

ReactXT: Understanding Molecular “Reaction-ship” via Reaction-Contextualized Molecule-Text Pretraining

Anonymous ACL submission

Abstract

Molecule-text modeling, which aims to facilitate molecule-relevant tasks with a textual interface and textual knowledge, is an emerging research direction. Beyond single molecules, studying reaction-text modeling holds promise for helping the synthesis of new materials and drugs. However, previous works mostly neglect reaction-text modeling: they primarily focus on modeling individual molecule-text pairs or learning chemical reactions without texts in context. Additionally, one key task of reaction-text modeling – experimental procedure prediction – is less explored due to the absence of an open-source dataset. The task is to predict step-by-step actions of conducting chemical experiments and is crucial to automating chemical synthesis. To resolve the challenges above, we propose a new pretraining method, **ReactXT**, for reaction-text modeling, and a new dataset, **OpenExp**, for experimental procedure prediction. Specifically, ReactXT features three types of input contexts to incrementally pretrain LMs. Each of the three input contexts corresponds to a pretraining task to improve the text-based understanding of either reactions or single molecules. ReactXT demonstrates consistent improvements in experimental procedure prediction and molecule captioning and offers competitive results in retrosynthesis. Our code is available at <https://anonymous.4open.science/r/ReactXT>.

1 Introduction

Multi-modal large language models (LMs) have recently attracted extensive research attention. Remarkably, in the vision-language domain, LMs enhanced with visual encoders show impressive results in visual question-answering and image captioning (Liu et al., 2023a; Li et al., 2023). Inspired by their successes, molecule-text modeling (MTM)

becomes an emerging research field (Zeng et al., 2022; Su et al., 2022; Liu et al., 2023b), aiming to build the natural language interface for molecular tasks, including text-guided molecule generation, molecule captioning, and molecule-text retrieval (Edwards et al., 2022; Liu et al., 2022).

Building upon these MTM works, we study reaction-text modeling (RTM), aiming to improve LMs’ performance on reaction-relevant tasks. Chemical reactions, involving the transformation of reactants into products, are fundamental to advancing drug discovery and material science (Schwaller et al., 2022). Revisiting prior works, we identify key research gaps in both the learning paradigm and the evaluation benchmark for RTM:

- **Learning Paradigm.** Most prior works either focus on generating the textual description of a single molecule (*cf.* Figure 1a) (Liu et al., 2023b; Edwards et al., 2022; Su et al., 2022), or apply LMs for chemical reaction prediction without including the textual descriptions of molecules/reactions in the context (*cf.* Figure 1b) (Christofidellis et al., 2023; Fang et al., 2023; Born and Manica, 2023). Such methods overlook the potential knowledge in textual descriptions to improve performance. Pioneer works (Guo et al., 2023; Shi et al., 2023) include labels of molecular roles and experimental conditions when prompting ChatGPT, but achieve suboptimal performances for being limited to prompt engineering.
- **Evaluation Benchmark.** An open-source dataset for experimental procedure prediction is notably missing. As illustrated in Figure 2, experimental procedure prediction aims to deduce the step-by-step actions for experimental execution through interpreting chemical reactions (Vaucher et al., 2021), which has a significant value for automating chemical synthesis processes (Vaucher et al., 2020; Zeng et al., 2023). This task aligns

chosen samples, 50 samples could be directly used without any human intervention, and 90 samples required only minor modifications for experimental execution (*cf.* Figure 5).

Our contributions can be summarized as follows:

- We propose ReactXT, a method that incorporates three types of input contexts to incrementally pretrain an LM. These contexts are tailored to enhance LMs’ understanding of chemical reactions and individual molecules.
- We curate an open-source experimental procedure prediction dataset OpenExp, a new benchmark for automating chemical synthesis research.
- ReactXT achieves state-of-the-art performances for experimental procedure prediction on the OpenExp dataset, highlighting its superior RTM ability. It also outperforms baselines by 3.2% for molecule captioning on the PubChem324k dataset. ReactXT has competitive performances for retrosynthesis, and we are refining it to surpass the current state-of-the-art method.

2 Related Works

Molecule-Text Modeling (MTM). MTM aims to jointly model molecules and texts to address text-related molecular tasks (Edwards et al., 2022, 2021). Molecules can be represented by 1D sequences of SMILES (Weininger, 1988) and SELFIES (Krenn et al., 2020), making it feasible to pretrain unified LMs on mixed 1D sequences of texts and molecules (Taylor et al., 2022; Edwards et al., 2022; Chithrananda et al., 2020; Zeng et al., 2022). Further, these LMs can be aligned to human preference via instruction tuning (Christofidellis et al., 2023; Fang et al., 2023). In parallel to 1D LMs, multi-modal methods are also studied, using graph neural networks (GNNs) (Hu et al., 2020) to encode 2D molecular graphs. Notably, CLIP-style (Radford et al., 2021) cross-modal contrastive learning and BLIP2-style (Li et al., 2023) cross-modal projector are both investigated to facilitate molecule-text retrieval (Su et al., 2022; Liu et al., 2022), and molecule-to-text generation (Liu et al., 2023b), respectively. However, prior works mainly focus on individual molecules rather than chemical reactions. To bridge the gap, ReactXT explores reaction-text modeling, facilitating reaction-relevant tasks with a text interface and textual knowledge.

Experimental Procedure Prediction. Synthesizing complex compounds requires detailed plan-

ning of synthetic pathways and intermediate steps, a process that is both labor-intensive and complex. Machine learning (ML) can potentially automate the process by predicting experimental procedures. Prior works have explored predicting reaction conditions (*e.g.*, catalyst and solvent) (Gao et al., 2018) and sequences of synthesis steps (Vaucher et al., 2021) by reading chemical reactions. Given known experimental procedures, ML is also explored to empower chemical lab robots (Burger et al., 2020), and automated lab pipelines (Coley et al., 2019; Nicolaou et al., 2020). Notably, tool-augmented GPT4 (OpenAI, 2023) is explored to plan and execute known chemical experiments (Boiko et al., 2023). Unlike prior works, our OpenExp dataset is the first open-source dataset to facilitate the procedure prediction of unseen chemical experiments.

Retrosynthesis and Chemical Reaction Prediction. Given a chemical reaction, retrosynthesis is to predict reactants from products and reaction prediction is to predict products from reactants (Schwaller et al., 2022). They can be formalized as sequence-to-sequence translation represented by SMILES strings (Liu et al., 2017; Irwin et al., 2022; Zhong et al., 2022; Tetko et al., 2020; Ucak et al., 2022). Concurrently, 2D molecular graphs are explored for reaction prediction: selection-based methods focus on classifying the most suitable reaction templates (Chen and Jung, 2021; Dai et al., 2019); and graph-based generative models directly synthesize target molecules (Shi et al., 2020; Sacha et al., 2021; Yan et al., 2020). However, the methods above leverage only reactions without texts. While notably two pioneer works apply ChatGPT for reaction prediction (Shi et al., 2023; Bran et al., 2023), their performances are limited to exploring only prompt engineering.

3 ReactXT: Reaction-Contextualized Molecule-Text Pretraining

ReactXT consists of two key components: 1) the method of creating input contexts to incrementally pretrain an LM, and 2) a balanced sampling strategy for the reaction contexts. We begin by introducing our multi-modal LM backbone, then proceed to elaborate on ReactXT’s two components.

Multi-Modal Language Model Backbone. Molecules can be represented by their 1D SMILES or 2D molecular graphs (Wells, 2012). We employ MolCA (Liu et al., 2023b) as our primary LM backbone to effectively harness both the 1D and

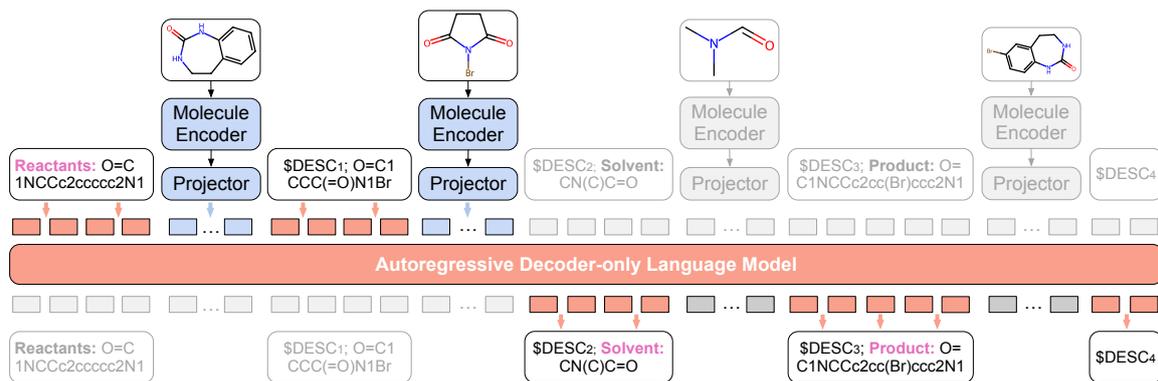


Figure 3: Illustration of Reaction-Contextualized Molecule-Text Pretraining. Example uses forward reaction context.

Context Type	Prompt Template
Forward reaction	Reactants: $\underbrace{\$SMI_1 \langle Mo1_1 \rangle \$DESC_1}_{\times n: \text{Number of reactants}}$; Solvent: $\$SMI_{n+1} \langle Mo1_{n+1} \rangle \$DESC_{n+1}$; Product: $\$SMI_{n+2} \langle Mo1_{n+2} \rangle \$DESC_{n+2} \langle STOP \rangle$
Backward reaction	Product: $\$SMI_1 \langle Mo1_1 \rangle \$DESC_1$; Solvent: $\$SMI_2 \langle Mo1_2 \rangle \$DESC_2$; Reactants: $\underbrace{\$SMI_3 \langle Mo1_3 \rangle \$DESC_3}_{\times n: \text{Number of reactants}} \langle STOP \rangle$
Random molecule	$\$SMI_1 \langle Mo1_1 \rangle \$DESC_1$; $\$SMI_2 \langle Mo1_2 \rangle \$DESC_2$; $\$SMI_3 \langle Mo1_3 \rangle \$DESC_3$; $\$SMI_4 \langle Mo1_4 \rangle \$DESC_4 \langle STOP \rangle$

Table 1: Prompt templates for creating input contexts. $\langle Mo1_i \rangle$ is the placeholder for the 2D graph embedding of the i -th molecule; $\$SMI_i$ and $\$DESC_i$ is the SMILES and textual description for the i -th molecule, respectively.

[Abstract] The invention relates to indole acetic acid compounds which function as antagonists of the CRTH2 receptor. The invention also relates to the use of these compounds to inhibit the binding of prostaglandin D2 and its metabolites or certain thromboxane metabolites to the CRTH2 receptor and to treat disorders responsive to such inhibition. [Properties] Molecular Weight: 547.60; XLogP3: 6.10; Hydrogen Bond Donor Count: 0; Hydrogen Bond Acceptor Count: 7; Rotatable Bond Count: 8; Exact Mass: 547.19; Monoisotopic Mass: 547.19; Topological Polar Surface Area: 89.40; Heavy Atom Count: 39; Formal Charge: 0; Complexity: 1020; Isotope Atom Count: 0; Defined Atom Stereocenter Count: 0; Undefined Atom Stereocenter Count: 0; Defined Bond Stereocenter Count: 0; Undefined Bond Stereocenter Count: 0; Covalently-Bonded Unit Count: 1; Compound Is Canonicalized: Yes.

Table 2: Molecule description example, including the patent abstract and the computed/experimental properties. The described molecule is Cc1c(C2=NN(CCC3CCCCC3)S(=O)(=O)C3CCCC32)c2cc(F)ccc2n1CC(=O)OC(C)(C)C.

225 2D molecular modalities. Specifically, MolCA in-
 226 corporates a GNN encoder (You et al., 2020)
 227 for encoding 2D molecular graphs. This GNN’s output
 228 then is mapped to an LM’s (i.e., Galactica; Taylor
 229 et al. (2022)) input space via a cross-modal projec-
 230 tor, thereby enabling the LM to perceive 2D molec-
 231 ular graphs. Both the cross-modal projector and the
 232 GNN have been pretrained for molecule-text align-
 233 ment (Li et al., 2023). MolCA shows promising
 234 performances when finetuned for molecule caption-
 235 ing and IUPAC name prediction.

3.1 Creating Input Contexts

236 Addressing the core challenges of LMs hinges on
 237 the careful selection of the input data. As shown in
 238 Table 1, ReactXT incorporates three types of input
 239 contexts to incrementally pretrain LMs: forward
 240 reaction context, backward reaction context, and
 241 random molecule context. These contexts are tai-
 242 lored for a text-based understanding of chemical
 243 reactions and individual molecules:

245 • **Forward Reaction Context.** As Figure 3 il-
 246 lustrates, the forward reaction context labels
 247 molecules according to their roles – Reactant,

248 Catalyst, Solvent, and Product – in the reac-
 249 tion, and arranges them in this specific sequen-
 250 tial order. Note, not every reaction has a Catalyst
 251 or Solvent. For each molecule, we append its
 252 2D molecular graph embeddings (e.g., $\langle Mo1_1 \rangle$;
 253 Liu et al. (2023b)) after its SMILES to enhance
 254 the LM’s understanding of molecular structures;
 255 and append molecular descriptions (e.g., $\$DESC_1$)
 256 following the 2D molecular graph embeddings
 257 to align molecules with texts.

258 • **Backward Reaction Context.** Similar to the
 259 forward context but with the order of molecular
 260 roles reversed, this context aims to combat the
 261 Reversal Curse (Berglund et al., 2023) of LMs:
 262 LMs trained on “A is B” fail to generalize to “B
 263 is A”. The reversal generalization is crucial be-
 264 cause downstream applications include backward
 265 retrosynthesis (Schwaller et al., 2022).

266 • **Random Molecule Context.** Introduced to en-
 267 sure LMs retain the capability to describe indi-
 268 vidual molecules outside chemical reactions.

269 **Context Length.** In each input context, we use
 270 up to k molecules and their descriptions, where

271 k is a hyperparameter. For reactions with over k
272 molecules, we apply weighted molecule sampling,
273 as explained in Section 3.2.

274 **Molecule Descriptions.** One crucial component
275 of the input contexts is the molecule description,
276 whose quality and comprehensiveness are vital for
277 molecule-text alignment. We collect molecular de-
278 scriptions and properties from multiple sources,
279 encompassing three types of content:

280 • **Molecule Patent Abstracts.** We source patent
281 abstracts from PubChem’s Patent View¹. These
282 abstracts typically describe molecular structures,
283 properties, or applications, but may also in-
284 clude irrelevant information if the molecule is
285 merely mentioned in passing rather than be-
286 ing the central subject. Despite the noise,
287 patent abstracts are indispensable for RTM: they
288 cover $\sim 95\%$ molecules in our employed reaction
289 databases (Lowe, 2017; Kearnes et al., 2021). In
290 contrast, the molecule-text datasets (Liu et al.,
291 2022, 2023b) derived from PubChem’s descrip-
292 tion section only cover $\sim 1\%$ of these molecules.

293 • **Computed and Experimental Properties.** We
294 retrieve these numerical properties from Pub-
295 Chem, aiming to enhance the understanding of
296 molecular structures through predictive learning.
297 Certain properties are also helpful for reaction
298 prediction. For example, knowing the solubility
299 helps determine concentrations when preparing
300 solutions; the knowledge of melting and boiling
301 points helps identify the states of matter at given
302 temperatures. Table 2 shows an example of a
303 patent abstract and computed/experimental prop-
304 erties. Table 12 includes detailed statistics of our
305 collected molecule properties.

306 • **PubChem Descriptions.** Following (Liu et al.,
307 2022, 2023b), we employ molecular descrip-
308 tions from PubChem. Due to their limited
309 coverage ($\sim 1\%$) for molecules in reaction
310 databases (Lowe, 2017; Kearnes et al., 2021),
311 we incorporate them exclusively for the random
312 molecule context.

313 **Autoregressive Language Modeling for Inter-**
314 **leaved Molecule-Text Sequences.** Given the input
315 contexts above of interleaved molecules and texts,
316 we apply language modeling loss to incrementally
317 pretrain the LM, molecule encoder, and projector.
318 We compute loss only for text tokens, excluding
319 2D molecular graph embeddings.

¹<https://pubchem.ncbi.nlm.nih.gov/docs/patents>

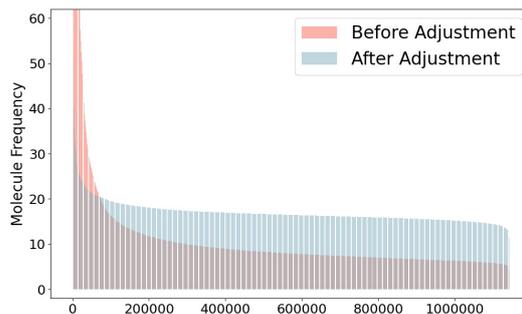


Figure 4: Distribution of molecules in the pretraining chemical reactions. For after adjustment, we conduct weighted sampling of chemical reactions matching the size of the pretraining dataset.

3.2 Balanced Sampling of Reaction Contexts

Figure 4 reveals a skewed distribution of molecules in chemical reactions (the red bars), with a small group of molecules appearing far more frequently than others. To address this imbalance, we develop a sampling strategy that promotes a fairer representation of molecules across reactions. This method reduces the dominance of commonly occurring molecules by adjusting 1) the sampling weight of each reaction r : $W(r)$, and 2) the sampling weight of each molecule m within a chosen reaction r : $W(m|r)$, based on the equations below:

$$W(r) = \frac{\sum_{m \in r} \text{Count}(m)}{\sum_{r' \in \mathcal{R}} \sum_{m \in r'} \text{Count}(m)}, \quad (1)$$

$$W(m|r) = \frac{1/\text{Count}(m)}{\sum_{m' \in r} 1/\text{Count}(m')}, \quad (2)$$

where \mathcal{R} denotes the dataset of chemical reactions; $\text{Count}(m)$ denotes molecule m ’s count in \mathcal{R} .

Equation (1) sets a reaction’s sampling weight inversely to the total occurrences of its molecules, favoring reactions with rare molecules; Equation (2) boosts the weights of rarer molecules within a given reaction. These weights are then applied for weighted random sampling without replacement (Efraimidis and Spirakis, 2006). The blue bars in Figure 4 present the sampling frequency of molecules after adjustment, showing a flatter distribution. Implementation details are in Appendix B.

4 OpenExp: An Open-Source Dataset for Experimental Procedure Prediction

Here we briefly introduce OpenExp’s curation process and defer the details to Appendix A.1. OpenExp is sourced from chemical reaction databases of USPTO-Applications (Lowe, 2017)

Total reactions	2262637	100%
Too large perplexity score	329160	14.55%
More than one product	105577	4.67%
Incomplete mapping of molecules (from chemical equation)	1034908	45.74%
Incomplete mapping of molecules (from action sequence)	178689	7.90%
Remove duplicate reactions	254099	11.23%
Filter out too short actions	14022	0.62%
Other errors	71743	3.16%
Remaining reactions	274439	12.13%

Table 3: Preprocessing steps and the number of samples removed at each step.

Dataset	Total	Train	Valid	Test	Open Source
Vaucher et al. (2021)	693k	555k	69k	69k	No
OpenExp, Ours	274k	220k	27k	27k	Yes

Table 4: Dataset statistics and comparison to prior work.

and ORD (Kearnes et al., 2021). As illustrated in Figure 2, these databases include chemical reactions and the corresponding unstructured descriptions of experimental procedures. To convert these unstructured descriptions into structured action sequences, we first run the pragraph2action model from (Christofidellis et al., 2023), and then conduct preprocessing following (Vaucher et al., 2021). The preprocessing is to remove low-quality data, eliminate duplicates, and construct molecule mapping between reactions and experimental procedures. Specific preprocessing steps are summarized in Table 3. Table 10 shows an example of the final dataset.

As shown in Table 4, the final OpenExp dataset includes 274k reaction-procedure pairs. It is randomly divided into train/valid/test sets by the 8:1:1 ratio. Compared to the prior work (Vaucher et al., 2021), which is closed-source for using the commercial Pistachio database², we open-source this dataset to assist future research.

To obtain insights on dataset quality, we invite two graduate students in chemistry to rate the alignment between the action sequences and their original descriptions, on a scale from 1 (lowest) to 5 (highest), as depicted in Figure 5. Briefly, of the 100 samples evaluated, 50 action sequences are deemed directly executable (scores above 4), and 90 are considered executable with slight manual adjustments (scores above 3).

²<https://www.nextmovesoftware.com/pistachio>

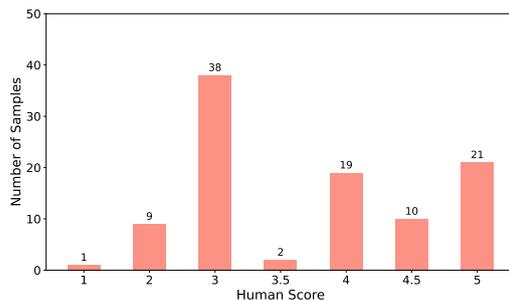


Figure 5: Human evaluations on OpenExp.

5 Experiment

We empirically evaluate ReactXT across three downstream tasks, including experimental procedural prediction, molecule captioning, and retrosynthesis. Further, we include ablation studies showcasing the contributions of individual components.

5.1 Experimental Setting

ReactXT is initialized by the stage-2 checkpoint of MolCA_{1.3B} (Liu et al., 2023b), if not specially noted. It is then pretrained using our proposed method, and subsequently finetuned for each downstream dataset separately. The context length k is 4. We employ full-parameter tuning for pretraining and finetuning. More details are in Appendix B.

ReactXT’s Pretraining Dataset. Our pretrain dataset includes PubChem324k’s pretrain subset (Liu et al., 2023b), which includes 298k molecule-text pairs, and 1.11 million chemical reactions from the USPTO-Applications (Lowe, 2017) and ORD (Kearnes et al., 2021) databases. For molecules in reactions, we obtain their patent abstracts and molecular properties following Section 3.1. To prevent information leakage, we have excluded 54k reactions that appear in the valid/test sets of the downstream datasets (*i.e.*, OpenExp, USPTO-50K (Schneider et al., 2016)) from the initial collection of 1.16 million reactions. See Appendix A.2 for more details.

Baselines. We compare ReactXT with the state-of-the-art LMs in science domain, including Galactica (Taylor et al., 2022), MolT5 (Edwards et al., 2022), TextChemT5 (Christofidellis et al., 2023), and MolCA (Liu et al., 2023b). For retrosynthesis and forward reaction prediction tasks, we also compare with task-specific LMs: R-SMILES (Zhong et al., 2022), AT (Tetko et al., 2020), MEGAN (Sacha et al., 2021), and Chemformer (Irwin et al., 2022). For captioning, we additionally compare against MoMu (Su et al., 2022).

Method	Validity	BLEU-2	BLEU-4	100%LEV	90%LEV	75%LEV	50%LEV	ROUGE-1	ROUGE-2	ROUGE-L
Random, among all reactions	63.2	34.5	19.1	0.0	0.0	0.0	13.6	46.6	18.1	36.4
Random, compatible pattern	100.0	37.8	22.1	0.0	0.0	0.1	16.5	47.8	21.0	38.4
Nearest neighbor	76.0	45.0	30.7	0.6	6.5	13.0	38.4	55.7	29.2	47.0
TextChemT5 _{220M}	99.3	54.1	40.6	0.4	4.6	13.7	61.2	61.5	40.3	56.4
MolT5-Large _{780M}	99.6	54.5	41.0	0.6	6.6	16.6	63.7	62.5	40.9	57.2
Galactica _{1.3B}	99.9	53.5	39.5	0.4	5.7	13.4	60.5	60.9	38.6	55.2
MolCA, Galac _{1.3B}	99.9	54.9	41.5	1.0	9.2	18.9	65.3	62.5	40.4	57.0
ReactXT, Galac _{1.3B} , Ours	100.0	57.4	44.0	1.0	9.5	22.6	70.2	64.4	42.7	58.9

Table 5: Comparison of experimental procedure prediction performances (%) on the OpenExp dataset. The subscript denotes each model’s parameter size. We conduct full-parameter fine-tuning for all models.

Method	BLEU-2	BLEU-4	ROUGE-1	ROUGE-2	ROUGE-L	METEOR
MolT5-Small _{80M}	14.8	8.5	26.5	13.5	23.6	18.5
MolT5-Base _{250M}	30.1	20.9	40.3	25.1	33.8	35.6
MolT5-Large _{780M}	30.2	22.2	41.5	25.9	34.8	36.6
Galactica _{1.3B} , LoRA ft	34.6	26.9	46.3	32.3	41.5	41.1
MoMu-Small _{82M}	19.1	12.0	29.7	16.3	26.7	21.8
MoMu-Base _{252M}	30.2	21.5	40.5	25.1	34.4	34.2
MoMu-Large _{782M}	31.1	22.8	41.8	25.7	36.7	36.2
MolCA, MolT5-Large _{877M}	32.9	26.3	49.8	35.7	44.2	42.4
MolCA, Galac _{125M}	31.9	24.3	47.3	33.9	43.2	41.6
MolCA, Galac _{1.3B} , LoRA ft	38.7	30.3	50.2	35.9	44.5	45.6
MolCA, Galac _{1.3B} , full ft*	39.4	32.2	52.7	39.4	47.6	49.2
ReactXT, Galac _{1.3B} , Ours	42.6	35.2	54.7	41.7	49.6	51.2

(a) PubChem324k dataset.

Method	BLEU-2	BLEU-4	ROUGE-1	ROUGE-2	ROUGE-L	METEOR
MolT5-Small _{80M}	51.9	43.6	62.0	46.9	56.3	55.1
MolT5-Base _{250M}	54.0	45.7	63.4	48.5	57.8	56.9
MolT5-Large _{780M}	59.4	50.8	65.4	51.0	59.4	61.4
TextChemT5 _{60M}	56.0	47.0	63.8	48.8	58.0	58.8
TextChemT5 _{220M}	62.5	54.2	68.2	54.3	62.2	64.8
MoMu-Small _{82M}	53.2	44.5	-	-	56.4	55.7
MoMu-Base _{252M}	54.9	46.2	-	-	57.5	57.6
MoMu-Large _{782M}	59.9	51.5	-	-	59.3	59.7
MolCA, Galac _{125M}	61.2	52.6	67.4	52.1	60.6	63.6
MolCA, Galac _{1.3B} , LoRA ft	62.0	53.1	68.1	53.7	61.8	65.1
ReactXT, Galac _{1.3B}	62.9	55.0	69.2	56.0	63.4	66.4

(b) CheBI-20 dataset.

Table 6: Molecule captioning performance (%) on the PubChem324k and CheBI-20 datasets. * denotes our re-implementation. Other baseline results are borrowed from (Liu et al., 2023b; Christofidellis et al., 2023).

Method	Top-1	Top-3	Top-5	Top-10
MEGAN	48.1	70.7	78.4	86.1
AT	53.5	-	81.0	85.7
Chemformer	54.3	-	62.3	63.0
<i>Train with aug., test without aug.</i>				
R-SMILES	51.2	74.9	81.1	83.0
MolT5-Large _{780M} *	53.9	69.9	74.6	77.3
ReactXT, Galac _{1.3B} , Ours	54.2	<u>70.9</u>	<u>74.9</u>	<u>78.3</u>
<i>Train with aug., test with aug.</i>				
R-SMILES	56.3	79.2	86.2	91.0
MolT5-Large _{780M} *	56.0	<u>76.0</u>	80.7	85.1
ReactXT, Galac _{1.3B} , Ours	<u>56.2</u>	<u>75.8</u>	<u>81.4</u>	<u>86.1</u>

Table 7: Retrosynthesis accuracies (%) on USPTO-50K. * denotes our re-implementation. Other baselines are from (Zhong et al., 2022). In each part, **bold** denotes the best result, and underline denotes the second best.

5.2 Experimental Procedure Prediction

Following (Vaucher et al., 2021), we employ the following evaluation metrics: Validity, which checks the syntactical correctness of the action sequence; machine-translation metrics BLUE (Papineni et al., 2002) and ROUGE (Lin, 2004); and the normalized Levenshtein similarity (Levenshtein et al., 1966). Specifically, 90%LEV denotes the proportion of predictions with a normalized Levenshtein score larger than 0.9. The three naive baselines based on random sampling and nearest neighbor are borrowed from (Vaucher et al., 2021). See Appendix B for details.

Table 5 presents the performances. We can observe that ReactXT consistently outperforms baselines across all metrics. Specifically, it sur-

Pretrain Input Context	Pretrain Data Type	BLEU-2	BLEU-4	75%LEV	50%LEV	ROUGE-1	ROUGE-2	ROUGE-L
No incremental pretrain	-	54.9	41.5	18.9	65.3	62.5	40.4	57.0
Random molecules	reaction, sing. mol.	56.6	43.2	20.9	69.4	63.8	41.9	58.3
Reactions w/o bal. samp.	reaction	56.8	43.3	21.3	69.2	64.0	42.1	58.5
Reactions	reaction	57.1	43.8	22.2	70.1	64.3	42.6	58.9
ReactXT	reaction, sing. mol.	57.4	44.0	22.6	70.2	64.4	42.7	58.9

Table 8: Ablation study of input contexts for incrementally pretrain MolCA, Galac_{1.3B}. Results are for experimental procedure prediction. Reactions denote both the forward reaction context and the backward reaction context.

passes baselines by 2.2% for BLEU-2 and 3.3% for 75%LEV, demonstrating ReactXT’s effectiveness for text-based reaction understanding.

5.3 Molecule Captioning

To evaluate ReactXT’s ability to understand single-molecules, we present its performances of molecule captioning on the PubChem324k (Liu et al., 2023b) and CheBI-20 (Edwards et al., 2022) datasets. We report metrics of BLEU (Papineni et al., 2002), ROUGE (Lin, 2004), and METEOR (Banerjee and Lavie, 2005).

Table 6 presents the captioning performances. We can observe that ReactXT consistently outperforms the baselines. Specifically, ReactXT shows improvements of 3.2% BLEU-2 and 2.3% ROUGE-2 scores on PubChem324k, and 1.7% ROUGE-2 on CheBI-20. These improvements underscore the effectiveness of our pretraining method for enhancing understanding of individual molecules.

5.4 Retrosynthesis

Retrosynthesis is to predict the reactant molecules given the product molecules. For this task, we employ the evaluation metrics of top-k accuracy, which measures the percentage of exact match to the ground truth in the top-k predictions. Following (Zhong et al., 2022), we use the root-aligned augmentations of SMILES during training and testing. Additionally, we report performances of testing without these augmentations.

Table 7 presents the results. We can observe that ReactXT outperforms MolT5-Large, which is also a multi-modal LM, in most metrics. This highlights the effectiveness of our approach among multi-modal methods. Further, we observe that ReactXT and MolT5 outperform R-SMILES in top-1 accuracy when testing without augmentation, but underperform R-SMILES for top-3, top-5, and top-10 accuracies. We conjecture that this discrepancy arises from a distribution shift between pretraining and finetuning: unlike R-SMILES, which uses root-aligned augmentations during pretraining, ReactXT and MolT5 do not. To address this, we are

pretraining a new ReactXT model that includes root-aligned augmentations.

5.5 Ablation Study

In this study, we ablate the key components of ReactXT, using the baseline of MolCA, Galac_{1.3B} without incremental pretraining. Table 8 presents the results. Specifically, we compare three variants of ReactXT: 1) pretraining with solely the random molecule contexts using the same pretrain dataset; 2) pretraining with forward and backward reaction contexts without the random molecule context; and 3) applying uniform sampling on reaction contexts instead of balanced sampling.

We can observe that 1) ReactXT’s full model shows the best performance, showing its performance is the integrated contribution of all components; 2) applying random molecule contexts alone improves upon the baseline, underscoring the valuable textual knowledge from our meticulously crafted pretraining dataset; 3) incorporating reaction contexts yields better results than random molecule contexts, highlighting the benefits of learning reaction knowledge during pretraining; and 4) balanced sampling improves the performance upon uniform sampling.

6 Conclusion and Future Works

In this work, we explore reaction-text modeling to empower reaction-relevant tasks with textual interfaces and knowledge. We present ReactXT, a pretraining method to learn chemical reactions within the context of the corresponding molecular textual descriptions. Additionally, we propose a new dataset OpenExp to support open-source research for experimental procedure prediction. ReactXT establishes the best performances across tasks of experimental procedure prediction and molecule captioning. It presents competitive performances for retrosynthesis. In future work, we plan to apply LMs to learn the interactions among large molecules (*e.g.*, proteins and nucleic acids), focusing on their dynamics and 3D spatial structures.

520 Limitations

521 In this and also the previous work (Vaucher et al.,
522 2021), the evaluation for experimental procedure
523 prediction is constrained to the comparison be-
524 tween the predictions and the reference action se-
525 quences. While improving this metric does reflect
526 the improvement in experimental design, it should
527 be acknowledged that the evaluation of real-world
528 chemical experiments is preferred for the devel-
529 oped models in future. For this purpose, the meth-
530 ods on automated chemistry pipelines (Boiko et al.,
531 2023; Coley et al., 2019; Nicolaou et al., 2020) can
532 be potentially considered.

533 Another limitation or future direction is improv-
534 ing the action space defined in our proposed Open-
535 Exp dataset, aiming to cover a wider range of chem-
536 ical experiments. For example, the action of ‘Pu-
537 rify’ is absent; and the action of ‘Concentration’
538 can be refined into operations such as ‘Evapora-
539 tion’ and ‘Pressurize’ for clearer instructions of
540 chemical experiments.

541 Potential Ethics Impact

542 In this study, the proposed method and dataset fo-
543 cus on chemical reactions and molecules, and in-
544 clude no human subjects. Consequently, we believe
545 this study presents no direct ethical concerns. How-
546 ever, the inclusion of LMs in our study does raise
547 potential issues, as LMs can be misused to produce
548 incorrect or biased information. Therefore, the
549 ethical implications of our work align with those
550 common to LM research, emphasizing the need for
551 responsible use and application of LMs.

552 References

553 Peter Atkins and Loretta Jones. 2007. *Chemical princi-*
554 *ples: The quest for insight*. Macmillan.

555 Satanjeev Banerjee and Alon Lavie. 2005. METEOR:
556 an automatic metric for MT evaluation with improved
557 correlation with human judgments. In *IEEvalu-*
558 *ation@ACL*, pages 65–72. Association for Computa-
559 tional Linguistics.

560 Lukas Berglund, Meg Tong, Max Kaufmann, Mikita
561 Balesni, Asa Cooper Stickland, Tomasz Korbak, and
562 Owain Evans. 2023. The reversal curse: Lms trained
563 on "a is b" fail to learn "b is a". *arXiv preprint*
564 *arXiv:2309.12288*.

565 Daniil A Boiko, Robert MacKnight, Ben Kline, and
566 Gabe Gomes. 2023. Autonomous chemical research
567 with large language models. *Nature*, 624(7992):570–
568 578.

Jannis Born and Matteo Manica. 2023. Regression
569 transformer enables concurrent sequence regression
570 and generation for molecular language modelling.
571 *Nat. Mac. Intell.*, 5(4):432–444. 572

Andres M Bran, Sam Cox, Oliver Schilter, Carlo Baldas-
573 sari, Andrew White, and Philippe Schwaller. 2023.
574 Augmenting large language models with chemistry
575 tools. In *NeurIPS 2023 AI for Science Workshop*. 576

Benjamin Burger, Phillip M Maffettone, Vladimir V
577 Gusev, Catherine M Aitchison, Yang Bai, Xiaoyan
578 Wang, Xiaobo Li, Ben M Alston, Buyi Li, Rob
579 Clowes, et al. 2020. A mobile robotic chemist. *Na-*
580 *ture*, 583(7815):237–241. 581

Shuan Chen and Yousung Jung. 2021. Deep retrosyn-
582 thetic reaction prediction using local reactivity and
583 global attention. *JACS Au*, 1(10):1612–1620. 584

Seyone Chithrananda, Gabriel Grand, and Bharath Ram-
585 sundar. 2020. Chemberta: large-scale self-supervised
586 pretraining for molecular property prediction. *arXiv*
587 *preprint arXiv:2010.09885*. 588

Dimitrios Christofidellis, Giorgio Giannone, Jannis
589 Born, Ole Winther, Teodoro Laino, and Matteo Man-
590 ica. 2023. Unifying molecular and textual represen-
591 tations via multi-task language modelling. In *ICML*. 592

Connor W Coley, Dale A Thomas III, Justin AM Lum-
593 miss, Jonathan N Jaworski, Christopher P Breen, Vic-
594 tor Schultze, Travis Hart, Joshua S Fishman, Luke
595 Rogers, Hanyu Gao, et al. 2019. A robotic platform
596 for flow synthesis of organic compounds informed
597 by ai planning. *Science*, 365(6453):eaax1566. 598

Hanjun Dai, Chengtao Li, Connor Coley, Bo Dai, and
599 Le Song. 2019. Retrosynthesis prediction with con-
600 ditional graph logic network. *Advances in Neural*
601 *Information Processing Systems*, 32. 602

Carl Edwards, Tuan Manh Lai, Kevin Ros, Garrett
603 Honke, Kyunghyun Cho, and Heng Ji. 2022. Trans-
604 lation between molecules and natural language. In
605 *EMNLP*, pages 375–413. Association for Computa-
606 tional Linguistics. 607

Carl Edwards, ChengXiang Zhai, and Heng Ji. 2021.
608 Text2mol: Cross-modal molecule retrieval with natu-
609 ral language queries. In *EMNLP (1)*, pages 595–607.
610 Association for Computational Linguistics. 611

Pavlos S Efrimidis and Paul G Spirakis. 2006.
612 Weighted random sampling with a reservoir. *Infor-*
613 *mation processing letters*, 97(5):181–185. 614

Yin Fang, Xiaozhuan Liang, Ningyu Zhang, Kangwei
615 Liu, Rui Huang, Zhuo Chen, Xiaohui Fan, and Hua-
616 jun Chen. 2023. Mol-instructions: A large-scale
617 biomolecular instruction dataset for large language
618 models. *CoRR*, abs/2306.08018. 619

Hanyu Gao, Thomas J Struble, Connor W Coley, Yu-
620 ran Wang, William H Green, and Klavs F Jensen.
621 2018. Using machine learning to predict suitable
622

623	conditions for organic reactions. <i>ACS central science</i> , 4(11):1465–1476.	677
624		678
625	Taicheng Guo, Kehan Guo, Bozhao Nan, Zhengwen Liang, Zhichun Guo, Nitesh V. Chawla, Olaf Wiest, and Xiangliang Zhang. 2023. What indeed can GPT models do in chemistry? A comprehensive benchmark on eight tasks. <i>CoRR</i> , abs/2305.18365.	679
626		680
627		681
628		682
629		683
630	Weihua Hu, Bowen Liu, Joseph Gomes, Marinka Zitnik, Percy Liang, Vijay Pande, and Jure Leskovec. 2020. Strategies for pre-training graph neural networks. In <i>ICLR</i> .	684
631		685
632		686
633		687
634	Ross Irwin, Spyridon Dimitriadis, Jiazhen He, and Esben Jannik Bjerrum. 2022. Chemformer: a pre-trained transformer for computational chemistry. <i>Machine Learning: Science and Technology</i> , 3(1):015022.	688
635		689
636		690
637		691
638		692
639	Steven M Kearnes, Michael R Maser, Michael Wleklin-ski, Anton Kast, Abigail G Doyle, Spencer D Dreher, Joel M Hawkins, Klavs F Jensen, and Connor W Coley. 2021. The open reaction database. <i>Journal of the American Chemical Society</i> , 143(45):18820–18826.	693
640		694
641		695
642		696
643		697
644	Mario Krenn, Florian Häse, AkshatKumar Nigam, Pascal Friederich, and Alán Aspuru-Guzik. 2020. Self-referencing embedded strings (SELFIES): A 100% robust molecular string representation. <i>Mach. Learn. Sci. Technol.</i> , 1(4):45024.	698
645		699
646		700
647		701
648		702
649	Vladimir I Levenshtein et al. 1966. Binary codes capable of correcting deletions, insertions, and reversals. In <i>Soviet physics doklady</i> , volume 10, pages 707–710. Soviet Union.	703
650		704
651		705
652		706
653	Junnan Li, Dongxu Li, Silvio Savarese, and Steven C. H. Hoi. 2023. BLIP-2: bootstrapping language-image pre-training with frozen image encoders and large language models. <i>CoRR</i> , abs/2301.12597.	707
654		708
655		709
656		710
657	Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In <i>Text summarization branches out</i> , pages 74–81.	711
658		712
659		713
660	Bowen Liu, Bharath Ramsundar, Prasad Kawthekar, Jade Shi, Joseph Gomes, Quang Luu Nguyen, Stephen Ho, Jack Sloane, Paul Wender, and Vijay Pande. 2017. Retrosynthetic reaction prediction using neural sequence-to-sequence models. <i>ACS central science</i> , 3(10):1103–1113.	714
661		715
662		716
663		717
664		718
665		719
666	Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023a. Visual instruction tuning. <i>arXiv preprint arXiv:2304.08485</i> .	720
667		721
668		722
669	Shengchao Liu, Weili Nie, Chengpeng Wang, Jiarui Lu, Zhuoran Qiao, Ling Liu, Jian Tang, Chaowei Xiao, and Anima Anandkumar. 2022. Multi-modal molecule structure-text model for text-based retrieval and editing. <i>CoRR</i> , abs/2212.10789.	723
670		724
671		725
672		726
673		727
674	Zhiyuan Liu, Sihang Li, Yanchen Luo, Hao Fei, Yixin Cao, Kenji Kawaguchi, Xiang Wang, and Tat-Seng Chua. 2023b. Molca: Molecular graph-language modeling with cross-modal projector and uni-modal adapter. In <i>EMNLP</i> , pages 15623–15638. Association for Computational Linguistics.	728
675		729
676		730
		731
		732
		733
		734
		735
		736
		737
		738
		739
		740
		741
		742
		743
		744
		745
		746
		747
		748
		749
		750
		751
		752
		753
		754
		755
		756
		757
		758
		759
		760
		761
		762
		763
		764
		765
		766
		767
		768
		769
		770
		771
		772
		773
		774
		775
		776
		777
		778
		779
		780
		781
		782
		783
		784
		785
		786
		787
		788
		789
		790
		791
		792
		793
		794
		795
		796
		797
		798
		799
		800

733	Yaorui Shi, An Zhang, Enzhi Zhang, Zhiyuan Liu, and Xiang Wang. 2023. ReLM: Leveraging language models for enhanced chemical reaction prediction . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023</i> , pages 5506–5520. Association for Computational Linguistics.	Zheni Zeng, Yuan Yao, Zhiyuan Liu, and Maosong Sun. 2022. A deep-learning system bridging molecule structure and biomedical text with comprehension comparable to human professionals. <i>Nature communications</i> , 13(1):862.	789
734			790
735			791
736			792
737			793
738			
739			
740	Bing Su, Dazhao Du, Zhao Yang, Yujie Zhou, Jiangmeng Li, Anyi Rao, Hao Sun, Zhiwu Lu, and Ji-Rong Wen. 2022. A molecular multimodal foundation model associating molecule graphs with natural language. <i>CoRR</i> , abs/2209.05481.	Zipeng Zhong, Jie Song, Zunlei Feng, Tiantao Liu, Lingxiang Jia, Shaolun Yao, Min Wu, Tingjun Hou, and Mingli Song. 2022. Root-aligned smiles: a tight representation for chemical reaction prediction. <i>Chemical Science</i> , 13(31):9023–9034.	794
741			795
742			796
743			797
744			798
745	Ross Taylor, Marcin Kardas, Guillem Cucurull, Thomas Scialom, Anthony Hartshorn, Elvis Saravia, Andrew Poulton, Viktor Kerkez, and Robert Stojnic. 2022. Galactica: A large language model for science. <i>CoRR</i> , abs/2211.09085.		
746			
747			
748			
749			
750	Igor V Tetko, Pavel Karpov, Ruud Van Deursen, and Guillaume Godin. 2020. State-of-the-art augmented nlp transformer models for direct and single-step retrosynthesis. <i>Nature communications</i> , 11(1):5575.		
751			
752			
753			
754	Umit V Ucak, Islambek Ashyrmamatov, Junsu Ko, and Juyong Lee. 2022. Retrosynthetic reaction pathway prediction through neural machine translation of atomic environments. <i>Nature communications</i> , 13(1):1186.		
755			
756			
757			
758			
759	Alain C Vaucher, Philippe Schwaller, Joppe Geluykens, Vishnu H Nair, Anna Iuliano, and Teodoro Laino. 2021. Inferring experimental procedures from text-based representations of chemical reactions. <i>Nature communications</i> , 12(1):2573.		
760			
761			
762			
763			
764	Alain C Vaucher, Federico Zipoli, Joppe Geluykens, Vishnu H Nair, Philippe Schwaller, and Teodoro Laino. 2020. Automated extraction of chemical synthesis actions from experimental procedures. <i>Nature communications</i> , 11(1):3601.		
765			
766			
767			
768			
769	David Weininger. 1988. Smiles, a chemical language and information system. 1. introduction to methodology and encoding rules. <i>J. Chem. Inf. Comput. Sci.</i> , 28(1):31–36.		
770			
771			
772			
773	Alexander Frank Wells. 2012. <i>Structural inorganic chemistry</i> . Oxford university press.		
774			
775	Chaochao Yan, Qianggang Ding, Peilin Zhao, Shuangjia Zheng, Jinyu Yang, Yang Yu, and Junzhou Huang. 2020. Retroxpert: Decompose retrosynthesis prediction like A chemist. In <i>NeurIPS</i> .		
776			
777			
778			
779	Yuning You, Tianlong Chen, Yongduo Sui, Ting Chen, Zhangyang Wang, and Yang Shen. 2020. Graph contrastive learning with augmentations. In <i>NeurIPS</i> .		
780			
781			
782	Zheni Zeng, Yi-Chen Nie, Ning Ding, Qian-Jun Ding, Wei-Ting Ye, Cheng Yang, Maosong Sun, E Weinan, Rong Zhu, and Zhiyuan Liu. 2023. Transcription between human-readable synthetic descriptions and machine-executable instructions: an application of the latest pre-training technology. <i>Chemical Science</i> , 14(35):9360–9373.		
783			
784			
785			
786			
787			
788			

A Dataset Details

A.1 Collection and Preprocessing of OpenExp

OpenExp is compiled from the raw data from the two following sources:

- **USPTO-Applications** (Lowe, 2017). This dataset comprises records of 1.94 million reactions and their corresponding applications from the United States Patent and Trademark Office (USPTO) published between 2001 and September 2016. We download the raw XML files from the Figshare website³. For each reaction in this dataset, we extract its key information from four elements: `<productList>`, which contains the products of the reaction; `<reactantList>`, detailing the reactants; `<spectatorList>`, encompassing the catalysts and solvents; and `<dl:paragraphText>`, which provides a textual description of the experimental procedures.
- **Open Reaction Database** (Kearnes et al., 2021). The ORD⁴ dataset contains over 2 million chemical reactions, which include detailed records of reaction conditions and experimental procedures. It includes data from the USPTO applications (2001-2016 Sep), USPTO-granted patents (1976-2016 Sep), and experimental records from chemical literature.

Paragraph2Action. As illustrated in Figure 2, these databases include chemical reactions and the corresponding unstructured descriptions of experimental procedures. The unstructured nature of these descriptions poses a significant challenge to 1) automate chemical synthesis with robots (Vaucher et al., 2020; Burger et al., 2020); and 2) apply ML methods to predict experimental procedures of unseen reactions. To address this, the task of paragraph2action (Vaucher et al., 2020; Zeng et al., 2023) is proposed, aiming to convert unstructured experimental procedure descriptions into structured, step-by-step instructions with pre-defined actions. In this study, we leverage the action space defined by (Vaucher et al., 2020, 2021), and the paragraph2action model released by (Christofidellis et al., 2023).

Preprocessing. Following (Vaucher et al., 2021), we conduct preprocessing after the paragraph2action conversion. The preprocessing has

³https://figshare.com/articles/dataset/Chemical_reactions_from_US_patents_1976-Sep2016_5104873?file=8664370

⁴<https://open-reaction-database.org>

Action	Occurrence	Action	Occurrence
Add	744,533	Wait	38,211
Stir	287,413	Recrystal.	25,600
Concentrate	276,551	PhaseSepa.	24,141
Yield	274,439	PH	21,756
MakeSolution	272,537	Quench	18,699
Filter	247,625	Partition	16,045
Wash	224,286	Triturate	13,390
DrySolution	178,248	DrySolid	6,435
CollectLayer	146,379	Degas	4,789
Extract	114,855	Microwave	2,237
SetTemp.	44,126	Sonicate	450
Reflux	43,296		

Table 9: Action space and actions’ occurrences in the OpenExp dataset.

two purposes: 1) extracting the important entities (*i.e.*, molecules) in experimental procedures and mapping all molecules to their precursors in the chemical reaction; 2) applying a rule-based filtration to improve the dataset quality. Our preprocessing strategy is inspired by (Vaucher et al., 2020), augmented with additional 2 steps: perplexity filtering and similar action aggregation. The complete preprocessing steps are listed below:

- **Perplexity Filtering.** To ensure the quality of the above translation step, we compute a perplexity score for each output and exclude samples with a score larger than 1.0. These perplexity scores are calculated using the TextChemT5 model.
- **Entity Recognition.** We extract all the molecules (either by name or SMILES) from the action sequences using the source codes of (Vaucher et al., 2020). Then, we conduct string matching of IUPAC names between the extracted molecules and those in the chemical reactions. STOUT (Rajan et al., 2021) and PubChemPy⁵ are used for the translation between IUPAC names and SMILES. If any molecule cannot be matched with its counterpart in the chemical reactions, we consider the reaction data invalid and remove it from the dataset. However, we permit the inclusion of certain common substances, such as common organic solvents, in every reaction. The names and SMILES expressions of the 134 common substances are included in our code. After entity recognition, we assign each entity a unique ID and update the experimental procedures by replacing the entity mentions with the corresponding entity IDs.

⁵<https://github.com/mcs07/PubChemPy>

Field	Value
Reactant	\$1\$: OC(CCC1CCCCN1)C(F)(F)F \$3\$: CC(C)(C)[Si](C)(C)Cl \$4\$: c1c[nH]cn1
Solvent	\$2\$: ClCCl
Catalyst	\$5\$: CN(C)c1cncnc1
Product	\$-1\$: CC(C)(C)[Si](C)(C)OC(CCC1CCCCN1)C(F)(F)F
Experimental Procedures	MAKESOLUTION with \$1\$ and \$2\$ (10 mL) ; ADD \$3\$ (616 mg, 4.1 mmol, 1.2 eq) at 0°C ; ADD \$4\$ (697 mg, 10.2 mmol, 3.0 eq) at 0°C ; ADD \$5\$ (415 ng, 3.4 mmol) at 0°C ; STIR for 36 hours ; CONCENTRATE ; YIELD \$-1\$ (970 mg, 89%).
Source	A solution of 700 mg (3.4 mmol) of 1,1,1-trifluoro-4-pyridin-2-ylbutan-2-ol in 10 mL of dichloromethane was treated with 616 mg (4.1 mmol, 1.2 eq.) of tert-butyldimethylsilyl chloride, 697 mg (10.2 mmol, 3.0 eq.) of imidazole and 415 ng (3.4 mmol) of 4-dimethylaminopyridine at 0° C. The resulting mixture was allowed to warm to room temperature and as stirred for 36 hours. Then the mixture w was concentrated and the residue was purified by flash chromatography to give 970 mg (89%) of 2-[3-(tert-butyldimethylsilyloxy)-4,4,4-trifluorobutyl]pyridine as a colorless oil.

Table 10: Illustrative example of the OpenExp dataset. **BOLDED BLUE** indicates pre-defined action.

- Common Substance Renaming. We standardized the nomenclature for common substances that are known by multiple names (*e.g.*, water may also be referred to as H₂O, pure water, water (aq.), *etc.*) to improve the dataset’s precision. Using PubChemPy, we align the different names to their standardized SMILES representations, allowing us to identify when different terms refer to the same molecule by comparing their SMILES expressions.
- Similar Action Aggregation. If two adjacent operations are highly similar (*e.g.*, *STIR* and *STIR for 5 min*), they are merged together.
- Ensuring Single Product. This dataset focuses on the preparation of a single material, hence we remove reactions that yield multiple products.
- Action Filtering. We remove action sequences that have fewer than five actions or contain invalid actions.
- Reaction Deduplication. We remove the duplicated reactions from the dataset.

Table 11 presents the number of samples removed at each preprocessing step. Further, Table 10 provides an example from the final OpenExp dataset, we can observe that it encompasses:

Total reactions	2262637	100%
Too large perplexity score	329160	14.55%
More than one product	105577	4.67%
Incomplete mapping of molecules (from chemical reaction)	1034908	45.74%
Incomplete mapping of molecules (from action sequence)	178689	7.90%
Remove duplicate reactions	254099	11.23%
Filter out too short actions	14022	0.62%
Other errors	71743	3.16%
Remaining reactions	274439	12.13%

Table 11: Number of samples removed at each preprocessing step.

- Structured, step-by-step instructions of experimental procedures; 904 905
 - All molecules in the reaction and their roles (*i.e.*, reactant, solvent, catalyst, product). 906 907
 - The mapping between the recognized entities (*i.e.*, molecules) and their IDs. 908 909
 - The original unstructured experimental procedures. 910 911
- Discussion on License.** The ORD database is accessible under the CC-BY-SA license, and the USPTO-Applications dataset is available under the CC0 license. We have used codes from TextChemT5 (Christofidellis et al., 2023) and Paragraph2Actions (Vaucher et al., 2021), which are 912 913 914 915 916 917

918 both licensed under the MIT license. Therefore, we
919 will release OpenExp under the CC-BY-SA license
920 to comply with the most restrictive license of these
921 resources. This license permits content distribution
922 and sharing, provided the same license is applied.

923 **Human Evaluation.** We invite two graduate stu-
924 dents majoring in chemistry to evaluate the quality
925 of the OpenExp dataset. They are compensated
926 by the local average hourly rate. Specifically, we
927 randomly sample 100 data points, the evaluators
928 are then asked to rate the quality of each data point
929 on a scale from 1 (lowest) to 5 (highest). Our in-
930 structions to the evaluators are shown below:

Instructions to human evaluators.

We are curating a dataset partially generated by an AI model and want to seek feedback on its quality from human experts. During the evaluation process, we will provide both machine language sequences (the machine-generated operational sequences of experimental actions) and the corresponding natural language sequences (descriptions of experimental procedures in their original free texts).

You should rate these samples based on how well the operational sequences align with the original descriptions. Please use a rating scale of 1 (low alignment) to 5 (high alignment). Molecular skeletal formulas are provided as images for reference during evaluation. All original data for this dataset come from the United States Patent and Trademark Office (USPTO), ensuring the viability of the reactions.

The following are the detailed scoring guidelines, with a maximum score of 5:

- **5:** The AI model’s output captures key operations and experimental details present in the original description.
- **4:** No key experimental steps are missing compared to the original description. Minor discrepancies of in experimental details may exist, but they do not impede the execution of the experiment.
- **3:** There are discrepancies in key steps compared to the original description, yet these can be rectified with minor manual modifications to successfully carry out the experiment.
- **2:** There are significant differences in key experimental steps compared to the original description, requiring manual corrections on more than 50% of the sequence.
- **1:** The AI model’s output differs substantially from the original description, rendering it ineffective.

931
932 Each data point is evaluated by a single evalua-
933 tor. Figure 5 presents the human evaluation results.
934 In certain cases, evaluators are undecided about
935 assigning a lower or higher score, leading to the
936 assignment of decimal scores (*e.g.*, 3.5 and 4.5).

A.2 Collection and Preprocessing of ReactXT’s Pretraining Dataset

937
938
939 In Section 3, we collect and compile a dataset to
940 incrementally pretrain an LM for improved un-
941 derstanding of chemical reactions and individual
942 molecules. Here we elaborate on the details of this
943 dataset, which includes the following contents:

- A total of 1,162,551 chemical reactions; 944
- Patent abstracts and computed/experimental 945
properties of 1,254,157 molecules, which are all 946
from the chemical reactions. 947

948 We extract chemical reactions from ORD and 949
USPTO datasets. Then, we source patent ab- 950
stracts from PubChem’s Patent View⁶ and obtain 951
molecular properties using the PubChem’s Pub- 952
View API⁷. For each molecule, the abstract text de- 953
rives from the abstracts of patent documents where 954
the molecule is mentioned, and its properties in- 955
clude both computational and experimental ones. 956
Table 12 shows a complete list of these properties. 957

958 In Table 13, we compare the statistics of our pre- 959
training dataset with that of PubChem324k. We 960
can observe that ReactXT’s pretraining dataset in- 961
cludes more molecules and additionally includes 962
chemical reactions. 963

964 To prevent information leakage, we exclude a 965
total of 54,403 reactions that appear in the vali- 966
dation and test sets of the downstream datasets 967
(*i.e.*, OpenExp and USPTO-50K (Schneider et al., 968
2016)) from the pretraining dataset. The remaining 969
1,108,148 reactions are used for pretraining. 970

971 **Discussion on License.** The ORD database is 972
accessible under the CC-BY-SA license, and the 973
USPTO-Applications dataset is available under the 974
CC0 license. The patent abstracts from PubChem 975
are provided by Google Patent⁸, which is released 976
under the CC-BY-4.0 license. To comply with the 977
strictest license terms, we will release our dataset 978
under the CC-BY-SA license. 979

980 Additionally, we have utilized textual descrip- 981
tions, computed properties, and experimental prop- 982
erties from the PubChem website for pretraining.
Given that this data is aggregated from various
sources by PubChem, determining a single appro-
priate license is challenging. To support future
research while avoiding licensing complexities, we

⁶pubchem.ncbi.nlm.nih.gov/docs/patents

⁷pubchem.ncbi.nlm.nih.gov/docs/pug-view

⁸patents.google.com

Computed Properties			Experimental Properties				
Property	Count	Property	Count	Property	Count	Property	Count
Molecular Weight	1244109	Physical Description	8368	Vapor Density	1043	Enthalpy of Sublimation	9
Hydrogen Bond Donor Count	1244109	Kovats Retention Index	6878	Autoignition Temperature	771	Acid Value	4
Hydrogen Bond Acceptor Count	1244109	Solubility	5909	Heat of Vaporization	583	Dielectric Constant	2
Rotatable Bond Count	1244109	Chemical Classes	5726	Viscosity	550	Dispersion	1
Exact Mass	1244109	Melting Point	4468	Taste	514	Hydrophobicity	1
Monoisotopic Mass	1244109	Vapor Pressure	3032	Henry’s Law Constant	502		
Topological Surface Area	1244109	Boiling Point	2996	Surface Tension	448		
Heavy Atom Count	1244109	Color/Form	2927	pH	444		
Formal Charge	1244109	Density	2862	Odor Threshold	442		
Complexity	1244109	LogP	2763	Corrosivity	410		
Isotope Atom Count	1244109	Other Experimental Properties	2393	Heat of Combustion	405		
Defined Atom Stereocenter Count	1244109	Decomposition	2033	Ionization Efficiency	332		
Undefined Atom Stereocenter Count	1244109	Refractive Index	1777	Optical Rotation	265		
Defined Bond Stereocenter Count	1244109	Collision Cross Section	1634	Ionization Potential	253		
Undefined Bond Stereocenter Count	1244109	Odor	1512	LogS	166		
Covalently-Bonded Unit Count	1244109	Stability/Shelf Life	1506	Polymerization	134		
Compound Is Canonicalized	1244109	Flash Point	1479	Relative Evaporation Rate	101		
XLogP3	1184175	Dissociation Constants	1250	Caco2 Permeability	79		

Table 12: Statistics of the collected molecule properties, including computed properties and experimental properties.

	Our Dataset	Pubchem324k
Num of Molecules	1, 254, 157	313, 083
Num of Reactions	1, 162, 551	-
Avg. Molecule Weight	362.4	502.4
Avg. Atom Count	24.9	35.2
Avg. Bond Count	26.8	37.6
Avg. Ring Count	2.9	3.5
Avg. Text Length	517.8	120.4
Avg. Property Count	17.8	-

Table 13: Statistics of ReactXT’s pretraining dataset and Pubchem324k.

will provide the scripts for downloading and pre-processing this data, rather than distributing the data directly.

B Experimental Details

B.1 Hyperparameters

Here we detail the hyperparameters for ReactXT’s pretraining and finetuning across three downstream tasks. Due to the prohibitive costs associated with training large LMs, finetuning on downstream datasets is limited to a single run.

ReactXT Pretrain. The pretraining stage of ReactXT has 5 million steps, with the number of

molecules per reaction being $k = 4$. Following MolCA’s (Liu et al., 2023b) experimental setup, we employ a Q-former with 8 query tokens. We use AdamW as the optimizer, with a weight decay set to 0.05. The optimizer’s peak learning rate is set to 1×10^{-4} , scheduled by linear warmup with cosine decay. The warmup has 1000 steps and starts at a learning rate of 1×10^{-6} .

Experimental Procedure Prediction. We fully finetune all the baseline methods and ReactXT for 20 epochs, with a batch size of 32. The optimizer and learning rate settings are consistent with the pretraining phase.

Retrosynthesis. Following (Zhong et al., 2022), we sample 20 root-aligned augmentations for the training and testing subsets. We train MolT5 for 20 epochs and ReactXT for 10 epochs on the augmented training set using a batch size of 32. During testing, we conduct a beam search with a beam size of 20 for both models and return the top ten results as the model’s predictions. The optimizer and learning rate settings are kept consistent with the pretraining phase.

Pretrain Input Context	Pretrain Data Type	BLEU-2	BLEU-4	ROUGE-1	ROUGE-2	ROUGE-L	METEOR
No incremental pretrain	-	39.4	32.2	52.7	39.4	47.6	49.2
Reactions	reaction	37.3	29.9	50.3	36.5	45.0	46.7
ReactXT	reaction, sing. mol.	42.6	35.2	54.7	41.7	49.6	51.2

Table 14: Ablation study. Performances (%) for molecule captioning on the PubChem324k dataset.

Molecule Captioning. On both datasets, we full finetune MolCA and ReactXT 20 epochs, with a batch size of 32. The optimizer and learning rate settings are consistent with the pretraining phase.

B.2 Other Implementation Details

Baselines. We briefly introduce the baselines:

- **Galactica** (Taylor et al., 2022). Galactica is a scientific language model which is pretrained on 2 million compounds from PubChem. It has a decent understanding of SMILES formulas.
- **MolT5** (Edwards et al., 2022). MolT5 is developed based on the T5 model. Its training corpora include both natural language and SMILES data, making it suitable for both molecule captioning and text-based molecular generation tasks.
- **TextChemT5** (Christofidellis et al., 2023). TextChemT5 is a T5-based multi-domain LM, which is tuned on various text-molecule tasks.
- **MolCA** (Liu et al., 2023b). MolCA is a multi-modal language model finetuned on Galactica. It includes both graph encoder and LM, where a Querying Transformer is applied to align their latent spaces.
- **AT** (Tetko et al., 2020). AT trains transformers with data augmentation for retrosynthesis. The data augmentation is achieved by rearranging the order of characters in SMILES strings in both the training and test sets.
- **MEGAN** (Sacha et al., 2021). MEGAN represents chemical reactions as a sequence of graph edits and performs retrosynthesis by sequentially modifying the target molecule.
- **MoMu** (Su et al., 2022). Momu contrastively pre-trains a GNN and an LM with paired molecular graph-text data, and can be adapted to retrieval and generation tasks.
- **Chemformer** (Irwin et al., 2022). Chemformer is a Transformer-based molecule LM that is self-supervised pretrained on a SMILES corpus. It

can be applied to both generation and property prediction tasks.

- **Random, among all reactions** (Vaucher et al., 2021). Randomly pick an action sequence from the training set.
- **Random, compatible pattern** (Vaucher et al., 2021). Randomly pick an action sequence from the training subset of reactions that have the same number of molecules as the current reaction.
- **Nearest Neighbor** (Vaucher et al., 2021). Pick the action sequence from the training set with the reaction most similar to the current one, as determined by reaction fingerprints (Schwaller et al., 2019).

C More Experimental Results

C.1 Ablation Study

Table 14 presents an ablation study examining the impact of input contexts on molecule captioning. The removal of the random molecule context results in diminished captioning performance. This observation can be attributed to two factors: 1) including the PubChem324k dataset, which is used for creating random molecule contexts, is important to maintain molecule captioning performance; and 2) without random molecule contexts, the LM becomes overly dependent on reaction contexts, compromising its capability to accurately caption individual molecules. This finding underscores the significance of incorporating random molecule contexts in training.

C.2 Case Studies and Error Analysis

In this section, we present case studies from the experimental procedure prediction task to inform future research. We include examples of accurate predictions (see Table 15), inaccurate predictions (see Tables 16), and predictions that are different from the annotations but may also work (see Table 17 and Table 18). Our selection criteria prioritizes the accuracy of action sequences and the correct identification of primary materials, while

Field	Value			
Reactant	\$1\$: OCCCCCCCc1ccccc1 \$2\$: C#CC(=O)O \$4\$: c1ccccc1			
Catalyst	\$3\$: Cc1ccc(S(=O)(=O)O)cc1			
Product	\$-1\$: C#CC(=O)OCCCCCCCc1ccccc1			
Source	A mixture of 0.5 g of 7-phenylheptanol, 0.27 g of propiolic acid, 0.005 g of p-toluenesulfonic acid and 25 ml of benzene was refluxed with stirring for six hours while water formed was removed by a Dean-Stark water separator. After the reaction was completed, the reaction solution was washed successively with a 5% aqueous sodium bicarbonate solution and a saturated sodium chloride solution, and dried over anhydrous magnesium sulfate. After removal of the solvent under reduced pressure, the obtained residue was subjected to silicagel column chromatography to yield 0.368 g of 7-phenylheptyl propiolate (compound 3).			
Annotated Actions	<table border="0"> <tr> <td style="vertical-align: top;"> <p>MAKESOLUTION with \$1\$ (1.1 g) and \$2\$ (0.005 g) and \$3\$ (25 ml) and \$4\$; REFLUX for 6 hours ; CONCENTRATE ; WASH with NaHCO₃ ; WASH with sodium chloride ; DRYSOLUTION over magnesium sulfate ; FILTER keep filtrate ; YIELD \$-1\$ (1.15 g).</p> </td> <td style="vertical-align: top; text-align: center;"> <p>Predicted Actions</p> </td> <td style="vertical-align: top;"> <p>MAKESOLUTION with \$1\$ (0.27 g) and \$2\$ (0.005 g) and \$3\$ (25 ml) and \$4\$; REFLUX for 10 hours ; CONCENTRATE ; WASH with NaHCO₃ ; WASH with sodium chloride ; DRYSOLUTION over magnesium sulfate ; FILTER keep filtrate ; YIELD \$-1\$ (0.368 g).</p> </td> </tr> </table>	<p>MAKESOLUTION with \$1\$ (1.1 g) and \$2\$ (0.005 g) and \$3\$ (25 ml) and \$4\$; REFLUX for 6 hours ; CONCENTRATE ; WASH with NaHCO₃ ; WASH with sodium chloride ; DRYSOLUTION over magnesium sulfate ; FILTER keep filtrate ; YIELD \$-1\$ (1.15 g).</p>	<p>Predicted Actions</p>	<p>MAKESOLUTION with \$1\$ (0.27 g) and \$2\$ (0.005 g) and \$3\$ (25 ml) and \$4\$; REFLUX for 10 hours ; CONCENTRATE ; WASH with NaHCO₃ ; WASH with sodium chloride ; DRYSOLUTION over magnesium sulfate ; FILTER keep filtrate ; YIELD \$-1\$ (0.368 g).</p>
<p>MAKESOLUTION with \$1\$ (1.1 g) and \$2\$ (0.005 g) and \$3\$ (25 ml) and \$4\$; REFLUX for 6 hours ; CONCENTRATE ; WASH with NaHCO₃ ; WASH with sodium chloride ; DRYSOLUTION over magnesium sulfate ; FILTER keep filtrate ; YIELD \$-1\$ (1.15 g).</p>	<p>Predicted Actions</p>	<p>MAKESOLUTION with \$1\$ (0.27 g) and \$2\$ (0.005 g) and \$3\$ (25 ml) and \$4\$; REFLUX for 10 hours ; CONCENTRATE ; WASH with NaHCO₃ ; WASH with sodium chloride ; DRYSOLUTION over magnesium sulfate ; FILTER keep filtrate ; YIELD \$-1\$ (0.368 g).</p>		

(a) Example 1.

Field	Value			
Reactant	\$1\$: C[Si]1(C)CC[Si](C)(C)N1c1ccc(C(O)c2cn(S(=O)(=O)c3ccccc3)c3ncc(Cl)cc23)cn \$2\$: Nc1ccc(C(O)c2cn(S(=O)(=O)c3ccccc3)c3ncc(Cl)cc23)cn1 \$4\$: CC[SiH](CC)CC \$5\$: O=C(O)C(F)(F)F			
Solvent	\$3\$: ClCCl			
Product	\$-1\$: Nc1ccc(Cc2cn(S(=O)(=O)c3ccccc3)c3ncc(Cl)cc23)cn1			
Source	To (1-benzenesulfonyl-5-chloro-1H-pyrrolo[2,3-b]pyridin-3-yl)-[6-(2,2,5,5-tetramethyl-[1,2,5]azadisilolidin-1-yl)-pyridin-3-yl]-methanol and (6-amino-pyridin-3-yl)-(1-benzenesulfonyl-5-chloro-1H-pyrrolo[2,3-b]pyridin-3-yl)-methanol (118, 119, 1.70/1.25 g mix, 2.41 mmol) in 25.0 mL of dichloromethane, triethylsilane (3.00 mL, 18.8 mmol) and trifluoroacetic acid (1.50 mL, 19.5 mmol) were added and the reaction stirred at room temperature overnight. The reaction was concentrated under vacuum, combined with aqueous potassium carbonate and extracted with ethyl acetate. The organic layer was dried over sodium sulfate, filtered and the filtrate concentrated under vacuum. The resulting material was purified by silica gel column chromatography eluting with 20-100% ethyl acetate in hexane to provide the desired compound (120, 0.70 g).			
Annotated Actions	<table border="0"> <tr> <td style="vertical-align: top;"> <p>MAKESOLUTION with \$1\$ and \$2\$ and \$3\$ (25.0 mL) ; ADD \$4\$ (3.00 mL, 18.8 mmol) ; ADD \$5\$ (1.50 mL, 19.5 mmol) ; STIR for overnight at room temperature ; CONCENTRATE ; ADD K₂CO₃ ; EXTRACT with ethyl acetate ; COLLECTLAYER organic ; DRYSOLUTION over sodium sulfate ; FILTER keep filtrate ; CONCENTRATE ; YIELD \$-1\$ (0.70 g).</p> </td> <td style="vertical-align: top; text-align: center;"> <p>Predicted Actions</p> </td> <td style="vertical-align: top;"> <p>MAKESOLUTION with \$1\$ (1.00 g, 1.91 mmol) and \$2\$ (0.69 g, 1.72 mmol) and \$3\$ (35 mL) ; ADD \$4\$ (1.35 mL, 7.84 mmol) ; ADD \$5\$ (1.90 mL, 26.7 mmol) ; STIR for 8 h at room temperature ; CONCENTRATE ; EXTRACT with K₂CO₃ ; EXTRACT with ethyl acetate ; COLLECTLAYER organic ; DRYSOLUTION over sodium sulfate ; FILTER keep filtrate ; CONCENTRATE ; YIELD \$-1\$ (0.13 g, 19%).</p> </td> </tr> </table>	<p>MAKESOLUTION with \$1\$ and \$2\$ and \$3\$ (25.0 mL) ; ADD \$4\$ (3.00 mL, 18.8 mmol) ; ADD \$5\$ (1.50 mL, 19.5 mmol) ; STIR for overnight at room temperature ; CONCENTRATE ; ADD K₂CO₃ ; EXTRACT with ethyl acetate ; COLLECTLAYER organic ; DRYSOLUTION over sodium sulfate ; FILTER keep filtrate ; CONCENTRATE ; YIELD \$-1\$ (0.70 g).</p>	<p>Predicted Actions</p>	<p>MAKESOLUTION with \$1\$ (1.00 g, 1.91 mmol) and \$2\$ (0.69 g, 1.72 mmol) and \$3\$ (35 mL) ; ADD \$4\$ (1.35 mL, 7.84 mmol) ; ADD \$5\$ (1.90 mL, 26.7 mmol) ; STIR for 8 h at room temperature ; CONCENTRATE ; EXTRACT with K₂CO₃ ; EXTRACT with ethyl acetate ; COLLECTLAYER organic ; DRYSOLUTION over sodium sulfate ; FILTER keep filtrate ; CONCENTRATE ; YIELD \$-1\$ (0.13 g, 19%).</p>
<p>MAKESOLUTION with \$1\$ and \$2\$ and \$3\$ (25.0 mL) ; ADD \$4\$ (3.00 mL, 18.8 mmol) ; ADD \$5\$ (1.50 mL, 19.5 mmol) ; STIR for overnight at room temperature ; CONCENTRATE ; ADD K₂CO₃ ; EXTRACT with ethyl acetate ; COLLECTLAYER organic ; DRYSOLUTION over sodium sulfate ; FILTER keep filtrate ; CONCENTRATE ; YIELD \$-1\$ (0.70 g).</p>	<p>Predicted Actions</p>	<p>MAKESOLUTION with \$1\$ (1.00 g, 1.91 mmol) and \$2\$ (0.69 g, 1.72 mmol) and \$3\$ (35 mL) ; ADD \$4\$ (1.35 mL, 7.84 mmol) ; ADD \$5\$ (1.90 mL, 26.7 mmol) ; STIR for 8 h at room temperature ; CONCENTRATE ; EXTRACT with K₂CO₃ ; EXTRACT with ethyl acetate ; COLLECTLAYER organic ; DRYSOLUTION over sodium sulfate ; FILTER keep filtrate ; CONCENTRATE ; YIELD \$-1\$ (0.13 g, 19%).</p>		

(b) Example 2.

Table 15: Examples of accurate experimental procedure predictions.

1097 overlooking specifics like material quantities and
1098 temperatures. All the examples are from the test
1099 set of OpenExp.

1100 Table 15 displays two examples where experi-
1101 mental procedures are accurately predicted, show-
1102 ing close alignment between predicted and anno-
1103 tated actions, albeit with slight variances in mate-
1104 rial quantities and experiment times. These cases
1105 highlight the capability of LMs to predict experi-
1106 mental procedures, suggesting a path toward au-
1107 tomating chemical synthesis.

1108 Table 16 displays two failed examples of experi-
1109 mental procedure prediction. The predicted action
1110 sequences significantly deviate from the annotated
1111 sequences, making them impractical. Additionally,
1112 we can observe one common error of repetition,
1113 with the same or similar actions being duplicated.

1114 Tables 17 and Table 18 showcase three exam-
1115 ples where the predictions, while different from
1116 the annotations, could still be viable. In Example
1117 5, as an alternative to the annotated 'EXTRACT
1118 with ethyl acetate', the model proposes a series of
1119 actions ('COLLECT LAYER', 'WASH with ethyl
1120 acetate', 'DRY SOLUTION', and 'FILTER'), serv-
1121 ing a similar function. In Example 6, instead of
1122 the specified 'SET TEMPERATURE' and 'STIR',
1123 the model recommends 'STIR for 1h at 0 °C', serv-
1124 ing the same purpose. In Example 7, the model
1125 suggests adding components ('ADD \$4\$', 'ADD
1126 \$5\$', 'ADD \$6\$') sequentially rather than making
1127 a single solution as annotated, which could also be
1128 effective.

Field	Value
Reactant	\$1\$: <chem>Nc1ccc(C(=O)N[C@H](CO)Cc2ccccc2)c(/C=C/c2ccccc2)c1</chem> \$3\$: <chem>CC(=O)OC(C)=O</chem>
Solvent	\$2\$: <chem>C1CCOC1</chem>
Product	\$-1\$: <chem>CC(=O)Nc1ccc(C(=O)N[C@H](CO)Cc2ccccc2)c(/C=C/c2ccccc2)c1</chem>
Source	1 g (2.7 mmol) of (S)-4-amino-2-(E-2-phenylethen-1-yl)-N-(3-phenylpropan-1-ol-2-yl)benzamide (intermediate 43f) was suspended in 50 ml of tetrahydrofuran and mixed with 0.25 ml (2.7 mmol) of acetic anhydride at 100° C. The mixture was stirred for 16 h. The reaction was then concentrated under reduced pressure and the residue was recrystallized from ethanol. 0.78 g (71%) of the product was obtained.

Annotated Actions	<p>ADD \$1\$ (1 g, 2.7 mmol) ; ADD \$2\$ (50 ml) ; ADD \$3\$ (0.25 ml, 2.7 mmol) ; STIR for 16 h at 100 °C ; CONCENTRATE ; RECRYSTALLIZE from ethanol (0.78 g) ; YIELD \$-1\$.</p>	Predicted Actions	<p>MAKESOLUTION with \$1\$ and \$2\$ (15 ml) ; ADD \$3\$ (0.27 ml, 2.9 mmol) ; STIR for 10 h at 100 °C ; WAIT for 10 h at 100 °C ; STIR for 10 h at 100 °C ; CONCENTRATE ; ADD ethyl acetate ; ADD water ; PHASESEPARATION ; COLLECTLAYER aqueous ; EXTRACT with ethyl acetate ; COLLECTLAYER organic ; WASH with brine ; DRYSOLUTION over magnesium sulfate ; FILTER keep filtrate ; CONCENTRATE ; YIELD \$-1\$ (144 mg, 75%).</p>
-------------------	------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	-------------------	-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

(a) Example 3.

Field	Value
Reactant	\$1\$: <chem>BrC1ccc2noc(-c3ccccc3)c2c1</chem> \$2\$: <chem>O</chem>
CATALYST	\$3\$: <chem>[Zn]</chem>
SOLVENT	\$4\$: <chem>CC(=O)O</chem>
PRODUCT	\$-1\$: <chem>Nc1ccc(Br)cc1C(=O)c1ccccc1</chem>
Source	5-Bromo-3-phenyl-2,1-benzisoxazole (7.5 g, 28.6 m mol), water (14.6 ml), and zinc dust (9.3 g, 143 m mol) were combined. Acetic acid (8.6 ml, 143 m mol) was added and the mixture was stirred and heated at 80° C. for 90 minutes. After cooling to room temperature, both the liquid and solid portion of the reaction were extracted with methylene chloride. The combined methylene chloride solutions was washed once with sodium hydroxide solution (10%) and several times with water. Drying (sodium sulfate) and evaporation yielded the desired product (7.42 g) of mp. 92°-102° C.

Annotated Actions	<p>ADD \$1\$ (7.5 g) ; ADD \$2\$ (14.6 ml) ; ADD \$3\$ (9.3 g) ; ADD \$4\$ (8.6 ml) ; STIR for 90 minutes at 80 °C ; EXTRACT with methylene chloride ; EXTRACT with CH₂ Cl₂ ; WASH with NaOH ; WASH with \$2\$; DRYSOLUTION over sodium sulfate ; FILTER keep filtrate ; YIELD \$-1\$ (7.42 g).</p>	Predicted Actions	<p>MAKESOLUTION with \$1\$ and \$2\$ (150 ml) ; ADD \$3\$ (24 g) ; ADD \$4\$ (60 ml) ; STIR for 120 minutes at 80 °C ; WAIT for 120 minutes at 80 °C ; FILTER ; ADD carbon ; COLLECTLAYER organic ; WASH with \$2\$ (200 ml) 3 x ; DRYSOLUTION over sodium sulfate ; FILTER keep filtrate ; CONCENTRATE ; RECRYSTALLIZE from 2-amino-5-bromo-benzophenone (20.7 g) ; YIELD \$-1\$ (20.7 g, 57.9%).</p>
-------------------	------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	-------------------	------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

(b) Example 4.

Table 16: Examples of inaccurate experimental procedure predictions. Green denotes error of repetition.

Field	Value
Reactant	\$1\$: COc1ccc(-c2cccc(CC(=O)O)c2)cc1C(C)=NO \$3\$: FC(F)(F)c1ccc(CBr)cc1 \$4\$: [H-].[Na+] \$5\$: Cl
SOLVENT	\$2\$: C1CCOC1
PRODUCT	\$-1\$: COc1ccc(-c2cccc(CC(=O)O)c2)cc1C(C)=NOCc1ccc(C(F)(F)F)cc1
Source	9 mg of [3'-(1-Hydroxyiminoethyl)-4'-methoxybiphenyl-3-yl]acetic acid was dissolved in 0.2 mg of tetrahydrofuran, then 25 mg of 4-trifluoromethylbenzyl bromide and 10 mg of sodium hydride were added thereto, and the mixture was stirred at room temperature for 3 days. The reaction mixture was neutralized with 1N hydrochloric acid and then extracted with ethyl acetate. The organic layer was concentrated, and the residue was purified by reverse-phase high performance liquid chromatography to give 5.21 mg of the title compound.
Annotation Actions	<p>ADD \$1\$ (9 mg) ; ADD \$2\$ (0.2 mg) ; ADD \$3\$ (25 mg) ; ADD \$4\$ (10 mg) ; STIR for 3 days at room temperature ; PH with \$5\$ (1 normal) to pH neutral at room temperature ; EXTRACT with ethyl acetate ; CONCENTRATE ; YIELD \$-1\$ (5.21 mg).</p>
Predicted Actions	<p>MAKESOLUTION with \$1\$ (25 mL) and \$2\$; ADD \$3\$ (292 mg) ; ADD \$4\$ (104 mg) ; STIR for overnight at room temperature ; PH with \$5\$ (1 normal) to pH neutral at room temperature ; COLLECTLAYER organic ; WASH with ethyl acetate ; COLLECTLAYER organic ; WASH with water ; WASH with brine ; DRYSOLUTION over magnesium sulfate ; FILTER keep filtrate ; CONCENTRATE ; YIELD \$-1\$ (204 mg).</p>

(a) Example 5.

Field	Value
REACTANT	\$3\$: Cc1ccc2c(N)ccc(O)c2n1 \$4\$: O=N[O-].[Na+] \$5\$: [N-]=[N+]=[N-].[Na+]
Solvent	\$1\$: Cl \$2\$: O
PRODUCT	\$-1\$: Cc1ccc2c(N=[N+]=[N-])ccc(O)c2n1
Source	5-Amino-8-hydroxy-2-methylquinoline (12; 723 mg, 4.2 mmol) was dissolved in a solution of concentrated hydrochloric acid (0.4 mL) and water (5 mL), cooled to -3° C. in a salt-ice bath, stirred for 10 min, then treated dropwise with a cold solution of sodium nitrite (0.50 g, 7.2 mmol) in water (5 mL). The mixture was stirred for 20 min, then treated dropwise with sodium azide (0.60 g, 9.2 mmol) in water (40 mL), stirred at 0° C. for a further 1.5 h, then allowed to warm to room temperature over 24 h in the dark. Isolation by extraction with diethyl ether gave a dark brown solid, which was recrystallized from light petroleum to yield 13 as light brown crystals (554 mg, 66%). 1H NMR (500 MHz, CD3OD) δ ppm: 3.31 (s, 3H), 7.06 (d, J=8.5 Hz, 1H), 7.15 (d, J=8.5 Hz, 1H), 7.36 (d, J=8.5 Hz, 1H), 8.21 (d, J=8.5 Hz, 1H). 13C NMR (125 MHz, CD3OD) δ ppm: 24.9, 111.8, 115.6, 121.6, 123.7, 127.9, 132.6, 139.6, 151.1, 159.7. HRMS (ESI): calcd for C10H8N4O [M+H]+: 201.0771, found: 201.0773.
Annotation Actions	<p>MAKESOLUTION with \$1\$ (0.4 mL) and \$2\$ (5 mL) ; ADD \$3\$ (723 mg, 4.2 mmol) ; SETTEMPERATURE -3 °C ; STIR for 10 min ; MAKESOLUTION with \$4\$ (0.50 g, 7.2 mmol) and \$2\$ (5 mL) ; STIR for 20 min ; MAKESOLUTION with \$5\$ (0.60 g, 9.2 mmol) and \$2\$ (40 mL) ; STIR for 1.5 h at 0 °C ; RECRYSTALLIZE from petroleum ; YIELD \$-1\$ (554 mg, 66%).</p>
Predicted Actions	<p>MAKESOLUTION with \$1\$ (3 ml) and \$2\$ (1 ml) ; ADD \$3\$ (210 mg, 1.24 mmol) ; STIR for 1 h at 0 °C ; MAKESOLUTION with \$4\$ (80 mg, 1.24 mmol) and \$2\$ (1 ml) ; STIR for 1 h at 0 °C ; MAKESOLUTION with \$5\$ (88 mg, 1.24 mmol) and \$2\$ (1 ml) ; STIR for 2 h at 0 °C ; RECRYSTALLIZE ; YIELD \$-1\$ (120 mg, 47%).</p>

(b) Example 6.

Table 17: Examples of experimental procedure predictions that are different from the annotation but might be viable.

Field	Value
Reactant	\$1\$: <chem>CNC(=O)c1cn(CCCCc2ccc(N)nn2)nn1</chem> \$2\$: <chem>O=C(O)Cc1cc(Br)ccn1</chem> \$4\$: <chem>CCCP1(=O)OP(=O)(CCC)OP(=O)(CCC)O1</chem> \$6\$: <chem>CCN(C(C)C)C(C)C</chem>
Solvent	\$3\$: <chem>CN(C)C=O</chem> \$5\$: <chem>CCOC(C)=O</chem> \$7\$: <chem>O</chem>
Product	\$-1\$: <chem>CNC(=O)c1cn(CCCCc2ccc(NC(=O)Cc3cc(Br)ccn3)nn2)nn1</chem>
Source	To a solution of 1-(4-(6-aminopyridazin-3-yl)butyl)-N-methyl-1H-1,2,3-triazole-4-carboxamide (2.0 g, 7.2 mmol) and 2-(4-bromopyridin-2-yl)acetic acid (1.6 g, 7.2 mmol) in DMF (5 ml) were added T3P (50 wt. % in EtOAc, 9.2 g, 14 mmol) and DIEA (2.7 g, 22 mmol). The mixture was stirred at RT for 3 h, then water (50 mL) was added, the mixture was stirred for 30 min, and precipitate was isolated by filtration to give the title compound as a yellow solid (1 g, 40%). MS (ES+) C19H21BrN8O2 requires: 473. found: 474 [M+H]+.
Annotated Actions	<div style="display: flex; justify-content: space-between;"> <div style="width: 45%;"> <p>MAKESOLUTION with \$1\$ (2.0 g, 7.2 mmol) and \$2\$ (1.6 g, 7.2 mmol) and \$3\$ (5 ml) ;</p> <p>MAKESOLUTION with \$4\$ and \$5\$ (9.2 g, 14 mmol) and \$6\$ (2.7 g, 22 mmol) ;</p> <p>STIR for 3 h at RT ;</p> <p>ADD \$7\$ (50 mL) ;</p> <p>STIR for 30 min ;</p> <p>YIELD \$-1\$ (1 g, 40%).</p> </div> <div style="width: 45%; border-left: 1px solid black; padding-left: 10px;"> <p>Predicted Actions</p> <p>MAKESOLUTION with \$1\$ (30 mg, 0.109 mmol) and \$2\$ (23 mg, 0.109 mmol) and \$3\$ (0.5 ml) ;</p> <p>ADD \$4\$ (67 l, 0.109 mmol, 0.109 mmol) ;</p> <p>STIR for 3 h at room temperature ;</p> <p>ADD \$5\$ (5 ml) ;</p> <p>ADD \$6\$ (67 l, 0.38 mmol, 0.38 mmol) ;</p> <p>STIR for 30 min at room temperature ;</p> <p>ADD \$7\$ (0.5 ml) ;</p> <p>CONCENTRATE ;</p> <p>YIELD \$-1\$ (36 mg, 70% yield).</p> </div> </div>

Table 18: Examples of experimental procedure predictions that are different from the annotation but might be viable. Example 7.