

Machine Learning in Nanoscience: Big Data at Small Scales

Keith A. Brown,^{*,†} Sarah Brittman,[‡] Nicolò Maccaferri,[§] Deep Jariwala,^{||} and Umberto Celano[⊥]

[†]Department of Mechanical Engineering, Physics Department, and Division of Materials Science and Engineering, Boston University, Boston, Massachusetts 02215, United States

[‡]U.S. Naval Research Laboratory, Washington, DC 20375, United States

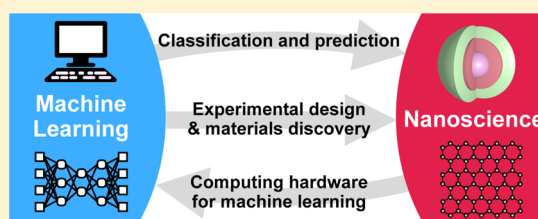
[§]Department of Physics and Materials Science, University of Luxembourg, 162a avenue de la Faïencerie, L-1511 Luxembourg, Luxembourg

^{||}Department of Electrical and Systems Engineering, University of Pennsylvania, Philadelphia, Pennsylvania 19104, United States

[⊥]imec, Kapeldreef 75, B-3001 Heverlee (Leuven), Belgium

ABSTRACT: Recent advances in machine learning (ML) offer new tools to extract new insights from large data sets and to acquire small data sets more effectively. Researchers in nanoscience are experimenting with these tools to tackle challenges in many fields. In addition to ML's advancement of nanoscience, nanoscience provides the foundation for neuromorphic computing hardware to expand the implementation of ML algorithms. In this Mini Review, we highlight some recent efforts to connect the ML and nanoscience communities by focusing on three types of interaction: (1) using ML to analyze and extract new insights from large nanoscience data sets, (2) applying ML to accelerate material discovery, including the use of active learning to guide experimental design, and (3) the nanoscience of memristive devices to realize hardware tailored for ML. We conclude with a discussion of challenges and opportunities for future interactions between nanoscience and ML researchers.

KEYWORDS: Machine learning, data-driven research, active learning, materials discovery, design of experiments



In nanoscience, high-throughput experimentation enabled by the small size of nanoscale samples and rapid, high-resolution imaging tools are becoming increasingly widespread.^{1,2} For example, in nanophotonics^{3,4} and catalysis^{2,5} material properties have been varied systematically across the same wafer-sized substrate and characterized locally using high-resolution scanning probe and optical or electron microspectroscopy techniques. These or similar methods can generate data sets that are too vast and complex for researchers to mentally parse without computational assistance; yet, these data are rich in relationships that the researchers would like to understand. Machine learning (ML) enables researchers to analyze large data sets by training models that can be used to classify observations into discrete groups, learn which features determine a metric of performance, or predict the outcome of new experiments. Furthermore, even in fields where such data-intensive methods are not typical, ML can assist researchers in designing experiments to optimize performance or test hypotheses more effectively. From nano-optoelectronics, to catalysis, to the bionano interface, ML is reshaping how researchers collect, analyze, and interpret their data. These methods will likely evolve into new standards tailored for each field complementary to the role statistics currently plays in scientific research.⁶ In return, nanoscience has the potential to benefit ML by developing electronic or photonic hardware that can implement algorithms more efficiently than conventional computing architectures. Deepening this unique relationship (Figure 1) has much to offer both research communities.

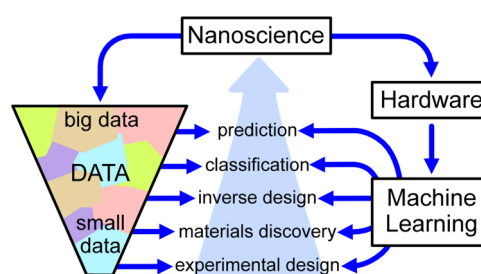


Figure 1. Scheme showing how the fields of nanoscience and machine learning interact via data and hardware. On the left, nanoscience generates data that, combined with machine learning, feeds back a variety of functions to advance nanoscience research. On the right, nanoscience contributes the nanoscale hardware components that can advance the field of machine learning by enabling new processing architectures.

Broadly speaking, ML encompasses algorithmic approaches for classifying data, identifying empirical correlations within the data, and predicting the consequences of these correlations in new data. These algorithms learn from the data itself to refine the accuracy of their predictions, typically by minimizing a mathematical error function. As learning can be defined as modifying behavior based on past experience, which in the case

Received: October 3, 2019

Revised: November 27, 2019

Published: December 5, 2019

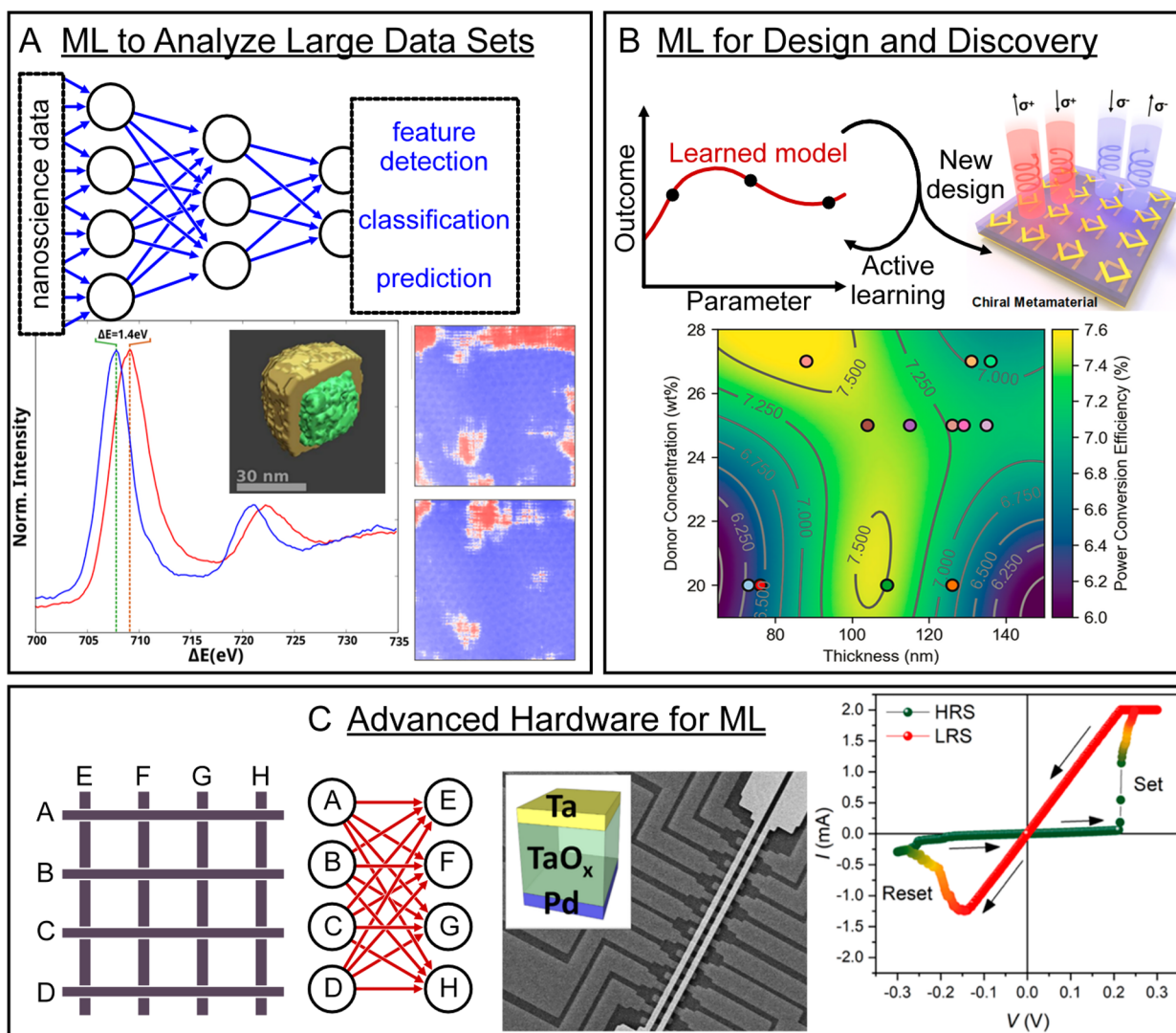


Figure 2. Schematics and representative images showing the categories of interaction between machine learning and nanoscience. (A) Large sets of data can be used to train models that facilitate analysis (top). Typical data sets include images (bottom right, reproduced from ref 7. Copyright 2018 American Chemical Society) and spectra (bottom left, reproduced from ref 8. Copyright 2016 American Chemical Society). (B) Learned models can aid in the selection of new experiments via active learning (bottom, reproduced from ref 9. Copyright 2018 American Chemical Society) or the design/discovery of novel materials or structures such as chiral metasurfaces (top, reproduced from ref 10. Copyright 2018 American Chemical Society). (C) Nanoscience-enabled hardware, such as memristor arrays (reproduced from ref 11. Copyright 2017 American Chemical Society) can function as physical embodiments of machine learning algorithms such as artificial neural networks (ANN). Specifically, a crossbar array of memristors connecting inputs A–D to output E–H has a strong analogy to a fully connected neural network (top). By allowing one to tune electrically between a high resistance state (HRS) and a low resistance state (LRS), memristors can perform both storage and computation functions (reproduced from ref 12. Copyright 2019 American Chemical Society).

of ML is presented in the form of training data, ML is defined very broadly from basic regression techniques to state-of-the-art approaches. Algorithms are written to address classes of problems and then trained for a specific task based on which type of data is available. To train an algorithm by supervised learning, the data must be labeled, which means that each piece of data comprises an input (e.g., parameters of an experiment or design of a material) and an output (e.g., outcome of an experiment or material property of interest). The algorithm takes in the input features of the data set and builds a model (based on internal assumptions) that produces the output of the data set with as little error as possible. This model, which is just a mapping of input to output, can then be used to predict the outputs when it is given inputs that are not included in the training data. For supervised learning, the input features of the

data must be predetermined and are not selected by the algorithm. Prediction and classification are two tasks commonly performed by supervised ML. Alternatively, if there is no information about output in the data, unsupervised learning can still be used to uncover relationships. Clustering and component analysis are tasks that are commonly performed via unsupervised learning. Learning can also be semisupervised, in which the model is initially produced by supervised learning on labeled data and then refined by unsupervised learning on unlabeled data. Through the extensive work of mathematicians and computer scientists, algorithms for classes of common problems can be rapidly applied using off-the-shelf and often open-source platforms. Even artificial neural networks (ANNs), a versatile and recently

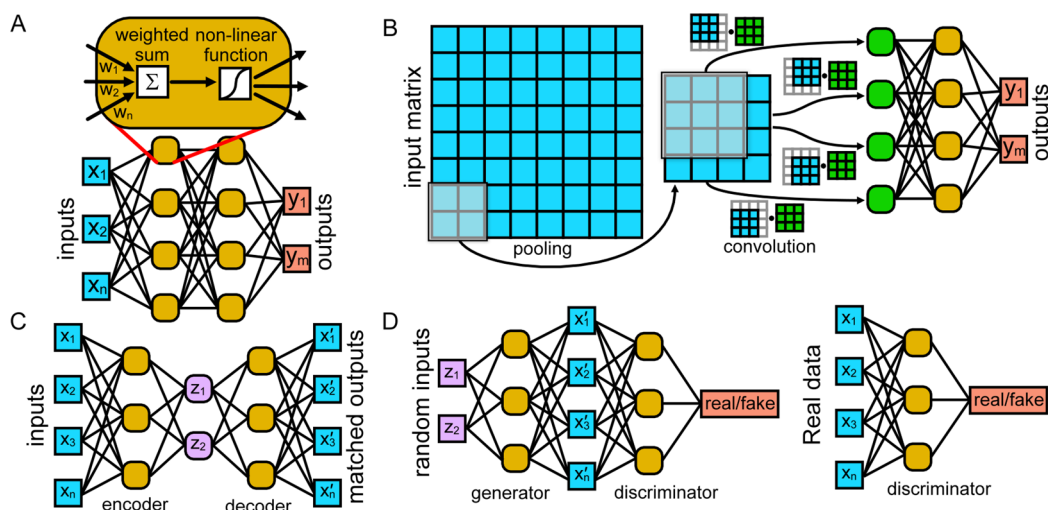


Figure 3. Schematics showing simplified architectures of commonly applied neural networks. (A) An ANN in which yellow nodes depict nonlinear operations on weighted sums of inputs to yield outputs. The nodes collectively form layers of computation within the ANN. (B) A CNN in which pooling and convolution operations on local subsets of the data effectively prioritize spatial correlations within the model. Pooling replaces a local region of the data by its maximum value (or other summarizing statistic) and therefore reduces the dimensionality of the data. (C) An autoencoder, which is trained to match its output to its input and in the process, finds a lower-dimensional representation of the input data (purple layer). (D) A GAN in which real data is used to train a discriminator network to differentiate between real and fake data. A generator therein produces fake data from random inputs and tries to fool this discriminator. Both networks are trained jointly in opposition to each other to improve at their respective tasks.

popular class of algorithms, are now available for researchers in nanoscience to bring to bear on problems in their own fields.

In this Mini Review, which is not able to be comprehensive, we explore the intersection between ML and nanoscience and highlight three expanding classes of interactions (Figure 2): (1) *The use of ML tools in nanoscience to analyze data.* Manufacturing and metrology of nanoelectronics for integrated circuits (ICs) are highlighted here as a case study. (2) *The discovery of new materials at the nanoscale using ML for inverse design and design of experiments.* Novel two-dimensional materials are discussed as a case study. (3) *How nanoscience can empower ML through the development of novel hardware.* In particular, nanoscale memristors, along with other emerging architectures including nanophotonics, have the potential to provide a hardware platform tailored for ML. Finally, we conclude by discussing challenges and opportunities arising at the interface between nanoscience and ML.

1. Machine Learning to Analyze Large Data Sets.

When many simulations or measurements are available as training data, ML can be used to identify features in data from nanoscale systems. Commonly, scientists interpret these features with physical models that suggest further experiments. This interaction illustrates the most basic role that ML can play in the scientific method.

Analysis of Spectra, Images, and Biological Outcomes. Machine Learning is widely contributing to the recognition and classification of key features in nanoscience data sets. For example, in X-ray spectroscopy an ANN (Figure 3A) was trained on simulated X-ray absorption fine structure (EXAFS) spectra generated by molecular dynamics simulations. Then, using this trained model, researchers extracted partial radial distribution functions that yielded insight into the chemical structure beyond the first atomic coordination shell.¹³ This new analysis indicated that surface effects change the atomic ordering of PdAu nanocrystals and therefore their catalytic properties. In an example from photoluminescence spectroscopy, supervised learning on experimental data was used to

extract the distribution of decay rates in CsPbBr₃ nanocrystals without being forced to assume a functional form of this distribution. The resulting distribution was well explained by three types of emissive species proposed to exist in these nanocrystals.¹⁴

In the analysis of images, which have some degree of spatial correlation, convolutional neural networks (CNNs) trained by supervised learning have been very successful (Figure 3B).^{15,16} These are ANNs whose internal operations are restricted so that they learn primarily local correlations within the data, and their models are invariant to small translations. They are therefore well suited to identify image features, which depend on local spatial correlations. For example, skyrmions were studied using labeled Lorentz transmission electron microscope (TEM) images,⁷ and phases of matter were identified in Monte Carlo simulations of Ising systems or square spin-ice models.¹⁷ Although a fully connected ANN could identify phases with a simple order parameter, a CNN was required to correctly analyze more complex spin models that had no order parameter.¹⁷ It is possible that the CNN's constrained focus on local spatial correlations allowed it to learn the distinguishing features between the phases more efficiently than a fully connected ANN; however, such a hypothesis still remains to be validated. The image classification capabilities of CNNs and their utility for pattern recognition have also been applied in data storage. Information encoded in deep subwavelength structures was recovered by training a CNN on the observed color of the structures to achieve high-information density and robustness against fabrication defects.¹⁸

The bio-nano interface has also been a major focus for the application of ML. An early example was the training of an ANN to study uptake of 109 types of nanoparticles into cells in order to predict their toxicity by their chemical composition.¹⁹ More recently, a similar process was implemented where nanoparticles were injected into a rat and then isolated to determine the makeup of the proteins that decorate their surface. This proteomic information, together with the fate of

the nanoparticles within the rat, was then used to train an ANN to predict nanoparticle fate.²⁰ In combination with a nanofluidic chip that can detect a variety of biomarkers, ML has been used to interpret gene expression data from exosomes to classify the disease state of patients.²¹ More generally, ML-enabled analysis of the content of liquid biopsies has experienced a rapid acceleration in recent years powered by the high degree of clinical relevance and inherent complexity of the systems.²²

Deconvoluting Components in Mixed Signals. Experiments often measure the collective effect of many individual components, so deconvoluting complex mixed signals given little initial information is a common analytical task. Approaches to this problem, known as blind source separation (BSS), are based on statistics of the data, such as the covariance matrix. For example, principal component analysis (PCA) identifies the eigenvectors of the covariance matrix to determine which linear combinations of input variables encode the most unique information. Such processes have been used in scanning probe microscopy to rapidly denoise and compress data using all available data channels.²³ In addition to PCA, which does not explicitly classify data, linear discriminant analysis (LDA) identifies linear combinations of parameters that optimally classify data. LDA has been used with a nanofluidic system to interpret gene expression data from exosomes to classify the disease state of patients.²¹ A related technique, independent component analysis (ICA), identifies nonorthogonal basis vectors that best account for the correlations in the data, including higher order statistics than covariances. For instance, electron microscopy-based energy-dispersive X-ray spectroscopy (EDS) has long been used to determine elemental maps and recently been expanded to identify phases. In nanoscience, ICA has been used to interpret EDS data and identify multielemental phases such as Fe₂O₃ and a PtFe alloy²⁴ in core-shell nanoparticles and to identify AuI₂ nanowires grown on an InSb substrate.²⁵ Similar techniques have also been applied to interpret electron energy loss spectra of iron oxide nanocubes to achieve three-dimensional chemical mapping of iron oxidation states.⁸ While PCA, LDA, and ICA are relatively simple forms of unsupervised learning,²⁶ their increasingly widespread use constitutes an important step of popularizing advanced statistics in nanoscience research. A deeper comparison of PCA, LDA, and ICA, along with the mathematical formalism for these approaches, can be found in textbooks on statistical learning²⁶ and has been recently provided in the literature related to their use in electroencephalogram classification.²⁷

Case Study: Machine Learning for Metrology of Nanoelectronics. Recently, ML has been applied to address major challenges in the large-scale manufacturing of nanoscale devices for integrated circuits (ICs). Miniaturization still drives the evolution of IC technology with transistor density consistently increasing as the technology continues to mature. However, transistor fabrication is currently facing some of the most intense challenges yet in the area of photolithography, new materials, defects analysis, and device architectures.²⁸ Semiconductor ICs are based on increasingly complex manufacturing processes realized at the nanoscale by many interconnected tools that generate enormous sets of data, particularly metrology data related to the properties and performance of fabricated structures.

Given the high-dimensional and vast nature of this data, optimally acting on the available data in a timely fashion is a

major challenge with profound implications to device cost.²⁸ For example, the front-end of the line (FEOL) represents when the finest features on a wafer are manufactured and metrology at this stage is employed in order to inspect indicators of a successful execution. Here, uncertainty requirements for fault detection and classification are beyond the capability of inline metrology techniques. ANNs trained on the readout of inline inspection tools can offer predictive modeling of complex measurements when the metrology system does not have sufficient resolution. For example, ANN-based accurate prediction of critical dimensions (CD) has been demonstrated in the sub-40 nm trenches for extreme ultraviolet resist, where an ANN was used to predict the CD error, resist shrinkage, and metal line resistance in the early process pipeline.^{29,30} Taking advantage of spatial correlations inherent to images, CNNs have been successfully applied to improve the accuracy in the classification of wafer maps³¹ and fault detection during growth of thin films made using chemical vapor deposition.³² Rather than inspecting every location on a sample, a recent trend in semiconductor processing is virtual metrology (VM) in which a selection of wafers are sampled and used to train models that correlate process sensor data to the performance metric of interest. Since process sensor data are always available, this trend potentially represents a rapid acceleration in manufacturing of complex nanostructures if such relationships can be rapidly and accurately learned.³³ Such metrology using ML algorithms can pinpoint defects and measure structures in advanced chips. As the technology advances, inspection and metrology based on ML will become even more important for correlation of process flows, predictive metrology, and yield analysis.

2. Machine Learning for Design and Discovery.

Predicting the properties of a material is a central challenge in materials science and chemistry.³⁴ Nanomaterials are even more complex: structuring materials on the nanoscale leads them to adopt different properties than their bulk counterparts and allows the construction of heterostructures or metamaterials that include multiple materials. In the face of this vast parameter space, ML can help predict novel materials, optimize structures, and even plan experiments.³⁵ A particularly impactful area of design and discovery is inverse design or finding a set of parameters that produces a desired outcome.

Inverse Design and Adversarial Networks in Nanophotonics. Optical metamaterials and nanophotonics are fields in which the experimental design space is vast due to the availability of high-resolution lithographic tools to construct intricate structures. Although a conventional approach of measuring the optical properties of a given nanostructure is conceptually straightforward, inverse design is made extremely difficult because the existence or uniqueness of an acceptable design cannot be guaranteed.

One approach to address this challenge is to train an ANN in a supervised fashion using known input/output combinations and then to use the ANN to iterate through unknown input parameters until a desired outcome is predicted. Such a trial-and-error approach was applied to design tailored optical responses of multilayer nanoparticles.⁴ This approach takes advantage of the fact that evaluating a trained ANN is typically much faster than running a brute-force optimization algorithm that evaluates possible combinations of parameters using a physics-based process. Here, this acceleration was achieved because the gradient can be found analytically for ANN's whereas it must be computed numerically for the optimization

algorithm. However, as this study notes, agreement between the ANN and the numerical solution, calculated here using the transfer matrix method, must always be verified. While this approach can be construed as relying heavily on regression to produce inverse design solutions, other approaches that rely more strongly upon physical principles to simplify or explore parameter space are also being used to great effect.³⁶

A different approach that has been explored in the design of metasurfaces utilizes a cascading deep neural network (DNN) in which two networks are trained, one that maps each design to an outcome and one that maps the outcome to designs,³ in a manner that mirrors the use of an encoder-decoder pairs in DNNs for parameter reduction.³⁷ This architecture is commonly known as an autoencoder and is used to reduce the dimensionality of the data (Figure 3C).³⁸ These networks produced candidate patterns that matched the desired spectra with high fidelity, expediting the discovery of metasurfaces with tailored optical responses.³ In addition, a deep-learning-based model, comprising two bidirectional neural networks assembled by a partial stacking strategy, was used to optimize three-dimensional chiral metamaterials with strong chiroptical responses at predesignated wavelengths.¹⁰ In principle, this type of model can help provide insight into the physical underpinnings that connect structures and their properties by elucidating the intricate relationships between metamaterial structures and their optical response.¹⁰

Another way to address inverse problems is through the use of generative adversarial networks (GANs), which is a recently invented unsupervised learning strategy (Figure 3D).³⁹ The GAN comprises two networks, a generator that guesses distributions of parameters and a discriminator that evaluates the quality of each guess by comparing it to existing unlabeled data. GANs have been used to design nanophotonic structures that have precise user-defined spectral responses.⁴⁰ In this case, the use of a GAN was motivated by the desire to allow the designed structures to not require input from an expert scientist. The rapid application of GANs after their invention signals that the nanoscience community is motivated to adopt novel learning approaches to rapidly meet pressing challenges such as inverse design. However, there is a need to benchmark ML-based approaches to inverse design against existing methods in order to elucidate the acceleration and improvements that are possible using ML.

Active Learning, Automated Experimentation, and Autonomous Researchers. There have been a number of innovations that allow researchers to conduct experiments more efficiently and to explore more of parameter space. For example, ML has been used to aid researchers in selecting combinations of experimental parameters that reduce the number of total experiments necessary to optimize multilayer organic solar cells.⁹ Further, experiments can be selected without input from humans based on insight from ML. Specifically, active learning describes the use of ML to select experiments to most efficiently achieve a goal. For example, by iteratively performing density functional theory (DFT) calculations that had been selected by machine learning, promising intermetallic surfaces for catalysis were identified.⁵ This approach bears some similarity to directed evolution, which is often applied in the space of protein design. Here, candidate protein structures are selected, modified, and tested in sequential generations. Directed evolution has been highly successful because structurally similar proteins often share similar properties.⁴¹

Perhaps most exciting are systems in which active learning is combined with automated experimentation to realize fully autonomous researchers or robot scientists (i.e., research systems that select and carry out experiments without a human in the loop).⁴² In the first example of this paradigm in nanoscience, a system was built that allowed for the automated growth and characterization of carbon nanotubes on the surface of micropillars by locally heating them with laser illumination. In addition to this automated experimental system, logical regression analysis was used to autonomously pick the next experimental conditions (in terms of temperature and precursor partial pressure) to offer the best chance of achieving the experimental goal. This analysis is a form of stepwise regression in which linear terms connecting inputs (e.g., synthesis conditions) and growth outcomes were sequentially tested to see if they improved the regression quality. By realizing a fully autonomous research system (ARES), the pace of information generation and analysis was accelerated substantially.⁴³ Autonomy has also been realized in scanning tunneling microscopy where a learning system was constructed to determine the state of the tip using CNN of images and recondition the tip when needed to maintain high imaging performance.⁴⁴

Case Study: Prediction of new 2D Materials and Heterostructures. An important driving force for innovations in materials discovery has been the isolation of van der Waals bonded layered materials into atomically thin, two-dimensional (2D) sheets.⁴⁵ The stacking of 2D materials with various compositions and rotational orientations has led to heterostructures with novel properties.⁴⁶ The number of possible combinations has become experimentally intractable and hence ML techniques have been critical in identifying new compounds and structures and classifying them by properties. Pairing first-principles models with active learning in a Bayesian framework, van der Waals heterostructures with desired electronic band gaps and thermoelectric properties have been proposed.⁴⁷ Further, vast material databases with functional descriptors such as bonding directionality, packing factor, and interlayer gap, have enabled new van der Waals 2D and 1D (chainlike) compounds to be identified.^{48,49} These physics-based ML models can be further improved with additional descriptors such as structure and composition⁵⁰ to predict unknown magnetic phases in known materials.⁵¹ Furthermore, by tuning the weights attributed to different features of these physics-based models, these algorithms can even predict the “synthesizability” of new 2D materials from their bulk counterparts, as has been demonstrated for the MXene family of 2D materials.⁵²

3. Nanoscience To Advance Hardware for Machine Learning. The technological revolution enabled by computation relies on devices made by nanoscience and, correspondingly, a tremendous amount of basic science and engineering in nanoscience has emerged from the study of ICs or related systems. However, the von Neumann architecture of computing is not the most efficient for implementing the myriad ML algorithms that have been developed.⁵³ The lack of efficiency of the von Neumann architecture stems from the time and energy required to transfer data between spatially separated memory and processing units,⁵⁴ as well as its failure to take advantage of analog operations that arise naturally from the physics of hardware components.⁵³ The result is that a von Neumann computer requires many thousands of transistors and memory elements to compute the action of a single

artificial neuron.⁵³ This observation raises the question of whether novel nanoscale architectures could provide a new neuromorphic hardware platform for advanced machine learning algorithms.

Novel Materials for Memristors. While neuromorphic architectures can be approximated using digital logic in ANNs, neuromorphic computing is an emerging field in which the hardware-level system is designed to more closely represent the architecture of the brain. A memristor, which is defined to be a two-terminal device in which conductance is a function of the prior voltages experienced by the device, is often considered to be the fundamental unit of such neuromorphic computation. These devices have been only recently (2008) experimentally realized through the observation that many nanoscale materials exhibit memristive properties through ionic motion.⁵⁵ For example, tantalum oxide films exhibit a conductivity that is tunable based upon the local density of oxygen vacancies, leading to complex and useful memristive dynamics.¹¹ Such devices can be as small as a single conductive filament.^{56,57} Sheets of MoS₂ have been found to exhibit a frequency-dependent memristive effect due to Joule heating, making them amenable to tasks such as localizing sound while dissipating exceptionally low power.⁵⁸ Such 2D materials can also be directly used as memristors, as was recently shown by diffusing copper atoms vertically between a MoS₂ bilayer.¹² Conductance modulation has also been achieved using ferroelectric thin-film transistors, which are three-terminal analogues of memristors,⁵⁹ as well as by engineering defects such as grain boundaries⁶⁰ and dislocations⁶¹ in 2D and 3D semiconductors, respectively. Finally, nanophotonic systems have also recently been explored as candidates for neuromorphic computing with successes including the realization of deep learning networks⁶² and adsorption-based photonic neural networks.⁶³

Design and Realization of Arrays of Memristors. Memristors are commonly fabricated in crossbar arrays so that the number of devices scales favorably with the number of electrical connections. For example, an 18-memristor array was used to realize a PCA algorithm that analyzed nine metrics of cell geometry to determine whether a cell was benign or malignant.¹¹ More recently, a 30-memristor array with a tantalum oxide active layer was utilized to implement a *k*-means analysis to accurately categorize flowers according to their geometry.⁶⁴ Nanoscale films of poly(1,3,5-trivinyl-1,3,5-trimethyl cyclotrisiloxane) (pV3D3) have been used to form arrays with 64 independent memristors.⁶⁵ Simulations of larger arrays of these devices were shown to be able to perform common image processing tasks including facial recognition.⁶⁵ Memristor arrays can be very compact: a recent work showed a functioning 3 × 3 array that was constructed with a 6 nm pitch.⁶⁶ Although making high density arrays has been a major push in the nanoscience community, efforts to realize arrays with larger dimensions have been very successful as well with 128 × 64 member arrays of microscale components.⁶⁷

4. Challenges and Opportunities for Machine Learning and Nanoscience. As this is a Mini Review, we are unable to cover all of the ways in which ML and nanoscience enrich each other. Both disciplines will benefit from closer interactions between data scientists and nanoscience researchers. For example, discussions of the quantity, quality, and labeling of data that is realistically available can inform the development of new algorithms. Also, incorporating data scientists into the review process will ensure that sound

technical practices from ML make their way into nanoscience. Although this Mini Review has principally focused on the advantages of using ML in nanoscience, there are pitfalls that practitioners can encounter such as overtraining that warrant caution.²²

The following six areas deserve special mention as challenges and opportunities for the future:

(1) For many applications, efficient, and confident use of ML requires large data sets that sample broadly from the target problem's parameter space. A key avenue for transformatively increasing the data available to researchers is enabling cross-study comparison through standardization of data format and conventions for metadata. Metadata is critical for nanomaterial research because it is unclear which of the numerous structural or compositional descriptors are most valuable for a given property. Although there are significant practical hurdles to be overcome, such standardization could lead to a "nanomaterials genome" that could be utilized to springboard diverse research initiatives.⁶⁸

(2) Simulation is often used to generate training data for ML algorithms because calculations are typically faster than experiments and can be performed with total control over the initial conditions. However, when a ML model is trained using simulated data, its accuracy depends on the ability of simulation to reproduce experiment, which can vary dramatically from field to field. How to assign uncertainty to simulated data and use it for multifidelity learning approaches, which try to account for differences in the accuracy of subsets of training data, is a rich and unsolved problem.

(3) Forefront ML algorithms must continue to be rapidly applied in nanoscience research to fully realize their impact. These types of collaborations are often difficult as they require nuanced communication across disciplines and go beyond widely utilized open source tools. Advances in these areas are expected to range from advanced generative algorithms for inverse design to new processes for learning how the structure of ANN provides deeper insights into the physical system. For example, recent studies from biology use "visible" neural networks, where existing knowledge of the hierarchical biological system is used to generate a structured neural network. Its substructures can be directly interpreted because they were designed to mimic the connectivity in a real cell.⁶⁹ Along similar lines, ML approaches that have physics and other domain knowledge built into them could provide a rapid approach to reducing the volume of training data and increasing the physical intuition that can be drawn from the outcome.⁷⁰

(4) Autonomous research systems have been demonstrated to be powerful tools for accelerating discoveries but the hardware infrastructure is highly application specific, which often makes investing in the development of such systems difficult to justify. As a field, combinatorial chemistry has overcome this by producing a unified set of tools and form factors (namely microtiter plates) that allow standardized experiments to be carried out without human intervention. Similar hardware standardization would lower the barrier to apply automation to address any given system. Candidates for such massive parallelization include inkjet printing and scanning probe lithography.^{71,72}

(5) While memristors and neuromorphic computing are progressing toward the level of a scalable practical technology, several challenges remain in terms of large area uniformity, reproducibility of the components, switching speed/efficiency,

and total lifetime in terms of cycles.⁷³ These challenges can be addressed through either the development of novel memristive systems or improvements to existing systems. However, in all cases, integration with existing CMOS platforms and competitive performance advantage over CMOS neurons must be realized to have a long-term technological impact. Lastly, although such analog networks have the potential to be highly efficient once they are trained, training such networks is not yet as flexible or efficient as digital logic. Thus, integrating such systems with traditional digital logic is still of high importance.

(6) Finally, integration of ML techniques and quantum computing has the potential to address currently untenable problems.^{74,75} Recently, quantum algorithms that can solve a problem of supervised learning, such as constructing a classifier, were proposed.⁷⁶ It was also recently demonstrated that superconducting quantum circuits can be used to realize quantum generative adversarial learning.⁷⁷ These successes indicate that systems that combine quantum science and machine learning might have applications that extend far beyond classification. In this context, the role of nanoscience is 2-fold. First, the performance of ML will be boosted beyond the current limits through the implementation of quantum algorithms, and this might have a huge impact on how we approach quantum chemistry problems. Second, nanoscience will play a crucial role in developing hardware components that can be used to build quantum computers capable of implementing such quantum machine learning algorithms.^{78,79}

AUTHOR INFORMATION

Corresponding Author

*E-mail: brownka@bu.edu.

ORCID

Keith A. Brown: 0000-0002-2379-2018

Sarah Brittman: 0000-0002-7307-7565

Nicolò Maccaferri: 0000-0002-0143-1510

Deep Jariwala: 0000-0002-3570-8768

Umberto Celano: 0000-0002-2856-3847

Author Contributions

The manuscript was written through contributions of all authors. All authors have given approval to the final version of the manuscript.

Notes

The authors declare no competing financial interest.

ACKNOWLEDGMENTS

We acknowledge very helpful discussions with Aldair Gongora, Chong Liu, and Mario Miscuglio. K.A.B. acknowledges support from the NSF (CMMI-1661412) and AFOSR (FA9550-16-1-0150). S.B. holds an NRC Research Associate award at the U.S. Naval Research Laboratory. D.J. thanks support from startup funds from Penn Engineering and NSF (DMR-1905853) and University of Pennsylvania Materials Research Science and Engineering Center (MRSEC) (DMR-1720530). N.M. acknowledges the financial support from the FEDER program (Grant 2017-03-022-19 Lux-Ultra-Fast).

REFERENCES

- (1) Lignos, I.; Stavarakis, S.; Nedelcu, G.; Protesescu, L.; deMello, A. J.; Kovalenko, M. V. Synthesis of Cesium Lead Halide Perovskite Nanocrystals in a Droplet-Based Microfluidic Platform: Fast Parametric Space Mapping. *Nano Lett.* **2016**, *16* (3), 1869–1877.
- (2) Kluender, E. J.; Hedrick, J. L.; Brown, K. A.; Rao, R.; Meckes, B.; Du, J.; Moreau, L.; Maruyama, B.; Mirkin, C. A. Catalyst discovery through megalibraries of nanomaterials. *Proc. Natl. Acad. Sci. U. S. A.* **2019**, *116* (1), 40–45.
- (3) Liu, D.; Tan, Y.; Khoram, E.; Yu, Z. Training deep neural networks for the inverse design of nanophotonic structures. *ACS Photonics* **2018**, *5* (4), 1365–1369.
- (4) Peurifoy, J.; Shen, Y.; Jing, L.; Yang, Y.; Cano-Renteria, F.; DeLacy, B. G.; Joannopoulos, J. D.; Tegmark, M.; Soljačić, M. Nanophotonic particle simulation and inverse design using artificial neural networks. *Science advances* **2018**, *4* (6), No. eaar4206.
- (5) Tran, K.; Ulissi, Z. W. Active learning across intermetallics to guide discovery of electrocatalysts for CO₂ reduction and H₂ evolution. *Nature Catalysis* **2018**, *1* (9), 696.
- (6) Bzdok, D.; Altman, N.; Krzywinski, M. Points of significance: statistics versus machine learning. *Nat. Methods* **2018**, *15*, 233–234.
- (7) Shibata, K.; Tanigaki, T.; Akashi, T.; Shinada, H.; Harada, K.; Niitsu, K.; Shindo, D.; Kanazawa, N.; Tokura, Y.; Arima, T.-h. Current-Driven Motion of Domain Boundaries between Skyrmion Lattice and Helical Magnetic Structure. *Nano Lett.* **2018**, *18* (2), 929–933.
- (8) Torruella, P.; Arenal, R.; de la Peña, F.; Saghi, Z.; Yedra, L.; Eljarrat, A.; López-Conesa, L.; Estrader, M.; López-Ortega, A.; Salazar-Alvarez, G.; et al. 3D visualization of the iron oxidation state in FeO/Fe₃O₄ core-shell nanocubes from electron energy loss tomography. *Nano Lett.* **2016**, *16* (8), 5068–5073.
- (9) Cao, B.; Adutwum, L. A.; Oliyynyk, A. O.; Luber, E. J.; Olsen, B. C.; Mar, A.; Buriak, J. M. How to optimize materials and devices via design of experiments and machine learning: demonstration using organic photovoltaics. *ACS Nano* **2018**, *12* (8), 7434–7444.
- (10) Ma, W.; Cheng, F.; Liu, Y. Deep-learning-enabled on-demand design of chiral metamaterials. *ACS Nano* **2018**, *12* (6), 6326–6334.
- (11) Choi, S.; Shin, J. H.; Lee, J.; Sheridan, P.; Lu, W. D. Experimental Demonstration of Feature Extraction and Dimensionality Reduction Using Memristor Networks. *Nano Lett.* **2017**, *17* (5), 3113–3118.
- (12) Xu, R.; Jang, H.; Lee, M.-H.; Amanov, D.; Cho, Y.; Kim, H.; Park, S.; Shin, H.-J.; Ham, D. Vertical MoS₂ double layer memristor with electrochemical metallization as an atomic-scale synapse with switching thresholds approaching 100 mV. *Nano Lett.* **2019**, *19* (4), 2411–2417.
- (13) Timoshenko, J.; Wrasman, C. J.; Luneau, M.; Shirman, T.; Cargnello, M.; Bare, S. R.; Aizenberg, J.; Friend, C. M.; Frenkel, A. I. Probing atomic distributions in mono- and bimetallic nanoparticles by supervised machine learning. *Nano Lett.* **2019**, *19* (1), 520–529.
- (14) Đorđević, N.; Beckwith, J. S.; Yarema, M.; Yarema, O.; Rosspeintner, A.; Yazdani, N.; Leuthold, J.; Vauthey, E.; Wood, V. Machine Learning for Analysis of Time-Resolved Luminescence Data. *ACS Photonics* **2018**, *5* (12), 4888–4895.
- (15) Chua, L. O.; Roska, T. The CNN paradigm. *IEEE Transactions on Circuits and Systems I: Fundamental Theory and Applications* **1993**, *40* (3), 147–156.
- (16) Egmont-Petersen, M.; de Ridder, D.; Handels, H. Image processing with neural networks—a review. *Pattern recognition* **2002**, *35* (10), 2279–2301.
- (17) Carrasquilla, J.; Melko, R. G. Machine learning phases of matter. *Nat. Phys.* **2017**, *13* (5), 431.
- (18) Wiecha, P. R.; Lecestre, A.; Mallet, N.; Larrieu, G. Pushing the limits of optical information storage using deep learning. *Nat. Nanotechnol.* **2019**, *14* (3), 237.
- (19) Epa, V. C.; Burden, F. R.; Tassa, C.; Weissleder, R.; Shaw, S.; Winkler, D. A. Modeling biological activities of nanoparticles. *Nano Lett.* **2012**, *12* (11), 5808–5812.
- (20) Lazarovits, J.; Sindhvani, S.; Tavares, A. J.; Zhang, Y.; Song, F.; Audet, J.; Krieger, J. R.; Syed, A. M.; Stordy, B.; Chan, W. C. W. Supervised Learning and Mass Spectrometry Predicts the in Vivo Fate of Nanomaterials. *ACS Nano* **2019**, *13* (7), 8023–8034.
- (21) Ko, J.; Bhagwat, N.; Yee, S. S.; Ortiz, N.; Sahnoud, A.; Black, T.; Aiello, N. M.; McKenzie, L.; O'Hara, M.; Redlinger, C.; Romeo, J.

Carpenter, E. L.; Stanger, B. Z.; Issadore, D. Combining Machine Learning and Nanofluidic Technology To Diagnose Pancreatic Cancer Using Exosomes. *ACS Nano* **2017**, *11* (11), 11182–11193.

(22) Ko, J.; Baldassano, S. N.; Loh, P.-L.; Kording, K.; Litt, B.; Issadore, D. Machine learning to detect signatures of disease in liquid biopsies—a user's guide. *Lab Chip* **2018**, *18* (3), 395–405.

(23) Jesse, S.; Kalinin, S. V. Principal component and spatial correlation analysis of spectroscopic-imaging data in scanning probe microscopy. *Nanotechnology* **2009**, *20* (8), 085714.

(24) Rossouw, D.; Burdet, P.; de la Peña, F.; Ducati, C.; Knappett, B. R.; Wheatley, A. E.; Midgley, P. A. Multicomponent signal unmixing from nanoheterostructures: Overcoming the traditional challenges of nanoscale x-ray analysis via machine learning. *Nano Lett.* **2015**, *15* (4), 2716–2720.

(25) Jany, B. R.; Janas, A.; Krok, F. Retrieving the Quantitative Chemical Information at Nanoscale from Scanning Electron Microscope Energy Dispersive X-ray Measurements by Machine Learning. *Nano Lett.* **2017**, *17* (11), 6520–6525.

(26) Hastie, T.; Tibshirani, R.; Friedman, J.; Franklin, J. The elements of statistical learning: data mining, inference and prediction. *Mathematical Intelligencer* **2005**, *27* (2), 83–85.

(27) Subasi, A.; Gursoy, M. I. EEG signal classification using PCA, ICA, LDA and support vector machines. *Expert Systems with Applications* **2010**, *37* (12), 8659–8666.

(28) Orji, N. G.; Badaroglu, M.; Barnes, B. M.; Beitia, C.; Bunday, B. D.; Celano, U.; Kline, R. J.; Neisser, M.; Obeng, Y.; Vladar, A. Metrology for the next generation of semiconductor devices. *Nature electronics* **2018**, *1* (10), 532.

(29) Rana, N.; Zhang, Y.; Kagalwala, T.; Bailey, T. Leveraging advanced data analytics, machine learning, and metrology models to enable critical dimension metrology solutions for advanced integrated circuit nodes. *J. Micro/Nanolithogr., MEMS, MOEMS* **2014**, *13* (4), 041415.

(30) Breton, M.; Chao, R.; Muthinti, G. R.; Abraham, A.; Simon, J.; Cepler, A. J.; Sendelbach, M.; Gaudiello, J.; Emans, S.; Shifrin, M. In *Electrical test prediction using hybrid metrology and machine learning, Metrology, Inspection, and Process Control for Microlithography XXXI*; International Society for Optics and Photonics, 2017; p 1014504.

(31) Nakazawa, T.; Kulkarni, D. V. Wafer map defect pattern classification and image retrieval using convolutional neural network. *IEEE Transactions on Semiconductor Manufacturing* **2018**, *31* (2), 309–314.

(32) Lee, K. B.; Cheon, S.; Kim, C. O. A convolutional neural network for fault classification and diagnosis in semiconductor manufacturing processes. *IEEE Transactions on Semiconductor Manufacturing* **2017**, *30* (2), 135–142.

(33) Kang, P.; Kim, D.; Cho, S. Semi-supervised support vector regression based on self-training with label uncertainty: An application to virtual metrology in semiconductor manufacturing. *Expert Systems with Applications* **2016**, *51*, 85–106.

(34) Butler, K. T.; Davies, D. W.; Cartwright, H.; Isayev, O.; Walsh, A. Machine learning for molecular and materials science. *Nature* **2018**, *559* (7715), 547.

(35) Wang, M.; Wang, T.; Cai, P.; Chen, X. Nanomaterials Discovery and Design through Machine Learning. *Small Methods* **2019**, *3* (5), 1900025.

(36) Leuchs, G.; Sondermann, M. Time-reversal symmetry in optics. *Phys. Scr.* **2012**, *85* (5), 058101.

(37) Mao, X.; Shen, C.; Yang, Y.-B. In *Image restoration using very deep convolutional encoder-decoder networks with symmetric skip connections*; Advances in neural information processing systems; Neural Information Processing Systems Foundation, Inc.: 2016; pp 2802–2810.

(38) Kingma, D. P.; Welling, M. Auto-encoding variational bayes. **2013**, arXiv:1312.6114.

(39) Goodfellow, I.; Pouget-Abadie, J.; Mirza, M.; Xu, B.; Warde-Farley, D.; Ozair, S.; Courville, A.; Bengio, Y. In *Generative adversarial nets*; Advances in neural information processing systems; Neural

Information Processing Systems Foundation, Inc.: 2014; pp 2672–2680.

(40) Liu, Z.; Zhu, D.; Rodrigues, S. P.; Lee, K.-T.; Cai, W. Generative model for the inverse design of metasurfaces. *Nano Lett.* **2018**, *18* (10), 6570–6576.

(41) Arnold, F. H. Design by directed evolution. *Acc. Chem. Res.* **1998**, *31* (3), 125–131.

(42) King, R. D.; Whelan, K. E.; Jones, F. M.; Reiser, P. G.; Bryant, C. H.; Muggleton, S. H.; Kell, D. B.; Oliver, S. G. Functional genomic hypothesis generation and experimentation by a robot scientist. *Nature* **2004**, *427* (6971), 247.

(43) Nikolaev, P.; Hooper, D.; Perea-López, N.; Terrones, M.; Maruyama, B. Discovery of Wall-Selective Carbon Nanotube Growth Conditions via Automated Experimentation. *ACS Nano* **2014**, *8* (10), 10214–10222.

(44) Rashidi, M.; Wolkow, R. A. Autonomous scanning probe microscopy in situ tip conditioning through machine learning. *ACS Nano* **2018**, *12* (6), 5185–5189.

(45) Butler, S. Z.; Hollen, S. M.; Cao, L.; Cui, Y.; Gupta, J. A.; Gutiérrez, H. R.; Heinz, T. F.; Hong, S. S.; Huang, J.; Ismach, A. F.; et al. Progress, challenges, and opportunities in two-dimensional materials beyond graphene. *ACS Nano* **2013**, *7* (4), 2898–2926.

(46) Geim, A. K.; Grigorieva, I. V. Van der Waals heterostructures. *Nature* **2013**, *499* (7459), 419.

(47) Bassman, L.; Rajak, P.; Kalia, R. K.; Nakano, A.; Sha, F.; Sun, J.; Singh, D. J.; Aykol, M.; Huck, P.; Persson, K.; et al. Active learning for accelerated design of layered materials. *npj Computational Materials* **2018**, *4* (1), 74.

(48) Cheon, G.; Duerloo, K.-A. N.; Sendek, A. D.; Porter, C.; Chen, Y.; Reed, E. J. Data mining for new two-and one-dimensional weakly bonded solids and lattice-commensurate heterostructures. *Nano Lett.* **2017**, *17* (3), 1915–1923.

(49) Ashton, M.; Paul, J.; Sinnott, S. B.; Hennig, R. G. Topology-scaling identification of layered solids and stable exfoliated 2D materials. *Phys. Rev. Lett.* **2017**, *118* (10), 106101.

(50) Cheon, G.; Cubuk, E. D.; Antoniuk, E. R.; Blumberg, L.; Goldberger, J. E.; Reed, E. J. Revealing the spectrum of unknown layered materials with superhuman predictive abilities. *J. Phys. Chem. Lett.* **2018**, *9* (24), 6967–6972.

(51) Miyazato, I.; Tanaka, Y.; Takahashi, K. Accelerating the discovery of hidden two-dimensional magnets using machine learning and first principle calculations. *J. Phys.: Condens. Matter* **2018**, *30* (6), 06LT01.

(52) Frey, N. C.; Wang, J.; Vega Bellido, G. I.; Anasori, B.; Gogotsi, Y.; Shenoy, V. B. Prediction of Synthesis of 2D Metal Carbides and Nitrides (MXenes) and Their Precursors with Positive and Unlabeled Machine Learning. *ACS Nano* **2019**, *13* (3), 3031–3041.

(53) Mead, C. Neuromorphic electronic systems. *Proc. IEEE* **1990**, *78* (10), 1629–1636.

(54) Backus, J. Can programming be liberated from the von Neumann style? A functional style and its algebra of programs. *Commun. ACM* **1978**, *21*, 613–641.

(55) Strukov, D. B.; Snider, G. S.; Stewart, D. R.; Williams, R. S. The missing memristor found. *Nature* **2008**, *453* (7191), 80.

(56) Berco, D.; Zhou, Y.; Gollu, S. R.; Kalaga, P. S.; Kole, A.; Hassan, M.; Ang, D. S. Nanoscale conductive filament with alternating rectification as an artificial synapse building block. *ACS Nano* **2018**, *12* (6), 5946–5955.

(57) Celano, U.; Goux, L.; Degraeve, R.; Fantini, A.; Richard, O.; Bender, H.; Jurczak, M.; Vandervorst, W. Imaging the three-dimensional conductive channel in filamentary-based oxide resistive switching memory. *Nano Lett.* **2015**, *15* (12), 7970–7975.

(58) Sun, L.; Zhang, Y.; Hwang, G.; Jiang, J.; Kim, D.; Eshete, Y. A.; Zhao, R.; Yang, H. Synaptic computation enabled by joule heating of single-layered semiconductors for sound localization. *Nano Lett.* **2018**, *18* (5), 3229–3234.

(59) Kim, M.-K.; Lee, J.-S. Ferroelectric Analog Synaptic Transistors. *Nano Lett.* **2019**, *19* (3), 2044–2050.

(60) Sangwan, V. K.; Jariwala, D.; Kim, I. S.; Chen, K.-S.; Marks, T. J.; Lauthon, L. J.; Hersam, M. C. Gate-tunable memristive phenomena mediated by grain boundaries in single-layer MoS₂. *Nat. Nanotechnol.* **2015**, *10* (5), 403.

(61) Choi, S.; Tan, S. H.; Li, Z.; Kim, Y.; Choi, C.; Chen, P.-Y.; Yeon, H.; Yu, S.; Kim, J. SiGe epitaxial memory for neuromorphic computing with reproducible high performance based on engineered dislocations. *Nat. Mater.* **2018**, *17* (4), 335.

(62) Shen, Y.; Harris, N. C.; Skirlo, S.; Prabhu, M.; Baehr-Jones, T.; Hochberg, M.; Sun, X.; Zhao, S.; Larochelle, H.; Englund, D.; et al. Deep learning with coherent nanophotonic circuits. *Nat. Photonics* **2017**, *11* (7), 441.

(63) George, J. K.; Mehrabian, A.; Amin, R.; Meng, J.; De Lima, T. F.; Tait, A. N.; Shastri, B. J.; El-Ghazawi, T.; Prucnal, P. R.; Sorger, V. J. Neuromorphic photonics with electro-absorption modulators. *Opt. Express* **2019**, *27* (4), 5181–5191.

(64) Jeong, Y.; Lee, J.; Moon, J.; Shin, J. H.; Lu, W. D. K-means data clustering with memristor networks. *Nano Lett.* **2018**, *18* (7), 4447–4453.

(65) Jang, B. C.; Kim, S.; Yang, S. Y.; Park, J.; Cha, J.-H.; Oh, J.; Choi, J.; Im, S. G.; Dravid, V. P.; Choi, S.-Y. Polymer Analog Memristive Synapse with Atomic-Scale Conductive Filament for Flexible Neuromorphic Computing System. *Nano Lett.* **2019**, *19* (2), 839–849.

(66) Pi, S.; Li, C.; Jiang, H.; Xia, W.; Xin, H.; Yang, J. J.; Xia, Q. Memristor crossbar arrays with 6-nm half-pitch and 2-nm critical dimension. *Nat. Nanotechnol.* **2019**, *14* (1), 35.

(67) Hu, M.; Graves, C. E.; Li, C.; Li, Y.; Ge, N.; Montgomery, E.; Davila, N.; Jiang, H.; Williams, R. S.; Yang, J. J.; et al. Memristor-based analog computation and neural network classification with a dot product engine. *Adv. Mater.* **2018**, *30* (9), 1705914.

(68) Qian, C.; Siler, T.; Ozin, G. A. Exploring the possibilities and limitations of a nanomaterials genome. *Small* **2015**, *11* (1), 64–69.

(69) Yu, M. K.; Ma, J.; Fisher, J.; Kreisberg, J. F.; Raphael, B. J.; Ideker, T. Visible machine learning for biomedicine. *Cell* **2018**, *173* (7), 1562–1565.

(70) Stewart, R.; Ermon, S. In *Label-free supervision of neural networks with physics and domain knowledge*; Thirty-First AAAI Conference on Artificial Intelligence, San Francisco, California, USA, February 04–09, 2017.

(71) Park, J.-U.; Hardy, M.; Kang, S. J.; Barton, K.; Adair, K.; Mukhopadhyay, D. k.; Lee, C. Y.; Strano, M. S.; Alleyne, A. G.; Georgiadis, J. G.; Ferreira, P. M.; Rogers, J. A. High-resolution electrohydrodynamic jet printing. *Nat. Mater.* **2007**, *6* (10), 782–789.

(72) Brown, K. A.; Hedrick, J. L.; Eichelsdoerfer, D. J.; Mirkin, C. A. Nanocombinatorics with Cantilever-Free Scanning Probe Arrays. *ACS Nano* **2019**, *13* (1), 8–17.

(73) Zidan, M. A.; Strachan, J. P.; Lu, W. D. The future of electronics based on memristive systems. *Nature Electronics* **2018**, *1* (1), 22–29.

(74) Dunjko, V.; Briegel, H. J. Machine learning & artificial intelligence in the quantum domain: a review of recent progress. *Rep. Prog. Phys.* **2018**, *81* (7), 074001.

(75) Schuld, M.; Killoran, N. Quantum machine learning in feature Hilbert spaces. *Phys. Rev. Lett.* **2019**, *122* (4), 040504.

(76) Havlíček, V.; Córcoles, A. D.; Temme, K.; Harrow, A. W.; Kandala, A.; Chow, J. M.; Gambetta, J. M. Supervised learning with quantum-enhanced feature spaces. *Nature* **2019**, *567* (7747), 209.

(77) Hu, L.; Wu, S.-H.; Cai, W.; Ma, Y.; Mu, X.; Xu, Y.; Wang, H.; Song, Y.; Deng, D.-L.; Zou, C.-L.; Sun, L. Quantum generative adversarial learning in a superconducting quantum circuit. *Science Advances* **2019**, *5* (1), No. eaav2761.

(78) von Lilienfeld, O. A. Quantum machine learning in chemical compound space. *Angew. Chem., Int. Ed.* **2018**, *57* (16), 4164–4169.

(79) Krylov, A. I. The quantum chemistry of open-shell species. *Reviews in Computational Chemistry* **2017**, *30*, 151–224.