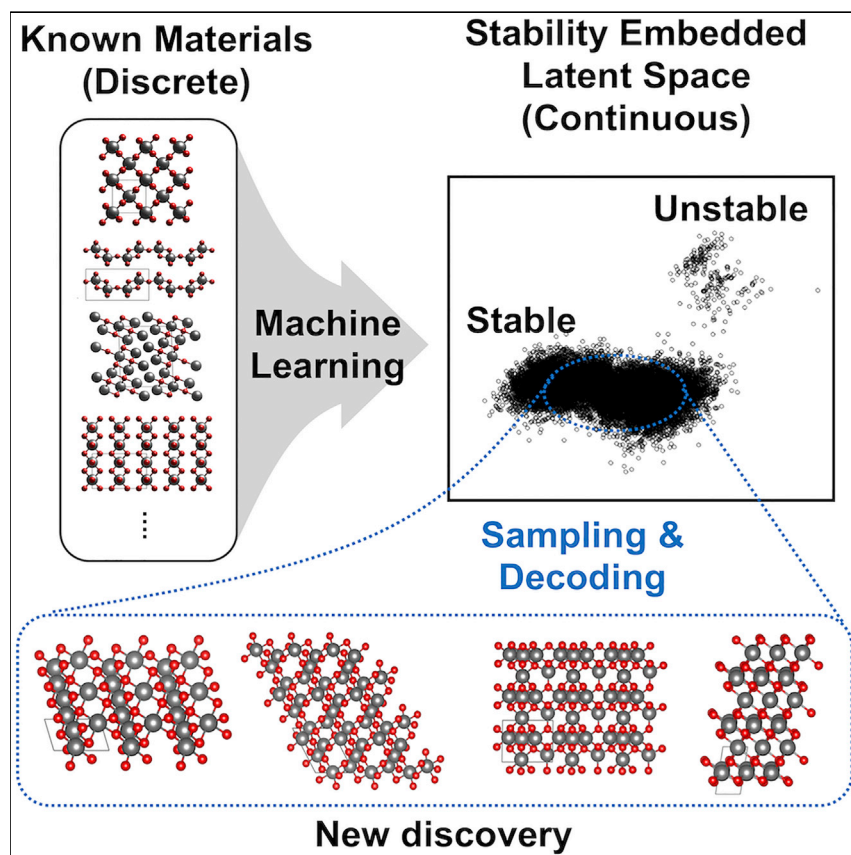


## Article

# Inverse Design of Solid-State Materials via a Continuous Representation



The inverse design of new materials with desired properties is the ultimate goal of materials research, but demonstrating such a possibility for inorganic solid-state materials has been challenging, due partly to the invertibility of representation. Here, we demonstrate that the generative model using invertible image-based representation yields accurate reconstruction performance and can successfully rediscover experimentally known vanadium oxides. The model predicts several completely new compositions and polymorphs of vanadium oxides that are metastable and may be synthesizable.

Juhwan Noh, Jaehoon Kim, Helge S. Stein, Benjamin Sanchez-Lengeling, John M. Gregoire, Alan Aspuru-Guzik, Yousung Jung

alan@aspuru.com (A.A.-G.)  
ysjn@kaist.ac.kr (Y.J.)

## HIGHLIGHTS

Invertible image-based representation is introduced to encode solid-state materials

Inverse design framework is constructed using a continuous materials-latent space

The model is shown to rediscover experimental materials that have been synthesized

New metastable V-O compounds that may be synthesizable are generated



## Benchmark

First qualification/assessment of material properties and/or performance

Noh et al., Matter 1, 1370–1384  
November 6, 2019 © 2019 Elsevier Inc.  
<https://doi.org/10.1016/j.matt.2019.08.017>



## Article

# Inverse Design of Solid-State Materials via a Continuous Representation

Juhwan Noh,<sup>1</sup> Jaehoon Kim,<sup>2</sup> Helge S. Stein,<sup>3</sup> Benjamin Sanchez-Lengeling,<sup>4</sup> John M. Gregoire,<sup>3</sup> Alan Aspuru-Guzik,<sup>5,6,7,\*</sup> and Yousung Jung<sup>1,2,8,\*</sup>

## SUMMARY

The non-serendipitous discovery of materials with targeted properties is the ultimate goal of materials research, but to date, materials design lacks the incorporation of all available knowledge to plan the synthesis of the next material. This work presents a framework for learning a continuous representation of materials and building a model for new discovery using latent space representation. The ability of autoencoders to generate experimental materials is demonstrated with vanadium oxides via rediscovery of experimentally known structures when the model was trained without them. Approximately 20,000 hypothetical materials are generated, leading to several completely new metastable  $V_xO_y$  materials that may be synthesizable. Comparison with genetic algorithms suggests computational efficiency of generative models that can explore chemical compositional space effectively by learning the distributions of known materials for crystal structure prediction. These results are an important step toward machine-learned inverse design of inorganic functional materials using generative models.

## INTRODUCTION

In the pursuit of improving everyday life through novel materials, the materials science community reported about  $2 \times 10^5$  inorganic materials within the past century for various purposes and applications<sup>1,2</sup> ranging from energy storage and production to healthcare. Considering the breadth of inorganic chemistries offered by the periodic table, the size of the entire materials space is almost unlimited, and the charted territory by various databases is only a tiny fraction. Previous efforts in exploring this combinatorially large space have relied mainly on chemical intuition and empirical rules obtained from prior knowledge to devise plausible materials. To cover pressing material needs, however, one must accelerate efforts toward the inverse design of materials, which seeks to discover a new material given a desired functionality.<sup>3,4</sup>

In this context, data-driven approaches such as high-throughput screening (HTS), both experimental and computational, have shown success in designing new materials for a myriad of applications such as organic photovoltaic materials (OPV),<sup>5</sup> manganese-based solar fuel photoanodes by combining experimental and theoretical approaches,<sup>6</sup> metal chalcogenide semiconductor materials,<sup>7</sup> BaTiO<sub>3</sub>-based piezoelectric materials by using statistical learning and experimental test,<sup>8</sup> metastable nitrides,<sup>9,10</sup> and other applications.<sup>11–13</sup> In particular, a few success stories for OPV and photoanodes in which data-driven HTS led to successful synthesis lead us to believe that this pathway is promising to reach the goal of accelerated non-serendipitous discovery of new and improved materials.

## Progress and Potential

While a traditional strategy for materials design has been to use chemical intuition and empirical rules, combining it with data science and machine learning can significantly expand the search space and accelerate the new discovery. Machine-learning models in materials science have been most extensively developed to predict properties of candidate materials, which still requires the selection of candidates. Inverting the role of machine learning to generate a candidate material with selected properties requires development of generative models for materials, as demonstrated herein. The inverse design pipeline for inorganic solids presented here is based on an invertible image-based featurization and is applied to find new crystal polymorphs of vanadium oxides. This proof-of-concept demonstration opens a great possibility of inverse designing new inorganic solid-state functional materials with desired properties.

Recent advances in deep learning, availability of large high-quality datasets, and more affordable computation have propelled the inverse design of materials. Previous work in this direction lies with the discovery of new alloy materials. For crystal structure prediction, Fischer et al.<sup>14</sup> proposed data-mining structure prediction, a probability-based model for binary alloys later extended by Hautier et al.<sup>15</sup> to ternary materials. Unlike the latter probabilistic learning model, Ryan et al.<sup>16</sup> proposed normalized atomic fingerprints to predict crystal structures for given alloy compositions. There are other software packages using genetic algorithm-based models, such as CALYPSO,<sup>17</sup> USPEX,<sup>18</sup> and XtalOpt,<sup>19</sup> which predict thermodynamically stable or metastable crystal structures for given chemical compositions and external conditions. All the aforementioned models, however, are mapping from an input (e.g., composition) to an output (e.g., structure and/or stability), and only leverage patterns in its input representation while not leveraging knowledge about the underlying structure of chemical space.

Generative models (GMs),<sup>3</sup> the focus of the present work, are a particularly promising new approach to find novel functional materials in this context. Contrary to existing models<sup>14–17,20–22</sup> for exploring the materials space, GMs tackle the inverse problem by learning to model the distribution of a dataset, which for our purposes will be a set of materials and properties. Just as one can sample random numbers from a Gaussian distribution, new compositions and structures are sampled from this learned distribution. One example of a GM is a variational autoencoder (VAE),<sup>23</sup> composed of two deep neural networks, an encoder and a decoder. The encoder maps data points to a low-dimensional continuous vector space, the latent space, and the decoder maps latent vectors back to data points. Both models are trained to deconstruct (encoder) and reconstruct (decoder) data points, learning to compress and decompress data via a meaningful intermediate representation. Gomez-Bombarelli et al.<sup>24</sup> demonstrated the use of VAE for the generation of small drug-like molecules. By including properties in the training procedure, the latent space would organize itself around the molecular properties, enabling global optimization to find novel materials that extremize such a property. Several other GMs using an adversarial autoencoder,<sup>25</sup> a recurrent neural network,<sup>26</sup> and reinforcement learning models<sup>27–29</sup> were also proposed to generate drug-like molecules.

Unlike molecules, however, the application of GMs to inorganic solid-state materials has been challenging due to the limited available data for inorganic solids, the enormous chemical breadth offered by the periodic table, and lack of invertible representations. For molecules, there are several large databases with millions of molecules (e.g., ZINC,<sup>30</sup> ChEMBL,<sup>31</sup> GDB-13,<sup>32</sup> GDB-17<sup>33</sup>), as well as several representations such as SMILES, InChI, and molecular-graph<sup>34</sup> that can be reversibly encoded. For inorganic solids, databases such as the International Crystal Structure Database (ICSD),<sup>1,2</sup> Materials Project (MP),<sup>35</sup> Open Quantum Materials Database,<sup>36</sup> and Atomic-FLOW for materials discovery<sup>37</sup> have been developed and used for prediction of formation energy, band gap, and various materials properties using machine learning. For the latter models, fingerprinting of the materials using properties of materials, crystal graphs,<sup>38</sup> Atom2Vec,<sup>39</sup> Ewald sum matrix,<sup>40</sup> Generalized Coulomb matrix,<sup>40</sup> and partial radial distribution function<sup>41</sup> have been proposed with a promising accuracy of prediction. Notably, however, no fingerprint (representation) for inorganic solid-state materials has been demonstrated to be invertible from representation to a material, hence the lack of generative inverse design for inorganic materials.

In this work, we present an invertible materials encoding/decoding scheme, and based on it construct an inverse materials design framework (iMatGen: Image-based

<sup>1</sup>Department of Chemical and Biomolecular Engineering, Korea Advanced Institute of Science and Technology (KAIST), 291 Daehakro, Daejeon 34141, Korea

<sup>2</sup>Graduate School of EEWS, Korea Advanced Institute of Science and Technology (KAIST), 291 Daehakro, Daejeon 34141, Korea

<sup>3</sup>Joint Center for Artificial Photosynthesis, California Institute of Technology, Pasadena, CA 91125, USA

<sup>4</sup>Department of Chemistry and Chemical Biology, Harvard University, 12 Oxford Street, Cambridge, MA 02138, USA

<sup>5</sup>Department of Chemistry and Department of Computer Science, University of Toronto, Toronto, ON M5S 3H6, Canada

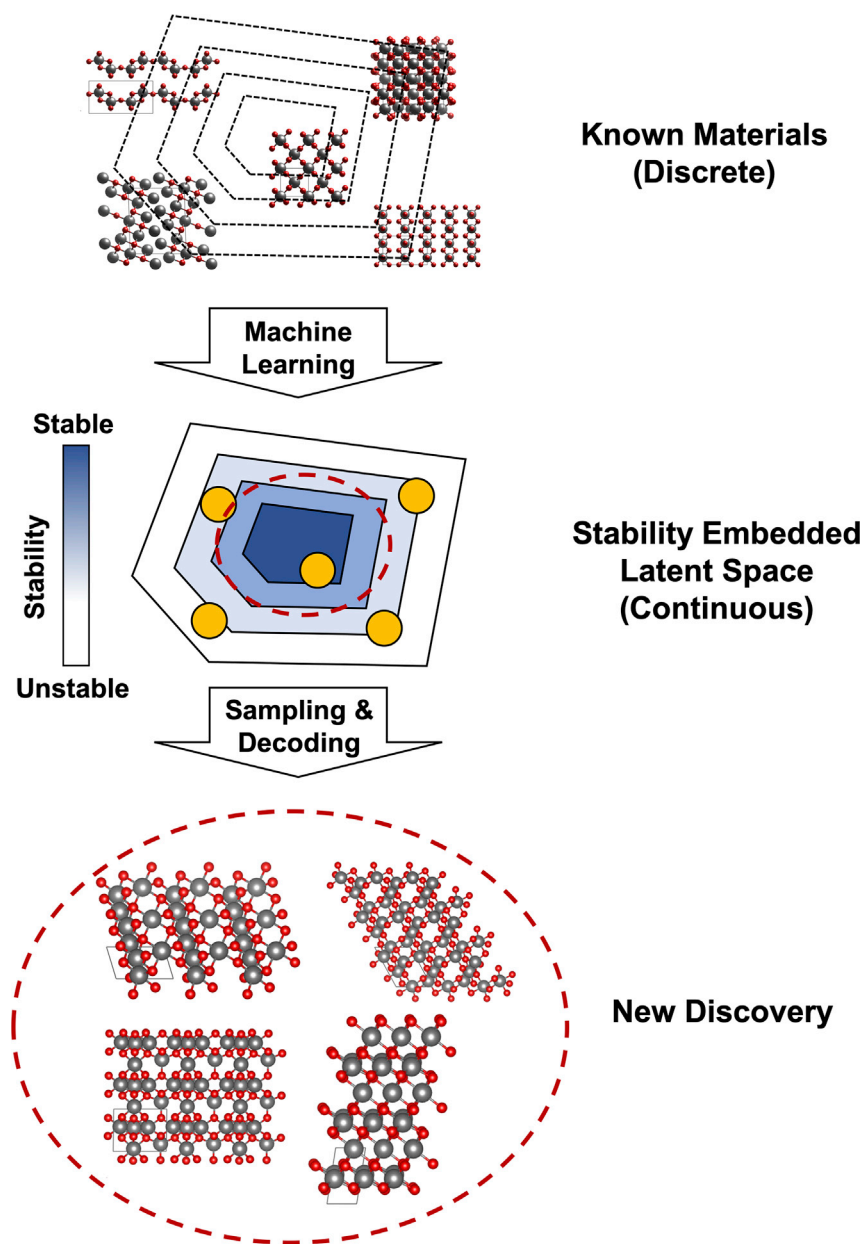
<sup>6</sup>Vector Institute for Artificial Intelligence, Toronto, ON M5S 1M1, Canada

<sup>7</sup>Canadian Institute for Advanced Research (CIFAR) Senior Fellow, Toronto, ON M5S 1M1, Canada

<sup>8</sup>Lead Contact

\*Correspondence: [alan@aspuru.com](mailto:alan@aspuru.com) (A.A.-G.), [ysjn@kaist.ac.kr](mailto:ysjn@kaist.ac.kr) (Y.J.)

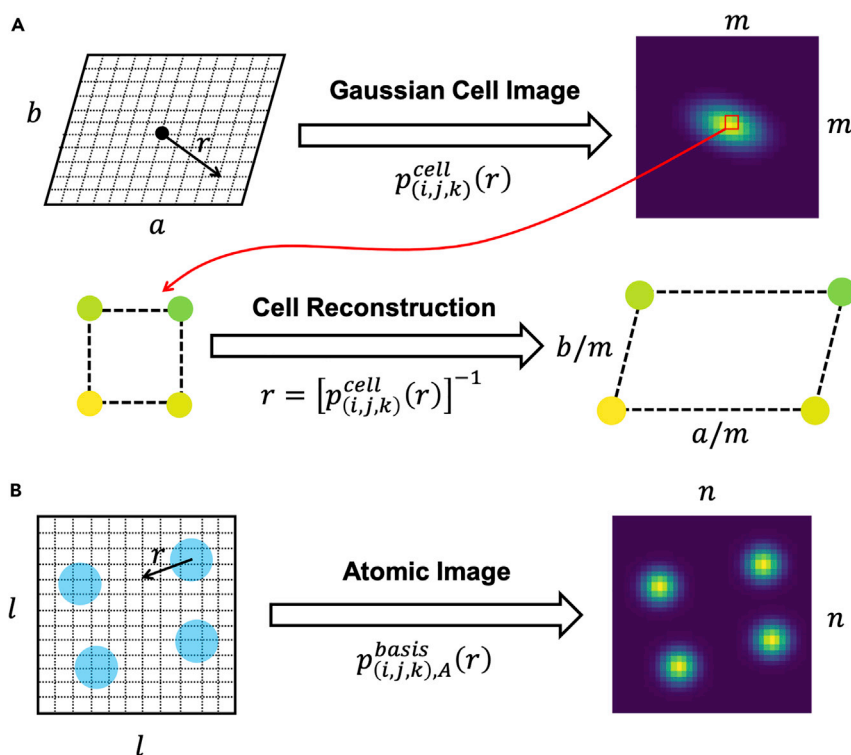
<https://doi.org/10.1016/j.matt.2019.08.017>



**Figure 1. Schematic Flow of the Inverse Materials Design via Latent Space Construction**

By building a stability embedded continuous materials space using the proposed crystal representation and generative model, stable materials can be sampled and decoded for a new discovery.

Materials Generator) for the first time to computationally predict new compositions and structural polymorphs of vanadium oxide (see Figure 1 for schematic diagram). The V-O system was chosen for the study largely because of the extent of training data available in the MP database due to a wide range of possible oxidation states of V (between +2 and +5), but the size of this dataset is also emblematic of the detailed investigation of this system to date, making it a particularly challenging space for discovery and motivating our inspection of how the iMatGen-discovered materials compare with known structures. We show that the current iMatGen can generate experimentally known structures when the model is trained without



**Figure 2. Image Representation of Crystal Structure and Inverse Transform**

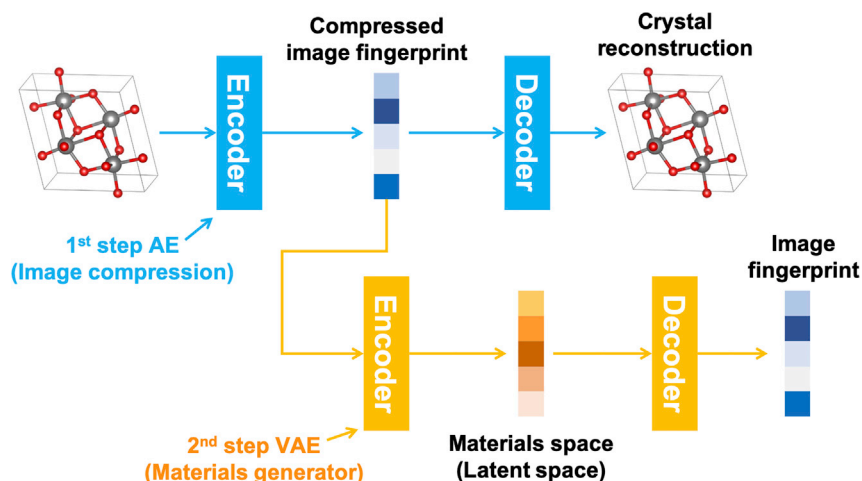
For simplicity, we show the representation method for a slice of the cubic box. Crystal structure is encoded using the two images, (A) cell (in the grid space of  $n^3$  grid points) and (B) basis (in the grid space of  $m^3$  grid points). Crystal lattice parameters can be exactly recovered from the cell image using a simple inverse transform scheme, and atom positions are recovered from the basis images using an image filter. See also [Experimental Procedures](#) and [Section S3.2](#) for details.

them, and in addition discover new compositions of V-O as well as new polymorphs of known V-O compositions with low energy above convex hull that might be synthesizable thermodynamically.

## RESULTS AND DISCUSSION

### Continuous Representation of Inorganic Materials

The first building block in the invertible materials framework is the continuous representation of inorganic materials that includes crystal structure and composition. To represent inorganic crystal structures, we use a three-dimensional (3D) grid-based image representation as initial input, similar to earlier proposals.<sup>42–44</sup> Particular to our representation, we decompose the crystal structure into a unit cell image (length of the cell edges and angles between them) and basis image (atomic positions within a unit cell), as shown in [Figure 2](#). We note that the representation scheme in [Figure 2](#) is shown for a two-dimensional (2D) slice of the 3D crystal unit cell for simplicity, but actual implementation of the method for real crystals in the current study used 3D grids. Since image-based tasks such as image generation or object recognition are widely developed in the deep learning community, we leverage these technologies by using an image-based representation for crystal structures.<sup>45</sup> Computational details of how to encode the materials data into image-based representations and decode back to original materials data (coordinates and unit cell parameters) are described in [Experimental Procedures](#) and [Section S3.2](#).



**Figure 3. The Proposed Hierarchical Two-Step Image-Based Materials Generator**

Image compression for basis and unit cell followed by materials generator (MG) using a variational autoencoder (VAE) model. For detailed model structure, see [Section S1](#).

### Image-Based Materials Generator

We constructed the iMatGen in two steps (see [Figure 3](#)), the first step to reduce the dimension of images (image compression) for both basis and cell, and the second step to encode the materials (materials generator) utilizing the element information from the first step. As shown in [Figure 3](#), the role of convolutional autoencoders in the first step is to reduce the materials input image data to intermediate vectors of smaller dimension that can be used as a fingerprint for the second materials generation step. In the second step, the materials generator (MG), the actual V-O materials space is then constructed using the VAE.<sup>23</sup> To additionally organize the latent space around a property, we follow an approach proposed by Snoek et al.<sup>46</sup> by adding a neural network mapping from the latent space to a binary label relating the formation energies to stability of materials (see [Figures S3](#) and [S4](#)). We labeled a crystal structure satisfying  $E_f \leq 0.5$  eV/atom as a stable material (+1), and all others as unstable (0). The 0.5 eV/atom buffer was used not to miss any metastable yet potentially interesting materials as well as considering inherent errors of the present density functionals used to compute  $E_f$  (see [Section S2.1](#)). This classification task was incorporated into the training procedure of VAE by introducing an additional loss function to a usual expression for VAE loss function as shown in [Section S1.3](#). The detailed structure of the model and its hyper-parameters are described in [Section S1](#) and [Figure S1](#).

### Dataset

We used the MP database,<sup>35</sup> one of the largest open source databases including 83,989 inorganic structures. In searching for a new composition or polymorphs of V-O system using autoencoder, we note that there are 112 known vanadium oxide materials in the MP database, among which there are 25 unique V-O compositions. Neither of these 112 structures nor 25 compositions are sufficient for training an autoencoder, and thus we expanded the data for V-O binary compounds by the substitution of existing binary materials. To this end, we constructed a database of 10,981  $V_xO_y$  compounds to train and test the proposed generative model. We call this the VO dataset henceforth throughout this paper. Since this VO dataset consists of a small number of known  $V_xO_y$  compounds as well as virtual data generated via substitution, we denote the dataset containing only those generated via the substitution the VO virtual dataset. This will be used for validating the model. Other



**Table 1. Reconstruction Performances of the Proposed iMatGen Framework**

	Unit Cell						Basis (Å)
	a (Å)	b (Å)	c (Å)	$\alpha$ (°)	$\beta$ (°)	$\gamma$ (°)	
Intrinsic error	0.00	0.00	0.00	0.0	0.0	0.0	0.11
Training set (n = 9,792)	0.08	0.09	0.10	2.1	2.1	2.0	0.19
Test set (n = 1,189)	0.09	0.09	0.11	2.1	2.1	2.0	0.19
Reconstructed structures (n = 10,079)	0.07	0.08	0.09	1.9	1.9	1.9	0.13

Statistical accuracy of the proposed iMatGen measured by RMSE of atomic positions and cell parameters for the reconstructed materials relative to the reference input materials. Overall, 10,079 of 10,981 materials are fully reconstructed after autoencoding. The autoencoder is trained and tested with VO dataset.

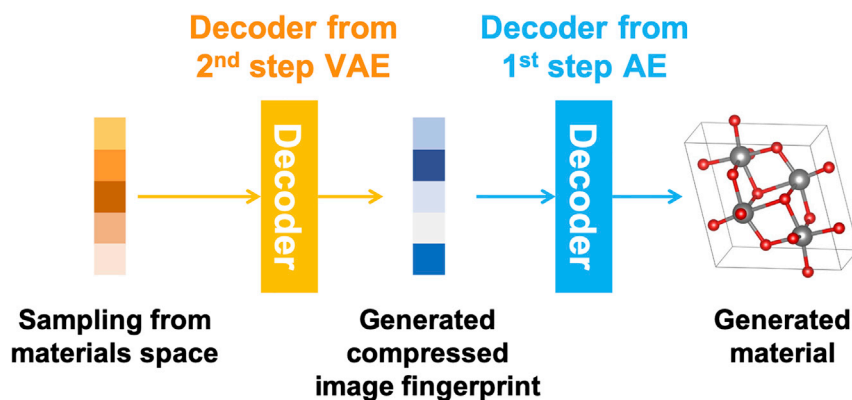
computational details for the generation of the datasets are described in [Experimental Procedures](#) and [Section S2](#).

### Performance of Autoencoders

[Table 1](#) shows statistical accuracy, in terms of root-mean-square error (RMSE), of the proposed iMatGen by comparing the inverse transformed lattice parameters and atomic positions with those of the input VO dataset structures (see [Figures S9](#) and [S10](#) for the detailed error distributions). The first row of [Table 1](#) corresponds to the intrinsic errors arising from the inverse transform of the input images themselves, that is, the inverse transform of the cell image using [Equation 2](#) and an image filter to locate the atomic positions from the basis images just after obtaining representation vectors. This analysis allows one to assess the accuracy of the representation-to-material back conversion. As discussed in [Experimental Procedures](#), the cell information is indeed fully recovered, while the determination of atomic positions from the basis images shows an intrinsic average error of 0.11 Å due to a finite grid spacing (0.23 Å) and the possibility of Gaussian peaks (atomic positions) appearing off the grid centers. After accounting for these small intrinsic errors of the inverse transform due to a finite grid spacing, the performance of iMatGen looks promising; overall, the cell lengths and angles are recovered within 0.1 Å and 2°, and the atomic positions within the unit cell are recovered within 0.2 Å on average. In terms of the number of structures that are fully reconstructed to the reference structure (after post-processing described in [Section S3.3](#)), 10,079 materials out of 10,981 were successfully reconstructed (reproducing simultaneously the correct cell and basis information within 0.1 Å and 2°, the last row in [Table 1](#)) by the present iMatGen.

### Materials Generation

To generate new materials using the proposed iMatGen, we first sample the materials vector ( $z$ ) from the materials latent space, and by applying the two decoders (one for MG and the other for image compression) consecutively, we obtain materials images for basis and cell in grid space. These images in grid space are then back-transformed to real-space atomic positions and cell parameters as described in [Experimental Procedures](#) (see [Figure 4](#) for scheme of materials generation). Because there is a possibility of generating blurred images from the materials-latent space, it may cause an unfavorable overlap between atomic positions after applying the image filter to the basis images. Therefore, post-processing was applied for the basis images in two steps as described in [Section S3.3](#); one for the position overlap under the periodic boundary condition for each element type (within V/O images) and the other for the position overlap between V and O (across V/O images). Conceptually, these post-processing steps are the application of our prior knowledge of atomic spacings to provide the user with the nearest meaningful structure corresponding to the MG output.



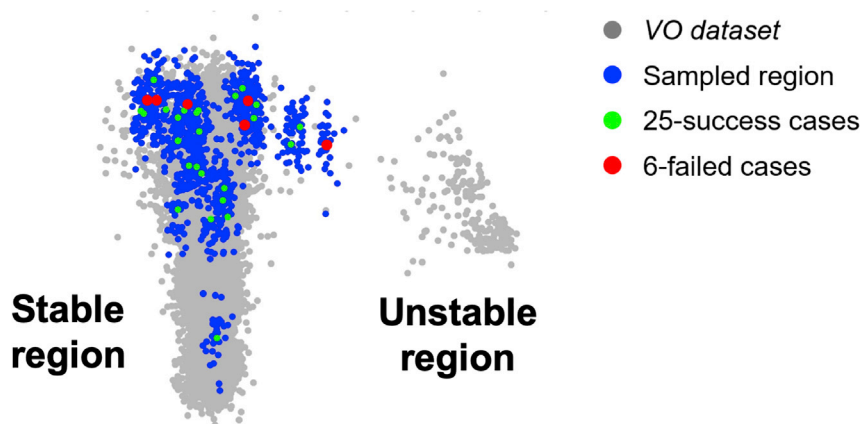
**Figure 4. Hierarchical Decoding for New Materials after Sampling in the Latent Space**

Before predicting new structures, first we have validated the proposed iMatGen scheme by attempting to rediscover experimentally known  $V_xO_y$  structures when these experimental structures were not included in the training data. To do so, we removed 31 known  $V_xO_y$  structures from the original training set of 10,981 data, and optimized the entire model set using the remaining set of 10,950 data points (denoted here as *VO virtual dataset* to distinguish it from *VO dataset*). The latent space of the latter MG is visualized in Figure 5 using the two main eigenvectors taken from principal component analysis. Two clusters can be identified, corresponding to stable and unstable materials, due to additional classification network mapping of the latent space to the binary stability labels. This clear separation of regions allows one to sample effectively the stable regions for the generation of stable materials, with most of the known materials indeed residing in the stable region. For each known  $V_xO_y$  structure, we sampled 30 random Gaussian distributed vectors around their latent vector with a fixed distance of 0.001, as shown in Figure 5 (blue circles). The latter constraints ensure that any sampled points do not correspond precisely to the 31 known  $V_xO_y$  structures (green circles in Figure 5). We then decoded all 930 ( $= 31 \times 30$ ) sampled points in latent space into the real-space crystal structures using the decoding and back-transformation procedures described above (Figure 4). To estimate the similarity between the reference and generated structures, we used a simple measure that accounts for the coordination environments of each atom in the unit cell and computes the vector distance between the two structures,<sup>47</sup> as shown in the Figure S8. A dissimilarity value of 0 thus means an identical match. We further inspected visually the lowest dissimilarity-value structures and optimized the structures using density functional theory to confirm that they are indeed identical to the reference structures. As shown in Figures 5 and S8, iMatGen successfully rediscovered 25 out of 31  $V_xO_y$  structures in MP, although none of these 31 structures were included in the training. The six unsuccessful cases correspond to a region with sparse data in the latent space for the machine to learn, and are hence not sampled reliably. It is notable that 8 out of 10 experimentally reported  $V_xO_y$  structures in the ICSD repository<sup>1,2</sup> (among the 31  $V_xO_y$  structures in MP) are fully generated and reconstructed. This result clearly demonstrates that the present iMatGen scheme has an ability to generate and decode into experimental materials that have been discovered over the course of the last century, and not just any computer-generated imaginary materials.

### Generation of Entirely New Structures

For the purpose of predicting new  $V_xO_y$  structures by utilizing all available data, we trained iMatGen using the entire *VO dataset* (10,981). The latent space of the resulting MG is shown in Figure 5 using the two main principal components, again showing





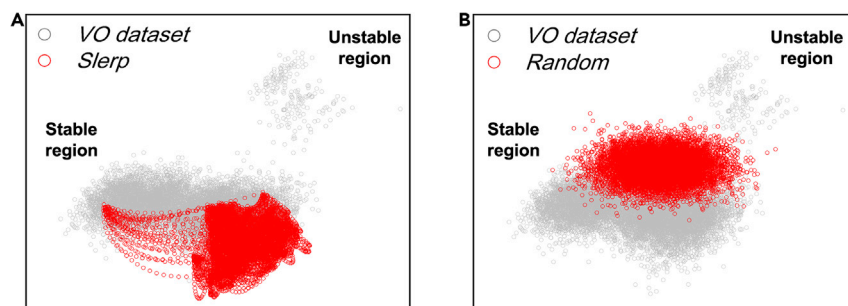
**Figure 5. Visualizations of the Learned Latent Space**

Latent space constructed using the VO virtual dataset (gray circle) and rediscovery of the known  $V_xO_y$  materials using the proposed iMatGen framework. A total of 930 latent vectors were randomly sampled (blue circles) around 31 known structures from Gaussian distribution (green and red circles). Twenty-five of 31 known  $V_xO_y$  structures in the Materials Project were successfully generated and rediscovered by the current iMatGen (green circles). See also Figure S8.

two distinct clusters of stable and unstable regions (the effect of inclusion of classification task is shown in Figure S3) where we further analyzed the latent space in terms of structural diversity of data as shown in Figure S4. We used two latent space sampling methods to generate new materials, as described in Section S3.1, one based on spherical linear interpolation (*Slerp*)<sup>48</sup> and another based on Gaussian random sampling (*Random*). In the case of *Slerp*, 31 known structures were used as initial guesses (two end points) and 19 structures were interpolated between each pair, yielding a total of  $31C_2 \times 19 = 8,835$  new materials. Similarly, we generated 10,000 structures using the Gaussian random sampling around the mean of the data distribution. Sampled data points in the *Slerp* and *Random* sampling methods are plotted in Figures 6A and 6B, respectively. Unlike *Slerp* that would mainly sample stable structures, in the case of *Random* sampling, all the sampled latent vectors are centered around a boundary region between the stable and unstable clusters, since the VAE learns the likelihood of data distribution using variational inference and tries to cover the whole data distribution based on a prior Gaussian (see Figure S2). One may randomly sample around the stable region only by shifting the center of the Gaussian, but we chose to randomly sample for the whole data to include the possibility of generating metastable and potentially more interesting structures at the boundary as well.

The newly generated structures were then post-processed and further refined as described in Section S3.2. In addition, we discarded the structures that do not match a V:O stoichiometry corresponding to a valid oxidation state of V ( $OS_V$ , i.e., between +2 and +5). The statistics of the latter results are summarized in Table 2. The *Slerp* sampling generates structures satisfying the valid  $OS_V$  condition with a much higher ratio ( $\sim 89\%$ ) than the *Random* sampling ( $\sim 17\%$ ), which can be attributed to the nature of *Slerp* that samples mainly in the stable region of the latent space (hence more reasonable structures), and that of *Random* that samples more at the boundary of stable and unstable regions. This shows the importance of sampling and complementary roles of *Slerp* and *Random*.

As shown in Table 2, the MP database has 25 unique compositions and 112 known  $V_xO_y$  structures, and the VO dataset includes 83 compositions and 10,981  $V_xO_y$



**Figure 6. Latent Vector Sampling for New Discovery**

The region of MG latent space trained with the VO dataset is shown in gray circles, and the new latent vectors sampled are shown in red circles using two methods: (A) *Spherical linear interpolation (Slerp)* and (B) *Random*. See also [Section S3.1](#) for sampling details.

structures as a result of V/O substitution for  $A_xB_y$ , among which 7 compositions and 31 structures correspond to existing  $V_xO_y$  materials from the MP database. For *Slerp* sampling, among the 8,835 structures generated, 7,845 structures meet the  $OS_V$  range, yielding 36 unique compositions in total (among which 16 are the known  $V_xO_y$  compositions from the MP database and 20 are from  $A_xB_y$ ). For *Random* sampling, among the 10,000 structures generated, 1,697 structures meet the  $OS_V$  range, yielding 47-unique compositions in total. The fact that the proposed iMatGen can generate  $V_xO_y$  compositions not in the MP database is partly due to the VO dataset, which includes the corresponding  $V_xO_y$  compositions by the substitution of  $A_xB_y$  binary materials with V and O as described above. However, it is interesting to note that many of the new and valid compositions generated by the current iMatGen are actually neither in the MP database nor the VO dataset (and quite stable as described below), demonstrating a possibility of our iMatGen to predict new and experimentally synthesizable vanadium oxide compositions. This subset of newly generated compositions corresponds to the column “New” in [Table 2](#). These results lead to a conclusion that the current iMatGen can effectively explore a large V-O binary chemical space for a new discovery.

To analyze the newly discovered materials further in detail, we selected the top four most frequently generated compositions for each sampling method for the new V-O compositions not in the MP database (i.e., the combined area of sections 3 and 4 in [Figure S11](#)):  $V_5O_8$ ,  $V_6O_7$ ,  $V_4O_5$ , and  $V_3O_4$  from *Slerp* ([Figure 7A](#)) and  $V_6O_7$ ,  $V_5O_6$ ,  $V_4O_5$ , and  $V_3O_4$  from *Random* ([Figure 7B](#)). It is interesting to note that  $V_6O_7$ ,  $V_4O_5$ , and  $V_3O_4$  commonly appear as the most frequently generated compositions in both sampling methods, and, in particular,  $V_6O_7$  is a completely new composition that is neither in the MP database nor the VO dataset. We then performed DFT geometry optimizations for these generated materials and estimated their phase stabilities using formation energies (see [Sections S2.2](#) and [S2.3](#)). As shown in [Figure 6](#), all of the newly generated structures of the four most frequent compositions are calculated to be metastable<sup>49,50</sup> with negative formation energies for both *Slerp* and *Random* sampling methods. In particular, a previous statistic that 80% of the experimentally known verified sulfides and oxides were within the energy above hull ( $E_{\text{hull}}$ )  $\leq 0.08$  eV/atom<sup>51</sup> suggests that 40 newly generated structures in this study could be considered synthesizable after full DFT calculation analysis.

To show examples of newly generated polymorphs of a given composition, we visualized the structures with the lowest  $E_{\text{hull}}$  values for each composition in [Figure 8](#). Interestingly, all of these lowest-energy polymorphs show porous structures

**Table 2. Database for Learning and the Newly Generated Data for Vanadium Oxides**

	MP and VO <sup>a</sup>	MP only <sup>b</sup>	VO only <sup>c</sup>	New <sup>d</sup>	Total
Database					
MP	31 (7)	81 (18)	–	–	112 (25)
VO dataset	31 (7)	–	10,950 (76)	–	10,981 (83)
Newly generated materials using iMatGen					
<i>Slerp</i>	6,099 (12)	35 (4)	1,409 (10)	302 (10)	7,845 (36)
<i>Random</i>	1,399 (11)	8 (4)	238 (20)	52 (12)	1,697 (47)

Data distribution (the number of materials) for the original database used in training (Materials Project [MP] and VO dataset [VO]) and for newly generated materials satisfying the valid oxidation state of V ( $OS_V$ ) range using spherical linear interpolation (*Slerp*) and random sampling (*Random*). The number of unique  $V_xO_y$  compositions are also listed in parentheses. For more detailed data distribution, see [Figure S11](#).

<sup>a</sup> $V_xO_y$  compositions known in both the MP and VO dataset.

<sup>b</sup> $V_xO_y$  compositions known in the MP but not in the VO dataset.

<sup>c</sup> $V_xO_y$  compositions known in the VO dataset but not in the MP.

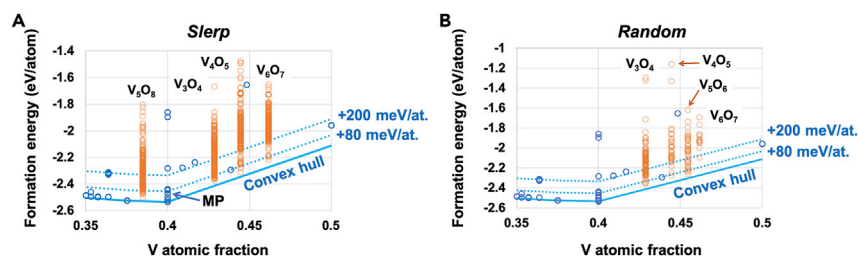
<sup>d</sup>New  $V_xO_y$  compositions neither in the MP nor in the VO dataset.

with non- $P1$  space group. Since many of these newly generated structures observed from [Figure 7](#) have compositions of the  $V_{n-1}O_n$ -type ( $n = 4-7$ ), we have examined whether these structures ([Figures 8C and 8D](#)) originate from the existing structures in VO dataset, in which one vanadium atom is removed from the unit cell. In other words, we added one V atom to  $V_{n-1}O_n$  and tested whether the  $V_nO_n$  structures existing in VO dataset can be recovered. Interestingly, indeed both cases ([Figures 8C and 8D](#)) converged to the same  $V_nO_n$  structures existing in VO dataset, suggesting that the present iMatGen has an ability to generate structures with long-range ordered vacancies that are similar in energy to the defect-free structures and thus likely to occur experimentally. While the training data did not include any V vacancy structures, the presence of V vacancies in V-O materials has been demonstrated both experimentally<sup>52</sup> and theoretically,<sup>53</sup> and the iMatGen discovery of these materials is likely enabled via training with a broad range of V:O stoichiometries from the  $A_xB_y$  structures. While crystal symmetries were not explicitly provided to the iMatGen, its structuring of the latent space from the symmetries of the training structures results in the prediction of new materials with various symmetries, which is demonstrated by the various space-group labels for the structures in [Figure 8](#).

The proposed model can easily be extended to generate higher-order compounds (i.e., ternary, quaternary, and so forth), since the input features can handle up to quinary compounds ([Section S1](#)). The most time-consuming procedure would be, as in many machine-learning-based approaches, to construct the data to learn from a significant combinatorial complexity embedded in higher-order compounds. Recently, many machine-learning-based force-field models have been proposed and increasingly used to describe potential energy surface of the given chemical space to compute energy, force, and stress of a material system.<sup>54–56</sup> One thus expects that utilizing such machine-learning-based force fields can be an effective way to construct the large materials database for higher-order compounds.

### Comparison with Genetic Algorithm

Since we applied the materials generative model to predicting new crystal structures as a proof of concept, we compared the results with a genetic-algorithm-based method<sup>17</sup> that has been widely and successfully used to predict new



**Figure 7. Phase Stabilities of the Newly Generated Materials**

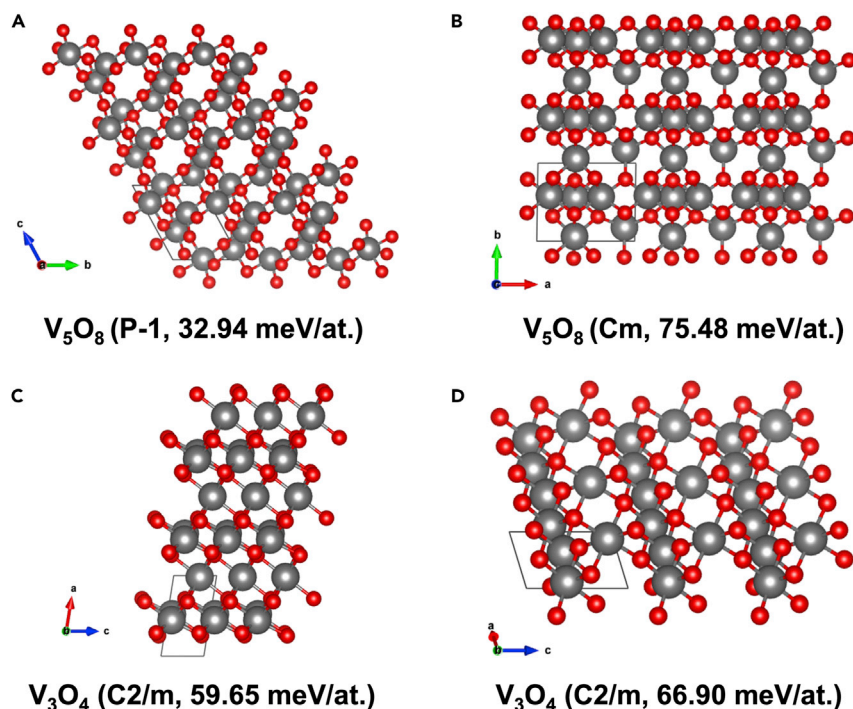
DFT calculated formation energies for the newly generated materials (orange circles) for the four most frequently generated compositions: (A)  $V_5O_8$ ,  $V_6O_7$ ,  $V_4O_5$ , and  $V_3O_4$  from spherical linear interpolation (*Slerp*) and (B)  $V_6O_7$ ,  $V_5O_6$ ,  $V_4O_5$ , and  $V_3O_4$  from *Random*. Blue circles correspond to the materials in the MP database, and the blue solid line is the minimum energy convex hull. Computational details related to the calculations of formation energies are shown in Sections S2.2 and S2.3.

(meta-)stable materials in various applications. Since these existing algorithms require as an input the preselected composition to be explored, here we chose  $V_3O_4$  (one and two formula units) and  $V_5O_8$ , since these compositions are discovered to have stable new structures from iMatGen. We used three different pressure conditions (0, 10, and 100 GPa) to predict new crystal structures, and at each pressure we performed a total of 400 structure generations. Thus, we performed a total of  $9 \times 400 = 3,600$  DFT geometry optimizations to predict new  $V_3O_4$  and  $V_5O_8$  materials using a genetic algorithm-based model. Surprisingly, the most stable structures predicted in a genetic algorithm-based model are consistent with the structures obtained using iMatGen; the lowest-energy  $V_5O_8$  and  $V_3O_4$  structures predicted by the genetic algorithm-based model were converged to the same  $V_5O_8$  structure predicted by iMatGen (Figure 8A) and  $V_3O_4$  structure in the VO dataset, respectively.

In terms of prediction efficiency, we note that the current iMatGen generated 52 unique new compositions and 40 relatively stable structures ( $E_{\text{hull}} \leq 80$  meV/atom) using 10,981 DFT calculations for the VO dataset construction (training) plus  $\sim 3,000$  DFT optimizations for most frequently generated compositions (post-processing), to total about 14,000 DFT optimizations. In contrast, the genetic algorithm-based model explored two new compositions,  $V_5O_8$  and  $V_3O_4$ , and discovered the same crystal structures as in iMatGen using 3,600 DFT calculations. Of course, the latter two preselected compositions were provided in this study by iMatGen but generally are a priori unknown in practice, and hence one would have to explore all 52 or more compositions (approximate total  $3,600 \times 26 = 93,600$  estimated DFT calculations) to make the same discovery using genetic algorithms. This demonstrates the potential advantage of GMs that can explore chemical compositional space effectively when there is no guide for the possible composition by learning the data distribution of known materials.

## Conclusions

We proposed a general framework for materials inverse design using hierarchical two-step autoencoder models (image compression followed by MG), called iMatGen, and as a proof of concept we applied the suggested scheme to discover new vanadium oxide materials. Our data-driven materials generation scheme is achieved by utilizing 3D image-based invertible input representation for crystal structures for both cell and basis information. We showed the invertibility of our representation scheme and the efficiency of the present iMatGen to generate



**Figure 8. New Polymorphs with Low Energy above a Convex Hull ( $E_{\text{hull}} \leq 80$  meV/atom)**  
Generated from the Present iMatGen

For each structure (A–D), the space group is also shown in parentheses.

experimental vanadium oxide materials by successfully rediscovering 26 out of 31 vanadium oxide materials in the MP when the model was trained without them. In addition to these materials in the database, more than 40 entirely new  $\text{V}_x\text{O}_y$  structures with energy above hull  $\leq 80$  meV/atom were generated. These discovered materials had not been identified by materials searches to date and include ordered-defect structures that are particularly promising for validation by experiments. Comparison of iMatGen with the genetic algorithm-based crystal structure prediction model also suggests an advantage of the proposed generative model to effectively explore the chemical compositional space and efficiently design unknown materials. The current framework is easily extendable for functional materials design with desired properties by modifying the training database and placing a separate property optimization task on the latent space,<sup>24</sup> in addition to utilizing the pretrained parameters obtained in this study.

## EXPERIMENTAL PROCEDURES

### Invertible Image-Based 3D Representations for Crystal Structures

We encode the crystal unit cell information using a single Gaussian function represented on a grid (Figure 2A), defined as

$$p_{(i,j,k)}^{\text{cell}} = \exp\left(-\frac{r_{(i,j,k)}^2}{2\sigma^2}\right), \text{ for } i, j, k \in [1, m], \quad (\text{Equation 1})$$

where  $p_{(i,j,k)}^{\text{cell}}$  is a Gaussian value at a grid index  $(i, j, k)$ ,  $r_{(i,j,k)}$  is a real-space distance between  $(i, j, k)$  and the center of the unit cell,  $\sigma$  is a Gaussian width, and  $m$  is the number of grid points in each direction. Due to a mapping of generally non-cubic crystal unit cell to a cubic grid, the resulting Gaussian image that represents the cell is generally distorted on a grid representation as shown in Figure 2A. The latter

distortion then makes the inverse transformation from images to the unit cell parameters (lengths and angles) fully invertible (as shown in Figure S7) using Equation 2,

$$r_{(i,j,k)} = \sqrt{-2\sigma^2 \ln p_{(i,j,k)}^{\text{cell}}}. \quad (\text{Equation 2})$$

To represent the basis (Figure 2B), all the atoms in the unit cell are put in a cubic box with length  $l$  Å in each direction and one Gaussian function is assigned for each atom before putting the resulting images on the  $n^3$  grid space, where  $n$  is the number of grid points in each dimension. The pixel value  $p_{(i,j,k),A}^{\text{basis}}$  at a grid index  $(i, j, k)$  is then the sum of contributions from all atoms of the atom type  $A$  ( $N_A$ ) as described by Equation 3,

$$p_{(i,j,k),A}^{\text{basis}} = \sum_{\alpha=1}^{N_A} \exp\left(-\frac{r_{(i,j,k),\alpha}^2}{2\sigma^2}\right), \text{ where } i, j, k \in [1, n], \quad (\text{Equation 3})$$

and  $r_{(i,j,k),\alpha}$  is a real-space distance between a grid index  $(i, j, k)$  and the position of atom  $\alpha$  for a given atom type  $A$ . The same value for  $\sigma$  as in Equation 2 can be used in Equation 3. This image-generation task for basis is repeated for each atom type of the material. For example, for a binary material consisting of two different elements V and O, one would generate two 3D images (one for V and one for O). Once the basis images are generated for each elemental type, the representation of a material is completed by combining the image of a unit cell and images of each atom type of a material. For the recovery of atom positions from the basis images during decoding, an image filter widely used in image-processing technologies can be used. In this work, an image filter implemented in SciPy<sup>57</sup> was used.

### Construction of VO Dataset

We generated 14,944 ( $=2 \times 7,472$ ) vanadium oxide structures by replacing the A or/and B in 7,472  $A_xB_y$ -type binary materials in MP to V or/and O. We then further removed the structures that do not satisfy the following two conditions for computational efficiency of cubic scaling grid generation: the maximum number of atoms in unit cell  $N_{\text{sites}} \leq 20$  and the maximum length of unit cell  $\leq 10$  Å. This yielded 10,981  $V_xO_y$  possible crystal structures of which 90% was used in training and 10% for testing the autoencoders in Figure 8. Among 10,981 entries, 31 materials (7 unique compositions) correspond to known vanadium oxides in MP. We denote this set of data as VO dataset. We used a cubic box of dimension  $l = 15$  Å to generate Gaussian images for the basis. To represent the cell and basis images in grid space, we used  $m = 32$  (cell) and  $n = 64$  (basis) in each dimension. The latter grid spacing for the basis corresponds to a resolution of 0.23 Å in real space. We note that there can be multiple computational unit cells definable for a given periodic material, and systematic inclusion of these additional data and even supercells in the training set would improve the prediction. In this study, we used the unit cells provided by MP without such data augmentation for computational efficiency.

### DATA AND CODE AVAILABILITY

The datasets used to train the model and the generated crystal structures are available at <https://github.com/kaist-amsg/imatgen.git>. Source codes and trained parameters are available at <https://github.com/kaist-amsg/imatgen.git>.

### SUPPLEMENTAL INFORMATION

Supplemental Information can be found online at <https://doi.org/10.1016/j.matt.2019.08.017>.



## ACKNOWLEDGMENTS

We acknowledge the support from the National Research Foundation of Korea (NRF-2017R1A2B3010176) and Korea Institute of Energy Technology Evaluation and Planning (KETEP-20188500000440) grants from the Korean Government, and a generous supercomputing time from Korea Institute of Science and Technology Information (KISTI). H.S.S. and J.M.G. are supported through the Office of Science of the U.S. Department of Energy under award no. DE-SC0004993. A.A.-G. thanks the Canada 150 Research Chairs Program, Natural Resources Canada, and the Vancouver Bush Faculty Fellowship Program for support. A.A.-G. acknowledges the generous support of Anders G. Frøseth.

## AUTHOR CONTRIBUTIONS

J.N., J.K., A.A.-G., and Y.J. designed the project. J.N. performed the machine-learning simulations, DFT calculations, and analyses. J.N. and Y.J. analyzed the results and wrote the manuscript. H.S.S. and J.M.G. assisted with data analysis and interpretation of the generated materials. B.S.-L. and A.A.-G. assisted with the machine-learning model construction. All authors contributed to the discussion and editing of the manuscript. Y.J. supervised the project.

## DECLARATION OF INTERESTS

The authors declare no competing interests.

Received: July 22, 2019

Revised: August 6, 2019

Accepted: August 17, 2019

Published: October 2, 2019

## REFERENCES

1. Belsky, A., Hellenbrandt, M., Karen, V.L., and Luksch, P. (2002). New developments in the Inorganic Crystal Structure Database (ICSD): accessibility in support of materials research and design. *Acta Crystallogr. B* 58, 364–369.
2. Allmann, R., and Hinek, R. (2007). The introduction of structure types into the inorganic crystal structure database ICSD. *Acta Crystallogr. A* 63, 412–417.
3. Sanchez-Lengeling, B., and Aspuru-Guzik, A. (2018). Inverse molecular design using machine learning: generative models for matter engineering. *Science* 361, 360–365.
4. Gu, G.H., Noh, J., Kim, I., and Jung, Y. (2019). Machine learning for renewable energy materials. *J. Mater. Chem. A* 7, 17096–17117.
5. Pyzer-Knapp, E.O., Li, K., and Aspuru-Guzik, A. (2015). Learning from the harvard clean energy project: the use of neural networks to accelerate materials discovery. *Adv. Funct. Mater.* 25, 6495–6502.
6. Shinde, A., Suram, S.K., Yan, Q., Zhou, L., Singh, A.K., Yu, J., Persson, K.A., Neaton, J.B., and Gregoire, J.M. (2017). Discovery of manganese-based solar fuel photoanodes via integration of electronic structure calculations, pourbaix stability modeling, and high-throughput experiments. *ACS Energy Lett.* 2, 2307–2312.
7. Davies, D.W., Butler, K.T., Skelton, J.M., Xie, C., Oganov, A.R., and Walsh, A. (2018). Computer-aided design of metal chalcogenide semiconductors: from chemical composition to crystal structure. *Chem. Sci.* 9, 1022–1030.
8. Yuan, R., Liu, Z., Balachandran, P.V., Xue, D., Zhou, Y., Ding, X., Sun, J., Xue, D., and Lookman, T. (2018). Accelerated discovery of large electrostrains in BaTiO<sub>3</sub>-based piezoelectrics using active learning. *Adv. Mater.* 30, <https://doi.org/10.1002/adma.201702884>.
9. Sun, W., Holder, A., Orvañanos, B., Arca, E., Zakutayev, A., Lany, S., and Ceder, G. (2017). Thermodynamic routes to novel metastable nitrogen-rich nitrides. *Chem. Mater.* 29, 6936–6946.
10. Sun, W., Bartel, C., Arca, E., Bauers, S., Matthews, B., Orvañanos, B., Chen, B.-R., Toney, M.F., Schelhas, L.T., Tumas, W., et al. (2019). A map of the inorganic ternary metal nitrides. *Nat. Mater.* 18, 732–739.
11. Sparks, T.D., Gaultois, M.W., Oliynyk, A., Brgoch, J., and Meredig, B. (2016). Data mining our way to the next generation of thermoelectrics. *Scr. Mater.* 111, 10–15.
12. Hinuma, Y., Hatakeyama, T., Kumagai, Y., Burton, L.A., Sato, H., Muraba, Y., Iimura, S., Hiramatsu, H., Tanaka, I., Hosono, H., et al. (2016). Discovery of earth-abundant nitride semiconductors by computational screening and high-pressure synthesis. *Nat. Commun.* 7, 11962.
13. Pandey, M., Vojvodic, A., Thygesen, K.S., and Jacobsen, K.W. (2015). Two-dimensional metal dichalcogenides and oxides for hydrogen evolution: a computational screening approach. *J. Phys. Chem. Lett.* 6, 1577–1585.
14. Fischer, C.C., Tibbetts, K.J., Morgan, D., and Ceder, G. (2006). Predicting crystal structure by merging data mining with quantum mechanics. *Nat. Mater.* 5, 641–646.
15. Hautier, G., Fischer, C.C., Jain, A., Mueller, T., and Ceder, G. (2010). Finding nature's missing ternary oxide compounds using machine learning and density functional theory. *ChemInform* 41, 3762–3767.
16. Ryan, K., Lengyel, J., and Shatruk, M. (2018). Crystal structure prediction via deep learning. *J. Am. Chem. Soc.* 140, 10158–10168.
17. Wang, Y., Lv, J., Zhu, L., and Ma, Y. (2012). CALYPSO: a method for crystal structure prediction. *Comput. Phys. Commun.* 183, 2063–2070.
18. Glass, C.W., Oganov, A.R., and Hansen, N. (2006). USPEX—evolutionary crystal structure prediction. *Comput. Phys. Commun.* 175, 713–720.
19. Lonie, D.C., and Zurek, E. (2011). XtalOpt: an open-source evolutionary algorithm for crystal

- structure prediction. *Comput. Phys. Commun.* **182**, 372–387.
20. Oliynyk, A.O., Antono, E., Sparks, T.D., Ghadbeigi, L., Gaultois, M.W., Meredig, B., and Mar, A. (2016). High-throughput machine-learning-driven synthesis of full-Heusler compounds. *Chem. Mater.* **28**, 7324–7331.
21. Oliynyk, A.O., Adutwum, L.A., Rudyk, B.W., Pisavadia, H., Lotfi, S., Hlukhyi, V., Harynuk, J.J., Mar, A., and Brgoch, J. (2017). Disentangling structural confusion through machine learning: structure prediction and polymorphism of equiatomic ternary phases ABC. *J. Am. Chem. Soc.* **139**, 17870–17881.
22. Legrain, F., Carrete, J., van Roekeghem, A., Madsen, G.K.H., and Mingo, N. (2018). Materials screening for the discovery of new Half-Heuslers: machine learning versus ab initio methods. *J. Phys. Chem. B* **122**, 625–632.
23. Kingma, D.P., and Welling, M. (2013). Auto-encoding variational bayes. *arXiv*, preprint [arXiv:1312.6114](https://arxiv.org/abs/1312.6114).
24. Gomez-Bombarelli, R., Wei, J.N., Duvenaud, D., Hernandez-Lobato, J.M., Sanchez-Lengeling, B., Sheberla, D., Aguilera-Iparraguirre, J., Hirzel, T.D., Adams, R.P., and Aspuru-Guzik, A. (2018). Automatic chemical design using a data-driven continuous representation of molecules. *ACS Cent. Sci.* **4**, 268–276.
25. Kadurin, A., Nikolenko, S., Khrabrov, K., Aliper, A., and Zhavoronkov, A. (2017). druGAN: an advanced generative adversarial autoencoder model for de novo generation of new molecules with desired molecular properties in silico. *Mol. Pharm.* **14**, 3098–3104.
26. Segler, M.H.S., Kogej, T., Tyrchan, C., and Waller, M.P. (2018). Generating focused molecule libraries for drug discovery with recurrent neural networks. *ACS Cent. Sci.* **4**, 120–131.
27. Guimaraes, G.L., Sanchez-Lengeling, B., Outeiral, C., Farias, P.L.C., and Aspuru-Guzik, A. (2017). Objective-reinforced generative adversarial networks (ORGAN) for sequence generation models. *arXiv*, preprint [arXiv:1705.10843](https://arxiv.org/abs/1705.10843).
28. Sanchez-Lengeling, B., Outeiral, C., Guimaraes, G.L., and Aspuru-Guzik, A. (2017). Optimizing distributions over molecular space. An objective-reinforced generative adversarial network for inverse-design chemistry (ORGANIC). *ChemRxiv*. <https://doi.org/10.26434/chemrxiv.5309668.v3>.
29. Putin, E., Asadulaev, A., Ivanenkov, Y., Aladinskiy, V., Sanchez-Lengeling, B., Aspuru-Guzik, A., and Zhavoronkov, A. (2018). Reinforced adversarial neural computer for de novo molecular design. *J. Chem. Inf. Model.* **58**, 1194–1204.
30. Irwin, J.J., and Shoichet, B.K. (2005). ZINC—a free database of commercially available compounds for virtual screening. *J. Chem. Inf. Model.* **45**, 177–182.
31. Gaulton, A., Hersey, A., Nowotka, M., Bento, A.P., Chambers, J., Mendez, D., Mutowo, P., Atkinson, F., Bellis, L.J., Cibrián-Uhalte, E., et al. (2016). The ChEMBL database in 2017. *Nucleic Acids Res.* **45**, D945–D954.
32. Blum, L.C., and Raymond, J.-L. (2009). 970 million druglike small molecules for virtual screening in the chemical universe database GDB-13. *J. Am. Chem. Soc.* **131**, 8732–8733.
33. Ruddigkeit, L., Van Deursen, R., Blum, L.C., and Raymond, J.-L. (2012). Enumeration of 166 billion organic small molecules in the chemical universe database GDB-17. *J. Chem. Inf. Model.* **52**, 2864–2875.
34. Kearnes, S., McCloskey, K., Berndl, M., Pande, V., and Riley, P. (2016). Molecular graph convolutions: moving beyond fingerprints. *J. Comput. Aided Mol. Des.* **30**, 595–608.
35. Jain, A., Ong, S.P., Hautier, G., Chen, W., Richards, W.D., Dacek, S., Cholia, S., Gunter, D., Skinner, D., Ceder, G., et al. (2013). Commentary: the Materials Project: a materials genome approach to accelerating materials innovation. *APL Mater.* **1**, 011002.
36. Kirklin, S., Saal, J.E., Meredig, B., Thompson, A., Doak, J.W., Aykol, M., Rühl, S., and Wolverton, C. (2015). The Open Quantum Materials Database (OQMD): assessing the accuracy of DFT formation energies. *NPJ Comput. Mater.* **1**, <https://doi.org/10.1038/npjcompumats.2015.10>.
37. Curtarolo, S., Setyawan, W., Hart, G.L.W., Jahnatek, M., Chepulskii, R.V., Taylor, R.H., Wang, S., Xue, J., Yang, K., Levy, O., et al. (2012). AFLOW: an automatic framework for high-throughput materials discovery. *Comput. Mater. Sci.* **58**, 218–226.
38. Xie, T., and Grossman, J.C. (2018). Crystal graph convolutional neural networks for an accurate and interpretable prediction of material properties. *Phys. Rev. Lett.* **120**, 145301.
39. Zhou, Q., Tang, P., Liu, S., Pan, J., Yan, Q., and Zhang, S.C. (2018). Learning atoms for materials discovery. *Proc. Natl. Acad. Sci. U S A* **115**, E6411–E6417.
40. Faber, F., Lindmaa, A., von Lilienfeld, O.A., and Armiento, R. (2015). Crystal structure representations for machine learning models of formation energies. *Int. J. Quan. Chem.* **115**, 1094–1101.
41. Schütt, K.T., Glawe, H., Brockherde, F., Sanna, A., Müller, K.R., and Gross, E.K.U. (2014). How to represent crystal structures for machine learning: towards fast prediction of electronic properties. *Phys. Rev. B* **89**, <https://doi.org/10.1103/PhysRevB.89.205118>.
42. Kajita, S., Ohba, N., Jinnouchi, R., and Asahi, R. (2017). A Universal 3D voxel descriptor for solid-state material informatics with deep convolutional neural networks. *Sci. Rep.* **7**, 16991.
43. Ryczko, K., Mills, K., Luchak, I., Homenick, C., and Tamblyn, I. (2018). Convolutional neural networks for atomistic systems. *Comput. Mater. Sci.* **149**, 134–142.
44. Jimenez, J., Skalic, M., Martinez-Rosell, G., and De Fabritiis, G. (2018). KDEEP: protein-ligand absolute binding affinity prediction via 3D-convolutional neural networks. *J. Chem. Inf. Model.* **58**, 287–296.
45. Goodfellow, I., Bengio, Y., Courville, A., and Bengio, Y. (2016). *Deep Learning, Vol. 1* (MIT Press).
46. Snoek, J., Adams, R., and Larochelle, H. (2012). On nonparametric guidance for learning autoencoder representations. Paper presented at: Artificial Intelligence and Statistics.
47. Zimmermann, N.E.R., Horton, M.K., Jain, A., and Haranczyk, M. (2017). Assessing local structure motifs using order parameters for Motif recognition, interstitial identification, and diffusion path characterization. *Front. Mater.* **4**, <https://doi.org/10.3389/fmats.2017.00034>.
48. White, T. (2016). Sampling generative networks. *arXiv*, preprint [arXiv:1609.04468](https://arxiv.org/abs/1609.04468).
49. Singh, A.K., Zhou, L., Shinde, A., Suram, S.K., Montoya, J.H., Winston, D., Gregoire, J.M., and Persson, K.A. (2017). Electrochemical stability of metastable materials. *Chem. Mater.* **29**, 10159–10167.
50. Zakutayev, A., Allen, A.J., Zhang, X., Vidal, J., Cui, Z., Lany, S., Yang, M., DiSalvo, F.J., and Ginley, D.S. (2014). Experimental synthesis and properties of metastable CuNbN<sub>2</sub> and theoretical extension to other ternary copper nitrides. *Chem. Mater.* **26**, 4970–4977.
51. Singh, A.K., Montoya, J.H., Gregoire, J.M., and Persson, K.A. (2019). Robust and synthesizable photocatalysts for CO<sub>2</sub> reduction: a data-driven materials discovery. *Nat. Commun.* **10**, 443.
52. Chamberland, B. (1973). New defect vanadium dioxide phases. *J. Solid State Chem.* **7**, 377–384.
53. Galy, J., and Miehle, G. (1999). Ab initio structures of (M2) and (M3) VO<sub>2</sub> high pressure phases. *Solid State Sci.* **1**, 433–448.
54. Artrith, N., Urban, A., and Ceder, G. (2017). Efficient and accurate machine-learning interpolation of atomic energies in compositions with many species. *Phys. Rev. B* **96**, <https://doi.org/10.1103/PhysRevB.96.014112>.
55. Behler, J., and Parrinello, M. (2007). Generalized neural-network representation of high-dimensional potential-energy surfaces. *Phys. Rev. Lett.* **98**, 146401.
56. Khorshidi, A., and Peterson, A.A. (2016). Amp: a modular approach to machine learning in atomistic simulations. *Comput. Phys. Commun.* **207**, 310–324.
57. Jones, E., Oliphant, T., and Peterson, P. (2014). *SciPy: open source scientific tools for Python*. <http://www.scipy.org/>.