

# TEXT-AUGMENTED MULTIMODAL LLMs FOR CHEMICAL REACTION CONDITION RECOMMENDATION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

High-throughput reaction condition (RC) screening is fundamental to chemical synthesis. However, current RC screening suffers from laborious and costly trial-and-error workflows. Traditional computer-aided synthesis planning (CASP) tools fail to find suitable RCs due to data sparsity and inadequate reaction representations. Nowadays, large language models (LLMs) are capable of tackling chemistry-related problems, such as molecule design, and chemical logic Q&A tasks. However, LLMs have not yet achieved accurate predictions of chemical reaction conditions. **Here, we present Chemma-RC, a text-augmented multimodal LLM that responds to task-specific questions by generating answers about reaction conditions. It learns a unified reaction representation via modality alignment from a corpus of reactions and question prompts, molecular structures in SMILES format, and graphical representations of chemical reactions.** We construct a 1.2 million pair-wised Q&A instruction dataset to train Chemma-RC and design a projection module for modality alignment. Our experimental results demonstrate that Chemma-RC achieves state-of-the-art performance on two open benchmark datasets and exhibits strong generalization capabilities on out-of-domain (OOD) and High-Throughput Experimentation (HTE) datasets. Chemma-RC has the potential to accelerate high-throughput condition screening in chemical synthesis.

## 1 INTRODUCTION

Chemical synthesis is a crucial step for the discovery of transformative molecules in multiple fields, including drug design, materials, renewable energy, etc. In chemical synthesis, reaction conditions are usually optimized to maximize the yield of each target molecule or minimize the cost of the corresponding process (Shields et al., 2021; Taylor et al., 2023). Despite significant advancements in chemical synthesis over the past few decades, discovering suitable reaction conditions from the extensive substrates combined with high-dimensional conditions renders exhaustive experimental impractical (Angello et al., 2022). Chemists have focused on building reliable and convenient computer-aided synthesis planning (CASP) tools to facilitate chemical synthesis (Corey & Wipke, 1969; Mikulak-Klucznik et al., 2020; Schwaller et al., 2021). However, few efforts have been made to solve the problem of reaction condition screening due to the low sparsity of chemical data, and the lack of effective reaction representation (Mehr et al., 2020; Rohrbach et al., 2022). *In summary, to realize efficient synthesis in chemistry, there is an urgent need to realize high-efficiency reaction condition recommendations.*

There are various types of data in the field of chemistry, including simplified molecular-input line-entry system (SMILES) (Weininger et al., 1989), graphs, and textual corpus of reaction (Schlichtkrull et al., 2018), which encompasses the descriptions of reaction processes, and reaction mechanisms. Traditional methods tackling the reaction condition recommendation (RCR) task typically rely on sequence-based SMILES data for end-to-end training (Gao et al., 2018; Schwaller et al., 2019; Andronov et al., 2023). However, training exclusively on sequence-based SMILES representations may hinder the model’s ability to capture the difference between similar reactions, as the feature distances encoded by transformers may be too close in the representation space. The capability to encode different reactions is critical for prediction, as even minor variations in a substrate’s functional group can result in fundamentally different reaction conditions. Therefore, it is necessary to incorporate additional information into reaction representations for RCR tasks. Given that the tex-

054 tual corpus contains chemical knowledge, which is invaluable for a comprehensive understanding of  
055 reactions, we aim to leverage cross-modality data to predict reaction conditions precisely.

056  
057 Nowadays, the emergence of generative large language models (LLMs), typified by GPT-4, has  
058 sparked significant interest in the field of AI for chemistry (Baum et al., 2021; Achiam et al., 2023;  
059 Boiko et al., 2023; Guo et al., 2023; M. Bran et al., 2024). Large multimodal models (LMMs) have  
060 demonstrated remarkable predictive capabilities in integrating modalities such as vision, text, and  
061 speech (Li et al., 2023; Zhu et al., 2024; Liu et al., 2024a). Therefore, we hypothesize that LMMs  
062 endowed with LLMs’ foundational capabilities in chemistry can deal with various modalities of  
063 chemical data, thereby enhancing the predictive performance in chemical tasks. However, it presents  
064 a significant challenge in designing modules to integrate various modalities effectively. Hence, *it is*  
065 *imperative to develop an effective prediction model that can incorporate different chemical data into*  
066 *LLMs to achieve a more comprehensive understanding of reaction processes, facilitating the task of*  
*chemical reaction condition recommendation.*

067 In view that molecules can be expressed as sequences, and reactions are described as natural lan-  
068 guage, e.g. text corpus, LMMs can be a potential solution due to the following advantages: (i) foun-  
069 dational LLMs can learn relationships between molecules in reactions, thereby acquiring chemical  
070 knowledge akin to the learning process of chemists (Achiam et al., 2023); (ii) via learning the joint  
071 representation of chemical reactions from different modalities, including graphs, SMILES, and cor-  
072 pus, LLMs might be empowered to understand the mechanism of reactions, which facilitates the task  
073 of RCR. To this end, we fine-tune general-purpose LLMs with domain-specific reaction data. Specif-  
074 ically, we present Chemma-RC, a multimodal LLM that jointly learns from the SMILES, graphs,  
075 and textual corpus of reactions. The contributions of this work can be summarized as follows:

- 076 1. We propose a multimodal LLM, a.k.a. Chemma-RC, to jointly learn representation from  
077 SMILES, graphs, and textual corpus of reactions for condition recommendation tasks. We  
078 further develop two distinct types of condition prediction modules, a classification mod-  
079 ule, and a generation module for Chemma-RC to enhance its compatibility with different  
080 reaction condition combinations.
- 081 2. We design text-augmented instruction prompts to construct a 1.2 million pair-wised Q&A  
082 dataset for training. We propose the Perceiver (Jaegle et al., 2021) module for modality  
083 alignment, which utilizes latent queries to align graphs and SMILES tokens with text-  
084 related tokens.
- 085 3. Through experimental validation on benchmark datasets, Chemma-RC achieves compet-  
086 itive results comparable to state-of-the-art models. Furthermore, Chemma-RC exhibits  
087 strong generalization capabilities on out-of-domain (OOD) and high-throughput experi-  
088 mentation (HTE) datasets.

## 091 2 RELATED WORK

092  
093 In chemical synthesis, reaction conditions are usually developed and optimized to maximize the  
094 yield of each target molecule or minimize the cost of the corresponding process (Shields et al.,  
095 2021; Taylor et al., 2023). High-throughput reaction condition (RC) screening, as an important tool  
096 in synthesizing molecules, exerts an important influence on chemical synthesis. However, discover-  
097 ing suitable reaction conditions from the extensive matrix of substrates combined with the high-  
098 dimensional reaction conditions renders exhaustive experimental impractical. (Angello et al., 2022).  
099 For decades, chemists have focused on building reliable and convenient computer-aided synthesis  
100 planning (CASP) tools to facilitate chemical synthesis (Corey & Wipke, 1969; Mikulak-Klucznik  
101 et al., 2020). For instance, Coley et al. built a multiway classification model based on a two-step  
102 graph convolutional network (GCN) for the reaction prediction task (Coley et al., 2017; 2019). Due  
103 to the effectiveness of a simplified molecular-input line-entry system (SMILES) (Weininger et al.,  
104 1989), as strings of a context-free, Nam et al. proposed the first sequence-to-sequence model for  
105 forward prediction using the SMILES representations of molecules (Nam & Kim, 2016). Inspired  
106 by attention-based transformer model (Vaswani et al., 2017), Schwaller et al. proposed molecular  
107 transformers (Schwaller et al., 2019; Ding et al., 2024), which were applied in forward prediction  
and reaction condition recommendation (RCR) tasks (Schwaller et al., 2019; Andronov et al., 2023).

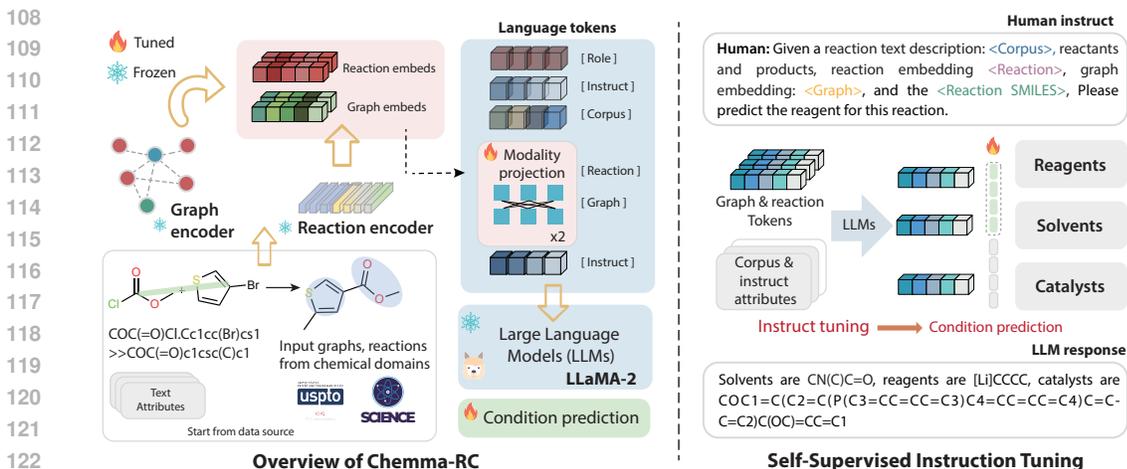


Figure 1: **Architecture of Chemma-RC.** Chemma-RC processes task-specific questions constructed by text-augmented multimodal instruction prompts and generates answers. Specifically, it takes three modalities of data as inputs: text (a textual corpus of reactions and question prompts), molecular SMILES, and reaction graphs. Two distinct types of prediction modules, a classification module, and a generation module are proposed to predict chemical reaction conditions.

Chemical reaction condition recommendation tasks aim to recommend catalysts, reagents, solvents, or other conditions for a specific reaction. The exploration of a suitable condition is crucial for the realization of CASP, as it dictates the expected outcomes, including reaction yields and rates (Schnitzer et al., 2024). Gao et al. developed a neural network model to predict the chemical context as well as the temperature for any particular organic reaction (Gao et al., 2018); Maser et al. proposed a machine-learned ranking model to predict the set of conditions used in a reaction as a binary vector (Maser et al., 2021); Wang et al. proposed Parrot, a powerful and interpretable transformer-based model for the prediction of reaction condition (Wang et al., 2023a); In the meantime, in order to enhance the representation of reactions, Qian et al. (Qian et al., 2023) designed TextReact, which introduced relevant corpus retrieved from literature to enhance the molecular representation of the reaction based on SMILES. Nevertheless, these methods rely on manual feature selection by experts’ knowledge and lack a general prediction model with powerful reaction representation.

Nowadays, the emergence of generative pre-trained transformer-based large language models (LLMs), typified by GPT-4, has triggered keen interest in leveraging such techniques to tackle chemistry challenges (Baum et al., 2021; Achiam et al., 2023). Several works focus on chemical agents for the exploration of chemical conditions. Boiko et al. (Boiko et al., 2023) proposed a GPT-4 driven scientific agent system to plan and perform complex experiments, which accelerates reaction condition screening and experimental automation in chemistry; Bran et al. developed ChemCrow, which augmented LLMs with chem-expert-designed tools (M. Bran et al., 2024); However, for tasks demanding a precise understanding of molecular SMILES representation, such as reaction prediction, and retrosynthesis, LLMs exhibited a less competitive performance than traditional machine learning baselines (Guo et al., 2023). Partially, the reason is that, without an in-depth understanding of the SMILES strings, and the reaction process that transforms reactants into products, it will be difficult for LLMs to generate accurate responses.

Besides SMILES strings, there are various types of data such as molecule graphs and the reactions’ external textual corpus in the chemistry synthesis field. By synergizing the strengths of multiple modalities, large multimodal models (LMMs) can achieve higher accuracy, and perform more effectively in a wide range of applications (Edwards et al., 2022; Li et al., 2023; Zhu et al., 2024; Liu et al., 2024a; Li et al., 2024; Liu et al., 2024b).

### 3 METHODS

#### 3.1 PROBLEM SETUP

For a task of reaction condition recommendation, we define the  $X$  as the input for the chemical reaction  $R$ ,  $T$  as the reaction corpus,  $G$  as the graph representations of reactions, and the output

162  $Y$  as a list of reaction conditions including the catalyst, solvent, and reagent. Thus, we define the  
 163 prediction model  $\mathcal{F}$ , i.e.,  $Y = \mathcal{F}(X, G, T)$ .

164  
 165 In this paper, we incorporate three types of data for the training of model  $\mathcal{F}$ :

- 166 1. **SMILES of a reaction**  $X$ : each example in the training set is presented by chemical  
 167 SMILES, i.e., "CC(C)O.O=C(n1ccnc1)ncnc1 >> CC(C)OC(=O)n1ccnc1".
- 168 2. **Graphs of a reaction**  $G$ : each SMILES representation of the reactants and the product is  
 169 encoded using a graph neural network (GNN). All compounds are integrated to generate a  
 170 comprehensive reaction representation.
- 171 3. **An unlabeled reaction corpus**: a paragraph describing a chemical reaction, e.g., "To a  
 172 solution of CDI (2 g, 12.33 mmol), in DCM (25 mL) was added isopropyl alcohol (0.95  
 173 mL, 12.33 mmol) at 0° C."

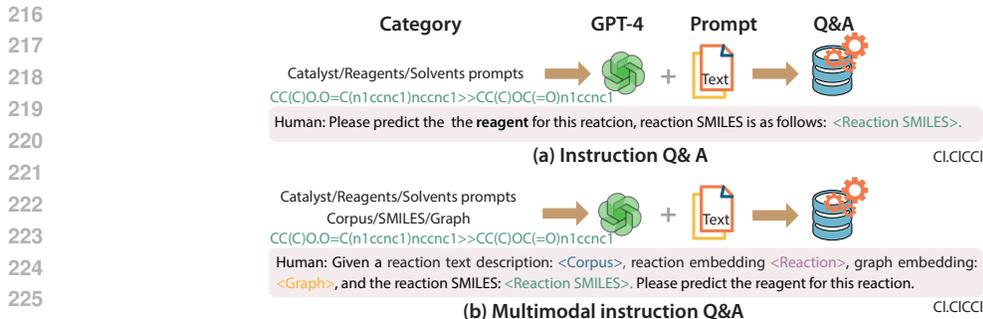
### 175 3.2 MODEL STRUCTURE

177 Here we first introduce the **Chemma-RC**, a multimodal LLM designed for reaction condition rec-  
 178 ommendation (RCR). An overview of Chemma-RC is illustrated in Figure. 1. Chemma-RC responds  
 179 to task-specific questions constructed by instruction prompts such as "*Please predict the reagent for*  
 180 *this reaction.*", and generates answers about reaction conditions. The Chemma-RC model accepts  
 181 three different data modalities as inputs. This includes text from a corpus of reactions and question  
 182 prompts, molecular structures in SMILES format, and graphical representations of chemical reac-  
 183 tions. We employ both transformer-based reaction encoder and GCN models to learn reaction rep-  
 184 resentations from SMILES and graph structure jointly. Subsequently, the modality projection trans-  
 185 forms the graph and SMILES embeddings into language tokens compatible with LLM space. These  
 186 learnable tokens, defined as graph and reaction tokens, along with tokens of instruction prompts, are  
 187 then input into the LLM to predict chemical reaction conditions. Note that, we develop two distinct  
 188 types of condition prediction modules, a **classification** and a **generation** prediction module to en-  
 189 hance its compatibility with different chemical reaction conditions. **On the one hand, the reason for**  
 190 **performing classification tasks is to select the most suitable reaction conditions from commercially**  
 191 **available libraries, as it is common practice to prioritize purchasable molecules. On the other hand,**  
 192 **the generation module can assist in designing novel molecules, which can be obtained by synthesis**  
 193 **experiments conducted. Therefore, we define two distinct tasks including classification and gener-**  
 194 **ation modules to address these objectives. Furthermore, existing baseline methods treat RCR as a**  
 195 **classification task for the USPTO-Condition datasets. To ensure a fair comparison, we conduct a**  
 196 **classification module for prediction and evaluation.**

#### 197 3.2.1 CONSTRUCTION OF TEXT-AUGMENTED INSTRUCTION PROMPTS

198 Instruction prompt datasets refer to format structured or unstructured data as natural language in-  
 199 structions so that LLMs can respond properly (Reynolds & McDonell, 2021; Wang et al., 2023b).  
 200 Compared to creating language instruction datasets for fine-tuning LLMs, constructing multimodal  
 201 instruction datasets requires a thorough understanding of domain-specific tasks. Recent advance-  
 202 ments indicate that the other data modalities, such as images, and graphs, can be transformed as the  
 203 prefix of prompts thereby facilitating effective reasoning based on inputs (Tsimpoukelli et al., 2021;  
 204 Zhu et al., 2024; Liu et al., 2024a).

205 Toward reaction condition recommendation task in chemical synthesis, we design a tailored instruc-  
 206 tion prompt system for better cross-modality alignment and instruction tuning (Figure. 2). Compared  
 207 to instruction prompts for natural language instruction tuning (Figure. 2(a)), we introduce augmented  
 208 text tokens and multimodal tokens into instruction prompts (Figure. 2(b)). **To be specific, given a**  
 209 **reaction, we retrieve a relevant corpus—a paragraph containing contextual information that closely**  
 210 **resembles the reaction—and populate the <Corpus>placeholder with this data. Next, the reaction**  
 211 **is converted into its corresponding SMILES representation, which is then inserted into the <Re-**  
 212 **action SMILES>placeholder. Finally, we introduce two additional placeholders, <Reaction>and**  
 213 **<Graph>, designed to accommodate the reaction and graph-based representations, respectively. In**  
 214 **instruction fine-tuning, all reaction embedding representations are extracted by reaction encoders.**  
 215 **Via the modality alignment module, all embeddings are inserted into token placeholders to align**  
**text-related tokens in language space. We also give pseudo-code as follows to explain this integra-**  
**tion process, which can be found in the Appendix. C Algorithm 1.**



227 Figure 2: Instruction of text-augmented prompts. (a) Traditional instruction prompts for natural language  
228 instruction tuning; (b) Our proposed text-augmented multimodal instruction Q&A prompts.

### 229 3.2.2 ENCODER AND DECODER

230  
231 Given a reaction  $R$ , we adapt a pioneering transformer-based encoder, Parrot (Wang et al., 2023a)  
232 to produce the reaction embeddings  $\mathbf{X}_R \in \mathbb{R}^{N \times C}$ . Here,  $N$  and  $C$  indicate the length of text tokens  
233 and embedding channels, respectively. During training, the encoder computes a contextual vector  
234 representation of the reactions by performing self-attention on the masked canonicalized SMILES  
235 string of molecules. We denote reaction embeddings as SMILES embedding in the following section.  
236

237 In the meantime, we leverage a GNN (Schlichtkrull et al., 2018) to model the relationship between  
238 atoms in molecules. We denote directed and labeled multi-graphs as  $G = (\mathcal{V}, \mathcal{E}, \mathcal{R})$  with nodes  
239 (atom entities),  $v_i \in \mathcal{V}$  and labeled edges (atom relations)  $(v_i, r, v_j) \in \mathcal{E}$ , where  $r \in \mathcal{R}$  is a  
240 relation type. GNN can be understood as special cases of a simple differentiable message-passing  
241 framework:

$$242 h_i^{(l+1)} = \sigma \left( \sum_{m \in \mathcal{M}_i} g_m(h_i^{(l)}, h_j^{(l)}) \right) \quad (1)$$

243  
244 where  $h_i^{(l)} \in \mathbb{R}^{d^{(l)}}$  is the hidden state of node  $v_i$  in the  $l$ -th layer of the neural network, with  $d^{(l)}$   
245 being the dimensionality of this layer’s representations. Incoming messages of the form  $g_m(\cdot, \cdot)$   
246 are accumulated and passed through an element-wise activation function  $\sigma(\cdot)$ , such as the  $\text{ReLU}(\cdot) =$   
247  $\max(0, \cdot)$ ,  $\mathcal{M}_i$  denotes the set of incoming messages for node  $v_i$  and is often chosen to be identical to  
248 the set of incoming edges.  $g_m(\cdot, \cdot)$  is typically chosen to be a (message-specific) neural network-like  
249 function or simply a linear transformation  $g_m(h_i, h_j) = Wh_j$  with a weight matrix  $W$ . Motivated  
250 by this architecture, GCNN (Schlichtkrull et al., 2018) proposed a refined propagation model for the  
251 forward-pass update of an entity or node:  
252

$$253 h_i^{(l+1)} = \sigma \left( \sum_{r \in \mathcal{R}} \sum_{j \in \mathcal{N}_i^r} \frac{1}{c_{i,r}} W_r^{(l)} h_j^{(l)} + W_0^{(l)} h_i^{(l)} \right) \quad (2)$$

254  
255 where  $\mathcal{N}_i^r$  denotes the set of neighbor indices of node  $i$  under relation  $r \in \mathcal{R}$ .  $c_{i,r}$  is a problem-  
256 specific normalization constant that can either be learned or chosen in advance (such as  $c_{i,r} = |\mathcal{N}_i^r|$ ).  
257

258 We develop two distinct types of prediction modules, a classification module and a generation mod-  
259 ule for Chemma-RC to enhance its compatibility with different chemical reaction conditions. Pre-  
260 diction modules are used to generate probability distributions over potential tokens, and we define  
261 two types of loss for this:

$$262 \text{Prediction : } \begin{cases} (1) X, G, T \xrightarrow{\text{classifier}} (c_i, \hat{c}_i) : \mathcal{L} = \sum_{i \in I} \text{CrossEntropyLoss}(c_i, \hat{c}_i) \\ (2) X, G, T \xrightarrow{\text{generate}} (C, \hat{C}) : \mathcal{L} = - \sum_{l=1}^L \sum_{v=1}^V y_l^v \log P_\theta(y_l^v | y_{<l}, (x, g, t)) \end{cases} \quad (3)$$

263  
264 where *classifier* refers to classification head,  $I$  is the chemical context condition number,  $c_i$  is  
265 the predicted label of the  $i$ -th condition,  $\hat{c}_i$  is the ground truth label of the  $i$ -th condition; *generate*  
266  
267  
268  
269

refers to generation head,  $C$  and  $\widehat{C}$  are the combination of predicted and the ground truth conditions, respectively.  $L$  is the sequence length,  $V$  is the vocabulary size.  $y_l$  is the one-hot encoded target token at position  $l$ ,  $y_l^v$  is the  $v$ -th element of the one-hot encoded target token at position  $l$ ;  $y_{<l}$  represents all previous tokens before position  $l$ ;  $(x, g, t)$  is the input context tokens representing SMILES, graphs, and corpus.

### 3.2.3 MODALITY PROJECTION

For the reaction condition recommendation task, the representation of the reaction is extracted by encoders (see in section 3.2.2), and LLMs tokenize the text representation. However, fusing two types of representation introduces inductive biases issues (Baltrušaitis et al., 2018; Jaegle et al., 2021). To effectively fuse representations from multiple modalities, we propose the Perceiver (Jaegle et al., 2021) module for modality projection, seen ‘modality projection’ in Figure 1. This module employs latent tokens to align graphs and SMILES embeddings with text-related tokens extracted from question prompts and a text-augmented corpus. During training, we employ two transformer-based Perceivers as projectors. Although these modules share an identical model architecture, they are distinguished by their unique weights. Consequently, learnable tokens contain highlighted reaction cues that are most related to the text tokens. We show the pseudo-code for modality projection in Appendix. C.

## 4 EXPERIMENTS AND RESULTS

### 4.1 DATA

We curate two large datasets, named USPTO-Condition and USPTO\_500MT\_Condition for evaluation. Data volumes are presented in Table. 1. The visualization of data distribution is depicted in Figure. 4. As depicted in Table. 8, for the USPTO-Condition dataset, five conditions categories are separated by commas in order. For the USPTO\_500MT\_Condition dataset, all conditions are combined by dot as strings. The detailed data description can be seen in Appendix. B.

Table 1: Data description of USPTO-Condition and USPTO\_500MT\_Condition.

Dataset	Sample of conditions	Prediction type	Training set
USPTO-Condition	[Zn],C1CCOC1,O,CO,[Cl-].[NH4+]	classification	546,728
USPTO_500MT_Condition	CO.[Na+].CC(=O)O.[BH3-]C#N	generation	88,410

### 4.2 EXPERIMENT SETUP

In our work, the reaction encoder is implemented based on Wang et al. (Wang et al., 2023a). A pre-trained graph model proposed by (Schlichtkrull et al., 2018) encodes the molecules in the reaction. We utilize LLaMA-2 (Touvron et al., 2023) as a text decoder. Each reaction has the corresponding corpus, a paragraph describing a chemical reaction with an average length of 190 tokens. During the training process, we fix the weight parameters of GCN, reaction encoder, and LLaMA-2. The modality projection and condition prediction layer is trainable. **The trainable parameters constitute approximately 0.3 billion out of the total 7 billion parameters. The training process is conducted with a batch size of 16 for fewer than 6 epochs over 48 hours, utilizing a GPU configuration of 2x48 GB NVIDIA A6000 GPUs. The inference process is highly efficient and can be performed using a single 80 GB NVIDIA A800 GPU. The detailed training setting can be seen in Appendix. A.**

### 4.3 PERFORMANCE COMPARISON

We assess the performance of our proposed Chemma-RC for reaction condition recommendation. The top- $N$  accuracy of condition recommendation on the combined test datasets of USPTO-Condition and USPTO\_500MT\_Condition are presented in Table. 2 and Table. 3, respectively. Compared methods include rxnfp LSTM (Gao et al., 2018), Reaction GCNN (Maser et al., 2021), TextReact (Qian et al., 2023), and Reagent Transformer (Andronov et al., 2023). The details of the baselines are present in Appendix. D.

Table 2: Results of reaction condition recommendation on USPTO-Condition dataset. The best performance is in **bold**.

Model	Top-k Accuracy (%)														
	Catalyst			Solvent 1			Solvent 2			Reagent 1			Reagent 2		
	1	3	5	1	3	5	1	3	5	1	3	5	1	3	5
rxnfp LSTM	91.6	94.1	95.2	48.3	64.4	70.2	81.4	83.4	84.6	48.2	64.4	70.8	76.5	84.1	86.4
Parrot	89.9	96.4	97.7	35.2	60.9	72.2	81.2	93.7	96.7	40.4	62.3	71.7	<b>80.6</b>	90.6	93.6
TextReact <sub>s</sub>	92.1	98.0	99.1	51.4	68.5	79.3	81.6	93.4	96.9	51.1	69.6	79.1	77.9	91.1	94.9
Chemma-RC <sub>s</sub>	<b>92.8</b>	<b>98.6</b>	<b>99.3</b>	<b>54.7</b>	<b>76.5</b>	<b>84.9</b>	<b>81.9</b>	<b>94.8</b>	<b>97.6</b>	<b>53.4</b>	<b>75.9</b>	<b>83.9</b>	78.6	<b>93.2</b>	<b>96.2</b>

For the USPTO-Condition dataset, we calculate top- $k$  accuracy with a strict matching policy. As depicted in Table. 2, TextReact<sub>s</sub> refers that we utilize *similar text* (Qian et al., 2023) paired with the corresponding reaction for training. To avoid label leak issues, we do not use *gold text* mentioned in his work for training or testing. Chemma-RC<sub>s</sub> refers that we use a similar corpus paired with each reaction as input to construct Q&A instruction datasets for training. Thanks to the work of Qian et al., we can retrieve the most similar corpus for each reaction from the literature or patents using their pre-trained model.

Table 3: Results of reaction condition recommendation on USPTO\_500MT\_Condition dataset. The best performance is in **bold**.

Model	Top-k Accuracy (%)			
	1	3	5	10
Reagent Transformer	17.5	27.5	31.6	35.6
Reaction GCNN	16.1	27.5	33.0	40.2
Parrot	13.8	25.3	31.4	37.9
nach0	13.1	-	-	-
<b>Chemma-RC</b>	<b>25.9</b>	<b>47.2</b>	<b>67.8</b>	<b>79.2</b>

tion, akin to the learning process of chemists (Achiam et al., 2023).

Unlike the USPTO-Condition dataset which includes three types of chemical condition data—catalysts, solvents, and reagents—the USPTO\_500MT\_Condition dataset categorizes all conditions as ‘reagents’. The performance of comparative methods on the USPTO\_500MT\_Condition dataset is shown in Table. 3. We have broadened several sets of baseline models to illustrate the feasibility of Chemma-RC, including nach0 (Livne et al., 2024), transformer-based models (Andronov et al., 2023), and other methods. The visualization of performance is shown in Appendix Figure. 6. We examine top-1, top-3, top-5, and top-10 predictive results. Notably, for USPTO\_500MT\_Condition datasets (Table. 3), we can see that Chemma-RC demonstrates the most favorable performance, where achieves 25.9% top-1 accuracy when compared with other baseline methods such as Reagent Transformer (17.5%), Reaction GCNN (16.1%), nach0 (13.1%). All SMILES conditions in the USPTO\_500MT\_Condition dataset are concatenated with dots, resulting in challenges due to the lengthy token sequences. However, Chemma-RC, pre-trained on a vast natural language corpus, effectively manages and accurately generates these long tokens.

## 4.4 ABLATION STUDY

### 4.4.1 MODEL STRUCTURE

In Chemma-RC, SMILES strings provide a textual representation of molecular structures, concisely encoding vital connectivity and stereochemistry details. Structural graphs of molecules offer a topological view of molecules in two-dimensional space, where atoms are nodes and bonds are edges. The textual corpus introduces a natural language context into the model to enhance the chemical interpretation capability of LLMs.

First, to examine the effect of different modalities on the performance of Chemma-RC, we evaluate the performance under the different combinations of mono-domain data including SMILES, graph, and corpus on the USPTO-Condition dataset. As indicated in Table. 4, from the results, we can see that different mono-domain data have different contributions for the entire performance. For the

Table 4: Performance evaluation of Chemma-RC under different combinations of mono-domain data on the USPTO-Condition Dataset.

SMILES	Graph	Corpus	Top- <i>k</i> Accuracy (%)														
			Catalyst			Solvent 1			Solvent 2			Reagent 1			Reagent 2		
			1	3	5	1	3	5	1	3	5	1	3	5	1	3	5
✓	✗	✗	90.3	97.5	98.7	37.1	64.5	75.7	80.8	92.9	96.8	37.1	63.5	74.7	73.7	89.9	94.1
✗	✓	✗	87.1	93.3	95.5	15.3	40.5	58.2	80.7	91.9	95.5	34.6	56.8	67.5	75.4	86.6	90.6
✗	✗	✓	87.1	87.4	87.8	14.1	26.1	44.9	80.7	88.1	92	26.0	32.1	37.3	75.1	76.6	77.9
✓	✗	✓	92.6	98.5	<b>99.3</b>	54.0	76.0	84.4	<b>81.8</b>	94.7	<b>97.6</b>	52.8	75.4	83.3	78.6	93.1	96.1
✓	✓	✗	91.3	98.1	99.1	42.1	68.8	79.4	80.1	93.5	97.1	45.2	70.4	79.9	76.7	91.4	95.1
✓	✓	✓	<b>92.7</b>	<b>98.6</b>	99.2	<b>54.6</b>	<b>76.4</b>	<b>84.8</b>	<b>81.8</b>	<b>94.8</b>	<b>97.6</b>	<b>53.4</b>	<b>75.8</b>	<b>83.9</b>	<b>78.7</b>	<b>93.2</b>	<b>96.2</b>

prediction of solvent 1, which is the most challenging task, the model enhanced with SMILES representation (first row) outperforms the models trained solely on graph-based features (second row) and corpus data (third row), achieving 21.8% and 23.0% higher top-1 accuracy, respectively. Subsequently, we investigate how chemical mono-domain data combination affects model performance compared to individual types of data (fourth row to sixth row). **By incorporating a corpus into the model already trained with SMILES representations, we achieve a 16.9% improvement in solvent 1 top-1 prediction accuracy.** Similarly, **integrating graph features into the SMILES-based model results in a 5.0% improvement in solvent 1 top-1 accuracy.** The effectiveness of incorporating additional corpus data and SMILES representations can be attributed to the LLM’s pre-training on extensive SMILES sequences and reaction data, which equips it with a more comprehensive understanding of chemical reactions and enhances its performance on RCR tasks. In a word, experimental results substantiate that integrating different modalities of chemical data including SMILES, graphs, and natural corpus, presents an effective representation of reactions, which is effective for RCR scenarios.

#### 4.4.2 DATA SPLIT STRATEGY

We include the other baseline methods for comparison on the USPTO-Condition dataset. We also evaluate Chemma-RC’s performance under different dataset splitting strategies, including random split (RS) and time-based split (TS), to further demonstrate its robustness across diverse conditions. A detailed introduction of each method and experiment settings are illustrated in the Appendix. D. TextReact (gr) refers to the TextReact model without retrieving gold texts for testing. From the results, we can see that the performance of other baseline models such as rxnfp LSTM (Gao et al., 2018), rxnfp retrieval, Transformer, and ChemBERTa (Chithrananda et al., 2020) shows moderate success. However, these models consistently deliver lower accuracy rates compared to TextReact (gr) and Chemma-RC. Chemma-RC significantly outperforms all baseline methods across both RS and TS settings. Notably, it achieves a Top-1 (RS) accuracy of 72.3%, which is substantially higher than the second-best approach, TextReact (gr), at 47.2%.

Table 5: Evaluation results for reaction condition recommendation (RCR). RS: random split; TS: time split. Scores are accuracy in %.

	RCR (RS)				RCR (TS)			
	Top-1	Top-3	Top-10	Top-15	Top-1	Top-3	Top-10	Top-15
rxnfp LSTM	20.5	30.7	41.7	45.3	15.2	26.2	40.7	45.4
rxnfp retrieval	27.2	37.5	47.9	51.1	7.8	15.2	27.3	31.5
Transformer	30.0	43.8	56.7	60.5	18.7	31.8	47.6	52.7
ChemBERTa	30.3	44.7	58.0	62.0	18.7	31.9	47.6	52.8
TextReact(gr)	47.2	59.9	65.0	71.4	36.3	50.4	56.2	63.8
Chemma-RC	<b>72.3</b>	<b>87.8</b>	<b>92.4</b>	<b>96.5</b>	<b>69.6</b>	<b>86.7</b>	<b>91.7</b>	<b>96.2</b>

#### 4.4.3 MODALITY PROJECTION

By leveraging the strengths of multiple modalities, multimodal LLMs can achieve higher accuracy in a wide range of applications. However, aligning representations among different modalities remains a challenging task. In our proposed Chemma-RC, we employ the Perceiver module (Jaegle et al.,

2021) to integrate molecular SMILES tokens and graphs tokens into text-related language space, where text tokens are augmented by the reaction corpus, as illustrated in Figure 1. This modality projection module maps the embeddings of reactions to a latent vector and enhances this representation using a Transformer tower. Consequently, learnable queries contain highlighted reaction contents that are most related to the text tokens. We compared three typical methods for modality projection, including Perceiver (Jaegle et al., 2021), Reprogramming (Jin et al., 2024), and MLP.

Table 6: Performance evaluation of Chemma-RC under different modality projections, the best performance are in bold.

Projection Layer	Top-k Accuracy (%)														
	Catalyst			Solvent 1			Solvent 2			Reagent 1			Reagent 2		
	1	3	5	1	3	5	1	3	5	1	3	5	1	3	5
MLP	90.9	97.8	98.9	51.1	73.3	82.2	81.1	93.9	97.1	47.4	71.0	79.9	77.0	91.7	95.2
Reprogramming	92.1	98.3	99.1	52.8	75.1	83.7	81.3	94.3	97.4	50.2	73.5	81.9	77.7	92.5	95.7
Perceiver	<b>92.7</b>	<b>98.6</b>	<b>99.2</b>	<b>54.6</b>	<b>76.4</b>	<b>84.8</b>	<b>81.8</b>	<b>94.8</b>	<b>97.6</b>	<b>53.4</b>	<b>75.8</b>	<b>83.9</b>	<b>78.7</b>	<b>93.2</b>	<b>96.2</b>

As depicted in Table. 6, the Perceiver module achieves significant gains in the prediction of all categories. Compared with Chemma-RC (with Reprogramming), Chemma-RC (with Perceiver) can be further enhanced and attains peak performance in all predicted categories with 7.1% significant gain. Specifically, For the solvent 1 prediction, a hard case, the Perceiver module stands out with a top-1 accuracy of 54.6%, significantly surpassing MLP (51.1%) and Reprogramming (52.8%). Its ability to consistently achieve high accuracy in both top-1 and top-k evaluations suggests a robust and versatile approach for reaction condition recommendation.

#### 4.5 TRANSFERABILITY EVALUATION ON HIGH-THROUGHPUT EXPERIMENTATION REACTION

Discovering effective reaction conditions precisely for high-throughput reaction condition screening is very important, as it has the potential to release chemists from laborious and costly trial-and-error workflows. Thus, we illustrate the transferability of our models through zero-shot evaluation on distinct high-throughput experimentation (HTE) datasets. We expect that Chemma-RC recommends conditions that yield high-product outputs. We select the Imidazole C-H arylation dataset extracted from the work proposed by Shields et al. in 2021 (Shields et al., 2021) for evaluation, where the substrate scope contains 8 imidazoles and 8 aryl bromides associated with conditions including ligands, bases, and solvents.

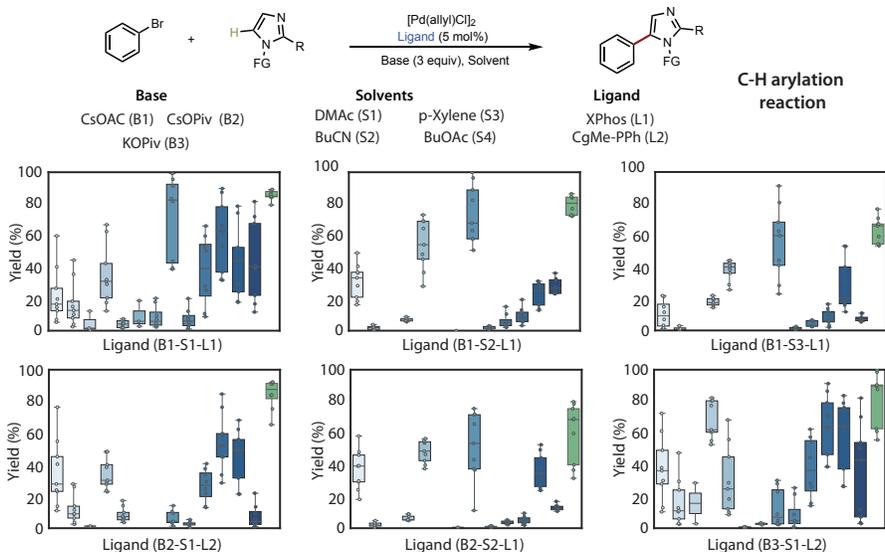


Figure 3: Boxplot of the performance for ligand recommendation on C-H arylation reaction.

Catalysts are vital compounds in chemical reactions, as they play a crucial role in determining both reactivity and yield. The catalyst used in imidazole C-H arylation comprises a metal (Pd) and ligands. Thus, we evaluate the performance of ligand recommendations. First, we ensure that reaction

486 data of imidazole C–H functionalization is excluded from the test set of the USPTO-Condition  
487 dataset to prevent data leakage issues. Chemma-RC recommends a ligand under a pre-defined  
488 solvent-base combination of conditions. As shown in Figure. 3, we randomly select six cases for  
489 performance evaluation. The referenced bases, solvents, and ligands can be found in the reaction  
490 formula, which has been annotated by ‘B’, ‘S’, ‘L’. For example, in Figure. 3, under the combination  
491 of CsOAc and DMAc, Chemma-RC identifies the XPhos ligand, which results in a higher yield.

492 For recommended results (Figure. 10, Figure. 11) we can observe that, for **15** of the 16 base-solvent  
493 combinations, the recommended ligand performs best in terms of the median value of reaction yields,  
494 suggesting that Chemma-RC can recommend ligands with higher yields.

495 Moreover, we can conclude that the capability of Chemma-RC to recommend suitable conditions  
496 for chemical reactions has the potential to accelerate high-throughput reaction condition screening  
497 in the future.

## 499 5 CONCLUSION AND LIMITATIONS

501 **Conclusions** In this paper, we present a multimodal LLM, a.k.a. Chemma-RC for chemical reac-  
502 tion condition recommendation. Trained with 1.2 million pair-wised Q&A instruction datasets that  
503 integrate with multimodal reaction representations and corpus in natural language, Chemma-RC  
504 effectively answers questions regarding reaction conditions through either a classification head or  
505 sequence generation.

506 **Limitations** Further, we will focus on how the token length of each modality improves its perfor-  
507 mance across various chemical reaction tasks in future work.

## 510 6 REPRODUCIBILITY STATEMENT

511 To ensure the reproducibility of our work, we have used datasets which have been pub-  
512 lished in (Wang et al., 2023a; Lu & Zhang, 2022), and the data links are as follows:  
513 USPTO\_500MT\_Condition and USPTO-Condition. Additionally, we commit to releasing the full  
514 implementation of our code, including model architectures, training pipelines, and evaluation scripts,  
515 upon acceptance and publication of this paper. Detailed instructions and necessary dependencies are  
516 provided in the Appendix to facilitate easy reproduction of our results.

## 519 REFERENCES

- 520 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-  
521 man, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. GPT-4 Techni-  
522 cal Report. *arXiv preprint arXiv:2303.08774*, 2023.
- 523
- 524 Mikhail Andronov, Varvara Voinarovska, Natalia Andronova, Michael Wand, Djork-Arné Clevert,  
525 and Jürgen Schmidhuber. Reagent prediction with a molecular transformer improves reaction data  
526 quality. *Chemical Science*, 14(12):3235–3246, 2023.
- 527
- 528 Nicholas H Angello, Vandana Rathore, Wiktor Beker, Agnieszka Wołos, Edward R Jira, Rafał  
529 Roszak, Tony C Wu, Charles M Schroeder, Alán Aspuru-Guzik, Bartosz A Grzybowski, et al.  
530 Closed-loop optimization of general reaction conditions for heteroaryl Suzuki-Miyaura coupling.  
531 *Science*, 378(6618):399–405, 2022.
- 532 Tadas Baltrušaitis, Chaitanya Ahuja, and Louis-Philippe Morency. Multimodal Machine Learning:  
533 A Survey and Taxonomy. *IEEE transactions on pattern analysis and machine intelligence*, 41(2):  
534 423–443, 2018.
- 535 Zachary J Baum, Xiang Yu, Philippe Y Ayala, Yanan Zhao, Steven P Watkins, and Qiongqiong  
536 Zhou. Artificial Intelligence in Chemistry: Current Trends and Future Directions. *Journal of*  
537 *Chemical Information and Modeling*, 61(7):3197–3212, 2021.
- 538
- 539 Daniil A Boiko, Robert MacKnight, Ben Kline, and Gabe Gomes. Autonomous chemical research  
with large language models. *Nature*, 624(7992):570–578, 2023.

- 540 Seyone Chithrananda, Gabriel Grand, and Bharath Ramsundar. Chemberta: large-scale self-  
541 supervised pretraining for molecular property prediction. *arXiv preprint arXiv:2010.09885*, 2020.  
542
- 543 Connor W Coley, Regina Barzilay, Tommi S Jaakkola, William H Green, and Klavs F Jensen. Pre-  
544 diction of Organic Reaction Outcomes Using Machine Learning. *ACS central science*, 3(5):434–  
545 443, 2017.
- 546 Connor W Coley, Wengong Jin, Luke Rogers, Timothy F Jamison, Tommi S Jaakkola, William H  
547 Green, Regina Barzilay, and Klavs F Jensen. A graph-convolutional neural network model for the  
548 prediction of chemical reactivity. *Chemical science*, 10(2):370–377, 2019.  
549
- 550 Elias James Corey and W Todd Wipke. Computer-Assisted Design of Complex Organic Syntheses:  
551 Pathways for molecular synthesis can be devised with a computer and equipment for graphical  
552 communication. *Science*, 166(3902):178–192, 1969.  
553
- 554 Yuheng Ding, Bo Qiang, Qixuan Chen, Yiqiao Liu, Liangren Zhang, and Zhenming Liu. Exploring  
555 Chemical Reaction Space with Machine Learning Models: Representation and Feature Perspec-  
556 tive. *Journal of Chemical Information and Modeling*, 2024.
- 557 Carl Edwards, Tuan Lai, Kevin Ros, Garrett Honke, Kyunghyun Cho, and Heng Ji. Translation  
558 between Molecules and Natural Language. In *Proceedings of the 2022 Conference on Empirical  
559 Methods in Natural Language Processing*, pp. 375–413, Abu Dhabi, United Arab Emirates, De-  
560 cember 2022. Association for Computational Linguistics. URL <https://aclanthology.org/2022.emnlp-main.26>.  
561  
562
- 563 Hanyu Gao, Thomas J Struble, Connor W Coley, Yuran Wang, William H Green, and Klavs F  
564 Jensen. Using Machine Learning To Predict Suitable Conditions for Organic Reactions. *ACS  
565 central science*, 4(11):1465–1476, 2018.
- 566 Taicheng Guo, Bozhao Nan, Zhenwen Liang, Zhichun Guo, Nitesh Chawla, Olaf Wiest, Xiangliang  
567 Zhang, et al. What can Large Language Models do in chemistry? A comprehensive benchmark  
568 on eight tasks. *Advances in Neural Information Processing Systems*, 36:59662–59688, 2023.  
569
- 570 Ross Irwin, Spyridon Dimitriadis, Jiazhen He, and Esben Jannik Bjerrum. Chemformer: a pre-  
571 trained transformer for computational chemistry. *Machine Learning: Science and Technology*, 3  
572 (1):015022, 2022.  
573
- 574 Andrew Jaegle, Felix Gimeno, Andy Brock, Oriol Vinyals, Andrew Zisserman, and Joao Carreira.  
575 Perceiver: General Perception with Iterative Attention. In *International conference on machine  
576 learning*, pp. 4651–4664. PMLR, 2021.
- 577 Ming Jin, Shiyu Wang, Lintao Ma, Zhixuan Chu, James Y Zhang, Xiaoming Shi, Pin-Yu Chen,  
578 Yuxuan Liang, Yuan-Fang Li, Shirui Pan, and Qingsong Wen. Time-LLM: Time Series Fore-  
579 casting by Reprogramming Large Language Models. In *International Conference on Learning  
580 Representations (ICLR)*, 2024.  
581
- 582 Jiatong Li, Yunqing Liu, Wenqi Fan, Xiao-Yong Wei, Hui Liu, Jiliang Tang, and Qing Li. Em-  
583 powering Molecule Discovery for Molecule-Caption Translation with Large Language Models:  
584 A ChatGPT Perspective. *IEEE Transactions on Knowledge and Data Engineering*, 2024.  
585
- 586 Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping Language-Image Pre-  
587 training with Frozen Image Encoders and Large Language Models. In *International conference  
588 on machine learning*, pp. 19730–19742. PMLR, 2023.
- 589 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual Instruction Tuning. *Advances  
590 in neural information processing systems*, 36, 2024a.  
591
- 592 Pengfei Liu, Yiming Ren, Jun Tao, and Zhixiang Ren. GIT-Mol: A Multi-modal Large Language  
593 Model for Molecular Science with Graph, Image, and Text. *Computers in Biology and Medicine*,  
171:108073, 2024b.

- 594 Micha Livne, Zulfat Miftahutdinov, Elena Tutubalina, Maksim Kuznetsov, Daniil Polykovskiy, An-  
595 nika Brundyn, Aastha Jhunjunwala, Anthony Costa, Alex Aliper, Alán Aspuru-Guzik, et al.  
596 nach0: Multimodal natural and chemical languages foundation model. *Chemical Science*, 15(22):  
597 8380–8389, 2024.
- 598 Jieyu Lu and Yingkai Zhang. Unified Deep Learning Model for Multitask Reaction Predictions with  
599 Explanation. *Journal of chemical information and modeling*, 62(6):1376–1387, 2022.
- 600 Andres M. Bran, Sam Cox, Oliver Schilter, Carlo Baldassari, Andrew D White, and Philippe  
601 Schwaller. Augmenting large language models with chemistry tools. *Nature Machine Intelli-*  
602 *gence*, pp. 1–11, 2024.
- 603 Michael R Maser, Alexander Y Cui, Serim Ryou, Travis J DeLano, Yisong Yue, and Sarah E Reis-  
604 man. Multi-Label Classification Models for the Prediction of Cross-Coupling Reaction Condi-  
605 tions. *Journal of Chemical Information and Modeling*, 61(1):156–166, 2021.
- 606 S Hessam M Mehr, Matthew Craven, Artem I Leonov, Graham Keenan, and Leroy Cronin. A  
607 universal system for digitization and automatic execution of the chemical synthesis literature.  
608 *Science*, 370(6512):101–108, 2020.
- 609 Barbara Mikulak-Klucznik, Patrycja Gołebiewska, Alison A Bayly, Oskar Popik, Tomasz Klucznik,  
610 Sara Szymkuć, Ewa P Gajewska, Piotr Dittwald, Olga Staszewska-Krajewska, Wiktor Beker, et al.  
611 Computational planning of the synthesis of complex natural products. *Nature*, 588(7836):83–88,  
612 2020.
- 613 Juno Nam and Jurae Kim. Linking the Neural Machine Translation and the Prediction of Organic  
614 Chemistry Reactions. *arXiv preprint arXiv:1612.09529*, 2016.
- 615 Yujie Qian, Zhening Li, Zhengkai Tu, Connor Coley, and Regina Barzilay. Predictive Chemistry  
616 Augmented with Text Retrieval. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceed-*  
617 *ings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 12731–  
618 12745, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/  
619 2023.emnlp-main.784. URL <https://aclanthology.org/2023.emnlp-main.784>.
- 620 Laria Reynolds and Kyle McDonell. Prompt Programming for Large Language Models: Beyond  
621 the Few-Shot Paradigm. In *Extended Abstracts of the 2021 CHI Conference on Human Factors*  
622 *in Computing Systems*, pp. 1–7, 2021.
- 623 Simon Rohrbach, Mindaugas Šiaučius, Greig Chisholm, Petrisor-Alin Pirvan, Michael Saleeb,  
624 S Hessam M Mehr, Ekaterina Trushina, Artem I Leonov, Graham Keenan, Aamir Khan, et al.  
625 Digitization and validation of a chemical synthesis literature database in the ChemPU. *Science*,  
626 377(6602):172–180, 2022.
- 627 Michael Schlichtkrull, Thomas N Kipf, Peter Bloem, Rianne Van Den Berg, Ivan Titov, and Max  
628 Welling. Modeling relational data with graph convolutional networks. In *The semantic web: 15th*  
629 *international conference, ESWC 2018, Heraklion, Crete, Greece, June 3–7, 2018, proceedings*  
630 *15*, pp. 593–607. Springer, 2018.
- 631 Tobias Schnitzer, Martin Schnurr, Andrew F Zahrt, Nader Sakhaee, Scott E Denmark, and Helma  
632 Wennemers. Machine Learning to Develop Peptide Catalysts- Successes, Limitations, and Op-  
633 portunities. *ACS Central Science*, 2024.
- 634 Philippe Schwaller, Teodoro Laino, Théophile Gaudin, Peter Bolgar, Christopher A Hunter, Costas  
635 Bekas, and Alpha A Lee. Molecular Transformer: A Model for Uncertainty-Calibrated Chemical  
636 Reaction Prediction. *ACS central science*, 5(9):1572–1583, 2019.
- 637 Philippe Schwaller, Daniel Probst, Alain C Vaucher, Vishnu H Nair, David Kreutter, Teodoro Laino,  
638 and Jean-Louis Reymond. Mapping the Space of Chemical Reactions Using Attention-Based  
639 Neural Networks. *Nature machine intelligence*, 3(2):144–152, 2021.
- 640 Benjamin J Shields, Jason Stevens, Jun Li, Marvin Parasram, Farhan Damani, Jesus I Martinez  
641 Alvarado, Jacob M Janey, Ryan P Adams, and Abigail G Doyle. Bayesian reaction optimization  
642 as a tool for chemical synthesis. *Nature*, 590(7844):89–96, 2021.

- 648 Connor J Taylor, Alexander Pomberger, Kobi C Felton, Rachel Grainger, Magda Barecka,  
649 Thomas W Chamberlain, Richard A Bourne, Christopher N Johnson, and Alexei A Lapkin. A  
650 Brief Introduction to Chemical Reaction Optimization. *Chemical Reviews*, 123(6):3089–3126,  
651 2023.
- 652 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-  
653 lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open Founda-  
654 tion and Fine-Tuned Chat Models. *arXiv preprint arXiv:2307.09288*, 2023.
- 655 Maria Tsimpoukelli, Jacob L Menick, Serkan Cabi, SM Eslami, Oriol Vinyals, and Felix Hill. Mul-  
656 timodal Few-Shot Learning with Frozen Language Models. *Advances in Neural Information*  
657 *Processing Systems*, 34:200–212, 2021.
- 658 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,  
659 Łukasz Kaiser, and Illia Polosukhin. Attention Is All You Need. *Advances in neural information*  
660 *processing systems*, 30, 2017.
- 661 Xiaorui Wang, Chang-Yu Hsieh, Xiaodan Yin, Jike Wang, Yuquan Li, Yafeng Deng, Dejun Jiang,  
662 Zhenxing Wu, Hongyan Du, Hongming Chen, et al. Generic Interpretable Reaction Condition  
663 Predictions with Open Reaction Condition Datasets and Unsupervised Learning of Reaction Cen-  
664 ter. *Research*, 6:0231, 2023a.
- 665 Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and  
666 Hannaneh Hajishirzi. Self-Instruct: Aligning Language Models with Self-Generated Instructions.  
667 In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings of the 61st Annual*  
668 *Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 13484–  
669 13508, Toronto, Canada, July 2023b. Association for Computational Linguistics. doi: 10.18653/  
670 v1/2023.acl-long.754. URL <https://aclanthology.org/2023.acl-long.754>.
- 671 David Weininger, Arthur Weininger, and Joseph L Weininger. SMILES. 2. Algorithm for generation  
672 of unique SMILES notation. *Journal of chemical information and computer sciences*, 29(2):  
673 97–101, 1989.
- 674 Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. MiniGPT-4: Enhancing  
675 Vision-Language Understanding with Advanced Large Language Models. In *The Twelfth Interna-*  
676 *tional Conference on Learning Representations*, 2024. URL [https://openreview.net/  
677 forum?id=1tzZbq88f27](https://openreview.net/forum?id=1tzZbq88f27).

## 682 APPENDIX

### 683 A TRAINING SETTINGS

684 To realize peak efficiency within our Chemma-RC model, we carefully design the training phases.  
685 This section offers a comprehensive summary of the training settings and the hyperparameter values.  
686 Through the detailed orchestration of these parameters, we ensure that Chemma-RC is capable of  
687 fully leveraging its capabilities in the application contexts.

- 688 • **Optional Settings:** There are alternatives for modification in the Chemma-RC framework,  
689 such as the replacement of the Perceiver-based modality projection layer with other archi-  
690 tectures like Reprogramming and MLP.
- 691 • **Reaction Condition Recommendation task:** Within the framework, the model takes the  
692 32-layer LLaMA-2-7b as the LLM backbone. Besides, we utilize a pre-trained SMILES-  
693 to-text retriever proposed by Qian et al. (Qian et al., 2023) and extract the most similar  
694 unpaired corpus as the reaction text. Meanwhile, we introduce Parrot, a Bert-like model  
695 to encode the reaction SMILES. We leverage R-GCN (Schlichtkrull et al., 2018) to encode  
696 the molecules in the reaction, and the combination of reactant and product embeddings  
697 is considered as the reaction representation. In the training process, the encoders in all  
698 modalities are frozen. After the alignment of the representation space, the SMILES and  
699 the graph-based tokens have a length of 128 and 3, respectively. Additionally, the model  
700  
701

702 employs the OneCycleLR as the learning rate scheduler, initializing the learning rate as  
 703 3e-5. The batch size is set to 16, with less than 6 epochs 48 hours in training. The GPU  
 704 configuration is  $8 \times 80\text{G A800}$ .  
 705

## 706 B DATA DESCRIPTION

707 We curate two large datasets, named USPTO-Condition and USPTO\_500MT\_Condition, with the  
 708 data volumes presented in Table. 8. Both datasets are split with the ratio of train:validation:test as  
 709 8:1:1 in our work. For USPTO-Condition dataset, all molecules including reactants, products, and  
 710 conditions are collected in canonical SMILES. Each reaction entry contains five condition labels,  
 711 including one catalyst, two solvents, two reagents, and an additional “none” category is introduced  
 712 to illustrate that the reaction does not require this type of reaction condition (Gao et al., 2018).  
 713 The visualization of data distribution is depicted in Figure. 4 (left). From Figure. 4 we can see  
 714 that this dataset covers a vast variety of reaction types, characterized by a substantial proportion  
 715 of heteroatom alkylation, arylation, and acylation reactions, while C-C formation reactions are less  
 716 included. We also introduce the corpus of reaction descriptions proposed by Qian et al. (Qian et al.,  
 717 2023) into the USPTO-Condition dataset. Each reaction is associated with a corpus of reaction  
 718 descriptions. It should be noted that the corpus will not be utilized directly for training. Instead, we  
 719 employ the corpus as an input for the pre-trained retrieval module proposed by (Qian et al., 2023).  
 720 This approach allows us to obtain similar embeddings necessary for the multimodal representation  
 721 learning of our Chemma-RC, and avoid data leaking issues. For USPTO\_500MT\_Condition datasets,  
 722 it collects top-500 types of reactions from the USPTO-MIT datasets (Coley et al., 2017), in which  
 723 the top-100 types of reactions make up 59% of the entire dataset, which can be seen in Figure. 4  
 724 (right). In order to calculate the predicted accuracy on the USPTO\_500MT\_Condition dataset, it  
 725 is necessary to separate all reagents in an appropriate manner. However, separating reagents using  
 726 the dot as a delimiter is challenging, as compounds like  $[\text{Na}^+].[\text{OH}^-]$  constitutes a single reagent  
 727 and cannot be split. Besides, to have a comprehensive knowledge of the datasets, we do sparsity  
 728 analyses. We calculate the non-empty count and density of every condition in the USPTO-Condition  
 729 dataset, which is presented in Table. 9. From the table, we can see that some conditions, such as  
 730 ‘Catalyst’, ‘Solvent 2’, and ‘Reagent 2’ show a high extent of sparsity, with a non-empty density  
 731 of fewer than 30%. For the USPTO\_500MT\_Condition, as it only covers the condition of non-split  
 732 reagents, all of the reaction entries have their corresponding non-empty condition label.

733 Furthermore, we make an investigation on the condition categories in the USPTO-Condition and  
 734 USPTO\_500MT\_Condition dataset, which is illustrated in Figure. 5. The visualization of the most  
 735 common chemical contexts of the reagents, catalysts, and solvents in USPTO-Condition, and separate  
 736 reagents in USPTO\_500MT\_Condition is depicted in Figure. 5 (A-D), respectively. From the  
 737 figures, we learn that reaction conditions have a property of diversity and imbalance. Besides, we  
 738 count categories of every condition, as is presented in Figure. 5 (E). Reagents in both datasets consist  
 739 of more than 200 categories, which highlights the difficulty of the reaction condition recommenda-  
 740 tion task. Additionally, we prove that reagents in the USPTO\_500MT\_Condition dataset follow the  
 741 power-law distribution, which indicates the condition keeps the long-tail feature in distribution and  
 742 a small number of categories account for the majority of the data size.

743 Table 7: Question templates generated by GPT-4.

745 Task	746 Description
747 Solvent prediction	748 Could you suggest potential solvents that could have been used 749 in the given chemical reaction, taking into consideration their polarity and compatibility with the reactants?
750 Reagent prediction	751 Please suggest some possible reagents that could have been used in the following chemical reaction.
752 Catalyst prediction	753 Considering the chemical reaction in question, which catalysts could be effective?
754 Condition prediction (all)	755 Given the current chemical reaction, what would be the appropriate conditions to consider?

Table 8: Data volume of USPTO-Condition and USPTO\_500MT\_Condition datasets.

Dataset	Training set	Validation set	Testing set
USPTO-Condition	546,728	68,341	68,341
USPTO_500MT_Condition	88,410	9,778	10,828

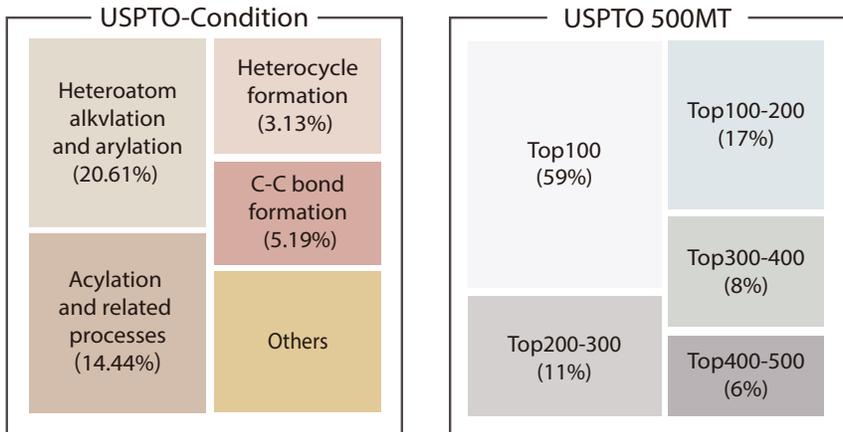


Figure 4: Left: The reaction distribution of USPTO-Condition. Right: The reaction distribution of USPTO\_500MT\_Condition.

## C DETAILS OF MODALITY ALIGNMENT

For the reaction condition recommendation task, the representation of the reaction is extracted by encoders (see in section 3.2.2), and the text representation is tokenized by LLMs. However, fusing two types of representation introduces inductive biases issues (Baltrušaitis et al., 2018; Jaegle et al., 2021). To effectively fuse representations from multiple modalities, we propose the use of a projection module, the Perceiver (Jaegle et al., 2021), for modality alignment (Figure 1). This module employs latent queries to align graph and SMILES tokens with text-related tokens, such as question prompts and a text-augmented corpus. We show the pseudo-code for modality projection in Algorithm. 1.

## D MODEL PERFORMANCE

A chemical reaction can be represented as the transformation of a sequence of characters (reactants, conditions) into another sequence (products), with compounds connected by special characters, such as '>>'. This structure makes sequence-to-sequence models, such as the Transformer, well-suited for predictive modeling of reaction representation (Schwaller et al., 2019; Irwin et al., 2022). However, existing SMILES-based Transformer models for reaction representation encounter limitations in various aspects, particularly with respect to atom permutations and the interpretability of reaction mechanisms. Consequently, our proposed Chemma-RC fuses data from diverse sources including corpus, SMILES and graphs of molecules to present a comprehensive view of the reaction. We assess the performance of our proposed Chemma-RC and the aforementioned baseline methods for reaction condition recommendation. The top- $N$  accuracy of condition recommendation on the combined test datasets of USPTO-Condition and USPTO\_500MT\_Condition are presented in Table. 2

Table 9: Sparsity analysis of the USPTO-Condition dataset.

USPTO-Condition	Catalyst	Solvent 1	Solvent 2	Reagent 1	Reagent 2
<b>Non-empty count</b>	89,756	673,634	130,326	504,169	170,752
<b>Non-empty density</b>	13%	99%	19%	74%	25%

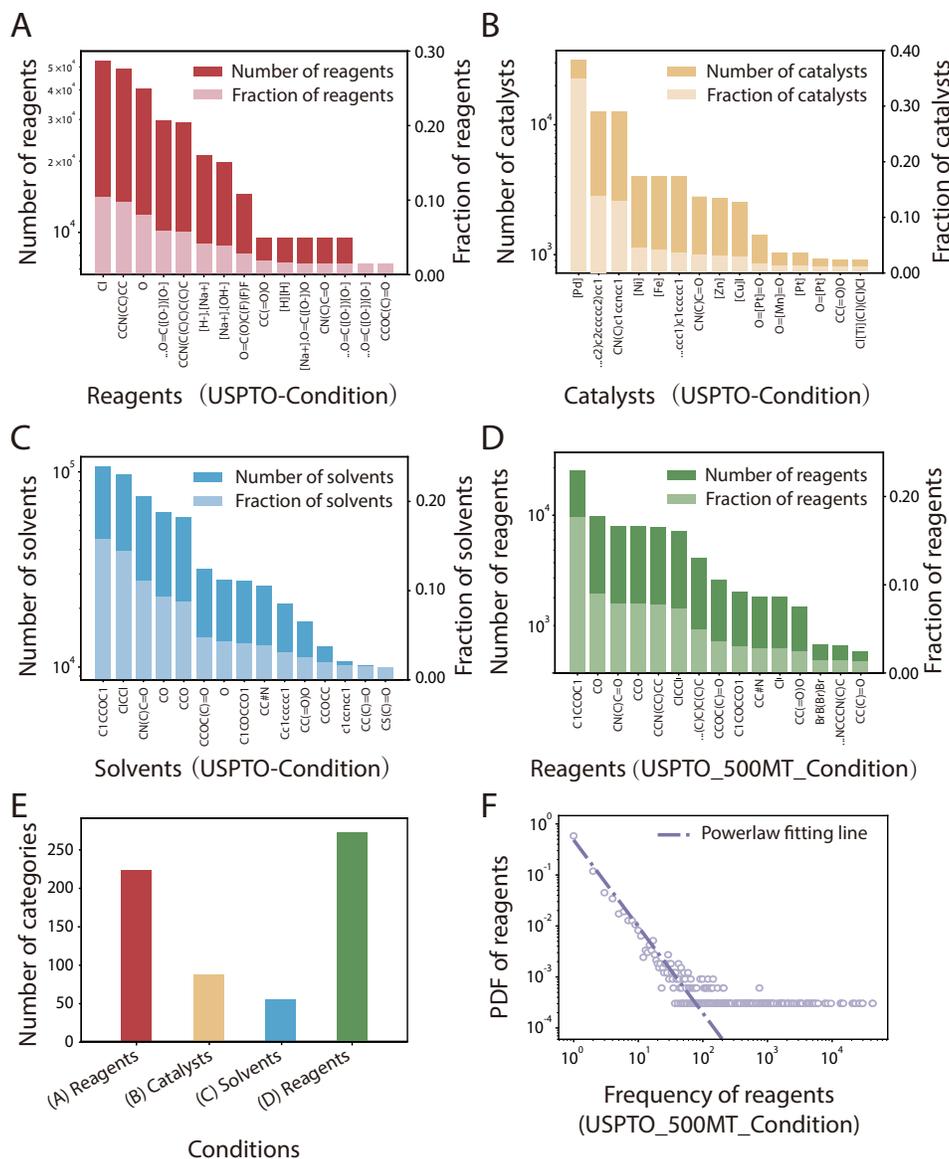


Figure 5: **Distribution of types of reactions in the USPTO-Condition and USPTO\_500MT-Condition.** (A-D) The bar charts of the fifteen most common reagents, catalysts, and solvents in the USPTO-Condition and reagents in the USPTO\_500MT-Condition, respectively, where the shallow color presents the decimal-scale proportion and the deep color presents the log-scale count. (E) The bar charts of the total category count of the conditions illustrated in (A-D). (F) Power law fitting of the reagent distribution in the USPTO\_500MT-Condition, where the shallow points show the probability density and the deep dashed-line shows the ideal power-law fitting, respectively.

---

```

864 Algorithm 1 Pseudo code for modality projection.
865 word_proj, perceiver_proj: predefined linear and transformer-based projectors, respec-
866 tively.
867 # B: batch size; C: channel size; n: content shape
868 # M: query length; N: shape of flatten reaction tokens;
869 # text_q: text query in shape (B, M, C)
870 # react_embed: reaction embedding in shape (B, N, C)
871 # word_embed: word embedding in shape (B, vocab_size, C)
872
873 # Key part 1: map transformer-based reaction feature
874 word_embed = self.word_proj(word_embed)
875 word_embed = word_embed.repeat(react_embed.size()[0], 1, 1)
876 react_embed = torch.cat([react_embed, word_embed], dim=1)
877 smiles_react_tokens = linear_layer(react_embed) # to make 128
878     tokens
879
880 # Key part 2: map graph-based reaction features
881 graph_embed = self.word_proj(graph_embed)
882 graph_react_tokens = linear_layer(graph_embed) # to make 3 tokens
883
884 # Key part 3:
885 reaction_tokens = torch.cat([smiles_react_tokens,
886     graph_react_tokens], dim=1)
887
888 # Key part 4: modality projection
889 reaction_tokens_from_smiles = self.perceiver_proj_smiles(
890     smiles_react_tokens)
891 reaction_tokens_from_graphs = self.perceiver_proj_graphs(
892     graph_react_tokens)
893
894 # concat token
895 final_token = torch.cat([reaction_tokens_from_smiles,
896     reaction_tokens_from_graphs, text_q], dim=1)

```

---

and Table. 3, respectively. We introduce several comparative methods to illustrate the performance of Chemma-RC.

1. rxnfp LSTM (Gao et al., 2018). This method proposes a reaction fingerprint to represent the difference between the product and reactant fingerprints.
2. rxnfp retrieval. It uses the conditions of the most similar reactions in the training set as the prediction. Similar reactions are determined based on the  $L_2$  distance of reaction fingerprints.
3. Transformer. It uses the same architecture as the TextReact predictor. This baseline represents the state-of-the-art model that only takes chemistry input.
4. ChemBERTa Chithrananda et al. (2020). It is same as the Transformer baseline except that the encoder is pre-trained on external SMILES data.
5. Reaction GCNN (Maser et al., 2021). This method proposes a machine-learned ranking model to predict the set of conditions used in a reaction as a binary vector.
6. Parrot (Wang et al., 2023a). This method leverages the attention-based model architecture to encode the reaction and design a training methodology specifically to enhance the reaction center.
7. TextReact (Qian et al., 2023). It aims to enhance the molecular representation of the reaction by introducing relevant corpus retrieved from literature into sequence-to-sequence Transformers.
8. Reagent Transformer (Andronov et al., 2023). This method leverages Molecular Transformer, (Schwaller et al., 2019) a state-of-the-art model to tackle the task of reagent prediction.

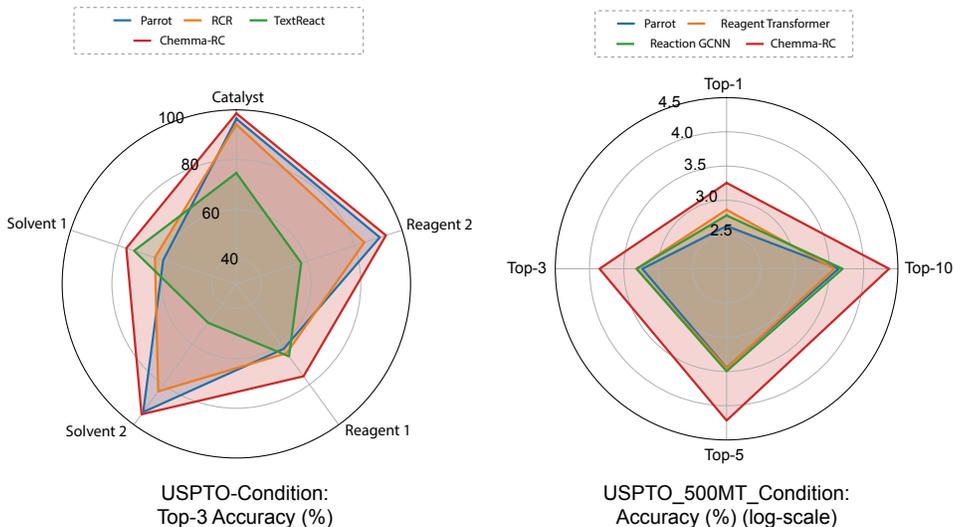


Figure 6: **Left:** Radar plot of top-3 prediction accuracy of conditions on the USPTO-Condition dataset. The classification performance consists of comparative methods such as Parrot, RCR, TextReact, and our methods with similar corpus. **Right:** Radar chart of log-scale accuracy of reagents in the USPTO\_500MT\_Condition dataset.

To have a comprehensive overview of the recommendation performance, we visualize the prediction results of USPTO-Condition and USPTO\_500MT\_Condition datasets, as described in Table. 2, 3. Specifically, we draw radar charts of our model and other competitive models, which are presented in Figure. 6. For the USPTO-Condition dataset, we reproduce Parrot, RCR, and TextReact. Then, we plot the top-3 predicting accuracy of different conditions (catalyst, solvent 1, solvent 2, reagent 1, and reagent 2), as depicted in Figure. 6 (left). For the USPTO\_500MT\_Condition dataset, we recommend reagents in SMILES sequence and take Parrot, Reagent Transformer, and Reaction GCNN as comparative methods. For more intuition, we visualize top-1, 3, 5, and 10 exactly matched accuracy in log scale, which is shown in Figure. 6 (right). From the charts, we can see that our model covers the largest area of the performance circle in both datasets, indicating that Chemma-RC markedly outperforms other competitive models.

#### D.1 GERALIZATION PERFORMANCE

In order to validate the out-of-domain performance of Chemma-RC, we employ Chemma-RC trained on the USPTO\_500MT\_Condition to test on the USPTO-Condition. The evaluation strategy includes three specific training conditions: reagents, catalysts, and solvents. We adopt a metric of **partial matched accuracy** to illustrate the generalization capability of Chemma-RC. **Different from the complete matched accuracy that requires perfect matching between predictions and labels, the partial matched accuracy is more suitable to test the generalization capacity, which focuses more on whether the predicted results match a substitutable part of the ground truth. For example, if the predicted result is '[Na+].[OH-]' and the condition label is 'CO.[Na+].[OH-]', we consider that the prediction partially matches the ground truth, but not completely.** The evaluation strategy includes three specific training conditions: reagents, catalysts, and solvents. Table. 10 reports the top-1 partial match accuracy for each condition prediction. From the results we can see that, Chemma-RC achieves a top-1 partial matched accuracy of 67.1% and 58.1%, respectively. This relatively high accuracy indicates that solvents and reagents have more consistent characteristics that the model can learn effectively from USPTO\_500MT\_Condition and apply to USPTO-Condition. In contrast, The model’s performance in predicting catalysts demonstrates a lower top-1 partial match accuracy at 89.9%.

Chemma-RC can successfully distinguish reagents from the combination of all conditions in a reaction. Additionally, training Chemma-RC on USPTO-Condition, a larger chemical reaction dataset, further enhances its ability to akin chemical knowledge.

Table 10: The top-1 partial matched accuracy of Chemma-RC under OOD setting.

Evaluation strategy (train $\rightarrow$ test)	Acc (%)
USPTO_500MT_Condition $\rightarrow$ USPTO-Condition (reagent)	67.1
USPTO_500MT_Condition $\rightarrow$ USPTO-Condition (catalyst)	89.9
USPTO_500MT_Condition $\rightarrow$ USPTO-Condition (solvent)	58.1

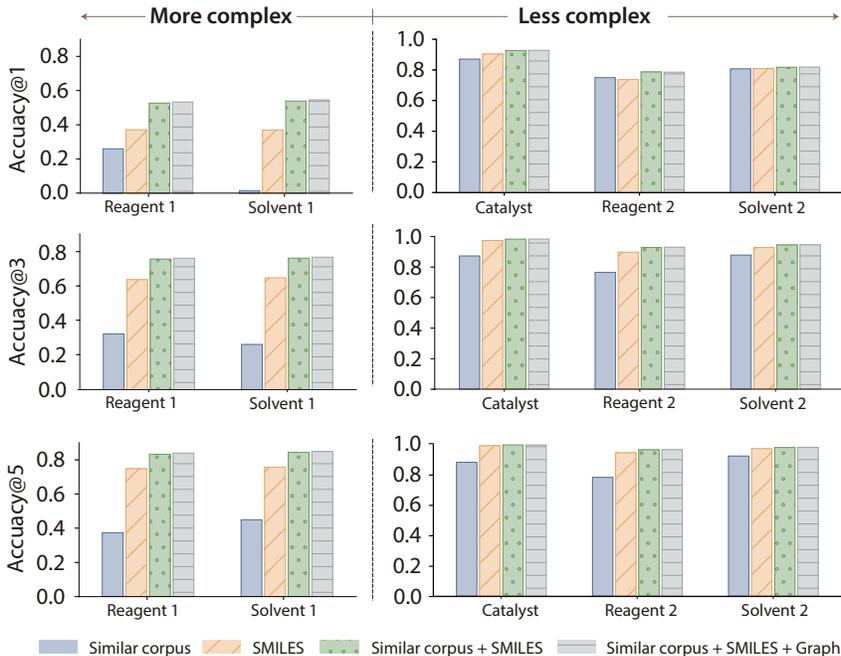


Figure 7: Bar charts demonstrating the ablation study of modalities including similar corpus, SMILES and graph. The classification performance is assessed on the conditions in the USPTO-Condition dataset, which are split into two groups according to data sparsity.

## D.2 ABLATION STUDY ON MODALITY

Besides, we visualize the results of the ablation study on modality on the USPTO-Condition dataset, which can be seen in Table 4. Specifically, we categorize the conditions of the USPTO-Condition into two types: more complex and less complex. According to the data sparsity, reagent 1 and solvent 1 are considered more complex, while catalyst, reagent 2, and solvent 2 are considered less complex. Then, the investigation on the effectiveness of modalities comprising similar corpus, SMILES, graph is depicted in Figure 7. From the results, we can see that compared with the model with multiple modalities, the model with single one modality degrades dramatically. Moreover, Chemma-RC with three modalities combined achieves the best performance, which demonstrates the vital importance of capturing the reaction representations from different dimensions.

## D.3 CASE STUDY

In this section, we select four cross-coupling reactions from USPTO-Condition datasets for performance validation. We visualize the predicted results in Figure 9. As depicted in Figure 9, the reaction centers and leaving groups are highlighted in different colors. For C–N cross-coupling reactions (the first and the third row), Chemma-RC can predict all conditions precisely. For C–C bond formation and Formylation reactions (the second and the fourth row), Chemma-RC fails to predict Ethyl Acetate (the second case) and THF (the fourth case). The reason why Chemma-RC is less effective for these reactions is that the data volume of C–C bond formation reactions in the USPTO-Condition dataset is only 5%, as shown in Figure 4. This limited representation constrains the

model’s ability to learn the patterns associated with C–C bond formation reactions. Consequently, Chemma-RC lacks sufficient training examples to capture and generalize the underlying reaction mechanisms accurately. The scarcity of diverse and representative data hampers its effectiveness, leading to a lower precision in predicting these types of reactions.

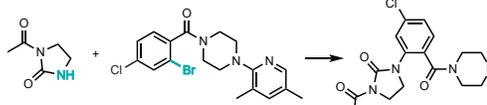
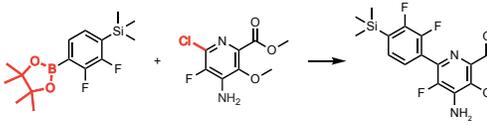
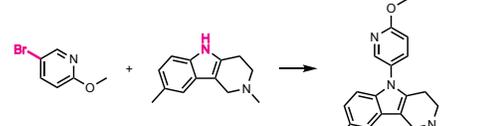
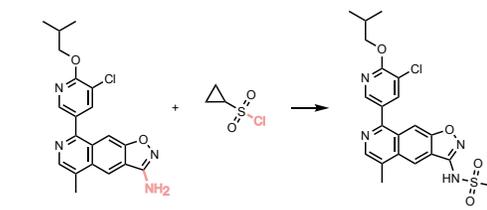
Reactions	First line: label;	Second line: prediction	Catalyst 1	Solvent 1	Solvent 2	Reagent 1	Reagent 2
	Cu–I	Cu–I	Cu–I	1,4-Dioxane	H <sub>2</sub> O	DMEN	K <sub>3</sub> PO <sub>4</sub>
	Dichlorobis (tricyclohexylphosphine) palladium(II)	Dichlorobis (tricyclohexylphosphine) palladium(II)	Dichlorobis (tricyclohexylphosphine) palladium(II)	Ethyl Acetate	H <sub>2</sub> O	MeCN	Na <sub>2</sub> CO <sub>3</sub>
	Dichlorobis (tricyclohexylphosphine) palladium(II)	Dichlorobis (tricyclohexylphosphine) palladium(II)	Dichlorobis (tricyclohexylphosphine) palladium(II)	H <sub>2</sub> O	None	MeCN	Na <sub>2</sub> CO <sub>3</sub>
				✗ Ethyl Acetate has not been predicted			✗
	Cu–I	Cu–I	Cu–I	DMF	H <sub>2</sub> O	L-Proline	K <sub>3</sub> PO <sub>4</sub>
	Cu–I	Cu–I	Cu–I	DMF	H <sub>2</sub> O	L-Proline	K <sub>3</sub> PO <sub>4</sub>
	DMAP	DMAP	DMAP	Ethyl Acetate	1,10-phenanthroline	H <sub>2</sub> O	THF
	DMAP	DMAP	DMAP	1,10-phenanthroline	1,10-phenanthroline	H <sub>2</sub> O	H <sub>2</sub> O
				✗ Ethyl Acetate has been predicted to 1,10-phenanthroline			✗
				THF has been predicted to H <sub>2</sub> O			✗

Figure 8: Visualization of recommended conditions on four reactions. We select four Suzuki–Miyaura cross-coupling reactions to present the performance of condition recommendation. The reaction centers and leaving groups are highlighted in different colors.

Further, we visualize the predicted results on OOD datasets in Figure. 9. We select two reaction cases for analysis. In case 1, Toluene is not predicted by Chemma-RC. In case 2, 1,4-Dioxane and 1-(diphenylphosphaneyl)cyclopenta-2,4-dien-1-ide are predicted. However, it is confirmed that Toluene and 1,4-Dioxane are common solvents, and 1-(diphenylphosphaneyl)cyclopenta-2,4-dien-1-ide is frequently used as a ligand. Therefore, we do not categorize these as failed cases because the model successfully predicts all the reagents in the labels and avoids predicting other conditions.

1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
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

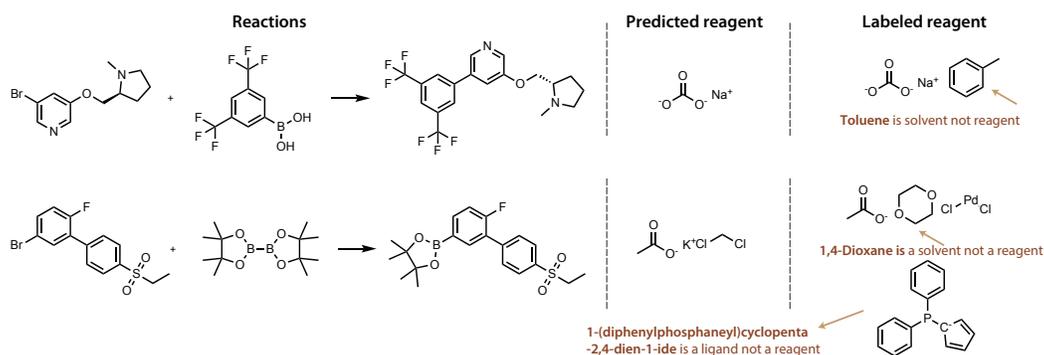


Figure 9: Visualization of recommended conditions on two reactions. In case 1, Toluene was not predicted by Chemma-RC. In case 2, 1,4-Dioxane and 1-(diphenylphosphanyl)cyclopenta-2,4-dien-1-ide were predicted. However, it is confirmed that Toluene and 1,4-Dioxane are common solvents, and 1-(diphenylphosphanyl)cyclopenta-2,4-dien-1-ide is frequently used as a ligand. Therefore, we do not categorize these as failed cases because the model successfully predicts all the reagents in the labels and avoids predicting other conditions.

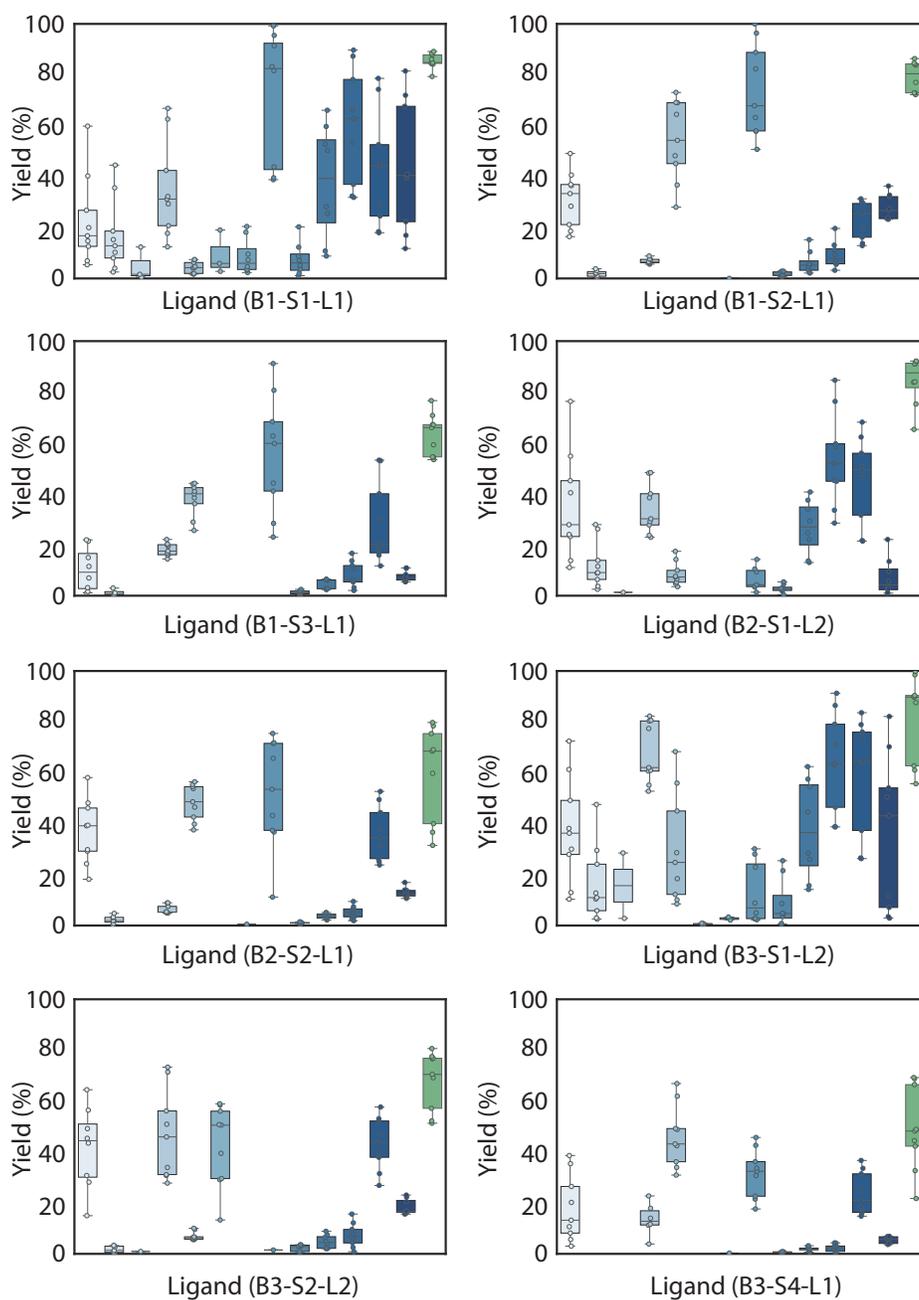


Figure 10: Boxplot of the performance for ligand recommendation (1).

1188  
1189  
1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197  
1198  
1199  
1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1240  
1241

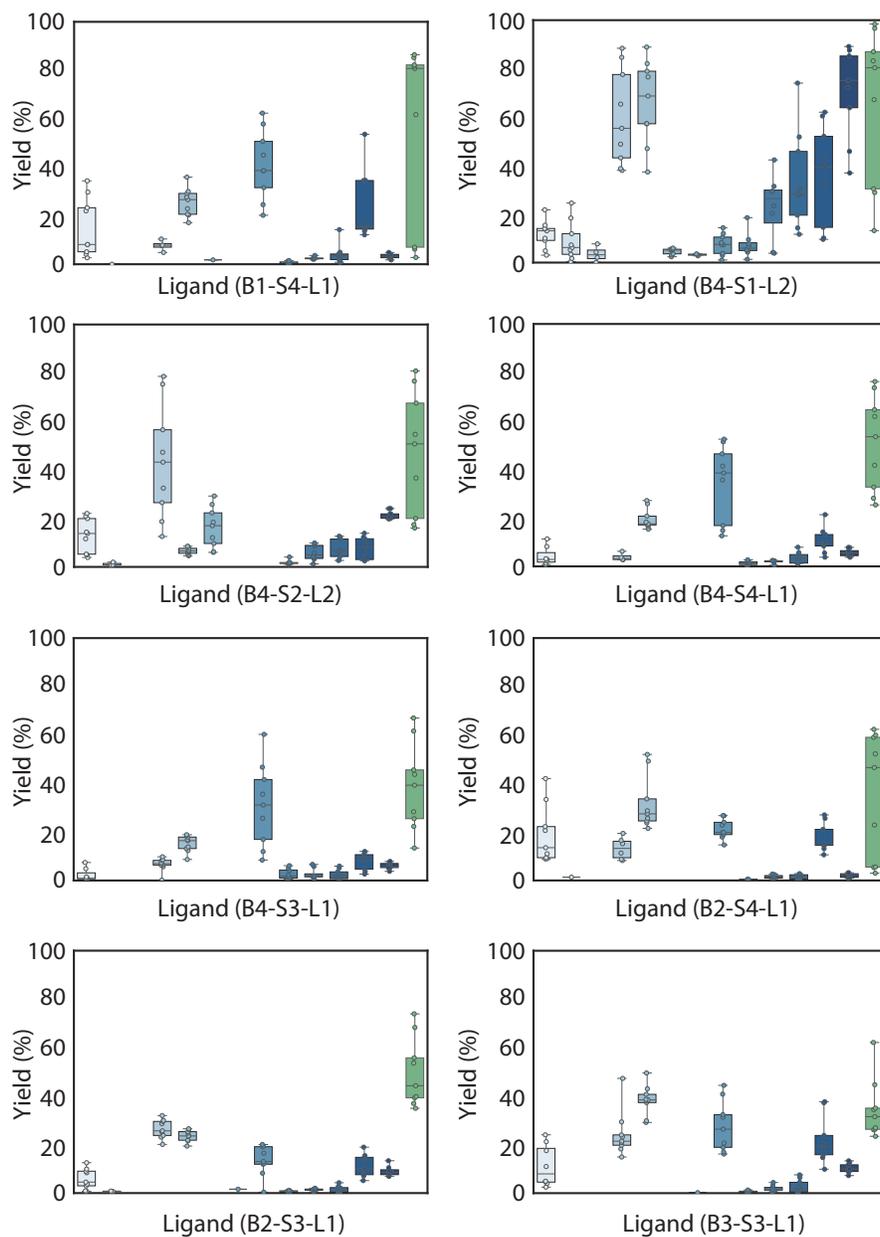


Figure 11: Boxplot of the performance for ligand recommendation (2).