

1 Transfer learning based on atomic feature extraction for  
2 the prediction of experimental  $^{13}\text{C}$  chemical shifts<sup>†</sup>

3 Žarko Ivković<sup>1,2</sup>, Jesús Jover<sup>1</sup>, and Jeremy Harvey<sup>2</sup>

4 <sup>1</sup>Institut de Química Teòrica i Computacional (IQTC), Faculty of Chemistry,  
5 University of Barcelona, Spain

6 <sup>2</sup>Department of Chemistry, KU Leuven, Belgium

7 June 19, 2024

---

<sup>†</sup> Electronic Supplementary Information (ESI) available.

## 8 Abstract

9 Forecasting experimental chemical shifts of organic compounds is a long-standing challenge  
10 in organic chemistry. Recent advances in machine learning (ML) have led to routines that  
11 surpass the accuracy of ab initio Density Functional Theory (DFT) in estimating experi-  
12 mental  $^{13}\text{C}$  shifts. The extraction of knowledge from other models, known as transfer learn-  
13 ing, has demonstrated remarkable improvements, particularly in scenarios with limited data  
14 availability. However, the extent to which transfer learning improves predictive accuracy in  
15 low-data regimes for experimental chemical shift predictions remains unexplored.

16 This study indicates that atomic features derived from a message passing neural network  
17 (MPNN) forcefield are robust descriptors for atomic properties. A dense network utilizing  
18 these descriptors to predict  $^{13}\text{C}$  shifts achieves a mean absolute error (MAE) of 1.68 ppm.  
19 When these features are used as node labels in a simple graph neural network (GNN),  
20 the model attains a better MAE of 1.34 ppm. On the other hand, embeddings from a self-  
21 supervised pre-trained 3D aware transformer are not sufficiently descriptive for a feedforward  
22 model but show reasonable accuracy within the GNN framework, achieving an MAE of 1.51  
23 ppm. Under low-data conditions, all transfer-learned models show a significant improvement  
24 in predictive accuracy compared to existing literature models, regardless of the sampling  
25 strategy used to select from the pool of unlabeled examples.

26 We demonstrated that extracting atomic features from models trained on large and di-  
27 verse datasets is an effective transfer learning strategy for predicting NMR chemical shifts,  
28 achieving results on par with existing literature models. This method provides several ben-  
29 efits, such as reduced training times, simpler models with fewer trainable parameters, and  
30 strong performance in low-data scenarios, without the need for costly ab initio data of the  
31 target property. This technique can be applied to other chemical tasks opening many new  
32 potential applications where the amount of data is a limiting factor.

33 Keywords: Machine Learning, Atomic Representation, Transfer Learning, Graph Neural  
34 Networks, NMR, Chemical Shifts, Feature Extraction, Low-data, Atomic Embeddings

## 35 Introduction

### 36 NMR Chemical Shifts

37 NMR chemical shifts are valuable in the structure elucidation of organic compounds within  
38 classical and computer-assisted frameworks.<sup>1-5</sup> Carbon chemical shifts have been used to  
39 elucidate reaction products<sup>6</sup>, metabolites<sup>7</sup>, and natural products, including in the revision  
40 of the structures.<sup>8-10</sup> Furthermore, chemical shifts carry information about the local chem-  
41 ical environments of atoms and have been used as descriptors for predicting chemical re-  
42 activity<sup>11,12</sup> and in QSAR/QSPR models<sup>13</sup>. Prediction of carbon chemical shifts from the  
43 molecular structure has been extensively studied and many methods have been developed,  
44 ranging from ab initio to fully data-driven methods.<sup>14,15</sup>

45 Predicting carbon NMR shifts from molecular structures from the first principles is com-  
46 putationally intensive. First, the geometry is optimized, followed by calculating the electronic  
47 structure. In addition to errors from the electronic structure calculations, treatment of solva-  
48 tion, conformational flexibility, and rovibronic effects introduce further errors.<sup>16</sup> Considering  
49 all these factors comprehensively is computationally impractical at any level of theory that  
50 ensures reasonable accuracy. For example, even a basic DFT calculation of chemical shifts  
51 on an inexpensive geometry is too resource-intensive for large-scale rapid structure elucida-  
52 tion. The chosen functional, basis set, and solvation model influences the precision of DFT  
53 predictions for NMR shifts.<sup>17,18</sup> Although different results in the literature are reported on  
54 different sets for the same computational protocols, the best-reported protocol achieves a  
55 root mean square error (RMSE) of 3.68 ppm when compared to experimental shifts.<sup>17</sup> This  
56 is insufficient for typical applications, as an initial investigation has shown that an accu-

57 racy of 1.1-1.2 ppm of MAE is necessary for correctly identifying 99% of molecules in the  
58 metabolomic database.<sup>19</sup>

59 The errors of DFT-predicted shifts have a systematic component that can be corrected  
60 using available experimental data. Lodewyk et al.<sup>16</sup> developed a linear scaling protocol for  
61 different combinations of levels of theory, solvents, and solvation models, and their findings  
62 were compiled in the CHESHIRE repository.<sup>20</sup> This became the standard for chemical shift  
63 prediction using DFT. Gao et al.<sup>21</sup> went beyond linear interpolation and constructed a  
64 deep neural network that takes molecular structure and descriptors derived from calculated  
65 DFT shielding constants as input to predict experimental chemical shifts. Their method  
66 demonstrated superior performance, achieving an RMSE of 2.10 ppm, which is a significant  
67 notable improvement over the 4.77 ppm RMSE the authors report from linear regression on  
68 the same small test set.

69 The Exp5K dataset, developed as part of the CASCADE project,<sup>12</sup> is the largest dataset  
70 that compares empirically scaled DFT chemical shifts with experimental shifts. The authors  
71 excluded structures where DFT significantly disagreed with experimental results to avoid  
72 introducing noise from potential misassignments in the experimental data. This exclusion  
73 inevitably removes challenging examples where the disagreement arises from DFT's inability  
74 to accurately predict shifts due to molecular complexity. Additionally, the atom ordering was  
75 altered when comparing DFT with experimental shifts, leading to the unjustified exclusion  
76 of some examples from the dataset. After correcting the atom order, the calculated shifts  
77 deviate from the experiments with an MAE of 2.21 ppm and an RMSE of 3.31 ppm.<sup>†</sup> This  
78 should be considered the most realistic measure of the accuracy of DFT-calculated shifts  
79 corrected with linear scaling. These correction methods, along with others reported in the  
80 literature,<sup>22,23</sup> enhance the accuracy of predictions but do not reduce their computational  
81 cost.

82 On the other hand, data-driven methods are significantly faster by several orders of

83 magnitude. The efficiency of machine learning in predicting carbon chemical shifts arises  
84 from the avoidance of expensive geometry optimizations or electronic structure computations.  
85 Nevertheless, the top models in the literature explicitly include geometrical data of the  
86 lowest energy conformers in their predictions.<sup>12,24–26</sup> The compromise is achieved by utilizing  
87 inexpensive forcefield geometries instead of costly DFT-optimized geometries.

88 The accuracy of predictions in data-driven models is influenced by the quality and quan-  
89 tity of the training data.<sup>27,28</sup> By using experimental data for training, common errors in ab  
90 initio methods can be avoided. The most extensive open NMR shift database with fully as-  
91 signed spectra is nmrshiftdb2.<sup>29,30</sup> User-contributed databases like this often face issues such  
92 as missing solvent and temperature details, peak misassignments, measurement noise, and  
93 incorrect structure identification. A model’s performance is limited not only by the quantity  
94 but also by the quality of data. Thus, models that perform well in low-data scenarios are  
95 necessary when data is scarce and when prioritizing high-quality data over quantity.

## 96 **Transfer Learning**

97 Transfer learning involves using a model trained on one task as a foundation for training  
98 on another task, known as a downstream task.<sup>31</sup> Generally, pre-training is performed on a  
99 similar task with a much larger dataset, followed by training on a smaller dataset for the  
100 specific task of interest. Feature extraction and fine-tuning are two main implementations of  
101 transfer learning.\* The choice of method depends on task similarity, the size and architecture  
102 of the pre-trained model, and the amount of available data. Feature extraction is commonly  
103 used in computer vision,<sup>32,33</sup> while fine-tuning is widely used in language models.<sup>34,35</sup>

104 One of the major challenges for machine learning in chemistry is the scarcity of train-  
105 ing data.<sup>36,37</sup> Acquiring experimental and high-quality ab initio data is costly, and more

---

\*In the literature, the term fine-tuning is not well-defined; it can refer to the second phase of training in general or to training models with weights initialized from other models. Here, we refer to the latter and simply call the second phase of training ‘training,’ as opposed to the ‘pre-training’ in the first phase.

106 affordable ab initio data often comes with substantial errors. Complex models, which are  
107 generally necessary to represent intricate chemical phenomena, demand a large amount of  
108 data for training. Integrating chemical and physical knowledge and intuition into the model  
109 architecture is one strategy to lessen the required training data.<sup>38</sup> Transfer learning provides  
110 an alternative method to enhance models and can be used alongside other techniques to  
111 address issues related to limited data for chemical problems.

112 Most previous studies employ transfer learning for chemical models by initially training  
113 models on data generated from ab initio methods and then fine-tuning them on experimental  
114 data.<sup>12,39,40</sup> This quasi-transfer approach is effective if a significantly larger amount of ab  
115 initio data compared to the available experimental data can be produced. However, certain  
116 experimental properties like the smell, catalytic activity, and reaction yield are difficult  
117 or impossible to model using ab initio methods, while calculating others such as NMR  
118 properties, free energies, and absorption spectra can be prohibitively costly. In such cases,  
119 pre-training must be conducted on less relevant tasks where it is feasible to generate large-  
120 scale datasets.

## 121 **Related work**

122 In the notable CASCADE study,<sup>12</sup> graph neural networks (GNN) were employed to pre-  
123 dict experimental chemical shifts. The ExpNN-ff model takes 3D structures optimized using  
124 MMFF forcefield as the way to incorporate geometrical information while maintaining rel-  
125 atively low computational cost. The authors implemented an interesting double-transfer  
126 learning training. First, the model was trained on DFT-optimized geometries and scaled  
127 DFT shifts. Second, the model was retrained on DFT-optimized geometries and experimen-  
128 tal shifts, keeping the interaction layers frozen. Finally, the model was retrained again on  
129 forcefield geometries and experimental shifts, keeping the readout layers frozen. It is unclear

130 what advantage this approach has over doing single-step transfer learning, updating all layers  
131 in the model simultaneously. Still, the ExpNN-ff model with an MAE of 1.43 ppm on a 500  
132 hold-out test set performs better than the DFT with empirical scaling which has an MAE  
133 of 2.21 ppm on the whole training dataset of around 5000 compounds.

134 To avoid the costly DFT calculations for large molecules during the generation of the pre-  
135 training dataset, Han and Choi<sup>39</sup> pretrained a GNN using the QM9 dataset of DFT shielding  
136 constants. They subsequently fine-tuned the model using an experimental chemical shifts  
137 database that includes larger molecules and atoms such as P, Cl, and S, which are absent in  
138 the QM9 dataset. The authors evaluated the model in low data scenarios, achieving an MAE  
139 of approximately 2.3 ppm with 2112 training examples. Nonetheless, the authors pre-trained  
140 on ab initio NMR data on a dataset comparable to the size of the experimental dataset used  
141 to fine-tune the model, similar to the approach used in CASCADE.

142 The first example of adopting true transfer learning for predicting chemical shifts was  
143 done in a recent work by El Samman et al.<sup>41</sup> The authors extracted atomic embeddings from  
144 the last interaction layer from the SchNet model<sup>42</sup> trained to predict molecular energies on  
145 the QM9 dataset. The authors tested linear and feedforward network models for different  
146 chemical tasks, including predicting carbon chemical shifts calculated by HOSE codes.<sup>43</sup>  
147 However, the dataset for the chemical shifts consisted of only 200 examples of shifts predicted  
148 by the HOSE code, so the performance relative to the literature models trained from scratch  
149 could not be assessed.

150 To tackle low-data scenarios without resorting to transfer learning, Rull et al.<sup>44</sup> modified  
151 a GNN architecture to enhance its efficiency in such conditions. While the modified archi-  
152 tecture performed better in low-data scenarios than a similar GNN model, it significantly  
153 underperformed in high-data scenarios. This underscores the importance of considering the  
154 volume of training data when evaluating model performance and designing model architec-  
155 tures.

## 156 Approach

157 In an ideal situation, pre-training is performed on a highly similar task for which either more  
158 data is available or it is significantly cheaper to generate. However, such tasks are rarely  
159 available for any downstream chemical task, necessitating some form of compromise. Many  
160 of the latest pre-trained chemical models employ self-supervised pre-training tasks on huge  
161 unlabeled datasets of 2D chemical structures.<sup>45–48</sup> Conversely, there are numerous instances  
162 of quasi-transfer learning, involving pre-training on datasets of ab initio calculated properties  
163 of the size comparable to the available experimental datasets.<sup>12,39</sup> We propose the atomic  
164 feature extraction from the models pre-trained for different chemical tasks on larger datasets,  
165 and we evaluate it by predicting experimental <sup>13</sup>C chemical shifts. The proposed approach  
166 is illustrated in Figure 1.

## 167 Choice of pre-training task and model

168 The downstream task in this study is to predict the chemical shifts of carbon atoms. Pre-  
169 dicting other atomic properties influenced by the chemical environment of the atom is the  
170 most relevant task. However, no other atomic properties have as extensive experimental data  
171 as chemical shifts. Fortunately, many models designed for predicting molecular properties  
172 incorporate atomic representations within their architectures.<sup>49,50</sup> Moreover, the pre-trained  
173 model must consider geometrical information since chemical shifts are influenced by molecu-  
174 lar conformation. Therefore, most pre-trained models based on 2D molecular structures are  
175 not suitable candidates. This leads us to neural network forcefields, whose architectures are  
176 designed to sum atomic energy contributions.\* We selected the MACE-OFF23 transferable  
177 organic forcefield<sup>51,52</sup>, which is state-of-the-art for predicting DFT molecular energies, open-  
178 source, and trained on a reasonably large dataset. Since we are not concerned with inference

---

\*This architecture design is not mandatory. The only requirement for architecture is the presence of atomic embeddings within the model

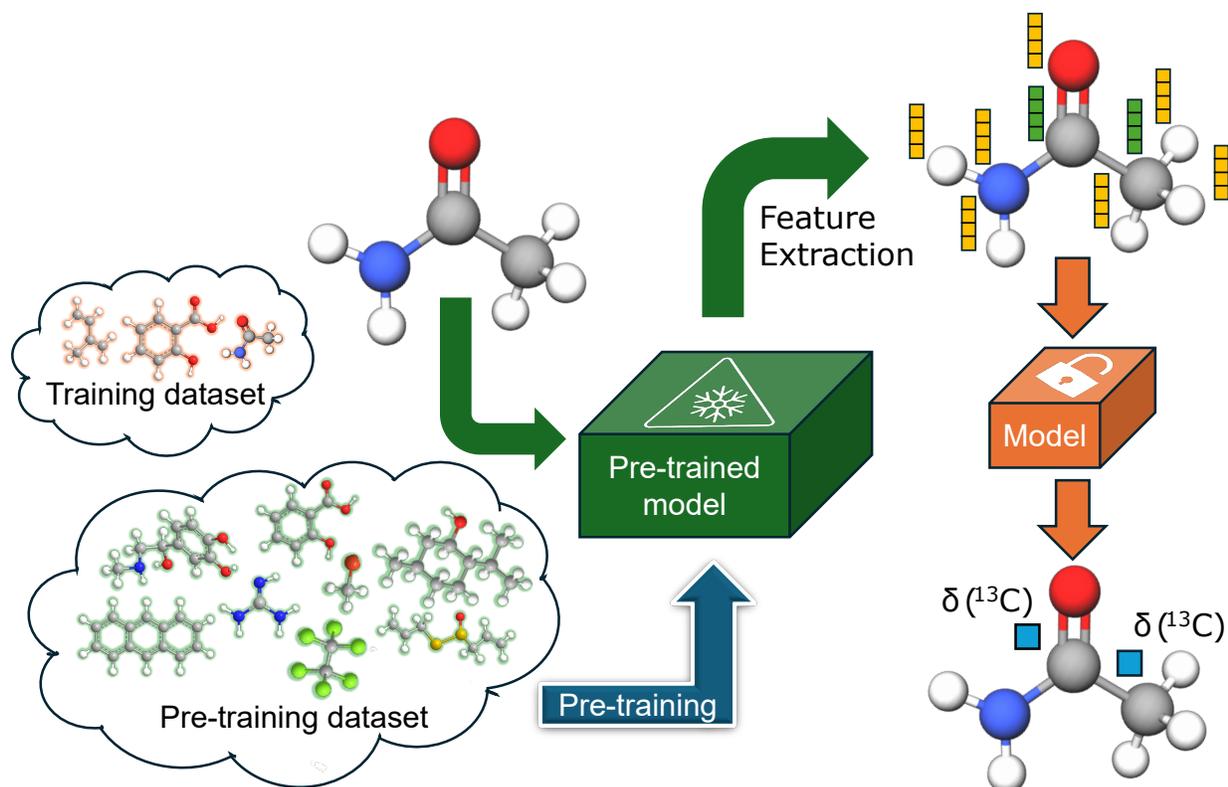


Figure 1: Transfer Learning based on atomic feature extraction.

179 time, we chose the large variant of the forcefield. The other model we tested is Uni-Mol<sup>53</sup>, a  
 180 3D-aware self-supervised pre-trained transformer known for its performance in downstream  
 181 molecular property prediction tasks. Although self-supervised pre-training is less directly  
 182 related to atomic property prediction, it is done on an even larger dataset. The model in-  
 183 cludes atomic representation in its architecture, and integrates geometrical information in  
 184 its embeddings, making it appropriate for this transfer learning approach.

## 185 Feature extraction

186 We extract atomic embeddings from the first of two interaction layers in the large variant  
 187 of the MACE-OFF23 forcefield. This approach contrasts with the method of El Samman et

188 al.<sup>41</sup>, where embeddings are extracted from the final interaction layer of the SchNet model.<sup>42</sup>  
189 We retain only the invariant portion of the embedding to ensure rotational and translational  
190 invariance, resulting in a 244-dimensional vector atomic embedding. Given that Uni-Mol is  
191 intended as a backbone pre-trained model for various downstream tasks, we directly extract  
192 the atomic representation from the output of the backbone, yielding a 512-dimensional vector  
193 per atom, invariant to translation and rotation. Both models use atomic coordinates and  
194 identities as inputs, akin to the input used by typical ab initio codes, and produce atomic  
195 embeddings for each atom as outputs.

## 196 **Models architecture**

197 We evaluated two distinct types of downstream models: a feedforward network (FFN) and a  
198 graph neural network (GNN). For the feedforward network, we assume that the pre-trained  
199 model has captured all necessary information regarding the chemical environment of each  
200 carbon atom. We use the embeddings of carbon atoms as input and train the network to  
201 predict chemical shifts. Additionally, we tested the GNN based on the GraphSAGE<sup>54</sup> archi-  
202 tecture, which facilitates the exchange of information between different atomic environment  
203 embeddings. This leads to a more robust model as it can learn more relevant embeddings  
204 for NMR shifts. Unlike the other methods where fully connected graphs with a cutoff dis-  
205 tance or graphs with implicitly represented hydrogens have been used, we used a chemical  
206 graph where all atoms are explicitly included. Consequently, GNN models require atomic  
207 connectivity as input, whereas FFN models only need atomic coordinates. Finally, after  
208 the message passing layers, the atomic embeddings of carbon atoms are fed into a readout  
209 feedforward network to predict chemical shifts. Both methodologies are illustrated in Figure  
210 2.

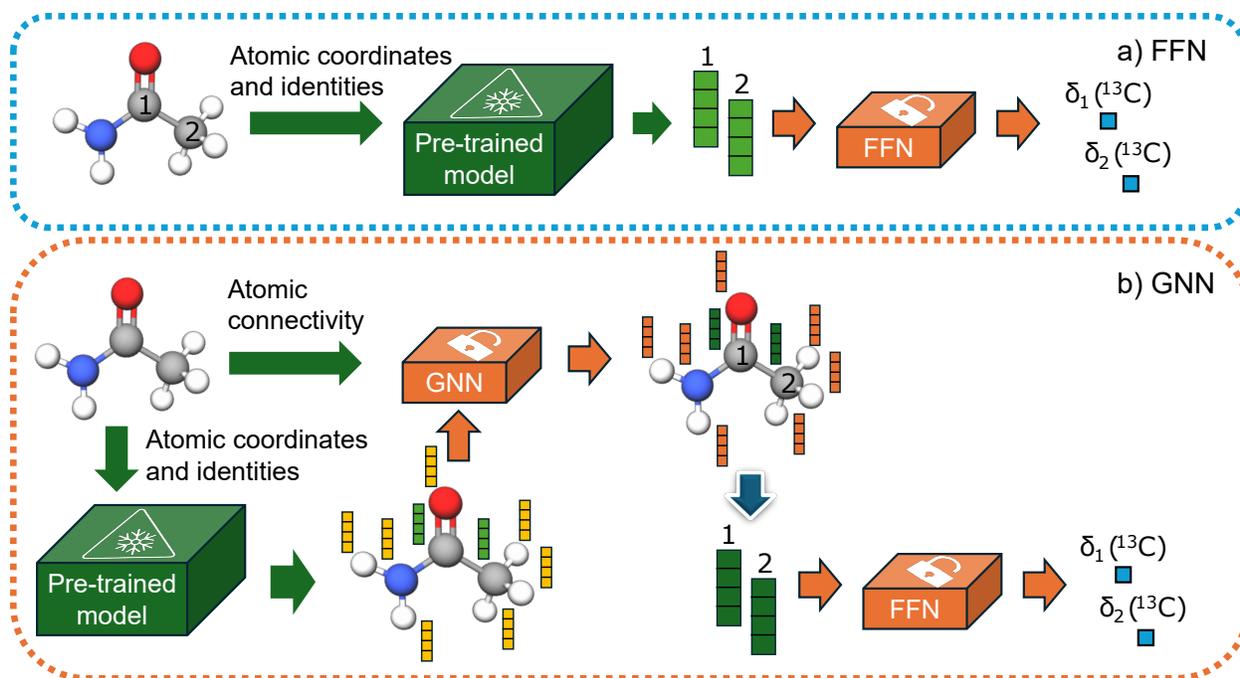


Figure 2: a) FNN model b) GNN model.  
 Only orange models are trained, while the green models' weights are frozen.

## 211 Low-data regimes

212 To evaluate model performance with fewer training examples, we selected varying quantities  
 213 of samples from the original dataset, treating it as a pool of unlabeled examples. Although  
 214 this dataset is smaller than the typical molecular datasets of unlabeled molecules, it is suffi-  
 215 ciently large to compare different sampling methods. We examined three sampling strategies:  
 216 random sampling, MaxMin<sup>55</sup> sampling based on the Tanimoto distance<sup>56</sup> between Morgan  
 217 fingerprints<sup>57</sup>, and MaxMin sampling based on the undirected Hausdorff distance<sup>58</sup> between  
 218 sets of transferred embeddings of all carbon atoms in two molecules. The directed Hausdorff  
 219 distance between two sets of vectors  $A$  and  $B$  is defined as:

$$h(A, B) = \max_{a \in A} \min_{b \in B} d(a, b)$$

220 where  $d(a, b)$  is any distance metric between two vectors. However, the directed Hausdorff  
221 distance is not symmetric, so we use the undirected Hausdorff distance, employing the Eu-  
222 clidean distance as the distance metric  $d$ :

$$H(A, B) = \max(h(A, B), h(B, A))$$

$$h(A, B) = \max_{a \in A} \min_{b \in B} \|a - b\|^2$$

223 In our scenario, sets of vectors represent sets of transferred embeddings of carbon atoms.  
224 While we could have used embeddings of all atoms, the carbon atom embeddings also convey  
225 information about their neighboring atoms. Since our primary interest lies in the differences  
226 in carbon atom environments between two molecules, we used only the embeddings of carbon  
227 atoms, which also reduces the computational cost, a crucial factor when sampling large pools  
228 of examples.

## 229 Results

230 The mean absolute error (MAE), root mean square error (RMSE), and Pearson correlation  
231 coefficient ( $\rho$ ) for all models are presented in Table 1. The results are based on a modified  
232 test set, where we excluded a couple of broken examples from the original test set. Additional  
233 details, including more performance metrics for each model and examples of molecules where  
234 models fail, can be found in SI.<sup>†</sup> The ensemble of two independently trained GNN models  
235 performs the best, with the lowest MAE and RMSE. MACE models outperform their Uni-  
236 Mol equivalents significantly, indicating that the forcefield is an excellent option for the  
237 pre-training task. Even though the Uni-Mol GNN has a lower MAE than the MACE FFN  
238 model, its RMSE is higher, highlighting the necessity to report at least both MAE and RMSE

239 when reporting the model’s performance. Regarding parameter efficiency, MACE GNN is  
240 by far the best model.

Table 1: Performance on a test set and number of trainable parameters

Model	MAE [ppm]	RMSE [ppm]	$\rho$	$N^\circ$ params
MACE FFN	1.68	2.74	0.9986	$1.3 \times 10^6$
Uni-Mol FFN	2.07	3.40	0.9978	$1.8 \times 10^6$
Ensemble MACE & Uni-Mol FFN	1.65	2.68	0.9986	$3.1 \times 10^6$
MACE GNN	1.34	2.38	0.9989	$1.9 \times 10^6$
Uni-Mol GNN	1.51	2.81	0.9985	$9.3 \times 10^6$
Ensemble MACE & Uni-Mol GNN	<b>1.28</b>	<b>2.37</b>	<b>0.9989</b>	$1.0 \times 10^7$

241 A comparison with relevant literature models that take forcefield geometries as input is  
242 shown in Figure 3. The ensemble of two GNNs and MACE GNN performs equally well as the  
243 best-reported literature models. Comparison with models trained using the same train/test  
244 split is more reliable, and the FullSSPrUCe model is trained on the larger portion of the  
245 nmrshiftdb2 database, which explains its slightly better performance. In any case, since  
246 all reported models are solvent agnostic, it is clear that the accuracy has reached its limit  
247 because it is not unusual for  $^{13}\text{C}$  shifts to differ by more than 1 ppm in different solvents.

248 The distinct advantages of our models are their simpler architectures<sup>†</sup> and fewer trainable  
249 parameters, which result in significantly reduced training time. We do not consider the  
250 parameters of pre-trained models because the entire training dataset can be encoded by  
251 pre-trained models before training, making the training time independent of the number of  
252 parameters of the pre-trained model. However, the complexity of pre-trained models affects  
253 inference speed. Fortunately, the bottleneck in inference is conformer generation, so our  
254 models are faster to train and equally fast for inference.

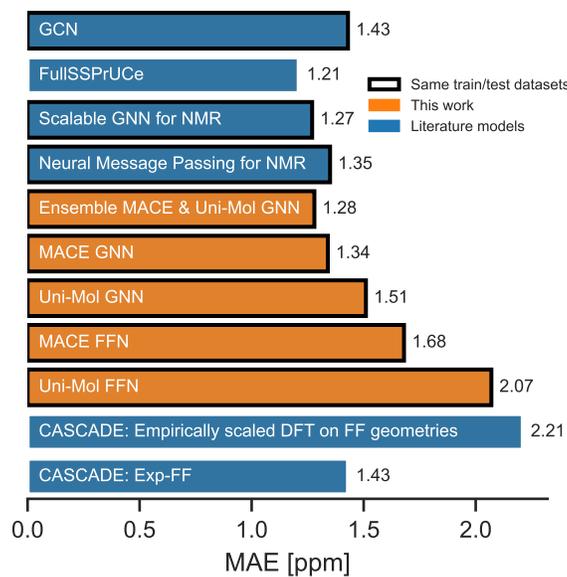


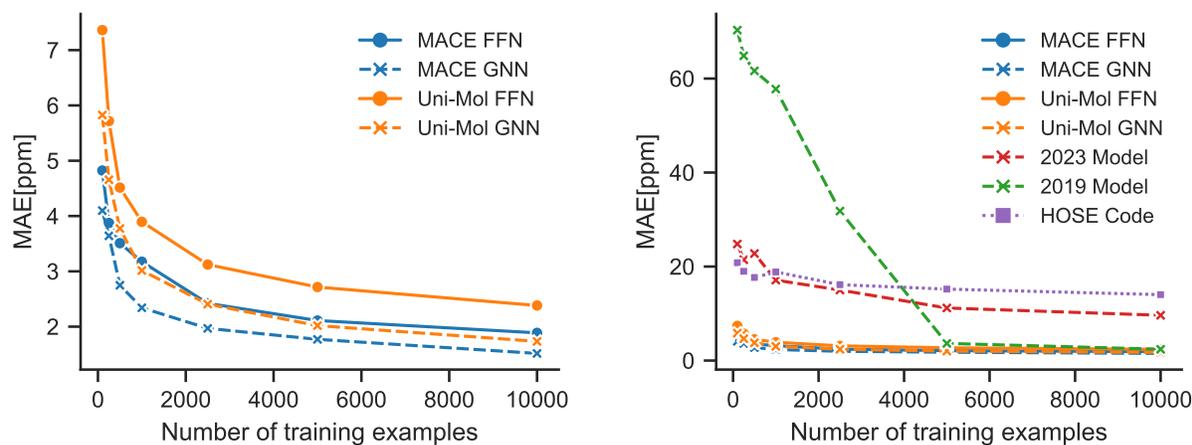
Figure 3: Comparison with the literature models.<sup>12,24–26,59</sup>

## 255 Low-data regimes

256 To simulate low-data regimes, we sampled data points from the training dataset, maintaining  
 257 the same model architectures<sup>†</sup> as used in the full data scenario to emphasize the effectiveness  
 258 of transfer learning. Nonetheless, the performance can be enhanced by optimizing hyperpa-  
 259 rameters for low-data regimes, especially by reducing model complexity and the dropout rate.  
 260 Furthermore, an additional molecule was excluded from the test set because MACE-based  
 261 models gave erroneous predictions for that molecule.<sup>†</sup>

262 Figure 4a illustrates that the performance of all models is improved with an increased  
 263 number of training examples. Notably, the MACE FFN model outperforms the Uni-Mol  
 264 GNN model in extremely low-data scenarios, whereas the reverse is true in high-data sce-  
 265 narios. The varying complexities of the models can explain this difference, as smaller models  
 266 need less training data. Figure 4b compares models in this paper with a model that per-  
 267 forms similarly on the full dataset, a model specifically designed for low-data scenarios, and  
 268 a classical HOSE Code model.<sup>43,44</sup> Transfer learning significantly boosts accuracy in low-

269 data scenarios compared to models trained from scratch. Furthermore, there is no trade-off  
270 between performance in high-data and low-data scenarios, unlike in the 2019 model.<sup>44</sup>



(a) This work.

(b) Comparison with literature models.<sup>43,44</sup>

Figure 4: Low-data regimes simulated using random sampling

## 271 Tautomer identification

272 In contrast to other outliers that possess uncommon functional groups or complex bonding  
273 and geometrical configurations,<sup>†</sup> one simple molecule yielded unsatisfactory results across all  
274 models developed in this study. Detailed examination reveals that the structure listed in  
275 the dataset, 1,3-cyclopentanedione, does not correspond to the tautomer present in solution  
276 under the conditions where the experimental chemical shifts were obtained. The tautomeric  
277 equilibrium that takes place for this molecule is illustrated in Figure 6.

278 Experimental findings on a similar compound<sup>60</sup> indicate that the two tautomers on the  
279 right-hand side of Fig. 6 predominate in solution, with rapid interconversion between them on  
280 the NMR time scale. Consequently, the NMR chemical shift of this compound represents an  
281 average of the chemical shifts of these two structures. The predicted shifts by the Ensemble  
282 MACE & Unimol GNN model for the diketo form (structure **a**) and the averaged prediction

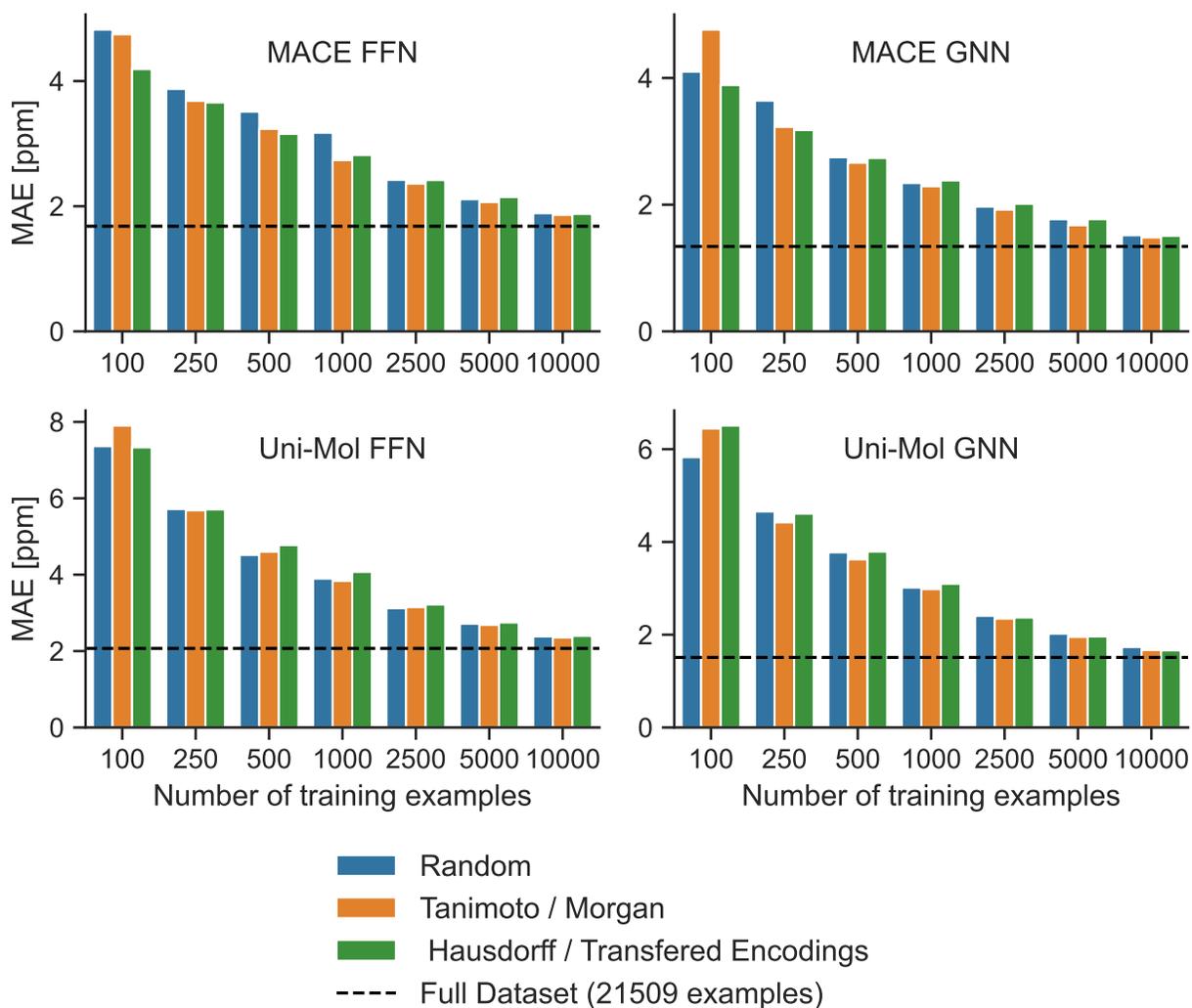


Figure 5: The effect of three different sampling strategies

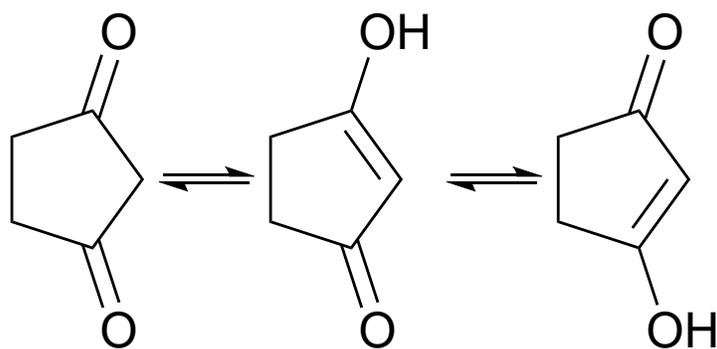


Figure 6: Equilibrium of different tautomers of 1,3-cyclopentanedione (**a**, **b** and **c**)

283 for the keto-enol forms (structures **b** and **c**) are illustrated in Figures 7a and 7b. The  
284 comparison of structure **a**, structure **b**, and the averaged prediction for structures **b** and **c**  
285 with observed shifts is shown in Table 2. The good match with experiment when using the  
286 prediction for the mixture of tautomers **b** and **c** is consistent with the rapid interconversion  
287 between two tautomeric structures, and demonstrates the ability of the model to assist in  
288 typical organic chemistry problems.

Table 2: Mean absolute errors of shifts predicted by the Ensemble GNN model

	Structure <b>a</b>	Structure <b>b</b>	Structures <b>b</b> and <b>c</b>
MAE [ppm]	19.03	3.42	0.34

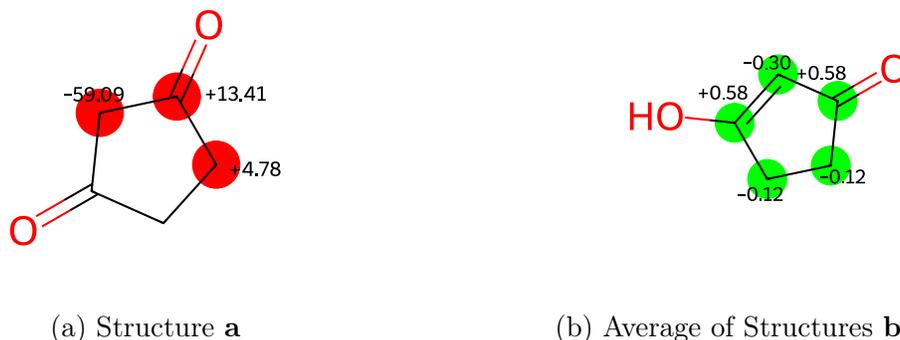


Figure 7: Errors [ppm] in predictions by Ensemble GNN model

## 289 Conclusion

290 We introduced atomic feature extraction as a transfer learning method applicable to both  
291 atomic and molecular-level prediction tasks. Unlike previous quasi-transfer methods, this  
292 approach does not require generating ab initio data for the target property. Moreover, the  
293 only information needed are atomic coordinates and atomic connectivity.

294 We evaluated this method on the prediction of experimental  $^{13}\text{C}$  chemical shifts, a well-  
295 studied atomic property prediction task. Our method performs on par with the best models

296 trained from scratch and surpasses them in low-data scenarios. When using this transfer  
297 learning approach, we demonstrated that the details of the sampling strategy used to se-  
298 lect from the pool of unlabeled examples don't matter. Lastly, we identified the MPNN  
299 forcefield as a superior candidate for pre-trained models for transfer learning compared to  
300 self-supervised pre-trained models.

301 The proven efficacy in low-data scenarios reveals new potential uses for this transfer  
302 learning approach in chemical problems with limited experimental data and in tasks where  
303 plenty of data exists but predictions are limited by data quality. For chemical shifts, em-  
304 ploying more precise geometries and data with recorded solvents and peaks assigned through  
305 multiple spectra will enhance the accuracy of data-driven models. This enhancement is fea-  
306 sible only if models can be trained on less data, which can be achieved through the transfer  
307 learning method described here.

## 308 **Methods**

### 309 **Data**

310 The dataset utilized in this work is taken from Kwon et al.<sup>26</sup>, and is derived from the original  
311 dataset published by Jonas and Kun.<sup>59</sup> It includes a predefined train/test split. This dataset  
312 comprises molecules with experimental spectra from nmrshiftdb2, which contain elements  
313 H, C, O, N, P, S, and F, and have no more than 64 atoms. The molecular geometries  
314 are obtained as the lowest energy conformers found in EDTKG conformer search<sup>61</sup> followed  
315 by MMFF minimization<sup>62</sup>. Molecules that failed rdkit sanitization, likely due to version  
316 discrepancies, were excluded. A detailed summary of the resulting dataset is available in the  
317 supplementary information.<sup>†</sup>

## 318 **Models**

319 FFN models consist of simple fully connected layers with exponential linear unit (ELU)  
320 activation functions.<sup>63</sup> The final layer is linear without any activation function. GNN models  
321 employ GraphSAGE message passing layers with ELU activation function, followed by a  
322 readout feedforward network of the same type as FFN models. Dropout was applied after  
323 each layer in all models.<sup>64</sup> The models were trained using L1 loss (mean absolute error) as the  
324 cost function and the AdamW optimizer with a weight decay of 0.01.<sup>65</sup> Hyperparameters were  
325 optimized through automated hyperparameter tuning and manual adjustments. Additional  
326 training and model architecture details can be found in the SI.<sup>†</sup>

## 327 **Computational details**

328 We accessed the pre-trained models using code from the associated repositories. Rdkit<sup>66,67</sup>  
329 (version 2023.09.5) was employed to process data, extract atomic connectivity from molec-  
330 ular structures, and perform MaxMin sampling. PyTorch<sup>68</sup> (version 2.2.1) and PyTorch  
331 Lightning<sup>69</sup> (version 2.2.1) were used for constructing and training FFN models, while Py-  
332 Torch Geometric<sup>70</sup> (version 2.5.2) was used for GNN models. All models were trained on a  
333 single Nvidia L4 Tensor core GPU. MaxMin sampling and Morgan fingerprints with a radius  
334 of 3 were implemented using rdkit. The Hausdorff distance was calculated using the scipy  
335 package<sup>71,72</sup>. Training for low-data examples continued until the validation loss ceased to  
336 decrease or until 800 epochs were reached. We sampled 120% of training data points for each  
337 regime, then randomly divided the data into train and validation sets. This ensured that the  
338 validation dataset size was always 20% of the training dataset size, and the train/validation  
339 split was performed as usual, making the conditions closer to a real low-data regime. Con-  
340 versely, testing was conducted on the entire test set for a realistic performance evaluation.  
341 Note that this approach differs from the work we compared low-data performance to, where

342 the test set size was proportional to the training dataset size.

## 343 **Code and Data availability**

344 The code used in the paper is publicly available in the repository

345 <https://github.com/zarkoivkovicc/AFE-TL-for-13C-NMR-chemical-shifts> under the ASL li-

346 cense, including the transfer learned models' weights. Pre-trained models and original

347 datasets can be downloaded from the code repositories of the corresponding publications.

## 348 **Author contributions**

349 Ž.I.: conceptualization, investigation, methodology, software, visualization, writing - original

350 draft J.J.: funding acquisition, supervision, writing - review and editing J.H.: resources,

351 supervision, writing - review and editing.

## 352 **Conflicts of interest**

353 There are no conflicts to declare.

## 354 **Acknowledgements**

355 Ž.I. acknowledge the IQTC-UB Master grant. J.J. acknowledges Maria de Maeztu grant

356 (code: CEX2021-001202-M).

## References

- [1] J. Stothers, *Carbon-13 NMR Spectroscopy: Organic Chemistry, A Series of Monographs, Volume 24*, Elsevier, 2012, vol. 24.
- [2] U. Sternberg, R. Witter and A. S. Ulrich, *Annual Reports on NMR Spectroscopy*, Academic Press, 2004, vol. 52, pp. 53–104.
- [3] A. Bagno and G. Saielli, *Theoretical Chemistry Accounts*, 2007, **117**, 603–619.
- [4] A. Wu, Q. Ye, X. Zhuang, Q. Chen, J. Zhang, J. Wu and X. Xu, *Precision Chemistry*, 2023, **1**, 57–68.
- [5] Z. Huang, M. S. Chen, C. P. Woroch, T. E. Markland and M. W. Kanan, *Chem. Sci.*, 2021, **12**, 15329–15338.
- [6] T. D. Michels, M. S. Dowling and C. D. Vanderwal, *Angewandte Chemie International Edition*, 2012, **51**, 7572–7576.
- [7] M. DiBello, A. R. Healy, H. Nikolayevskiy, Z. Xu and S. B. Herzon, *Acc. Chem. Res.*, 2023, **56**, 1656–1668.
- [8] S. D. Rychnovsky, *Org. Lett.*, 2006, **8**, 2895–2898.
- [9] H. A. Sánchez-Martínez, J. A. Morán-Pinzón, E. del Olmo Fernández, D. L. Eguiluz, J. F. Adserias Vistué, J. L. López-Pérez and E. G. De León, *J. Nat. Prod.*, 2023, **86**, 2294–2303.
- [10] D. J. Tantillo, *Nat. Prod. Rep.*, 2013, **30**, 1079–1086.
- [11] C. P. Gordon, C. Raynaud, R. A. Andersen, C. Copéret and O. Eisenstein, *Acc. Chem. Res.*, 2019, **52**, 2278–2289.

- 378 [12] Y. Guan, S. V. Shree Sowndarya, L. C. Gallegos, P. C. St. John and R. S. Paton, *Chem.*  
379 *Sci.*, 2021, **12**, 12012–12026.
- 380 [13] R. P. Verma and C. Hansch, *Chem. Rev.*, 2011, **111**, 2865–2899.
- 381 [14] E. Jonas, S. Kuhn and N. Schlörer, *Magnetic Resonance in Chemistry*, 2022, **60**, 1021–  
382 1031.
- 383 [15] I. Cortés, C. Cuadrado, A. Hernández Daranas and A. M. Sarotti, *Front. Nat. Prod.*,  
384 2023, **2**, 1122426.
- 385 [16] M. W. Lodewyk, M. R. Siebert and D. J. Tantillo, *Chem. Rev.*, 2012, **112**, 1839–1862.
- 386 [17] E. Benassi, *Journal of Computational Chemistry*, 2017, **38**, 87–92.
- 387 [18] P. Cimino, L. Gomez-Paloma, D. Duca, R. Riccio and G. Bifulco, *Magnetic Resonance*  
388 *in Chemistry*, 2004, **42**, S26–S33.
- 389 [19] Y. Yesiltepe, N. Govind, T. O. Metz and R. S. Renslow, *Journal of Cheminformatics*,  
390 2022, **14**, 64.
- 391 [20] T. Cheshire, P. Ramblenm, D. J. Tantillo, M. R. Siebert and M. W. Lodewyk, *CHEmical*  
392 *SHift REpository with Coupling Constants Added Too*, <http://cheshirenmr.info/>.
- 393 [21] P. Gao, J. Zhang, Q. Peng, J. Zhang and V.-A. Glezakou, *J. Chem. Inf. Model.*, 2020,  
394 **60**, 3746–3754.
- 395 [22] A. M. Sarotti and S. C. Pellegrinet, *J. Org. Chem.*, 2009, **74**, 7254–7260.
- 396 [23] D. Xin, C. A. Sader, O. Chaudhary, P.-J. Jones, K. Wagner, C. S. Tautermann, Z. Yang,  
397 C. A. Busacca, R. A. Saraceno, K. R. Fandrick, N. C. Gonnella, K. Horspool, G. Hansen  
398 and C. H. Senanayake, *J. Org. Chem.*, 2017, **82**, 5135–5145.

- 399 [24] J. Williams and E. Jonas, *Chem. Sci.*, 2023, **14**, 10902–10913.
- 400 [25] J. Han, H. Kang, S. Kang, Y. Kwon, D. Lee and Y.-S. Choi, *Phys. Chem. Chem. Phys.*,  
401 2022, **24**, 26870–26878.
- 402 [26] Y. Kwon, D. Lee, Y.-S. Choi, M. Kang and S. Kang, *J. Chem. Inf. Model.*, 2020, **60**,  
403 2024–2030.
- 404 [27] L. Budach, M. Feuerpfeil, N. Ihde, A. Nathansen, N. Noack, H. Patzlaff, F. Naumann  
405 and H. Harmouch, *The Effects of Data Quality on Machine Learning Performance*,  
406 2022.
- 407 [28] F. J. Fan and Y. Shi, *Bioorganic & Medicinal Chemistry*, 2022, **72**, 117003.
- 408 [29] S. Kuhn and N. E. Schlörer, *Magnetic Resonance in Chemistry*, 2015, **53**, 582–589.
- 409 [30] S. Kuhn, H. Kolshorn, C. Steinbeck and N. Schlörer, *Magnetic Resonance in Chemistry*,  
410 2024, **62**, 74–83.
- 411 [31] A. Farahani, B. Pourshojae, K. Rasheed and H. R. Arabnia, *A Concise Review of*  
412 *Transfer Learning*, 2021.
- 413 [32] G. Kumar and P. K. Bhatia, 2014 Fourth International Conference on Advanced Com-  
414 puting & Communication Technologies, Rohtak, India, 2014, pp. 5–12.
- 415 [33] E. d. S. Puls, M. V. Todescato and J. L. Carbonera, *An Evaluation of Pre-Trained*  
416 *Models for Feature Extraction in Image Classification*, 2023.
- 417 [34] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Nee-  
418 lakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger,  
419 T. Henighan, R. Child, A. Ramesh, D. M. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen,

- 420 E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford,  
421 I. Sutskever and D. Amodei, *Language Models Are Few-Shot Learners*, 2020.
- 422 [35] B. Weng, *Navigating the Landscape of Large Language Models: A Comprehensive Review  
423 and Analysis of Paradigms and Fine-Tuning Strategies*, 2024.
- 424 [36] D. van Tilborg, H. Brinkmann, E. Criscuolo, L. Rossen, R. Özçelik and F. Grisoni,  
425 *Current Opinion in Structural Biology*, 2024, **86**, 102818.
- 426 [37] S. G. Espley, E. H. E. Farrar, D. Buttar, S. Tomasi and M. N. Grayson, *Digital Discov-  
427 ery*, 2023, **2**, 941–951.
- 428 [38] G. E. Karniadakis, I. G. Kevrekidis, L. Lu, P. Perdikaris, S. Wang and L. Yang, *Nature  
429 Reviews Physics*, 2021, **3**, 422–440.
- 430 [39] H. Han and S. Choi, *J. Phys. Chem. Lett.*, 2021, **12**, 3662–3668.
- 431 [40] F. H. Vermeire and W. H. Green, *Chemical Engineering Journal*, 2021, **418**, 129307.
- 432 [41] A. El-Samman, S. De Castro, B. Morton and S. De Baerdemacker, *Can. J. Chem.*, 2024,  
433 **102**, 275–288.
- 434 [42] K. Schütt, P.-J. Kindermans, H. E. Saucedo Felix, S. Chmiela, A. Tkatchenko and K.-R.  
435 Müller, *Advances in Neural Information Processing Systems*, 2017.
- 436 [43] W. Bremser, *Analytica Chimica Acta*, 1978, **103**, 355–365.
- 437 [44] H. Rull, M. Fischer and S. Kuhn, *Journal of Cheminformatics*, 2023, **15**, 114.
- 438 [45] W. Ahmad, E. Simon, S. Chithrananda, G. Grand and B. Ramsundar, *ChemBERTa-2:  
439 Towards Chemical Foundation Models*, 2022.

- 440 [46] J. Ross, B. Belgodere, V. Chenthamarakshan, I. Padhi, Y. Mroueh and P. Das, *Large-*  
441 *Scale Chemical Language Representations Capture Molecular Structure and Properties,*  
442 2022.
- 443 [47] J. Xia, Y. Zhu, Y. Du and S. Z. Li, *A Systematic Survey of Chemical Pre-trained Models,*  
444 2022.
- 445 [48] Y. Rong, Y. Bian, T. Xu, W. Xie, Y. Wei, W. Huang and J. Huang, *Self-Supervised*  
446 *Graph Transformer on Large-Scale Molecular Data,* 2020.
- 447 [49] E. Heid, K. P. Greenman, Y. Chung, S.-C. Li, D. E. Graff, F. H. Vermeire, H. Wu,  
448 W. H. Green and C. J. McGill, *J. Chem. Inf. Model.*, 2024, **64**, 9–17.
- 449 [50] J. Gilmer, S. S. Schoenholz, P. F. Riley, O. Vinyals and G. E. Dahl, *Neural Message*  
450 *Passing for Quantum Chemistry,* 2017.
- 451 [51] D. P. Kovács, J. H. Moore, N. J. Browning, I. Batatia, J. T. Horton, V. Kapil, W. C.  
452 Witt, I.-B. Magdău, D. J. Cole and G. Csányi, *MACE-OFF23: Transferable Machine*  
453 *Learning Force Fields for Organic Molecules,* 2023.
- 454 [52] I. Batatia, S. Batzner, D. P. Kovács, A. Musaelian, G. N. C. Simm, R. Drautz, C. Ort-  
455 ner, B. Kozinsky and G. Csányi, *The Design Space of  $E(3)$ -Equivariant Atom-Centered*  
456 *Interatomic Potentials,* 2022.
- 457 [53] G. Zhou, Z. Gao, Q. Ding, H. Zheng, H. Xu, Z. Wei, L. Zhang and G. Ke, *Uni-Mol: A*  
458 *Universal 3D Molecular Representation Learning Framework,* 2023.
- 459 [54] W. L. Hamilton, R. Ying and J. Leskovec, *Inductive Representation Learning on Large*  
460 *Graphs,* 2017.

- 461 [55] M. Ashton, J. Barnard, F. Casset, M. Charlton, G. Downs, D. Gorse, J. Holliday,  
462 R. Lahana and P. Willett, *Quantitative Structure-Activity Relationships*, 2002, **21**, 598–  
463 604.
- 464 [56] D. Bajusz, A. Rácz and K. Héberger, *Journal of Cheminformatics*, 2015, **7**, 20.
- 465 [57] D. Rogers and M. Hahn, *J. Chem. Inf. Model.*, 2010, **50**, 742–754.
- 466 [58] T. Birsan and D. Tiba, *System Modeling and Optimization*, Kluwer Academic Publish-  
467 ers, Boston, 2006, vol. 199, pp. 35–39.
- 468 [59] E. Jonas and S. Kuhn, *Journal of Cheminformatics*, 2019, **11**, 50.
- 469 [60] V. Lacerda, M. G. Constantino, G. V. J. da Silva, Á. C. Neto and C. F. Tormena,  
470 *Journal of Molecular Structure*, 2007, **828**, 54–58.
- 471 [61] S. Riniker and G. A. Landrum, *J. Chem. Inf. Model.*, 2015, **55**, 2562–2574.
- 472 [62] T. A. Halgren, *Journal of Computational Chemistry*, 1996, **17**, 490–519.
- 473 [63] D.-A. Clevert, T. Unterthiner and S. Hochreiter, *Fast and Accurate Deep Network Learn-*  
474 *ing by Exponential Linear Units (ELUs)*, 2015.
- 475 [64] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever and R. Salakhutdinov, *Journal of*  
476 *Machine Learning Research*, 2014, **15**, 1929–1958.
- 477 [65] I. Loshchilov and F. Hutter, *Decoupled Weight Decay Regularization*, 2017.
- 478 [66] *RDKit*.
- 479 [67] G. Landrum, P. Tosco, B. Kelley, Ric, D. Cosgrove, sriniker, R. Vianello, gedec, Nadi-  
480 neSchneider, G. Jones, E. Kawashima, D. N, A. Dalke, B. Cole, M. Swain, S. Turk,

481 A. Savelev, A. Vaucher, M. Wójcikowski, I. Take, V. F. Scalfani, D. Probst, K. Uji-  
482 hara, guillaume godin, A. Pahl, R. Walker, J. Lehtivarjo, F. Berenger, strets123 and  
483 jasondbiggs, *Rdkit/Rdkit: Release\_2023.09.5*, Zenodo, 2024.

484 [68] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin,  
485 N. Gimeshein, L. Antiga, A. Desmaison, A. Köpf, E. Yang, Z. DeVito, M. Raison,  
486 A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai and S. Chintala, *PyTorch: An  
487 Imperative Style, High-Performance Deep Learning Library*, 2019.

488 [69] W. Falcon and The PyTorch Lightning team, *PyTorch Lightning*, 2019.

489 [70] M. Fey and J. E. Lenssen, *Fast Graph Representation Learning with PyTorch Geometric*,  
490 2019.

491 [71] P. Virtanen, R. Gommers, T. E. Oliphant, M. Haberland, T. Reddy, D. Cournapeau,  
492 E. Burovski, P. Peterson, W. Weckesser, J. Bright, S. J. Van Der Walt, M. Brett,  
493 J. Wilson, K. J. Millman, N. Mayorov, A. R. J. Nelson, E. Jones, R. Kern, E. Larson,  
494 C. J. Carey, Í. Polat, Y. Feng, E. W. Moore, J. VanderPlas, D. Laxalde, J. Perktold,  
495 R. Cimrman, I. Henriksen, E. A. Quintero, C. R. Harris, A. M. Archibald, A. H. Ribeiro,  
496 F. Pedregosa, P. Van Mulbregt, SciPy 1.0 Contributors, A. Vijaykumar, A. P. Bardelli,  
497 A. Rothberg, A. Hilboll, A. Kloeckner, A. Scopatz, A. Lee, A. Rokem, C. N. Woods,  
498 C. Fulton, C. Masson, C. Häggström, C. Fitzgerald, D. A. Nicholson, D. R. Hagen,  
499 D. V. Pasechnik, E. Olivetti, E. Martin, E. Wieser, F. Silva, F. Lenders, F. Wilhelm,  
500 G. Young, G. A. Price, G.-L. Ingold, G. E. Allen, G. R. Lee, H. Audren, I. Probst, J. P.  
501 Dietrich, J. Silterra, J. T. Webber, J. Slavič, J. Nothman, J. Buchner, J. Kulick, J. L.  
502 Schönberger, J. V. De Miranda Cardoso, J. Reimer, J. Harrington, J. L. C. Rodríguez,  
503 J. Nunez-Iglesias, J. Kuczynski, K. Tritz, M. Thoma, M. Newville, M. Kümmerer,  
504 M. Bolingbroke, M. Tartre, M. Pak, N. J. Smith, N. Nowaczyk, N. Shebanov, O. Pavlyk,

505 P. A. Brodtkorb, P. Lee, R. T. McGibbon, R. Feldbauer, S. Lewis, S. Tygier, S. Sievert,  
506 S. Vigna, S. Peterson, S. More, T. Pudlik, T. Oshima, T. J. Pingel, T. P. Robitaille,  
507 T. Spura, T. R. Jones, T. Cera, T. Leslie, T. Zito, T. Krauss, U. Upadhyay, Y. O.  
508 Halchenko and Y. Vázquez-Baeza, *Nat Methods*, 2020, **17**, 261–272.

509 [72] A. A. Taha and A. Hanbury, *IEEE Trans. Pattern Anal. Mach. Intell.*, 2015, **37**, 2153–  
510 2163.