



CECAM UQ, June 2022



Predictive accuracy for real materials

- First-principles calculations promise quantitatively accurate simulations that make no use of experimental data
- Emergent physics from first principles: still a tremendous challenge
- Machine learning to the rescue



Grabowski et al., PRB (2009); Kapil, Engel, Rossi, MC, JCTC (2019)

Predictive accuracy for real materials

- First-principles calculations promise quantitatively accurate simulations that make no use of experimental data
- Emergent physics from first principles: still a tremendous challenge
- Machine learning to the rescue

2



Uncertainty and errors



We don't talk about errors

UQ for atomistic simulations

- Errors from model approximations:
 - 50 shades of DFT (and implementations!)
 - Classical vs quantum nuclei, Born-Oppenheimer vs non-adiabatic

Statistical convergence of averages, accelerated sampling

Scorecard

	0.3	0.3	0.6	1.0	0.9	0.3	1.5	0.6	0.9	0.4	1.0	0.4	0.4	6.3	13.5	1.1	2.1	0.7	1.4
0.3		0.1	0.5	0.9	0.8	0.2	1.5	0.6	0.8	0.4	1.0	0.5	0.3	6.3	13.4	1.1	2.1	0.7	1.4
0.3	0.1		0.5	0.9	0.8	0.2	1.5	0.6	0.8	0.4	0.9	0.5	0.3	6,3	13.4	1.1	2.1	0.7	1.4
0.6	0.5	0.5		0.8	0.6	0.4	1.5	0.6	0.8	0.6	1.0	0.7	0.5	6.3	13.2	1.0	1.9	0.6	1.3
1.0	0.9	0.9	0.8		0.9	0.9	1.8	0.9	1.3	1.0	1.4	1.0	0.9	6.4	13.0	1.2	1.8	1.0	1.6
0.9	0.8	0.8	0.6	0.9		0.8	1.7	0.7	1.1	0.8	1.3	1.0	0.8	6.5	13.2	1.1	1.8	0.8	1.5
0.3	0.2	0.2	0.4	0.9	0.8		1.5	0.5	0.8	0.3	1.0	0.5	0.3	6.2	13.4	1.0	2.0	0.6	1.4

Lejaeghere et al., Science (2016)

UQ for atomistic simulations

- Errors from model approximations:
 - 50 shades of DFT (and implementations!)
 - Classical vs quantum nuclei, Born-Oppenheimer vs non-adiabatic
- Statistical convergence of averages, accelerated sampling



UQ for atomistic simulations

- Errors from model approximations:
 - 50 shades of DFT (and implementations!)
 - Classical vs quantum nuclei, Born-Oppenheimer vs non-adiabatic
- Statistical convergence of averages, accelerated sampling



MC, Brain, Riordan, Manolopoulos Proc. Royal Soc. A (2012)

UQ for ML: the classics

- NN sanity check comparing the outputs of multiple fits
- Gaussian process regression provides built-in uncertainty measure
- Active learning driven by error estimates



Behler, J. Phys. Cond. Mat (2014)

UQ for ML: the classics

- NN sanity check comparing the outputs of multiple fits
- Gaussian process regression provides built-in uncertainty measure
- Active learning driven by error estimates



Deringer et al. Chem. Rev. (2021)

UQ for ML: the classics

- NN sanity check comparing the outputs of multiple fits
- Gaussian process regression provides built-in uncertainty measure
- Active learning driven by error estimates



Jinnouchi et al., Phys. Rev. B (2019)

What we do, what we need

- Most frameworks can be expressed in terms of *n*-body correlations of atom positions. Only difference the choice of basis
- Extension to a fully equivariant framework (NICE)
- ... to features to describe long-range interactions (LODE)
- ... and to message-passing, N-center features (MP-ACDC)



Introductory review: Musil et al., Chem. Rev. (2021)

- Most frameworks can be expressed in terms of *n*-body correlations of atom positions. Only difference the choice of basis
- Extension to a fully equivariant framework (NICE)
- ... to features to describe long-range interactions (LODE)
- ... and to message-passing, N-center features (MP-ACDC)



Introductory review: Musil et al., Chem. Rev. (2021)

- Most frameworks can be expressed in terms of *n*-body correlations of atom positions. Only difference the choice of basis
- Extension to a fully equivariant framework (NICE)
- ... to features to describe long-range interactions (LODE)
- ... and to message-passing, N-center features (MP-ACDC)



Introductory review: Musil et al., Chem. Rev. (2021)

- Most frameworks can be expressed in terms of *n*-body correlations of atom positions. Only difference the choice of basis
- Extension to a fully equivariant framework (NICE)
- ... to features to describe long-range interactions (LODE)
- ... and to message-passing, N-center features (MP-ACDC)



Latest & greatest: Nigam, Pozdnyakov, Fraux, MC, JCP (2022)

Want to learn vectors or general tensors?
 Need features that are *equivariant* to rotations

$$d_{\alpha}\left(\hat{R}A_{i}\right) = \sum_{q} \langle d|q\rangle \langle q|\hat{R}A; \overline{\rho_{i}^{\otimes\nu}; \alpha}\rangle$$



Grisafi, Wilkins, Csányi, & MC, PRL (2018)

Want to learn vectors or general tensors?
 Need features that are *equivariant* to rotations

$$d_{\alpha}\left(\hat{R}A_{i}\right) = \sum_{\boldsymbol{q}} \langle \boldsymbol{d} | \boldsymbol{q} \rangle \sum_{\boldsymbol{\alpha}'} R_{\boldsymbol{\alpha}\boldsymbol{\alpha}'} \langle \boldsymbol{q} | \boldsymbol{A}; \overline{\rho_{i}^{\otimes \nu}; \boldsymbol{\alpha}'} \rangle = \sum_{\boldsymbol{\alpha}'} R_{\boldsymbol{\alpha}\boldsymbol{\alpha}'} d_{\boldsymbol{\alpha}'} \left(\boldsymbol{A}_{i}\right)$$



Grisafi, Wilkins, Csányi, & MC, PRL (2018)

Want to learn vectors or general tensors?
 Need features that are *equivariant* to rotations



Machine learning you can trust

Want to learn vectors or general tensors?
 Need features that are *equivariant* to rotations



Lewis, Grisafi, MC, Rossi, JCTC (2021)

Want to learn vectors or general tensors?
 Need features that are *equivariant* to rotations



Nigam, Willatt, MC, JCP (2022)

Want to learn vectors or general tensors?
 Need features that are *equivariant* to rotations

$$\boldsymbol{y}_{\mu}^{\lambda}\left(\hat{\boldsymbol{R}}\boldsymbol{A}_{i}\right)=\sum_{\boldsymbol{q}}\left\langle\boldsymbol{d}|\boldsymbol{q}\right\rangle\sum_{\mu'}\boldsymbol{D}_{\mu\mu'}^{\lambda}\left(\hat{\boldsymbol{R}}\right)\left\langle\boldsymbol{q}|\boldsymbol{A};\overline{\rho_{i}^{\otimes\nu};\lambda\mu}\right\rangle$$



Nigam, Willatt, MC, JCP (2022)

- Predicting any property accessible to quantum calculations
- Realistic time and size scales, with first-principles accuracy and mapping of structural and functional properties



Kapil, Wilkins, Lan, MC, JCP (2020)

- Predicting any property accessible to quantum calculations
- Realistic time and size scales, with first-principles accuracy and mapping of structural and functional properties



N. Lopanitsyna, C. Ben Mahmoud, MC, Phys. Rev. Mater. (2021)

- Predicting any property accessible to quantum calculations
- Realistic time and size scales, with first-principles accuracy and mapping of structural and functional properties



V. Deringer et al., Nature (2021)

- Predicting any property accessible to quantum calculations
- Realistic time and size scales, with first-principles accuracy and mapping of structural and functional properties



Gigli et al., NPJ Comp. Mat. in press (2022)

Uncertainty quantification made easy

• Ensemble of N_{RS} models, trained on subsets of the train set

$$\bar{\mathbf{y}}(\mathbf{A}) = \frac{1}{N_{RS}} \sum_{j} \tilde{\mathbf{y}}_{j}(\mathbf{A}), \qquad \sigma^{2}(\mathbf{A}) = \frac{1}{N_{RS} - 1} \sum_{j} \left(\tilde{\mathbf{y}}_{j}(\mathbf{A}) - \bar{\mathbf{y}}(\mathbf{A}) \right)^{2}$$

Verify accuracy by the distribution of errors P (|ȳ (A) - y_A (A)| |σ (A))
 Adjust the spread by maximum likelihood ỹ_j ← ȳ + α (ỹ_j - ȳ)

Calibrated model can be used for easy uncertainty propgation

1. train set $\{(A, y_A)\}$



Musil et al., JCTC (2019)

• Ensemble of *N_{RS}* models, trained on subsets of the train set

$$\bar{\mathbf{y}}(\mathbf{A}) = \frac{1}{N_{RS}} \sum_{j} \tilde{\mathbf{y}}_{j}(\mathbf{A}), \qquad \sigma^{2}(\mathbf{A}) = \frac{1}{N_{RS} - 1} \sum_{j} \left(\tilde{\mathbf{y}}_{j}(\mathbf{A}) - \bar{\mathbf{y}}(\mathbf{A}) \right)^{2}$$

- Verify accuracy by the distribution of errors $P\left(\left|\bar{y}\left(\mathcal{A}
 ight)-y_{\mathcal{A}}\left(\mathcal{A}
 ight)\right|\right.\left|\sigma\left(\mathcal{A}
 ight)
 ight)$
- Adjust the spread by maximum likelihood $ilde{y}_i \leftarrow ar{y} + lpha \left(ilde{y}_i ar{y}
 ight)$
- Calibrated model can be used for easy uncertainty propgation
- **1.** train set $\ \left\{ (A,y_A)
 ight\}$



Machine learning you can trust

• Ensemble of *N_{RS}* models, trained on subsets of the train set

$$\bar{y}(A) = \frac{1}{N_{RS}} \sum_{j} \tilde{y}_{j}(A), \qquad \sigma^{2}(A) = \frac{1}{N_{RS} - 1} \sum_{j} \left(\tilde{y}_{j}(A) - \bar{y}(A) \right)^{2}$$

- Verify accuracy by the distribution of errors $P\left(\left|\bar{y}\left(\mathcal{A}\right)-y_{\mathcal{A}}\left(\mathcal{A}\right)\right| \; \left|\sigma\left(\mathcal{A}\right)\right)$
- Adjust the spread by maximum likelihood $ilde{y}_j \leftarrow ar{y} + lpha \left(ilde{y}_j ar{y}
 ight)$

Calibrated model can be used for easy uncertainty propgation

 (A, y_A)



Musil et al., JCTC (2019)

Machine learning you can trust

• Ensemble of *N_{RS}* models, trained on subsets of the train set

$$\bar{y}(A) = \frac{1}{N_{RS}} \sum_{j} \tilde{y}_{j}(A), \qquad \sigma^{2}(A) = \frac{1}{N_{RS} - 1} \sum_{i} \left(\tilde{y}_{i}(A) - \bar{y}(A) \right)^{2}$$

- Verify accuracy by the distribution of errors $P\left(\left|\bar{y}\left(\mathcal{A}
 ight)-y_{\mathcal{A}}\left(\mathcal{A}
 ight)
 ight|\,\left|\sigma\left(\mathcal{A}
 ight)
 ight)$
- Adjust the spread by maximum likelihood $\tilde{y}_j \leftarrow \bar{y} + \alpha \left(\tilde{y}_j \bar{y} \right)$
- Calibrated model can be used for easy uncertainty propgation



Uncertainty estimation on a shoestring

- Most expensive step is usually features/kernels evaluation. Generating committee by resampling allows inexpensive error estimation
- Same idea applies to linear & (sparse) kernel models
- Extension to NNs by only re-training output layer



Uncertainty estimation on a shoestring

- Most expensive step is usually features/kernels evaluation. Generating committee by resampling allows inexpensive error estimation
- Same idea applies to linear & (sparse) kernel models
- Extension to NNs by only re-training output layer



Uncertainty estimation on a shoestring

- Most expensive step is usually features/kernels evaluation. Generating committee by resampling allows inexpensive error estimation
- Same idea applies to linear & (sparse) kernel models
- Extension to NNs by only re-training output layer



Uncertainty estimation in action

- Error vs uncertainty for the atomization energies of QM9
 - Uncertainty \neq error! Spread of predictions is what matters
 - "Rejected" QM9 structures are highly uncertain
- Raman spectra with uncertainty propagation



Musil, Willatt, MC JCTC (2019)

Uncertainty estimation in action

- Error vs uncertainty for the atomization energies of QM9
 - Uncertainty \neq error! Spread of predictions is what matters
 - "Rejected" QM9 structures are highly uncertain
- Raman spectra with uncertainty propagation



Musil, Willatt, MC JCTC (2019)

Uncertainty estimation in action

- Error vs uncertainty for the atomization energies of QM9
 - Uncertainty \neq error! Spread of predictions is what matters
 - "Rejected" QM9 structures are highly uncertain
- Raman spectra with uncertainty propagation



N. Raimbault, A. Grisafi, MC, M. Rossi, New J. Phys. (2019)
Uncertainty estimation in action

- Error vs uncertainty for the atomization energies of QM9
 - Uncertainty \neq error! Spread of predictions is what matters
 - "Rejected" QM9 structures are highly uncertain
- Raman spectra with uncertainty propagation



N. Raimbault, A. Grisafi, MC, M. Rossi, New J. Phys. (2019)

• ML potential as a correction to a (semi)empirical baseline

Use baseline error (constant σ_b) and uncertainty (structure-dependent σ (A)) to get a weighted-baseline model that avoids instability

$$V(A) = V_b(A) + \overline{V}_{\delta}(A)$$



• ML potential as a correction to a (semi)empirical baseline

Use baseline error (constant σ_b) and uncertainty (structure-dependent σ (A)) to get a weighted-baseline model that avoids instability

 $V(A) = V_b(A) + \bar{V}_{\delta}(A)$



- ML potential as a correction to a (semi)empirical baseline
- Use baseline error (constant σ_b) and uncertainty (structure-dependent σ (A)) to get a weighted-baseline model that avoids instability



- ML potential as a correction to a (semi)empirical baseline
- Use baseline error (constant σ_b) and uncertainty (structure-dependent σ (A)) to get a weighted-baseline model that avoids instability



- Errors in ML-based thermodynamic averages combine effects on the observable σ_a^2 and those from sampling σ_{aV}^2
- A committee of predictions can be obtained from a single trajectory!

$$\langle a \rangle_{V^{(i)}} = \left\langle a e^{\beta \left(\overline{V} - V^{(i)} \right)} \right\rangle_{\overline{V}}$$

• Statistically stable estimates with a Cumulant Expansion Approximation



Machine learning you can trust

- Errors in ML-based thermodynamic averages combine effects on the observable σ_a^2 and those from sampling σ_{aV}^2
- A committee of predictions can be obtained from a single trajectory!

$$\langle a \rangle_{V^{(i)}} = \left\langle a \, e^{\beta \left(\bar{V} - V^{(i)} \right)} \right\rangle_{\bar{V}}$$

• Statistically stable estimates with a Cumulant Expansion Approximation



Machine learning you can trust

- Errors in ML-based thermodynamic averages combine effects on the observable σ_a^2 and those from sampling σ_{aV}^2
- A committee of predictions can be obtained from a single trajectory!

$$\langle a \rangle_{V^{(i)}} = \left\langle a \, e^{\beta \left(\bar{V} - V^{(i)} \right)} \right\rangle_{\bar{V}}$$

• Statistically stable estimates with a Cumulant Expansion Approximation



Machine learning you can trust

- Errors in ML-based thermodynamic averages combine effects on the observable σ²_a and those from sampling σ²_{aV}
- A committee of predictions can be obtained from a single trajectory!

$$\langle a \rangle_{\mathcal{V}^{(i)}} = \left\langle a \, e^{\beta \left(\bar{\mathcal{V}} - \mathcal{V}^{(i)} \right)} \right\rangle_{\bar{\mathcal{V}}}$$

Statistically stable estimates with a Cumulant Expansion Approximation



- Compute the phase diagram of Ga_xAs_{1-x}: interface-pinning simulations for a 2-component system
- DFT-accurate ML potential, i-PI+LAMMPS+PLUMED setup
- Estimate uncertainty in melting points



- Compute the phase diagram of Ga_xAs_{1-x}: interface-pinning simulations for a 2-component system
- DFT-accurate ML potential, i-PI+LAMMPS+PLUMED setup

Estimate uncertainty in melting points



- Compute the phase diagram of Ga_xAs_{1-x}: interface-pinning simulations for a 2-component system
- DFT-accurate ML potential, i-PI+LAMMPS+PLUMED setup
- Estimate uncertainty in melting points



- Compute the phase diagram of Ga_xAs_{1-x}: interface-pinning simulations for a 2-component system
- DFT-accurate ML potential, i-PI+LAMMPS+PLUMED setup
- Estimate uncertainty in melting points



- Compute the phase diagram of Ga_xAs_{1-x}: interface-pinning simulations for a 2-component system
- DFT-accurate ML potential, i-PI+LAMMPS+PLUMED setup
- Estimate uncertainty in melting points



- Compute the phase diagram of Ga_xAs_{1-x}: interface-pinning simulations for a 2-component system
- DFT-accurate ML potential, i-PI+LAMMPS+PLUMED setup
- Estimate uncertainty in melting points



- Compute the phase diagram of Ga_xAs_{1-x}: interface-pinning simulations for a 2-component system
- DFT-accurate ML potential, i-PI+LAMMPS+PLUMED setup
- Estimate uncertainty in melting points



Regio-selective catalysis

- Acid-catalyzed oxidation of phenol yields catechol and hydroquinone. Regioselectivity is poorly understood
- Very complex solution: phenol, H₂O₂, methanesulphonic acid. Explicit simulations require hybrid DFT accuracy



Rossi et al., JCTC (2020)

Regio-selective catalysis

- Acid-catalyzed oxidation of phenol yields catechol and hydroquinone. Regioselectivity is poorly understood
- Very complex solution: phenol, H₂O₂, methanesulphonic acid. Explicit simulations require hybrid DFT accuracy



Machine learning you can trust

Regio-selective catalysis

- Acid-catalyzed oxidation of phenol yields catechol and hydroquinone. Regioselectivity is poorly understood
- Very complex solution: phenol, H₂O₂, methanesulphonic acid. Explicit simulations require hybrid DFT accuracy



UQ for acid-base equilibria

- Energetics described with a DFTB baseline and a MLP correction. Accelerated by multiple time stepping
- PIGLET thermostatting & ring-polymer contraction for the quantum sampling, Plumed-driven metadynamics applied to the centroid
- Metadynamics sampling of the dissociation of CH₃SO₂OH.

• UQ for the free-energy profile!



Kapil, Behler, MC, JCP (2016); Zamani et al., Adv. Mater. (2020); Rossi et al., JCTC (2020)

UQ for acid-base equilibria

- Energetics described with a DFTB baseline and a MLP correction. Accelerated by multiple time stepping
- PIGLET thermostatting & ring-polymer contraction for the quantum sampling, Plumed-driven metadynamics applied to the centroid
- Metadynamics sampling of the dissociation of CH₃SO₂OH.
- UQ for the free-energy profile!



Kapil, Behler, MC, JCP (2016); Zamani et al., Adv. Mater. (2020); Rossi et al., JCTC (2020)

UQ for acid-base equilibria

- Energetics described with a DFTB baseline and a MLP correction. Accelerated by multiple time stepping
- PIGLET thermostatting & ring-polymer contraction for the quantum sampling, Plumed-driven metadynamics applied to the centroid
- Metadynamics sampling of the dissociation of CH₃SO₂OH.
- UQ for the free-energy profile!



Outlook

A software stack for atomistic machine learning

- Integrating ML and atomistic simulations: from representations to models to advanced MD
- Interoperability and data sharing with the rest of the ecosystem



Computational science & modeling @ EPFL

cosmo.epfl.ch







Building trust and understanding for ML

- ML is a purely inductive approach: dangerous in extrapolative contexts
- Physics-based constraints, and understanding of mathematical underpinnings, provide useful inductive biases
- Uncertainty quantification: easy, cheap, and universally applicable



B. Russel, The problems of philosophy, Chapter VI

Building trust and understanding for ML

- ML is a purely inductive approach: dangerous in extrapolative contexts
- Physics-based constraints, and understanding of mathematical underpinnings, provide useful inductive biases
- Uncertainty quantification: easy, cheap, and universally applicable



Pozdnyakov, MC, arXiv:2201.07136 (2022)

Building trust and understanding for ML

- ML is a purely inductive approach: dangerous in extrapolative contexts
- Physics-based constraints, and understanding of mathematical underpinnings, provide useful inductive biases
- Uncertainty quantification: easy, cheap, and universally applicable



Imbalzano et al., JCP (2021); Imbalzano & MC, Phys. Rev. Materials (2021)

Backup slides

Want to learn vectors or general tensors?
 Need features that are *equivariant* to rotations

$$d_{\alpha}(A_{i}) = \sum_{q} \langle d|q \rangle \langle q|A; \overline{\rho_{i}^{\otimes \nu}; \alpha} \rangle$$
$$d_{\alpha}(\hat{R}A_{i}) = \sum_{q} \langle d|q \rangle \langle q|\hat{R}A; \overline{\rho_{i}^{\otimes \nu}; \alpha} \rangle$$



Glielmo, Sollich, De Vita, PRB (2017); Grisafi, Wilkins, Csányi, & MC, PRL (2018); Veit et al., JCP (2020)

 Want to learn vectors or general tensors? Need features that are *equivariant* to rotations

$$d_{\alpha}(A_{i}) = \sum_{q} \langle d|q \rangle \langle q|A; \overline{\rho_{i}^{\otimes \nu}}; \alpha \rangle$$
$$d_{\alpha}(\hat{R}A_{i}) = \sum_{q} \langle d|q \rangle \sum_{\alpha'} R_{\alpha\alpha'} \langle q|A; \overline{\rho_{i}^{\otimes \nu}}; \alpha' \rangle = \sum_{\alpha'} R_{\alpha\alpha'} d_{\alpha'}(A_{i})$$



Glielmo, Sollich, De Vita, PRB (2017); Grisafi, Wilkins, Csányi, & MC, PRL (2018); Veit et al., JCP (2020)

 Want to learn vectors or general tensors? Need features that are *equivariant* to rotations

$$d_{\alpha}(A_{i}) = \sum_{q} \langle d|q \rangle \langle q|A; \overline{\rho_{i}^{\otimes \nu}; \alpha} \rangle$$

$$d_{\alpha}(\hat{R}A_{i}) = \sum_{q} \langle d|q \rangle \sum_{\alpha'} R_{\alpha\alpha'} \langle q|A; \overline{\rho_{i}^{\otimes \nu}; \alpha'} \rangle = \sum_{\alpha'} R_{\alpha\alpha'} d_{\alpha'}(A_{i})$$

$$\xrightarrow{5A}$$
Partial charge / e



Machine learning you can trust

 Want to learn vectors or general tensors? Need features that are *equivariant* to rotations

$$\begin{aligned} d_{\alpha}_{\alpha}\left(\boldsymbol{A}_{i}\right) &= \sum_{\boldsymbol{q}} \left\langle \boldsymbol{d} | \boldsymbol{q} \right\rangle \left\langle \boldsymbol{q} | \boldsymbol{A}; \overline{\rho_{i}^{\otimes \nu}; \alpha} \right\rangle \\ \boldsymbol{y}_{\mu}^{\lambda}\left(\hat{\boldsymbol{R}}\boldsymbol{A}_{i}\right) &= \sum_{\boldsymbol{q}} \left\langle \boldsymbol{d} | \boldsymbol{q} \right\rangle \sum_{\mu'} \boldsymbol{D}_{\mu\mu'}^{\lambda}\left(\hat{\boldsymbol{R}}\right) \left\langle \boldsymbol{q} | \boldsymbol{A}; \overline{\rho_{i}^{\otimes \nu}; \lambda\mu} \right\rangle \end{aligned}$$



Grisafi, Wilkins, Csányi, & MC, PRL (2018); Willatt, Musil, & MC, JCP (2019)

Molecular polarizabilities at the CCSD level

- Symmetry-adapted tensorial model of polarizabilities
- Training on high-end CCSD calculations of small molecules
- Extrapolate to larger molecules (tested up to aciclovir C₈H₁₁N₅O₃)
- Better-than-DFT accuracy try it on alphaml.org



Wilkins, Grisafi, Yang, Lao, DiStasio, MC, PNAS (2019);

Transferable model of the electron density

- Expand the charge density on atom-centered basis $\phi_k \equiv R_n Y_l^m$
- Learning of coefficients with symmetry-adapted kernels
- Highly transferable: train on small molecules, predict on polypeptides
- Condensed phase implementation. Scaling up complexity!



Grisafi, Wilkins, Meyer, Fabrizio, Corminboeuf, MC, ACS Central Science (2019);

Meyer, Grisafi, Fabrizio, MC, Corminboeuf, Chem. Sci., (2019)

Transferable model of the electron density

- Expand the charge density on atom-centered basis $\phi_k \equiv R_n Y_I^m$
- Learning of coefficients with symmetry-adapted kernels
- Highly transferable: train on small molecules, predict on polypeptides
- Condensed phase implementation. Scaling up complexity!



Grisafi, Wilkins, Meyer, Fabrizio, Corminboeuf, MC, ACS Central Science (2019);

Meyer, Grisafi, Fabrizio, MC, Corminboeuf, Chem. Sci., (2019)
Transferable model of the electron density

- Expand the charge density on atom-centered basis $\phi_k \equiv R_n Y_I^m$
- Learning of coefficients with symmetry-adapted kernels
- Highly transferable: train on small molecules, predict on polypeptides
- Condensed phase implementation. Scaling up complexity!



Lewis, Grisafi, MC, Rossi, JCTC (2021)

- Molecular Hamiltonian in an atomic orbital basis
- Decompositon into irreducible symmetric blocks
- Learn with a fully equivariant model
- Consistent with molecular orbital theory by design



- Molecular Hamiltonian in an atomic orbital basis
- Decompositon into irreducible symmetric blocks
- Learn with a fully equivariant model
- Consistent with molecular orbital theory by design



- Molecular Hamiltonian in an atomic orbital basis
- Decompositon into irreducible symmetric blocks
- Learn with a fully equivariant model
- Consistent with molecular orbital theory by design



- Molecular Hamiltonian in an atomic orbital basis
- Decompositon into irreducible symmetric blocks
- Learn with a fully equivariant model
- Consistent with molecular orbital theory by design



- Molecular Hamiltonian in an atomic orbital basis
- Decompositon into irreducible symmetric blocks
- Learn with a fully equivariant model
- Consistent with molecular orbital theory by design



ML, ML everywhere (and all the errors did shrink)

Ab initio thermodynamics of water

- Thermodynamics of water liquid/ice-Ih/ice-Ic.
- Hybrid DFT+D3 level, with quantum nuclei
- ML, for sampling, promote to full-DFT with free-energy perturbation
- Melting point of water within ~5K. $\Delta G_{lh/lc}$ within 0.1 meV/molecule



Cheng, Engel, Behler, Dellago, MC, PNAS (2019)

Ab initio thermodynamics of water

- Thermodynamics of water liquid/ice-Ih/ice-Ic.
- Hybrid DFT+D3 level, with quantum nuclei
- ML, for sampling, promote to full-DFT with free-energy perturbation
- Melting point of water within ~5K. $\Delta G_{lh/lc}$ within 0.1 meV/molecule



Cheng, Engel, Behler, Dellago, MC, PNAS (2019)

Ab initio thermodynamics of water

- Thermodynamics of water liquid/ice-Ih/ice-Ic.
- Hybrid DFT+D3 level, with quantum nuclei
- ML, for sampling, promote to full-DFT with free-energy perturbation
- Melting point of water within ~5K. $\Delta G_{lh/lc}$ within 0.1 meV/molecule



Cheng, Engel, Behler, Dellago, MC, PNAS (2019)

Liquid-liquid phase transition in dense H

- Atomic/molecular fluid for H at giant-planets conditions
- Difficult to characterize the transition experimentally
- First-order transition predicted by ab initio MD: finite-size effect?
- ML potential and polyamorphic solution model finds critical point to lie on the solid-liquid phase line



Cheng, Mazzola, Pickard, MC, Nature (2020)

Liquid-liquid phase transition in dense H

- Atomic/molecular fluid for H at giant-planets conditions
- Difficult to characterize the transition experimentally
- First-order transition predicted by ab initio MD: finite-size effect?
- ML potential and polyamorphic solution model finds critical point to lie on the solid-liquid phase line



Cheng, Mazzola, Pickard, MC, Nature (2020)

- Synthesis of GaAs nanostructures: solid-liquid-gas process, Ga_l/GaAs interface based on droplets of molten Ga
- Problem: controlling the zinc-blende/wurtzite polymorphism
- Role of surface polarity and liquid structure on defect rate
- MD simulation with DFT/NN to study liquid ordering at the interfaces
- More to come: transferable model of the Ga_xAs_{1-x} binary



- Synthesis of GaAs nanostructures: solid-liquid-gas process, Ga_l/GaAs interface based on droplets of molten Ga
- Problem: controlling the zinc-blende/wurtzite polymorphism
- Role of surface polarity and liquid structure on defect rate
- MD simulation with DFT/NN to study liquid ordering at the interfaces
 More to come: transferable model of the Ga_xAs_{1-x} binary



- Synthesis of GaAs nanostructures: solid-liquid-gas process, Ga_l/GaAs interface based on droplets of molten Ga
- Problem: controlling the zinc-blende/wurtzite polymorphism
- Role of surface polarity and liquid structure on defect rate
- MD simulation with DFT/NN to study liquid ordering at the interfaces
 More to come: transferable model of the GavAs1 with binary



- Synthesis of GaAs nanostructures: solid-liquid-gas process, Ga_l/GaAs interface based on droplets of molten Ga
- Problem: controlling the zinc-blende/wurtzite polymorphism
- Role of surface polarity and liquid structure on defect rate
- MD simulation with DFT/NN to study liquid ordering at the interfaces
- More to come: transferable model of the Ga_xAs_{1-x} binary



- Synthesis of GaAs nanostructures: solid-liquid-gas process, Ga_l/GaAs interface based on droplets of molten Ga
- Problem: controlling the zinc-blende/wurtzite polymorphism
- Role of surface polarity and liquid structure on defect rate
- MD simulation with DFT/NN to study liquid ordering at the interfaces
- More to come: transferable model of the Ga_xAs_{1-x} binary



Imbalzano, MC, Phys. Rev. Mat. (2021)

Machine learning you can trust

ML-powered NMR crystallograpy

- Solid-state NMR relies on GIPAW-DFT to determine crystal structure
- ML model trained on 2000 CSD structures. DFT accuracy of chemical shielding predictions (RMSE H: 0.5, C: 5, N: 13, O: 18 ppm)
- Precise enough to do structure determination try it on shiftml.org



Paruzzo, Hofstetter, Musil, De, MC, Emsley, Nature Comm. (2018) [data: CSD-500]

Quantum nuclei in NMR experiments

- Combining a ML potential and NMR shielding model \rightarrow full ab initio modeling of NMR at constant temperature
- Thermal and quantum fluctuations affect the NMR *shifts*: comparable to typical errors of reference DFT calculations



Engel, Kapil, MC J. Phys. Chem. Lett. (2021)

Consolidating the mathematical foundations

Are representations complete?

- More fundamental: are representations complete (injective)?
- Well-known: 2-body correlations (distances) are ambiguous
- Surprise: neither are 3 (angles) and 4 (dihedrals) body features!

• Limits the asymptotic accuracy of models



Boutin, Kemper, Ann. Adv. Math. (2004); Figure from Bartók, Kondor, Csányi, PRB (2013)

Are representations complete?

- More fundamental: are representations complete (injective)?
- Well-known: 2-body correlations (distances) are ambiguous
- Surprise: neither are 3 (angles) and 4 (dihedrals) body features!

• Limits the asymptotic accuracy of models



Pozdniakov, Willatt, Bartók, Ortner, Csányi, MC PRL (2020)

Are representations complete?

- More fundamental: are representations complete (injective)?
- Well-known: 2-body correlations (distances) are ambiguous
- Surprise: neither are 3 (angles) and 4 (dihedrals) body features!
- Limits the asymptotic accuracy of models



Pozdniakov, Willatt, Bartók, Ortner, Csányi, MC PRL (2020)

How about graph convolution?

- Atoms: nodes in a fully-connected network. Edges are decorated by (functions of) interatomic distances *r_{ij}*
- Each node is augmented with information on its neighbors and their distance: $h(A_i) = (a_i, \{(a_j, r_{ij})\})$
- The multiset of neighbors and edges is hashed, and used as a label to describe the nodes. The process can be iterated



SchNET: Schütt et al., JCP (2018)

How about graph convolution?

- Atoms: nodes in a fully-connected network. Edges are decorated by (functions of) interatomic distances *r_{ij}*
- Each node is augmented with information on its neighbors and their distance: h (A_i) = (a_i, {(a_j, r_{ij})})
- The multiset of neighbors and edges is hashed, and used as a label to describe the nodes. The process can be iterated



How about graph convolution?

- Atoms: nodes in a fully-connected network. Edges are decorated by (functions of) interatomic distances *r_{ij}*
- Each node is augmented with information on its neighbors and their distance: h (A_i) = (a_i, {(a_j, r_{ij})})
- The multiset of neighbors and edges is hashed, and used as a label to describe the nodes. The process can be iterated



Graph convolution, pros and cons

- Bad news: there are known *discrete* graphs that cannot be distinguished by this procedure (W-L test)
- Good news: things look good for molecular graphs (*fully-connected*, distance-decorated 3D point clouds);
- Distance-GNN resolve all known ACDC counterexamples



Sato, arxiv:2003.04078

Graph convolution, pros and cons

- Bad news: there are known *discrete* graphs that cannot be distinguished by this procedure (W-L test)
- Good news: things look good for molecular graphs (*fully-connected*, distance-decorated 3D point clouds);
- Distance-GNN resolve all known ACDC counterexamples



Bartók et al. PRB (2013); Pozdnyakov et al. PRL (2020)

- A family of 3D point clouds with degenerate pairs for GNN. Key idea: the distance matrix is identical, except for a swap
- Can be folded to give finite 3D structures
- Hard limit to the accuracy for plausible molecular geometries
- Modern architectures that use angular/directional information (and simple models based on |p^{⊗2}/_i) are immune



Pozdnyakov, MC, arXiv:2201.07136 (2022)

- A family of 3D point clouds with degenerate pairs for GNN. Key idea: the distance matrix is identical, except for a swap
- Can be folded to give finite 3D structures
- Hard limit to the accuracy for plausible molecular geometries
- Modern architectures that use angular/directional information (and simple models based on $|\overline{\rho_i^{\otimes 2}}\rangle$) are immune



Pozdnyakov, MC, arXiv:2201.07136 (2022)

- A family of 3D point clouds with degenerate pairs for GNN. Key idea: the distance matrix is identical, except for a swap
- Can be folded to give finite 3D structures
- Hard limit to the accuracy for plausible molecular geometries
- Modern architectures that use angular/directional information (and simple models based on $\overline{|\rho_i^{\otimes 2}\rangle}$) are immune



Pozdnyakov, MC, arXiv:2201.07136 (2022)

- A family of 3D point clouds with degenerate pairs for GNN. Key idea: the distance matrix is identical, except for a swap
- Can be folded to give finite 3D structures
- Hard limit to the accuracy for plausible molecular geometries
- Modern architectures that use angular/directional information (and simple models based on $|\overline{\rho_i^{\otimes 2}}\rangle$) are immune



Nigam, Fraux, MC, arXiv:2202.01566 (2022)

The charged elephant in the other room

Understanding the range of interactions

- Representations are built for different cutoff radii
- Dimensionality/accuracy tradeoff: a measure of the range of interactions
- Multi-scale kernels $K(A, B) = \sum_{i} w_i K_i(A, B)$ yield the best of all worlds



Bartók, De, Poelking, Kermode, Bernstein, Csányi, MC, Science Advances (2017) [data: QM9, von Lilienfeld&C]

Understanding the range of interactions

- Representations are built for different cutoff radii
- Dimensionality/accuracy tradeoff: a measure of the range of interactions
- Multi-scale kernels $K(A, B) = \sum_{i} w_i K_i(A, B)$ yield the best of all worlds



Understanding the range of interactions

- Representations are built for different cutoff radii
- Dimensionality/accuracy tradeoff: a measure of the range of interactions
- Multi-scale kernels $K(A, B) = \sum_{i} w_i K_i(A, B)$ yield the best of all worlds


The problem with electrostatics

• 1/r decay \rightarrow pathological convergence of with interaction cutoff

• Capturing true long-range effects with local models is hopeless



The problem with electrostatics

- 1/r decay \rightarrow pathological convergence of with interaction cutoff
- Capturing true long-range effects with local models is hopeless



• Idea: *local* representation that reflects long-range *asymptotics*

- Atom-density potential $\langle a\mathbf{r} | V \rangle = \int \langle a\mathbf{r}' | \rho \rangle / |\mathbf{r}' \mathbf{r}| d\mathbf{r}'$
- ② Usual gig: symmetrize, decompose locally, learn!
- Efficient evaluation in reciprocal space



Grisafi, MC, JCP (2019)

- Idea: local representation that reflects long-range asymptotics
 - Atom-density potential $\langle a\mathbf{r} | \mathbf{V} \rangle = \int \langle a\mathbf{r}' | \rho \rangle / |\mathbf{r}' \mathbf{r}| d\mathbf{r}'$

Efficient evaluation in reciprocal space





Grisafi. MC. JCP (2019)

Machine learning you can trust

- Idea: *local* representation that reflects long-range *asymptotics*
 - Atom-density potential $\langle a\mathbf{r}|\mathbf{V}\rangle = \int \langle a\mathbf{r}'|\rho\rangle / |\mathbf{r}' \mathbf{r}| d\mathbf{r}'$
 - Usual gig: symmetrize, decompose locally, learn!
- Efficient evaluation in reciprocal space



Grisafi, MC, JCP (2019)

- Idea: local representation that reflects long-range asymptotics
 - Atom-density potential $\langle a\mathbf{r}|\mathbf{V}\rangle = \int \langle a\mathbf{r}'|\rho\rangle / |\mathbf{r}' \mathbf{r}| d\mathbf{r}'$
 - Osual gig: symmetrize, decompose locally, learn!
- Efficient evaluation in reciprocal space





Grisafi, MC, JCP (2019)

42 Michele Ceriotti cosmo.epfl.ch

Machine learning you can trust

Binding of charged molecules, and beyond

- A challenging test: rigid-molecule binding curves of charged dimers from the BioFragmentsDB
- "Multi-scale" LODE features | p_i ⊗ V_i ⟩ map to multipole electrostatics but enable learning all sorts of long-range physics



Grisafi, MC, JCP (2019); Grisafi, Nigam, MC, Chem. Sci. (2021)

Machine learning you can trust

Integrated ML models beyond size and time limits

- Predicting any property accessible to quantum calculations
- Realistic time and size scales, with first-principles accuracy and mapping of structural and functional properties



Kapil, Wilkins, Lan, MC, JCP (2020)

- Predicting any property accessible to quantum calculations
- Realistic time and size scales, with first-principles accuracy and mapping of structural and functional properties



N. Lopanitsyna, C. Ben Mahmoud, MC, Phys. Rev. Mater. (2021)

- Predicting any property accessible to quantum calculations
- Realistic time and size scales, with first-principles accuracy and mapping of structural and functional properties



V. Deringer et al., Nature (2021)

- Predicting any property accessible to quantum calculations
- Realistic time and size scales, with first-principles accuracy and mapping of structural and functional properties



Gigli et al., arxiv: 2111.05129

A software stack for atomistic machine learning

- Integrating ML and atomistic simulations: from representations to models to advanced MD
- Interoperability and data sharing with the rest of the ecosystem





A Dirac notation for ML



- A representation maps a structure A (or one environment A_i) to a vector discretized by a feature index Q
- Bra-ket notation (Q|A; rep.) indicates in an abstract way this mapping, leaving plenty of room to express the details of a representation
- Dirac-like notation reflects naturally a change of basis, the construction of a kernel, or a linear model

$$\langle \mathbf{Y} | \mathbf{A}
angle = \int \mathrm{d} \mathbf{Q} \left< \mathbf{Y} | \mathbf{Q}
ight> \left< \mathbf{Q} | \mathbf{A}
ight>$$

Willatt, Musil, MC, JCP (2019); https://tinyurl.com/dirac-rep

A Dirac notation for ML



- A representation maps a structure A (or one environment A_i) to a vector discretized by a feature index Q
- Bra-ket notation (Q|A; rep.) indicates in an abstract way this mapping, leaving plenty of room to express the details of a representation
- Dirac-like notation reflects naturally a change of basis, the construction of a kernel, or a linear model

$$k(A,A') = \langle A|A'
angle pprox \int \mathrm{d} Q \langle A|Q
angle \, \langle Q|A'
angle$$

Willatt, Musil, MC, JCP (2019); https://tinyurl.com/dirac-rep

A Dirac notation for ML



- A representation maps a structure A (or one environment A_i) to a vector discretized by a feature index Q
- Bra-ket notation (Q|A; rep.) indicates in an abstract way this mapping, leaving plenty of room to express the details of a representation
- Dirac-like notation reflects naturally a change of basis, the construction of a kernel, or a linear model

$${\it E}({\it A}) = \langle {\it E} | {\it A}
angle pprox \int {
m d} {\it Q} \, \langle {\it E} | {\it Q}
angle \, \langle {\it Q} | {\it A}
angle$$

Willatt, Musil, MC, JCP (2019); https://tinyurl.com/dirac-rep

- Understanding what goes into a representation is key to achieve meaningful results from automated data analytics
- Example: you don't always want to have rotational invariance



- Understanding what goes into a representation is key to achieve meaningful results from automated data analytics
- Example: you don't always want to have rotational invariance



- Understanding what goes into a representation is key to achieve meaningful results from automated data analytics
- Example: you don't always want to have rotational invariance



- Understanding what goes into a representation is key to achieve meaningful results from automated data analytics
- Example: you don't always want to have rotational invariance



- Understanding what goes into a representation is key to achieve meaningful results from automated data analytics
- Example: you don't always want to have rotational invariance



Variations on a theme

- Most of the existing density-based representations and kernels emerge as special cases of this framework
 - Basis set choice e.g. plane waves basis for $|\overline{\rho_i^{\otimes 2}}\rangle$ (Ziletti et al. N.Comm 2018)
 - Projection on symmetry functions (Behler-Parrinello, DeepMD)

$$\langle \mathbf{k} | \mathbf{A}; \overline{\rho^{\otimes 2}} \rangle = \sum_{ij \in \mathbf{A}} e^{i \mathbf{k} \cdot \mathbf{r}_{ij}}$$



Willatt, Musil, MC, JCP (2019), https://arxiv.org/pdf/1807.00408

Variations on a theme

- Most of the existing density-based representations and kernels emerge as special cases of this framework
 - Basis set choice e.g. plane waves basis for $|\rho_i^{\otimes 2}\rangle$ (Ziletti et al. N.Comm 2018)
 - Projection on symmetry functions (Behler-Parrinello, DeepMD)



Willatt, Musil, MC, JCP (2019), https://arxiv.org/pdf/1807.00408

Machine learning you can trust

- Quantitative comparison of relative information content of different features, metrics & kernels
- Feature space Reconstruction Error (FRE): linearly-embeddable mutual information



Goscinski, Fraux, MC, MLST (2021)

- Quantitative comparison of relative information content of different features, metrics & kernels
- Feature space Reconstruction Error (FRE): linearly-embeddable mutual information



 $\operatorname{GFRE}(\mathcal{F}, \mathcal{F}')$

Goscinski, Fraux, MC, MLST (2021)

- Quantitative comparison of relative information content of different features, metrics & kernels
- Feature space Reconstruction Error (FRE): linearly-embeddable mutual information



Goscinski, Fraux, MC, MLST (2021)

- Quantitative comparison of relative information content of different features, metrics & kernels
- Feature space Reconstruction Error (FRE): linearly-embeddable mutual information

$$\mathsf{GFRE}(\mathcal{F} \to \mathcal{F}') = \min_{\boldsymbol{P} \in \mathbb{R}^{n_{\mathcal{F}} \times n_{\mathcal{F}'}}} \|\boldsymbol{X}_{\mathcal{F}'} - \boldsymbol{X}_{\mathcal{F}} \boldsymbol{P}\|$$



Density expansion and SOAP

• What if we use radial functions and spherical harmonics?

- Symmetrized tensor product ightarrow SOAP power spectrum!
- Easily generalized to higher body order. δ -distribution limit \rightarrow atomic cluster expansion



Bartók, Kondor, Csányi, PRB (2013); Willatt, Musil, MC, JCP (2019); Drautz, PRB (2019)

Density expansion and SOAP

- What if we use radial functions and spherical harmonics?
- Symmetrized tensor product \rightarrow SOAP power spectrum!
- Easily generalized to higher body order. δ -distribution limit \rightarrow atomic cluster expansion



$$\langle nn'l | \overline{\rho_i^{\otimes 2}} \rangle = \sum_m \langle nlm | \rho_i \rangle^* \langle n'lm | \rho_i \rangle$$
$$p_{nn'l} = \sum_m c_{nlm}^* c_{n'lm}$$

Bartók, Kondor, Csányi, PRB (2013); Willatt, Musil, MC, JCP (2019); Drautz, PRB (2019)

Density expansion and SOAP

- What if we use radial functions and spherical harmonics?
- Symmetrized tensor product \rightarrow SOAP power spectrum!
- Easily generalized to higher body order. δ -distribution limit \rightarrow atomic cluster expansion



Bartók, Kondor, Csányi, PRB (2013); Willatt, Musil, MC, JCP (2019); Drautz, PRB (2019)

- Construction of a three-body (u = 2) invariant atomic descriptor
 - O Define relative position of neighbors (translation-invariant)
 - Positions are transformed in a neighbor density (permutation invariant)
 - Symmetrize over rotations a tensor product of the neighbor densities
 - O This is equivalent to a function of two distances and one angle

 - Linear model \Rightarrow 3-body potential!



Bartók, Kondor, Csányi, PRB (2013)

- Construction of a three-body (u = 2) invariant atomic descriptor
 - O Define relative position of neighbors (translation-invariant)
 - Positions are transformed in a neighbor density (permutation invariant)
 - Symmetrize over rotations a tensor product of the neighbor densities
 - O This is equivalent to a function of two distances and one angle
 - **3** $g \to \delta$ limit \Rightarrow list of 2-neighbors tuples $(\mathbf{r}_{j_1i}, \mathbf{r}_{j_2i}, \hat{\mathbf{r}}_{j_1i} \cdot \hat{\mathbf{r}}_{j_2i})$
 - Linear model \Rightarrow 3-body potential!



Bartók, Kondor, Csányi, PRB (2013)

- Construction of a three-body (u = 2) invariant atomic descriptor
 - O Define relative position of neighbors (translation-invariant)
 - Ositions are transformed in a neighbor density (permutation invariant)
 - Symmetrize over rotations a tensor product of the neighbor densities
 - O This is equivalent to a function of two distances and one angle
 - **3** $g \to \delta$ limit \Rightarrow list of 2-neighbors tuples $(\mathbf{r}_{j_1 i}, \mathbf{r}_{j_2 i}, \hat{\mathbf{r}}_{j_1 i} \cdot \hat{\mathbf{r}}_{j_2 i})$
 - Linear model \Rightarrow 3-body potential!



- Construction of a three-body (u = 2) invariant atomic descriptor
 - O Define relative position of neighbors (translation-invariant)
 - Ositions are transformed in a neighbor density (permutation invariant)
 - Symmetrize over rotations a tensor product of the neighbor densities
 - O This is equivalent to a function of two distances and one angle
 - **5** $g \to \delta$ limit \Rightarrow list of 2-neighbors tuples $(r_{j_1 i}, r_{j_2 i}, \hat{\mathbf{r}}_{j_1 i} \cdot \hat{\mathbf{r}}_{j_2 i})$
 - Linear model \Rightarrow 3-body potential!



Willatt, Musil, MC, JCP (2019)
Two-neighbors descriptors

- Construction of a three-body (u = 2) invariant atomic descriptor
 - O Define relative position of neighbors (translation-invariant)
 - Ositions are transformed in a neighbor density (permutation invariant)
 - Symmetrize over rotations a tensor product of the neighbor densities
 - O This is equivalent to a function of two distances and one angle
 - **3** $g \to \delta$ limit \Rightarrow list of 2-neighbors tuples $(\mathbf{r}_{j_1 i}, \mathbf{r}_{j_2 i}, \hat{\mathbf{r}}_{j_1 i} \cdot \hat{\mathbf{r}}_{j_2 i})$

• Linear model \Rightarrow 3-body potential!



Willatt, Musil, MC, JCP (2019)

Two-neighbors descriptors

- Construction of a three-body (u = 2) invariant atomic descriptor
 - O Define relative position of neighbors (translation-invariant)
 - Positions are transformed in a neighbor density (permutation invariant)
 - Symmetrize over rotations a tensor product of the neighbor densities
 - O This is equivalent to a function of two distances and one angle
 - $\mathbf{9} \ \mathbf{g} \to \delta \text{ limit} \Rightarrow \text{list of 2-neighbors tuples } (\mathbf{r}_{j_1i}, \mathbf{r}_{j_2i}, \hat{\mathbf{r}}_{j_1i} \cdot \hat{\mathbf{r}}_{j_2i})$
 - Linear model \Rightarrow 3-body potential!



Willatt, Musil, MC, JCP (2019)

Density trick in an $\langle nlm |$ basis

- The symmetrized correlations can be computed in closed form using a discrete basis
 - The neighbor density can be expanded on a basis of radial functions $\langle x|n \rangle \equiv R_n(x)$ and spherical harmonics $\langle \hat{\mathbf{x}}|lm \rangle \equiv Y_l^m(\hat{\mathbf{x}})$
 - Spherical harmonics transform linearly under rotations based on Wigner rotation matrices $\mathbf{D}^{l}\left(\hat{R}\right)$
 - Orthogonality of Wigner matrices yields the SOAP powerspectrum



Density trick in an $\langle nlm |$ basis

- The symmetrized correlations can be computed in closed form using a discrete basis
 - The neighbor density can be expanded on a basis of radial functions $\langle x|n \rangle \equiv R_n(x)$ and spherical harmonics $\langle \hat{\mathbf{x}}|lm \rangle \equiv Y_l^m(\hat{\mathbf{x}})$
 - Spherical harmonics transform linearly under rotations based on Wigner rotation matrices $\mathbf{D}^l\left(\hat{R}\right)$
 - Orthogonality of Wigner matrices yields the SOAP powerspectrum



Density trick in an $\langle nlm |$ basis

- The symmetrized correlations can be computed in closed form using a discrete basis
 - The neighbor density can be expanded on a basis of radial functions $\langle x|n \rangle \equiv R_n(x)$ and spherical harmonics $\langle \hat{\mathbf{x}}|lm \rangle \equiv Y_l^m(\hat{\mathbf{x}})$
 - Spherical harmonics transform linearly under rotations based on Wigner rotation matrices $\mathbf{D}^l\left(\hat{R}\right)$
 - Orthogonality of Wigner matrices yields the SOAP powerspectrum

$$\begin{split} &\int \mathrm{d}\hat{R} \sum_{kk'} D_{mk}^{l}(\hat{R}) D_{m'k'}^{l'}(\hat{R}) \propto \\ &\delta_{ll'} \delta_{mm'} \delta_{kk'} \\ &\langle nn'l | A; \overline{\rho_{i}^{\otimes 2}} \rangle = \\ &\sum_{m} \langle nlm | A; \rho_{i} \rangle \langle n'lm | A; \rho_{i} \rangle \end{split}$$

A hierarchy of equivariant features

• A generalization of the definition yields *N*-body features that transform like angular momenta

$$\langle \boldsymbol{X} | \overline{\rho_{\boldsymbol{i}}^{\otimes \nu}; \sigma; \lambda \mu} \rangle$$

 Recursive construction based on sums of angular momenta and an expansion of the atom density

$$\langle \mathbf{n}_{1} \mathbf{l}_{1} \mathbf{k}_{1} | \overline{\rho_{j}^{\otimes 1}; \lambda \mu} \rangle \equiv \langle \mathbf{n}_{1} \lambda (-\mu) | \rho_{i} \rangle \, \delta_{\mathbf{l}_{1} \lambda} \delta_{\mathbf{k}_{1} \lambda} \delta_{\sigma 1} \equiv \langle \mathbf{n}_{1} | \overline{\rho_{j}^{\otimes 1}; \lambda \mu} \rangle$$

$$\langle \dots; \mathbf{n}_{\nu} \mathbf{l}_{\nu} \mathbf{k}_{\nu}; \mathbf{n} \mathbf{l} \mathbf{k} | \overline{\rho_{i}^{\otimes (\nu+1)}}; \sigma; \lambda \mu \rangle = \delta_{\sigma((-1)^{l+k+\lambda} \mathbf{s})} \mathbf{c}_{k\lambda} \times \\ \sum_{qm} \langle \mathbf{l}m; \mathbf{k}q | \lambda \mu \rangle \langle \mathbf{n} | \overline{\rho_{i}^{\otimes 1}}; \mathbf{l}m \rangle \langle \dots; \mathbf{n}_{\nu} \mathbf{l}_{\nu} \mathbf{k}_{\nu} | \overline{\rho_{i}^{\otimes \nu}}; \mathbf{s}; \mathbf{k}q \rangle$$

• Can be used to compute efficiently *invariant* features $|
ho_i^{\otimes
u};0;00
angle$

Nigam, Pozdnyakov, MC, JCP (2020)

Machine learning you can trust

NICE features for ML

- Problem: number of features grows exponentially with u
- Solution: an N-body iterative contraction of equivariants (NICE) framework
 - After each body order increase, the most relevant features are selected and used for the next iteration
 - Systematic convergence with ν and contraction truncation



Nigam, Pozdnyakov, MC, JCP (2020)

NICE features for ML

- Problem: number of features grows exponentially with u
- Solution: an N-body iterative contraction of equivariants (NICE) framework
 - After each body order increase, the most relevant features are selected and used for the next iteration
 - Systematic convergence with ν and contraction truncation



Nigam, Pozdnyakov, MC, JCP (2020)

- How to learn with multiple species? Decorate atomic Gaussian with elemental kets $|H\rangle, |O\rangle, \ldots$
- Expand each ket in a finite basis, $|lpha
 angle = \sum_J u_{lpha J} |J
 angle$. Optimize coefficients
- Dramatic reduction of the descriptor space, more effective learning . . .
- ... and as by-product get a data-driven version of the periodic table!



- $\bullet\,$ How to learn with multiple species? Decorate atomic Gaussian with elemental kets $|H\rangle, |O\rangle, \ldots$
- Expand each ket in a finite basis, $|lpha
 angle = \sum_J u_{lpha J} |J
 angle$. Optimize coefficients
- Dramatic reduction of the descriptor space, more effective learning . . .
 . . . and as by-product get a data-driven version of the periodic table!

$$\begin{aligned} |\mathbf{H}\rangle &= 0.5 |\bigstar\rangle + 0.1 |\bigstar\rangle + 0.2 |\bigstar\rangle \\ |\mathbf{C}\rangle &= 0.2 |\bigstar\rangle + 0.8 |\bigstar\rangle + 0.3 |\bigstar\rangle \\ |\mathbf{O}\rangle &= 0.1 |\bigstar\rangle + 0.1 |\bigstar\rangle + 0.6 |\bigstar\rangle \end{aligned}$$

Empedocles et al. (ca 360BC). Metaphor courtesy of Albert Bartók

- How to learn with multiple species? Decorate atomic Gaussian with elemental kets $|H\rangle, |O\rangle, \ldots$
- Expand each ket in a finite basis, $|lpha
 angle = \sum_J u_{lpha J} |J
 angle$. Optimize coefficients
- Dramatic reduction of the descriptor space, more effective learning . . .
- ... and as by-product get a data-driven version of the periodic table!



Elpasolite dataset. Reference curve (red) from Faber et al. JCP (2018)

- How to learn with multiple species? Decorate atomic Gaussian with elemental kets $|H\rangle, |O\rangle, \ldots$
- Expand each ket in a finite basis, $|lpha
 angle = \sum_J u_{lpha J} |J
 angle$. Optimize coefficients
- Dramatic reduction of the descriptor space, more effective learning . . .
- . . . and as by-product get a data-driven version of the periodic table!



Willatt, Musil, Ceriotti, PCCP (2018)

Recognizing active protein ligands

- A SOAP-REMatch-based KSVM classifies active and inactive ligands with 99% accuracy; non-additive model is crucial!
- Sensitivity analysis help identify the active "warhead" and could guide drug design and optimization



Bartok, De, Poelking, Kermode, Bernstein, Csanyi, MC, Science Advances (2017) [data: DUD-E, Shoichet]

Structure-property landscapes

• Clustering/sketch-maps based on REMatch-SOAP correlate well with qualitative classification of packing motifs, and with properties (ex.: azapentacene structure-energy-property landscape maps)



Musil, De, Yang, Campbell, Day, MC, Chemical Science (2018);http://interactive.sketchmap.org

Structure-property landscapes

• Clustering/sketch-maps based on REMatch-SOAP correlate well with qualitative classification of packing motifs, and with properties (ex.: azapentacene structure-energy-property landscape maps)



Musil, De, Yang, Campbell, Day, MC, Chemical Science (2018);http://interactive.sketchmap.org

Structure-property landscapes

 Clustering/sketch-maps based on REMatch-SOAP correlate well with qualitative classification of packing motifs, and with properties (ex.: azapentacene structure-energy-property landscape maps)



Musil, De, Yang, Campbell, Day, MC, Chemical Science (2018); http://interactive.sketchmap.org

Principal Covariates Regression

 Very simple idea to combine PCA and latent-space LR to find a dimensionality reduction that preserves variance and predicts well

$$\ell = \alpha \|\mathbf{X} - \mathbf{X}\mathbf{P}_{XT}\mathbf{P}_{TX}\|^2 + (1 - \alpha) \|\mathbf{Y} - \mathbf{X}\mathbf{P}_{XT}\mathbf{P}_{TY}\|^2$$

• Solution can be found working in sample space (looking for the eigenvectors of a modified Gram matrix)

$$\tilde{\mathbf{K}} = \alpha \mathbf{X} \mathbf{X}^{\mathsf{T}} + (\mathbf{1} - \alpha) \mathbf{X} \mathbf{P}_{\mathbf{X}\mathbf{Y}} \mathbf{P}_{\mathbf{X}\mathbf{Y}}^{\mathsf{T}} \mathbf{X}^{\mathsf{T}}$$

• ... or in feature space by diagonalizing a modified covariance

$$\tilde{\mathbf{C}} = \alpha \mathbf{X}^{\mathsf{T}} \mathbf{X} + (1 - \alpha) \left(\mathbf{X}^{\mathsf{T}} \mathbf{X} \right)^{-1/2} \mathbf{X}^{\mathsf{T}} \mathbf{Y} \mathbf{Y}^{\mathsf{T}} \mathbf{X} \left(\mathbf{X}^{\mathsf{T}} \mathbf{X} \right)^{-1/2}$$



S. de Jong and HAL Kiers, Scandinavian Symposium on Chemometrics (1992)

Kernel PCovR

• Kernel versions of PCovR can be obtained with a modified kernel $\tilde{\mathbf{K}} = \alpha \mathbf{K} + (1 - \alpha) \hat{\mathbf{Y}} \hat{\mathbf{Y}}^{T}$, diagonalizing it and finding the projector

$$\mathbf{P}_{\mathbf{K}T} = \left(\alpha \mathbf{I} + (1 - \alpha) \left(\mathbf{K} + \lambda \mathbf{I}\right)^{-1} \mathbf{Y} \hat{\mathbf{Y}}\right) \mathbf{U}_{\tilde{\mathbf{K}}} \Lambda_{\tilde{\mathbf{K}}}^{1/2}$$



Machine learning you can trust

Where unsupervised meets supervised

 Using KPCovR to reveal structure-property relations in databases of materials structures



Machine learning you can trust

A Generalized Convex Hull Construction



Anelli, Engel, Pickard & MC, PRM (2019); Engel, Anelli, MC, Pickard & Needs, Nature Comm. (2018)