

001
002
003
004
005
006
007 AutoTemplate: Enhancing Chemical Reaction Datasets for
008
009
010 Machine Learning Applications in Organic Chemistry
011

012
013 Lung-Yi Chen¹ and Yi-Pei Li^{1,2*}

014
015 ¹Department of Chemical Engineering, National Taiwan University, No. 1, Sec. 4,
016 Roosevelt Road, Taipei, 10617, Taiwan.

017 ²Taiwan International Graduate Program on Sustainable Chemical Science and
018 Technology (TIGP-SCST), No. 128, Sec. 2, Academia Road, Taipei, 11529, Taiwan.
019

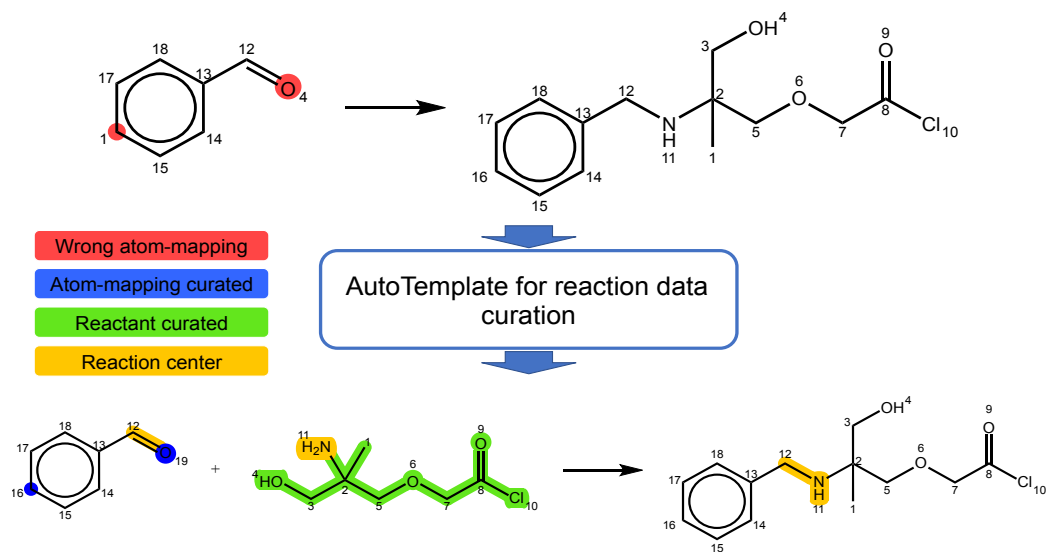
020
021 *Corresponding author: Yi-Pei Li. E-mail: yipeili@ntu.edu.tw;
022

023
024 **Abstract**

025 This paper presents AutoTemplate, an innovative data preprocessing protocol, addressing the
026 crucial need for high-quality chemical reaction datasets in the realm of machine learning
027 applications in organic chemistry. Recent advances in artificial intelligence have expanded the
028 application of machine learning in chemistry, particularly in yield prediction, retrosynthesis, and
029 reaction condition prediction. However, the effectiveness of these models hinges on the integrity
030 of chemical reaction datasets, which are often plagued by inconsistencies like missing reactants,
031 incorrect atom mappings, and outright erroneous reactions. AutoTemplate introduces a two-
032 stage approach to refine these datasets. The first stage involves extracting meaningful reaction
033 transformation rules and formulating generic reaction templates using a simplified SMARTS
034 representation. This simplification broadens the applicability of templates across various chem-
035 ical reactions. The second stage is template-guided reaction verification, where these templates
036 are systematically applied to validate and correct the reaction data. This process effectively
037 amends missing reactant information, rectifies atom-mapping errors, and eliminates incorrect
038 data entries. A standout feature of AutoTemplate is its capability to concurrently identify and
039 correct false chemical reactions. It operates on the premise that most reactions in datasets are
040 accurate, using these as templates to guide the correction of flawed entries. The protocol demon-
041 strates its efficacy across a range of chemical reactions, significantly enhancing dataset quality.
042 This advancement provides a more robust foundation for developing reliable machine learning
043 models in chemistry, thereby improving the accuracy of forward and retrosynthetic predictions.
044 AutoTemplate marks a significant progression in the preprocessing of chemical reaction datasets,
045 bridging a vital gap and facilitating more precise and efficient machine learning applications in
046 organic synthesis. *Scientific contribution:* The proposed automated preprocessing tool for chem-
047 ical reaction data aims to identify errors within chemical databases. Specifically, if the errors
048 involve atom mapping or the absence of reactant types, corrections can be systematically applied
049 using reaction templates, ultimately elevating the overall quality of the database.

049 **Keywords:** Reaction template, Data preprocessing, Atom-to-atom mapping, Reaction data curation
050
051
052
053
054
055

Graphical Abstract



1 Introduction

079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104

Recent advancements in artificial intelligence have greatly expanded its applications in the field of chemistry. Machine learning techniques have been integrated into various aspects of organic synthesis, including yield prediction [1–4], forward prediction [5–9], retrosynthesis [10–15] and reaction condition prediction [16–19]. These predictive models rely on extensive and reliable chemical reaction datasets, enabling the development of robust machine learning solutions for real-world scenarios [20–24].

Chemical reaction databases commonly utilized in the literature can be broadly categorized as open-source datasets such as the United States Patent and Trademark Office (USPTO) [25] and open reaction database (ORD) [26], or proprietary datasets like Pistachio [27], Reaxys [28], SciFinder [29], and Spresi [30]. These datasets are compiled through text-mining or manual recording, both of which can introduce errors in the chemical reaction data. Fig. 1 illustrates common data deficiencies observed in chemical databases, including missing reactants, inexplicable extra atoms in products, and even entirely erroneous reactions. Detecting and rectifying these data inconsistencies often require human intervention to ensure the quality of machine learning models.

To address these issues, Gimadiev et al. [31] employed atom-to-atom mapping toolkits [32–35] and the CGRTools [36] python library for preprocessing chemical transformations. They used a condensed graph of reaction (CGR), representing the superposition of the reactants and products, to remove duplicate reactions and balance reaction equations, particularly in cases where

105
106
107
108
109
110

111 simple reagents like amine and water were unspecified. In contrast, Vaucher et al. [37] developed a
112 transformer-based model [38] to complete reaction equations by filling in missing parts of molecules
113 in partial reactions using a sequence-to-sequence approach. Although the model exhibited versatil-
114 ity in handling retrosynthesis, forward prediction, and data curation tasks, it achieved an accuracy
115 of approximately 30% for exact matches, which may pose limitations in its application for extensive
116 preprocessing of external chemical reaction datasets. More recently, Toniato et al. [39] employed the
117 concept of catastrophic forgetting [40] to monitor the learning progress of molecular transformer
118 [9] during training. Data points with difficulty in learning were assumed to be associated with
119 errors and were subsequently removed from the dataset. However, the extent of data removal using
120 this approach significantly depended on the model used, its learning capacity, and hyperparameter
121 selection, rendering it less deterministic.

129 To the best of our knowledge, existing data-preprocessing methods have limited capacity to
130 detect and correct false chemical reactions simultaneously. This gap has motivated us to develop an
131 advanced data-preprocessing protocol called AutoTemplate in this work. AutoTemplate establishes
132 clear criteria for identifying and removing erroneous data while effectively recovering missing reac-
133 tants. It operates under the assumption that the majority of reactions in datasets are correct and
134 uses these reactions as templates to guide the curation of incorrect data. The proposed method can
135 successfully identify incorrect reactions, correct faulty atom mapping, and complete missing reac-
136 tants, providing a solid foundation for the development of data-driven machine learning models,
137 thereby enhancing the performance of forward and retrosynthetic predictions.

145 2 Method

148 The data cleaning methodology presented in this work is divided into two stages: generic template
149 extraction and template-guided reaction verification. In the generic template extraction stage, we
150 first identify meaningful reaction transformation rules within the dataset of interest. These rules are
151 then expressed as generic reaction templates using a simplified version of the SMARTS representa-
152 tion [43]. This simplification ensures that the templates can be applied to a wide range of reactions
153 with the same transformation. In the template-guided reaction verification stage, we leverage the
154 list of generic reaction templates to systematically validate the reaction data. This involves applying
155 retro templates to the product. If the original reactants are indeed a subset of the results obtained
156 through template application, the template-applied outcomes replace the original data. This process
157 effectively rectifies any missing reactant information and simultaneously corrects potential atom-
158 mapping errors. However, in situations where none of the templates match the reaction, indicating

166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220

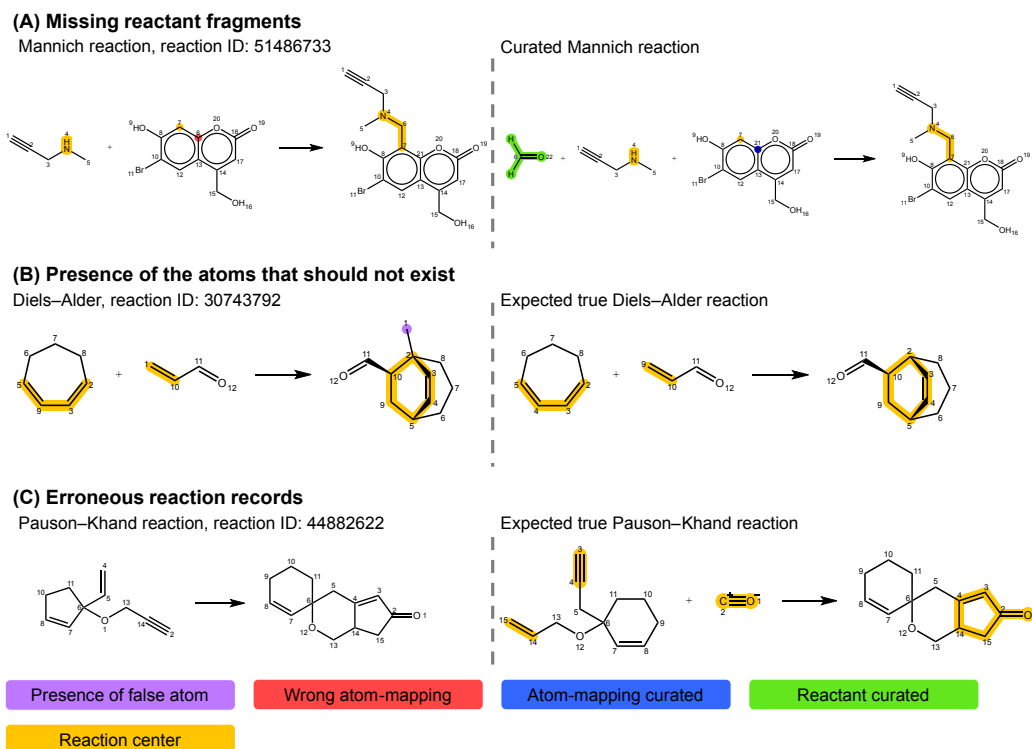


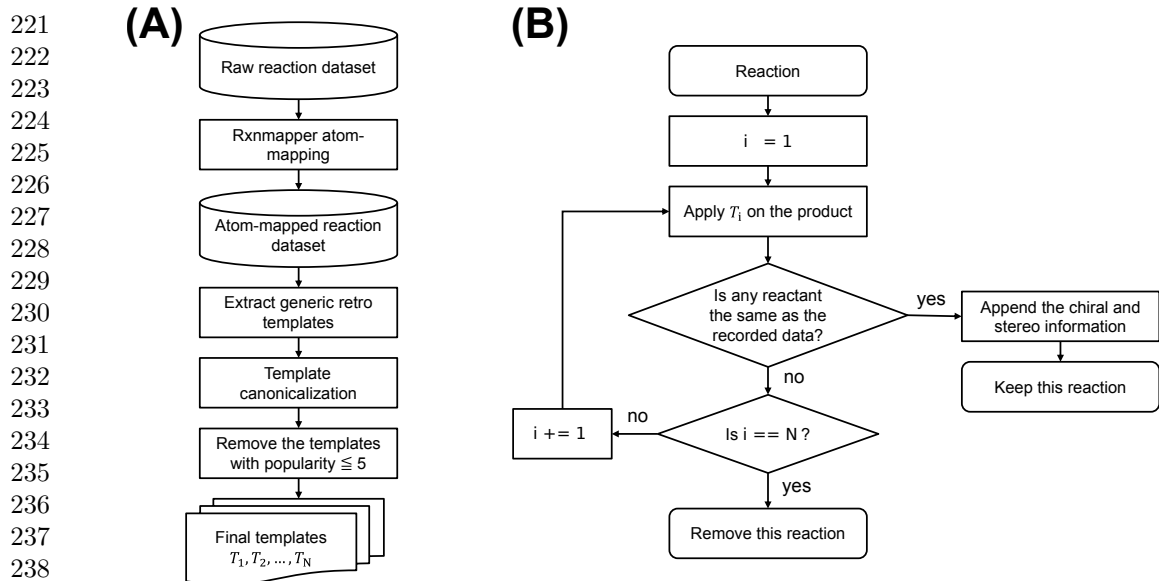
Fig. 1 Illustrations of deficiencies in reaction datasets: (A) The selected Mannich reaction omits formaldehyde in the reactants. (B) The presence of the carbon atom (labeled as purple) violates the law of conservation of matter, and the accurate product based on the work by De Nino et al. [41] is shown on the right. (C) The reactant and product do not match, and the correct chemical reaction is depicted on the right side, extracted from the study by Özdemirhan [42]. These examples are sourced from the Reaxys database [28], but it is important to note that similar errors exist in other databases. Notably, the original Reaxys dataset lacks atom-mapping information, and the atom-mapping labels in the left half of this figure were generated using the RXNMapper software [35].

an unusual chemical transformation and potentially incorrect data entry, we opt to remove that specific reaction from our dataset. The overall procedure is visually depicted in Fig. 2, with detailed step-by-step explanations provided in the following subsections.

2.1 Generic template extraction

2.1.1 Reaction data collection

To evaluate the effectiveness of our data cleaning protocol, we applied it to reaction data derived from the Reaxys database [28], a well-established resource in the field of computational chemistry that, like any large database, may contain some errors [31]. To demonstrate the broad applicability of our data preprocessing approach, we retrieved datasets for 20 different reaction types from Reaxys. These datasets were obtained by searching for specific reaction names, and they encompassed a variety of reactions, including Adams decarboxylation, Baylis–Hillman reaction, Buchwald–Hartwig cross coupling, Chan–Lam coupling, Diels–Alder, Fischer indole synthesis, Friedel–Crafts acylation, Friedel–Crafts alkylation, Grignard reaction, Hiyama coupling, Huisgen cycloaddition, Hydrogenation, Kabachnik–Fields reaction, Kumada coupling, Mannich reaction,



239 **Fig. 2** Overview of the two-stage data cleaning protocol of AutoTemplate for processing chemical reaction data. Panel
240 (A) illustrates the generic template extraction procedure. Panel (B) shows the template-guided reaction verification
241 process, which systematically validates the reaction data using a list of generic reaction templates.

242
243 Negishi coupling, Pauson–Khand reaction, reductive amination, Suzuki coupling, and Wittig reac-
244 tion. The reaction IDs for each reaction used in our study are provided in the GitHub repository
245 for reference [44]. We removed any reactions involving reactants or products that could not be
246 parsed by RDKit [45]. In addition, we eliminated isotope labels from the molecules since they do not
247 impact the chemical transformation. It is worth noting that the labels denoting reaction types in the
248 Reaxys database may not always align accurately with the actual reaction types. Therefore, despite
249 our efforts to collect data based on the 20 specified reaction names, there were instances where the
250 recorded reaction entries did not correspond precisely to these 20 designated reaction types.

257 2.1.2 Atom-to-atom mapping

258
259 The original reaction data obtained from Reaxys lacked information on atom mapping, a crucial ele-
260 ment for establishing correspondence between the atoms of reactants and products. This information
261 is essential to identify the reaction center where the connectivity of atoms has changed, a prerequisite
262 for extracting the reaction template. The accuracy of common atom-to-atom mapping toolkits has
263 been assessed in the study by Lin et al. [33]. According to their findings, the open-source tool RXN-
264 Mapper [35] demonstrated state-of-the-art performance, processing each reaction within one second.
265 Due to these advantages, we selected RXNMapper as our preprocessing toolkit for atom-to-atom
266 mapping. With atom-mapping information available, we can distinguish spectator molecules—those
267
268
269
270
271
272
273
274
275

276 that do not actively participate in the reaction or contribute any non-hydrogen atoms to the prod-
277 uct. These spectator molecules were removed because our data preprocessing framework focuses on
278 curating the chemical transformation itself, rather than the spectator molecules.
279
280

281

282 **2.1.3 Generic template definition and extraction**

283

284 Upon obtaining the atom-mapped reactions, the next step is to retrieve all the reaction templates
285 from the dataset using the RDChiral [46] template extractor. It is important to note that RDChi-
286 ral primarily focuses on generating retrosynthetic templates, which are designed for developing
287 computer-aided retrosynthesis models. Because chemical reaction datasets often focus on the major
288 product while not necessarily comprehensively documenting the reactants needed to produce that
289 product, our study utilizes retrosynthetic templates to verify and curate the reaction data.
290
291
292
293

294 The default templates generated by RDChiral provide highly detailed information around the
295 reaction center. This results in an excessive number of templates for the same type of chemical
296 transformation, particularly when there are minor variations in neighboring functional groups. It also
297 extends the time required for the subsequent template application process. The specificity of these
298 templates can make it challenging to apply a template from one reaction entry to curate another
299 entry, unless both entries have identical neighboring functional groups near the reaction center. To
300 overcome these challenges, we made modifications to the RDChiral functions. Our aim was to create
301 generic reaction templates that include only essential information concerning atom types and bond
302 types within the reaction centers, while excluding extraneous details. Table 1 provides a comparison
303 between the default and modified template extraction functions.
304
305

311 Consider the Grignard reaction in Fig. 3A as an example, the corresponding reaction template
312 generated by default RDChiral is `[OH;D1;+0:4]-[CH;D3;+0:5](-[c:6])-[c;H0;D3;+0:1](:[c:2]):[c:3]>>Br-[c;H0;D3;+0:1](:[c:2]):[c:3].[0;H0;D1;+0:4]=[CH;D2;+0:5]-[c:6]`. On
313 the other hand, its generic template reduces to `[#6:1]-[#6:2]-[#8:3]>>Br-[#6:1].[#6:2]=[`
314 `#8:3]`. In the generic template, details related to atomic aromaticity, degree of freedom, number of
315 hydrogen atoms, charge, and extra atoms are all discarded. The meanings of the notations used in
316 the template can be found in the reaction SMARTS documentation [47]. This simplification effec-
317 tively documents the chemical transformation for most cases. Nevertheless, there are special cases
318 that require unique treatment. The first exception involves specifying the number of connected
319 hydrogens in the generic template to accurately represent species involved in radical reactions, as
320 shown in Fig. 3B. The second exception is the inclusion of the number of charges in the template
321 when the reaction involves charge transfer, as illustrated in Fig. 3C. The third exceptional case
322
323
324
325
326
327
328
329
330

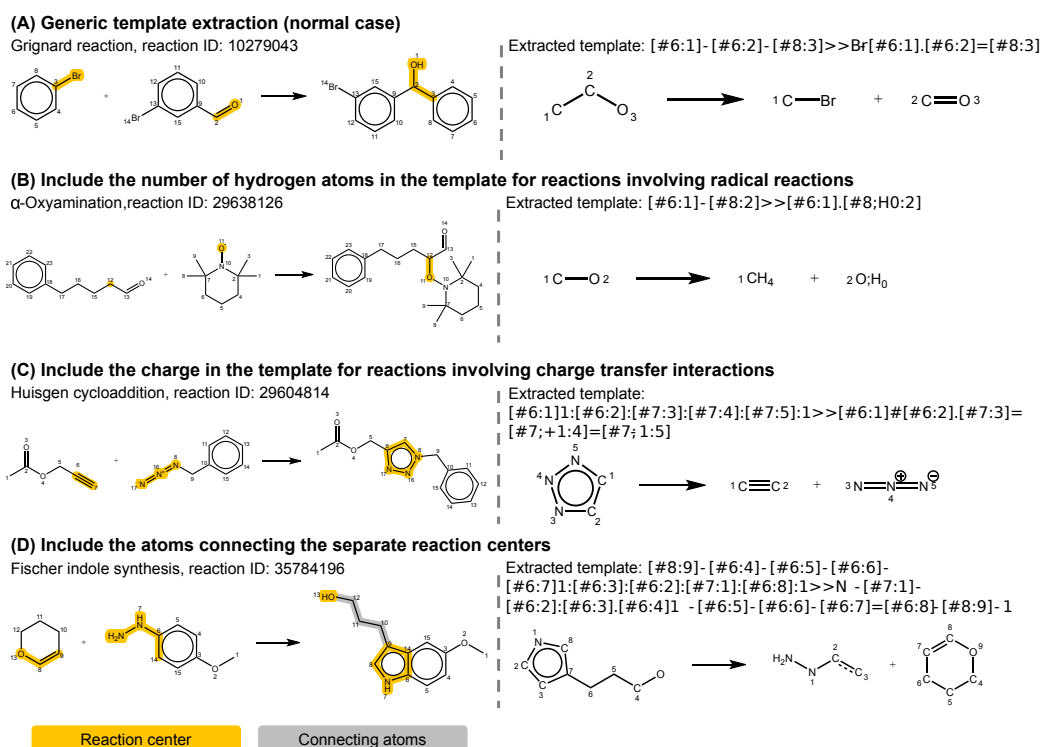
331 **Table 1** The features specified in default RDChiral and generic reaction templates.

Level	Features	RDChiral	Generic
334	Reactant radius ¹	1	0
335	Product radius ¹	0	0
336	Aliphatic or aromatic	Yes	No
337	Degree of freedom ²	Yes	No
338	Chirality	Yes	No
339	No. of hydrogen atoms	Yes	No, except for radical reactions
340	Charge	Yes	No, except for charge transfer reactions
341	Bond type	Yes	Yes
342	Cis-trans isomerism	Yes	No
343	Functional groups	Yes	Yes
344	Predefined groups	Yes	No

345 ¹Radius denotes the extending distance of the neighbor atoms around the reaction center.

346 ²Degree of freedom here represents the number of connecting non-hydrogen atoms.

347
348
349 arises when separate reaction centers occur in the product (Fig. 3D). In such cases, the connect-
350 ing atoms between the reaction centers should be incorporated into the generic template. These
351 connecting atoms can be identified using Dijkstra's algorithm [48], which finds the shortest path
352 between given nodes. This approach ensures that no redundant atoms are included in the template
353 and is effectively applicable to extracting templates for ring-opening reactions.
354
355
356



381 **Fig. 3** Illustration of generic template extraction with the normal and special cases.

382

383

384

385

386 2.1.4 Template canonicalization

387

388 To address the issue of having multiple generic reaction templates representing the same chemical
389 transformation but with different text representations, we employed a graph isomorphism check to
390 confirm whether the reactants and products in pairwise templates were identical. If both reactant and
391 product SMARTS patterns were graph isomorphic, we combined the two templates. Additionally, we
392 calculated the number of bond changes in the templates and keep the one with fewer changes. Fig.
393 4 illustrates this scenario with two Diels–Alder reaction templates that share identical subgraphs of
394 reactants and products but differ in reaction transformations due to mapping errors from the atom-
395 mapping tool. Such errors can lead to incorrect atom swaps, resulting in additional and incorrect
400 formation and breaking of chemical bonds. Therefore, we retained the template with fewer bond
401 changes.
402
403 changes.

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

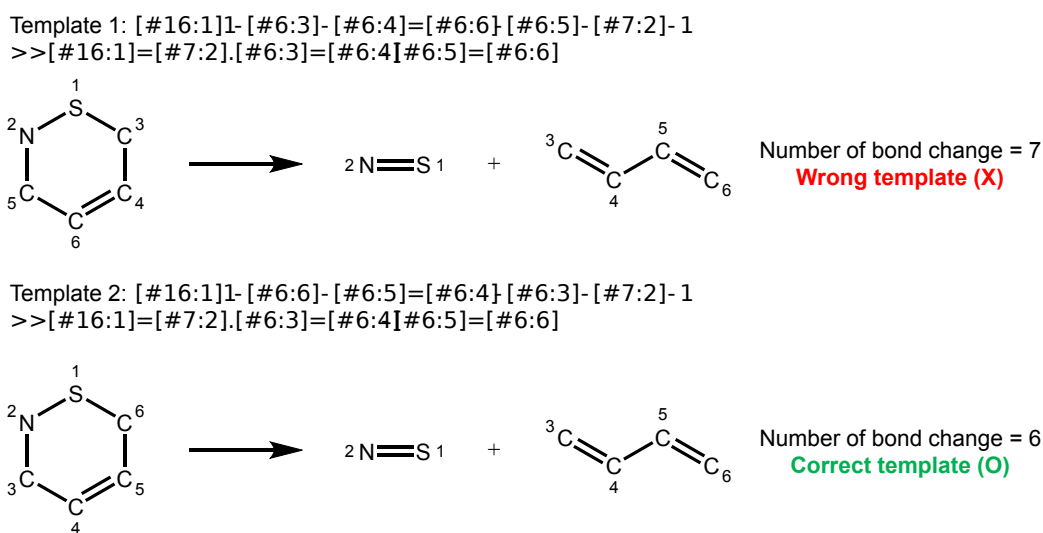
436

437

438

439

440



422 **Fig. 4** Examples of two generic templates extracted from Diels–Alder reactions.

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

426 2.1.5 Removal of rare templates

427

428

429

430

431

432

433

434

435

436

437

438

439

440

Generic templates are designed to be broadly applicable to reaction instances with similar chemical transformations. If a generic template matches only a few reaction entries, it suggests an unusual chemical transformation, possibly indicating that the template may have been derived from a reaction entry with errors. To address this, we monitored the occurrence frequency of each generic template during the template extraction process. Templates with a popularity of 5 or less were removed. This process resulted in the final set of generic templates $\{T_1, T_2, \dots, T_N\}$ for subsequent template-guided verification.

441 2.2 Template-guided verification

442

443 2.2.1 Template application procedure

444

445 This procedure primarily involves the iterative application of generic reaction templates to the
446 products of each reaction entry. When the reactants in the original data entry form a subset of
447 the reactants resulting from the applied template, we replace the original data's reactants with
448 those from the applied template. This rectifies any missing reactant information and simultaneously
449 corrects potential atom-mapping errors. In cases where none of the templates match the reaction,
450 indicating an unusual chemical transformation and potentially incorrect data entry, we choose to
451 remove that specific reaction entry from the dataset.

452 Throughout the template application process, the reactants are automatically supplemented with
453 the appropriate number of hydrogen atoms based on their charge state and the number of bonds
454 connected to them. For instance, neutral sulfur atoms are assigned either two or six bonds, resulting
455 in two possible configurations for a neutral sulfur atom with a connected chemical bond, acquiring
456 either one or five hydrogen atoms. Exceptions to this rule only occur when the template explicitly
457 specifies the number of hydrogen atoms connected to the reaction center.

468 2.2.2 Append atomic chirality and bond stereochemistry

469

470 We note that the reactants generated from template application lack annotations for atomic chirality
471 and bond stereochemistry at the reaction centers. Therefore, an additional step is necessary to rein-
472 troduce this information into the reactants, but only if this information was included in the original
473 dataset. This process involves establishing a one-to-one atom correspondence between the original
474 reactants and template-generated reactants. This can be achieved by initially converting both sets
475 of reactants into undirected graphs, followed by utilizing the exact graph matching algorithm [49]
476 to establish a strict one-to-one node correspondence between the two graphs.

480

483 3 Results and Discussion

484

486 3.1 Analysis of overall results

487

488 Table 2 provides information on the number of reactions in the dataset, the number of templates
489 extracted from these reactions, and the residual proportion after data processing. The variation
490 in the number of templates for each type of reaction is due to the unique characteristics of their
491 reaction mechanisms. For example, coupling reactions that involve multiple possible leaving groups
492 often result in a higher template count. Conversely, reductive amination, where the carbonyl group

496 **Table 2** The data preprocessing results for the chemical reaction datasets.

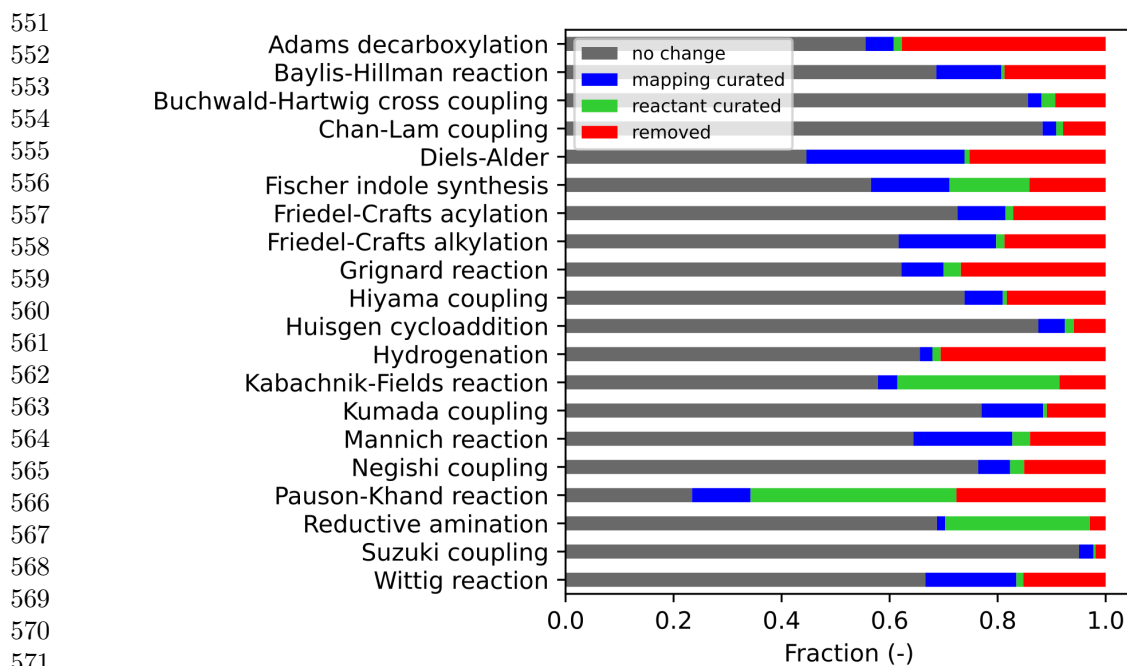
497	Reaction type	No. of reactions	No. of generic templates	Residual rate
498				
499	Adams decarboxylation	2,636	54	62.3%
500	Baylis–Hillman reaction	7,507	84	81.3%
501	Buchwald–Hartwig cross coupling	18,341	96	90.7%
502	Chan–Lam coupling	6,885	43	92.1%
503	Diels–Alder	18,757	258	74.8%
504	Fischer indole synthesis	6,841	28	85.9%
505	Friedel–Crafts acylation	10,095	118	82.9%
506	Friedel–Crafts alkylation	17,248	164	81.3%
507	Grignard reaction	13,530	154	73.2%
508	Hiyama coupling	4,089	106	81.7%
509	Huisgen cycloaddition	54,183	144	94.1%
510	Hydrogenation	41,217	306	69.4%
511	Kabachnik–Fields reaction	5,575	14	91.4%
512	Kumada coupling	16,371	82	89.1%
513	Mannich reaction	29,698	271	86.0%
514	Negishi coupling	10,909	146	84.9%
515	Pauson–Khand reaction	2,703	19	72.4%
516	Reductive amination	50,406	16	97.1%
517	Suzuki coupling	184,219	216	98.2%
518	Wittig reaction	16,337	94	84.8%

516 is reduced to an amine, has a large number of reaction entries, but only 16 reaction templates are
517 extracted, indicating less variation in its reaction transformation.

519 Fig. 5 displays curated reaction results, addressing issues such as false atom-mapping, reactant
520 omissions, and the identification and removal of incorrect reaction records. Notably, the Diels–Alder
521 reactions exhibited a high atom-mapping correction rate of 29.3%. This is likely attributed to the
522 complexity of Diels–Alder reactions, which involve numerous bond transformations and instances of
523 intramolecular or fused ring formation, making them challenging for accurate atom-mapping predic-
524 tions. Conversely, coupling reactions generally showed relatively fewer atom-mapping errors, likely
525 because they involve fewer bond changes. Accurate atom-mapping data can significantly improve
526 reaction prediction quality, particularly for graph-based models. Regarding the issue of missing reac-
527 tants, Fischer indole synthesis, Kabachnik–Fields reaction, Pauson–Khand reaction, and reductive
528 amination display a noteworthy proportion of data with absent reactants. In the case of the Pau-
529 son–Khand reaction, most instances systematically omit carbon monoxide as a reactant. However,
530 there is no clear pattern indicating which reactants may be omitted in the data for the other three
531 types of reactions. Further discussions on specific data errors and curated results are provided in
532 the following subsections for selected examples.

544 3.2 Visualized results of selected mapping curated examples

546 Currently, there is no package available that can generate atom-mapping information perfectly for
547 all reactions [33]. In this study, the data-driven neural network RXNMapper [35] was utilized to
548 predict atom mapping. However, it is important to note that even for reactions considered relatively
549
550



572 **Fig. 5** Distribution of the proportion of repaired reactions after data processing.

575 straightforward for humans, there can still be instances of incorrect atom mapping, as shown in
576 Fig. 6A. This example of the Baylis-Hillman reaction incorrectly assigns the atom-mapping number
577 (6 and 14) at the position of the carbon-carbon double bond, which would lead to the incorrect
578 reaction template during template extraction. Applying the data processing procedure proposed
579 in this work can recover this reaction with the true atom-mapping labels. Another example is the
580 Buchwald-Hartwig cross-coupling reaction illustrated in Fig. 6B, which has the same issue at the
581 reaction center where the carbon atoms are labeled incorrectly in the intramolecular ring-closing
582 reaction. We note that false atom-mapping issues occur more frequently at the reaction centers, and
583 systematically addressing this problem would benefit downstream template-based and graph-based
584 model applications.

592 3.3 Visualized results of selected reactant curated examples

595 The data processing procedure proposed in this work primarily focuses on addressing omitted reac-
596 tants rather than products, as byproducts and leaving groups are typically not the main focus and
597 are not specified in reaction datasets. The issue of missing reactants can be identified by comparing
598 the atom counts between reactants and products, with reactions having fewer atoms on the reactant
599 side categorized as this type of error. To the best of our knowledge, there is no existing approach
600 tailored for adding missing reactants. However, with the template-guided verification method pro-
601 posed in this work, erroneous reaction entries can be recovered along with the omitted reactants.

606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660

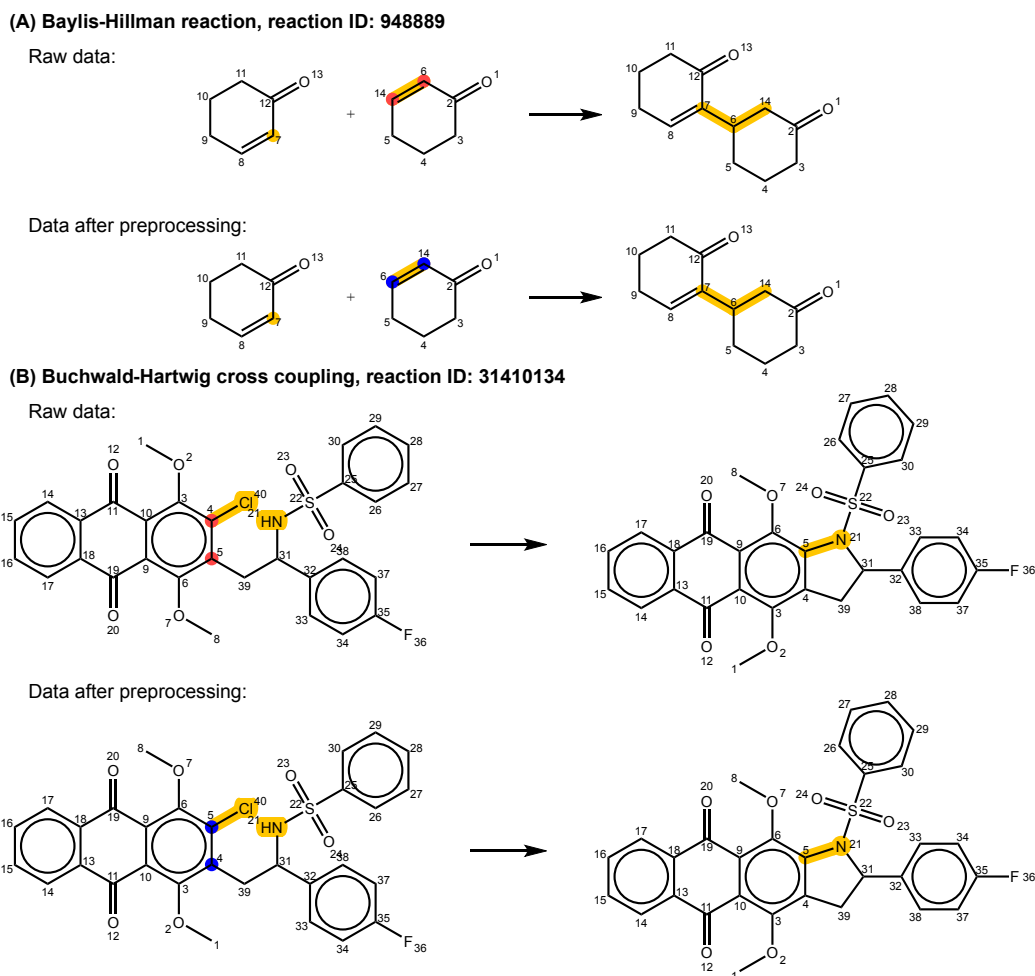


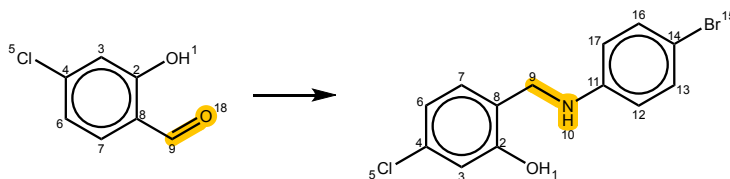
Fig. 6 Two selected examples of (A) Baylis–Hillman reaction and (B) Buchwald–Hartwig cross coupling to demonstrate the curated results of the reaction entries with incorrect atom-mapping. Yellow highlights indicate the reaction centers, red highlights denote atoms with incorrect atom mapping, and blue highlights represent atoms with curated mapping.

Fig. 7A illustrates a typical example from the reductive amination dataset, where the missing reactant with an amine functional group was generated by applying the generic template to the product, thus balancing the reaction equation. In the case of the second instance of the Kabachnik–Fields reaction shown in Fig. 7B, which involves three molecules in the reaction, the two missing fragments were successfully recovered from the template. It is worth noting that the chirality of the phosphorus atom cannot be inferred because the generic template does not specify chiral and cis-trans stereoisomerism at the reaction center. Including such detailed information in templates would lead to an excessive number of templates, reducing the chances of applying a template from one reaction entry to curate another entry.

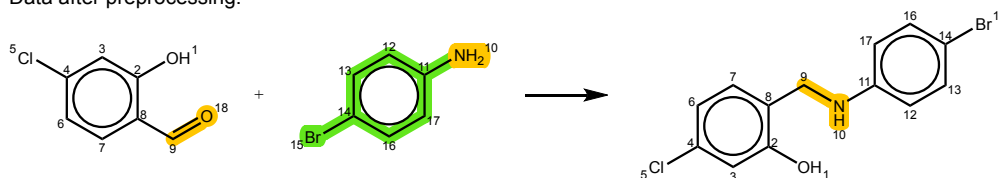
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715

(A) Reductive amination, reaction ID: 32187214

Raw data:

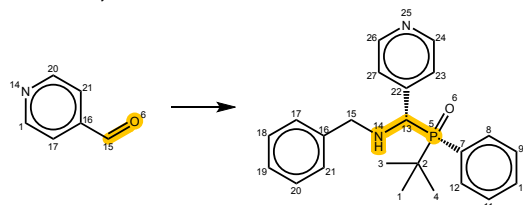


Data after preprocessing:



(B) Kabachnik–Fields reaction, reaction ID: 12318568

Raw data:



Data after preprocessing:

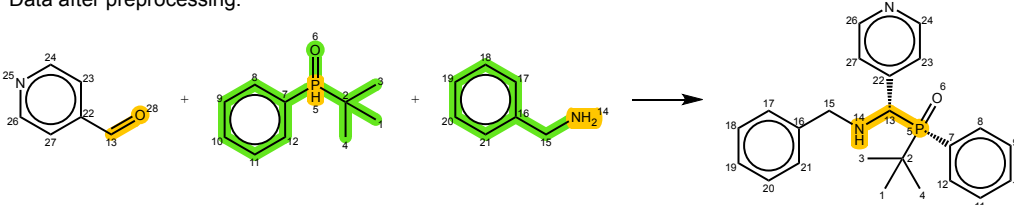


Fig. 7 Two selected examples of (A) reductive amination and (B) Kabachnik–Fields reaction to demonstrate the curated results of the reaction entries with incomplete reactant information. Yellow highlights represent the reaction centers, while green highlights indicate molecular fragments added through the data curation process.

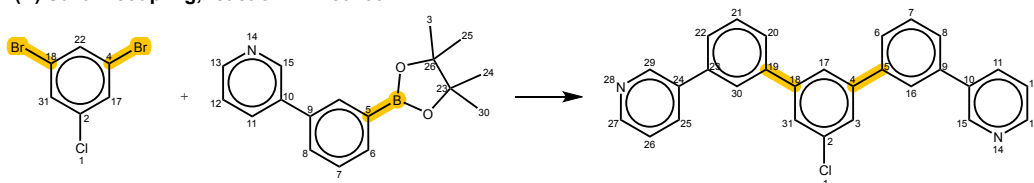
3.4 Visualized results of selected removed reactions

In cases where none of the templates matched the reaction, indicating an unusual chemical transformation or potential data entry errors, the specific reaction entry was removed from the dataset. Several examples of such removals are presented in Fig. 8 and discussed below.

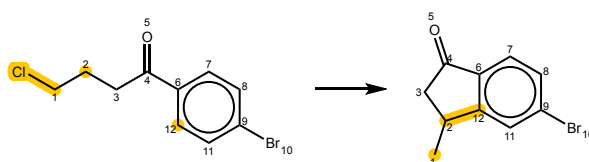
Fig. 8A illustrates a two-step Suzuki coupling reaction. To automatically identify multi-step reactions like this, one would need to repetitively validate them using all the single-step reaction templates, which becomes increasingly time-consuming as the number of steps allowed grows. Because most reaction prediction models focus on single-step reactions, the accommodation of multi-step reactions is less critical in this study. The reactions shown in Fig. 8B and 8C are actually correct reactions, but none of the generic templates in the final list match them. This occurred because the templates extracted from these reactions did not match a sufficient number of reaction entries, leading to their exclusion from the final list of generic templates. As discussed in the method section, templates with low matching frequencies may indicate errors in the template source data. While

716 this approach effectively removes erroneous reaction entries, it can also inadvertently exclude rare
 717 but valid reactions, as demonstrated in Fig. 8B and 8C. The reaction depicted in Fig. 8D belongs
 718 to the category of Huisgen cycloaddition. In this reaction, the atom highlighted in purple (number
 719 10) in the product is identified as a carbon atom. However, at the same position in the reactant,
 720 an oxygen atom is indicated. Rectifying this type of error is challenging because it is difficult to
 721 determine whether the correct structure should be attributed to the reactant or the product. This
 722 particular entry originates from a study by McNitt et al. [50], where atom number 10 was labeled as
 723 an oxygen atom, suggesting a potential error in the recorded product information in the database.
 724
 725
 726
 727
 728

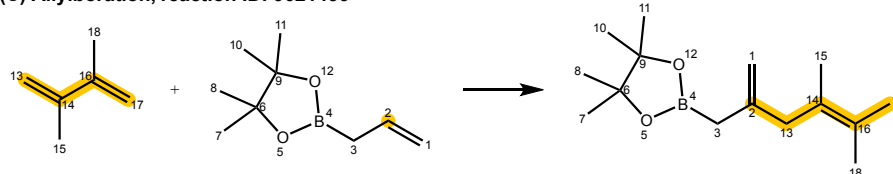
729
 730 **(A) Suzuki coupling, reaction ID: 28976014**



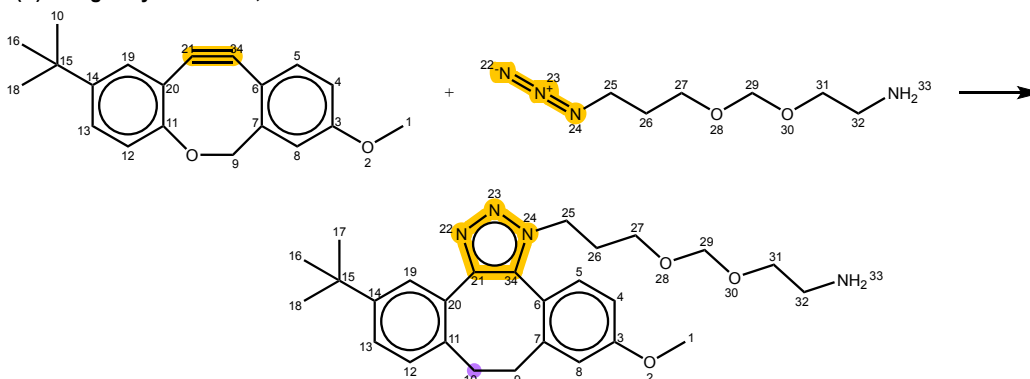
737 **(B) Friedel–Crafts alkylation, reaction ID: 3808146**



743 **(C) Allylboration, reaction ID: 9021436**



750 **(D) Huisgen cycloaddition, reaction ID: 34297252**



763 **Fig. 8** Four selected examples of (A) Suzuki coupling, (B) Friedel–Crafts alkylation, (C) allylboration, and (D) Huisgen cycloaddition to illustrate reactions that did not match any of the final generic templates and were consequently removed during the data processing procedure.
 764
 765

766
 767
 768
 769
 770

771 4 Conclusions

772
773 Recent advancements in artificial intelligence have significantly impacted the field of organic chem-
774 istry. The reliability of predictive models in chemistry, essential for applications such as yield
775 prediction, retrosynthesis, and reaction condition prediction, is heavily contingent on the quality
776 of chemical reaction datasets. However, these datasets, sourced from both open-source and propri-
777 etary databases, often contain inconsistencies like missing reactants, incorrect atom mappings, or
778 erroneous reactions, necessitating rigorous data preprocessing.
779

783 This work introduces a novel data preprocessing protocol called AutoTemplate, designed to
784 enhance the quality of chemical reaction datasets. AutoTemplate employs a two-stage approach:
785 generic template extraction and template-guided reaction verification. The process begins with the
786 extraction of meaningful reaction transformation rules from a dataset, which are then expressed as
787 generic reaction templates using a simplified version of the SMARTS representation. This simplifi-
788 cation ensures broad applicability across various reactions. In the subsequent stage, these generic
789 templates are systematically applied to validate and correct reaction data. This involves rectifying
790 missing reactant information, correcting atom-mapping errors, and removing incorrect data entries.
791

796 Our method stands out by its ability to simultaneously identify and correct false chemical reac-
797 tions, leveraging the assumption that the majority of reactions in datasets are correct. By using
798 these reactions as templates for data curation, AutoTemplate not only rectifies existing errors but
799 also aids in the recovery of missing reactants. The protocol's effectiveness is demonstrated through
800 its application to diverse chemical reactions, highlighting significant improvements in dataset qual-
801 ity. This refined data provides a more reliable foundation for developing machine learning models
802 in chemistry, enhancing the accuracy of forward and retrosynthetic predictions.
803

808 This study represents a significant step forward in preprocessing chemical reaction datasets,
809 addressing a critical gap in the field and paving the way for more accurate and efficient machine
810 learning applications in organic synthesis.
811

815 Abbreviations

817 CGR: Condensed Graph of Reaction; ORD: Open Reaction Database; SMILES: Simplified Molecular
818 Input Line Entry Specification; SMARTS: SMILES Arbitrary Target Specification; USPTO: United
819 States Patent and Trademark Office
820

826 **Declarations**

827

828

829 **Availability of data and materials**

830

831 Full code and reaction IDs for searching the reactions are available at: <https://github.com/Lung->

832

833 [Yi/AutoTemplate](https://github.com/Lung-Yi/AutoTemplate)

834

835

836 **Acknowledgements**

837

838 We are grateful to the National Center for High-performance Computing (NCHC) and the Com-

839

840 puter and Information Networking Center at NTU for the support of computing facilities. AI tools

841

842 were utilized in the process of correcting grammatical mistakes and enhancing the fluency of the

843

844 manuscript.

845

846

847 **Competing interests**

848

849 The authors declare no competing financial interest.

850

851

852 **Funding**

853

854 Y.P.L. is supported by Taiwan NSTC Young Scholar Fellowship Einstein Program (112-2636-E-002-

855

856 005) and the Higher Education Sprout Project by the Ministry of Education in Taiwan (113L891305).

857

858

859 **Authors' contributions**

860

861 LYC: Methodology, Formal Analysis, Writing - Original Draft. YPL: Funding Acquisition, Supervi-

862

863 sion, Writing - Review & Editing.

864

865

866 **References**

867

868 [1] Jiang, S., Zhang, Z., Zhao, H., Li, J., Yang, Y., Lu, B.-L., Xia, N.: When smiles smiles, practical-

869

870 ity judgment and yield prediction of chemical reaction via deep chemical language processing.

871

872 IEEE Access **9**, 85071–85083 (2021)

873

874 [2] Probst, D., Schwaller, P., Reymond, J.-L.: Reaction classification and yield prediction using

875

876 the differential reaction fingerprint drfp. Digital discovery **1**(2), 91–97 (2022)

877

878 [3] Saebi, M., Nan, B., Herr, J.E., Wahlers, J., Guo, Z., Zurański, A.M., Kogej, T., Norrby, P.-

879

880 O., Doyle, A.G., Chawla, N.V.: On the use of real-world datasets for reaction yield prediction.

881

882 Chemical Science **14**(19), 4997–5005 (2023)

- 881 [4] Schwaller, P., Vaucher, A.C., Laino, T., Reymond, J.-L.: Prediction of chemical reaction yields
882 using deep learning. *Machine learning: science and technology* **2**(1), 015016 (2021)
883
884
- 885 [5] Coley, C.W., Barzilay, R., Jaakkola, T.S., Green, W.H., Jensen, K.F.: Prediction of organic
886 reaction outcomes using machine learning. *ACS central science* **3**(5), 434–443 (2017)
887
888
- 889 [6] Coley, C.W., Jin, W., Rogers, L., Jamison, T.F., Jaakkola, T.S., Green, W.H., Barzilay, R.,
890 Jensen, K.F.: A graph-convolutional neural network model for the prediction of chemical
891 reactivity. *Chemical science* **10**(2), 370–377 (2019)
892
893
- 894 [7] Do, K., Tran, T., Venkatesh, S.: Graph transformation policy network for chemical reaction
895 prediction. In: *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge*
896 *Discovery & Data Mining*, pp. 750–760 (2019)
897
898
- 899 [8] Fooshee, D., Mood, A., Gutman, E., Tavakoli, M., Urban, G., Liu, F., Huynh, N., Van Vranken,
900 D., Baldi, P.: Deep learning for chemical reaction prediction. *Molecular Systems Design &*
901 *Engineering* **3**(3), 442–452 (2018)
902
903
- 904 [9] Schwaller, P., Laino, T., Gaudin, T., Bolgar, P., Hunter, C.A., Bekas, C., Lee, A.A.: Molecu-
905 lar transformer: a model for uncertainty-calibrated chemical reaction prediction. *ACS central*
906 *science* **5**(9), 1572–1583 (2019)
907
908
- 909 [10] Coley, C.W., Green, W.H., Jensen, K.F.: Machine learning in computer-aided synthesis
910 planning. *Accounts of chemical research* **51**(5), 1281–1289 (2018)
911
912
- 913 [11] Coley, C.W., Rogers, L., Green, W.H., Jensen, K.F.: Computer-assisted retrosynthesis based
914 on molecular similarity. *ACS central science* **3**(12), 1237–1245 (2017)
915
916
- 917 [12] Dong, J., Zhao, M., Liu, Y., Su, Y., Zeng, X.: Deep learning in retrosynthesis planning: datasets,
918 models and tools. *Briefings in Bioinformatics* **23**(1), 391 (2022)
919
920
- 921 [13] Schreck, J.S., Coley, C.W., Bishop, K.J.: Learning retrosynthetic planning through simulated
922 experience. *ACS central science* **5**(6), 970–981 (2019)
923
924
- 925 [14] Tu, Z., Coley, C.W.: Permutation invariant graph-to-sequence model for template-free retrosyn-
926 thesis and reaction prediction. *Journal of chemical information and modeling* **62**(15), 3503–3513
927
928
- 929 [15] Zhong, W., Yang, Z., Chen, C.Y.-C.: Retrosynthesis prediction using an end-to-end graph
930
931
932
933
934
935

- 936 generative architecture for molecular graph editing. *Nature Communications* **14**(1), 3009 (2023)
937
- 938 [16] Chen, L.-Y., Li, Y.-P.: Enhancing chemical synthesis: a two-stage deep neural network for
939 predicting feasible reaction conditions. *Journal of Cheminformatics* **16**(1), 1–14 (2024)
940
941
- 942 [17] Gao, H., Struble, T.J., Coley, C.W., Wang, Y., Green, W.H., Jensen, K.F.: Using machine
943 learning to predict suitable conditions for organic reactions. *ACS central science* **4**(11), 1465–
944 1476 (2018)
945
946
- 947 [18] Kwon, Y., Kim, S., Choi, Y.-S., Kang, S.: Generative modeling to predict multiple suitable
948 conditions for chemical reactions. *Journal of Chemical Information and Modeling* **62**(23), 5952–
949 5960 (2022)
950
951
- 952 [19] Maser, M.R., Cui, A.Y., Ryou, S., DeLano, T.J., Yue, Y., Reisman, S.E.: Multilabel classi-
953 fication models for the prediction of cross-coupling reaction conditions. *Journal of Chemical*
954 *Information and Modeling* **61**(1), 156–166 (2021)
955
956
- 957 [20] Ahneman, D.T., Estrada, J.G., Lin, S., Dreher, S.D., Doyle, A.G.: Predicting reaction
958 performance in *c*–*n* cross-coupling using machine learning. *Science* **360**(6385), 186–190 (2018)
959
960
- 961 [21] Chen, Y., Zhang, L.: How much can deep learning improve prediction of the responses to drugs
962 in cancer cell lines? *Briefings in bioinformatics* **23**(1), 378 (2022)
963
964
- 965 [22] Li, B., Su, S., Zhu, C., Lin, J., Hu, X., Su, L., Yu, Z., Liao, K., Chen, H.: A deep learning frame-
966 work for accurate reaction prediction and its application on high-throughput experimentation
967 data. *Journal of Cheminformatics* **15**(1), 1–12 (2023)
968
969
- 970 [23] Panteleev, J., Gao, H., Jia, L.: Recent applications of machine learning in medicinal chemistry.
971 *Bioorganic & medicinal chemistry letters* **28**(17), 2807–2815 (2018)
972
973
- 974 [24] Chen, L.-Y., Li, Y.-P.: *Machine Learning Applications in Chemical Kinetics and Thermochem-*
975 *istry*, pp. 203–226. Springer, ??? (2023)
976
977
- 978 [25] Lowe, D.: Chemical reactions from US patents (1976-Sep2016) (2017). [https://figshare.](https://figshare.com/articles/dataset/Chemical_reactions_from_US_patents_1976-Sep2016_/5104873)
979 [com/articles/dataset/Chemical_reactions_from_US_patents_1976-Sep2016_/5104873](https://figshare.com/articles/dataset/Chemical_reactions_from_US_patents_1976-Sep2016_/5104873) Accessed
980 September 05, 2023
981
982
- 983 [26] Kearnes, S.M., Maser, M.R., Wlekliński, M., Kast, A., Doyle, A.G., Dreher, S.D., Hawkins,
984 J.M., Jensen, K.F., Coley, C.W.: The open reaction database. *Journal of the American Chemical*
985
986
987
988
989
990

- 991 Society **143**(45), 18820–18826 (2021)
992
993 [27] Nextmove Software Pistachio (2023). <https://www.nextmovesoftware.com/pistachio.html>
994
995 Accessed September 05, 2023
996
997 [28] Reaxys (2023). <https://www.reaxys.com/> Accessed September 05, 2023
998
999
1000 [29] CAS, SciFinder-n (2023). <https://scifinder-n.cas.org/> Accessed September 05, 2023
1001
1002 [30] Roth, D.L.: SPRESIweb 2.1, a selective chemical synthesis and reaction database. ACS
1003 Publications (2005)
1004
1005
1006 [31] Gimadiev, T.R., Lin, A., Afonina, V.A., Batyrshin, D., Nugmanov, R.I., Akhmetshin, T.,
1007 Sidorov, P., Duybankova, N., Verhoeven, J., Wegner, J.: Reaction data curation i: chemical
1008 structures and transformations standardization. *Molecular Informatics* **40**(12), 2100119 (2021)
1009
1010
1011 [32] Chen, W.L., Chen, D.Z., Taylor, K.T.: Automatic reaction mapping and reaction center
1012 detection. *Wiley Interdisciplinary Reviews: Computational Molecular Science* **3**(6), 560–593
1013
1014 (2013)
1015
1016 [33] Lin, A., Dyubankova, N., Madzhidov, T.I., Nugmanov, R.I., Verhoeven, J., Gimadiev, T.R.,
1017 Afonina, V.A., Ibragimova, Z., Rakhimbekova, A., Sidorov, P.: Atom-to-atom mapping: a
1018 benchmarking study of popular mapping algorithms and consensus strategies. *Molecular*
1019 *Informatics* **41**(4), 2100138 (2022)
1020
1021
1022 [34] Nugmanov, R., Dyubankova, N., Gedich, A., Wegner, J.K.: Bidirectional graphormer for reac-
1023 tivity understanding: neural network trained to reaction atom-to-atom mapping task. *Journal*
1024 *of Chemical Information and Modeling* **62**(14), 3307–3315 (2022)
1025
1026 [35] Schwaller, P., Hoover, B., Reymond, J.-L., Strobelt, H., Laino, T.: Extraction of organic chem-
1027 istry grammar from unsupervised learning of chemical reactions. *Science Advances* **7**(15), 4166
1028 (2021)
1029
1030 [36] Nugmanov, R.I., Mukhametgaleev, R.N., Akhmetshin, T., Gimadiev, T.R., Afonina, V.A.,
1031 Madzhidov, T.I., Varnek, A.: Cgrtools: Python library for molecule, reaction, and condensed
1032 graph of reaction processing. *Journal of chemical information and modeling* **59**(6), 2516–2521
1033 (2019)
1034
1035 [37] Vaucher, A.C., Schwaller, P., Laino, T.: Completion of partial reaction equations. Chemrxiv
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045

1046 (2020)
1047
1048 [38] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L.,
1049 Polosukhin, I.: Attention is all you need. arXiv preprint arXiv:1706.03762 (2017)
1050
1051
1052 [39] Toniato, A., Schwaller, P., Cardinale, A., Geluykens, J., Laino, T.: Unassisted noise reduction
1053 of chemical reaction datasets. *Nature Machine Intelligence* **3**(6), 485–494 (2021)
1054
1055
1056 [40] Goodfellow, I.J., Mirza, M., Xiao, D., Courville, A., Bengio, Y.: An empirical investigation
1057 of catastrophic forgetting in gradient-based neural networks. arXiv preprint arXiv:1312.6211
1058 (2013)
1059
1060
1061 [41] De Nino, A., Bortolini, O., Maiuolo, L., Garofalo, A., Russo, B., Sindona, G.: A sustainable pro-
1062 cedure for highly enantioselective organocatalyzed diels–alder cycloadditions in homogeneous
1063 ionic liquid/water phase. *Tetrahedron letters* **52**(13), 1415–1417 (2011)
1064
1065
1066 [42] Özdemirhan, D.: Optically active tertiary alcohols by biocatalysis. *Synthetic Communications*
1067 **47**(7), 629–645 (2017)
1068
1069 [43] Dolfus, U., Briem, H., Rarey, M.: Visualizing generic reaction patterns. *Journal of Chemical*
1070 *Information and Modeling* **62**(19), 4680–4689 (2022)
1071
1072 [44] Chen, L.-Y.: AutoTemplate. <https://github.com/Lung-Yi/AutoTemplate> Accessed September
1073 05, 2023
1074
1075 [45] RDKit: Open-Source Cheminformatics Software. <http://www.rdkit.org/> Accessed September
1076 05, 2023
1077
1078 [46] Coley, C.W., Green, W.H., Jensen, K.F.: Rdcchiral: An rdkit wrapper for handling stereochem-
1079 istry in retrosynthetic template extraction and application. *Journal of chemical information*
1080 *and modeling* **59**(6), 2529–2537 (2019)
1081
1082 [47] Daylight SMARTS Documentation. [https://www.daylight.com/dayhtml/doc/theory/theory.](https://www.daylight.com/dayhtml/doc/theory/theory.smarts.html)
1083 [smarts.html](https://www.daylight.com/dayhtml/doc/theory/theory.smarts.html) Accessed September 05, 2023
1084
1085 [48] Dijkstra, E.W.: A note on two problems in connexion with graphs. In: Edsger Wybe Dijkstra:
1086 *His Life, Work, and Legacy*, pp. 287–290 (2022)
1087
1088 [49] Riesen, K., Jiang, X., Bunke, H.: Exact and inexact graph matching: Methodology and
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100

1101 applications. Managing and mining graph data, 217–247 (2010)
1102
1103 [50] McNitt, C.D., Popik, V.V.: Photochemical generation of oxa-dibenzocyclooctyne (odibo) for
1104 metal-free click ligations. *Organic & biomolecular chemistry* **10**(41), 8200–8202 (2012)
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133
1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155