# ChemBART: A Pre-trained BART Model Assisting Organic Chemistry Analysis

Kenan Li [a,1], Yijian Zhang [a,1], Jin Wang [b,1], Haipeng Gan [a], Zeying Sun [b], Xiaoguang Lei [b,*], & Hao Dong [a,*]

[a] State Key Laboratory of Analytical Chemistry for Life Science, Kuang Yaming Honors School, Chemistry and Biomedicine Innovation Centre (ChemBIC), ChemBioMed Interdisciplinary Research Centre at Nanjing University, Engineering Research Centre of Protein and Peptide Medicine of the Ministry of Education, Institute for Brain Sciences, Nanjing University, Nanjing 210023, China.
[b] Beijing National Laboratory for Molecular Sciences, Key Laboratory of Bioorganic Chemistry and Molecular Engineering of Ministry of Education, College of Chemistry and Molecular Engineering, and Peking-Tsinghua Center for Life Sciences, Peking University, Beijing 100871, China

[1] these authors contributed equally

* corresponding authors: xglei@pku.edu.cn (X.G.L), donghao@nju.edu.cn (H.D.)

**Abstract:** Recent advances in large language models (LLMs) have demonstrated transformative potential across diverse fields. While LLMs have been applied to molecular simplified molecular input line entry system (SMILES) in computer-aided synthesis planning (CASP), existing methodologies typically address single tasks, such as precursor prediction. We introduce ChemBART, a SMILES-based LLM pre-trained on chemical reactions, which enables a unified model for multiple downstream chemical tasks—achieving the paradigm of "one model, one pre-training, multiple tasks." By leveraging outputs from a mask-filling pre-training task on reaction expressions, ChemBART effectively solves a variety of chemical problems, including precursor/reagent generation, temperature-yield regression, molecular property classification, and optimizing the policy and value functions within a reinforcement learning framework, integrated with Monte Carlo tree search for multi-step synthesis route design. Unlike single-molecule pre-trained LLMs constrained to specific applications, ChemBART addresses broader chemical challenges and integrates them for comprehensive synthesis planning. Crucially, ChemBART-designed multi-step synthesis routes and reaction conditions directly inspired wet-lab validation, which confirmed shorter pathways with ~30% yield improvement over literature benchmarks. Our work validates the power of reaction-focused pre-training and showcases the broad utility of ChemBART in advancing the complete synthesis planning cycle.

# Introduction

The field of Natural Language Processing (NLP) has witnessed remarkable achievement with the advent of large language models (LLMs) based on the Transformer architecture[1]. Prominent examples include GPT[2], BERT[3], and BART[4].

Recently, the simplified molecular input line entry system (SMILES) notation[5], a standard representation for chemical molecules, has been widely seen as a language of the chemical world. Consequently, LLMs have been increasingly applied to acquire chemical knowledge by learning the textual information embedded within SMILES strings.

A particularly successful application based on this concept is precursor prediction, exemplified by models such as AutoSynRoute[6], Molecular Transformer[7], GTA[8], and Graph2SMILES[9], etc.. These systems typically input a target product molecule (often as its SMILES string) into an LLM, which then generates the synthetic precursors. This methodology is usually referred to as "template-free," as reaction rules are implicitly embedded within the parameters learned by the LLMs during training.

Modern LLMs have widely adopted the idea of pre-training – training models to recover a large amount of unlabeled data can boost downstream tasks with only a small amount of labeled data. Some language models based on SMILES notation, such as Chemformer[10], X-Mol[11], SMILES-BERT[12], MolBert[13], T5Chem[14], ChemLM[15], CleanMol,[16] and ReactSeq[17], have also conducted similar pre-training on datasets of single molecules. For example, Chemformer[10] undergoes pre-training on mask-filling and synonym transformations of molecule SMILES expressions, followed by fine-tuning for downstream tasks, such as precursor prediction, molecular structure optimization, and regression of molecular properties and biological activities.

However, chemical reactions are unquestionably essential in chemical contexts, as they encompass rules governing atomic and functional group interactions. Presumably, pre-training exclusively on single-molecule datasets allows the model to learn SMILES syntax and valency rules, yet it is likely insufficient for acquiring comprehensive knowledge of chemical systems. Moreover, when mask-filling parts of single molecules, diverse substituents and carbon chain structures in molecules create multiple plausible answers, thereby complicating model optimization and hindering the model's ability to learn definitive representations. The limitations of LLMs pre-trained on single molecules and the success of LLMs on precursor prediction thus inspire us to move pre-training forward to chemical reactions.

PMSR[18] introduced a pre-training approach incorporating reaction information. While offering insights into chemical knowledge from reactions, the specificity of the pre-training objective may limit generalization. Evaluation is currently limited to precursor prediction, leaving broader applicability untested.

Recently, RSGPT[19], a retrosynthesis model built upon LLaMA2, generated over 1 billion reaction data points. Using pre-training, RLAIF, and fine-tuning, it achieves a state-of-the-art Top-1 accuracy of 63.4% with 3.2 billion parameters, significantly outperforming prior models.

Some works in the literature have attempted to expand retrosynthesis tasks beyond

single precursor prediction. For example, the Triple Transformers Loop (TTL)[20] model employs separate Transformer models to predict precursors, reagents and products separately, with product prediction validating the others. This approach provided an important reference for our research; however, it required training distinct models for each task.
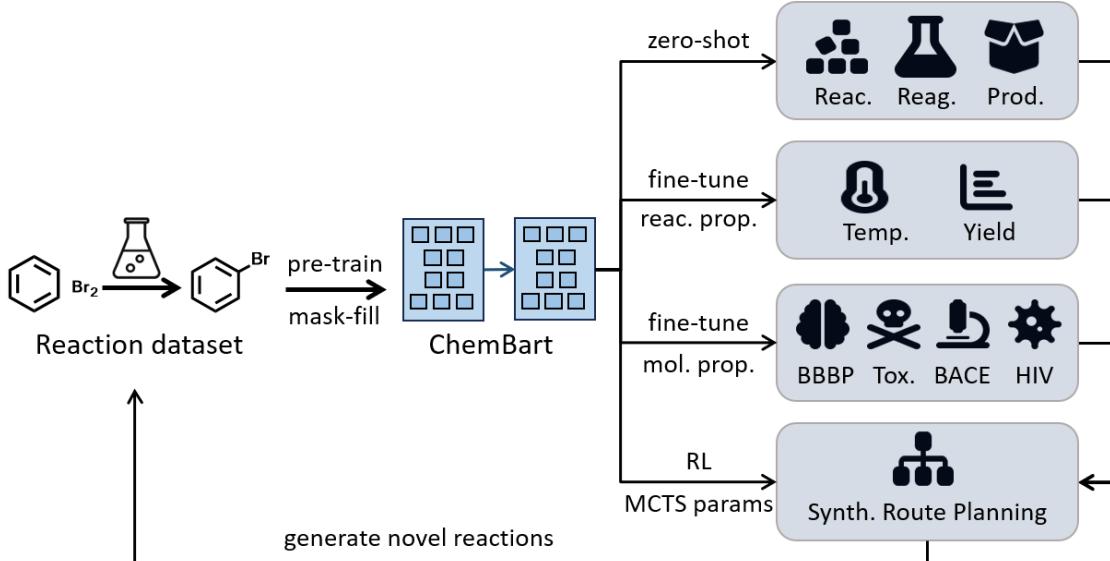


**Fig. 1 | Overview of ChemBART, a pre-trained language model from mask-filling of chemical reaction data.** The model predicts reaction precursors and reagents directly based on pre-trained parameters. Upon fine-tuning with additional data, the model can also be applied to predict detailed chemical reaction-related information and molecular-level information, further aiding chemical synthesis. By combining reinforcement learning (RL) with Monte Carlo tree search (MCTS), ChemBART can perform multi-step synthesis planning and provide feedback on experimental data.

Our goal is to achieve simultaneous multi-task learning using a single pre-trained model on chemical reaction data, thereby achieving information gains, generalization, and resource savings. Here, we introduce ChemBART, an open-source pre-trained SMILES language model. It is designed for predicting chemical reaction pathways and various chemical tasks using a unified framework pre-trained on reaction SMILES. Unlike current pre-trained SMILES LLMs, ChemBART can handle diverse chemical tasks and integrate them for full synthesis planning (**Fig. 1**). Specifically, we used mask-filling of precursors, reagents, and products in reactions as the pre-training task, enabling single-step reaction predictions. Building on our ReSynZ algorithm[21], we incorporated Monte Carlo tree search (MCTS) [22] and reinforcement learning to allow multi-step planning by fine-tuning the model for policy and value functions. We further fine-tuned it for reaction-relevant data, including temperature, yield, and molecular properties like blood-brain barrier penetration (BBBP), aiding synthetic path planning.

After 7 rounds of pre-training (**Supplementary Figs. S1-S2**), ChemBART achieved convergence, resulting in 84.7% top-10 accuracy in precursor prediction and 85.8% in reagent prediction. For reaction data, it had a 10% mean absolute error (MAE) for temperature and 23% for yield across categories, with an 8% MAE for yield on the Suzuki-Miyaura cross-coupling dataset[23]. In biochemical property prediction,

ChemBART outperformed state-of-the-art models on four out of five tasks. For MCTS parameter fitting, ChemBART showed improved performance over traditional "template-based" networks. Perhaps most compellingly, ChemBART-designed multi-step synthesis routes and reaction conditions were validated experimentally, confirming shorter pathways with over 20% yield improvement compared to established literature benchmarks. These results highlight ChemBART's effectiveness and potential utility.

## Results
### 1. Pre-training and fine-tuning workflow

ChemBART is a BART-based language model pre-trained on chemical reactions, each represented as a SMILES sentence ("reactant > reagent > product"). As illustrated in **Supplementary Tables. S1-S3**, tokenization is performed at the atom/mapping index level, with masked components reconstructed during training. This mask-filling task captures SMILES syntax and chemical principles, providing versatile representations for diverse chemical tasks. We utilized USPTO-full[24] and USPTO-MIT[25] as the pre-training data source. Hereafter, we denote ChemBART pre-trained on USPTO-full as ChemBART-F, the one pre-trained on USPTO-MIT as ChemBART-M, and the randomly initialized model as ChemBART-R.

After pre-training, ChemBART is fine-tuned for multiple downstream tasks using task-specific tokens and heads, supporting reaction temperature and yield regression, molecular property classification, and MCTS parameter estimation. The integrated model enables multi-step retrosynthetic route planning and associated reaction implementation information prediction, offering a unified framework for a broad range of chemistry applications.

### 2. Interpretability of ChemBART.

After pre-training, ChemBART embeds chemical knowledge within its parameters. To illuminate its learning process, we analyze the embedding vectors and attention matrices, revealing connections between LLM computation representations and fundamental chemical properties.

**Interpretability of embedding token vector.** To interpret the chemical meaning encoded in the token vectors from ChemBART's embedding layer, we calculated the Euclidean distances between vectors representing specific element tokens (Cl, Br, O, N, C, Na, K).

The model captured chemical concepts including electronegativity and bond covalency within its embeddings, where smaller distances (darker colors in **Fig. 2a**) denote higher similarity. Intra-group pairs (Cl/Br: 1.009, Na/K: 0.906) were more similar than cross-group pairs (e.g., Cl/Na: 1.064), reflecting electronegativity trends. Similarly, O/Cl similarity (0.993) exceeded O/Br (1.080), consistent with their higher electronegativity and anion basicity. C/O/N formed a distinct cluster reflecting covalent bonding, unlike the ionic-leaning Cl/Br/Na/K group. The lower C/Cl distance versus C/Br further supports stronger C-Cl covalency. A 3D visualization, using two PCA-reduced dimensions and atomic number as the third axis, further indicates that ChemBART learns periodic element trends (**Supplementary Fig. S3**).

Model interpretability was further validated via alcohol oxidation experiments. Dimensionality reduction of end token outputs (**Fig. 2b**) revealed distinct mechanistic patterns, showing consistent vectors for alcohol to aldehyde and aldehyde to acid oxidations. This confirms ChemBART's understanding of SMILES syntax and chemical principles, effectively encoding chemical laws during pre-training.

**Interpretability of the attention matrix.** We further investigated ChemBART's reaction understanding using attention mechanisms, exemplified by a substitution reaction between a Grignard reaction **1** and a bromo compound **2**. After product prediction via masking, the encoder's last-layer attention matrix (**Fig. 2c**) revealed strong $C^-{\rightarrow}C\text{-}Br$ interaction in the alkyl bromide, indicating recognition of the reactive C-Br bond. Enhanced $C^-{\rightarrow}Mg^{2+}$ attention reflected charge attraction. Attention decay along the Br-originating alkyl chain suggested diminishing electronic influence. In contrast, attention of $C^-$ to C atoms in the chain in the Grignard reagent increased with chain length from C4 to C1, possibly due to steric effects. These patterns demonstrate capture of reaction mechanisms. Strong local token interactions, particularly diagonal ones, indicated neighboring atomic group influences, also seen in cross-attention (**Supplementary Figs. S4-S7**).
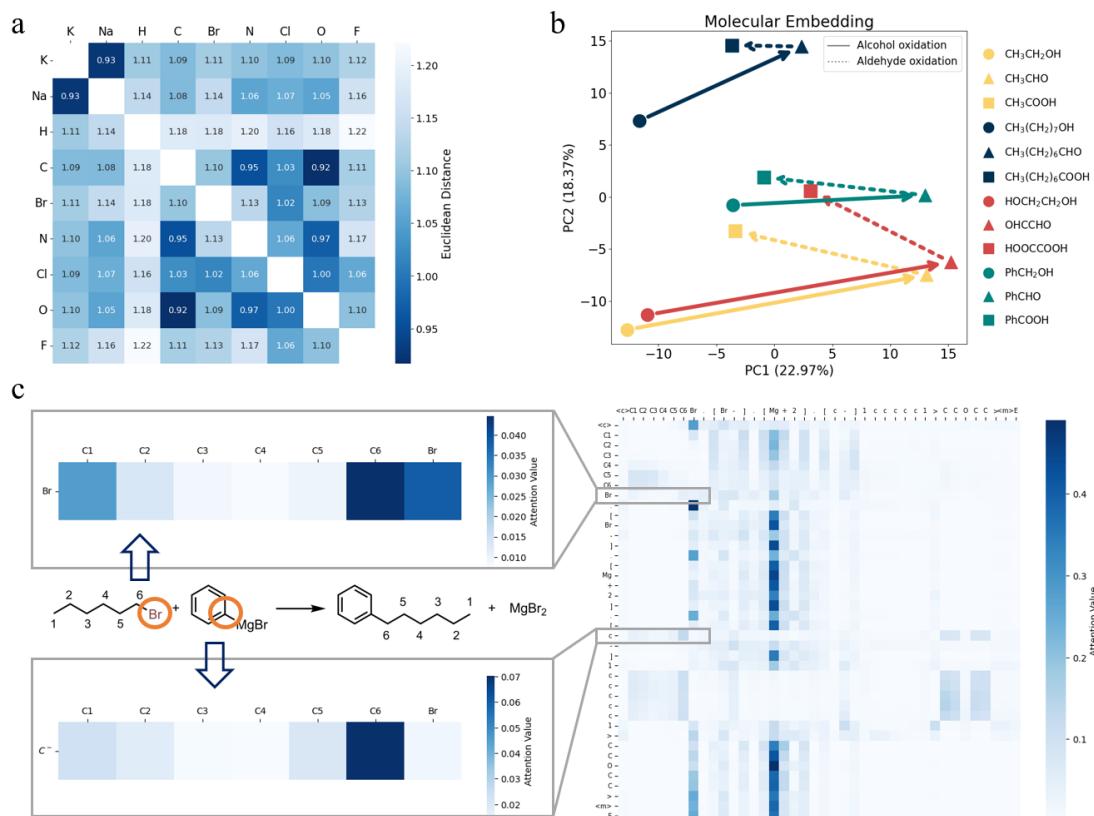


**Fig. 2 | Interpretability analysis of the ChemBART model.** (a) Euclidean distances between characteristic element representations in the embedding layer. (b) 2D projections of the reaction pathway. Alcohols (circle), aldehydes (triangle), and acids (square) of the same series are represented by the same color; arrows indicate oxidation pathways. (c) The encoder's last-layer attention matrix during masked product prediction for a Grignard reaction.

**3. Single-step retrosynthesis prediction.**

Our pre-training task of mask-filling reactions directly enables retrosynthesis precursor and reagent prediction. We employed a beam search strategy to give multiple predictions. We calculated top-$k$ accuracy (accurate results in top-$k$ predictions) and syntax error rate (invalid SMILES in top-$k$ predictions).

**Table 1. Top-$k$ accuracies and syntax error rates for the prediction of precursors and reagents. The unit for the parameter size of models is million.**

| Criterion | Task | Model | Parameter size | Top-1 | Top-3 | Top-5 | Test set |
|---|---|---|---|---|---|---|---|
| Top-k Accuracy | Precursor Prediction | ChemBART-F | 350 | 45.0 | 62.6 | 68.8 | USPTO-full |
| | | PMSR (8-layer)[18] | 110 | 45.5 | 60.9 | 65.5 | |
| | | RSGPT[26] | 3228 | 59.2 | 74.2 | 78.2 | |
| | | ChemBART-M | 350 | 56.8 | 74.3 | 80.0 | USPTO-MIT |
| | | Segler-Coley retrained[27] | >110 | 47.8 | 67.6 | 74.1 | |
| | | AutoSynRoute - character[6] | 44.4 | 54.1 | 71.8 | 76.9 | |
| | | RSGPT[26] | 3228 | 63.9 | 86.1 | 88.7 | |
| | | ReactSeq[28] | 17.4 | 58.9 | 80.5 | 86.4 | USPTO-50k |
| | | Karpov Transformer[29] | 1.88 | 42.7 | 63.9 | 69.8 | |
| | | AutoSynRoute - token[6] | 44.4 | 42.0 | 64.0 | 71.3 | |
| | | Chemformer[10] | 230 | 54.3 | - | 62.3 | |
| | | T5Chem[14] | 6.67 | 46.5 | 64.4 | 70.5 | |
| | | PMSR (8-layer)[18] | 110 | 67.1 | 82.1 | 85.2 | |
| | | RSGPT[26] | 3228 | 63.4 | 84.2 | 89.2 | |
| | Reagent Prediction | ChemBART-F | 350 | 54.5 | 67.8 | 74.4 | USPTO-full |
| | | ChemBART-M | 350 | 50.3 | 65.2 | 71.0 | USPTO-MIT |
| Syntax Error Rate | Precursor Prediction | ChemBART-F | 350 | 0.388 | 0.657 | 0.869 | USPTO-full |
| | | ChemBART-M | 350 | 0.204 | 0.353 | 0.494 | USPTO-MIT |
| | Reagent Prediction | ChemBART-F | 350 | 0.0243 | 0.0179 | 0.0151 | USPTO-full |
| | | ChemBART-M | 350 | 0.0102 | 0.0171 | 0.0212 | USPTO-MIT |

As shown in **Table 1**, our model achieves performance comparable to state-of-the-art template-free tools, (e.g., AutoSynRoute[6]) and other pre-trained SMILES language models (e.g., Chemformer,[10] T5Chem,[14] PMSR[18] and ReactSeq[17]). These findings validate the effectiveness of the ChemBART model, offering three key advantages. Its multi-task pre-training concurrently predicts precursors, reagents, and products, enhancing chemical reasoning via implicit task correlation and tripling training data. Unlike precursor-focused template-free models, ChemBART's simplified design improves generalization and allows fine-tuning across various tasks. Low overall syntax error rates, with minimal impact on generation, indicate robust SMILES syntax acquisition. Reagent predictions show near-zero errors, likely due to their simpler structures and higher training frequency.

## 4. ChemBART fine-tuning: predicting reaction features and molecular properties.

For classification and regression tasks, BART[4] adds special task tokens to the input at the end due to its decoder's masked attention; a linear layer then maps each token's output to the target label. This method, leveraging the model's pre-trained capabilities, is believed to capture overall input information. ChemBART follows this practice. We fine-tuned it on reaction temperature/yield prediction and biochemical property classification.

**Prediction of reaction temperature and yield.** We evaluated ChemBART-F, ChemBART-M, and ChemBART-R for predicting reaction temperature and yield on the ORD dataset[30]. ChemBART-M showed the best individual performance for property prediction (**Supplementary Tables S4-S5**). We therefore fine-tuned ChemBART-M with task-specific tokens for joint temperature and yield regression, selecting the checkpoint with the lowest validation root mean squared error (RMSE) to mitigate overfitting.

The model showed robust temperature prediction (MAE 10% for -150 to 250°C, **Fig. 3a**) and yield estimation (MAE 23% for 0–110%, **Fig. 3b**). For the Suzuki-Miyaura dataset[23], our model yielded the $R^2$ values ranging from 0.77 to 0.82 across 10 random splits (**Fig. 3c & Supplementary Table S6**), rivaling the state-of-the-art rxnfp model designed explicitly for this dataset $(0.81 \pm 0.01)$[31].

Our exploration of fine-tuning strategies confirmed that appending task tokens at the reaction's end and using a linear layer on the task token output yields optimal performance. For regression, the BART model's encoder-decoder combination outperformed the Bert model's encoder-only approach, likely due to more parameters. Pre-training significantly boosted regression performance, validating our design.
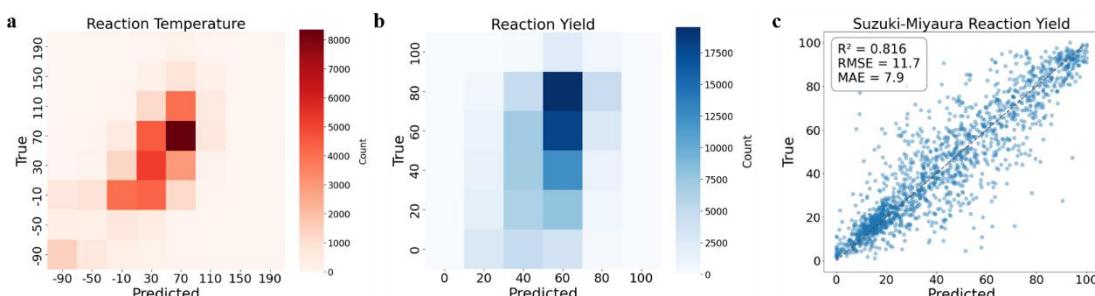


**Fig. 3 | Performance of joint regression of temperature and yield.** (a) Temperature prediction ($R^2$=0.55, MAE=10%), (b) Yield prediction ($R^2$=0.16, MAE=23%), with heatmap intensity reflecting data density per 40°C/10% interval. (c) Yield prediction on Suzuki-Miyaura dataset ($R^2$=0.816, MAE=7.9%) using one data split.

**Prediction of molecular properties.** For molecular property prediction, we used the ChemBART-M model with an added binary classification layer. We fine-tuned it on five benchmark datasets, BBBP (bioactivity, 2,053), HIV (antiviral, 41,127), BACE (binding affinity, 1,514), Tox21 (toxicology, 8,014), and ClinTox (clinical toxicity, 1,484)[32], using an 8:1:1 split. As shown in **Table 2**, ChemBART achieved competitive

performance against state-of-the-art methods when evaluated using the area under the curve (AUC). This suggests that chemical reaction pre-training effectively captures organic molecule interaction principles, validating our paradigm's generalization across biochemical tasks. We also conducted LoRA (Low-Rank Adaptation) fine-tuning [33] in addition to full-parameter fine-tuning, which showed further improvement, likely due to its regularization effect against overfitting.

Visualizing pre-training embeddings showed poor separation (**Supplementary Figs. 8a and 8c**), due to the lack of property supervision. Fine-tuning, however, induced apparent clustering (**Supplementary Figs. 8b and 8d**), demonstrating the distillation of task-specific features from the pre-trained model.

**Table 2. AUC index for ChemBART and other models for classification tasks on molecular properties.**

| | Datasets | | | | |
|---|---|---|---|---|---|
| | **BBBP** | **HIV** | **BACE** | **Tox21** | **ClinTox** |
| SMILES-Transformers[34] | 0.704 | 0.729 | 0.701 | 0.802 | 0.954 |
| ECFP4[35] | 0.729 | 0.792 | 0.867 | 0.822 | 0.799 |
| GraphConv[36] | 0.690 | 0.763 | 0.793 | 0.829 | 0.807 |
| Weave[37] | 0.671 | 0.703 | 0.806 | 0.820 | 0.832 |
| ChemBERTa[38] | 0.643 | 0.622 | - | 0.728 | 0.733 |
| D-MPNN[39] | 0.708 | 0.752 | - | 0.688 | 0.906 |
| CDDD[40] | 0.761 | 0.753 | 0.833 | - | - |
| MolBERT[13] | 0.762 | 0.783 | 0.866 | - | - |
| GraphCL[41] | 0.685±0.005 | 0.776±0.009 | 0.782±0.012 | 0.754±0.009 | 0.701±0.019 |
| GraphLoG[42] | 0.725±0.008 | 0.778±0.008 | 0.835±0.012 | 0.757±0.005 | 0.767±0.033 |
| Mol2vec[43] | 0.872±0.021 | 0.769±0.021 | 0.862±0.027 | 0.803±0.041 | 0.841±0.062 |
| MolR-GCN[44] | 0.890±0.032 | 0.802±0.024 | 0.882±0.019 | 0.818±0.023 | 0.916±0.039 |
| MolR-GAT[44] | 0.887±0.026 | 0.794±0.022 | 0.863±0.026 | 0.839±0.039 | 0.908±0.039 |
| MolR-SAGE[44] | 0.879±0.032 | 0.793±0.026 | 0.859±0.029 | 0.811±0.039 | 0.890±0.058 |
| MolR-TAG[44] | 0.895±0.031 | 0.801±0.023 | 0.875±0.023 | 0.820±0.028 | 0.913±0.043 |
| Xmol | 0.962 | 0.798 | 0.872 | - | 0.984 |
| ChemBART-M | 0.910±0.015 | 0.809±0.001 | 0.881±0.029 | 0.844±0.012 | 0.866±0.024 |
| ChemBART-R | 0.839±0.052 | 0.752±0.004 | 0.780±0.047 | 0.738±0.019 | 0.639±0.074 |
| ChemBART-LoRA | 0.915±0.012 | 0.789±0.009 | 0.872±0.015 | 0.830±0.014 | 0.920±0.016 |

## 5. Synthetic path planning

In previous work, we developed ReSynZ [21], which integrates reinforcement learning and MCTS methods for template-based chemical retrosynthetic analysis. Based on the predictive capability of ChemBART for single-step reaction, we will next leverage it to fit the policy and value functions. By combining these functions with MCTS, we aim to achieve efficient planning of multi-step synthesis pathways.

Using 12,000 MCTS-generated data points from ReSynZ[21], we evaluated three models: the template-based neural network from ReSynZ[21], the pre-trained ChemBARTF/M, and the randomly initialized ChemBART-R. Model performance was quantified by mean RMSE. Training dynamics (**Supplementary Figs. S9-S10**) show

ChemBART achieved balanced policy and value prediction, outperforming ReSynZ's poor policy generalization. Seemingly, pre-training enhanced multi-task coordination without overfitting as in the case of template-based networks[21,45-48]. Even without pre-training, ChemBART-R outperforms traditional template-based neural networks in predicting *p*. Therefore, we saved model parameters only when validation RMSE for both policy and value hit new minima, with final test results in **Supplementary Table S7**.
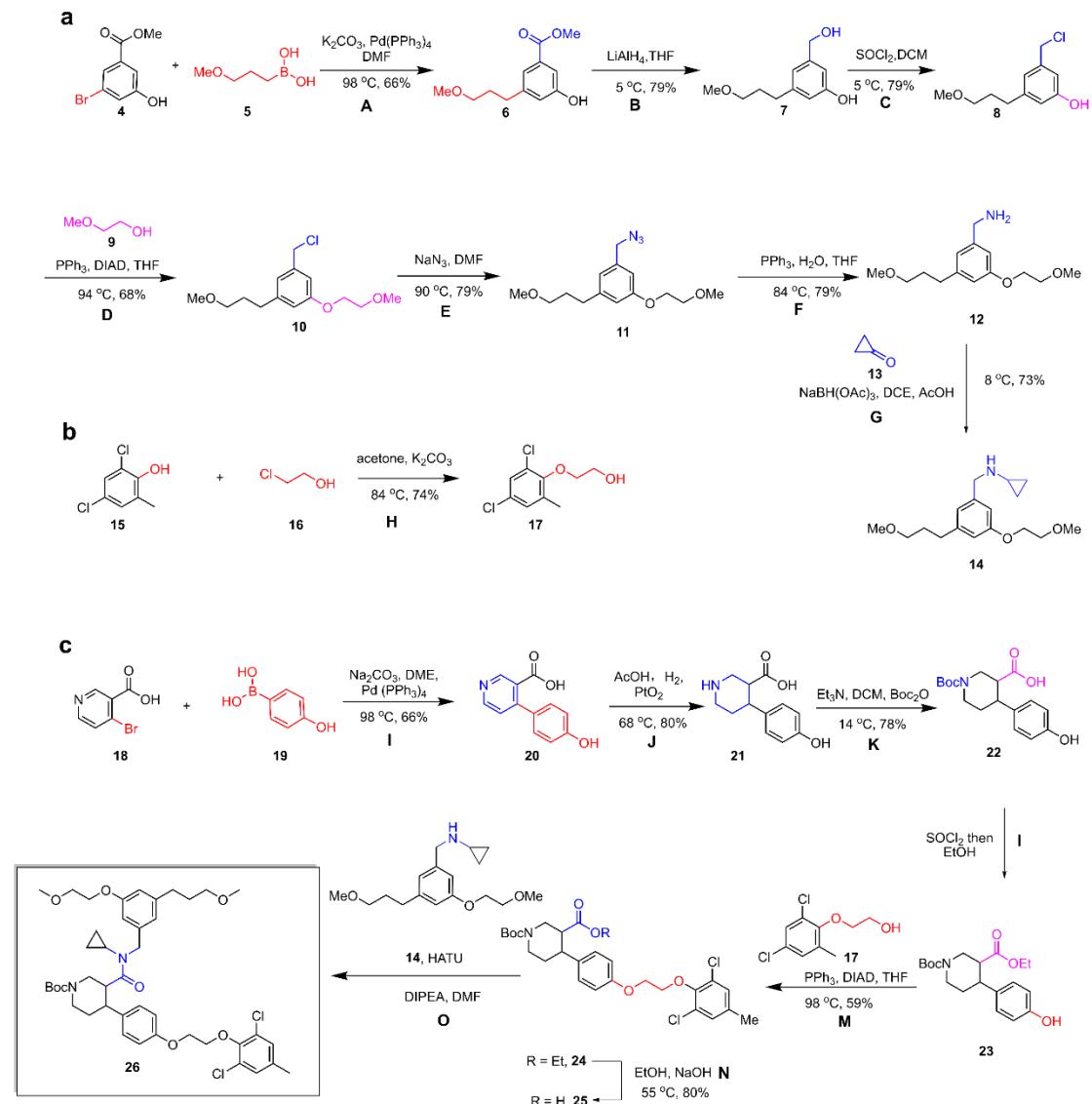


**Fig. 4 | A 15-step synthetic pathway for compound 26**[49] **designed by ChemBART.** Panels (a), (b), and (c) show the planned paths for intermediates **14** and **17**, and the target molecule **26**, respectively.

We then evaluated its predictive performance for multistep synthetic routes across various datasets. At the node of the search tree, ChemBART employs beam search (with num_sample=50 throughout this work) to predict potential precursors and their corresponding reagents. It scores these candidates for chemical feasibility using probabilities from the pre-trained model and synthetic difficulty through the policy and

value functions of the fine-tuned model.

ChemBART-M/F, a template-free method, achieved success rates of 64.9% and 70.1% on the Retro*-190 (also known as USPTO-190) benchmark[45] (**Supplementary Table S8**). Impressively, to the best of our knowledge, the success rate of multi-step synthetic route design via template-free Transformer-based retrosynthesis remains unreported. Our work achieves success rates nearly matching template-based methods, alongside consideration for route quality, underscoring the power of ChemBART's pre-training and fine-tuning strategy. **Fig. 4** presents a representative example selected from this dataset. In this case, ChemBART planned a 15-step synthetic route to compound **26** (PubChem CID: 86639854, ID 40 in **Supplementary Table S8**), a designed renin inhibitor[49], with detailed reaction conditions and yields. This example highlights ChemBART's capability for long-path planning, showcasing innovative strategies such as the reductive selectivity (step-J in **Supplementary Fig. S11**), which could potentially inspire new pathways. In other cases, ChemBART adopted the same bond-breaking strategy as patents but shortened pathways via reagent selection.

To further assess predictive performance on contemporary organic molecules, we created a novel multi-step retrosynthesis dataset containing 53 highly diverse and distinct molecules from 30 recent articles in the *Journal of Medicinal Chemistry* in 2025 (denoted as JMC2025 dataset, **Supplementary Figs. S12-S13, & Table S9**). Using this dataset, we systematically investigated different generation strategies. Beam search, top-$k$ ($k$=10), and top-$p$ ($p$=0.9) achieved success rates of 88.68%, 84.91%, and 83.02%, respectively (**Supplementary Fig. S14**). Evidently, beam search tends to produce high-confidence yet chemically conservative routes with limited reaction diversity (e.g., ChemBART paths matching literature examples), while top-$k$ and top-$p$ sampling may introduce greater synthetic novelty at the cost of increased risk of hallucination. Seemingly, balancing novelty and reliability remains a challenge.

## 6. Wet-lab validation

To further evaluate the applicability of ChemBART in multistep synthetic planning, we selected a representative compound **P1** (ID 3 in **Supplementary Table S8**) from the JMC2025 dataset for wet-lab validation (**Supplementary Figs. S15-S22**). **P1** was designed as a PD-L1/VISTA dual inhibitor and had been previously synthesized in six steps with an overall yield of 6.5% (**Fig. 5a**)[50].

A more convergent route was designed by ChemBART featuring the late-stage coupling of amine **35** and alcohol **37** (**Fig. 5b**). Following this route, amine **35** was obtained under Suzuki coupling conditions in 70% yield. Meanwhile, ester **28** was reduced with DABAL-H to give aldehyde **36**, which subsequently underwent reductive amination to afford alcohol **37**. Finally, the $S_N$Ar reaction between amine **35** and alcohol **37** was achieved under acidic conditions (HCl), albeit with a yield of only 10%. After screening various conditions (see supporting information), we identified *p*-TsOH as the optimal acid, achieving a substitution efficiency of 67%. Overall, the ChemBART-designed route delivered **P1** in an overall yield of 35%, representing an improvement of 28.5% compared with the original route.

Notably, most reaction conditions proposed by ChemBART were readily applicable

in the laboratory, and the incorporation of expert knowledge and further optimization enabled us to achieve the best overall performance, underscoring both the practical utility and the inspirational value of ChemBART's designs.
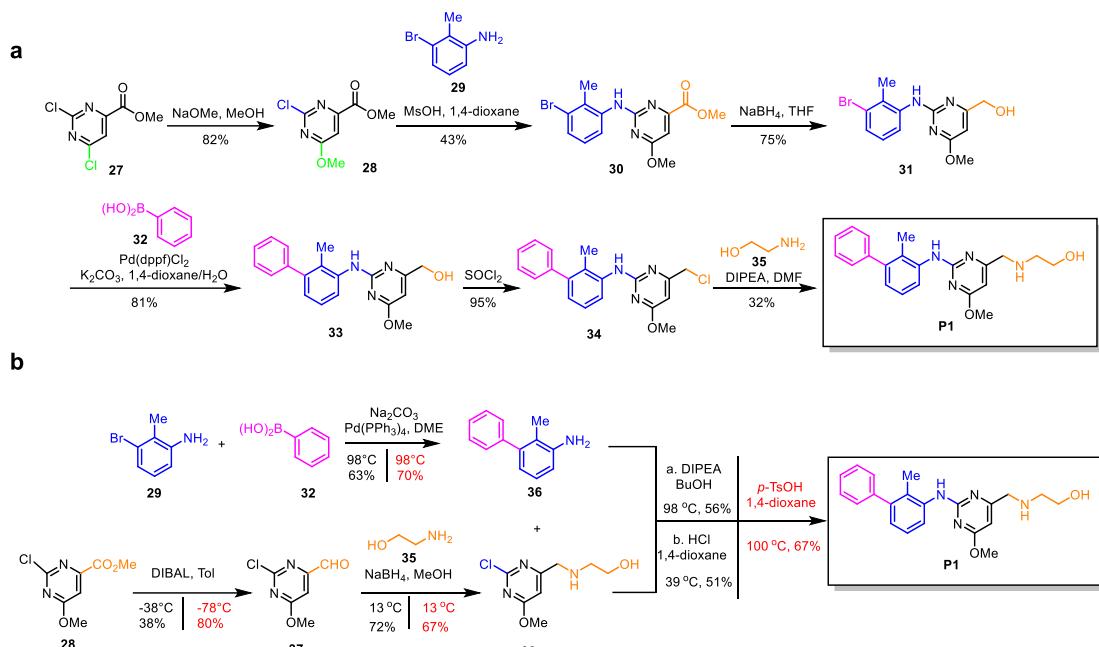


**Fig. 5 | Comparison of synthetic routes for compound P1: Literature-reported vs. ChemBART-predicted and experimentally validated.** (a) The literature-reported experimental route[50]. (b) The ChemBART-predicted synthetic route, including predicted reaction conditions and yields, alongside experimental validation results. Within panel (b), optimized reaction conditions and corresponding yields, determined experimentally, are highlighted in red.

## Discussion

**The mask-filling of reaction components for pre-training is effective for the model to learn general chemical knowledge.** Various experiments have demonstrated that abundant chemical knowledge captured within the pre-trained parameters significantly enhances performance across various downstream fine-tuning tasks[10-17]. Our analysis reveals that during pre-training, model embedding and attention layers acquire SMILES syntax and chemical knowledge, embedding this directly into parameters and thus improving fine-tuning outcomes on pre-trained models.

However, recent SMILES LLMs[10-17] often pre-train on individual molecules using tasks like mask-filling or synonym conversion. These methods primarily capture SMILES syntax, failing to deeply embed essential reaction knowledge needed for molecular interconversion. Furthermore, the variability within molecules complicates conventional pre-training tasks (e.g., predicting masked segments), as they lack definitive answers, hindering the learning of consistent chemical rules.

In contrast, ChemBART employs a novel three-directional pre-training task with exact answers. This allows the model to learn clear chemical laws, yielding robust pre-training results effectively. Unlike single-molecule LLMs limited to specific tasks,

ChemBART can tackle more versatile chemical challenges and integrate these capabilities to accomplish the full synthesis planning cycle, highlighting the advantage of reaction-level pre-training.

**A single copy of the model parameters can be used for multiple tasks, saving computational resources.** ChemBART can utilize a single set of parameters for both multi-task generation and regression. Its pre-training design allows the learned parameters to be directly applied across diverse generation tasks—predicting precursors, reagents, and products—and regression tasks—optimizing temperature and yield. Additionally, parameter fitting provides policy and value estimates crucial for synthesis planning. This multi-task capability is particularly economical, as it avoids the need to load and maintain separate models, making efficient use of GPU resources, a critical advantage given the high computational demands of large language models.

**Directions to be further studied.** While our model demonstrates strength in reaction-related tasks, its performance in predicting molecular physical-chemical properties, such as activation energy on the QM9 dataset[51], remains limited, suggesting sequence information alone is insufficient. Potential improvements include integrating 3D structural data via graph neural networks, combining it with the language model output, as done by KFLM2[52], and incorporating sequence attention and structural contact maps, inspired by Uni-Mol[53], to better link sequence and structure representations is another option.

We have also observed efforts to adapt general-purpose LLMs pre-trained on natural languages, such as LLaMA[54] and Qwen[55], using both SMILES strings and chemistry literature[16,52,56-60], or just simply prompt these LLMs using their general knowledge.[61,62] While their advantage over SMILES-only models is yet to be clearly demonstrated, their potential to leverage vast general knowledge and learn from natural language descriptions holds significant promise for enhancing downstream performance and enabling broader application through natural language input.

Further reinforcement learning optimization, including additional MCTS cycles and retraining, is planned to improve performance. We also aim to implement conditional precursor generation, guided by reaction type and site, to enhance precursor diversity and tree search flexibility, ultimately boosting the success rate of multi-step synthesis planning.

# Method

**Pre-training of the model.**

**Definition of chemical sentence and vocabulary list.** We define a complete chemical reaction "reactant > reagent > product" as a sentence in the SMILES[5] language. Then, tokenization and index conversion are required to transform a chemical reaction into a digital representation readable by the language model. In tokenization, the SMILES sentence is segmented into sequences of tokens according to a predefined word list (**Supplementary Tables S1-S3**). In this study, we take a complete atom or atom mapping index as a token. Although AutoSynRoute[6] pointed out after experiments that there is no significant difference in the training effectiveness of taking atoms or separate characters as tokens, we choose to use complete atoms or atom mapping indices as the token units to improve interpretability and reduce computational overhead. Given that the SMARTS[63] language serves as an extension of SMILES, the word list has been designed to accommodate SMARTS syntax as well. This compatibility enables researchers to extend the application of this pre-trained model to address challenges expressed within the SMARTS language framework. In index conversion, each token in the resulting sequence is then mapped to its corresponding index based on the word list.

After passing through the embedding layer of LLM, each token index is transformed into a $d$-dimensional vector, composed of the sum of a lexical vector and a positional vector. In our problem, the lexical vector encodes the chemical semantics of the token. For example, if the token represents a chemical element, its lexical vector captures general chemical properties associated with that element. Although such high-dimensional representations are not directly readable to humans, the vectors corresponding to different tokens are comparable. Specifically, the distance between two representation vectors reflects the similarity of their chemical properties. To quantify this similarity, we used the Euclidean distance $D(x,y)$ between the two vectors, $x = [x_1, x_2\ldots, x_d]$, $y = [y_1, y_2\ldots, y_d]$, where a smaller Euclidean distance indicates greater similarity between the two vectors.

$$D(\mathrm{x}, \mathrm{y}) = \sqrt{\sum_{i=1}^{d} (x_i - y_i)^2}$$

In our setting, dimension $d$ is set to 1024.

**Model structure.** Our language model adopts the framework of the BART-large model,[4] featuring approximately 0.4 billion parameters. Each token is represented by a 1024-dimensional vector, and the model accommodates a maximum input length of 1024 tokens. The attention mechanism comprises 16 attention heads per layer, with the architecture structured into 12 encoder layers and 12 decoder layers. The inclusion of both encoder and decoder components is specifically suited to addressing the diverse requirements in chemical tasks, as the encoder is suited for understanding the semantic meaning, which is the molecular structure in our problem, while the decoder is suited for the generation of molecular structures. Moreover, the decoder can be extended to

support classification and regression tasks by integrating specific heads. In addition, the trained encoder can be extracted separately to build a BERT model (referred to as ChemBert in the following context) for dedicated labeling, classification and regression.

As a key component of Transformers,[1] the attention mechanism enables the model to perform far beyond earlier models dealing with sequential problems, such as RNN[64] and LSTM[65]. The basic algorithm of the attention mechanism is that, for the hidden layer representation of a sentence $X$ input into the attention module,

$$X = [T_1, T_2, \ldots T_n]^T, \text{where } T_i = [x_1, x_2, \ldots x_d], x_i \in R$$

where $T_i$ is the representation of token $i$ in this hidden layer and $n$ is the number of tokens in the sentence, an attention matrix $A$ is calculated by

$$K = XM_1$$
$$Q = XM_2$$
$$A = \text{softmax}(\frac{QK^T}{\sqrt{n}})$$

where $M_1$, $M_2$ are the matrices with trainable parameters, and the value of $A_{ij}$ represents the attention of token $i$ to $j$, i.e., the degree of relevance of $j$ for $i$. The representation of token $i$ in the new round output by the module, $T_i'$, is a weighted summation of the vector representations of all the tokens concerning the attention of $i$ to each

$$V = XM_3 = [V_1, V_2, \ldots V_n]^T$$
$$T_i' = \sum_{j=1}^{n} A_{ij} * V_j$$

where $M_3$ is the matrix with trainable parameters, $V_i$ still represents the information of the token $i$, which is linearly transformed on the basis of $T_i$. Therefore, in the next layer, the hidden layer representation of the input sentence $X'$ is

$$X' = AV$$

**Mask-filling pre-training.** In our framework, each "reaction sentence" consists of a triple of reactant, reagent, and product, any of which can be masked. Based on chemical principles, the masked component remains predictable given the remaining parts. Thus, the pre-training task in this study is designed such that an input reaction sentence with one masked component (reactant, reagent, or product) is processed by the decoder to reconstruct the original complete sentence. This approach facilitates the model's acquisition of both fundamental SMILES syntax and essential chemical reaction principles.

The model was pre-trained on the USPTO-full reaction dataset [24] (containing approximately 1.5 million reactions), primarily derived from medicinal chemistry synthesis patents. According to Schneider et al.,[66] the main reaction types in this dataset, ordered by frequency, include: heteroatom alkylation and arylation, acylation, protection/deprotection, C–C bond formation, redox reactions, heterocycle formation, and functional group interconversion or addition. As R-smiles[67] has recently reported that the increase in reaction data quality will significantly increase the performance of language models, we also conducted a pre-training on the USPTO-MIT dataset,[25] a filtered USPTO dataset excluding unreliable reaction records, and compared its performance with the model trained on USPTO-full. A 9:1 split was applied to separate

the training and test sets. The USPTO-50k dataset,[66,68,69], a much smaller data set with 50k pieces of data filtered from USPTO-full, was also used for testing.

The training was performed using the AdamW 70 optimizer with the Binary Cross Entropy loss function, a learning rate of 1e-5, a weight decay value of 1e-4, and a batch size of 256.
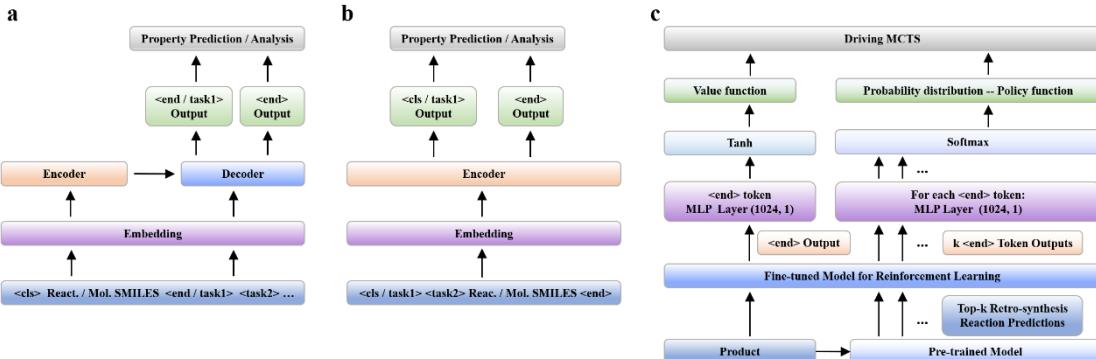
**Fine-tuning of the ChemBART model.**



**Fig. 6 | Frameworks for classification, regression, and generation tasks using Transformers.** (a) BART's classification/regression scheme: Task-specific tokens are appended to the input sequence. The complete sequence is processed by both the encoder and decoder. The decoder's output corresponding to each task token is passed through a linear layer to obtain the classification or regression label. (b) BERT's classification/regression scheme: Task-specific tokens are prepended to the input sequence. The encoder processes the entire sequence, and the output corresponding to each task token is passed through a linear layer to obtain the classification or regression label. (c) Fitting policy and value functions: A pre-trained model generates multiple possible precursors for a specific product. The value function $v$ is obtained by inputting the product's expression into a fine-tuned regression model. The policy function $p$ is derived by inputting each corresponding synthetic reaction for all precursor options into a fine-tuned regression model, then normalizing the output probabilities.

**Reaction temperature and yield prediction.** We employed our pre-trained model to analyze the temperature and yield data of organic reactions provided by the Open Reaction Database (ORD)[30]. The dataset comprises approximately 700k entries associated with reaction temperatures and yields. We partitioned the dataset into training, validation, and test sets with an 8:1:1 ratio.

To evaluate different fine-tuning approaches, we implemented multiple training schemes. Specifically, we compared the strategy of combining both encoder and decoder for regression tasks (following the BART model strategy) with a strategy that uses only the encoder (aligned with the BERT model strategy). Additionally, we investigated the potential benefits of joint regression across multiple tasks. This involved appending several task-specific tokens at the end of the reaction expression and applying a linear head for each token, to enhance information flow to improve task performance. These approaches are illustrated in **Fig. 6a & 6b**, respectively. Further, we examined whether incorporating a Sigmoid layer to constrain the output range of

the regression head could enhance prediction accuracy. We also compared the effectiveness of a Bi-LSTM layer, which integrates the output vectors of all tokens before feeding them into a linear head for regression, against the "one task token for one task" approach.

**Biochemical molecular property classification.** In this task, ChemBART is fine-tuned on five benchmark datasets: BBBP (bioactivity, 2,053), HIV (antiviral, 41,127), BACE (binding affinity, 1,514), Tox21 (toxicology, 8,014), and ClinTox (clinical toxicity, 1,484)[32], all structured as 8:1:1 training/validation/test splits.

We also explore the effectiveness of Low-Rank Adaptation (LoRA)[33] in fine-tuning on these datasets. LoRA offers the advantages of faster training and reduced memory consumption:

$$W' = W + \alpha BA$$

Where $W$ is the weight matrix of attention, frozen in LoRA fine-tuning, and $A \in R^{r*d}, B \in R^{d*r}$ are very small trainable weight matrices. For each task, two hyperparameters were tuned: the rank $r$ governs the model's capacity, and the scaling factor $\alpha$ controls the importance of the low-rank updates applied to the original weights.

We applied LoRA fine-tuning using the PEFT library from Hugging Face[71]. For the ClinTox dataset, we set 6 for both parameters, and for the other tasks, we used 8 for both, following hyperparameter optimization.

**MCTS Policy and Value Regression.** In the search tree of MCTS in our problem, a node is a single molecule with a value function. A node has multiple child nodes, which are choices of possible synthetic precursors. These child nodes have a probability distribution of being chosen, called the policy. Traditional multi-step synthesis approach like RetrosynthesisZero (ReSynZ)[21], Retro*[45], PDVN[46] and EG-MCTS[47] to fitting the policy and value function in MCTS was to use ECFP algorithm[72] to transform a molecular structure into a vector, and then use fully connected layers to yield policy and value with mlp_retrosyn library[48].

We propose a novel method, utilizing our pre-trained ChemBART model, to effectively fit the policy and value function in MCTS. We input the SMILES of the molecule at a node into ChemBART to fit the value function.

$$v(n) = v_\Theta(m)$$

where n denotes a node, m denotes the SMILES expression of the molecule at the node, and $\Theta$ denotes the fine-tuned model parameters.

forget. We then input each child node with its corresponding reagents and parent node into ChemBART to obtain a success score for the reaction. Next, we normalize the scores across all child nodes to derive the policy function.

$$p(n) = [p_1, \dots p_n], \text{ where } p_i = p_\Theta(r_i)$$

$$p(n) = \frac{p}{\sum_{i=1}^{n} p_i}$$

where $r_i$ denotes the SMILES expression of the $i^{th}$ candidate retrosynthesis reaction.

The algorithm is depicted in **Fig. 6c**. Compared with the traditional regression method of MCTS parameters, the embedding of Transformers serves a similar functionality to ECFP, converting molecule structure to a digital representation, while

the encoder and decoder of Transformers are also similar to the fully connected layers, extracting high-level information of molecular structures. Our experimental results have proved the effectiveness of our novel model structure on this task.

Since running MCTS with LLMs is computationally expensive, requiring hundreds of LLM calls per instance, we did not run the reinforcement learning algorithm of ReSynZ directly on ChemBART through iterative MCTS execution and retraining. Instead, we adopted a cold-start approach by leveraging the ReSynZ program that had finished reinforcement learning to generate MCTS data efficiently, with 6,000 data points from successful retrosynthesis routes and 6,000 from failed ones. The dataset was split into training, validation, and test sets in a 90:5:5 ratio. We then trained a template-based neural network from scratch, as well as ChemBART-F, -M and -R, respectively, using this dataset to compare performance. Subsequent reinforcement learning on ChemBART with MCTS yielded only marginal improvement in success rate (1~2%), indicating that the model had already absorbed most of the relevant knowledge from the ReSynZ-generated data.

**Single Step Generation.**

Our pre-training approach involves masking one of the three components—precursors, reagents, or products—within the reaction SMILES expression, prompting the model to output the complete reaction. This setup inherently allows for various predictive tasks, such as deriving precursors from products, inferring reagents based on both reactants and products, and predicting products given reactants and reagents, without requiring additional fine-tuning. We employed a beam search strategy for LLMs to generate k alternative answers for each input across these three tasks. In beam search, at each round of predicting the following tokens, only the generated sequences with top-k cumulative probability are retained.

**Synthesis Route Planning with MCTS and Reinforcement Learning.**

The MCTS[22] algorithm plays a key role in efficiently exploring potential synthetic pathways in large chemical reaction spaces. Each node in the search tree, representing a molecule, is initialized with two parameters: policy $p$ and value $v$, from neural networks. Each round of exploration tries to find a synthesis route, with each step decided by the parameters of the nodes. Through several times of exploration, the policy and value of each node would be gradually updated and improved to point to the optimal route. Finally, MCTS outputs the probability distribution $\Pi$ of all the child nodes of the root node, i.e. the target molecule, by the following formula

$$\Pi = \frac{[N_1^\tau, N_2^\tau, \dots]}{N^\tau}$$

where $N_i$ is the number of visits of the $i^{th}$ child node, $N$ is the total number of visits conducted, and $\tau$ is the temperature coefficient. The $\Pi$ can be seen as an improved policy $p$ of the root. If only the best route is needed, in each step the precursors with the highest probability in $\Pi$ are greedily picked. For alternative synthesis routes for a target molecule, we also adopted beam search to pick the top-$k$ routes in each step.

To drive MCTS, the algorithms to generate all the child nodes, as well as the policy and value for each node are required.

**Generating candidate child nodes for a parent node.** We used beam search to

generate top-k possible synthetic precursor candidates. To ensure the quality of the overall synthesis pathway, we filtered out unreliable predictions at each step. Using the RDKit[73] module, we canonicalized the SMILES of all generated child nodes and discard any that cannot be canonicalized, which indicates structurally invalid predictions. We then deduplicated child nodes that share the same canonical SMILES, summing their associated probabilities. Next, we used ChemBART to predict the reagents for each candidate retrosynthesis reaction, which are essential implementation details, and subsequently predict the product from each precursor-reagent pair. If the predicted product does not match the original parent node, the predicted reaction is considered incorrect and the corresponding child node is discarded. Finally, we normalized the probabilities of the remaining child nodes to get a probability distribution $p_1$.

**Predicting value and policy functions.** As described in Section 2.2, we obtained the value function by inputting the target molecule at a node into our fine-tuned model. The policy $p_2$ is derived by inputting each filtered candidate reaction at that node into the model to obtain a score, then normalized. However, the resulting policy is derived purely from reinforcement learning, which optimizes only the likelihood of successfully completing a synthesis pathway, without accounting for the practical feasibility of the individual reactions, which can only come from literature. To address this, we incorporated the generation probability $p_1$ reflecting confidence based on literature, with the policy from the fine-tuned ChemBART model reflecting the ease of achieving complete synthesis, by multiplying them elementwise and then normalizing the resulting probability to yield the final policy $p^*$.

$$p^* = (p_1 \odot p_2)$$
$$p^* = \frac{p^*}{\sum_{i=1}^{n} p_i^*}$$

It should also be noted that the methodology of combining the knowledge from literature and exploration is also the popular trend in recent CASP works like PDVN[46], Retro*+[74], EG-MCTS[47] and so on.

To optimize computing efficiency, a separate process is assigned to each target molecule while planning synthesis pathways, while ChemBART parameters (pre-trained and fine-tuned) on GPUs are shared across these processes.

After a complete synthesis route is designed, we used other fine-tuned ChemBART models stated above to add additional information to each step of the route, such as reaction temperature, reaction yield and molecular property of the molecules involved.

# Acknowledgements

Science Foundation through the XPLORER PRIZE.

## Author Contributions
K.N.L and H.D. conceptualized the work and designed the methodology. K.N.L, Y.J.Z., and H.P.G. performed the computational studies, J.W., Z.Y.S., and X.G.L. performed wet-lab validations. All authors discussed the data and contributed to the drafting. H.D. supervised the project and acquired funding.

## Competing Interests
The authors declare no competing financial interests.

## Additional information
**Supplementary information.** The online version contains supplementary material available at https://doi.org/XXXXXX
**Correspondence and requests for materials** should be addressed to Hao Dong.

## Data Availability
All the data used in this paper are from open-source datasets. Please check the reference in the paper for all of them.

## Code Availability
The source code is available on GitHub: https://github.com/njukenanli/ChemBART_Planner. The model weights would be available on Huggingface.

## References

1       Vaswani, A. *et al.* in *31st Annual Conference on Neural Information Processing Systems (NIPS).* (2017).

2       Alec Radford, K. N., Tim Salimans, Ilya Sutskever. *Improving language understanding by generative pre-training*, <https://cdn.openai.com/research-covers/language-unsupervised/language_understanding_paper.pdf> (2018).

3       Devlin, J., Chang, M.-W., Lee, K. & Toutanova, K. in *North American Chapter of the Association for Computational Linguistics.*

4       Lewis, M. *et al.* in *Annual Meeting of the Association for Computational Linguistics.*

5       Weininger, D. SMILES, A CHEMICAL LANGUAGE AND INFORMATION-SYSTEM .1. INTRODUCTION TO METHODOLOGY AND ENCODING RULES. *Journal of Chemical Information and Computer Sciences* **28**, 31-36, doi:10.1021/ci00057a005 (1988).

6       Lin, K., Xu, Y., Pei, J. & Lai, L. Automatic retrosynthetic route planning using template-free models. *Chemical Science* **11**, 3355-3364, doi:10.1039/c9sc03666k (2020).

7       Schwaller, P. *et al.* Molecular Transformer: A Model for Uncertainty-Calibrated Chemical Reaction Prediction. *Acs Central Science* **5**, 1572-1583, doi:10.1021/acscentsci.9b00576 (2019).

8       Seo, S.-W. *et al.* in *35th AAAI Conference on Artificial Intelligence / 33rd Conference on Innovative Applications of Artificial Intelligence / 11th Symposium on Educational Advances in Artificial Intelligence.* 531-539 (2021).

9       Tu, Z. & Coley, C. W. Permutation Invariant Graph-to-Sequence Model for Template-Free Retrosynthesis and Reaction Prediction. *Journal of Chemical Information and Modeling* **62**, 3503-3513, doi:10.1021/acs.jcim.2c00321 (2022).

10      Irwin, R., Dimitriadis, S., He, J. & Bjerrum, E. J. Chemformer: a pre-trained transformer for computational chemistry. *Machine Learning-Science and Technology* **3**, doi:10.1088/2632-2153/ac3ffb (2022).

11      Xue, D. *et al.* X-MOL: large-scale pre-training for molecular understanding and diverse molecular analysis. *bioRxiv* (2020).

12      Wang, S., Guo, Y., Wang, Y., Sun, H. & Huang, J. in *Proceedings of the 10th ACM International Conference on Bioinformatics, Computational Biology and Health Informatics*    429–436 (Association for Computing Machinery, Niagara Falls, NY, USA, 2019).

13      Fabian, B. *et al.* Molecular representation learning with language models and domain-relevant auxiliary tasks.    (2020).

14      Lu, J. & Zhang, Y. Unified Deep Learning Model for Multitask Reaction Predictions with Explanation. *Journal of Chemical Information and Modeling* **62**, 1376-1387, doi:10.1021/acs.jcim.1c01467 (2022).

15      Kallergis, G. *et al.* Domain adaptable language modeling of chemical compounds identifies potent pathoblockers for Pseudomonas aeruginosa. *Communications Chemistry* **8**, 114 (2025).

16      Jang, Y., Kim, J. & Ahn, S. Improving Chemical Understanding of LLMs via SMILES Parsing. *arXiv preprint arXiv:2505.16340* (2025).

17      Xiong, J. *et al.* Bridging chemistry and artificial intelligence by a reaction description language. *Nature Machine Intelligence* **7**, 782-793, doi:10.1038/s42256-025-01032-8 (2025).

18      Jiang, Y. *et al.* in *Proceedings of the AAAI Conference on Artificial Intelligence* Vol. 37(4) 5113-5121 (2023).

19      Deng, Y. *et al.* RSGPT: a generative transformer model for retrosynthesis planning pre-trained on ten billion datapoints. *Nature communications* **16**, 7012-7012, doi:10.1038/s41467-025-62308-6 (2025).

20      Kreutter, D. & Reymond, J.-L. Multistep retrosynthesis combining a disconnection aware triple transformer loop with a route penalty score guided tree search. *Chemical Science* **14**, 9959-9969, doi:10.1039/d3sc01604h (2023).

21      Jiasheng Guo, C. Y., Kenan Li, Yijian Zhang, Guoqiang Wang, Shuhua Li, Hao Dong. Retrosynthesis Zero: Self-Improving Global Synthesis Planning Using Reinforcement Learning. *Journal of Chemical Theory and Computation* (2024).

22      Coulom, R. in *5th International Conference on Computers and Games.*    72-83 (2007).

23      Perera, D. *et al.* A platform for automated nanomole-scale reaction screening and micromole-scale synthesis in flow. *Science* **359**, 429-434, doi:10.1126/science.aap9112 (2018).

24      *USPTO*, <https://www.uspto.gov/> (

25      *USPTO-MIT*, <https://github.com/wengong-jin/nips17-rexgen/blob/master/USPTO/> (

26      Deng, Y. *et al.* RSGPT: a generative transformer model for retrosynthesis planning pre-trained on ten billion datapoints. *Nature Communications* **16**, 7012 (2025).

27      Segler, M. H. & Waller, M. P. Neural‑symbolic machine learning for retrosynthesis and reaction prediction. *Chemistry–A European Journal* **23**, 5966-5971 (2017).

28      Xiong, J. *et al.* Bridging chemistry and artificial intelligence by a reaction description language. *Nature Machine Intelligence*, 1-12 (2025).

29      Karpov, P., Godin, G. & Tetko, I. V. in *International conference on artificial neural networks.* 817-830 (Springer).

30      <https://docs.open-reaction-database.org/en/latest/> (

31      Schwaller, P., Vaucher, A. C., Laino, T. & Reymond, J.-L. Prediction of chemical reaction yields using deep learning. *Machine learning: science and technology* **2**, 015016 (2021).

32      Wu, Z. *et al.* MoleculeNet: a benchmark for molecular machine learning††Electronic supplementary information (ESI) available. See DOI: 10.1039/c7sc02664a. *Chemical Science* **9**, 513-530, doi:https://doi.org/10.1039/c7sc02664a (2018).

33      Hu, E. J. *et al.* Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685* (2021).

34      Honda, S., Shi, S. & Ueda, H. R. SMILES Transformer: Pre-trained Molecular Fingerprint for Low Data Drug Discovery. *CoRR* **abs/1911.04738** (2019).

35      Rogers, D. & Hahn, M. Extended-connectivity fingerprints. *Journal of chemical information and modeling* **50**, 742-754 (2010).

36      Duvenaud, D. *et al.* in *Proceedings of the 28th International Conference on Neural Information Processing Systems - Volume 2.*    2224–2232-2224–2232 (MIT Press).

37      Kearnes, S., McCloskey, K., Berndl, M., Pande, V. & Riley, P. Molecular graph convolutions: moving beyond fingerprints. *Journal of Computer-Aided Molecular Design* **30**, 595-608, doi:10.1007/s10822-016-9938-8 (2016).

38      Chithrananda, S., Grand, G. & Ramsundar, B. ChemBERTa: Large-Scale Self-Supervised Pretraining for Molecular Property Prediction. *CoRR* **abs/2010.09885** (2020).

39      Yang, K. *et al.* Analyzing Learned Molecular Representations for Property Prediction. *Journal of chemical information and modeling* **59**, 3370-3388 (2019).

40      Winter, R., Montanari, F., Noé, F. & Clevert, D.-A. Learning continuous and data-driven molecular descriptors by translating equivalent chemical representations. *Chemical science* **10**, 1692-1701 (2019).

41      You, Y. *et al.* in *Proceedings of the 34th International Conference on Neural Information Processing Systems.*    (Curran Associates Inc.).

42      Xu, M., Wang, H., Ni, B., Guo, H. & Tang, J. Self-supervised Graph-level Representation Learning with Local and Global Structure. *CoRR* **abs/2106.04113** (2021).

43      Jaeger, S., Fulle, S. & Turk, S. Mol2vec: Unsupervised Machine Learning Approach with Chemical Intuition. *Journal of Chemical Information and Modeling* **58**, 27-35, doi:10.1021/acs.jcim.7b00616 (2018).

44      Wang, H. *et al.* Chemical-Reaction-Aware Molecule Representation Learning. *arXiv e-prints*, arXiv:2109.09888-arXiv:02109.09888 (2021).

45      Chen, B., Li, C., Dai, H. & Song, L. Retro*: Learning Retrosynthetic Planning with Neural Guided A* Search. *ArXiv* **abs/2006.15820** (2020).

46      Liu, G. *et al.* in *International Conference on Machine Learning.*

47      Hong, S., Zhuo, H. H., Jin, K., Shao, G. & Zhou, Z. Retrosynthetic planning with experience-guided Monte Carlo tree search. *Communications Chemistry* **6**, doi:10.1038/s42004-023-00911-8 (2023).

48      Segler, M. H. S. & Waller, M. P. Neural-Symbolic Machine Learning for Retrosynthesis and Reaction Prediction. *Chemistry-a European Journal* **23**, 5966-5971, doi:10.1002/chem.201605499 (2017).

49       Wu, D. D. G. L. G. J. L. Y.-H. RENIN INHIBITORS. WO (World Patent via PCT) patent (2008).

50       Sun, C. *et al.* Novel PD-L1/VISTA Dual Inhibitor as Potential Immunotherapy Agents. *Journal of Medicinal Chemistry* **68**, 156-173, doi:10.1021/acs.jmedchem.4c01640 (2025).

51       quantum-machine.org. *QM9 Dataset*, <https://paperswithcode.com/dataset/qm9> (2025).

52       Xie, L., Jin, Y., Xu, L., Chang, S. & Xu, X. Fusing Domain Knowledge with a Fine-Tuned Large Language Model for Enhanced Molecular Property Prediction. *Journal of Chemical Theory and Computation* (2025).

53       Zhou, G. *et al.* in *International Conference on Learning Representations.*

54       Dubey, A. *et al.* The llama 3 herd of models. *arXiv e-prints*, arXiv: 2407.21783 (2024).

55       Yang, A. *et al.* Qwen3 technical report. *arXiv preprint arXiv:2505.09388* (2025).

56       Zhang, Y. *et al.* Large language models to accelerate organic chemistry synthesis. *Nature Machine Intelligence*, doi:10.1038/s42256-025-01066-y (2025).

57       Zhang, Y. *et al.* A Large Language Model for Chemistry and Retrosynthesis Predictions. *arXiv preprint arXiv:2507.01444* (2025).

58       Zhang, C. *et al.* SynAsk: unleashing the power of large language models in organic synthesis. *Chemical science* **16**, 43-56 (2025).

59       Zhao, Z. *et al.* Developing ChemDFM as a large language foundation model for chemistry. *Cell Reports Physical Science* **6** (2025).

60       Lin, X., Liu, Q., Xiang, H., Zeng, D. & Zeng, X. Enhancing Chemical Reaction and Retrosynthesis Prediction with Large Language Model and Dual-task Learning. *arXiv preprint arXiv:2505.02639* (2025).

61       Wang, H. *et al.* LLM-Augmented Chemical Synthesis and Design Decision Programs. *arXiv preprint arXiv:2505.07027* (2025).

62       Zheng, Y. *et al.* Large language models for scientific discovery in molecular property prediction. *Nature Machine Intelligence*, 1-11 (2025).

63       Daylight. *A Language for Describing Molecular Patterns*, <https://www.daylight.com/dayhtml/doc/theory/theory.smarts.html> (

64       Lipton, Z. C. A Critical Review of Recurrent Neural Networks for Sequence Learning. *ArXiv* **abs/1506.00019** (2015).

65       Hochreiter, S. & Schmidhuber, J. Long short-term memory. *Neural Computation* **9**, 1735-1780, doi:10.1162/neco.1997.9.8.1735 (1997).

66       Schneider, N., Stiefl, N. & Landrum, G. A. What's What: The (Nearly) Definitive Guide to Reaction Role Assignment. *Journal of Chemical Information and Modeling* **56**, 2336-2346, doi:10.1021/acs.jcim.6b00564 (2016).

67       Zhong, Z. *et al.* Root-aligned SMILES: a tight representation for chemical reaction prediction. *Chemical Science* **13**, 9023-9034, doi:10.1039/d2sc02763a (2022).

68       Coley, C. W., Rogers, L., Green, W. H. & Jensen, K. F. Computer-Assisted Retrosynthesis Based on Molecular Similarity. *Acs Central Science* **3**, 1237-1245, doi:10.1021/acscentsci.7b00355 (2017).

69       Liu, B. *et al.* Retrosynthetic Reaction Prediction Using Neural Sequence-to-Sequence Models. *Acs Central Science* **3**, 1103-1113, doi:10.1021/acscentsci.7b00303 (2017).

70       Loshchilov, I. & Hutter, F. Fixing Weight Decay Regularization in Adam. *ArXiv* **abs/1711.05101** (2017).

71       Mangrulkar, S. *et al.* Peft: State-of-the-art parameter-efficient fine-tuning methods. *URL:*

*https://github. com/huggingface/peft* (2022).

72      Hu, Y., Lounkine, E. & Bajorath, J. Filtering and Counting of Extended Connectivity Fingerprint Features Maximizes Compound Recall and the Structural Diversity of Hits. *Chemical Biology & Drug Design* **74**, 92-98, doi:10.1111/j.1747-0285.2009.00830.x (2009).

73      *RDKit*, <https://www.rdkit.org/> (2024).

74      Kim, J., Ahn, S., Lee, H. & Shin, J. Self-Improved Retrosynthetic Planning. *ArXiv* **abs/2106.04880** (2021).

# Supplementary Information

## ChemBART: A Pre-trained BART Model Assisting Organic Chemistry Analysis

Kenan Li [a,1], Yijian Zhang [a,1], Jin Wang [b,1], Haipeng Gan [a], Zeying Sun [b], Xiaoguang Lei [b,*], & Hao Dong [a,*]
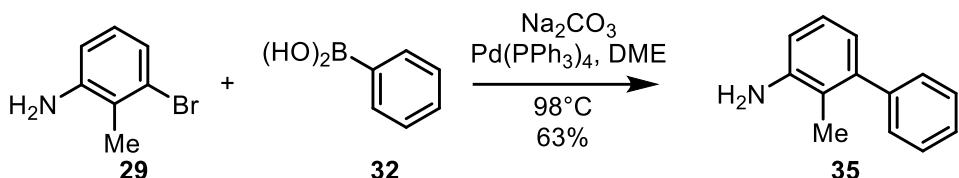
[a] State Key Laboratory of Analytical Chemistry for Life Science, Kuang Yaming Honors School, Chemistry and Biomedicine Innovation Centre (ChemBIC), ChemBioMed Interdisciplinary Research Centre at Nanjing University, Engineering Research Centre of Protein and Peptide Medicine of Ministry of Education, Institute for Brain Sciences, Nanjing University, Nanjing 210023, China.
[b] Beijing National Laboratory for Molecular Sciences, Key Laboratory of Bioorganic Chemistry and Molecular Engineering of Ministry of Education, College of Chemistry and Molecular Engineering, and Peking-Tsinghua Center for Life Sciences, Peking University, Beijing 100871, China; Academy for Advanced Interdisciplinary Studies, Peking University, Beijing 100871, China; Institute for Cancer Research, Shenzhen Bay Laboratory, Shenzhen 518107, China

[1] these authors contributed equally
* corresponding authors: xglei@pku.edu.cn (X.G.L), donghao@nju.edu.cn (H.D.)

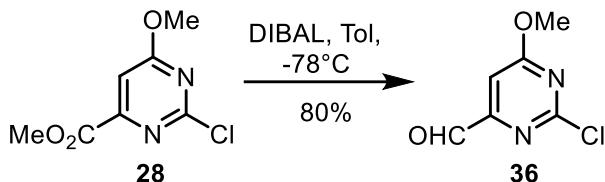## Experimental procedures for the synthesis of P1



To a stirred solution of amine **29** (186 mg, 1 mmol, 1.0 equiv.), phenylboronic acid **32** (145 mg, 1.2 mmol, 1.2 equiv.), and $Na_2CO_3$ (296 mg, 2.8 mmol, 2.8 equiv.) in DME and water (1.8 mL, volume ratio = 5:1) was added $Pd(PPh_3)_4$ (57 mg, 0.05 mmol, 0.05 equiv.). The reaction mixture was then heated to 100 °C and maintained for 16 h under an argon atmosphere. After completion, the residue was extracted three times with ethyl acetate. The combined organic phases were washed with brine (10 mL), dried over $Na_2SO_4$, filtered and concentrated in vacuo to afford a crude oil, which was purified by by flash column chromatography on silica gel (Petroleum ether/Ethyl acetate =2:1) to afford **35** (110 mg, 60%) as a white solid.

[1]H NMR (400 MHz, CDCl$_3$) δ 7.47 – 7.39 (m, 2H), 7.38 – 7.31 (m, 3H), 7.10 (t, $J$ = 7.7 Hz, 1H), 6.73 (dd, $J$ = 7.7, 2.0 Hz, 2H), 3.66 (brs, 2H), 2.09 (s, 3H).
[13]C NMR (101 MHz, CDCl$_3$) δ 145.1, 143.2, 142.5, 129.5, 128.1, 126.7, 126.2, 120.7, 120.0, 114.2, 14.65.
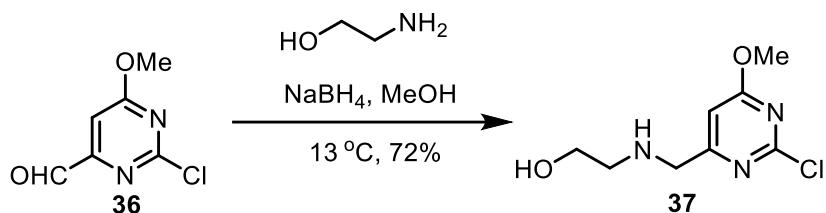HRMS(ESI) [M + H]+ calculated for C13H14N: 184.1126, found: 184.1124;



To a solution of methyl 2-chloro-6-methoxypyrimidine-4-carboxylate **28** (56 mg, 0.27 mmol, 1.0 equiv.) in toluene (10 mL) at -78°C under argon atmosphere was added DIBAL-H (1.0 M solution in hexane, 0.8 mL, 0.8 mmol, 3.0 equiv.) dropwise. The reaction was stirred at -78 °C for 3 h, then quenched with MeOH and saturated potassium sodium tartrate. The reaction was warmed to room temperature gradually and stirred overnight. The aqueous layer was extracted with ethyl acetate for three times, then the combined organic layers was dried over $Na_2SO_4$, and concentrated under reduced pressure to afford a crude oil, which was purified by by flash column chromatography on silica gel (Petroleum ether/CH$_2$Cl$_2$ =2:1) to afford **36** (37 mg, 80%) as a white solid.

$^1$H NMR (400 MHz, CDCl$_3$) δ 9.92 (s, 1H), 7.17 (s, 1H), 4.08 (s, 3H).
$^{13}$C NMR (101 MHz, CDCl$_3$) δ 190.6, 172.1, 161.3, 160.4, 104.1, 55.5.
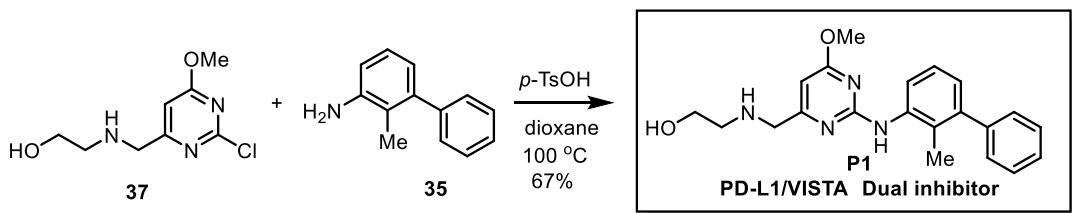HRMS(ESI) [M + H]+ calculated for C6H6N2O2Cl: 173.0118, found: 173.0119;



To a solution of the aldehyde **36** (17 mg, 0.1 mmol, 1.0 equiv.) and ethanolamine (10 μL, 0.65 mmol, 1.5 equiv.) in MeOH (5 mL) at 0°C under argon atmosphere was added ethanolamine (10 μL, 0.16 mmol, 1.6 equiv.). The solution was stirred at 0°C for 60 min. Then NaBH$_4$ (4.5 mg, 0.12 mmol, 1.2 equiv.) was added at 0°C. The reaction mixture was gradually warmed to room temperature and stirred overnight. The solution was concentrated under reduced pressure and the residue was purified by flash column chromatography on silica gel (CH2Cl2/MeOH = 60 /1- 20/1)to afford the title compound **37** (15 mg, 67%) as a white solid.

$^1$H NMR (400 MHz, DMSO-$d_6$) δ 7.07 (s, 1H), 4.08 (s, 2H), 3.96 (s, 3H), 3.60 (t, $J$ = 5.5 Hz, 2H), 2.88 (t, $J$ = 5.5 Hz, 2H).
$^{13}$C NMR (101 MHz, DMSO-$d_6$) δ 171.0, 167.6, 158.6, 105.2, 57.8, 54.9, 50.4, 49.8, 39.5.
HRMS(ESI) [M + H]+ calculated for C8H13N3O2Cl: 218.0696, found: 218.0700;



Alcohol **37** (8 mg, 0.036 mmol, 1.0 equiv.), 2-methyl-[1,1'-biphenyl]-3-amine **35** (10 mg, 0.054 mmol, 1.5 equiv.) and *p*-toluenesulfonic acid (8 mg, 0.046 mmol, 1.3 equiv.) was dissolved in 1,4-dioxane (0.5 mL) under argon atmosphere. The reaction was heated to 100°C and stirred overnight. After that, the solution was cooled to room temperature and concentrated under reduced pressure. The residue was purified by flash column chromatography on silica gel (CH2Cl2/MeOH = 60 /1- 20/1) to afford the title compound **P1** (9 mg, 67%) as a white solid.

[1]H NMR (400 MHz, DMSO-$d_6$) δ 8.78 (s, 1H), 7.53 (d, $J$ = 7.4 Hz, 1H), 7.45 (t, $J$ = 7.3 Hz, 2H), 7.39 – 7.34 (m, 1H), 7.33 – 7.28 (m, 2H), 7.22 (t, $J$ = 7.8 Hz, 1H), 7.00 (d, $J$ = 7.5 Hz, 1H), 6.28 (s, 1H), 3.82 (s, 3H), 3.71 (s, 2H), 3.52 (t, $J$ = 5.5 Hz, 2H), 2.73 (t, $J$ = 5.6 Hz, 2H), 2.09 (s, 3H).

[13]C NMR (101 MHz, DMSO) δ 170.3, 160.5, 142.3, 141.6, 138.3, 129.6, 129.1, 128.2, 126.9, 125.6, 125.3, 124.3, 94.9, 59.3, 53.1, 52.5, 50.7, 15.9.

HRMS(ESI) [M + H]+ calculated for C21H25N4O2: 365.1978, found: 365.1968;

**Supplementary Figure S1 | The top-*k* success rate for reactant prediction on USPTO-MIT dataset for ChemBART training.** After the $7^{th}$ epoch, model convergence observed.

**Supplementary Figure S2 | The top-*k* syntax error rate for reactant prediction on USPTO-MIT dataset for ChemBART training.**

**Supplementary Figure S3 | 3D visualization using PCA and ChemBART embeddings for the top 20 most frequent elements in the USPTO dataset.** The resulting visualization suggests the model has learned certain periodic trends: clustering is observed for organic backbone elements from the same period C, N, O; elements from the same group H, Li, Na, K, Cs and F, Cl, Br, I are nearer in the PCA space; additionally, elements with lower electronegativity B, Al, Si, Cu, Pd, Sn also show a tendency to cluster.

**Supplementary Figure S4 | The cross-attention matrix generated by the ChemBART model for a Grignard reaction, $CH_3(CH_2)_5Br + PhMgBr \xrightarrow{C_2H_5OC_2H_5} CH_3(CH_2)_5Ph + MgBr_2$, by masking the product.** (A) The entire reaction; (B) Bromine atom and related atoms; (C) Carbon anion and related atoms.

**Supplementary Figure S5 | The decoder attention matrix generated by the ChemBART model for a Grignard reaction, $CH_3(CH_2)_5Br + PhMgBr \xrightarrow{C_2H_5OC_2H_5} CH_3(CH_2)_5Ph + MgBr_2$, by masking the product.** (A) The entire reaction; (B) Bromine atom and related atoms; (C) Carbon anion and related atoms.

**Supplementary Figure S6 | Attention matrices generated by the ChemBART model for a reaction, $PhO(CH_2)_3CH_3 + Cl_2 \xrightarrow{FeCl_3} o\text{-}ClPhO(CH_2)_3CH_3 + p\text{-}ClPhO(CH_2)_3CH_3$, with products masked.** (a-c) Encoder, cross, and decoder attention matrices, respectively. The top panel is about entire reaction, and the bottom panel is about n-Butyl phenyl ether and $Cl_2$.

**Supplementary Figure S7 | Encoder attention matrix generated by masking the reagents for products to reactants (a), and for reactants to products (b).**

**Supplementary Figure S8 | Fine-tuning the pre-trained model on specific datasets enhances the distinction in molecular feature vector distributions.** (a-b) The BBBP dataset, before (a) and after (b) fine-tuning, where dark and light shades represent non-penetrating and penetrating molecules, respectively. (c-d) The BACE dataset, before (c) and after (d) fine-tuning, where dark and light shades correspond to non-binders and binders, respectively.

**Supplementary Figure S9 | Comparison between ReSynZ and ChemBART in the training of value function.** Pre-trained ChemBART-F/M fits value function data well with a lower minimum validation RMSE compared to ReSynZ, while ChemBART-R performs significantly worse.

**Supplementary Figure S10 | Comparison between ReSynZ and ChemBART in the training of policy function.** Pre-trained ChemBART-F/M show good fitting performance on policy function, while ChemBART-R and ReSynZ exhibit early over-fitting.

**Patent 1:**



**ChemBart 1:**



**Patent 2:**



**ChemBart 2:**



**Supplementary Figure S11 | Comparison of ChemBART and Patent Synthetic Routes for two compounds.** ChemBART adopted the same bond-breaking scheme as the patent but shortened the synthetic pathway by selecting appropriate reagents.

**Supplementary Figure S12 | Probability distributions for 53 molecules from the JMC2025 dataset across six properties.** (a) LogP, the logarithm of the partition coefficient. (b) MW, molecular weight. (c) QED, quantitative estimate of drug-likeness, quantifies drug-likeness by taking into account the main molecular properties, ranging from 0 (unfavorable) to 1 (favorable). (d) TPSA, topological polar surface area, the sum of surface area over all polar atoms. (e) SAS, synthetic accessibility score, the measurement of the difficulty of synthesizing a compound. It is a score between 1 (easy to make) and 10 (very difficult to make). (f) Heavy atom count.

**Supplementary Figure S13 | t-SNE plot of 53 molecular Morgan fingerprints from the JMC20225 dataset.**

**Supplementary Figure S14 | Histograms, KDE curves, and boxplots of the route step distributions of different strategies.** (a) Beam search, (b) top-*k* sampling, and (c) top-*p* sampling.

# NMR spectra



**Supplementary Figure S15 | $^1$H NMR (400 MHz, CDCl$_3$) of 35**

**Supplementary Figure S16 |** $^{13}$C NMR (100 MHz, CDCl$_3$) of **35**

**Supplementary Figure S17 | $^1$H NMR (400 MHz, CDCl$_3$) of 36**

**Supplementary Figure S18 | $^{13}$C NMR (100 MHz, CDCl$_3$) of 36**

**Supplementary Figure S19 | $^1$H NMR (400 MHz, DMSO-$d_6$) of 37**

**Supplementary Figure S20 |** $^{13}$**C NMR (101 MHz, DMSO-***d*$_6$**) of 37**

**Supplementary Figure S21 | $^1$H NMR (400 MHz, DMSO-$d_6$) of P1**

**Supplementary Figure S22 |** $^{13}$C NMR (100 MHz, DMSO-$d_6$) of P1

**Supplementary Table S1 | Classification of all possible tokens in the SMILES language**

| Category | Tokens |
|---|---|
| location label | <pad>, <cls>, . (sep), >(generate), <end>, <msk> |
| punctuation mark | [, ], @, +, -, (, ), /, \. =, #, :, % |
| non-aromatic atoms | H, He, Li, ... , Bi, Po, U, Th |
| aromatic atoms | c, o, n, s, se, p, te |
| numbers (rings, atom index, charge, atom mapping, isotope atomic weight) | 0, 1, ... , 83 (numbers > 83 would be broken up) |
| unused tokens to be set by user | <n00> ... <n04> |
| unknown tokens / wildcard | ? |

**Supplementary Table S2 | The vocabulary defined by this model, used to specify a one-to-one mapping from tokens to their indices.**

| | | | |
|---|---|---|---|
| { | "32": 50, | "83": 101, | "Zr": 152, |
| "<pad>": 0, | "33": 51, | "?": 102, | "Nb": 153, |
| "<cls>": 1, | "34": 52, | "Th": 103, | "Mo": 154, |
| ">": 2, | "35": 53, | ":": 104, | "Tc": 155, |
| "<end>": 3, | "36": 54, | "te": 105, | "Ru": 156, |
| "<msk>": 4, | "37": 55, | "p": 106, | "Rh": 157, |
| "%": 5, | "38": 56, | "<n00>": 107, | "Pd": 158, |
| "[": 6, | "39": 57, | "c": 108, | "Ag": 159, |
| "]": 7, | "40": 58, | "o": 109, | "Cd": 160, |
| "@": 8, | "41": 59, | "n": 110, | "In": 161, |
| "+": 9, | "42": 60, | "s": 111, | "Sn": 162, |
| "-": 10, | "43": 61, | "se": 112, | "Sb": 163, |
| ".": 11, | "44": 62, | "H": 113, | "Te": 164, |
| "(": 12, | "45": 63, | "He": 114, | "I": 165, |
| ")": 13, | "46": 64, | "Li": 115, | "Xe": 166, |
| "/": 14, | "47": 65, | "Be": 116, | "Cs": 167, |
| "\\": 15, | "48": 66, | "B": 117, | "Ba": 168, |
| "=": 16, | "49": 67, | "C": 118, | "La": 169, |
| "#": 17, | "50": 68, | "N": 119, | "Ce": 170, |
| "0": 18, | "51": 69, | "O": 120, | "Pr": 171, |
| "1": 19, | "52": 70, | "F": 121, | "Nd": 172, |
| "2": 20, | "53": 71, | "Ne": 122, | "Pm": 173, |
| "3": 21, | "54": 72, | "Na": 123, | "Sm": 174, |
| "4": 22, | "55": 73, | "Mg": 124, | "Eu": 175, |
| "5": 23, | "56": 74, | "Al": 125, | "Gd": 176, |
| "6": 24, | "57": 75, | "Si": 126, | "Tb": 177, |
| "7": 25, | "58": 76, | "P": 127, | "Dy": 178, |
| "8": 26, | "59": 77, | "S": 128, | "Ho": 179, |
| "9": 27, | "60": 78, | "Cl": 129, | "Er": 180, |
| "10": 28, | "61": 79, | "Ar": 130, | "Tm": 181, |
| "11": 29, | "62": 80, | "K": 131, | "Yb": 182, |
| "12": 30, | "63": 81, | "Ca": 132, | "Lu": 183, |
| "13": 31, | "64": 82, | "Sc": 133, | "Hf": 184, |
| "14": 32, | "65": 83, | "Ti": 134, | "Ta": 185, |
| "15": 33, | "66": 84, | "V": 135, | "W": 186, |
| "16": 34, | "67": 85, | "Cr": 136, | "Re": 187, |
| "17": 35, | "68": 86, | "Mn": 137, | "Os": 188, |
| "18": 36, | "69": 87, | "Fe": 138, | "Ir": 189, |
| "19": 37, | "70": 88, | "Co": 139, | "Pt": 190, |
| "20": 38, | "71": 89, | "Ni": 140, | "Au": 191, |
| "21": 39, | "72": 90, | "Cu": 141, | "Hg": 192, |
| "22": 40, | "73": 91, | "Zn": 142, | "Tl": 193, |
| "23": 41, | "74": 92, | "Ga": 143, | "Pb": 194, |
| "24": 42, | "75": 93, | "Ge": 144, | "Bi": 195, |
| "25": 43, | "76": 94, | "As": 145, | "Po": 196, |
| "26": 44, | "77": 95, | "Se": 146, | "U": 197, |
| "27": 45, | "78": 96, | "Br": 147, | "<n01>": 198, |
| "28": 46, | "79": 97, | "Kr": 148, | "<n02>": 199, |
| "29": 47, | "80": 98, | "Rb": 149, | "<n03>": 200, |
| "30": 48, | "81": 99, | "Sr": 150, | "<n04>": 201 |
| "31": 49, | "82": 100, | "Y": 151, | } |

**Supplementary Table S3 | SMILES notation and corresponding digital sequence for a halogenation reaction, conceptualized as a "reaction sentence" in this framework.** Each sentence is prefixed with a start token <cls> and suffixed with an

end token <end> for compatibility within the language model.

| Reaction |  |
|---|---|
| SMILES expression | <cls>c1ccccc1C(C)(C)(C).ClCl>[Fe+3].[Cl-].[Cl-].[Cl-]>c1cc(Cl)ccc1C(C)(C)(C)<end> |
| Digital sequence | **1**, 108, 19, 108, 108, 108, 108, 108, 19, 118, 12, 118, 13, 12, 118, 13, 12, 118, 13, **11**, 129, 129, **2**, 6, 138, 9, 21, 7, 11, 6, 129, 10, 7, 11, 6, 129, 10, 7, 11, 6, 129, 10, 7, **2**, 108, 19, 108, 108, 12, 129, 13, 108, 108, 108, 19, 118, 12, 118, 13, 12, 118, 13, 12, 118, 13, **3** |

**Supplementary Table S4 | RMSE of Temperature-Yield Regression on the ORD Dataset Across Different Training Schemes.** The pre-trained model outperforms the randomly initialized model, demonstrating the effectiveness of pre-training. The units of temperature and yield are Celsius and percentage, respectively.

| Scheme | Temperature | | | Yield | | |
|---|---|---|---|---|---|---|
| | Training | Validation | Test | Training | Validation | Test |
| ChemBART-F (c) (e) | 39.50 | 36.91 | 37.11 | - | - | - |
| ChemBART-F (d) (e) | 39.77 | 36.68 | 36.99 | - | - | - |
| ChemBert-F (c) (e) | 41.33 | 38.33 | 38.80 | - | - | - |
| ChemBART-F (a) (c) (f) | 42.11 | 39.45 | 40.42 | 25.03 | 25.71 | 26.08 |
| ChemBART-F (b) (c) (f) | 41.25 | 39.43 | 39.38 | 23.67 | 25.18 | 25.13 |
| ChemBert-F (a) (c) (f) | 44.43 | 42.78 | 42.97 | 25.34 | 25.99 | 26.23 |
| ChemBART-M (b) (c) (f) | 44.61 | 39.82 | 38.80 | 24.02 | 24.85 | 26.02 |
| ChemBART-R (b) (c) (f) | 50.06 | 44.75 | 45.07 | 24.89 | 25.98 | 26.51 |

(a) The output head of the yield uses a sigmoid layer to control the output range;
(b) the output head of the yield does not use a sigmoid layer to control the output range;
(c) Use only one task token per attribute;
(d) use a Bi-LSTM layer that integrates the outputs from all the input tokens before inputting the output vector into the linear head;
(e) The regression for the single task of temperature;
(f) The joint regression for temperature and yield.

**Supplementary Table S5 | Accuracy of Different Models on Temperature-Yield Classification for the ORD Dataset.** For the temperature classification, the ranges for each class are defined as (-150, 0), (0, 15), (15, 30), (30, 80), (80, 150), and (150, 400) Celsius. For the yield classification, the ranges are (0, 25), (25, 50), (50, 75), and (75, 110) percentage. The units of accuracy is percentage.

| | | Temperature | | | Yield | | |
|---|---|---|---|---|---|---|---|
| | | **Training** | **Validation** | **Test** | **Training** | **Validation** | **Test** |
| Top1 | ChemBART-F | 55.1 | 55.42 | 64.17 | 45.38 | 43.36 | 41.52 |
| | ChemBART-M | 57.1 | 55.08 | 64.69 | 45.99 | 43.75 | 41.55 |
| | ChemBART-R | 46.62 | 46.25 | 57.49 | 43.40 | 42.49 | 40.4 |
| Top3 | ChemBART-F | 55.1 | 55.42 | 86.70 | 45.38 | 43.36 | 78.22 |
| | ChemBART-M | 57.19 | 55.08 | 86.61 | 45.99 | 43.75 | 77.59 |
| | ChemBART-R | 46.62 | 46.25 | 81.64 | 43.40 | 42.49 | 70.88 |

**Supplementary Table S6 | 10-fold validation of ChemBART fine-tuned on the Suzuki-Miyaura dataset [1].** Data points with yields exceeding 1% (n=5469) were retained from the total dataset (n=5760) to account for experimental errors.

| Split | $R^2$ | RMSE (%) | MAE (%) |
|---|---|---|---|
| Split 0 | 0.802 | 12.0 | 8.4 |
| Split 1 | 0.816 | 11.7 | 7.9 |
| Split 2 | 0.822 | 11.5 | 7.7 |
| Split 3 | 0.766 | 13.1 | 9.0 |
| Split 4 | 0.817 | 11.6 | 8.0 |
| Split 5 | 0.824 | 11.4 | 7.9 |
| Split 6 | 0.795 | 12.4 | 8.7 |
| Split 7 | 0.808 | 11.8 | 8.1 |
| Split 8 | 0.795 | 12.5 | 8.3 |
| Split 9 | 0.793 | 12.4 | 8.7 |
| average | 0.804 | 12.0 | 8.3 |
| STDEV | 0.0166 | 0.516 | 0.402 |

**Supplementary Table S7 | Comparison of the Performance (Mean RMSE) of ChemBART and Traditional Fully-connected-layer Networks on the Task of Fitting MCTS Data.**

| | training set | | validation set | | test set | |
|---|---|---|---|---|---|---|
| | $v$ | $p$ | $v$ | $p$ | $v$ | $p$ |
| ChemBART-F | 0.011 | 0.092 | 0.091 | 0.154 | 0.126 | 0.160 |
| ChemBART-M | 0.006 | 0.077 | 0.099 | 0.156 | 0.106 | 0.158 |
| ChemBART-R | 0.054 | 0.170 | 0.293 | 0.181 | 0.323 | 0.180 |
| ReSynZ$^2$ | 0.015 | 0.204 | 0.085 | 0.202 | 0.111 | 0.200 |

**Supplementary Table S8 | The Results for the Retro\*-190 Dataset.**

| ID | Mol | Retro\*-190 step | ChemBART-F | | ChemBART-M | |
|---|---|---|---|---|---|---|
| | | | step | success | step | success |
| 0 | C[C@H](c1ccccc1)N1C[C@]2(C(=O)OC(C)(C)C)C=CC[C@@H]2C1=S | 7 | 1 | 0 | 0 | 0 |
| 1 | CCCC[C@@H](C(=O)N1CCC[C@H]1C(=O)O)[C@@H](F)C(=O)OC | 8 | 0 | 0 | 4 | 1 |
| 2 | CC(C)c1ccc(-n2nc(O)c3c(=O)c4ccc(Cl)cc4[nH]c3c2=O)cc1 | 2 | 0 | 0 | 6 | 0 |
| 3 | CCOC(=O)c1nc(N2CC[C@H](NC(=O)c3nc(C(F)(F)F)c(CC)[nH]3)[C@H](OC)C2)sc1C | 9 | 5 | 1 | 6 | 1 |
| 4 | NC1=N[C@@]2(c3ccc(F)cc3F)CO[C@@H](c3nnco3)C[C@H]2CS1 | 14 | 14 | 0 | 16 | 0 |
| 5 | COc1ccc2c(Nc3c(Cl)cncc3Cl)cc(=O)[nH]c2c1OCCCCCCN1CCCNCC1 | 10 | 7 | 1 | 7 | 1 |
| 6 | COCCCc1cc(CN(C(=O)[C@H]2CN(C(=O)OC(C)(C)C)CC[C@@H]2c2ccc(OCCOc3c(Cl)cc(C)cc3Cl)cc2)C2CC2)cc(OCc2ccc(C#N)cc2)c1 | 14 | 11 | 1 | 13 | 1 |
| 7 | O=C(OCc1ccccc1)N1CC[C@H]2CCCN(CCc3ccccc3)C[C@H]21 | 3 | 5 | 1 | 11 | 1 |
| 8 | N[C@H]1CC[C@H]1c1ccc(Cl)cc1 | 4 | 1 | 1 | 1 | 1 |
| 9 | CCc1nc(C)c(C(=O)Nc2ccc(-c3ccc(C45CCC(CC(=O)OC)(CC4)OC5)cc3)cc2)o1 | 7 | 9 | 1 | 3 | 1 |
| 10 | O=C(Cc1cccs1)NC1C(=O)N2C(C(=O)O)C(CBr)=CS[C@H]12 | 3 | 16 | 0 | 16 | 0 |
| 11 | O=[N+]([O-])c1ccc(OC2(S(=O)(=O)c3ccccc3)CC2)nc1Cl | 3 | 1 | 1 | 1 | 1 |
| 12 | CCCN(CCC)Cc1cc(CO)c(=O)n2c(-c3c(C)cc(C)cc3C)cccc12 | 4 | 14 | 0 | 5 | 1 |
| 13 | C/N=C(\C)C1=C(O)C=C2Oc3c(C(=O)NCc4c(C)ccc5ccccc45)c(OC)cc(O)c3[C@]2(C)C1=O | 3 | 1 | 0 | 3 | 0 |
| 14 | CCOC(=O)c1cc2c(F)cccc2nc1[C@H](C)NC(=O)OC(C)(C)C | 2 | 2 | 0 | 9 | 1 |
| 15 | COc1ccc2nc(-c3ccc(OC)c(F)c3)nc(C(=O)NC(CO)Cc3c[nH]c4cccnc34)c2c1 | 7 | 6 | 1 | 11 | 1 |
| 16 | COC(=O)CC12CCC(c3ccc(Br)cc3)(CC1)CO2 | 5 | 0 | 1 | 0 | 1 |
| 17 | C#CC1(O)C(C)=CC2(CC1(C)C(F)(F)F)OC(C)C(C)O2 | 5 | 0 | 0 | 0 | 0 |
| 18 | CCC/C=C/[C@H]1CC[C@H](C(C(=O)O)C(=O)O)CC1 | 5 | 3 | 1 | 1 | 0 |
| 19 | CC(C)(C)OC(=O)NC1(C2CCc3cc(Sc4cccc(OCc5ccccc5)c4)ccc3C2)COC(C)(C)OC1 | 7 | 16 | 0 | 16 | 1 |
| 20 | C[Si](C)(C)CCOCn1cc(C2CCc3c(C(=O)O)nn(COCC[Si](C)(C)C)c3C2)cn1 | 4 | 5 | 1 | 10 | 1 |
| 21 | COCCCc1cc(CN(C(=O)[C@H]2CNCC[C@@H]2c2ccc(OCCOc3c(Cl)cc(C)cc3Cl)cc2)C2CC2)cc(OCCOC)c1 | 15 | 14 | 1 | 16 | 0 |
| 22 | CC(C)c1ccc2c(c1)OC1(O)c3ccccc3C(=O)C21NC(=O)c1cc(-c2ccccc2)n[nH]1 | 4 | 3 | 1 | 4 | 1 |
| 23 | COCCCc1cc(CN(C(=O)[C@H]2CN(C(=O)OC(C)(C)C)CC[C@@H]2c2ccc(OCCOc3c(Cl)cc(C)cc3Cl)cc2)C2CC2)cc(OCC2(CC(=O)O)CC2)c1 | 15 | 16 | 1 | 0 | 0 |
| 24 | COC[C@H](C)COCc1ccc([C@@]2(O)CCN(C(=O)OC(C)(C)C)C[C@@H]2c2noc(-c3ccccc3CCNC(C)=O)c2Br)cc1 | 9 | 16 | 0 | 13 | 1 |
| 25 | CC1(C)Oc2ccc(CC(=O)NC3CCCc4ccccc43)cc2C(N)C1O | 3 | 6 | 1 | 14 | 1 |
| 26 | O=C(Nc1cccc(Cl)c1)N1CCc2[nH]nc(C(=O)N3CC(F)CO3)c2C1 | 3 | 2 | 1 | 2 | 1 |
| 27 | C[C@@H](O)C[C@H]1OC[C@@H](C2CCCCC2)N(c2cc(C#CC(C)(C)C)sc2C(=O)O)C1=O | 5 | 8 | 1 | 8 | 1 |
| 28 | FC(F)(F)Cn1ncnc1-c1cc2n(n1)-c1cc(C3CCNCC3)ccc1OCC2 | 7 | 11 | 1 | 10 | 1 |
| 29 | C[C@]1(F)CC(F)(F)[C@@](C)(c2cc(N)ccc2F)N=C1N | 6 | 7 | 0 | 8 | 0 |
| 30 | COc1ccc2nccc(C(N)CC[C@@H]3CCN(CCSc4cccs4)C[C@@H]3C(=O)O)c2c1 | 4 | 10 | 1 | 16 | 0 |

| | | | | | | |
|---|---|---|---|---|---|---|
| 31 | N#Cc1cc(F)cc([C@H]2C[C@H](F)CN2c2ccn3ncc(C(N)=O)c3c2)c1 | 9 | 5 | 1 | 7 | 1 |
| 32 | CCN1CCN(CCCCCCOc2c(OC)ccc3c(Nc4c(Cl)cncc4Cl)cc(=O)[nH]c23)CC1 | 10 | 7 | 1 | 11 | 1 |
| 33 | CC(=O)Nc1nc2c(s1)-c1c(c(C3CC3)nn1C1CCN(C(=O)N3CCN(C(C)C)CC3)CC1)CC2 | 3 | 10 | 1 | 9 | 1 |
| 34 | CC(C)(C)OC(=O)N[C@@H]1c2cccnc2[C@H](N)CC[C@H]1c1cccc(F)c1F | 12 | 14 | 0 | 15 | 0 |
| 35 | COc1cc2c(=O)[nH]c(=O)n([C@@H]3O[C@H](CO)[C@H]4OC(C)(C)O[C@H]43)c2cc1OC | 3 | 11 | 1 | 5 | 0 |
| 36 | Cc1cc(Nc2nccc(C(F)(F)F)n2)cc(-c2cnc([C@](C)(O)[C@H]3CC[C@H](C(=O)O)CC3)s2)c1 | 4 | 5 | 1 | 3 | 1 |
| 37 | CC(C)(C)OC(=O)N[C@@H]1c2cccnc2C(=O)CC[C@H]1c1cccc(F)c1F | 10 | 11 | 0 | 16 | 0 |
| 38 | CC(C)S(=O)(=O)N[C@H]1CCOC[C@H]1c1ccc(-c2ccc(C#N)s2)cc1 | 8 | 6 | 1 | 3 | 1 |
| 39 | COc1cc2ncc3c(N)nc(-c4cncc(OCCNCc5ccc(F)cc5)c4)cc3c2cc1OC | 4 | 7 | 1 | 7 | 0 |
| 40 | COCCCc1cc(CN(C(=O)[C@H]2CN(C(=O)OC(C)(C)C)CC[C@@H]2c2ccc(OCCOc3c(Cl)cc(C)cc3Cl)cc2)C2CC2)cc(OCCOC)c1 | 14 | 14 | 1 | 15 | 1 |
| 41 | CN1CCN(c2ccc3c(c2)[nH]c2c(C(N)=O)cc(-c4ccc(O)c(Cl)c4)nc23)CC1 | 8 | 7 | 1 | 7 | 1 |
| 42 | Cn1nccc1[C@]1(O)CCCC[C@H]1O | 2 | 1 | 1 | 1 | 1 |
| 43 | COCCCc1cc(CN(C(=O)[C@H]2CNCC[C@@H]2c2ccc(OCCOc3c(Cl)cc(C)cc3Cl)cc2)C2CC2)cc(OCC(C)(C)O)c1 | 15 | 13 | 1 | 23 | 0 |
| 44 | COc1ccc2c(Nc3c(Cl)cncc3Cl)cc(=O)[nH]c2c1OCCCCCCN1CCCN(C)CC1 | 10 | 7 | 1 | 8 | 1 |
| 45 | COC[C@H](C)COCc1ccc([C@@]2(O)CCN(C(=O)OC(C)(C)C)C[C@@H]2C=NO)cc1 | 6 | 3 | 1 | 8 | 1 |
| 46 | CN1CC[C@@H](c2c(O)cc(O)c3c(=O)cc(-c4ccc(C(F)(F)F)cc4)oc23)[C@@H]1CO | 7 | 13 | 1 | 12 | 0 |
| 47 | C[C@H](c1ccccc1)N1C[C@]2(C(=O)OC(C)(C)C)C=CC[C@@H]2C1=O | 6 | 0 | 0 | 0 | 0 |
| 48 | Cn1oc(=O)nc1/C(=N\OCc1cccc(NC(=O)OCCc2ccccc2)n1)c1ccccc1 | 5 | 3 | 1 | 3 | 1 |
| 49 | COc1ccccc1C[C@H](C[C@H](O[Si](C)(C)C(C)(C)C)[C@H](Cc1ccccc1)NC(=O)OC(C)(C)C)C(=O)O | 6 | 6 | 1 | 2 | 0 |
| 50 | CCC[C@@H](CCO)Nc1nc(N)nc(C)c1Cc1ccc(CN2CCC[C@@H]2C(=O)O)cc1OC | 8 | 6 | 1 | 4 | 1 |
| 51 | CCc1[nH]c(C(=O)N[C@H]2CCN(c3cccc(C(=O)O)c3)C[C@H]2OC)nc1C(F)(F)F | 12 | 5 | 1 | 6 | 1 |
| 52 | O=C(Nc1ccc(-c2cnc(C34CC5CC(CC(C(=O)O)(C5)C3)C4)s2)cc1)Nc1ccccc1F | 5 | 5 | 1 | 8 | 0 |
| 53 | COc1cc(-c2cccc(C(F)(F)F)c2)c(F)cc1-c1nncc2cc(S(=O)(=O)Nc3ccon3)ccc12 | 5 | 4 | 1 | 6 | 1 |
| 54 | Cc1cc(C#N)cnc1C(=O)Nc1ccc(F)c([C@@]2(C)N=C(N)[C@@](C)(F)CC2(F)F)c1 | 7 | 8 | 0 | 10 | 0 |
| 55 | Cc1ccc(C(C)(C)C)cc1S[C@@H]1O[C@H](CO)[C@@H](O)[C@H](O)[C@H]1O | 2 | 1 | 1 | 1 | 1 |
| 56 | Oc1ccc2c3c(ccc2c1)Cc1ccccc1OC3c1ccc(OCCN2CCCCC2)cc1 | 6 | 11 | 0 | 7 | 0 |
| 57 | C[C@@H](n1ncn(-c2ccc(Cl)cc2)c1=O)[C@@]1(c2ccc(F)cc2F)CO1 | 7 | 4 | 1 | 8 | 1 |
| 58 | CCCc1c(Cc2ccc(-c3ccccc3C#N)cc2F)c(=O)n([C@H]2CC[C@H](OCC3(C(C)=O)CCC3)CC2)c2ncnn12 | 8 | 14 | 1 | 6 | 1 |
| 59 | CS(=O)(=O)CCN1CCC(c2ccc3c(c2)-n2nc(-c4ncnn4CC(F)(F)F)cc2CCO3)CC1 | 8 | 15 | 0 | 11 | 1 |
| 60 | CC(C)(CO)n1c(CO)nc2cnc(Br)cc21 | 6 | 2 | 1 | 2 | 1 |
| 61 | COCCCc1cc(O)cc(CN(C(=O)[C@H]2CN(C(=O)OC(C)(C)C)CC[C@@H]2c2ccc(OCCOc3c(Cl)cc(C)cc3Cl)cc2)C2CC2)c1 | 13 | 11 | 1 | 13 | 1 |
| 62 | CC(C)N1CCC(Oc2ccc3c(c2)cc2n3[C@H](C)CN(CCO)C2=O)CC1 | 6 | 6 | 0 | 8 | 1 |
| 63 | COCCCc1cc(CN(C(=O)[C@H]2CN(C(=O)OC(C)(C)C)CC[C@@H]2c2ccc(OCCOc3c(Cl)cc(C)cc3Cl)cc2)C2CC2)cc(OC[C@@H]2C[C@H]2C(=O)OCC(=O)N(C)C)c1 | 16 | 3 | 0 | 3 | 0 |
| 64 | CC(=O)NC[C@H]1CN(c2ccc3c(c2)CCCc2c(C(C)C)n[nH]c2-3)C(=O)O1 | 3 | 3 | 1 | 3 | 1 |
| 65 | COC[C@H](C)COCc1ccc([C@@]2(O)CCN(C(=O)OC(C)(C)C)C[C@@H]2CO)cc1 | 4 | 3 | 1 | 3 | 1 |
| 66 | O=C1C=C(N2CCNCC2)C(C2CCCCC2)CC1 | 2 | 3 | 1 | 1 | 0 |
| 67 | C[Si](C)(C)CCOCn1c(O[C@@H]2CO[C@@H]3[C@H](O)CO[C@H]23)nc2cc(Cl)c(-c3ccc(- | 5 | 5 | 0 | 1 | 0 |

| | | | | | | |
|---|---|---|---|---|---|---|
| | c4ccc(N=S(C)(=O)N5CCCC5)cc4)cc3)nc21 | | | | | |
| 68 | C/C=C/c1c(N)nc(-c2ccc(Cl)c(OC)c2F)nc1C(=O)OC | 3 | 2 | 1 | 2 | 1 |
| 69 | COCCCc1cc(CN(C(=O)[C@H]2CN(C(=O)OC(C)(C)C)CC[C@@H]2c2ccc(OCCOc3c(Cl)cc(C)cc3Cl)cc2)C2CC2)cc(OCc2ccc(-c3nnn[nH]3)cc2)c1 | 15 | 16 | 1 | 14 | 1 |
| 70 | CC[C@@H](OC(=O)c1ccccc1)[C@H]1CCCN(C(=O)OC(C)(C)C)C1 | 2 | 1 | 1 | 2 | 1 |
| 71 | ClCc1ccc2c(c1)Nc1nccnc1S2 | 4 | 3 | 1 | 4 | 1 |
| 72 | CC(=O)NC[C@H]1CC[C@@H](C(=O)Nc2cccc(OC(F)(F)F)c2)N1 | 10 | 2 | 1 | 2 | 1 |
| 73 | COC(=O)[C@H]1CN(C(=O)OC(C)(C)C)CC[C@@H]1c1ccc(OCCOc2c(Cl)cc(C)cc2Cl)cc1 | 6 | 6 | 1 | 6 | 1 |
| 74 | CN(C1CC1)[C@H]1CC[C@H]2[C@@H]3CC=C4C[C@@H](O)CC[C@]4(C)[C@H]3CC[C@@]21C | 4 | 3 | 0 | 2 | 0 |
| 75 | Nc1ncnc2c1c(-c1ccc3ccc(-c4ccccc4)nc3c1)cn2C1CCC1 | 3 | 3 | 1 | 4 | 1 |
| 76 | CC(=O)OCc1nc2cnc(Br)cc2n1C(C)(C)COC(C)=O | 5 | 2 | 1 | 2 | 1 |
| 77 | CCOC(=O)/C(N)=N/Nc1cc(Cl)ccc1[N+](=O)[O-] | 2 | 1 | 0 | 0 | 0 |
| 78 | C=C(C[C@@H](Cc1ccc(-c2ccccc2)cc1)NC(=O)OC(C)(C)C)C(=O)OC | 3 | 1 | 0 | 0 | 0 |
| 79 | COC[C@H](C)COCc1ccc([C@@]2(O)CCN(C(=O)OC(C)(C)C)C[C@@H]2C=O)cc1 | 5 | 2 | 1 | 7 | 1 |
| 80 | O=S(=O)(C#Cc1ccc(Cl)cc1)N1CCNCC1 | 3 | 1 | 1 | 2 | 1 |
| 81 | OC[C@H]1O[C@](O)(c2ccc(Cl)c(Cc3ccc(C#Cc4cnccn4)cc3)c2)[C@H](O)[C@@H](O)[C@@H]1O | 6 | 5 | 0 | 16 | 0 |
| 82 | O[C@H]1C[C@H](c2cnn3c(N[C@H]4CCc5ccccc54)ncnc23)C=C1COCc1ccccc1 | 3 | 12 | 1 | 12 | 0 |
| 83 | C[C@@H](n1ccn(-c2ccc(Cl)cc2)c1=O)[C@@]1(c2ccc(F)cc2F)CO1 | 7 | 4 | 1 | 8 | 1 |
| 84 | CN(C)CC(OCC1(c2ccc(F)cc2)CCN(C(=O)OC(C)(C)C)CC1)c1cc(Cl)cc2cn(COCC[Si](C)(C)C)nc12 | 11 | 7 | 1 | 3 | 1 |
| 85 | CC(=O)N1c2ccc(N3CCNCC3)nc2[C@H](Nc2ccccc2)[C@@H](C)[C@@H]1C1CC1 | 8 | 9 | 0 | 11 | 0 |
| 86 | COC(=O)c1ccc2c(c1)C=CC(=CCl)CO2 | 6 | 6 | 1 | 0 | 0 |
| 87 | CCOP(=O)(Cc1ccc(Nc2ncc(C(F)(F)F)c(Nc3ccc([C@H]4CC[C@H](N5CCN(C)CC5)CC4)c4c3C(=O)N(C)C4)n2)c(OC)c1)OCC | 8 | 6 | 1 | 7 | 1 |
| 88 | CCCC[Sn](/C=C/C1(O)C(C)=CC2(CC1(C)C(F)(F)F)OC(C)C(C)O2)(CCCC)CCCC | 6 | 0 | 0 | 0 | 0 |
| 89 | CC(C)(C)OC(=O)N[C@@H]1c2cccnc2[C@H](O)CC[C@H]1c1cccc(F)c1F | 9 | 12 | 0 | 9 | 0 |
| 90 | C[C@H](O[Si](C)(C)C(C)(C)C)[C@@H]1CC(=O)CC(C)(C)N1 | 5 | 9 | 1 | 9 | 1 |
| 91 | CC1=NC2(N=C1N)c1cc(-c3cc(F)cc(C#N)c3)ccc1CCC21CC1 | 6 | 8 | 0 | 4 | 0 |
| 92 | CC(C)(C)OC(=O)NC1(c2nc(NCc3ccccn3)c3c(Cl)ccn3n2)CC1 | 4 | 10 | 1 | 9 | 0 |
| 93 | CC(=O)c1ccc2c(c1)C=CC(O)(CO)CO2 | 4 | 8 | 1 | 0 | 0 |
| 94 | COC[C@H](C)COCc1ccc([C@@]2(O)CCNC[C@@H]2c2noc(-c3ccccc3CCNC(C)=O)c2Br)cc1 | 10 | 15 | 1 | 14 | 1 |
| 95 | CC(=O)NC[C@H]1CC[C@@H](C(=O)Nc2cccc(OC(F)(F)F)c2)N1C(=O)OCc1ccccc1 | 9 | 3 | 1 | 7 | 1 |
| 96 | COC(=O)C(C)(C)CC(CCCCO)CCCc1cccnc1 | 12 | 11 | 1 | 4 | 1 |
| 97 | CCCCCCC(C)(C)c1ccc([C@H]2C[C@@H](O)CCN2C=O)c(OCc2ccccc2)c1 | 4 | 12 | 1 | 17 | 1 |
| 98 | COC(=O)[C@@]12C[C@H]1C=CCCCC[C@H](NC(=O)CC1CC1)C(=O)N1C[C@H](Oc3nc(-c4ccccn4)nc4ccsc34)C[C@H]1C(=O)N2 | 9 | 17 | 1 | 15 | 0 |
| 99 | CC1(C)OCC(C)(CO)N(Cc2ccccc2)C1=O | 5 | 5 | 0 | 3 | 0 |
| 100 | COC(=O)CC12CCC(c3ccc(-c4ccc(NC(=O)c5nc(C)oc5C(F)(F)F)cc4)cc3)(CC1)CO2 | 7 | 9 | 1 | 2 | 1 |
| 101 | COCCCc1cc(CN(C(=O)[C@H]2CN(C(=O)OC(C)(C)C)CC[C@@H]2c2ccc(OCCOc3c(Cl)cc(C)cc3Cl)cc2)C2CC2)cc(OCCC(C)(C)C(=O)OC)c1 | 14 | 11 | 1 | 18 | 1 |

| | | | | | | |
|---|---|---|---|---|---|---|
| 102 | COc1ccc(C[C@H](C[C@H](O[Si](C)(C)C(C)(C)C)[C@H](CC2CCCCC2)NC(=O)OC(C)(C)C)C(=O)O)c(OC)c1OC | 8 | 10 | 1 | 2 | 0 |
| 103 | CC(=O)N1c2ccc(N3CCNCC3)cc2[C@H](Nc2ccccc2)[C@@H](C)[C@@H]1C | 5 | 9 | 0 | 6 | 1 |
| 104 | CC1(C)OCC(C)(C=O)N(Cc2ccccc2)C1=O | 6 | 3 | 1 | 11 | 1 |
| 105 | CNC(=O)c1ccc2ncc(C(=O)OCn3ccnc3)c(Nc3ccc(OC)cc3)c2c1 | 4 | 1 | 1 | 1 | 1 |
| 106 | CCCCCN1C(=O)C2(CNC(=O)c3cc4c(cc32)OCO4)c2ccccc21 | 6 | 0 | 0 | 0 | 0 |
| 107 | CC1=NC2(N=C1N)c1cc(-c3cc(Cl)cc(C#N)c3)ccc1CCC21CC1 | 6 | 11 | 0 | 7 | 0 |
| 108 | COCC1=CS[C@@H]2C(NC(=O)Cc3cccs3)C(=O)N2C1C(=O)O | 4 | 16 | 0 | 13 | 0 |
| 109 | C[C@@H](O)c1nc2cnc3ccsc3c2n1[C@H]1CC[C@H](CO)CC1 | 3 | 6 | 1 | 6 | 1 |
| 110 | CC1(C)OCC(C)(CO[Si](C)(C)C(C)(C)C)N(Cc2ccccc2)C1=O | 4 | 6 | 0 | 11 | 1 |
| 111 | CN1CCN(c2ccc3c(c2)[nH]c2c(C(N)=O)cc(-c4cnn(C)c4)nc23)CC1 | 8 | 7 | 1 | 16 | 0 |
| 112 | COc1cc2ncc3c(N)nc(-c4cncc(OCCN(Cc5ccc(F)cc5)C(=O)OC(C)(C)C)c4)cc3c2cc1OC | 3 | 17 | 0 | 5 | 1 |
| 113 | OC1CCCc2sc(-c3ccncc3)cc2C1c1ccc(Cl)cc1 | 7 | 5 | 1 | 3 | 1 |
| 114 | O=C(O)[C@H]1CCOC[C@H]1c1ccc(I)cc1 | 5 | 5 | 1 | 1 | 1 |
| 115 | CC(C)(C)OC(=O)[C@@]12C=CC[C@@H]1CN(C(=O)OCc1ccccc1)C2 | 9 | 1 | 0 | 4 | 0 |
| 116 | Cn1oc(=O)nc1/C(=N\OCc1cccc(N)n1)c1ccccc1 | 4 | 3 | 1 | 3 | 1 |
| 117 | C[C@@]1(O)[C@H](O)[C@@H](CO)O[C@H]1n1cc(-c2ccccc2)c2c(N)ncnc21 | 3 | 5 | 1 | 4 | 1 |
| 118 | COCCCc1cc(CN(C(=O)[C@H]2CN(C(=O)OC(C)(C)C)CC[C@@H]2c2ccc(OCCOc3c(Cl)cc(C)cc3Cl)cc2)C2CC2)cc(OCCC(C)(C)C(=O)O)c1 | 15 | 13 | 1 | 0 | 0 |
| 119 | C[Si](C)(C)CCOCn1c(O[C@@H]2CO[C@@H]3[C@H](O)CO[C@H]23)nc2cc(Cl)c(-c3ccc(-c4ccc(N=S(C)(=O)N5CCC5)cc4)cc3)nc21 | 5 | 6 | 1 | 1 | 0 |
| 120 | COc1cc2c(Oc3cc(C)c(C)nc3-c3cccc(C)n3)ccnc2cc1OCCNCCCO | 7 | 7 | 1 | 7 | 1 |
| 121 | COC(=O)[C@@H]1CCCC2(CCCCC2)[C@H]1O | 2 | 3 | 1 | 1 | 1 |
| 122 | COC(=O)c1ccc2c(c1)C=CC(=C(Cl)Cl)CO2 | 6 | 0 | 0 | 3 | 0 |
| 123 | COCCCc1cc(CN(C(=O)[C@H]2CN(C(=O)OC(C)(C)C)CC[C@@H]2c2ccc(OCCOc3c(Cl)cc(C)cc3Cl)cc2)C2CC2)cc(OCC2(CC#N)CC2)c1 | 15 | 14 | 1 | 13 | 1 |
| 124 | CC1(C)COCc2nc3cnc(Br)cc3n21 | 7 | 1 | 1 | 1 | 1 |
| 125 | C=CC[C@@H]1C(=O)N([C@H](C)c2ccccc2)C[C@@]1(C=C)C(=O)OC(C)(C)C | 5 | 1 | 0 | 1 | 0 |
| 126 | COCCCc1cc(CN(C(=O)[C@H]2CN(C(=O)OC(C)(C)C)CC[C@@H]2c2ccc(OCCOc3c(Cl)cc(C)cc3Cl)cc2)C2CC2)cc(OC[C@@H]2C[C@H]2C(=O)Oc2ccc3c(c2)CCC3)c1 | 16 | 17 | 1 | 0 | 0 |
| 127 | Nc1ccc(OC2(S(=O)(=O)c3ccccc3)CC2)nc1Cl | 4 | 1 | 1 | 1 | 1 |
| 128 | CC1=NC2(N=C1N)c1cc(-c3cc(F)cc(F)c3)ccc1CCC21CC1 | 6 | 1 | 0 | 7 | 0 |
| 129 | CC#CCn1c(Br)nc(C=O)c1C(=O)OC | 2 | 2 | 1 | 2 | 1 |
| 130 | NS(=O)(=O)OC[C@H]1C[C@@H](c2cnn3c(N[C@H]4CCc5ccccc54)ncnc23)C[C@@H]1O | 5 | 10 | 1 | 16 | 0 |
| 131 | CC(C)S(=O)(=O)N[C@H]1CCOC[C@H]1c1ccc(I)cc1 | 7 | 3 | 1 | 3 | 1 |
| 132 | O=C(Nc1ccc(-c2cnc(C34CC5CC(CC(C(=O)O)(C5)C3)C4)s2)cc1)Nc1cccc(C(F)(F)F)c1 | 5 | 5 | 1 | 6 | 1 |
| 133 | NC(=O)CN1CCC(c2ccc3c(c2)-n2nc(-c4ncnn4CC(F)(F)F)cc2CCO3)CC1 | 8 | 16 | 0 | 11 | 1 |
| 134 | C[C@H](O[Si](C)(C)C(C)(C)C)[C@@H]1CC(O)CC(C)(C)N1 | 6 | 8 | 1 | 15 | 0 |
| 135 | Cc1ccn2cc(CO)n(-c3ccccc3)c(=O)c12 | 8 | 13 | 1 | 6 | 1 |

| | | | | | | |
|---|---|---|---|---|---|---|
| 136 | COC(=O)c1ccc(-c2ccc(O[C@H]3O[C@H](CO)[C@@H](O)[C@H](O)[C@@H]3O)c(Cl)c2)cc1 | 3 | 1 | 1 | 12 | 0 |
| 137 | CC(=O)Nc1cccc(N2CCN(CCc3nn(C)c(=O)n3CC3CCCCC3)CC2)c1 | 8 | 3 | 1 | 3 | 1 |
| 138 | C=C(C[C@@H](Cc1ccc(-c2ccccc2)cc1)NC(=O)OC(C)(C)C)C(=O)O | 2 | 4 | 1 | 6 | 1 |
| 139 | COCCCCc1cc(CN(C(=O)[C@H]2CN(C(=O)OC(C)(C)C)CC[C@@H]2c2ccc(OCCOc3c(Cl)cc(C)cc3Cl)cc2)C2CC2)cc(OC[C@@H]2C[C@H]2CO)c1 | 15 | 12 | 1 | 14 | 1 |
| 140 | Cc1cc(Nc2nccc(C(F)(F)F)n2)cc(-c2cnc([C@@](C)(O)[C@H]3CC[C@H](C(=O)O)CC3)s2)c1 | 4 | 5 | 1 | 3 | 1 |
| 141 | COC(=O)CCc1cc2cc(-c3noc(-c4ccc(OC(C)C)c(Cl)c4)n3)ccc2n1C | 5 | 5 | 1 | 4 | 1 |
| 142 | Cn1nccc1[C@]12CCCC[C@H]1OC(=O)O2 | 3 | 0 | 0 | 0 | 0 |
| 143 | C[C@@]1(O)[C@H](O)[C@@H](CO)O[C@H]1n1ccc2c(N)nc(N)nc21 | 3 | 6 | 1 | 3 | 1 |
| 144 | O=S(=O)(c1cc(C(F)(F)F)ccc1Br)C1CCOC(c2ccc(Cl)cc2)C1 | 4 | 4 | 1 | 3 | 1 |
| 145 | CCc1[nH]c(C(=O)N[C@H]2CCN(c3nc(C(=O)O)c(C)s3)C[C@H]2OC)nc1C(F)(F)F | 10 | 4 | 1 | 7 | 1 |
| 146 | CC(C)(C)OC(=O)N1CC=C(c2ccc3c(c2)-n2nc(-c4ncnn4CC(F)(F)F)cc2CCO3)CC1 | 6 | 0 | 0 | 0 | 0 |
| 147 | CC(C)(C)OC(=O)N[C@@H]1c2cccnc2[C@H](N2C(=O)c3ccccc3C2=O)CC[C@H]1c1cccc(F)c1F | 11 | 14 | 0 | 14 | 0 |
| 148 | CC1=NC2(N=C1N)c1cc(Br)ccc1CCC21CC1 | 5 | 5 | 0 | 3 | 0 |
| 149 | COCCCCc1cc(CN(C(=O)[C@H]2CN(C(=O)OC(C)(C)C)CC[C@@H]2c2ccc(OCCOc3c(Cl)cc(C)cc3Cl)cc2)C2CC2)cc(OCCCC#N)c1 | 14 | 11 | 1 | 15 | 1 |
| 150 | CCOC(=O)[C@@H]1C[C@H]1COc1cc(CCCCOC)cc(CN(C(=O)[C@H]2CN(C(=O)OC(C)(C)C)CC[C@@H]2c2ccc(OCCOc3c(Cl)cc(C)cc3Cl)cc2)C2CC2)c1 | 14 | 12 | 1 | 16 | 1 |
| 151 | OCCN1CCC(c2ccc3c(c2)-n2nc(-c4ncnn4CC(F)(F)F)cc2CCO3)CC1 | 8 | 15 | 0 | 11 | 1 |
| 152 | CCOC(=O)[C@@H]1C[C@H]1COc1cc(CCCCOC)cc(CN(C(=O)[C@H]2CNCC[C@@H]2c2ccc(OCCOc3c(Cl)cc(C)cc3Cl)cc2)C2CC2)c1 | 15 | 12 | 1 | 16 | 1 |
| 153 | CC(C)c1ccc2c(c1)OC1(O)c3ccccc3C(=O)C21NC(=O)C(=O)c1cccs1 | 4 | 2 | 1 | 3 | 1 |
| 154 | CC(=O)OCC1=CS[C@@H]2C(NC(=O)Cc3cccs3)C(=O)N2C1C(=O)O | 2 | 10 | 1 | 13 | 1 |
| 155 | NC(=O)c1cc(-c2ccc3c(cnn3CCN3CCOCC3)c2)nc2c1[nH]c1cc(N3CCOCC3)ccc12 | 7 | 14 | 1 | 9 | 1 |
| 156 | O=C(NC[C@@H]1CC[C@@H](COc2ccccc2)OC1)c1ccc(O)cc1 | 9 | 5 | 1 | 9 | 1 |
| 157 | COCCCCc1cc(CN(C(=O)[C@H]2CN(C(=O)OC(C)(C)C)CC[C@@H]2c2ccc(OCCOc3c(Cl)cc(C)cc3Cl)cc2)C2CC2)cc(OS(=O)(=O)C(F)(F)F)c1 | 14 | 11 | 1 | 15 | 1 |
| 158 | CCOC(=O)/C=C1\CC[C@@H](c2cccc(F)c2F)[C@H](NC(=O)OC(C)(C)C)c2cccnc21 | 11 | 0 | 0 | 16 | 0 |
| 159 | C[C@H](c1ccccc1)N1C[C@H]2CC=C[C@@]2(C(=O)OC(C)(C)C)C1 | 8 | 1 | 0 | 3 | 0 |
| 160 | COCCCCc1cc(CN(C(=O)[C@H]2CN(C(=O)OC(C)(C)C)CC[C@@H]2c2ccc(OCCOc3c(Cl)cc(C)cc3Cl)cc2)C2CC2)cc(OC[C@@H]2C[C@H]2C(=O)O)c1 | 15 | 0 | 0 | 0 | 0 |
| 161 | CCc1[nH]c(C(=O)N[C@H]2CCN(c3cccc(C(=O)OC(C)(C)C)c3)C[C@H]2OC)nc1C(F)(F)F | 11 | 5 | 1 | 6 | 1 |
| 162 | C=CCC1C=C(C)CC(C)CC(OC)C2OC(O)(C(=O)C(=O)N3CCCCC3C(=O)OC(C(C)=CC3CCC(O)C(OC)C3)C(C)C=CC1O)C(C)C2OC | 2 | 0 | 0 | 0 | 0 |
| 163 | Cc1nnc(SCC2=C(C(=O)O)N3C(=O)C(NC(=O)C(O)c4csc(N)n4)[C@H]3SC2)s1 | 7 | 1 | 1 | 1 | 1 |
| 164 | CCOC(=O)CCc1cc2cc(-c3noc(-c4ccc(OC(C)C)c(Cl)c4)n3)ccc2[nH]1 | 4 | 4 | 1 | 4 | 1 |
| 165 | CCOC(=O)[C@H]1[C@H](c2ccc3c(c2)OC(F)(F)O3)C1(C)C | 2 | 4 | 1 | 4 | 1 |
| 166 | C[C@@H]1CNC(=O)c2cc3cc(OCCCN4CCCCC4)ccc3n21 | 5 | 8 | 1 | 4 | 1 |

| | | | | | | |
|---|---|---|---|---|---|---|
| 167 | CC(C)(C)OC(=O)NC1(c2nc(O)c3c(Cl)ccn3n2)CC1 | 3 | 0 | 0 | 0 | 0 |
| 168 | C[C@@H]1CCCN1CCc1nnc2cc(Br)ccc2c1O | 2 | 2 | 1 | 3 | 1 |
| 169 | CC(C)c1ccc(OC2CCC3(CC2)N[C@@H](C(=O)O)C(C)(C)S3)cc1 | 2 | 9 | 0 | 4 | 1 |
| 170 | CC1=NC2(N=C1N)c1cc(-c3cncc(Cl)c3)ccc1CCC21CC1 | 6 | 9 | 0 | 7 | 0 |
| 171 | CC(=O)Nc1nc2c(s1)-c1c(c(C3CC3)nn1C1CCNCC1)CC2 | 2 | 9 | 1 | 7 | 1 |
| 172 | COCCCCC1(CNC(=O)C2CCCNC2)c2ccccc2Oc2ccccc21 | 12 | 3 | 1 | 3 | 1 |
| 173 | N[C@H]1CCOC[C@H]1c1ccc(I)cc1 | 6 | 2 | 1 | 2 | 1 |
| 174 | COCCCc1cc(CN(C(=O)[C@H]2CN(C(=O)OC(C)(C)C)CC[C@@H]2c2ccc(OCCOc3c(Cl)cc(C)cc3Cl)cc2)C2CC2)cc(OCC2(CC(=O)OC)CC2)c1 | 14 | 12 | 1 | 16 | 1 |
| 175 | C[C@H](C(=O)NCCF)N(C)C(=O)c1ccc2c(c1)c1c(n2C)CCC(C2CCOCC2)C1 | 7 | 4 | 1 | 4 | 1 |
| 176 | CCC/C=C/[C@H]1CC[C@H](C(CO)CO)CC1 | 6 | 5 | 1 | 9 | 0 |
| 177 | O=C(OCc1ccccc1)N1CC[C@H]2CCCNC[C@H]21 | 4 | 4 | 1 | 10 | 1 |
| 178 | OC[C@H]1O[C@@H](SC(c2ccccc2)(c2ccccc2)c2ccccc2)[C@H](O)[C@@H](O)[C@@H]1O | 2 | 3 | 1 | 7 | 1 |
| 179 | CC(C)(C)OC(=O)N[C@@H]1c2cccnc2[C@@H](O)CC[C@H]1c1cccc(F)c1F | 10 | 12 | 0 | 9 | 0 |
| 180 | CCc1ccc(S(=O)(=O)NC2c3cc(CC(=O)NC4CCCc5ccccc54)ccc3OC(C)(C)C2O)cc1 | 4 | 7 | 1 | 9 | 1 |
| 181 | CCCCCOC(=O)N1CCN(C(=O)[C@H](CCCO[Si](c2ccccc2)(c2ccccc2)C(C)(C)C)NC(=O)OC(C)(C)C)CC1 | 3 | 3 | 1 | 2 | 1 |
| 182 | CC(=O)OC1CCC2C(CO)CCC12 | 7 | 2 | 1 | 3 | 1 |
| 183 | C[C@H]1OC[C@@H]2CC[C@@H](c3nc(-c4ccc(C(=O)NC5cc(C(F)(F)F)ccn5)cc4)c4c(N)nccn34)CN2C1=O | 10 | 16 | 0 | 16 | 1 |
| 184 | CCCC(C)(C)OC(=O)Nc1nc(C(O)C(=O)NC2C(=O)N3C(C(=O)O)=C(CSc4nnc(C)s4)CS[C@H]23)cs1 | 6 | 2 | 1 | 2 | 1 |
| 185 | NC1=Nc2ccc(F)cc2C2CCCC12 | 4 | 3 | 1 | 2 | 1 |
| 186 | OC[C@H]1C[C@@H](c2cnn3c(N[C@H]4CCc5ccccc54)ncnc23)C[C@@H]1O | 4 | 9 | 1 | 14 | 0 |
| 187 | O=C(CO)N1CCC(c2ccc3c(c2)-n2nc(-c4ncnn4CC(F)(F)F)cc2CCO3)CC1 | 8 | 15 | 0 | 11 | 1 |
| 188 | COCCCCc1cc(CN(C(=O)[C@H]2CN(C(=O)OC(C)(C)C)CC[C@@H]2c2ccc(OCCOc3c(Cl)cc(C)cc3Cl)cc2)C2CC2)cc(OCC(C)(C)O)c1 | 14 | 12 | 1 | 15 | 1 |
| 189 | Oc1ccc2c(c1)[C@H](c1ccc(OCCN3CCCC3)cc1)[C@H](c1ccccc1)CO2 | 7 | 6 | 1 | 4 | 1 |
| | **Average** | | | 0.705 | | 0.647 |

**Supplementary Table S9 | The JMC2025 dataset[3-32].**

| ID | SMILES | Name | JMC step | Beam Search | | Top-*k* | | Top-*p* | |
|---|---|---|---|---|---|---|---|---|---|
| | | | | step | success | step | success | step | success |
| 1 | O=C1C2=CC(NC3CCN(C(C)=O)CC3)=CC=C2CCN1C[C@H](O)CN4CC5=C(C=CC=C5)CC4 | A1 | 6 | 7 | 1 | 6 | 1 | 7 | 1 |
| 1 | O=C1C2=CC(NC3CCN(C(C4CC4)=O)CC3)=CC=C2C5(CCCC5)CN1C[C@H](O)CN6CC7=C(C=CC=C7)CC6 | C9 | 3 | 7 | 1 | 8 | 0 | 8 | 1 |
| 2 | O=C(N(C1=CC=C(N2CCC(C)CC2)C(C(F)(F)F)=C1)C)C3=CC=C(CN4CCOCC4)C=C3 | 40 | 3 | 2 | 1 | 2 | 1 | 2 | 1 |
| 3 | CC1=C(C=CC=C1C2=CC=CC=C2)NC3=NC(OC)=CC(CNCCO)=N3 | P1 | 6 | 2 | 1 | 4 | 1 | 4 | 1 |
| 3 | CC1=C(C=CC=C1C2=CC=CC=C2)NC3=NC(OC)=CC(CNCC(O)=O)=N3 | P13 | 7 | 3 | 1 | 7 | 1 | 6 | 1 |
| 4 | ClC1=C(C2=CC=C1CNC(C3CC4(C3)CNC4)=O)N(C=C2C5=CC(C(N)=S)=CC=C5)C6CCOCC6 | 22 | 13 | 6 | 1 | 6 | 1 | 8 | 0 |
| 4 | ClC1=C(C2=CC=C1CNC(C3CC4(C3)CCN(CC4)CCNC(C)=O)=O)N(C=C2C5=CC(C(N)=S)=CC=C5)C6CCOCC6 | 50 | 16 | 7 | 1 | 8 | 1 | 8 | 1 |
| 5 | O=C(CN1CCN(CCN(CCN(CC1)CC(O)=O)CC(O)=O)CC(N[C@@H](CCC(NCCCN2C=CN=C2[N+]([O-])=O)=O)C(N3CCN(CCCOC4=CC=C(C5=C4)N=CC=C5C(NCC(N6[C@@H](CC(F)(C6)F)C#N)=O)=O)CC3)=O)=O)O | DOTA-NI-FAPI-04 | 7 | 1 | 0 | 1 | 0 | 1 | 0 |
| 6 | CN1C=C(C2=CC=CC=C21)C3=CC=NC(N(C4=CC=CC(N)=C4)C)=N3 | A35 | 4 | 3 | 1 | 3 | 1 | 4 | 1 |
| 6 | NC1=CC(NC2=NC=CC(OC3=CC=C4NC=CC4=C3)=N2)=CC=C1 | A76 | 3 | 1 | 1 | 1 | 1 | 1 | 1 |
| 7 | OC(C1=C(N=C(C2=CC=C(C=C2)OCC)S1)C)=O | CIB-L25 | 4 | 0 | 1 | 0 | 1 | 0 | 1 |
| 7 | CC1=C(SC(CC2=CC=C(C=C2)OCC3CC3)=N1)C(O)=O | CIB-L63 | 7 | 3 | 1 | 3 | 1 | 3 | 1 |
| 8 | OC(C1=CC=C(C(OC2CCN(CC2)C)=C1)CC3=CC=CC=C3)=O | 2a | 4 | 3 | 1 | 3 | 1 | 3 | 1 |
| 8 | CN1CCC(OC2=CC(CC(O)=O)=CC=C2CC3=CC=CC=C3)CC1 | 2b | 9 | 3 | 1 | 7 | 1 | 2 | 1 |
| 9 | O=C1NC2=CC(C(OC)=O)=CC=C2/C1=C(C3=CC=CC=C3)/NC4=CC=C(C=C4)N(C)C(C(C)N5CCN(CC5)C)=O | 10a | 7 | 3 | 1 | 3 | 1 | 3 | 1 |
| 10 | CN(C1=CC=C(OCC2=C(C=CC=C2)[N+]([O-])=O)C=C1)N=O | N4 | 2 | 2 | 1 | 2 | 1 | 3 | 1 |
| 10 | CN(C1=CC=C(OCC2=CC=C(C=C2)[N+]([O-])=O)C=C1)N=O | N1 | 2 | 2 | 1 | 3 | 1 | 3 | 1 |
| 11 | CCCCC1=CC(C2=CC=CC=C2)=NN1CC3=CC=C(C=C3)C(NO)=O | 4 | 4 | 2 | 1 | 3 | 1 | 3 | 1 |
| 11 | O=C(C1=CC=C(CN2N=C(C=C2C(F)(F)F)C(NC3=CC=CC=C3)=O)C=C1)NO | 9 | 5 | 2 | 1 | 3 | 1 | 3 | 1 |
| 12 | CN1N=C2C(N(C=CC2=C1CC3=CC=C(N=C3)C(F)(F)F)CC)=O | 5 | 4 | 5 | 1 | 8 | 1 | 7 | 1 |
| 12 | CN1N=C2C(N(C=C(C2=C1CC3=CC=C(N=C3)C(F)(F)F)CC))=O | 12 | 6 | 9 | 1 | 8 | 1 | 8 | 1 |
| 13 | CC1=C(C2=CC=C(CNC([C@@H]3C[C@H](CN3C([C@H](C(C)(C)C)NC(CN(CCCCCCOC4=CC5=C(N=C4)C(C6CCN(CC6)C(C7=CC=C(C=C7)OC(F)(F)F)=O)=NC=N5)C)=O)=O)O)=O)C=C2)SC=N1 | 7 | 3 | 10 | 1 | 9 | 1 | 2 | 0 |

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| 13 | CC1=C(C2=CC=C(CNC([C@@H]3C[C@H](CN3C([C@H](C(C)(C)C)NC(CCOCCNC(COC4=CC5=C(N=C4)C(C6CCN(CC6)C(C7=CC=C(C=C7)OC(F)(F)F)=O)=NC=N5)=O)=O)=O)O)=O)C=C2)SC=N1 | 3 | 3 | 8 | 0 | 9 | 0 | 0 | 0 |
| 14 | O=S(C1=CC=C(C=C1)C(N(CC2=CC=C(C=C2)NCCO)CC3CC3)=O)(NC4=CC=CC=C4)=O | 85 | 3 | 2 | 1 | 3 | 1 | 2 | 1 |
| 14 | O=S(C1=CC=C(C=C1)CN(CC2=CC=CC=C2)C(CCC)=O)(NC3=CC=CC=C3)=O | 9 | 3 | 2 | 1 | 2 | 1 | 2 | 1 |
| 15 | N#CC1=C2C(C3=CC=C(N=C3)N4CC5CC(C4)N5CC6=CN=C(C=C6)OC)=CC(C7=CN(N=C7)CCCCC(N[C@@H](C(C)(C)C)C(N8C[C@H](O)C[C@H]8C(N[C@@H](C)C9=CC=C(C%10=C(C)N=CS%10)C=C9)=O)=O)=O)=CN2N=C1 | 6 | 9 | 0 | 0 | 0 | 0 | 0 | 0 |
| 15 | O=C(NC1=O)CCC1N2C(C3=CC=CC(C#CCCN4C=C(C=N4)C5=CN6N=CC(C#N)=C6C(C7=CC=C(N8CC9CC(N9CC%10=CN=C(OC)C=C%10)C8)N=C7)=C5)=C3C2)=O | 12 | 3 | 0 | 0 | 0 | 0 | 0 | 0 |
| 16 | O=C1N(CC2=CC=CC=C2)CCC3=C1N(C4=CN(C(C)C)N=C4)N=C3C(NC5=CC=CC(C(O)=O)=C5)=O | 10 | 4 | 8 | 0 | 7 | 1 | 8 | 1 |
| 16 | OC(C1=CC=CC=C1NC(C2=NN(CC3=CC=C(C=C3)OC)C4=C2CN(CC5=CC=CC=C5)CC4)=O)=O | 45 | 3 | 5 | 1 | 6 | 1 | 2 | 1 |
| 17 | FC(C=C1)=CC=C1CNC2=C3C(N(N=C3CC)=NC(C4=CC=C5C(COB5O)=C4)=N2 | 63 | 4 | 6 | 1 | 5 | 1 | 11 | 1 |
| 17 | FC1=CC=C(C=C1)CNC2=NC(C3=CC=C(C=C3)P(O)(O)=O)=NC4=C2C(CC)=NN4C | 53 | 4 | 8 | 1 | 6 | 1 | 8 | 1 |
| 18 | ClC1=CC=CC=C1OC2=NC=CC(CNS(C3=CC=C4C(CCN4C(C)=O)=C3)(=O)=O)=C2 | MR49915 | 3 | 2 | 1 | 4 | 1 | 6 | 1 |
| 19 | O=C(C1=CC(Cl)=CC=C1O)NC2=C(Cl)C=C([N+]([O-])=O)C=C2 | A | 1 | 0 | 1 | 0 | 1 | 0 | 1 |
| 20 | CCCCN(C(C1)=O)C2=C1C=C(C3=CC=CC=C3C4=NN=NN4)C=C2 | 6 | 5 | 3 | 1 | 4 | 1 | 4 | 1 |
| 20 | O=C(OC)C1=CC=CC=C1C2=CC=C(/C=N/C3=CC=CC=C3)C=C2 | 42 | 2 | 2 | 1 | 1 | 1 | 1 | 1 |
| 21 | ClC1=NC2=C(N=C(NCC3CC3)N=C2)N(C4=CC=C(OC(F)F)C=C4)C1=O | 51b | 5 | 7 | 1 | 4 | 1 | 5 | 1 |
| 21 | O=C1C(C2=CC=N(C(C)C)C=N3)=C3C=C2)=NC4=C(C=C(C5CC5)C=C4)N1C6=CC=C(OC(F)F)C=C6 | 32 | 6 | 4 | 1 | 8 | 1 | 6 | 1 |
| 22 | CN1CCN(C2=C(NC(C3=CC(N)=C(F)C(C)=C3Cl)=O)C=C(C4=CC=C(NC(CCCC(NO)=O)=O)C=C4)C=C2)CC1 | 27a | 11 | 5 | 1 | 5 | 1 | 7 | 1 |
| 22 | CN1CCN(C2=C(NC(C3=CC(N)=C(F)C(C)=C3Cl)=O)C=C(C4=CC=C(C(NCCCC(NO)=O)=O)C=C4)C=C2)CC1 | 24a | 12 | 5 | 1 | 5 | 1 | 4 | 1 |
| 23 | CCC1=C(NC(C(O)=O)=C1C2=CC=CC3=C2C=CN3)C4=CC=CC5=C4C=CN5 | 1j | 5 | 6 | 1 | 5 | 1 | 5 | 1 |
| 23 | CC(C)C1=C(NC(C(O)=O)=C1C2=CC=CC3=C2C=CN3)C4=CC=CC5=C4C=CN5 | 2a | 5 | 5 | 1 | 9 | 1 | 7 | 1 |
| 24 | CN(C1=CC=C(C2=C1)N=CC=C2C(NCC(N3[C@@H](CC(F)(C3)F)C#N)=O)=O)CCCN4CCN(CC4)C(OC(C)(C)C)=O | 17 | 15 | 4 | 1 | 1 | 1 | 1 | 1 |
| 25 | FC1=C(C2=CC=C3C(CC4(OC3=C2)CCC4)=O)N=C(NC5=CC=C(C=N5)N6CCN(CC6)C)N=C1 | 1 | 4 | 5 | 1 | 3 | 1 | 5 | 1 |
| 25 | FC1=C(C2=CC=C3C(CC4(OC3=C2)CCN(CC4)C(C)=O)=O)N=C(NC5=CC=C(C=N5)N6CCN(CC6)C)N=C1 | 14 | 6 | 3 | 1 | 5 | 1 | 5 | 1 |
| 26 | O=C(C1=CC(COCOC)=C(C(COCOC)=C1)[Sn](CCCC)(CCCC)CCCC)NCC(O)=O | 13 | 9 | 8 | 0 | 8 | 0 | 8 | 0 |
| 26 | O=C(C1=CC(C(NCC(O)=O)=O)=CC(C(N(C)C)=O)=C1[Sn](CCCC)(CCCC)CCCC)N(C)C | 19 | 9 | 6 | 1 | 8 | 1 | 7 | 1 |

| # | SMILES | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| 27 | CC1=NOC2=CC(OCC(NCCNC(COC3=CC=CC(C(N4C5C(NC(CC5)=O)=O)=O)=C3C4=O)=O)=O)=C(C=C21)NS(C6=C(C=CC(Br)=C6)OC)(=O)=O | 6a | 6 | 4 | 1 | 3 | 1 | 5 | 1 |
| 27 | CC1=NOC2=CC(OCC(NCCNC3=CC=CC(C(N4C5CCC(NC5=O)=O)=O)=C3C4=O)=O)=C(C=C21)NS(C6=C(C=CC(Br)=C6)OC)(=O)=O | 7a | 4 | 4 | 1 | 4 | 1 | 4 | 1 |
| 28 | O=C(C1=C(SC2=C1CCOC2)CC3=CC=C(C=C3F)C4=CC=C5OC=CC5=C4)NCC6=CC=C(C=C6)C(O)=O | 28 | 2 | 7 | 1 | 8 | 0 | 8 | 0 |
| 28 | O=C(C1=C(SC2=C1CCOC2)CC3=CC=C(C=C3)C4=CC=C(C=C4)OC)NC5(C6=CC=C(C=C6)C(O)=O)CC5 | 2 | 4 | 6 | 1 | 4 | 1 | 5 | 1 |
| 29 | O=C1COC2=CC=C(C=C2N1CC3=NOC(C4CCC4)=N3)S(NC5=CC=CC(C)=C5C)(=O)=O | 19 | 3 | 0 | 1 | 0 | 1 | 0 | 1 |
| 30 | FC1=CC(F)=C(C=C1)C(O)(CN2C=NC=N2)CN(C)CC3=CC=C(OCCCC(N4CCN(CC4)CC5=CC=C(C=C5)N(CC)C(C6=C(O)C=C(O)C(C7=CC=C(C)C=C7)=C6)=O)=O)C=C3 | MM3 | 12 | 10 | 1 | 8 | 1 | 10 | 1 |
| 30 | FC1=CC(F)=C(C=C1)C(CN2CCN(CC3=CC=C(N(C(C4=C(OCC5=CC=CC=C5)C=C(OCC6=CC=CC=C6)C(C7=CC=C(C)C=C7)=C4)=O)CC)C=C3)CC2)(CN8C=NC=N8)O | Bn-MM4 | 8 | 8 | 1 | 11 | 0 | 1 | 0 |
| | **Average** | | | | **0.887** | | **0.849** | | **0.830** |

Reference:

1    Perera, D. *et al.* A platform for automated nanomole-scale reaction screening and micromole-scale synthesis in flow. *Science* **359**, 429–434 (2018). https://doi.org/10.1126/science.aap9112

2    Jiasheng Guo, C. Y., Kenan Li, Yijian Zhang, Guoqiang Wang, Shuhua Li, Hao Dong. Retrosynthesis Zero: Self-Improving Global Synthesis Planning Using Reinforcement Learning. *Journal of Chemical Theory and Computation* (2024).

3    Abeje, Y. E. *et al.* Impact of Linker Composition on VHL PROTAC Cell Permeability. *Journal of Medicinal Chemistry* **68**, 638–657 (2025). https://doi.org/10.1021/acs.jmedchem.4c02492

4    Ackloo, S. *et al.* A Target Class Ligandability Evaluation of WD40 Repeat-Containing Proteins. *Journal of Medicinal Chemistry* **68**, 1092–1112 (2025). https://doi.org/10.1021/acs.jmedchem.4c02010

5    Argade, M. D. *et al.* Guardians at the Gate: Optimization of Small Molecule Entry Inhibitors of Ebola and Marburg Viruses. *Journal of Medicinal Chemistry* **68**, 135–155 (2025). https://doi.org/10.1021/acs.jmedchem.4c01646

6    Baldwin, A. G. *et al.* Discovery of MDI-114215: A Potent and Selective LIMK Inhibitor To Treat Fragile X Syndrome. *Journal of Medicinal Chemistry* **68**, 719–752 (2025). https://doi.org/10.1021/acs.jmedchem.4c02694

7    Chen, P. *et al.* Monomeric or Homodimer Conjugates of Fibroblast Activation Protein Inhibitor and Cyanine 7 Bearing a Meso-Chloride as Near-Infrared Fluorescence Probes: Design, Synthesis, and Comparative In Vivo Imaging of Distinct Breast Cancer Subtypes. *Journal of Medicinal Chemistry* **68**, 1417–1432 (2025). https://doi.org/10.1021/acs.jmedchem.4c01968

8    Cheng, Z. *et al.* Subtle Structural Modifications Spanning from EP4 Antagonism to EP2/EP4 Dual Antagonism: A Novel Class of Thienocyclic-Based Derivatives. *Journal of Medicinal Chemistry* **68**, 1587–1607 (2025). https://doi.org/10.1021/acs.jmedchem.4c02241

9    Cui, H. *et al.* Discovery of N-Phenyl-5-propyl-1H-pyrazole-3-carboxamide, with Selective Inhibition and Degradation of HDAC6 for the Treatment of Acute Liver Injury. *Journal of Medicinal Chemistry* **68**, 531–554 (2025). https://doi.org/10.1021/acs.jmedchem.4c02341

10   Etxebeste-Mitxeltorena, M. *et al.* Modulation of Mitochondria–Endoplasmic Reticulum Contacts (MERCs) by Small Molecules as a New Strategy for Restoring Lipid Metabolism in an Amyotrophic Lateral Sclerosis Model. *Journal of Medicinal Chemistry* **68**, 1179–1194 (2025). https://doi.org/10.1021/acs.jmedchem.4c01368

11   Gasonoo, M. *et al.* Discovery of an Orally Efficacious Pyrazolo[3,4-d]pyrimidine Benzoxaborole as a Potent Inhibitor of Cryptosporidium. *Journal of Medicinal Chemistry* **68**, 832–849 (2025). https://doi.org/10.1021/acs.jmedchem.4c02805

12   Hernandez-Olmos, V. *et al.* Identification of a Chemical Probe for BLT2 Activation by Scaffold Hopping. *Journal of Medicinal Chemistry* **68**, 1195–1221 (2025). https://doi.org/10.1021/acs.jmedchem.4c01617

13   Hippman, R. S. *et al.* Structure–Activity Relationship Studies of an Autophagy Inhibitor That Targets the ATG14L-Beclin1 Protein–Protein Interaction Reveal Modifications That Improve Potency and Solubility and Maintain Selectivity. *Journal of Medicinal Chemistry* **68**, 1645–1667 (2025). https://doi.org/10.1021/acs.jmedchem.4c02312

14   Hirata, S. *et al.* In Vivo Stability Improvement of Astatobenzene Derivatives by Introducing Neighboring Substituents. *Journal of Medicinal Chemistry* **68**, 1540–1552 (2025). https://doi.org/10.1021/acs.jmedchem.4c02188

15    Huang, G. *et al.* Structure-Based Development of Novel Spiro-Piperidine ASH1L Inhibitors. *Journal of Medicinal Chemistry* **68**, 174–195 (2025). https://doi.org/10.1021/acs.jmedchem.4c01673

16    Kuhne, S. *et al.* Probing the Histamine H1 Receptor Binding Site to Explore Ligand Binding Kinetics. *Journal of Medicinal Chemistry* **68**, 448–464 (2025). https://doi.org/10.1021/acs.jmedchem.4c02043

17    Li, C. *et al.* Discovery of 2(1H)-Quinoxalinone Derivatives as Potent and Selective MAT2A Inhibitors for the Treatment of MTAP-Deficient Cancers. *Journal of Medicinal Chemistry* **68**, 1222–1244 (2025). https://doi.org/10.1021/acs.jmedchem.4c01635

18    Li, Q.-Q. *et al.* Design, Synthesis, and Biological Evaluation of 3,4-Dihydroisoquinolin-1(2H)-one Derivatives as Protein Arginine Methyltransferase 5 Inhibitors for the Treatment of Non-Hodgkin's Lymphoma. *Journal of Medicinal Chemistry* **68**, 108–134 (2025). https://doi.org/10.1021/acs.jmedchem.4c01548

19    Long, G. *et al.* Discovery of WDR5–MLL1 and HDAC Dual-Target Inhibitors for the Treatment of Acute Myeloid Leukemia. *Journal of Medicinal Chemistry* **68**, 1260–1279 (2025). https://doi.org/10.1021/acs.jmedchem.4c01720

20    Luo, Y. *et al.* Development of [68Ga]Ga/[177Lu]Lu-DOTA-NI-FAPI-04 Containing a Nitroimidazole Moiety as New FAPI Radiotracers with Improved Tumor Uptake and Retention. *Journal of Medicinal Chemistry* **68**, 348–360 (2025). https://doi.org/10.1021/acs.jmedchem.4c02015

21    Ma, D. *et al.* Discovery of Novel 4,5,6,7-Tetrahydro-7H-pyrazolo[3,4-c]pyridin-7-one Derivatives as Orally Efficacious ATX Allosteric Inhibitors for the Treatment of Pulmonary Fibrosis. *Journal of Medicinal Chemistry* **68**, 792–818 (2025). https://doi.org/10.1021/acs.jmedchem.4c02719

22    Nie, Z. *et al.* Structure-Based Discovery and Development of Highly Potent Dihydroorotate Dehydrogenase Inhibitors for Malaria Chemoprevention. *Journal of Medicinal Chemistry* **68**, 590–637 (2025). https://doi.org/10.1021/acs.jmedchem.4c02394

23    Shen, C. *et al.* Design, Synthesis, and Biological Evaluation of 2-Arylaminopyrimidine Derivatives as Dual Cathepsin L and JAK Inhibitors for the Treatment of Acute Lung Injury. *Journal of Medicinal Chemistry* **68**, 361–386 (2025). https://doi.org/10.1021/acs.jmedchem.4c02030

24    Shen, H. *et al.* Discovery of a Potent and Selective GSPT1 Molecular Glue Degrader for the Treatment of Castration-Resistant Prostate Cancer. *Journal of Medicinal Chemistry* **68**, 1553–1571 (2025). https://doi.org/10.1021/acs.jmedchem.4c02205

25    Shi, H. *et al.* Design, Synthesis, and Antitumor Activity Evaluation of 2-Phenylthiazole-5-Carboxylic Acid Derivatives Targeting Transactivation Response RNA-Binding Protein 2. *Journal of Medicinal Chemistry* **68**, 421–447 (2025). https://doi.org/10.1021/acs.jmedchem.4c02041

26    Shi, W. *et al.* Discovery of Potent and Selective CDK4/6 Inhibitors for the Treatment of Chemotherapy-Induced Myelosuppression. *Journal of Medicinal Chemistry* **68**, 1446–1472 (2025). https://doi.org/10.1021/acs.jmedchem.4c02080

27    Sun, C. *et al.* Novel PD-L1/VISTA Dual Inhibitor as Potential Immunotherapy Agents. *Journal of Medicinal Chemistry* **68**, 156–173 (2025). https://doi.org/10.1021/acs.jmedchem.4c01640

28    Yang, W., Zhou, W. & Gou, S. Hypoxia Activated Nitric Oxide Donor Compounds for the Prevention and Treatment of Myocardial Hypoxia-Induced Injury. *Journal of Medicinal Chemistry* **68**, 491–505 (2025). https://doi.org/10.1021/acs.jmedchem.4c02132

29      Yao, H. *et al.* The Discovery of a Novel AXL/Triple Angiokinase Inhibitor Based on 6-Chloro-Substituted Indolinone and Side Chain Methyl Substitution Inhibiting Pancreatic Cancer Growth and Metastasis. *Journal of Medicinal Chemistry* **68**, 465–490 (2025). https://doi.org/10.1021/acs.jmedchem.4c02130

30      Yao, Y. *et al.* Development of PRC1 Inhibitors Employing Fragment-Based Approach and NMR-Guided Optimization. *Journal of Medicinal Chemistry* **68**, 1382–1396 (2025). https://doi.org/10.1021/acs.jmedchem.4c01955

31      Zhang, M. *et al.* Lanosterol 14α-Demethylase (CYP51)/Heat Shock Protein 90 (Hsp90) Dual Inhibitors for the Treatment of Invasive Candidiasis. *Journal of Medicinal Chemistry* **68**, 1668–1681 (2025). https://doi.org/10.1021/acs.jmedchem.4c02305

32      Zhang, Q. *et al.* Discovery of an Efficacious RET PROTAC Degrader with Enhanced Antiproliferative Activity against Resistant Cancer Cells Harboring RET Solvent-Front Mutations. *Journal of Medicinal Chemistry* **68**, 753–775 (2025). https://doi.org/10.1021/acs.jmedchem.4c02692