

# Machine Learning Models Capable of Chemical Deduction for Identifying Reaction Products

Tianfan Jin, Qiyuan Zhao, Andrew B. Schofield, and Brett M. Savoie\*

*Davidson School of Chemical Engineering, Purdue University, West Lafayette, IN, 47906*

E-mail: bsavoie@purdue.edu

## Abstract

Deductive solution strategies are required in prediction scenarios that are under determined, when contradictory information is available, or more generally wherever one-to-many non-functional mappings occur. In contrast, most contemporary machine learning (ML) in the chemical sciences is inductive learning from example, with a fixed set of features. Chemical workflows are replete with situations requiring deduction, including many aspects of lab automation and spectral interpretation. Here, a general strategy is described for designing and training machine learning models capable of deduction that consists of combining individual inductive models into a larger deductive network. The training and testing of these models is demonstrated on the task of deducing reaction products from a mixture of spectral sources. The resulting models are capable of distinguishing between intended and unintended reaction outcomes and identifying starting material based on a mixture of spectral sources. The models are also capable of performing well on tasks that they were not directly trained on, like predicting minor products from named organic chemistry reactions, identifying reagents and isomers as plausible impurities, and handling missing or conflicting information. A new

dataset of 1,124,043 simulated spectra that were generated to train these models is also distributed with this work. These findings demonstrate that deductive bottlenecks for chemical problems are not fundamentally insuperable for ML models.

1 Product identification is a central task in every reaction development workflow.<sup>1-5</sup> There is no  
2 standardized solution to this problem, with practices ranging from separation and crystallization  
3 for unequivocal identification, to using a mixture of analytical information sources (e.g., mass spec-  
4 trometry (MS), nuclear magnetic resonance (NMR), infrared spectroscopy (IR), etc.) and general  
5 reactivity knowledge to distinguish between plausible products. The lack of standardization reflects  
6 that product identification is typically underdetermined by simple knowledge of the reactants and  
7 conditions. For example, a new reaction may yield a complex product mixture that requires several  
8 iterations of characterization and interpretation to fully identify, and even putatively established  
9 reactions can yield unexpected products if a hot-plate fails or a starting material has an impurity.  
10 Underdetermination also occurs because most analytical characterizations only provide partial or  
11 indirect structural information, and a particular analytical method may yield decisive information  
12 for identifying one product but not another.<sup>6-9</sup> For these reasons, the state-of-the-art for general  
13 product identification remains manual expert interpretation of multiple information sources.

14 Product identification is a member of a larger group of deduction problems that are common  
15 in the chemical sciences (Fig. 1A). In deductive scenarios, external information is used to restrict  
16 the potential solution space when making a prediction. Deduction is required for underdetermined  
17 problems or when there is a mixture of competing information sources. In contrast, most ma-  
18 chine learning (ML) in chemistry is inductive, learning from example, with a fixed set of input  
19 features.<sup>10-13</sup> In the case of product identification, deduction takes the form of using established  
20 reactivity relationships to narrow the solution space to a small number of potential products that  
21 can then be inductively distinguished using one or more analytical spectra. More generally, deduc-  
22 tion is needed whenever a non-functional one-to-many relationship exists between input features

23 and prediction targets. In the context of ML, this distinction is critical, because regardless of their  
 24 complexity, neural networks are incapable of circumventing the information limitations posed by  
 25 non-functional mappings.

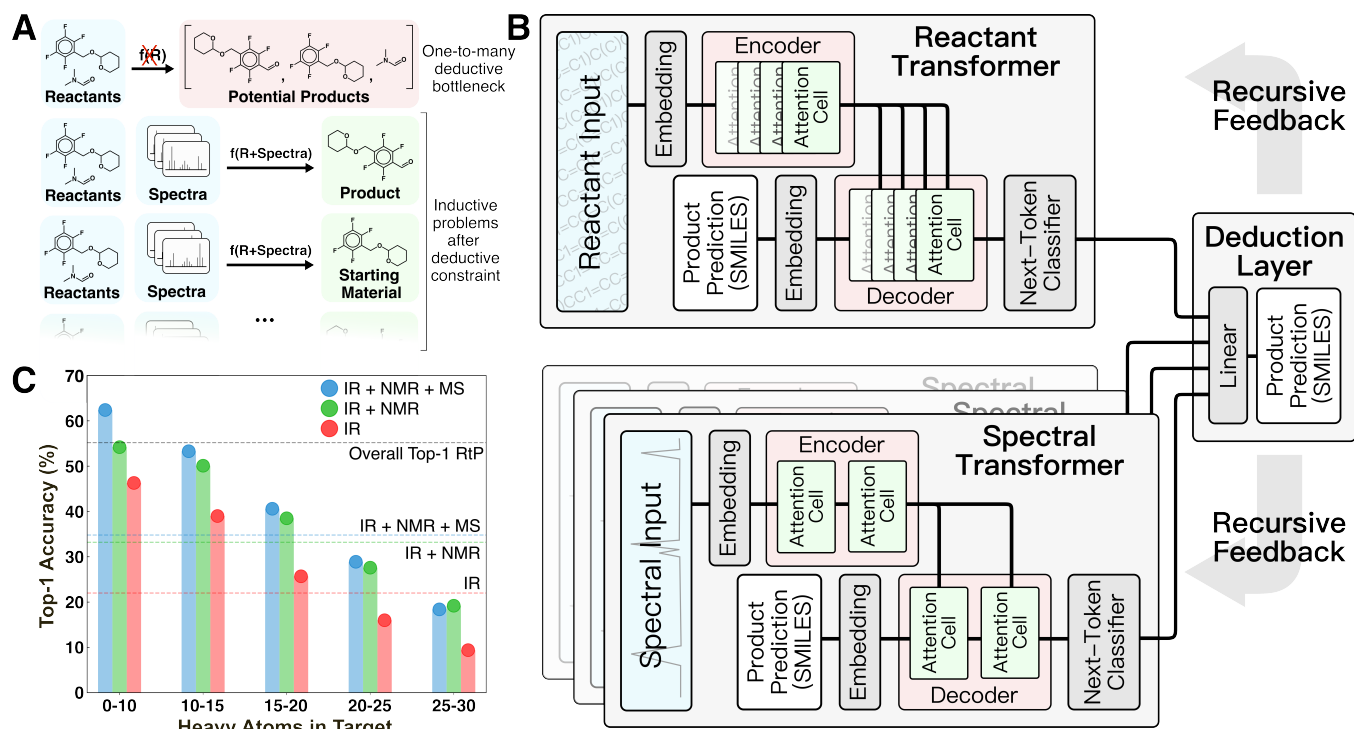


Figure 1: Overview of deductive architecture and bottleneck for product identification. (A) Illustration of the general non-functional one-to-many relationship between reactant information and some potential species that can be found as intended and unintended products. (B) Deductive super-network consisting of a reactant to product (RtP) transformer and one or more spectrum to structure (StS) transformers combined by a terminal linear layer. The model predicts product SMILES in probabilistic token-by-token fashion. (C) Top-1 accuracy of StS models in predicting structures from the testing set with an increasing number of heavy atoms. The dotted lines indicate the overall top-1 accuracy of each model on the whole testing set.

26 The motivation for the current study was to develop a ML-framework capable of emulating  
 27 expert deduction to perform product identification based on a flexible mixture of spectral input  
 28 sources. We hypothesized that deduction would be an emergent property of a super-network  
 29 composed of individual task-specific inductive neural networks and a decision-making layer for

30 weighing task-specific evidence (Fig. 1B). This idea was directly motivated by the manual analog  
31 of interpreting individual spectra to obtain derived information (e.g., identifying the presence of  
32 certain functional groups from IR or a probable chemical formula from MS) then forming structural  
33 hypotheses from comparisons of this derived information.

34 Here, we experimented with combining up to four task-specific transformers for ingesting re-  
35 actant/reagent information and IR, H-NMR, and electron-ionization (EI) MS spectra, respec-  
36 tively. The overall architecture takes reactant/reagent simplified molecular-input line-entry sys-  
37 tem (SMILES) strings and one or more analytical spectra as inputs and probabilistically decodes  
38 the product SMILES as an output in recursive token-by-token fashion.<sup>14</sup> Each task-specific trans-  
39 former provides a probabilistic prediction of the next token in the product that informs a final  
40 linear deduction layer (see Methods). This architecture provides two sources of deductive coupling  
41 between the transformers. The first is the straightforward probability reweighting that happens  
42 in the final linear deduction layer, which provides the opportunity for one or more of the trans-  
43 formers to form a consensus over the other transformer(s). The second is through the recursive  
44 token-by-token decoding by which the product prediction is made. Because the partially decoded  
45 product string is used as an input to each transformer during inference, it is possible for control  
46 to shift between transformers for different portions of the decoding (e.g., one may dominate the  
47 scaffold, while another dominates predictions of certain functional groups). In this way, the trans-  
48 formers can dynamically provide deductive constraints on each other during different portions of  
49 the decoding.

50 The deduction models were trained and tested on 299,658 reactions taken from the Lowe patent  
51 dataset after filtering (see Methods).<sup>15,16</sup> Artificial EI-MS, H-NMR, and IR spectra were generated  
52 for all products, reactants, and reagents due to the unavailability of suitable training data for this  
53 task. To turn this into a deductive product identification task, the dataset was augmented with  
54 null reactions that corresponded to obtaining starting material from the reaction instead of the

55 expected product. The final dataset consisted of 299,658 real reactions and 146,672 null reactions,  
56 that were split using a 80:10:10 training, validation, testing split while ensuring that there were  
57 no prediction targets shared between the splits. All accuracies are reported for the testing set.

58 Prediction baselines for this task were set by training analogous transformer models on the  
59 reactant-to-product (RtP) and spectrum-to-structure (StS) prediction tasks (Fig. 1C). The RtP  
60 model exhibits an obvious deductive bottleneck in this task, since a given reactant can map to  
61 either the expected product or starting material(s). The RtP model was trained only to predict  
62 the expected products, because attempts to train with null reactions in the training data led to  
63 confusion due to the one-to-many relationship between inputs and targets. The RtP model's top-1  
64 accuracy of  $\sim 55\%$  reflects a combined top-1 accuracy of  $\sim 0.6\%$  on null reactions and  $\sim 84.5\%$  on  
65 real reactions in the testing set. The latter result is comparable to the state-of-the-art RtP mod-  
66 els.<sup>17,18</sup> Several StS models were trained with different combinations of spectral transformers (IR,  
67 IR+NMR, and IR+NMR+MS models in Fig. 1B). The StS models exhibit lower overall perfor-  
68 mance than the RtP model, with a top-1 accuracy of  $\sim 35\%$  for the best model (IR+NMR+MS).  
69 The accuracies monotonically increase with the number of spectral sources used in the prediction  
70 and monotonically decrease with the molecular size of the prediction target. Although the de-  
71 ductive bottleneck is less obvious, it is qualitatively expected that spectral uniqueness decreases  
72 with molecular size (e.g., the structural isomers of large molecules often cannot be distinguished  
73 by these spectra). These accuracies favorably compare with recently published StS models that  
74 also exhibit relatively low performance for large molecules.<sup>7,19,20</sup> Notably, groups have reported  
75 StS accuracies that significantly improve when the molecular formula is supplied to the model  
76 in addition to the spectra.<sup>20</sup> Although it has not been identified as such, this is an elementary  
77 deductive constraint.

78 To test the hypothesis that combining a RtP transformer with one or more StS transformers  
79 circumvents the deductive bottleneck in the product identification task, the top-1 and top-5 testing

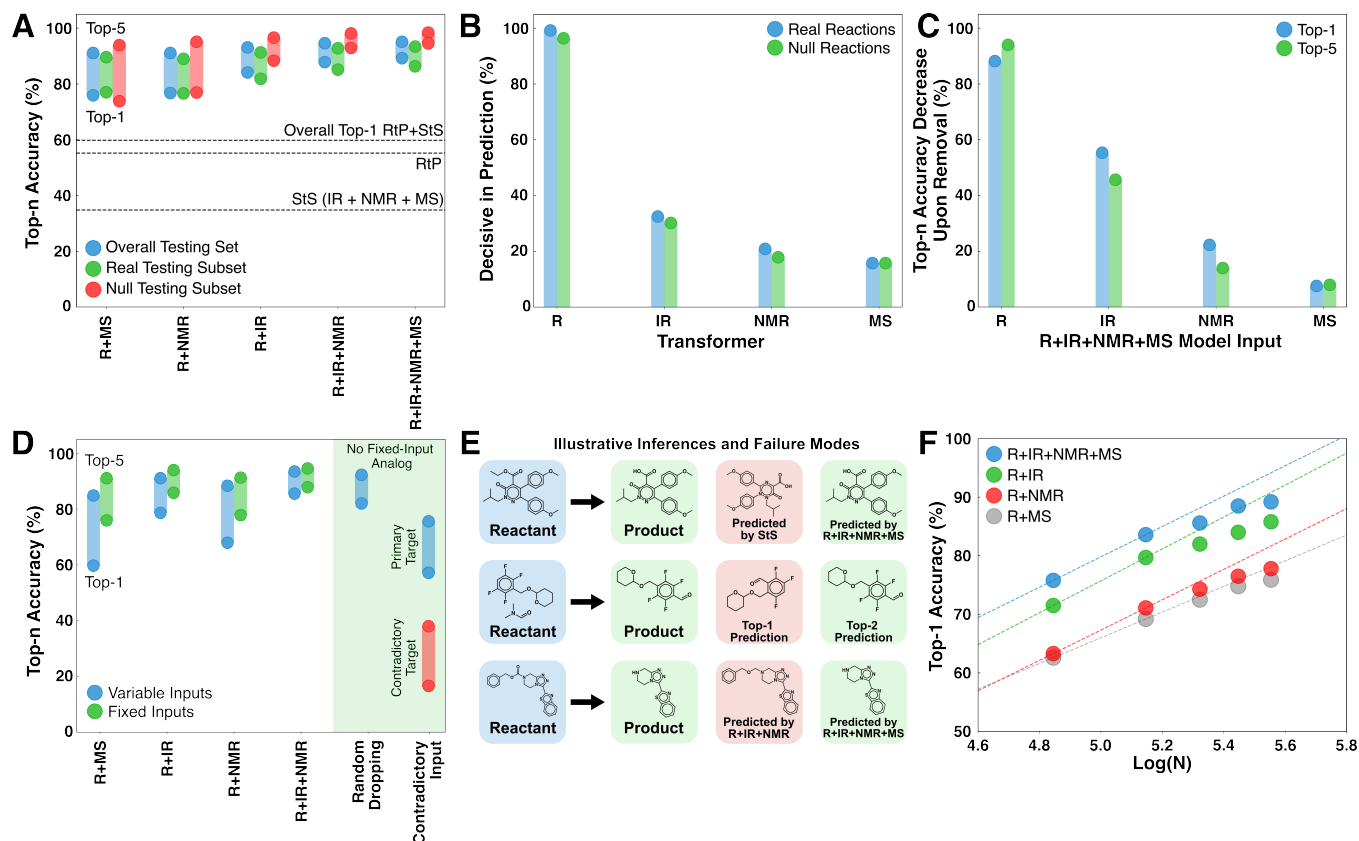


Figure 2: Overview of deductive performance in product identification tasks. (A) Comparison of several reactant+spectrum deductive models with RtP and StS models. The RtP+StS result corresponds to the accuracy obtained by combining the correct predictions from both models. (B) The fraction of products for which each transformer provides decisive input on at least one token. Multiple transformers can provide decisive contributions to a given product and a consensus results in no transformer being decisive, so the sum does not equal unity. (C) The reduction in top-n accuracy on the testing set upon zeroing out the input to the indicated transformer. (D) Comparison of a R+IR+NMR+MS model trained with missing spectra (blue) with the corresponding fixed input models (green). The cases on the right correspond to the performance with random dropping of one spectral input and supplying a contradictory spectrum (i.e., of starting material or a real product) to one of the spectral transformers. The red bars correspond to the fraction of cases where the contradictory species corresponding to the supplied spectrum was predicted in the top-n structures. (E) Three illustrative comparisons of the inferences of different models. (F) The convergence of the accuracy with respect to the number of training data on each of the deduction models.

80 accuracies of the deduction models were compared with the RtP and StS results (Fig. 2A). All  
81 the deduction models (even those with fewer spectral inputs) outperform the RtP and StS models  
82 by  $\sim 20\%$ , showing a qualitative difference between the inductive and deductive architectures. To  
83 clearly illustrate the non-linear impact of combining general reaction knowledge and the spectral  
84 information within a single model, we also calculated the top-1 accuracy of a hypothetical RtP+StS  
85 model that combines the correct predictions of the two separate models (line in Fig. 2A). Despite  
86 this generous accuracy calculation, the best deduction model still outperforms the RtP+StS model  
87 by 29%, illustrating the non-additive coupling between the reactant and spectral transformers.  
88 The deductive models also show no significant accuracy difference between predicting starting  
89 material versus expected products. This confirms that the reactant knowledge provided by the  
90 RtP transformer also assists with identifying starting material when incorporated within the larger  
91 deductive network.

92 The deductive architecture was motivated by the hypothesis that predictive control might  
93 switch between transformers during the token-by-token product decoding. To directly test this,  
94 the probability vectors produced by the transformers were individually zeroed out during inference  
95 to test whether the most probable overall token predicted by the model changed. If such a swap  
96 occurred for at least one token in a product, then the transformer was considered decisive in that  
97 decoding (Fig. 2B). The reactant transformer was found to be decisive for at least one token in  
98 over 95% of products, followed by the IR transformer at  $\sim 30\%$ . The lower decisiveness of the  
99 spectral transformers at least partially reflects their tendency to form a consensus and therefore  
100 not be individually decisive. For example, the decisiveness of the IR rises in the R+IR model to  
101 58% and 78% on real and null testing reactions, respectively. Approximately half of the products in  
102 the testing set had two or more decisive transformers from the R+IR+NMR+MS model involved  
103 in their decoding (Fig. S3). The mode decoding behavior is to switch between a consensus for the  
104 majority of the tokens (60-80%) and one or more decisive predictions for a minority of the tokens

105 (20-40%) (Fig. S4). This is strong support for the mechanism of dynamic deductive constraints  
106 being supplied by the different transformers during the token-by-token inference cycle.

107 To investigate the overall importance of the different input sources, the accuracy loss upon  
108 zeroing out each feature was averaged across the testing data (Fig. 2C). Given the stochastic  
109 nature of the decoding, a given input can influence a prediction even if it is not decisive for any  
110 particular token. Conversely, even if a transformer is decisive for a particular token, the flexibility  
111 of SMILES in decoding the same structure multiple ways means that a correct prediction may still  
112 be possible absent that transformer. The accuracy contributions roughly mirror the decisiveness  
113 of each transformer (Fig. 2B). In the case of IR, the influence on accuracy is  $\sim 20\%$  larger than  
114 the decisiveness measure, whereas for R, NMR, and MS it is marginally smaller. We interpret the  
115 relative contributions of the different spectra to reflect the simulation accuracy rather than the  
116 intrinsic information content of each spectral source. Nevertheless, there are many cases where  
117 even EI-MS makes decisive contributions to top predictions.

118 Several additional tests were performed to interrogate the ability of the deductive models to  
119 operate in scenarios of incomplete information and even contradictory information (Fig. 2D). For  
120 these trials, a version of the R+IR+NMR+MS model was trained from scratch using a ten percent  
121 random chance of dropping each spectral input based on the hypothesis that this would reduce  
122 the model reliance on consensus formation (see Methods). First, we tested the performance of  
123 this model in situations where one or more spectral inputs were unavailable. The performance of  
124 the model monotonically decreases on the testing set as spectral information is removed, but the  
125 top-1 and top-5 performance remain comparable to the models with fixed inputs (e.g., comparing  
126 R+IR+NMR+MS when deprived of IR and NMR data against the R+MS model). The perfor-  
127 mance remains comparably high in the case where the spectrum being removed is randomized, and  
128 for which there is no analog among the fixed input models. These trials show that the deductive  
129 architecture is capable of basing predictions on a flexible number of input sources, analogous to



130 the situation in product identification when spectra arrive asynchronously or may be unavailable  
131 for a given analyte (e.g., EI-MS may not be available for large molecules).

132 The R+IR+NMR+MS model trained with missing spectra was also tested in situations with  
133 contradictory information by supplying one of the spectral transformers at random with a con-  
134 tradictory spectrum (either starting material or real product) from the others (Fig. 2D, right).  
135 The performance in this case is lower than the situation where the model is simply deprived of a  
136 spectrum; nevertheless, the model shows the capacity to form a consensus that overrules the pre-  
137 dictions of the misinformed transformer. Remarkably, the model still predicts the contradictory  
138 species in the top-5 in nearly 40% cases. Although unanticipated, this behavior is more consis-  
139 tent with the supplied evidence than if the model never predicted the contradictory species. This  
140 also provides encouraging evidence that this architecture might be extended to predicting product  
141 mixtures. For example, a binary mixture of species with large differences in ionization efficiency  
142 or oscillator strengths could present similarly to the contradictory use case.

143 Inspection of some specific testing set examples illustrates the various ways that information is  
144 being used by the model (Fig. 2E). The first example shows a case where the IR+NMR+MS StS  
145 model fails for a relatively large product molecule, whereas the R+IR+NMR+MS model correctly  
146 predicts the product. This improvement reflects the transferable knowledge about organic reactions  
147 imparted by the reactant transformer. The second example shows a case where the deduction model  
148 fails to predict a product as top-1, but includes it as a top-5 prediction. This example is typical of  
149 many of the inaccurate predictions, where the model predicts structural isomers or molecules with  
150 similar scaffolds that are difficult to distinguish spectrally.  $\sim 18\%$  of the R+IR+NMR+MS top-  
151 1 mispredictions are structural isomers of the target. The third example shows a case where the  
152 R+IR+NMR model fails to predict a product as top-5 but the R+IR+NMR+MS model predicts it  
153 as a top choice. This case illustrates the complementary information supplied by the MS, despite  
154 it exhibiting the lowest overall decisiveness and accuracy contribution among the investigated

155 spectra.

156 A major data curation effort was required to train these models; nevertheless the accuracy  
157 versus training data size curves for the various models make it clear that there is additional scope  
158 for improvement (Fig. 2F). All of the models show clear evidence of saturation that we attribute  
159 to two factors. The first is that the performance of the models in identifying real products is  
160 already approaching the probable irreducible error of the underlying patent-sourced reaction data  
161 (i.e., many of the expected product labels are likely incorrect and cannot be accurately predicted  
162 regardless of having more data). The second potential source of saturation is the use of simulated  
163 spectra for these models. It is possible that real spectra would exhibit more information and  
164 saturate later.

165 Because these models were only trained on predicting starting material and major products,  
166 it was unclear how their performance would translate to predicting the products of side-reactions  
167 or other off-target species. We curated two external testing datasets, REAGENT and MULTI  
168 (see Methods), to test this (Fig. 3). The REAGENT dataset is made of 3262 reactions where  
169 the prediction target was a reagent rather than the starting material or expected product, as in  
170 the training data (see Methods). Reagent identification was an untrained task for these models  
171 and all reagents were unseen as prediction targets during training. The performance trend for  
172 reagent prediction is similar to the main testing cases, with a monotonic decrease in accuracy  
173 as spectral sources are removed and a baseline accuracy that is above the best StS model (Fig.  
174 3A). The accuracy is still reduced overall, as is expected given the difference between the training  
175 task and this task, but nevertheless the transferability to an unseen task is excellent. The RtP  
176 model is not compared here because it has  $\sim 0\%$  accuracy on this task, which is a reminder of the  
177 qualitative difference between the deductive and inductive architectures despite the decisiveness  
178 of the reactant transformer in the deductive architecture.

179 The capacity of the models to predict minor products was tested on the MULTI dataset of

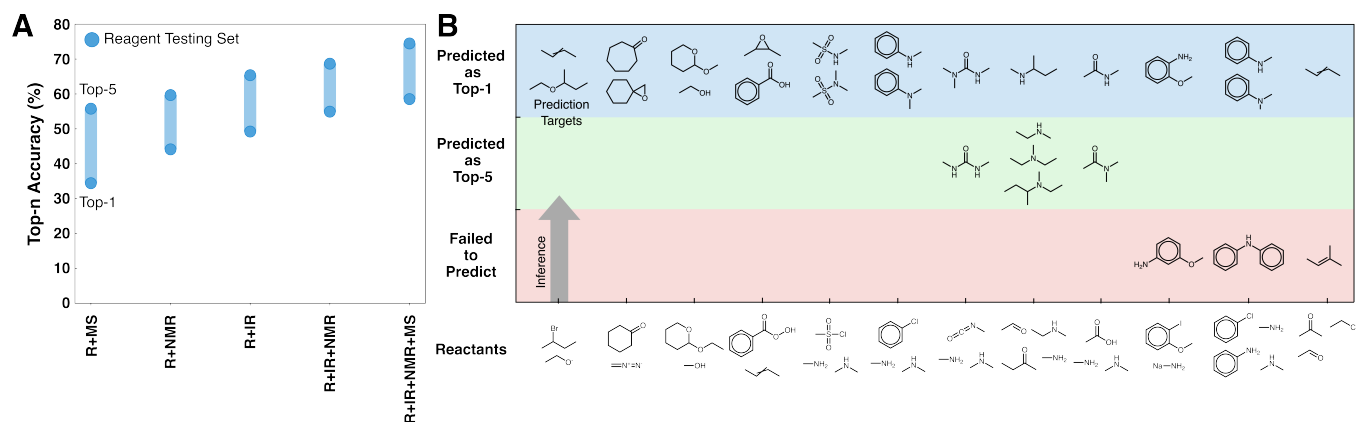


Figure 3: Performance of the deduction models on external testing sets. (A) Comparison of top-n performance in identifying reagents that were unseen as prediction targets during training. (B) Performance of the R+IR+NMR+MS model in predicting major and minor products of unseen reactions involving 18 reactants. The 15 products for the 7 reactants that are not shown were not predicted in the top-5 by the model.

180 18 organic reactants, each with two or more possible products producing a total of 40 distinct  
 181 reactions, curated from published and textbook sources (see Methods).<sup>21,22</sup> None of these reactions  
 182 existed in the training data, and predicting side-products (as opposed to starting material) was  
 183 not a task that was directly trained for. The R+IR+NMR+MS model can identify the major and  
 184 minor products in the top-1 for 19/40 of the reactions for 11/18 of the distinct reactants (Fig. 3B).  
 185 For reference, the IR+NMR+MS model, for which this is an on-target task, correctly identifies  
 186 19/30 of the major and minor products after excluding those seen during its training. Several of  
 187 the failure cases are also illuminating. For example, the structural isomers of anisidine are largely  
 188 indistinguishable using the limited analytical sources provided to the model. Nevertheless, the  
 189 transferability to this unseen task suggests that when provided with additional spectral sources  
 190 and task-specific training, this architecture is also capable of side-product identification.

191 The deductive super-networks studied here were designed to weight evidence from inductive  
 192 sub-models responsible for digesting individual information sources. This concept was loosely  
 193 inspired by human deduction, whereby training occurs on specific inductive tasks (e.g., certain

194 types of math, physics, or organic synthesis problems) that are consulted to construct and weight  
195 hypotheses and reject solutions in practical scenarios. This idea is also consistent with deductive  
196 behavior being an emergent capability of sufficiently expansive inductive subsystems or training  
197 datasets. For example, large language models show emergent deductive behavior as evidenced by  
198 their ability to respond to non sequiturs, questions that assume certain knowledge, and questions  
199 with false premises that contradict established knowledge.<sup>23</sup> Similarly, the surprising versatility  
200 of language models in generative chemical applications and general chemical problem solving has  
201 been documented by several groups.<sup>17,24,25</sup> The initial version of this architecture demonstrated  
202 surprising transferability to off-target tasks and in prediction scenarios with partial and even con-  
203 tradictory information. Additional variations on this architecture for product prediction and other  
204 deductive problems are immediately possible. Among the most obvious that were left unexplored  
205 are finding the optimal manner of combining the inductive sub-models (e.g., more sophisticated  
206 couplings beyond the linear reweighting used here) and training the super-network (e.g., training  
207 on multiple tasks or contrasting examples).

208 There are many opportunities for further improving these models and for applications beyond  
209 product identification. For example, the current work has not addressed the problem of product  
210 identification when the spectra contain product mixtures. Knowledge about the number of species  
211 is a powerful deductive constraint that was provided here implicitly through the training data  
212 curation; however, this too could be treated as a learnable deduction using an additional classifier  
213 or spectral segmentation model to deconvolute spectra for the spectral transformers. This is beyond  
214 the current scope, other than to acknowledge the opportunity. Deductive architectures should find  
215 application more generally in any prediction scenario where a non-functional one-to-many mapping  
216 occurs. These include predictions of materials aging, predictive maintenance, reaction planning,  
217 and inverse materials design, among others where missing variables, stochastic factors, or extra  
218 degrees of freedom make the prediction problem underdetermined. Such scenarios require deductive

219 reasoning, for which the state-of-the-art is often manual expert analysis of disparate information  
220 sources. Deductive ML models of the kind demonstrated here should find use in a multitude of  
221 similar applications.

## 222 **Acknowledgments**

223 The work was made possible by the Office of Naval Research (ONR) through support provided by  
224 the Energetic Materials Program (MURI grant number:N00014-21-1-2476, Program Manager: Dr.  
225 Chad Stoltz). B.M.S also acknowledges partial support for this work from the Dreyfus Program  
226 for Machine Learning in the Chemical Sciences and Engineering

## 227 **Data and Materials Availability**

228 All results are summarized in the main text or supplementary figures. Jupyter notebooks for gen-  
229 erating main text and supplementary figures, trained models, model training scripts, and training  
230 datasets have been uploaded to [XXX figshare link to be populated upon publication XXX].

## 231 **Methods**

### 232 **Dataset Curation**

#### 233 **Dataset Summary**

234 The final product identification dataset curated here consists of 446,330 samples, split between  
235 299,658 samples corresponding to real product prediction and 146,672 samples corresponding to  
236 starting material prediction. Each sample in the dataset is composed of the reactant and reagent  
237 SMILES, the simulated EI-MS, IR, and <sup>1</sup>H-NMR of the prediction target as available features,

238 and the product SMILES as the prediction target. Two versions of the dataset were used, one  
239 with reagents distinguished from other reactants using a special token, “>”, and one without. A  
240 80:10:10 training:validation:testing split was used for all model development. The curation details  
241 of this dataset and the data splits are summarized in the remaining sections.

## 242 **Dataset Curation**

243 The USPTO reaction dataset originally curated by Derek Lowe then filtered and split by Jin et al  
244 served as the starting point for data curation.<sup>15,16</sup> This dataset provided reactant:product pairs in  
245 the form of SMILES strings that needed to be augmented with spectral data (i.e., EI-MS, IR, and  
246 H-NMR) for each species for use in the product identification learning task. Filtering the reactions  
247 for compatibility with the spectral generation workflow (described next) resulted in 299,658 distinct  
248 reactions involving 374,681 distinct molecules (counting distinct reactants, reagents, and products).

## 249 **Simulated Spectra**

250 Spectra were simulated for all 374,681 distinct molecules in the dataset, because open-source  
251 spectral databases are insufficiently large and have limited overlap with the Lowe species to be  
252 useful for training a practical product identification model. IR spectra with 4 cm<sup>-1</sup> resolution from  
253 400-4000 cm<sup>-1</sup> were generated from the SMILES string of each molecule using the message-passing  
254 neural network model published by McGill et al.<sup>26</sup> EI-MS spectra with 1 m/z resolution from 1-  
255 999 m/z were generated using bidirectional neural network model (NEIMS) and rapid approximate  
256 subset-based spectra prediction (rassp) model published by Wei et al and Zhu et al respectively.<sup>27,28</sup>  
257 In general, the rassp spectra are more accurate but have size limitations, so NEIMS spectra were  
258 used as substitutions wherever rassp spectra were unavailable (about half of the spectra). 1H-  
259 NMR spectra with 0.0121 ppm resolution from -2ppm - 10ppm were generated using Mestrenova  
260 v14.3.0.<sup>29</sup> Spectral generation for both EI-MS and 1H-NMR required optimized geometries of each

261 species that were generated using Auto3D.<sup>30</sup> Reactions from the Jin et al USPTO dataset involving  
262 species with more than 30 heavy atoms or elements besides H, B, C, Si, N, P, O, S, Se, F, Cl, Br,  
263 and I were discarded to conform to the current constraints of Auto3D.<sup>16</sup> These exclusions resulted  
264 in the final set of 299,658 reactions with real products as prediction targets.

## 265 **Null Reactions**

266 To test the model's deductive capability, a set of "null reactions" was generated that share the same  
267 reactants and reagents as real reactions but with products and input spectra corresponding to one of  
268 the reactants. Predicting the product of such reactions corresponds to identifying starting material  
269 as an unintended product using the information provided by the spectra. The introduction of null  
270 reactions also creates an underdetermined scenario for a RtP model, since a given reactant can yield  
271 multiple potential products. Null reactions were generated for each of the 299,658 real reactions.  
272 All possible null reactions were generated for reactions with multiple reactants. Null reactions  
273 were discarded if their prediction target matched a product of a real reaction in the dataset. This  
274 was done to avoid accidental information leakage between null reactions and real reactions and  
275 also because it yielded a useful 2:1 data balance between real and null reactions without further  
276 filtering. A total of 146,672 null reactions satisfied this criteria, resulting in a combined dataset of  
277 446,330 reactions (i.e., 146,672 null and 299,658 real) for the product identification task.

## 278 **Dataset Splitting**

279 An 80:10:10 training:validation:testing split was used for model development. The splitting was  
280 performed so that all reactions that shared a prediction target were partitioned to the same split.  
281 This was done to ensure that the testing and validation sets correspond to unseen prediction targets.  
282 For example, if ibuprofen was a product of five different real reactions and two null reactions in the  
283 dataset, then all seven would be partitioned to the same split (at random) since they all share the

284 same prediction target (i.e., ibuprofen). This avoids information exchange between tasks, where  
285 the model would potentially see the same prediction spectra during training and testing. The total  
286 number of real and null reactions, together with their training-validation-test split is summarized  
287 in Table1.

Table 1: Dataset Split Used for Deduction Model Training

	Training set	Validation set	Test set
Real reactions	249766	24969	24923
Null reactions	108212	11000	13349

## 288 External Testing Datasets

289 Two additional datasets, MULTI and REAGENT, were curated to test the performance of the  
290 deduction models when predicting reactions with side products and identifying reagents as potential  
291 products, respectively. The MULTI dataset consists of a set of organic reactions with known side-  
292 products curated from Grossman’s textbook and the dataset compiled by Hartenfeller et al.<sup>21,22</sup>  
293 These reactions were combined to produce a total 18 reactants involved in reactions yielding 40  
294 distinct products. The REAGENT dataset was curated by identifying all unique reagent species  
295 from the main dataset and excluding any that overlapped with targets in the training set or that  
296 were incompatible with the spectral generation workflow. This resulted in 2707 distinct reagents.  
297 Up to three reactions, if available, from the main dataset involving each reagent was selected at  
298 random and the prediction target and input spectra were swapped for the reagent to yield a total  
299 3262 reactions. This dataset tests whether the models are able to identify reagents as a potential  
300 product. The spectra of all species in the MULTI and REAGENT datasets were simulated using  
301 the same protocol as the main training dataset.



## 302 **Neural Network Architecture**

### 303 **Architecture Summary**

304 All product identification models used an architecture composed of a reaction transformer, one or  
305 more spectral transformers, and a single linear deduction layer. The transformers were adapted  
306 from those now typical of neural machine translation (NMT) tasks,<sup>31</sup> using hyperparameter tun-  
307 ing based on the validation set accuracy. Both reactant and spectral data were pre-processed  
308 beforehand and then fed into the attention score calculation module of each transformer through  
309 the trainable embedding network. Inference was performed by these models in recursive token-  
310 by-token fashion until encountering an end token. An illustration of the R+IR+NMR+MS model  
311 architecture is shown in Figure S1. The largest model trained here, R+IR+NMR+MS, has  $\sim 30M$   
312 weights.

### 313 **Input Embedding**

314 The raw reactant input data were represented as SMILES strings, because this is currently the  
315 most reliable representation in reaction prediction tasks.<sup>32</sup> The SMILES strings were tokenized  
316 using a standard SMILES vocabulary of 284 possible tokens in addition to a special > symbol  
317 used (when present) to separate the reactants and reagents (e.g., solvents or catalysts), a padding  
318 token, and special start and end tokens (only present in the decoded product strings). Reactant  
319 inputs were converted to fixed 276-length ( $d_{seq}$ ) input vectors using padding tokens before being  
320 passed to a linear token embedding layer that converted each token to a 256-length vector ( $d_{emb}$ ).  
321 The dimensions of the reactant input after embedding were [276,256] (i.e.,  $d_{seq}$  by  $d_{emb}$ ). The batch  
322 dimension is omitted for clarity from all reported sizes.

323 The raw simulated 1H-NMR, EI-MS, and IR spectra were represented as intensity versus ppm,  
324 m/z, and  $\text{cm}^{-1}$  vectors, respectively. To prepare the 1H-NMR and EI-MS spectra for embedding,

325 the intensity values were normalized to a range between 0 and 1, binned by percentile (lower  
326 range exclusive, upper range inclusive), then tokenized based on the 100 possible percentile ranges  
327 and a special bin for zero (i.e., the percentiles served as a vocabulary for tokenization). The  
328 embedding of the IR spectra was identical except that intensities less than 1% were zeroed out to  
329 eliminate potential background noise, resulting in 100 total possible tokens rather than 101 (i.e.,  
330 the zero token for IR includes the first bin in the 1H-NMR and EI-MS cases, so there is one less  
331 token). The preprocessed input vectors for the IR, 1H-NMR, and EI-MS spectra were of length  
332 900 (representing 400-4000  $\text{cm}^{-1}$  with a 4  $\text{cm}^{-1}$  resolution), 993 (representing -2ppm - 10ppm with  
333  $\sim 0.0121$  ppm resolution), and 999 (representing 1-999  $m/z$  with 1  $m/z$  resolution). The input  
334 vectors were then embedded using a linear layer (specific to each transformer but with  $d_{emb} = 256$   
335 in all cases) in the same manner as the reactants, resulting in embedded inputs of size [900,256],  
336 [993,256], and [999,256] for the the IR, 1H-NMR, and EI-MS transformers, respectively.

337 To retain the spatial information of the inputs for use by the models (i.e., token position for  
338 the reactants and peak location for the spectra), standard trigonometric positional embedding (P)  
339 was added to the token-based embeddings according to

$$\begin{aligned} P(k, 2i) &= \sin\left(\frac{k}{n^{2i/d}}\right) \\ P(k, 2i + 1) &= \cos\left(\frac{k}{n^{2i/d}}\right) \end{aligned} \quad (1)$$

340 where  $k$  is the position of the input token,  $i$  is the position in the embedding dimension,  $d$  is the  
341 hidden dimension ( $d_{emb}$ ), and  $n$  is a convenient constant for determining the relative frequency  
342 shift between the sequentially sampled periodic functions (taken to be  $10^4$ , here).

### 343 Attention Cells

344 Each transformer is composed of a task-specific encoder and decoder that use two to four attention  
345 cells. Each encoder attention cell consists of a sequence of layer norm, multi-head self-attention

346 layer, residual connection, layer norm, feed-forward layer, and residual connection (Fig. S2). The  
347 layer norm is performed before other attention and feed-forward operations with an  $\epsilon$  value of  $10^{-6}$ .  
348 Eight attention heads were used, using linear projections of the input embedding dimension to form  
349 key and query vectors of length 256 ( $d_k = d_q = 256$ ) and value vectors of length  $d_v = d_{emb}/8 = 32$ ,  
350 and the dot-product attention mechanism calculated according to

$$\text{Score}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (2)$$

351 where Q, K, and V are matrices containing the queries, keys, and values for each embedded token  
352 (for the first cell, afterwards the derived feature of the previous cell) in the sequence with sizes  
353 of  $[d_{seq}, d_k]$ ,  $[d_{seq}, d_k]$ , and  $[d_{seq}, d_v]$ , respectively, and  $\sqrt{d_k}$  is a normalization factor. The outputs  
354 of each head are catenated along the value dimension to recover a matrix of the same size as the  
355 input to the attention layer. The catenated output from the multi-head attention layer is added  
356 to the input of the attention cell via a residual connection, then passed to a second layer norm  
357 and fed to a feed-forward block that consists of a linear layer to project the  $d_{emb}$ -dimension into  
358 a 2048-length vector, followed by a ReLU activation layer, and a second linear layer to project  
359 the hidden dimension from 2048 back to  $d_{emb}$ . Two drop-out layers with drop-out rate of 0.1 were  
360 applied after each linear transformation during training. Finally, the input to the attention cell is  
361 mixed with the output via another residual connection.

362 The decoder attention cells used in these models are identical to the encoder attention cells,  
363 with the exceptions that the target SMILES embedding is used as an input to the first cell,  
364 the multi-head self-attention layer uses masking to restrict non-zero attention calculations to later  
365 tokens, and a multi-head cross-attention layer is inserted after the masked multi-head self-attention  
366 layer (Fig. S2). The embedding layer used for the predicted product SMILES is shared across  
367 transformers and determined by training. The self-attention masking is identical to that used by

368 Vaswani et al.<sup>31</sup> The multi-head cross-attention layer is identical to the unmasked multi-head self-  
369 attention layer in the encoder attention cells, except that the key and value inputs are obtained as  
370 linear projections of the embedding dimension of the encoder output and the queries are obtained  
371 as linear projections of the embedding dimension of the output of the masked self-attention layer.  
372 Layer norms are used before each attention layer and residual connections are used after each  
373 attention layer (the same as for the encoder, there is just an extra one of each); all other details  
374 (sizes, sequence, number of heads, the final feed-forward layer, etc.) are identical to the encoder  
375 attention cells.

## 376 **Transformers**

377 All models were constructed from one or more transformers, with each consisting of an encoder,  
378 decoder, and terminal linear softmax classifier to predict the next token in the sequence. The  
379 encoder and decoder of each transformer were composed of a series of the attention cells described  
380 in the previous section. In the case of the reactant transformer, four attention cells were used in  
381 the encoder and decoder; whereas, for all spectral transformers only two attention cells were used  
382 in the encoder and decoder. A minimal loss in validation accuracy was observed upon reducing  
383 the number of attention cells in the spectral transformers and this expedited model training. More  
384 transformers might be useful when training on different data sources or other spectral inputs.

385 The RtP model consists of a single reactant transformer; the various StS models consist of one or  
386 more spectral transformers and no reactant transformer; and the various deduction models consist  
387 of a reactant transformer and one or more spectral transformers. For each case, the  $[d_{seq}, d_{emb}]$   
388 output of each transformer is linearly projected along the embedding-dimension to a 288-length  
389 vector (i.e., the number of SMILES plus special tokens) with a softmax to predict the probability  
390 of the next token.

## 391 **Deductive Layer**

392 The models that combine more than one transformer (i.e., the various StS and R+spectra models)  
393 are linked together by a single linear layer that projects the  $288 \times N$  token-probabilities outputted  
394 by the  $N$  individual transformers to predict the next token. Specifically, the outputs of the trans-  
395 formers are catenated to a  $288 \times N$ -length vector that is linearly projected to a 288-length vector  
396 with a softmax to predict the probability of the next token. Because the weights of this linear pro-  
397 jection layer are static after training and independent of the input, this layer represents a simple  
398 weighting of the evidence from the different transformers that potentially also accounts for any  
399 average linear correlations in the token-predictions observed during training.

400 The linear linkage of the transformers provides two mechanisms by which the task-specific  
401 transformers can act as deductive constraints on each other. The first is through the formation of  
402 a consensus prediction of the next token. This simple mechanism allows the more confident trans-  
403 formers to potentially overrule one or more less confident transformers in predicting a particular  
404 token. The second is through the recursive token-by-token manner in which the product predic-  
405 tion is made. At each step of this process, the prediction string, updated with the token from the  
406 last inference, is passed to all transformers to make their individual next-token predictions. This  
407 creates a mechanism by which the transformers can perform inference on prediction strings that  
408 they never would have encountered via a greedy decoding. For example, a particular transformer  
409 may be overruled by the others for several tokens, such that it is now performing inference on a  
410 partially decoded product scaffold that it would not have predicted on its own. In such a case, the  
411 other transformers have acted as a deductive constraint on the transformer.

412 Other deductive connections are likely useful but have not been significantly explored due to  
413 the immediate success of the current architecture for these prediction tasks. The only alternative  
414 that was significantly tested was an architecture that terminated in an additive layer rather than

415 a linear projection, which resulted in a marginal reduction in validation set accuracy.

## 416 **Training**

417 All models were trained using the Adam optimizer and a batch size of 20. The learning rate,  $\eta$ ,  
418 was linearly increased each update step followed by an exponential decay according to

$$\eta = \frac{1}{\sqrt{d_{\text{emb}}}} * \min\left(\frac{1}{\sqrt{s}}, \frac{s}{s_{\text{warm}}^{3/2}}\right) \quad (3)$$

419 where,  $s$ , is the step,  $s_{\text{warm}}$  is the number of steps within the warmup phase, and  $d_{\text{emb}}$  is the  
420 embedding dimension length.  $s_{\text{warm}}$  was set to 37500 steps, roughly 4% of the overall training  
421 steps, which is consistent with Vaswani et al.<sup>31</sup> No label smoothing was used during training.  
422 Early stopping was applied to terminate training if the validation loss did not decrease in the  
423 consecutive 30 epochs.

424 One R+IR+NMR+MS model was trained with random dropping of the spectral sources for  
425 use in Figure 2D of the main text. All other results are for models trained without dropping. For  
426 the model trained with dropping, a 10% probability of dropping was separately applied to each  
427 input spectrum during training (i.e., on average 1/1000 training samples had no input spectra).

## 428 **Inference**

429 During the inference cycle, all models' top-k outputs are determined by a beam search with beam  
430 size set to five. The beam search algorithm is consistent with the previous implementation pub-  
431 lished by Schwaller et al.<sup>17</sup> The inference cycle is initiated by feeding the target input with a  
432 dummy string only containing the start token "<". This replaces the target product's SMILES  
433 that is used in the training cycle. The model then selects the five most probable tokens decoded  
434 from the start string to form five new beams. At each decoding step, each of the beams produces

435 another five candidate strings, and the five candidates with the highest overall probability are se-  
436 lected from the pool of 25 strings, which are then assigned to the new beams for the next decoding  
437 step. The decoding of each beam terminates if the end token “\$” is predicted as the top-1 or the  
438 string length reaches the upper limit of 67.

## 439 **Transformer Decisiveness and Input Accuracy Reduction**

440 The decisiveness measure was implemented by zeroing out the final probability prediction of each  
441 transformer before it was passed to the linear deduction layer. If this caused a change in the top-1  
442 predicted token compared with the unmodified inference, then the transformer was classified as  
443 being decisive for that token. According to this definition, one or more transformers can be decisive  
444 for a token, and also no transformer can be decisive if a sufficiently strong consensus exists. If a  
445 transformer was decisive for at least one token in a given product decoding, then it was classified  
446 as being decisive for that product.

447 The overall accuracy reduction is an alternative measure of input importance that simply  
448 reports the reduction in overall top-n accuracy when each of the input sources are individually  
449 zeroed out. This was implemented by supplying a single padding token to the reactant transformer,  
450 and three zero intensity tokens as inputs to the spectral transformers, respectively. The overall  
451 accuracy reduction is not necessarily equivalent to the decisiveness of each transformer, because of  
452 the flexibility of the SMILES language, which allows the same molecule to be decoded in multiple  
453 ways, and the important role of consensus formation in the decoding.

## 454 **References**

- 455 (1) Bubliauskas, A.; Blair, D. J.; Powell-Davies, H.; Kitson, P. J.; Burke, M.; Cronin, L. A  
456 practical approach to combine modular reactions and reactionware for the digitization of

- 457 chemical synthesis. *Angew. Chem. Int. Ed.* **2022**, *61*.
- 458 (2) Lin, Y.; Zhang, R.; Wang, D.; Cernak, T. Computer-aided key step generation in alkaloid  
459 total synthesis. *Sci.* **2023**, *379*, 453–457.
- 460 (3) Manzano, J. S.; Hou, W.; Zalesskiy, S. S.; Frei, P.; Wang, H.; Kitson, P. J.; Cronin, L.  
461 An autonomous portable platform for universal chemical synthesis. *Nat. Chem.* **2022**, *14*,  
462 1311–1318.
- 463 (4) Zahrt, A. F.; Mo, Y.; Nandiwale, K. Y.; Shprints, R.; Heid, E.; Jensen, K. F. Machine-  
464 Learning-Guided Discovery of Electrochemical Reactions. *J. Am. Chem. Soc.* **2022**, *144*,  
465 22599–22610.
- 466 (5) Lumley, J. A.; Sharman, G.; Wilkin, T.; Hirst, M.; Cobas, C.; Goebel, M. A KNIME Workflow  
467 for Automated Structure Verification. *SLAS Discov.* **2020**, *25*, 950–956.
- 468 (6) Fine, J. A.; Rajasekar, A. A.; Jethava, K. P.; Chopra, G. Spectral deep learning for prediction  
469 and prospective validation of functional groups. *Chem. Sci.* **2020**, *11*, 4618–4630.
- 470 (7) Huang, Z.; Chen, M. S.; Worocho, C. P.; Markland, T. E.; Kanan, M. W. A framework for  
471 automated structure elucidation from routine NMR spectra. *Chem. Sci.* **2021**, *12*, 15329–  
472 15338.
- 473 (8) Yao, L.; Yang, M.; Song, J.; Yang, Z.; Sun, H.; Shi, H.; Liu, X.; Ji, X.; Deng, Y.; Wang, X.  
474 Conditional Molecular Generation Net Enables Automated Structure Elucidation Based on  
475 <sup>13</sup>C NMR Spectra and Prior Knowledge. *Anal. Chem.* **2023**, *95*, 5393–5401.
- 476 (9) Jung, G.; Jung, S. G.; Cole, J. M. Automatic materials characterization from infrared spectra  
477 using convolutional neural networks. *Chem. Sci.* **2023**, *14*, 3600–3609.



- 478 (10) Lemm, D.; von Rudorff, G. F.; von Lilienfeld, O. A. Machine learning based energy-free  
479 structure predictions of molecules, transition states, and solids. *Nat. Commun.* **2021**, *12*,  
480 4468.
- 481 (11) Heinen, S.; von Rudorff, G. F.; von Lilienfeld, O. A. Toward the design of chemical reactions:  
482 Machine learning barriers of competing mechanisms in reactant space. *J. Chem. Phys.* **2021**,  
483 *155*, 064105.
- 484 (12) Krenn, M.; Pollice, R.; Guo, S. Y.; Aldeghi, M.; Cervera-Liarta, A.; Friederich, P.; dos  
485 Passos Gomes, G.; Häse, F.; Jinich, A.; Nigam, A., et al. On scientific understanding with  
486 artificial intelligence. *Nat. Rev. Phys.* **2022**, *4*, 761–769.
- 487 (13) Anstine, D. M.; Isayev, O. Generative Models as an Emerging Paradigm in the Chemical  
488 Sciences. *J. Am. Chem. Soc.* **2023**, *145*, 8736–8750.
- 489 (14) Weininger, D. SMILES, a chemical language and information system. 1. Introduction to  
490 methodology and encoding rules. *J. Chem. Inf. Comput. Sci.* **1988**, *28*, 31–36.
- 491 (15) Lowe, D. M. Extraction of chemical structures and reactions from the literature. Ph.D. thesis,  
492 University of Cambridge, 2012.
- 493 (16) Jin, W.; Coley, C.; Barzilay, R.; Jaakkola, T. Predicting organic reaction outcomes with  
494 weisfeiler-lehman network. *Adv. Neural Inf. Process. Syst.* **2017**, *30*.
- 495 (17) Schwaller, P.; Laino, T.; Gaudin, T.; Bolgar, P.; Hunter, C. A.; Bekas, C.; Lee, A. A. Molecular  
496 transformer: a model for uncertainty-calibrated chemical reaction prediction. *ACS Cent. Sci.*  
497 **2019**, *5*, 1572–1583.
- 498 (18) Tu, Z.; Coley, C. W. Permutation invariant graph-to-sequence model for template-free ret-  
499 rosynthesis and reaction prediction. *J. Chem. Inf. Model.* **2022**, *62*, 3503–3513.

- 500 (19) Ji, H.; Deng, H.; Lu, H.; Zhang, Z. Predicting a molecular fingerprint from an electron  
501 ionization mass spectrum with deep neural networks. *Anal. Chem.* **2020**, *92*, 8649–8653.
- 502 (20) Alberts, M.; Laino, T.; Vaucher, A. Leveraging Infrared Spectroscopy for Automated Struc-  
503 ture Elucidation. **2023**,
- 504 (21) Grossman, R. B.; Grossman, R. *The art of writing reasonable organic reaction mechanisms*;  
505 Springer, 2003.
- 506 (22) Hartenfeller, M.; Eberle, M.; Meier, P.; Nieto-Oberhuber, C.; Altmann, K.-H.; Schneider, G.;  
507 Jacoby, E.; Renner, S. A collection of robust organic synthesis reactions for in silico molecule  
508 design. *J. Chem. Inf. Model.* **2011**, *51*, 3093–3098.
- 509 (23) OpenAI, non-sequitur; assumed outside knowledge; false premise. 2023; [https://chat.  
510 openai.com/share/e678c670-2ec8-44fb-bcd0-056d993c4192](https://chat.openai.com/share/e678c670-2ec8-44fb-bcd0-056d993c4192).
- 511 (24) Flam-Shepherd, D.; Zhu, K.; Aspuru-Guzik, A. Language models can learn complex molecular  
512 distributions. *Nat. Commun.* **2022**, *13*, 3293.
- 513 (25) White, A. D.; Hocky, G. M.; Gandhi, H. A.; Ansari, M.; Cox, S.; Wellawatte, G. P.; Sasmal, S.;  
514 Yang, Z.; Liu, K.; Singh, Y.; Peña Ccoa, W. J. Assessment of chemistry knowledge in large  
515 language models that generate code. *Dig. Discov.* **2023**, *2*, 368–376.
- 516 (26) McGill, C.; Forsuelo, M.; Guan, Y.; Green, W. H. Predicting infrared spectra with message  
517 passing neural networks. *J. Chem. Inf. Model.* **2021**, *61*, 2594–2609.
- 518 (27) Wei, J. N.; Belanger, D.; Adams, R. P.; Sculley, D. Rapid prediction of electron–ionization  
519 mass spectrometry using neural networks. *ACS Cent. Sci.* **2019**, *5*, 700–708.
- 520 (28) Zhu, R. L.; Jonas, E. Rapid approximate subset-based spectra prediction for electron  
521 ionization–mass spectrometry. *Anal. Chem.* **2023**, *95*, 2653–2663.

- 522 (29) Willcott, M. R. MestRe Nova. *J. Am. Chem. Soc.* **2009**, *131*, 13180–13180.
- 523 (30) Liu, Z.; Zubatiuk, T.; Roitberg, A.; Isayev, O. Auto3D: Automatic Generation of the Low-  
524 Energy 3D Structures with ANI Neural Network Potentials. *J. Chem. Inf. Model.* **2022**, *62*,  
525 5373–5382.
- 526 (31) Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A. N.; Kaiser, Ł.;  
527 Polosukhin, I. Attention is all you need. *Adv. Neural Inf. Process. Syst.* **2017**, *30*.
- 528 (32) Jaume-Santero, F.; Bornet, A.; Valery, A.; Naderi, N.; Vicente Alvarez, D.; Proios, D.; Yaz-  
529 dani, A.; Bournez, C.; Fessard, T.; Teodoro, D. Transformer Performance for Chemical Reac-  
530 tions: Analysis of Different Predictive and Evaluation Scenarios. *J. Chem. Inf. Model.* **2023**,  
531 *63*, 1914–1924.