



# COMPUTING CATALYTIC REACTION TIMES AND PATHS WITH MACHINE LEARNING AND RARE EVENTS SAMPLING METHODS

---

THOMAS PIGEON

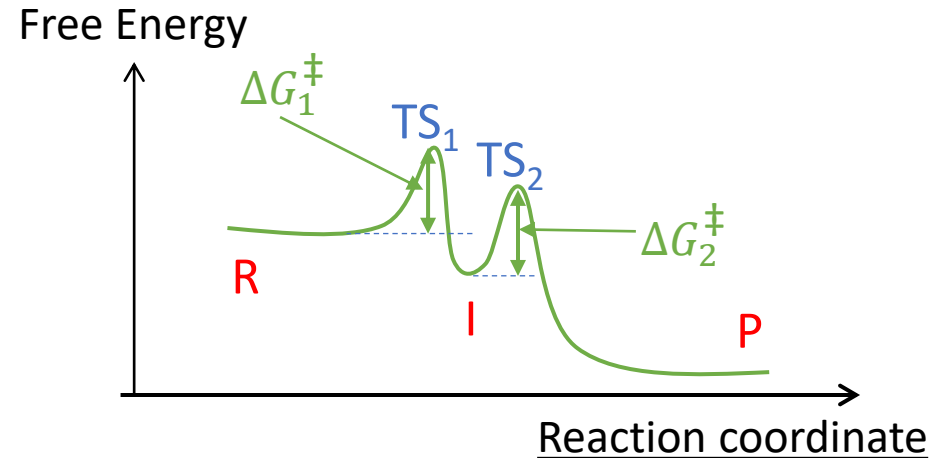
Pascal Raybaud<sup>1</sup>, Manuel Corral Valero<sup>1</sup>,  
Ani Anciaux-Sedrakian<sup>2</sup>, Maxime Moreaud<sup>2</sup>  
Gabriel Stoltz<sup>3</sup>, Tony Lelièvre<sup>3</sup>

<sup>1</sup> IFPEN: Catalysis, Biocatalysis and Separation

<sup>2</sup> IFPEN: Digital Science and Technology

<sup>3</sup> CERMICS, Ecole des Ponts ParisTech, and Equipe-projet MATHEMATICALS, Inria Paris,

Targets : compute reaction rates (average times)  
identify reaction paths (how atoms rearranges themselves)



Different methods exist:

- Transition State Theory (TST)  
→ Need to estimate free energy
- Directly from the time evolution of the system ?  
→ Need Molecular Dynamics (MD)

Simulates the dynamic of the system by adding a thermostat to newton equations of motion

ex. Langevin formalism<sup>1</sup>

$$\begin{cases} dq_t = M^{-1}p_t dt \\ dp_t = -\nabla V(q_t)dt - \gamma p_t dt + \sqrt{2\gamma M k_B T} dW_t \end{cases}$$

NVE ensemble

NVT ensemble

Preserves energy

Dissipate energy

Provides energy

Newton equation

Langevin part

**Not efficient for the simulation of rare events due to high energy barriers and entropic bottlenecks**

**Time scales:** integration time step :  $\sim 10^{-15} s$  rare event rate  $\sim 10^{-9} s^{-1}$  to  $10^3 s^{-1}$

MD based approaches to overcome barriers:

→ biased MD such as Metadynamics<sup>2</sup>, Blue-Moon sampling<sup>3</sup> ...

**Dynamics is lost but we can estimate free energy**

→ rare events sampling methods such as Adaptive multi-level splitting<sup>4</sup>

**Dynamics preserved and rates can be “directly” be computed**

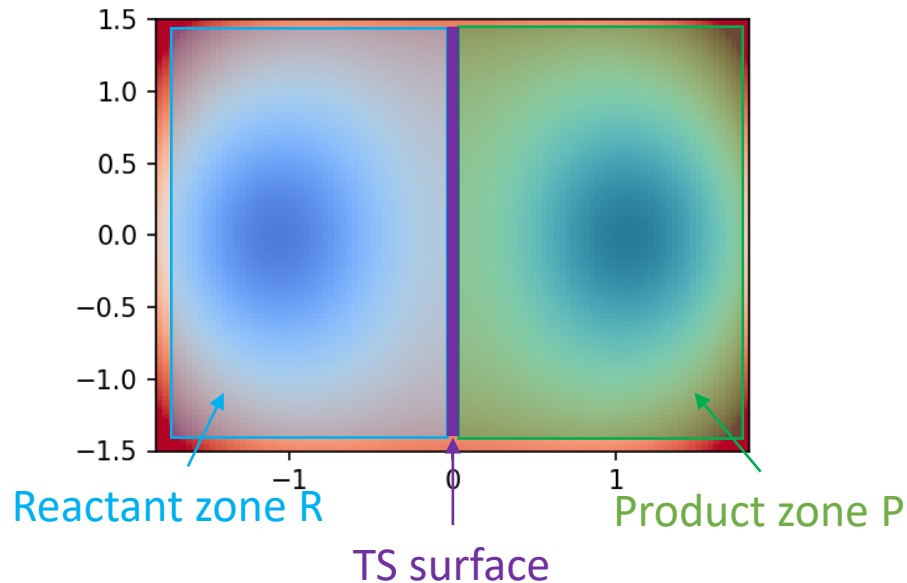
<sup>1</sup> P. Langevin P. (1908), Comptes-Rendus de l'Académie des Sciences, 146, 530

<sup>2</sup> A. Laio, M. Parrinello, (2002) PNAS, 99, 20, 12562

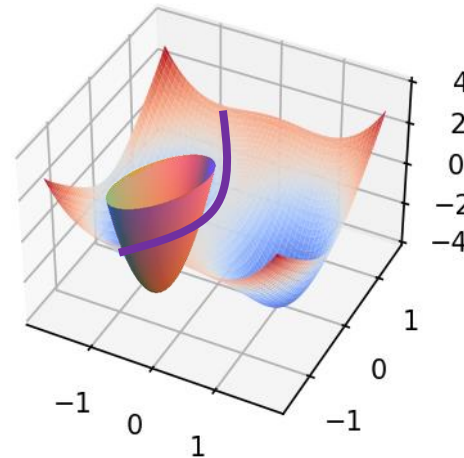
<sup>3</sup> E. A. Carter, G. Ciccotti, J. T. Hynes, R. Kapral, (1989). Chem. Phys. Lett., 156, 5, 472

<sup>4</sup> Cérou, F., & Guyader, A. (2007) Stoch. Anal. Appl., 25, 2, 417

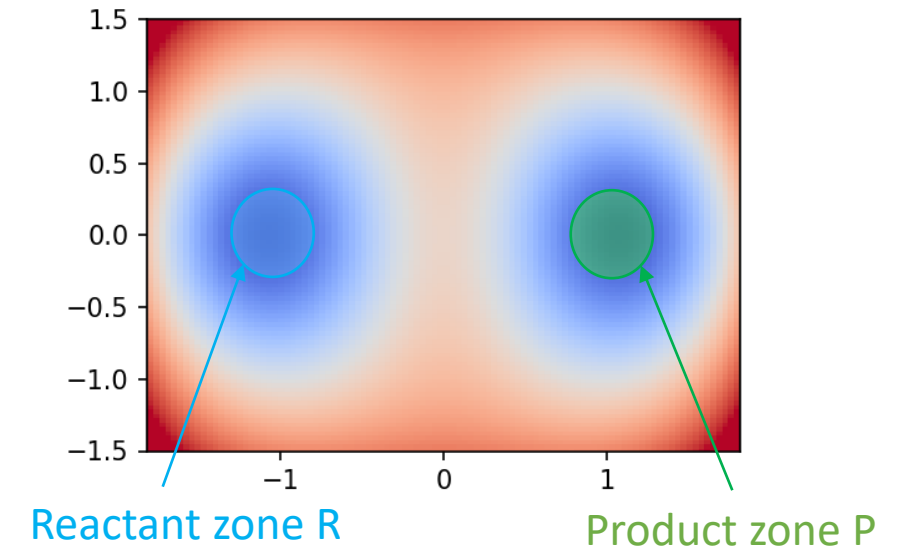
## Transition State Theory



## 2-dimensional potential



## Hill relation



**Rate** = probability of being in TS with respect to R  
 × frequency of decomposition to P

$$k^{TST} = p(TS | R) \phi_{TS \rightarrow P} \quad k^{hTST} = e^{-\frac{\Delta G^\ddagger}{k_B T}} \frac{k_B T}{h}$$

Sensitive to the TS definition  
 TST overestimates rates ( $\kappa$ )  
 hTST poorly captures entropy

**Rate** = probability of reaching P before R  
 starting from  $\partial R$  × frequency of exits of R

$$k^{Hill} = p_{R \rightarrow P}(\partial R) \phi_R$$

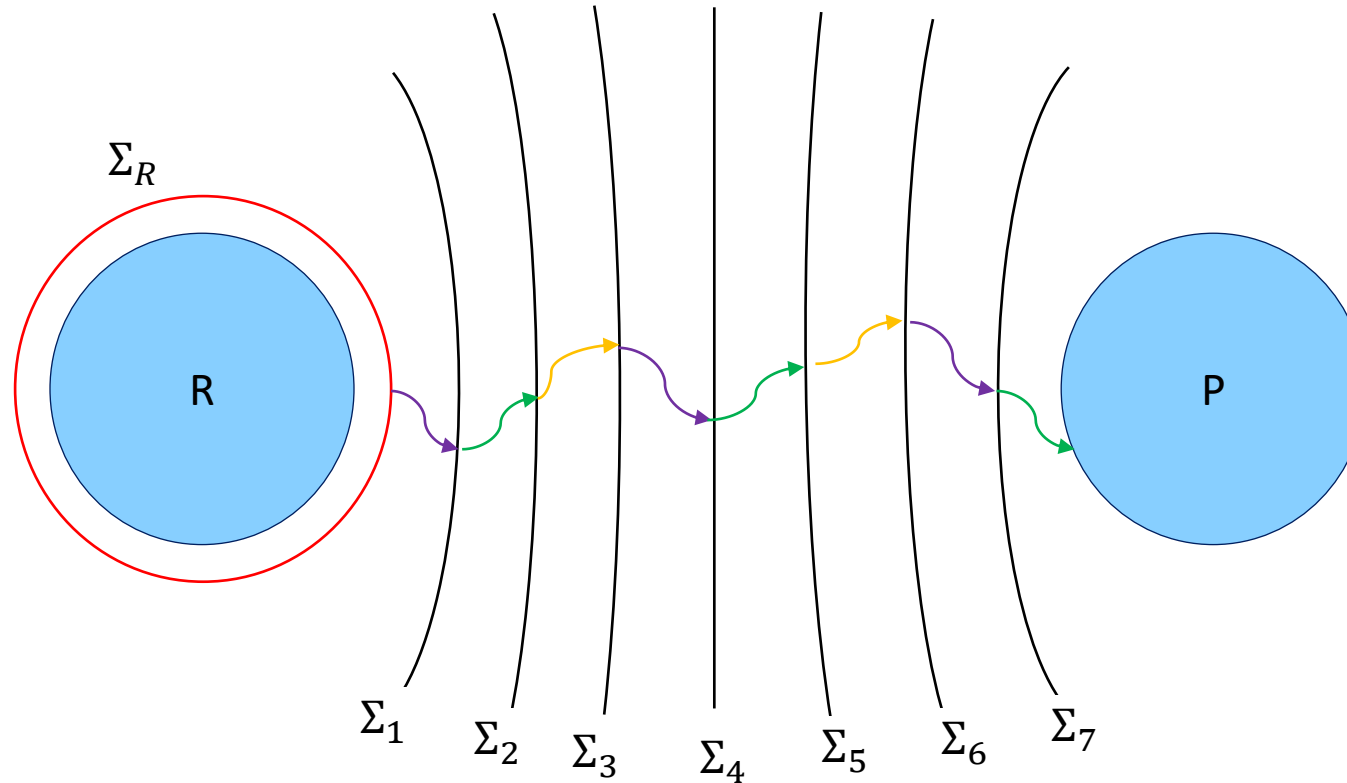
Not extremely sensitive to the  
 definition of R and P

<sup>1</sup> H. Eyring, (1935). J. Chem. Phys., 3, 2, 107

<sup>2</sup> P. Hänggi, P. Talkner, M. Borkovec, (1990) Rev. Mod. Phys., 62, 2, 251

<sup>2</sup> T. Hill, (2012) Free energy transduction in biology: The steady-state kinetic and thermodynamic formalism. Elsevier Science and Technology Books

What is a Multilevel Splitting estimator:



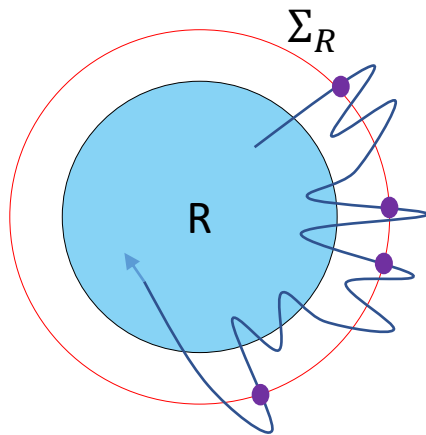
$$p_{R \rightarrow \Sigma_1}(\Sigma_R) p_{R \rightarrow \Sigma_2}(\Sigma_1) p_{R \rightarrow \Sigma_3}(\Sigma_2) p_{R \rightarrow \Sigma_4}(\Sigma_3) p_{R \rightarrow \Sigma_5}(\Sigma_4) p_{R \rightarrow \Sigma_6}(\Sigma_5) p_{R \rightarrow \Sigma_7}(\Sigma_6) p_{R \rightarrow P}(\Sigma_7) \\ = p_{R \rightarrow P}(\Sigma_R)$$

How to place  $\Sigma_i$  and compute  $p_{R \rightarrow \Sigma_{i+1}}(\Sigma_i)$  ?

● AMS aims at estimating  $p_{\Sigma \rightarrow P}^{1,2}$ . It can be split in 3 steps:

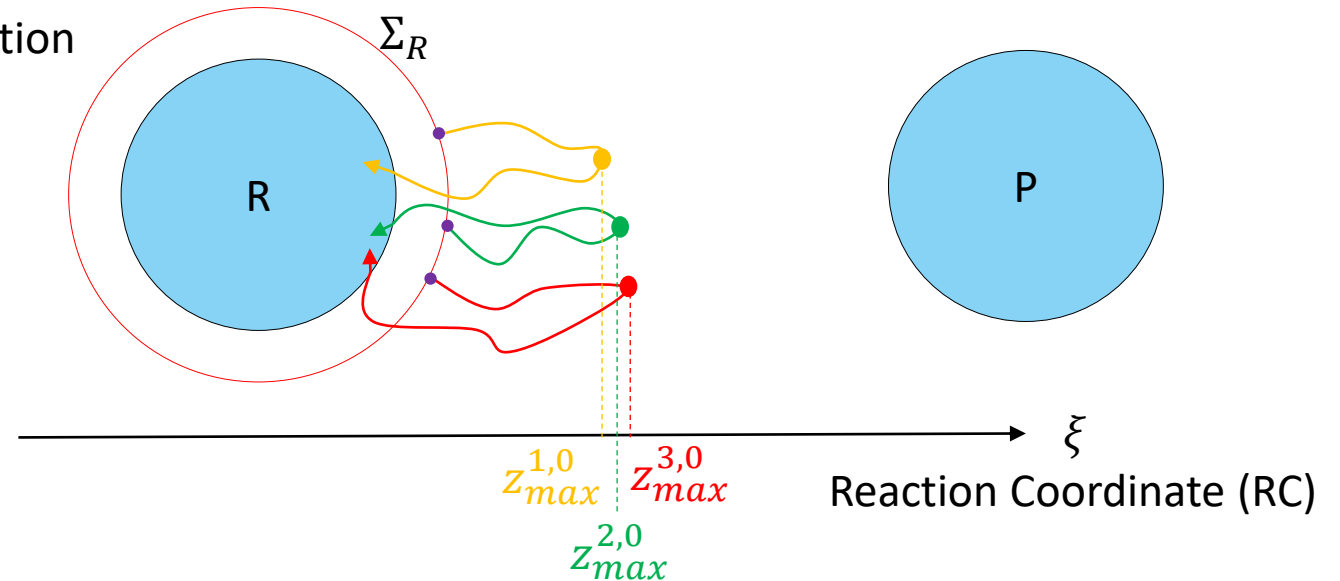
1. Generating initial conditions on  $\Sigma$  and estimate  $t_{R-\Sigma-R} = \frac{1}{\phi_R}$
2. Initialize N replicas by running an unbiased dynamics until it reaches R or P. Set  $p = 1$ . Classify all the replicas by increasing  $\xi_{\max}$ .
3. Apply the AMS loop until all replicas have reached P.

### 1. Initial conditions and flux



MD for Initial conditions.

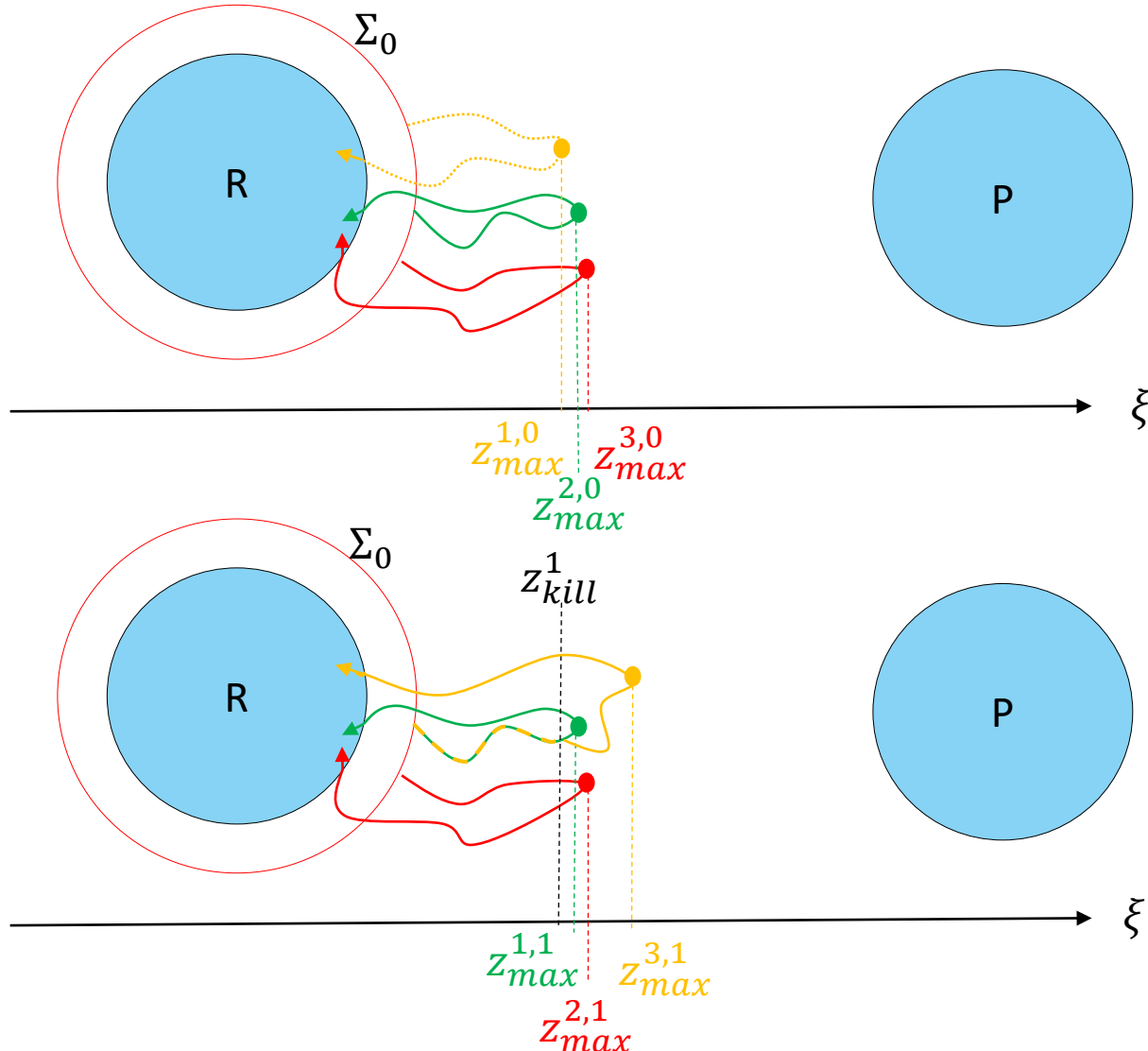
### 2. Initialization



<sup>1</sup> F. Cérou, A. Guyader, (2007) Stoch. Anal. and Appl. 25, 2, 417.

<sup>2</sup> L. J. S. Lopes, T. Lelièvre, (2019) J. Comput. Chem. 40, 1198

# I. ADAPTIVE MULTI-LEVEL SPLITTING METHOD FOR REACTION RATES



## 3. AMS iterations: $i \geq 0$

a) Save the smallest ( $z_{max}^{1,i}$ ) as  $z_{kill}^{i+1}$  and delete all the trajectories that did not “go above”  $z_{kill}^{i+1}$

b) Randomly select one trajectory within the remaining ones. Copy it until it reaches  $z_{kill}^{i+1}$  and continue it until it reaches R or P.

c) Classify all the replicas by increasing  $z_{max}$ .

$$\tilde{p} = \prod_{i=0}^{i_{max}} \tilde{p}_{\Sigma_{z_{kill}^i} \rightarrow \Sigma_{z_{kill}^{i+1}}} = \left(1 - \frac{1}{N}\right)^{i_{max}}$$

Unbiased estimator:  $\mathbb{E}[\tilde{p}] = p_{R-P}(\Sigma_R)$

Variance depends on RC:

$$\text{Var}[\tilde{p}] = f(\xi)$$

Implemented with VASP software<sup>1,2</sup>

<sup>1</sup> G. Kresse, J. Hafner, (1993) J. Phys. Rev. B, 47, 558–561.

<sup>2</sup> G. Kresse, D. Joubert, (1999) Phys. Rev. B, 59, 1758–1775.

## II. AIMD METHOD APPLIED TO WATER DISSOCIATION ON (100) SURFACE

### Multistate problem

With

$$R = A_1$$

$$\Sigma_R = \Sigma_{A_1}$$

$$P = A_2A_3 \cup A_4 \cup D_1D_3 \cup D_2D_4$$

AMS can sample :

$$A_1 \rightarrow A_2A_3$$

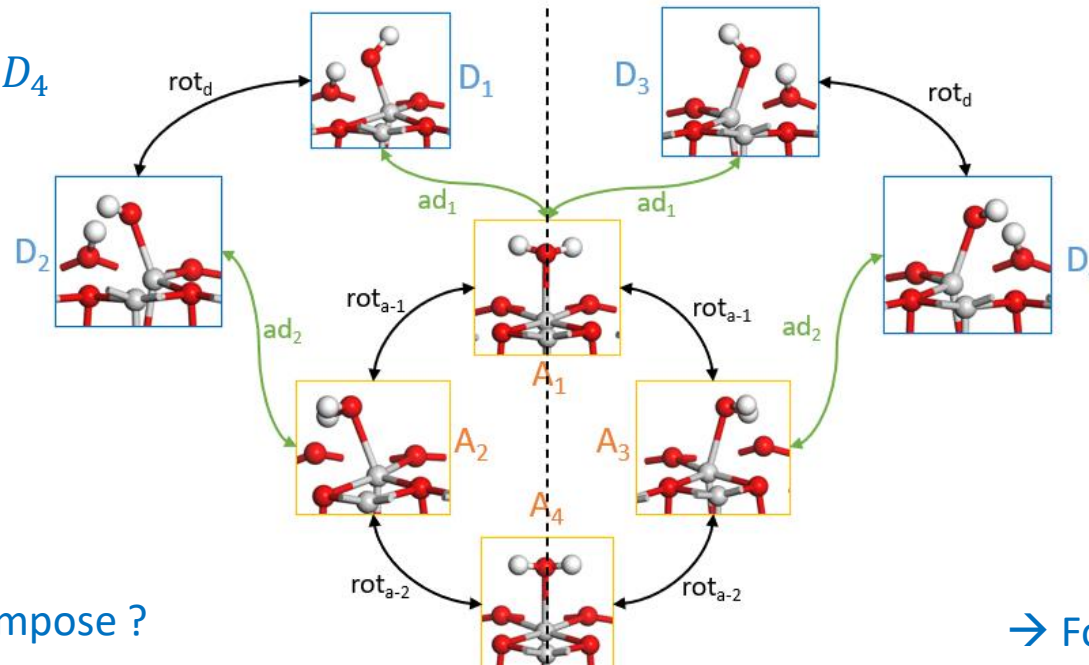
$$A_1 \rightarrow A_4$$

$$A_1 \rightarrow D_1D_3$$

$$A_1 \rightarrow D_2D_4$$

→ Answers how  $A_1$  can decompose ?

The most probable transition will be sampled, with precision conditioned by  $\xi$



Metastable states of  $H_2O$  on the (100) surface of  $\gamma$ -alumina

With

$$R = A_1 \cup A_2A_3 \cup A_4 \cup D_2D_4$$

$$\Sigma_R = \Sigma_{A_1}$$

$$P = D_1D_3$$

AMS can sample :

$$A_1 \rightarrow D_1D_3$$

→ Focus specifically on one event

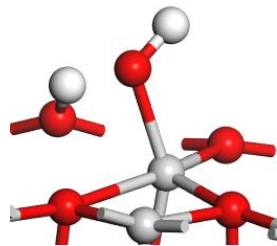
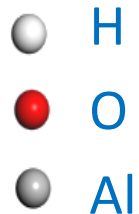
Quality of the sampling depends on  $\xi$



## II. AIMD METHOD APPLIED TO WATER DISSOCIATION ON (100) SURFACE

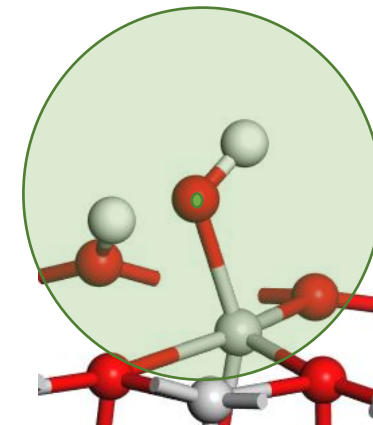
Method:

1. Identify the various metastable states (intermediates)  
→ dissociated ( $D_i$ ) or associated ( $A_i$ )
2. Run short dynamics in these states to sample Potential Energy Surface (PES) around the minima
3. SOAP<sup>1</sup> atom centered descriptors to numerically encode the structure for training the MLCV.



Cartesian coordinates  
(vector of  $3N$  lines)

Select central atom



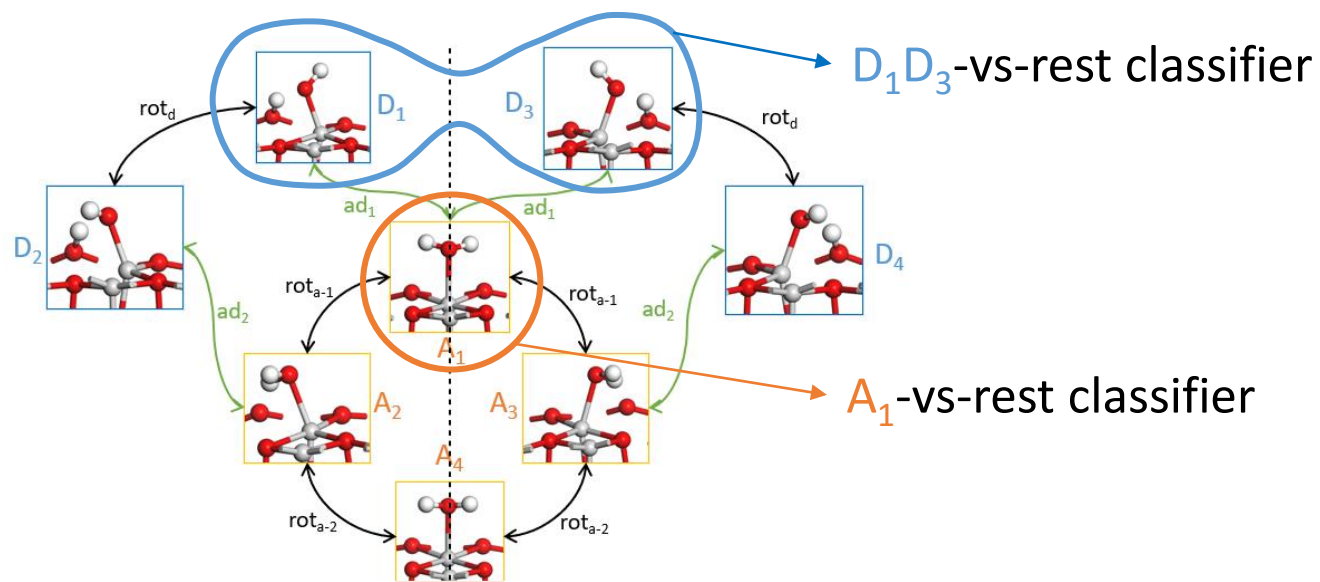
Atom centered descriptor of the structure  
(vector of  $\sim 10^3$ - $10^4$  lines)

<sup>1</sup> A. P. Bartók, R. Kondor, .. G. Csányi, (2013) Phys. Rev. B, 87, 18, 184115.

## II. AIMD METHOD APPLIED TO WATER DISSOCIATION ON (100) SURFACE

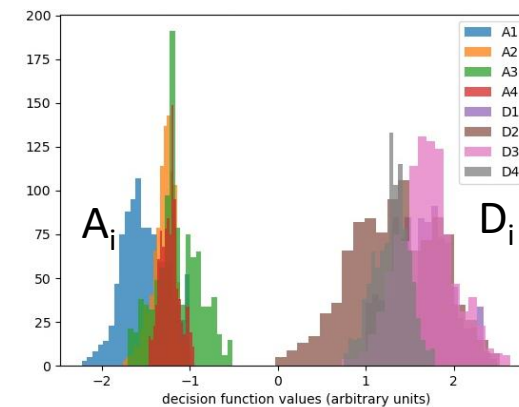
SVM classifiers separate two sets of points by the highest margin plane<sup>1</sup>

SOAP-SVM CV : classifier decision function ( $f_X$ ): algebraic distance to the plane.

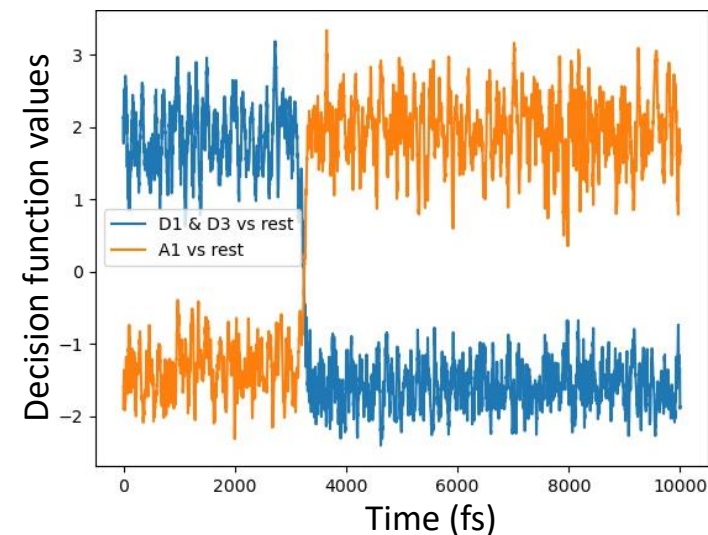


$$f_X(\mathbf{q}) \in (-\infty, -1] \Leftrightarrow \mathbf{q} \in X$$

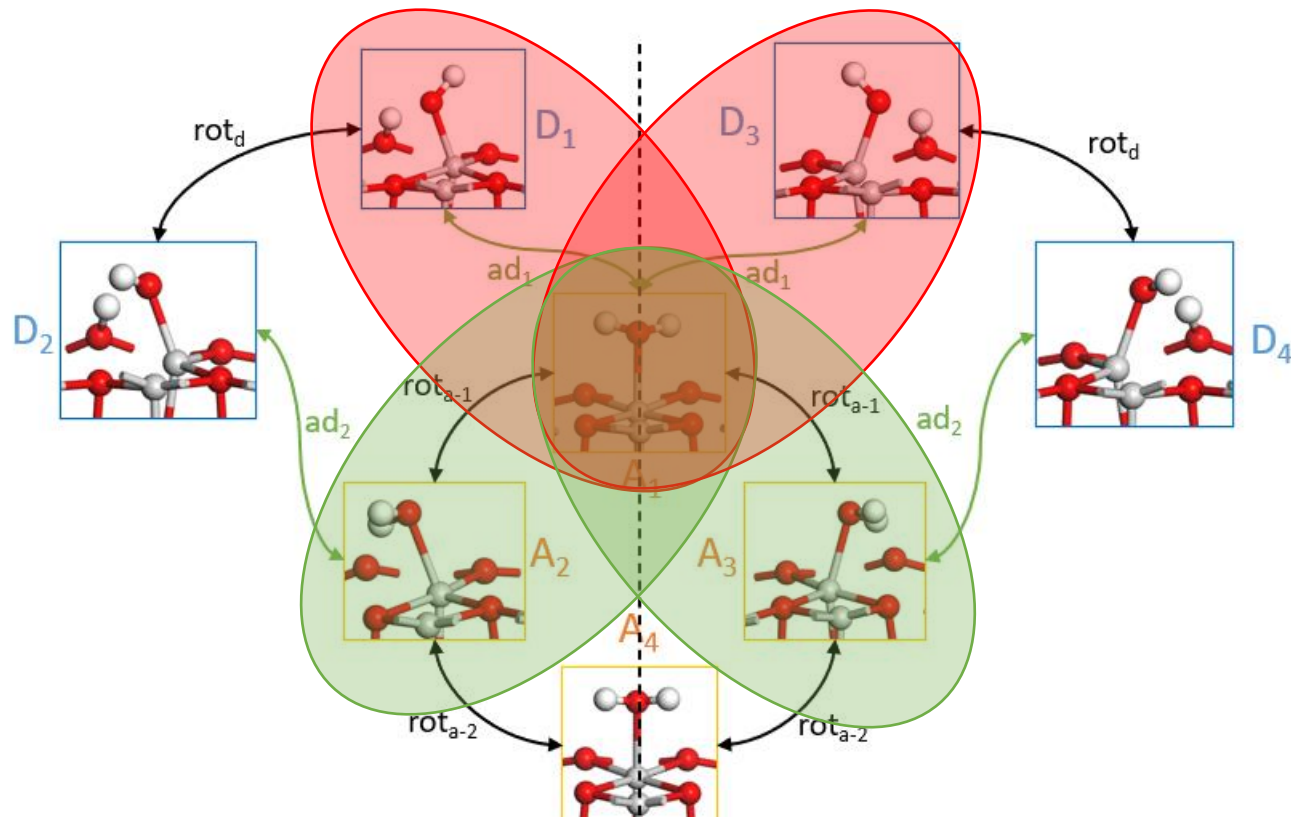
Histogram of SVM classifier decision function values



Training set : A1 = 0, D1 = 1



<sup>1</sup>K. P. Murphy, (2022) Probabilistic Machine Learning: An introduction; MIT Press: , 2022



Dissociation

$$k_{A_1 \rightarrow D_1 D_3} =$$

$$k_{D_1 D_3 \rightarrow A_1} =$$

Hill

$$1.6 \cdot 10^9 \text{ s}^{-1}$$

$$2.3 \cdot 10^{10} \text{ s}^{-1}$$

hTST

$$3.4 \cdot 10^{11} \text{ s}^{-1}$$

$$1.1 \cdot 10^{12} \text{ s}^{-1}$$

Rotation

$$k_{A_1 \rightarrow A_2 A_3} =$$

$$k_{A_2 A_3 \rightarrow A_1} =$$

Hill

$$3.8 \cdot 10^{10} \text{ s}^{-1}$$

$$1.5 \cdot 10^{11} \text{ s}^{-1}$$

hTST

$$7.6 \cdot 10^{10} \text{ s}^{-1}$$

$$2.1 \cdot 10^{12} \text{ s}^{-1}$$

hTST rates are larger

Might come from entropy estimation / recrossing

~ 2 · 10<sup>6</sup> CPU Hours

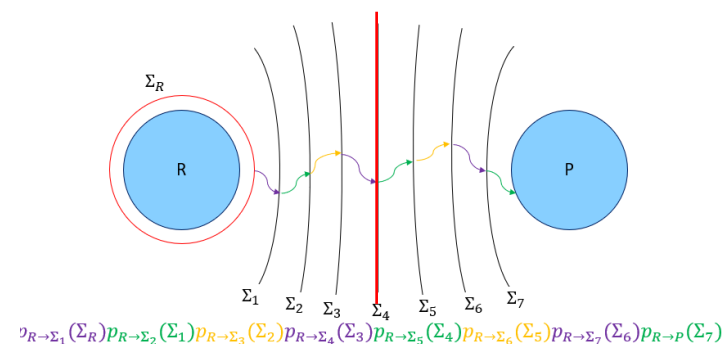
## II. AIMD METHOD APPLIED TO WATER DISSOCIATION ON (100) SURFACE

### Identify TS structures

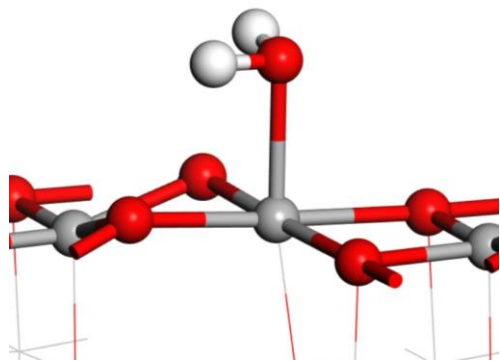
TS in the sense of committor function  $p_{R \rightarrow P}$  (probability of reaching P before R)<sup>1</sup>

Find the level of the RC  $z_{kill}^n$  such that  $p_{R \rightarrow P}(\Sigma_{z_{kill}^n}) = 0.5$

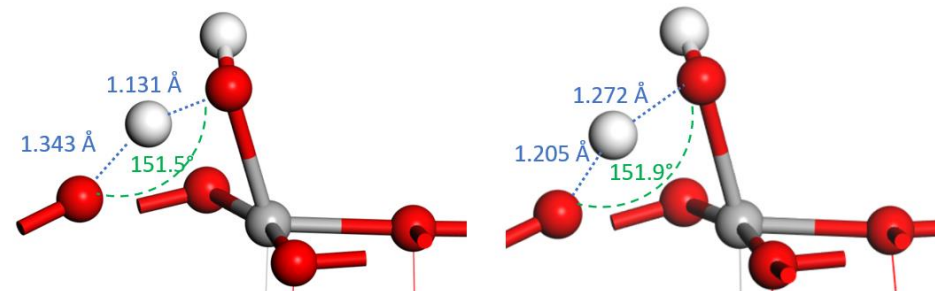
$$\prod_{i=n}^{i_{max}} \tilde{p}_{R \rightarrow \Sigma_{z_{kill}^{i+1}}}(\Sigma_{z_{kill}^i}) = 0.5$$



Along each trajectory, take the structure right after the level  $\Sigma_{z_{kill}^n}$  is crossed, then find the average structure



Example for the  $A_1 \rightarrow D_1 D_3$  reaction



Saddle point

AMS estimated

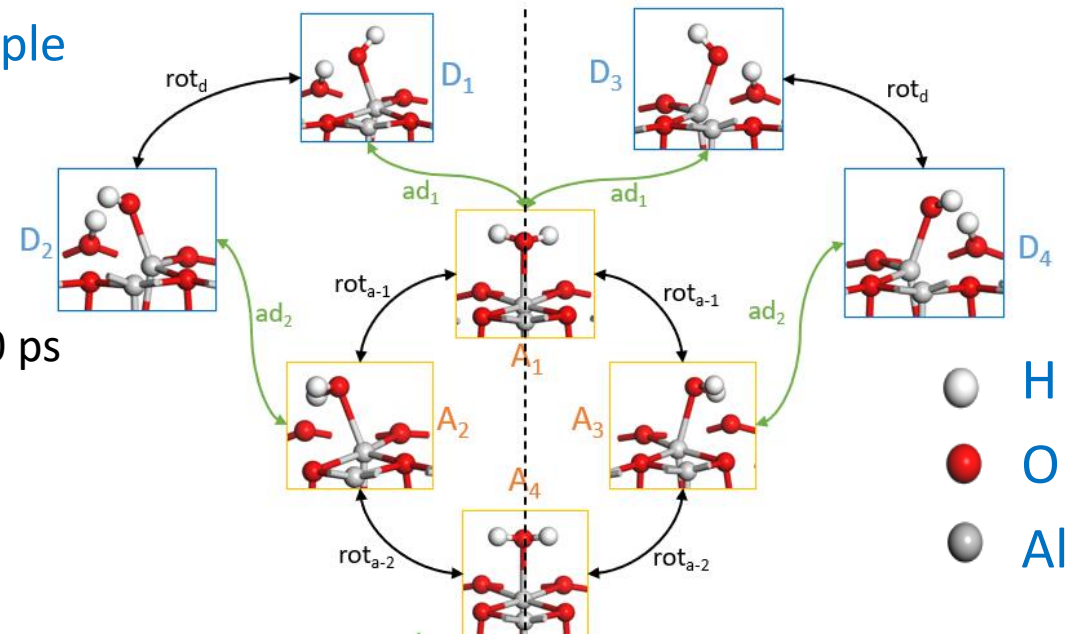
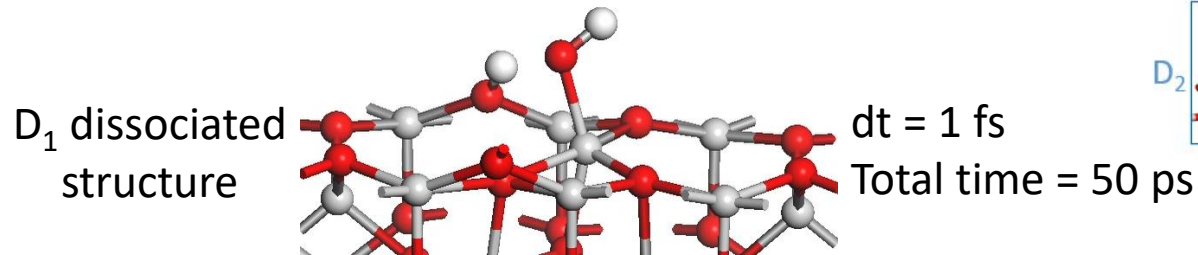
$$p_{A_1 \rightarrow D_1 D_3} = 0.5$$

<sup>1</sup>E. Vanden-Eijnden, E. Transition Path Theory (2006) in Computer Simulations in Condensed Matter Systems: From Materials to Chemical Biology. p 453-493

# III. REDUCING AMS COMPUTATIONAL COST WITH MACHINE LEARNING FORCE FIELD

Method to train a MLFF<sup>1,2</sup>:

1. Identify the various metastable states (intermediates)  
→ dissociated ( $D_i$ ) or associated ( $A_i$ )
2. Run active learning dynamics in these states to sample Potential Energy Surface (PES) around the minima



3. Concatenate the dataset (containing E, Forces, and positions)
  4. AMS with Active learning
  5. Re-fit the force field
- Rate can be estimated with AMS with MLFF without active learning

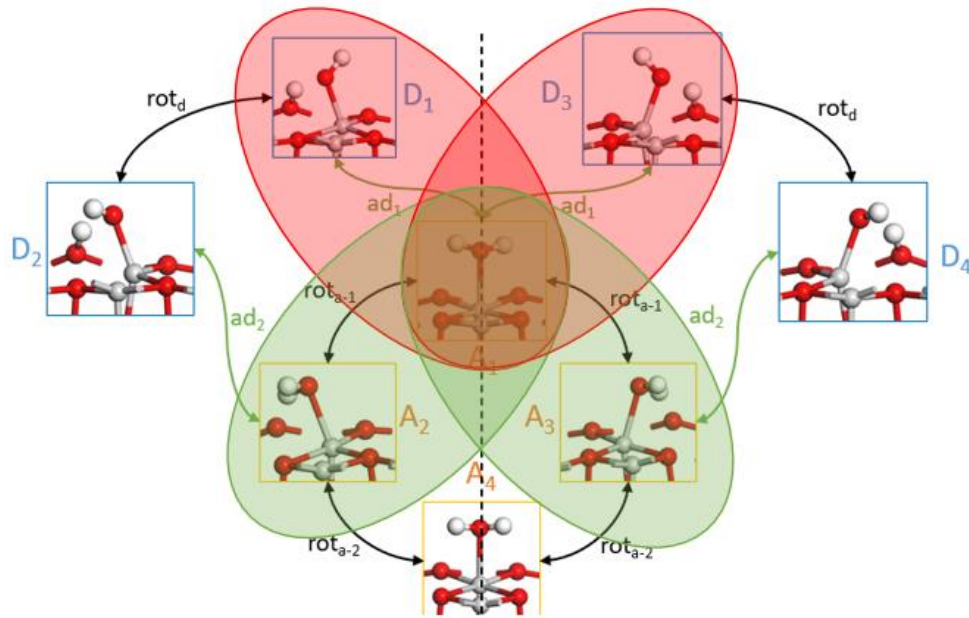
<sup>1</sup> R. Jinnouchi, F. Karsai, G. Kresse, G. (2019) Phys. Rev. B, 100, 014105.

<sup>2</sup> R. Jinnouchi, K. Miwa, F. Karsai, G. Kresse, R. Asahi, (2020) J. Phys. Chem. Lett., 11, 6946–6955.

# III. REDUCING AMS COMPUTATIONAL COST WITH MACHINE LEARNING FORCE FIELD

with  $N_{\text{rep}} = 200$  and  $M_{\text{real}} = 10$

$A_1 \rightarrow D_1 D_3$  Dissociation



DFT

$$k_{A_1 \rightarrow D_1 D_3} = (1.64 \pm 1.59) 10^9 \text{ s}^{-1}$$

MLFF

$$k_{A_1 \rightarrow D_1 D_3} = (2.76 \pm 3.81) 10^9 \text{ s}^{-1}$$

with  $N_{\text{rep}} = 800$  and  $M_{\text{real}} = 10$

$$k_{A_1 \rightarrow D_1 D_3} = (2.43 \pm 1.15) 10^9 \text{ s}^{-1}$$

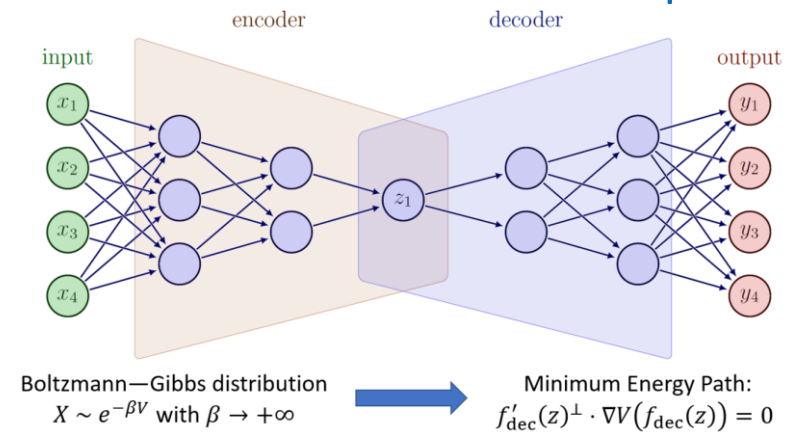
Test MAE forces = 76 meV/Å

Test MAE energy = 20 meV

1000 configurations randomly drawn  
from sampled reactive trajectories

# 15 CONCLUSION AND PERSPECTIVES

- hTST overestimate the DFT-MD rate estimated using AMS
- MLFF-MD and DFT-MD rates are consistent
- MLFF used in prediction mode drastically reduces de computational cost
- Current implementation of AMS with VASP limits the application of active learning  
→ Restart does have an important cost for the active learning.
- Using D-optimality criterion active learning with VASP as calculator of ab-initio with ACE potential seems a good opportunity<sup>1</sup>
- Active learning of RC  $\xi$  can be included in the workflow<sup>2</sup>



<sup>1</sup> Y. Lysogorskiy, A. Bochkarev, M. Mrovec, R. Drautz, (2023) Phys. Rev. Mater., 7, 4, 043801

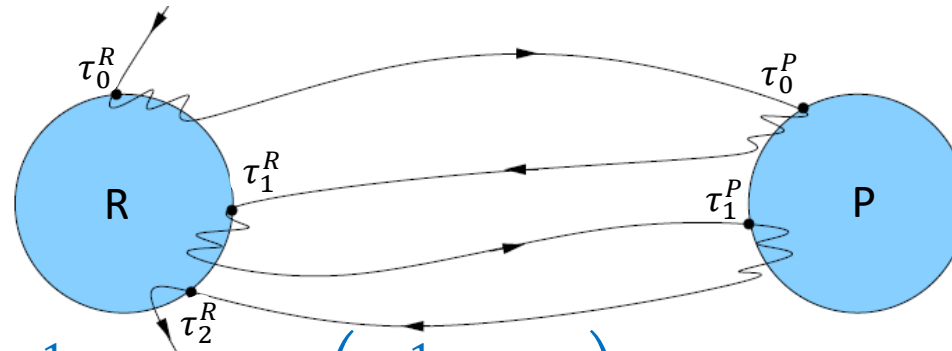
<sup>2</sup> T. Lelièvre, TP, G. Stoltz, W. Zhang, (2024) J. Phys. Chem. B, 128, 11, 2607

Thank you for you attention



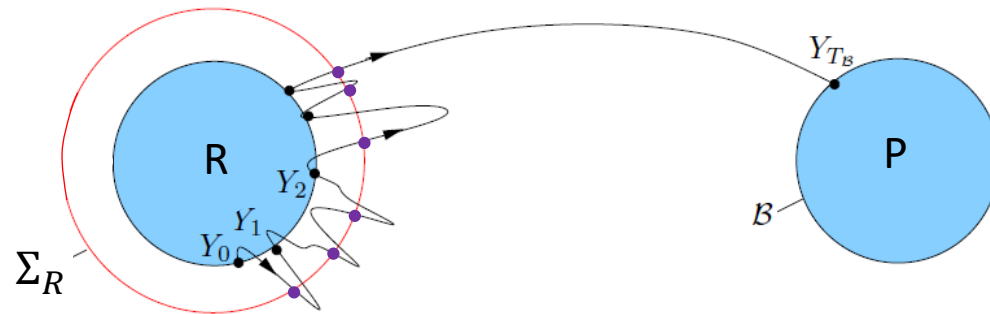
# I. ADAPTIVE MULTI-LEVEL SPLITTING METHOD FOR REACTION RATES

Transition time:  $\frac{1}{k_{RP}} = t_{RP} = \text{mean}(\tau_i^R - \tau_i^P)$



We model the reaction time as:  $\frac{1}{k_{RP}} = t_{RP} = \left( \frac{1}{p_{\Sigma_{R \rightarrow P}}} - 1 \right) (t_{R \rightarrow \Sigma} + t_{\Sigma \rightarrow R}) + t_{R \rightarrow \Sigma}^\dagger + t_{\Sigma \rightarrow P} \approx \frac{t_{R \rightarrow \Sigma \rightarrow R}}{p_{\Sigma_{R \rightarrow P}}} =$

$$\frac{1}{p_{\Sigma_{R \rightarrow P}} \phi_R}$$



$p_{\Sigma_{R \rightarrow P}}$ : probability of reaching P before R when starting from  $\Sigma_R$ .

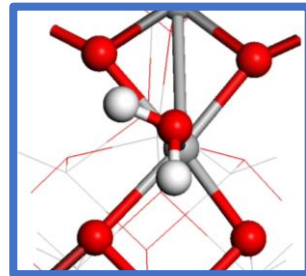
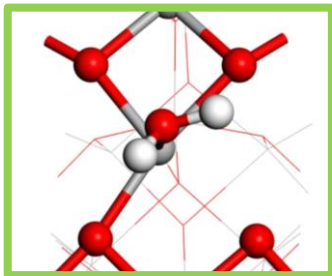
<sup>1</sup> Baudel, M., Guyader, A., & Lelièvre, T. (2020). On the Hill relation and the mean reaction time for metastable processes. *arXiv preprint, arXiv:2008.09790*.

## II. AIMD METHOD APPLIED TO WATER DISSOCIATION ON (100) SURFACE

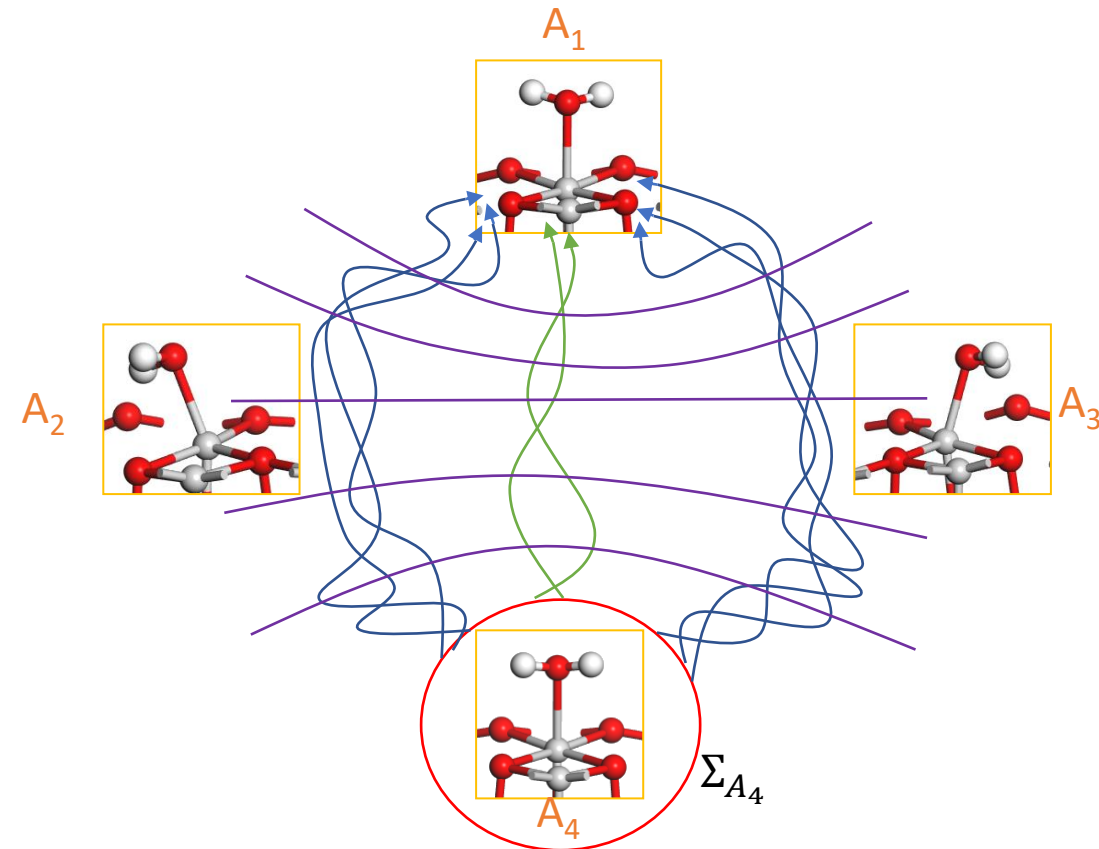
Use K-means clustering method to identify groups of trajectories.

Based on SOAP descriptor + PCA to describe 5 structures per trajectory.

5 Structures = First time trajectory cross RC iso-levels



- Reactive trajectories
- Iso-levels of a reaction coordinate



# III. REDUCING AMS COMPUTATIONAL COST WITH MACHINE LEARNING FORCE FIELD

Wall clock time	Total number of structures	Methodology	Final force mean Bayesian error
~0.5 days	2600	50 ps active learning in each metastable states $A_i$ and $D_i$ starting from scratch	10 to 60 meV/Å
~0.5 days	4000	50 ps active learning starting in $A_1$ with the merged datasets	25 meV/Å
~10 days	4010	10 $A_1$ -vs-all AMS runs using active learning	25 meV/Å

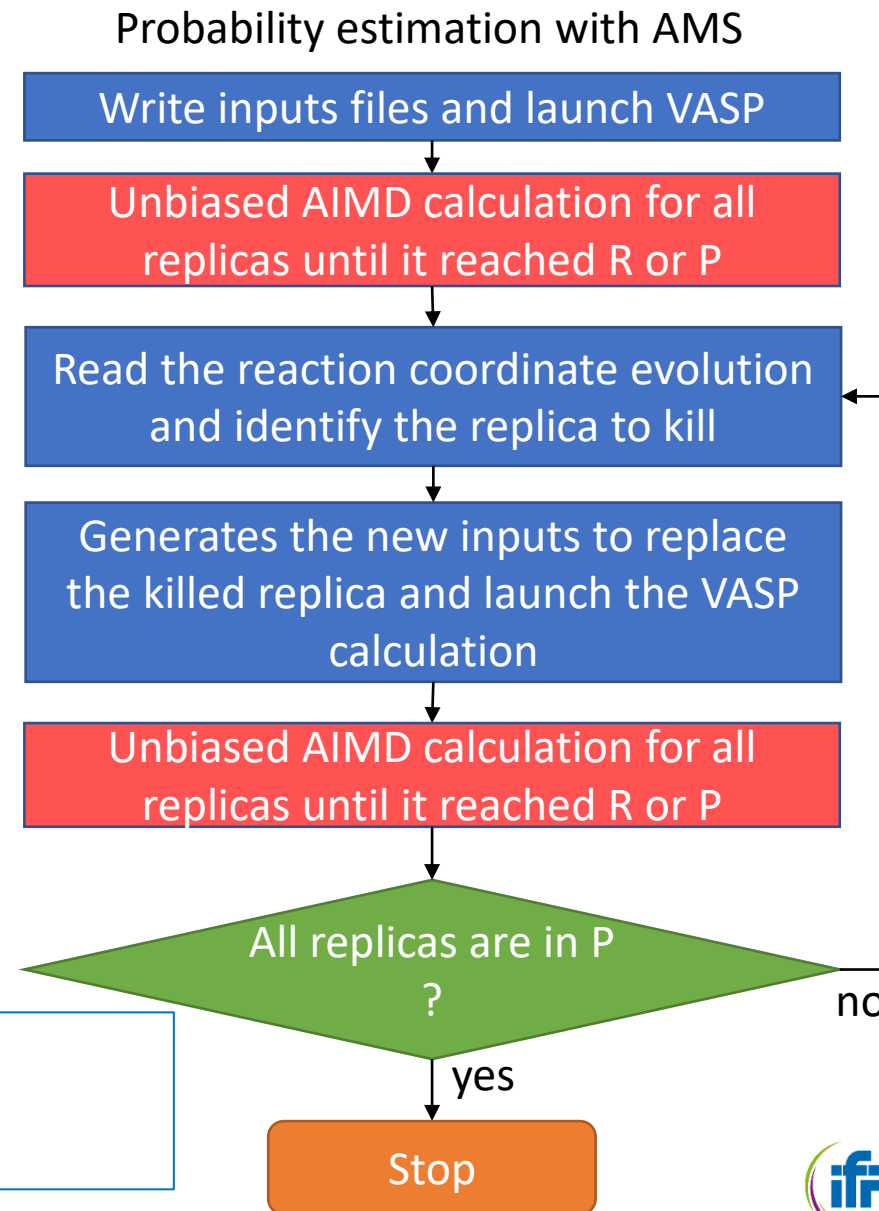
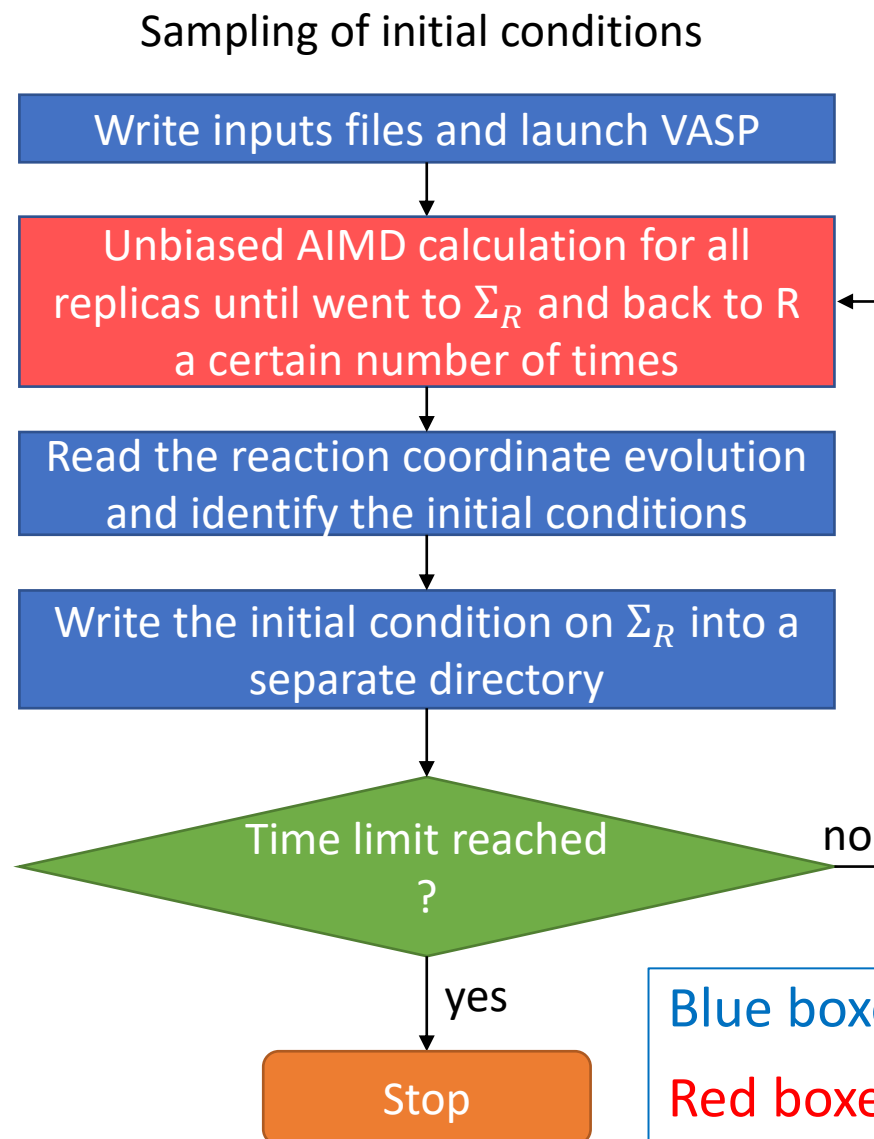
Threshold updated using the stored Bayesian errors<sup>1,2</sup>

→ Transitions already sampled during first two steps?

<sup>1</sup> Jinnouchi, R.; Karsai, F.; Kresse, G. (2019) *Phys. Rev. B* 100, 014105

<sup>2</sup> Jinnouchi, R.; Miwa, R.; Karsai, F.; Kresse, G.; Asahi, R. (2020) *J Phys. Chem. Lett.* 11, 6946

# I. ADAPTIVE MULTI-LEVEL SPLITTING METHOD FOR REACTION RATES



Blue boxes: python code

Red boxes: VASP code<sup>1,2</sup>

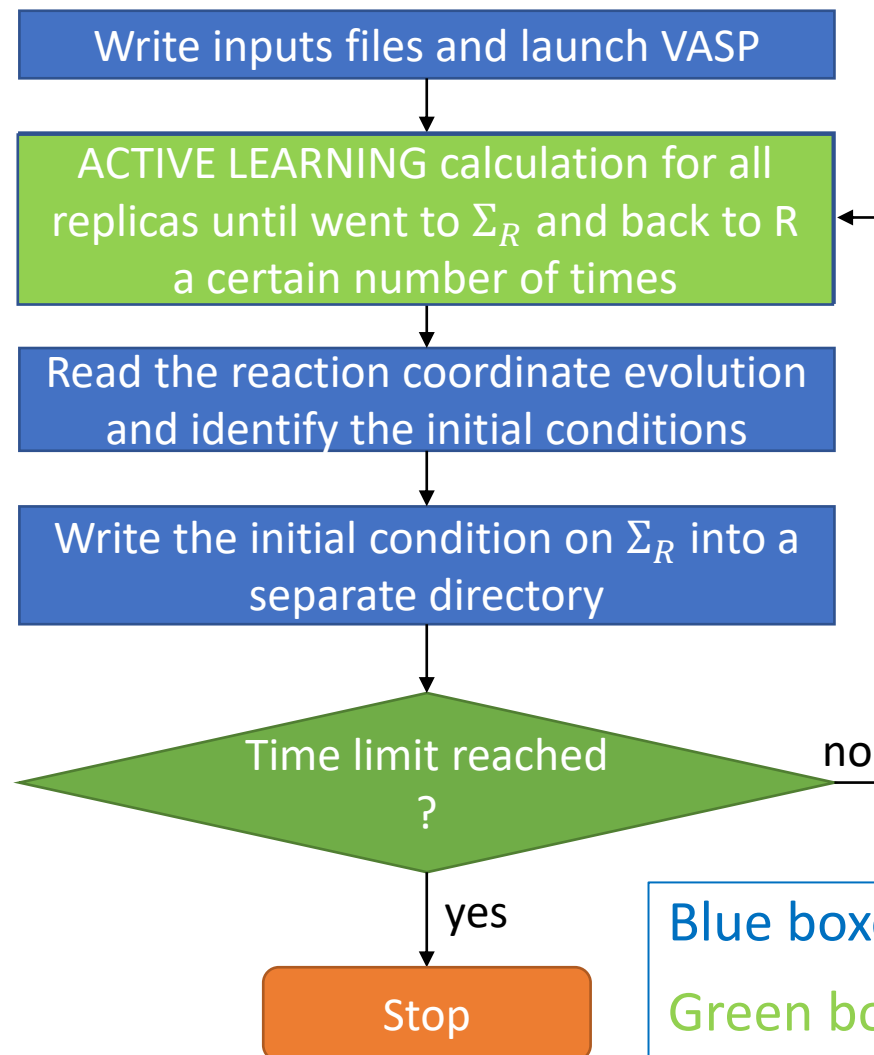
<sup>1</sup>Kresse, G.; Hafner, J. *Phys. Rev. B* **1993**, 47, 558–561.

<sup>2</sup>Kresse, G.; Joubert, D. *Phys. Rev. B* **1999**, 59, 1758–1775.

### III. USING AMS WITH ACTIVE LEARNING

#### AMS IMPLEMENTATION WITH VASP (PLANE WAVE DFT)

##### Sampling of initial conditions



Blue boxes: python code

Green boxes: VASP active learning

##### Probability estimation with AMS

