

# BioTool: A Comprehensive Tool-Calling Dataset for Enhancing Biomedical Capabilities of Large Language Models

Anonymous ACL submission

## Abstract

Despite the success of large language models (LLMs) on general-purpose tasks, their performance in highly specialized domains such as biomedicine remains unsatisfactory. A key limitation is the inability of LLMs to effectively leverage biomedical tools, which clinical experts and biomedical researchers rely on extensively in daily workflows. While recent general-domain tool-calling datasets have substantially improved the capabilities of LLM agents, existing efforts in the biomedical domain largely rely on in-context learning and restrict models to a small set of tools. To address this gap, we introduce BIOTOOL, a comprehensive biomedical tool-calling dataset designed for fine-tuning LLMs. BIOTOOL comprises 34 frequently used tools collected from the NCBI, Ensembl, and UniProt databases, along with 7,040 high-quality, human-verified query-API call pairs spanning variation, genomics, proteomics, evolution, and general biology. Fine-tuning a 4-billion-parameter LLM on BIOTOOL yields substantial improvements in biomedical tool-calling performance, outperforming state-of-the-art commercial LLMs such as GPT-5.1. Furthermore, human expert evaluations demonstrate that integrating a BIOTOOL-fine-tuned tool caller significantly improves downstream answer quality compared to the same LLM without tool usage, highlighting the effectiveness of BIOTOOL in enhancing the biomedical capabilities of LLMs.

## 1 Introduction

The rapid advancement of large language models (LLMs) has revolutionized natural language processing, enabling unprecedented performance across a wide range of general-purpose tasks (OpenAI, 2023; Bai et al., 2023). However, their capabilities in biomedical domains remain limited, which hinders their deployment in high-stakes, real-world biomedical applications (Chen et al., 2025; Li et al., 2025a). A key reason for this limitation

is the insufficient ability of LLMs to effectively leverage specialized biomedical tools (Jin et al., 2024). Unlike commonsense questions that can often be answered directly, biomedical problems typically require even expert researchers to consult external tools and databases before drawing reliable conclusions (NCBI, 2017). For instance, even for human biologists, the biological function of a raw nucleotide sequence cannot be reliably inferred without the aid of computational tools, such as BLAST or other sequence similarity-based methods (Altschul et al., 1990). As shown in Figure 1, LLMs that lack access to or integration with such tools are therefore prone to hallucinations and imprecise generalizations, undermining their reliability for scientific discovery.

Given these challenges, early attempts have integrated biomedical and chemistry tools into LLMs via in-context learning (Jin et al., 2024; Bran et al., 2024). Although these approaches show improvements, they are constrained to a small set of available tools due to limited context length. Moreover, biomedical research tools often support diverse and complex usage scenarios that cannot be fully captured by a few lines of textual prompts, which hinders LLMs from fully realizing their potential in biomedical tool usage. Inspired by the success of instruction-tuning-based tool-calling datasets in the general NLP domain (Liu et al., 2024; Patil et al., 2024), we address this gap by curating a comprehensive biomedical tool-calling dataset, BIOTOOL.

BIOTOOL is an instruction fine-tuning-style biomedical tool-calling dataset consisting of 7,040 high-quality, human-verified query-API call pairs. It includes 34 frequently used tools from the NCBI (NCBI, 2017), Ensembl (Hubbard et al., 2002), and UniProt (The UniProt Consortium, 2017) databases, spanning multiple subdomains such as variation, genomics, proteomics, evolution, and general biology. To construct the dataset, we

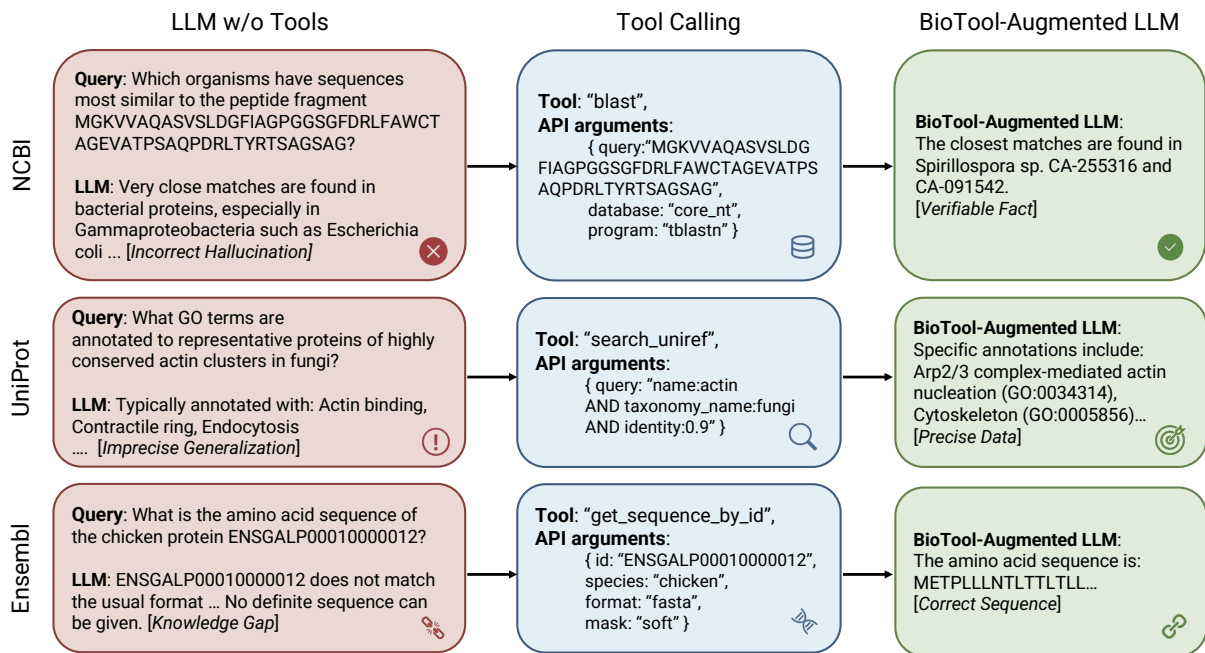


Figure 1: Comparison between answers generated by LLMs without tools and BIOTOOL-augmented LLMs for biomedical queries. LLMs without tools often hallucinate or produce imprecise answers (left), whereas BIOTOOL-augmented LLMs (right) generate API calls and retrieve critical information from biomedical databases, leading to higher-quality responses.

085 first manually select 34 tools from NCBI, Ensembl, 112  
086 and UniProt that are widely used in biomedical 113  
087 research. We then collect official documentation for 114  
088 these tools from their respective websites and use 115  
089 them to generate diverse combinations of API pa- 116  
090 rameters with the assistance of LLMs. The synthe- 117  
091 sized API calls are executed and filtered to remove 118  
092 cases with unavailable or uninformative responses, 119  
093 resulting in 3,829 unique API calls. Next, we 120  
094 prompt state-of-the-art reasoning models (OpenAI, 121  
095 2025) with these API calls and their corresponding 122  
096 responses to generate potential user queries. These 123  
097 queries are subsequently evaluated by an LLM- 124  
098 based judge to assess whether the API responses 125  
099 meaningfully support answering the queries, fol- 126  
100 lowed by a final round of human expert review 127  
101 focusing on biological relevance and correctness. 128  
102 This process yields 7,040 high-quality query-API 129  
103 call pairs, which is the final BIOTOOL dataset. 130

104 We evaluate the quality and effectiveness of 131  
105 BIOTOOL through two sets of experiments. First, 132  
106 we fine-tune several open-source LLMs ranging 133  
107 from 4B to 20B parameters on the BIOTOOL train- 134  
108 ing split and compare them with state-of-the-art 135  
109 commercial LLMs, including GPT-5.1, Gemini-3 136  
110 Pro, and Claude-4.5-Sonnet, using in-context learn- 137  
111 ing. Results on the test split show that smaller

112 LLMs fine-tuned with BIOTOOL significantly out- 113  
114 perform commercial LLMs with hundreds of times 114  
115 more parameters in terms of tool-calling quality. 115  
116 For example, a BIOTOOL-fine-tuned 4B Qwen-3 116  
117 model outperforms the best-performing Claude- 117  
118 4.5-Sonnet by 14.8% in overall API-calling qual- 118  
119 ity. Second, we conduct human evaluations to as- 119  
120 sess whether BIOTOOL-enhanced LLMs produce 120  
121 higher-quality answers from the perspective of 121  
122 biomedical researchers. On 1,048 test queries, a 122  
123 GPT-5.1 model augmented with oracle BIOTOOL 123  
124 API calls achieves 86% higher normalized answer 124  
125 quality compared to the same model without tool 125  
126 usage, demonstrating the intrinsic quality of the 126  
127 BIOTOOL dataset. Moreover, a GPT-5.1 model 127  
128 augmented with a BIOTOOL-fine-tuned API caller 128  
129 achieves 65% higher normalized answer quality 129  
130 compared to the raw GPT-5.1 model, highlighting 130  
131 the effectiveness of BIOTOOL in training tool-using 131  
132 LLMs and enhancing their biomedical capabilities. 132

## 2 Related Works 132

133 Early general-purpose tool-calling models, such as 133  
134 Toolformer (Schick et al., 2023) and Gorilla (Patil 134  
135 et al., 2024), established that LLMs can be trained 135  
136 to invoke external APIs, thereby grounding re- 136  
137 sponses in retrieved data to mitigate hallucinations. 137

Subsequent frameworks like ToolBench (Qin et al., 2023) and APiGen (Liu et al., 2024) advanced this capability by introducing scalable pipelines for generating synthetic instruction-tuning data. Despite these advancements, generalist models often struggle with specialized scientific domains like biomedicine because they rely on broad datasets that include only a negligible fraction of corresponding tools and frequently fail to adhere to the rigorous schema constraints of scientific databases. To address these limitations, domain-specific agents have emerged. GeneGPT (Jin et al., 2024) pioneered this shift by utilizing in-context learning (Wei et al., 2023) to enable access to NCBI Web APIs. Similarly, systems such as Sci-Agent (Li et al., 2025b) and ChemCrow (Bran et al., 2024) have successfully integrated tool-augmented agents for complex reasoning in scientific and chemical research. While more recent entries like Biomni (Huang et al., 2025) have introduced general-purpose agents for biomedical tasks, they primarily focus on a restricted subset of tools. Consequently, they lack the comprehensive, full-list interface to primary authoritative biomedical databases.

### 3 The BioTool Dataset

**Example of BioTool Data Entry**

**User Query**  
 Could you provide concise definitions for the major severe immunodeficiency disorders?

**Tool Information**  
 Database: "UniProt"  
 Tool: "human\_diseases"  
 API Endpoint: "search\_human\_diseases"

**API Arguments**  
 query: "name: immunodeficiency AND name: severe"  
 fields: "definition"  
 sort: "id asc"

**Observation**  
 (id: "DI-00171", definition: "An autosomal recessive immunologic disorder characterized by the loss of expression of MHC class II antigens on antigen-presenting cells..."),  
 (id: "DI-00305", definition: "A form of chronic granulomatous disease..."),  
 ...

This section details the development and composition of BIOTOOL. We first present an example

data entry from BIOTOOL to illustrate the structure of a query-API call pair. Each entry includes a *user query* field, which contains a realistic clinical or biomedical question expressed in free-form text. The *tool information* field provides descriptions of the tools required to answer the query, while the *API arguments* specify the input parameters for the corresponding API endpoint. Executing the API endpoint with these arguments returns an *observation*, which contains information used to augment the LLM’s response. We note that the observation is fully determined by the API endpoint and its arguments; it is included in the dataset for completeness and user convenience.

Next, we describe the sequential construction pipeline used to generate and verify biomedical tool calling pairs in Section 3.1, illustrated in Figure 2. We then provide a quantitative analysis of the resulting dataset, highlighting its functional utility and biological diversity in Section 3.2.

#### 3.1 Dataset Construction Pipeline

**Tool Selection** We select three major online API providers: the National Center for Biotechnology Information (NCBI), UniProt, and Ensembl as the tool source for BIOTOOL, motivated by their roles as the authoritative repositories within the global biomedical research infrastructure (Sayers, 2010; Ahmad et al., 2025; Yates et al., 2014). These three platforms are widely considered the definitive standard because they offer expansive and highly interoperable data spanning the entire central dogma of biology, encompassing the full spectrum from raw genomic sequences to functional protein annotations.

Across the three databases, we comprehensively review their websites and manually select tools that are critical for answering biomedical and clinical questions. During this process, we exclude tools with limited biomedical relevance (e.g., APIs that only return service or versioning information) as well as deprecated or unstable tools. As a result, we curate a diverse set of 34 tools comprising 124 API endpoints, each of which is frequently used in biomedical research workflows. The complete list of selected tools is provided in Appendix D. In addition, we collect the official documentation for each API endpoint from the corresponding website. These documents specify API usage, input arguments, constraints, and example calls, and serve as essential resources for subsequent stages of API call synthesis and user query generation.

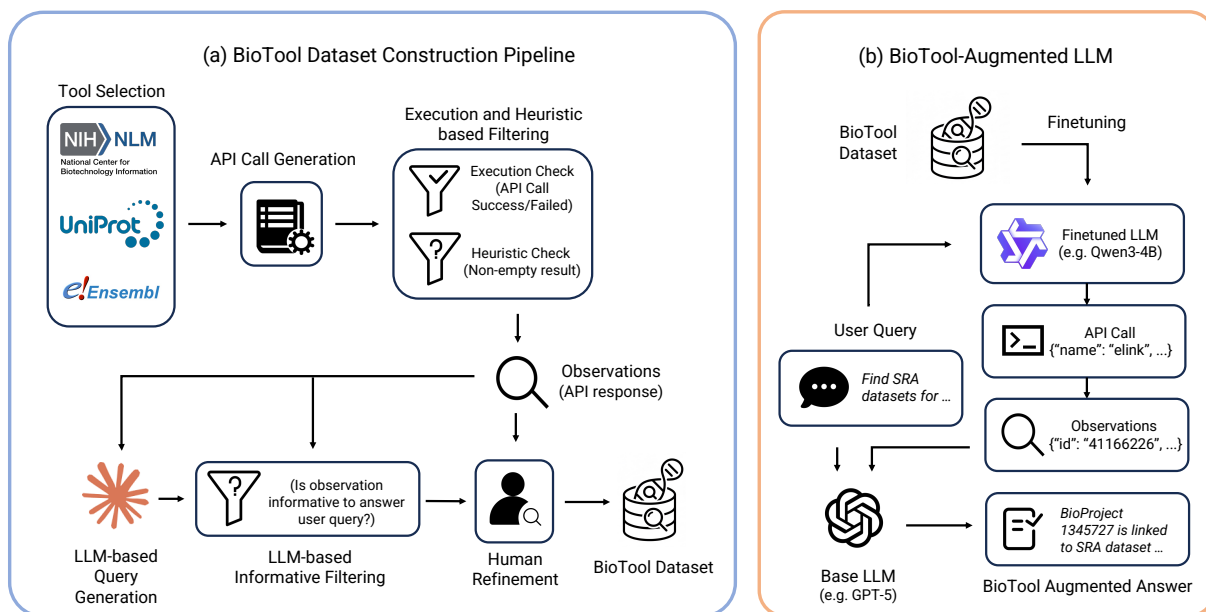


Figure 2: The systematic workflow of BIO TOOL spans from automated dataset construction to downstream application. Panel (a) illustrates the multi-stage construction pipeline, which includes initial tool selection from primary databases, automated API call generation, and a rigorous filtering process involving execution checks, heuristic validation, and LLM-based informativeness assessment. Panel (b) depicts the inference-time application, where specialized API-calling models fine-tuned on BIO TOOL enable base LLMs to retrieve grounded observations and generate verifiable biological answers.

**API Call Synthesis and Verification** Based on the curated tool set and associated documentation, we manually select critical API arguments corresponding to biologically meaningful identifiers for each API endpoint. These arguments, such as taxon IDs, gene symbols, and UniProt accession numbers, ensure that the synthesized API calls are biologically diverse and scientifically plausible. Given the selected arguments, we follow prior work (Liu et al., 2024) to randomly sample a large set of candidate API calls. These candidates are then executed to filter out cases that result in client errors, timeouts, or empty responses. To further improve data quality, we design a novel heuristic-based filtering strategy to remove API calls that are overly similar to existing ones, as well as those whose returned observations lack biological significance. Details of this heuristic filter are provided in Appendix A. After this verification process, we obtain a collection of 6,391 unique API calls.

**User Query Generation** Given the synthesized API calls, we leverage state-of-the-art LLMs to generate corresponding user queries, following a self-instruct-style paradigm established in prior work (Wang et al., 2022; Patil et al., 2024; Liu et al., 2024). Specifically, LLMs are prompted with an API call, its documentation, and its corresponding

observation, together with a small set of human-crafted in-context query-API call pairs, to generate realistic user queries.

To further improve the quality and biological relevance of BIO TOOL, we introduce two novel adaptations to ensure both the *necessity* and *sufficiency* of the API observations. First, to enforce *necessity*, we apply Chain-of-Thought (CoT) prompting (Wei et al., 2023) using a strong reasoning model (OpenAI o3 (OpenAI, 2025)) when generating user queries. The model is first prompted to summarize the technical details of the API observation into a natural-language description, which is then used to generate the final user query. This procedure ensures that the observation is required to answer the query, while keeping the query realistic and avoiding explicit references to specific tools or API calls. The detailed system and user prompts for this process are provided in Appendix B.1. Second, to ensure *sufficiency*, we employ another state-of-the-art LLM (Claude Haiku 4.5 (Anthropic, 2025)) to perform informativeness-based filtering, inspired by the LLM-as-a-judge framework (Zheng et al., 2023). The model is prompted to follow a structured rubric and classify a query-API call pair as informative if the observation contains at least one relevant fact or a partial summary that supports

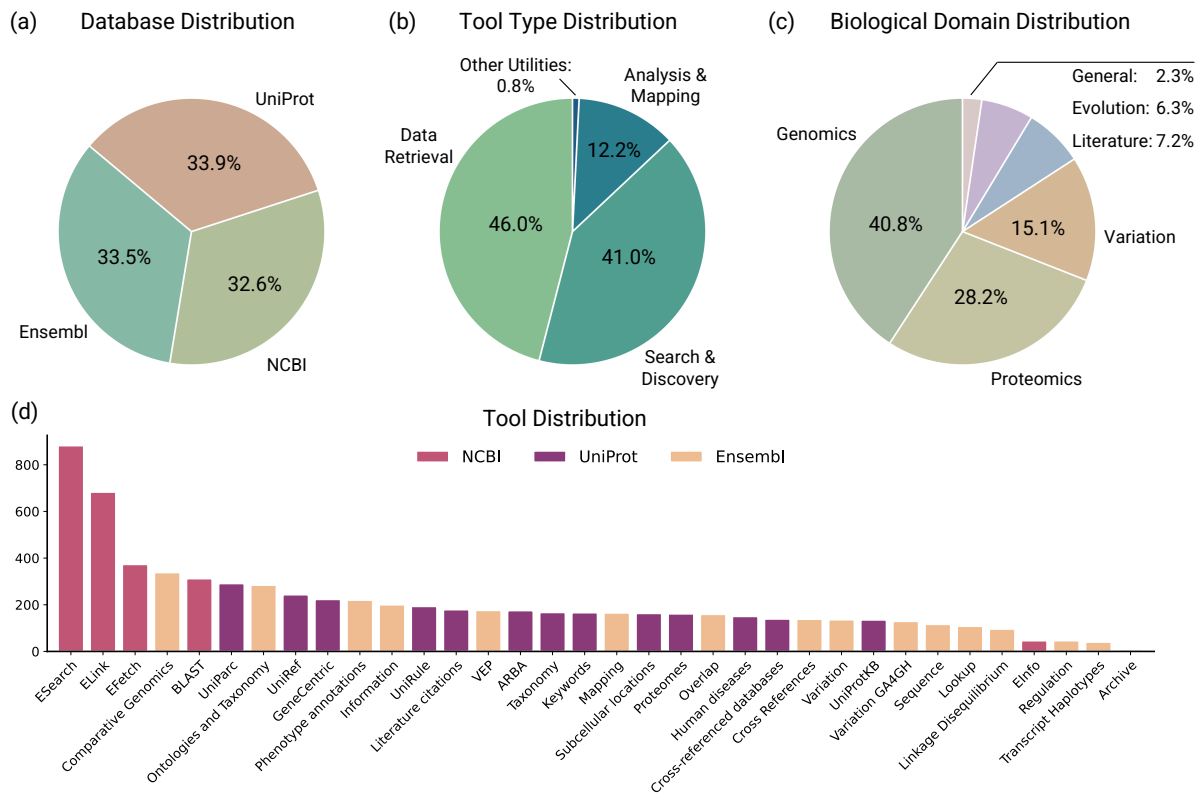


Figure 3: Distribution analysis of the 7,040 samples within BIOTOOL across four dimensions. Panel (a) shows the distribution across source databases. Panel (b) illustrates the distribution of samples by tool type. Panel (c) presents the distribution across various biological domains. Panel (d) delineates the distribution of user queries across the 34 distinct biological tools.

the user’s intent. Pairs in which the observation is unrelated to the query or too vague to support a concrete response are discarded. The specific judge prompts are provided in Appendix B.2.

**Human Refinement** The final stage involves a comprehensive manual review conducted by human evaluators with at least a college-level background in bioinformatics. The evaluators first identify and remove low-quality queries. For the remaining samples, they refine pedantic or unnatural phrasing and ensure the accuracy of biological terminology and nomenclature. After this round of filtering and correction, the final BIOTOOL dataset comprises 7,040 high-quality samples.

This instruction fine-tuning–style dataset is primarily used to train open-source LLMs as API-calling models, following training paradigms established in general-domain tool-calling datasets (Patil et al., 2024; Liu et al., 2024). A BIOTOOL-trained LLM can assist state-of-the-art LLMs in generating grounded and scientifically accurate responses, as illustrated in the right panel of Figure 2.

### 3.2 Data Statistics

The BIOTOOL dataset is derived from 34 distinct biological tools and 124 unique API endpoints, encompassing a wide array of scientific content categorized across several key dimensions. As shown in Figure 3(a), the distribution of tools across databases is well balanced, with comparable proportions from NCBI, UniProt, and Ensembl. Figure 3(b) illustrates the diversity of tool types included in BIOTOOL, ranging from data retrieval (e.g., nucleotide identifiers fetching) and search and discovery (e.g., phenotype-based gene discovery) to biological analysis and mapping (e.g., cross-referencing SNP identifiers). Figure 3(c) highlights the dataset’s broad scientific scope, covering domains such as genomics (e.g., gene tree querying), proteomics (e.g., protein sequence alignment), variation analysis (e.g., linkage disequilibrium analysis), and evolutionary biology (e.g., species-level taxonomy identification). Finally, Figure 3(d) shows that BIOTOOL includes both frequently accessed general-purpose tools and a long tail of specialized tools, all of which are

essential for complex scientific discovery across the central dogma.

## 4 Experimental Results

To evaluate the effectiveness of BIOTOOL, we first compare the API-calling capabilities of small open-source LLMs fine-tuned on BIOTOOL against state-of-the-art LLMs using in-context learning. We then conduct human expert evaluations to compare the answer quality of baseline LLMs with that of BIOTOOL-augmented LLMs.

### 4.1 Experimental Setup

**BioTool score** We define a BioTool performance score to automatically evaluate the capability of an LLM as an API caller on the BIOTOOL dataset, especially the alignment of retrieved information with the user’s intent. Specifically, assume we have the test set  $D = \{(q_1, o_1), \dots, (q_n, o_n)\}$ , where  $q_i$  is the  $i^{\text{th}}$  user query and  $o_i$  is the observation obtained from ground-truth API calling in the dataset. The BioTool score on this test set  $S(D)$  for a LLM API caller  $f$  is then defined as follows:

$$S(D) = \sum_{i=1}^n \text{Sim}(f(q_i), o_i) \quad (1)$$

where  $\text{Sim}(\hat{o}, o)$  computes the semantic embedding similarity of two text strings: the ground truth observation  $o$  and the corresponding observation  $\hat{o}$  from LLM API caller prediction. In practice, we use a MedCPT model (Jin et al., 2023) to get a sentence embedding for an observation. API calls may fail due to incorrect model generation, yielding an empty string  $\hat{o} = \varepsilon$ . In this case, we set  $\text{Sim}(\varepsilon, o) = 0$ . Intuitively, this score determines model performance by measuring whether the retrieved biological facts remain semantically similar to the required information, even when the technical implementation of the call differs from the reference.

**Additional Metrics** Based on the BioTool score, we define two additional metrics to further characterize model performance. Similar metrics have been widely adopted in existing API-calling benchmarks (Patil et al., 2025). Firstly, we define API calling success rate AS as follows:

$$AS(D) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}[\text{Sim}(f(q_i), o_i) > 0] \quad (2)$$

where  $\mathbf{1}[\cdot]$  is the indicator function. A zero similarity indicates API calling failure due to incorrect formatting, invalid API names, or improper parameter values. Conceptually, this metric focuses on the model’s capability to generate API calls that execute correctly and return a valid response containing data. Secondly, we define an exact match score EM as follows:

$$EM(D) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}[\text{Sim}(f(q_i), o_i) = 1] \quad (3)$$

which measures the proportion of predictions whose resulting observations exactly match the ground-truth reference observation, requiring the model to correctly identify the API endpoint and provide all required parameters with values that exactly match the reference.

**Models** In this study, we use four cutting-edge proprietary models, including GPT-5.1, GPT-5.1-Codex, Gemini 3 Pro, and Claude 4.5 Sonnet (OpenAI, 2025b,a; Google, 2025; Anthropic, 2025) under an in-context learning scheme. We use five open source models, which are gpt-oss-20b, Llama3.1-8B-Instruct, Qwen3-8B, Qwen2.5-7B-Instruct, and Qwen3-4B-Instruct (OpenAI et al., 2025; Grattafiori et al., 2024; Yang et al., 2025; Qwen et al., 2025), for BioTool-based fine-tuning.

### 4.2 Results on Tool Calling Capability

In this section, we first fine-tune small open-source models on the training split of the BIOTOOL dataset, which is randomly split under a four-to-one ratio. We use the state-of-the-art proprietary as baselines, and the evaluation for all models was conducted equally on the held-out test set consisting of 1,408 samples in terms of BioTool score. As shown in Table 1, there is a clear performance advantage for BIOTOOL-fine-tuned models over much larger LLMs under in-context learning. The fine-tuned 4B model achieved the highest overall BioTool score, representing a 14.8% improvement over the strongest proprietary model, Claude 4.5 Sonnet, and 68.8% higher performance than GPT-5.1. This gap suggests that the general-purpose pre-training of frontier LLMs together with in-context learning is insufficient to navigate the specialized technical constraints and precise parameter mappings of biological repositories. Instead, the high-density training signals within the BIOTOOL

Model	NCBI	UniProt	Ensembl	Overall
Proprietary Models				
GPT-5.1	15.5	78.7	72.8	55.4
GPT-5.1-Codex	14.8	80.5	74.5	56.3
Gemini 3 Pro	72.5	87.8	77.3	79.2
Claude 4.5 Sonnet	79.8	87.3	77.0	81.4
BioTool-fine-tuned Open-Source Models				
gpt-oss-20b	89.7	80.9	86.8	85.8
Llama3.1-8B-Ins	92.1	89.1	93.0	91.4
Qwen3-8B	90.7	85.5	81.3	85.9
Qwen2.5-7B-Ins	92.0	88.3	84.7	88.4
Qwen3-4B-Ins	<b>93.4</b>	<b>91.2</b>	<b>95.8</b>	<b>93.5</b>

Table 1: Comparative evaluation of models on the BIO TOOL dataset, measured by the BioTool score (higher is better). Scores are reported for each constituent database (NCBI, UniProt, Ensembl) and overall. Model names with the suffix *Ins* denote instruction-tuned variants. Bold values indicate the best performance in each column.

dataset allow significantly smaller models to acquire the necessary domain expertise that remains elusive to even the largest proprietary models.

### 4.3 Human Evaluation of Answer Quality

The ultimate criterion for assessing the usefulness of a tool-calling dataset is its ability to improve the quality of LLM-generated answers. To evaluate this, we use GPT-5.1 as the base model and compare its performance under three settings: (1) no tool augmentation, (2) augmentation with ground-truth BIO TOOL API calls, and (3) augmentation with a BioTool-fine-tuned Qwen3-4B-Instruct tool-calling model. We evaluate these three settings on 700 randomly sampled test queries using side-by-side human judgments by two annotators with college-level bioinformatics backgrounds. Annotators compare settings (1) vs. (2) and (1) vs. (3), selecting the better answer based on informativeness and task fulfillment, while rejecting vague or scientifically incorrect responses. The normalized win rates for the two comparisons are shown in Figure 4. Raw preference results and normalization procedures are detailed in Appendix C.

We observe that tool augmentation substantially improves the quality of biomedical answers, demonstrating that grounding LLMs in verifiable data from NCBI, Ensembl, and UniProt effectively mitigates domain-specific hallucinations and imprecise generalizations. The oracle con-

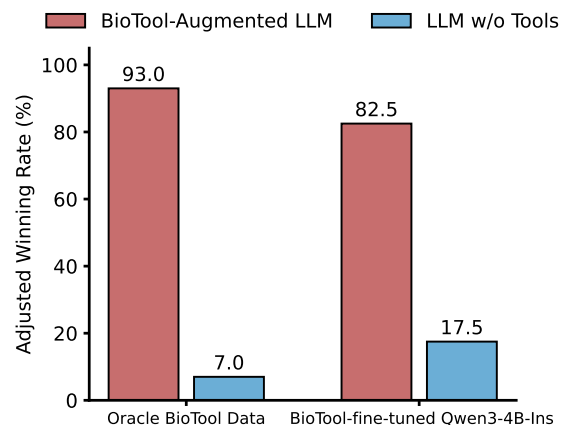


Figure 4: Human evaluation results comparing answer quality between BIO TOOL-augmented LLMs and LLMs without tool usage, using GPT-5.1 as the base model. The augmented settings include GPT-5.1 with oracle BIO TOOL data (left) and GPT-5.1 with a BIO TOOL-fine-tuned Qwen3-4B-Instruct tool caller (right).

figuration achieves a 93.0% win rate over the base model, highlighting the high quality of the BIO TOOL dataset. Similarly, the BIO TOOL-fine-tuned Qwen3-4B-Instruct model attains an 82.5% win rate, indicating that a small, fine-tuned model can improve the correctness and helpfulness of large commercial LLMs as judged by human evaluators, further demonstrating the practical utility of BIO TOOL.

### 4.4 Additional Results

We report results for the additional metrics under the same experimental settings in Table 1 to provide further insights into model behavior and dataset characteristics. As shown in Table 2, there is a clear divergence between Exact Match (EM) and API Success (AS), particularly for proprietary models. Although models such as Claude 4.5 Sonnet and Gemini 3 Pro achieve high AS scores, their EM remains extremely low (often below 5%), indicating difficulty in producing parameterizations that exactly match reference specifications. In contrast, the BIO TOOL-fine-tuned Qwen3-4B-Instruct achieves an EM more than six times higher than the best proprietary model, highlighting the necessity of fine-tuning for learning the precise syntax of biological APIs. The EM-AS gap also reflects the varying complexity of biological repositories. On the NCBI subset, proprietary models such as GPT-5.1 fail to achieve any exact matches and frequently encounter execution errors, likely due to strict iden-

Model	NCBI		UniProt		Ensembl		Overall (Avg)	
	EM	AS	EM	AS	EM	AS	EM	AS
Proprietary Models								
GPT-5.1	0.0	16.6	1.7	97.9	8.9	79.4	3.5	64.4
GPT-5.1-Codex	0.0	16.2	1.7	98.3	6.3	80.5	2.6	64.7
Gemini 3 Pro	3.2	76.7	2.1	99.8	16.3	82.0	7.1	86.2
Claude 4.5 Sonnet	3.6	85.5	1.1	<b>100.0</b>	15.0	82.4	6.5	89.4
BioTool-fine-tuned Open-Source Models								
gpt-oss-20b	30.7	94.5	4.0	98.5	30.4	92.6	21.7	95.2
Llama3.1-8B-Ins	42.2	94.3	7.2	99.8	43.0	95.9	30.8	96.7
Qwen3-8B	43.7	93.5	4.0	99.6	37.7	85.9	28.5	93.0
Qwen2.5-7B-Ins	45.6	93.9	5.5	<b>100.0</b>	42.3	87.9	31.1	94.0
Qwen3-4B-Ins	<b>65.6</b>	<b>94.5</b>	<b>9.3</b>	99.8	<b>51.8</b>	<b>98.1</b>	<b>42.3</b>	<b>97.4</b>

Table 2: Comparative evaluation of EM and AS metrics (higher is better). Model names with the suffix *Ins* denote instruction-tuned variants. Bold values indicate the best performance in each column.

tifier formats and nested parameters. Fine-tuned models, however, maintain high execution success, demonstrating that BIOTOOL trains functionally robust models that produce valid and biologically meaningful API calls even without exact string matches.

#### 4.5 Case Study

Comparison of API Call and Response to an Example User Query

##### User Query

Within 29:707234-757234, which variant pairs show strong linkage in the NextGen goat population?

##### BioTool-Fine-tuned Qwen3-4B API Call

```
get_ld_region(species="capra_hircus",
region="29:707234-757234",
population_name="NextGen:All",
d_prime=0.8)
```

##### BioTool-Fine-tuned Qwen3-4B Response

```
[ {"variation1": "rs661133063",
"variation2": "rs668584442", "d_prime":
1.0 }, ... ]
```

##### Claude API Call

```
get_ld_region(species="goat",
region="29:707234-757234",
population_name="NextGen")
```

##### Claude Response

```
{ "error": "Can not find internal name for
species 'goat' " }
```

##### Oracle API Call

```
get_ld_region(species="capra_hircus",
region="29:707234-757234",
population_name="NextGen:All", r2=0.5)
```

To illustrate how BIOTOOL enhances LLMs’ biomedical capabilities, we present an example comparing the API-calling behavior and tool-augmented answers of a BIOTOOL-fine-tuned Qwen3-4B-Instruct model and Claude 4.5 Sonnet on a complex genomic linkage query. The ICL-based Claude model exhibits low precision in handling domain-specific tool usage rules, directly passing literal natural-language terms such as “goat” and “NextGen” from the user query into the API arguments. In contrast, the BIOTOOL-fine-tuned model correctly maps these concepts to the required API parameters. This example demonstrates that BIOTOOL provides high-quality supervision for learning the implicit rules of biomedical tool usage, enabling reliable natural-language-to-tool mapping that state-of-the-art ICL-based LLMs often fail to achieve.

#### 5 Conclusion

In this work, we introduce BIOTOOL, a comprehensive biomedical tool-calling dataset comprising 7,040 human-verified query-API call pairs spanning 124 biomedical tools. Fine-tuning 4-billion-parameter LLMs on BIOTOOL leads to substantial improvements in API-calling performance, surpassing state-of-the-art commercial LLMs. Furthermore, human evaluations confirm that BIOTOOL-augmented LLMs generate more helpful, informative, and scientifically accurate answers compared to the same base models without tool usage, shedding light on the development of reliable biomedical agents in the future.

## 505 Limitations

506 Despite the performance gains observed with  
507 BIOTOOL, several limitations remain. Our current  
508 framework focuses exclusively on one-hop tool  
509 calling responses. This ignores more complex bio-  
510 logical problems that cannot be solved with a sin-  
511 gle API interaction and instead require multi-hop  
512 search results or iterative reasoning across multiple  
513 tools. Furthermore, we did not fine-tune an inde-  
514 pendent, specialized biomedical agent. This archi-  
515 tectural choice was necessitated by the extreme con-  
516 text length of raw biological observations, which  
517 frequently exceed our resource limitations even  
518 after post-processing and summarization. Future  
519 work should explore long-context architectures and  
520 multi-step reasoning trajectories to better support  
521 the most intricate clinical and research workflows.

## 522 References

523 Shadab Ahmad, Leonardo Jose da Costa Gonzales,  
524 Emily H Bowler-Barnett, Daniel L Rice, Minjoon  
525 Kim, Supun Wijerathne, Aurélien Luciani, Swaathi  
526 Kandasamy, Jie Luo, Xavier Watkins, Edd Turner,  
527 Maria J Martin, and the UniProt Consortium. 2025.  
528 [The uniprot website api: facilitating programmatic  
529 access to protein knowledge](#). *Nucleic Acids Research*,  
530 53(W1):W547–W553.

531 Stephen F. Altschul, Warren Gish, Webb Miller, Eu-  
532 gene W. Myers, and David J. Lipman. 1990. [Basic  
533 local alignment search tool](#). *Journal of Molecular  
534 Biology*, 215(3):403–410.

535 Anthropic. 2025. [System card: Claude haiku 4.5](#). Sys-  
536 tem Card.

537 Anthropic. 2025. [System card: Claude sonnet 4.5](#).

538 Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang,  
539 Xiaodong Deng, Yang Fan, Wenhang Ge, Yu Han,  
540 Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang  
541 Lin, Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang  
542 Lu, K. Lu, and 31 others. 2023. [Qwen technical  
543 report](#). *ArXiv*, abs/2309.16609.

544 Andres M. Bran, Sean Cox, Oliver Schilter, and 1 oth-  
545 ers. 2024. [Augmenting large language models with  
546 chemistry tools](#). *Nature Machine Intelligence*, 6:525–  
547 535.

548 Qiang Chen, Yifan Hu, Xiaohan Peng, and 1 others.  
549 2025. [Benchmarking large language models for  
550 biomedical natural language processing applications  
551 and recommendations](#). *Nature Communications*,  
552 16:3280.

553 Google. 2025. [Gemini 3 pro model card](#).

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, 554  
Abhinav Pandey, Abhishek Kadian, Ahmad Al- 555  
Dahle, Aiesha Letman, Akhil Mathur, Alan Schel- 556  
ten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh 557  
Goyal, Anthony Hartshorn, Aobo Yang, Archi Mi- 558  
tra, Archie Sravankumar, Artem Korenev, Arthur 559  
Hinsvark, and 542 others. 2024. [The llama 3 herd of  
560 models](#). *Preprint*, arXiv:2407.21783. 561

Kexin Huang, Serena Zhang, Hanchen Wang, Yuanhao 562  
Qu, Yingzhou Lu, Yusuf Roohani, Ryan Li, Lin Qiu, 563  
Junze Zhang, Yin Di, and 1 others. 2025. [Biomni: A  
564 general-purpose biomedical ai agent](#). *bioRxiv*, pages  
565 2025–05. 566

Tim Hubbard, David Barker, Ewan Birney, Graham 567  
Cameron, Yong Chen, Lucy Clark, Tony Cox, 568  
James Cuff, Val Curwen, Thomas Down, Richard 569  
Durbin, Eduardo Eyras, James Gilbert, Matthew 570  
Hammond, Lukasz Huminiecki, Arek Kasprzyk, 571  
Heikki Lehvaslaiho, Peter Lijnzaad, Chris Melsopp, 572  
and 16 others. 2002. [The ensembl genome database  
573 project](#). *Nucleic Acids Research*, 30(1):38–41. 574

Qiao Jin, Won Kim, Qingyu Chen, Donald C Comeau, 575  
Lana Yeganova, W John Wilbur, and Zhiyong Lu. 576  
2023. [Medcpt: Contrastive pre-trained transformers  
577 with large-scale pubmed search logs for zero-shot  
578 biomedical information retrieval](#). *Bioinformatics*,  
579 39(11):btad651. 580

Qiao Jin, Yifan Yang, Qingyu Chen, and Zhiyong Lu. 581  
2024. [Genegpt: Augmenting large language models  
582 with domain tools for improved access to biomedical  
583 information](#). *Bioinformatics*, 40(2):btae075. 584

Mingchen Li, Zaifu Zhan, Han Yang, Yongkang 585  
Xiao, Huixue Zhou, Jiatan Huang, and Rui Zhang. 586  
2025a. [Benchmarking retrieval-augmented large lan-  
587 guage models in biomedical nlp: Application, ro-  
588 bustness, and self-awareness](#). *Science Advances*,  
589 11(47):eadr1443. 590

Xuchen Li, Ruitao Wu, Xuanbo Liu, Xukai Wang, Jinbo 591  
Hu, Zhixin Bai, Bohan Zeng, Hao Liang, Leheng 592  
Chen, Mingrui Chen, Haitian Zhong, Xuanlin Yang, 593  
Xu-Yao Zhang, Liu Liu, Jia Li, Kaiqi Huang, Jiahao 594  
Xu, Haitao Mi, Wentao Zhang, and Bin Dong. 2025b. 595  
[Sciagent: A unified multi-agent system for generalis-  
596 tic scientific reasoning](#). *Preprint*, arXiv:2511.08151. 597

Zuxin Liu, Thai Hoang, Jianguo Zhang, Ming Zhu, 598  
Tian Lan, Shirley Kokane, Juntao Tan, Weiran Yao, 599  
Zhiwei Liu, Yihao Feng, and 1 others. 2024. [Api-  
600 gen: Automated pipeline for generating verifiable  
601 and diverse function-calling datasets](#). *arXiv preprint  
602 arXiv:2406.18518*. 603

Quinn McNemar. 1947. [Note on the sampling error  
604 of the difference between correlated proportions or  
605 percentages](#). *Psychometrika*, 12(2):153–157. 606

NCBI. 2017. [Database resources of the national cen-  
607 ter for biotechnology information](#). *Nucleic Acids  
608 Research*, 46(D1):D8–D13. 609

610	OpenAI, :, Sandhini Agarwal, Lama Ahmad, Jason	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	664
611	Ai, Sam Altman, Andy Applebaum, Edwin Arbus,	Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and	665
612	Rahul K. Arora, Yu Bai, Bowen Baker, Haiming Bao,	Denny Zhou. 2023. <a href="#">Chain-of-thought prompting elic-</a>	666
613	Boaz Barak, Ally Bennett, Tyler Bertao, Nivedita	<a href="#">its reasoning in large language models</a> . <i>Preprint</i> ,	667
614	Brett, Eugene Brevdo, Greg Brockman, Sebastien	arXiv:2201.11903.	668
615	Bubeck, and 108 others. 2025. <a href="#">gpt-oss-120b &amp; gpt-</a>		
616	<a href="#">oss-20b model card</a> . <i>Preprint</i> , arXiv:2508.10925.		
617	OpenAI. 2023. <a href="#">Gpt-4 technical report</a> .	An Yang, Anfeng Li, Baosong Yang, Beichen Zhang,	669
618	OpenAI. 2025a. <a href="#">Gpt-5.1-codex-max system card</a> .	Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao,	670
619	OpenAI. 2025b. <a href="#">Gpt-5.1 instant and gpt-5.1 thinking</a>	Chengen Huang, Chenxu Lv, Chujie Zheng, Day-	671
620	<a href="#">system card addendum</a> .	iheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao	672
621	OpenAI. 2025. <a href="#">Openai o3 and o4-mini system card</a> .	Ge, Haoran Wei, Huan Lin, Jialong Tang, and 41	673
622	System Card.	others. 2025. <a href="#">Qwen3 technical report</a> . <i>Preprint</i> ,	674
623	Shishir G Patil, Huanzhi Mao, Fanjia Yan, Char-	arXiv:2505.09388.	675
624	lie Cheng-Jie Ji, Vishnu Suresh, Ion Stoica, and		
625	Joseph E. Gonzalez. 2025. <a href="#">The berkeley function</a>	Andrew Yates, Kathryn Beal, Stephen Keenan, William	676
626	<a href="#">calling leaderboard (BFCL): From tool use to agentic</a>	McLaren, Miguel Pignatelli, Graham R. S. Ritchie,	677
627	<a href="#">evaluation of large language models</a> . In <i>Forty-second</i>	Magali Ruffier, Kieron Taylor, Alessandro Vullo, and	678
628	<i>International Conference on Machine Learning</i> .	Paul Flicek. 2014. <a href="#">The ensembl rest api: Ensembl</a>	679
629	Shishir G. Patil, Tianjun Zhang, Xin Wang, and	<a href="#">data for any language</a> . <i>Bioinformatics</i> , 31(1):143–	680
630	Joseph E. Gonzalez. 2024. Gorilla: Large language	145.	681
631	model connected with massive apis.	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan	682
632	Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan	Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin,	683
633	Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang,	Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang,	684
634	Bill Qian, Sihan Zhao, Runchu Tian, Ruobing Xie,	Joseph E. Gonzalez, and Ion Stoica. 2023. <a href="#">Judg-</a>	685
635	Jie Zhou, Mark Gerstein, Dahai Li, Zhiyuan Liu,	<a href="#">ing llm-as-a-judge with mt-bench and chatbot arena</a> .	686
636	and Maosong Sun. 2023. <a href="#">Toolllm: Facilitating large</a>	<i>Preprint</i> , arXiv:2306.05685.	687
637	<a href="#">language models to master 16000+ real-world apis</a> .		
638	<i>Preprint</i> , arXiv:2307.16789.	<b>A Heuristic Filter Detail</b>	688
639	Qwen, :, An Yang, Baosong Yang, Beichen Zhang,	In this section, we provide a more granular ex-	689
640	Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan	planation of the heuristic filtering strategies em-	690
641	Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan	ployed during the API call synthesis and verifica-	691
642	Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin	tion phase.	692
643	Yang, Jiayi Yang, Jingren Zhou, and 25 oth-	The specific filtering logic varies across the three	693
644	ers. 2025. <a href="#">Qwen2.5 technical report</a> . <i>Preprint</i> ,	integrated databases to account for differences in	694
645	arXiv:2412.15115.	their API architectures and the nature of the bi-	695
646	Eric Sayers. 2010. A general introduction to the e-	ological data they provide. For UniProt, which	696
647	utilities. <i>Entrez Programming Utilities Help [Inter-</i>	primarily provides functional protein annotations	697
648	<i>net]</i> . Bethesda (MD): National Center for Biotech-	and sequence data, we implement a strict dedupli-	698
649	<i>nology Information (US)</i> .	cation process by filtering out all API calls target-	699
650	Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta	ing the same unique identifier, such as a UniRef	700
651	Raileanu, Maria Lomeli, Luke Zettlemoyer, Nicola	entry ID or keyword entry ID, within the same	701
652	Cancedda, and Thomas Scialom. 2023. <a href="#">Toolformer:</a>	tool to prevent the over-representation of specific	702
653	<a href="#">Language models can teach themselves to use tools</a> .	proteins. Furthermore, we validate every execu-	703
654	<i>Preprint</i> , arXiv:2302.04761.	tion result by discarding any responses that return	704
655	The UniProt Consortium. 2017. <a href="#">Uniprot: the univer-</a>	empty lists or "null" search results, thereby ensur-	705
656	<a href="#">sal protein knowledgebase</a> . <i>Nucleic Acids Research</i> ,	ing that every retained API call contains at least	706
657	45(D1):D158–D169.	one valid, non-empty biological observation. En-	707
658	Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa	sembl requires a more nuanced dual-path approach	708
659	Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh	to balance diversity and validity when handling	709
660	Hajishirzi. 2022. <a href="#">Self-instruct: Aligning language</a>	complex genomic coordinates. For endpoints with	710
661	<a href="#">models with self-generated instructions</a> . In <i>Annual</i>	a restricted set of valid parameter combinations	711
662	<i>Meeting of the Association for Computational Lin-</i>	(defined as fewer than 20), where strict ID dedupli-	712
663	<i>guistics</i> .	cation would yield insufficient data, we selectively	713
		retain entries where the optional parameters, such	714
		as species or variants, are not identical, while query	715
		IDs are the same. Conversely, for "rich" APIs with	716

an expansive range of possible inputs, we apply a strategy similar to UniProt by filtering out any samples where the combination of required parameters is identical to an existing entry to prevent the model from over-fitting to specific genomic regions. For NCBI, the strategy is optimized for high-throughput tools and general metadata retrieval. We apply specialized heuristics to the BLAST tool, only retaining parameter combinations that involve unique query sequences and return at least one significant alignment hit, while removing matchless queries that cannot support downstream scientific reasoning. Other NCBI APIs are filtered using a standard heuristic method that removes identical identifier calls and verifies that the retrieved observations remain biologically informative.

## B Prompts

### B.1 Prompt for creating user queries

The following prompt is used to generate natural language user queries. It requires four distinct input streams: (1) the source document context, (2) API function specifications, (3) retrieved biological observations, and (4) in-context few-shot demonstrations.

#### System Prompt

You generate realistic biomedical questions that researchers naturally ask.

#### TASK OVERVIEW: TWO-PHASE REASONING

**Phase 1 – ANALYZE:** Map technical parameters to natural language concepts (Qualitative Mapping).

**Phase 2 – GENERATE:** Create TWO questions (one Broad/Implicit, one Specific/Qualitative).

#### PHASE 1: PARAMETER MAPPING & ABSTRACTION

Do not simply list parameters. You must translate *Data* into *Language*.

- **VERBATIM (Keep Exact):** Unique identifiers (gene symbols, rsIDs, accessions), Raw sequences (FASTA format), and Coordinates (e.g., “chr1:100-200”).
- **QUALITATIVE MAPPING (Translate Numbers/Codes):**
  - **Thresholds:** Map high numbers to adjectives like “strong” or “significant” (e.g.,  $d\_prime=0.8 \rightarrow$  “strong linkage”).
  - **Complex Codes:** Simplify technical strings to common terms (e.g., 1000GENOMES...  $\rightarrow$  “1000 Genomes data”).
- **IMPLICIT DEFAULTS (Selectively Omit):** If a parameter just ensures usable results (e.g.,

format=json), OMIT it. The user implies “good results” by asking.

#### PHASE 2: QUESTION GENERATION STRATEGY

Your goal is **Tool Bias:** The question should be specific enough that *this tool* is the logical choice, without naming it.

**Question 1: The “Implicit” Question (Natural & Broad).** A question a biologist asks a colleague. Hide strict parameters; assume the tool’s filters represent the broad intent.

**Question 2: The “Qualitative” Question (Specific Demand).** A researcher asking for a specific *quality*. Use adjectives to reflect parameter values (e.g., “strong LD”).

#### RULES:

- Ask only **one** question at a time.
- Keep the question concise and to the point.
- Do not use parentheses for supportive information.

#### OUTPUT FORMAT

Return JSON only with keys: param\_analysis, observation\_check, and questions.

742

#### User Prompt

#### GENERATE TWO NATURAL BIOMEDICAL QUESTIONS

#### API DOCUMENTATION

(Understand the tool’s specific bias and domain.)

[API\_DOC\_TEXT]

#### STEP 1 – Classify parameters

#### PARAMS:

[PARAMS\_JSON]

#### STEP 2 – Check observation

(Distinguish between AVAILABLE data and MISSING/EMPTY data)

#### OBSERVATION:

[OBSERVATION\_JSON]

#### STEP 3 – Write TWO questions

- **Question 1:** Broad intent (Natural tone, implies need for this specific tool).
- **Question 2:** Specific feature (Focus on a field that HAS data).
- **MUST include these identifiers:** [IDENTIFIERS\_BLOCK]

#### REFERENCE EXAMPLES

[FEW\_SHOTS\_TEXT]

Output JSON with param\_analysis, observation\_check, and questions.

743

## B.2 Prompt for informative check

The following prompt is used to evaluate whether an observation is informative enough to answer a specific user query. It requires (1) the natural language user query and (2) the JSON representation of the biological observation.

### System Prompt

You are an evaluator for dataset filtering.

**Goal:** Decide whether the OBSERVATION is informative (useful) for answering the USER QUERY.

**IMPORTANT CONTEXT:** Observations are POST-PROCESSED SUMMARIES (often partial). This is NOT a strict completeness check.

Be concise and deterministic.

### User Prompt

**USER QUERY:**  
[USER\_QUERY\_TEXT]

**OBSERVATION (tool output / retrieved data):**  
[OBSERVATION]

### RUBRIC

- **informative=true** if the observation contains at least ONE relevant, non-trivial fact that can be used to answer part of the query without inventing details.
  - Examples/partial lists still count.
  - Counts/aggregates/summaries still count.
  - If the query asks “which X” but the observation only gives a count or a few examples, that is still informative=true (partial answer).
- **informative=false ONLY** when:
  - Observation is an error / empty / placeholder, OR
  - Observation content is clearly unrelated to the query intent, OR
  - Observation is too vague to support even a single concrete statement relevant to the query.

### When writing the reason:

- Focus on what CAN be answered using the observation (even partially).
- If partial, put the missing parts into limitations, but do NOT flip to false just because it’s incomplete.
- Be concise and specific (name the fields/signals you used).

### OUTPUT FORMAT

(do NOT output JSON; output exactly these lines)

INFORMATIVE: true|false

REASON: <short reason>

LIMITATIONS: <optional; if none, write “none”>

## B.3 Prompt for generating answers

The following prompts are used to generate the final natural language responses for the human expert evaluation. These prompts require the original user query, the generated api call, and the corresponding biological observations as input.

### System Prompt (Base Model)

You are a concise, accurate biomedical assistant. Answer the user question as directly as possible in 2–5 sentences, using your general biomedical knowledge and reasonable domain assumptions. Do NOT mention tools, APIs, databases, internet access, or that you cannot look things up. Do NOT tell the user how to get the information. Answer directly. If the question asks for record-level details you cannot know exactly, give the best plausible answer in a natural, helpful way without refusals or meta statements (avoid phrasing like “I can’t”, “I don’t know”, “without risk of error”).

### User Prompt (Base Model)

[USER\_QUERY]

### System Prompt (Tool-Augmented)

You are a concise, accurate biomedical assistant. You are given a user question plus a tool call and its observation output. Answer in 2–6 sentences. Use the observation as primary evidence and your general knowledge. Write the answer directly as if you already know the facts. If the observation is insufficient, you may add general biomedical context, but do not invent record-level facts that should come from the observation.

### User Prompt (Tool-Augmented)

**User question:**  
[USER\_QUERY]

**API call (for context):**  
[API\_CALL\_JSON]

**Observation (tool output):**  
[OBSERVATION\_TEXT]

## C Human Evaluation Details

This section details the manual side-by-side assessment process and provides the raw preference data used to derive the winning rates reported in Section 4.3. A total of 700 samples were evaluated by researchers with biological backgrounds to compare the performance of tool-augmented models against the base GPT-5.1 generator. Table 3 sum-

770 marizes the distribution of these outcomes, includ-  
 771 ing cases where both models performed well, or  
 772 both failed to provide a satisfactory answer.

Model	Qwen3-4B	Oracle
Total Samples	700	700
Model A Wins	496 (70.9%)	610 (87.1%)
Model B Wins	41 (5.9%)	8 (1.1%)
Both Bad	99 (14.1%)	31 (4.4%)
Both Good	64 (9.1%)	51 (7.3%)

Table 3: Raw human preference distribution for pair-wise model evaluations. Model A refers to the tool-augmented configuration (Qwen3-4B or Oracle), and Model B refers to the base GPT-5.1 model without tool access.

773 To provide a balanced comparison that accounts  
 774 for samples where neither model showed a distinct  
 775 advantage, we calculated the adjusted winning rate  
 776 reported in Figure 4 based on the logic of McNe-  
 777 mar’s test (McNemar, 1947). In this framework,  
 778 “Both Good” and “Both Bad” responses are col-  
 779 lectively treated as ties ( $n_c = n_{good} + n_{bad}$ ) and  
 780 adjusted by splitting them evenly between the two  
 781 conditions. Specifically, given the raw preference  
 782 counts  $n_a$  and  $n_b$ , the adjusted preference numbers  
 783  $n'_a$  and  $n'_b$  were calculated as  $n'_a = n_a + \frac{1}{2}n_c$  and  
 784  $n'_b = n_b + \frac{1}{2}n_c$ .

## 785 D Tool and API List

786 The following part enumerates all tools and their  
 787 corresponding APIs used in this work, grouped by  
 788 data source.

### NCBI Tools

- **ESearch**
  - esearch
- **ELink**
  - elink
- **EFetch**
  - efetch
- **EInfo**
  - einfo
- **BLAST**
  - blast

### UniProt Tools

- **UniProtKB**
  - get\_uniprotkb\_entry
  - search\_uniprotkb
  - stream\_uniprotkb
- **UniRef**
  - get\_uniref\_by\_id
  - get\_uniref\_light

- get\_uniref\_members
- search\_uniref
- stream\_uniref
- **UniParc**
  - get\_uniparc\_by\_upi
  - get\_uniparc\_databases
  - get\_uniparc\_light
  - search\_uniparc
  - stream\_uniparc
- **GeneCentric**
  - get\_genecentric\_by\_accession
  - get\_genecentric\_by\_proteome\_id
  - search\_genecentric
  - stream\_genecentric
- **Proteomes**
  - get\_proteome\_by\_upid
  - search\_proteomes
  - stream\_proteomes
- **Literature citations**
  - get\_citation\_by\_id
  - search\_literature\_citations
  - stream\_literature\_citations
- **Keywords**
  - get\_keyword\_by\_id
  - search\_keywords
  - stream\_keywords
- **Human diseases**
  - get\_disease\_by\_id
  - search\_human\_diseases
  - stream\_human\_diseases
- **Subcellular locations**
  - get\_location\_by\_id
  - search\_subcellular\_locations
  - stream\_subcellular\_locations
- **Cross-referenced databases**
  - get\_crossref\_database\_by\_id
  - search\_crossref\_databases
  - stream\_crossref\_databases
- **Taxonomy**
  - get\_taxonomy\_by\_id
  - search\_taxonomy
  - stream\_taxonomy
- **UniRule**
  - get\_unirule\_by\_id
  - search\_unirule
  - stream\_unirule
- **ARBA**
  - get\_arba\_by\_id
  - search\_arba
  - stream\_arba
- **Archive**
  - get\_archive\_id

### Ensembl Tools

- **Comparative Genomics**
  - get\_alignment\_region
  - get\_cafe\_genetree\_by\_id
  - get\_cafe\_genetree\_by\_member\_id
  - get\_cafe\_genetree\_by\_member\_symbol
  - get\_genetree\_by\_id
  - get\_genetree\_member\_by\_id
  - get\_genetree\_member\_by\_symbol
  - get\_homology\_by\_id
  - get\_homology\_by\_symbol
- **Cross References**

791

792

- get\_xrefs\_by\_id
- get\_xrefs\_by\_symbol
- lookup\_xref\_name
- **Information**
  - get\_info\_analysis
  - get\_info\_assembly
  - get\_info\_assembly\_region
  - get\_info\_biotypes
  - get\_info\_biotypes\_groups
  - get\_info\_biotypes\_name
  - get\_info\_compara\_methods
  - get\_info\_compara\_species\_sets
  - get\_info\_external\_dbs
  - get\_info\_genomes
  - get\_info\_genomes\_accession
  - get\_info\_genomes\_assembly
  - get\_info\_genomes\_division
  - get\_info\_genomes\_taxonomy
  - get\_info\_species
  - get\_info\_variation\_population\_name
  - get\_info\_variation\_populations
  - get\_info\_variation\_sources
- **Linkage Disequilibrium**
  - get\_ld\_around\_variant
  - get\_ld\_pairwise
  - get\_ld\_region
- **Lookup**
  - lookup\_by\_id
  - lookup\_by\_symbol
- **Mapping**
  - map\_assembly
  - map\_cdna\_to\_genome
  - map\_cds\_to\_genome
  - map\_translation\_to\_genome
- **Ontologies and Taxonomy**
  - get\_ontology\_ancestors
  - get\_ontology\_ancestors\_chart
  - get\_ontology\_descendants
  - get\_ontology\_id
  - get\_ontology\_name
  - get\_taxonomy\_classification
  - get\_taxonomy\_id
  - get\_taxonomy\_name
- **Overlap**
  - overlap\_by\_id
  - overlap\_by\_region
  - overlap\_translation
- **Phenotype annotations**
  - get\_phenotype\_by\_accession
  - get\_phenotype\_by\_gene
  - get\_phenotype\_by\_region
  - get\_phenotype\_by\_term
- **Regulation**
  - get\_binding\_matrix
- **Sequence**
  - get\_sequence\_by\_id
  - get\_sequence\_by\_region
- **Transcript Haplotypes**
  - get\_transcript\_haplotypes
- **VEP**
  - vep\_by\_hgvs
  - vep\_by\_id
  - vep\_by\_region
- **Variation**
  - get\_variation
  - get\_variation\_by\_pmcid
  - get\_variation\_by\_pmids
  - variant\_recoder

- **Variation GA4GH**
  - get\_ga4gh\_callsets
  - get\_ga4gh\_datasets
  - get\_ga4gh\_features
  - get\_ga4gh\_featuresets
  - get\_ga4gh\_references
  - get\_ga4gh\_referencesets
  - get\_ga4gh\_variantannotationsets
  - get\_ga4gh\_variants
  - get\_query\_beacon