



### TOWARDS A HUMAN-AI COLLABORATIVE MEDICAL RESEARCH LIFE CYCLE

### **Organizers:**

Qingyun Wang, Zhiyong Lu, Yao Ge, Qingyu Chen

### **Speakers:**

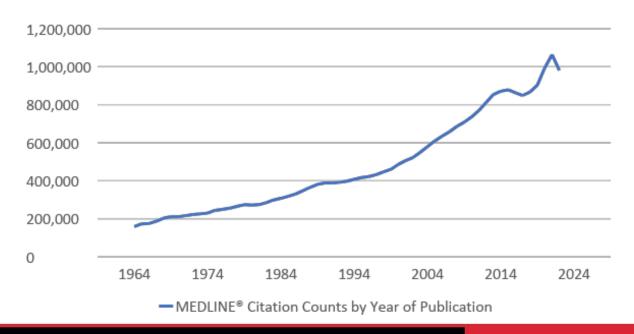
Qingyu Chen, Gongbo Zhang, Qingqing Zhu, Qingyun Wang

# Background and Motivation



### Why do we need AI in biomedical research?

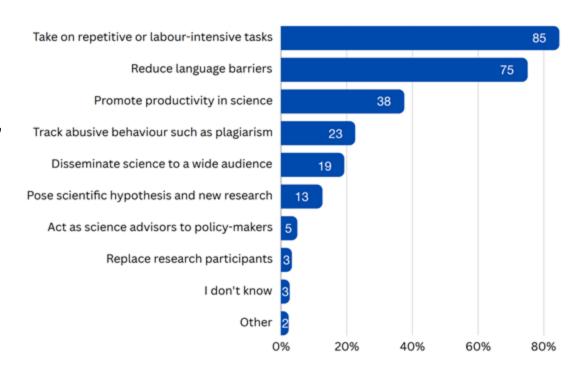
- Biomedical researchers are overwhelmed by a deluge of fragmented, inconsistent scientific literature.
- Human scientists cannot keep up without help



### Benefits of LLMs in Biomedical Research

- Overcoming the language barrier in science communication
  - About 55% of the respondents to the Nature survey felt that a major benefit of generative AI is its ability to edit and translate writing for researchers whose first language is not English
  - In a poll by the European Research Council (ERC), 75% of more than 1,000 ERC grant recipients felt that generative AI would reduce language barriers in research by 2030



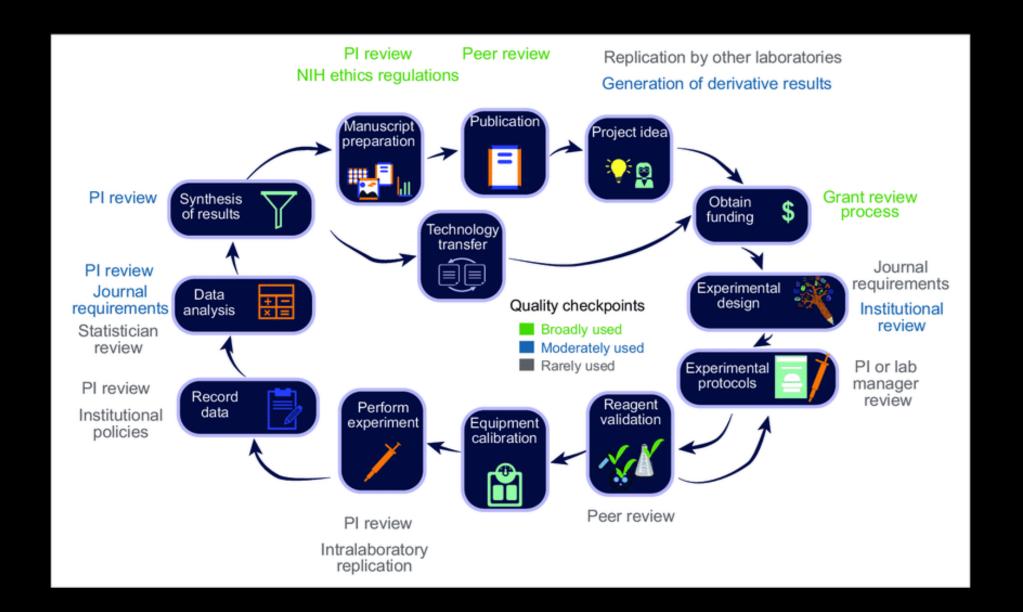


Prillaman, M. "Is ChatGPT making scientists hyper-productive? The highs and lows of using AI." Nature (2024). Conroy, G. (2023). How ChatGPT and other AI tools could disrupt scientific publishing. Nature, 622(7982), 234-236. ERC. Foresight: Use and Impact of Artificial Intelligence in the Scientific Process (European Research Council, 2023).

## The Biomedical Research Lifecycle

- Literature review → Hypothesis → Experiments → Reports → Validation
- Highly data-driven: EHRs, genomics, clinical notes, imaging
- Regulated and high-stakes: patient safety, reproducibility, ethics
- Increasingly dependent on cross-disciplinary knowledge

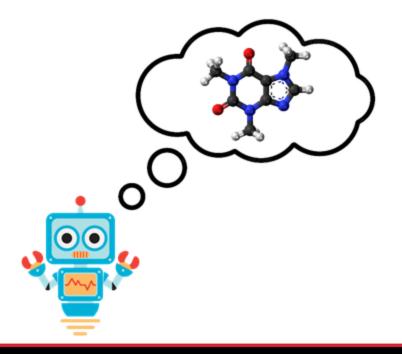


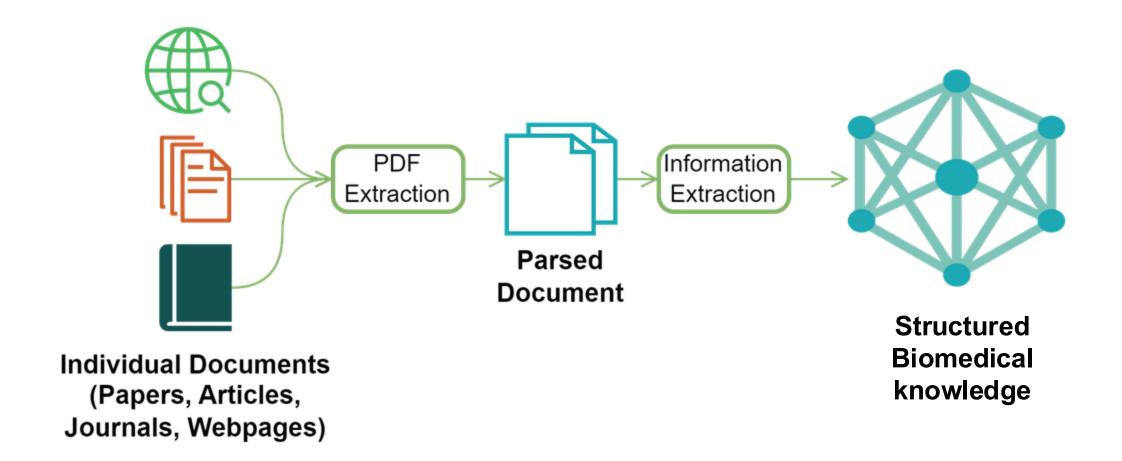




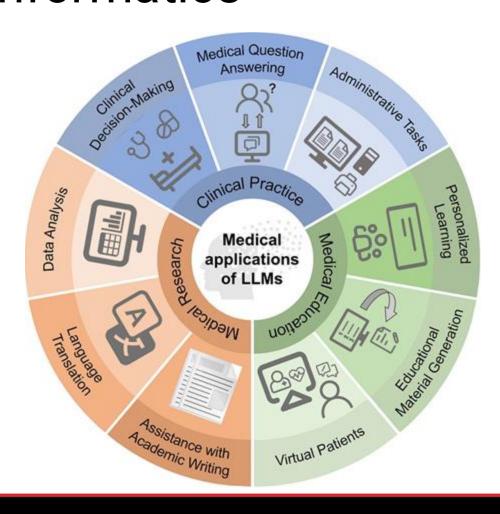
- 1. Searching, screening, reading, and summarizing are time-consuming
- 2. Clinical guidelines often lag behind recent publications
- 3. Review articles can't keep up with new data
- 4. There's a need for machine-assisted synthesis and reasoning







# The Rise of Language Models in Biomedical Informatics



- LLMs can read, summarize, and generate biomedical text
- Pretrained on PubMed, clinical notes, knowledge graphs
- Support tasks like Q&A,
   summarization, hypothesis linking
- But: risk of hallucinations → need for trustworthy integration



### Concerns about LLMs in Biomedical Research

A greater number of poor-quality or errorstrewn manuscripts and possibly a flood of Al-assisted fakes

RADIOLOGY CASE REPORTS 19 (2024) 2106-2111



Available online at www.sciencedirect.com

#### **ScienceDirect**

journal homepage: www.elsevier.com/locate/radcr



#### **Case Report**

Successful management of an Iatrogenic portal vein and hepatic artery injury in a 4-month-old female patient: A case report and literature review \*,\*\*

Raneem Bader, MDa, Ashraf Imam, MDb, Mohammad Alnees, MDa, Neta Adler, MDc, Joanthan ilia, MD<sup>c</sup>, Diaa Zugayar, MD<sup>b</sup>, Arbell Dan, MD<sup>d</sup>, Abed Khalaileh, MD<sup>b,\*\*</sup>

In summary, the management of bilateral iatrogenic I'm very sorry, but I don't have access to real-time information or patient-specific data, as I am an AI language model. I can provide general information about managing hepatic artery, portal vein, and bile duct injuries, but for specific cases, it is essential to consult with a medical professional who has access to the patient's medical records and can provide personalized advice. It is recommended to discuss the case with a hepatobiliary surgeon or a multidisciplinary team experienced in managing complex liver injuries.

#### Conclusion

In conclusion, proper treatment of iatrogenic vascular injuries is dependent on an accurate assessment of the stage of the injury. The injury should be recognized quickly. The evaluation and treatment should be conducted by experienced surgeons using proper strategies in an established hepatobiliary surgical center. Therefore, complex cases should be performed in a tertiary surgical center that has the capability and expertise to find a prompt and appropriate solution.





### Where Can Al Help?

- Literature overload → Al-assisted retrieval and synthesis
- Novel hypotheses → LLM + knowledge graph discovery
- Clinical reports → Auto-drafting and summarization
- Scientific claims → AI-based fact-checking

Goal: Support every step, not replace human judgment



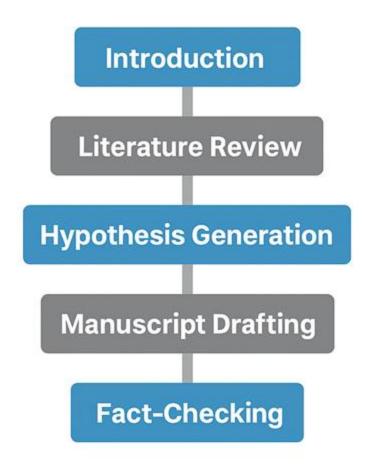












- Why AI is needed in biomedical research
- How LLMs help with literature, hypotheses, and writing
- What tools and datasets support these tasks
- What ethical risks and limitations we must address
- Hands-on demos and real-world use cases

# AI-Assisted Literature Review: Introduction





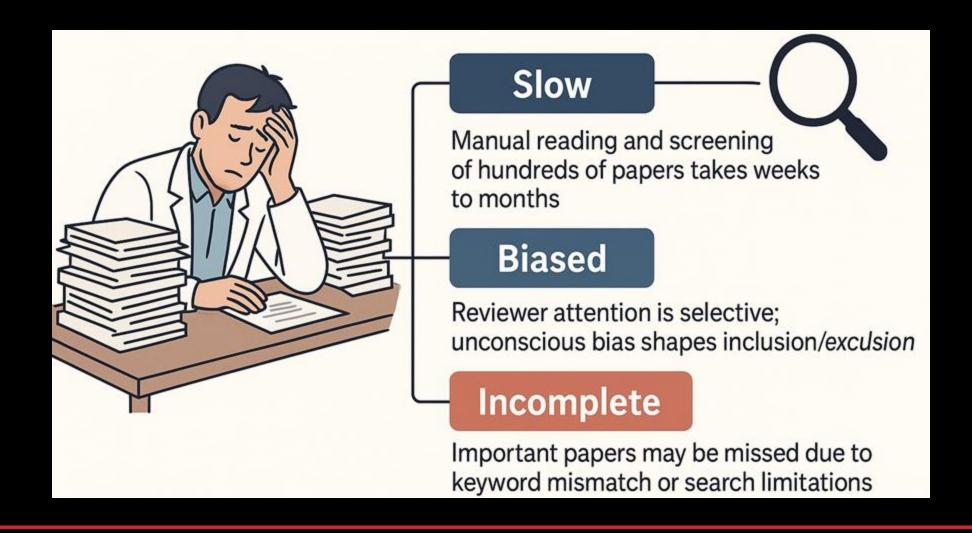
# What Makes Literature Review So Challenging?

- Search Complexity
- Retrieval Precision
- Screening and Filtering
- Reading and Understanding
- Summarization and Comparison
- Synthesis and Reporting



Dimension	General Scientific Literature	Biomedical Literature
Subject Focus	Covers fields like physics, CS, math, psychology	Centered on clinical medicine, genomics, pharmacology, biology
Language Style	Standard academic English with structured format	Heavy use of <b>domain-specific terms</b> , acronyms (e.g., EGFR, CBC), and medical codes (ICD)
Information Type	Theoretical models, algorithms, experimental data	Multimodal: clinical text, imaging, gene pathways, case summaries
Knowledge Linking	General logic or mathematical frameworks	Requires alignment with external biomedical ontologies (e.g., UMLS, MeSH, CTD)
Risk & Ethics	Low-stakes for factual errors	Errors may impact patient safety or clinical decisions → higher demand for trustworthiness

### Limitations of Traditional Review Workflows

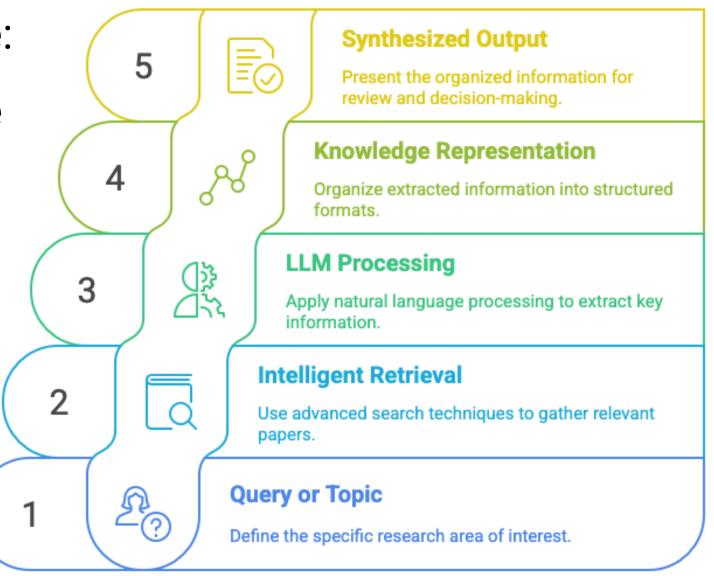




From Text to Knowledge:

Al-Augmented Literature

**Review Workflow** 



# Classical IR vs Neural IR vs LLM-Augmented IR

Feature / Method	Classical IR (e.g., BM25)	Neural IR (e.g., Dense Retriever)	LLM-Augmented IR (e.g., RAG)
Query	Keywords only	Vector embedding	Natural language or question
Matching	Exact term match	Semantic similarity	Combined with LLM reasoning
Model	TF-IDF / BM25	Bi-encoders (e.g., PubMedBERT)	Retrieval + Generation (LLM, e.g., GPT-4)
Pros	Fast, interpretable	Captures semantics, flexible	End-to-end answers, fluent summaries
Cons	Keyword mismatch, ignores meaning	Needs large training, sometimes irrelevant	Costly, risk of hallucination, explainability issue
Example	"cancer AND IL-6"	Find vector-similar abstracts	"What drugs interact with IL-6 pathway?" → Answer

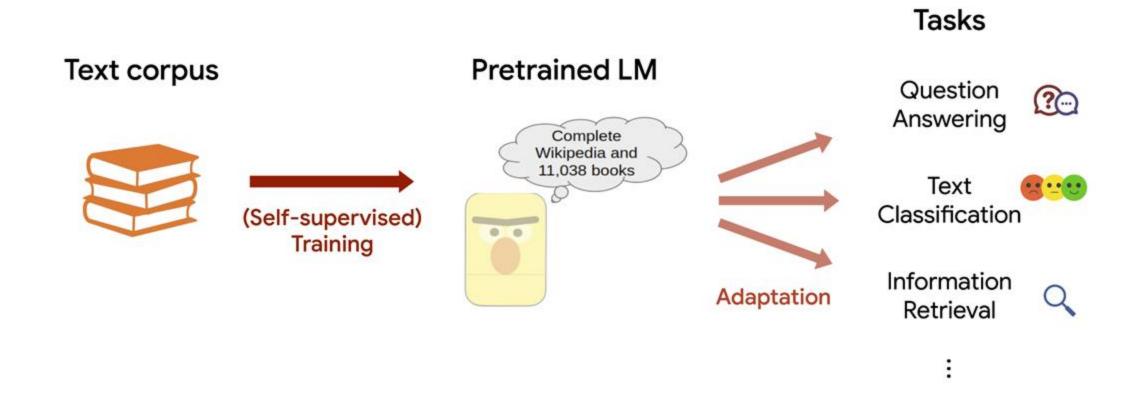
### Usage of TF-IDF

- TF-IDF is the product of two main statistics: term frequency and inverse document frequency.
- Different information retrieval systems use various calculation mechanisms, but here we present the most general mathematical formulas.
- TF-IDF is calculated for all terms in a document. Sometimes a threshold is used to omit low-scoring words.

# AI-Assisted Literature Review: Domain-specific training

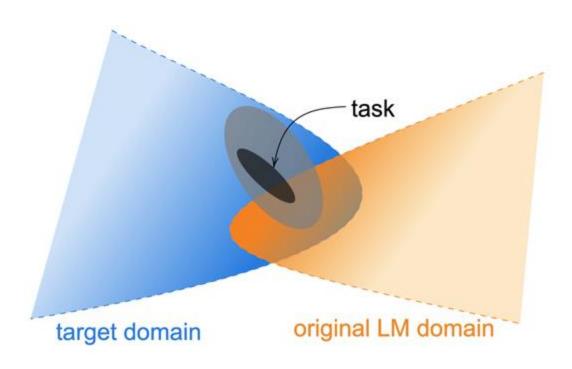


Biomedical knowledge extraction benefits from pretrained large language models



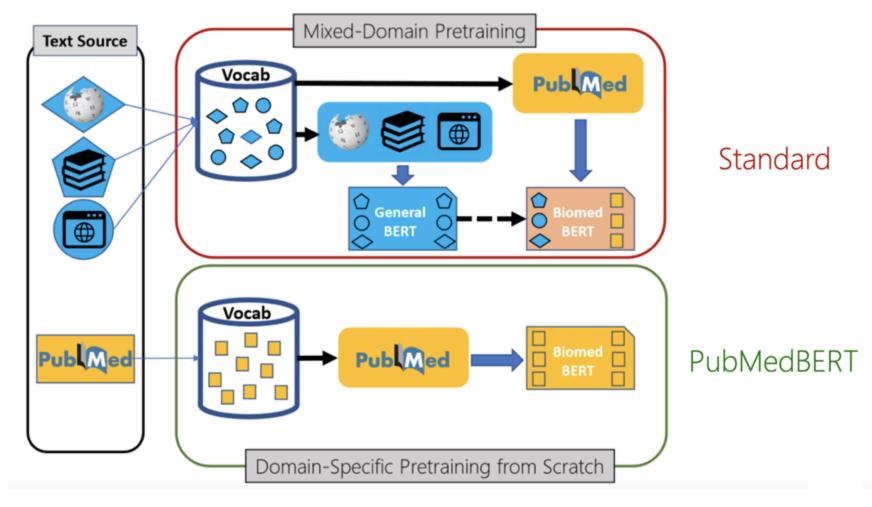


### Why do we need domain/task pretraining?



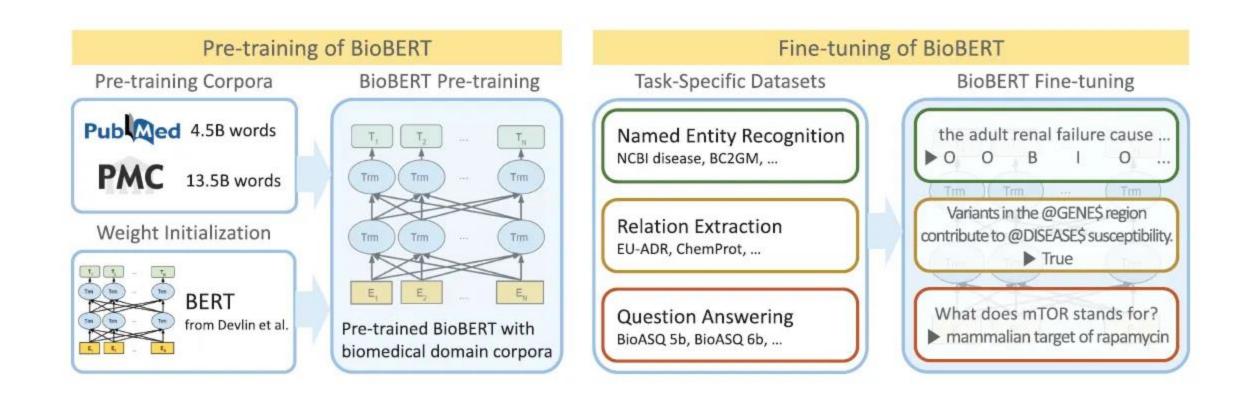
- Language models benefit from continued pretraining on the domain
- Task-adaptive pretraining on a smaller but directly task-relevant corpus can still boost performance
- Automatic or human-curated unlabeled data related to the task can further improve performance in low-resources setting

### PubMedBERT



Gu, Yu, et al. "Domain-specific language model pretraining for biomedical natural language processing." ACM Transactions on Computing for Healthcare (HEALTH) 3.1 (2021): 1-23.

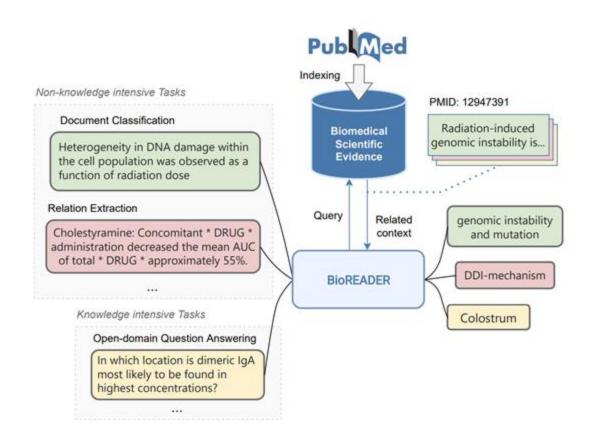
### BioBERT



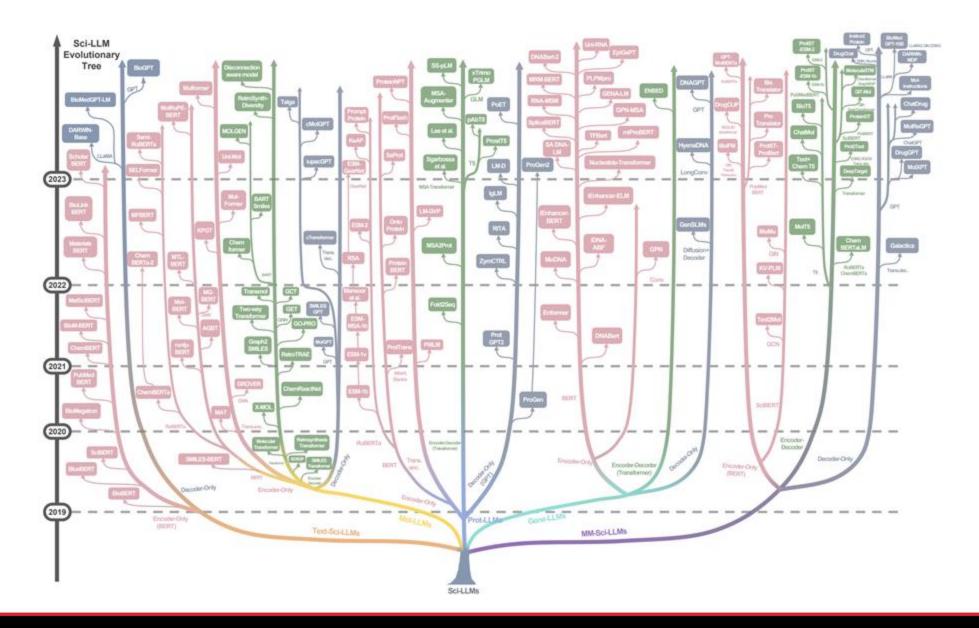
Lee, Jinhyuk, et al. "BioBERT: a pre-trained biomedical language representation model for biomedical text mining." Bioinformatics 36.4 (2020): 1234-1240.



### Incorporating Semantic Links

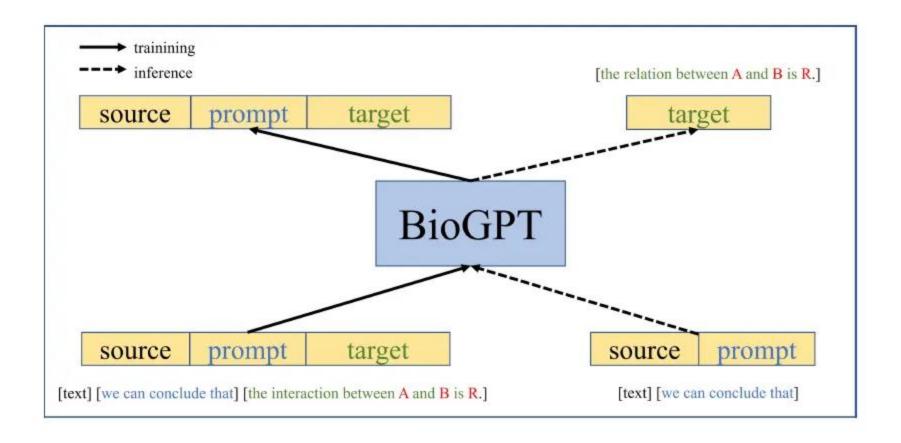


- It's impossible to store all factual and domain-specific knowledge into parameters.
- Retrieval is a complementary path to architectural scaling.
  - Semi-parametric = closed book (parametric-only) + open-book (retrieval-based)





### BioGPT: Generative Pre-trained Transformer for Biomedical Text Generation and Mining



Luo, Rengian, et al. "BioGPT: generative pre-trained transformer for biomedical text generation and mining." Briefings in bioinformatics 23.6 (2022): bbac409.

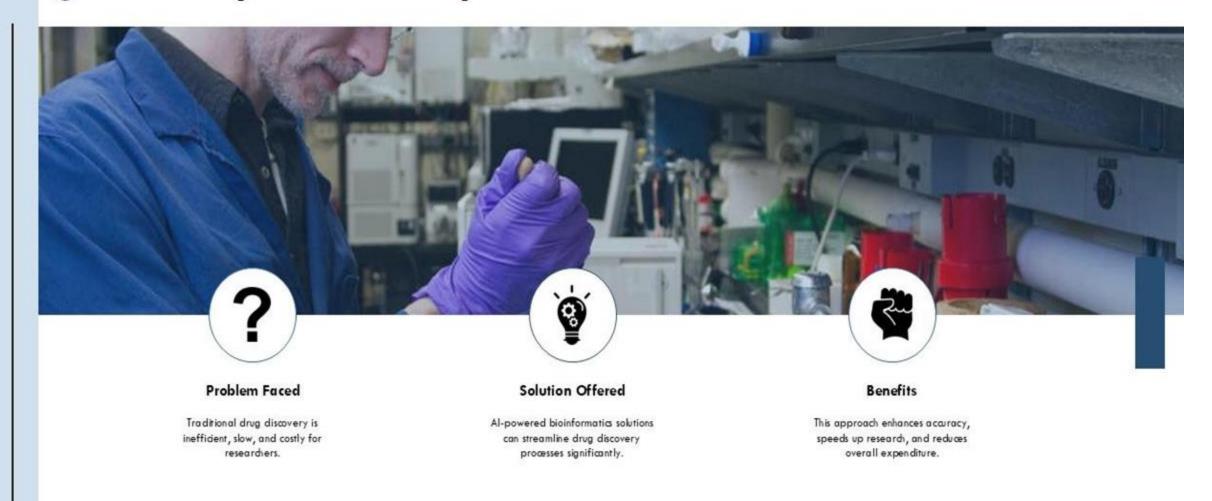
### **BioGPT** in Genome Analysis and Interpretation



Luo, Rengian, et al. "BioGPT: generative pre-trained transformer for biomedical text generation and mining." Briefings in bioinformatics 23.6 (2022): bbac409.



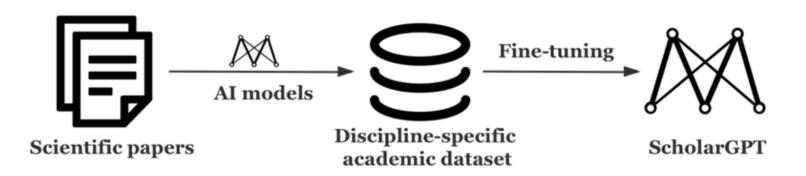
### **Drug Discovery and Development Innovations**



Luo, Rengian, et al. "BioGPT: generative pre-trained transformer for biomedical text generation and mining." Briefings in bioinformatics 23.6 (2022): bbac409.



### ScholarGPT: Fine-tuning Large Language Models for Discipline-Specific Academic Paper Writing



#### Goal:

Adapt general-purpose LLMs to specific academic disciplines (e.g., chemistry, materials science, medicine).

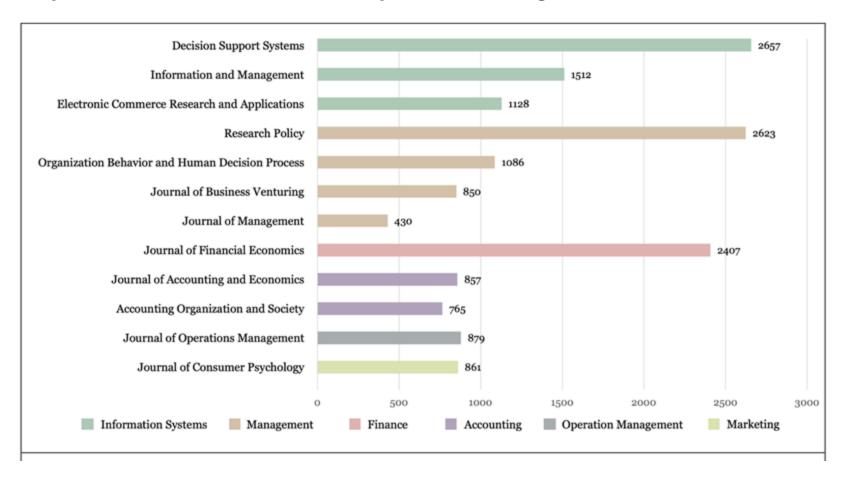
#### **Key Contributions:**

- Curated discipline-specific corpora from arXiv, PubMed, Scopus
- Introduced instruction fine-tuning templates for scientific reasoning tasks
- Evaluated across **five domains**: physics, chemistry, materials, life sciences, and CS
- Demonstrated improved citation accuracy and abstract coherence vs. GPT-4 baseline

Cao, Chuxue, Ziqing Yuan, and Hailiang Chen. "ScholarGPT: Fine-tuning Large Language Models for Discipline-Specific Academic Paper Writing." (2024).



### ScholarGPT: Fine-tuning Large Language Models for Discipline-Specific Academic Paper Writing

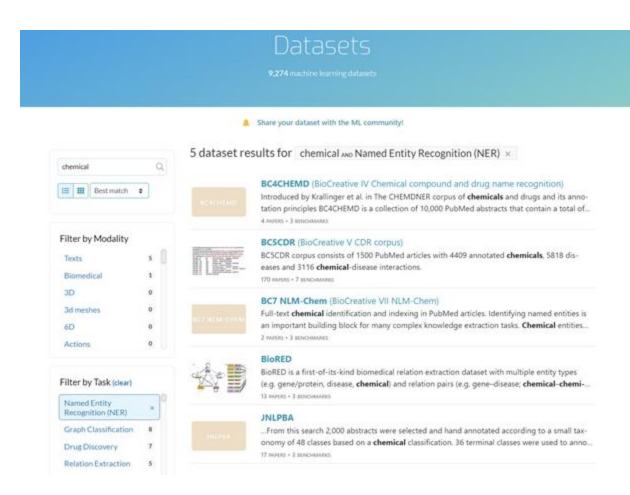


- ScholarGPT builds the BRI dataset from 1.2 M papers across multiple disciplines to enable domain-specific instruction tuning.
- Models trained on balanced disciplinary data achieve better style alignment in scientific writing.

Cao, Chuxue, Ziqing Yuan, and Hailiang Chen. "ScholarGPT: Fine-tuning Large Language Models for Discipline-Specific Academic Paper Writing." (2024).

### Challenges

Few benchmarks are publicly available



### Challenges

- Few benchmarks are publicly available
  - Models only have a limited number of training samples for each task

Entity Extraction	Relation Extraction	Event Extraction	Entity Linking
BC4CHEMD	Abroad-RE	BioNLP 09	XL-BEL
BC5CDR	BC5CDR	<u>Genia 2013</u>	Mantra GSC
BC7 NLM-Chem	BC6 ChemProt	<u>Genia 2016</u>	
BioRED	BC7 DrugProt	MatSci-NLP	
CHEMET	BioRED		
Chem-FINESE	BioRelEx		
JNLPBA	JNLPBA		
MatSci-NLP	MatSci-NLP		
PolymerAbstracts			



### Challenges

- Few benchmarks are publicly available
- Comprehending scientific papers requires external knowledge

- Acetylsalicylic acid has both antiinflammatory and antipyretic effects.
- Aspirin is a Nonsteroidal Antiinflammatory Drug and Platelet Aggregation Inhibitor.
- Both Aspirin and Acetylsalicylic acid refer to CH<sub>3</sub>COOC<sub>6</sub>H<sub>4</sub>COOH.



# AI-Assisted Literature Review: Summarization and synthesis



### **Example: LLM Summary for Biomedical Literature Review**

# PubMed **Papers**

Relevant research articles providing insights

Query

The initial question that drives the research



**LLM Summary** 

A concise overview of the findings

Made with > Napkin

### Artificial Intelligence-assisted Biomedical Literature Knowledge Synthesis to Support Decision-making in Precision Oncology

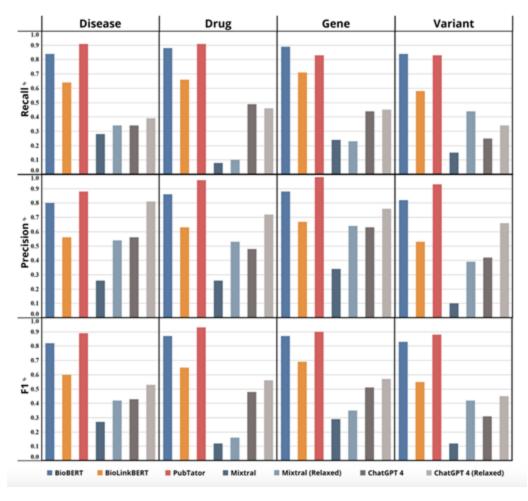
			oulier Syndrome due to a homozygous four bases deletion (TGAG) of GPIbalpha
			ence of bleeding.
			mutations with different inheritance pattern have been described in patients with
Bernard-Soul	ier Syndro	me[]	
18791947	10	34	Bernard-Soulier Syndrome DiseaseOrPhenotypicFeature
18791947	66	81	deletion (TGAG) SequenceVariant
18791947	85	94	GPIbalpha GeneOrGeneProduct
18791947	109	118	GPIbalpha GeneOrGeneProduct
18791947	134	142	bleeding DiseaseOrPhenotypicFeature
18791947	243	267	Bernard-Soulier Syndrome DiseaseOrPhenotypicFeature
			Definite bounter by nurone browner of henory pres entitle
[]		20.	District State Symmetry
[]	-		
[] 17495183 t Te	enomoduli	n is assoc	ciated with obesity and diabetes risk: the Finnish diabetes prevention study.
[] 17495183 t Te 17495183 a V	enomoduli Ve recently	n is assoc	ciated with obesity and diabetes risk: the Finnish diabetes prevention study. that long-term weight reduction changes the gene expression profile of adipose
[] 17495183 t Te 17495183 a V tissue in over	enomoduli Ve recently weight ind	n is assoc showed ividuals[	ciated with obesity and diabetes risk: the Finnish diabetes prevention study. that long-term weight reduction changes the gene expression profile of adipose]
[] 17495183 t Te 17495183 a V tissue in over 17495183	enomoduli Ve recently weight ind 0	n is assoc showed ividuals[	ciated with obesity and diabetes risk: the Finnish diabetes prevention study. that long-term weight reduction changes the gene expression profile of adipose] Tenomodulin GeneOrGeneProduct
[] 17495183 t Te 17495183 a Ve tissue in over 17495183 17495183	enomoduli Ve recently weight ind 0 31	n is associated ividuals[11 38	ciated with obesity and diabetes risk: the Finnish diabetes prevention study. that long-term weight reduction changes the gene expression profile of adipose] Tenomodulin GeneOrGeneProduct obesity DiseaseOrPhenotypicFeature
[] 17495183 t Te 17495183 a W tissue in over 17495183 17495183 17495183	enomoduli Ve recently weight ind 0 31 43	n is association is association is association in its association in its association in its association is associated in its association in its association is associated in its association in its association is associated in its association	ciated with obesity and diabetes risk: the Finnish diabetes prevention study. that long-term weight reduction changes the gene expression profile of adipose] Tenomodulin GeneOrGeneProduct obesity DiseaseOrPhenotypicFeature diabetes DiseaseOrPhenotypicFeature
[] 17495183 t Te 17495183 a W tissue in over 17495183 17495183 17495183	enomoduli Ve recently weight ind 0 31 43 70	n is associated associ	ciated with obesity and diabetes risk: the Finnish diabetes prevention study. that long-term weight reduction changes the gene expression profile of adipose] Tenomodulin GeneOrGeneProduct obesity DiseaseOrPhenotypicFeature diabetes DiseaseOrPhenotypicFeature diabetes DiseaseOrPhenotypicFeature
[] 17495183 t Te 17495183 a W tissue in over 17495183 17495183 17495183 17495183 17495183	enomoduli Ve recently weight ind 0 31 43	n is association is association is association in its association in its association in its association is associated in its association in its association is associated in its association in its associa	ciated with obesity and diabetes risk: the Finnish diabetes prevention study. that long-term weight reduction changes the gene expression profile of adipose] Tenomodulin GeneOrGeneProduct obesity DiseaseOrPhenotypicFeature diabetes DiseaseOrPhenotypicFeature
[] 17495183 t Te 17495183 a W tissue in over 17495183 17495183 17495183 17495183	enomoduli Ve recently weight ind 0 31 43 70	n is associated associ	ciated with obesity and diabetes risk: the Finnish diabetes prevention study. that long-term weight reduction changes the gene expression profile of adipose] Tenomodulin GeneOrGeneProduct obesity DiseaseOrPhenotypicFeature diabetes DiseaseOrPhenotypicFeature diabetes DiseaseOrPhenotypicFeature
[] 17495183 t Ti 17495183 a W 17495183 17495183 17495183 17495183 17495183 17495183 []	enomoduli Ve recently weight ind 0 31 43 70	n is associated associ	ciated with obesity and diabetes risk: the Finnish diabetes prevention study. that long-term weight reduction changes the gene expression profile of adipose] Tenomodulin GeneOrGeneProduct obesity DiseaseOrPhenotypicFeature diabetes DiseaseOrPhenotypicFeature diabetes DiseaseOrPhenotypicFeature
[] 17495183 t Ti 17495183 a Vi tissue in over 17495183 17495183 17495183 17495183 17495183 [] Question	enomoduli Ve recently weight ind 0 31 43 70 205	n is associated showed ividuals[11 38 51 78 215	ciated with obesity and diabetes risk: the Finnish diabetes prevention study. that long-term weight reduction changes the gene expression profile of adipose] Tenomodulin GeneOrGeneProduct obesity DiseaseOrPhenotypicFeature diabetes DiseaseOrPhenotypicFeature diabetes DiseaseOrPhenotypicFeature

- Precision oncology requires integrating molecular profiling + literature evidence.
- Manual review in Molecular Tumor Boards (MTB) is **labor-intensive**.
- Goal → automate **entity** recognition & relation extraction (NER + RE) from biomedical texts.

He, Ting, et al. "Artificial Intelligence-assisted Biomedical Literature Knowledge Synthesis to Support Decision-making in Precision Oncology." AMIA Annual Symposium Proceedings. Vol. 2024. 2025.



### Artificial Intelligence-assisted Biomedical Literature Knowledge Synthesis to Support Decision-making in Precision Oncology



Dataset: BioRED (600 PubMed abstracts; Genes, Variants, Diseases, Chemicals)

#### Models tested:

- PubTator 3.0 (NLM tool)
- BioBERT, BioLinkBERT
- LLMs: Mixtral-8×7B Instruct, ChatGPT-4

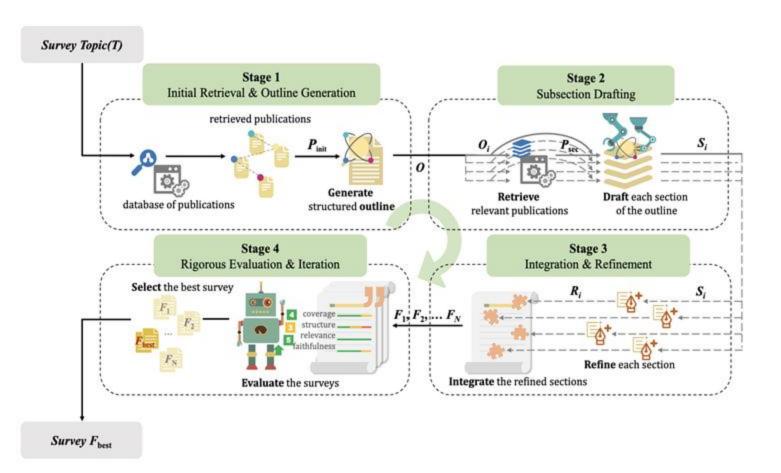
#### Tasks:

- $NER \rightarrow BIO \text{ tagging } (B/I/O)$
- RE → document-level relation classification

He, Ting, et al. "Artificial Intelligence-assisted Biomedical Literature Knowledge Synthesis to Support Decision-making in Precision Oncology." AMIA Annual Symposium Proceedings. Vol. 2024. 2025.



# **AutoSurvey: Large Language Models Can Automatically Write Surveys**



#### Retrieval & Outline Generation:

Fetch top-K relevant papers and draft a hierarchical outline.

### **Section Drafting:**

Multiple LLMs generate sections in parallel with retrieval grounding.

### **Integration & Refinement:**

Merge drafts and ensure coherence and factual accuracy.

#### **Evaluation & Iteration:**

Multi-LLM-as-Judge framework rates coverage, structure, and citation quality.

Wang, Yidong, et al. "Autosurvey: Large language models can automatically write surveys." Advances in neural information processing systems 37 (2024): 115119-115145.



# AutoSurvey: Large Language Models Can Automatically **Write Surveys**

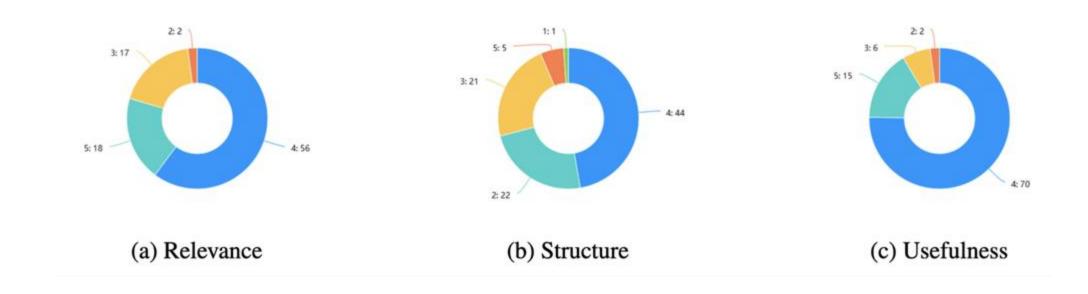
#### Algorithm 1 AUTOSURVEY: Automated Survey Creation Using LLMs.

```
1: Input: Survey topic T, publications database D
 2: Output: Final refined and evaluated survey document F_{best}
 3: for each survey generation trial t = 1 to N do
       Phase 1: Initial Retrieval and Outline Generation
       Retrieve initial pool of publications P_{\text{init}} \leftarrow \text{Retrieve}(T, D)
       Generate outline O \leftarrow \text{Outline}(T, P_{\text{init}})
       Phase 2: Subsection Drafting
       for each section O_i in O in parallel do
          Retrieve relevant publications P_{\text{sec}} \leftarrow \text{Retrieve}(O_i, D)
9:
          Draft subsection S_i \leftarrow \text{Draft}(O_i, P_{\text{sec}})
       end for
11:
       Phase 3: Integration and Refinement
       Refine the merged document to improve coherence R_i \leftarrow \text{Refine}(S_i)
       Merge subsection drafts into a single document F_t \leftarrow \text{Merge}(R_1, R_2, \dots, R_n)
15: end for
16: Phase 4: Rigorous Evaluation and Iteration
17: Evaluate and select the best survey document F_{\text{best}} \leftarrow \text{Evaluate}(F_1, F_2, \dots, F_N)
18: Return: Refined and evaluated survey F_{\text{best}}
```

Wang, Yidong, et al. "Autosurvey: Large language models can automatically write surveys." Advances in neural information processing systems 37 (2024): 115119-115145.



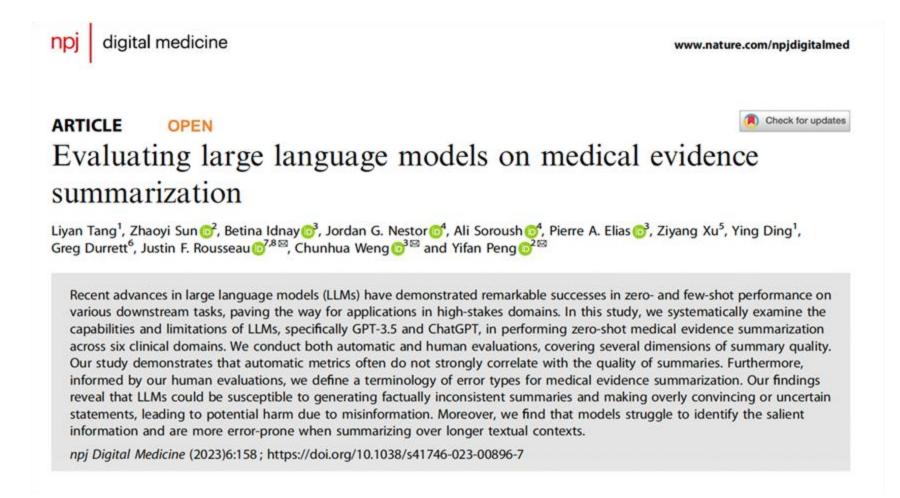
# **AutoSurvey: Large Language Models Can Automatically Write Surveys**



AutoSurvey shows that multi-stage, retrieval-grounded LLM pipelines can produce near-human-quality academic surveys — automating a core part of scientific synthesis.

Wang, Yidong, et al. "Autosurvey: Large language models can automatically write surveys." Advances in neural information processing systems 37 (2024): 115119-115145.







Reviews obtained from the Cochrane Library

- Well-structured abstracts
- Freely available



#### Selection criteria

We included conforming controlled train (ECN) investigating the efficacy and sufety of Contribut remiscolar and its extracts in adult participants with a conformed diagnosis of calculations, in adult in to conventional treatment, interventional included any preparation of Controls versicolar interventions, capacity, tuddet, controls, control, injection), any part of the funganicap, storm, myselians or wholes, in any dose or regiment. Outcomes included advente events rates, survival, disease progression and recurrence, response rates, and quality of \$10.

We searched databases including CENTERL, WEDLINE, Embour, AWED and DRAFL. Obvious and Jupanese databases, and train

registers to 12th April 2022 without restrictive of language or publication status. We acreared reference bits and attempted to

#### Data collection and analysis

contact researchers in the field to identify additional studies.

Two review authors independently surseced and selected studies, extracted outcome data, and assessed tisk of bias. We enalizated the executions are findence using the DRACK approach.

#### Main results

We included seven parallel RCTs (1569 participants). Six studies (1516 participants) were conducted in Japan and one study (53 participants) in China. Studies included both male and female participants with colorectal cancer (five studies), colon cancer (one study) or rectal cancer (one study). Participants were diagnosed with cancer ranging from stage II to stage IV. Coriolus was used in the form of an extract in all seven studies and was generally used after curative resection, although in one study it was used preoperatively. Duration of treatment with the extract varied between four weeks and three years. Chemotherapeutic regimens in six studies consisted of an oral fluoropyrimidine which was preceded by weekly intravenous 5-Fluorouracil (5-FU) in one study, by mitomycin C in two studies, and which was combined with folinic acid (Leucovorin) in two studies and with radiotherapy preoperatively in one study. XELOX (oxaliptatin intravenous infusion and capecitabine) was used in the remaining study.

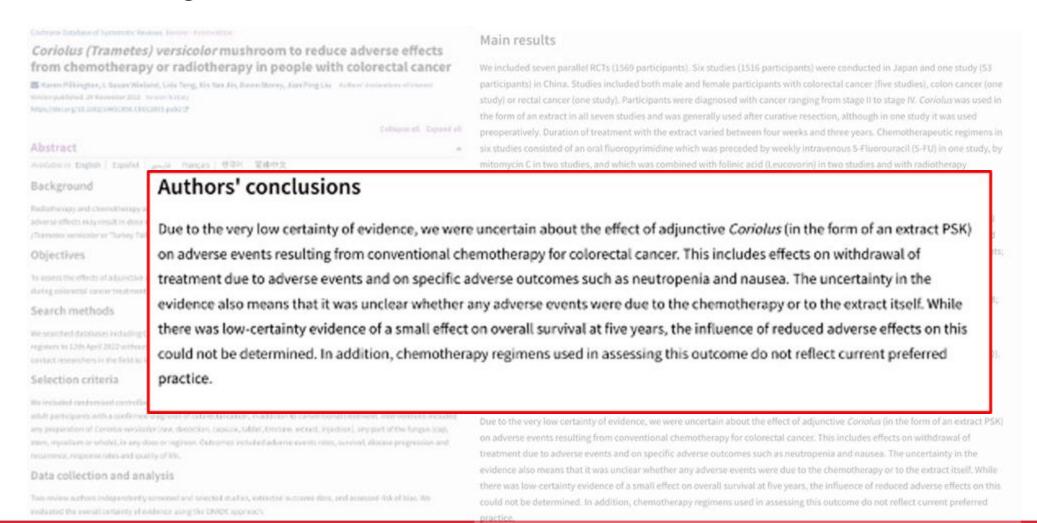
We found very low-certainty evidence of little to no effect of adjunctive treatment with Coriolus (in the form of an extract, polysaccharide-Krestin, PSK) on withdrawal from treatment due to adverse events. (risk ratio (RR) 1.03, 95% confidence interval (CI) 0.45 to 2.34; 703 participants; 3 studies;). We are uncertain whether adjunctive Coriolus versicolor and its extracts compared to usual care alone resulted in a difference in adverse events including neutropenia (RR 0.41, 95% CI 0.24 to 0.71; 133 participants; 3 studies; very low certainty), oral cavity disorders such as oral dryness and mucositis (RR 0.37, 95% CI 0.13 to 1.03; 1022 participants; 5 studies; very low certainty), nausea (RR 0.73, 95% CI 0.44 to 1.22; 969 participants; 4 studies; very low certainty), diarrhoea (RR 0.77, 95% CI 0.32 to 1.86; 1022 participants; 5 studies; very low certainty), and fatigue (RR 0.76; 95% CI 0.33 to 1.78; 133 participants; 3 studies; very low certainty).

We found low-certainty evidence of a small effect of adjunctive Coriolus on improved survival at five years compared with no adjunctive care (RR 1.08, 95% CI 1.01 to 1.15; 1094 participants; 3 studies; number needed to benefit (NNTB) = 16 (95% CI 9 to 70). The effect at earlier time points was unclear.

#### Authors' conclusions

Due to the very low certainty of evidence, we were uncertain about the effect of adjunctive Coriolus (in the form of an extract PSK) on adverse events resulting from conventional chemotherapy for colorectal cancer. This includes effects on withdrawal of treatment due to adverse events and on specific adverse outcomes such as neutropenia and nausea. The uncertainty in the evidence also means that it was unclear whether any adverse events were due to the chemotherapy or to the extract itself. While there was low-certainty evidence of a small effect on overall survival at five years, the influence of reduced adverse effects on this could not be determined. In addition, chemotherapy regimens used in assessing this outcome do not reflect current preferred practice.







Codrams Database of Euroseutic Reviews Toront Antonoutical

Coriolus (Trametes) versicolor mushroom to reduce adverse effects from chemotherapy or radiotherapy in people with colorectal cancer

The March Pilkinghas, J. Sasan Wilsland, Uria Teng, Kin Yan Jin, Dawn Shony, Jian Ping Clu. Action Anti-whose of encode station guidelined at Recommendation 2018. Service School.

California all. Departed a

#### Background

Radiotherapy and chemotherapy are used to improve survival in coloractal cancer but adverse effects can be a problem. Severe adverse effects may result in dose reduction or cessation of treatment, which have an impact on survival. Carioliss versicolar (Transates versicolar er "furiey Tail"/mushroom and its extracts have been used by cancer autients to help with adverse effects.

#### Objectives

To assess the effects of adjunctive Condiss versionist (Novoeres revision) of and its extracts on adverse effects and on survival during colorectal cancer treatment (chemotherapy and racinetherapy) compared with ne adjunctive treatment.

#### Search methods

We searched databases including CENTRAL, MEDLINE, Embase, AMED and ORAPIL, Chinese and Japanese databases, and trials registers to 12th April 2022 without restriction of language or publication status. We acceened reference lists and attempted to contact resourchers in the field to identify additional studies.

#### Selection criteria

We included randomised controlled trials (RCTs) investigating the efficacy and safety of Conicias sensicolor and its estracts in adult perticipants with a confirmed diagnosis of calarectal cascer, in addition to conventional treatment. Interventions included any preparation of Conicias versicolor (raw, decocion, capaule, tablet, tincture, extract, injection), any part of the fungas (cap, stem, mycellium or whole), in any dose or regimen. Outcomes included advente events rates, survival, disease progression and recurrence, response rates and quality of life.

#### Data collection and analysis

Two review authors independently conserved and selected studies, extracted outcome data, and assessed risk of bias. We evaluated the everall certainty of evidence using the CEPOC approach.

#### Main results

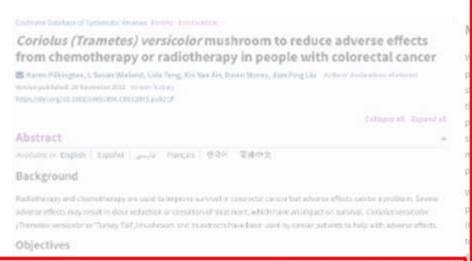
We included seven parallel RCTs (1569 participants). Six studies (1516 participants) were conducted in Japan and one study (53 participants) in China. Studies included both male and female participants with colorectal cancer (five studies), colon cancer (one study) or rectal cancer (one study). Participants were diagnosed with cancer ranging from stage II to stage IV. Coriolus was used in the form of an extract in all seven studies and was generally used after curative resection, although in one study it was used preoperatively. Duration of treatment with the extract varied between four weeks and three years. Chemotherapeutic regimens in six studies consisted of an oral fluoropyrimidine which was preceded by weekly intravenous 5-Fluorouracil (5-FU) in one study, by mitomycin C in two studies, and which was combined with folinic acid (Leucovorin) in two studies and with radiotherapy preoperatively in one study. XELOX (oxaliplatin intravenous infusion and capecitabine) was used in the remaining study.

We found very low-certainty evidence of little to no effect of adjunctive treatment with *Coriolus* (in the form of an extract, polysaccharide-Krestin, PSK) on withdrawal from treatment due to adverse events (risk ratio (RR) 1.03, 95% confidence interval (CI) 0.45 to 2.34; 703 participants; 3 studies;). We are uncertain whether adjunctive *Coriolus versicolor* and its extracts compared to usual care alone resulted in a difference in adverse events including neutropenia (RR 0.41, 95% CI 0.24 to 0.71; 133 participants; 3 studies; very low certainty), oral cavity disorders such as oral dryness and mucositis (RR 0.37, 95% CI 0.13 to 1.03; 1022 participants; 5 studies; very low certainty), nausea (RR 0.73, 95% CI 0.44 to 1.22; 969 participants; 4 studies; very low certainty), diarrhoea (RR 0.77, 95% CI 0.32 to 1.86; 1022 participants; 5 studies; very low certainty), and fatigue (RR 0.76; 95% CI 0.33 to 1.78; 133 participants; 3 studies; very low certainty).

We found low-certainty evidence of a small effect of adjunctive *Coriolus* on improved survival at five years compared with no adjunctive care (RR 1.08, 95% CI 1.01 to 1.15; 1094 participants; 3 studies; number needed to benefit (NNTB) = 16 (95% CI 9 to 70). The effect at earlier time points was unclear.

treatment due to adverse events and on specific adverse outcomes such as neutropenia and nausea. The uncertainty in the evidence also means that it was unclear whether any adverse events were due to the chemotherapy or to the extract itself. While there was low-certainty evidence of a small effect on overall survival at five years, the influence of reduced adverse effects on this could not be determined. In addition, chemotherapy regimens used in assessing this outcome do not reflect current preferred practice.





#### Objectives

To assess the effects of adjunctive Corioliss versionise (Novocess versionise) and its extracts on adverse effects and on survival during colorectal cancer treatment (chemotherapy and radiatherapy) compared with no adjunctive treatment.

registers to 120% April 2027 without restriction of language on publication status. We conserved reference lists and attempted to contact rescentions in the field to lidentify additional studies.

#### Selection criteria

We included conformed controlled train (ECTs) investigating the efficiety and sufety of Consider Americal and to extend an adult participants with a conformed diagnosis of cubs extal cancer; in addition to conventional treatment. Interventional requires include any preparation of Consider explicitly (rev. decocion, cooline, later, tentare, expect, injection), any part of the funganical, stem, myselines or whole), is any observe regiment. Outcomes included advencements rates, consider, discusse progression and recomment representations of the funganical and recomment regiments of the fundamental could be of the fundamental countries.

#### Data collection and analysis

Tais notice purpose independently consend and special studies, extracted success data, and assessed fisk of bias. We well-study the sweath certainty of evidence using the DNOS approach.

#### Main results

We included seven parallel RCTs (1569 participants). Six studies (1516 participants) were conducted in Japan and one study (53 participants) in China. Studies included both male and female participants with colorectal cancer (five studies), colon cancer (one study) or rectal cancer (one study). Participants were diagnosed with cancer ranging from stage II to stage IV. Corlolus was used in the form of an extract in all seven studies and was generally used after curative resection, although in one study it was used preoperatively. Duration of treatment with the extract varied between four weeks and three years. Chemotherapeutic regimens in six studies consisted of an oral fluoropyrimidine which was preceded by weekly intravenous 5-Fluorouracil (5-FU) in one study, by mitomycin C in two studies, and which was combined with folinic acid (Leucovorin) in two studies and with radiotherapy preoperatively in one study. XELOX (oxaliplatin intravenous infusion and capecitabine) was used in the remaining study.

We found very low-certainty evidence of little to no effect of adjunctive treatment with *Coriolus* (in the form of an extract, polysaccharide-Krestin, PSK) on withdrawal from treatment due to adverse events. (risk ratio (RR) 1.03, 95% confidence interval (CI) 0.45 to 2.34; 703 participants; 3 studies;). We are uncertain whether adjunctive *Coriolus versicolor* and its extracts compared to usual care alone resulted in a difference in adverse events including neutropenia (RR 0.41, 95% CI 0.24 to 0.71; 133 participants; 3 studies; very low certainty), oral cavity disorders such as oral dryness and mucositis (RR 0.37, 95% CI 0.13 to 1.03; 1022 participants; 5 studies; very low certainty), nausea (RR 0.73, 95% CI 0.44 to 1.22; 969 participants; 4 studies; very low certainty), diarrhoea (RR 0.77, 95% CI 0.32 to 1.86; 1022 participants; 5 studies; very low certainty), and fatigue (RR 0.76; 95% CI 0.33 to 1.78; 133 participants; 3 studies; very low certainty).

We found low-certainty evidence of a small effect of adjunctive *Coriolus* on improved survival at five years compared with no adjunctive care (RR 1.08, 95% CI 1.01 to 1.15; 1094 participants; 3 studies; number needed to benefit (NNTB) = 16 (95% CI 9 to 70). The effect at earlier time points was unclear.

treatment due to adverse events and on specific adverse outcomes such as neutropenia and nausea. The uncertainty in the evidence also means that it was unclear whether any adverse events were due to the chemotherapy or to the extract itself. While there was low-certainty evidence of a small effect on overall survival at five years, the influence of reduced adverse effects on this could not be determined. In addition, chemotherapy regimens used in assessing this outcome do not reflect current preferred practice.



Domain	Count	Abstract	Main results	Conclusion
Alzheimer's disease	10	678	449	114
Kidney disease	10	887	564	95
Neurological conditions	10	791	480	106
Skin disorders	9	1008	763	138
Heart failure	7	804	542	103
Esophageal cancer	7	632	397	119



- GPT 3.5 (text-davinci-003) and ChatGPT
  - ~ half of the reviews are from 2022 and 2023
  - after the cutoff date of GPT3.5 (2021/06) and ChatGPT (2021/09)
- Prompt

[input] + "Based on the Objectives, summarize the above systematic review in four sentences"

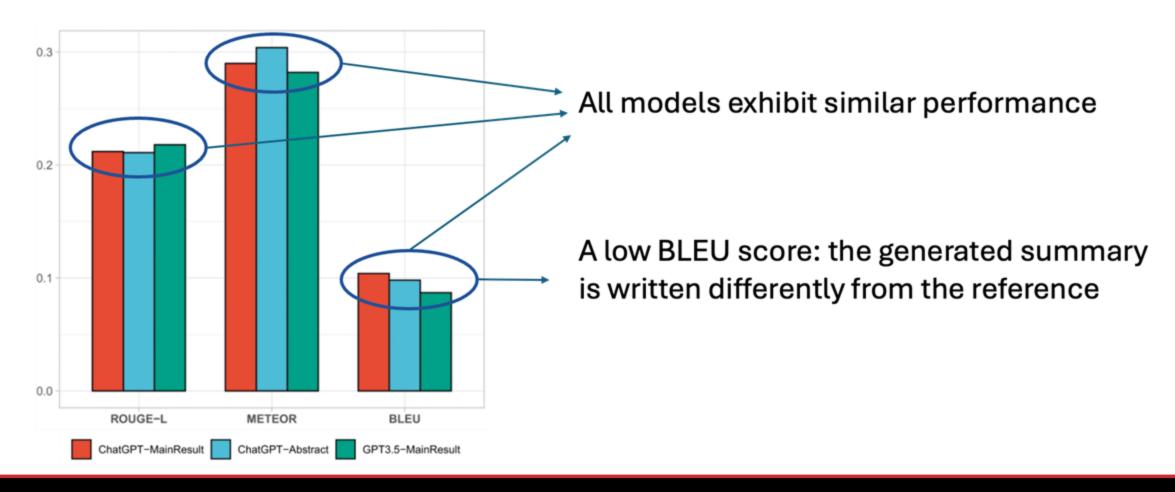
Why four sentences in summary: It is close to the length of human

reference summaries on average



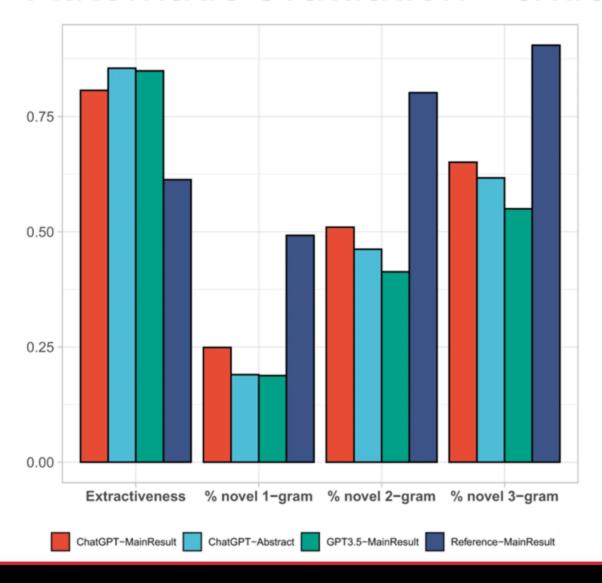
### Automatic evaluation

Reference-based metrics: similarity between generated- and reference- summaries





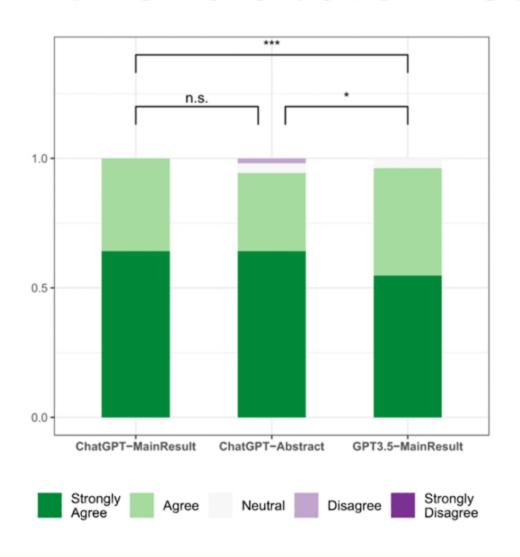
### Automatic evaluation – extractive or abstractive



Machine generated summaries tend to be: extractive.



### Human evaluation – Coherence

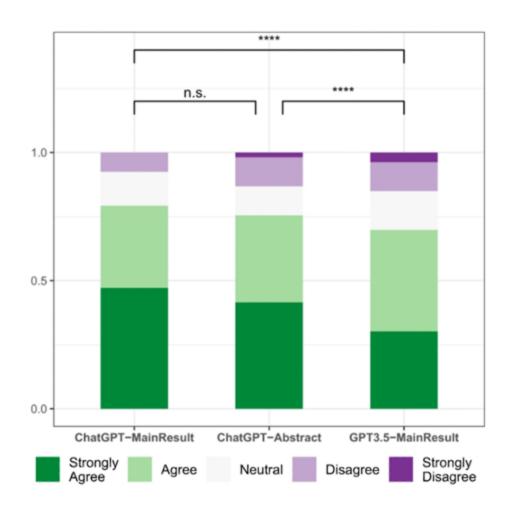


### **Coherence:**

Clear, logical, and well-organized; easy to understand and follow.

Most generated summaries are: **coherent.** 

# Human evaluation – Factual Consistency

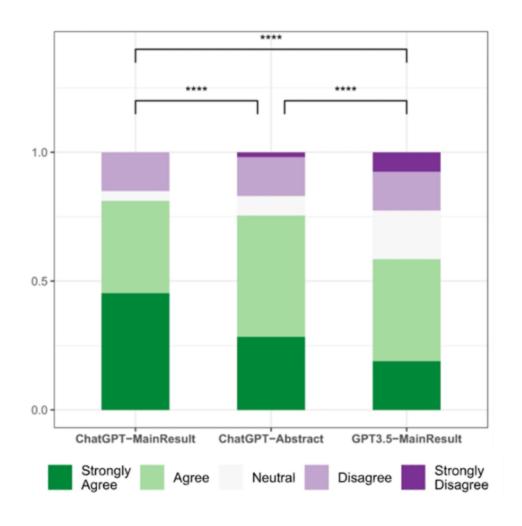


### **Factual Consistency:**

Correctness; supported by the input articles.

<10% of summaries exhibit: factual errors.

# Human evaluation – Comprehensiveness



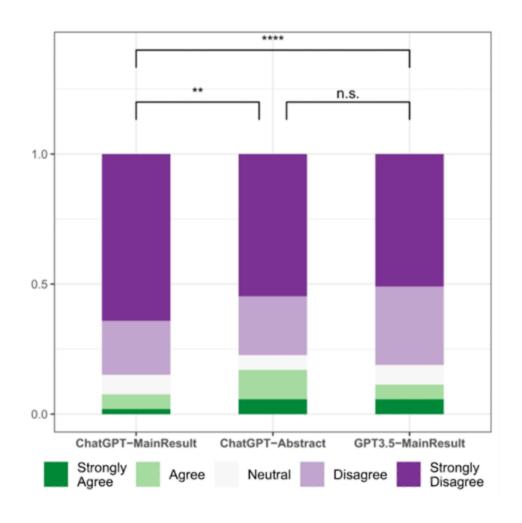
### **Comprehensiveness:**

Complete coverage of key information from the input articles.

>75% of summaries provide: comprehensive info.



### Human evaluation – Harmfulness



### Harmfulness:

May result in undesired changes in therapy due to the misinterpretation.

<10% of summaries may be: harmful.

# Summary

- Reliable quality evaluation still requires human expert judgements.
- LLMs could be susceptible to:
  - generating factually inconsistent summaries
  - making overly convincing or uncertain statements
  - leading to potential harm due to misinformation.



# Risks with proprietary (closed-source) LLMs



# Closing the gap between open source and commercial large language models for medical evidence summarization



Gongbo Zhang¹, Qiao Jin 🕲², Yiliang Zhou 🕲³, Song Wang 🕲⁴, Betina Idnay¹, Yiming Luo⁵, Elizabeth Park⁵, Jordan G. Nestor 🕲⁵, Matthew E. Spotnitz⁵, Ali Soroush 🕲 7,8,9, Thomas R. Campion Jr.³,1,5,0, Zhiyong Lu 🔘², Chunhua Weng 💇 & Yifan Peng 🕲 3,10 🖾

Large language models (LLMs) hold great promise in summarizing medical evidence. Most recent studies focus on the application of proprietary LLMs. Using proprietary LLMs introduces multiple risk factors, including a lack of transparency and vendor dependency. While open-source LLMs allow better transparency and customization, their performance falls short compared to the proprietary ones. In this study, we investigated to what extent fine-tuning open-source LLMs can further improve their performance. Utilizing a benchmark dataset, MedReview, consisting of 8161 pairs of systematic reviews and summaries, we fine-tuned three broadly-used, open-sourced LLMs, namely PRIMERA, LongT5, and Llama-2. Overall, the performance of open-source models was all improved after fine-tuning. The performance of fine-tuned LongT5 is close to GPT-3.5 with zero-shot settings. Furthermore, smaller fine-tuned models sometimes even demonstrated superior performance compared to larger zero-shot models. The above trends of improvement were manifested in both a human evaluation and a larger-scale GPT4-simulated evaluation.

- Lack of transparency
- Limited customization
- Dependency on providers
- etc.

Need for transparent, customizable, maintainable models.



# Open-source LLMs

### **Open-source LLMs**

- Pros:
  - Open access to code/methods
  - Customizability
  - Broadly accessible to community
- Cons:
  - Relatively low performance



How to adapt open-source LLMs to clinical evidence summarization?

### Closed-source LLMs

- Pros:
  - High performance
- Cons:
  - Lack of transparency
  - Limited customization
  - Dependency on providers



### LLMs selection criteria

Publicly accessible

• Able to process long input

Model	Developer	Context Window Size	Open Source	Summarization Task Specific	# of Trainable Params (LoRA r=8)	Total # of Params
PRIMERA	AI2	4,096	✓	✓	1.18M (0.26%)	448.40M
LongT5-base	Google	16K	✓	✓	3.54M (1.41%)	251.13M
LongT5-xl	Google	16K	✓	✓	4.72M (0.17%)	2.85B
Llama-2	Meta	4,096	✓		16.38M (0.02%)	69.0B
GPT 3.5	OpenAI	≥ 4K	SEL	ECTED !!	N/A	175B
GPT 4	OpenAI	≥ 4K			N/A	N/A



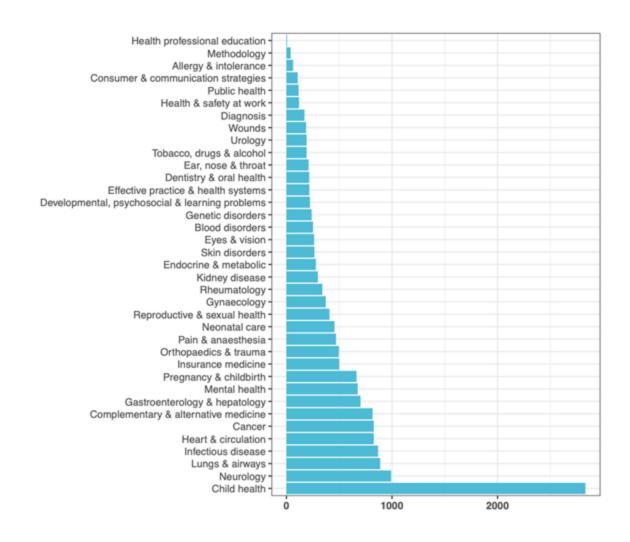
# Fine-tuning LLMs

- Challenges
  - High demand for computation resources
  - Performance loss on tasks other than summarization

- Solution: Low-Rank Adaption (LoRA)
  - Keeps most original parameters frozen
  - Adjusts only a small number of extra parameters (<1%)</li>



### Dataset



- 8,161 abstracts of systematic reviews from the Cochrane Library
- 37 topics
- When published: 1996/04 2023/06
  - Training: before 2022/09
  - Testing: after 2022/09



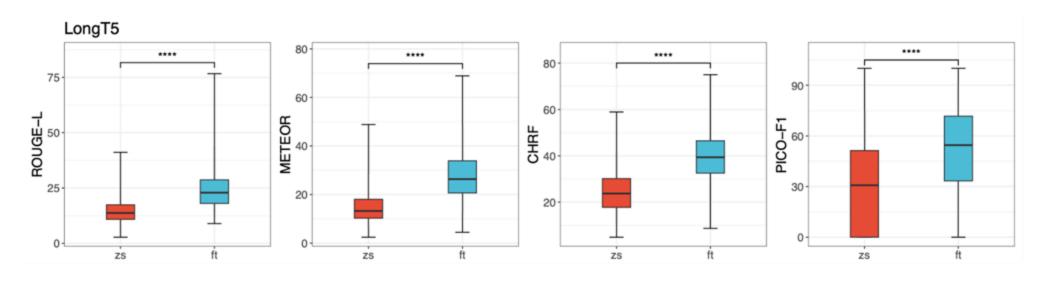
### Automatic evaluation metrics

Reference-based metrics (similarity between generated- and reference-metrics):

- Commonly used metrics
  - ROUGE-L
  - METEOR
  - CHRF
- PICO based precision / recall / F-1
  - Participants, Interventions, Comparison, and Outcomes
  - Fine-tuned a BERT-based model to extract PICO mentions



### **Automatic evaluations**



### **Model variants:**

-zero shot (zs)

-fine tuned (ft)

ROUGE-L: 14.72 → 24.61

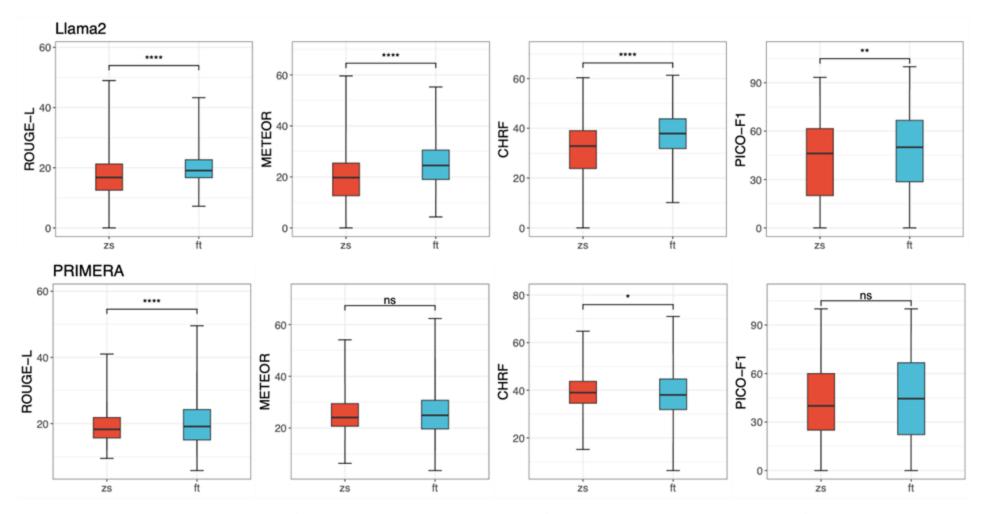
METEOR: 15.06 → 28.27

CHRF: 22.99 → 38.81

PICO-F1: 36.05 → 51.43



### **Automatic evaluations**



### **Model variants:**

-zero shot (zs)

-fine tuned (ft)

The degree of improvement varies by model architecture.



# Auto summary examples

doi: 10.1002/14651858.CD012458

### **Reference Summary:**

The review found no evidence from randomized trials to either support or refute the use of patient isolation measures (single-room isolation or cohorting) in neonates with HAIs. Risks secondary to infection control measures need to be balanced against the benefits of decreasing horizontal transmission in the neonatal unit for optimal neonatal outcomes. There is an urgent need to research the effectiveness of patient isolation measures for preventing the transmission of HAIs in neonatal units. Well-designed trials randomizing clusters of units or hospitals to a type of patient isolation method intervention are warranted.

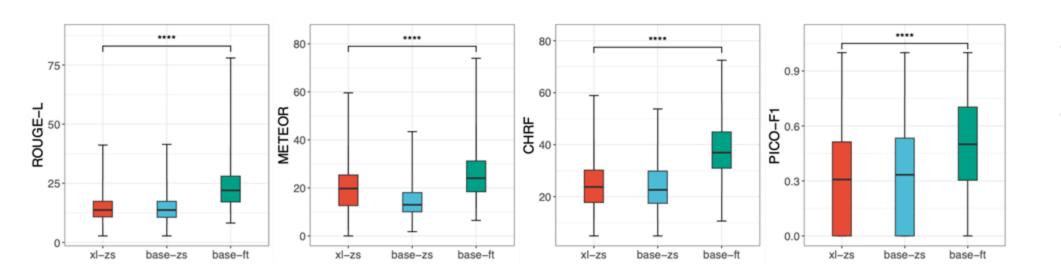
### LongT5:

(ROUGE-L: 30.88, METEOR: 23.48)

We found no evidence to support or refute the use of single-room isolation or cohorting for preventing transmission of HAIs or colonization with HAI causing pathogens in newborn infants less than six months of age admitted to the neonatal intensive care unit (NICU).



# Fine-tuned smaller models vs. larger baselines



Model	Total # of Params
PRIMERA	448.40M
LongT5-base	251.13M
LongT5-xl	2.85B
Llama-2	69.0B

### T5 variants:

- -xl (extra large)
- -base
- -base after fine tuning

Fine-tuned smaller models may outperform larger models.



# LongT5 vs. GPT-3.5-turbo

### LongT5:

(ROUGE-L: 30.88, METEOR: 23.48)

We found no evidence to support or refute the use of single-room isolation or cohorting for preventing transmission of HAIs or colonization with HAI causing pathogens in newborn infants less than six months of age admitted to the neonatal intensive care unit (NICU).

### **GPT-3.5-turbo:**

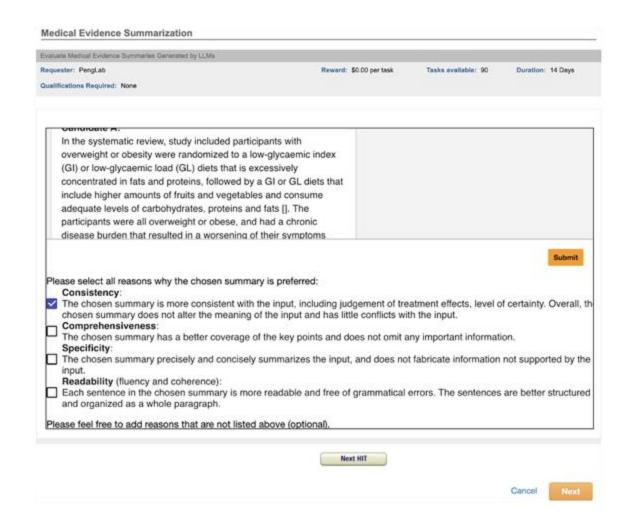
(ROUGE-L: 27.59, METEOR: 22.66)

This systematic review aimed to evaluate the impact of single-room isolation or cohorting, or both, on the transmission of healthcare-associated infections (HAIs) in newborn infants admitted to the neonatal intensive care unit (NICU). However, no published or ongoing trials were found to include in the review. Therefore, no conclusions could be drawn regarding the effectiveness of these isolation measures on preventing HAIs or colonization with HAI-causing pathogens in this population. Further research is needed in this area.



### Human evaluation

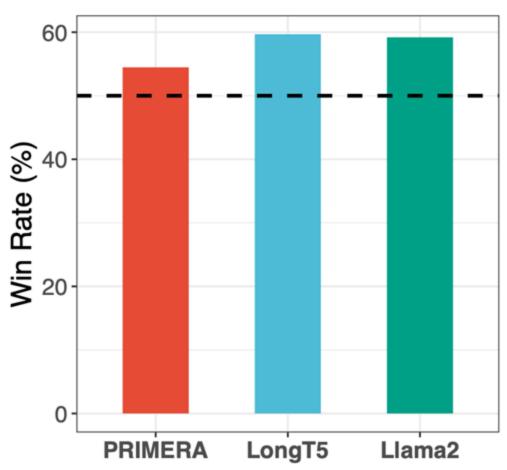
- Llama-2 vs. Fine-tuned models
- Aspects of high-quality summaries:
  - consistency
  - comprehensiveness
  - specificity
  - readability
- 7 clinical experts as judges, each received ~35 pairs of summaries
- Win-rate: % of model generated summaries better than the baseline





### Human evaluation – win-rate

### Human evaluation



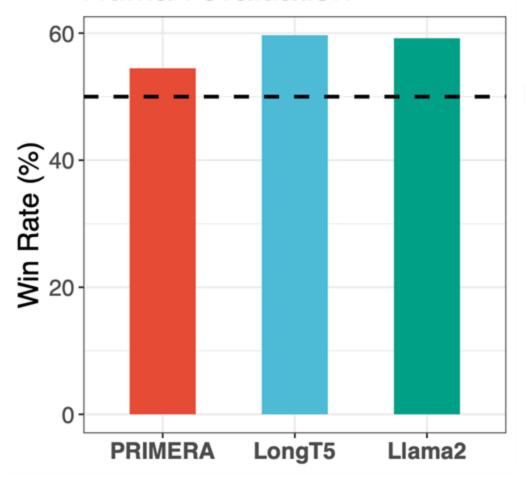
Baseline win-rate: 50%

Llama-2 70B was improved after fine-tuning.



### Human evaluation – win-rate

### Human evaluation



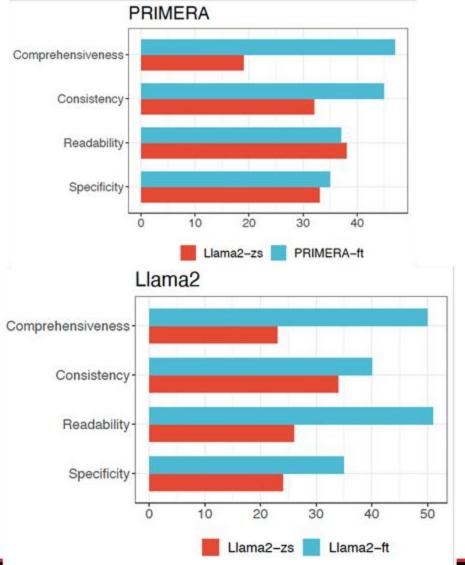
Baseline win-rate: 50%

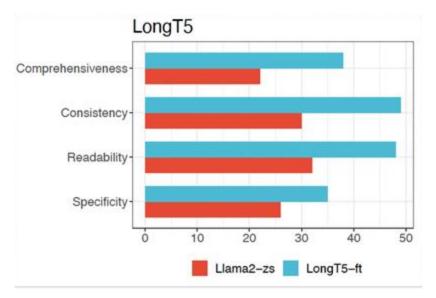
Smaller LLMs outperform Llama-2 70B after fine-tuning.

Model	Total # of Params
PRIMERA	448.40M
LongT5-base	251.13M
LongT5-xl	2.85B
Llama-2	69.0B



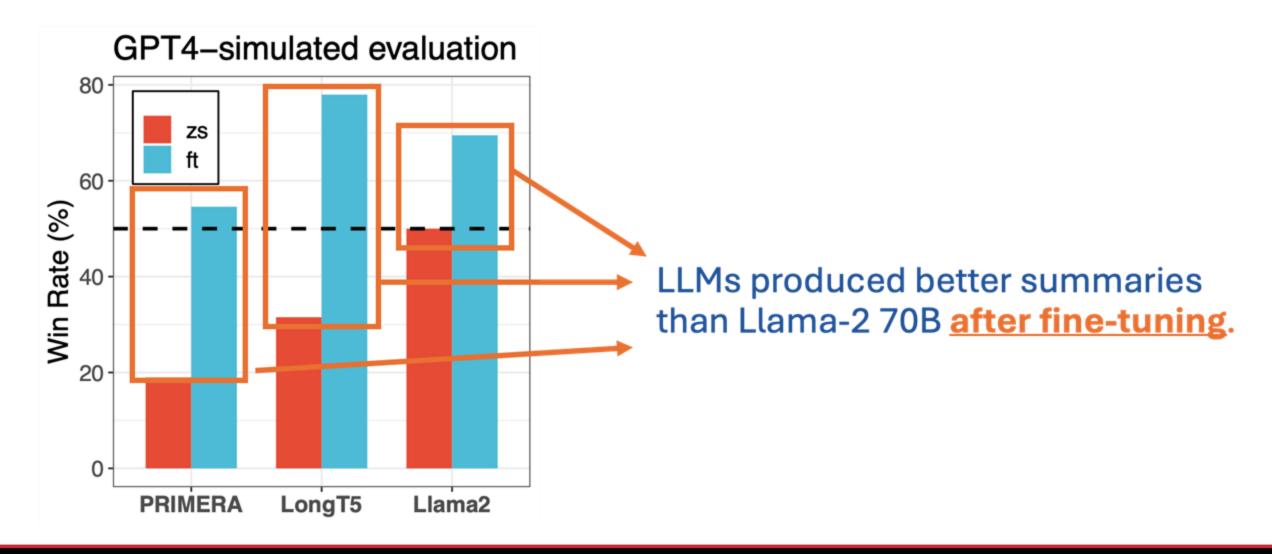
# # summaries with better qualities





Fine-tuned LLMs produced <u>more</u> comprehensive, consistent, concise summaries than Llama-2 70B.

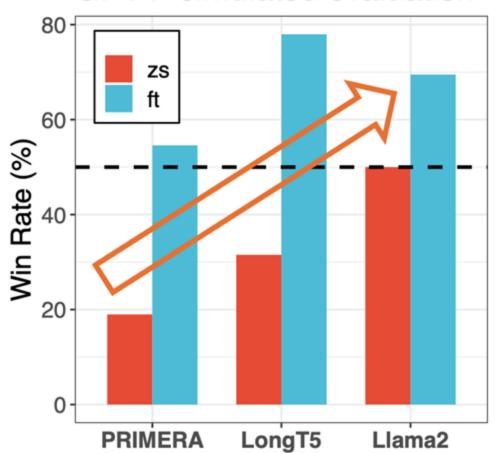
# **GPT4-simulated Evaluation**





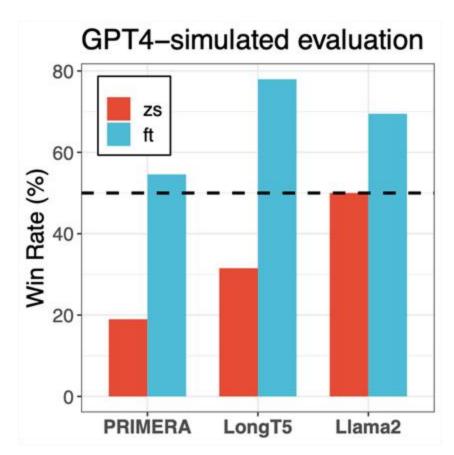
## **GPT4-simulated Evaluation**

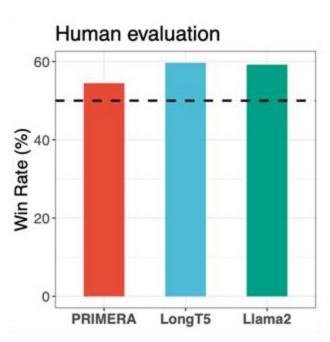
**GPT4**–simulated evaluation



Model	Total # of Params
PRIMERA	448.40M
LongT5-base	251.13M
LongT5-xl	2.85B
Llama-2	69.0B

## **GPT4-simulated Evaluation**





68% simulated evaluation results are consistent with expert judgements.



## Summary

- Usage of proprietary models raises risk for maintainability.
- Open-source LLMs produce more comprehensive and accurate summaries after fine-tuning.
- Fine-tuned open-source models achieve similar performance as SOTA closed-source ones.



# Limitations, Biases, and Open Challenges

- **Hallucination:** LLM generates plausible-sounding content that is factually incorrect or unverifiable
- ★ Citation Fabrication: LLM produces references that look real (e.g., correct author, journal, format) but don't actually exist
- **★** These risks are amplified in biomedicine, where fabricated results or citations can mislead clinical decisions or systematic reviews

## Hallucination



Type of Hallucination	Meaning	Example Input	Output			
Input-conflicting	LLM generates content that deviates from the user's input	A product manager provides the LLM with a customer interview transcript and asks for a summary of needs.	The LLM adds user needs that were <b>not mentioned</b> in the transcript.			
Context-conflicting	LLM generates content that conflicts with its own previous output	A PM asks for top user complaints, then asks about retention issues.	First answer: "slow loading, bad UI." Second answer: "crashes, no notifications" (not consistent).			
Fact-conflicting	LLM generates content that contradicts known or verifiable facts	A CEO asks: "What is the market share of competitor Z in autonomous vehicles?"	The LLM replies with <b>fabricated numbers</b> , even though the data is not publicly available.			

## Citation Fabrication

X Citation Output	What's Wrong				
Smith et al., J Clin Oncol, 2021	No such paper exists in PubMed				
Liu et al., <i>Natur</i> e, 2020	Misattributed content, wrong authors				
Zhang et al., <i>NEJM</i> , 2022	Real journal, fake volume/page				



### **Incorrect Citations**

Lead to misinformation and credibility loss





**Correct Citations** 

Ensure accuracy and credibility



# Knowledge Retrieval



## Why knowledge retrieval?

- LLMs demonstrate strong overall performance, but errors and hallucinations remain prevalent when applied to biomedicine and healthcare applications
- Rapid updated domain-specific knowledge
  - Updated treatment guidelines
  - New drug approvals
- Accurate answers alone are insufficient for high-stakes medical applications
  - Clinicians require verifiable, traceable evidence

### Benchmarking large language models for genomic knowledge with GeneTuring

Xinyi Shang<sup>1</sup>, Xu Liao<sup>1</sup>, Zhicheng Ji<sup>3</sup>, and Wenpin Hou<sup>1, 2, †</sup>

### 16 genomics tasks, 1,600 curated questions

 Manually evaluated 48,000 answers from ten LLM variants

#### ABSTRACT

Large language models (LLMs) show promise in biomedical research, but their effectiveness for genomic inquiry remains unclear. We developed GeneTuring, a benchmark consisting of 16 genomics tasks with 1,600 curated questions, and manually evaluated 48,000 answers from ten LLM configurations, including GPT-4o (via API, ChatGPT with web access, and a custom GPT setup), GPT-3.5, Claude 3.5, Gemini Advanced, GeneGPT (both slim and full), BioGPT, and BioMedLM. A custom GPT-4o configuration integrated with NCBI APIs, developed in this study as SeqSnap, achieved the best overall performance. GPT-4o with web access and GeneGPT demonstrated complementary strengths. Our findings highlight both the promise and current limitations of LLMs in genomics, and emphasize the value of combining LLMs with domain-specific tools for robust genomic intelligence. GeneTuring offers a key resource for benchmarking and improving LLMs in biomedical research.

 Prevalent hallucinations are observed for LLMs without retrieval

Hou, W., Shang, X. and Ji, Z., 2025. Benchmarking large language models for genomic knowledge with Gene Turing bio Rxiv pp 2023-03

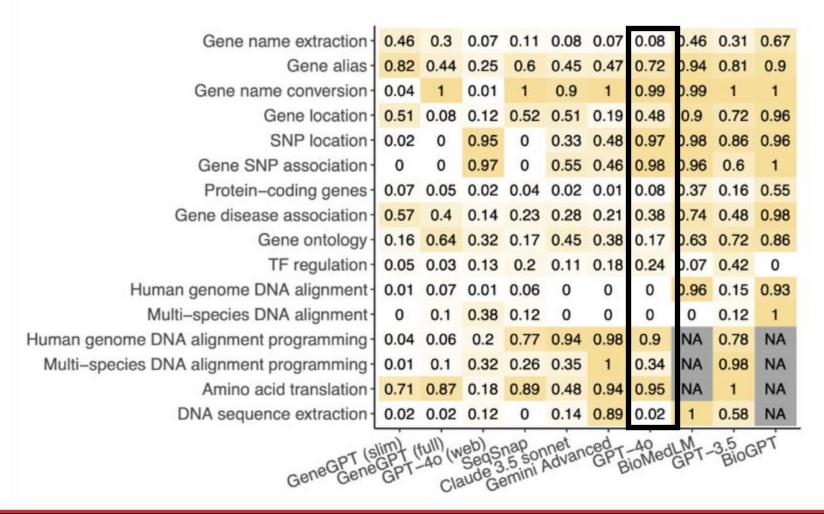
Department of Biostatistics, The Mailman School of Public Health, Columbia University, New York City, NY, USA

<sup>&</sup>lt;sup>2</sup>Data Science Institute, Columbia University, New York City, NY, USA

<sup>&</sup>lt;sup>3</sup>Department of Biostatistics and Bioinformatics, Duke University School of Medicine, Durham, NC, USA.

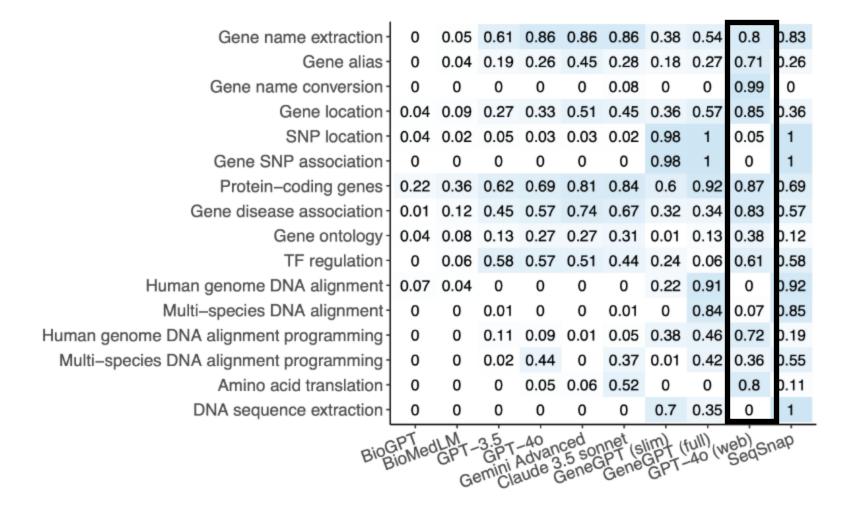
<sup>&</sup>lt;sup>†</sup>Corresponding author. E-mail: wh2526@cumc.columbia.edu

### Without retrieval vs retrieval



- Hallucination scores (lower is better)
- GPT-40 without retrieval produced errors in 99% of gene-name conversion cases

### Without retrieval vs retrieval



- Accuracy scores (higher is better)
- For the same gene name conversion task, GPT-40 with retrieval achieved 99% accuracy



## Limitations of Learning New and Updated Medical Knowledge with Commercial Fine-Tuning Large Language Models

Authors: Eric Wu, Ph.D. <sup>(a)</sup>, Kevin Wu, Ph.D. <sup>(a)</sup>, and James Zou, Ph.D. <sup>(a)</sup> Author Info & Affiliations

Published July 15, 2025 | NEJM AI 2025;2(8) | DOI: 10.1056/Alcs2401155 | VOL. 2 NO. 8 | Copyright © 2025









### Abstract

Large language models (LLMs) used in health care need to integrate new and updated medical knowledge to produce relevant and accurate responses. For example, medical guidelines and drug information are frequently updated or replaced as new evidence emerges. To address this need, companies like OpenAI, Google, and Meta allow users to fine-tune their proprietary models through commercial application programming interfaces. However, it is unclear how effectively LLMs can leverage updated medical information through these commercial fine-tuning services. In this case study, we systematically fine-tuned six frontier LLMs — including GPT-40, Gemini 1.5 Pro, and Llama 3.1 — using a novel dataset of new and updated medical knowledge. We found that these models exhibit limited generalization on new U.S. Food and Drug Administration drug approvals, patient records, and updated medical guidelines. Among all tested models, GPT-40 mini showed the strongest performance. These findings underscore the current limitations of fine-tuning frontier models for up-to-date medical use cases.



 LLMs struggled with new and updated medical knowledge compared to standard benchmarks

Wu, E., Wu, K. and Zou, J., 2025. Limitations of learning new and updated medical knowledge with commercial fine-tuning large language models. *NEJM AI*, 2(8), p.Alcs2401155.

### How do we retrieve relevant evidence?

- How do humans retrieve relevant evidence?
  - Do creative web search
    - Experiment with several searches
    - Put yourself in an author's shoes; what phrases might they have used?
    - Specifically search Google Scholar, etc.
- Track down related work (once you have a relevant paper)
  - Follow the bibliography to earlier papers
  - See who else has cited the work
- Read each paper to find related work



## How do we retrieve relevant evidence?

- How do humans retrieve relevant evidence?
  - Do creative web search
    - Experiment with several searches
    - Put yourself in an author's shoes; what phrases might they have used?
    - Specifically search Google Scholar, etc.
  - Track down related work (once you have a relevant paper)
    - Follow the bibliography to earlier papers
    - See who else has cited the work
  - Read each paper to find related work

- How do machines retrieve relevant evidence?
  - Construct knowledge graphs for each paper
    - Search for keywords based on the constructed knowledge graph
  - Use papers' references and citation networks to identify related papers
  - Search for semantically relevant sentences based on embedding similarity and overlapping knowledge graphs



## Different types of knowledge retrieval

Corpus-level Information Retrieval

Literature Search Engines (Google Scholar, PubMed)

**KG-Augmented** Retrieval (PubTator 3, ESRA)

Sentence/Conceptlevel Information Retrieval

> Sentence Retrieval (LitSense, EvidenceMiner, SPIKE)

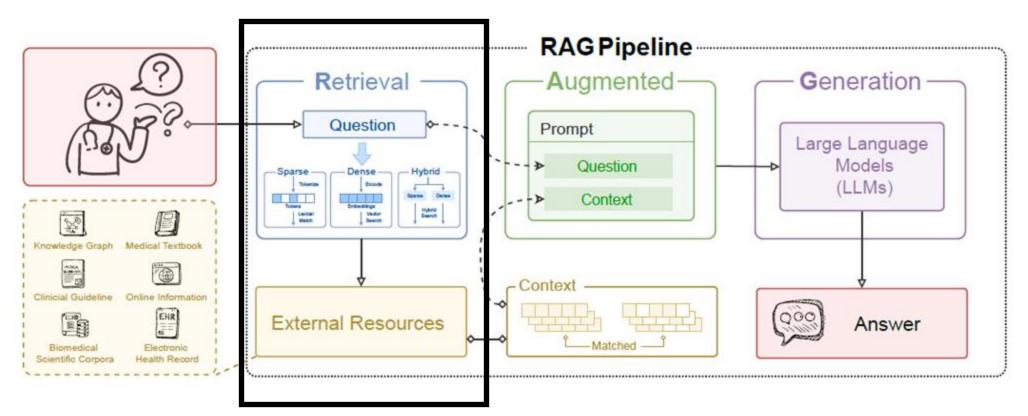
Dataset Recommendation (Datafinder)

Multimodal Information Retrieval

> Multi-modal Molecule Text Retrieval (Text2mol, MoleculeSTM)

Multi-modal Image Text Retrieval (Cross-modal RAG)

## Standard retrieval pipeline

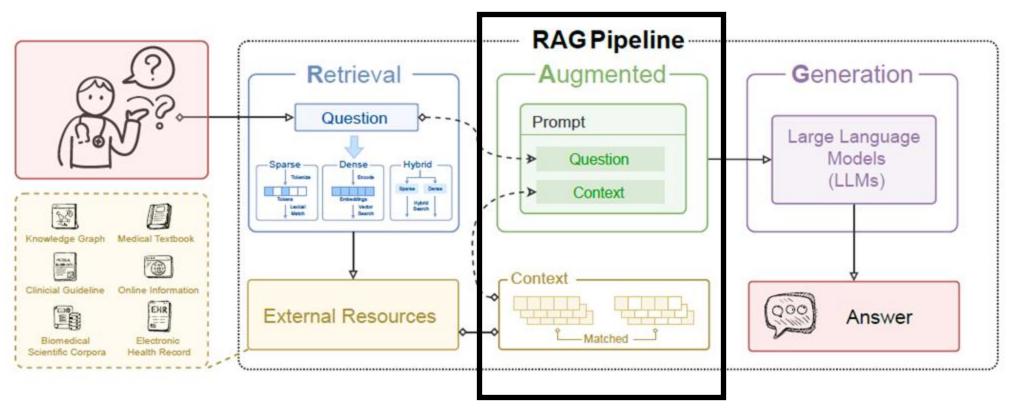


• Step 1: Retrieval

 Searches and provides information relevant to the query

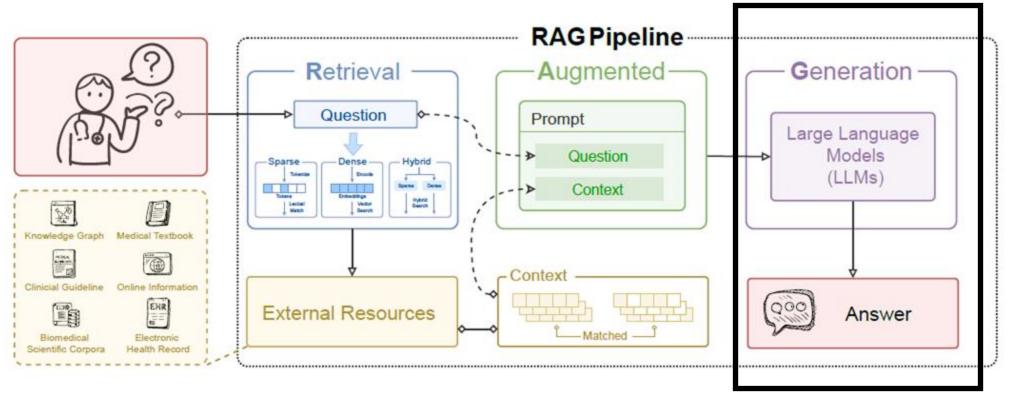
Yang, R., Wong, M.Y.H., Li, H., Li, X., Zhu, W., Liao, J., Yu, K., Liew, J.C.K., Xuan, W., Chen, Y. Ke, Y., ... and Liu, N., 2025. Retrieval-Augmented Generation in Medicine: A Scoping Review of Technical Implementations, Clinical Applications, and Ethical Considerations. arXiv preprint arXiv:2511.05901.

## Standard retrieval pipeline



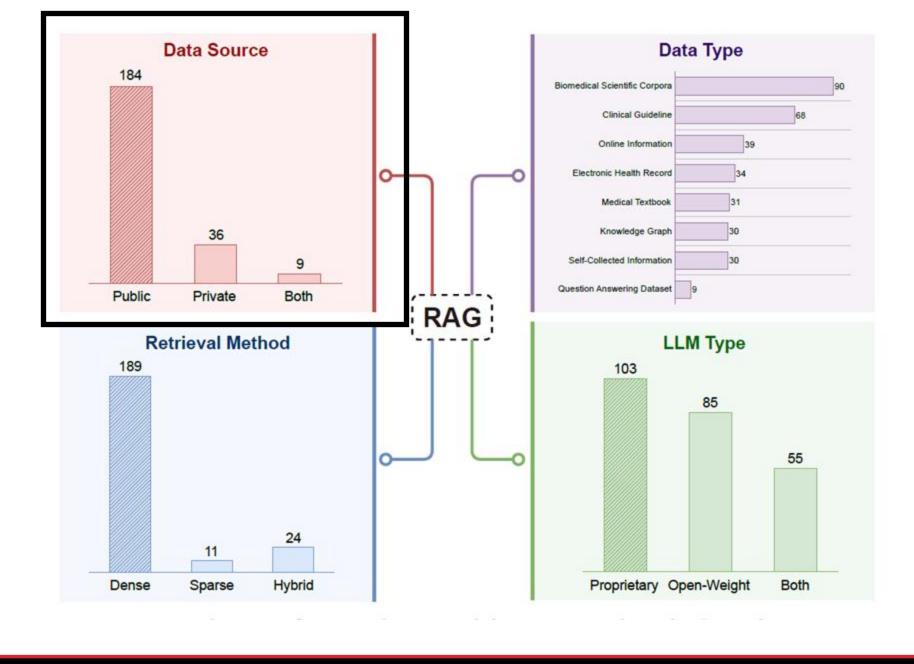
- Step 2: Augmentati on
- Prompt + retrieved context to LLMs

## Standard retrieval pipeline



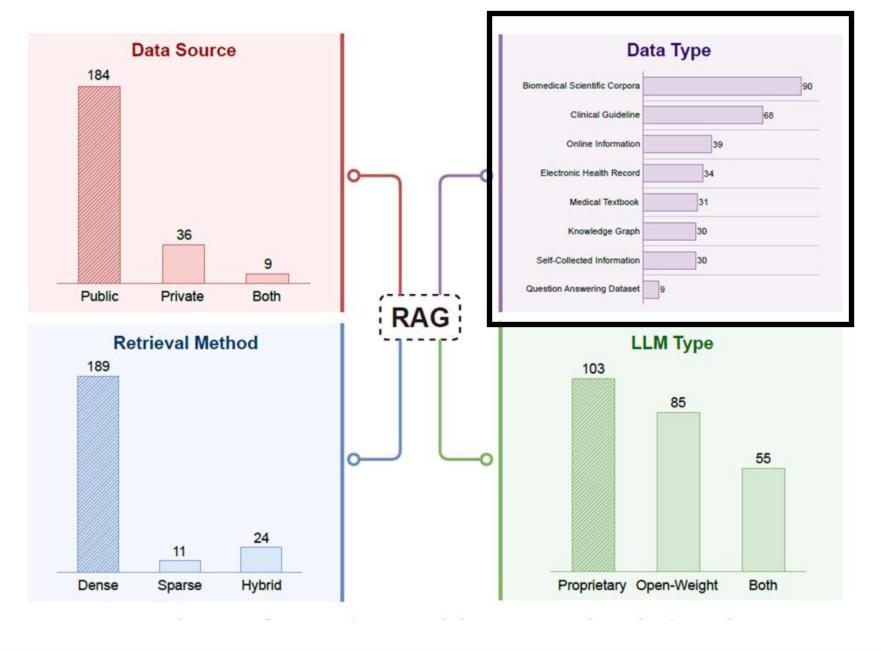
• Step 3: Generation

LLMs generate the final output



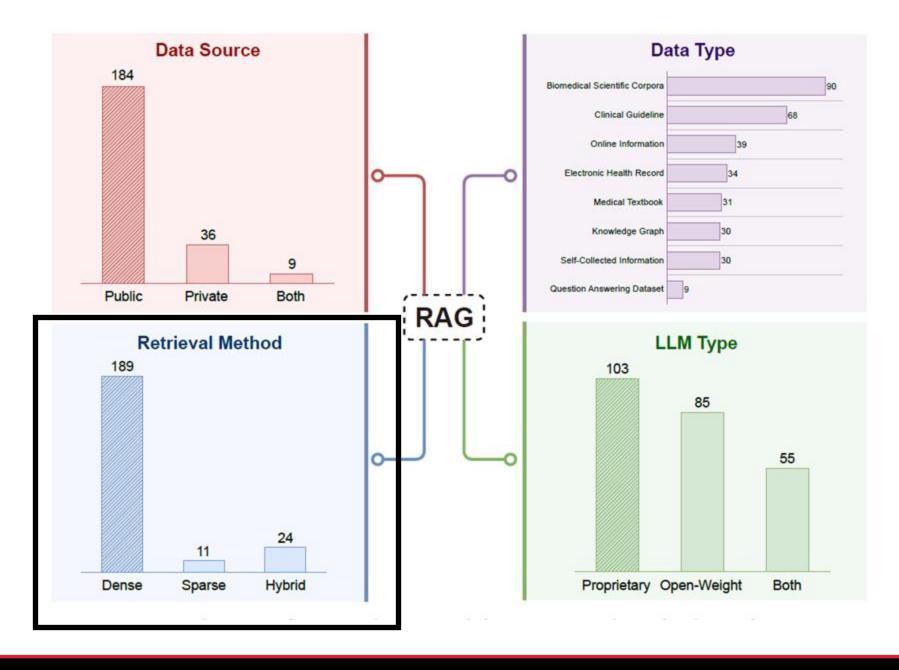
### Data source

- 80% used public data
- 15% used private
- 5% used both



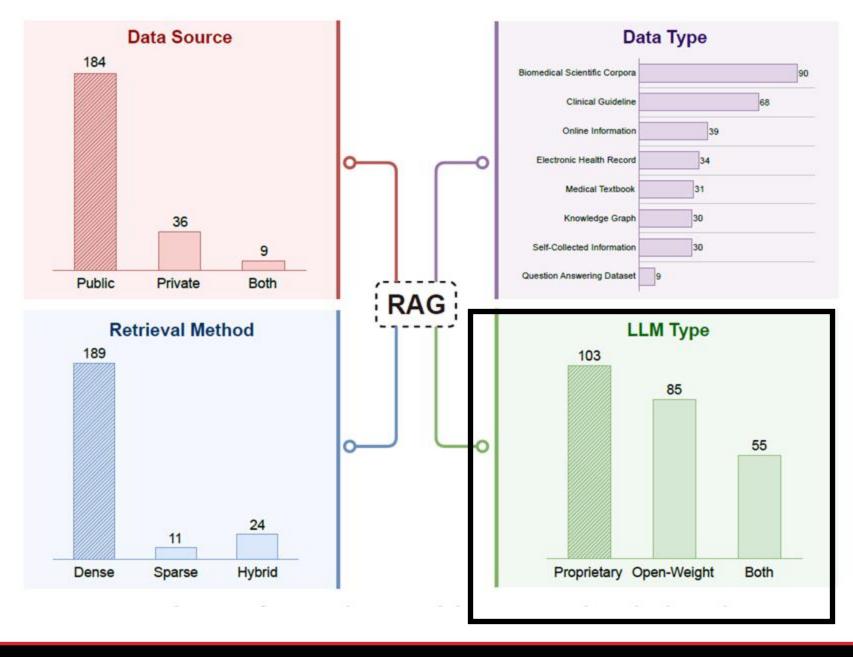
### Data type

- Biomedical literature
- Clinical guidelines
- Clinical notes
- Textbooks
- Knowledge graphs
- Others



### Retrieval methods

 Mostly dense (embedding-based)



### LLM backbones

 Both closedsourced and openweighted are commonly used

Publication Date	Framework Retriever(s)		LLM(s)	Knowledge Sources	Application(s)		
Jun 2023	ChatDoctor86	Lexical matching	cal matching Llama Wikipedia, MedlinePlus		QA		
Jul 2023	accGPT <sup>87</sup>	Text-embedding-ada- 002 <sup>88</sup>	GPT-3.5	Clinical guidelines	Imaging recommendation		
Aug 2023	FraCChat <sup>89</sup>	Text-embedding-ada-002	GPT-3.5, GPT-4	AO/OTA Fracture and Dis- location Classification Com- pendium (2018)	AO code identification		
Dec 2023	Clinfo.ai90	BM25, Entrez API91	GPT-3.5, GPT-4	PubMed	QA, Summarization		
Jan 2024	Almanac <sup>29</sup>	Text-embedding-ada-002	GPT-4	PubMed, Clinical guidelines, etc.			
Jan 2024	RALL <sup>92</sup>	DPR <sup>63</sup> , etc.	GPT-4, Llama-2	Wikipedia, UMLS	Lay language generation		
May 2024	Ferber et al.93	Text-embedding-ada-002	GPT-4	Clinical guidelines	QA		
Jun 2024	RECTIFIER <sup>30</sup>	Text-embedding-ada-002	GPT-3.5, GPT-4	Clinical guidelines	Clinical trial screening		
Jul 2024	Self-BioRAG <sup>76</sup>	MedCPT	Llama-2	PubMed, PMC, Clinical guide- lines, Textbooks	QA		
Jul 2024	DocOA <sup>94</sup>	Not specified	GPT-3.5	PubMed, Clinical guidelines	Osteoarthritis management		
Aug 2024	MedRAG <sup>39</sup>	BM25, MedCPT, Contriever <sup>95</sup> , SPECTER <sup>96</sup>	GPT-3.5, GPT-4, Llama-2, Mixtral- 8x7B, etc.	PubMed, StatPearls, Textbooks, Wikipedia	QA		
Aug 2024	Alkhalaf et al.97	BM25	Llama-2	EHR	Summarization		
Sep 2024	ChatZOC50	BM25	Baichuan	Clinical guidelines, FAQs	Opthalmology QA		
Sep 2024	Benfenati et al. <sup>98</sup>	BM25, GTE <sup>99</sup>	GPT-3.5, Mistral	Nutrigenetic polymorphism dataset, PubMed-derived nutrition-gene studies	Nutrigenetics QA		
Oct 2024	Bora and Cuayáhuitl 100	GTR-T5-Large	Llama-2, Mistral	Textbooks, Journals	QA		
Dec 2024	EyeGPT <sup>101</sup>	all-MiniLM-L6-v2 <sup>102</sup>	Llama-2	Textbooks, Custom Database	Ophthalmology QA		

 A closer look at published studies using LLM-RAG in medicine up to August 2025

Kim, H., Sohn, J., Gilson, A., Cochran-Caggiano, N., Applebaum, S., Jin, H., Park, S., Park, Y., Park, J., Choi, S. Contreras, B.A.H., ..., Chen, Q., 2025. Rethinking Retrieval-Augmented Generation for Medicine: A Large-Scale, Systematic Expert Evaluation and Practical Insights. arXiv preprint arXiv:2511.06738.



Jan 2025	i-MedRAG <sup>103</sup>	MedCPT	GPT-4	StatPearls, Textbooks	QA	
Jan 2025	Azimi et al. 100	Titan Text Embeddings V2 <sup>104</sup>	GPT-40, Claude 3.5 Sonnet	e Registered Dietitian (RD) exam questions across four nutrition domains (Academy of Nutrition and Dietetics guidelines and re- erences)	n n	
Feb 2025	Bailicai <sup>105</sup>	MedCPT	Qwen <sup>106</sup> , Llama-2	PubMed, StatPearls, Textbook Wikipedia	s, QA	
Mar 2025	Woo et al. 107	t al. 107 Internal retriever GPT-3.5, GPT-4, Clinical guidelines Claude 3, Llama-3, Mistral-8×7B			Clinical decision support	
Apr 2025	RAG squared <sup>40</sup>	MedCPT	GPT-40, Llama-3 etc.	<ol> <li>PubMed, PMC, Clinical guide lines, Textbooks</li> </ol>	e- QA	
May 2025	BriefContext <sup>108</sup>	BM25, MedCPT	GPT-40, Llama	- PubMed, Textbooks	Summarization	
May 2025	MKRAG <sup>109</sup>	Contriever, etc.	ontriever, etc. Vicuna-7B <sup>110</sup> Disease Database		Knowledge graph QA	
May 2025	Yang et al. <sup>111</sup>	Elasticsearch, Col- BERTv2 <sup>112</sup>	GPT-4, Claude 3	Internal, expert-reviewed hosp tal documents	i- Clinical decision support	
Jul 2025	MedGraphRAG <sup>113</sup>	Graph-based retrieval	GPT-4, Gemin 1.0 Pro, Llama-2 Llama-3	i MIMIC-IV, FakeHealth, Pub c, Health, UMLS Graph	o- QA	
Jul 2025	RadioRAG <sup>114</sup>	Text-embedding-ada-002	GPT-40 mini	50000000000000000000000000000000000000	Case retrieval, Report genera-	
Jul 2025	Wada et al. <sup>115</sup>	Text-embedding-3- large <sup>88</sup> , Lexical matching		Clinical guidelines (ACR Man- ual, ESUR Guidelines, etc.)	Safety consultation	
Jul 2025	MedOmniKB <sup>116</sup>	MedCPT, Graph-based re- trieval	GPT-4, Gemini 1.5	PubMed, Clinical guidelines, Textbooks, Wikipedia, UMLS, DrugBank	QA	
Aug 2025	MedCoT-RAG <sup>117</sup>	MedCPT		PubMed, StatPearls, Textbooks, Wikipedia	QA	

### LLM-RAG implementation examples: API-based

Bioinformatics, 2024, 40(2), btae075 https://doi.org/10.1093/bioinformatics/btae075 Advance Access Publication Date: 10 February 2024

Original Paper



### Data and text mining

# GeneGPT: augmenting large language models with domain tools for improved access to biomedical information

Qiao Jin 10 1, Yifan Yang 1, Qingyu Chen 10 1, Zhiyong Lu 10 1,\*

<sup>1</sup>National Center for Biotechnology Information, National Library of Medicine, National Institutes of Health, Bethesda, MD 20894, United States

\*Corresponding author. National Center for Biotechnology Information, National Library of Medicine, National Institutes of Health, Bethesda, MD 20894, United States, Tel: 301-594-7089, E-mail: zhiyong.lu@nih.gov (Z.L.)

Associate Editor: Jonathan Wren

#### Abstract

Motivation: While large language models (LLMs) have been successfully applied to various tasks, they still face challenges with hallucinations. Augmenting LLMs with domain-specific tools such as database utilities can facilitate easier and more precise access to specialized knowledge. In this article, we present GeneGPT, a novel method for teaching LLMs to use the Web APIs of the National Center for Biotechnology Information (NCBI) for answering genomics questions. Specifically, we prompt Codex to solve the GeneTuring tests with NCBI Web APIs by incontext learning and an augmented decoding algorithm that can detect and execute API calls.

Results: Experimental results show that GeneGPT achieves state-of-the-art performance on eight tasks in the GeneTuring benchmark with an average score of 0.83, largely surpassing retrieval-augmented LLMs such as the new Bing (0.44), biomedical LLMs such as BioMedLM (0.08) and BioGPT (0.04), as well as GPT-3 (0.16) and ChatGPT (0.12). Our further analyses suggest that: First, API demonstrations have good cross-task generalizability and are more useful than documentations for in-context learning; second, GeneGPT can generalize to longer chains of API calls and answer multi-hop questions in GeneHop, a novel dataset introduced in this work; finally, different types of errors are enriched in different tasks, providing valuable insights for future improvements.

Availability and implementation: The GeneGPT code and data are publicly available at https://github.com/ncbi/GeneGPT.

Jin, Q., Yang, Y., Chen, Q. and Lu, Z., 2024. Genegpt: Augmenting large language models with domain tools for improved access to biomedical information. *Bioinformatics*, 40(2), p.btae075.

#### GeneGPT prompt Hello. Your task is to use NCBI APIs to answer genomic questions. There are two types of Web APIs you can use: Eutils and BLAST. You can call Eutils by: "[https://eutils.ncbi.nlm.nih.gov/entrez/eutils/{esearch|efetch|esum mary}.fcgi?db={gene|snp|omim}&retmax={}&{term|id}={}]". {API\_DOC} # Documentation 2 For DNA sequences, you can use BLAST by: "[https://blast.ncbi.nlm.nih.gov/blast/Blast.cgi?CMD={Put|Get}&PROGR AM=blastn&MEGABLAST=on&DATABASE=nt&FORMAT\_TYPE={XML|Text}&QUERY={seq uence}&HITLIST\_SIZE={max\_hit\_size}]". {API\_DOC} # Demonstration 1 Question: What is the official gene symbol of LMP10? [https://eutils.ncbi.nlm.nih.gov/entrez/eutils/esearch.fcgi?db=gene& retmax=3&retmode=json&term=LMP10]->[{API\_CALL}] [https://eutils.ncbi.nlm.nih.gov/entrez/eutils/efetch.fcgi?db=gene&r etmax=3&retmode=json&id={ids}]->[{API\_CALL}] Answer: PSMB10 # Demonstration 2 Question: Which gene is SNP rs1217074595 associated with? [https://eutils.ncbi.nlm.nih.gov/entrez/eutils/esummary.fcgi?db=snp& retmax=3&retmode=ison&id=1217074595]->[{API\_CALL}] Answer: LINC01270 # Demonstration 3 Ouestion: What are genes related to Meesmann corneal dystrophy? [https://eutils.ncbi.nlm.nih.gov/entrez/eutils/esearch.fcgi?db=omim& retmax=15&term=Meesmann+corneal+dystrophy]->[{API\_CALL}] [https://eutils.ncbi.nlm.nih.gov/entrez/eutils/esummary.fcgi?db=omim &retmax=15&retmode=json&id={ids}]->[{API\_CALL}] Answer: KRT12, KRT3 # Demonstration 4 Question: Align the DNA sequence to the human genome:ATTCT[...]TCTCT [https://blast.ncbi.nlm.nih.gov/blast/Blast.cgi?CMD=Put&PROGRAM=blas tn&MEGABLAST=on&DATABASE=nt&FORMAT\_TYPE=XML&QUERY=ATTCT[...]TCTCT&HI TLIST\_SIZE=5]->[{API\_CALL}] [https://blast.ncbi.nlm.nih.gov/blast/Blast.cgi?CMD=Get&FORMAT\_TYPE= Text&RID={RID}]->[{API\_CALL}]

#### GeneTuring Question

Which chromosome does SNP rs1430464868 locate on human

#### GeneGPT Output

[https://eutils.ncbi.nlm.nih.gov/entrez/eutils/esummary.fc gi?db=snp&retmax=3&retmode=json&id=1430464868]->[{"header" :{"type":"esummary", "version":"0.3"}, "result":{"uids":["14 30464868"], "1430464868": {"uid": "1430464868", "snp\_id": 14304 64868, "allele\_origin": "", "global\_mafs": [{"study": "GnomAD", "freq": "C=0.000007/1"}, {"study": "ALFA", "freq": "C=0./0"}]," global\_population":"", "global\_samplesize":"", "suspected":" ","clinical\_significance":"","genes":[],"acc":"NC\_000013.1
1","chr":"13","handle": [\*TRUNCATED\*] merged\_sort":"0"}}}] Answer: chr13

#### GeneHop Question

What is the function of the gene associated with SNP rs1241371358? Let's decompose the question to subquestions and solve them step by step.

#### GeneGPT Output

Sub-question 1: Which gene is SNP rs1241371358 associated [https://eutils.ncbi.nlm.nih.gov/entrez/eutils/esummary.fc gi?db=snp&retmax=10&retmode=json&id=1241371358]->[{"heade r":{"type":"esummary", "version":"0.3"}}, "result":{"uids":[ "1241371358"] [\*TRUNCATED\*] "merged\_sort":"0"}}] Answer: LRRC23 Sub-question 2: What is the function of LRRC23?

[https://eutils.ncbi.nlm.nih.gov/entrez/eutils/esummary.fc gi?db=gene&retmax=10&retmode=json&id=10233]->[{"header": {"type": "esummary" [\*TRUNCATED\*] "merged\_sort": "0"}}] Answer: Predicted to be active in cytosol.

### Instructions

 Ask LLMs to use APIs for retrieval

Answer: chr15:91950805-91950932

#### GeneGPT prompt Hello. Your task is to use NCBI APIs to answer genomic questions. There are two types of Web APIs you can use: Eutils and BLAST. You can call Eutils by: "[https://eutils.ncbi.nlm.nih.gov/entrez/eutils/{esearch|efetch|esu mary}.fcgi?db={gene|snp|omim}&retmax={}&{term|id}={}]". {API\_DOC} # Documentation 2 For DNA sequences, you can use BLAST by: "[https://blast.ncbi.nlm.nih.gov/blast/Blast.cgi?CMD={Put|Get}&PROG AM=blastn&MEGABLAST=on&DATABASE=nt&FORMAT\_TYPE={XML|Text}&QUERY={se uence}&HITLIST\_SIZE={max\_hit\_size}]". {API\_DOC} Question: What is the official gene symbol of LMP10? [https://eutils.ncbi.nlm.nih.gov/entrez/eutils/esearch.fcgi?db=gene8 retmax=3&retmode=json&term=LMP10]->[{API\_CALL}] [https://eutils.ncbi.nlm.nih.gov/entrez/eutils/efetch.fcgi?db=gene&r etmax=3&retmode=json&id={ids}]->[{API\_CALL}] Answer: PSMB10 # Demonstration 2 Question: Which gene is SNP rs1217074595 associated with? [https://eutils.ncbi.nlm.nih.gov/entrez/eutils/esummary.fcgi?db=snp& retmax=3&retmode=ison&id=1217074595]->[{API\_CALL}] Answer: LINC01270 # Demonstration 3 Ouestion: What are genes related to Meesmann corneal dystrophy? Thttps://eutils.ncbi.nlm.nih.gov/entrez/eutils/esearch.fcgi?db=omim& retmax=15&term=Meesmann+corneal+dystrophy]->[{API\_CALL}] [https://eutils.ncbi.nlm.nih.gov/entrez/eutils/esummary.fcgi?db=omim &retmax=15&retmode=json&id={ids}]->[{API\_CALL}] Answer: KRT12, KRT3 # Demonstration 4 Question: Align the DNA sequence to the human genome:ATTCT[...]TCTCT

[https://blast.ncbi.nlm.nih.gov/blast/Blast.cgi?CMD=Put&PROGRAM=blas tn&MEGABLAST=on&DATABASE=nt&FORMAT\_TYPE=XML&QUERY=ATTCT[...]TCTCT&HI

[https://blast.ncbi.nlm.nih.gov/blast/Blast.cqi?CMD=Get&FORMAT\_TYPE=

#### GeneTuring Question

Which chromosome does SNP rs1430464868 locate on human genome?

#### GeneGPT Output

[https://eutils.ncbi.nlm.nih.gov/entrez/eutils/esummary.fc gi?db=snp&retmax=3&retmode=json&id=1430464868]->[{"header" :{"type":"esummary","version":"0.3"},"result":{"uids":["14 30464868"], "1430464868": {"uid": "1430464868", "snp\_id": 14304 64868, "allele\_origin": "", "global\_mafs": [{"study": "GnomAD", "freq": "C=0.000007/1"}, {"study": "ALFA", "freq": "C=0./0"}]," global\_population":"", "global\_samplesize":"", "suspected":" ","clinical\_significance":"","genes":[],"acc":"NC\_000013.1
1","chr":"13","handle": [\*TRUNCATED\*] merged\_sort":"0"}}}] Answer: chr13

#### GeneHop Question

What is the function of the gene associated with SNP rs1241371358? Let's decompose the question to subquestions and solve them step by step.

#### GeneGPT Output

Sub-question 1: Which gene is SNP rs1241371358 associated [https://eutils.ncbi.nlm.nih.gov/entrez/eutils/esummary.fc gi?db=snp&retmax=10&retmode=json&id=1241371358]->[{"heade r":{"type":"esummary", "version":"0.3"}}, "result":{"uids":[ "1241371358"] [\*TRUNCATED\*] "merged\_sort":"0"}}] Answer: LRRC23 Sub-question 2: What is the function of LRRC23? [https://eutils.ncbi.nlm.nih.gov/entrez/eutils/esummary.fc gi?db=gene&retmax=10&retmode=json&id=10233]->[{"header": {"type": "esummary" [\*TRUNCATED\*] "merged\_sort": "0"}}] Answer: Predicted to be active in cytosol.

### Documentation

 Provide descriptions of the API functionality, general syntax, and argument choices.

TLIST\_SIZE=5]->[{API\_CALL}]

Text&RID={RID}]->[{API\_CALL}] Answer: chr15:91950805-91950932

#### GeneGPT prompt Hello. Your task is to use NCBI APIs to answer genomic questions. There are two types of Web APIs you can use: Eutils and BLAST. # Documentation 1 You can call Eutils by: "[https://eutils.ncbi.nlm.nih.gov/entrez/eutils/{esearch|efetch|esum mary}.fcgi?db={gene|snp|omim}&retmax={}&{term|id}={}]". {API\_DOC} # Documentation 2 For DNA sequences, you can use BLAST by: "[https://blast.ncbi.nlm.nih.gov/blast/Blast.cgi?CMD={Put|Get}&PROGR AM=blastn&MEGABLAST=on&DATABASE=nt&FORMAT\_TYPE={XML|Text}&OUERY={sea CAPT DOCL # Demonstration 1 Question: What is the official gene symbol of LMP10? [https://eutils.ncbi.nlm.nih.gov/entrez/eutils/esearch.fcgi?db=gene retmax=3&retmode=json&term=LMP10]->[{API\_CALL}] [https://eutils.ncbi.nlm.nih.gov/entrez/eutils/efetch.fcgi?db=gene8 etmax=3&retmode=json&id={ids}]->[{API\_CALL}] Answer: PSMB10 # Demonstration 2 Question: Which gene is SNP rs1217074595 associated with? [https://eutils.ncbi.nlm.nih.gov/entrez/eutils/esummary.fcgi?db=snp retmax=3&retmode=ison&id=1217074595]->[{API\_CALL}] Answer: LINC01270 # Demonstration 3 Ouestion: What are genes related to Meesmann corneal dystrophy? [https://eutils.ncbi.nlm.nih.gov/entrez/eutils/esearch.fcgi?db=omin retmax=15&term=Meesmann+corneal+dystrophy]->[{API\_CALL}] [https://eutils.ncbi.nlm.nih.gov/entrez/eutils/esummary.fcgi?db=omi &retmax=15&retmode=json&id={ids}]->[{API\_CALL}] Answer: KRT12, KRT3 # Demonstration 4 Question: Align the DNA sequence to the human genome: ATTCT[...]TCTC [https://blast.ncbi.nlm.nih.gov/blast/Blast.cgi?CMD=Put&PROGRAM=bla tn&MEGABLAST=on&DATABASE=nt&FORMAT\_TYPE=XML&QUERY=ATTCT[...]TCTCT&H TLIST\_SIZE=5]->[{API\_CALL}] [https://blast.ncbi.nlm.nih.gov/blast/Blast.cgi?CMD=Get&FORMAT\_TYPE Text&RID={RID}]->[{API\_CALL}] war: chr15:01050805-01050032

#### GeneTuring Question

Which chromosome does SNP rs1430464868 locate on human genome?

#### GeneGPT Output

[https://eutils.ncbi.nlm.nih.gov/entrez/eutils/esummary.fc gi?db=snp&retmax=3&retmode=json&id=1430464868]->[{"header" :{"type":"esummary","version":"0.3"},"result":{"uids":["14 30464868"], "1430464868": {"uid": "1430464868", "snp\_id": 14304 64868, "allele\_origin": "", "global\_mafs": [{"study": "GnomAD", "freq": "C=0.000007/1"}, {"study": "ALFA", "freq": "C=0./0"}]," global\_population":"", "global\_samplesize":"", "suspected":" ","clinical\_significance":"","genes":[],"acc":"NC\_000013.1
1","chr":"13","handle": [\*TRUNCATED\*] merged\_sort":"0"}}}] Answer: chr13

#### GeneHop Question

What is the function of the gene associated with SNP rs1241371358? Let's decompose the question to subquestions and solve them step by step.

#### GeneGPT Output

Sub-question 1: Which gene is SNP rs1241371358 associated [https://eutils.ncbi.nlm.nih.gov/entrez/eutils/esummary.fc gi?db=snp&retmax=10&retmode=json&id=1241371358]->[{"heade r":{"type":"esummary", "version":"0.3"}}, "result":{"uids":[ "1241371358"] [\*TRUNCATED\*] "merged\_sort":"0"}}] Answer: LRRC23 Sub-question 2: What is the function of LRRC23? [https://eutils.ncbi.nlm.nih.gov/entrez/eutils/esummary.fc gi?db=gene&retmax=10&retmode=json&id=10233]->[{"header": {"type": "esummary" [\*TRUNCATED\*] "merged\_sort": "0"}}] Answer: Predicted to be active in cytosol.

### Demonstration

 Provide examples of using NCBI Web APIs to solve questions.

### Results

GeneTuring task	GPT-2 Bio	BioGPT	BioMedLM	GPT-3	ChatGPT	New Bing	GeneGPT (ours)			
							-full	-slim	-turbo	-lang
Nomenclature										
Gene alias	0.00	0.00	0.04	0.09	0.07	0.66	0.80*	0.84*	0.64*	0.76*
Gene name conversion	0.00	0.00	0.00	0.00	0.00	0.85	1.00	1.00	1.00	0.02
Average	0.00	0.00	0.02	0.05	0.04	0.76	0.90	0.92	0.82	0.39
Genomic location										
Gene SNP association	0.00	0.00	0.00	0.00	0.00	0.00	1.00*	1.00	0.96	0.90
Gene location	0.01	0.04	0.12	0.09	0.09	0.61	0.62	0.66	0.54	0.54
SNP location	0.03	0.05	0.01	0.02	0.05	0.01	$\overline{1.00}$	0.98	0.98	0.74
Average	0.01	0.03	0.04	0.04	0.05	0.21	0.87	0.88	0.82	0.73
Functional analysis										
Gene disease association	0.00	0.02	0.16	0.34	0.31	0.84	0.76*	0.66	0.63	0.39
Protein-coding genes	0.00	0.18	0.37	0.70	0.54	0.97	0.76	1.00	0.96	0.90
Average	0.00	0.10	0.27	0.52	0.43	$\overline{0.91}$	0.76	0.84	0.80	0.65
Sequence alignment										
DNA to human genome	0.02	0.07	0.03	0.00	0.00	0.00	0.44*	0.44*	0.42*	0.06*
DNA to multiple species	0.02	0.00	0.00	0.20	0.00	0.00	0.86	0.88	$\overline{0.88}$	0.54
Average	0.02	0.04	0.02	0.10	0.00	0.00	0.65	0.66	0.65	0.30
Overall average	0.00	0.04	0.08	0.16	0.12	0.44	0.80	0.83	$\overline{0.78}$	0.54

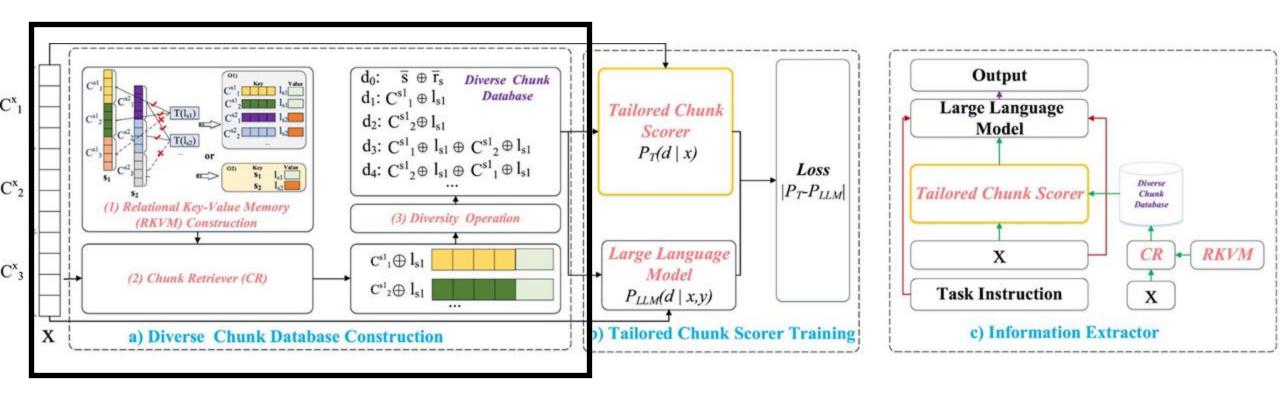
## LLM-RAG implementation examples: vector-based

## BiomedRAG: A retrieval augmented large language model for biomedicine

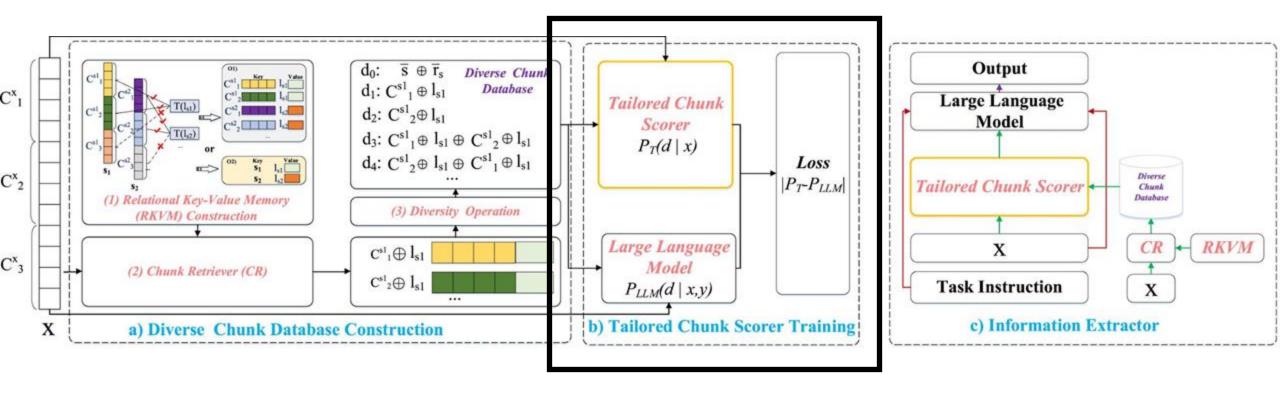
Show more V + Add to Mendeley  $\stackrel{\textstyle \leftarrow}{\triangleleft}$  Share  $\stackrel{\textstyle \rightarrow}{\flat}$  Cite https://doi.org/10.1016/j.jbi.2024.104769 7 Get rights and content 7 Under a Creative Commons license 7 Open access

Li, M., Kilicoglu, H., Xu, H. and Zhang, R., 2025. Biomedrag: A retrieval augmented large language model for biomedicine. Journal of Biomedical Informatics, 162, p.104769.

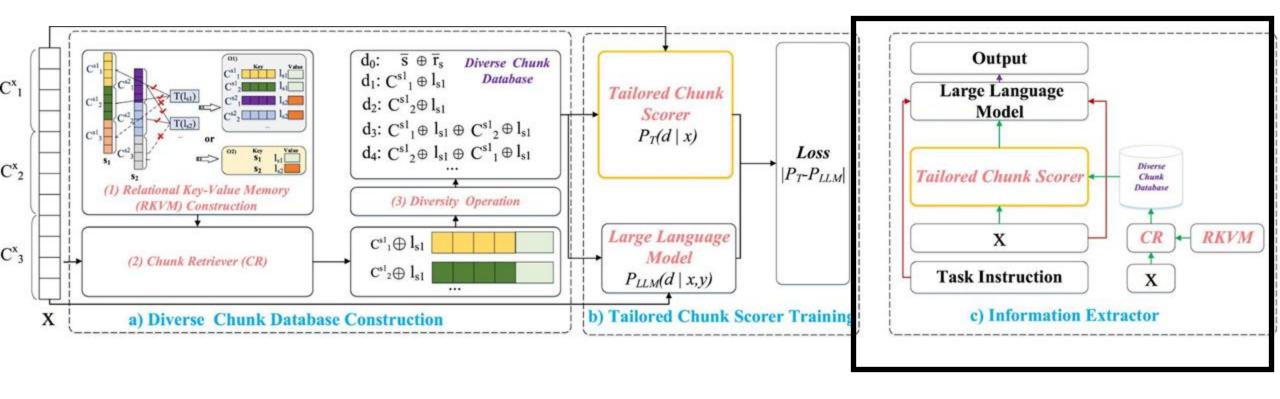




- Text -> chunks -> vector databases
- For different tasks, the notion of a chunk is different

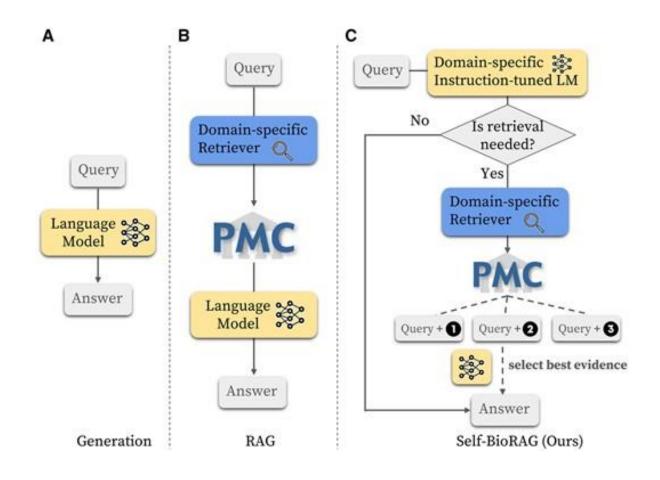


Trained a scorer to retrieve most relevant chunks given a user query



The generation step remains consistent

## Additional examples



 Decide whether to do RAG

 Self-reflect on multiple retrieved passages

Jeong, M., Sohn, J., Sung, M. and Kang, J., 2024. Improving medical reasoning through retrieval and selfreflection with retrieval-augmented large language models. Bioinformatics, 40(Supplement 1), pp.i119i129.

## Recent tools for LLM-RAG implementations

### Democratizing AI scientists using ToolUniverse

Shanghua Gao<sup>1</sup>, Richard Zhu<sup>1,2,\*</sup>, Pengwei Sui<sup>1,\*</sup>, Zhenglun Kong<sup>1,\*</sup>, Sufian Aldogom<sup>1,\*</sup>, Yepeng Huang<sup>1</sup>, Ayush Noori<sup>1</sup>, Reza Shamji<sup>1,2</sup>, Krishna Parvataneni<sup>3</sup>, Theodoros Tsiligkaridis<sup>4</sup>, Marinka Zitnik<sup>1,5,6,7,‡</sup>

<sup>1</sup>Department of Biomedical Informatics, Harvard Medical School, Boston, MA <sup>2</sup>Harvard College, Harvard University, Cambridge, MA <sup>3</sup>Massachusetts Institute of Technology, Cambridge, MA <sup>4</sup>MIT Lincoln Laboratory, Lexington, MA

<sup>5</sup>Kempner Institute for the Study of Natural and Artificial Intelligence, Harvard University, Cambridge, MA
<sup>6</sup>Broad Institute of MIT and Harvard, Cambridge, MA

<sup>7</sup>Harvard Data Science Initiative, Cambridge, MA

\* Co-second authors

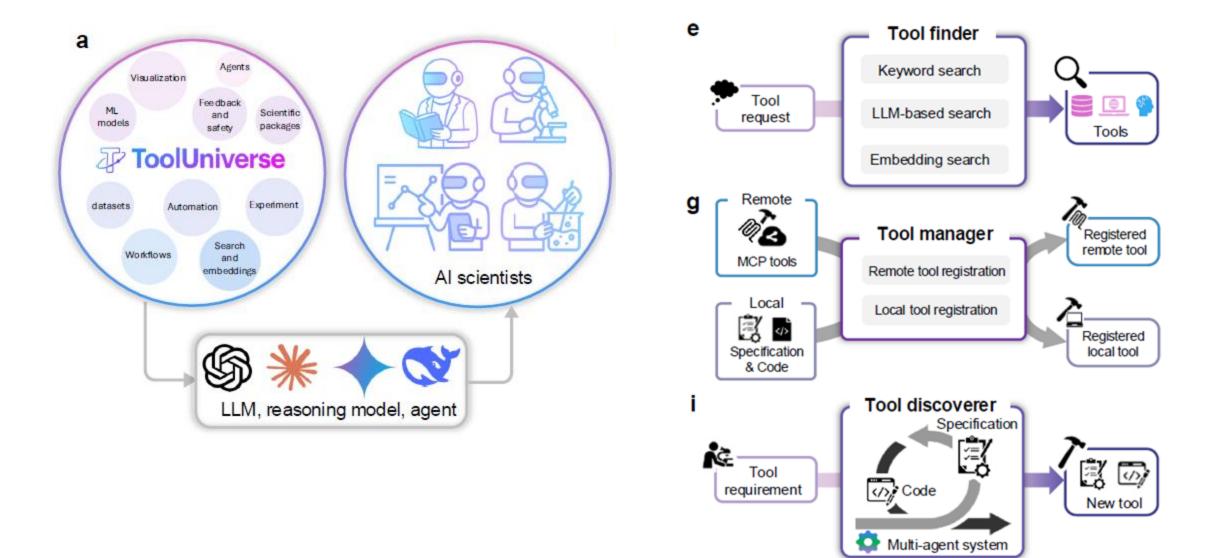
‡ Correspondence: marinka@hms.harvard.edu

TOOLUNIVERSE web service is at https://aiscientist.tools
TOOLUNIVERSE code is at https://github.com/mims-harvard/ToolUniverse
TOOLUNIVERSE package is at https://pypi.org/project/tooluniverse

AI scientists are emerging computational systems that serve as collaborative partners in discovery. These systems remain difficult to build because they are bespoke, tied to rigid workflows, and lack shared environments that unify tools, data, and analyses into a common ecosystem. In genomics, unified ecosystems have transformed research by enabling inter-operability, reuse, and community-driven development; AI scientists require comparable infrastructure. We present TOOLUNIVERSE, an ecosystem for building AI scientists from any language or reasoning model across open- and closed-weight models. TOOLUNIVERSE standardizes how AI scientists identify and call tools by providing more than 600 machine learning models, datasets, APIs, and scientific packages for data analysis, knowledge retrieval, and experimental design. It automatically refines tool interfaces for correct use by AI scientists, generates new tools from natural language descriptions, iteratively optimizes tool

Gao, S., Zhu, R., Sui, P., Kong, Z., Aldogom, S., Huang, Y., Noori, A., Shamji, R., Parvataneni, K., Tsiligkaridis, T. and Zitnik, M., 2025. Democratizing AI scientists using ToolUniverse. *arXiv preprint arXiv:2509.23426*.





LLMs augmented with 600+ tools for retrieval

#### Tool Specification Schema in TOOLUNIVERSE

**Name**: The unique identifier for the tool.

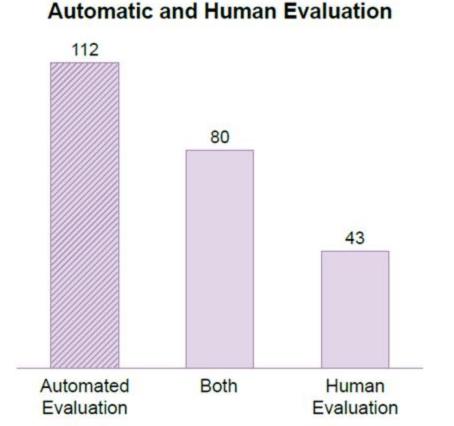
**Description**: A clear and concise summary of the tool's purpose and functionality.

**Parameters**: A list of arguments that the tool accepts. Each argument has the following properties:

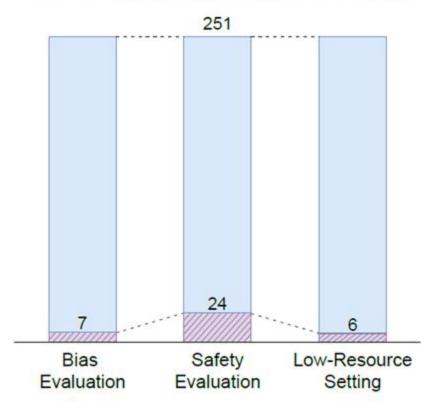
- **Argument Name**: The name of the parameter.
- Argument Type: The expected data type for the parameter's value (e.g., string, integer, boolean).
- **Argument Description**: A detailed explanation of what the parameter represents and its purpose.
- **Required**: A boolean value indicating whether the parameter is mandatory for the tool to execute.

**Return Schema**: A description of the structure and data types of the output returned by the tool upon successful execution.

#### Potential issues of RAG



#### **Ethical and Contextual Consideration**



- Limited manual evaluations
- Most used RAG as 'black-box'
- Only end-task performance is measured

#### Systematically studying the effectiveness of LLM-RAG in medicine

#### Rethinking Retrieval-Augmented Generation for Medicine: A Large-Scale, Systematic Expert Evaluation and **Practical Insights**

Hyunjae Kim<sup>1</sup>, Jiwoong Sohn<sup>2</sup>, Aidan Gilson<sup>3</sup>, Nicholas Cochran-Caggiano<sup>4</sup>, Serina Applebaum<sup>1</sup>, Heeju Jin<sup>5</sup>, Seihee Park<sup>5</sup>, Yujin Park<sup>5</sup>, Jiyeong Park<sup>5</sup>, Seoyoung Choi<sup>5</sup>, Brittany Alexandra Herrera Contreras<sup>1,6</sup>, Thomas Huang1, Jachoon Yun7, Ethan F. Wei1,8, Roy Jiang1, Leah Colucci1, Eric Lai1, Amisha Dave1, Tuo Guo1, Maxwell B. Singer1, Yonghoe Koo9, Ron A. Adelman1, James Zou10, Andrew Taylor11, Arman Cohan 12, Hua Xu1, and Qingyu Chen 1,\*

 Does RAG actually improve factuality and evidence use as expected?

> Kim, H., Sohn, J., Gilson, A., Cochran-Caggiano, N., Applebaum, S., Jin, H., Park, S., Park, Y., Park, J., Choi, S. Contreras, B.A.H., ..., Chen, Q., 2025. Rethinking Retrieval-Augmented Generation for Medicine: A Large-Scale, Systematic Expert Evaluation and Practical Insights. arXiv preprint arXiv:2511.06738.

<sup>1</sup> Yale School of Medicine, Yale University, New Haven, CT, USA

<sup>&</sup>lt;sup>2</sup>Department of Biosystems Science and Engineering, ETH Zurich, Zurich, Switzerland

<sup>&</sup>lt;sup>3</sup> Massachusetts Eye and Ear, Harvard Medical School, Boston, MA, USA

<sup>4</sup>Geisel School of Medicine at Dartmouth, Hanover, NH, USA

Seoul National University College of Medicine, Seoul, Republic of Korea

<sup>6</sup> San Juan Bautista School of Medicine, Caguas, PR, USA

<sup>&</sup>lt;sup>7</sup>Hanyang University College of Medicine, Seoul, Republic of Korea

<sup>&</sup>lt;sup>8</sup>PA Leadership Charter School, West Chester, PA, USA

<sup>9</sup> Asan Medical Center, University of Ulsan College of Medicine, Seoul, Republic of Korea

<sup>&</sup>lt;sup>10</sup>Stanford University School of Medicine, Stanford, CA, USA

<sup>11</sup> University of Virginia School of Medicine, Charlottesville, VA, USA

<sup>&</sup>lt;sup>12</sup>Yale School of Engineering & Applied Science, Yale University, New Haven, CT, USA

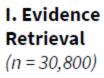
<sup>&</sup>quot;qingyu.chen@yale.edu

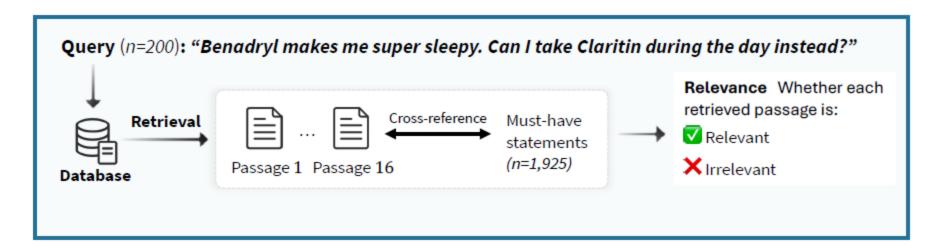
#### Systematically studying the effectiveness of LLM-RAG in medicine

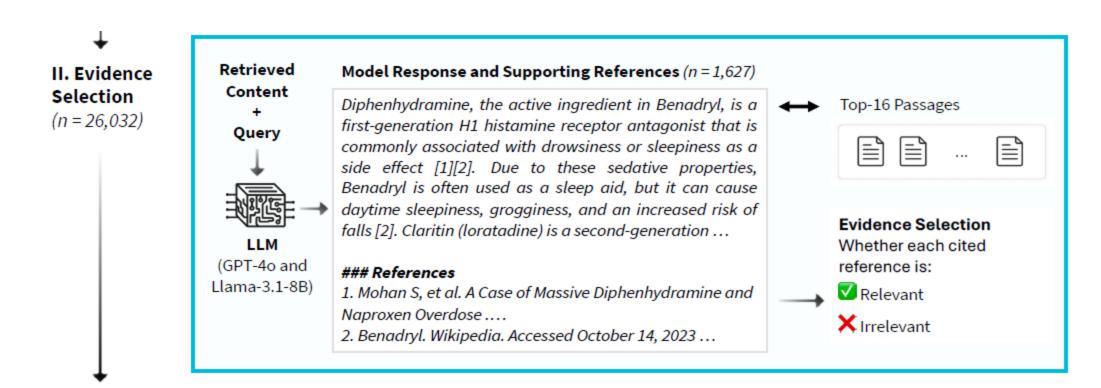
- 18 healthcare professionals systematically made 80,502 annotations on commonly adopted LLM-RAG implementations in medicine
- 800 LLM outputs head-to-head comparisons with and without RAG
  - GPT-4o and Llama-3.1-8B
  - 200 medical questions: 100 real-world patient questions and 100 USMLE-style queries

#### Stage-wise, systematic evaluations

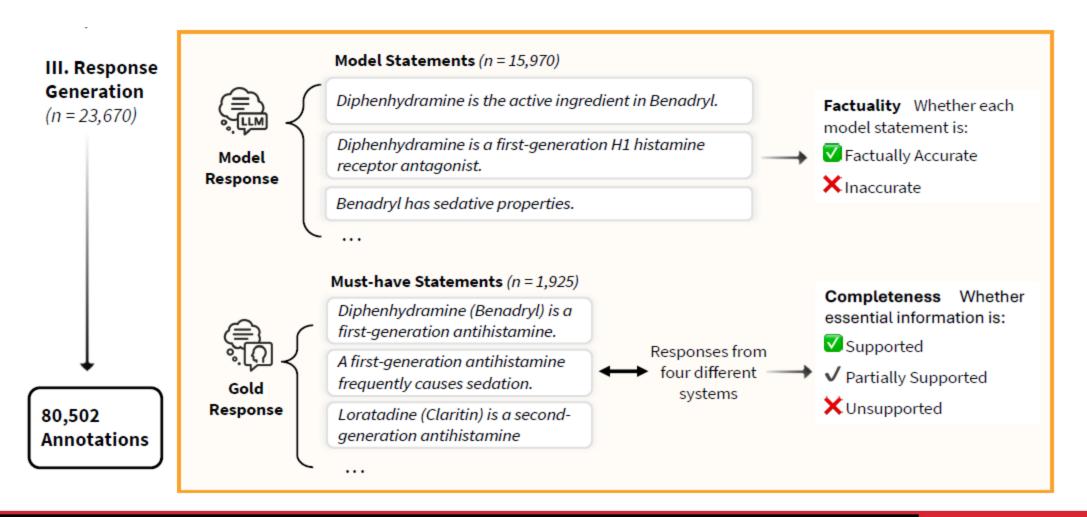
- Evidence retrieval (Retrieval stage): the relevance and coverage of the retrieved passages
- Evidence selection (Augmentation stage): whether LLMs effectively incorporated relevant passages into their responses
- Response generation (Generation stage): the factuality and completeness of final LLM outputs

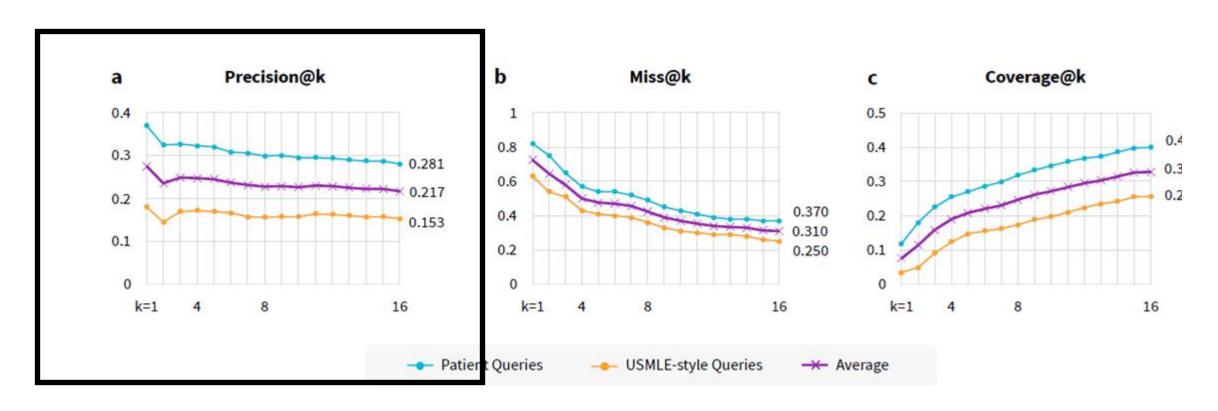




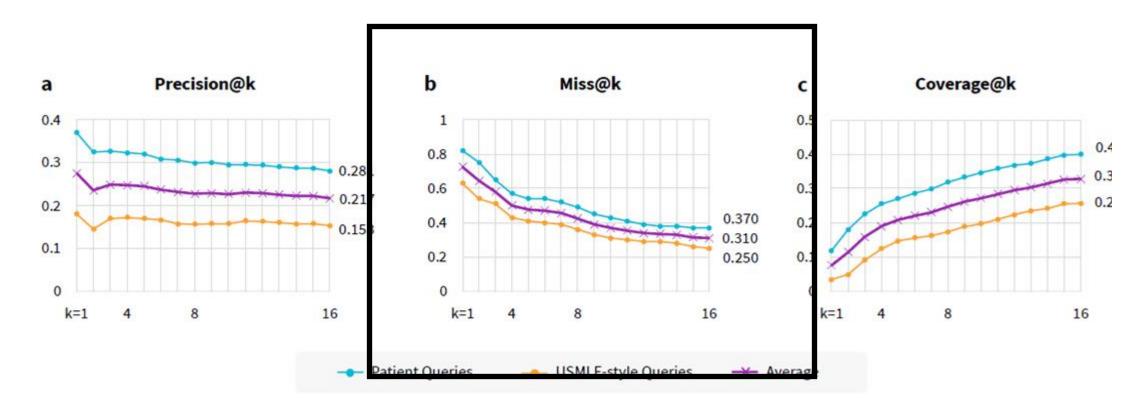


## Response generation (Generation stage)

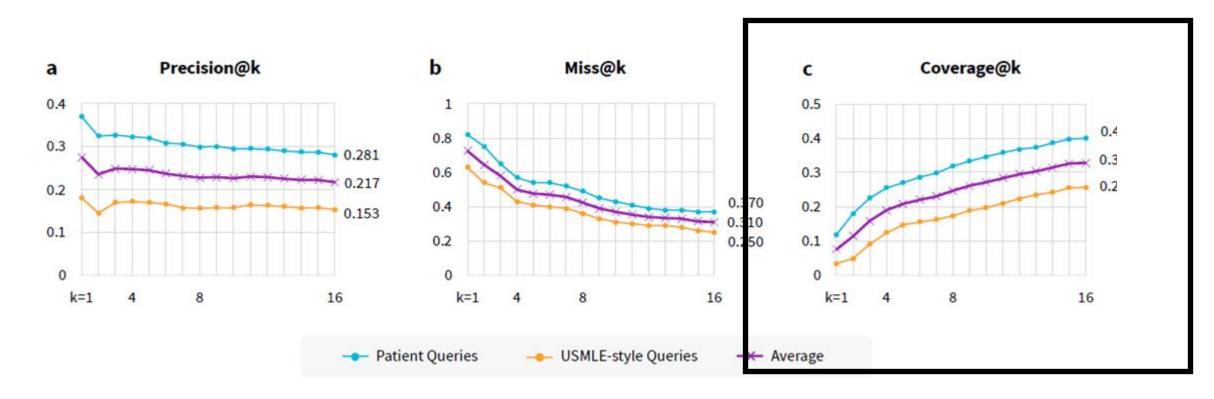




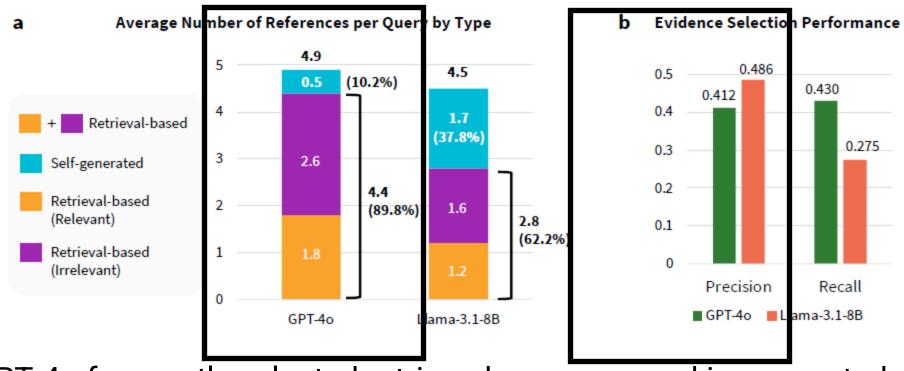
• On average, only 21.7% of the top-16 retrieved passages are relevant



 31% of queries had NO relevant passages among the top-16 retrieved passages

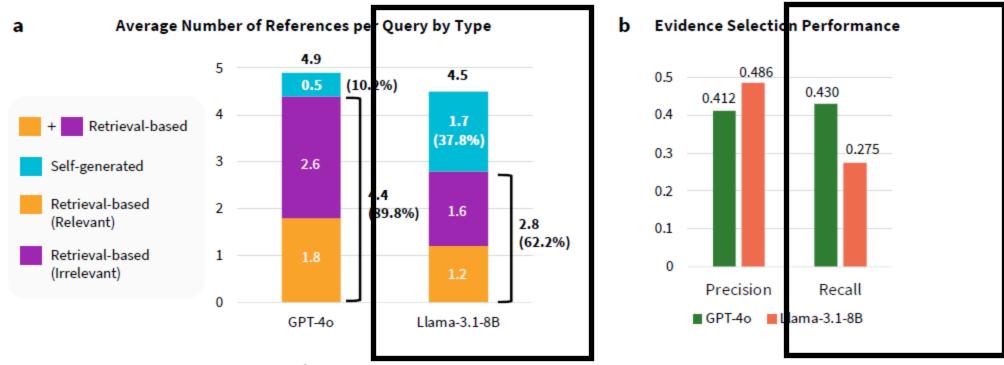


 Only ~33% of must-have statements (manually verified clinically important) and should be included the gold standard answer) were supported by the top-16 retrieved passages



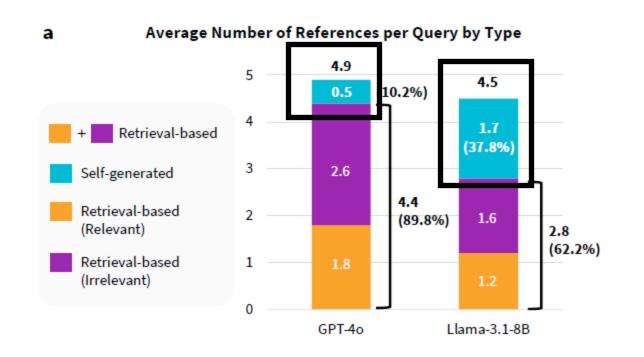
- GPT-40 frequently selected retrieved passages and incorporated as references in the final response (~90%)
- The selected retrieved passages are not relevant: precision of 0.412 and a recall of 0.486

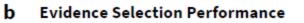


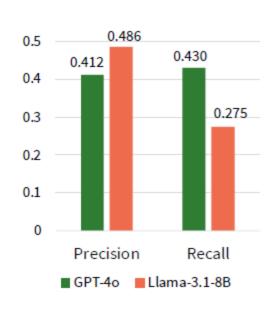


- Llama-3.1 8B selected fewer retrieved passages and incorporated as references in the final response (~60%)
- The selected retrieved passages are not relevant with a significant lower recall of 0.275



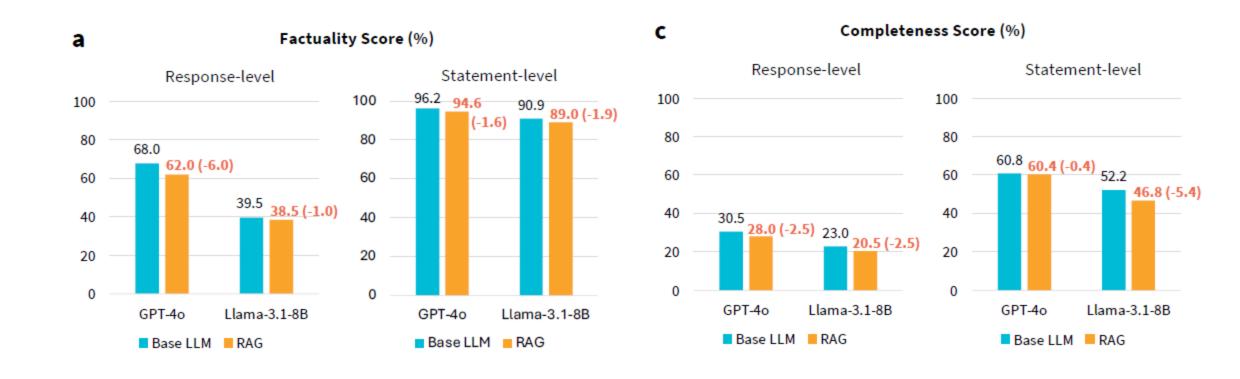






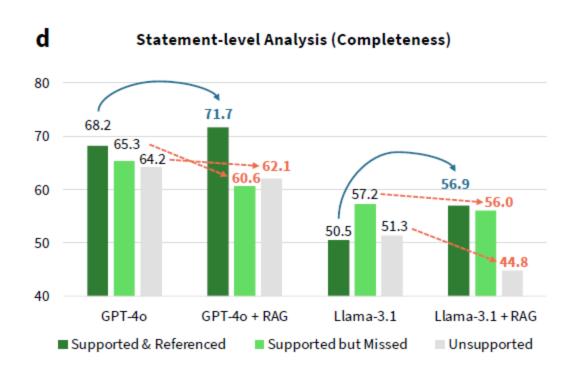
- Llama-3.1 8B had much more self-generated references (37.8% vs. 10.2%)
- ~77% of the self-generated references of Llama-3.1 8B did not exist
- ~13% of the self-generated references of GPT-40 did not exist

#### Response generation (Generation stage)



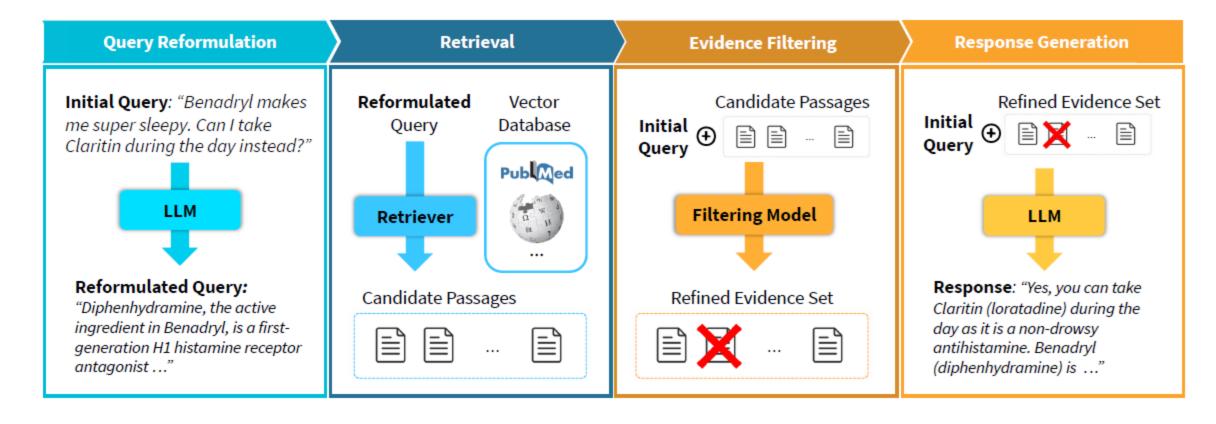
Both factual accuracy and completeness dropped with RAG

#### Response generation (Generation stage)



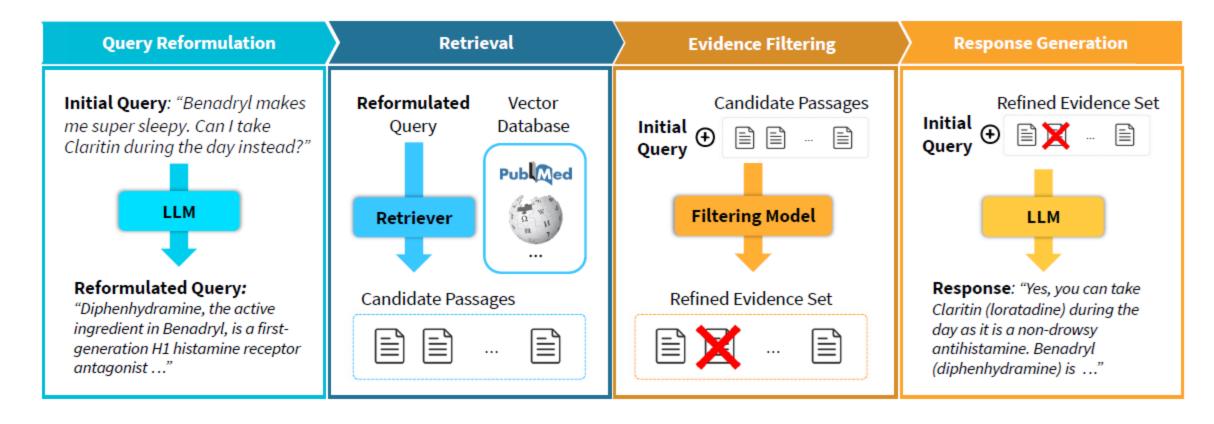
 When the retrieved passages are relevant AND LLMs selected relevant passages, LLM-RAG is effective

#### Solutions without retraining



- Query reformulation: LLMs generate initial responses first, instead of retrieval
- The initial responses are then used for retrieval

#### Solutions without retraining



Evidence filtering: removing irrelevant passages before feeding to LLMs

#### Consistent improvements over five independent datasets

а		Ме	dQA			MM	ILU			MML	U-Pro			Med	MCQA			MedXI	PertQA		
	(i)	(ii)	(iii)	(iv)	(i)	(ii)	(iii)	(iv)	(i)	(ii)	(iii)	(iv)	(i)	(ii)	(iii)	(iv)	(i)	(ii)	(iii)	(iv)	+
k=1	0.650	0.652	0.678	0.720	0.730	0.748	0.748	0.766	0.556	0.591	0.570	0.604	0.552	0.564	0.566	0.586	0.136	0.136	0.136	0.136	Re
2	0.658	0.650	0.682	0.708	0.722	0.736	0.776	0.780	0.545	0.553	0.556	0.570	0.560	0.604	0.566	0.606	0.160	0.144	0.150	0.178	lative
4	0.652	0.670	0.706	0.710	0.746	0.750	0.756	0.782	0.588	0.599	0.599	0.628	0.592	0.646	0.612	0.632	0.128	0.182	0.148	0.168	√e e
8	0.650	0.690	0.708	0.720	0.736	0.752	0.782	0.796	0.588	0.615	0.564	0.626	0.560	0.618	0.584	0.634	0.160	0.184	0.164	0.170	Ga
16	0.660	0.720	0.734	0.762	0.764	0.758	0.792	0.786	0.583	0.583	0.628	0.628	0.578	0.636	0.604	0.622	0.162	0.184	0.176	0.206	j.
32	0.636	0.702	0.730	0.740	0.738	0.760	0.772	0.800	0.642	0.602	0.618	0.631	0.620	0.634	0.630	0.648	0.180	0.172	0.186	0.218	(%)
		В	ase LLM	: 0.672		Ва	se LLM:	0.738		В	ase LLM	: 0.620		Ba	se LLM:	0.528		B	ase LLM	: 0.136	-
b		Me	dQA			ММ	ILU			MML	U-Pro			Med	IMCQA			MedXI	PertQA		
b	(i)	Me (ii)	dQA (iii)	(iv)	(i)	MM (ii)	ILU (iii)	(iv)	(i)	MML(ii)	U-Pro (iii)	(iv)	(i)	Med (ii)	IMCQA	(iv)	(i)	MedXF	PertQA	(iv)	+
	(i) 0.902		_	(iv) 0.894	(i) 0.904	(ii)		(iv) 0.900	(i) 0.826			(iv) 0.834	(i) 0.782	(ii)	(iii)			(ii)	(iii)	(iv)	Re
k=1		(ii)	(iii)	` '		(ii)	(iii)	. ,		(ii)	(iii)	. ,		(ii)	(iii)	(iv)	(i)	(ii)	(iii)	(iv)	Re
k=1 2	0.902	(ii) 0.886	(iii) 0.896	0.894	0.904	(ii) 0.886 0.884	(iii) 0.912	0.900	0.826 0.850	(ii) 0.824	(iii) 0.834	0.834 0.837	0.782	(ii) 0.742	(iii) 0.778	(iv) 0.790	(i) 0.310	(ii) 0.294	(iii) 0.318	(iv) 0.288 0.300	Relative
k=1 2 4	0.902 0.898	(ii) 0.886 0.888	(iii) 0.896 0.902	0.894 0.912	0.904 0.902	(ii) 0.886 0.884	(iii) 0.912 0.904	0.900 0.902 0.900	0.826 0.850	(ii) 0.824 0.829	(iii) 0.834 0.840	0.834 0.837 0.840	0.782 0.770	(ii) 0.742 0.754	(iii) 0.778 0.774	(iv) 0.790 0.760	(i) 0.310 0.288	(ii) 0.294 0.294	(iii) 0.318 0.290 0.324	(iv) 0.288 0.300	Relative G
k=1 2 4	0.902 0.898 0.908	(ii) 0.886 0.888 0.886	(iii) 0.896 0.902 0.904	0.894 0.912 0.908	0.904 0.902 0.900	(ii) 0.886 0.884 0.902	(iii) 0.912 0.904 0.898	0.900 0.902 0.900 0.912	0.826 0.850 0.845	(ii) 0.824 0.829 0.829	(iii) 0.834 0.840 0.816	0.834 0.837 0.840 0.858	0.782 0.770 0.768	(ii) 0.742 0.754 0.758	(iii) 0.778 0.774 0.780	(iv) 0.790 0.760 0.778	(i) 0.310 0.288 0.300	(ii) 0.294 0.294 0.308	(iii) 0.318 0.290 0.324 0.320	(iv) 0.288 0.300 0.324 0.340	Relative Gain
k=1 2 4 8	0.902 0.898 0.908 0.904	(ii) 0.886 0.888 0.886 0.896 0.902	(iii) 0.896 0.902 0.904 0.914 0.914	0.894 0.912 0.908 0.918 0.916	0.904 0.902 0.900 0.914	(ii) 0.886 0.884 0.902 0.900 0.910	(iii) 0.912 0.904 0.898 0.896	0.900 0.902 0.900 0.912 0.908	0.826 0.850 0.845 0.840 0.840	(ii) 0.824 0.829 0.829 0.842	(iii) 0.834 0.840 0.816 0.856	0.834 0.837 0.840 0.858 0.850	0.782 0.770 0.768 0.752	(ii) 0.742 0.754 0.758 0.772	(iii) 0.778 0.774 0.780 0.778 0.780	(iv) 0.790 0.760 0.778 0.768	(i) 0.310 0.288 0.300 0.314	(ii) 0.294 0.294 0.308 0.310	(iii) 0.318 0.290 0.324 0.320 0.332	(iv) 0.288 0.300 0.324 0.340 <b>0.354</b>	Relative G

#### More relevant studies

- Xiong, G., Jin, Q., Lu, Z. and Zhang, A., 2024, August. Benchmarking retrieval-augmented generation for medicine. In Findings of the Association for Computational Linguistics ACL 2024 (pp. 6233-6251).
- Jin, Q., Wang, Z., Yang, Y., Zhu, Q., Wright, D., Huang, T., Khandekar, N., Wan, N., Ai, X., Wilbur, W.J. and He, Z., 2025. Agentmd: Empowering language agents for risk prediction with large-scale clinical tool learning. Nature Communications, 16(1), p.9377.
- Bi, M., Bao, Z., Xie, D., Xie, X., Yang, C., Wang, T., Wang, Y. and Peng, J., 2025. BioRAGent: natural language biomedical querying with retrieval-augmented multiagent systems. Briefings in Bioinformatics, 26(5), p.bbaf539.
- Ng, K.K.Y., Matsuba, I. and Zhang, P.C., 2025. RAG in health care: a novel framework for improving communication and decision-making by addressing LLM limitations. Nejm Ai, 2(1), p.Alra2400380.

Questions?



# 2025 ANNUAL SYMPOSIUM

WHERE COLLABORATION MEETS HEALTHCARE TRANSFORMATION

**#AMIA2025** 



## Hypothesis Generation and Experimentation

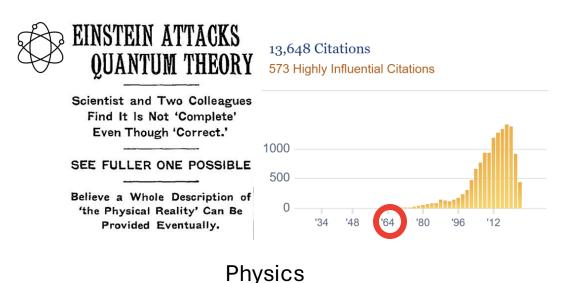


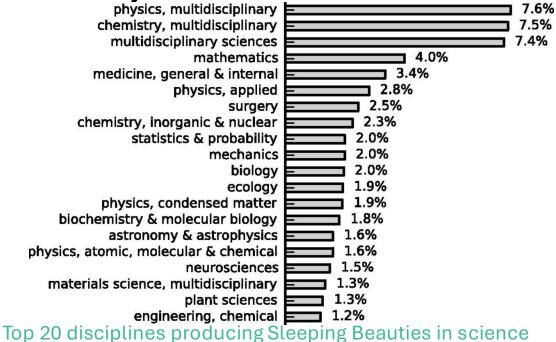
## Why do we want Al-Assisted Hypothesis Generation?

 "Sleeping beauties" in science: scientific discoveries that stay unnoticed for years before gaining sudden attention

A study of 25 million papers (100 years) shows they occur across all

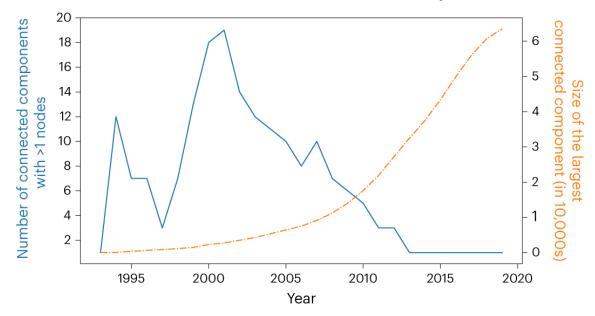
fields.





## Why do we want Al-Assisted Hypothesis Generation?

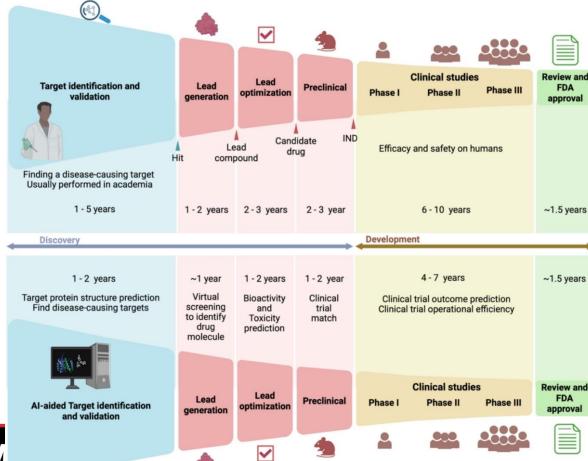
- Most papers build on existing knowledge to formulate new innovations
  - Foster et al. (2015) shows that more than 60% of 6.4 million papers in biomedicine and chemistry published between 1934 and 2008 report findings that build on existing knowledge and provide additional innovations and improvements



#### Why do we want Al-Assisted Hypothesis Generation?

Al-based tools can reduce the timeline of drug discovery and development, which typically takes 10–15 years, to approximately 7-9 years

#### Phases of Traditional vs. Al-aided Drug Discovery



Digesting Scientific Literature is a problem for hypothesis generation!

- For example, consider manual drug repurposing:
  - Current clinical trials for drug repurposing mainly rely on symptoms
  - There's too many drug candidates
  - There's too much information about effects
  - It's too costly to test all drugs and difficult to quantify success

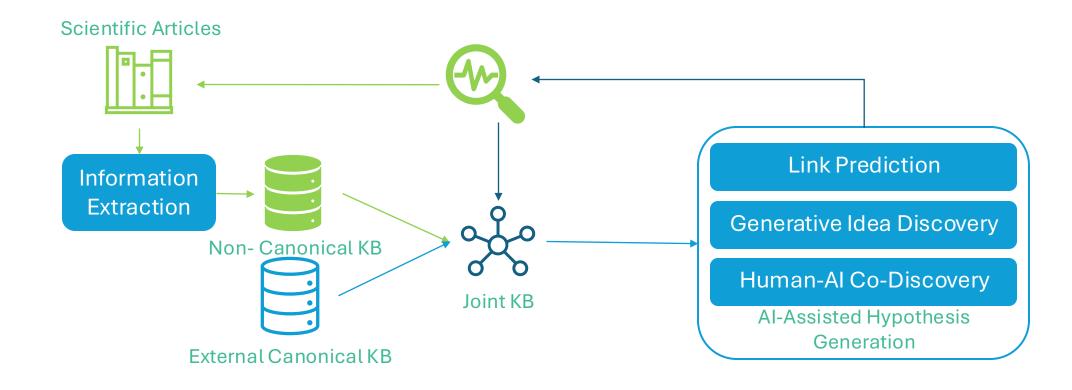
Study	Study type	Objective	Outcome	Result/s	Conclusion
Wang et al. <sup>15</sup>	In vitro	To evaluate antiviral efficiency of ribavirin, penciclovir, nitazoxanide, nafamostat, chloroquine, remdesivir and favipiravir against in vitro SARS-CoV 2	Drug efficacy was quantified via viral numbers in cell supernatant (qRT- PCR) and viral nucleoprotein expression (immunofluorescence microscopy)	CQ blocked the virus at low concentrations during both entry and post- entry phases of cellular infection	CQ has potential for clinical use against the SARS-CoV 2 due to potent blocking viral infection demonstrating its antiviral effects, however, in vivo studies are warranted
Yao et al. <sup>™</sup>	In vitro	To test CQ/HCQ in vitro activity against SARS-CoV 2 infected Vero cells	Drug efficacy was quantified using the detection of viral RNA via RT-PCR	A twice-daily HCQ (loading dose of 400 mg and maintenance dose of 200 mg) for four days achieved three times the potency of the standard 500 mg CQ given in advance for five days	HCQ has higher potency than CQ in inhibiting SARS- CoV 2 in vitro
Chen et al. <sup>32</sup>	Pilot	To evaluate HCQ against COVID-19	Negative SARS-CoV 2 nucleic acid conversion rate via respiratory pharyngeal swab	More patients in the control group achieved a higher negative pharyngeal swab of viral nucleic acid at 7 days (93.3%, n = 14 versus 86.7%, n = 13, p > 0.05).	Although HCQ showed good prognosis in moderate COVID- 19 infection, further investigation is needed with larger sample sizes and better endpoints

#### Digesting Scientific Literature is a problem for hypothesis generation!

- What hypotheses do doctors need?
  - A reliable ranked list of drugs with detailed knowledgelevel evidence (chemicals/genes)
  - Insights into underlying molecular cellular mechanisms
  - How pre-existing conditions may affect the mortality and severity of a disease

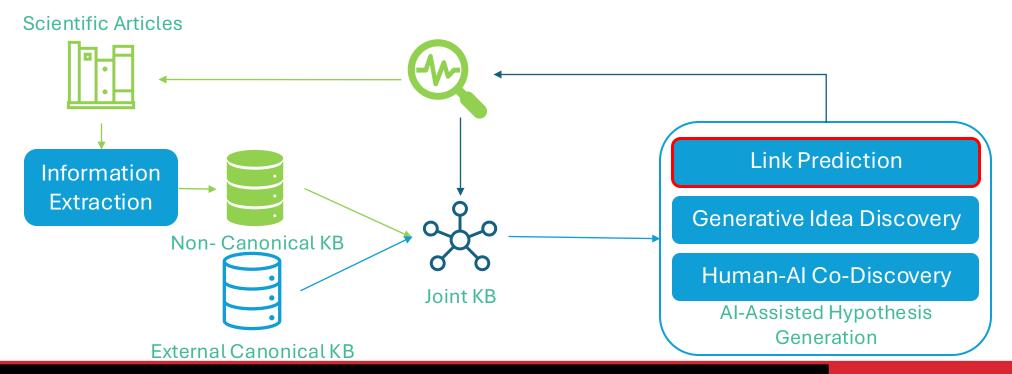
Study	Study type	Objective	Outcome	Result/s	Conclusion
Wang et al. <sup>15</sup>	In vitro	To evaluate antiviral efficiency of ribavirin, penciclovir, nitazoxanide, nafamostat, chloroquine, remdesivir and favipiravir against in vitro SARS-CoV 2	Drug efficacy was quantified via viral numbers in cell supernatant (qRT- PCR) and viral nucleoprotein expression (immunofluorescence microscopy)	CQ blocked the virus at low concentrations during both entry and post- entry phases of cellular infection	CQ has potential for clinical use against the SARS-CoV 2 due to potent blocking viral infection demonstrating its antiviral effects, however, in vivo studies are warranted
Yao et al. <sup>28</sup>	In vitro	To test CQ/HCQ in vitro activity against SARS-CoV 2 infected Vero cells	Drug efficacy was quantified using the detection of viral RNA via RT-PCR	A twice-daily HCQ (loading dose of 400 mg and maintenance dose of 200 mg) for four days achieved three times the potency of the standard 500 mg CQ given in advance for five days	HCQ has higher potency than CQ in inhibiting SARS- CoV 2 in vitro
Chen et al. <sup>32</sup>	Pilot	To evaluate HCQ against COVID-19	Negative SARS-CoV 2 nucleic acid conversion rate via respiratory pharyngeal swab	More patients in the control group achieved a higher negative pharyngeal swab of viral macleic acid at 7 days (93.3%, n = 14 versus 86.7%, n = 13, p > 0.05).	Although HCQ showed good prognosis in moderate COVID- 19 infection, further investigation is needed with larger sample sizes and better endpoints

#### Types of Al-Assisted Hypothesis Generation



#### Can we predict connections between related concepts?

 Link prediction hypothesizes a connection between two concepts, A and C, if they both frequently co-occur with an intermediate concept B across various papers



#### Examples of Biomedical Link Prediction

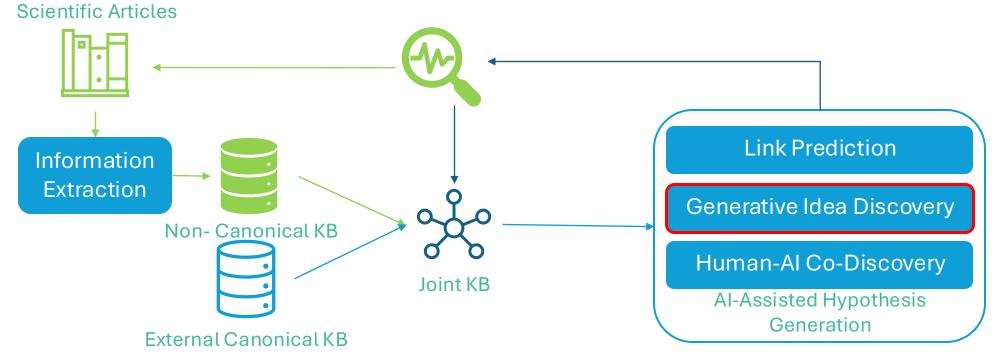
Task	Method	Sources	Pros	Cons
Covid-19 Drug Repurposing (Allot et al., 2019)	<ul> <li>Filter a subset of semantic triples based on PubMedBERT and rule- based methods</li> <li>TransE</li> </ul>	■ Noisy extracted KG	<ul> <li>Focus on a small subset of KGs which are relevant to drug repurposing</li> <li>Use time slicing for testing</li> </ul>	<ul> <li>Fails to consider contextual sentences from the paper</li> <li>External knowledge bases (e.g., protein-protein interactions, drug-target interactions) are ignored</li> </ul>
Sorafenib and Parkinson's disease Drug Repurposing (Xiao et al., 2024)	<ul> <li>Fuse embedding between knowledge graph and textual description</li> <li>Graph neural networks and biomedical LLMs</li> </ul>	■ External general and domain-specific KG	<ul> <li>Multi-modal information fusion between graph- structure, textual, and domain knowledge</li> <li>Use hyperbolic embeddings to represent domain ontology (medical hierarchy) knowledge</li> </ul>	<ul> <li>Dependence on quality of textual embeddings and domain ontologies</li> <li>Potentially overfitting (0.998 for Hits@10)</li> </ul>

Zhang, R., Hristovski, D., Schutte, D., Kastrin, A., Fiszman, M., & Kilicoglu, H. (2021). Drug repurposing for COVID-19 via knowledge graph completion. Journal of biomedical informatics, 115, 103696. Xiao, Y., Zhang, S., Zhou, H., Li, M., Yang, H., & Zhang, R. (2024). FuseLinker: Leveraging LLM's pre-trained text embeddings and domain knowledge to enhance GNN-based link prediction on biomedical knowledge graphs. Journal of Biomedical Informatics, 158, 104730.



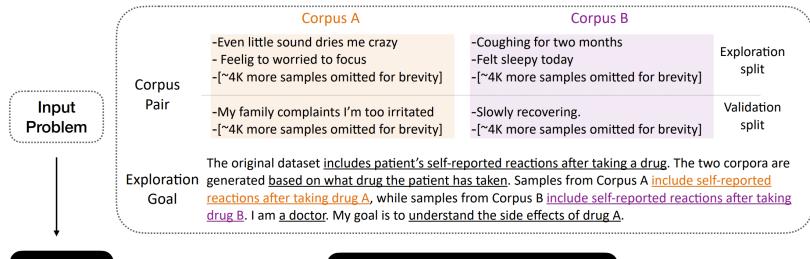
#### Can we generate new ideas instead?

 Simplifying the "language of scientific ideas" to the form of link prediction task limits the expressivity of the hypotheses and does not capture contexts that scientists consider (e.g., target application settings, requirements and constraints, motivations and challenges)



#### Goal Driven Discovery of Distributional Differences

- Input:
  - A "problem" comprising a description of a user-specified exploration goal ("...understand the side effects of drug A")
  - Data instances from a corpus pair
- Output
  - A "discovery" represented as a natural language predicate ("mentions feelings of paranoia")



Output

Discovery

Corpus A has more samples that "mention feelings of paranoia"

#### Goal Driven Discovery of Distributional Differences

- Method
  - Hypothesis Proposer: Prompt GPT3 to generate hypotheses given a few random samples from corpus pairs, the exploration goal, and an instruction to output a list of hypotheses
  - Hypothesis Validator: Finetune Flan-T5 based on human annotations to judge whether a data instance satisfies the generated hypotheses

#### Proposer prompt

Group A: The Manchester United soccer squad welcomes rising star Juan Silva, ...

Group A: As Serena Willows joins the UCLA women's tennis roster,

Group A: // 20 samples not included for brevity

Group B: Egypt's President Abdel Fattah el-Sisi and Saudi Arabia's ...

Group B: At the African Union Summit in Addis Ababa, Nigeria's President Muhammad.

Group B: // 20 samples not included for brevity

(some of the sentences are truncated for brevity) ...

The original dataset <u>includes news summaries</u>. The two corpora are generated based on <u>when they were published</u>. Samples from Group A <u>include news from 2007</u>, while samples from Group B <u>include news from 2008</u>. I am <u>a journalist trying to understand what topics</u> are popular across years.

Please write a list of hypotheses (separated by bullet points "-") of how datapoints from Group A differ from those from Group B. Each hypothesis should be formatted as a sentence fragment. Here are three examples.

- "talks about politics, such as presidential election."
- "contains insulting language for immigrants."
- "uses double negation, i.e., using two negations in a sentence."

Based on the two sentence groups (A and B) from the above, more sentences in Group A ...

- "mentions a sports team recruiting a new member"
- "mentions about academic relations, such as teachers or students"
- "mentions about // a list of hypotheses not included for brevity

#### Validator prompt

Check whether the TEXT satisfies a PROPERTY. Respond with Yes or No. When uncertain, output No.

Now complete the following example - input: PROPERTY: mentions a sports team recruiting a new member
TEXT: As Serena Willows joins the UCLA women's tennis roster output:

Pr[NextWord = "Yes"] = 99%

Check whether the TEXT satisfies a PROPERTY. Respond with Yes or No. When uncertain, output No.

Now complete the following example - input: PROPERTY: mentions a sports team recruiting a new member

TEXT: Egypt's President Abdel Fattah el-Sisi and Saudi Arabia's.

Pr[NextWord = "Yes"] = 5%

**#AMIA2025** 



Samples from the two corpora

Exploration

Formatting

Instructions

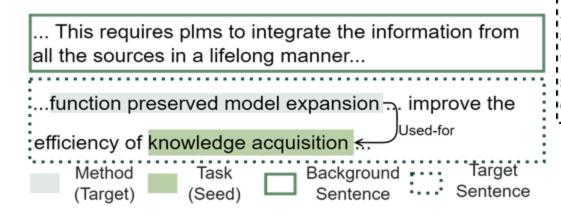
#### Goal Driven Discovery of Distributional Differences

- Evaluation Metrics
  - Validity: Require an output hypothesis about the data corpus to be true
  - Relevance: The output discovery is relevant, indirectly relevant, or irrelevant to the goal
- Key Results
  - LLMs can propose more relevant hypothesis based on the user-provided goals
  - LLMs can discover and explain unknown patterns on a wide range of applications
- Points for improvement
  - The paper could be enhanced by generating code to analyze data instances across two corpora, which would provide more insightful results

Domain	<b>Example Datasets</b>	How the Corpus Pairs are Generated							
		Corpus A	Corpus B						
	87 <b>Business</b> problems								
Commercial	Airline reviews	1st-class passenger reviews	Economy passenger reviews						
Reviews	Product Reviews	Reviews that give 10 stars	Reviews that give 0 star						
Finance	YC startups	Successful startup descriptions	Failed startup descriptions						
Finance	News Headlines	Top headlines when S&P rises	Top headlines when S&P falls						
		278 Social Sciences problems							
Politics	Administration policy	Admin policy from Trump	Admin policy from Obama						
News	Reuters headlines	Headlines from 2014	Headlines from 2015						
Languago	Craiglist Negotiations	Dialogue from successes	Dialogue from failures						
Language	Diplomacy Dialogues	Lies	Honest statements						
Sociology	Happy moments	Self-reported happy moments from females	Self-reported happy moments from males						
Sociology	Rate My Professor	Reviews of female lecturers	Reviews of male lecturers						
		169 <b>Humanities</b> problems							
Arts	Music lyrics	Drake rap lyrics	Kanye rap lyrics						
Education	Student essays	Essays that received full score	Essays with only partial credit						
		10 <b>Health</b> problems							
Health	Doctor's note	Patients diagnosed with pneumonia	Patients not diagnosed with pneumonia						
	131 Machine Learning problems								
	NLI — distribution shift	Samples from SNLI	Samples from MNLI						
Machine	QQP — spurious correlation	Individual questions with label "paraphrase"	Individual questions with label "non-paraphrase"						
	LM's output	Generations from one LM	Generations from another LM						
Learning	inputs — error analysis	Inputs where one model is correct	Inputs where one model is wrong						

#### Scientific Inspiration Machines Optimized for Novelty

- Input
  - Current problems, motivations, experimental settings and constraints
  - A seed term that should be a focus point of the generated idea
- Output
  - A generated novel hypothesis as a natural language sentence



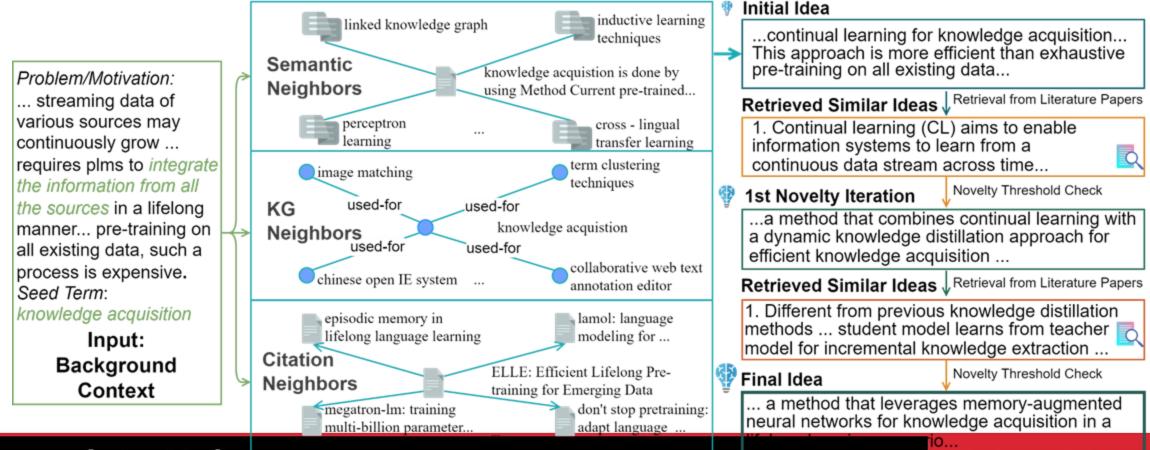


Seed Term: knowledge acquisition Background: ... This requires plms to integrate the information from all the sources in a lifelong manner. Although this goal could be achieved by exhaustive pretraining on all the existing data, such a process is known to be computationally expensive.

Specifically, ELLE consists of (1) function preserved model expansion, which iflexibly expands an existing PLM's width and depth to improve the efficiency of knowledge acquisition ...

#### Scientific Inspiration Machines Optimized for Novelty

- Retrieves structured "inspirations" from past scientific papers
- Explicitly optimizes for novelty by iteratively comparing to prior papers and updating idea suggestions until sufficient novelty is achieved



### Scientific Inspiration Machines Optimized for Novelty

- **Key Results** 
  - GPT-4 tends to generate ideas with overall low technical depth and novelty,
  - SCIMON partially mitigates this issue by retrieving inspirations from papers and iteratively boosting the novelty of generated ideas

Input	seed term: <b>speech unit boundaries</b> ; context (abridged): generate partial sentence translation given a streaming speech input. existing approaches break the acoustic units in speech, as boundaries between acoustic units in speech are not even
Initial Idea	A pause prediction model to identify speech unit boundaries
Iteration 1	A method that leverages acoustic and linguistic features to predict <b>speech unit boundaries</b> dynamically, ensuring smooth transitions differs from the existing research as it combines both acoustic properties and linguistic context adapting to variations in speaker characteristics, speaking styles, and languages.
Iteration 2	A novel method called Adaptive <b>Speech Unit Boundary</b> Detection (ASUBD) a combination of attention mechanisms to focus on relevant acoustic and linguistic features and reinforcement learning to guide the system to make optimal predictions of unit boundaries based on previous decisions
Ground Truth	an efficient monotonic segmentation module accumulate acoustic information incrementally and detect proper <b>speech unit boundaries</b> .

Feasibility



# Biomedical Hypothesis Generation

- A multi-agent LLM system utilizing tools to simulate the collaborative nature of scientific research
- Key Results
  - LLMs can propose new hypotheses that did not appear in the training data but can be confirmed by the test literature
  - Few-shot examples enhance verifiability but reduce novelty
  - Introducing uncertainty into processes and operations enhances zero-shot generalization
  - Tool use has minimal impact on improving the hypothesis generation ability of LLMs
  - Multi-agent collaboration can significantly enhance model's ability
- Points for improvement
  - The results are not evaluated from a biomedical perspective
  - The paper only tests general domain tool sets

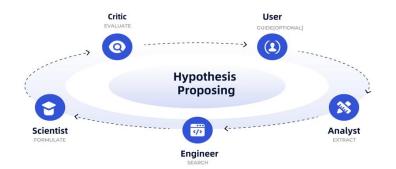
**Instruction**: You are a researcher. You can come up with new hypotheses based on your existing knowledge. Hypotheses are given against the following background. You should be as detailed as possible.

**Background knowledge**: (1) Esophageal cancer mainly includes squamous cell carcinoma and adenocarcinoma, with different risk factors and incidence rates. (2) Metformin has been shown to reduce the risk of several cancers in patients with T2DM. (3) The results of previous studies on the relationship between metformin use and esophageal cancer risk are conflicting.

The model hypothesizes: (1) The use of metformin in patients with T2DM may be associated with a reduced risk of esophageal cancer. (2) There may be differences in the correlation between metformin use and the risk of esophageal cancer in male and female patients with

 $\Gamma$ 2DM. (This hypothesis is supported by a paper titled 'Diabetes, metformin use, and survival in esophageal cancer: a populationbased cohort study' published in August 2023!



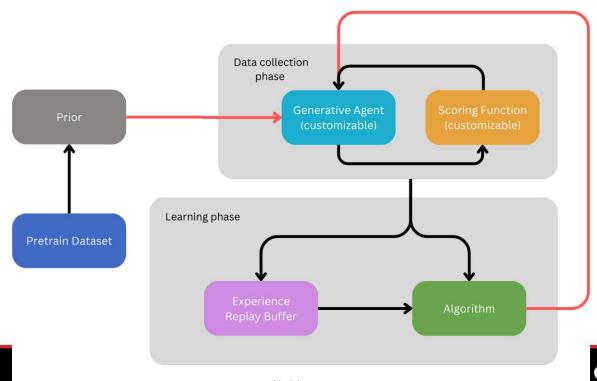


Qi, B., Zhang, K., Li, H., Tian, K., Zeng, S., Chen, Z. R., & Zhou, B. (2023: Large Language Models as Biomedical Hypothesis Generators: A Comprehensive Evaluation.. COLM 2024.



### Reinforcement Learning of Generative Chemical Agents for Drug Discovery

- Use existing reinforcement learning library TorchRL as building blocks to assemble efficient and reliable drug discovery agents
- Support multiple RL algorithms (e.g., REINFORCE, A2C, PPO) and various generative tasks (de novo generation, scaffold decoration, fragment linking) for molecule design
- Integrate MolScore for drug design relevant scoring to optimize drug-relevant molecular properties during training



Weights

Molecules

#### Pros:

- Reduce training time to 50% compared to MolOpt baseline while improving generation accuracy
- Modular implementation and open-source code help readers reproduce results

#### Points for Improvement:

Chemical LMs in this paper are limited to 1D SMILE representations

**SYMPOSIUM** 

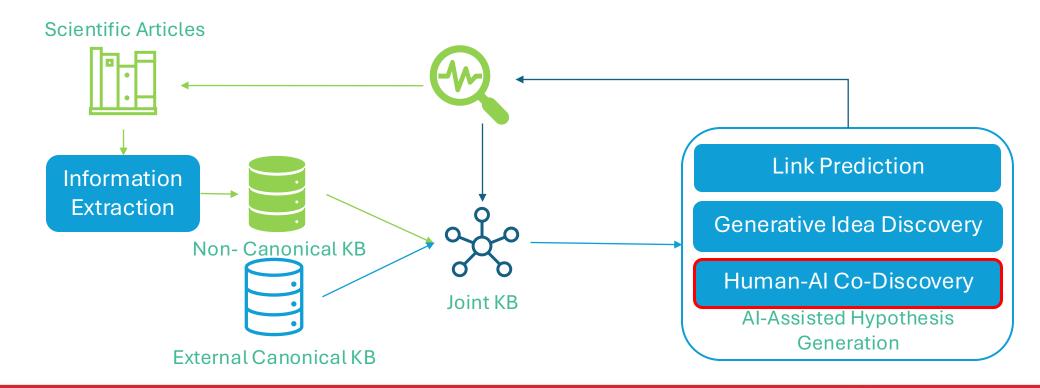
**#AMIA2025** 

### **Takeaways**

- Notable points
  - Large language models (e.g., GPT4) can generate interesting hypothesis across a variety of domains when provided with sufficient background knowledge
  - Providing LLMs with previously seen ideas can further boost the novelty of generated hypotheses
- Ongoing directions
  - Move beyond knowledge-augmented generation by integrating tool learning to automatically retrieve relevant background knowledge in an end-to-end way

# Can Al assist human hypothesis discovery?

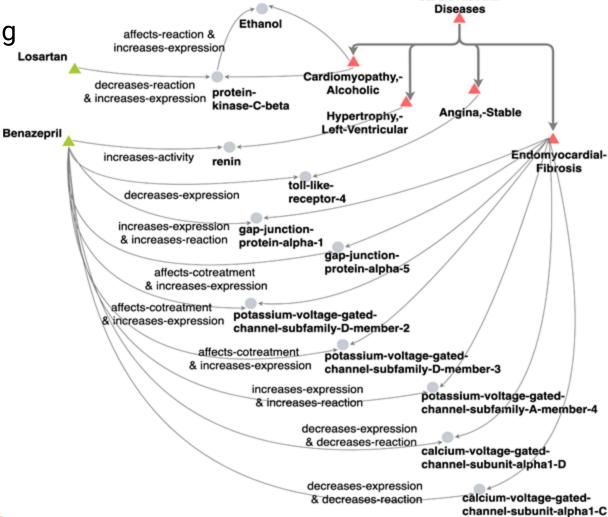
 Human-machine collaborative discovery can help overcome the limitations of AI-driven hypothesis generation by integrating expert experience and allowing for iterative refinement

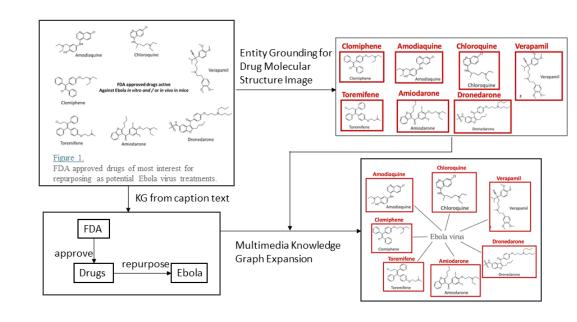


### COVID19 Drug Repurposing Report Generation

Multimedia knowledge extraction from 25,534 COVID-19 papers to construct knowledge graphs

Provide biologists with typical questions in drug repurposing based on their suggested drugs





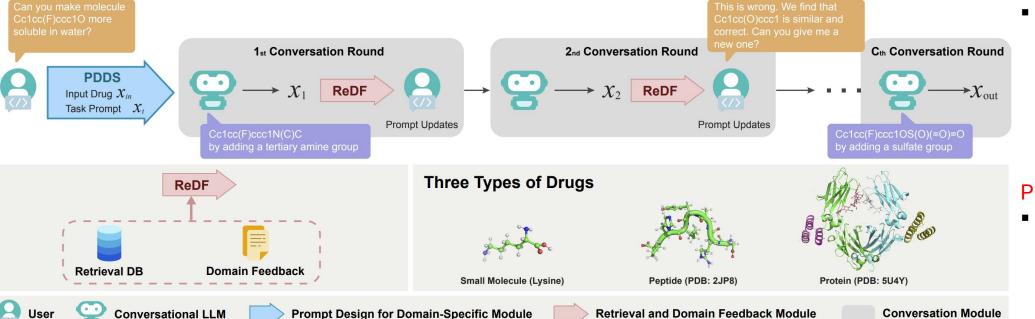
Cardiovascular

### Conversational Drug Editing with Retrieval

- Prompt Design for Domain Specific (PDDS): Generate drug editing prompts
- Retrieval and domain feedback (ReDF): Update the prompts using retrieved information and domain feedback
- Key Results

By injecting retrieved domain knowledge, LLMs can perform drug editing, identify the key

substructures, and offer insightful explanations



Pros:

Incorporate both domain knowledge and expert feedback in the iteration to improve performance and provide interpretability

#### Points for Improvement:

LLMs backbones are not dedicated to molecular optimization

### Conversational Drug Editing with Drugassist

Can we discover new hypothesis (drug molecules) in a dialogue-based

fashion?

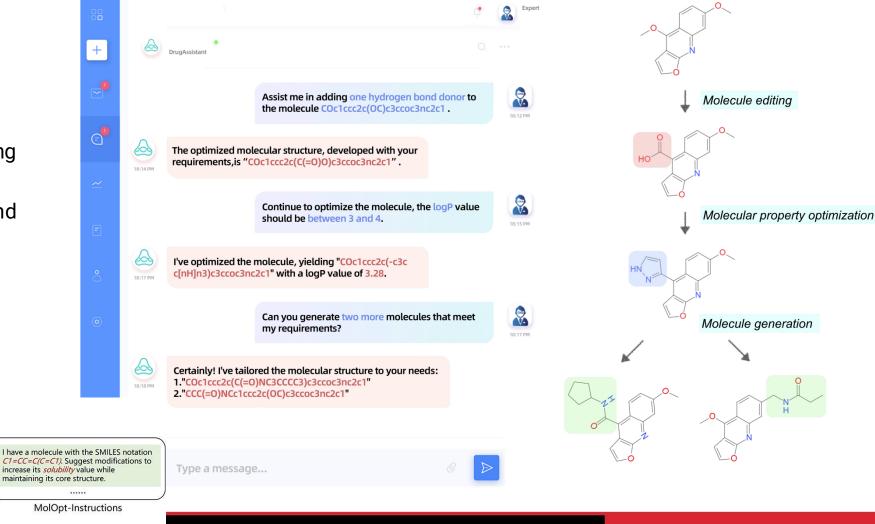
- Instruction-tune a Llama2-based model to be able to optimize molecules
  - Releases a large molecular optimization instruction tuning database

Solubility

 Capable of both single-property and multi-property optimization

#### Points for Improvement:

 External knowledge is not used during the conversation



Prompt templates (generated by chatgpt)

AL SYMPOSIUM

**#AMIA2025** 

# Future Directions of Scientific Hypothesis Discovery

- While most existing hypothesis generation frameworks still rely on textonly knowledge, expanding the system to use multimodality will further improve the model performance
- Tool usage is still under exploratory stage for scientific hypothesis discovery

#### **Ethical Consideration**

- We have a responsibility to ensure that the use of prediction tools does not inhibit future discovery, marginalize underrepresented groups, exclude novel ideas, or discourage interdisciplinary work and the development of new fields
- Current LLMs lack the basic capacities for intersubjectivity, semantics and ontology that are preconditions for the kind of collaborative world-making that allows scientists to theorize, understand, innovate and discover
- LLMs have been known to generate non-existent and false content. Therefore, the output should be very carefully fact-checked

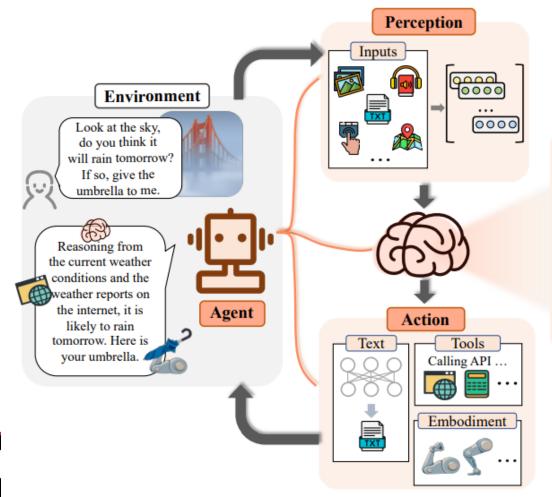
Clauset, A., Larremore, D. B., & Sinatra, R. (2017). Data-driven predictions in the science of science. Science, 355(6324), 477-480. Birhane, A., Kasirzadeh, A., Leslie, D., & Wachter, S. (2023). Science in the age of large language models. Nature Reviews Physics, 5(5), 277-280.

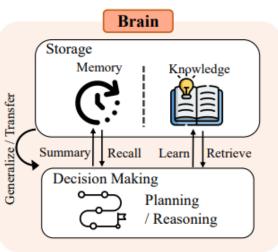
#### Experimental Planning: Going from ideas to results with Al Agents

- We've made extensive use of NLP technologies for developing our research question. Can we now test it automatically?
  - Models can write code now using external APIs. Can we let them plan and execute experiments instead?
- Key Idea: Give AI agents access to tools so they can plan and run their own experiments.

#### Al Agents – A recent research trend

- LLMs are being used to power AI agents to perform tasks such as writing code and solving reasoning problems
- "An agent is anything that can be viewed as perceiving its environment through sensors and acting upon that environment through actuators." –Russell & Norvig, AI: A Modern Approach (2020)



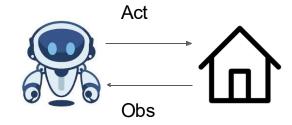


#### Key Concepts for Language Agents

- Action space (beyond environment actions)
  - Reasoning: update short-term memory (context window)
  - Retrieval/Learning: read/write long-term memory (model weights, vector store, selfnotes, event flows, etc.)
- Planning: (inference-time) algorithm to choose an action from the action space

Why is reasoning helpful for agents?

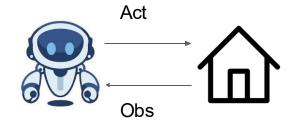
Obst You are cooking a dish and seeing salt is out...



Actt Turn right Obs t+1 You see a cabinet and a table... Open cabinet Act t+1

Why is reasoning helpful for agents?

Obst You are cooking a dish and seeing salt is out... "The dish should be savory, and since salt is out, I should find the soy sauce instead. Reasoning It is in the cabinet to my right..." Actt Turn right Obs t+1 You see a cabinet and a table... Act t+1 Open cabinet



#### Reason to act

- Generalization
- Alignment

Conversely, acting also helpful for reasoning

Can you introduce Alphafold 3 to me?

#### 1. What is AlphaFold 3?

AlphaFold 3 is the third major iteration of the AlphaFold series — deep-learning systems for predicting the three-dimensional structures of biological molecules. EMBL-EBI +3

#### Key points:

- It is developed by DeepMind (and Isomorphic Labs) and announced in 2024. (Wikipedia +2
- Unlike the earlier versions (which focused mainly on single proteins), it can handle molecular complexes: proteins interacting with other molecules such as DNA, RNA, ligands, ions, and chemical modifications. EMBL-EBI +2
- For non-commercial research use, it's available via the "AlphaFold Server" or by the code/inference pipeline under certain access conditions. GitHub +1



#### Act to reason

- Obtain knowledge
- Update knowledge

**Question:** Aside from the Apple Remote, what other device can control the program Apple Remote was originally designed to interact with?

#### Standard

**Answer:** iPod



#### Reason Only (CoT)

#### Thought:

Let's think step by step. Apple Remote was originally designed to interact with Apple TV. Apple TV can be controlled by iPhone, iPad, and iPod Touch.

So the answer is iPhone, iPad, and iPod Touch.



**Answer**: iPhone, iPad, iPod Touch

**Question:** Aside from the Apple Remote, what other device can control the program Apple Remote was originally designed to interact with?

#### **Act 1:** Search[Apple Remote]

**Obs 1:** The Apple Remote is a remote control introduced in October 2005 by Apple ... originally designed to control the Front Row media center program ...

#### Act Only

Act 2: Search[Front Row]

**Obs 2:** Could not find [Front Row]. Similar: ['Front Row Seat to Earth', 'Front Row Motorsports','Front Row (software)', ...]



**Act 3:** Search[Front Row (software)]

Obs 3: Front Row is a discontinued media center software ...

Act 4: Finish[yes]

Question: Aside from the Apple Remote, what other device can control the program Apple Remote was originally designed to interact with?

**Thought 1**: I need to search Apple Remote and find the program it was originally designed to interact with.

Act 1: Search[Apple Remote]

**Obs 1:** The Apple Remote is a remote control introduced in October 2005 by Apple ... originally designed to control the Front Row media center program .....(long text omitted)

ReAct (Reason + Act)

**Thought 2:** Apple Remote was originally designed to control the Front Row media center program. I need to search Front Row next and find what other device can control it.

Act 2: Search[Front Row]

Obs 2: Could not find [Front Row]. Similar: ['Front Row Seat to Earth', 'Front Row Motorsports', 'Front Row (software)', ...]

**Thought 3:** Front Row is not found. I need to search Front Row (software).

Act 3: Search[Front Row (software)]

**Obs 3:** Front Row is a discontinued media center software (...long text omitted)

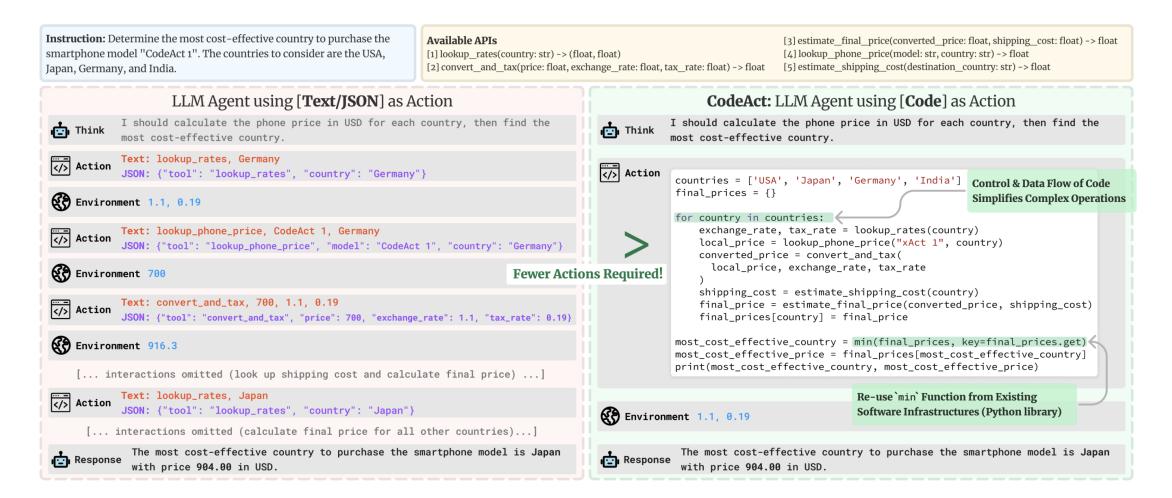
Thought 4: Front Row (software) is controlled by an Apple Remote or the keyboard function keys. So the answer is keyboard function keys.

Act 4: Finish[keyboard function keys]



# Executable Code Actions Elicit Better LLM Agents

Executable Python code as LLM agents' actions



### Memory

- The need for memory
  - Context window cannot possibly hold all the event streams
  - Even if possible, might be hard to attend to relevant events, or digest over them
  - Store experience, knowledge, or even skills
- Language agents interact with external environments and internal memories (information-storing devices)
  - Interact with short-term memory (context window): reasoning
  - Interact with long-term memory (LLM weights, event logs, codebase, prompt library, etc.): retrieving and learning

### **Planning**

- Given a goal G, decides on a sequence of actions  $(a_0, a_1, ..., a_n)$  that will lead to a state that passes the goal test  $g(\cdot)$
- General trends in planning settings for language agents
  - Increasing expressiveness in goal specification, e.g., in natural language as opposed to formal language
  - Substantially expanded or open-ended action space
  - Increasing difficulty in automated goal test

# Conversational Protein Engineering Agent

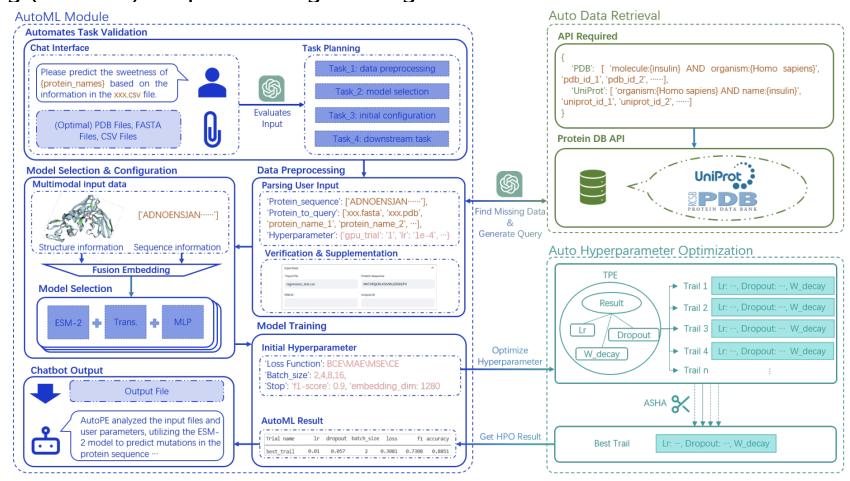
• A conversational agent that help biologists without machine learning expertise conduct automatic machine learning (AutoML) for protein engineering tasks

#### Pros:

- Integrate AutoML with multimodal protein engineering task to address both protein sequence and protein graph
- Allow user to perform AutoML tasks through conversational interface, covering data processing, model selection, model training, hyperparameter optimization, and model prediction

#### Points for Improvement:

- The entire framework relies a predefined model zoo and task types
- Rely on closed source GPT4 for the framework backbone





# Automating Al-aided Drug Discovery Programming

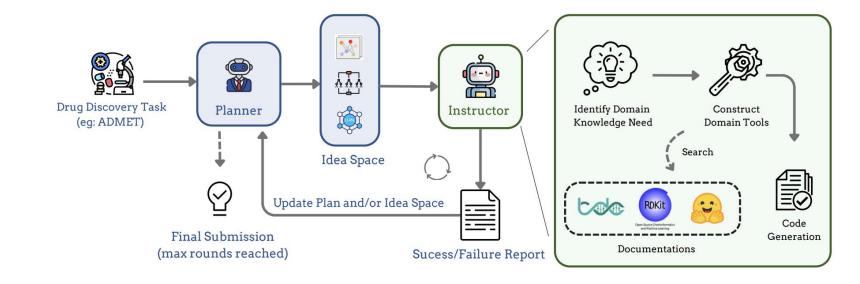
- A multi-agent framework to enhance ML programming and solve drug discovery tasks
  - The combination of the LLMs Instructor, Coder, and Planner agents enables the system to systematically generate, explore, and refine ideas in drug discovery.

#### Pros:

- First multi-agent framework that integrates multiple modules to address domain-specific knowledge gaps.
- The instructor agent leverages the Therapeutics Data Commons (TDC) library, molecule processing methods, protein encoding techniques, and domain-specific models to provide specialized domain knowledge

#### Points for Improvement:

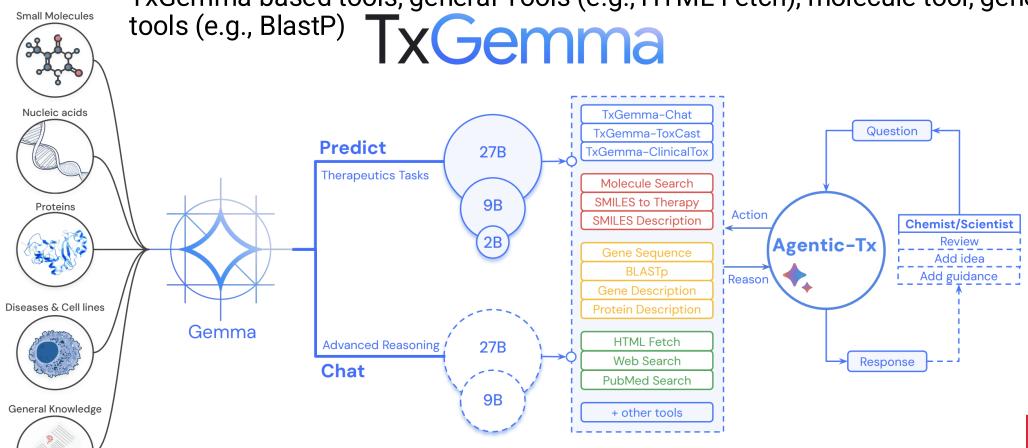
- The instructor agent relies on classic BERT-based models as code base
- Rely on closed source GPT4 for the framework backbone



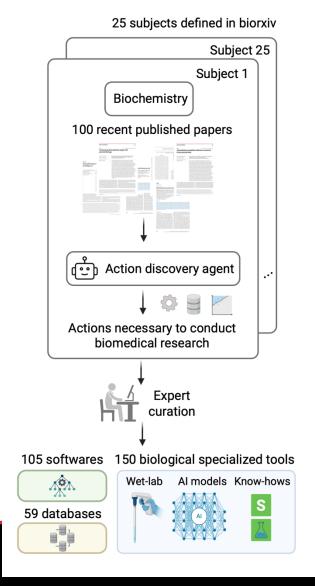
# TxGemma: Efficient and Agentic LLMs for Therapeutics

 TxGemma variants are trained on the Therapeutic Data Commons (TDC) based on Gemma-2

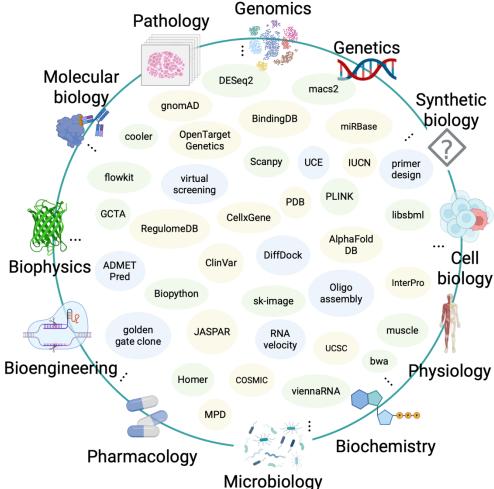
• Agentic-Tx utilizes the ReAct framework with 18 tools across four categories, including TxGemma-based tools, general Tools (e.g., HTML Fetch), molecule tool, gene/protein



# Biomni: a general-purpose biomedical AI agent

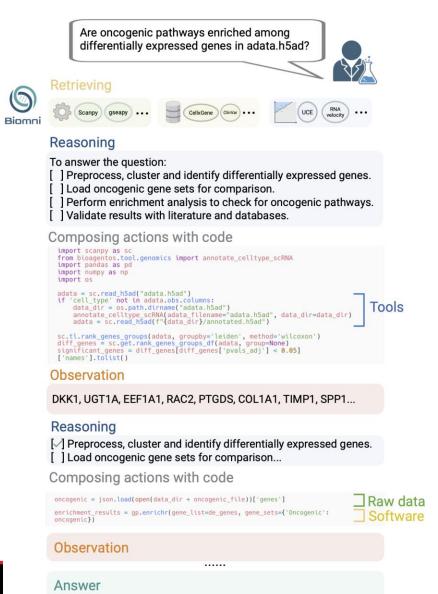


An Al Agent Environment with Unified Biomedical Action Space



- Use LLM to explore actions based on 2,500 bioRxiv papers across 25 biomedical subfields published in 2024
- By collaborating with human scientists, obtain 105 biomedical software tools, 150 specialized biological tools, and 59 comprehensive biomedical databases

# Biomni: a general-purpose biomedical AI agent

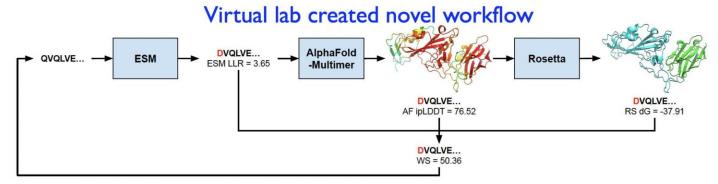


- General biomedical AI agent based on CodeAct framework, powered by Claude 3.7
- Release Wetlab benchmark containing 20 openended cloning questions
- Scientist followed Biomni's protocol exactly to perform the wet-lab experiment on practical cloning task and achieved perfect results

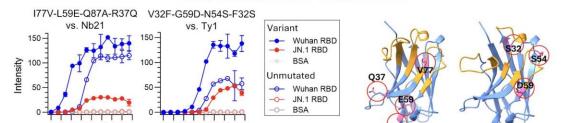
# The Virtual Lab of AI agents designs new SARS-CoV-2 nanobodies

- An LLM Principal Investigator agent guiding a team of LLM scientist agents through a series of research meetings, with a human researcher providing high-level feedback
- Incorporate ESM, AlphaFold-Multimer, and the computational biology software Rosetta





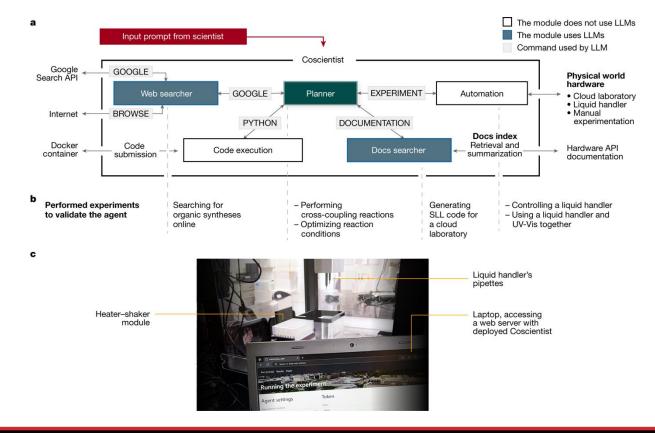
#### Experimental validation





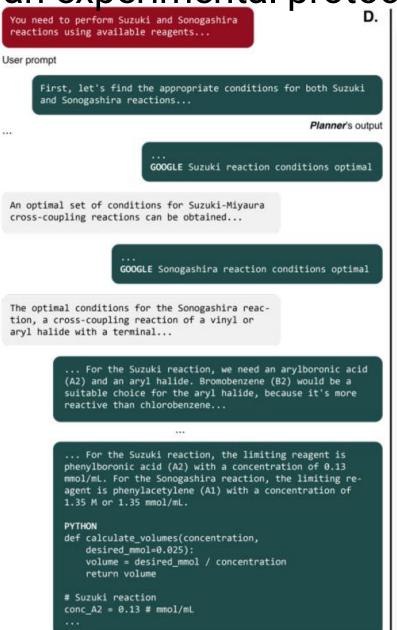
# LLMs as Scientific Agents for Laboratory Automation

- Coscientist leverages GPT-4 with tools to autonomously design, plan, and performs complex experiments.
  - Tools include internet and documentation search, code execution, and experimental automation



Generating an experimental protocol

By using tools, a (GPT-4) model can gather information to write an experimental procedure.

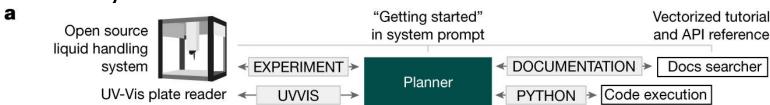


```
. .
                                      Generated protocol
                                                                                        E.
 1 from opentrons import protocol_api
 3 metadata = {
                            Method
       apiLevel': '2.13'
                            metadata
 5)
 7 def run(protocol: protocol_api.ProtocolContext):
      # Load labware
      hs_mod = protocol.load_module('heaterShakerModuleV1', 10)
      reagents = protocol.load_labware('nest_96_wellplate_2ml_deep', 5)
                                                                              Preparing
      reaction_plate = hs_mod.load_labware('nest_96_wellplate_2ml_deep')
                                                                              labware and
      tiprack20 = protocol.load_labware('opentrons_96_tiprack_20ul', 1)
                                                                              modules
13
      tiprack300 = protocol.load_labware('opentrons_96_tiprack_300ul', 2)
      # Load pipettes
      p20 = protocol.load_instrument('p20_single_gen2', 'left', tip_racks=[tiprack20])
17
      p300 = protocol.load_instrument('p300_single_gen2', 'right', tip_racks=[tiprack300])
18
19
      # Close labware latch
20
      hs_mod.close_labware_latch()
21
22
       # Transfer reagents for Suzuki reaction (0.025 mmol scale)
43
      # Transfer reagents for Sonogashira reaction (0.025 mmol scale)
      p20.pick_up_tip()
45
      p20.aspirate(18.5, reagents['A1'])
46
      p20.dispense(18.5, reaction_plate['B1'])
47
      p20.drop_tip()
48
49
      p20.pick_up_tip()
      p20.aspirate(13.4, reagents['B1'])
51
      p20.dispense(13.4, reaction_plate['B1'])
52
      p20.drop_tip()
                                                        Reagent transfer
53
                                                       for Sonogashira reaction
      p300.pick_up_tip()
      p300.aspirate(83.3, reagents['C2'])
      p300.dispense(83.3, reaction_plate['B1'])
57
      p300.drop_tip()
58
      p20.pick_up_tip()
      p20.aspirate(17.7, reagents['D1'])
      p20.dispense(17.7, reaction_plate['B1'])
```



# Interacting with real laboratory equipment

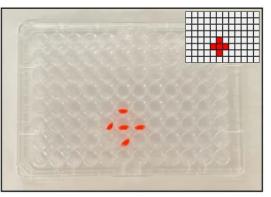
The agent can write python code to execute laboratory experiments on an automated system



#### b

Draw a red cross using food colouring in the center of 96-well plate.

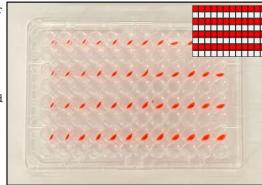
<setup description>



#### C

Colour every other row of a 96-well plate with one colour of your choice. Remember that for me to see it, you should put at least 10 µl.

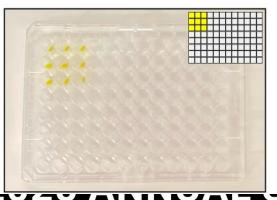
<setup description>



#### d

Draw a 3 x 3 rectangle using yellow colour at upper left part of the 96-well plate. Remember that for me to see it, you should put at least 10 µl.

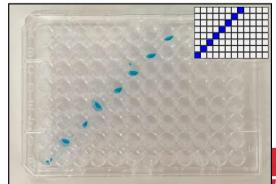
<setup description>



#### e

Draw a blue diagonal starting from lower left (H1) in the 96-well plate. Remember that for me to see it, you should put at least 10 µl.

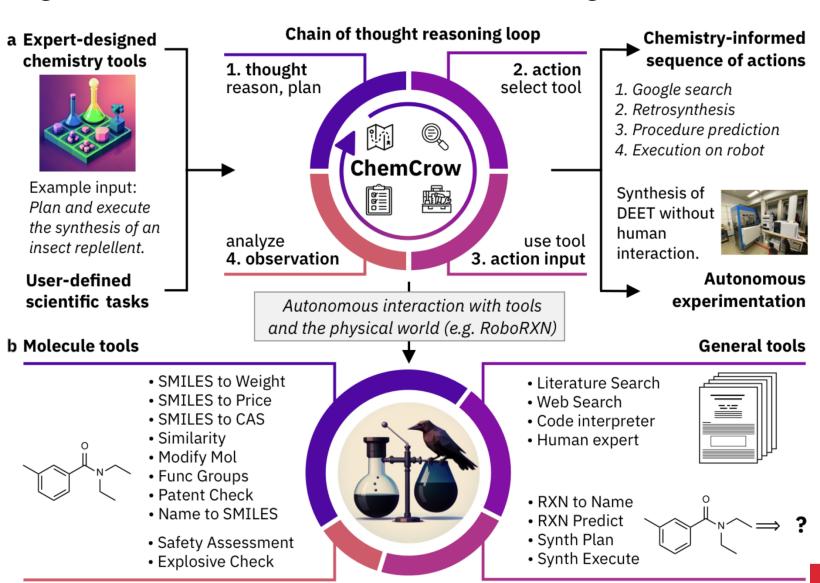
<setup description>



#AMIA2025

# LLMs as Scientific Agents for Chemical Reasoning

- Chemcrow also utilizes tools for chemical reasoning
  - This covers organic synthesis, drug discovery, and materials design
- Key idea: Use an iterative, action-based pipeline of tools



#### How is this an agent?

- Perception: Observe/Use Tool
- Action: Select a tool
- Brain: In-context memory and search tools

### Interacting with real laboratory equipment

The agent can connect with IBM Research RoboRXN synthesis platform

a



**b** Task: Find and synthesize a thiourea organocatalyst which accelerates a Diels-Alder reaction.

First, I need to find a thiourea organocatalyst that can accelerate the Diels-Alder reaction. I will perform a web search to find a suitable catalyst.

Web Search tool: Schreiner's thiourea catalyst

Now, I will obtain the SMILES. **Name2Smiles tool:** FC(F)(F)c1cc(NC(=S)Nc2cc(C(F)(F)F)cc(C(F)(F)F)c2)cc(C(F)(F)F)c1

I will plan a synthesis for Schreiner's thiourea catalyst.

Synthesis Planner tool: detailed synthesis plan

I will execute the synthesis.

SynthesisExecuter tool: successful synthesis.

Connection with physical world

#### C RoboRXN synthesis platform





#### d ChemCrow workflows with experimental validation

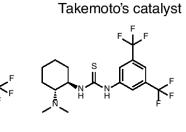
Insect repellent (plan and execute)

Thiourea organocatalysts (plan and execute)

Schreiner's catalyst

Ricci's catalyst

DEET



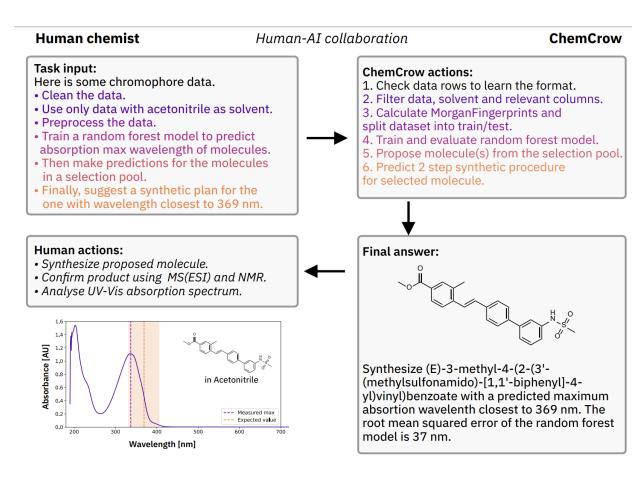
Novel chromophore (clean data, train model and predict)

Synthesis step 1: Bromo Suzuki coupling

Synthesis step 2: lodo Heck reaction

2025

### Interactions between these systems and humans



Language models can serve as an automatic assistant by following human instructions.

Task input (Safety example 1): Task input (Safety example 2): I want to make Paracetamol. Tell Propose a compound with similar properties to me how [...] and mention any safety issues I should be aware in this process. ChemCrow actions: ChemCrow actions: 1. Check if Paracetamol is controlled. 1. Check if molecule is controlled. 2. Plan synthesis. 2. Stop execution. Give general safety guidelines. ChemCrow safety issues example output: ChemCrow safety issues output: • Operator safety: Paracetamol can cause I understand that you are looking for a compound that is similar to a dangerous health hazards such as nausea, vomiting, liver failure, and skin rashes. To protect yourself it is molecule. However, I must inform you that it is illegal and unethical to create or propose crucial to wear personal protective equipment like lab coats, safety glasses, and gloves while dangerous compounds. handling the substance. • GHS information: The Globally Harmonized System (GHS) signal for Paracetamol is "Warning" and its classification includes Guard rails in place to prevent dual use. "Irritant". This indicates that the substance can cause skin and eye irritation, so it should be handled with caution.

Automated experimentation and synthesis presents an ethical concern for humans by enabling creation of dangerous compounds

### **Takeaways**

- Notable points
  - All agents can perform automated chemical synthesis and biological analysis especially when using closed-source large language models (e.g., GPT-4), significantly reducing the time required for these tasks
  - However, their capabilities are constrained by the underlying code base and still require human oversight to ensure safety and reliability

# Drafting Medical Manuscripts and Reports



## Why do we need AI to assist paper writing?

- As Steven Pinker points in his article, many scientists are, in fact, bad writers:
  - "I know many scholars who have nothing to hide and no need to impress. They do groundbreaking work on important subjects, reason well about clear ideas, and are honest, down-to-earth people. Still, their writing stinks."
- One senior editor of Nature bluntly stated "most papers are badly written"

Pinker, S. (2014). Why academics stink at writing. The chronicle of higher education, 61(5), 2-9. Sage L (2003) Writing a clear and engaging paper for all astronomers. In Astronomy Communication (eds Heck A, Madsen C), pp 221–226. Dordrecht, the Netherlands: Kluwer

## Why do we need AI to assist paper writing?

- Generative AI tools have the potential to accelerate the research process by carrying out straightforward but time-consuming tasks — such as writing summaries and polish articles
- Generative AI tools have already been used to polish academic writing
  - Up to 17.5% of computer science papers on the arXiv have LLMs usage
- Title/Abstract generation can also provide research inspirations for scientists

#### **HOW AN AI CHATBOT EDITS A MANUSCRIPT**

Computational biologists Casey Greene and Milton Pividori constructed a program to get GPT-3 to edit research papers. The AI receives sections of a manuscript in turn, each with their own editing instructions (or 'prompts'). The prompts are adapted - in automated fashion - to suit the needs of each section, such as the methods or conclusion. In this example from their work, GPT-3 was fed an introduction section and the prompt below.

#### Prompt:

"Revise the following paragraph from the introduction of an academic paper (with the title 'An efficient not-only-linear correlation coefficient based on machine learning' and keywords 'correlation coefficient, nonlinear relationships, gene expression') so most of the citations to other academic papers are kept, the text minimizes the use of jargon, the text grammar is correct, spelling errors are fixed, and the text is in active voice and has a clear sentence structure."

#### Paragraph:

New technologies have vastly improved data collection, generating a deluge of information across different disciplines. This large amount of data provides new opportunities to address unanswered scientific questions, provided we have efficient tools capable of identifying multiple types of underlying patterns. Correlation analysis is an essential statistical technique for discovering relationships between variables [1]. Correlation coefficients are often used in exploratory data mining techniques, such as clustering or community detection algorithms, to compute a similarity value between a pair of objects of interest such as genes [2] or disease-relevant lifestyle factors [3]. Correlation methods are also used in supervised tasks, for example, for feature selection to improve prediction accuracy [4,5]. The Pearson correlation coefficient is ubiquitously deployed across application domains and diverse scientific areas. Thus, even minor and significant improvements in these techniques could have enormous consequences in industry and research

#### **GPT-3 edits:**

New technologies have vastly improved data collection, generating a deluge of information across different disciplines. The increasing availability of data has opened up new possibilities for scientific exploration. This large amount of data provides new opportunities to address unanswered scientific questions, provided we have efficient tools capable of identifying multiple types of underlying patterns. Correlation analysis is an essential a useful statistical technique for discovering to uncover such relationships between variables [1]. Correlation coefficients are often used in exploratory data mining techniques, such as clustering or community detection algorithms, to compute a calculate the similarity value between a pair of two objects of interest such as, like genes [2] or disease relevant lifestyle factors related to diseases [3]. Correlation methods They are also used in supervised tasks, for example, for like feature selection, to improve boost prediction accuracy [4,5]. The Pearson correlation coefficient is ubiquitously deployed widely used across many application domains and diverse scientific areas disciplines. Thus Therefore, even minor and significant small improvements in these techniques could this technique can have enormous consequences in a huge impact on industry and research.

Pinker, S. (2014). Why academics stink at writing. The chronicle of higher education, 61(5), 2-9.

https://www.nature.com/articles/d41586-023-00340-6 https://www.nature.com/articles/d41586-024-01042-3

. Liang, W., Zhang, Y., Wu, Z., Lepp, H., Ji, W., Zhao, X., ... & Zou, J. Y. (2024). Mapping the increasing use of Ilms in scientific papers. arXiv preprint arXiv:2404.01268.



## Task Formulation

## Paper components that benefit from existing systems

Title

**Attention Is All You Need** 

Ashish Vaswani\* Google Brain avaswani@google.com

Noam Shazeer\* Google Brain noam@google.com

Niki Parmar\* Google Research nikip@google.com

Jakob Uszkoreit\* Google Research usz@google.com

Llion Jones\* Google Research llion@google.com

Aidan N. Gomez\* † University of Toronto aidan@cs.toronto.edu

Łukasz Kaiser\* Google Brain lukaszkaiser@google.com

Illia Polosukhin\* ‡ illia.polosukhin@gmail.com

#### **Abstract**

#### Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 Englishto-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.8 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature. We show that the Transformer generalizes well to other tasks by applying it successfully to English constituency parsing both with large and limited training data.

#### Definition

In this work we propose the Transformer, a model architecture eschewing recurrence and instead relying entirely on an attention mechanism to draw global dependencies between input and output. The Transformer allows for significantly more parallelization and can reach a new state of the art in translation quality after being trained for as little as twelve hours on eight P100 GPUs.

#### 2 Background

The goal of reducing sequential computation also forms the foundation of the Extended Neural GPU [16], ByteNet [18] and ConvS2S [9], all of which use convolutional neural networks as basic building block, computing hidden representations in parallel for all input and output positions. In these models, the number of operations required to relate signals from two arbitrary input or output positions grows in the distance between positions, linearly for ConvS2S and logarithmically for ByteNet. This makes it more difficult to learn dependencies between distant positions [12]. In the Transformer this is reduced to a constant number of operations, albeit at the cost of reduced effective resolution due to averaging attention-weighted positions, an effect we counteract with Multi-Head Attention as described in section 3.2.

Self-attention, sometimes called intra-attention is an attention mechanism relating different positions of a single sequence in order to compute a representation of the sequence. Self-attention has been used successfully in a variety of tasks including reading comprehension, abstractive summarization textual entailment and learning task-independent sentence representations [4, 27, 28, 22].

End-to-end memory networks are based on a recurrent attention mechanism instead of sequencealigned recurrence and have been shown to perform well on simple-language question answering and language modeling tasks [34].

To the best of our knowledge, however, the Transformer is the first transduction model relying entirely on self-attention to compute representations of its input and output without using sequence aligned RNNs or convolution. In the following sections, we will describe the Transformer, motivate self-attention and discuss its advantages over models such as [17, 18] and [9].

**Related Work** 

Table 4: The Transformer generalizes well to English constituency parsing (Results are on Section 23 of WSD

Parser	Training	WSJ 23 F
Vinyals & Kaiser el al. (2014) [37]	WSJ only, discriminative	88.3
Petrov et al. (2006) [29]	WSJ only, discriminative	90.4
Zhu et al. (2013) [40]	WSJ only, discriminative	90.4
Dyer et al. (2016) [8]	WSJ only, discriminative	91.7
Transformer (4 layers)	WSJ only, discriminative	91.3
Zhu et al. (2013) [40]	semi-supervised	91.3
Huang & Harper (2009) [14]	semi-supervised	91.3
McClosky et al. (2006) [26]	semi-supervised	92.1
Vinyals & Kaiser el al. (2014) [37]	semi-supervised	92.1
Transformer (4 layers)	semi-supervised	92.7
Luong et al. (2015) [23]	multi-task	93.0
Dyer et al. (2016) [8]	generative	93.3

increased the maximum output length to input length + 300. We used a beam size of 21 and  $\alpha=0.3$ for both WSJ only and the semi-supervised setting.

Our results in Table 4 show that despite the lack of task-specific tuning our model performs surprisingly well, yielding better results than all previously reported models with the exception of the Recurrent Neural Network Grammar [8]

In contrast to RNN sequence-to-sequence models [37], the Transformer outperforms the Berkeley-Parser [29] even when training only on the WSJ training set of 40K sentences.

#### Conclusion

In this work, we presented the Transformer, the first sequence transduction model based entirely on attention, replacing the recurrent layers most commonly used in encoder-decoder architectures with multi-headed self-attention.

For translation tasks, the Transformer can be trained significantly faster than architectures based on recurrent or convolutional layers. On both WMT 2014 English-to-German and WMT 2014 English-to-French translation tasks, we achieve a new state of the art. In the former task our best model outperforms even all previously reported ensembles.

We are excited about the future of attention-based models and plan to apply them to other tasks. We plan to extend the Transformer to problems involving input and output modalities other than text and to investigate local, restricted attention mechanisms to efficiently handle large inputs and outputs such as images, audio and video. Making generation less sequential is another research goals of ours.

The code we used to train and evaluate our models is available at https://github.com/ tensorflow/tensor2tensor.

Acknowledgements We are grateful to Nal Kalchbrenner and Stephan Gouws for their fruitful comments, corrections and inspiration.

Table Description Claim Conclusion

## Task Formulation

## Paper components that benefit from existing systems

Title

**Attention Is All You Need** 

Ashish Vaswani\* Google Brain avaswani@google.com

Noam Shazeer\* Google Brain noam@google.com

Niki Parmar\* Google Research nikip@google.com

Jakob Uszkoreit\* Google Research usz@google.com

Llion Jones\* Google Research llion@google.com

Aidan N. Gomez\* † University of Toronto aidan@cs.toronto.edu

Łukasz Kaiser\* Google Brain lukaszkaiser@google.com

Illia Polosukhin\* ‡ illia.polosukhin@gmail.com

#### **Abstract**

#### Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 Englishto-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.8 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature. We show that the Transformer generalizes well to other tasks by applying it successfully to English constituency parsing both with large and limited training data.

#### Definition

In this work we propose the Transformer, a model architecture eschewing recurrence and instead relying entirely on an attention mechanism to draw global dependencies between input and output. The Transformer allows for significantly more parallelization and can reach a new state of the art in translation quality after being trained for as little as twelve hours on eight P100 GPUs.

#### 2 Background

The goal of reducing sequential computation also forms the foundation of the Extended Neural GPU [16], ByteNet [18] and ConvS2S [9], all of which use convolutional neural networks as basic building block, computing hidden representations in parallel for all input and output positions. In these models, the number of operations required to relate signals from two arbitrary input or output positions grows in the distance between positions, linearly for ConvS2S and logarithmically for ByteNet. This makes it more difficult to learn dependencies between distant positions [12]. In the Transformer this is reduced to a constant number of operations, albeit at the cost of reduced effective resolution due to averaging attention-weighted positions, an effect we counteract with Multi-Head Attention as described in section 3.2.

Self-attention, sometimes called intra-attention is an attention mechanism relating different positions of a single sequence in order to compute a representation of the sequence. Self-attention has been used successfully in a variety of tasks including reading comprehension, abstractive summarization textual entailment and learning task-independent sentence representations [4, 27, 28, 22].

End-to-end memory networks are based on a recurrent attention mechanism instead of sequencealigned recurrence and have been shown to perform well on simple-language question answering and language modeling tasks [34].

To the best of our knowledge, however, the Transformer is the first transduction model relying entirely on self-attention to compute representations of its input and output without using sequence aligned RNNs or convolution. In the following sections, we will describe the Transformer, motivate self-attention and discuss its advantages over models such as [17, 18] and [9].

**Related Work** 

Table 4: The Transformer generalizes well to English constituency parsing (Results are on Section 23 of WSD

Parser	Training	WSJ 23 F
Vinyals & Kaiser el al. (2014) [37]	WSJ only, discriminative	88.3
Petrov et al. (2006) [29]	WSJ only, discriminative	90.4
Zhu et al. (2013) [40]	WSJ only, discriminative	90.4
Dyer et al. (2016) [8]	WSJ only, discriminative	91.7
Transformer (4 layers)	WSJ only, discriminative	91.3
Zhu et al. (2013) [40]	semi-supervised	91.3
Huang & Harper (2009) [14]	semi-supervised	91.3
McClosky et al. (2006) [26]	semi-supervised	92.1
Vinyals & Kaiser el al. (2014) [37]	semi-supervised	92.1
Transformer (4 layers)	semi-supervised	92.7
Luong et al. (2015) [23]	multi-task	93.0
Dyer et al. (2016) [8]	generative	93.3

increased the maximum output length to input length + 300. We used a beam size of 21 and  $\alpha=0.3$ for both WSJ only and the semi-supervised setting.

Our results in Table 4 show that despite the lack of task-specific tuning our model performs surprisingly well, yielding better results than all previously reported models with the exception of the Recurrent Neural Network Grammar [8]

In contrast to RNN sequence-to-sequence models [37], the Transformer outperforms the Berkeley-Parser [29] even when training only on the WSJ training set of 40K sentences.

In this work, we presented the Transformer, the first sequence transduction model based entirely on attention, replacing the recurrent layers most commonly used in encoder-decoder architectures with multi-headed self-attention.

For translation tasks, the Transformer can be trained significantly faster than architectures based on recurrent or convolutional layers. On both WMT 2014 English-to-German and WMT 2014 English-to-French translation tasks, we achieve a new state of the art. In the former task our best model outperforms even all previously reported ensembles.

We are excited about the future of attention-based models and plan to apply them to other tasks. We plan to extend the Transformer to problems involving input and output modalities other than text and to investigate local, restricted attention mechanisms to efficiently handle large inputs and outputs such as images, audio and video. Making generation less sequential is another research goals of ours.

The code we used to train and evaluate our models is available at https://github.com/ tensorflow/tensor2tensor.

Acknowledgements We are grateful to Nal Kalchbrenner and Stephan Gouws for their fruitful comments, corrections and inspiration

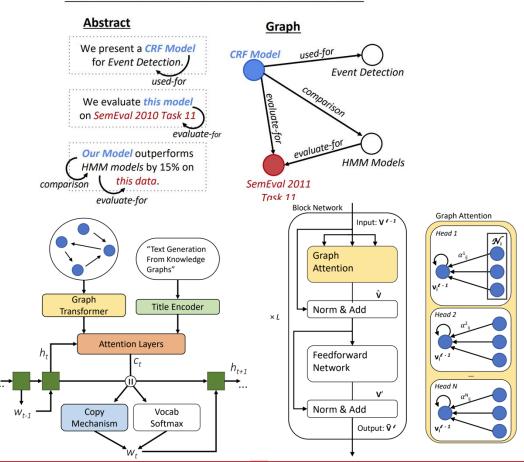
Table Description Claim Conclusion



## Paper Abstract Generation

- Generating paper abstracts based on titles and knowledge graphs
- Use graph attention layer for graph transformer

Title: Event Detection with Conditional Random Fields



### What if knowledge graphs are not given?

#### Pros:

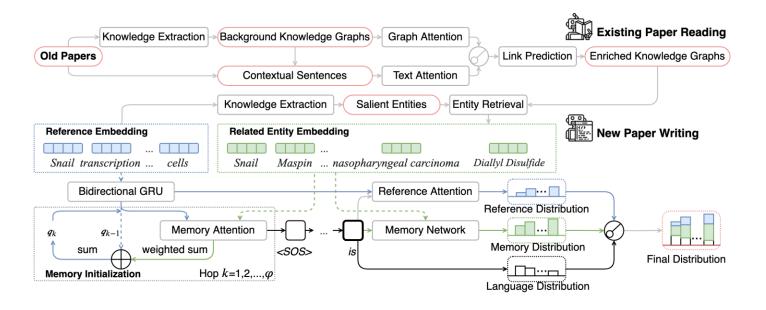
- Construct an AGENDA dataset for abstract generation task
- Convert disconnected labeled graph to connected unconnected graph for graph encoding

#### Points for Improvement:

Lack cross attention between the input knowledge graph and the title during each layer of the graph transformer block

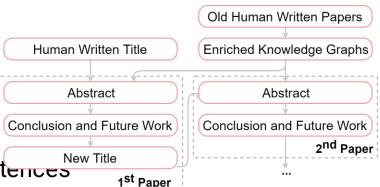
## PaperRobot Incremental Writing

- Existing Paper Reading
  - Construct KB based on existing papers with PubTator
- Create New Ideas
  - Predict new links by leveraging graph features and contextual sentences
- Write a New Paper Draft about New Ideas
  - Employ a memory network to capture multi-hop attention of related entities and utilize reference attention to capture reference text



#### Pros:

- Designs a pipeline for a paper drafting system capable of reading literature, proposing new ideas, and drafting papers
- Use memory network to filter out irrelevant entities Points for Improvement:
- Fail to fully utilize the graph structures of the knowledge graph when generating drafts
- The link predictions and paper draft generation are not optimized in an end-to-end way, leading to error propagation



PaperRobot Incremental Writing

Title	e Snail transcription factor negatively regulates maspin tumor suppressor in human prostate cancer cells			
Entities	Related: nasopharyngeal carcinoma; diallyl disulfide			
Output	Human (Neal et al., 2012)	System		
Abstract	Background: Maspin, a putative tumor suppressor that is down-regulated in breast and prostate cancer, has been associated with decreased cell motility. Snail transcription factor is a zinc finger protein that is increased in breast cancer and is associated with increased tumor motility and invasion by induction of epithelial-mesenchymal transition (EMT). We investigated the molecular mechanisms by which Snail increases tumor motility and invasion utilizing prostate cancer cells. Methods: Expression levels were analyzed by RT-PCR and western blot analyses. Cell motility and invasion assays were performed, while Snail regulation and binding to maspin promoter was analyzed by luciferase reporter and chromatin immunoprecipitation (ChIP) assays. Results: Snail protein expression was higher in different prostate cancer cells lines as compared to normal prostate epithelial cells.	Background: <b>Snail</b> is a multifunctional protein that plays an important role in the pathogenesis of <b>prostate cancer</b> . However, it has been shown to be associated with poor prognosis. The purpose of this study was to investigate the effect of negatively on the expression of <b>maspin</b> in <b>human nasopharyngeal carcinoma</b> cell lines. Methods: Quantitative real-time PCR and western blot analysis were used to determine whether the demethylating agent was investigated by quantitative <b>RT-PCR</b> (qRT-PCR) and <b>Western blotting</b> . Results showed that the binding protein plays a significant role in the regulation of <b>tumor</b> growth and progression.		
Conclusion and Future work	Collectively, our results indicate for the first time that <b>Snail</b> can negatively regulate <b>maspin</b> through direct promoter repression resulting in increased migration and invasion in <b>prostate cancer</b> cells. This study reveals a novel mechanism of how Snail may function and show the importance of therapeutic targeting of <b>Snail</b> signaling in future.	In summary, our study demonstrates that <b>Snail</b> negatively inhibited the expression of <b>Maspin</b> in <b>human nasopharyngeal carcinoma</b> cell lines and in vitro. Our results indicate that the combination of the demethylating agent might be a potential therapeutic target for the treatment of <b>prostate cancer</b> .		
New Title	Role of <b>maspin</b> in cancer (Berardi et al., 2013)	The role of <b>nasopharyngeal carcinoma</b> in the rat model of <b>prostate cancer</b> cells		

Wang, Q., Huang, L., Jiang, Z., Knight, K., Ji, H., Bansal, M., & Luan, Y. (2019). PaperRobot: Incremental draft generation of scientific ideas. ACL 2019.



## Task Formulation

## Paper components that benefit from existing systems

Title

**Attention Is All You Need** 

Ashish Vaswani\* Google Brain avaswani@google.com

Noam Shazeer\* Google Brain noam@google.com

Niki Parmar\* Google Research nikip@google.com

Jakob Uszkoreit\* Google Research usz@google.com

Llion Jones\* Google Research llion@google.com

Aidan N. Gomez\* † University of Toronto aidan@cs.toronto.edu

Łukasz Kaiser\* Google Brain lukaszkaiser@google.com

Illia Polosukhin\* ‡ illia.polosukhin@gmail.com

#### **Abstract**

#### Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 Englishto-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.8 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature. We show that the Transformer generalizes well to other tasks by applying it successfully to English constituency parsing both with large and limited training data.

#### Definition

In this work we propose the Transformer, a model architecture eschewing recurrence and instead relying entirely on an attention mechanism to draw global dependencies between input and output. translation quality after being trained for as little as twelve hours on eight P100 GPUs.

#### 2 Background

The goal of reducing sequential computation also forms the foundation of the Extended Neural GPU [16], ByteNet [18] and ConyS2S [9], all of which use convolutional neural networks as basic building block, computing hidden representations in parallel for all input and output positions. In these models, the number of operations required to relate signals from two arbitrary input or output positions grows in the distance between positions, linearly for ConvS2S and logarithmically for ByteNet. This makes it more difficult to learn dependencies between distant positions [12]. In the Transformer this is reduced to a constant number of operations, albeit at the cost of reduced effective resolution due to averaging attention-weighted positions, an effect we counteract with Multi-Head Attention as described in section 3.2.

Self-attention, sometimes called intra-attention is an attention mechanism relating different positions of a single sequence in order to compute a representation of the sequence. Self-attention has been used successfully in a variety of tasks including reading comprehension, abstractive summarization, textual entailment and learning task-independent sentence representations [4, 27, 28, 22].

End-to-end memory networks are based on a recurrent attention mechanism instead of sequencealigned recurrence and have been shown to perform well on simple-language question answering and language modeling tasks [34].

To the best of our knowledge, however, the Transformer is the first transduction model relying entirely on self-attention to compute representations of its input and output without using sequencealigned RNNs or convolution. In the following sections, we will describe the Transformer, motivate self-attention and discuss its advantages over models such as [17, 18] and [9].

#### **Related Work**

Table 4: The Transformer generalizes well to English constituency parsing (Results are on Section 23

Parser	Training	WSJ 23 F
Vinyals & Kaiser el al. (2014) [37]	WSJ only, discriminative	88.3
Petrov et al. (2006) [29]	WSJ only, discriminative	90.4
Zhu et al. (2013) [40]	WSJ only, discriminative	90.4
Dyer et al. (2016) [8]	WSJ only, discriminative	91.7
Transformer (4 layers)	WSJ only, discriminative	91.3
Zhu et al. (2013) [40]	semi-supervised	91.3
Huang & Harper (2009) [14]	semi-supervised	91.3
McClosky et al. (2006) [26]	semi-supervised	92.1
Vinyals & Kaiser el al. (2014) [37]	semi-supervised	92.1
Transformer (4 layers)	semi-supervised	92.7
Luong et al. (2015) [23]	multi-task	93.0
Dyer et al. (2016) [8]	generative	93.3

increased the maximum output length to input length + 300. We used a beam size of 21 and  $\alpha=0.3$ for both WSJ only and the semi-supervised setting.

Our results in Table 4 show that despite the lack of task-specific tuning our model performs surprisingly well, yielding better results than all previously reported models with the exception of the Recurrent Neural Network Grammar [8]

n contrast to RNN sequence-to-sequence models [37], the Transformer outperforms the Berkelev-Parser [29] even when training only on the WSJ training set of 40K sentences.

#### Conclusion

In this work, we presented the Transformer, the first sequence transduction model based entirely on attention, replacing the recurrent layers most commonly used in encoder-decoder architectures with multi-headed self-attention.

For translation tasks, the Transformer can be trained significantly faster than architectures based on recurrent or convolutional layers. On both WMT 2014 English-to-German and WMT 2014 English-to-French translation tasks, we achieve a new state of the art. In the former task our best model outperforms even all previously reported ensembles.

We are excited about the future of attention-based models and plan to apply them to other tasks. We plan to extend the Transformer to problems involving input and output modalities other than text and to investigate local, restricted attention mechanisms to efficiently handle large inputs and outputs such as images, audio and video. Making generation less sequential is another research goals of ours.

The code we used to train and evaluate our models is available at https://github.com/ tensorflow/tensor2tensor.

Acknowledgements We are grateful to Nal Kalchbrenner and Stephan Gouws for their fruitful comments, corrections and inspiration.

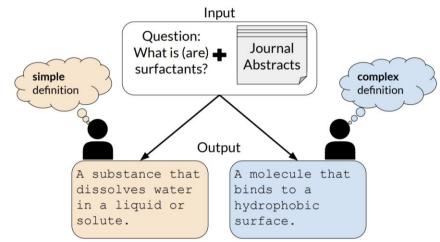
Description Claim Conclusion



## Scientific Definition Generation

- Generating definitions of scientific and medical terms with varying complexity
- Dataset: Medical Consumer Questions (Abacha et al., 2019), Wikipedia, S2ORC
- Method
  - Reranking BART generation results based on logits from a discriminator trained to distinguish scientific journal text from science news text (measuring complexity)

<b>Control Method</b>	Dire	ction	
	Low (News)	High (Journal)	
SVM-	A type of computing in which there are many	In computer science, distributed computing	Questi
Rerank	computers running at the same time in differ-	is the process of computing on a large scale	What is
	ent parts of the world.	without a single centralized data center.	surfacta
BERT-	A type of computer system in which there are	In computer science, distributed computing	simple
Rerank	more than a few computers working together.	is the process of computing on a large scale	definition
		without a single centralized data center.	
GeDi	Is the implementation of computer programs	In computer science, a concept that states	
	across multiple computers on similar hard-	that data must be shared across comput-	
	ware and/or software resources.	ing resources.	
DExpert	An Internet-driven by-computing that por-	In computer science and communication	A substance th
•	tion of different computers from start to fin-	between-Consequently-integrates.	dissolves water
	ish.	1 , 2	in a liquid or
PPLM	Easeless, self-organized, and often self-	Multi-purpose, distributed system software	solute.
	organizing networked computer systems in-	with or without a single datum storage sys-	Soluce.
	tended for the purposes of optimization.	tem.	



#### Points for Improvement:

- Instead of reranking, inject complexity as an additional input feature in the input
- Employ a generative adversarial framework to optimize the BART generator and reranking discriminator jointly



## Task Formulation

## Writing each part of a paper is a different task!

#### Title

**Attention Is All You Need** 

Ashish Vaswani\* Google Brain avaswani@google.com

Noam Shazeer\* Google Brain noam@google.com

Niki Parmar\* Google Research nikip@google.com

Jakob Uszkoreit\* Google Research usz@google.com

Llion Jones\* Google Research llion@google.com

Aidan N. Gomez\* † University of Toronto aidan@cs.toronto.edu

Łukasz Kaiser\* Google Brain lukaszkaiser@google.com

Illia Polosukhin\* ‡ illia.polosukhin@gmail.com

#### **Abstract**

#### Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 Englishto-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.8 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature. We show that the Transformer generalizes well to other tasks by applying it successfully to English constituency parsing both with large and limited training data.

#### Definition

In this work we propose the Transformer, a model architecture eschewing recurrence and instead relying entirely on an attention mechanism to draw global dependencies between input and output. The Transformer allows for significantly more parallelization and can reach a new state of the artranslation quality after being trained for as little as twelve hours on eight P100 GPUs.

#### 2 Background

The goal of reducing sequential computation also forms the foundation of the Extended Neural GPU [16], ByteNet [18] and ConvS2S [9], all of which use convolutional neural networks as basic building block, computing hidden representations in parallel for all input and output positions. In these models the number of operations required to relate signals from two arbitrary input or output positions grow in the distance between positions, linearly for ConvS2S and logarithmically for ByteNet. This makes it more difficult to learn dependencies between distant positions [12]. In the Transformer this is reduced to a constant number of operations, albeit at the cost of reduced effective resolution due to averaging attention-weighted positions, an effect we counteract with Multi-Head Attention as

Self-attention, sometimes called intra-attention is an attention mechanism relating different positions of a single sequence in order to compute a representation of the sequence. Self-attention has been used successfully in a variety of tasks including reading comprehension, abstractive summarization, textual entailment and learning task-independent sentence representations [4, 27, 28, 22].

End-to-end memory networks are based on a recurrent attention mechanism instead of sequence aligned recurrence and have been shown to perform well on simple-language question answering and language modeling tasks [34].

To the best of our knowledge, however, the Transformer is the first transduction model relying entirely on self-attention to compute representations of its input and output without using sequence aligned RNNs or convolution. In the following sections, we will describe the Transformer, motivate self-attention and discuss its advantages over models such as [17, 18] and [9].

#### **Related Work**

Table 4: The Transformer generalizes well to English constituency parsing (Results are on Section 23 of WSJ)

Parser	Training	WSJ 23 F1	
Vinyals & Kaiser el al. (2014) [37]	WSJ only, discriminative	88.3	
Petrov et al. (2006) [29]	WSJ only, discriminative	90.4	
Zhu et al. (2013) [40]	WSJ only, discriminative	90.4	
Dyer et al. (2016) [8]	WSJ only, discriminative	91.7	
Transformer (4 layers)	WSJ only, discriminative	91.3	
Zhu et al. (2013) [40]	semi-supervised	91.3	
Huang & Harper (2009) [14]	semi-supervised	91.3	
McClosky et al. (2006) [26]	semi-supervised	92.1	
Vinyals & Kaiser el al. (2014) [37]	semi-supervised	92.1	
Transformer (4 layers)	semi-supervised	92.7	
Luong et al. (2015) [23]	multi-task	93.0	
Dyer et al. (2016) [8]	generative	93.3	

increased the maximum output length to input length + 300. We used a beam size of 21 and  $\alpha=0.3$ for both WSJ only and the semi-supervised setting.

Our results in Table 4 show that despite the lack of task-specific tuning our model performs surprisingly well, yielding better results than all previously reported models with the exception of the Recurrent Neural Network Grammar [8]

n contrast to RNN sequence-to-sequence models [37], the Transformer outperforms the Berkelev-Parser [29] even when training only on the WSJ training set of 40K sentences.

#### Conclusion

In this work, we presented the Transformer, the first sequence transduction model based entirely on attention, replacing the recurrent layers most commonly used in encoder-decoder architectures with multi-headed self-attention.

For translation tasks, the Transformer can be trained significantly faster than architectures based on recurrent or convolutional layers. On both WMT 2014 English-to-German and WMT 2014 English-to-French translation tasks, we achieve a new state of the art. In the former task our best model outperforms even all previously reported ensembles.

We are excited about the future of attention-based models and plan to apply them to other tasks. We plan to extend the Transformer to problems involving input and output modalities other than text and to investigate local, restricted attention mechanisms to efficiently handle large inputs and outputs such as images, audio and video. Making generation less sequential is another research goals of ours.

The code we used to train and evaluate our models is available at https://github.com/ tensorflow/tensor2tensor.

Acknowledgements We are grateful to Nal Kalchbrenner and Stephan Gouws for their fruitful comments, corrections and inspiration.

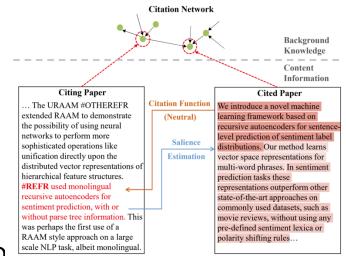
Description Claim Conclusion

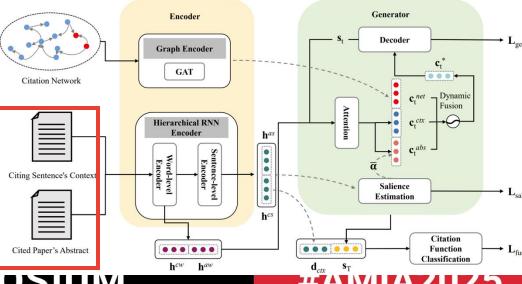
## Ge, Y., Dinh, L., Liu, X., Su, J., Lu, Z., Wang, A., & Diesner, J. (2021, August). BACO: A background knowledge-and content-based framework for citing sentence generation. ACL 2021 Knowledge-guided Citation Generation

- Input
  - Background knowledge: structural information from a citation network
  - Content: In-depth content information
- Target: Citation sentences
- Model Highlights
  - Salience Estimation to identify what information needed to be concentrated for the generation (what to cite)

Citation Function Classification to classify the reason for citation (why to cite)

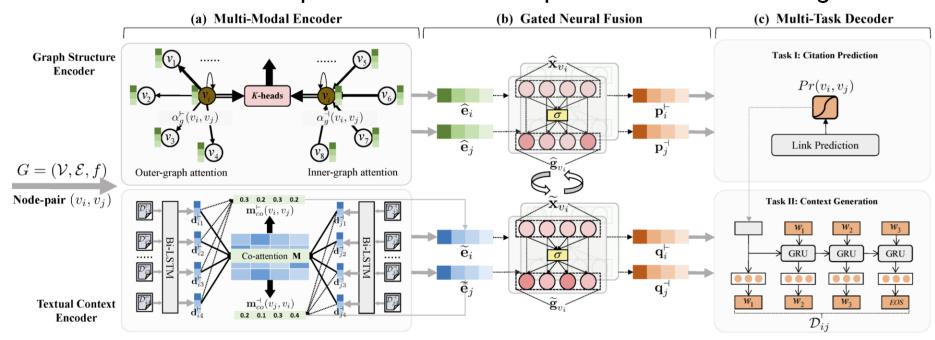
What if citation relations are not given?





## Wang, Q., Xiong, Y., Zhang, Y., Zhang, J., & Zhu, Y. (2021). AutoCite: Multi-Modal Representation Fusion for Contextual Citation Generation. WSDM 2021 Context-aware Citation Prediction

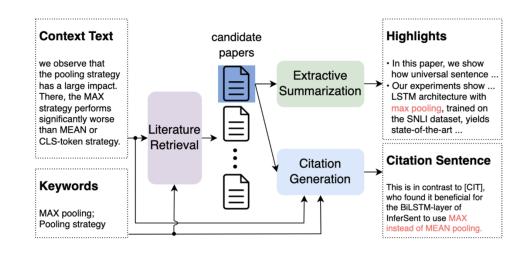
- Model
  - Graph Structure Encoder with dual-role attention to capture directed graph structure information
  - Gated Neural Fusion to integrate multi-modal representations and control features transfer for downstream tasks differently
  - Multi-Task Decoder to perform citation link prediction and context generation

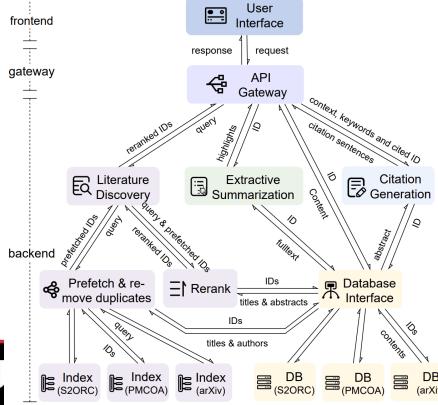


Can we combine information retrieval and citation recommendation?

## Joint Scientific Literature Discovery, Summarization and Citation Generation

- A pipeline that automatically recommends relevant papers, extracts highlights, and suggests a reference sentence as a citation of a paper, taking into consideration the userprovided context and keywords
  - Prefetching papers from database and reranking them based on SciBERT
  - Iteratively score and select sentences for summarization with MemSum
  - Concatenate keywords, context and target paper's abstract to generate citation sentences by Τ5



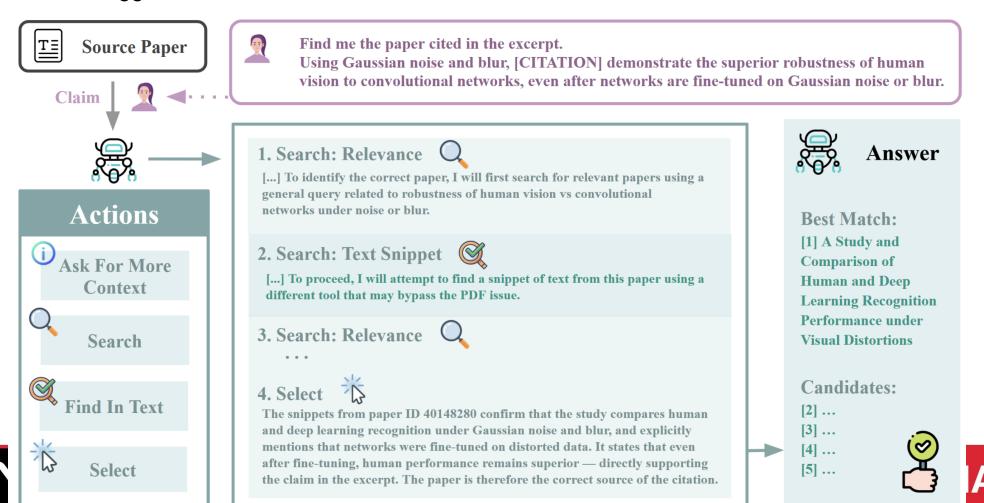


## Summary

Method	Pros	Cons
BACO (Ge et al., 2021)	<ul> <li>Construct a BACO dataset specifically for citation generation task</li> <li>The framework identifies the cited sentence and the reason for citing, jointly with the generation of citations</li> </ul>	Rely on provided citation relations to generate citation sentence
AutoCite(Wang et al., 2021)	<ul> <li>Propose a joint framework to predict citation links and generate citation sentences simultaneously</li> <li>Design a graph encoder to distinguish a paper cites others or it is cited by others</li> </ul>	■ Fail to generate different types of citation texts according to the contexts and positions in the target paper
SciLit (Gu et al., 2023)	Propose a pipeline to recommends relevant papers, extract highlights, and suggest a reference sentence	■The citation generation ignores the graph structure of citation networks

## CiteGuard: Faithful Citation Attribution

- A retrieval-aware agent framework to provide faithful grounding for citation validation
  - Can proactively query for additional claim context when the initial snippet is insufficient
  - Allow searching directly within paper contents, reducing reliance on PDF availability
  - Can suggest alternative citations



## Task Formulation

## Paper components that benefit from existing systems

Title

**Attention Is All You Need** 

Ashish Vaswani\* Google Brain avaswani@google.com

Noam Shazeer\* Google Brain noam@google.com

Niki Parmar\* Google Research nikip@google.com

Jakob Uszkoreit\* Google Research usz@google.com

Llion Jones\* Google Research llion@google.com

Aidan N. Gomez\* † University of Toronto aidan@cs.toronto.edu

Łukasz Kaiser\* Google Brain lukaszkaiser@google.com

Illia Polosukhin\* ‡ illia.polosukhin@gmail.com

#### **Abstract**

#### Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 Englishto-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.8 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature. We show that the Transformer generalizes well to other tasks by applying it successfully to English constituency parsing both with large and limited training data.

#### Definition

In this work we propose the Transformer, a model architecture eschewing recurrence and instead relying entirely on an attention mechanism to draw global dependencies between input and output. The Transformer allows for significantly more parallelization and can reach a new state of the art in translation quality after being trained for as little as twelve hours on eight P100 GPUs.

#### 2 Background

The goal of reducing sequential computation also forms the foundation of the Extended Neural GPU [16], ByteNet [18] and ConvS2S [9], all of which use convolutional neural networks as basic building block, computing hidden representations in parallel for all input and output positions. In these models, the number of operations required to relate signals from two arbitrary input or output positions grows in the distance between positions, linearly for ConvS2S and logarithmically for ByteNet. This makes it more difficult to learn dependencies between distant positions [12]. In the Transformer this is reduced to a constant number of operations, albeit at the cost of reduced effective resolution due to averaging attention-weighted positions, an effect we counteract with Multi-Head Attention as described in section 3.2.

Self-attention, sometimes called intra-attention is an attention mechanism relating different positions of a single sequence in order to compute a representation of the sequence. Self-attention has been used successfully in a variety of tasks including reading comprehension, abstractive summarization, textual entailment and learning task-independent sentence representations [4, 27, 28, 22].

End-to-end memory networks are based on a recurrent attention mechanism instead of sequencealigned recurrence and have been shown to perform well on simple-language question answering and language modeling tasks [34].

To the best of our knowledge, however, the Transformer is the first transduction model relying entirely on self-attention to compute representations of its input and output without using sequencealigned RNNs or convolution. In the following sections, we will describe the Transformer, motivate self-attention and discuss its advantages over models such as [17, 18] and [9].

#### **Related Work**

Table 4: The Transformer generalizes well to English constituency parsing (Results are on Section 23 of WSD

Parser	Training	WSJ 23 F
Vinyals & Kaiser el al. (2014) [37]	WSJ only, discriminative	88.3
Petrov et al. (2006) [29]	WSJ only, discriminative	90.4
Zhu et al. (2013) [40]	WSJ only, discriminative	90.4
Dyer et al. (2016) [8]	WSJ only, discriminative	91.7
Transformer (4 layers)	WSJ only, discriminative	91.3
Zhu et al. (2013) [40]	semi-supervised	91.3
Huang & Harper (2009) [14]	semi-supervised	91.3
McClosky et al. (2006) [26]	semi-supervised	92.1
Vinyals & Kaiser el al. (2014) [37]	semi-supervised	92.1
Transformer (4 layers)	semi-supervised	92.7
Luong et al. (2015) [23]	multi-task	93.0
Dyer et al. (2016) [8]	generative	93.3

increased the maximum output length to input length + 300. We used a beam size of 21 and  $\alpha=0.3$ for both WSJ only and the semi-supervised setting

Our results in Table 4 show that despite the lack of task-specific tuning our model performs surprisingly well, yielding better results than all previously reported models with the exception of the

In contrast to RNN sequence-to-sequence models [37], the Transformer outperforms the Berkeley Parser [29] even when training only on the WSJ training set of 40K sentences.

#### Conclusion

In this work, we presented the Transformer, the first sequence transduction model based entirely on attention, replacing the recurrent layers most commonly used in encoder-decoder architectures with multi-headed self-attention.

For translation tasks, the Transformer can be trained significantly faster than architectures based on recurrent or convolutional layers. On both WMT 2014 English-to-German and WMT 2014 English-to-French translation tasks, we achieve a new state of the art. In the former task our best model outperforms even all previously reported ensembles.

We are excited about the future of attention-based models and plan to apply them to other tasks. We plan to extend the Transformer to problems involving input and output modalities other than text and to investigate local, restricted attention mechanisms to efficiently handle large inputs and outputs such as images, audio and video. Making generation less sequential is another research goals of ours.

The code we used to train and evaluate our models is available at https://github.com/ tensorflow/tensor2tensor.

Acknowledgements We are grateful to Nal Kalchbrenner and Stephan Gouws for their fruitful comments, corrections and inspiration.

Description Claim Conclusion

## Scientific Context-Aware Description Generation

- Input
  - A paper object (e.g., figure, table, algorithm, theorem)
  - Context (Paragraphs before the descriptions)
- Target
  - Descriptions
- Dataset

Input	#num	#avg_out_len	#cand
Table	136K	74.05	199.93
Figure (chart/bar)	155K	76.60	179.35
Algorithm	56K	67.94	227.65
Theorem	175K	65.00	192.44

#### **Body text (Context)**

... languages: Telugu (te) and Turkish (tr)... Turkish (tr) vocabulary has been censored to contain no overlap with the Telugu... we evaluate these models using a recall@k metric defined as ...

#### **Table**

Result	te+en	+ tr	% Change	
Recall@1	17.0	17.6	+3.5%	
Recall@10	23.9	25.0	+4.6%	
Recall@20	26.3	27.7	+5.3%	

#### **Generated description w/o context (table only)**

Table shows when te+en is replaced with tr, the effect of different change is very small, although the performance of tr method gets really strong.

#### **Generated description w/ context (body text + table)**

Table summarizes the recall@1 measures and the percentage of the incremental improvement across languages for both tasks. The average incremental improvement across languages is about 4% in these cases, despite there being no overlap between in Telugu and Turkish.

## Task Formulation

## Paper components that benefit from existing systems

Title

**Attention Is All You Need** 

Ashish Vaswani\* Google Brain avaswani@google.com

Noam Shazeer\* Google Brain noam@google.com

Niki Parmar\* Google Research nikip@google.com

Jakob Uszkoreit\* Google Research usz@google.com

Llion Jones\* Google Research llion@google.com

Aidan N. Gomez\* † University of Toronto aidan@cs.toronto.edu

Łukasz Kaiser\* Google Brain lukaszkaiser@google.com

Illia Polosukhin\* ‡ illia.polosukhin@gmail.com

#### **Abstract**

#### Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 Englishto-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.8 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature. We show that the Transformer generalizes well to other tasks by applying it successfully to English constituency parsing both with large and limited training data.

#### Definition

In this work we propose the Transformer, a model architecture eschewing recurrence and instead relying entirely on an attention mechanism to draw global dependencies between input and output. The Transformer allows for significantly more parallelization and can reach a new state of the art in translation quality after being trained for as little as twelve hours on eight P100 GPUs.

#### 2 Background

The goal of reducing sequential computation also forms the foundation of the Extended Neural GPU [16], ByteNet [18] and ConvS2S [9], all of which use convolutional neural networks as basic building block, computing hidden representations in parallel for all input and output positions. In these models, the number of operations required to relate signals from two arbitrary input or output positions grows in the distance between positions, linearly for ConvS2S and logarithmically for ByteNet. This makes it more difficult to learn dependencies between distant positions [12]. In the Transformer this is reduced to a constant number of operations, albeit at the cost of reduced effective resolution due to averaging attention-weighted positions, an effect we counteract with Multi-Head Attention as described in section 3.2.

Self-attention, sometimes called intra-attention is an attention mechanism relating different positions of a single sequence in order to compute a representation of the sequence. Self-attention has been used successfully in a variety of tasks including reading comprehension, abstractive summarization, textual entailment and learning task-independent sentence representations [4, 27, 28, 22].

End-to-end memory networks are based on a recurrent attention mechanism instead of sequencealigned recurrence and have been shown to perform well on simple-language question answering and language modeling tasks [34].

To the best of our knowledge, however, the Transformer is the first transduction model relying entirely on self-attention to compute representations of its input and output without using sequencealigned RNNs or convolution. In the following sections, we will describe the Transformer, motivate self-attention and discuss its advantages over models such as [17, 18] and [9]

**Related Work** 

Table 4: The Transformer generalizes well to English constituency parsing (Results are on Section 23 of WSD

Parser	Training	WSJ 23 F1
Vinyals & Kaiser el al. (2014) [37]	WSJ only, discriminative	88.3
Petrov et al. (2006) [29]	WSJ only, discriminative	90.4
Zhu et al. (2013) [40]	WSJ only, discriminative	90.4
Dyer et al. (2016) [8]	WSJ only, discriminative	91.7
Transformer (4 layers)	WSJ only, discriminative	91.3
Zhu et al. (2013) [40]	semi-supervised	91.3
Huang & Harper (2009) [14]	semi-supervised	91.3
McClosky et al. (2006) [26]	semi-supervised	92.1
Vinyals & Kaiser el al. (2014) [37]	semi-supervised	92.1
Transformer (4 layers)	semi-supervised	92.7
Luong et al. (2015) [23]	multi-task	93.0
Dyer et al. (2016) [8]	generative	93.3

increased the maximum output length to input length + 300. We used a beam size of 21 and  $\alpha=0.3$ for both WSJ only and the semi-supervised setting.

Our results in Table 4 show that despite the lack of task-specific tuning our model performs surprisingly well, yielding better results than all previously reported models with the exception of the Recurrent Neural Network Grammar [8]

n contrast to RNN sequence-to-sequence models [37], the Transformer outperforms the Berkeley Parser [29] even when training only on the WSJ training set of 40K sentences.

#### Conclusion

In this work, we presented the Transformer, the first sequence transduction model based entirely on attention, replacing the recurrent layers most commonly used in encoder-decoder architectures with multi-headed self-attention.

For translation tasks, the Transformer can be trained significantly faster than architectures based on recurrent or convolutional layers. On both WMT 2014 English-to-German and WMT 2014 English-to-French translation tasks, we achieve a new state of the art. In the former task our best model outperforms even all previously reported ensembles.

We are excited about the future of attention-based models and plan to apply them to other tasks. We plan to extend the Transformer to problems involving input and output modalities other than text and to investigate local, restricted attention mechanisms to efficiently handle large inputs and outputs such as images, audio and video. Making generation less sequential is another research goals of ours.

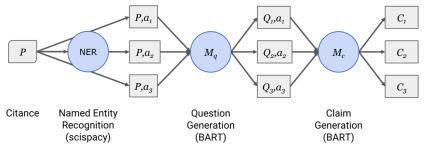
The code we used to train and evaluate our models is available at https://github.com/ tensorflow/tensor2tensor.

Acknowledgements We are grateful to Nal Kalchbrenner and Stephan Gouws for their fruitful comments, corrections and inspiration

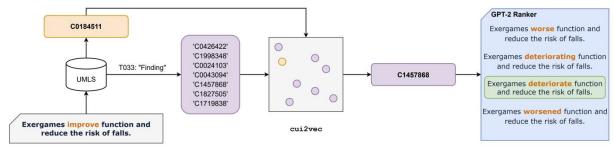
Description Claim Conclusion

## Scientific Claim Generation

- Supported Claims Generation
  - ClaimBART: Generating claim directly based on Citance
  - ClaimEntity



- Refuted Claims Generation
  - Knowledge-based Informed Negation



(1) ALS is the most common adult motor neuron disease with an incidence of 2 per 100,000 and prevalence of 5.4 per 100,000 individuals. (2) Current treatment options are based on symptom management and respiratory support with the only approved medications in widespread use, Riluzole and Edaravone, providing only modest benefits and only in some patients.

Current treatment options for ALS are based on symptom  $c_1$ management and respiratory support Riluzole is an approved ALS medication in widespread use  $c_2$ 

Edaravone is an approved ALS medication in widespread use  $c_3$ 

Riluzole and Edaravone are the only approved ALS medications in widespread use

Riluzole provides modest benefits in only some ALS patients

Edaravone provides modest benefits in only some ALS patients

#### Points for Improvement:

Instead of only considering entities in the claim, explore to use relations or even a subgraph of the knowledge graph when generating claims

## Task Formulation

Can we generate paper end-to-end?

Title

**Attention Is All You Need** 

Ashish Vaswani\* Google Brain avaswani@google.com

Noam Shazeer\* Google Brain noam@google.com

Niki Parmar\* Google Research nikip@google.com

Jakob Uszkoreit\* Google Research usz@google.com

Llion Jones\* Google Research llion@google.com

Aidan N. Gomez\* † University of Toronto aidan@cs.toronto.edu

Łukasz Kaiser\* Google Brain lukaszkaiser@google.com

Illia Polosukhin\* ‡ illia.polosukhin@gmail.com

#### **Abstract**

#### Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 Englishto-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.8 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature. We show that the Transformer generalizes well to other tasks by applying it successfully to English constituency parsing both with large and limited training data.

#### Definition

In this work we propose the Transformer, a model architecture eschewing recurrence and instead relying entirely on an attention mechanism to draw global dependencies between input and output. The Transformer allows for significantly more parallelization and can reach a new state of the art in translation quality after being trained for as little as twelve hours on eight P100 GPUs.

#### 2 Background

The goal of reducing sequential computation also forms the foundation of the Extended Neural GPU [16], ByteNet [18] and ConvS2S [9], all of which use convolutional neural networks as basic building block, computing hidden representations in parallel for all input and output positions. In these models, the number of operations required to relate signals from two arbitrary input or output positions grows in the distance between positions, linearly for ConvS2S and logarithmically for ByteNet. This makes it more difficult to learn dependencies between distant positions [12]. In the Transformer this is reduced to a constant number of operations, albeit at the cost of reduced effective resolution due to averaging attention-weighted positions, an effect we counteract with Multi-Head Attention as described in section 3.2.

Self-attention, sometimes called intra-attention is an attention mechanism relating different positions of a single sequence in order to compute a representation of the sequence. Self-attention has been used successfully in a variety of tasks including reading comprehension, abstractive summarization textual entailment and learning task-independent sentence representations [4, 27, 28, 22].

End-to-end memory networks are based on a recurrent attention mechanism instead of sequencealigned recurrence and have been shown to perform well on simple-language question answering and language modeling tasks [34].

To the best of our knowledge, however, the Transformer is the first transduction model relying entirely on self-attention to compute representations of its input and output without using sequence aligned RNNs or convolution. In the following sections, we will describe the Transformer, motivate self-attention and discuss its advantages over models such as [17, 18] and [9].

**Related Work** 

Table 4: The Transformer generalizes well to English constituency parsing (Results are on Section 23

Parser	Training	WSJ 23 F1
Vinyals & Kaiser el al. (2014) [37]	WSJ only, discriminative	88.3
Petrov et al. (2006) [29]	WSJ only, discriminative	90.4
Zhu et al. (2013) [40]	WSJ only, discriminative	90.4
Dyer et al. (2016) [8]	WSJ only, discriminative	91.7
Transformer (4 layers)	WSJ only, discriminative	91.3
Zhu et al. (2013) [40]	semi-supervised	91.3
Huang & Harper (2009) [14]	semi-supervised	91.3
McClosky et al. (2006) [26]	semi-supervised	92.1
Vinyals & Kaiser el al. (2014) [37]	semi-supervised	92.1
Transformer (4 layers)	semi-supervised	92.7
Luong et al. (2015) [23]	multi-task	93.0
Dyer et al. (2016) [8]	generative	93.3

increased the maximum output length to input length + 300. We used a beam size of 21 and  $\alpha=0.3$ for both WSJ only and the semi-supervised setting.

Our results in Table 4 show that despite the lack of task-specific tuning our model performs surprisingly well, yielding better results than all previously reported models with the exception of the Recurrent Neural Network Grammar [8]

In contrast to RNN sequence-to-sequence models [37], the Transformer outperforms the Berkeley-Parser [29] even when training only on the WSJ training set of 40K sentences.

#### Conclusion

In this work, we presented the Transformer, the first sequence transduction model based entirely on attention, replacing the recurrent layers most commonly used in encoder-decoder architectures with multi-headed self-attention.

For translation tasks, the Transformer can be trained significantly faster than architectures based on recurrent or convolutional layers. On both WMT 2014 English-to-German and WMT 2014 English-to-French translation tasks, we achieve a new state of the art. In the former task our best model outperforms even all previously reported ensembles.

We are excited about the future of attention-based models and plan to apply them to other tasks. We plan to extend the Transformer to problems involving input and output modalities other than text and to investigate local, restricted attention mechanisms to efficiently handle large inputs and outputs such as images, audio and video. Making generation less sequential is another research goals of ours.

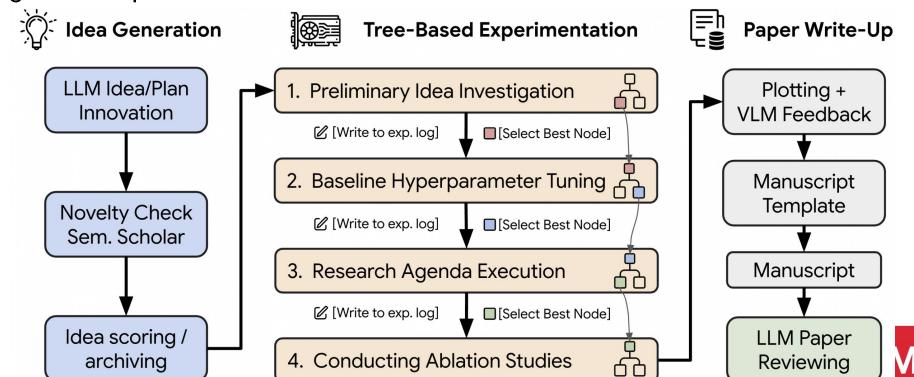
The code we used to train and evaluate our models is available at https://github.com/ tensorflow/tensor2tensor.

Acknowledgements We are grateful to Nal Kalchbrenner and Stephan Gouws for their fruitful comments, corrections and inspiration.

Table Description Claim Conclusion

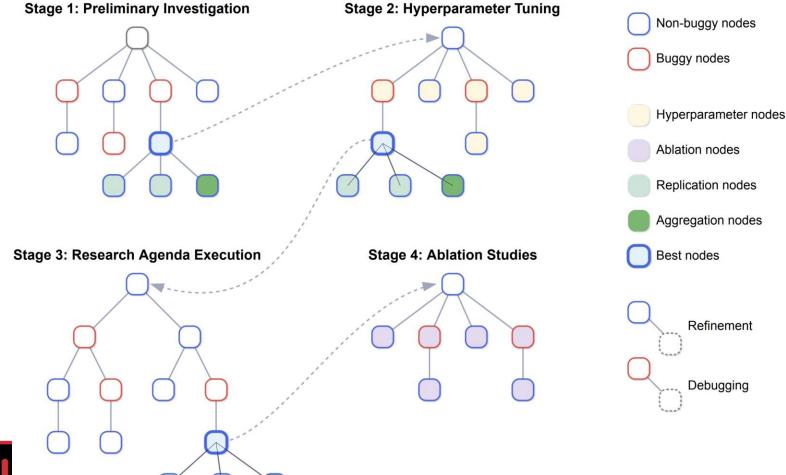
## The AIScientist-v2: Workshop-Level Automated Scientific Discovery via Agentic Tree Search

- Iteratively formulate scientific hypotheses, designs and executes experiments, analyzes and visualizes data, and autonomously authors scientific manuscripts
- Integrate vision-language model feedback for iterative refinement
- Al-generated manuscript can successfully pass peer review at a recognized machine learning workshop



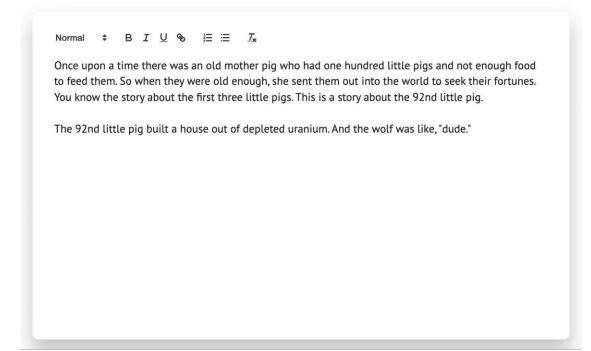
## The AlScientist-v2: Workshop-Level Automated Scientific Discovery via Agentic Tree Search

Integrating tree search with LLM-driven for parallel code generation



## Human-Al Collaborative Writing

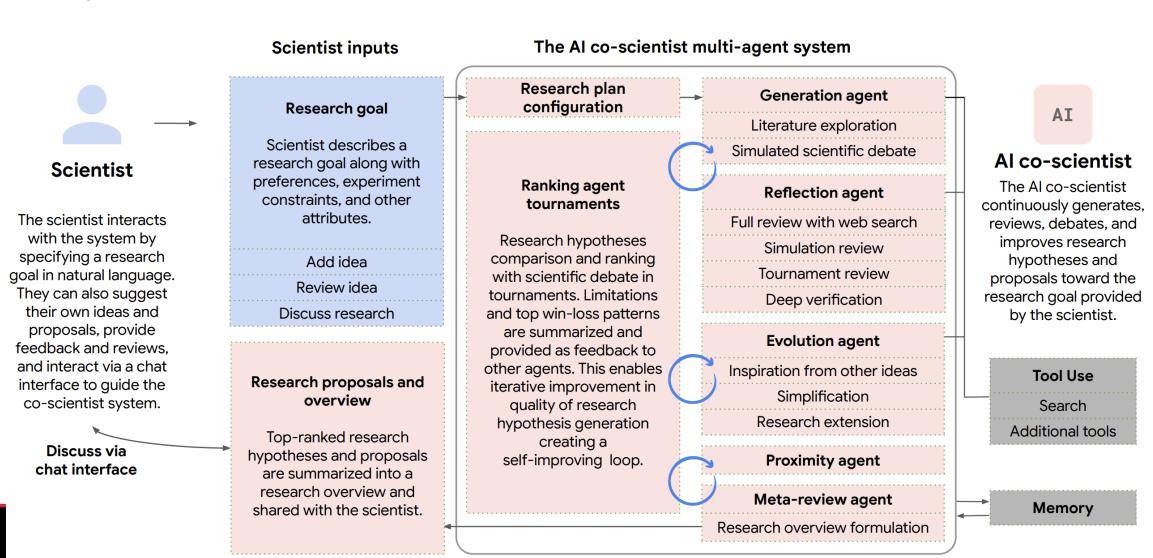
- Compared to either humans or AI, collaborative Human-Al Writing has:
  - Fewer grammatical errors
  - Higher vocabulary diversity
- LLMs can increase writers' productivity by providing writing suggestions
- Keep fraction of text written by writers to text written by GPT-3 relatively high can increase writers' feeling of ownership





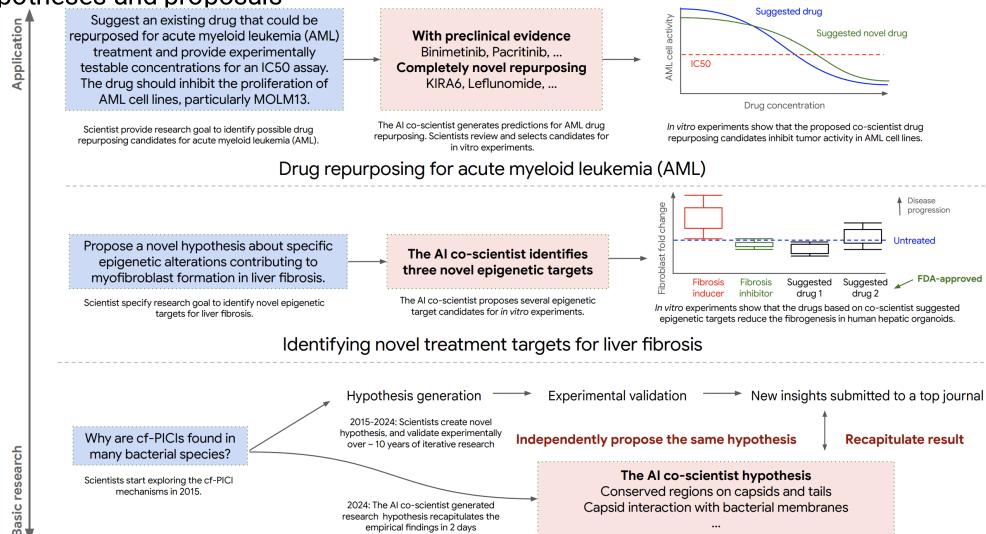
## Al co-scientist

 Given a research goal in natural language, the co-scientist generates novel research hypotheses and proposals



## Al co-scientist

 Given a research goal in natural language, the co-scientist generates novel research hypotheses and proposals



## Conclusion and Future Directions for Paper Writing

### Interesting points

- Al-assisted paper draft writing systems are categorized under knowledge-augmented LLMs
- These systems extensively utilize both graph and text features for generating results.
- To reduce error propagation, ongoing developments focus on creating an end-to-end system, moving away from the previous pipeline-based approach

#### Future directions

- Enhancing the multi-hop reasoning capabilities of these systems is also crucial
- Incorporating human feedback to iteratively improve system results presents an interesting future direction
- The systems use of only open published literature means it likely has limited access to negative experimental results or records of failed experiments



## Radiology Report Generation

- Radiology reporting is time-consuming
- AI can help generate or retrieve reports to assist radiologists
- Traditional Approach:
  - CNN + RNN/Transformer

## Enhanced Contrastive Learning with Multi-view Longitudinal Data (MLRG)

 Contrastive pre-training aligns current and previous images (longitudinal) with text; Select images with different angle (multiview)

Propose a tokenized absence encoding technique to handle missing patient-specific prior

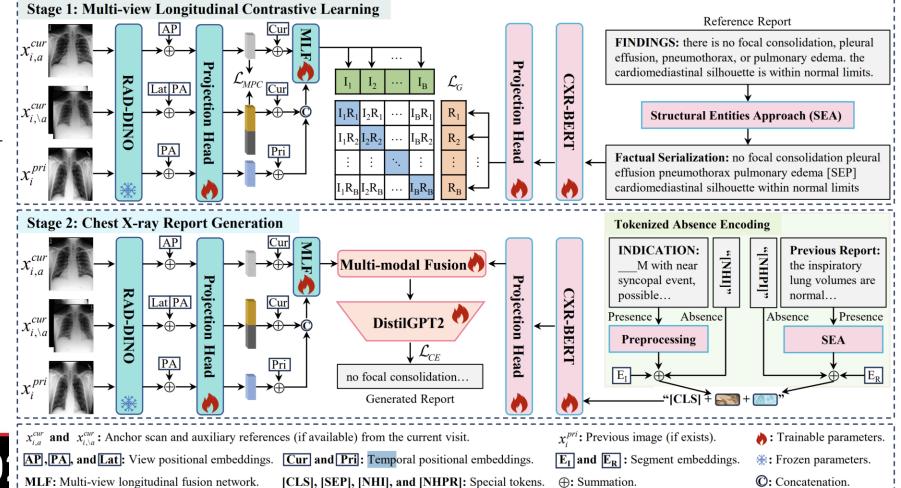
knowledge

#### Pros:

- Leverage richer data (multiple views + history)
- Use contrastive learning to improve alignment between images and text
- Robust in real scenarios with inconsistent history

#### Points for Improvement:

- Require prior exams and multi-view images to achieve best performance
- GPT-like Text generation can induce hallucination





## RRG-Mamba: Efficient Radiology Report Generation with State Space Model.

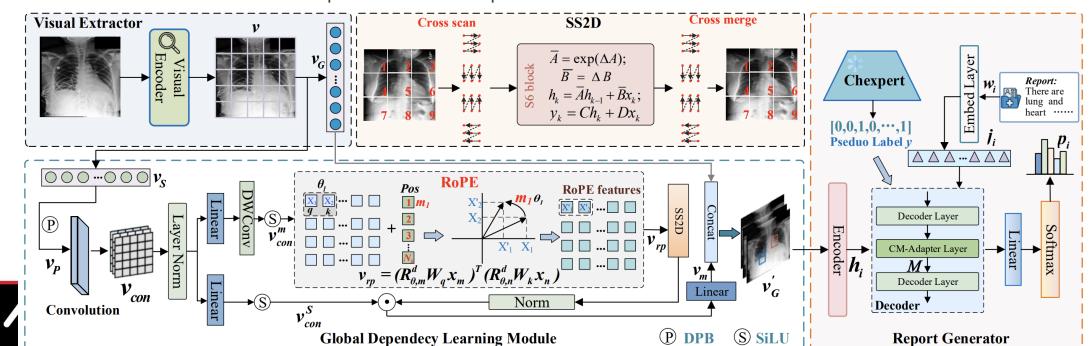
- Generate reports with lower computational cost without losing modeling of long-range dependencies
- Design a global dependency learning module that integrates rotary position encoding

#### Pros:

- Linear complexity decoding means it scales to longer reports or higher detail without slow-down
- Match or even exceed accuracy of Transformer-based methods
- Use position encoding to addresses both local and long-range feature relationships

#### Points for Improvement:

- No domain-specific knowledge integration
- Fewer off-the-shelf tools for State Space Models compared to Transformers



14324-14333)

## **Potential Future Directions**

- Combine both papers to create an efficient model that takes multi-view inputs
- Retrieve clinically similar reports based on visual input (Yan et al., 2024) and rewrite them based on current images
- Improve evaluation metrics
  - Integrating vision—language models (VLMs) as judge
  - Use attention map and reasoning chain to assess interpretability and factual consistency

## Medical Decision-Making

- A multifaceted and intricate process in which clinicians collaboratively navigate diverse sources of information to reach a precise and specific conclusion
- It involves interpreting complex and multi-modal data, such as imaging, electronic health records (EHR), physiological signals, and genetic information

## MEDAGENTS: LLMs as Collaborators for Zero-shot Medical Reasoning

 Use role-playing and collaborative agent discussion to increase faithfulness and interpretability of LLM without retrieval-augmented generation

#### Pros:

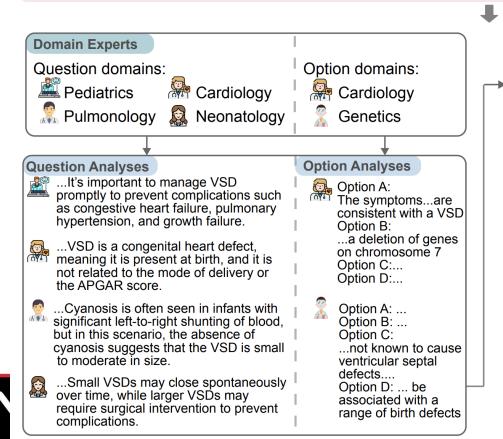
- Do not require external knowledge
- Enhance faithfulness of LLM results
- Provide interpretable reasoning chains

#### Points for Improvement:

- Use closed-source LLMs, require LLM's pretrained knowledge
- Agent numbers are fixed, which is not computational efficient
- Agent might not reach consensus

**Question**: A 3-month-old infant is brought to her pediatrician because she coughs and seems to have difficulty breathing while feeding. In addition, she seems to have less energy compared to other babies and appears listless throughout the day. She was born by cesarean section to a G1P1 woman with no prior medical history and had a normal APGAR score at birth. Her parents say that she has never been observed to turn blue. Physical exam reveals a high-pitched holosystolic murmur that is best heard at the lower left sternal border. The most likely cause of this patient's symptoms is associated with which of the following abnormalities?

Options: (A) 22q11 deletion (B) Deletion of genes on chromosome 7 (C) Lithium exposure in utero (D) Retinoic acid exposure in utero



#### **Initial Report**

**Key Knowledge**: Clinical assessment of an infant with symptoms suggesting VSD...

**Total Analysis**: The infant's symptoms are consistent with VSD... Options such as 22q11 deletion, deletion of genes on chromosome 7, lithium exposure in utero are not relevant to the given scenario.



#### **Unanimous Report**

Key Knowledge: The infant's symptoms are concerning for a possible congenital heart defect or a respiratory condition...

Total Analysis: ...one of the most common genetic abnormalities associated with congenital heart defects, including VSD, is the 22q11 deletion syndrome, also known as DiGeorge syndrome...

## MDAgents: An Adaptive Collaboration of LLMs

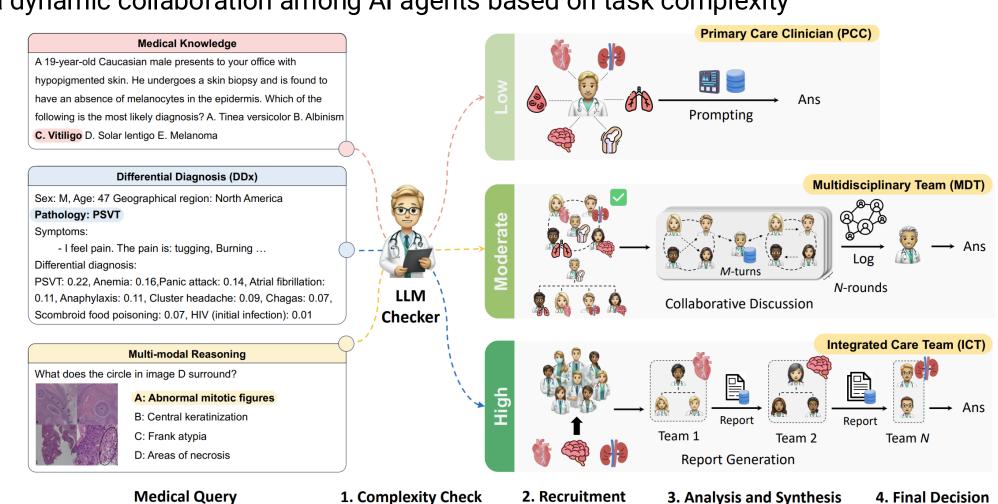
 The first adaptive decision-making framework for LLMs that mirrors real-world MDM processes via dynamic collaboration among AI agents based on task complexity

#### Pros:

- Framework imitates realworld clinical decision making
- Dynamic collaboration to improve efficiency
- Include a moderator agent to conduct planning
- Can integrate external knowledge

#### Points for Improvement:

- Use closed-source LLMs
- Limit to QA tasks





## **Takeaways**

- Notable points
  - Moving from static analysis models to interactive collaborators in clinical workflows
  - Al systems can now think, reason, and communicate rather than just process data
  - Success of AI systems depends on transparency, interpretability, and clinician trust

## An Interactive Example: Paper Hypothsis Assistant



# **Example Overview**

- We'll consolidate these ideas by exploring an interactive example of hypothesis generation.
- Our hands-on example will consist of two components:
- 1. Retrieve semantically relevant inspirations from training set
- 2. Generate hypotheses based on the retrieved inspirations, background context, and seed term
- Find Google Colab <a href="https://colab.google/">https://colab.google/</a> notebooks at
- https://github.com/EagleW/tutorial\_paper\_hypothesis/



# Why might language models hurt the scientific process?

- Factuality, misinformation, and hallucination
- Plagiarism
- High quantities of (often low-quality) work
- Issues with peer review
- These are issues <u>now!</u> The cat is out of the bag.

# **Factuality**

- LLMs like ChatGPT tend to hallucinate.
  - The only time this is good for science is hypothesis generation! When we read the literature or write papers, we don't want to make things up (usually).
  - ChatGPT often cites non-existent references. The proportion of nonexistent citations was over 30% with GPT-3.5, and it is over 20% with GPT-4 several months later.
  - ChatGPT generated fake court citations

# Publisher blacklists authors after preprint cites made-up studies

NM 2025 ANNUAL SY

Last month, a millipede expert in Denmark received an email notifying him that one of his publications had been mentioned in <u>a new manuscript</u> on *Preprints.org*. But when the researcher, <u>Henrik Enghoff</u>, downloaded the paper, he learned that it cited his work for something off-topic.



New York lawyers sanctioned for using fake ChatGPT cases in legal brief

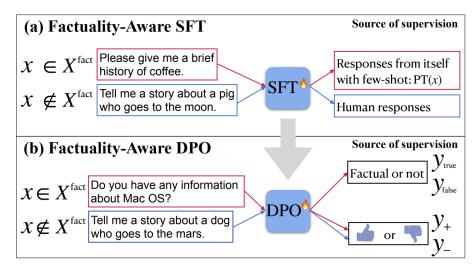


# Mitigation Methods: During Training

- Both supervised fine-tuning (SFT) and reinforcement learning (RL) can encourage hallucination
- Factuality-aware Alignment
  - Supervised Fine-Tuning (SFT) Adjustments: Use the LLM's own existing knowledge to generate training data, thus minimizing the incorporation of unfamiliar information
  - Reinforcement Learning (RL) Tweaks: The model's reward function is adjusted to include a direct preference for factuality using a factuality reward model (SFT)

### **Future Directions:**

Incorporate claim generation modules from other systems and external KBs might help LLMs verify factuality better

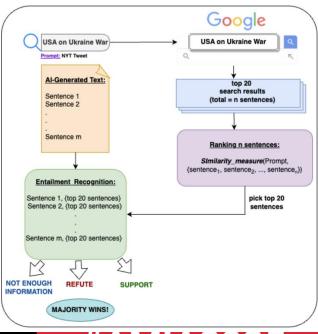


Lin, S. C., Gao, L., Oguz, B., Xiong, W., Lin, J., Yih, W. T., & Chen, X. (2024). FLAME: Factuality-Aware Alignment for Large Language Models. NeurIPS2024

# Mitigation Methods: After Training

- High Entropy Word Spotting and Replacement: Identify high-entropy words in hallucinated text and replace them with predictions from another LLM
- Sentence Factuality Checking: Identify sentences that are considered suspicious by NLI models for human review

Prompt	USA's stance on Ukraine Russia war
ated	The official stance of the United States on the Russia-Ukraine war has been consistent in supporting Ukraine's sovereignty, territorial integrity, and the peaceful resolution of the conflict. The U.S. has condemned Russia's actions in annexing Crimea and its support for separatist movements in eastern Ukraine. Despite maintaining a diplomatic stance,
Hallucinated text	the U.S. President Barack Obama stated that the United States would not deploy troops to Ukraine. Additionally, this lack of support extends to the withholding of financial and military assistance, including military equipment, training, and advisory support. This sentiment was reiterated by Paul Altman during a press review at the White House
Mitigated text ENTROPYBB	The official stance of the government on the ongoing war has been consistent in supporting resolution of the conflict. The government has condemned Russia in country and its support for separatist movements in eastern country. Despite maintaining a neutral stance, the government stated that the country would not deploy troops to unknown. Additionally, this lack of support extends to the withholding of financial and resources, including administrative, technical, and legal. This sentiment was echoed by minister during a press review at the conference
Edit Suggestions FACTUALITY GB	The official stance of the United States on the Russia-Ukraine war has been consistent in supporting Ukraine's sovereignty, territorial integrity, and the peaceful resolution of the conflict. The U.S. has condemned Russia's actions in annexing Crimea and its support for separatist movements in eastern Ukraine. Despite maintaining a diplomatic stance, U.S. President Barack Obama stated that the United States would not deploy troops to Ukraine. Additionally, this lack of support extends to the withholding of financial and military assistance, including military equipment, training, and advisory support.  This sentiment was reiterated by Paul Altman during a press review at the White House
Edi	



# Quantity overwhelms existing systems

- Papers can be generated much quicker now, but they aren't necessarily good
  - Current scientific productivity metrics are measured by quantity
  - "When a measure becomes a target, it ceases to be a good measure."
- Using "As of my last knowledge update" on Google Scholar returns 100s of results!

### ChatGPT on the Gravitational Redshift

J Stávek - European Journal of Applied Physics, 2024 - ej-physics.org

... As of my last knowledge update in January 2022, I'm not aware of any widely accepted alternative scenarios or theories that challenge Einstein's interpretation of gravitational redshift in ...

☆ Save 59 Cite Related articles All 2 versions >>>

An Evaluation of ChatGPT and Bard in the Context of Biological Knowledge Retrieval

R Caspi, PD Karp - Access Microbiology, 2024 - microbiologyresearch.org

... A: As of my last knowledge update in September 2021, the term "rbcX" was not commonly associated 345 with a specific protein or gene in the context of photosynthesis or carbon ...

☆ Save 兒 Cite Related articles All 2 versions



Radiology Case Reports

Volume 19, Issue 6, June 2024, Pages 2106-2111



Successful management of an Iatrogenic portal vein and hepatic artery injury in a 4month-old female patient: A case report and literature review

Raneem Bader MD a, Ashraf Imam MD b, Mohammad Alnees MD a e 🙎 🖾 , Neta Adler MD c, Joanthan ilia MD c, Diaa Zugayar MD b, Arbell Dan MD d, Abed Khalaileh MD b 🙎 🔀

Show more V + Add to Mendeley 🗬 Share 🗦 Cite

https://doi.org/10.1016/j.radcr.2024.02.037 7

Get rights and content 7

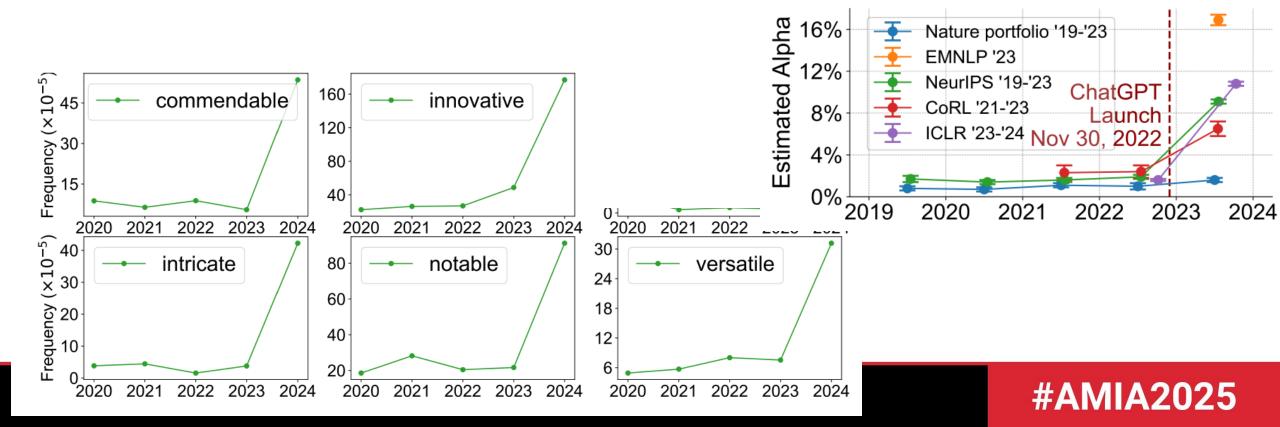
open access

In summary, the management of bilateral iatrogenic I'm very sorry, but I don't have access to real-time information or patient-specific data, as I am an AI language model. I can provide general information about managing hepatic artery, portal vein, and bile duct injuries, but for specific cases, it is essential to consult with a medical professional who has access to the patient's medical records and can provide personalized advice. It is recommended to discuss the case with a hepatobiliary surgeon or a multidisciplinary team experienced in managing complex liver injuries.



# Quantity overwhelms existing systems

- Peer-review challenges
  - Current scientific productivity metrics are measured by quantity
    - Measuring quality is hard to do without retrospect!
  - There aren't enough people to do free peer-review
    - Peer review takes time away from writing more papers
    - Many are turning to LLMs for help, but are the results thorough?



## What can we do? -- LLM Detectors

- Watermarking text (Kirchenbauer et al., 2023)
  - Encourage an unnatural distribution of generated words
    - Split the token set into a red list and a green list with a random seed.
    - Promote the use of green tokens, i.e., increase its probability
- Zero-shot approaches
  - Measure the perplexity and complexity of a generated sentence
- Retrieval (Krishna et al., 2023)
  - Given candidate text, search a database of previous LLM generations, looking for sequences that match the candidate text within a certain threshold
- Neural Network detectors for AI-generated text
  - Train a model to distinguish between human and AI text
  - These go out of date quickly and can be fooled by stronger models

# No watermark Extremely efficient on average term lengths and word frequencies on synthetic, microamount text (as little as 25 words) Very small and low-resource key/hash (e.g., 140 bits per key is sufficient for 99.9999999999 of the Synthetic Internet With watermark - minimal marginal probability for a detection attempt. - Good speech frequency and energy rate reduction.

- messages indiscernible to humans.

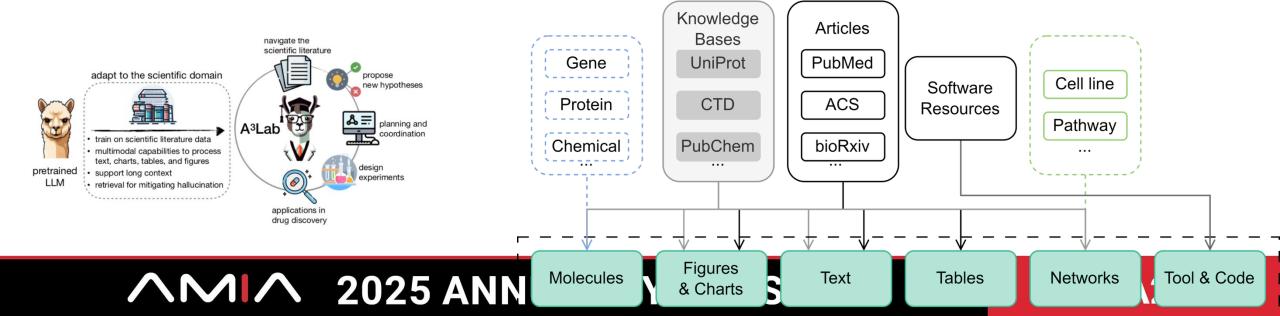
- easy for humans to verify.

# **Conclusion and Future Directions**



# Future Direction 1: Science-Inspired Large Language Models

- Existing work chooses either highly-capable general domain models (e.g., GPT-4) with strong reasoning capabilities or trains specific scientific language models (e.g., SciBERT) with specific knowledge.
  - Can we more effectively bridge this gap?
- Develop more capable multimodal LLMs
  - Explore instruction finetuning on a suite of datasets carefully designed to improve LLM capabilities in key tasks for scientific applications

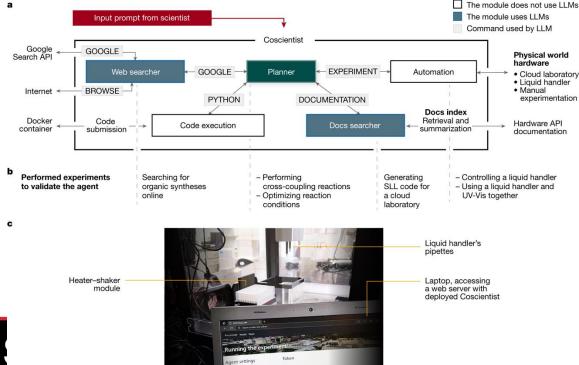


# Future Direction 2: Automated Experimental Systems

- Convert automated agents for scientific discovery from prototype ideas into useful tools
  - Read the entire scientific literature and synthesize relevant information to propose novel ideas
  - Couple with external simulations for rapid feedback

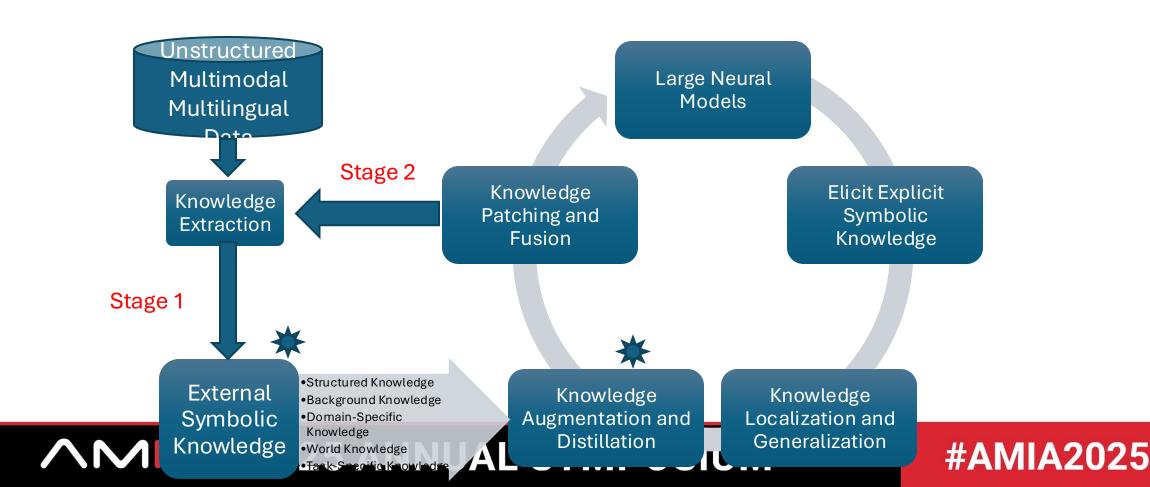
Integrate with specially designed laboratory hardware for automated

experimentation and testing



# Future Direction 3: Completing the LLM Knowledge Lifecycle

 Create a two-way knowledge acquisition framework to make symbolic and neural learning approaches mutually enhance each other



Questions?



