

BIOME-Bench: A Benchmark for Biomolecular Interaction Inference and Multi-Omics Pathway Mechanism Elucidation from Scientific Literature

Anonymous ACL submission

Abstract

Multi-omics studies often rely on pathway enrichment to interpret heterogeneous molecular changes, but pathway enrichment (PE)-based workflows inherit structural limitations of pathway resources, including curation lag, functional redundancy, and limited sensitivity to molecular states and interventions. Although recent work has explored using large language models (LLMs) to improve PE-based interpretation, the lack of a standardized benchmark for end-to-end multi-omics pathway mechanism elucidation has largely confined evaluation to small, manually curated datasets or ad hoc case studies, hindering reproducible progress. To address this issue, we introduce **BIOME-Bench**, constructed via a rigorous four-stage workflow, to evaluate two core capabilities of LLMs in multi-omics analysis: **Biomolecular Interaction Inference** and end-to-end **Multi-Omics Pathway Mechanism Elucidation**. We develop evaluation protocols for both tasks and conduct comprehensive experiments across multiple strong contemporary models. Experimental results demonstrate that existing models still exhibit substantial deficiencies in multi-omics analysis, struggling to reliably distinguish fine-grained biomolecular relation types and to generate faithful, robust pathway-level mechanistic explanations. Code and datasets are available on GitHub.¹

1 Introduction

Multi-omics profiling, including proteomics, metabolomics, and single-cell transcriptomics, among others, has become central to studying complex biological systems and disease mechanisms (Sanches et al., 2024; Baião et al., 2025). By capturing coordinated molecular changes across regulatory layers, multi-omics can support mechanistic hypotheses that connect perturbations to pathway-level phenotypes. However, translating heterogeneous multi-omics signals into coherent, causally

¹<https://anonymous.4open.science/r/BIOME-Bench>

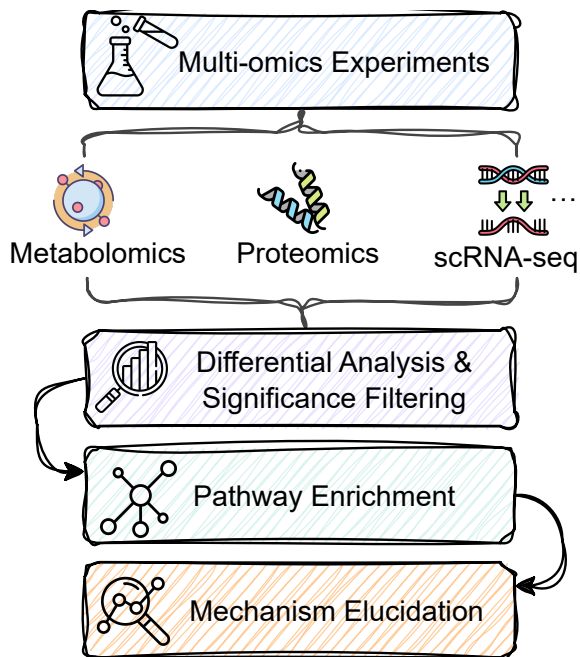


Figure 1: Overview of a pathway enrichment-based multi-omics mechanism elucidation workflow. Multi-omics experiments (e.g., metabolomics, proteomics, and scRNA-seq) are followed by differential analysis with significance filtering to identify perturbed entities, which are then mapped to biological pathways via enrichment analysis to support downstream mechanistic interpretation.

grounded explanations remains a major bottleneck (Mohr et al., 2024).

A common approach is pathway enrichment (PE) interpretation (Zhao and Rhee, 2023; Ryan V et al., 2025; Mubeen et al., 2022). As shown in Figure 1, practitioners typically perform differential analysis and significance filtering to identify perturbed entities, map them to pathways through enrichment, and interpret the enriched pathways to derive mechanistic narratives. PE therefore serves as a widely used interface between entity-level changes and higher-level biological processes (Elizarraras et al., 2024).

056 However, many limitations of PE-based interpretation arise from the structure of pathway re- 108
057 sources. Pathway knowledge bases are curated and 109
058 versioned, which introduces curation lag and can 110
059 omit newly discovered or context-specific mecha- 111
060 nisms (Agrawal et al., 2024). Enrichment results 112
061 also exhibit substantial functional redundancy due 113
062 to overlapping gene sets, often yielding long lists of 114
063 near-duplicate terms that are difficult to prioritize 115
064 (Balestra et al., 2023; Ge et al., 2025; Ozisik et al., 116
065 2022). More fundamentally, enrichment scores are 117
066 largely context-insensitive. They do not represent 118
067 molecular states (e.g., phosphorylation or activa- 119
068 tion), intervention directionality, or the multi-hop 120
069 causal structure required to connect perturbed en- 121
070 tities to pathway-level phenotypes. Consequently, 122
071 downstream interpretation often relies on post hoc 123
072 narrative stitching rather than state-aware mecha- 124
073 nistic reasoning (Slobodyanyuk et al., 2024). In 125
074 practice, researchers may obtain an interpretable 126
075 list of pathway names, but often still need addi- 127
076 tional analysis to determine how the observed per- 128
077 turbations give rise to the phenotype. 129

078
079 Recent advances in large language models 130
080 (LLMs) provide opportunities to improve PE-based 131
081 interpretation, for example by reducing redundancy 132
082 in enrichment outputs and using LLMs to gener- 133
083 ate mechanistic explanations (Ge et al., 2025; 134
084 Zhou et al., 2024). However, progress in this di- 135
085 rection is difficult to measure reliably because ex- 136
086 isting biomedical benchmarks rarely evaluate the 137
087 end-to-end capability required to produce mecha- 138
088 nistic interpretations from multi-omics observa- 139
089 tions. Many literature-grounded datasets focus 140
090 on question answering (Jin et al., 2019) or lo- 141
091 cal relation extraction (Zhang et al., 2019; Luo 142
092 et al., 2022), where relations are often assessed 143
093 in isolation from pathway context and phenotype- 144
094 level consequences. Meanwhile, existing pathway- 145
095 reasoning benchmarks typically emphasize naviga- 146
096 tion or subgraph operations rather than generating 147
097 state-aware mechanistic chains that connect per- 148
098 turbed entities to pathway-level phenotypes (Zhao 149
099 et al., 2025). In the absence of such end-to-end 150
100 benchmarks, evaluations often fall back on small 151
101 manually curated datasets or ad hoc case studies 152
102 (Ge et al., 2025). This limits coverage and repro- 153
103 ducibility, and it leaves open how well current 154
104 models can directly generate coherent pathway- 155
105 mechanism explanations on perturbed entities and 156
106 pathway context. 157

107 Motivated by these limitations, we consider 158

a complementary formulation that more closely 108
aligns with how scientists reason from observa- 109
tions to hypotheses: *end-to-end multi-omics path- 110
way mechanism elucidation*. Under this formula- 111
tion, a model receives only (i) a set of significantly 112
perturbed entities derived from multi-omics mea- 113
surements and (ii) a pathway context describing the 114
relevant biological process. The model must then 115
generate a coherent mechanistic explanation with- 116
out relying on explicit pathway retrieval and graph 117
traversal. This end-to-end perspective can reduce 118
redundancy by shifting the unit of interpretation 119
from overlapping pathway labels to mechanistic 120
hypotheses. It also better reflects research prac- 121
tice, where the goal is to elucidate mechanisms 122
from perturbed entities conditioned on pathway 123
context. The task is nontrivial because it requires 124
rich biological knowledge and multi-step causal 125
reasoning to connect perturbed entities to pathway- 126
level phenotypes (see Appendix B), which moti- 127
vates standardized, instance-level supervision for 128
reliable evaluation. 129

To address this gap, we introduce BIOME- 130
Bench, a literature-grounded benchmark designed 131
to evaluate LLMs on two core capabilities: 132
biomolecular interaction inference and end-to-end 133
multi-omics pathway mechanism elucidation. To 134
construct high-quality evaluation instances, we de- 135
velop a dedicated data construction workflow. As 136
shown in Figure 2, the workflow transforms path- 137
way information and evidence from biomedical 138
literature into structured, validated knowledge rep- 139
resentations through three sequential phases: (i) 140
Literature Search and Relevance Filtering, (ii) In- 141
formation Extraction and Standardization, and (iii) 142
Knowledge Structuring and Validation. The re- 143
sulting knowledge representations are then used to 144
formulate the benchmark tasks. Our contributions 145
are as follows: 146

- We formulate end-to-end multi-omics pathway 147
mechanism elucidation as a benchmarkable 148
task. Given perturbed entities and pathway 149
context, models must generate coherent, state- 150
aware, and intervention-consistent mechanistic 151
explanations without external retrieval or graph 152
traversal. 153
- We construct BIOME-Bench, a literature- 154
grounded benchmark with instance-level super- 155
vision that captures key mechanistic elements, 156
including entities, molecular states, biomolec- 157
ular relations, and pathway-level phenotypes, 158

159	enabling evaluation beyond surface narrative	206
160	quality.	207
161	• We design evaluation protocols that measure	208
162	mechanistic correctness at multiple granulari-	209
163	ties, including structured knowledge graph eval-	210
164	uation and holistic mechanism evaluation, en-	211
165	abling diagnosis of common failure modes in	212
166	current LLMs.	213
167	• We benchmark a diverse set of LLMs and	214
168	analyze their deficiencies, motivating future	215
169	research toward reliable end-to-end pathway	216
170	mechanism elucidation from multi-omics ob-	217
171	servations.	218
172	2 Related Work	219
173	2.1 Pathway Knowledge Bases	220
174	KEGG, Reactome, and WikiPathways provide cu-	221
175	rated pathway representations as molecular in-	222
176	teraction and reaction networks (Kanehisa et al.,	223
177	2025; Milacic et al., 2024; Agrawal et al., 2024).	224
178	KEGG organizes pathway maps using the KO	225
179	system to support cross-species mapping (Kane-	226
180	hisa et al., 2025). Reactome provides a con-	227
181	sistent, manually curated human-centric pathway	228
182	model with disease and drug context (Milacic et al.,	229
183	2024). WikiPathways supports community-driven	230
184	pathway curation (Agrawal et al., 2024). How-	231
185	ever, these databases are not designed as end-to-	232
186	end benchmarks for multi-omics mechanism elu-	233
187	cidation. They typically lack benchmark-ready,	234
188	instance-level supervision, such as explicit inter-	235
189	ventions, molecular states, relations, and phenotype	236
190	consequences, and they may lag behind newly pub-	237
191	lished mechanistic findings. Accordingly, we start	238
192	from KEGG pathways and ground benchmark in-	239
193	stances in scientific literature to derive state-aware	240
194	structured supervision and mechanistic explana-	241
195	tions. Notably, our method is applicable to other	242
196	pathway resources.	243
197	2.2 Literature-Grounded Benchmarks for	244
198	Pathway Reasoning and Relation	245
199	Extraction	246
200	Many biomedical benchmarks are derived from	247
201	scientific literature, but they do not evaluate	248
202	end-to-end pathway mechanism elucidation from	249
203	multi-omics observations. PubMedQA focuses on	250
204	evidence-based question answering over PubMed	251
205	abstracts (Jin et al., 2019). DrugProt (Luo et al.,	252
	2022) and ChemProt (Zhang et al., 2019) evalu-	253
	ate local relation extraction from abstracts or sen-	254
	tences. However, these benchmarks do not assess	
	biomolecular interaction inference under pathway	
	context.	
	More recently, BioMaze benchmarks	
	intervention-centric pathway reasoning and	
	introduces pathway subgraph navigation agents	
	such as PathSeeker (Zhao et al., 2025). However,	
	it emphasizes navigation and subgraph operations	
	rather than generating state-aware mechanistic	
	chains that connect perturbed entities to pathway-	
	level phenotypes. In contrast, BIOME-Bench	
	provides literature-grounded, instance-level su-	
	pervision with explicit molecular states, relations,	
	and pathway phenotypes, enabling fine-grained	
	evaluation of both interaction types and end-to-end	
	mechanistic explanations.	
	2.3 LLM-Based Systems for Multi-Omics and	
	Pathway Interpretation	
	Recent LLM-based systems such as AutoBA	
	(Zhou et al., 2024) and MAPA (Ge et al., 2025)	
	demonstrate how LLMs can support practical multi-	
	omics workflows, including analysis planning, tool	
	execution, and interpretation of pathway analysis	
	outputs. However, due to the lack of end-to-end	
	benchmarks for multi-omics pathway mechanism	
	elucidation, model performance in this setting is of-	
	ten evaluated through ad hoc case studies or small	
	manually curated datasets. This limits coverage	
	and reproducibility, and it makes it difficult to diag-	
	nose whether models can, given only perturbed enti-	
	ties and pathway context, directly generate accurate	
	pathway-mechanism explanations. To address this	
	gap, BIOME-Bench provides literature-grounded,	
	instance-level supervision and standardized evalua-	
	tion for this end-to-end setting.	
	3 Methodology	
	As shown in Figure 2, we propose a data con-	
	struction workflow that converts pathway infor-	
	mation and supporting evidence from biomedical	
	literature into structured, validated knowledge rep-	
	resentations in three sequential phases (detailed	
	prompts are provided in Appendix C). Building	
	on these representations, we introduce BIOME-	
	Bench , a literature-grounded benchmark for evalu-	
	ating LLMs on two tasks: Biomolecular Interaction	
	Inference and Multi-Omics Pathway Mechanism	
	Elucidation.	

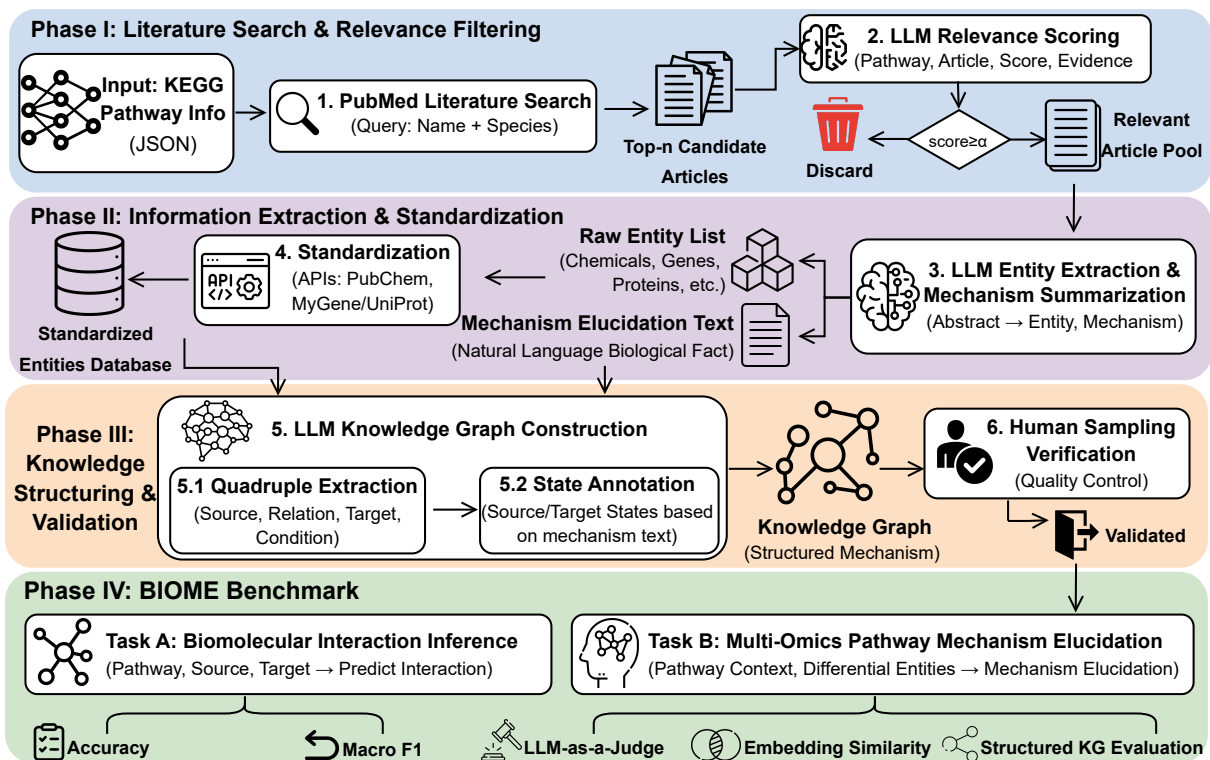


Figure 2: Workflow for constructing BIOME-Bench: (I) MeSH-guided PubMed retrieval with LLM relevance filtering; (II) LLM-based entity extraction and standardization; (III) state-aware knowledge graph construction with human sampling verification; and (IV) benchmark formulation with two tasks—biomolecular interaction inference and multi-omics pathway mechanism elucidation.

3.1 Phase I: Literature Retrieval and Relevance Filtering

To ensure strong biological validity, the construction process begins with a rigorous literature acquisition stage. Let

$$\mathcal{P} = \{p_1, p_2, \dots, p_n\}$$

denote a predefined set of target KEGG pathways. Each pathway p_i is characterized by its pathway name N_{p_i} and associated species S_{p_i} .

MeSH-guided Literature Retrieval. For each pathway p_i , we perform a structured literature search on the PubMed (White, 2020) database using Medical Subject Headings (MeSH) (National Library of Medicine (US)) to improve recall precision and semantic consistency. Specifically, the pathway name N_{p_i} is mapped to a set of MeSH descriptors, denoted as $\text{MeSH}(N_{p_i})$, and the species S_{p_i} is mapped to the corresponding MeSH organism term, denoted as $\text{MeSH}(S_{p_i})$.

The final PubMed query is constructed as a conjunction of pathway-related MeSH terms and species constraints:

$$Q(p_i) = \text{MeSH}(N_{p_i}) \wedge \text{MeSH}(S_{p_i}). \quad (1)$$

Executing $Q(p_i)$ yields an initial candidate document set:

$$D_{\text{cand}}(p_i) = \{d_1, d_2, \dots, d_m\}, \quad (2)$$

where each document d_j is indexed by PubMed and annotated with curated MeSH terms.

LLM-based Semantic and Mechanistic Relevance Scoring. MeSH-guided retrieval yields a controlled, high-recall candidate set, but MeSH annotations alone do not ensure that an article contains pathway-specific mechanistic evidence. To identify literature suitable for mechanism-level benchmarking, we use an LLM-based semantic evaluator with parameters θ .

Given a document and pathway pair (d, p_i) , the evaluator assigns a relevance score

$$s = f_{\theta}(d, p_i), \quad s \in [0, 10], \quad (3)$$

where f_{θ} aggregates multiple biologically motivated dimensions:

$$f_{\theta}(d, p_i) = g_{\theta}(\mathbf{S}), \quad \mathbf{S} = \begin{pmatrix} S_{\text{subj}} \\ S_{\text{spec}} \\ S_{\text{mol}} \\ S_{\text{ctx}} \end{pmatrix}. \quad (3)$$

S_{subj} measures whether the pathway’s biological process is the primary focus of the article, as opposed to a background mention. S_{spec} measures whether the organism studied matches the pathway species, and it also credits appropriate model organisms used to study human physiology or disease while penalizing biologically unrelated species. S_{mol} measures whether the article mentions pathway-defined molecular entities, such as genes, enzymes, or metabolites. S_{ctx} measures whether the article describes pathway regulation, such as activation, inhibition, or other modulatory effects, rather than merely reporting pathway presence.

We retain a document only if it exceeds a strict relevance threshold:

$$D_{\text{rel}}(p_i) = \{d \in D_{\text{cand}}(p_i) \mid f_{\theta}(d, p_i) \geq \alpha\}. \quad (4)$$

In this work, we set $\alpha = 8$ to prioritize articles in which the target pathway is central and supported by explicit molecular and regulatory evidence.

3.2 Phase II: Information Extraction and Entity Standardization

3.2.1 LLM-based Mechanistic Extraction

For each document $d \in D_{\text{rel}}(p_i)$, we process the abstract using a LLM to produce two complementary outputs:

1. **Raw Entity Set** E_{raw} : a collection of mentioned biological entities categorized into Chemicals, Genes/Proteins, and Phenotypes.
2. **Mechanism Description** M_{text} : a coherent natural language explanation describing the molecular interactions and regulatory mechanisms reported in the document. This text later serves as ground truth for generative evaluation.

3.2.2 Entity Normalization and Ontology Mapping

To ensure interoperability with external biological resources, we normalize each raw entity $e \in E_{\text{raw}}$ to a canonical identifier using an ontology resolution function $\phi(e)$. Specifically,

- **Chemical entities** are mapped to PubChem (Kim et al., 2016) compound identifiers (CIDs) using PubChemPy.
- **Genes and proteins** are mapped to NCBI Gene (Brown et al., 2015) identifiers or

UniProt (Consortium, 2015) accessions via MyGene.info.

To improve benchmark quality, we discard a candidate document if any entity cannot be resolved to a valid identifier. Only documents for which all entities are successfully normalized are retained, yielding the standardized entity set:

$$E_{\text{std}} = \{\phi(e) \mid e \in E_{\text{raw}} \wedge \forall e' \in E_{\text{raw}}, \phi(e') \neq \emptyset\} \quad (5)$$

3.3 Phase III: Knowledge Structuring and Validation

This phase leverages an LLM to convert the extracted mechanistic information and standardized entities into a fine-grained, state-aware knowledge graph representation.

3.3.1 Interaction Quadruple Extraction

We first extract the core interaction structure from M_{text} . Each interaction is represented as a quadruple:

$$T_{\text{core}} = (e_s, r, e_t, c), \quad (6)$$

where $e_s, e_t \in E_{\text{std}}$ denote the source and target entities, $r \in \mathcal{R}$ is a relation type drawn from a controlled biological vocabulary, and c specifies the biological condition under which the interaction occurs.

3.3.2 Biological State Annotation

To capture dynamic molecular behavior, we further annotate entity-specific biological states. Let σ_s and σ_t denote the states of the source and target entities, respectively (e.g., mutated, overexpressed). Incorporating state information yields a state-aware hexaplet representation:

$$T_{\text{final}} = (e_s, \sigma_s, r, e_t, \sigma_t, c). \quad (7)$$

This formulation enables the benchmark to distinguish subtle yet critical mechanistic differences, such as changes in protein abundance versus post-translational modifications.

3.3.3 Human Expert Verification

To establish a high-confidence gold standard, we perform human-in-the-loop validation. A randomly sampled subset of the constructed knowledge graph entries is reviewed by domain experts in molecular biology and systems biology, who cross-check each entry against the supporting literature to verify its accuracy and grounding (see Appendix A for details).

3.4 Phase IV: BIOME-Bench Task Formulation

Based on the curated and validated knowledge representations, BIOME-Bench defines two complementary evaluation tasks.

3.4.1 Task A: Biomolecular Interaction Inference

This task evaluates an LLM’s ability to infer precise molecular relationships within a pathway context. Given a pathway p_i , a source entity e_s with state σ_s , a target entity e_t with state σ_t , and a biological condition c , the model is required to predict the correct interaction relation from a finite controlled vocabulary \mathcal{R} :

$$\hat{r} = \arg \max_{r \in \mathcal{R}} P(r \mid p_i, e_s, \sigma_s, e_t, \sigma_t, c). \quad (8)$$

Model performance is evaluated using Accuracy and Macro-F1 over relation labels in \mathcal{R} (with invalid or unrecognized predictions treated as errors for the corresponding ground-truth label).

3.4.2 Task B: Multi-Omics Pathway Mechanism Elucidation

This task simulates realistic omics-driven pathway analysis scenarios. The model is provided with a pathway context p_i and a set of differentially observed entities

$$E_{\text{diff}} \subseteq E_{\text{std}},$$

and is required to generate a coherent mechanistic explanation \hat{Y} that elucidates the biological interactions, regulatory relationships, and molecular processes connecting these entities within the given pathway context.

We adopt a multi-dimensional evaluation strategy to comprehensively assess the quality of the generated explanations:

- **LLM-as-a-Judge:** Given the model-generated explanation \hat{Y} and the literature-derived ground truth M_{text} , judge model evaluates the output across four biologically motivated dimensions (scale: 1-5): Phenotype Coverage, Causal Reasoning, Factuality, and Hallucination.
- **Structured Knowledge Graph Evaluation:** Based on the literature-derived knowledge graph, we adopt a closed-set evaluation protocol. Specifically, an extraction model is provided with the standardized knowledge graph

and is only allowed to select supporting knowledge tuples from this graph based on the generated explanation \hat{Y} . As a result, the predicted tuple set satisfies

$$\mathcal{T}_{\text{pred}} \subseteq \mathcal{T}_{\text{GT}},$$

ensuring that no out-of-graph or hallucinated knowledge can be introduced.

Under this constraint, factual completeness is quantified using coverage, defined as:

$$\text{Coverage} = \frac{|\mathcal{T}_{\text{pred}}|}{|\mathcal{T}_{\text{GT}}|}.$$

- **Semantic Embedding Similarity:** We compute the cosine similarity between vector representations of \hat{Y} and the ground truth mechanism text M_{text} using an LLM-based embedding model, providing a complementary measure of semantic alignment.

4 Experiments

4.1 Benchmark Statistics and Characteristics

Species	Number of Pathways	Number of Entities	Number of Processes and Phenotypes	Task A Biomolecular Interaction Inference	Task B Multi-Omics Pathway Mechanism Elucidation
hsa	80	1,349	1,781	4,032	490
mmu	80	1,356	1,860	4,162	496
rno	80	1,141	1,265	3,384	361
Total	240	3,846	4,906	11,578	1,347

Table 1: Benchmark statistics of BIOME-Bench across species.

BIOME-Bench is a multi-species benchmark that covers three commonly used organisms: hsa (human), mmu (mouse), and rno (rat). Table 1 summarizes the core statistics, including the numbers of curated pathways, standardized entities, process and phenotype terms, mechanism analysis instances, and knowledge graph relations. Overall, it includes 1,347 instances for multi-omics mechanism elucidation and 11,578 instances for biomolecular interaction inference, both evaluated under consistent pathway contexts.

4.2 Experimental Setup

We evaluate a range of LLMs, including Qwen3-14B, 32B and 235B (Yang et al., 2025), DeepSeek-V3.2-R1 (Liu et al., 2025), GLM-4.6 (GLM et al., 2024), Gemini3-Pro (Google DeepMind), GPT-5.2 (OpenAI), Doubao-Seed-1.8 (Seed et al., 2025),

Model	Biomolecular Interaction Inference		Multi-Omics Pathway Mechanism Elucidation				Similarity	Coverage	Avg.
	Acc	Macro-F1	LLM-as-a-Judge						
			Phenotype Coverage	Causal Reasoning	Factuality	Hallucination			
Qwen3-14B	47.43	43.72	3.12	3.31	3.97	4.64	78.73	42.38	64.13
Qwen3-32B	41.84	40.51	3.00	3.26	3.89	4.79	78.98	45.43	63.20
Qwen3-235B	51.41	46.21	3.66	4.32	4.54	4.40	77.34	42.22	69.45
DeepSeek-V3.2-R1	53.10	47.52	3.28	4.31	4.20	4.10	75.12	40.76	66.79
GLM-4.6	53.60	50.08	3.50	4.14	4.32	4.18	76.89	39.95	67.92
Gemini3-Pro	52.34	46.54	3.60	4.57	4.59	4.54	77.21	41.13	69.74
GPT-5.2	54.66	50.70	3.68	4.58	4.69	4.62	71.38	37.49	70.70
Doubao-Seed-1.8	55.42	50.40	3.81	4.69	4.69	4.57	74.92	39.72	71.96
Intern-S1-235B	54.15	50.36	3.96	4.28	4.75	4.92	78.71	44.49	73.24
S1-Base-671B	54.68	50.41	4.02	4.48	4.76	4.83	77.36	44.45	73.59

Table 2: Performance comparison of large language models on BIOME-Bench. Similarity refers to the cosine similarity between the embeddings of the generated answer and the reference answer. Coverage refers to the knowledge graph coverage derived using Structured Knowledge Graph Evaluation. **Avg.** is the arithmetic mean of all metrics after normalizing each score to the 0–100 range.

469 Intern-S1-235B (Bai et al., 2025) and S1-Base-
470 671B (ScienceOne). We deploy Qwen3-14B, 32B
471 and 235B with vLLM (Kwon et al., 2023), and
472 access the remaining models through their respec-
473 tive APIs. For all models, we conducted a sin-
474 gle experimental run with temperature = 0 and
475 max_tokens = 10240.

476 For LLM-as-a-Judge and information extrac-
477 tion for Structured Knowledge Graph Evaluation,
478 we use Qwen3-32B with temperature = 0 and
479 max_tokens = 10240. For embedding-based sim-
480 ilarity, we use Qwen3-Embedding-8B (Zhang et al.,
481 2025). All experiments run on a server with 8×
482 NVIDIA A100-SXM4-80GB GPUs.

483 4.3 Main Results

484 Table 2 summarizes the overall performance of con-
485 temporary LLMs on BIOME-Bench. For biomolec-
486 ular interaction inference, most models fall within
487 a relatively narrow range (Acc 41.84% to 55.42%,
488 Macro-F1 40.51% to 50.70%), indicating that the
489 task remains challenging and that current progress
490 is largely incremental. Doubao-Seed-1.8 achieves
491 the highest accuracy (55.42%), whereas GPT-5.2
492 attains the best Macro-F1 (50.70%). The discrep-
493 ancy suggests that leading systems are comparable
494 in overall correctness but differ in robustness on
495 minority classes.

496 For multi-omics pathway mechanism elucida-
497 tion, the LLM-as-a-Judge dimensions exhibit a
498 consistent pattern. Factuality and hallucination
499 control are generally strong, typically above 4,
500 while phenotype coverage is lower, around 3 to
501 4. This suggests that current models often produce
502 plausible and reasonably grounded narratives, but
503 they frequently omit required process or pheno-

504 type elements, making phenotype-level reasoning
505 a primary bottleneck. Doubao-Seed-1.8 and GPT-
506 5.2 score highest on causal reasoning (4.69 and
507 4.58) while maintaining high factuality (both 4.69).
508 Intern-S1 performs particularly well on evidence
509 alignment, achieving the best hallucination control
510 (4.92) together with reliable phenotype coverage
511 (3.96). In contrast, S1-Base-671B attains the high-
512 est phenotype coverage (4.02) and factuality (4.76),
513 indicating an advantage in producing more com-
514 plete and reliable mechanistic narratives.

515 In addition, similarity and knowledge graph cov-
516 erage do not fully reflect judge-assessed mechanis-
517 tic quality. Qwen3-32B attains the highest similar-
518 ity (0.7898) and KG coverage (45.43), but it does
519 not lead on causal reasoning or factuality. This
520 indicates that similarity and KG coverage primar-
521 ily capture proximity to reference phrasing and the
522 breadth of entity linking, rather than the correct-
523 ness and sufficiency of multi-step causal explana-
524 tions. In addition, similarity can be sensitive to
525 verbosity: stronger models often produce longer
526 responses with auxiliary background or contextual
527 text beyond the core mechanistic trace, which can
528 reduce surface-form similarity despite improving
529 explanatory content. Overall, strong BIOME per-
530 formance requires balancing mechanistic complete-
531 ness, causal coherence, and evidence grounding.
532 Improvements in similarity or KG coverage there-
533 fore do not necessarily translate into better judge-
534 assessed mechanistic reasoning.

535 Finally, as scientific foundation models, Intern-
536 S1-235B and S1-Base-671B achieve stronger aver-
537 age performance than general-purpose models, sug-
538 gesting that continued scientific-domain training
539 has a positive effect on multi-omics pathway mech-

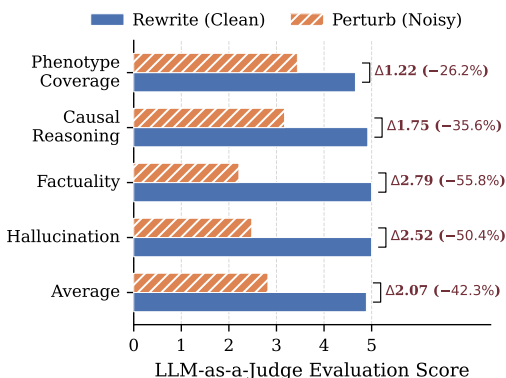
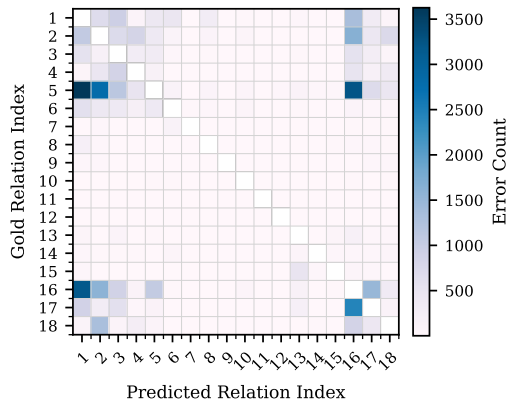


Figure 3: Sensitivity of Qwen3-32B judge to semantic perturbations. Scores are reported for rewrite and perturb. Drop% denotes the relative score decrease from rewrite to perturb.



- | | |
|-----------------------------|---------------------|
| 1: activates | 10: ubiquitinates |
| 2: inhibits | 11: glycosylates |
| 3: upregulates expression | 12: methylates |
| 4: downregulates expression | 13: produces |
| 5: regulates | 14: consumes |
| 6: binds | 15: converts to |
| 7: dissociates from | 16: leads to |
| 8: phosphorylates | 17: increases level |
| 9: dephosphorylates | 18: decreases level |

Figure 4: Error confusion matrix for Biomolecular Interaction Inference. Rows are gold relation types and columns are predicted types. Color encodes the count of misclassified gold→predicted relations.

540 anism elucidation. Nevertheless, the two BIOME-
 541 Bench tasks remain highly challenging for current
 542 models. Existing models still struggle to reliably
 543 distinguish fine-grained relation types and to de-
 544 rive robust mechanistic explanations directly from
 545 perturbed entities.

546 4.4 Validity of LLM-as-a-Judge

547 To ensure the reliability of automated evaluation in
 548 BIOME, we conducted experiments to assess the
 549 effectiveness of the LLM-as-a-judge. For each test
 550 instance, we construct two candidate answers from
 551 the gold reference. The first is a rewrite version that
 552 paraphrases the reference while preserving entities,
 553 relations, and causal semantics. The second is a
 554 perturb version that introduces targeted semantic er-
 555 rors by replacing key entities and/or interaction re-
 556 lations, while keeping the text fluent. A valid judge
 557 should remain insensitive to semantics-preserving
 558 rewrites, yet substantially penalize perturbations
 559 that degrade mechanistic correctness.

560 Figure 3 shows that Qwen3-32B consistently
 561 distinguishes semantics-preserving rewrites from
 562 mechanistically corrupted perturbations. The
 563 rewrite answers obtain near-ceiling scores of
 564 4.66/4.92/5.00/5.00 on Phenotype Coverage,
 565 Causal Reasoning, Factuality, and Hallucination,
 566 with an overall average of 4.89. In contrast, the per-
 567 turb answers drop to 3.44/3.17/2.21/2.48 and an
 568 average of 2.82. This corresponds to relative score
 569 decreases of 26.2%, 35.6%, 55.8%, and 50.4% for
 570 the four dimensions, and 42.3% on average. The
 571 large and consistent drops indicate that the judge
 572 is highly sensitive to entity and relation perturba-
 573 tions even when the text remains fluent, supporting

574 the validity of Qwen3-32B as an LLM-as-a-Judge
 575 under our reference-grounded rubric.

576 4.5 Interaction Type Confusion in 577 Biomolecular Inference

578 To analyze failure modes in *Biomolecular Inter-*
 579 *action Inference*, we aggregate all misclassified
 580 relations and summarize the confusion matrix in
 581 Figure 4. Errors are dominated by coarse regulatory
 582 and causal labels. Across all mistakes, predictions
 583 most often fall into leads_to (relation 16, 11,079
 584 cases), activates (1, 10,469), and inhibits (2,
 585 7,930), indicating a strong tendency to default to
 586 generic causality or signed regulation under uncer-
 587 tainty.

588 Two confusions account for most of this mass.
 589 First, the underspecified label regulates (5)
 590 is frequently polarized or rewritten as causa-
 591 tion, most often to activates (3,626), inhibits
 592 (2,783), or leads_to (3,222). Second, models
 593 blur pathway-level causation and direct regulation:
 594 gold activates/inhibits are often predicted as
 595 leads_to (1,325/1,643), while gold leads_to
 596 is over-interpreted as activates or inhibits
 597 (3,174/1,573).

598 Overall, the results indicate that current models
 599 still fall short in distinguishing fine-grained inter-
 600 action types and separating direct regulation from
 601 pathway-level causation, making it difficult to re-
 602 cover complete, verifiable mechanistic chains.

603 Limitations

604 BIOME-Bench has three main limitations. First,
605 each instance conditions on a single pathway con-
606 text; extending to multi-pathway settings with ex-
607 plicit crosstalk and compositional reasoning is an
608 important direction. Second, our supervision is
609 largely organized as a one-to-one mapping between
610 a pathway and a supporting paper; future work
611 should build multi-pathway, multi-document mech-
612 anistic graphs that integrate dispersed evidence and
613 enable reasoning over interconnected mechanisms.
614 Third, our current instruction set has limited stylistic
615 and structural diversity, while modern LLMs
616 can be highly prompt-sensitive due to training-data
617 distribution effects; as a result, benchmark perfor-
618 mance may vary with prompt phrasing and intro-
619 duce evaluation bias. Expanding prompt templates
620 and reporting robustness across prompts would
621 help mitigate this issue.

622 Ethics Statement

623 Our benchmark construction relies on multiple
624 public biomedical resources, including KEGG,
625 PubMed, PubChem, NCBI Gene, and UniProt. For
626 all resources, we strictly adhere to their respective
627 usage policies and access them solely for academic
628 research purposes. We also control API request
629 rates to avoid excessive traffic and ensure responsi-
630 ble use. In addition, we provide explicit citations
631 for all public biomedical resources used in this
632 work.

633 Large language models are used primarily to
634 support data processing tasks. The generated con-
635 tent is limited to biomedical mechanistic analysis
636 grounded in the provided literature evidence and
637 pathway context, and is not intended to produce
638 offensive, hateful, or otherwise harmful content.
639 We further incorporate validation procedures and
640 human checks to mitigate the risk of unverified or
641 misleading statements.

642 We did not employ third-party annotators. The
643 two validators involved in expert verification are
644 collaborators on this project based in China, and
645 their verification efforts constitute one component
646 of their contribution to this work. Both validators
647 are domain experts in Bioinformatics with Ph.D.
648 degrees, ensuring the professional quality of the
649 annotations.

650 This study does not involve human-subject re-
651 search and does not collect personal data. All in-
652 formation used is derived exclusively from public

biomedical databases and published scientific liter- 653
ature. 654

References 655

- 656 Ayushi Agrawal, Hasan Balci, Kristina Hanspers, Su- 657
san L Coort, Marvin Martens, Denise N Slen- 658
ter, Friederike Ehrhart, Daniela Digles, Andra 659
Waagmeester, Isabel Wassink, and 1 others. 2024. 660
WikiPathways 2024: next generation pathway 661
database. *Nucleic acids research*, 52(D1):D679– 662
D689.
- 663 Lei Bai, Zhongrui Cai, Yuhang Cao, Maosong Cao, 664
Weihan Cao, Chiyu Chen, Haojiong Chen, Kai Chen, 665
Pengcheng Chen, Ying Chen, and 1 others. 2025. 666
Intern-S1: A scientific multimodal foundation model. 667
arXiv preprint arXiv:2508.15763.
- 668 Ana R Baião, Zhaoxiang Cai, Rebecca C Poulos, 669
Phillip J Robinson, Roger R Reddel, Qing Zhong, Su- 670
sana Vinga, and Emanuel Gonçalves. 2025. A tech- 671
nical review of multi-omics data integration meth- 672
ods: from classical statistical to deep generative ap- 673
proaches. *Briefings in bioinformatics*, 26(4):bbaf355.
- 674 Chiara Balestra, Carlo Maj, Emmanuel Müller, and An- 675
dreas Mayr. 2023. Redundancy-aware unsupervised 676
ranking based on game theory: Ranking pathways in 677
collections of gene sets. *Plos one*, 18(3):e0282699.
- 678 Garth R Brown, Vichet Hem, Kenneth S Katz, Michael 679
Ovetsky, Craig Wallin, Olga Ermolaeva, Igor Tolstoy, 680
Tatiana Tatusova, Kim D Pruitt, Donna R Maglott, 681
and 1 others. 2015. Gene: a gene-centered infor- 682
mation resource at NCBI. *Nucleic acids research*, 683
43(D1):D36–D42.
- 684 UniProt Consortium. 2015. UniProt: a hub for protein 685
information. *Nucleic acids research*, 43(D1):D204– 686
D212.
- 687 John M Elizarraras, Yuxing Liao, Zhiao Shi, Qian Zhu, 688
Alexander R Pico, and Bing Zhang. 2024. We- 689
bGestalt 2024: faster gene set analysis and new sup- 690
port for metabolomics and multi-omics. *Nucleic 691
acids research*, 52(W1):W415–W421.
- 692 Yifei Ge, Feifan Zhang, Yijiang Liu, Chao Jiang, Peng 693
Gao, Nguan Soon Tan, Sai Zhang, Yuchen Shen, 694
Qianyi Zhou, Xin Zhou, and 1 others. 2025. Lever- 695
aging large language models for redundancy-aware 696
pathway analysis and deep biological interpretation. 697
bioRxiv, pages 2025–08.
- 698 Team GLM, Aohan Zeng, Bin Xu, Bowen Wang, Chen- 699
hui Zhang, Da Yin, Dan Zhang, Diego Rojas, Guanyu 700
Feng, Hanlin Zhao, and 1 others. 2024. ChatGLM: A 701
family of large language models from GLM-130B to 702
GLM-4 all tools. *arXiv preprint arXiv:2406.12793*.
- 703 Google DeepMind. [Gemini 3](#).

704	Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William Cohen, and Xinghua Lu. 2019. Pubmedqa: A dataset for biomedical research question answering. In <i>Proceedings of the 2019 conference on empirical methods in natural language processing and the 9th international joint conference on natural language processing (EMNLP-IJCNLP)</i> , pages 2567–2577.	760
705		761
706		762
707		763
708		764
709		765
710		766
711	Minoru Kanehisa, Miho Furumichi, Yoko Sato, Yuriko Matsuura, and Mari Ishiguro-Watanabe. 2025. KEGG: biological systems database as a model of the real world. <i>Nucleic acids research</i> , 53(D1):D672–D677.	767
712		768
713		769
714		770
715		771
716	Sunghwan Kim, Paul A Thiessen, Evan E Bolton, Jie Chen, Gang Fu, Asta Gindulyte, Lianyi Han, Jane He, Siqian He, Benjamin A Shoemaker, and 1 others. 2016. PubChem substance and compound databases. <i>Nucleic acids research</i> , 44(D1):D1202–D1213.	772
717		773
718		774
719		775
720		776
721	Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention. In <i>Proceedings of the 29th symposium on operating systems principles</i> , pages 611–626.	777
722		778
723		779
724		780
725		781
726		782
727		783
728	Aixin Liu, Aoxue Mei, Bangcai Lin, Bing Xue, Bingxuan Wang, Bingzheng Xu, Bochao Wu, Bowei Zhang, Chaofan Lin, Chen Dong, and 1 others. 2025. DeepSeek-v3.2: Pushing the frontier of open large language models. <i>arXiv preprint arXiv:2512.02556</i> .	784
729		785
730		786
731		787
732		788
733	Ling Luo, Po-Ting Lai, Chih-Hsuan Wei, and Zhiyong Lu. 2022. A sequence labeling framework for extracting drug-protein relations from biomedical literature. <i>Database</i> , 2022:baac058.	789
734		790
735		791
736		792
737	Marija Milacic, Deidre Beavers, Patrick Conley, Chuqiao Gong, Marc Gillespie, Johannes Griss, Robin Haw, Bijay Jassal, Lisa Matthews, Bruce May, and 1 others. 2024. The Reactome pathway knowledgebase 2024. <i>Nucleic acids research</i> , 52(D1):D672–D678.	793
738		794
739		795
740		796
741		797
742		798
743	Alex E Mohr, Carmen P Ortega-Santos, Corrie M Whisner, Judith Klein-Seetharaman, and Paniz Jasbi. 2024. Navigating challenges and opportunities in multi-omics integration for personalized healthcare. <i>Biomedicines</i> , 12(7):1496.	799
744		800
745		801
746		802
747		803
748	Sarah Mubeen, Alpha Tom Kodamullil, Martin Hofmann-Apitius, and Daniel Domingo-Fernandez. 2022. On the influence of several factors on pathway enrichment analysis. <i>Briefings in bioinformatics</i> , 23(3):bbac143.	804
749		805
750		806
751		807
752		808
753	National Library of Medicine (US). Medical subject headings (mesh) 2025 .	809
754		810
755	OpenAI. GPT-5.2 .	811
756	Ozan Ozisik, Morgane T��rezol, and Ana��s Baudot. 2022. orsum: a Python package for filtering and comparing enrichment analyses using a simple principle. <i>BMC bioinformatics</i> , 23(1):293.	
757		
758		
759		
	William G Ryan V, Smita Sahay, John Vergis, Corey Weistuch, Jarek Meller, and Robert E McCullumsmith. 2025. Pathway analysis interpretation in the multi-omic era. <i>BioTech</i> , 14(3):58.	
	Pedro H Godoy Sanches, Nicolly Clemente de Melo, Andreia M Porcari, and Lucas Miguel de Carvalho. 2024. Integrating molecular perspectives: Strategies for comprehensive multi-omics integrative data analysis and machine learning applications in transcriptomics, proteomics, and metabolomics. <i>Biology</i> , 13(11):848.	
	ScienceOne. S1-Base-671B .	
	ByteDance Seed, Jiaze Chen, Tiantian Fan, Xin Liu, Lingjun Liu, Zhiqi Lin, Mingxuan Wang, Chengyi Wang, Xiangpeng Wei, Wenyuan Xu, and 1 others. 2025. Seed1.5-thinking: Advancing superb reasoning models with reinforcement learning. <i>arXiv preprint arXiv:2504.13914</i> .	
	Mykhaylo Slobodyanyuk, Alexander T Bahcheli, Zoe P Klein, Masroor Bayati, Lisa J Strug, and J��ri Reimand. 2024. Directional integration and pathway enrichment analysis for multi-omics data. <i>Nature Communications</i> , 15(1):5690.	
	Jacob White. 2020. Pubmed 2.0. <i>Medical reference services quarterly</i> , 39(4):382–387.	
	An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, and 1 others. 2025. Qwen3 technical report. <i>arXiv preprint arXiv:2505.09388</i> .	
	Yanzhao Zhang, Mingxin Li, Dingkun Long, Xin Zhang, Huan Lin, Baosong Yang, Pengjun Xie, An Yang, Dayiheng Liu, Junyang Lin, and 1 others. 2025. Qwen3 Embedding: Advancing text embedding and reranking through foundation models. <i>arXiv preprint arXiv:2506.05176</i> .	
	Yijia Zhang, Hongfei Lin, Zhihao Yang, Jian Wang, and Yuanyuan Sun. 2019. Chemical-protein interaction extraction via contextualized word representations and multihead attention. <i>Database</i> , 2019:baz054.	
	Haiteng Zhao, Chang Ma, Lingpeng Kong, and Zhi-Hong Deng. 2025. Benchmarking and enhancing large language models for biological pathway reasoning. <i>OpenReview</i> .	
	Kangmei Zhao and Seung Yon Rhee. 2023. Interpreting omics data with pathway enrichment analysis. <i>Trends in Genetics</i> , 39(4):308–319.	
	Juexiao Zhou, Bin Zhang, Guowei Li, Xiuying Chen, Haoyang Li, Xiaopeng Xu, Siyuan Chen, Wenjia He, Chencheng Xu, Liwei Liu, and 1 others. 2024. An ai agent for fully automated multi-omic analyses. <i>Advanced Science</i> , 11(44):2407094.	

812 Appendices

813 A Human Expert Verification Protocol

814 To evaluate the reliability of our data construction
815 workflow and to establish a high-confidence gold
816 standard, we performed a human-in-the-loop veri-
817 fication study with two domain experts in molecular
818 biology and systems biology.

819 **Sampling and Materials.** We randomly sam-
820 pled 50 *pathway–paper* pairs from the constructed
821 dataset. For each pair, we provided (i) the paper
822 abstract and the intermediate artifacts produced
823 by our pipeline, including (ii) the LLM-extracted
824 mechanistic text, (iii) the LLM-based relevance
825 score with its abstract-grounded evidence and ratio-
826 nale, and (iv) the normalized entity list. For each
827 normalized entity, we additionally supplied multi-
828 ple database-derived aliases and synonyms to fa-
829 cilitate matching to the abstract. We also provided
830 the structured outputs derived from the abstract,
831 including annotated entity states and relations in
832 the resulting knowledge graph.

833 **Verification Steps and Criteria.** Experts evalu-
834 ated each instance using a four-step checklist. They
835 were instructed to use the abstract as the evidence
836 source, which matches the evidence constraints of
837 our construction pipeline:

- 838 1. **Pathway–paper relevance.** Based on the ab-
839 stract, determine whether the paper is highly
840 relevant to the specified pathway context.
- 841 2. **Hallucination check for mechanistic text.** De-
842 termine whether the LLM-extracted mechanistic
843 text contains unsupported claims, defined as
844 statements not substantiated by the abstract.
- 845 3. **Entity presence validation.** Verify whether
846 each normalized entity is present in the abstract,
847 allowing matches via the provided aliases and
848 synonyms.
- 849 4. **State and relation grounding.** Verify whether
850 the annotated molecular states and inter-entity
851 relations in the knowledge graph are consistent
852 with the abstract.

853 We labeled an instance as **reliable** only if it sat-
854 isfied all four criteria. Otherwise, it was labeled as
855 unreliable.

856 **Results.** All 50 sampled instances satisfied the
857 checklist criteria (100% pass rate). This result sup-
858 ports that our workflow produces robust intermedi-
859 ate artifacts and accurate structured supervision.

860 B Case Study: Multi-Omics Pathway 861 Mechanism Elucidation

862 **Legend.** **Entity** *Relation* **Phenotype/Process**
863 [State/Observation].

864 We present a representative end-to-end instance
865 of multi-omics pathway mechanism elucidation
866 for the KEGG pathway *Osteoclast differentiation*
867 (hsa04380), grounded in a linked mechanistic study
868 (PMID: 41266497). The instance is specified as
869 follows:

870 **Input. Pathway Name:** Osteoclast differentia-
871 tion – Homo sapiens (human).

872 **Pathway Description:** The osteoclasts, multin-
873 ucleated cells originating from the hematopoietic
874 monocyte-macrophage lineage, are responsible
875 for bone resorption. Osteoclastogenesis is mainly
876 regulated by signaling pathways activated by
877 RANK and immune receptors, whose ligands are
878 expressed on the surface of osteoblasts. Signaling
879 from RANK changes gene expression patterns
880 through transcription factors like NFATc1 and
881 characterizes the active osteoclast.

882 **Observed Entities:** Gsk-650394, SGK1, STAT3
883 (dephosphorylated), MYCL (expression reduced),
884 CTSK (transcription promoted), NFATC1.

885 **Instruction:** Based on the pathway context and
886 biological principles, generate a coherent mecha-
887 nism analysis explaining how these entities interact
888 to produce the observed phenotype.
889

890 **Output.** Mechanistically, **SGK1 phosphorylates**
891 **STAT3** at Tyr705, and Tyr705-phosphorylated
892 **STAT3 upregulates MYCL** expression. **MYCL**
893 directly **binds** the **CTSK** promoter and **promotes**
894 **CTSK** transcription, thereby supporting **osteoclast**
895 **differentiation** and **bone resorption**. Consistent
896 with this signaling axis, pharmacological inhibition
897 of **SGK1** by **GSK650394** reduces **STAT3** Tyr705
898 phosphorylation, leading to reduced **MYCL** ex-
899 pression and impaired osteoclastogenesis. Import-
900 tantly, **MYCL** overexpression can **rescue** the osteo-
901 clast differentiation defect induced by **SGK1** inhi-
902 bition, providing functional support for the causal
903

904 **SGK1**→**STAT3**→**MYCL**→**CTSK** pathway in os-
905 teoclastogenesis.

906 This instance specifies a chemical perturbation
907 (**GSK650394**) and heterogeneous molecular alter-
908 ations spanning signaling and transcriptional reg-
909 ulation, including **SGK1**, **STAT3** (dephosphory-
910 lated at Tyr705), **MYCL** (reduced expression), and
911 **CTSK** (increased transcription), together with the
912 pathway-context transcription factor **NFATC1**. It is
913 challenging because the model must integrate these
914 observations into a single, causally constrained
915 explanation that links a targeted intervention to
916 pathway-level phenotypes **osteoclast differentia-**
917 **tion** and **bone resorption**. In particular, the model
918 must (i) perform state-aware grounding to align
919 entities with explicit molecular states (e.g., Tyr705
920 phosphorylation of **STAT3**) rather than treating
921 them as state-free symbols, (ii) propagate interven-
922 tion effects with correct directionality under **SGK1**
923 inhibition, (iii) construct a multi-step causal chain
924 from upstream kinase signaling to transcriptional
925 control and phenotype rather than enumerating iso-
926 lated facts, and (iv) incorporate pathway priors such
927 as **NFATC1** as contextual constraints without as-
928serting unsupported causal edges. Although these
929 requirements reflect routine mechanistic reason-
930 ing in biology, contemporary models often fail to
931 produce end-to-end explanations that are simultane-
932 ously state-consistent, intervention-consistent, and
933 evidence-grounded from the provided perturbed
934 entities and pathway context.

C Prompts for Data Construction Workflow

Prompt for Semantic and Mechanistic Relevance Scoring

Role

You are an expert Biocurator and Molecular Biologist specializing in pathway analysis and literature mining. Your task is to evaluate the relevance between a specific KEGG Pathway and a scientific article (PubMed).

Task

Analyze the provided "KEGG Pathway Info" and "Literature Info" to determine if the article provides meaningful evidence, context, or experimental data related to the pathway.

Input Data

1. KEGG Pathway Info

```
{{KEGG_PATHWAY_JSON}}
```

2. Literature Info

```
{{PAPER_JSON}}
```

Evaluation Criteria & Steps

- 1. Subject Matching:** Does the article primarily discuss the biological process described in the pathway (e.g., Glycolysis, Gluconeogenesis)? Differentiate between core focus vs. background mention.
- 2. Species Consistency:** Check if the species in the pathway (e.g., "hsa" for Human) matches the organism model in the paper.
 - Note: If the pathway is Human but the paper uses a model organism (e.g., Mouse/Murine) to simulate human physiology/disease, this is considered **Relevant**.
 - Note: If the species are completely unrelated (e.g., Plant pathway vs. Human study), penalize the score.
- 3. Molecular Evidence:** Look for specific mentions of the genes, enzymes, or metabolites described in the pathway description.
- 4. Directionality/Context:** Does the paper discuss the activation, inhibition, or regulation of this pathway?

Scoring Standard (0-10)

- **0-1 (Irrelevant):** No meaningful connection. The terms might appear only in references or unrelated contexts.
- **2-4 (Low):** Pathway is mentioned as a keyword or broad concept, but not investigated. Major species mismatch without translational value.
- **5-7 (Medium):** The pathway is part of the results (e.g., "we observed changes in glycolysis"). Valid species match or relevant model.
- **8-10 (High):** The pathway is the central topic. The paper elucidates mechanisms, regulation, or disease implications of this specific pathway. High species alignment.

Output Format

Provide the result in a valid JSON object strictly adhering to the following structure. Do not output markdown backticks (```) or extra text.

```
{
  "relevance_score": <int, 0-10>,
  "relevance_level": "<High/Medium/Low/Irrelevant>",
  "species_check": "<Briefly state if species match or if a valid model organism is used>",
  "evidence_summary": "<Extract 1-2 key sentences from the abstract that support the link>",
  "reasoning": "<Concise explanation of the score, focusing on biological mechanisms and study focus>"
}
```

Prompt for Entity Extraction and Mechanism Summarization

Role

You are an expert Systems Biologist and Biomedical Literature Curator. Your task is to extract, categorize, and standardize key molecular information to construct a high-quality “Ground Truth” benchmark for multi-omics pathway analysis.

Goal

You will be provided with a **Target Pathway** and a scientific paper’s details. Your goal is to generate two structured outputs:

1. **Significant Entities (Classified & Normalized)**: A structured list of metabolites, genes, proteins, or phenotypes explicitly mentioned in the text. You must classify them by type and normalize them (expand abbreviations) to facilitate downstream database mapping.
2. **Mechanism Analysis**: A concise, coherent biological explanation of how these entities interact. Crucially, this must be written as a direct fact or expert interpretation, not as a summary of a study.

Instructions

• Entity Classification & Normalization

- **Chemicals/Metabolites**: Target for PubChem mapping. Provide the original text found in the abstract and a `standard_name` (full chemical name, expand abbreviations like 5-ALA to 5-aminolevulinic acid).
- **Genes/Proteins**: Target for UniProt/NCBI mapping. Provide the original text and a `standard_name` (official symbol or full protein name).
- **Processes/Phenotypes**: Target for GO/MeSH. Biological outcomes or processes (e.g., Ferroptosis, Oxidative Stress). Keep these separate as they are not valid targets for chemical/gene databases.

- **Tone & Style (Critical)**: Write the `mechanism_analysis` as an objective biological fact. Do not use phrases such as “This study reveals”, “The authors found”, or “We observed”. Start directly with the biological subject (e.g., “Elevated levels of X cause . . .”).

- **Contextualize**: Use the Pathway Description and Evidence Summary to filter for relevance.

- **Precision**: Extract only entities explicitly mentioned in the text.

Constraints

1. **No Hallucinations**: Do not invent standard names if an abbreviation is ambiguous. Use the most likely biological expansion based on context.
2. **Relevance**: Focus on the **Target Pathway**.
3. **No Meta-Language**: Strictly ban words referring to the source material in the mechanism analysis.
4. **Format**: Return strict JSON only.

Input

Target Pathway: {{pathway_name}}
Pathway Description: {{pathway_description}}
Title: {{title}}
Abstract: {{abstract}}
Evidence Summary: {{evidence_summary}}

Output JSON Format

Do not output markdown backticks (```) or extra text.

```
{
  "significant_entities": {
    "chemicals": [
      { "original": "Text in abstract", "standard_name": "Full Standardized Name" }
    ],
    "genes_proteins": [
      { "original": "Text in abstract", "standard_name": "Official Symbol/Name" }
    ],
    "processes_phenotypes": [
      "Phenotype 1",
      "Phenotype 2"
    ]
  },
  "mechanism_analysis": "A direct biological explanation describing the interaction of these entities and the resulting outcome."
}
```

Prompt for Interaction Quadruple Extraction

Role

You are an expert Biological Knowledge Graph Constructor. Your task is to convert a biological mechanism text into a structured list of interactions (quadruplets) using a strict vocabulary.

Task

Extract all biological interactions from the "Mechanism Text" and map them to the following JSON structure: {"source": "Entity A", "relation": "Relation_Type", "target": "Entity B", "condition": "Context/Prerequisite"}.

Constraints

1. **Entity Mapping:** Use the exact names provided in the "Standardized Entities List" for source and target.
 - Note the [Type] tags provided in the input (e.g., [Chemical], [Gene], [Phenotype]) to understand the biological context of each entity.
 - If an entity is missing from the list but critical for the logic, use its name from the text.
2. **Relation Vocabulary (Strict):** Use only the specific relation strings defined in the lists within the following JSON schema. Use the inline comments (//) as guidance.

```
{
  "Regulatory": [
    "activates",           // covers: activation
    "inhibits",           // covers: inhibition
    "upregulates_expression", // covers: expression
    "downregulates_expression", // covers: repression
    "regulates"           // covers: indirect effect, state change
  ],
  "Physical_Interaction": [
    "binds",              // covers: binding/association
    "dissociates_from"    // covers: dissociation
  ],
  "Modification": [
    "phosphorylates",     // covers: phosphorylation
    "dephosphorylates",  // covers: dephosphorylation
    "ubiquitinates",      // covers: ubiquitination
    "glycosylates",       // covers: glycosylation
    "methylates"          // covers: methylation
  ],
  "Metabolic": [
    "produces",           // covers: compound (enzyme -> product)
    "consumes",           // covers: metabolic_reaction (substrate -> enzyme)
    "converts_to"         // covers: metabolic_reaction (substrate -> product)
  ],
  "Causal (Phenotypic)": [
    "leads_to",           // e.g., leads to Apoptosis
    "increases_level",    // e.g., increases ROS levels
    "decreases_level"
  ]
}
```

3. **Compound Handling:** If "Gene A increases Metabolite B to activate Gene C", split it into:

- (Gene A, increases_level, Metabolite B)
- (Metabolite B, activates, Gene C)

Do not use compound as a relation.

4. **Condition Extraction:** Extract a short, specific phrase describing when, where, or how the interaction happens (e.g., "upon light activation", "in HO-1 deficient cells", "under hypoxia", "when combined with DCA"). If the interaction is a general biological fact with no specific constraint, use "General".

Input

Mechanism Text: {{mechanism_text}}

Standardized Entities List: {{formatted_entity_str}}

Output Format

Return a strict JSON object following this exact schema:

```
[
  {
    "source": "Entity A",
    "relation": "Relation_Type_From_Vocabulary",
    "target": "Entity B",
    "condition": "Context string or 'General'"
  }
]
```

Prompt for Biological State Annotation

Role

You are an expert Biological Graph Annotator.

Task

You will be provided with:

1. **Mechanism Text:** A biological description.
2. **Existing Interactions:** A list of structured interactions (quadruplets) extracted from the text.

Your job is to annotate the biological state (e.g., increased, mutated, added, phosphorylated) for the source and target entities in each interaction, based strictly on the text.

Instructions

1. For each interaction in the list, add two new fields: `source_state` and `target_state`.
2. Extract the state as a short phrase from the text.
 - Examples: "elevated levels", "administration", "mutated", "deficiency", "accumulation", "overexpression".
3. **Treatment/Drug Handling:** If an entity is a drug or treatment added to the system, use terms such as "administration", "added", or "treated with".
4. **Default:** If the text does not specify a change or state for that entity in that specific context (i.e., it is simply a pathway component), use "Present" or "Endogenous".
5. **Constraint:** Do not change the original source, relation, target, or condition values. Keep them exactly as provided.

Input

Mechanism Text: `{{mechanism_text}}`

Existing Interactions to Annotate: `{{existing_graph_json}}`

Output Format

Return a strict JSON object with a single key "annotated_kg" containing the updated list.

Example Output Schema

```
{
  "annotated_kg": [
    {
      "source": "Entity A",
      "source_state": "elevated levels",
      "relation": "increases_level",
      "target": "Entity B",
      "target_state": "accumulation",
      "condition": "General"
    }
  ]
}
```