

Crystal Structure Generation with Autoregressive Large Language Modeling

Luis M. Antunes^{*1}, Keith T. Butler², and Ricardo Grau-Crespo^{*1}

¹*Department of Chemistry, University of Reading, Whiteknights, Reading RG6 6DX, United Kingdom.
l.m.antunes@pgr.reading.ac.uk; r.grau-crespo@reading.ac.uk*

²*Department of Chemistry, University College London, WC1H 0AJ, United Kingdom.*

Abstract

The generation of plausible crystal structures is often the first step in predicting the structure and properties of a material from its chemical composition. Quickly generating and predicting inorganic crystal structures is important for the discovery of new materials, which can target applications such as energy or electronic devices. However, most current methods for crystal structure prediction are computationally expensive, slowing the pace of innovation. Seeding structure prediction algorithms with quality generated candidates can overcome a major bottleneck. Here, we introduce CrystaLLM, a methodology for the versatile generation of crystal structures, based on the autoregressive large language modeling (LLM) of the Crystallographic Information File (CIF) format. Trained on millions of CIF files, CrystaLLM focuses on modeling crystal structures through text. CrystaLLM can produce plausible crystal structures for a wide range of inorganic compounds unseen in training, as demonstrated by *ab initio* simulations. The integration with predictors of formation energy permits the use of a Monte Carlo Tree Search algorithm to improve the generation of meaningful structures. Our approach challenges conventional representations of crystals, and demonstrates the potential of LLMs for learning effective ‘world models’ of crystal chemistry, which will lead to accelerated discovery and innovation in materials science.

1 Introduction

The *in silico* search for new materials often involves the exploration of a space of compositions in a chemical system, and the investigation of various predicted structural phases in that space (see [1], [2] and [3] for examples). To elucidate the structures of unknown materials, a Crystal Structure Prediction (CSP) approach is often employed, which attempts to derive the ground state crystal structure for a given chemical composition under specific physical conditions. [4] CSP approaches are relatively computationally expensive, typically involving *ab initio* techniques. [5] They often begin with the generation of candidate structures. Examples are the AIRSS [6, 7] and USPEX [8] approaches. Initializing the search space with sensible structures increases the likelihood of success, and decreases the

amount of computation required. It is therefore expected that effective crystal structure generation tools would help accelerate the prediction of structures using CSP methods.

Increasingly, techniques from machine learning (ML) and data science are being used to solve problems in materials science. [9,10] In particular, generative modeling approaches based on autoencoder architectures and generative adversarial networks (GANs) [11] have been used to generate crystal structures. [12–15] Indeed, generative modeling has become commonplace, an outcome catalyzed by astounding advancements in the computational generation of images, audio and natural language over the last several years. [16] The Large Language Model (LLM), backed by the Transformer architecture [17], is the approach behind state-of-the-art performance on natural language processing tasks. This approach begins with a generative pre-training step, which is autoregressive in nature, involving the unsupervised task of predicting the next token given a sequence of preceding tokens. [18] When such models are scaled to billions of parameters, their effectiveness becomes quite remarkable, as tools such as ChatGPT [19] demonstrate.

LLMs have recently been used in the context of materials science. [20–25] These attempts have been focused on using existing and publicly accessible LLMs, training and tuning LLMs for natural language generation tasks involving chemical subject matter, or training LLMs on a corpus of expanded chemical compositions for the purposes of generating unseen compositions. However, the potential of training LLMs on textual representations of crystal structures has not been considered. A sole exception is a recent pre-print by Flam-Shepherd and Aspuru-Guzik, where the idea of generating the structures of molecules, materials, and protein binding sites with LLMs has been preliminarily explored [26].

Here, we report the first LLM specifically designed for crystal generation. This model is distinctively trained on textual representations of inorganic crystal structures, specifically in the Crystallographic Information File (CIF) format [27], instead of relying solely on natural language corpora, or chemical compositions alone. The motivation for this approach originates from two conjectures: The first states that a sequence of symbols (i.e. tokens) is an appropriate representation modality for many predictive tasks, including those involving chemical structure. The idea of representing any domain with a sequence of tokens may at first seem counter-intuitive. However, consider that even images can be represented this way, and be subject to the autoregressive language modeling of pixels [28]. This challenges the notion that domain-specific representations, such as graphs for chemical structure [29], are necessary for superior performance. The second conjecture states that LLMs learn more than simply *surface statistics* and the conditional probability distribution of tokens. Indeed, autoregressive pre-training involving next-token prediction may result in learning an effective *world model*: an internalized causal model of the processes generating the target phenomena. A model which simply learns spurious correlations in the data is less desirable, as it may have greater difficulty in generalizing beyond the training distribution. Recent studies have demonstrated that LLMs trained on sequences of board game play (e.g. chess and Othello) do indeed track the state of the board, and probes of the internal activations of the model reveal the existence of representations of various abstract concepts specific to the domain. [30,31] We therefore asked whether a model trained to predict the 3-dimensional coordinates of atoms, digit-by-digit, could learn the chemistry implicit in crystal structures, and generate unseen structures, borrowing from its model of the world of atoms.

As such, we herein describe the CrystaLLM model, a tool for crystal structure generation trained on an extensive corpus of CIF files representing the structures of millions of inorganic solid-state materials. Unlike small molecule organic compounds, the generative modeling of inorganic crystals presents unique challenges: the structures are complex and periodic, are not readily described by simple graphs, are imbued with different forms of symmetry, and can be constructed from more than 100 different elements. Even so, the model is capable of reliably generating correct CIF syntax and physically plausible crystal structures for many classes of inorganic compounds. Moreover, we demonstrate how sampling from the model can be improved using the Monte Carlo Tree Search (MCTS) algorithm [32, 33] together with a pre-trained graph-based neural network predictor of formation energy.

2 Results

CrystaLLM is a Transformer-based, decoder-only language model of the CIF file format, trained autoregressively on a corpus of millions of CIF files (Figure 1a). Rather than training on structural representations derived from the CIF files, the model is directly trained on the standardized and tokenized text contents of the CIF files. During training, the model is given a sequence of tokens from the corpus of CIF files, and is tasked with predicting the tokens which follow each of the given tokens. Once the model is trained, it can be used to generate new CIF files, conditioned on some starting sequence of tokens. Generating a CIF file involves repeatedly sampling tokens from the model, conditioning on the accumulated generated content, until a terminating condition is reached (Figure 1b).

To assess the ability of the model to generate structures, a test set of approximately 10,000 randomly chosen CIF files is withheld from a training set of approximately 2.2 million CIF files, and the model is tasked with generating CIF files beginning from prompts constructed from the test set. Moreover, we assemble what we call a *challenge set*, which consists of 70 structures, 58 of which were obtained from the recent literature, and were not in the training set. The remaining 12 structures are from the training set, and are included as representatives of different structural classes. They serve to assess the model’s ability to recover what it has seen in training, and as a means of comparing the model’s generations of seen and unseen structures. (Supplementary Table 1 contains the full list of the challenge set compounds, and their sources.) The permutative nature of the dataset, with many structures having been derived by substituting atoms into pre-defined templates, results in a test set with the potential for some structures to closely resemble those of the training set. The challenge set provides a source of structures that are guaranteed to have been produced through a different process. Moreover, the challenge set constitutes a manageable set of compounds that reflects a variety of solid-state structural classes, allowing for a fine-grained picture of the model’s capabilities. The test set, on the other hand, is better suited for a bulk assessment, and originates from the same distribution as the training set.

The following terminology is used in the remainder of this article: A *formula*, *reduced formula*, or *reduced composition*, refers to the empirical formula, or formula unit, which is the simplest, whole-number ratio of atoms in the compound. An example of a formula is Ba_2MnCr . A *cell composition* is a chemical formula referring to the total number of

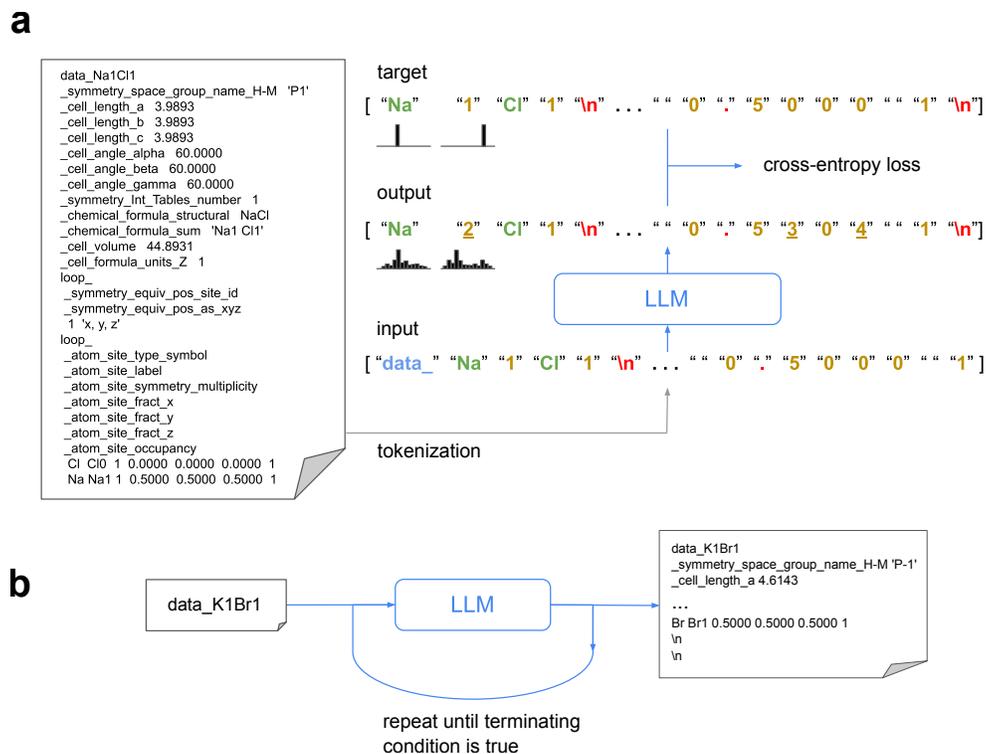


Figure 1: **a** Core concepts in training a Large Language Model of CIF files: A CIF file (left) is converted into a sequence of symbols, through tokenization. The sequence is processed by the model, which produces a list of probability distributions over the vocabulary, for each corresponding symbol in the input. The resulting predicted probability distributions are evaluated against the target distributions (which contain the entire probability mass on the correct subsequent token), using the cross-entropy loss metric. The target tokens are the input tokens shifted one spot to the left, as the objective is to predict the next token given a sequence of preceding tokens. The tokens are categorized as CIF tags (blue), atoms (green), numeric digits (gold), and punctuation (red). Output tokens (not actually sampled during training) represent the tokens assigned the highest probability by the model. Underlined tokens represent predicted distributions assigning a relatively low probability to the correct next token. **b** Generation of a CIF file: First, a prompt is constructed by concatenating the symbol `data_` with the desired cell composition, which is then tokenized and processed by the model. Next, a token is sampled from the predicted distribution for the upcoming token in the sequence. Finally, the sampled token is added to the accumulating contents of the CIF file. This procedure continues iteratively until a predefined terminating condition is met (e.g. two consecutive newline tokens are sampled).

atoms of each type in the unit cell of a crystal. It represents the chemical formula of the compound as it would appear in the crystallographic unit cell, which might contain Z formula units. An example of a cell composition is $\text{Ba}_6\text{Mn}_3\text{Cr}_3$, with a Z of 3.

2.1 Training and Learned Representations

Training consists of iteratively sampling sequences of tokens, of fixed length, and adjusting the model’s parameters so that it becomes progressively better at predicting which token should follow a preceding sequence. (See the Methods, and Supplementary Note 2, for more information on the model architecture and training.) Since it has been observed that LLM performance improves as the number of model parameters is increased [34], we train a *small* model, consisting of 25 million parameters, and a *large* model, consisting of 200 million parameters.

To monitor the progress of training, we withhold a validation set that constitutes 10% of the set held-out for training. Over the course of training, the model continues to improve in terms of its total cross-entropy loss on the validation set, even after 90,000 iterations (see Supplementary Figure 2). We note, however, that improvements appear to become smaller with more training time.

As a consequence of the model’s architecture, each token in a processed sequence is mapped to a distinct learned vector representation using an embedding table, whose parameters are adjusted during training. The result is that, through autoregressive training, distributed representations are learned for each symbol in the vocabulary. The vocabulary consists of symbols for atoms, space groups, and numeric digits. (See Supplementary Note 1 for a detailed description of the vocabulary and the tokenization procedure.) The training process appears to result in sensible representations of these various symbols. Plots of dimensionally-reduced atom and space group vectors demonstrate a logical structure, where similar entities cluster together, indicating that intrinsic properties and relationships are captured. (See Supplementary Figure 3 for plots of the learned atom vectors, and Supplementary Figure 4 for a plot of the learned space group vectors.) Moreover, examination of the learned numeric digit vectors reveals that numerical relationships are captured in the representations, as measurements of cosine and Euclidean distances between the learned digit vectors demonstrate a logical spatial relationship. (See Supplementary Figure 5.) While not explored further in this work, we note that distributed representations of chemical entities, such as atoms, are useful for the prediction of materials properties [35].

2.2 Generalizing to Unseen Structures

To evaluate the ability of the model to generate an unseen structure, the model is prompted with the structure’s cell composition, and allowed to generate up to 3,000 tokens. The prompt includes the first line of the CIF file, which consists of the data block header, containing the cell composition of the structure. Subsequently, the model is prompted with both the structure’s cell composition and space group, and again allowed to generate up to 3,000 tokens. The prompt includes the first several lines of the pre-processed CIF file, up to the line containing the specification of the space group. Prompting the model with both the cell composition and space group allows us to assess

how reliant the model is on the space group. This process is repeated for all CIF files of the held-out test set (10,286 in total).

The generated CIF files are then assessed for correctness and quality. Any syntactically incorrect CIF files are declared invalid. Syntactically correct CIF files are subjected to further analysis, and are considered to be valid only if specific criteria are met, such as being consistent in terms of generated structure and declared space group, and having reasonable bond lengths (see Supplementary Note 3 for further details on the validation of generated CIF files). The results of evaluating the generation of the CIF files of the test set using the small model are presented in Table 1.

Table 1: Performance of the small model on the held-out test set. The percentages represent the fraction of test set compounds which meet the corresponding criteria. For example, the first row represents the percentage of test set compounds where the declared space group in the generated CIF file is consistent with the generated structure. Valid generated length refers to the length of a valid generated CIF file in terms of the number of tokens.

	No Space Group	With Space Group
Space Group Consistent	98.8%	99.1%
Atom Site Multiplicity Consistent	99.4%	99.4%
Bond Length Reasonableness Score	0.9878 ± 0.0686	0.9878 ± 0.0671
Bond Lengths Reasonable	94.6%	94.6%
Valid	93.8%	94.0%
Longest Valid Generated Length	1145	970
Average Valid Generated Length	331.885 ± 42.567	339.002 ± 41.361

The CIF files generated by prompting the model with the cell composition and space group were compared to the corresponding CIF files of the test set using a structure matching algorithm. The fraction of matching structures is presented in Table 2.

Table 2: Structure matching results for the test set when the space group is included in the prompt. The *Reduced Unseen* column represents the results for formulas that were not seen in training with any *Z*.

	All	Reduced Unseen
At least 1 match within 3 attempts	88.1%	86.3%
All 3 attempts matching	67.4%	70.0%
Matched on 1st attempt	78.4%	78.7%

We further examined how closely the generated cell parameters resembled the actual cell parameters, for the cases where there was a structural match. We took the first matching structure for samples that had at least one generated structure matching the test set structure, and measured the R^2 and mean absolute error (MAE) for the true versus generated cell lengths, the true versus generated (i.e. printed) volume, and the implied (from cell parameters) versus generated volume. The results are presented in Figure 2.

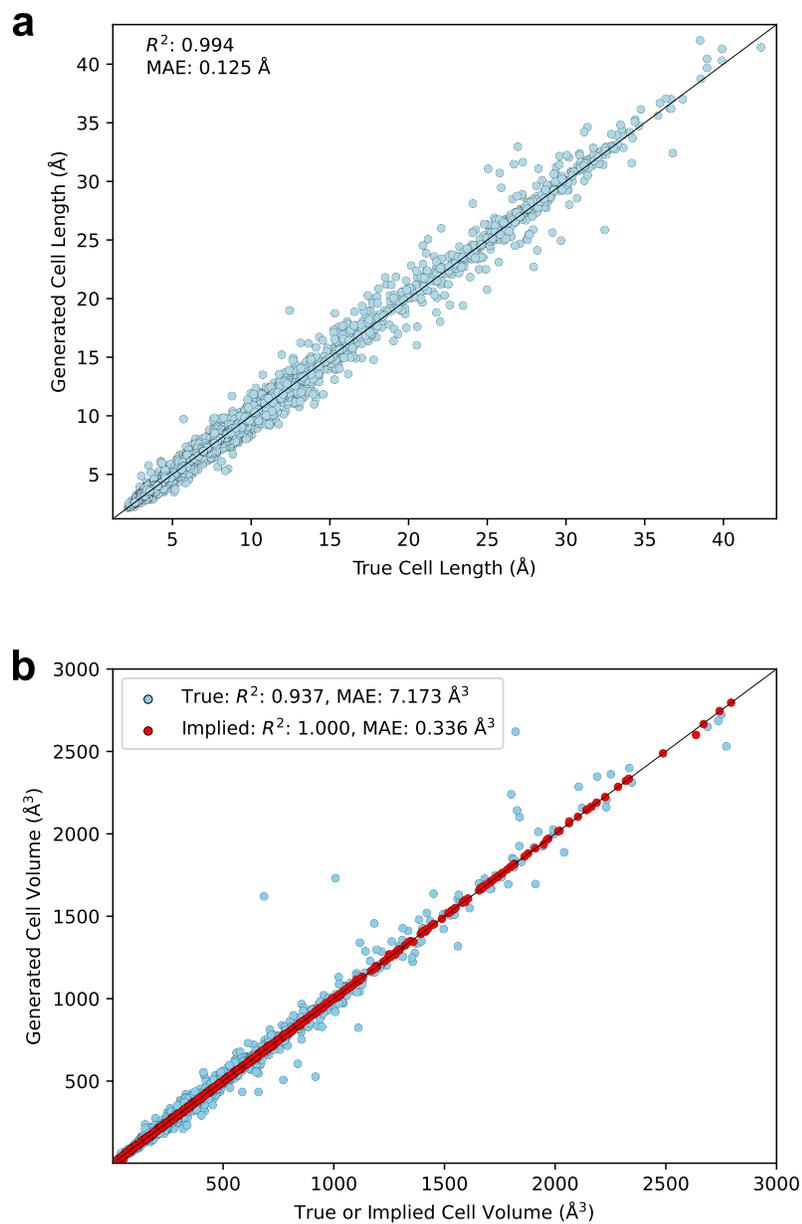


Figure 2: **a** The generated cell lengths for matching structures of the test set vs. the true cell lengths, when space group is included. **b** The generated cell volumes for matching structures of the test set vs. either the true cell volumes, or the cell volumes implied from the generated cell parameters, when space group is included.

To further assess the model’s ability to generalize to unseen structures, we prompted the model with the cell compositions of the challenge set. The challenge set contains 58 structures not seen in training. These structures were all manually sourced from the recent literature, and represent experimentally characterized materials. Crucially, these compounds originate through a process different from the process which generated the training set (namely, a high-throughput DFT analysis of hypothetical materials). They also represent a variety of different structural classes, such as intermetallics, silicates, sulfides and selenides, borates, phosphates, carbonates, and complex mixed-anion compounds.

Both the small and large models were prompted with the cell compositions of the challenge set, both with and without the space group. A total of 100 attempts were made to generate a structure from the given cell composition (and optionally space group). We record the successful generation rate, representing the fraction of compounds where at least one valid CIF file was generated in the 100 attempts, and the true match rate, representing the fraction of compounds where there was a structural match between a valid generated structure and the true structure reported in the literature. The results are presented in Table 3 and Supplementary Tables 2 to 5.

Table 3: Results of the small and large models on the challenge set, both with a space group (‘s.g.’) and without. The first row represents the percentage of cases where the model was able to generate a valid structure within 100 attempts. The second row represents the percentage of cases where a generated structure matched the true structure, for the compounds seen in training. The last row represents the percentage of cases where a generated structure matched the true structure, for unseen compounds only.

	Small model		Large model	
	no s.g.	with s.g.	no s.g.	with s.g.
Successful Generation Rate	85.7%	88.6%	87.1%	91.4%
Match Rate (Seen)	50.0%	50.0%	83.3%	83.3%
Match Rate (Unseen)	25.9%	34.5%	37.9%	41.4%

The results in Table 3 indicate that inclusion of the space group in the prompt increases the likelihood of generating a valid structure, and of generating a match with the true structure. The large model appears to be superior to the small model in all categories. While the models can recover the reported structure more often when the structure was seen in training, it is noteworthy that they are able to generate unseen structures which match the reported structure in up to 40% of the cases.

2.3 Comparison with Other ML-based Approaches

Generative models of materials based on advanced ML techniques have recently been developed. CDVAE [13] and DiffCSP [36] are two such examples. Both models use diffusion-based methods for generating materials, with DiffCSP focusing on crystal structure prediction through an equivariant diffusion process, while CDVAE uses a diffusion-based approach within a variational autoencoder framework for generating periodic materials. We compare CrystaLLM to these models on four benchmarks: Perov-5 [37, 38], Carbon-24 [39], MP-20 [40], and MPTS-52 [41]. The Perov-5 dataset consists of 18,928

perovskites, Carbon-24 consists of 10,153 carbon allotropes, MP-20 consists of 45,231 stable inorganic materials of various classes, while MPTS-52 consists of 40,476 various inorganic materials. MPTS-52 is by far the most complex dataset, with up to 52 atoms in the unit cells of the constituent structures.

The benchmark datasets have each been split into training, validation and test sets. Models are trained on the training set, and then are used to generate 20 structures for each of the cell compositions of the test set. The models are evaluated in terms of the *match rate*, which is the fraction of compositions for which the true structure was generated within n attempts (we tried $n=1$ and 20), and the average root mean squared error (RMSE) of the closest candidate for each test set structure. The results are presented in Table 4.

We present results for three different versions of CrystaLLM. Versions *a* and *b* are trained on the benchmark data only and differ in the size of the model used. Version *c* is trained on the full 2.3M training points minus the test set of MPTS-52 and is included to demonstrate how the results improve with the size of training data, but is not directly comparable to other models due to the different training data sets.

Table 4: Benchmark CSP results. Numbers in bold indicate the best $n=20$ result, while italicized numbers represents the best $n=1$ result, amongst the models trained only on the benchmark training sets, where n represents the number of samples generated for each structure of the benchmark test set. *a* Results for the small model architecture trained only on the benchmark training sets. *b* Results for the large model architecture trained only on the benchmark training sets. *c* Results for the small model architecture trained on the original 2.3M-structure dataset without the structures of the MPTS-52 validation or test sets. The CDVAE and DiffCSP results are taken from Jiao *et al.* [36].

Model	n	Perov-5		Carbon-24		MP-20		MPTS-52	
		Match Rate	RMSE						
CDVAE	1	45.31	0.1138	17.09	0.2969	33.90	0.1045	5.34	0.2106
CDVAE	20	88.51	0.0464	88.37	0.2286	66.95	0.1026	20.79	0.2085
DiffCSP	1	<i>52.02</i>	<i>0.0760</i>	17.54	0.2759	51.49	0.0631	12.19	0.1786
DiffCSP	20	98.60	0.0128	88.47	0.2192	77.93	0.0492	34.02	0.1749
CrystaLLM ^a	1	47.95	0.0966	<i>21.13</i>	<i>0.1687</i>	55.85	0.0437	17.47	0.1113
CrystaLLM ^a	20	98.26	0.0236	83.60	0.1523	75.14	0.0395	32.98	0.1197
CrystaLLM ^b	1	46.10	0.0953	20.25	0.1761	58.70	<i>0.0408</i>	<i>19.21</i>	<i>0.1110</i>
CrystaLLM ^b	20	97.60	0.0249	85.17	0.1514	73.97	0.0349	33.75	0.1059
CrystaLLM ^c	1	-	-	-	-	-	-	28.30	0.0850
CrystaLLM ^c	20	-	-	-	-	-	-	47.45	0.0780

CrystaLLM outperforms DiffCSP on three out of four benchmarks in terms of RMSE for both $n=20$ and $n=1$, and in terms of match rate when constrained to only a single generation attempt. This is achieved even in the most challenging of the benchmarks, MPTS-52, which contains structures with larger unit cells and more atoms.

CrystaLLM has other important advantages when compared to the other models. Notably, it supports the conditioning of structure generation on specific symmetry space groups, a capability unique to CrystaLLM. The flexibility of its inputs suggests that CrystaLLM may be conditioned on other properties of the structure as well, including those not traditionally included in the CIF format. Moreover, as a large language model, it can leverage the established practice of fine-tuning, allowing the pre-trained model to be adapted for the prediction of materials properties. There is far less precedent in

fine-tuning models based on diffusion and variational autoencoder architectures for tasks involving regression or classification.

2.4 Examples of Generated Structures

To further examine the model’s ability to generalize to unseen scenarios, we prompted the model with various formulas, and examined its output. The results are presented in Figure 3.

An example of the model generalizing to a formula that had been seen in training, but with different space groups, is presented in Figure 3a. The formula, Ba_2MnCr , was in the held-out test set, with the $R\bar{3}m$ space group. That combination of formula and space group had not been seen in training. The model generated a structure matching the one in the test set on the first attempt, when the space group was provided.

The model also demonstrated the ability to generate plausible structures for formulas not seen in training with any Z . An example is the quaternary compound CsCuTePt . This compound was not in the training set, but was in the held-out test set (with $Z=4$). The model generated a structure matching the one in the test set, in the $F\bar{4}3m$ space group, on the third attempt when the space group was provided. The generated structure is presented in Figure 3b.

Finally, in Figure 3c is the generated structure of YbMn_6Sn_6 [42], an example of the model generalizing to structural motifs with atoms not seen in training. This formula was not seen in training for any Z , and was not in the held-out test set. However, ZrMn_6Sn_6 was seen in training, in the $P6/mmm$ space group. The model generated a structure in the same space group on the first attempt, without the space group being provided. The generated structure matched the ZrMn_6Sn_6 structure, with Yb substituted for Zr, and with cell parameters and atomic coordinates adjusted accordingly. This demonstrates the model performing a structure prediction by analogy procedure, as commonly used by materials scientists for discovery [43, 44], despite never having been provided with the procedure to do this.

2.4.1 Rutiles

Rutiles are a class of binary compounds that adopt a tetragonal unit cell, in the $P4_2/mnm$ space group ($Z=2$), as is seen in TiO_2 , from which this class of materials adopts its name. The general formula for rutile oxides is MO_2 , where M is a metallic species in the +4 oxidation state. Rutile fluorides are also known, where the metal is in the +2 oxidation state.

The model’s training dataset consisted of essentially all of the rutiles one might expect to be able to find in nature. Therefore, to test the model’s ability to generate unseen rutiles, we requested the generation of theoretically possible, but unlikely compounds, such as AuO_2 . With gold in a highly unlikely +4 oxidation state, AuO_2 is not expected to be formed under most conditions. However, the model was able to imagine what the structure of such a compound might be (when the space group is provided). While TiO_2 has cell parameters $a=4.594\text{\AA}$, $c=2.959\text{\AA}$, the generated rutile gold variant has $a=4.838\text{\AA}$ $c=3.429\text{\AA}$, reflecting the increased volume occupied by the larger gold atoms (Figure 3d).

2.4.2 Spinel

Spinel is a group of ternary compounds with general formula AB_2X_4 . The most common combination of elements in spinel is one where A is a cation in the +2 oxidation state, B is a cation in the +3 oxidation state, and X, normally a chalcogen, is a -2 anion. Spinel forms cubic close-packed structures, with eight tetrahedral, and four octahedral sites, normally in the $Fd\bar{3}m$ space group.

To explore the model’s ability to generate unseen spinels, we selected the thiospinel Sm_2BS_4 , which was absent from both the training and test sets. The model was able to generate the expected spinel structure when the cell composition and space group were provided (Figure 3e). During training, the model encountered a number of different oxy-, thio-, and selenospinel, and this likely contributed to its ability to generate this compound.

2.4.3 Elpasolites

Elpasolites are quaternary compounds with the general formula ABC_2X_6 . The A and C species are typically alkali metal cations in the +1 oxidation state, B is usually a transition metal cation in the +3 oxidation state, and X is a halogen anion. The elpasolites are often referred to as “double perovskites”, since their structures are related to perovskites by the doubling of their unit cell dimensions, and the replacement of the M^{2+} cation with alternating M^+ and M^{3+} cations. Elpasolites crystallize in the $Fm\bar{3}m$ space group, and are the most common quaternary crystal system reported in the Inorganic Crystal Structure Database (ICSD) [45]. We wondered if the CrystaLLM model could generate elpasolites not seen during training.

We selected an elpasolite from the held-out test, that was not seen in training: the fluoride KRb_2TiF_6 . The model was able to generate the correct elpasolite structure when the cell composition and space group was provided (Figure 3f).

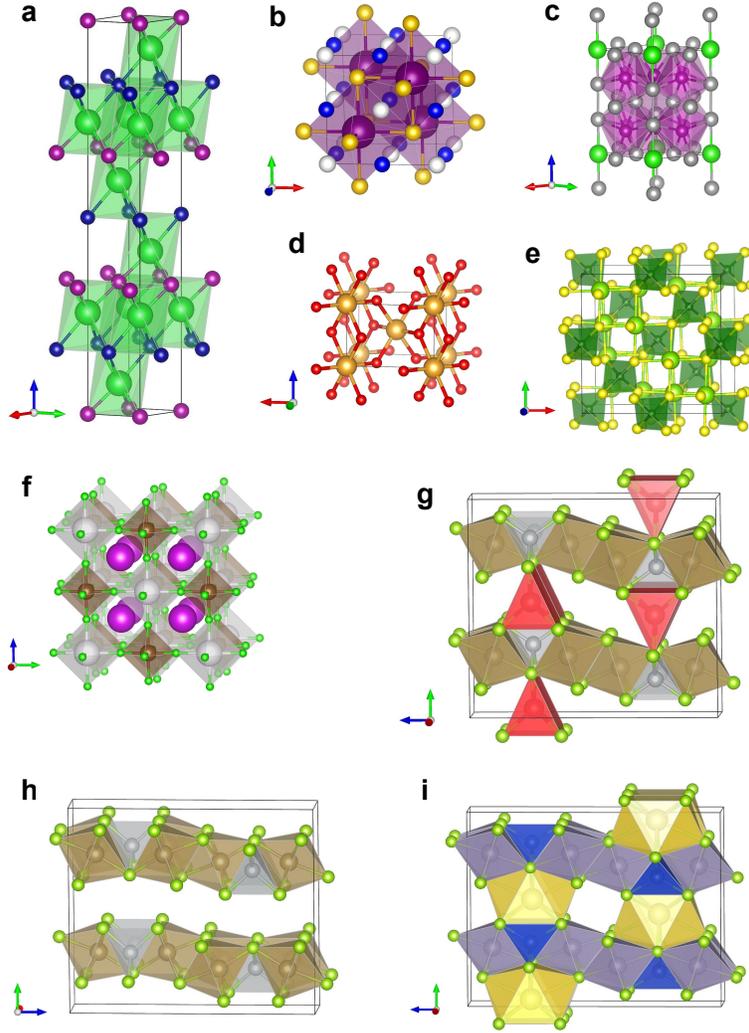


Figure 3: The generated structures of various inorganic compounds. **a** Ba_2MnCr . Cell parameters: a, b, c : 3.778 Å, 27.503 Å, α, β, γ : 90.0°, 90.0°, 120.0°. Color scheme: Ba: green, Mn: purple, Cr: blue. **b** CsCuTePt . Cell parameters: a, b, c : 7.153 Å, α, β, γ : 90.0°. Color scheme: Cs: purple, Cu: blue, Te: gold, Pt: white. **c** YbMn_6Sn_6 . Cell parameters: a, b, c : 5.488 Å, 8.832 Å, α, β, γ : 90.0°, 90.0°, 120.0°. ZrMn_6Sn_6 , in the training set, possessed the same structure, but with the following cell parameters: a, b, c : 5.364 Å, 8.933 Å, α, β, γ : 90.0°, 90.0°, 120.0°. Color scheme: Yb: green, Mn: magenta, Sn: grey. **d** AuO_2 . Cell parameters: a, b, c : 4.838 Å, 3.429 Å, α, β, γ : 90.0°. Color scheme: Au: yellow, O: red. **e** Sm_2BS_4 . Cell parameters: a, b, c : 10.884 Å, α, β, γ : 90.0°. Color scheme: Sm: light green, B: green, S: yellow. **f** KRb_2TiF_6 . Cell parameters: a, b, c : 8.688 Å, α, β, γ : 90.0°. Color scheme: K: white, Rb: purple, Ti: brown, F: green. **g** $\text{LiTa}_2\text{NiSe}_5$ (a : 3.517 Å, b : 13.362 Å, c : 15.156 Å, $Z=4$), which resembles the recently reported structure in [46]. **h** Ta_2NiSe_5 , seen in training. **i** $\text{NaSn}_2\text{CuSe}_5$, seen in training.

2.4.4 Pyrochlores

The general formula for the pyrochlores is $A_2B_2O_7$, where A, a trivalent cation, and B, a tetravalent cation, are either rare-earths or transition metals (other oxidation states, e.g. combining monovalent and pentavalent cations, are also possible, but we focus here on the trivalent/tetravalent pyrochlores). Pyrochlores crystallize in the $Fd\bar{3}m$ space group ($Z=8$). There are many combinations of A and B that are possible for this structure, by using lanthanide ions, actinide ions, and Y(III) for the A species, and various transition metal ions, as well as Ti(IV), Zr(IV), and Hf(IV) for the B species. We investigated whether CrystaLLM could generate valid pyrochlore structures for any unseen combinations, and whether it could estimate reasonable cell parameters in line with the trends observed for the pyrochlore series, as the cell parameters are expected to be correlated with the ionic radii of the A and B cations.

We created a space of pyrochlores consisting of 144 compounds by producing different combinations of A and B species. Of these, 54 were seen in training. We selected 10 compounds from among the 90 not seen in training, and attempted 3 generations with the model, for each. The cell composition and space group were included in the prompt. All generations resulted in valid pyrochlore structures (Table 5).

Table 5: Values of mean generated cell length for the selected pyrochlores not seen in training, over 3 generation attempts.

Formula	Cell Length (Å)
$Ce_2Hf_2O_7$	10.75 ± 0.07
$Ce_2Mn_2O_7$	10.50 ± 0.22
$Ce_2V_2O_7$	10.53 ± 0.09
$La_2Mn_2O_7$	10.21 ± 0.07
$La_2V_2O_7$	10.48 ± 0.06
$Lu_2Hf_2O_7$	10.30 ± 0.08
$Lu_2Zr_2O_7$	10.45 ± 0.12
$Pr_2Mn_2O_7$	10.40 ± 0.08
$Pr_2V_2O_7$	10.51 ± 0.06
$Pr_2Hf_2O_7$	10.80 ± 0.06

We subsequently performed DFT relaxation calculations on the first generated structure for each of the 10 compounds. One case, $Ce_2V_2O_7$, posed challenges in calculation under the generalized gradient approximation and was thus excluded from further analysis. The DFT-derived value of the cell parameter for each of the remaining compounds is plotted against the mean value generated by CrystaLLM in Figure 4. A good agreement exists between the DFT-derived and generated cell lengths, with an R^2 of 0.62 and MAE of 0.08 Å being exhibited. This example illustrates CrystaLLM’s capability to accurately estimate cell parameters of compounds not seen in training with any structure.

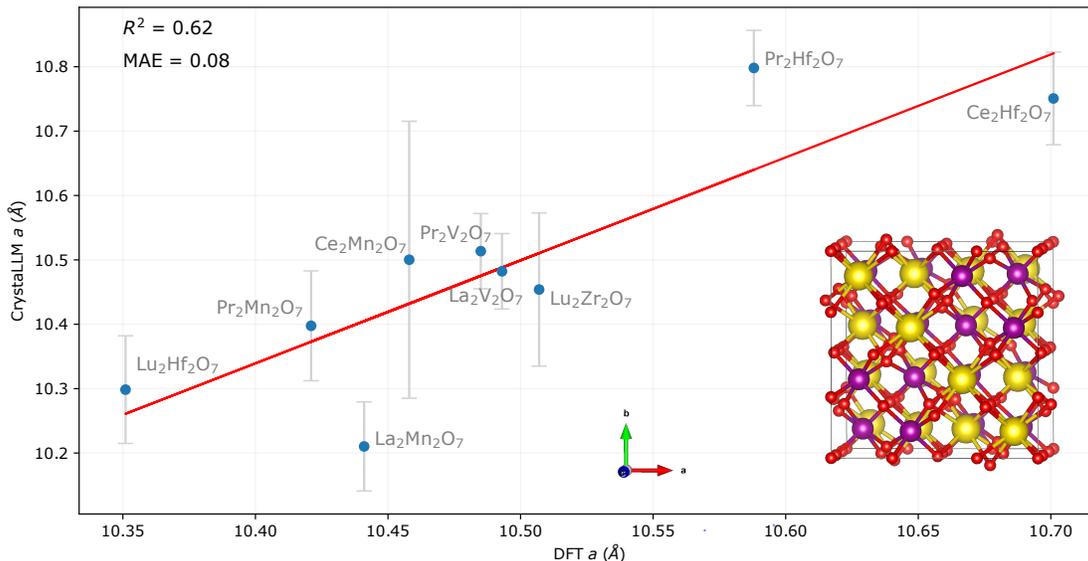


Figure 4: The generated vs. DFT-derived value of the cell parameter a for selected pyrochlores not in the training dataset. The error bars represent the \pm standard deviation of the value of the a cell parameter for the three generation attempts (all of which resulted in the pyrochlore structure), while the y -coordinate of the points represents the mean value of the cell parameter across the three attempts. The inset represents the structure of the generated pyrochlore $\text{Pr}_2\text{Mn}_2\text{O}_7$, with cell parameters a, b, c : 10.34 \AA , α, β, γ : 90.0°. Color scheme: Pr = yellow, Mn = purple, O = red.

2.4.5 Problematic Cases

While the model seems capable of generating structures for many different classes of inorganic crystals, it does nonetheless have difficulty in certain cases. All of the cases appear to involve systems that are rare, and under-represented in the training dataset, or missing from the training set altogether. More precisely, we define a *template* as a unique combination of the reduced composition ratio, the space group, and Z . For example, the combination of the reduced composition ratio 1:1:3:4, space group $Cmcm$, and $Z = 4$, represents a unique template. There are 25,921 unique templates in the dataset.

The problematic cases in the challenge set are largely represented by unseen templates, and templates for which there are few examples. For example, validation rates were low for $\text{Mg}_7\text{Pt}_4\text{Ge}_4$, the structure of which was reported recently to exist in the $P6_3mc$ space group ($Z=2$). [47] In this case, there were only 38 examples of 7:4:4 systems in the training dataset, none contained Mg or Pt, and none were in the $P6_3mc$ space group.

The small version of the model also seems to struggle with generating phosphates, sulfates, carbonates, and organic-inorganic hybrid structures. Examples include carbonate hydroxide minerals, such as $\text{Co}_2\text{CO}_3(\text{OH})_2$ [48] and $\text{Cu}_2\text{CO}_3(\text{OH})_2$ (malachite). While present in the dataset, they belong to a group of analogous structures for which there are

only a handful of examples. While both the small and large versions of the model can generate $\text{Mn}_4(\text{PO}_4)_3$, they generally fail to generate a valid structure for $\text{Ca}_5(\text{PO}_4)_3(\text{OH})$ (hydroxyapatite). A common theme is the appearance of multiple oxyanions, which can give rise to more complex arrangements of atoms, for which the model may not have seen enough examples. In contrast, the model can generate compounds of the perovskite class reliably. However, over 5,000 examples of the ABX_3 ($\text{X}=\text{O},\text{F}$) system in the $Pm\bar{3}m$ space group were seen in training. Finally, structures represented by CIF files with a relatively large number of tokens also pose challenges for the models.

Future versions of the model will consider strategies for addressing these occurrences of class imbalance.

2.5 Heuristic Search for Low-Energy Structures

The examples generated in the previous section were produced through top- k random sampling of the model. Essentially, as the CIF file is generated, each subsequent token is sampled randomly from amongst the top k tokens, according to their probabilities. (See Supplementary Note 2.3 for a detailed description of top- k sampling.) However, random sampling may not necessarily result in the most desirable sequence, and consequently, there are more strategic approaches for constructing sequences that incorporate the probability distributions produced by the model, along with additional heuristics. An example of a heuristic search is Beam Search [49], which is commonly used in natural language contexts to improve the quality of generated sequences. Another popular heuristic search algorithm is MCTS, which has traditionally been used in the context of planning and games, but has recently also been used to increase the quality of generated natural language, through incorporation with LLMs. [50]

Here, we employ the MCTS algorithm, informed by CrystaLLM, to generate a collection of sequences, which is expected to progressively yield sequences of increasingly higher quality as the search advances. In this implementation, each node in the tree represents a cumulative context of tokens. The algorithm operates through a series of steps, including selection, expansion, rollout, evaluation, and backpropagation. The search tree is constructed iteratively, as the search proceeds (Figure 5). In the selection phase, nodes are chosen using the PUCT algorithm (Predictor-Upper Confidence bound applied to Trees) [51, 52], which is a principled means of obtaining a balance between exploring untried nodes, and exploiting promising nodes. The expansion involves adding child nodes based on predicted probabilities. During the rollout step, the CrystaLLM model is prompted with token sequences until a terminating condition is met, leading to the evaluation of the completed sequence. Evaluation is conducted using the ALIGNN (Atomistic Line Graph Neural Network) model of formation energy per atom [53], while the backpropagation step accumulates outcomes in the tree nodes, scoring each based on the quality of the generated structure. (See Supplementary Note 4 for a more detailed description of the algorithm.) The objective is to produce structures with lower formation energy per atom, E_f , and the incorporation of the ALIGNN model allows for a fast and sufficiently accurate estimate of the target property.

When compared to random sampling, MCTS improves the overall validity rate for a compound, and also generally produces lower energy structures. To evaluate the MCTS decoding procedure, we took the 20 most problematic cases of the challenge set where the

validity rate was greater than 0, and performed 1,000 generation attempts using random top- k sampling, and 1,000 iterations of MCTS. The results are presented in table 6.

Table 6: Results of MCTS decoding for the 20 most problematic cases of the challenge set. The percentages represent the fraction of cases with the corresponding improvement after using MCTS decoding, when compared to random sampling. The first row represents the percentage of cases where the validity rate improved. The second row represents the percentage of cases where the minimum E_f obtained was improved. The third row represents the percentage of cases where the mean E_f was improved.

	No Space Group	With Space Group
Validity Rate Improvement	95.0%	60.0%
Minimum E_f Improvement	85.0%	65.0%
Mean E_f Improvement	70.0%	65.0%

When no space group is provided in the prompt, the validity rate improves in 95% of the cases, and the minimum E_f attained improves in 85% of cases. (See Supplementary Tables 6 and 7 for more detailed results.) In some cases, the validity rate increases as the search proceeds when using MCTS (see Supplementary Figure 6).

2.6 Beyond Element Substitution

Although CrystaLLM appears to be very effective at finding appropriate template systems for a given cell composition, and making the necessary adjustments of cell parameters to substitute different atoms, it appears capable of going further, synthesizing information from different template systems. An example is the selenide $\text{LiTa}_2\text{NiSe}_5$, which is obtained by lithium intercalation into Ta_2NiSe_5 [46].

The compound $\text{LiTa}_2\text{NiSe}_5$ was not present in the training set, however, the layered material Ta_2NiSe_5 was (Figure 3g,h). As $\text{LiTa}_2\text{NiSe}_5$ was included in the challenge set, we performed 100 generation attempts with the model. While the model was not able to recover the lowest energy structure reported, it did produce structures with close resemblance to low-energy polymorphs. Upon closer examination of the dataset, we found that $\text{NaSn}_2\text{CuSe}_5$ was present (Figure 3i), which likely provided some precedent for the intercalation of atoms between layered structures. It thus appears that the model is capable of integrating information from different template systems to form new structural predictions.

2.7 The CrystaLLM.com Web Application

To allow for easy and open access to the CrystaLLM model, we make it available through a web application, published at <https://crystallm.com>. The application allows users to enter in a reduced formula, and optionally a value for Z and the desired space group. The option to select the model size is also provided. The request is sent to the model, and the resulting structure (or the CIF contents, if the structure is invalid) is presented to the user. By making the model easily accessible, we hope to contribute a potentially useful tool to the materials structure research community. We also hope to receive feedback from users that may help improve future versions of the model.

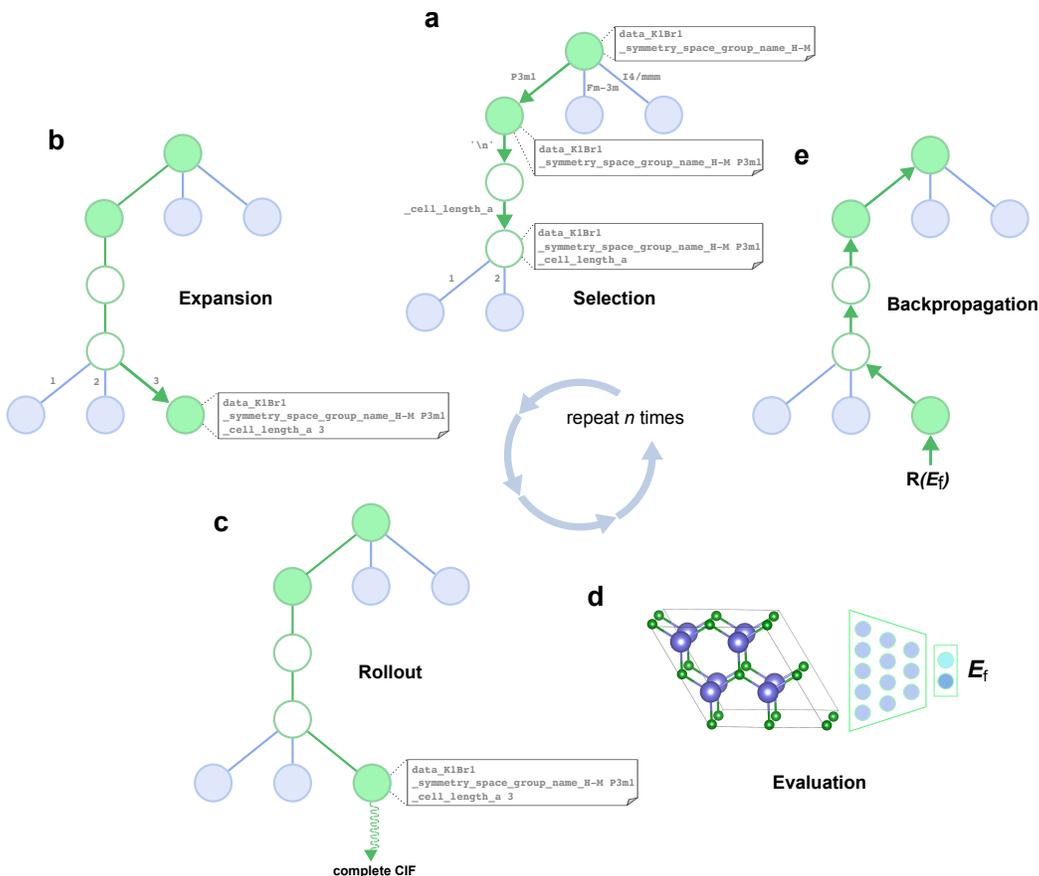


Figure 5: Schematic depiction of the Monte Carlo Tree Search decoding procedure. CIF files are generated as a tree is iteratively constructed, with each iteration guiding the generation of subsequent structures towards more desirable parameters (e.g. lower formation energy per atom). The nodes in the tree represent the cumulative contents of a CIF file at various points. **a** The Selection step involves descending the tree by choosing the most promising node at each level, using a variant of the PUCT algorithm. **b** During Expansion, an unexplored child node is randomly selected and added to the tree. If a node has only one highly probable child (represented as empty nodes), the child node bypasses the Rollout step. **c** The Rollout step involves prompting the model with the contents of the selected node, and sampling from the model until a terminal condition is met, so as to obtain a complete CIF file and an estimate of the value of a node. **d** The generated structure is validated and scored, incorporating the prediction of the structure’s formation energy per atom, as given by a pre-trained neural network. **e** Finally, the score is backpropagated through the selected nodes, which store the accumulated results of each iteration. The resulting generated CIF file, if valid, is returned.

3 Discussion

Here, we have shown that LLMs of the CIF format are able to generate inorganic crystal structures for a variety of known classes. Indeed, the model is able to produce valid and sensible arrangements of atoms in 3-dimensional space by generating *xyz* coordinates digit-by-digit. The model also seems to have captured the relationship between space group symbols and the symmetries inherent in the structures it generates.

We chose to build a language model of the CIF format (instead of a simplified format, for example, which might include a minimal vocabulary) for several reasons. First, the CIF format is not particularly verbose. The model learns the grammatical structure of the format fairly quickly. We can thus avoid having to devise an intermediate format that requires inter-conversion between more common formats, which could also be error prone. Second, we believe that having the model learn to generate the more redundant parts of the CIF format, such as the cell volume, and *Z*, which are inferable from prior inputs, helps the model to perform better overall.

A number of approaches for crystal structure generation have been reported. [54–57] These approaches generally require the existence of pre-defined structural templates, and are followed by the procedural or machine learning-assisted substitution of atoms and adjustment of cell parameters, under the constraint of a specified space group. These types of approaches can also be enhanced to increase the structural diversity of generated materials, by allowing partial substitutions and adjusting substitution probabilities [58]. Conversely, CrystaLLM automatically selects the templates which can be applied to a given composition, utilizing the implicit templates it has absorbed through autoregressive training. Moreover, the model can automatically adjust cell parameters to accommodate the atoms in the unit cell. It can also produce structures based on templates it has not explicitly encountered in training, borrowing from its internalized concepts of chemical structure. In comparison with recently reported diffusion-based ML methods for crystal generation (CDVAE [13] and DiffCSP [36]), not only does CrystaLLM outperform them on established benchmarks in several aspects, but it also offers additional advantages in terms of flexibility (e.g. in using symmetry as input) and the potential for fine-tuning.

While the CrystaLLM model can generate sensible structures, this does not by itself make it suitable, as is, for CSP. Just as natural language LLMs, such as GPT-3 and -4, are not suitable chatbots without further fine-tuning and alignment, the CrystaLLM model will also need to be fine-tuned for more advanced tasks. Fine-tuning involves an additional and separate training step, where the model’s parameters are adjusted in the context of a different task. This may also involve altering the model’s output layer, such as to make it suitable for a regression task. Models can be fine-tuned using a variety of techniques, but supervised learning and reinforcement learning [59] are most common. One might use reinforcement learning, for example, when a task is not clearly defined as a supervised learning problem. When fine-tuning natural language LLMs for chatbot applications, it is common to use Reinforcement Learning from Human Feedback (RLHF) [60, 61]. With RLHF, the idea is to gather data from human annotators to be used to train a reward model, which scores generated text according to its desirability. The reward model is then used as part of a reinforcement learning-based tuning of the LLM. In CSP, one would like to produce ground-state structures (for some given physical conditions). One could thus imagine an analogous procedure where CrystaLLM is fine-

tuned for the goal of generating low-energy structures, via feedback from an external evaluator of the generated structure’s energy, resulting in *Reinforcement Learning from Thermodynamic Feedback*. This procedure would also require a reward model, and such a model should ideally provide a timely estimate of a structure’s energy. This excludes time-consuming approaches such as DFT. A viable approach could make use of a separate machine learning-based model of formation energy, such as one based on ALIGNN. Indeed, neural network potentials have been used to accelerate the prediction of crystal structures, and the identification of potentially stable materials. [62, 63]

There are several limitations with the current approach. First, none of the structures of the dataset have site-occupancy disorder (fractional site occupancies). Therefore, CrystaLLM cannot generate disordered structures, and may not successfully generate structures for combinations of cell composition and space group that imply a disordered structure. An example is $\text{K}_2\text{NaTiOF}_5$, which is reported to be an elpasolite, in the $Fm\bar{3}m$ space group ($Z=4$), with F and O species sharing the same crystal site [64]. Another limitation is that the CIF files of the dataset were not all created using the same level of theory. The training set is derived from a combination of DFT sources using different settings, functionals, etc., which may make it difficult for the model, in some instances, to learn a consistent relationship between cell composition and detailed structure [65].

Nevertheless, we believe that CrystaLLM will be a useful tool for crystal structure generation, which is quickly becoming a critical step in large scale materials discovery [58, 66], and materials informatics. We plan to explore fine-tuning the model for physical property prediction tasks, such as the prediction of lattice thermal conductivity, where experimental data is relatively scarce. [67] The architecture of the model allows it to be fine-tuned for either composition-based or structure-based prediction tasks. This implies that CrystaLLM may be the basis for a general-purpose materials informatics model, which can be used for generative tasks, and fine-tuned for property prediction tasks that require either composition or structure. If the model is able to transfer what it has learned about the world of atoms to these various predictive problems, it may prove to be a quite flexible tool relevant to many aspects of materials chemistry.

4 Methods

4.1 Dataset Curation

The dataset was assembled by obtaining structures from the Materials Project [40], the OQMD [68], and NOMAD [69], which were originally optimized using density functional theory (DFT) simulations. Specifically, the structures from the Materials Project were downloaded in April 2022, and from NOMAD in April 2023. We use version 1.5 of the OQMD, which was released in October 2021. In total, approximately 3.6 million structures were obtained. This dataset consists of compounds containing anywhere from 1 to 10 elements, with most consisting of 3 or 4 elements. The elements up to and including atomic number 94 are present, with the exception of polonium, astatine, radon, francium, and radium. The dataset contains roughly 800,000 unique formulas, and 1.2 million unique cell compositions. When paired with space groups, there are 2.3 million unique cell composition-space group pairs. (See Supplementary Figure 1.) To choose between duplicate structures containing the same cell composition and space group, the

structure with the lowest volume per formula unit was selected. The 2.3 million structures in this dataset were converted to CIF files using the pymatgen library [70], and were used for training. The CIF files were created with the pymatgen option for symmetry finding tolerance set to 0.1 Å. All floating point numbers in the files were rounded to 4 decimal places. The dataset was split randomly into train, validation, and test sets, such that the training set consisted of 2,047,889 CIF files, the validation set 227,544 CIF files, and the test set 10,286 CIF files.

4.2 CIF Syntax Standardization and Tokenization

The dataset of CIF files was standardized and tokenized prior to training. The vocabulary consisted of CIF tags, space group symbols, element symbols, numeric digits, and various punctuation symbols, for a total of 371 symbols. After tokenization, the training set consisted of 768 million tokens. See Supplementary Note 1 for further details.

4.3 Generative Pre-training

The generative pre-training step requires a vocabulary, \mathcal{V} , and an ordered list of tokens $\mathcal{U} = (u_1, \dots, u_n)$, with $u_i \in \mathcal{V}$. We want to maximize the following likelihood:

$$\mathcal{L}(\theta; \mathcal{U}) = \sum_i \log P(u_i | u_{i-c}, \dots, u_{i-1}; \theta) \quad (1)$$

where c is the size of a context window, P is the conditional probability distribution to be modelled, and θ the parameters of a neural network. We therefore minimize $\mathcal{J}(\theta; \mathcal{U}) = -\mathcal{L}$, using stochastic gradient descent to adjust the parameters. We use a multi-layer Transformer decoder [71] for the neural network, as described in [18]. Our model consists of 25 million parameters, with 8 layers, 8 attention heads, and an embedding size of 512. We decay the learning rate from 10^{-3} to 10^{-4} over the course of training, and use a batch size of 32. For further details, see Supplementary Note 2.

4.4 Evaluation of Generated Structures

A CIF file is said to be *valid* if: 1) the declared space group is consistent with the generated structure, 2) the generated bond lengths are reasonable, and 3) the declared atom site multiplicity is consistent with the cell composition. To check if the generated structure is consistent with the printed space group, we use the `SpacegroupAnalyzer` class of the pymatgen library, which uses the `spglib` library [72]. To check if bond lengths are reasonable, we first use a Voronoi-based nearest-neighbour algorithm in pymatgen to identify bonded atoms; then, we establish expected bond lengths based on the electronegativity difference between the bonded atoms, and their ionic or covalent radii. We classify a structure as having reasonable bond lengths if all the detected bond lengths are within 30% of the corresponding expected bond lengths. See Supplementary Note 3 for more details on how the validity of a generated CIF file is established.

In some scenarios, we wish to determine whether a generated structure matches a target structure, which typically represents a ground-truth structure. To determine whether two structures are a match, we use the pymatgen `StructureMatcher` class, which

performs a structural similarity assessment of two crystals. We use a fractional length tolerance of 0.2, a site tolerance of 0.3 Å, and an angle tolerance of 5 degrees, which are the default values in pymatgen. Both structures are reduced to primitive cells before matching, and are scaled to equivalent volume.

4.5 Benchmark Evaluations

To evaluate CrystaLLM on the Perov-5, Carbon-24, MP-20 and MPTS-52 benchmarks, we consider two different scenarios: 1) the model is trained only on the benchmark training sets, and 2) the model is trained on the full 2.3 million-structure dataset minus the validation and test set structures of the MPTS-52 dataset. For the first scenario, both the small and large model architectures are used. We use the same 60-20-20 train/validation/test splits used in the CDVAE study [13] for the Perov-5, Carbon-24, and MP-20 datasets, and we use the same 27,380/5,000/8,096 train/validation/test split used in the DiffCSP study for the MPTS-52 dataset. These models are trained for a fixed number of iterations: the Perov-5 model is trained for 1,750 iterations, the Carbon-24 model is trained for 8,000 iterations, the MP-20 model is trained for 5,000 iterations, and the MPTS-52 model is trained for 3,500 iterations. For the second scenario, we train a model with the small model architecture on the full 2.3 million-structure dataset minus the structures of the MPTS-52 validation and test sets. The model is trained for 100,000 iterations. We decay the learning rate from 10^{-3} to 10^{-4} over the course of training, and use a batch size of 32, for all models. For both scenarios, we take the structures of the test set(s), and prompt the models with only the cell compositions of these structures. Models are given 20 attempts to generate a structure. We use top- k sampling with $k = 10$ and a temperature of 1.0 for all models and in both scenarios.

To establish the match rate and RMSE, we use the same procedure defined in the DiffCSP study. Specifically, we use the pymatgen `StructureMatcher` class, with a fractional length tolerance of 0.3, a site tolerance of 0.5 Å, and an angle tolerance of 10 degrees, to determine if a generation attempt matches the ground truth structure. The RMSE, normalized by $\sqrt[3]{V/N}$ (where V is the volume of the lattice and N is the number of sites), is computed between the corresponding ground truth structure and each matching generated structure. The test set’s average RMSE is computed by taking the lowest RMSE for each entry’s matching generated structure.

4.6 Monte Carlo Tree Search Decoding

The MCTS search tree is constructed iteratively, as the search proceeds. We maintain a tree width of 5, and maximum tree depth of 1,000. The PUCT constant c_{puct} is set at 1.0. The expansion involves adding child nodes based on predicted probabilities. When a node has a probability of 0.99 or greater, it becomes the only child node, and bypasses the rollout step. During the rollout step, the CrystaLLM model is prompted with token sequences until a terminating condition is met, up to a maximum of 1,000 tokens. Evaluation is conducted using the ALIGNN model of formation energy per atom. The ALIGNN model is given the generated CIF file, and the predicted formation energy per atom (in eV) is used to compute the reward. The backpropagation step accumulates outcomes in the tree nodes, scoring each based on the quality of the generated structure,

with a reward constant λ of 2.0. For all compounds, we perform 1,000 search iterations. See Supplementary Note 4 for a more detailed description of the algorithm.

4.7 DFT Calculations

For the pyrochlore case study, a small number of DFT calculations were performed using VASP, following as closely as possible the settings used in the OQMD project (where most of the pyrochlore structures seen in training were taken from). For example, the recommended PAW potential was used for each element: Zr_sv for zirconium, Hf_pv for hafnium, Lu_3 for lutetium, Pr_3 for praseodymium, Ce_3 for cerium (for the remaining elements, the name of the PAW potential simply matched the element’s symbol). The Perdew-Burke- Ernzerhof (PBE) exchange-correlation functional [73], in the generalized-gradient approximation, was used in all calculations. Hubbard (PBE+U) corrections were applied for transition metal elements with unfilled d levels ($U_{\text{eff}}=3.8$ eV for Mn and 3.1 eV for V). Although the cell parameters reported here correspond to the conventional cubic cell with 8 formula units, the DFT calculations were performed using the primitive cell with two formula units, and sampling of the reciprocal space corresponding to that primitive cell was performed using a 7x7x7 grid, as done for all pyrochlore calculations in the OQMD project.

4.8 Web Application

The web application is made available at <https://crystallm.com>. The user of the application is presented with a text field requiring a formula to be entered. Optionally, they may provide the number of formula units (Z), the desired space group, and the size of the model. Once they press the `Generate` button, a request is sent to a GPU server which has the model in memory. The request is converted into a prompt, and the generated contents are returned to the user. If no Z is provided, we scan through Z values of 1, 2, 3, 4, 6, and 8, and return the first valid structure generated by the model. We validate the generated structure using the same procedure described previously, checking that the generated structure is consistent in terms of the printed space group, and other elements of the CIF file. If no valid structure can be found, the user is presented with an informative error message, including the option to view the generated content. Requests typically take several seconds to process, but can take longer if no Z is provided and the model has trouble generating a valid structure for the attempted Z values. Generated structures are displayed in a web browser-based 3D structure viewer provided by the Crystal Toolkit framework, upon which the front-end of the web application is built [74].

5 Code Availability

The code for training and using the CrystaLLM model is open source, released under the MIT License. The code repository is accessible online, at: <https://github.com/lantunes/CrystaLLM>.

6 Data Availability

The structures used in the experiments described in this work were obtained from the Materials Project (<https://materialsproject.org/>), the OQMD (<https://oqmd.org/>), and NOMAD (<https://nomad-lab.eu/>). All structures were made available by those sources under the Creative Commons Attribution 4.0 License [75].

All trained models, training sets, and artifacts generated by the models have been deposited to Zenodo. The files are publicly accessible at: <https://zenodo.org/records/10642388>. All files are released under the CC-BY 4.0 license.

7 Acknowledgements

This work was partially supported by computational resource donations from Amazon Web Services through the AWS Activate program, obtained with assistance from the Communtech Hub. For the DFT calculations, we used the Young supercomputer facility via the UK Materials and Molecular Modelling Hub, which is partially funded by EPSRC (EP/T022213/1, EP/W032260/1).

8 Author Contributions

L.M.A. conceived the project, performed the experiments, and drafted the manuscript. L.M.A. and R.G.-C. designed the experiments. R.G.-C. carried out the DFT calculations. R.G.-C. and K.T.B. supervised and guided the project. All authors reviewed, edited and approved the manuscript.

References

- [1] Cerqueira, T. F. *et al.* Identification of Novel Cu, Ag, and Au Ternary Oxides from Global Structural Prediction. *Chemistry of Materials* **27**, 4562–4573 (2015).
- [2] Zhu, B. & Scanlon, D. O. Predicting Lithium Iron Oxysulfides for Battery Cathodes. *ACS Applied Energy Materials* **5**, 575–584 (2022).
- [3] Harper, A. F., Evans, M. L. & Morris, A. J. Computational Investigation of Copper Phosphides as Conversion Anodes for Lithium-Ion Batteries. *Chemistry of Materials* **32**, 6629–6639 (2020).
- [4] Oganov, A. R., Pickard, C. J., Zhu, Q. & Needs, R. J. Structure prediction drives materials discovery. *Nature Reviews Materials* **4**, 331–348 (2019).
- [5] Oganov, A. R. *Modern Methods of Crystal Structure Prediction* (John Wiley & Sons, 2011).
- [6] Pickard, C. J. & Needs, R. High-Pressure Phases of Silane. *Physical Review Letters* **97**, 045504 (2006).

- [7] Pickard, C. J. & Needs, R. Ab initio random structure searching. *Journal of Physics: Condensed Matter* **23**, 053201 (2011).
- [8] Oganov, A. R. & Glass, C. W. Crystal structure prediction using *ab initio* evolutionary techniques: Principles and applications. *The Journal of Chemical Physics* **124**, 244704 (2006).
- [9] Butler, K. T., Davies, D. W., Cartwright, H., Isayev, O. & Walsh, A. Machine learning for molecular and materials science. *Nature* **559**, 547–555 (2018).
- [10] Podryabinkin, E. V., Tikhonov, E. V., Shapeev, A. V. & Oganov, A. R. Accelerating crystal structure prediction by machine-learning interatomic potentials with active learning. *Physical Review B* **99**, 064114 (2019).
- [11] Goodfellow, I. *et al.* Generative Adversarial Nets. *Advances in Neural Information Processing Systems* **27** (2014).
- [12] Court, C. J., Yildirim, B., Jain, A. & Cole, J. M. 3-D Inorganic Crystal Structure Generation and Property Prediction via Representation Learning. *Journal of Chemical Information and Modeling* **60**, 4518–4535 (2020).
- [13] Xie, T., Fu, X., Ganea, O.-E., Barzilay, R. & Jaakkola, T. Crystal Diffusion Variational Autoencoder for Periodic Material Generation. *arXiv preprint arXiv:2110.06197* (2021).
- [14] Yan, D., Smith, A. D. & Chen, C.-C. Structure prediction and materials design with generative neural networks. *Nature Computational Science* 1–3 (2023).
- [15] Alverson, M. *et al.* Generative adversarial networks and diffusion models in material discovery. *Digital Discovery* (2024).
- [16] Cao, Y. *et al.* A Comprehensive Survey of AI-Generated Content (AIGC): A History of Generative AI from GAN to ChatGPT. *arXiv preprint arXiv:2303.04226* (2023).
- [17] Vaswani, A. *et al.* Attention Is All You Need. *Advances in Neural Information Processing Systems* **30** (2017).
- [18] Radford, A., Narasimhan, K., Salimans, T., Sutskever, I. *et al.* Improving Language Understanding by Generative Pre-Training (2018).
- [19] Introducing ChatGPT. <https://openai.com/blog/chatgpt>. Accessed: 2023-06-21.
- [20] Bran, A. M., Cox, S., White, A. D. & Schwaller, P. ChemCrow: Augmenting large-language models with chemistry tools. *arXiv preprint arXiv:2304.05376* (2023).
- [21] Jablonka, K. M., Schwaller, P., Ortega-Guerrero, A. & Smit, B. Is GPT-3 all you need for low-data discovery in chemistry? *ChemRxiv* 1–32 (2023).
- [22] Xie, T. *et al.* Large Language Models as Master Key: Unlocking the Secrets of Materials Science with GPT. *arXiv preprint arXiv:2304.02213* (2023).

- [23] Fu, N. *et al.* Material transformers: deep learning language models for generative materials design. *Machine Learning: Science and Technology* **4**, 015001 (2023).
- [24] Jablonka, K. M. *et al.* 14 examples of how LLMs can transform materials science and chemistry: a reflection on a large language model hackathon. *Digital Discovery* **2**, 1233–1250 (2023).
- [25] Boiko, D. A., MacKnight, R., Kline, B. & Gomes, G. Autonomous chemical research with large language models. *Nature* **624**, 570 (2023).
- [26] Flam-Shepherd, D. & Aspuru-Guzik, A. Language models can generate molecules, materials, and protein binding sites directly in three dimensions as XYZ, CIF, and PDB files. *arXiv preprint arXiv:2305.05708* (2023).
- [27] Hall, S. R., Allen, F. H. & Brown, I. D. The crystallographic information file (CIF): a new standard archive file for crystallography. *Acta Crystallographica Section A: Foundations of Crystallography* **47**, 655–685 (1991).
- [28] Chen, M. *et al.* Generative Pretraining from Pixels. In *International Conference on Machine Learning*, 1691–1703 (PMLR, 2020).
- [29] Chen, C., Ye, W., Zuo, Y., Zheng, C. & Ong, S. P. Graph Networks as a Universal Machine Learning Framework for Molecules and Crystals. *Chemistry of Materials* **31**, 3564–3572 (2019).
- [30] Toshniwal, S., Wiseman, S., Livescu, K. & Gimpel, K. Chess as a Testbed for Language Model State Tracking. In *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 36, 11385–11393 (2022).
- [31] Li, K. *et al.* Emergent World Representations: Exploring a Sequence Model Trained on a Synthetic Task. In *The Eleventh International Conference on Learning Representations* (2023). URL https://openreview.net/forum?id=DeG07_TcZvT.
- [32] Coulom, R. Efficient Selectivity and Backup Operators in Monte-Carlo Tree Search. In *International Conference on Computers and Games*, 72–83 (Springer, 2006).
- [33] Browne, C. B. *et al.* A Survey of Monte Carlo Tree Search Methods. *IEEE Transactions on Computational Intelligence and AI in games* **4**, 1–43 (2012).
- [34] Brown, T. *et al.* Language Models are Few-Shot Learners. *Advances in Neural Information Processing Systems* **33**, 1877–1901 (2020).
- [35] Antunes, L. M., Grau-Crespo, R. & Butler, K. T. Distributed representations of atoms and materials for machine learning. *npj Computational Materials* **8**, 44 (2022).
- [36] Jiao, R. *et al.* Crystal Structure Prediction by Joint Equivariant Diffusion. *arXiv preprint arXiv:2309.04475* (2023).
- [37] Castelli, I. E. *et al.* New cubic perovskites for one- and two-photonwater splitting using the computational materials repository. *Energy & Environmental Science* **5**, 9034–9043 (2012).

- [38] Castelli, I. E. *et al.* Computational screening of perovskite metal oxides for optimal solar light capture. *Energy & Environmental Science* **5**, 5814–5819 (2012).
- [39] Pickard, C. J. AIRSS Data for Carbon at 10GPa and the C+N+H+O System at 1GPa. <https://archive.materialscloud.org/record/2020.0026/v1> (2020).
- [40] Jain, A. *et al.* Commentary: The Materials Project: A materials genome approach to accelerating materials innovation. *APL Materials* **1**, 011002 (2013).
- [41] Baird, S. mp-time-split. <https://github.com/sparks-baird/mp-time-split> (Accessed in 2024).
- [42] Mazet, T., Welter, R. & Malaman, B. A study of the new ferromagnetic YbMn_6Sn_6 compound by magnetization and neutron diffraction measurements. *Journal of Magnetism and Magnetic Materials* **204**, 11–19 (1999).
- [43] Pamplin, B. A systematic method of deriving new semiconducting compounds by structural analogy. *Journal of Physics and Chemistry of Solids* **25**, 675–684 (1964).
- [44] Davies, D. W. *et al.* Computational Screening of All Stoichiometric Inorganic Materials. *Chem* **1**, 617–627 (2016).
- [45] Zagorac, D., Müller, H., Ruehl, S., Zagorac, J. & Rehme, S. Recent developments in the Inorganic Crystal Structure Database: theoretical crystal structure data and related features. *Journal of Applied Crystallography* **52**, 918–925 (2019).
- [46] Hyde, P. *et al.* Lithium Intercalation into the Excitonic Insulator Candidate Ta_2NiSe_5 . *Inorganic Chemistry* **62**, 12027–12037 (2023).
- [47] Ponou, S., Lidin, S. & Mudring, A.-V. Optimization of Chemical Bonding through Defect Formation and Ordering—The Case of $\text{Mg}_7\text{Pt}_4\text{Ge}_4$. *Inorganic Chemistry* (2023).
- [48] González-López, J., Cockcroft, J. K., Fernández-González, A., Jimenez, A. & Grau-Crespo, R. Crystal structure of cobalt hydroxide carbonate $\text{Co}_2\text{CO}_3(\text{OH})_2$: density functional theory and X-ray diffraction investigation. *Acta Crystallographica Section B: Structural Science, Crystal Engineering and Materials* **73**, 868–873 (2017).
- [49] Science, C.-M. U. P. P. D. O. C. *Speech Understanding Systems. Summary of Results of the Five-Year Research Effort at Carnegie-Mellon University* (1977).
- [50] Chaffin, A., Claveau, V. & Kijak, E. PPL-MCTS: Constrained Textual Generation Through Discriminator-Guided MCTS Decoding. In Carpuat, M., de Marneffe, M. & Ruíz, I. V. M. (eds.) *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL 2022, Seattle, WA, United States, July 10-15, 2022*, 2953–2967 (Association for Computational Linguistics, 2022).
- [51] Rosin, C. D. Multi-armed Bandits with Episode Context. *Annals of Mathematics and Artificial Intelligence* **61**, 203–230 (2011).

- [52] Silver, D. *et al.* Mastering the game of Go with deep neural networks and tree search. *Nature* **529**, 484–489 (2016).
- [53] Choudhary, K. & DeCost, B. Atomistic Line Graph Neural Network for improved materials property predictions. *npj Computational Materials* **7**, 185 (2021).
- [54] Kusaba, M., Liu, C. & Yoshida, R. Crystal structure prediction with machine learning-based element substitution. *Computational Materials Science* **211**, 111496 (2022).
- [55] Wei, L. *et al.* TCSP: a Template-Based Crystal Structure Prediction Algorithm for Materials Discovery. *Inorganic Chemistry* **61**, 8431–8439 (2022).
- [56] Fredericks, S., Parrish, K., Sayre, D. & Zhu, Q. PyXtal: A Python library for crystal structure generation and symmetry analysis. *Computer Physics Communications* **261**, 107810 (2021).
- [57] Avery, P. & Zurek, E. RandSpg: An open-source program for generating atomistic crystal structures with specific spacegroups. *Computer Physics Communications* **213**, 208–216 (2017).
- [58] Merchant, A. *et al.* Scaling deep learning for materials discovery. *Nature* 1–6 (2023).
- [59] Sutton, R. S. & Barto, A. G. *Reinforcement Learning: An Introduction* (MIT press, 2018).
- [60] Ziegler, D. M. *et al.* Fine-Tuning Language Models from Human Preferences. *arXiv preprint arXiv:1909.08593* (2019).
- [61] Illustrating Reinforcement Learning from Human Feedback (RLHF). <https://huggingface.co/blog/rlhf>. Accessed: 2023-07-05.
- [62] Kang, S. *et al.* Accelerated identification of equilibrium structures of multicomponent inorganic crystals using machine learning potentials. *npj Computational Materials* **8**, 108 (2022).
- [63] Chen, C. & Ong, S. P. A universal graph deep learning interatomic potential for the periodic table. *Nature Computational Science* **2**, 718–728 (2022).
- [64] Pausewang, G. & Rüdorff, W. Über Alkali-oxofluorometallate der Übergangsmetalle. $A'_3\text{MeO}_x\text{F}_{6-x}$ -Verbindungen mit $x = 1, 2, 3$. *Zeitschrift für anorganische und allgemeine Chemie* **364**, 69–87 (1969).
- [65] Hegde, V. I. *et al.* Quantifying uncertainty in high-throughput density functional theory: A comparison of AFLOW, Materials Project, and OQMD. *Physical Review Materials* **7**, 053805 (2023).
- [66] Ye, W., Lei, X., Aykol, M. & Montoya, J. H. Novel inorganic crystal structures predicted using autonomous simulation agents. *Scientific Data* **9**, 302 (2022).

- [67] Antunes, L. M. *et al.* Machine Learning Approaches for Accelerating the Discovery of Thermoelectric Materials. In *Machine Learning in Materials Informatics: Methods and Applications*, 1–32 (ACS Publications, 2022).
- [68] Saal, J. E., Kirklin, S., Aykol, M., Meredig, B. & Wolverton, C. Materials Design and Discovery with High-Throughput Density Functional Theory: The Open Quantum Materials Database (OQMD). *JOM* **65**, 1501–1509 (2013).
- [69] Draxl, C. & Scheffler, M. The NOMAD laboratory: from data sharing to artificial intelligence. *Journal of Physics: Materials* **2**, 036001 (2019).
- [70] Ong, S. P. *et al.* Python Materials Genomics (pymatgen): A robust, open-source python library for materials analysis. *Computational Materials Science* **68**, 314–319 (2013).
- [71] Liu, P. J. *et al.* Generating Wikipedia by Summarizing Long Sequences. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings* (2018).
- [72] Togo, A. & Tanaka, I. Spglib: a software library for crystal symmetry search. *arXiv preprint arXiv:1808.01590* (2018).
- [73] Perdew, J. P., Burke, K. & Ernzerhof, M. Generalized gradient approximation made simple. *Physical review letters* **77**, 3865 (1996).
- [74] Horton, M. *et al.* Crystal Toolkit: A Web App Framework to Improve Usability and Accessibility of Materials Science Research Algorithms. *arXiv preprint arXiv:2302.06147* (2023).
- [75] Creative Commons Attribution 4.0 License. <https://creativecommons.org/licenses/by/4.0/>. Accessed: 2023-06-26.

Supplementary Information for the article “Crystal Structure Generation with Autoregressive Large Language Modeling”

Luis M. Antunes¹, Keith T. Butler², and Ricardo Grau-Crespo¹

¹*Department of Chemistry, University of Reading, Whiteknights, Reading RG6 6DX, United Kingdom.
l.m.antunes@pgr.reading.ac.uk*

²*Department of Chemistry, University College London, WC1H0AJ, United Kingdom.*

Supplementary Notes

1. CIF Syntax Standardization and Tokenization

The CIF format is flexible in terms of the sequence of tags in the file. Moreover, not all tags are required to be present in the file. While a large language model could, in principle, learn to process variable arrangements of the tags, we chose to restrict the CIF file syntax, such that every CIF file in the dataset is structured identically. Furthermore, we added several tags that are not part of the CIF specification.

To ensure consistency, and enhance the model’s ability to learn from the data, we standardized the CIF files using a sequence of pre-processing steps. The steps were designed to not only normalize the format of the CIF files, but also to incorporate additional information beneficial for the model’s training. The pre-processing steps are as follows:

1. Each structure in the dataset was first converted into a `pymatgen Structure` object.
2. The `pymatgen CifWriter` class was used to generate CIF files from the `Structure` objects.
3. In the generated CIF files, we replaced the content of the `data_` tag, which contains the reduced formula, with the cell composition of the structure. The atoms of the cell composition appended to `data_` are sorted by electronegativity.
4. We removed the symmetry operators from the CIF files.
5. We introduced a custom block in the CIF files to include specific atomic properties, namely, the electronegativity, the radius, and ionic radius. These properties are not part of the standard CIF specification.
6. All numerical values in the CIF files were rounded to four decimal places.

An example of a pre-processed CIF file from the training dataset is given below:

Pre-processed CIF file for PbTe (Z=2, Pmma)

```

data_Te2Pb2
loop_
_atom_type_symbol
_atom_type_electronegativity
_atom_type_radius
_atom_type_ionic_radius
Te 2.1000 1.4000 1.2933
Pb 2.3300 1.8000 1.1225
_symmetry_space_group_name_H-M Pmma
_cell_length_a 5.6440
_cell_length_b 4.0012
_cell_length_c 5.6807
_cell_angle_alpha 90.0000
_cell_angle_beta 90.0000
_cell_angle_gamma 90.0000
_symmetry_Int_Tables_number 51
_chemical_formula_structural TePb
_chemical_formula_sum 'Te2 Pb2'
_cell_volume 128.2864
_cell_formula_units_Z 2
loop_
_symmetry_equiv_pos_site_id
_symmetry_equiv_pos_as_xyz
1 'x, y, z'
loop_
_atom_site_type_symbol
_atom_site_label
_atom_site_symmetry_multiplicity
_atom_site_fract_x
_atom_site_fract_y
_atom_site_fract_z
_atom_site_occupancy
Te Te0 2 0.2500 0.5000 0.7357 1
Pb Pb1 2 0.2500 0.0000 0.2691 1

```

After the pre-processing step, the CIF files are tokenized; that is, a CIF file is parsed and converted into a sequence of numbers, where each number represents a particular token. Tokenization is necessary for converting the structured, text-based data of CIF files into a format that the model can process. The selection of tokens is guided by a custom vocabulary. The vocabulary defines the distinct, irreducible elements of the CIF file syntax that are relevant to the problem. In constructing this vocabulary, we chose to represent numeric digits, atomic symbols, space group symbols, and CIF tags with distinct tokens. Specifically, the vocabulary consists of digits: 0 1 2 3 4 5 6 7 8 9, as well as various symbols: x y z . () ' , _ (space) \n (newline). A complete enumeration of the supported atom, CIF tag, and space group symbols follows:

Supported atom tokens.

```

Ac Ag Al Ar As Au B Ba Be Bi Br C Ca Cd Ce Cl Co Cr Cs Cu Dy Er Eu F Fe Ga Gd Ge
H He Hf Hg Ho I In Ir K Kr La Li Lu Mg Mn Mo N Na Nb Nd Ne Ni Np O Os P Pa Pb Pd
Pm Pr Pt Pu Rb Re Rh Ru S Sb Sc Se Si Sm Sn Sr Ta Tb Tc Te Th Ti Tl Tm U V W Xe
Y Yb Zn Zr

```

Supported CIF tag tokens.

_cell_length_b	_atom_site_occupancy
_atom_site_attached_hydrogens	_cell_length_a
_cell_angle_beta	_symmetry_equiv_pos_as_xyz
_cell_angle_gamma	_atom_site_fract_x
_symmetry_space_group_name_H-M	_symmetry_Int_Tables_number
_chemical_formula_structural	_chemical_name_systematic
_atom_site_fract_y	_atom_site_symmetry_multiplicity
_chemical_formula_sum	_atom_site_label
_atom_site_type_symbol	_cell_length_c
_atom_site_B_iso_or_equiv	_symmetry_equiv_pos_site_id
_cell_volume	_atom_site_fract_z
_cell_angle_alpha	_cell_formula_units_Z
loop_	data_
_atom_type_symbol	_atom_type_electronegativity *
_atom_type_radius *	_atom_type_ionic_radius *
_atom_type_oxidation_number	

*(Tags with * do not exist in the official CIF specification.)*

Supported space group tokens.

Aea2	Aem2	Ama2	Amm2	C2	C2/c	C2/m
C222	C222_1	Cc	Ccc2	Ccce	Cccm	Cm
Cmc2_1	Cmce	Cmcm	Cmm2	Cmme	Cmmm	F-43c
F-43m	F222	F23	F432	F4_132	Fd-3	Fd-3c
Fd-3m	Fdd2	Fddd	Fm-3	Fm-3c	Fm-3m	Fmm2
Fmmm	I-4	I-42d	I-42m	I-43d	I-43m	I-4c2
I-4m2	I222	I23	I2_12_12_1	I2_13	I4	I4/m
I4/mcm	I4/mmm	I422	I432	I4_1	I4_1/a	I4_1/acd
I4_1/amd	I4_122	I4_132	I4_1cd	I4_1md	I4cm	I4mm
Ia-3	Ia-3d	Iba2	Ibam	Ibca	Im-3	Im-3m
Ima2	Imm2	Imma	Immm	P-1	P-3	P-31c
P-31m	P-3c1	P-3m1	P-4	P-42_1c	P-42_1m	P-42c
P-42m	P-43m	P-43n	P-4b2	P-4c2	P-4m2	P-4n2
P-6	P-62c	P-62m	P-6c2	P-6m2	P1	P2
P2/c	P2/m	P222	P222_1	P23	P2_1	P2_1/c
P2_1/m	P2_12_12	P2_12_12_1	P2_13	P3	P312	P31c
P31m	P321	P3_1	P3_112	P3_121	P3_2	P3_212
P3_221	P3c1	P3m1	P4	P4/m	P4/mbm	P4/mcc
P4/mmm	P4/mnc	P4/n	P4/nbm	P4/ncc	P4/nmm	P4/nnc
P422	P42_12	P4_1	P4_122	P4_12_12	P4_132	P4_2
P4_2/m	P4_2/mbc	P4_2/mcm	P4_2/mmc	P4_2/mnm	P4_2/n	P4_2/nbc
P4_2/ncm	P4_2/nmc	P4_2/nnm	P4_22_12	P4_232	P4_2bc	P4_2cm
P4_2mc	P4_2nm	P4_3	P4_322	P4_32_12	P4_332	P4bm
P4cc	P4mm	P4nc	P6/m	P6/mcc	P6/mmm	P622
P6_1	P6_122	P6_2	P6_222	P6_3	P6_3/m	P6_3/mcm
P6_3/mmc	P6_322	P6_3cm	P6_3mc	P6_4	P6_422	P6_5
P6_522	P6cc	P6mm	Pa-3	Pba2	Pbam	Pban
Pbca	Pbcm	Pbcn	Pc	Pca2_1	Pcc2	Pcca
Pccm	Pccn	Pm	Pm-3	Pm-3m	Pm-3n	Pma2
Pmc2_1	Pmm2	Pmma	Pmmm	Pmmn	Pmn2_1	Pmna
Pn-3	Pn-3m	Pn-3n	Pna2_1	Pnc2	Pnma	Pnn2
Pnna	Pnnm	Pnnn	R-3	R-3c	R-3m	R3
R32	R3c	R3m				

After the model has generated a sequence of tokens representing a CIF file, we perform a post-processing step in which the custom `loop_` section with atomic properties is

removed, and the symmetry equivalent site IDs and positions which match the printed space group are introduced.

2. Model Architecture and Generative Pre-training

The generative pre-training step consists of training a GPT-style transformer model autoregressively. The implementation is based on the nanoGPT project [1]. The model consists of a series of transformer blocks, each consisting of multi-head self-attention and a feed-forward neural network. The input to the model is a sequence of token indices representing the token sequence. The tokens are embedded using a learned embedding table. The token embeddings are combined with learned positional embeddings, to which dropout is applied. The result is passed through a series of transformer blocks. A transformer block consists of causal self-attention [2] and a feed-forward network containing a non-linear layer with GELU activation [3], and dropout. A linear output layer transforms the features produced by the transformer blocks into a vector of logits. A softmax operation is applied to convert the logits into the probabilities of the tokens of the vocabulary, for each position in the output sequence. Weight tying [4] is used: the output layer and the input embedding layer share the same weights. The objective is to minimize the cross-entropy loss between the predicted probability distribution over the vocabulary and the actual next token in the sequence, for all the tokens in the sequence.

Training consists of iteratively sampling sequences from the dataset, performing a forward-pass through the model, computing the loss, and backpropagating the error. We use the AdamW optimizer [5], and apply a cosine decay schedule to the learning rate, from 10^{-3} to 10^{-4} , over the course of training. Gradients were clipped to have a norm of at most 1.0. During each training iteration, we perform 40 gradient accumulation steps, and in each step we randomly sample a batch of 32 sequences. The dataset consists of a single list of all tokens from all CIF files, concatenated together, and the beginning of each sequence is a randomly sampled token from the list. The number of tokens in each sequence is equal to the block size of the model, which is the maximum length of the input sequence the model can process. All models were trained on a single A100 GPU with 80 GB of memory.

2.1 Small Model

The small model consists of 25 million parameters, with 8 transformer blocks, each with 8 attention heads, an embedding size of 512, a block size of 1,024, and dropout with probability $p = 0.1$. To determine the optimal number of training iterations, we train the model using 10% of the dataset as a validation set, and monitor the model’s performance on the validation set, in terms of the cross-entropy loss. It was determined that the model continues to improve beyond 90,000 iterations. Therefore, the final model was trained on the entire dataset for 100,000 iterations (due to computational resource and time constraints).

2.2 Large Model

The large model consists of 200 million parameters, with 16 transformer blocks, each with 16 attention heads, an embedding size of 1,024, a block size of 2,048, and dropout

with probability $p = 0.1$. Due to computational resource and time constraints, we train the large model on the entire dataset for 48,000 iterations. Additionally, the starting point for each sequence is sampled from a pre-compiled list of tokens, each known to be the starting token of a CIF file in the dataset. This approach ensures that each sequence begins at the start of a distinct CIF file.

2.3 CIF File Generation via Random Sampling

CIF files are generated using top- k random sampling. Top- k sampling involves randomly selecting the next token from the top k most likely candidates as predicted by the model. We first apply temperature scaling \mathbf{x}/τ to the logits $\mathbf{x} \in \mathbb{R}^{|\mathcal{V}|}$ at the final position, and keep only the top k logits, where $|\mathcal{V}|$ represents the size of the vocabulary. The top k logits are then converted into normalized probabilities through application of the softmax operation. Finally, the next token is sampled using the given probabilities. More formally,

$$\text{token} \sim \text{softmax} \left(\text{top}_k \left(\frac{\mathbf{x}}{\tau} \right) \right) \quad (1)$$

$$\text{softmax}(x)_i = \frac{e^{x_i}}{\sum_{j=1}^k e^{x_j}} \quad (2)$$

where i represents the i -th element of the vector \mathbf{x} . Tokens are sampled iteratively, each conditioned on the progressively growing sequence of previously sampled tokens, until two consecutive newline tokens are sampled.

3. Validation of Generated CIF Files

To ensure the consistency of the printed information and the chemical sensibility of the implied structure, we conduct a series of validations on the generated CIF file. The procedure is described in Algorithm 1.

First, we require that the chemical formula, which is printed in several locations in the file, is consistent everywhere. Specifically, we ensure that the formula associated with the `_chemical_formula_sum` tag, and the (reduced) formula associated with the `_chemical_formula_structural` tag are consistent with the cell composition in the first line of the file, and with each other.

Next, we require that the printed atom site multiplicity is consistent with the cell composition. This information is printed at the end of the file, and the values are associated with the `_atom_site_type_symbol` and `_atom_site_symmetry_multiplicity` tags.

We also check that the structure’s bond lengths are reasonable. To check if bond lengths are reasonable, we first use a Voronoi-based nearest-neighbour algorithm in `py-matgen` to define which atoms are bonded together; then, we establish expected bond lengths based on the electronegativity difference between the bonded atoms, and their ionic or covalent radii. We compute a bond length reasonableness score, $B \in [0, 1]$, which represents the fraction of bonds which are within 30% of the corresponding expected bond lengths. We classify a structure as having reasonable bond lengths if $B \geq c_{\text{bond}}$, where

$c_{\text{bond}} \in [0, 1]$ is a bond length acceptability score minimum, which in this work is set to 1.0 (i.e. all bond lengths must be within 30% of the expected bond lengths).

Finally, we check that the generated CIF file is consistent in terms of space group. To check if the generated structure is consistent with the printed space group, we use the `SpacegroupAnalyzer` class of the `pymatgen` library, which uses the `spglib` library [6].

Algorithm 1 Check Validity of Generated CIF File

```

1: Input:  $S$ , the contents of the generated CIF file
2: Input:  $c_{\text{bond}}$ , the bond length acceptability score minimum
3: Output: True or False, indicating whether  $S$  is valid
4: if not is_formula_consistent(S) then
5:   return False
6: end if
7: if not is_atom_site_multiplicity_consistent(S) then
8:   return False
9: end if
10:  $B \leftarrow \text{bond\_length\_reasonableness\_score}(S)$ 
11: if  $B < c_{\text{bond}}$  then
12:   return False
13: end if
14: if not is_space_group_consistent(S) then
15:   return False
16: end if
17: return True

```

4. Monte Carlo Tree Search Decoding

To improve the efficiency and quality of sampling from the model, we use the Monte Carlo Tree Search (MCTS) algorithm [7, 8]. Typically, the MCTS algorithm is used in the context of games, and similar decision processes, where the aim is to select an optimal action to perform. Here, we use MCTS to generate a collection of sequences, which should improve (according to some measure of quality) as the algorithm proceeds.

A sequence of tokens can be considered the outcome of following a specific path during the traversal of a tree of tokens, starting from the root and progressing to a leaf. In this framework, each node in the tree represents a token, and each edge denotes the transition from one token to the next in the sequence. By systematically exploring and expanding the most promising paths, MCTS balances exploitation of well-performing token sequences with exploration of new, potentially better sequences. The trade-off between exploitation and exploration is achieved through a principled selection strategy, which guides the search towards areas of the tree that either have high potential or have not been sufficiently explored. As the search progresses, the algorithm builds a more informed representation of the tree, enabling more efficient and higher-quality sampling of token sequences.

The MCTS algorithm is comprised of a sequence of steps performed for a fixed number of iterations. Our implementation is described in Algorithm 2, and a detailed explanation of each step follows.

4.1 Selection

Each node, i , in the tree represents the cumulative context up to that point, akin to constructing a sentence word by word. The first step in every iteration involves descending the tree, from the root node to a leaf node, by selecting the most promising node, i_t , at each level t . To select the node, we use a variant of the PUCT (Predictor-Upper Confidence bound applied to Trees) algorithm [9, 10].

The selection of a node at each level is guided by the statistics accumulated in the tree. The specific node i_t is chosen by maximizing the PUCT value, expressed as $i_t = \underset{i}{\operatorname{argmax}}(\operatorname{PUCT}(i_t))$, where the PUCT value is calculated as:

$$\operatorname{PUCT}(i) = \frac{w_i}{N_i} + c_{\text{puct}} P_i \frac{\sqrt{N_h}}{1 + N_i} \tag{3}$$

where w_i represents the total score accumulated at node i , indicating the node’s past performance. N_i is the number of times node i has been visited, reflecting its exploitation level. P_i is the prior probability of selecting the token leading to node i , given its parent node h , which is provided by the CrystaLLM model. N_h is the total number of visits to the parent node h , and c_{puct} is a constant determining the level of exploration in the PUCT algorithm.

4.2 Expansion

If a node has children that haven’t been added to the tree, a child node is selected randomly and added to the tree. This newly added node is the selected node for the remainder of the iteration. To determine what children a node contains, the CrystaLLM model’s predicted probabilities are used to select the top k tokens. If a child node’s probability exceeds 0.99, it becomes the sole child node. In this case, where a node has only a single child node, the child node is *bypassed*, foregoing the Rollout step. The process then proceeds directly to the Selection step for the child nodes of the bypassed node, continuing the iteration from that point onwards.

4.3 Rollout

The Rollout step involves prompting the CrystaLLM model with the sequence of tokens represented by the selected node in the tree, and then sampling from the model repeatedly, until a terminating condition is reached. The aim is to arrive at a completed structure, which can then be further evaluated and scored.

4.4 Evaluation

Once a sequence has been completed, either by reaching a terminal node through Selection, or through Rollout, it represents the contents of a completed CIF file. The generated CIF contents are then validated (see Supplementary Note 3), and if the generated CIF file is valid, the structure is evaluated using the ALIGNN model of formation energy per atom, to produce a prediction of the structure’s energy, E_f .

4.5 Backpropagation

The outcomes of iterations are accumulated in the tree nodes. All nodes selected during a simulation increase in their visit count, and a score is added to each. The score, $R \in [-1, 1]$, represents the quality of the generated structure. A more positive R represents a better structure.

Because scores are required to be between -1 and 1, and since the range of formation energies is not known for a composition *a priori*, the score for valid structures ($R_{\text{valid}} \in [0, 1]$) is computed using the statistics of the predicted energies over the course of the search:

$$R_{\text{valid}} = \frac{1}{1 + e^{\lambda((E_f - \mu)/\sigma)}} \tag{4}$$

where E_f is the formation energy per atom (eV) according to ALIGNN (Atomistic Line Graph Neural Network) [11], μ is the mean over all of the obtained E_f , σ is the standard deviation over all the obtained E_f , and λ is a constant that determines how responsive the reward is to E_f .

The overall score is computed piecewise:

$$R = \begin{cases} R_{\text{valid}} & \text{if valid,} \\ B - 1 & \text{if bond lengths unreasonable,} \\ -1 & \text{otherwise} \end{cases} \tag{5}$$

where B is the bond length reasonableness score. An invalid structure receives a score of -1, unless it is invalid because of unreasonable bond lengths. In cases where a CIF file is otherwise valid, but the structure contains unreasonable bond lengths, a negative score is assigned that is proportional to the number of unreasonable bonds.

Algorithm 2 Monte Carlo Tree Search Decoding

```

1: Input: trained large language model, LLM
2: Input: number of simulations,  $n$ 
3: Input: tree width,  $k$ 
4: Input: PUCT exploration constant,  $c_{\text{puct}}$ 
5: Input: text prompt,  $P$ 
6: Output: list of valid sequences
7: Initialize tree with root node based on  $P$ 
8: valid_sequences  $\leftarrow$  []
9: for simulation = 1 to  $n$  do
10:   current_node  $\leftarrow$  root
11:   // Select
12:   while not current_node.has_untried_children() and current_node.has_children() do
13:     current_node  $\leftarrow$  select_node(current_node.children, LLM,  $c_{\text{puct}}$ )
14:   end while
15:   // Expand
16:   if current_node.has_untried_children() then
17:     untried_child  $\leftarrow$  select_untried_child_randomly(current_node,  $k$ )
18:     current_node.add_child(untried_child)
19:     current_node  $\leftarrow$  untried_child
20:   end if
21:   // Rollout
22:   complete_sequence  $\leftarrow$  sample_randomly(current_node, LLM)
23:   // Evaluate
24:   score  $\leftarrow$  evaluate_sequence(complete_sequence)
25:   if is_valid(complete_sequence) then
26:     valid_sequences.append(complete_sequence)
27:   end if
28:   // Backpropagate
29:   while current_node is not null do
30:     current_node.visits  $\leftarrow$  current_node.visits + 1
31:     current_node.wins  $\leftarrow$  current_node.wins + score
32:     current_node  $\leftarrow$  current_node.parent
33:   end while
34: end for
35: return valid_sequences

```

Supplementary Tables

Supplementary Table 1: The compounds of the Challenge Set, their sources, and their formation energies per atom, as predicted by the ALIGNN model.

Formula	Source	ALIGNN E_f (eV/atom)
Ba ₂ MnCr	training set	0.906
Ca ₁₀ (PO ₄) ₆ (OH) ₂	training set	-3.029
CH ₃ NH ₃ PbI ₃	training set	-0.358
Co ₂ CO ₃ (OH) ₂	training set	-1.019
CsCuTePt	training set	0.137
Cu ₂ C ₁ O ₅ H ₂	training set	-0.923
Cu ₃ (CO ₃) ₂ (OH) ₂	training set	-0.999
K ₂ AgMoI ₆	training set	-0.639
MgF ₂	training set	-3.782
Mn ₄ (PO ₄) ₃	training set	-2.010
PbCu(OH) ₂ SO ₄	training set	-1.160
Sm ₂ BO ₄	training set	-2.993
AlCu ₂ As(HO) ₁₂	ref. [12]	-1.185
Ba ₂ AuIO ₆	ref. [13]	-1.329
Ba ₂ Fe ₂ F ₉	ref. [14]	-3.073
Ba ₆ Fe ₂ Te ₃ S ₇	ref. [15]	-1.593
Ba ₂ Gd(BO ₃) ₂ F	ref. [16]	-3.325
Ba ₄ GeSb ₂ Se ₁₁	ref. [17]	-1.136
Ba ₃ GeTeS ₄	ref. [18]	-1.721
Ba ₂ HfF ₈	ref. [19]	-4.177
BaY ₁₆ Si ₄ O ₃₃	ref. [20]	-3.702
Ba ₉ Yb ₂ (SiO ₄) ₆	ref. [21]	-3.245
Ca ₂ Bi ₂ O ₇	ref. [22]	-1.952
CaFe ₆ Ge ₆	ref. [23]	-0.200
CaHPO ₃	ref. [24]	-2.445
CaPt ₄ P ₆	ref. [25]	-0.826
Ca ₂ Te ₃ O ₈	ref. [26]	-1.911
CaZnV ₂ O ₆	ref. [27]	-2.573
Ce ₆ Cd ₂₃ Te	ref. [28]	-0.295
Cs ₂ Al ₂ O ₃ F ₂	ref. [29]	-3.116
Cs ₈ Cu ₃ Si ₁₄ O ₃₅	ref. [30]	-2.561
Cs ₃ LuSi ₃ O ₉	ref. [31]	-2.970
Cu ₄ FeGe ₂ S ₇	ref. [32]	-0.368
Eu ₂ FeGe ₂ OS ₆	ref. [33]	-1.265
HgB ₂ S ₄	ref. [34]	-0.340
Ho ₂ Ir ₃ Si ₅	ref. [35]	-0.885
KScP ₂ O ₇	ref. [36]	-2.801
K ₂ Sr ₄ (PO ₃) ₁₀	ref. [37]	-2.626
K ₆ Zn(CO ₃) ₄	ref. [38]	-1.817
La ₄ Ga ₂ S ₈ O ₃	ref. [39]	-2.171
LaScSe ₃	ref. [40]	-1.912
Li ₉ Al ₄ Sn ₅	ref. [41]	-0.173
LiBa ₂ AlO ₄	ref. [42]	-2.965
Li ₂ GeS ₃	ref. [43]	-0.989
LiMnBi	ref. [44]	-0.004
LiTa ₂ NiSe ₅	ref. [45]	-0.842
Mg ₇ Pt ₄ Ge ₄	ref. [46]	-0.676
NaGdSi ₂ O ₆	ref. [47]	-3.074
Na ₂ Hf(BO ₃) ₂	ref. [48]	-2.834
Na ₆ Li ₄ WO ₄ (CO ₃) ₄	ref. [49]	-2.010
NaMgV ₅ (H ₅ O ₆) ₄	ref. [50]	-1.540
Na ₅ Mn ₄ P ₄ H ₄ (O ₉ F ₂) ₂	ref. [51]	-2.140
NaSbSe ₂ O ₇	ref. [52]	-1.213
NaSb ₂ TeO ₇	ref. [52]	-1.502
Na ₄ Sn ₂ Ge ₅ O ₁₆	ref. [53]	-1.862
Na ₃ Te ₂ (FeO ₄) ₃	ref. [54]	-1.430
Nd ₃ BSi ₂ O ₁₀	ref. [55]	-3.352
Ni ₃ Te ₂ O ₂ (PO ₄) ₂ (OH) ₄	ref. [56]	-1.302
RbNiFe(PO ₄) ₂	ref. [57]	-1.930
Rb ₃ SnCl ₇	ref. [58]	-1.434
Sr ₂ Bi ₂ O ₇	ref. [22]	-1.963
SrCo ₄ (OH)(PO ₄) ₃	ref. [59]	-1.861
Sr(ClO ₄) ₂	ref. [60]	-0.593
Sr ₆ Ge ₃ OSe ₁₁	ref. [61]	-1.145
Tb ₃ S ₃ BO ₃	ref. [62]	-2.857
Tb ₃ TeBO ₉	ref. [63]	-2.433
YbMn ₆ Sn ₆	ref. [64]	-0.070
Zn ₂ (HTeO ₃)(AsO ₄)	ref. [65]	-1.290
Zn ₂ BS ₃ Br	ref. [66]	-0.692
Zn ₄ CuH ₆ (CO ₆) ₂	ref. [67]	-1.251

Supplementary Table 2: Performance of the small model on the Challenge Set. No space group was included in the prompt.

Composition	Mean E_f	Min E_f	% Valid	Any match true?
AlCu ₂ As(HO) ₁₂	-0.366	-0.568	19	no
Ba ₂ AuIO ₆	-1.546	-1.599	96	no
Ba ₂ Fe ₂ F ₉	-2.324	-2.598	14	no
Ba ₂ Gd(BO ₃) ₂ F	-2.091	-2.578	12	no
Ba ₂ HfF ₈	-3.427	-3.940	35	yes
Ba ₂ MnCr	0.985	0.611	100	yes
Ba ₃ GeTeS ₄	-0.985	-1.366	33	no
Ba ₄ GeSb ₂ Se ₁₁	-0.707	-0.859	18	no
Ba ₆ Fe ₂ Te ₃ S ₇	-0.743	-0.989	15	no
Ba ₉ Yb ₂ (SiO ₄) ₆	-1.524	-1.524	1	no
BaY ₁₆ Si ₄ O ₃₃	-	-	0	no
CH ₃ NH ₃ PbI ₃	0.339	0.110	11	no
Ca ₁₀ (PO ₄) ₆ (OH) ₂	-	-	0	no
Ca ₂ Bi ₂ O ₇	-1.759	-1.889	100	yes
Ca ₂ Te ₃ O ₈	-	-	0	no
CaFe ₆ Ge ₆	0.232	-0.207	26	yes
CaHPO ₃	-1.471	-2.030	18	no
CaPt ₄ P ₆	-0.330	-0.671	29	no
CaZnV ₂ O ₆	-1.981	-2.551	32	yes
Ce ₆ Cd ₂₃ Te	-0.309	-0.369	91	yes
Co ₂ CO ₃ (OH) ₂	-0.180	-0.527	29	no
Cs ₂ Al ₂ O ₃ F ₂	-1.959	-2.455	16	no
Cs ₃ LuSi ₃ O ₉	-	-	0	no
Cs ₈ Cu ₃ Si ₁₄ O ₃₅	-	-	0	no
CsCuTePt	0.136	0.092	100	yes
Cu ₂ C ₁ O ₅ H ₂	-0.279	-0.708	48	no
Cu ₃ (CO ₃) ₂ (OH) ₂	-0.164	-0.479	22	no
Cu ₄ FeGe ₂ S ₇	-0.060	-0.266	39	no
Eu ₂ FeGe ₂ OS ₆	-0.924	-1.265	45	yes
HgB ₂ S ₄	0.336	0.051	18	no
Ho ₂ Ir ₃ Si ₅	-0.883	-0.890	98	yes
K ₂ AgMol ₆	-0.638	-0.643	100	yes
K ₂ Sr ₄ (PO ₃) ₁₀	-	-	0	no
K ₆ Zn(CO ₃) ₄	-	-	0	no
KScP ₂ O ₇	-2.626	-2.794	81	yes
La ₄ Ga ₂ S ₈ O ₃	-1.024	-1.288	5	no
LaScSe ₃	-1.879	-1.978	98	yes
Li ₂ GeS ₃	-0.554	-0.960	44	no
Li ₆ Al ₄ Sn ₅	-0.122	-0.189	2	no
LiBa ₂ AlO ₄	-1.837	-2.053	3	no
LiMnBi	0.130	0.075	100	no
LiTa ₂ NiSe ₅	-0.693	-0.847	71	no
Mg ₇ Pt ₄ Ge ₄	-0.263	-0.521	23	no
MgF ₂	-3.512	-3.811	93	yes
Mn ₄ (PO ₄) ₃	-1.750	-2.014	16	yes
Na ₂ Hf(BO ₃) ₂	-2.766	-2.835	69	yes
Na ₃ Te ₂ (FeO ₄) ₃	-1.362	-1.455	97	yes
Na ₄ Sn ₂ Ge ₅ O ₁₆	-0.741	-0.867	2	no
Na ₅ Mn ₄ P ₄ H ₄ (O ₉ F ₂) ₂	-0.917	-0.948	2	no
Na ₆ Li ₄ WO ₄ (CO ₃) ₄	-	-	0	no
NaGdSi ₂ O ₆	-2.721	-3.060	63	no
NaMgV ₅ (H ₅ O ₆) ₄	-	-	0	no
NaSb ₂ TeO ₇	-0.855	-1.049	3	no
NaSbSe ₂ O ₇	-0.528	-0.813	10	no
Nd ₃ BSi ₂ O ₁₀	-	-	0	no
Ni ₃ Te ₂ O ₂ (PO ₄) ₂ (OH) ₄	-0.538	-0.824	28	no
PbCu(OH) ₂ SO ₄	-0.362	-0.731	51	no
Rb ₃ SnCl ₇	-1.329	-1.506	52	yes
RbNiFe(PO ₄) ₂	-0.896	-1.315	9	no
Sm ₂ BO ₄	-2.978	-3.011	92	yes
Sr(ClO ₄) ₂	-0.044	-0.357	20	no
Sr ₂ Bi ₂ O ₇	-1.729	-1.931	97	yes
Sr ₆ Ge ₃ OSe ₁₁	-0.768	-1.017	4	no
SrCo ₄ (OH)(PO ₄) ₃	-0.868	-0.868	1	no
Tb ₃ S ₃ BO ₃	-1.173	-2.026	45	no
Tb ₃ TeBO ₉	-2.274	-2.477	73	yes
YbMn ₆ Sn ₆	-0.042	-0.071	100	yes
Zn ₂ (HTeO ₃)(AsO ₄)	-0.556	-1.209	26	no
Zn ₂ BS ₃ Br	-0.077	-0.602	77	no
Zn ₄ CuH ₆ (CO ₆) ₂	-0.207	-0.641	19	no

Supplementary Table 3: Performance of the small model on the Challenge Set. The space group was included in the prompt.

Composition	Mean E_f	Min E_f	% Valid	Any match true?
AlCu ₂ As(HO) ₁₂	-0.403	-0.735	50	no
Ba ₂ AuIO ₆	-1.410	-1.608	74	yes
Ba ₂ Fe ₂ F ₉	-2.163	-2.416	8	no
Ba ₂ Gd(BO ₃) ₂ F	-1.869	-2.266	12	no
Ba ₂ HfF ₈	-3.565	-4.082	40	yes
Ba ₂ MnCr	1.020	0.833	100	yes
Ba ₃ GeTeS ₄	-1.148	-1.589	42	yes
Ba ₄ GeSb ₂ Se ₁₁	-0.707	-0.888	16	no
Ba ₆ Fe ₂ Te ₃ S ₇	-0.563	-1.020	6	no
Ba ₉ Yb ₂ (SiO ₄) ₆	-	-	0	no
BaY ₁₆ Si ₄ O ₃₃	-	-	0	no
CH ₃ NH ₃ PbI ₃	0.523	0.068	47	no
Ca ₁₀ (PO ₄) ₆ (OH) ₂	-	-	0	no
Ca ₂ Bi ₂ O ₇	-1.758	-1.887	100	yes
Ca ₂ Te ₃ O ₈	-	-	0	no
CaFe ₆ Ge ₆	0.013	-0.175	8	yes
CaHPO ₃	-1.175	-1.917	18	no
CaPt ₄ P ₆	-0.346	-0.645	31	yes
CaZnV ₂ O ₆	-2.197	-2.583	51	yes
Ce ₆ Cd ₂₃ Te	-0.315	-0.357	89	yes
Co ₂ CO ₃ (OH) ₂	-0.086	-0.304	21	no
Cs ₂ Al ₂ O ₃ F ₂	-2.028	-2.813	20	no
Cs ₃ LuSi ₃ O ₉	-1.505	-1.927	3	no
Cs ₈ Cu ₃ Si ₁₄ O ₃₅	-	-	0	no
CsCuTePt	0.136	0.128	100	yes
Cu ₂ C ₁ O ₅ H ₂	-0.135	-0.541	29	no
Cu ₃ (CO ₃) ₂ (OH) ₂	-0.157	-0.581	29	no
Cu ₄ FeGe ₂ S ₇	-0.055	-0.314	57	no
Eu ₂ FeGe ₂ OS ₆	-1.151	-1.290	52	yes
HgB ₂ S ₄	0.272	-0.043	25	no
Ho ₂ Ir ₃ Si ₅	-0.881	-0.892	97	yes
K ₂ AgMol ₆	-0.638	-0.643	100	yes
K ₂ Sr ₄ (PO ₃) ₁₀	-	-	0	no
K ₆ Zn(CO ₃) ₄	-0.820	-0.820	1	no
KSP ₂ O ₇	-2.703	-2.798	77	yes
La ₄ Ga ₂ S ₈ O ₃	-1.028	-1.460	19	no
LaScSe ₃	-1.912	-1.975	99	yes
Li ₂ GeS ₃	-0.590	-0.962	24	yes
Li ₆ Al ₄ Sn ₅	-0.117	-0.216	9	no
LiBa ₂ AlO ₄	-1.392	-2.074	15	no
LiMnBi	0.123	-0.037	100	yes
LiTa ₂ NiSe ₅	-0.249	-0.555	48	no
Mg ₇ Pt ₄ Ge ₄	-0.412	-0.602	19	no
MgF ₂	-3.774	-3.810	100	yes
Mn ₄ (PO ₄) ₃	-1.574	-1.997	34	yes
Na ₂ Hf(BO ₃) ₂	-2.802	-2.839	99	yes
Na ₃ Te ₂ (FeO ₄) ₃	-1.363	-1.458	94	yes
Na ₄ Sn ₂ Ge ₅ O ₁₆	-0.974	-1.215	3	no
Na ₅ Mn ₄ P ₄ H ₄ (O ₉ F ₂) ₂	-1.053	-1.053	1	no
Na ₆ Li ₄ WO ₄ (CO ₃) ₄	-0.712	-1.012	3	no
NaGdSi ₂ O ₆	-1.696	-2.620	12	no
NaMgV ₅ (H ₅ O ₆) ₄	-	-	0	no
NaSb ₂ TeO ₇	-0.752	-0.752	1	no
NaSbSe ₂ O ₇	-0.543	-1.081	21	no
Nd ₃ BSi ₂ O ₁₀	-1.915	-1.915	1	no
Ni ₃ Te ₂ O ₂ (PO ₄) ₂ (OH) ₄	-0.582	-0.947	24	no
PbCu(OH) ₂ SO ₄	-0.378	-0.776	46	no
Rb ₃ SnCl ₇	-1.454	-1.538	76	yes
RbNiFe(PO ₄) ₂	-0.887	-1.082	4	no
Sm ₂ BO ₄	-2.981	-3.011	92	yes
Sr(ClO ₄) ₂	-0.064	-0.158	2	no
Sr ₂ Bi ₂ O ₇	-1.762	-1.968	100	yes
Sr ₆ Ge ₃ OSe ₁₁	-0.672	-0.912	5	no
SrCo ₄ (OH)(PO ₄) ₃	-	-	0	no
Tb ₃ S ₃ BO ₃	-1.281	-1.477	6	no
Tb ₃ TeBO ₉	-2.251	-2.466	84	yes
YbMn ₆ Sn ₆	-0.052	-0.066	99	yes
Zn ₂ (HTeO ₃)(AsO ₄)	-0.865	-1.210	50	no
Zn ₂ BS ₃ Br	0.111	-0.181	56	no
Zn ₄ CuH ₆ (CO ₆) ₂	-0.205	-0.302	4	no

Supplementary Table 4: Performance of the large model on the Challenge Set. No space group was included in the prompt.

Composition	Mean E_f	Min E_f	% Valid	Any match true?
AlCu ₂ As(HO) ₁₂	-0.333	-0.620	33	no
Ba ₂ AuIO ₆	-1.541	-1.589	65	no
Ba ₂ Fe ₂ F ₉	-2.292	-2.701	10	no
Ba ₂ Gd(BO ₃) ₂ F	-1.996	-2.207	5	no
Ba ₂ HfF ₈	-3.522	-4.109	23	yes
Ba ₂ MnCr	1.027	0.593	100	yes
Ba ₃ GeTeS ₄	-1.346	-1.681	71	yes
Ba ₄ GeSb ₂ Se ₁₁	-1.085	-1.137	82	yes
Ba ₆ Fe ₂ Te ₃ S ₇	-0.722	-1.036	12	no
Ba ₉ Yb ₂ (SiO ₄) ₆	-2.970	-3.244	41	yes
BaY ₁₆ Si ₄ O ₃₃	-	-	0	no
CH ₃ NH ₃ PbI ₃	-0.352	-0.358	100	yes
Ca ₁₀ (PO ₄) ₆ (OH) ₂	-2.999	-3.003	93	no
Ca ₂ Bi ₂ O ₇	-1.751	-1.888	100	yes
Ca ₂ Te ₃ O ₈	-	-	0	no
CaFe ₆ Ge ₆	0.207	-0.003	7	no
CaHPO ₃	-1.165	-1.504	14	no
CaPt ₄ P ₆	-0.437	-0.756	34	yes
CaZnV ₂ O ₆	-2.394	-2.585	64	yes
Ce ₆ Cd ₂₃ Te	-0.325	-0.342	93	yes
Co ₂ CO ₃ (OH) ₂	-1.009	-1.019	100	yes
Cs ₂ Al ₂ O ₃ F ₂	-2.114	-2.685	14	no
Cs ₃ LuSi ₃ O ₉	-2.815	-2.946	19	no
Cs ₈ Cu ₃ Si ₁₄ O ₃₅	-	-	0	no
CsCuTePt	0.140	0.131	100	yes
Cu ₂ C ₁ O ₅ H ₂	-0.924	-0.933	99	yes
Cu ₃ (CO ₃) ₂ (OH) ₂	-1.007	-1.018	95	yes
Cu ₄ FeGe ₂ S ₇	-0.073	-0.282	53	no
Eu ₂ FeGe ₂ OS ₆	-0.824	-1.264	31	yes
HgB ₂ S ₄	0.380	0.154	9	no
Ho ₂ Ir ₃ Si ₅	-0.880	-0.887	99	yes
K ₂ AgMol ₆	-0.637	-0.644	100	yes
K ₂ Sr ₄ (PO ₃) ₁₀	-	-	0	no
K ₆ Zn(CO ₃) ₄	-	-	0	no
KScP ₂ O ₇	-2.761	-2.795	100	yes
La ₄ Ga ₂ S ₈ O ₃	-1.174	-1.421	5	no
LaScSe ₃	-1.866	-1.940	96	yes
Li ₂ GeS ₃	-0.735	-0.917	23	no
Li ₆ Al ₄ Sn ₅	-0.112	-0.218	14	no
LiBa ₂ AlO ₄	-2.000	-2.628	9	no
LiMnBi	0.124	-0.055	99	yes
LiTa ₂ NiSe ₅	-0.436	-0.844	58	no
Mg ₇ Pt ₄ Ge ₄	-	-	0	no
MgF ₂	-3.380	-3.803	93	yes
Mn ₄ (PO ₄) ₃	-1.980	-2.020	81	no
Na ₂ Hf(BO ₃) ₂	-2.567	-2.829	75	yes
Na ₃ Te ₂ (FeO ₄) ₃	-1.430	-1.460	100	yes
Na ₄ Sn ₂ Ge ₅ O ₁₆	-1.023	-1.034	3	no
Na ₅ Mn ₄ P ₄ H ₄ (O ₉ F ₂) ₂	-1.094	-1.151	2	no
Na ₆ Li ₄ WO ₄ (CO ₃) ₄	-	-	0	no
NaGdSi ₂ O ₆	-2.962	-3.083	73	yes
NaMgV ₅ (H ₅ O ₆) ₄	-	-	0	no
NaSb ₂ TeO ₇	-0.689	-0.943	4	no
NaSbSe ₂ O ₇	-0.497	-0.698	9	no
Nd ₃ BSi ₂ O ₁₀	-3.338	-3.374	86	yes
Ni ₃ Te ₂ O ₂ (PO ₄) ₂ (OH) ₄	-0.502	-1.010	31	no
PbCu(OH) ₂ SO ₄	-1.153	-1.160	98	yes
Rb ₃ SnCl ₇	-1.438	-1.481	6	yes
RbNiFe(PO ₄) ₂	-1.011	-1.557	9	no
Sm ₂ BO ₄	-2.986	-3.006	95	yes
Sr(ClO ₄) ₂	-0.170	-0.348	12	yes
Sr ₂ Bi ₂ O ₇	-1.672	-1.874	100	yes
Sr ₆ Ge ₃ OSe ₁₁	-0.870	-0.870	1	no
SrCo ₄ (OH)(PO ₄) ₃	-	-	0	no
Tb ₃ S ₃ BO ₃	-1.792	-2.195	13	no
Tb ₃ TeBO ₉	-1.994	-2.347	71	yes
YbMn ₆ Sn ₆	-0.056	-0.067	100	yes
Zn ₂ (HTeO ₃)(AsO ₄)	-0.470	-0.671	22	no
Zn ₂ BS ₃ Br	0.024	-0.510	56	no
Zn ₄ CuH ₆ (CO ₆) ₂	-0.181	-0.550	18	no

Supplementary Table 5: Performance of the large model on the Challenge Set. The space group was included in the prompt.

Composition	Mean E_f	Min E_f	% Valid	Any match true?
AlCu ₂ As(HO) ₁₂	-0.343	-0.540	30	no
Ba ₂ AuIO ₆	-1.411	-1.572	65	yes
Ba ₂ Fe ₂ F ₉	-2.175	-2.503	4	no
Ba ₂ Gd(BO ₃) ₂ F	-2.505	-2.739	4	no
Ba ₂ HfF ₈	-3.466	-4.058	47	yes
Ba ₂ MnCr	1.066	0.862	100	yes
Ba ₃ GeTeS ₄	-1.520	-1.681	70	yes
Ba ₄ GeSb ₂ Se ₁₁	-1.100	-1.138	75	yes
Ba ₆ Fe ₂ Te ₃ S ₇	-0.724	-1.068	31	no
Ba ₉ Yb ₂ (SiO ₄) ₆	-2.949	-3.241	38	yes
BaY ₁₆ Si ₄ O ₃₃	-	-	0	no
CH ₃ NH ₃ PbI ₃	-0.358	-0.358	100	yes
Ca ₁₀ (PO ₄) ₆ (OH) ₂	-1.634	-1.634	1	no
Ca ₂ Bi ₂ O ₇	-1.745	-1.864	100	yes
Ca ₂ Te ₃ O ₈	-	-	0	no
CaFe ₆ Ge ₆	0.079	-0.063	12	no
CaHPO ₃	-1.159	-1.857	22	no
CaPt ₄ P ₆	-0.475	-0.813	33	yes
CaZnV ₂ O ₆	-2.310	-2.571	79	yes
Ce ₆ Cd ₂₃ Te	-0.327	-0.360	99	yes
Co ₂ CO ₃ (OH) ₂	-1.009	-1.020	100	yes
Cs ₂ Al ₂ O ₃ F ₂	-2.190	-2.831	18	no
Cs ₃ LuSi ₃ O ₉	-	-	0	no
Cs ₈ Cu ₃ Si ₁₄ O ₃₅	-	-	0	no
CsCuTePt	0.139	0.131	100	yes
Cu ₂ C ₁ O ₅ H ₂	-0.923	-0.929	99	yes
Cu ₃ (CO ₃) ₂ (OH) ₂	-0.260	-0.634	27	no
Cu ₄ FeGe ₂ S ₇	-0.104	-0.215	60	no
Eu ₂ FeGe ₂ OS ₆	-1.123	-1.300	91	yes
HgB ₂ S ₄	0.191	0.040	10	no
Ho ₂ Ir ₃ Si ₅	-0.881	-0.886	99	yes
K ₂ AgMol ₆	-0.637	-0.643	100	yes
K ₂ Sr ₄ (PO ₃) ₁₀	-	-	0	no
K ₆ Zn(CO ₃) ₄	-0.857	-0.997	2	no
KScP ₂ O ₇	-2.759	-2.803	98	yes
La ₄ Ga ₂ S ₈ O ₃	-1.235	-1.290	4	no
LaScSe ₃	-1.880	-1.960	97	yes
Li ₂ GeS ₃	-0.629	-0.930	15	yes
Li ₆ Al ₄ Sn ₅	-0.149	-0.249	37	no
LiBa ₂ AlO ₄	-1.621	-2.328	37	no
LiMnBi	0.066	0.001	100	yes
LiTa ₂ NiSe ₅	-0.166	-0.493	50	no
Mg ₇ Pt ₄ Ge ₄	-0.448	-0.591	50	no
MgF ₂	-3.783	-3.808	100	yes
Mn ₄ (PO ₄) ₃	-1.903	-2.021	84	yes
Na ₂ Hf(BO ₃) ₂	-2.802	-2.857	98	yes
Na ₃ Te ₂ (FeO ₄) ₃	-1.430	-1.456	100	yes
Na ₄ Sn ₂ Ge ₅ O ₁₆	-0.946	-1.214	15	no
Na ₅ Mn ₄ P ₄ H ₄ (O ₉ F ₂) ₂	-0.902	-1.123	2	no
Na ₆ Li ₄ WO ₄ (CO ₃) ₄	-0.914	-1.307	3	no
NaGdSi ₂ O ₆	-3.011	-3.083	79	yes
NaMgV ₅ (H ₅ O ₆) ₄	-	-	0	no
NaSb ₂ TeO ₇	-0.717	-0.892	4	no
NaSbSe ₂ O ₇	-0.575	-0.983	21	no
Nd ₃ BSi ₂ O ₁₀	-3.363	-3.372	76	yes
Ni ₃ Te ₂ O ₂ (PO ₄) ₂ (OH) ₄	-0.588	-0.851	21	no
PbCu(OH) ₂ SO ₄	-1.138	-1.161	99	yes
Rb ₃ SnCl ₇	-1.431	-1.521	51	yes
RbNiFe(PO ₄) ₂	-0.966	-1.218	6	no
Sm ₂ BO ₄	-2.986	-3.005	96	yes
Sr(ClO ₄) ₂	-0.184	-0.337	11	yes
Sr ₂ Bi ₂ O ₇	-1.671	-1.861	100	yes
Sr ₆ Ge ₃ OSe ₁₁	-0.698	-0.836	6	no
SrCo ₄ (OH)(PO ₄) ₃	-0.506	-0.506	1	no
Tb ₃ S ₃ BO ₃	-1.246	-1.410	7	no
Tb ₃ TeBO ₉	-2.035	-2.443	78	yes
YbMn ₆ Sn ₆	-0.055	-0.070	100	yes
Zn ₂ (HTeO ₃)(AsO ₄)	-0.389	-0.889	34	no
Zn ₂ BS ₃ Br	0.167	-0.245	21	no
Zn ₄ CuH ₆ (CO ₆) ₂	-0.065	-0.126	3	no

Supplementary Table 6: MCTS results for the small model. No space group was included in the prompt.

Composition	Algorithm	Best E_f	Best Iter.	Mean E_f	% Valid
Ba ₂ Fe ₂ F ₉	Random	-2.787	10	-2.273	14.50
	MCTS	-2.812	570	-2.359	42.10
Ba ₂ Gd(BO ₃) ₂ F	Random	-2.950	943	-2.121	18.50
	MCTS	-2.992	699	-2.104	27.20
Ba ₄ GeSb ₂ Se ₁₁	Random	-0.928	571	-0.704	16.70
	MCTS	-0.925	509	-0.735	37.40
Ba ₆ Fe ₂ Te ₃ S ₇	Random	-1.123	110	-0.689	16.60
	MCTS	-1.216	510	-0.730	17.40
Ba ₉ Yb ₂ (SiO ₄) ₆	Random	-2.712	834	-1.907	1.40
	MCTS	-2.777	534	-1.815	2.10
CH ₃ NH ₃ PbI ₃	Random	-0.027	257	0.431	10.60
	MCTS	-0.199	611	0.448	52.80
CaHPO ₃	Random	-2.048	55	-1.397	18.10
	MCTS	-2.247	367	-1.596	59.30
Cs ₂ Al ₂ O ₃ F ₂	Random	-2.825	283	-1.972	22.90
	MCTS	-2.922	651	-2.051	37.90
HgB ₂ S ₄	Random	-0.142	765	0.284	16.90
	MCTS	-0.212	282	0.270	34.20
La ₄ Ga ₂ S ₈ O ₃	Random	-1.404	696	-1.016	3.30
	MCTS	-1.495	27	-1.094	5.70
Li ₉ Al ₄ Sn ₅	Random	-0.225	409	-0.147	1.50
	MCTS	-0.231	123	-0.143	4.60
LiBa ₂ AlO ₄	Random	-2.683	667	-1.728	4.40
	MCTS	-2.504	281	-1.854	61.20
Mn ₄ (PO ₄) ₃	Random	-2.029	122	-1.787	22.00
	MCTS	-2.045	632	-1.946	68.90
Na ₄ Sn ₂ Ge ₅ O ₁₆	Random	-1.126	40	-0.863	2.50
	MCTS	-1.264	848	-0.915	8.70
Na ₅ Mn ₄ P ₄ H ₄ (O ₉ F ₂) ₂	Random	-1.510	660	-1.020	2.40
	MCTS	-1.531	335	-0.979	1.10
NaSb ₂ TeO ₇	Random	-1.292	64	-0.851	4.50
	MCTS	-1.391	787	-0.875	25.50
NaSbSe ₂ O ₇	Random	-0.969	795	-0.473	11.50
	MCTS	-1.108	792	-0.658	42.70
RbNiFe(PO ₄) ₂	Random	-1.465	197	-0.835	4.70
	MCTS	-1.599	699	-1.044	6.80
Sr ₆ Ge ₃ OSe ₁₁	Random	-0.974	658	-0.716	1.50
	MCTS	-1.214	207	-0.910	31.10
SrCo ₄ (OH)(PO ₄) ₃	Random	-1.245	493	-0.845	2.20
	MCTS	-1.223	102	-0.674	9.90

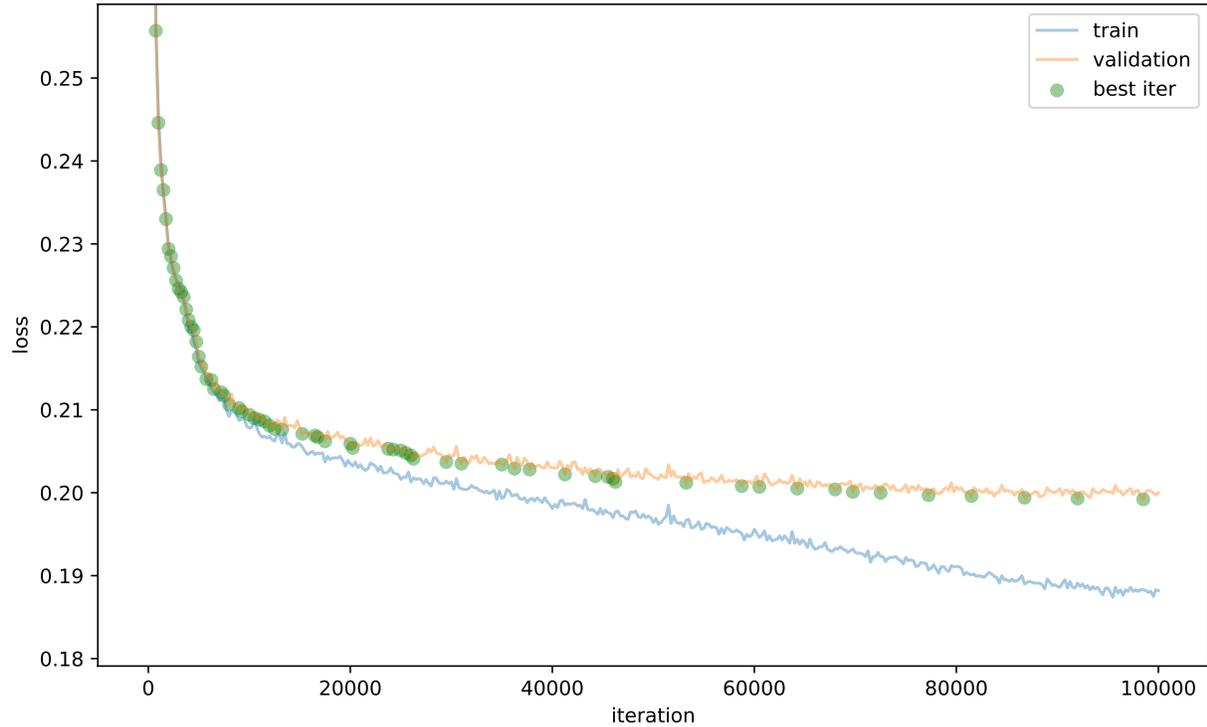
Supplementary Table 7: MCTS results for the small model. The space group was included in the prompt.

Composition	Algorithm	Best E_f	Best Iter.	Mean E_f	% Valid
Ba ₂ Fe ₂ F ₉	Random	-2.523	822	-2.213	8.10
	MCTS	-2.642	188	-2.204	8.30
Ba ₂ Gd(BO ₃) ₂ F	Random	-2.761	906	-1.922	7.60
	MCTS	-2.644	461	-1.956	10.20
Ba ₄ GeSb ₂ Se ₁₁	Random	-0.904	490	-0.677	15.30
	MCTS	-0.923	450	-0.725	25.30
Ba ₆ Fe ₂ Te ₃ S ₇	Random	-1.110	297	-0.632	6.20
	MCTS	-1.137	160	-0.711	24.60
CaFe ₆ Ge ₆	Random	-0.205	353	0.018	9.30
	MCTS	-0.209	675	0.065	8.70
Cs ₃ LuSi ₃ O ₉	Random	-1.967	640	-1.491	1.60
	MCTS	-1.953	325	-1.514	1.40
K ₆ Zn(CO ₃) ₄	Random	-1.308	143	-0.694	1.50
	MCTS	-0.958	487	-0.673	1.40
Li ₉ Al ₄ Sn ₅	Random	-0.247	926	-0.115	10.50
	MCTS	-0.290	559	-0.134	42.50
LiBa ₂ AlO ₄	Random	-2.395	365	-1.370	15.60
	MCTS	-2.380	125	-1.319	16.40
Na ₄ Sn ₂ Ge ₅ O ₁₆	Random	-1.271	849	-0.907	5.80
	MCTS	-1.432	119	-0.920	10.50
Na ₅ Mn ₄ P ₄ H ₄ (O ₉ F ₂) ₂	Random	-1.567	463	-1.105	2.10
	MCTS	-1.472	671	-1.076	3.50
Na ₆ Li ₄ WO ₄ (CO ₃) ₄	Random	-1.488	675	-0.677	6.80
	MCTS	-1.203	694	-0.722	2.90
NaGdSi ₂ O ₆	Random	-2.486	128	-1.686	12.80
	MCTS	-2.591	893	-1.669	9.60
NaSb ₂ TeO ₇	Random	-0.902	616	-0.571	1.40
	MCTS	-1.025	881	-0.667	5.40
Nd ₃ BSi ₂ O ₁₀	Random	-2.910	740	-2.234	1.60
	MCTS	-2.777	411	-2.372	0.50
RbNiFe(PO ₄) ₂	Random	-1.477	104	-0.926	4.70
	MCTS	-1.555	73	-0.961	8.50
Sr(ClO ₄) ₂	Random	-0.278	125	-0.043	2.80
	MCTS	-0.314	633	-0.051	3.10
Sr ₆ Ge ₃ OSe ₁₁	Random	-1.155	679	-0.763	3.80
	MCTS	-1.358	158	-0.874	3.00
Tb ₃ S ₃ BO ₃	Random	-1.877	487	-1.358	6.10
	MCTS	-1.915	199	-1.327	5.70
Zn ₄ CuH ₆ (CO ₆) ₂	Random	-0.350	834	-0.139	3.70
	MCTS	-0.575	815	-0.255	4.00

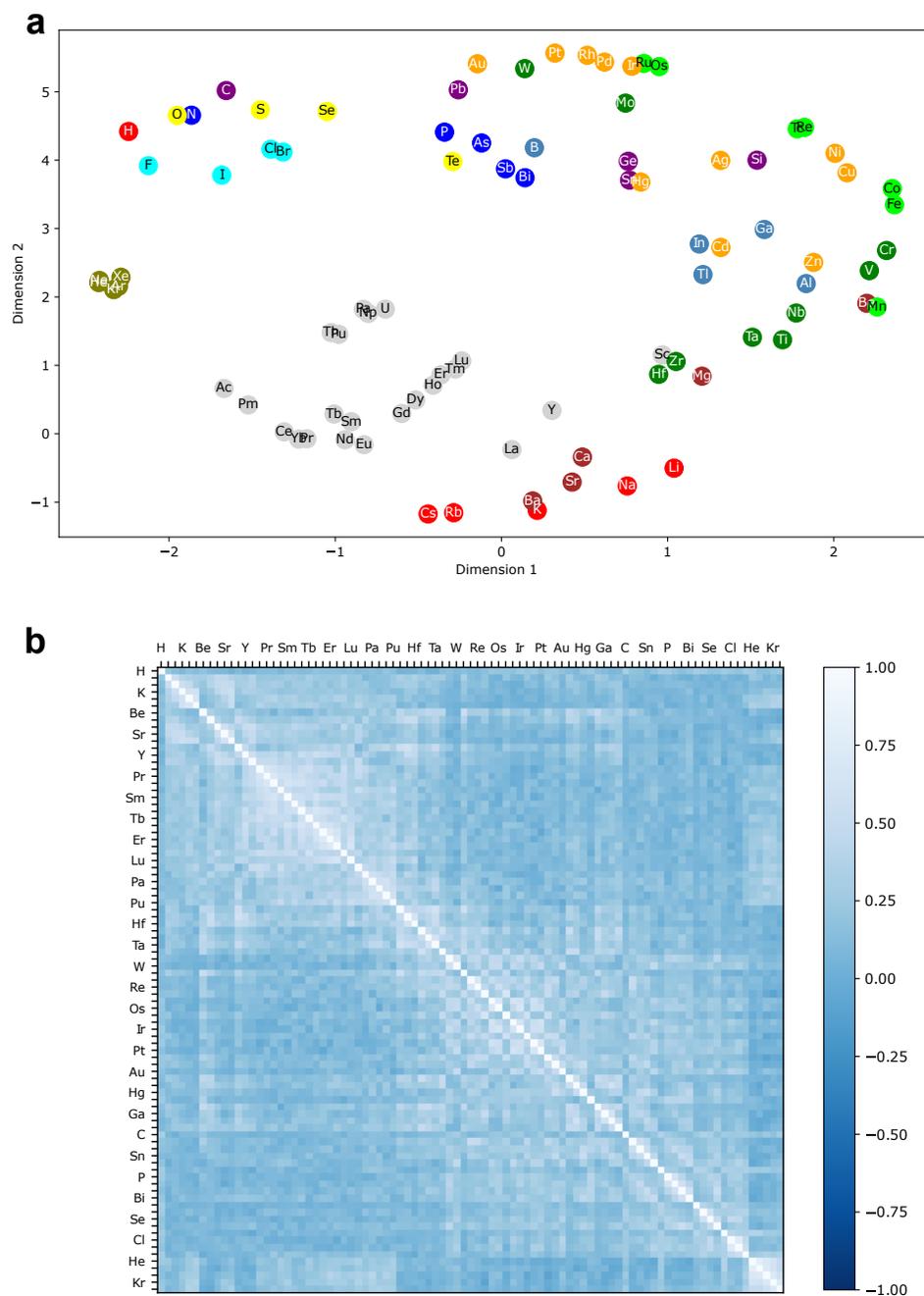
Supplementary Figures



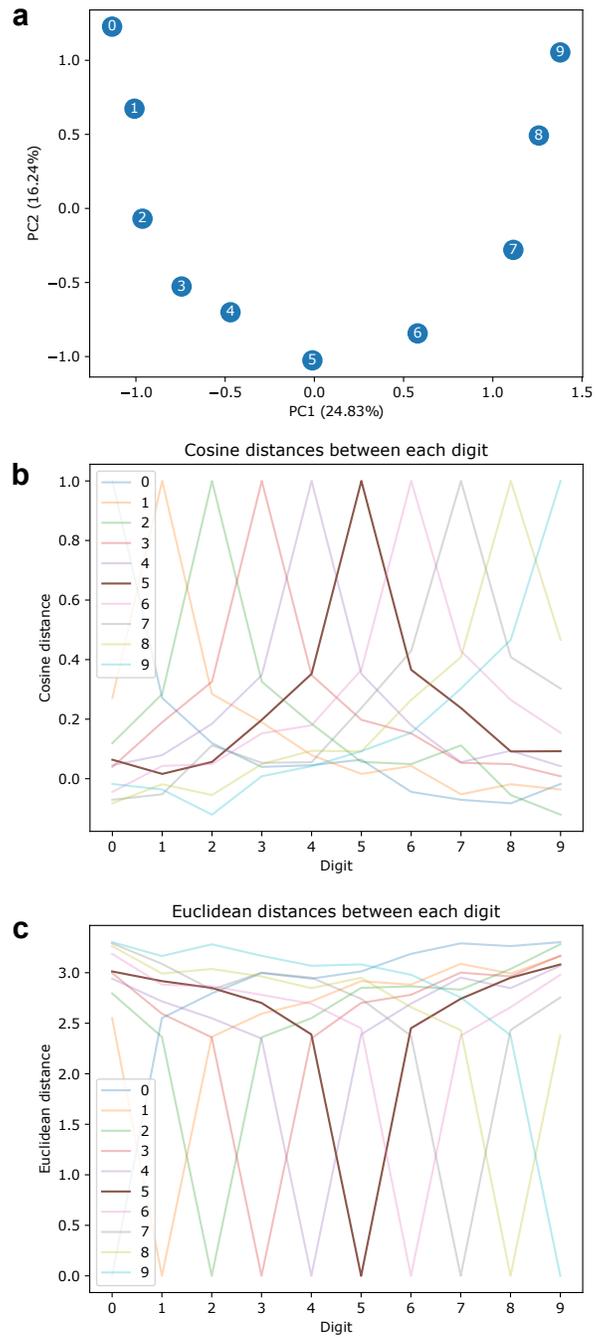
Supplementary Figure 1: Various plots describing the contents of the CIF file dataset. **a** The distribution of CIF files containing the atoms indicated (by atomic number) on the x-axis. The most abundant element in the dataset is oxygen, followed by copper, then lithium. **b** The distribution of Z values (i.e. the number of formula units in the unit cell) in the dataset. The majority of structures have Z of 1-4. **c** The distribution of compositions by the number of constituent elements in the formula. Most formulas are ternary or quaternary. Inset: The Venn diagram illustrates the numbers of unique reduced compositions obtained from each of the publicly accessible materials databases used to create the training dataset. **d** The distribution of space groups occurring in the CIF files of the dataset.



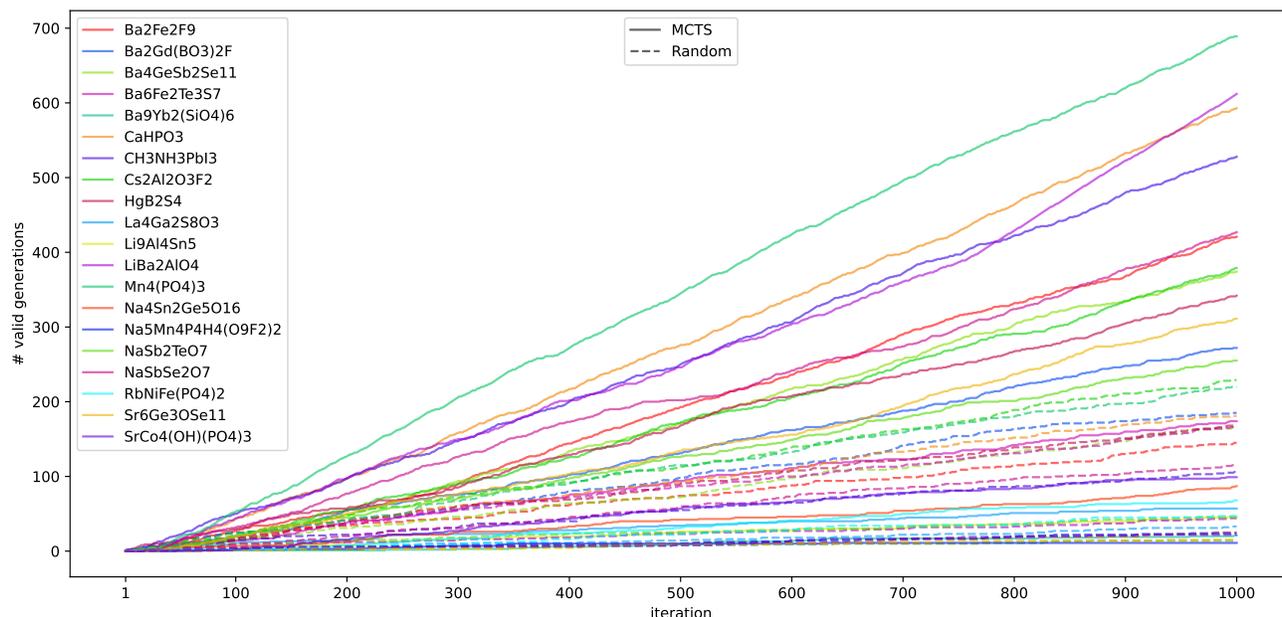
Supplementary Figure 2: A plot of the training set and validation set losses for the small model over the course of training. The green points represent an improved loss on the validation set, demonstrating that the model continues to improve its performance on the validation set even after 90,000 iterations. While the gap between the training set loss and the validation set loss appears to grow over the course of training, this is not necessarily indicative of overfitting. The growing gap could be more indicative of the differences between the distributions of the training and validation sets. On absolute terms, the difference between the curves is less than 0.02 units. Moreover, the performance on the validation and challenge sets indicate that the model trained for more iterations is superior.



Supplementary Figure 3: Plots depicting the small model's learned atom vectors. **a** A t-SNE [68] plot of the small model's dimensionally reduced learned atom vectors. **b** A heatmap of the cosine similarities between the small model's learned atom vectors.



Supplementary Figure 5: Plots depicting the small model's learned numeric digit vectors. a A plot of the small model's learned numeric digit vectors, dimensionally reduced using PCA. b A plot of the cosine similarities between the small model's learned numeric digit vectors. c A plot of the Euclidean distances between the small model's learned numeric digit vectors.



Supplementary Figure 6: Plots of the number of valid generations over the course of 1,000 iterations for the MCTS experiment (with no space group). The plot illustrates the finding that MCTS produces more valid generations than sampling randomly, and that, in some cases, the validation rate increases over time.

Supplementary References

- [1] Karpathy, A. nanoGPT: The simplest, fastest repository for training/finetuning medium-sized GPTs. <https://github.com/karpathy/nanoGPT> (2023).
- [2] Radford, A., Narasimhan, K., Salimans, T., Sutskever, I. *et al.* Improving Language Understanding by Generative Pre-Training (2018).
- [3] Hendrycks, D. & Gimpel, K. Gaussian Error Linear Units (GELUs). *arXiv preprint arXiv:1606.08415* (2016).
- [4] Press, O. & Wolf, L. Using the Output Embedding to Improve Language Models. *arXiv preprint arXiv:1608.05859* (2016).
- [5] Loshchilov, I. & Hutter, F. Decoupled Weight Decay Regularization. *arXiv preprint arXiv:1711.05101* (2017).
- [6] Togo, A. & Tanaka, I. Spglib: a software library for crystal symmetry search. *arXiv preprint arXiv:1808.01590* (2018).
- [7] Coulom, R. Efficient Selectivity and Backup Operators in Monte-Carlo Tree Search. In *International Conference on Computers and Games*, 72–83 (Springer, 2006).

- [8] Browne, C. B. *et al.* A Survey of Monte Carlo Tree Search Methods. *IEEE Transactions on Computational Intelligence and AI in Games* **4**, 1–43 (2012).
- [9] Rosin, C. D. Multi-armed Bandits with Episode Context. *Annals of Mathematics and Artificial Intelligence* **61**, 203–230 (2011).
- [10] Silver, D. *et al.* Mastering the game of Go with deep neural networks and tree search. *Nature* **529**, 484–489 (2016).
- [11] Choudhary, K. & DeCost, B. Atomistic Line Graph Neural Network for improved materials property predictions. *npj Computational Materials* **7**, 185 (2021).
- [12] Plumhoff, A. *Thermodynamic properties, crystal structures, phase relations and isotopic studies of selected copper oxyalts*. Ph.D. thesis (2020). Friedrich-Schiller-Universität Jena.
- [13] Pogue, E. A. *et al.* A gold (I) oxide double perovskite: Ba₂AuIO₆. *Journal of the American Chemical Society* **143**, 19033–19042 (2021).
- [14] Huang, Q. *et al.* Investigation of Charge-Ordered Barium Iron Fluorides with One-Dimensional Structural Diversity and Complex Magnetic Interactions. *Inorganic Chemistry* **62**, 14044–14054 (2023).
- [15] Frøen, E. H., Adler, P. & Valldor, M. Synthesis and Properties of Ba₆Fe₂Te₃S₇, with an Fe Dimer in a Magnetic Singlet State. *Inorganic Chemistry* **62**, 12548–12556 (2023).
- [16] Chen, Y. *et al.* Improvement on Magnetocaloric Effect through Structural Evolution in Gadolinium Borate Halides Ba₂Gd(BO₃)₂X (X= F, Cl). *Inorganic Chemistry* **62**, 15584–15592 (2023).
- [17] Yuan, F.-Y. *et al.* Ba₄GeSb₂Se₁₁: an infrared nonlinear optical crystal with a V-shaped Se₃₂-group possessing a large contribution to the SHG response. *Inorganic Chemistry* **60**, 15593–15598 (2021).
- [18] Yadav, S., Panigrahi, G., Niranjana, M. K. & Prakash, J. Ba₃GeTeS₄: A new quaternary heteroanionic chalcogenide semiconductor. *Journal of Solid State Chemistry* **323**, 124028 (2023).
- [19] Keerthisinghe, N., Ayer, G. B., Smith, M. D. & Zur Loye, H.-C. Comparative Study on Crystal Structures and Synthetic Techniques of Ternary Hafnium/Zirconium Fluorides. *Inorganic Chemistry* **62**, 12089–12098 (2023).
- [20] Motozawa, S., Kimura, H., Takahashi, J., Simura, R. & Yamane, H. BaY₁₆Si₄O₃₃ containing Ba(SiO₄)₄ orthosilicates. *Acta Crystallographica Section E: Crystallographic Communications* **78**, 1249–1252 (2022).
- [21] Liu, A. *et al.* Ba₉RE₂(SiO₄)₆ (RE= Ho–Yb): A Family of Rare-Earth-Based Honeycomb-Lattice Magnets. *Inorganic Chemistry* **62**, 13867–13876 (2023).
- [22] Saiduzzaman, M. *et al.* Hydrothermal synthesis of pyrochlore-type pentavalent bismuthates Ca₂Bi₂O₇ and Sr₂Bi₂O₇. *Inorganic Chemistry* **58**, 1759–1763 (2019).

- [23] Braun, T. & Hlukhyy, V. Structural order-disorder in CaFe_6Ge_6 and $\text{Ca}_{1-x}\text{Co}_x\text{Ge}_6$. *Journal of Solid State Chemistry* **318**, 123742 (2023).
- [24] Phillips, M. L. & Harrison, W. T. Synthesis and crystal structure of calcium hydrogen phosphite, CaHPO_3 . *Acta Crystallographica Section E: Crystallographic Communications* **75**, 997–1000 (2019).
- [25] Makhaneva, A. Y., Zakharova, E. Y., Nesterenko, S. N., Lyssenko, K. A. & Kuznetsov, A. N. CaPt_4P_6 , first calcium-containing representative of the ternary pyrite-derived pnictides of the BaPt_4As_6 type: Synthesis, crystal, and electronic structure. *Journal of Solid State Chemistry* **322**, 123969 (2023).
- [26] Weil, M. $\text{Ca}_2\text{Te}_3\text{O}_8$, a new phase in the $\text{CaO}\text{--}\text{TeO}_2$ system. *Acta Crystallographica Section E: Crystallographic Communications* **75**, 26–29 (2019).
- [27] Fukuda, M. *et al.* A-Site Columnar-Ordered Perovskite CaZnV_2O_6 as a Pauli-Paramagnetic Metal. *Inorganic Chemistry* (2023).
- [28] Desroches, G. & Bobev, S. Synthesis and structure determination of $\text{Ce}_6\text{Cd}_{23}\text{Te}$: a new chalcogen-containing member of the $\text{RE}_6\text{Cd}_{23}\text{T}$ family (RE is a rare-earth metal and T is a late group 14, 15 and 16 element). *Acta Crystallographica Section C: Structural Chemistry* **73**, 121–125 (2017).
- [29] Šimko, F. *et al.* Cesium Oxo-fluoro-aluminates in the $\text{CsF}\text{--}\text{Al}_2\text{O}_3$ System: Synthesis and Structural Characterization. *Inorganic Chemistry* **62**, 15651–15663 (2023).
- [30] Morrison, G., Jones, V. G., Zamorano, K. P., Greedan, J. E. & Zur Loye, H.-C. Flux Synthesis, UV–vis Absorbance, and Magnetism of Cesium Copper Silicates with an Isolated Super-Super Exchange Spin Dimer in $\text{Cs}_6\text{Cu}_2\text{Si}_9\text{O}_{23}$. *Inorganic Chemistry* **62**, 11682–11689 (2023).
- [31] Kimura, H. & Yamane, H. Crystal structure of chain silicate $\text{Cs}_3\text{LuSi}_3\text{O}_9$. *Acta Crystallographica Section E: Crystallographic Communications* **77**, 1239–1242 (2021).
- [32] Craig, A. J., Stoyko, S. S., Bonnoni, A. & Aitken, J. A. Syntheses and crystal structures of the quaternary thiogermanates $\text{Cu}_4\text{FeGe}_2\text{S}_7$ and $\text{Cu}_4\text{CoGe}_2\text{S}_7$. *Acta Crystallographica Section E: Crystallographic Communications* **76**, 1117–1121 (2020).
- [33] Zhang, N. *et al.* $\text{Eu}_2\text{MGe}_2\text{OS}_6$ (M= Mn, Fe, Co): Three Melilite-Type Rare-Earth Oxythiogermanates Exhibiting Balanced Nonlinear-Optical Behaviors. *Inorganic Chemistry* **62**, 16299–16303 (2023).
- [34] Huang, Y. *et al.* HgB_2S_4 : A d^{10} Metal Thioborate with Giant Birefringence and Wide Band Gap. *Chemistry of Materials* (2023).
- [35] Ramakrishnan, S. *et al.* Coupling between Charge Density Wave Ordering and Magnetism in $\text{Ho}_2\text{Ir}_3\text{Si}_5$. *Chemistry of Materials* **35**, 1980–1990 (2023).
- [36] Redhammer, G. J. & Tippelt, G. The crystal structure of KScP_2O_7 . *Acta Crystallographica Section E: Crystallographic Communications* **76**, 1412–1416 (2020).

- [37] Dong, W. *et al.* $\text{K}_2\text{Sr}_4(\text{PO}_3)_{10}$: A Polyphosphate with Deep-UV Cutoff Edge and Enlarged Birefringence. *Inorganic Chemistry* **62**, 16215–16221 (2023).
- [38] Eder, F. & Weil, M. Crystal structure of $\text{K}_6[\text{Zn}(\text{CO}_3)_4]$. *Acta Crystallographica Section E: Crystallographic Communications* **79**, 718–721 (2023).
- [39] Yan, H. *et al.* $\text{La}_4\text{Ga}_2\text{S}_8\text{O}_3$: A Rare-Earth Gallium Oxysulfide with Disulfide Ions. *Inorganic Chemistry* (2023).
- [40] Zhang, H. *et al.* Prediction and Synthesis of a Selenide Perovskite for Optoelectronics. *Chemistry of Materials* (2023).
- [41] Pavlyuk, V., Dmytriv, G., Tarasiuk, I. & Ehrenberg, H. $\text{Li}_9\text{Al}_4\text{Sn}_5$ as a new ordered superstructure of the $\text{Li}_{13}\text{Sn}_5$ type. *Acta Crystallographica Section C: Structural Chemistry* **73**, 337–342 (2017).
- [42] Nishita, Y., Simura, R., Inaguma, Y. & Yamane, H. $\text{LiBa}_2\text{AlO}_4$: A new lithium barium aluminate having an oxygen tetrahedral framework. *Journal of Solid State Chemistry* **317**, 123654 (2023).
- [43] Roh, J., Do, N., Manjón-Sanz, A. & Hong, S.-T. Li_2GeS_3 : Lithium Ionic Conductor with an Unprecedented Structural Type. *Inorganic Chemistry* **62**, 15856–15863 (2023).
- [44] Gvozdetzkyi, V. *et al.* From Layered Antiferromagnet to 3D Ferromagnet: LiMnBi -to- MnBi Magneto-Structural Transformation. *Chemistry of Materials* **35**, 3236–3248 (2023).
- [45] Hyde, P. *et al.* Lithium Intercalation into the Excitonic Insulator Candidate Ta_2NiSe_5 . *Inorganic Chemistry* **62**, 12027–12037 (2023).
- [46] Ponou, S., Lidin, S. & Mudring, A.-V. Optimization of Chemical Bonding through Defect Formation and Ordering—The Case of $\text{Mg}_7\text{Pt}_4\text{Ge}_4$. *Inorganic Chemistry* (2023).
- [47] Kamutzki, F. *et al.* $\text{NaGdSi}_2\text{O}_6$ – A novel antiferromagnetically coupled silicate with Vierer chain structure. *Journal of Solid State Chemistry* **317**, 123677 (2023).
- [48] Nagai, T. & Kimura, T. Chemical Switching of Ferroaxial and Nonferroaxial Structures Based on Second-Order Jahn–Teller Activity in $(\text{Na K})_2\text{Hf}(\text{BO}_3)_2$. *Chemistry of Materials* **35**, 4109–4115 (2023).
- [49] Galven, C. *et al.* $\text{Na}_6\text{Li}_4\text{MO}_4(\text{CO}_3)_4$ (M= W and Mo): An Alternative Electrolyte for High-Temperature Electrochemical Cells. *Inorganic Chemistry* **62**, 15367–15374 (2023).
- [50] Hughes, J. M., Wise, W. S., Gunter, M. E., Morton, J. P. & Rakovan, J. Lasalite, $\text{Na}_2\text{Mg}_2[\text{V}_{10}\text{O}_{28}] \cdot 20\text{H}_2\text{O}$, a new decavanadate mineral species from the Vanadium Queen Mine, La Sal District, Utah: Description, atomic arrangement, and relationship to the pascoite group of minerals. *The Canadian Mineralogist* **46**, 1365–1372 (2008).

- [51] Luo, Q., Li, N., Zhao, Z., Cui, M. & He, Z. A new compound $\text{Na}_5\text{Mn}_4(\text{PO}_4)_4\text{F}_4 \cdot 2\text{H}_2\text{O}$ with a rarely mixed valence spin chain showing multiple magnetic transitions. *Inorganic Chemistry Frontiers* **10**, 6303–6307 (2023).
- [52] Robert, R., Mangalassery, S., Rao, D. N. & Vidyasagar, K. Syntheses and characterization of quaternary selenites and tellurite of antimony, $\text{NaSbSe}_2\text{O}_7$, $\text{AgSbSe}_2\text{O}_7$ and $\text{Na}_2\text{Sb}_4\text{Te}_2\text{O}_{14}$. *Journal of Solid State Chemistry* **327**, 124228 (2023).
- [53] Novikov, S. *et al.* $\text{Na}_{4-x}\text{Sn}_{2-x}\text{Sb}_x\text{Ge}_5\text{O}_{16}$, an Air-Stable Solid-State Na-Ion Conductor. *Inorganic Chemistry* **62**, 16068–16076 (2023).
- [54] Eder, F. & Weil, M. Garnet-type $\text{Na}_3\text{Te}_2(\text{FeO}_4)_3$. *Acta Crystallographica Section E: Crystallographic Communications* **79** (2023).
- [55] Chong, S., Kroll, J. O., Crum, J. V. & Riley, B. J. Synthesis and crystal structure of a neodymium borosilicate, $\text{Nd}_3\text{BSi}_2\text{O}_{10}$. *Acta Crystallographica Section E: Crystallographic Communications* **75**, 700–702 (2019).
- [56] Eder, F. & Weil, M. $\text{Ni}_3\text{Te}_2\text{O}_2(\text{PO}_4)_2(\text{OH})_4$, an open-framework structure isotypic with $\text{Co}_3\text{Te}_2\text{O}_2(\text{PO}_4)_2(\text{OH})_4$. *Acta Crystallographica Section E: Crystallographic Communications* **76**, 625–628 (2020).
- [57] Badri, A., Bembli, M., Alvarez-Serrano, I., López, M. L. & Amara, M. B. Synthesis, single crystal structure, optical and magnetic properties of a new rubidium nickel iron phosphate $\text{RbNiFe}(\text{PO}_4)_2$. *Journal of Solid State Chemistry* 124141 (2023).
- [58] Huang, D. *et al.* Metal Halide Single Crystals RbCdCl_3 : Sn^{2+} and Rb_3SnCl_7 with Blue and White Emission Obtained via a Hydrothermal Process. *Inorganic Chemistry* **62**, 15943–15951 (2023).
- [59] Cherif, F.-Z. *et al.* Crystal structure of $\text{SrCo}_4(\text{OH})(\text{PO}_4)_3$, a new hydroxyphosphate. *Acta Crystallographica Section E: Crystallographic Communications* **76**, 1022–1026 (2020).
- [60] Hyung, J., Lee, H. W., Kim, S. J., Shin, H. R. & Hong, S.-T. Crystal structure of strontium perchlorate anhydrate, $\text{Sr}(\text{ClO}_4)_2$, from laboratory powder X-ray diffraction data. *Acta Crystallographica Section E: Crystallographic Communications* **75**, 447–450 (2019).
- [61] Menezes, L. T. *et al.* $\text{Sr}_6\text{Ge}_3\text{OSe}_{11}$: A Rationally Designed Noncentrosymmetric Oxyselenide with Polar $[\text{GeOSe}_3]$ Building Blocks. *Chemistry of Materials* **35**, 3033–3040 (2023).
- [62] Xie, Y. *et al.* A series of new rare-earth sulfide borates $\text{RE}_3\text{S}_3\text{BO}_3$ (RE= Nd, Tb, Dy): Syntheses, structures and optical properties. *Journal of Solid State Chemistry* **327**, 124277 (2023).
- [63] Zhou, C. & Li, R. Large Difference in Nonlinear Optical Activity of Rare Earth Ion Substitution of Bi^{3+} in A_3TeBO_9 (A= Bi, La, Pr, Nd, Sm-Dy). *Inorganic Chemistry* **62**, 11265–11270 (2023).

- [64] Mazet, T., Welter, R. & Malaman, B. A study of the new ferromagnetic YbMn_6Sn_6 compound by magnetization and neutron diffraction measurements. *Journal of Magnetism and Magnetic Materials* **204**, 11–19 (1999).
- [65] Eder, F. & Weil, M. Crystal structure of $\text{Zn}_2(\text{HTeO}_3)(\text{AsO}_4)$. *Acta Crystallographica Section E: Crystallographic Communications* **77**, 555–558 (2021).
- [66] Hu, C.-L., Han, Y.-X., Fang, Z. & Mao, J.-G. $\text{Zn}_2\text{BS}_3\text{Br}$: An Infrared Nonlinear Optical Material with Significant Dual-Property Enhancements Designed through a Template Grafting Strategy. *Chemistry of Materials* **35**, 2647–2654 (2023).
- [67] Santamaría-Pérez, D., Chuliá-Jordán, R., Otero-de-la Roza, A., Oliva, R. & Popescu, C. High-Pressure Experimental and DFT Structural Studies of Aurichalcite Mineral. *Minerals* **13**, 619 (2023).
- [68] Van der Maaten, L. & Hinton, G. Visualizing Data using t-SNE. *Journal of Machine Learning Research* **9** (2008).