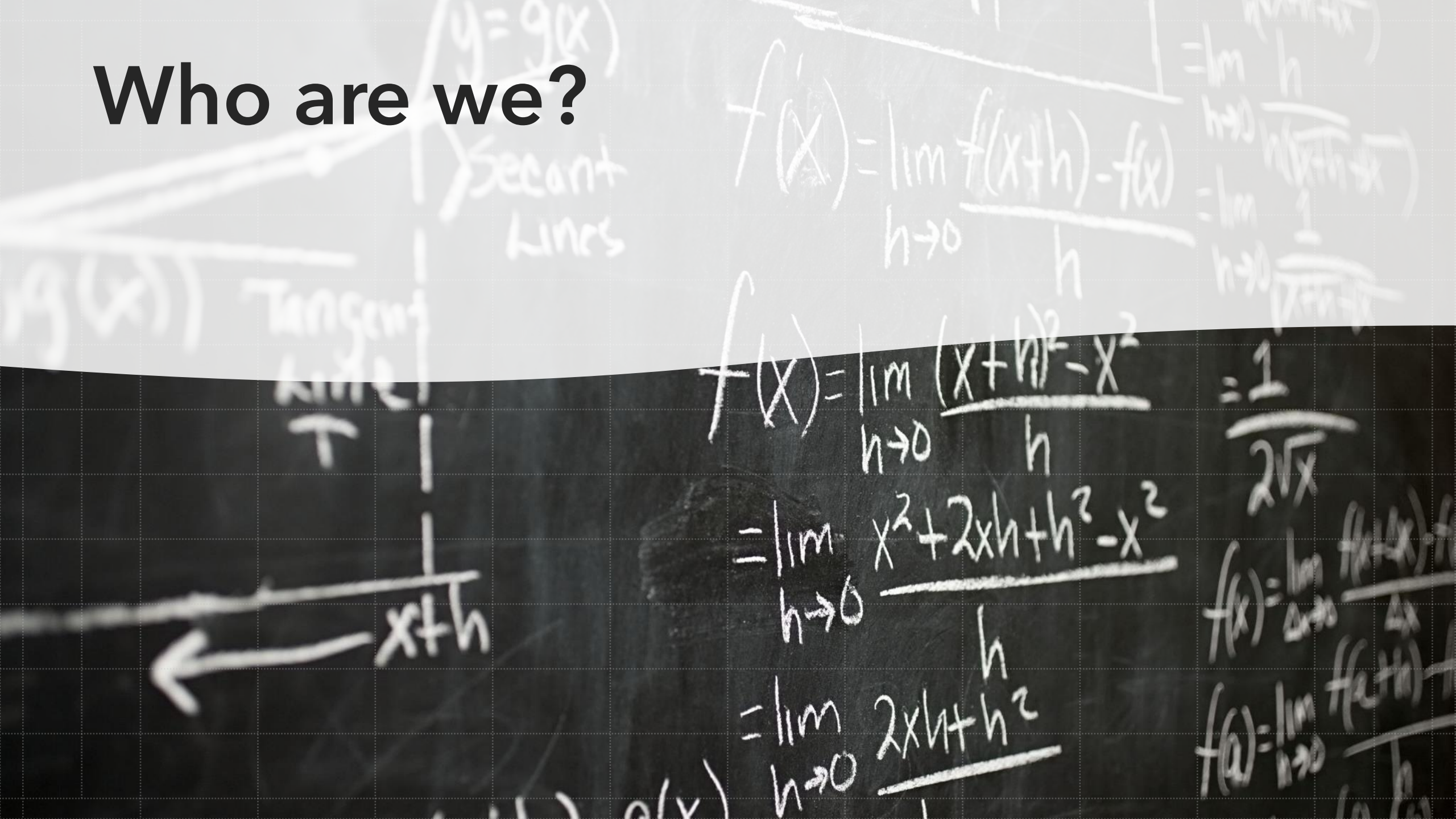


# Machine learning for theoretical chemistry

**Mario Barbatti**

*Aix Marseille Université, CNRS, Institut de Chimie Radicalaire  
Institut Universitaire de France*

# Who are we?



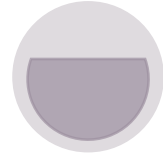
# The Light & Molecules group

# The Light & Molecules Group



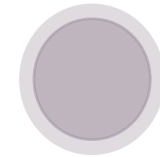
## Methods

Nonadiabatic dynamics  
Nuclear ensembles



## Software

NEWTON-X  
platform



## Applications

Photoprocesses in

- Fundamental PhysChem
- Molecular biology
- Organic devices
- Environment

# The Light & Molecules Group

## Current members

Mario Barbatti (PI)

Baptiste Demoulin (IT researcher)

Josene M Toldo (postdoc)

Saikat Mukherjee (postdoc)

Bidhan Garain (postdoc) \*

Rafael Mattos (PhD candidate)

Matheus Bispo (PhD candidate) \*

## Recent past members

Mariana T do Casal

Ritam Mansour

Shuming Bai

Lijljana Stojanovic

Carlos E de Moura

Fabris Kossoski

Prateek Goel

Max Pinheiro Jr \*

Moumita Kar

\* Dedicated to ML projects

**Marseille-Xiamen consortium**



---

## Marseille (ICR)

Mario Barbatti



Nuclear ensembles  
Nonadiabatic dynamics  
Unsupervised ML

---

## Xiamen Univ

Pavlo Dral



Atomistic supervised ML

Pinheiro Jr et al. *Sci Data* **2023**, 10, 95

Zhang et al., *In Quantum Chemistry in the Age of ML*, **2023**

Pinheiro Jr and Dral., *In Quantum Chemistry in the Age of ML*, **2023**

Barbatti et al. *JCTC* **2022**, 18, 6851

Dral; Barbatti. *Nat Rev Chem* **2021**, 5, 388

Pinheiro Jr et al. *Chem Sci* **2021**, 12, 14396

Dral et al. *Top Curr Chem* **2021**, 379, 27

Xue; Barbatti; Dral. *J Phys Chem A* **2020**, 124, 7199

Dral; Barbatti; Thiel. *J Phys Chem Lett* **2018**, 9, 5660

---

**Marseille (ICR)**

Mario Barbatti



Nuclear ensembles  
Nonadiabatic dynamics  
Unsupervised ML

---

**Xiamen Univ**

Pavlo Dral



Atomistic supervised ML

---

**Marseille (LIS)**

Thierry Artières

Hachem Kadri



Data science expertise



# The Newton-X platform



# Newtonian Dynamics Close to the X-Seam

LIGHT AND  
MOLECULES

- Surface hopping & Nuclear ensemble simulations
- Freeware
- Open source



Barbatti *et al.* *JCTC* **2022**, 18, 6851

Baptiste Demoulin



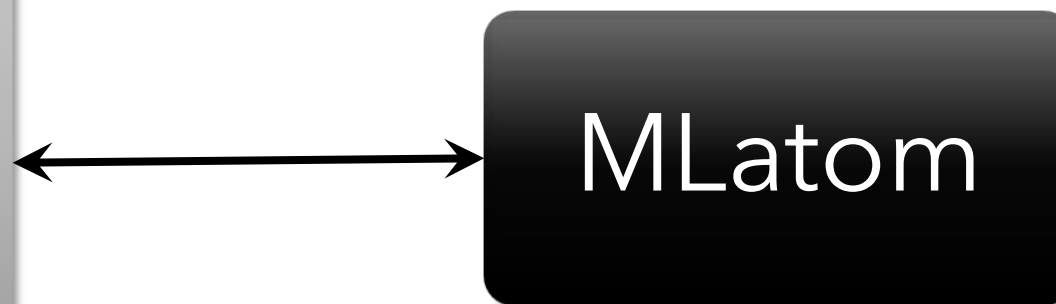
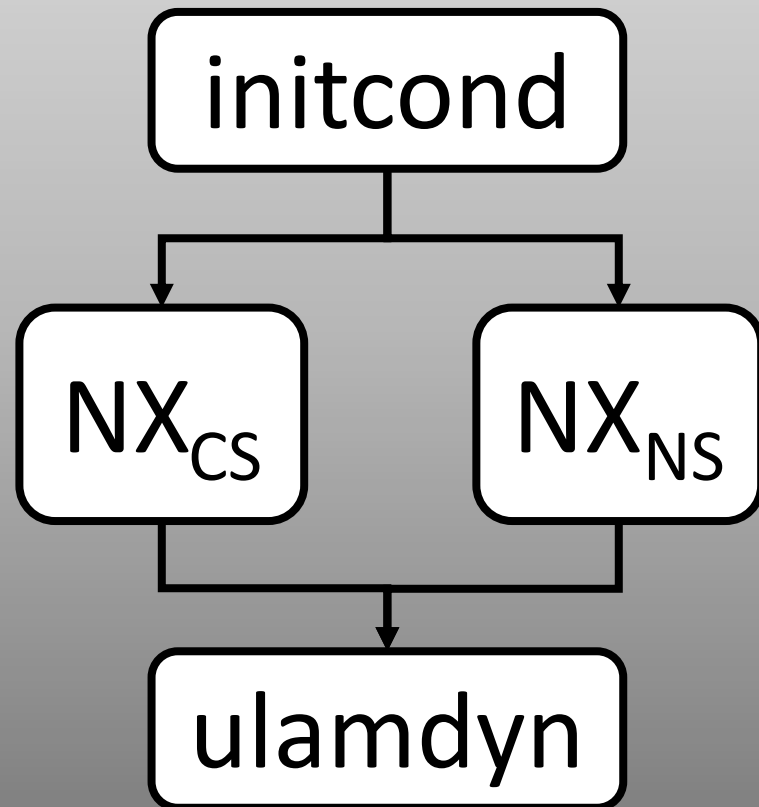
# Newtonian Dynamics

## Close to the X-Seam

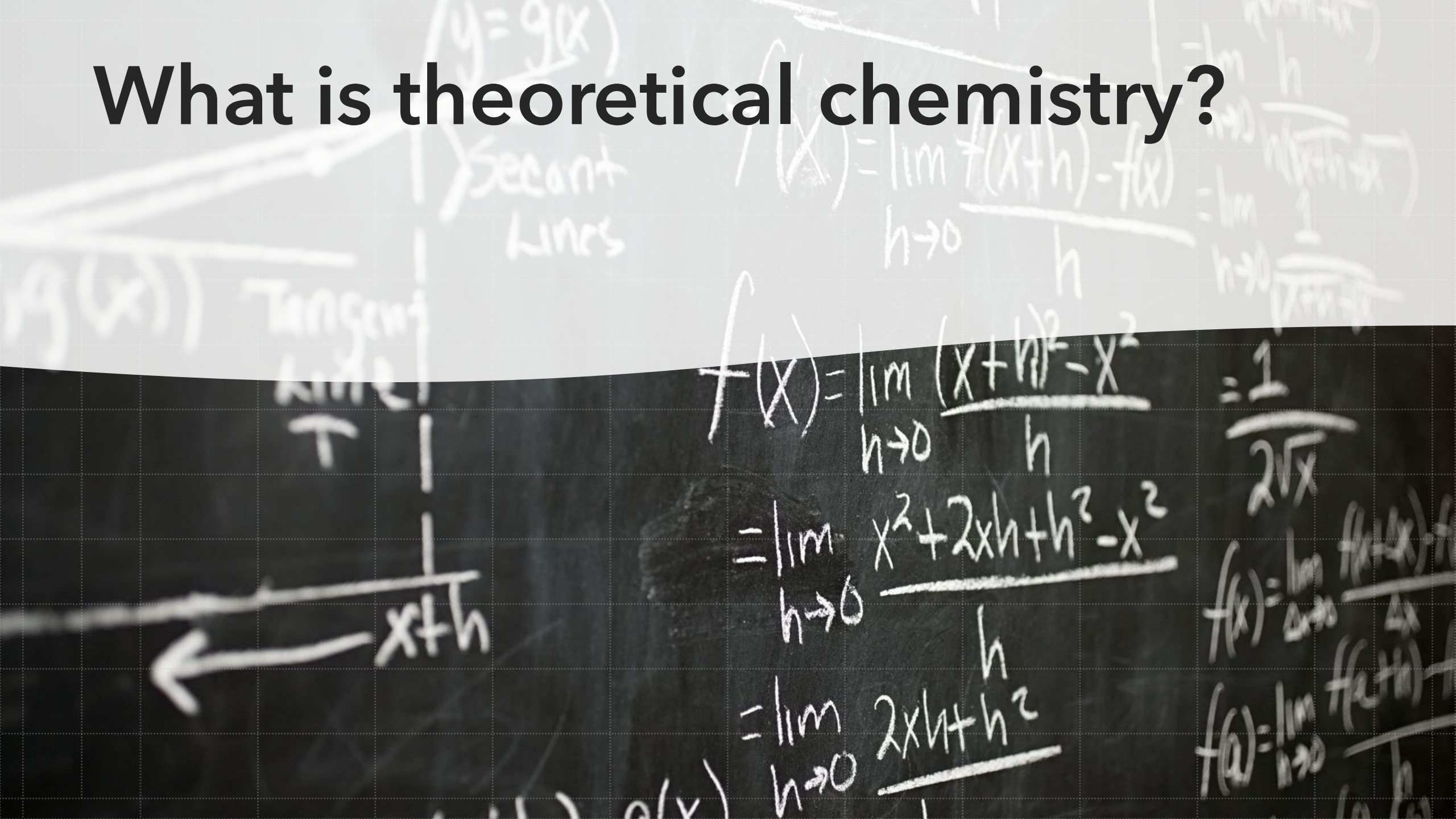
LIGHT AND  
MOLECULES

- Simulations with MRCI, MCSCF, CASPT2, ADC(2), TDDFT, TD-DFTB, Semiempirical/CI, Analytical models, ML potentials
- Interfaces to Columbus, Turbomole, Gaussian, Bagel, Gamess, CP2K, DFTB+, Mopac (Pisa), ORCA, Open Molcas, MNDO, MLatom

# Newton-X Platform



# What is theoretical chemistry?



# Quantum mechanics of molecules

Schrödinger equation for the molecule (including electrons and nuclei)

$$i\hbar \frac{\partial \Psi}{\partial t} = \hat{H} \Psi$$

Following Born and Oppenheimer's approach, this problem simplifies to

### **Electrons ( $\mathbf{r}$ )**

Electronic Schrödinger equation  
(adiabatic approximation)

$$(T_{elec}(\mathbf{r}) + V(\mathbf{r}, \mathbf{R}))\varphi(\mathbf{r}; \mathbf{R}) = E(\mathbf{R})\varphi(\mathbf{r}; \mathbf{R})$$

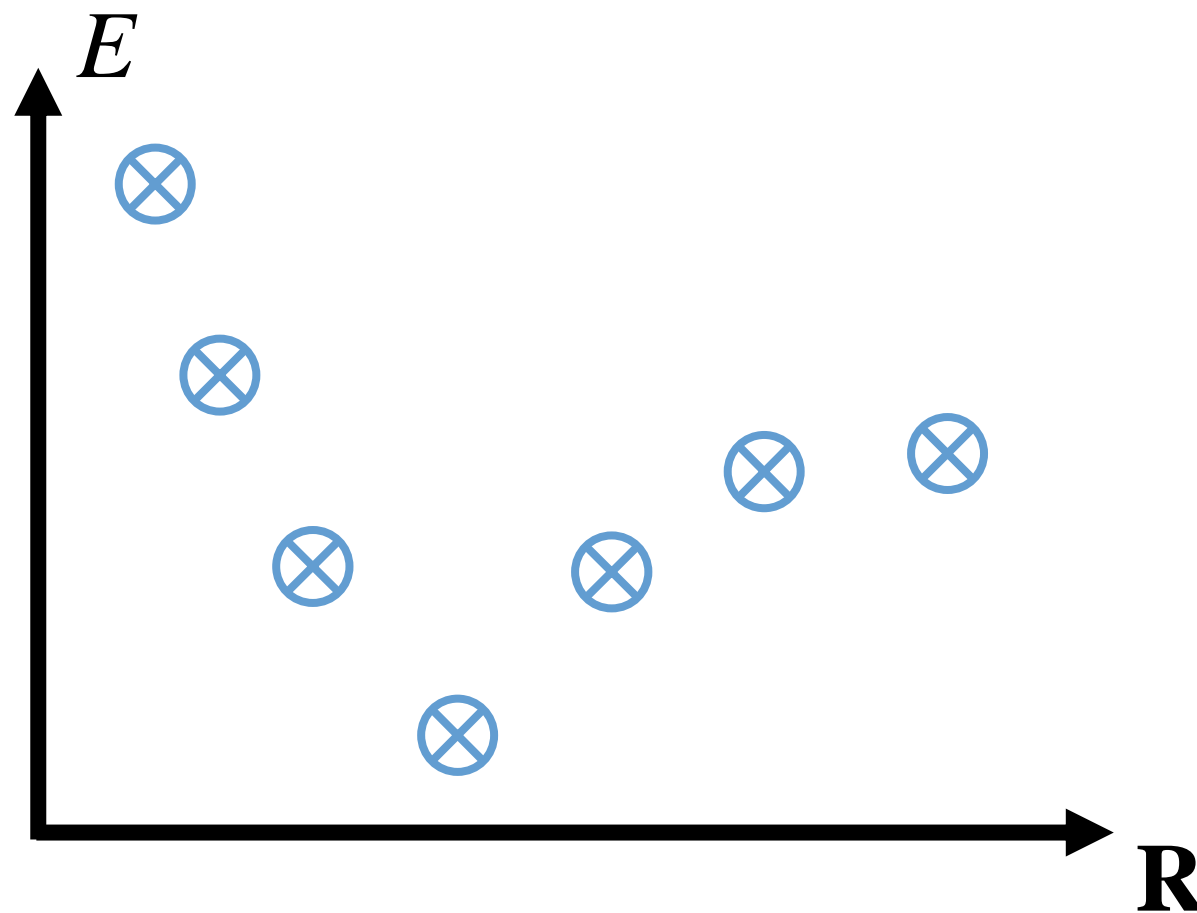
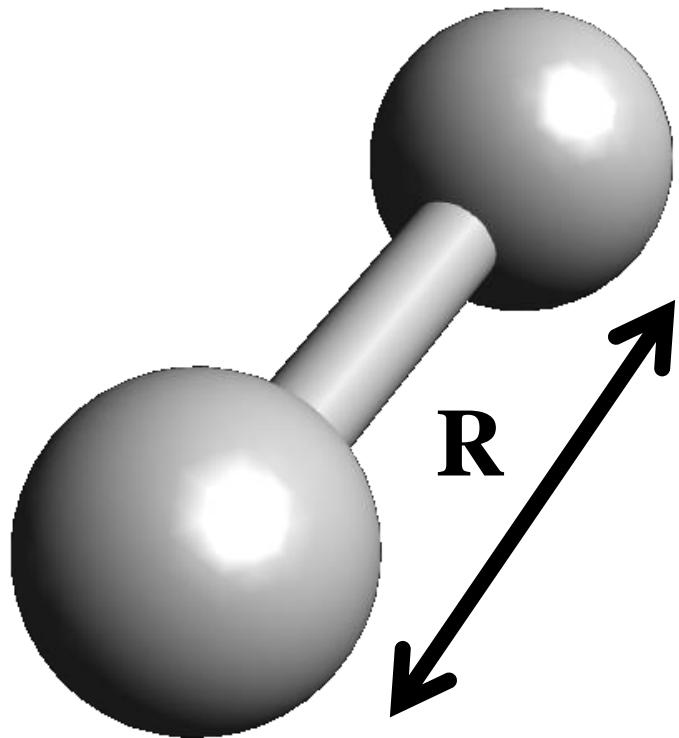
### **Nuclei ( $\mathbf{R}$ )**

Nuclear Newton's equation  
(Classical approximation)

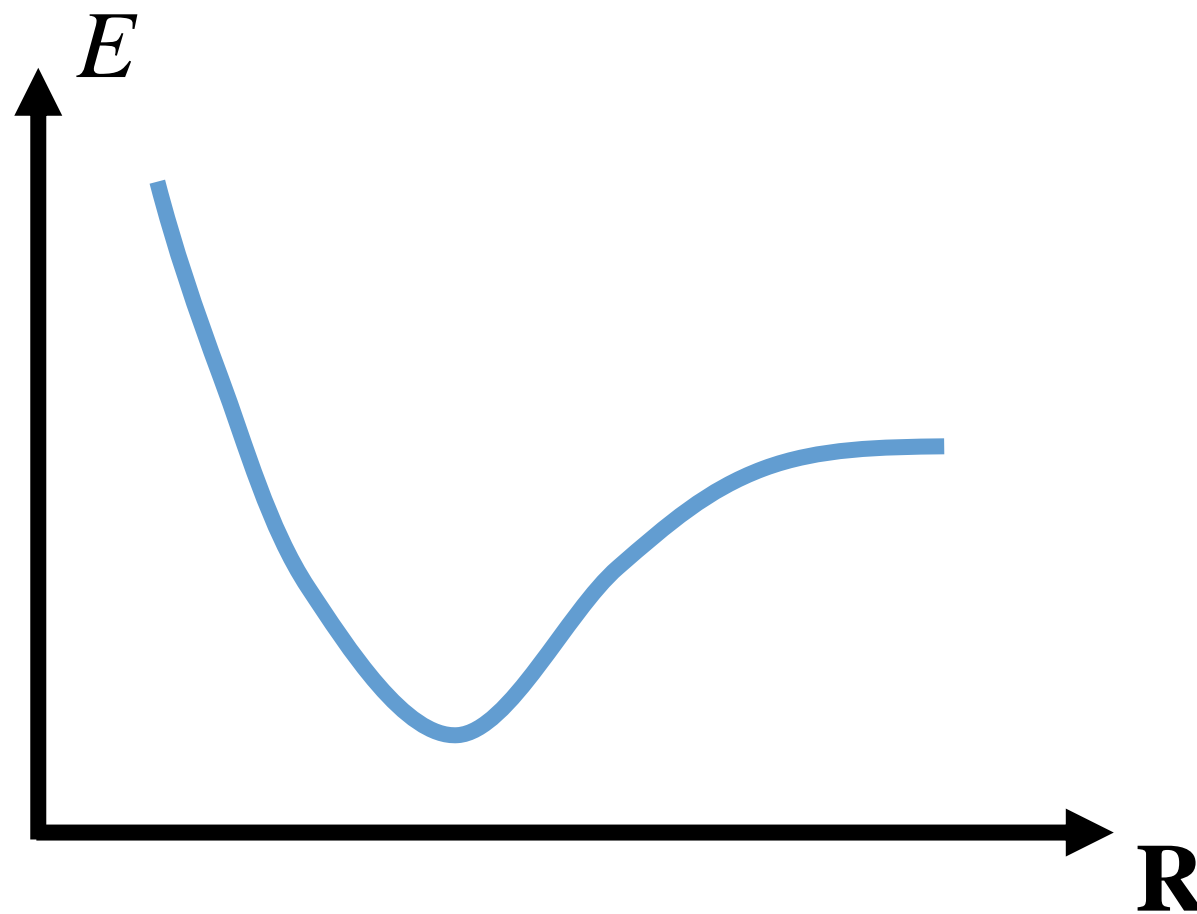
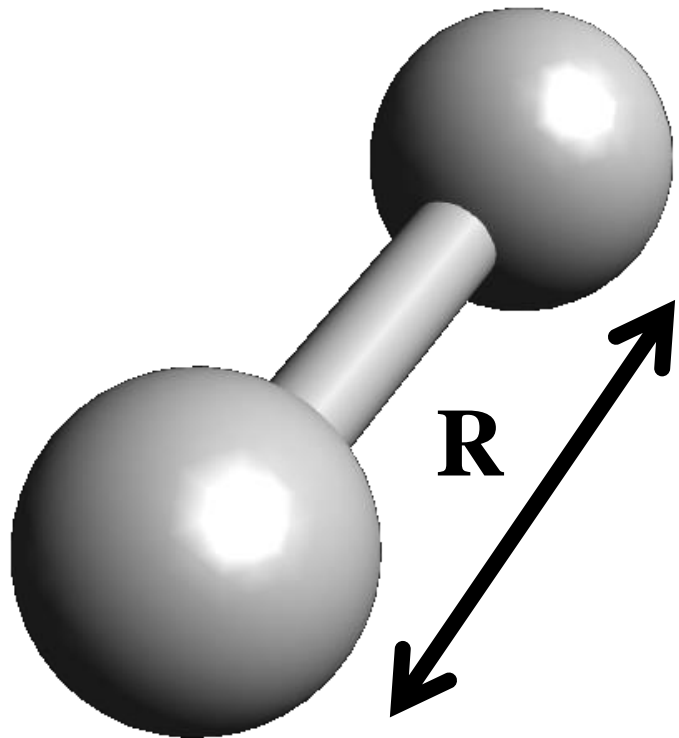
$$M_{\alpha} \frac{d^2 \mathbf{R}_{\alpha}}{dt^2} = -\nabla_{\alpha} E(\mathbf{R})$$

The core quantity is the potential energy surface  $E(\mathbf{R})$

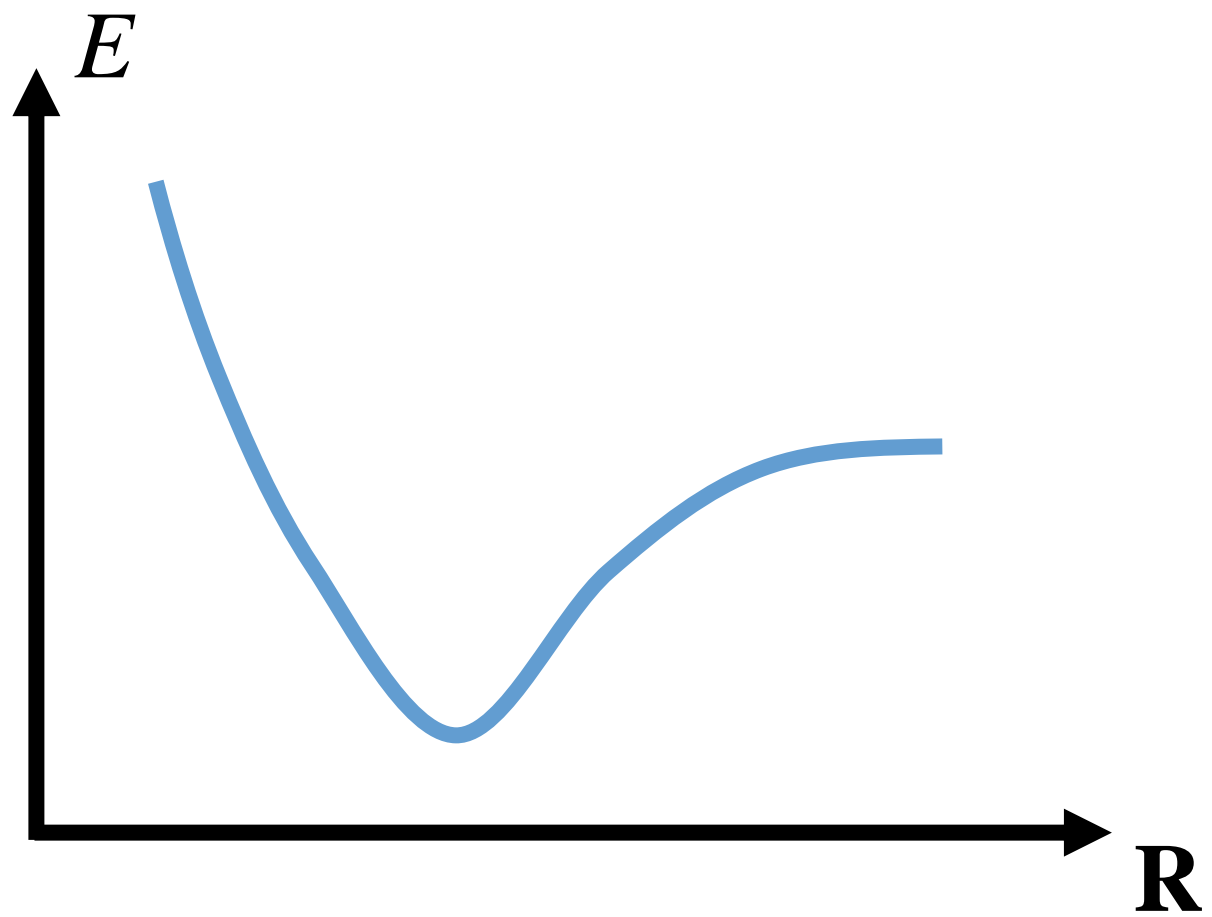


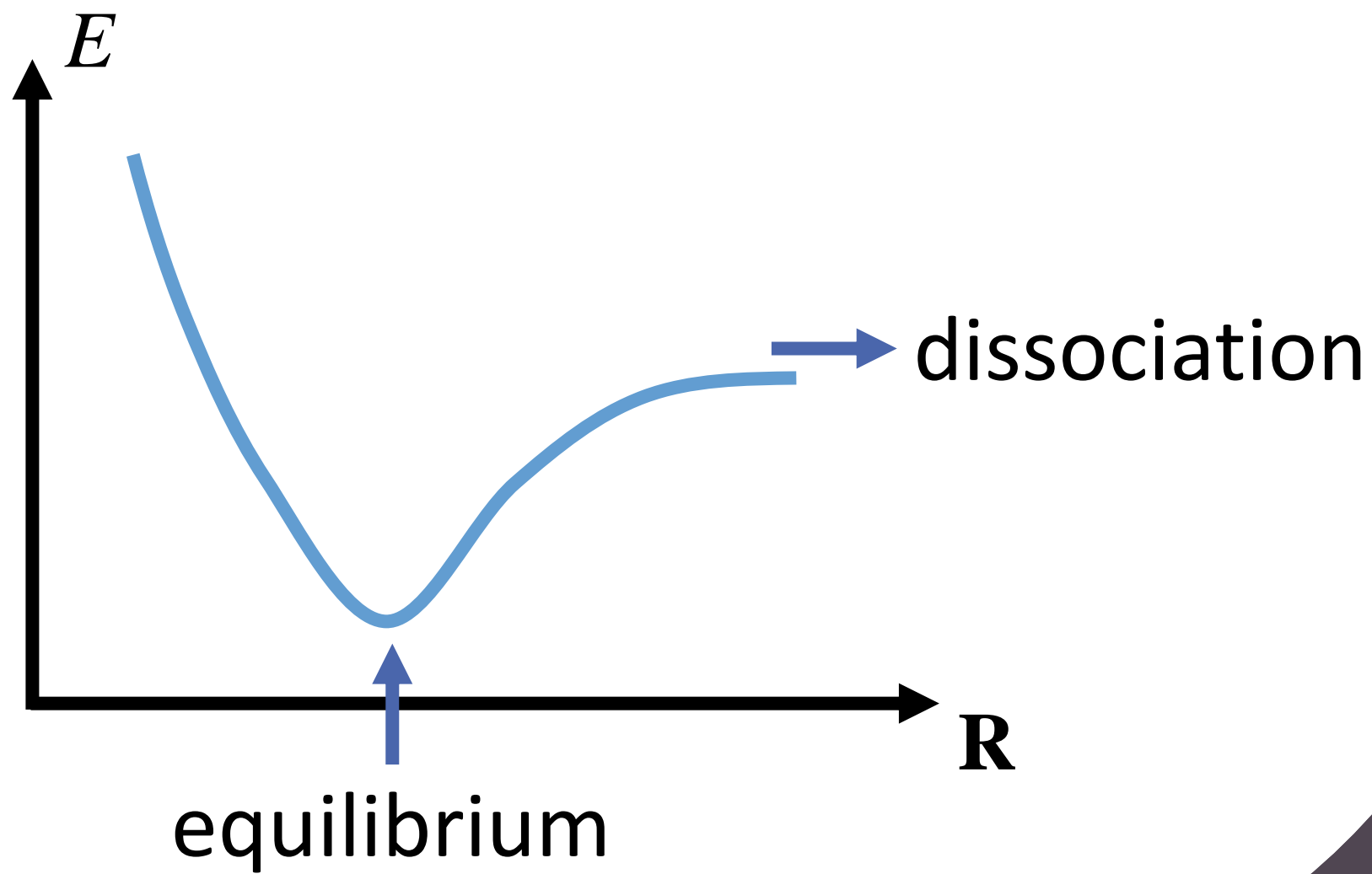


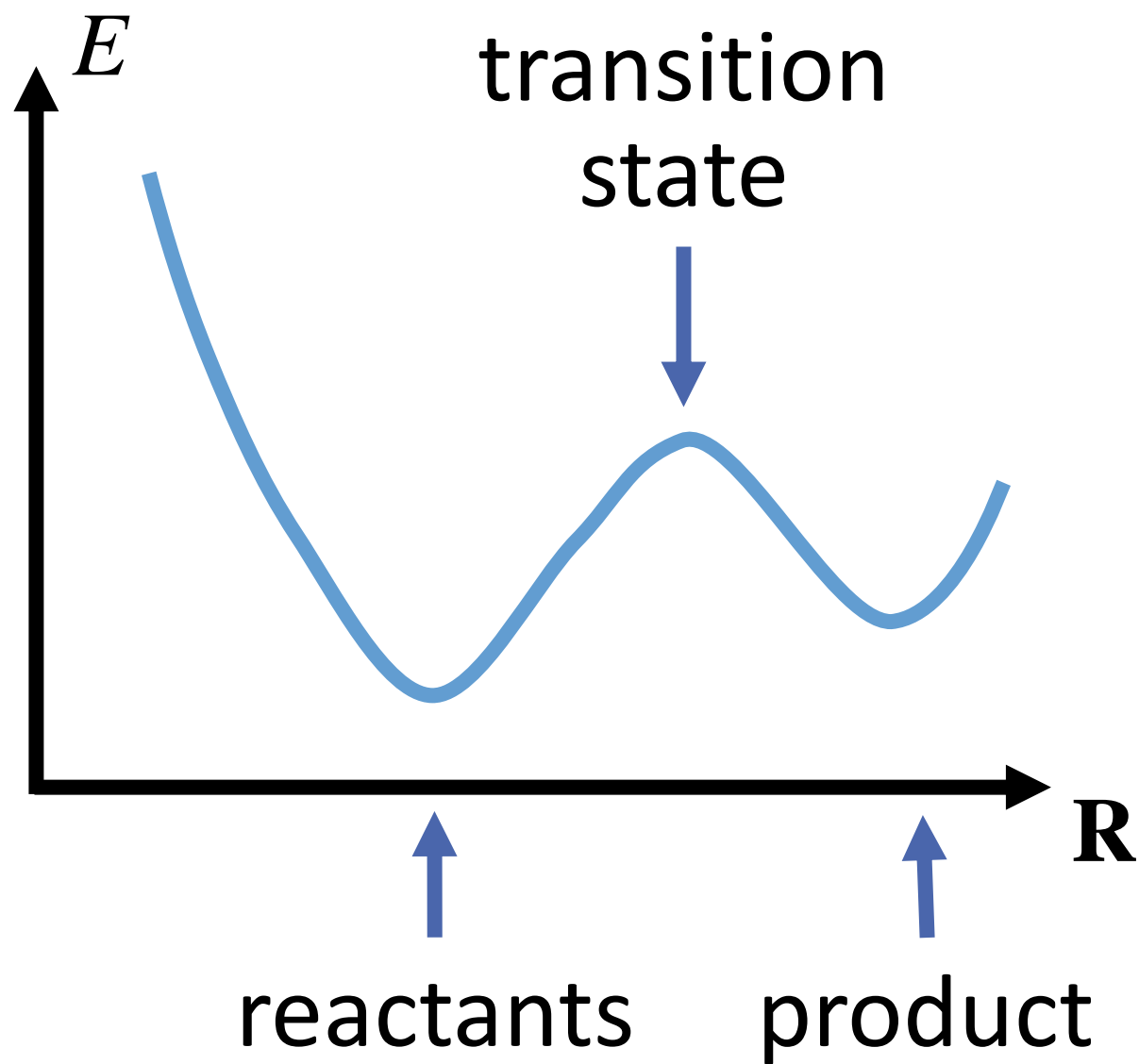
$$H_{elec} \varphi(\mathbf{r}_1, \mathbf{r}_2; \mathbf{R}) = E(\mathbf{R}) \varphi(\mathbf{r}_1, \mathbf{r}_2; \mathbf{R})$$

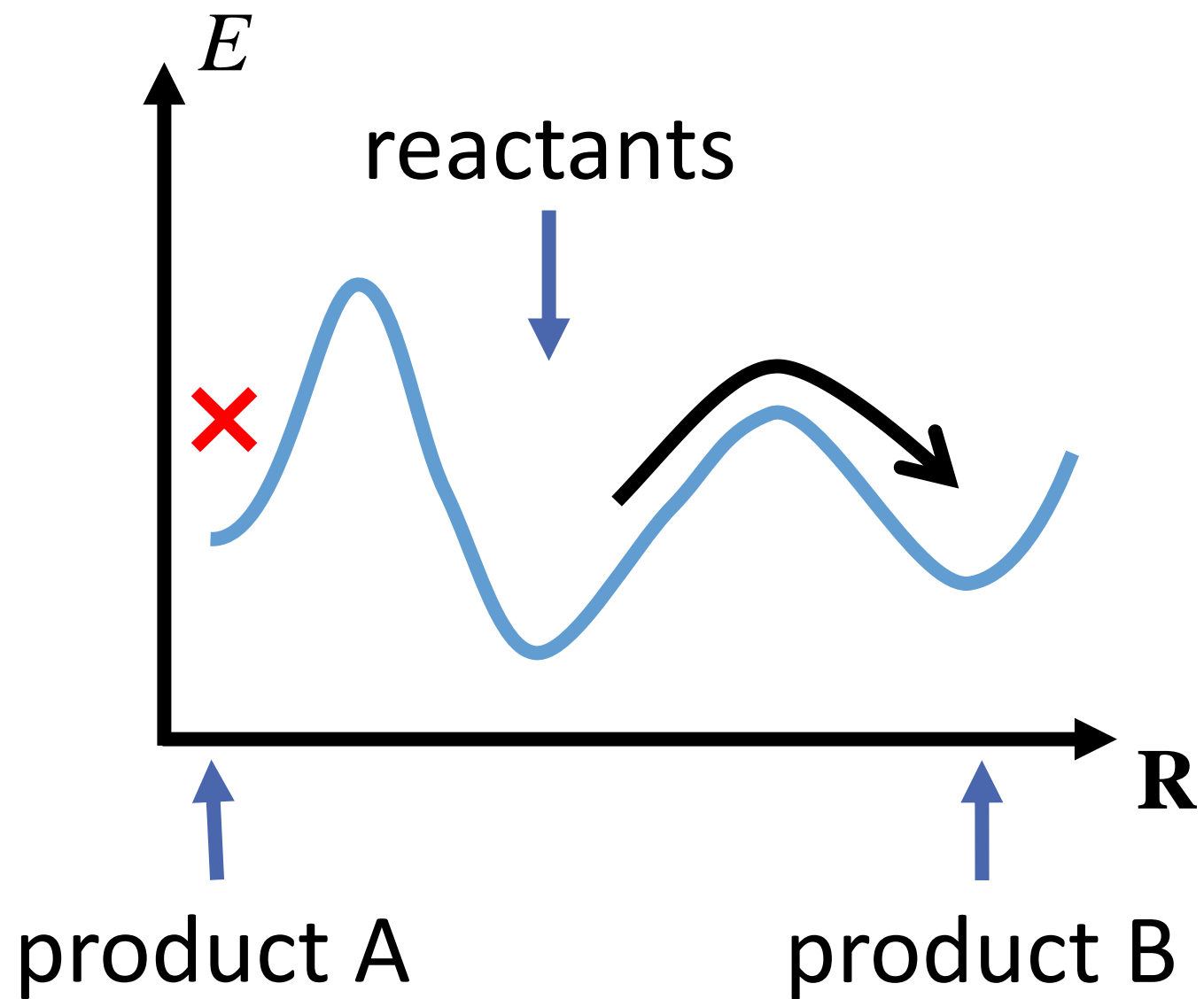


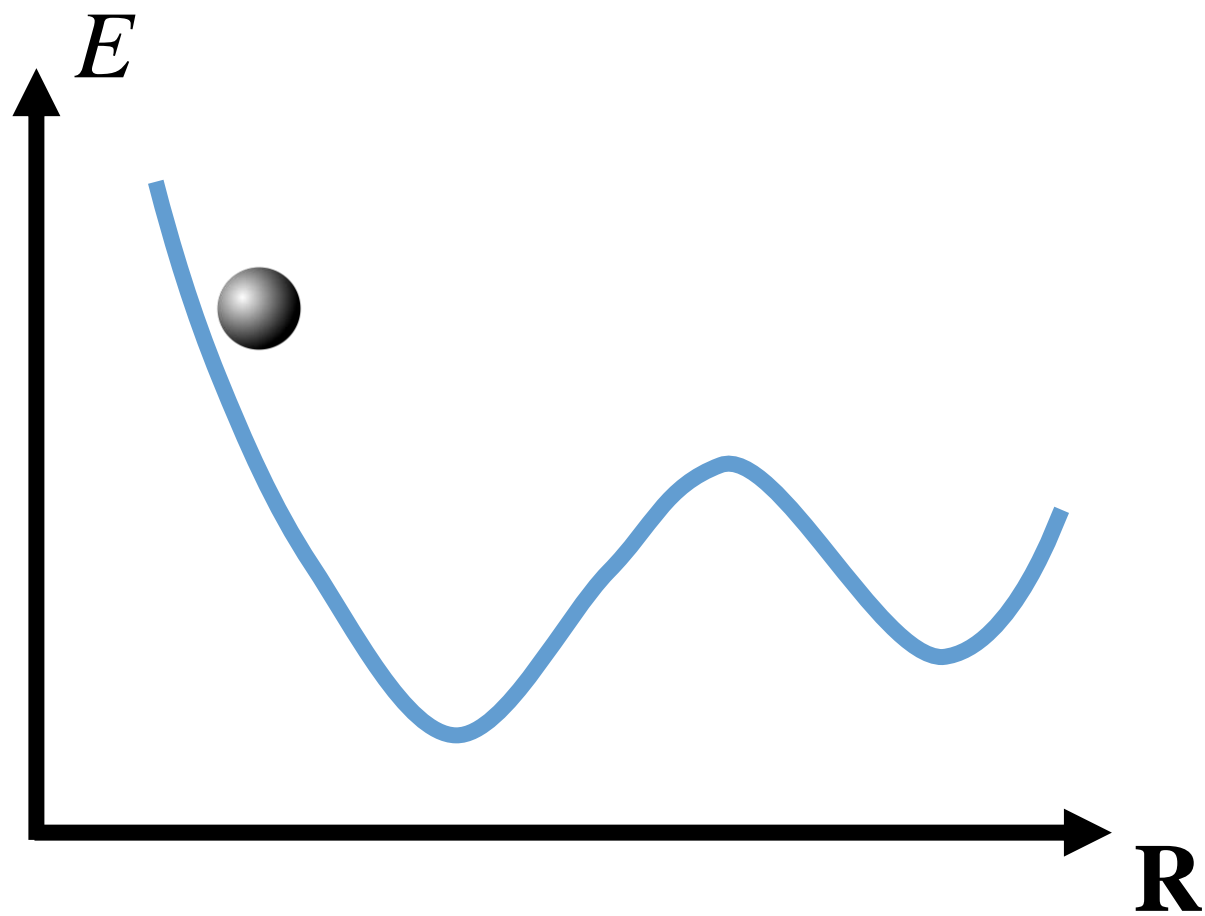
$$H_{elec} \varphi(\mathbf{r}_1, \mathbf{r}_2; \mathbf{R}) = E(\mathbf{R}) \varphi(\mathbf{r}_1, \mathbf{r}_2; \mathbf{R})$$







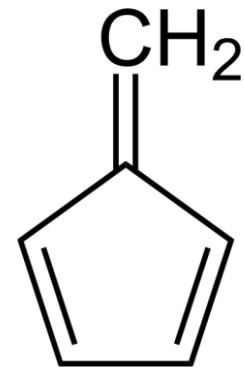






Potential energy surfaces  $E(\mathbf{R})$  have  $3N_{at}-6$  dimensions

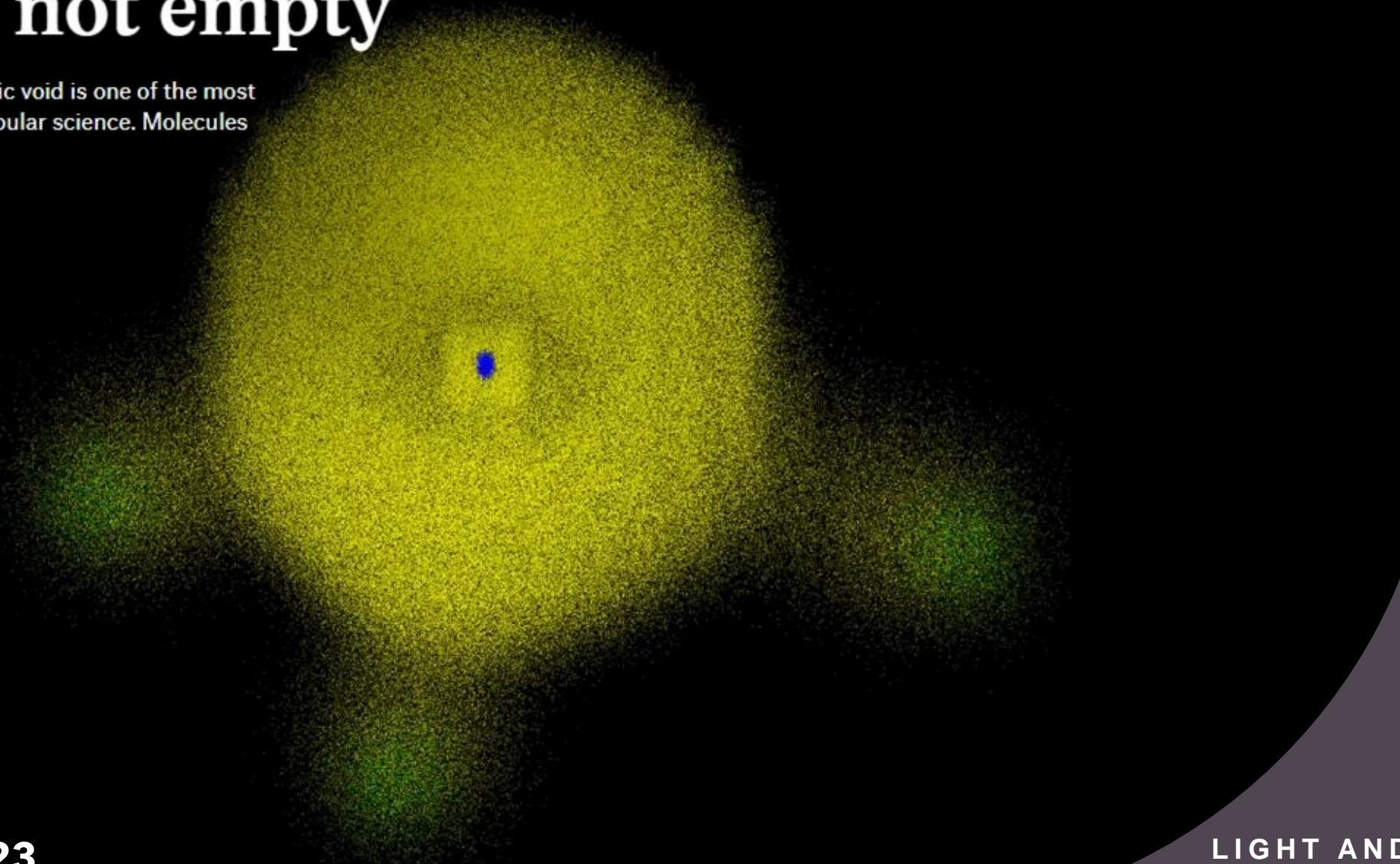
For fulvene,  $N_{at} = 12$ ,  $E(\mathbf{R})$  has 30 dimensions



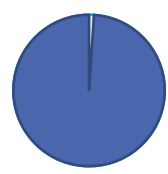
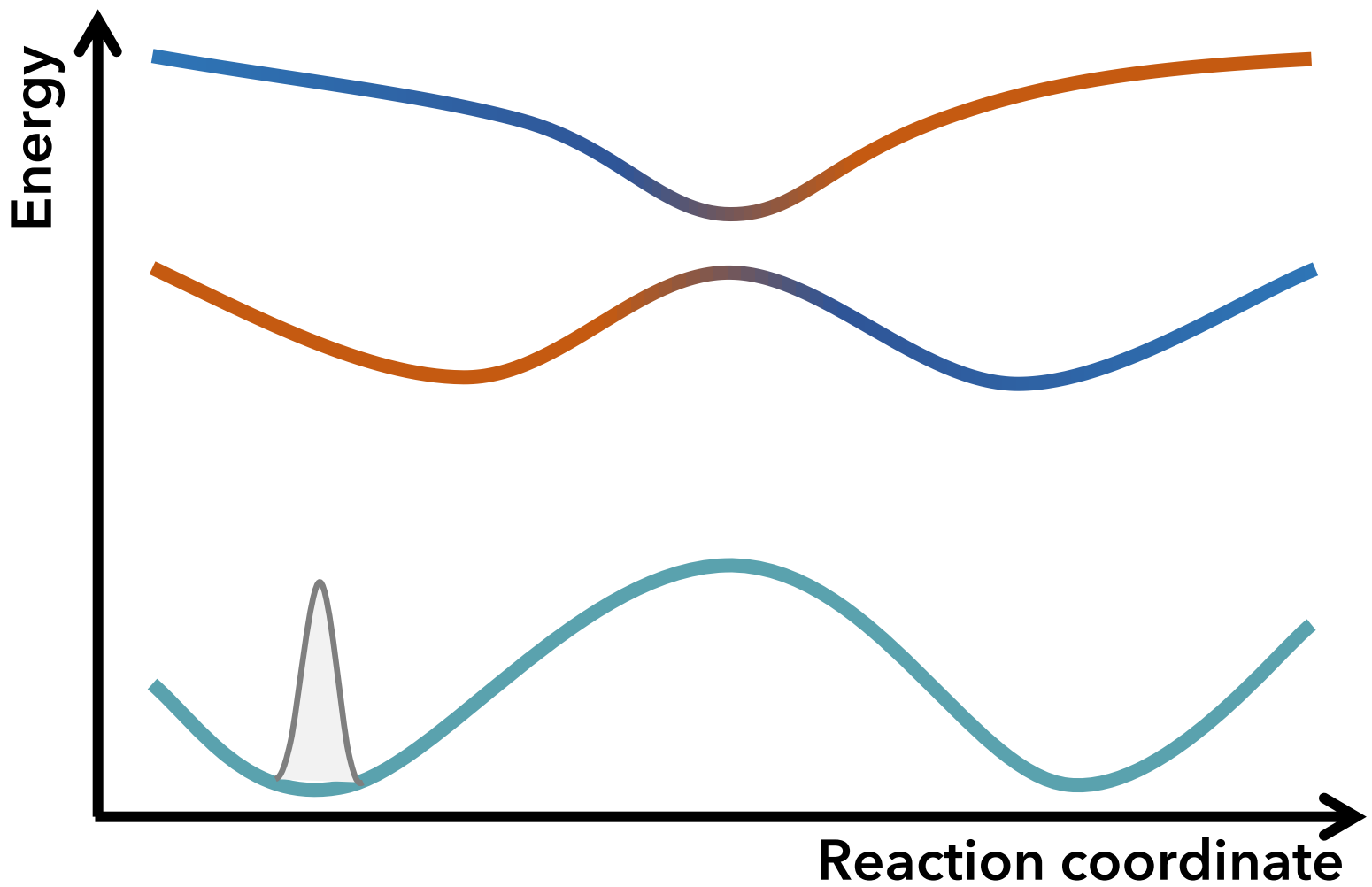
fulvene

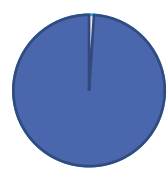
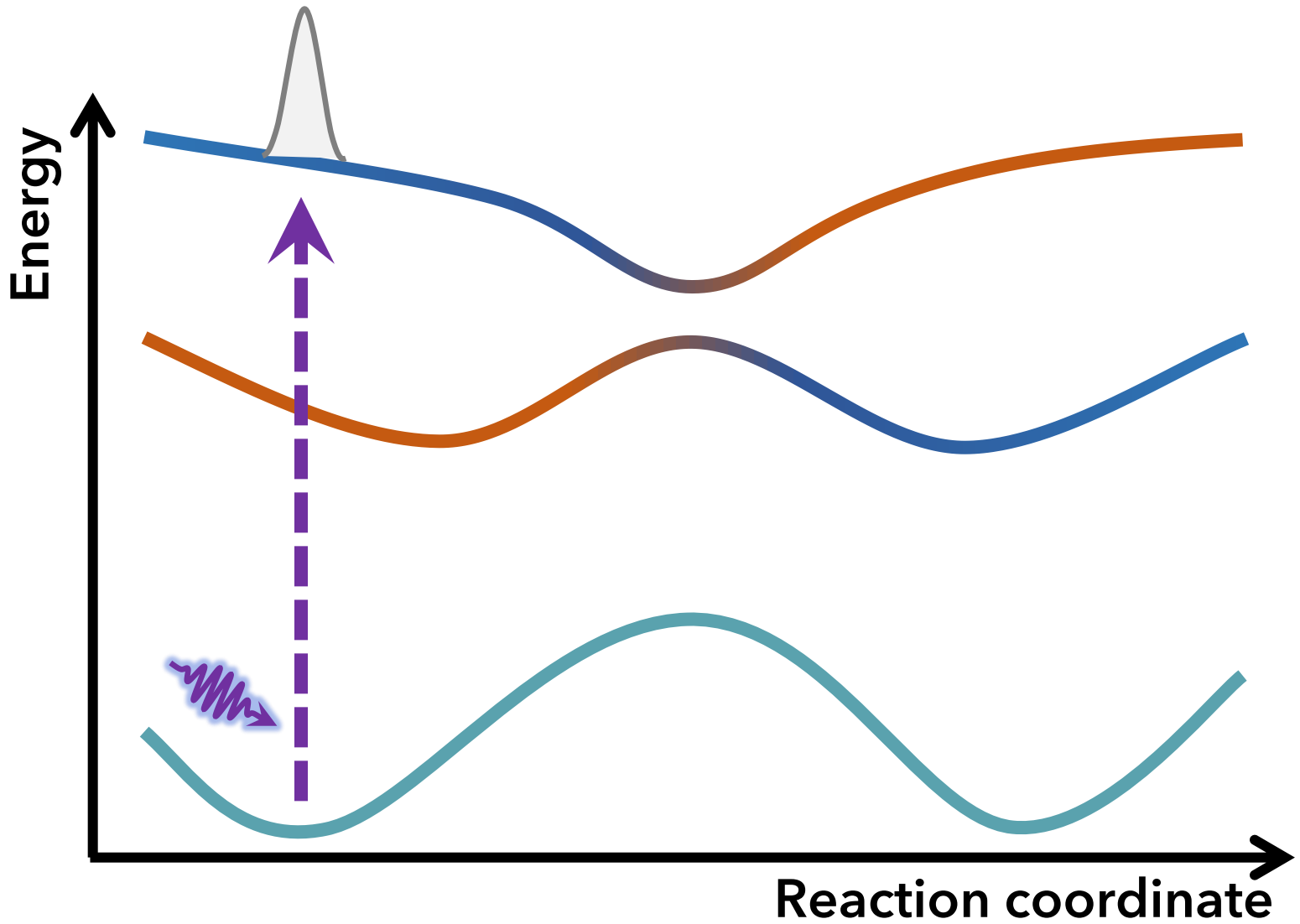
# We are not empty

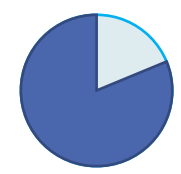
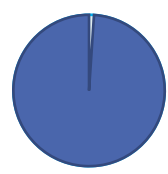
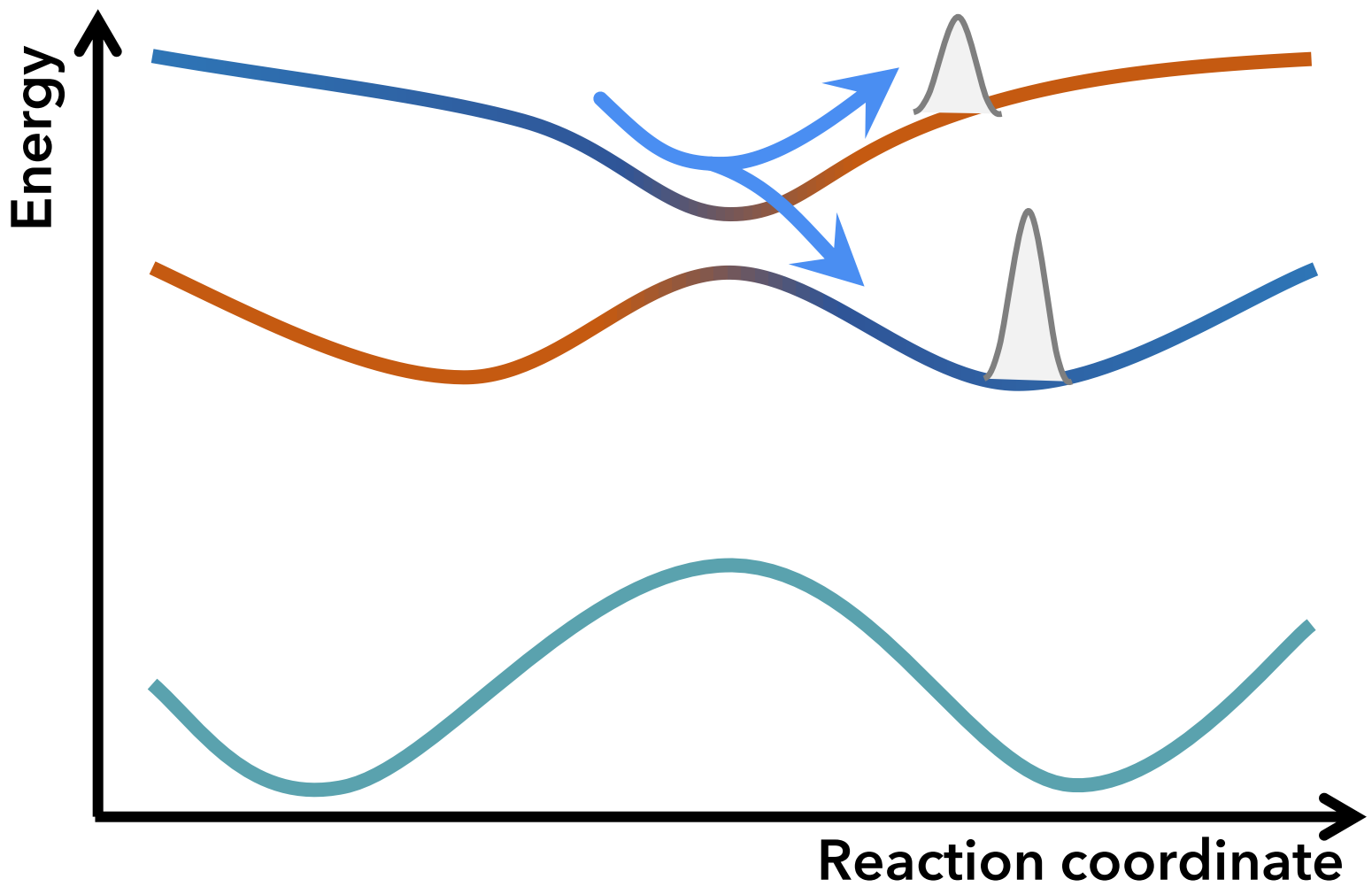
The concept of the atomic void is one of the most repeated mistakes in popular science. Molecules are packed with stuff

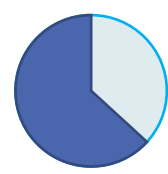
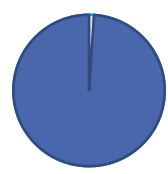
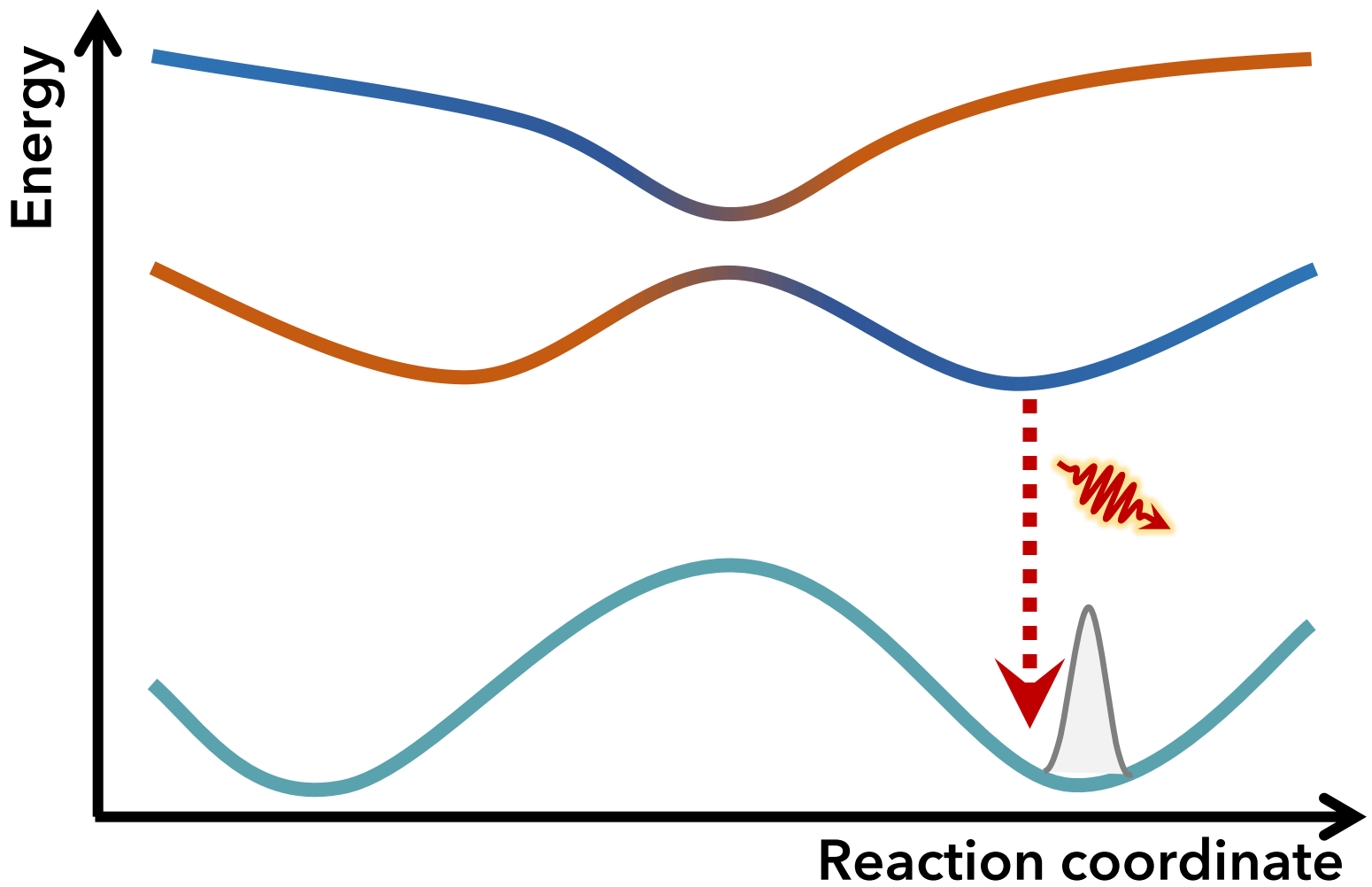


# Nonadiabatic molecular dynamics (NAMD)







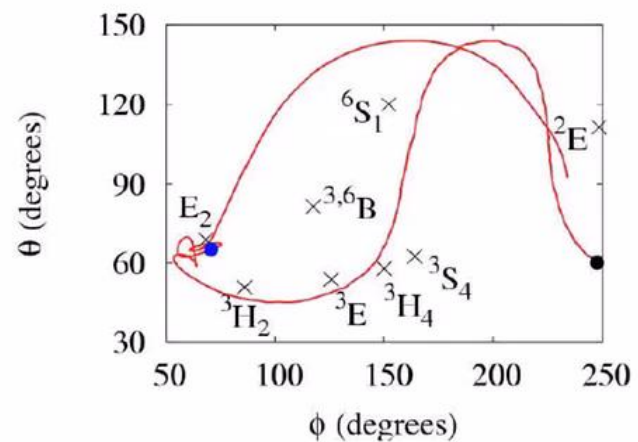
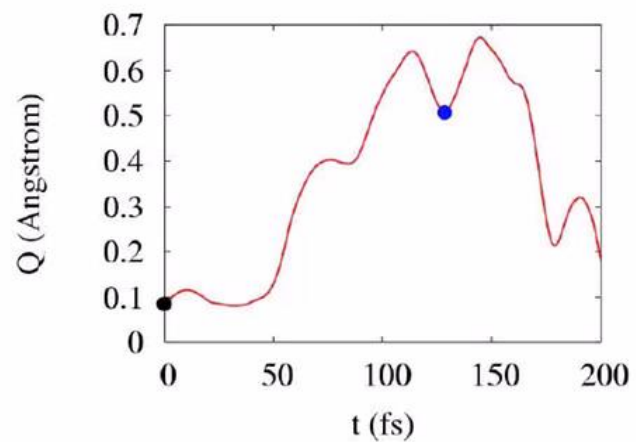
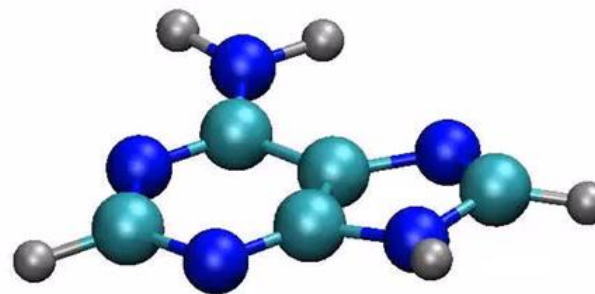
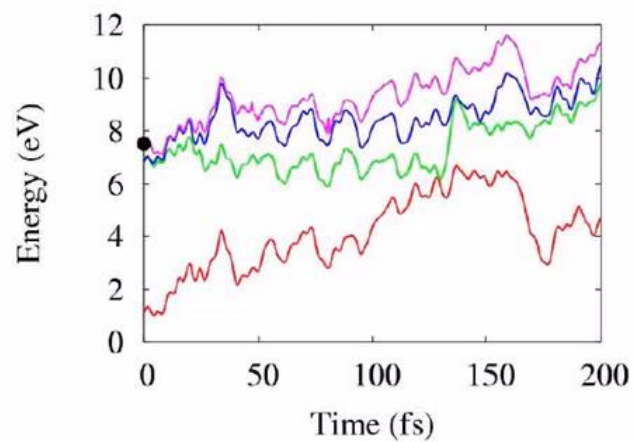




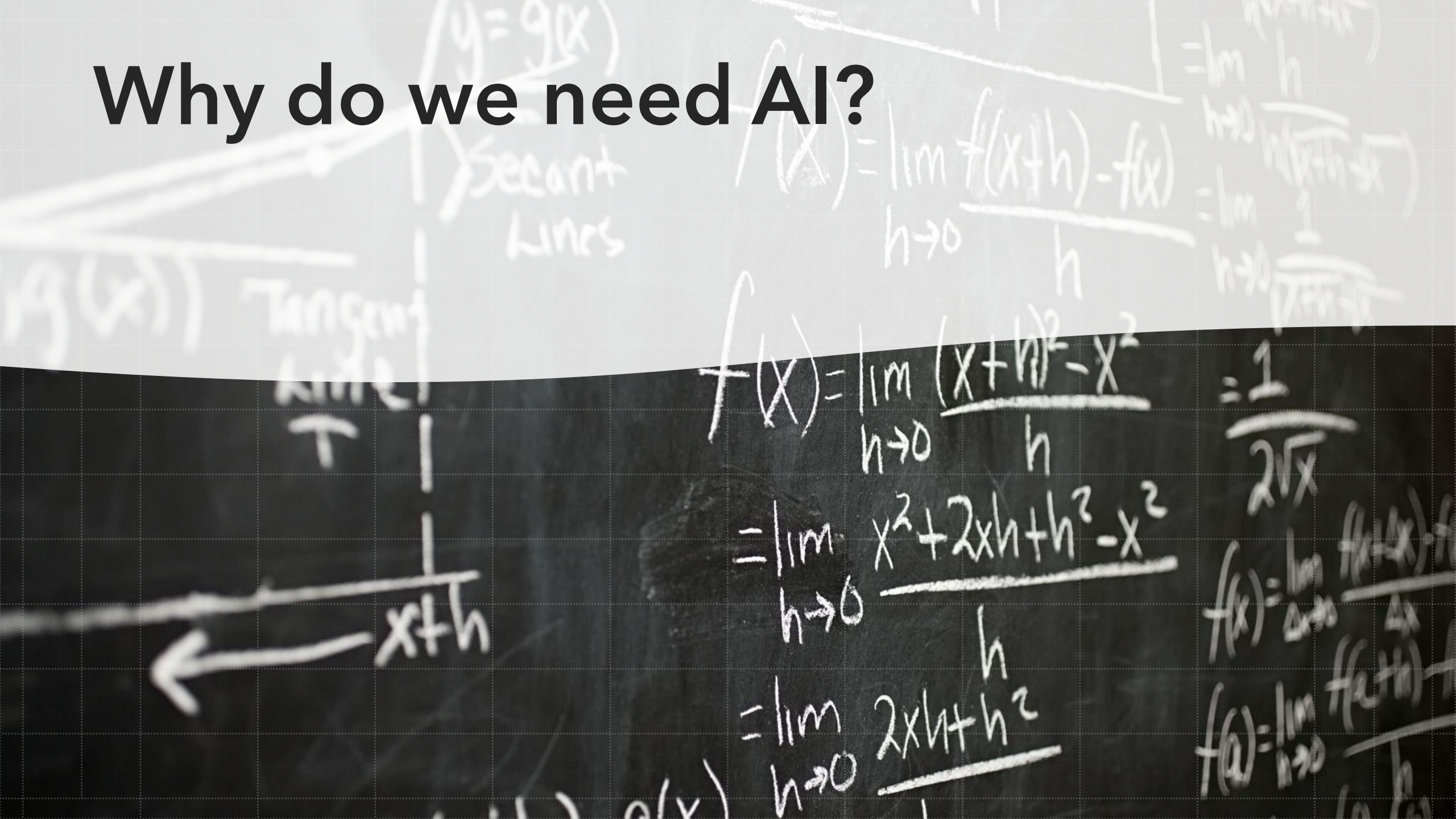
How can we simulate  
nonadiabatic molecular dynamics?

# Mixed quantum-classical methods

1. Nuclei are treated via *classical trajectories*
2. Electrons are treated *quantum mechanically*
3. A nonadiabatic algorithm introduces *post Born-Oppenheimer effects*



# Why do we need AI?



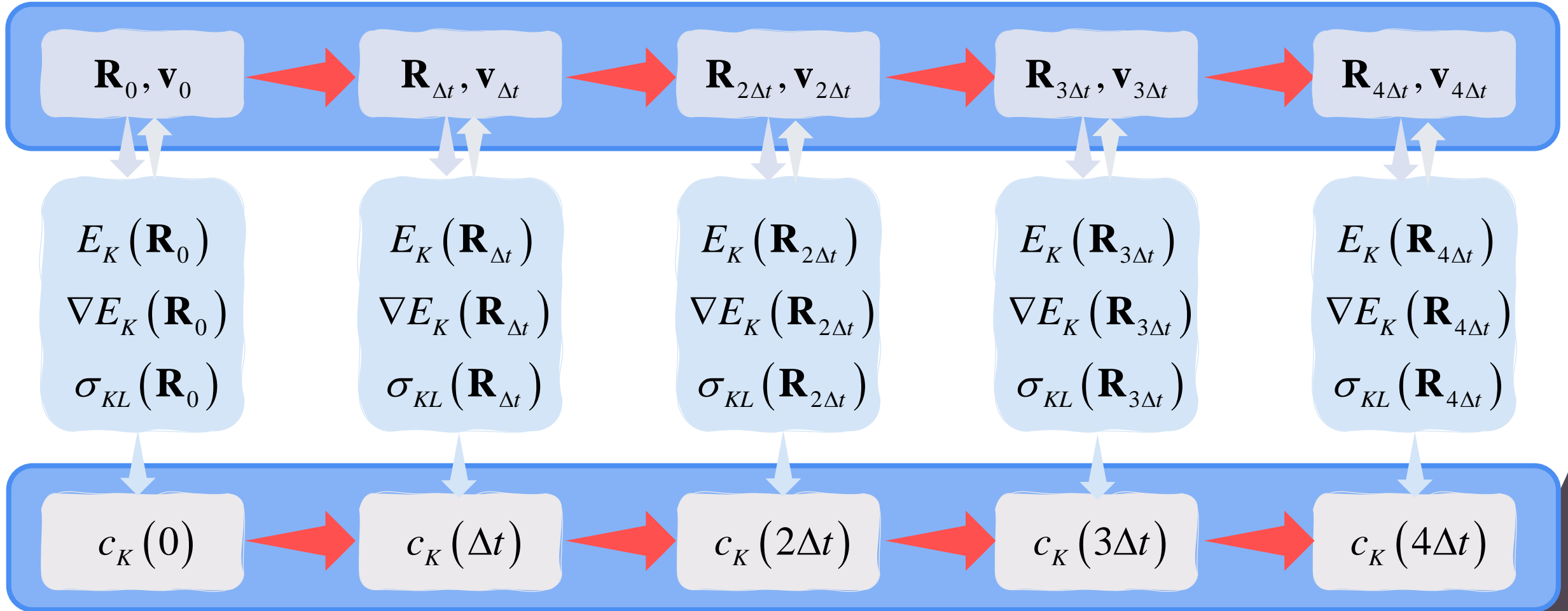
AI for theoretical chemistry has been used to

- Search the chemical space of compounds
- Perform dimensionality reduction, clustering, and pattern recognition
- Improve or accelerate quantum chemical methods
- Predict properties as a surrogate approach

# Costs of dynamics

# Nonadiabatic dynamics

## Classical EOM



## Quantum EOM

# Dynamics may be expensive

$$T_{total} \approx N_{\text{Trajectories}} \times N_{\text{Single Points}} \times T_{\text{Single Point}}$$

How much does dynamics cost? [tinyurl.com/dyncost](https://tinyurl.com/dyncost)  
How many trajectories should we run? [tinyurl.com/trajs](https://tinyurl.com/trajs)



# Dynamics may be expensive

$$T_{total} \approx N_{\text{Trajectories}} \times \left( \frac{\tau_{\text{chem process}}}{\Delta \tau} \right) \times T_{\text{Single Point}}$$

$N_{\text{Trajectories}}$	= 100 trajectories
$T_{\text{Single Point}}$	= 6 min = 0.1 CPUh
$\tau_{\text{chem process}}$	= 500,000 fs = 0.5 ns
$\Delta \tau$	= 0.5 fs

$$T_{total} \approx 10 \text{ MCPUh}$$


Price 1 CPUh	= 0.02 € (France)
--------------	-------------------

Price 10 MCPUh	= 200 k€
----------------	----------

How much does dynamics cost? [tinyurl.com/dyncost](https://tinyurl.com/dyncost)

How many trajectories should we run? [tinyurl.com/trajs](https://tinyurl.com/trajs)






# Dynamics leaves a huge carbon footprint

1 CPUh @ 32 GB = 1.3 g CO<sub>2</sub>e 

10 MCPUh = 13 tCO<sub>2</sub>e



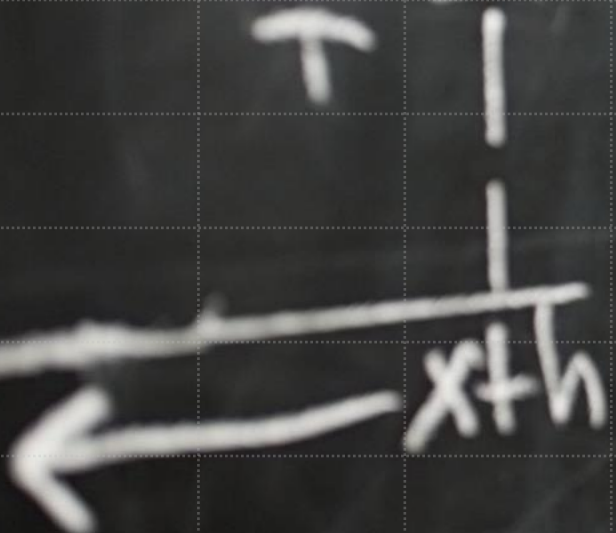
11.5 tCO<sub>2</sub>e/year

	→ × 2
	→ × 7
	→ × 10
	→ × 12
	→ × 14

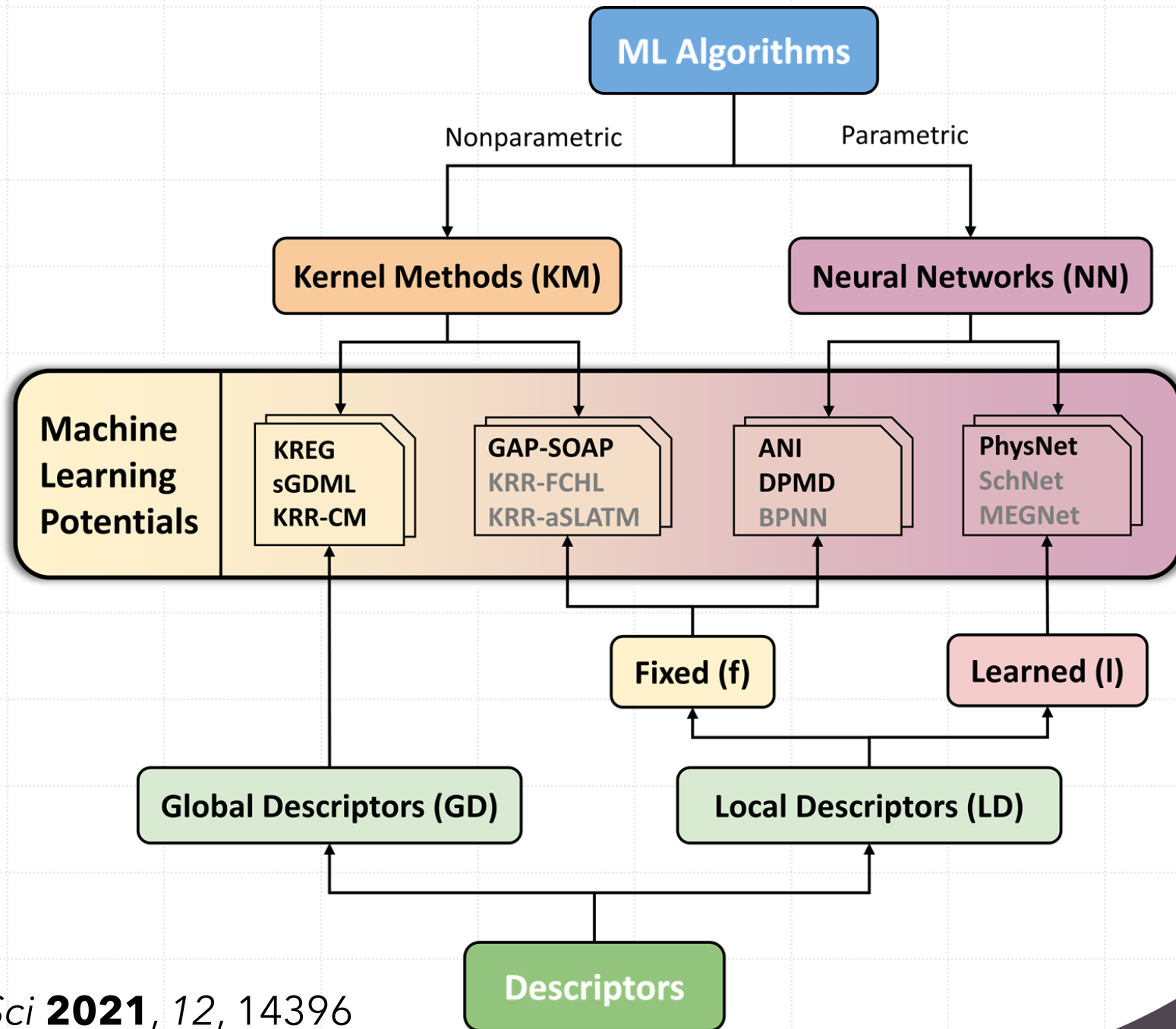
# What is the current status of ML for chemistry?

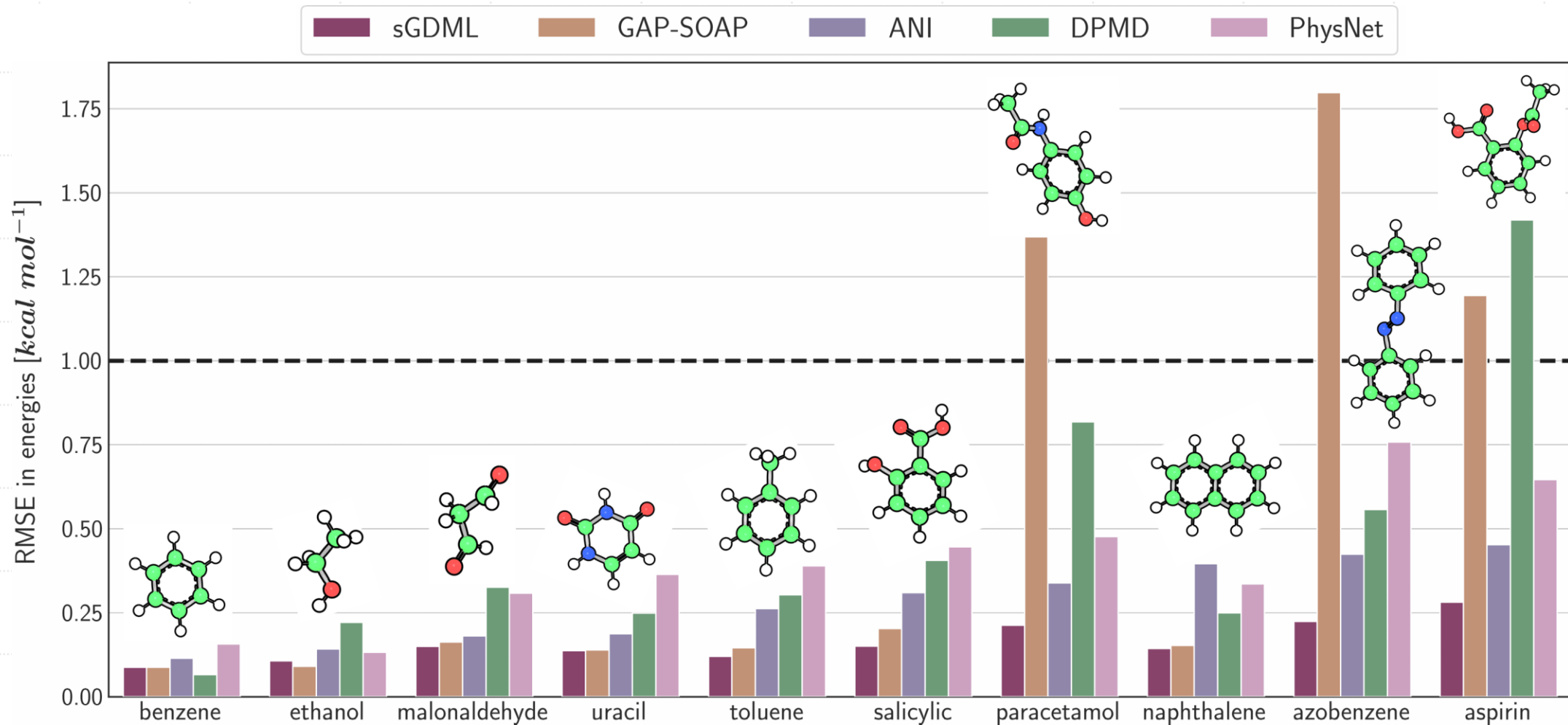
$$\begin{aligned} f'(x) &= \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h} \\ f(x) &= \lim_{h \rightarrow 0} \frac{(x+h)^2 - x^2}{h} \\ &= \lim_{h \rightarrow 0} \frac{x^2 + 2xh + h^2 - x^2}{h} \\ &= \lim_{h \rightarrow 0} \frac{2xh + h^2}{h} \end{aligned}$$

$$\begin{aligned} &= \lim_{h \rightarrow 0} \frac{1}{2\sqrt{x}} \\ f(x) &= \lim_{\Delta x \rightarrow 0} \frac{f(x+\Delta x) - f(x)}{\Delta x} \\ f(a) &= \lim_{h \rightarrow 0} f(a+h) \end{aligned}$$

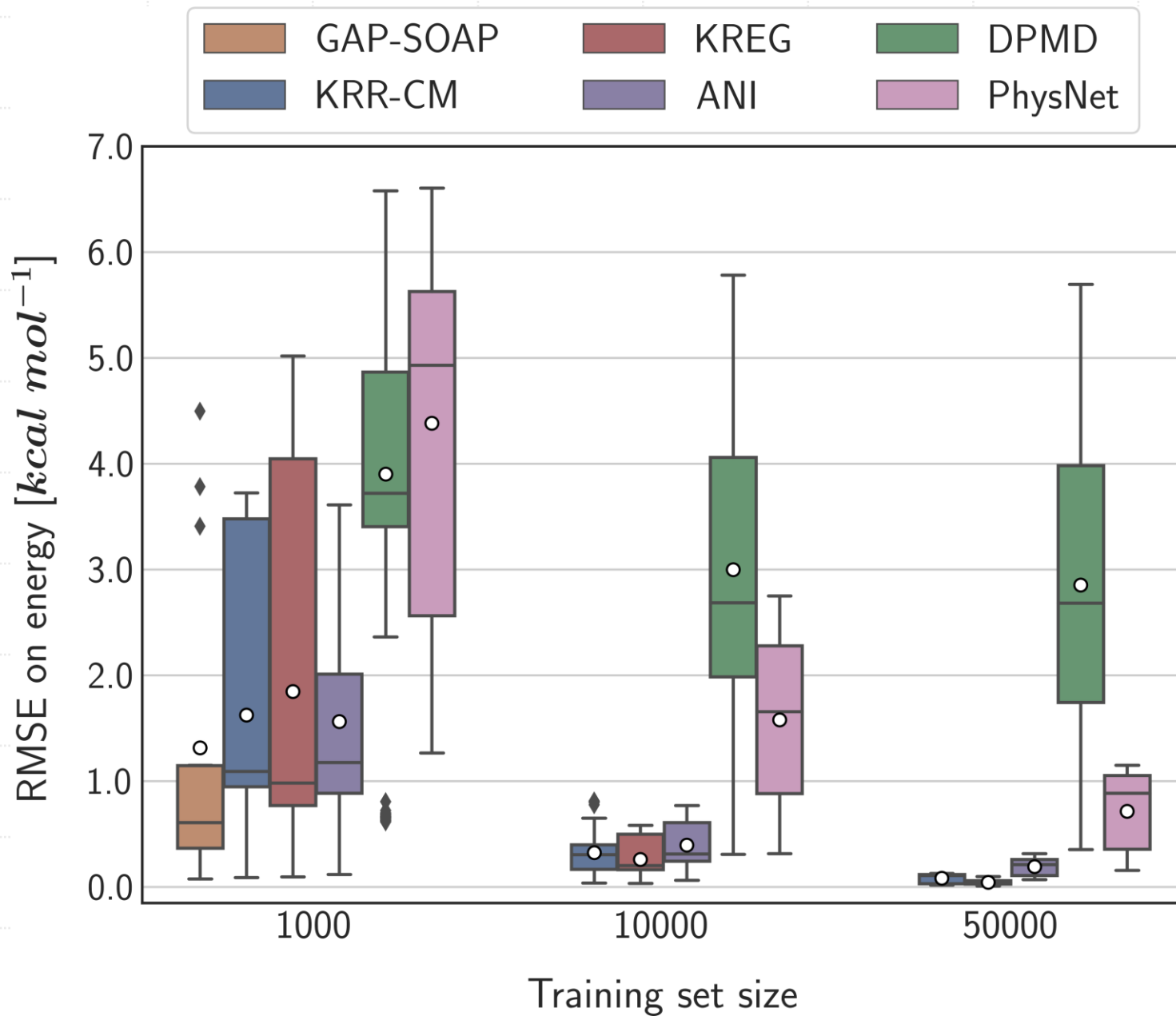


The ground state





- MD17 Database
- Energy + Force
- $N_{train} = 1k$ ;  $N_{model} = 20$ ;  $N_{test} = 20k$



- MD17 Database
- Energy only
- $N_{test} = 20k$



# The excited states



Simulating excited states is much more challenging:

1. They usually correspond to electronic densities that are difficult to compute
2. They are strongly anharmonic
3. They cluster in state bundles, mixing with each other

ML can simulate excited states for the ground state equilibrium

If a dataset spanning the  $3N_{at}-6$  dimensions is available, ML can deliver excellent fittings

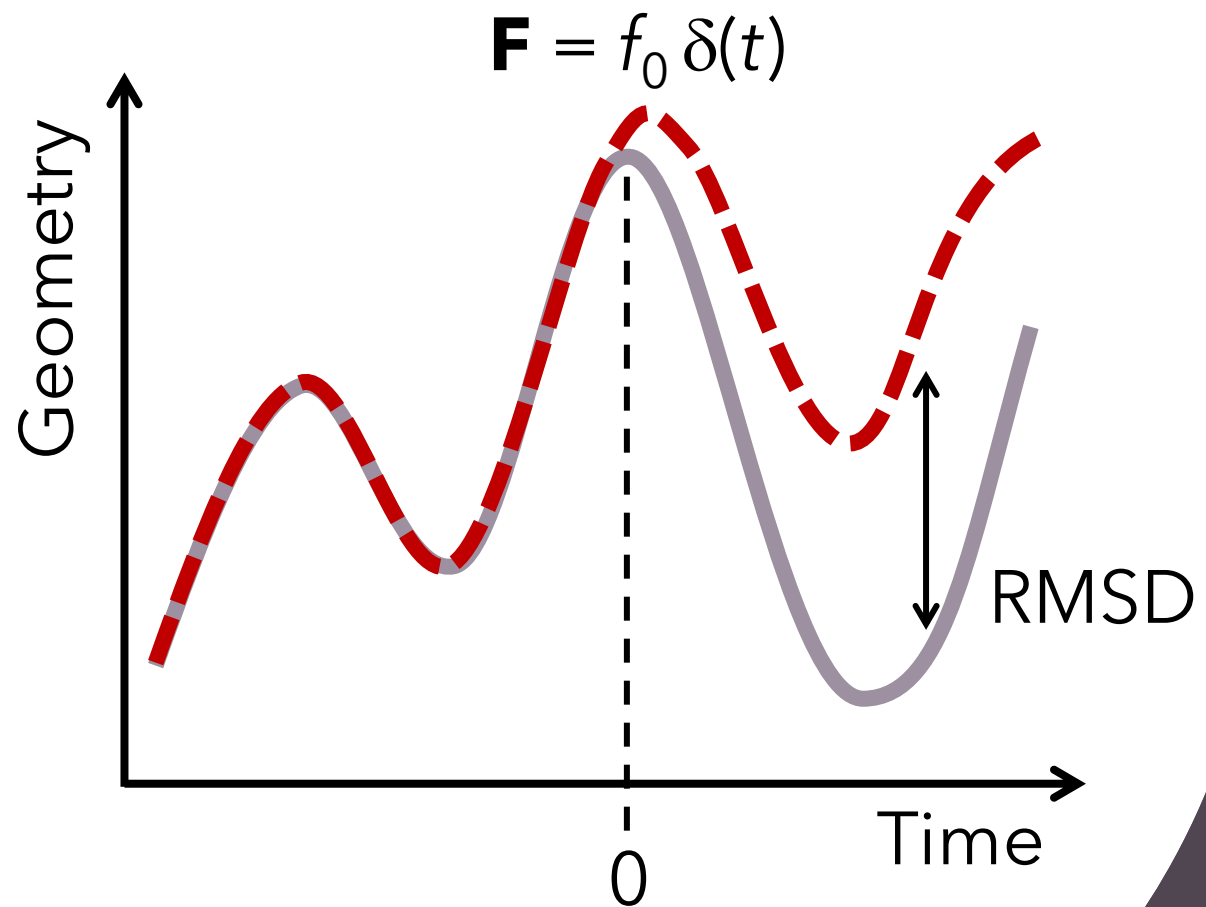
However, for sparse datasets, robust ML protocols are still missing



Matheus Bispo

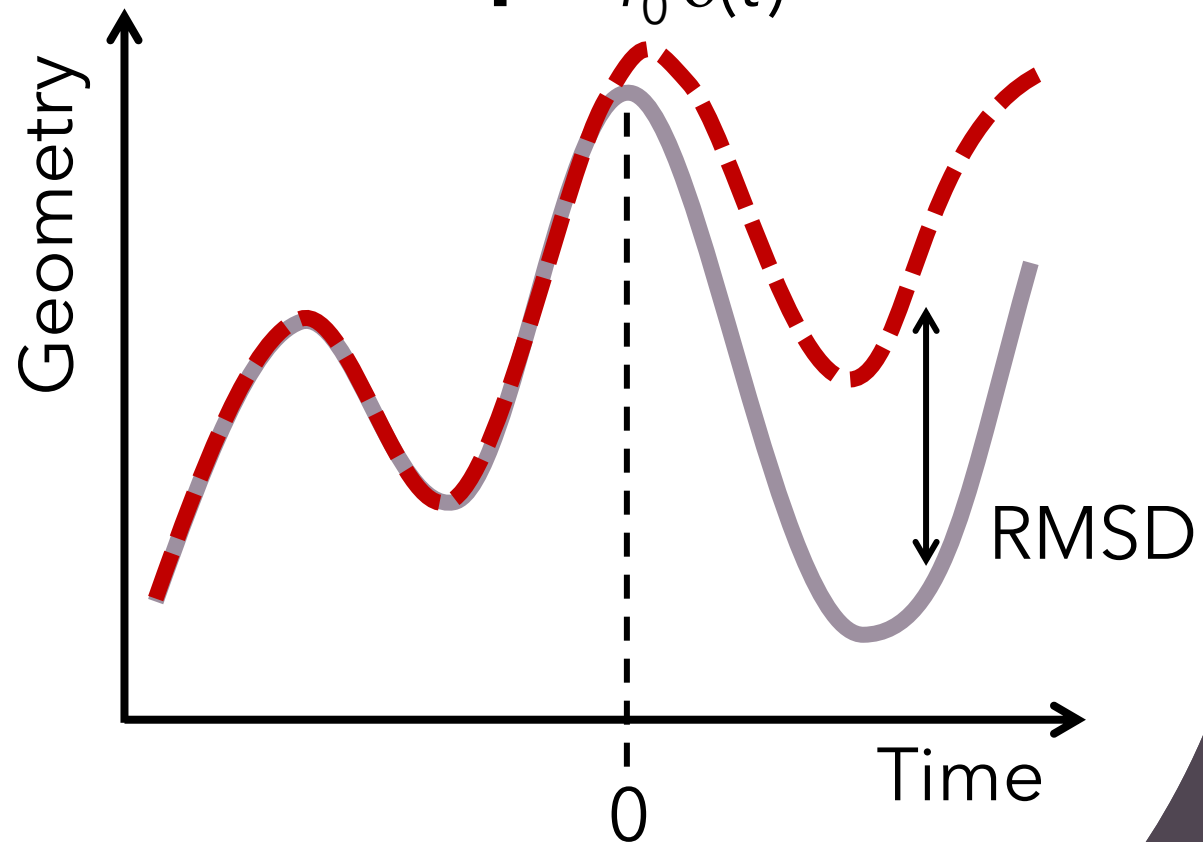
# The challenge

# Effect of force uncertainty



# Effect of force uncertainty

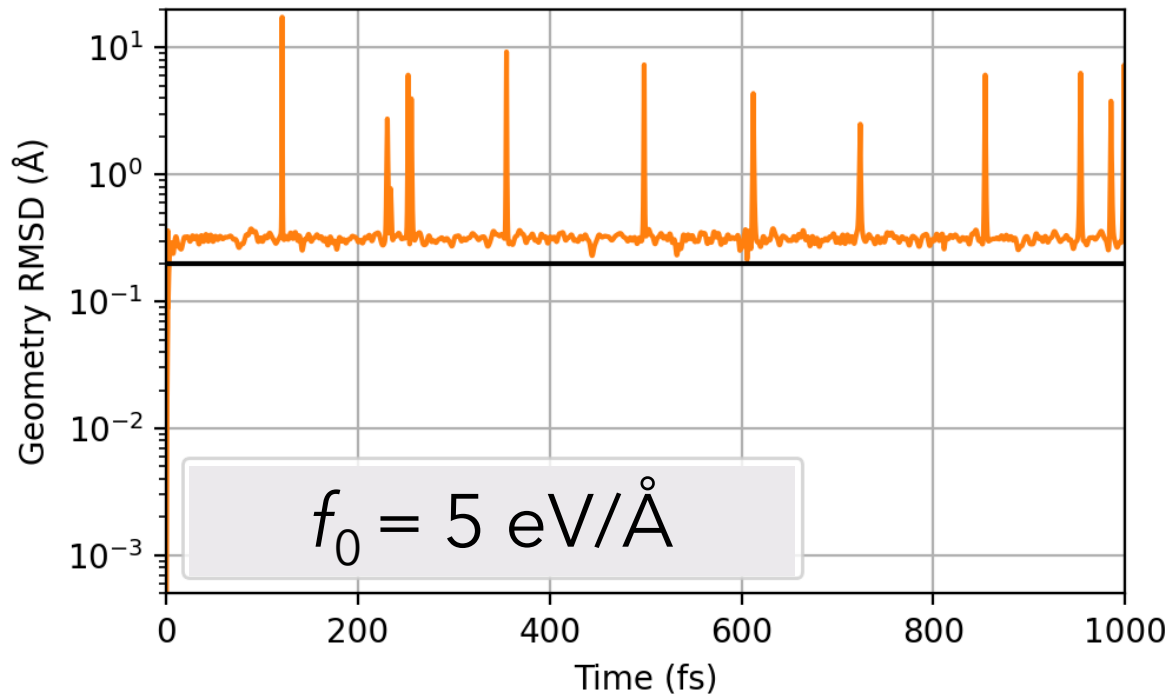
$$\mathbf{F} = f_0 \delta(t)$$



## Geometrical accuracy

We want results better than  $0.2 \text{ \AA}$

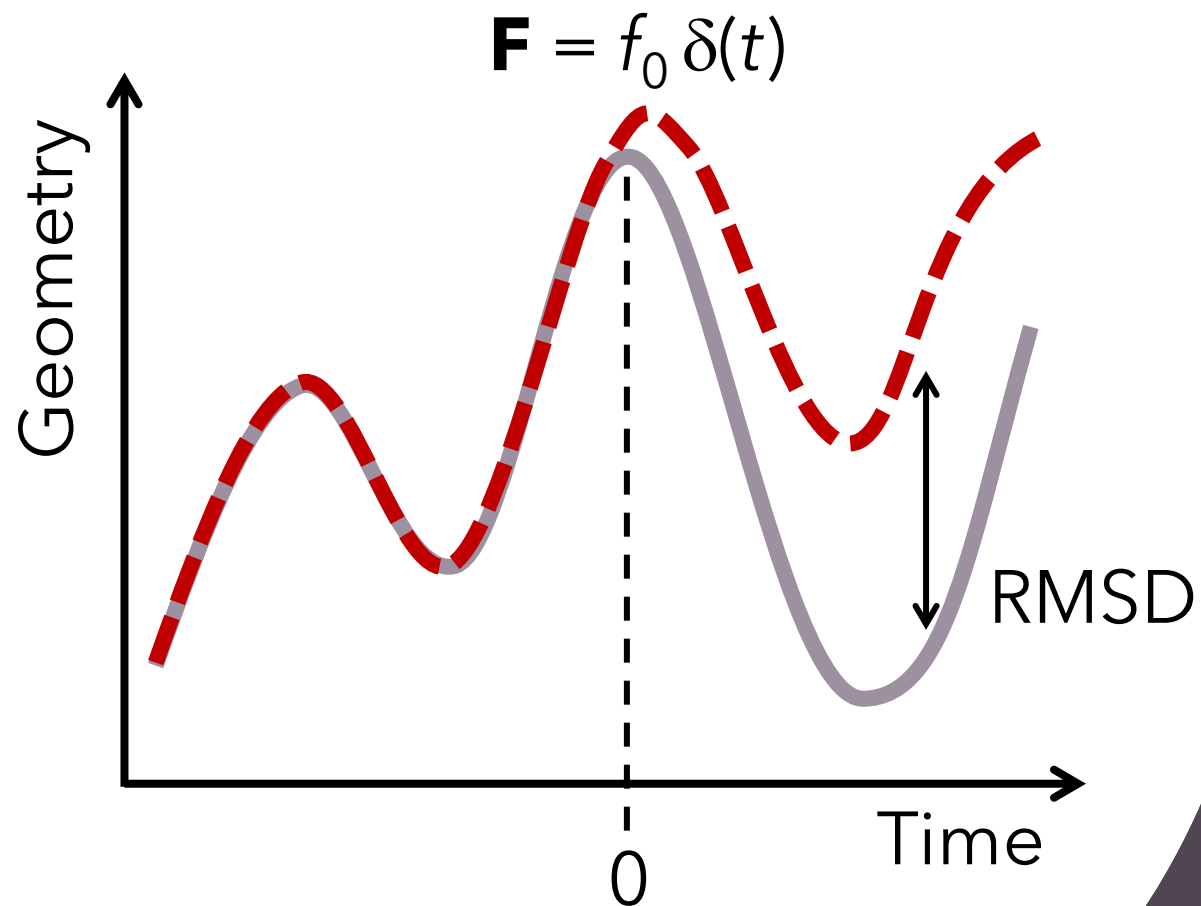
- A-SBH 33D
- dynamics on  $E_1$



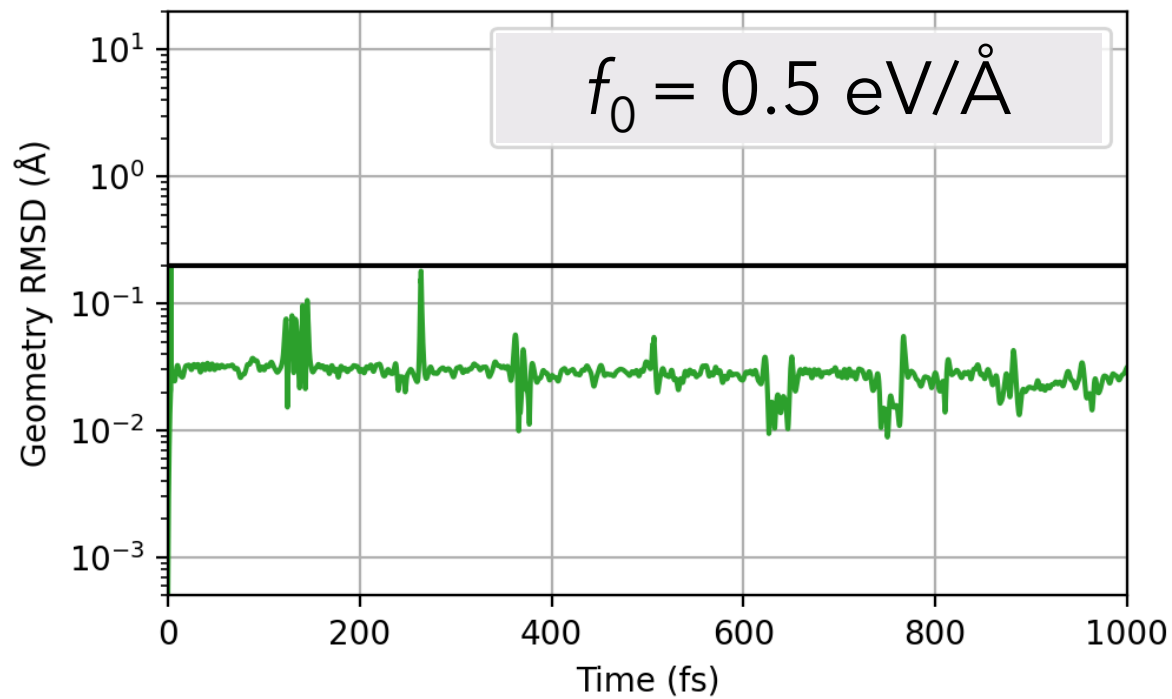
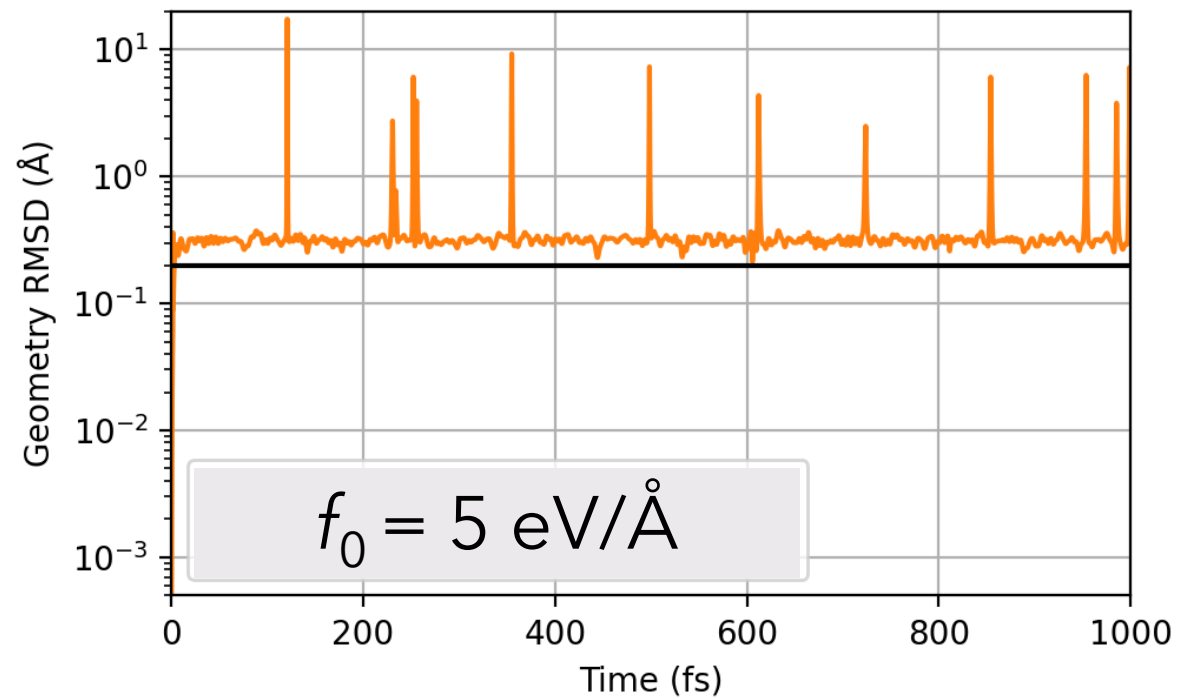
## Geometrical accuracy

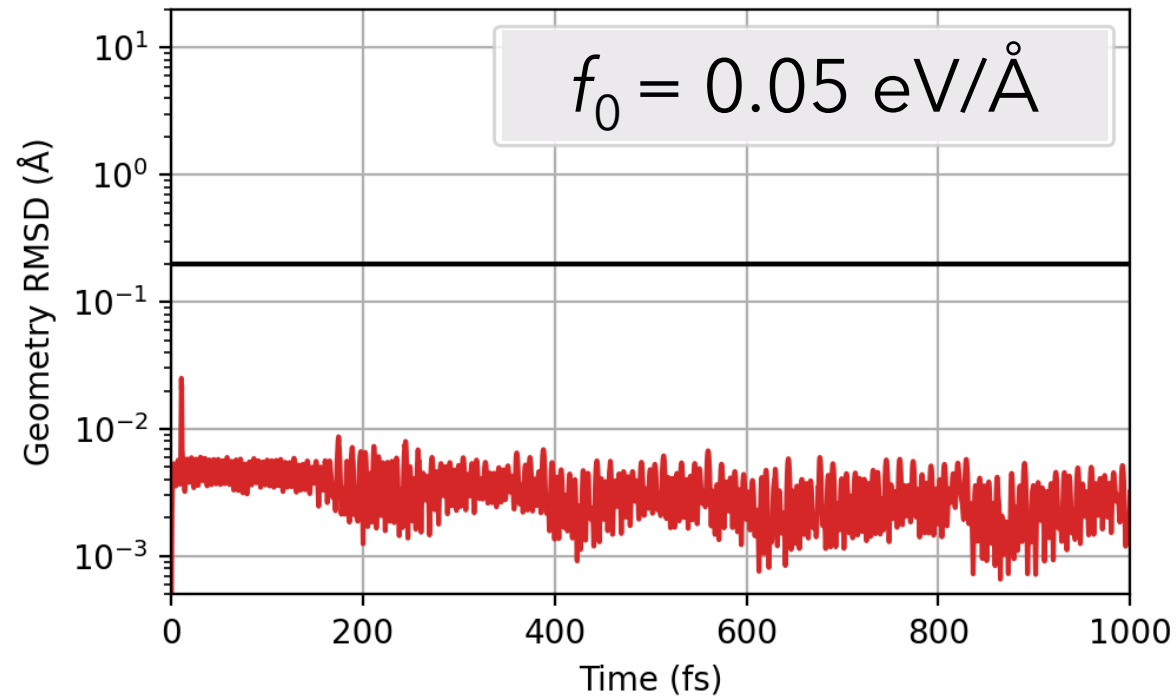
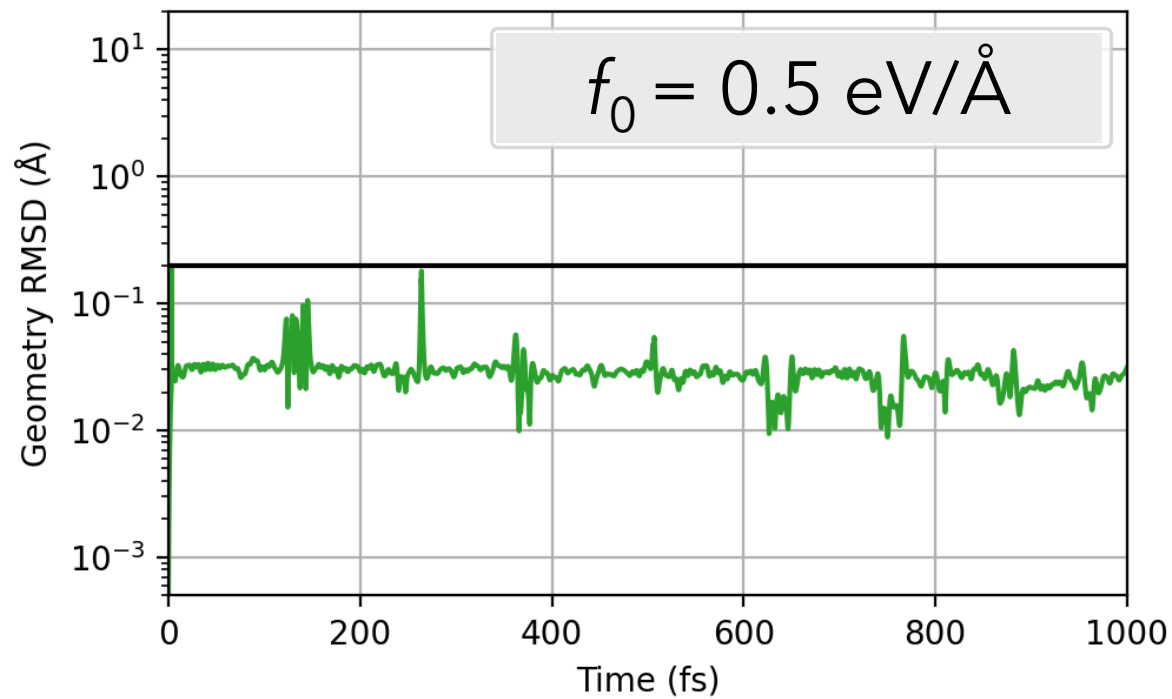
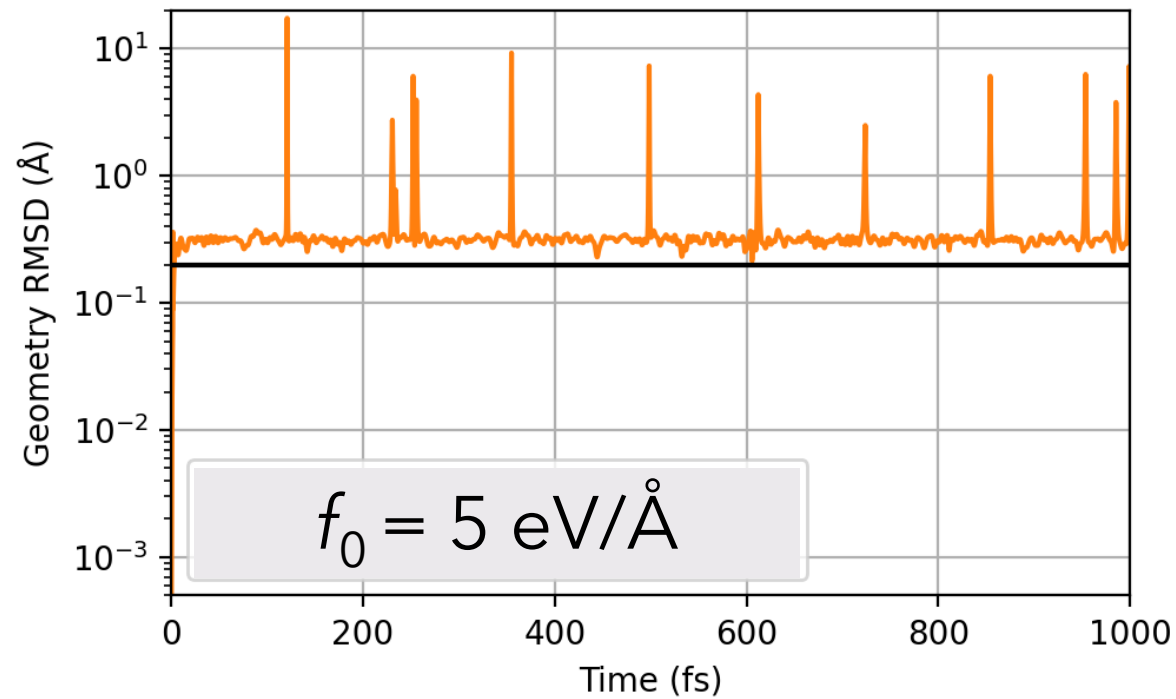
We want results better than  $0.2 \text{ Å}$

## Effect of force uncertainty

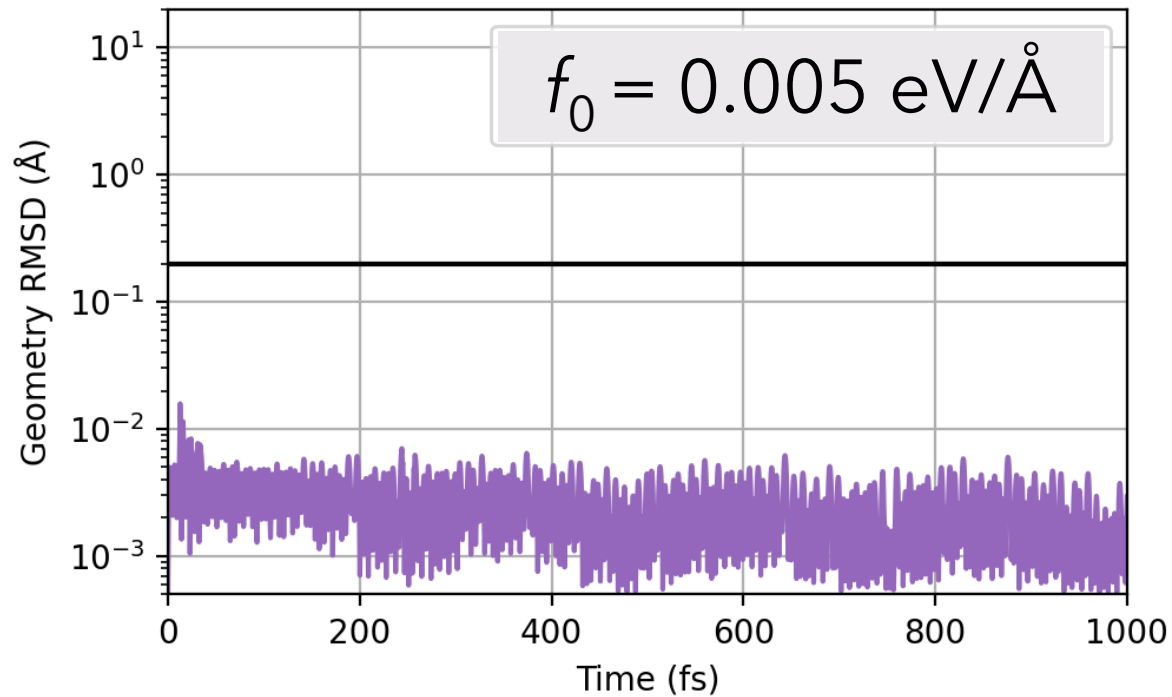
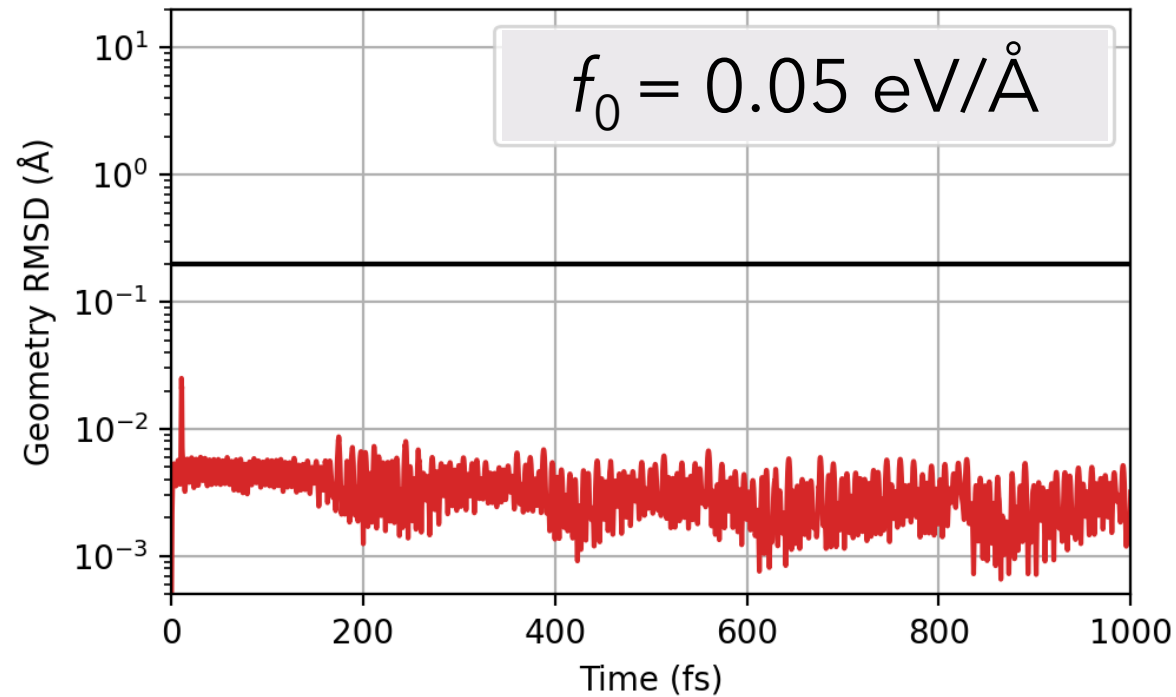
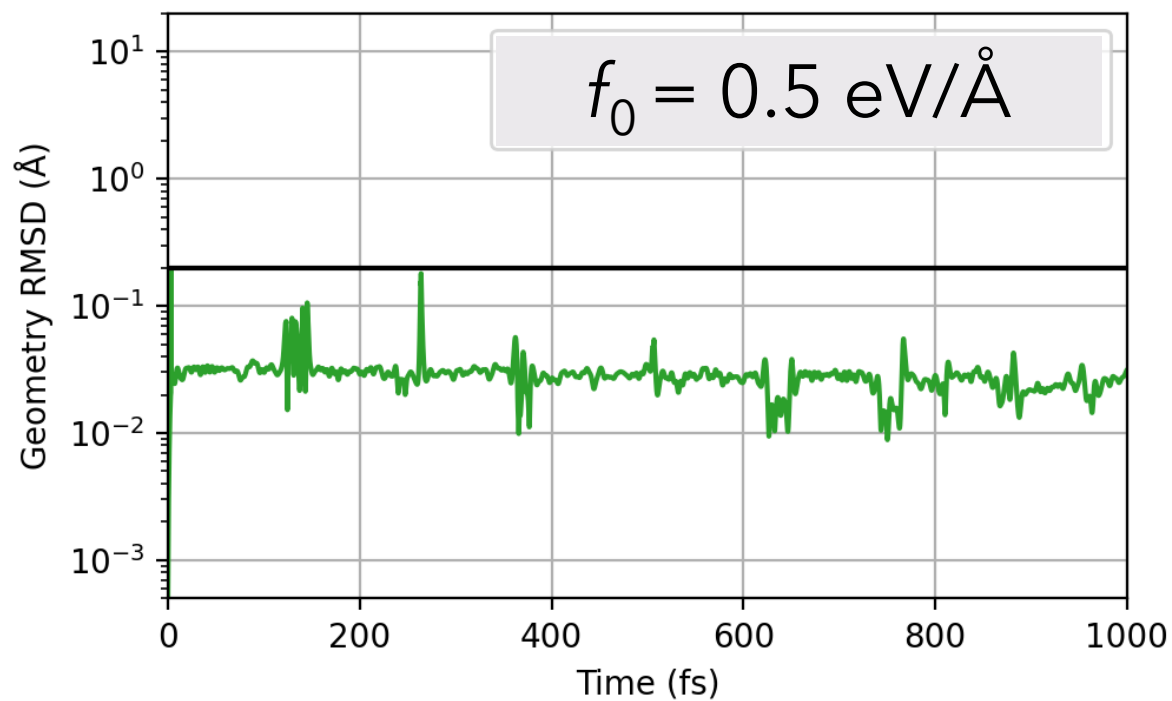
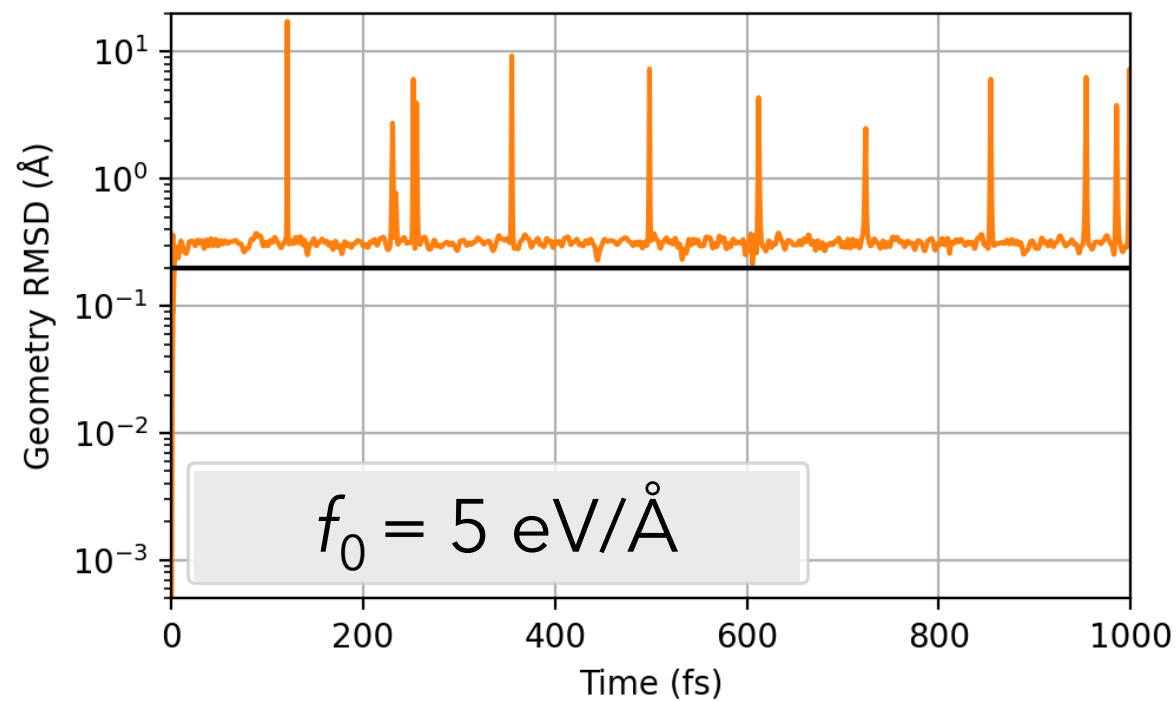


- A-SBH 33D
- dynamics on  $E_1$









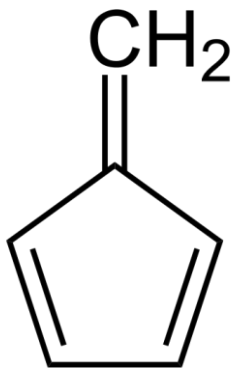
We must predict forces better than  
 $0.5 \text{ eV/\AA}$  ( $0.001 \text{ Hartree/Bohr}$ )

(Maximum absolute error)



Max Pinheiro Jr

# ML-NAMD test cases



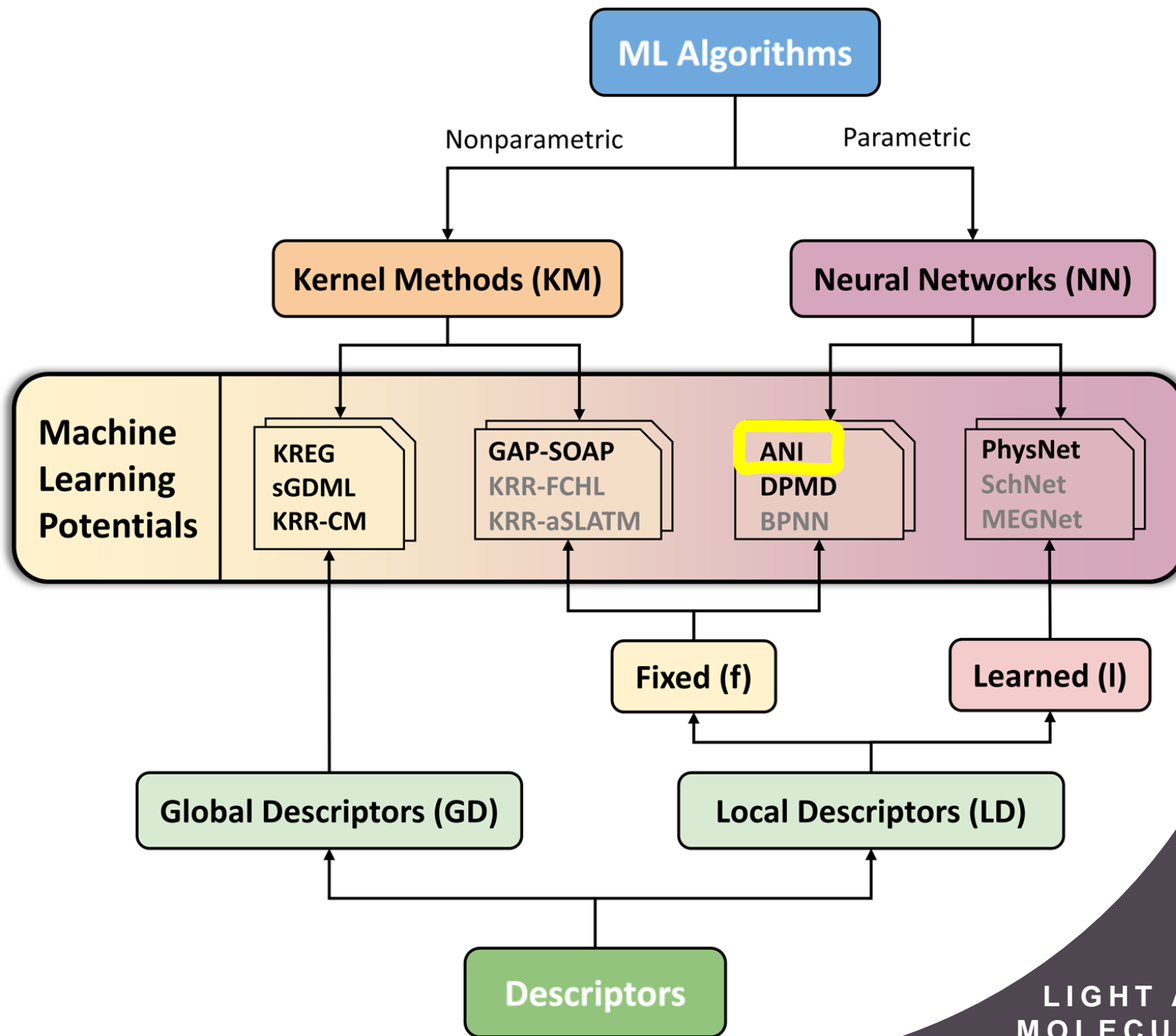
# Fulvene test

Dataset:

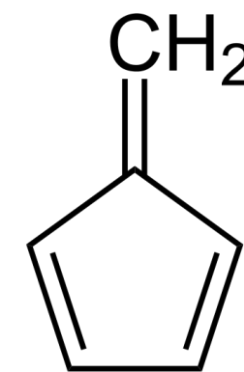
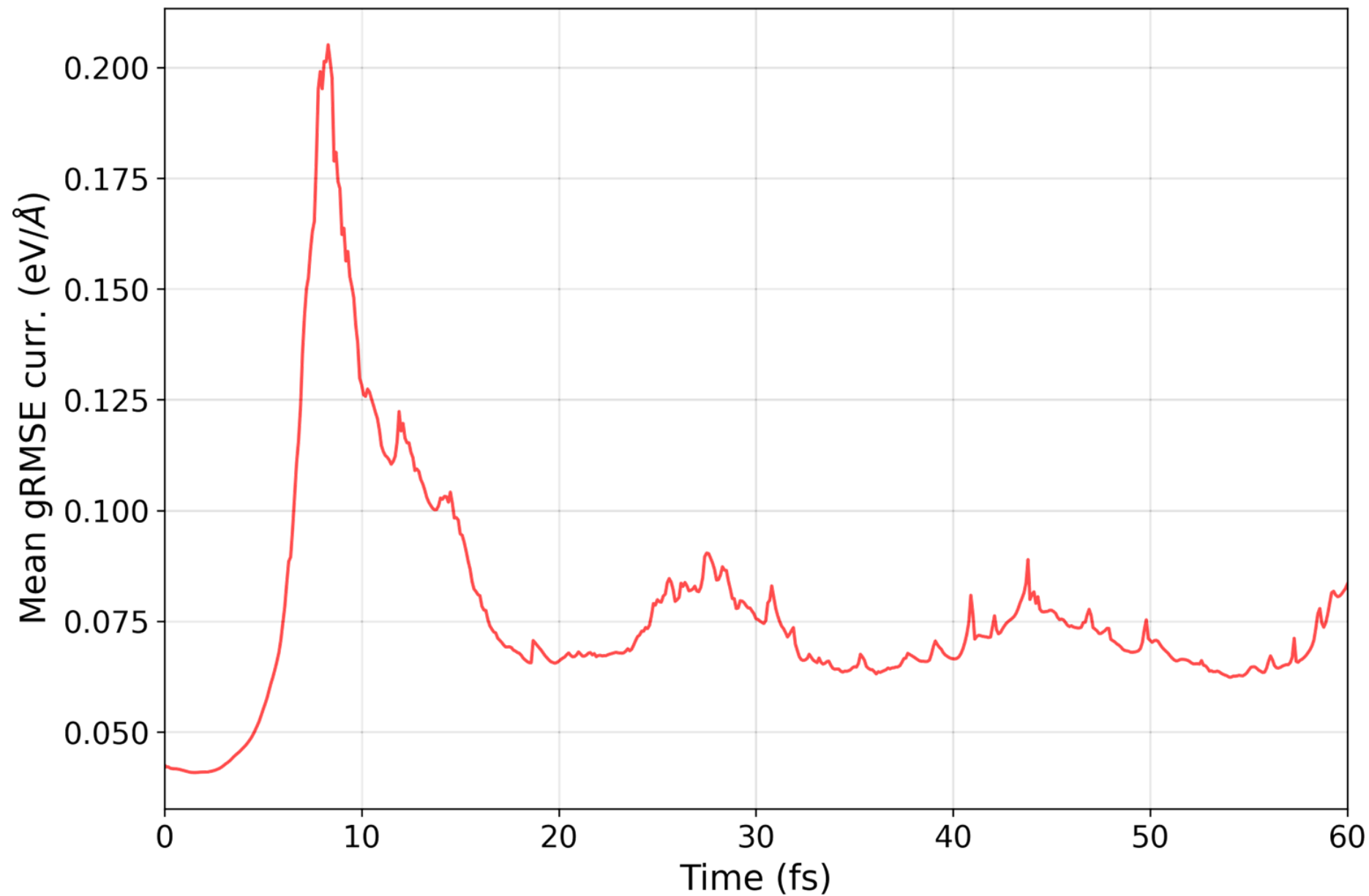
- DC-FSSH / CASSCF
- 200 traj; 2 states
- $t_{\max} = 60$  fs;  $\Delta t = 0.1$  fs

ML potential:

- ANI on 40k
- energy + forces



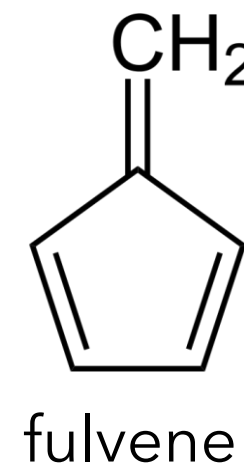
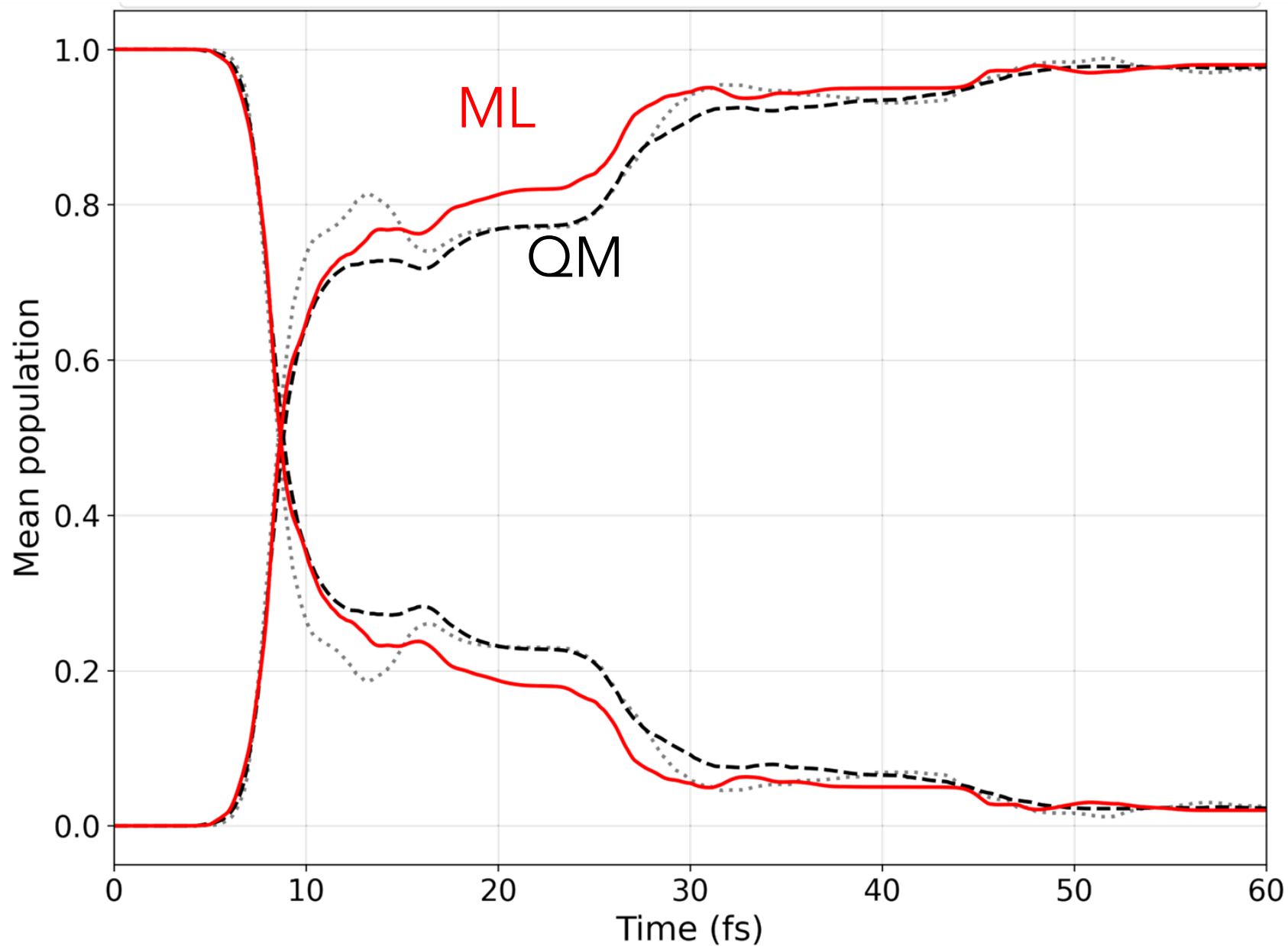
# Gradient (current state)



fulvene

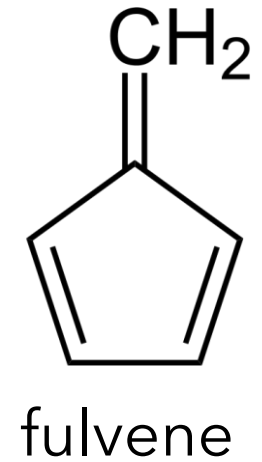
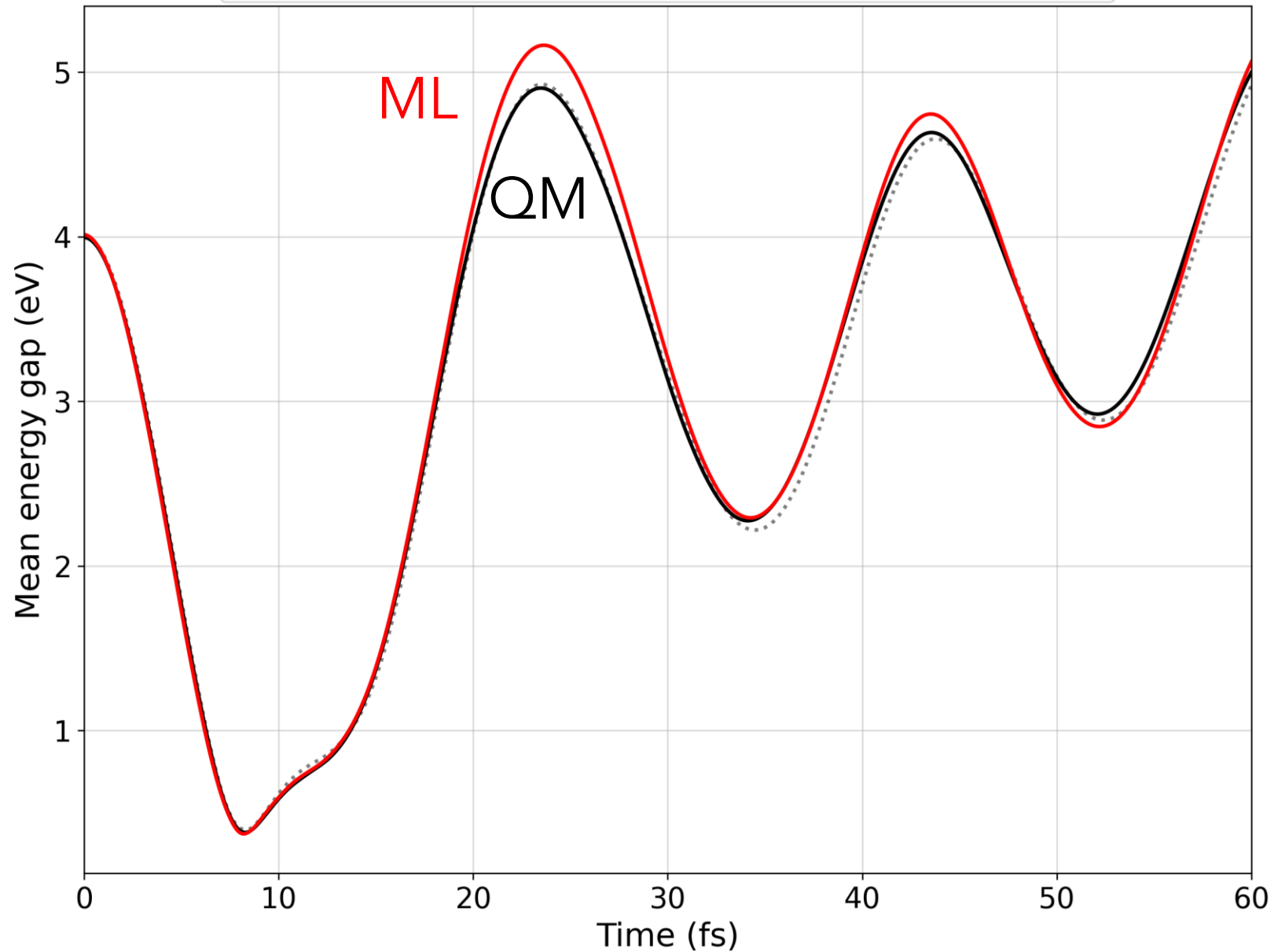
Error for the gradient norm (current state)

# Population



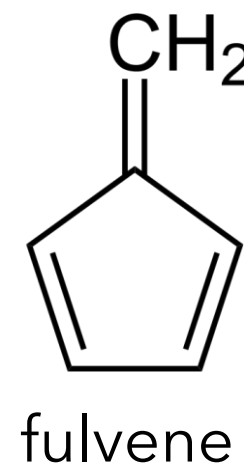
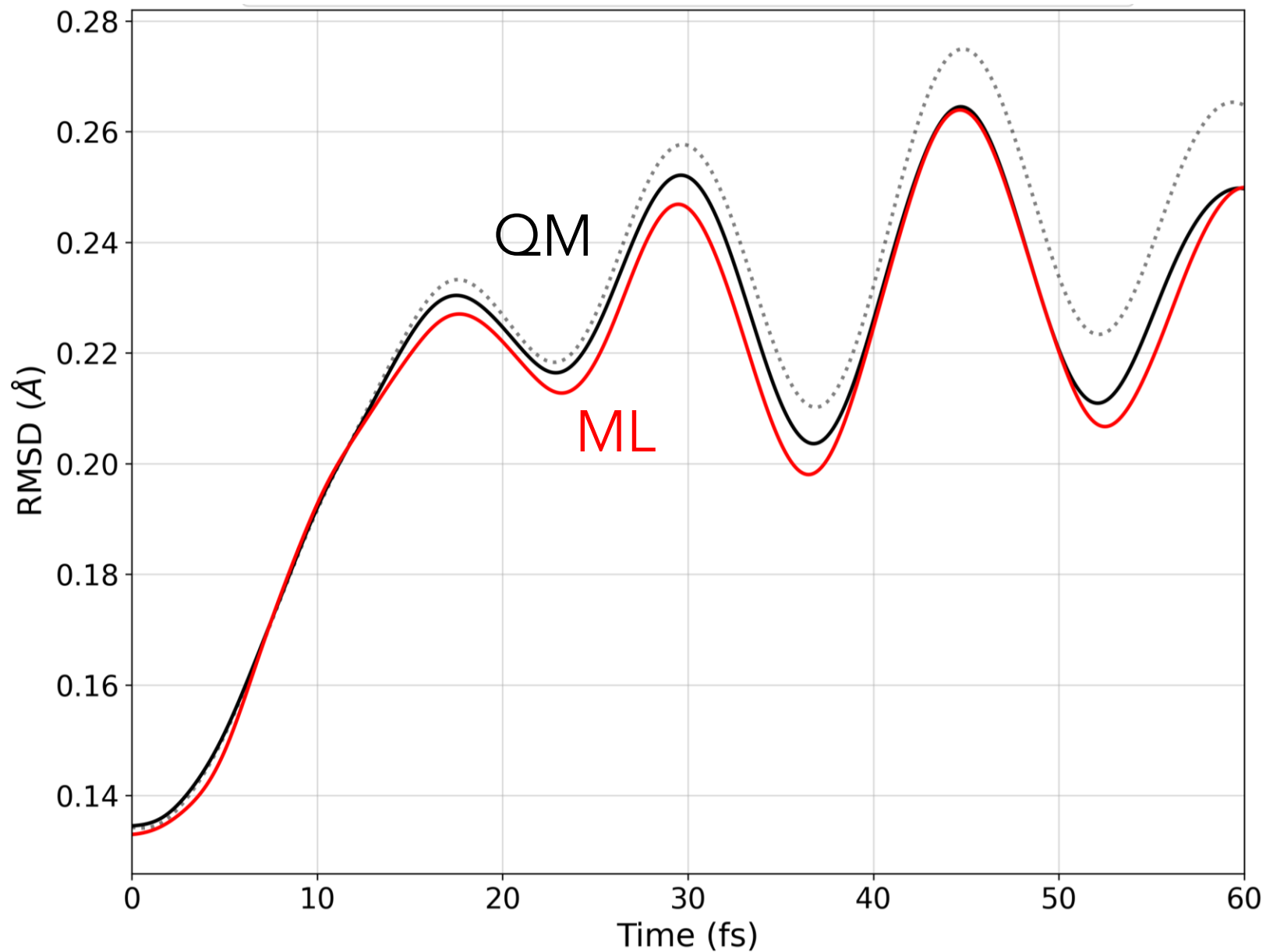
Population averaged over  
100 new trajectories

# Energy gap



$\Delta E_{10}$  evolution averaged over 100 new trajectories

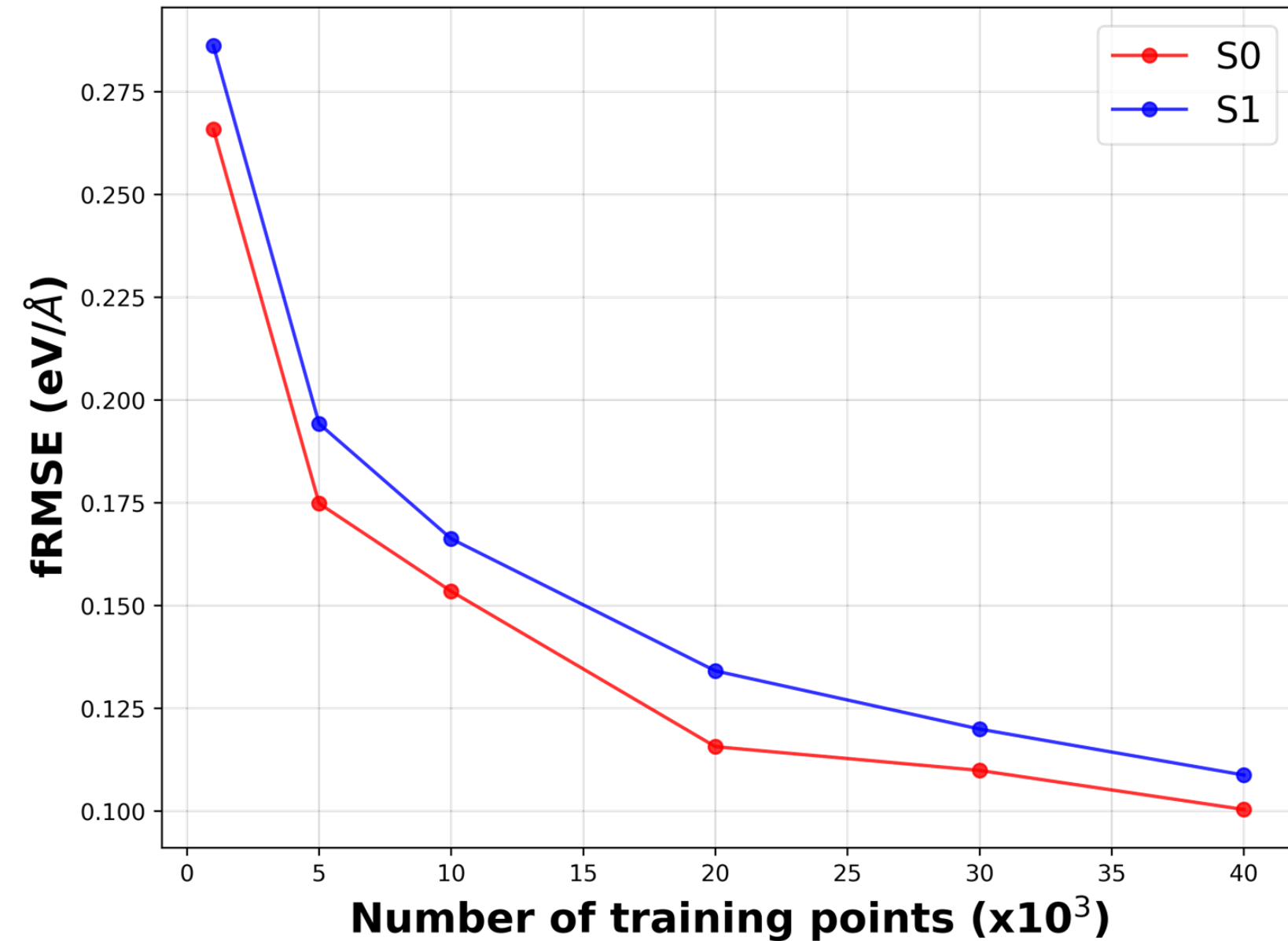
# Geometry



Geometry evolution averaged over 100 new trajectories

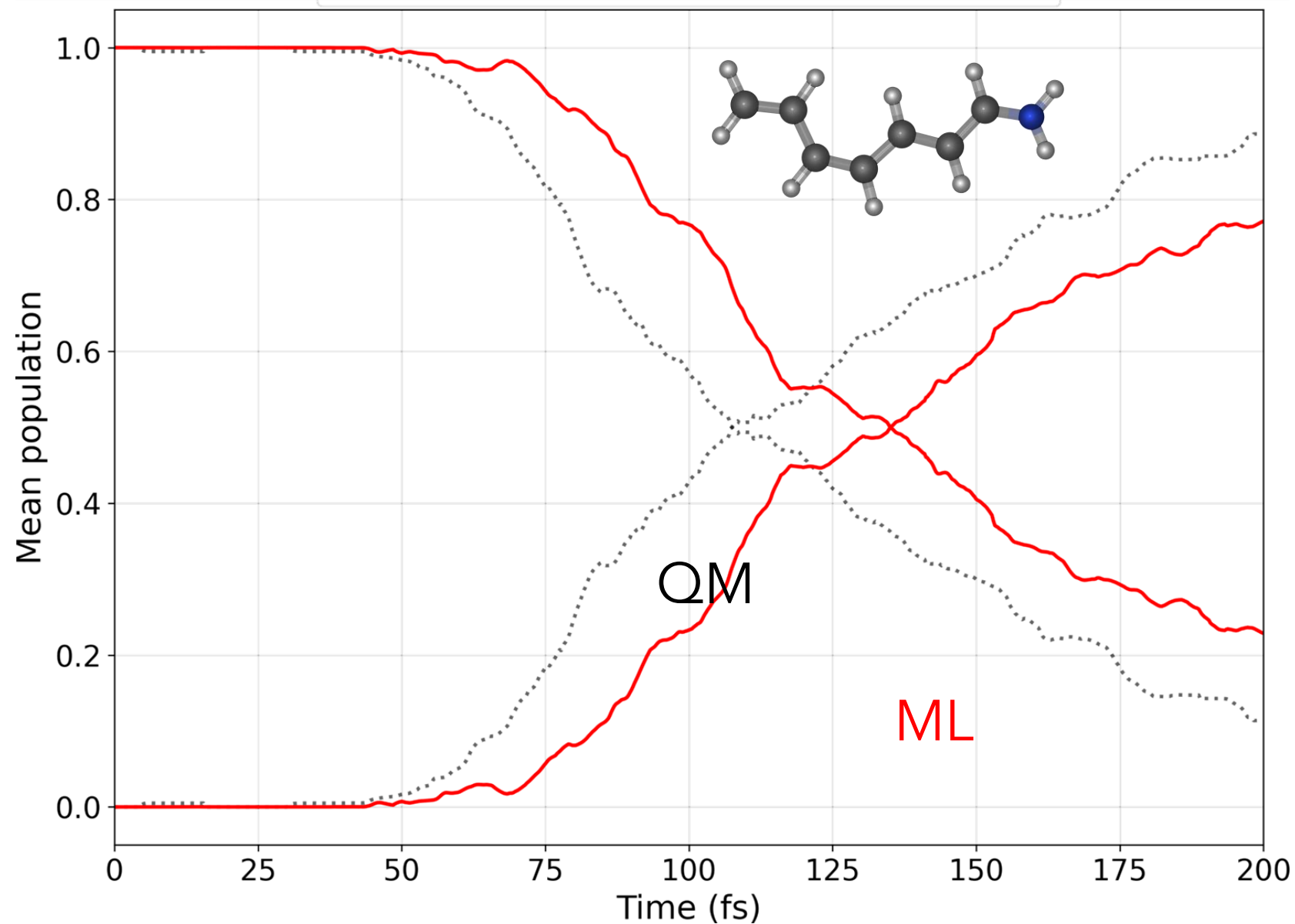


# Problem 1: Training set size



40k points are too expensive to adopt as routine protocol

# Problem 2: Lack of robustness

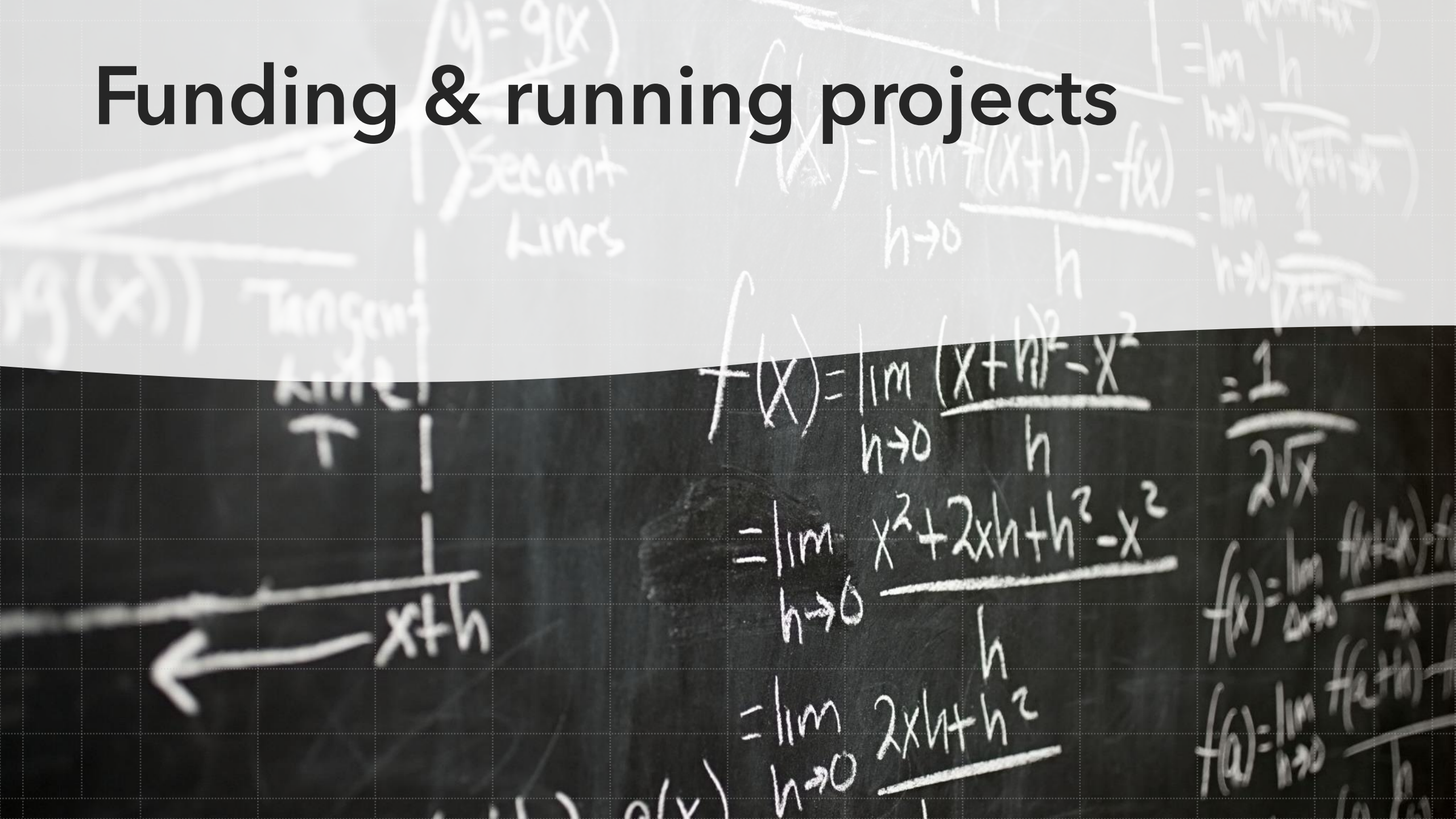


Using the same ML model with the same training set size for another molecule may not give adequate results

To address both problems, we are testing **active learning**:

1. We train multiple machines with an initially small number of points
2. Then, we run dynamics until the prediction between machines diverges
3. The geometries showing divergence are included in the training set, and the machine is retrained
4. Repeat the procedure until no divergence is observed

# Funding & running projects





Sep 19

Aug 24

## Subnano (ERC AdG)

**Goal:** ML potentials for dynamics in long timescale

- 72 months postdoc
- 1 PhD

Apr 24

Mar 25

## Sony (Award)

**Goal:** ML potentials NAMD of large chromophores

- 12 months postdoc

Jan 24

Dec 27

## MLChem (A\*Midex 2022 Res. & Training)

**Goal:** Establish an international AI center for chemistry at AMU

- 3x36 months postdoc
- 1 PhD

**A\*Midex MLChem**

ICR  
Institut Chimie Radicalaire

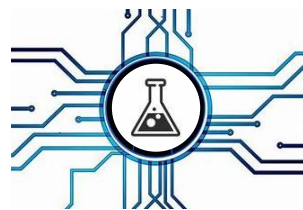


- Barbatti
- Ferré

LIS  
LABORATOIRE  
D'INFORMATIQUE  
& SYSTÈMES  
UMR 7020



- Artières
- Kadri



Center  
of AI for  
Chemistry

MLChem aims to establish  
an **international AI center  
for chemistry** at AMU



廈門大學化學化工學院  
College of Chemistry and Chemical Engineering, Xiamen University

- Dral



# MLChem mission 1: ML showcase

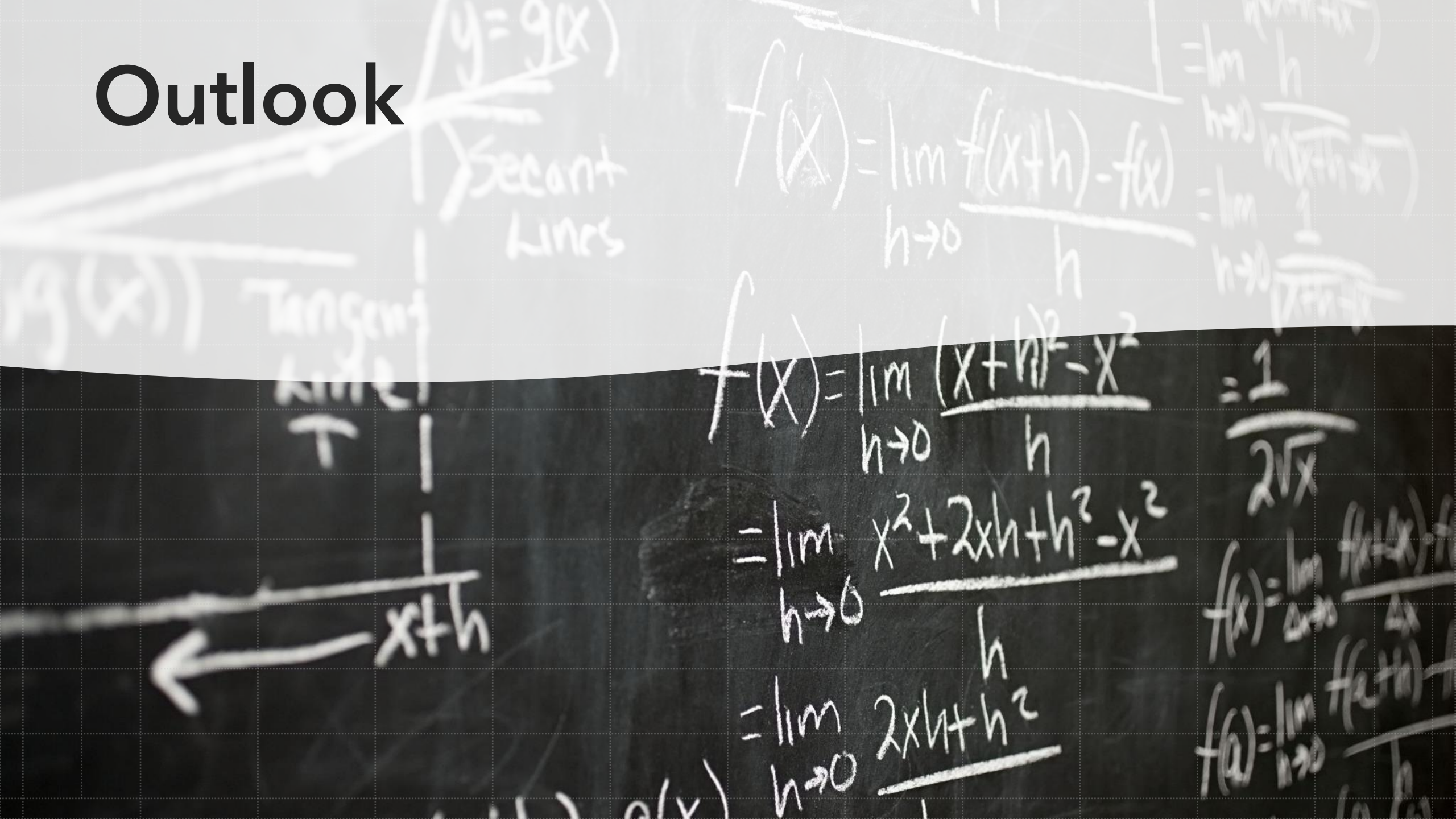
- Developing new methodologies for computational chemistry using ML for energy and charge propagation in organic crystals
- Implementing unsupervised-learning approaches for analyzing computational chemistry data
- Creating open data protocols for sharing ML results in a dedicated system server



# MLChem mission 2: AI culture

- Offering project-tailored consulting on using AI solutions for chemistry groups at AMU
  1. help to identify preexistent software and hardware available
  2. help writing grant proposals to secure funds to implement those solutions.
- Offering courses and tutorials on ML for chemistry at several levels
  1. Two virtual ML schools
  2. Courses for ED250
  3. ML at master programs

# Outlook



Supervised ML has great potential for accelerating excited-state potential energy predictions.

However, no robust protocol has been published  
(If the model and training procedure worked for a molecule, it does not mean they will also work for another one)

Problems boil down to the high accuracy required in the predictions in highly multidimensional spaces

Our experience developing methods, benchmarks, and programs for ML for theoretical chemistry can be a seed for creating a culture of AI for chemistry at AMU

Our main challenge may not be scientific.

It is our need for career attractiveness for young researchers.

# SUBNANO



erc



# A\*Midex

Initiative d'excellence Aix-Marseille



institut  
universitaire  
de France



[www.barbatti.org](http://www.barbatti.org)



[mario.barbatti@univ-amu.fr](mailto:mario.barbatti@univ-amu.fr)



[@MarioBarbatti](#)