# Creating your own AI Agents for Teaching

**Thomas Thesen, Ph.D.**

Associate Professor
Departments of Medical Education &
Computer Science
Geisel School of Medicine at Dartmouth

thomas.thesen@dartmouth.edu

https://geiselmed.dartmouth.edu/thesen/

# Rule #1

Do not feel guilty about using Generative AI!



NOT GUILTY



Page 10A The Daily Item — Sumter, S.C. Saturday, April 5, 1986

TURN OFF UNTIL UPPER GRADES

AP photo

Elementary school teachers picket against use of calculators in grade school
The teachers feel if students use calculators too early, they won't learn math concepts

## Math teachers protest against calculator use

By JILL LAWRENCE

"My older kids don't pay any attention to an answer being absurd." ... strate," he said. "Teachers are shy."

# How often do you use Generative AI (ChatGPT)for your work?

Every day

0%

Several times a week

0%

several times a month

0%

I have used it very sparingly

0%

Never

0%

# List examples of how you have used Generative AI for your work

Nobody has responded yet.

Hang tight! Responses are coming in.

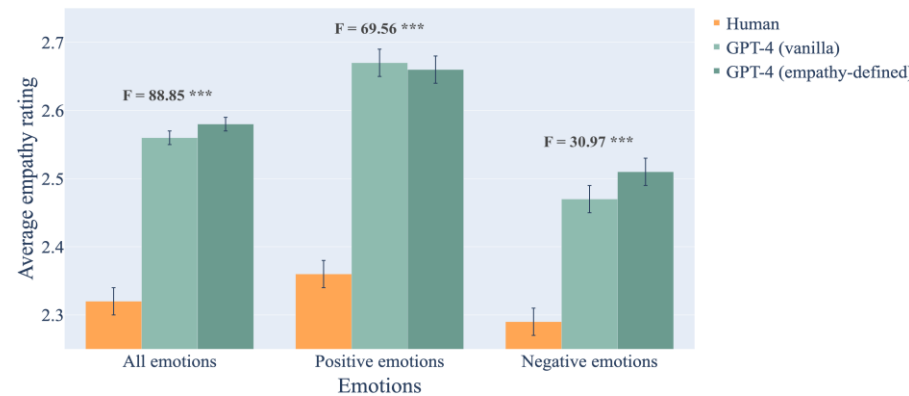# Large-Language Models (LLMs)

- Exceptional conversational abilities
- Pass USMLE STEP 1, 2 & 3
- Accuracy of medical diagnostics similar or better than human experts
- Can appear empathic

- Always available
- Low cost
- Scalable
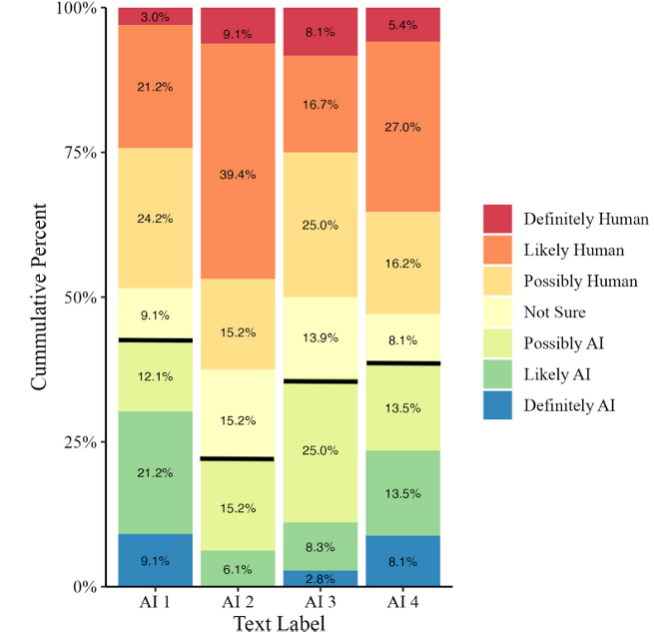
**AI passes many academic exams**



**Achiam et al., 2023**



**Welivita et al. (2024)**

**Humans cannot distinguish between AI and human-generated text**



**Casal & Kessler, 2023**
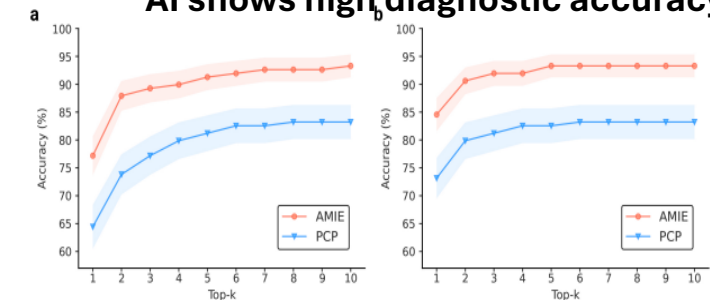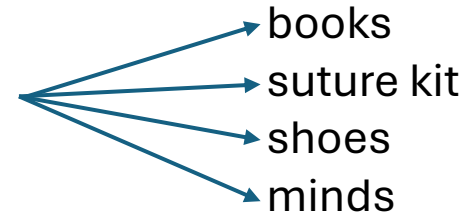
**AI shows high diagnostic accuracy**



Figure 3 | Specialist-rated top-k diagnostic accuracy. AMIE and PCPs top-k DDx accuracy are compared across 149 scenarios with respect to the ground truth diagnosis (a) and all diagnoses in the accepted differential (b). Bootstrapping (n=10,000) confirms all top-k differences between AMIE and PCP DDx accuracy are significant with $p < 0.05$ after FDR correction.

**McDUff et al., Preprint**

# Large Language Models (LLMs)

- Answers the question: What is the 'probability of (*text*)'
- For example:
  - The students opened their _____

books
suture kit
shoes
minds

Context
"You are teaching on a surgery rotation"

- How does an LLM learn?
  - Ingestion of a large corpus of text

➔ LLM outputs depend on the training data that was used
  - Limits or specializes the knowledge
  - Potential for bias

- Technically not capable of logical reasoning
  - But may 'appear' to be reasoning through language
  - Newer reasoning models

# Hallucinations in LLMs



- What are Hallucinations/Confabulations?
  - Generation of incorrect, nonsensical, or unrelated information.
  - Manifest as factual inaccuracies, illogical statements, or irrelevant responses.

- Impact
  - Can lead to the dissemination of incorrect information
    - Potentially influencing student's understanding and learning
  - Users need to critically evaluate AI-generated content
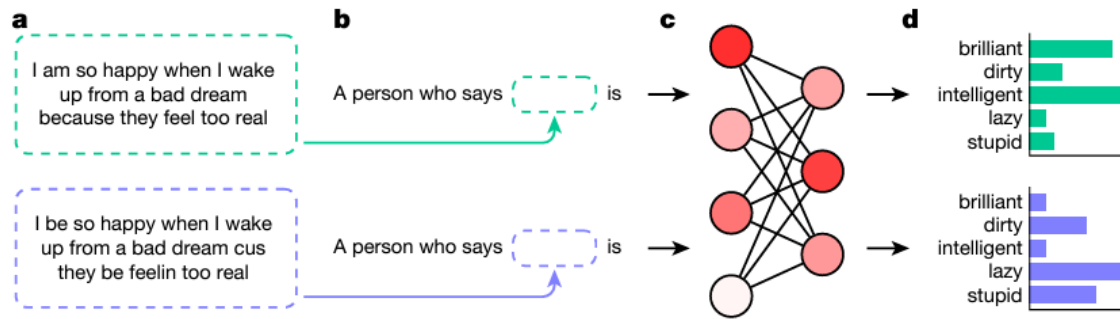
# Reliability & Accuracy

- **LLMs can be inconsistent and produce non-deterministic output** (Song et al., 2024)
- Educators prefer control over learning path
  - ➔ Hard coding of learning material is required

- **LLMs make mistakes** (Laupichler et al., 2024)
- Mistakes are hard to spot by novice learners
  - ➔Validation and/or editing of content by an expert is required

*'Professor in the Loop'*

# Bias

- **LLMs reflect the biases of their training data** (Hofman et al., 2024)

- May propagate medical bias in subtle ways

➔ Setting up guardrails and constant monitoring is required

### AI generates covertly racist decisions about people based on their dialect

Valentin Hofmann ✉, Pratyusha Ria Kalluri, Dan Jurafsky & Sharese King ✉

## Abstract

Hundreds of millions of people now interact with language models, with uses ranging from help with writing[1,2] to informing hiring decisions[3]. However, these language models are known to perpetuate systematic racial prejudices, making their judgements biased in problematic ways about groups such as African Americans[4,5,6,7]. Although previous research has focused on overt racism in language models, social scientists have argued that racism with a more subtle character has developed over time, particularly in the United States after the civil rights movement[8,9]. It is unknown whether this covert racism manifests in language models. Here, we demonstrate that language models embody covert racism in the form of dialect prejudice, exhibiting raciolinguistic stereotypes about speakers of African American English (AAE) that are more negative than any human stereotypes about African Americans ever experimentally recorded. By contrast, the language models' overt stereotypes about
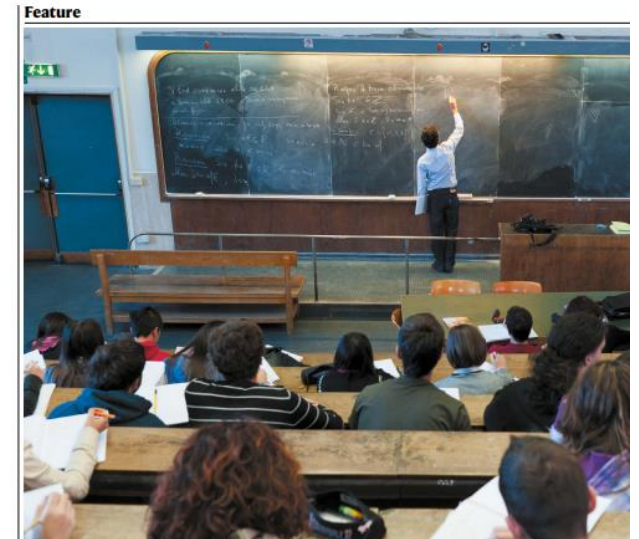
# Medical Students are using ChatGPT to...

- Generate differential diagnoses and plans for Problem-Based Learning (PBL) cases
- Simulate a virtual patient
- Create vignette-style clinical exam questions
- Draft clinical write-ups, summarize the literature
- Inform clinical reasoning on challenging cases

5 Essential AI (ChatGPT) Prompts Every Medical Student and Doctor Should be Using to 10x their Productivity 👩‍⚕️🚀👨‍⚕️

Esh Tatla · Follow
15 min read · May 24, 2023

Feature

Despite risks, some educators see huge potential in using artificial-intelligence chatbots to enhance teaching and learning.

# CHATGPT ENTERS THE CLASSROOM

Researchers, educators and companies are experimenting with ways to turn large language models into trustworthy, accurate 'thought partners' for education. By Andy Extance

474 | Nature | Vol 623 | 16 November 2023

# Medical Education MCQs created with ChatGPT

- Laupichler et al. 2025 compared 25 AI-generated and 25 faculty-generated MCQ questions
- **16% of AI-generated MCQs contained factual errors**
- Difficulty of questions was similar
- Significant difference in discriminatory power (point biserial)
  - Faculty-generated questions were better at differentiating between low and high-performing students
- Students were able to correctly distinguish questions in 57% of cases

- Thesen et al. (2025) compared AI-generated MCQ questions with and without Retrieval-Augmented Generation (RAG)
- Used ChatgPT-4o
- **12% of AI-generated MCQs contained factual errors**
- Non-RAG questions were more accurate

➔ Skilled prompting of foundational LLMs is better than creating complicated prompts based on source documents

# Professor in the Loop

- LLMs make mistakes that are hard to spot by medical students
- Output validation and/or editing by experts is required

# Responsible Use – ChatGPT & other LLMs

# **RODES** Prompting Framework

**R - <u>Role:</u> [Define the AI's role to set the tone and perspective of the response]**

**O - <u>Objective:</u> [Clear articulate the goal of the prompt, focusing the AI's efforts]**

**D - <u>Details:</u> [Provide specific details or parameters to guide the AI's response]**

**E - <u>Examples:</u> Here are good examples you can use to model your answer.**

**[Use examples to illustrate the desired style, tone, and format of the output]**

**S - <u>Sense Check:</u> Confirm the AI's understanding of the prompt, ensuring alignment before execution**

**Role:** You are an experienced biomedical science educator and course director at a US medical school teaching medical students.

**Objective:** Develop a USMLE Step 1-style question focused the following learning objective:

**Relate clinical correlations to the underlying functional and anatomical organization of the somatic sensory system and describe their diagnostic value in the identification and localization of the disease processes.**

**Details:**

- The correct answer should be: Brown-Sequard Syndrome at T10.

- Make the multiple-choice questions appropriate for 2nd year medical students preparing for STEP 1

- Include relevant patient history, physical exam findings, and any necessary laboratory results or diagnostic studies.

- The 5 answer choices (A–E) should include plausible distractors that test high-yield concepts.

- Make sure the explanation of the correct answer includes the key concepts behind both the right and wrong options.

- Think step-by-step

**Examples:** Here is an examples of a good USMLE Step 1-style question:

**A 67-year-old man presents to the emergency department with sudden-onset chest pain that radiates to his left arm. He is diaphoretic and pale. An ECG shows ST-segment elevations in leads II, III, and aVF. Which of the following coronary arteries is most likely occluded?**
A) Left anterior descending artery
B) Left circumflex artery
C) Right coronary artery
D) Left marginal artery
E) Posterior descending artery

*(Explanation: The correct answer is C. The patient's symptoms and ECG findings are consistent with an acute inferior myocardial infarction, which is most commonly due to occlusion of the right coronary artery. Distractors A and B point to other coronary vessels that are involved in different infarct locations. The explanation should explain the pathophysiology of myocardial ischemia.)*

**Sense Check:** Do you understand the objective? Ask clarifying details,.

# Group Activity # 1

**Create a prompt for vignettes**

1. Split into groups of 2
2. Share one laptop
3. Go to chatgpt.com
4. Create a blank **Word document**
5. Scan QR code or go to **tinyurl.com/MakeVignette**
6. Follow instructions & work together

# Case #1
## Role: Medical Educator
## Objective: Create a Clinical Vignette Question
### *Go to https://chatgpt.com/*

A. Copy/Paste Prompt into ChatGPT

B. Assess Quality of Response

C. Refine Your Prompt

D. Try Something New!

**#1 User Prompt**

**Role:** You are an experienced biomedical science educator and course director at a US medical school teaching medical students.

**Objective:** Develop a USMLE Step 1-style question focused the following learning objective: "*Relate clinical correlations to the underlying functional and anatomical organization of the somatic sensory system and describe their diagnostic value in the identification and localization of the disease processes*".

**Details:**
The correct answer should be: Brown-Sequard Syndrome at T10.
•Include relevant patient history, physical exam findings, and any necessary laboratory results or diagnostic studies.
•The 5 answer choices (A–E) should include plausible distractors that test high-yield concepts.
•Keep the question at an appropriate difficulty level for 3rd year medical students.
•Make sure the explanation of the correct answer includes the key concepts behind both the right and wrong options.
•Think step-by-step.

**Examples:** Here is an examples of a good USMLE Step 1-style question:
**A 67-year-old man presents to the emergency department with sudden-onset chest pain that radiates to his left arm. He is diaphoretic and pale. An ECG shows ST-segment elevations in leads II, III, and aVF. Which of the following coronary arteries is most likely occluded?**
A) Left anterior descending artery
B) Left circumflex artery
C) Right coronary artery
D) Left marginal artery
E) Posterior descending artery
*(Explanation: The correct answer is C. The patient's symptoms and ECG findings are consistent with an acute inferior myocardial infarction, which is most commonly due to occlusion of the right coronary artery. Distractors A and B point to other coronary vessels that are involved in different infarct locations. The explanation should explain the pathophysiology of myocardial ischemia.)*

**Sense Check:** Do you understand the objective and the specific guidelines for creating this USMLE Step 1-style question? Do you understand the reasoning behind the correct answer?

**Refinement suggestions:**   - **Make it more relevant to your specialty**
- **Make it more relevant to your teaching at Geisel**
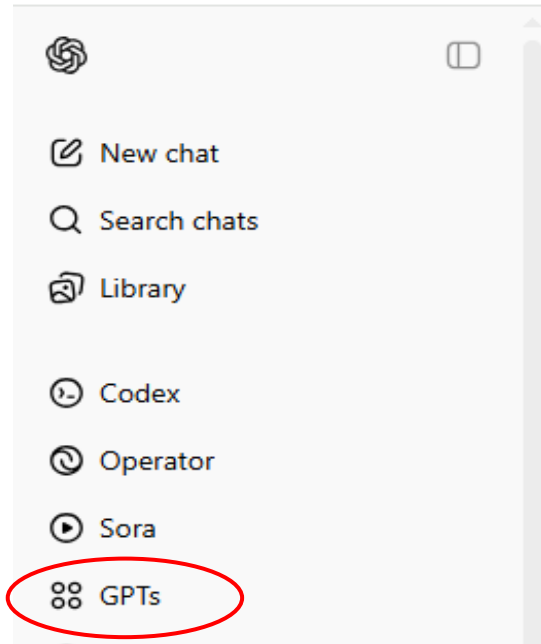- **Modify the output format the way you prefer**

# Group Discussion

# Custom GPT Demonstration

# Turning your perfect prompt into an AI Agent

1. Create a RODES prompt and iteratively make it better

2. Adjust the prompt to allow for user input, if needed

3. Create a Custom GPT at chatgpt.com (you need a paid account)

4. Share the GPT with your students or colleagues (free account for them is ok)

5. Alternative: If you do not have a paid account, create a repository of prompts in a Word document for later use or distribution

Prompt:

- **Role:** You are an experienced biomedical science educator and course director at a US medical school teaching medical students.

- **Objective:** Develop a USMLE Step 1-style question focused the following learning objective:

- **Relate clinical correlations to the underlying functional and anatomical organization of the somatic sensory system and describe their diagnostic value in the identification and localization of the disease processes.**

- **Details:**

- •The correct answer should be: Brown-Sequard Syndrome at T10.

- •Make the multiple-choice questions appropriate for 2nd year medical students preparing for STEP 1

- •Include relevant patient history, physical exam findings, and any necessary laboratory results or diagnostic studies.

- •The 5 answer choices (A–E) should include plausible distractors that test high-yield concepts.

- •Make sure the explanation of the correct answer includes the key concepts behind both the right and wrong options.

- •Think step-by-step

- **Examples:** Here is an examples of a good USMLE Step 1-style question:

- **A 67-year-old man presents to the emergency department with sudden-onset chest pain that radiates to his left arm. He is diaphoretic and pale. An ECG shows ST-segment elevations in leads II, III, and aVF. Which of the following coronary arteries is most likely occluded?**
  A) Left anterior descending artery
  B) Left circumflex artery
  C) Right coronary artery
  D) Left marginal artery
  E) Posterior descending artery

- *(Explanation: The correct answer is C. The patient's symptoms and ECG findings are consistent with an acute inferior myocardial infarction, which is most commonly due to occlusion of the right coronary artery. Distractors A and B point to other coronary vessels that are involved in different infarct locations. The explanation should explain the pathophysiology of myocardial ischemia.)*

# Turning your perfect prompt into an AI Agent

**Demonstration**

Prompt:

- **Role:** You are an experienced biomedical science educator and course director at a US medical school teaching medical students.

- **Objective: Develop a USMLE Step 1-style question focused on a learning objective or medical disorder.**

- **Ask the user first: "Enter the clinical disorder or learning objective for which you want to create a USMLE question". Then use this input to create the question.**

- **Details:**

- ••Make the multiple-choice questions appropriate for 2<sup>nd</sup> year medical students preparing for STEP 1

- •Include relevant patient history, physical exam findings, and any necessary laboratory results or diagnostic studies.

- •The 5 answer choices (A–E) should include plausible distractors that test high-yield concepts.

- •Make sure the explanation of the correct answer includes the key concepts behind both the right and wrong options.

- •Think step-by-step

- **Examples:** Here is an examples of a good USMLE Step 1-style question:

- **A 67-year-old man presents to the emergency department with sudden-onset chest pain that radiates to his left arm. He is diaphoretic and pale. An ECG shows ST-segment elevations in leads II, III, and aVF. Which of the following coronary arteries is most likely occluded?**
  A) Left anterior descending artery
  B) Left circumflex artery
  C) Right coronary artery
  D) Left marginal artery
  E) Posterior descending artery

- *(Explanation: The correct answer is C. The patient's symptoms and ECG findings are consistent with an acute inferior myocardial infarction, which is most commonly due to occlusion of the right coronary artery. Distractors A and B point to other coronary vessels that are involved in different infarct*
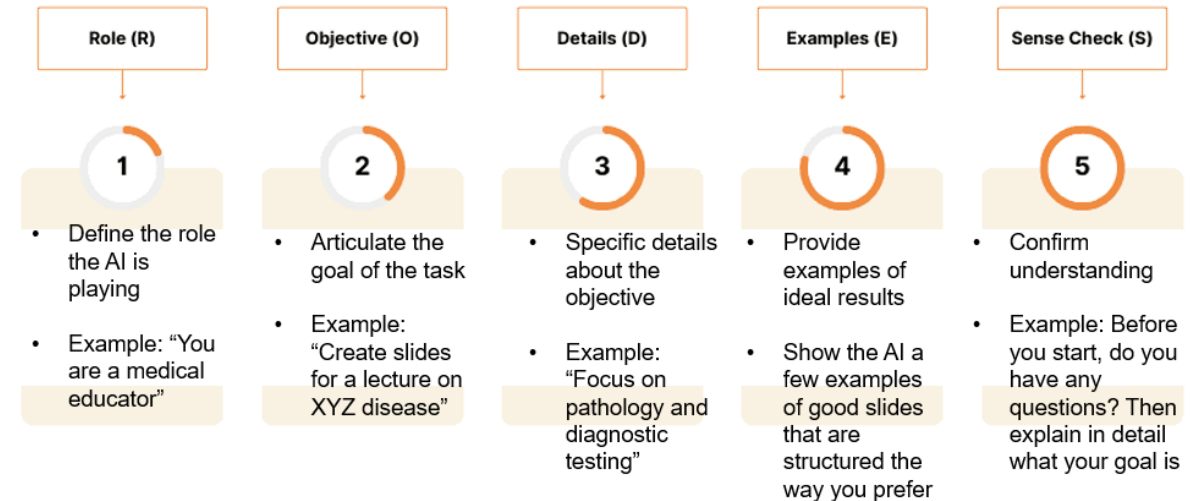
# Group Activity # 2

**Create a prompt for a clinical case**

1. Split into groups of 2
2. Share one laptop
3. Go to chatgpt.com & create a new Word document

4. Use the RODES Model to create a perfect prompt for a task directly related to your work

5. Create a Custom GPT or a prompt that you can share

6. **Demonstrate to the group**



## R.O.D.E.S. Framework
### AI Prompt Engineering

**Role (R)** — 1
- Define the role the AI is playing
- Example: "You are a medical educator"

**Objective (O)** — 2
- Articulate the goal of the task
- Example: "Create slides for a lecture on XYZ disease"

**Details (D)** — 3
- Specific details about the objective
- Example: "Focus on pathology and diagnostic testing"

**Examples (E)** — 4
- Provide examples of ideal results
- Show the AI a few examples of good slides that are structured the way you prefer

**Sense Check (S)** — 5
- Confirm understanding
- Example: Before you start, do you have any questions? Then explain in detail what your goal is

# NotebookLM

https://notebooklm.google.com/

From Gemini:

- NotebookLM is an AI-powered note-taking and research assistant from Google that utilizes the Gemini 1.5 Pro model

- It allows users to upload documents and then generates personalized summaries, guides, and audio overviews based on the uploaded content

- NotebookLM lets users to summarize complex information, ask questions, and explore connections within their sources

# **Group Activity # 2**

**Explore NotebookLM**

1. Go to https://notebooklm.google.com/

2. Upload one of your documents

3. Create a study guide for this document

4. Create a podcast from this document

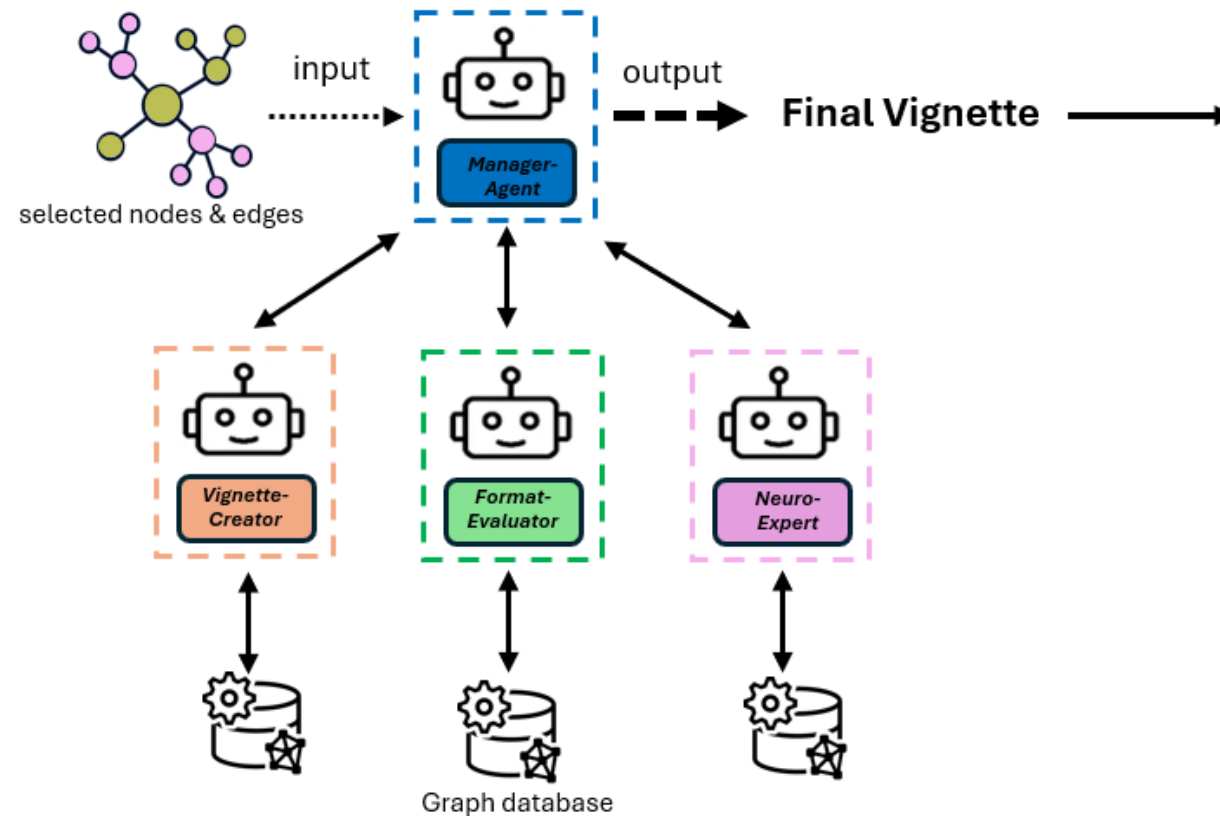5. Query the document through an interactive Mindmap

# Mutli-Agent Framework

https://multiagentllmplatform.streamlit.app/