

# Exploring the polymorphism of dicalcium silicates using transfer learning enhanced machine learning atomic potentials

1 Jon López-Zorrilla,\* Xabier M. Aretxabaleta, and Hegoi Manzano\*

*Physics department, University of the Basque Country UPV/EHU, 48940 Leioa, Basque Country, Spain*

E-mail: jon.lopezz@ehu.eus; hegoi.manzano@ehu.eus

## 2 Abstract

3 Belitic cements are a greener alternative to Ordinary Portland Cements due to the  
4 lower CO<sub>2</sub> associated to their production. However, their low reactivity with water  
5 is currently a drawback, resulting in longer setting times. In this study, we utilize a  
6 combination of evolutionary algorithms and machine learning atomic potentials (MLPs)  
7 to identify previously unreported belite polymorphs that may exhibit higher hydraulic  
8 reactivity than the known phases. To address the high computational demand of this  
9 methodology, we propose a novel transfer learning approach for generating MLPs. First,  
10 the models are pre-trained on a large set of classical data (ReaxFF) and then re-  
11 trained with Density Functional Theory (DFT) data. We demonstrate that the transfer  
12 learning enhanced potentials exhibit higher accuracy, require less training data, and are  
13 more transferable than those trained exclusively on DFT data. The generated machine  
14 learning potential enables a fast, exhaustive, and reliable exploration of the dicalcium  
15 silicate polymorphs. This includes studying their stability through phonon analysis and  
16 calculating their structural and elastic properties. Overall, we identify ten new belite

17 polymorphs within the energy range of the existing ones, including a layered phase with  
18 potentially high reactivity.

## 19 **1 Introduction**

20 Cement is the most extensively manufactured product globally in terms of mass. In 2022,  
21 its global production reached an astounding 4.2 billion tons, well over 500 kg per capita.<sup>1</sup>  
22 Cement and concrete, characterized by their versatility, cost-effectiveness, abundance, and  
23 local availability, are indispensable components in modern societies, playing a fundamental  
24 role in various construction applications. However, despite their ubiquity, the prevalent  
25 technologies associated with their production contribute significantly to carbon dioxide (CO<sub>2</sub>)  
26 emissions, accounting for 8% of total per capita emissions. Most of the emissions, up to 60%,  
27 are intrinsic to the material.<sup>2</sup> The raw minerals, mainly clays and limestone, are melted at  
28 high temperatures to form the clinker phases: calcium silicate and aluminate phases, which  
29 mixed with other components such as gypsum and additives, form the cement. During  
30 melting, the calcination of limestone (CaCO<sub>3</sub>) releases a considerable amount of CO<sub>2</sub>, which  
31 is unavoidable.<sup>3</sup> Therefore, the strategies for cement's environmental impact reduction are  
32 based on Carbon Capture and Utilization or a reduction of the clinker in cement through  
33 the use of Supplementary Cementitious Materials. But there is a third alternative that  
34 researchers have long pursued: the so-called *belite cements*.<sup>4</sup>

35 Belite cements (BCs) are, as the name indicates, cements in which the main component  
36 is dicalcium silicate (also known as belite or C<sub>2</sub>S in cement chemistry notation), in contrast  
37 to Ordinary Portland Cement (OPC) in which tricalcium silicate or alite predominates. Due  
38 to its lower Ca content, BCs require less limestone, reducing by up to 1/3 the OPC CO<sub>2</sub>  
39 emissions upon calcination. Furthermore, the sintering temperature is lower, also reducing  
40 fuel consumption. The resulting cement paste after BC hydration is equal to or even out-  
41 performs OPC pastes in terms of durability and mechanical properties.<sup>5,6</sup> However, BCs are

42 unpractical for most industrial applications due to the low dissolution rate of belite. There-  
43 fore, a practical transition from OPC to BC requires an acceleration of belite dissolution,  
44 often called activation.

45 Dicalcium silicate,  $\text{Ca}_2\text{SiO}_4$ , can be found in nature as an orthosilicate called larnite, the  
46 Ca end member of the olivine mineral group. In cement chemistry, this stable polymorph  
47 is denoted as  $\gamma$ , and it is not desirable due to its low hydraulic activity. Four additional  
48 polymorphs are found during clinker production, named  $\beta$ ,  $\alpha$ ,  $\alpha'_L$ , and  $\alpha'_H$ . The  $\beta$  form is  
49 predominant in cement, stabilized by the presence of guest ions in the structure, mainly Mg.<sup>7</sup>  
50 A partial activation of belite has been achieved by chemical and mechanical means, as well  
51 as by the use of additives.<sup>6,8</sup> Besides these strategies, a new polymorph denoted as  $X$  has  
52 been recently obtained by thermal decomposition of a hydrated calcium silicate  $\alpha\text{-C}_2\text{SH}$ .<sup>9</sup>  
53 The  $X$  polymorph is obtained together with a considerable amount of amorphous phase, and  
54 the mixture is reported to hydrate faster than the conventional  $\beta$  and  $\alpha$  forms.<sup>10</sup> Large-scale  
55 production of amorphous and  $X$ -belite is currently impractical, although laboratory-scale  
56 synthesis is feasible.

57 The discovery of the  $X$  polymorph motivates the current work: could other metastable  
58 and highly reactive polymorphs of belite exist? To answer this question, we have used  
59 Evolutionary Algorithms (EA) to explore the configurational space of dicalcium silicate and  
60 search for unreported belite polymorphs. Performing an exhaustive search using EA requires  
61 thousands of Density Functional Theory (DFT) simulations, which can be prohibitive due  
62 to their high computational cost. The recent outburst of machine learning atomic poten-  
63 tials<sup>11,12</sup> (MLPs) provides a new alternative, enabling simulations with DFT precision at a  
64 significantly reduced computational cost. However, training a MLP demands, in turn, a sub-  
65 stantial number of DFT calculations to build the database. To break the deadlock, we have  
66 used the transfer learning (TL) methodology,<sup>13,14</sup> which involves pre-training the models on  
67 low-quality data before training on the DFT data to enhance the performance of machine  
68 learning potentials. In practice, the re-training on the smaller set of high-quality data can be

69 accomplished by either fine-tuning all the parameters or keeping some of the layers frozen.

70 Transfer learning is a valuable tool in computational materials science for predicting var-  
71 ious properties, where models pre-trained on different levels of computational data are used  
72 to improve performance when only a few data points are available.<sup>13–15</sup> In particular, in the  
73 field of machine learning potentials, a common approach involves fine-tuning Density Func-  
74 tional Theory models to achieve post-Hartree-Fock accuracy,<sup>16,17</sup> mainly by training on the  
75 difference between the two methods (commonly known as  $\Delta$ -learning).<sup>18</sup> Although transfer  
76 learning models are typically used to achieve coupled cluster accuracy from DFT,<sup>19–21</sup> some  
77 research has been conducted to reduce the amount of DFT data required by pre-training on  
78 more primitive DFT approximations.<sup>22</sup> Nevertheless, to the best of our knowledge, none of  
79 the published works demonstrate the feasibility of transferring the physical knowledge from  
80 classical potentials to *ab initio* quantum methods. Even low-quality data based on empirical  
81 potentials contains significant, even if not very accurate, physical information about the sys-  
82 tems at the atomic scale that can be used to enhance the MLP while minimizing the required  
83 amount of data. Thus, our approach involves exploiting the speed of empirical potentials  
84 to thoroughly sample the phase space and pre-train the machine learning models. We then  
85 select a small subset of those configurations to include in the DFT training database. In this  
86 work, we choose to pre-train the models using the ReaxFF reactive force field,<sup>23,24</sup> which is  
87 itself fitted to reproduce *ab initio* calculations. As for the high-level method, we consider  
88 DFT under the PBE exchange-correlation functional<sup>25</sup> sufficient for our purpose.

## 89 2 Methods

### 90 Reference data generation

91 For each phase in the data set, the same sampling technique was followed, consisting of  
92 different cell deformations and molecular dynamics simulations. First, several MD runs  
93 were performed using the ReaxFF forcefield<sup>24</sup> with the Ca/Si/O/H set of parameters from

94 Refs.<sup>26,27</sup> in LAMMPS<sup>28</sup> under the NVT ensemble. Various simulations were performed at  
95 different combinations of temperatures ( $T = 300K, 600K, 900K$ ) and cell volumes ( $\Delta V/V =$   
96  $0.9, 1, 1.1$ ). A time step of 0.2 fs was used, and snapshots of the trajectory were saved every  
97 500 steps. Second, all the non-symmetric axes were deformed from 10% compression to 10%  
98 expansion, including hydrostatic deformations with a maximum of 10% variation in volume,  
99 and angles were varied from  $-10^\circ$  to  $10^\circ$ . The DFT data set was generated by randomly  
100 selecting structures from the ReaxFF data set, and evaluating their energy and forces.

## 101 **Density Functional Theory**

102 DFT calculations were performed using the quantum ESPRESSO software<sup>29,30</sup> using ON-  
103 CVSP pseudopotentials<sup>31</sup> from pseudodojo,<sup>32</sup> under the PBE exchange-correlation func-  
104 tional,<sup>25</sup> and with the empirical dispersion by Grimme.<sup>33</sup> The plane wave energy cutoff was  
105 set to 100 Ry, and calculations were converged to  $10^{-6}$  eV. Geometry optimizations were  
106 converged to  $10^{-5}$  eV and  $10^{-4}$  eV $\text{\AA}^{-1}$  for energy and forces, respectively. Taking into ac-  
107 count that systems of very different sizes have been studied, the number of k points was  
108 systematically selected such that the distance between points in the reciprocal space was  
109 about  $0.25\text{\AA}^{-1}$ .

## 110 **MLP architecture and training**

111 The machine learning atomic potentials used in this work are based on artificial neural  
112 networks.<sup>11,12</sup> All of them were trained using the  $\text{\ae}net$ -PyTorch software,<sup>34,35</sup> using all the  
113 energies and 50% of the atomic forces. Chebyshev polynomials were used as descriptors for  
114 the atomic environments,<sup>36</sup> with a  $N_{\text{rad}} = 18$  and  $N_{\text{ang}} = 6$  order expansion for the radial and  
115 angular basis respectively. The radial cutoff distance was  $6.5\text{\AA}$ , while the angular distance  
116 was  $4\text{\AA}$ . This leads to a fingerprint with 52 components for each atomic environment. The  
117 MLP architecture for all models was  $52 - 10 - 10 - 1$ , with hyperbolic tangent as activation  
118 function. The only exception is the initial toy model of the calcium ion and the silicon

119 dioxide molecule, where the architecture was reduced to 40 – 3 – 3 – 1. Transfer learning is  
120 performed by fine-tuning all the parameters of the pre-trained models.

121 The Supplementary Information contains a detailed analysis of the transferring method-  
122 ology by freezing all combinations of the layers, and a comparison between different network  
123 architecture and descriptor sizes, showing that fine-tuning all layers is the optimal choice in  
124 our case.

125 Note that for every training data set of each experiment presented throughout this work,  
126 several MLPs have been trained, and the results shown correspond to the average of all  
127 MLPs.

## 128 **Evolutionary algorithms**

129 The exploration of the dicalcium silicate polymorphs was done using evolutionary algorithms  
130 as implemented in the USPEX code<sup>37–39</sup> (version 10.5). For each system size, EA runs were  
131 performed for enough generations until all experimentally known phases were found. Each  
132 structure was relaxed using the `ænet-LAMMPS` interface,<sup>40,41</sup> first minimizing the energy  
133 using ReaxFF to avoid random structures far from the included in the training data, and  
134 then using MLPs.

## 135 **Phonon and elastic properties**

136 Phonons were computed under the finite difference approximation to build the dynamical  
137 matrix. The `phonopy`<sup>42,43</sup> software was employed to generate the appropriate atomic dis-  
138 placements for each crystal structure, to build the dynamical matrix, and to compute the  
139 force constants and phonon dispersion along the high-symmetry path of the corresponding  
140 space group. The atomic displacements were set to 0.1Å, and supercells of at least 13Å  
141 (twice the cutoff distance of the descriptors) along each crystallographic axis were used in  
142 order to guarantee convergence.

143 The elastic tensor of all the structures was computed fitting the stress-strain relationship

144  $\sigma_i = C_{ij}\epsilon_j$ . Each crystal parameter (cell-vector lengths and angles) was deformed indepen-  
145 dently 10 times, with a deformation in a range  $\epsilon_0 \in (-0.01, 0.01)$ . The elastic properties are  
146 computed under the Hill scheme. The forces and stress tensors for each of the structures  
147 were evaluated using the LAMMPS interface of `ænet`.

## 148 **Annealing and amorphous dicalcium silicate**

149 All MLP-based molecular dynamics simulations were performed with the `ænet`-LAMMPS  
150 interface.<sup>40,41</sup> The annealing to refine the polymorphs was performed under the NPT ensem-  
151 ble with a time step of 0.5 fs, both heating from 0 K to 400 K and cooling back at a rate of  
152  $2 \cdot 10^{12} \text{ Ks}^{-1}$ .

153 The amorphous dicalcium silicate models were generated by heating a  $4 \times 2 \times 4$  supercell  
154 of  $\gamma\text{-C}_2\text{S}$  up to 2000 K. The three different amorphous models were obtained by cooling the  
155 heated structure at three rates:  $2 \cdot 10^{12}$ ,  $2 \cdot 10^{13}$ , and  $2 \cdot 10^{14} \text{ Ks}^{-1}$ . In all three cases, the  
156 time step was lowered to 0.1 fs to ensure the stability of the high-temperature molecular  
157 dynamics.

## 158 **3 Results and discussion**

159 Our results are organized as follows: first, the advantages of transfer learning are qualitatively  
160 introduced with a simple toy model. Second, we consider a more complex and realistic  
161 dataset to train a MLP for dicalcium silicates, while quantifying the benefits of the transfer  
162 learning approach. Finally, we use the trained MLP to explore the polymorphism of dicalcium  
163 silicates.

### 164 **A simple transfer learning model**

165 Let us first consider a simplified scenario to illustrate the capabilities of the transfer learning  
166 methodology: a system formed by a calcium ion and a silicon dioxide molecule, with the Si

167 and Ca atoms fixed at a distance of  $5\text{\AA}$ . We will explore the potential energy surface (PES)  
 168 of the system as a function of the distance  $d_{\text{Si-O}}$  from the silicon atom to a mobile oxygen,  
 169 which breaks its bond to move toward the calcium atom. Sampling that PES using ReaxFF  
 170 and DFT yields two similar landscapes, with two possible bound states for the oxygen and an  
 171 energy barrier for the oxygen transfer. According to DFT, being bonded to the calcium ion  
 172 is the lowest-energy configuration, while ReaxFF predicts the Si-O bond to be more stable,  
 173 as shown in Figure 1 (a).

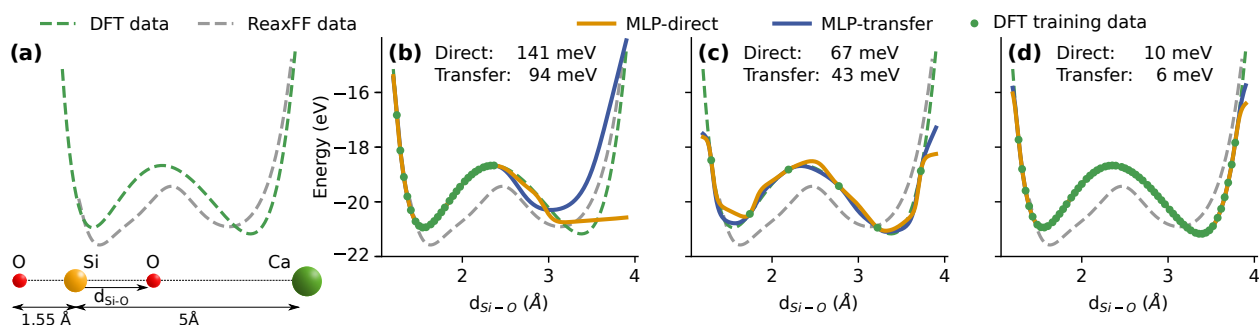


Figure 1: **Simple transfer learning model.** (a) Potential energy surface of the calcium ion and the silicon dioxide molecule computed with ReaxFF (grey dashed lines) and DFT (green dashed lines) as a function of the Si-O distance. The rest of the figure shows machine-learned potential energy surface for both direct learning (orange lines) and transfer learning (blue lines), for different distributions of DFT training data (green dots). The mean absolute error of the energy of both transfer and direct learning is displayed in each case.

174 We evaluate the advantage of the TL strategy in three scenarios with different distribu-  
 175 tions of the DFT training data. Two independent models are trained for each of them: one  
 176 trained on all the ReaxFF data and subsequently re-trained on the selected DFT data points  
 177 (MLP-transfer), and the second trained exclusively on that DFT data (MLP-direct).

178 First, we focus on a scenario where one of the bound states is correctly sampled by DFT  
 179 (i.e. the Si-O bound state) while no DFT training data about the second state is included,  
 180 see Figure 1(b). In this case, the direct training leads to an incorrect representation of the  
 181 Ca-O region, even predicting an unphysical PES. On the contrary, the transferred model  
 182 does predict a bound state resembling that of the ReaxFF data. Second, we explore the case  
 183 where the DFT training points cover both regions of the PES but they are sparse (with only



184 3 points·Å<sup>-1</sup>). In this case, both the direct and transfer learning protocols give a reasonable  
185 answer, but the transfer learning results are clearly smoother and the error with respect  
186 to the DFT data is considerably lower, see Figure 1(c). By increasing the DFT training  
187 configurations to 30 points·Å<sup>-1</sup> [Figure 1(d)], both direct and transfer models yield similar  
188 results, but the error on the validation set is still lower for TL.

189 Overall, this simple model illustrates the ability of TL to reach accurate predictive capa-  
190 bility with a reduced DFT dataset by employing empirical potentials to pre-train the MLP,  
191 and even predict energies for unsampled areas of the phase-space.

## 192 **Transfer Learning MLP for dicalcium silicates**

193 We now focus on the construction of a large dataset to train the MLP for dicalcium silicates,  
194 containing the 12 polymorphs available in Materials Project<sup>44</sup> as of December 2022, includ-  
195 ing experimental and theoretical phases. To pre-train the model, we sample a total of 10000  
196 configurations using ReaxFF, by performing molecular dynamics simulations for each poly-  
197 morph under different conditions and deforming the equilibrium cell along all independent  
198 crystallographic axes, as detailed in the Methods section. The DFT dataset is built from  
199 this data, by randomly selecting structures to be evaluated by DFT.

200 First, we study the impact of the transfer learning protocol on a realistic system like this.  
201 We consider several subsets of the database with an increasing amount of data and train  
202 models within both MLP-direct and MLP-transfer approaches. Figures 2(a) and (c) show  
203 the mean absolute error (MAE) of the energy and force as a function of the amount of DFT  
204 training data. Very interestingly, the transferred model outperforms its direct counterpart  
205 for any given amount of training data. The improvement (shown in the lower panels) is  
206 most significant with only a few hundred DFT data available for training, reaching up to  
207 a 40% reduction in both energy and force errors. With approximately 2000 DFT training  
208 data (half the total available set), the transfer learning model already reaches the same level  
209 of accuracy as the model trained directly on the full data set. Additionally, the model pre-

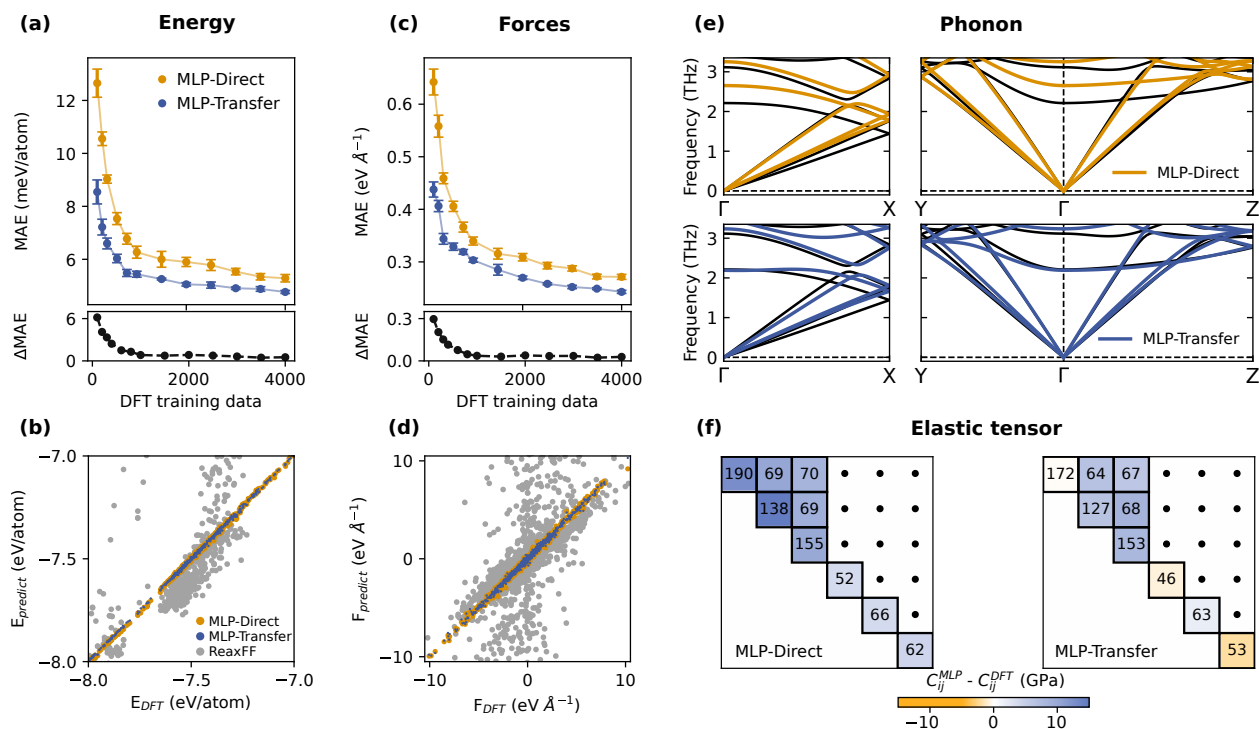


Figure 2: **Transfer learning results for dicalcium silicates.** (a) Energy error of transfer (blue) and direct learning (orange), as a function of the amount of DFT data used for training. The bottom panel shows the decrease of the error due to transfer learning. (b) Energy of a test set of 700 structures evaluated using transfer and direct learning and ReaxFF (grey), compared to the DFT value. (c)-(d) Same as the previous figure, in the case for the error of the forces. (e) Phonon dispersion of  $\gamma$ -belite computed using DFT (black), MLP-direct (orange) and MLP-transfer (blue). (f) Elastic tensor for the same polymorph. The color indicates the deviation with respect to the DFT value.

210 trained on ReaxFF data and trained on the full DFT data set is about 10% more accurate  
 211 than the best-performing MLP-direct model.

212 Let us further characterize the performance of the transferred and direct MLPs trained  
 213 on all the DFT data. Figures 2 (b) and (d) display the correlation of energy and forces  
 214 with the actual DFT calculations, for both machine learning approaches and for the ReaxFF  
 215 potential used for pre-training. Both MLP models outperform the ReaxFF potential and  
 216 demonstrate high accuracy, resulting in a nearly perfect correlation with the DFT data.  
 217 However, the MLP-transfer is still more accurate than the MLP-direct, by approximately  
 218 1 meV/atom in energies and  $0.05 \text{ eV}\text{\AA}^{-1}$  in forces, as quantified in Figure 2 (a) and (c).  
 219 This is a clear indication of the capabilities of the transfer learning enhanced MLPs and

220 their suitability for the exploration of the  $C_2S$  phase space at a similar computational cost  
221 to empirical potentials. A simple efficiency check with a supercell of  $4 \times 2 \times 4$  of  $\gamma$ -belite  
222 containing 896 atoms shows that the MLP is as fast as ReaxFF for 4 cores, but it scales  
223 better and is twice as fast for 32 cores, as shown in the Supplementary Information.

224 Finally, we explore the ability of our models to describe magnitudes related to higher-  
225 order derivatives of the PES, in particular, phonons [Figure 2(e)] and elastic properties  
226 [Figure 2(f)], which are of paramount importance for the discussion of our main results  
227 in the subsequent sections. As an illustrative example, we compute those magnitudes for  
228 the most stable belite polymorph,  $\gamma$ , within both training approaches, and compare the  
229 results to DFT calculations. The phonon dispersion curves are closer to DFT within the TL  
230 approach (blue lower panel) than directly training (orange upper panel). This is particularly  
231 noticeable for the optical modes, but it is also significant for acoustic modes. Moreover, the  
232 TL model excels at describing the phonon dispersion near the selected high symmetry points.  
233 Regarding the elastic tensor, direct MLPs are relatively accurate, within a 15% deviation  
234 from the DFT reference values. Additionally, the TL model further reduces the error in all  
235 the elastic constants, with a mean absolute error of 8.9 GPa on the elastic tensor, compared  
236 to the original 13.4 GPa of the direct model.

237 Thus, the TL approach improves the prediction of the PES and its first and second-order  
238 derivatives over directly training on all the available data. The Supplementary Information  
239 includes results from a similar study where only some phases are undersampled in the train-  
240 ing data, demonstrating that our transfer learning approach also improves the performance  
241 in such scenarios. This is also interesting for the exploration of the phase space using evolu-  
242 tionary algorithms, where many atomic arrangements not included in the database are likely  
243 to be encountered. Hence, pre-training the MLPs in a ReaxFF dataset as diverse as possible  
244 will enhance the predictive power of the resulting potentials.

245 Therefore, to further improve the flexibility of the MLP on those high-energy regions of  
246 the energy landscape, we incorporated several new structures into the previous training set:

247 several polymorphs of silicon dioxide ( $\text{SiO}_2$ ), calcium oxide ( $\text{CaO}$ ), calcium silicates (CS),  
 248 and tricalcium silicates ( $\text{C}_3\text{S}$ ), as well as 1000 data points taken from a preliminary DFT-  
 249 EA run to account for pseudorandom and high-energy conformations. Finally, to prevent  
 250 the system from collapsing if the interatomic distances are too small, we included dimers for  
 251 each pair of chemical elements in the system. The distance was reduced until repulsive forces  
 252 exceeded  $20 \text{ eV}\text{\AA}^{-1}$  and expanded up to  $4\text{\AA}$ . The final ReaxFF training database comprises  
 253 20000 structures, while the DFT database is a subset of 8000 data points, as detailed in the  
 254 S.I. Given the previous performance analysis, all the subsequent results are computed only  
 255 with the MLP generated by the TL approach.

## 256 Belite polymorph search and computational screening

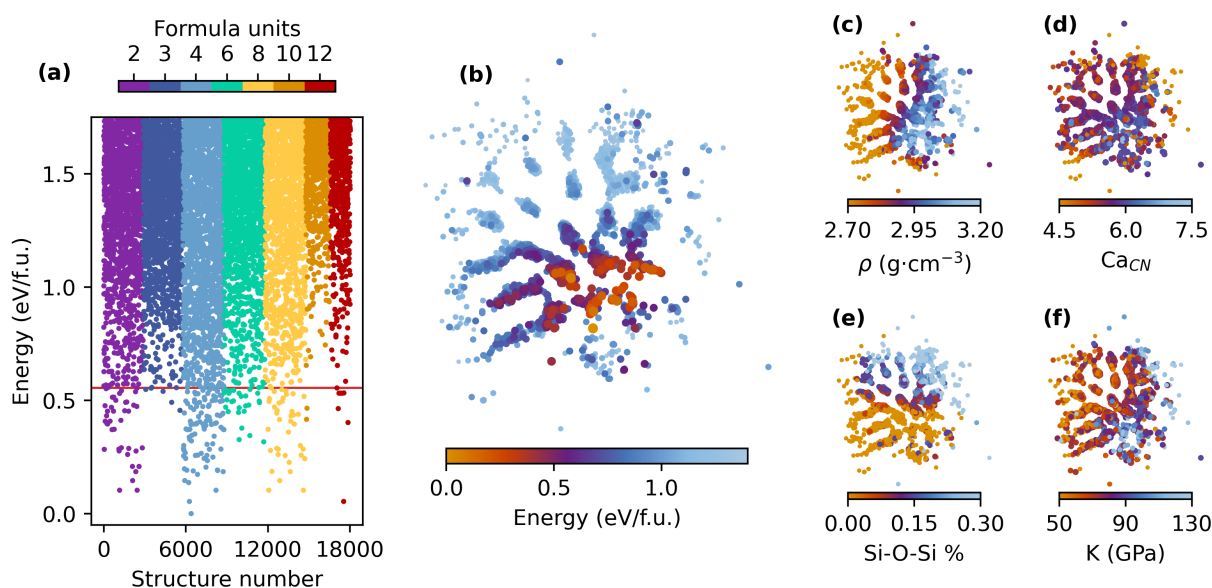


Figure 3: **Sketch-map of the generated structures.** (a) Energy of the found structures for all the considered numbers of formula units. (b) Sketch-map of the lowest-energy 5000 structures. Each point represents one of the structures, and the distance between them represents their structural similarity, i.e., the closer the points, the more similar the structures. (c) - (f) Several structural and mechanical properties represented using the same sketch-map: density ( $\rho$ ), average coordination number of calcium ( $\text{Ca}_{CN}$ ), percentage of oxygen atoms bonded to more than one silicon atom (Si-O-Si), and bulk modulus (K).

257 To discover new metastable polymorphs of belite, we conducted several independent EA

258 searches for different system sizes, ranging from 2 to 12 formula units per cell, at least until  
259 the experimental phases are found. This led to about 18000 potential structures, as displayed  
260 in Figure 3(a). Additional information regarding the number of structures for each size can  
261 be found in the Supplementary Information.

262 As shown in Figure 3(a), most of the generated structures have an energy considerably  
263 above the range of the experimental structures, delimited by the energy of the  $\alpha$  polymorph,  
264 indicated by a red line. Considering the large amount of structures generated, we focus on  
265 the 5000 with the lowest energy, which are more likely to contain metastable phases. For  
266 these phases, we calculated their structural dissimilarities as detailed in the Supplementary  
267 Material. We then reduced their dimensionality using a sketch-map,<sup>45,46</sup> leading to the  
268 representation displayed in Figures 3(b)-(f). Each point of the figure represents one structure,  
269 and the distance between points indicates the similarity between structures: the closer two  
270 points are, the more similar the structures are. Over the maps in Figures 3(b)-(f), we  
271 projected several structural properties and elastic properties obtained from the elastic tensor.  
272 The low-energy structures generally exhibit high density and high bulk modulus values. In  
273 particular, all the structures with a density above  $2.86 \text{ gcm}^{-3}$ , the lowest experimental density  
274 corresponding to  $X\text{-C}_2\text{S}$ , have energies below the higher energy experimental polymorph  $\alpha$ -  
275 belite. The calcium coordination number ( $\text{Ca}_{\text{CN}}$ ) distribution is centered at 6. Structures  
276 with  $\text{Ca}_{\text{CN}}$  below 5.5 are generally high-energy phases, while higher  $\text{Ca}_{\text{CN}}$  are more favorable.  
277 Finally, most of the low-energy polymorphs are orthosilicates, i.e. they have isolated silicate  
278 monomers, like the already known experimental phases.

279 The number of structures found in our initial search is too large, so we devise a compu-  
280 tational screening procedure to systematically filter the unique and most stable polymorphs.

- 281 • In the first step, we identified duplicate structures and superstructures, by examining  
282 the structures with matching energies, densities, and space groups. Furthermore, the  
283 dissimilarity analysis described in the previous section was used to discard phases with  
284 a structural distance lower than 0.05. At this step, we identified 3000 unique structures.

- 285 • A large number of structures were still considerably above the highest energy of the  
286 known experimental polymorphs. To narrow down the searching space, we discarded all  
287 the phases with cohesive energy 7.5% over the energy of the highest-energy experimen-  
288 tal phase, i.e.  $\alpha$ , reducing the number of polymorphs to 215. The phonon dispersion  
289 was computed for each of them along the high-symmetry path corresponding to their  
290 crystal symmetry, rejecting any phase displaying imaginary modes. Although these  
291 phases may give rise to lower energy structures after undergoing phase transitions, we  
292 did not explore such possibilities due to the complexity of the problem. Instead, we as-  
293 sume that the initial set of 18000 structures already includes any of those lower-energy  
294 structures. This leaves 70 dynamically stable phases.
- 295 • Finally, an annealing process was performed, which involved increasing the temperature  
296 up to 400K and then cooling it down to 0K, followed by a geometry optimization.  
297 Since MD simulations might break the crystal symmetry, a lousy symmetry check was  
298 performed to identify the symmetry group using the ASE interface of `Spglib`.<sup>47,48</sup> The  
299 process concludes with one last structure optimization with fixed symmetry. After this  
300 refinement, the structural dissimilarity analysis was performed again, removing similar  
301 and identical structures.

302 After the computational screening, only 12 possible candidates remain from the initial  
303 18000 structures. All the candidates are orthosilicates with IV-coordinated silicon, consistent  
304 with the experimental phases. It is worth noting that two structures ( $S_{12}$  and  $S_5$ ) present  
305 only translational symmetry (P1 space group) and have large unit cells; therefore, we argue  
306 that they could be classified as glasses. These structures were generated in the EA stage and  
307 survived the annealing stage and the stability checks. Furthermore, none of the candidate  
308 structures were included in the Materials Projects database and therefore were not part  
309 of the training set. As a matter of fact, all non-experimental phases in the training set  
310 have energies above the  $\alpha$  phase and all of the polymorphs found by EA, except for the  $S_{12}$   
311 phase. However, this phase has already been discarded. For instance, the energy of the next

312 polymorph with the highest energy ( $S_{11}$ ) predicted by the MLP is 0.64 eV/f.u. above  $\gamma$ ,  
 313 while the remaining non-experimental phases in the database are approximately 0.7 eV/f.u.  
 314 above that reference.

## 315 Reactivity analysis of the $C_2S$ candidates

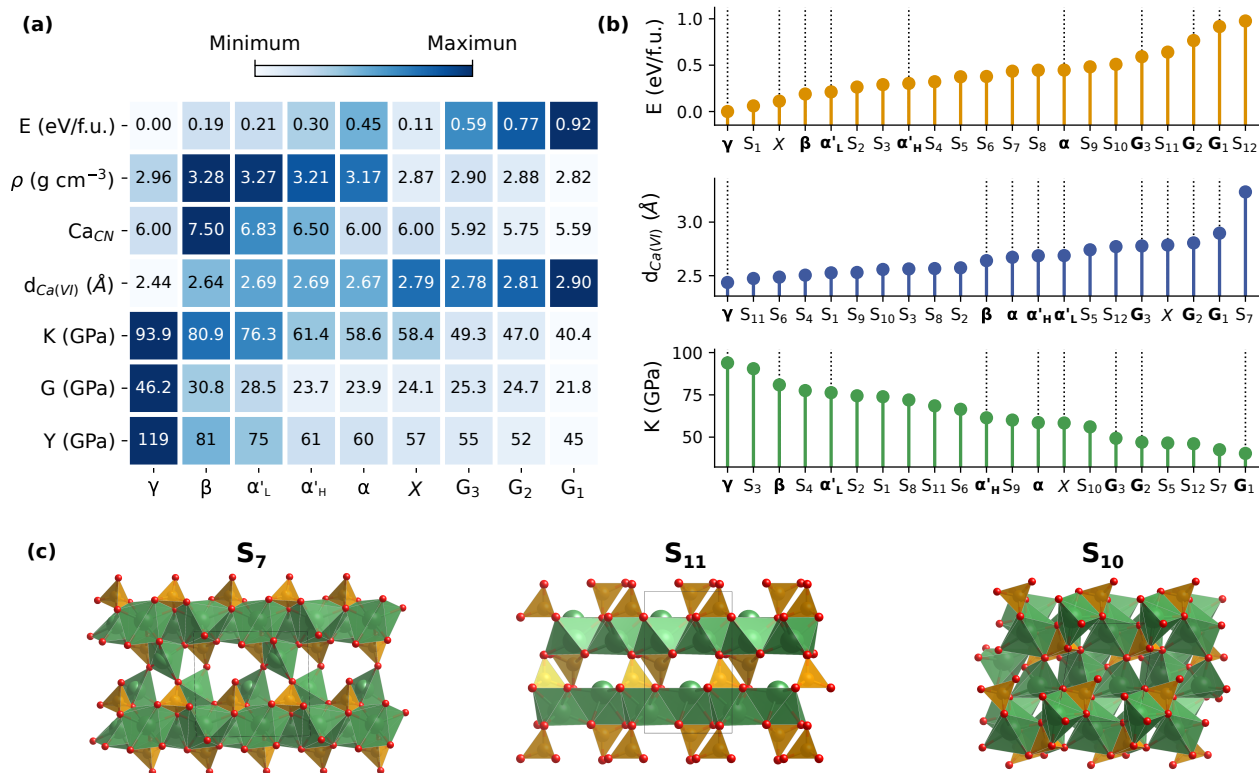


Figure 4: **Reactivity analysis of the candidates.** (a) Several properties computed for the known polymorphs and the three generated glassy structures: energy per formula unit (E), density ( $\rho$ ), averaged coordination number for calcium ( $Ca_{CN}$ ), average Ca-O distance to fulfill the Ca coordination shell ( $d_{Ca-O(VI)}$ ), bulk modulus (K), shear modulus (G) and Young's modulus (Y). The color ranges from the minimum to the maximum value of each magnitude, and the polymorphs are ordered from least to most reactive. (b) Energy, average Ca-O distance of the Ca coordination shell and bulk modulus of all the candidates. (c) Selected candidates with the potential to display high reactivity.

316 The final aim of this work is to find potential  $C_2S$  polymorphs with high dissolution  
 317 rates, and for that, we need an atomic scale reactivity indicator. Unfortunately, a quantita-  
 318 tive prediction of mineral dissolution rates based on atomistic simulations is a complicated  
 319 task. Electronic structure calculations are suitable for surface chemisorption reactions, for

320 example, to predict the catalytic properties of materials.<sup>49,50</sup> However, they are less appro-  
321 priate to predict dissolution, as the individual water chemisorption at the mineral surface do  
322 not correlate with the dissolution rates. For instance, it has been shown that the water dis-  
323 sociation reaction in  $\gamma$ -C<sub>2</sub>S can be barrier-less despite being the polymorph with the lowest  
324 dissolution rate. In contrast, water dissociation at the fast-dissolving  $\beta$ -C<sub>2</sub>S surface presents  
325 energy barriers between 2.5 and 47.2 kJ mol<sup>-1</sup>. Furthermore, electronic properties are not  
326 accessible to the MLP, so we need a structure or energy-based descriptor.<sup>51</sup>

327 A key factor that determines mineral dissolution rates is the ligand-cation exchange abil-  
328 ity.<sup>52-54</sup> There is a clear correlation between the dissolution rates of oxides, orthosilicates,  
329 and carbonates, with the water exchange rate of the forming cation in solution.<sup>52,54</sup> We as-  
330 sume that the ligand-cation exchange mechanism applies to phases with the same cation,  
331 and the exchange tendency is related to how strongly bonded is the cation in the structure,  
332 reflected in factors such as coordination and flexibility. To find correlations we plotted in  
333 Figure 4 (a) the binding energy and several structural and elastic properties of the experi-  
334 mental polymorphs in increasing order of hydraulic reactivity. Unfortunately, there are no  
335 quantitative sample-independent values for the dissolution rate of these polymorphs, but it  
336 is well-known that they follow the order (from least to most reactive)  $\gamma < \beta < \alpha$ .<sup>6</sup> More-  
337 over, the high dissolution rate of X-C<sub>2</sub>S samples has been suggested to correspond actually  
338 to an amorphous coexisting phase. Therefore, three amorphous dicalcium silicate models  
339 have been constructed from MD simulations by heating a  $\gamma$ -C<sub>2</sub>S supercell to 700K and rapid  
340 cooling at different cooling rates (see Methods for details). These models are denoted as G1,  
341 G2 and G3, from lowest to highest cooling rate. The final expected order for the dissolution  
342 rate is  $\gamma < \beta < \alpha < X < G3 < G2 < G1$ , with uncertainty about the actual ranking of the  
343 X phase.

344 Figure 4 (a) shows that there is a correlation between the cohesive energy and the re-  
345 activity: the more energetically stable the polymorph is, the less reactive they are. The X  
346 phase is an exception to the trend, but it has already been discussed that its reactivity may



347 not have been properly quantified. In fact, all three amorphous models have higher energy  
348 than the rest of the polymorphs. Nevertheless, the  $X$  phase correlates with the remaining  
349 properties, while the  $\gamma$  phase, undoubtedly the least reactive phase, breaks the trends in  
350 density. For all the other phases a lower density, higher Ca-O distance to fulfill the Ca  
351 coordination shell  $d_{\text{Ca(VI)}}$ , lower  $\text{Ca}_{\text{CN}}$ , and lower elastic properties correlate with a higher  
352 reactivity. It is interesting to note that the amorphous structures ranked as a function of  
353 the cooling rate follow the correlations, which suggests that amorphous  $\text{C}_2\text{S}$  could be indeed  
354 responsible for the high dissolution rate in samples with  $X\text{-C}_2\text{S}$ .

355 Considering the previous trends, we propose three magnitudes as potential reactivity  
356 descriptors of the dicalcium silicates: the cohesive energy, the Ca-O distance ( $d_{\text{Ca(VI)}}$ ), and  
357 the bulk modulus ( $K$ ). Figure 4 (b) ranks the 12 new polymorphs according to these three  
358 reactivity indicators. In theory, the phases located at the right end (high energy, high  $d_{\text{Ca(VI)}}$ ,  
359 and low  $K$ ) should be the most reactive polymorphs, and the objective of our search. Despite  
360 being in the range of interest, we will not consider the  $S_5$  and  $S_{12}$  because they are amorphous  
361 systems as mentioned before, and therefore it is natural that they lie close to the amorphous  
362 models. Out of the remaining candidates, the  $S_7$  structure shows the most promise. It has  
363 the lowest bulk modulus among the crystalline structures, half the value of  $\gamma\text{-C}_2\text{S}$ , and the  
364 largest  $d_{\text{Ca(VI)}}$ . The  $S_7$  is a layered structure, with CaO forming a central sheet and silicate  
365 groups at both sides with three of their four oxygen atoms coordinated to the CaO central  
366 sheet. The interlayer space contains a Ca atom that links consecutive layers. The  $S_{11}$  is  
367 a similar layered structure, without Ca in the interlayer space. Its energy is in the upper  
368 range, only surpassed by the amorphous structures. However, its overall potential is limited  
369 due to low values of other indicators. Other possible candidates that may display higher  
370 reactivity than the known phases include  $S_{10}$ , a bulk phase with a monoclinic axis, which  
371 also exhibits high energy and low bulk modulus.

## 372 4 Conclusions

373 In this work, we have introduced a hitherto unexplored approach to efficiently generating  
374 accurate MLPs based on transfer learning (TL) from the ReaxFF reactive force field to  
375 DFT. Previous attempts to reuse a lower-quality training set to reduce the amount of high-  
376 quality data were limited to using quantum methods as both low- and high-quality data.  
377 The present study demonstrates that TL from a classical force fields to DFT is both feasible  
378 and effective. In particular, we find that building the MLPs from models pre-trained on  
379 ReaxFF data can boost their accuracy from 10% to 40% in both energy and forces. Very  
380 importantly, the generation of the data for pre-training is virtually free of computational cost,  
381 and the methodology has no drawback: the TL-enhanced MLPs outperform those trained  
382 exclusively on DFT data in every tested scenario. In addition, MLPs pre-trained on large  
383 datasets made by empirical potentials can cover larger regions of the configurational space,  
384 providing flexibility and generality to the potential.

385 The TL methodology has been applied to build a MLP for calcium silicates. First,  
386 the MLPs were pre-trained on a dataset of 20000 ReaxFF configurations, followed by a  
387 refinement on 8000 DFT data points. The resulting MLP can successfully reproduce the  
388 DFT energies and forces with a mean absolute error of 4.8 meV/atom and  $0.25 \text{ eV}\text{\AA}^{-1}$   
389 respectively, as well as phonon spectra and elastic properties of calcium silicate crystals.  
390 This potential has been used to search for new dicalcium silicate polymorphs, aiming to  
391 find new (and hopefully highly reactive) belite phases. The combination of the DFT-like  
392 accuracy with the efficiency of classical potentials permits to examine and sieve thousands  
393 of polymorphs. In particular, we generated 18000 structures using EA, which were filtered  
394 using a computational screening protocol to discard duplicates, supercells, and dynamically  
395 unstable structures according to their phonon spectra and annealing at 400K. From the  
396 initial 18000, we identified 10 new crystalline  $\text{C}_2\text{S}$  polymorphs that are potentially stable.  
397 Based on our mechanical and structural descriptors of reactivity, a layered structure, denoted  
398 as  $\text{S}_7$  in this work, is particularly promising for displaying higher hydraulic activity than the

399 currently known belite phases.

400 The next step will be to investigate the hydration of these structures by performing  
401 molecular dynamics simulations at the crystal/water interfaces. If the proposed polymorphs  
402 are indeed highly reactive phases, it will be essential to test their thermodynamic stabiliza-  
403 tion by guest ions, in order to guide the synthesis and eventual production of highly reactive  
404 belitic cements. To conduct these studies, the computational work should focus on exploiting  
405 the presented TL methodology to include large and complex systems beyond the DFT ca-  
406 pabilities. This could include belite/water interfaces,<sup>55</sup> complexes and clusters in solution,<sup>56</sup>  
407 an extension of the MLP to new chemical species etc., allowing quantitative studies of the  
408 C<sub>2</sub>S reactivity and stability.

## 409 **Acknowledgement**

410 The authors thank Professor Iñigo Etxebarria for his valuable comments on the manuscript.  
411 This work was supported by the “Departamento de Educación, Política Lingüística y Cul-  
412 tura del Gobierno Vasco” (Grant No. IT1458-22), the "Ministerio de Ciencia e Innovación"  
413 (TED2021-130860B-I00), the University of the Basque Country UPV/EHU (Colab22/06)  
414 and the Transnational Common Laboratory “Aquitaine-Euskadi Network in Green Concrete  
415 and Cement-based Materials” (LTC-Green Concrete). The authors thank for technical and  
416 human support provided by SGIker (UPV/EHU/ ERDF, EU). J.L.-Z. acknowledges the  
417 financial support from the Basque Country Government (PRE\_2019\_1\_0025).

## 418 **Supporting Information Available**

419 The database created for this work is available in Zenodo [https://doi.org/10.5281/](https://doi.org/10.5281/zenodo.10796241)  
420 [zenodo.10796241](https://doi.org/10.5281/zenodo.10796241), together with the generated MLPs and the structures of the final can-  
421 didates. All the software used in this work is open source. The modifications to the  
422 `ænet-PyTorch` code to perform transfer learning can be found at <https://github.com/>

423 jlopez141/TL\_aenet-PyTorch. The Supplementary Information contains:

- 424 • Chebyshev descriptors and structural similarity.
- 425 • Performance of machine learning potentials.
- 426 • Training data set details.
- 427 • Hyperparameters for training and transfer learning.
- 428 • Transfer Learning with different compositions.
- 429 • Details about the structures generated by EA.
- 430 • Validation of the MLPs.

## 431 References

- 432 (1) Activity Report 2022. [https://cembureau.eu/media/m3jcyfre/  
433 cembureau-activity-report-2022-light.pdf](https://cembureau.eu/media/m3jcyfre/cembureau-activity-report-2022-light.pdf), 2023; Cembureau: The European  
434 Cement Association.
- 435 (2) Environment, U.; Scrivener, K. L.; John, V. M.; Gartner, E. M. Eco-efficient cements:  
436 Potential economically viable solutions for a low-CO<sub>2</sub> cement-based materials industry.  
437 *Cement and concrete Research* **2018**, *114*, 2–26.
- 438 (3) Barcelo, L.; Kline, J.; Walenta, G.; Gartner, E. Cement and carbon emissions. *Materials  
439 and structures* **2014**, *47*, 1055–1065.
- 440 (4) Ghosh, S. N.; Rao, P. B.; Paul, A.; Raina, K. The chemistry of dicalcium silicate  
441 mineral. *Journal of Materials Science* **1979**, *14*, 1554–1566.
- 442 (5) Wang, L.; Yang, H.; Zhou, S.; Chen, E.; Tang, S. Hydration, mechanical property and  
443 CSH structure of early-strength low-heat cement-based materials. *Materials Letters*  
444 **2018**, *217*, 151–154.

- 445 (6) Cuesta, A.; Ayuela, A.; Aranda, M. A. Belite cements and their activation. *Cement and*  
446 *Concrete Research* **2021**, *140*, 106319.
- 447 (7) Zhao, Y.; Lu, L.; Wang, S.; Gong, C.; Huang, Y. Modification of dicalcium silicates  
448 phase composition by BaO, SO<sub>3</sub> and MgO. *Journal of Inorganic and Organometallic*  
449 *Polymers and Materials* **2013**, *23*, 930–936.
- 450 (8) Kim, Y.-M.; Hong, S.-H. Influence of minor ions on the stability and hydration rates  
451 of  $\beta$ -dicalcium silicate. *Journal of the American Ceramic Society* **2004**, *87*, 900–905.
- 452 (9) Link, T.; Bellmann, F.; Ludwig, H.; Haha, M. B. Reactivity and phase composition of  
453 Ca<sub>2</sub>SiO<sub>4</sub> binders made by annealing of alpha-dicalcium silicate hydrate. *Cement and*  
454 *Concrete Research* **2015**, *67*, 131–137.
- 455 (10) Pirvan, A. A.; Haha, M. B.; Boehm-Courjault, E.; Scrivener, K. L. Calcium-silicate-  
456 hydrates from reactive dicalcium silicate binder. 39 th Cement and Concrete Science  
457 Conference 2019. 2019; p 215.
- 458 (11) Behler, J.; Parrinello, M. Generalized neural-network representation of high-  
459 dimensional potential-energy surfaces. *Physical review letters* **2007**, *98*, 146401.
- 460 (12) Behler, J. Neural network potential-energy surfaces in chemistry: a tool for large-scale  
461 simulations. *Physical Chemistry Chemical Physics* **2011**, *13*, 17930–17955.
- 462 (13) Yamada, H.; Liu, C.; Wu, S.; Koyama, Y.; Ju, S.; Shiomi, J.; Morikawa, J.; Yoshida, R.  
463 Predicting materials properties with little data using shotgun transfer learning. *ACS*  
464 *central science* **2019**, *5*, 1717–1730.
- 465 (14) Jha, D.; Choudhary, K.; Tavazza, F.; Liao, W.-k.; Choudhary, A.; Campbell, C.;  
466 Agrawal, A. Enhancing materials property prediction by leveraging computational and  
467 experimental data using deep transfer learning. *Nature communications* **2019**, *10*, 5316.

- 468 (15) Cai, C.; Wang, S.; Xu, Y.; Zhang, W.; Tang, K.; Ouyang, Q.; Lai, L.; Pei, J. Transfer  
469 learning for drug discovery. *Journal of Medicinal Chemistry* **2020**, *63*, 8683–8694.
- 470 (16) Smith, J. S.; Nebgen, B. T.; Zubatyuk, R.; Lubbers, N.; Devereux, C.; Barros, K.;  
471 Tretiak, S.; Isayev, O.; Roitberg, A. E. Approaching coupled cluster accuracy with a  
472 general-purpose neural network potential through transfer learning. *Nature communi-*  
473 *cations* **2019**, *10*, 2903.
- 474 (17) Kaser, S.; Boittier, E. D.; Upadhyay, M.; Meuwly, M. Transfer learning to CCSD (T):  
475 Accurate anharmonic frequencies from machine learning models. *Journal of Chemical*  
476 *Theory and Computation* **2021**, *17*, 3687–3699.
- 477 (18) Ramakrishnan, R.; Dral, P. O.; Rupp, M.; Von Lilienfeld, O. A. Big data meets quan-  
478 tum chemistry approximations: the  $\Delta$ -machine learning approach. *Journal of chemical*  
479 *theory and computation* **2015**, *11*, 2087–2096.
- 480 (19) Bogojeski, M.; Vogt-Maranto, L.; Tuckerman, M. E.; Müller, K.-R.; Burke, K. Quantum  
481 chemical accuracy from density functional approximations via machine learning. *Nature*  
482 *communications* **2020**, *11*, 5223.
- 483 (20) Ruth, M.; Gerbig, D.; Schreiner, P. R. Machine learning of coupled cluster (T)-energy  
484 corrections via delta ( $\Delta$ )-Learning. *Journal of Chemical Theory and Computation* **2022**,  
485 *18*, 4846–4855.
- 486 (21) Maier, S.; Collins, E. M.; Raghavachari, K. Quantitative Prediction of Vertical Ioniza-  
487 tion Potentials from DFT via a Graph-Network-Based Delta Machine Learning Model  
488 Incorporating Electronic Descriptors. *The Journal of Physical Chemistry A* **2023**, *127*,  
489 3472–3483.
- 490 (22) Sun, G.; Sautet, P. Toward fast and reliable potential energy surfaces for metallic Pt  
491 clusters by hierarchical delta neural networks. *Journal of chemical theory and compu-*  
492 *tation* **2019**, *15*, 5614–5627.

- 493 (23) Van Duin, A. C.; Dasgupta, S.; Lorant, F.; Goddard, W. A. ReaxFF: a reactive force  
494 field for hydrocarbons. *The Journal of Physical Chemistry A* **2001**, *105*, 9396–9409.
- 495 (24) Senftle, T. P.; Hong, S.; Islam, M. M.; Kylasa, S. B.; Zheng, Y.; Shin, Y. K.; Junker-  
496 meier, C.; Engel-Herbert, R.; Janik, M. J.; Aktulga, H. M.; others The ReaxFF reactive  
497 force-field: development, applications and future directions. *npj Computational Mate-*  
498 *rials* **2016**, *2*, 1–14.
- 499 (25) Perdew, J. P.; Burke, K.; Ernzerhof, M. Generalized gradient approximation made  
500 simple. *Physical review letters* **1996**, *77*, 3865.
- 501 (26) Van Duin, A. C.; Strachan, A.; Stewman, S.; Zhang, Q.; Xu, X.; Goddard, W. A.  
502 ReaxFFSiO reactive force field for silicon and silicon oxide systems. *The Journal of*  
503 *Physical Chemistry A* **2003**, *107*, 3803–3811.
- 504 (27) Manzano, H.; Pellenq, R. J.; Ulm, F.-J.; Buehler, M. J.; Van Duin, A. C. Hydration of  
505 calcium oxide surface predicted by reactive force field molecular dynamics. *Langmuir*  
506 **2012**, *28*, 4187–4197.
- 507 (28) Thompson, A. P.; Aktulga, H. M.; Berger, R.; Bolintineanu, D. S.; Brown, W. M.;  
508 Crozier, P. S.; in 't Veld, P. J.; Kohlmeyer, A.; Moore, S. G.; Nguyen, T. D.; Shan, R.;  
509 Stevens, M. J.; Tranchida, J.; Trott, C.; Plimpton, S. J. LAMMPS - a flexible simulation  
510 tool for particle-based materials modeling at the atomic, meso, and continuum scales.  
511 *Comp. Phys. Comm.* **2022**, *271*, 108171.
- 512 (29) Giannozzi, P.; Baroni, S.; Bonini, N.; Calandra, M.; Car, R.; Cavazzoni, C.; Ceresoli, D.;  
513 Chiarotti, G. L.; Cococcioni, M.; Dabo, I.; others QUANTUM ESPRESSO: a modu-  
514 lar and open-source software project for quantum simulations of materials. *Journal of*  
515 *physics: Condensed matter* **2009**, *21*, 395502.
- 516 (30) Giannozzi, P.; Andreussi, O.; Brumme, T.; Bunau, O.; Nardelli, M. B.; Calandra, M.;  
517 Car, R.; Cavazzoni, C.; Ceresoli, D.; Cococcioni, M.; others Advanced capabilities for

- 518 materials modelling with Quantum ESPRESSO. *Journal of physics: Condensed matter*  
519 **2017**, *29*, 465901.
- 520 (31) Hamann, D. Optimized norm-conserving Vanderbilt pseudopotentials. *Physical Review*  
521 *B* **2013**, *88*, 085117.
- 522 (32) van Setten, M. J.; Giantomassi, M.; Bousquet, E.; Verstraete, M. J.; Hamann, D. R.;  
523 Gonze, X.; Rignanese, G.-M. The PseudoDojo: Training and grading a 85 element  
524 optimized norm-conserving pseudopotential table. *Computer Physics Communications*  
525 **2018**, *226*, 39–54.
- 526 (33) Grimme, S. Semiempirical GGA-type density functional constructed with a long-range  
527 dispersion correction. *Journal of computational chemistry* **2006**, *27*, 1787–1799.
- 528 (34) López-Zorrilla, J.; Aretxabaleta, X. M.; Yeu, I. W.; Etxebarria, I.; Manzano, H.; Ar-  
529 trith, N. ænet-PyTorch: A GPU-supported implementation for machine learning atomic  
530 potentials training. *The Journal of Chemical Physics* **2023**, *158*.
- 531 (35) Artrith, N.; Urban, A. An implementation of artificial neural-network potentials for  
532 atomistic materials simulations: Performance for TiO<sub>2</sub>. *Computational Materials Sci-*  
533 *ence* **2016**, *114*, 135–150.
- 534 (36) Artrith, N.; Urban, A.; Ceder, G. Efficient and accurate machine-learning interpolation  
535 of atomic energies in compositions with many species. *Physical Review B* **2017**, *96*,  
536 014112.
- 537 (37) Oganov, A. R.; Glass, C. W. Crystal structure prediction using ab initio evolutionary  
538 techniques: Principles and applications. *The Journal of chemical physics* **2006**, *124*.
- 539 (38) Oganov, A. R.; Lyakhov, A. O.; Valle, M. How Evolutionary Crystal Structure Predic-  
540 tion Works and Why. *Accounts of chemical research* **2011**, *44*, 227–237.



- 541 (39) Lyakhov, A. O.; Oganov, A. R.; Stokes, H. T.; Zhu, Q. New developments in evolution-  
542 ary structure prediction algorithm USPEX. *Computer Physics Communications* **2013**,  
543 *184*, 1172–1182.
- 544 (40) Chen, M. S.; Morawietz, T.; Mori, H.; Markland, T. E.; Artrith, N. AENET–LAMMPS  
545 and AENET–TINKER: Interfaces for accurate and efficient molecular dynamics simu-  
546 lations with machine learning potentials. *The Journal of Chemical Physics* **2021**, *155*.
- 547 (41) Mori, H.; Tsuru, T.; Okumura, M.; Matsunaka, D.; Shiihara, Y.; Itakura, M. Dynamic  
548 interaction between dislocations and obstacles in bcc iron based on atomic potentials  
549 derived using neural networks. *Physical Review Materials* **2023**, *7*, 063605.
- 550 (42) Togo, A.; Chaput, L.; Tadano, T.; Tanaka, I. Implementation strategies in phonopy  
551 and phono3py. *J. Phys. Condens. Matter* **2023**, *35*, 353001.
- 552 (43) Togo, A. First-principles Phonon Calculations with Phonopy and Phono3py. *J. Phys.*  
553 *Soc. Jpn.* **2023**, *92*, 012001.
- 554 (44) Jain, A.; Ong, S. P.; Hautier, G.; Chen, W.; Richards, W. D.; Dacek, S.; Cholia, S.;  
555 Gunter, D.; Skinner, D.; Ceder, G.; others Commentary: The Materials Project: A  
556 materials genome approach to accelerating materials innovation. *APL materials* **2013**,  
557 *1*.
- 558 (45) Oganov, A. R.; Valle, M. How to quantify energy landscapes of solids. *The Journal of*  
559 *chemical physics* **2009**, *130*.
- 560 (46) Ceriotti, M.; Tribello, G. A.; Parrinello, M. Simplifying the representation of com-  
561 plex free-energy landscapes using sketch-map. *Proceedings of the National Academy of*  
562 *Sciences* **2011**, *108*, 13023–13028.
- 563 (47) Larsen, A. H.; Mortensen, J. J.; Blomqvist, J.; Castelli, I. E.; Christensen, R.;  
564 Dułak, M.; Friis, J.; Groves, M. N.; Hammer, B.; Hargus, C.; others The atomic sim-

- 565        ulation environment—a Python library for working with atoms. *Journal of Physics:*  
566        *Condensed Matter* **2017**, *29*, 273002.
- 567 (48) Togo, A.; Tanaka, I. Spglib: a software library for crystal symmetry search. *arXiv*  
568        *preprint arXiv:1808.01590* **2018**,
- 569 (49) Nørskov, J. K.; Abild-Pedersen, F.; Studt, F.; Bligaard, T. Density functional theory  
570        in surface chemistry and catalysis. *Proceedings of the National Academy of Sciences*  
571        **2011**, *108*, 937–943.
- 572 (50) Nørskov, J. K.; Bligaard, T.; Rossmeisl, J.; Christensen, C. H. Towards the computa-  
573        tional design of solid catalysts. *Nature chemistry* **2009**, *1*, 37–46.
- 574 (51) Rejmak, P.; Dolado, J. S.; Aranda, M. A.; Ayuela, A. First-principles calculations on  
575        polymorphs of dicalcium silicate—Belite, a main component of Portland cement. *The*  
576        *Journal of Physical Chemistry C* **2019**, *123*, 6768–6777.
- 577 (52) Ohlin, C. A.; Villa, E. M.; Rustad, J. R.; Casey, W. H. Dissolution of insulating oxide  
578        materials at the molecular scale. *Nature materials* **2010**, *9*, 11–19.
- 579 (53) Rustad, J. R.; Casey, W. H. Metastable structures and isotope exchange reactions in  
580        polyoxometalate ions provide a molecular view of oxide dissolution. *Nature materials*  
581        **2012**, *11*, 223–226.
- 582 (54) Casey, W. H.; Swaddle, T. W. Why small? The use of small inorganic clusters to under-  
583        stand mineral surface and dissolution reactions in geochemistry. *Reviews of Geophysics*  
584        **2003**, *41*.
- 585 (55) Qi, C.; Manzano, H.; Spagnoli, D.; Chen, Q.; Fourie, A. Initial hydration process of  
586        calcium silicates in Portland cement: A comprehensive comparison from molecular  
587        dynamics simulations. *Cement and Concrete Research* **2021**, *149*, 106576.

588 (56) Aretxabaleta, X. M.; López-Zorrilla, J.; Etxebarria, I.; Manzano, H. Multi-step nu-  
589 cleation pathway of CSH during cement hydration from atomistic simulations. *Nature*  
590 *Communications* **2023**, *14*, 7979.