GeneAgent: self-verification language agent for gene-set analysis using domain databases

Zhiyong Lu, PhD FACMI FIAHSI Senior Investigator, NLM, NIH Adjunct Professor, CS UIUC

> Bioinformatics Community Fair NIH Research Festival September 9, 2025

Our research at NIH IRP

- Research Areas
 - AI & Machine Learning, LLMs
 - Natural Language Processing (NLP)
 - Medical Image Analysis
- Multimodal data analysis
 - Biomedical Literature
 - Clinical notes, EHRs
 - CT, CXR & retinal images







- > LLM-powered applications in biomedicine
 - 1. TrialGPT: Assisting patient-to-trial matching with LLMs (2024)
 - GeneAgent: Al agent for gene set analysis (2025)

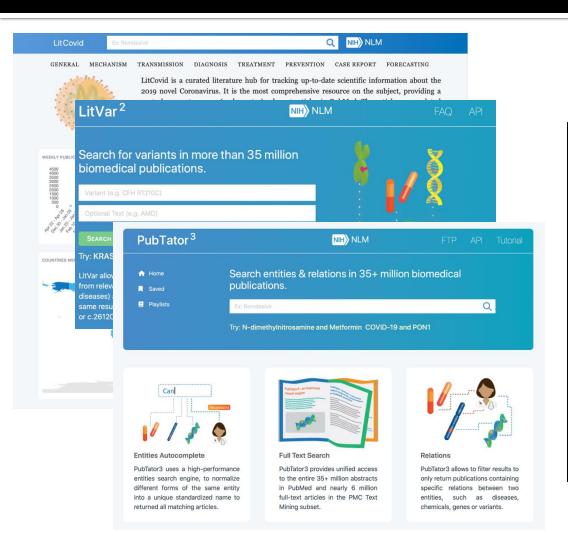
AI in PubMed - serving millions each day

- Related articles
- Spell checker
- Query autosuggest
- Semantic query understanding
- Citation sensor
- Author name disambiguation
- Query expansion
- Best Match: Sort by Relevance

- ...



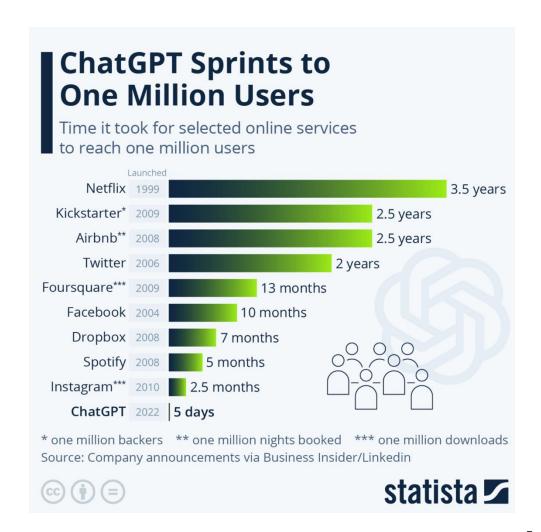
Biomedical Literature Mining





ChatGPT: revolution or hype?





WIRED

Dr. ChatGPT Will See You Now

24 days ago

E (M) Health Life, But Better Fitness Food Sleep Mindfulness More

Al may be as effective as medical specialists at diagnosing disease

By Jack Guy, CNN



The Guardian

First NHS AI-run physio clinic in England halves back-pain waiting list

3 days ago



wsj Wall Street Journal

Companies Bring Al Agents to Healthcare

Feb 27, 2025

The Washington Post

ChatGPT is little help for doctors in diagnosing diseases, study finds

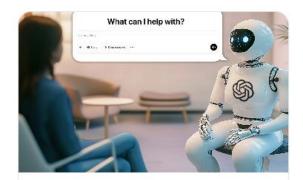
The research, conducted with 50 physicians last year, found that using ChatGPT did not significantly improve doctors' diagnostic reasoning.



CNN · 11d

FDA's artificial intelligence is supposed to revolutionize drug approvals. It's making up studies

Insiders tell CNN the FDA's AI is "hallucinating" studies and can't access ...



New York Post

Harmful AI therapy: Chatbots endanger users with suicidal thoughts, delusions, researchers warn

Jun 28, 2025



News & Events

Recent News Releases



NIH researchers develop AI agent that improves accuracy of gene set analysis by leveraging expertcurated databases

July 28, 2025 — The AI agent could help lead to a better understanding of how different diseases and conditions affect groups of genes individually and together.

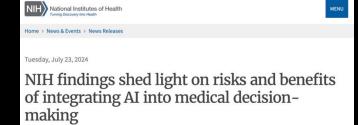


NEWS RELEASES

Monday, November 18, 2024

NIH-developed AI algorithm matches potential volunteers to clinical trials

Such an algorithm may save clinicians time and accelerate clinical enrollment and research.



Researchers at the National Institutes of Health (NIH) found that an artificial intelligence (AI) model solved medical quiz questions—designed to test health professionals' ability to diagnose patients based on clinical images and a brief text summary—with high accuracy. However, physician-graders found the AI model made mistakes when describing images and explaining how its decision-making led to the correct answer. The findings, which shed light on AI's potential in the clinical setting, were published in npj Digital Medicine a: The study was led by researchers from NIH's National Library of Medicine (NLM) and Weill Cornell Medicine, New York City.

"Integration of AI into health care holds great promise as a tool to help medical professionals diagnose

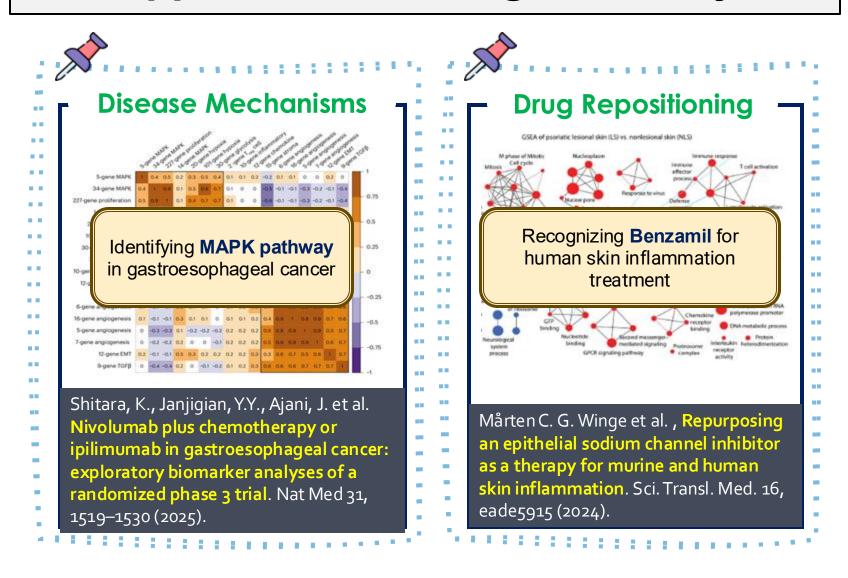


GPT-4V, an AI model, orten made mistakes when describing the medical image and explaining its reasoning behind the diagnosis—even in cases where it made the correct final choice. NIH/NLM

Intro: Gene Set Analysis

- Gene sets from high-throughput experiments
- Differentially expressed genes under different conditions
- Data analysis goal: determine the collective functions by a group of genes

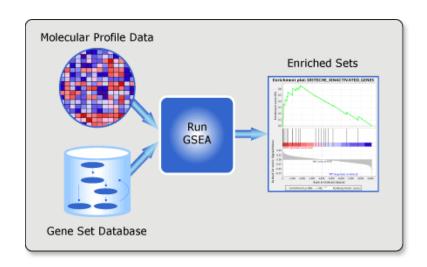
GSA applications in drug discovery



Existing GESA methods & their limitations

 Limited to curated knowledge and/or predefined gene sets

No explanation to support predictions



Related works

[Submitted on 21 May 2023 (v1), last revised 25 May 2023 (this version, v2)]

Gene Set Summarization using Large Language Models

Marcin P. Joachimiak, J. Harry Caufield, Nomi L. Harris, Hyeongsik Kim, Christopher J. Mungall

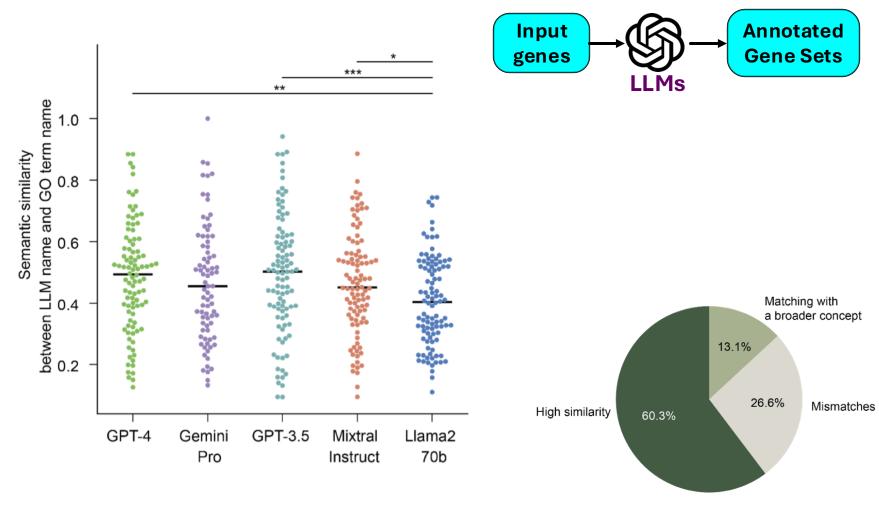
Molecular biologists frequently interpret gene lists derived from high-throughput experiments and computational analysis. This is typically done as a statistical enrichment analysis that measures the over- or under-representation of biological function terms associated with genes or their properties, based on curated assertions from a knowledge base (KB) such as the Gene Ontology (GO). Interpreting gene lists can also be framed as a textual summarization task, enabling the use of Large Language Models (LLMs), potentially utilizing scientific texts directly and avoiding reliance on a KB. We developed SPINDOCTOR (Structured Prompt Interpolation of Natural Language Descriptions of Controlled Terms for Ontology Reporting), a method tl [Submitted on 7 Sep 2023 (v1), last revised 1 Apr 2024 (this version, v2)] summarization as a complement to standard enrichme annotations, (2) ontology-free narrative gene summari function We demonstrate that these methods are able to genera values and often return terms that are not statistically able to recapitulate the most precise and informative t inability to generalize and reason using an ontology. R variations in prompt resulting in radically different terr based methods are unsuitable as a replacement for sta curation of ontological assertions remains necessary.

sources of gene functional information: (1) structured Evaluation of large language models for discovery of gene set

term lists for gene sets. However, GPT-based approach Mengzhou Hu, Sahar Alkhairy, Ingoo Lee, Rudolf T. Pillich, Dylan Fong, Kevin Smith, Robin Bachelder, Trey Ideker, Dexter Pratt

> Gene set analysis is a mainstay of functional genomics, but it relies on curated databases of gene functions that are incomplete. Here we evaluate five Large Language Models (LLMs) for their ability to discover the common biological functions represented by a gene set, substantiated by supporting rationale, citations and a confidence assessment. Benchmarking against canonical gene sets from the Gene Ontology, GPT-4 confidently recovered the curated name or a more general concept (73% of cases), while benchmarking against random gene sets correctly yielded zero confidence. Gemini-Pro and Mixtral-Instruct showed ability in naming but were falsely confident for random sets, whereas Llama2-70b had poor performance overall. In gene sets derived from 'omics data, GPT-4 identified novel functions not reported by classical functional enrichment (32% of cases), which independent review indicated were largely verifiable and not hallucinations. The ability to rapidly synthesize common gene functions positions LLMs as valuable 'omics assistants.

Key results in Hu et al.,



LLM as a judge

- Pros: scalable; consistent
- Cons: Selfpreference; limited factual checking
- We propose an AI
 agent that performs
 automatic self verification grounded
 in domain knowledge



Gene-centric information in expertcurated biological databases



















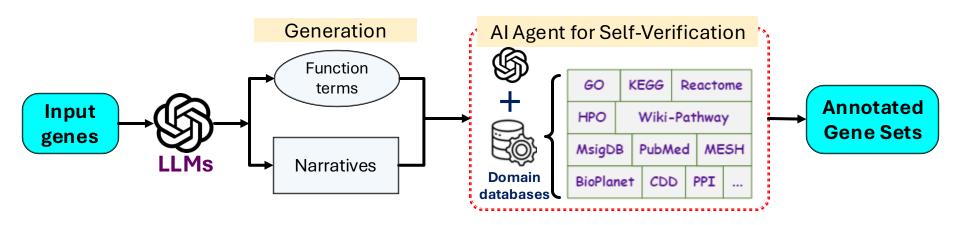




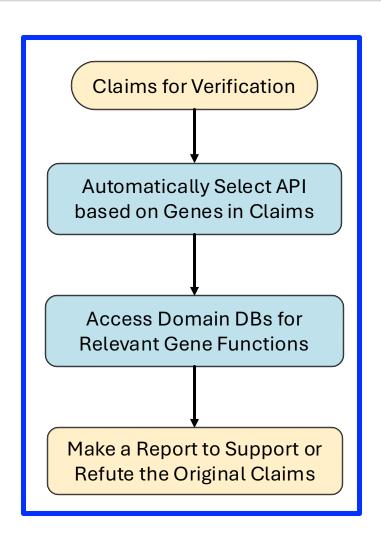
Our Gene Agent Approach

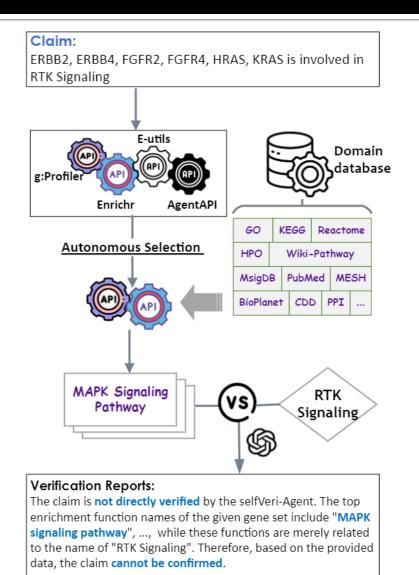


Zhizheng Wang, PhD



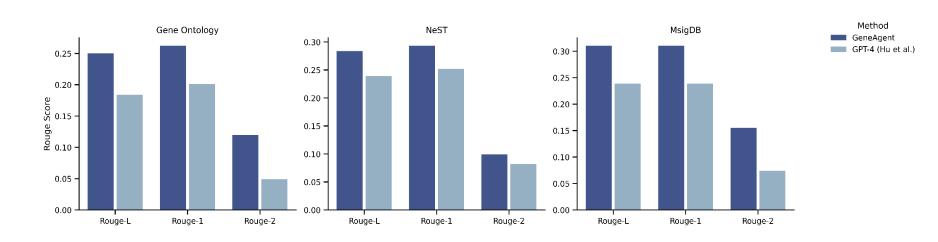
Al Agent for self-Verification





GeneAgent vs. standard GPT-4

Dataset	#gene sets	#genes	Avg. genes		
Gene Ontology	1,000	3 to 456	48.32		
NeST	50	5 to 323	18.96		
MsigDB	56	4 to 200	112.00		
All	1,106	3 to 456	50.67		



Pilot study with novel gene sets

- To assess its potential utility in real-world applications
- Worked with domain experts from NCI/NLM
- Novel gene sets from mouse
 B2905 melanoma cell line



Chi-Ping Day PhD Lab of Cancer Biology and Genetics, Cancer Data Science Lab, NCI



Christina Ross, Ph.D. NCI Lab of Cancer Biology and Genetics

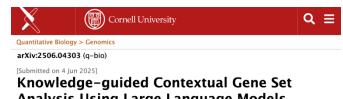
Wang, et al. GeneAgent: Self-verification Language Agent for Gene Set Knowledge Discovery using Domain Databases. *Nature Methods*, 2025.

Evaluation results by 2 domain experts

ID	Generated by GPT-4	Generated by GeneAgent	Gene Coverage	Better Output Annotated by Genomic Experts									
				Relevance		Readability		Consistency		Comprehensive		Final Decision	
				GPT-4	GeneAgent	GPT-4	GeneAgent	GPT-4	GeneAgent	GPT-4	GeneAgent	GPT-4 C	eneAgent
mmu05171 (HA-R)	Ribosomal Protein Synthesis	Cytosolic Ribosome and Protein Synthesis	33/36		0	0	0	0	0	0	0		✓
mmu03010 (HA-R)	Ribosomal Protein Synthesis and Assembly	Cytosolic Ribosome	34/35		0	0	0	0	0	0	0		✓
mmu03010 (HA-S)	Ribosomal Protein Synthesis	Cytosolic Ribosome	13/49									×	×
mmu05171 (HA-S)	Ribosomal Protein Synthesis	Cytosolic Ribosome Assembly and Protein Synthesis	47/47		0	0	0		0		0		~
mmu04015 (HA-S)	MAPK/ERK Pathway Regulation	Rap1 Signaling Pathway	27/27		0	0	0	0	0	0	0		✓
mmu05100 (HA-S)	Caveolae- Mediated Endocytosis and Actin Remodeling	Bacterial Invasion of Epithelial Cells	19/19	0		0	0	0		0		>	
mmu05022 (LA-S)	Oxidative Phosphorylation and Neurodegeneration	Neurodegener ation and Respiratory Chain Complex	23/24		0	0		0	0		0		~

Limitations & Future Directions

- Human & mouse genes only
- Very large gene sets
- Add contextual information (study designs, experimental conditions)



Analysis Using Large Language Models

Zhizheng Wang, Chi-Ping Day, Chih-Hsuan Wei, Qiao Jin, Robert Leaman, Yifan Yang, Shubo Tian, Aodong Qiu, Yin Fang, Qingging Zhu, Xinghua Lu, Zhiyong Lu

View PDF

Gene set analysis (GSA) is a foundational approach for interpreting genomic data of diseases by linking genes to biological processes. However, conventional GSA methods overlook clinical context of the analyses, often generating long lists of enriched pathways with redundant, nonspecific, or irrelevant results. Interpreting these requires extensive, ad-hoc manual effort, reducing both reliability and reproducibility. To address this limitation, we introduce cGSA, a novel Al-driven framework that enhances GSA by incorporating context-aware pathway prioritization, cGSA integrates gene cluster detection, enrichment analysis, and large language models to identify pathways that are not only statistically significant but also biologically meaningful. Benchmarking on 102 manually curated gene sets across 19 diseases and ten disease-related biological mechanisms shows that cGSA outperforms baseline methods by over 30%, with expert validation confirming its increased precision and interpretability. Two independent case studies in melanoma and breast cancer further demonstrate its potential to uncover context-specific insights and support targeted hypothesis generation.

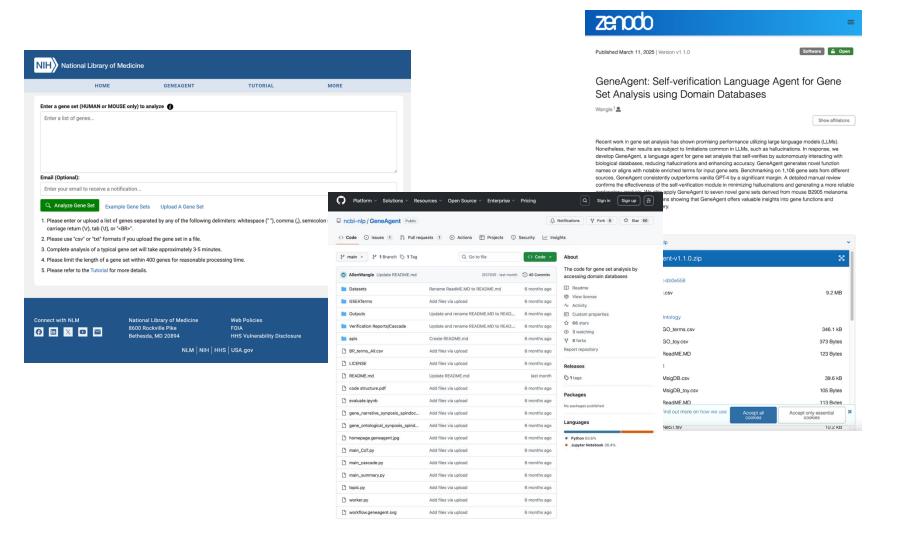
Comments: 56 pages, 9 figures, 1 table

Subjects: Genomics (q-bio.GN); Artificial Intelligence (cs.AI); Machine Learning (cs.LG)

arXiv:2506.04303 [q-bio.GN]

(or arXiv:2506.04303v1 [q-bio.GN] for this version) https://doi.org/10.48550/arXiv.2506.04303

We welcome your feedback



Current & Future Work



Benchmarking stock models on medical tasks

- MedCal-Bench: evaluating LLMs for medical calculations (NeurIPS oral, 24)
- > Evaluating LLMs on various BioNLP tasks (*Nature Communications*. 2025)

zhiyong.lu@nih.gov

> Enhancing standard LLMs with domain-specific data or tools

- GeneGPT: domain tool learning (Bioinformatics, 2024)
- MedRAG: Best practices with retrieval augmented generation (ACL 2024, PSB 2025)

Multimodal LLMs in Healthcare

- Assessing Progress of Multimodal LLM Performance on Clinical Cases (Radiology, 2025)
- Developing MLLMs with CT images and clinical notes (*Nature Biomedical Engineering*, in press)

Agentic Al systems:

- > TrialGPT: patient-trial matching (*Nature Communications*, 2024)
- > GeneAgent: Gene set analysis with domain knowledge, Nαture Methods, 2025
- > AgentMD: automating medical risk calculation (Nature Communications, to appear)

Challenges & Limitations:

- Hidden flaws in multi-modal GPT4-V (npj Digital Medicine, 2024)
- Evaluation beyond Multiple-Choice Accuracy (*Annual Review of Biomedical Data Science*, 2025)
- Risks of Al Scientists: Prioritizing Safeguarding Over Autonomy (Nature Communications, in press)
- > AI model vulnerability under adversarial attacks (*Nαture Communicαtions*, to appear)