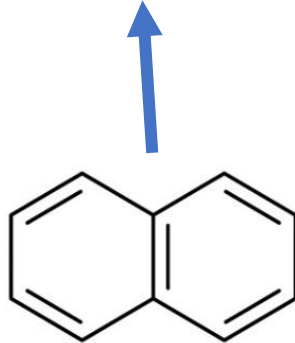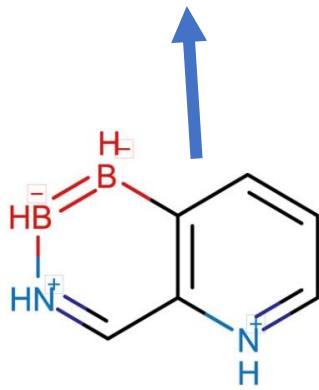  interpolate coordinates, nuclear charges (not: number of electrons)

- Force field describes material/molecule
- No electronic structure, but physical effects



Hydrophobic effect is roughly proportional to surface area

water
Continuum solvent model

torsion angle

Lennard Jones

Distance

bond length or 3-atom angle

https://commons.wikimedia.org/wiki/File:MM_PEF.png

- Hamiltonian describes electronic structure
- Mixing Hamiltonians creates paths

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

Physical meaning?

$$-\sum_i \frac{\hbar^2}{2M_i} \nabla^2_{\mathbf{R}_i}$$

$$-\sum_i \frac{\hbar^2}{2m_e} \nabla^2_{\mathbf{r}_i}$$

$$-\sum_i \sum_j \frac{Z_i e^2}{4\pi\epsilon_0 |\mathbf{R}_i - \mathbf{r}_j|}$$

$$\frac{1}{2}\sum_i \sum_{j\neq i} \frac{e^2}{4\pi\epsilon_0 |\mathbf{r}_i - \mathbf{r}_j|} = \sum_i \sum_{j>i} \frac{e^2}{4\pi\epsilon_0 |\mathbf{r}_i - \mathbf{r}_j|}$$

$$\frac{1}{2}\sum_i \sum_{j\neq i} \frac{Z_i Z_j e^2}{4\pi\epsilon_0 |\mathbf{R}_i - \mathbf{R}_j|} = \sum_i \sum_{j>i} \frac{Z_i Z_j e^2}{4\pi\epsilon_0 |\mathbf{R}_i - \mathbf{R}_j|}.$$

**What could possibly go wrong?**

**How hard can it be?**
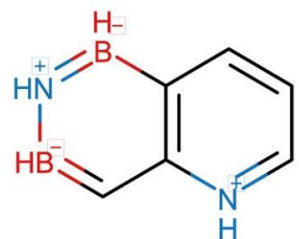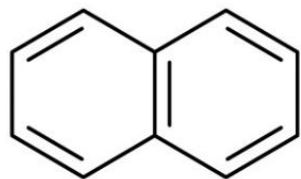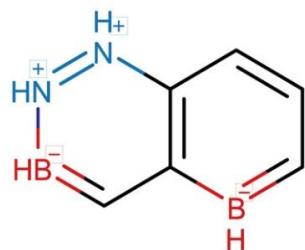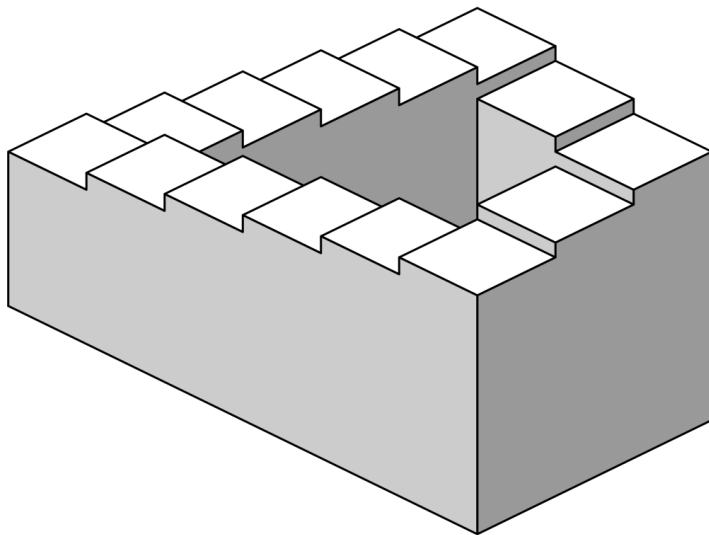
- State function: uniquely defined for a system = $Z_I, \mathbf{R}_I, N_e, \sigma$

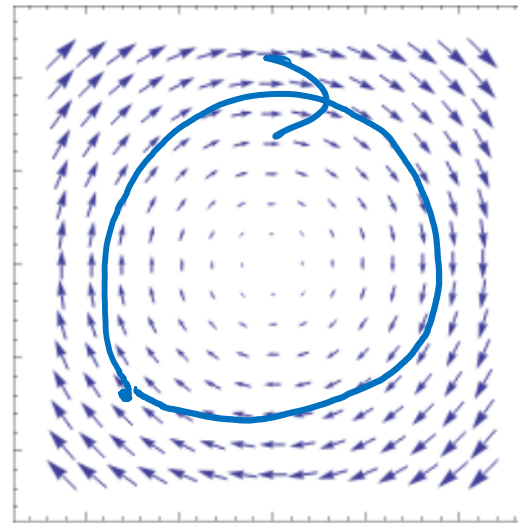Energy, and any other observable (!)  Other examples?

- State function: uniquely defined for a system
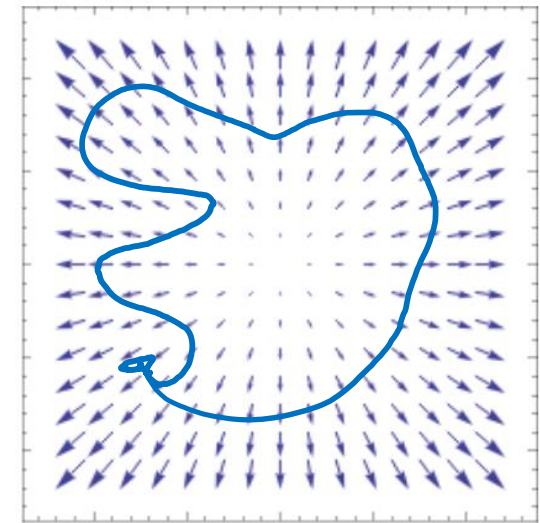- Integrals being path-independent*



non-conservative        conservative

$$\int_{C_1} \mathbf{F} \cdot d\mathbf{s} \neq \int_{C_2} \mathbf{F} \cdot d\mathbf{s}$$

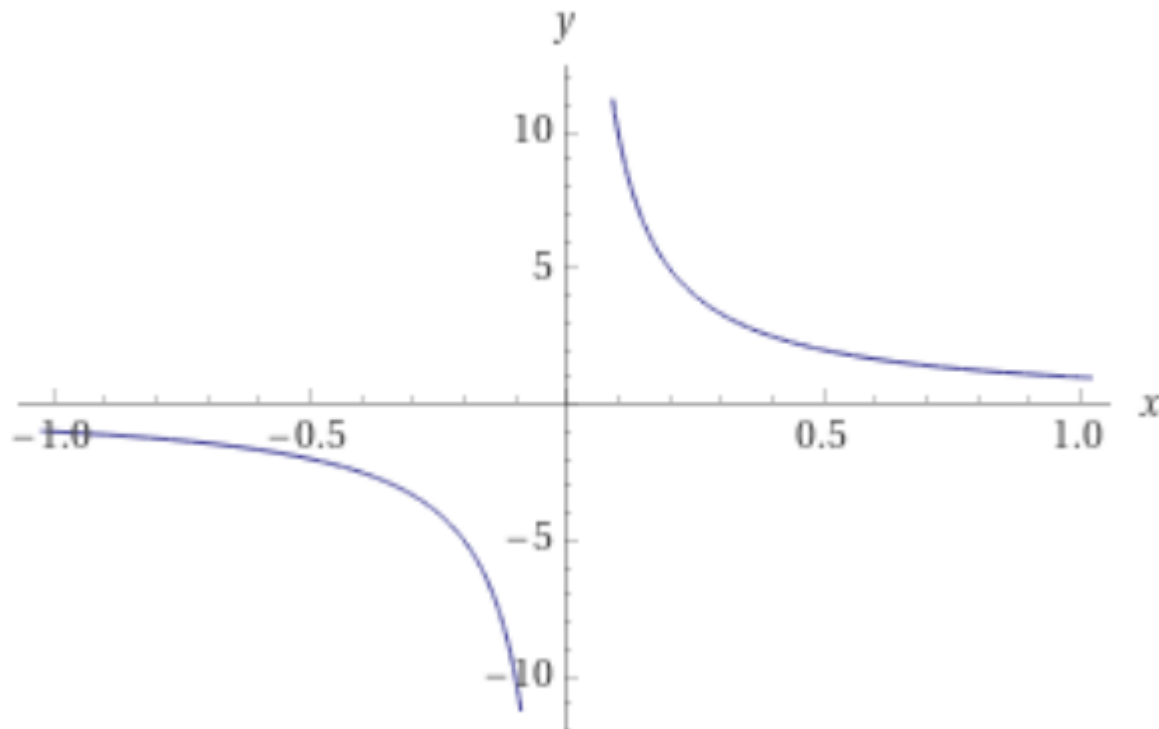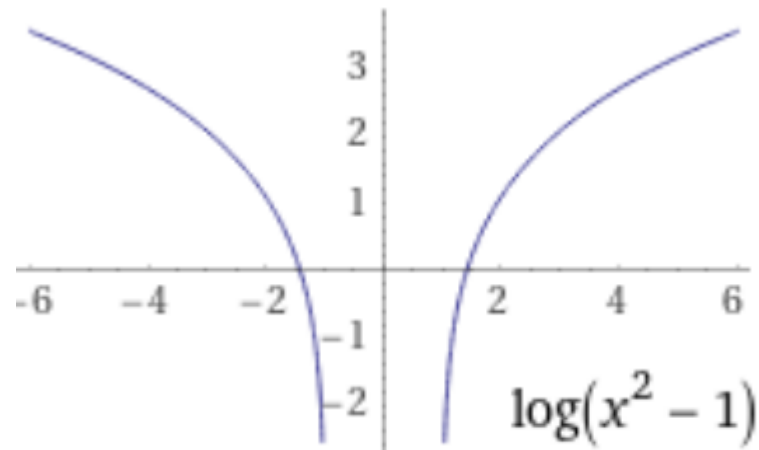- State function: uniquely defined for a system
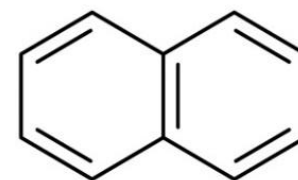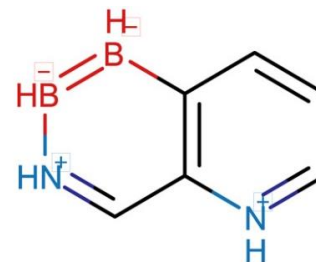- Integrals being path-independent*
- Smoothness

- State function: uniquely defined for a system
- Integrals being path-independent*
- Smoothness
- Defined for the relevant domain



$$\log(x^2 - 1)$$

Non-integer values are no problem

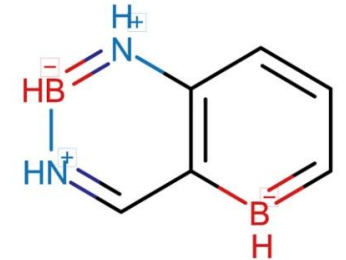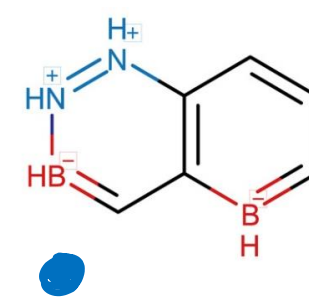Fundamental theorem of calculus

$$F(x) = \int_a^x f(t)\, dt.$$

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

Fundamental theorem of calculus

$$F(x) = \int_a^x f(t)\, dt.$$
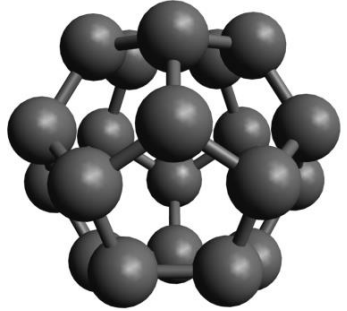
$$\hat{H}(\lambda) \equiv \lambda \hat{H}_\mathrm{t} + (1 - \lambda)\hat{H}_\mathrm{r}$$

- Experimentally observable
- Energetically stable
- Simple relationships / easy paths

Covalent Energies



Non-covalent Interactions



Deprotonation Energies



Energy Decomposition

# Example
26

Alchemical Perturbation Density Functional Theory (APDFT)
Uses calculations of *one* molecule to estimate *many* molecules

$$\frac{\partial E}{\partial z_1} = \frac{\partial E}{\partial z_3}$$

$$\frac{\partial^2 E}{\partial z_1 \partial z_3}$$



$$E, \rho, \{\partial_\lambda^i \rho\} \rightarrow \{E_i\}, \{\rho_i\}, \{F_i\}, \{\mu_i\}, \{Q_i\}, \dots$$

1 system $\rightarrow$ Millions of systems

- BN-doped $C_{20}$
- APDFT2: one SCF, one derivative
- Key trick to scale with chemical space:

$$\frac{\partial \rho}{\partial \lambda} = \sum_I \frac{\partial \rho}{\partial Z_I} \frac{\partial Z_I}{\partial \lambda}$$

$C_{20}$

$3.1 \cdot 10^6$ targets





GFvR, O. A. von Lilienfeld, *Phys. Rev. Res.* **2020** (*arXiv* 1809.01647).

- BN-doped coronene dimer
- APDFT2: one SCF, 24 derivatives
- SCAN+VV10

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



2.8 · $10^{10}$ targets

GFvR, O. A. von Lilienfeld, *Phys. Rev. Res.* **2020** (*arXiv* 1809.01647).

1    A tweet used to be up to 140 characters. Summarize the key idea and relevance of computational alchemy in 140 characters or less.

2    a) The fundamental theorem of calculus allows us to obtain one function value by integrating the derivative of that function along a path. In practise, this is done by evaluating the function derivative at a finite number of equidistant points $n$ along the path.
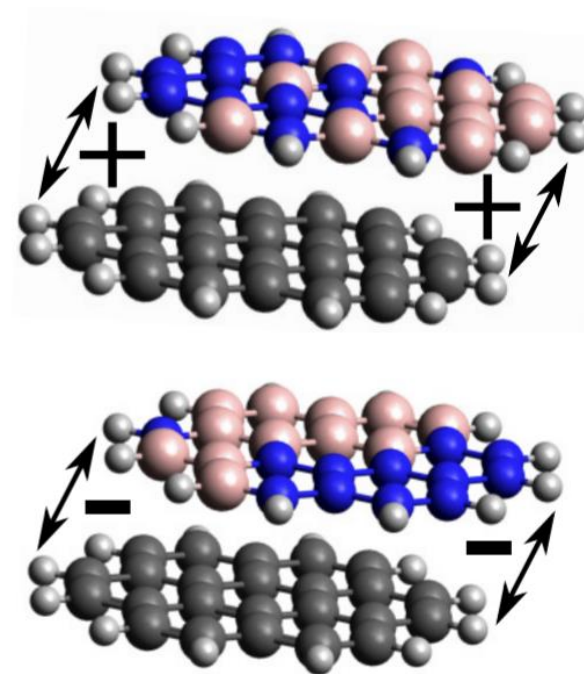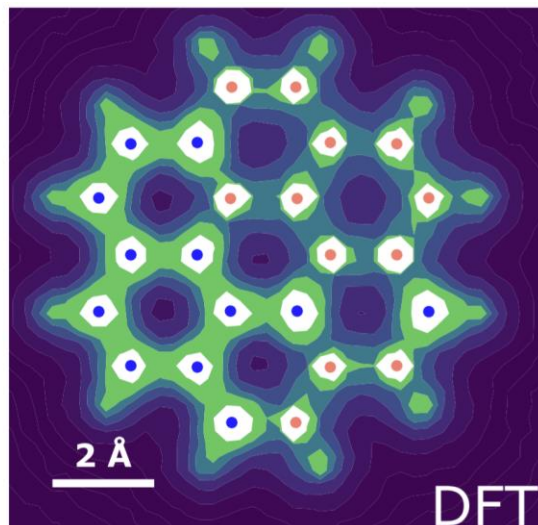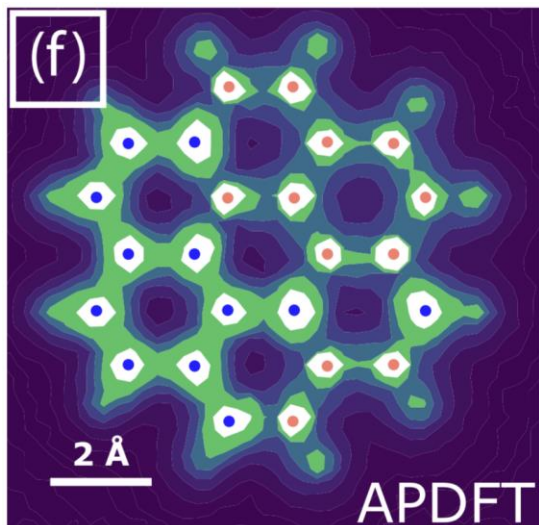
Why can this method yield wrong results?
Can you sketch a function where the integral as obtained from $n=3$ and $n=5$ differs in sign?
Which features in an integrand would be problematic for this approach and why?

b) In classical calculations, a common application of computational alchemy is to grow a molecule (e.g. benzene) inside an existing environment (e.g. water or a membrane). This is done by turning on the interactions between a static benzene and environment molecules, while keeping the interactions between the environment molecules unchanged.
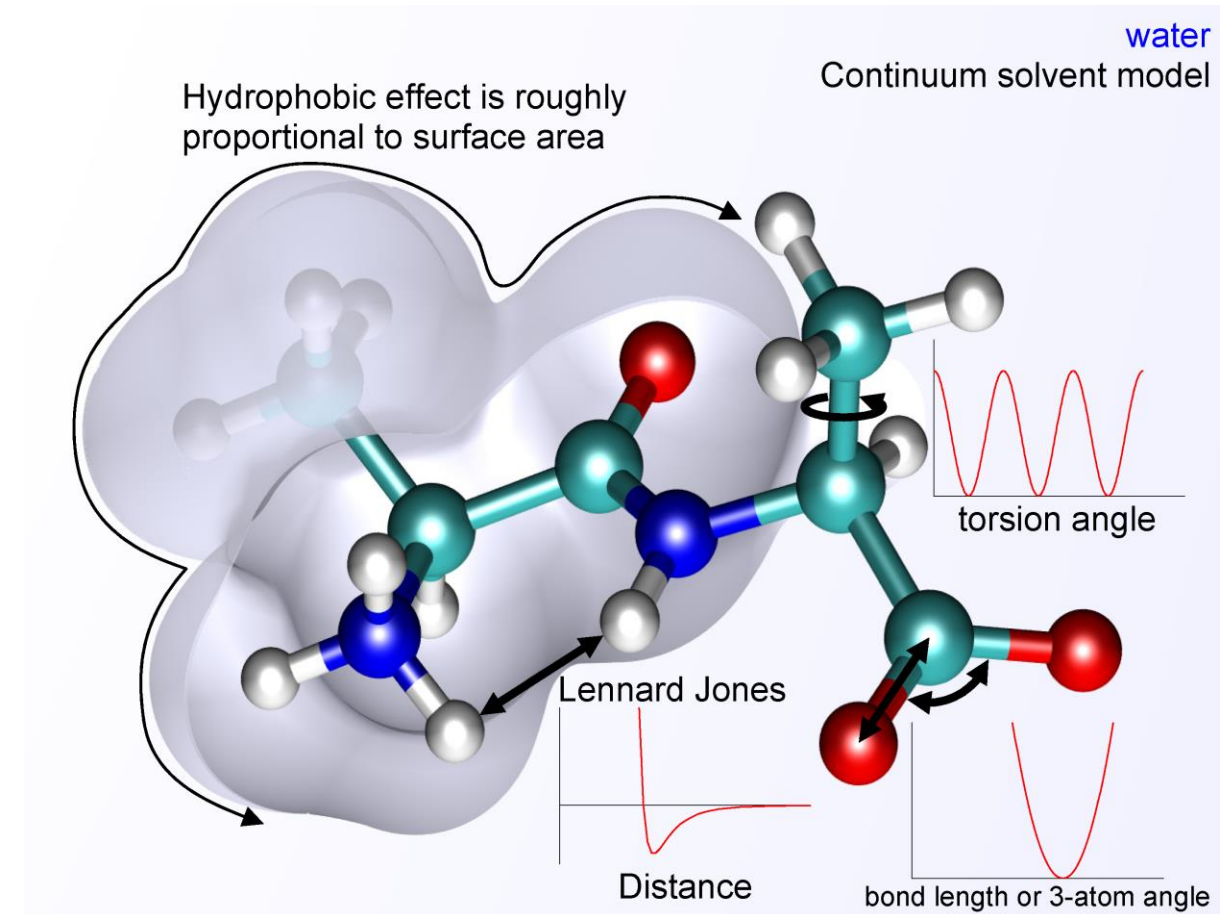
Why does this have to be done slowly?

# Classical Alchemy: Theory

- Free energy is a state function
- Free energy differences
  … drive action
  … are hard to calculate (if at all)
- Following alchemical paths yields such differences

**Context**
- Classical force fields (all-classical description)
  - Bond, angle, dihedral
  - Non-bonded: Lennard Jones
- Interactions get scaled
  - Masses stay constant
  - Charges are scaled
- No need to be linear in interpolation
- Often seen within a
  molecular dynamics trajectory
- No charge conservation required!



Hydrophobic effect is roughly proportional to surface area

water
Continuum solvent model

torsion angle

Lennard Jones

Distance

bond length or 3-atom angle

https://commons.wikimedia.org/wiki/File:MM_PEF.png

alchemical

cheap

slow, expensive

slow, expensive

cheap

- Two states: i, j
- NVT ensemble, equilibrated

$$\Delta A_{ij} \equiv -k_B T \ln \frac{Q_j}{Q_i}$$

$$Q_i \equiv \int_{\Gamma_i} \exp\left[-\frac{U_i(\vec{q})}{k_B T}\right] d\vec{q}$$

Helmholtz free energy

Temperature

Partition function

Phase space

Potential energy

- Need: ensemble average
- Ergodicity: Time average = Ensemble average
- Propagate in small steps (~fs) explicit positions from repeated force evaluations

- Challenges
  - Numerically stable
  - Time reversibility
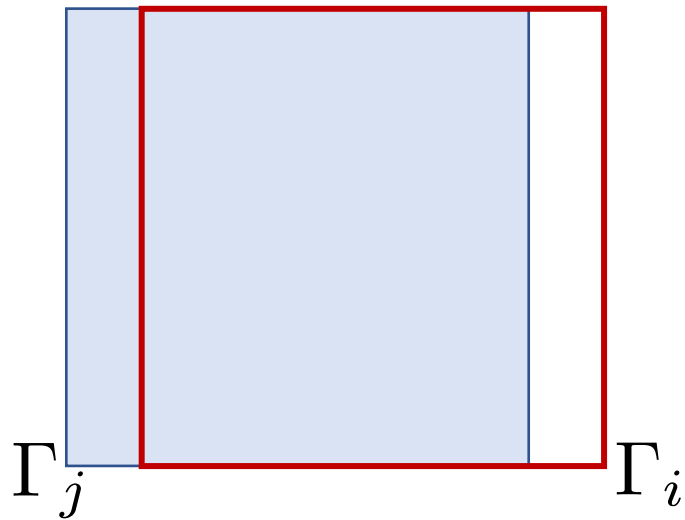  - Thermostats / Barostats
  - Equilibration

Questions:
- Why are classical calculations typically time-reversible but quantum mechanical calculations are not?
- What does it mean if a setup is lacking time-reversibility?

- Two states: i, j
- NVT ensemble

$$Q_i \equiv \int_{\Gamma_i} \exp\left[-\frac{U_i(\vec{q})}{k_B T}\right] d\vec{q}$$

$\Gamma_j$   $\Gamma_i$

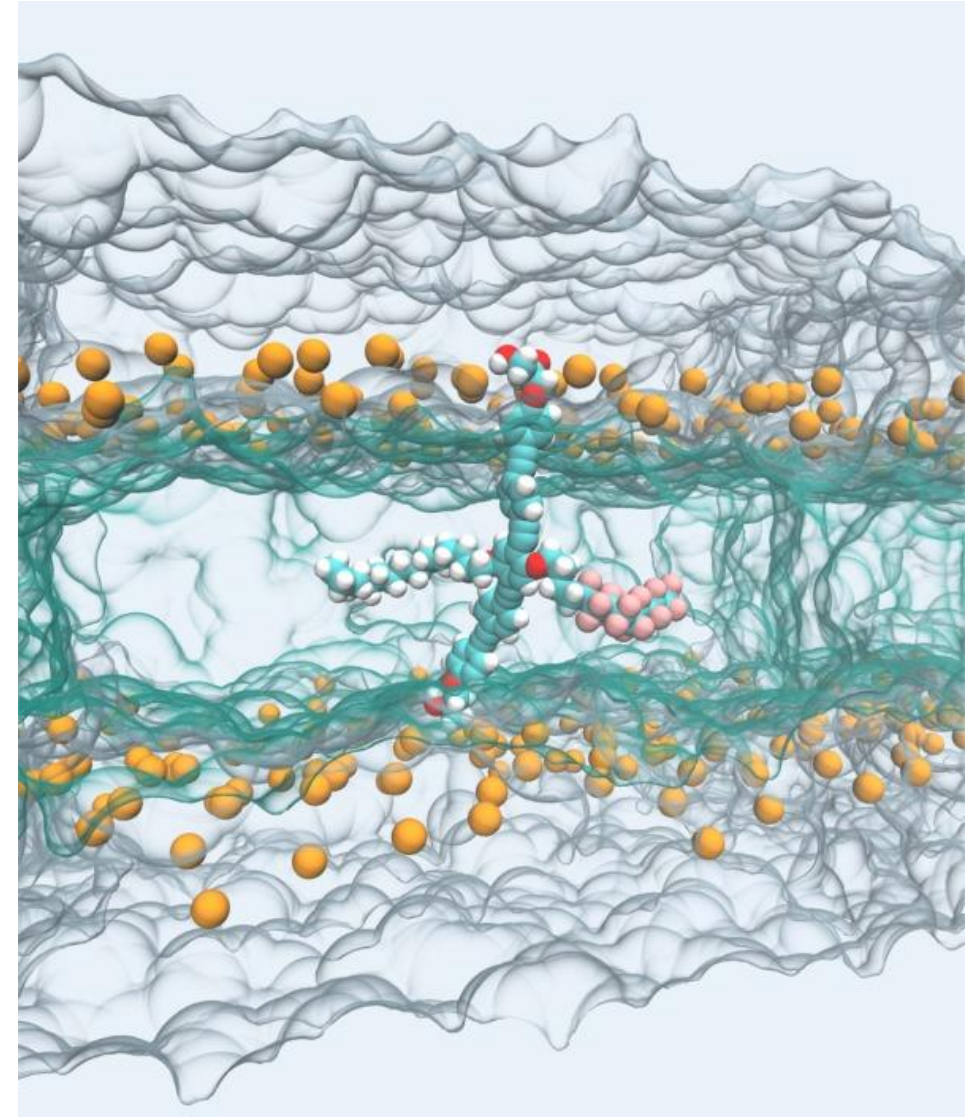Partition function    Phase space    Potential energy

Example: hard spheres with different radii: close interaction never happens

Insert molecule in membrane: turn on interactions

$$V_{\mathrm{LJ}}(r) = 4\varepsilon \left[ \left(\frac{\sigma}{r}\right)^{12} - \left(\frac{\sigma}{r}\right)^{6} \right]$$

Questions:

1. How to "turn on" interactions?
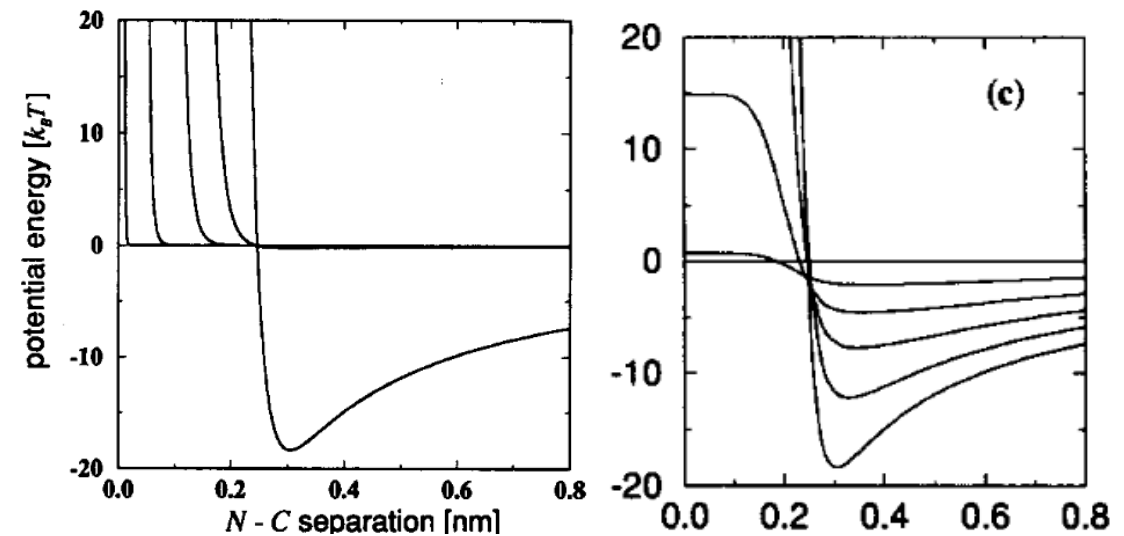2. Why might that fail?
3. What happens if that fails?



GFvR, T. Watermann, X.-Y. Guo, D. Sebastiani, *J. Comput. Chem.* **2017.**

- No matter the scaling: unbounded energy

$$V_{\mathrm{LJ}}(r) = 4\varepsilon \left[ \left( \frac{\sigma}{r} \right)^{12} - \left( \frac{\sigma}{r} \right)^{6} \right]$$

- Unbounded energy: no reliable derivatives
  - Question: Why is that an issue in molecular dynamics?
- Solution: soft-core potentials

*T.C. Beutler et al. / Chemical Physics Letters 222 (1994) 529–539*

$$U(\lambda, r) = 4\epsilon\lambda^{n} \left[ \left( \alpha(1-\lambda)^{m} + \left( \frac{r}{\sigma} \right)^{6} \right)^{-2} - \left( \alpha(1-\lambda)^{m} + \left( \frac{r}{\sigma} \right)^{6} \right)^{-1} \right]$$
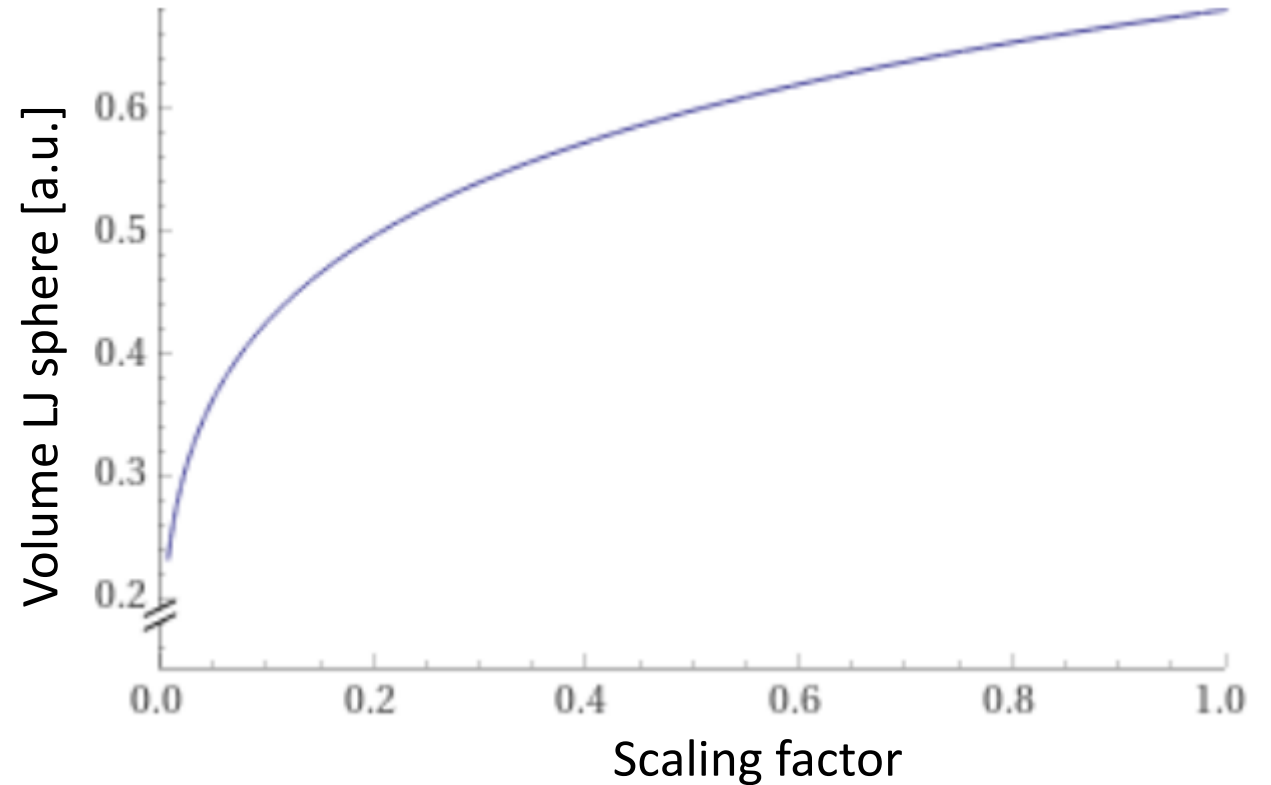
- Typical molecules: effective charges on each site + Lennard Jones

$$V_{\mathrm{LJ}}(r) = 4\varepsilon \left[ \left( \frac{\sigma}{r} \right)^{12} - \left( \frac{\sigma}{r} \right)^{6} \right]$$

- Can be scaled independently
    - Energies remain state function of parameters
- Caveat:
    - If LJ is scaled: charges can get closer to each other. If charges are of opposite sign: trapping
    - Therefore: electrostatics first, LJ second
- Question: Would separate paths be acceptable and if so, why?

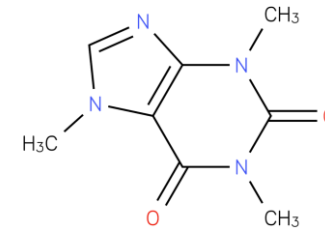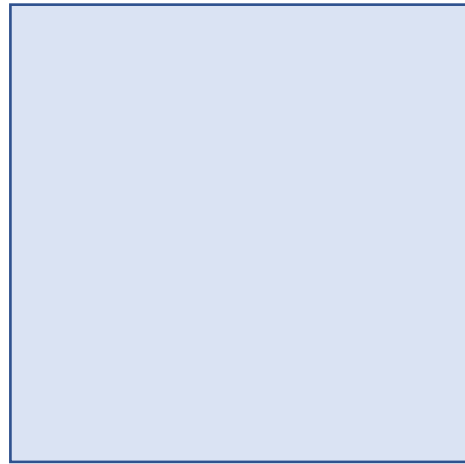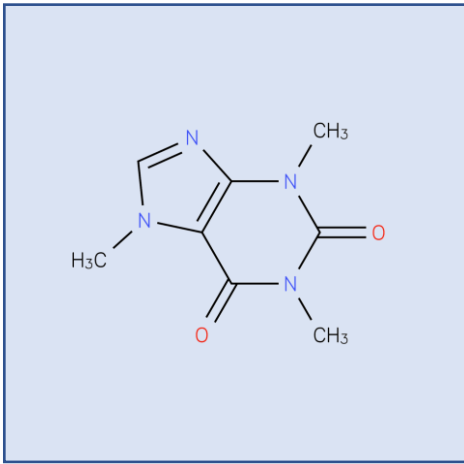Linear pathway not necessarily efficient / rarely "effectively linear"

$$V_{\mathrm{LJ}}(r) = 4\varepsilon \left[ \left(\frac{\sigma}{r}\right)^{12} - \left(\frac{\sigma}{r}\right)^{6} \right]$$

- Avoid constrained/restrained configurations                                    Why?


- Choose low-change path: large changes mean large derivatives          Why is that bad?


- Change parameters to create effectively linear results


- Restrict number of intermediates (=mixed states)


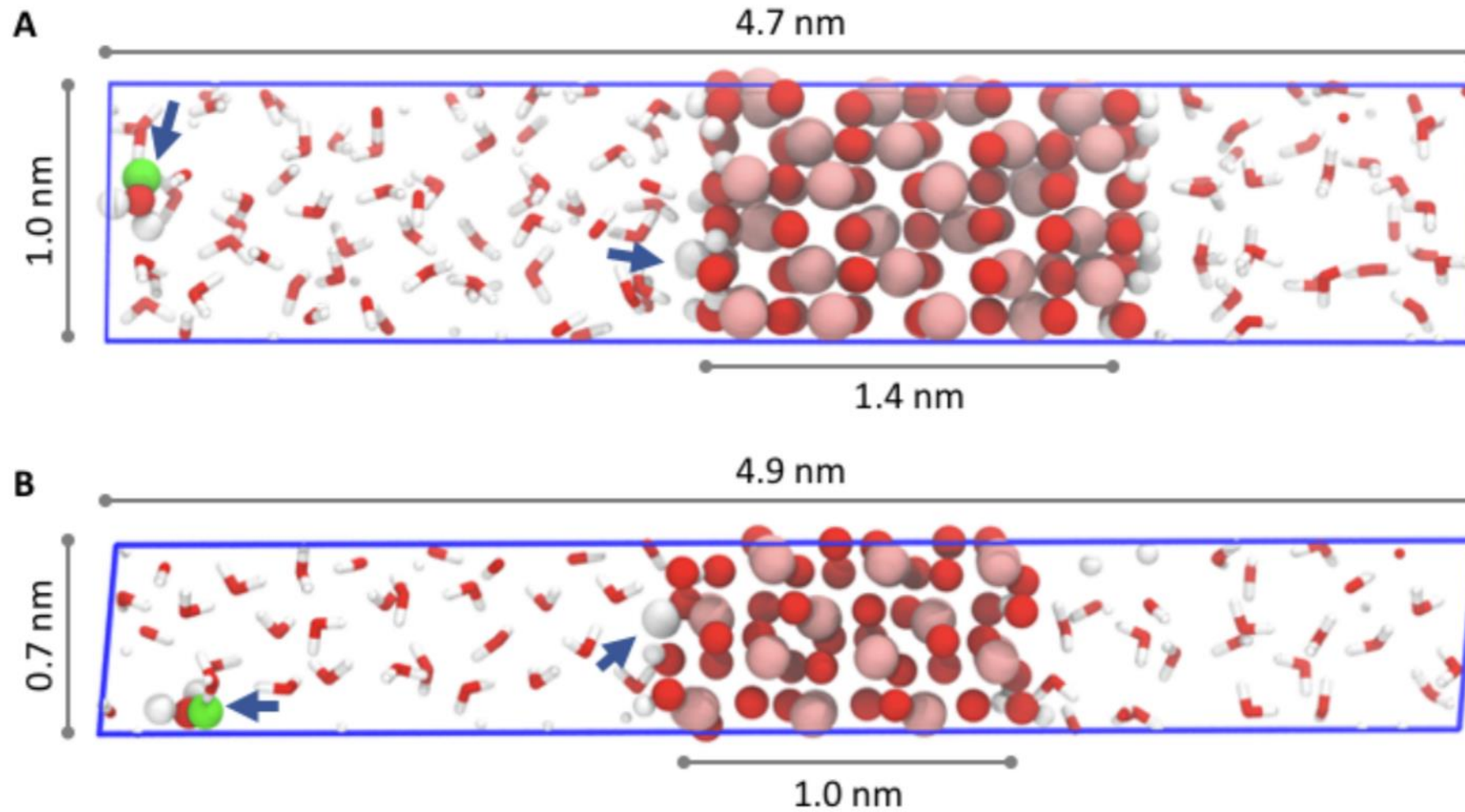- Beware electrostatics: keep net charge                                          Why?

Consider NVT:
Why is the free energy of solvation NOT simply the free energy differences with solute-solvent interactions turned off?
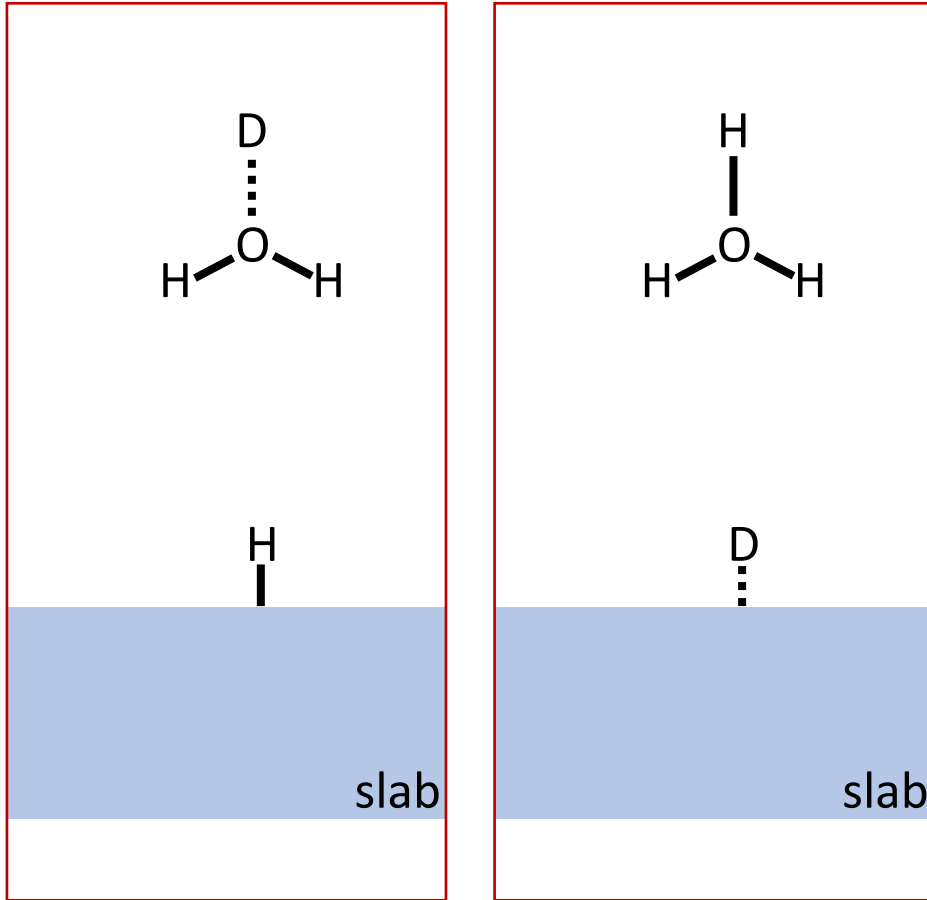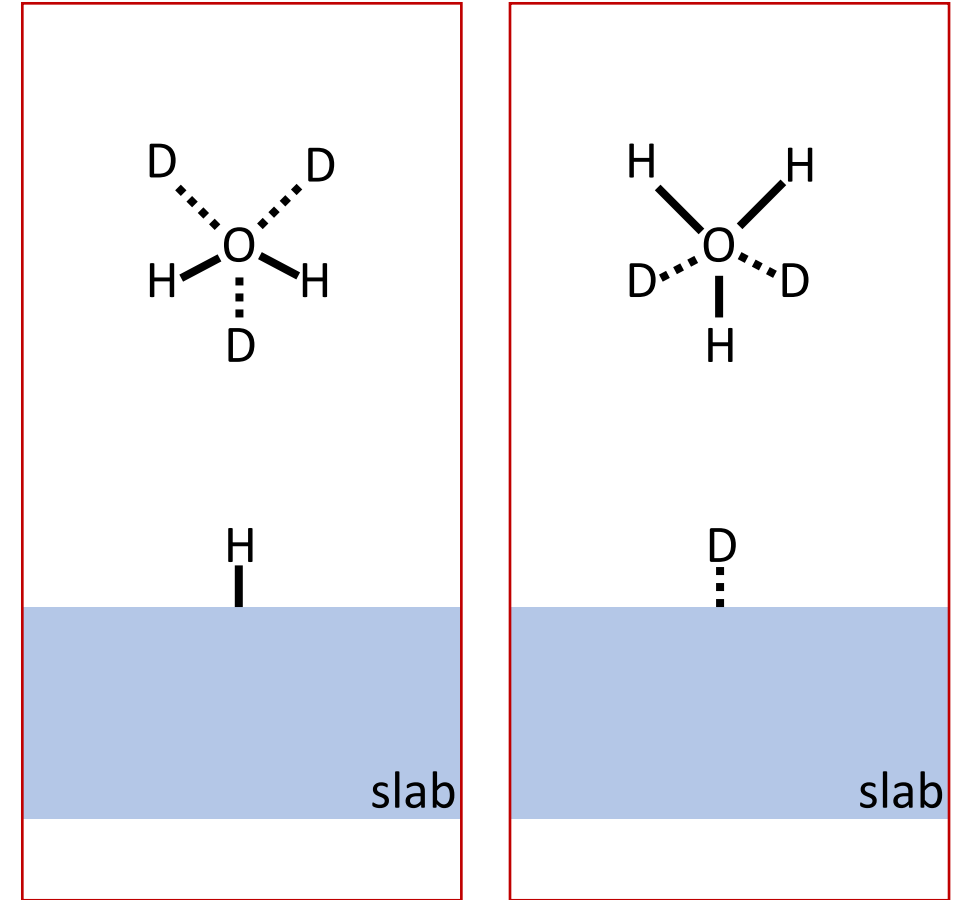
O. Gittus, GFvR, K. Rosso, J. Blumberger, *J. Phys. Chem. Lett.* **2018.**
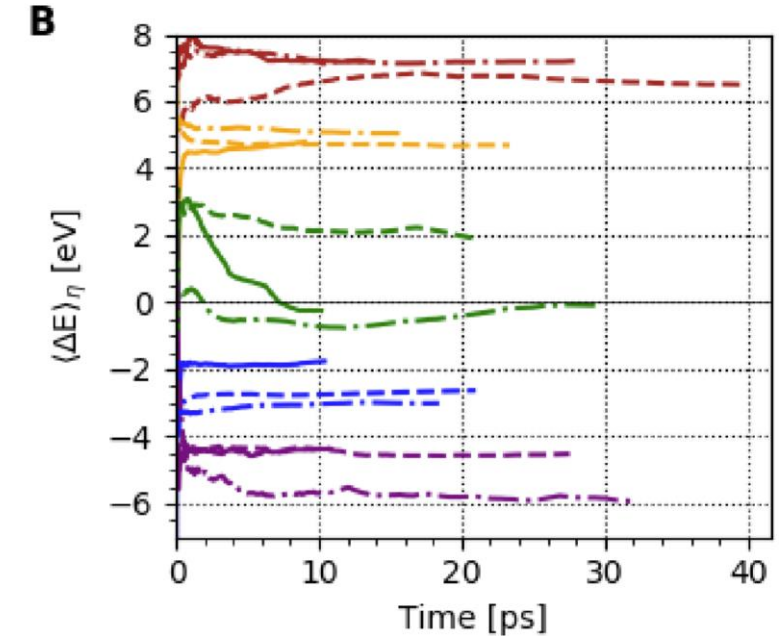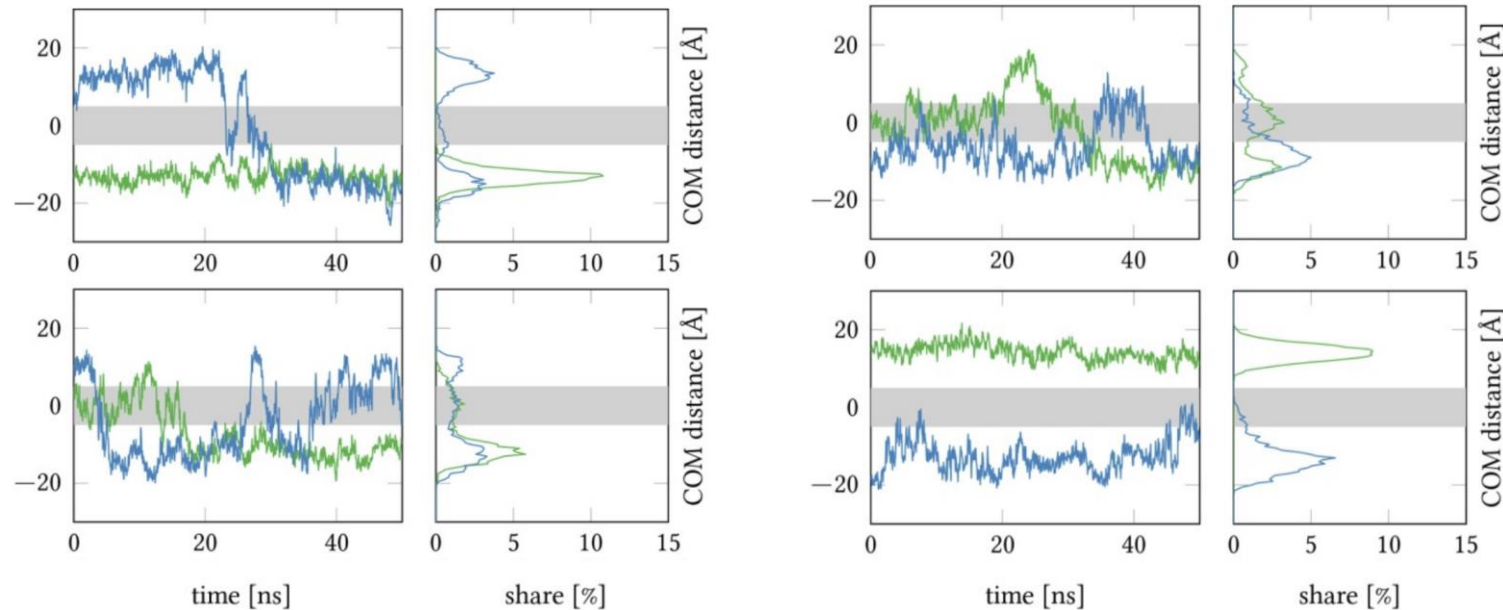
**Single topology approach**



**Double topology approach**

Conformational sampling slow – four independent trajectories
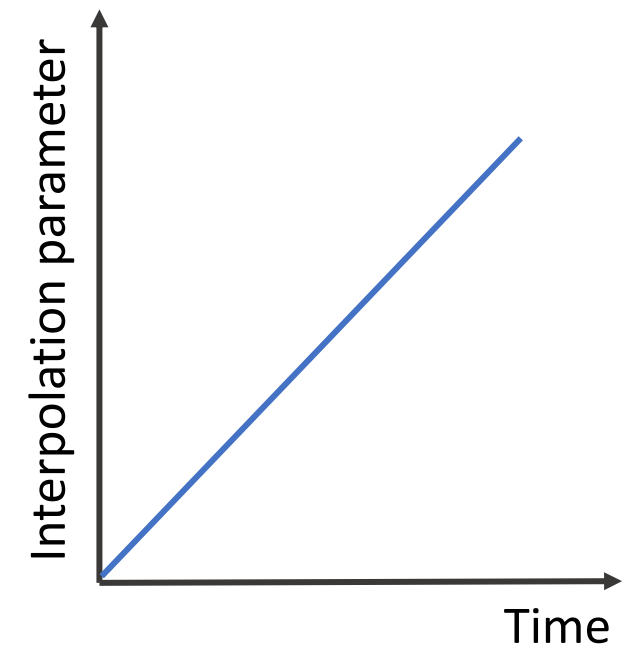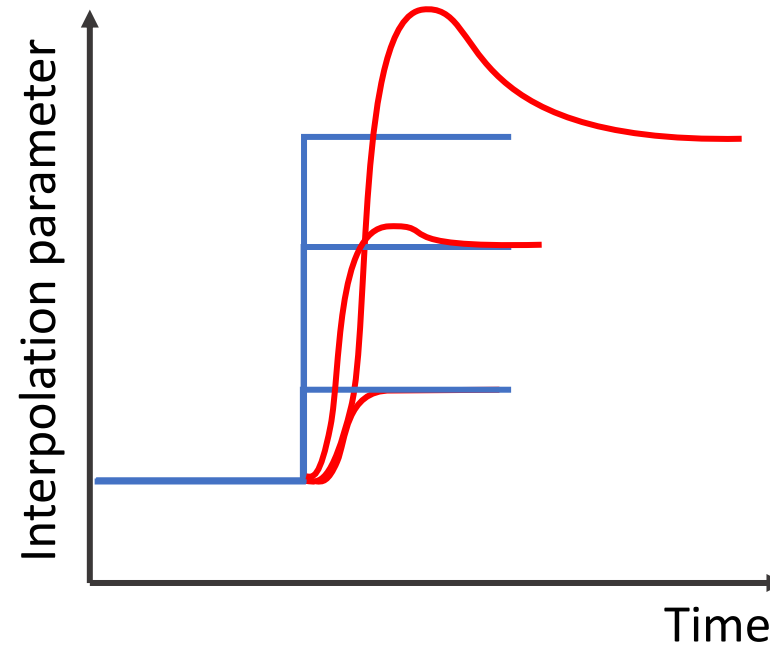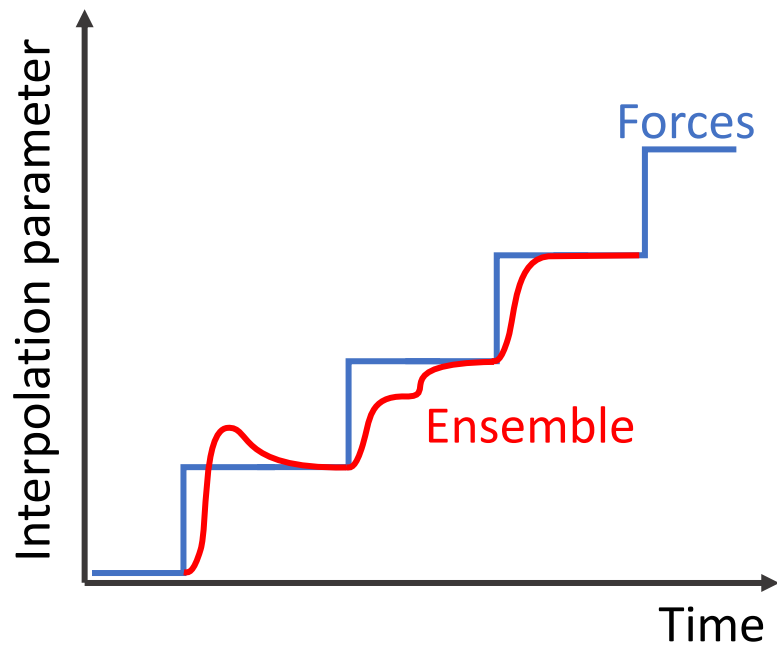


High auto-correlation

O. Gittus, GFvR, K. Rosso, J. Blumberger, *J. Phys. Chem. Lett.* **2018.**
GFvR, T. Watermann, X.-Y. Guo, D. Sebastiani, *J. Comput. Chem.* **2017.**

Molecular dynamics: propagation of configurations: auto-correlated, delayed



Questions:
1. Which limit is required?
2. Third panel: What does the ensemble do?

- Convergence hard to achieve in practise
  - How about being wrong twice?



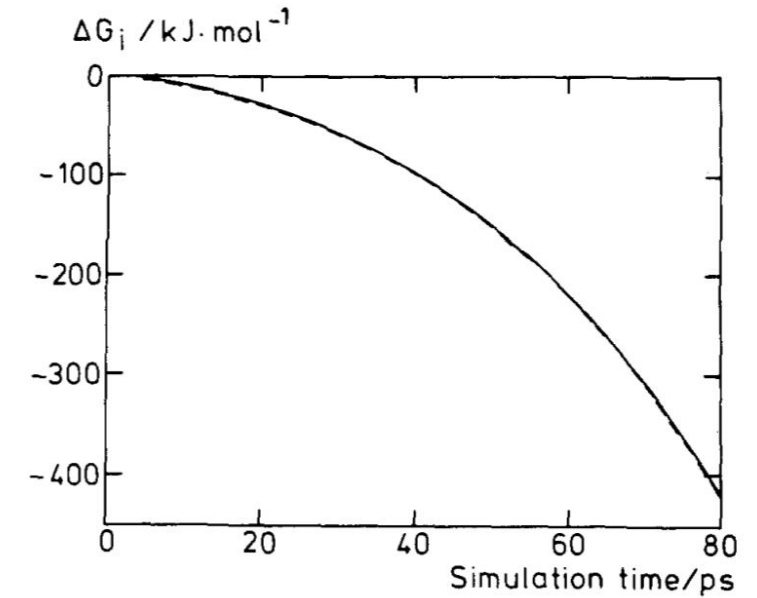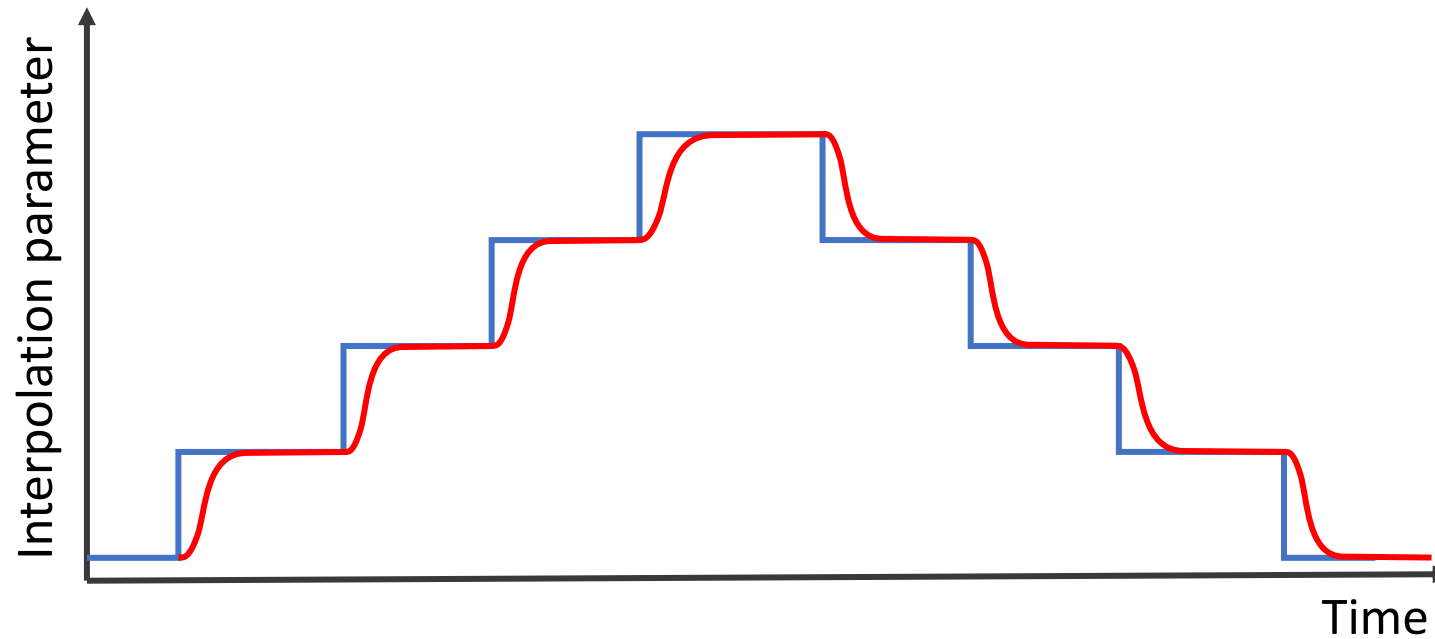$\Delta G_i / kJ \cdot mol^{-1}$

FIG. 1. Free energy of hydration change in a thermodynamic integration from neon to a sodium cation in 80 ps. The dashed curve gives the free energy change of the reverse process.

Interpolation parameter

Time

More accurate: slower change. Low hysteresis does not mean accurate: why?

T.P. Straatsma, H.J.C. Berendsen, *J. Chem. Phys.* **1988.**

Main sources: Correlation and sampling uncertainty
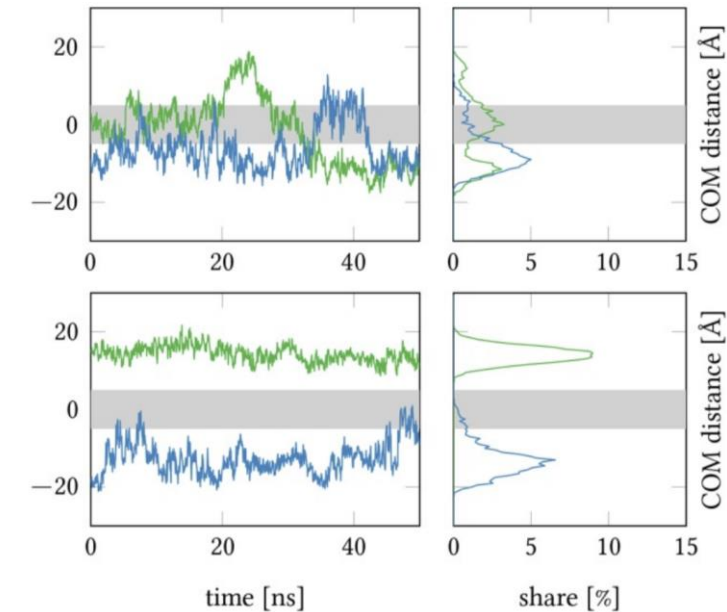
Method 1: **Bootstrapping** (What if we knew less?)
- Sub-sample observations, estimate uncertainty from variance
- Risks:
  - Rare events
  - Biased sampling

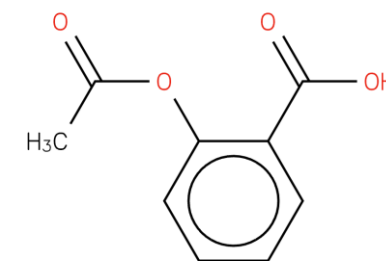Method 2: **Autocorrelation** (What if we measured independently?)
- Observable over time, quantify lagging
- Risks:
  - Not necessarily static: fast/slow process
  - Complex interplay with constraints

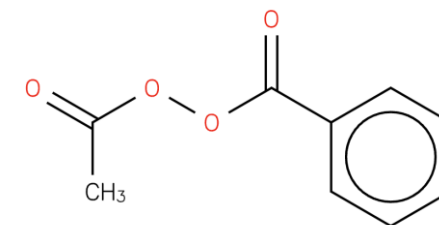Method 3: **Block averages** (What if time domains are representative?)
- Slice time in blocks and treat them as observables
- Risks:
  - Hard to figure the duration out

1  Sketch the single topology and double topology approach for one specific alchemical path change from aspirin to acetozone. Try to find the shortest path in terms of sites alchemically changed. Which method is likely to be preferable and why?

aspirin

2  Consider the surface deprotonation case: One question you might have is why not to vanish a surface proton instead of the complicated way of placing the proton in the liquid.

   a) List reasons why vanishing a proton is not the right thing to do.

   b) Using the dimensions from the slides and approximating the unit cell to be rectangular: what change in pH is expected if you add/remove one proton? (Pick one slab only)

acetozone