Can Agents Judge Systematic Reviews Like Humans? Evaluating SLRs with LLM-based Multi-Agent System

Anonymous ACL submission

Abstract

Systematic Literature Reviews (SLRs) are foundational to evidence-based research but remain labor-intensive and prone to inconsis-004 005 tency across disciplines. We present an LLMbased SLR evaluation copilot built on a Multi-007 Agent System (MAS) architecture to assist researchers in assessing the overall quality of the systematic literature reviews. The system automates protocol validation, methodological assessment, and topic relevance checks using a scholarly database. Unlike conventional singleagent methods, our design integrates a specialized agentic approach aligned with PRISMA 015 guidelines to support more structured and interpretable evaluations. We conducted an initial study on five published SLRs from diverse 017 domains, comparing system outputs to expertannotated PRISMA scores, and observed 84% agreement. While early results are promising, this work represents a first step toward scalable and accurate NLP-driven systems for interdisciplinary workflows and reveals their capacity for rigorous, domain-agnostic knowledge aggregation to streamline the review process.

1 Introduction

041

Systematic Literature Reviews are foundational to evidence-based research, offering a structured, protocol-driven approach for identifying, analyzing, and synthesizing prior work. In contrast to narrative reviews, SLRs follow a predefined methodology to promote transparency, reproducibility, and reduced bias (Grant and Booth, 2009; Booth et al., 2021). They are widely employed to surface research trends, identify knowledge gaps, and establish a grounded basis for future inquiry across disciplines.

However, the exponential growth of scholarly publications, varying in quality and relevance, has made it increasingly challenging to maintain rigor and comprehensiveness in the SLR process. This



Figure 1: An overview of the proposed multi-agent system LLM-based SLR evaluation framework.

overload slows down the review pipeline and introduces risks of redundancy and overlooks literature.

043

044

046

047

050

051

052

054

060

061

062

063

064

065

To address these challenges, we propose an AIaugmented SLR system that leverages a multiagent LLM architecture. Inspired by humancentered co-pilot design principles (Sellen and Horvitz, 2024), our system supports the SLR workflow: from protocol critique and methodological assessment to relevance checking, duplication detection, and collaborative drafting. Figure 1 provides a high-level overview of the proposed system.

The tool is designed with an interdisciplinary lens, recognizing the cross-domain nature of SLRs, from health sciences to software engineering, while leveraging LLM-based agentic architectures to enhance quality and efficiency using the NLP architecture proposed by (Vaswani et al., 2017). By combining automation with human oversight, our approach takes an initial step toward improving accessibility and robustness in evidence synthesis. To guide our study, we focus on the practical capabilities and evaluation of the proposed system within the SLR workflow.

Our research questions (RQs) are as follows:

- 067 068

- 073
- 074
- 077

094

100

102

104

105

107

108

109

110

111

112

113

114

115

069

tency during SLR evaluation?

Background and Related Work 2

of the SLR process?

RQ1: How can a multi-agent LLM system support

RQ2: How well does the system's output align with

the protocol validation and compliance steps

PRISMA standards, based on initial expert

evaluations, and does the system offer mea-

surable improvements in efficiency or consis-

SLRs are essential for synthesizing research findings, identifying gaps, and guiding future studies. The PRISMA guidelines (Moher et al., 2009; Page et al., 2021) offer a structured approach to ensure transparency and reproducibility. Tools like Covidence (Covidence Systematic Review Software, 2024) and Rayyan (Ouzzani et al., 2016) assist in screening and data extraction but rely heavily on manual effort, making them prone to fatigue, selection bias, and inconsistency.

Recent work explores leveraging LLMs for automating SLR stages such as study identification, summarization, and quality assessment (Susnjak, 2023; Ge and Others, 2024; Smith and Others, 2023; Jones and Others, 2022). State-of-the-art LLMs like GPT-4, Claude 3.7, Llama, and Gemini 2.5 (OpenAI, 2025b; DeepMind, 2025; Anthropic-AI, 2025; Meta-AI, 2025) demonstrate impressive few-shot learning capabilities (Brown et al., 2020), making them suitable for structured tasks with minimal supervision.

To enhance performance on complex, multi-step tasks, MAS have emerged as powerful frameworks capable of decomposing problems, enabling cooperative reasoning, and outperforming single-agent models on structured benchmarks (Park et al., 2025; Zhang et al., 2025; Wang et al., 2024). Recent MAS-driven tools illustrate this shift: Google's AI Co-Scientist (Gottweis et al., 2025) leverages Gemini-based agents to generate hypotheses and propose experiments; OpenAI's Deep Research (OpenAI, 2025a) conducts autonomous literature synthesis via web search; and SciSpace (SciSpace, 2025) offers interactive document parsing and drafting.

While existing tools support general scientific exploration and offer basic workflow automation, they often lack the methodological rigor required for systematic reviews. Similarly, recent studies underscore the limitations of core monolithic LLMs in structured tasks: Lieberum et al. reviewed 37 GPT-

based SR prototypes and found them largely unval-116 idated; Penzo et al. demonstrated that LLaMA2-117 13B exhibits prompt and structure sensitivity; and 118 Wei et al.; Wang et al. reported only modest 119 gains, ranging from 3.9% to 17.9%, from tech-120 niques like chain-of-thought and self-consistency 121 across various benchmarks. These findings sug-122 gest inherent performance ceilings in end-to-end 123 LLM pipelines for SRs. To address the limitations 124 posed by monolithic LLMs, we introduce a mod-125 ular, multi-agent LLM system explicitly aligned 126 with PRISMA guidelines, where each checklist 127 item is handled by a specialized agent under expert 128 oversight. Early results indicate improved consis-129 tency and reliability. Our open-source implemen-130 tation¹ promotes transparency, supports scalable, 131 high-quality reviews, and provides a foundation for 132 robust, end-to-end SLR support. 133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

3 Methodology & System Design

3.1 Architecture

Our MAS-LLM SLR Evaluation Framework consists of 27 specialized agents organized into six PRISMA-aligned societies (see Table 1) along with two utility agents (PDF Parsing and Follow-up Conversation). This structure (number of agents and division in societies) mirrors the PRISMA checklist which clearly elicits the checklist items for each part of the systematic reviews (Page et al., 2021; Susnjak, 2023), with one agent per checklist item. Sections like Methods have more agents (11) due to their detailed protocol requirements (greater number of checklist items), while sections like Discussion require only one agent. Early experiments with multi-item agents resulted in unstable behavior-overloaded agents, degraded performance, and unpredictable agent spawning, so we adopted a one-agent-per-item design with oneshot detailed prompting along with examples for robustness and clarity. All agents use GPT-4.1 (OpenAI, 2025b), a state-of-the-art LLM built for agentic workflows with a 1M token context window. Upon SLR PDF upload, an OCR-enabled Vision-Language Model (Unstructured Technologies, Inc., 2025) converts it into structured text.

A Coordinator Agent and Task Agent decompose the PRISMA checklist into modular evaluation tasks, dispatching them via few-shot prompts to specialized agents. Each agent uses the arXiv

¹GitHub link will be added here upon publication



Figure 2: Architecture of the MAS-LLM SLR evaluation framework. The user-uploaded SLR document is processed by multiple agentic societies designed to evaluate it according to PRISMA guidelines. The results are displayed on the web user interface, which also supports follow-up questions and interactive conversations.

Toolkit to retrieve relevant research as needed, assigns a 0–5 score, and provides qualitative feedback. If outputs fall below thresholds, the Coordinator reallocates tasks or spawns new agents. As shown in Figure 2, agent outputs are synthesized into a unified format, accessible via web-interface and provided to the Follow-up Conversation Agent.

Society	Function	Agents	
Abstract &	Evaluate title clarity and	2	
Title	abstract completeness	2	
Introduction	Assess rationale, scope, and	2	
	objectives		
Methods	Check eligibility criteria, search	11	
	strategy, and bias assessment		
	Verify result reporting,		
Results	visualizations, and statistical	7	
	summaries		
Discussion	Examine interpretation,	1	
	limitations, and implications	1	
Other Infor- mation	Review registration, funding,		
	conflicts of interest, and data	4	
	policies		
Standalone	PDF parsing and follow-up	2	
	dialogue	2	

Table 1: Agent societies, their responsibilities, and the number of agents in each society.

The follow-up conversation agent named *SLR-GPT Agent* provides co-pilot style research support through professional interactions. Using the same arXiv Toolkit available to all agents, it sug-

gests new papers, verifies citations, cross-checks literature results, and recommends editorial refinements to maximize PRISMA compliance. By combining structured evaluation outputs from agents from societies with in-context retrieval for both the original manuscript and PRISMA protocol, it transforms assessments into actionable guidance for manuscript improvement. Open source implementation ² is also shared for reproducibility.

175

176

177

178

179

180

181

183

184

185

186

187

188

190

191

192

194

195

196

197

199

200

3.2 Evaluation Methodology

We assess our framework on five published SLRs from diverse fields (Medical, E-commerce, AI, Metaverse, IoT), comparing agent outputs against ratings by three expert SLR reviewers. Both agents and human experts score each PRISMA item on a 0-5 scale (0 = Not Addressed, ..., 5 = ThoroughlyAddressed), enabling a standardized, ordinal evaluation across sections. Agent prompts incorporate PRISMA guidelines and one-shot exemplars to standardize evaluation; human reviewers assess the original manuscripts along with PRISMA guidelines. We quantify alignment using Mean Absolute Error (MAE), agreement level using MAE and additional statistical analysis to pinpoint areas where multi-agent collaboration aligns best with human experts. This benchmark demonstrates the tech-

164

165

166

167

168

169

170

172 173

²GitHub link will be added here upon publication

237 238

239 240

241

242

243

244

245

246

247

248

249

250

252

253

254

255

256

257

258

260

261

262

263

264

265

266

267

268

269

270

271

272

273

274

275

276

277

278

279

282

4.2 Paper-wise & Inter-expert Analysis

minutes based on length and complexity, offering

early-stage insights that can guide and accelerate

subsequent human reviews, effectively reducing

Per-paper analysis (Appendix Fig.4) shows agent scores consistently track human evaluations across all five SLRs. Expert comparison (Appendix Fig.5) indicates the highest inter-rater variability. High inter-expert agreement—Intraclass Correlation Coefficient = 0.924, Krippendorff's $\alpha = 0.889$, Pearson $\rho = 0.898$ (Appendix Table 2), validates our human benchmark and suggests remaining divergences highlight targets for system refinement.

5 Future Directions

overall turnaround time.

To ensure practical utility, we plan to deploy an interactive, browser-based interface that allows users to ask questions, revise summaries, and re-score sections, enabling both quantitative and qualitative evaluation of UI design and agent responsiveness. We will test the system with real-world users, systematic reviewers, and authors to assess its collaborative effectiveness. Structured feedback will be collected using Likert scales to evaluate clarity, usefulness, and trust, following best practices in human–AI interaction (Zhao et al., 2024; Rong et al., 2022). This feedback will be integrated into a preference-tuning loop to better align agent behavior with user preferences (Shao et al., 2025).

6 Conclusion

We introduced a multi-agent system for evaluating Systematic Literature Reviews aligned with the PRISMA protocol. By leveraging specialized agents for protocol validation, topic relevance, structural assessment, and ArXiv integration, the system aims to reduce the manual burden of interdisciplinary SLR evaluation. The built-in SLR-GPT Assistant supports citation checks and editorial feedback. An initial small-scale empirical evaluation on five SLRs from diverse domains showed 84% agreement with expert SLR judgments. While promising, these results are preliminary. This work offers a first step toward scalable, accurate MAS for SLRs, with future efforts focused on broader evaluations and system refinement with a focus on human-AI collaboration.

nical viability and interdisciplinary applicability of agentic LLMs in streamlining SLR workflows across diverse scientific domains.

4 Results

201

202

206

207

208

210

211

212

213

214

215

216

217

219

222

223

227

236

4.1 MAE & Agreement Level

We evaluated our system on five published SLRs from diverse domains (Medical, AI, Ecommerce, IoT, Metaverse), comparing its PRISMA-based scores against expert human reviewers (Section 3.2). Agreement percentages, computed as $100\% - (MAE/5 \times 100)$, are shown in Figure 3.

Mean Absolute Error (MAE): Agents vs. Human Experts



Figure 3: Agreement between agents and humans. Avg. agreement is 84%, with the highest alignment in Introduction (97%) and lower in Other Information (81%).

According to our results, the overall agreement is 84%, with strongest alignment in Introduction (97%), Discussion (94%), and Methods (93%). Alignment dips only slightly to Results (84%) and Other Information (81%). In absolute terms, the highest agreement exceeds the overall agreement by 13 points, while the lowest is just 3 points below, indicating that all section-level agreements fall within a narrow window around the mean (close alignment overall). These agreement levels demonstrate that our system faithfully reproduces expert judgments on core review components by aligning itself with the PRISMA protocol while highlighting opportunities for future improvement. This alignment with experts also shows that our proposed MAS supports the protocol validation and compliance mainly through its system design and architectural choices mentioned in Section 3.

To address the significant delays in traditional peer review, averaging 15 weeks for the first round, with 10 weeks in medical sciences, 14 in natural sciences, and 17 in social sciences (SciRev, 2014), and up to 25 weeks in fields like Economics (Huisman and Smits, 2017), we present a automated MAS review system. It analyzes papers in just 15–20

Limitations

Our MAS-LLM framework shows early promise, but several limitations should be acknowledged. The current evaluation spans five SLRs from distinct domains. In subsequent studies, we are looking towards increasing the number of papers and hence the credibility of our results. Agent perfor-289 mance is bounded by the capabilities of current LLMs, which may overlook fine-grained domain knowledge in technical contexts. The system's integration with arXiv enhances open-access coverage but excludes other key databases like PubMed or Scopus, creating potential gaps. Moreover, the system currently supports only evaluation, not realtime drafting or collaboration. Nonetheless, this 297 study demonstrates the feasibility of structured, agentic LLM support for SLRs and lays the groundwork for more scalable and interactive systems in future work. 301

References

305

310

311

312

313

314

315

316

317

318 319

320

321

323

324

325

326

327

330

- Anthropic-AI. 2025. Claude 3.7 Sonnet and Claude Code. Technical report, Anthropic AI. Retrieved from https://www.anthropic.com/news/ claude-3-7-sonnet.
- Andrew Booth, Martyn-St James, Mark Clowes, Anthea Sutton, and 1 others. 2021. Systematic approaches to a successful literature review.
 - Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, and 1 others. 2020. Language models are few-shot learners. Advances in neural information processing systems, 33:1877–1901.
 - Covidence Systematic Review Software. 2024. Covidence: Streamlining systematic review workflow. Online resource. Available at https://www. covidence.org/.
 - DeepMind. 2025. Gemini: Revolutionizing AI with multimodal capabilities. Technical report, Deep-Mind. Retrieved from https://deepmind.google/ technologies/gemini/.
- X Ge and Others. 2024. AI-driven systematic reviews in health research. *Medical Informatics Journal*.
- Juraj Gottweis, Wei-Hung Weng, Alexander Daryin, Tao Tu, Anil Palepu, Petar Sirkovic, Artiom Myaskovsky, Felix Weissenberger, Keran Rong, Ryutaro Tanno, and 1 others. 2025. Towards an AI coscientist. *arXiv preprint arXiv:2502.18864*.
- Maria J Grant and Andrew Booth. 2009. A typology of 331 reviews: an analysis of 14 review types and associ-332 ated methodologies. Health information & libraries journal, 26(2):91-108. 334 Jeroen Huisman and Jeroen Smits. 2017. Duration and 335 quality of the peer review process: the author's per-336 spective. Scientometrics, 113(1):633-650. 337 L Jones and Others. 2022. Automating sentiment analysis in systematic reviews. AI Review Journal. Judith-Lisa Lieberum, Markus Töws, Maria-Inti Metzendorf, Felix Heilmeyer, Waldemar Siemens, Chris-341 tian Haverkamp, Daniel Böhringer, Joerg J. Meerpohl, and Angelika Eisele-Metzger. 2025. Large lan-343 guage models for conducting systematic reviews: on 345 the rise, but not yet ready for use—a scoping review. Journal of Clinical Epidemiology, 181:111746. 346 Meta-AI. 2025. The llama 4 herd: The begin-347 ning of a new era of natively multimodal AI 348 innovation. Technical report, Meta AI. Re-349 trieved from https://ai.meta.com/blog/ 350 llama-4-multimodal-intelligence/. 351 David Moher, Alessandro Liberati, Jennifer Tetzlaff, 352 and Douglas G Altman. 2009. Preferred reporting 353 items for systematic reviews and meta-analyses: The 354 PRISMA statement. BMJ, 339:b2535. 355 Introducing OpenAI. 2025a 356 deep research. https://openai.com/index/ 357 introducing-deep-research/. Accessed: 358 2025-05-06. 359 OpenAI. 2025b. Introducing GPT-4.1 in the API. 360 Mourad Ouzzani, Hossam Hammady, Zbys Fedorow-361 icz, and Ahmed Elmagarmid. 2016. Rayyan-a web and mobile app for systematic reviews. Systematic 363 Reviews, 5:210. 364 Matthew J Page, Joanne E McKenzie, Patrick M 365 Bossuyt, Isabelle Boutron, Tammy C Hoffmann, Cyn-366 thia D Mulrow, Larissa Shamseer, Jennifer M Tet-367 zlaff, Elie A Akl, Sue E Brennan, and 1 others. 2021. The PRISMA 2020 statement: an updated guideline 369 for reporting systematic reviews. BMJ, 372:n71. 370 Sooyoung Park, Franziska Roesner, Tadayoshi Kohno, 371 and Daniel Weld. 2025. Why do multi-agent LLM 372 systems fail? arXiv preprint arXiv:2503.13657. 373 Nicolò Penzo, Maryam Sajedinia, Bruno Lepri, Sara 374 Tonelli, and Marco Guerini. 2024. Do LLMs suffer 375 from multi-party hangover? a diagnostic approach to 376 addressee recognition and response selection in con-377 versations. In Proceedings of the 2024 Conference on 378 Empirical Methods in Natural Language Processing, 379 pages 11210–11233. Association for Computational 380 Linguistics. 381 Yao Rong, Tobias Leemann, Thai-Trang Nguyen, Lisa 382 Fiedler, Peizhu Qian, Vaibhav Unhelkar, Tina Seidel, 383

- 390 394 396 397 400 401 402 403 404 405 406 407 408 409 410 411 412 413 414 415 416 417 418 419 420 421 422 423 494 425 426 427 428 429 430

- 431
- 432 433

434 435

436 437

Gjergji Kasneci, and Enkelejda Kasneci. 2022. Towards human-centered explainable ai: A survey of user studies for model explanations. arXiv preprint arXiv:2210.11584.

- SciRev. 2014. Average duration first review round 15 weeks. Accessed: 2025-05-18.
- SciSpace. 2025. Scispace: AI chat for scientific pdfs. https://typeset.io/. Accessed: 2025-05-06.
- Abigail Sellen and Eric Horvitz. 2024. The rise of the AI co-pilot: Lessons for design from aviation and beyond. Communications of the ACM, 67(7):18-23.
- Zekai Shao, Yi Shan, Yixuan He, Yuxuan Yao, Junhong Wang, Xiaolong Zhang, Yu Zhang, and Siming Chen. 2025. Do language model agents align with humans in rating visualizations? an empirical study. arXiv preprint arXiv:2505.06702.
- J Smith and Others. 2023. Sentiment analysis in systematic literature reviews. Computational Linguistics.
- Teo Susnjak. 2023. PRISMA-DFLLM: An extension of PRISMA for systematic literature reviews using domain-specific finetuned large language models. arXiv preprint arXiv:2306.14905.
- Unstructured Technologies, Inc. 2025. Unstructured your GenAI has a data problem. Accessed: February 27, 2025.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. Advances in neural information processing systems, 30.
- Xuezhi Wang, Yixin Zhou, Dale Schuurmans, and 1 others. 2023. Self-consistency improves chain of thought reasoning in language models. In Proceedings of the 2023 EMNLP, pages 4116-4128. Association for Computational Linguistics.
- Zihan Wang, Xiang Ren, and 1 others. 2024. Llm multiagent systems: Challenges and open problems. arXiv preprint arXiv:2402.03578.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Harm Le, and 1 others. 2022. Chain of thought prompting elicits reasoning in large language models. In Proceedings of the 2022 EMNLP, pages 2480-2493. Association for Computational Linguistics.
- Ke Zhang, Yali Zhang, Jinlong Li, and 1 others. 2025. A comprehensive survey on multi-agent cooperative decision-making. Information Fusion. ArXiv:2503.13415.
- Lingjun Zhao, Nguyen X. Khanh, and Hal Daumé III. 2024. Successfully guiding humans with imperfect instructions by highlighting potential errors and suggesting corrections. In Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, pages 719-736.

Appendix Α

Papers (SLR)-Wise Analysis A.1

Figure 4 presents a paper-wise comparison of av-440 erage PRISMA scores assigned by our MAS-LLM 441 system versus human experts across five SLRs. The 442 mean absolute error (MAE) by paper ranges from 443 0.05 (Paper 3) to 0.44 (Paper 2), demonstrating con-444 sistently strong alignment and identifying specific 445 instances for targeted model refinement. 446

438

439

447

448

449

450

451

452

453

Paper-wise Analysis: Agents vs. Human Experts Average Scores Per Paper





Figure 4: Paper-wise comparison of average scores (Agents vs. Humans).

A.2 Human Experts' Scores

Figure 5 illustrates inter-expert variability across PRISMA checklist societies. Methods sections exhibit the highest reviewer disagreementreflecting the inherent complexity of methodological assessments-whereas Title & Abstract sections achieve the greatest consensus.

Human Experts Comparison



Figure 5: Variation in human expert scores by society.

A.3 Inter-Expert Reliability Analysis

454

455

456

457

458

459

460

461

Table 2 summarizes inter-expert reliability metrics, including Intraclass Correlation Coefficient (ICC), Krippendorff's Alpha, and average Pearson ρ , all exceeding 0.88. These high values confirm the robustness of our expert benchmark and substantiate the validity of comparing agent outputs against human ratings.

Metric	Value	Interpretation
Intraclass Correlation Coefficient (ICC)	0.924	Excellent reliabil- ity
Krippendorff's Alpha	0.889	Strong agreement
Avg. Inter-Human Pearson ρ	0.898	Strong correlation

Table 2: Inter-expert agreement metrics.