

CHEMBOMAS: ACCELERATED BO FOR SCIENTIFIC DISCOVERY IN CHEMISTRY WITH LLM-ENHANCED MULTI-AGENT SYSTEM

Dong Han^{1*} Zhehong Ai^{1*} Pengxiang Cai^{1*} Shanya Lu²

Jianpeng Chen¹ Zihao Ye⁴ Shuzhou Sun¹ Ben Gao¹ Lingli Ge¹ Weida Wang¹
 Xiangxin Zhou¹ Xihui Liu³ Mao Su¹ Wanli Ouyang¹ Lei Bai¹ Dongzhan Zhou¹
 Tao Xu^{2†} Yuqiang Li^{1†} Shufei Zhang^{1†}

¹ Shanghai Artificial Intelligence Laboratory, Shanghai, China

² Tongji University, Shanghai, China

³ The University of Hong Kong, Hong Kong, China

⁴ Fudan University, Shanghai, China

ABSTRACT

Bayesian optimization (BO) is a powerful tool for scientific discovery in chemistry, yet its efficiency is often hampered by the sparse experimental data and vast search space. Here, we introduce **ChemBOMAS**: a large language model (LLM)-enhanced multi-agent system that accelerates BO through synergistic data- and knowledge-driven strategies. Firstly, the data-driven strategy involves an 8B-scale LLM regressor fine-tuned on a mere 1% labeled samples for pseudo-data generation, robustly initializing the optimization process. Secondly, the knowledge-driven strategy employs a hybrid Retrieval-Augmented Generation approach to guide LLM in dividing the search space while mitigating LLM hallucinations. An Upper Confidence Bound algorithm then identifies high-potential subspaces within this established partition. Across the LLM-refined subspaces and supported by LLM-generated data, BO achieves the improvement of effectiveness and efficiency. Comprehensive evaluations across multiple chemical benchmarks demonstrate that ChemBOMAS set a new state-of-the-art, accelerating optimization efficiency by up to 5-fold compared to baseline methods.

1 INTRODUCTION

Manual experimentation and traditional control variable methods have long underpinned chemical discovery, yet they remain labor-intensive and time-consuming, slowing the generation of new scientific insights Xie et al. (2023); Tom et al. (2024). To address these constraints, automated or self-driving laboratories integrate robotic execution with AI algorithms, delivering high throughput, precision, and efficiency Seifrid et al. (2022); Chen et al. (2024); Ai et al. (2024a). Within these experimental platforms, Bayesian Optimization (BO) algorithms are widely recognized as a crucial decision-making tool for experiment design Guo et al. (2023); Abolhasani and Kumacheva (2023); Chen et al. (2023); Ai et al. (2024b). BO enables efficient navigation of complex experimental variable spaces and converges toward optimal reaction conditions or material compositions by integrating prior data, constructing probabilistic surrogate models, quantifying uncertainty, and iteratively selecting the most informative subsequent experiments Shields et al. (2021a).

Despite BO achieving remarkable success in complex scientific domains, particularly chemistry, it still contends with two major obstacles: (I) the scarcity and high cost of experimental observations during the early optimization stages, and (II) the multitude of reaction parameters that inflate the search into high-dimensional design spaces Shahriari et al. (2015); Wang et al. (2023). The two obstacles exacerbate the limitations of vanilla BO, also known as the "cold start" problem and the

*Equal contribution.

†Correspondence to taoxu@tongji.edu.cn, liyuqiang@pjlab.org.cn, zhangshufei@pjlab.org.cn.

"curse of dimensionality", frequently leading to slow convergence Guo et al. (2023). Without effective acceleration strategies, the protracted optimization process may yield only marginal improvements, which could cause researchers to abandon the search before discovering the optimal conditions.

Several strategies have been proposed to accelerate BO, including search space partitioning Wang et al. (2020a), specialized encoding embeddings Tripp et al. (2020); Moriconi et al. (2020); Nayeibi et al. (2019), pseudo-data generation Yin et al. (2023), and transfer across similar tasks Swersky et al. (2013). However, when these acceleration strategies are applied to the intricate chemical reactions, two critical shortcomings emerge. First, most approaches employ a single acceleration technique, which might be insufficient for the chemical optimization problems with multiple demands, such as exploration of diverse reaction parameter combinations while overcoming data scarcity in the early-stage. Second, current acceleration methods are predominantly data-driven. Because chemical reaction pathways differ widely in their underlying kinetics and thermodynamics, a purely statistical BO framework frequently expends resources in chemically implausible regions of the search space, missing opportunities to leverage mechanistic insight that could guide the search more efficiently.

To overcome these limitations, we propose **ChemBOMAS**, an LLM-Enhanced **Multi-Agent System** specifically designed for accelerated **Bayesian Optimization** in chemistry. ChemBOMAS synergistically integrates two LLM-powered modules: a **knowledge-driven search space decomposition module** and a **data-driven pseudo-data generation module**. The knowledge-driven module employs an LLM-powered agent to reason over existing chemical knowledge (e.g., literature, databases), intelligently decompose the vast search space and identify promising candidate regions, dynamically pruning the search space for better BO efficiency. Simultaneously, the data-driven module utilizes a fine-tuned LLM to generate informative pseudo-data points across the entire search space. These pseudo-data not only warm-start the BO process but also inform the knowledge-driven module’s sub-space selection. This closed-loop interaction enables ChemBOMAS to achieve superior optimization efficiency and convergence speed even under extreme data scarcity.

The effectiveness of ChemBOMAS was rigorously evaluated. We conducted extensive experiments on four chemical performance optimization benchmarks, demonstrating consistent improvements in optimal results, convergence speed, initialization performance, and robustness compared to various baseline methods. Crucially, ablation studies confirmed that the synergy between the knowledge-driven and data-driven strategies is essential for creating a highly efficient and robust optimization framework. Additionally, the practical utility and real-world applicability of ChemBOMAS were validated through a previously unreported wet-lab experiment.

Our main contributions are summarized as follows:

1. We systematically investigated how LLM-based approaches could address two key limitations in BO for scientific discovery in chemistry: data scarcity and inefficiency in the vast search spaces.
2. We propose ChemBOMAS, a novel framework that synergistically combines LLM-powered knowledge and data strategies, which enables efficient BO in chemical performance optimization tasks. It consists of a knowledge-driven module to decompose the search space and a data-driven module to generate pseudo-data.
3. ChemBOMAS significantly accelerates BO under limited data (1% of all search space) on four chemical benchmarks. It achieves accelerated convergence and superior final performance compared to four relevant baseline methods, improving optimal results by approximately 3–10% and converging about $2\text{--}5\times$ faster.

2 RELATED WORK

Large Language Models (LLMs) offer synergistic potential with Bayesian Optimization (BO) to address traditional BO limitations (e.g., sample inefficiency, cold starts) by providing prior knowledge Souza et al. (2021), enhancing surrogate models Liu et al. (2024); Nguyen et al. (2024); Ramos et al. (2023a), automating acquisition function design Austin et al. (2024), and enabling optimization in novel problem representations. Prior work has explored LLM-driven BO improvements in warm-starting, surrogate modeling, candidate generation, acquisition function design, and search space understanding.

However, directly replacing core BO modules with LLMs introduces significant challenges. LLM "hallucinations" can mislead optimization, compromising reliability. Furthermore, the direct suitability of LLMs as surrogates or acquisition functions is limited by concerns regarding uncertainty quantification, theoretical guarantees, computational cost, efficiency in low-data regimes, adaptability to specific numerical tasks, and interpretability.

On another front, some techniques such as LA-MCTS Wang et al. (2020a) was proposed, which employ tree structures to decompose the search space Wang et al. (2023; 2024; 2019). Some works propose hierarchical Bayesian optimization Moriconi et al. (2020); Reker et al. (2020). These approaches offer valuable strategies for managing and navigating complex optimization landscapes. Unlike previous works, we focus on robustly integrating LLM knowledge to enhance BO, leveraging their strengths as auxiliary tools while mitigating weaknesses such as hallucinations, to achieve this synergy over substitution.

3 METHODOLOGY

3.1 PROBLEM SETUP

This work aims to significantly improve the efficiency of searching a task’s variable space for the optimal combination that maximizes the objective function.

3.2 THE FRAMEWORK OF CHEMBOMAS

As illustrated in Figure 1, we propose ChemBOMAS, an LLM-enhanced multi-agent optimization framework that systematically integrates data-driven and knowledge-driven strategies. First, the data-driven strategy utilizes a pre-trained and fine-tuned LLM regressor to generate pseudo-data, thereby robustly initializing the optimization process. Second, the knowledge-driven strategy employs a hybrid Retrieval-Augmented Generation (RAG) approach, which guides an LLM to partition the search space based on variable impact ranking and property similarity. Third, an Upper Confidence Bound (UCB) algorithm then identifies the most promising subspaces from this partition. Finally, BO is performed within the selected subspaces, supported by the LLM-generated pseudo-data, leading to enhanced effectiveness and efficiency. The complete algorithm process can be seen in Appendix F. The two strategies are detailed below.

3.3 DATA-DRIVEN STRATEGY: LLM-GENERATED PSEUDO DATA

An LLM-based regression model was constructed and utilized in three sequential steps to generate pseudo-data for optimization initialization.

Step 1: Pre-training. The base LLaMA 3.1 model Grattafiori et al. (2024) was pre-trained on the Pistachio dataset $\mathcal{D}_{\text{chem}}$ to enhance its representational ability for chemical reactions. The dataset was formatted as Q&A pairs where, given reactants \mathbf{R} and products \mathbf{P} , the model learns to predict the corresponding reaction conditions $\mathbf{c} = (c_1, c_2, \dots, c_T)$, thereby avoiding direct exposure to objective performance labels. Pre-training employed a Causal Language Modeling loss: $\mathcal{L}_{\text{pre-train}} = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{\text{chem}}} \left[\sum_{t=1}^T \log p(x_t | x_{<t}) \right]$, where t denotes the token index, x_t represents the t -th token, and T is the sequence length.

Step 2: Fine-tuning. The pre-trained model was subsequently fine-tuned on a small labeled dataset $\mathcal{D}_{\text{labeled}} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$, which only comprises 1% of the data, by integrating a regression head. A reaction configuration \mathbf{x} (including \mathbf{R} , \mathbf{P} , and \mathbf{c}) is fed into the LLM via prompt engineering. The final hidden state $\mathbf{h}_T = \text{LLM}(\mathbf{x})^{[T]}$ is then projected to a reaction performance prediction \hat{y} via an MLP: $\hat{y} = f_{\theta_{\text{MLP}}}(\mathbf{h}_T) = f_{\theta_{\text{MLP}}}(\text{LLM}(\mathbf{x})^{[T]})$. Fine-tuning used Low-Rank Adaptation (LoRA) Hu et al. (2022) with rank $r = 8$, introducing adaptable parameters ϕ_{LoRA} alongside the frozen pre-trained weights θ_{LLM} . The MLP parameters θ_{MLP} were fully trained to minimize an L2-loss with regularization:

$$\mathcal{L}_{\text{fine-tune}} = \frac{1}{|\mathcal{D}_{\text{labeled}}|} \sum_{(\mathbf{x}, y) \in \mathcal{D}_{\text{labeled}}} \left\| f_{\theta_{\text{MLP}}}(\text{LLM}_{\theta_{\text{LLM}}, \phi_{\text{LoRA}}}^{[T]}(\mathbf{x})) - y \right\|_2^2 + \lambda \|\theta_{\text{MLP}}\|_2^2 \quad (1)$$

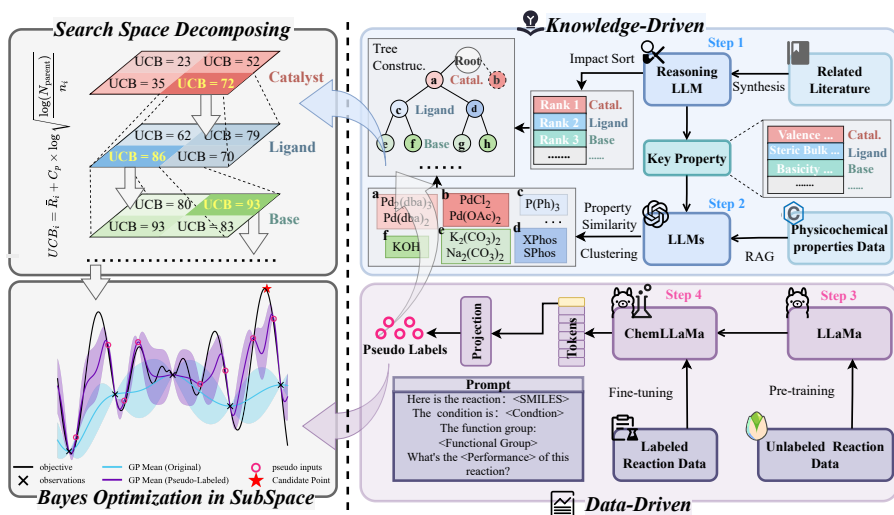


Figure 1: ChemBOMAS: A synergistic knowledge- and data-driven framework for efficient Bayesian Optimization. The framework operates as a closed-loop system: the **knowledge-driven module** decomposes the search space into subspaces using LLM-extracted chemical insights, followed by a UCB algorithm to select promising subspaces; the **data-driven module** generates pseudo-data to initialize both the subspace selection and the Bayesian Optimization process within the selected subspaces. The two modules interact iteratively, with real data from optimization feedback refining subsequent search directions.

Step 3: Pseudo-data Generation and Utilization. The fine-tuned LLM regressor was used to generate pseudo-data for all unsampled data points, forming a pseudo-dataset $\mathcal{D}_{\text{pseudo}} = \{(\mathbf{x}_k, \hat{y}_k)\}_{k=1}^M$, where $\hat{y}_k = f_{\theta_{\text{MLP}}}(\text{LLM}_{\theta_{\text{LLM}}, \phi_{\text{LoRA}}}^{[T]}(\mathbf{x}_k))$ was used to initialize a UCB algorithm. The UCB algorithm (Section 3.4) then identifies the high-potential subspaces. BO is then conducted within these subspaces (Section 3.5), leveraging both the selected pseudo-data and limited real data to accelerate surrogate model fitting. To mitigate the influence of the noise in the pseudo-data, a refinement strategy based on data similarity and reverse-order removal is applied (see Appendix D for details).

3.4 KNOWLEDGE-DRIVEN STRATEGY: LLM-DIVIDED SEARCH SPACE

To efficiently identify high-potential regions in the vast chemical reaction parameter space, we construct a Subspace Tree Search module using the GPT-4o API in three steps.

Step 1: LLM-Guided Space Partitioning. The n -dimensional variable space is defined as $\mathcal{X} = \{C_i\}_{i=1}^n$, containing n categories of chemicals involved in the reaction, where each $C_i = \{x_{i,1}, \dots, x_{i,k}\}$ represented a category including k candidate substances. A hybrid Retrieval-Augmented Generation (RAG) approach integrates multi-source information (literature, professional databases, web search) to facilitate the LLM’s decisions and minimize hallucination. The LLM first ranks the chemical categories C_i by their importance to the chemical reaction, generating an ordered sequence $\mathcal{O} = (o_1, \dots, o_n)$. Subsequently, for each chemical category C_i , the LLM identifies key influencing physicochemical properties $p_{i,1}, p_{i,2}, \dots$ and clusters the candidates based on these property values. This partitions each C_i into a collection of separate subspaces $\Pi_i = \{S_{i,1}, \dots, S_{i,q_i}\}$, where candidate substances within each subspace $S_{i,l}$ share similar properties.

Step 2: Hierarchical Search Tree Construction. A hierarchical tree is built based on the category importance order \mathcal{O} and the clustering results Π_i . The l -th layer of the tree corresponds to the l -th most important category C_l and contains nodes representing its q_l clusters. Each path from the root to a leaf node defines a unique search subspace as the Cartesian product of n clusters, resulting in a total of $\prod_{i=1}^n q_i$ disjoint subspaces that comprehensively partition the original space.

Step 3: UCB-based Subspace Selection. A UCB algorithm is employed to explore the tree and identify promising subspaces. Starting from the root, UCB selects child nodes layer-by-layer until reaching a leaf node. The UCB value for a child node i is computed as: $UCB_i = \bar{R}_i + C_p \times \log \sqrt{\frac{\log(N_{parent})}{n_i}}$, where \bar{R}_i is the average performance value (exploitation), N_{parent} is the parent’s visit count, n_i is the child’s visit count, and C_p is an exploration constant. At each layer, the top-5 nodes by UCB value are selected for further exploration. This path traversal pinpoints high-value subspaces for subsequent BO. The UCB values are updated dynamically as BO progresses and new samples are acquired.

3.5 BAYESIAN OPTIMIZATION IN CHEMBOMAS

BO is performed within the promising subspaces identified by the preceding modules, leveraging the LLM-generated pseudo-data for initialization. The procedure consists of two main steps.

Step 1: Surrogate Modeling. A Gaussian Process (GP) surrogate model, using a Matérn kernel with constant scaling and a white noise kernel, is fitted to the combined set of actual observations and pseudo-data points, which serve as an informative prior. This model provides, for any unsampled point \mathbf{x} in the target subspaces, a posterior distribution over the performance value, characterized by a mean function $\mu(\mathbf{x})$ and a variance function $\sigma^2(\mathbf{x})$.

Step 2: Acquisition Function Optimization. An acquisition function $\alpha(\mathbf{x})$, such as UCB or Expected Improvement (EI), is used to recommend the next sample by balancing exploration (high uncertainty) and exploitation (high predicted mean). The next query point is selected by maximizing $\alpha(\mathbf{x})$ over the unsampled points within the high-potential subspaces: $\mathbf{x}_{next} = \arg \max \mathbf{x} \in \mathcal{X}_{sub} \alpha(\mathbf{x})$. This point is then evaluated to obtain a new real observation, which updates the GP surrogate model for the next iteration.

4 EXPERIMENT

4.1 DATA

The pre-training phase employed a subset of the Pistachio dataset containing approximately 50,000 chemical reaction entries, none of which contained objective performance labels.

For the fine-tuning and Bayesian Optimization (BO) phases, three benchmark datasets were used: Suzuki Perera et al. (2018), Arylation Shields et al. (2021b), and Buchwald Ahneman et al. (2018). In each case, only 1% of the labeled data was randomly selected for fine-tuning the LLM regressor; the effectiveness of this data volume is analyzed in Appendix I.2. For the Buchwald dataset, which exhibits inconsistencies in reactants and products across entries, two consistent subsets were constructed, denoted $\text{Buchwald}_{\text{sub-1}}$ and $\text{Buchwald}_{\text{sub-2}}$, to serve as rational benchmarks for the optimization task. Detailed statistics for all benchmarks are provided in Appendix C.

4.2 EXPERIMENT SETUP

The LLM regressor in data-driven module was trained on $2 \times$ NVIDIA A800 GPUs. For fine-tuning the LLM regressor, the hyperparameters were set as follows: learning rate of 1×10^{-4} , batch size of 24, and 100 training epochs.

In the knowledge-driven module, the search tree was constructed using a UCB policy with an exploration constant $\kappa = 20$. The BO process was run for 40 iterations, initialized with 1% of the data as the prior and sampling 0.1% of the dataset in each iteration. The acquisition function and other BO configurations were kept consistent with the baseline methods for a fair comparison. Each optimization experiment was independently repeated 5 times with different random seeds, and the average performance across these runs is reported as the final result. (The prompts for LLM clustering and further implementation details are provided in the Appendix G.)

4.3 PERFORMANCE COMPARISON

4.3.1 REGRESSION MODELS

The quality of the pseudo-data is directly influenced by the prediction accuracy of the regression model. We evaluated the performance prediction accuracy of ChemBOMAS against three categories of existing regression models on three chemical datasets: 1) general-purpose LLMs (GPT series) with zero-shot inference; 2) open LLMs fine-tuned on 1% labeled data; and 3) scientific LLMs with molecule pre-training fine-tuned on 1% labeled data. The prediction metrics for each model are summarized in Table 1, from which two key observations can be drawn.

Table 1: Comparative performance of various LLM-based regression models on the chemical performance prediction task.

Model	Suzuki			Arylation			Buchwald		
	MSE↓	MAE↓	R ² ↑	MSE↓	MAE↓	R ² ↑	MSE↓	MAE↓	R ² ↑
<i>General-purpose LLMs with zero-shot inference</i>									
GPT-4o	2207.17	40.02	-1.80	2702.58	44.86	-2.63	1512.44	33.60	-1.03
GPT-5	1218.93	30.34	-0.55	1515.81	33.68	-1.04	1516.62	33.55	-1.04
<i>Open LLMs fine-tuned on 1% labeled data</i>									
Qwen3-7B	820.48	22.10	-0.04	848.51	22.52	-0.14	998.22	25.25	-0.34
GLM4-9B	593.49	18.94	0.25	739.20	20.78	0.01	719.72	20.77	0.00
LLaMa-3.1-8B	685.55	20.50	0.13	679.72	19.57	0.09	739.27	20.57	0.01
<i>Scientific LLMs with molecule pre-training and fine-tuned on 1% labeled data</i>									
MolT5-Large	1081.23	25.13	-0.37	1094.86	25.40	-0.47	1098.16	25.37	-0.47
Galactica-1.3B	727.18	22.23	0.08	785.01	21.83	-0.05	857.79	22.54	-0.15
ChemBOMAS	<u>633.68</u>	<u>19.47</u>	<u>0.20</u>	650.00	19.55	0.13	593.76	18.52	0.20

First, ChemBOMAS demonstrated superior effectiveness and versatility in chemical performance regression. As shown in Table 1, ChemBOMAS achieved the highest prediction accuracy on the Arylation and Buchwald datasets, with R² scores exceeding the second-best model by 2000% and 140%, respectively. On the Suzuki dataset, ChemBOMAS outperformed six of the seven compared models, ranking second only to GLM4-9B. However, the poor generality of GLM4-9B is evident from its near-zero R² scores on the other two datasets.

Second, task-specific fine-tuning proved essential. Despite their general capabilities, the off-the-shelf general-purpose LLMs GPT-4o OpenAI et al. (2024) and GPT-5 OpenAI (2025) performed poorly on this specialized regression task, consistently yielding strongly negative R² scores, lower than most fine-tuned models, which also confirms that these chemical datasets were not part of their training data. Among the open-source models, LLaMa-3.1-8B Grattafiori et al. (2024) exhibited a favorable balance of prediction accuracy and generalization, justifying its selection as the base model for the data-driven module of ChemBOMAS.

Furthermore, we investigated the impact of fine-tuning data volume on pseudo-data quality and BO performance (see Tables 7 and 8). The predictive performance improved gradually as data volume increased from 0.25% to 4%. Notably, the R² value first turned positive and consistently exceeded 0.1 across all datasets at the 1% volume. Table 8 indicates that BO’s optimization performance does not linearly correlate with the regression model’s R²; a value between 0.1 and 0.2 is sufficient for BO to identify high-performing conditions. Therefore, using 1% data volume for fine-tuning represents a rational and effective trade-off between cost and performance.

4.3.2 CLUSTER METHODS

We evaluated the impact of the search tree structure on BO by comparing scenarios with and without a tree, as well as trees constructed using three distinct strategies: expert guidance (Expert Guided), data-driven approach (ChemBOMAS_{d-d}), and knowledge-driven approach (ChemBOMAS_{k-d}). All methods were initialized with an identical, fixed set of pseudo-data to ensure a fair comparison.

Table 2: Comparison of Bayesian Optimization performance using different search space decomposition strategies: standard BO without a tree (BO w/o tree), expert-guided clustering (Expert Guided), data-driven clustering (ChemBOMAS_{d-d}) and knowledge-driven clustering (ChemBOMAS_{k-d}).

Method	Suzuki				Arylation			
	Best Found (%)	Initial (%)	95% Max Iter	Best Iter	Best Found (%)	Initial (%)	95% Max Iter	Best Iter
BO _{w/o tree}	83.16	54.04	16	38	80.45	64.64	8	35
Expert-Guided	96.15	61.96	3	3	79.67	78.71	1	36
ChemBOMAS _{d-d}	89.51	72.98	29	29	79.67	67.75	2	28
ChemBOMAS _{k-d}	96.15	72.98	3	3	82.23	78.71	1	39

Method	Buchwald _{sub-1}				Buchwald _{sub-2}			
	Best Found (%)	Initial (%)	95% Max Iter	Best Iter	Best Found (%)	Initial (%)	95% Max Iter	Best Iter
BO _{w/o tree}	78.86	45.39	26	31	56.29	27.96	13	34
Expert-Guided	79.71	75.55	2	32	56.81	53.33	2	2
ChemBOMAS _{d-d}	79.68	79.04	1	6	56.81	56.81	1	1
ChemBOMAS _{k-d}	79.77	75.55	2	11	56.81	56.81	1	1

The results, summarized in Table 2, demonstrate that employing a subspace tree search significantly accelerates Bayesian Optimization, yielding up to a 34-fold improvement in optimization efficiency compared to the standard BO baseline without a tree. Notably, the clustering methods derived from ChemBOMAS—both ChemBOMAS_{d-d} and ChemBOMAS_{k-d}—consistently matched or surpassed the performance of the expert-guided approach across all benchmarks. This underscores the robustness and reliability of our automated framework for variable space decomposition. Furthermore, the knowledge-driven clustering strategy achieved superior optimization performance on more benchmarks compared to its data-driven counterpart, highlighting the value of incorporating structured chemical knowledge.

4.3.3 OPTIMIZATION

As shown in Figure 2, ChemBOMAS demonstrates consistent and superior performance over all baseline methods across the four benchmark datasets in terms of optimal result, convergence rate, initialization performance, and robustness.

In terms of final performance and convergence speed, ChemBOMAS identified the highest objective values—96.15% (Suzuki), 81.26% (Arylation), 80.00% (Buchwald_{sub-1}) and 56.80% (Buchwald_{sub-2})—achieving convergence in just 3, 39, 23, and 2 iterations, respectively. This represents the fastest convergence among all methods.

Regarding initialization performance, ChemBOMAS attained the highest initial performance on the Arylation, Buchwald_{sub-1}, and Buchwald_{sub-2} datasets. Although it started 5.12% lower than LA-MCTS and BO-ICL on the Suzuki dataset, it surpassed all baselines by the third iteration and proceeded to converge, highlighting its strong optimization capability.

Two additional observations further underscore the robustness of ChemBOMAS. First, it exhibited the lowest variance across five independent optimization runs, indicating high stability. Second, its performance remained consistently effective and was largely unaffected by the sampling batch size (see Appendix I.1 for details). These findings collectively confirm the reliability of our method.

To evaluate the generality of ChemBOMAS beyond chemistry, we assessed its optimization performance on a materials science benchmark. As shown in Table 9 (see Appendix J for dataset details), ChemBOMAS maintains competitive performance, demonstrating its applicability to other scientific domains.

To further validate the practical applicability of ChemBOMAS and preclude the possibility of knowledge leakage in the LLM, we conducted wet-lab optimization for a previously unreported chemical reaction (see Appendix H for details). As shown in Figure 5, ChemBOMAS identified the optimal reaction condition with a yield of 96% after evaluating only 43 samples in 2 iterations. This result markedly outperforms the 15% yield obtained by a chemist using the traditional control variable method, demonstrating the framework’s effectiveness in real-world scenarios.

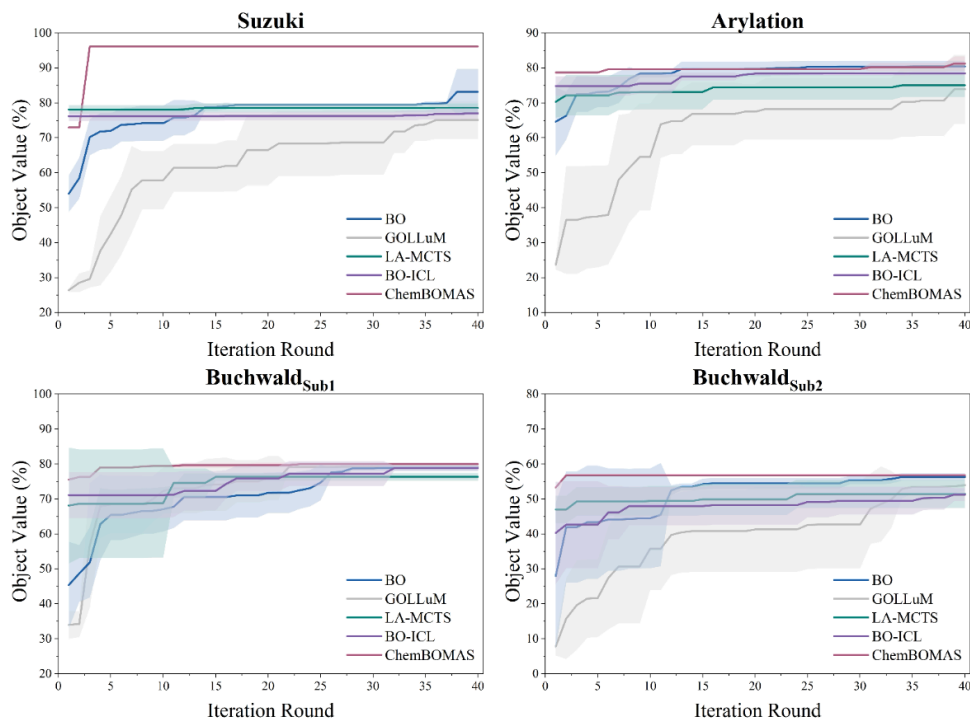


Figure 2: Optimization performance comparison between ChemBOMAS and baseline methods on the four benchmark datasets: (a) Suzuki, (b) Arylation, (c) Buchwald_{sub-1}, and (d) Buchwald_{sub-2}. ChemBOMAS exhibits accelerated convergence and achieves superior final performance with lower variance across all tasks, demonstrating its enhanced efficiency and robustness.

4.4 ABLATION STUDY

We first evaluate the impact of pre-training and fine-tuning on the regression performance of ChemBOMAS. The prediction accuracy is measured on the Suzuki, Arylation, and Buchwald datasets using Mean Squared Error (MSE), Mean Absolute Error (MAE), and the coefficient of determination (R^2).

Table 3: Impact of pre-training and fine-tuning strategies on the regression performance of ChemBOMAS across the Suzuki, Arylation, and Buchwald datasets.

Model Configuration	Suzuki			Arylation			Buchwald		
	MSE↓	MAE↓	R^2	MSE↓	MAE↓	R^2	MSE↓	MAE↓	R^2
w/o Pre & SFT	2338.02	39.43	-1.966	1797.88	32.54	-1.413	1881.96	33.73	-1.527
Pre-train Only (w/o SFT)	2407.22	40.27	-2.054	1853.70	33.24	-1.488	1795.72	32.52	-1.411
SFT Only (w/o Pre-train)	685.55	20.50	0.130	679.72	19.57	0.088	739.27	20.57	0.007
Pre-train & SFT	633.68	19.47	0.196	650.00	19.55	0.128	667.16	19.51	0.104

The results in Table 3 indicate that the combined use of pre-training and supervised fine-tuning (SFT) yields the best predictive performance across all benchmarks. Notably, SFT alone (without pre-training) achieves the second-best performance and substantially outperforms models using only pre-training or those without any training. This strongly suggests that supervised fine-tuning is the most critical component for adapting large models to chemical performance prediction tasks.

We further evaluate the contribution of each module by comparing the complete ChemBOMAS framework against three ablated versions: (i) without the data-driven module, (ii) without the knowledge-driven module, and (iii) without both modules. The results in Table 4 demonstrate

that both modules are critical to the framework’s performance. Ablating either module leads to a significant degradation in both optimization efficiency and final effectiveness.

For instance, on the Suzuki dataset, the full ChemBOMAS achieves the optimal value of 96.15% within only three iterations. In contrast, removing the data-driven module reduces the optimum to 82.65%, while disabling the knowledge-driven module limits it to 88.98%. The performance of the single-module ablations is comparable to or only marginally better than the version lacking both modules, indicating that neither strategy alone is sufficient. These results underscore that the synergy between the knowledge-driven and data-driven strategies is essential for creating a highly efficient and robust optimization framework.

Table 4: Optimization performance of the full ChemBOMAS framework compared to its ablated variants: without the data-driven module (w/o data module), without the knowledge-driven module (w/o knowledge module), and without both modules (w/o both).

Method	Suzuki				Arylation			
	Best Found (%)	Initial (%)	95% Max Iter	Best Iter	Best Found (%)	Initial (%)	95% Max Iter	Best Iter
Full ChemBOMAS	96.15	72.98	3	3	81.26	78.71	1	39
w/o data module	82.65	60.44	13	29	80.65	45.40	13	40
w/o knowledge module	88.98	65.95	21	37	79.67	45.98	10	40
w/o both	83.16	54.04	16	38	80.45	64.64	8	35

Method	Buchwald _{sub-1}				Buchwald _{sub-2}			
	Best Found (%)	Initial (%)	95% Max Iter	Best Iter	Best Found (%)	Initial (%)	95% Max Iter	Best Iter
Full ChemBOMAS	80.00	75.55	2	23	56.81	53.33	2	2
w/o data module	79.90	59.32	10	40	55.53	33.12	17	32
w/o knowledge module	79.63	54.07	9	31	56.81	12.87	9	11
w/o both	78.86	45.39	26	31	56.29	27.96	13	34

5 LIMITATIONS

While ChemBOMAS represents a significant advancement in accelerating Bayesian Optimization for chemical reactions, its performance remains constrained by several factors. Most notably, the framework is inherently dependent on the accuracy and scope of the underlying LLM and its knowledge base; inference errors in literature parsing or incomplete corpora can lead to suboptimal search-space decompositions. In addition, the absence of explicit safety and feasibility constraints raises the risk of recommending theoretically optimal yet practically hazardous or infeasible conditions, underscoring the need for expert oversight or integration of safety-aware modules in future implementations.

6 CONCLUSION

ChemBOMAS presents an LLM-enhanced multi-agent framework designed to accelerate Bayesian Optimization in the context of chemical reactions. Through a synergistic combination of knowledge-driven search space decomposition and data-driven pseudo-data generation, this approach seeks to mitigate common challenges like data scarcity and complex reaction mechanisms. Results from benchmark evaluations, along with encouraging outcomes from wet-lab validation on a demanding, previously unreported industrial reaction—where ChemBOMAS showed improved performance compared to domain expert methods—suggest its potential for practical application. ChemBOMAS offers a promising direction for facilitating chemical discovery and enhancing the optimization of complex chemical processes.

REFERENCES

- Yunchao Xie, Kianoosh Sattari, Chi Zhang, and Jian Lin. Toward autonomous laboratories: Convergence of artificial intelligence and experimental automation. *Progress in Materials Science*, 132: 101043, 2023.
- Gary Tom, Stefan P Schmid, Sterling G Baird, Yang Cao, Kourosh Darvish, Han Hao, Stanley Lo, Sergio Pablo-García, Ella M Rajaonson, Marta Skreta, et al. Self-driving laboratories for chemistry and materials science. *Chemical Reviews*, 124(16):9633–9732, 2024.

- Martin Seifrid, Robert Pollice, Andres Aguilar-Granda, Zamyla Morgan Chan, Kazuhiro Hotta, Cher Tian Ser, Jenya Vestfrid, Tony C Wu, and Alan Aspuru-Guzik. Autonomous chemical experiments: Challenges and perspectives on establishing a self-driving lab. *Accounts of Chemical Research*, 55(17):2454–2466, 2022.
- Yangguan Chen, Longhan Zhang, Zhehong Ai, Yifan Long, Ji Qi, Pengxiao Bao, and Jing Jiang. Robot-assisted optimized array design for accurate multi-component gas quantification. *Chemical Engineering Journal*, 496:154225, 2024.
- Zhehong Ai, Longhan Zhang, Yangguan Chen, Yu Meng, Yifan Long, Julin Xiao, Yao Yang, Wei Guo, Yueming Wang, and Jing Jiang. Customizable colorimetric sensor array via a high-throughput robot for mitigation of humidity interference in gas sensing. *ACS sensors*, 9(8):4143–4153, 2024a.
- Jeff Guo, Bojana Ranković, and Philippe Schwaller. Bayesian optimization for chemical reactions. *Chimia*, 77(1/2):31–38, 2023.
- Milad Abolhasani and Eugenia Kumacheva. The rise of self-driving labs in chemical and materials sciences. *Nature Synthesis*, 2(6):483–492, 2023.
- Yangguan Chen, Longhan Zhang, Zhehong Ai, Yifan Long, Temesgen Muruts Weldengus, Xubin Zheng, Di Wang, Haowen Wang, Yiteng Zhai, Yuqing Huang, et al. Robot-accelerated development of a colorimetric co₂ sensing array with wide ranges and high sensitivity via multi-target bayesian optimizations. *Sensors and Actuators B: Chemical*, 390:133942, 2023.
- Zhehong Ai, Longhan Zhang, Yangguan Chen, Yifan Long, Boyuan Li, Qingyu Dong, Yueming Wang, and Jing Jiang. On-demand optimization of colorimetric gas sensors using a knowledge-aware algorithm-driven robotic experimental platform. *ACS sensors*, 9(2):745–752, 2024b.
- Benjamin J Shields, Jason Stevens, Jun Li, Marvin Parasram, Farhan Damani, Jesus I Martinez Alvarado, Jacob M Janey, Ryan P Adams, and Abigail G Doyle. Bayesian reaction optimization as a tool for chemical synthesis. *Nature*, 590(7844):89–96, 2021a.
- Bobak Shahriari, Kevin Swersky, Ziyu Wang, Ryan P Adams, and Nando De Freitas. Taking the human out of the loop: A review of bayesian optimization. *Proceedings of the IEEE*, 104(1):148–175, 2015.
- Xilu Wang, Yaochu Jin, Sebastian Schmitt, and Markus Olhofer. Recent advances in bayesian optimization. *ACM Computing Surveys*, 55(13s):1–36, 2023.
- Linnan Wang, Rodrigo Fonseca, and Yuandong Tian. Learning search space partition for black-box optimization using monte carlo tree search. *Advances in Neural Information Processing Systems*, 33:19511–19522, 2020a.
- Austin Tripp, Erik Daxberger, and José Miguel Hernández-Lobato. Sample-efficient optimization in the latent space of deep generative models via weighted retraining. *Advances in Neural Information Processing Systems*, 33:11259–11272, 2020.
- Riccardo Moriconi, Marc Peter Deisenroth, and KS Sesh Kumar. High-dimensional bayesian optimization using low-dimensional feature spaces. *Machine Learning*, 109:1925–1943, 2020.
- Amin Nayebi, Alexander Munteanu, and Matthias Poloczek. A framework for bayesian optimization in embedded subspaces. In *International Conference on Machine Learning*, pages 4752–4761. PMLR, 2019.
- Yuxuan Yin, Yu Wang, and Peng Li. High-dimensional bayesian optimization via semi-supervised learning with optimized unlabeled data sampling. *arXiv preprint arXiv:2305.02614*, 2023.
- Kevin Swersky, Jasper Snoek, and Ryan P Adams. Multi-task bayesian optimization. *Advances in neural information processing systems*, 26, 2013.

- Artur L. F. Souza, Luigi Nardi, Leonardo B. Oliveira, Kunle Olukotun, Marius Lindauer, and Frank Hutter. Bayesian optimization with a prior for the optimum. In Nuria Oliver, Fernando Pérez-Cruz, Stefan Kramer, Jesse Read, and José Antonio Lozano, editors, *Machine Learning and Knowledge Discovery in Databases. Research Track - European Conference, ECML PKDD 2021, Bilbao, Spain, September 13-17, 2021, Proceedings, Part III*, volume 12977 of *Lecture Notes in Computer Science*, pages 265–296. Springer, 2021. doi: 10.1007/978-3-030-86523-8_17. URL https://doi.org/10.1007/978-3-030-86523-8_17.
- Tennison Liu, Nicolás Astorga, Nabeel Seedat, and Mihaela van der Schaar. Large language models to enhance bayesian optimization. *arXiv preprint arXiv:2402.03921*, 2024.
- Tung Nguyen, Qiuyi Zhang, Bangding Yang, Chansoo Lee, Jörg Bornschein, Yingjie Miao, Sagi Perel, Yutian Chen, and Xingyou Song. Predicting from strings: Language model embeddings for bayesian optimization. *CoRR*, abs/2410.10190, 2024. doi: 10.48550/ARXIV.2410.10190. URL <https://doi.org/10.48550/arXiv.2410.10190>.
- Mayk Caldas Ramos, Shane S Michtavy, Marc D Porosoff, and Andrew D White. Bayesian optimization of catalysts with in-context learning. *arXiv preprint arXiv:2304.05341*, 2023a.
- David Eric Austin, Anton Korikov, Armin Toroghi, and Scott Sanner. Bayesian optimization with llm-based acquisition functions for natural language preference elicitation. In Tommaso Di Noia, Pasquale Lops, Thorsten Joachims, Katrien Verbert, Pablo Castells, Zhenhua Dong, and Ben London, editors, *Proceedings of the 18th ACM Conference on Recommender Systems, RecSys 2024, Bari, Italy, October 14-18, 2024*, pages 74–83. ACM, 2024. doi: 10.1145/3640457.3688142. URL <https://doi.org/10.1145/3640457.3688142>.
- Shukuan Wang, Ke Xue, Lei Song, Xiaobin Huang, and Chao Qian. Monte carlo tree search based space transfer for black-box optimization. *arXiv preprint arXiv:2412.07186*, 2024.
- Linnan Wang, Saining Xie, Teng Li, Rodrigo Fonseca, and Yuandong Tian. Sample-efficient neural architecture search by learning action space. *arXiv preprint arXiv:1906.06832*, 2019.
- Daniel Reker, Emily A Hoyt, Gonçalo JL Bernardes, and Tiago Rodrigues. Adaptive optimization of chemical reactions with minimal experimental information. *Cell Reports Physical Science*, 1(11), 2020.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, et al. The llama 3 herd of models, 2024. URL <https://arxiv.org/abs/2407.21783>.
- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net, 2022. URL <https://openreview.net/forum?id=nZeVKeeFYf9>.
- Damith Perera, Joseph W Tucker, Shalini Brahmabhatt, Christopher J Helal, Ashley Chong, William Farrell, Paul Richardson, and Neal W Sach. A platform for automated nanomole-scale reaction screening and micromole-scale synthesis in flow. *Science*, 359(6374):429–434, 2018.
- Benjamin J Shields, Jason Stevens, Jun Li, Marvin Parasram, Farhan Damani, Jesus I Martinez Alvarado, Jacob M Janey, Ryan P Adams, and Abigail G Doyle. Bayesian reaction optimization as a tool for chemical synthesis. *Nature*, 590(7844):89–96, 2021b.
- Derek T Ahneman, Jesús G Estrada, Shishi Lin, Spencer D Dreher, and Abigail G Doyle. Predicting reaction performance in c–n cross-coupling using machine learning. *Science*, 360(6385):186–190, 2018.
- OpenAI et al. Gpt-4o system card, 2024. URL <https://arxiv.org/abs/2410.21276>.
- OpenAI. Introducing gpt-5, 2025. URL <https://openai.com/gpt-5>.
- Connor J. Taylor, Alexander Pomberger, Kobi C. Felton, Rachel Grainger, Magda Barecka, Thomas W. Chamberlain, Richard A. Bourne, Christopher N. Johnson, and Alexei A. Lapkin. A brief introduction to chemical reaction optimization. *Chemical Reviews*, 123(6):3089–3126, 2023.

Mayk Caldas Ramos, Shane S Michtavy, Marc D Porosoff, and Andrew D White. Bayesian optimization of catalysts with in-context learning. *arXiv preprint arXiv:2304.05341*, 2023b.

Bojana Ranković and Philippe Schwallier. Gollum: Gaussian process optimized llms—reframing llm finetuning through bayesian optimization. In *ICLR 2025 Workshop on World Models: Understanding, Modelling and Scaling*, 2025.

Linnan Wang, Rodrigo Fonseca, and Yuandong Tian. Learning search space partition for black-box optimization using monte carlo tree search. *Advances in Neural Information Processing Systems*, 33:19511–19522, 2020b.

A ILLUSTRATIVE EXAMPLE OF LLM ASSISTED CONSTRUCTION OF A REACTION OPTIMIZATION TREE

In the following illustrative example, we demonstrate how an LLM can assist in constructing an optimization tree for the reaction $A + B \rightarrow C$ in two steps, enabling efficient optimization of reaction conditions (e.g., catalyst, ligand, solvent, base), given that reactants A and B are fixed. In the first step, an LLM is used to infer the possible reaction type for $A + B \rightarrow C$. Based on the inferred reaction type and specific optimization objectives (e.g., improving yield or selectivity), relevant scientific literature is retrieved. Literature acquisition can be done through manual downloads or by using publisher-provided APIs (noting that not all APIs are openly accessible). The collected literature is then used to construct a vector database to support the subsequent retrieval process. Using the information from the literature in the vector database, the LLM is queried via analyzing literature to determine the relative importance of different reaction conditions (variables) on the reaction objects, generating a ranked list. For instance, the LLM might determine the order of influence as: Catalyst > Ligand > Solvent > Base. Further queries to the LLM identify the key physicochemical properties within each category that significantly influence the chemical reaction performance. For example, within the ligand category, the LLM may highlight "steric and electronic effects" as crucial physicochemical properties. Subsequently, detailed information regarding the key physicochemical properties of each ligand candidate is retrieved from online databases, after which the LLM clusters these ligand candidates into subsets based on similarities in "steric and electronic effects".

In the second step, the optimization tree is constructed based on the variable importance ranking and clustering results. The first level of the tree corresponds to the most important variable—the catalyst. At the first level, several child nodes can be established, representing different subsets of catalyst candidates clustered by property similarity. The second level of the tree corresponds to the next most important variable—the ligand. Under each catalyst subset node at the first level, additional child nodes branch out, representing various subsets of ligand candidates categorized by their physicochemical properties. This process continues iteratively, layer by layer, incorporating additional variables (e.g., solvent, base) until the complete optimization tree is constructed.

B UPDATE ON CHEMBOMAS DURING OPTIMIZATION

After receiving the observation feedback on each round of the experiment, ChemBOMAS would update. First, the data module would be retrained with the prior and newly acquired data points, and then infer the unsampled data points to generate pseudo-labels. Second, the optimization tree would recount the visit number and value of each node to refine the identified hot regions. Third, with the updated observations, pseudo-labels, and refined hot regions, the BO module would recommend next-round reaction conditions, targeting potentially higher object values.

C BENCHMARK DETAIL

This section provides further details on the benchmark datasets used for evaluating ChemBOMAS.

C.1 DATASET FOR LLM PRE-TRAIN

The Pistachio dataset employed during the pre-training phase is a large-scale reaction information repository. Its core data was systematically extracted from the full texts of US patents (USPTO) and European patents (EPO) through automated text mining techniques. To enhance data diversity and accuracy, the dataset integrates information from multiple sources, including: - Structured data parsed from ChemDraw (CDX) files embedded directly within patent documents - Records sourced from specialized chemical databases such as Reaxys - Exported data from select electronic laboratory notebooks (ELNs) The dataset contains a total of 19.17 million chemical reactions. In this project, we primarily utilize the reaction SMILES strings for model pre-training.

C.2 DATASET FOR LLM FINE-TUNE

To conduct a rigorous and unbiased evaluation of model performance, we selected a series of publicly available benchmark datasets widely used in the field of chemical reaction optimization. The core strength of these datasets lies in their completeness: all were generated via high-throughput automated experimental platforms and encompass experimental results for every variable combination within a clearly defined chemical space (full factorial design). This exhaustive coverage effectively eliminates sampling bias, enabling deterministic quantitative evaluation of algorithmic recommendation performance against known experimental ground truth.

Specifically, we employed three recognized benchmark datasets: Suzuki, Arylation, and Buchwald reactions. During fine-tuning, we randomly sampled 1% of data from each dataset as training samples to adjust the pre-trained model.

Suzuki originates from the automated nanomolar-scale flow screening study reported by Perera et al. in 2018. The chemical space of the experiments comprised a full factorial combination of 4 halogenated quinolines, 3 boronic acid derivatives, 11 phosphine ligands, 7 bases, and 4 solvents. All reactions were conducted under uniform conditions (100 °C, 1-minute residence time, 9:1 organic/aqueous phase). Reaction yields were detected via dual UPLC-MS online detection and uniformly calibrated. The data is comprehensive and highly consistent, making it one of the widely adopted validation standards in the field.

Arylation was reported by Shields et al. in 2021 for Bayesian optimization studies. Its chemical space was generated via a full factorial design comprising 12 phosphine ligands, 4 bases, 4 solvents, 3 temperature gradients, and 3 concentration gradients. All experiments were conducted at high throughput in 96-well plates, with yields precisely quantified via UHPLC-MS coupled with internal standard methods. This dataset features no duplicates or missing data, exhibits uniform variable distribution, and has been validated by 50 practicing chemists, establishing it as a critical benchmark for optimizing C-H functionalization reactions.

Buchwald was published by Ahneman et al. in 2018, this dataset aims to predict yields of C-N coupling reactions via machine learning. Experiments were conducted in nanomolar-scale high-throughput format using 1536-well plates, systematically examining all combinations of 15 aryl halides, 4 ligands, 3 bases, and 23 isoxazole additives. All reactions proceeded under standard conditions (60 °C, DMSO, 16 hours), with yields quantified by LC-MS. This dataset is complete with no missing values, serving as an authoritative open-access resource for studying additive effects and modeling complex reaction systems.

C.3 DATASET FOR BAYESIAN OPTIMIZATION

Table 5: Descriptive statistics of the four reaction datasets. The table summarizes key statistical measures for the reaction yields, including measures of central tendency, dispersion, and distribution shape.

Statistic	Suzuki	Arylation	Buchwald _{sub-1}	Buchwald _{sub-2}
Total data points (N)	5030	3678	629	765
Maximum Yield (%)	96.15	84.65	80.91	56.81
Minimum Yield (%)	0.00	0.00	0.00	0.00
Mean (%)	33.04	29.05	42.24	18.71
Median (%)	26.86	25.53	42.21	11.34
Standard Deviation (%)	22.47	23.79	22.86	18.98
25% Quantile (%)	15.26	6.87	23.14	0.72
75% Quantile (%)	51.27	47.14	63.01	38.77

For Bayesian optimization tasks, we employed four benchmark datasets: Suzuki, Arylation, Buchwald_{sub-1}, and Buchwald_{sub-2}. The latter two originate from partitions of the aforementioned Buchwald-Hartwig dataset. To ensure consistency of target products within the optimization space, the original dataset was first divided into five independent subsets based on product molecular structures. We observed distinct high-yield and low-yield patterns in the reaction yields of these subsets.

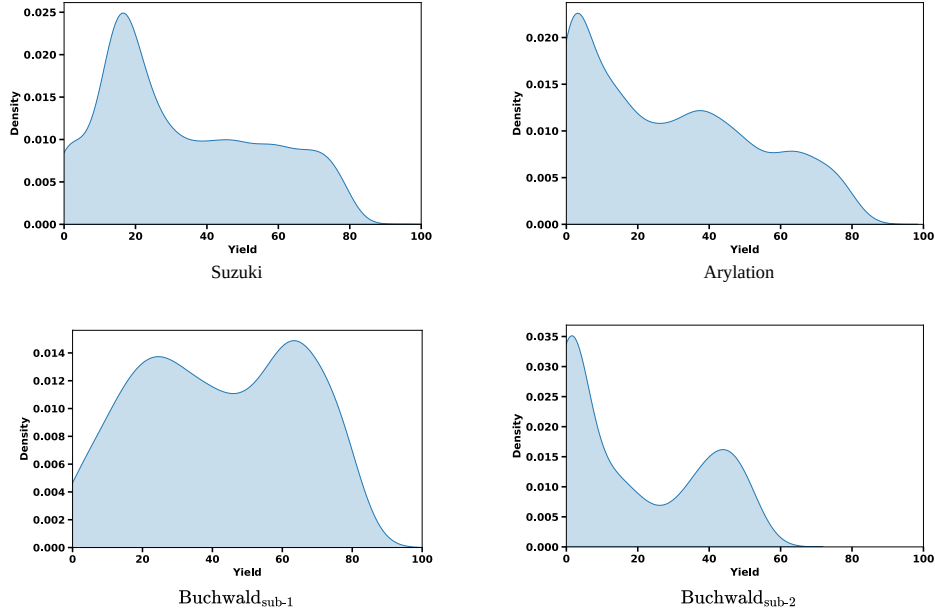


Figure 3: KDE plots illustrating the yield distributions for the four benchmark datasets.

To ensure comprehensive evaluation, we selected one representative subset from each category, naming them $\text{Buchwald}_{\text{sub-1}}$ and $\text{Buchwald}_{\text{sub-2}}$, respectively. Table 5 summarizes key descriptive statistics for these four datasets, while Figure 3 visually depicts their respective yield distributions via kernel density estimation (KDE) plots. These datasets exhibit distinct statistical characteristics, with average yields ranging from 18.71% to 42.24% and diverse distribution shapes. Collectively, they form a challenging optimization problem that effectively tests algorithm performance across varying data environments.

C.4 COMPARATIVE ALGORITHMS

To evaluate the efficacy of our proposed method, we benchmark it against four algorithms that represent diverse approaches to black-box optimization.

Traditional Bayesian Optimization (**BO**) is a sequential optimization method that utilizes a surrogate model to approximate the objective function and an acquisition function to guide subsequent sampling. The baseline implementation in this work employs a Gaussian Process (GP) with a Matérn kernel as the surrogate, a qLogEI acquisition function, and one-hot encoding for inputs. BO did not adopt more complex encoding schemes because prior research indicated their benefits were negligible Taylor et al. (2023); Shields et al. (2021a).

Bayesian Optimization with In-Context Learning (**BO-ICL**) is a method that integrates a frozen Large Language Model (LLM) with BO, as proposed by Ramos et al. (2023b). It leverages the LLM’s in-context learning ability to form a surrogate model by translating experimental parameters into textual prompts, thereby predicting outcomes and uncertainty without model fine-tuning.

Gaussian Process Optimized LLMs (**GOLLM**) is a method that achieves a deeper fusion of LLMs and Bayesian optimization, introduced by Ranković and Schwaller (2025). It utilizes the LLM’s embedding space as a deep kernel for a GP and jointly optimizes the embedding and GP hyperparameters to learn a task-specific representation space.

Latent Action Monte Carlo Tree Search (**LA-MCTS**) is a meta-algorithm for high-dimensional optimization, developed by Wang et al. (2020b). It employs Monte Carlo Tree Search to dynamically partition the search space into high- and low-performance regions and subsequently deploys a local optimizer within promising subregions.

D DUAL-STRATEGY REFINEMENT FOR ENHANCED OPTIMIZATION

To mitigate the detrimental influence of noise and redundancy inherent in generated pseudo-data, we introduced a dual-pronged refinement strategy. This approach was designed to dynamically curate the pseudo-dataset, ensuring its quality and diversity throughout the optimization process. The strategy combined a local, similarity-based removal mechanism with a global, performance-driven pruning policy. This ensured that the pseudo-dataset remained a reliable and informative asset for guiding the optimization, particularly in complex search spaces.

Data Similarity (Local Removal): We utilized the final token embedding, $\mathbf{e}(\mathbf{x}) = \text{LLM}_{\theta_{\text{LLM}}, \phi_{\text{LoRA}}}^{[T]}(\mathbf{x})$, to calculate cosine similarity. Upon acquiring a new real data point $(\mathbf{x}_{\text{new}}, y_{\text{new}})$, the pseudo-dataset was updated by removing points that were too similar:

$$\mathcal{D}_{\text{pseudo}} \leftarrow \mathcal{D}_{\text{pseudo}} \setminus \left\{ (\mathbf{x}_j, \hat{y}_j) \in \mathcal{D}_{\text{pseudo}} \mid \frac{\mathbf{e}(\mathbf{x}_j) \cdot \mathbf{e}(\mathbf{x}_{\text{new}})}{\|\mathbf{e}(\mathbf{x}_j)\| \|\mathbf{e}(\mathbf{x}_{\text{new}})\|} > \tau \right\} \quad (2)$$

where τ was a predefined similarity threshold.

Observed Performance (Global Removal): As the optimization progresses, the model should be encouraged to explore more broadly. Therefore, based on the predicted performance values \hat{y} of the pseudo-points, we randomly discarded a proportion of pseudo-data, starting from those with high predicted performance downwards. The probability of discarding a pseudo-point $(\mathbf{x}_j, \hat{y}_j)$ was a monotonically increasing function of its predicted performance.

These generated pseudo-points could also provide further support for the construction of the knowledge-guided optimization tree in Stage 1. By adjusting the LLM’s temperature parameter during generation, we could produce a set of candidate tree structures, $\{\mathcal{T}_1, \mathcal{T}_2, \dots, \mathcal{T}_K\}$. Using the pseudo-points, we quantitatively evaluated these candidates. Let $\mathcal{N}(\mathcal{T}_k)$ be the set of leaf nodes of tree \mathcal{T}_k , and let $\mathcal{D}_{\text{pseudo}}^{(j)}$ be the subset of pseudo-points belonging to node $j \in \mathcal{N}(\mathcal{T}_k)$. The tree structure that minimized the weighted average of intra-node variances was selected as the optimal one:

$$\mathcal{T}^* = \arg \min_{\mathcal{T}_k} \sum_{j \in \mathcal{N}(\mathcal{T}_k)} \frac{|\mathcal{D}_{\text{pseudo}}^{(j)}|}{|\mathcal{D}_{\text{pseudo}}|} \text{Var}(\{\hat{y} \mid (\mathbf{x}, \hat{y}) \in \mathcal{D}_{\text{pseudo}}^{(j)}\}) \quad (3)$$

This ensured the selection of a tree that best partitions the search space into regions of homogeneous performance, guiding the subsequent optimization more effectively.

E QUALITATIVE COMPARISON OF OPTIMIZATION TRAJECTORIES

To qualitatively assess how well the automated clustering strategies of ChemBOMAS emulate expert-level reasoning, we visualized the optimization progress. Figure 4 provides a comparative heatmap of the "Best Found" objective value over 40 iterations for three different search tree configurations: one guided by human experts, one by our knowledge-driven module (ChemBOMAS_{k-d}), and one by our data-driven module (ChemBOMAS_{d-d}).

The visual evidence strongly suggests that both automated ChemBOMAS strategies produce optimization trajectories that are remarkably consistent with the expert-guided approach. The color progression—from blue (lower values) to red (higher values)—is highly similar across all three methods. This indicates that the subspaces identified as promising by the LLM-driven modules align well with those selected by human domain experts. The ability of both the knowledge-driven and data-driven variants to rapidly progress towards high-yield regions in a manner analogous to the expert baseline underscores the effectiveness of our framework in automatically structuring the search space in a chemically meaningful way. This qualitative alignment provides further confidence in the robustness and practical utility of ChemBOMAS for real-world chemical optimization tasks.

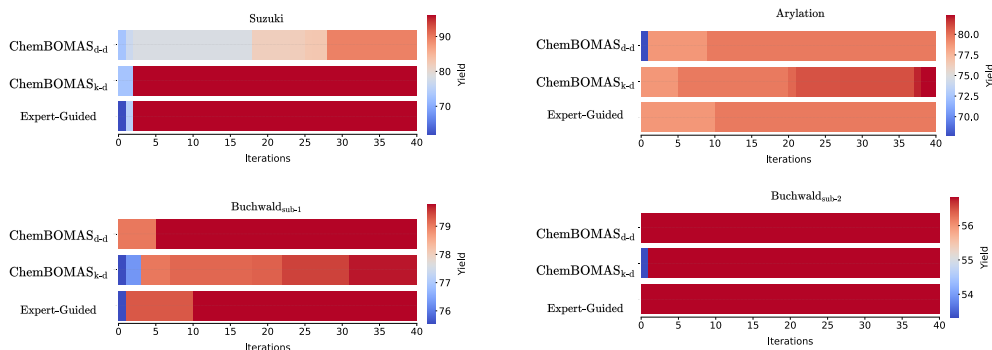


Figure 4: Heatmap of the best-found objective value over 40 iterations on the Suzuki dataset for three different tree-building strategies. Each colored block represents the highest value discovered up to that iteration, with the color scale progressing from blue (low) to red (high). The visual similarity in the optimization trajectories demonstrates that both the knowledge-driven (ChemBOMAS_{k-d}) and data-driven (ChemBOMAS_{d-d}) methods closely mirror the performance progression of the expert-guided approach.

F COMPLETE ALGORITHM PROCESS

To provide a comprehensive and formal description of the ChemBOMAS framework, we present its complete algorithmic process in Algorithm F. This pseudocode encapsulates the synergistic, two-stage optimization strategy detailed in Section 3.

The algorithm begins by initializing a hierarchical search tree using the LLM-guided knowledge-driven approach as shown in Section 3.4. The main loop then iterates through the coarse-grained optimization phase, where a UCB policy navigates the tree to select a promising subspace. Within this selected subspace, the algorithm transitions to the fine-grained, data-driven optimization phase. Here, a standard Bayesian Optimization procedure is executed, but it is significantly accelerated by an informative prior constructed from both real experimental data and pseudo-data generated by the fine-tuned LLM regressor (Section 3.3).

After each experimental evaluation, the results are backpropagated to update the value estimates of the nodes in the search tree, refining the knowledge-driven search for subsequent iterations. This process continues until the predefined budget of evaluations is exhausted, ultimately returning the best-performing experimental configuration found.

G DETAILS OF PROMPTS

As outlined in the main text, our methodology leverages LLMs to support several critical tasks in reaction optimization, such as analyzing literature, assessing parameter significance, and understanding physicochemical properties to inform the construction of a hierarchical optimization tree. This appendix section presents a detailed overview of the specific prompts designed to guide the LLM in executing these crucial Tasks.

G.1 PROMPT OF DATA MODULE

As detailed in the Section 3.3, our pre-training phase employs a conditional prediction task. Given the reactants and products, the model’s objective is to predict the corresponding reaction conditions. This process utilizes a Causal Language Modeling (CLM) loss, where the model learns to predict the next token in the sequence of reaction conditions.

To provide concrete examples of the input format for this task, this appendix section presents a selection of prompts utilized during the pre-training phase. These prompts typically consist of the reactants, products, and the target reaction condition sequence that the model is trained to predict.

Algorithm 1 The Complete Algorithm Process of ChemBOMAS

Input: Search space \mathcal{X} , black-box objective function $h(\cdot)$, coarse iterations N_{coarse} , fine iterations per evaluation N_{fine} , exploration constant C_p , fine-tuned LLM regressor $f_{\theta_{MLP}}(\text{LLM}_{\theta_{LLM}, \phi_{LoRA}}(\cdot))$.

Initialize:

Construct hierarchical search tree via LLM-guided space partitioning (see Section 3.4).

Construct hierarchical search tree \mathcal{T} by partitioning \mathcal{X} via LLM-driven analysis.

Initialize value estimate $Q_v \leftarrow 0$, visit count $n_v \leftarrow 0$ for all nodes $v \in \mathcal{T}$.

Initialize global set of real experimental data $\mathcal{D}_{real} \leftarrow \emptyset$.

procedure MAIN LOOP

for $i = 1$ to N_{coarse} **do**

$v_{current} \leftarrow \text{root}(\mathcal{T})$

$\text{path} \leftarrow [v_{current}]$

// Stage 1: Knowledge-driven Strategy

while $v_{current}$ is not a leaf node **do**

$v_{current} \leftarrow \arg \max_{v_k \in \text{children}(v_{current})} \left(\frac{Q_{v_k}}{n_{v_k}} + C_p \sqrt{\frac{\ln n_{v_{current}}}{n_{v_k}}} \right)$

 Append $v_{current}$ to path.

end while

 Let \mathcal{S}_j be the promising subspace corresponding to the leaf node $v_{current}$.

// Stage 2: Data-driven Strategy

$(y_{new}, \mathbf{x}_{new}) \leftarrow \text{BO}(\mathcal{S}_j, N_{fine}, \mathcal{D}_{real}, f_{\theta_{MLP}}(\text{LLM}_{\theta_{LLM}, \phi_{LoRA}}(\cdot)))$

$\mathcal{D}_{real} \leftarrow \mathcal{D}_{real} \cup \{(\mathbf{x}_{new}, y_{new})\}$.

for v in path **do**

$n_v \leftarrow n_v + 1$

$Q_v \leftarrow Q_v + y_{new}$

end for

end for

end procedure

▷ Backpropagation

function $\text{BO}(\mathcal{S}_j, N_{fine}, \mathcal{D}_{real}, \text{LLM_regressor})$

 //Initialize surrogate model with LLM-generated pseudo-data.

 Generate pseudo-dataset $\mathcal{D}_{pseudo} = \{(\mathbf{x}_k, \hat{y}_k)\}_{k=1}^M$ for $\mathbf{x}_k \in \mathcal{S}_j$ using LLM_regressor.

 Let $\mathcal{D}_{real}^{(j)} = \{(\mathbf{x}, y) \in \mathcal{D}_{real} \mid \mathbf{x} \in \mathcal{S}_j\}$.

 ▷ Fit Gaussian Process (GP) on combined data to serve as an informative prior.

 Initialize GP surrogate model \mathcal{M} on $\mathcal{D}_{pseudo} \cup \mathcal{D}_{real}^{(j)}$.

for $k = 1$ to N_{fine} **do**

 ▷ Select next point by maximizing the acquisition function $\alpha(\cdot)$.

$\mathbf{x}_{next} \leftarrow \arg \max_{\mathbf{x} \in \mathcal{S}_j} \alpha(\mathbf{x} | \mathcal{M})$

$y_{next} \leftarrow h(\mathbf{x}_{next})$

 ▷ Perform real experiment to get objective value.

$\mathcal{D}_{real}^{(j)} \leftarrow \mathcal{D}_{real}^{(j)} \cup \{(\mathbf{x}_{next}, y_{next})\}$

// Apply refinement strategy

 Update \mathcal{D}_{pseudo} by removing points based on similarity and performance rules.

 Update GP surrogate model \mathcal{M} with new data $\{(\mathbf{x}_{next}, y_{next})\}$ and pruned \mathcal{D}_{pseudo} .

end for

return $(y_{next}, \mathbf{x}_{next})$

▷ Return the result of the last experiment.

end function

Output: The configuration \mathbf{x}^* with the highest observed objective value $h(\mathbf{x}^*)$ from \mathcal{D}_{real} .

Furthermore, in line with the methodology described in the main text, these input sequences are augmented with functional group annotations (generated via RDKit) to enhance the model’s chemical awareness; the augmentation of the prompt is also reflected in the examples provided below.

Prompt of Condition Prediction Pretraining: For the condition prediction pre-training, the input prompts are structured to provide the model with comprehensive reaction information. Typically, a prompt is formatted as: [Reactants_SMILES]; [Products_SMILES]; [Reaction Type];[Target_Reaction_Conditions]. Prior to constructing these prompts, the SMILES strings for both reactants and products are canonicalized using RDKit. This normalization step ensures a standardized and consistent representation of molecular structures, which is vital for robust model training. The model then processes this complete sequence, aiming to predict the [Target_Reaction_Conditions] segment token by token, guided by the Causal Language Modeling objective and conditioned on the preceding reaction type, reactants, and products.

To further clarify the input structure for this prediction task, the following examples demonstrate the format used:

Condition Prediction Pre-training Prompts

"reaction": "Here is a chemical reaction.
Reactants are: C1=CC=CC=2C3=CC=CC=C3N(C12)CC#C,BrC#CCCCCO.
Product is: C1=CC=CC=2C3=CC=CC=C3N(C12)CC#CC#CCCCCO.
Reaction type is Cadiot-Chodkiewicz coupling.",
"condition": "The reaction conditions of this reaction are:
Solvent: O,CN(C=O)C,CN(C=O)C. Catalyst: Cl[Cu]. Atmosphere: N#N. Additive: C(C)N,[Na]Cl,Cl.NO.", "reaction_type": "Cadiot-Chodkiewicz coupling",

Prompt of Yield Prediction Fine-tuning: To fine-tune LLM for precise prediction of chemical reaction yields, we combine key chemical information—including reaction type, products, reactants, and reaction conditions—into structured prompts. This approach guides the model to learn the complex relationships between these variables and reaction outcomes, enabling it to output a specific numerical prediction.

Below is an example prompt for yield prediction fine-tuning from the reactants in the Suzuki coupling dataset.

An example prompt for yield prediction fine-tuning

Here is a chemical reaction:
Reactants are: CCc1cccc(CC)c1.Clc1ccc2ncccc2c1,Cc1ccc2c(cnn2C2CCCCO2)c1B(O)O.
Product is: Cc1ccc2c(cnn2C2CCCCO2)c1-c1ccc2ncccc2c1.
Reaction type is Suzuki Miyaura.
The reaction conditions of this reaction are:
Solvent: CC#N.O
Ligand: CC(C)(C)P(C(C)(C)C)C(C)(C)C
Base: [Na+].[OH-]
What is the yield of this reaction?

G.2 PROMPTS OF KNOWLEDGE MODULE

The Knowledge Module, as described in Section 3.4, employs the LLM to systematically analyze chemical literature and physicochemical data. This involves ranking the impact of various reaction parameters and classifying components based on their physicochemical properties.

Variable Candidates Clustering Prompt: The prompt guides the LLM to identify key physicochemical properties of each variable and cluster variable candidates based on their similarity in the physicochemical properties. Below is an example of the prompt for variable candidates classification.

Prompt for Variable Candidates Classification Based on Physicochemical Data

Objective:

Classify the provided list of candidate chemical substances into NO MORE THAN THREE groups according to the [Specified_physicochemical_Properties], or place them all in ONE class if justified.. Your primary method for classification must be the utilization of quantitative data that would typically be found in a comprehensive physicochemical property database.

Crucial Instructions:

Prioritize Quantitative Data: For each substance and property, you should first attempt to classify it based on specific, measurable, quantitative values (e.g., pKa for basicity/acidity, dielectric constant for polarity, boiling point for volatility, specific functional group counts).

Minimize General Knowledge/Intuition: Avoid relying on your general, unquantified chemical knowledge or intuition. If a quantitative value from the "database" directly supports a classification, state that. If a direct value isn't typically used for a category but strong structural indicators (which could be quantified, e.g., number of H-bond donors) point to it, explain this as an inference based on data-like principles.

Adhere to Provided Categories: Classify substances strictly into the categories provided for each property. If a substance does not clearly fit or straddles categories based on (assumed) data, note this ambiguity.

Candidate Substances to Classify:

[TYPE] : [CANDIDATE_SUBSTANCES_LIST]

Provided Literature:

[LITERATURE_1]

[LITERATURE_2]

...

Available Tools:

[PubMedToolkit], [PubChemToolkit], [GoogleSearchToolkit]

H WET EXPERIMENTS

To further validate the efficacy and applicability of our proposed method, an algorithm-driven **wet laboratory experiment** was conducted. Guided by ChemBOMAS, this study aimed to maximize product yield in a challenging chemical reaction optimization—the palladium-catalyzed cross-coupling of boronic esters with aryl chlorides. This demanding optimization task, originating from a pharmaceutical enterprise, was governed by four stringent practical constraints: (1) a **previously-unreported** chemical reaction, resulting in the complete absence of reference data; (2) a six-dimensional process parameter space, reportedly **exceeding seventy times** the scale of those in comparable published studies, thus posing a considerable exploration challenge; (3) a cost-saving imperative requiring a **tenfold reduction** in catalyst loading relative to conventional levels, substantially hampering product formation; and (4) a restriction on approximately **60 experimental runs** to curtail labor intensity. Detailed experiment settings could be found in the Supplementary Material.

H.1 EXPERIMENT RESULTS

As shown in Figure 5, during the wet experiment task, ChemBOMAS successfully identified the optimal reaction condition with a yield of 96%, markedly outperforming the 15% yield achieved by a chemist employing the traditional control variable method. Additionally, three noteworthy phenomena emerged. First, in the initial round, ChemBOMAS had attained a maximum product yield of 90%, surpassing the target threshold of 75%. Second, the optimal reaction condition yielding 96% was discovered in the early stage of the optimization process, specifically in the second iteration. Third, as the optimization progressed, ChemBOMAS increasingly recommended reaction conditions with yields exceeding the 75% target threshold, indicating a continuous refinement of the surrogate model. The number of high-yielding conditions ($\geq 75\%$) identified in rounds one through five was

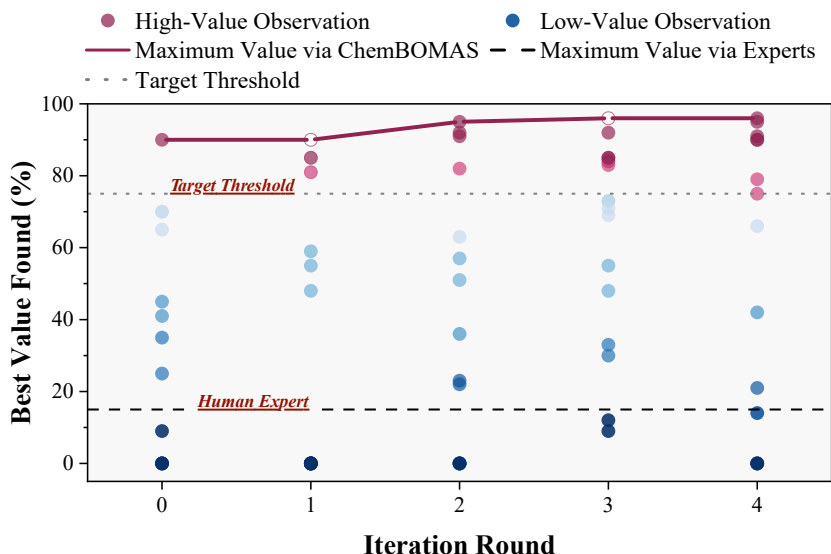


Figure 5: **Wet laboratory experiment result.** Comparison of ‘Best Value Found (%)’ over ‘Iteration Rounds’, showing individual high and low-value observations. Lines indicate maximum values achieved via ChemBOMAS, human experts, and a target threshold.

one, two, three, five, and five, respectively. The strong initialization performance, rapid convergence, and progressive model improvement collectively demonstrate the effectiveness and efficiency of ChemBOMAS in accelerating chemical reaction optimization.

H.2 WET EXPERIMENT DETAIL PROTOCOL

Additionally, to validate the robustness of ChemBOMAS’s initialization performance, the initial-round sampling was repeated ten times with the fixed experimental configurations, as detailed in Supplementary Materials. In the ten repeated initialization tests using ChemBOMAS, each run consistently identified at least two reaction conditions with yields exceeding 60%. Moreover, reaction conditions achieving yields above 80% appeared in 70% of the validation tests, totaling 11 such high-yield conditions across all trials. These results demonstrate that ChemBOMAS reliably mitigates the “cold-start” problem inherent to BO optimization.

General Procedure for Reaction Optimization For the wet experiment involving palladium-catalyzed coupling of boronic esters with aryl chlorides, first, an oven-dried 10 mL Schlenk tube fitted with a Teflon-coated magnetic stir bar was charged inside an N_2 -filled glovebox with Pd-catalyst (0.002 mmol), Phosphine ligand (0.008 mmol), and base (0.30 mmol, 1.5 equiv). Then, the tube was sealed with a septum, removed from the glovebox, and placed under a positive flow of N_2 . The Mixture of organic solvent and water (2 mL) was introduced via a syringe. Next, pinacol boronic ester 2 (Reactant 1, 0.20 mmol, 1 equiv) and Aryl chloride 1 (Reactant 1, 0.25 mmol, 1.25 equiv) were added sequentially by syringe. The tube was capped tightly, placed in a pre-heated aluminum heating block maintained at 80 °C, 100 °C, or 120 °C, and the mixture was stirred (approximately 1500 rpm) for 24 hours. After cooling to room temperature, the mixture was diluted with ethyl acetate (3 mL) and quenched with water (3 mL). Finally, GC yields were determined directly from the crude mixture against the n-dodecane standard.

ChemBOMAS Configuration Some configurations of ChemBOMAS described in the Experiment Section of the main text were adjusted for the wet experiment task. First, in the Knowledge Module, the additional process parameters (here, water usage and temperature) were divided into multiple subsets automatically by the LLM using RAG, and these subsets were grouped by the similarity of physical properties, which is the same as the category variables. For instance, temperature conditions

were categorized into three distinct subsets corresponding to low, intermediate, and high activation energy levels. Moreover, during the Bayesian Optimization (BO), considering the relatively high experimental throughput, multiple acquisition functions (here, EI and UCB) were applied to generate fourteen samples per round. Apart from the aforementioned adjustments, all other configurations within ChemBOMAS remained consistent with those used in the dry-lab experiments.

Sample in The Initial Round The initial experiment was only designed by Knowledge module due to the lack of prior data. Specifically, after the Knowledge Module partitioned the variables into subsets, a sampling function that can select variables from different subsets evenly was applied to generate fourteen diverse reaction conditions. The generated reaction conditions were then sent to the experiment operators for actual observation, which facilitated providing data to inform the experimental design in the next round.

Sample in The Iterated Round As illustrated in Section B of the Supplementary Material, after receiving the observation feedback on each round of the wet experiment, all ChemBOMAS modules would update based on the feedback from each round of the wet-lab experiments. Following the update of ChemBOMAS, the BO module would recommend fourteen reaction conditions with potentially higher yields for the subsequent round.

I ADDITIONAL RESULTS ANALYSIS

I.1 IMPACT OF BATCH SIZE

To investigate the influence of the number of samples per iteration on the performance of ChemBOMAS, we conducted an ablation study by varying the batch size. Specifically, we configured the batch sizes to be 0.05%, 0.1%, 0.2%, and 0.4% of the total dataset volume for each reaction. It is important to note that for the Buchwald_{sub-1} and Buchwald_{sub-2} datasets, which have a significantly smaller number of data points, a batch percentage of 0.05% resulted in a batch size of less than 1. Consequently, experiments for this specific setting were omitted for these two datasets.

Table 6: Bayesian Optimization Performance with Vary Batch Size Per Iteration

Dataset	Batch %	Batch Size	Best Found	Initial Value	Time	Iter. of Best
Suzuki	0.05	3	96.15	72.98	105.18	3
	0.10	5	96.15	72.98	158.03	3
	0.20	10	96.15	72.98	318.11	3
	0.40	20	96.15	72.98	646.02	3
Arylaton	0.05	2	80.67	78.71	59.24	37
	0.10	3	80.64	78.71	95.01	39
	0.20	6	81.65	78.71	240.63	39
	0.40	12	81.25	78.71	755.86	29
Buchwald _{sub-1}	0.05	-	-	-	-	-
	0.10	1	79.58	75.55	10.51	25
	0.20	2	79.39	75.55	24.43	23
	0.40	4	79.60	75.55	56.30	25
Buchwald _{sub-2}	0.05	-	-	-	-	-
	0.10	1	56.81	53.34	8.77	2
	0.20	2	56.81	53.34	20.57	2
	0.40	4	56.81	53.34	49.02	2

The experimental results, as detailed in Table I.1, revealed that while there was a marginal improvement in the best-found values with an increasing batch size, the overall impact on optimization performance was minimal. This observation underscored a critical strength of ChemBOMAS: its ability to efficiently navigate the variable space and identify optimal or near-optimal parameter combinations using a remarkably small subset of data in each iteration.

Conversely, a clear trend emerged regarding computational cost: the runtime increased substantially with larger batch sizes. Given the negligible gains in the optimal solution found, a larger batch size presented an unfavorable trade-off. Therefore, to balance computational efficiency and performance,

we selected a batch percentage of 0.1% for our main experiments, as it demonstrated the capability of ChemBOMAS to achieve excellent results with minimal data sampling per round.

I.2 IMPACT OF PRIOR DATA VOLUME

To ascertain the influence of prior data volume on the fine-tuning process, we conducted a comprehensive ablation study. We fine-tuned the pre-trained Large Language Model on five distinct proportions of the four datasets: 0.25%, 0.5%, 1.0%, 2.0%, and 4.0%. The model’s predictive performance was evaluated with MSE, MAE, and R^2 as key metrics. Here, to more clearly demonstrate the impact of pseudo-data quality on Bayesian optimization, we have removed the knowledge-driven module.

Table 7: Performance Evaluation of the Fine-tuned LLM with Varying Amounts of Prior Data.

Prior Data	Suzuki			Arylation			Buchwald		
	MSE↓	MAE↓	R^2 ↑	MSE↓	MAE↓	R^2 ↑	MSE↓	MAE↓	R^2 ↑
0.25%	1205.19	27.80	-0.53	793.92	24.20	-0.07	737.64	23.50	0.01
0.5%	774.70	21.13	0.02	1016.37	26.53	-0.36	617.28	19.81	0.17
1.0%	633.68	19.47	0.20	650.00	19.55	0.13	593.76	18.52	0.20
2.0%	479.09	15.92	0.39	462.52	15.75	0.38	365.60	13.98	0.51
4.0%	360.02	13.44	0.54	286.56	11.97	0.62	248.44	11.28	0.67

The empirical results, as presented in Table 7, demonstrate a clear and positive correlation between the volume of fine-tuning data and the model’s predictive accuracy. Specifically, we observed a monotonic improvement across all metrics as the data percentage increased. A critical threshold was identified at the 1.0% data level. At this point, the R^2 values for both the Suzuki and Arylation datasets became positive for the first time, signifying that the model’s predictions had surpassed the explanatory power of a simple mean-based model. Given that the R^2 scores for all three datasets converged to a reasonable performance level (approximately 0.2) at this stage, we concluded that 1.0% represents an efficient and effective baseline for prior data quantity, balancing model fidelity with data utilization. Therefore, this level was adopted for subsequent experiments.

We further investigated how the volume of the fine-tuned model, when leveraged as a surrogate for generating pseudo-data, affects the performance of a downstream Bayesian optimization task. For this, the models trained with 1%, 2%, and 4% of the prior data were utilized to guide the optimization process on four distinct reaction datasets: Suzuki, Arylation, Buchwald_{sub-1}, and Buchwald_{sub-2}.

Table 8: Bayesian Optimization Performance Using Pseudo-data from Models Fine-tuned with Different Prior Data Scales.

Dataset	Prior Data (%)	Volume Size	Best Found	Initial Value	Iteration of Best↓
Suzuki	1.0	50	88.98	65.95	37
	2.0	100	88.99	70.25	36
	4.0	200	91.41	52.13	40
Arylation	1.0	30	79.67	45.98	40
	2.0	60	81.65	48.95	40
	4.0	120	83.81	37.11	31
Buchwald _{sub-1}	1.0	7	79.63	54.07	31
	2.0	14	79.56	41.14	39
	4.0	28	80.91	32.36	15
Buchwald _{sub-2}	1.0	7	56.81	12.87	11
	2.0	14	53.95	15.59	38
	4.0	28	53.93	16.22	40

The outcomes, summarized in Table 8, reveal a nuanced and somewhat counter-intuitive trend. While increasing the fine-tuning data from 1% to 4% (e.g., from 50 to 200 samples for Suzuki) does yield modest improvements in the "Best Found" values for Suzuki and Arylation, the enhancement is not proportional to the four-fold increase in data. More strikingly, for the Buchwald_{sub-1} and

Buchwald_{sub-2} datasets, this trend does not hold. For instance, on Buchwald_{sub-2}, the model fine-tuned with 1% of the data achieved a superior result (56.81) compared to the models trained on 2% (53.95) and 4% (53.93) of the data. Similarly, for Buchwald_{sub-1}, the 1% model’s performance (79.63) was marginally better than the 2% model (79.56).

These findings suggest that simply increasing the volume of data for training the surrogate model does not guarantee enhanced performance in Bayesian optimization. Beyond a certain point, it appears to yield diminishing returns and can even be detrimental to the optimization outcome. This may indicate a complex interplay between the surrogate model’s accuracy and its ability to generalize across the search space, suggesting that a more moderately-sized, yet sufficiently representative, prior dataset may be optimal for guiding the exploration-exploitation balance in our optimization framework.

J GENERALIZATION TO BROADER SCIENTIFIC DOMAINS

To assess the universality of the ChemBOMAS framework across different scientific domains, we extended our evaluation to a materials science benchmark. This expansion aims to validate a core hypothesis: that the fundamental principle of combining knowledge-driven decomposition with data-driven fine-tuning can be extended beyond chemistry to complex black-box optimization problems. We selected a publicly available dataset representing a unique challenge in scientific discovery:

LNP3 originates from the field of materials science, aiming to address a critical challenge in nanomedicine: optimizing the formulation of lipid nanoparticles (LNPs) for drug delivery. The task involves optimizing LNP composition to effectively encapsulate cannabidiol (CBD). The dataset comprises 768 experimental formulations defined by a 5-dimensional parameter space, encompassing the type and quantity of solid lipids, liquid lipids, and surfactants. This constitutes a complex multi-objective optimization problem requiring the simultaneous achievement of three competing goals: maximizing drug loading and encapsulation efficiency while minimizing particle size. Its parameter space combines categorical and discrete variables, forming a non-trivial search space that represents the challenges encountered in real-world material formulation.

The performance comparison between ChemBOMAS and baseline methods on this dataset is summarized in Table 9.

Table 9: Performance Comparison on a Non-Chemical Scientific Benchmark.

Dataset	Method	Best Found	Initial Value	95% Max Iter↓	Iteration of Best↓
LNP3	ChemBOMAS	0.62	0.23	12	28
	Gollum	0.62	0.21	13	33
	BO	0.62	0.25	12	38
	LA-MCTS	0.47	0.44	4	15
	BO-ICL	0.60	0.15	24	38

In the LNP3 material formulation benchmark, ChemBOMAS demonstrated highly competitive performance. It successfully identified an optimal value of 0.62, matching the final performance achieved by GoLLuM and traditional Bayesian optimization (BO). More importantly, ChemBOMAS demonstrated higher sample efficiency, locating this optimal solution in just 28 iterations, compared to 33 for GoLLuM and 38 for BO. This result indicates that the framework’s structured exploration mechanism can effectively accelerate convergence even in non-chemical optimization scenarios. Notably, while LA-MCTS delivered strong initial performance, it prematurely converged to a suboptimal solution, highlighting the risks of overly aggressive early exploration.

Overall, testing results in the field of materials science demonstrate that ChemBOMAS’s fundamental architecture—namely, the synergistic integration of knowledge-based search space partitioning with data-driven model optimization—holds potential as a universal strategy. It has proven that beyond core chemical domains, this framework possesses equally robust applicability and competitiveness in accelerating black-box optimization across diverse scientific discovery tasks.