

# From Atoms to Accelerated Discovery: Machine Learning Pathways for Materials Innovation

## Discuss: Descriptors / Unsupervised Learning

Nong Artrith

Debye Institute for Nanomaterials Science  
Faculty of Science, Utrecht University, The Netherlands

Email: [n.artrith@uu.nl](mailto:n.artrith@uu.nl)



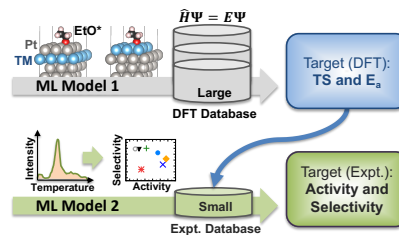
2026 CaMML - Chemistry and Materials Machine Learning School  
Daresbury Laboratory, UK, April 14<sup>th</sup>, 2026

1

### Computational Discovery of Energy Materials and Interpretation of Experimental Observations with Atomistic First-Principles Methods and Machine Learning

#### Research Activities

Nong Artrith and Team  
Computational Materials Science  
and Machine Learning (ML)

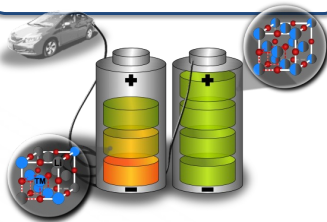


#### 1. Machine Learning & Data Science

- Data from simulation and experiment (with collaborators)
- Accelerating DFT calculations
- Predicting materials properties
- Data mining for materials discovery

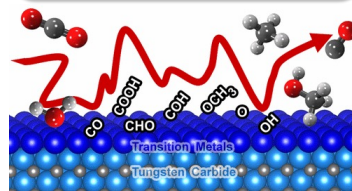
#### 2. Energy Storage

- Storage for renewable energy
- Enable electric vehicles
- Improve portable electronics



#### 3. Energy Conversion

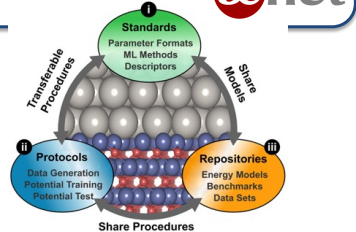
- Replace fossil fuels (oil/gas)
- Produce synthetic fuels using renewable energy



#### 4. Method Development

- **Open-source** ML modeling tools

<http://nartrith.atomistic.net>



2

## We know the laws of physics that govern (most) materials properties

- Quantum mechanics describes extremely accurately how chemical bonds form
- Already developed 100 years ago!
- If we can simulate QM, we can predict new materials *in silico*

Nong Artrith, n.artrith@uu.nl

2026 CaMML - Chemistry and Materials Machine Learning School

3

## But length and time scales are severely limited!

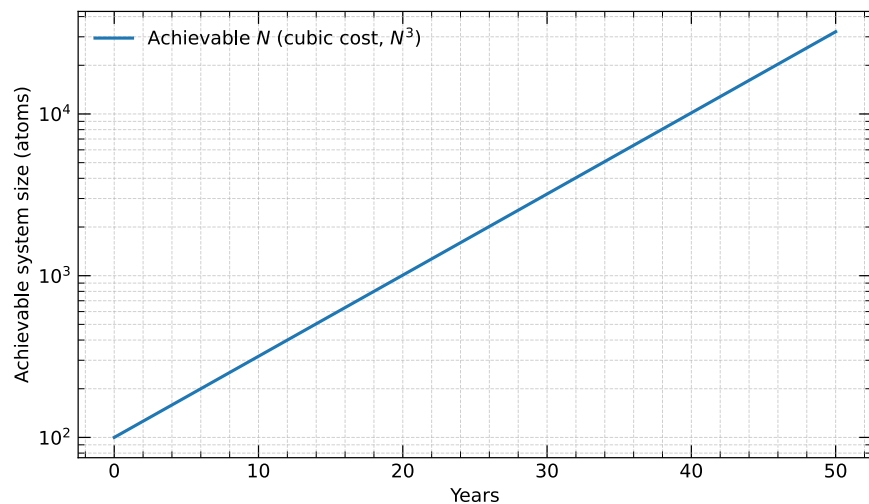
Cost of DFT

$$\propto (N_{\text{atoms}})^3$$

Moore's law

$$\propto 2^{N_{\text{years}}}$$

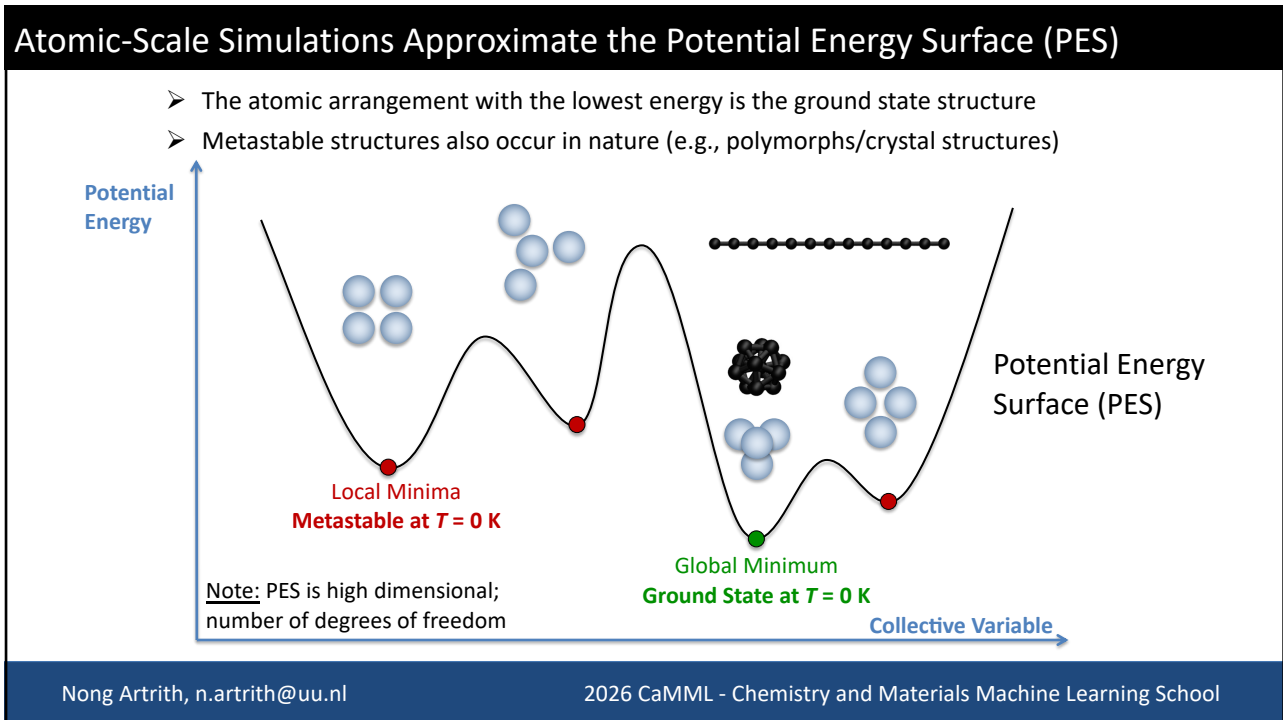
Practical limit:  
Tens of thousands  
of atoms.



Nong Artrith, n.artrith@uu.nl

2026 CaMML - Chemistry and Materials Machine Learning School

4



5

### How can Machine Learning be used for Materials Chemistry?

- Machine Learning can be used as:  
**A method for data analysis that automates analytical model building**
- Learn from data, identify patterns, and make decisions with minimal human intervention
- Data can be from experimental measurement or from computational modeling

6

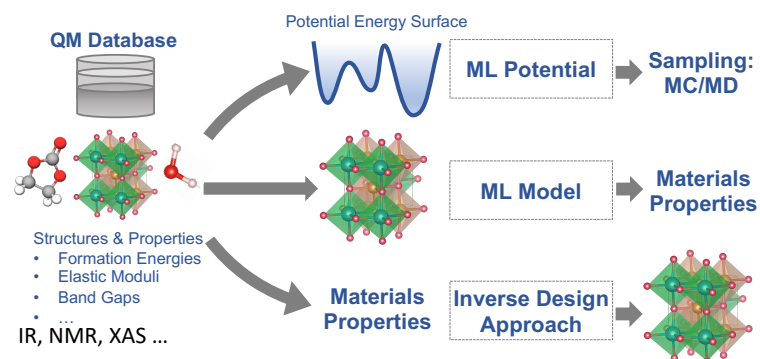
## Textbooks

- Christopher Bishop,  
*Pattern Recognition and Machine Learning*,  
Springer **2006**
- Kevin P. Murphy,  
*Machine Learning: A Probabilistic Perspective*,  
MIT Press **2012**
- Grégoire Montavon, Geneviève Orr, Klaus-Robert Müller,  
*Neural Networks: Tricks of the Trade*,  
Springer **2012**

7

## Machine Learning for Accelerated Modeling

- Machine learning (ML) can be used to accelerate atomic-scale simulations
- ML models can also be trained on measured data to predict materials
- ML is also useful for the interpretation of measurements (e.g., from spectroscopy)



H. Guo, Q. Wang, A. Stuke, A. Urban, and **N. Arith**, *Frontiers in Energy Research* **9**, 2021, 695902

8

## Outline: Principles of Machine Learning

1. Popular Learning Types and Their Applications
  - Supervised learning
  - Unsupervised learning
  - Semi-supervised learning
  - Reinforcement learning
2. Models and Model Training
3. Model Validation
4. Materials Science Examples

9

## 1.1 Supervised Learning

### **Supervised learning: Learn from labeled data**

- Data pairs  $(x, y)$  are available to train the model
- $x$  is the *feature*,  $y$  is the *label*
- ML model learns to predict  $y$  for unseen  $x$
- Typical applications: Classification, regression

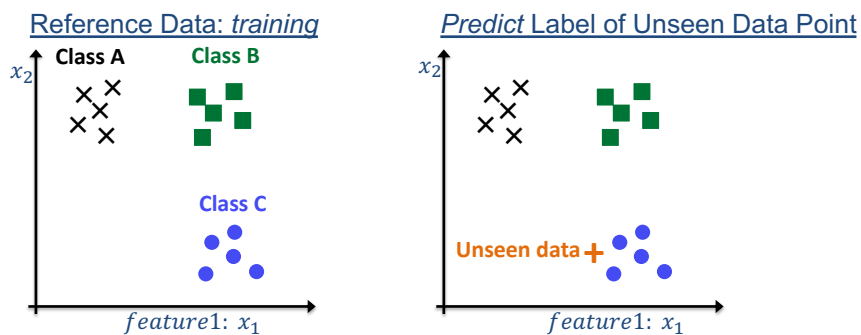
10

## Classification: Predict Label of New Data Points

- **Labeled** data points used for *training*



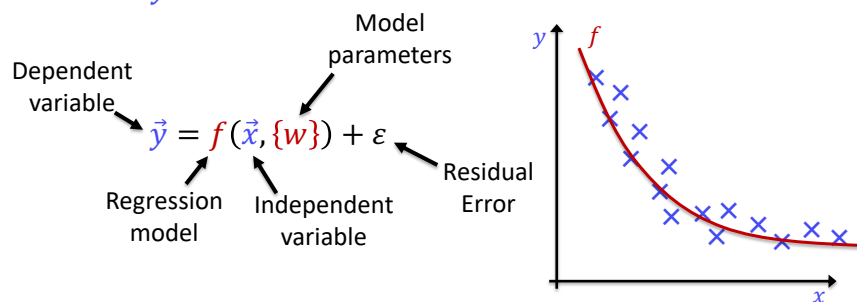
- *Prediction* of **discrete** labels for **new** data points



11

## Regression: Estimating Variable Relationships

- Given data points  $(\vec{x}, \vec{y})$  determine a **continuous** function that relates the independent variables  $\vec{x}$  to the dependent variables  $\vec{y}$



- Machine learning provides *models* and methods for determining the model parameters by *training*

12

## 1.2 Unsupervised Learning

### Unsupervised learning: Detect structure in data

- Only features given, but *no labels*
- Given data points  $\{x\}$ ; determine relationship between data points
- Typical applications: Clustering, dimension reduction

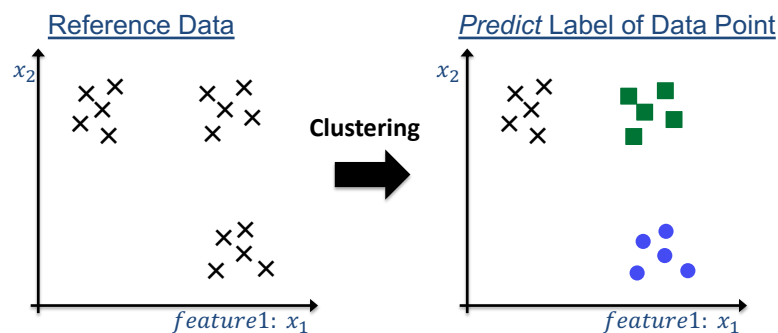
13

## Clustering: Determine Which Data Belongs Together

- Only **features** used for *training*



- *Prediction* of the labels for **original** data points



14

## 2. Models and Model Training

15

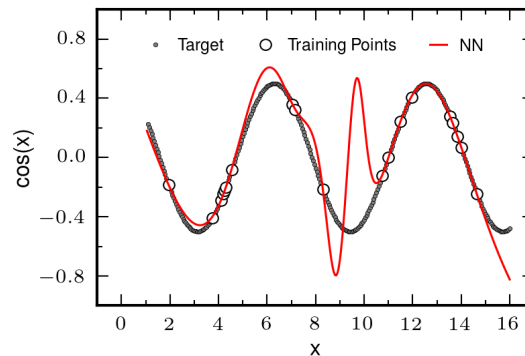
### The Right Choice of ML Model Depends on the Data

- **Size of the reference data set**
  - Billions of Google search requests
  - vs.**
  - Hundreds of materials compositions
- **Dimension of the learning problem**
  - Pixels in a photo (1600 x 800)
  - vs.**
  - Free energy  $G(T,p)$
- ...

16

## Complex Models can be *Overfitted*

- Model accurately reproduces data used for training
- But model interpolates incorrectly → **poor generalization**



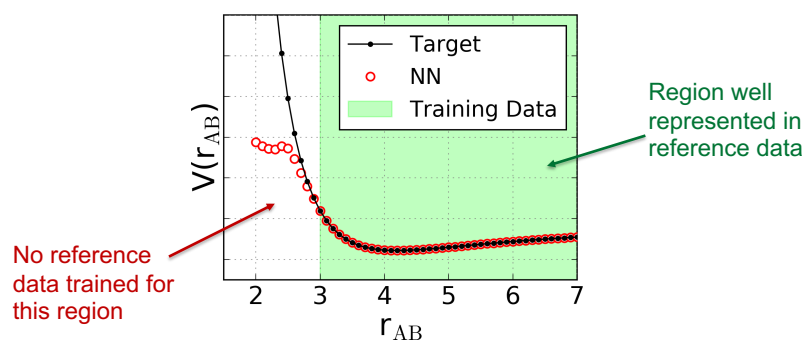
A.M. Miksch, T. Morawietz, J. Kästner, A. Urban, **N. Artrith\***,  
*Machine Learning: Science and Technology*, **2**, 031001(2021), DOI: <https://doi.org/10.1088/2632-2153/abfd96>.

Code: <https://github.com/atomisticnet/MLP-beginners-guide>

17

## Machine Learning is Poor for *Extrapolation*

- ANNs are flexible functions
- Can well interpolate between reference data points
- But **no good for extrapolation** beyond trained region



A.M. Miksch, T. Morawietz, J. Kästner, A. Urban, **N. Artrith\***,  
*Machine Learning: Science and Technology*, **2**, 031001(2021), DOI: <https://doi.org/10.1088/2632-2153/abfd96>.

Code: <https://github.com/atomisticnet/MLP-beginners-guide>

18

### 3. Model Validation

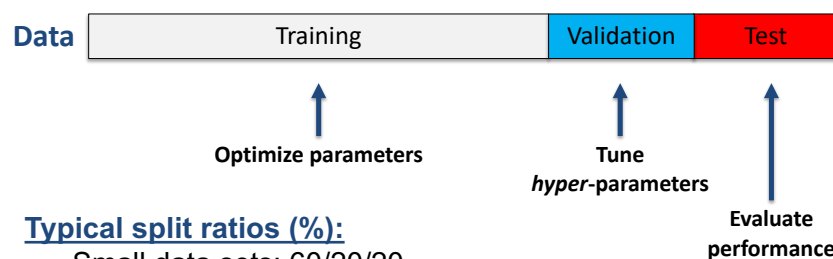
Machine Learning models have to be validated.

19

## ML Models Have to be Validated

**Solution:** Split reference data set into different parts

- **Training set:** data samples that are used for training
- **Validation set:** used to tune hyper parameters
- **Test set:** used to evaluate performance (*error estimate*)



**Typical split ratios (%):**

Small data sets: 60/20/20

Large data sets: 98/1/1

20

## Cross-Validation

- **Problem:** hyper-parameters can be overfit to test data
- **Solution:** Use many splits
  - Example: Nested cross-validation

optimize parameters
Tune hyper-parameters

**Outer Loop:**

- Train with optimal set
- Average test errors

↙

**Inner Loop:**  
tune hyper-parameters

21

## Early Stopping: Prevent Overfitting

- Monitor mean error of samples in test set during training
- Test set error does no longer improve when the model is overfitted

Training error

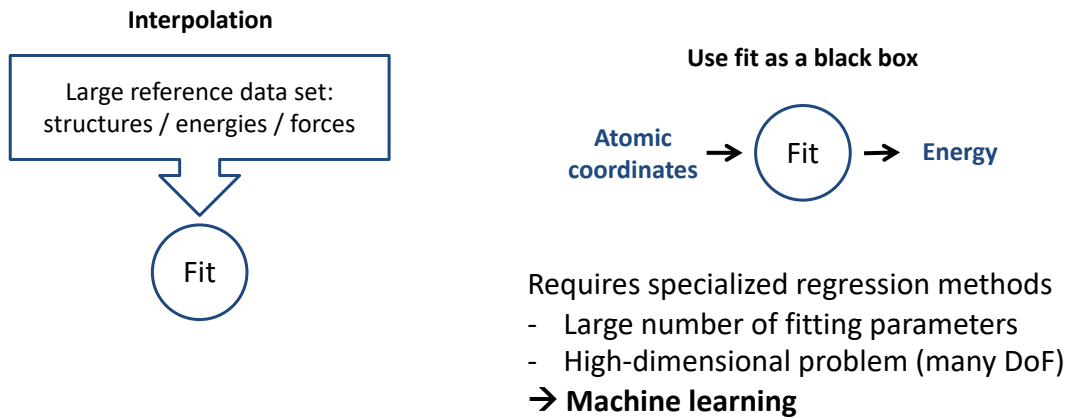
Training iteration

Test set error remains same or increases

Training set error continues to decrease

22

## Machine-Learning Interatomic Potentials: Learn from QM, Inference with ML



Nong Artrith, n.artrith@uu.nl

2026 CaMML - Chemistry and Materials Machine Learning School

23

## Not a Novel Idea!

- Machine learning is used for classification (e.g., image recognition) and regression (function fitting)
- Well suited for regression:
  - *Gaussian Process Regression (or Kriging)*
  - *Artificial Neural Networks*
- Both methods have been used for potential fitting

Nong Artrith, n.artrith@uu.nl

2026 CaMML - Chemistry and Materials Machine Learning School

24

## Early ANN Potentials Were Published in the 1990s



6 June 1997

**CHEMICAL  
PHYSICS  
LETTERS**

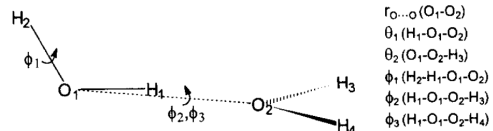
Chemical Physics Letters 271 (1997) 152–156

### Description of the potential energy surface of the water dimer with an artificial neural network

 Kyoung Tai No <sup>a,1</sup>, Byung Ha Chang <sup>a</sup>, Su Yeon Kim <sup>a</sup>, Mu Shik Jhon <sup>b,1</sup>,  
Harold A. Scheraga <sup>c</sup>
<sup>a</sup> Department of Chemistry, Sogang University, Seoul 156-743, South Korea  
<sup>b</sup> Department of Chemistry, Korea Advanced Institute of Science and Technology, Taejeon 305-701, South Korea  
<sup>c</sup> Baker Laboratory of Chemistry, Cornell University, Ithaca, NY 14853-1301, USA

#### ANN Input:

Most molecular degrees of freedom (distances, angles, dihedrals); fixed O-H bond length



#### Drawbacks:

- Not symmetric wrt. exchange of equivalent atoms.
- Potential can **only** be used for the water dimer
- Extension to other systems is not straightforward

Nong Artrith, n.artrith@uu.nl

2026 CaMML - Chemistry and Materials Machine Learning School

25

## The Breakthrough: Learn the Atomic Energy

 Inspiration from bond-order potentials:  
 Decompose total energy into **atomic contributions**

$$E = \sum_i^{\text{atoms}} E_i(\sigma_i)$$

$E_i$  is an ANN       $\sigma_i$  describes the local environment of atom  $i$

#### Note:

- The atomic energy is neither known nor uniquely defined.
- The target for the training is the total energy of the atomic structure.
- The model *learns* one possible atomic energy decomposition.

J. Behler, and M. Parrinello, *Phys. Rev. Lett.* **98** (2007) 146401.

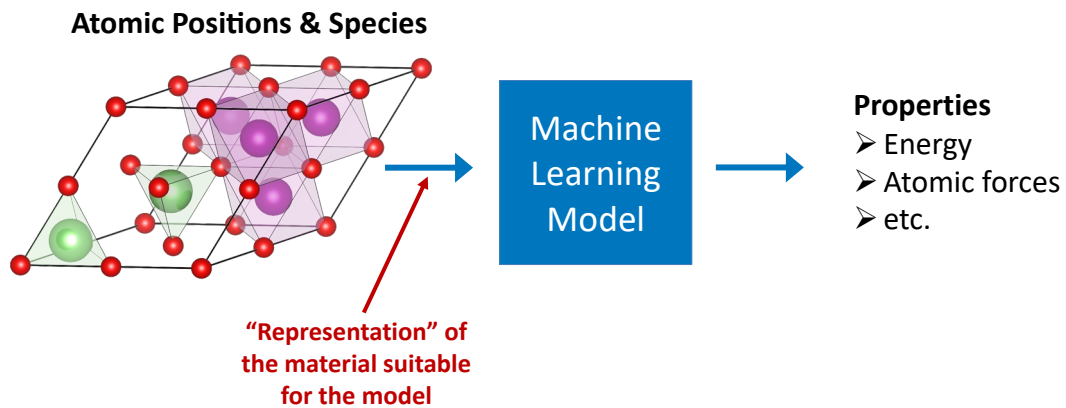
Nong Artrith, n.artrith@uu.nl

2026 CaMML - Chemistry and Materials Machine Learning School

26

## How do you input a chemical compound into a machine-learning model?

- Our model should work with any arbitrary atomic structure in a defined chemical space



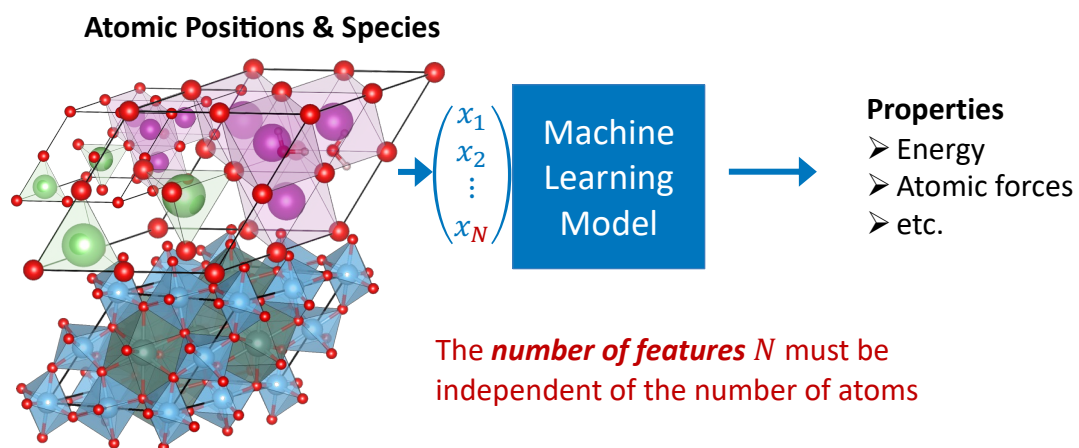
Nong Artrith, n.artrith@uu.nl

2026 CaMML - Chemistry and Materials Machine Learning School

27

## Requirement 1: (Descriptor) *Feature vector* with constant dimension

- Many ML models expect an input vector with constant dimension



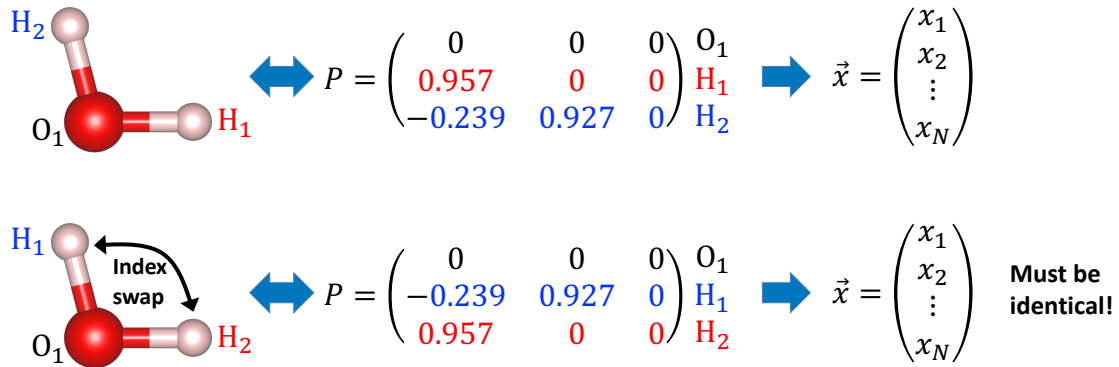
Nong Artrith, n.artrith@uu.nl

2026 CaMML - Chemistry and Materials Machine Learning School

28

## Requirement 2: Exchanging equivalent atoms must not affect the result

- Many ML models expect an input vector with constant dimension



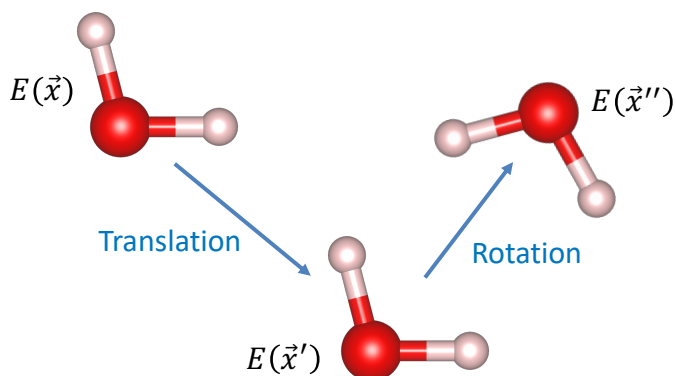
Nong Artrith, n.artrith@uu.nl

2026 CaMML - Chemistry and Materials Machine Learning School

29

## Requirement 3: Scalar properties must be invariant wrt. translation & rotation

- The energy and other scalar properties, such as the charge, do not change upon translation, rotation, or reflection



$$E(P + \vec{t}) = E(P)$$

- $P = (\vec{r}_1, \vec{r}_2, \dots)^T$  is a matrix with the atomic positions
- $\vec{t}$  is an arbitrary translation
- $P + \vec{t} = (\vec{r}_1 + \vec{t}, \vec{r}_2 + \vec{t}, \dots)^T$

$$E(RP) = E(P)$$

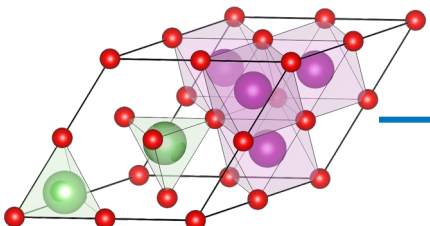
- $R$  is an arbitrary rotation matrix

Nong Artrith, n.artrith@uu.nl

2026 CaMML - Chemistry and Materials Machine Learning School

30

### Summary of the requirements for materials features



$$\vec{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{pmatrix}$$

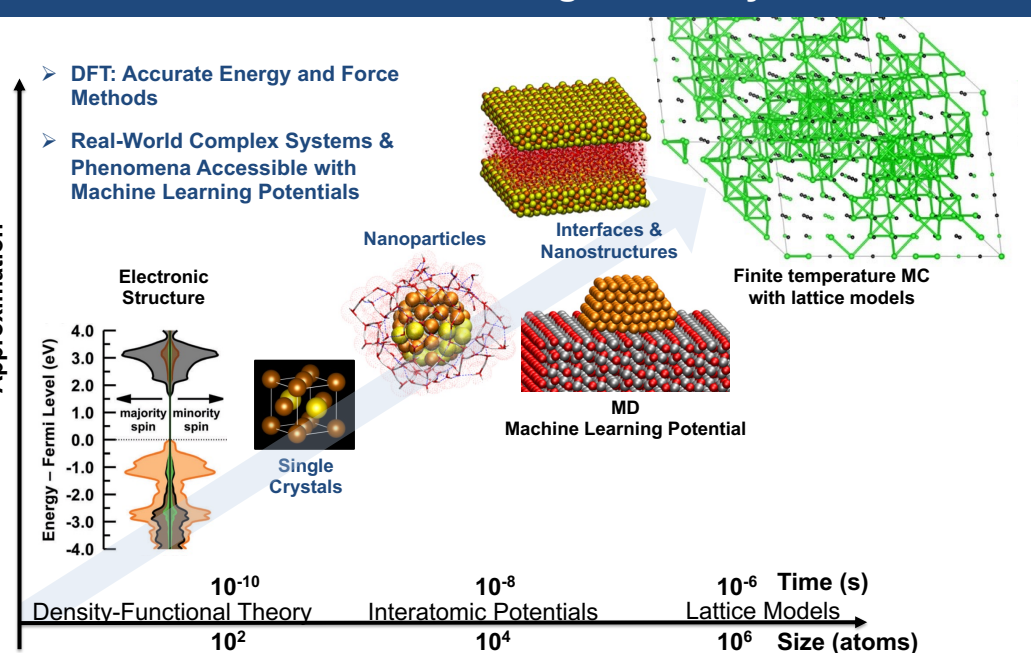
- Dimension independent of the number of atoms
- Invariant to the exchange of equivalent atoms
- Support for model E(3) equivariance; = energy invariant with respect to
  - Rotation,
  - Reflection, and
  - Translation

Nong Artrith, n.artrith@uu.nl
2026 CaMML - Chemistry and Materials Machine Learning School

31

### Accurate Simulation of Large-Scale Systems

- DFT: Accurate Energy and Force Methods
- Real-World Complex Systems & Phenomena Accessible with Machine Learning Potentials



**10<sup>-10</sup>**  
Density-Functional Theory

**10<sup>2</sup>**

**10<sup>-8</sup>**  
Interatomic Potentials

**10<sup>4</sup>**

**10<sup>-6</sup>** Time (s)  
Lattice Models

**10<sup>6</sup>** Size (atoms)

32



## Featurization: Representation Methods for Translating Materials into Vectors

### Option 1 (relatively easy): Ignore the atomic details of the crystal structure

- Elements can be represented by their known properties, *i.e.*, from the periodic table
- Group & period, electronegativity, ionization energies, mass, quantum numbers, etc.
- Element-wise features need to be combined to represent compositions
  - weighted mean, standard deviation, other averages, etc.

## Featurizing the structure (lattice & sites) as well is harder

### Option 2 (harder): Also featurized the structure, *i.e.*, atomic positions/sites and lattice

- Need to (should) maintain invariants (rotation, translation, permutation of equivalent atoms)
- Ideally works with **different numbers of atoms**

#### Two principal directions:

- “Hand-crafted” representations = recipes that a human has come up with. For example, based on basis-set expansions.
- Learned representations = specialized ANNs such as E(3)-equivariant **graph neural networks**

## Selection of hand-crafted (physics-based) material representations

### Some hand-crafted representations (selection only)

#### Atomic sites (→ local representation)

- Atom-centered Symmetry functions, ACSF (Behler, Parrinello, *Phys. Rev. Lett.* **98**, 2007, 146401)
- Smooth Overlap of Atomic Positions, SOAP (Bartók, Kondor, Csányi, *Phys. Rev. B* **87**, 2013, 184115)
- Moment tensors (Shapeev, *Multiscale Model. Simul.* **14**, 2016, 1153–1173)
- Atomic Cluster Expansion, ACE (Drautz, *Phys. Rev. B* **99**, 2019, 014104)

#### Entire structure (→ global representation)

- Diffraction patterns
- Coulomb matrices (Rupp ... von Lilienfeld, *Phys. Rev. Lett.* **108**, 2012, 058301)
- Crystal-symmetry and sites (Gómez-Bombarelli and coworkers)

#### Recent reviews:

Himanan, ..., Rinke, Foster, *Comput. Phys. Commun.* **247**, 2020, 106949.

Musli, ..., Csányi, Ceriotti, *Chem. Rev.* **121**, 2021, 9759–9815.

Damewood, ..., Gómez-Bombarelli, *Annu. Rev. Mater. Res.* **53**, 2023, 12.1–12.28.

Nong Artrith, n.artrith@uu.nl

2026 CaMML - Chemistry and Materials Machine Learning School

37

## Local representation: Expansion of the RDF, ADF, etc. in Orthonormal Basis Set

$$\text{RDF}_i(r) = \sum_{\alpha} c_{\alpha}^{(2)} \phi_{\alpha}(r) \text{ for } 0 \leq r \leq R_c$$

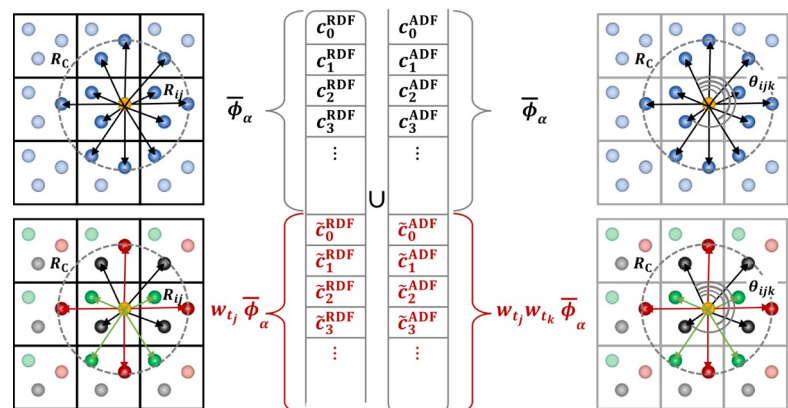
$$\text{ADF}_i(\theta) = \sum_{\alpha} c_{\alpha}^{(3)} \phi_{\alpha}(\theta) \text{ for } 0 \leq \theta \leq \pi$$

$\{\phi_i\}$  is a complete orthonormal basis.

Features: Expansion coefficients **with and without** element-specific weights.

$$c_{\alpha}^{(2)} = \sum_{j \neq i} w_{t_j} \bar{\phi}_{\alpha}(R_{ij}) f_c(R_{ij})$$

$$c_{\alpha}^{(3)} = \sum_{j, k \neq i} w_{t_j} w_{t_k} \bar{\phi}_{\alpha}(\theta_{ijk}) f_c(R_{ij}) f_c(R_{ik})$$



Method: N. Artrith, A. Urban, G. Ceder, *Phys. Rev. B* **96**, 2017, 014112.

Figure: A.M. Miksch, T. Morawietz, J. Kästner, A. Urban, N. Artrith, *Mach. Learn. Sci. Tech.* **2**, 2021, 031001.

Nong Artrith, n.artrith@uu.nl

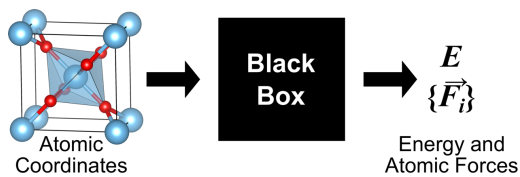
2026 CaMML - Chemistry and Materials Machine Learning School

38

## Our Model $\text{\ae}net$ : Actively Developed Since ~2011. First Public Version in 2016.

### Machine-learned interatomic potentials (MLIPs) as "drop-in replacement" for DFT

- Trained on first principles reference data
- Nearly as accurate as the reference method
- But very data hungry; can require thousands of reference calculations



$\text{\ae}net$

<https://github.com/atomisticnet>

N. Artrith and A. Urban, *Comput. Mater. Sci.* **114** (2016) 135-150.

N. Artrith, A. Urban, G. Ceder, *Phys. Rev. B* **96** (2017) 014112.

A. Cooper, J. Kästner, A. Urban, N. Artrith, *npj Comput. Mater.* **6** (2020) 54.

I. W. Yeu et al., *npj Comput. Mater.* **11** (2025) 156.

Nong Artrith, n.artrith@uu.nl

2026 CaMML - Chemistry and Materials Machine Learning School

39

## The $\text{\ae}net$ Package

First publicly available implementation of the Machine Learning Potential (MLP) method!

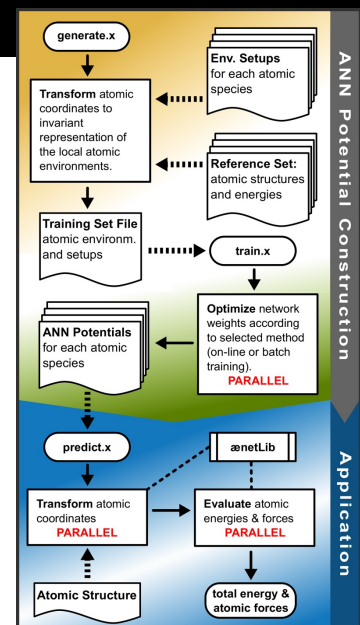
The Atomic Energy Network

$\text{\ae}net$

<https://github.com/atomisticnet>

Modules for potential fit and for application

- C compatible library & Python interface
- Interfaces with simulation software: ASE, LAMMPS, Tinker, DL\_POLY
- Efficient (approximate) force pre-training

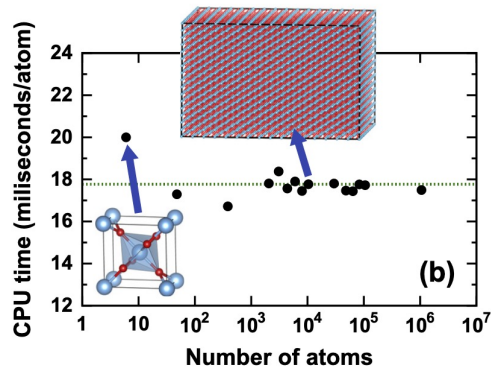


Nong Artrith, n.artrith@uu.nl

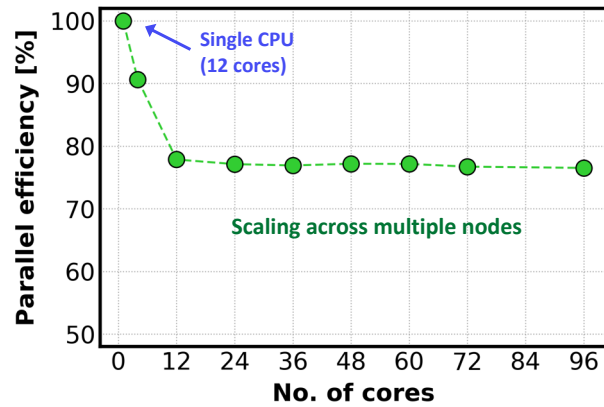
2026 CaMML - Chemistry and Materials Machine Learning School

40

## $\mathcal{A}$ enet Potentials Scale Linearly and are Highly Parallelizable



N. Artrith and A. Urban,  
*Comput. Mater. Sci.* **114** (2016) 135-150.



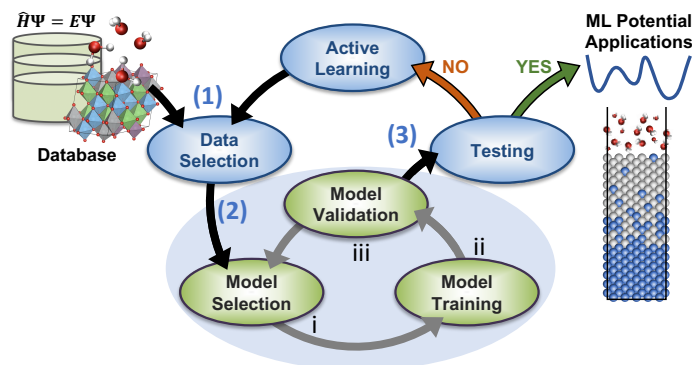
M.S. Chen, T. Morawietz, H. Mori, T.E. Markland, N. Artrith,  
*J. Chem. Phys.* **155**, 074801 (2021).

Nong Artrith, n.artrith@uu.nl

2026 CaMML - Chemistry and Materials Machine Learning School

41

## MLIPs are Typically Constructed in an Iterative (Active Learning) Fashion



A.M. Miksch, T. Morawietz, J. Kästner, A. Urban, N. Artrith,  
*Mach. Learn.: Sci. Tech.* **2** (2021) 031001.  
Code: <https://github.com/atomisticnet/MLP-beginners-guide>

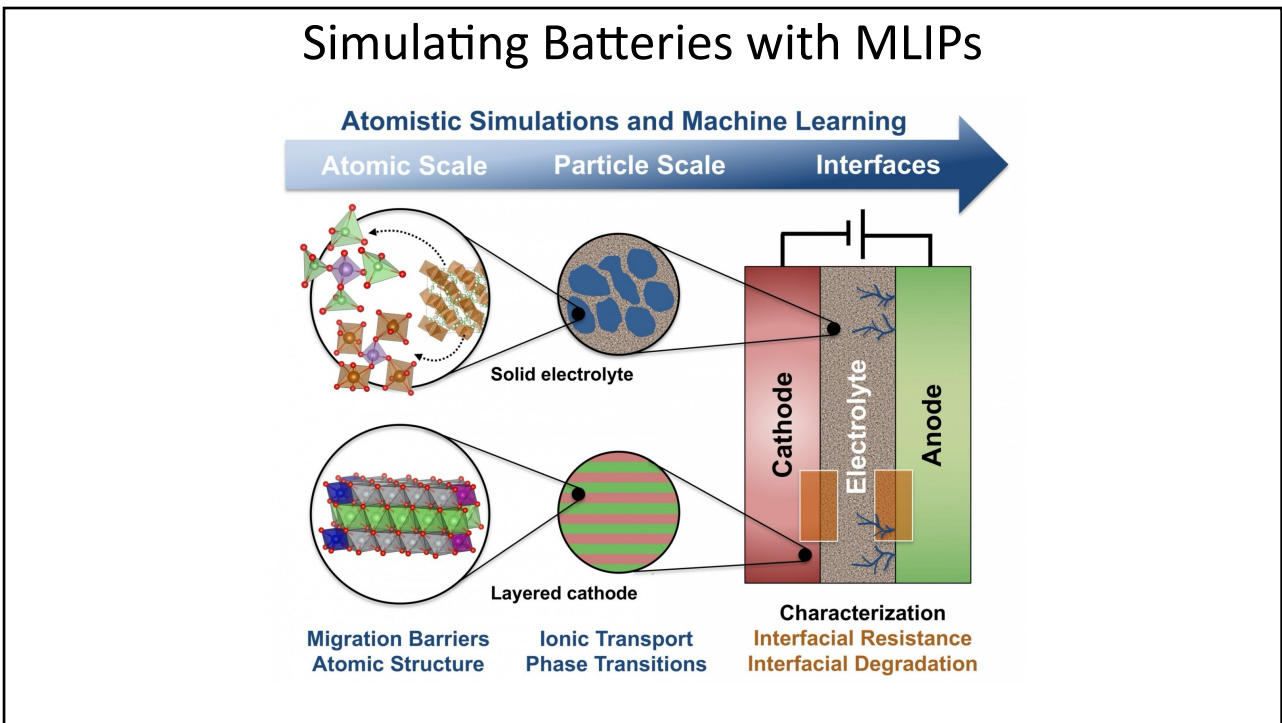
Nong Artrith, n.artrith@uu.nl

2026 CaMML - Chemistry and Materials Machine Learning School

42

# Some Applications

43



44

### Enable the Sampling of Amorphous Materials

Describing amorphous phases requires large numbers of atoms due to the lack of long-range order

1 String Representation  
2 Initial Population of  $N$  Delithiated Configurations  
3 Energy from ANN Potential  
4a Select Low-Energy Configurations  
4b Converged? Select Lowest-Energy Configurations  
5 Random Crossing  
6 Next Generation  
DFT Optimization of Vacancy Structures

$\text{Li}_{15}\text{Si}_4$

$E_1$   
 $E_2$   
 $E_3$

Lowest-Energy DFT Structure as Input for the Next Step

ae net

N. Artrith, A. Urban, G. Ceder, *J. Chem. Phys.* **148** (2017) 241711.

Nong Artrith, n.artrith@uu.nl 2026 CaMML - Chemistry and Materials Machine Learning School

45

### Quantitative Modeling of Amorphous Materials Feasible

ANN potential predicts phase diagram and voltage profiles in agreement with experiment

N. Artrith, A. Urban, and G. Ceder, *J. Chem. Phys.* **148**, (2018) 241711.

**Computed Phase Diagram**

Formation Energy (eV/ $\text{Li}_x\text{Si}$ ) vs  $x$  in  $\text{Li}_x\text{Si}$

Crystalline  $\text{Li}_x\text{Si}$  (DFT)  
GA  $\alpha\text{-Li}_{15-x}\text{Si}_4$  (DFT)

$\text{Li}_8\text{Si}_3$   
 $\text{Li}_{12}\text{Si}_7$   
 $\text{Li}_7\text{Si}_3$   
 $\text{Li}_{13}\text{Si}_4$   
 $\text{Li}_{15}\text{Si}_4$   
 $\text{Li}_{21}\text{Si}_5$

**Voltage Profile**

Potential vs.  $\text{Li}^+/\text{Li}$  (V) vs  $x$  in  $\text{Li}_x\text{Si}$

Huggins 1981 (415°C)  
Dahn 2004 (Thin Film RT)  
Crystalline  $\text{Li}_x\text{Si}$  (DFT)  
GA  $\alpha\text{-Li}_{15-x}\text{Si}_4$  (DFT)

Please see this reference for how to calculate a DFT phase diagram:  
Computational understanding of Li-ion batteries. *npj Comput Mater* **2**, 16002 (2016). <https://doi.org/10.1038/npjcompumats.2016.2>

Nong Artrith, n.artrith@uu.nl 2026 CaMML - Chemistry and Materials Machine Learning School

46

### Makes it Possible to Model Entire Nanoparticles with Tens of Thousands of Atoms

Isolated Si atoms      Shell: Si chains  
Core: Isolated Si atoms      Gradual Si clustering

$\text{Li}_{3.28}\text{Si}$        $\text{Li}_{2.89}\text{Si}$        $\text{Li}_{2.50}\text{Si}$        $\text{Li}_{1.98}\text{Si}$        $\text{Li}_{1.71}\text{Si}$        $\text{Li}_{1.44}\text{Si}$

● Si  
● Li

Nong Artrith, n.artrith@uu.nl      2026 CaMML - Chemistry and Materials Machine Learning School

47

### Predictions are in Quantitative Agreement with Experiment

$\text{Li}_{3.50}\text{Si}$        $\text{Li}_{2.25}\text{Si}$        $\text{Li}_{1.00}\text{Si}$

● Si  
● Li

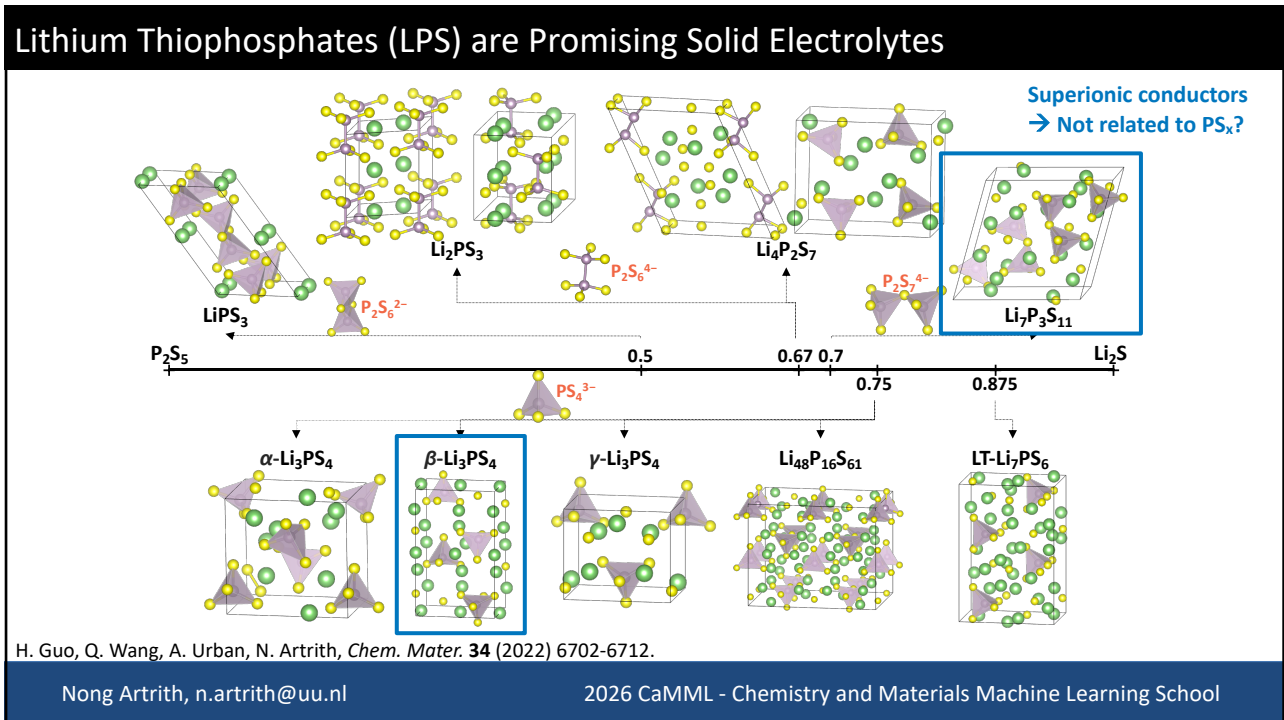
Li Transport (Arrhenius Plot)

$x_{\text{Li}}$	$E_a$ (eV)	$D$ ( $\text{cm}^2\text{s}^{-1}$ )
0.75	0.789	$1.154 \times 10^{-14}$
1.00	0.500	$5.986 \times 10^{-11}$
2.25	0.483	$9.607 \times 10^{-11}$
3.50	0.682	$3.820 \times 10^{-13}$

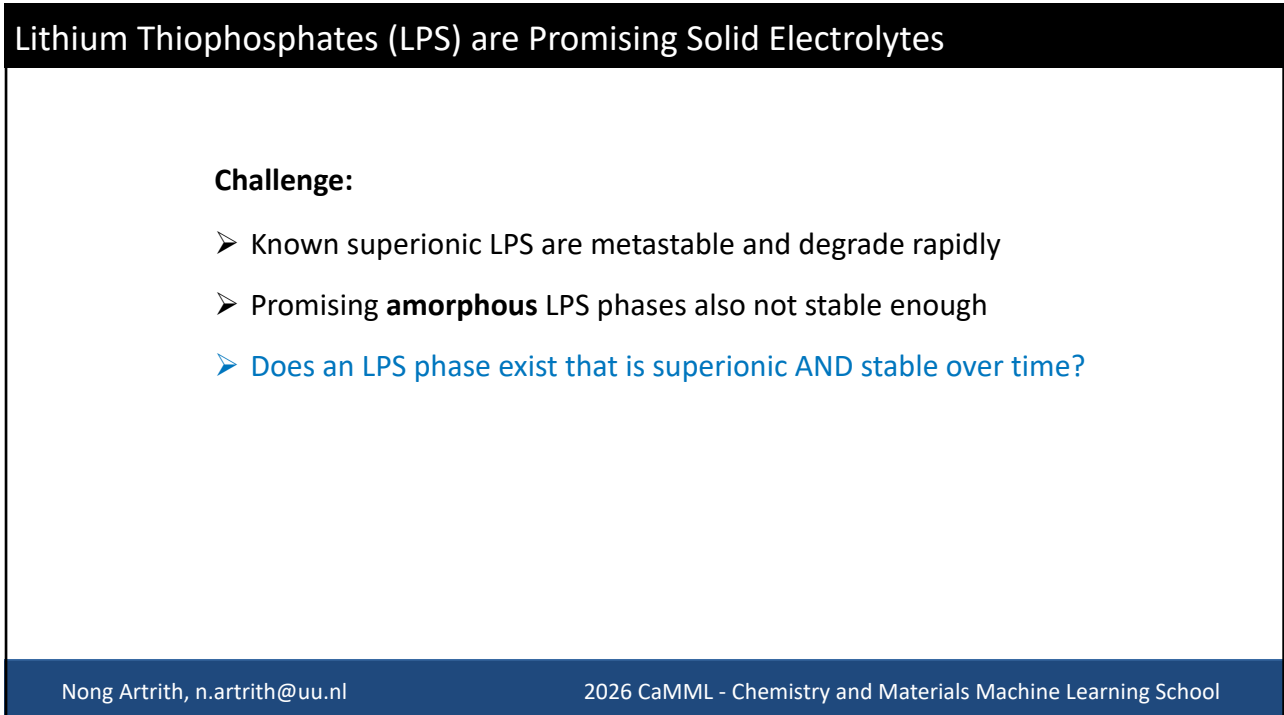
Experimental References		
$D$ ( $\text{cm}^2/\text{s}$ )	Method	Reference
$10^{-10}$	EIS	[1]
$10^{-12}$	CV, EIS, GITT	[2]
$10^{-14}$	EIS, PITT	[3]
$10^{-14}$ – $10^{-13}$	PITT	[4]

Nong Artrith, n.artrith@uu.nl      2026 CaMML - Chemistry and Materials Machine Learning School

48

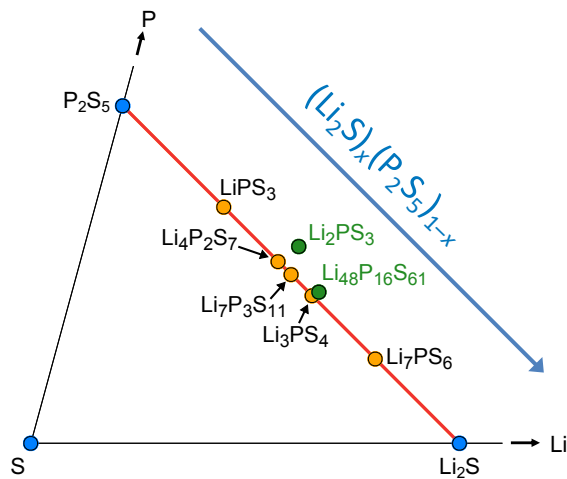


49



50

## Approach: First Map the Phase Diagram of Crystalline and Amorphous LPS



### ML-accelerated sampling

Vary  $x$  in  $(\text{Li}_2\text{S})_x(\text{P}_2\text{S}_5)_{1-x}$

1. Start with supercell of crystalline LPS
2. 2 Li and 1 S (or 2 P and 5 S)
3. Search for most stable arrangement of the remaining atoms
4. Repeat steps 2 & 3

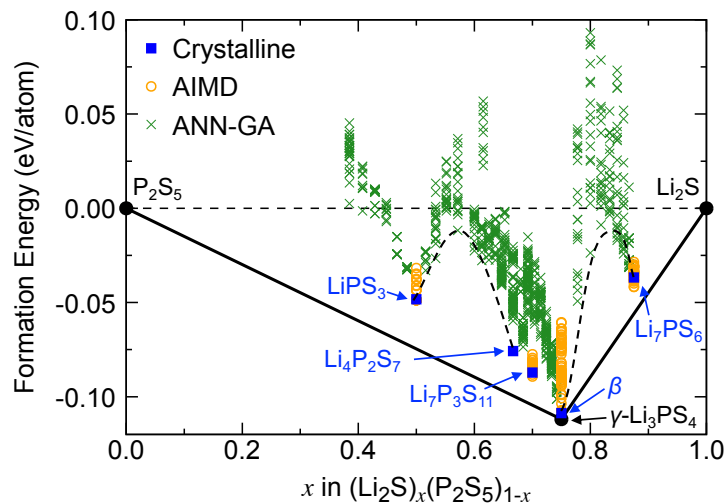
Nong Artrith, n.artrith@uu.nl

2026 CaMML - Chemistry and Materials Machine Learning School

51

## Only $\gamma\text{-Li}_3\text{PS}_4$ is Thermodynamically Stable at Zero Temperature

OK LPS phase diagram: which compositions can likely be made  $\rightarrow$  "Energy above hull"



Which compositions are best for solid electrolytes?

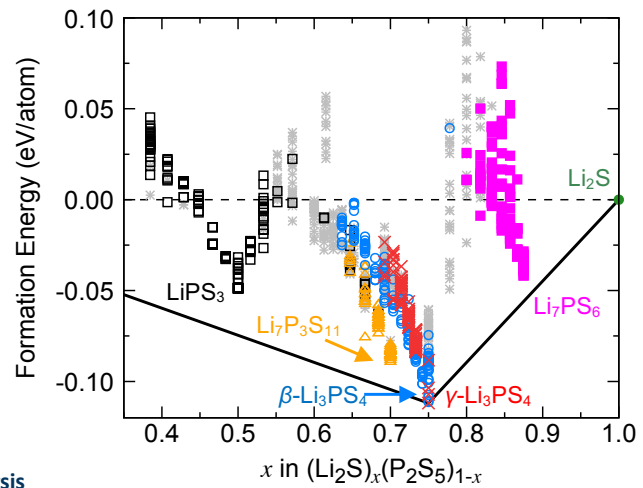
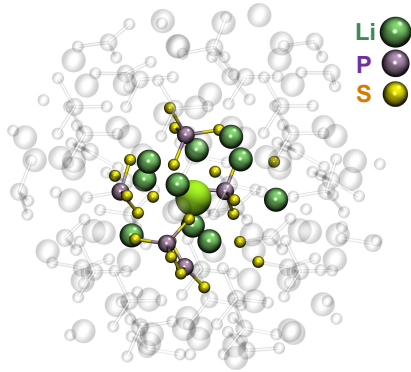
H. Guo, Q. Wang, A. Urban, N. Artrith, *Chem. Mater.* **34** (2022) 6702-6712.

Nong Artrith, n.artrith@uu.nl

2026 CaMML - Chemistry and Materials Machine Learning School

52

## Classifying Amorphous Phases based on their Local Li Site Environment



► Unsupervised learning for similarity classification using our  $\text{\ae}$ net descriptors and k-means clustering analysis of Li-site environments in our structural databases.

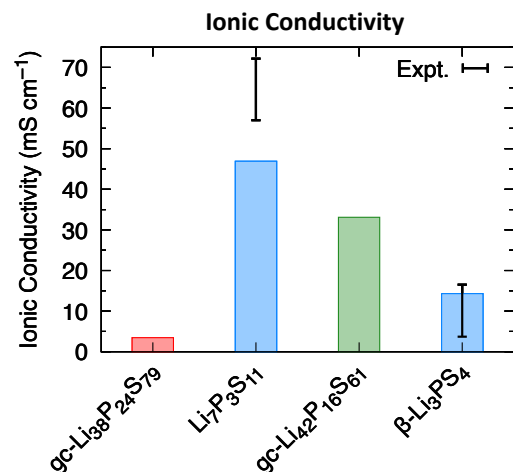
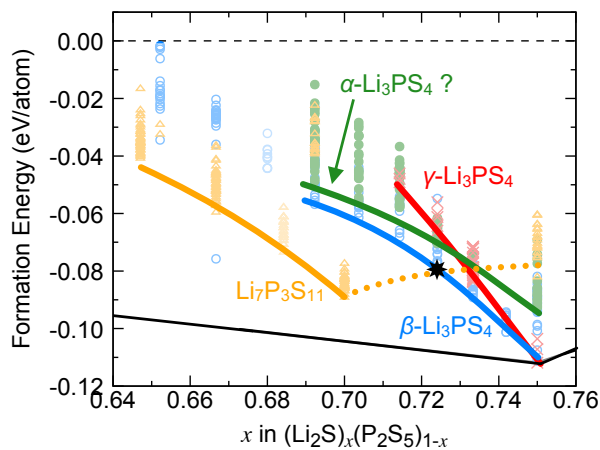
Nong Artrith, n.artrith@uu.nl

2026 CaMML - Chemistry and Materials Machine Learning School

53

## LPS Compositions Between $\text{Li}_7\text{P}_3\text{S}_{11}$ and $\beta\text{-Li}_3\text{PS}_4$ are Likely Stable Conductors

Molecular dynamics simulations indicate that  $\text{gc-Li}_{42}\text{P}_{16}\text{S}_{61}$  is a stable superionic conductor



H. Guo, Q. Wang, A. Urban, N. Artrith, *Chem. Mater.* **34** (2022) 6702-6712.

Nong Artrith, n.artrith@uu.nl

2026 CaMML - Chemistry and Materials Machine Learning School

54



55

## Universal MLIPs

**Matbench Discovery** [↗](#)

<https://matbench-discovery.materialsproject.org>

Full Test Set Unique Prototypes 10k Most Stable

Model	CPS ↑	Acc	F1	DAF	Prec	MAE	R <sup>2</sup>	K <sub>SRME</sub>	RMSD	Training Set	Params	Targets	Date Added
EquiformerV3+DeNS-OAM	0.902	0.978	0.931	6.074	0.928	0.018	0.868	0.118	0.059	6.6M (113M) OMat24+MPtrj+sAlex	30.3M	EF <sub>S<sub>G</sub></sub>	2026-04-07
PET-OAM-XL	0.898	0.977	0.924	6.075	0.929	0.019	0.864	0.119	0.060	6.6M (113M) OMat24+sAlex+MPtrj	730M	EF <sub>S<sub>G</sub></sub>	2026-01-10
TACE-OAM-L	0.889	0.972	0.910	5.898	0.902	0.020	0.868	0.126	0.061	6.6M (113M) OMat24+sAlex+MPtrj	82.9M	EF <sub>S<sub>G</sub></sub>	2026-04-09
eSEN-30M-OAM	0.888	0.977	0.925	6.069	0.928	0.018	0.866	0.170	0.061	6.6M (113M) OMat24+MPtrj+sAlex	30.2M	EF <sub>S<sub>G</sub></sub>	2025-03-17
EquiFlash	0.888	0.975	0.919	5.983	0.915	0.019	0.871	0.158	0.060	6.6M (113M) OMat24+MPtrj+sAlex	28.7M	EF <sub>S<sub>G</sub></sub>	2025-06-23
Nequip-OAM-XL	0.886	0.971	0.906	5.869	0.897	0.020	0.872	0.125	0.063	6.6M (113M) OMat24+sAlex+MPtrj	32.1M	EF <sub>S<sub>G</sub></sub>	2025-11-30
MatRIS-10M-OAM	0.877	0.976	0.921	6.039	0.923	0.019	0.871	0.218	0.060	6.6M (113M) OMat24+sAlex+MPtrj	10.4M	EF <sub>S<sub>G</sub>M</sub>	2025-10-29
SevenNet-Omni-112	0.873	0.971	0.906	5.954	0.910	0.021	0.868	0.192	0.062	243M COSMOSDataset	54.9M	EF <sub>S<sub>G</sub></sub>	2026-01-12
Nequip-OAM-L	0.870	0.967	0.893	5.823	0.890	0.022	0.865	0.166	0.065	6.6M (113M) OMat24+sAlex+MPtrj	9.6M	EF <sub>S<sub>G</sub></sub>	2025-09-08
GRACE-2L-OAM-L	0.865	0.964	0.883	5.840	0.893	0.022	0.862	0.169	0.064	6.6M (113M) OMat24+sAlex+MPtrj	26.4M	EF <sub>S<sub>G</sub></sub>	2025-09-09
ORB v3	0.860	0.971	0.905	5.912	0.904	0.024	0.821	0.210	0.075	6.47M (133M) MPtrj+Alex+OMat24	25.5M	EF <sub>S<sub>G</sub></sub>	2025-04-05
Allegro-OAM-L	0.840	0.966	0.895	5.674	0.867	0.022	0.868	0.319	0.065	6.6M (113M) OMat24+sAlex+MPtrj	9.7M	EF <sub>S<sub>G</sub></sub>	2025-09-08
GRACE-2L-OAM	0.837	0.963	0.880	5.774	0.883	0.023	0.862	0.294	0.067	6.6M (113M) OMat24+sAlex+MPtrj	12.6M	EF <sub>S<sub>G</sub></sub>	2025-02-06
EquiformerV3+DeNS-MP	0.830	0.956	0.863	5.479	0.838	0.029	0.840	0.275	0.070	146k (1.58M) MPtrj	30.3M	EF <sub>S<sub>G</sub></sub>	2026-04-07
DPA-3.1-3M-FT	0.802	0.963	0.884	5.667	0.866	0.023	0.869	0.469	0.069	163M OpenLAM	3.27M	EF <sub>S<sub>G</sub></sub>	2025-06-05
eSEN-30M-MP	0.797	0.946	0.831	5.260	0.804	0.033	0.822	0.340	0.075	146k (1.58M) MPtrj	30.1M	EF <sub>S<sub>G</sub></sub>	2025-03-17
MACE-MPA-0	0.795	0.954	0.852	5.582	0.853	0.028	0.842	0.412	0.073	3.37M (12M) MPtrj+sAlex	9.06M	EF <sub>S<sub>G</sub></sub>	2024-12-09
MatRIS-10M-MP	0.778	0.951	0.847	5.422	0.829	0.031	0.824	0.489	0.072	146k (1.58M) MPtrj	10.4M	EF <sub>S<sub>G</sub>M</sub>	2025-10-29
AlphaNet-v1-OAM	0.769	0.968	0.901	5.747	0.879	0.024	0.831	0.643	0.079	6.6M (113M) OMat24+sAlex+MPtrj	4.65M	EF <sub>S<sub>G</sub></sub>	2025-05-12

56

## Can we Train a Universal MLIP to Replace DFT in All Simulations?

nature computational science

Article

<https://doi.org/10.1038/s43588-022-00349-3>

### A universal graph deep learning interatomic potential for the periodic table

Received: 18 March 2022

Accepted: 5 October 2022

Published online: 28 November 2022

Check for updates

Chi Chen<sup>1</sup> & Shyue Ping Ong<sup>1</sup>

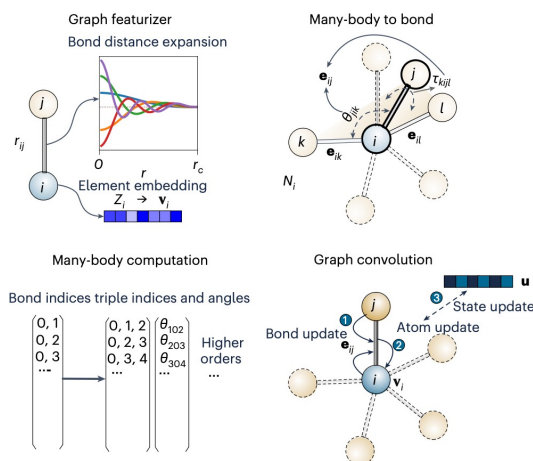
Interatomic potentials (IAPs), which describe the potential energy surface of atoms, are a fundamental input for atomistic simulations. However, existing IAPs are either fitted to narrow chemistries or too inaccurate for general applications. Here we report a universal IAP for materials based on graph neural networks with three-body interactions (M3GNet). The M3GNet IAP was trained on the massive database of structural relaxations performed by the Materials Project over the past ten years and has broad applications in structural relaxation, dynamic simulations and property prediction of materials across diverse chemical spaces. About 1.8 million materials from a screening of 31 million hypothetical crystal structures were identified to be potentially stable against existing Materials Project crystals based on M3GNet energies. Of the top 2,000 materials with the lowest energies above the convex hull, 1,578 were verified to be stable using density functional theory calculations. These results demonstrate a machine learning-accelerated pathway to the discovery of synthesizable materials with exceptional properties.

C. Chen and S. P. Ong, *Nat. Comput. Sci.* **2** (2022) 718–728.

Nong Artrith, n.artrith@uu.nl

2026 CaMML - Chemistry and Materials Machine Learning School

Trained on 1.8 million DFT calculations!



57

## Tech Industry Shows What is Possible When Money is Unlimited

### UMA: A Family of Universal Models for Atoms

Brandon M. Wood<sup>1,†,\*</sup>, Misko Dzamba<sup>1,†,\*</sup>, Xiang Fu<sup>1,\*</sup>, Meng Gao<sup>1,\*</sup>, Muhammed Shuaibi<sup>1,\*</sup>, Luis Barroso-Luque<sup>1</sup>, Kareem Abdelmaqsood<sup>2</sup>, Vahe Gharakhanyan<sup>1</sup>, John R. Kitchin<sup>2</sup>, Daniel S. Levine<sup>1</sup>, Kyle Michel<sup>1</sup>, Anuroop Sriram<sup>1</sup>, Taco Cohen<sup>1</sup>, Abhishek Das<sup>1</sup>, Ammar Rizvi<sup>1</sup>, Sushree Jagriti Sahoo<sup>1</sup>, Zachary W. Ulissi<sup>1</sup>, C. Lawrence Zitnick<sup>1,†</sup>

<sup>1</sup>FAIR at Meta, <sup>2</sup>Department of Chemical Engineering, Carnegie Mellon University  
\*Co-first Author, <sup>†</sup>Co-corresponding Author

The ability to quickly and accurately compute properties from atomic simulations is critical for advancing a large number of applications in chemistry and materials science including drug discovery, energy storage, and semiconductor manufacturing. To address this need, Meta FAIR presents a family of Universal Models for Atoms (UMA), designed to push the frontier of speed, accuracy, and generalization. UMA models are trained on half a billion unique 3D atomic structures (the largest training runs to date) by compiling data across multiple chemical domains, e.g. molecules, materials, and catalysts. We develop empirical scaling laws to help understand how to increase model capacity alongside dataset size to achieve the best accuracy. The UMA small and medium models utilize a novel architectural design we refer to as mixture of linear experts that enables increasing model capacity without sacrificing speed. For example, UMA-medium has 1.4B parameters but only ~50M active parameters per atomic structure. We evaluate UMA models on a diverse set of applications across multiple domains and find that, remarkably, a single model without any fine-tuning can perform similarly or better than specialized models. We are releasing the UMA code, weights, and associated data to accelerate computational workflows and enable the community to continue to build increasingly capable AI models.

Models: <https://huggingface.co/facebook/UMA>

Code: <https://github.com/facebookresearch/fairchem>

Correspondence: B.M.W. ([bmwood@meta.com](mailto:bmwood@meta.com)), C.L.Z. ([zitnick@meta.com](mailto:zitnick@meta.com))

Meta

- Trained on ½ billion DFT calculations
- Our benchmarks show: Reproduces DFT for **molecular simulations** with great reliability
- Not yet reliable for inorganic materials

#### Other “players”

- Bytedance (Bytedance AI Molecular BOOster, BAMBOO 2024)
- Microsoft (e.g., SimPoly, 2025)
- Schrödinger

Nong Artrith, n.artrith@uu.nl

2026 CaMML - Chemistry and Materials Machine Learning School

58

## Is ænet a Relic of the Past? — Model Distillation



- Graph neural nets are complex and computationally demanding
- No strict linear scaling with the number of atoms
- Require GPU hardware

### Renewed interest in ænet

- Get domain-specialized models for efficient simulations
- Distillation in Python, but inference in Fortran
- ~ 1000 times more efficient

Figure source: <https://www.thoughtco.com/what-is-distillation-601964> (2/26)

Nong Artrith, n.artrith@uu.nl

2026 CaMML - Chemistry and Materials Machine Learning School

59

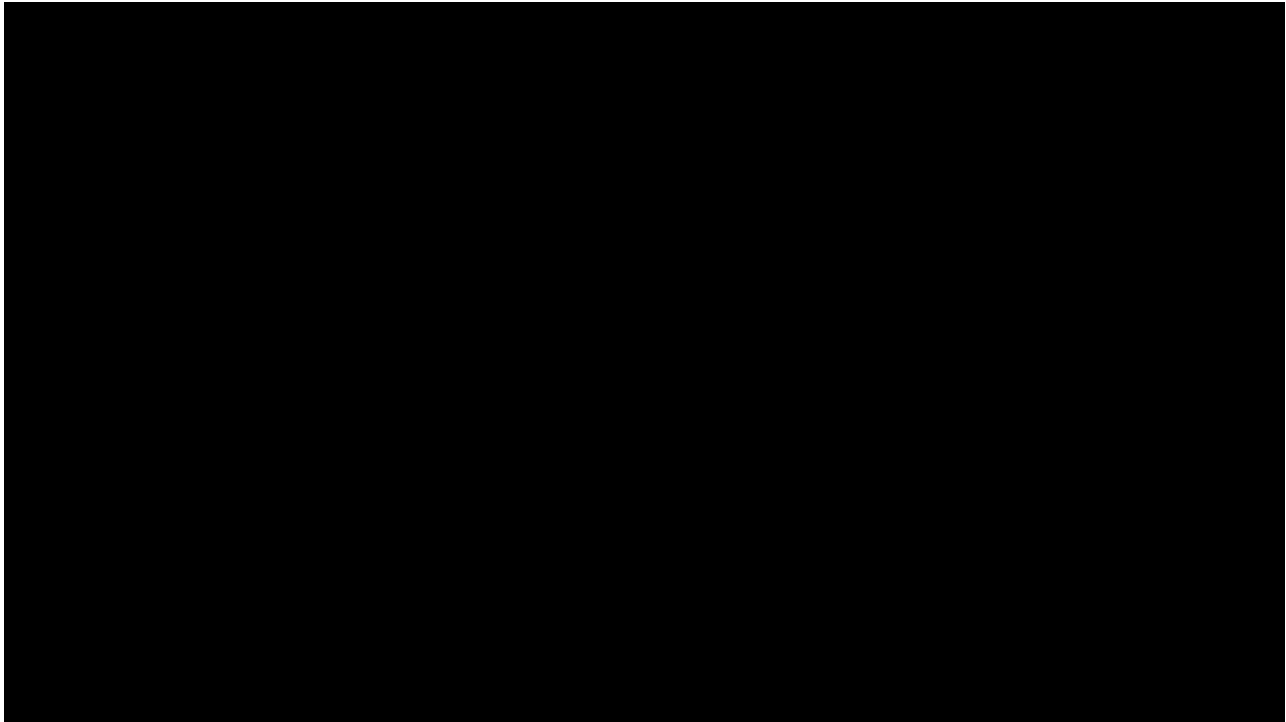
## Summary

- The laws of quantum mechanics determine material properties
- Often, DFT is a good numerical approximation
- (Although sometimes, it is not yet good enough.)
- The computational scaling of DFT prevents routine materials discovery
- ML interatomic potentials can reduce the scaling to linear with atoms
- Recent enormous advances through model scaling (in industry)

Nong Artrith, n.artrith@uu.nl

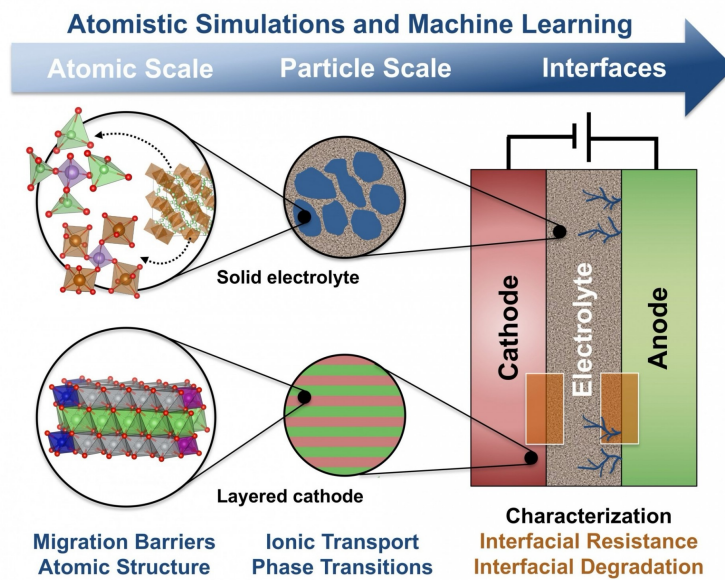
2026 CaMML - Chemistry and Materials Machine Learning School

60



61

## Simulating Batteries with MLIPs



62

## Practical: SOAP descriptor

[https://colab.research.google.com/github/atomisticnet/aenet-gpr/blob/main/tutorial/tutorial\\_3\\_Li-EC.ipynb](https://colab.research.google.com/github/atomisticnet/aenet-gpr/blob/main/tutorial/tutorial_3_Li-EC.ipynb)



**Note on performance**

In section 3-1, the use of the SOAP descriptor in GPR can be computationally intensive. On high-performance computing (HPC) or GPU systems, the SOAP-based GPR runs significantly faster. In this Colab environment, you can switch to a GPU by going to `Runtime → Change runtime type → Hardware accelerator → T4 GPU`

```
try:
    import torch
    print("successfully imported torch")
    print(torch.__version__)
except ImportError:
    ! pip install torch torchvision torchaudio --user --index-url https://download.pytorch.org/whl/cpu
    print("completed installing torch")
```

successfully imported torch  
2.6.0+cu124

**Homework: SOAP descriptor tutorials**

- [GPR Tutorial: H<sub>2</sub>](#)
- [GPR Tutorial: EC-EC](#)
- [GPR Tutorial: Li-EC](#)

 **aenet-gpr**

IW Yeu, A Urban, N Artrith et al., *npj Comput. Mater.* **11**, 156 (2025)

GPR-ANN and GPR-aided NEB Code: <https://github.com/atomisticnet/aenet-gpr>

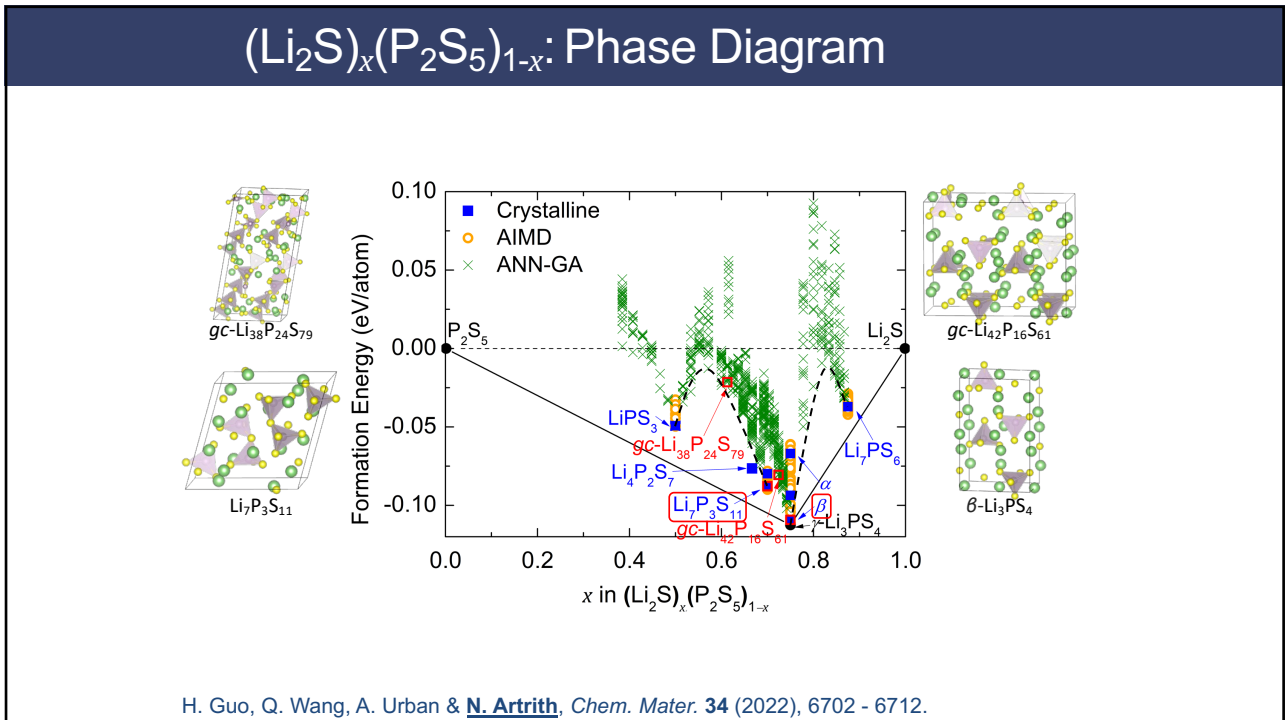
63

## Local Geometry and/or Electronic Structure

K-means is an unsupervised learning algorithm

Nong Artrith (Email: [n.artrith@uu.nl](mailto:n.artrith@uu.nl))

64



65

## X-ray Absorption Spectroscopy

**Transmission XAS**  
 $\mu(E) = \ln(I_0/I_t)$ , i.e. Beers' Law for X-rays

**Fluorescence XAS**  
 $\mu(E) \propto I_f/I_0$

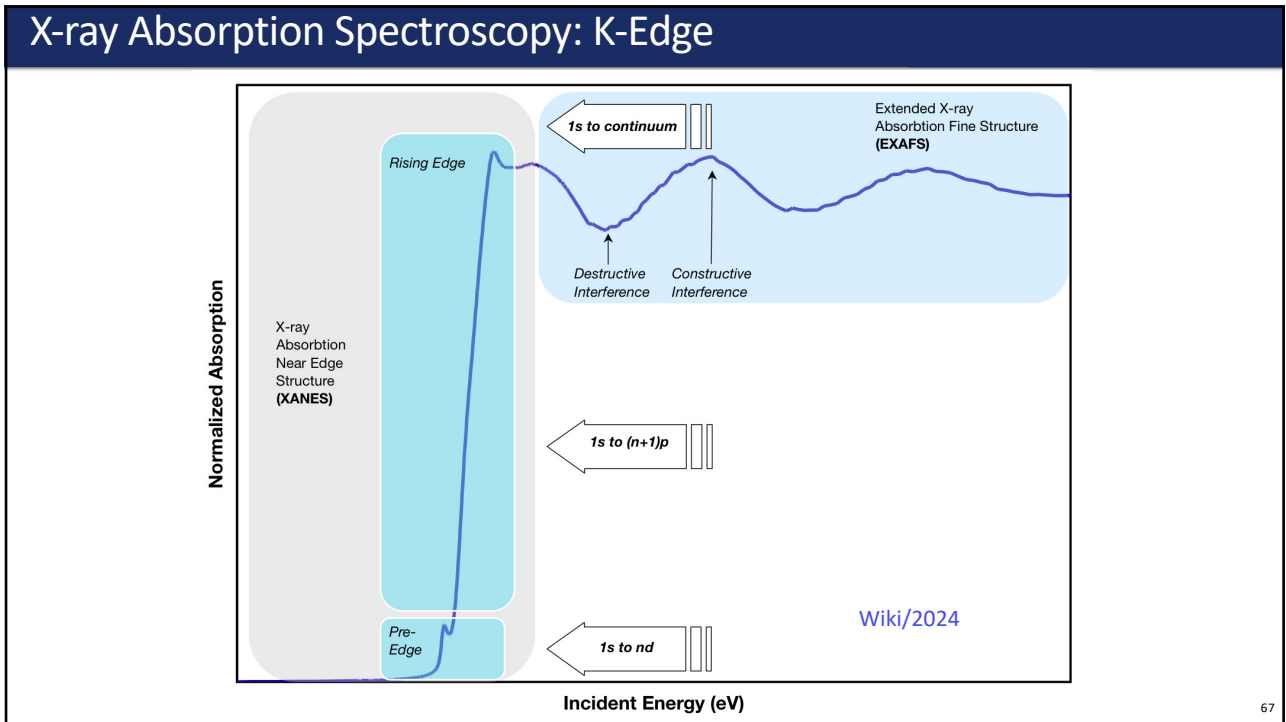
Email: n.artrith@uu.nl

1. An incoming photon interacts with a deep-core electron. Here, a 1s electron is excited for a K-edge spectrum.
2. The deep-core electron is promoted to some unoccupied state above the Fermi energy, propagates away, and leaves behind a core-hole.
3. A short time later (1 or 2 femtoseconds), a higher-lying electron decays into the core-hole and emits a photon.
4. Alternately, the energy from the higher-lying electron can be used to emit an Auger electron.

- Relationship between formal valence of a metal and the position of the edge in XANES spectrum: The shift to higher energy is, to first order, a Coulomb effect. Less charge on the atom means less screening of the core.

[1] Bruce Ravel. *Introduction to X-ray Absorption Spectroscopy* (2015).

66



67

## XAS Simulation

Email: n.artrith@uu.nl

(b)

**Fermi's golden rule:**

$$\sigma(\omega) = 4\pi^2 \frac{\omega}{c} \sum_F |\langle 0|d|F\rangle|^2 \delta(\omega + E_0 - E_F).$$

**Dipole app. :**  $d = \hat{e} \cdot \mathbf{r}$

**Single particle app. :**

$$M_{i \rightarrow f} = \langle \Psi_f | \hat{e} \cdot \mathbf{R} | \Psi_i \rangle \approx S \langle \psi_f | \hat{e} \cdot \mathbf{r} | \psi_i \rangle,$$

**Core hole final state effect:**

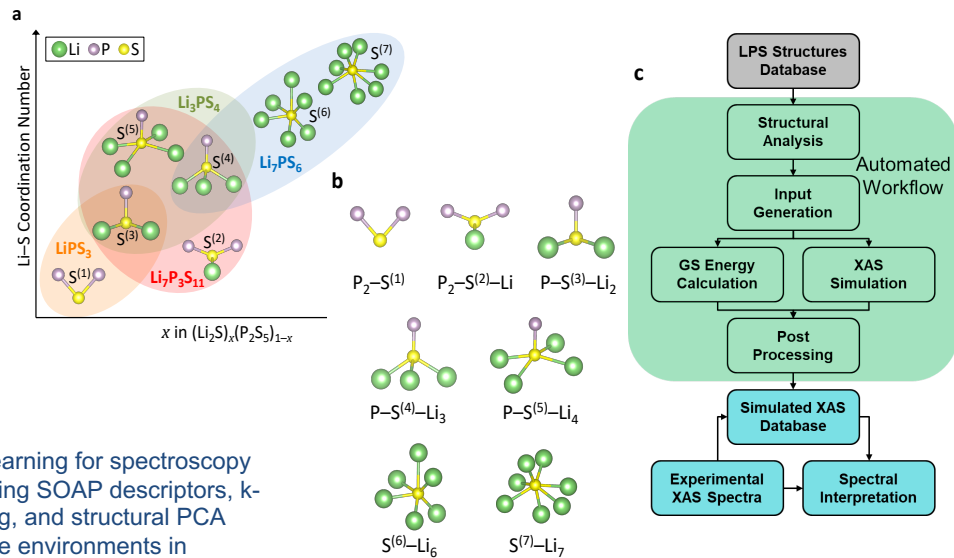
- > Core hole potential approach: SC relaxation of valence electron (Xspectra, VASP 6.1.1)
- > Perturbation approach: linear response of valence electrons (OCEAN, Exciting)

[1] Phil. Trans. R. Soc. A 371.1995 (2013).

68

## Local Coordination: S atoms wt Li & P atoms in *gc*-LPS

Guo, Carbone, N. Artrith, Urban, Lu, et al. *Nature Scientific Data* 10, 349 (2023)



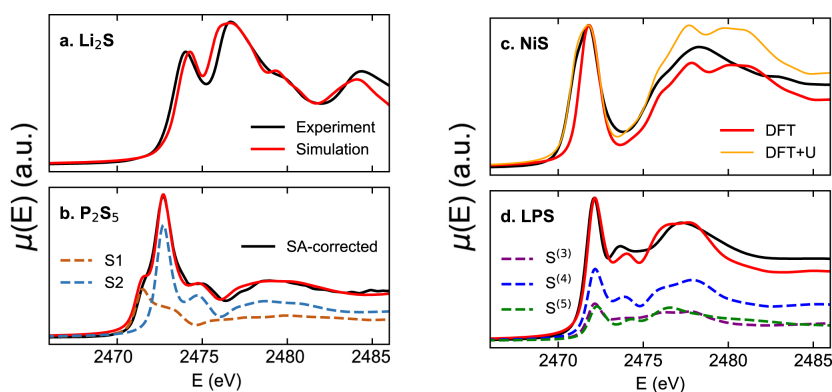
- Unsupervised learning for spectroscopy classification using SOAP descriptors, k-means clustering, and structural PCA analysis of S-site environments in structural databases.

<https://github.com/atomisticnet/xas-tools>

69

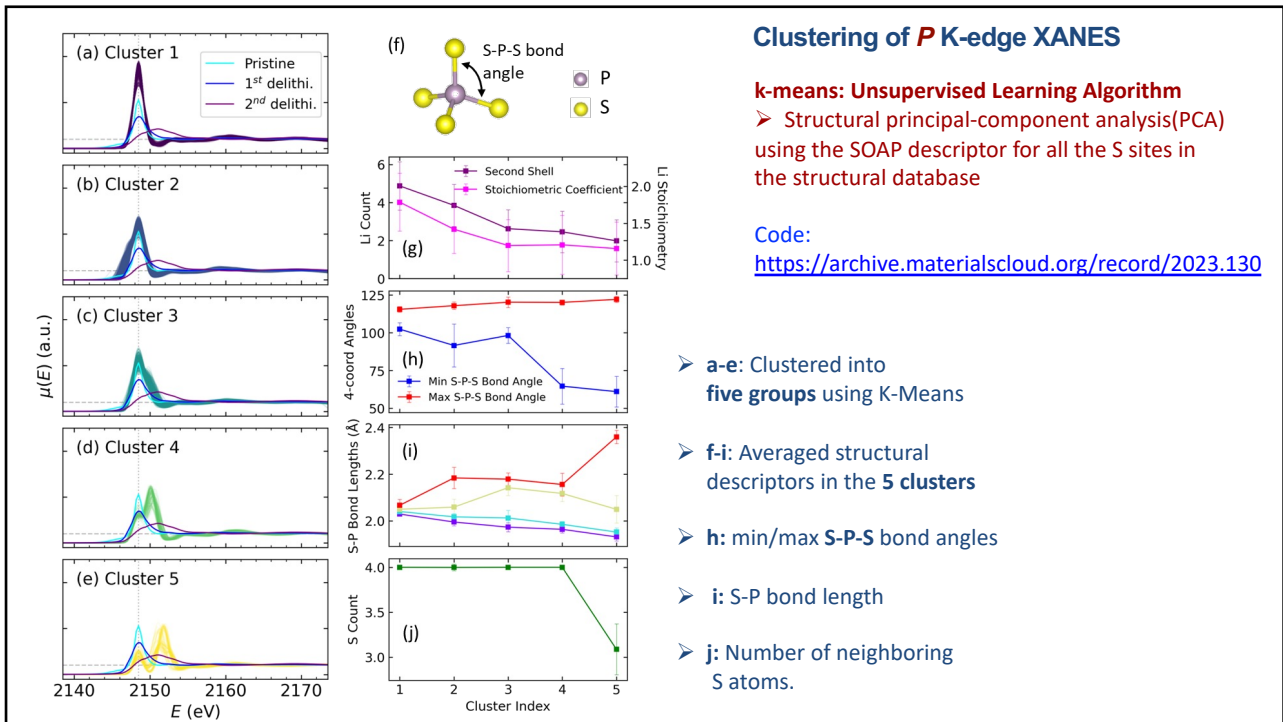
## Validation: Benchmark of the XAS (DFT) simulations

- Simulated spectra reproduce the main features in reference systems



Guo, Carbone, N. Artrith, Urban, Lu, et al. *Nature Scientific Data* 10, 349 (2023)


70



71

## β-Li<sub>3</sub>PS<sub>4</sub> Structures are Stable Conductors

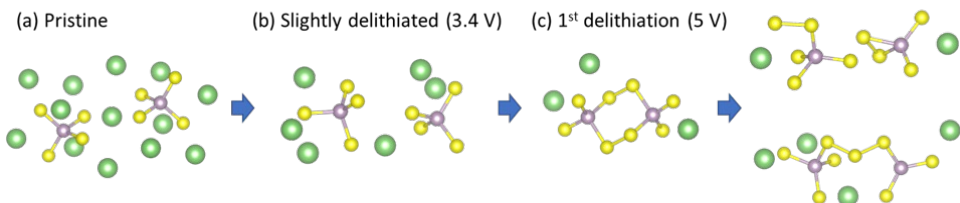
**Chuntian Cao**



➤ **Atomic Insights into the Oxidative Degradation Mechanisms of Sulfide Solid Electrolytes**

<https://github.com/AI-multimodal/Lightshow>

(a) Pristine      (b) Slightly delithiated (3.4 V)      (c) 1<sup>st</sup> delithiation (5 V)      (d) 2<sup>nd</sup> delithiation (5 V)



● Li    ● P    ● S

➤ **Unsupervised learning for spectroscopy classification using SOAP descriptors, k-means clustering, and structural PCA analysis of S-site environments in structural databases.**

Chuntian Cao et al., *Cell Rep. Phys. Sci*, 5, 101909 (2024) <https://doi.org/10.1016/j.xcrp.2024.101909>

72

## Summary

Guo, Carbone, N. Artrith, Urban, Lu, *et al.* *Nature Scientific Data* **10**, (2023) 349.  
Cao, Carbone, N. Artrith, Urban, Lu, Wang, *et al.* *Cell Rep. Phys. Sci.* **5**, (2024) 101909.

**(a)**

**Data-driven Spectral Interpretation**

Li<sub>3-x</sub>PS<sub>4</sub> structure database

X-ray absorption spectroscopy database

**Experimental Spectra**

**Atomic-scale delithiation process**

● Li   ● P   ● S

**(b)**

**Automated Workflow**

LPS Structures Database  
↓  
Structural Analysis  
↓  
Input Generation  
↓  
GS Energy Calculation   XAS Simulation  
↓  
Post Processing  
↓  
Simulated XAS Database  
↕  
Experimental XAS Spectra   Spectral Interpretation

<https://github.com/atomisticnet/xas-tools>

➤ **ML-Assisted XAS Analysis Combines Structural & Spectral Databases**

73

## Practical: SOAP descriptor

[https://colab.research.google.com/github/atomisticnet/aenet-gpr/blob/main/tutorial/tutorial\\_3\\_Li-EC.ipynb](https://colab.research.google.com/github/atomisticnet/aenet-gpr/blob/main/tutorial/tutorial_3_Li-EC.ipynb)

tutorial\_3\_Li-EC.ipynb

File Edit View Insert Runtime Tools Help

Q Commands + Code + Text ▶ Run all Copy to Drive

Connect

**Note on performance**

In section 3-1, the use of the SOAP descriptor in GPR can be computationally intensive.

On **high-performance computing (HPC)** or **GPU** systems, the SOAP-based GPR runs significantly faster.

In this Colab environment, you can switch to a **GPU** by going to

Runtime → Change runtime type → Hardware accelerator → T4 GPU

```

try:
    import torch
    print("successfully imported torch")
    print(torch.__version__)
except ImportError:
    !! pip install torch torchvision torchaudio --user --index-url https://download.pytorch.org/whl/cpu
    print("completed installing torch")

```

successfully imported torch  
2.6.0+cu124

**Homework: SOAP descriptor tutorials**

[GPR Tutorial: H<sub>2</sub>](#)

[GPR Tutorial: EC-EC](#)

[GPR Tutorial: Li-EC](#)

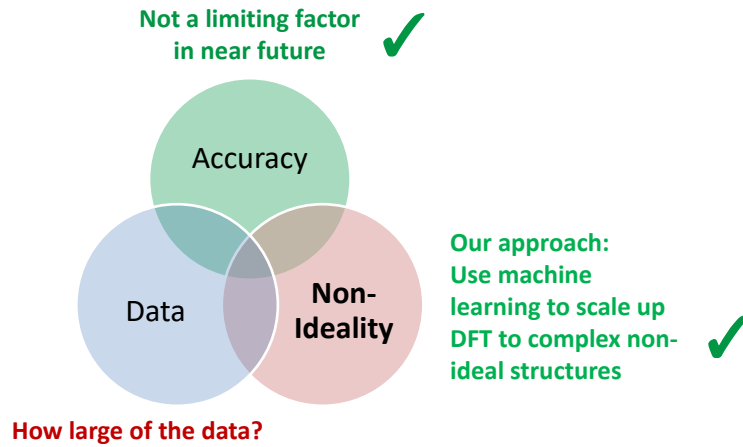
**aenet-gpr**

IW Yeu, A Urban, N Artrith *et al.*, *npj Comput. Mater.* **11**, 156 (2025)

GPR-ANN and GPR-aided NEB Code: <https://github.com/atomisticnet/aenet-gpr>

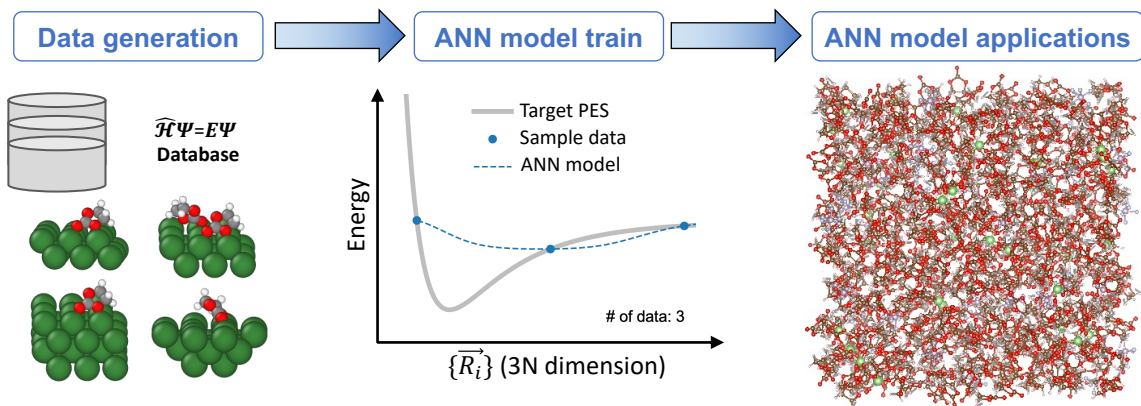
74

## What's Holding Back Computational Materials Discovery?



75

## Practical MLIP applications are hindered mainly by three obstacles



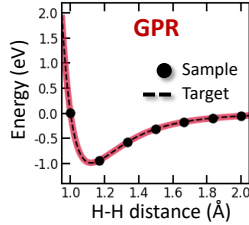
### Problems

- **(Database)** High-fidelity data is lacking
- **(Sampling)** Data generation is arbitrary and inefficient
- **(Cost)** Training is not scalable

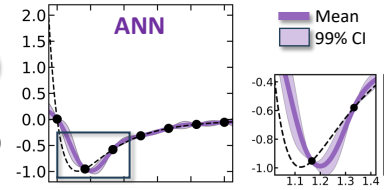
76

### GPR vs ANN: Which is better MLIP?

- **GPR** excels for interpolation in low-data regimes
- **GPR** uncertainty can be used for active sampling
- But **GPR's** scalability is limited



Train data: 7  
 Test data: 200



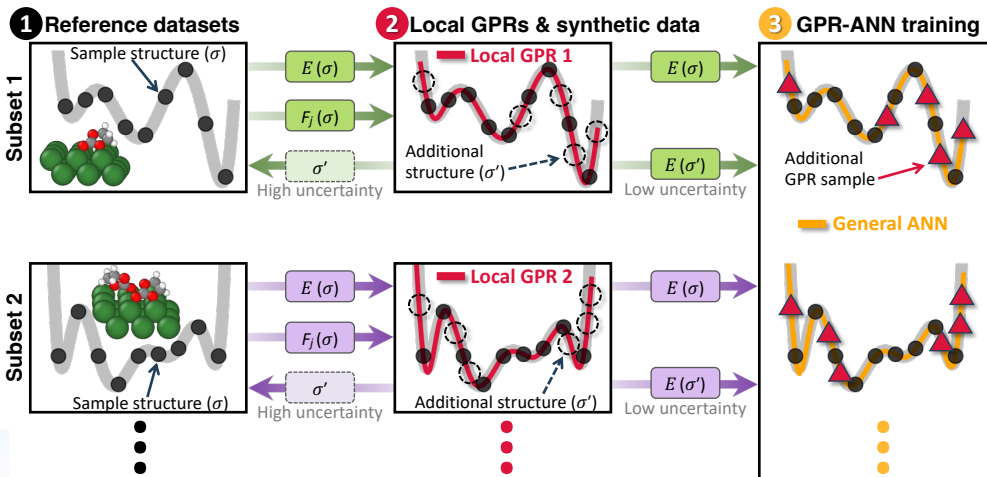
	<b>GPR</b> Non-parametric kernel-based	<b>ANN</b> Parametric model defined by {w}
<b>Model uncertainty</b>	O	X
<b>Data requirements</b>	Small	Large
<b>Training cost</b>	$\mathcal{O}(N^{2\sim 3})$	$\mathcal{O}(N^{1\sim 2})$
<b>Transferability to large-scale</b>	Not Great	Good

77

### GPR-ANN: Combine the Best of Both Worlds

**GPR** for low-data interpolation, uncertainty-based data refinement

**ANN** obtained from efficient energy training, final model for large-scale evaluation



InWon Yue

**aenet-gpr**

<https://github.com/atomicnet/aenet-gpr>

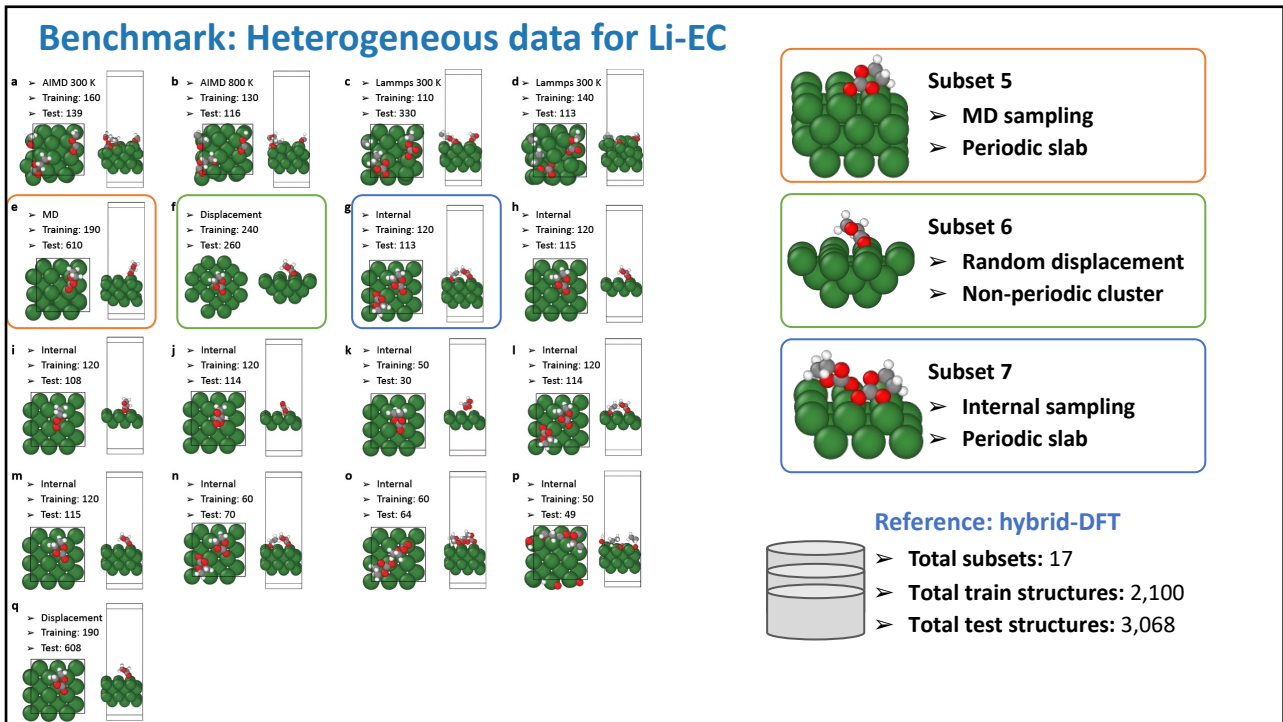
Yeu<sup>†</sup> et al., *npj Comput. Mater.* 11, 156, (2025)

**Atomic Energy Network**

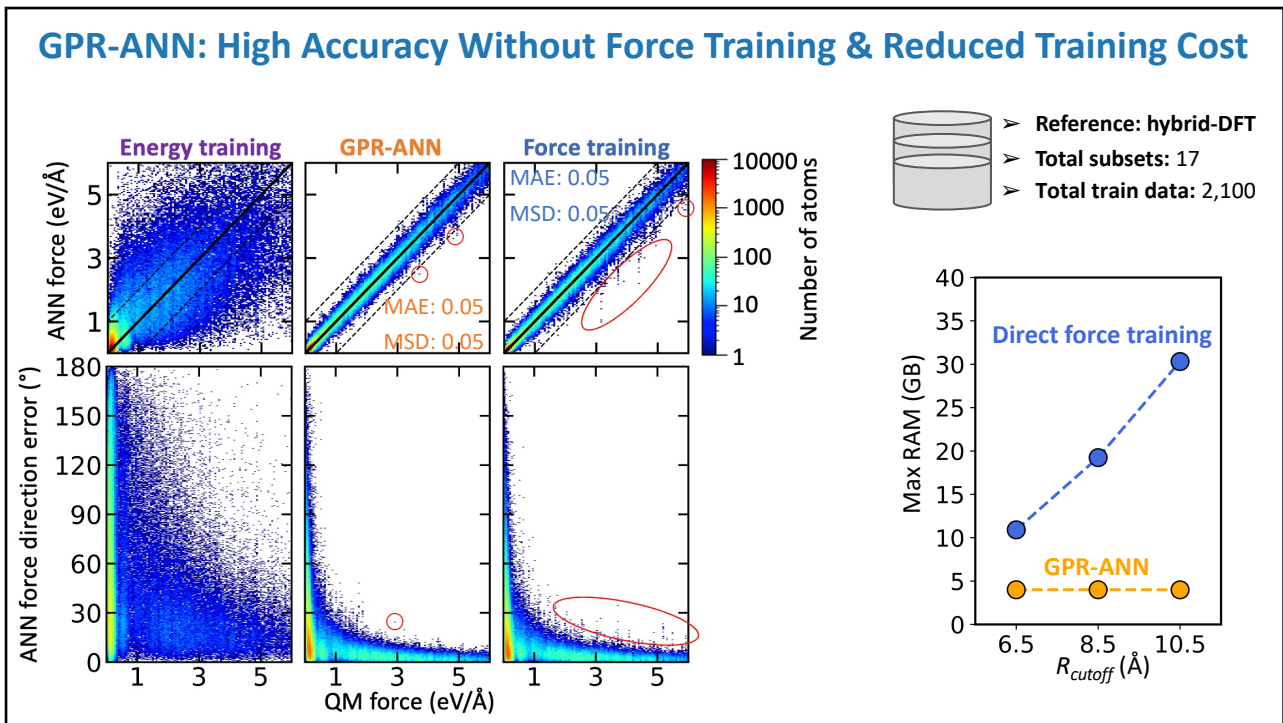
**aenet**

<https://github.com/atomicnet/aenet>

78

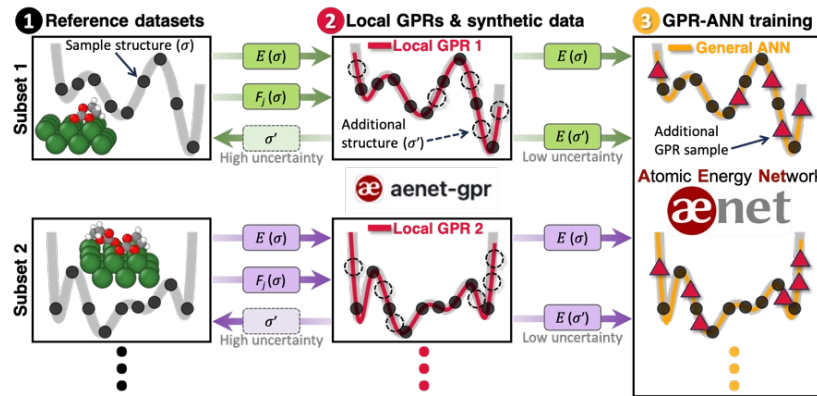


79



80

### Solved the common issues!



#### Problems solved

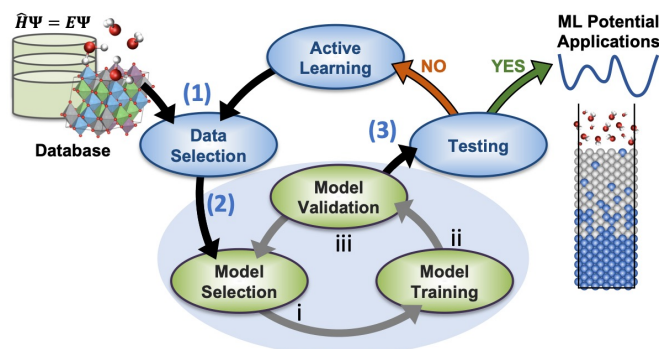
- > **(Database)** Reference data requirement is reduced
- > **(Sampling)** Uncertainty-aware surrogate guides targeted sampling
- > **(Cost)** Training is scalable



InWon Yue

81

### Can we systematically sample missing data that greatly improve MLIPs?



A.M. Miksch, T. Morawietz, J. Kästner, A. Urban, N. Artrith\*,  
*Machine Learning: Science and Technology*, 2, 031001 (2021), DOI: <https://doi.org/10.1088/2632-2153/abfd96>.

Code: <https://github.com/atomisticnet/MLP-beginners-guide>

- **(Model Selction)** kernel based GPR + neural networks
- **(Model Training)** direct force training → indirect force training (distill force into energy)
- **(Active learning)** Query of committee → GPR uncertainty

82

## Practical: SOAP descriptor

[https://colab.research.google.com/github/atomisticnet/aenet-gpr/blob/main/tutorial/tutorial\\_3\\_Li-EC.ipynb](https://colab.research.google.com/github/atomisticnet/aenet-gpr/blob/main/tutorial/tutorial_3_Li-EC.ipynb)

**Note on performance**

In section 3-1, the use of the SOAP descriptor in GPR can be computationally intensive. On high-performance computing (HPC) or GPU systems, the SOAP-based GPR runs significantly faster. In this Colab environment, you can switch to a GPU by going to Runtime → Change runtime type → Hardware accelerator → T4 GPU

```
try:
    import torch
    print("successfully imported torch")
    print(torch.__version__)
except ImportError:
    ! pip install torch torchvision torchaudio --user --index-url https://download.pytorch.org/whl/cpu
    print("completed installing torch")
```

successfully imported torch  
2.6.0+cu124

**aenet-gpr**

**Homework: SOAP descriptor tutorials**

[GPR Tutorial: H<sub>2</sub>](#)  
[GPR Tutorial: EC-EC](#)  
[GPR Tutorial: Li-EC](#)

IW Yeu, A Urban, N Artrith et al., *npj Comput. Mater.* **11**, 156 (2025)  
 GPR-ANN and GPR-aided NEB Code: <https://github.com/atomisticnet/aenet-gpr>

83

## Acknowledgements: Collaborators and Facilities

### Thank You for Your Attention

**Chuntian Cao**

**Matt Carbone**

**Qian Wang**

**Haoyue Guo**

DOE VTO: DE-SC0012704 (2020-2024)

**aenet**

**pymatgen**

84

## Acknowledgements: Collaborators and Facilities



Thank You for Your Attention



Sharing science,  
*shaping tomorrow*

Advanced Research Center  
Chemical Building Blocks Consortium

Dutch National Supercomputer Snellius



Dutch Sector Plan Project

