

PERSPECTIVE • OPEN ACCESS

Artificial intelligence for advanced functional materials: exploring current and future directions

To cite this article: Cristiano Malica *et al* 2025 *J. Phys. Mater.* **8** 021001

View the [article online](#) for updates and enhancements.

You may also like

- [The 2025 roadmap to ultrafast dynamics: frontiers of theoretical and computational modeling](#)

Fabio Caruso, Michael A Sentef, Claudio Attaccalite et al.

- [Rigid-dipole magnetic nanoparticles for sub-second 3D viscosity imaging](#)

J M Costa, G F Resende, Y Gu et al.

- [Towards high loading cesium lead halide nanocomposites for radiation detection](#)

Jan Král, Kateina Dcká, Vojtch Zabloudil et al.

Journal of Physics: Materials



OPEN ACCESS

RECEIVED
4 February 2025

REVISED
8 March 2025

ACCEPTED FOR PUBLICATION
19 March 2025

PUBLISHED
23 April 2025

Original content from
this work may be used
under the terms of the
[Creative Commons
Attribution 4.0 licence](#).

Any further distribution
of this work must
maintain attribution to
the author(s) and the title
of the work, journal
citation and DOI.



PERSPECTIVE

Artificial intelligence for advanced functional materials: exploring current and future directions

Cristiano Malica^{1,*} , Kostya S Novoselov^{2,3} , Amanda S Barnard⁴ , Sergei V Kalinin^{5,6} , Steven R Spurgeon^{7,8} , Karsten Reuter⁹ , Maite Alducin^{10,11} , Volker L Deringer¹² , Gábor Csányi¹³ , Nicola Marzari^{1,14} , Shirong Huang¹⁵ , Gianaurelio Cuniberti¹⁵ , Qiushi Deng¹⁶ , Pablo Ordejón¹⁷ , Ivan Cole¹⁸ , Kamal Choudhary¹⁹ , Kedar Hippalgaonkar^{20,21} , Ruiming Zhu^{20,21}, O Anatole von Lilienfeld^{22,23,24} , Mohamed Hibat-Allah^{23,25} , Juan Carrasquilla²⁶ , Giulia Cisotto²⁷ , Alberto Zancanaro²⁸ , Wolfgang Wenzel²⁹ , Andrea C Ferrari³⁰ , Andrey Ustyuzhanin^{2,31} and Stephan Roche^{17,32,*}

¹ U Bremen Excellence Chair, Bremen Center for Computational Materials Science, and MAPEX Center for Materials and Processes, University of Bremen, Bremen, D-28359, Germany

² Institute for Functional Intelligent Materials, National University of Singapore, Singapore, 117544, Singapore

³ Department of Materials Science and Engineering, National University of Singapore, Singapore, 117575, Singapore

⁴ School of Computing, Australian National University, Acton, Australia

⁵ Department of Materials Science and Engineering, University of Tennessee, Knoxville, TN, United States of America

⁶ Pacific Northwest National Laboratory, Richland, WA, United States of America

⁷ National Renewable Energy Laboratory, Golden, CO, United States of America

⁸ University of Colorado, Boulder, CO, United States of America

⁹ Fritz-Haber-Institut der Max-Planck-Gesellschaft, Berlin, Germany

¹⁰ Centro de Física de Materiales (CFM/MPC), Donostia-San Sebastián, Spain

¹¹ Donostia International Physics Center, Donostia-San Sebastián, Spain

¹² Inorganic Chemistry Laboratory, Department of Chemistry, University of Oxford, Oxford, United Kingdom

¹³ Engineering Laboratory, University of Cambridge, Cambridge, United Kingdom

¹⁴ Theory and Simulation of Materials (THEOS), and National Centre for Computational Design and Discovery of Novel Materials (MARVEL), École Polytechnique Fédérale de Lausanne (EPFL), Lausanne, Switzerland

¹⁵ Institute for Materials Science and Max Bergmann Center for Biomaterials TUD Dresden University of Technology, Dresden, Germany

¹⁶ School of Engineering, RMIT University, Melbourne, Australia

¹⁷ Catalan Institute of Nanoscience and Nanotechnology ICN2 (CSIC and BIST), Campus UAB, Bellaterra, Spain

¹⁸ School of Engineering, Australia National University, Acton, Australia

¹⁹ Materials Science and Engineering Division, National Institute of Standards and Technology, Gaithersburg, MD, United States of America

²⁰ School of Materials Science and Engineering, Nanyang Technological University, Singapore, 639798, Singapore

²¹ Institute of Materials Research and Engineering, Agency for Science, Technology and Research (A*STAR), Singapore, 138634, Singapore

²² Departments of Chemistry, Materials Science & Engineering, Physics, and the Acceleration Consortium, University of Toronto, Canada

²³ Vector Institute, Toronto, Ontario, Canada

²⁴ ML Group, Technische Universität Berlin and Institute for the Foundations of Learning and Data, Berlin, Germany

²⁵ Department of Applied Mathematics, University of Waterloo, Waterloo, ON, Canada

²⁶ Institute for Theoretical Physics, ETH Zürich, Switzerland

²⁷ Department of Mathematics, Informatics and Geosciences, University of Trieste, Italy

²⁸ Department of Information Engineering, University of Padova, Italy

²⁹ Karlsruhe Institute of Technology, Eggenstein-Leopoldshafen, Germany

³⁰ Cambridge Graphene Centre, University of Cambridge, Cambridge, United Kingdom

³¹ Constructor University, Bremen, Germany

³² ICREA—Institució Catalana de Recerca i Estudis Avançats, Barcelona, Spain

* Authors to whom any correspondence should be addressed.

E-mail: cmalica@uni-bremen.de and stephan.roche@icn2.cat

Keywords: artificial intelligence (AI), machine learning (ML), materials science

Abstract

This perspective addresses the topic of harnessing the tools of artificial intelligence (AI) for boosting innovation in functional materials design and engineering as well as discovering new materials for targeted applications in energy storage, biomedicine, composites, nanoelectronics or quantum technologies. It gives a current view of experts in the field, insisting on challenges and

opportunities provided by the development of large materials databases, novel schemes for implementing AI into materials production and characterization as well as progress in the quest of simulating physical and chemical properties of realistic atomic models reaching the trillion atoms scale and with near *ab initio* accuracy.

1. Introduction

1.1. Innovative material design and engineering

The design and engineering of innovative advanced materials (IAMs) is facing a variety of challenges in today's industries, including the access to proper design strategies for reaching upper performances of materials for targeted applications, the discovery of alternatives and the search for more functional intelligent materials which can help solving health, energy or environmental issues. Such new functionalities often require complex material structure, such as complex alloys, composites, heterostructures, etc. Traditional methods of material modelling cannot cope with such demands. In this context the use of artificial intelligence (AI) tools has become cornerstone for boosting innovation strategies and ensuring sustainability and safe-by-design approaches.

The development and availability of material databases is a fundamental part for further training AI models (machine learning (ML) and so on) able to cope with diversity and complexity and extract hidden information which could ultimately offer further intelligent guidance in optimisation and also materials (property) discovery. However, the multiplicity of databases also calls for efforts in improving universality in development languages, interoperability as well as integrated workflows which can connect information concerning the structure to the end physical or chemical properties of materials of concern. More, the needs for more and more predictive modelling and capability of simulation tools to cope with systems reaching the trillion atoms-scale limit (while keeping a near *ab initio* accuracy) presents grand challenges and demands for novel workflows to be developed and more synergies between academic research and industrial developments. To this end, property and functions-oriented databases are required, especially when aiming at solving the inverse problem of finding the material with predetermined properties.

On the other side, intelligent materials are defined as structural materials with advanced functionalities and can be classified as structure-mimetic (mimicking the structure of organisms) and function-mimetic (mimicking the function of organisms). Intelligent materials target self-sensing of the material during its use (e.g. damage, loads, shape, temperature, pressure, etc) and/or target adaptive actuation (e.g., changing deformation, colour, shape, inner stresses, stiffness, temperature, etc) which depends on the biological or environmental conditions (humidity, pH, temperature, etc). The quest for more innovative intelligent materials depends on the capability of AI tools to provide a proper booster in benchmarking, fast and precise analysis and extrapolation for materials design.

As a result, the synergy between the activities of computer scientists, material engineers and AI developers with the experimental activities and elaboration of novel types of functional intelligent materials has become key for advanced development in innovative materials design and engineering.

In this perspective, we provide snapshots about efforts made in a variety of different fields and visions of international experts, searching for the same common objective, that is the deployment of AI tools for an accelerated development of materials design and deeper access to hidden dimensions of materials growth, structure–property correlations and reverse engineering strategies. Such a vast field of research calls for structuring the exploding amount of information and also for implementing chains of tools able to communicate information and extract essential parameters that can be ultimately accessible to the largest possible public. In that perspective, international events such as AI4AM (www.ai4am.net) are enabling platforms to gather communities, facilitate networking, roadmapping and ultimately enhance our knowledge and methodologies.

2. Higher-order pattern recognition for materials informatics using explainable AI

Explainable AI (XAI) is an emerging field in computer science based in statistics that can augment materials informatics workflows. XAI can be used as a forensic analysis technique to understand the consequences of data, model, and application decisions, or as a model refinement method capable of distinguishing important information [1, 2]. This approach is often used to explain how the structural characteristics of materials (features) contribute to a target property prediction using tools such as feature importance rankings that highlight useful or nuisance variables. However, an alternative approach is to apply similar methods to the instance space and identify influential or unproductive data instances (materials). Data sets

contain a range of special cases such as outliers (unusual types), archetypes (pure types), stereotypes (those assumed to be representative) and prototypes (those that actually are representative), and groups of data instances (clusters) that are similar in the high dimensional feature space. The amount of influence these special types of data instances have on a pattern, cluster or prediction is rarely explored or quantified, but they can also have a profound effect on model architecture and predictive ability.

Recent work has shown that it is possible to decompose the residuals of ML loss functions to better understand how individual materials contribute to model predictions [3]. This has been used to explain how including certain materials in a data set can improve the ability to accurately predict the properties of others [4]. This research has now been expanded to explain unsupervised patterns in data and identify special subsets of materials worthy of detailed consideration [5]. The first step is to represent materials using Shapley values, which are a solution in cooperative game theory where each game is assigned a unique distribution of a total surplus generated by the coalition of all players [1]. A popular tool for studying cost-sharing, market analytics and voting, in materials informatics the game is usually the model, and the players are the materials. By testing the impact of removing individual instances or features, and aggregating across the feature space or instance space, respectively, Shapley values quantify how much the inclusion or exclusion of a particular material (or a structural feature) affects the result. The second step is to transform the data, represented by its Shapley values, in different ways to reveal hidden groups or patterns. This two-step process aids the data analysis process, and acts as a precursor to the residual decomposition; simultaneously finding influential materials in the data set and quantifying how they are impacting the prediction of other materials.

The novelty in this new model-agnostic approach is that the cooperative game is the underlying data distribution, not a model, which opens up the opportunity for explainable unsupervised learning. This enables researchers to better understand how ML methods use the latent information captured in the data, informing better decisions about what kind of materials to make or simulate, what kind of characterisation or analysis to perform, and how these choices impact the outcome.

3. ML for autonomous microscopy: from physics discovery to atomic fabrication

Electron and scanning probe microscopies are now one of the foundational methods for characterization of structure and functional properties of matter on the nanometre and atomic scales. Scanning probe microscopy (SPM) enables rapid characterization of surface topography and mechanical, magnetic, ferroelectric, and electrochemical properties. Electron microscopy now provides comprehensive probe of structure, chemical composition, and vibrational properties at nanometre and atomic scales.

For most domain areas, microscopies traditionally represent downstream characterization methods in materials discovery cycle yielding the qualitative data. Recent progress in quantitative SPMs and scanning transmission electron microscopy (STEM) is challenging this paradigm, delivering large volumes of quantitative structural and high-veracity property data. However, the sheer volume of data has necessitated very complex analyses, minimizing the impact. The recent progress in ML and rapid data analytics for post acquisition analyses and particularly active learning methods that can be operationalized on active microscopes offer to change this paradigm [6]. On the data analytic side, ML provides the flexibility and speed necessary to analyse large volumes of multidimensional imaging and spectroscopy data for building low-dimensional representations and, in many, cases extraction of relevant materials parameters.

A fundamentally new spectrum of opportunities emerges in the context of active learning, where ML based workflows not only inform human-based decision making, but directly return control commands to the instrument. Operationalized on the SPM and STEM machines, these methods can be used for rapid mapping of the structure–property relationships. This knowledge can further be used for the discovery of generative physical models such as microstructure evolution equations, free energy functionals and Hamiltonians, and learning processing mechanisms. By combining zero-shot [7] and predictive [8] ML models with *in situ* particle beam, heating, or other processing, it is possible to learn materials responses and impart desired metastable states. These models are especially well suited to discovery scenarios, where they can reveal latent features to scientists, informing synthesis or degradation mechanisms.

These approaches create new opportunities for materials discovery. The last 20 years have seen exponential growth of the theoretical predictive capability for crystalline materials and small molecules. The last 5 years have seen the exponential growth in the capability to accelerate materials synthesis via laboratory robotics and microfluidic synthesis. However, the lesson of the past two decades is that scaling computation or synthesis individually by many orders of magnitude is insufficient to expedite materials discovery. Rather, the key is accelerating the feedback loop between theory and hypothesis making, experiment planning, synthesis, and characterization with subsequent update of theoretical models. Currently, characterization is the bottleneck—while synthesis can be scaled to 1000s compositions per day, the sequential structural, functional, and chemical probing outside of fast optical/photoluminescent methods still require hours and

days. Closing these characterization loops requires scaling down the probing volume and reducing measurement times, tasks ideally matched to microscopy capabilities. Here, microscopy offers the natural tool for exploration of multidimensional composition and processing spaces via strong (i.e. matching target macroscopic functionalities) and weak proxies [9]. There is also an opportunity to leverage ML-based adaptive sampling and intelligently select modalities based on uncertainty metrics, shortcircuiting the time to discovery.

A fundamentally new space of opportunities for materials discovery emerges based on controlled interventions in microscopy. In SPM, these include local polarization switching and electrochemical reactions that can now be studied at the time- and length scale well outside of conventional characterization methods, but very close to the intrinsic length scales of these phenomena. For electron microscopy, unique opportunities are the result of the electron beam's power to break local chemical bonds, enabling controlled fabrication of atomic defects [8], beam controlled atomic motion, and building homo- and heteroatomic artificial molecules atom by atom [10]. The rapid exploration of materials synthesis and degradation pathways at spatial, chemical, and temporal scales commensurate with fundamental physical interactions is now more viable than ever before.

Incorporation of ML methods both in real time and post-acquisition data analysis offers the compelling case to greatly increase the efficiency of instrument utilization by orders of magnitude and close the materials characterization gap, ushering the new era of materials and physics discovery and atomic fabrication.

4. Beyond crystallinity and throughput: AI for working interfaces in energy conversion technologies

The urgency with which mankind needs to accomplish the transition to a sustainable energy economy dictates a drastic acceleration of established research and development cycles toward ever improved energy conversion devices like solar cells, catalysts, electrolyzers or batteries. With respect to materials discovery much prospect to this end is seen in data centric approaches, which harness the powerful algorithms of ML or AI. In many areas of materials science, corresponding techniques ranging from high-throughput screening to inverse design are already most successfully employed to search the vast materials spaces for promising candidates at unprecedented efficiency [11]. The discovery is thereby often conducted entirely *in silico*, exploiting the predictive quality of first-principles computational data.

Unfortunately, such developments are at present still largely stalled for the energy conversion context as the functionality of corresponding devices is generally limited by interfacial problems and interfacial data is much more difficult to come by than the bulk data that suffices for many other application areas. This relates to the involved (experimental or computational) costs for the generation of such data, but even more so to the lack of best practice protocols to do this reliably and reproducibly. A crucial component here is the strong structural, compositional and morphological evolution that the interfaces in energy conversion devices undergo *operando* [12]. These working surfaces or interfaces are thus anything but simple truncations or ideal junctions of known bulk materials, respectively. Instead, they extend over a finite width, and exhibit novel purely interface-stabilized phases with often a low degree of crystalline order. For the experimental data generation, this *operando* evolution dictates not only stringent protocols for the initial synthesis, but a seamless and exhaustive documentation of the entire history of environmental operation conditions to which the interfaces were subject to. As this is rarely reached, data is not comparable and interoperable, preventing a community-wide build-up of large-scale data bases. At the same time, the *operando* evolution also precludes the generation of pertinent first-principles data. There are essentially no established *operando*-aware descriptors, and even if there were, there are in general no established structural models for working interfaces that could be used to compute them.

In this situation, there are two major strands in which AI and ML is presently employed toward an accelerated discovery. On the computational side, data-centric approaches are used to gain a deeper mechanistic understanding into working interfaces, with the long-term goal to use this insight to formulate *operando*-aware descriptors that could then be used for an efficient exploration of materials spaces [13]. A dominant development to this end is ML surrogate models, and there in particular ML interatomic potentials (MLIPs), which allow to generate first-principles quality data at orders of magnitude reduced computational costs. Appropriate for the data-scarce regime, the MLIP training is thereby done in agile active learning loops, with automated workflows being developed that ideally fully interlace this with the actual simulations to ensure consistent reliability [14]. In cutting-edge studies, the unprecedented capabilities are presently used to conduct the simulations in much larger simulation cells (therefore also allowing to address disorder) or perform significantly increased and therewith powerful samplings.

On the experimental side, AI and ML is increasingly employed to reach a deeper analysis of (*operando*) characterization data, either to also reach an improved mechanistic understanding or to identify structure

and correlations in the data that would enable improved workflows (proxy experiments, multi-fidelity experiments, etc). Notably, AI and ML are employed within emerging self-driving laboratories (SDLs). Here they complement lab automation and robotics to reach higher throughputs, but foremost they take over the experiment planning. SDLs operate in active-learning loops, in which data from executed experiments is fed back into the ML model to refine it and design subsequent experiments. Current methodological frontiers in employed Bayesian optimization or adaptive Design of Experiment approaches concern significant or varying noise levels (e.g. in case of multi-fidelity measurements), the design of larger numbers of data points (to meet batch-type operation in increasingly parallelized workflows), or agility to either autonomously adapt the shape and dimensions of the search spaces across loops or react to corresponding changes imposed by human scientists [15].

5. Accessing photoinduced reaction dynamics on surfaces with neural networks (NNs)

Laser-induced reactions at surfaces are particularly interesting because this kind of excitation mechanism can increase significantly the reaction probability with respect to ordinary thermal activation and, importantly, even open new reaction channels. Still, despite the impressive technical advances, experiments alone cannot fully determine with atomistic space and time resolution all the elementary steps involved in the reaction as well as the properties determining each of these steps. It is at this point that molecular-dynamics (MD) simulations become crucial to dissect the reaction dynamics.

Modelling the ultrafast photo-induced dynamics and reactivity of adsorbates on metals requires including the effect of the laser-excited electrons and also the effect of the concomitantly excited surface lattice. All these features can be effectively achieved by solving Langevin equations of motion, in which the coupling of each adsorbate and surface atom nuclei to the excited electrons is modelled in terms of electronic friction and stochastic forces that depend on the time-dependent electronic temperature that characterizes the excited Fermi-Dirac distribution, while the rest of interactions are described with the adiabatic potential energy surface (PES) that must account for all the system degrees of freedom. In spite of the apparent simplicity of the model, such simulations are highly demanding. Low energy molecules/atoms are particularly sensitive to energy differences of tens of meV that they experience in the proximity of a solid surface. And this sensitivity is even amplified when measuring, for instance, photoinduced desorption and reactivity probabilities and final-state distributions of the scattered gas species, such as, kinetic energies, scattering angles, and rovibrational quantum state distributions to cite some. Thus, any reliable description of gas/surface interactions requires the knowledge of the corresponding accurate first-principles multidimensional PES. First-principles molecular dynamics (FPMD) with electronic friction and thermostats (T_e, T_1)-FPMDEF, which calculate on-the-fly the adiabatic forces with density-functional theory (DFT) [16, 17], do enable such a complex modelling [18], but, unfortunately, these simulations come with a very large computational expense that severely limits any statistical analysis of the reaction and, also, it restricts the simulation time to just a few picoseconds that might be insufficient to guarantee well-converged reaction yields.

In the last years, the use of NN-generated multidimensional PESs and, in particular, the use of atomistic NNs (AtNNs), is becoming the accurate alternative to FPMD studies of diverse gas-surface reactions [19, 20]. Aside AtNN methods, the newer message passing NN potentials, which are also discussed in forthcoming sections, are certainly promising in terms of accuracy and efficiency when using the capabilities of graphical processing units (GPUs). However, it must be emphasized that the requirements imposed to a NN-PES capable of describing photoinduced reactions are even more demanding than those required in usual elementary gas-surface processes. A reliable NN-PES must be able to model large and out-of-phase movements of multiple and different adsorbates and also surface atoms and it must describe accurately the very distinct and changing adsorbate coverages that may exist during the photoinduced dynamics because of desorption events. This means that it is necessary to assure a precise description of both adsorbate–substrate and inter adsorbate interactions under very different and changing conditions, including local variations in the number of neighbour adsorbates and strong lattice distortions, since the lattice temperature may vary rapidly in the range of tens to thousands Kelvin.

The AtNN-like embedded atom NN (EANN) method, which uses descriptors inspired in the embedded-atom electron density, demonstrates to be impressively accurate and flexible to account for all these requirements [21]. EANN PESs allowed us to reproduce and understand the experimental strong coverage dependence of CO phodesorption in Pd(111) [22], the large branching ratio between CO photo-desorption and CO photo-oxidation in Ru(0001) [23], and reveal the dynamical nature of the CO physisorption well that so far was only found in XPS experiments. But there are additional open questions that we can now think in treating by exploiting NN capabilities. Besides the general challenges faced in gas/surface dynamics [19, 20], specific challenges for photoinduced surface chemistry are related to

performing nonadiabatic dynamics, either by advancing in orbital-based electronic friction coefficients adapted to the highly dynamic surface or even more challenging by developing excited state NN-PESs that could contribute to clarify the role on the initial nonthermal distribution of excited states.

6. Data-driven advances in modelling and understanding amorphous materials

ML has transformed atomic-scale materials modelling: rather than building simplified models of reality, we can now describe ‘the real thing’ in increasingly accurate simulations [24]. MLIPs are trained to reproduce quantum-mechanical energy and force data for this very purpose. In the domain of inorganic materials, MLIPs are typically based on DFT ground-truth data; once they have been fitted and properly validated, they therefore enable very-large-scale MD simulations of bulk and nanostructured materials, all while retaining DFT-like accuracy. MLIPs have evolved from specialised tools to increasingly widely available (and visibly popular) simulation methods, and their development has been documented in numerous review articles [19, 25]. The architectures used to train MLIPs have been advanced over many years and, thanks to these efforts, have now reached impressive accuracy. There are still many future research directions: among them are improved strategies for dataset construction, and for MLIPs that can be distilled for downstream tasks [26].

Looking from the development of MLIPs onwards to their (current and expected) impact on materials chemistry research, MLIPs are particularly promising tools in the area of amorphous materials—non-crystalline solids, whose complex atomic structures and structure–property relationships are now increasingly exploited for practical applications. Indeed, amorphous materials are of growing interest for energy storage, computing, catalysis, and many other fields (see [27] and references therein). Accordingly, amorphous materials are a frontier research challenge in computational materials design, and MLIPs are well placed to help to address this challenge [27].

A recent ML-driven study of graphene oxide (GO) exemplifies several aspects related to MLIPs and their applications to disordered and amorphous materials [28]. Formally, GO is a sheet of graphene modified by the presence of various functional groups (say, hydroxyl, carbonyl, and so on). In laboratory experiments, this modification is achieved using chemical reactions; in simulations, one can now quickly construct atomistic structural models over a wide-ranging parameter space of compositions and functional groups, and yet only the subsequent comparison with experiment will ultimately validate a given structural model. The study in [28] takes a two-step approach: first, exploring structures with ML-accelerated FPMD; second, using a graph-neural-network (GNN) architecture for fitting increasingly accurate MLIPs that iteratively ‘learn’ about 2D extended and subsequently about 1D edge structures—details and methodological references may be found in [28]. With the final MLIP model available, MD simulations were carried out, exploring the gradual thermal reduction of a GO sheet.

ML-driven simulations have already begun to have a major impact in materials chemistry and related fields. In the future, together with other emerging AI/ML approaches, they might enable the discovery and design of amorphous functional materials for a variety of practical applications [27].

7. Foundation models for atomistic materials chemistry

DFT and its associated methods have become the standard toolkit of computational materials science and also to a large extent computational chemistry. As such, DFT constitutes the pinnacle of the *Dirac programme* of first-principles modelling [29]: start with the fundamental equations of quantum mechanics that describe the electrons and atomic nuclei (the latter represented as point charges to an exceedingly good approximation), and derive the consequences for the behaviour of crystals, molecules, currents, etc. The resulting sequence of approximations over the past ~ 50 years have enabled the description of known—and prediction of new material properties and underpins our understanding of the material world at the microscopic scale.

The computational cost and scaling of DFT in practice limits its general usefulness to the treatment of hundreds of atoms and picoseconds of time scale. While these limitations are being challenged and push all the time by the progress in computational hardware and also algorithmic efficiency, but the extension of first-principles modelling to significantly larger length scales require a change in the modelling framework. Just as DFT eliminates the degrees of freedom inherent in the full many-body wave function and just retains the one-particle operator corresponding to the electron density, we can go further and eliminate electronic degrees of freedom altogether by writing the total energy as a function of just the atomic coordinates: a force field. This function is very complicated, but advances in parametrising functions using a very large number (typically millions) of parameters based on fits to a large amount of data (widely known as *machine learning*) has enabled useful approximations that allow the simulation of tens of thousands of atoms for millions of time steps, i.e. *nanometres* of material for *nanoseconds*.

The past decade or so has seen incredible progress, and was spent mostly understanding how to build datasets for fitting ML force fields for particular systems, and how to achieve the accuracy that is required for the model to be usefully predictive of interesting properties [25]. Indeed even just characterising the relationship between the ‘pointwise’ accuracy of the potential energy of a force field to the error in its prediction of any particular material property is highly nontrivial, and turns out to be critical for success. Simulations of phase transitions both under equilibrium and nonequilibrium conditions, heterogeneous catalysis, study of diffusion and spatiotemporal correlation are now routinely possible for complex materials.

Just very recently it was discovered that when the training set is diverse enough, a force field can be made that covers most of periodic table, and despite only having been fitted to DFT calculations of small inorganic periodic crystals, is capable of running stable molecular dynamics on essentially any chemical system [30]. Such extreme generalisation goes somewhat counter to the conventional wisdom in ML research, which has made tremendous progress recently by using ever larger data sets. There is currently little understanding of what gives rise to such generalisation, but it raises the tantalising possibility of a *universal* force field. There is no doubt that further accuracy for a wide range of systems will be gained by training on large databases, and the construction of numerically consistent DFT data is currently the limiting factor.

8. ML electrochemistry

The accurate description of redox reactions from the perspective of first-principles calculations still represents a challenge. Standard DFT approximations to the exchange-correlation functionals suffers from the so-called self interaction errors (SIEs), leading to an unphysical delocalization of electrons and thereby limiting its ability to accurately study processes where changes in oxidation states are critical. Hybrid functionals, and even more extended Hubbard functionals (DFT + U + V) can provide a successful solution to this challenge [31, 32]. As recently shown, DFT + U + V provides a robust framework to mitigate SIE in materials with strongly localized *d* or *f* electrons, especially for systems where the electronic localization occurs with substantial hybridization. Recently, it has been shown how the use of DFT + U + V along with FPMD is capable of following the adiabatic evolution of oxidation states over time in representative cathode materials for Li-ion batteries [32]. In addition, this opens the door to incorporating the concept of redox-aware into ML potentials. Starting from the physical rationale that atoms with different oxidation states behave like distinct species, it has been shown that a NN training that considers atoms with different oxidation states (obtained through DFT + U + V FPMD) as distinct species can identify the correct ground state and pattern of oxidation states for the redox elements present [32]. This can be achieved, e.g. through a combinatorial search for the lowest-energy configuration, among all possible patterns, and is shown to recover correctly the DFT + U + V ground state. This brings the advantages of ML potentials to central technological applications (e.g., rechargeable batteries), which require the correct description of redox states.

The predictive accuracy of DFT + U + V heavily depends on the precise determination of the onsite U and inter-site V Hubbard parameters, which describe localization and hybridization, respectively. While in the simplest cases these parameters could be obtained through semiempirical tuning (but then negating the predictive power of the approach, and the capability to deal with complex and very diverse local environments, that require atom-specific U and V), unbiased predictions identify Hubbard parameters self-consistently through linear response calculations, particularly efficient when density-functional perturbation theory (DFPT) is deployed. This approach has now been fully automated [33], enabling high throughput calculations of Hubbard parameters that can provide extensive datasets for further investigations.

In particular, it becomes even possible to build a ML model to predict U and V bypassing the DFPT step. For example, the ML method of [34] has been recently devised to this goal. The model is based on equivariant NNs, and uses electronic occupation matrices as descriptors, capturing the electronic structure, local chemical environment, and oxidation states of the system in question. The model significantly speeds up the prediction of Hubbard parameters, while approaching the accuracy of DFPT. The model uses two DFT-based calculations: first, a DFT + U + V ground-state calculation with initial guesses for U and V (which can be set to zero) to obtain atomic occupation matrices; second, a structural optimization using the model-predicted self-consistent (SC) Hubbard parameters to obtain the SC structural-electronic ground state. Furthermore, thanks to its strong transferability, it enables accelerated materials discovery and design via high-throughput calculations, with relevance for various technological applications.

Another key topic in computational electrochemistry is the accurate calculation of molecular ion solvation energies, crucial for controlling electrochemical reactions. In particular, this information is essential for the characterization of relative phase stability in different environments, and thus of major interest to advanced materials and manufacturing. First and foremost, first-principles accuracy is needed to determine the solvation energies of ions and small molecules in arbitrary solvents. The NN potentials discussed in the previous sections make these calculations viable, and overcome the computational bottlenecks of FPMD. A

recently developed NN-based workflow [35] has shown the capability to compute ion solvation energies for alkaline(-earth) cations with chemical accuracy. Future directions will involve developing active learning schemes to automate the calculations' workflows. Moreover, electrostatic interactions that have been treated directly through the NN for the short range and analytically for the long range might need to take into account the complex nature of the electrochemical potential across all length scales.

9. ML for molecular sensing

ML is fundamentally transforming molecular sensing, particularly in gas sensing field, by revolutionizing the screening of sensing materials and the enhancement of sensor performance through advanced signal processing techniques. By integrating ML with theoretical tools (e.g., DFT), researchers have unlocked a powerful methodology for designing selective gas sensing materials and decoding complex sensor signals. This synergistic approach accelerates material discovery and sensor optimization, paving the way for molecular sensing devices that are more sensitive, selective, and reliable.

Traditional methods for designing and screening gas sensing materials rely on trial-and-error experimentation, which is often labour-intensive and time-consuming. By contrast, ML, when combined with computational tools, facilitates the efficient prediction of material properties, significantly streamlining the process. For instance, ML models can correlate key material descriptors, such as adsorption energy, surface reactivity and electronic properties, with their responses to specific gases. This enables rapid screening and selection of materials without the need for exhaustive experimental validation. For instance, in a recent study, ML combined with DFT successfully predicted the sensitivity of $\text{Cs}_3\text{Cu}_2\text{I}_5$ to hydrogen sulfide, achieving a remarkable 92 % accuracy in predictions, which were later validated experimentally [36]. Beyond accelerating material discovery, this integration also provides valuable mechanism insights into gas adsorption and sensitivity, enabling a deeper understanding of the materials' functionality. Similarly, for metal oxide materials-based sensors, ML models have been instrumental in identifying critical descriptors that dictate their sensing capabilities, guiding the targeted selection of materials for diverse applications ranging from industrial safety to environmental monitoring [37].

Following the selection of sensing materials, ML continues to play a pivotal role in optimizing sensor performance by fine tuning critical parameters such as sensitivity, selectivity, and response time. It is reported that ML techniques have been applied to gas-sensing platforms based on copper phthalocyanine functionalized graphene, enhancing their ability to detect trace amounts of gases like ammonia and phosphine [38]. By analysing the sensor's responses, ML improves accuracy and specificity, even in complex gas mixtures. Furthermore, ML has proven invaluable in the design of sensor arrays capable of detecting multiple gases simultaneously. Algorithms are employed to analyse interactions between sensor elements and their collective responses, enabling the identification of the most effective configurations. This optimization is particularly critical for applications like air quality monitoring, where the simultaneous and accurate identification of various pollutants is essential.

In molecular sensing, interpreting gas sensing signals is crucial. ML techniques are extensively utilized in signal processing to extract meaningful features from raw sensor data while minimizing noise, a common challenge in real-world sensing environments. For instance, ML has been used in electronic noses to extract transient kinetic features from sensor response profiles. These features act as distinct fingerprints of odorants, enabling the accurate classification of volatile organic compounds and addressing one of the field's primary challenges [39].

Despite its transformative potential, the application of ML in molecular sensing faces several challenges, including improving the interpretability of ML models, reducing dependence on large datasets, and enhancing the real-time performance of sensing systems, as well as energy-consuming. Overcoming these hurdles will require continued advancements in ML techniques, such as the integration of deep learning and reinforcement learning, development of more accurate adaptive sensing systems, as well as development of brain-inspired neuromorphic computing system [40]. Future developments could enable gas sensors that not only detect and classify gases but also predict environmental changes or potential hazards. Such advancements will pave the way for smart, autonomous sensing systems across diverse domains, including healthcare, environmental monitoring, and industrial safety.

To sum up, the integration of ML with theoretical tools is revolutionizing the design and optimization of molecular sensors. By expediting material discovery, refining sensor configurations, and enhancing signal processing, ML stands at the forefront of developing next-generation molecular sensing technologies. These innovations promise enhanced sensitivity, selectivity, and performance, ensuring their pivotal role in addressing the challenges of modern sensing applications.

10. Refining molecular characterization for robust machine-guided corrosion inhibitor discovery

A particularly urgent and industrially significant case where ML is being applied to discover effective molecular materials is in the discovery of corrosion inhibitors. Such inhibitors would be embedded in the primer of a paint system or used in initial metal passivation. Traditionally inhibitors have been based on chromate or other toxic compounds that are being banned by legislation worldwide. Small hetero-cyclic compounds are a promising alternative, yet to determine the exact molecular structure with high-efficient corrosion inhibition from tens of thousands of possibilities remains a challenge. Various methods including high-throughput experimentation and computational modelling have been developed to select or design the optimal molecular structure. A very promising approach is the use of inverse design, in which high throughput experiments defining electrochemical performance and computation methods deriving inhibitor characteristics and attributes are linked by a ML method to define the molecular attributes essentially important for inhibition performance. These critical features are then used to search molecular databases and select promising candidate inhibitors that are then subject to testing for verification.

Early work was able to use quantitative structural activity relationships (QSARs) methods based on a NN approach to obtain reasonable models of the features controlling inhibition [41]. However, while these models represent successfully the existing dataset, further generalization ability was still to be enhanced. Challenges may come in twofold: (1) the relevance of molecular characterization and attribute definitions; and (2) the deficiency of computationally/experimentally generated datasets. Since then, large programs have been undertaken, where the databases were significantly enhanced, and great care was cast to ensure that both experimental and computational data were accurate and reproducible. Further the molecular attributes were refined to better represent molecular interactions with solvent and metal surfaces. Ranges of statistical and QSAR techniques have been used to define the relationships between the molecular attributes and electrochemical performance. These models demonstrated an enhanced ability to represent existing data, but their predictive ability could still be enhanced [42]. Recent experimental work reveals the complexity of corrosion inhibition process, of which peak performance (both electrochemically and mechanically) was reached by short-term inhibitor treatment, but subsequent voids appeared within the inhibitor film with time expansion [43]. MD models indicated that the inhibitor film may be subject to electroporation where charge at the metal surface causes the inhibitors to clump together allowing water to again reach the metal surface [44].

Thus, it is evident that additional factors involved in the inhibitor adsorption control the overall performance and stability of inhibitor layer, entailing a deeper understanding of inhibitor layer formation and lifetime. A recent review [45] highlighted the limitations of previous models: inadequate or no representation of solvent, lack of potential effects and relatively small models. New methodology was proposed based on a combined quantum mechanics/molecular mechanics/non equilibrium greens function approach. This approach enables the simulation of larger models that include both solvent and voltage effects. The system has been applied to the inhibition study of both copper and zinc surfaces by 2-mercaptobenzimidazole (MBI). A major result of the study is that, when MBI binds to the surface, a major electronic re-alignment across the inhibitor assembly rearranges the dipole moment at the exterior of the molecule. The traditional theory of inhibitors is that they form a barrier to both water and solvents against charge transfer. This study proposed that while MBI acts as an effective barrier against water, it cannot be regarded as a charge barrier. In fact, charge realignment and the formation of the dipole will have a profound influence on the deposition of the subsequent inhibitor layer. The relevant molecular attributes contributing to the dynamics of corrosion inhibition processes are potentially important descriptors that were previously overlooked for ML development. As highlighted in our studies, the molecular attributes that can represent these processes are quite different from those that reflect surface bonding and thus our datasets used in ML approaches need significant redesign.

11. Exploring new frontiers in inverse materials design through GNNs and large language models (LLMs)

Finding new materials with suitable properties has been a challenging task due to the computational and experimental costs. AI/ML techniques have been successfully used for both forward (structure to property) and inverse (property to structure) tasks in materials design [46]. Inverse design approaches can surpass traditional funnel-like materials screening methods and facilitate the computational discovery of next-generation materials. Since no explicitly available physics-based methods exist for inverse design tasks, AI/ML is an obvious choice.

To accomplish inverse materials design tasks, we require the following: (1) a well-curated and diverse dataset, (2) an AI/ML model and architecture that can establish a mapping between the properties of materials and material structures, and (3) suitable metrics and a benchmarking strategy to guide the design process. While there are numerous material properties—such as electronic bandgap, bulk modulus, refractive index, etc, or their combinations—that can be used as target properties, we can start with a specific property, such as superconducting transition temperature (T_c). Superconductors are one of the most celebrated classes of materials in materials science, but there are very few such materials known experimentally. As mentioned above, we require a dataset for superconductors. While many materials databases exist, they lacked superconducting properties until JARVIS-DFT.

JARVIS-DFT [47] consists of more than 80 000 materials and millions of material properties, with around 1000 superconducting materials in the dataset. Note that predicting T_c is computationally expensive compared to other properties, such as formation energy, when using DFT. As the next step for inverse design, we require AI/ML methods suitable for this task. There are a variety of AI/ML methods, such as fingerprint-based traditional methods, deep learning techniques like convolutional NNs, GNNs, and generative pre-trained transformers (GPTs).

GNNs, in particular, have been successful recently for atomistic materials design tasks. In these models, atoms are represented as nodes, bonds as edges, and angles as edges of the corresponding line graphs, for instance. GNNs such as atomistic line GNNs, combined with diffusion models like the crystal diffusion variational autoencoder, have enabled the generation of superconducting atomic structures [48]. The dataset was split into training and testing sets, and the metric for performance was the interatomic distances between target and predicted structures in the test dataset. After model development, more than 50 candidate superconductors were computationally discovered and later characterized with DFT to validate AI predictions. Another AI approach used was GPT models.

In GPT models such as AtomGPT, both the atomic structure and the target property can be represented as text [49]. These texts are converted into tokens, and GPT models establish the relationship between the atomic structure and property/prompt tokens. Such GPT models have shown remarkable promise for both forward and inverse materials design tasks. For the superconducting dataset, we followed similar train-test splits as in GNN methods and measured performance based on the interatomic bond distance comparison metric between target and predicted materials in the test dataset. We found that GPT-based models surpass GNN models in terms of this metric, and new candidate superconductors were computationally discovered and later validated with DFT. Additionally, GPT models are much faster and easier to implement than GNN models. These comparisons are hosted on the JARVIS-Leaderboard [50] open-source platform to enhance reproducibility, transparency, and allow others to contribute their models as well.

12. Property directed generative design of inorganic materials

Property-directed generative design presents a unique opportunity in modern materials discovery, shifting from large-scale data-driven screening to precise, generative approaches aimed at discovering novel compounds with tailored properties. Recent advancements in computational science for materials discovery have made significant progress in addressing one of the key challenges in crystal structure prediction (CSP)—identifying stable and metastable structures efficiently [51]. Traditional CSP methods, such as evolutionary algorithms and particle swarm optimization, however, are computationally expensive and limited in their ability to explore vast chemical spaces. Generative models, by contrast, provide a promising alternative by efficiently targeting structures that are near ground-state configurations when trained on existing data. These generative frameworks also enable inverse design, where the desired targets guide the generation of materials, making them particularly valuable for property-directed generative design.

However, a major challenge in applying generative models is in ensuring that generated structures obey the symmetry and periodicity essential for physical plausibility. Symmetry considerations are fundamental for determining key properties of inorganic materials, including electronic band structures, optical behaviour, and mechanical strength. We summarize some recent frameworks, such as DiffSCP++ [52], CrystalFormer [53], WyCryst [54], MatterGen [55], and physics guided crystal generative model (PGCGM) [56], which have demonstrated the importance of integrating symmetry constraints into a generative framework to ensure the physical plausibility of generated materials. These models utilize a variety of AI-driven models, such as symmetry-based representations, diffusion-based methods, and GNNs, to generate stable and diverse crystal structures that satisfy specific property requirements. By embedding symmetry into the generative process, these frameworks enhance the efficiency of materials discovery, reduce reliance on trial-and-error experimentation, and open new avenues for the design of materials with applications in energy, electronics, and catalysis.

Among these models, WyCryst enables symmetry-constrained structure generation through three key components: a Wyckoff position-based representation to enforce symmetry constraints, a property-directed variational autoencoder for generating novel crystal structures, and an automated DFT workflow for validating the stability and properties of the generated materials. By embedding symmetry constraints, WyCryst efficiently generates materials that adhere to space group symmetries while meeting desired property criteria. Similarly, DiffSCP++ employs a symmetry-constrained diffusion model to refine atom types, positions, and lattice parameters, ensuring that the generated structures maintain realistic symmetry and periodicity. This approach enhances the diversity and stability of generated materials, opening new possibilities for discovering synthesizable inorganic compounds. CrystalFormer employs a transformer-based architecture to generate crystal structures by predicting symmetry-inequivalent Wyckoff positions in the unit cell ensuring compliance with space group symmetries: this produces thermodynamically stable materials with various symmetries. CrystalFormer is also capable of performing property guided exploration with probabilistic modelling, facilitating the discovery of inorganic compounds with targeted properties. MatterGen employs an SE(3) equivariant diffusion approach to generate crystal structures by iteratively refining random initial configurations until they conform to a targeted distribution. MatterGen, as a base pre-trained model, can be finetuned towards stability or functional properties, facilitating the discovery of materials tailored to specific applications. The PGCGM achieves symmetry-based generative design by incorporating physics-oriented losses related to physics and space group symmetry. The model training emphasizes thermodynamic stability ensuring the generation of low energy compounds. PGCGM also enables the generation of crystal structures with specific space group symmetries, allowing further discovery of functional materials.

In conclusion, property-directed generative design frameworks represent a significant advancement in the field of materials science. By embedding symmetry-based constraints into the generative process, these models enhance the validity and stability of predicted materials, thereby accelerating the discovery of inorganic compounds with desired properties. A key bottleneck is the generation of experimental or high-quality computational data to train such generative models. This approach, however, promises not only a streamlined materials design process but also new avenues for the development of advanced materials tailored for specific applications. The synergy between AI models and physics-based property-directed design holds immense promise for revolutionizing the way materials are discovered and optimized for real-world use.

13. Physics based ML for materials and compound space

The virtual navigation of chemical compound space has been significantly constrained by the prohibitive computational demand associated with numerically solving approximations to Schrödinger's equation with satisfying accuracy for an exponentially growing number of possible systems. Over the last decade, considerable progress has been realized thanks to the application of statistical techniques commonly referred to as AI, as recently documented in an entire issue in *Chemical Reviews* dedicated to ML at the atomic scale [57]. Due to the colossal number of potential and costly training compounds, the central inquiry has been on how to improve training efficiency—as quantified by scaling laws (or learning curves). This question has persisted ever since it was first demonstrated that ML models of quantum properties can be applied throughout chemical compound space, i.e. for out-of-sample systems (not part of training), with prediction errors that decay systematically with training set size [58]. Subsequent applications have highlighted the promise of ML for the atomistic sciences by systematically surpassing the accuracy of hybrid DFT approximations for various quantum properties [59], estimating formation energies for millions of quaternary crystals [60] or reaching the accuracy of explicitly correlated electronic structure theory methods through Δ learning [61], or multi-level learning [62]. Further breakthroughs in training efficiency, scalability, and transferability were achieved by virtue of similarity-based query aware models, trained on the fly, and decomposition of training and testing systems into fragments, based on atoms-in-molecule-ONs (AMONs) [63].

Most recent contributions indicate that meaningful combinations of these techniques are possible, often via intricate combinations with DFT. As such, DFT has assumed an outstanding role for the use of AI in chemistry and materials not only for merely generating data sets for training and testing but also for informing superior ML model architectures and workflows [64]. Specific examples include the combination of Δ learning and AMONs to enable quantum Monte Carlo level of accuracy [65], similarity-based learning and ridge regression identifying potentially superconducting candidates [66], or adaptive hybrid DFT which reaches superior accuracies when it comes to singlet-triplet spin gaps or other quantum observables [67].

The work mentioned only represents a small glimpse of recent activities in the entire and rapidly growing field. Overall, remarkable progress has been made towards the generic goal of reaching EAST, i.e. Efficiency,

Accuracy, Scalability, and Transferability [64]. Remaining challenges include the generation of more and sufficient data that is universally representative not only for minima but also for barriers, foundational ML models that can be used to estimate any quantum mechanical observable in any electronic state, and the possibility to account for multi-reference, as well as nuclear quantum and relativistic effects.

14. Language models for many-body physics

Now is an exciting time for research on quantum physics due to the opportunities and significant advances in the application of ML and AI to fundamental problems in physics, chemistry, and materials science. In particular, the transformative power of language models like recurrent NNs (RNNs) and Transformers [68], originally designed for natural language processing (NLP), has opened a new frontier across a wide array of technologically and scientifically relevant disciplines, including classical and quantum many-body physics.

Although historically these models originally demonstrated breakthrough performances in NLP, such as in ChatGPT [68], they have in principle little to do with ‘language’ itself. From a broader perspective, these models constitute powerful statistical modelling and information processing machines that can process a wide array of data types exhibiting correlations of different nature not limited to language. Tokens traditionally understood as pieces of words or phrases could also represent physical or chemical degrees of freedom, namely, spins in a lattice, lattice occupation numbers, atomic coordinates, or generally any sequence of inputs that are statistically mutually dependent. By expanding the token universe to encompass states from any other degrees of freedom relevant to a physical system, LLMs can allow physicists to simulate many-body interactions with unprecedented precision and efficiency.

While the origin of these token streams is disparate, the statistical correlations in datasets commonly used in NLP, computer vision, and other popular tasks in ML, display striking similarities with data from physical systems. Key similarities include symmetries, high dimensionality, and correlation functions. For example, spatial symmetries are present in natural datasets and classical and quantum systems simultaneously improve the sample complexity and learnability of models in computer vision as well as enriches our understanding of physical systems in classical and quantum mechanics. The behaviour of the correlations among the constituent elements in the token streams in computer vision and NLP display strikingly analogous behaviour to classical and quantum systems in thermal equilibrium near a critical point [69]. These commonalities make it natural to attempt to use these models to study classical and quantum many-body systems and are an important reason behind their rise and success in quantum many-body physics research.

Recent research has begun to capitalize on this potential. Techniques such as RNN wave functions [70] and language model-based quantum state tomography offer flexible and powerful representations of quantum states than conventional approaches. These studies have been extended to the task of finding ground states of quantum many-body systems, e.g., ground states of frustrated magnets, Rydberg atoms arrays, and fermionic systems, as well as to simulate the time evolution of quantum states and to solve combinatorial optimization problems [71]. Quantum chemistry, a field crucial for understanding molecular interactions and reactions, has also benefited from these advances. Transformer-based models can predict molecular ground state energies with comparable accuracy to traditional methods. Such advancements hold immense promise for rapid simulations in quantum chemistry, offering a pathway toward scalable tools that handle large basis sets and dense electron correlations that are challenging for standard quantum chemistry methods.

One potential application of LLMs in physics is the simulation of many-body fermion systems, such as those encountered in Rydberg arrays, exotic material phases, or molecules. Models such as ‘RydbergGPT’ [72] could enable simulations, potentially influencing quantum computing and materials science in the long term. By offering a scalable and adaptable approach to many-body physics, LLMs could present a complementary method to state-of-the-art algorithms such as quantum Monte Carlo and density matrix renormalization group, especially when dealing with high-dimensional frustrated or out-of-equilibrium systems.

Looking ahead, the development of efficient and environmentally conscious LLMs is critical. The computational costs of training these models using GPUs are significant, and reducing the environmental impact of large-scale simulations remains a pressing concern even in physics simulations based on language models. Innovations in model architecture design could help address this issue, aligning with the broader push for sustainable computing. In conclusion, language models can become versatile tools for scientists by bridging the gap between language processing and physical simulations, impacting fields beyond NLP [71]. As research in this space advances, LLMs may catalyse breakthroughs across many-body physics, quantum chemistry, and beyond, unlocking a new era of data- and physics-driven, efficient, and scalable many-body systems simulations.

15. Variational autoencoders-enabled high-fidelity reconstruction and effective anomaly detection in time-series data

Robust modelling of multi-channel biological time-series data, such as EEG, across different individuals is crucial in numerous applications. Most often, identifying common patterns (*biomarkers*) is as relevant as distinguishing them from individual behaviours (*fingerprints*). However, achieving accurate modelling involves tackling three primary challenges: intersubject variability, intra-subject variability, and ensuring data quality and fairness, including the automatic detection of artifacts [73].

We used the well-known *BCI dataset 2a*, a very popular EEG dataset collected from nine subjects performing motor imagery of hand and feet movements, to test both classification and reconstruction using various deep learning models [74]. First, *vEEGNet-ver1* served as the baseline model upon which we built subsequent versions. *vEEGNet-ver1* is a variational autoencoder with the encoder inspired by the popular EEGNet architecture, with three main convolutional layers. The decoder is the mirrored version of the encoder. By enhancing its encoder architecture, we developed *vEEGNet-ver2*, which offered improved performance over the first version in terms of reconstruction. Then, we decided to focus on the reconstruction task, in line with previous literature indicating a trade-off between classification and reconstruction learning abilities [75]. This led to the creation of *vEEGNetver3*, which targets a single task, i.e., the reconstruction. In *vEEGNet-ver3*, we defined the reconstruction loss as the (soft) dynamic time warping distance between the original and the reconstructed time-series. This approach significantly improved the model's performance, suggesting the importance of concentrating on specific tasks to achieve better results. Finally, by employing a hierarchical variational autoencoder architecture [76], we transformed *vEEGNet-ver3* into the *hvEEGNet* model. This advanced architecture demonstrated high fidelity reconstruction performance and provides three distinct latent representations, extracted from the three latent spaces of the model. As the reconstruction performance on the *dataset 2a* were very high, we tested *hvEEGNet* as an automatic artifact detector, enabling the identification of artifacts that had not been previously detected in the wellknown public dataset.

One of the key insights from our work is the crucial role of domain knowledge that allowed us to recognize that poor reconstruction results were linked to acquisition problems, such as signal saturation, or physiological artefacts, such as eye blinking. Ensuring high data quality is essential for the successful and reliable learning of ML models. Without high quality data, even the most advanced algorithms can produce misleading or suboptimal results.

Moreover, the latent representations extracted by *hvEEGNet* can be further investigated to develop new physics-informed smaller and more effective latent space structures [77]. Such advancements could pave the way for more robust and informative deep learning models for time-series modelling and anomaly detection. By improving the effectiveness and interpretability of latent representations, future research could address the challenge of distinguishing common patterns from individual ones and better quantify inter- and intrasubject variability. Also, improved interpretability will enable a higher degree of interaction with domain experts, who can help drive the development of deep learning models tailored to their research and clinical questions.

The significance of this research extends beyond this immediate application, as the above challenges are common to other domains where complex living systems are under investigation. Moreover, *hvEEGNet* is a versatile model which can be adjusted to other types time-series data, with different dynamics, and different applications.

16. Multiscale materials science: tasks, challenges, and cross-domain synergies

Materials science is a field driven by its multiscale nature, where phenomena at vastly different spatial and temporal scales interact to define the properties and behaviours of materials. From atomic vibrations that dictate thermal conductivity to macroscopic structures determining mechanical strength, understanding and predicting material behaviour requires bridging these scales.

Traditionally, physics has provided a robust set of mathematical tools to address multiscale problems. Methods such as renormalization groups, effective field theories, and closure coordinates have been used to study specific properties like critical points or ground states. While these approaches have been immensely successful in understanding phase transitions and other fundamental phenomena, they were typically designed to address narrowly focused problems.

Today, the scope of problems in materials science has expanded significantly. Researchers are not only interested in understanding ground states or critical phenomena but also in exploring broader challenges like finding meta-stable states, analysing mechanical properties under various conditions, and even generating entirely new structures.

Addressing these challenges requires a paradigm shift from traditional analytical methods to new data-driven approaches. ML and AI have emerged as powerful tools to augment classical methods, enabling scientists to model, predict, and design materials across multiple scales with unprecedented efficiency.

This shift toward data-driven methodologies is transforming materials science, creating opportunities to solve problems that were previously intractable and broadening the field's potential impact across domains.

16.1. Integrated multiscale tasks in materials science

In materials science, tasks such as property prediction, conditional structure generation, automated synthesis, and physical law discovery are inherently multiscale. Each of these tasks requires understanding the interplay between small-scale phenomena and large-scale outcomes.

Predicting material properties often involves connecting atomistic interactions to macroscopic behaviours. For example, multiscale deep learning models can predict the elastic properties of woven composites by analysing data from simulations at the microstructural level [78]. These models provide valuable insights into how small changes at the microscale influence the overall performance of a material.

Generating structures with specific properties is a complex inverse problem. Advanced generative models, like those used to predict domain boundaries in potassium sodium niobate thin films, can reveal previously unobserved structural motifs that emerge from simple local rules [79]. This work highlights how structural complexity arises naturally from underlying physical principles.

Automating synthesis processes accelerates material discovery by optimizing experimental conditions. For instance, the LeapFrog framework combines adaptive mesh refinement with ML to simulate the solidification of alloys, offering insights into how synthesis parameters affect microstructure formation [80].

Discovering physical laws and principles requires connecting diverse scales of phenomena. Compression theory, for example, identifies relevant degrees of freedom in complex systems like quasicrystals, uncovering new critical behaviours that were previously hidden [81, 82].

16.2. Universality of multiscale methods

One of the most exciting aspects of multiscale methods is their universality. Once developed for a specific domain, these methods can often be applied to entirely different fields. For example, techniques for coarse-graining molecular dynamics with GNNs reduce computational costs while generating transferable representations applicable across molecular systems [83]. Similarly, data-driven models used to study DNA methylation patterns have uncovered thermodynamic variables that govern healthspan and lifespan across species, demonstrating the potential for cross-domain applications [84].

By leveraging the universality of multiscale approaches, we can accelerate discoveries not only in materials science but also in fields like biology, chemistry, and even social systems.

16.3. The core challenge: balancing detail and holism

A central challenge in multiscale modelling is determining the appropriate level of detail. Too much detail can make models computationally infeasible, while too little can lead to inaccuracies. For example, the FE^{ANN} framework balances accuracy and computational efficiency by using physics-constrained NNs to model fibre-reinforced composites [85].

When a single level of detail is insufficient, integrating multiple scales into a cohesive framework becomes essential. Flow-matching, a novel method for coarse-grained molecular dynamics, combines generative modelling with force-matching to efficiently capture key interactions across scales [86]. Such approaches demonstrate how multiscale frameworks can provide a holistic view without overwhelming computational resources.

16.4. Outlook: toward a unified multiscale ecosystem

The future of materials science lies in creating a unified ecosystem that integrates multiscale simulations, AI, and experimental data. Graph-enhanced deep material networks, for example, unify the modelling of diverse microstructures, enabling predictions across families of materials [87]. These tools not only improve accuracy but also pave the way for entirely new material designs.

Furthermore, integrating theoretical principles with data-driven models offers powerful opportunities. For instance, a platform based on the Onsager principle creates reduced thermodynamic coordinates for stochastic systems, allowing for a more profound understanding of complex material behaviours [88].

By developing interoperable, scalable, and transferable tools, we can accelerate innovation, enabling faster discoveries and broader applications of multiscale methodologies.

Multiscale materials science stands at the intersection of computation, experimentation, and theory. From predicting properties to discovering universal principles, multiscale methods allow us to tackle challenges across domains. By balancing detail and holism and leveraging the universality of these

approaches, we can push the boundaries of what is possible, not only in materials science but in many other fields as well.

17. Conclusion and perspective

Digitalization of materials is a strategic action in the frame of emerging twin green & digital transition which aim at a more sustainable and resilient world economy. New digitalization solutions are needed covering the whole materials value chain and interconnecting all phases of the materials life cycle, from materials design and development, production, optimal usage, to maintenance, re-use, and recycling. Principally these efforts can be broadly categorized into two domains: 'digital twins' and materials models in a very broad sense, provided digital representations of materials in the context of their application independent of a specific measurement. Secondly, curated dataspaces provide access to experimental data that forms the basis for the generation of model-based digital representations.

Digital twins and materials models need to extrapolate beyond the limited number of available data points. Given the vast dimension of the materials space this problem will not be solved by high-throughput experiments. For this reason, one of the major challenges of developing accurate and functional modelling of IAMs is to reach the accuracy of first-principles approaches over a very large volume of systems. Moreover, a generic modelling procedure at an experimental scale where material imperfections such as defects, disordered chemical composition, rough interfaces, etc, play a major role, is hardly possible with conventional first-principles techniques. The development of ML strategies that allow obtaining atomistic models from a large dataset of small and accurate first-principles calculations could enable achieving unprecedented time and length scales. The elaboration of novel types of datasets required to train ML models is also of major concern, but while open-access material databases offer valuable information on thousands of crystalline materials, they overlook the nature and impact of the variety of possible atomic imperfections as usually observed in experiments. Finally, a substantial challenge persists in extracting meaningful physical insights from the vast amount of data (raw images and spectra) generated during experimental analysis. The use of AI-driven methodologies associated with (S)TEM data analysis for instance could boost the automation of experiments and data analysis of IAMs.

It is therefore urgent to develop AI and ML based models embedded into workflows that can accelerate the design of IAMs, and optimize their compositions and structures for enhancing their application performances. Efforts need to be focused on the generalization of AI-driven ML techniques to cope with realistic modelling of IAMs, as well as on the development of workflows connecting the generation of atomistic models to the simulation of their electronic, transport, thermal and optical properties. Critically, the impact of disorder, interface symmetries or chemical composition on their physical properties (electronic, optical, magnetic, etc), in limiting the use of IAMs for optimization of devices and achieving device metrics' upper limits need to be considered. Additionally, AI-enhanced characterization workflow should be developed to facilitate breakthroughs in data analysis methodologies of IAMs.

Developing physically informed AI-based models can allow material scientists, engineers, and companies to determine the physical properties (electronic, optical, transport, magnetic...) of IAMs in significantly less time than through conventional modelling, hence accelerating the path to innovation and new discoveries. These approaches will boost the exploration of complex structures relevant to energy, electronic, photonic quantum and composites applications. Moreover, properly trained models will enable to quickly test multiple experimental conditions with minimal modelling effort, which will help conventional fab and lab metrologist to access a comprehensive analysis of intricate architectures and compositions of IAMs, serving as a solution to the lack of sufficient statistical sampling for understanding performance variability among individual devices.

The transition from traditional data-centric approaches to more sophisticated foundation models is essential to address these challenges. Foundation models, which generalize across diverse datasets and tasks, will help in scaling up the modelling of IAMs and extracting more actionable insights from data. Thus, it is increasingly important to develop AI- and ML-based models embedded in workflows that not only accelerate the design of IAMs but also optimize their compositions and structures for enhanced application performance. These models should be capable of handling a broader range of tasks, from the generation of atomistic models to the simulation of electronic, transport, thermal, and optical properties. Moreover, the impact of material disorder, interface symmetries, and chemical composition on the physical properties (e.g., electronic, optical, magnetic) must be considered, especially when aiming to push the limits of device performance and metrics.

Materials data spaces (such as material digital, www.materialdigital.de/, FAIRmat www.fairmat-nfdi.eu/fairmat/, DIADEM, PSDI, CAPeX or NIMS-MPDF) [89] constitute an asset for establishing a materials commons infrastructure in which federated data repositories with trusted data

management, access, and exchange are provided. Building on the experience of several national, European and international initiatives, harmonized semantic data documentation according to FAIR principles should be developed to support interoperability and AI-readiness of produced IAM data. EU and National initiatives provided a huge and ever-increasing body of materials data that needs to be curated and made accessible to ensure maximal exploitation. To this end the data must be findable and accessible independent of its original format, i.e. semantically. Here the development of core and domain-specific ontologies opens significant long-term opportunities. While sharing the meta-data is uncritical for many stakeholders, the efforts to organize data-spaces must take into account the need for data-provenance for certain datasets. Given the enormous increase of the power of foundational models, the data needed for training emerges as the bottleneck and the potential competitive advantage for Europe.

Beyond modelling and data sharing, the role of platforms that facilitate not only the exchange of materials data but also the processing, automated analysis, and on-the-fly literature analysis is crucial. Such platforms would enable seamless integration of various stages of data flow, from acquisition to interpretation, while simultaneously providing relevant, up-to-date research knowledge that can inform and hasten the experimental or computational task at hand. This functionality would enable the acceleration of innovation in IAMs by ensuring that researchers have access to both empirical data and the latest theoretical insights. AI enhanced characterization workflows are particularly important for automating and refining data analysis techniques, allowing for rapid testing of experimental conditions with minimal manual effort. This will help experimental researchers access comprehensive analyses of complex IAM architectures and compositions, solving the problem of insufficient statistical sampling and enabling a better understanding of performance variability across breadth of domains of material science.

Combined efforts in the digitalization of materials and the validation of predictive AI models for materials enable the establishment of materials acceleration platforms, or self-driving labs, with enormous potential to revolutionize and accelerate the development of IAM both in industrial and academic settings [89]. In that sense the emerging initiatives in Europe and elsewhere stand as an opportunity but also great challenge and will demand sustained efforts and funding for the decade to come.

Data availability statement

No new data were created or analysed in this study.

Acknowledgments

C M acknowledges the support by the European Commission through the MaX Centre of Excellence for supercomputing applications (grant number 101093374). C M and N M acknowledge support from the Deutsche Forschungsgemeinschaft (DFG) under Germany's Excellence Strategy (EXC 2077, No. 390741603, University Allowance, University of Bremen) and Lucio Colombi Ciacchi, the host of the 'U Bremen Excellence Chair Program'. S R acknowledges funding from 2021 SGR 00997, funded by Generalitat de Catalunya and Grant PID2022-138283NB-I00 funded by MICIU/AEI/ 10.13039/501100011033 and by 'ERDF/EU'. ICN2 is funded by the CERCA Programme/Generalitat de Catalunya and supported by the Severo Ochoa Centres of Excellence programme, Grant CEX2021-001214-S, funded by MCIN/AEI/10.13039.501100011033. K S N and A U acknowledge support by the Ministry of Education, Singapore, under its Research Centre of Excellence award to the Institute for Functional Intelligent Materials (I-FIM, project No. EDUNC-33-18-279-V12) and National Research Foundation, Singapore under its AI Singapore Programme (AISG Award No: AISG3-RP-2022-028). M A acknowledges financial support from the Spanish MCIN/AEI/10.13039/501100011033 (Grant No. PID2022-140163NB-I00), Gobierno Vasco-UPV/EHU (Project No. IT1569-22), and the Basque Government Education Departments' IKUR program, also co-funded by the European NextGenerationEU action through the Spanish Plan de Recuperación, Transformación y Resiliencia (PRTR). V L D acknowledges support from UK Research and Innovation [grant number EP/X016188/1]. O A v L acknowledges the support by the Natural Sciences and Engineering Research Council of Canada (NSERC), [funding reference number RGPIN-2023-04853], the University of Toronto's Acceleration Consortium via the Canada First Research Excellence Fund, grant number: CREF-2022-00042, the Ed Clark Chair of Advanced Materials, a Canada CIFAR AI Chair, the European Research Council (ERC) under the European Union's Horizon 2020 research and innovation programme (grant agreement No. 772834). W W acknowledges the support by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) under Germany's Excellence Strategy for the Excellence Cluster '3D Matter Made to Order' (Grant No. EXC-2082/1-390761711) and by the Carl Zeiss Foundation.

Author contributions statement

C M conceived the perspective and coordinated it with S R. The authors' contributions within the sections are listed as follows. Sec. 1: C M, K S N, A U, S R; Sec. 2: A S B; Sec. 3: S V K, S R S; Sec. 4: K R, Sec. 5: M A, Sec. 6: V L D; Sec. 7: G Csányi; Sec. 8: C M, N M; Sec. 9: S H; G Cuniberti; Sec. 10: Q D, P O, I C; Sec. 11: K C; Sec. 12: K H, R Z; Sec. 13: O A v L; Sec. 14: M H-A, J C; Sec. 15: G Cisotto, A Z; Sec. 16: A U; Sec. 17: W W, A C F, A U, S R.

ORCID iDs

Cristiano Malica  <https://orcid.org/0000-0003-1159-4468>
Kostya S Novoselov  <https://orcid.org/0000-0003-4972-5371>
Amanda S Barnard  <https://orcid.org/0000-0002-4784-2382>
Sergei V Kalinin  <https://orcid.org/0000-0001-5354-6152>
Steven R Spurgeon  <https://orcid.org/0000-0003-1218-839X>
Karsten Reuter  <https://orcid.org/0000-0001-8473-8659>
Maite Alducin  <https://orcid.org/0000-0002-1264-7034>
Volker L Deringer  <https://orcid.org/0000-0001-6873-0278>
Gábor Csányi  <https://orcid.org/0000-0002-8180-2034>
Nicola Marzari  <https://orcid.org/0000-0002-9764-0199>
Shirong Huang  <https://orcid.org/0000-0002-4349-793X>
Gianaurelio Cuniberti  <https://orcid.org/0000-0002-6574-7848>
Qiushi Deng  <https://orcid.org/0000-0001-8118-1998>
Pablo Ordejón  <https://orcid.org/0000-0002-2353-2793>
Ivan Cole  <https://orcid.org/0000-0001-6582-1457>
Kamal Choudhary  <https://orcid.org/0000-0001-9737-8074>
Kedar Hippalgaonkar  <https://orcid.org/0000-0002-1270-9047>
O Anatole von Lilienfeld  <https://orcid.org/0000-0001-7419-0466>
Mohamed Hibat-Allah  <https://orcid.org/0000-0002-5298-8589>
Juan Carrasquilla  <https://orcid.org/0000-0001-7263-3462>
Giulia Cisotto  <https://orcid.org/0000-0002-9554-9367>
Alberto Zancanaro  <https://orcid.org/0000-0002-5276-7030>
Wolfgang Wenzel  <https://orcid.org/0000-0001-9487-4689>
Andrea C Ferrari  <https://orcid.org/0000-0003-0907-9993>
Andrey Ustyuzhanin  <https://orcid.org/0000-0001-7865-2357>
Stephan Roche  <https://orcid.org/0000-0003-0323-4665>

References

- [1] Liu T and Barnard A S 2023 *Cell Rep. Phys. Sci.* **4** 101630
- [2] Li S, Wang R, Deng Q and Barnard A S 2024 *Proc. 12th Int. Conf. on Learning Representations*
- [3] Liu T and Barnard A S 2023 *Proc. 40th Int. Conf. on Machine Learning* vol 202 p 21375
- [4] Liu T, Tho Z Y and Barnard A S 2024 *Digit. Discov.* **3** 422–35
- [5] Liu T and Barnard A S 2025 *Mach. Learn. : Eng.* accepted (<https://doi.org/10.1088/3049-4761/adaaf6>)
- [6] Liu Y T, Kelley K P, Vasudevan R K, Funakubo H, Ziatdinov M A and Kalinin S V 2022 *Nat. Mach. Intell.* **4** 341–50
- [7] Ter-Petrosyan A H, Billbrey J A, Doty C M, Matthews B E, Wang L, Du Y, Lang E, Hattar K and Spurgeon S R 2023 *Proc. Machine Learning and the Physical Sciences Workshop, NeurIPS 2023*
- [8] Lewis N R, Jin Y, Tang X, Shah V, Doty C, Matthews B E, Akers S and Spurgeon S R 2022 *npj Comput. Mater.* **8** 252
- [9] Ziatdinov M A, Liu Y, Morozovska A N, Eliseev E A, Zhang X, Takeuchi I and Kalinin S V 2022 *Adv. Mater.* **34** e2201345
- [10] Dyck O, Kim S, Jimenez-Izal E, Alexandrova A N, Kalinin S V and Jesse S 2018 *Small* **14** e1801771
- [11] Peng J *et al* 2022 *Nat. Rev. Mater.* **7** 991
- [12] Reuter K 2016 *Catal. Lett.* **146** 541
- [13] Bruix A, Margraf J T, Andersen M and Reuter K 2019 *Nat. Catal.* **2** 659
- [14] Margraf J T, Jung H, Scheurer C and Reuter K 2023 *Nat. Catal.* **6** 112
- [15] Scheurer C and Reuter K 2024 *Nat. Catal.* **8** 13–19
- [16] Hohenberg P and Kohn W 1964 *Phys. Rev.* **136** B864
- [17] Kohn W and Sham L J 1965 *Sham. Phys. Rev.* **140** A1133
- [18] Alducin M, Camillone N, Hong S-Y and Juaristi J I 2019 *Phys. Rev. Lett.* **123** 246802
- [19] Behler J 2021 *Chem. Rev.* **121** 10037
- [20] Omranpour A, Elsner J, Lausch K N and Behler J 2025 *ACS Catal.* **15** 1616–34
- [21] Serrano Jiménez A, Muzas A S, Zhang Y, Ovčar J, Jiang B, Lončarić I, Juaristi J I and Alducin M 2021 *J. Chem. Theory Comput.* **17** 4648
- [22] Muzas A S, Serrano Jiménez A, Zhang Y, Jiang B, Juaristi J I and Alducin M 2024 *J. Phys. Chem. Lett.* **15** 2587
- [23] Žugec I, Tetenoire A, Muzas A S, Zhang Y, Jiang B, Alducin M and Juaristi J I 2024 *J. Am. Chem. Soc. Att.* **4** 1997

[24] Chang C, Deringer V L, Katti K S, Van Speybroeck V and Wolverton C M 2023 Simulations in the era of exascale computing *Nat. Rev. Mater.* **8** 309

[25] Deringer V L, Bartók A P, Bernstein N, Wilkins D M, Ceriotti M and Csányi G 2021 Gaussian process regression for materials and molecules *Chem. Rev.* **121** 10073–141

[26] Ben Mahmoud C, Gardner J L A and Deringer V L 2024 Data as the next challenge in atomistic machine learning *Nat. Comput. Sci.* **4** 384

[27] Liu Y, Madanchi A, Anker A S, Simine L and Deringer V L 2025 The amorphous state as a frontier in computational materials design *Nat. Rev. Mater.* **10** 228–41

[28] El-Machachi Z, Frantzov D, Nijamudheen A, Zarrouk T, Caro M A and Deringer V L 2024 Accelerated first-principles exploration of structure and reactivity in graphene oxide *Angew. Chem., Int. Ed.* **63** e202410088

[29] Dirac P A M 1929 Quantum mechanics of many-electron systems *Proc. R. Soc. Lond. A* **123** 714–33

[30] Batatia *et al* 2023 A foundation model for atomistic materials chemistry (arXiv:2401.00096)

[31] Cococcioni M and Marzari N 2019 *Phys. Rev. Mater.* **3** 033801

[32] Malica C and Marzari N 2024 Teaching oxidation states to neural networks (arXiv:2412.01652)

[33] Bastonero L, Malica C, Macke E, Berce M, Huber S, Timrov I and Marzari N 2025 First-principles Hubbard parameters with automated and reproducible workflows (arXiv:2503.01590v1)

[34] Uhrin M, Zadoks A, Binci L, Marzari N and Timrov I 2024 (arXiv:2406.02457)

[35] Bonnet N and Marzari N 2023 *J. Chem. Theory Comput.* **20** 4820

[36] Gao S, Cheng Y, Chen L and Huang S 2024 Rapid discovery of gas response in materials via density functional theory and machine learning *Energy Environ. Mater.* **8** e12816

[37] Yang Z, Sun Y, Gao S, Yu Q, Zhao Y, Huo Y, Wan Z, Huang S, Wang Y and Gu X 2024 General model for predicting response of gas-sensitive materials to target gas based on machine learning *ACS Sens.* **9** 2509–19

[38] Huang S *et al* 2022 Machine learning-enabled smart gas sensing platform for identification of industrial gases *Adv. Intell. Syst.* **4** 2200016

[39] Huang S *et al* 2023 Machine learning-enabled graphene-based electronic olfaction sensors and their olfactory performance assessment *Appl. Phys. Rev.* **10** 021406

[40] Baek E *et al* 2020 Intrinsic plasticity of silicon nanowire neurotransistors for dynamic memory and learning functions *Nat. Electron.* **3** 398–408

[41] Chen F F, Breedon M, White P, Chu C, Mallick D, Thomas S, Sapper E and Cole I 2016 Correlation between molecular features and electrochemical properties using an artificial neural network *Mater. Des.* **112** 410–8

[42] Winkler D A, Breedon M, White P, Hughes A E, Sapper E D and Cole I 2016 Using high throughput experimental data and in silico models to discover alternatives to toxic chromate corrosion inhibitors *Corros. Sci.* **106** 229–35

[43] Deng Q, Rafiuddin Jakeria M, Elbourne A, Chen X-B and Cole I S 2025 Revising inhibiting stability of 2-mecaptobenzimidazole as corrosion inhibitor against saline corrosive media: a combined *in-situ* and *ex-situ* investigation *Appl. Surf. Sci.* **681** 161558

[44] Jeschke S, Eiden P, Deng Q, Cole I S and Keil P 2024 Structure and dynamics of aqueous 2-Aminothiazole/NaCl electrolytes at electrified interfaces *J. Phys. Chem. B* **128** 6189–96

[45] Castillo-Robles J M, de Freitas Martins E, Ordejón P and Cole I 2024 Molecular modeling applied to corrosion inhibition: a critical review *npj Mater. Degrad.* **8** 72

[46] Choudhary K *et al* 2022 *npj Comput. Mater.* **8** 59

[47] Choudhary K *et al* 2020 *npj Comput. Mater.* **6** 173

[48] Wines D, Xie T and Choudhary K 2023 *J. Phys. Chem. Lett.* **14** 6630–38

[49] Choudhary K J 2024 *Phys. Chem. Lett.* **15** 27

[50] Choudhary K *et al* 2024 *npj Comput. Mater.* **10** 93

[51] Yan D, Smith A D and Chen C-C 2023 Structure prediction and materials design with generative neural networks *Nat. Comput. Sci.* **3** 572–4

[52] Jiao R *et al* 2024 Space group constrained crystal generation (arXiv:2402.03992)

[53] Cao Z, Luo X, Lv J and Wang L 2024 Space group informed transformer for crystalline materials generation (arXiv:2403.15734) pp 1–26

[54] Zhu R, Nong W, Yamazaki S and Hippalgaonkar K 2024 WyCryst: Wyckoff inorganic crystal generator framework *Matter* **7** 1–20

[55] Zeni C *et al* 2023 MatterGen: a generative model for inorganic materials design pp 1–56 (arXiv:2312.03687)

[56] Zhao Y, Siriwardane E M D, Wu Z, Fu N, Al-Fahdi M, Hu M and Hu J 2023 Physics guided deep learning for generative design of crystal materials with symmetry constraints *npj Comput. Mater.* **9** 38

[57] Ceriotti M, Clementi C and von Lilienfeld O A 2021 *Chem. Rev.* **121** 9719

[58] von Lilienfeld O A 2018 *Angew. Chem., Int. Ed.* **57** 4164

[59] Faber F, Hutchison L, Huang B, Gilmer J, Schoenholz S S, Dahl G E, Vinyals O, Kearnes S, Riley P F and von Lilienfeld O A 2017 *J. Chem. Theory Comput.* **13** 5255

[60] Faber F, Lindmaa A, von Lilienfeld O A and Armiento R 2016 *Phys. Rev. Lett.* **117** 135502

[61] Ramakrishnan R *et al* 2015 *J. Chem. Theory Comput.* **11** 2087

[62] Heinen S *et al* 2024 *Mach. Learn.: Sci. Technol.* **5** 025058

[63] von Lilienfeld O A 2023 *Mach. Learn.: Sci. Technol.* **4** 045043

[64] Huang B, von Rudorff G F and von Lilienfeld O A 2023 *Science* **381** 170

[65] Huang B *et al* 2023 *J. Chem. Theory Comput.* **19** 1711

[66] Lee S *et al* 2024 (arXiv:2406.14524)

[67] Khan D *et al* 2025 *Sci. Adv.* **11**

[68] Naveed H, Ullah Khan A, Qiu S, Saqib M, Anwar S, Usman M, Akhtar N, Barnes N and Mian A 2024 A comprehensive overview of large language models (arXiv:2307.06435)

[69] Stephens G J, Mora T, Tkacik G and Bialek W 2013 Statistical thermodynamics of natural images *Phys. Rev. Lett.* **110** 018701

[70] Hibat-Allah M, Ganahl M, Hayward L E, Melko R G and Carrasquilla J 2020 Recurrent neural network wave functions *Phys. Rev. Res.* **2** 023358

[71] Melko R G and Carrasquilla J 2024 Language models for quantum simulation *Nat. Comput. Sci.* **4** 11–18

[72] Fitzek D, Hong Teoh Y, Pok Fung H, Dagnew G A, Ejaz Merali M S M, MacLellan B and Melko R G 2024 Rydberggpt (arXiv:2405.21052)

[73] Gyori N G, Palombo M, Clark C A, Zhang H and Alexander D C 2022 Training data distribution significantly impacts the estimation of tissue microstructure with machine learning *Magn. Reson. Med.* **87** 932–47

[74] Cisotto G, Zancanaro A, Zoppis I F and Manzoni S L 2024 hvEEGNet: a novel deep learning model for high-fidelity EEG reconstruction *Front. Neuroinform.* **18** 1459970

[75] Rathjens J and Wiskott L 2024 Classification and reconstruction processes in deep predictive coding networks: antagonists or allies (arXiv:2401.09237)

[76] Vahdat A and Kautz J 2020 NVAE: a deep hierarchical variational autoencoder *Advances in Neural Information Processing Systems* vol 33 pp 19667–79

[77] Karniadakis G E, Kevrekidis I G, Lu L, Perdikaris P, Wang S and Yang L 2021 Physics-informed machine learning *Nat. Rev. Phys.* **3** 422–40

[78] Ghane E, Fagerström M and Mirkhalaf S M 2023 A multiscale deep learning model for elastic properties of woven composites *Int. J. Solids Struct.* **282** 112452

[79] Jiadong D, Waqar M, Erofeev I, Yao K, Wang J, Pennycook S J and Duane Loh N 2023 A multiscale generative model to understand disorder in domain boundaries *Sci. Adv.* **9** eadj0904

[80] Pinto D, Greenwood M and Provatas N 2024 LeapFrog: getting the jump on multi-scale materials simulations using machine learning (arXiv:2406.15326)

[81] Gökmen D E, Biswas S, Huber S D, Ringel Z, Flickerr F and Koch-Janusz M 2023 Compression theory for inhomogeneous systems *Nat. Commun.* **15** 10214

[82] Gökmen D E, Ringel Z, Huber S D and Koch-Janusz M 2021 Symmetries and phase diagrams with real-space mutual information neural estimation *Phys. Rev. E* **104** 064106

[83] Husic B E *et al* 2020 Coarse graining molecular dynamics with graph neural networks *J. Chem. Phys.* **153** 194101

[84] Denisov K and Fedichev P 2024 Discovery of thermodynamic control variables that independently regulate healthspan and maximum lifespan (available at: www.biorxiv.org/content/10.1101/2024.12.01.626230v1.full.pdf)

[85] Kalina K A, Linden L, Brummund J and Kästner M 2023 FE^{ANN}: An efficient data-driven multiscale approach based on physics-constrained neural networks and automated data mining *Comput. Mech.* **71** 827–51

[86] Köhler J, Chen Y, Krämer A, Clementi C and Noé F 2023 FlowMatching—efficient coarse-graining of molecular dynamics without forces (arXiv:2203.11167)

[87] Jean J G, Tung-Huan S, Huang S-J, Cheng-Tang W and Chen C-S 2025 Graph-enhanced deep material network: multiscale materials modeling with microstructural informatics *Comput. Mech.* **75** 113–36

[88] Chen X, Soh B W, Ooi Z-E, Vissol-Gaudin E, Haijun Y, Novoselov K S, Hippalgaonkar K and Qianxiao L 2023 Constructing custom thermodynamics using deep learning (arXiv:2308.04119)

[89] Stier S P *et al* 2024 *Adv. Mater.* **36** 2407791