

# AIDD Spring School: Advanced Machine Learning for Drug Discovery Constructing accurate machine learning force fields for flexible molecules

Leonardo Medrano Sandonas

Department of Physics and Materials Science, University of Luxembourg, Luxembourg

May 10, 2022



O.T. Unke et al., Chem. Rev. 121, 10142-10186, (2021).

Email: leonardo.medrano@uni.lu

# The role of system dimensionality **Potential energy surfaces (PES)**

Describe the energy of a system (molecule) in terms of certain parameters (positions, bonds,





Constructing accurate ML force fields for flexible molecules AIDD Spring School Lugano-Switzerland // 10.05.2022

Protein (secondary

# Potential energy surface (PES): what's good for?



M.A. M.-Drumel et al., Phys. Chem. Chem. Phys. 16, 22062, (2014).

**Rotational energy profiles** 



Y. Ali, Sci. Rep. 10, 10995, (2020).

#### Thermodynamic binding constants



H. Isla et al., J. Org. Chem. 84, 5383, (2019).



# **Potential energy surface (PES)**

"Original" energy function for molecules: Molecular mechanics (MM)



 $U_{\mathrm{bond}}$  = Oscillations about the equilibrium bond length

= Oscillations of 3 atoms about an equilibrium bond angle

 $U_{\rm dihedral}$  = Torsional rotation of 4 atoms about a central bond

 $U_{\rm non-bond}$  = Non-bonded energy terms (Lennard-Jones and electrostatics)



 $U_{\text{angle}}$ 

# Potential energy surface (PES)

"Original" energy function for molecules: Molecular mechanics (MM)

Analysis of the conformational sampling of B3 domain of Protein G (GB3)



Conformations and energies highly depend on the chosen force field.

#### F. Martin-Garcia et al., PLoS One **10**, 3, (2015).



## **Potential energy surface (PES)** Hierarchies in atomistic modeling methods



http://quantum-machine.org
J. Hoja *et al., Sci. Data* 8, 43, (2021).
J.S. Smith *et al., Sci. Data* 4, 170193, (2017).
J. Rezac *et al., J. Chem. Theory Comput.* 7, 8, (2011).
Z.M. Sparrow *et al., J. Chem. Phys.* 155, 184303, (2021).



## **Potential energy surface (PES)** Hierarchies in atomistic modeling methods



http://quantum-machine.org
J. Hoja *et al., Sci. Data* 8, 43, (2021).
J.S. Smith *et al., Sci. Data* 4, 170193, (2017).
J. Rezac *et al., J. Chem. Theory Comput.* 7, 8, (2011).
Z.M. Sparrow *et al., J. Chem. Phys.* 155, 184303, (2021).





# QM dataset of small molecules

#### Molecular representations (3D geometric descriptors)

#### Coulomb matrix



Spectrum of London and Axilrod-Teller-Muto potential (SLATM)



Two-body term:

$$rac{1}{2}\sum_{J
eq I}Z_J\delta({f r}-{f R}_{IJ})g({f r})$$

Three-body term:

$$\frac{1}{3}\sum_{J\neq K\neq I}Z_J Z_K \delta(\theta - \theta_{IJK}) h(\theta, \mathbf{R}_{IJ}, \mathbf{R}_{IK})$$

K. Hansen *et al.*, *J. Chem. Theory Comput.* **9**, (2013). W. Pronobis *et al.*, *J. Chem. Theory Comput.* **14**, (2018).

> UNIVERSITÉ DU LUXEMBOURG

Constructing accurate ML force fields for flexible molecules AIDD Spring School Lugano-Switzerland // 10.05.2022

FCHL19 (previous FCHL18) Two-body term:

$$G^{2-\text{body}} = \xi_2(r_{IJ}) f_{\text{cut}}(r_{IJ}) \frac{1}{R_s \sigma(r_{ij}) \sqrt{2\pi}} \exp\left(-\frac{(\ln R_s - \mu(r_{ij}))^2}{2\sigma(r_{ij})^2}\right),$$

Three-body term:

$$G^{3\text{-body}} = \xi_3 G^{3\text{-body}}_{\text{Radial}} G^{3\text{-body}}_{\text{Angular}} f_{\text{cut}}(r_{IJ}) f_{\text{cut}}(r_{JK}) f_{\text{cut}}(r_{KI}).$$

- Neural Network representation (SchNet)
  - Distances are expanded with radial basis functions,

$$e_k(\mathbf{r}_i - \mathbf{r}_j) = \exp(-\gamma \|d_{ij} - \mu_k\|^2)$$



Many-body atomic interactions.

# **QM dataset of small molecules** Molecular representations (3D geometric descriptors)



#### O.A. von Lilenfeld *et al.*, *Nat. Rev. Chem.* **4**, 347, (2020). F. Faber *et al.*, *J. Chem. Phys.* **148**, 241717, (2018).



Constructing accurate ML force fields for flexible molecules AIDD Spring School Lugano-Switzerland // 10.05.2022

□ Atomic forces prediction (MD17 dataset)



Shortcomings for application in **large molecules:** 

- > High-dimensional geometric descriptors.
- ➤ Large degrees of freedom.
- Strong non-covalent interaction.

#### Problems to consider when studying large systems **Conformational sampling** Long range interactions

#### **UniLu-Janssen dataset**

- 60,082 molecular conformations (1673 unique compositions). Elements: H, C, N, O, S, Cl, P, F.
- Structures containnig up to 92 atoms (54 nonhydrogen/heavy atoms).
- ~43 QM properties: PBE0(tight)+many body dispersion (MBD) with MPB implicit solvent





Representations must consider long range terms (vdW, electrostatics).

1.6

0.0

1.5

Conformer identification is more challenging by only considering energies.

-E<sub>MBD</sub> [eV]





# Methods for developing machine learning force fields



# **Method 1: Kernel ridge regression**

#### Isomerization: Azobenzene reaction paths





#### Transition paths: reference data





# **Method 1: Kernel ridge regression**

**Isomerization:** Azobenzene reaction paths

### Training/Testing scheme

- > Training sets from 100 up to 1000 points.
- Subsets of size equal to five times the number of training points.
- 5-fold cross-validation on each subset.
- Used one fold for training and the rest for testing.

### **KRR procedures**



S. Chmiela *et al., Sci. Adv.* **3**, 1603015, (2017).

V. Vassilev-Galindo et al., J. Chem. Phys. 154, 094119, (2021).





# **Method 1: Kernel ridge regression**

**Isomerization:** Azobenzene reaction paths



A single descriptor is not able to optimally resolve all different states on a PES.

#### V. Vassilev-Galindo et al., J. Chem. Phys. 154, 094119, (2021).



# Method 2: First neural network architectures

#### SchNetPack: end-to-end NN with cut-off

> Distances are expanded with radial basis functions,  $e_k(\mathbf{r}_i - \mathbf{r}_j) = \exp(-\gamma ||d_{ij} - \mu_k||^2)$ 



> Many-body atomic interactions.

#### Custom loss function: energies and forces

$\ell(\hat{E}, (E, \mathbf{F}_1, \dots, \mathbf{F}_n)) = \rho \ E - \hat{E}\ ^2 + \frac{1}{n} \sum_{i=0}^n \left\ \mathbf{F}_i - (\mathbf{F}_i)^2 - \hat{E}_i^2 - \hat{E}_i^2\right\ ^2$	$\left(-\frac{\partial \hat{E}}{\partial \mathbf{R}_i}\right)$	)
--	--	---

Т	F	ho	Energy	Forces
3	64	0.010	0.228	0.401
6	64	0.010	0.202	0.217
3	128	0.010	0.188	0.197
6	128	0.010	0.1002	0.120
6	128	0.100	0.027	0.171
6	128	0.010	0.100	0.120
6	128	0.001	0.238	0.061
6	128	0.000	0.260	0.058

K.T. Schütt et al., J. Chem. Phys. 148, 241722, (2018).



Constructing accurate ML force fields for flexible molecules AIDD Spring School Lugano-Switzerland // 10.05.2022

 $\mathbf{2}$ 





Slide 16

#### Method 2: First neural network architectures **View Article Online**

## **ANI potentials**





$$\begin{split} G_m^{\mathbf{R}} &= \sum_{j \neq i} e^{-\eta \left(R_{ij} - R_s\right)} f_{\mathbf{C}} \left(R_{ij}\right) \\ G_m^{\mathbf{A}_{\text{mod}}} &= 2^{1-\zeta} \sum_{j,k \neq i}^{\text{all atoms}} \left(1 + \cos(\theta_{ijk} - \theta_s)\right)^{\zeta} \\ &\times \exp\left[-\eta \left(\frac{R_{ij} + R_{ik}}{2} - R_s\right)^2\right] f_{\mathbf{C}} \left(R_{ij}\right) f_{\mathbf{C}} \left(R_{ik}\right) \end{split}$$

Each different chemical symbol has a distinct NN.

target property

$$L^{\text{tot}} = \frac{1}{N} \sum_{t}^{T} \sum_{p} \sum_{i}^{N} w_{t} w_{p} (y_{tpi} - \hat{y}_{tpi})^{2}$$

Atomic forces of the ANI potentials calculated by using the are automatic differentiation feature of PyTorch library.

#### **Rotational energy profiles**



J.S. Smith et al., Chem. Sci. 8, 3192, (2017).



# Method 3: Recent Physics-inspired NN potentials Scalable and accurate ML force field

<u>SpookyNet</u>: It models electronic degrees of freedom and non-local interactions using attention in a transformer architecture.



#### **Physics-based components**





# Method 3: Recent Physics-inspired NN potentials Scalable and accurate ML force field

<u>SpookyNet</u>: It models electronic degrees of freedom and non-local interactions using attention in a transformer architecture.



## SpookyNet model: QM7-X dataset

- Training set: ~4 M of equilibrium and nonequilibrium molecules up to 23 atoms.
- Level of theory: PBE0+MBD.
- Features: 128.
- Cut-off = 5.29 Å.
- Model parameters: 3 630 142.

O. T. Unke *et al., Nat. Commun.* **12**, 7273, (2021).





# Method 3: Recent Physics-inspired NN potentials Scalable and accurate ML force field

- □ Four molecules: 30, 40, 50, and 60 atoms.
- References geometries optimized at PBE0(tight)+MBD.



P. Thölke and G. De Fabritiis, (2022). arXiv:2202.02541 O. T. Unke et al., Nat. Commun. 12, 7273, (2021).



## Method 4: Hybrid ML/Molecular-Mechanics potentials Binding free energy calculations

 $U_{\rm ML/MM}\left(X_{\rm P}, X_{\rm L}\right) =$ 

 $U_{\mathrm{MM}}\left(X_{\mathrm{P}},X_{\mathrm{L}}
ight)$  -  $U_{\mathrm{MM}}^{^{vacuum}}\left(X_{\mathrm{L}}
ight)$  +  $U_{\mathrm{ML}}^{^{vacuum}}\left(X_{\mathrm{L}}
ight)$ 



#### D.M. Rufa et al., (2020). bioRxiv: 2020.07.29.227959.



## Method 4: Hybrid ML/Molecular-Mechanics potentials **Binding free energy calculations**



 $U_{\mathrm{MM}}\left(X_{\mathrm{P}},X_{\mathrm{L}}
ight)$  -  $U_{\mathrm{MM}}^{^{vacuum}}\left(X_{\mathrm{L}}
ight)$  +  $U_{\mathrm{ML}}^{^{vacuum}}\left(X_{\mathrm{L}}
ight)$ 





**MM: openFF** (+solvent)

-8

Α



-8







#### D.M. Rufa et al., (2020). bioRxiv: 2020.07.29.227959.



# Method 4: Hybrid ML/Molecular-Mechanics potentials BuRNN: Buffer region NN for polarizable embedding



$$V_{tot} = V_{\mathbb{I}+\mathbb{B}}^{QM} - V_{\mathbb{B}}^{QM} + V_{\mathbb{B}}^{MM} + V_{\mathbb{O}(\mathbb{I}+\mathbb{B}+\mathbb{O})}^{MM}$$

$$\cong V^{NN}_{\mathbb{I}+\Delta\mathbb{B}} + V^{MM}_{\mathbb{B}+\mathbb{O}}$$



B. Lier et al., J. Phys. Chem. Lett. 13, 3812, (2022).



Constructing accurate ML force fields for flexible molecules AIDD Spring School Lugano-Switzerland // 10.05.2022

Features:

- Predict the difference between two QM regions.
- Polarizable embedding of buffer region at full QM level.



## Method 4: Hybrid ML/Molecular-Mechanics potentials System specific MD simulations with HDNNP

- □ Challenges
  - Long range interactions (electrostatic, vdW)
  - Large phase space
  - Long time scales
- □ Approach
  - Symmetry functions as in ANI-x models
  - $\Delta$ -learning scheme with DFTB as baseline

#### **Molecules in water**



L. Böselt et al., J. Chem. Theory Comput. 17, 2641, (2021).



 $L = \frac{1}{N} \cdot \sum_{i=1}^{N} (E_i - \tilde{E}_i)^2$ 

N<sub>QM</sub> 3

Loss functions:





# Method 5: Hybrid ML/Quantum-Mechanics potentials Density functional tight-binding (DFTB) method



$$E_{ ext{rep}} = E_{ ext{DFT}} - E_{ ext{DFTB}}^{( ext{el})} 
ot\approx \sum_{A < B} \mathcal{V}_{ ext{rep}}^{(AB)}ig(R_{AB}ig)$$

Problems of accuracy and general validity.



#### M. Stöhr et al., J. Phys. Chem. Lett. 11, 6835, (2020).



Constructing accurate ML force fields for flexible molecules AIDD Spring School Lugano-Switzerland // 10.05.2022

#### Rotational energy profile:



Slide 25

#### Many-body NN repulsive potentials

# Summary

#### Increasing the system dimensionality







**Drug-protein binding** 

#### Increasing the system complexity

- ✓ 2-and 3-body geometric descriptors
- ✓ Few degrees of freedom
- ✓ Weak non-covalent interaction

- High-dimensional geometric descriptors
- Large degrees of freedom
- Strong non-covalent interaction

energy

- Strong intermolecular interaction
- Dependence on the chemical environment
- > Hybrid QM/ML models



# Acknowledgments

## **University of Luxembourg**



#### **External collaborators**

- Prof. Dr. Klaus-Robert Müller (TU Berlin)
- Prof. Dr. Robert DiStasio Jr. (Cornell univ.)
- Dr. Hugo Ceulemans (Janssen Phar.)
- Dr. Dries van Rompaey (Janssen Phar.)

# Funding and computational resources















#### Many thanks for your attention.

#### Email: leonardo.medrano@uni.lu





