# Latterview: A Multi-Agent Framework for Systematic Review Automation using Large Language Models

Pouria Rouzrokh, MD, MPH, MHPE, Yale New Haven Hospital, Yale University

## Introduction/Background

Systematic literature reviews are vital for evidence-based medicine but are time-consuming and labor-intensive. We introduce and evaluate an open-source, multi-agent framework utilizing large language models (LLMs) to automate and accelerate key steps in the systematic medical literature review process, aiming to enhance efficiency while maintaining rigor.

## Methods/Intervention

We developed LatteReview, a Python framework employing configurable LLM agents (e.g., for title/abstract screening, relevance scoring, data abstraction) orchestrated within customizable multi-round workflows. The framework supports various LLM providers (cloud/local), Retrieval-Augmented Generation, and structured outputs. We evaluated a representative workflow using two 'junior' LLM agents (Gemini-1.5-flash, GPT-4O-mini) for initial title/abstract screening based on inclusion/exclusion criteria, with a 'senior' LLM agent (GPT-4O) resolving disagreements. Evaluation used six diverse public systematic review datasets (SYNERGY collection) and a custom dataset derived from prior cardiothoracic imaging AI reviews, featuring progressively complex criteria. Performance was measured by Area Under the Curve (AUC) and accuracy, evaluated across various decision thresholds.

## Results/Outcome

Across the six SYNERGY datasets, the automated workflow demonstrated substantial discriminative capability with AUCs ranging from 0.71 to 0.95. Performance variability was observed, linked to heterogeneity in criteria clarity and dataset inclusion rates (0.9% to 12.4%). On the custom dataset with well-defined criteria, AUCs ranged from 0.79 (most complex criteria) to 0.94 (simplest criteria). Furthermore, by refining agent prompts and calibrating the decision threshold using a 150-record validation set per task, performance was significantly enhanced, achieving accuracy of 99% on subsequent unseen data. Efficiency testing showed processing of 1000 abstracts took approximately 1 minute costing ~$1.20 using the selected models.

## Conclusion

LLM-driven multi-agent systems, as implemented in our framework, can effectively automate significant portions of the systematic review process, particularly title and abstract screening. Performance is promising and improves markedly with well-defined review criteria and thoughtful workflow design.

## Statement of Impact

Automating systematic reviews using AI tools like LatteReview can drastically accelerate the synthesis of medical research findings. This enables faster development and updating of evidence-based clinical

guidelines, quicker assessment of new diagnostic or therapeutic advancements, and helps clinicians and researchers stay current with the rapidly expanding body of literature.
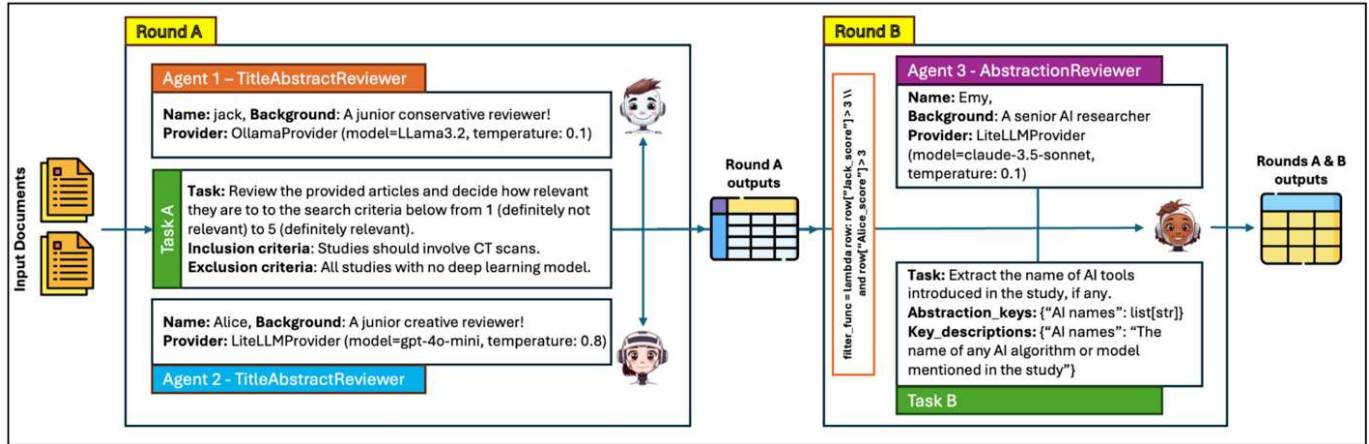


Figure 1. An example multiagent review workflow for title and abstract screening with two junior reviewer agents (round A) followed by concept extraction by a senior and more powerful reviewer agent (round B).



Figure 2. An example workflow for systematic review of some articles using the TitleAbstractReviewer agents.
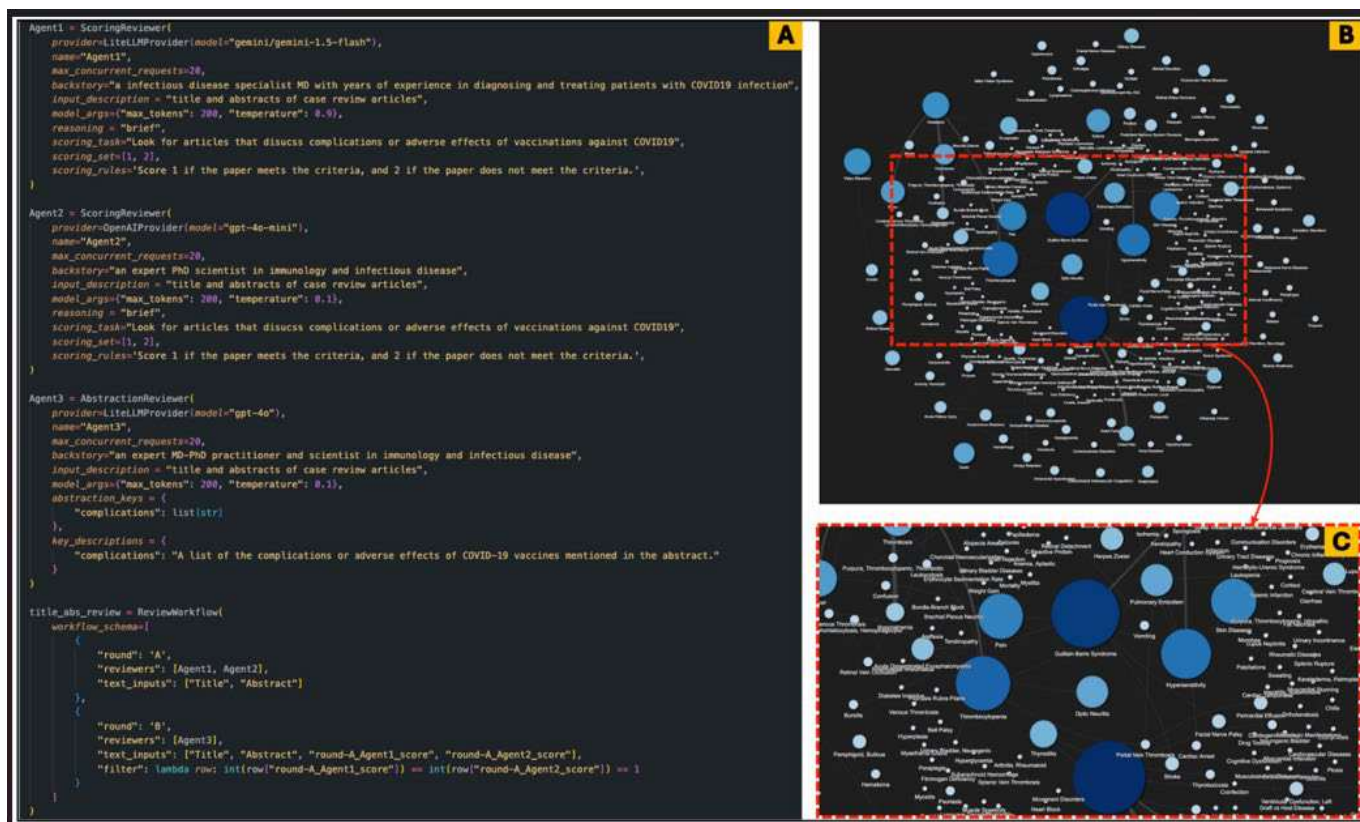
Figure 3. (A) A review workflow for screening titles and abstracts from a dataset created by querying PubMed for case report articles on complications of COVID-19 vaccines. The workflow consists of two agents that detect relevant articles and a third agent that abstracts the complications mentioned in those articles into a list. (B and C) A knowledge graph of complications extracted by the reviewers, where each node represents a complication. The size of the nodes correlates with the number of articles mentioning that complication, with larger and darker nodes indicating higher mention counts.