
DeepRoot: A KG-Coordinated Multi-Agent System for Therapeutic Reasoning over Historical Medical Texts

Zijian (Carl) Ma^{*†}
Department of Bioengineering
Stanford University

Sean J. Wang[†]
Department of Bioengineering
Stanford University

Sijbren Kramer[†]
Department of Bioengineering
Stanford University

Li Erran Li
Amazon AWS AI

Abstract

Historical medical archives and traditional medicines hold immense potential for drug discovery and remain a primary source for current drug development. However, pre-ontological prose and idiosyncratic taxonomies prevent the standardization and medical modernization of the data for use in current biomedical pipelines. Furthermore, no existing LLM agent system, whether tool-calling, retrieval-augmented, or agentic deep-research, can convert such text into verifiable drug-discovery leads at scale. We close this gap with DeepRoot, a multi-agent LLM system that jointly builds and utilizes a verified knowledge graph, showing that grounding and reasoning—often conflated—are separable axes the system can compose for therapeutic reasoning. Applied to the *Shen Nong Ben Cao Jing*, DeepRoot recovers 10 of 21 held-out compound–disease treatment pairs at R@20 (47.6% vs. 4.8% for a raw corpus LLM and $\sim 2.4\%$ random) and dominates an LLM-as-judge audit for reasoning quality over baseline LLMs and LLMs with direct tool-call access to the same APIs DeepRoot itself queries. Tool-using LLMs hallucinate evidence on 87% of claims, versus 7–10% for DeepRoot. Graph-only inference hallucinates 0% but ranks lowest on reasoning coherence; DeepRoot KG + LLM is the only condition to win on both axes, pointing toward a route for systematic mining and repurposing of historical medical knowledge.

1 Introduction

Natural products—chemical compounds synthesized by living organisms—remain the leading source of approved drugs and provide scaffolds for developing more potent derivatives [Newman and Cragg, 2020, Koehn, 2012]. Many natural products have been uncovered through mining traditional medicines, including morphine from opium poppies and the antimalarial artemisinin, with the latter isolated by Tu Youyou after consulting a 4th-century Chinese medical text [Tu, 2011, Brook et al., 2017].

ML, DL, and LLM approaches for mining historical medical texts at scale have been reported before but treat the text as pure input-classification problems without a reasoning trace grounded in verified biological evidence or mechanism ontologies [Li et al., 2024, Hui et al., 2020, Liu et al.,

^{*}Corresponding author: maziujian@stanford.edu. [†]Equal contribution.

2025, Dai et al., 2024]. In parallel, multi-agent LLM systems leverage a shared knowledge graph (KG) for coordinated reasoning [Ghafarollahi and Buehler, 2025, Rasmussen et al., 2025], but only qualitatively: they neither ablate the graph against agent decomposition, nor evaluate on the regimes we target—historical clinical cases where traditional text lacks clean ontological anchors, and discovery problems with sparse ground truth.

Building on these advances, we introduce **DeepRoot** (Figure 1), a multi-agent LLM pipeline where agents collectively construct and reason over a shared KG (Neo4j). Closest to our work is OpenTCM [He et al., 2025], which uses a Graph-RAG architecture for LLM reasoning. However, its construction relied on expert oversight and pure LLM-generated outputs. **DeepRoot Assembly** agentically populates the knowledge graph via seven specialized agents that combine LLM canonicalization with strict verification against curated biomedical databases. **DeepRoot Discovery** then employs critic and discovery agents, leveraging Neo4j Cypher walks for subgraph traversal to evaluate therapeutic claims and identify potential therapeutics that are mechanistically grounded in the KG.

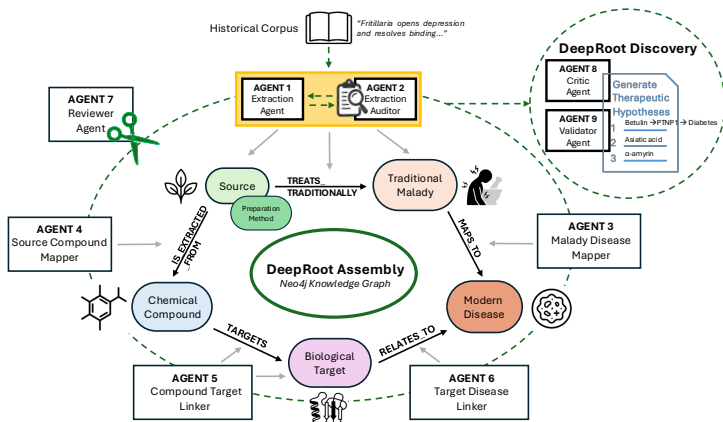


Figure 1: Schematic of DeepRoot. Graph nodes and edges are represented by rounded rectangles and black arrows. Gray arrows indicate creation of specific nodes and edges by particular agents.

2 Methods and KG construction

2.1 Dataset and grounding sources

Corpus. We evaluate on the *Shen Nong Ben Cao Jing* materia medica, segmented into 71 chunks. The corpus catalogues plants, animals, and minerals (*sources*), maladies, preparation methods, and claimed therapeutic uses.

External grounding. Every entity is verified against curated biomedical databases. Sources are linked to compounds via COCONUT2.0 [Chandrasekhar et al., 2025] (natural products) and PubChem [Kim et al., 2023] (chemicals); compounds are linked to molecular targets and clinical indications via ChEMBL [Mendez et al., 2019]; protein targets are linked to diseases via Open Targets [Ochoa et al., 2021], and pathogenic-organism targets via NCBI Taxonomy [Schoch et al., 2020] with OLS4 [McLaughlin et al., 2025]. Modern disease nodes are anchored to ICD-10, MeSH, SNOMED, MONDO, and DOID identifiers via NLM and EBI lookup services.

2.2 Knowledge graph schema

The graph has six node types and seven edge types (Figure 1). A therapeutic claim is *verifiable* when its mechanistic loop closes: a *Source* treats a *Traditional_Malady* that maps to a *Modern_Disease*; the source contains a *Chemical_Compound* that targets a *Biological_Target* which itself relates to that same *Modern_Disease*. Identity for compounds is the RDKit-computed InChIKey and identity for targets is the curated ChEMBL ID, so equivalent entities arriving from different routes collapse onto the same node. Full schema is tabulated in Table S1.

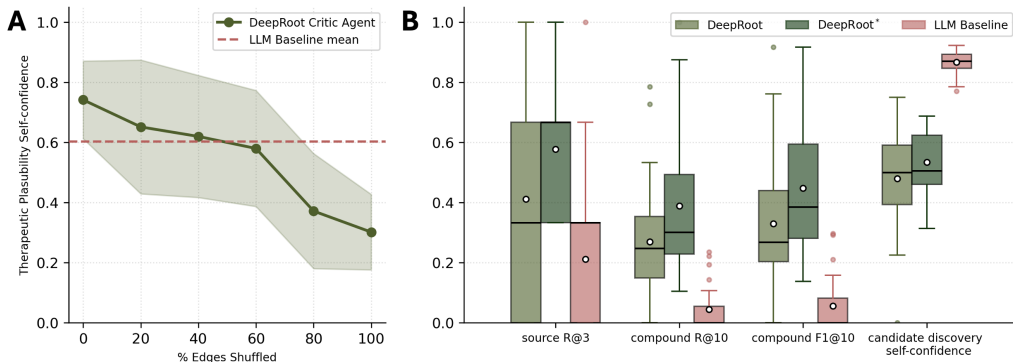


Figure 2: **(A)** Critic agent self-reported confidence for the therapeutic plausibility of each source-text pair vs. KG edge-shuffle fraction ($n=30$ source-text pairs). **(B)** Source-and-compound recovery of DeepRoot Discovery, LLM baseline (both using Gemini3.1 Flash Lite). *Batch evaluation by processing all 30 mini-corpora in a single invocation. Candidate discovery self-confidence refers to the mean self-reported confidence in each model’s proposed compound candidates per mini-corpus ($n=30$ mini-corpora).

Assembly. Seven specialized agents populate the graph in dependency order: an *extractor* emits Source, Malady, and Preparation nodes from raw text; an *auditor* canonicalizes sources and archives evidence spans that fail substring verification against their source chunk; three *linkers* ground audited entities to compounds, molecular targets, and target-to-disease associations using the databases above; a *malady-to-disease mapper* follows a generate-then-verify protocol in which LLMs propose canonical names and ontology codes are recovered only by tolerant exact match, eliminating hallucinated identifiers; and a *reviewer* archives orphans and off-domain entities.

Resulting graph. On *Shen Nong Ben Cao Jing*, Assembly yields **21,111 active nodes** (415 sources, 294 maladies, 129 modern diseases, 18,012 compounds, 2,211 targets, 50 preparations) and **52,467 active edges** (32,909 IS_EXTRACTED_FROM, 16,696 TARGETS, 1,841 RELATES_TO, 431 TREATS_TRADITIONALLY, 257 MAPS_TO, 301 KNOWN_TREATS, 32 PREPARED_AS). A visual example of nodes originating from a single extracted source is presented in Figure S1.

3 Results

3.1 Knowledge-graph ablation: edge perturbation tests structural dependence

First, to verify that DeepRoot Discovery genuinely relies on graph structure, we progressively shuffled the graph edges and tasked the critic agent with evaluating 30 extracted closed-loop source-malady claims. As expected, the Critic’s self-confidence in the therapeutic plausibility of the source based on the text decreases as edge perturbation increases, demonstrating responsiveness to the KG’s integrity (Figure 2A). Around 50% perturbation, the critic’s confidence converges with the raw LLM baseline, suggesting that the KG signal has been degraded enough that the critic behaves similarly to an LLM without structured graph support. Furthermore, past 50%, the score continues to decrease to 0.30, reflecting KG-dependent scoring.

3.2 KG-guided recovery of mechanistically supported candidates

Next, we tested whether DeepRoot can use the KG to recover mechanistically grounded candidates from noisy historical text. For this, we synthesized evaluation cases by selecting sets of 3 closed-loop and 7 non-closed-loop distractor sources. The associated paragraphs of those sources were then interweaved into a mini-corpus and fed to different models to rank the sources and candidate chemical compounds (Figure 2B). We report source recall@3, compound recall@10 (the fraction of closed-loop compounds recovered within the top-10 candidates), and mean self-confidence (0–1) related to the therapeutic plausibility of candidate compounds. Because each passage may contain ~ 500 compounds, compound recall@10 directly tests whether KG-grounded scoring concentrates the likely leads.

Over 30 mini-corpora, DeepRoot Discovery outperforms the LLM baseline, achieving $1.95\times$ higher source recall and $6.11\times$ higher compound recall (Figure 2B). Surprisingly, despite DeepRoot Discovery being theoretically capable of fully traversing the KG, recovery was not perfect. This is likely due to two factors: framing evaluation as an inference task [Edwards and Camacho-Collados, 2024] and incomplete subgraph traversal via Cypher walks, which could explain why batch evaluation across all mini-corpora improved all metrics since it could indirectly surface shared relevant paths for the Critic.

Nevertheless, the recall and F1 gains validate that KG augmentation meaningfully enhances parsing and ranking. Notably, the LLM baseline overstates therapeutic relevance, with a self-confidence of 0.87, versus DeepRoot Discovery’s 0.48, which closely aligns with the latter’s source recall@3 of 0.41 (Figure 2A). This alignment suggests that self-confidence in a KG-augmented system is effectively bounded by retrieval accuracy. In contrast, the other modalities demonstrated high self-confidence hallucinations, which is a phenomenon previously reported for both LLMs and agents [Lin et al., 2022, 2025]. Together with the KG ablation study, we establish that KG contributes meaningfully to the reasoning capabilities of an underlying LLM.

3.3 Blind rediscovery of held-out validated treatments

Whereas Section 3.2 tested whether DeepRoot surfaces *mechanistically grounded* candidates—compounds for which the graph itself closes a compound→target→disease loop—this experiment tests the complementary, harder question: can the system blind-rediscover *empirically validated* compound–disease treatments after we hide them? Concretely, for each held-out pair (a KNOWN_TREATS edge sourced from ChEMBL clinical indications) we delete the edge and all stereochemical siblings (planar-InChIKey prefix) from the validator, then ask DeepRoot Discovery to re-rank candidates for the disease. We evaluate on a 21-pair historical set, and compare against a raw-corpus LLM given the full *Shen Nong Ben Cao Jing* and asked to rank the same top- K . (Table 1)

DeepRoot Discovery recovers 10 of 21 held-out pairs ($R@20 = 47.6\%$), compared with 1 of 21 for the raw-corpus LLM baseline ($R@20 = 4.8\%$). Per-disease candidate pools span 87–1,954 compounds (median 835), so random $R@20 \approx 2.4\%$, suggesting that the result is far above random retrieval.

Table 1: Held-out KNOWN_TREATS recovery on 21 historically reachable ChEMBL indication pairs. $R@k$ in %, MRR unitless; higher is better for all metrics.

Method	R@1	R@5	R@10	R@20	MRR
DeepRoot Discovery	9.5	28.6	33.3	47.6	0.161
Raw-corpus LLM	0.0	4.8	4.8	4.8	0.012

3.4 Benchmarking DeepRoot’s therapeutic reasoning against diverse baselines

We audit critic-agent outputs with an independent LLM judge (Claude Sonnet 4.6, cross-family from the graded systems) on 30 stratified source–malady claims across seven conditions (Table 2): the DeepRoot Discovery at three LLM tiers (Gemini 3.1 Pro / 2.5 Flash / 3.1 Flash Lite), a graph-only baseline (no LLM), an LLM-only baseline given corpus passages, and a tool-call LLM with direct access to the same APIs (ChEMBL, Open Targets, PubMed, MeSH) that DeepRoot Assembly itself queries. The judge scores six dimensions on [1, 5] and flags hallucinated evidence per claim.

All three KG-augmented configurations outperform every baseline on overall score. Even DeepRoot-Lite (3.70) exceeds both the graph-only condition (3.55) and the tool-calling LLM (2.47). This contrast highlights a tradeoff between grounding and synthesis. The tool-calling agent triggers the judge’s hallucinated-evidence flag on 87% of claims, despite having on-the-fly access to the same set of APIs. By contrast, the graph-only condition produces no hallucinated evidence by construction, but exhibits the weakest reasoning coherence (2.69). KG-augmented LLMs therefore occupy a favorable middle ground: they maintain low hallucination rates (7–10%) while preserving the reasoning and synthesis capacity absent from graph-only scoring. Biomni/Phylo was evaluated using the same LLM judge suite, but its access to recently updated papers and external validation evidence may exceed the judges’ closed-book biomedical knowledge, particularly for claims grounded in post-cutoff literature.

Table 2: Reasoning-quality evaluation: seven conditions graded by Claude Sonnet 4.6 over a stratified sample of 30 source→malady claims. Scores are means on [1, 5]; Hallu. is the rate of the judge’s hallucinated_evidence flag in [0, 1]. **Bold** = best per column.

System	Components	Overall↑	EF↑	VA↑	RC↑	CM↑	UC↑	AC↑	Hallu.↓
DeepRoot — Gemini 3.1 Pro	graph + LLM	3.83	4.53	4.47	3.97	4.07	3.73	3.67	0.10
DeepRoot — Gemini 2.5 Flash	graph + LLM	3.77	4.67	4.37	3.83	3.63	3.57	3.63	0.07
DeepRoot — Gemini 3.1 Flash Lite	graph + LLM	3.70	4.60	4.27	3.73	3.70	3.60	3.67	0.07
Graph-only	graph, no LLM	3.55	4.55	4.55	2.69	3.21	3.31	2.93	0.00
Text + LLM (G3.1 FL)	corpus + LLM	3.17	3.10	2.80	3.47	3.67	3.27	3.17	0.13
Tool-call + LLM (G3.1 FL)	ChEMBL/OT/PubMed/MeSH	2.47	2.30	2.97	2.80	3.30	2.63	2.70	0.87

EF: evidence fidelity; VA: verdict alignment; RC: reasoning coherence; CM: clinical mapping; UC: uncertainty calibration; AC: actionability; Hallu.: hallucination rate.

3.5 Human-expert evaluation of DeepRoot reasoning as a qualitative case study

To evaluate the reasoning quality of DeepRoot Discovery as a traditional medicine knowledge assistant, we constructed a mini-corpus of 50 randomly sampled source–malady pairs from *Shen Nong Ben Cao Jing*, each paired with its associated textual evidence. This setting reflects a potential use case in which a scientist seeks to assess whether observed historical claims about a source’s therapeutic potential are grounded in modern biological evidence.

Table 3: Classified modern disease agreement and verdict agreement with reference to DeepRoot. Verdict refers to system-interpreted therapeutic signal of the source–malady pair: Previously Reported, Very Plausible, Plausible, Weak, and Unsupported.

System	Disease Agreement	Verdict Agreement	# Previously Reported	# Plausible or Better	# Weak	# Unsupported
DeepRoot — Gemini 3.1 Flash Lite	–	–	3	39	4	4
Biomni (w/ full agent env.)	42	30	0	32	12	6
Text + LLM (G3.1 FL)	33	13	1	11	19	19

Comparing DeepRoot with Biomni (a biomedical reasoning agent environment that is our gold-standard comparison) [Huang et al., 2025], we see good alignment of disease classification and verdicts (Table 3). We provide full responses from DeepRoot and Biomni Appendix A.10. Overall, DeepRoot Discovery clearly leverages the KG for reasoning. Biomni also shows strong reasoning and tool-calling capabilities. Notably, both DeepRoot and Biomni cite the same biological targets for their reasoning in example 1. But example 8 highlights a key limitation of KG over-reliance. Although DeepRoot’s underlying LLM identified that the hydrolyzed version of certain compounds were bioactive, those downstream products were absent from the KG, preventing KG path completion and leading to low confidence despite biological plausibility.

4 Discussion and Conclusion

DeepRoot shows that historical materia medica can be converted from pre-ontological prose into an auditable biomedical knowledge graph that supports mechanistic therapeutic reasoning. On the *Shen Nong Ben Cao Jing*, this construction pass enables held-out treatment recovery and substantially lower hallucinated-evidence rates than LLM-only, tool-calling, and biomedical-agent baselines; with broader historical corpora, the same framework could support larger-scale drug repurposing, de novo therapeutic candidate nomination, and prioritization of experimentally testable natural-product hypotheses.

It is also important to highlight the gap between DeepRoot (LLM + KG) and LLMs given direct access to the same biomedical APIs. Our results suggest that building a verified knowledge graph suppresses hallucination in a way that querying those resources at inference time does not. We show that building a verified knowledge graph suppresses hallucination in a way where querying those same APIs at inference time does not. For corpora that predate modern ontologies, retrieval-augmented and tool-using agents need a construction pass first, rather than on-demand calling. The same pattern may

transfer to other historical materia medica, including Ayurvedic and broader ethnopharmacological archives, as well as structured ranking problems beyond traditional medicine. In these settings, an agentically constructed KG offers the additional advantage that new curated claims, user submissions, and external evidence can be incorporated over time, expanding coverage while preserving traceability for future candidate ranking.

Limitations. (i) *Corpus*: a single 71-chunk materia medica; transfer to other historical corpora is unverified. (ii) *Sample size*: the held-out slice is $N = 21$ pairs, single seed, no bootstrap CIs. (iii) *Priors*: flat, face-validity, uncalibrated (Appendix A.5). (iv) *Coverage*: Open Targets is human-disease-only, leaving non-modern indications unscored. (v) *DeepRoot Discovery reasoning*: rationale quality is bounded by the underlying LLM. (vi) *Comparison evaluations with other LLM modalities*: While we reported the best result for LLM from prompt engineering, limited analysis was placed in this area, but previous studies have demonstrated that such interventions seldom provide substantial improvements [Qian et al., 2024, Wu et al., 2024].

Broader impacts. DeepRoot offers a low-cost (\sim \$0.25/corpus), auditable route to surface mechanism-grounded leads from historical medical knowledge modern pipelines ignore, and extends to other underrepresented corpora (Ayurvedic, ethnopharmacological). Outputs are leads for *in vitro* validation, not clinical advice, and commercial repurposing of communal traditional knowledge warrants attribution consideration. Historical claims may be ineffective, toxic, or culturally specific, and graph-supported plausibility does not establish safety or efficacy. Any downstream use requires expert review, provenance tracking, toxicity assessment, and experimental validation.

References

- David J. Newman and Gordon M. Cragg. Natural products as sources of new drugs over the nearly four decades from 01/1981 to 09/2019. *Journal of Natural Products*, 83(3):770–803, 2020. doi: 10.1021/acs.jnatprod.9b01285. URL <https://pubs.acs.org/doi/10.1021/acs.jnatprod.9b01285>.
- Frank E. Koehn. Biosynthetic medicinal chemistry of natural product drugs. *MedChemComm*, 3(8):854–865, 2012. doi: 10.1039/C2MD00316C. URL <https://pubs.rsc.org/en/content/articlelanding/2012/md/c2md00316c>.
- Youyou Tu. The discovery of artemisinin (Qinghaosu) and gifts from Chinese medicine. *Nature Medicine*, 17(10):1217–1220, 2011. doi: 10.1038/nm.2471. URL <https://www.nature.com/articles/nm.2471>.
- K. Brook, J. Bennett, and S. P. Desai. The chemical history of morphine: An 8000-year journey, from resin to de-novo synthesis. *Journal of Anesthesia History*, 3(2):50–55, 2017. doi: 10.1016/j.janh.2017.02.001. URL <https://www.sciencedirect.com/science/article/pii/S2352452916301293>.
- Yizhen Li, Zhongzhi Luan, Yixing Liu, Heyuan Liu, Jiaying Qi, and Dongran Han. Automated information extraction model enhancing traditional Chinese medicine RCT evidence extraction (Evi-BERT): Algorithm development and validation. *Frontiers in Artificial Intelligence*, 7:1454945, 2024. doi: 10.3389/frai.2024.1454945. URL <https://www.frontiersin.org/journals/artificial-intelligence/articles/10.3389/frai.2024.1454945/full>.
- Ye Hui, Lin Du, Shuyuan Lin, Yiqian Qu, and Dong Cao. Extraction and classification of TCM medical records based on BERT and Bi-LSTM with attention mechanism. In *2020 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*, pages 1626–1631, 2020. doi: 10.1109/BIBM49941.2020.9313359. URL <https://ieeexplore.ieee.org/document/9313359>.
- Yu Liu, Yishan Yuan, Keming Yan, Yuanyuan Li, Valeria Sacca, Sierra Hodges, Mattia Cannistra, Pauline Jeong, Jiani Wu, and Jian Kong. Evaluating the role of large language models in traditional Chinese medicine diagnosis and treatment recommendations. *npj Digital Medicine*, 8:466, 2025. doi: 10.1038/s41746-025-01845-2. URL <https://www.nature.com/articles/s41746-025-01845-2>.
- Yizheng Dai, Xin Shao, Jinlu Zhang, Yulong Chen, Qian Chen, Jie Liao, Fei Chi, Junhua Zhang, and Xiaohui Fan. TCMChat: A generative large language model for traditional Chinese medicine. *Pharmacological Research*, 210:107530, 2024. doi: 10.1016/j.phrs.2024.107530. URL <https://www.sciencedirect.com/science/article/pii/S1043661824004754>.
- Alireza Ghafarollahi and Markus J. Buehler. SciAgents: Automating scientific discovery through bioinspired multi-agent intelligent graph reasoning. *Advanced Materials*, 37(22):2413523, 2025. doi: 10.1002/adma.202413523. URL <https://doi.org/10.1002/adma.202413523>. Published online 2024.

- Preston Rasmussen, Pavlo Paliychuk, Travis Beauvais, Jack Ryan, and Daniel Chalef. Zep: A temporal knowledge graph architecture for agent memory. *arXiv preprint arXiv:2501.13956*, 2025. doi: 10.48550/arXiv.2501.13956. URL <https://arxiv.org/abs/2501.13956>.
- Jinglin He, Yunqi Guo, Lai Kwan Lam, Waikai Leung, Lixing He, Yuanan Jiang, Chi Chiu Wang, Guoliang Xing, and Hongkai Chen. Opentcm: A graphrag-empowered llm-based system for traditional chinese medicine knowledge retrieval and diagnosis, 2025. URL <https://arxiv.org/abs/2504.20118>.
- Venkata Chandrasekhar, Kohulan Rajan, Sri Ram Sagar Kanakam, Nisha Sharma, Viktor Weißenborn, Jonas Schaub, and Christoph Steinbeck. COCONUT 2.0: A comprehensive overhaul and curation of the collection of open natural products database. *Nucleic Acids Research*, 53(D1):D634–D643, 2025. doi: 10.1093/nar/gkae1063. URL <https://academic.oup.com/nar/article/53/D1/D634/7908792>. Published online 2024.
- Sunghwan Kim, Jie Chen, Tiejun Cheng, Asta Gindulyte, Jia He, Siqian He, Qingliang Li, Benjamin A. Shoemaker, Paul A. Thiessen, Bo Yu, Leonid Zaslavsky, Jian Zhang, and Evan E. Bolton. PubChem 2023 update. *Nucleic Acids Research*, 51(D1):D1373–D1380, 2023. doi: 10.1093/nar/gkac956. URL <https://academic.oup.com/nar/article/51/D1/D1373/6777787>.
- David Mendez, Anna Gaulton, A. Patrícia Bento, Jon Chambers, Marleen De Veij, Eloy Félix, María Paula Magariños, Juan F. Mosquera, Prudence Mutowo, Michał Nowotka, María Gordillo-Marañón, Fiona Hunter, Laura Junco, Grace Mugumbate, Milagros Rodriguez-Lopez, Francis Atkinson, Nicolas Bosc, Chris J. Radoux, Aldo Segura-Cabrera, Anne Hersey, and Andrew R. Leach. ChEMBL: Towards direct deposition of bioassay data. *Nucleic Acids Research*, 47(D1):D930–D940, 2019. doi: 10.1093/nar/gky1075. URL <https://academic.oup.com/nar/article/47/D1/D930/5162468>.
- David Ochoa, Andrew Hercules, Miguel Carmona, Daniel Suveges, Asier Gonzalez-Urriarte, Cinzia Malangone, Alfredo Miranda, Luca Fumis, Denise Carvalho-Silva, Michaela Spitzer, Jarrod Baker, Javier Ferrer, Arwa Raies, Olga Razuvayevskaya, Adam Faulconbridge, Evangelia Petsalaki, Prudence Mutowo, Sandra Machlitt-Northen, Gareth Peat, Elaine McAuley, Chuang Kee Ong, Edward Mountjoy, Maya Ghousaini, Andrea Pierleoni, Eliseo Papa, Miguel Pignatelli, Gautier Koscielny, Mohd Karim, Jeremy Schwartzentruber, David G. Hulcoop, Ian Dunham, and Ellen M. McDonagh. Open Targets platform: Supporting systematic drug-target identification and prioritisation. *Nucleic Acids Research*, 49(D1):D1302–D1310, 2021. doi: 10.1093/nar/gkaa1027. URL <https://academic.oup.com/nar/article/49/D1/D1302/6024045>.
- Conrad L. Schoch, Stacy Ciuffo, Mikhail Domrachev, Carol L. Hutton, Sivakumar Kannan, Rogneda Khovanskaya, Detlef Leipe, Richard McVeigh, Kathleen O’Neill, Barbara Robbertse, Shobha Sharma, Vladimir Soussov, John P. Sullivan, Lu Sun, Seán Turner, and Ilene Karsch-Mizrachi. NCBI taxonomy: A comprehensive update on curation, resources and tools. *Database*, 2020:baaa062, 2020. doi: 10.1093/database/baaa062. URL <https://academic.oup.com/database/article/doi/10.1093/database/baaa062/5881509>.
- James McLaughlin, Josh Lagrimas, Haider Iqbal, Helen Parkinson, and Henriette Harmse. OLS4: A new ontology lookup service for a growing interdisciplinary knowledge ecosystem. *Bioinformatics*, 41(5):btaf279, 2025. doi: 10.1093/bioinformatics/btaf279. URL <https://academic.oup.com/bioinformatics/article/41/5/btaf279/8125017>.
- Aleksandra Edwards and Jose Camacho-Collados. Language models for text classification: Is in-context learning enough? In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 10058–10072, Torino, Italia, 2024. ELRA and ICCL. URL <https://aclanthology.org/2024.lrec-main.879/>.
- Stephanie Lin, Jacob Hilton, and Owain Evans. TruthfulQA: Measuring how models mimic human falsehoods. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3214–3252, Dublin, Ireland, 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.acl-long.229. URL <https://aclanthology.org/2022.acl-long.229/>.
- Xixun Lin, Yucheng Ning, Jingwen Zhang, Yan Dong, Yilong Liu, Yongxuan Wu, Xiaohua Qi, Nan Sun, Yanmin Shang, Kun Wang, Pengfei Cao, Qingyue Wang, Lixin Zou, Xu Chen, Chuan Zhou, Jia Wu, Peng Zhang, Qingsong Wen, Shirui Pan, Bin Wang, Yanan Cao, Kai Chen, Songlin Hu, and Li Guo. LLM-based agents suffer from hallucinations: A survey of taxonomy, methods, and directions. *arXiv preprint arXiv:2509.18970*, 2025. doi: 10.48550/arXiv.2509.18970. URL <https://arxiv.org/abs/2509.18970>.
- Kexin Huang, Serena Zhang, Hanchen Wang, Yuanhao Qu, Yingzhou Lu, Yusuf Roohani, Ryan Li, Lin Qiu, Gavin Li, Junze Zhang, Di Yin, Shruti Marwaha, Jennefer N. Carter, Xin Zhou, Matthew Wheeler, Jonathan A. Bernstein, Mengdi Wang, Peng He, Jingtian Zhou, Michael Snyder, Le Cong, Aviv Regev, and Jure Leskovec. Biomni: A general-purpose biomedical AI agent. *bioRxiv*, June 2025. doi: 10.1101/2025.05.30.656746. URL <https://www.biorxiv.org/content/10.1101/2025.05.30.656746v1>. Preprint.

Chen Qian, Wei Liu, Hongzhang Liu, Nuo Chen, Yufan Dang, Jiahao Li, Cheng Yang, Weize Chen, Yusheng Su, Xin Cong, Juyuan Xu, Dahai Li, Zhiyuan Liu, and Maosong Sun. Chatdev: Communicative agents for software development, 2024. URL <https://arxiv.org/abs/2307.07924>.

Yiran Wu, Tianwei Yue, Shaokun Zhang, Chi Wang, and Qingyun Wu. Stateflow: Enhancing llm task-solving through state-driven workflows, 2024. URL <https://arxiv.org/abs/2403.11322>.

A Technical appendices and supplementary material

A.1 Knowledge graph schema

Edge types. Seven typed edges, each carrying a numeric `confidence_score` (flat prior, see Appendix A.5), an `evidence_type` tag (see Appendix A.6), and a `source_db` provenance field where applicable.

- **TREATS_TRADITIONALLY** (Source \rightarrow Traditional_Malady): evidence span quoted from the source chunk text.
- **MAPS_TO** (Traditional_Malady \rightarrow Modern_Disease): `is_primary`, `mapping_role` \in {primary, syndrome_component}, `mapping_source` \in {gemini+icd10_exact, gemini+mesh_exact, gemini+snomed_exact, gemini_unverified}, `mapping_alternatives` (JSON).
- **IS_EXTRACTED_FROM** (Chemical_Compound \rightarrow Source): evidence type encodes COCONUT/PubChem provenance and resolution level (canonical vs. alias vs. formula).
- **TARGETS** (Chemical_Compound \rightarrow Biological_Target): `pchembl_score`, `assay_id`, `assay_type` \in {B, F}, `assay_description`, `mechanism_action`.
- **RELATES_TO** (Biological_Target \rightarrow Modern_Disease): `ot_overall_score`, `match_tier` \in {efo_id, mondo_id, doid_id, mesh_id, norm_name}.
- **KNOWN_TREATS** (Chemical_Compound \rightarrow Modern_Disease): `clinical_phase` \in {1, 2, 3, 4}, materialized from ChEMBL drug indications; held-out evaluation slice (§3.3).
- **PREPARED_AS** (Source \rightarrow Preparation_Method).

Identity and write semantics. All writes are idempotent MERGE-on-identity. Compound identity is the RDKit-computed InChIKey, which is invariant to canonical-SMILES variants and to naming differences across COCONUT and PubChem. Target identity is the ChEMBL ID, which unifies SINGLE PROTEIN, PROTEIN COMPLEX, PROTEIN FAMILY, and ORGANISM target types under one key. Modern disease identity is the canonical name, with ontology codes coalesce-backfilled as later agents verify them against additional services.

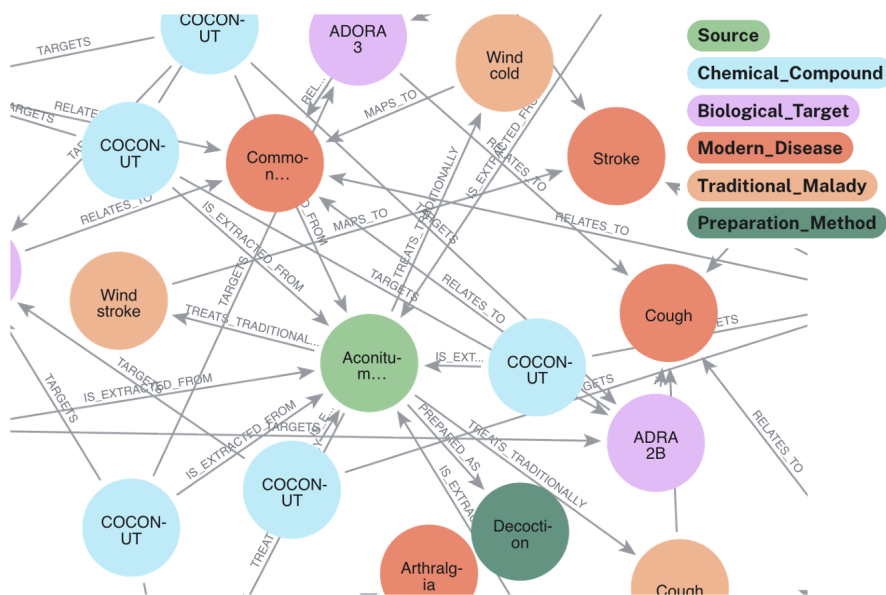


Figure S1: Example of a node cluster in Neo4j.

Table S1: Complete node schema for the DeepRoot knowledge graph.

Node	Property	Description
Source	name	Canonical Latin binomial or common name (primary key)
	aliases	Alternative names from source text
	evidence_span	Verbatim passage from which the node was extracted
	source_document	Origin corpus file identifier
	canonical_name	Auditor-resolved canonical name
	canonical_type	Taxonomic category (herb, mineral, animal, fungus, ...)
	canonical_part	Plant/animal part used (root, bark, seed, whole, ...)
	canonical_source	Database used to resolve canonical form
	canonical_raw_response	Raw LLM response from canonicalization step
	linker_status	Compound-linker outcome (ok, skipped, failed)
	linker_attempted_at	ISO timestamp of last linker run
	linker_compound_count	Number of compounds linked from this source
linker_evidence_type	Evidence type used (coconut, chEMBL)	
Chemical_Compound	name	IUPAC or common compound name (primary key)
	smiles	Canonical SMILES string
	inchikey	Standard InChIKey identifier
	molecular_formula	Molecular formula (e.g. C ₂₁ H ₂₃ N ₀₄)
	np_likeness	Natural-product likeness score from COCONUT (-5 to +5)
	annotation_level	Structural confidence tier (1 = MS2 confirmed, 5 = predicted)
	source_db	Source database (COCONUT, ChEMBL, ...)
	coconut_row	COCONUT row index for traceability
	pubchem_cid	PubChem Compound ID
	target_linker_status	ChEMBL target-linker outcome
	linker_attempted_at	ISO timestamp of target-linker run
	linker_chembl_id	ChEMBL molecule ID used for target lookup
	linker_lookup_method	Match method (inchikey, smiles, name)
	linker_target_count	Targets written after filtering
	linker_dropped_count	Targets dropped below pChEMBL floor
	linker_pchembl_floor	pChEMBL activity threshold applied
	linker_max_targets	Cap on targets written per compound
	kt_linker_status	KNOWN_TREATS linker outcome
kt_linker_attempted_at	ISO timestamp of KNOWN_TREATS linker run	
kt_linker_indication_count	Drug indications written	
kt_linker_dropped_count	Indications dropped below phase threshold	
kt_linker_min_phase	Minimum clinical trial phase accepted	
Biological_Target	name	Target protein name (primary key)
	target_pref_name	ChEMBL preferred target name
	gene_symbol	HGNC gene symbol
	uniprot_id	UniProt accession
	target_chembl_id	ChEMBL target identifier
	target_type	Target class (SINGLE PROTEIN, PROTEIN COMPLEX, ...)
	ncbi_tax_id	NCBI taxonomy ID of the target organism
	td_linker_status	Target-disease linker outcome
	td_linker_attempted_at	ISO timestamp of target-disease linker run
	td_linker_association_count	Disease associations written
	td_linker_dropped_count	Associations dropped below score threshold
	td_linker_min_score	Open Targets association score floor applied

Continued on next page

Table S1 (continued)

Node	Property	Description
Modern_Disease	name	Disease name (primary key)
	doid_id	Disease Ontology identifier
	mondo_id	MONDO disease ontology identifier
	mesh_id	MeSH descriptor
	efo_id	Experimental Factor Ontology (EFO) identifier
	icd10_code	ICD-10 classification code
	snomed_id	SNOMED CT concept identifier
	verified_by	Agent or curator that confirmed the mapping
Traditional_Malady	name	TCM ailment name (primary key)
	description	Classical definition from source text
	evidence_span	Verbatim passage supporting the malady
	source_document	Origin corpus file identifier
	mapper_status	Malady-to-disease mapper outcome
	mapper_classification	Ontology mapping confidence class
	mapper_attempted_at	ISO timestamp of mapper run
	mapper_raw_response	Raw LLM response from mapping step
	archived	Reviewer flag: duplicate or low-quality node
archive_reason	Free-text reason for archiving	
	reviewed_by	Reviewer agent identifier
Preparation_Method	name	Preparation name (decoction, pill, powder, ...)
	route	Administration route (oral, topical, inhaled, ...)
	evidence_span	Verbatim passage describing the preparation

A.2 Assembly agent protocols

Table S2: The seven DeepRoot Assembly agents in dependency order. Each agent combines LLM proposal with deterministic verification against the listed grounding source.

#	Agent	Role	Grounding source
i	Extraction	emit Source / Malady / Preparation; — TREATS, PREPARED_AS	
ii	Auditor	canonicalize; verify evidence spans; merge duplicates	COCONUT, Pub- Chem
iii	Malady→Disease	generate-then-verify ontology mapping (MAPS_TO)	MeSH, ICD-10, SNOMED
iv	Source→Compound	natural-product / chemical lookup (IS_EXTRACTED_FROM)	COCONUT, Pub- Chem
v	Compound→Target	molecular targets and bioactivities (TARGETS)	ChEMBL
vi	Target→Disease	mechanism-to-disease, dispatched on target type (RELATES_TO)	Open Targets, NCBI Tax, OLS4
vii	Reviewer	rules + LLM archival pass (orphans, OCR artifacts, off-domain entities)	—

Extraction. A Gemini call per text chunk emits Source nodes with aliases and species hints, Traditional_Malady nodes with descriptions, Preparation_Method nodes, and the TREATS_TRADITIONALLY and PREPARED_AS edges among them. Each edge records the literal evidence span from the source chunk. Confidence is the LLM’s self-assessed score, used only as a soft signal for downstream auditing (the auditor verifies span and identity independently).

Auditor. Three deterministic post-extraction jobs. (1) *Canonicalization.* Gemini Flash-Lite (temperature 0, batched 20 sources/call, structured-JSON schema) labels each source organism (Latin binomial + part), chemical (name or formula), or uncanonicalized. Per-type external lookup verifies: organisms hit a local COCONUT inverted index over 62,792 species keys, chemicals hit PubChem REST. Canonical labels with external-DB hits are tagged `gemini+coconut` / `gemini+pubchem`; the rest are tagged `gemini_unverified_*`, `uncanonicalized`, or `error`. (2) *Source merge.* Sources sharing the same (`canonical_name`, `canonical_part`) collapse onto a keeper (highest-degree, alphabetical tiebreak). Merged-from nodes are soft-archived with reason `merged_into:<keeper>`; outgoing edges are re-targeted (parallel edges take maximum confidence) and aliases are unioned. (3) *Evidence-span verification.* Substring check (whitespace-normalized) of every TREATS_TRADITIONALLY evidence span against its source chunk. Hallucinated spans (typically LLM-introduced ellipsis) trigger soft-archival with reason `hallucinated_evidence`.

Malady→Disease (generate-then-verify). One Gemini call per malady (temperature 0, six-shot system prompt) emits a typed exit: `disease`, `symptom`, `syndrome`, `ambiguous`, or `tcm_no_equivalent`. Crucially, the LLM emits canonical *names only*—never codes—along with an ontology hint. Each proposed name is then verified in parallel against ICD-10 (NLM Clinical Tables), MeSH (NLM RDF Lookup), and SNOMED (EBI OLS4). Verification accepts only *tolerant exact match* (case- and punctuation-insensitive); fuzzy hits, even from the API itself, are rejected. Syndromes can produce one primary plus up to two `syndrome_component` edges, but only components that pass exact-match verification are written (unverified components are dropped and logged in the primary edge’s `mapping_alternatives`). Default mode rejects unverified mappings entirely; an `-allow-unverified` flag stores them with `requires_review=true` and a degraded prior.

Source→Compound. Routes by `canonical_type`. Organisms hit a local COCONUT inverted index (in-memory, ~700,000 structures with species provenance) by exact normalized species name and the “first two tokens” (genus species) prefix. Chemicals hit a disk-cached PubChem REST client (4 RPS), with formula fallback via `compound/fastformula` for formula-shaped canonicals. Compound identity is the RDKit-computed InChIKey from canonical SMILES, which collapses canonicalization variants and unifies COCONUT/PubChem name splits onto a single node. Edges are written with `evidence_type` encoding the resolution path (`coconut_organism_canonical`, `pubchem_chemical_canonical`, `alias`, `formula`, `unverified`). Per-source `linker_status` stamping makes the agent fully resumable across network interruptions.

Compound→Target. ChEMBL queries by InChIKey (91.4% resolution rate), falling back to canonical SMILES then preferred name. For each resolved compound, the agent retrieves mechanism-of-action records and bioactivity records. Activity filtering: `pchembl_value` ≥ 5.0 (10 μM floor) for quantitative tier, `assay_type` $\in \{B, F\}$, `standard_relation` $\in \{=, \sim\}$; data validity flags must be empty. Target type is *not* restricted to SINGLE PROTEIN: PROTEIN COMPLEX (subunit fan-out), PROTEIN FAMILY (broad-spectrum inhibitors), and ORGANISM (anti-pathogen evidence, e.g. *Plasmodium falciparum* for antimalarials) are all admitted. A `-include-phenotypic` flag additionally retrieves phenotypic activities (`pchembl_value` null, `assay_type` F/B) at a lower confidence prior; for terpenes, sterols, and other natural products tested phenotypically rather than against named molecular targets, this opt-in is necessary to avoid silent loss of ~2,000 compounds. Salt and tautomer parents are aggregated via the ChEMBL molecule hierarchy; without hierarchy expansion, 30–50% of activities are missed for multi-form compounds.

Target→Disease. Four-way dispatch on `target_type`. SINGLE PROTEIN: Open Targets GraphQL via UniProt→Ensembl, returning disease associations with overall scores binned to confidence tiers. PROTEIN COMPLEX: subunit fan-out via ChEMBL, then per-subunit Open Targets, with per-disease max-score deduplication. PROTEIN FAMILY: intentionally skipped (`td_linker_status = skipped_protein_family`), as family-level evidence is too coarse for clinical association. ORGANISM: an EFO/DOID walk over OLS4 starting from the NCBI `tax_id`, emitting both the specific disease class and its ancestors, falling back to a parallelized LLM safety net (Gemini proposes candidate disease names, NLM MeSH verifies via tolerant exact match) when the ontology walk returns empty. Three-tier disease matching: exact ontology-ID match first, then normalized-name exact match, then MeSH-synonym expansion. A plan-then-apply phase performs all writes via UNWIND-batched Cypher transactions (~12 transactions for ~5,500 rows), with DELETE restricted to terminal-status rows so transient API failures cannot wipe live edges.

Reviewer. Two-pass deterministic-then-LLM archival. Pass 1 catches OCR artifacts (single-character entities, mojibake), orphans (degree 0), generic categories (“herb”, “compound”), and metaphysical concepts that escaped extraction. Pass 2 batches the residual ambiguous nodes (~20) to Gemini for biomedical-relevance filtering. Cascade archival propagates to incident edges. Archival is soft (archived=true with reason); no records are deleted.

A.3 Discovery agent protocols

DeepRoot Discovery comprises two agent roles operating over the typed graph: a *critic agent* that scores existing Source→Malady claims, and a *discovery agent* that nominates novel compound candidates for a target Modern_Disease. Both consume the tier-bucketed path-scoring layer described in Appendix A.4.

Critic agent. For each (Source, Malady) claim, the agent receives a structured payload assembled from the KG: the claim itself (Source name + aliases, Malady description, primary mapped Modern_Disease, mapping rationale), deterministic Pass-1 signals (path bucket distribution, loop-closure counts, top-bucket score), the top- K mechanistic chains (Source→Compound→Target→Disease) with edge metadata, and four cross-cutting enrichments—compound profiles (KNOWN_TREATS for other diseases, target spectrum), target genericity (number of associated diseases per target), source-level target convergence (multi-compound hits on the same target), and sibling verdicts (Pass-1 verdicts of other claims on the same source). The model returns a structured JSON CriticVerdict with: a verdict on the four-tier ladder (VALIDATED / PLAUSIBLE / WEAK / UNSUPPORTED); biological_plausibility and evidence_coherence scores in $[0, 1]$, defensively clamped to that range; a key_evidence list of cited compound–target–disease triples; a concerns list with typed enum values (generic_target, weak_evidence_only, indirect_mechanism, wrong_disease_mapping, syndrome_underutilized, promiscuous_compound, unverified_evidence); a free-form rationale; and a requires_human_review flag (auto-set when the LLM and deterministic Pass-1 verdicts disagree by ≥ 2 rungs). The prompt instructs the model to quote specific input fields and never speculate beyond the provided evidence; numeric ranges are enforced via post-hoc clamping rather than relying on the model to obey them.

Discovery (nominator) agent. Given a target Modern_Disease query d^* , the agent walks the KG backward (disease ← malady ← source ← compound) to enumerate all corpus-supported candidate compounds, then walks forward (compound → target → disease) to score each candidate’s mechanistic plausibility. A novelty filter drops compounds whose KNOWN_TREATS edge already reaches d^* (supporting an in-memory mask for held-out evaluation without mutating the graph). The remaining candidates are ranked lexicographically by (i) has_loop_closure (does at least one forward chain reach d^*), (ii) forward_bucket (T1 > T2 > T3 > T4; weakest-link tier of the strongest loop-closing path), (iii) unique_sources_count, (iv) unique_maladies_count, and (v) forward_max_score (multiplicative product of edge confidences along the strongest path). The output is an ordered list of CompoundCard entries containing top historical paths (source, malady, evidence span), top forward paths (target, assay description, OT score), and KNOWN_TREATS for other diseases as polypharmacology context. The discovery agent is fully deterministic—no LLM is in the loop—making the ranking auditable and stable across re-runs.

A.4 Tier-bucket path scoring

A *path* is a sequence of typed edges connecting a source node to a disease node through compound and target intermediaries. Each edge carries an evidence_type tag (e.g., chembl_mechanism, ot_association_strong, coconut_organism_canonical) which maps to one of four tiers $T \in \{T1, T2, T3, T4\}$ (Appendix A.5, Table S3) and a flat numeric prior $c \in [0, 1]$.

For a path p with edges e_1, \dots, e_n , the *path bucket* $B(p)$ and *path score* $S(p)$ are

$$B(p) = \min_{i=1\dots n} T(e_i), \quad S(p) = \prod_{i=1\dots n} c(e_i).$$

Paths are ordered lexicographically by $(B(p), S(p))$ with the bucket as the primary key (highest-tier bucket wins) and the multiplicative score as tiebreak within a bucket. The bucket captures the

qualitative claim “a chain is only as strong as its weakest edge” (weakest-link), while the score gives a continuous ordering inside each tier.

The same scoring layer is consumed by both Discovery agents (Appendix A.3) and by the deterministic Pass-1 signals fed to the critic. Because priors are flat (Table S3) rather than learned, raw external scores (`ot_overall_score`, `pchembl_value`, `np_likeness`) are preserved as edge attributes so downstream consumers can recalibrate without re-running Assembly.

A.5 Confidence priors

Tier ladder (used for path scoring). $T1 > T2 > T3 > T4$. Path bucket is the minimum tier across edges (weakest-link); within a bucket, ranking uses the multiplicative product of edge confidences as tiebreak.

Table S3: Per-edge-type confidence priors. Priors are flat (not learned), chosen on biomedical face validity, and never replaced by self-reported LLM confidence. Raw external scores (e.g., `ot_overall_score`, `pchembl_value`, `np_likeness`) are preserved as edge properties so downstream consumers can recalibrate without re-running Assembly.

Edge	Evidence type	Tier (prior)
IS_EXTRACTED_FROM	coconut_organism_canonical / pubchem_chemical_canonical	T1 (0.70–0.80)
	coconut_organism_alias	T2 (0.55)
	coconut_organism_unverified / pubchem_chemical_unverified	T3 (0.50–0.55)
	pubchem_chemical_formula	T4 (0.50)
TARGETS	chembl_mechanism	T1 (0.95)
	chembl_activity_strong ($pchembl \geq 7$)	T2 (0.75)
	chembl_activity_moderate ($pchembl \geq 6$)	T3 (0.60)
	chembl_activity_weak ($pchembl \geq 5$) / chembl_phenotypic	T4 (0.40)
RELATES_TO	ncbi_pathogen_consensus	T1 (0.92)
	ot_association_strong ($OT \geq 0.7$)	T1 (0.85)
	ot_association_moderate ($OT \geq 0.4$) / complex_aggregate / pathogen_llm_verified	T2 (0.65–0.75)
	ot_association_weak ($OT \geq 0.2$)	T3 (0.45)
MAPS_TO	icd10/mesh/snomed exact (primary)	T1 (0.80–0.85)
	syndrome_component / symptom	T2 (0.65–0.75)
KNOWN_TREATS	clinical_phase = 4 (approved)	T1 (0.95)
	clinical_phase = 3	T1 (0.85)
	clinical_phase = 2	T2 (0.65)
	clinical_phase = 1	T3 (0.45)

A.6 Graph statistics

Node and edge totals (active, post-Assembly). 21,111 nodes active, 94 archived. By type: 415 Source, 294 Traditional_Malady, 129 Modern_Disease, 18,012 Chemical_Compound, 2,211 Biological_Target, 50 Preparation_Method. Edges: 52,467 active. By type: 32,909 IS_EXTRACTED_FROM, 16,696 TARGETS, 1,841 RELATES_TO, 431 TREATS_TRADITIONALLY, 301 KNOWN_TREATS, 257 MAPS_TO (208 primary +49 syndrome_component), 32 PREPARED_AS.

Per-evidence-type breakdown. IS_EXTRACTED_FROM: 32,885 organism_canonical, 21 chemical_canonical, 3 formula. TARGETS: 60 mechanism, 1,148 strong, 1,185 moderate, 2,264 weak, 12,039 phenotypic. RELATES_TO: 936 ot_weak, 666 ot_moderate, 14 ot_strong, 203 com-

plex_aggregate, 17 llm_verified, 3 pathogen_consensus, 2 efo. KNOWN_TREATS: 60 phase 4, 92 phase 3, 84 phase 2, 65 phase 1. The phenotypic-heavy distribution of TARGETS reflects the natural-product corpus: terpenes, sterols, and flavonoids are predominantly characterized by phenotypic bioassays rather than named molecular targets.

Convergence. 4,605 compounds appear in ≥ 2 sources (classic phytochemistry pattern: β -sitosterol 112 \times , quercetin 87 \times , kaempferol 66 \times). Across the 129 Modern_Disease nodes, 257 MAPS_TO edges resolve to an average of 1.99 maladies per disease, indicating strong canonical convergence rather than fragmentation.

Coverage. 504 of 2,211 targets (22.8%) link to at least one disease; 88 of 129 disease nodes (68%) are reached by at least one target. 3,221 of 18,012 compounds (17.9%) have at least one TARGETS edge; the remainder either lack ChEMBL records or have no admissible target-class data, a documented limitation of the underlying databases rather than of the pipeline.

A.7 Evaluation protocols

Eval 1: edge-perturbation sensitivity. A fixed test set of closed-loop Source–Malady claims is sampled from the KG. For each perturbation level $p \in \{0\%, 20\%, 40\%, 60\%, 80\%, 100\%\}$, a fraction p of edges across four mechanistic edge types (TARGETS, RELATES_TO, KNOWN_TREATS, MAPS_TO) is selected uniformly at random and their target endpoints are shuffled among themselves. A single perturbation is applied per level; all test claims are then evaluated by the Critic against the same perturbed graph. Self-confidence (mean biological plausibility over all claims) is reported per level.

Eval 2: source and compound recovery on mini-corpora. To build each mini-corpus, build_recovery_eval_corpus.py first queries the knowledge graph for two disjoint source pools: **closed-loop sources** (those with at least one complete Source \rightarrow Compound \rightarrow Target \rightarrow Disease chain where the source also TREATS_TRADITIONALLY \rightarrow Malady \rightarrow MAPS_TO the same disease) and **distractor sources** (all other non-archived KG sources). The full *Shen Nong Ben Cao Jing* text is split into paragraphs, each tagged by keyword regex against every source name in the KG. Gemini Flash then verifies which tagged sources a paragraph actually *describes therapeutically* (rather than merely cross-referencing in a compatibility list). Verified single-source paragraphs are banked into the two pools. Each synthetic eval case is assembled by a diversity-maximising greedy algorithm: it picks $K=3$ least-used closed-loop source paragraphs and $N=7$ least-used distractor source paragraphs, ensures no overlap between the two sets, then deterministically shuffles all 10 paragraphs into an interleaved mini-corpus. The label set for compound recovery includes both **closed-loop compounds** (retrieved via a Cypher walk confirming Compound \rightarrow TARGETS \rightarrow Target \rightarrow RELATES_TO \rightarrow Disease for a disease the source already treats; typically 1–5 per source) and **distractor compounds** (IS_EXTRACTED_FROM compounds of the distractor sources). Compound recall@ k is reported separately against each label set so the closed-loop and broad-coverage signals can be distinguished. Thirty such mini-corpora are generated, each with a distinct closed-loop source signature enforced by deduplication.

Eval 3: positive-control recovery of hidden known treatments. The 301 KNOWN_TREATS edges are filtered to those whose Modern_Disease has a backward chain $d \leftarrow \text{malady} \leftarrow \text{source} \leftarrow \text{compound}$ in the KG (the *historical-reachability* subset), yielding 21 (compound, disease) pairs across 10 diseases. For each test pair (c^*, d^*) , we compute c^* 's planar InChIKey prefix (first 14 characters, dropping stereochemistry) and mask every KNOWN_TREATS edge from any compound sharing that prefix to d^* (in-memory mask only; the graph is not mutated). The discovery agent is then run on d^* with top- $K = 20$. A trial succeeds at rank r if any nominee within the top r shares c^* 's planar prefix. The candidate pool—all compounds reachable via the backward chain from d^* —is recorded per-disease (range 87–1,954, median 835), giving a uniform random recall@20 baseline of $\approx 2.4\%$.

Eval 4: LLM-as-judge reasoning quality. A stratified sample of 30 closed-loop Source–Malady claims is drawn from the 431 candidate claims, with the strata chosen to exercise distinct verdict regimes (representative balance of *unsupported*, *strong_support*, *T1-bucket-without-loop-closure*, *mechanistic-only*, and *traditional-only*). The same 30 claims are scored by six conditions: DeepRoot Discovery at three LLM tiers (Gemini 3.1 Pro, 2.5 Flash, 3.1 Flash-Lite), a graph-only deterministic

baseline (Pass-1 verdict only, no LLM), an LLM baseline given just the corpus passages, and a tool-call LLM baseline given direct API access to ChEMBL, Open Targets, PubMed, and MeSH. Outputs are graded by Claude Sonnet 4.6 (cross-family from the graded systems) on six dimensions in [1, 5]: Evidence Fidelity (does the critic cite evidence present in the payload?), Verdict Alignment (does the verdict follow from the visible evidence?), Reasoning Coherence (does the rationale explain *this* claim’s chain?), Clinical Mapping (does the critic responsibly handle the malady→disease mapping?), Uncertainty Calibration (are the scores, concerns, and review flag calibrated?), and Actionability (would a curator know what to inspect next?). The judge additionally returns six binary flags (hallucinated_evidence, unsupported_verdict_jump, ignored_loop_closure_status, overclaims_strength, contradictory_scores, needs_human_review) and a recommended status in {pass, weak_pass, fail, human_review}. The judge sees only the critic’s visible artifacts (verdict, scores, key_evidence, concerns, rationale) plus, where applicable, the structured payload that the critic was given; it does not have access to ground truth and grades the quality of the critic’s argument rather than its absolute correctness.

A.8 Prompt templates

This section summarizes the four critical prompts in DeepRoot. Each template is described as a tuple of (role, input schema, output schema, key instructions); the verbatim text is in the released codebase.

Extraction prompt. *Role:* extract typed entities from a single text chunk of the source corpus. *Input:* the chunk text plus a source_document identifier. *Output schema (structured JSON):* sources [], maladies [], preparations [], treats_edges [], prepared_as_edges []; each entity carries name, aliases, evidence_span (verbatim quote from the chunk), and a self-assessed confidence. *Key instructions:* evidence spans must be substring-matchable to the chunk text; identifiers (binomials, ChEMBL IDs) are never to be invented; the LLM emits names, not codes.

Malady→Disease mapping prompt (generate-then-verify). *Role:* classify each Traditional_Malady into one of {disease, symptom, syndrome, ambiguous, tcm_no_equivalent} and propose canonical English name(s). *Input:* the malady’s name, description, evidence_span, and source classical context. *Output schema:* a top-level classification field plus a mappings [] list, where each mapping carries name (canonical English), an ontology_hint ∈ {icd10, mesh, snomed}, a role ∈ {primary, syndrome_component}, and a one-sentence rationale. *Key instructions:* the model emits canonical names only, never codes; for syndrome maladies, up to two syndrome_component entries may be returned in addition to the primary; six in-context examples cover the typical TCM patterns (wind heat, gu toxin, counterflow, etc.). Codes (ICD-10, MeSH, SNOMED) are recovered downstream by deterministic exact-match verification against the corresponding ontology services—never trusted from the LLM.

Critic agent prompt. *Role:* biomedical reasoning expert evaluating whether a historical Source plausibly treats a Modern_Disease via known mechanisms. *Input:* the structured payload described in Appendix A.3 (claim, Pass-1 signals, evidence paths, compound profiles, target profiles, source-level target convergence, sibling verdicts), serialized as a single JSON object. *Output schema:* the CriticVerdict object described in Appendix A.3. *Key instructions (highlighted in the system prompt):* (i) assess TARGET QUALITY via the per-target disease count (pleiotropy of ≥ 50 associated diseases is flagged generic_target); (ii) ground COMPOUND PHARMACOLOGY in the compound’s KNOWN_TREATS record and target spectrum; (iii) scrutinize the DISEASE MAPPING for clinical plausibility against the historical evidence span; (iv) upweight POLYPHARMACOLOGY / CONVERGENCE when multiple compounds from the same source hit the same target; (v) discount SPECIFICITY when the source has many unrelated claims (kitchen-sink remedy). The prompt instructs the model to quote specific input fields and never speculate beyond the provided payload.

Baseline LLM prompt. *Role:* Candidate nominator. You are a pharmaceutical and natural products expert with broad knowledge of traditional Chinese medicine, pharmacognosy, and bioactive plant compounds. You will be given a passage from the Shen Nong Ben Cao Jing and will be asked to discover and rank plausibility of therapeutic compounds. Read this passage from a historical Chinese herbal text. List each medicinal source discussed; for each, propose up to 10 therapeutic chemical compounds it likely contains, with a plausibility score 0.0-1.0 (0.0 = incoherent or no

known mechanism, 1.0 = biologically obvious) and a one-sentence reasoning grounded in known biochemistry or pharmacology.

Biomni prompt. *Role:* Biomedical reasoning expert. *Input:* You are a pharmaceutical and natural-products expert with broad knowledge of traditional Chinese medicine, pharmacognosy, bioactive plant/mineral/animal compounds, and modern clinical pharmacology. You will be given an extracted corpus of historical medicinal entries. Each block is one claim: a SOURCE (a medicinal substance) and a MALADY (the traditional ailment the historical text says it treats), followed by the original passage. You must judge to your fullest capability, using any databases and subagents available necessitating @ChemBL, @PubChem, @PubMed, @OpenTargets as places for you to find relevant compounds or links to mechanistically evaluate the plausibility that the source provided can treat the malady listed. Rules: Score each claim independently; do not let one claim bias another. Do not invent citations or database IDs. Reason from known pharmacology. If the source is an inert mineral or has no plausible bioactive route to the disease, say so and mark Unsupported.

LLM-judge prompt. *Role:* independent biomedical evaluation judge grading the visible reasoning of an automated critic, not deciding whether the underlying therapeutic claim is true. *Input:* condition under judgment (one of the six in Eval 4), the claim, the Pass-1 deterministic signals, the visible payload the critic received (or empty for graph-free conditions), and the critic’s full structured output. *Output schema:* per-dimension scores in [1, 5] for the six rubric dimensions (Appendix A.7), six binary flags, a recommended status, and a free-text justification referencing specific input fields. *Key instructions:* treat Pass-1 as an input signal, not ground truth (a critic can be good even when it agrees with an imperfect Pass-1 verdict, if it explains the limitation correctly); penalize hallucinated citations (every compound, target, and disease named by the critic must trace to a payload field, except for widely accepted biomedical facts in text-only conditions); penalize rationales that confuse global top-bucket evidence with loop-closing disease support (a claim with 76 T1-bucket paths but zero of those paths reaching the mapped disease is *not* well-supported).

A.9 Implementation notes

Models. Extraction, Auditor canonicalization, Malady mapper, and Reviewer Pass 2 use Gemini 3.1-Flash-Lite at temperature 0 with structured-JSON schemas (rate-limited 1 request/sec, retried up to 4 times on transient failure). The Discovery LLM stages use Gemini Pro under the same rate-limit envelope.

Storage. Neo4j AuraDB. All writes are idempotent MERGE statements with status fields (e.g., linker_status, td_linker_status, mapper_status) written *last* so transient failures retry safely. Default re-run filters select only nodes with status IS NULL; -retry-misses reprocesses error and no_* branches; -force-relink deletes prior edges from terminal-status branches only, never from in-flight nodes. This makes the pipeline crash-safe end-to-end.

Caching. PubChem responses are disk-cached by query; COCONUT is loaded once into an in-memory inverted index. Re-runs after network interruptions complete in seconds.

End-to-end cost. Full Assembly over the 71-chunk corpus runs in ~30–40 minutes wall-clock, dominated by per-edge Neo4j MERGE latency on AuraDB. Total LLM spend is ~\$0.25 (Flash-Lite); ChEMBL, COCONUT, PubChem, Open Targets, NLM, and EBI are all free.

A.10 Qualitative scoring of representative critic agent responses

We assembled fifty source-malady pairs, constrained so that every pair has both a traditional TREATS_TRADITIONALLY edge and a primary mapping to a modern disease. For each pair we extracted the corresponding historical passage from the *Shen Nong Ben Cao Jing* to form a mini-corpus. The mini-corpus was given to DeepRoot Discovery, Biomni, or baseline LLM and tasked to score each pair based on the plausibility that the particular source can treat the traditional malady it is claimed to treat. Shown below are the first 10 examples for DeepRoot and Biomni.

Table S4: DeepRoot versus Biomni (baseline): malady→modern-disease mappings and plausibility verdicts across ten (Source, Malady) pairs.

#	Source → Malady	DeepRoot		Biomni	
		Mapped disease	Verdict	Mapped disease	Verdict
1	<i>Epimedium</i> → impotence	Erectile Dysfunction	Very Plausible	Erectile Dysfunction	Very Plausible
2	<i>Trichosanthes cucumeroides</i> → menstrual block	Amenorrhea	Unsupported	Amenorrhea	Plausible
3	<i>Rhizoma Arisaematis</i> → heart pain	Angina Pectoris	Previously Reported	Angina Pectoris	Weak
4	<i>Calcareaous Spar</i> → generalized fever	Fever	Unsupported	Pyrexia (Fever)	Unsupported
5	<i>Aster tataricus</i> → cough	Cough	Plausible	Bronchitis; Asthma	Very Plausible
6	<i>Fructus Ailanthi Altissimi</i> → impotence	Erectile Dysfunction	Plausible	Erectile Dysfunction	Weak
7	<i>Ge Gen</i> → toxicity	Poisoning	Plausible	Alcohol Poisoning	Plausible
8	<i>Sanguisorba officinalis</i> → vaginal discharge	Leukorrhea	Unsupported	Leukorrhea; Candidiasis	Plausible
9	<i>Hedgehog Pelt</i> → hemorrhoids	Hemorrhoids	Plausible	Hemorrhoidal Disease	Unsupported
10	<i>Calcareaous Spar</i> → cough	Cough	Unsupported	Cough; Bronchitis	Unsupported

Each box below shows one verdict for a (Source, Malady) pair. **VERY PLAUSIBLE** **PREVIOUSLY REPORTED** **PLAUSIBLE** **WEAK** **UNSUPPORTED** badges reflect the reported plausibility tier (green = strongest; PREVIOUSLY REPORTED denotes a source compound with an already-known therapeutic edge to the mapped disease).

Example 1: *Epimedium* (Source) and *Impotence* (Malady)

DeepRoot Verdict: **VERY PLAUSIBLE**

Weakest Link: The KG contains several non-closing paths for icariin that distract from the primary, well-established PDE5 inhibitory mechanism.

Reasoning: Epimedium is well-documented for its primary active constituent, icariin, a known phosphodiesterase-5 (PDE5) inhibitor. The KG confirms the icariin→PDE5A→erectile dysfunction path, which is the canonical mechanism for treating impotence by enhancing nitric oxide-mediated vasodilation. Although other paths in the graph do not close the loop, the primary mechanism is highly robust and clinically validated.

Biomni Verdict: **VERY PLAUSIBLE**

Weakest Link: Icariin's PDE5A IC₅₀ (~5,900 nM) is ~1,000-fold weaker than sildenafil (~3.5 nM); oral bioavailability of the native glycoside in humans remains the key uncertainty.

Reasoning: Icariin (CID 5318997; ChEMBL553204) is a confirmed PDE5A inhibitor (IC₅₀ = 5,900 nM, ChEMBL; active in PubChem AID 362785 — human recombinant PDE5A1 assay), sharing the same mechanism as sildenafil: cGMP preservation → corpus cavernosum smooth muscle relaxation. In vivo rat models confirm efficacy via GRK2 downregulation in penile tissue (PMID 39905744) and protection of cavernous tissue under hyperglycemia (PMID 40161546). The metabolite icariside II is significantly more potent. Icariin has reached ChEMBL Max Phase 3, making this the best-validated claim in the dataset.

Example 2: *Trichosanthes cucumeroides* (Source) and *Menstrual block* (Malady)

DeepRoot Verdict: **UNSUPPORTED**

Weakest Link: The lack of any mechanistic path connecting cucurbitacins to reproductive hormonal regulation or uterine function.

Reasoning: *Trichosanthes cucumeroides* contains cucurbitacins, which are potent cytotoxic and anti-inflammatory triterpenoids. While these compounds modulate pathways like IGF2BP1 and ITGAL, there is no clear mechanistic link to the regulation of the menstrual cycle or the treatment of amenorrhea. The KG paths converge on unrelated conditions like diabetes and malaria, failing to provide a coherent basis for this traditional claim.

Biomni Verdict: PLAUSIBLE

Weakest Link: The primary active agent, β -trichosanthin, is a large protein (~28 kDa) with an abortifacient rather than emmenagogue mechanism — it does not address functional amenorrhea unrelated to pregnancy.

Reasoning: *Trichosanthes* root contains β -trichosanthin, a Type I ribosome-inactivating protein (RIP) that depurinates 28S rRNA, arresting protein synthesis in trophoblast cells and inducing abortion in vivo (PMID 3596896; 198 PubMed results for ‘trichosanthin ribosome inactivating’). This coherently explains the TCM use for pregnancy-related amenorrhea. Cucurbitacin B (CID 5281316; ChEMBL2180736) contributes cytotoxic activity (IDH1 IC₅₀ = 640 nM; 79/122 active PubChem assays). The mechanism is abortifacient rather than broadly emmenagogue, limiting applicability to non-pregnancy amenorrhea.

Example 3: *Rhizoma Arisaematis* (Source) and heart pain (Malady)

DeepRoot Verdict: PREVIOUSLY REPORTED

Weakest Link: The evidence tier for the adenosine-receptor interactions is limited to bronze, suggesting a need for higher-confidence binding data.

Reasoning: *Rhizoma Arisaematis* contains adenosine, which acts on adenosine receptors (ADORA1, ADORA2A, ADORA2B, ADORA3) to modulate cardiac rhythm and coronary blood flow. These paths close the loop to angina pectoris, providing a plausible mechanism for alleviating heart pain through vasodilation and metabolic protection. The evidence is consistent across multiple receptor subtypes, though the overall tier remains bronze.

Biomni Verdict: WEAK

Weakest Link: *Arisaema* species (the actual botanical source) have no documented cardiac pharmacology in any modern database; the claim’s plausibility depends entirely on resolving ‘Hu Zhang’ as *Polygonum cuspidatum* (resveratrol), which is a different plant.

Reasoning: *Rhizoma Arisaematis* (Tian Nan Xing) derives from *Arisaema* species; PubMed and ChEMBL searches returned no cardiac or anti-anginal activity for any *Arisaema* compound. The classical text’s reference to ‘Hu Zhang’ maps to *Polygonum cuspidatum* in modern pharmacopoeia — a rich source of resveratrol (CID 445154; ChEMBL165, Max Phase 3) with COX-1 (IC₅₀ = 380–2,840 nM), COX-2 (IC₅₀ = 3,490 nM), and SIRT1 inhibition relevant to anti-platelet and cardioprotective activity (PMID 41978708). If the source is confirmed as knotweed, the rating rises to Plausible; as *Arisaema*, it remains Weak.

Example 4: *Calcareous Spar* (Source) and Generalized fever (Malady)

DeepRoot Verdict: UNSUPPORTED

Weakest Link: The substance is an inert mineral with no pharmacological activity relevant to systemic fever.

Reasoning: *Calcareous Spar* is primarily composed of calcium carbonate (CaCO₃), an inert mineral in the context of systemic fever reduction. It lacks bioactive organic compounds capable of modulating thermoregulatory centers or inflammatory cytokines. The KG correctly shows no paths, as there is no plausible pharmacological mechanism for this substance to treat fever.

Biomni Verdict: UNSUPPORTED

Weakest Link: CaCO₃ has no antipyretic target in any database; the antipyretic TCM mineral is Gypsum (CaSO₄), a distinct compound with documented TLR4/NF- κ B inhibition.

Reasoning: Calcareous Spar is calcium carbonate (CaCO₃; CID 10112; ChEMBL1200539, Approved). Its 31 ChEMBL and 22 OpenTargets indications cover hypocalcemia, renal osteodystrophy, GERD, and osteoporosis — none include fever. No antipyretic bioactivity was found in ChEMBL or PubMed. The TCM antipyretic mineral is Gypsum (CaSO₄, Shi Gao), which inhibits TLR4/NF- κ B signaling, reduces IL-6/PGE₂/TNF- α , and upregulates AVP in LPS-fever models (PMID 37769495). The TCM ‘cold nature clears heat’ rationale has no pharmacological correlate for CaCO₃.

Example 5: *Aster tataricus* (Source) and Cough (Malady)

DeepRoot Verdict: PLAUSIBLE

Weakest Link: The low-confidence wood-tier evidence for the strobopinin-PTGS interaction limits the certainty of this specific mechanism.

Reasoning: *Aster tataricus* is traditionally used as an antitussive, and its constituents like strobopinin target PTGS1/2 (cyclooxygenases). Inhibition of these enzymes can reduce airway inflammation and prostaglandin-mediated cough reflexes. While the KG confirms the loop to cough, the evidence is wood-tier, and other paths for quercetin and kaempferol do not close the loop to respiratory conditions.

Biomni Verdict: VERY PLAUSIBLE

Weakest Link: Shionone (CID 12315507) has zero PubChem bioassay data and is absent from ChEMBL; specific antitussive IC₅₀ values for any *Aster* compound against a defined respiratory target have not been reported.

Reasoning: *Aster tataricus* (Zi Wan) root contains shionone, epifriedelanol (CID 119242; 2 active PubChem assays), and triterpenoid saponins. Multiple in vivo mouse studies confirm mechanistically specific anti-inflammatory activity in respiratory disease: PI3K/AKT/NF- κ B inhibition in asthma (PMID 38081397), HIF-1 α /VEGFA pathway suppression in chronic bronchitis, and iNOS/LPS-mediated mediator suppression in acute bronchitis. Saponins additionally provide expectorant activity via mucociliary stimulation. Comprehensive reviews (PMIDs 40271061, 36080154) confirm the traditional antitussive use is well-corroborated by modern pharmacology.

Example 6: *Fructus Ailanthi Altissimi* (Source) and impotence (Malady)

DeepRoot Verdict: PLAUSIBLE

Weakest Link: The reliance on wood-tier evidence for the identified targets makes the specific pharmacological contribution of these compounds to erectile function speculative.

Reasoning: *Fructus Ailanthi Altissimi* contains flavonoids like quercetin and rutin, which exhibit antioxidant and vasorelaxant properties. The KG identifies paths such as quercetin→AKR1B1 and rutin→ADRA2A/C, which are linked to erectile dysfunction. While these targets are relevant to vascular health and smooth muscle tone, the evidence is primarily from wood-tier sources, limiting the strength of the mechanistic link.

Biomni Verdict: WEAK

Weakest Link: Ailanthone and related quassinoids have no demonstrated activity at PDE5, androgen receptors, eNOS, or any pro-erectile target; the pharmacological profile is anticancer and antimalarial only.

Reasoning: Ailanthus altissima fruit contains quassinoids (ailanthone, CID 72965; ChEMBL161015) with documented anticancer (HepG2 IC₅₀ = 240 nM), antiparasitic (P. falciparum IC₅₀ = 0.003 μg/mL), and NF-κB/PKM2 inhibitory activity. PubMed returned zero results for ‘Ailanthus altissima erectile dysfunction.’ No PDE5A, androgen receptor, or nitric oxide pathway activity was found in ChEMBL (8 activities) or PubChem (11 assays, 1 active). The TCM bark use is for spermatorrhea — a distinct condition — and even that lacks modern mechanistic validation.

Example 7: *Ge Gen* (Source) and *Toxicity* (Malady)

DeepRoot Verdict: PLAUSIBLE

Weakest Link: The term ‘toxicity’ is too non-specific to definitively map to the ACHE-inhibition mechanism identified in the KG.

Reasoning: Ge Gen (*Pueraria lobata*) contains isoflavones like daidzin and flavonoids like rutin and apigenin. The KG identifies paths where these compounds target ACHE (acetylcholinesterase), which is a critical enzyme in managing certain types of neurotoxicity and organophosphate poisoning. While the loop closes, the clinical application of Ge Gen for general ‘toxicity’ is broad and requires more specific evidence regarding the type of toxin.

Biomni Verdict: PLAUSIBLE

Weakest Link: The detoxification evidence is specific to alcohol-related hepatotoxicity and does not extend to the broader traditional claim of resolving ‘various toxins.’

Reasoning: Ge Gen (*Radix Puerariae*) contains puerarin (CID 5281807; ChEMBL486386, Max Phase 2) and daidzein (CID 5281708; ChEMBL8145). Puerarin modulates macrophage polarization (M1→M2) and activates PI3K/AKT signaling to suppress alcohol-induced hepatocellular cytotoxicity (PMID 42010989; 14 PubMed papers). Daidzein inhibits ALDH2 (IC₅₀ = 9,000 nM) and acts as a phytoestrogen via ERβ (IC₅₀ = 303 nM). OpenTargets confirms Phase 2 clinical registration for alcohol abuse/dependence. The mechanism is coherent for alcohol detoxification specifically, but the broader ‘various toxins’ claim lacks support.

Example 8: *Sanguisorba officinalis* (Source) and *Vaginal discharge* (Malady)

DeepRoot Verdict: UNSUPPORTED

Weakest Link: The absence of any KG paths linking the known astringent tannins of *Sanguisorba* to the treatment of vaginal discharge.

Reasoning: *Sanguisorba officinalis* is rich in hydrolyzable tannins and ellagic acid, which provide astringent and antimicrobial properties useful for treating inflammatory exudates. However, the KG paths for its constituents like astragaloside and genistein do not close the loop to vaginal discharge or related gynecological infections. The current graph evidence focuses on unrelated targets like ESR2 and CA12, which do not explain the traditional use.

Biomni Verdict: PLAUSIBLE

Weakest Link: The most direct anti-Candida evidence uses the charred/processed form (not raw root), and no human clinical trial for vaginal discharge specifically has been identified.

Reasoning: *Sanguisorba officinalis* contains ellagic acid (CID 5281855; ChEMBL6246; Src IC₅₀ = 300 nM, DNA gyrase IC₅₀ = 3.3–47 μg/mL), sanguinolone H-6 (CID 16130897), and ziyuglycoside I/II (CID 71609288). Direct anti-Candida albicans activity is confirmed for the

whole extract (PMID 12169406). Ziyuglycoside II synergizes with tetracycline against MDR *S. aureus* (MIC = 32 $\mu\text{g}/\text{mL}$; FICI = 0.25–0.375; PMID 41385864) via membrane disruption and ROS generation. Astringent ellagitannins reduce mucosal secretion and inflammation, providing a dual antimicrobial-astringent mechanism supporting the traditional use.

Example 9: *Hedgehog Pelt* (Source) and *Hemorrhoids* (Malady)

DeepRoot Verdict: PLAUSIBLE

Weakest Link: The reliance on wood-tier evidence for the interaction between the pelt's constituents and the identified receptors.

Reasoning: Hedgehog Pelt contains various organic acids and signaling molecules that the KG links to ADRA1A and NR3C1 (glucocorticoid receptor). These targets are involved in vasoconstriction and anti-inflammatory responses, which are relevant to reducing the swelling and pain associated with hemorrhoids. While the mechanism is plausible, the source material is complex and the evidence tier is low.

Biomni Verdict: UNSUPPORTED

Weakest Link: No bioactive compounds from hedgehog pelt (*Erinaceus europaeus* skin) have been characterized in any database; the entire mechanistic chain is unestablished.

Reasoning: Hedgehog pelt (Wei Pi, Pellis Erinacei) from *Erinaceus europaeus* returned zero PubMed results for 'hedgehog pelt hemorrhoids' or 'Ci Wei Pi TCM hemorrhoids.' No ChEMBL entries exist for *E. europaeus* skin compounds. Compounds found in searches (erinacines, hericenones) derive from the unrelated hedgehog mushroom (*Hericium erinaceus*). The traditional use likely reflects sympathetic magic or empirical use of charred animal material as a non-specific hemostatic, neither of which constitutes a credible pharmacological mechanism.

Example 10: *Calcareous Spar* (Source) and *Cough* (Malady)

DeepRoot Verdict: UNSUPPORTED

Weakest Link: The substance is an inert mineral with no pharmacological activity relevant to cough.

Reasoning: As a mineral composed of calcium carbonate, Calcareous Spar possesses no antitussive properties or mechanisms to modulate the cough reflex. It does not interact with respiratory receptors or inflammatory pathways involved in cough. The lack of KG paths reflects the biological reality that this substance is inert for this indication.

Biomni Verdict: UNSUPPORTED

Weakest Link: CaCO_3 has no antitussive, bronchodilatory, or anti-inflammatory mechanism in any database; PubMed returned zero results for 'calcium carbonate antitussive cough.'

Reasoning: Calcareous Spar is the same CaCO_3 (CID 10112) as Example 4, with no antitussive indication in ChEMBL or OpenTargets. While calcium signaling is mechanistically relevant to cough — Ca^{2+} -activated Cl^- channels (TMEM16A) regulate airway smooth muscle tone and cough reflex sensitivity (PMID 39608849) — these are endogenous calcium-signaling mechanisms unrelated to exogenous CaCO_3 supplementation. The TCM rationale of minerals 'descending rebellious qi' to suppress cough has no pharmacological correlate for this compound.

NeurIPS Paper Checklist

The checklist is designed to encourage best practices for responsible machine learning research, addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove the checklist: **The papers not including the checklist will be desk rejected.** The checklist should follow the references and follow the (optional) supplemental material. The checklist does NOT count towards the page limit.

Please read the checklist guidelines carefully for information on how to answer these questions. For each question in the checklist:

- You should answer [Yes], [No], or [N/A].
- [N/A] means either that the question is Not Applicable for that particular paper or the relevant information is Not Available.
- Please provide a short (1–2 sentence) justification right after your answer (even for [N/A]).

The checklist answers are an integral part of your paper submission. They are visible to the reviewers, area chairs, senior area chairs, and ethics reviewers. You will also be asked to include it (after eventual revisions) with the final version of your paper, and its final version will be published with the paper.

The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation. While [Yes] is generally preferable to [No], it is perfectly acceptable to answer [No] provided a proper justification is given (e.g., error bars are not reported because it would be too computationally expensive” or “we were unable to find the license for the dataset we used”). In general, answering [No] or [N/A] is not grounds for rejection. While the questions are phrased in a binary way, we acknowledge that the true answer is often more nuanced, so please just use your best judgment and write a justification to elaborate. All supporting evidence can appear either in the main paper or the supplemental material, provided in appendix. If you answer [Yes] to a question, in the justification please point to the section(s) where related material for the question can be found.

IMPORTANT, please:

- **Delete this instruction block, but keep the section heading “NeurIPS Paper Checklist”.**
- **Keep the checklist subsection headings, questions/answers and guidelines below.**
- **Do not modify the questions and only use the provided macros for your answers.**

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope?

Answer: [Yes]

Justification: The abstract’s three quantitative claims map directly to results: $R@20 = 47.6\%$ vs. 4.8% for the raw-corpus LLM and $\sim 2.4\%$ random (§3.3, Table 1); reasoning-quality dominance across all KG-augmented configurations with hallucination rates 7–10% (§3.4, Table 2); and the 87% tool-call hallucination rate is the audited value from the same table. Scope is restricted to *Shen Nong Ben Cao Jing* as stated, with generalization flagged in §4.

Guidelines:

- The answer [N/A] means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A [No] or [N/A] answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: §4 discusses single-corpus evaluation, the small ($N = 21$) held-out KNOWN_TREATS slice, flat (face-validity) priors that have not been calibration-characterized, Open Targets' human-disease-only coverage truncating non-modern indications, and the dependence of critic rationale quality on the underlying LLM. Coverage caveats on the natural-product target distribution (phenotypic-heavy TARGETS) are also discussed in Appendix A.6.

Guidelines:

- The answer [N/A] means that the paper has no limitation while the answer [No] means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [N/A]

Justification: The paper makes no theoretical claims requiring proof.

Guidelines:

- The answer [N/A] means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: The full graph schema (Appendix A.1), each Assembly agent’s grounding source, prompt structure, and verification protocol (Appendix A.2), the Discovery agents’ deterministic ranking key (Appendix A.3), the four evaluation protocols including held-out construction and InChIKey planar-prefix masking (Appendix A.7), the prompt templates (Appendix A.8), and the flat tier priors (Appendix A.5) are reported. The corpus is public and external grounding databases (COCONUT, PubChem, ChEMBL, Open Targets, NCBI Taxonomy, OLS4, MeSH, ICD-10, SNOMED) are all freely accessible.

Guidelines:

- The answer [N/A] means that the paper does not include experiments.
- If the paper includes experiments, a [No] answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [No]

Justification: Code is not released at submission to preserve anonymity; we commit to a public release (agent code, Cypher schema, prompt files, evaluation harness, and the constructed Neo4j KG dump) upon acceptance. In the interim, the full agent protocols (Appendix A.2), prompt templates (Appendix A.8), evaluation protocols (Appendix A.7), and tier priors (Appendix A.5) are documented in detail, and all grounding databases are public.

Guidelines:

- The answer [N/A] means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://neurips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so [No] is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://neurips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer) necessary to understand the results?

Answer: [Yes]

Justification: No model training is performed—the LLM stages call frozen Gemini 3.1-Flash-Lite/2.5 Flash/3.1 Pro at temperature 0 with structured-JSON schemas (Appendix A.9). Test settings are specified per evaluation: 30 mini-corpora with 3+7 loop-closing/distractor splits (Eval 1), 11-step perturbation grid (Eval 2), 21-pair held-out slice with planar-InChIKey masking and per-disease pool sizes 87–1,954 (Eval 3, Appendix A.7), and 30 stratified claims graded by Claude Sonnet 4.6 across six rubric dimensions (Eval 4). Activity thresholds ($pchembl \geq 5.0$, assay type $\in \{B, F\}$), Open Targets score floors, and tier priors (Table S3) are all reported.

Guidelines:

- The answer [N/A] means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: We report the structural-significance test in Figure 2B / §3.1: critic self-confidence over 30 closed-loop claims is monotonically decreasing across an 11-level edge-perturbation grid, crossing the LLM-only baseline at $\sim 50\%$ corruption and continuing to ~ 0.30 at full shuffle—ruling out both the “graph-as-decoration” null and the “silent fallback to LLM-only” null. Recall@ K in Eval 3 is reported on the complete historically reachable held-out slice ($N = 21$) rather than a sub-sample, and per-disease pool sizes (87–1,954) anchor the random-recall reference. We do not report bootstrap CIs around point estimates; given the small fixed evaluation sets and single-seed runs, this is acknowledged in §4.

Guidelines:

- The answer [N/A] means that the paper does not include experiments.

- The authors should answer [Yes] if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g., negative error rates).
- If error bars are reported in tables or plots, the authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: No GPUs are required; all LLM calls are made to hosted APIs (Gemini, Claude). Appendix A.9 reports end-to-end Assembly wall-clock (~ 30–40 minutes over the 71-chunk corpus, dominated by Neo4j AuraDB MERGE latency) and total LLM spend (~ \$0.25 at Flash-Lite). Table 1 reports nominator wall time (18 s) vs. raw-corpus LLM (1,469 s). External grounding services (ChEMBL, COCONUT, PubChem, Open Targets, NLM, EBI) are free. The local COCONUT inverted index (~700,000 structures) fits in standard laptop memory.

Guidelines:

- The answer [N/A] means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [Yes]

Justification: The work uses a public millennia-old corpus and curated public biomedical databases under their respective terms of use; no human subjects, no scraped private data, no PII. Anonymity is preserved (no author identifiers in the submission). Attribution and benefit-sharing concerns specific to traditional knowledge are flagged in §4.

Guidelines:

- The answer [N/A] means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer [No], they should explain the special circumstances that require a deviation from the Code of Ethics.

- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: §4 discusses positive impact (auditable, mechanism-grounded drug-discovery leads from underused historical corpora; provenance preserved end-to-end) and two non-trivial negative impacts: (i) stripping leads from their KG provenance to present them as therapeutic claims, and (ii) attribution and benefit-sharing concerns when commercializing compounds mined from communal traditional knowledge. We also note that outputs are leads for downstream *in vitro* validation, not patient-facing inference.

Guidelines:

- The answer [N/A] means that there is no societal impact of the work performed.
- If the authors answer [N/A] or [No], they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate Deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pre-trained language models, image generators, or scraped datasets)?

Answer: [N/A]

Justification: We release no pretrained models, no generative models, and no scraped data. The new asset planned for release on acceptance is a constructed knowledge graph over a public corpus and public biomedical databases, with provenance fields on every edge. The misuse channel relevant to this work—reading therapeutic leads as clinical advice—is addressed in §4 rather than via gated release.

Guidelines:

- The answer [N/A] means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.

- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: All grounding databases and the corpus are cited at first use in §2.1: COCONUT [Chandrasekhar et al., 2025], PubChem [Kim et al., 2023], ChEMBL [Mendez et al., 2019], Open Targets [Ochoa et al., 2021], NCBI Taxonomy [Schoch et al., 2020], OLS4 [McLaughlin et al., 2025]; ICD-10, MeSH, and SNOMED are accessed via NLM/EBI lookup services. We use each within its public terms of use (COCONUT CC BY 4.0; ChEMBL CC BY-SA 3.0; PubChem and NCBI as public-domain US Government works; Open Targets CC0 1.0; SNOMED via the EBI OLS4 API as permitted for research use). The *Shen Nong Ben Cao Jing* is a public-domain text. RDKit (BSD-3) is used for InChIKey computation.

Guidelines:

- The answer [N/A] means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset’s creators.

13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [Yes]

Justification: The new asset is the DeepRoot knowledge graph constructed by Assembly over *Shen Nong Ben Cao Jing* (21,111 active nodes, 52,467 active edges; per-type counts in Appendix A.6). Documentation is in-paper: the full node and edge schema with property-level descriptions (Appendix A.1, Table S1), per-edge evidence-type breakdown (Appendix A.6), tier priors per evidence type (Table S3), and the Assembly agent that produced each edge type (Table S2). The graph dump and code will be released alongside this paper on acceptance; the paper version of submission is anonymized.

Guidelines:

- The answer [N/A] means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [N/A]

Justification: The paper involves no crowdsourcing and no human subjects.

Guidelines:

- The answer [N/A] means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [N/A]

Justification: No human subjects research is conducted, so IRB review does not apply.

Guidelines:

- The answer [N/A] means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does *not* impact the core methodology, scientific rigor, or originality of the research, declaration is not required.

Answer: [Yes]

Justification: LLMs are a core methodological component and their usage is fully declared. The Assembly agents (extraction, auditor canonicalization, malady-to-disease mapping, reviewer Pass 2) call Gemini 3.1-Flash-Lite at temperature 0 with structured-JSON schemas (Appendix A.2, Appendix A.9). The Discovery critic agent uses Gemini 3.1 Pro by default, with 2.5 Flash and 3.1 Flash-Lite reported as ablation tiers in Table 2; the discovery (nominator) agent contains *no* LLM in the loop and is fully deterministic (Appendix A.3). The reasoning-quality judge in Eval 4 is Claude Sonnet 4.6, deliberately cross-family from the graded systems. Prompt structure for each role is described in Appendix A.8.

Guidelines:

- The answer [N/A] means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy in the NeurIPS handbook for what should or should not be described.