

# Gold standard finite temperature simulations of materials via machine learning

---

**Basile Herzog**<sup>1</sup>, Alejandro Gallo<sup>2</sup>, Mauricio Chagas da Silva<sup>1</sup>, Andreas Irmeler<sup>2</sup>, Felix Hummel<sup>2</sup>, Michael Badawi<sup>1</sup>, Tomas Bucko<sup>3</sup>, Sébastien Lebegue<sup>1</sup>, Andreas Grüneis<sup>2</sup> and Dario Rocca<sup>1</sup>

<sup>1</sup> LPCT, Université de Lorraine & CNRS, Nancy (France)

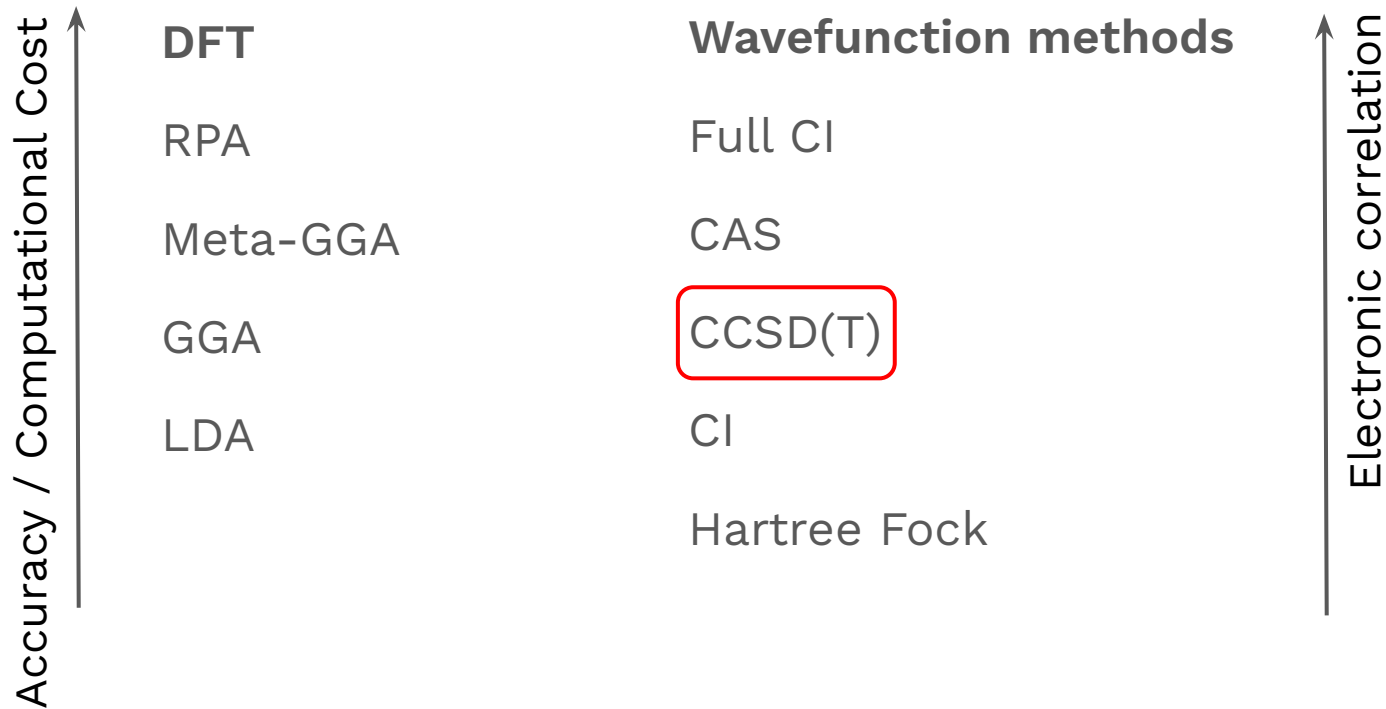
<sup>2</sup> Institute for Theoretical Physics, TU Wien, Vienna (Austria)

<sup>3</sup> Comenius University in Bratislava and Slovak Academy of Sciences, Bratislava (Slovakia)

# Outline

- Introduction & Motivation
- Machine Learning Perturbation Theory (MLPT)
- Adsorption enthalpy of CO<sub>2</sub> in Protonated Chabazite
- Possible limitations and solution: Machine Learning Monte-Carlo (MLMC)
- Conclusion

# Methods vs computational cost



# Coupled Cluster Theory

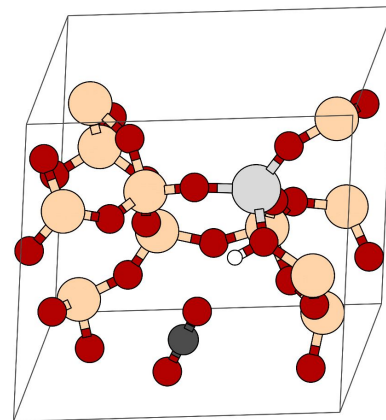
- CCSD:  $|\Psi\rangle = e^{T_1+T_2}|HF\rangle$
- $T_1, T_2$ : single and double substitutions  
triples are treated perturbatively
- exponential ansatz: allows systematic inclusion of highest degree of correlations (Taylor expansion)
- **Gold standard** method in the quantum chemistry community: reaches chemical accuracy for most applications
- $O(N^6)$  complexity

# Case study: adsorption enthalpy of CO<sub>2</sub> in protonated chabazite

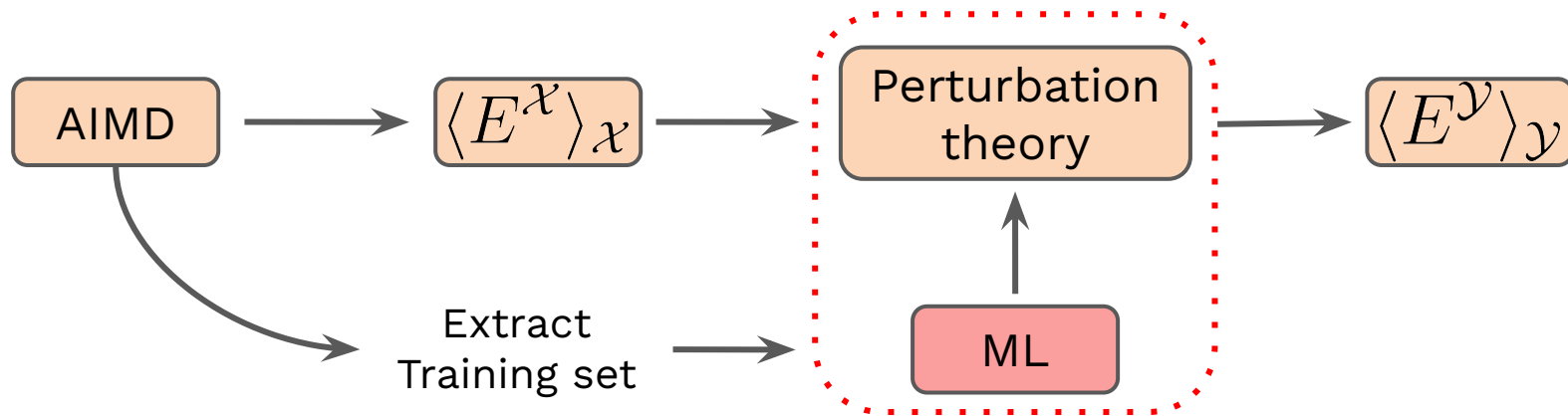
$$\Delta H = \langle E(\text{CO}_2@HCHAB) \rangle - \langle E(HCHAB) \rangle - \langle E(\text{CO}_2) \rangle - k_B T$$

- *Ab-initio* Molecular Dynamics (AIMD) - 200K steps  
-CCSD(T) → 10B CPU Hours → 1000 human years

- Machine learning thermodynamic perturbation theory (MLPT)  
a “cheap” AIMD (e.g. PBE) + only a few number (~100) of CCSD(T)  
calculations → 1 human week



# Machine Learning Thermodynamic Perturbation Theory



1. AIMD is performed using the cheap theory  $\mathcal{X} \rightarrow \langle E^{\mathcal{X}} \rangle_{\mathcal{X}}$
2.  $\Delta E_i = E_i^{\mathcal{X}} - E_i^{\mathcal{Y}}$  is learned on  $N_{\text{train}} \sim 100$  configurations evenly spaced from the MD
3. Perturbation theory is applied to  $\{E_i^{\mathcal{X}}\} \rightarrow \langle E^{\mathcal{Y}} \rangle_{\mathcal{Y}}$

# Thermodynamic Perturbation Theory

- Computationally cheap Hamiltonian (e.g. PBE)

$$H^x = T + V^x$$

- Average in canonical ensemble of  $H^x$

$$\langle O \rangle_x$$

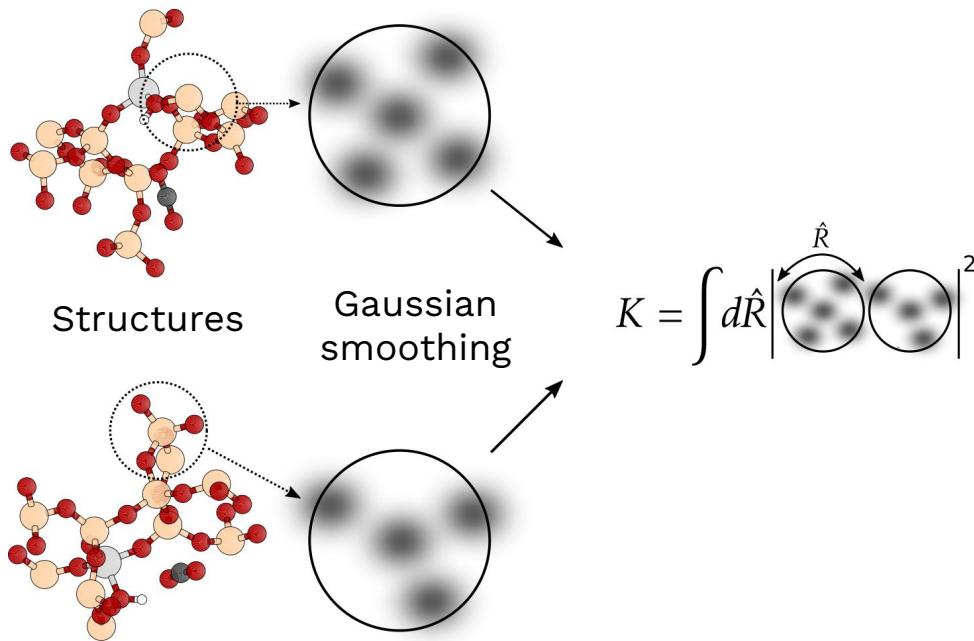
- Computationally expensive Hamiltonian (e.g. CCSD(T))

$$\begin{aligned} H^y &= H^x + V^y - V^x \\ &= H^x - \Delta V \end{aligned}$$

- Average in canonical ensemble of  $H^y$

$$\langle O \rangle_y = \frac{\langle O e^{\beta \Delta V} \rangle_x}{\langle e^{\beta \Delta V} \rangle_x}$$

# Machine Learning model : Descriptor & KRR



- Smooth Overlap of Atomic Positions (SOAP)

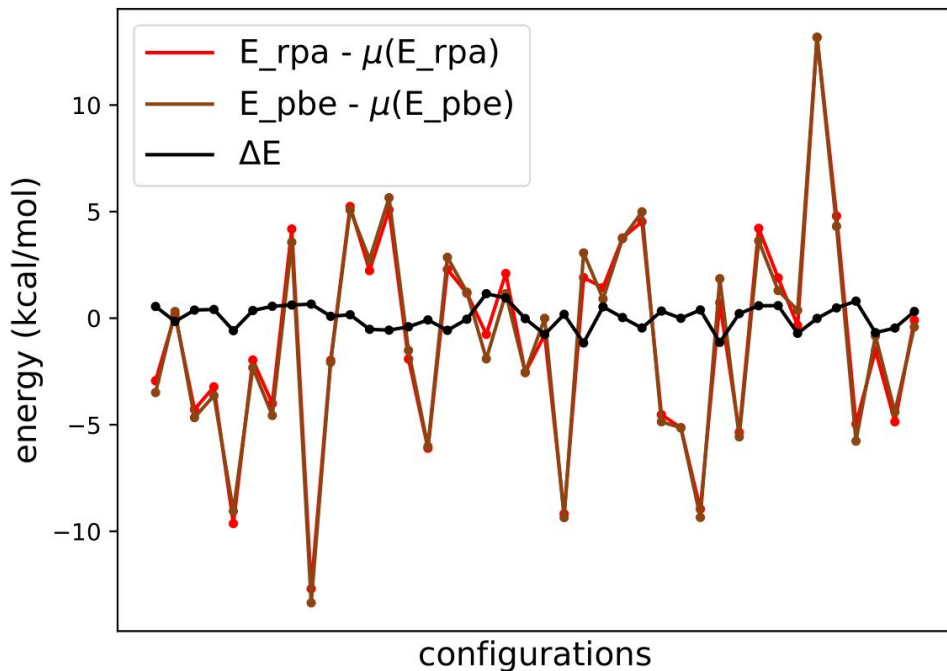
- Kernel Ridge Regression (KRR)

$$\alpha = (K_{train} + \lambda I)^{-1} y_{train}$$

$$y_{pred} = K_{pred} \alpha$$



# Machine Learning model : $\Delta$ -Machine Learning



$\Delta E$ : smooth function, easy to learn

Small training set:  
Ntrain=100/Nval=10

System	HCHAB	CO2@HCHAB
RMSE (kcal/mol)	0.50	0.63

## Results : CO<sub>2</sub>@HCHAB

Method	PBE-D2	MP2	CCSD(T)	Exp.
Enthalpy (kcal/mol)	-9.72	-9.50	-7.69	-8.41

- Good agreement with experiment, within chemical accuracy
- First finite temperature result at this level of theory

# Possible sources of error

$$\langle E^{\mathcal{Y}} \rangle_{\mathcal{Y}} = \frac{\langle E^{\mathcal{Y}} e^{\beta \Delta E} \rangle_{\mathcal{X}}}{\langle e^{\beta \Delta E} \rangle_{\mathcal{X}}} = \sum_i^M \frac{w_i E_i^{\mathcal{Y}}}{\sum_j w_j}$$

$$w_i = \exp(\beta \Delta E_i)$$

Production and target methods might have insufficient statistical overlap: averages in target space would be dominated by few configurations.

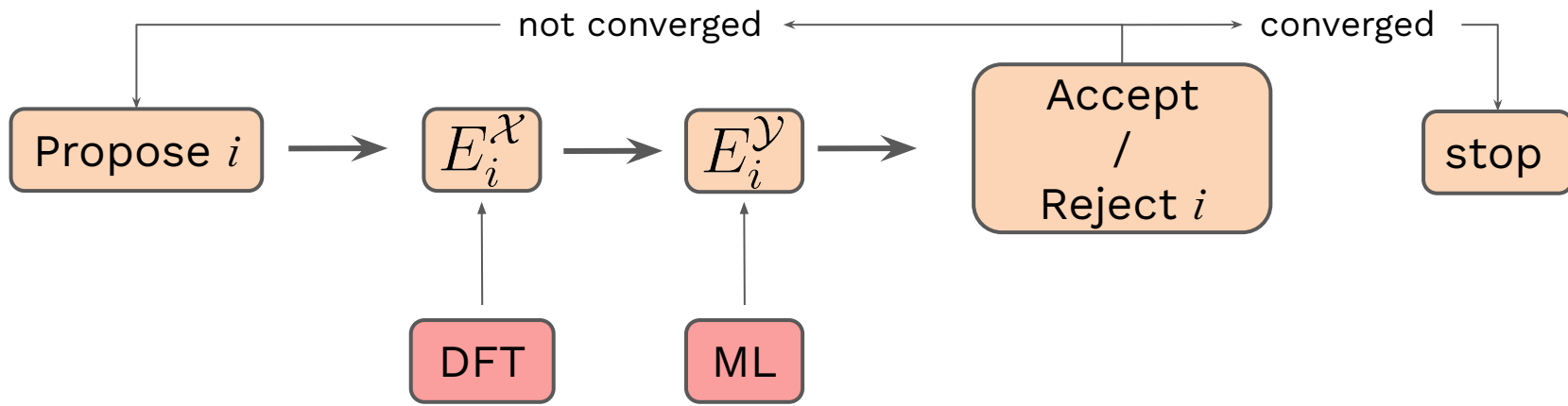
$$I_w = \frac{(M - N)}{M} \in [0, 0.5]$$

$$N, M \text{ such that } \frac{\sum_i^N w_i}{\sum_j^M w_j} \geq 0.5$$

$I_w = 0$ : MLPT will fail

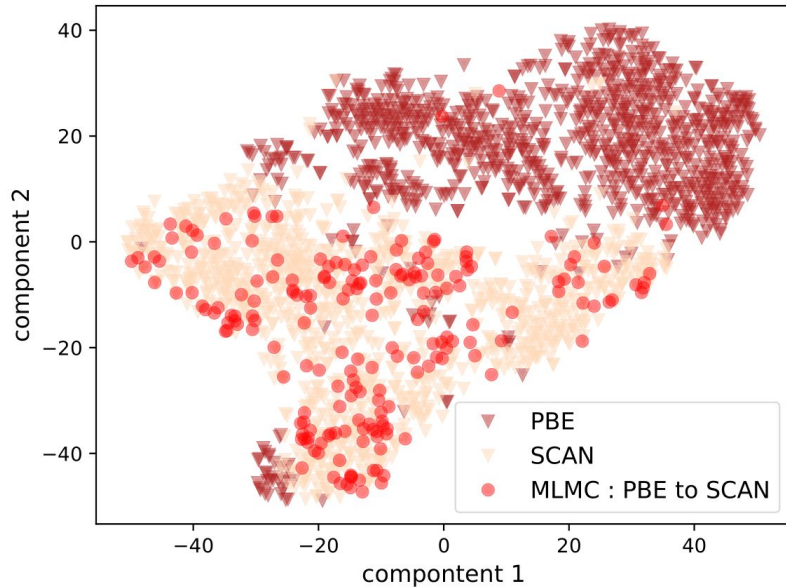
Solution: Machine Learning Monte-Carlo (MLMC)

# Machine Learning Monte-Carlo



1. Energy of the proposed configuration is computed using the cheap theory
2. Correction to the expensive theory is done using the ML model of the previous MLPT
3. Sampling is done in the expensive configurational space until convergence

# Example: MLPT estimate of SCAN energy from PBE



$\Delta E(\text{MLPT})$

$\Delta E(\text{MLMC})$

-4.38

0.64

- $\text{CH}_4$  in protonated chabazite
- Case study with DFT functionals SCAN & PBE. MD reference is known
- MLPT  $\rightarrow I_w=0$ , bad overlap, big deviation compared to reference
- MLMC trajectory lies in correct configurational space, good agreement with reference

# Conclusion

- The MLPT method allows computation of highly accurate thermodynamic property at CCSD(T) level for the first time
- Monte-Carlo resampling is currently running to confirm this result
- Future directions: other porous materials applications, surface adsorption problems, activation energies

# Acknowledgements

Group leaders: Dario Rocca<sup>\*</sup>,  
Sébastien Lebegue & Andreas  
Grüneis



Michael Badawi for the porous  
materials applications and  
computational time



Collaborators: Mauricio Chagas da Silva, Alejandro Gallo,  
Andreas Irmler, Felix Hummel



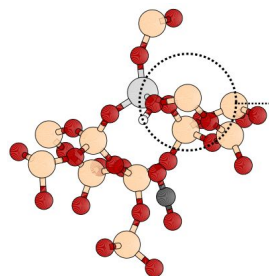
Tomáš Bučko for the fruitful discussions and  
collaboration that led to the development of MLPT



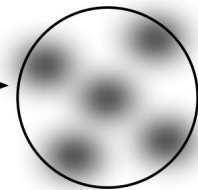
<sup>\*</sup> [dario.rocca@univ-lorraine.fr](mailto:dario.rocca@univ-lorraine.fr)

# SOAP Kernel

$$\rho^{\chi_a}(\mathbf{r}) = \sum_{i \in \chi_a} \exp(-|\mathbf{x}_i - \mathbf{r}|^2 / 2\sigma^2)$$



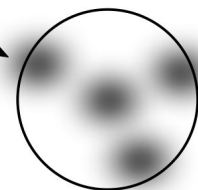
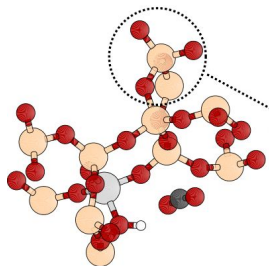
Structures



Gaussian Smoothing

$$K = \int d\hat{\mathbf{R}} \left| \int d\mathbf{r} \rho^{\chi_a}(\mathbf{r}) \rho^{\chi_b}(\hat{\mathbf{R}}\mathbf{r}) \right|^2$$

$$K(\chi_a, \chi'_b) = \int d\hat{\mathbf{R}} \left| \int d\mathbf{r} \rho^{\chi_a}(\mathbf{r}) \rho^{\chi'_b}(\hat{\mathbf{R}}\mathbf{r}) \right|^2$$



$$K(\chi, \chi') = \sum_{a=1}^{N_\chi} \sum_{b=1}^{N_{\chi'}} \frac{K(\chi_a, \chi'_b)}{\sqrt{K(\chi_a, \chi_a) K(\chi'_b, \chi'_b)}}$$

$$\alpha = (K_{train} + \lambda I)^{-1} y_{train}$$

$$y_{pred} = K_{pred} \alpha$$



# Electronic Correlation

Chemical bonds

Non covalent interactions

Free energy profiles

Strong correlations ?