

# Graph to Activation Energy Models Easily Reach Irreducible Errors but Show Limited Transferability

Sai Mahit Vadaddi,<sup>†</sup> Qiyuan Zhao,<sup>\*,‡</sup> and Brett M. Savoie<sup>\*,†</sup>

<sup>†</sup>*Davidson School of Chemical Engineering, Purdue University, West Lafayette, IN, 47906*

<sup>‡</sup>*Department of Medicinal Chemistry, University of Michigan, Ann Arbor, MI, 48109*

E-mail: zhaoqiy@med.umich.edu; bsavoie@purdue.edu

## Abstract

Activation energy characterization of competing reactions is a costly, but crucial step for understanding the kinetic relevance of distinct reaction pathways, product yields, and myriad other properties of reacting systems. The standard methodology for activation energy characterization has historically been a transition state search using the highest level of theory that can be afforded. However recently, several groups have popularized the idea of predicting activation energies directly, based on nothing more than the reactant and product graphs, a sufficiently complex neural network, and a broad enough dataset. Here, we have revisited this task using the recently developed Reaction Graph Depth 1 (RGD1) transition state dataset and several newly developed graph attention architectures. All of these new architectures achieve similar state-of-the-art results of  $\sim 4$  kcal/mol mean absolute error on withheld testing sets of reactions but poor performance on external testing sets composed

of reactions with differing mechanisms, reaction molecularity, or reactant size distribution. Limited transferability is also shown to be shared by other contemporary graph to activation energy architectures through a series of case-studies. We conclude that an array of standard graph architectures can already achieve results comparable to the irreducible error of available reaction datasets but that out-of-distribution performance remains poor.

## 1 Introduction

Reaction activation energies ( $E_a$ ) and heats of reaction ( $\Delta H_r$ ) are essential to understanding of reactivity in applications ranging from biofuel utilization,<sup>1-6</sup> drug design,<sup>7-10</sup> and materials stability.<sup>11-13</sup> Collecting this information from experiments is costly,<sup>14-18</sup> making it highly desirable to develop predictive methods that can be used prior to synthesis to expedite hypothesis formation and optimization.<sup>19-25</sup> Over the past several decades, quantum chemistry has delivered many algorithms for localizing transition states (TSs) and characterizing activation energies;<sup>26-31</sup> however, finding transition states remains relatively expensive for on-the-fly and high-throughput applications. It would be a qualitative advance if reaction properties like activation energy could be directly calculated without first localizing a transition state.

In the past few years several groups have shown the feasibility of predicting activation energies from only the reactant and product graphs (Fig. 1A). Early examples focused on summarizing changes between reactants and products using expert-generated features<sup>32-36</sup> and molecular fingerprints.<sup>37-40</sup> With the advent of larger reaction datasets,<sup>41-43</sup> several standard practices have been identified, including the necessity to avoid including reverse reactions in testing datasets (i.e., having an example in the testing set that was seen during training with the reactant and product switched),<sup>44</sup> the potential usefulness of the heat of reaction ( $\Delta H_r$ ) as an input feature,<sup>38,44</sup> and the advantage of learnable reaction fingerprints over pre-defined expert fingerprints.<sup>44-46</sup> In paral-

19 lel, prediction strategies are also being developed based on three-dimensional featurizations of the  
 20 reactant and product,<sup>47,48</sup> relatively inexpensive information from approximate levels of theory,<sup>49</sup>  
 21 and additional features from quantum chemistry.<sup>39,50–55</sup> The cost of additional features can nullify  
 22 the advantage of using a graph-based approach (albeit, with potentially higher transferability). In  
 23 this sense, the ideal model would be able to achieve high accuracy based solely on the reactant  
 24 and product graphs. The negotiation of these trade-offs remains a live issue.<sup>44,56,57</sup>

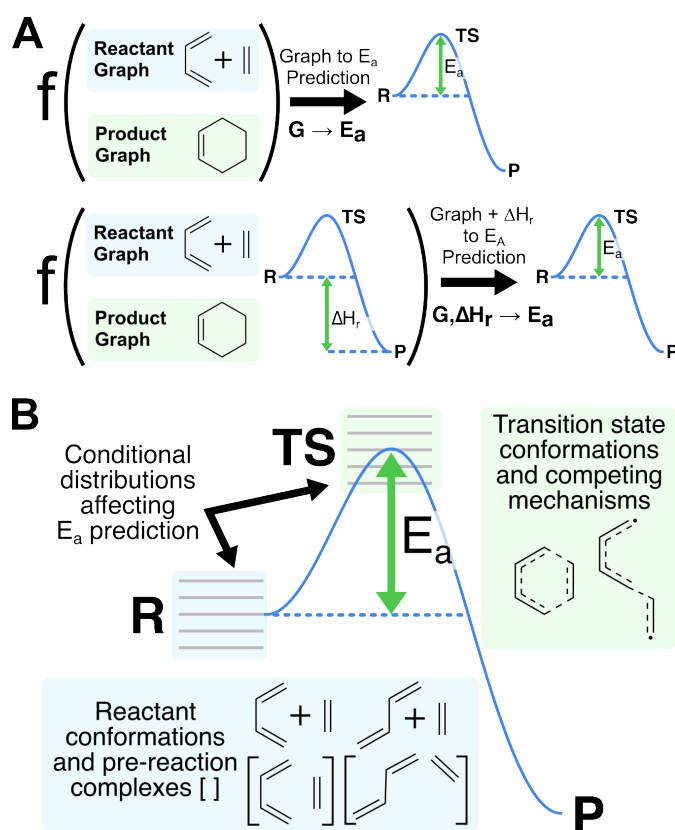


Figure 1: Overview of the graph to activation energy (G2Ea) prediction task. (a) A minimal featurization of this problem consists of only using the reactant and product graphs, while sometimes relatively inexpensive but informative features like heat of reaction ( $\Delta H_r$ ) are also used. (b)  $E_a$  depends on the reference initial state and the sampling distribution used to localize the transition state. Illustrative examples are shown for the sequential versus concerted Diels-Alder mechanisms. Comparisons across datasets and predictions on unseen reactions will be out-of-distribution if such factors are not consistently sampled.

25 Despite many practical demonstrations of the graph-to-activation-energy (G2Ea) concept, sev-

26 eral challenges persist that limit the usefulness of these models as drop-in replacements for quantum-  
27 chemistry based TS searches. One challenge is that the scarcity of large reaction datasets has  
28 limited convincing out-of-distribution tests of the transferability of G2Ea models. Under the as-  
29 sumption that reaction mechanisms are conserved, G2Ea models should be capable of extrapolating  
30 to unseen substrates involved in classes of reactions that have been directly trained on, but this  
31 remains a largely untested hypothesis.

32 A more fundamental challenge that is specific to the G2Ea learning task is that the prediction  
33 problem is underdetermined with respect to the manner in which computational  $E_a$  datasets are  
34 currently generated (Fig. 1B). While experiments measure an effective activation energy based  
35 on the Boltzmann average over all accessible conformations and transition states, computationally  
36 derived activation energies are typically extracted from a single pair of energies corresponding  
37 to a particular reactant and TS conformation that are not always uniquely defined. Available  
38 datasets have no guarantee of conformational completeness or of having found a globally minimum  
39 energy barrier for the observed reactions. This will lead to an irreducible error for any G2Ea  
40 model trained on available quantum-chemistry derived datasets, since predictions are conditioned  
41 on the conformational distribution of reactants and TSs used during curation. For example, the  
42 Conformer-Rotamer Ensemble Sampling Tool (CREST)<sup>58</sup> algorithm was used during the curation  
43 of the Reaction Graph Depth 1 (RGD1) dataset<sup>41</sup> to find the minimum energy conformer of the  
44 isolated reactants, and a protocol specific to Yet Another Reaction Program (YARP) was used  
45 to select up to three conformations for double-ended TS searches.<sup>59</sup> These choices will show up  
46 as inductive biases in models trained on RGD1 and any incomplete conformational sampling will  
47 show up as an irreducible error when predicting on unseen reactions.

48 Here, the transferability of G2Ea models has been revisited using the recently developed RGD1  
49 dataset and an adaptation of the graph attention architecture.<sup>60-62</sup> The motivation for this study  
50 was that the size and mechanistic diversity of the RGD1 dataset potentially allows for the training

51 of more data-demanding architectures with better transferability. Using a graph attention archi-  
52 tecture we are able to train models that reliably approach the estimated irreducible error of RGD1  
53 and perform well on withheld reactions drawn from the same distribution. Nevertheless, these  
54 models show minimal transferability in external testing scenarios, oftentimes performing worse  
55 than naïve mean-predicting models. Through several comparative case-studies it is shown that  
56 this behavior is shared by models using closely related featurizations, and other contemporary  
57 graph to property architectures recently published for use in the G2Ea prediction task.

## 58 **2 Methods**

59 Several variations of the Edge-featured Graph Attention Network<sup>62</sup> (EGAT) architecture are de-  
60 veloped here for the prediction of activation energies from reactant and product graphs. The main  
61 elements shared by all the architectures are described in the Model Overview, Input Features,  
62 and Description subsections (Sections 2.1-2.3) and the differences between models are described in  
63 Learning Tasks subsection (Section 2.4). These models are compared with Chemprop,<sup>63</sup> a directed  
64 message passing neural network (D-MPNN) model that has previously been trained for the G2Ea  
65 task. The training and implementation details of Chemprop in this work followed those supplied  
66 by the developers through their distributed code,<sup>63</sup> additional can be found in the SI (Section 2).

### 67 **2.1 Model Overview**

68 Edge-featured Graph Attention Networks (EGAT) are a subset of Graph Attention Networks  
69 (GAT).<sup>60-62</sup> The basic idea behind EGATs is to use the features of each edge to create an attention  
70 score that is used to weight the information mixing between nodes. Each pass through an EGAT  
71 layer results in the mixing of information between nodes as determined by learnable attention  
72 scores. Thus, the use of  $n$  EGAT layers results in the mixing of information from nodes up to  $n$

73 edges away into the fingerprint of each node. A fingerprint of the whole reaction graph is obtained  
74 by pooling the node and edge fingerprints, which can be used for reaction property prediction.

75 Here, the EGAT architecture is adapted for  $E_a$  prediction by featurizing each atom as a node  
76 and each bond as an edge. The reactant and product are separately passed to the model, converted  
77 to individual fingerprints, and then the fingerprint difference is used to predict  $E_a$  subject to several  
78 small architectural variations (See Section 2.4). The relative atom-mapping between reactants and  
79 products affects the prediction of the model, but the architecture includes an intermediate pooling  
80 operation to ensure the prediction is invariant to the absolute numbering of the atom sequence.  
81 Differences in atom mapping reflect distinct bond-changes and so the atom-mapping awareness  
82 is an important feature of the model that is not captured by models based on simpler Morgan  
83 fingerprints.

84 The edge topology that is used for the reactant and product corresponds to the union of the  
85 bonds that are presented in either of the species. This results in molecular fingerprints for the  
86 reactant and product that have the same number of edges when taking the fingerprint difference.  
87 Thus, when processed by the EGAT model, bonds that are only physically present in the reactant  
88 or product (i.e., they are formed or broken in the reaction) are still present as edges in the product  
89 and reactant graphs, respectively, albeit with a special feature indicating that the bond was broken  
90 or formed.

## 91 **2.2 Input Features**

92 The input to the EGAT architecture is an ordered set of features for each atom and bond in the  
93 reactant and product. Table 1 and 2 list the input atom (i.e., node) features,  $\mathbf{n}$ , and bond (i.e.,  
94 edge) features,  $\mathbf{e}$ , fed into the EGAT model, respectively, and their comparisons with input features  
95 of Chemprop. Many of the features could be incorporated as either distinct numeric values (i.e.,  
96 integers) or categorical values (i.e., one-hot vectors). In such cases we elected to use integers over

97 one-hot vectors wherever possible to simplify the representation. For example, hybridization in  
98 EGAT is represented using a four bit one-hot vector (i.e., this is clearly a categorical feature), but  
99 the number of bonded hydrogens are featurized with integers (i.e., this is a case that could be  
100 treated either numerically or categorically). Wherever one-hots are used, they are transposed and  
101 concatenated with the other features during the input preprocessing. The resulting dimensions for  
102 the featurized node and edge vectors are 17x1 and 14x1, respectively.

103 The rationale for some of the new features is as follows. A node's proximity to a bond that  
104 is broken or formed is novel to the EGAT model, which denotes the effect it may have on the  
105 reaction as atoms closer to the reaction center exhibit more local changes than those farther away.  
106 Node features corresponding to the number of bonded CHON atoms were also directly featurized,  
107 even though this is an implicit feature that might be learned from the convolution. The novel edge  
108 features include a set of binary features corresponding to whether the bond was broken, formed,  
109 if the bond order changed, or was unchanged. These features are necessary to use the union of  
110 the reactant and product graphs as the inputted graph topology to the model. For example, the  
111 reactant may have some edges with bond order zero, because these might only exist in the product,  
112 but it is nonetheless potentially useful for the model to mix information along these edges given  
113 the fact that a bond forms or breaks between such atoms in the reaction.

## 114 **2.3 Model Description**

115 The overall model architecture consists of four EGAT layers that yield compressed fingerprints of  
116 the reactant(s) and product(s), which are then used to predict the  $E_a$  using a feed-forward stack  
117 (Fig. 2). Even though there are four total EGAT layers, only two are distinct, with the first being  
118 unique and the last three sharing weights. The first is responsible for embedding the raw node  
119 and edge features in a higher dimensional space so it has different internal dimensions from the  
120 other three. The three-fold application of the second EGAT layer results in mixing of information

Table 1: Input Node Features in the EGAT and Chemprop Models

Features	Chemprop	EGAT
Atom Type	100x1 one-hot vector	Integer
Number of Bonds	6x1 one-hot vector	-
Charge	5x1 one-hot vector	Integer
Neighboring H	5x1 one-hot vector	Integer
Neighboring C	-	Integer
Neighboring N	-	Integer
Neighboring O	-	Integer
Distance to Nearest Reacting Atom	-	Float
Hybridization	5x1 one-hot vector	4x1 one-hot vector <sup>1</sup>
Aromaticity	One-hot Value	One-hot Value
Atomic Mass	Mass/100	Float
Chirality	4x1 one-hot vector	3x1 one-hot vector
Ring Atom	-	One-hot value

1. The bits correspond to whether the atom is s,sp,sp<sup>2</sup>,or sp<sup>3</sup> hybridized. These should be expanded if applying the model beyond second row chemistry.

Table 2: Input Edge Features in the EGAT and Chemprop Models

Features	Chemprop	EGAT
Bond Type	4x1 one-hot vector	5x1 one-hot vector <sup>1</sup>
Conjugated	One-hot Value	One-hot Value
In Ring	One-hot Value	One-hot Value
Stereochemistry	4x1 one-hot vector	3x1 one-hot vector
Bond Change	-	4x1 one-hot vector <sup>2</sup>
Change in Bond Order	-	One-hot value

1. Four of the vector values determine the order of the bond, while the additional value determines whether the bond is aromatic.

2. One of the bits denotes whether there is a bond break, a second denotes whether a bond forms, a third denotes whether the bond order changes, a fourth denotes whether the bond is unchanged.



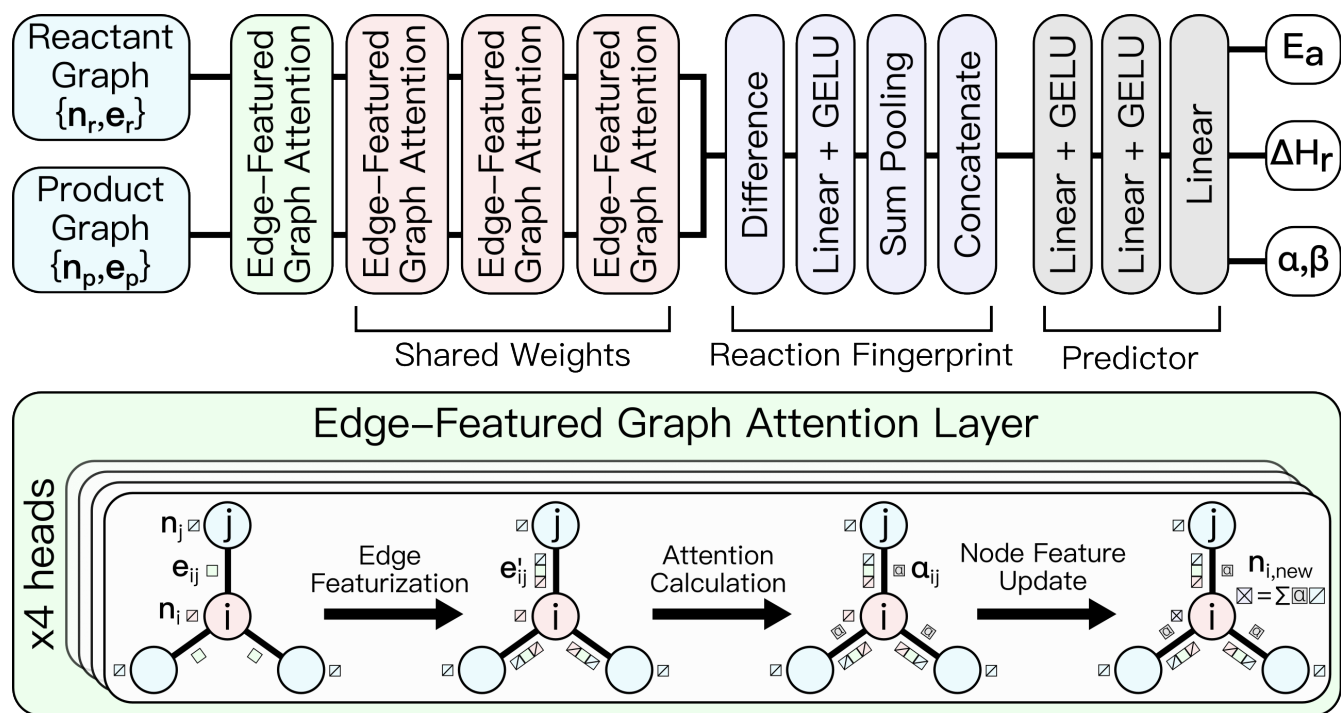


Figure 2: Overview of the EGAT model architecture developed in this work. The model consists of a stack of edge-featured graph attention layers (green and red) that generate the derived features for each edge and node in the reactant and product. The differences in these features are then transformed and pooled to generate a global fixed-size reaction fingerprint (purple) that is used as an input for the  $E_a$  prediction task (gray). Depending on the formulation of the learning task, the final predictor may be trained to solely predict  $E_a$ ,  $E_a$  and  $\Delta H_r$  in multi-task fashion, or  $\alpha$  and  $\beta$  the linear free energy parameters associated with  $E_a$ . The lower inset illustrates the edge-featured graph attention mechanism applied to a single node and its edges.

121 from nodes separated by up to four edges away. Experiments on validation data showed no further  
122 benefit from using more EGAT layers or distinct weights in consecutive layers.

123 Within each EGAT layer, each edge is embedded based on both the edge and node features of  
124 the atoms connected by the edge. For an edge connecting nodes,  $i$  and  $j$ , the edge features,  $\mathbf{e}_{ij}$ ,  
125 and connected node features,  $\mathbf{n}_i$  and  $\mathbf{n}_j$ , are embedded as a concatenated vector

$$\mathbf{e}'_{ij} = \mathbf{A} (\mathbf{n}_i || \mathbf{e}_{ij} || \mathbf{n}_j) + \mathbf{a} \quad (1)$$

126 where  $||$  refers to the catenation operation,  $\mathbf{A}$  is a learnable linear transformation, and  $\mathbf{a}$  is  
127 a bias term. In the first EGAT layer,  $\mathbf{A}$  has dimensions of  $512 \times 48$ , where 48 is the size of the  
128 catenated inputted node and edge features (Tables 1 and 2) and 512 is the embedding dimension  
129 such that  $\mathbf{e}'_{ij}$  has dimensions of  $512 \times 1$ . In subsequent EGAT layers,  $\mathbf{A}$  has dimensions of  $512 \times 1536$ ,  
130 where 1536 is the size of the catenated embedded node and edge features, and  $\mathbf{e}'_{ij}$  has unchanged  
131 dimensions of  $512 \times 1$ . Recall that edges are defined based on the union of the bonds in the reactant  
132 and product graphs, such that an edge may exist in the reactant or product where no physical  
133 bond exists. The embedded edges are directional (i.e.,  $\mathbf{e}'_{ij} \neq \mathbf{e}'_{ji}$ ) and so there are twice as many  
134 embedded edges as connected nodes.

135 Each node is embedded into a size  $512 \times 1$  space using a learnable linear projection of its node  
136 features according to

$$\mathbf{n}'_i = \mathbf{B} \mathbf{n}_i + \mathbf{b} \quad (2)$$

137 where  $\mathbf{B}$  is a learnable linear transformation and  $\mathbf{b}$  is a  $512 \times 1$  bias vector. In the first EGAT layer,  
138  $\mathbf{B}$  has dimensions of  $512 \times 17$ , where 17 is the size of the inputted node features (Table 1) and 512  
139 is the embedding dimension such that  $\mathbf{n}'_i$  has dimensions of  $512 \times 1$ . In subsequent EGAT layers,  
140  $\mathbf{B}$  has dimensions of  $512 \times 512$ , where 512 is the size of the embedded node features after the first  
141 EGAT layer, and  $\mathbf{n}'_i$  has unchanged dimensions of  $512 \times 1$ .

142 Each graph edge attention layer performs a non-linear transform and mixing operation between  
143 nodes connected by edges. The embedded edges,  $\mathbf{e}'_{ij}$ , are transformed to create attention scores that  
144 determine the degree of mixing between nodes along each attention head. Here, each EGAT layer  
145 has four heads, meaning that the embedded 512x1 node feature vector,  $\mathbf{n}'_i$ , is reshaped into four  
146 128x1 vectors,  $\mathbf{n}'_{i,h}$ , each of which is mixed with the corresponding reshaped vectors of neighboring  
147 nodes,  $\mathbf{n}'_{j,h}$ , based on the attention values of each head,  $\alpha_{ij}^h$ . The mathematical description of these  
148 steps is as follows.

149 The edge vectors are reshaped into a 128x1x4 tensor,  $\mathbf{e}'_{ij,h}$ , where each column vector corresponds  
150 to a query for each head in the attention mechanism. This matrix is non-linearly transformed using  
151 the LeakyReLU function and projected along learnable directions by each head according to

$$\epsilon_{ij} = \mathbf{C} \text{LeakyReLU} \left( \mathbf{e}'_{ij,h} \right) \quad (3)$$

152 where  $\mathbf{C}$  is a 1x128x4 learnable tensor, with each row vector corresponding to the key for each  
153 head in the attention mechanism, and  $\epsilon_{ij}$  is a 1x1x4 tensor holding the logits for each head. The  
154  $\epsilon_{ij}$  values for each head are then softmax normalized over all edges originating at node i to obtain  
155 the attention scores

$$\alpha_{ij} = \text{softmax}(\epsilon_{ij}) = \frac{e^{\epsilon_{ij}}}{\sum_{k \in N(i)} e^{\epsilon_{ik}}} \quad (4)$$

156 where  $N(i)$  refers to all edges originating at node i, the operations are performed per-element,  
157 and  $\alpha_{ij}$  is a 1x1x4 tensor. The softmax guarantees that  $\sum_{k \in N(i)} \alpha_{ik} = \mathbf{1}$  such that the attention  
158 scores associated with each head can be interpreted as mixing probabilities along each edge. Finally,  
159 the inputted 512x1 node feature vectors,  $\mathbf{n}'_i$ , are reshaped into 128x1x4 tensors,  $\mathbf{n}'_{i,h}$ , where each  
160 row vector plays the role of a value in the attention mechanism. The updated node features are  
161 calculated as attention-weighted mixtures of the neighboring node features according to

$$\mathbf{n}_{i,\text{out}}^h = \sum_{k \in N(i)} \alpha_{ik} \mathbf{n}_k^{',h} \quad (5)$$

162 where  $\mathbf{n}_{i,\text{out}}^h$  is a 128x1x4 tensor that is reshaped to a 512x1 vector,  $\mathbf{n}_{i,\text{out}}$ , before being returned  
 163 by the layer. Each EGAT layer returns both the edge features,  $\mathbf{e}'_{ij}$ , and node features,  $\mathbf{n}_{i,\text{out}}$ , for  
 164 use by subsequent layers.

165 The nodes and edges of the reactant and product graphs are separately transformed by the  
 166 EGAT layers to yield a set of 512x1 node and edge vectors for the reactant and product. The same  
 167 stack of layers and weights are used for the reactant and product graphs. Reaction node and edge  
 168 features are obtained based on the differences between these vectors,

$$\mathbf{n}_{i,\text{rxn}} = \mathbf{n}_{i,\text{product}} - \mathbf{n}_{i,\text{reactant}} \quad (6)$$

$$\mathbf{e}_{ij,\text{rxn}} = \mathbf{e}_{ij,\text{product}} - \mathbf{e}_{ij,\text{reactant}} \quad (7)$$

169 Where the product and reactant labels refer to the vector features outputted by the final EGAT  
 170 layer. The reaction features are then subjected to a final projection and non-linear transform

$$\mathbf{n}'_{i,\text{rxn}} = \text{GeLU}(\mathbf{D}\mathbf{n}_{i,\text{rxn}} + \mathbf{d}) \quad (8)$$

$$\mathbf{e}'_{ij,\text{rxn}} = \text{GeLU}(\mathbf{F}\mathbf{e}_{ij,\text{rxn}} + \mathbf{f}) \quad (9)$$

171 where  $\mathbf{D}$  and  $\mathbf{F}$  are 512x512 learnable matrices and  $\mathbf{d}$  and  $\mathbf{f}$  are 512x1 bias vectors. At this  
 172 stage there are a variable number of node and edge features depending on the size and topology of  
 173 the reactant and product graphs. A fixed size 1024x1 reaction fingerprint,  $\mathbf{FP}$ , that is invariant to

174 the absolute atomic numbering is obtained sum pooling the reaction node and edge features and  
175 concatenating the result

$$\mathbf{FP} = \sum_i \mathbf{n}'_{i,\text{rxn}} || \sum_{ij} \mathbf{e}'_{ij,\text{rxn}}. \quad (10)$$

176 **FP** is used as an input to a stack of three feed-forward layers, with hidden dimensions of  
177 [256,128,1], a GeLU activation function after the first and second layers, and the last linear layer  
178 mapping to one output in the simplest model. Several small variations of the **FP** predictor stack  
179 were explored, depending on the formulation of the learning task as described next.

## 180 2.4 Learning Tasks

181 Four models were developed using the shared architecture for generating reaction fingerprints  
182 described in the previous section but with minor variations in the predictor stack. The basic  
183 model predicts  $E_a$  as a scalar output of the model using only the reactant and product graphs.  
184 The **FP** predictor stack for this model has dimensions of [256,128,1], with the final layer predicting  
185  $E_a$ . We refer to this model in the results as the graph to  $E_a$  model ( $G \rightarrow E_a$ ).

186 Three other models were developed that use  $\Delta H_r$  as an additional input feature or as an  
187 additional prediction target. The first uses  $\Delta H_r$  as an additional input to the **FP** predictor stack,  
188 so that it has dimensions of [257,128,1], with the final layer predicting  $E_a$ . We refer to this model in  
189 the results section as the graph and  $\Delta H_r$  to  $E_a$  model ( $G, \Delta H_r \rightarrow E_a$ ). The second uses  $\Delta H_r$  as an  
190 indirect feature by having the model predict  $E_a$  using the Bell-Evans-Polanyi<sup>64</sup> (BEP) relationship

$$E_a = \alpha \Delta H_r + \beta \quad (11)$$

191 where  $\alpha$  and  $\beta$  are constants predicted by the model. The predictor stack for this model  
192 has dimensions of [256,128,2], where  $\alpha$  and  $\beta$  are predicted in the final layer and  $\Delta H_r$  is used

193 to calculate  $E_a$ . The rationale for this physics-inspired formulation is that it might show better  
194 transferability due to the well established linear free-energy relationships expected for reactions  
195 sharing a common mechanism. We refer to this model in the results as the graph to BEP model  
196 (G→BEP).

197 The last model uses a multi-task predictor stack with dimensions of [256,128,2] that predicts  
198 both  $\Delta H_r$  and  $E_a$  in the last layer. Whereas both of the previous models require  $\Delta H_r$  to make  
199 a prediction, the rationale for multi-task training is that it indirectly informs the model of the  
200 underlying mechanistic information associated with  $\Delta H_r$  without requiring it at the time of pre-  
201 diction. This is a common form of transfer learning. We refer to this model in the results as the  
202 graph to  $\Delta H_r$  and  $E_a$  model (G→ $\Delta H_r, E_a$ ).

## 203 2.5 Training Details

204 The loss function for training the single-task models (G→ $E_a$ ; G, $\Delta H_r$ → $E_a$ ; G→BEP) was the  
205 mean absolute error in  $E_a$  prediction

$$L = \frac{1}{N_{\text{batch}}} \sum_i^{N_{\text{batch}}} |E_{a,i,0} - E_{a,i,p}| \quad (12)$$

206 where  $i$  runs over all samples in the batch,  $N_{\text{batch}}$  is the number of samples per batch,  $E_{a,i,0}$   
207 refers to the reference activation energy, and  $E_{a,i,p}$  refers to the predicted activation energy. For  
208 the multi-task model (G→ $\Delta H_r, E_a$ ) trained to predict both  $\Delta H_r$  and  $E_a$ , the loss function was  
209 weighted to prioritize  $E_a$  accuracy according to

$$L = \frac{1}{N_{\text{batch}}} \sum_i^{N_{\text{batch}}} 0.8|E_{a,i,0} - E_{a,i,p}| + 0.2|\Delta H_{r,i,0} - \Delta H_{r,i,p}| \quad (13)$$

210 where  $\Delta H_{r,i,0}$  refers to the reference heat of reaction,  $\Delta H_{r,i,p}$  refers to the predicted heat of  
211 reaction, and all other symbols have the same meaning as in Eq. 12.

212 All EGAT models were trained using the Adam optimizer and a batch size of 50. The learning  
213 rate,  $\eta$ , was initially set to 5e-4, linearly increased each update step for 10 epochs to 1e-3, followed  
214 by an exponential decay to a minimum of 1e-5. Early stopping was applied to terminate train-  
215 ing if the validation loss did not decrease in 30 consecutive epochs. Further information on the  
216 hyperparameters for each model can be found in the Supplemental Information. All Chemprop  
217 models were trained using the optimized parameters found via the hyperopt package listed in the  
218 Supplemental Information (SI) section of the Heid et. al.<sup>46</sup> paper.

## 219 **2.6 Data**

220 The Chemprop and EGAT models used in this paper were trained on the RGD1 dataset,<sup>41</sup> which  
221 contains around 177,000 reactions with up to ten heavy atoms consisting of carbon, hydrogen,  
222 nitrogen and oxygen. In brief, the RGD1 dataset was generated by a graph-based enumeration  
223 of  $\sim 700,000$  reactions involving reactants sampled from PubChem.<sup>20,65</sup> A reaction conformational  
224 sampling strategy<sup>20,58</sup> was applied to generate up to three conformations for each reaction that  
225 were used to initialize double-ended TS searches,<sup>29</sup> followed by Berny optimization,<sup>66</sup> and in-  
226 trinsic reaction coordinate<sup>67</sup> (IRC) validation at the GFN2-xTB<sup>68</sup> level of theory. The GFN2-  
227 xTB optimized TSs were further refined using the Gaussian16<sup>69</sup> quantum chemistry engine at the  
228 B3LYP-D3/TZVP<sup>70,71</sup> level of theory with D3 dispersion.<sup>72</sup> The DFT-level TS were classified as  
229 intended or unintended using an XGBoost<sup>73</sup> model that uses geometric features of the TS and the  
230 GFN2-xTB level information to classify the TS as intended or unintended (i.e., whether the TS  
231 corresponds to the reaction that was used to seed the TS search). This model exhibits a testing  
232 set accuracy of  $\sim 95\%$  in a previous smaller testing set.<sup>20</sup> For a detailed description of the RGD1  
233 database, we direct readers to our previous publication.<sup>41</sup>

234 Several data processing steps were applied to prepare the RGD1 data for training the activation  
235 energy models. Firstly, reactions were pruned by the uniqueness of the reactant and product

236 InChIKeys<sup>74</sup> to eliminate the presence of cases where reactions may have the same reactant and  
237 product but only differ by tautomerization. This was done to prevent possible data leakage by  
238 training on one tautomer while including another in the testing set. Secondly, only the minimum  
239 activation energy and its corresponding heat of reaction was listed for each reaction. This step is  
240 required because RGD1 contains multiple conformationally distinct TSs for many of the reactions.  
241 In total 135,455 distinct reactions passed these filters. Random splitting was used to generate  
242 80/10/10 training (101,236 reactions), validation (12,139), and testing splits (12,139). The splits  
243 included only one direction of each reaction (i.e., the data was not augmented by including the  
244 reverse versions of reactions in training or testing). Unless otherwise stated, all accuracies are  
245 reported for prediction accuracy on the RGD1 testing set.

## 246 **3 Results and Discussion**

### 247 **3.1 Overall RGD1 Performance**

248 The overall performance of the various EGAT architectures was tested by predicting the activation  
249 energy of the 12.3k reactions in the RGD1 testing split (Fig. 3A). There is remarkably little  
250 variation across the different EGAT architectures; all models show uniformly low mean absolute  
251 prediction errors (MAE) of  $\sim 4$  kcal/mol. None of the models show a significant systematic bias  
252 based on the negligible mean signed errors (Fig. S1A). The  $G \rightarrow E_a$  model that doesn't use the  
253  $\Delta H_r$  information intuitively shows the largest MAE, but only performs  $\sim 0.5$  kcal/mol worse on  
254 average than  $G \rightarrow \text{BEP}$ , the best EGAT model. The models were tested as ensembles of five  
255 independently trained models of each class. The standard deviation in the mean testing split  
256 performance of the individual models comprising the ensemble provides an estimate of the testing  
257 performance uncertainty (error bars in Fig. 3A). The mean performance uncertainty is also within



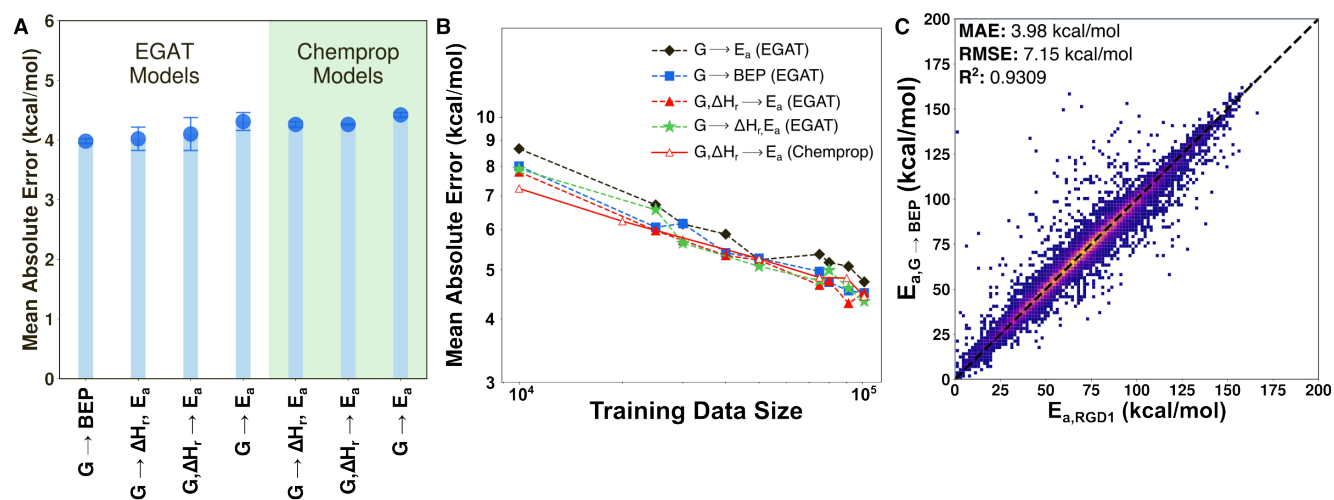


Figure 3: Performance of the G2Ea architectures on the RGD1 dataset. (A) Mean performance of each EGAT and Chemprop architecture on the RGD1 testing set ranked by accuracy. Error bars correspond to the standard deviation in mean performance across the ensemble. Each datapoint reflects the mean over five models trained with independent starting weights but the same training data. (B) Testing set error versus training dataset size for the EGAT architectures and best overall Chemprop model. Linear scaling is expected on a log-log plot. Each datapoint represents the best performance of a single model. (C) Parity plot showing the performance of best EGAT model ( $G \rightarrow \text{BEP}$ ) on the testing set reactions. Individual pixels are colored by the density of datapoints (purple to yellow is low to high).

258  $\sim 0.5$  kcal/mol in all cases.

259 The performance of analogous Chemprop models trained and tested on RGD1 provide a useful  
260 reference (Fig. 3A). Similar to the EGAT architectures, the Chemprop model that eschews  $\Delta H_r$   
261 information performs the worst on average, but only by  $\sim 0.5$  kcal/mol compared with the multitask  
262 Chemprop architecture ( $G \rightarrow \Delta H_r, E_a$ ). The best EGAT architecture marginally outperforms the  
263 best Chemprop architecture in this comparison, but we do not consider this difference significant  
264 for reasons that are further explored below.

265 How should we interpret all of the architectures approaching a MAE of  $\sim 4$  kcal/mol on the  
266 RGD1 testing split? To answer this the reader should consider that there are several known sources  
267 of irreducible error in RGD1—indeed to some extent all available computational TS datasets suffer  
268 from these issues. These sources include the inaccuracy of the underlying DFT method (B3LYP-  
269 D3/TZVP), incomplete conformational sampling of TSs and reactants, and the potential inclusion  
270 of TSs that correspond to unintended reactions (i.e., they connect a reactant and product that  
271 are different than the label). Because we are testing on DFT-level  $E_a$  values and not experimental  
272 values, the absolute DFT errors should not contribute to the irreducible error for this particular  
273 learning task; but incomplete conformational sampling and mislabeled TSs still represent sources  
274 of irreducible error. Below we will provide some lower bound estimates of the latter errors, but  
275 cumulatively they are expected to easily amount to an uncertainty of 4 kcal/mol across the entire  
276 dataset. As such, we interpret the prediction errors of all the models as effectively approaching  
277 the irreducible error of this prediction task.

278 If the models are approaching the irreducible error of this particular learning task, then that  
279 should be evident in the known reducible sources of error, such as dataset size, approaching zero.  
280 To test this, the training data size versus testing set accuracy learning curves were generated by  
281 training individual models for each architecture on subsets of the training data and evaluating  
282 their performance (Fig. 3B). Between training sample sizes of 10k and 40k, the models show

283 large error reductions (MAE reductions from 8 kcal/mol to 5.5 kcal/mol in the largest cases), but  
284 for sizes between 75k and 100k, the models show average improvements of less than 1 kcal/mol.  
285 Extrapolating the power law scaling evident in these curves, we estimate that increasing training  
286 data from order 100k to 1m samples would only further reduce errors by  $\sim 0.1$  kcal/mol for the  
287 best EGAT and Chemprop architectures. The minimal errors associated with training data size  
288 are consistent with the interpretation that the models are approaching the prediction accuracy  
289 limit associated with the irreducible error of this task.

290 Only looking at the mean performance obscures the number of outliers predicted by all the  
291 models. The presence of outliers are already evidenced by the significantly lower median absolute  
292 errors—between 2-3 kcal/mol for all EGAT and chemprop models—compared with the MAEs  
293 (Fig S1A). To illustrate some individual outliers, a parity plot is presented that shows the per  
294 sample accuracy of the  $G \rightarrow$  BEP model on the testing set (Fig. 3B). The median absolute  
295 deviation ( $MAD = \text{median}(|X_i - \text{median}(X_i)|)$ ) has been used as a robust estimator for the  
296 width of the error distribution that is minimally affected by outliers. Assuming normal statistics  
297 the corresponding estimate for the standard deviation in the testing errors for this model is 3.40  
298 kcal/mol. The estimated standard deviation in prediction errors is similar when calculated for the  
299 underlying reaction classes (Fig. S4A), reactions with distinct molecularity (Fig. S4B), and across  
300 model architectures. The interested reader can also find violin plots of the MAEs of different  
301 populations of the testing set reactions differentiated by reaction type (Fig S5c) and molecularity  
302 (Fig S5D) in the SI. As an estimate of the number of outliers, there are a total of 796 testing  
303 samples (out of 12,139 total, or  $\sim 7\%$ ) with absolute errors greater than 10.2 kcal/mol (3x the  
304 standard deviation as estimated from the median absolute deviation). Based on the testing set  
305 size and assuming normal statistics only  $\sim 14$  samples (i.e.,  $\sim 0.1\%$  of the testing split) would be  
306 expected with errors this large. An analogous calculation using the absolute percent errors rather  
307 than the absolute errors produces a large number of percentage based outliers that is driven mainly

308 by low barrier reactions and not necessarily unphysical predictions (Fig. S6), and these are not  
309 further analyzed here.

## 310 **3.2 Sources of Error**

311 Although the EGAT architectures achieve mean accuracies on RGD1 comparable to the best G2EA  
312 models previously published, we consider it useful spend some time elaborating the obstacles to  
313 further accuracy improvements (Fig. 4). First we analyzed the conservation of poorly predicted  
314 outliers across the different architectures by histogramming testing samples by their error percentile  
315 for each EGAT model and calculating the overlapping samples in each bin (Fig. 4A). This analysis  
316 reveals that the architectures tended to perform poorly on the same samples, whereas the other  
317 bins show overlaps that are more consistent with chance (the expected overlap is  $\sim 1\%$  for this  
318 number of bins and models).

319 It is possible that the outliers represent unusually hard samples for all models, but the more  
320 likely scenario is that the outliers are cases affected by incomplete conformational sampling or  
321 with unintended transition states (i.e., the known sources of irreducible error in RGD1). To test  
322 this hypothesis, we took the 100 lowest error testing samples and 100 highest error testing samples  
323 based on the  $G \rightarrow$  BEP model performance and re-investigated their transition states. For these 200  
324 reactions from the RGD1 testing set, transition state searches were re-performed using 10 rather  
325 than 3 reaction conformers (i.e., the RGD1 curation protocol) and DFT-level IRC calculations  
326 were performed on the RGD1 TSs.

327 The additional conformational sampling was done to estimate the error associated with incom-  
328 plete sampling by calculating the average reduction in  $E_a$  for the top 100 and bottom 100 samples  
329 (Fig. 4B). For example, if the additional conformations merely rediscovered the same TS or higher  
330 energy TSs, then the reduction would be zero, but if they led to the discovery of lower barrier TSs,  
331 then the reduction would be positive. This experiment revealed that incomplete conformational

332 sampling disproportionately affects the poorly predicted samples, with the most and least accurate  
333 populations showing mean  $E_a$  reductions of  $\sim 2.5$  kcal/mol and  $\sim 6.5$  kcal/mol, respectively. The  
334 reaction with the largest  $E_a$  reduction of  $\sim 30$  kcal/mol occurs in the poorly predicted samples  
335 versus a maximum reduction of  $\sim 13$  kcal/mol in the top 100 predictions (See Fig. S6 for other  
336 relevant analyses). For the bottom 100 predictions that were intended after the IRC calculation  
337 (i.e., the 27 reactions where unintended TSs are not a confounding factor as revealed by the anal-  
338 ysis in the next paragraph) the G $\rightarrow$ BEP model underestimates  $E_a$  by 21.8 kcal/mol on average  
339 compared with a negligible mean signed error on the whole dataset (Fig. S1A), meaning that the  
340 model recognizes these reactions as being conformationally undersampled. The  $E_a$  reduction of  
341  $\sim 2.5$  kcal/mol observed in the best predictions can be considered a lower bound on the irreducible  
342 error of RGD1, assuming no other factors contribute to the error and that unintended reactions  
343 can be perfectly filtered. Thus while conformational sampling errors significantly contribute to the  
344 irreducible error for G2Ea prediction on RGD1, alone they cannot explain the majority of outliers  
345 that comprise the worst predictions.

346 IRC calculations on the top 100 and bottom 100 samples were performed to investigate the  
347 prevalence of unintended reactions in the two populations (Fig. 4C). Recall that IRC calculations  
348 are expensive at the DFT level and so RGD1 was curated using a machine learning model to filter  
349 intended and unintended TSs based on IRC calculations performed at the GFN2-xTB level.<sup>41</sup>  
350 Despite the high accuracy of this model, presumably some unintended reactions remain in the  
351 dataset. This experiment revealed that unintended reactions are disproportionately represented in  
352 the bottom 100 predictions (73%) compared with the best predictions (1%). This suggests that  
353 a majority of the outliers in Fig. 3C are in fact unintended reactions. This result motivated us  
354 to perform a larger IRC study on the 1000 worst predicted samples from the training split, which  
355 returned a similarly high proportion (63.8%) of unintended reactions (Figs. S8). Based on these  
356 tests, unintended reactions likely comprise only a few percent of RGD1, but they are the cause of

357 most poorly predicted outliers. Recalculating the MAE without the three sigma outliers results in  
 358 an improvement of  $\sim 1 - 2$  kcal/mol across the models, which can be taken as an estimate of the  
 359 irreducible error associated with imperfect filtering of unintended reactions from RGD1.

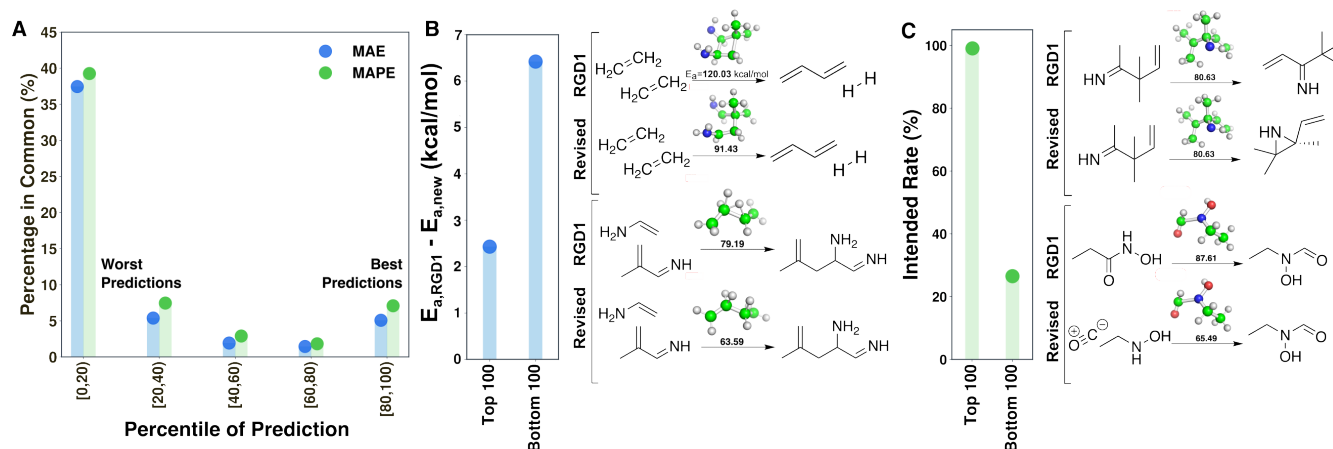


Figure 4: Investigating the origin of outlier behavior in G2Ea models. (A) Comparison of the fraction of testing set reactions that are consistently predicted accurately and inaccurately by each EGAT model. The rank of each reaction was determined by accuracy, binned by performance pentile, then the membership in each pentile was compared across models (percentage in common). (B) The mean  $E_a$  reduction after expanding the conformational sampling for the 100 best and 100 worst predicted testing set reactions by the G $\rightarrow$ BEP model. Two examples showing newly discovered TSs with lower barriers (revised) are shown. (C) The intended rates determined by IRC calculation for the RGD1 TSs of the 100 best and 100 worst predicted testing set reactions by the G $\rightarrow$ BEP model. Two illustrative examples are shown where the TS connected a different product (top) and different reactant (bottom) after the IRC calculation (revised). In the latter case,  $E_a$  is also incorrect because of the misidentified reactant.

### 360 3.2.1 Model Transferability

361 Graph-based neural-network models have become notorious in many contexts for overfitting and  
 362 poor out-of-distribution performance.<sup>44,56,57</sup> Although the models trained on RGD1 show excellent  
 363 testing performance on unseen reactions, this is a large dataset and reactions typically involve a  
 364 small number of bond changes and conserved mechanisms. This means that even if the testing  
 365 set involves unseen reactions in terms of reactants or products, it is not expected to necessarily

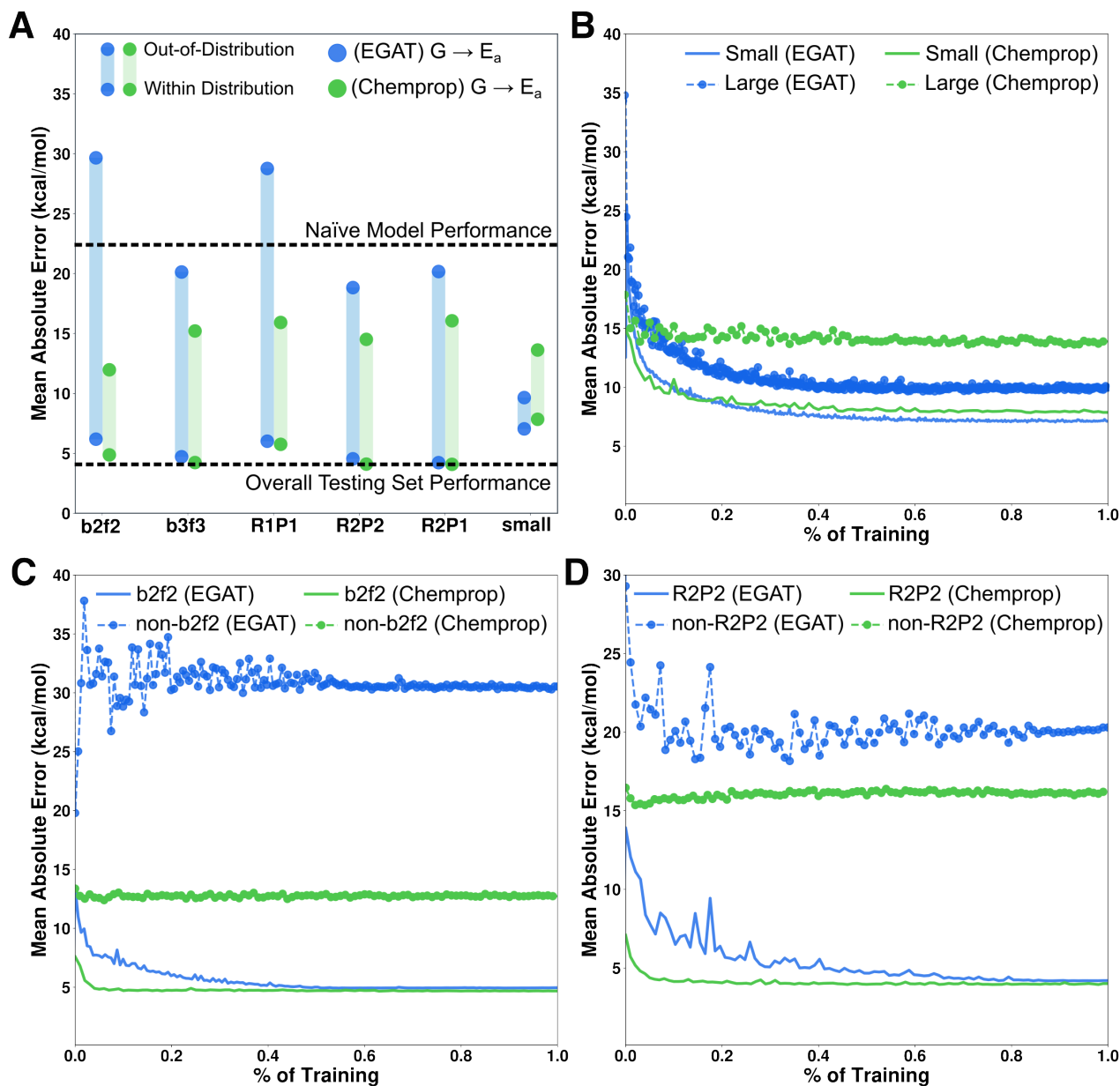


Figure 5: Transferability case-studies for  $G \rightarrow E_a$  models. (A) Mean performance of models trained on subsets of the RGD1 dataset indicated by the x-axis label. These models were tested on the RGD1 testing set reactions with the same reaction type as in training (within distribution) and reactions types excluded from training (out-of-distribution). (B) Learning curves for EGAT and Chemprop  $G \rightarrow E_a$  models trained on small reaction (3-7 heavy atoms) and evaluated on both small and large (8-10 heavy atoms) reactions from the validation split. (C) Learning curves for EGAT and Chemprop  $G \rightarrow E_a$  models trained on b2f2 data and evaluated on both b2f2 and non-b2f2 reactions from the validation split. (D) Learning curves for EGAT and Chemprop  $G \rightarrow E_a$  models trained on R2P2 data and evaluated on both R2P2 and non-R2P2 reactions from the validation split.

366 present novel reactivity (e.g., in terms of new types of bonds being broken and formed) that is not  
367 seen elsewhere in the training data.

368 To interrogate these architectures with a more rigorous test of out-of-distribution performance,  
369 several case studies were performed where the EGAT and Chemprop architectures were trained  
370 on subsets of the training data split that excluded one or more classes of reactions (Fig. 5).  
371 Three factors were used to classify different subpopulations within RGD1—the number of bonds  
372 broken and formed in the reaction (bnfn), the molecularity (number of species) of the reactant  
373 and product (RnPn), and the size of the reactant (more than seven heavy atoms was considered  
374 large). The manner in which RGD1 was generated resulted in reactions involving breaking two  
375 bonds and forming two bonds (b2f2) being the most abundant class of reactions and b3f3 reactions  
376 being the second most abundant with 22,538 and 68,615 reactions in the training dataset, respec-  
377 tively. In terms of molecularity, R1P1 (i.e., unimolecular) and R2P2 (i.e., bimolecular reactant  
378 and bimolecular product) reactions are the most abundant with 38,240 and 16,817 reactions in the  
379 training dataset, respectively. Similarly, the heavy atom cutoff was selected to make the small and  
380 large subpopulations approximately evenly matched with 40,018 and 57,746 training set reactions,  
381 respectively. In these case studies, the models were trained on one subpopulation of the training  
382 split but tested on the original testing split. The performance on the class of reactions included in  
383 training is referred to as “within distribution” and the class of reactions excluded from training is  
384 referred to as “out-of-distribution” in the discussed comparisons. All performance is reported for  
385 the testing set of reactions unless indicated otherwise.

386 The gap between the within distribution and out-of-distribution performance is large in all of  
387 the cases studies (Fig. 5A). The accuracy drop for within distribution performance is small in  
388 most cases and generally mirrors the reduced amount of training data available. In contrast, the  
389 accuracy drop for the out-of-distribution reaction classes is much larger than can be explained from  
390 reduced training data and reflects the qualitative failure of these architectures to predict  $E_a$  for



391 reactions that differ from the training distribution in either the number and kind of bonds changes,  
392 molecularity, or size. We had hypothesized that the models would extrapolate better across some  
393 of these classes than others. This was incorrect—the transferability to unseen reaction populations  
394 is consistently poor, with accuracy often worse than the naïve mean-predicting model. We had also  
395 hypothesized that the multitask or transfer learning architectures might show better transferability,  
396 but this also wasn't the case (see Fig. S8 for the combined results with all architectures). All the  
397 EGAT architectures perform poorly in the out-of-distribution prediction tasks with only a marginal  
398 benefit from the architectures with access to  $\Delta H_r$ . With the exception of large reactants, Chemprop  
399 generally shows significantly lower errors than EGAT in out-of-distribution tasks. But neither  
400 architecture achieves out-of-distribution performance that could be relied on for applications.

401 The learning curves for these case studies also illustrate the limited transferability of the learned  
402 reaction representation to out-of-distribution prediction (Figs. 5B-D). These learning curves show  
403 the performance of the  $G \rightarrow E_a$  EGAT and Chemprop architectures on the validation split as they  
404 are trained. A normalized x-axis is used because the Chemprop and EGAT architectures train  
405 over differing numbers of epochs. The within distribution validation curves look typical, with  
406 an asymptotic approach to a saturation value that is near the testing set performance shown in  
407 Figure 5A. However, the out-of-distribution samples in the validation split show remarkably little  
408 reduction in error throughout the training. In some of the cases, a marginal initial reduction  
409 occurs that can be interpreted as the models learning generic reaction fingerprint features relevant  
410 to all reactions, but these curves still plateau at much larger values. Additional case studies were  
411 performed where models that were pretrained on one reaction class were then retrained on an  
412 excluded reaction class (Fig. S9). In all cases, these models show catastrophic forgetfulness after  
413 a single training epoch in predicting validation reactions from their original training class.

414 These case studies highlight the extremely limited transferability of these models to classes  
415 of reactions that are nevertheless expected to share many essential mechanistic features with the

416 training reactions. For example, many reaction mechanisms are expected to be conserved regardless  
417 of molecularity. Similarly, many small and large molecules should exhibit comparable activation  
418 energies for reactions with conserved mechanisms. The lack of transferability in these cases studies  
419 speaks to a gap in the current architectures. Based on the class-based reaction modeling of experts,  
420 this gap should be addressable, but it will require better regularization or additional architectural  
421 developments to promote mechanistic transferability.

## 422 4 Conclusions

423 This study has revisited the chemical graph to activation energy (G2Ea) prediction problem using  
424 an edge-featured graph attention (EGAT) architecture. This was motivated by the recent devel-  
425 opment of the relatively large ( $\sim 176\text{k}$ ) RGD1 organic chemistry reaction dataset, which enabled  
426 benchmarking across a broad swathe of reaction types and interrogation of out-of-distribution per-  
427 formance using case-studies with partial training on subsets of RGD1. Another motivation was to  
428 contribute an additional open-source architecture for other research groups to experiment with on  
429 their own. To the credit of the Chemprop developers, they have opened up their code for adapta-  
430 tion and comparison. Expanding the pool of publicly available models is critical to resolve where  
431 fundamental obstacles exist. The summary observations from these experiments are that it is  
432 relatively straightforward for G2Ea models to approach the irreducible error of RGD1—estimated  
433 to be  $\sim 4$  kcal/mol—but that out-of-distribution performance is often worse than naïve models.

434 Sources of irreducible error for the G2Ea prediction task warrant more detailed consideration  
435 moving forward. Incomplete conformational sampling results in the curation of TSs that are not  
436 the lowest barrier possible, and thus out-of-distribution for training and prediction. Likewise,  
437 unintended TSs (i.e., true saddle points that nonetheless correspond to reactions that are different  
438 from the labeled reactants and products) are commonly produced by TS search algorithms. Due to

439 the incorrect input features, unintended samples are intrinsically unpredictable by G2Ea models.  
440 Both sources of irreducible error were observed in the RGD1 dataset and are expected to be found in  
441 any other computational dataset of conformationally complex reactants without explicit protocols  
442 for mitigation. Although these errors are “irreducible” from the perspective of performance in  
443 the G2Ea task for a given dataset, they are reducible from the perspective of curating better  
444 datasets. For example, there are many possibilities for developing better filters for unintended  
445 transition states—from more complex models to using more informative featurizations for the  
446 unintended/intended classification task—but they have yet to be implemented. We estimated that  
447 the majority of prediction outliers in the current study were in fact unintended reactions and that  
448 perfect filtering would reduce the irreducible error for the G2Ea task by  $\sim 1 - 2$  kcal/mol with the  
449 remaining  $\sim 2 - 3$  kcal/mol being accounted for by incomplete conformational sampling.

450 The discussion of errors also bears on the relative advantage of formulating  $E_a$  prediction as  
451 a G2Ea task. For example, unintended transition states and biased conformational sampling are  
452 not intrinsic problems for some three-dimensional formulations of  $E_a$  prediction. For example, a  
453 model with a three dimensional reactant featurization could plausibly learn the contribution to  
454  $E_a$  associated with a particular conformer and thus reduce that source of error. Or consider an  
455  $E_a$  model that performs  $E_a$  estimation identically to quantum chemical approaches but instead  
456 uses ML atomic potentials.<sup>75</sup> Because such a model is formulated to predict reactive potential  
457 energy surfaces, it can learn from both intended and unintended TSs and all reactant and product  
458 conformations. This formulation may sound much more expensive than directly predicting  $E_a$   
459 from a pair of chemical graphs, but the advent of GPU-compatible routines for performing TS  
460 searches on ML-potentials may render such cost differences moot. These differing formulations of  
461 the  $E_a$  prediction problem can also be expected to affect model transferability. Better benchmarks  
462 will ultimately be required to resolve the accuracy vs cost Pareto front of different  $E_a$  prediction  
463 formulations being gestured toward by this discussion.

464 Within the context of G2Ea models, these results also suggest several pathways for improve-  
465 ment. First, the fact that the overall accuracy of all models uniformly approached the irreducible  
466 error of RGD1 shows that these models have spare complexity to learn broader classes of reactions.  
467 RGD1 only contains closed-shell neutral reactants containing CHON elements, but extensions to  
468 ionic, radical, and other elements are expected to be successful. Second, the current featurization  
469 does not directly consider the reaction conditions (e.g., solvent) or the availability of a catalyst.  
470 It is plausible that a condition fingerprint could be catenated to the reaction fingerprint to pre-  
471 dict how  $E_a$  would be modulated by environment, but to our knowledge the requisite data for  
472 a convincing attempt at this does not yet exist. Lastly, our current experiments found marginal  
473 advantage from using complementary information sources like  $\Delta H_r$  during training and prediction  
474 in the large data limit. More advantages might become apparent for datasets possessing a broader  
475 range of reaction mechanisms, datasets with lower irreducible error for the G2Ea task, or by using  
476 auxiliary information sources beyond  $\Delta H_r$ . These and adjacent opportunities suggest that the  
477 field is far from determining the ultimate performance of G2Ea models.

## 478 References

- 479 (1) Demirbas, A.; Arin, G. An Overview of Biomass Pyrolysis. *Energy Sources* **2002**, *24*, 471–482.
- 480 (2) Akbar Ali, M.; Violi, A. Reaction Pathways for the Thermal Decomposition of Methyl Bu-  
481 tanoate. *J. Org. Chem.* **2013**, *78*, 5898–5908.
- 482 (3) Kang, P.-L.; Shang, C.; Liu, Z.-P. Glucose to 5-Hydroxymethylfurfural: Origin of Site-  
483 Selectivity Resolved by Machine Learning Based Reaction Sampling. *J. Am. Chem. Soc.*  
484 **2019**, *141*, 20525–20536.

- 485 (4) Mohan, D.; Pittman Jr, C. U.; Steele, P. H. Pyrolysis of wood/biomass for bio-oil: A critical  
486 review. *Energy Fuels* **2006**, *20*, 848–889.
- 487 (5) Corma, A.; Iborra, S.; Velty, A. Chemical routes for the transformation of biomass into  
488 chemicals. *Chem. Rev.* **2007**, *107*, 2411–2502.
- 489 (6) Zhao, Q.; Savoie, B. M. Deep reaction network exploration of glucose pyrolysis. *Proc. Natl.*  
490 *Acad. Sci.* **2023**, *120*, e2305884120.
- 491 (7) Reymond, J.-L.; Awale, M. Exploring Chemical Space for Drug Discovery Using the Chemical  
492 Universe Database. *ACS Chem. Neurosci.* **2012**, *3*, 649–657.
- 493 (8) Virshup, A. M.; Contreras-García, J.; Wipf, P.; Yang, W.; Beratan, D. N. Stochastic Voyages  
494 into Uncharted Chemical Space Produce a Representative Library of All Possible Drug-Like  
495 Compounds. *J. Am. Chem. Soc.* **2013**, *135*, 7296–7303.
- 496 (9) Blakemore, D. C.; Castro, L.; Churcher, I.; Rees, D. C.; Thomas, A. W.; Wilson, D. M.;  
497 Wood, A. Organic synthesis provides opportunities to transform drug discovery. *Nat. Chem.*  
498 **2018**, *10*, 383–394.
- 499 (10) Montgomery, J. High-Throughput Discovery of New Chemical Reactions. *Science* **2011**, *333*,  
500 1387–1388.
- 501 (11) Heiskanen, S. K.; Kim, J.; Lucht, B. L. Generation and Evolution of the Solid Electrolyte  
502 Interphase of Lithium-Ion Batteries. *Joule* **2019**, *3*, 2322–2333.
- 503 (12) Takenaka, N.; Suzuki, Y.; Sakai, H.; Nagaoka, M. On Electrolyte-Dependent Formation of  
504 Solid Electrolyte Interphase Film in Lithium-Ion Batteries: Strong Sensitivity to Small Struc-  
505 tural Difference of Electrolyte Molecules. *J. Phys. Chem. C* **2014**, *118*, 10874–10882.

- 506 (13) Chamas, A.; Moon, H.; Zheng, J.; Qiu, Y.; Tabassum, T.; Jang, J. H.; Abu-Omar, M.;  
507 Scott, S. L.; Suh, S. Degradation Rates of Plastics in the Environment. *ACS Sustain. Chem.*  
508 *Eng.* **2020**, *8*, 3494–3511.
- 509 (14) Huang, Z.; Chen, M. S.; Wroch, C. P.; Markland, T. E.; Kanan, M. W. A framework for  
510 automated structure elucidation from routine NMR spectra. *Chem. Sci.* **2021**, *12*, 15329–  
511 15338.
- 512 (15) Firet, N. J.; Smith, W. A. Probing the Reaction Mechanism of CO<sub>2</sub> Electroreduction over  
513 Ag Films via Operando Infrared Spectroscopy. *ACS Catal.* **2017**, *7*, 606–612.
- 514 (16) Lansford, J. L.; Vlachos, D. G. Infrared spectroscopy data- and physics-driven machine learn-  
515 ing for characterizing surface microstructure of complex materials. *Nat. Commun.* **2020**, *11*,  
516 1513.
- 517 (17) Fine, J.; Kuan-Yu Liu, J.; Beck, A.; Alzarieni, K. Z.; Ma, X.; Boulos, V. M.; Kenttämaa, H. I.;  
518 Chopra, G. Graph-based machine learning interprets and predicts diagnostic isomer-selective  
519 ion–molecule reactions in tandem mass spectrometry. *Chem. Sci.* **2020**, *11*, 11849–11858.
- 520 (18) Fang, Y.; Li, J.; Chen, Y.; Lu, Q.; Yang, H.; Wang, X.; Chen, H. Experiment and Modeling  
521 Study of Glucose Pyrolysis: Formation of 3-Hydroxy- $\gamma$ -butyrolactone and 3-(2H)-Furanone.  
522 *Energy Fuels* **2018**, *32*, 9519–9529.
- 523 (19) Zhao, Q.; Savoie, B. M. Simultaneously improving reaction coverage and computational cost  
524 in automated reaction prediction tasks. *Nat. Comput. Sci.* **2021**, *1*, 479–490.
- 525 (20) Zhao, Q.; Savoie, B. M. Algorithmic explorations of unimolecular and bimolecular reaction  
526 spaces. *Angew. Chem., Int. Ed.* **2022**, *61*, e202210693.

- 527 (21) Barter, D.; Clark Spotte-Smith, E. W.; Redkar, N. S.; Khanwale, A.; Dwaraknath, S.; Pers-  
528 son, K. A.; Blau, S. M. Predictive stochastic analysis of massive filter-based electrochemical  
529 reaction networks. *Digital Discovery* **2023**, *2*, 123–137.
- 530 (22) Coley, C. W.; Green, W. H.; Jensen, K. F. Machine learning in computer-aided synthesis  
531 planning. *Acc. Chem. Res.* **2018**, *51*, 1281–1289.
- 532 (23) Suleimanov, Y. V.; Green, W. H. Automated Discovery of Elementary Chemical Reaction  
533 Steps using Freezing String and Berny Optimization Methods. *J. Chem. Theory Comput.*  
534 **2015**, *11*, 4248–4259.
- 535 (24) Zimmerman, P. M. Automated discovery of chemically reasonable elementary reaction steps.  
536 *J. Comput. Chem.* **2013**, *34*, 1385–1392.
- 537 (25) Vernuccio, S.; Broadbelt, L. J. Discerning complex reaction networks using automated gen-  
538 erators. *AIChE J.* **2019**, *65*, e16663.
- 539 (26) Ásgeirsson, V.; Birgisson, B. O.; Bjornsson, R.; Becker, U.; Neese, F.; Riplinger, C.;  
540 Jónsson, H. Nudged Elastic Band Method for Molecular Reactions Using Energy-Weighted  
541 Springs Combined with Eigenvector Following. *J. Chem. Theory Comput.* **2021**, *17*, 4929–  
542 4945.
- 543 (27) Henkelman, G.; Uberuaga, B. P.; Jónsson, H. A Climbing Image Nudged Elastic Band Method  
544 for Finding Saddle Points and Minimum Energy Paths. *J. Chem. Phys.* **2000**, *113*, 9901–9904.
- 545 (28) Zimmerman, P. M. Growing string method with interpolation and optimization in internal  
546 coordinates: Method and examples. *J. Chem. Phys.* **2013**, *138*, 184102.
- 547 (29) Zimmerman, P. Reliable transition state searches integrated with the growing string method.  
548 *J. Chem. Theory Comput.* **2013**, *9*, 3043–3050.

- 549 (30) Maeda, S.; Harabuchi, Y.; Takagi, M.; Taketsugu, T.; Morokuma, K. Artificial force induced  
550 reaction (AFIR) method for exploring quantum chemical potential energy surfaces. *Chem.*  
551 *Rec.* **2016**, *16*, 2232–2248.
- 552 (31) Simm, G. N.; Reiher, M. Context-driven exploration of complex chemical reaction networks.  
553 *J. Chem. Theory Comput.* **2017**, *13*, 6108–6119.
- 554 (32) Abdelfatah, K.; Yang, W.; Vijay Solomon, R.; Rajbanshi, B.; Chowdhury, A.; Zare, M.;  
555 Kundu, S. K.; Yonge, A.; Heyden, A.; Terejanu, G. Prediction of Transition-State Energies  
556 of Hydrodeoxygenation Reactions on Transition-Metal Surfaces Based on Machine Learning.  
557 *J. Phys. Chem. C* **2019**, *123*, 29804–29810.
- 558 (33) Wang, S. et al. Universal transition state scaling relations for (de)hydrogenation over transi-  
559 tion metals. *Phys. Chem. Chem. Phys.* **2011**, *13*, 20760–20765.
- 560 (34) Wang, S.; Vorotnikov, V.; Sutton, J. E.; Vlachos, D. G. Brønsted–Evans–Polanyi and Tran-  
561 sition State Scaling Relations of Furan Derivatives on Pd(111) and Their Relation to Those  
562 of Small Molecules. *ACS Catal.* **2014**, *4*, 604–612.
- 563 (35) Van Santen, R. A.; Neurock, M.; Shetty, S. G. Reactivity Theory of Transition-Metal Surfaces:  
564 a Brønsted-Evans-Polanyi Linear Activation Energy-Free-Energy analysis. *Chem. Rev.* **2010**,  
565 *110*, 2005–2048.
- 566 (36) Bragato, M.; von Rudorff, G. F.; von Lilienfeld, O. A. Data enhanced Hammett-equation:  
567 reaction barriers in chemical space. *Chem. Sci.* **2020**, *11*, 11859–11868.
- 568 (37) Rogers, D.; Hahn, M. Extended-Connectivity Fingerprints. *J. Chem. Inf. Model.* **2010**, *50*,  
569 742–754.



- 570 (38) Ismail, I.; Robertson, C.; Habershon, S. Successes and challenges in using machine-learned  
571 activation energies in kinetic simulations. *J. Chem. Phys.* **2022**, *157*, 014109.
- 572 (39) Li, X.; Zhang, S.-Q.; Xu, L.-C.; Hong, X. Predicting Regioselectivity in Radical C-H Func-  
573 tionalization of Heterocycles through Machine Learning. *Angew. Chem., Int. Ed.* **2020**, *59*,  
574 13253–13259.
- 575 (40) Choi, S.; Kim, Y.; Kim, J. W.; Kim, Z.; Kim, W. Y. Feasibility of Activation Energy Predic-  
576 tion of Gas-Phase Reactions by Machine Learning. *Chem. Eur. J.* **2018**, *24*, 12354–12358.
- 577 (41) Zhao, Q.; Vaddadi, S. M.; Woulfe, M.; Ogunfowora, L. A.; Garimella, S. S.; Isayev, O.;  
578 Savoie, B. M. Comprehensive exploration of graphically defined reaction spaces. *Sci. Data*  
579 **2023**, *10*, 145.
- 580 (42) Rudorff, G. F.; Heinen, S. N.; Bragato, M.; Lilienfeld, O. A. Thousands of reactants and  
581 transition states for competing E2 and SN2 reactions. *Mach. Learn.: Sci. Technol* **2020**, *1*,  
582 045026.
- 583 (43) Grambow, C. A.; Pattanaik, L.; Green, W. H. Reactants, products, and transition states of  
584 elementary chemical reactions based on quantum chemistry. *Sci. Data* **2020**, *7*, 137.
- 585 (44) Spiekermann, K. A.; Pattanaik, L.; Green, W. H. Fast Predictions of Reaction Barrier Heights:  
586 Toward Coupled-Cluster Accuracy. *J. Phys. Chem. A* **2022**, *126*, 3976–3986.
- 587 (45) Grambow, C. A.; Pattanaik, L.; Green, W. H. Deep Learning of Activation Energies. *J. Phys.*  
588 *Chem. Lett.* **2020**, *11*, 2992–2997.
- 589 (46) Heid, E.; Green, W. H. Machine Learning of Reaction Properties via Learned Representations  
590 of the Condensed Graph of Reaction. *J. Chem. Inf. Model.* **2022**, *62*, 2101–2110.

- 591 (47) Heinen, S.; von Rudorff, G. F.; von Lilienfeld, O. A. Toward the Design of Chemical Reactions:  
592 Machine Learning Barriers of Competing Mechanisms in Reactant Space. *J. Chem. Phys.*  
593 **2021**, *155*, 064105.
- 594 (48) Heinen, S.; von Rudorff, G. F.; von Lilienfeld, O. A. Transition state search and geometry  
595 relaxation throughout chemical compound space with quantum machine learning. *J. Chem.*  
596 *Phys.* **2022**, *157*, 221102.
- 597 (49) Zhao, Q.; Anstine, D. M.; Isayev, O.; Savoie, B. M.  $\Delta^2$  machine learning for reaction property  
598 prediction. *Chem. Sci.* **2023**,
- 599 (50) Marques, E.; de Gendt, S.; Pourtois, G.; van Setten, M. J. Improving Accuracy and Trans-  
600 ferability of Machine Learning Chemical Activation Energies by Adding Electronic Structure  
601 Information. *J. Chem. Inf. Model.* **2023**, *63*, 1454–1461.
- 602 (51) Friederich, P.; dos Passos Gomes, G.; De Bin, R.; Aspuru-Guzik, A.; Balcells, D. Machine  
603 learning dihydrogen activation in the chemical space surrounding Vaska's complex. *Chem.*  
604 *Sci.* **2020**, *11*, 4584–4601.
- 605 (52) Migliaro, I.; Cundari, T. R. Density Functional Study of Methane Activation by Frustrated  
606 Lewis Pairs with Group 13 Trihalides and Group 15 Pentahalides and a Machine Learning  
607 Analysis of Their Barrier Heights. *J. Chem. Inf. Model.* **2020**, *60*, 4958–4966.
- 608 (53) Mikami, K. Interactive-quantum-chemical-descriptors enabling accurate prediction of an ac-  
609 tivation energy through machine learning. *Polymer* **2020**, *203*, 122738.
- 610 (54) Döntgen, M.; Fenard, Y.; Heufer, K. A. Atomic Partial Charges as Descriptors for Barrier  
611 Heights. *J. Chem. Inf. Model.* **2020**, *60*, 5928–5931.

- 612 (55) Jorner, K.; Brinck, T.; Norrby, P.-O.; Buttar, D. Machine learning meets mechanistic mod-  
613 elling for accurate prediction of experimental activation energies. *Chem. Sci.* **2021**, *12*, 1163–  
614 1175.
- 615 (56) Gerwen, P. v.; Wodrich, M. D.; Laplaza, R.; Corminboeuf, C. Reply to Comment on ‘Physics-  
616 based representations for machine learning properties of chemical reactions’. *Mach. learn.:  
617 Sci. Technol.* **2023**, *4*, 048002.
- 618 (57) Spiekermann, K. A.; Stuyver, T.; Pattanaik, L.; Green, W. H. Comment on ‘Physics-based  
619 representations for machine learning properties of chemical reactions’. *Mach. learn.: Sci.  
620 Technol.* **2023**, *4*, 048001.
- 621 (58) Pracht, P.; Bohle, F.; Grimme, S. Automated exploration of the low-energy chemical space  
622 with fast quantum chemical methods. *Phys. Chem. Chem. Phys.* **2020**, *22*, 7169–7192.
- 623 (59) Zhao, Q.; Hsu, H.-H.; Savoie, B. Conformational Sampling for Transition State Searches on  
624 a Computational Budget. *J. Chem. Theory Comput.* **2022**, *18*, 3006–3016.
- 625 (60) Veličković, P.; Cucurull, G.; Casanova, A.; Romero, A.; Liò, P.; Bengio, Y. Graph Attention  
626 Networks. International Conference on Learning Representations. 2018.
- 627 (61) Kamiński, K.; Ludwiczak, J.; Jasiński, M.; Bukala, A.; Madaj, R.; Szczepaniak, K.; Dunin-  
628 Horkawicz, S. Rossmann-toolbox: a deep learning-based protocol for the prediction and design  
629 of cofactor specificity in Rossmann fold proteins. *Briefings Bioinf.* **2021**, *23*, bbab371.
- 630 (62) Wang, Z.; Chen, J.; Chen, H. EGAT: Edge-Featured Graph Attention Network. Artificial  
631 Neural Networks and Machine Learning – ICANN 2021. Cham, 2021; pp 253–264.
- 632 (63) Yang, K.; Swanson, K.; Jin, W.; Coley, C.; Eiden, P.; Gao, H.; Guzman-Perez, A.; Hopper, T.;  
633 Kelley, B.; Mathea, M.; Palmer, A.; Settels, V.; Jaakkola, T.; Jensen, K.; Barzilay, R. Analyz-

- 634 ing Learned Molecular Representations for Property Prediction. *J. Chem. Inf. Model.* **2019**,  
635 *59*, 3370–3388.
- 636 (64) Evans, M. G.; Polanyi, M. Inertia and driving force of chemical reactions. *Trans. Faraday*  
637 *Soc.* **1938**, *34*, 11–24.
- 638 (65) Kim, S. Pubchem in 2021: new data content and improved web interfaces. *Nucleic Acids Res.*  
639 **2021**, *49*, D1388–D1395.
- 640 (66) Schlegel, H. B. Optimization of equilibrium geometries and transition structures. *J. Comput.*  
641 *Chem.* **1982**, *3*, 214–218.
- 642 (67) Maeda, S.; Harabuchi, Y.; Ono, Y.; Taketsugu, T.; Morokuma, K. Intrinsic reaction coor-  
643 dinate: Calculation, bifurcation, and automated search. *Int. J. Quantum Chem.* **2015**, *115*,  
644 258–269.
- 645 (68) Bannwarth, C.; Ehlert, S.; Grimme, S. GFN2-xTB—An accurate and broadly parametrized  
646 self-consistent tight-binding quantum chemical method with multipole electrostatics and  
647 density-dependent dispersion contributions. *J. Chem. Theory Comput.* **2019**, *15*, 1652–1671.
- 648 (69) Frisch, M. J. et al. Gaussian 16 Revision C.01. **2016**, Gaussian Inc. Wallingford CT.
- 649 (70) Stephens, P. J.; Devlin, F. J.; Chabalowski, C. F.; Frisch, M. J. Ab initio calculation of  
650 vibrational absorption and circular dichroism spectra using density functional force fields. *J.*  
651 *Phys. Chem.* **1994**, *98*, 11623–11627.
- 652 (71) Schäfer, A.; Huber, C.; Ahlrichs, R. Fully optimized contracted Gaussian basis sets of triple  
653 zeta valence quality for atoms Li to Kr. *J. Chem. Phys.* **1994**, *100*, 5829–5835.

- 654 (72) Grimme, S.; Antony, J.; Ehrlich, S.; Krieg, H. A consistent and accurate ab initio parametriza-  
655 tion of density functional dispersion correction (DFT-D) for the 94 elements H-Pu. *J. Chem.*  
656 *Phys* **2010**, *132*, 154104.
- 657 (73) Chen, T.; Guestrin, C. XGBoost: A Scalable Tree Boosting System. Proceedings of the 22nd  
658 ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. 2016;  
659 pp 785–794.
- 660 (74) Heller, S. R.; McNaught, A.; Pletnev, I.; Stein, S.; Tchekhovskoi, D. InChI, the IUPAC  
661 International Chemical Identifier. *J. ChemInf.* **2015**, *7*, 1–34.
- 662 (75) Dylan Anstine, O. I., Roman Zabatyuk AIMNet2: A Neural Network Potential to Meet your  
663 Neutral, Charged, Organic, and Elemental-Organic Needs. *ChemRxiv* **2023**,

## 664 **5 Data and Code Availability**

665 The authors declare that the data supporting the findings of this study are available within the  
666 paper and its supplementary information files. Pretained EGAT models and corresponding training  
667 scripts are available on Github (XXX To be filled upon publication XXX). Training and testing  
668 set splits and all relevant data for the figures are detailed on Figshare (XXX To be filled upon  
669 publication XXX).

## 670 **Acknowledgements**

671 The work performed by S.M.V., Q.Z., and B.M.S was made possible by the Office of Naval Research  
672 (ONR) through support provided by the Energetic Materials Program (MURI grant number:

673 N00014-21-1-2476, Program Manager: Dr. Chad Stoltz). B.M.S also acknowledges partial support  
674 for this work from the Purdue Process Safety and Assurance Center.

## 675 **Author contributions statement**

676 Q.Z. and B.M.S. conceived and designed the study. Q.Z. developed the initial iterations of the  
677 model. Q.Z. and S.M.V. pre-processed the data into training, testing, and validation sets. S.M.V.  
678 performed the tests and analyzed the data. All authors prepared the manuscript.

## 679 **Competing interests**

680 The authors declare no competing interests.