
ORDERly: Datasets and benchmarks for chemical reaction data

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Machine learning has the potential to provide tremendous value to the life sciences
2 by providing models that aid in the discovery of new molecules and reduce the
3 time for new products to come to market. Chemical reactions play a significant role
4 in these fields, but there is a lack of high-quality open-source chemical reaction
5 datasets for training ML models. Herein, we present ORDERly, an open-source
6 Python package for customizable and reproducible preparation of reaction data
7 stored in accordance with the increasingly popular Open Reaction Database (ORD)
8 schema. We use ORDERly to clean US patent data stored in ORD and generate
9 datasets for forward prediction, retrosynthesis, as well as the first benchmark for
10 reaction condition prediction. We train neural networks on datasets generated with
11 ORDERly for condition prediction and show that datasets missing key cleaning
12 steps can lead to silently overinflated performance metrics. Additionally, we
13 train transformers for forward and retrosynthesis prediction and demonstrate how
14 non-patent data can be used to evaluate model generalisation. By providing a
15 customizable open-source solution for cleaning and preparing large chemical
16 reaction data, ORDERly is poised to push forward the boundaries of machine
17 learning applications in chemistry.

18 1 Introduction

19 Advancements in chemistry and material science hinge on the availability of high-quality chemical
20 reaction data, and the advent of machine learning (ML) for science has highlighted the value that data
21 can bring to chemistry. One important application is in the pharmaceutical industry, where figuring
22 out *how* to make novel molecules remains a significant bottleneck, causing delays in the "make" step
23 of the "design, make, test" cycle [1]. Making a molecule (product) includes predicting the reaction
24 pathway (retrosynthesis) and suitable reaction conditions (e.g. solvents and reagents), and optimising
25 for one or more outcomes such as reaction yield, selectivity, and conversion. ML is well suited to
26 assist with these tasks, with a range of tools being developed for forward reaction prediction [2, 3, 4],
27 retrosynthesis [5, 6, 7, 8, 9], condition prediction [10, 11, 12], yield prediction [13, 14, 15], and
28 closed-loop optimisation [16, 17, 18].

29 Building reaction prediction tools requires access to large datasets for training. Historically, re-
30 searchers have accessed proprietary in-house datasets or acquired the data through commercial
31 databases such as Reaxys [19]. The advantage of commercial databases is both the scale of the
32 datasets available (often millions of reactions) and the annotation already completed by the publish-
33 ers. Yet, these datasets are not freely available to ML practitioners, stymieing advances in reaction
34 condition prediction in both academia and industry.

35 Recently, efforts have been made to create openly-accessible databases for chemical reaction data. In
 36 particular, the Open Reaction Database (ORD) [20] is promising due to its exhaustive schema for
 37 describing chemical reaction data and breadth of data already incorporated. Yet, many of the datasets
 38 in ORD require further processing before they can be used in ML pipelines, preventing practical use.
 39 This is especially true for the largest dataset in ORD extracted from the US patent literature (the
 40 "USPTO dataset" [21]). In this work, we endeavor to close this gap.

41 Herein, we present ORDERly, a new framework for extracting and cleaning data from ORD, accompa-
 42 nied by datasets for three reaction related tasks: retrosynthesis, forward, and condition prediction.
 43 By offering an open-source and customizable solution for cleaning chemical reaction data, ORDERly
 44 aims to contribute to the development of advanced ML models in chemistry and material science.

45 2 Problem formulation

46 As noted by Meng *et al.* [22], reaction related tasks operate on molecules. There are numerous
 47 machine readable molecular representations [23], including molecular graphs and strings, and in this
 48 work molecules are represented as SMILES strings. Each character m_i in a SMILES string represents
 49 an atom or a molecular feature (bond, branch, ring closure): $\mathcal{M} := m_1, m_2, m_3, \dots, m_L$, where L is
 50 the total number of characters in the string. Molecules can take on one of three roles in a reaction:
 51 reactant, product, or agent. A reaction \mathcal{R} transforms N reactant molecules (sometimes called educts)
 52 $\{\mathcal{M}_i^{\mathcal{E}}\}_{i=1}^N$ by breaking and forming bonds to form M new product molecules $\{\mathcal{M}_i^{\mathcal{P}}\}_{i=1}^M$ using K
 53 agent molecules $\{\mathcal{M}_i^{\mathcal{A}}\}_{i=1}^K$. Agents are helper molecules that enable the reaction to proceed (e.g.,
 54 solvents, catalysts).

$$\mathcal{R} : \{\mathcal{M}_i^{\mathcal{E}}\}_{i=1}^N, \{\mathcal{M}_i^{\mathcal{A}}\}_{i=1}^K \rightarrow \{\mathcal{M}_i^{\mathcal{P}}\}_{i=1}^M, \{\mathcal{M}_i^{\mathcal{A}}\}_{i=1}^K \quad (1)$$

55 Given this view of reactions, we define three different reaction related tasks in this work.

56 **Forward prediction** is the task of predicting the product of a reaction $\mathcal{M}^{\mathcal{P}}$ given its reactants
 57 $\{\mathcal{M}_i^{\mathcal{E}}\}_{i=1}^N$ and, potentially, agents $\{\mathcal{M}_i^{\mathcal{A}}\}_{i=1}^K$. Probabilistically, the task is to predict the distribution
 58 $p(\mathcal{M}^{\mathcal{P}} | \{\mathcal{M}_i^{\mathcal{E}}\}_{i=1}^N)$. While experimental evaluation in a wet lab requires expert chemists and is a time
 59 intense task, reaction outcome prediction can help as a tool to evaluate the quality of a predicted
 60 retrosynthetic route (i.e., the probability that the reaction predicted by the single-step retrosynthesis
 61 model leads to the desired product) [24].

62 **Retrosynthesis** is the task of designing a sequence of Z reactions $\mathcal{R}_1, \mathcal{R}_2, \mathcal{R}_3, \dots, \mathcal{R}_Z$ that trans-
 63 form a set of readily available reactant molecules $\{\mathcal{M}_i^{\mathcal{E}_1}\}_{i=1}^N$ to a desired product(s) $\{\mathcal{M}_i^{\mathcal{P}_Z}\}_{i=1}^{M_Z}$.
 64 Retrosynthesis is done in the reverse direction by starting with the desired product(s) $\{\mathcal{M}_i^{\mathcal{P}_Z}\}_{i=1}^{M_Z}$
 65 and predicting reactants $\{\mathcal{M}_i^{\mathcal{E}_Z}\}_{i=1}^{N_Z}$ that would react to form the desired product(s). The predicted

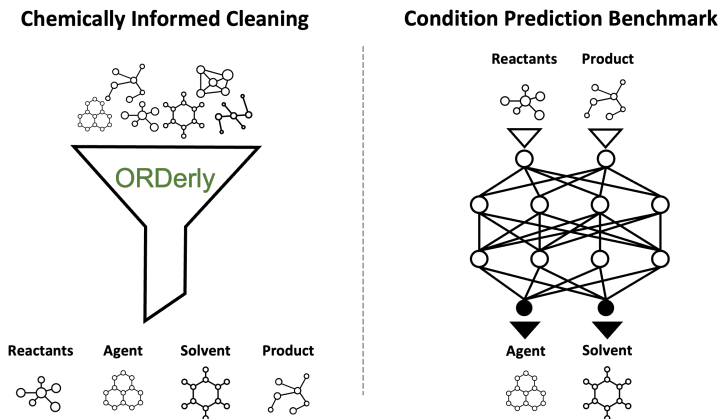


Figure 1: Overview of ORDERly.

66 reactants $\{\mathcal{M}_i^{\mathcal{E}_z}\}_{i=1}^{N_z}$ then become the products of the next reaction to be predicted $\{\mathcal{M}_i^{\mathcal{P}_{z-1}}\}_{i=1}^{M_{z-1}}$.
67 This process is repeated until a readily available set of starting reactant molecules are predicted
68 $\{\mathcal{M}_i^{\mathcal{E}_1}\}_{i=1}^N$. Therefore, the key machine learning task, often called single-step retrosynthesis, is
69 predicting the distribution $p(\{\mathcal{M}_i^{\mathcal{E}_j}\}_{i=1}^{N_j}|\mathcal{M}^{\mathcal{P}_j})$ or the set of reactants that could lead to a given
70 product(s) $\{\mathcal{M}_i^{\mathcal{P}_j}\}_{i=1}^{M_j}$. Single-step retrosynthesis can be seen as the inverse of forward prediction.
71 **Condition prediction** is the task of predicting the distribution $p(\{\mathcal{M}_i^{\mathcal{A}}\}_{i=1}^K|\{\mathcal{M}_i^{\mathcal{E}}\}_{i=1}^N, \mathcal{M}^{\mathcal{P}})$ (i.e.,
72 the agents for a reaction given reactants and product). In addition to agents, some models can predict
73 continuous variables such as reaction temperature and concentrations of reactants and agents [10].

74 3 Related work

75 3.1 Chemical reaction cleaning tools

76 Existing tools for cleaning reaction data are primarily targeted at retrosynthesis and forward prediction
77 tasks [25, 26, 27, 28] and have somewhat limited extensibility, given that they are built to take as
78 inputs CSV files or the stationary XML files of the US patent (USPTO) dataset [21] instead of
79 the outputs of continuously updated databases such as ORD [20]. Furthermore, in the original
80 publications, there is little to no discussion of how decisions made during cleaning (e.g. restricting
81 the number of components in a reaction or the minimum frequency of occurrence) impact the datasets
82 being cleaned or performance of models trained on the datasets. We believe that this is in part due to
83 data cleaning historically being viewed as a "low value" task, and therefore not adequately discussed
84 and published on.

85 USPTO, being the largest open-source chemical reaction dataset, has been cleaned a number of times
86 for different learning tasks. For example, the USPTO-50K [29, 30] and USPTO-MIT datasets [31]
87 are commonly used for benchmarking single-step retrosynthesis and forward predictions models¹,
88 and these benchmarks are available in aggregate benchmarking sets such as the Therapeutics Data
89 Commons (TDC) [32]. However, the code used to process the raw data to generate the aforementioned
90 USPTO benchmarks was not published and, there is no publicly available benchmark for reaction
91 condition prediction extracted from these datasets.

92 4 Dataset generation

93 ORDERly extracts data directly from ORD [20]. Even though the data in ORD is stored in accordance
94 with a structured schema, we found that further effort is required to transform the labeled data into
95 ML-ready datasets. Therefore, ORDERly is centered around a data extraction script and a data cleaning
96 script, both of which take numerous arguments that customize the operations being performed.

97 4.1 Extraction and cleaning methodology

98 The extraction script allows the user to choose whether reaction roles should be assigned using the
99 labeling in ORD or using chemically-informed logic on the atom-mapped reaction string (if available).
100 It also enables specification of data source (e.g., USPTO or non-USPTO), allowing users to train
101 models with data from one source and test the performance with data from another source. Creating
102 test sets from different data sources is a robust way to evaluate generalization performance.

103 We chose cleaning operations motivated by first-principles understanding of chemistry. Cleaning
104 operations on the chemical reaction data include: (1) Restricting the number of reactants and product,
105 preventing multi-step reactions being included in the dataset; (2) Ensuring that all molecules can
106 be sanitized by the cheminformatics package RDKit [33]; (3) Restricting the maximum number of
107 unique catalysts, solvents, and reagents in a reaction based on commonly used experimental amounts;
108 (4) Frequency filtering to remove outliers; (5) Sanity checking the yield ($0\% \leq yield \leq 100\%$),
109 temperature, and pressure; (6) Removing duplicates, and finally; (7) Applying a random split to create

¹We discuss the difference between these datasets and our dataset in Appendix A.3.2

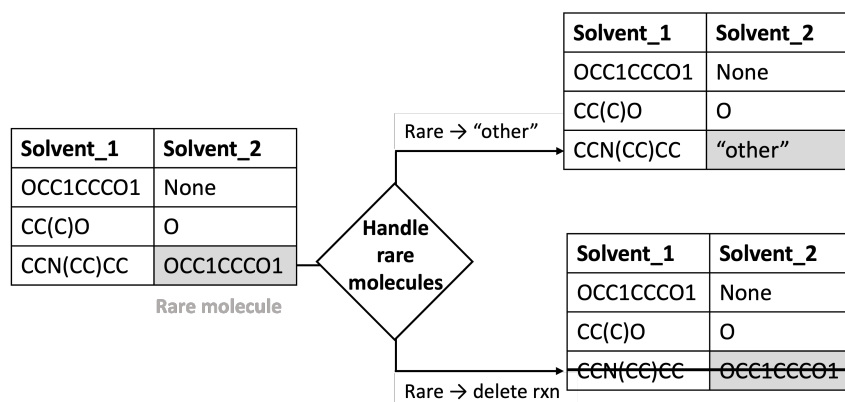


Figure 2: We present two different approaches for handling rare molecules. Rare \rightarrow "other" is investigated as a strategy to avoid deleting reactions with rare molecules.

110 training/validation/test sets, carefully ensuring that any inputs present in the train set (i.e. reactants
 111 and products for reaction condition prediction) are not also present in the test set.

112 **Computational details:** All extraction/cleaning operations described in this section were performed
 113 using a 2022 Mac Studio with an Apple M1 Max chip and 32GB memory. In ORD there are roughly
 114 1.7 million reactions from US patents (USPTO) and 91k reactions that are not from US patents. For
 115 the USPTO dataset extraction and sanitation took roughly 35 minutes, while the cleaning steps took 8
 116 minutes.

117 4.2 Reaction role assignment

118 We experimented with two approaches to assigning roles to the molecules found in a reaction (e.g.,
 119 whether a molecule is a reactant or an agent): trusting the labeling of molecules in ORD (referred to
 120 as "labeling") or applying chemical reaction logic to identify the role of different molecules from
 121 the reaction string (referred to as "rxn string" or "reaction string"). Our reaction logic identified
 122 reactants (molecules that contribute atoms to the product(s)) and spectator molecules (molecules
 123 that do not contribute atoms to the product(s)) based on the atom-mapping and their position in
 124 the reaction SMILES string. Solvents were identified in the list of spectator molecules by cross
 125 checking against a list of solvents compiled from prior research (see Appendix A.1.1), while all other
 126 spectator molecules are marked as agents. Catalysts were not separated into their own category since
 127 identifying catalysts can be quite subtle (especially with organocatalysis), and few reactions in the
 128 reaction string datasets contained transition metals.

129 4.3 Frequency filtering

130 Removing rare molecules can increase the signal to noise ratio in a dataset by removing outliers.
 131 In this work, we investigated two different strategies for filtering spectator molecules based on
 132 their frequency: deleting reactions with rare spectator molecules (rare \rightarrow delete rxn) or keeping the
 133 reactions but mapping the rare molecules to an "other" category (rare \rightarrow "other") (see Figure 2). We
 134 conducted experiments with both the rare \rightarrow delete rxn and rare \rightarrow "other" strategies for the task of
 135 condition prediction. The frequency threshold was set at 100 in line with previous research [10],
 136 though the sensitivity of dataset size to frequency threshold was still investigated (see Appendix C.2).
 137 Deleting reactions with rare molecules may create a more cohesive dataset by removing outliers,
 138 while renaming rare molecules "other" allows more reactions to be kept, offering more training data
 139 for the model.

Table 1: Number of reactions left in each dataset after cleaning. A description of each dataset can be found in section 4. Note that the actual number of reactions used for training will differ from the dataset size shown below due to train/test splits and augmentation. Non-USPTO-retro had a final dataset size of 20,830 and was cleaned in the same way as ORDerly-retro.

Dataset name:	ORDerly- condition (labeling)	ORDerly- condition (rxn string)	ORDerly- forward	ORDerly- retro	Non- USPTO- forward
Full dataset	1,771,032	1,771,032	1,771,032	1,771,032	91,067
Too many reactants	1,470,060	1,631,394	1,743,585	1,631,394	43,845
Too many products	1,329,399	1,593,196	1,740,655	1,593,196	40,770
Too many solvents	1,222,381	1,388,312	1,689,445	NA	36,522
Too many agents	1,202,790	1,279,833	1,550,800	NA	31,187
No reactants/products	1,202,758	1,262,333	1,533,680	1,567,697	31,095
No solvents	870,888	950,189	NA	NA	NA
No agents	135,139	690,234	NA	NA	NA
Inconsistent yields	126,948	658,071	NA	NA	NA
Dropping duplicates	76,634	392,996	919,231	941,566	28,496
Frequency filtering	75,033	356,906	NA	NA	NA

140 4.4 Dataset composition

141 Datasets generated with ORDerly have the following column groups: Reaction SMILES (string),
 142 is_mapped (bool), Reactants & products (SMILES strings), Solvents and agents (rxn string data), or
 143 solvents, catalysts, and reagents (labeling data) (SMILES strings), Temperature, reaction time, yield
 144 (floats), Procedure details (string), Grant date (datetime), date of experiment (datetime), file name
 145 (string).

146 Three new benchmarks were created from the USPTO dataset: ORDerly-forward for forward
 147 prediction, ORDerly-retro for retrosynthesis prediction, and ORDerly-condition for reaction
 148 condition prediction. Several additional datasets were created, including datasets from non-USPTO
 149 data in ORD and datasets to investigate data labeling and frequency filtering. An overview of the
 150 datasets and benchmarks showing how each cleaning step impacted the dataset size can be found
 151 in Table 1. The datasets are freely available and can be downloaded immediately from FigShare or
 152 regenerated using the code in the ORDerly Github repository.

153 5 Results and discussion

154 Experimental evaluation of the ORDerly-forward and ORDerly-retro benchmarks was performed
 155 using the Molecular Transformer architecture built by Schwaller *et al.* [2]. To switch from forward
 156 prediction to retrosynthesis prediction no changes to the transformer architecture were necessary,
 157 only the data was changed. The ORDerly-condition benchmark was evaluated together with the
 158 impact of different approaches to reaction role assignment and frequency filtering using the neural
 159 network architecture built by Gao *et al.* [10] with only minor modifications.

160 5.1 Forward and retrosynthesis prediction with transformers

161 Transformers were applied to two tasks: forward prediction (predicting products given reactants,
 162 solvents, and agents) and retrosynthesis (predicting reactants given a product). For the task of forward
 163 reaction prediction two different modes were tested: mixing the reactants, solvents, and agents, or
 164 weakly separating the reactants from the solvents and agents with a ">" token. Forward prediction
 165 with mixed inputs is a more difficult task, since it is less obvious which atoms (characters) will appear
 166 in the product.

167 For both forward and retrosynthesis prediction the order of the molecules was randomized, and the
 168 dataset was augmented by replacing each SMILES string in the reaction with a random equivalent

Table 2: Test performance with Molecular Transformer on forward prediction and retrosynthesis (%). The first column shows the percentage of invalid SMILES strings produced by the transformer (lower is better), while the second and third column show the top-1 accuracy with and without consideration of stereochemistry (SC), respectively (higher is better).

Test sets:	Random split from USPTO			Non USPTO		
	Invalid SMILES	Accuracy (with SC)	Accuracy (w/o SC)	Invalid SMILES	Accuracy (with SC)	Accuracy (w/o SC)
Forward (separated)	0.46	82.18	84.31	0.31	82.61	83.62
Forward (mixed)	0.47	80.79	82.86	0.31	82.61	83.62
Retrosynthesis	0.25	49.96	50.99	0.09	42.28	42.47

SMILES string (thus doubling the dataset size), before finally being tokenized [2]. Performance metrics are reported in Table 2, showing that across all tasks only a small percentage of the generated SMILES strings are invalid.

On the forward prediction tasks, the accuracies achieved are similar (albeit slightly lower) to the accuracies reported by Schwaller et al. [2] (88-90% top-1 accuracy when trained on the USPTO_MIT [31] dataset), though the accuracies are not directly comparable since different subsets of USPTO were used. As expected, the performance with separated agents is higher than mixed, since it is an easier task, and it is encouraging to see that the models get stereochemical information correct most of the time. Accuracy with the retrosynthesis model on the held out test set was roughly 50%, which is similar previous work on retrosynthesis [34]. It is interesting that prediction accuracy on the non-USPTO data was similar on the forward prediction tasks, but markedly worse on the retrosynthesis task.

Computational details: The transformer models were trained for around 35 hours (roughly 600 epochs) on a T4 cloud GPU instance provided by lightning.ai. Evaluation was done with the final model checkpoint.

5.2 Reaction condition prediction with neural networks

The reaction condition prediction model used in this work predicts five categorical variables: two solvents and three agents. These five molecules form a set (order invariant), though the loss function in the model used to predict the molecules considers them sequentially (with order) since this was found to work better in practice [10]. The metric used to evaluate the accuracy of the model should be order invariant, since the problem is order invariant, and for this reason the accuracy metrics used are top-1 (see appendix B) and top-3 (see Table 3) exact match combination accuracy for each type of component (i.e., solvent, agent). Beam search was used to identify the top-3 highest probability sets of reaction conditions. The top-3 accuracy was compared to the baseline predictive accuracy of simply predicting on the test set the most common molecules found in the train set.

Additionally, we define a metric inspired by Maser *et al.* [12] called the average improvement over baseline (AIB%):

$$AIB\% = \frac{A_m - A_b}{1 - A_b} * 100 \quad (2)$$

where A_m is the exact match combination accuracy of the model and A_b is the exact match combination accuracy of choosing the top 3 most common values of a component in the respective train set.

Table 3 shows the predictive performance on the test set using four different flavours of the ORDerly-condition benchmark. All models show an improvement over the frequency informed baseline.

Table 3: Top-3 metrics on condition prediction with the model architecture of Gao et al. [10]: frequency informed guess accuracy // model prediction accuracy // AIB%.

Datasets:	labeling	labeling	reaction string	reaction string
	rare→"other"	rare→delete rxn	rare→"other"	rare→delete rxn
Solvents	47 // 58 // 21%	50 // 61 // 22%	23 // 42 // 26%	24 // 45 // 28%
Agents	54 // 70 // 35%	58 // 72 // 32%	19 // 39 // 25%	21 // 42 // 27%
Solvents & Agents	31 // 44 // 19%	33 // 47 // 21%	4 // 21 // 18%	5 // 24 // 21%

The performance of the labeling datasets at first appears to be better than those that use our custom logic to extract reaction components from the reaction string. However, as shown in Figure 5, many of the reactions in datasets where we trust the labeling in ORD have more than three reactants, while most reactions in organic chemistry only have two reactants. Upon manual inspection, we found that many agents were mislabeled as reactants and, therefore, the prediction problem was made significantly easier. This insight is confirmed in Table 4; there are fewer unique solvents and agents and a higher density of null components when using the ORD labeling instead of the reaction string. This discrepancy demonstrates that naive creation of datasets based on ORD can lead to inflated performance metrics. In dealing with rare spectator molecules to avoid sparse OHE (see Table 1) we found that rare → delete rxn strategy performed better in practice. Therefore the ORDerly-condition benchmark uses the reaction string to assign reaction roles with the rare → delete rxn strategy.

For the datasets that extract the components from the reaction string, overall top-3 accuracy is less than 25% across solvents and agents. While not directly comparable, our overall accuracy is lower than what Gao *et al.* [10] achieved with 50.1% top-3 accuracy across catalysts, solvents and agents. However, Gao *et al.* trained on approximately ten million reactions, while we train on less than four percent of that (~350k). As shown in Figure 3, we see consistent increases in AIB (%) with the number of data points for the dataset which uses reaction strings and deletes rare reactions, and this scaling performance indicates that as ORD grows, better performance could be achieved, even with potentially fewer data points than used in the paper by Gao *et al.*

Computational details: These models were trained on an A10G cloud GPU instance provided by lightning.ai for 100 epochs to minimize cross entropy loss for each reaction component. The best model by validation loss was chosen for evaluation.

6 Technical limitations

6.1 Component labeling

Identifying the role of molecules in a reaction provides crucial context to machine learning models, and this identification could be improved with better atom-mapping [35]. However, an atom-mapping algorithm was not integrated into ORDERly to keep ORDERly lightweight. Even with perfect atom-mapping reaction role identification [29] can be challenging since the role of a molecule depends

Table 4: Diversity in the datasets. Frequency filtering was applied for the solvents and agents to create a more dense one-hot encoding. Columns: Number of unique molecules with a frequency above the threshold; number of unique molecules with a frequency below the threshold; percentage of the dataset that is None.

	labeling			reaction string		
	above	below	% None	above	below	% None
Reactants	40,020	0	25.7%	317,184	0	18.4%
Products	38,816	0	0.0%	382,850	0	0.0%
Solvents	29	204	40.0%	85	313	28.0%
Agents	48	447	56.2%	255	11,945	37.0%

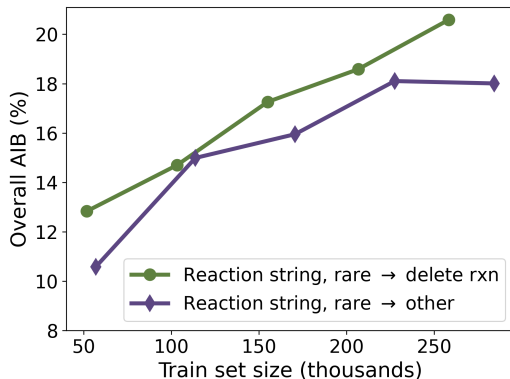


Figure 3: Scaling behaviour of different datasets with respect to overall top-3 AIB (%) for all solvents and agents (third row from Table 3.)

231 on the context. Reaction roles can more easily be identified when only considering one reaction
 232 class at [12], since this allows the mechanistic details of the reaction class [36, 37, 38] to be
 233 considered. Handling large and diverse datasets inevitably requires generalizations that may result in
 234 contradictions upon a more fine-grained inspection.

235 6.2 Order invariance

236 Although order of addition may play a role in wet lab chemistry, reaction prediction tasks are
 237 often cast as order invariant, where the goal is to predict a set of molecules. However, both of the
 238 architectures used for experimental validation of the ORDERly datasets are not agnostic to the ordering
 239 of the targets, since the neural networks used predict one molecule at a time in the OHE, and the
 240 transformers used predict one token at a time. Incorporating order invariance (and canonicalization)
 241 of the molecules into the loss calculation during training may allow for better generalisability of the
 242 predictive models, and is an exciting area for further study. It is worth noting that the evaluation
 243 metrics used throughout are order invariant.

244 7 Conclusions

245 In this work, we presented ORDERly, an open-source framework for preparing chemical reaction
 246 data stored in the Open Reaction Database (ORD) for machine learning applications. ORDERly was
 247 used to generate benchmark datasets for forward prediction (ORDERly-forward), retrosynthesis
 248 (ORDERly-retro), and condition prediction (ORDERly-condition) based on US patent data. Trans-
 249 former models were trained on the forward prediction and retrosynthesis datasets, and they were
 250 found to only generate invalid SMILES strings very infrequently, while also achieving similar test
 251 accuracy to that found in the literature on a held-out set of US patents. To further investigate model
 252 generalisation ORDERly was used to generate test sets from all non-patent data from ORD, and for
 253 the forward prediction task the accuracy was comparable, while the accuracy was slightly lower for
 254 the retrosynthesis task. The condition prediction task was used to investigate different strategies for
 255 assigning reaction roles and frequency filtering of the spectator molecules. When building datasets for
 256 condition prediction using the labeling in ORD, we found contamination of the inputs (reactants) with
 257 the outputs (agents), resulting in a problem that was unrealistically easy. We therefore chose to use
 258 chemically informed logic to better assign reaction roles for the ORDERly-condition benchmark.

259 All benchmarks and datasets experimented with in this work, as well as the code used to generate
 260 them, are freely available online, and we hope the benchmarks will make reaction prediction tasks
 261 more accessible. ORDERly presents a fully open-source pipeline to go from raw ORD data to a fully
 262 trained condition prediction model, allowing for an avenue to leverage the growing contributions to
 263 open source chemistry.

264 References

- 265 [1] Klavs F. Jensen, Connor W Coley, and Natalie S Eyke. Autonomous discovery in the chemical
266 sciences part I: Progress. *Angewandte Chemie International Edition*, September 2019.
- 267 [2] Philippe Schwaller, Teodoro Laino, Théophile Gaudin, Peter Bolgar, Christopher A. Hunter,
268 Costas Bekas, and Alpha A. Lee. Molecular Transformer: A Model for Uncertainty-Calibrated
269 Chemical Reaction Prediction. *ACS Central Science*, 5(9):1572–1583, September 2019.
- 270 [3] Connor W. Coley, Regina Barzilay, Tommi S. Jaakkola, William H. Green, and Klavs F. Jensen.
271 Prediction of Organic Reaction Outcomes Using Machine Learning. *ACS Central Science*,
272 3(5):434–443, May 2017.
- 273 [4] Zhengkai Tu and Connor W. Coley. Permutation Invariant Graph-to-Sequence Model for
274 Template-Free Retrosynthesis and Reaction Prediction. *Journal of Chemical Information and
275 Modeling*, 62(15):3503–3513, August 2022.
- 276 [5] Connor W. Coley, William H. Green, and Klavs F. Jensen. RDChiral: An RDKit Wrapper for
277 Handling Stereochemistry in Retrosynthetic Template Extraction and Application. *Journal of
278 Chemical Information and Modeling*, 59(6):2529–2537, June 2019.
- 279 [6] Alpha A. Lee, Qingyi Yang, Vishnu Sresht, Peter Bolgar, Xinjun Hou, Jacquelyn L. Klug-
280 McLeod, and Christopher R. Butler. Molecular Transformer unifies reaction prediction and
281 retrosynthesis across pharma chemical space. *Chemical Communications*, 55(81):12152–12155,
282 October 2019.
- 283 [7] Igor V. Tetko, Pavel Karpov, Ruud Van Deursen, and Guillaume Godin. State-of-the-art
284 augmented NLP transformer models for direct and single-step retrosynthesis. *Nature Communi-
285 cations*, 11(1):5575, November 2020.
- 286 [8] Umit V. Ucak, Islambek Ashyrmamatov, Junsu Ko, and Juyong Lee. Retrosynthetic reac-
287 tion pathway prediction through neural machine translation of atomic environments. *Nature
288 Communications*, 13(1):1186, March 2022.
- 289 [9] Yijia Sun and Nikolaos V Sahinidis. Computer-aided retrosynthetic design: fundamentals, tools,
290 and outlook. *Current Opinion in Chemical Engineering*, 35:100721, March 2022.
- 291 [10] Hanyu Gao, Thomas J. Struble, Connor W. Coley, Yuran Wang, William H. Green, and Klavs F.
292 Jensen. Using Machine Learning To Predict Suitable Conditions for Organic Reactions. *ACS
293 Central Science*, 4(11):1465–1476, November 2018.
- 294 [11] Philippe Schwaller, Riccardo Petraglia, Valerio Zullo, Vishnu H. Nair, Rico Andreas Haeusel-
295 mann, Riccardo Pisoni, Costas Bekas, Anna Iuliano, and Teodoro Laino. Predicting ret-
296 rosynthetic pathways using transformer-based models and a hyper-graph exploration strategy.
297 *Chemical Science*, 11(12):3316–3325, March 2020.
- 298 [12] Michael R. Maser, Alexander Y. Cui, Serim Ryou, Travis J. DeLano, Yisong Yue, and Sarah E.
299 Reisman. Multilabel Classification Models for the Prediction of Cross-Coupling Reaction
300 Conditions. *Journal of Chemical Information and Modeling*, 61(1):156–166, January 2021.
- 301 [13] Daniel Probst, Philippe Schwaller, and Jean-Louis Reymond. Reaction classification and yield
302 prediction using the differential reaction fingerprint DRFP. *Digital Discovery*, 1(2):91–97,
303 2022.
- 304 [14] Martin Fitzner, Georg Wuitschik, Raffael Koller, Jean-Michel Adam, and Torsten Schindler.
305 Machine Learning C–N Couplings: Obstacles for a General-Purpose Reaction Yield Prediction.
306 *ACS Omega*, 8(3):3017–3025, January 2023.

- 307 [15] Philippe Schwaller, Alain C. Vaucher, Teodoro Laino, and Jean-Louis Reymond. Prediction
308 of chemical reaction yields using deep learning. *Machine Learning: Science and Technology*,
309 2(1):015016, March 2021.
- 310 [16] A. Pomberger, A. A. Pedrina McCarthy, A. Khan, S. Sung, C. J. Taylor, M. J. Gaunt, L. Colwell,
311 D. Walz, and A. A. Lapkin. The effect of chemical representation on active machine learning
312 towards closed-loop optimization. *Reaction Chemistry & Engineering*, 7(6):1368–1379, May
313 2022.
- 314 [17] Nicholas H. Angello, Vandana Rathore, Wiktor Beker, Agnieszka Wołos, Edward R. Jira, Rafał
315 Roszak, Tony C. Wu, Charles M. Schroeder, Alán Aspuru-Guzik, Bartosz A. Grzybowski,
316 and Martin D. Burke. Closed-loop optimization of general reaction conditions for heteroaryl
317 Suzuki-Miyaura coupling. *Science*, 378(6618):399–405, October 2022.
- 318 [18] Connor J. Taylor, Kobi C. Felton, Daniel Wigh, Mohammed I. Jeraal, Rachel Grainger, Gianni
319 Chessari, Christopher N. Johnson, and Alexei A. Lapkin. Accelerated Chemical Reaction
320 Optimization Using Multi-Task Learning. *ACS Central Science*, April 2023.
- 321 [19] Elsevier. Reaxys, January 2009. <https://www.reaxys.com>.
- 322 [20] Steven M. Kearnes, Michael R. Maser, Michael Wlekinski, Anton Kast, Abigail G. Doyle,
323 Spencer D. Dreher, Joel M. Hawkins, Klavs F. Jensen, and Connor W. Coley. The Open Reaction
324 Database. *Journal of the American Chemical Society*, 143(45):18820–18826, November 2021.
- 325 [21] Daniel Lowe. Chemical reactions from US patents (1976-Sep2016), June 2017.
- 326 [22] Ziqiao Meng, Peilin Zhao, Yang Yu, and Irwin King. A Unified View of Deep Learning for
327 Reaction and Retrosynthesis Prediction: Current Status and Future Challenges. In *Electronic*
328 *proceedings of IJCAI 2023*, volume 6, pages 6723–6731, August 2023.
- 329 [23] Daniel S. Wigh, Jonathan M. Goodman, and Alexei A. Lapkin. A review of molecular represen-
330 tation in the age of machine learning. *WIREs Computational Molecular Science*, 12(5):e1603,
331 2022.
- 332 [24] Connor W. Coley, Dale A. Thomas, Justin A. M. Lummiss, Jonathan N. Jaworski, Christopher P.
333 Breen, Victor Schultz, Travis Hart, Joshua S. Fishman, Luke Rogers, Hanyu Gao, Robert W.
334 Hicklin, Pieter P. Plehiers, Joshua Byington, John S. Piotti, William H. Green, A. John Hart,
335 Timothy F. Jamison, and Klavs F. Jensen. A robotic platform for flow synthesis of organic
336 compounds informed by AI planning. *Science*, 365(6453), August 2019.
- 337 [25] Christos Kannas, Amol Thakkar, Esben Bjerrum, and Samuel Genheden. rxnutils – A Chemin-
338 formatics Python Library for Manipulating Chemical Reaction Data, August 2022. ChemRxiv,
339 2022. DOI: 10.26434/CHEMRXIV-2022-WT440-V2.
- 340 [26] Alain Vaucher and Hélder Lopes. RXN reaction preprocessing, May 2023.
341 <https://github.com/rxn4chemistry/rxn-reaction-preprocessing>.
- 342 [27] Amol Thakkar, Thierry Kogej, Jean-Louis Reymond, Ola Engkvist, and Esben Jannik Bjerrum.
343 Datasets and their influence on the development of computer assisted synthesis planning tools
344 in the pharmaceutical domain. *Chemical Science*, 11(1):154–168, 2020.
- 345 [28] Samuel Genheden, Per-Ola Norrby, and Ola Engkvist. AiZynthTrain: Robust, Reproducible,
346 and Extensible Pipelines for Training Synthesis Prediction Models. *Journal of Chemical*
347 *Information and Modeling*, 63(7):1841–1846, April 2023.
- 348 [29] Nadine Schneider, Nikolaus Stiefl, and Gregory A. Landrum. What’s What: The (Nearly)
349 Definitive Guide to Reaction Role Assignment. *Journal of Chemical Information and Modeling*,
350 56(12):2336–2346, December 2016.

- 351 [30] Bowen Liu, Bharath Ramsundar, Prasad Kawthekar, Jade Shi, Joseph Gomes, Quang
352 Luu Nguyen, Stephen Ho, Jack Sloane, Paul Wender, and Vijay Pande. Retrosynthetic Reaction
353 Prediction Using Neural Sequence-to-Sequence Models. *ACS central science*, 3(10):1103–1113,
354 October 2017.
- 355 [31] Wengong Jin, Connor Coley, Regina Barzilay, and Tommi Jaakkola. Predicting Organic Reaction
356 Outcomes with Weisfeiler-Lehman Network. In *Advances in Neural Information Processing*
357 *Systems*, volume 30, 2017.
- 358 [32] Kexin Huang, Tianfan Fu, Wenhao Gao, Yue Zhao, Yusuf Roohani, Jure Leskovec, Connor
359 Coley, Cao Xiao, Jimeng Sun, and Marinka Zitnik. Therapeutics Data Commons: Machine
360 Learning Datasets and Tasks for Drug Discovery and Development. *Proceedings of the Neural*
361 *Information Processing Systems Track on Datasets and Benchmarks*, 1, December 2021.
- 362 [33] Greg Landrum. The RDKit Documentation (accessed January 10 2020), 2006.
363 <https://www.rdkit.org/docs/>.
- 364 [34] Shuan Chen and Yousung Jung. Deep Retrosynthetic Reaction Prediction using Local Reactivity
365 and Global Attention. *JACS Au*, 1(10):1612–1620, October 2021.
- 366 [35] Arkadii Lin, Natalia Dyubankova, Timur I. Madzhidov, Ramil I. Nugmanov, Jonas Verhoeven,
367 Timur R. Gimadiev, Valentina A. Afonina, Zarina Ibragimova, Assima Rakhimbekova, Pavel
368 Sidorov, Andrei Gedich, Rail Suleymanov, Ravil Mukhametgaleev, Joerg Wegner, Hugo Ceule-
369 mans, and Alexandre Varnek. Atom-to-atom Mapping: A Benchmarking Study of Popular
370 Mapping Algorithms and Consensus Strategies. *Molecular Informatics*, 41(4):2100138, 2022.
- 371 [36] A. A. Thomas and S. E. Denmark. Pre-transmetalation intermediates in the Suzuki-Miyaura
372 reaction revealed: The missing link. *Science*, 352(6283):329–332, April 2016.
- 373 [37] Daniel S. Wigh, Matthieu Tissot, Patrick Pasau, Jonathan M. Goodman, and Alexei A. Lapkin.
374 Quantitative In Silico Prediction of the Rate of Protodeboronation by a Mechanistic Density
375 Functional Theory-Aided Algorithm. *The Journal of Physical Chemistry A*, March 2023.
- 376 [38] Donna G. Blackmond. Reaction Progress Kinetic Analysis: A Powerful Methodology for
377 Mechanistic Studies of Complex Catalytic Reactions. *Angewandte Chemie International*
378 *Edition*, 44(28):4302–4320, 2005.
- 379 [39] Yehia Amar, Artur M. Schweidtmann, Paul Deutsch, Liwei Cao, and Alexei Lapkin. Machine
380 learning and molecular descriptors enable rational solvent selection in asymmetric catalysis.
381 *Chemical Science*, 10(27):6697–6706, July 2019.
- 382 [40] Louis J. Diorazio, David R. J. Hose, and Neil K. Adlington. Toward a More Holistic Framework
383 for Solvent Selection. *Organic Process Research & Development*, 20(4):760–773, April 2016.
- 384 [41] Kobi C. Felton, Jan G. Rittig, and Alexei A. Lapkin. Summit: Benchmarking Machine Learning
385 Methods for Reaction Optimisation. *Chemistry–Methods*, 1(2):116–122, 2021.
- 386 [42] Adarsh Arun, Zhen Guo, Simon Sung, and Alexei A. Lapkin. Reaction Impurity Prediction
387 using a Data Mining Approach**. *Chemistry–Methods*, n/a(n/a):e202200062, March 2023.
- 388 [43] Connor W. Coley, Luke Rogers, William H. Green, and Klavs F. Jensen. Computer-Assisted Ret-
389 rosynthesis Based on Molecular Similarity. *ACS Central Science*, 3(12):1237–1245, December
390 2017.
- 391 [44] Monya Baker. 1,500 scientists lift the lid on reproducibility. *Nature*, 533(7604):452–454, May
392 2016.

393 Appendix (ORDERly: Datasets and benchmarks for chemical reaction data)

394 A Dataset extraction and cleaning

395 In the main paper, we describe the "labeling" and "reaction string" datasets; in the code this
396 is denoted by `trust_labeling=True`, and `trust_labeling=False`, respectively. We also
397 presented two different strategies for dealing with rare molecules, either "rare→"other" and
398 "rare→delete rxn", these are denoted in the code as `map_rare_molecules_to_other=True`, and
399 `map_rare_molecules_to_other=False`, respectively. There are a number of other tuneable pa-
400 rameters in the scripts, and below we explain how default values were chosen for each of these.

401 A.1 Extraction script

402 There are three fields in the Open Reaction Database schema to extract molecules from: the input, the
403 outcome and the reaction string. Molecules in the reaction string are represented as SMILES, while
404 molecules in the input and outcome field can be represented with a number of different representations,
405 including SMILES, InChI, and plain text English names. When extracting molecules from the input
406 or outcome field, the preferred representation was SMILES. However, how should the situation where
407 only an English name exists be dealt with? It is tempting to check whether the representation is
408 interpretable by RDKit (potentially implying that the molecular representation was mislabeled as
409 a name rather than SMILES), however, this can lead to unexpected behaviour. As an example, the
410 string, "1400C", was encountered as the name for a molecule, should this be interpreted as a graphene
411 structure, a typo for carbon-14, a typo for 1400°C, or simply carbon? Another situation which was
412 encountered was BOC; this *is* a resolvable smiles string, representing boron oxygen and carbon
413 bonded together, however, in context, it was actually referring to a BOC group (tert-butyloxycarbonyl
414 protecting group). Another example of unintended behaviour is the case of II, which could mean
415 diiodine, but also mark the second step/item when counting. Therefore, when the user decides not to
416 trust the labeling of the molecules, molecules only represented with a plain text name were ignored,
417 to avoid ambiguity.

418 The extraction script generates the relevant data from each ORD file, and allows for the following
419 customization Note that we only mention the arguments that materially affect the science/logic of
420 how cleaning is done.

- 421 • `trust_labeling`: If True, maintain the labeling of the data in ORD. If False: chemical
422 logic (described extensively in the paper) is applied to the reaction string to determine the
423 reaction role of molecules.
- 424 • `solvents_path`: If the user does not trust the labeling, all agent molecules are cross-
425 checked against a set of industrially relevant solvents, and any matches are re-labeled as
426 solvents. See section A.1.1 for how this set of solvents was constructed.
- 427 • `name_contains_substring`: Only extract filenames from ORD that includes this string.
428 If left empty will not search for anything, and if set to None it will extract data from all ORD
429 files in the designated folder. For example, setting `name_contains_substring="uspto"`
430 will grab all files that have "uspto" in the file name (i.e. the USPTO data).
- 431 • `inverse_substring`: The inverse of `name_contains_substring`, e.g. setting
432 `inverse_substring="uspto"` will grab everything *except* the USPTO data.

433 A.1.1 Building a set of solvents

434 The solvents set can be found in `orderly/data/solvents.csv` in the ORDERly GitHub repository.
435 The set was created from the intersection of solvents from the following three sources:

- 436 • Machine learning and molecular descriptors enable rational solvent selection in asymmetric
437 catalysis (458 solvents) [39]

- 438 • ACS Green Chemistry: Solvent Selection Tool (272 solvents) [40]
- 439 • Summit GitHub Repository (115 solvents) [41]

440 After the data from these three sources were concatenated into a new CSV files, the solvents
441 were filtered by: making all solvent names lower case, stripping spaces, and then removing du-
442 plicate names. (Before removing duplicates: 458+272+115=845 solvents. After removing dupli-
443 cates: 615 solvents.) Then Pura was run to resolve the solvent name where no SMILES string
444 was available. Each solvent (with no SMILES string) was represented with up to four different
445 names: three English solvent names (synonym names) and one CAS number. Pura was used with
446 `services=[PubChem(autocomplete=True), Opsin(), CIR()]` and `agreement=2` on each En-
447 glish name, and `services=[CAS()]` with `agreement=1` on the CAS number. This yielded up to four
448 different SMILES strings for each solvent. SMILES strings with full agreement for a solvent were
449 trusted, and any rows with disagreement between the SMILES strings (≈ 40 solvents) were resolved
450 by hand. The final solvents set is a CSV file with seven columns: up to three different English solvent
451 names (synonyms), a CAS number, a chemical formula, SMILES, and finally the source.

452 An obvious drawback of identifying solvents by crosschecking against a curated set is that the set
453 naturally will be incomplete; there are unfathomably many different organic molecules, and it is
454 unclear how many of these could act as solvents. However, not distinguishing between solvents and
455 agents may make the learning task more difficult for machine learning models, and using the labeling
456 that already exists in ORD was routinely found to be inaccurate. In practice, the vast majority of
457 solvents used in industry and academia are inspired by what has previously proven successful, and
458 thus the solvents set curated for this work is likely going to capture a majority of solvents. Another
459 difficulty is that the role of solvent molecules may depend on the context (e.g. polar protic solvents
460 may contribute protons to the product, in which case the role of the molecule becomes murky (i.e.
461 is it a reactant since it contributed atoms to the product, is it a solvent since it dissolved the (other)
462 reactants, or is it a reagent since it acts like an acid?).

463 A.2 Cleaning script

- 464 • `remove_reactions_with_no_reactants [bool]`: Self-explanatory
- 465 • `remove_reactions_with_no_products [bool]`: Self-explanatory
- 466 • `consistent_yield [bool]`: If True, removes reactions that have yields that do not make
467 sense, e.g. if any individual yields, or the sum of yields, is outside of [0%; 100%] (reactions
468 with no yields are kept).
- 469 • `num_reactant, num_product, num_solv, num_agent, num_cat, num_reag [int]`:
470 The maximum number of components allowed of the specified type in a reaction. E.g.
471 if `num_solv=2` any reactions with 3 or more solvents will be dropped from the DataFrame.
472 See section C for how the default values were chosen.
- 473 • `min_frequency_of_occurrence [int]`: The frequency of molecules across all columns
474 of the same type (e.g. solvents) are counted, and any reactions containing molecules below
475 the frequency cutoff are dealt with in accordance with `map_rare_molecules_to_other`.
476 See section C for how the default values were chosen.
- 477 • `map_rare_molecules_to_other [bool]`: If False, any reactions containing molecules
478 that fall below the threshold will be deleted. If True, the rare molecules will be mapped to a
479 string "other", allowing us to keep the reactions in the dataset. This behaviour can be shut
480 off simply by setting `min_frequency_of_occurrence=0`.
- 481 • `set_unresolved_names_to_none_if_mapped_rxn_str_exists_else_del_rxn,`
482 `remove_rxn_with_unresolved_names, set_unresolved_names_to_none [bool]`:
483 These three bools control the handling of unresolvable names (i.e. names that are
484 unresolvable by RDKit, and do not exist in our manually curated name resolution dictionary,
485 and at most one of them can be True (if all are set to False, unresolvable names are kept in
486 the dataset.) While the second and third bool are self-explanatory, this is the logic applied

Table 5: Comparison between different datasets for retrosynthesis and forward prediction.

Dataset	Size	Split	Reference
USPTO-50K	50 016	Random	[29]
USPTO-MIT	479 035	Random	[31]
USPTO-full	997 415 ²	Random	[32]
ORDerly-retro	941 566	Random	This work
ORDerly-forward	919 231	Random	This work

487 if the first bool is True: if a reaction contains a mapped reaction, the reaction is seen as
 488 quite trustworthy, and therefore the unresolvable names can safely be set to None, while the
 489 remaining data associated with that reaction is maintained; if a reaction does not have an
 490 associated mapped reaction, the presence of an unresolvable name is a red flag casting
 491 doubt on the veracity of that reaction, and thus the whole reaction (a row in the DataFrame)
 492 is removed).

493 A.3 Further justification for cleaning thresholds

494 A.3.1 Condition prediction benchmark

- 495 • **Reactant filtering:** Reactions with more than two reactants were filtered out, since they are
 496 likely to be multi-step reactions or complex one-pot reactions (tri-molecular mechanisms
 497 are exceedingly rare in chemistry).
- 498 • **Product filtering:** Reactions with multiple products were also filtered out since nearly
 499 all reactions in USPTO only report one product (see Figure A5); predicting reaction side
 500 products and impurities remains an active area of research [42], and thus fell beyond the
 501 scope of ORDerly.
- 502 • **Solvent and agent filtering:** Thresholds for the number of spectator molecules was set at
 503 two solvents and three agents to have the same number of categorical variables as in the
 504 model of Gao et al. [10].
- 505 • **No conditions filtering:** Reactions will not work without a solvent, and will usually require
 506 an agent. There are exceptions to this (e.g. the Diels-Alder reaction), however, the number
 507 of reactions with an erroneous recording of no agents is likely going to outnumber the
 508 amount of genuine exceptions. These filtering steps imply that a model trained on the
 509 ORDerly-condition benchmark may be ill-equipped to deal with reaction impurities or make
 510 predictions for reactions with no agents. It is worth noting that these drawbacks may be
 511 relatively inconsequential, since a skilled chemist is unlikely to query a model to predict
 512 agents for a class of reaction that requires no agents.
- 513 • **Not predicting temperature:** Only 192k out of 323k reactions in the ORDerly-condition
 514 training set contain a temperature, of which over half report 25C. Filtering away reactions
 515 without a temperature would leave a much smaller dataset, and we do not believe that it is
 516 reasonable to assume that reactions without a reported temperature were performed at room
 517 temperature.

518 A.3.2 Forward prediction and single-step retrosynthesis benchmarks

519 The ORDerly-retro dataset is compared to other standard forward prediction and retrosynthesis
 520 datasets in Table 5. USPTO-50K was created by Schneider *et al.* for testing reaction role assignment
 521 [29]. They used NameRxn to assign reaction classes to all the reactions in the dataset. Liu *et al.* [30]
 522 then used the USPTO-50K for benchmarking their retrosynthesis model, however, they did not use
 523 the reaction classes to create a split, and instead opted for a random split. Coley *et al.* [43] is often
 524 cited for their train/test split of USPTO-50K. USPTO-MIT is a larger set that was introduced by Jin
 525 *et al.* [31].

- 526 • **Forward prediction:** A small number of reactions in USPTO reported two products, and for
527 the forward prediction dataset we allowed up to two products and three reactants, solvents,
528 and agents.
- 529 • **Retrosynthesis prediction:** In retrosynthesis prediction the goal is to predict reactants
530 that can be used to form a desired product. To ensure that the difficulty of the task was
531 reasonable, we limit reactions to having one product and two reactants, such that the models
532 only have to learn how to break one molecule into two, and not consider e.g. multi-product
533 or multi-step reactions. Only product and reactant molecules were used in the retrosynthesis
534 dataset, so there were no restrictions in the number of solvents and agents.

535 B Further experimental details

536 B.1 Condition prediction with neural networks

537 The code from *Gao et al.* [10] was used for training condition prediction models. The hyperparameters
538 in Table 6 were used, which reflect those used in the original paper. Training on an A10G required 30
539 minutes or less for a full training run.

Table 6: Hyperparameters used for training condition prediction models

batch size	512
learning rate	0.01
hidden size 1	1024
hidden size 2	100
dropout	0.2
fingerprint size	2048

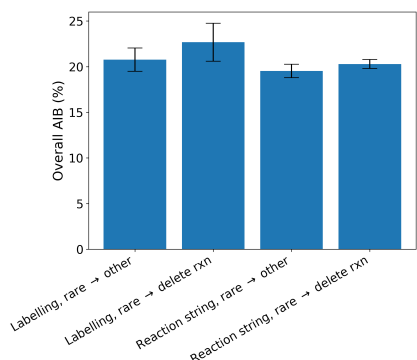


Figure 4: AIB (%) on the test sets for each training dataset. Error bars are with respect to the random seed in splitting the training and validation data (test data stayed the same).

540 B.2 Forward prediction and retrosynthesis prediction with transformers

541 Most of the hyperparameters used in the Molecular Transformer architecture (see Table 7) were the
542 defaults suggested by Schwaller *et al.* [2] (GitHub). The transformer models were trained for around
543 35 hours (approximately 600 epochs).

Table 7: Hyperparameters for Molecular Transformer.
Training

seed	42
param_init	0
param_init_glorot	
max_generator_batches	32
batch_size	4096
batch_type	tokens
normalization	tokens
max_grad_norm	0
accum_count	4
optim	adam
adam_beta1	0.9
adam_beta2	0.998
decay_method	noam
warmup_steps	8000
learning_rate	2
label_smoothing	0.0
layers	4
rnn_size	256
word_vec_size	256
encoder_type	transformer
decoder_type	transformer
dropout	0.1
position_encoding	
share_embeddings	
global_attention	general
global_attention_function	softmax
self_attn_type	scaled-dot
heads	8
transformer_ff	2048
Inference	
batch_size	512
replace_unk	
max_length	200
beam_size	5

544 C ORDERly benchmark statistics

545 C.1 Number of components

546 Figure 5 shows the distribution in the number of components of the unfiltered datasets, allowing us
 547 to compare the reaction string datasets to the labeling datasets. The distributions look quite similar
 548 for products and solvents. However, the distributions are different for reactants and agents/catalysts,
 549 which can be explained by reagents routinely being labelled as reactants in ORD.

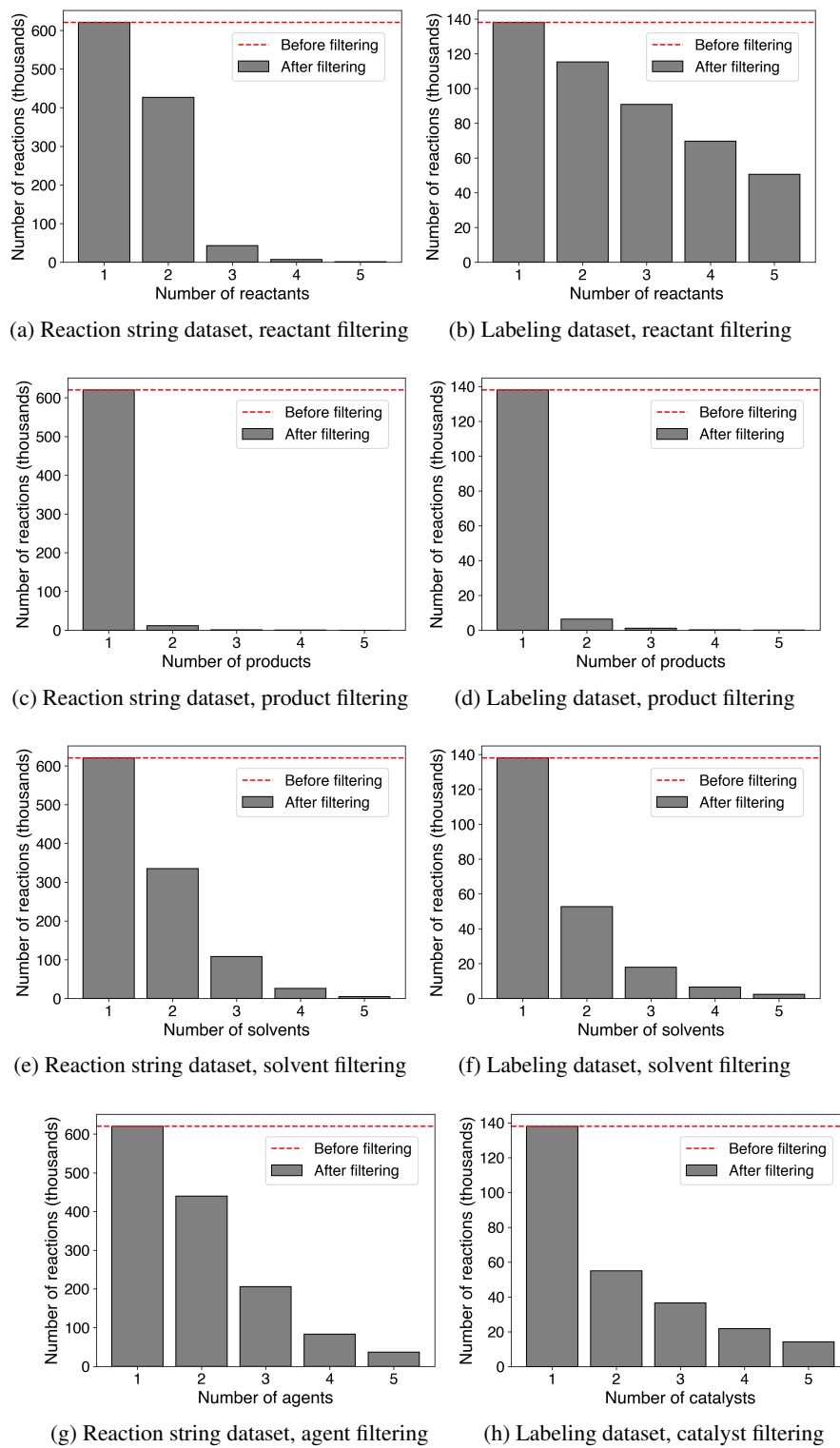


Figure 5: Distribution of the number of components between the reaction string and labeling datasets. There are no reagents in the labeling dataset, so after filtering excess catalysts were re-labelled as reagents.

550 C.2 Minimum frequency of occurrence

551 Figure 6 shows how many reactions would be left in the reaction string and labeling datasets as
552 a function of the minimum frequency of occurrence. The minimum frequency of occurrence is
553 the threshold applied to the spectator molecules (solvents, agents, reagents, agents, catalysts) to be
554 considered rare, and any reactions containing a rare molecule will be deleted if (rare→delete rxn).

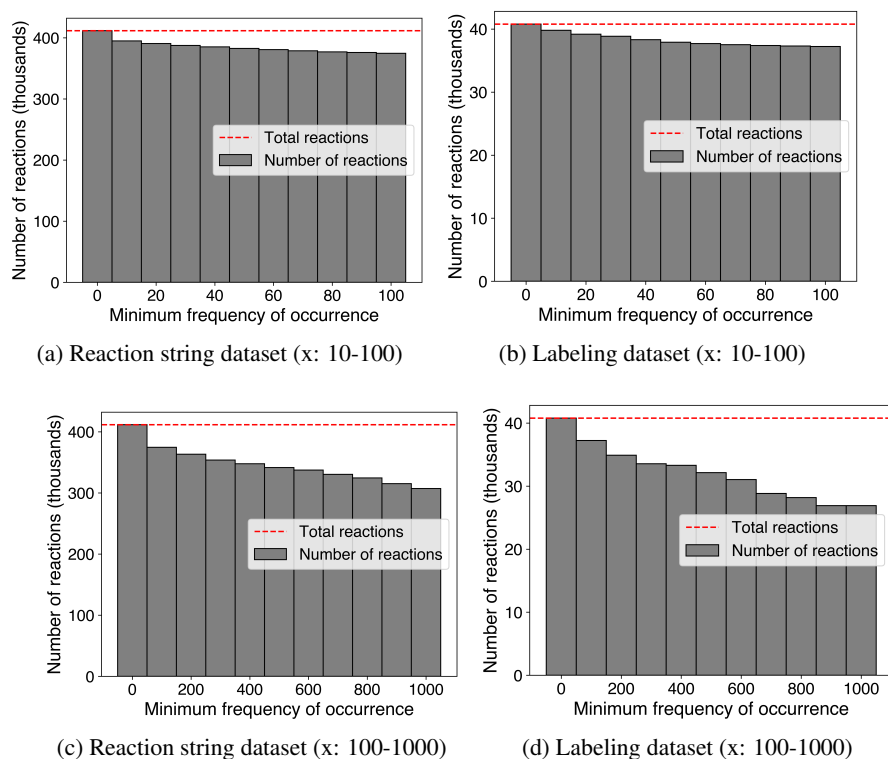


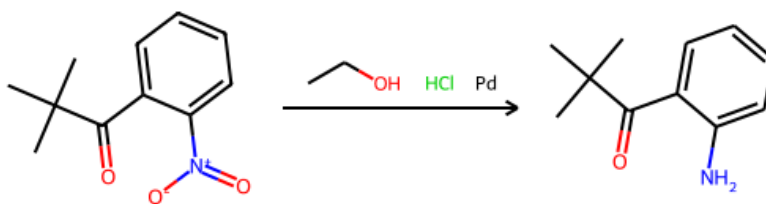
Figure 6: Impact on dataset size by changing the minimum frequency of occurrence.

555 C.3 Molecule popularity

556 Figure 7 shows the distribution of occurrence of the top 100 most popular molecules across the
557 different categories of molecules for the labeling and rxn string datasets. Across categories, the
558 reaction string dataset is more diverse and not as heavily dominated by the most popular component.
559 It is also interesting that the most popular molecules between the datasets are not the same, despite
560 being based on the same raw data.

561 D Example reaction instances and predictions

562 In this section, we give examples of reactions that are in both the trust labeling and reaction string
563 datasets (Table 3) to demonstrate the differences between the different cleaning methodologies.



Reaction string input:

CC(C)(C)C(=O)c1cccc1[N+](=O)[O-]>CCO.Cl.[Pd]>CC(C)(C)C(=O)c1cccc1N

	Reaction string dataset (This work)	Trust labelling dataset
Reactants		
Products		
Ground truth solvents		HCl
Ground truth agents	Pd, HCl	Pd
Predicted solvents	✓	✗
Predicted agents	Pd ✓	Pd ✓

Reaction type: Reduction (of a nitro group)

Comment: Trust labelling predicted the wrong solvent, which, however, can still serve as a solvent due to similar properties (polar and protic). It must be noted that the agent prediction was incomplete - no strategy predicted HCl as an agent which is crucial to enable the reaction and serve as a proton source.



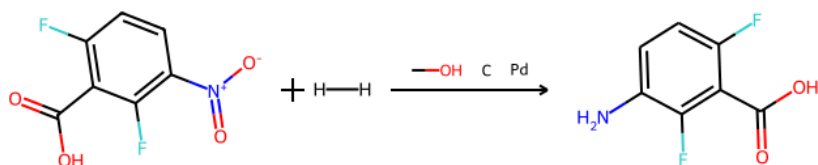
Correct prediction



Incorrect prediction



Partially correct prediction



Reaction string input:

O=C(O)c1c(F)ccc([N+](=O)[O-])c1F.[H][H]>CO.[C].[Pd]>Nc1ccc(F)c(C(=O)O)c1F

	Reaction string dataset (This work)	Trust labelling dataset
Reactants		H-H
Products		
Ground truth solvents		
Ground truth agents	C, Pd, H-H	C, Pd
Predicted solvents		
Predicted agents	Pd	Pd

Reaction type: Hydrogenation (of a nitro group)

Comment: Hydrogen gas can be either categorized as reactant or as agent – here the approaches vary depending on the dataset. In both cases ethanol is predicted which, however, can still serve as a solvent due to similar properties (polar and protic). It must be noted that the agent prediction was incomplete - no strategy predicted H₂ as an agent which is crucial to enable the reaction and serve as a hydrogen source.



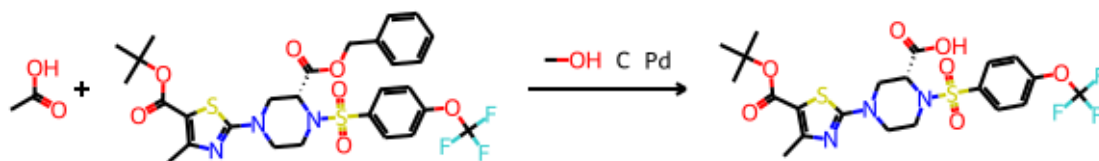
Correct prediction



Incorrect prediction



Partially correct prediction



Reaction string input:

```
CC(=O)O.Cc1nc(N2CCN(S(=O)(=O)c3ccc(OC(F)(F)F)cc3)[C@@H](C(=O)OCc3cccc3)C2)sc1C(=O)OC(C)(C)C>CO.[C].[Pd]>Cc1nc(N2CCN(S(=O)(=O)c3ccc(OC(F)(F)F)cc3)[C@@H](C(=O)O)C2)sc1C(=O)OC(C)(C)C
```

	Reaction string dataset (This work)	Trust labelling dataset
Reactants		
Products		
Ground truth solvents		
Ground truth agents	C, Pd	C, Pd
Predicted solvents		
Predicted agents	Pd	Pd

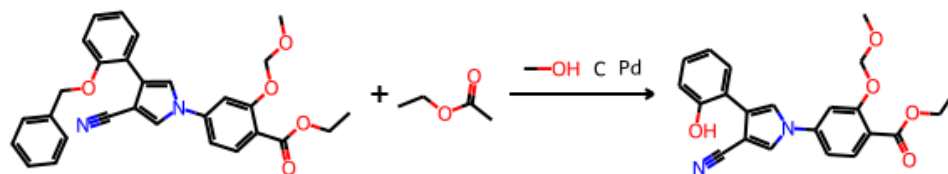
Reaction type: Acidic ester cleavage

Comment: We observed differences in categorizing the acetic acid as either reactant or solvent. Chemically, it should be considered a reactant or agent. In both cases ethanol is predicted which can still serve as a solvent due to similar properties (polar and protic). In the case that acetic acid is not passed as reactant the model should also predict it as agent.

Correct prediction

Incorrect prediction

Partially correct prediction



Reaction string input:

CCOC(=O)c1ccc(-n2cc(C#N)c(-c3ccccc3OCc3ccccc3)c2)cc1OCOC.CCOC(C)=O
>CO.[C].[Pd]>CCOC(=O)c1ccc(-n2cc(C#N)c(-c3ccccc3O)c2)cc1OCOC

	Reaction string dataset (This work)	Trust labelling dataset
Reactants		
Products		
Ground truth solvents		
Ground truth agents	C, Pd	C, Pd
Predicted solvents	✗	✓
Predicted agents	Pd ✓	Pd ✓

Reaction type: Ether cleavage (cleaving an Obn protection group)

Comment: Acetyl acetate is categorized either as solvent or reactant. Here both roles makes sense chemically. For the prediction using reaction string dataset it must be noted that while EtOH is predicted, the ground truth solvent is ethylacetate. However, under acidic conditions acetyl acetate can fall apart into acetic acid and EtOH.



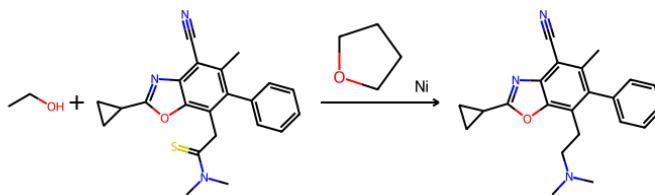
Correct prediction



Incorrect prediction



Partially correct prediction



Reaction string input:

CCO.Cc1c(-c2ccccc2)c(CC(=S)N(C)C)c2oc(C3CC3)nc2c1C#N>C1CCOC1.[Ni]>
Cc1c(-c2ccccc2)c(CCN(C)C)c2oc(C3CC3)nc2c1C#N

	Reaction string dataset (This work)	Trust labelling dataset
Reactants		
Products		
Ground truth solvents		
Ground truth agents	Ni	Ni
Predicted solvents	✓	✗
Predicted agents	Pd ✗	Pd ✗

Reaction type: Corey Seebach reaction

Comment: Ethanol is categorized either as solvent or reactant - both roles makes sense chemically. Within the prediction, trust labelling predicted ethyl acetate which is uncommon for this transformation. Using the reaction string dataset, THF was predicted which is correct, however, the initial data also contained EtOH. Pd has been predicted in both cases as agent which is incorrect.

✓ Correct prediction ✗ Incorrect prediction ✓ Partially correct prediction

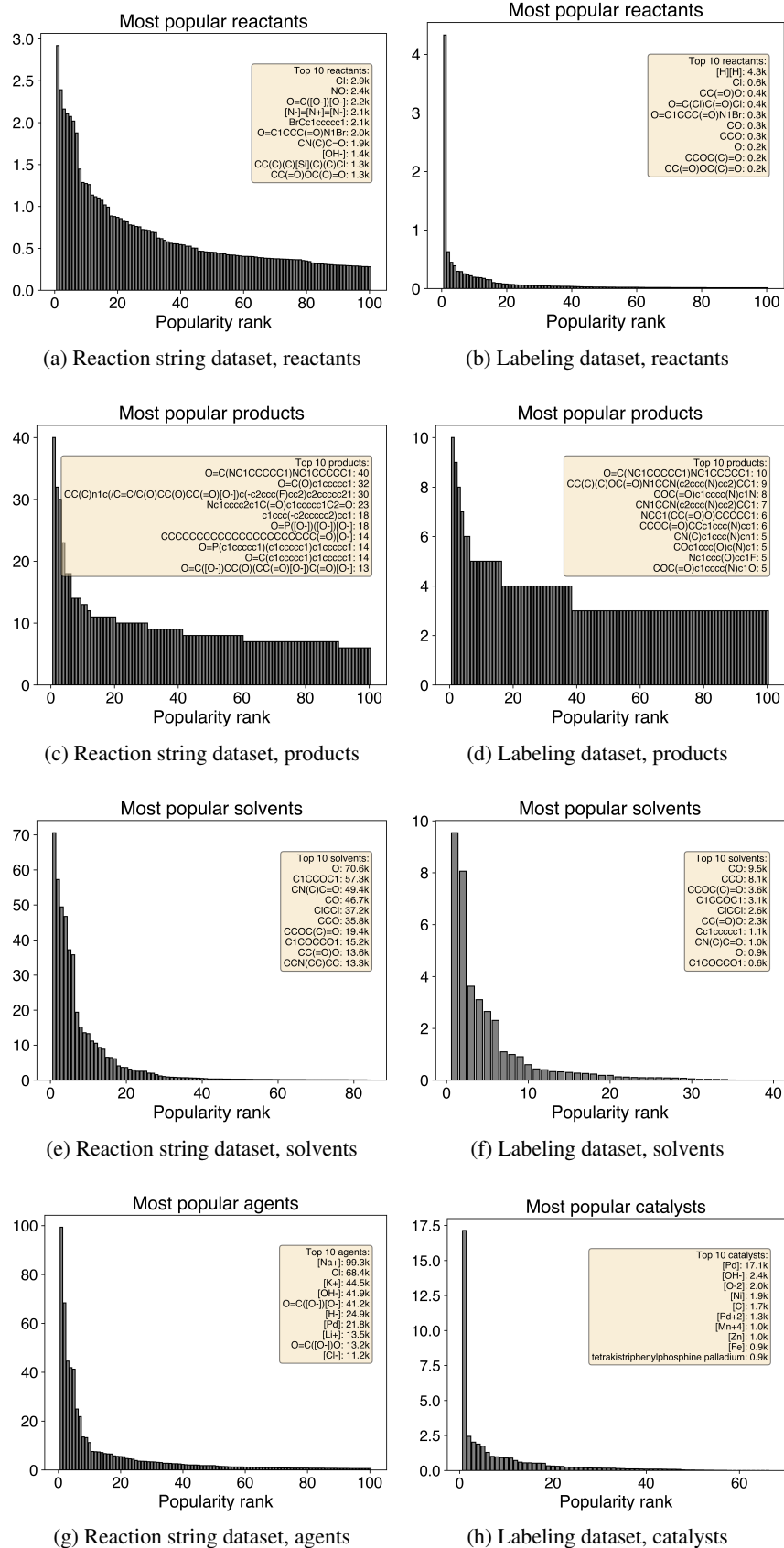


Figure 7: Frequency of occurrence of the most popular molecules. NULL has been removed, reagents and catalysts have been merged.

569 E Datasheet for ORDERly dataset

570 E.1 Motivation

571 Q1: **For what purpose was the dataset created?** Was there a specific task in mind? Was there
572 a specific gap that needed to be filled? Please provide a description.

573 – The datasets were created to facilitate building machine learning models for prediction
574 of reaction products, retrosynthesis, and reaction conditions in chemical synthesis,
575 particularly in the context of the pharmaceutical industry. There was a need of a
576 clean, high-quality reaction condition benchmark dataset, in addition to a need for an
577 open-source repository for cleaning reactions, and an investigation of how decisions
578 made during cleaning impact the usefulness of the model that is trained on the datasets.
579 ORDERly solves all three of these issues. The code for ORDERly, and the raw data
580 used to generate the ORDERly benchmark datasets, are both open-source, making the
581 benchmark generation accessible and reproducible.

582 Q2: **Who created the dataset (e.g., which team, research group) and on behalf of which**
583 **entity (e.g., company, institution, organization)?**

584 – ORDERly was built by researchers from the group of Anon at Institution.

585 Q3: **Who funded the creation of the dataset?** If there is an associated grant, please provide the
586 name of the grantor and the grant name and number.

587 – This work is co-funded by UCB Pharma and Engineering and Physical Sciences Re-
588 search Council via project EP/S024220/1 EPSRC Centre for Doctoral Training in
589 Automated Chemical Synthesis Enabled by Digital Molecular Technologies. This
590 project was co-funded by European Regional Development Fund via the project "Inno-
591 vation Centre in Digital Molecular Technologies".

592 Q4: **Any other comments?**

593 – No.

594 E.2 Composition

595 Q5: **What do the instances that comprise the dataset represent (e.g., documents, photos,**
596 **people, countries)?** *Are there multiple types of instances (e.g., movies, users, and ratings;*
597 *people and interactions between them; nodes and edges)? Please provide a description.*

598 – Eight datasets were presented in this work. Each dataset was saved in Apache Parquet
599 format, and has the following column groups:

- 600 * Reaction SMILES string (string), is_mapped (bool)
- 601 * Reactants & products (SMILES strings)
- 602 * Solvents and agents (rxn string data), or solvents, catalysts, and reagents (labeling
603 data) (SMILES strings)
- 604 * Temperature, reaction time, yield (floats)
- 605 * Procedure details (string)
- 606 * Grant date (datetime), date of experiment (datetime), file name (string)

607 Q6: **How many instances are there in total (of each type, if appropriate)?**

608 – The number of reactions in each dataset is outlined in detail in Table 1.

609 Q7: **Does the dataset contain all possible instances or is it a sample (not necessarily random)**
610 **of instances from a larger set?** *If the dataset is a sample, then what is the larger set? Is the*
611 *sample representative of the larger set (e.g., geographic coverage)? If so, please describe*
612 *how this representativeness was validated/verified. If it is not representative of the larger set,*
613 *please describe why not (e.g., to cover a more diverse range of instances, because instances*
614 *were withheld or unavailable).*

615 – All the data in ORD was used to generate the datasets presented in this paper. Datasets
616 A-F were built from the subset of ORD belonging to USPTO (1.7m reactions in total),
617 while datasets G-H were built on the subset of data from ORD that do not belong to
618 USPTO (91k reactions in total, as of August 2023).

619 **Q8: What data does each instance consist of? “Raw” data (e.g., unprocessed text or images)**
620 *or features? In either case, please provide a description.*

621 – Chemical reaction data stored in ORD is structured like a json/dictionary, with strings
622 and floats as the values. The values that are relevant to ORDERly were discussed in
623 response to Q5. A full description of the data stored in ORD is available elsewhere
624 [20].

625 **Q9: Is there a label or target associated with each instance? If so, please provide a description.**

626 – There is a label associated with molecules in ORD, and in this work we show the
627 pitfalls of relying on this label, and present ORDERly to more robustly assign labels.
628 The targets are the reaction conditions (solvents, agents, catalysts, reagents).

629 **Q10: Is any information missing from individual instances? If so, please provide a description,**
630 *explaining why this information is missing (e.g., because it was unavailable). This does not*
631 *include intentionally removed information, but might include, e.g., redacted text.*

632 – Many reactions were missing temperature, reaction time, and yield data; this is likely
633 due to this information not being recorded by the experimentalist, or not extracted
634 when the information was scraped from a patent/paper.

635 **Q11: Are relationships between individual instances made explicit (e.g., users’ movie ratings,**
636 **social network links)? If so, please describe how these relationships are made explicit.**

637 – Each row contains information for a single step chemical reaction. The only explicit link
638 between reactions is the year they were performed or the year that the corresponding
639 patent was granted. The year a chemical reaction was performed may imply some
640 degree of chemical information, since chemical reactions of a certain type obviously
641 could not have been performed before they were invented. Furthermore, "hype" around
642 a particular type of reaction may influence how often certain reaction classes are used
643 through time. For these reasons, a time-based split can be viewed as a (somewhat poor)
644 proxy for a reaction class split. There is a column in the dataset containing the year
645 that the grant was awarded, and another column for time of experiment.

646 **Q12: Are there recommended data splits (e.g., training, development/validation, testing)? If**
647 *so, please provide a description of these splits, explaining the rationale behind them.*

648 – We recommend using a random split of the ORDERly benchmarks, and provide pre-split
649 data to ensure that ML researchers using the benchmark use the same train/test split.
650 There are three data splits that would make sense on a chemical reactions dataset: a
651 random split, a time split, and the reaction class split. A reaction class split would
652 require models to generalise to unseen reaction classes (as opposed to unseen reactions
653 of the same class), making the prediction task much more difficult. As explained above
654 (Q16), using a time split would effectively just serve as a proxy for a reaction class split,
655 and is therefore not desirable. There are a number of reasons for the random split being
656 preferred over the reaction class split: 1) A reaction class split would need to either use
657 an ML clustering algorithm (which usually work quite well, but cannot be viewed as a
658 ground-truth split), or using proprietary software based on manually curated chemistry
659 rules (which would mean that the full pipeline is no longer fully open source and
660 reproducible). 2) The reaction prediction task is already difficult enough with a random
661 split (e.g. considering our top-3 accuracy of sub 50%, and models trained on a random
662 split are still able to provide value even if they can only make predictions on reaction
663 classes that they have seen before - the reaction classes represented in the dataset will
664 likely be the most popular reaction classes, and therefore also those most likely to be
665 queried by the end user.

- 666 Q13: **Are there any errors, sources of noise, or redundancies in the dataset?** *If so, please*
667 *provide a description.*
- 668 – The ORDerly-condition, ORDerly-forward, and ORDerly-retro datasets are
669 generated from the USPTO dataset, which is a dataset made from chemical reac-
670 tions from US Patents. When a molecule is patented, it is also a requirement to publish
671 the synthesis pathway to produce the molecule, and it is from these synthesis pathways
672 that reactions are extracted. To avoid giving away proprietary information there is an
673 incentive to use already published "industry standard" reaction conditions in the patent
674 application; furthermore, the "first to file" nature of the US patent system means there
675 is an incentive to apply for patents as soon as possible. These two factors may bias
676 the reactions in the USPTO dataset towards being unoptimized, low-yielding reactions
677 that can also be found elsewhere. In fact, we observed that $\approx 40\%$ of reactions were
678 dropped because they were duplicates (see Table 1), indicating that many reactions
679 are executed at "standard conditions" for a particular class of reaction instead of being
680 optimized for the specific reactants.
 - 681 – Reproducibility is known to be difficult in chemistry[44], which implies a base-level of
682 noise in the dataset.
- 683 Q14: **Is the dataset self-contained, or does it link to or otherwise rely on external resources**
684 **(e.g., websites, tweets, other datasets)?** *If it links to or relies on external resources, a) are*
685 *there guarantees that they will exist, and remain constant, over time; b) are there official*
686 *archival versions of the complete dataset (i.e., including the external resources as they*
687 *existed at the time the dataset was created); c) are there any restrictions (e.g., licenses, fees)*
688 *associated with any of the external resources that might apply to a future user? Please*
689 *provide descriptions of all external resources and any restrictions associated with them, as*
690 *well as links or other access points, as appropriate.*
- 691 – The ORDerly datasets are self-contained. To be able to reproduce cleaning of ORD
692 data, the ORD data will naturally need to continue to exist. ORD was built to be an
693 open-source tool, so there should not be any restrictions on its use in the future.
- 694 Q15: **Does the dataset contain data that might be considered confidential (e.g., data that is**
695 **protected by legal privilege or by doctor–patient confidentiality, data that includes the**
696 **content of individuals’ non-public communications)?** *If so, please provide a description.*
- 697 – No.
- 698 Q16: **Does the dataset contain data that, if viewed directly, might be offensive, insulting,**
699 **threatening, or might otherwise cause anxiety?** *If so, please describe why.*
- 700 – No.
- 701 Q17: **Does the dataset relate to people?** *If not, you may skip the remaining questions in this*
702 *section.*
- 703 – No.
- 704 Q18: **Does the dataset identify any subpopulations (e.g., by age, gender)?**
- 705 – No.
- 706 Q19: **Is it possible to identify individuals (i.e., one or more natural persons), either directly or**
707 **indirectly (i.e., in combination with other data) from the dataset?** *If so, please describe*
708 *how.*
- 709 – No.
- 710 Q20: **Does the dataset contain data that might be considered sensitive in any way (e.g., data**
711 **that reveals racial or ethnic origins, sexual orientations, religious beliefs, political**
712 **opinions or union memberships, or locations; financial or health data; biometric or**
713 **genetic data; forms of government identification, such as social security numbers;**
714 **criminal history)?** *If so, please provide a description.*

715 – No.

716 **Q21: Any other comments?**

717 – No.

718 **E.3 Collection process**

719 **Q22: How was the data associated with each instance acquired?** Was the data directly ob-
720 servable (e.g., raw text, movie ratings), reported by subjects (e.g., survey responses), or
721 indirectly inferred/derived from other data (e.g., part-of-speech tags, model-based guesses
722 for age or language)? If data was reported by subjects or indirectly inferred/derived from
723 other data, was the data validated/verified? If so, please describe how.

724 – The raw data of each instance (reaction) was extracted from United States Patents to
725 form the "USPTO dataset" [21]. The USPTO dataset was parsed into ORD format [20],
726 where we extracted it from. ORD does contain additional data, beyond the USPTO
727 dataset. Other reactions in ORD are contributed by chemists in academia and industry.

728 **Q23: What mechanisms or procedures were used to collect the data (e.g., hardware apparatus
729 or sensor, manual human curation, software program, software API)?** How were these
730 mechanisms or procedures validated?

731 – Data in the ORD database is readily downloadable through the GitHub repository.

732 **Q24: If the dataset is a sample from a larger set, what was the sampling strategy (e.g.,
733 deterministic, probabilistic with specific sampling probabilities)?**

734 – See Q7.

735 **Q25: Who was involved in the data collection process (e.g., students, crowdworkers, contrac-
736 tors) and how were they compensated (e.g., how much were crowdworkers paid)?**

737 – N/A.

738 **Q26: Over what timeframe was the data collected? Does this timeframe match the creation
739 timeframe of the data associated with the instances (e.g., recent crawl of old news
740 articles)?** *If not, please describe the timeframe in which the data associated with the
741 instances was created.*

742 – The reactions in the USPTO dataset are from patents which were published between
743 1976 and September 2016. The USPTO dataset was parsed into ORD in 2020. Addi-
744 tional reactions not from patents have since been added to ORD. ORDERly was built in
745 2023.

746 **Q27: Were any ethical review processes conducted (e.g., by an institutional review board)?**
747 *If so, please provide a description of these review processes, including the outcomes, as well
748 as a link or other access point to any supporting documentation.*

749 – No.

750 **Q28: Does the dataset relate to people?** *If not, you may skip the remaining questions in this
751 section.*

752 – No.

753 **Q29: Did you collect the data from the individuals in question directly, or obtain it via third
754 parties or other sources (e.g., websites)?**

755 – N/A.

756 **Q30: Were the individuals in question notified about the data collection?** *If so, please describe
757 (or show with screenshots or other information) how notice was provided, and provide a link
758 or other access point to, or otherwise reproduce, the exact language of the notification itself.*

759 – N/A.

760 Q31: **Did the individuals in question consent to the collection and use of their data?** *If so,*
761 *please describe (or show with screenshots or other information) how consent was requested*
762 *and provided, and provide a link or other access point to, or otherwise reproduce, the exact*
763 *language to which the individuals consented.*

764 – N/A.

765 Q32: **If consent was obtained, were the consenting individuals provided with a mechanism to**
766 **revoke their consent in the future or for certain uses?** *If so, please provide a description,*
767 *as well as a link or other access point to the mechanism (if appropriate).*

768 – N/A.

769 Q33: **Has an analysis of the potential impact of the dataset and its use on data subjects (e.g.,**
770 **a data protection impact analysis) been conducted?** *If so, please provide a description*
771 *of this analysis, including the outcomes, as well as a link or other access point to any*
772 *supporting documentation.*

773 – N/A.

774 Q34: **Any other comments?**

775 – No.

776 E.4 Preprocessing, cleaning, and/or labeling

777 Q35: **Was any preprocessing/cleaning/labeling of the data done (e.g., discretization or bucket-**
778 **ing, tokenization, part-of-speech tagging, SIFT feature extraction, removal of instances,**
779 **processing of missing values)?** *If so, please provide a description. If not, you may skip the*
780 *remainder of the questions in this section.*

781 – Yes, this is described in detail in section 4 and A.

782 Q36: **Was the “raw” data saved in addition to the preprocessed/cleaned/labeled data (e.g., to**
783 **support unanticipated future uses)?** *If so, please provide a link or other access point to*
784 *the “raw” data.*

785 – The raw structured data is stored in the ORD GitHub repository.

786 Q37: **Is the software used to preprocess/clean/label the instances available?** *If so, please*
787 *provide a link or other access point.*

788 – This paper is for the software used to preprocess, clean, and label the instances.

789 Q38: **Any other comments?**

790 – No.

791 E.5 Uses

792 Q39: **Has the dataset been used for any tasks already?** *If so, please provide a description.*

793 – Yes, in section 5 we train a previously published neural network model for reaction
794 condition prediction and a previously published transformer for forward prediction and
795 retrosynthesis.

796 Q40: **Is there a repository that links to any or all papers or systems that use the dataset?** *If*
797 *so, please provide a link or other access point.*

798 – No.

799 Q41: **What (other) tasks could the dataset be used for?**

800 – As described in section 2, other key problems in chemical synthesis include reaction
801 outcome prediction, retrosynthesis, and reaction condition prediction. An important
802 task which was not described is reaction yield prediction. Successful reaction yield
803 models are predominantly trained on high-throughput experimentation (HTE) datasets
804 [15], and is known to be difficult (if not impossible) with patent data (e.g. USPTO)

805 [13, 14]. As long as ORD primarily consists of USPTO data, ORDERly will probably
806 not be very useful for yield prediction, but it could be in the future.

807 **Q42: Is there anything about the composition of the dataset or the way it was collected**
808 **and preprocessed/cleaned/labeled that might impact future uses?** *For example, is there*
809 *anything that a future user might need to know to avoid uses that could result in unfair*
810 *treatment of individuals or groups (e.g., stereotyping, quality of service issues) or other*
811 *undesirable harms (e.g., financial harms, legal risks) If so, please provide a description. Is*
812 *there anything a future user could do to mitigate these undesirable harms?*

813 – Yes, ORDERly relies on the ORD schema, and changes to the ORD schema or ORD
814 database may require updates to ORDERly. ORD may change in the future, as the it
815 becomes more clear how the community wishes to use ORD (e.g. which classes of
816 information are stored).

817 **Q43: Are there tasks for which the dataset should not be used?** *If so, please provide a*
818 *description.*

819 – The ORDERly datasets were generated to make it easier to train models that can predict
820 how to make small molecules. The intended usage is to predict synthesis pathways
821 for therapeutics, however, within this category of small molecules is also energetic
822 materials, such as explosives.

823 **Q44: Any other comments?**

824 – No.

825 **E.6 Distribution**

826 **Q45: Will the dataset be distributed to third parties outside of the entity (e.g., company,**
827 **institution, organization) on behalf of which the dataset was created?** *If so, please*
828 *provide a description.*

829 – Yes, the datasets will be open-source.

830 **Q46: How will the dataset be distributed (e.g., tarball on website, API, GitHub)?** *Does the*
831 *dataset have a digital object identifier (DOI)?*

832 – The data is available through FigShare. (<https://doi.org/10.6084/m9.figshare.23298467>)

833 – It can also reliably be recreated using the instructions in the ORDERly GitHub repository
834 (<https://github.com/sustainable-processes/ORDERly>).

835 **Q47: When will the dataset be distributed?**

836 – It is already publicly available.

837 **Q48: Will the dataset be distributed under a copyright or other intellectual property (IP)**
838 **license, and/or under applicable terms of use (ToU)?** *If so, please describe this license*
839 *and/or ToU, and provide a link or other access point to, or otherwise reproduce, any relevant*
840 *licensing terms or ToU, as well as any fees associated with these restrictions.*

841 – CC-BY-4.0

842 **Q49: Have any third parties imposed IP-based or other restrictions on the data associated**
843 **with the instances?** *If so, please describe these restrictions, and provide a link or other*
844 *access point to, or otherwise reproduce, any relevant licensing terms, as well as any fees*
845 *associated with these restrictions.*

846 – No.

847 **Q50: Do any export controls or other regulatory restrictions apply to the dataset or to**
848 **individual instances?** *If so, please describe these restrictions, and provide a link or other*
849 *access point to, or otherwise reproduce, any supporting documentation.*

850 – No,

851 **Q51: Any other comments?**

852 – No.

853 **E.7 Maintenance**

854 **Q52: Who will be supporting/hosting/maintaining the dataset?**

- 855 – The dataset is hosted on FigShare, the code to generate the dataset is hosted on GitHub.
- 856 – The group of Anon will be maintaining ORDERly.

857 **Q53: How can the owner/curator/manager of the dataset be contacted (e.g., email address)?**

- 858 – Anon.

859 **Q54: Is there an erratum?** If so, please provide a link or other access point.

- 860 – N/A.

861 **Q55: Will the dataset be updated (e.g., to correct labeling errors, add new instances, delete instances)?** *If so, please describe how often, by whom, and how updates will be communicated to users (e.g., mailing list, GitHub)?*

- 864 – ORDERly will be maintained by the group of Anon, updates will be tracked through GitHub. ORDERly is built to be extensible, such that as the ORD dataset grows, users can run ORDERly to create new, larger, datasets. The ORDERly benchmark datasets are unlikely to change (to ensure model accuracy is comparable).

868 **Q56: If the dataset relates to people, are there applicable limits on the retention of the data associated with the instances (e.g., were individuals in question told that their data would be retained for a fixed period of time and then deleted)?** *If so, please describe these limits and explain how they will be enforced.*

- 872 – N/A.

873 **Q57: Will older versions of the dataset continue to be supported/hosted/maintained?** If so, please describe how. If not, please describe how its obsolescence will be communicated to users.

- 876 – The datasets are small enough to easily be versioned and hosted on FigShare (350k-1m reactions, 200MB-500MB).

878 **Q58: If others want to extend/augment/build on/contribute to the dataset, is there a mechanism for them to do so?** *If so, please provide a description. Will these contributions be validated/verified? If so, please describe how. If not, why not? Is there a process for communicating/distributing these contributions to other users? If so, please provide a description*

- 883 – All contributions to ORDERly will be managed through the ORDERly GitHub repository. Pull requests into main will need to be verified by a member of Anon’s group.

885 **Q59: Any other comments?**

- 886 – No.