



# Machine Learning Interatomic Potentials I

### Ioan-Bogdan Magdău ioan.magdau@ncl.ac.uk

**CAMMLS Spring School 2025** 

### Few words about myself.

ត៍ពីពី

fitt

i ii i

F

Lecturer in Computational Chemistry (3<sup>rd</sup> April 2023) Newcastle University

### Few words about myself.

Modelling Excited States Dynamics





Adam O'Hare



Julien Eng Ashley

Leon

Foundational MLIPs for complex interfaces

Samuel Niblett

Panagiotis Kourtis

#### Modelling ion transport in electrolytes



Zitong Wu

Andrea Gardin

### Lecture Outline

- Brief Intro to MLIPs: Context and Overview
- Anatomy of a potential: locality, EOs, forces
- Atomic Descriptors: Symmetry, Smoothness, Completeness
- Models Architectures: Linear, Kernel, MPNNs (MACE)
- Fitting and testing MLIPs

### Multiscale Molecular Modelling - the Battle of Approximations



# Force-field, Potential Energy, Molecular Dynamics



Molecular Dynamics is the method we use to explore the potential energy surface of the nuclei

We obtain forces as the gradient of the potential energy surface:

$$F = -\frac{dE(x)}{dx}$$

We integrate Newton's second low of motion to update positions  $d^2 \mathbf{r}$ 

$$\mathbf{F}_i = m_i \mathbf{a}_i = m \frac{a \mathbf{r}_i}{dt^2}$$

Solve very **expensive Schrodinger Equation** to obtain energies, forces Update positions **very slightly** (1.0 femtosecond at a time)



**Discard** old Schrodinger solution, solve again for very similar geometry.

# Machine Learning Interatomic Potentials (MLIPs)

#### MLIPs learn the mapping from **geometry** to **energy, forces** Instead of using simple functions (MM), they use **universal approximators**



Can speed up the dynamics significantly! **1000-1M X** 

# MLIPs generalize Molecular Mechanics (MM) Forcefields

The PES: 
$$E = E(r_1, r_2, ...)$$

**In QM**, PES from solving the Schrodinger Equation (Born-Oppenheimer approximation)

In MM, PES using empirical parameters: bond, angles, charges (generally non-reactive)

In MLIPS, PES interpolated from know geometries to new geometries (capture chemical reactions)



#### Not all QM methods are created equal:



MLIP cost is constant regardless of accuracy of PES! Training is more expensive of course!

### **Body-order in Molecular Mechanics**



# Is There Anything Missing?



# MLIPs generalize MM Force Fields

**a** Conventional *N*-body terms



Bonds (2-body) Angles (3-body) Dihedrals / Impropers (4-body) enumerated from the outset

**b** General 3-body descriptor



all Pairs within cutoff (2-body) all Triangles (3-body) all Tetrahedra (4-body) ... determined on the fly: inferred from data

Maximum Body-order of MLIP:

$$\frac{\partial^N P}{\partial x_{i_1}, \dots, \partial x_{i_N}} \neq 0$$

### General Anatomy of the MLIP



**Sort-range ML** models rely on energy decomposition and locality:

• total energy can be written as a sum of atomic energies:



**Sort-range ML** models rely on energy decomposition and locality:

• total energy can be written as a sum of atomic energies:

$$E_{tot} = \sum_{i} \varepsilon_{i}$$

• atomic energies are a function of local environment (within a cutoff)



 $\varepsilon_i = \varepsilon_i (r < R_{cut})$ 

#### Central Concepts in MLIPs: Isolated energies



This ensures transferability across chemical compounds and a physically meaningful baseline in the absence of data for chemical reactions.

### Central Concept in MLIPs: Calculating forces

• forces are derivatives of **total** energy:

$$F_{i} = \frac{\partial E_{tot}}{\partial r_{i}} = \sum_{i} \frac{\partial \varepsilon_{i}}{\partial r_{i}} + \sum_{j \neq i} \frac{\partial \varepsilon_{j}}{\partial r_{i}}$$

$$F_{i} = \sum_{i} \frac{\partial \varepsilon_{i}}{\partial r_{i}} (r < R_{cut}) + \sum_{j \neq i} \frac{\partial \varepsilon_{j}}{\partial r_{i}} (r < 2R_{cut})$$

 $\mathcal{E}_i$ 

 $\mathcal{E}_{i}$ 

This ensures:

- energy conservation
- force equivariance
- extends effective interaction

# **Representation and Atomic Descriptors**

Use Atomic Descriptors to represent Structure/Geometry



#### Large and growing family tree of atomic descriptors!

#### **Properties**

- symmetry: translational, rotational, permutational
- smoothness
- completeness
- correlation order (2, 3, 4-body)

### Descriptor Symmetries: Permutation, Rotation, Translation



PERMUTATION

**Permutation** is a <u>discrete</u> <u>symmetry</u> of the energy!

TRANSLATION

**Translation** is a <u>continuous</u> <u>symmetry</u> of the energy!

Infinitely many environments map to the same energy!

ROTATION

Rotation is also a <u>continuous symmetry</u> of the energy!

Descriptors should be **Invariant** to **Permutations:** 

 $P(x_i) = P(x_{\sigma(i)})$ 

Difficult to `learn`, our descriptors require **Translational Symmetry**:

 $P(x_i) = P(x_i + T)$ 

Descriptors should be **Invariant** to **Rotations**:

 $P(x_i) = P(x_i \times R)$ 

#### **Rotational Equivariance**



What about properties that do change with rotations?

# E.g. Forces, Dipole vectors rotate with the coordinate frame!

We want to avoid learning the transformation, instead we can make atomic descriptors **Equivariant with Rotations**:

 $P(x_i \times R) = P(x_i) \times R$ 

### Question: can we fit equivariant forces with an invariant MLIP?

Let's define a rotated coordinate system (i.e. Rotate molecule in Ovito):

Energy is rotationally invariant:

$$E(x) = E(x')$$

x' = Rx



### Question: can we fit equivariant forces with an invariant MLIP?

Let's define a rotated coordinate system (i.e. Rotate molecule in Ovito):

Energy is rotationally invariant:

$$E(x) = E(x')$$

x' = Rx

Forces in the original and rotated frame:

The relation between the forces:

$$\frac{\partial E(x)}{\partial x} = \frac{\partial E(x)}{\partial x'} \frac{\partial x'}{\partial x} = \frac{\partial E(x')}{\partial x'} R$$
$$F(x) = RF(x')$$

 $F(x) = -\frac{\partial E(x)}{\partial x}$   $F(x') = -\frac{\partial E(x')}{\partial x'}$ 

YES! If we obtain forces as the **gradient of an invariant energy**, we ensure they are equivariant! This will work for all vectors that can be expressed as gradients of a scalar field.



# Other Properties of a 'Good' Descriptor: Completeness



#### **Descriptor Space** >> **Geometry Space**

- **many** P **do not** map onto real  $x_i$
- **all**  $x_i$  have to map onto a P
- **different**  $x_i$  must map onto different P





Bijection (One-to-One and Onto)



### Other Properties of a 'Good' Descriptor: Smoothness





descriptors need to change smoothly with geometry:

#### energy is smooth!

so should their arbitrary derivatives:

forces, Hessians are smooth!

### Forces and energy conservation

The total energy must be conserved:

$$E = \frac{1}{2}mv^2 + U = const$$

Which means its time derivative must zero:

$$\frac{dE}{dt} = mv \cdot \frac{dv}{dt} + \frac{dU}{dt} = F \cdot v + \frac{\partial U}{\partial x} \cdot \frac{dx}{dt} = F \cdot v + \frac{\partial U}{\partial r} \cdot v = 0$$

It follows that forces must obtained as the gradient of the total energy:

$$F = -\frac{\partial U}{\partial r}$$

### Examples: SOAP

Atomic density **centered** on atom *i*:

$$\rho^{i}(\boldsymbol{r}) = \sum_{j} \exp\left(\frac{-|\boldsymbol{r} - \boldsymbol{r}_{ij}|^{2}}{2\sigma_{atom}^{2}}\right) f_{cut}(|\boldsymbol{r}_{ij}| < R_{cut})$$

permutational invariance
 translational invariance

Expand into radial / angular basis:

$$\rho^{i\alpha}(\boldsymbol{r}) = \sum_{nlm} c_{nlm}^{i\alpha} R_n(r) Y_{lm}(\hat{r})$$

**Power spectrum** (the 'feature' vector):

$$P = p_{nn'l}^{i} = \frac{1}{\sqrt{2l+1}} \sum_{m} c_{nlm}^{i} c_{n'lm}^{i}$$

rotational invariance

# SOAP is a 3-body Descriptor

#### <u>3-body Degeneracy:</u>

# Two different environments (4 neighbors): same SOAP SOAP is incomplete!



body order ↔ minimum neighborhood degeneracy!

4-body: minimum 7 neighbors proven!

ACE: arbitrary bodyorder expansion!

SN Pozdnyakov, MJ Willatt, AP Bartók, C Ortner, G Csányi, and M Ceriotti PRL **125**, 166001, 2020

Atomic Cluster Expansion (ACE) as a systematic approach:

One-particle basis, atom-centred (translational invariance)

$$\varphi_{nlm,z_iz_j}(\mathbf{r}) = R_{nl,z_iz_j}(r_{ji})Y_l^m(\hat{\mathbf{r}})$$

Evaluate for each neighbour and <u>sum</u> (permutation invariance)

$$A_{iv} = \sum_{j \in \mathcal{N}(i)} \varphi_v(\sigma_j, \sigma_i) \quad \text{(2-body)}$$

<u>Phys. Rev. B **100**</u>, 249901 (2019)

 $arphi_v(m{r_{20}})$ 

 $arphi_v(\dot{r_{10}})$ 

# ACE: Systematic Body-ordered Expansion

Construct the product basis ( $\nu + 1$  body terms):

$$\mathbf{A}_{i,\mathbf{v}} = \prod_{\xi=1}^{\nu} A_{i,v_{\xi}}$$
 ( $u$  + 1 body)



 $\phi(\mathbf{r}_{10})$ 

 $\phi(\mathbf{r}_{20})$ 

Symmetrize the product atomic-basis using the

Clebsch–Gordan coefficients (rotational invariance)

$$oldsymbol{B}_{i,oldsymbol{v}} = \sum_{oldsymbol{w}} \mathcal{C}_{oldsymbol{v}oldsymbol{w}} oldsymbol{A}_{i,oldsymbol{v}}$$

<u>Phys. Rev. B 100,</u> 249901 (2019)

 $\phi(\mathbf{r}_{10})$ 

 $\phi(\mathbf{r}_{20})$ 



A

ACE inspired from cluster-expansions of the energy function:

$$E_{i} = E_{i}^{(0)} + \sum_{v} \tilde{c}_{v}^{(1)} A_{iv} + \sum_{v_{1}v_{2}}^{v_{1} \ge v_{2}} \tilde{c}_{v_{1}v_{2}}^{(2)} A_{iv_{1}} A_{iv_{2}} + \sum_{v_{1}v_{2}v_{3}}^{v_{1} \ge v_{2} \ge v_{3}} \tilde{c}_{v_{1}v_{2}v_{3}}^{(3)} A_{iv_{1}} A_{iv_{2}} A_{iv_{3}} + \dots$$

#### The problem becomes linear:

$$E = \sum_{i, v} c_{i, v} \boldsymbol{B}_{i, v}$$

... and we know how to solve it:

$$L = \sum_{j} \| E(X_{j}) - \mathbf{B}(X_{j})^{T} \mathbf{w} \|^{2}$$
  
CE is a precursor to MACE



# Dealing with chemical species: embedding

#### From categorical space to continuous space

**Embedding vs one-hot encoding representations** 

#### **4-Dimensional Embedding**

dense, low-dimensional representations of data
each word (e.g., "cat," "mat," "on") is represented by a 4dimensional vector with <u>learned values</u>

#### **One-Hot Encoding**

- categorical data as binary vector
- each word is represented by a vector where:
  - One position is "1" (indicating the word).
  - All others are "0."
  - "the"  $\rightarrow$  [0, 0, 0, 0, 1]
  - "cat"  $\rightarrow$  [1, 0, 0, 0, 0]
- high dimensionality: one vector per word
- no notion of similarity between "cat" and "dog."

MACE starts with one-hot encoding and learns an embedding

#### A 4-dimensional embedding











# MPNN+ACE: MACE

Molecules are graphs



- each node is an atom
- messages are interactions between neighbor atoms
- MACE uses ACE basis as initial features





# 1. Embeddings: initializing the graph

Node level



Edge level

Computed once in the beginning.

For each atom/node expand distances/ angles into radial/angular basis.

Encode species one-hot and project into a learnable embedding.

#### Important hypers:

- $R_{cut}$  (r\_max)
- $L_{max}$  (max\_ell)
- N<sub>max</sub> (num\_radial\_basis)

# 2. Interaction: pooling across neighbours



Updated every time the model evaluates the interaction layer (*S* times)

The one-particle basis ( $\phi$ ) is constructed using **e3nn**, and unlike in ACE, it is equivariant (there are  $\eta_1$  ways to construct equivariances)

The atomic basis (A) is obtained by pooling over neighbors (permutational invariance)

Everything is multiplied by additional learnable weights
### 3. Product: updating the node features



•  $\nu$  correlation order

Updated every time the model evaluates the product layer (*S* times)

The atomic basis is tensor-multiplied with itself to form the product basis (main difference to NequIP: short-cut to obtaining higher body orders with fewer layers)

The fully symmetric basis (B) is constructed by contracting the atomic basis (A) using **e3nn**. There are  $\eta_{\nu}$  ways to achieve this.

Finally, the message is obtained by multiplying with learnable weights and the node feature is updated.

#### 4. Readout from each layer



In an S-layer MACE, the readout from the features at layer s is:

Forces obtained via autodiff. Loss in each batch:

$$\mathcal{L} = \frac{\lambda_E}{B} \sum_{b=1}^{B} \left( \frac{E_b - \hat{E}_b}{N_b} \right)^2 + \frac{\lambda_F}{3B} \sum_{b=1}^{B} \sum_{i_b, \alpha=1}^{N_b, 3} \left( -\frac{\partial E_b}{\partial r_{i_b, \alpha}} - \hat{F}_{i_b, \alpha} \right)^2$$

Total field-of view:

 $R_{cut} \times S$ 

Total body-order (within the first )

 $(\nu + 1) \times S + 1$ 

#### Training the Model: Minimizing the Loss

**Input:** atomic environments p**Output:** observables F(p)

Learning: minimize loss on train set  $\{p, F(p)\}$ 

$$\mathcal{L} = \sum_{n=1}^{N_{data}} [F_n - F(p_n)]^2$$

to determine  $w_{ij}$  etc

Data can be different molecules, or different geometries or both ...



### Case study: MLIP for an organic solvent

Li-ion battery







- 6 atomic species: add Li, P, F (quadratic scaling)
- charged long-range effects  $(\sim 1/r)$

#### Case study: MLIP for an organic solvent



# Schrödinger

#### High-Dimensional Neural Network Potential for Liquid Electrolyte Simulations

Steven Dajnowicz\*, Garvit Agarwal\*, James M. Stevenson, Leif D. Jacobson, Farhad Ramezanghorbani, Karl Leswing, Richard A. Friesner, Mathew D. Halls, and Robert Abel

- training data type: molecular clusters
- Neural Network + long-range
- we will fit a MACE model on a subset of this data



#### Intra-/Inter- Decomposition as a Test









$$L = \left\| F_{total}^{DFT} - F_{total}^{ML} \right\|^{2}$$
$$L = \left\| \left( F_{intra}^{DFT} - F_{intra}^{ML} \right) + \left( F_{inter}^{DFT} - F_{inter}^{ML} \right) \right\|^{2}$$

dominates loss

drives dynamics

#### Intra/Inter test:

- split liquid configs into molecular
- recompute molecules (same geometry) in vacuum: intra contributions
- compute inter- as difference between total and intra-

#### Trans-/Rot-/Vib- Decomposition of Forces







## Machine Learning Interatomic Potentials II

Ioan-Bogdan Magdău ioan.magdau@ncl.ac.uk

**CAMMLS Spring School 2025** 

#### Lecture Outline

- Iterative Training: improving stability and accuracy
- Error estimation: committee models
- Active learning: unsupervised iterative training
- Foundational models: out-of-the-box MLIPs
- Fine-tuning on new data and new labels

#### MACE-MP-0: a First Generation Foundational Model

#### Batatia, I et al, arXiv preprint, arXiv:2401.00096 (2023)



to fully appreciate foundational models

let's first look at how we used to fit MLIPs a few years ago

these concepts: **iterative training**, **active learning** will remain relevant, hopefully needed less often

### Stability and Accuracy in Molecular Dynamics

Stability isn't everything, but without it we can't do anything

Most of us are interested in applications that involve MD beyond AIMD.



### Test Error, MD Stability and MD Accuracy

#### 1. Root Mean Square Errors (RMSE)



### Test Error, MD Stability and MD Accuracy

**1. Root Mean Square Errors (RMSE)** 

low test errors do no guarantee stability!

**2. MD Stable Potentials** 

stability does not guarantee MD accuracy!

3. Prediction accuracy: thermodynamics and kinetics

Goal: correct probability distribution! difficult to check



### GAP Potential for the Organic Electrolyte Solvent

Li-ion battery





3 atomic species: H, C, O LP57 solvent: 33% EC : 67% EMC neutral molecules ( $\sim 1/r^3$ ) simple GAP model on total energies **Initial SEI Evolved SEI** Dimethyl carbonate, DMC Ethylmethyl carbonate, Diethyl carbonate, EMC DEC Anod Anod Carbonate Solvents Sraphite LIPF LIOR Graphi Ethylene carbonate, Propylene carbonate, Carbonate Solvents EC PC LIPE Li<sup>+</sup> **Electrolyte-Interface** +Salt:

- 6 atomic species: add Li, P, F (quadratic scaling)
- charged long-range effects  $(\sim 1/r)$

### Training on Fixed / Pre-generated Data



**OPLS** sampling: decorrelated samples, low computational cost DFT calculations **PBE+G06** 12-molecule configurations Wide range of **densities and temperatures**, maintaining **diffusive** behavior



### The Density Problem: Bubble Formation



- <u>well-behaved molecules</u>
- liquid density collapse
- bubbles formation



- Coverage of phase space?
- Transferability across molecular compositions?
- Difference in scale and dimensionality of Intra- / Inter-?

#### Iterative Training: model learns from mistakes



### Improving the models – Iterative Training



### **Iterative Training**



50 OPLS configs Iterations with **12-molecules** (accessible to ab initio) Multiple Temperatures: 500K → 1200K (extrapolation)
Multiple Pressure: 1bar → 25kbar
Proxy for folding back: density instabilities
Gen 4: 50 OPLS + 35 GAP-MD configs: stable densities!

### Testing on Larger System (48 molecules)



#### **Target Composition (33:67)**

- GAP densities stable at all temperatures!
- Densities not well reproduced at 400K
- Same density at 400K and 500K

#### Transferability to other Molecular Compositions



#### poor transferability between compositions!

#### Transferability to other Molecular Compositions





#### Scan inter-molecular PES:

- thermalized configurations from GAP-MD
- > energy binding curve with frozen molecules











#### <u>GAP</u>

- describes bottom of the well
- stable densities, but values too high



#### <u>GAP</u>

- describes bottom of the well
- stable densities, but values too high

#### GAP + Single Molecules (SM)

- reproduces the free molecule limit
- density values still too high







#### <u>GAP</u>

- describes bottom of the well
- stable densities, but values too high

#### GAP + Single Molecules (SM)

- reproduces the free molecule limit
- density values still too high

#### GAP + SM + Volume Scans (VS)

- reproduces the entire E(V) curve
- captures the strong repulsion limit
- still some non-smooth behavior

### 2 Postdoc Years: GAP-MD vs AIMD on PBE-D3



### Iterative Training: How to best choose 'failed' configs?



-1.0

-0.5

Active Learning is a way to automate the Iterative Training by choosing new configs automatically, for example based on true error.

In practice, hard to know **true error**, because it can be very expensive to compute!

0.0

Х

0.5

1.0

#### **Uncertainty Prediction – Bayesian View**

Linear/Kernel models:

 $\mathbf{y} = \mathbf{\Psi}\mathbf{c} + \boldsymbol{\epsilon}$ 

 $Posterior = \frac{Likelihood * Prior}{Normalization}$ 

We can use Bayesian Regression to compute **both mean and variance**:

 $\overline{\mathbf{c}} = \lambda \mathbf{\Sigma} \mathbf{\Psi}^T \mathbf{y}$ 

$$\sigma = \sqrt{rac{1}{\lambda} + oldsymbol{B}^T oldsymbol{\Sigma} oldsymbol{B}}$$

#### Which configs to add back to training?



A good read: https://towardsdatascience.com/introduction-to-bayesian-linear-regression-e66e60791ea7

#### **Uncertainty Prediction – Committees**

In practice, analytical error estimation can be very expensive

 $O(N_{\rm basis}^2)$ 

Use committee models: 'basket' of models with different solutions

$$\tilde{\sigma} = \sqrt{\frac{1}{\lambda} + \frac{1}{K} \sum_{k=1}^{K} \left( E^{k} - \overline{E} \right)^{2}},$$

Advantage: committees work for highly-nonlinear Neural Networks (NN) as well: MACE



### Uncertainty Prediction guides Active Learning

## In practice, hard to know true error, assume it is predicted by the committee disagreement



### Uncertainty Prediction guides Active Learning

## In practice, hard to know true error, assume it is predicted by the committee disagreement



### Neural Network (MACE) Committee Models

more

diverse

# Strategies for building NN committees:

- $\succ$  different training seeds  $\rightarrow$  solutions
- different representative data subsets
- slightly different hyperparameters

MACE result on EC/EMC cluster configs labeled with XTB (semiempirical/fast)

Shows convincing model disagreement (uncertainty) when true error explodes

Good correlation between uncertainty and true error!



### **Beyond Active Learning**



Often it can take a long time until we find a failed config (useful new data). This is can amount to a lot of wasted computing time.

This problem becomes more pronounced as the models get better. Can we find a way to speed up the failure?

### Hyper-active Learning (HAL)



$$E = \sum_i E_i = \mathbf{c} \cdot \mathbf{B}$$

$$\tilde{\sigma} = \sqrt{\frac{1}{\lambda} + \frac{1}{K} \sum_{k=1}^{K} \left( E^{k} - \overline{E} \right)^{2}}$$

HAL uses the uncertainty prediction to bias the potential energy surface (PES)

This pushes the simulation to areas of high uncertainty

#### npj computational materials

Explore content v About the journal v Publish with us v

nature > npj computational materials > articles > article

#### Article Open access Published: 13 September 2023

## Hyperactive learning for data-driven interatomic potentials

Cas van der Oord <sup>™</sup>, Matthias Sachs, Dávid Péter Kovács, Christoph Ortner & Gábor Csányi

npj Computational Materials 9, Article number: 168 (2023) Cite this article

### Hyper-active Learning



HAL vs simple Active Learning (AL)

HAL requires fewer iterative training cycles

HAL requires orders of magnitude less dynamics

HAL uses the uncertainty prediction to bias the potential energy surface (PES)

This pushes the simulation to areas of high uncertainty

 $E_{HAL} := E - \tau \sigma$ 

npj computational materials

Explore content v About the journal v Publish with us v

nature > npj computational materials > articles > article

#### Article Open access Published: 13 September 2023

Hyperactive learning for data-driven interatomic potentials

Cas van der Oord 🖾, Matthias Sachs, Dávid Péter Kovács, Christoph Ortner & Gábor Csányi

npj Computational Materials 9, Article number: 168 (2023) Cite this article
### The Advent of Foundational Models

Training a MLIP for every application/system is not sustainable!

**Big Question:** 

Can we train one big foundational model that works for all systems? Message-passing Atomic Cluster Expansion (MACE)

GAP MACE Inter Energy (meV/atom) Inter Energy (meV/atom) -10Inter-energy -15-15MACE (2x6A) 52- 52-6AP (TS) 05--25 RMSE = 1.1351 RRMSE = 0.1770 RMSE = 0.4020 RRMSE = 0.0627 -30 -30 -25 -20 -15 -10 -30 -25 -20 -15 -30 DFT DFT Inter Forces (eV/A) Inter Forces (eV/A) 1.0 1.0 nter-forces 0.5 0.5 MACE (2×6A) GAP (TS) 0.0 0.0 -0.5 -0.5-1.0RMSE = 0.0644RMSE = 0.0190 RRMSE = 0.1604RMSE = 0.5430DFT DFT



- species embedding →
  constant scaling with
  elements, transferability
- body-ordered expansion → accuracy, large capacity
- regularity/smoothness  $\rightarrow$  stable MD
- availability of Big Data

### MACE-MP-0: a First Generation Foundational Model

- Materials Project (Berkeley Lab): 89 elements, ~150K materials (90% < 70 atoms)</li>
- MPtraj: ~1.5M configs + structural relaxations
- XC-functional: PBE+U





https://next-gen.materialsproject.org/

MACE: trained on MPtraj, no iterative training

Model 🛈	CPS ↑	F1 ↑	DAF ↑	Prec ↑	Acc ↑	MAE ↓	R² ↑	κ <sub>srme</sub> ↓	RMSD ↓	Training Set	Params	Targets	Date Adde
eSEN-30M-OAM	0.896	0.925	6.069	0.928	0.977	0.018	0.866	0.1704	0.0096	6.6M (113M) (OMat24+MPtrj+sAlex)	30.2M	$EFS_{G}$	2025-03-1
SevenNet-MF-ompa	0.849	0.901	5.825	0.89	0.969	0.021	0.867	0.317	0.0115	6.6M (113M) (OMat24+sAlex+MPtrj)	25.7M	$EFS_{G}$	2025-03-1
GRACE-2L-OAM	0.841	0.88	5.774	0.883	0.963	0.023	0.862	0.294	0.0121	6.6M (113M) (OMat24+sAlex+MPtrj)	12.6M	$EFS_{G}$	2025-02-0
eSEN-30M-MP	0.800	0.831	5.26	0.804	0.946	0.033	0.822	0.3398	0.0142	146k (1.58м) <b>(MPtrj)</b>	30.1M	$EFS_G$	2025-03-1
MACE-MPA-0	0.796	0.852	5.582	0.853	0.954	0.028	0.842	0.412	0.0143	3.37M (12M) (MPtrj+sAlex)	9.06M	$EFS_G$	2024-12-0
DPA3-v2-OpenLAM	0.765	0.89	5.747	0.879	0.966	0.022	0.869	0.687	0.0127	163M (OpenLAM)	7.02M	$EFS_{G}$	2025-03-14
GRACE-1L-OAM	0.762	0.824	5.255	0.803	0.944	0.031	0.842	0.516	0.0139	6.6M (113M) (OMat24+sAlex+MPtrj)	3.45M	$EFS_{G}$	2025-02-0
MatterSim v1 5M	0.716	0.862	5.852	0.895	0.959	0.024	0.863	0.574	0.0733	17M (MatterSim)	4.55M	$EFS_{G}$	2024-12-16
SevenNet-I3i5	0.709	0.76	4.629	0.708	0.92	0.044	0.776	0.55	0.0182	146k (1.58м) <b>(MPtrj)</b>	1.17M	$EFS_G$	2024-12-10
MatRIS v0.5.0 MPtrj	0.680	0.809	5.049	0.772	0.938	0.037	0.803	0.861	0.0156	146k (1.58м) <b>(MPtrj)</b>	5.83M	$EFS_{G}M$	2025-03-1
GRACE-2L-MPtrj	0.678	0.691	4.163	0.636	0.896	0.052	0.741	0.525	0.0186	146k (1.58м) <b>(MPtrj)</b>	15.3M	$EFS_{G}$	2024-11-21
DPA3-v2-MPtrj	0.647	0.786	4.822	0.737	0.929	0.039	0.804	0.959	0.0164	146k (1.58м) <b>(MPtrj)</b>	4.92M	$EFS_{G}$	2025-03-14
MACE-MP-0	0.640	0.669	3.777	0.577	0.878	0.057	0.697	0.647	0.0194	146k (1.58м) <b>(MPtrj)</b>	4.69M	$EFS_{G}$	2023-07-14
AlphaNet-MPTrj	0.562	0.799	4.863	0.743	0.933	0.041	0.745	1.31	0.0227	146k (1.58м) <b>(MPtrj)</b>	16.2M	$EFS_{G}$	2025-03-0
eqV2 M	0.558	0.917	6.047	0.924	0.975	0.02	0.848	1.771	0.0138	3.37M (102M) (OMat24+MPtrj)	86.6M	$EFS_D$	2024-10-18
ORB v2	0.540	0.88	6.041	0.924	0.965	0.028	0.824	1.732	0.016	3.25M (32.1M) (MPtrj+Alex)	25.2M	EFS <sub>D</sub>	2024-10-11
eqV2 S DeNS	0.526	0.815	5.042	0.771	0.941	0.036	0.788	1.676	0.0138	146k (1.58м) <b>(MPtrj)</b>	31.2M	$EFS_D$	2024-10-18
ORB v2 MPtrj	0.476	0.765	4.702	0.719	0.922	0.045	0.756	1.725	0.0185	146k (1.58м) <b>(MPtrj)</b>	25.2M	EFS <sub>D</sub>	2024-10-14
M3GNet	0.430	0.569	2.882	0.441	0.813	0.075	0.585	1.412	0.0217	62.8k (188k) (MPF)	228k	$EFS_{G}$	2022-09-2
CHGNet	0.391	0.613	3.361	0.514	0.851	0.063	0.689	1.717	0.0216	146k (1.58M) (MPtrj)	413k	EFS <sub>G</sub> M	2023-03-0
GNoME	n/a	0.829	5.523	0.844	0.955	0.035	0.785	n/a	n/a	6M (89M) (GNoME)	16.2M	$EF_{G}$	2024-02-0

### The Emergence of Foundational Models

#### Batatia, I et al, arXiv preprint, arXiv:2401.00096 (2023)

A foundation model for atomistic materials chemistry

Ilyes Batatia<sup>†1</sup>, Philipp Benner<sup>†2</sup>, Yuan Chiang<sup>†3,4</sup>, Alin M. Elena<sup>†17</sup> Dávid P. Kovács<sup>†1</sup>, Janosh Riebesell<sup>†4,13</sup>, Xavier R. Advincula<sup>12,13</sup>, Mark Asta<sup>3,4</sup>, Matthew Avaylon<sup>30</sup>, William J. Baldwin<sup>1</sup>, Fabian Berger<sup>12</sup>, Noam Bernstein<sup>11</sup>, Arghya Bhowmik<sup>25</sup>, Samuel M. Blau<sup>10</sup>, Vlad Cărare<sup>1,13</sup>, James P. Darby<sup>1</sup>, Sandip De<sup>18</sup>, Flaviano Della Pia<sup>12</sup>, Volker L. Deringer<sup>16</sup>, Rokas Elijošius<sup>1</sup>, Zakariya El-Machachi<sup>16</sup>, Fabio Falcioni<sup>31</sup>, Edvin Fako<sup>18</sup>, Andrea C. Ferrari<sup>26</sup>, Annalena Genreith-Schriever<sup>12</sup>, Janine George<sup>2,6</sup>, Rhys E. A. Goodall<sup>15</sup>, Clare P. Grey<sup>12</sup>, Petr Grigorev<sup>27</sup>, Shuang Han<sup>18</sup>, Will Handley<sup>13,19</sup>, Hendrik H. Heenen<sup>9</sup>, Kersti Hermansson<sup>23</sup>, Christian Holm<sup>22</sup>, Stephan Hofmann<sup>1</sup>, Jad Jaafar<sup>1</sup>, Konstantin S. Jakob<sup>9</sup>, Hyunwook Jung<sup>9</sup>, Venkat Kapil<sup>12, 21</sup>, Aaron D. Kaplan<sup>4</sup>, Nima Karimitari<sup>20</sup>, James R. Kermode<sup>28</sup>, Namu Kroupa<sup>13,19,1</sup>, Jolla Kullgren<sup>23</sup>, Matthew C. Kuner<sup>3,4</sup>, Domantas Kuryla<sup>12</sup>, Guoda Liepuoniute<sup>1,26</sup>, Johannes T. Margraf<sup>8</sup>, Ioan-Bogdan Magdău<sup>24</sup>, Angelos Michaelides<sup>12</sup>, J. Harry Moore<sup>1</sup>, Aakash A. Naik<sup>2,6</sup>, Samuel P. Niblett<sup>12</sup>, Sam Walton Norwood<sup>25</sup>, Niamh O'Neill<sup>12,13</sup>, Christoph Ortner<sup>5</sup>, Kristin A. Persson<sup>3,4,7</sup>, Karsten Reuter<sup>9</sup>, Andrew S. Rosen<sup>3,4</sup>, Lars L. Schaaf<sup>1</sup>, Christoph Schran<sup>13</sup>, Benjamin X. Shi<sup>12</sup>, Eric Sivonxay<sup>10</sup>, Tamás K. Stenczel<sup>1</sup>, Viktor Svahn<sup>23</sup>, Christopher Sutton<sup>20</sup>, Thomas D. Swinburne<sup>27</sup>, Jules Tilly<sup>31</sup>, Cas van der Oord<sup>1</sup>, Santiago Vargas<sup>29</sup>, Eszter Varga-Umbrich<sup>1</sup>, Teis Vegge<sup>25</sup>, Martin Vondrák<sup>8,9</sup>, Yangshuai Wang<sup>5</sup>, William C. Witt<sup>14</sup>, Fabian Zills<sup>22</sup>, and Gábor Csánvi<sup>\*1</sup>

- MD-stable out-of-box
- New frontiers in molecular modelling
- vast range of applications: organic, inorganic, interfaces, full devices



### Now possible to compute with ab initio accuracy!

Water - Benzene

#### IR spectra of water



#### Free energy landscape: CO<sub>2</sub> in MOF



#### Formation of SEI in Li-ion battery





NaCl dissolution

#### Million atom perovskite

x - direction tilt

Ethanol - Benzene



### Nudge-Elastic-Band

multistep catalysis on surfaces



#### 0.0-1.50 1.75 2.00 2.25 2.50 2.75 3.00 3.25 C - O distance / Å



Water - Ethanol

Water - Heptane

### Back to the Electrolyte: stability

#### EC/EMC LiPF<sub>6</sub>, 500K NPT, MACE-MP-0 (PBE+D3) EC/EMC, 500K NPT, bespoke MACE model (PBE+D2)



### Out-of-box MACE-MP-0

### no iterative training

simulation is stable at 500K, NPT, 150ps:

- molecules remain intact!
- density is stable!
- electrolyte is flowing, exploring phasespace



### Back to the Electrolyte: Accuracy

0

0



- independent PBE test set ullet
- 0.1-2.5 g/cm<sup>3</sup>, all EC/EMC ٠ LiPF<sub>6</sub> compositions
- largest errors at unphysically\_10 ٠ high-density

#### C. Rigid-molec. Volume Scans

- neat solvent
- full electrolyte (incl. ions) ۲
- surprisingly good ۲ performance
- notice the potential ٠ reproduces the 1/r behavior without formal charges



B. Intra/inter errors for electrolyte

### Full Battery 'baby' Model: can we Break the Simulation?





Cu | H-capped graphite+Li | EC/EMC+LiPF<sub>6</sub> | NMC+Li

# Simulating the Li-ion battery environments

- ~ 1800 atoms
- 9 chemical elements
- 4 different materials: metal Copper, Li-loaded graphite, EC/EMC + LiPF6, Li-loaded (partial) NMC
- 4 different chemical interfaces

all local atomic environments are represented to some extend in MPtraj

Will this simulation be stable with the out-of-the-box model?

### Full Battery 'baby' Model

#### **MLIP** performance:

- ~1000s atoms
- ~100ps a day/GPU



#### **DFT cost/performance**

(rough estimate, linear-scaling DFT can certainly do much better) this system: ~6000 electrons ~single point 2 day / 20 CPUs



## Beginning of SEI Formation? Is the Science Right?



#### New Compounds in the Electrolyte / Interface

- bare NMC (no H): CO<sub>2</sub> forming, some EC, EMC become radical
- hydrogenated NMC (surface O atoms): less CO<sub>2</sub>, more H<sub>2</sub>O forming, mostly EMC become radical



#### **Final MD snapshot**

- gray: products of reacted solvent molecules
- red: oxygen atoms originating from cathode, forming new CO<sub>2</sub> and H<sub>2</sub>O molecules
- solvent molecules chemisorbed on cathode surface





### Fine Tuning: Pretrain on Big Data, Transfer to Small Data



 working assumption: parts of the model are transferable (feature construction)

 pretrain model on available large dataset

 fine-tune on small dataset for specific application

### Fine Tuning Strategies



 freeze parts of the network (same features)

#### Loss Landscape

 unfreeze all network: start solution is closer to minimum on Loss Landscape \_\_\_\_\_



multi-head: one head is the old data (keeps weights constrained), new head is new data

### Fine Tuning in MLIPs: Early Days

• out-of-train systems (here ice lh), same level of theory (PBE)

new level of theory

both out-of-train system and new level of theory (RPA, revPBE-D3)





Harveen Kaur, et al, arXiv preprint arXiv:2405.20217 (2024)

https://github.com/venkatkapil24/fine-tuning-MLPs-ice-polymorphs