# LA-UR-24-26124

**Approved for public release; distribution is unlimited.**

| | |
|---|---|
| **Title:** | Improving Atomistic Simulations With Machine Learning |
| **Author(s):** | Allen, Alice Elisabeth Anastasia |
| **Intended for:** | Seminar at Heidelberg Institute for Theoretical Studies |
| **Issued:** | 2024-06-24 |

# Improving Atomistic Simulations With Machine Learning
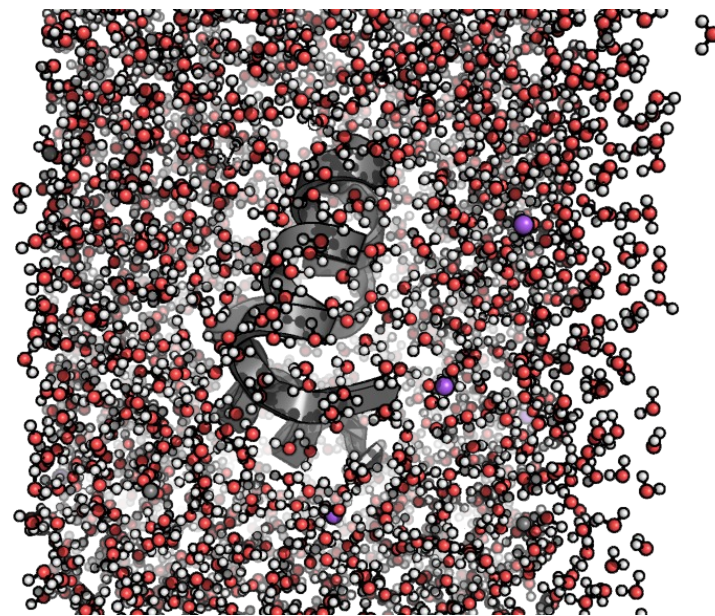
## Alice Allen

# Introduction

What are atomistic simulations and interatomic potentials?

How does machine learning fit in?

How can interatomic potentials be improved?

# Atomistic Simulation

Atomistic simulations can be used to determine material and chemical properties
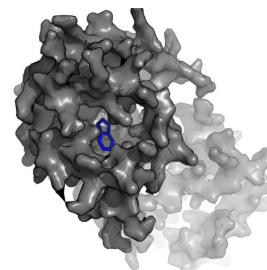
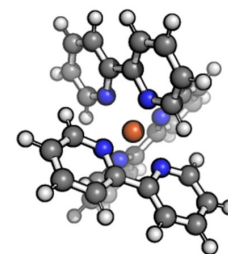Simulating how atoms and molecules behave in a given system

For example:

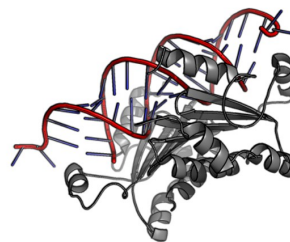 Understand reactive systems

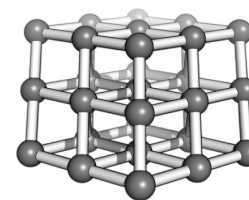 Exploring phases of materials

Protein binding to a drug   Transition Metal Complex

## Atomistic Simulations

Biological    Material

# Interatomic Potentials

Potentials provide forces and energies

Trade off between speed and accuracy

Classical models use well defined functional form



$$E(\mathbf{x}) = \sum_{bonds} k_b(r - r_e)^2 + \sum_{angle} k_\theta(\theta - \theta_e)^2$$

$$+ \sum_{dihedrals} k_\chi(1 + \cos n\chi - \delta) + \sum_{impropers} k_{imp}(\phi - \phi_0)^2$$

$$+ \sum_{\substack{nonbonded \\ pairs}} \left[ \frac{A_{ij}}{r_{ij}^{12}} - \frac{B_{ij}}{r_{ij}^6} \right] + \frac{q_i q_j}{\epsilon_l r_{ij}}$$

# Classical Examples

$$E(\mathbf{x}) = \sum_{bonds} k_b(r - r_e)^2 + \sum_{angle} k_\theta(\theta - \theta_e)^2$$

$$+ \sum_{dihedrals} k_\chi(1 + \cos n\chi - \delta) + \sum_{impropers} k_{imp}(\phi - \phi_0)^2$$

$$+ \sum_{\substack{nonbonded \\ pairs}} \left[ \frac{A_{ij}}{r_{ij}^{12}} - \frac{B_{ij}}{r_{ij}^6} \right] + \frac{q_i q_j}{\epsilon_l r_{ij}}$$
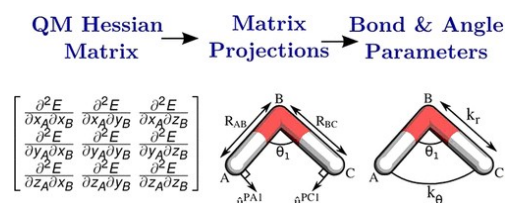
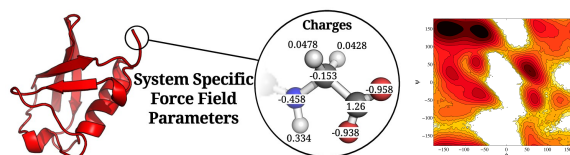# How can classical potentials be improved?

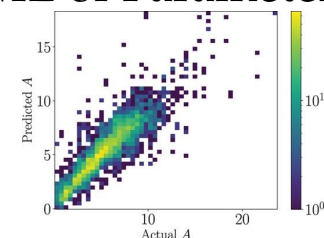## Improve Parameters

### Better bond and angle terms



Harmonic force constants for molecular mechanics force fields via Hessian matrix projection, AEA Allen, MC Payne, DJ Cole, 2018, JCTC

### System specific parameters



Development and validation of the quantum mechanical bespoke protein force field, AEA Allen, MJ Robertson, MC Payne, DJ Cole, 2019, ACS Omega

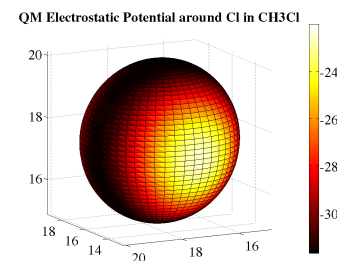### ML of Parameters



Toward transferable empirical valence bonds: Making classical force fields reactive, AEA Allen, Gabor Csanyi;, J. Chem. Phys, 2024

## Improve Functional Form

### Improve description of electrostatics



QUBEKit: Automating the derivation of force field parameters from quantum mechanics, JT Horton, AEA Allen, LS Dodda, DJ Cole, 2019, Journal of chemical information and modeling

# Limitations of Classical Force Fields



QUBEKit: Automating the derivation of force field parameters from quantum mechanics, JT Horton, AEA Allen, LS Dodda, DJ Cole, 2019, Journal of chemical information and modeling

# Machine Learning

# Why do we need Machine Learning Potentials?

Potentials can provide forces and energies

Problem suits machine learning as the potential energy surface is highly complex and requires a flexible functional form

Lots of data can be supplied for the problem

# Machine Learning Potentials

The flexible functional form allows for increased accuracy.

Multiple considerations for these forms of models.
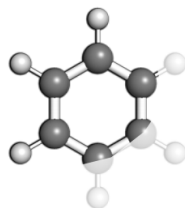
<u>Atom Centred</u>
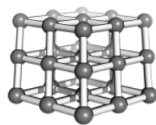
$$E = \sum_{i=1}^{N} E_i$$

Machine Learning Potentials

<u>Cutoffs</u>

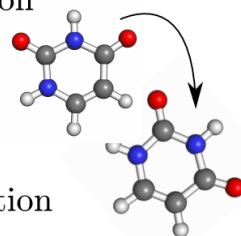<u>Training Data</u>

Automated

Manual Construction
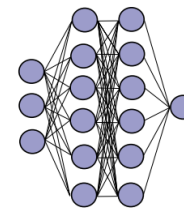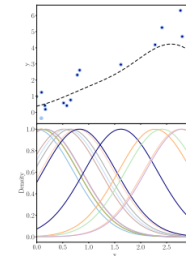
<u>Symmetries</u>

Translation

Rotation

Permutation

Rotate

<u>Model Choice</u>

Kernels

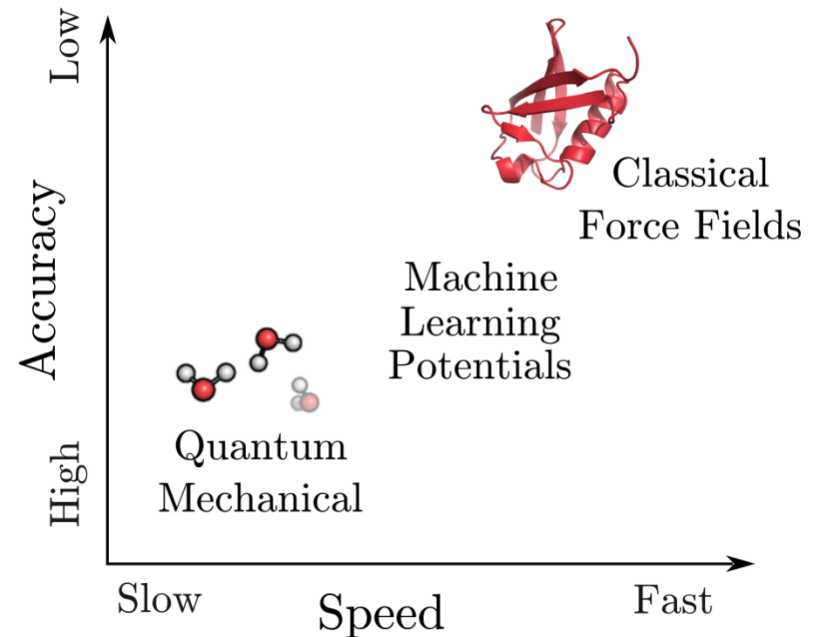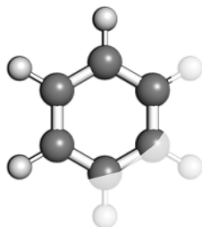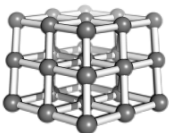Neural Networks

Linear

$$Loss = \sum_{i=0}^{N} (y_i - \sum_{j=0}^{M} x_{ij} W_j)^2$$

# Improving ML Models

## Architecture

Atomic permutationally invariant polynomials for fitting molecular force fields. Alice E. A. Allen, Geneviève Dusson, Christoph Ortner, and Gábor Csányi. Machine Learning: Science and Technology, 2021

Linear Atomic Cluster Expansion Force Fields for Organic Molecules: beyond RMSE., Kovacs, D. P., van der Oord, C., Kucera, J., Allen, A., Cole, D., Ortner, C., & Csanyi, G., JCTC, 2021

## Training Data

Machine learning potentials with Iterative Boltzmann Inversion: training to experiment, Sakib Matin, Alice Allen, Justin S Smith, Nicholas Lubbers, Ryan B Jadrich, Richard A Messerly, Benjamin T Nebgen, Ying Wai Li, Sergei Tretiak, Kipton Barros, JCTC, 2024

## Fitting Methods

Learning together: Towards foundation models for machine learning interatomic potentials with meta-learning, AEA Allen, N Lubbers, S Matin, J Smith, R Messerly, S Tretiak, K Barros, npj Computational Materials, 2024

## Interpretability

Alice E. A. Allen and Alexandre Tkatchenko, Machine Learning of Material Properties: Predictive and Interpretable Multilinear Models, Science Advances, 2022

## Architecture

Atomic permutationally invariant polynomials for fitting molecular force fields. Alice E. A. Allen, Geneviève Dusson, Christoph Ortner, and Gábor Csányi. Machine Learning: Science and Technology, 2021

Linear Atomic Cluster Expansion Force Fields for Organic Molecules: beyond RMSE., Kovacs, D. P., van der Oord, C., Kucera, J., Allen, A., Cole, D., Ortner, C., & Csanyi, G., JCTC, 2021



## Training Data

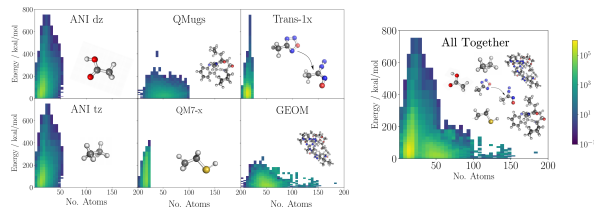Machine learning potentials with Iterative Boltzmann Inversion: training to experiment, Sakib Matin, Alice Allen, Justin S Smith, Nicholas Lubbers, Ryan B Jadrich, Richard A Messerly, Benjamin T Nebgen, Ying Wai Li, Sergei Tretiak, Kipton Barros, JCTC, 2024
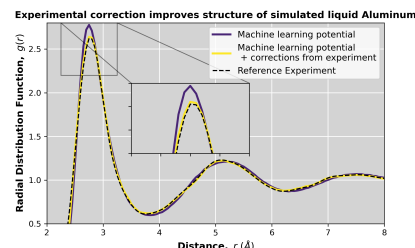


## Fitting Methods

Learning together: Towards foundation models for machine learning interatomic potentials with meta-learning, AEA Allen, N Lubbers, S Matin, J Smith, R Messerly, S Tretiak, K Barros, npj Computational Materials, 2024
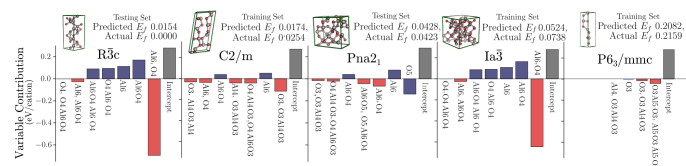


## Interpretability

Alice E. A. Allen and Alexandre Tkatchenko, Machine Learning of Material Properties: Predictive and Interpretable Multilinear Models, Science Advances, 2022

# Architecture

- Interatomic potentials are a regression problem with a regression model and representation chosen

- The representation chosen must satisfy the following properties:

    - Rigid translations

# Architecture

- Interatomic potentials are a regression problem with a regression model and representation chosen

- The representation chosen must satisfy the following properties:

    - Rigid translations

# Architecture

- Interatomic potentials are a regression problem with a regression model and representation chosen

- The representation chosen must satisfy the following properties:

    - Rigid translations

    - Rotations and reflections

# Architecture

- Interatomic potentials are a regression problem with a regression model and representation chosen

- The representation chosen must satisfy the following properties:

    - Rigid translations

    - Rotations and reflections

# Architecture

- Interatomic potentials are a regression problem with a regression model and representation chosen

- The representation chosen must satisfy the following properties:

    - Rigid translations

    - Rotations and reflections

    - Permutation of like atoms

# Architecture

- Interatomic potentials are a regression problem with a regression model and representation chosen

- The representation chosen must satisfy the following properties:

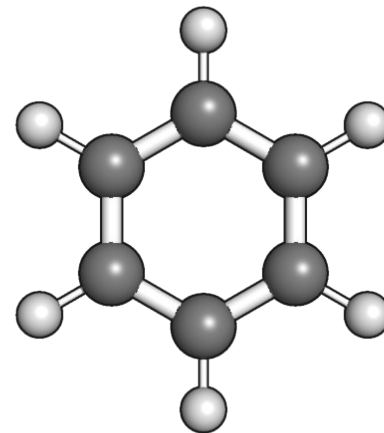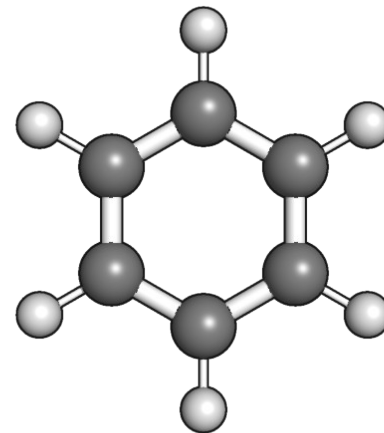    - Rigid translations

    - Rotations and reflections

    - Permutation of like atoms

- Three identical atoms example for PIPs:



Rotational Invariance    Permutational Invariance

$$E(R_1, R_2, R_3) = E(r_{12}, r_{13}, r_{23}) = E(I_1, I_2, I_3)$$

$$I_1 = r_{12} + r_{13} + r_{23}$$

$$I_2 = r_{12}r_{13} + r_{12}r_{23} + r_{13}r_{23}$$

$$I_3 = r_{12}r_{13}r_{23}$$

$$E(R_1, R_2, R_3) = P(I_1, I_2, I_3) = \sum_k c_k I_1^{k_1} I_2^{k_2} I_3^{k_3}$$

- These ideas can then be extended to more atoms and to the non-identical atoms case.

- PIPs are a polynomial construction with the underlying invariances present in the basis used.

- aPIPs retains the polynomial construction

$$E = \sum_{i=1}^{N} E_i$$

- But moves to an atom centred model

- Uses a body-ordered approached:

$$P(I_1, I_2, I_3) = \sum_k c_k I_1^{k_1} I_2^{k_2} I_3^{k_3}$$

$$E(\mathbf{R}) = \sum_i E_1 + \sum_{i<j} E_2(r_{ij}) + \sum_{i<j<k} E_3(r_{ij}, r_{ik}, r_{jk}) + \sum_{i<j<k<l} E_4(r_{ij}, r_{ik}, r_{il}, r_{jk}, r_{jl}, r_{kl}) + \ldots$$

Cas van der Oord, Geneviève Dusson, Gábor Csányi and Christoph Ortner, Regularised atomic body-ordered permutation-invariant polynomials for the construction of interatomic potentials, (2020). Machine Learning: Science and Technology

Allen, Alice E A, Dusson, Geneviève, Ortner, Christoph, Csányi, Gábor, Atomic permutationally invariant polynomials for fitting molecular force fields, (2021). Machine Learning: Science and Technology

The atomic basis is define as [3, 4]:

$$\alpha_{z_i, znlm} = \sum_{\substack{j \\ \text{where } z_j = z}} \phi_{nlm}^{z_i z}(\mathbf{r}_{ji}) \qquad (1)$$

$$\phi_{nlm}^{z_i z_j}(\mathbf{r}) = R_{nl}^{z_i z_j}(r) Y_l^m(\hat{\mathbf{r}}). \qquad (2)$$

where $R_{nl}$ is a radial basis function and $Y_l^m$ is a spherical harmonic function. A permutation-invariant basis function is constructed by forming products of the atomic basis:

$$A_{z_i \mathbf{v}} = \prod_{t=1}^{\nu} \alpha_{z_i v_t}, \quad \mathbf{v} = (v_1, \ldots, v_\nu). \qquad (3)$$

where $v = znlm$, defining the atomic number of the atom and the properties of the radial function and spherical harmonics. Rotational invariance is then achieved by the following stage:

$$B_{z_i \mathbf{v}} := \int_{\hat{R} \in O(3)} \prod_{t=1}^{\nu} A_{z_i v_t}(\{\hat{R}\mathbf{r}_{ij}\}) \, d\hat{R} \qquad (4)$$

$$= \sum_{\mathbf{v}'} C_{\mathbf{v}\mathbf{v}'} A_{z_i \mathbf{v}'}, \qquad (5)$$

where the coefficients $C_{\mathbf{v}\mathbf{v}'}$ are the Clebsch-Gordan coupling coefficients. The energy at site $i$ is then given by:

$$E_i = \sum c_{z_i \mathbf{v}} B_{z_i \mathbf{v}} = \mathbf{c} \cdot \mathbf{B}. \qquad (6)$$



Kovacs, D. P., van der Oord, C., Kucera, J., Allen, A., Cole, D., Ortner, C., & Csanyi, G., Linear Atomic Cluster Expansion Force Fields for Organic Molecules: beyond RMSE. JCTC, 2022

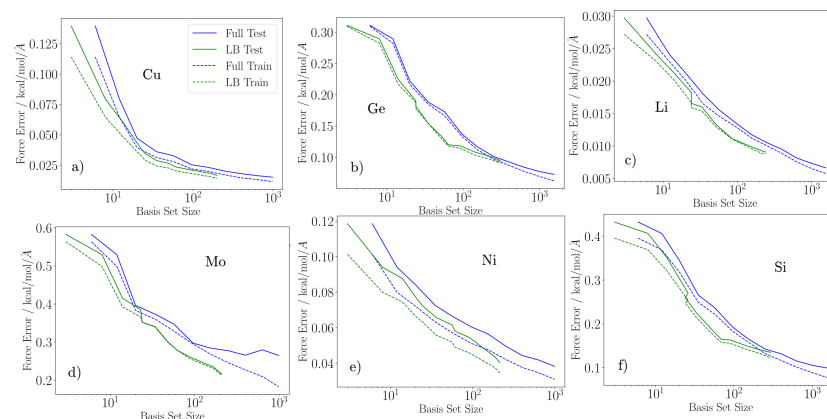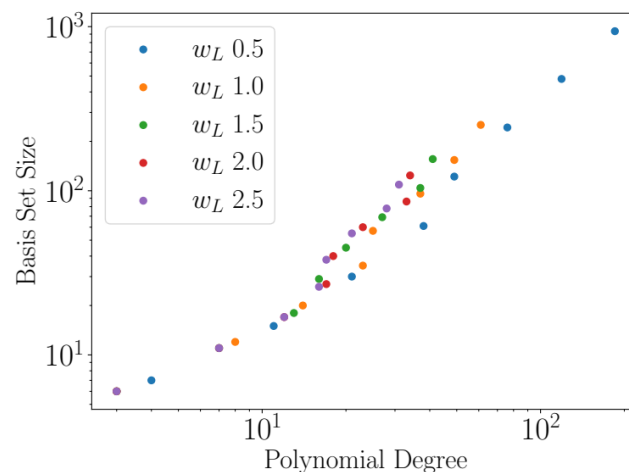- Construct matrix of $\frac{\partial B_i}{\partial \alpha_j}$ for one atom (either from training dataset or assigned random neighbours). The matrix $\frac{\partial B_i}{\partial \alpha_j}$ is calculated numerically.

- Initialize a matrix with columns equivalent to the number of atomic basis functions ($\alpha_{z_i, znlm}$) and no rows.

- For each $B_i$ until the utmost number of $B_i$ terms is reached: append the row vector of the partial derivatives of the subsequent invariant into the matrix if it augments the rank of the matrix. From this, constructing an index of $B_i$ components that increase rank.

- Repeat this $N$ times, constructing an index of all $B_i$ components that increase rank for any datapoint present.
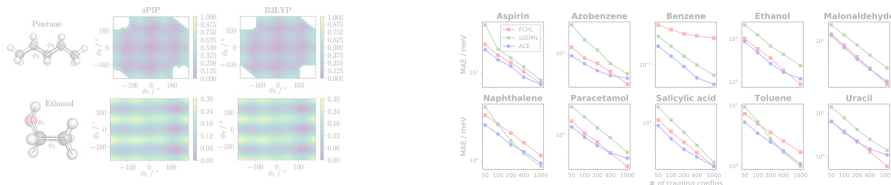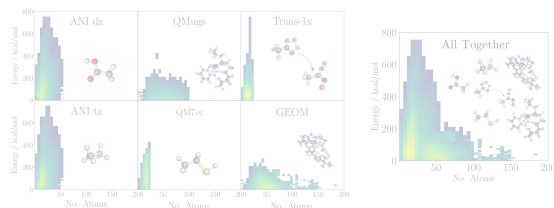
# Improving ML Models

## Architecture

Atomic permutationally invariant polynomials for fitting molecular force fields. Alice E. A. Allen, Geneviève Dusson, Christoph Ortner, and Gábor Csányi. Machine Learning: Science and Technology, 2021

Linear Atomic Cluster Expansion Force Fields for Organic Molecules: beyond RMSE., Kovacs, D. P., van der Oord, C., Kucera, J., Allen, A., Cole, D., Ortner, C., & Csanyi, G., JCTC, 2021



## **Training Data**

Machine learning potentials with Iterative Boltzmann Inversion: training to experiment, Sakib Matin, Alice Allen, Justin S Smith, Nicholas Lubbers, Ryan B Jadrich, Richard A Messerly, Benjamin T Nebgen, Ying Wai Li, Sergei Tretiak, Kipton Barros, JCTC, 2024
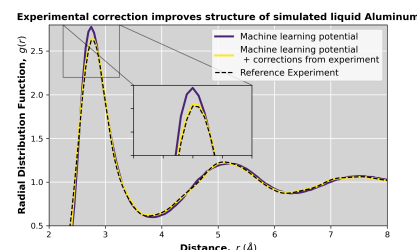


## Fitting Methods

Learning together: Towards foundation models for machine learning interatomic potentials with meta-learning, AEA Allen, N Lubbers, S Matin, J Smith, R Messerly, S Tretiak, K Barros, npj Computational Materials, 2024
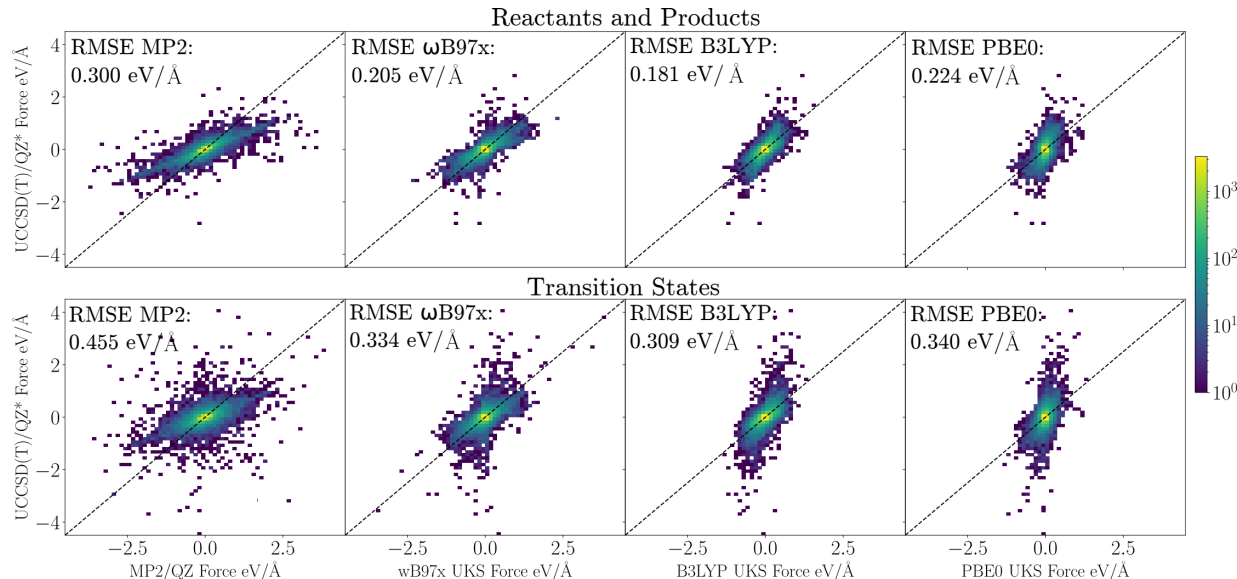


## Interpretability

Alice E. A. Allen and Alexandre Tkatchenko, Machine Learning of Material Properties: Predictive and Interpretable Multilinear Models, Science Advances, 2022
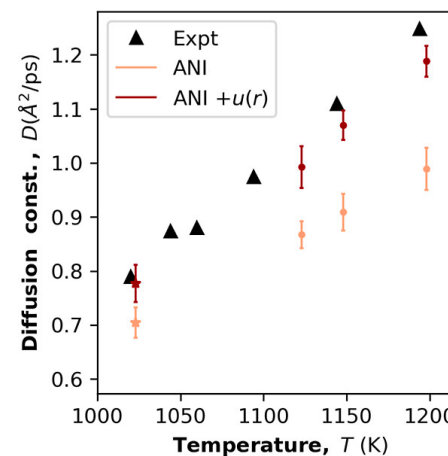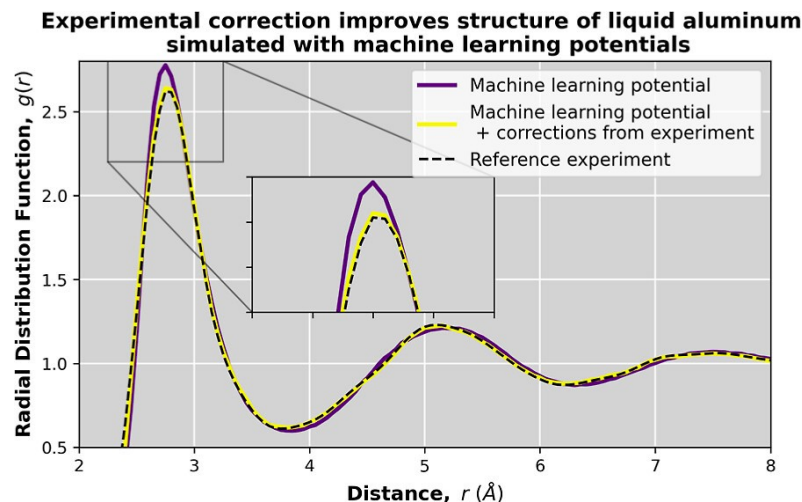
# Fitting Data – QM Data

- Improving QM datasets can include improving the level of theory used

- We have been building a unrestricted CCSD(T) dataset for gas phase reactions

- Energies and forces are included

# Fitting Data - Experimental

- QM datasets have fundamental limitations

- How can experimental data be incorporated into ML models?

- Investigated using a pair potential correction to an existing MLIP using  radial distribution function data

- Improvements then seen in diffusion coefficients when correction present



Machine learning potentials with Iterative Boltzmann Inversion: training to experiment, Sakib Matin, Alice Allen, Justin S Smith, Nicholas Lubbers, Ryan B Jadrich, Richard A Messerly, Benjamin T Nebgen, Ying Wai Li, Sergei Tretiak, Kipton Barros,  JCTC, 2024

# Improving ML Models

## Architecture

Atomic permutationally invariant polynomials for fitting molecular force fields. Alice E. A. Allen, Geneviève Dusson, Christoph Ortner, and Gábor Csányi. Machine Learning: Science and Technology, 2021

Linear Atomic Cluster Expansion Force Fields for Organic Molecules: beyond RMSE., Kovacs, D. P., van der Oord, C., Kucera, J., Allen, A., Cole, D., Ortner, C., & Csanyi, G., JCTC, 2021
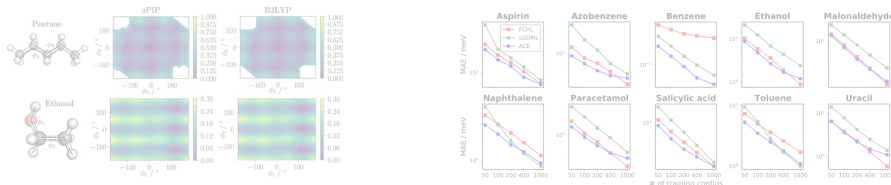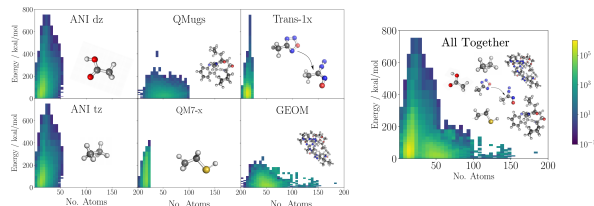


## Training Data

Machine learning potentials with Iterative Boltzmann Inversion: training to experiment, Sakib Matin, Alice Allen, Justin S Smith, Nicholas Lubbers, Ryan B Jadrich, Richard A Messerly, Benjamin T Nebgen, Ying Wai Li, Sergei Tretiak, Kipton Barros, JCTC, 2024
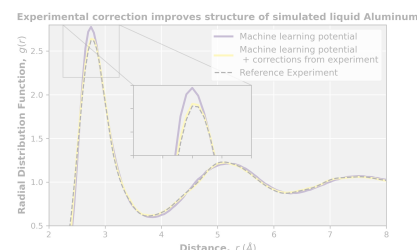


## Fitting Methods

Learning together: Towards foundation models for machine learning interatomic potentials with meta-learning, AEA Allen, N Lubbers, S Matin, J Smith, R Messerly, S Tretiak, K Barros, npj Computational Materials, 2024



## Interpretability

Alice E. A. Allen and Alexandre Tkatchenko, Machine Learning of Material Properties: Predictive and Interpretable Multilinear Models, Science Advances, 2022
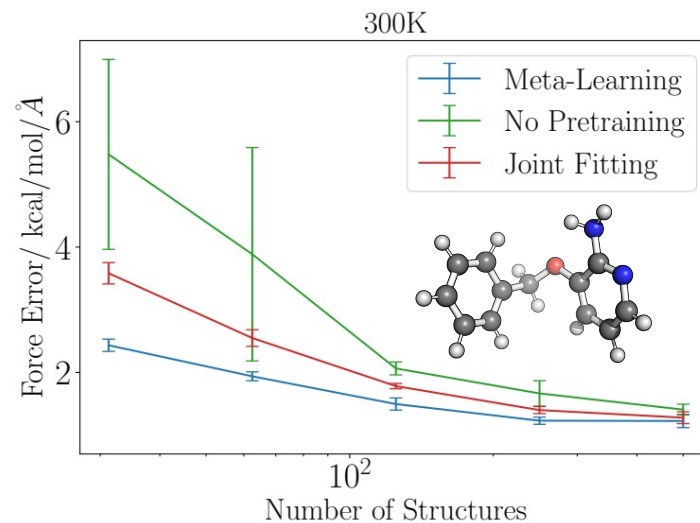
# Fitting Methods – Meta Learning

- There is an abundance of datasets containing quantum mechanical calculations for molecular and material systems.

- However, using the information from different datasets together remains a challenge due to the varying levels of theory employed.

- We have shown that meta-learning can be used to pre-train models to multiple datasets.

**Algorithm 1** Reptile

1: Initialize $\Phi$, the initial parameter vector
2: **for** iteration 1,2,3,... **do**
3:      Randomly sample a task T
4:      Perform k > 1 steps in task T, starting with parameters $\Phi$, resulting in parameters W
5:      Update: $\Phi \leftarrow \Phi + \varepsilon(W - \Phi)$
6: Return $\Phi$

300K

Meta-Learning
No Pretraining
Joint Fitting

Force Error / kcal/mol/Å

Number of Structures
$10^2$

Learning together: Towards foundation models for machine learning interatomic potentials with meta-learning, AEA Allen, N Lubbers, S Matin, J Smith, R Messerly, S Tretiak, K Barros, npj Computational Materials, 2024
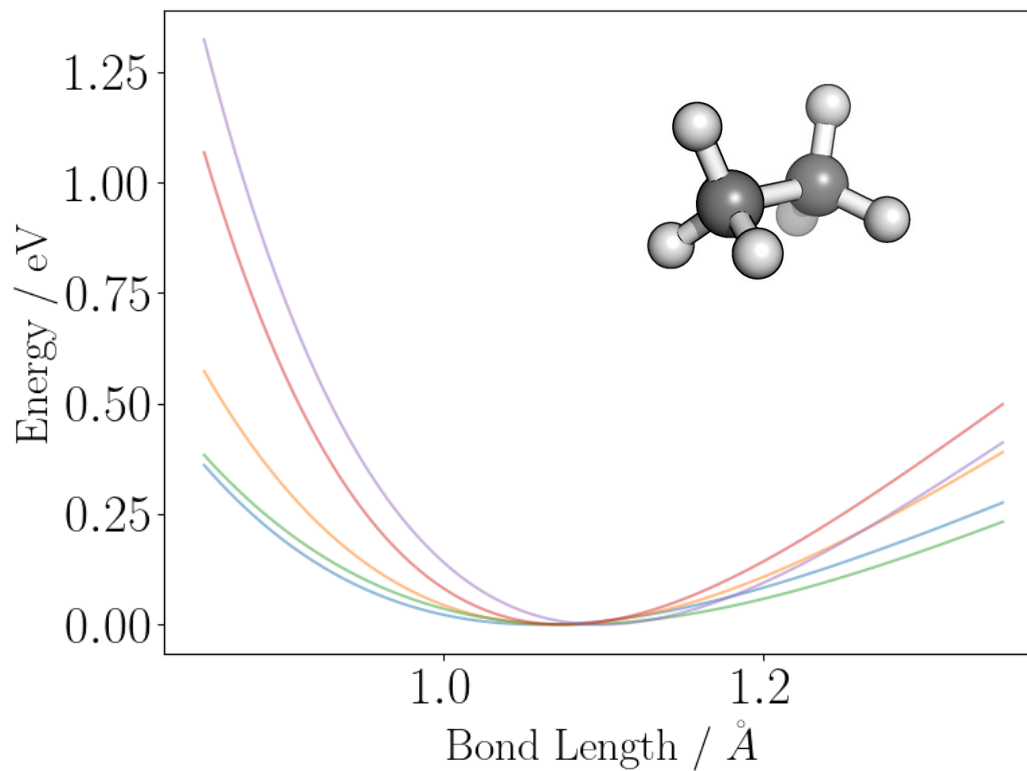
# Problem

How do we fit multiple datasets together that use different QM approximations?
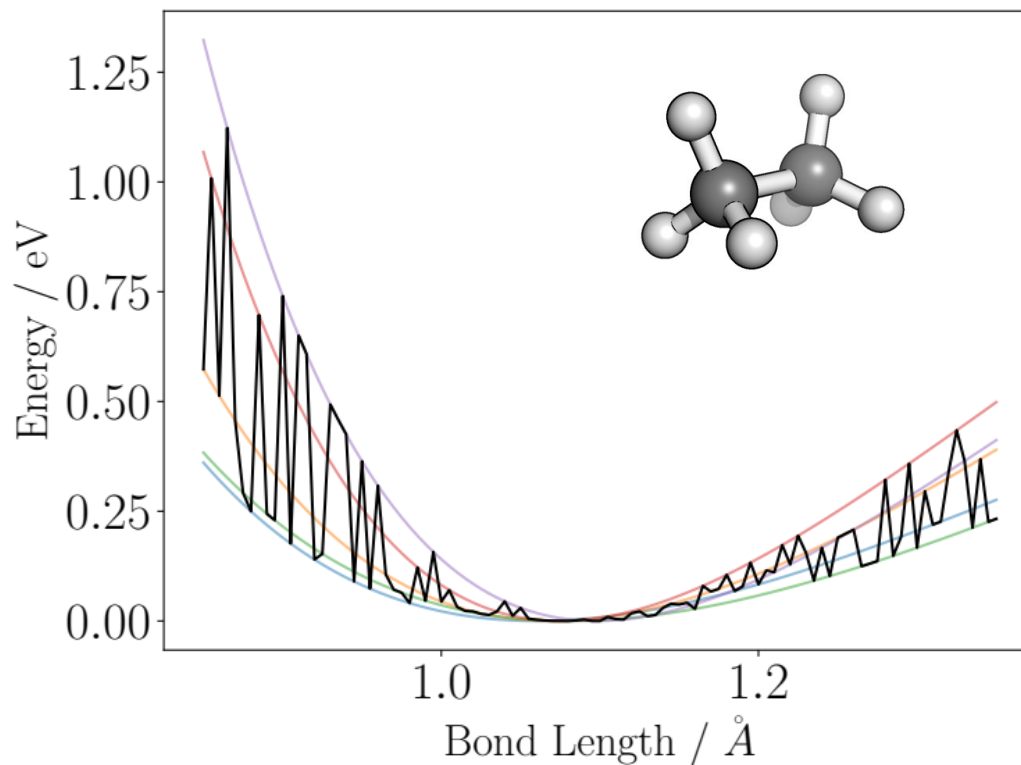
# Multiple Existing Datasets Available

| Dataset | Unique Compounds | Total Conformers | Heavy Atoms Max | Conformer Generation | Method | Dispersion | Transition Paths |
|---|---|---|---|---|---|---|---|
| QM9 | 133,885 | 133,885 | 9 | None | 76 DFT Functionals | Yes | No |
| AN1-1x | ~64,000 | 4,956,005 and 4,617,229 | 8 | Normal Mode Sampling, MD sampling, Torsional Sampling, Active Learning with QBC | $\omega$B97x/ 6-31G* and $\omega$B97x/ def2-TZVPP | No | No |
| QMugs | 665,911 | 1,992,984 | 100 | Meta Dynamics with xTB | $\omega$B97X-D/ def2-SVP | Yes | No |
| GEOM | 437,724 | 32,657,609 | 91 | Meta Dynamics with CREST | r2scan-3c/ mTZVPP | Yes | No |
| QM7-x | 41,537* | 4,195,237 | 7 | Normal-mode Sampling with DFTB | PBE0+MBD | Yes | No |
| Transition1x | 10,073 Reactions | 9,644,740 | 7 | Nudged Elastic Band | $\omega$B97x/ 6-31G(d) | No | Yes |
| ANI-1ccx | ~64,000? | 489,571 | 8 | Active Learning from ANI-1x | CCSD(T)*/ CBS | Yes | No |

Over 90 million calculations included, with more than 200 different levels of theory.

# Transfer Learning

## Approaching coupled cluster accuracy with a general-purpose neural network potential through transfer learning

Justin S. Smith[1,2,3,7], Benjamin T. Nebgen [2,4,7], Roman Zubatyuk[2,5,7], Nicholas Lubbers[2,3], Christian Devereux[1], Kipton Barros[2], Sergei Tretiak [2,4], Olexandr Isayev [6] & Adrian E. Roitberg[1]

Computational modeling of chemical and biological systems at atomic resolution is a crucial tool in the chemist's toolset. The use of computer simulations requires a balance between cost and accuracy: quantum-mechanical methods provide high accuracy but are computationally expensive and scale poorly to large systems, while classical force fields are cheap and scalable, but lack transferability to new systems. Machine learning can be used to achieve the best of both approaches. Here we train a general-purpose neural network potential (ANI-1ccx) that approaches CCSD(T)/CBS accuracy on benchmarks for reaction thermochemistry, isomerization, and drug-like molecular torsions. This is achieved by training a network to DFT data then using transfer learning techniques to retrain on a dataset of gold standard QM calculations (CCSD(T)/CBS) that optimally spans chemical space. The resulting potential is broadly applicable to materials science, biology, and chemistry, and billions of times faster than CCSD(T)/CBS calculations.

Trained ANI potential to DFT and then retrained to CCSD(T) dataset with frozen parameters.

Freezing parameters and other transfer learning techniques have limitations and can't readily train to multiple datasets.

Problem: How could we extend this to multiple datasets/level of theory?

# Meta-learning

**The problem:**

How can we combine multiple QM datasets?

**The solution**

Meta-learning trains models in such a way that the solution can easily be refit to a new task

Not about finding a model that works well for one task

But finding a model that can be easily retrained to new tasks with limited data

On First-Order Meta-Learning Algorithms

Alex Nichol and Joshua Achiam and John Schulman
OpenAI
{alex, jachiam, joschu}@openai.com

**Abstract**

This paper considers meta-learning problems, where there is a distribution of tasks, and we would like to obtain an agent that performs well (i.e., learns quickly) when presented with a previously unseen task sampled from this distribution. We analyze a family of algorithms for learning a parameter initialization that can be fine-tuned quickly on a new task, using only first-order derivatives for the meta-learning updates. This family includes and generalizes first-order MAML, an approximation to MAML obtained by ignoring second-order derivatives. It also includes Reptile, a new algorithm that we introduce here, which works by repeatedly sampling a task, training on it, and moving the initialization towards the trained weights on that task. We expand on the results from Finn et al. showing that first-order meta-learning algorithms perform well on some well-established benchmarks for few-shot classification, and we provide theoretical analysis aimed at understanding why these algorithms work.

Learning Together: Towards foundation models for machine learning interatomic potentials with meta-learning, Alice E. A. Allen, Nicholas Lubbers, Sakib Matin, Justin Smith, Richard Messerly, Sergei Tretiak, Kipton Barros, npj Computational Materials, 2024

# Meta-learning

What is meta-learning?

Trying to train models in such a way that the solution can easily be refit to a new task

Not about finding a model that works well for one task

But finding a model that can be easily retrained to new tasks with limited data
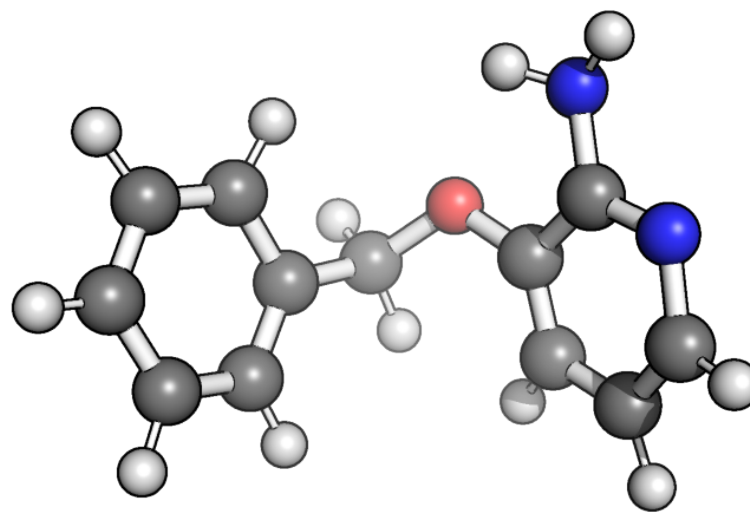
Training

Task 1

Task 2

Test

New Task

Learning Together: Towards foundation models for machine learning interatomic potentials with meta-learning, Alice E. A. Allen, Nicholas Lubbers, Sakib Matin, Justin Smith, Richard Messerly, Sergei Tretiak, Kipton Barros, npj Computational Materials, 2024
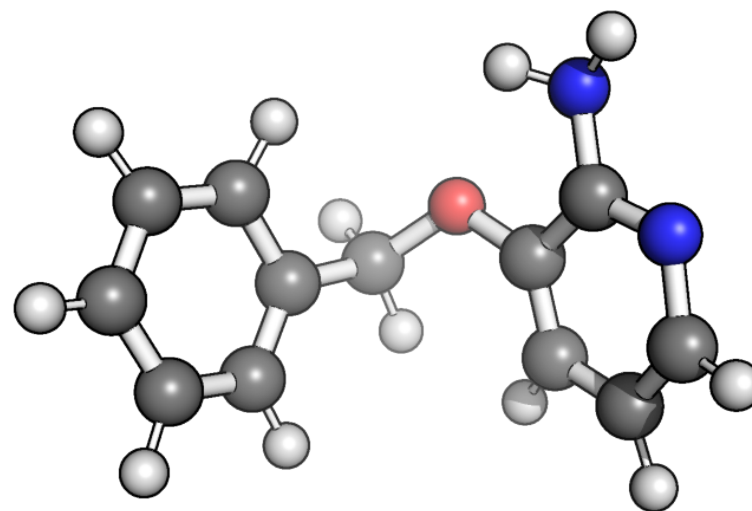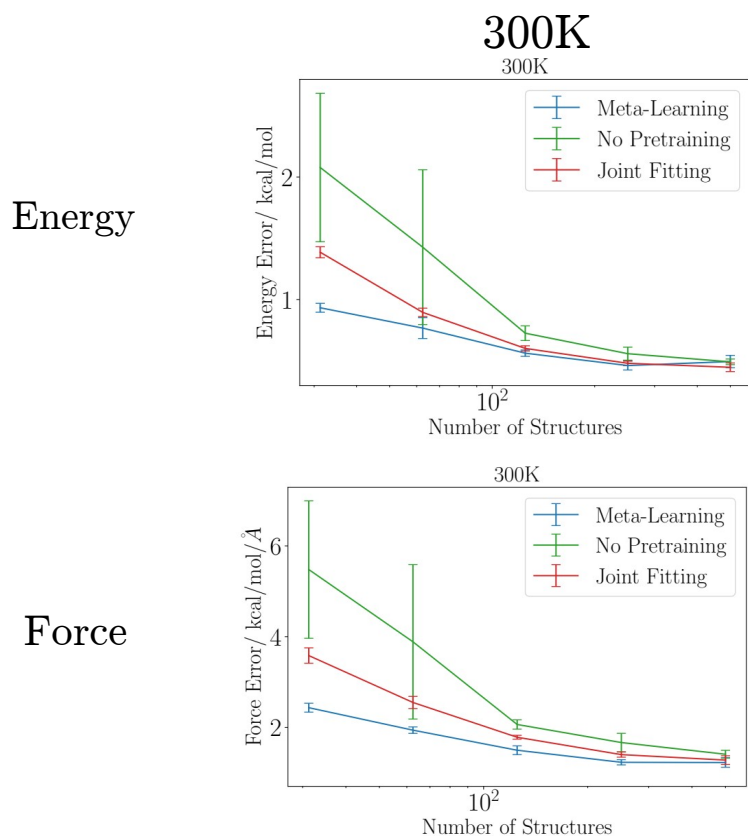
# Meta-learning

Previously implemented the Reptile meta-learning algorithm.

This divides a problem into multiple tasks and then performs multiple optimization steps on individual tasks to find a model that generalizes well to new tasks.

**Algorithm 1** Reptile

1: Initialize $\Phi$, the initial parameter vector
2: **for** iteration 1,2,3,... **do**
3:     Randomly sample a task T
4:     Perform k > 1 steps in task T, starting with parameters $\Phi$, resulting in parameters W
5:     Update: $\Phi \leftarrow \Phi + \varepsilon(W - \Phi)$
6: Return $\Phi$

Reptile Algorithm

$$\theta_i = \theta_{i-1} + \epsilon(W_{i-1} - \theta_{i-1})$$

$\theta_1$

Trained to Task 1

k steps

$W_1$

$\theta_2$

k steps

Trained to Task 2

$W_2$

$\theta_3$

Trained to Task 3

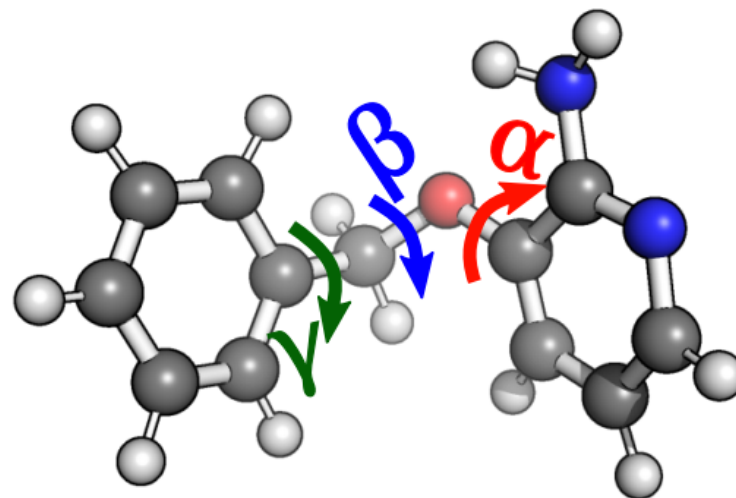$\theta_4$

k steps

k steps

$W_3$

Learning Together: Towards foundation models for machine learning interatomic potentials with meta-learning, Alice E. A. Allen, Nicholas Lubbers, Sakib Matin, Justin Smith, Richard Messerly, Sergei Tretiak, Kipton Barros, npj Computational Materials, 2024
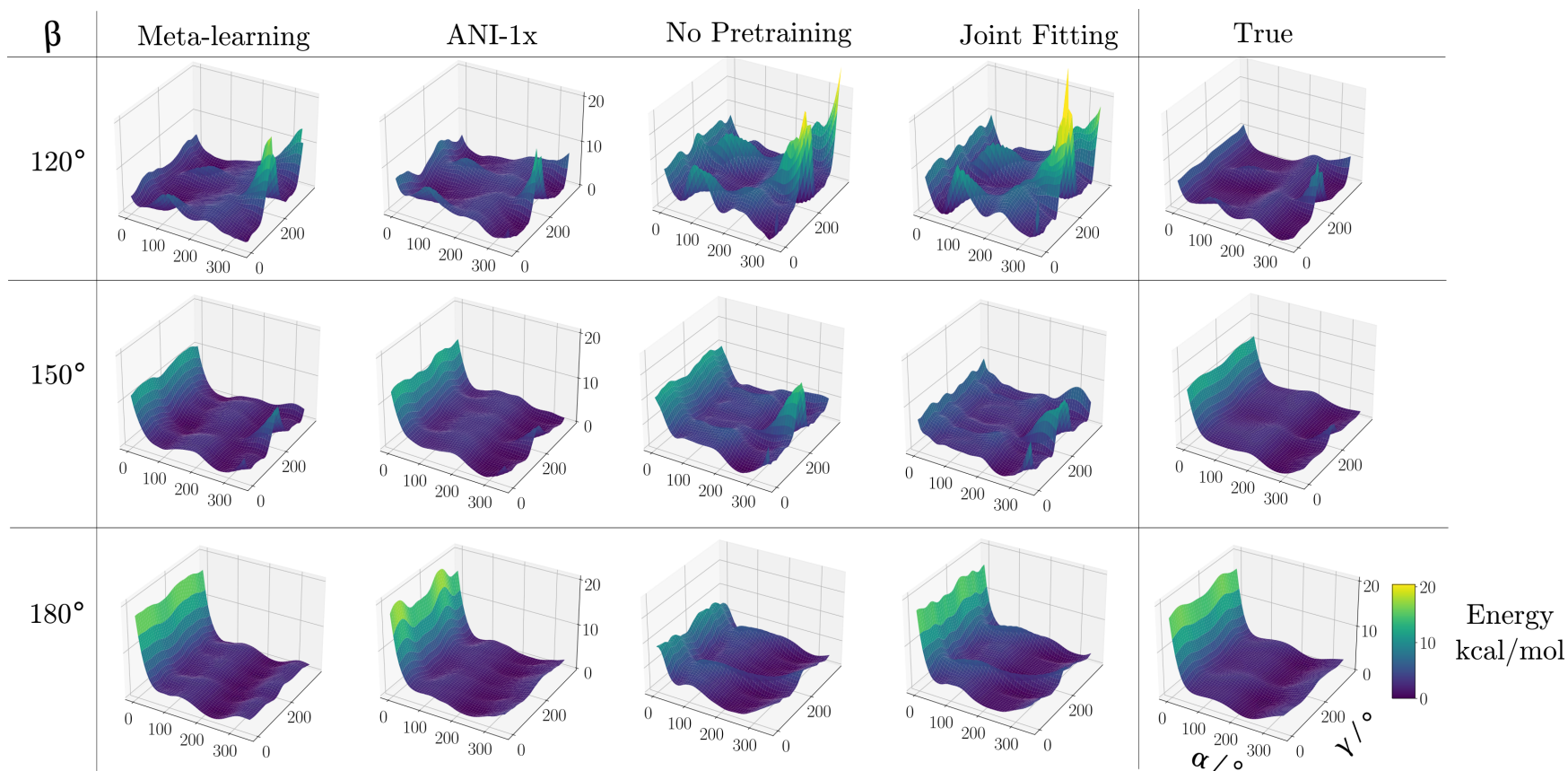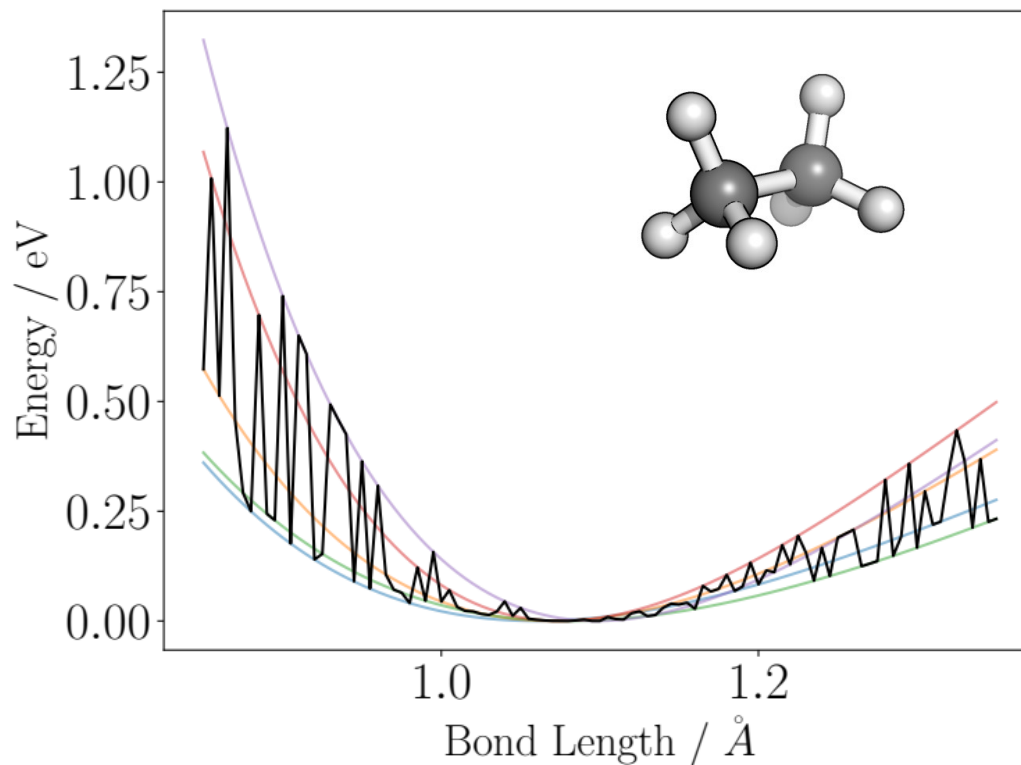
# Reptile

We have implemented the Reptile meta-learning algorithm.

This divides a problem into multiple tasks and then performs multiple optimization steps on individual tasks to find a model that generalizes well to new tasks.

NOT:

    Trying to find parameters that are good for one specific task

    Trying to find parameters that work for all tasks at once

# Aspirin

Aspirin was pre-trained to three datasets from 300K, 600K and 900K MD simulations

Three levels of theory were used to pre-train the potential

The potential was then retrained to 400 structures at MP2.



$$\theta_i = \theta_{i-1} + \epsilon(W_{i-1} - \theta_{i-1})$$

Reptile Algorithm

# QM9

Meta-learning potential was fit to 150 different levels of theory from the QM9 dataset

Refitting was then performed on 15 different levels of theory that had not previous been encountered

The results of this are shown to the right.

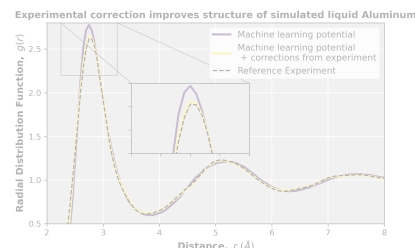S. Nandi, T. Vegge, and A. Bhowmik, ChemRxiv (2022), 10.26434/chemrxiv-2022-fs70n.



Learning Together: Towards foundation models for machine learning interatomic potentials with meta-learning, Alice E. A. Allen, Nicholas Lubbers, Sakib Matin, Justin Smith, Richard Messerly, Sergei Tretiak, Kipton Barros, npj Computational Materials, 2024

# Combining Datasets

Previously implemented the Reptile meta-learning algorithm.



Learning Together: Towards foundation models for machine learning interatomic potentials with meta-learning, Alice E. A. Allen, Nicholas Lubbers, Sakib Matin, Justin Smith, Richard Messerly, Sergei Tretiak, Kipton Barros, npj Computational Materials, 2024

# 3BPA



Start by investigating the potential for a single molecule

3BPA – shown to the right

Kovacs, D. P., van der Oord, C., Kucera, J., Allen, A., Cole, D., Ortner, C., & Csanyi, G., Linear Atomic Cluster Expansion Force Fields for Organic Molecules: beyond RMSE., (2021). JCTC.

# Meta-learning



Energy

Force

300K

Learning Together: Towards foundation models for machine learning interatomic potentials with meta-learning, Alice E. A. Allen, Nicholas Lubbers, Sakib Matin, Justin Smith, Richard Messerly, Sergei Tretiak, Kipton Barros, npj Computational Materials, 2024

# 3BPA

Next lets look at the 2D torsional energy scans

# Meta-learning

Learning Together: Towards foundation models for machine learning interatomic potentials with meta-learning, Alice E. A. Allen, Nicholas Lubbers, Sakib Matin, Justin Smith, Richard Messerly, Sergei Tretiak, Kipton Barros, npj Computational Materials, 2024

# Improving ML Models

## Architecture

Atomic permutationally invariant polynomials for fitting molecular force fields. Alice E. A. Allen, Geneviève Dusson, Christoph Ortner, and Gábor Csányi. Machine Learning: Science and Technology, 2021

Linear Atomic Cluster Expansion Force Fields for Organic Molecules: beyond RMSE., Kovacs, D. P., van der Oord, C., Kucera, J., Allen, A., Cole, D., Ortner, C., & Csanyi, G., JCTC, 2021



## Training Data

Machine learning potentials with Iterative Boltzmann Inversion: training to experiment, Sakib Matin, Alice Allen, Justin S Smith, Nicholas Lubbers, Ryan B Jadrich, Richard A Messerly, Benjamin T Nebgen, Ying Wai Li, Sergei Tretiak, Kipton Barros, JCTC, 2024



## Fitting Methods

Learning together: Towards foundation models for machine learning interatomic potentials with meta-learning, AEA Allen, N Lubbers, S Matin, J Smith, R Messerly, S Tretiak, K Barros, npj Computational Materials, 2024



## Interpretability

Alice E. A. Allen and Alexandre Tkatchenko, Machine Learning of Material Properties: Predictive and Interpretable Multilinear Models, Science Advances, 2022

# Interpretability

- Machine learning models can provide fast and accurate predictions of material properties but often lack transparency.

- Interpretability techniques can be used with black box solutions, or alternatively, models can be created that are directly interpretable.

- We revisited several works and demonstrate that simple linear combinations of nonlinear basis functions can be created, which have comparable accuracy to the kernel and neural network approaches originally used.



Machine learning of material properties: Predictive and interpretable multilinear models, Allen and Tkatchenko, Sci. Adv. 8, eabm7185 (2022)

- A predictive model for the formation energy of 10,000 elpasolite structures ($ABC_2D_6$ in the *Fm3m* space group) was produced using kernel ridge regression.

## Machine Learning Energies of 2 Million Elpasolite ($ABC_2D_6$) Crystals

Felix A. Faber,[1] Alexander Lindmaa,[2] O. Anatole von Lilienfeld,[1,3,*] and Rickard Armiento[2,†]

[1]*Institute of Physical Chemistry and National Center for Computational Design and Discovery of Novel Materials, Department of Chemistry, University of Basel, 4056 Basel, Switzerland*
[2]*Department of Physics, Chemistry and Biology, Linköping University, SE-581 83 Linköping, Sweden*
[3]*General Chemistry, Free University of Brussels, Pleinlaan 2, 1050 Brussels, Belgium*
(Received 24 August 2015; published 20 September 2016)

Elpasolite is the predominant quaternary crystal structure ($AlNaK_2F_6$ prototype) reported in the Inorganic Crystal Structure Database. We develop a machine learning model to calculate density functional theory quality formation energies of all $\sim 2 \times 10^6$ pristine $ABC_2D_6$ elpasolite crystals that can be made up from main-group elements (up to bismuth). Our model's accuracy can be improved systematically, reaching a mean absolute error of 0.1 eV/atom for a training set consisting of $10 \times 10^3$ crystals. Important bonding trends are revealed: fluoride is best suited to fit the coordination of the *D* site, which lowers the formation energy whereas the opposite is found for carbon. The bonding contribution of the elements *A* and *B* is very small on average. Low formation energies result from *A* and *B* being late elements from group II, *C* being a late (group I) element, and *D* being fluoride. Out of $2 \times 10^6$ crystals, 90 unique structures are predicted to be on the convex hull—among which is $NFAl_2Ca_6$, with a peculiar stoichiometry and a negative atomic oxidation state for Al.
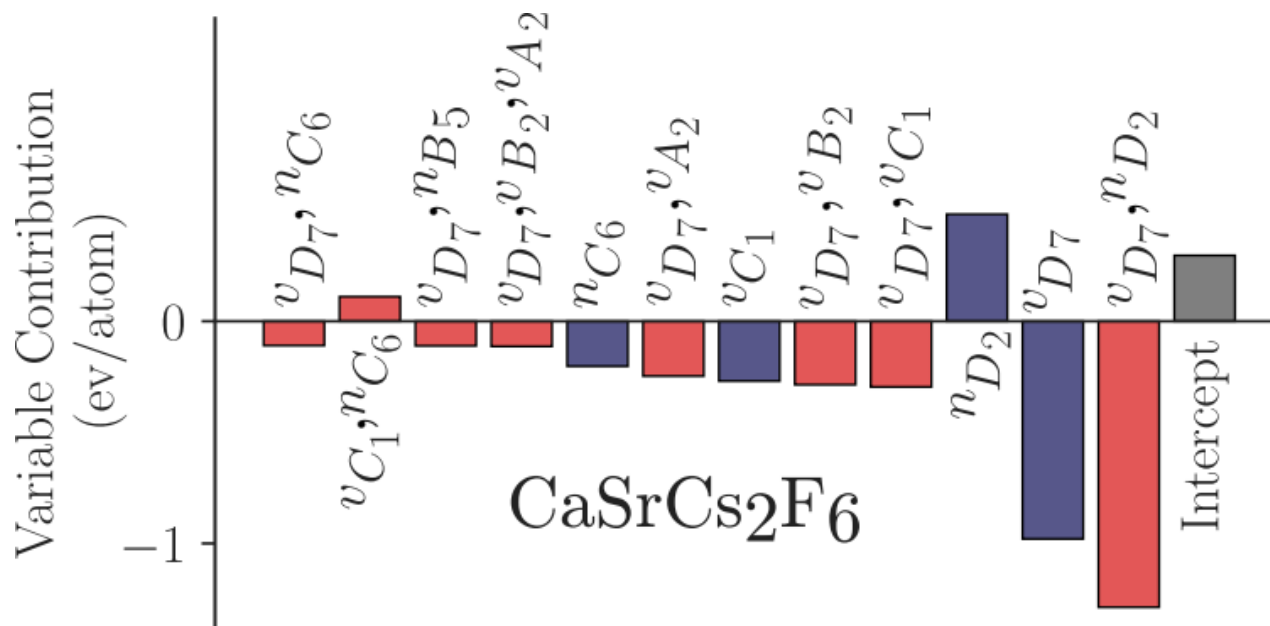
# Interpretability

- The features used to describe the structures were the principle quantum number (n) and number of valence electrons (v) at each site A,B,C or D.

- Reach accuracy of 0.11 eV/atom, compared to 0.10 eV/atom for the KRR.

- The accuracy of the DFT data has been stated as between 0.10 – 0.19 eV/atom

- So we have comparable accuracy to KRR and DFT

$$
\begin{aligned}
E(n,v) = & \sum_i \alpha_i n_i + \sum_i \beta_i v_i + \sum_{j<i} \alpha_{ij} n_i n_j + \sum_{j<i} \beta_{ij} v_i v_j \\
& + \sum_{i,j} \gamma_{ij} n_i v_j + \sum_{k<j<i} \alpha_{ijk} n_i n_j n_k + \sum_{k<j<i} \beta_{ijk} v_i v_j v_k \\
& + \sum_{k<i,j} \gamma_{ijk} n_i v_j n_k + \sum_{k<j,i} \lambda_{ijk} n_i v_j v_k \quad (3)
\end{aligned}
$$
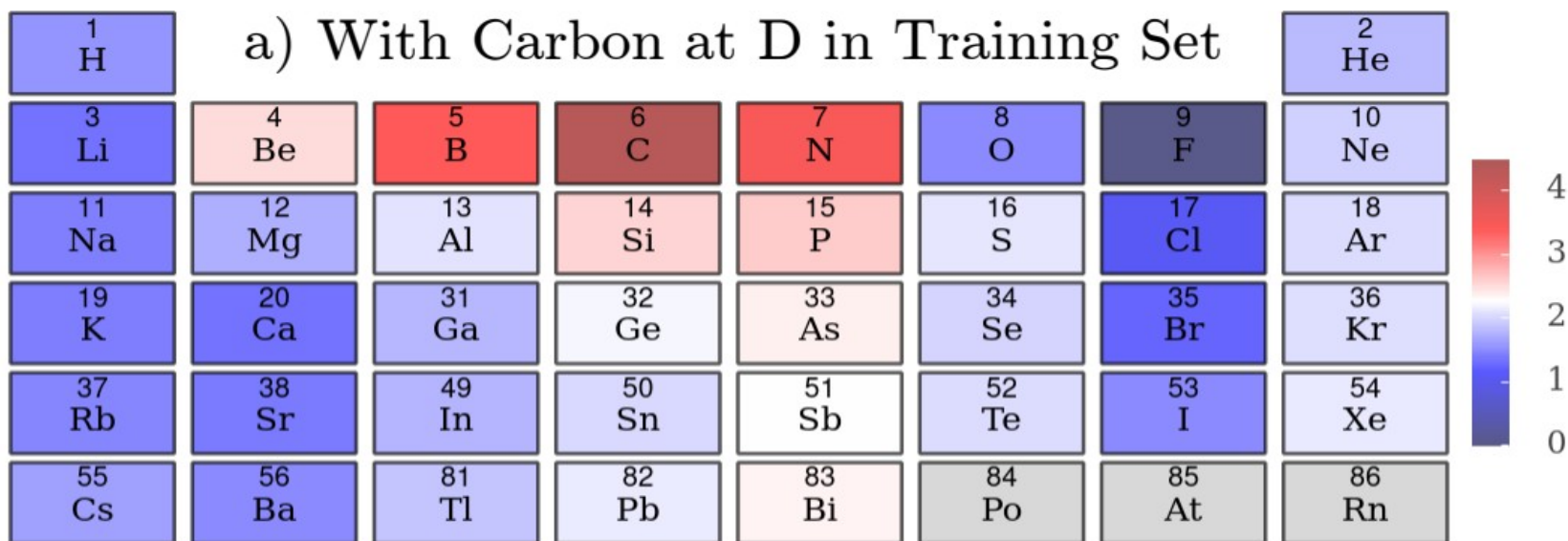
- Search 2 million possible structures

- Individual predictions can be broken down into variable and interaction contributions:

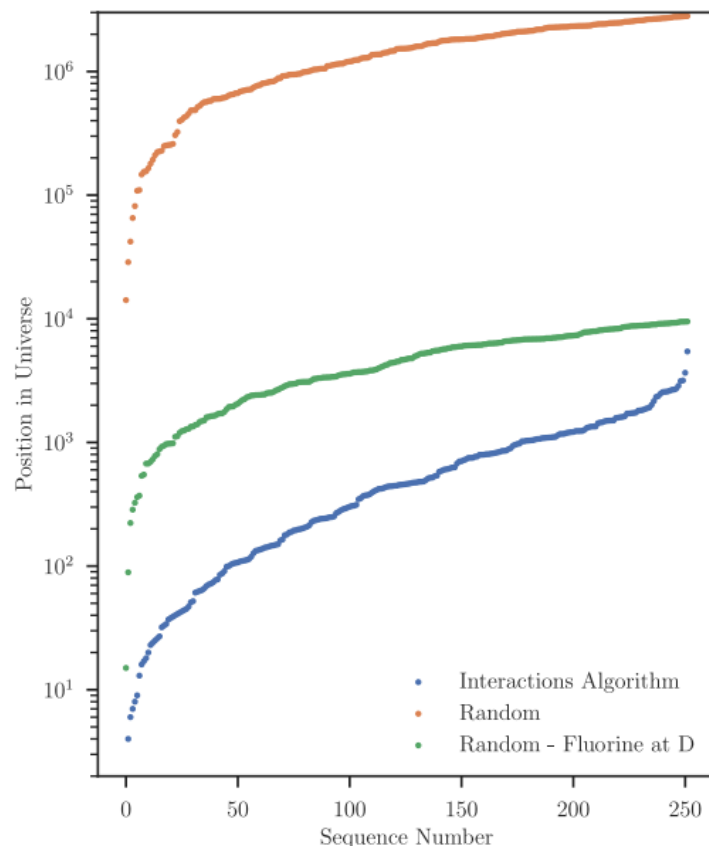- Contributions to the formation energy at position D can be found directly from the coefficients of the model:



a) With Carbon at D in Training Set

# Interpretability

- With the linear model the coefficients can be used to guide predictions.

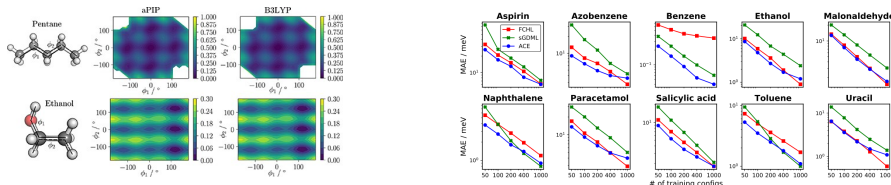- If we want low formation energy structures we can focus on this region.



Machine learning of material properties: Predictive and interpretable multilinear models, Allen and Tkatchenko, Sci. Adv. 8, eabm7185 (2022)
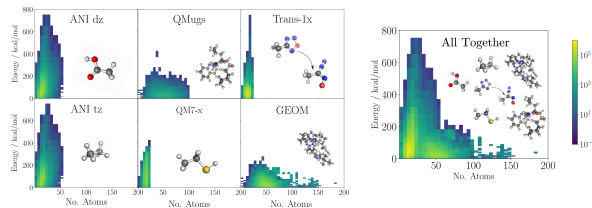
# Improving ML Models

## Architecture

Atomic permutationally invariant polynomials for fitting molecular force fields. Alice E. A. Allen, Geneviève Dusson, Christoph Ortner, and Gábor Csányi. Machine Learning: Science and Technology, 2021

Linear Atomic Cluster Expansion Force Fields for Organic Molecules: beyond RMSE., Kovacs, D. P., van der Oord, C., Kucera, J., Allen, A., Cole, D., Ortner, C., & Csanyi, G., JCTC, 2021



## Training Data

Machine learning potentials with Iterative Boltzmann Inversion: training to experiment, Sakib Matin, Alice Allen, Justin S Smith, Nicholas Lubbers, Ryan B Jadrich, Richard A Messerly, Benjamin T Nebgen, Ying Wai Li, Sergei Tretiak, Kipton Barros, JCTC, 2024
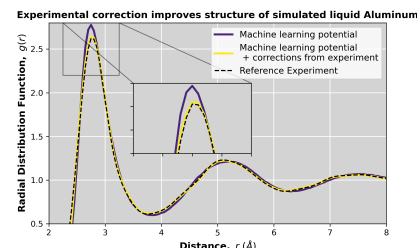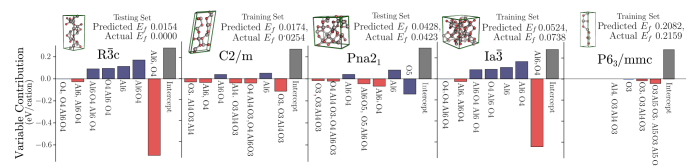


## Fitting Methods

Learning together: Towards foundation models for machine learning interatomic potentials with meta-learning, AEA Allen, N Lubbers, S Matin, J Smith, R Messerly, S Tretiak, K Barros, npj Computational Materials, 2024
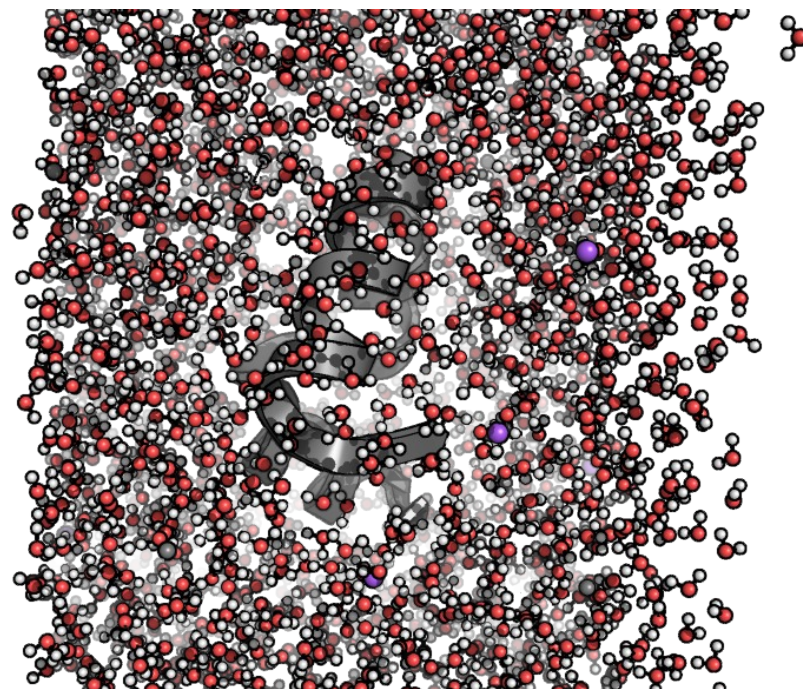


## Interpretability

Alice E. A. Allen and Alexandre Tkatchenko, Machine Learning of Material Properties: Predictive and Interpretable Multilinear Models, Science Advances, 2022

# Conclusion

- Atomistic simulations and the description of interactions between atoms and molecules can be greatly improved with ML

- Both building novel machine learning models and exploiting techniques from ML community is important for constructing effective models

- This can help us improve ML through better architecture, fitting data, fitting methods and interpretability

# Acknowledgements

**Classical Force Fields:**

Daniel Cole

Mike Payne

Joshua Horton

Michael Robertson

**aPIPs/ACE**:

Geneviève Dusson

Gábor Csányi

Cas van der Oord

Christoph Ortner

Dávid Péter Kovács

Emily Shinkle

Roxana Bujack

**Meta-learning & Reactive Potentials**

Kipton Barros

Nicholas Lubbers

Sakib Matin

Richard Messerly

Ben Nebgen

Sergei Tretiak

Justin Smith

Garnet Chan

**Interpretability**

Alexandre Tkatchenko

UNIVERSITY OF CAMBRIDGE

Los Alamos NATIONAL LABORATORY

Center for Nonlinear Studies

# Questions?