

Dynamic Embedding Representation for Graph Neural Networks to Enhance Materials Property Prediction with Limited Datasets

Vishu Gupta^{1,2,3}, Kamal Choudhary⁴, Youjia Li¹, Muhammed Nur Talha Kilic⁵, Daniel Wines⁴, Wei-keng Liao¹, Alok Choudhary¹, Ankit Agrawal¹

¹*Department of Electrical and Computer Engineering, Northwestern University, Evanston, IL, U.S.A.*

²*Lewis-Sigler Institute for Integrative Genomics, Princeton University, Princeton, NJ, U.S.A.*

³*Ludwig Institute for Cancer Research, Princeton University, Princeton, NJ, U.S.A.*

⁴*Material Measurement Laboratory, National Institute of Standards and Technology, Gaithersburg, MD, U.S.A.*

⁵*Department of Computer Science, Northwestern University, Evanston, IL, U.S.A.*

1 Graph neural networks (GNNs) have proven effective in understanding and predicting
2 diverse material properties, even when working with limited datasets. An important step in
3 training GNN is to use an appropriate and informative graph embedding that can adequately
4 represent the structural and compositional information in the chemical space. Current graph
5 embeddings consist of composition and structure-agnostic element-level encodings, which are
6 static in nature. This makes it challenging to differentiate between different compounds on
7 the element level, especially for datasets with limited data size, thereby relying more on
8 the complex input and architecture for model training. Here, we present a novel framework

9 for GNN-based prediction tasks that use dynamic embedding to significantly improve the
10 models’ predictive ability on materials properties with limited data size. We evaluated the
11 proposed framework on multiple materials datasets across various domains to find that the
12 model trained using dynamic embedding outperforms the models trained using conventional
13 static embedding and features obtained using a pre-trained model. The proposed framework
14 holds significant potential for expediting artificial intelligence (AI)-driven materials discovery.

15 **Introduction**

16 The field of materials science has experienced an increasing utilization of artificial intelligence
17 (AI) and machine learning (ML) methods, boosting the development of various data-driven
18 models for forward ¹⁻³ and inverse design problems. ⁴⁻⁶ Among these methodologies, graph
19 neural networks (GNNs), in particular, have become increasingly sought-after owing to their
20 ability to aid the discovery process of novel materials and molecules for various applica-
21 tions ⁷⁻¹⁰. In recent years, the development of fast and efficient GNN models for forward
22 predictive modeling of materials properties with desired properties has arguably received
23 the most interest for its potential to accelerate materials design ^{11,12}. Broadly speaking, the
24 quality of the predictive model trained using GNN depends on several factors, including the
25 capability of the model architecture to effectively extract relevant information from the input
26 data, the amount of data available for training the model, and the meaningful information
27 contained within the graph embedding used to represent the chemical compound. There has
28 been ongoing research to improve the model’s performance by taking one or more of these

29 factors into consideration ¹³⁻¹⁵.

30 Various works have attempted to improve the performance of the GNN model by
31 modifying the architectures used for predicting material properties such as Crystal Graph
32 Convolutional Neural Networks (CGCNN) ¹⁶, Representation Learning from Stoichiome-
33 try (Roost) ¹⁷, SchNet ¹⁸, MatErials Graph Network (MEGNet) ¹⁹, DimeNet++ ²⁰, and
34 Atomistic Line Graph Neural Network (ALIGNN) ²¹. Several works have also attempted
35 to improve the performance of the model when dealing with limited datasets²², either via
36 materials property specific feature engineering performed before training the model ²³⁻²⁷
37 or utilizing the knowledge learned from a trained model on a large dataset to boost the
38 performance of the small dataset using advanced data mining techniques such as feature
39 extraction ²⁸⁻³⁰ and architectural optimization ³¹⁻³³. Additionally, graph-based descriptors
40 and topological index methods have been widely developed for molecular systems, including
41 benzenoid hydrocarbons and molecular trees, with well-known indices such as the Zagreb
42 and Szeged indices. However, these approaches are primarily designed for specific classes of
43 organic molecules characterized by well-defined molecular graphs. Consequently, there is less
44 visibility of works that focus on developing informative input descriptors in the form of graph
45 embeddings for the GNN designed for solids with well-defined structures. Conventionally, a
46 crystalline material, when represented as a graph, comprises of nodes corresponding to con-
47 stituent atom features and edges corresponding to bond features. As the fundamental node
48 and edge information available to use are limited ³⁴, recent works have tried to improve the
49 model performance by incorporating more information (such as angle-based information ^{21,35})

50 or implementing complex components in the model architecture. Moreover, publicly used
51 and available graph embeddings often use atom-level encodings that are composition and
52 structure agnostic. This means that the embeddings represent individual atoms (such as Co
53 or O) the same way, regardless of the molecule’s overall structure or complexity ^{36,37}. For
54 example, the atom-level embedding for cobalt (Co) in a simple molecule like Co_3O_4 (cobalt
55 oxide) would be identical to that in a more complex structure like LiCoO_2 (lithium cobalt
56 oxide). This uniform representation does not take into account how the atom’s bonding
57 environment or the molecule’s composition might affect its chemical properties, which can
58 limit the expressiveness of the embeddings. Additionally, these uniform representations can
59 be problematic, particularly when working with properties that depend heavily on the pre-
60 cise chemical environment, such as electronic, mechanical, or thermal properties. This issue
61 becomes more pronounced when large, diverse datasets are unavailable for model training,
62 as the models have fewer opportunities to effectively extract relevant information. Thus, the
63 problem is not just the structure-agnostic nature of the embeddings, but also the challenge
64 of training effective models in low-data regimes, where simple model training frameworks
65 may not be able to compensate for these limitations.

66 In this work, we present a novel framework that uses composition and structure aware
67 element-level encodings to represent a chemical compound in a graph neural network (GNN)
68 to improve the predictive performance of the model for materials properties with limited
69 data size. The workflow comparison of the traditional and proposed approach for training
70 graph neural networks is shown in Figure 1. Here, we first apply a GNN architecture that

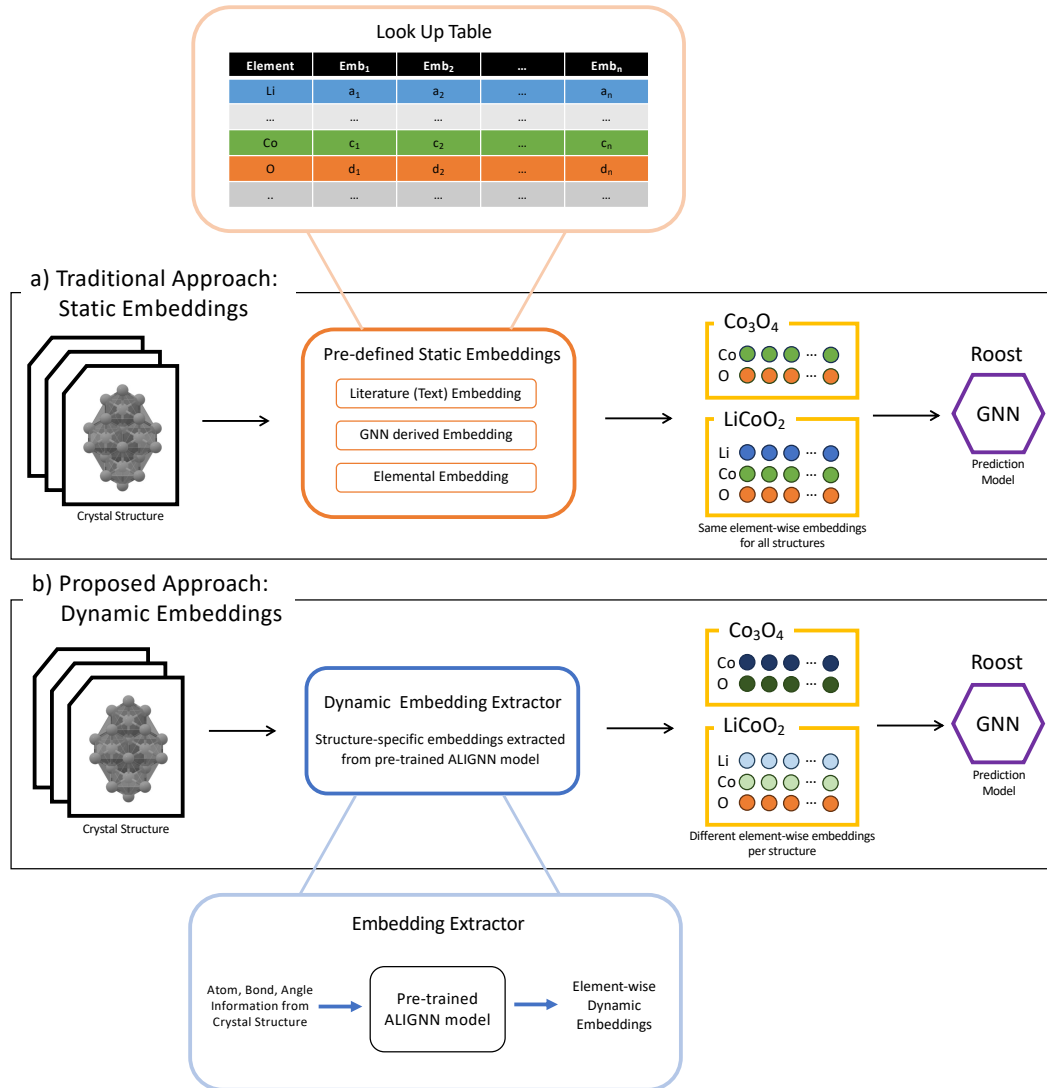


Figure 1: Comparison of traditional static embedding and proposed dynamic embedding approaches for materials property prediction. (a) Traditional approach: Pre-defined static embeddings are retrieved from pre-calculated look up table. Each element (e.g., Co, Li, O) has identical embedding vectors regardless of its structural environment, resulting in the same representations despite their different crystal structures. (b) Proposed approach: Dynamic embeddings are extracted from a pre-trained structure-aware ALIGNN model that processes atomic positions, bond lengths, and bond angles. Each element receives different embedding vectors depending on its local structural environment, enabling the same element (e.g., Co) to have distinct representations in different materials (Co₃O₄ vs. LiCoO₂).

71 incorporates composition and structure information as model input to capture the under-
72 lying chemistry associated with the existing large data containing crystal structures. The
73 resulting model is then used as an embedding extractor to extract dynamic graph embed-
74 dings to represent each input of the target dataset. We use an ALIGNN²¹ model as the
75 embedding extractor, as it takes atom, bond, and angle-based information as the model
76 input and has been shown to learn meaningful information from large materials datasets³⁸.
77 The ALIGNN-based embedding extractor consists of one embedding layer, followed by four
78 ALIGNN layers and four GCN layers, all connected in a linear sequence. This structure re-
79 sults in a total of nine embedding layers. While the ALIGNN²¹ work focuses on learning the
80 chemistry between elements in compounds using conventional CGCNN embeddings (static),
81 our method leverages this learned chemistry from ALIGNN to extract more meaningful rep-
82 resentations (dynamic). These representations are expected to better capture interactions
83 and be more effective for downstream tasks. For training the target model, we use Roost¹⁷
84 as it takes element-level graph embedding as the model input and has been shown to out-
85 perform various deep learning (DL) models for training materials properties with limited
86 data points. We compare models trained using the proposed framework with models trained
87 using conventional static embedding and methodologies that use embedding and/or features
88 obtained using a pre-trained model. As, the performance difference between the models
89 trained on the GNN and the conventional DL architectures are already explored in³⁸, here
90 we mainly compare the results with models trained on the conventional static embedding
91 and methodologies that use embedding and/or features obtained using a pre-trained model.

92 We observe that the proposed framework can readily accommodate the continuously ex-
93 panding datasets and evolving model training techniques, facilitating further enhancements
94 to the models. These improvements are anticipated to aid materials science researchers in
95 effectively leveraging data mining techniques³⁹⁻⁴¹. This can enable more reliable and accu-
96 rate screening and identification of potential material candidates, thereby accelerating the
97 materials discovery process.

98 **Results**

99 **Datasets** We use two datasets of density functional theory (DFT)-computed properties in
100 this work: Materials Project (MP)⁴² and Joint Automated Repository for Various Integrated
101 Simulations (JARVIS)^{43,44}. MP dataset was downloaded from²¹ and the JARVIS dataset
102 from the following figshare link https://figshare.com/collections/ALIGNN_data/5429274.

103 The MP dataset with Formation Energy²¹ is used to train embedding extractor model,
104 which is then used to extract dynamic embedding and train small target datasets, which
105 correspond to materials properties with limited data points, to improve their predictive
106 performance. Formation energy is selected as the materials property as it has shown to learn
107 to meaningful representations from large source datasets^{28,38}. Materials properties in the
108 JARVIS datasets are used to train target model. Once trained, the target model is used to
109 make predictions and evaluate material properties for compounds in holdout test set. We use
110 dataset without structural polymorphs for a few of the analyses by keeping the most stable

111 structure available in the database, i.e., data entry corresponding to the lowest formation
112 energy among all compounds with the same composition. The target datasets are randomly
113 split into training, validation, and holdout test sets in the ratio of 80:10:10 with a fixed
114 random seed. The data size corresponding to each materials property in the two datasets
115 are shown in Supplementary Table 1, modifications applied to certain materials properties
116 within the target dataset to align with the model input are shown in Supplementary Table
117 2, and data distribution for training and test sets for all properties of JARVIS dataset are
118 shown in Supplementary Figure 1 (a) to (ar). We use mean absolute error (MAE) as the
119 primary evaluation metric for all models. Due to the extensive range of materials properties
120 examined in this study and the constraints posed by limited computational resources, we
121 did not explore the aleatoric uncertainty resulting from the random initialization of the
122 models. Additionally, since the performance differences between models trained directly
123 on pre-trained GNNs and those trained via transfer learning—using either fine-tuning or
124 extracted features—have already been explored in prior work^{38,45,46}, our focus is primarily on
125 comparing downstream models. Specifically, we evaluate models trained using conventional
126 static embeddings against those utilizing embeddings or features derived from pre-trained
127 models.

128 **Dynamic Embedding Extractor** We use a structure-aware GNN-based architecture,
129 ALIGNN²¹ as the base architecture to train the embedding extractor and extract dynamic
130 embedding as it takes atom, bond, and angle based information as the model input and has
131 been shown to learn meaningful information from the large materials datasets using crystal

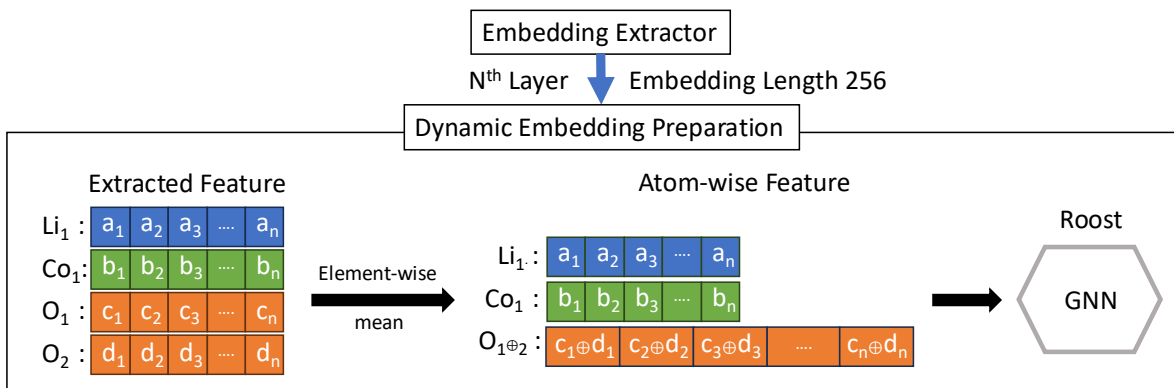


Figure 2: Outline of the operation performed to create dynamic embedding from ALIGNN-based embedding extractor. The embedding comprises of 256 vector dynamic representation of elements. We take element-wise mean (depicted with \oplus) of embeddings for the compound with multiple atoms for the same element. The final set of embeddings are used as an input to the graph neural network.

132 structure information as the model input ³⁸. For the initial set of input features used to
 133 train ALIGNN, please refer to the publication ²¹. To extract embeddings from ALIGNN, we
 134 design a novel ALIGNN-based Embedding Extractor, shown in Figure 2.

135 The structure file containing information on the atomic positions and lattice geometry
 136 of a chemical compound is divided into atom, bond, and angle based encodings before being
 137 fed into an ALIGNN-based embedding extractor, where we perform embedding extraction.
 138 The ALIGNN-based embedding extractor consists of one embedding layer, followed by four
 139 ALIGNN layers and four GCN layers, all connected in a linear sequence. This structure
 140 results in a total of nine embedding layers. We extract the element-level encoding before
 141 and after the ALIGNN and GCN (Graph Convolution Network) layers in the ALIGNN ar-

142 chitecture. We call these encodings dynamic embeddings as they vary depending on the
143 layer they are extracted from and for each element in every compound. After performing
144 embedding extraction, we obtain nine sets of dynamic embeddings, each with a unique 256-
145 vector atom-level representation of the elements of the chemical compound. Due to the
146 nature of the model architecture, all the embeddings extracted from the embedding extrac-
147 tor contain atom, bond, and angle based structure information. For a detailed explanation
148 of the pre-processing of the structure-based encoding associated with the embedding ex-
149 tractor, please refer to the methods section. Note that the concept of dynamic embedding
150 representations—where features extracted from a neural network are used for downstream
151 tasks—has been explored in recent years, existing methods predominantly employ frame-
152 works that extract one-dimensional embeddings which can only be used with feed-forward
153 neural networks ^{28,47} to train models for the downstream tasks. Our work distinguishes
154 itself by introducing a framework that extracts atom-wise dynamic embedding representa-
155 tions, which are then utilized as inputs for a GNN which requires a more complex form of
156 input for model training. To the best of our knowledge, no previous work has employed this
157 approach, making our framework novel in its application of atom-level embeddings to GNN
158 architectures for enhanced structural representation and improved predictive performance
159 on the downstream task. We show the difference between the traditional framework and
160 proposed framework for feature extraction in Figure 3.

161 Next, we perform model training using the above-defined set of embeddings as input
162 for the deep graph neural network where we use Roost ¹⁷ as the base architecture as it takes

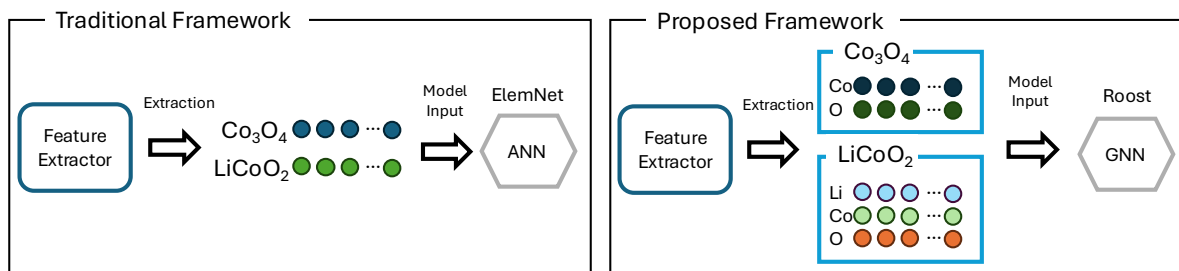


Figure 3: Workflow comparison of the traditional and proposed framework for feature extraction. Traditional methods predominantly employ frameworks with one-dimensional embeddings as the final output which can only be used with artificial neural networks (ANN) to train models for the downstream tasks. Our proposed framework extracts atom-wise dynamic embedding representations, which are then utilized as inputs for a GNN which requires a more complex form of input for model training.

163 element-level graph embedding as input for model training. We use the formation energy of
 164 the JARVIS dataset as the materials property for the property prediction task.

165 Table 1 shows that, in general, dynamic embeddings obtained from the last layer of
 166 the GNN architecture perform better than other layers. Hence, for the rest of the analysis,
 167 we only use the dynamic embedding extracted from the last layer to perform model training
 168 on the target datasets.

169 **Static Embedding vs. Dynamic Embedding** Here, we compare the performance of the
 170 model trained on dynamic embeddings against conventional static embeddings using multiple
 171 target properties in the JARVIS dataset. For static embeddings, we use OneHot, Magppie⁴⁸,
 172 CGCNN¹⁶, MEGNet¹⁹, Mat2Vec⁴⁹ and Matscholar⁵⁰ representations to train the target

Table 1: The table shows the test mean absolute error (MAE) when run on embeddings extracted from different layers of embedding extractor on formation energy of JARVIS dataset.

Layer Number	Validation MAE (eV/atom)	Test MAE (eV/atom)
1	0.107	0.108
2	0.068	0.068
3	0.061	0.059
4	0.057	0.055
5	0.051	0.049
6	0.047	0.046
7	0.046	0.045
8	0.047	0.045
9	0.045	0.043

173 model which are shown to be meaningful high-dimensional representations of chemical ele-
 174 ments³⁴. OneHot and Magpie⁴⁸ embeddings only take composition-based information into
 175 consideration. CGCNN¹⁶ and MEGNet¹⁹ embeddings are derived from crystal graph convo-
 176 lutional neural networks trained to predict materials properties using structure information
 177 contained in crystalline materials as input. Mat2Vec⁴⁹ and Matscholar⁵⁰ are literature word
 178 embeddings obtained by using over 3.27 million materials science related abstracts. While
 179 OneHot and Magpie utilize domain-driven feature engineering, CGCNN, MEGNet, Mat2Vec,
 180 and Matscholar use atom-level embeddings obtained from pre-trained models, representing a
 181 form of transfer learning. Note that all models use Roost as the base architecture for training
 182 with the respective embeddings (i.e., they correspond to OneHot+Roost, CGCNN+Roost,
 183 etc). Interested readers can refer to their respective publications for more details. Table 2
 184 presents the prediction accuracy for 44 target materials properties trained on different sets
 185 of embeddings.

186 Table 2 indicates that models trained using dynamic embedding significantly outper-

Table 2: The table shows the test MAE, and % error change between proposed and best performing static embedding for each of the target materials properties trained on different sets of embeddings using Roost for the prediction task of “Static Embedding v/s Dynamic Embedding”.

Property	Data Size	Onehot	Magpie [39]	CGCNN [16]	MEGNET [19]	Mat2Vec [40]	Matscholar [41]	Dynamic (proposed)	% Error Change
E _g OPT (eV)	39760	0.227	0.223	0.239	0.224	0.229	0.219	0.172	-21.5
E _f (eV/atom)	39760	0.090	0.086	0.092	0.093	0.087	0.086	0.043	-50.0
KLU (Å)	39578	12.19	11.86	12.25	12.24	12.08	11.65	10.47	-10.1
Ehull (eV/atom)	39487	0.051	0.052	0.062	0.054	0.051	0.052	0.097	90.2
Encut (eV)	39461	148.6	144.9	148.4	142.4	146.0	140.4	136.5	-2.8
Magoszi (μ_B)	37574	0.446	0.438	0.470	0.438	0.442	0.429	0.342	-20.3
Magout (μ_B)	37932	0.548	0.565	0.574	0.547	0.545	0.531	0.452	-14.9
Epsx	31227	29.34	28.68	29.28	28.94	28.72	27.60	22.85	-17.2
Epsy	31227	29.09	28.66	28.77	28.87	29.08	27.65	23.03	-16.7
Epsz	31227	28.10	27.75	27.87	28.47	28.18	27.02	22.26	-17.6
PPF ($\mu W m^{-1} K^{-2}$)	16443	568.8	525.0	545.5	527.6	553.9	521.9	470.2	-9.9
NPF ($\mu W m^{-1} K^{-2}$)	16443	570.9	552.8	543.1	543.4	558.4	534.8	496.8	-7.1
NSB ($\mu V/K$)	16387	57.31	57.28	56.75	54.92	55.49	53.91	42.53	-21.1
PSB ($\mu V/K$)	16400	64.81	63.04	63.38	60.88	64.21	59.04	49.01	-17.0
Nem300k (m_0)	15252	0.400	0.417	0.407	0.412	0.422	0.408	0.371	-7.3
Pem300k (m_0)	14364	0.686	0.656	0.645	0.624	0.671	0.617	0.475	-23.0
ETC33 (GPa)	13195	45.07	45.22	43.51	44.92	43.60	43.44	34.22	-21.2
ETC22 (GPa)	13098	43.59	41.78	42.71	40.95	40.86	41.34	32.79	-19.8
ETC11 (GPa)	13107	44.59	43.26	42.79	42.69	42.58	41.59	34.23	-17.7
BulkKV (GPa)	12789	18.60	18.21	18.46	19.22	17.46	17.31	13.74	-20.6
ETC13 (GPa)	13056	18.55	18.16	17.84	17.80	18.08	17.30	15.10	-12.7
ETC12 (GPa)	12895	22.11	20.85	21.91	21.61	20.67	20.62	17.61	-14.6
Poisson	12676	0.234	0.236	0.222	0.204	0.223	0.230	0.187	-8.3
ShearGV (GPa)	12179	13.75	12.82	12.86	12.87	13.00	12.77	9.89	-22.6
E _g MBJ (eV)	12281	0.410	0.379	0.401	0.369	0.371	0.345	0.285	-17.4
ETC44 (GPa)	11863	18.98	18.76	18.27	18.14	17.96	17.98	14.34	-20.2
AvgME (electron mass unit)	12637	0.100	0.103	0.098	0.111	0.099	0.097	0.083	-14.4
AvgMH (electron mass unit)	12637	0.164	0.162	0.170	0.161	0.176	0.159	0.130	-18.2
ETC55 (GPa)	11673	16.61	16.50	16.26	15.34	15.46	15.41	12.40	-19.2
ETC66 (GPa)	11538	17.14	16.51	16.49	16.36	16.48	16.31	12.79	-21.6
Mepsx	11380	32.79	32.26	31.59	30.94	32.00	30.62	27.36	-10.6
Mepsy	11380	32.65	30.93	31.25	30.61	31.14	31.50	26.62	-13.0
Mepsz	11380	30.59	29.59	30.25	31.30	30.37	28.91	25.02	-13.5
MaxM (cm^{-1})	9156	55.51	50.72	53.38	54.17	52.65	55.19	41.63	-17.9
MinM (cm^{-1})	8644	21.55	20.03	20.68	20.34	19.87	20.56	14.95	-24.8
MaxEFG ($10^{21} Vm^{-2}$)	7286	32.33	30.58	28.83	30.08	30.23	29.52	23.27	-19.3
Spillage	8755	0.415	0.398	0.407	0.394	0.402	0.406	0.343	-12.9
SLME (%)	5877	6.479	5.900	5.764	6.289	5.867	5.700	4.948	-13.2
PMEij (cm^{-2})	2711	0.134	0.116	0.133	0.118	0.129	0.127	0.112	-3.4
PMDi	2649	20.63	19.05	18.60	18.63	18.83	19.41	16.07	-13.6
PMDIel (ϵ_{11})	2592	2.959	3.167	3.410	3.184	2.866	2.948	2.631	-8.2
PMDIlo (ϵ_{11})	2496	3.864	3.516	3.399	3.841	3.434	3.465	3.037	-10.7
PMDij (cm^{-1})	1885	8.878	8.702	7.900	8.461	7.954	7.763	5.771	-25.7
Exfoli (meV/atom)	499	48.83	42.86	47.26	42.83	43.37	40.78	40.62	-0.4

187 form the models trained using conventional static embeddings in 43/44 cases, i.e., in $\approx 98\%$
188 of the cases. Although we observe improvement across all the target materials properties,
189 the most improvement (i.e., $\approx 50\%$) is shown with the formation energy as the target
190 property (0.086 to 0.043 eV/atom). The improvement can be attributed to the fact that
191 the embedding extractor was trained using the formation energy as the materials property.
192 The proposed model performed worse only for Ehull as the materials property. Ehull is
193 unique among our tested properties as it represents relative stability rather than an intrinsic
194 material property. It measures the energy difference between a compound and the thermo-
195 dynamically stable phase-separated state. This relative nature may explain why embeddings
196 trained on absolute formation energy do not transfer as effectively. It would be interesting to
197 explore why the proposed framework performed worse only for Ehull by analyzing its relation
198 with the embeddings in future work by training embedding extractors directly on convex-
199 hull-related objectives or by using multi-task learning frameworks that jointly optimize for
200 both absolute and relative thermodynamic properties. Overall, the results illustrate that the
201 embedding extractor can learn and extract useful and widely applicable embeddings during
202 model training on source data and improve the performance of the target model trained on
203 materials properties with limited data size.

204 **Comparison Against Other Methods** In this section, we investigate the performance
205 of the proposed framework against other well-known deep neural networks and methodolo-
206 gies that use embedding and/or features obtained using a pre-trained model to improve the
207 performance of the model for small datasets, i.e., CrabNet⁵¹, and AtomSets⁴⁵. CrabNet⁵¹

208 uses literature word embedding mat2vec ⁴⁹ as input representation with attention-based
209 network to train the model. AtomSets ⁴⁵ comprises compositional and structural descrip-
210 tors extracted from pre-trained MEGNet ¹⁹ model. These features serve as input to train
211 conventional DL models, thereby implementing transfer learning similar to the approach
212 proposed in this work. We also compare the proposed framework with conventional deep
213 learning (ElemNet ⁴⁷) and traditional machine learning (AutoML ⁵²) models. As they require
214 one-dimensional vector representation as input, we use elemental fraction (EF) ⁴⁷, physical
215 attributes (PA) ⁴⁸, and atom-wise averaged value of dynamic embedding as the inputs to
216 the model. As we have multiple trained models for AtomSets, ElemNet, and AutoML, we
217 use the model with the least test MAE to compare the performance in Table 3 and show the
218 rest of the results in Supplementary Tables 3, and 4.

219 Table 3 indicates that models trained using dynamic embedding outperform other
220 embedding and/or features based methods in 33/44 cases, i.e., in $\approx 75\%$ of the cases. For
221 the remaining 11 properties, the model trained on ElemNet using atom-wise averaged value
222 of dynamic embedding performed the best for 7/11 cases. It is interesting that the methods
223 that incorporate dynamic embedding in one form or another can outperform other methods
224 in 40/44 cases, i.e., in $\approx 91\%$ of the cases. We also observe close error values for some of
225 the materials properties between the proposed method and ElemNet trained using atom-
226 wise averaged value of dynamic embedding. However, as both the models are trained using
227 default hyperparameter settings, the proposed methods train for 250 epochs, and ElemNet
228 trains with patience of 100 epochs (i.e., stop the model training if the validation loss does not

Table 3: The table shows the test MAE, and % error change between proposed method and (1) best performing existing method and (2) method using embedding extracted from ALIGNN as input for each of the target materials properties trained on different models for the prediction task of “Comparison Against Other Methods”.

Property	Data Size	Existing		Embedding		Dynamic (proposed)	% Error Change (Existing)	% Error Change (Embedding)
		CrabNet [42]	AtomSets [36]	AutoML [43]	ElemNet [38]			
E _g OPT (eV)	39760	0.302	0.237	0.283	0.171	0.172	-27.4	0.6
E _f (eV/atom)	39760	0.108	0.066	0.075	0.043	0.042	-36.4	-2.3
KLU (Å)	39578	10.50	12.10	10.62	10.23	10.47	-0.3	2.3
Ehull (eV/atom)	39487	0.089	0.168	0.142	0.156	0.097	9.0	-31.7
Encut (eV)	39461	124.0	164.7	152.9	138.1	136.5	10.1	-1.2
Magoszi (μ_B)	37574	0.504	0.431	0.537	0.373	0.342	-20.6	-8.3
Magout (μ_B)	37932	0.641	0.614	0.665	0.466	0.452	-26.4	-3.0
Epsx	31227	25.10	31.32	25.98	23.96	22.85	-9.0	-4.6
Epsy	31227	25.40	31.69	26.90	23.61	23.03	-9.3	-2.5
Epsz	31227	24.60	31.74	25.84	22.98	22.26	-9.5	-3.1
PPF ($\mu W m^{-1} K^{-2}$)	16443	487.0	575.1	526.9	471.4	470.2	-3.4	-0.3
NPF ($\mu W m^{-1} K^{-2}$)	16443	520.0	601.9	541.4	484.7	496.8	-4.5	2.5
NSB ($\mu V/K$)	16387	50.50	57.70	50.28	45.67	42.53	-15.8	-6.9
PSB ($\mu V/K$)	16400	51.50	62.27	56.78	49.69	49.01	-4.8	-1.4
Nem300k (m_0)	15252	0.389	0.475	0.438	0.357	0.371	-4.6	3.9
Pem300k (m_0)	14364	0.658	0.642	0.633	0.501	0.475	-26.0	-5.2
ETC33 (GPa)	13195	38.60	47.62	40.82	35.74	34.22	-11.3	-4.3
ETC22 (GPa)	13098	36.40	44.58	40.05	34.79	32.79	-9.9	-5.7
ETC11 (GPa)	13107	37.70	45.67	40.67	35.07	34.23	-9.2	-2.4
BulkKV (GPa)	12789	14.50	19.50	17.83	14.15	13.74	-5.2	-2.9
ETC13 (GPa)	13056	15.30	20.41	16.87	15.63	15.10	-1.3	-3.4
ETC12 (GPa)	12895	18.80	23.53	19.72	18.02	17.61	-6.3	-2.3
Poisson	12676	0.168	0.234	0.181	0.173	0.187	11.3	8.1
ShearGV (GPa)	12179	12.30	14.60	12.39	10.27	9.89	-19.6	-3.7
E _g MBJ (eV)	12281	0.489	0.428	0.406	0.354	0.285	-33.4	-19.5
ETC44 (GPa)	11863	15.60	19.56	16.47	14.58	14.34	-8.1	-1.6
AvgME (electron mass unit)	12637	0.103	0.111	0.113	0.090	0.083	-19.4	-7.8
AvgMH (electron mass unit)	12637	0.171	0.178	0.183	0.142	0.130	-24.0	-8.5
ETC55 (GPa)	11673	14.20	16.61	15.25	12.50	12.40	-12.7	-0.8
ETC66 (GPa)	11538	14.60	17.20	14.61	13.21	12.79	-12.4	-3.2
Mepsx	11380	28.50	36.11	30.58	27.61	27.36	-4.0	-0.9
Mepsy	11380	27.40	35.07	29.69	26.68	26.62	-2.8	-0.2
Mepsz	11380	27.90	34.92	29.28	25.88	25.02	-10.3	-3.3
MaxM (cm^{-1})	9156	60.60	62.70	55.58	57.47	41.63	-31.3	-25.1
MinM (cm^{-1})	8644	20.60	24.92	23.12	15.70	14.95	-27.4	-4.8
MaxEFG ($10^{21} Vm^{-2}$)	7286	27.80	31.70	27.65	25.82	23.27	-16.3	-9.9
Spillage	8755	0.372	0.451	0.375	0.350	0.343	-7.8	-2.0
SLME (%)	5877	6.460	7.058	6.184	5.690	4.948	-23.4	-13.0
PMEij (cm^{-2})	2711	0.112	0.132	0.133	0.101	0.112	0.0	10.9
PMDi	2649	16.50	19.73	19.53	15.07	16.07	-2.6	6.6
PMDiEl (ϵ_{11})	2592	2.800	3.088	3.248	2.733	2.631	-6.0	-3.7
PMDiIo (ϵ_{11})	2496	3.480	3.935	16 568	3.330	3.037	-12.7	-8.8
PMDij (cm^{-1})	1885	6.560	9.407	8.305	5.685	5.771	-12.0	1.5
Exfoli (meV/atom)	499	40.10	30.37	36.92	42.07	40.62	33.8	10.0

229 improve after 100 epochs and save the model with the best validation error). This can cause
230 ElemNet to train for more epochs, leading to longer training time. It would be interesting
231 to see if the performance of the proposed method improves if we increase the number of
232 epochs to train the model. Additionally, as the AtomSets framework provides different
233 features for composition and structure based information, analyzing the results provided in
234 Supplementary Table 4 can provide insights on what information is beneficial to train which
235 materials properties. Overall, the results demonstrate that the proposed framework performs
236 better with respect to other embeddings and other methodologies that use embedding and/or
237 features obtained using a pre-trained model to improve the performance of the model for
238 small datasets.

239 **Including Structural Polymorphs in Dataset** In the previous sections, we used a
240 dataset without structural polymorphs entries. This is because the model used to train
241 the small target dataset was designed to work with static, structure-agnostic embeddings.
242 These embeddings cannot differentiate between various structural polymorphs of the same
243 composition, so polymorph entries had to be removed during model training. However, as
244 the proposed dynamic embedding framework uses different sets of element-level encodings
245 for each element in every compound, by using them as embeddings for training the model,
246 we can effectively convert the composition-based models into structure-based models and
247 thereby mitigate the shortcomings of models trained on composition-based inputs only. We
248 compare the performance of the proposed framework with AtomSets, which uses structure-
249 based features obtained using a pre-trained model along with ElemNet and AutoML trained

250 using atom-wise averaged value of dynamic embedding. We could not use structure-based
251 representations such as classical force-field inspired descriptors (CFID)⁵³ and voronoi tes-
252 sellations based descriptors⁵⁴ as model input to ElemNet and AutoML as their respective
253 libraries failed to featurize substantial amount of training, validation and test sets. More-
254 over, using only structure information as input to the deep learning and traditional machine
255 learning model has been shown not to improve the performance of the model⁵⁵. As we have
256 multiple trained models for AtomSets, we use the model with the least test MAE to compare
257 the performance in Table 4 and show the rest of the results in Supplementary Tables 5.

258 Table 4 indicates that models trained using dynamic embedding outperform other em-
259 bedding and/or features based methods in 34/44 cases, i.e., in $\approx 77\%$ of the cases. For the
260 remaining 10 properties, the model trained on ElemNet using atom-wise averaged value of
261 dynamic embedding performed the best for 10/10 cases. Similar to the observations made
262 in the previous section, the methods that incorporate dynamic embedding in one form or
263 another can outperform other methods, and for the materials properties with close error
264 values between the proposed method and ElemNet trained using atom-wise averaged value
265 of dynamic embedding, we can explore the effect of epochs for training the model with the
266 performance. We also observe a clear benefit of using dynamic embeddings as features when
267 working with datasets that include structural polymorphs. These embeddings capture more
268 meaningful structural differences between compounds with identical compositions but dis-
269 tinct crystal structures, as evidenced by improved model performance. It is quite encouraging
270 to observe that the model that was originally designed to take composition-based inputs only

Table 4: The table shows the test MAE, and % error change between proposed method and (1) best performing existing method and (2) method using embedding extracted from ALIGNN as input for each of the target materials properties trained on different models for the prediction task of “Including Structural Polymorphs in Dataset”.

Property	Data Size	Existing	Embedding		Dynamic (proposed)	% Error	% Error
		AtomSets [36]	AutoML [43]	ElemNet [38]		Change (Existing)	Change (Embedding)
E _g OPT (eV)	54934	0.227	0.265	0.158	0.140	-38.3	-11.4
E _f (eV/atom)	54934	0.075	0.073	0.047	0.046	-38.7	-2.1
KLU (Å)	54606	11.45	9.99	9.41	9.92	-4.1	-13.4
Ehull (eV/atom)	54546	0.191	0.231	0.159	0.075	-60.7	-52.8
Encut (eV)	54469	161.0	152.4	129.9	122.0	-24.2	-6.1
Magoszi (μ_B)	51443	0.431	0.504	0.297	0.289	-32.9	-2.7
Magout (μ_B)	52125	0.588	0.700	0.446	0.430	-26.9	-3.6
Epsx	43760	30.07	24.97	21.81	21.48	-28.6	-1.5
Epsy	43760	29.32	25.00	22.10	21.31	-27.3	-3.6
Epsz	43760	28.78	24.76	21.94	21.93	-23.8	0.0
PPF ($\mu W m^{-1} K^{-2}$)	22692	585.5	490.0	445.3	457.2	-21.9	2.7
NPF ($\mu W m^{-1} K^{-2}$)	22692	587.1	518.8	468.2	481.9	-17.9	2.9
NSB ($\mu V/K$)	22609	58.56	51.43	43.69	42.70	-27.1	-2.3
PSB ($\mu V/K$)	22634	59.62	55.05	45.84	45.35	-23.9	-1.1
Nem300k (m_0)	20899	0.439	0.460	0.318	0.310	-29.4	-2.5
Pem300k (m_0)	19730	0.630	0.660	0.454	0.497	-21.1	9.5
ETC33 (GPa)	19530	48.78	38.44	34.45	34.44	-29.4	0.0
ETC22 (GPa)	19416	44.83	38.21	34.35	32.73	-27.0	-4.7
ETC11 (GPa)	19381	47.58	38.18	35.03	33.89	-28.8	-3.3
BulkKV (GPa)	19003	20.46	19.14	14.30	12.90	-37.0	-9.8
ETC13 (GPa)	19023	22.06	18.73	15.38	14.90	-32.5	-3.1
ETC12 (GPa)	18945	25.21	20.92	18.61	18.12	-28.1	-2.6
Poisson	18902	0.218	0.161	0.143	0.160	-26.6	11.9
ShearGV (GPa)	18004	12.607	10.787	8.878	8.683	-31.1	-2.2
E _g MBJ (eV)	17578	0.416	0.445	0.273	0.243	-41.6	-11.0
ETC44 (GPa)	17231	19.62	15.83	14.40	14.03	-28.5	-2.6
AvgME (electron mass unit)	17181	0.106	0.136	0.081	0.080	-24.5	-1.2
AvgMH (electron mass unit)	17181	0.157	0.192	0.121	0.115	-26.8	-5.0
ETC55 (GPa)	16859	16.34	14.36	12.17	12.27	-24.9	0.8
ETC66 (GPa)	16650	15.76	13.63	11.88	11.71	-25.7	-1.4
Mepsx	16349	31.23	28.36	24.69	23.36	-25.2	-5.4
Mepsy	16349	34.46	28.72	23.91	23.24	-32.6	-2.8
Mepsz	16349	34.09	29.36	24.02	24.55	-28.0	2.2
MaxM (cm^{-1})	13379	52.58	81.74	46.05	35.12	-33.2	-23.7
MinM (cm^{-1})	12678	30.83	30.28	20.59	20.45	-33.7	-0.7
MaxEFG ($10^{21} Vm^{-2}$)	11506	31.51	26.47	22.66	22.05	-30.0	-2.7
Spillage	11203	0.464	0.408	0.349	0.351	-24.4	0.6
SLME (%)	8923	6.228	5.849	4.996	4.445	-28.6	-11.0
PMEij (cm^{-2})	4745	0.143	0.139	0.109	0.098	-31.5	-10.1
PMDi	4600	23.57	22.61	15.35	15.27	-35.2	-0.5
PMDiEl (ϵ_{11})	4495	3.798	4.035	3.173	3.585	-5.6	13.0
PMDiIo (ϵ_{11})	4307	4.249	3.5919	3.160	3.122	-26.5	-1.2
PMDij (cm^{-1})	3290	14.29	13.87	10.05	9.92	-30.6	-1.3
Exfoli (meV/atom)	801	44.14	44.66	38.83	45.42	2.9	17.0

271 as input can outperform other methods when including structure-based information. Ad-
272 ditionally, we performed model training by incorporating an additional set of meaningful
273 information in the dynamic embedding to train Roost to observe if the current set of embed-
274 ding works well with other chemical information associated with a compound. The dynamic
275 embedding uses different element-level encoding to represent each element of every com-
276 pound. However, as we know that the ALIGNN uses bond and angle based features to train
277 the model, it can also be extracted from a given layer in the architecture. As each compound
278 has a different number of bond and angle information associated with it, we take the aver-
279 age across all that information and create one-dimensional compound-level encoding, which
280 will be different for every compound. Next, we concatenate this compound-level encoding
281 with element-level encoding to create a more informative dynamic embedding. We then use
282 this large dynamic embedding that combines compound-level encoding with element-level
283 encoding as the model input to train the model. Additionally, we can also concatenate static
284 embeddings with dynamic embeddings to improve the results. Hence, we try combinations
285 for a couple of static and dynamic embeddings with element-level encoding and compound-
286 level encoding for the analysis. The prediction accuracy of the proposed framework with
287 different combinations of embeddings as shown in Supplementary Table 6 indicates that
288 adding compound-level encodings, i.e., bond and angle representation, to the element-level
289 encodings helps improve the model’s performance. On the other hand, training model using
290 the combination of dynamic embeddings with static embedding does not contribute towards
291 improving the performance of the model. This, along with the other results, demonstrates

292 that the proposed framework can significantly and consistently help improve the prediction
293 of the materials properties across various domains with limited datasets, thereby potentially
294 saving time and resources in the process of future materials discovery.

295 **Discussion**

296 In this paper, we presented a novel framework that uses composition and structure aware
297 element-level encodings to represent a compound in a graph neural network (GNN) to build
298 accurate models with limited datasets for improved materials property prediction across a
299 wide variety of materials properties. Additionally, our framework offers a principled method
300 for incorporating structural polymorphs during training, even in neural network architectures
301 originally designed to accept only a single structure per composition. To show the benefit of
302 the proposed approach, we built an embedding extractor using a GNN architecture ALIGNN
303 that incorporates composition and structure based information as input and formation en-
304 ergy of the MP dataset as the materials property. This model was used to extract dynamic
305 embeddings to train a target model on 44 different materials properties to find that the pro-
306 posed framework produces accurate and robust models for datasets with limited data size.
307 We compare the performance of the target model trained using dynamic embedding with
308 conventional static embedding and methodologies that use embedding and/or features ob-
309 tained using a pre-trained model. We also observed that adding compound-level encodings,
310 i.e., bond and angle representation, to the element-level encodings helps improve the model’s
311 performance. On the other hand, training model using the combination of dynamic embed-

dings with static embedding does not contribute towards improving the performance of the model. Hence, one can also explore different sets of embeddings to train the target dataset or use more sophisticated architectures for the target model in a bid to boost the performance of the target model for a specific materials property. To check the robustness of the proposed framework even further, we perform additional experiments using the formation energy of the JARVIS dataset to examine the other applications of the dynamic embedding obtained using the embedding extractor. First, we performed model training using other deep neural network architecture that uses embeddings as input, i.e., CrabNet⁵¹ to see if dynamic embedding produces similar performance improvements as observed in the Roost. The model trained on original mat2vec embeddings gives a test MAE of 0.108 eV/atom, whereas the model trained on proposed dynamic embeddings gives a test MAE of 0.066 eV/atom, which is $\approx 39\%$ improvement in accuracy. This shows that dynamic embedding can be applied to other types of deep neural network architectures that take embeddings as input for datasets with limited data size. Next, we evaluate the computational overhead associated with the proposed dynamic embeddings relative to conventional static embeddings (matscholar) when training the Roost model, using formation energy (a widely studied materials property) from the JARVIS dataset, using 10-fold cross-validation in Supplementary Table 7. We find that both the per-epoch training time and the total training time for Roost models employing dynamic embeddings are comparable to those using static embeddings, despite the substantial gains in predictive accuracy achieved with dynamic embeddings. The principal additional cost arises during the initial embedding extraction stage; however, this step is performed only

333 once and the resulting embeddings can be reused across multiple downstream tasks without
334 re-extraction. Consequently, dynamic embeddings provide a computationally efficient and
335 reusable representation that consistently enhances model performance across diverse target
336 material properties. At last, we present the training time per epoch (in seconds) and the
337 total number of epochs for each target material property trained using ElemNet or Roost
338 in Tables 3 and 4 in Supplementary Table 8. Although both Roost and ElemNet perform
339 similarly, with Roost showing slightly better average performance, it is notable that Roost
340 takes more time per epoch, while ElemNet requires more epochs to train.

341 The proposed framework is thus flexible and robust, and can leverage a wide variety
342 of deep neural networks to improve upon the performance using a rich set of embeddings.
343 This model is primarily designed for solids with well-defined structures. However, existing
344 studies ^{28,38} have demonstrated that transfer learning can be applied to various materials
345 properties across domains with limited datasets making such approaches agnostic to the
346 type of downstream application. This indicates potential transferability to materials with
347 less precise structures, such as polymer melts or ionic liquids. Previous works ^{56,57} that have
348 applied TL to materials with less precise structures have used simple ML/DL techniques to
349 train the models. Extending this approach to more recent models, such as the one proposed
350 in our work, could offer significant improvements and represents a promising direction for
351 future research. The proposed framework is expected to easily adapt to ever-increasing data,
352 ever-advancing model training techniques, and other scientific domains beyond materials
353 science. Moreover, the presented framework is conceptually easy to implement, understand,

354 use, and build upon. For future work, it would be interesting to explore the effect on the
355 performance of the target model when materials properties other than formation energy
356 are used as the material property and GNN architecture other than ALIGNN is used for
357 training the embedding extractor. It would also be interesting to see the effect of dynamic
358 embedding when used as input for the state-of-the-art GNN architecture or when training
359 even larger target datasets. Additionally, we can also explore the uncertainty associated
360 with the materials property prediction by incorporating neural network components that
361 help perform uncertainty estimation, such as dropout within the network architecture, or
362 by creating an ensemble model using multiple graph neural networks. One can also explore
363 different combinations of embeddings to train the neural network or use more sophisticated
364 neural network architectures for the target model in a bid to boost the performance of the
365 target model for a specific materials property.

366 **Methods**

367 **Embedding Extractor** In this work, we implement an embedding extractor to extract
368 element-level encoding from DFT-relaxed structures to train the target model. The model
369 pre-trained on ALIGNN using formation energy from the MP dataset is used as an embedding
370 extractor to extract atom-based information from a given layer, each containing a variable
371 number of rows depending on the number of the atom information present in the input file
372 and 256 columns as embedding for each row. For example, let us consider a hypothetical
373 compound $W_wX_xY_yZ_z$ where $w + x + y + z = Q$ and the number of unique elements is 4.

374 Next, we extract the embeddings before and after every ALIGNN and GCN layer, where the
375 dimensions of the extracted embeddings will be $(Q, 256)$. We then take the element-wise
376 mean of all embeddings across each column to create a $(4, 256)$ embedding representation
377 where each row represents the element-wise embedding for a given compound of the target
378 dataset. The extracted embeddings from a given layer can then be used as a dynamic
379 embedding for any GNN network that uses embedding as a form of model input.

380 **Network Settings and Model Architecture** Roost¹⁷ and ALIGNN²¹ were implemented
381 using Pytorch. The hyperparameters used in the Roost comprise of the following: Adam as
382 the optimizer with weight decay parameter of 10^{-6} , LeakyReLU as the base activation func-
383 tion, mini-batch size of 128, and learning rate as 0.0003. We train all Roost models for 250
384 epochs as done in the original work¹⁷. The hyperparameters used in the ALIGNN comprise
385 of the following: Adaptive Moment Estimation with decoupled weight decay (AdamW) as
386 the optimizer with normalized weight decay of 10^{-5} , Sigmoid Linear Unit (SiLU) as the base
387 activation function, mini-batch size of 64, and learning rate as 0.001. We train all ALIGNN
388 models for 300 epochs with a fixed random seed as done in the original work²¹. Readers
389 interested in in-depth hyperparameter settings for Roost and ALIGNN models are referred
390 to those publications^{17,21} for details. We use mean absolute error (MAE) as the loss function
391 as well as the primary evaluation metric for all models. Each model is trained on a single
392 Tesla V100-PCIE-16GB Graphics processing unit (GPU). For all other downstream models
393 we used the default hyperparameters and training procedures for training the models. For all
394 other downstream models, we retained the default hyperparameters and training procedures

395 recommended in the original publications and reference codebases referenced in the code
396 availability section.

397 **Data availability** The datasets used in this paper are publicly available from the corresponding
398 websites- MP ⁴² from <https://materialsproject.org/> and JARVIS ^{43,44} from [https://jarvis.](https://jarvis.nist.gov)
399 [nist.gov](https://jarvis.nist.gov).

400 **Code availability** The code for extracting dynamic embeddings from the pre-trained ALIGNN
401 model used in the proposed framework is available at [https://github.com/GuptaVishu2002/](https://github.com/GuptaVishu2002/DynamicEmbedding)
402 [DynamicEmbedding](https://github.com/GuptaVishu2002/DynamicEmbedding). For the downstream prediction models, we used publicly available imple-
403 mentations from the following repositories: Roost (<https://github.com/CompRhys/roost>), Elem-
404 Net (<https://github.com/NU-CUCIS/ElemNet>), CrabNet ([CrabNet](https://github.com/anthony-wang/
405 <a href=)), AtomSets (<https://github.com/materialyzeai/maml>), and AutoML ([sklearn](https://github.com/automl/au
406 <a href=)).

407 References

- 408 1. Dane Morgan and Ryan Jacobs. Opportunities and challenges for machine learning in
410 materials science. Annual Review of Materials Research, 50:71–103, 2020.
- 411 2. Pascal Friederich, Florian Häse, Jonny Proppe, and Alán Aspuru-Guzik. Machine-
412 learned potentials for next-generation matter simulations. Nature Materials, 20(6):750–
413 761, 2021.

- 414 3. Julia Westermayr, Michael Gastegger, Kristof T Schütt, and Reinhard J Maurer. Per-
415 spective on integrating machine learning into computational chemistry and materials
416 science. The Journal of Chemical Physics, 154(23):230903, 2021.
- 417 4. Yuwei Mao, Mahmudul Hasan, Arindam Paul, Vishu Gupta, Kamal Choudhary,
418 Francesca Tavazza, Wei-keng Liao, Alok Choudhary, Pinar Acar, and Ankit Agrawal.
419 An ai-driven microstructure optimization framework for elastic properties of titanium
420 beyond cubic crystal systems. npj Computational Materials, 9(1):111, 2023.
- 421 5. Ghanshyam Pilania. Machine learning in materials science: From explainable predictions
422 to autonomous design. Computational Materials Science, 193:110360, 2021.
- 423 6. Rui Jiao, Wenbing Huang, Peijia Lin, Jiaqi Han, Pin Chen, Yutong Lu, Yang Liu. Crys-
424 tal structure prediction by joint equivariant diffusion Advances in Neural Information
425 Processing Systems, 36:17464–17497, 2023.
- 426 7. Rampi Ramprasad, Rohit Batra, Ghanshyam Pilania, Arun Mannodi-Kanakkithodi,
427 and Chiho Kim. Machine learning in materials informatics: recent applications and
428 prospects. npj Computational Materials, 3(1):54, dec 2017.
- 429 8. Ankit Agrawal and Alok Choudhary. Deep materials informatics: Applications of deep
430 learning in materials science. MRS Communications, 9(3):779–792, 2019.
- 431 9. Vishu Gupta, Wei-keng Liao, Alok Choudhary, and Ankit Agrawal. Evolution of artificial
432 intelligence for application in contemporary materials science. MRS Communications,
433 pages 1–10, 2023.

- 434 10. Kamal Choudhary, Daniel Wines, Kangming Li, Kevin F. Garrity, Vishu Gupta, Aldo
435 H. Romero, Jaron T. Krogel, Kayahan Saritas, Addis Fuhr, Panchapakesan Ganesh,
436 Paul R. C. Kent, Keqiang Yan, Yuchao Lin, Shuiwang Ji, Ben Blaiszik, Patrick Reiser,
437 Pascal Friederich, Ankit Agrawal, Pratyush Tiwary, Eric Beyerle, Peter Minch, Trevor
438 David Rhone, Ichiro Takeuchi, Robert B. Wexler, Arun Mannodi-Kanakithodi, Elif
439 Ertekin, Avanish Mishra, Nithin Mathew, Mitchell Wood, Andrew Dale Rohskopf, Ja-
440 son Hattrick-Simpers, Shih-Han Wang, Luke E. K. Achenie, Hongliang Xin, Maureen
441 Williams, Adam J. Biacchi, Francesca Tavazza JARVIS-Leaderboard: a large scale
442 benchmark of materials design methods. npj Computational Materials, 10(1):93 , 2023.
- 443 11. Kamal Choudhary, Brian DeCost, Chi Chen, Anubhav Jain, Francesca Tavazza, Ryan
444 Cohn, Cheol Woo Park, Alok Choudhary, Ankit Agrawal, Simon J. L. Billinge, Elizabeth
445 Holm, Shyue Ping Ong, Chris Wolverton Recent advances and applications of deep
446 learning methods in materials science. npj Computational Materials, 8(1):59, 2022.
- 447 12. Xingyue Shi, Linming Zhou, Yuhui Huang, Yongjun Wu, Zijian Hong A review on the
448 applications of graph neural networks in materials science at the atomic scale Materials
449 Genome Engineering Advances, 2(2):e50, 2024.
- 450 13. Zhenyao Fang, Qimin Yan Towards accurate prediction of configurational disorder prop-
451 erties in materials using graph neural networks npj Computational Materials, 10(1):91,
452 2024.
- 453 14. Hongwei Du, Jiamin Wang, Jian Hui, Lanting Zhang, Hong Wang, DenseGNN: universal

- 454 and scalable deeper graph neural networks for high-performance property prediction in
455 crystals and molecules npj Computational Materials, 10(1):292, 2024.
- 456 15. Reshma Devi, Keith T Butler, Gopalakrishnan Sai Gautam, Optimal pre-train/fine-
457 tune strategies for accurate material property predictions npj Computational Materials,
458 10(1):300, 2024.
- 459 16. Tian Xie and Jeffrey C. Grossman. Crystal graph convolutional neural networks for
460 an accurate and interpretable prediction of material properties. Phys. Rev. Lett.,
461 120:145301, Apr 2018.
- 462 17. Rhys EA Goodall and Alpha A Lee. Predicting materials properties without crystal
463 structure: Deep representation learning from stoichiometry. Nature communications,
464 11(1):1–9, 2020.
- 465 18. Kristof Schütt, Pieter-Jan Kindermans, Huziel Enoc Saucedo Felix, Stefan Chmiela,
466 Alexandre Tkatchenko, and Klaus-Robert Müller. Schnet: A continuous-filter convolu-
467 tional neural network for modeling quantum interactions. Advances in neural information
468 processing systems, 30, 2017.
- 469 19. Chi Chen, Weike Ye, Yunxing Zuo, Chen Zheng, and Shyue Ping Ong. Graph networks
470 as a universal machine learning framework for molecules and crystals. Chemistry of
471 Materials, 31(9):3564–3572, 2019.

- 472 20. Johannes Klicpera, Shankari Giri, Johannes T Margraf, and Stephan Günnemann. Fast
473 and uncertainty-aware directional message passing for non-equilibrium molecules. arXiv
474 preprint arXiv:2011.14115, 2020.
- 475 21. Kamal Choudhary and Brian DeCost. Atomistic line graph neural network for improved
476 materials property predictions. npj Computational Materials, 7(1):1–8, 2021.
- 477 22. Dipendra Jha, Vishu Gupta, Wei-keng Liao, Alok Choudhary, and Ankit Agrawal. Mov-
478 ing closer to experimental level materials property prediction using ai. Scientific reports,
479 12, 2022.
- 480 23. Atsuto Seko, Atsushi Togo, Hiroyuki Hayashi, Koji Tsuda, Laurent Chaput, and Isao
481 Tanaka. Prediction of low-thermal-conductivity compounds with first-principles anhar-
482 monic lattice-dynamics calculations and bayesian optimization. Physical Review Letters,
483 115(20):205901, 2015.
- 484 24. Luca M Ghiringhelli, Jan Vybiral, Sergey V Levchenko, Claudia Draxl, and Matthias
485 Scheffler. Big data of materials science: Critical role of the descriptor. Physical Review
486 Letters, 114(10):105503, 2015.
- 487 25. Joohwi Lee, Atsuto Seko, Kazuki Shitara, Keita Nakayama, and Isao Tanaka. Prediction
488 model of band gap for inorganic compounds by combination of density functional theory
489 calculations and machine learning techniques. Physical Review B, 93(11):115104, 2016.
- 490 26. Austin D Sendek, Qian Yang, Ekin D Cubuk, Karel-Alexander N Duerloo, Yi Cui, and
491 Evan J Reed. Holistic computational structure screening of more than 12000 candidates

- 492 for solid lithium-ion conductor materials. Energy & Environmental Science, 10(1):306–
493 320, 2017.
- 494 27. Yuwei Mao, Shahriyar Keshavarz, Vishu Gupta, Andrew CE Reid, Wei-keng Liao, Alok
495 Choudhary, and Ankit Agrawal. Ai for learning deformation behavior of a material:
496 Predicting stress-strain curves 4000x faster than simulations. In 2023 International Joint
497 Conference on Neural Networks (IJCNN), pages 1–8. IEEE, 2023.
- 498 28. Vishu Gupta, Kamal Choudhary, Francesca Tavazza, Carelyn Campbell, Wei-keng Liao,
499 Alok Choudhary, and Ankit Agrawal. Cross-property deep transfer learning framework
500 for enhanced predictive analytics on small materials data. Nature communications,
501 12(1):1–10, 2021.
- 502 29. Vishu Gupta, Kamal Choudhary, Yuwei Mao, Kewei Wang, Francesca Tavazza, Carelyn
503 Campbell, Wei-keng Liao, Alok Choudhary, and Ankit Agrawal. Mppredictor: An ar-
504 tificial intelligence-driven web tool for composition-based material property prediction.
505 Journal of Chemical Information and Modeling, 63(7):1865–1871, 2023.
- 506 30. Vishu Gupta, Wei-keng Liao, Alok Choudhary, and Ankit Agrawal. Pre-activation based
507 representation learning to enhance predictive analytics on small materials data. In 2023
508 International Joint Conference on Neural Networks (IJCNN), pages 1–8. IEEE, 2023.
- 509 31. Vishu Gupta, Wei-keng Liao, Alok Choudhary, and Ankit Agrawal. Brnet: Branched
510 residual network for fast and accurate predictive modeling of materials properties. In

- 511 Proceedings of the 2022 SIAM International Conference on Data Mining (SDM), pages
512 343–351. SIAM, 2022.
- 513 32. Vishu Gupta, Alec Peltekian, Wei-keng Liao, Alok Choudhary, and Ankit Agrawal.
514 Improving deep learning model performance under parametric constraints for materials
515 informatics applications. Scientific Reports, 13(1):9128, 2023.
- 516 33. Vishu Gupta, Youjia Li, Alec Peltekian, Muhammed Nur Talha Kilic, Wei-keng Liao,
517 Alok Choudhary, and Ankit Agrawal. Simultaneously improving accuracy and computa-
518 tional cost under parametric constraints in materials property prediction tasks. Journal
519 of Cheminformatics, 16(1):17, 2024.
- 520 34. Anthony Onwuli, Ashish V Hegde, Kevin VT Nguyen, Keith T Butler, and Aron Walsh.
521 Element similarity in high-dimensional materials representations. Digital Discovery,
522 2(5):1558–1564, 2023.
- 523 35. Tim Hsu, Tuan Anh Pham, Nathan Keilbart, Stephen Weitzner, James Chapman, Peng-
524 hao Xiao, S Roger Qiu, Xiao Chen, and Brandon C Wood. Efficient and interpretable
525 graph network representation for angle-dependent properties applied to optical spec-
526 troscopy. npj Computational Materials, 8(1):151, 2022.
- 527 36. Yuxin You, Zhen Liu, Xiangchao Wen, Yongtao Zhang, Wei Ai, . Large language models
528 meet graph neural networks: a perspective of graph mining. Mathematics, 10(7):1147,
529 2022.

- 530 37. Jongmin Han, Youngchun Kwon, Youn-Suk Choi, Seokho Kang, . Improving chem-
531 ical reaction yield prediction using pre-trained graph neural networks. Journal of
532 Cheminformatics, 16(1):25, 2024.
- 533 38. Vishu Gupta, Kamal Choudhary, Brian DeCost, Francesca Tavazza, Carelyn Campbell,
534 Wei-keng Liao, Alok Choudhary, and Ankit Agrawal. Structure-aware graph neural
535 network based deep transfer learning framework for enhanced predictive analytics on
536 diverse materials datasets. npj Computational Materials, 10(1):1, 2024.
- 537 39. Zhenwen Sheng, Hui Zhu, Bo Shao, Yu He, Zhuang Liu, Suqin Wang, Ming Sheng,
538 . Accelerated Discovery of Energy Materials via Graph Neural Network. Inorganics,
539 13(12):395, 2025.
- 540 40. Dennis Delali Kwesi Wayo. Ensembles of Graph and Physics-Informed Machine Learn-
541 ing for Scientific Modeling in Materials Science: A Review: DDK Wayo. Archives of
542 Computational Methods in Engineering, 1–26, 2025.
- 543 41. Archit Anand, Priyanka Kumari, Ajay Kumar Kalyani. High throughput screening of
544 new piezoelectric materials using graph machine learning and knowledge graph approach.
545 Computational Materials Science, 246:113445, 2025.
- 546 42. Anubhav Jain, Shyue Ping Ong, Geoffroy Hautier, Wei Chen, William Davidson
547 Richards, Stephen Dacek, Shreyas Cholia, Dan Gunter, David Skinner, Gerbrand Ceder,
548 and Kristin a. Persson. The Materials Project: A materials genome approach to accel-
549 erating materials innovation. APL Materials, 1(1):011002, 2013.

- 550 43. Kamal Choudhary, Kevin F. Garrity, Andrew C. E. Reid, Brian DeCost, Adam J. Biacchi, Angela R. Hight Walker, Zachary Trautt, Jason Hattrick-Simpers, A. Gilad Kusne, 551 Andrea Centrone, Albert Davydov, Jie Jiang, Ruth Pachter, Gowoon Cheon, Evan 552 Reed, Ankit Agrawal, Xiaofeng Qian, Vinit Sharma, Houlong Zhuang, Sergei V. Kalinin, 553 Bobby G. Sumpter, Ghanshyam Pilania, Pinar Acar, Subhasish Mandal, Kristjan Haule, 554 David Vanderbilt, Karin Rabe, and Francesca Tavazza. JARVIS: An integrated infras- 555 tructure for data-driven materials design, 2020.
- 557 44. Daniel Wines, Ramya Gurunathan, Kevin F. Garrity, Brian DeCost, Adam J. Biacchi, 558 Francesca Tavazza, and Kamal Choudhary. Recent progress in the JARVIS infrastructure 559 for next-generation data-driven materials design. Applied Physics Reviews, 10(4):041302, 560 10 2023.
- 561 45. Chi Chen and Shyue Ping Ong. Atomsets as a hierarchical transfer learning framework 562 for small and large materials datasets. npj Computational Materials, 7(1):173, 2021.
- 563 46. Jian-Gang Kong, Ke-Lin Zhao, Jian Li, Qing-Xu Li, Yu Liu, Rui Zhang, Jia-Ji Zhu, 564 and Kai Chang. Self-supervised representations and node embedding graph neural net- 565 works for accurate and multi-scale analysis of materials. Machine Learning: Science and 566 Technology, 5(3):035018, 2024.
- 567 47. Dipendra Jha, Logan Ward, Arindam Paul, Wei-keng Liao, Alok Choudhary, Chris 568 Wolverton, and Ankit Agrawal. ElemNet: Deep learning the chemistry of materials 569 from only elemental composition. Scientific reports, 8(1):17593, 2018.

- 570 48. Logan Ward, Ankit Agrawal, Alok Choudhary, and Christopher Wolverton. A General-
571 Purpose Machine Learning Framework for Predicting Properties of Inorganic Materials.
572 npj Computational Materials, 2(August):16028, 2016.
- 573 49. Vahe Tshitoyan, John Dagdelen, Leigh Weston, Alexander Dunn, Ziqin Rong, Olga
574 Kononova, Kristin A Persson, Gerbrand Ceder, and Anubhav Jain. Unsupervised
575 word embeddings capture latent knowledge from materials science literature. Nature,
576 571(7763):95–98, 2019.
- 577 50. Leigh Weston, Vahe Tshitoyan, John Dagdelen, Olga Kononova, Amalie Trewartha,
578 Kristin A Persson, Gerbrand Ceder, and Anubhav Jain. Named entity recognition and
579 normalization applied to large-scale information extraction from the materials science
580 literature. Journal of chemical information and modeling, 59(9):3692–3702, 2019.
- 581 51. Anthony Yu-Tung Wang, Steven K Kauwe, Ryan J Murdock, and Taylor D Sparks.
582 Compositionally restricted attention-based network for materials property predictions.
583 npj Computational Materials, 7(1):1–10, 2021.
- 584 52. Matthias Feurer, Aaron Klein, Katharina Eggenberger, Jost Springenberg, Manuel
585 Blum, and Frank Hutter. Efficient and robust automated machine learning. In Advances
586 in Neural Information Processing Systems 28 (2015), pages 2962–2970, 2015.
- 587 53. Kamal Choudhary, Brian DeCost, and Francesca Tavazza. Machine learning with force-
588 field-inspired descriptors for materials: Fast screening and mapping energy landscape.
589 Physical review materials, 2(8):083801, 2018.

- 590 54. Logan Ward, Ruoqian Liu, Amar Krishna, Vinay I Hegde, Ankit Agrawal, Alok Choud-
591 hary, and Chris Wolverton. Including crystal structure attributes in machine learning
592 models of formation energies via Voronoi tessellations. Physical Review B, 96(2):024104,
593 2017.
- 594 55. Dipendra Jha, Vishu Gupta, Logan Ward, Zijiang Yang, Christopher Wolverton, Ian
595 Foster, Wei-keng Liao, Alok Choudhary, and Ankit Agrawal. Enabling deeper learning
596 on big data for materials informatics applications. Scientific reports, 11(1):1–12, 2021.
- 597 56. Alfred Yan, Tatiana Sokolinski, William Lane, Jinwang Tan, Kim Ferris, and Emily M
598 Ryan. Applying transfer learning with convolutional neural networks to identify
599 novel electrolytes for metal air batteries. Computational and Theoretical Chemistry,
600 1205:113443, 2021.
- 601 57. Zhan Ma, Shu Wang, Minhee Kim, Kaibo Liu, Chun-Long Chen, and Wenxiao Pan.
602 Transfer learning of memory kernels for transferable coarse-graining of polymer dynam-
603 ics. Soft Matter, 17(24):5864–5877, 2021.

604 **Acknowledgements** This work was performed under the following financial assistance award
605 70NANB24H136 from U.S. Department of Commerce, National Institute of Standards and Tech-
606 nology as part of the Center for Hierarchical Materials Design (CHiMaD). Partial support is also
607 acknowledged from NSF award OAC-2331329 and Northwestern Center for Nanocombinatorics.
608 Certain commercial equipment, instruments, software, or materials are identified in this paper in
609 order to specify the experimental procedure adequately. Such identifications are not intended to

610 imply recommendation or endorsement by NIST, nor it is intended to imply that the materials or
611 equipment identified are necessarily the best available for the purpose.

612 **Author Contributions** V.G. designed and carried out the implementation and experiments for
613 the dynamic embedding under the guidance of A.A., A.C., and W.L.. Y.L. and M.N.T.K performed
614 experiments to train some of the models and collect performance results. K.C. and D.W. provided
615 the necessary domain expertise for this work. V.G., A.A., D.W., and K.C. wrote the manuscript.
616 All authors discussed the results and reviewed the manuscript.

617 **Competing Interests** The authors declare that they have no competing interests.