

# Evidence Retrieval and Grounding in Medicine

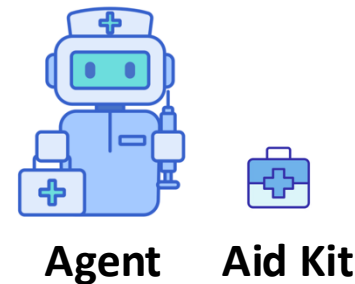
**Shuyue Jia**

Ph.D. Student

Advisor: Dr. Vijaya B. Kolachalama

Aug 9, 2025

Department of Electrical and Computer Engineering  
Boston University

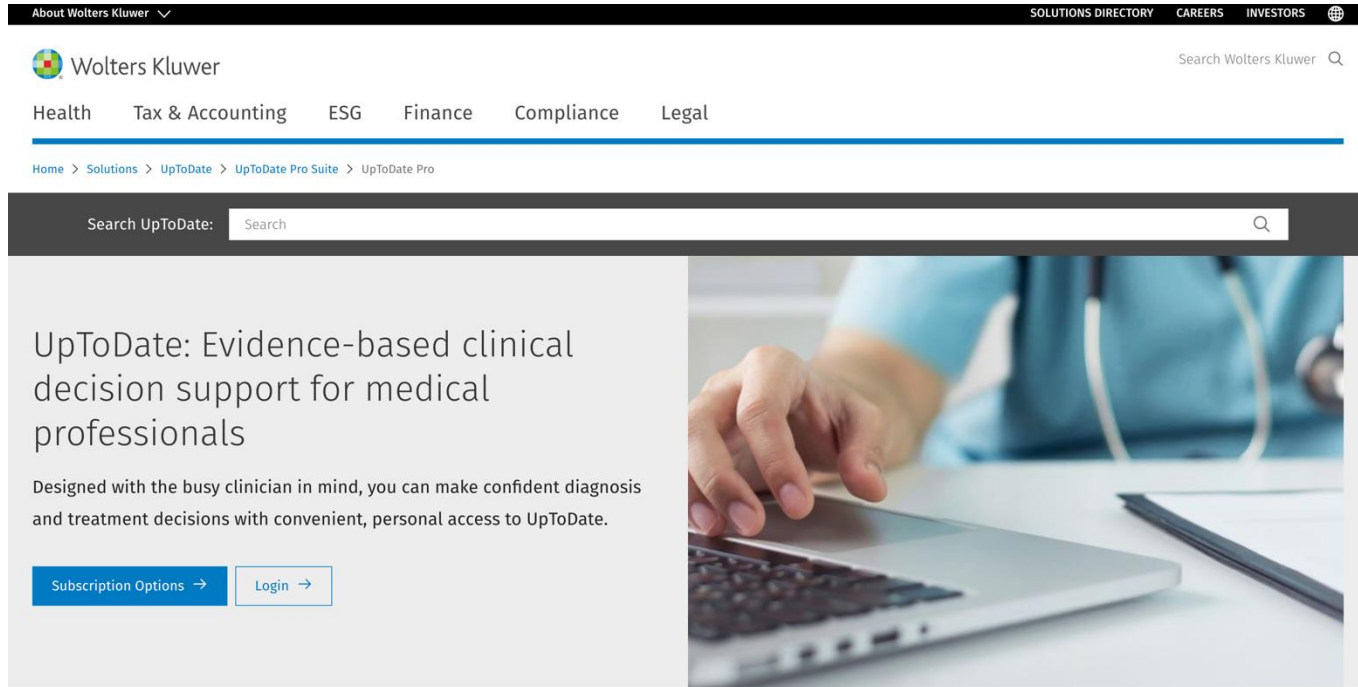


# Outline

- Background
- Related Works and Motivation
- MedPodGPT
- PodGPT
- Agentic System

# Part 1 – Background

- Traditional search engine: **UpToDate** ([\\$6,713 million revenue, 2024](#))



[Why UpToDate?](#)[Product](#)[Editorial](#)[Subscription Options](#)[All](#)[Adult](#)[Pediatric](#)[Patient](#)[Graphics](#)

Showing results for **what is the recommended treatment of chronic insomnia in adults**

## Overview of the treatment of insomnia in adults

... three active **treatment** arms (CBT-I, zolpidem, or both) in 63 young and middle-aged **adults** with **chronic insomnia**, there were no differences in total sleep time among active **treatment** groups and... every six months is **recommended**. In patients who fail insomnia **treatment**, it is important to discuss expectations of sleep, particularly for older **adults** and those with comorbidities....

## Pharmacotherapy for insomnia in adults

... the preferred first-line **treatment** for **chronic insomnia** in **adults** and has been endorsed as first-line **therapy** by multiple societies and guideline panels .... Pharmacotherapy should not be the sole **treatment** for insomnia. CBT-I is **recommended** as first-line **treatment** for **chronic insomnia** ....

## Cognitive behavioral therapy for insomnia in adults

... This approach is a **recommended treatment** for insomnia disorder and is a multistep process carried out over multiple sessions. During the **treatment**,... and delivery of CBT-I and other behavioral **treatments** for insomnia in **adults**. An overview of the **treatment** of insomnia and pharmacologic **therapies** for insomnia are presented separately.... behavioral **therapy** for insomnia (CBT-I) is a multicomponent **treatment** for **chronic insomnia** disorder that aims to identify and target the multiple...



[< Back](#)

## Pharmacotherapy for insomnia in adults

### Outline



#### SUMMARY AND RECOMMENDATIONS

#### INTRODUCTION

#### DRUG SELECTION

Our approach

Patients with isolated sleep-onset insomnia

Patients with sleep-maintenance or mixed insomnia

Special populations

#### PRESCRIBING AND MONITORING

Shared warnings and precautions

Safe prescribing practices

Initial doses and adjustments

Response assessment

Patients with inadequate response

Treatment duration

Select Language



**AUTHOR:** [David N Neubauer, MD](#)

**SECTION EDITORS:** [Ruth Benca, MD, PhD](#), [Joann G Elmore, MD, MPH](#)

**DEPUTY EDITOR:** [April F Eichler, MD, MPH](#)

Literature review current through: **Apr 2025**.

This topic last updated: **May 21, 2025**.

### INTRODUCTION

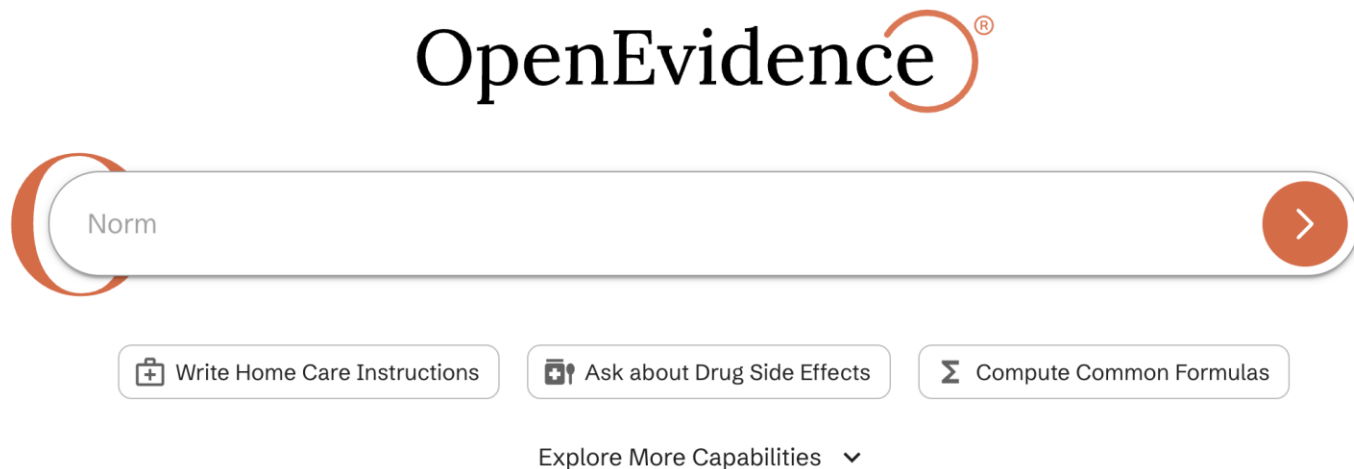
Insomnia is one of the most common symptoms for which adults seek medical advice. While pharmacotherapy should not be the sole treatment for insomnia, medication may be part of an integrated approach that includes behavioral strategies and treatment of relevant comorbidities.

A range of medications and substances are used to treat insomnia in adults, spanning the following categories:

- Medications with regulatory approval for treatment of insomnia
- Off-label prescription medications prescribed for insomnia
- Over-the-counter (OTC) sleep aids (eg, diphenhydramine and doxylamine, alone or in "PM" formulated medications)

# Part 1 – Background

- AI-powered search engine: **OpenEvidence** ([\\$1 billion valuation, 2025](#))





# Partnering with OpenEvidence: A Life-Saving Healthcare Revolution

This AI-powered search engine for doctors could be one of the most important companies of the next decade.



What is the latest diagnosis of dementia?

▼ Finished thinking

The **latest diagnostic criteria for dementia** are based on a combination of clinical assessment and, increasingly, biomarker evidence, with criteria evolving to reflect advances in the understanding of neurodegenerative diseases.

For all-cause dementia, the National Institute on Aging and the Alzheimer's Association (NIA-AA) recommend that dementia is diagnosed when there is evidence of cognitive or behavioral (neuropsychiatric) symptoms that interfere significantly with the ability to function at work or in usual daily activities. This determination is made through clinical judgment, based on history from the patient and a knowledgeable informant, and supported by cognitive testing. The diagnosis requires impairment in at least two cognitive domains (such as memory, executive function, language, visuospatial skills, or personality/behavior), and the deficits must represent a decline from previous levels of functioning and not be explained by delirium or major psychiatric disorder. The NIA-AA emphasizes that the distinction between mild cognitive impairment and dementia is based on the degree of functional impairment, with dementia requiring significant interference in daily life.<sup>[1]</sup>

Recent updates from the Alzheimer's Association Workgroup further incorporate biomarkers into the diagnostic process, particularly for Alzheimer's disease. The 2024 criteria define Alzheimer's disease as a biological process, with diagnosis possible in asymptomatic individuals based on abnormal amyloid and tau biomarkers (e.g., amyloid PET, CSF, or plasma phosphorylated tau). An abnormal "Core 1" biomarker is sufficient to establish a diagnosis of Alzheimer's disease, while "Core 2" biomarkers (such as tau PET) provide prognostic information and increase diagnostic confidence. The criteria also recognize the importance of integrating clinical and biological staging, acknowledging that comorbidities and cognitive reserve may influence the clinical presentation.<sup>[2-3]</sup>

In summary, the most up-to-date consensus is that **dementia is diagnosed clinically by evidence of acquired cognitive decline interfering with daily function, with increasing use of biomarkers for etiological specificity, especially in Alzheimer's disease, as recommended by the National Institute**

---

In summary, the most up-to-date consensus is that **dementia is diagnosed clinically by evidence of acquired cognitive decline interfering with daily function, with increasing use of biomarkers for etiological specificity, especially in Alzheimer's disease, as recommended by the National Institute on Aging and the Alzheimer's Association.**<sup>[2]</sup> ■



Helpful



Not Helpful



## References



1. [The Diagnosis of Dementia Due to Alzheimer's Disease: Recommendations From the National Institute on Aging-Alzheimer's Association Workgroups on Diagnostic Guidelines for Alzheimer's Disease.](#)



McKhann GM, Knopman DS, Chertkow H, et al.

Alzheimer's & Dementia : The Journal of the Alzheimer's Association. 2011;7(3):263-9.

doi:10.1016/j.jalz.2011.03.005.

☆ [Leading Journal](#) ⓘ

2. [Revised Criteria for Diagnosis and Staging of Alzheimer's Disease: Alzheimer's Association Workgroup.](#)



Jack CR, Andrews JS, Beach TG, et al.

Alzheimer's & Dementia : The Journal of the Alzheimer's Association. 2024;20(8):5143-5169.

doi:10.1002/alz.13859.

☆ [Leading Journal](#) ⓘ [New Research](#)

3. [Influence of Different Diagnostic Criteria on Alzheimer Disease Clinical Research.](#)



Bieger A, Brum WS, Borelli WV, et al.

Neurology. 2024;103(5):e209753. doi:10.1212/WNL.0000000000209753.

[New Research](#)

## USMLE Test Prep

## Create Study Session

BETA

Step 1

Step 2 CK

Step 3

Study specific topics with dynamically generated questions.

Main Topic ▼

Subtopic ▼

Start Study Session

A 29-year-old female presents to the clinic with a chief complaint of fatigue and weight gain over the past six months. She works as a high school teacher in a rural area with limited access to healthcare specialists. The patient reports feeling cold all the time and has experienced constipation and mild depression. She also notes her periods have become irregular, and she has had difficulty concentrating at work. She denies any recent illnesses or significant stressors.

Her past medical history is significant for a head injury from a motor vehicle accident two years ago, which required hospitalization but did not result in any lasting neurological deficits. She takes no regular medications and denies smoking, alcohol, or illicit drug use.

On physical examination, the patient appears tired and has a dry, pale skin texture. Her vital signs are: blood pressure 110/70 mmHg, heart rate 58 bpm, temperature 96.8°F (36°C), and respiratory rate 14 breaths per minute. Thyroid examination reveals no palpable goiter. Reflexes are delayed, and she has mild facial puffiness with periorbital edema. No other abnormalities are noted.

Laboratory results reveal:

- TSH: 0.4  $\mu$ IU/mL (normal: 0.5–4.5  $\mu$ IU/mL)
- Free T4: 0.6 ng/dL (normal: 0.8–1.8 ng/dL)
- Free T3: 1.8 pg/mL (normal: 2.3–4.2 pg/mL)
- Serum cortisol (8 AM): 5  $\mu$ g/dL (normal: 5–25  $\mu$ g/dL)
- Prolactin: 25 ng/mL (normal: 4–23 ng/mL)

Additional imaging with MRI of the brain reveals a 1 cm lesion in the pituitary region suggestive of a pituitary adenoma.

A 29-year-old female presents to the clinic with fatigue, weight gain, cold intolerance, constipation, mild depression, irregular periods, and difficulty concentrating. She has a history of a head injury from a motor vehicle accident two years ago. Physical examination reveals dry, pale skin, delayed reflexes, mild facial puffiness, and periorbital edema. Laboratory results show low TSH, free T4, and free T3 levels, with a slightly elevated prolactin level. An MRI of the brain reveals a 1 cm lesion in the pituitary region suggestive of a pituitary adenoma. What is the most appropriate next step in management?

- ☐ (A) Schedule the patient for surgical removal of the pituitary adenoma
- ☐ (B) Order a fine-needle aspiration biopsy of the pituitary lesion
- ☐ (C) Initiate levothyroxine therapy
- ☐ (D) Refer to an endocrinologist for further evaluation and management
- ☐ (E) Start dopamine agonist therapy

Submit

## USMLE-style multiple-choice question answering:



Medical Students



Residents

## Need harder questions: Fellows Doctors (Clinicians/Physicians)

mild depression, irregular periods, and difficulty concentrating. She has a history of a head injury from a motor vehicle accident two years ago. Physical examination reveals dry, pale skin, delayed reflexes, mild facial puffiness, and periorbital edema. Laboratory results show low TSH, free T4, and free T3 levels, with a slightly elevated prolactin level. An MRI of the brain reveals a 1 cm lesion in the pituitary region suggestive of a pituitary adenoma. What is the most appropriate next step in management?

- ☐ (A) Schedule the patient for surgical removal of the pituitary adenoma
- ☒ (B) Order a fine-needle aspiration biopsy of the pituitary lesion
- ☐ (C) Initiate levothyroxine therapy
- ☐ (D) Refer to an endocrinologist for further evaluation and management
- ☐ (E) Start dopamine agonist therapy

Submit

Incorrect. The correct answer is: Refer to an endocrinologist for further evaluation and management

Create Anki Card

Explain

Next



▼ Finished thinking

**Referral to an endocrinologist for further evaluation and management is the most appropriate next step** because pituitary adenomas presenting with central hypothyroidism and mild hyperprolactinemia require specialized assessment to confirm the diagnosis, evaluate for additional pituitary hormone deficiencies, and determine the functional status of the tumor. Central hypothyroidism, particularly in the context of a pituitary mass, is frequently accompanied by other pituitary hormone deficits, and its management—including the initiation and titration of levothyroxine—demands expertise to avoid precipitating adrenal insufficiency or other complications.<sup>[1-6]</sup>

Endocrinologists are uniquely qualified to distinguish between functioning and nonfunctioning adenomas, interpret laboratory artifacts such as the hook effect, and initiate appropriate first-line therapies (e.g., dopamine agonists for prolactinomas, surgical referral for nonfunctioning or other functioning adenomas).<sup>[2][6-8]</sup> The Endocrine Society and other expert guidelines recommend comprehensive endocrine evaluation in all patients with pituitary adenomas, including assessment for hormone hypersecretion and hypopituitarism, as well as formal visual field testing if there is concern for optic chiasm involvement.<sup>[1]-2[4]</sup> Early specialist involvement is critical to prevent morbidity from untreated hormone deficiencies, guide further imaging or surgical referral if indicated, and optimize long-term outcomes.<sup>[1-2][4]</sup>

In summary, the complexity of diagnosis and management, the need for tailored hormone replacement, and the potential for significant morbidity underscore the necessity of endocrinology referral in this clinical scenario. █

☰ References ▼



Large ▼

Ask a follow-up question...





Endocrinologists are uniquely qualified to distinguish between functioning and nonfunctioning adenomas, interpret laboratory artifacts such as the hook effect, and initiate appropriate first-line therapies (e.g., dopamine agonists for prolactinomas, surgical referral for nonfunctioning or other functioning adenomas).<sup>[2][6-8]</sup> The Endocrine Society and other expert guidelines recommend comprehensive endocrine evaluation in all patients with pituitary adenomas, including assessment for hormone hypersecretion and hypopituitarism, as well as formal visual field testing if there is concern for optic chiasm involvement.<sup>[1-2][4]</sup> Early specialist involvement is critical to prevent morbidity from untreated hormone deficiencies, guide further imaging or surgical referral if indicated, and optimize long-term outcomes.<sup>[1-2][4]</sup>

In summary, the complexity of diagnosis and management, the need for tailored hormone replacement, and the potential for significant morbidity underscore the necessity of endocrinology referral in this clinical scenario. ■

## References

1. [Pituitary Incidentaloma: An Endocrine Society Clinical Practice Guideline](#).  
Freda PU, Beckers AM, Katznelson L, et al.  
The Journal of Clinical Endocrinology and Metabolism. 2011;96(4):894-904. doi:10.1210/jc.2010-1048.
2. [Diagnosis and Management of Pituitary Adenomas: A Review](#).  
Tritos NA, Miller KK.  
Jama. 2023;329(16):1386-1398. doi:10.1001/jama.2023.5444.  
★ [Leading Journal](#)
3. [Central Hypothyroidism: Advances in Etiology, Diagnostic Challenges, Therapeutic Targets, and Associated Risks](#).  
Iglesias P.  
Endocrine Practice : Official Journal of the American College of Endocrinology and the American Association of Clinical Endocrinologists. 2025;31(5):650-659. doi:10.1016/j.eprac.2025.02.004.  
📄 [New Research](#)
4. [Hormonal Replacement in Hypopituitarism in Adults: An Endocrine Society Clinical Practice Guideline](#).  
Fleseriu M, Hashim IA, Karavitaki N, et al.  
The Journal of Clinical Endocrinology and Metabolism. 2016;101(11):3888-3921. doi:10.1210/jc.2016-2118.
5. [Interactions Between Hypothalamic Pituitary Thyroid Axis and Other Pituitary Dysfunctions](#).  
Feldt-Rasmussen U, Klose M, Benvenga S.  
Endocrine. 2018;62(3):519-527. doi:10.1007/s12020-018-1738-6.
6. [Approach to the Patient With Prolactinoma](#).  
Auriemma RS, Pirchio R, Pivonello C, et al.  
The Journal of Clinical Endocrinology and Metabolism. 2023;108(9):2400-2423. doi:10.1210/clinem/dgad174.
7. [Diagnosis and Treatment of Pituitary Adenomas: A Review](#).  
Molitch ME.  
Jama. 2017;317(5):516-524. doi:10.1001/jama.2016.19699.  
★ [Leading Journal](#)
8. [Pituitary-Tumor Endocrinopathies](#).  
Melmed S.  
The New England Journal of Medicine. 2020;382(10):937-950. doi:10.1056/NEJMra1810772.  
★ [Leading Journal](#)

# Part 1 – Background

## ■ Patient background

A 27-year-old man presents to the emergency department. He was brought in by staff from the homeless shelter when they found him unresponsive. The patient is a known IV drug abuser but otherwise has an unknown past medical history. He currently attends a methadone clinic. His temperature is 99.5°F (37.5°C), blood pressure is 97/48 mmHg, pulse is 140/min, respirations are 29/min, and oxygen saturation is 98% on room air. Initial laboratory values are shown below. Serum: Na<sup>+</sup>: 139 mEq/L. Cl<sup>-</sup>: 100 mEq/L. K<sup>+</sup>: 6.3 mEq/L. HCO<sub>3</sub><sup>-</sup>: 17 mEq/L. Glucose: 589 mg/dL. The patient is given treatment. After treatment, his temperature is 99.5°F (37.5°C), blood pressure is 117/78 mmHg, pulse is 100/min, respirations are 23/min, and oxygen saturation is 98% on room air. His laboratory values are seen below. Serum: Na<sup>+</sup>: 139 mEq/L. Cl<sup>-</sup>: 100 mEq/L. K<sup>+</sup>: 4.3 mEq/L. HCO<sub>3</sub><sup>-</sup>: 19 mEq/L. Glucose: 90 mg/dL.

# Part 1 – Background

- **Question and options**

- **Question**

What is the best next step in management?

- **Options**

- A. Insulin, potassium, IV fluids, and glucose
    - B. IV fluids only
    - C. Oral rehydration
    - D. Supportive therapy and close monitoring

- **Answer**

A. Insulin, potassium, IV fluids, and glucose

**(1) Multiple choice question-answering**

**(2) Open-ended question-answering**

## Part 2 – Related Works

- Question-answering without rationale → **lack explainability**

Dataset	Source	Description	n
<b>Finetuning</b>			
Medical Flash Cards	Anki Flashcards	Rephrased Q&A pairs derived from the front and back sides of medical flashcards	33,955
Stack Exchange	Academia	Q&A pairs generated from questions and their top-rated answers	39,633
	Biology		7,482
	Fitness		3,026
	Health		1,428
	Bioinformatics		906
Wikidoc	Living Textbook	Q&A pairs generated from paragraphs, where questions were formulated from rephrased paragraph titles, and answers were extracted from paragraph text	67,704
	Patient Information	Q&A pairs generated from paragraph headings and associated text content	5,942
<b>Evaluation</b>			
USMLE	Step 1	Multiple choice questions from the USMLE self-assessment with image-based questions excluded	119
	Step 2		120
	Step 3		135

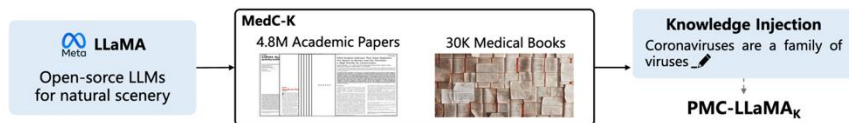
Instruction tuning

(Supervised fine-tuning)

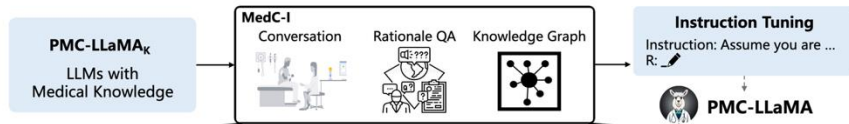
## Part 2 – Related Works

### ■ Question-answering with rationale → lack factuality

#### Step-I Data-centric Knowledge Injection



#### Step-II Medical-specific Instruction Tuning



#### MedC-I Samples

Conversation	Rationale QA	Knowledge Graph
<b>Instruction:</b> If you are a doctor, please answer the medical questions based on ... <b>Input:</b> Doctor, I have been experiencing ... What could be the problem? <b>Response:</b> It's possible that you have a vocal cord polyp. To confirm ...	<b>Instruction:</b> In your capacity as ... Answer the medical questions. <b>Input:</b> Question: Which of the following ... Options: A ... B ... C ... D ... <b>Response:</b> Option A is wrong because ... Answer: Option D is correct.	<b>Prompt-Description</b> <b>Instruction:</b> ... Explain the definition of ... <b>Input:</b> Question: What is the meaning ... <b>Response:</b> Answer: the entity denotes ... <b>Prompt-Relation</b> <b>Instruction:</b> ... Determine the relation between ... <b>Input:</b> Question: What is the relation ... <b>Response:</b> Answer: Mercaptopurine ... has ...

###Rationale: This patient's clinical presentation and laboratory abnormalities suggest a diagnosis of thrombotic thrombocytopenic purpura (TTP). TTP is a microangiopathic hemolytic anemia, resulting from severe deficiency or functional abnormality of the cleaving protease ADAMTS13, leading to accumulation of ultra-large multimers of von Willebrand factor, platelet adhesion, and widespread thrombi in microcirculation. Clinically, TTP is characterized by the pentad of thrombocytopenia, microangiopathic hemolytic anemia, neurologic abnormalities, renal insufficiency, and fever. This pentad is not always complete, and high clinical suspicion should be maintained even if not all five elements are present, as in this case. In terms of management, plasma exchange therapy is the mainstay and should be initiated as soon as the diagnosis is suspected, even before laboratory confirmation, given the severity and rapid progression of TTP. This therapy helps by removing the autoantibodies that inhibit ADAMTS13 and replenishing the deficient enzyme through plasma replacement. Options A and B, high-dose glucocorticoids and cyclophosphamide and rituximab, respectively, can be used in addition to plasma exchange in certain refractory or relapsed cases but are not first-line treatment. Option C, Vancomycin and cefepime, are antibiotics used to treat infections which doesn't align with this patient's presentation.

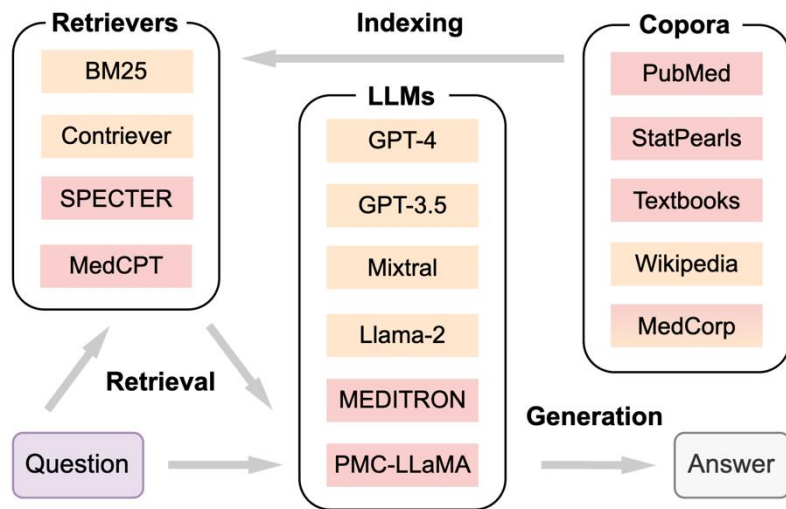
###Answer: OPTION D IS CORRECT.

Credit:

- [1] Wu et al., PMC-LLaMA: Towards Building Open-source Language Models for Medicine, In Journal of the American Medical Informatics Association'23.
- [2] Qiu et al., Towards Building Multilingual Language Model for Medicine, In Nature Communications'24.

## Part 2 – Related Works

- **Question-answering with retrieval-augmented generation (RAG)**
  - **Bias problem** ← no further evidence assessment
  - **Low efficiency** ← model doesn't decide whether to use evidence



## Part 2 – Motivation

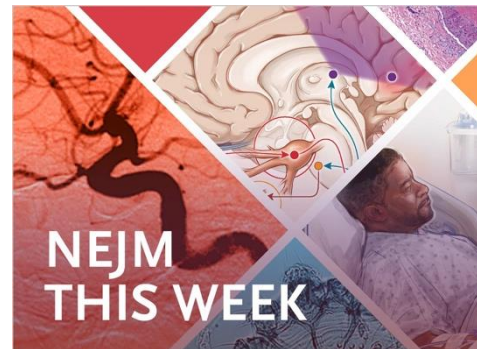
- **Question-answering without rationale**
  - Clinical decision-making with rationale → **explainability**
- **Question-answering with rationale**
  - Clinical decision-making with external evidence → **factuality**
- **Question-answering with RAG**
  - **Clinical decision-making with evidence assessment** ← bias problem
  - **Clinical decision-making with an automatic and dynamic process** ← low efficiency

Credit:

- [1] Han et al., MedAlpaca: An Open-Source Collection of Medical Conversational AI Models and Training Data, In arXiv'23.
- [2] Wu et al., PMC-LLaMA: Towards Building Open-source Language Models for Medicine, In Journal of the American Medical Informatics Association'23.
- [3] Qiu et al., Towards Building Multilingual Language Model for Medicine, In Nature Communications'24.
- [4] Xiong et al., Benchmarking Retrieval-Augmented Generation for Medicine, In ACL'24.

## Part 3 – MedPodGPT

- What: A **multilingual audio-augmented LLM** for medical research and education

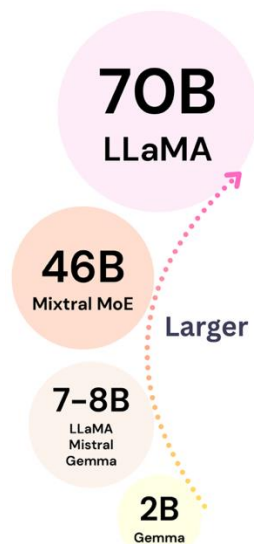


- Why: **Up to date and domain-relevant information** ← podcasts
- Where: Journals, exam preparation materials, and clinical practice
- When: Most recent (e.g., since 2023)
- How: Continual pretraining of large language models (LLMs)





## Part 3 – MedPodGPT



### Large language models

- 2B Gemma
- 7B Gemma, 7B Mistral, 8B LLaMA
- 46B Mixtral 8×7B MoE
- 70B LLaMA

$$\mathcal{L}_{\pi_{\theta}} = - \sum \log(\pi_{\theta}(x_i | \mathbf{x}_{<i})) .$$

$\pi_{\theta}$ : LLM, parameterized by  $\theta$

$\mathbf{x} = [x_1, x_2, \dots, x_t]$ : a sequence of texts



# Part 3 – MedPodGPT

## ■ Podcasts



The NEW ENGLAND  
JOURNAL of MEDICINE

CURRENT ISSUE ▾ SPECIALTIES ▾ TOPICS ▾ MULTIMEDIA ▾ LEARNING/CME ▾ AUTHOR CENTER 🔍

## Podcasts & Feeds

### INTENTION TO TREAT

Incisive analysis of critical and timely issues in medicine and health care.

[Recent episodes.](#)



FOLLOW:

### NEJM THIS WEEK

A summary of what's in this week's issue.

[Recent episodes.](#)



FOLLOW:

### NOT OTHERWISE SPECIFIED

Exploring some of health care's toughest challenges and areas of greatest promise.

[Recent episodes.](#)



FOLLOW:

### NEJM INTERVIEWS

Current topics in medicine and health care, with authors and editors.

[Recent episodes.](#)



FOLLOW:

### The AI Revolution in Medicine

[Download PDF Now](#)

Understanding  
How AI Can Advance  
Patient Care



### THE NEJM AI GRAND ROUNDS PODCAST

Informal conversations with a variety of experts exploring the deep issues at the intersection of artificial intelligence, machine learning, and medicine.

[Recent episodes.](#)



## Part 3 – MedPodGPT

### ■ Podcasts (journals)

- [OpenAI Whisper automatic speech recognition \(ASR\) model](#) (2022)
- Audio → texts (transcripts)

Journal Podcasts						
Podcast	Language	Episodes	Audio Time (min)	Mean Length Episode $\pm \sigma$ (min)	Number of Text Tokens	Mean Text Tokens per Episode $\pm \sigma$
NEJM	English	1974	39256.0	$19.89 \pm 9.74$	4,760,783	$1928.22 \pm 14.87$
JAMA	English	2235	32163.0	$14.39 \pm 8.66$	3,454,191	$1928.64 \pm 15.54$
The Lancet	English	2029	28279.0	$13.94 \pm 7.62$	3,300,982	$1925.89 \pm 20.88$
The BMJ	English	300	13264.2	$44.21 \pm 10.35$	2,235,458	$1897.67 \pm 75.07$
Annals Latest Highlights	English	396	6427.0	$16.23 \pm 8.09$	803,958	$1927.96 \pm 14.60$
Annals On Call	English	142	3440.0	$24.22 \pm 3.65$	522,547	$1928.22 \pm 15.64$
Pediatrics on Call	English	98	3299.0	$33.66 \pm 6.09$	565,781	$1930.99 \pm 15.00$
Procedure Ready: Ob/Gyn	English	20	383.7	$19.19 \pm 5.00$	63,667	$1929.30 \pm 13.41$
Revista Médica AFP Podcast	Spanish	40	1055.0	$26.38 \pm 3.70$	190,518	$1924.42 \pm 16.34$

Credit:

[1] <https://openai.com/index/whisper/>

[2] Jia et al., MedPodGPT: A Multilingual Audio-augmented Large Language Model for Medical Research and Education, In medRxiv'24.

## Part 3 – MedPodGPT

### ■ Podcasts (test preparations)

Test Preparation Podcasts						
Podcast	Language	Episodes	Audio Time (min)	Mean Length of Episode $\pm \sigma$ (min)	Number of Text Tokens	Mean Text Tokens per Episode $\pm \sigma$
Divine Intervention Podcasts	English	480	18363.8	$38.26 \pm 24.07$	2,269,153	$1931.19 \pm 13.53$
The Radiology Review Podcast	English	127	2517.7	$19.82 \pm 10.10$	292,949	$1927.29 \pm 26.81$
Crush Step 1: The Ultimate USMLE Step 1 Review	English	49	2176.2	$44.41 \pm 15.31$	328,194	$1930.55 \pm 13.03$
The USMLE Guys Podcast	English	31	1464.3	$47.24 \pm 43.47$	156,923	$1937.32 \pm 6.12$
Harrison's PodClass: Internal Medicine Cases and Board Prep	Spanish	95	905.2	$9.53 \pm 2.24$	101,574	$1916.49 \pm 22.77$
El Interno Desvelado	Spanish	4	99.13	$24.78 \pm 11.91$	17,121	$1902.33 \pm 25.54$
Curso MIR Asturias	Spanish	3	17.7	$5.89 \pm 4.53$	3,872	$1936.00 \pm 9.00$

# Part 3 – MedPodGPT

## ■ Podcasts (clinical experts)

Clinical Podcasts						
Podcast	Language	Episodes	Audio Time (min)	Mean Length of Episode $\pm \sigma$ (min)	Number of Text Tokens	Mean Text Tokens per Episode $\pm \sigma$
The Curbsiders Internal Medicine Podcast	English	485	28749.2	59.39 $\pm$ 16.08	5,772,083	1929.82 $\pm$ 17.06
This Podcast Will Kill You	English	168	18363.8	38.26 $\pm$ 24.07	2,269,153	1931.19 $\pm$ 13.53
The Clinical Problem Solvers	English	315	13500.1	42.86 $\pm$ 14.51	2,493,777	1927.18 $\pm$ 23.32
PsychEd: educational psychiatry podcast	English	62	3556.3	57.36 $\pm$ 17.52	607,237	1927.74 $\pm$ 16.36
Run the List	English	97	1973.0	20.34 $\pm$ 6.44	352,977	1928.84 $\pm$ 15.17
Goljan Pathology Lectures	English	37	1886.0	50.97 $\pm$ 4.58	412,086	1934.68 $\pm$ 13.45
Core IM: 5 Pearls	English	54	1847.1	34.21 $\pm$ 10.19	361,213	1931.62 $\pm$ 10.16
Neurology Clinical Pearls	English	27	333.2	12.34 $\pm$ 3.19	42,494	1931.54 $\pm$ 10.78
Tutorías Medicina Interna	Spanish	570	19834.9	34.80 $\pm$ 25.01	4,311,263	1898.39 $\pm$ 64.31
Leucocitos isotópicos	Spanish	68	2537.8	37.32 $\pm$ 9.55	481,676	1797.29 $\pm$ 154.42
Medicina Con Cabeza	Spanish	246	2457.8	9.99 $\pm$ 3.44	462,383	1902.81 $\pm$ 57.55
Medicina de impacto	Spanish	57	2406.5	42.22 $\pm$ 9.13	492,363	1915.81 $\pm$ 29.28
Ronda, El Podcast de Medicina Interna	Spanish	20	1084.4	54.22 $\pm$ 25.01	206,218	1891.91 $\pm$ 71.90
Medicina De Bolsillo   Hablando de Medicina	Spanish	45	958.3	21.30 $\pm$ 10.79	186,268	1844.24 $\pm$ 124.82
La Tertulia de Cajal	Spanish	27	876.3	32.46 $\pm$ 18.28	186,001	1897.97 $\pm$ 57.71
PedCast: Dos Pediatras y un Podcast	Spanish	14	458.5	32.75 $\pm$ 10.62	89,127	1896.32 $\pm$ 58.05
Neurobiologie et Immunité	French	21	1882.8	89.66 $\pm$ 14.77	383,189	1896.97 $\pm$ 40.12
Incubateur Néonate	French	25	1579.3	63.17 $\pm$ 21.28	391,475	1918.99 $\pm$ 24.15
Guideline.care	French	68	1369.1	20.13 $\pm$ 6.30	293,301	1917.0 $\pm$ 29.29
La Minute Rhumato	French	119	921.0	7.74 $\pm$ 2.19	132,354	1918.17 $\pm$ 23.59
Oncologie cellulaire et moléculaire - Hugues de Thé	French	11	852.9	77.53 $\pm$ 19.81	186,693	1905.03 $\pm$ 44.42
Le podcast des Conférenciers (UFR3S) by Université de Lille	French	65	768.4	11.82 $\pm$ 19.58	86,105	1913.44 $\pm$ 42.78
Super Docteur	French	47	676.3	14.39 $\pm$ 6.50	139,824	1915.40 $\pm$ 33.26
Médecine, Sciences et Recherche clinique	French	24	332.2	13.84 $\pm$ 4.58	63,314	1918.61 $\pm$ 26.60
NéphroDio	French	40	318.6	7.96 $\pm$ 2.58	55,716	1921.24 $\pm$ 19.59
La Minute Néonate	French	37	307.6	8.31 $\pm$ 1.93	57,435	1914.50 $\pm$ 31.58
Le Med G Eclairé	French	11	249.2	22.66 $\pm$ 16.76	51,988	1925.48 $\pm$ 12.57
La Minute du Pancréas	French	22	209.4	9.52 $\pm$ 2.34	38,376	1918.80 $\pm$ 23.34
L'essentiel des principales pathologies	French	14	151.3	10.81 $\pm$ 13.10	23,098	1924.83 $\pm$ 11.25
AR-Pod le Podcast de l'anesthésie-réanimation	French	12	139.0	11.59 $\pm$ 4.52	22,998	1916.50 $\pm$ 26.48



# Part 3 – MedPodGPT

## ■ Benchmarks (in-domain performance)

Language	Benchmark Datasets	Model											
		Gemma 2B		Gemma 7B		Mistral 7B		LLaMA 3 8B		Mixtral MoE		LLaMA 3 70B	
		Baseline	<i>Ours</i>	Baseline	<i>Ours</i>	Baseline	<i>Ours</i>	Baseline	<i>Ours</i>	Baseline	<i>Ours</i>	Baseline	<i>Ours</i>
English	MedExpQA	15.20	<b>23.40</b>	34.40	<b>45.20</b>	<b>47.20</b>	46.20	57.60	<b>61.40</b>	52.80	<b>61.20</b>	<b>78.40</b>	77.60
	MedMCQA	34.81	<b>35.24</b>	40.66	<b>44.86</b>	42.65	<b>45.50</b>	58.64	<b>58.82</b>	50.20	<b>53.54</b>	<b>71.12</b>	70.58
	MedQA	29.69	<b>33.38</b>	38.26	<b>44.56</b>	46.27	<b>47.80</b>	<b>61.12</b>	59.21	<b>54.05</b>	53.22	<b>77.85</b>	77.51
	PubMedQA	47.80	<b>55.50</b>	<b>63.40</b>	55.30	<b>51.60</b>	41.75	<b>59.40</b>	49.20	<b>42.80</b>	32.20	73.00	<b>75.75</b>
	Anatomy	<b>43.70</b>	42.04	49.63	<b>52.96</b>	56.30	<b>56.67</b>	68.89	<b>69.82</b>	64.44	<b>68.15</b>	77.04	<b>77.78</b>
	Clinical Knowledge	<b>41.51</b>	38.78	55.47	<b>62.17</b>	61.89	<b>62.26</b>	72.08	<b>73.68</b>	67.92	<b>74.90</b>	82.26	<b>83.40</b>
	College Biology	44.44	<b>47.05</b>	61.11	<b>68.06</b>	61.81	<b>64.93</b>	74.31	<b>77.43</b>	72.92	<b>77.95</b>	91.67	<b>92.36</b>
	College Medicine	36.99	<b>37.14</b>	50.29	<b>55.06</b>	57.80	<b>59.97</b>	67.05	<b>68.06</b>	63.58	<b>69.07</b>	<b>78.61</b>	78.18
	Medical Genetics	43.00	<b>44.75</b>	54.00	<b>66.00</b>	64.00	<b>65.50</b>	<b>80.00</b>	77.25	70.00	<b>78.00</b>	91.00	<b>91.00</b>
	Professional Medicine	29.78	<b>34.10</b>	50.37	<b>60.02</b>	56.99	<b>63.33</b>	<b>76.84</b>	75.64	72.06	<b>73.07</b>	<b>90.44</b>	90.26
	<b>Average</b>	36.69	<b>39.14</b>	49.76	<b>55.42</b>	54.65	<b>55.39</b>	<b>67.59</b>	67.05	61.08	<b>64.13</b>	81.14	<b>81.22</b>
French	FrenchMedMCQA	<b>29.91</b>	28.43	29.60	<b>40.27</b>	<b>45.48</b>	44.32	41.74	<b>44.63</b>	55.14	<b>58.02</b>	63.24	<b>73.05</b>
	MedExpQA	19.20	<b>20.60</b>	26.40	<b>39.20</b>	40.80	<b>41.20</b>	<b>48.00</b>	43.60	50.40	<b>56.00</b>	<b>76.80</b>	74.00
	Anatomy	<b>35.56</b>	35.18	48.15	<b>49.63</b>	33.33	<b>39.45</b>	45.19	<b>47.41</b>	55.56	<b>59.63</b>	67.41	<b>68.52</b>
	Clinical Knowledge	32.45	<b>36.51</b>	50.94	<b>57.92</b>	<b>55.47</b>	53.02	<b>61.89</b>	61.13	65.66	<b>71.51</b>	78.87	<b>80.56</b>
	College Biology	33.33	<b>38.02</b>	46.53	<b>52.78</b>	<b>53.47</b>	49.65	57.64	<b>62.50</b>	67.36	<b>72.92</b>	86.81	<b>87.67</b>
	College Medicine	32.95	<b>35.84</b>	43.93	<b>47.98</b>	<b>51.45</b>	48.56	57.80	<b>59.40</b>	57.80	<b>63.44</b>	69.94	<b>74.71</b>
	Medical Genetics	35.00	<b>40.00</b>	50.00	<b>57.25</b>	47.00	<b>59.00</b>	66.00	<b>67.00</b>	71.00	<b>72.00</b>	<b>90.00</b>	89.50
	Professional Medicine	24.26	<b>28.95</b>	33.09	<b>42.00</b>	43.38	<b>43.84</b>	51.47	<b>55.51</b>	59.56	<b>64.15</b>	72.79	<b>73.34</b>
	<b>Average</b>	30.33	<b>32.94</b>	41.08	<b>48.38</b>	46.30	<b>47.38</b>	53.72	<b>55.15</b>	60.31	<b>64.71</b>	75.73	<b>77.67</b>
Spanish	HeadQA	33.77	<b>34.32</b>	48.21	<b>54.47</b>	53.79	<b>55.54</b>	59.66	<b>61.24</b>	64.77	<b>68.00</b>	81.44	<b>82.44</b>
	MedExpQA	21.60	<b>23.00</b>	32.80	<b>38.40</b>	<b>46.40</b>	40.40	40.00	<b>43.00</b>	<b>52.80</b>	52.40	73.60	<b>76.60</b>
	Anatomy	37.78	<b>39.08</b>	42.22	<b>51.11</b>	45.93	<b>49.63</b>	48.15	<b>52.96</b>	60.74	<b>62.22</b>	71.11	<b>74.44</b>
	Clinical Knowledge	37.74	<b>38.78</b>	53.96	<b>55.47</b>	54.34	<b>56.13</b>	58.49	<b>62.08</b>	<b>68.68</b>	68.40	78.49	<b>80.00</b>
	College Biology	29.17	<b>35.94</b>	48.61	<b>50.35</b>	55.56	<b>56.25</b>	54.86	<b>55.04</b>	66.67	<b>69.10</b>	<b>85.42</b>	84.20
	College Medicine	32.37	<b>34.39</b>	43.93	<b>48.84</b>	<b>54.34</b>	48.99	49.71	<b>54.05</b>	<b>59.54</b>	58.24	69.94	<b>72.97</b>
	Medical Genetics	32.00	<b>34.75</b>	46.00	<b>59.50</b>	53.00	<b>57.25</b>	<b>72.00</b>	68.00	<b>67.00</b>	66.75	86.00	<b>86.75</b>
	Professional Medicine	26.47	<b>30.06</b>	38.24	<b>43.56</b>	<b>47.06</b>	45.68	<b>51.84</b>	50.74	53.68	<b>56.90</b>	<b>69.49</b>	68.94
	<b>Average</b>	31.36	<b>33.79</b>	44.25	<b>50.21</b>	<b>51.30</b>	51.23	54.34	<b>55.89</b>	61.74	<b>62.75</b>	76.94	<b>78.29</b>

## Part 3 – MedPodGPT

### ■ Benchmarks (zero-shot cross-lingual transfer performance)

Language	Benchmark Datasets	Model											
		Gemma 2B		Gemma 7B		Mistral 7B		LLaMA 3 8B		Mixtral MoE		LLaMA 3 70B	
		Baseline	<i>Ours</i>	Baseline	<i>Ours</i>	Baseline	<i>Ours</i>	Baseline	<i>Ours</i>	Baseline	<i>Ours</i>	Baseline	<i>Ours</i>
Chinese	MedQA-MCMLE	33.39	<b>33.43</b>	40.51	<b>45.98</b>	39.67	<b>39.25</b>	63.63	<b>66.32</b>	45.80	<b>47.14</b>	<b>84.68</b>	83.73
	Anatomy	<b>28.38</b>	23.98	25.00	<b>31.59</b>	25.00	<b>30.41</b>	33.78	<b>35.98</b>	<b>33.11</b>	26.35	63.51	<b>64.02</b>
	Clinical Knowledge	<b>29.11</b>	28.27	31.22	<b>39.87</b>	33.33	<b>32.60</b>	49.37	<b>50.95</b>	<b>39.24</b>	38.61	71.73	<b>71.94</b>
	College Medicine	28.94	<b>32.23</b>	33.70	<b>36.08</b>	<b>30.77</b>	30.68	52.01	<b>56.50</b>	38.46	<b>40.94</b>	75.82	<b>80.49</b>
	Medical Genetics	32.39	<b>32.39</b>	43.75	<b>45.17</b>	38.64	<b>42.33</b>	43.18	<b>44.60</b>	45.45	<b>45.88</b>	<b>61.36</b>	57.53
	Medical Nutrition	33.79	<b>35.69</b>	40.69	<b>44.66</b>	<b>42.07</b>	37.24	<b>53.10</b>	50.00	49.66	<b>51.90</b>	66.21	<b>68.28</b>
	Traditional Chinese Medicine	27.57	<b>28.52</b>	31.35	<b>36.35</b>	24.86	<b>28.52</b>	<b>43.24</b>	39.46	30.27	<b>30.94</b>	66.49	<b>67.98</b>
	Virology	<b>37.28</b>	36.98	46.15	<b>54.44</b>	43.79	<b>48.22</b>	<b>59.76</b>	58.88	<b>53.25</b>	50.15	76.33	<b>77.51</b>
	<b>Average</b>	31.36	<b>31.44</b>	36.55	<b>41.77</b>	34.77	<b>36.16</b>	49.76	<b>50.34</b>	<b>41.91</b>	41.49	70.77	<b>71.44</b>
Hindi	Anatomy	25.93	<b>32.22</b>	34.07	<b>36.86</b>	23.70	<b>30.00</b>	<b>40.00</b>	35.18	31.11	<b>34.44</b>	52.59	<b>57.78</b>
	Clinical Knowledge	26.42	<b>28.96</b>	<b>41.89</b>	41.04	24.91	<b>35.85</b>	<b>48.30</b>	46.70	<b>38.11</b>	36.70	63.40	<b>69.06</b>
	College Biology	26.39	<b>33.16</b>	26.39	<b>34.03</b>	19.44	<b>28.47</b>	32.65	<b>37.16</b>	<b>30.56</b>	32.81	58.33	<b>68.06</b>
	College Medicine	24.86	<b>27.60</b>	42.20	<b>43.35</b>	23.12	<b>33.09</b>	41.04	<b>43.64</b>	27.17	<b>33.24</b>	60.69	<b>64.74</b>
	Medical Genetics	<b>31.00</b>	30.50	36.00	<b>41.75</b>	28.00	<b>29.25</b>	<b>46.00</b>	45.75	40.00	<b>43.25</b>	71.00	<b>77.00</b>
	Professional Medicine	25.37	<b>26.19</b>	30.88	<b>41.08</b>	22.06	<b>28.67</b>	36.40	<b>39.34</b>	29.41	<b>29.50</b>	45.59	<b>64.70</b>
	<b>Average</b>	26.66	<b>29.77</b>	35.24	<b>39.69</b>	23.54	<b>30.89</b>	40.73	<b>41.29</b>	32.73	<b>34.99</b>	58.60	<b>66.89</b>



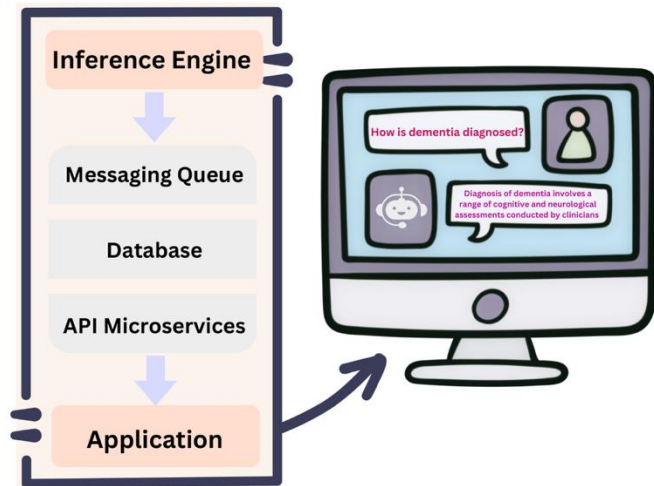
## Part 3 – MedPodGPT

### ■ Software development

- vLLM: inference engine
- Apache Cassandra: query optimization
- Flask: store chats and conversations in Cassandra and send text inference requests to a queue
- RabbitMQ: queue requests
- PostgreSQL: RAG database management
- OAuth 2.0: Authorization and user management

Credits:

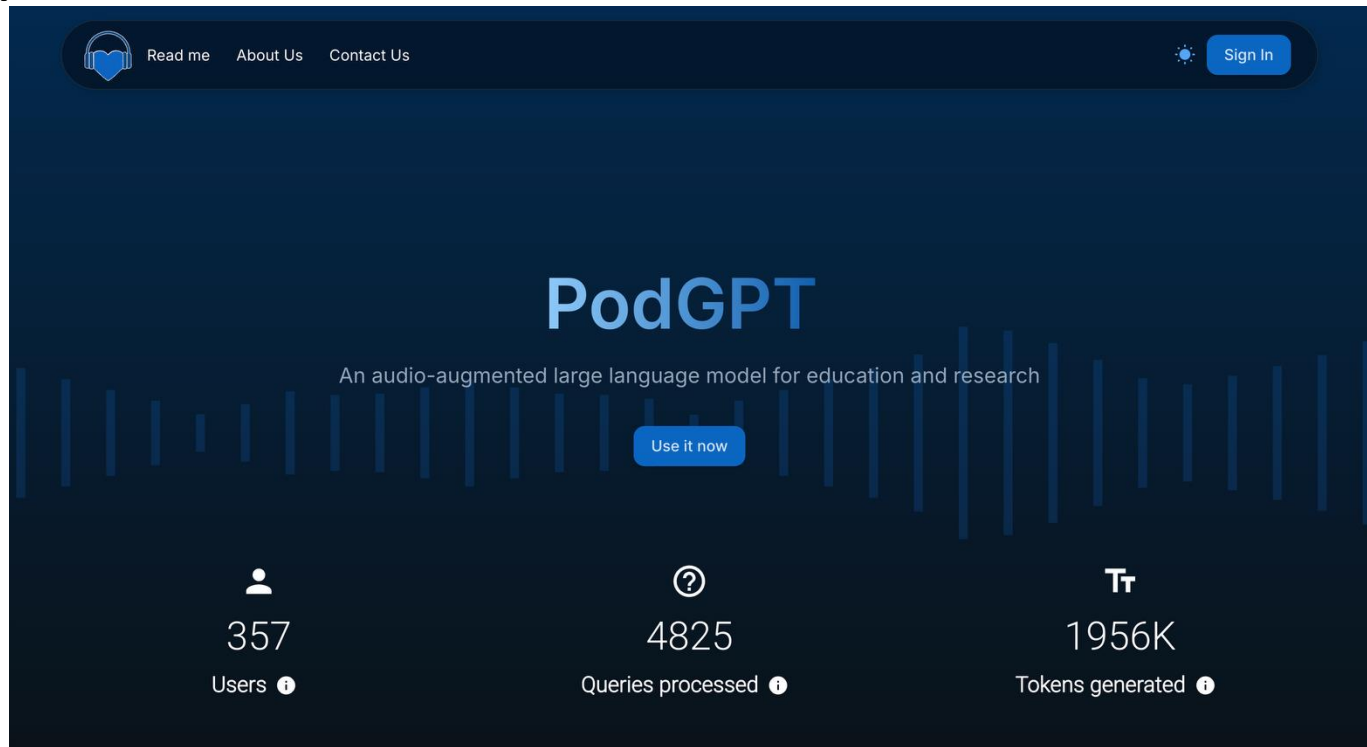
- [1] PodGPT: <https://podgpt.org/>
- [2] vLLM: <https://github.com/vllm-project/vllm>
- [3] Apache Cassandra: <https://github.com/apache/cassandra>
- [4] Flask: <https://github.com/pallets/flask>
- [5] RabbitMQ: <https://github.com/rabbitmq/rabbitmq-server>
- [6] PostgreSQL: <https://github.com/postgres/postgres>
- [7] OAuth 2.0: <https://github.com/postgres/postgres>



## Part 3 – MedPodGPT

- Software development

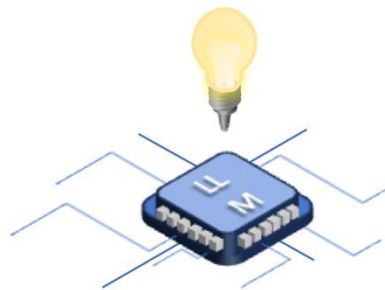
<https://podgpt.org>



### 1. Generating teaching cases



### 2. Leading morning report conferences



*Tip: Use LLMs to generate ideas and enhance your own creativity*

### 3. Preparing for journal club



### 4. Summarizing large volumes of feedback



### 5. Providing feedback on clinical documentation

- 1. Customized teaching cases:** “I’m teaching second-year medical students. Generate a case on a patient with a chronic obstructive pulmonary disease (COPD) exacerbation with diagnostic uncertainty about heart failure exacerbation.”
- 2. Expert discussion:** “I have a medical case in 4 parts. As our diagnostic expert, provide the problem representation, a prioritized differential, and your choice of next test for each part.”
- 3. Summarization and feedback:**
  - “Would my patient be included in this study?” or “Summarize methodological points of the study.”
  - “I need a 2-paragraph summary of these course evaluations. Include 2 direct quotes highlighting learner qualities.”
  - “I’ve written these patient instructions aiming for an 8th grade reading level. Determine their current reading level and offer feedback for improvement.”

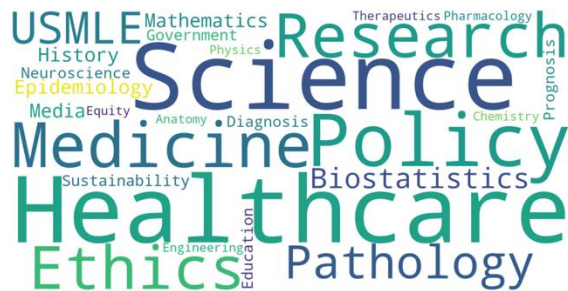
## Part 3 – MedPodGPT

### ■ Take away

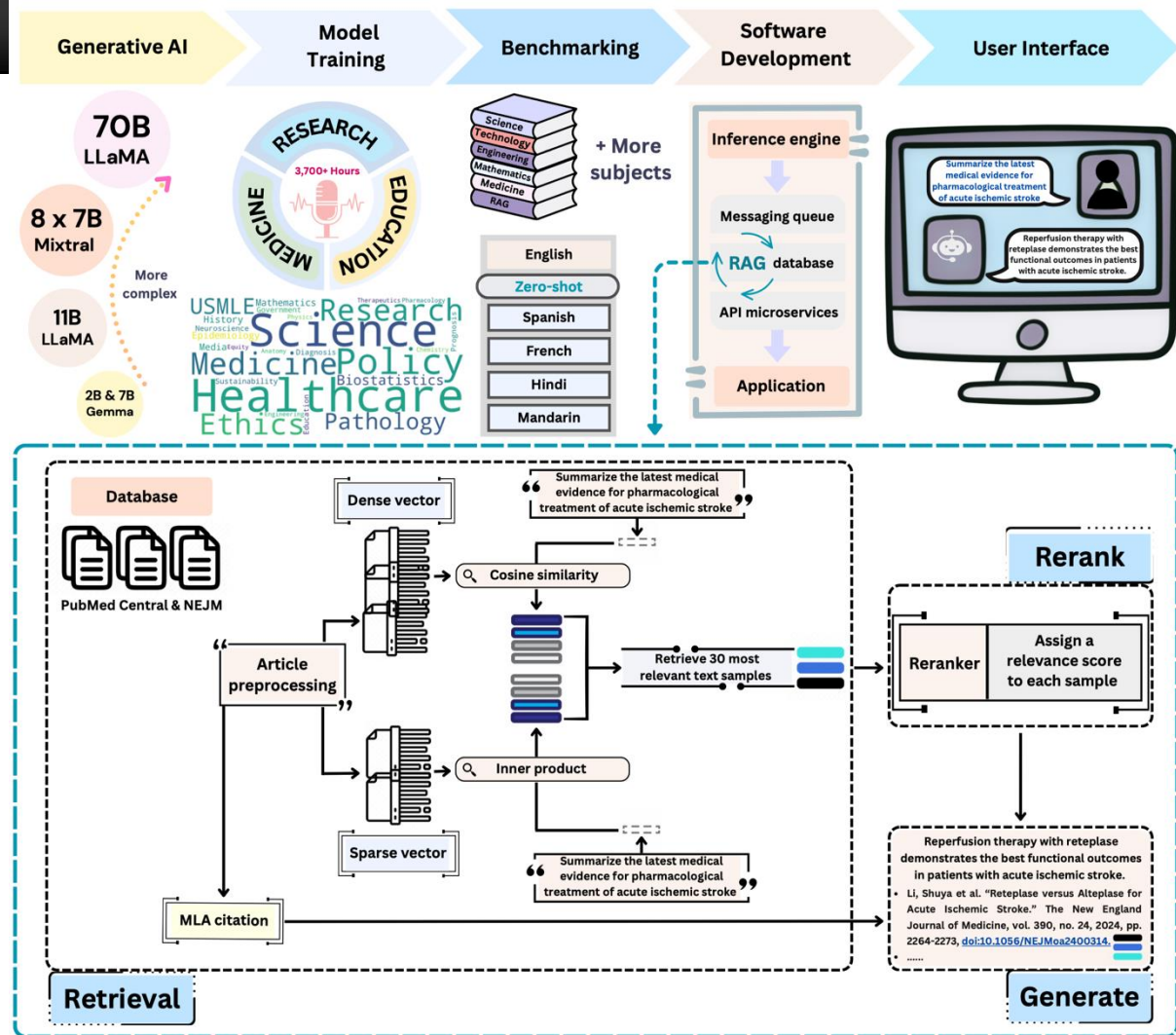
- **Podcasts** are valuable sources of **up to date** and **domain-relevant** information
- **Continual pretraining of LLMs on podcasts** can enhance overall **model performance**
- Improve **in-domain (language)** performance
- Improve **zero-shot multilingual transfer** capabilities

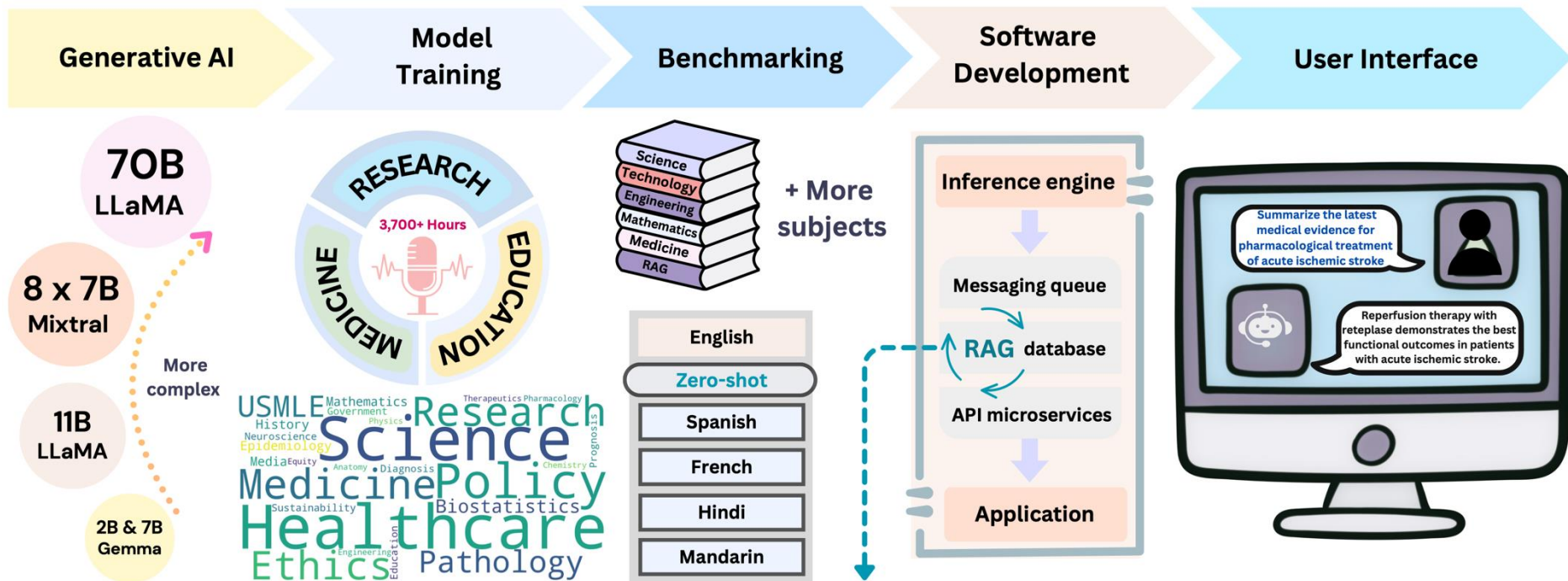
## Part 4 – PodGPT

- What: An **audio-augmented LLM** for research and education
  - Fields: Science, technology, engineering, mathematics, and medicine (STEMM)



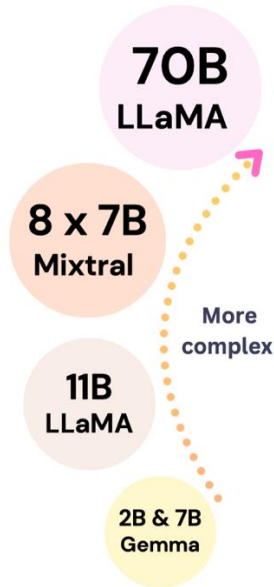
- Why: (1) **Up to date information** ← podcasts (2) **Factual information** ← journal articles
- Where: Creative Commons Attribution (CC BY) podcasts from journals/experts, and NEJM
- When: Most recent (e.g., since 2023)
- How: Continual pretraining of LLMs







## Part 4 – PodGPT



- Large language models
  - 2B Gemma
  - 7B Gemma, 7B Mistral
  - Quantized 70B LLaMA (11B)
  - 46B Mixtral 8×7B MoE
  - 70B LLaMA





# Part 4 – PodGPT

## ■ Podcasts

Podcast	Episodes	Audio time (min)	Mean length episode $\pm \sigma$ (min)	Number of text tokens	Mean text tokens per episode $\pm \sigma$
NEJM This Week	457	13,300.17	29.10 $\pm$ 2.92	2,029,219	4440.30 $\pm$ 419.35
NEJM Interviews	654	9,223.44	14.10 $\pm$ 8.15	1,732,879	2649.66 $\pm$ 1522.35
NEJM Core IM   Internal Medicine Podcast	170	5,285.72	31.09 $\pm$ 10.08	1,077,154	6336.20 $\pm$ 2093.18
NEJM Curbside Consults	74	1,977.46	26.72 $\pm$ 11.04	408,189	5516.07 $\pm$ 2522.46
NEJM Clinical Conversations	108	1,829.66	16.94 $\pm$ 4.13	320,968	2971.93 $\pm$ 774.89
NEJM Leadership Conversations	100	1,765.76	17.66 $\pm$ 4.38	306,490	3064.90 $\pm$ 759.09
NEJM AI Grand Rounds	24	1,459.11	60.80 $\pm$ 14.84	303,499	12645.79 $\pm$ 3107.35
NEJM Intention to Treat	40	997.22	24.93 $\pm$ 5.27	169,632	4240.80 $\pm$ 954.41
NEJM Not Otherwise Specified	20	836.03	41.80 $\pm$ 15.71	146,718	7335.90 $\pm$ 2880.36
TWiV: This Week In Virology	1,186	104,089.57	87.77 $\pm$ 30.30	20,188,268	17022.15 $\pm$ 5785.86
TWiP: This Week In Parasitism	245	21,129.03	86.24 $\pm$ 15.35	4,381,511	17883.72 $\pm$ 3814.75
TWiM: This Week in Microbiology	320	20,641.50	64.50 $\pm$ 10.36	3,667,519	11461.00 $\pm$ 2121.78
TWiEVO: This Week In Evolution	100	8,845.09	88.45 $\pm$ 10.99	1,756,480	17564.80 $\pm$ 2517.80
IMMUNE	93	6,918.84	74.40 $\pm$ 19.67	1,363,176	14657.81 $\pm$ 3888.47
TWiN: This Week In Neuroscience	53	3,581.13	67.57 $\pm$ 10.40	644,712	12164.38 $\pm$ 2078.49
Matters Microbial	62	3,418.33	55.13 $\pm$ 11.33	637,647	10284.63 $\pm$ 2357.42
Infectious Disease Puscast	65	2,298.71	35.36 $\pm$ 5.96	415,145	6386.85 $\pm$ 1126.15
Urban Agriculture	29	2,171.63	74.88 $\pm$ 17.58	443,859	15305.48 $\pm$ 4214.82
On The Wards: On The Pods Medical Podcast for Doctors	245	5,915.14	24.14 $\pm$ 8.00	1,175,307	4797.17 $\pm$ 1781.74
Digital Campus Podcast	64	3,169.44	49.52 $\pm$ 5.85	605,977	9322.72 $\pm$ 1528.99
emDOCs.net Emergency Medicine (EM) Podcast	112	1,576.72	14.08 $\pm$ 4.30	303,672	2711.35 $\pm$ 864.29
Policy in Plain English Podcast	73	1,089.20	14.92 $\pm$ 7.52	208,580	2857.26 $\pm$ 1509.80
Open Minds ... from Creative Commons	21	803.28	38.25 $\pm$ 12.79	145,624	6619.27 $\pm$ 2722.90
What is Global Health?	18	486.47	27.03 $\pm$ 10.02	87,653	4869.61 $\pm$ 2054.00
Consilience Sustainability In Progress (SIP) Podcast	9	403.95	44.88 $\pm$ 17.89	70,461	7829.00 $\pm$ 3535.75
Research Pulse: Future Focussed Health Insights	16	177.79	11.11 $\pm$ 2.29	33,672	2104.50 $\pm$ 476.85
Our People: Central to Healthcare	9	161.15	17.91 $\pm$ 7.42	31,545	3505.00 $\pm$ 1474.57

## Part 4 – PodGPT

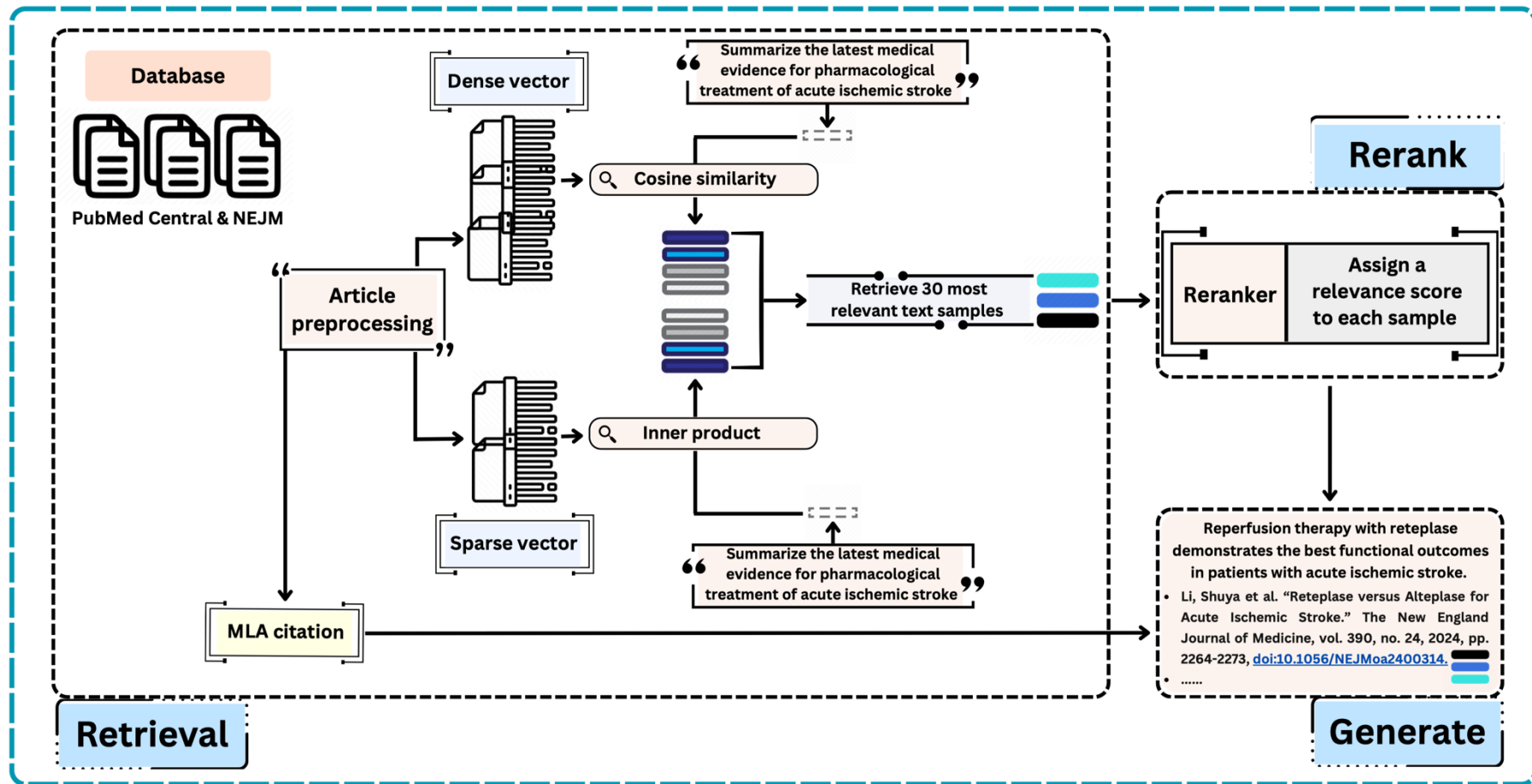
### ■ Benchmarks (in-domain performance)

Model	Gemma 2B		Gemma 7B		Quantized LLaMA 70B		Mixtral 8×7B MoE		LLaMA 70B	
MMLU benchmark	Baseline	<i>Ours</i>	Baseline	<i>Ours</i>	Baseline	<i>Ours</i>	Baseline	<i>Ours</i>	Baseline	<i>Ours</i>
Physics	29.85	<b>30.74 (0.12)</b>	42.38	<b>47.64 (0.64)</b>	76.79	<b>79.24</b>	56.08	<b>56.19 (0.71)</b>	78.63	<b>79.26</b>
Biology	44.82	<b>46.58 (0.45)</b>	63.32	<b>66.43 (0.58)</b>	90.97	<b>91.51</b>	72.14	<b>77.84 (0.90)</b>	93.53	<b>93.53</b>
Chemistry	30.04	<b>30.96 (0.62)</b>	43.42	<b>44.14 (1.31)</b>	69.94	<b>73.15</b>	46.60	<b>50.33 (1.34)</b>	<b>70.16</b>	69.66
Computer Science	40.59	<b>44.20 (0.11)</b>	53.58	<b>54.62 (0.29)</b>	80.24	<b>82.30</b>	59.25	<b>60.98 (1.25)</b>	78.72	<b>79.44</b>
Engineering	39.31	<b>42.07 (0.49)</b>	44.14	<b>47.94 (0.60)</b>	73.79	<b>73.79</b>	<b>57.24</b>	56.38 (0.90)	75.17	<b>75.17</b>
Mathematics	25.69	<b>26.44 (0.57)</b>	34.97	<b>39.23 (0.16)</b>	68.83	<b>71.34</b>	<b>47.42</b>	46.59 (1.01)	63.97	<b>64.18</b>
Medicine	40.62	<b>41.72 (0.15)</b>	55.22	<b>59.50 (0.14)</b>	86.11	<b>87.38</b>	67.38	<b>74.00 (0.71)</b>	88.65	<b>88.65</b>
<b>Average</b>	35.85	<b>37.53</b>	48.15	<b>51.36</b>	78.10	<b>79.82</b>	58.02	<b>60.33</b>	78.40	<b>78.56</b>

# ■ Benchmarks

## (zero-shot transfer)

Language	Benchmark datasets	Model									
		Gemma 2B		Gemma 7B		Quantized LLaMA 70B		Mixtral 8×7B MoE		LLaMA 70B	
		Baseline	<i>Ours</i>	Baseline	<i>Ours</i>	Baseline	<i>Ours</i>	Baseline	<i>Ours</i>	Baseline	<i>Ours</i>
Mandarin	MedQA-MCMLLE	<b>33.54</b>	33.14 (0.13)	40.81	<b>45.20 (0.06)</b>	82.69	<b>84.15</b>	<b>45.94</b>	45.49 (0.31)	<b>86.14</b>	86.02
	Physics	<b>32.75</b>	31.73 (0.37)	35.37	<b>40.13 (0.54)</b>	<b>61.08</b>	60.53	<b>42.72</b>	38.90 (1.42)	<b>62.42</b>	61.92
	Biology	24.26	<b>25.44 (0.00)</b>	33.14	<b>37.43 (0.64)</b>	62.72	<b>63.91</b>	<b>44.38</b>	39.94 (2.15)	59.76	<b>60.36</b>
	Chemistry	25.00	<b>27.08 (0.32)</b>	29.55	<b>37.12 (0.54)</b>	<b>45.45</b>	43.94	29.55	<b>34.47 (0.85)</b>	46.97	<b>47.73</b>
	Computer Science	32.24	<b>34.18 (0.32)</b>	40.78	<b>47.92 (0.59)</b>	73.60	<b>76.82</b>	<b>57.16</b>	56.20 (2.42)	78.46	<b>78.46</b>
	Engineering	<b>33.40</b>	32.05 (0.40)	41.72	<b>45.58 (0.20)</b>	62.48	<b>63.72</b>	42.30	<b>46.75 (0.93)</b>	<b>66.92</b>	66.16
	Mathematics	<b>26.82</b>	24.64 (0.65)	27.96	<b>30.73 (0.81)</b>	<b>53.72</b>	51.93	<b>38.21</b>	30.41 (0.90)	<b>58.89</b>	57.90
	Medicine	31.18	<b>31.62 (0.15)</b>	35.72	<b>39.65 (0.18)</b>	70.06	<b>71.21</b>	41.06	<b>41.77 (0.74)</b>	<b>73.42</b>	73.29
	<b>Average</b>	29.90	<b>29.98</b>	35.63	<b>40.47</b>	63.98	<b>64.53</b>	<b>42.66</b>	41.74	<b>66.62</b>	66.48
French	MedExpQA	19.20	<b>22.40 (0.00)</b>	28.00	<b>37.40 (1.18)</b>	<b>78.40</b>	77.60	49.60	<b>52.80 (1.50)</b>	77.60	<b>77.60</b>
	FrenchMedMCQA	<b>31.46</b>	28.04 (0.22)	33.64	<b>43.69 (0.34)</b>	81.62	<b>84.11</b>	57.01	<b>64.64 (0.47)</b>	87.23	<b>87.85</b>
	Physics	26.93	<b>28.88 (0.23)</b>	35.58	<b>42.65 (0.59)</b>	<b>71.28</b>	70.31	55.56	<b>56.51 (0.51)</b>	<b>72.40</b>	72.29
	Biology	34.51	<b>39.35 (0.39)</b>	50.16	<b>57.04 (0.34)</b>	89.29	<b>91.01</b>	70.03	<b>73.52 (1.01)</b>	90.90	<b>90.90</b>
	Chemistry	26.07	<b>29.09 (0.32)</b>	37.98	<b>40.20 (0.53)</b>	<b>64.94</b>	64.93	47.14	<b>49.73 (0.45)</b>	<b>65.18</b>	64.43
	Computer Science	35.46	<b>37.40 (0.69)</b>	44.20	<b>45.23 (0.70)</b>	75.85	<b>75.88</b>	<b>57.54</b>	56.43 (0.41)	<b>76.49</b>	76.24
	Engineering	<b>38.62</b>	36.20 (0.60)	46.90	<b>47.59 (0.00)</b>	68.28	<b>72.41</b>	53.79	<b>54.31 (0.57)</b>	71.72	<b>71.72</b>
	Mathematics	26.14	<b>28.01 (0.51)</b>	30.62	<b>32.92 (0.98)</b>	<b>64.87</b>	63.90	44.65	<b>46.68 (1.14)</b>	65.63	<b>65.92</b>
	Medicine	32.56	<b>35.48 (0.12)</b>	45.80	<b>51.55 (0.18)</b>	81.36	<b>83.44</b>	65.08	<b>66.49 (0.29)</b>	84.21	<b>84.29</b>
	<b>Average</b>	30.11	<b>31.65</b>	39.21	<b>44.25</b>	75.10	<b>75.95</b>	55.60	<b>57.90</b>	<b>76.82</b>	76.80
Hindi	Physics	25.49	<b>26.60 (0.56)</b>	29.32	<b>33.39 (0.46)</b>	<b>58.48</b>	55.90	<b>32.06</b>	31.14 (1.26)	58.68	<b>59.12</b>
	Biology	29.02	<b>32.58 (0.18)</b>	29.28	<b>39.08 (0.66)</b>	<b>66.55</b>	66.46	<b>35.69</b>	34.03 (0.74)	<b>72.74</b>	72.61
	Chemistry	<b>24.08</b>	20.88 (0.10)	35.26	<b>36.97 (0.90)</b>	50.84	<b>51.05</b>	<b>33.98</b>	30.25 (2.59)	<b>54.04</b>	53.29
	Computer Science	<b>32.15</b>	30.30 (0.40)	36.64	<b>41.49 (0.37)</b>	<b>65.78</b>	61.22	<b>37.20</b>	36.40 (2.60)	66.45	<b>66.45</b>
	Engineering	<b>43.45</b>	42.42 (0.34)	40.00	<b>41.72 (1.04)</b>	<b>59.31</b>	57.93	40.00	<b>43.28 (1.23)</b>	57.24	<b>58.62</b>
	Mathematics	<b>25.33</b>	24.87 (0.24)	29.33	<b>30.96 (0.50)</b>	<b>53.19</b>	48.35	<b>33.74</b>	33.10 (0.47)	<b>54.60</b>	52.82
	Medicine	26.77	<b>29.07 (0.15)</b>	34.00	<b>40.26 (0.21)</b>	64.60	<b>65.08</b>	<b>33.98</b>	33.43 (0.80)	71.04	<b>71.05</b>
	<b>Average</b>	29.47	<b>29.53</b>	33.40	<b>37.70</b>	<b>59.82</b>	58.00	<b>35.24</b>	34.52	<b>62.11</b>	61.99
Spanish	HEAD-QA	33.66	<b>34.38 (0.10)</b>	48.32	<b>52.87 (0.19)</b>	81.58	<b>82.97</b>	64.48	<b>66.76 (0.50)</b>	<b>84.28</b>	84.21
	MedExpQA	21.60	<b>23.20 (0.00)</b>	32.80	<b>37.40 (0.35)</b>	76.80	<b>80.00</b>	51.20	<b>54.00 (1.65)</b>	76.80	<b>78.40</b>
	Physics	28.06	<b>28.86 (0.26)</b>	40.64	<b>43.63 (0.11)</b>	70.20	<b>70.92</b>	53.32	<b>54.84 (0.67)</b>	71.59	<b>72.13</b>
	Biology	30.63	<b>39.50 (0.55)</b>	52.19	<b>57.26 (0.51)</b>	85.37	<b>85.46</b>	66.46	<b>71.11 (0.47)</b>	<b>86.73</b>	86.50
	Chemistry	<b>27.06</b>	25.94 (0.22)	35.98	<b>39.93 (0.36)</b>	<b>63.72</b>	61.72	48.36	<b>48.74 (1.00)</b>	<b>67.18</b>	65.93
	Computer Science	37.09	<b>40.43 (0.40)</b>	45.93	<b>48.04 (0.25)</b>	73.38	<b>74.30</b>	<b>56.74</b>	56.71 (0.78)	74.77	<b>75.72</b>
	Engineering	<b>43.45</b>	38.79 (0.75)	47.59	<b>49.14 (0.89)</b>	73.10	<b>73.10</b>	<b>55.17</b>	53.79 (1.89)	<b>72.41</b>	71.72
	Mathematics	<b>26.63</b>	25.80 (0.14)	31.14	<b>32.79 (0.39)</b>	<b>64.80</b>	61.97	<b>42.04</b>	41.70 (0.72)	<b>67.08</b>	66.47
	Medicine	31.89	<b>35.54 (0.15)</b>	45.94	<b>51.81 (0.44)</b>	78.28	<b>79.64</b>	61.17	<b>63.25 (0.90)</b>	<b>80.19</b>	80.09
	<b>Average</b>	31.12	<b>32.49</b>	42.28	<b>45.87</b>	74.14	<b>74.45</b>	55.44	<b>56.77</b>	75.67	<b>75.69</b>



# Part 4 – PodGPT

## ■ RAG Database

Journal	Number of articles	Number of text samples	Average length per sample
JAMA Network Open	9,367	119,672	389.61
The New England Journal of Medicine	2,013	33,344	400.33
Cell	497	24,977	417.86
British Medical Journal	601	10,307	399.27
The Lancet	458	9,425	403.77
The Lancet Global Health	539	9,003	399.56
Neurology	428	6,862	404.43
JAMA Health Forum	524	5,524	390.70
The Lancet Regional Health - Europe	337	5,435	398.94
The Lancet Infectious Diseases	272	4,720	400.31
The Lancet Regional Health – Western Pacific	243	3,267	397.28
The Lancet Public Health Home	179	3,059	401.88
JAMA Psychiatry	180	2,859	399.22
JAMA Neurology	141	2,462	397.50
The Lancet Oncology	116	2,304	400.62
The Lancet Regional Health – Americas	174	2,262	395.02
The Lancet Microbe	111	2,054	403.17
The Lancet Psychiatry	119	1,951	400.77
The Lancet Neurology	110	1,841	396.53
JAMA Oncology	138	1,810	394.91
The Lancet HIV	97	1,786	401.54
The Lancet Planetary Health	100	1,776	400.43
JAMA Pediatrics	135	1,773	391.99
JAMA Internal Medicine	128	1,766	395.25
The Lancet Regional Health – Southeast Asia	134	1,607	397.18
The Lancet Respiratory Medicine	72	1,360	404.70
The Lancet Child & Adolescent Health	68	1,205	402.95
The Lancet Diabetes & Endocrinology	64	1,082	399.89
The Lancet Gastroenterology and Hepatology	54	1,007	397.67
JAMA Cardiology	60	947	396.69
JAMA Surgery	68	941	399.31
The Lancet Healthy Longevity	49	855	403.10
JAMA Ophthalmology	56	810	398.01
The Lancet Rheumatology	35	713	403.76
The Lancet Haematology	33	604	398.36
JAMA Dermatology	30	413	399.73
JAMA Otolaryngology – Head & Neck Surgery	17	244	399.27
JAMA Facial Plastic Surgery	1	10	424.80

# Part 4 – PodGPT

## ■ Benchmarks (in-domain performance)

Benchmark datasets			MedExpQA	MedMCQA	MedQA	PubMedQA	MMLU Medicine	Average
Model	Gemma 2B	Baseline	19.20	<b>34.71</b>	29.54	46.80	40.62	34.17
		<i>Ours</i>	<b>21.20 (0.69)</b>	34.62 (0.02)	<b>32.91 (0.15)</b>	<b>54.25 (0.54)</b>	<b>41.72 (0.15)</b>	<b>36.94</b>
		<i>Baseline+RAG</i>	23.20	35.91	32.60	49.00	41.47	36.44
		<i>Ours+RAG</i>	<b>28.00 (0.00)</b>	<b>35.96 (0.07)</b>	<b>34.43 (0.12)</b>	<b>51.95 (0.78)</b>	<b>42.12 (0.13)</b>	<b>38.49</b>
	Gemma 7B	Baseline	34.40	40.69	37.78	<b>61.80</b>	55.22	45.98
		<i>Ours</i>	<b>42.00 (0.89)</b>	<b>44.64 (0.09)</b>	<b>44.14 (0.21)</b>	57.35 (1.37)	<b>59.50 (0.14)</b>	<b>49.53</b>
		<i>Baseline+RAG</i>	35.20	40.64	39.28	<b>61.40</b>	54.14	46.13
		<i>Ours+RAG</i>	<b>47.40 (1.18)</b>	<b>43.54 (0.07)</b>	<b>43.32 (0.25)</b>	55.70 (1.88)	<b>60.11 (0.39)</b>	<b>50.01</b>
	Quantized LLaMA 70B	Baseline	73.60	68.30	76.67	77.60	86.11	76.46
		<i>Ours</i>	<b>75.20</b>	<b>69.11</b>	<b>78.00</b>	<b>77.60</b>	<b>87.38</b>	<b>77.46</b>
		<i>Baseline+RAG</i>	69.60	66.51	66.77	76.80	83.91	72.72
		<i>Ours+RAG</i>	<b>76.00</b>	<b>68.13</b>	<b>69.60</b>	<b>78.80</b>	<b>85.62</b>	<b>75.63</b>
	Mixtral 8×7B MoE	Baseline	52.80	50.11	54.20	<b>42.00</b>	67.38	<b>53.30</b>
		<i>Ours</i>	<b>58.60 (1.54)</b>	<b>51.94 (0.31)</b>	<b>55.36 (0.81)</b>	25.35 (0.78)	<b>74.00 (0.71)</b>	53.05
		<i>Baseline+RAG</i>	49.60	48.58	47.84	<b>34.60</b>	65.45	49.21
		<i>Ours+RAG</i>	<b>57.12 (1.21)</b>	<b>51.29 (0.28)</b>	<b>51.14 (0.01)</b>	24.44 (1.11)	<b>72.01 (0.02)</b>	<b>51.20</b>
	LLaMA 70B	Baseline	75.20	70.31	<b>78.40</b>	76.20	88.65	77.75
		<i>Ours</i>	<b>75.20</b>	<b>70.62</b>	78.32	<b>76.40</b>	<b>88.65</b>	<b>77.84</b>
		<i>Baseline+RAG</i>	68.00	68.11	66.93	77.40	84.24	72.94
		<i>Ours+RAG</i>	<b>71.20</b>	<b>68.59</b>	<b>68.42</b>	<b>77.40</b>	<b>85.60</b>	<b>74.24</b>





Summarize the latest medical evidence for the effectiveness of crinecerfont in adult congenital adrenal hyperplasia.

Crinecerfont therapy allowed substantial and clinically meaningful glucocorticoid reduction to more physiological doses in adults with classic congenital adrenal hyperplasia and was associated with reduced adrenal androgen production. The treatment appeared to be acceptably tolerated, with similar incidences of treatment-emergent adverse events in both the crinecerfont and placebo groups, and achieved androstenedione control at week 24 in 74.6% of crinecerfont-treated participants compared with 52.6% with placebo.

- Auchus, Richard J., et al. "Phase 3 Trial of Crinecerfont in Adult Congenital Adrenal Hyperplasia." *The New England Journal of Medicine*, vol. 391, no. 6, 2024, pp. 504-514. [doi:10.1056/NEJMoa2404656](https://doi.org/10.1056/NEJMoa2404656). Score: 0.51



What is the association between gestational diabetes mellitus and cardiovascular disease?

Women with a history of gestational diabetes mellitus are at substantially higher risk of future cardiovascular and cerebrovascular diseases overall, with a pooled risk ratio of 1.72. This excess risk cannot be solely attributed to conventional cardiovascular risk factors, and is partially mediated by subsequent diabetes, highlighting the need for early intervention and continuous monitoring in women with gestational diabetes mellitus.

- Xie, Wenhui, et al. "Association of Gestational Diabetes Mellitus with Overall and Type Specific Cardiovascular and Cerebrovascular Diseases: Systematic Review and Meta-Analysis." *BMJ*, vol. 378, 21 Sep. 2022, e070244. [doi:10.1136/bmj-2022-070244](https://doi.org/10.1136/bmj-2022-070244). Score: 0.61
- Echouffo Tcheguigui, Justin B., et al. "Association of Concomitant Gestational Hypertensive Disorders and Gestational Diabetes With Cardiovascular Disease." *JAMA Network Open*, vol. 5, no. 11, 1 Nov. 2022, e2243618. [doi:10.1001/jamanetworkopen.2022.43618](https://doi.org/10.1001/jamanetworkopen.2022.43618). Score: 0.49



Summarize the latest medical evidence for monoclonal antibody L9LS to prevent malaria.

The monoclonal antibody L9LS has shown to be effective in preventing malaria, with 88% of participants protected after controlled human malaria infection, and protection was seen at serum concentrations as low as 9.2 µg per milliliter. A single low dose of L9LS administered subcutaneously may be a feasible strategy to limit malaria-associated morbidity and mortality among infants and young children in regions where malaria is endemic.

- Kayentao, Kassoum, et al. "Subcutaneous Administration of a Monoclonal Antibody to Prevent Malaria." *The New England Journal of Medicine*, vol. 390, no. 17, 2024, pp. 1549-1559. [doi:10.1056/NEJMoa2312775](https://doi.org/10.1056/NEJMoa2312775). Score: 0.52
- Wu, Richard L., et al. "Low-Dose Subcutaneous or Intravenous Monoclonal Antibody to Prevent Malaria." *The New England Journal of Medicine*, vol. 387, no. 5, 2022, pp. 397-407. [doi:10.1056/NEJMoa2207007](https://doi.org/10.1056/NEJMoa2207007). Score: 0.60



## Explain the role of GPC3-Unc5 receptors.

GPC3-Unc5 receptors play a crucial role in regulating collective migration and correct targeting, and their signaling must be finely balanced to achieve effective migration. The interaction between GPC3 and Unc5 receptors elicits a repulsive cellular response, and disrupting or enhancing this interaction can impact various cellular processes, including apoptosis, migration, and survival, as seen in embryonal tumors and neuroblastoma samples.

- Akkermans, Onno, et al. "GPC3-Unc5 Receptor Complex Structure and Role in Cell Migration." *Cell*, vol. 185, no. 21, 2022, pp. 3931-3949.e26, [doi:10.1016/j.cell.2022.09.025](https://doi.org/10.1016/j.cell.2022.09.025).

Score: 0.55



## Are there disparities in household air pollution exposure between urban and rural areas?

There are disparities in household air pollution exposure between urban and rural areas, with exposures due to polluting fuels generally being much higher in rural settings than in urban settings. The average measured 24-hour HAP-PM<sub>2.5</sub> personal exposure for users of biomass fuel with traditional stoves is 130  $\mu\text{g}/\text{m}^3$  in rural settings and 80  $\mu\text{g}/\text{m}^3$  in urban settings, indicating a significant difference in exposure levels between the two settings.

- Mohajeri, Nahid, et al. "Urban-Rural Disparity in Global Estimation of PM<sub>2.5</sub> Household Air Pollution and Its Attributable Health Burden." *The Lancet. Planetary Health*, vol. 7, no. 8, 2023, pp. e660-e672, [doi:10.1016/S2542-5196\(23\)00133-X](https://doi.org/10.1016/S2542-5196(23)00133-X).

Score: 0.55



## What is the role of amygdala neurons in social decision making in primates?

The amygdala neurons in primates play a crucial role in social decision making by encoding object-specific values learned from social observation and own experience, and translating these values into representations of the partner monkey's forthcoming choices. This is evident from the fact that amygdala damage profoundly impairs primates' social behavior, and that distinct "simulation neurons" in the amygdala dynamically translate values into representations of the partner monkey's choices, supporting understanding of others' mental states.

- Grabenhorst, Fabian et al. "Primate Amygdala Neurons Simulate Decision Processes of Social Partners." *Cell* vol. 177,4 (2019): 986-998.e15. [doi:10.1016/j.cell.2019.02.042](https://doi.org/10.1016/j.cell.2019.02.042)

Score: 0.59



## Part 4 – PodGPT

### ■ Performance on conversational dataset

Model		Gemma 2B		Gemma 7B		Quantized LLaMA 70B		Mixtral 8×7B MoE		LLaMA 70B	
		Baseline	<i>Ours</i>	Baseline	<i>Ours</i>	Baseline	<i>Ours</i>	Baseline	<i>Ours</i>	Baseline	<i>Ours</i>
Perplexity	Original transcripts	17.54	<b>10.15 (0.01)</b>	27.04	<b>8.75 (0.02)</b>	7.34	<b>7.27</b>	<b>6.11</b>	6.32 (0.04)	7.31	<b>7.10</b>
	Our transcripts	14.08	<b>7.64 (0.02)</b>	22.72	<b>6.67 (0.02)</b>	6.12	<b>5.51</b>	5.23	<b>5.05 (0.01)</b>	6.10	<b>5.55</b>

$$\text{PPL}(\mathbf{x}) = e^{-\frac{1}{t} \sum_i^t \log(p_{\theta}(x_i | \mathbf{x}_{<i}))}.$$

## Part 4 – PodGPT

### ■ Take away

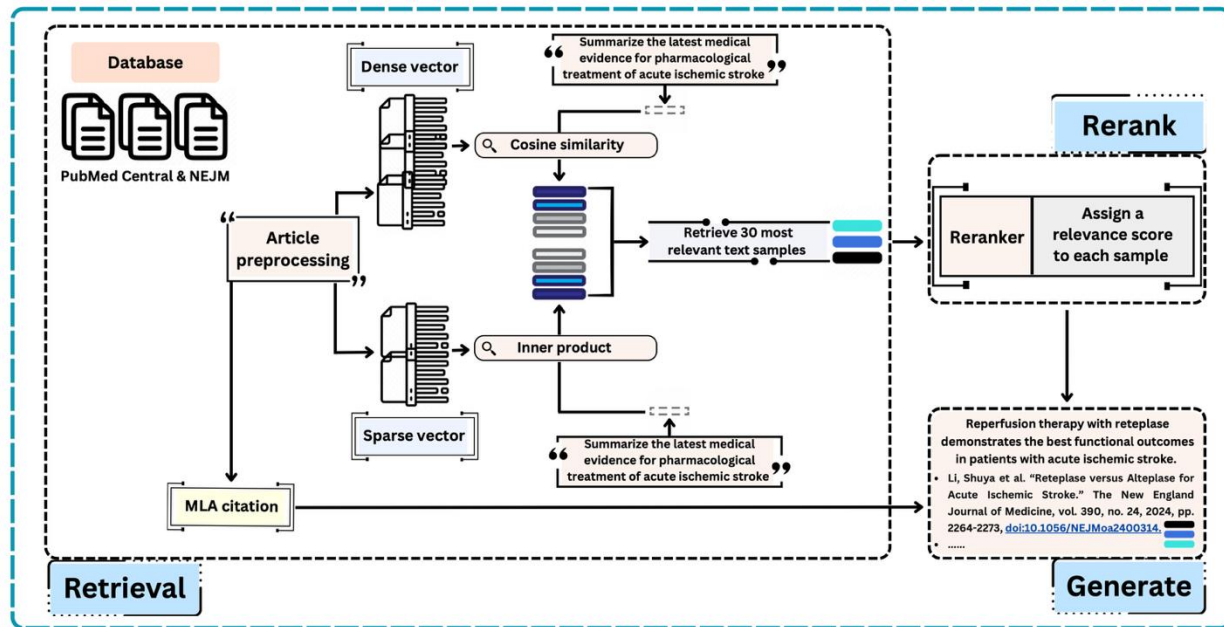
- **Podcasts** are valuable resources for **STEMM research and education**
- **Grounding LLMs with RAG** enhances **factual accuracy and reliability**
- **Continual pretraining of LLMs on podcasts** improves their **conversational capability**
- **Medical journals** are valuable resources for **evidence-based AI generation**

## Part 5 – Agentic System

### ■ Agentic memory-augmented retrieval and evidence grounding in medicine

- Retrieval
- Reranking
- Generation

- Unified
- Automatic
- Dynamic





## Multi-choice QA



## Open-ended QA

### Option Generation

Generate the most likely answers based on the patient's specific condition and needs.



### Option Comparison

A structured comparison of provided options.



### Evidence Search

Answer 'yes' or 'no' to indicate search necessity.

Yes

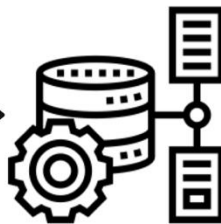
### Cache-and-prune Memory Bank Mechanism



#### Relevance Analysis



#### Evidence Grounding



Initial Memory  $M_1$

Grounded Evidence

Updated memory  $M_i$

$$\times \left\lceil \frac{R}{B} \right\rceil$$

Clinical Cases



Article Abstracts



Articles Textbooks



Clinical Trials



Encyclopedia



TopR  
Re-ranking



TopK  
Retrieval

No



## Final Diagnosis



### Aid Kit



generate\_options



perform\_comparison



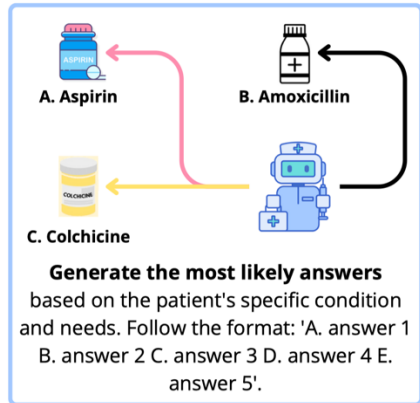
enable\_search



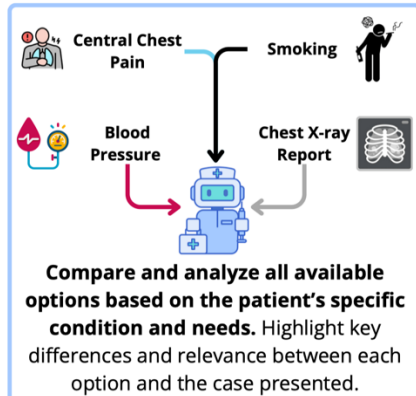
relevance\_analysis



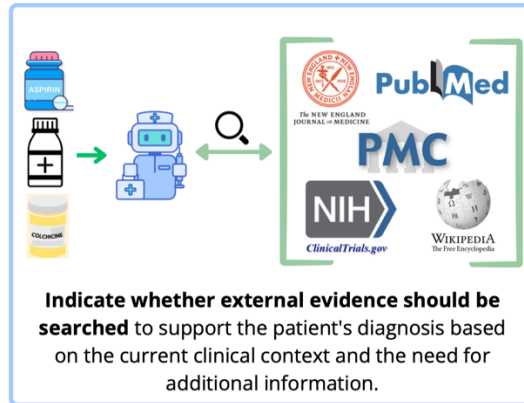
locate\_evidence



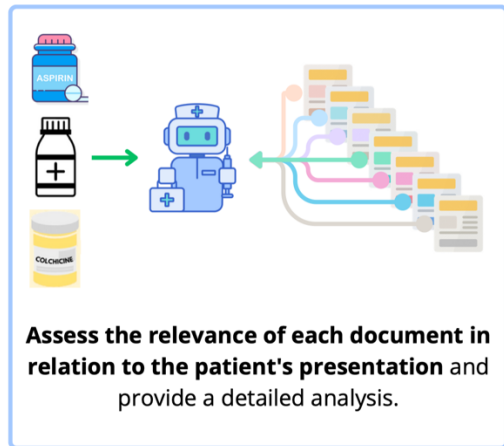
**generate\_options**



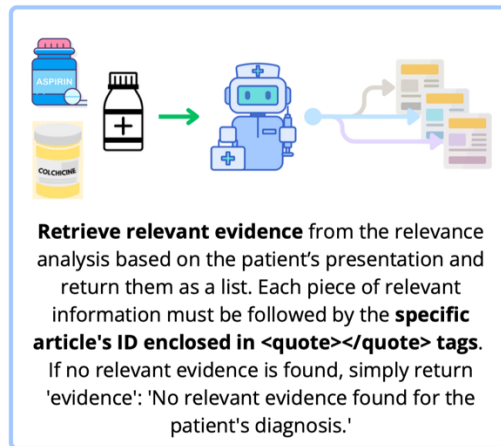
**perform\_comparison**



**enable\_search**



**relevance\_analysis**



**locate\_evidence**

---

**Algorithm 1** Agentic memory-augmented retrieval and evidence grounding system

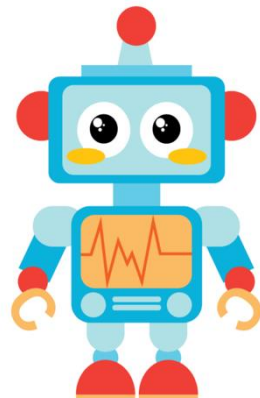
---

```
1: Initialize Document Retriever  $\phi$ , Evidence Reranker  $\psi$ 
2: Initialize AI Agent  $\pi$ , Conversation  $C$ , Memory Bank  $M_1$ 
3: Initialize Evidence database  $\mathcal{V}$ 
4: Given patient background and question  $Q$ , instructions  $I$ , tools  $T$ 
5: AI Agent  $\pi$  generates initial response  $\prod_t \pi(y_t | T, Q, I, \mathbf{y}_{<t})$ 
6: while tool calling do
7:   Retrieve content from the tool calling to update conversation  $C$ 
8:   if tool calling is enable_search then
9:     Retrieve TopK documents  $\arg \text{TopK}_{\mathbf{v} \in \mathcal{V}} - \|\phi(\mathbf{x}) - \phi(\mathbf{v})\|_2$ 
10:    Rerank TopR documents  $\arg \text{TopR}_{\mathbf{k} \in \mathcal{K}} \cos(\psi(\mathbf{x}), \psi(\mathbf{k}))$ 
11:    while  $i \leq \lceil R/B \rceil$  do
12:      Retrieve  $\mathcal{B}_i$  (a batch of  $\mathcal{R}$ ) to update conversation  $C$ 
13:      if tool calling is locate_evidence then
14:        if Relevant document is grounded within <quote></quote> tags then
15:          Update memory bank  $M_i = \text{Prune}(M_{i-1} \cup \mathcal{B}_i)$ 
16:        end if
17:      end if
18:      Remove  $\mathcal{B}_i$  from conversation  $C$ 
19:    end while until Sufficient information is gathered
20:  end if
21: end while
22: if  $M_i$  then
23:   return Final diagnosis  $\prod_t \pi(y_t | T, Q, I, C, M_i, \mathbf{y}_{<t})$ 
24: else
25:   return Final diagnosis  $\prod_t \pi(y_t | T, Q, I, C, \mathbf{y}_{<t})$ 
26: end if
```

---

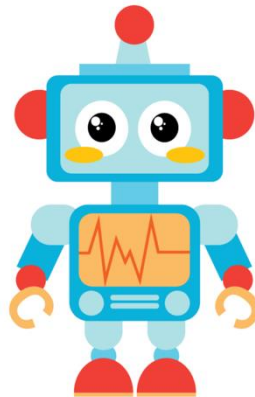
## Part 5 – Agentic System

- Tool using (functional call)
- **Enhanced interpretability and user trust**
  - Perform a **specific task** with **specific output constraints** (but with input perturbations).
  - Better information retrieval with reduced hallucinations.
- **Access to external tools, databases, and APIs**
  - **Knowledge Acquisition** ← External knowledge
- **Automation and efficiency**
  - Complete complex tasks without excessive prompting.
- **Interaction enhancement**
  - **Interact with external models**, e.g., MRI model, CT scan model, etc.



## Part 5 – Agentic System

- Tool using (functional call)
- **Task planning**
  - Which step should I solve?
- **Tool selection**
  - Should I use a tool(s)?
  - Which tool(s) should I use to solve this step problem?
- **Tool calling**
  - Calling and performing the specific action that the tool has defined.
- **Response generation**
  - Generating a response based on the tool's output description.





## Part 5 – Agentic System

### ■ Tools

Tool Name	Parameter	Parameter Description
<code>generate_options</code>	<code>answers</code>	The most likely answers based on the patient's specific condition and needs.
<code>perform_comparison</code>	<code>comparisons</code>	A structured comparison of all options, detailing their relevance to the patient's case.
<code>enable_search</code>	<code>search</code>	Answer 'yes' or 'no' to indicate search necessity.
<code>relevance_analysis</code>	<code>analysis</code>	A comprehensive analysis detailing the relevance of each document to the patient's presentation, highlighting key matches, inconsistencies, and important findings.
<code>locate_evidence</code>	<code>evidence</code>	Relevant evidence applicable to the patient's presentation, with article IDs in <code>&lt;quote&gt;&lt;/quote&gt;</code> tags.

## Part 5 – Agentic System

### ■ Database

Corpus	Number of Docs	Number of Snippets	Average Length
PubMed Abstracts	23,897,881	23,897,881	290.01
Wikipedia	6,458,670	29,642,311	166.47
Clinical Trials	156,887	4,177,121	268.33
PubMed Central Articles	123,194	8,155,929	202.46
Textbooks	8,226	2,224,013	207.95
Clinical Cases	1,479	17,821	215.61

## ■ Database

Journal Title	Article Count	Journal Title	Article Count
BMJ Open	37,488	JAMA Ophthalmol	434
Proc Natl Acad Sci U S A	16,619	Lancet HIV	387
JAMA Netw Open	10,824	BMJ Health Care Inform	366
Nature	8,148	JAMA Surg	287
Cell	4,811	BMJ Neurol Open	282
Science	4,660	JAMA Dermatol	279
BMJ	3,636	Lancet Psychiatry	270
BMJ Glob Health	3,460	Lancet Public Health	264
N Engl J Med	2,159	BMJ Support Palliat Care	262
BMJ Open Qual	1,569	BMJ Nutr Prev Health	254
JAMA	1,552	Lancet Respir Med	252
BMJ Open Diabetes Res Care	1,434	JAMA Cardiol	239
Lancet	1,344	Lancet Diabetes Endocrinol	225
Neurology	1,216	Lancet Microbe	167
BMJ Open Sport Exerc Med	1,201	BMJ Ment Health	167
Lancet Reg Health West Pac	1,196	JAMA Otolaryngol Head Neck Surg	164
BMJ Case Rep	1,190	Lancet Planet Health	162
BMJ Paediatr Open	1,145	Lancet Haematol	157
Lancet Reg Health Eur	1,077	BMJ Med	154
BMJ Open Respir Res	1,031	Lancet Child Adolesc Health	154
Lancet Reg Health Am	901	BMJ Evid Based Med	136
Ann Intern Med	881	Lancet Digit Health	124
Lancet Glob Health	805	BMJ Surg Interv Health Technol	120
JAMA Intern Med	797	Lancet Gastroenterol Hepatol	117
Lancet Infect Dis	676	BMJ Oncol	114
BMJ Open Ophthalmol	656	Lancet Healthy Longev	102
JAMA Neurol	639	BMJ Sex Reprod Health	100
JAMA Health Forum	628	BMJ Mil Health	64
BMJ Open Gastroenterol	625	Lancet Rheumatol	61
Lancet Oncol	613	BMJ Open Sci	49
BMJ Qual Saf	601	BMJ Innov	46
JAMA Psychiatry	597	BMJ Simul Technol Enhanc Learn	42
JAMA Pediatr	569	JAMA Facial Plast Surg	39
BMJ Qual Improv Rep	547	Ann Intern Med Clin Cases	6
JAMA Oncol	490	BMJ Outcomes	1
Lancet Reg Health Southeast Asia	464	BMJ Clin Evid	1
BMJ Public Health	453		
Lancet Neurol	444		
		<b>Total Number of Articles</b>	<b>123,194</b>

## Part 5 – Agentic System

- Benchmarks (medical analysis and diagnosis)

Benchmark	Number of Testing Cases	Number of Choices
USMLE Step 1 (40)	94	9
USMLE Step 2 (40)	109	6
USMLE Step 3 (40)	122	6
MedQA (41)	1,273	4
MedExpQA (6)	125	5

**Getting  
more  
challenging**

## Part 5 – Agentic System

### ■ Benchmarks (multiple choice question-answering)

Model	USMLE Step 1	USMLE Step 2	USMLE Step 3	MedQA	MedExpQA
GPT-4	<u>80.67</u>	<u>81.67</u>	<b>89.78</b>	<b>78.87</b>	N/A
ChatGPT	51.26	60.83	58.39	50.82	N/A
BioMistral (7B)	34.04	37.61	37.70	41.01	37.60
OpenBioLLM (8B)	47.87	44.04	50.00	47.84	43.20
UltraMedical (8B)	42.55	27.52	34.43	38.49	35.20
OpenBioLLM (70B)	69.15	70.64	68.85	69.13	<u>71.20</u>
UltraMedical (70B)	70.21	55.05	56.56	52.32	50.40
PodGPT (70B)	73.40	72.48	74.59	65.04	63.20
<b>Ours</b>	<b>82.98</b>	<b>86.24</b>	<u>88.52</u>	<u>73.29</u>	<b>78.40</b>

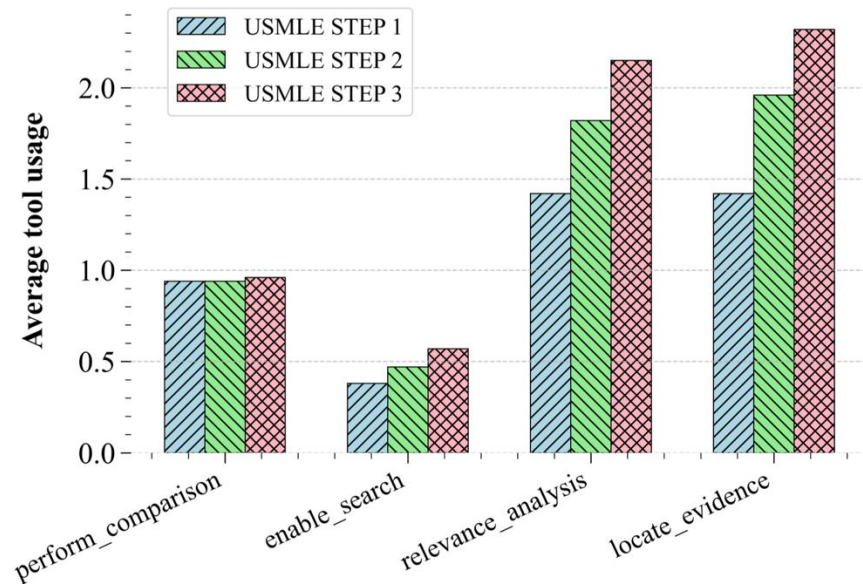
## Part 5 – Agentic System

### ■ Benchmarks (open-ended question-answering)

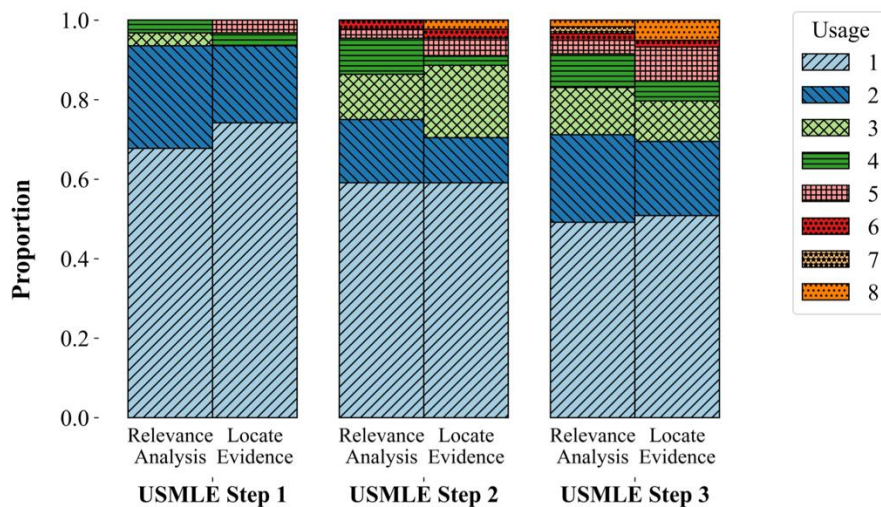
Benchmark	Model	BioMistral (7B)	OpenBioLLM (8B)	UltraMedical (8B)	OpenBioLLM (70B)	UltraMedical (70B)	PodGPT (70B)	Ours
USMLE Step 1	SFR	0.79 $\pm$ 0.09	0.70 $\pm$ 0.12	0.81 $\pm$ 0.13	0.85 $\pm$ 0.10	0.82 $\pm$ 0.11	<u>0.86</u> $\pm$ 0.11	<b>0.87</b> $\pm$ 0.09
	GTE	0.48 $\pm$ 0.17	0.38 $\pm$ 0.17	0.57 $\pm$ 0.21	0.60 $\pm$ 0.23	0.63 $\pm$ 0.23	<u>0.66</u> $\pm$ 0.24	<b>0.66</b> $\pm$ 0.22
	BERTScore	0.58 $\pm$ 0.12	0.51 $\pm$ 0.13	0.61 $\pm$ 0.16	0.66 $\pm$ 0.17	0.64 $\pm$ 0.17	<u>0.68</u> $\pm$ 0.20	<b>0.68</b> $\pm$ 0.17
USMLE Step 2	SFR	0.76 $\pm$ 0.11	0.71 $\pm$ 0.10	0.80 $\pm$ 0.11	0.82 $\pm$ 0.09	0.80 $\pm$ 0.10	<u>0.85</u> $\pm$ 0.10	<b>0.85</b> $\pm$ 0.09
	GTE	0.45 $\pm$ 0.19	0.38 $\pm$ 0.15	0.52 $\pm$ 0.19	0.54 $\pm$ 0.19	0.59 $\pm$ 0.22	<u>0.62</u> $\pm$ 0.21	<b>0.62</b> $\pm$ 0.22
	BERTScore	0.58 $\pm$ 0.11	0.56 $\pm$ 0.11	0.61 $\pm$ 0.13	0.64 $\pm$ 0.13	0.63 $\pm$ 0.14	<u>0.66</u> $\pm$ 0.15	<b>0.67</b> $\pm$ 0.15
USMLE Step 3	SFR	0.74 $\pm$ 0.10	0.70 $\pm$ 0.10	0.79 $\pm$ 0.12	<u>0.85</u> $\pm$ 0.11	0.80 $\pm$ 0.11	<u>0.84</u> $\pm$ 0.11	<b>0.86</b> $\pm$ 0.09
	GTE	0.41 $\pm$ 0.18	0.38 $\pm$ 0.14	0.53 $\pm$ 0.22	<u>0.63</u> $\pm$ 0.26	0.60 $\pm$ 0.23	0.63 $\pm$ 0.24	<b>0.65</b> $\pm$ 0.22
	BERTScore	0.57 $\pm$ 0.11	0.52 $\pm$ 0.14	0.60 $\pm$ 0.17	<b>0.71</b> $\pm$ 0.19	0.62 $\pm$ 0.15	0.67 $\pm$ 0.18	<u>0.70</u> $\pm$ 0.17
MedQA	SFR	0.76 $\pm$ 0.10	0.71 $\pm$ 0.12	0.80 $\pm$ 0.12	<b>0.86</b> $\pm$ 0.11	0.80 $\pm$ 0.11	0.84 $\pm$ 0.11	<u>0.85</u> $\pm$ 0.10
	GTE	0.43 $\pm$ 0.18	0.40 $\pm$ 0.17	0.53 $\pm$ 0.22	<b>0.63</b> $\pm$ 0.26	0.58 $\pm$ 0.23	0.60 $\pm$ 0.23	<u>0.61</u> $\pm$ 0.23
	BERTScore	0.56 $\pm$ 0.12	0.52 $\pm$ 0.15	0.60 $\pm$ 0.16	<b>0.70</b> $\pm$ 0.19	0.61 $\pm$ 0.16	0.65 $\pm$ 0.18	<u>0.67</u> $\pm$ 0.18
MedExpQA	SFR	0.76 $\pm$ 0.10	0.71 $\pm$ 0.11	0.78 $\pm$ 0.13	0.81 $\pm$ 0.11	0.77 $\pm$ 0.13	<u>0.83</u> $\pm$ 0.11	<b>0.84</b> $\pm$ 0.10
	GTE	0.47 $\pm$ 0.18	0.40 $\pm$ 0.18	0.52 $\pm$ 0.22	0.54 $\pm$ 0.24	0.55 $\pm$ 0.22	<b>0.61</b> $\pm$ 0.23	<u>0.60</u> $\pm$ 0.22
	BERTScore	0.58 $\pm$ 0.11	0.53 $\pm$ 0.12	0.58 $\pm$ 0.15	0.62 $\pm$ 0.17	0.60 $\pm$ 0.14	<u>0.65</u> $\pm$ 0.17	<b>0.65</b> $\pm$ 0.16

## Part 5 – Agentic System

### ■ Tool usage



(a)



(b)

## Part 5 – Agentic System

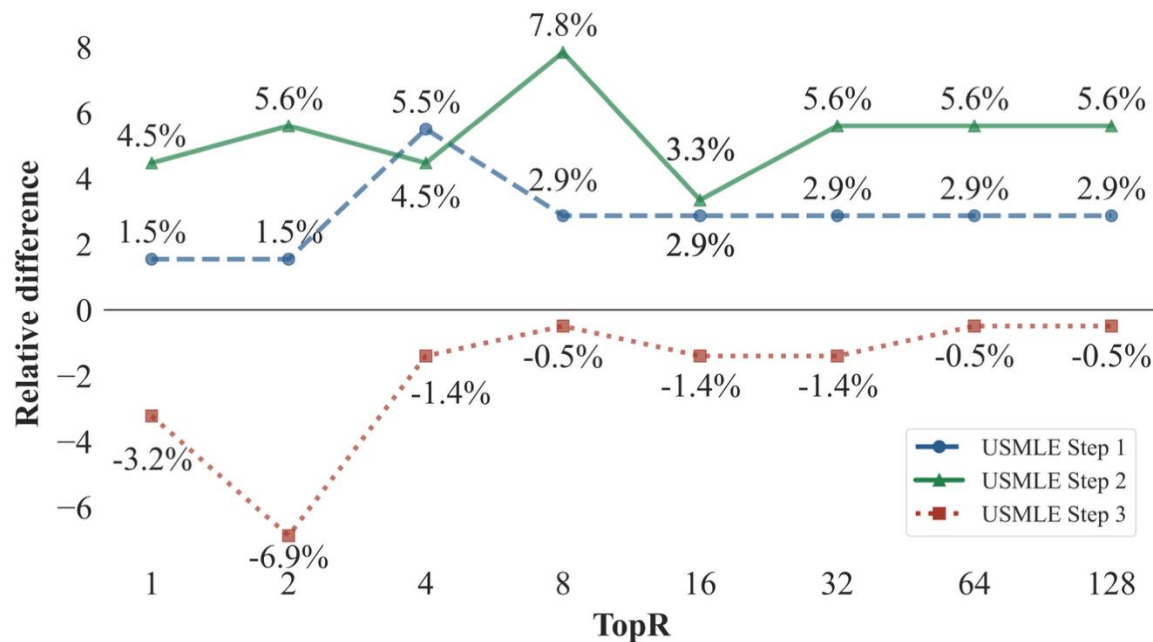
- Ablation study

Benchmark	USMLE Step 1	USMLE Step 2	USMLE Step 3	Average
Ours	<b>82.98</b>	<b>86.24</b>	<b>88.52</b>	<b>85.91</b>
w/o Tools	-1.07	-3.67	-4.91	-3.22
w/o Cache & Prune	-1.07	-2.75	-3.27	-2.36
w/o Evidence Search	-2.13	-3.67	-6.55	-4.12



## Part 5 – Agentic System

### ■ Ablation study



## Part 5 – Agentic System

### ■ Ablation study

Table 1: Performance evaluation of different LLMs. We report different LLMs' performance (accuracy) across the USMLE Steps 1 to 3.

Model	USMLE Step 1	USMLE Step 2	USMLE Step 3
GPT-4	80.67	81.67	<b>89.78</b>
ChatGPT	51.26	60.83	58.39
Ours (Qwen2.5 72B)	<u>82.98</u>	<u>86.24</u>	<u>88.52</u>
Ours (Qwen3 32B)	<b>91.49</b>	<b>87.16</b>	86.07

**Reasoning-based model**

## Part 5 – Agentic System

- Runtime performance metric

Table 1: Latency (in seconds), per-GPU seconds per query (in seconds), and GPU hours per query (in hours) of the agentic system deployed using 4 NVIDIA L40S GPUs.

<b>Qwen2.5 72B</b>	<b>USMLE Step 1</b>	<b>USMLE Step 2</b>	<b>USMLE Step 3</b>
Latency	1.12	1.48	1.80
per-GPU seconds per query	39.15	51.64	83.30
GPU hours per query	0.043	0.057	0.093

## Part 5 – Agentic System

### ■ Take away

- **Agentic framework unifies** document retrieval, evidence grounding, as well as AI generation, with an **automatic and dynamic process**
- **Tool-augmented LLM-based agent** enables **dynamic multistep tool use**, eliminating the need for manually engineered prompts or multi-stage pipelines
- **The cache-and-prune memory bank mechanism** efficiently extends the retention of relevant documents for **evidence grounding**, enhancing diagnostic accuracy and computational efficiency

# **Thank you very much for your attention!**