

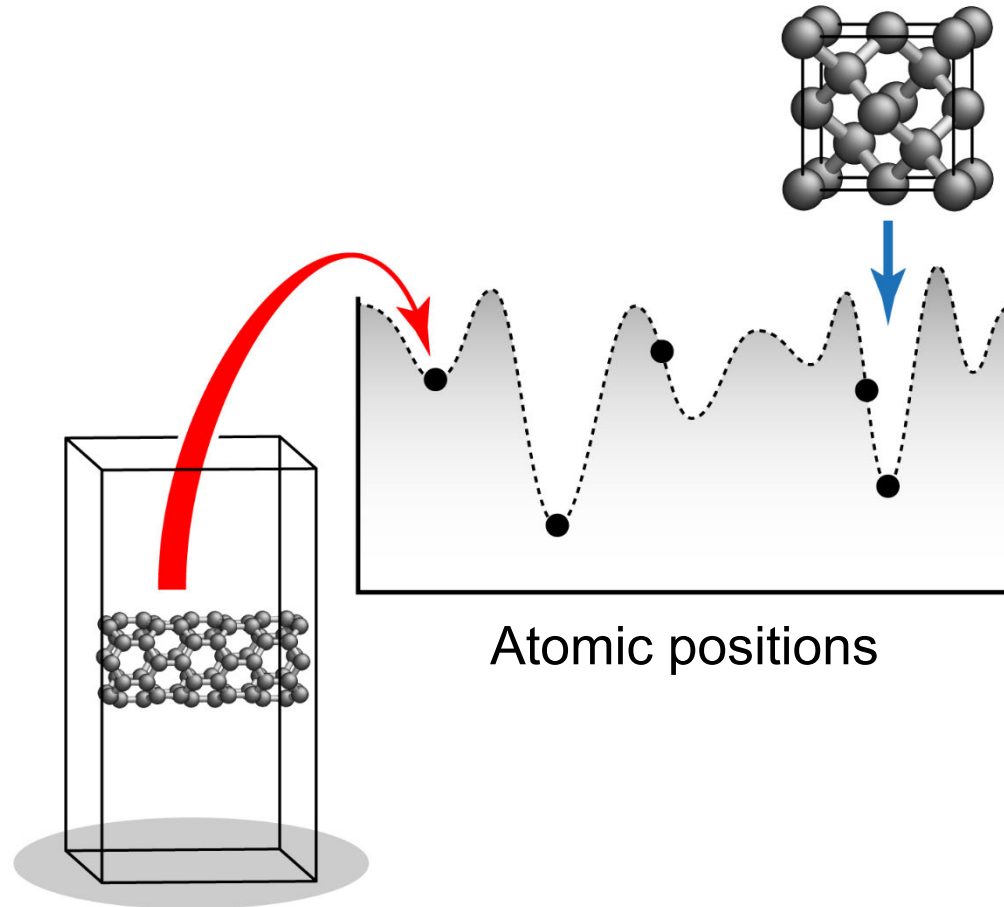


Gaussian process regression for atomistic materials modelling

Volker Deringer with tutorials by Joe Morrow

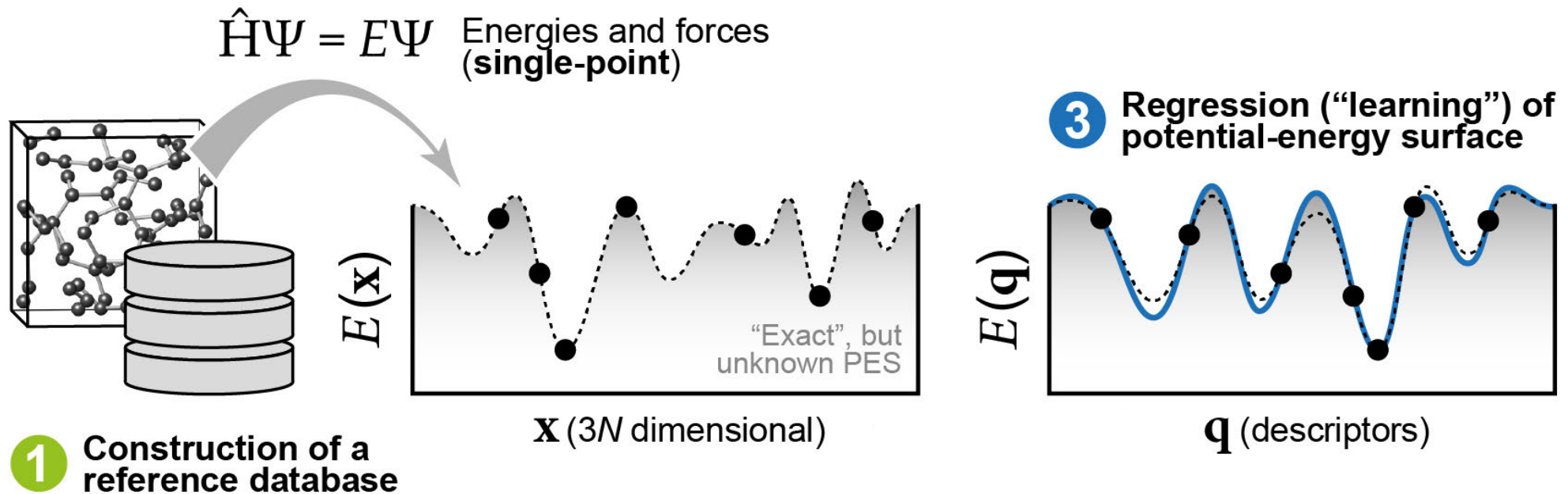
<http://deringer.chem.ox.ac.uk>

Machine-learning-driven materials modelling



Machine-learning-driven materials modelling

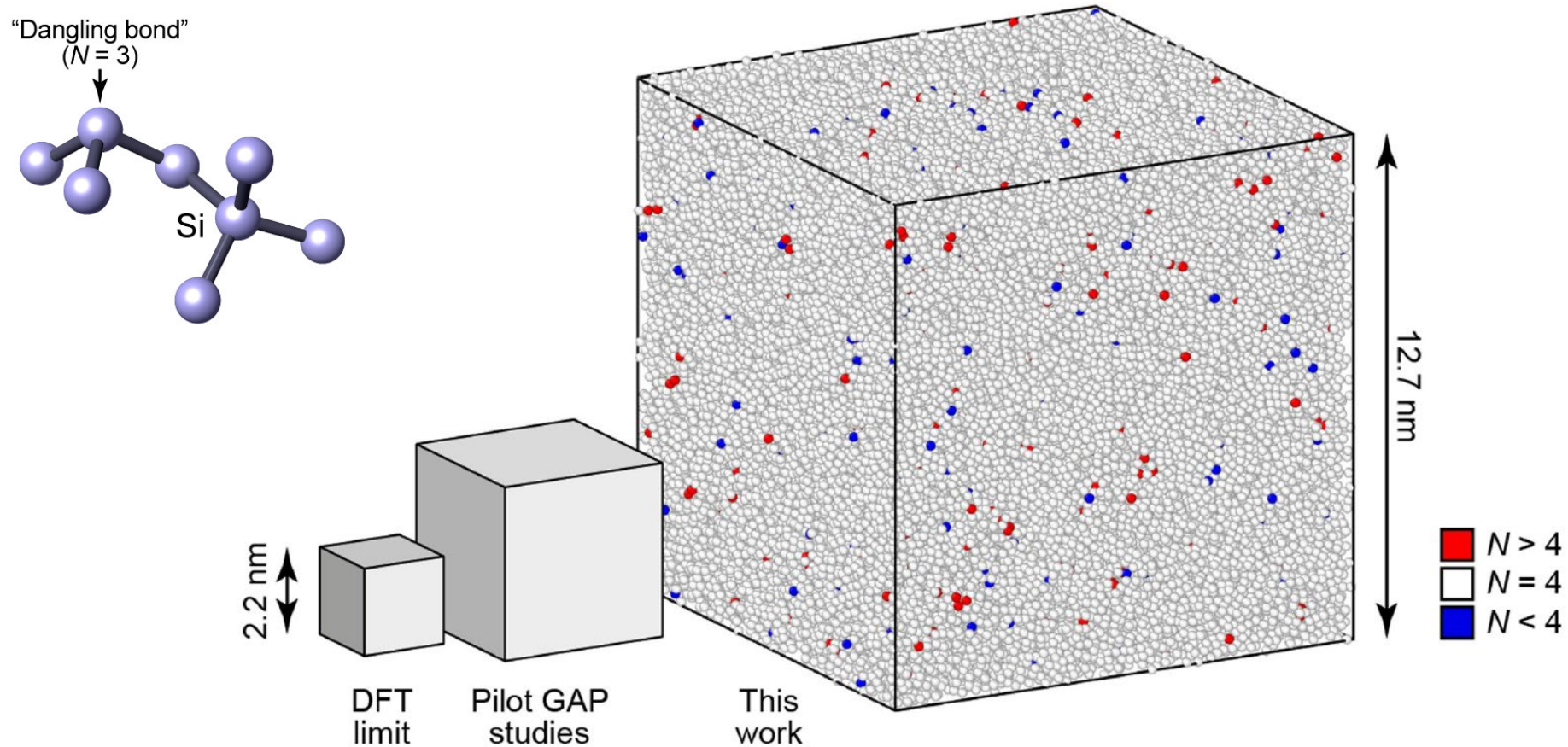
ML potentials approximate the quantum-mechanical potential-energy surface, using three main **ingredients**:



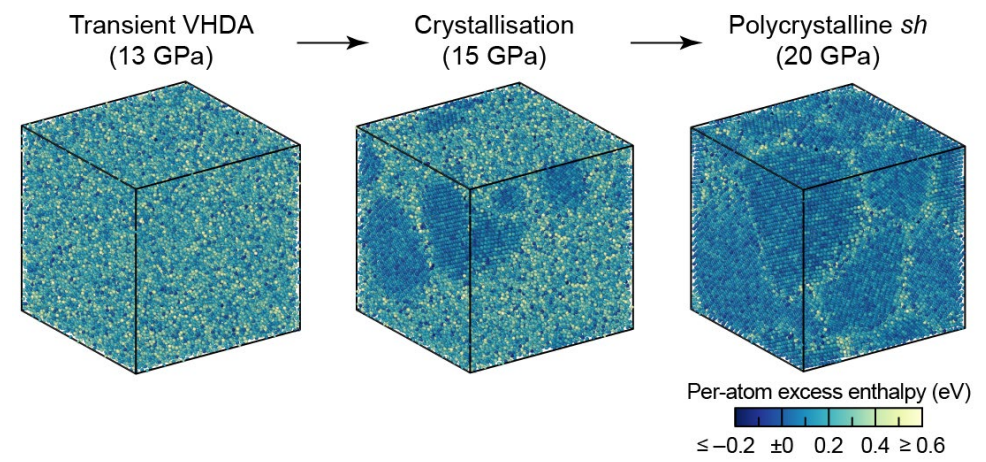
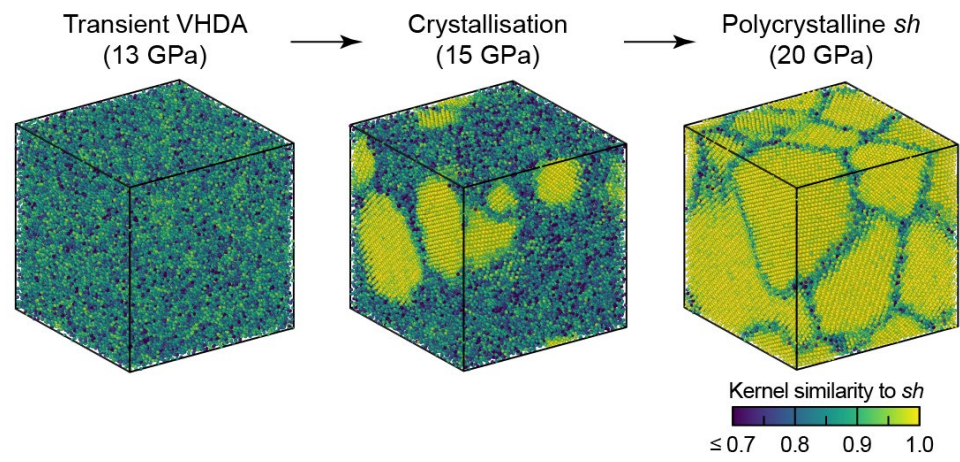
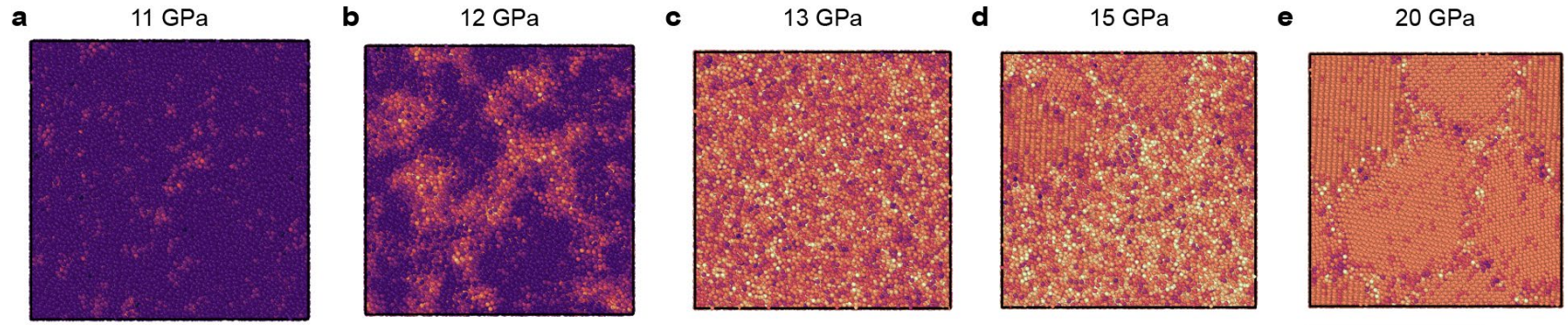
A “textbook case” for random networks: **Amorphous silicon**

Silicon “GAP-18”: Bartók et al., *Phys. Rev. X* **2018**, 8, 041048

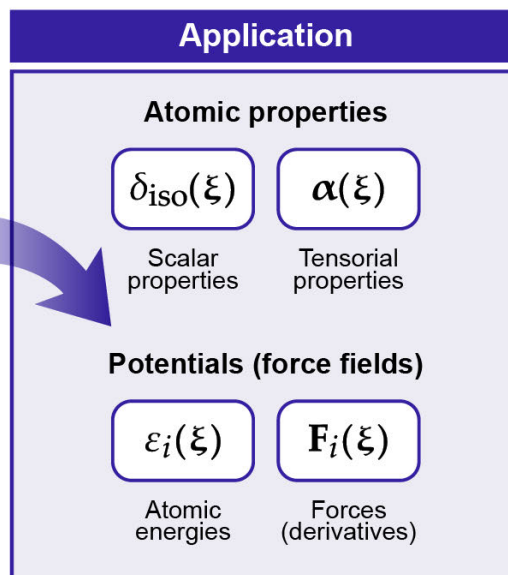
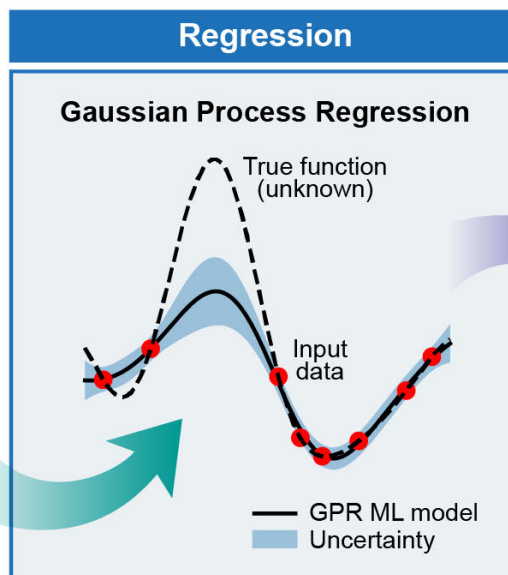
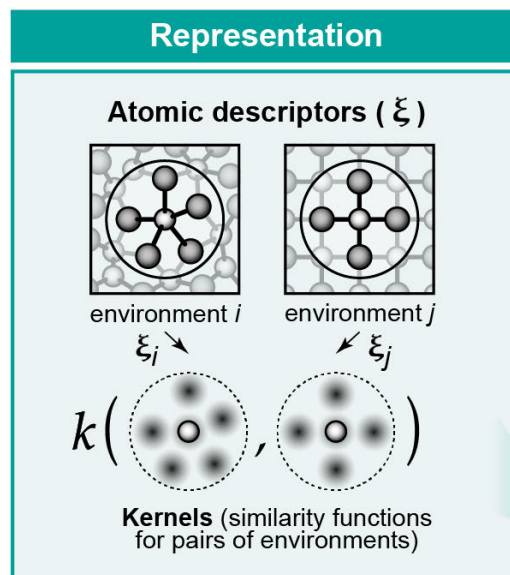
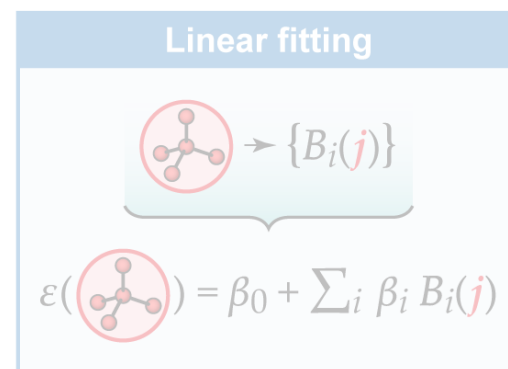
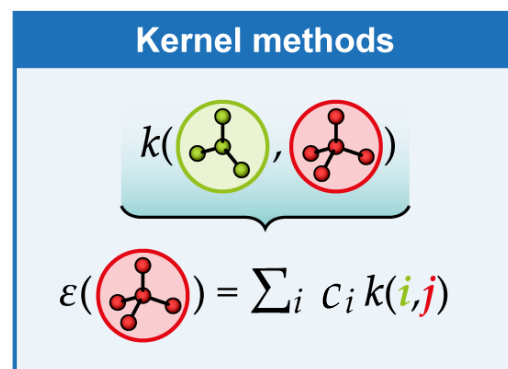
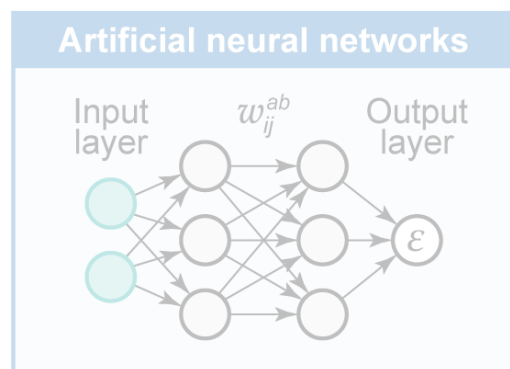
Validation for *a*-Si: VLD et al., *J. Phys. Chem. Lett.* **2018**, 9, 2879



Origins of structural transitions in disordered silicon



Gaussian process regression (GPR) for materials



CHEMICAL REVIEWS
pubs.acs.org/CR

Gaussian Process Regression for Materials and Molecules
Volker L. Deringer,* Albert P. Bartók,* Noam Bernstein, David M. Wilkins, Michele Ceriotti, and Gábor Csányi*

Cite This: Chem. Rev. 2021, 121, 10073–10141

ACCESS | Metrics & More | Article Recommendations

ABSTRACT: We provide an introduction to Gaussian process regression (GPR) machine-learning methods in computational materials science and chemistry. The focus of the present review is on the regression of atomistic properties: in particular, on the construction of interatomic potentials, or force fields, in the Gaussian Approximation Potential (GAP) framework; beyond this, we also discuss the fitting of arbitrary scalar, vectorial, and tensorial quantities. Methodological aspects of reference data generation, representation, and regression, as well as the question of how a data-driven model may be validated, are reviewed and critically discussed. A survey of applications to a variety of research questions in chemistry and materials science illustrates the rapid growth in the field. A vision is outlined for the development of the methodology in the years to come.

CONTENTS

1. Introduction	10074	4.5.1. Cutoff Radius	10102
2. Gaussian Process Regression	10074	4.5.2. Kernel Regularity	10102
2.1. Weight-Space View of GPR	10076	4.6. Regularization in GAPs	10103
2.2. Function-Space View of GPR	10077	4.6.1. Noise in the Input	10103
2.3. Explicit Construction of the Reproducing Kernel Hilbert Space	10079	4.6.2. Dealing with Inhomogeneous Data	10104
2.4. GPR Based on Linear Functional Observations	10080	4.6.3. Implementation	10104
2.5. Regularization	10083	5. Validation and Accuracy	10105
2.6. Hyperparameters	10083	5.1. Physical Behavior versus Numerical Errors	10105
3. Learning Atomistic Properties	10084	5.2. Predicted Errors in GPR	10107
3.1. Representing Atomic Structures	10087	5.3. Committee Models and Uncertainty Propagation	10108
3.2. Symmetry-Adapted Representation	10087	5.4. GPR Models for Isolated Molecules	10110
3.3. Ψ_0 Potential Energy: A Hands-On Example	10087	6. Applications (I): Force Fields	10111
3.4. Symmetry-Adapted GPR	10089	6.1. Transition Metals	10112
4. Gaussian Approximation Potential (GAP) Framework	10090	6.2. Complex Allotropy and Crystal-Structure Prediction	10113
4.1. Reference Data	10091	6.3. Structure of Amorphous Materials	10115
4.1.1. Hand-Built Databases	10092	6.3.1. Carbon Nanostructures	10115
4.1.2. Iterative and Active Learning	10092	6.3.2. Amorphous Silicon	10116
4.1.3. GAP-RSS	10093	6.3.3. Ge–Sb–Te Phase-Change Materials	10117
4.1.4. Automatic Training Set Selection	10095	6.4. Surface Chemistry	10118
4.1.5. General-Purpose Databases	10095	6.5. Functional Properties	10119
4.2. Hierarchical Models	10096	6.6. Molecular Materials	10121
4.3. Sparse GPR	10099	7. Applications (II): Beyond Force Fields	10123
4.4. Locality	10100	7.1. NMR Chemical Shieldings	10123
4.5. Practical Choices for Hyperparameters	10102		

Special Issue: Machine Learning at the Atomic Scale

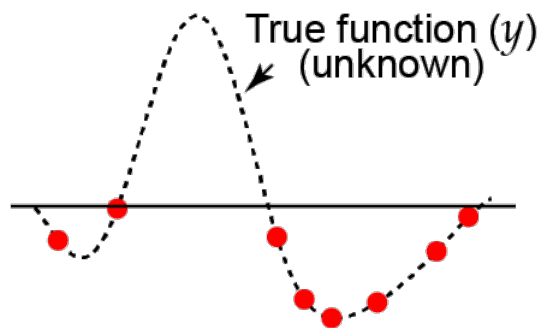
Received: January 8, 2021
Published: August 16, 2021

ACS Publications | 10073

Gaussian process regression (GPR)

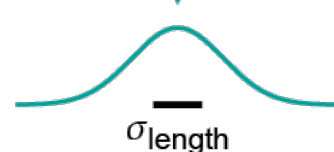
1 Observations

(\mathbf{x}_n, y_n)
at specific
locations (●)



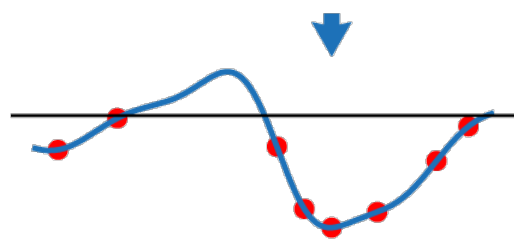
2 Basis function

$$k(\mathbf{x}, \mathbf{x}_m) = \exp\left(-\frac{|\mathbf{x} - \mathbf{x}_m|^2}{2\sigma_{\text{length}}^2}\right)$$



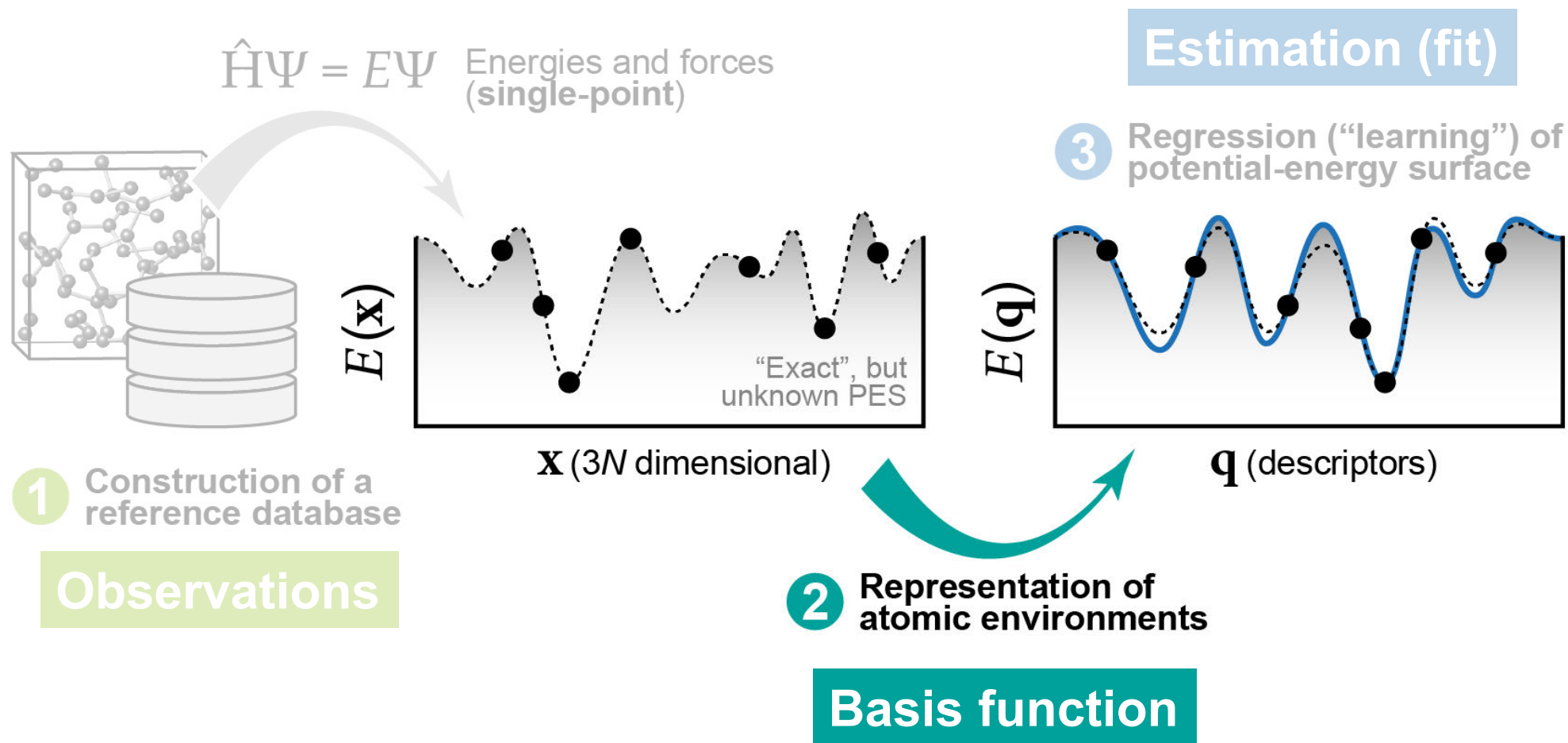
3 Estimation (fit)

$$\tilde{y}(\mathbf{x}) = \sum_{m=1}^M c_m k(\mathbf{x}, \mathbf{x}_m)$$

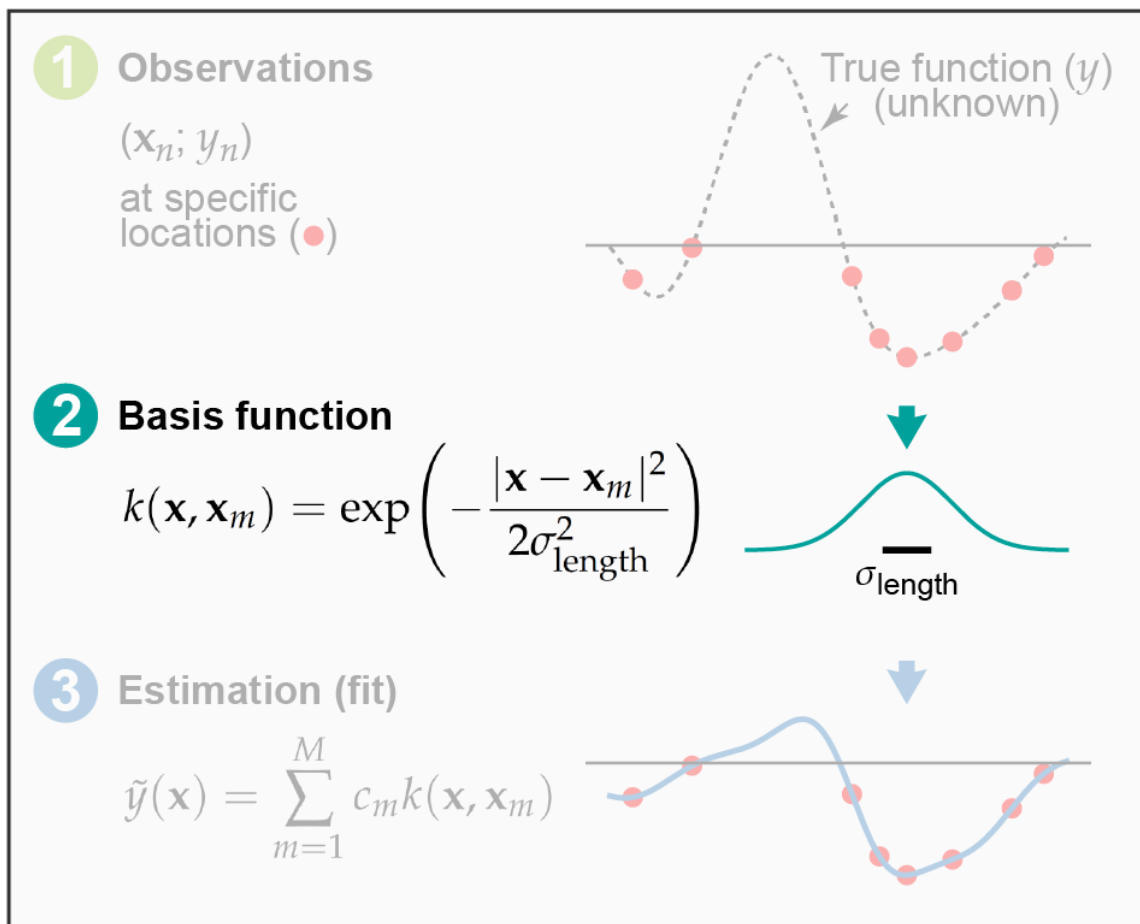


- **Observations**, using the “ground truth”, make up the training data
- In atomistic ML: \mathbf{x} encodes atomic environment, y can be any per-atom property
- **Basis functions** are built using a kernel similarity measure, k
- In atomistic ML: from simple pair distance to more complex forms
- **Estimation** means predicting y at a new, unknown location \mathbf{x}
- To do this, we need k and pre-determined fitting coefficients, c

Gaussian process regression (GPR) for materials

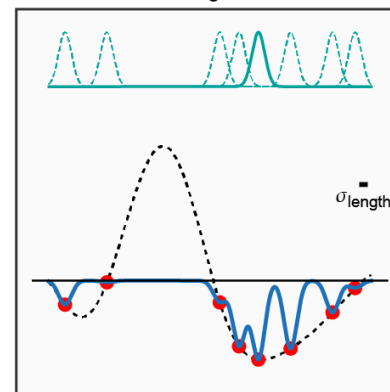


Gaussian process regression (GPR)

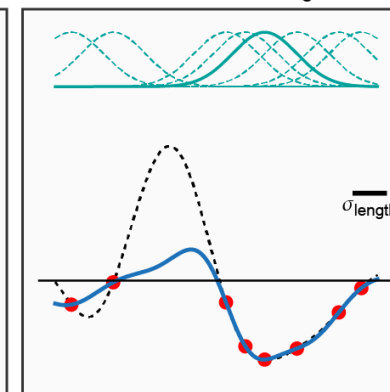


a Learning from function values

Too small σ_{length} (overfitting)

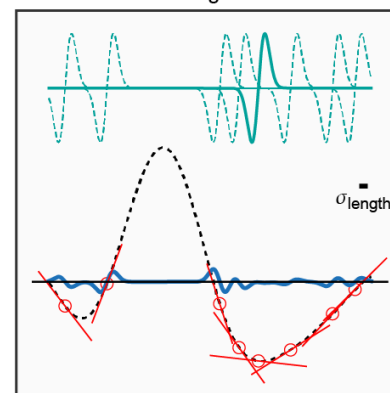


Appropriate σ_{length}

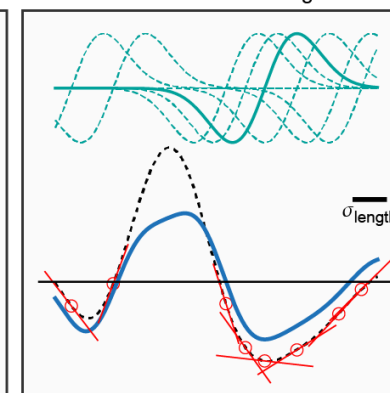


b Learning from derivative values

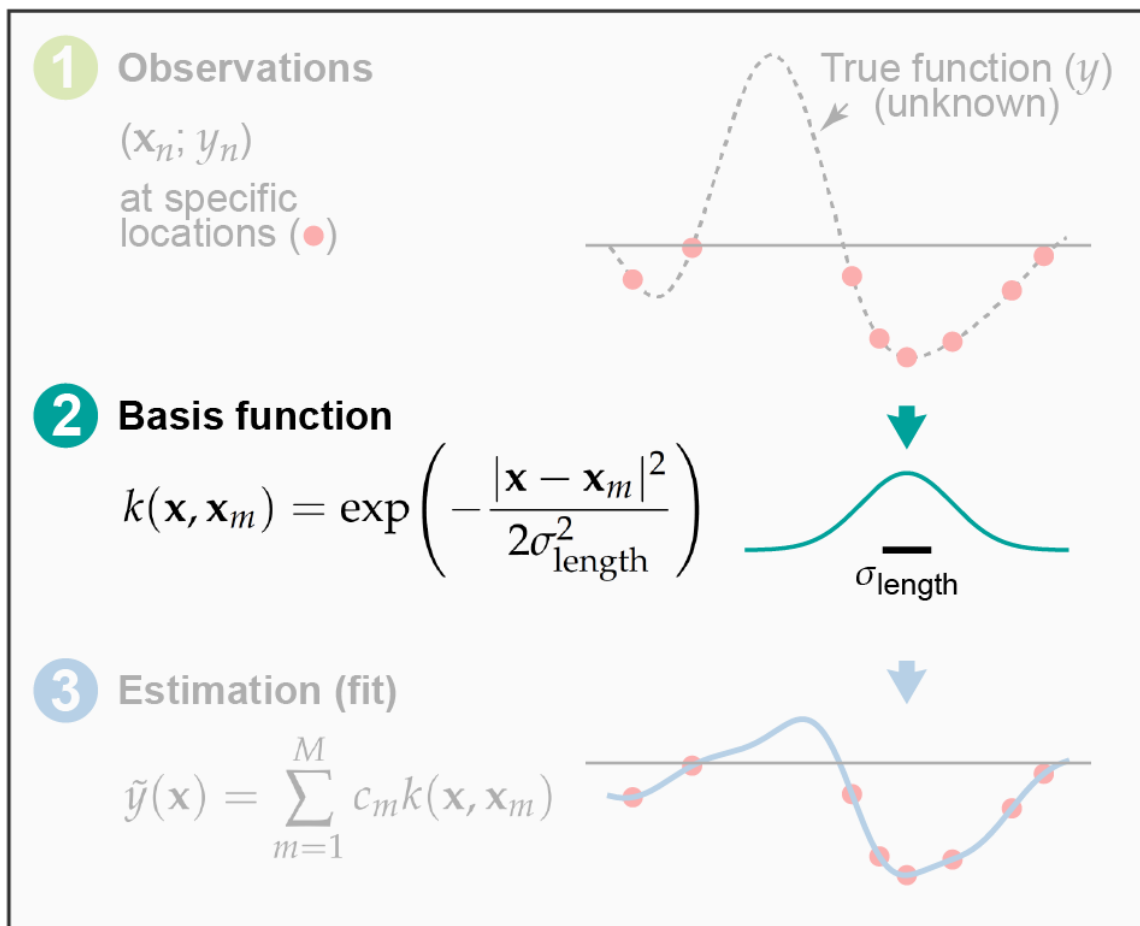
Too small σ_{length} (overfitting)



Appropriate σ_{length}

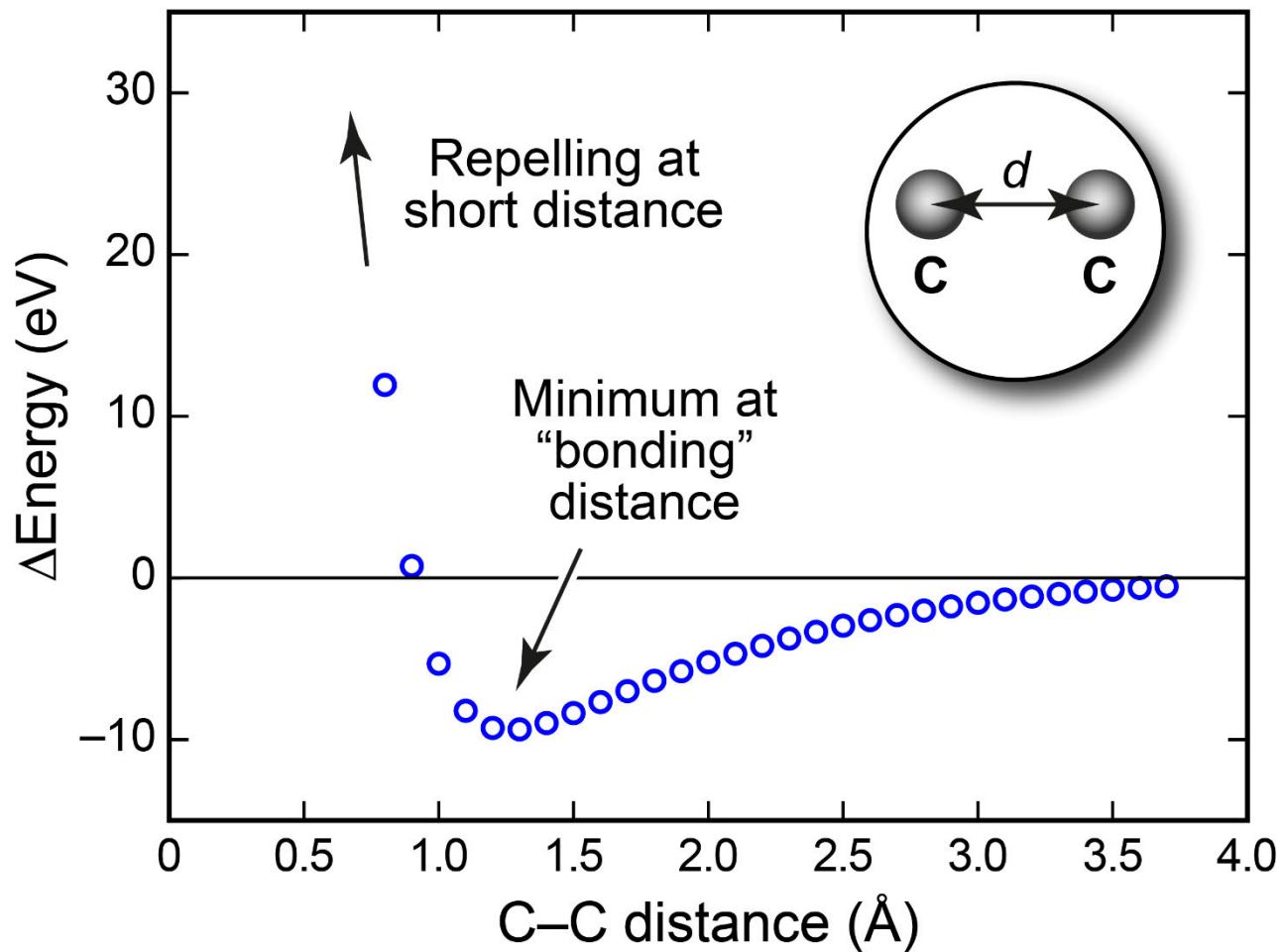


Gaussian process regression (GPR)

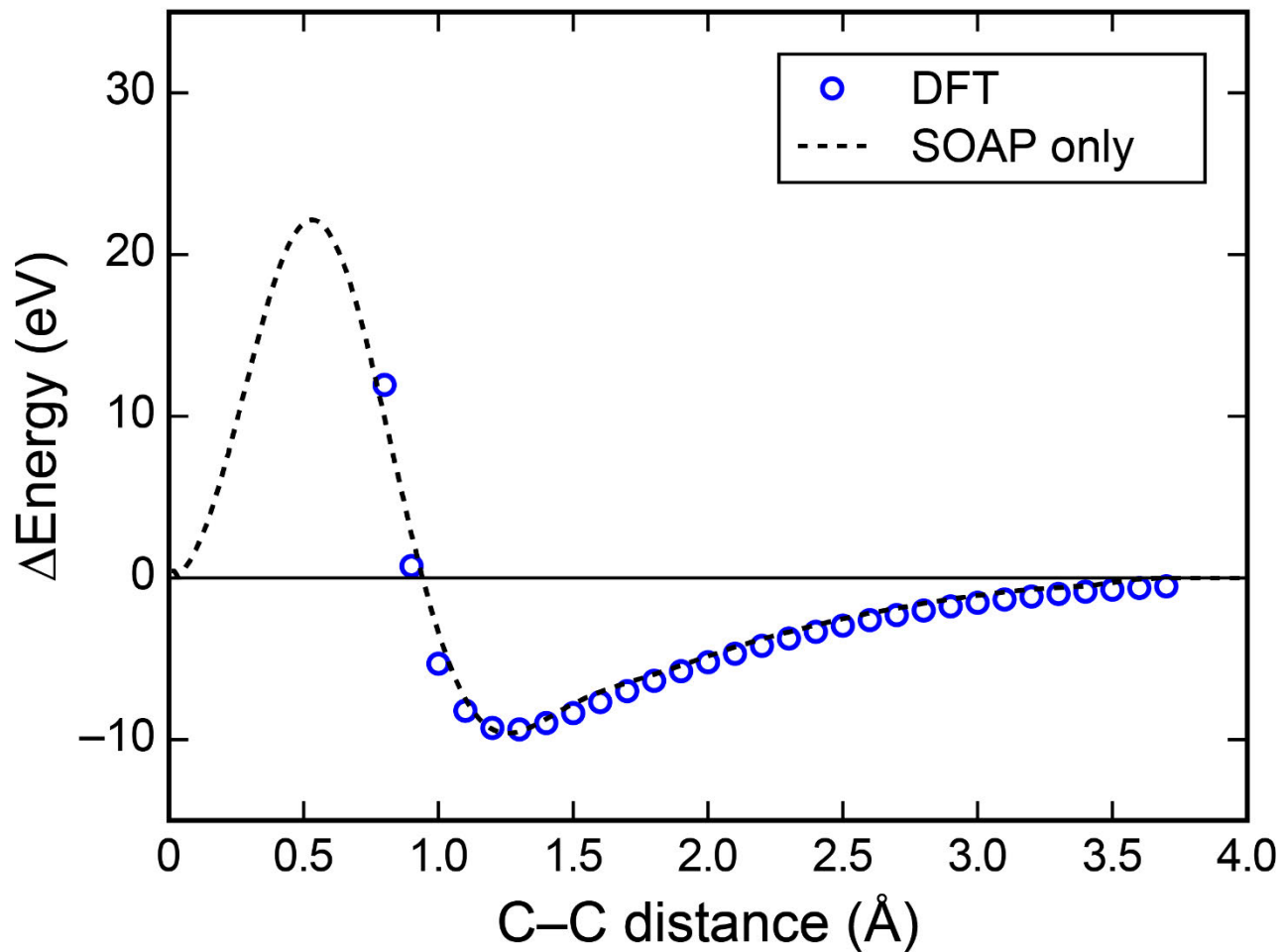


- **Observations**, using the “ground truth”, make up the training data
- In atomistic ML: \mathbf{x} encodes atomic environment, y can be any per-atom property
- **Basis functions** are built using a kernel similarity measure, k
- **In atomistic ML: from simple pair distance to more complex forms**
- **Estimation** means predicting y at a new, unknown location \mathbf{x}
- To do this, we need k and pre-determined fitting coefficients, c

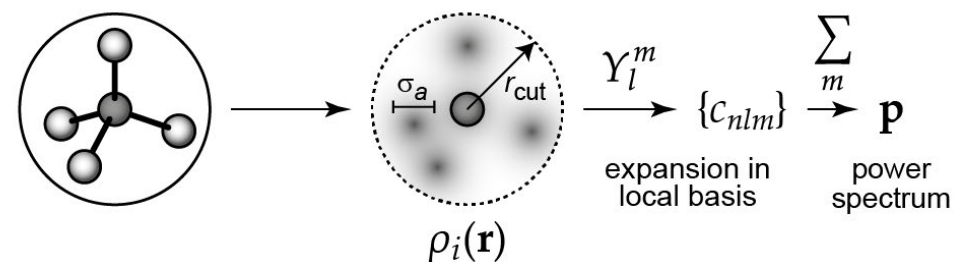
Representations (*descriptors*) in practice



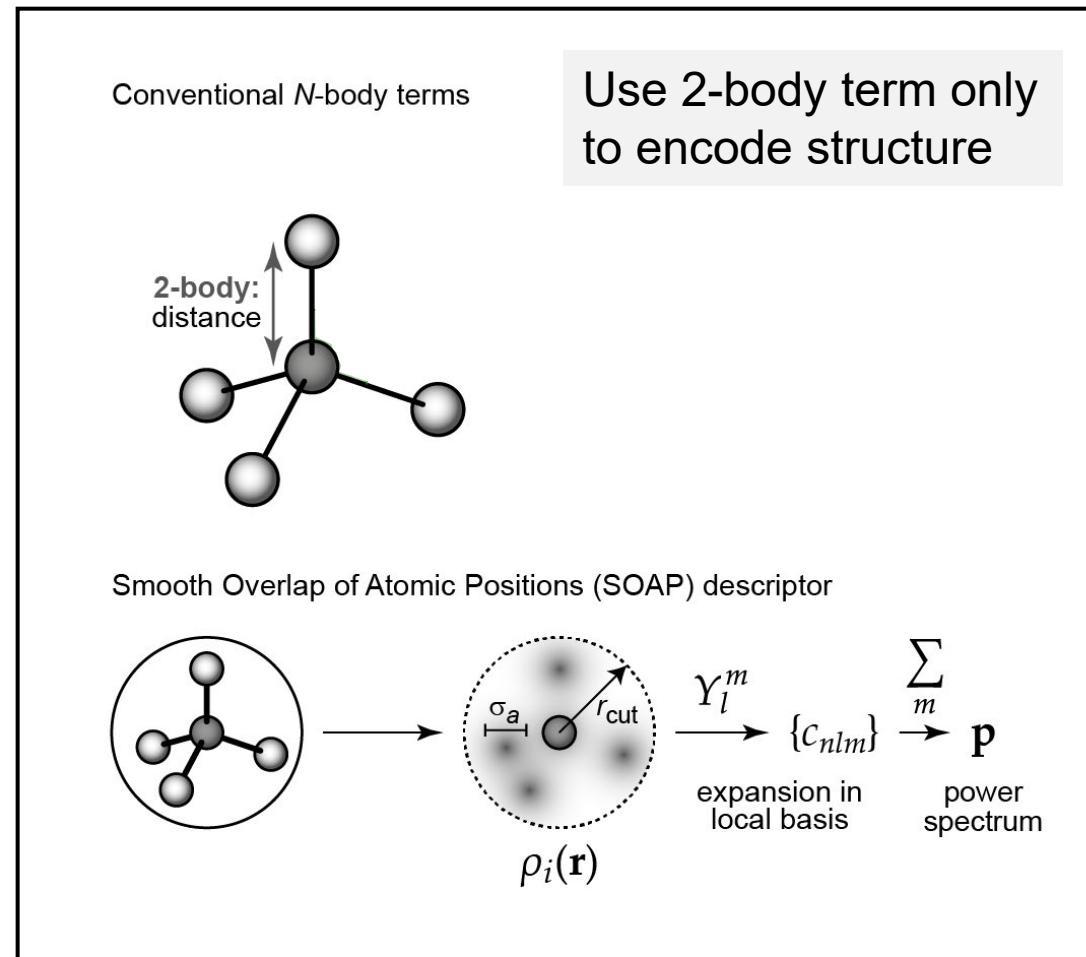
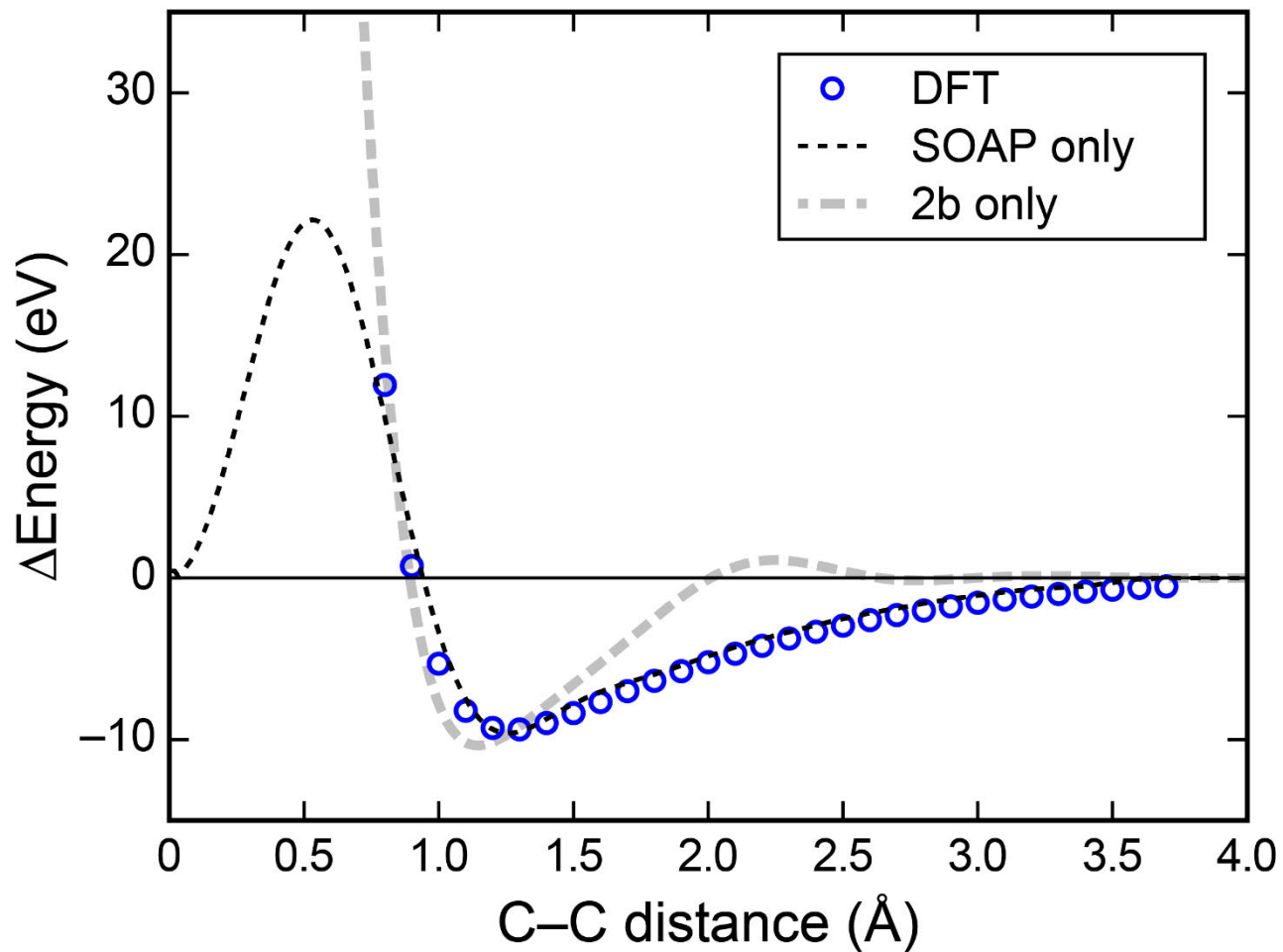
Representations (*descriptors*) in practice



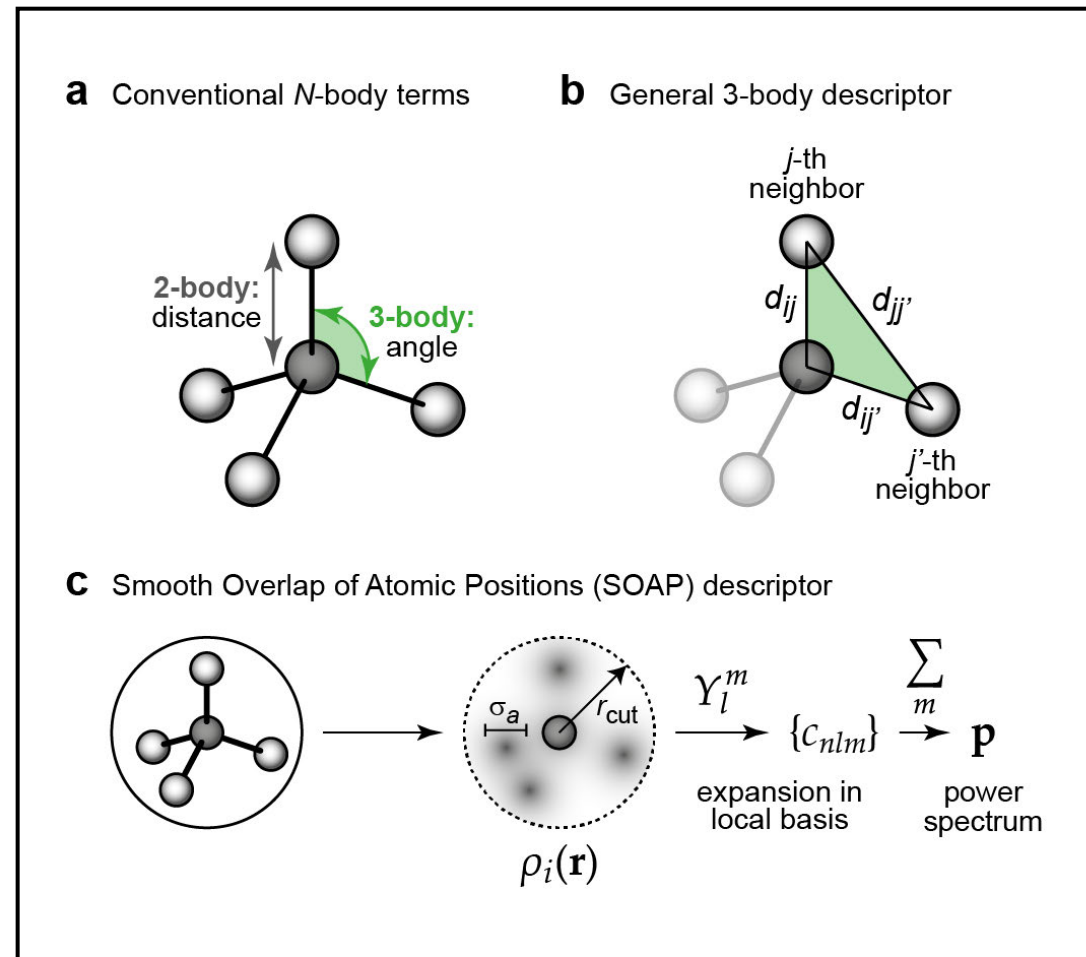
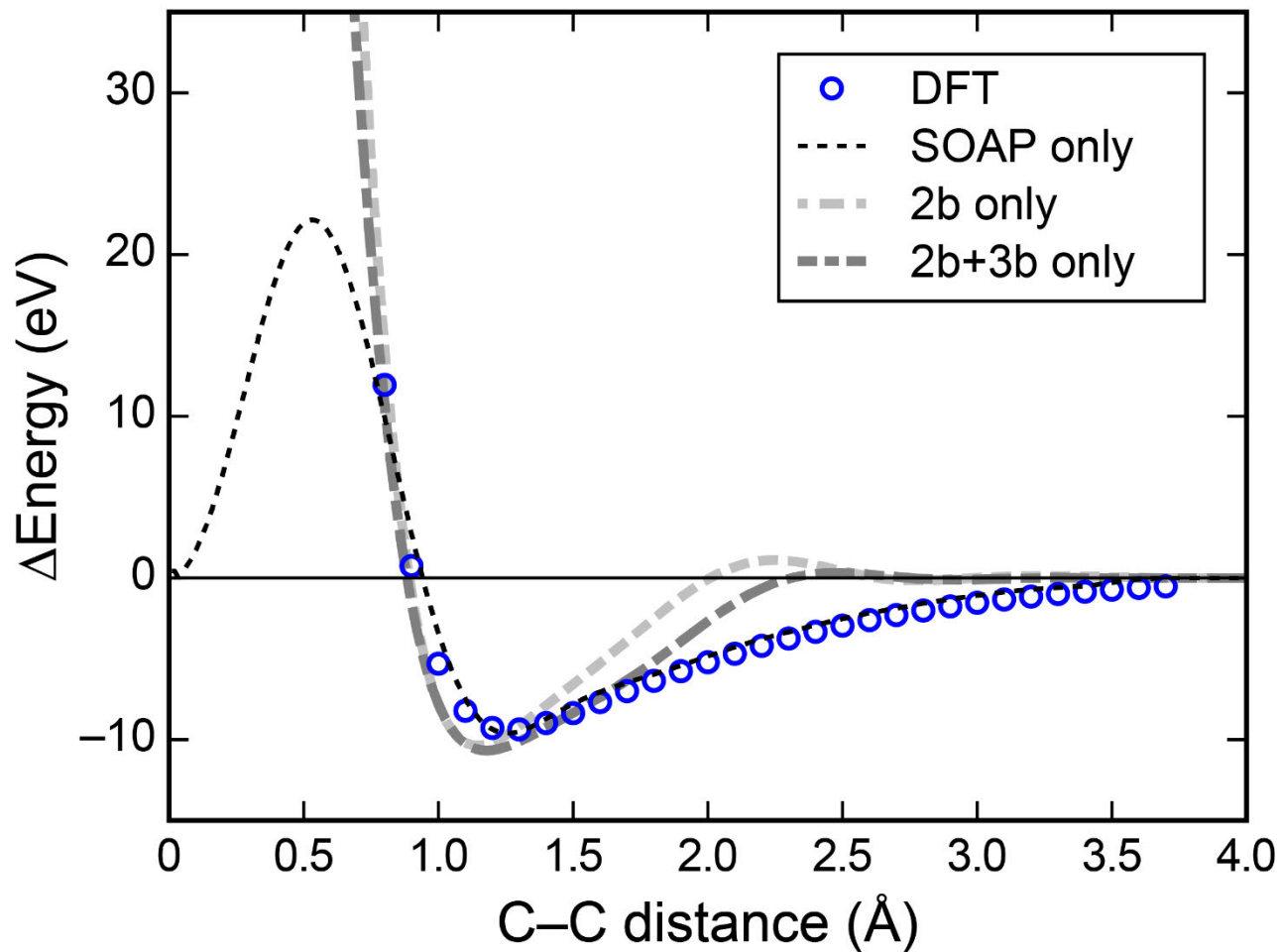
Use Smooth Overlap of Atomic Positions (**SOAP**) only to encode atomic structure



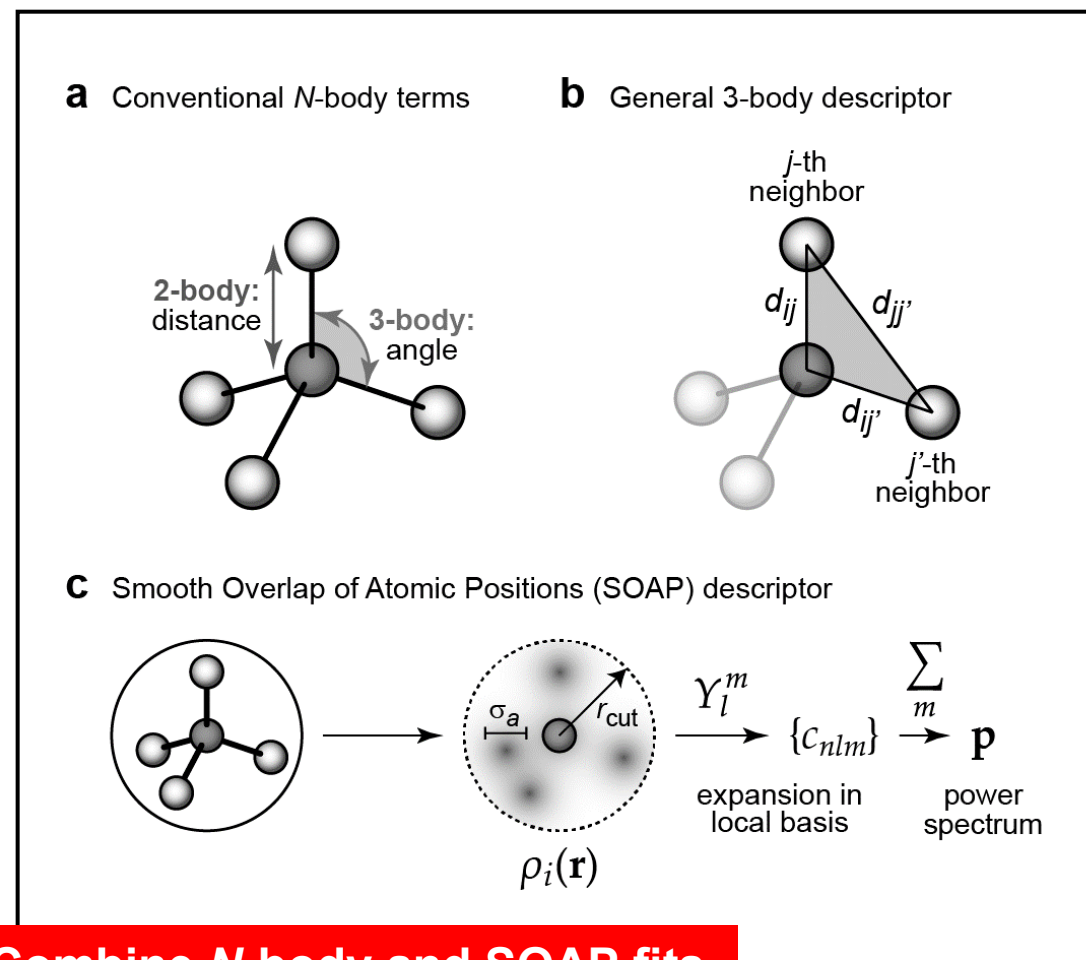
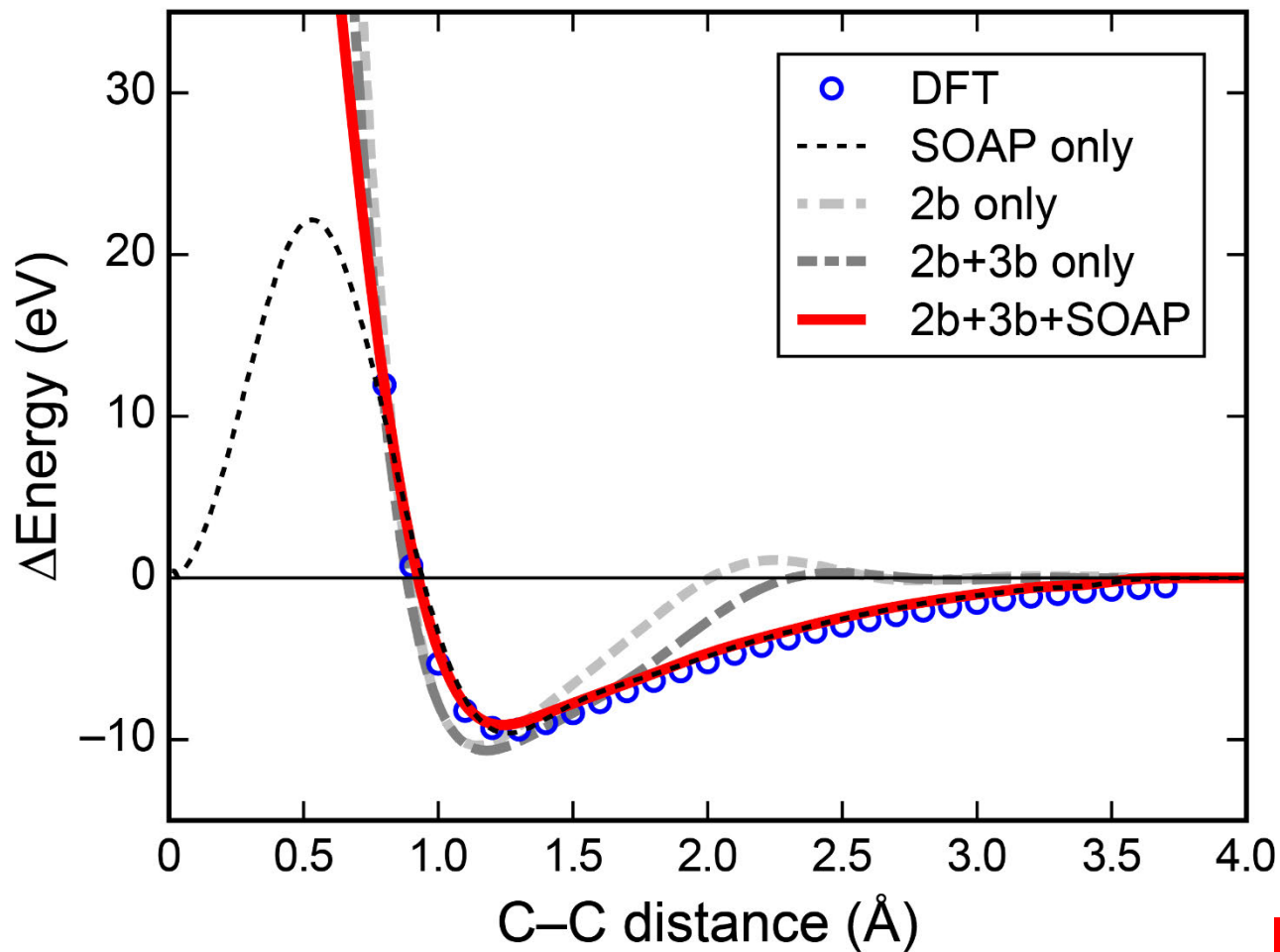
Representations (*descriptors*) in practice



Representations (*descriptors*) in practice



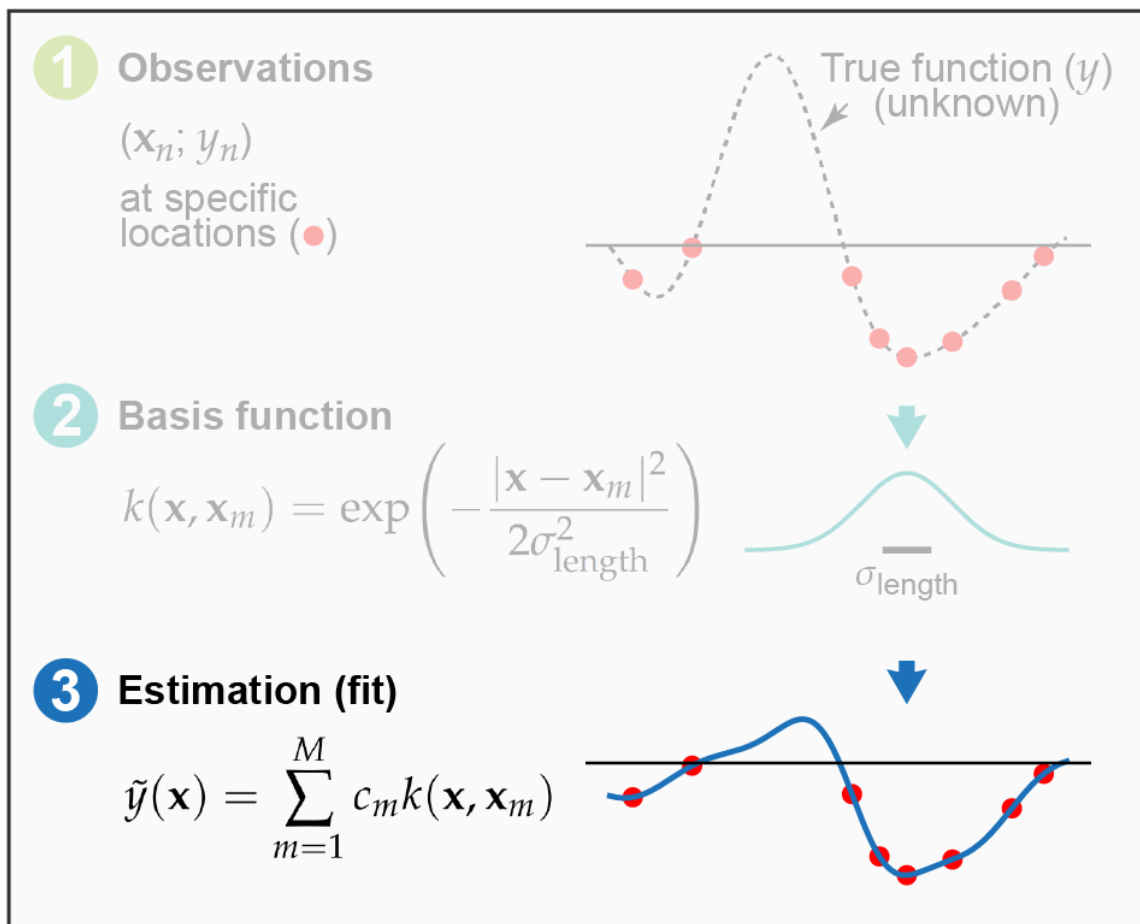
Representations (*descriptors*) in practice



Combine N -body and SOAP fits

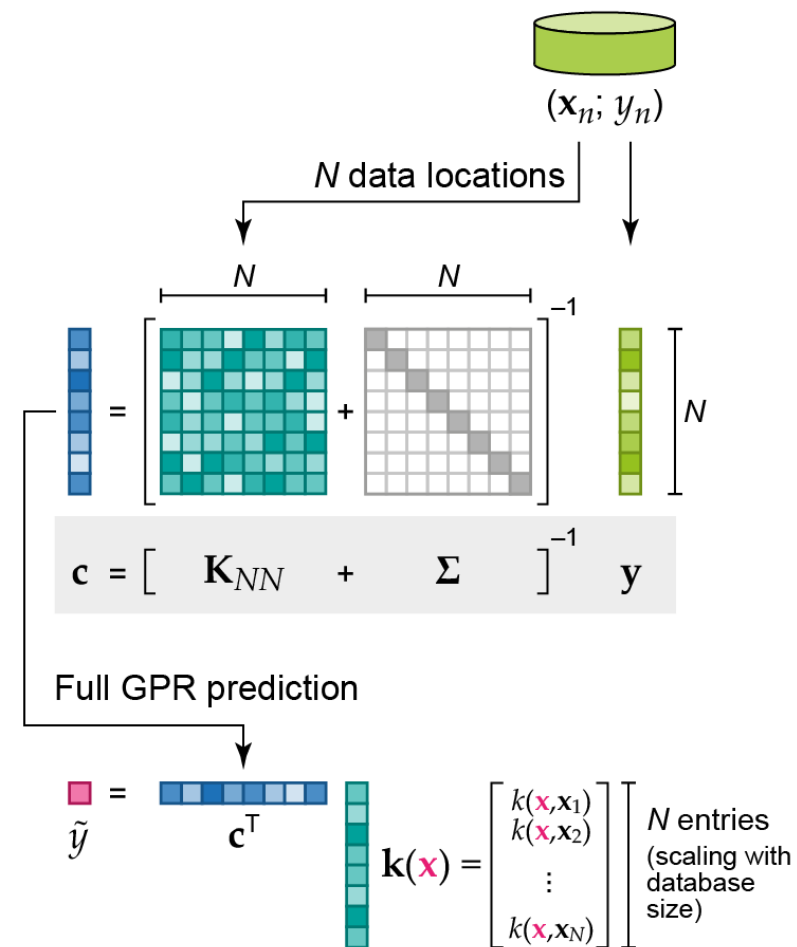
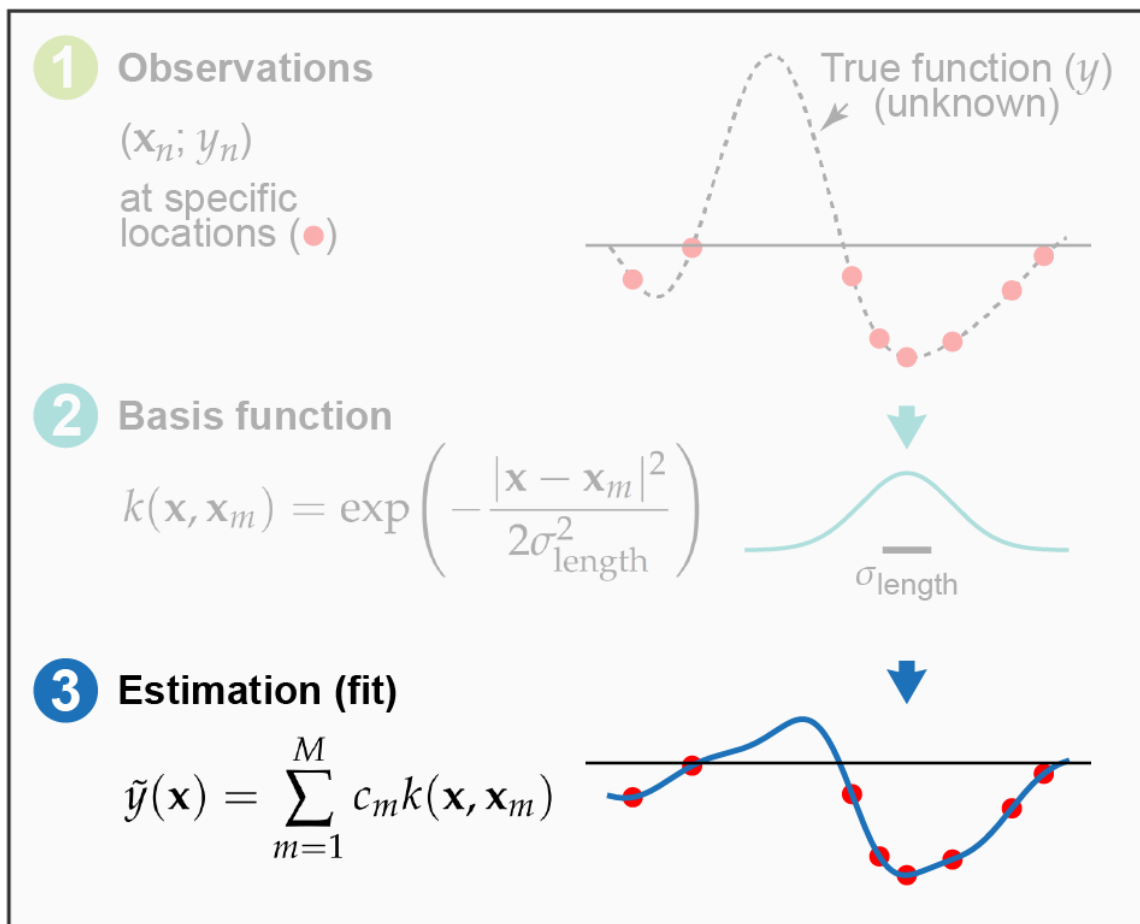
for **robustness** where needed, & **accuracy** where possible

Gaussian process regression (GPR)



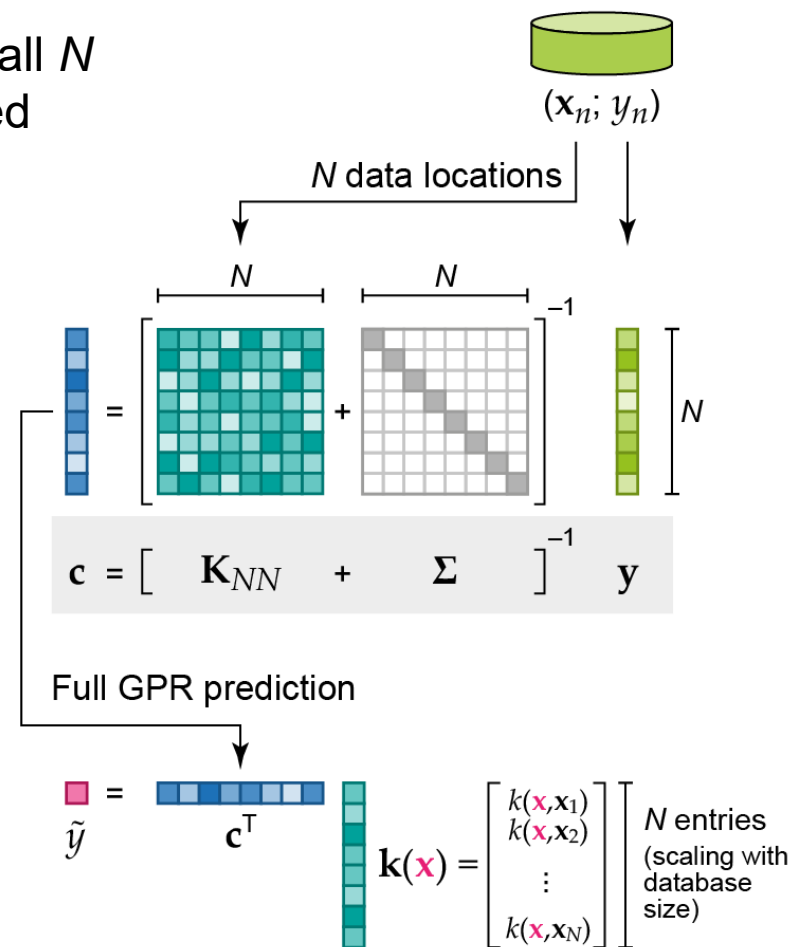
- **Observations**, using the “ground truth”, make up the training data
 - In atomistic ML: \mathbf{x} encodes atomic environment, y can be any per-atom property
 - **Basis functions** are built using a kernel similarity measure, k
 - In atomistic ML: from simple pair distance to more complex forms
 - **Estimation** means predicting y at a new, unknown location \mathbf{x}
- To do this, we need k and pre-determined fitting coefficients, c

Gaussian process regression (GPR)



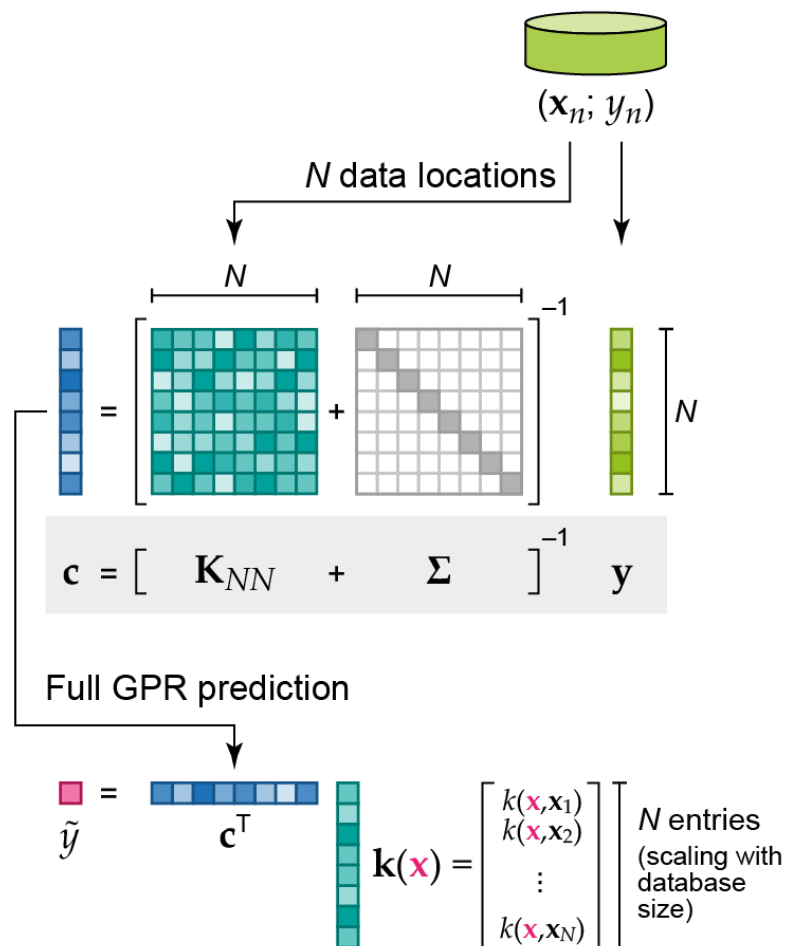
Gaussian process regression (GPR)

- Finding the coefficients, \mathbf{c} , means solving a linear algebra problem
 - Σ is a **regulariser** that encodes the expected noise in the input data
 - (effect of regularisation: see examples at <https://arxiv.org/abs/2211.16443>)
- This is “**full GPR**”: all N data points are used

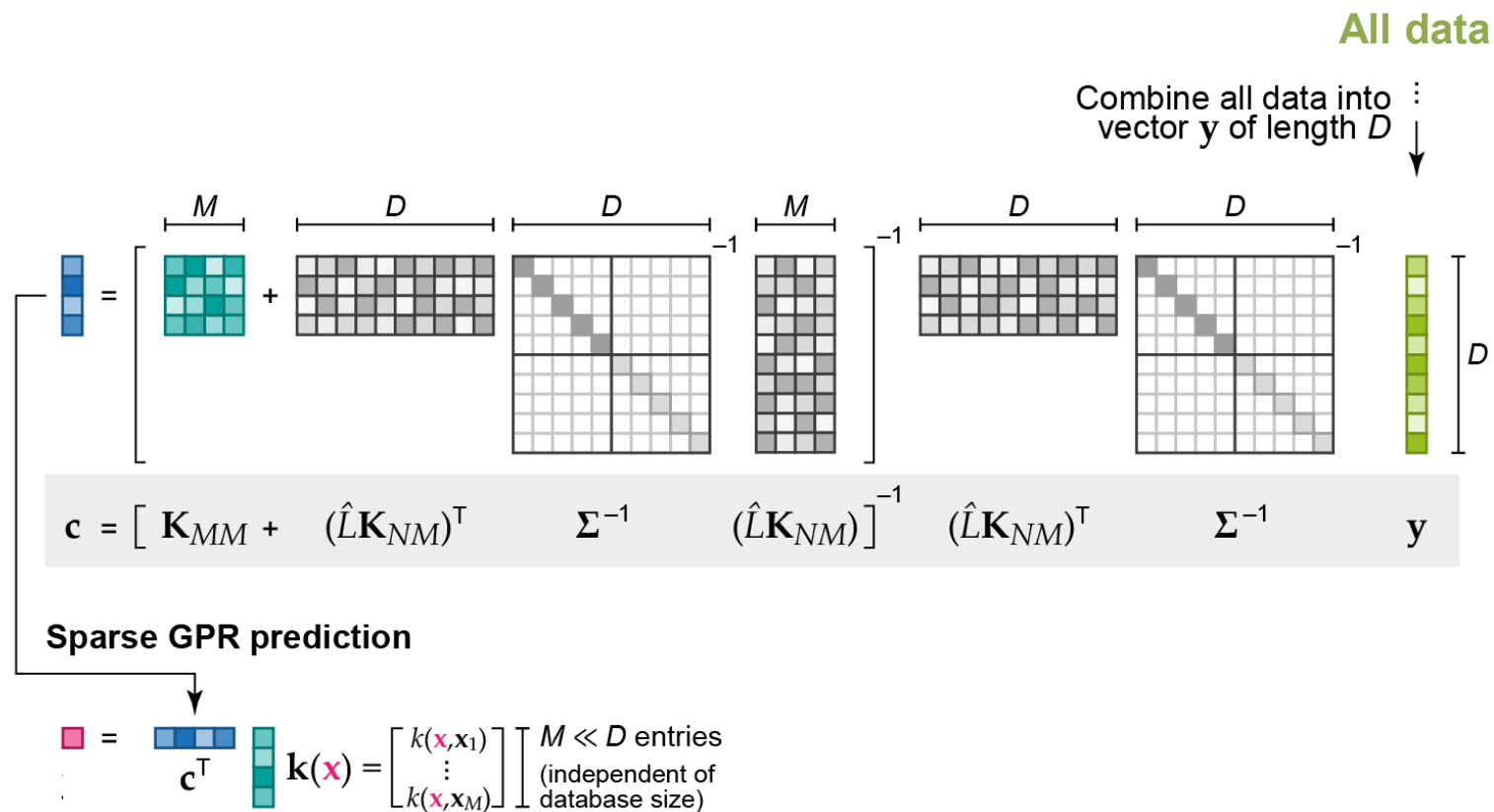


Gaussian process regression (GPR)

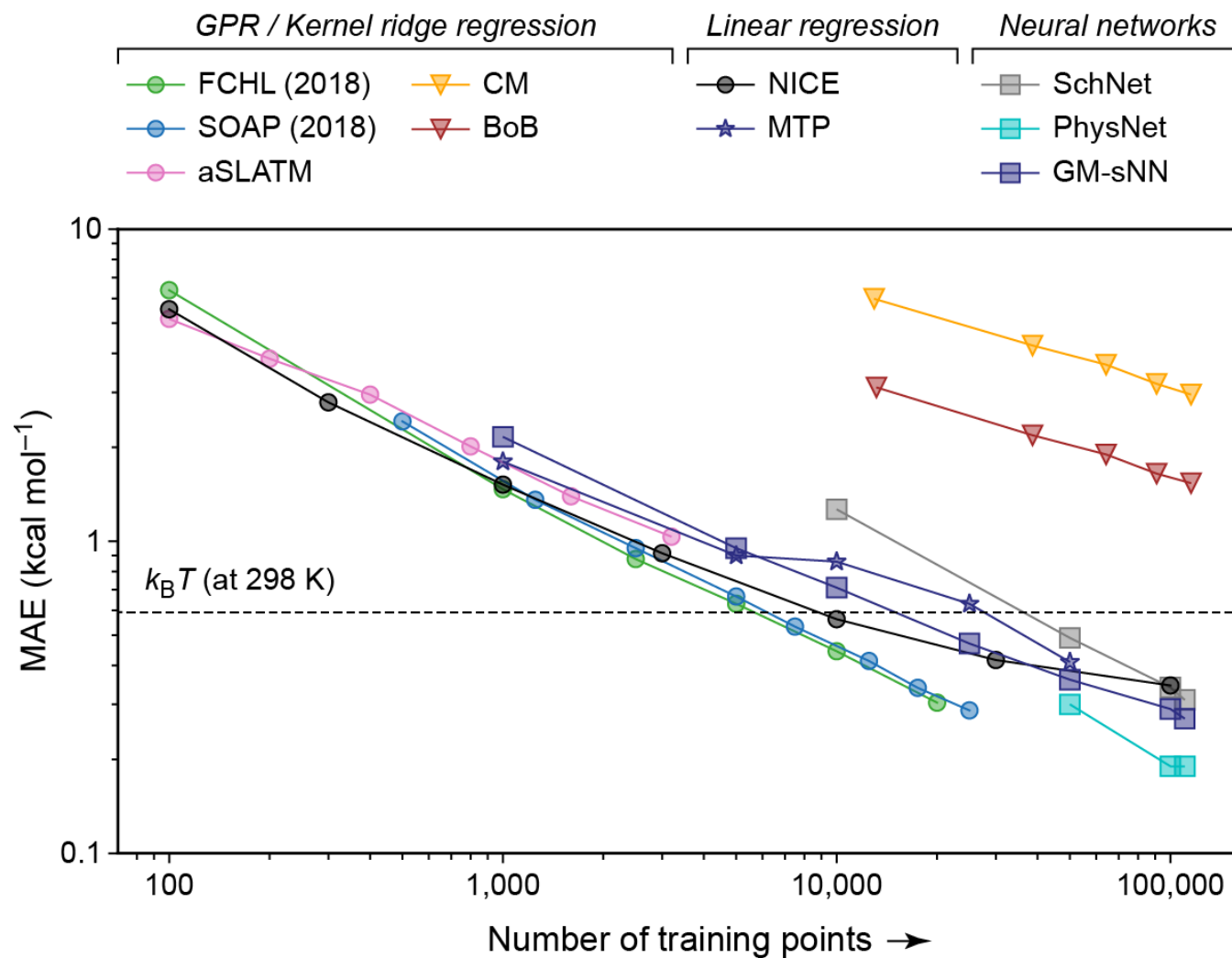
a Full GPR fitting



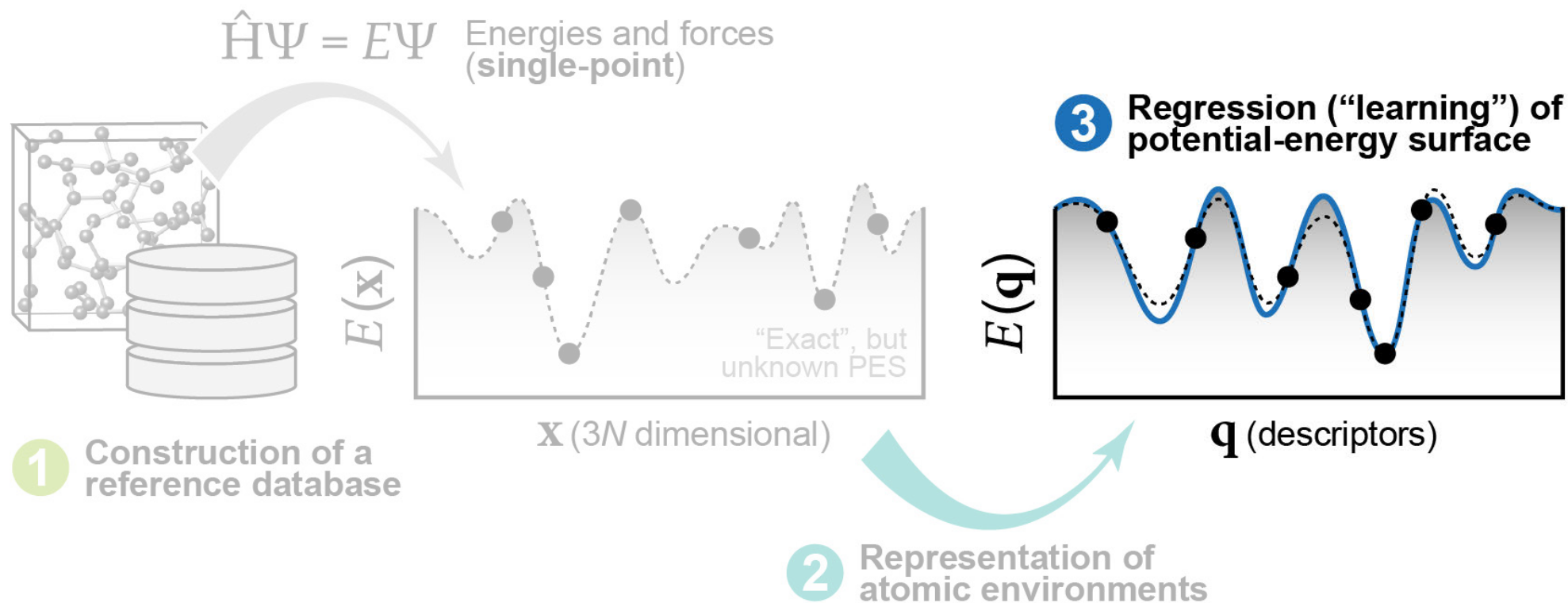
b Sparse GPR fitting: e.g., in **GAP** models



Gaussian process regression (GPR)

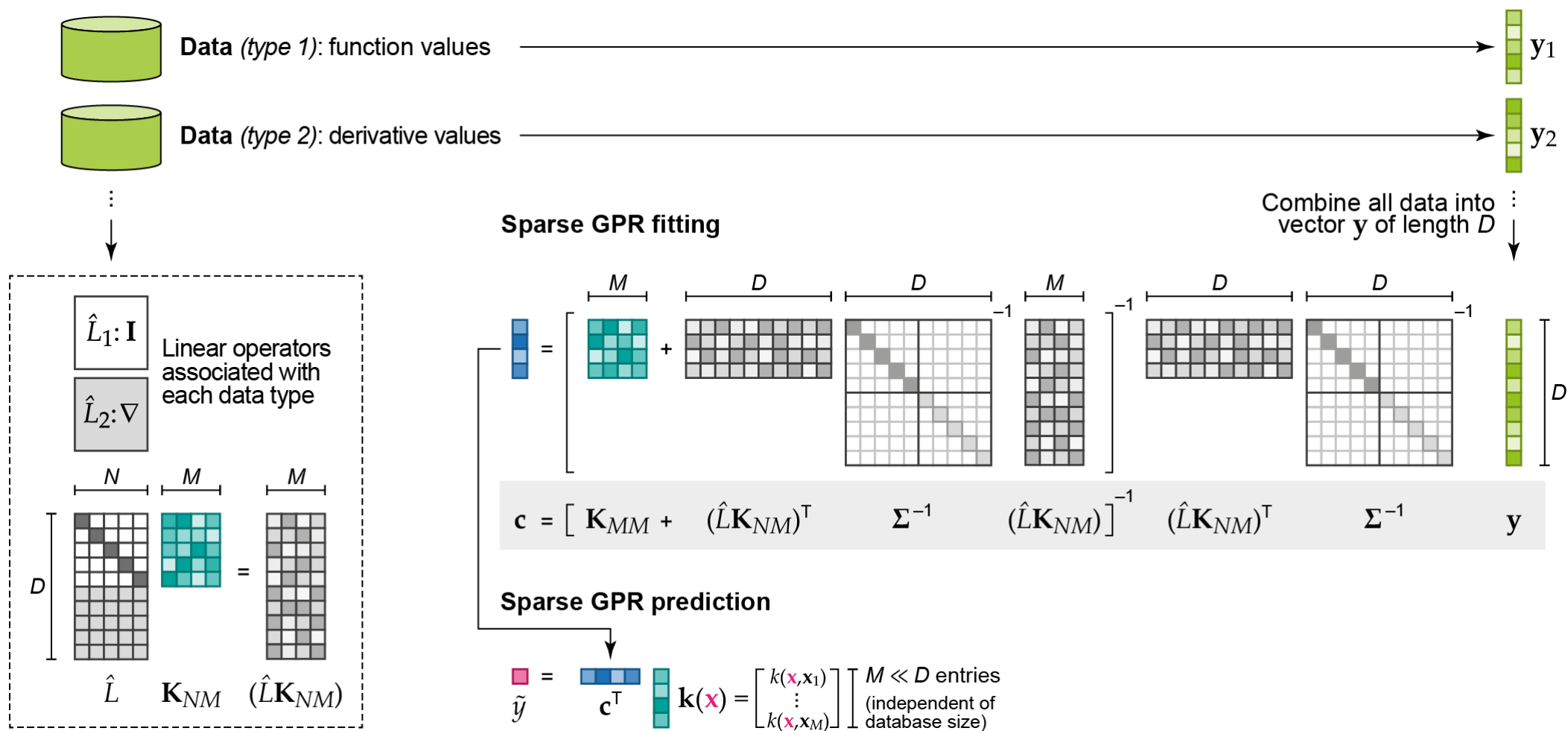


Gaussian process regression (GPR) for interatomic potentials

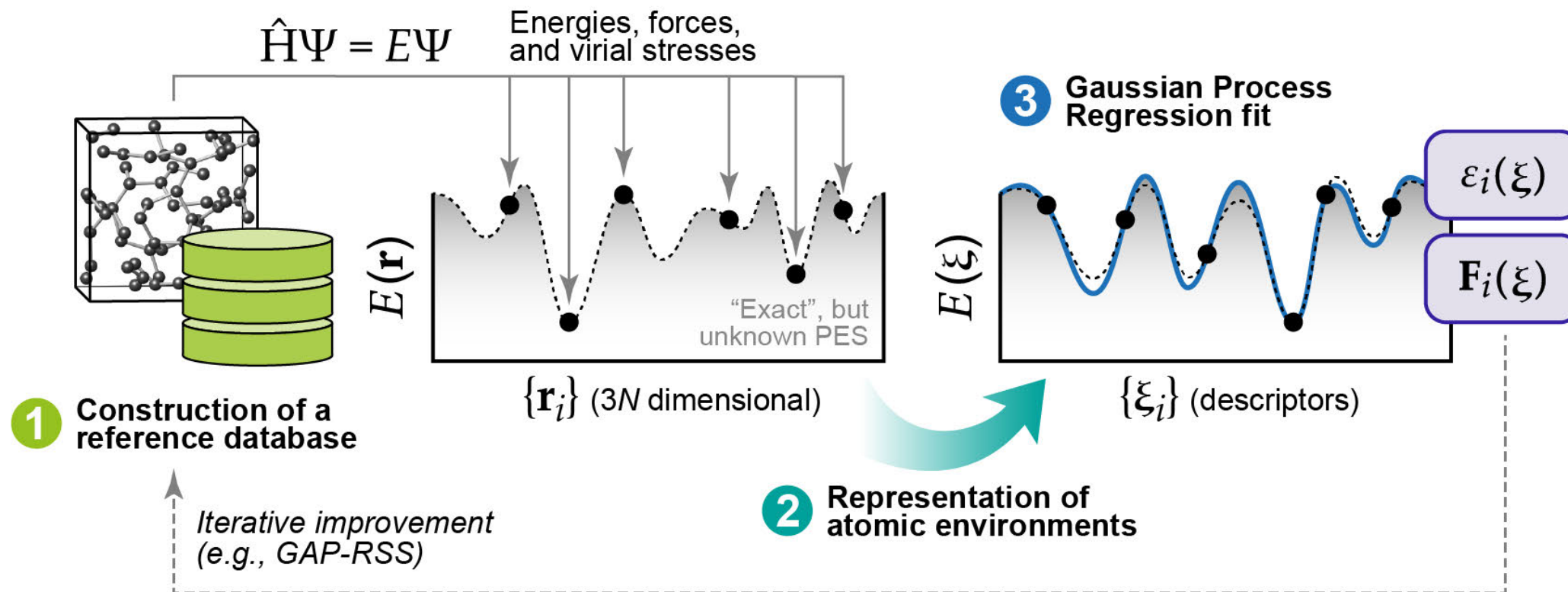


- So far, we have looked at learning an atomistic property, y , itself
- For interatomic potentials: much of the training data are forces (**derivatives**)

Gaussian process regression (GPR) for interatomic potentials



Gaussian approximation potential (GAP) models



Gaussian approximation potential (GAP) models

Reference data

- All available data are used: total energy, forces, and stresses (for periodic systems), combined into a single ML fit. The design of the input database is critical to the success of the model and has been a cornerstone of all presently available general-purpose GAPs. The selection of reference data is as much an area of ongoing methods development as is that of representation and regression (section 4.1).

The fit itself

- Hyperparameters of the GAP model are chosen and fixed a priori as much as possible and optimized only where required. The main hyperparameters are (i) the relevant length scales, which define the cutoff radius and the smoothness of the kernel, and (ii) the expected errors (arising both from noise in the input data and limitations of the model, e.g., due to the necessarily finite cutoff radius; section 4.4), which determine the regularization of the fit (section 4.6).

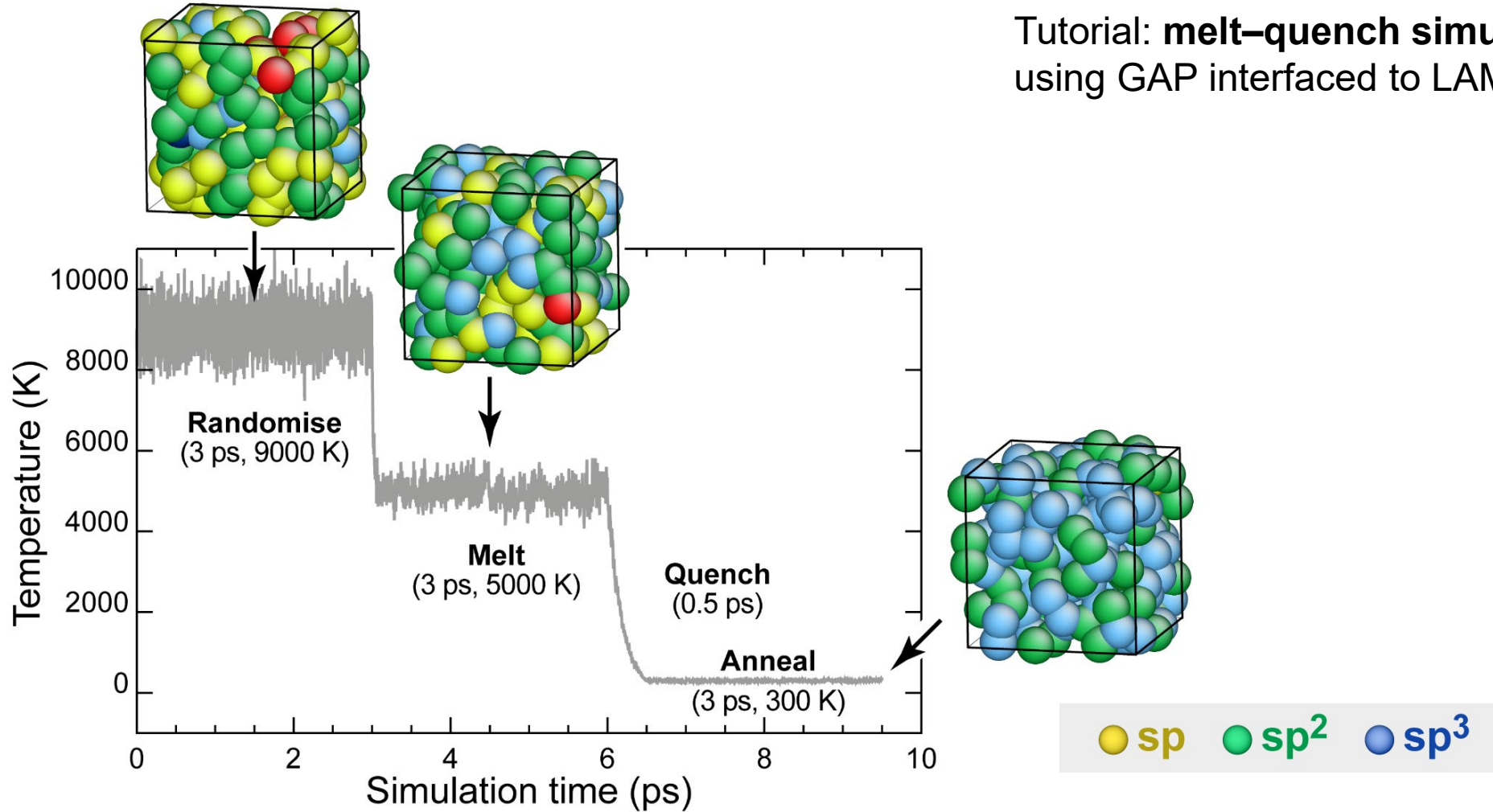
Representation

- The choice and specification of structural descriptors (representation) is tightly coupled with the choice of kernels, and both are an essential part of the user input. They incorporate prior knowledge about the nature of the potential-energy function—specifically, its regularity. Commonly used examples are distances and angles between atoms together with a squared exponential (Gaussian) kernel, or the many-body SOAP representation with a polynomial kernel. These are not mutually exclusive: low-dimensional kernel models can be fitted together with many-body ones, with appropriate weighting between them. All representations and kernels in GAP have finite distance cut-offs, typically about 5–6 Å, and therefore they represent the local environments of the atoms (section 4.2).

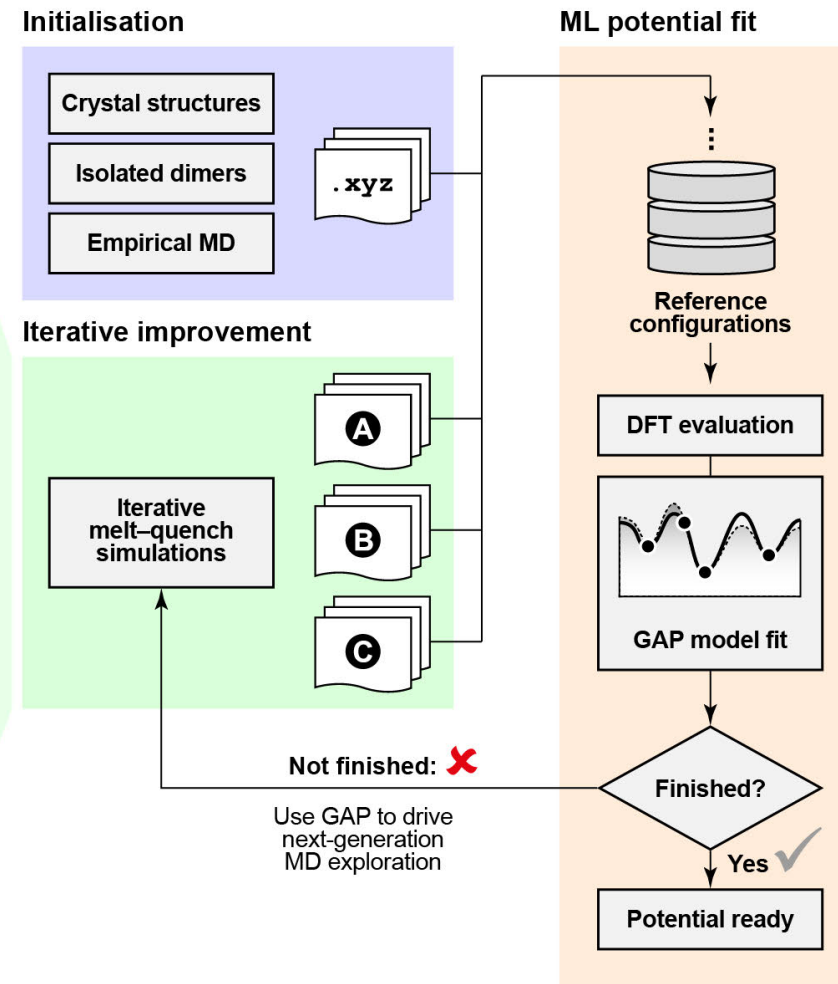
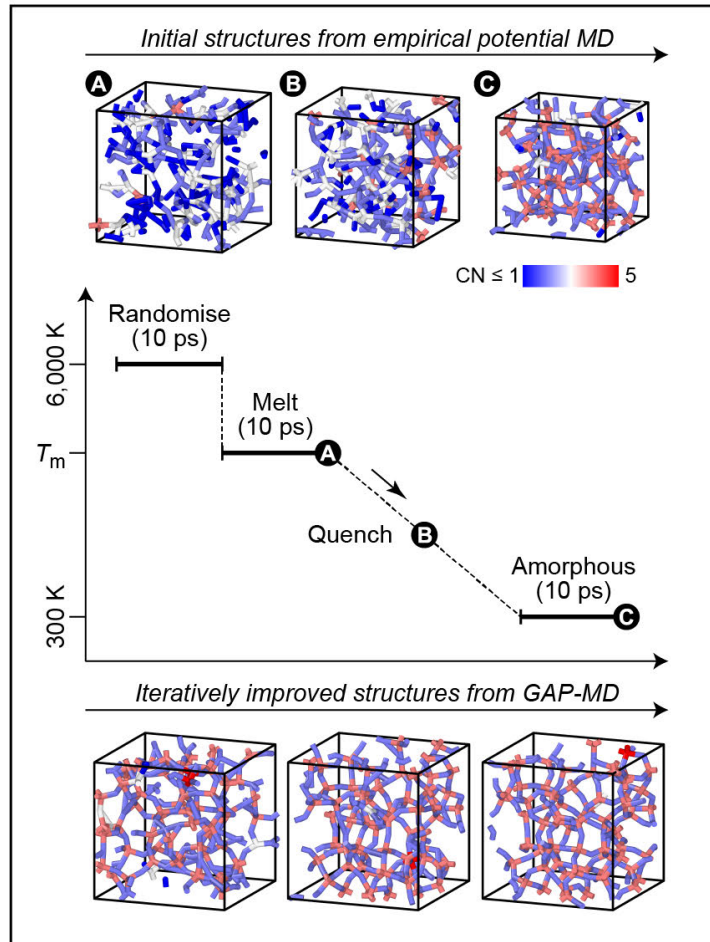
Available for non-commercial research at
<https://github.com/libAtoms/QUIP>

Gaussian approximation potential (GAP) models

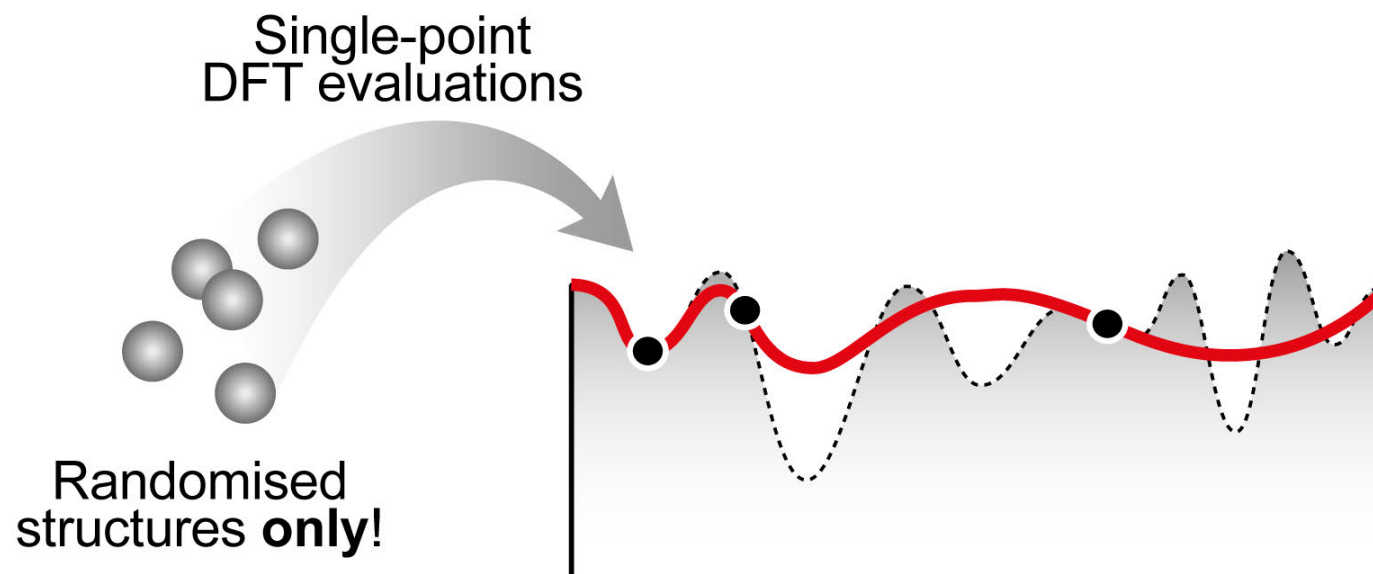
Tutorial: **melt-quench simulations**
using GAP interfaced to LAMMPS



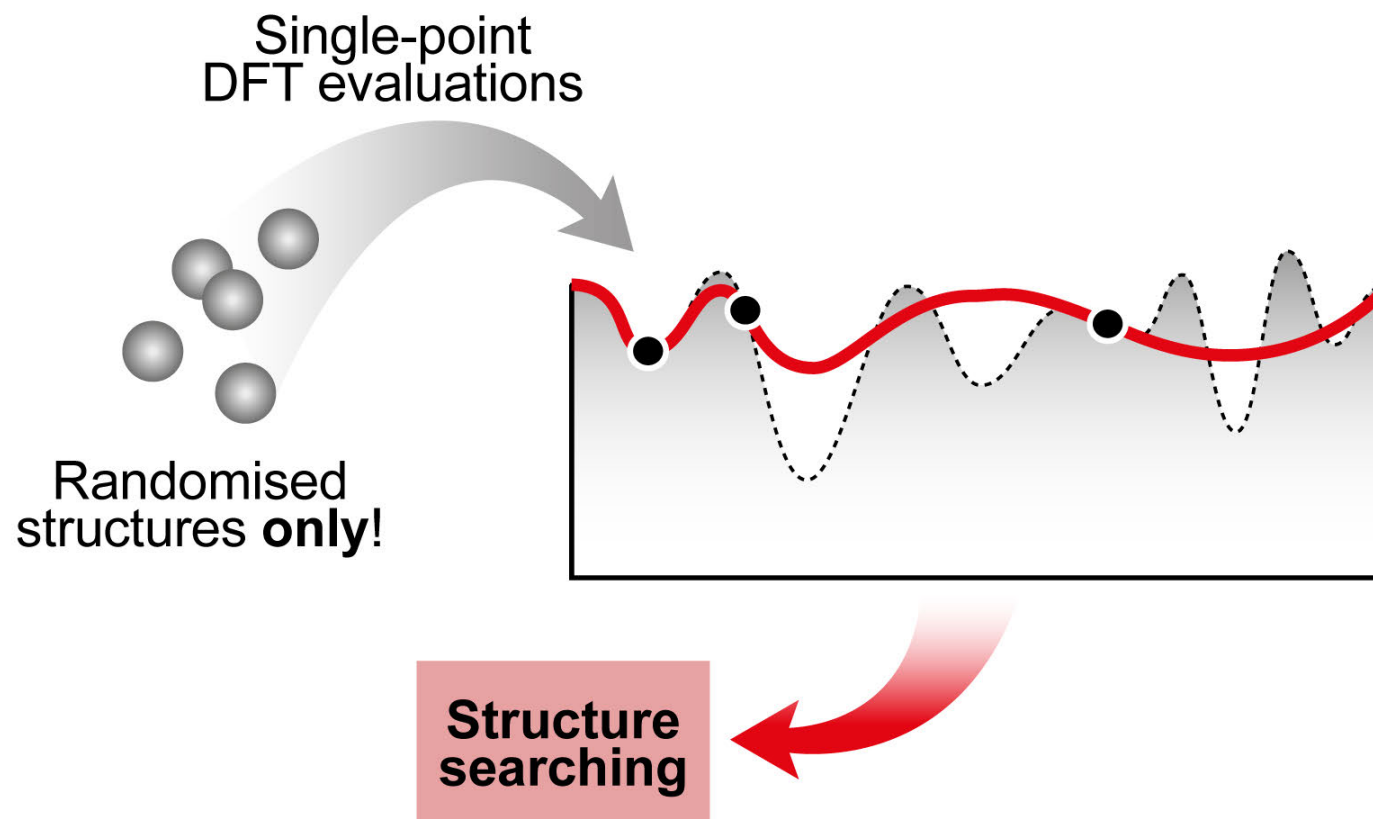
Reference data are critical



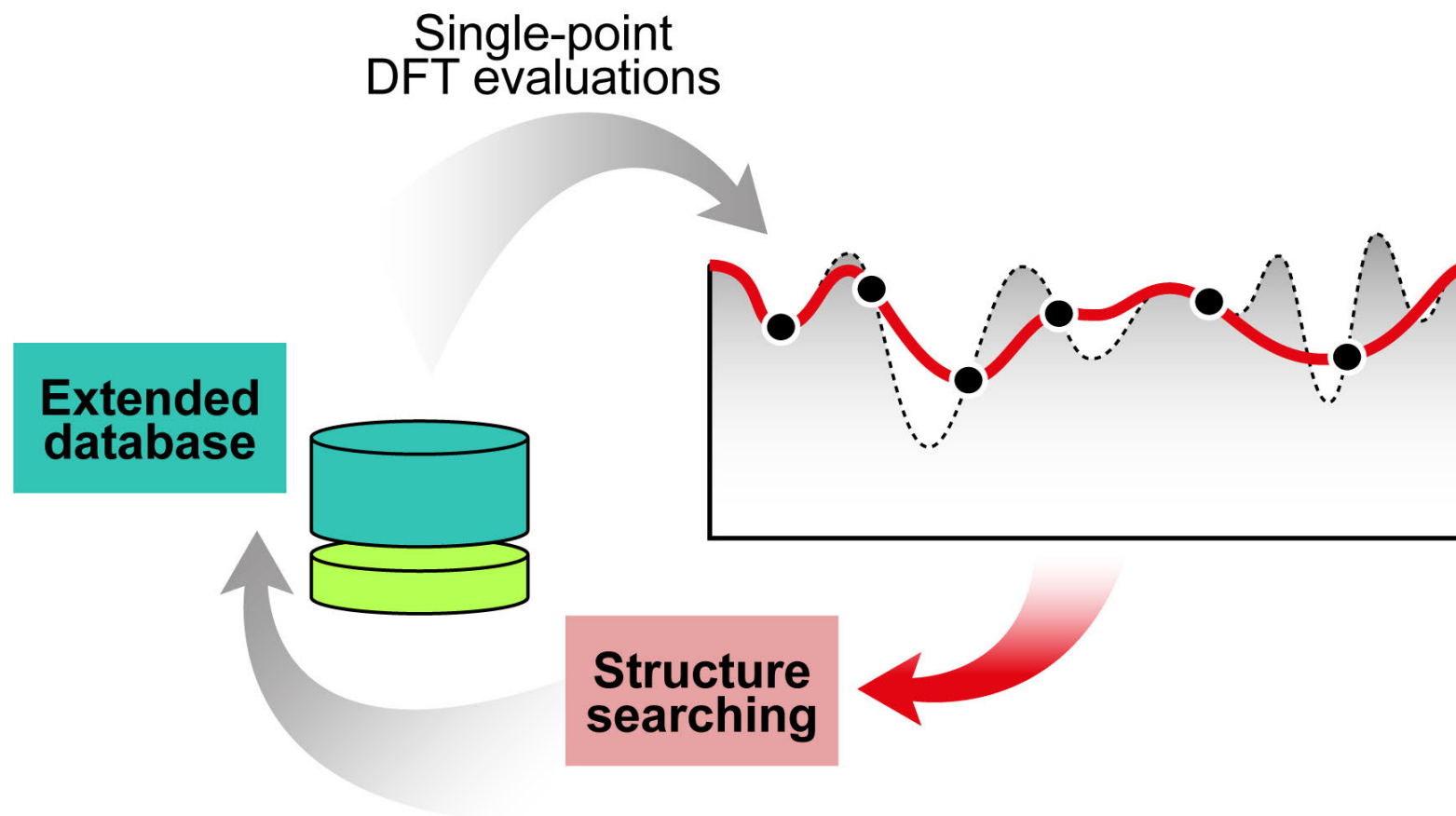
De novo exploration & fitting



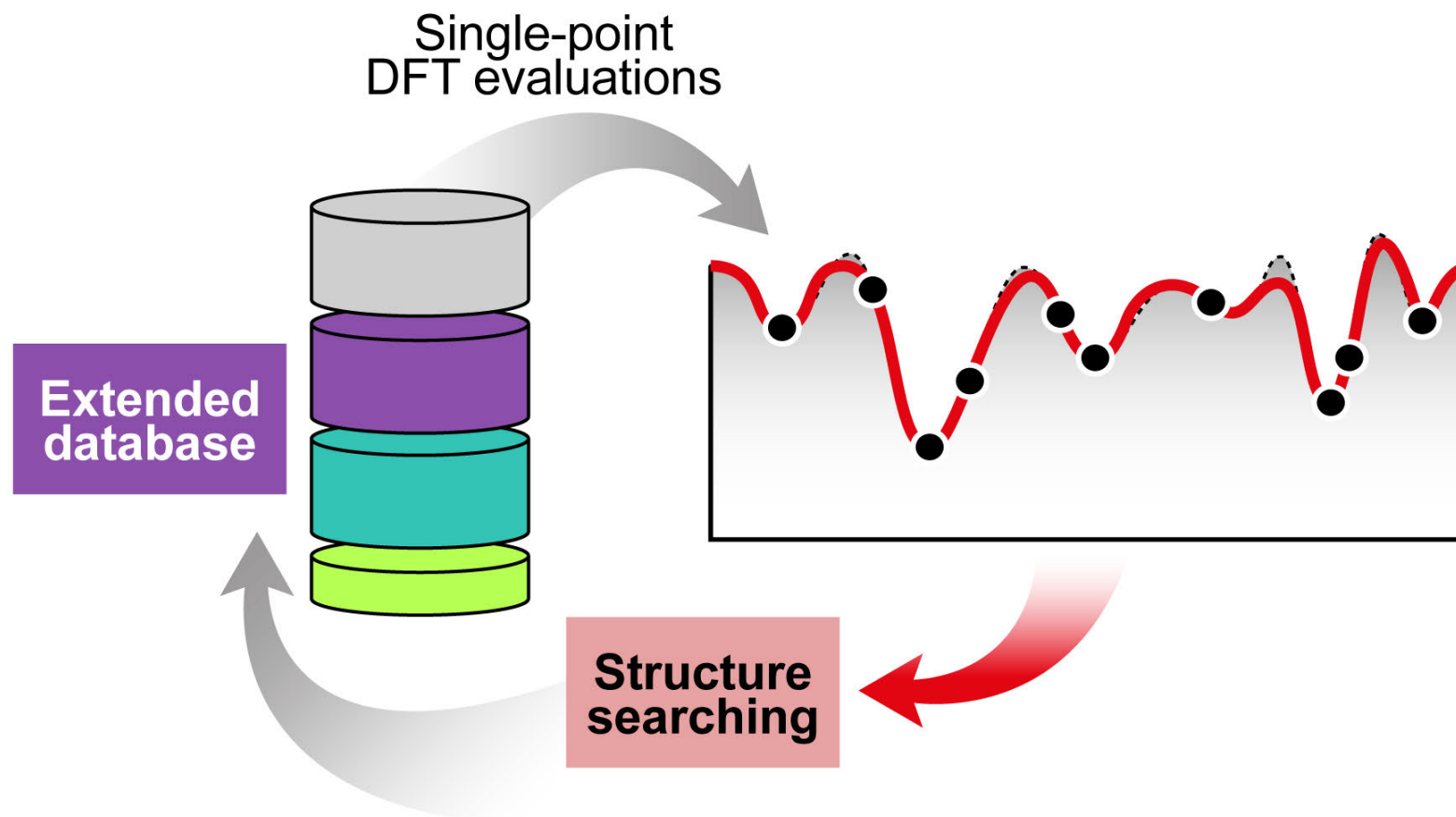
De novo exploration & fitting



De novo exploration & fitting



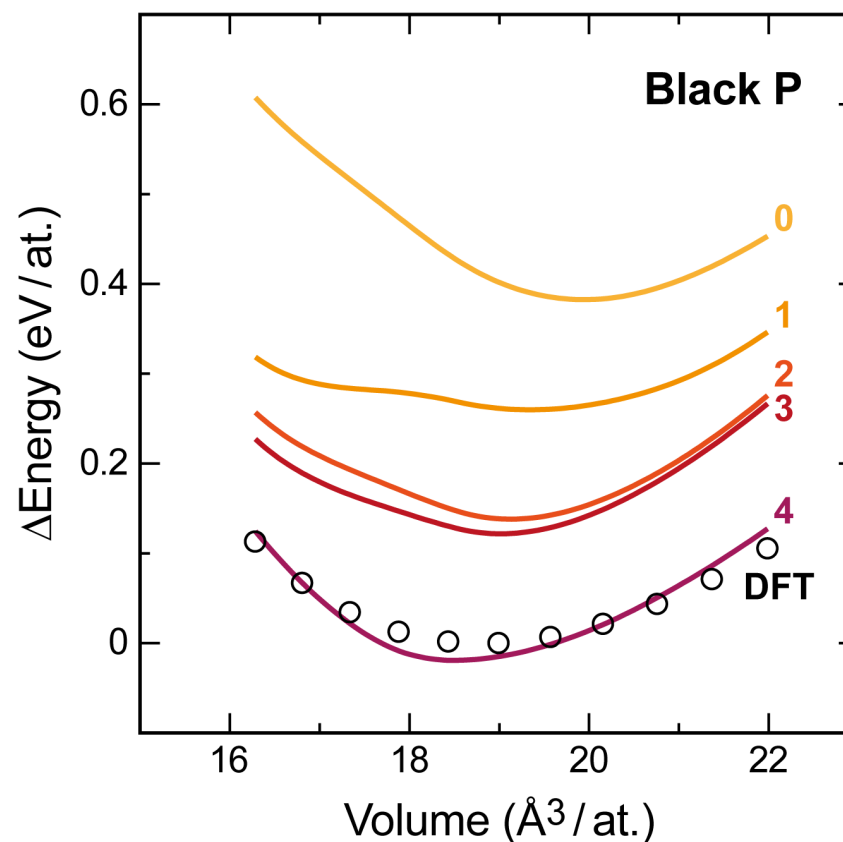
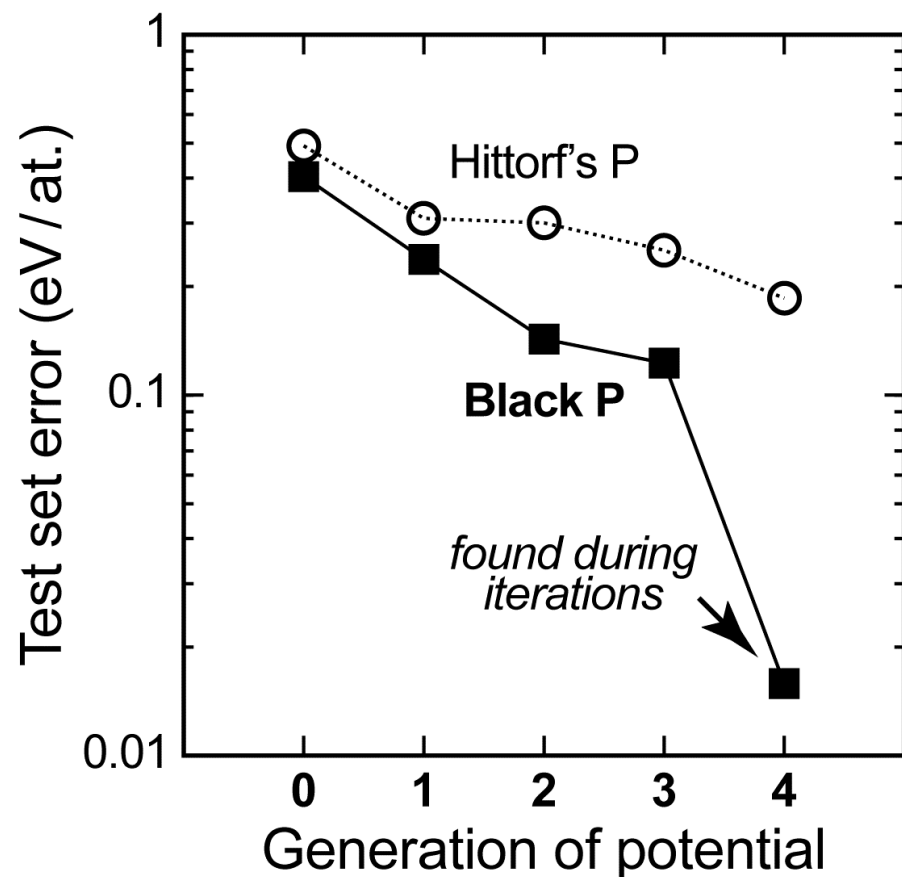
De novo exploration & fitting



VLD, C. J. Pickard, G. Csányi, *Phys. Rev. Lett.* **2018**, *120*, 156001

VLD, D. M. Proserpio, G. Csányi, C. J. Pickard, *Faraday Discuss.* **2018**, *211*, 45

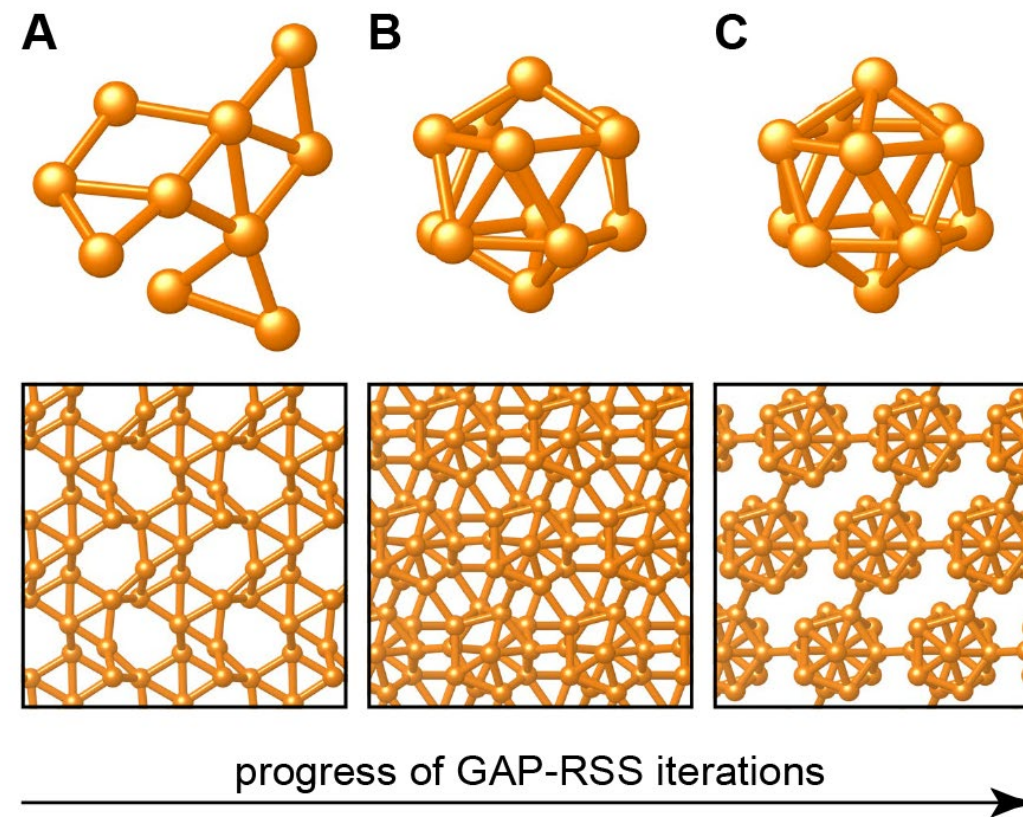
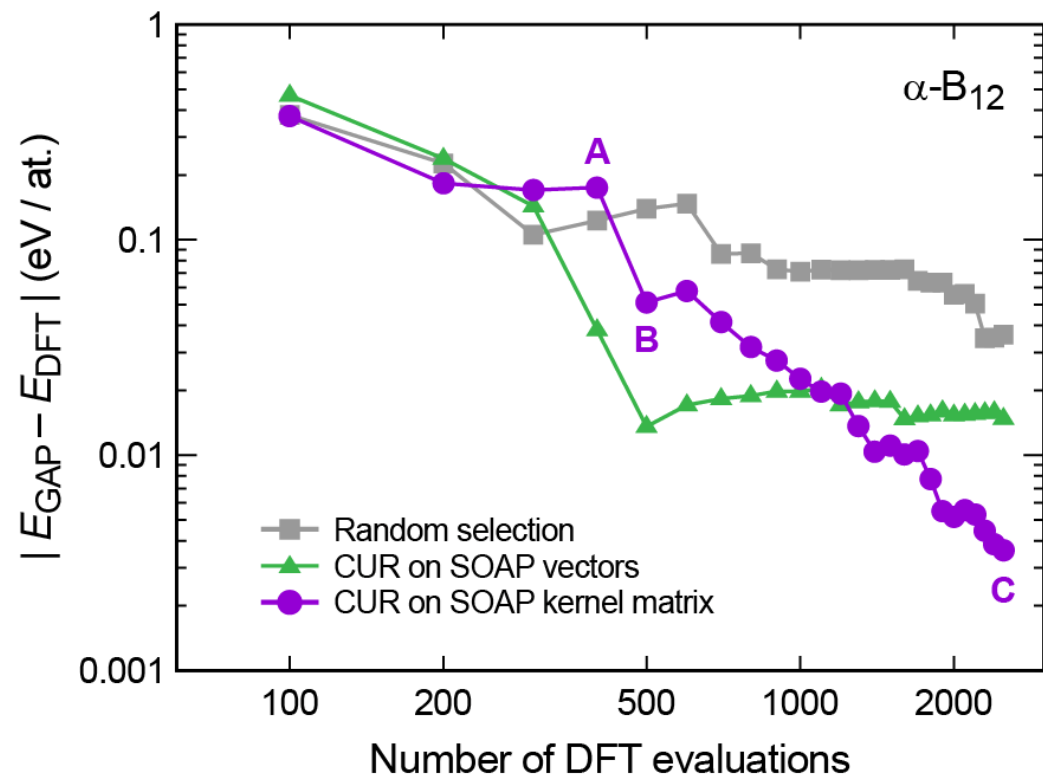
De novo exploration & fitting



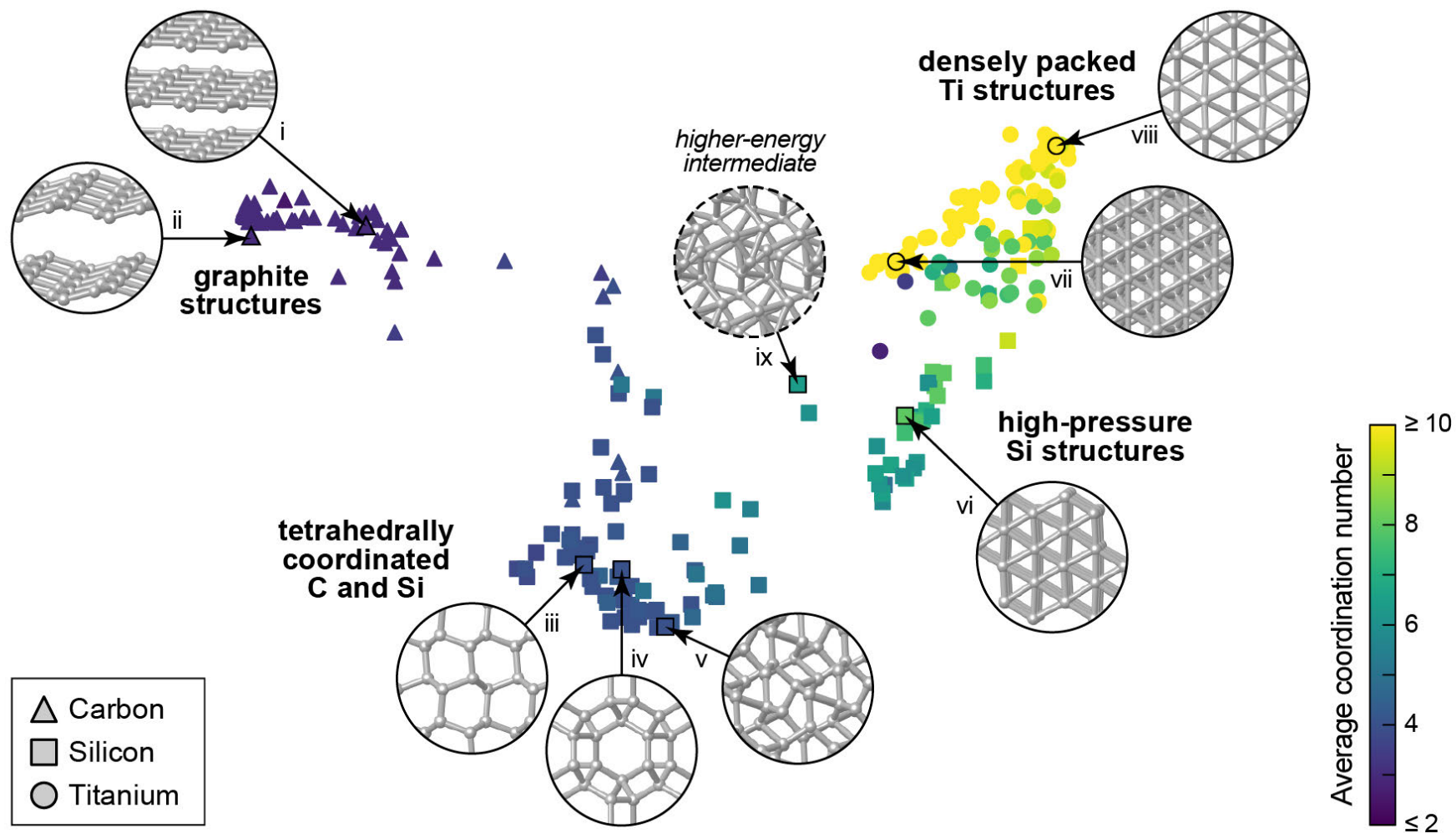
VLD, C. J. Pickard, G. Csányi, *Phys. Rev. Lett.* **2018**, 120, 156001

VLD, D. M. Proserpio, G. Csányi, C. J. Pickard, *Faraday Discuss.* **2018**, 211, 45

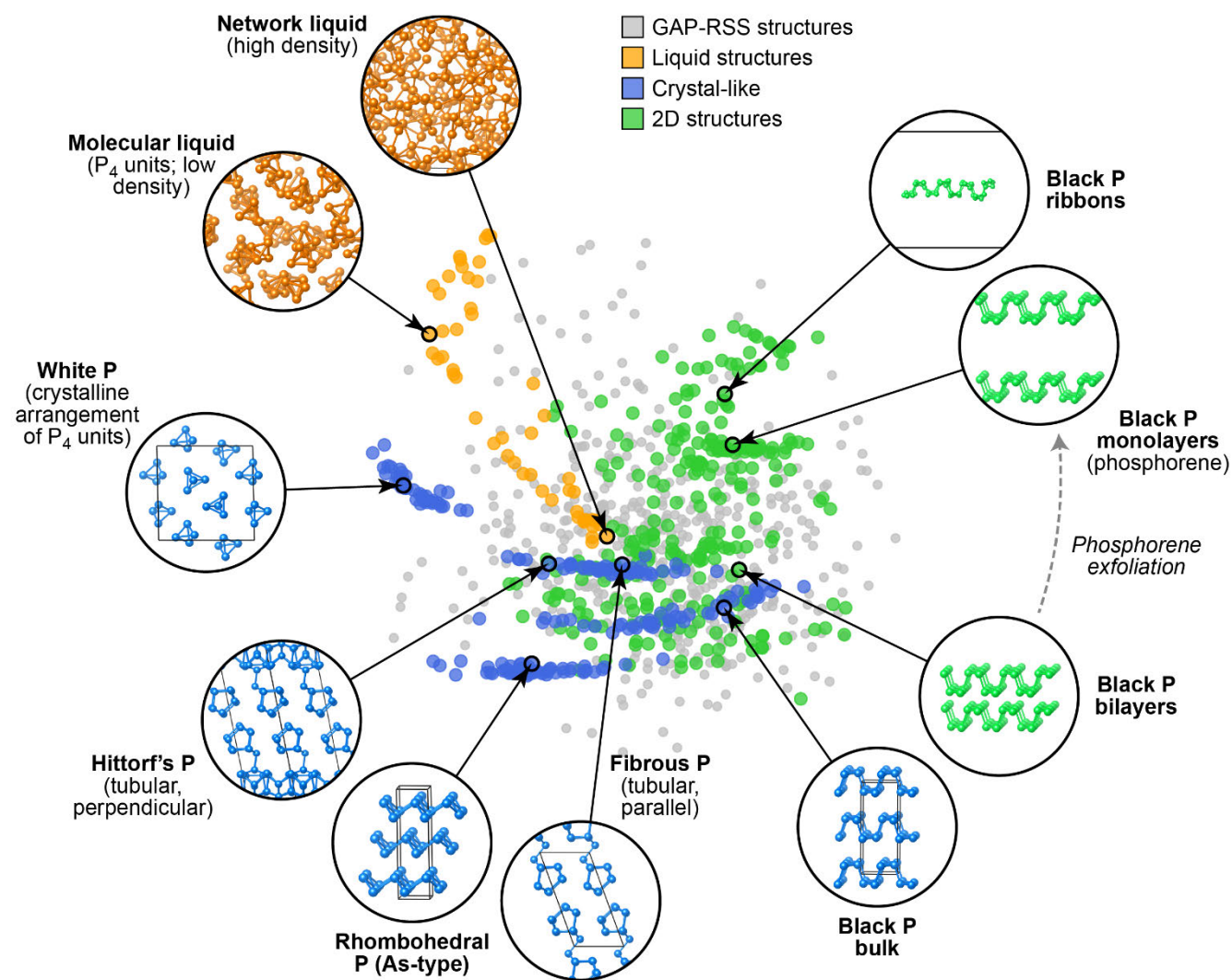
De novo exploration & fitting



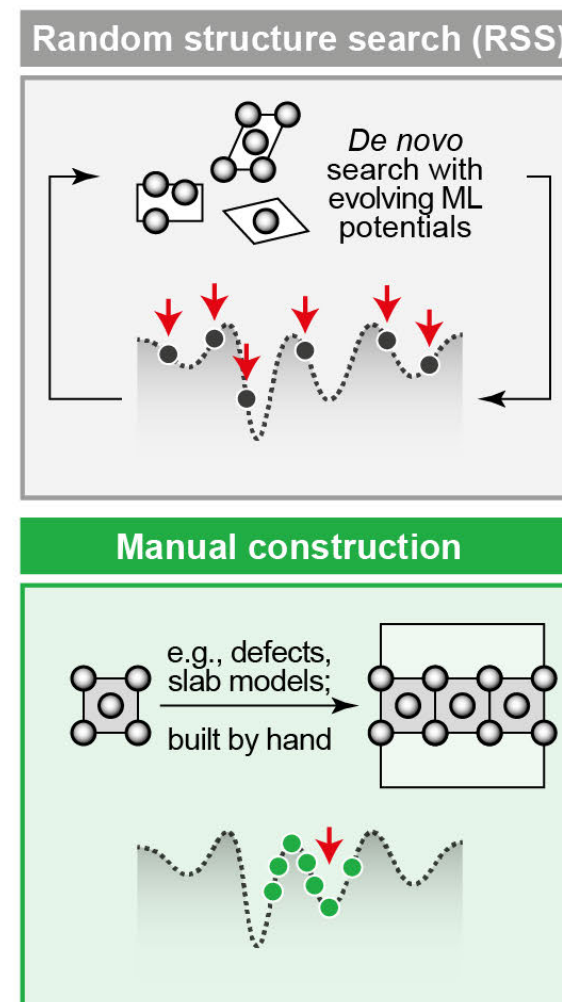
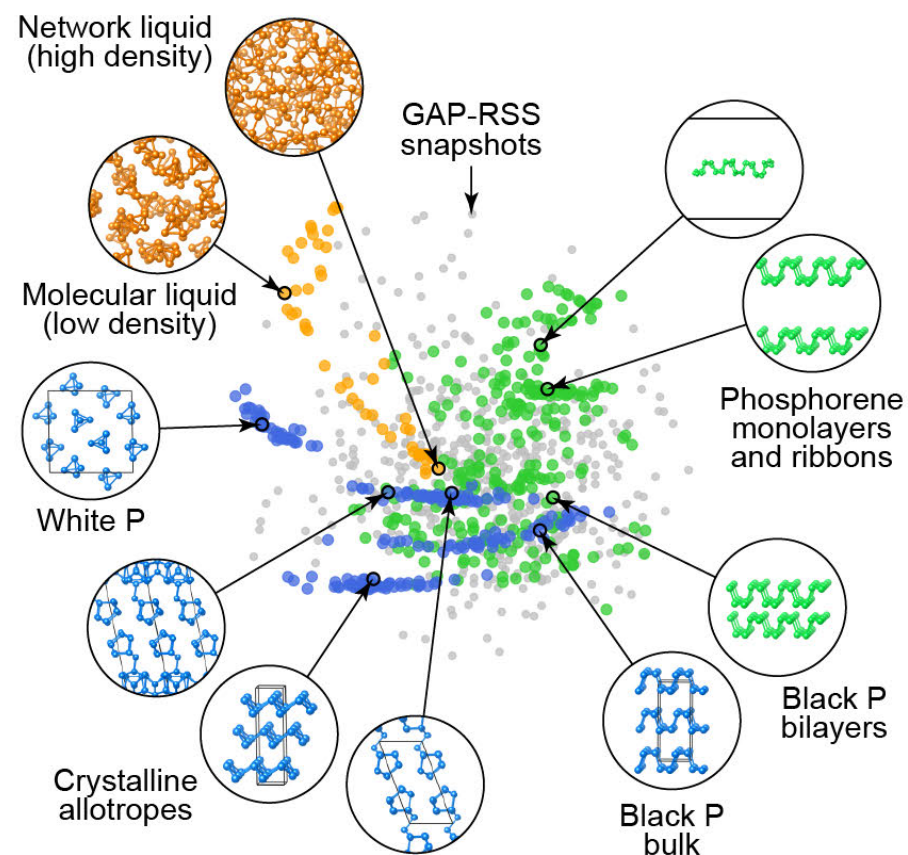
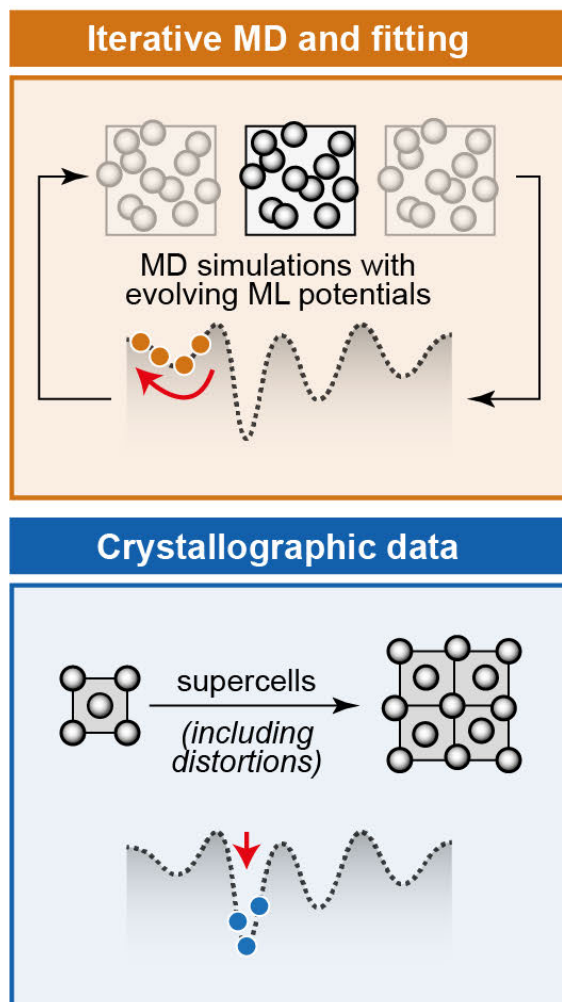
De novo exploration & fitting



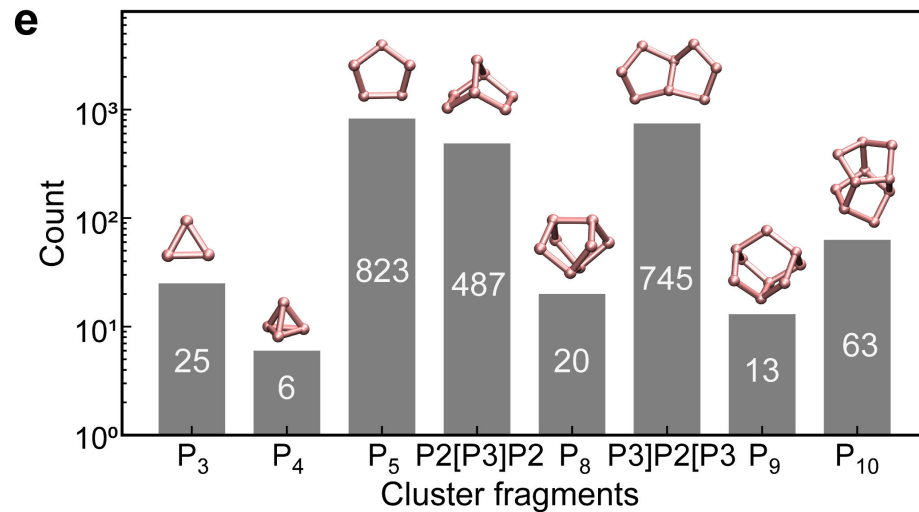
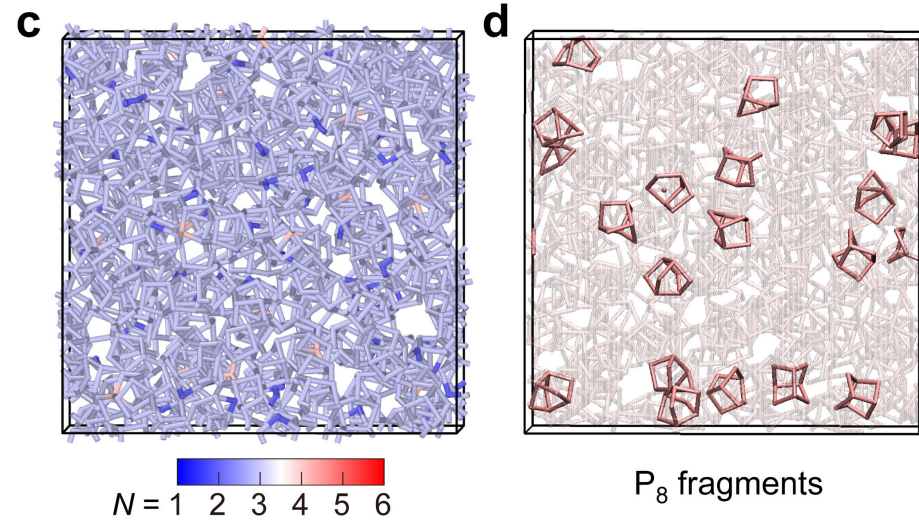
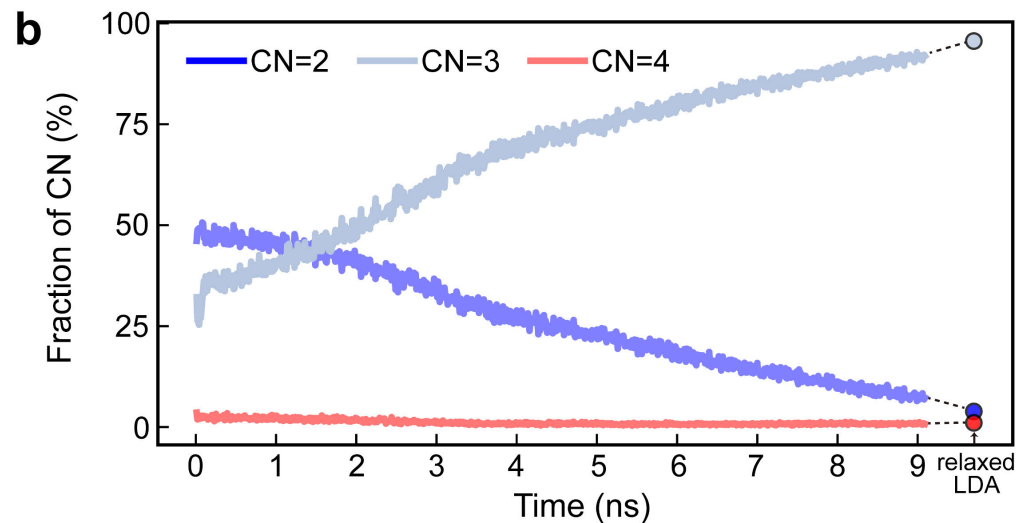
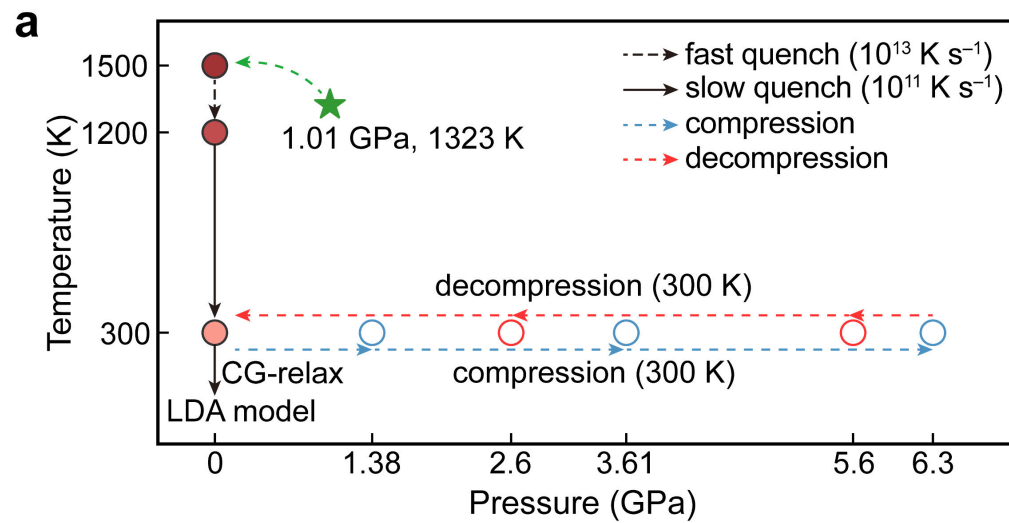
GAP-RSS as *starting point* for ML potential databases



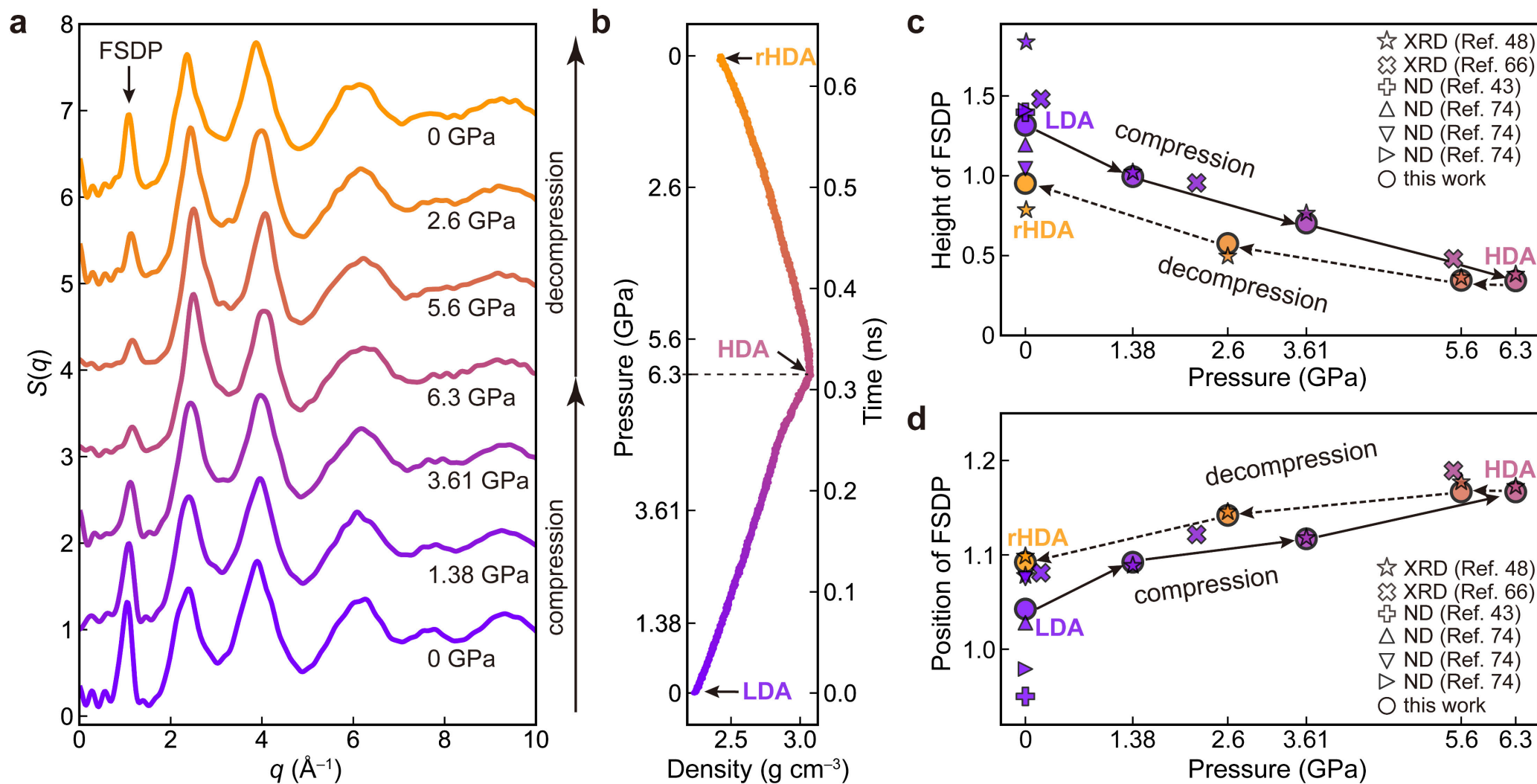
Towards automated & general ML potentials



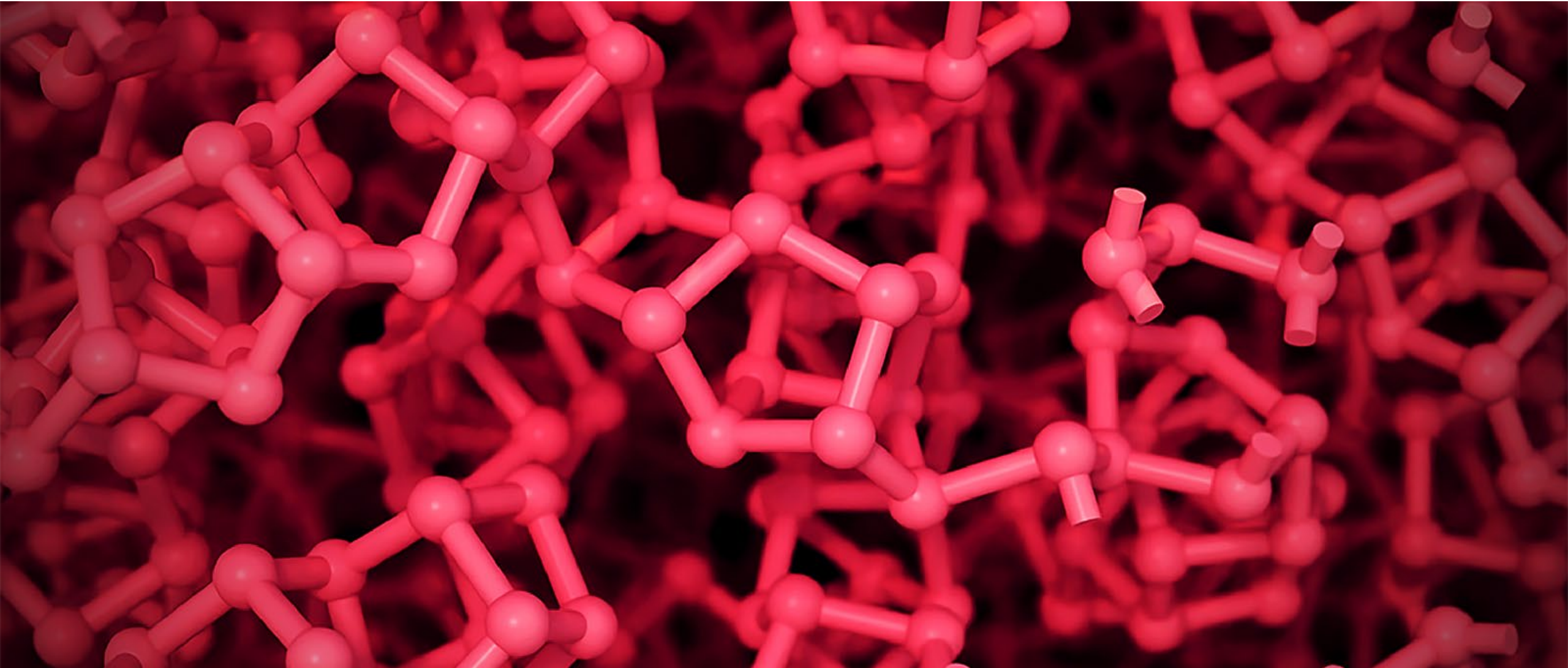
The structure of amorphous red phosphorus



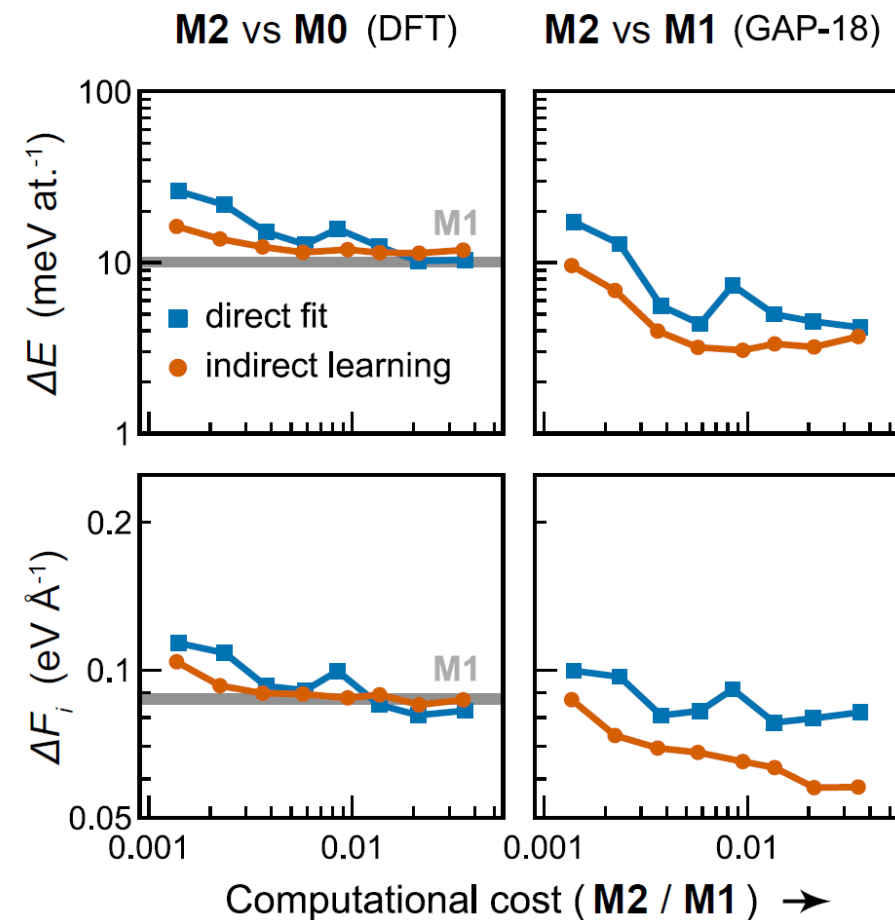
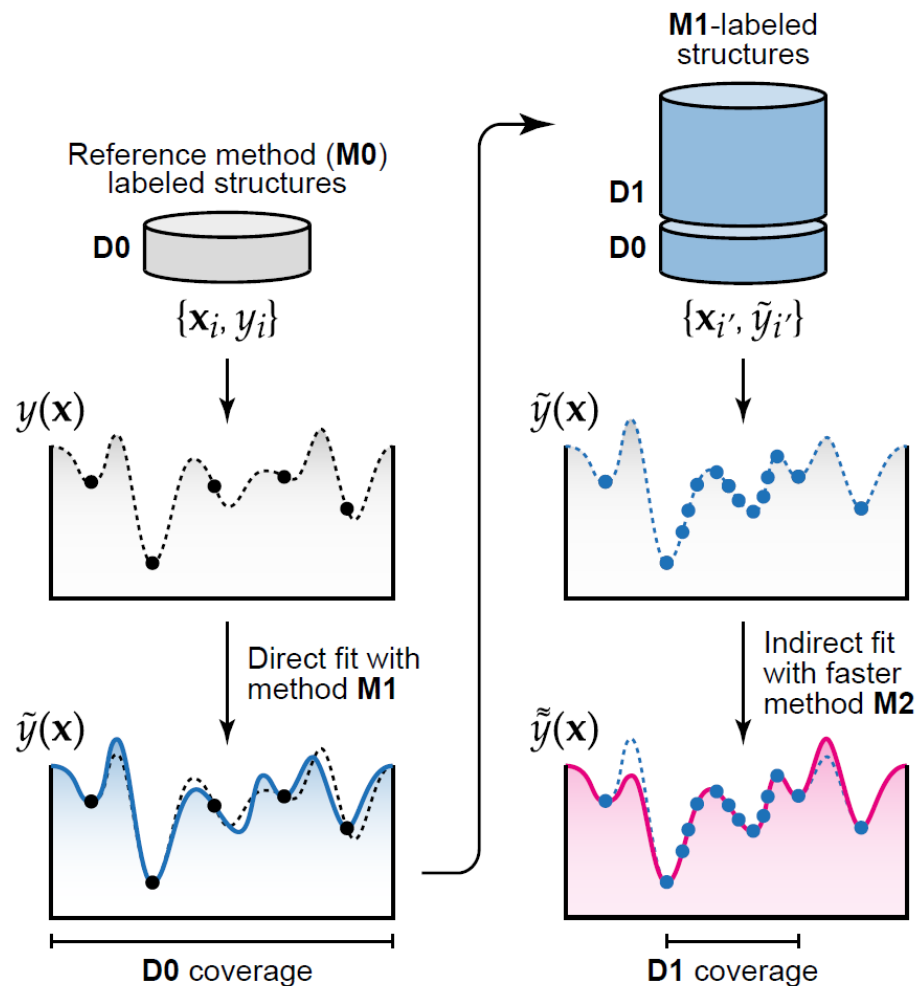
The structure of amorphous red phosphorus



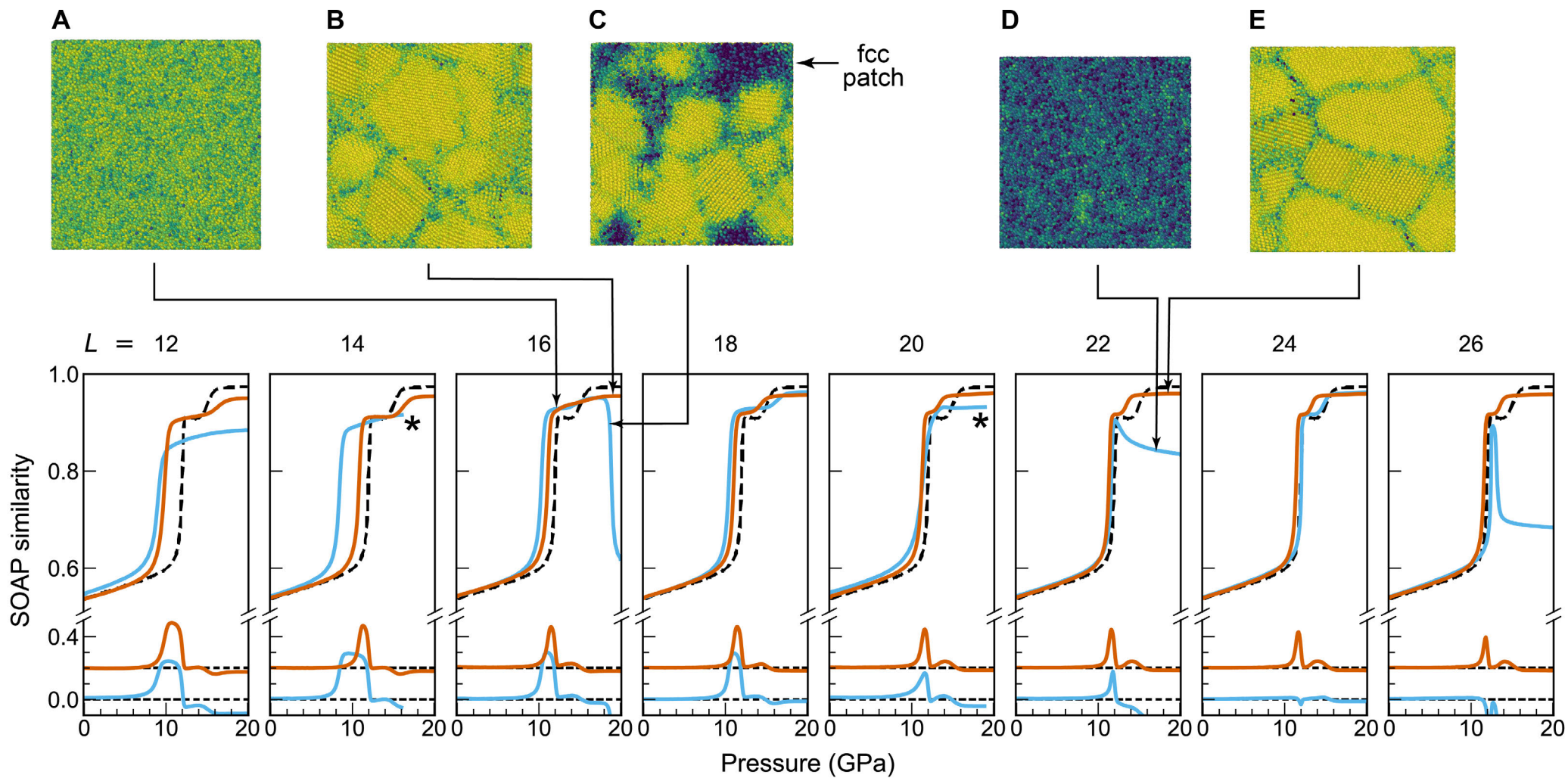
The structure of amorphous red phosphorus



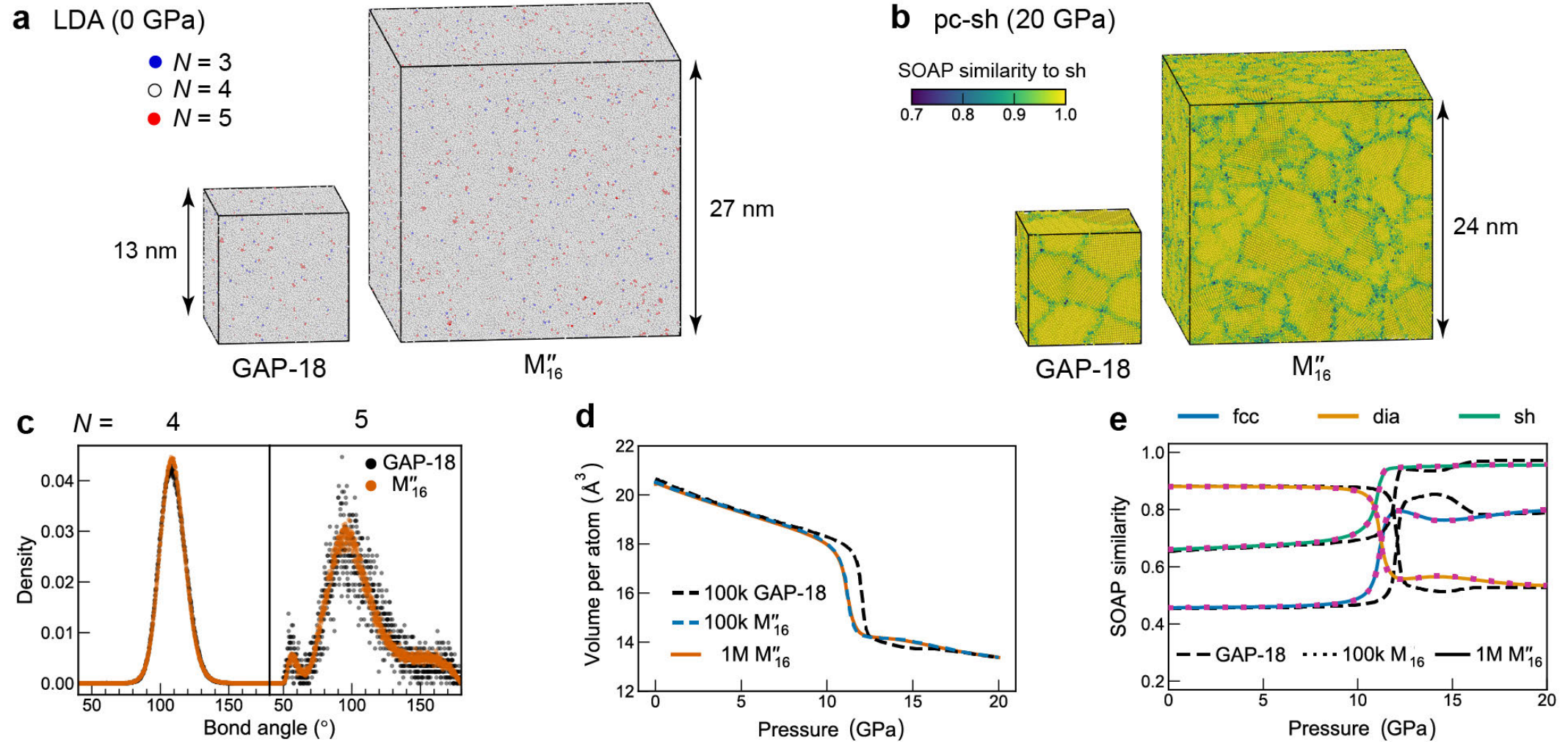
How to validate interatomic potentials



How to validate interatomic potentials



How to validate interatomic potentials



Executive summary

- **ML potentials** are increasingly popular simulation tools for materials modelling: accurate, flexible, and *fast*
(*Adv. Mater.* **2019**, 31, 1902765; *Chem. Rev.* **2021**, 121, 10073)
- They are therefore giving new insight into **real materials**, from fundamental questions to full device-scale simulations
(a-Si: *Nature* **2021**, 589, 59; phase-change materials: arXiv:2207.14228)
- **Validating potentials** is becoming increasingly important, and many ideas here will be transferable from GAP to other models
(arXiv:2211.12484 – [see also tutorial by Joe](#))

Thank you!

Collaborators

ML potential development:

Gábor Csányi (Cambridge)

Amorphous silicon:

Noam Bernstein (US NRL)

Michele Ceriotti (EPFL)

David Drabold (Ohio)

Stephen Elliott (Cambridge)

Mark Wilson (Oxford)

Phase-change materials:

En Ma (Xi'an Jiaotong)

Wei Zhang (Xi'an Jiaotong)

...and other key collaborators,
including:

Karsten Albe (Darmstadt)

Miguel Caro (Aalto)

Janine George (BAM Berlin)

Andrew Goodwin (Oxford)

Davide Proserpio (Milan)

VLD Group Oxford

Research Associates

Dr Nijamudheen

Abdulrahiman

DPHil (PhD) students

Zakariya El-Machachi

Zoé Faure Beaulieu

John Gardner

Joe Morrow

Tom Nicholas

Louise Rosset

Daniel Thomas du Toit

Yuxing Zhou

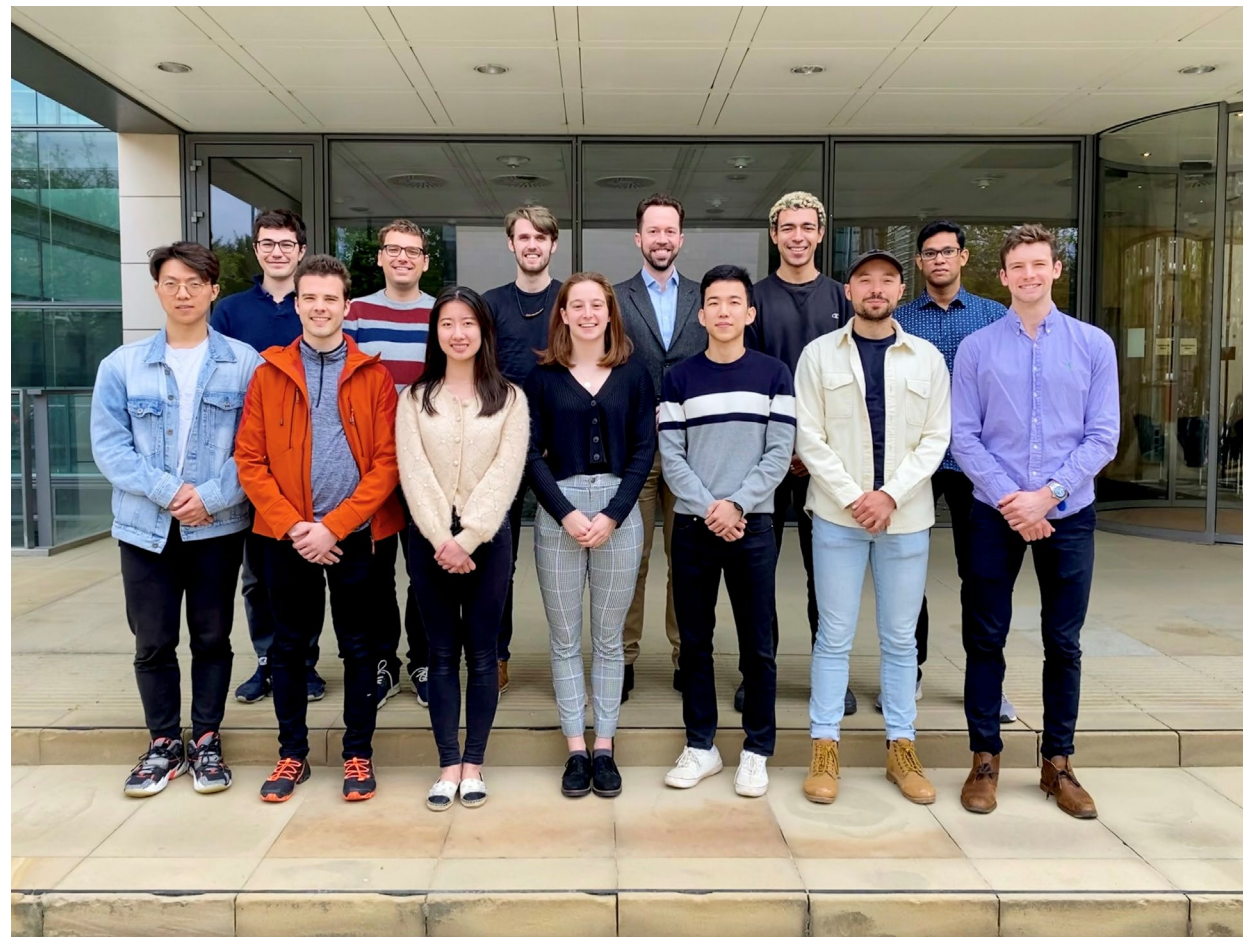
Part II (MChem) students

Kathryn Baker

Muxue Chen

Maximilian Weber

Aleksandra Zawadzka



(ERC guarantee funding)



Engineering and
Physical Sciences
Research Council

LEVERHULME
TRUST

