

An Agentic System for Rare Disease Diagnosis with Traceable Reasoning

Weike Zhao^{1,2,*}, Chaoyi Wu^{1,2,*}, Yanjie Fan^{3,*}, Xiaoman Zhang⁴, Pengcheng Qiu^{1,2}, Yuze Sun¹,
Xiao Zhou², Yanfeng Wang^{1,2}, Ya Zhang^{1,2,†}, Yongguo Yu^{3,†}, Kun Sun^{3,†} and Weidi Xie^{1,2,†}

¹Shanghai Jiao Tong University, Shanghai, China

²Shanghai Artificial Intelligence Laboratory, Shanghai, China

³Xinhua Hospital affiliated to Shanghai Jiao Tong University School of Medicine, Shanghai, China

⁴Department of Biomedical Informatics, Harvard Medical School, Boston, MA, USA

*Equal contributions †Corresponding author

Ya Zhang: ya_zhang@sjtu.edu.cn; Yongguo Yu: yuyongguo@shsmu.edu.cn;

Kun Sun: sunkun@xinhua.com.cn; Weidi Xie: weidi@sjtu.edu.cn

Rare diseases collectively affect over 300 million individuals worldwide, yet timely and accurate diagnosis remains a pervasive challenge. This is largely due to their clinical heterogeneity, low individual prevalence, and the limited familiarity most clinicians have with rare conditions. Here, we introduce **DeepRare**, the first rare disease diagnosis agentic system powered by a large language model (LLM), capable of processing heterogeneous clinical inputs: free-text clinical descriptions, structured Human Phenotype Ontology (HPO) terms, and genetic testing results in variant call format (VCF). The system generates ranked diagnostic hypotheses for rare diseases, each accompanied by a transparent chain of reasoning that links intermediate analytic steps to verifiable medical evidence. This interpretability is critical for clinical adoption, supporting human-AI collaboration in diagnostic workflows.

DeepRare comprises three key components: a central host with a long-term memory module; specialized agent servers responsible for domain-specific analytical tasks (*e.g.*, phenotype extraction, variant prioritization) integrating over 40 specialized tools and web-scale, up-to-date medical knowledge sources, ensuring access to the most current clinical information. This modular and scalable design enables complex diagnostic reasoning while maintaining traceability and adaptability. We evaluate DeepRare on eight datasets sourced from literature, case reports, and clinical centers across Asia, North America, and Europe, spanning 14 medical specialties, including neurology, cardiology, immunology, and genetics. The system demonstrates exceptional diagnostic performance among 2,919 diseases, achieving 100% accuracy for 1013 diseases. In HPO-based evaluations, DeepRare significantly outperforms other 15 methods, like traditional bioinformatics diagnostic tools, large language models, and other agentic systems, achieving an average Recall@1 score of 57.18% and surpassing the second-best method (Reasoning LLM) by a substantial margin of 23.79 percentage points. For multi-modal input scenarios, DeepRare achieves 70.60% at Recall@1 compared to Exomiser’s 53.20% in 109 cases. Manual verification of reasoning chains by clinical experts achieves 95.40% agreements, confirming that the system’s intermediate reasoning steps are both medically valid and traceable to authoritative sources, enhancing its potential as a trustworthy decision support tool in rare disease diagnostics. Furthermore, the DeepRare system has been implemented as a user-friendly web application <http://raredx.cn/doctor>.

1 Introduction

Rare diseases—defined as conditions affecting fewer than 1 in 2,000 individuals—collectively impact over 300 million people worldwide, with more than 7,000 distinct disorders identified to date, approximately 80% of which are genetic in origin [1, 2, 3, 4]. Despite their cumulative burden, rare diseases remain notoriously difficult to diagnose due to their clinical heterogeneity, low individual prevalence, and limited clinician familiarity [1, 2, 5, 6, 7, 8, 9, 10, 11, 12, 13]. Patients often experience a prolonged “diagnostic odyssey” averaging over five years, marked by repeated referrals, misdiagnoses, and unnecessary interventions, all of which contribute to delayed treatment and adverse outcomes [5, 14]. These challenges highlight the urgent

need for scalable, accurate, and interpretable diagnostic tools, an area where recent advances in multi-agent systems offer transformative potential.

Developing artificial intelligence (AI) systems for rare disease diagnosis presents several inherent challenges, (i) **multidisciplinary**: rare diseases often manifest with complex, heterogeneous, and multisystem symptoms, requiring diagnostic models to possess multidisciplinary medical knowledge and the ability to interpret diverse patient phenotypes [15, 16]; (ii) **limited cases**: the scarcity of cases for individual rare diseases limits the availability of training data, making it difficult to develop robust models and increasing the risk of overfitting and catastrophic forgetting; (iii) **dynamic knowledge updates**: the rare disease knowledge landscape is rapidly evolving, with approximately 260 to 280 new diseases added annually, according to the International Rare Diseases Research Consortium (IRDIRC) [17]. This dynamic nature demands AI systems that are not only updatable but also capable of integrating new knowledge efficiently; (iv) **transparency and traceability**: clinical deployment demands interpretability: diagnostic suggestions must be accompanied by transparent, traceable reasoning to support clinician trust and accountability.

Recent advances in agentic large language model (LLM) systems have opened new avenues for rare disease diagnosis [18, 19, 6, 20, 21, 22, 23, 24], which orchestrates multiple specialized tools and sub-agents [18, 19], enabling seamless integration of external knowledge bases, case repositories, and multimodal analytical components [25, 23]. Unlike conventional supervised learning approaches, these systems are typically training-free and excel in few-shot and zero-shot scenarios—an essential capability for rare disease applications where annotated data are scarce. Their modular and interpretable architectures further facilitate transparent, auditable, and clinically actionable diagnostic workflows.

Here, we present **DeepRare**, an agentic LLM-based system designed specifically for rare disease diagnosis. DeepRare is capable of processing heterogeneous patient inputs, including free-text clinical descriptions, structured Human Phenotype Ontology (HPO) terms, and genomic testing results. Based on the input, the system generates a ranked list of candidate diagnoses, each supported by a transparent chain of reasoning that directly references verifiable medical evidence. This design enhances interpretability and supports clinician trust in AI-assisted diagnostic decisions.

Specifically, our proposed system is inspired by the Model Context Protocol (MCP) [25] and consists of three hierarchical tiers. At its core, a central host backed by a memory bank and powered by a state-of-the-art LLM—coordinates the diagnostic process and retains contextual information. Surrounding this host are multiple agent servers, each dedicated to specialized analytical tasks, such as phenotype extractor, disease normalization, knowledge retrieval, case matching, phenotype analysis and genotype analysis. The outermost tier consists of curated and web-scale external data sources, ensuring access to the most current clinical evidence. To further improve diagnostic accuracy and robustness, **DeepRare** implements a self-reflective diagnostic loop, prompting the central host to iteratively re-evaluate intermediate hypotheses by gathering additional evidence. This reduces the risk of over-diagnosis and mitigates LLM hallucinations.

We evaluate **DeepRare** on 6,401 clinical cases collected from seven public datasets and one in-house dataset, sourced from diverse populations across Asia, North America, and Europe. Among these, we construct an in-house dataset comprising 975 patient cases representative of the Chinese rare disease population, including 109 cases with whole exome sequencing results. To the best of our knowledge, this is the only rare disease diagnosis benchmark featuring original gene testing data. All diagnoses in this cohort are rigorously validated by genetic testing, providing a high-quality standard for assessing diagnostic performance. **DeepRare** consistently achieves superior diagnostic accuracy across all 8 datasets of 2919 rare diseases spanning 14 medical specialties.

Notably, among the 2,919 evaluated rare diseases, **DeepRare** attains 100% accuracy for 1013 diseases. In HPO-based evaluations, compared with other 15 methods like traditional bioinformatics tools, large language models, and agentic systems, **DeepRare** achieves an average score of 57.18%, 65.25% at Recall@1, and Recall@3, respectively, surpassing the second-best method (Reasoning LLM) by substantial margins of 23.79%, 18.65%. In multi-modal input scenarios, **DeepRare** achieves a Recall@1 of 70.6%, outperforming Exomiser’s 53.2% in the 109 whole-exome cases. Additionally, we engage 10 rare disease physicians to manually verify the traceable reasoning chains generated by the system across 180 cases. DeepRare demonstrates high reliability in evidence factuality, achieving 95.4% agreement with clinical experts, thereby confirming that its intermediate

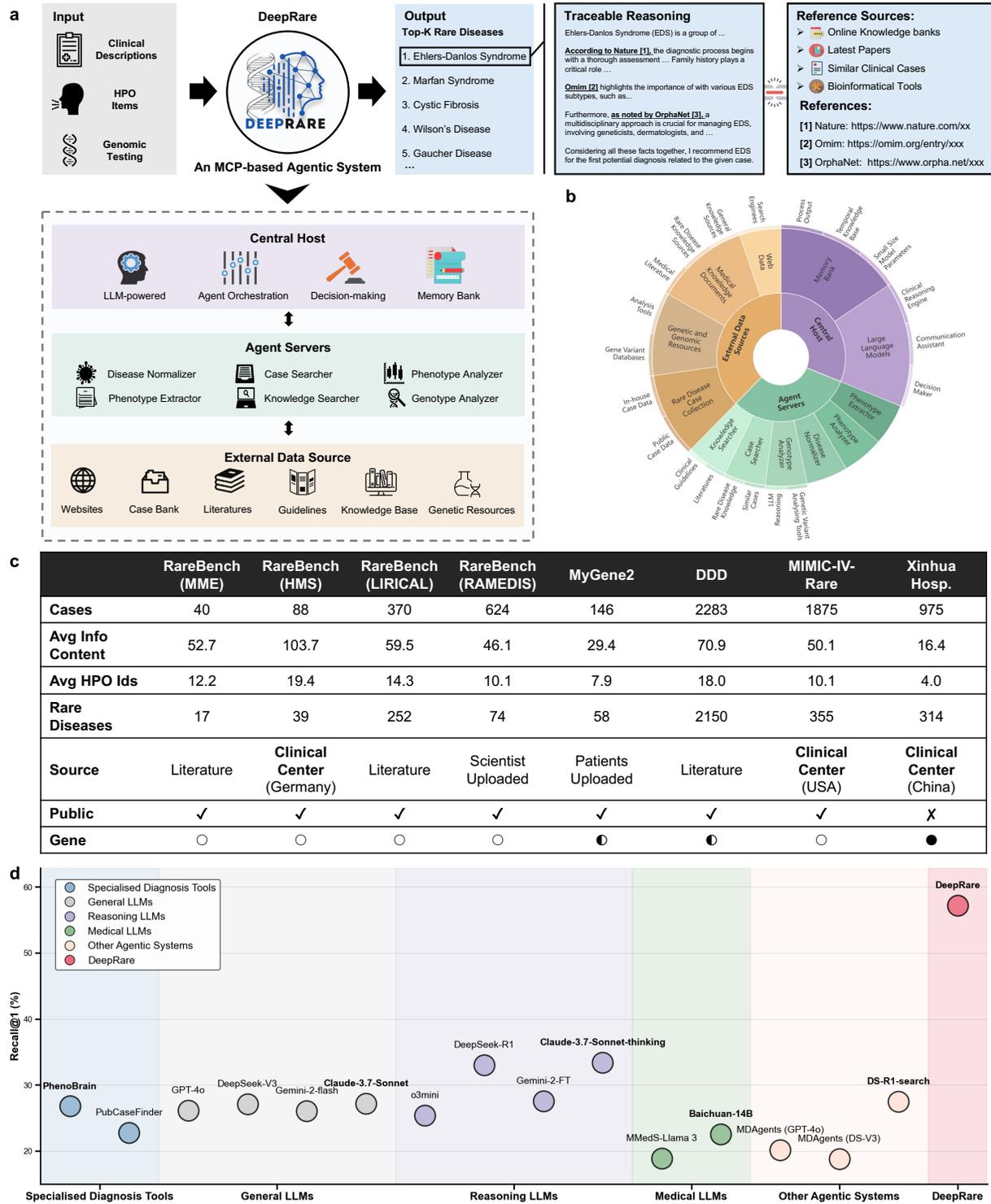


Figure 1 | DeepRare: An agentic framework for rare disease prioritization. (a) System workflow: Multi-modal patient data (HPO terms, genomic variants) are processed through a tiered MCP-inspired architecture, generating a ranked Top-K diagnosis list with evidence-supported reasoning chains. (b) Knowledge architecture: Sunburst visualization depicting hierarchical integration of diagnostic tools and biomedical knowledge sources within DeepRare. (c) Multi-center benchmark characteristics: Case distributions, phenotypic complexity (HPO metrics), disease spectrum, provenance, and genetic annotation status (solid: confirmed pathogenic variants; half-solid: candidate variants extracted; hollow: no genetic data). (d) Performance benchmarking: Comparative evaluation across diagnostic APIs, general-purpose LLMs, reasoning-enhanced LLMs, medically-tuned LLMs, and agentic systems.

reasoning steps are both medically valid and traceable to authoritative sources. To facilitate clinical adoption, we have deployed **DeepRare** as a user-friendly web application as a diagnostic copilot for rare disease physicians. Finally, we discuss the robustness of our agentic framework by evaluating different underlying LLMs and analyzing the contribution of each module, demonstrating the superiority of our system design.

2 Results

This section presents the results of our study, beginning with an overview of the proposed framework, **DeepRare**, and the evaluation settings, followed by a detailed analysis of the main findings.

2.1 System Overview

DeepRare is an LLM-powered agentic system for rare disease diagnosis. It features a three-tier architecture that synergistically integrates reasoning-enhanced large language models with a broad range of clinical knowledge sources as shown in Figure 1a and 1b. The system comprises: (i) a central host, powered by LLMs and equipped with a memory bank, which orchestrates the entire diagnostic workflow by synthesizing collected evidence; (ii) multiple specialized agent servers, each managing a local set of tools to perform various rare disease-related analytical tasks and interact with distinct resource environments; and (iii) heterogeneous web-scale medical sources, which provide essential and traceable diagnostic evidence—such as research articles, clinical guidelines, and existing patient cases, *etc.*

Upon receiving a clinical case—provided as free-text phenotypic descriptions, structured Human Phenotype Ontology (HPO) terms, raw VCF files, or any combination thereof—the central host systematically decomposes the diagnostic task. It first orchestrates the agent servers to retrieve relevant evidence and references from external data sources, tailored to the patient’s information. The host then synthesizes this evidence to generate preliminary diagnostic hypotheses, followed by a self-reflection phase in which it conducts additional searches to rigorously validate or refute these hypotheses. If no hypothesized diseases meet the self-reflection criteria, the system iteratively revisits earlier steps to acquire further patient-specific evidence, repeating this diagnostic loop until a satisfactory resolution is achieved. Ultimately, **DeepRare** outputs a ranked list of potential rare diseases, each accompanied by a transparent reasoning chain that directly links each inference step to trusted medical evidence. Further details on the system workflow are provided in the **Method** section.

2.2 Evaluation Settings

To evaluate the performance of **DeepRare**, we consider **three baseline approaches**:

- **Specialised rare disease diagnosis tools.** We consider the bioinformatics tools that are directly designed for rare disease diagnosis. PhenoBrain [26] is a tool for Human Phenotype Ontology (HPO) analysis that processes structured HPO terms and outputs potential rare disease candidates by ensembling the predicted probabilities from five classic diagnostic models. PubCaseFinder [27] performs HPO-based diagnostic analysis by identifying and matching the most similar public cases from PubMed case reports.
- **Latest large language models.** We compare different LLMs, including **general LLMs, reasoning-enhanced LLMs, and medical LLMs**. **General LLMs** denote the most commonly used LLMs without extra reasoning enhancement or domain alignment. Specifically, GPT-4o [28], DeepSeek-V3 [29], Gemini-2.0-flash [30], and Claude-3.7-Sonnet [31] are considered here. **Reasoning LLMs** denote the latest generation LLMs enhanced with an explicit reasoning chain, including o3mini [32], DeepSeek-R1 [29], Gemini-2.0-FT [30], Claude-3.7-Sonnet-thinking [31]. **Medical LLMs** refers to large language models specifically developed for the medical domain, with Baichuan-14B [33] and MMedS-Llama 3 [34] serving as notable representatives. All these LLMs are adapted to rare disease diagnosis, leveraging well-designed Prompt 1.
- **Other agentic systems.** We compare to existing agentic systems that build upon different LLMs. These systems are more powerful as they integrate LLMs with other external sources. However, inefficient agentic system design could also introduce unnecessary information (noise), which may result in a performance drop. We consider some open-source disease diagnosis methods, which include MDAgents [35], an agent that conducts discussion via multidisciplinary consultations, and the online

search capabilities of DeepSeek-V3, which augment LLMs with real-time internet information. To enable batch inference, we utilize the online interface provided by Volcano Engine [36].

More detailed descriptions on these baselines can be found in the **Method** section.

To demonstrate the effectiveness of our method, we perform a thorough cross-center evaluation. As shown in Figure 1c, we consider **eight rare disease diagnostic evaluation datasets**, with 6,401 clinical cases collected from seven public datasets and one in-house dataset. These comprised the standardized RareBench database (1,122 cases across 362 rare diseases) [8], MyGene2 (146 cases) [37, 38], the Deciphering Developmental Disorders Study (DDD) (2283 cases) [39], 1,875 public cases curated from MIMIC-IV [40, 41], and 975 cases from Xinhua Hospital affiliated to Shanghai Jiao Tong University.

These datasets can be categorized into three groups based on sources, denoting varying diagnostic difficulties:

- **From research papers:** RareBench-MME [42] (40 cases), RareBench-LIRICAL [43] (370 cases), DDD [39] (2283 cases). These cases are extracted from papers and manually verified. Considering that the cases present in the literature are often typical and well-documented, diagnosis on these datasets tends to be relatively easy.
- **From case reports:** RareBench-RAMEDIS [44] (624 cases), MyGene2 [37, 38] (146 cases). These rare disease cases are uploaded by scientists or patients through case reports. They offer authenticity but also have undergone extra manual filtering, thus presenting moderate diagnostic difficulty.
- **From real clinical centers:** RareBench-HMS [45] (88 cases), MIMIC-IV-Rare [46] (1,875 cases), Xinhua Hosp. (975 cases). These datasets are directly collected from three independent clinical centers of real patients in daily diagnostic procedures. Due to the complexity and diversity of real patients, these benchmarks are more challenging and aligned with real clinical practice. Specifically, RareBench-HMS is collected from the outpatient clinic at Hannover Medical School in Germany. MIMIC-IV-Rare is collected from Beth Israel Deaconess Medical Center in Boston, MA, USA by filtering for the rare-disease related cases. Xinhua Hosp. is a newly collected in-house data from Xinhua Hospital Affiliated to Shanghai Jiao Tong University School of Medicine, China. Notably, 109 cases within the Xinhua dataset include original gene VCF files generated from whole-exome sequencing (WES), representing the only test set in our collection with available genomic data. **Due to patient privacy considerations, the Xinhua hosp. dataset evaluates exclusively using local models without external API access.** These datasets cover three distinct regions, enabling cross-center evaluation of our methods on different population distributions in various countries.

As shown by Figure 1c from a statistical perspective, the number of rare diseases represented across various datasets ranges from 17 to 2150, while the average number of HPO items per patient varies between 4.0 and 19.4. Moreover, following [8, 47], we calculate the average information content (IC) for each dataset. Information content quantifies the specificity of a concept within an ontology by measuring its inverse frequency of occurrence, *i.e.*, concepts that appear less frequently in the corpus have higher IC values. Lower IC values typically correspond to more general terms within the ontology hierarchy [47]. This metric highlights the varying levels of diagnostic complexity across our evaluated datasets, enabling more nuanced interpretation of model performance in relation to task difficulty. To the best of our knowledge, this collection is the most comprehensive benchmark for rare disease diagnosis, covering **2919 diseases** from different case sources, and multiple independent clinical centers.

For each diagnostic task, we generated the top five most probable diagnostic predictions. The position of the correct diagnosis within these predictions was determined using GPT-4o under Prompt 2, and we subsequently calculated Recall@1, Recall@3, and Recall@5 metrics across the entire dataset.

To validate the reliability of our automated evaluation approach, we engaged eight rare disease specialists, each with over 10 years of clinical experience, to independently verify the accuracy of the LLM-based assessments. The Pearson correlation coefficient between physician rankings and LLM rankings was 0.8689, indicating strong agreement between human expert and automated evaluations. In our analysis of 240 cases, 88% demonstrated concordant assessments between physicians and the LLM evaluation system. In 10% of cases, physicians ranked the correct diagnosis higher (*i.e.*, assigned it a better position) than the LLM-based evaluation. This pattern

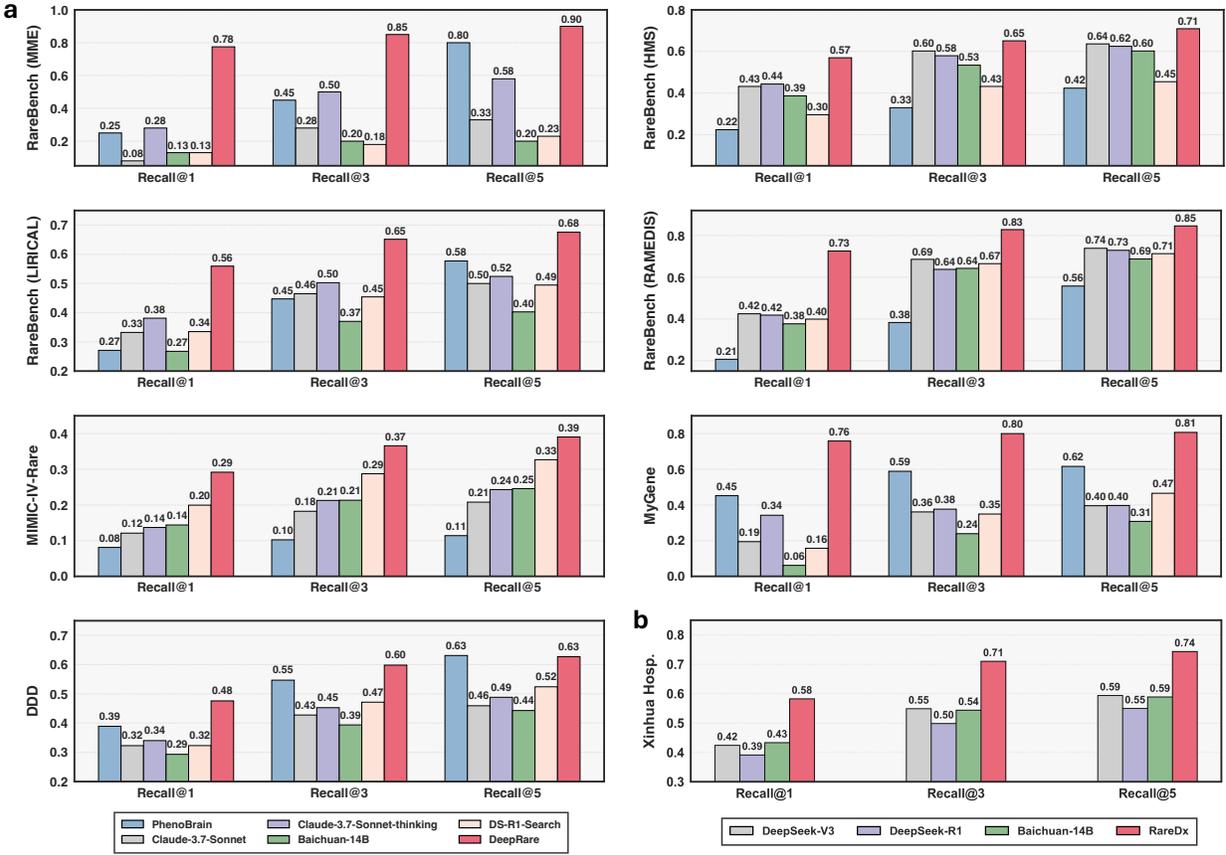


Figure 2 | HPO-wise cross-dataset evaluation and comparative performance of DeepRare. (a) Diagnostic accuracy on seven public rare disease registries, demonstrating DeepRare’s significant advantage over leading baselines – particularly in RareBench-MME (70.0% top-1 accuracy) and RareBench-RAMEDIS (72.6% top-1 accuracy). (b) Superior performance consistency on the Xinhua Hospital cohort (local model evaluation only due to privacy consideration).

suggests that physicians may consider diagnoses within the same broad disease category as sufficiently correct, potentially employing more nuanced clinical judgment that recognizes subtle indicators. Conversely, the LLM evaluation system appears to apply more conservative and stringent ranking criteria, tending to require a diagnosis that is completely and precisely correct to assign a high ranking. From a validation standpoint, these findings demonstrate that large-scale LLM-based evaluation can produce reliable and methodologically sound outcomes for automated assessment purposes, while acknowledging that human clinical expertise may capture additional diagnostic nuances. Additionally, these clinical specialists evaluated the validity of the models’ diagnostic reasoning processes to ensure clinical relevance and accuracy.

2.3 HPO-wise Diagnostic Performance Analysis across Datasets

Figure 1d presents a comparison on HPO-wise diagnosis of average Recall@1 across all benchmarks (except for Xinhua Hosp due to privacy issues). **Our proposed DeepRare clearly demonstrates superior performance across all method categories, achieving 57.18% top-1 diagnosis recall, and significantly outperforming the second-best method, Claude-3.7-Sonnet-thinking (33.39%).** Specifically, we can draw the following key observations: (i) LLM-supported approaches consistently outperform traditional rare disease diagnostic models (PhenoBrain, PubCaseFinder), demonstrating enhanced flexibility in handling diverse clinical presentations; (ii) Reasoning-enhanced LLMs systematically surpass their general-purpose counterparts without explicit reasoning, likely due to their transparent reasoning traces that improve diagnostic accuracy; (iii) General-purpose LLMs unexpectedly exceed medical domain-tuned LLMs in performance,

potentially reflecting parameter scale advantages and broader training diversity; (iv) Our multi-agent framework significantly advances beyond existing single-model approaches, highlighting the value of orchestrated specialist agents in complex diagnostic reasoning.

As shown in Figure 2a and 2b, we present detailed comparison on each dataset (only the top-performing models in each baseline category are shown). Complete results for all methods can be found in Supplementary Table 1. **Our proposed DeepRare system consistently outperforms all existing methods on all benchmarks.** Specifically, in the RareBench (MME) evaluation, **DeepRare** achieves exceptional scores of 78%, and 85% for Recall@1, Recall@3, respectively, surpassing the second-best baseline method (PubCaseFinder) by margins of 30%, 20%. The system demonstrates particularly strong results on the MyGene2 evaluation with scores of 74%, 81% surpassing second methods by substantial margins of 35%, 28%.

In clinical datasets, **DeepRare** maintains its performance edge on the MIMIC-IV-Rare test (29%, 37%). In addition to the public benchmarks, we also report its performance on the in-house clinical testset, Xinhua Hosp (Figure 2b). We mainly compare to the DeepSeek-V3, DeepSeek-R1, and MedIns, which can be implemented locally. **DeepRare** achieves 58%, 71% for Recall@1 and Recall@3, significantly surpassing the other methods.

2.4 HPO-wise Diagnostic Performance Analysis across Specialists

In addition to analyzing HPO-wise diagnostic performance across datasets, we also present results across different medical specialties, highlighting the system’s broad understanding of diverse medical knowledge.

Specifically, we categorized all test cases based on 14 body system specialists, following the taxonomy introduced by MedlinePlus¹ [48], namely, Blood, Heart and Circulation; Bones, Joints and Muscles; Brain and Nerves; Digestive System; Ear, Nose and Throat; Endocrine System; Eyes and Vision; Immune System; Kidneys and Urinary System; Lungs and Breathing; Mouth and Teeth; Skin, Hair and Nails; Female Reproductive System; and Male Reproductive System. The diseases are categorized by DeepSeek-V3 under Prompt 3.

We then present the performance comparison on various specialists. It is important to note that each case may involve multiple specialists. Similarly, due to privacy concerns regarding the in-house cases, we only evaluate methods that do not require uploading cases via online LLM APIs. Thus, DeepSeek-V3, DeepSeek-R1, and MedIns are retained to represent general LLMs, reasoning LLMs, and medical LLMs, respectively.

The results are illustrated in Figure 3a, **DeepRare** demonstrates substantial performance superiority across almost all specialties. For example, in the Endocrine System category, **DeepRare** achieves a top-1 diagnostic accuracy of 60%, significantly higher than the second-best method at 32%. Similarly, for 729 cases in the Digestive System category, **DeepRare**’s top-1 diagnostic accuracy reached 49%, substantially outperforming the second-best method at 34%. Notably, this analysis reveals that our **DeepRare** performs best in the Kidneys and Urinary System, achieving an accuracy of 66%, while showing relatively lower performance in the Lungs and Breathing System with an accuracy of 31%, reflecting its clinical application boundaries.

In Figure 3b, we further compare diagnostic performance across different methods on cases with varying number of medical specialists involved. **The results reveal a clear trend: as the number of specialists increases, DeepRare’s diagnostic accuracy improves significantly—a pattern not observed in other approaches.** Intuitively, cases involving multiple specialists tend to offer more distinctive and informative disease patterns, potentially reducing diagnostic ambiguity. However, they also present increased complexity, requiring a diagnostic system to integrate a broader spectrum of medical knowledge. This observation highlights **DeepRare**’s strength in synthesizing cross-domain clinical information, enabling it to effectively diagnose complex rare disease cases by leveraging external, diverse medical resources.

2.5 Diagnostic Performance Analysis with HPO and Gene Data

To comprehensively evaluate our system’s diagnostic capabilities, we investigate the performance when incorporating both HPO and genetic data as inputs. We conducted this evaluation on 109 cases from the Xinhua dataset, specifically selecting cases with complete whole-exome sequencing data to ensure robust comparative analysis. As shown in Figure 3c, the integration of genetic information yielded substantial

¹<https://medlineplus.gov/healthtopics.html>

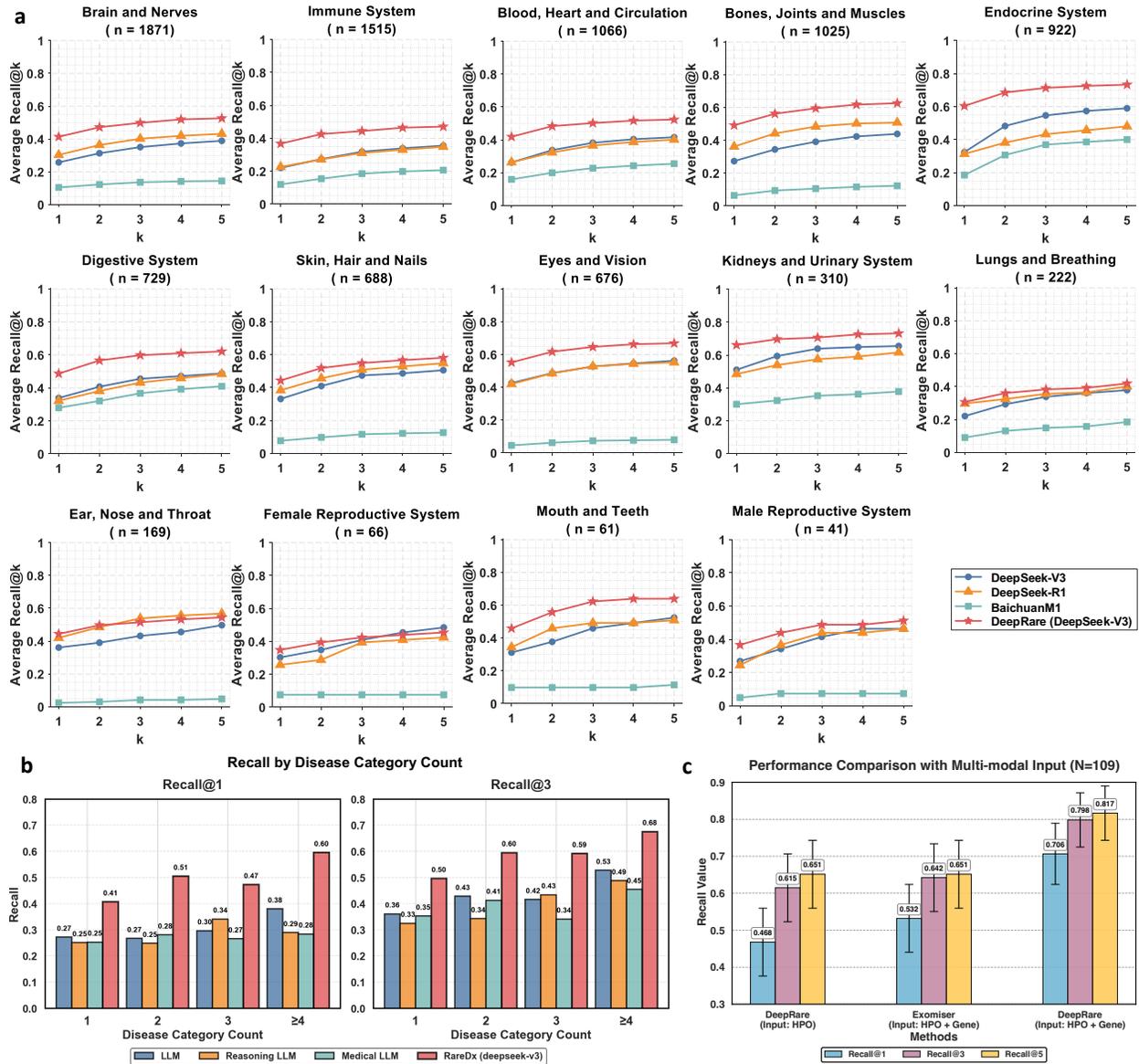


Figure 3 | DeepRare’s diagnostic performance. (a) Comparison of diagnostic accuracy across fourteen body systems: showing DeepRare’s superior performance in most specialties compared to LLM (DeepSeek-V3), Reasoning LLM (DeepSeek-R1), and Medical LLM (MedIns). (b) Diagnostic accuracy relative to disease complexity: demonstrating DeepRare’s increasing accuracy with greater disease complexity (measured by the number of specialties involved). (c) Diagnosis performance with HPO and gene data input compared with baseline method and only HPO input.

performance improvements, with Recall@1 increasing dramatically from 46.8% to 70.6%. In addition, we compared our approach with other bioinformatics tools that similarly process both HPO and genetic data, such as Exomiser [49], which is also utilized as a component within our system. As demonstrated in Figure 3c, our system achieved performance comparable to Exomiser even without genetic data input, and showed marked improvement when genetic information was incorporated. These results demonstrate that our agentic system significantly outperforms existing bioinformatics diagnostic tools in rare disease comprehensive analysis.

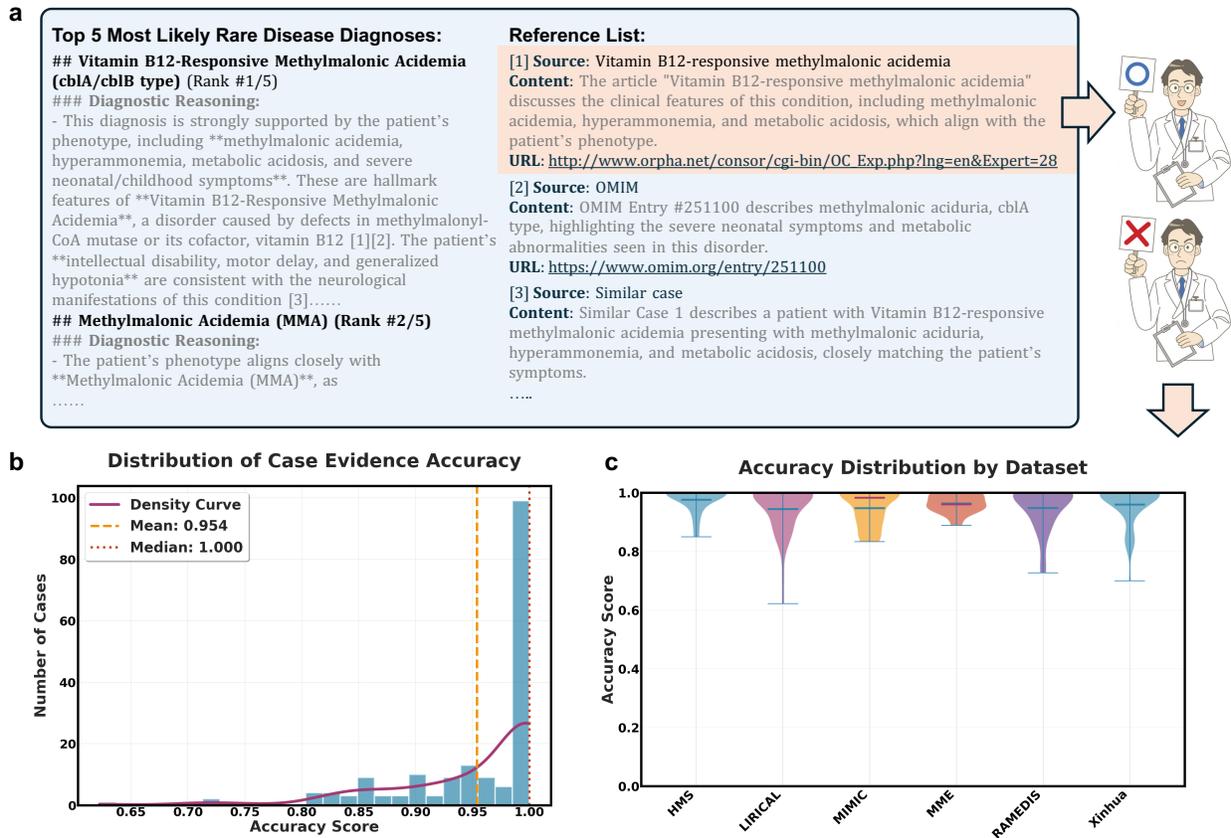


Figure 4 | Validation of traceable reasoning chain in DeepRare diagnostic system. (a) Representative case output demonstrating differential diagnosis with an evidence-based reference list. **(b)** Histogram of reference accuracy scores with density curve (mean = 0.954, median = 1.000). **(c)** Dataset-specific accuracy distributions showing robust performance across eight rare disease datasets.

2.6 Traceable Reasoning Chain Validation

To assess the reliability and clinical relevance of the reference lists generated by **DeepRare**, we enlisted 10 associate chief physicians specializing in rare diseases to evaluate the system's outputs on complex cases. A total of 180 cases were randomly sampled from **DeepRare**'s predictions across eight datasets. Each case was independently reviewed by three specialists, and the consensus was calculated as the mean score. We developed a dedicated annotation interface that presented experts with patient information, model-generated diagnostic results, and corresponding reference lists (see Figure 4a). Physicians were asked to assess the accuracy of each reference (including literature, case reports, and websites), with accuracy defined as the reference being both reliable and directly relevant to the model's final diagnostic decision.

Statistical results at the case level (Figure 4b) show an average reference accuracy of 95.4%. At the dataset level (Figure 4c), the system consistently demonstrates high performance across all datasets. Further analysis of references deemed incorrect by physicians revealed two main error categories: **(1) Hallucinated references**, where the system generated plausible but nonexistent URLs in the absence of actual literature links, leading to erroneous webpages; **(2) Irrelevant references**, resulting from incorrect diagnostic conclusions that caused the model to cite sources unrelated to the true disease.

Overall, physician validation confirms the robustness of **DeepRare**'s source attribution, highlighting its potential to substantially streamline the literature and case retrieval process during clinical diagnosis and to enhance diagnostic efficiency for healthcare professionals.

1 Clinical Data Entry

2 Systematic Clinical Inquiry

5 Clinical Report Generation

3 HPO Phenotype Mapping

4 Diagnostic Analysis and Output

Figure 5 | The five-stage DeepRare web application workflow. (1) Clinical data entry: Input of demographic parameters, family history, and clinical manifestations with optional file uploads (medical images, lab reports, VCF files). (2) Systematic clinical inquiry: AI-guided symptom refinement for organ involvement and disease progression. (3) HPO phenotype mapping: Automated terminology standardization with clinician-curated adjustment interface. (4) Diagnostic analysis: Integrated tool orchestration generating evidence-based recommendations. (5) Report downloading: Automated export of structured clinical reports (PDF/Word).

2.7 Web Application

To facilitate adoption by rare disease clinicians and patients, we developed a user-friendly web application interface for **DeepRare**². The platform enables users to input patient demographics, family history, and clinical presentations to obtain diagnostic predictions. The backend architecture processes and structures the model outputs, presenting results through an intuitive and interactive interface optimized for clinical workflow integration. As presented in Figure 5, the diagnostic workflow encompasses five sequential phases:

- **Clinical data entry:** Users input essential patient information, including age, sex, family history, and primary clinical manifestations. The platform supports the upload of supplementary materials such as case reports, diagnostic imaging, laboratory results, or raw genomic VCF files when available.
- **Systematic clinical inquiry:** The system first conducts "detailed symptom inquiry", further investigating information that helps clarify the scope of organ involvement, family genetic history, and symptom progression to narrow the diagnostic range. Users may also choose to skip this step and proceed directly to diagnosis.
- **HPO phenotype mapping:** The platform automatically maps clinical inputs to standardized Human Phenotype Ontology (HPO) terms, with manual curation capabilities allowing clinicians to refine, supplement, or remove assigned phenotypic descriptors.
- **Diagnostic analysis and output:** At this stage, the system executes a comprehensive analysis by invoking various tools and consulting medical literature and case databases to provide diagnostic recommendations and treatment suggestions for physicians. The web frontend renders the output results for user-friendly presentation.
- **Clinical report downloading:** Upon completion of the diagnostic analysis, users can generate comprehensive diagnostic reports that are automatically formatted and exported as PDF or Word documents for integration into electronic health records or clinical documentation.

This system has currently been deployed and tested in hospital settings, assisting rare disease physicians in diagnostic efficiency and clinical decision-making accuracy.

2.8 Ablation Study

In this section, we conduct a comprehensive ablation study on our system design, focusing on central host selection and the effectiveness of introducing various agent designs.

To begin, we evaluate various foundational models as central hosts for **DeepRare**. As illustrated in Figure 6a, we test Claude-3.5-Sonnet, DeepSeek-R1, DeepSeek-V3, GPT-4o, and Gemini-2-Flash across eight datasets, assessing their suitability as the core of our agentic system. According to the results, DeepSeek-V3 outperforms other models on most datasets, except for RareBench (MME), where the Gemini-2-flash-based agentic system achieves the best performance. Overall, the choice of central host LLMs has minimal impact on the results, highlighting the generalization of our system, which is not reliant on specific LLMs.

In Figure 6b, we compare the raw LLMs with their corresponding agentic systems. As shown in the figure, our agentic design significantly improves the performance of the original LLMs, highlighting the necessity of the proposed agentic workflow mechanisms. For instance, with GPT-4o, the average Recall@1 score across the five public datasets improves substantially from 25.60% to 54.67% with 28.56% performance gain, while for DeepSeek-V3, it increases from 26.18% to 56.94% with 29.95% performance gain. This enhancement is consistently observed across all tested LLMs, demonstrating the effectiveness of our approach.

Lastly, in Figure 6c, we show the effectiveness of different agentic components. It illustrates the Recall@1 improvements of each module, including the similar case retrieval, web knowledge, and self-reflection modules, in the DeepRare system (GPT-4o powered) compared to the baseline method (GPT-4o) on the RareBench dataset. As shown in the figure, each module in our agentic system design contributes varying degrees of improvement to diagnostic performance, with particularly pronounced effects on challenging datasets such as MME. Specifically, in MME dataset, the similar case module achieves a 40% improvement, the self-reflection

²<http://raredx.cn/doctor>

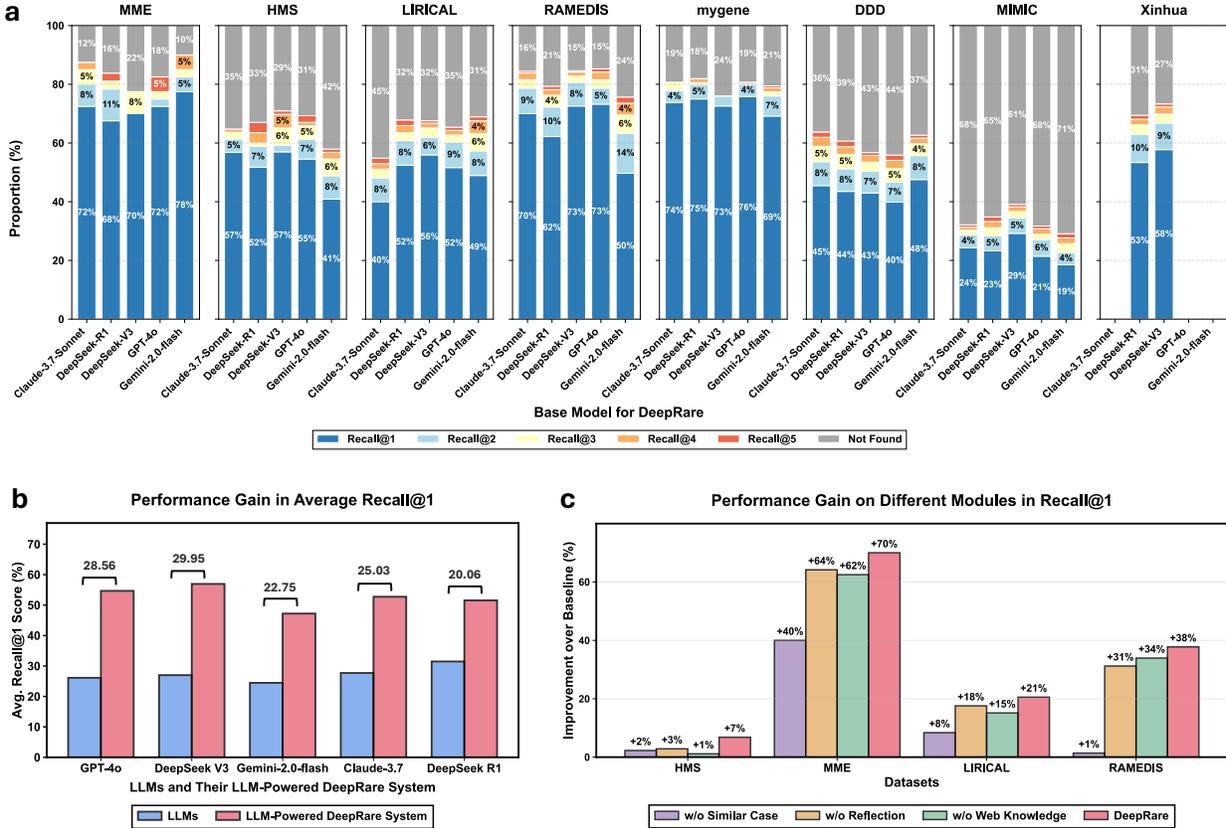


Figure 6 | Ablation Study of DeepRare System. (a) Performance comparison across different LLMs (Claude-3.7-Sonnet, DeepSeek-R1, DeepSeek-V3, GPT-4o, Gemini-2.0-flash) as central hosts on eight rare disease datasets. (b) Performance enhancement comparison between baseline large language models and their powered agentic DeepRare systems. (c) Module-wise contribution analysis on DeepRare system (GPT-4o powered) demonstrating the effectiveness of similar case retrieval, web knowledge integration, and self-reflection components compared to baseline GPT-4o performance.

module provides a 64% enhancement, the web knowledge module contributes a 62% boost, and the overall system demonstrates a 70% improvement in performance.

3 Discussion

In this study, we present **DeepRare**, an LLM-powered agentic system specifically designed for rare disease diagnosis. It can process a wide range of input types commonly encountered in clinical workflows, including chief complaints, genetic data, and detailed clinical phenotypes, thereby supporting clinicians in the timely and accurate diagnosis of rare diseases. A key feature of **DeepRare** is its ability to generate a comprehensive diagnostic reasoning chain, providing transparent and interpretable insights that enhance clinical decision-making. We rigorously evaluated **DeepRare** across diverse datasets spanning multiple sources, disease categories, and medical centers, where it consistently outperformed existing methods, demonstrating both effectiveness and generalizability.

Compared to existing diagnostic tools commonly used in clinical practice, **DeepRare** addresses several critical limitations: (i) Traditional HPO-based systems typically generate candidate disease lists without providing sufficient explanatory context or diagnostic rationale, thereby limiting their clinical applicability; (ii) Some tools that integrate both HPO phenotypes and genetic data remain highly dependent on genetic testing results, rendering them less suitable for initial patient assessments or first-line screening; (iii) Recent advances in large language models have improved clinician usability, but these models are still prone to hallucinations

that undermine diagnostic reliability. **DeepRare** overcomes these challenges by grounding its diagnostic reasoning in verifiable medical evidence, ensuring both interpretability and trustworthiness throughout the diagnostic process.

Our experimental results demonstrate two key achievements: (i) **Superior performance across benchmarks**: **DeepRare** achieved substantial improvements over existing methods across multiple benchmark datasets, including the publicly available RareBench, a rare disease subset of MIMIC-IV-Note, and our in-house Xinhua Hospital dataset, with comprehensive analyses showing consistent outperformance across different medical specialties and input modalities. (ii) **Evidence-based reasoning chains**: Beyond diagnostic accuracy, **DeepRare** provides transparent, step-by-step diagnostic reasoning with verifiable references that significantly reduce clinical decision-making time and minimize patient costs associated with misdiagnosis, as validated through expert clinical assessment.

The clinical implications of **DeepRare** extend beyond diagnostic accuracy to address fundamental challenges in rare disease care delivery. The system’s ability to provide evidence-based reasoning chains with verifiable references could significantly reduce the time required for literature review and case research, enabling clinicians to focus more on patient care rather than information gathering. Furthermore, the system’s consistent performance across different medical specialties suggests its potential as a valuable decision support tool for non-specialist physicians who may encounter rare diseases infrequently. This democratization of rare disease expertise could be particularly impactful in resource-limited settings or regions with limited access to specialized care, potentially reducing healthcare disparities in rare disease diagnosis.

4 Limitations

While **DeepRare** demonstrates strong performance and broad applicability, several areas offer opportunities for further enhancement. First, although the current tool architecture integrates a diverse set of specialized resources and databases, it does not yet encompass the full range of potentially valuable data sources. However, the flexible design of our agentic system and MCP-like interface readily allow for future expansion and integration of additional rare disease knowledge systems and bioinformatics tools, which may further strengthen diagnostic support.

Second, our present knowledge search agent processes phenotypic information in aggregate. While this approach has proven effective, future work could explore more refined and adaptive retrieval mechanisms to further optimize knowledge curation and potentially enhance diagnostic precision.

Finally, although we have developed modules for patient interaction to facilitate information gathering, the lack of suitable validation datasets has so far precluded experimental evaluation of this feature. As such datasets become available, further investigation and iterative improvement of patient interaction capabilities will be a natural next step.

Overall, these represent areas for ongoing development rather than fundamental limitations. Future work will focus on expanding the agentic system framework to encompass rare disease treatment and prognosis prediction, with the goal of evolving **DeepRare** into an even more flexible and comprehensive ecosystem for rare disease management.

5 Methods

We introduce **DeepRare**, an agentic framework designed to support rare disease diagnosis, structured upon a modular, multi-tiered architecture. The system comprises three core components: (i) a central host agent, equipped with a memory bank that integrates and synthesizes diagnostic information while coordinating system-wide operations; (ii) specialized local agent servers, each interfacing with specific diagnostic resource environments through tailored toolsets; and (iii) heterogeneous data sources that provide critical diagnostic evidence, including structured knowledge bases (*e.g.*, research literature, clinical guidelines) and real-world patient data. The architecture of **DeepRare** is described in a top-down manner, beginning with the central host’s core workflow and proceeding through the agent servers to the underlying data sources.

5.1 Problem Formulation

In this paper, we focus on rare disease diagnosis, where the input of a rare disease patient’s case typically consists of two components: **phenotype** and **genotype**, denoted as $\mathcal{I} = \{\mathcal{P}, \mathcal{G}\}$. Either \mathcal{P} or \mathcal{G} (but not both) may be an empty set \emptyset , indicating the absence of the corresponding input. Specifically, the input phenotype may consist of free-text descriptions \mathcal{T} , structured Human Phenotype Ontology (HPO) terms \mathcal{H} , or both. Formally, we define, $\mathcal{P} = (\mathcal{T}, \mathcal{H})$, where either \mathcal{T} or \mathcal{H} may be empty (*i.e.*, \emptyset), indicating the absence of that input modality. The **genotype input** denotes the raw Variant Call Format (VCF) file generated from Whole Exome Sequencing (WES).

Given \mathcal{P} , the goal of the system is to produce: a ranked list of the top K most probable rare diseases, $\mathcal{D} = \{d_1, d_2, \dots, d_K\}$, and a corresponding rationale \mathcal{R} , consisting of evidence-grounded explanations traceable to medical sources such as peer-reviewed literature, clinical guidelines, and similar patient cases. This can be formalized as:

$$\{\mathcal{D}, \mathcal{R}\} = \mathcal{A}(\mathcal{P}), \tag{1}$$

where $\mathcal{A}(\cdot)$ denotes the diagnostic model.

As shown in Figure 7b, our multi-agent system comprises of three main components:

- **A central host with a memory bank** serves as the coordinating brain of the system. The memory bank is initialized as empty and incrementally updated with information gathered by agent servers. Powered by a LLM, the central host integrates historical context from the memory bank to determine the system’s next actions.
- **Multiple agent servers** execute specialized tasks such as phenotype extraction and knowledge retrieval, enabling dynamic interaction with external data sources.
- **Diverse data sources** serve as the external environment, providing crucial diagnostic evidence from PubMed articles, clinical guidelines, publicly available case reports, and other relevant resources.

5.2 Main Workflow

The system operates in two primary stages, orchestrated by the central host: **information collection** and **self-reflective diagnosis**, as illustrated in Figure 7c. For clarity, the specific functionalities of the agent servers involved in each stage are detailed in the following section.

Information Collection

In the information collection stage, the system preprocesses the patient input and invokes specialized agent servers to gather relevant medical evidence from external sources. The process begins with two parallel steps: one focusing on phenotype inputs and the other on genotype data. Subsequently, the central host takes control to facilitate diagnostic decision-making and patient interaction.

Phenotype Information Collection: Given the phenotype input ($\mathcal{P} = (\mathcal{T}, \mathcal{H})$), the system performs the three main sub-steps to collect extra information: **HPO standardization**, **phenotype retrieval**, and **phenotype analysis**.

In HPO standardization, the phenotype extractor a_{hpo} agent server is called, to convert the given free-form

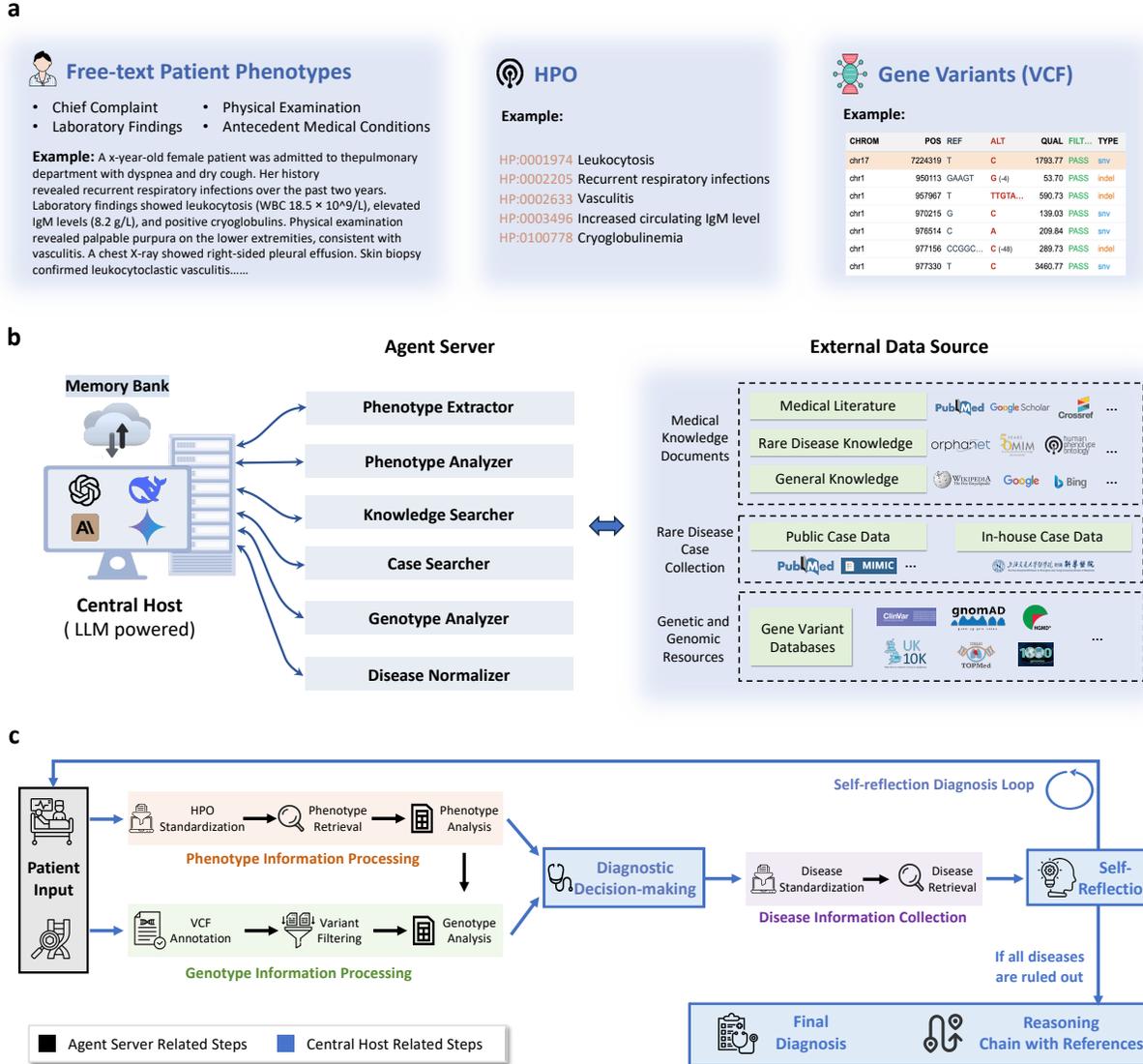


Figure 7 | Overview of the DeepRare system. **a** The input consists of patient free-text information, structured HPO IDs, or any combination of them. **b** The three-level components in RaraDx. Inspired by the MCP, our system can also be analogized to a personal computer system architecture, comprising: (1) a central host with a memory bank for centrally managing and coordinating the system, analogous to the main computer processing system; (2) multiple agent servers to organize tools, execute specific tasks, and interact with the external environment, analogous to auxiliary hardware assistant equipment; (3) comprehensive external data sources, representing a complete external rare-disease diagnostic environment, supporting the entire system by various medical reliable evidence, including medical knowledge and clinical cases. **c** The flowchart of the main workflow of our system illustrates two primary stages, i.e., the information collection stage and the self-reflection diagnosis stage. In the former, the central host actively collects medical support information relevant to the patient. In the latter, the central host performs self-reflection on its diagnostic results. Steps involving the central host are highlighted in blue boxes within the flowchart.

reports \mathcal{T} into a list of standardized entities \mathcal{H} , denoted as:

$$\hat{\mathcal{P}} = \begin{cases} a_{\text{hpo}}(\mathcal{T}), & \text{if } \mathcal{T} \neq \emptyset. \\ \mathcal{H}, & \text{otherwise.} \end{cases} \quad (2)$$

As a result, each patient is now denoted as a set of standardized HPO entities ($\hat{\mathcal{P}}$), that are further treated as the query for phenotype retrieval.

The knowledge searcher ($a_{\text{k-search}}$) and case searcher ($a_{\text{c-search}}$) agent servers are invoked to retrieve supporting documents from the web and relevant cases from an external database, respectively:

$$\mathcal{E}_{\text{hpo}} = a_{\text{k-search}}(\hat{\mathcal{P}}, \mathcal{M}, N) \cup a_{\text{c-search}}(\hat{\mathcal{P}}, \mathcal{M}, N), \quad (3)$$

where \mathcal{E}_{hpo} refers to a unified set of retrieved evidences and N denotes the search depth. Notably, the two search agents will also check the memory bank (\mathcal{M}) to avoid retrieving items that have already been recorded.

Lastly, in phenotype analysis, the agent server integrates various distinct bio-informatics tools to provide a set of diagnostic-related suggestions, for example, identifying diseases that are more likely to be associated with the patient based on their phenotype, denoted as:

$$\mathcal{Y}_{\text{hpo}} = a_{\text{hpo-analyzer}}(\hat{\mathcal{P}}). \quad (4)$$

Till here, we have gathered relevant information on the phenotype by exploring the web, or database with similar cases, and multiple existing bioinformatics analysis tools, collectively denoted as ($\mathcal{E}_{\text{hpo}}, \mathcal{Y}_{\text{hpo}}$) and update them into the system memory bank \mathcal{M} , denoted as:

$$\mathcal{M} \leftarrow \mathcal{M} \cup (\mathcal{E}_{\text{hpo}}, \mathcal{Y}_{\text{hpo}}), \quad (5)$$

Following the phenotype analysis, the central host generates a tentative diagnosis based on the available phenotype information:

$$\mathcal{D}' = \mathcal{A}_{\text{host}}(\hat{\mathcal{P}}, \mathcal{M} \mid \langle \text{prompt} \rangle_4) \quad (6)$$

where $\langle \text{prompt} \rangle_4$ is a diagnosis-related prompt instruction to drive the central host. The output \mathcal{D}' is an initialized rare disease list.

Genotype Information Collection: In parallel to the phenotype analysis, the system will also collect external information relevant to the genotypes, if provided ($\mathcal{G} \neq \emptyset$). It also consists of three main sub-steps: **VCF annotation**, **variant ranking**, and **synthetic analysis**. The first two steps are conducted by the genotype analyzer agent server, and the last one is processed by the central host.

In VCF annotation, the goal is to annotate the raw VCF input, which often contains thousands of gene variants using various genomic databases. This process enriches each variant with comprehensive functional annotations, population frequencies, and pathogenicity predictions. Subsequently, variant ranking is performed to prioritize variants based on their potential clinical significance. This step applies scoring algorithms that consider multiple factors, including functional impact, allele frequencies, conservation scores, and predicted pathogenicity:

$$\hat{\mathcal{G}} = a_{\text{geno-analyzer}}(\mathcal{G}) \quad (7)$$

where $\hat{\mathcal{G}}$ denotes the ranked set of variants ordered by their clinical relevance scores.

Finally, in synthetic analysis, the central host uses LLM to interpret the ranked variants and patient information, providing comprehensive variant interpretation, gene-phenotype association predictions, and inheritance pattern analysis:

$$\mathcal{D}'' = \mathcal{A}_{\text{host}}(\hat{\mathcal{G}}, \mathcal{M}, \mathcal{D}', N \mid \langle \text{prompt} \rangle_5) \quad (8)$$

where $\langle \text{prompt} \rangle_5$ is synthetic analysis prompt. If genotype data is not available ($\mathcal{G} = \emptyset$), the system uses the phenotype-only diagnosis:

$$\mathcal{D}'' = \mathcal{D}', \text{ if } \mathcal{G} = \emptyset \quad (9)$$

where \mathcal{D}'' represents the updated diagnosis list after incorporating genotype information (when available). The results are then consolidated and updated in the system memory bank:

$$\mathcal{M} \leftarrow \mathcal{M} \cup \mathcal{D}'' \quad (10)$$

Self-reflective Diagnosis

At this stage, the central host takes entire system control and proceeds to self-reflective diagnosis, which attempts to make a diagnosis based on all previously collected information.

Specifically, the central host will make a tentative diagnosis decision-making, defined as:

$$\mathcal{D}' = \mathcal{A}_{\text{host}}(\hat{\mathcal{I}}, \mathcal{M} \mid \langle \text{prompt} \rangle_4) \quad (11)$$

where $\hat{\mathcal{I}} = \{\hat{\mathcal{P}}, \hat{\mathcal{G}}\}$ denotes the preprocessed patient input and $\langle \text{prompt} \rangle_4$ is a diagnosis-related prompt instruction to drive the central host. The output \mathcal{D}' is an initialized rare disease list.

Then the system executes self-reflection that double-checks the disease list by collecting additional knowledge related to the predicted disease list from the Internet. This process involves two sub-steps: **disease standardization** and **disease retrieval**. In disease standardization, the disease normalizer agent server is invoked here, which converts the disease name list into items from Orphanet or OMIM. Subsequently, in disease retrieval, the knowledge searcher ($a_{\text{k-search}}$) is called, which treats the standardized disease items as queries and retrieves the relevant knowledge documents for each disease. The process can be formulated as:

$$\mathcal{M} \leftarrow \mathcal{M} \cup \mathcal{E}_{\text{disease}}, \quad \mathcal{E}_{\text{disease}} = a_{\text{k-search}}(a_{\text{d-norm}}(\mathcal{D}''), \mathcal{M}), \quad (12)$$

where $\mathcal{E}_{\text{disease}}$ denotes the collected disease-related knowledge and will also be stored into the memory bank.

After acquiring the external disease-specific knowledge, the central host self-reflects on the correctness of the predicted diseases, by synthesizing all collected information. This process is formulated as:

$$\mathcal{D} = \mathcal{A}_{\text{host}}(\mathcal{D}'', \hat{\mathcal{I}}, \mathcal{M} \mid \langle \text{prompt} \rangle_6), \quad (13)$$

where $\langle \text{prompt} \rangle_6$ is a self-reflection prompt, and $\mathcal{D} = \{d_1, d_2, \dots\}$ represents the ranked list of possible rare diseases, ordered by their likelihood. Notably, if $\mathcal{D} = \emptyset$, *i.e.*, all proposed rare diseases are ruled out during self-reflection, the system will return to the beginning and increase N by ΔN , re-collect new patient-wise information, and iterate through the entire program workflow until $\mathcal{D} \neq \emptyset$ is satisfied.

Once the system passes the former self-reflection step, the central host will further synthesize the collected information, and provide traceable transparent rationale explanations:

$$\{\mathcal{D}, \mathcal{R}\} = \mathcal{A}_{\text{host}}(\mathcal{D}, \hat{\mathcal{I}}, \mathcal{M} \mid \langle \text{prompt} \rangle_7), \quad (14)$$

where \mathcal{R} denotes the rationale explanation for each output rare disease, organized as free-text by the central host. This ensures that the final diagnosis is not only accurate but also interpretable, offering users a clear and auditable justification for the predicted diseases.

5.3 Agent Servers

Agent servers form the second tier of our **DeepRare system**, each manages one or multiple specific tools, interacting with a specialized working environment to gather evidence from external data sources. In specific, the following agent servers are utilized in our system: **phenotype extractor**, **disease normalizer**, **knowledge searcher**, **case searcher**, **phenotype analyzer**, and **genotype analyzer**.

Phenotype Extractor

The clinical rare disease diagnosis procedure requires converting patients' phenotype consultation records (\mathcal{T}) into standardized HPO items. Specifically, we extract potential phenotype candidates and modify the

phenotype name by prompting an LLM:

$$\mathcal{H} = \Phi_{\text{LLM}}(\Phi_{\text{LLM}}(\mathcal{T} \mid \langle \text{prompt} \rangle_8) \mid \langle \text{prompt} \rangle_9) \quad (15)$$

where $\mathcal{H} = \{h_1, h_2, \dots\}$ denotes the set of extracted, unnormalized HPO candidate entities, and $\langle \text{prompt} \rangle_8$ $\langle \text{prompt} \rangle_9$ represents the corresponding prompt instruction. Here, the two-step reasoning process is designed to extract more accurate phenotypic descriptions with LLM assistance, significantly reducing the probability of errors in the subsequent step.

Subsequently, we perform named-entity normalization leveraging BioLORD [50], a BERT-based text encoder, to map these candidate entities to standardized HPO terms. Specifically, we compute the cosine similarity between the text embeddings of the predicted entity name and all standardized HPO term names. The top-matching HPO term is then selected to represent the entity. Notably, if no HPO term achieves a cosine similarity of 0.8 or above, the entity is discarded.

Disease Normalizer

During the diagnostic process, free-text diagnostic diseases are mapped to standardized Orphanet or OMIM items, as more precise keywords for subsequent searches. Similar to that in the phenotype extractor, we use BioLORD [50] to perform named-entity normalization, by computing the cosine similarity between the text embeddings of the predicted disease name and all standardized disease names listed in the Orphanet or OMIM. The top-matched standardized disease name is then used, and if all standardized disease names cannot match the predicted term (cosine similarity less than 0.8), the predicted disease will be discarded.

Knowledge Searcher

The knowledge searcher is tasked with real-time knowledge document searching, interacting with external medical knowledge documents and the Internet, supporting the diagnosis system with latest rare disease knowledge.

While it is invoked, it will perform two distinct search modules with multiple searching tools, on a specific search query (\mathcal{Q}), for example, HPO or predicted diseases:

- **General web search:** This part executes the general search engines, including Bing [51], Google [52], and DuckDuckGo [53]. We will call them one by one, following the listed order. Each time, the top- N web pages (with a default value of $N = 5$) will be retrieved. Specifically, Bing³ is accessed through automated browser simulation using Selenium, while Google and DuckDuckGo are queried via their official APIs^{4 5}. If a search engine successfully completes the execution, the process will stop immediately.
- **Medical domain search:** Considering that some professional medical-specific web pages may not be ranked highly in general search engines, this part retrieves information from well-known medical databases. The following search engines are considered:
 - **Up-to-dated Academic literatures**, including PubMed [40](accessed via PubMedRetriever from langchain_community⁶), and Crossref [54] (queried through official API⁷);
 - **Rare disease-specific knowledge bases** such as Orphanet [55], OMIM [56], and HPO [57] (all utilizing offline knowledge bases and accessed through retrieval mechanisms);
 - **General medical knowledge repositories:** Wikipedia [58] (accessed via WikipediaRetriever from langchain_community⁸) and MedlinePlus⁹ [59] (accessed through automated browser simulation using Selenium).

³<https://www.bing.com>

⁴<https://developers.google.com/custom-search/v1/overview>

⁵<https://api.duckduckgo.com>

⁶https://python.langchain.com/api_reference/community/retrievers/langchain_community_retrievers_pubmed.PubMedRetriever.html

⁷<https://api.crossref.org/swagger-ui/index.html>

⁸https://python.langchain.com/api_reference/community/retrievers/langchain_community_retrievers_wikipedia.WikipediaRetriever.html

⁹<https://medlineplus.gov/>

Similarly, while searching academic papers and rare disease-specific knowledge bases, we retrieve the top- N web pages (with N defaulting to 5) from each source. The search engines are queried one by one, and the process stops upon successful execution.

These tools retrieve web pages and return them to the knowledge searcher, and are summarized by the agent server with a lightweight language model (GPT-4o-mini by default), to simultaneously extract key information and filter relevant content. This integrated processing pipeline can be formalized as:

$$\mathcal{R} = \Phi_{\text{LLM}}(\text{document}, \mathcal{Q} \mid \langle \text{prompt} \rangle_{10}), \quad (16)$$

where \mathcal{R} represents the processed output for each retrieved document, \mathcal{Q} denotes the given search queries, and $\langle \text{prompt} \rangle_{10}$ is the unified prompt instruction that simultaneously governs both summarization and relevance filtering. The system employs a binary classification approach: medical-related documents are retained and translated into the target language, while non-medical content is rejected with the output "*Not a medical-related page*".

Case Searcher

Inspired by clinical practice, where physicians often refer to publicly discussed cases when faced with rare or challenging patients, the case searcher agent is designed to explore an external case bank. Each patient in the database is represented as a list of HPO terms, transforming case search into an HPO similarity matching problem. Using an input HPO list (\mathcal{H}) from the query case, we implement a two-step retrieval method to interact with this external database: (i) **Initial retrieval**: we employ OpenAI’s text-embedding model (`text-embedding-3-small`) to encode both the query HPO list and each candidate patient’s HPO representation into dense vector embeddings. The embeddings for all candidate patients in the case database have been pre-computed and stored using the same embedding model. We then identify the top-50 candidate patients based on cosine similarity between these embeddings. (ii) **Re-ranking**: We further re-rank these candidates using MedCPT-Cross-Encoder [60], a BERT-based model specifically trained on PubMed search logs for biomedical information retrieval. This model computes refined cosine similarity scores between the query case’s HPO profile and each candidate’s HPO profile, leveraging domain-specific medical knowledge to improve matching accuracy.

We also evaluated alternative retrieval strategies, including single-stage methods with different embedding models such as BioLORD and MedCPT, as well as traditional approaches like BM25 [61], discussed in the results section. Experimental findings indicate that the two-stage retrieval approach outperforms all alternatives, optimizing both computational efficiency and clinical relevance of the retrieved cases.

Similar to the knowledge searcher, after receiving the similar cases, the case searcher will further assess their relevance to prevent misdiagnosis from irrelevant cases, powered by the lightweight language model:

$$r_{\text{case}} = \Phi_{\text{LLM}}(\text{Case}, \mathcal{H} \mid \langle \text{prompt} \rangle_{11}), \quad (17)$$

where $r_{\text{case}} \in \{\text{True}, \text{False}\}$, a binary scalar that indicates whether the case is related to the given HPO list, and $\langle \text{prompt} \rangle_{11}$ is the corresponding prompt instruction. Consistency is maintained by employing the same LLM architecture used in the diagnostic process for this assessment.

Phenotype Analyzer

This agent server controls various professional diagnosis tools, that are developed for phenotype analysis. By integrating the analysis results from these tools into the overall diagnostic pipeline, our system is enabled to incorporate more professional and comprehensive suggestions. Specifically, given the patient HPO list (\mathcal{H}), the following tools are used:

- **PhenoBrain** [26]: This is a tool for HPO analysis, that takes structured HPO items as input ($\hat{\mathcal{H}}$), and output 5 potential rare disease suggestions. We adopt it by calling its official API¹⁰.

¹⁰https://github.com/xiaohaomao/timgroup_disease_diagnosis/tree/main/PhenoBrain_Web_API

- **PubcaseFinder** [27]: It performs HPO-wise diagnostic analysis by matching the most similar public cases from PubMed case reports. Similarly, it takes the structured HPO items ($\hat{\mathcal{H}}$) as input and returns top-5 potential rare disease suggestions, each with a confidence score. We access it via its official API¹¹.
- **Zero-shot LLM Inference**: We additionally employ LLMs to perform zero-shot preliminary reasoning. Given the extensive knowledge base acquired during LLM training, these models can often suggest candidate diagnoses that conventional diagnostic tools might overlook. Specifically, this approach takes the structured HPO items ($\hat{\mathcal{H}}$) as input and returns the top-5 potential rare disease candidates under Prompt 12.

Genotype Analyzer

Similar to the phenotype analyzer, the genotype analyzer is tasked with performing professional genotype analysis by calling existing tools.

For patient genomic variant files (aligned to GRCh37 reference genome), we initially subjected the HPO phenotype terms \mathcal{H} and corresponding VCF files \mathcal{G} to comprehensive annotation and prioritization analysis using the **Exomiser** [49] framework, with configuration parameters detailed in the supplementary materials 11.3, which is configed to integrate multiple data sources and analytical steps: population frequency filtering utilizing databases including gnomAD [62], 1000 Genomes Project [63], TOPMed [64], UK10K [65], and ESP [66] across diverse populations; pathogenicity assessment through PolyPhen-2 [67], SIFT [68], and MutationTaster [69] prediction algorithms; variant effect filtering to retain coding and splice-site variants while excluding intergenic and regulatory variants; inheritance mode analysis supporting autosomal dominant/recessive, X-linked, and mitochondrial patterns; and gene-disease association prioritization through OMIM [56] and HiPhive [70] algorithms that leverage cross-species phenotype data.

The Exomiser output is ranked according to the composite `exomiser_score`, from which we selected the top-n candidate genes while preserving essential metadata including OMIM identifiers, `phenotype_score`, `variant_score`, statistical significance (`p_value`), detailed `variant_info`, ACMG pathogenicity classifications, ClinVar annotations, and associated disease phenotypes. The curated genomic annotations were subsequently transmitted to the Central Host for downstream processing and integration.

All outputs from the specialized tools are then transformed into free texts by the agent server, that can be seamlessly combined with the LLM-based central host or other LLM-driven tools. This is achieved by employing a predefined templates tailored to each tool’s specific output format, such as “[Tool Name] identified [Disease]” (with confidence scores included when available) for disease predictions.

5.4 External Data Sources

The external data sources form the third tier of our **DeepRare** framework, providing a comprehensive external environment for tool interaction. These diverse, rare disease-related information sources support the system with professional medical knowledge, we specifically consider the medical-focused databases.

Medical Literature. Scientific publications are essential for evidence-based diagnosis, especially for rapidly evolving rare diseases. DeepRare accesses peer-reviewed literature through:

- **PubMed database** [40]: The world largest database of biomedical literature containing over 34 million papers.
- **Google Scholar** [71]: A broad academic search engine covering publications across diverse sources.
- **Crossref** [54]: A comprehensive metadata database that enables seamless access to scholarly publications and related fields through persistent identifiers and open APIs.

Rare Disease Knowledge Sources. Curated repositories that aggregate structured information about rare diseases:

¹¹<https://pubcasefinder.dbcls.jp/api>

- **Orphanet** [55]: Comprehensive information for over 6,000 rare diseases, including descriptions, genetics, epidemiology, diagnostics, and treatments, etc.
- **OMIM** (Online Mendelian Inheritance in Man) [56]: A catalog of human genes and genetic disorders, documenting over 17,000 genes and their associated phenotypes.
- **Human Phenotype Ontology** [72]: A standardized vocabulary of phenotypic abnormalities in human diseases, containing over 18,000 terms and more than 156,000 hereditary disease annotations.

General Knowledge Sources. Broad clinical resources that provide contextual understanding:

- **MedlinePlus** [59]: A US National Library of Medicine resource providing reliable, up-to-date health information for patients and clinicians.
- **Wikipedia** [58]: General encyclopedia entries on all general knowledge, including medical conditions and rare diseases.
- **Online websites:** Resources accessible through search engines that provide up-to-date information, including medical news portals, patient advocacy groups, research institution websites, and clinical trial registries that may contain the latest developments not yet published in scholarly literature.

Case Collection. A large-scale case repository is constructed from multiple data sources to serve as the database for the case search agent server, with a subset of the data reserved as a test set to validate model performance. For large datasets such as MIMIC-IV-Note and Xinhua Hosp. Dataset, we perform data splitting at a 4:1 ratio, where earlier cases serve as reference similar cases; for smaller datasets such as RareBench, MyGene2 and DDD, we employ leave-one-out cross-validation.

- **RareBench** [8] is a benchmark designed to systematically evaluate LLM capabilities across four critical dimensions in rare disease analysis. We utilize Task 4 (Differential Diagnosis among Universal Rare Diseases), specifically its public subset comprising 1,114 patient cases collected from four open datasets: MME (The Matchmaker Exchange), HMS, LIRICAL, and RAMEDIS. MME and LIRICAL cases are extracted from published literature and manually verified. HMS contains data from the outpatient clinic at Hannover Medical School in Germany. RAMEDIS comprises rare disease cases autonomously submitted by researchers.
- **Mygene2** [37], a data-sharing platform connecting families with rare genetic conditions, clinicians, and researchers, provided additional data. We use preprocessed data (146 patients spanning 55 MONDO diseases) from [38], which extracted phenotype-genotype information as of May 2022, limited to patients with confirmed OMIM disease identifiers and single candidate genes to ensure diagnostic accuracy.
- **DDD** (the Deciphering Developmental Disorders Study) [73] data were obtained from the Gene2Phenotype (G2P) project, which curates gene-disease associations for clinical interpretation. We downloaded phenotype terms and associated gene sets from the G2P database¹² in May 2025. After preprocessing to remove cases with missing diagnostic results or phenotypes, the final DDD cohort comprised 2,283 cases.
- **MIMIC-IV-Note** [46] contains 331,794 de-identified discharge summaries from 145,915 patients admitted to Beth Israel Deaconess Medical Center in Boston, Massachusetts. Since our focus is exclusively on rare diseases, we first determine whether the case involved a rare disease, by prompting GPT-4o with the ICD-10 [74] codes associated with each note. Confirmed cases were mapped to our rare disease knowledge base using a methodology similar to disease normalization, while unmapped cases were discarded, resulting in a final dataset of 9,185 records.
- **Xinhua Hosp. Dataset (in-house)** encompasses all rare disease diagnostic records from 2014 to 2025, totaling 352,424 entries. Using a procedure similar to our MIMIC processing workflow, we applied GPT-4o and vector matching to eliminate records without definitive diagnoses or significant data gaps. We also consolidated multiple consultations for the same patient, resulting in a curated dataset of 5,820 records.

¹²https://ftp.ebi.ac.uk/pub/databases/gene2phenotype/G2P_data_downloads/2025_05_28/

- **PMC-Patients** [75, 41] comprises 167,000 patient summaries extracted from case reports in PubMed [40]. The RareArena GitHub [41] repository has processed this dataset with GPT-4o for rare disease screening. We therefore utilized their preprocessed dataset, which contains 69,759 relevant records.

Gene Variant Databases. Specialized repositories that support the analysis of genetic findings in rare disease diagnosis:

- **ClinVar** [76]: A freely accessible database containing 1.7 million interpretations of clinical significance for genetic variants, with particular value for identifying pathogenic mutations in rare disorders.
- **gnomAD** (Genome Aggregation Database) [62]: A resource of population frequency data for genetic variants from over 140,000 individuals, essential for distinguishing rare pathogenic variants from benign population polymorphisms.
- **HGMD** (Human Gene Mutation Database) [77]: A collection of published germline mutations in nuclear genes associated with human inherited disease.
- **1000 Genomes Project** [63]: A database of human genetic variation across diverse populations worldwide.
- **ExAC** (Exome Aggregation Consortium) [78]: A database of exome sequence data from over 60,000 individuals.
- **TOPMed** (Trans-Omics for Precision Medicine) [64]: A national heart, lung, blood, and sleep disorders research program providing whole-genome sequencing data from over 180,000 individuals.
- **UK10K** [65]: A British genomics project providing population-specific variant frequencies for the UK population through whole-genome and exome sequencing of approximately 10,000 individuals.
- **ESP** (NHLBI Exome Sequencing Project) [66]: A project focused on exome sequencing of individuals with heart, lung, and blood disorders, providing population frequency data stratified by ancestry.

5.5 Clinical Evaluation Dataset Curation

In this section, we will introduce the curation procedure of the two proposed evaluation datasets from the clinical centers, *i.e.*, **MIMIC-IV-Rare** and **Xinhua Hosp.** datasets.

The MIMIC-IV-Note dataset comprised 331,794 de-identified discharge summaries from 145,915 patients sourced from public repositories, while the Xinhua Hospital dataset contained 352,425 outpatient and emergency records from 42,248 patients specializing in genetic diseases, which is an in-house clinical data.

A systematic data preprocessing pipeline was implemented to ensure data quality and relevance. For the MIMIC-IV-Note dataset, we applied a two-stage exclusion process: first, cases without rare disease diagnoses were filtered out ($n = 318,976$ excluded), where rare disease classification was determined using a LLM under prompt X guidelines. Subsequently, records with incomplete patient information were removed ($n = 3,633$ excluded), with information completeness defined as the ability to correctly extract HPO entities that could be successfully matched to the HPO database. This filtering process resulted in 9,185 cases. Similarly, the Xinhua Hospital dataset underwent parallel filtering using identical criteria, excluding 28,150 cases without rare disease diagnoses and 8,278 cases with incomplete information, yielding 5,820 cases.

Subsequently, a time-based allocation strategy was employed to partition the data into evaluation and reference sets. Recent cases were designated for testing purposes, while historical cases were allocated to similar case libraries for retrieval-based analysis. This allocation resulted in the MIMIC-IV Test Set ($n = 1,875$ cases) and MIMIC-IV Similar Case Library ($n = 7,310$ cases), alongside the Xinhua Test Set ($n = 975$ cases) and Xinhua Similar Case Library ($n = 4,845$ cases).

It should be noted that these datasets, derived from authentic clinical records, inherently contain heterogeneous and potentially noisy phenotypic information, including patient-reported symptoms, post-operative complications, multiple consultation entries, and incomplete documentation. This real-world complexity significantly increases the diagnostic challenge compared to curated datasets, thereby providing a more rigorous evaluation framework for clinical decision support systems in rare disease diagnosis.

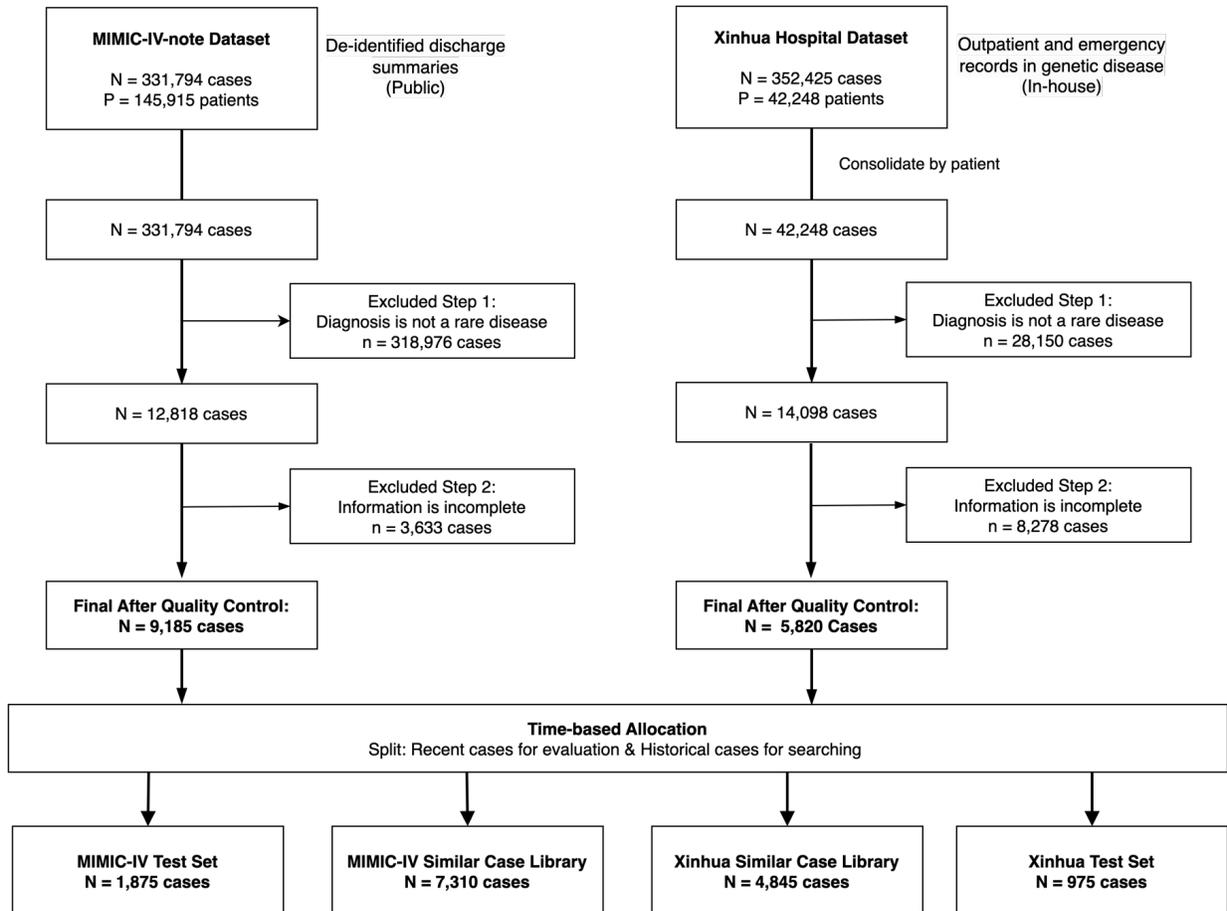


Figure 8 | Cohort curation pipeline and allocation strategy for MIMIC-IV-note and Xinhua Hospital datasets. Left: MIMIC-IV-note dataset including 331,794 cases, with 9,185 remaining after exclusions and divided into test ($n = 1,875$) and library ($n = 7,310$) sets. Right: Xinhua Hospital dataset including 352,425 cases, with 5,820 remaining after exclusions and divided into test ($n = 975$) and library ($n = 4,845$) sets. Both datasets underwent rare disease checks and information completeness filtering.

5.6 Baselines

In this section, we introduce the compared baselines in detail, covering specialized diagnostic methods, latest LLMs, and other agentic systems.

Specialized diagnostic methods:

- **PhenoBrain**¹³ [26]: Takes free-text or structured HPO items as input and suggests top potential rare diseases by an ensembling method integrating the result of a graph-based Bayesian method (PPO) and two machine learning methods (CNB and MLP) via its API.
- **PubcaseFinder**¹⁴ [27]: An website that can extract free-text input first and analyze HPO items by matching similar cases from PubMed reports, returning top potential rare disease suggestions with confidence scores, accessible via its API.

Latest LLMs:

- **GPT-4o** [79]: A closed-source model (version identifier: gpt-4o-2024-11-20) developed by OpenAI. The model was released in May 2024.

¹³<http://www.phenobrain.cs.tsinghua.edu.cn>

¹⁴<https://pubcasefinder.dbcls.jp/>

- **DeepSeek-V3** [80]: An open-source model (version identifier: deepseek-ai/DeepSeek-V3) with 671 billion parameters. It was trained on 14.8 trillion tokens and released in December 2024.
- **Gemini-2.0-flash** [81]: A closed-source model (version identifier: gemini-2.0-flash) developed by Google. This model was released in December 2024.
- **Claude-3.7-Sonnet** [31]: A closed-source model (version identifier: claude-3-7-sonnet) developed by Anthropic. It features a unique “hybrid reasoning” mechanism that allows it to switch between fast responses and extended thinking for complex tasks. This model was released in February 2025.
- **OpenAI-o3-mini** [32]: A closed-source model (version identifier: o3-mini-2025-01-31). This model was officially released in January 2025.
- **DeepSeek-R1** [29]: An open-source large language model (version identifier: deepseek-ai/DeepSeek-R1) with 671 billion parameters. The model was publicly released in January 2025.
- **Gemini-2.0-FT** [30]: A closed-source model (version identifier: gemini-2.0-flash-thinking-exp-01-21). Its training data encompasses information up to June 2024, and it was released in January 2025.
- **Claude-3.7-Sonnet-thinking** [31]: This is an extended version of Claude-3.7-sonnet that provides transparency into its step-by-step thought process. It is a closed-source reasoning model (version identifier: claude-3-7-sonnet-20250219-thinking), publicly released in January 2025.
- **Baichuan-M1** [33]: An open-source domain-specific model (version identifier: baichuan-inc/Baichuan-M1-14B-Instruct) designed specifically for medical applications, distinguishing it from the general-purpose LLMs above. This model comprises 14 billion parameters and was released in January 2025.
- **MMedS-Llama 3** [34]: An open-source domain-specific model (version identifier: Henrychur/MMedS-Llama-3-8B) specialized for the medical domain and released in January 2025. Built upon the Llama-3 architecture with extensive medical domain adaptation, this model comprises 8 billion parameters.

Other Agentic Systems:

- **MDAgents** [35]: A multi-agent architecture that adaptively orchestrates single or collaborative LLM configurations for medical decision-making through a five-phase methodology: Complexity Checking, Expert Recruitment, Initial Assessment, Collaborative Discussion, and Review and Final Decision.
- **DeepSeek-V3-Search** [36, 80]: An LLM agent framework augmented with internet search through Volcano Engine’s platform and web browser plugin.

Genetic Analysis Tools:

- **Exomiser** [49]: A variant prioritization tool (version identifier: 14.1.0 2024-11-14) that combines genomic variant data with HPO phenotype terms to identify disease-causing variants in rare genetic diseases. It integrates population frequency, pathogenicity prediction, and phenotype-gene associations to rank candidate variants.

References

- [1] Rodolfo Valdez, Lijing Ouyang, and Julie Bolen. Public health and rare diseases: oxymoron no more. *Preventing chronic disease*, 13:E05, 2016.
- [2] Stéphanie Nguengang Wakap, Deborah M Lambert, Annie Olry, Charlotte Rodwell, Charlotte Gueydan, Valérie Lanneau, Daniel Murphy, Yann Le Cam, and Ana Rath. Estimating cumulative point prevalence of rare diseases: analysis of the orphanet database. *European journal of human genetics*, 28(2):165–173, 2020.
- [3] Hope for rare diseases. *The Lancet*, 404(10464):1701, 2024.
- [4] The landscape for rare diseases in 2024. *The Lancet Global Health*, 12(3):e341, 2024.
- [5] Arrigo Schieppati, Jan-Inge Henter, Erica Daina, and Anita Aperia. Why rare diseases are an important medical and social issue. *The Lancet*, 371(9629):2039–2041, 2008.
- [6] Xuanzhong Chen, Ye Jin, Xiaohao Mao, Lun Wang, Shuyang Zhang, and Ting Chen. Rareagents: Autonomous multi-disciplinary team for rare disease diagnosis and treatment. *arXiv preprint arXiv:2412.12475*, 2024.
- [7] Xiaohong Liu, Hao Liu, Guoxing Yang, Zeyu Jiang, Shuguang Cui, Zhaoze Zhang, Huan Wang, Liyuan Tao, Yongchang Sun, Zhu Song, et al. A generalist medical language model for disease diagnosis assistance. *Nature Medicine*, pages 1–11, 2025.
- [8] Xuanzhong Chen, Xiaohao Mao, Qihan Guo, Lun Wang, Shuyang Zhang, and Ting Chen. Rarebench: Can llms serve as rare diseases specialists? In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 4850–4861, 2024.
- [9] Michael F Wangler, Shinya Yamamoto, Hsiao-Tuan Chao, Jennifer E Posey, Monte Westerfield, John Postlethwait, Undiagnosed Diseases Network (UDN), Philip Hieter, Kym M Boycott, Philippe M Campeau, et al. Model organisms facilitate rare disease diagnosis and therapeutic research. *Genetics*, 207(1):9–27, 2017.
- [10] Daniel McDuff, Mike Schaekermann, Tao Tu, Anil Palepu, Amy Wang, Jake Garrison, Karan Singhal, Yash Sharma, Shekoofeh Azizi, Kavita Kulkarni, et al. Towards accurate differential diagnosis with large language models. *Nature*, pages 1–7, 2025.
- [11] Sebastian Lunke, Sophie E Bouffler, Chirag V Patel, Sarah A Sandaradura, Meredith Wilson, Jason Pinner, Matthew F Hunter, Christopher P Barnett, Mathew Wallis, Benjamin Kamien, et al. Integrated multi-omics for rapid rare disease diagnosis on a national scale. *Nature medicine*, 29(7):1681–1691, 2023.
- [12] Kristin D Kernohan and Kym M Boycott. The expanding diagnostic toolbox for rare genetic diseases. *Nature Reviews Genetics*, 25(6):401–415, 2024.
- [13] The Lancet Digital Health. Exploring electronic health records to study rare diseases, 2025.
- [14] Genomics Education Programme. The diagnostic odyssey in rare disease — knowledge hub, February 2024. Last reviewed: 13 February 2024, Next review due: 13 February 2026.
- [15] William L Macken, Micol Falabella, Caroline McKittrick, Chiara Pizzamiglio, Rebecca Ellmers, Kelly Eggleton, Cathy E Woodward, Yogen Patel, Robyn Labrum, et al. Specialist multidisciplinary input maximises rare disease diagnoses from whole genome sequencing. *Nature communications*, 13(1):6324, 2022.
- [16] Leen Khalife, Rachel Gottlieb, Tara Daly, Xiaoting Ma, Asma Rashid, Bridget Funk, Emanuela Gussoni, Christina Hung, and Olaf Bodamer. Multidisciplinary clinical and translational approach for optimizing management for complex and rare conditions using kabuki syndrome as example. *Rare*, 1:100008, 2023.

- [17] Hugh JS Dawkins, Ruxandra Draghia-Akli, Paul Lasko, Lilian PL Lau, Anneliene H Jonker, Christine M Cutillo, Ana Rath, Kym M Boycott, Gareth Baynam, Hanns Lochmüller, et al. Progress in rare diseases research 2010–2016: an irdirc perspective. *Clinical and translational science*, 11(1):11, 2017.
- [18] LangChain. Langchain’s suite of products supports developers along each step of the llm application lifecycle, n.d. Accessed: 2025-04-22.
- [19] Anthropic. Building effective agents, n.d. Accessed: 2025-04-22.
- [20] Jianing Qiu, Kyle Lam, Guohao Li, Amish Acharya, Tien Yin Wong, Ara Darzi, Wu Yuan, and Eric J Topol. Llm-based agentic systems in medicine and healthcare. *Nature Machine Intelligence*, 6(12):1418–1420, 2024.
- [21] James Zou and Eric J Topol. The rise of agentic ai teammates in medicine. *The Lancet*, 405(10477):457, 2025.
- [22] Yongju Lee, Dyke Ferber, Jennifer E Rood, Aviv Regev, and Jakob Nikolas Kather. How ai agents will change cancer research and oncology. *Nature Cancer*, pages 1–3, 2024.
- [23] Yashar Talebirad and Amirhossein Nadiri. Multi-agent collaboration: Harnessing the power of intelligent llm agents. *arXiv preprint arXiv:2306.03314*, 2023.
- [24] Qiaoyu Zheng, Chaoyi Wu, Pengcheng Qiu, Lisong Dai, Ya Zhang, Yanfeng Wang, and Weidi Xie. Can modern llms act as agent cores in radiology environments? *arXiv preprint arXiv:2412.09529*, 2024.
- [25] Anthropic. Model context protocol, n.d. Accessed: 2025-04-22.
- [26] Xiaohao Mao, Yu Huang, Ye Jin, Lun Wang, Xuanzhong Chen, Honghong Liu, Xinglin Yang, Haopeng Xu, Xiaodong Luan, Ying Xiao, et al. A phenotype-based ai pipeline outperforms human experts in differentially diagnosing rare diseases using ehers. *npj Digital Medicine*, 8(1):68, 2025.
- [27] Toyofumi Fujiwara, Jae-Moon Shin, and Atsuko Yamaguchi. Advances in the development of pubcasefinder, including the new application programming interface and matching algorithm. *Human Mutation*, 43(6):734–742, 2022.
- [28] R OpenAI et al. Gpt-4 technical report. *ArXiv*, 2303:08774, 2023.
- [29] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- [30] Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, et al. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*, 2024.
- [31] Anthropic Team. Introducing the next generation of claude, 2024. Accessed on March 4, 2024.
- [32] OpenAI. Openai o3 mini, n.d. Accessed: 2025-02-23.
- [33] Bingning Wang, Haizhou Zhao, Huozhi Zhou, Liang Song, Mingyu Xu, Wei Cheng, Xiangrong Zeng, Yupeng Zhang, Yuqi Huo, Zecheng Wang, et al. Baichuan-m1: Pushing the medical capability of large language models. *arXiv preprint arXiv:2502.12671*, 2025.
- [34] Chaoyi Wu, Pengcheng Qiu, Jinxin Liu, Hongfei Gu, Na Li, Ya Zhang, Yanfeng Wang, and Weidi Xie. Towards evaluating and building versatile large language models for medicine. *npj Digital Medicine*, 8(1):58, 2025.
- [35] Yubin Kim, Chanwoo Park, Hyewon Jeong, Yik S Chan, Xuhai Xu, Daniel McDuff, Hyeonhoon Lee, Marzyeh Ghassemi, Cynthia Breazeal, and Hae W Park. Mdagents: An adaptive collaboration of llms for medical decision-making. *Advances in Neural Information Processing Systems*, 37:79410–79452, 2024.

- [36] Volcano Engine. Volcano engine, 2025. Accessed: 2025-05-08.
- [37] University of Washington. Mygene2. <https://mygene2.org>. Accessed: 2025-02-23.
- [38] Emily Alsentzer, Michelle M Li, Shilpa N Kobren, Ayush Noori, Undiagnosed Diseases Network, Isaac S Kohane, and Marinka Zitnik. Few shot learning for phenotype-driven diagnosis of patients with rare genetic diseases. *medRxiv*, pages 2022–12, 2022.
- [39] Helen V Firth, Shola M Richards, A Paul Bevan, Stephen Clayton, Manuel Corpas, Diana Rajan, Steven Van Vooren, Yves Moreau, Roger M Pettett, and Nigel P Carter. Decipher: database of chromosomal imbalance and phenotype in humans using ensembl resources. *The American Journal of Human Genetics*, 84(4):524–533, 2009.
- [40] MR Macleod. Pubmed: <http://www.pubmed.org>. *Journal of Neurology, Neurosurgery & Psychiatry*, 73(6):746–746, 2002.
- [41] Zhiyu Zhao et al. Rarearena: A comprehensive rare disease diagnostic dataset with nearly 50,000 patients covering more than 4000 diseases. <https://github.com/zhao-zy15/RareArena>, 2025.
- [42] Dongzhi Jiang, Renrui Zhang, Ziyu Guo, Yanwei Li, Yu Qi, Xinyan Chen, Lihui Wang, Jianhan Jin, Claire Guo, Shen Yan, et al. Mme-cot: Benchmarking chain-of-thought in large multimodal models for reasoning quality, robustness, and efficiency. *arXiv preprint arXiv:2502.09621*, 2025.
- [43] Peter N Robinson, Vida Ravanmehr, Julius OB Jacobsen, Daniel Danis, Xingmin Aaron Zhang, Leigh C Carmody, Michael A Gargano, Courtney L Thaxton, Guy Karlebach, Justin Reese, et al. Interpretable clinical genomics with a likelihood ratio paradigm. *The American Journal of Human Genetics*, 107(3):403–417, 2020.
- [44] Thoralf Töpel, Dagmar Scheible, Friedrich Trefz, and Ralf Hofestädt. Ramedis: a comprehensive information system for variations and corresponding phenotypes of rare metabolic diseases. *Human mutation*, 31(1):E1081–E1088, 2010.
- [45] Simon Ronicke, Martin C Hirsch, Ewelina Türk, Katharina Larionov, Daphne Tientcheu, and Annette D Wagner. Can a decision support system accelerate rare disease diagnosis? evaluating the potential impact of ada dx in a retrospective study. *Orphanet journal of rare diseases*, 14:1–12, 2019.
- [46] Alistair Johnson, Tom Pollard, Steven Horng, Leo Anthony Celi, and Roger Mark. MIMIC-IV-NOTE: Deidentified free-text clinical notes. <https://www.physionet.org/content/mimic-iv-note/2.2/>.
- [47] Yanjie Fan, Ying Zhou, Huili Liu, Xiaomei Luo, Ting Xu, Yu Sun, Tingting Yang, Linlin Chen, Xuefan Gu, and Yongguo Yu. Improving variant prioritization in exome analysis by entropy-weighted ensemble of multiple tools. *Clinical Genetics*, 103(2):190–199, 2023.
- [48] Naomi Miller, Eve-Marie Lacroix, and Joyce EB Backus. Medlineplus: building and maintaining the national library of medicine’s consumer health web service. *Bulletin of the Medical Library Association*, 88(1):11, 2000.
- [49] Damian Smedley, Julius OB Jacobsen, Marten Jäger, Sebastian Köhler, Manuel Holtgrewe, Max Schubach, Enrico Siragusa, Tomasz Zemojtel, Orion J Buske, Nicole L Washington, et al. Next-generation diagnostics and disease-gene discovery with the exomiser. *Nature protocols*, 10(12):2004–2015, 2015.
- [50] François Remy, Kris Demuyne, and Thomas Demeester. Biolord: Learning ontological representations from definitions (for biomedical concepts and their textual descriptions). *arXiv preprint arXiv:2210.11892*, 2022.
- [51] Microsoft Corporation. Bing. <https://www.bing.com>, 2009. Accessed: 2025-04-08.
- [52] Google LLC. Google search. <https://www.google.com>, 1998. Accessed: 2025-04-08.
- [53] Gabriel Weinberg. Duckduckgo. <https://duckduckgo.com>, 2008. Accessed: 2025-04-08.

- [54] Crossref. Crossref. <https://www.crossref.org/>, 2025. Accessed: April 2025.
- [55] Steffanie S Weinreich, R Mangon, JJ Sikkens, ME En Teeuw, and MC Cornel. Orphanet: a european database for rare diseases. *Nederlands tijdschrift voor geneeskunde*, 152(9):518–519, 2008.
- [56] Joanna S Amberger, Carol A Bocchini, François Schiettecatte, Alan F Scott, and Ada Hamosh. Omim.org: Online mendelian inheritance in man (omim®), an online catalog of human genes and genetic disorders. *Nucleic acids research*, 43(D1):D789–D798, 2015.
- [57] Michael A Gargano, Nicolas Matentzoglou, Ben Coleman, Eunice B Addo-Lartey, Anna V Anagnostopoulos, Joel Anderton, Paul Avillach, Anita M Bagley, Eduard Bakštein, James P Balhoff, Gareth Baynam, Susan M Bello, Michael Berk, Holli Bertram, Somer Bishop, Hannah Blau, David F Bodenstein, Pablo Botas, Kaan Boztug, Jolana Čady, Tiffany J Callahan, Rhiannon Cameron, Seth J Carbon, Francisco Castellanos, J Harry Caufield, Lauren E Chan, Christopher G Chute, Jaime Cruz-Rojo, Noémi Dahan-Oliel, Jon R Davids, Maud de Dieuleveult, Vinicius de Souza, Bert B A de Vries, Esther de Vries, J Raymond DePaulo, Beata Derfalvi, Ferdinand Dhombres, Claudia Diaz-Byrd, Alexander J M Dingemans, Bruno Donadille, Michael Duyzend, Reem Elfeky, Shahim Essaid, Carolina Fabrizzi, Giovanna Fico, Helen V Firth, Yun Freudenberg-Hua, Janice M Fullerton, Davera L Gabriel, Kimberly Gilmour, Jessica Giordano, Fernando S Goes, Rachel Gore Moses, Ian Green, Matthias Griese, Tudor Groza, Weihong Gu, Julia Guthrie, Benjamin Gyori, Ada Hamosh, Marc Hanauer, Kateřina Hanušová, Yongqun (Oliver) He, Harshad Hegde, Ingo Helbig, Kateřina Holasová, Charles Tapley Hoyt, Shangzhi Huang, Eric Hurwitz, Julius O B Jacobsen, Xiaofeng Jiang, Lisa Joseph, Kamyar Keramatian, Bryan King, Katrin Knoflach, David A Koolen, Megan L Kraus, Carlo Kroll, Maaïke Kusters, Markus S Ladewig, David Lagorce, Meng-Chuan Lai, Pablo Lapunzina, Bryan Laraway, David Lewis-Smith, Xiarong Li, Caterina Lucano, Marzieh Majd, Mary L Marazita, Victor Martinez-Glez, Toby H McHenry, Melvin G McInnis, Julie A McMurry, Michaela Mihulová, Caitlin E Millett, Philip B Mitchell, Veronika Moslerová, Kenji Narutomi, Shahrzad Nematollahi, Julian Nevado, Andrew A Nierenberg, Nikola Novák Čajbiková, Jr. Nurnberger, John I, Soichi Ogishima, Daniel Olson, Abigail Ortiz, Harry Pachajoa, Guiomar Perez de Nanclares, Amy Peters, Tim Putman, Christina K Rapp, Ana Rath, Justin Reese, Lauren Rekerle, Angharad M Roberts, Suzy Roy, Stephan J Sanders, Catharina Schuetz, Eva C Schulte, Thomas G Schulze, Martin Schwarz, Katie Scott, Dominik Seelow, Berthold Seitz, Yiping Shen, Morgan N Similuk, Eric S Simon, Balwinder Singh, Damian Smedley, Cynthia L Smith, Jake T Smolinsky, Sarah Sperry, Elizabeth Stafford, Ray Stefancsik, Robin Steinhaus, Rebecca Strawbridge, Jagadish Chandrabose Sundaramurthi, Polina Talapova, Jair A Tenorio Castano, Pavel Tesner, Rhys H Thomas, Audrey Thurm, Marek Turnovec, Marielle E van Gijn, Nicole A Vasilevsky, Markéta Vlčková, Anita Walden, Kai Wang, Ron Wapner, James S Ware, Addo A Wiafe, Samuel A Wiafe, Lisa D Wiggins, Andrew E Williams, Chen Wu, Margot J Wyrwoll, Hui Xiong, Nefize Yalin, Yasunori Yamamoto, Lakshmi N Yatham, Anastasia K Yocum, Allan H Young, Zafer Yüksel, Peter P Zandi, Andreas Zankl, Ignacio Zarante, Miroslav Zvolský, Sabrina Toro, Leigh C Carmody, Nomi L Harris, Monica C Munoz-Torres, Daniel Danis, Christopher J Mungall, Sebastian Köhler, Melissa A Haendel, and Peter N Robinson. The human phenotype ontology in 2024: phenotypes around the world. *Nucleic Acids Research*, 52(D1):D1333–D1346, 11 2023.
- [58] Wikimedia Foundation. Wikipedia, the free encyclopedia. <https://www.wikipedia.org/>, 2001. Accessed: 2025-04-08.
- [59] National Library of Medicine (US). Medlineplus [internet], 2020. [updated Jun 24; cited 2020 Jul 1].
- [60] Qiao Jin, Won Kim, Qingyu Chen, Donald C Comeau, Lana Yeganova, W John Wilbur, and Zhiyong Lu. Medcpt: Contrastive pre-trained transformers with large-scale pubmed search logs for zero-shot biomedical information retrieval. *Bioinformatics*, 39(11):btad651, 2023.
- [61] Stephen Robertson, Hugo Zaragoza, et al. The probabilistic relevance framework: Bm25 and beyond. *Foundations and Trends® in Information Retrieval*, 3(4):333–389, 2009.
- [62] Siwei Chen, Laurent C Francioli, Julia K Goodrich, Ryan L Collins, Masahiro Kanai, Qingbo Wang, Jessica Alföldi, Nicholas A Watts, Christopher Vittal, Laura D Gauthier, et al. A genomic mutational constraint map using variation in 76,156 human genomes. *Nature*, 625(7993):92–100, 2024.

- [63] Nayanah Siva. 1000 genomes project, 2008.
- [64] Daniel Taliun, Daniel N Harris, Michael D Kessler, Jedidiah Carlson, Zachary A Szpiech, Raul Torres, Sarah A Gagliano Taliun, André Corvelo, Stephanie M Gogarten, Hyun Min Kang, et al. Sequencing of 53,831 diverse genomes from the NHLBI TOPMed program. *Nature*, 590(7845):290–299, 2021.
- [65] Statistics group Ciampi Antonio 8 Greenwood Celia MT (co-chair) 7 8 14 19 Hendricks Audrey E. 1 12 Li Rui 7 13 14 Metrustry Sarah 5 Oualkacha Karim 80 Tachmazidou Ioanna 1 Xu ChangJiang 7 8 Zeggini Eleftheria (co-chair) 1 et al. The UK10K project identifies rare variants in health and disease. *Nature*, 526(7571):82–90, 2015.
- [66] Jacob A Tennessen, Abigail W Bigham, Timothy D O’connor, Wenqing Fu, Eimear E Kenny, Simon Gravel, Sean McGee, Ron Do, Xiaoming Liu, Goo Jun, et al. Evolution and functional impact of rare coding variation from deep sequencing of human exomes. *science*, 337(6090):64–69, 2012.
- [67] Ivan A Adzhubei, Steffen Schmidt, Leonid Peshkin, Vasily E Ramensky, Anna Gerasimova, Peer Bork, Alexey S Kondrashov, and Shamil R Sunyaev. A method and server for predicting damaging missense mutations. *Nature methods*, 7(4):248–249, 2010.
- [68] Pauline C Ng and Steven Henikoff. Sift: Predicting amino acid changes that affect protein function. *Nucleic acids research*, 31(13):3812–3814, 2003.
- [69] Jana Marie Schwarz, Christian Rödelsperger, Markus Schuelke, and Dominik Seelow. Mutationtaster evaluates disease-causing potential of sequence alterations. *Nature methods*, 7(8):575–576, 2010.
- [70] Fredrik Eriksson, Erik Fransson, and Paul Erhart. The hiphive package for the extraction of high-order force constants by machine learning. *Advanced Theory and Simulations*, 2(5):1800184, 2019.
- [71] Google LLC. Google scholar. <https://scholar.google.com>, 2004. Accessed: 2025-04-08.
- [72] Sebastian Köhler, Michael Gargano, Nicolas Matentzoglou, Leigh C Carmody, David Lewis-Smith, Nicole A Vasilevsky, Daniel Danis, Ganna Balagura, Gareth Baynam, Amy M Brower, et al. The human phenotype ontology in 2021. *Nucleic acids research*, 49(D1):D1207–D1217, 2021.
- [73] T Michael Yates, Morad Ansari, Louise Thompson, Sarah E Hunt, Elena Cibrian Uhalte, Rachel J Hobson, Joseph A Marsh, Caroline F Wright, and Helen V Firth. Curating genomic disease-gene relationships with gene2phenotype (g2p). *Genome Medicine*, 16(1):127, 2024.
- [74] ICD10. <https://www.icd10data.com/ICD10CM/Codes>, Accessed: 2025-04-08.
- [75] Zhengyun Zhao, Qiao Jin, Fangyuan Chen, Tuorui Peng, and Sheng Yu. A large-scale dataset of patient summaries for retrieval-based clinical decision support systems. *Scientific Data*, 10(1):909, 2023.
- [76] Melissa J Landrum, Jennifer M Lee, Mark Benson, Garth Brown, Chen Chao, Shanmuga Chitipiralla, Baoshan Gu, Jennifer Hart, Douglas Hoffman, Jeffrey Hoover, et al. ClinVar: public archive of interpretations of clinically relevant variants. *Nucleic acids research*, 44(D1):D862–D868, 2016.
- [77] Peter D Stenson, Matthew Mort, Edward V Ball, Molly Chapman, Katy Evans, Luisa Azevedo, Matthew Hayden, Sally Heywood, David S Millar, Andrew D Phillips, et al. The human gene mutation database (hgmd®): optimizing its use in a clinical diagnostic or research setting. *Human genetics*, 139:1197–1207, 2020.
- [78] Konrad J Karczewski, Ben Weisburd, Brett Thomas, Matthew Solomonson, Douglas M Ruderfer, David Kavanagh, Tymor Hamamsy, Monkol Lek, Kaitlin E Samocha, Beryl B Cummings, et al. The exac browser: displaying reference data information from over 60 000 exomes. *Nucleic acids research*, 45(D1):D840–D845, 2017.
- [79] Jungo Kasai, Yuhei Kasai, Keisuke Sakaguchi, Yutaro Yamada, and Dragomir Radev. Evaluating GPT-4 and ChatGPT on Japanese medical licensing examinations, 2023.

- [80] Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*, 2024.
- [81] Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.

6 Data availability

The preprocessed databases and datasets are being prepared for public release and will be made available at <https://huggingface.co/datasets/Angelakeke/DeepRare> under the licenses.

7 Code availability

The web application used in this study is accessible at <http://raredx.cn>. Access requires approval - please contact zwk0629@sjtu.edu.cn with your research details and intended use case.

8 Acknowledgments

This work is supported by the National Key R&D Program of China (No. 2022YFC2703400 to Yongguo Yu) and STCSM (No. 22511106101, No. 18DZ2270700, No. 21DZ1100100). Weidi would like to acknowledge the funding from Scientific Research Innovation Capability Support Project for Young Faculty (ZY-GXQNJSKYCXNLZCXM-I22). Additionally, we gratefully acknowledge the developers and contributors of publicly available rare disease datasets, foundational research works, bioinformatics tools, and large language models that have collectively enabled our research.

9 Author Contributions

All listed authors meet the ICMJE four criteria for authorship. W.Z., C.W., and Y.F. contributed equally to this work. Y.Z., Y.Y., K.S., and W.X. are the corresponding authors. All authors (W.Z., C.W., Y.F., X.Z., P.Q., Y.S., X.Z., Y.W., Y.Z., Y.Y., K.S., and W.X.) contributed to the conception and design of the study. W.Z. and C.W. led the computational algorithm design, while Y.F. led the clinical design and medical aspects. W.Z., C.W., and Y.F. performed data acquisition, analysis, and interpretation. W.Z., C.W., and Y.F. drafted the manuscript, while X.Z., P.Q., Y.S., X.Z., Y.W., Y.Z., Y.Y., K.S., and W.X. provided critical revisions for important intellectual content. All authors approved the final version for publication and agree to be accountable for all aspects of the work, ensuring that questions related to the accuracy or integrity of any part of the work are appropriately investigated and resolved.

10 Competing Interests

We declare that the authors have no competing interests as defined by Nature Portfolio, or other interests that might be perceived to influence the results and/or discussion reported in this paper.

11 Supplementary

11.1 Main Workflow of DeepRare System

Algorithm 1 Main Workflow of DeepRare System

```
1: Input:  $\mathcal{I} = \{\mathcal{T}, \mathcal{H}, \mathcal{G}\}$  ▷ Patient input: Free-texts, HPO Items, Genetic Variants
2: Output:  $\mathcal{D}, \mathcal{R}$  ▷ Final diagnosis list with rationale explanations
3:
4: Initialize memory bank  $\mathcal{M} \leftarrow \emptyset$ 
5:  $\mathcal{D} \leftarrow \emptyset$  ▷ Initialize diagnosis list as empty
6:  $N \leftarrow N_0$  ▷ Initialize search depth
7:
8: while  $\mathcal{D} = \emptyset$  do ▷ Iterate until valid diagnosis is obtained
9:
10: // Information Collection Stage
11:
12: Phenotype Information Collection:
13: if  $\mathcal{T} \neq \emptyset$  then
14:    $\hat{\mathcal{P}} \leftarrow a_{\text{hpo}}(\mathcal{T}, \mathcal{H})$  ▷ HPO standardization
15: else
16:    $\hat{\mathcal{P}} \leftarrow \mathcal{H}$ 
17: end if
18:  $\mathcal{E}_{\text{hpo}} \leftarrow a_{\text{k-search}}(\hat{\mathcal{P}}, \mathcal{M}, N) \cup a_{\text{c-search}}(\hat{\mathcal{P}}, \mathcal{M}, N)$  ▷ Phenotype retrieval with depth  $N$ 
19:  $\mathcal{Y}_{\text{hpo}} \leftarrow a_{\text{hpo-analyzer}}(\hat{\mathcal{P}})$  ▷ Phenotype analysis
20:  $\mathcal{M} \leftarrow \mathcal{M} \cup (\mathcal{E}_{\text{hpo}}, \mathcal{Y}_{\text{hpo}})$  ▷ Update memory bank
21:
22:  $\mathcal{D}' \leftarrow \mathcal{A}_{\text{host}}(\hat{\mathcal{P}}, \mathcal{M} | \langle \text{prompt} \rangle_4)$  ▷ Tentative diagnosis
23:
24: Genotype Information Collection:
25: if  $\mathcal{G} \neq \emptyset$  then
26:    $\hat{\mathcal{G}} \leftarrow a_{\text{geno-analyzer}}(\mathcal{G})$  ▷ VCF annotation & variant ranking
27:    $\mathcal{D}'' \leftarrow \mathcal{A}_{\text{host}}(\hat{\mathcal{G}}, \mathcal{M}, \mathcal{D}', N | \langle \text{prompt} \rangle_5)$  ▷ Synthetic analysis with depth  $N$ 
28:    $\mathcal{M} \leftarrow \mathcal{M} \cup \mathcal{D}''$  ▷ Update memory bank
29: else
30:    $\mathcal{D}'' \leftarrow \mathcal{D}'$  ▷ Use phenotype-only diagnosis when no genetic data
31: end if
32:
33: // Self-reflective Diagnosis Stage
34:  $\mathcal{E}_{\text{disease}} \leftarrow a_{\text{k-search}}(a_{\text{d-norm}}(\mathcal{D}''), \mathcal{M})$  ▷ Disease retrieval
35:  $\mathcal{M} \leftarrow \mathcal{M} \cup \mathcal{E}_{\text{disease}}$  ▷ Update memory bank
36:  $\mathcal{D} \leftarrow \mathcal{A}_{\text{host}}(\mathcal{D}'', \hat{\mathcal{I}}, \mathcal{M} | \langle \text{prompt} \rangle_6)$  ▷ Self-reflection
37:
38: if  $\mathcal{D} = \emptyset$  then ▷ If self-reflection rules out all diseases
39:    $N \leftarrow N + \Delta N$  ▷ Increase search depth
40: end if
41:
42: end while
43:
44:  $\{\mathcal{D}, \mathcal{R}\} \leftarrow \mathcal{A}_{\text{host}}(\mathcal{D}, \hat{\mathcal{I}}, \mathcal{M} | \langle \text{prompt} \rangle_7)$  ▷ Final diagnosis with rationale
45:
46: return  $\mathcal{D}, \mathcal{R}$ 
```

Table 1 | Rare Disease Diagnosis Result (HPO Input only)

Model	RareBench (MME)			RareBench (HMS)			RareBench (LIRICAL)			RareBench (RAMEDIS)			MyGene			DDD			MIMIC-IV -Rare			Avg.			Xinhua Hosp.					
	case=40			case=88			case=370			case=624			case=146			case=2283			case=1875			-			case=975					
	R@1	R@3	R@5	R@1	R@3	R@5	R@1	R@3	R@5	R@1	R@3	R@5	R@1	R@3	R@5	R@1	R@3	R@5	R@1	R@3	R@5	R@1	R@3	R@5	R@1	R@3	R@5	R@1	R@3	R@5
Diagnosis API																														
PhenoBrain	25.00	45.00	80.00	22.35	32.94	42.35	27.10	44.72	57.72	20.67	38.30	55.77	45.21	58.90	61.64	38.92	54.68	63.11	8.11	10.24	11.41	26.78	40.68	53.14	-	-	-	-	-	-
PubCaseFinder	40.00	57.50	57.50	7.32	13.41	15.85	37.77	47.01	48.91	25.84	32.58	39.97	20.55	26.71	28.08	26.92	32.32	34.05	1.06	1.39	1.60	22.78	30.13	32.28	-	-	-	-	-	-
LLM																														
GPT-4o	2.50	5.00	10.00	47.73	60.23	65.91	31.08	42.16	48.92	35.47	52.24	61.06	19.31	35.17	37.93	33.30	42.28	45.18	13.76	22.35	25.60	26.16	37.06	42.08	-	-	-	-	-	-
DeepSeek V3	5.00	20.00	27.50	42.05	56.82	62.50	33.24	42.97	45.68	40.06	66.51	72.11	20.83	39.58	47.22	34.01	45.27	48.83	14.61	23.04	26.18	27.11	42.03	47.15	42.43	54.91	59.41	-	-	-
Gemini 2	5.00	7.50	10.00	31.82	47.73	54.55	30.81	39.73	43.51	40.37	52.56	60.42	28.77	48.63	54.79	31.32	43.60	53.03	14.29	20.16	24.11	26.05	37.13	42.92	-	-	-	-	-	-
Claude-3.7-Sonnet	7.50	27.50	32.50	43.18	60.23	63.64	33.24	46.49	50.00	42.47	68.59	73.88	19.44	36.11	39.58	32.28	42.76	45.93	12.11	18.24	20.80	27.17	42.84	46.62	-	-	-	-	-	-
Reasoning LLM																														
o3mini	2.50	7.50	7.50	30.68	44.32	54.55	25.68	35.41	40.00	39.42	64.42	67.47	34.48	46.21	49.66	34.09	42.96	45.50	10.88	19.52	23.90	25.39	37.19	41.23	-	-	-	-	-	-
DeepSeek R1	35.00	60.00	65.00	38.64	53.41	60.23	36.76	45.14	47.57	33.17	47.60	53.21	35.62	46.58	47.95	38.24	49.43	53.01	13.81	19.92	22.93	33.03	46.01	49.99	39.16	49.85	54.98	-	-	-
Gemini 2 think	0.00	5.00	12.50	38.64	48.86	56.82	30.81	41.35	44.05	33.97	50.80	58.65	39.72	53.42	53.42	34.43	46.33	50.02	15.47	20.53	25.39	27.58	38.05	42.98	-	-	-	-	-	-
Claude 3.7 think	27.50	50.00	57.50	44.32	57.95	62.50	38.11	50.27	52.43	41.83	63.78	72.92	34.25	37.67	39.72	34.01	45.27	48.83	13.70	21.28	24.32	33.39	46.60	51.17	-	-	-	-	-	-
Medical LLM																														
MMedS-Llama 3	12.50	20.00	20.00	17.04	29.55	31.82	21.08	29.73	33.51	30.45	45.67	50.00	15.07	23.97	28.09	22.18	33.80	38.34	14.04	22.58	26.37	18.91	29.33	32.59	29.55	40.18	45.30	-	-	-
Baichuan-14B	5.00	10.00	10.00	38.64	53.41	60.23	26.77	37.03	40.27	37.66	64.26	68.75	6.16	23.91	30.82	29.32	39.37	44.29	14.40	21.33	24.59	22.56	35.62	39.85	43.35	54.40	58.89	-	-	-
Other Agentic System																														
MDAgents (GPT-4o)	2.25	2.25	5.00	27.27	40.91	46.59	21.89	28.92	31.08	37.82	58.01	62.98	15.75	22.60	24.66	25.01	33.33	36.01	10.98	17.12	20.27	20.14	29.02	32.37	-	-	-	-	-	-
MDAgents (DS-V3)	0.00	0.00	0.025	29.55	43.18	45.45	23.24	28.11	31.08	36.22	59.29	63.30	8.21	11.64	15.75	25.76	32.33	34.69	8.69	15.04	19.41	18.81	27.08	29.96	-	-	-	-	-	-
DS-R1-search	12.50	17.50	22.50	38.64	55.68	62.50	33.51	45.41	49.46	39.90	66.51	71.31	15.75	34.93	46.57	32.30	47.15	52.41	19.95	28.75	32.69	27.51	42.28	48.21	-	-	-	-	-	-
DeepRare System																														
DeepRare (GPT-4o)	72.50	77.50	82.50	54.54	65.90	69.31	51.62	62.70	65.40	73.23	81.57	85.26	75.86	80.00	80.69	39.85	51.36	55.88	21.44	28.89	31.71	55.58	63.99	67.25	-	-	-	-	-	-
DeepRare (DS-V3)	70.00	77.50	77.50	56.97	65.12	70.93	55.95	65.16	67.57	72.60	82.85	84.62	72.60	76.03	76.03	42.97	53.48	56.72	29.19	36.60	39.06	57.18	65.25	67.49	58.27	71.04	74.35	-	-	-
DeepRare (DS-R1)	67.57	81.08	83.78	51.76	60.00	67.06	52.43	63.51	67.84	62.17	76.28	79.32	75.00	80.55	81.94	44.27	57.09	61.27	23.27	31.16	34.79	53.78	64.24	68.00	54.30	66.90	70.20	-	-	-
DeepRare (claude-3.7)	72.50	85.00	87.50	56.82	63.64	64.77	40.00	51.08	54.86	70.03	81.57	84.46	73.79	80.68	80.69	45.49	58.85	63.76	24.32	30.30	32.12	54.71	64.45	66.88	-	-	-	-	-	-
DeepRare (gemini-2)	77.50	85.00	90.00	40.91	54.25	57.95	48.91	63.24	68.91	49.67	69.55	75.64	69.18	77.40	79.45	47.59	59.82	62.70	19.08	26.28	29.41	50.41	62.22	66.29	-	-	-	-	-	-

Note: All values reported in the table are percentages (%).

11.2 Prompt Sets

Prompt

Prompt 1. Prompt for Baseline LLM Diagnosis

You are a specialist in the field of rare diseases. You will be provided and asked about a complicated clinical case; read it carefully and then provide a diverse and comprehensive differential diagnosis.

Patient's {info_type}: {patient_info}

Enumerate the top 5 most likely diagnoses. Be precise, listing one diagnosis per line, and try to cover many unique possibilities (at least 5).

The top 10 diagnoses are:

Prompt

Prompt 2. Prompt for Diagnosis Result Evaluation

You are a specialist in the field of rare diseases.

I will now give you five predicted diseases if the predicted diagnosis is in the standard diagnosis. Please output the predicted rank; otherwise, output "No". Only output "No" or "1-5" numbers. If the predicted disease has multiple conditions, only output the top rank. Output only "No" or one number, no additional output.

Predicted diseases: {predict_diagnosis}

Standard diagnosis: {golden_diagnosis}

Prompt

Prompt 3. Prompt for Rare Disease Type Classification

You are a medical disease classifier specializing in categorizing diseases into predefined categories.

Your task is to classify the given disease into one or more of these categories:

```
[  
"Blood, Heart and Circulation",  
"Bones, Joints and Muscles",  
"Brain and Nerves",  
"Digestive System",  
"Ear, Nose and Throat",  
"Endocrine System",  
"Eyes and Vision",  
"Immune System",  
"Kidneys and Urinary System",  
"Lungs and Breathing",  
"Mouth and Teeth",  
"Skin, Hair and Nails",  
"Female Reproductive System",  
"Male Reproductive System"  
]
```

Note: If the input contains multiple disease names separated by "/" (slash), they are synonyms or alternate names for the same disease. Consider them as a single disease and classify accordingly.

Please output your results in JSON format ONLY as follows:

```
“‘json  
{  
"disease": "the original disease name(s) exactly as provided",  
"category": ["category1", "category2", ...]  
}  
“‘
```

Your response must be ONLY the JSON object - no additional text, explanations, or commentary. If a disease could belong to multiple categories, include all relevant categories in the "category" array. If you're unsure about the classification, make your best judgment based on the disease characteristics.

The following is a disease name. Please classify it according to the instructions:

Prompt

Prompt 4. Prompt for Tentative Diagnosis Decision-making

You are a specialist in the field of rare diseases.

You have access to the following context:

- **Online knowledge** (with titles and URLs): {web_diagnosis}
- **LLM-generated diagnoses**: {llm_response}
- **Diagnosis API results**: {diagnosis_api_response}
- **Similar cases**: {similar_case_detailed}
- **Prompt details**: {patient_info}

—

Based on the above and your knowledge, enumerate the **top 5 most likely rare disease diagnoses** for this patient.

For each diagnosis, use the following format:

DIAGNOSIS NAME (Rank #X/5)

Diagnostic Reasoning:

- Provide 2-3 concise sentences explaining why this rare disease fits the clinical picture.
- Integrate evidence from all available sources (online knowledge, similar cases, LLM outputs, and API results).
- Support your reasoning with specific, in-text citations in [X] format, referencing the most relevant sources (including specific similar cases, articles, or diagnostic tools).
- Briefly discuss the pathophysiological basis for the diagnosis, citing relevant literature or case evidence.

—

After listing all 5 diagnoses, include a reference section:

References:

- Number each reference in the order it is first cited ([1], [2], ...).
- Only include sources you directly cited in your diagnostic reasoning above.
- For each reference, should provide:
 - a. Source type (e.g., medical guideline, similar case, literature, diagnosis assisent tool...)
 - b. Use 3-4 sentences to describe of the content and its relevance.
 - c. For articles or literature, include the title and URL if provided.
- Every in-text citation [X] in your reasoning should correspond to a numbered entry in your reference list.
- Try to cover as many sources and references.
- Do not repeat!!

—

Key Instructions:

1. Always use in-text citations in [X] format, matching only the references you actually cite in your reasoning.
2. Each diagnosis must be a rare disease (**bolded** using markdown).
3. Rank from most (#1) to least (#5) likely.
4. Integrate information from all provided sources (medical literature, similar cases, and judgment analyses) wherever appropriate.
5. Do **not** copy or invent references—only include those present in the provided materials.

Prompt

Prompt 5. Prompt for Rare Disease Gene Analysis

You are a specialist in the field of rare diseases.

Here is a rare disease diagnosis case.

- **Exomiser gene/variant prioritization summary**: {exomiser_summary}
- **Phenotypic description (HPO terms)**: {hpo_terms}
- **Preliminary diagnosis based only on phenotype**: {pheno_only_diagnosis}

—

Based on the above information and your knowledge, enumerate the **top 5 most likely rare disease diagnoses** for this patient.

For each diagnosis, use the following format:

DIAGNOSIS NAME (Rank #X/5)

Diagnostic Reasoning:

- Provide 2-3 concise sentences explaining why this rare disease fits the clinical picture.
- Integrate evidence from all available sources, prioritizing the Exomiser gene/variant results while considering phenotypic data and preliminary diagnosis as supporting evidence.
- Support your reasoning with specific, in-text citations in [X] format, referencing the most relevant sources (Exomiser findings, HPO phenotype matches, or preliminary diagnostic considerations).
- Briefly discuss the pathophysiological basis for the diagnosis, relating gene function to observed phenotype.

—

After listing all 5 diagnoses, include a reference section:

References:

- Number each reference in the order it is first cited ([1], [2], ...).
- Only include sources you directly cited in your diagnostic reasoning above.
- For each reference, provide:
 - a. Source type (e.g., Exomiser gene prioritization, HPO phenotype analysis, preliminary phenotype-based diagnosis)
 - b. Use 3-4 sentences to describe the content and its relevance to the diagnostic reasoning.
- Every in-text citation [X] in your reasoning should correspond to a numbered entry in your reference list.
- Try to cover all relevant sources and references.
- Do not repeat!!

Prompt

Prompt 6. Prompt for Disease Reflection

Assume you are a doctor specialized in rare disease diagnosis.

Based on the patient information, similar case diagnoses, and disease knowledge, evaluate whether the proposed diagnosis is correct for this patient.

Begin with a clear "DIAGNOSIS ASSESSMENT: [Correct/Incorrect]" statement, followed by your reasoning.

Structure your analysis as follows:

1. **PATIENT SUMMARY:** Briefly summarize the patient's key symptoms
2. **PROPOSED DIAGNOSIS ANALYSIS:** Evaluate the proposed diagnosis (`{diagnosis_to_judge}`) in relation to the patient's symptoms
3. **REFERENCES:** Extract and number the most relevant evidence from the provided medical literature that supports your conclusion

Patient phenotype:
`{patient_info}`

Similar cases:
`{similar_case_detailed}`

Medical literature:
`{disease_knowledge}`

Prompt

Prompt 7. Prompt for Final Diagnosis

You have access to the following information:

- Patient presentation: {patient_info}
- Similar cases: {similar_case_detailed}
- Primary diagnosis results (with references): {tentative_result}
- Disease Reflection (with references): {judgements}

****Task:****

Based on all the above, enumerate the top 5 most likely rare disease diagnoses for this patient.

****For each diagnosis, follow this format exactly:****

****DISEASE NAME**** (Rank #X/5)

Diagnostic Reasoning:

- Provide 3-4 sentences explaining why this diagnosis fits the patient's presentation.
 - Specify which patient symptoms and findings support this diagnosis.
 - Clearly explain the underlying pathophysiological mechanisms (briefly).
 - Integrate and ****cite specific evidence**** from the provided references (including medical literature, similar cases, or judgement analyses), using in-text [X] citation style.
 - Try to cite as more sources and references but do not add hallucination content.
-

****After listing all 5 diagnoses, include a reference section:****

References:

- Number each reference in the order it is first cited ([1], [2], ...).
 - Only include sources you directly cited in your diagnostic reasoning above.
 - For each reference, should provide:
 - a. Source type (e.g., medical guideline, similar case, literature, diagnosis assisent tool...) (Do not use source type: "Judgement analysis", "Disease Reflection")
 - b. Use 3-4 sentences to describe of the content and its relevance.
 - c. For articles or literature, include the title and URL if provided.
 - Every in-text citation [X] in your reasoning should correspond to a numbered entry in your reference list.
 - Try to cover as more sources and references.
 - Do not repeat!!
-

****IMPORTANT GUIDELINES:****

1. Each diagnosis must be a rare disease (****bolded**** using markdown).
2. Rank from most (#1) to least (#5) likely.
3. Integrate information from all provided sources (medical literature, similar cases, and judgement analyses) wherever appropriate.
4. Do ****not**** copy or invent references—only include those present in the provided materials.
5. Remember to add the summary of the content, url for each reference.

Prompt

Prompt 8. Prompt for HPO Extraction

Given a paragraph of patient information from discharge note, please extract the phenotype about this patient only.

Check the Human Phenotype Ontology (HPO) database to determine the phenotype.

Only output the extracted phenotypes.

Use the format: {'HPO': 'HP:0000000', 'Phenotype': 'Phenotype description'}

Patient information: {case_report}

Prompt

Prompt 9. Prompt for HPO Extraction Modify

You are a medical terminology translator specializing in rare diseases. Your task is to convert patient phenotype description into standardized HPO (Human Phenotype Ontology) concept.

Input can be in either Chinese or English describing clinical phenotype. Analyze the description carefully, identify the phenotypes mentioned, and map them to standard English concepts in the HPO database.

Please output your results in JSON format as follows:

```
{  
  "original_term": "the original phenotype description",  
  "hpo_term": "standardized HPO term in English"  
}
```

If the phenotype doesn't exist in HPO, output:

```
{  
  "original_term": "the original phenotype description",  
  "hpo_term": "none"  
}
```

For each identified phenotype:

1. Do not include any phenotypes that are not present in the input.
2. Ensure the JSON is properly formatted and valid.
3. Your response must be ONLY the JSON object - no additional text, explanations, or commentary.
4. Provide the standard English name of the HPO term
5. If the phenotype doesn't have a corresponding concept in HPO, set the hpo_term to "none"

Example 1:

Input: "Metabolic dysfunction"

Output:

```
{  
  "original_term": "Metabolic dysfunction",  
  "hpo_term": "Abnormality of metabolism/homeostasis"  
}
```

Input: "Dark complexion"

Output:

```
{  
  "original_term": "Dark complexion",  
  "hpo_term": "Hyperpigmentation of the skin"  
}
```

The following is a patient phenotype description. Please convert it to HPO terms:

Prompt

Prompt 10. Prompt for Knowledge Summarization

Assume you are a doctor, please summarize these medical article into a paragraph.

Only keep key message, mainly focus on the phenotype and related disease.

If this article is not related to medical, please output "not a medical-related page".

Prompt

Prompt 11. Prompt for Case Summarization

Assume you are a doctor experienced in rare disease diagnosis.

Please judge if the two patient cases are likely to be the same disease based on the patient information.

Only output 'Yes' or 'No'.

Patient 1 phenotype: {patient_info}

Patient 2 phenotype: {retrieved_patient_case}

Prompt

Prompt 12. Prompt for Zero-shot LLM Inference

You are a specialist in the field of rare diseases.

You will be provided and asked about a complicated clinical case; read it carefully and then provide a diverse and comprehensive differential diagnosis.

Patient's {info_type}: {patient_info}

Enumerate the top 5 most likely diagnoses. Be precise, and try to cover many unique possibilities.

Each diagnosis should be a rare disease.

Use ** to tag the disease name.

Make sure to reorder the diagnoses from most likely to least likely.

Now, List the top 5 diagnoses.

11.3 Exomiser Config

```
1 CONFIG_TEMPLATE = {
2   "analysis": {
3     "genomeAssembly": "GRCh37",
4     "outputOptions": {
5       "outputFormat": ["TSV", "HTML"]
6       "outputPrefix": {output_prefix}
7     },
8     "frequencySources": [
9       "THOUSAND_GENOMES", "TOPMED", "UK10K", "ESP_AA", "ESP_EA",
10      "ESP_ALL", "GNOMAD_E_AFR", "GNOMAD_E_AMR", "GNOMAD_E_EAS",
11      "GNOMAD_E_NFE", "GNOMAD_E_SAS", "GNOMAD_G_AFR",
12      "GNOMAD_G_AMR", "GNOMAD_G_EAS", "GNOMAD_G_NFE", "GNOMAD_G_SAS"
13    ],
14    "pathogenicitySources": ["POLYPHEN", "MUTATION_TASTER", "SIFT"],
15    "analysisMode": "PASS_ONLY",
16    "inheritanceModes": {
17      "AUTOSOMAL_DOMINANT": 0.1,
18      "AUTOSOMAL_RECESSIVE_HOM_ALT": 0.1,
19      "AUTOSOMAL_RECESSIVE_COMP_HET": 2.0,
20      "X_DOMINANT": 0.1,
21      "X_RECESSIVE_HOM_ALT": 0.1,
22      "X_RECESSIVE_COMP_HET": 2.0,
23      "MITOCHONDRIAL": 0.2
24    },
25    "steps": [
26      {"failedVariantFilter": {}},
27      {"variantEffectFilter": {
28        "remove": [
29          "FIVE_PRIME_UTR_EXON_VARIANT",
30          "FIVE_PRIME_UTR_INTRON_VARIANT",
31          "THREE_PRIME_UTR_EXON_VARIANT",
32          "THREE_PRIME_UTR_INTRON_VARIANT",
33          "NON_CODING_TRANSCRIPT_EXON_VARIANT",
34          "NON_CODING_TRANSCRIPT_INTRON_VARIANT",
35          "CODING_TRANSCRIPT_INTRON_VARIANT",
36          "UPSTREAM_GENE_VARIANT",
37          "DOWNSTREAM_GENE_VARIANT",
38          "INTERGENIC_VARIANT",
39          "REGULATORY_REGION_VARIANT"
40        ]
41      }},
42      {"frequencyFilter": {"maxFrequency": 1.0}},
43      {"pathogenicityFilter": {"keepNonPathogenic": true}},
44      {"inheritanceFilter": {}},
45      {"omimPrioritiser": {}},
46      {"hiPhivePrioritiser": {}}
47    ]
48    "vcf": {vcf_path}
49    "hpoIds": {hpo_ids}
50  }
51 }
```

Listing 1 | Exomiser Configuration Template