



Information Synthesis Workflows with Large Language Models: A Review of Current Practices, Opportunities, and Challenges in Literature Review Tasks

Bryce Huffman (presenter),
Dr. Caitlin Grady & Dr. Zoe Szajnfarter

AI4SE & SE4AI Workshop 2025
September 17, 2025

Engineering Management and Systems Engineering
The George Washington University



Agenda



Introduction and motivations

Methods

Results and key findings

Wrap-up

Advancements in AI Large Language Models have led to rapid increases in commercial, purpose-built AI tools for unstructured text workflows

① GENERAL MODELS

- Public availability catalyzed in late-2022
- Ongoing advancements



② OPEN QUESTIONS

- How is their potential harnessed?
- Where could they be used?
- What are they good for?
- What are they not good for?

③ PROPOSAL WRITING?

- Common across industries
- Core challenge is unstructured text processing and synthesis
- Direct similarities to SE artifact generation and other related tasks

Exploding market of commercial LLM tools claiming to support proposal writing workflows

AutogenAI

AutoRFP.ai

govdash.com

GOVEAGLE

GovSignals

grant assistant
by FREEMILL

Grantable

grant boost

inventive

loopio

RS AI

responsive

sweetspot

/ IDENTIFIED SET

Agenda

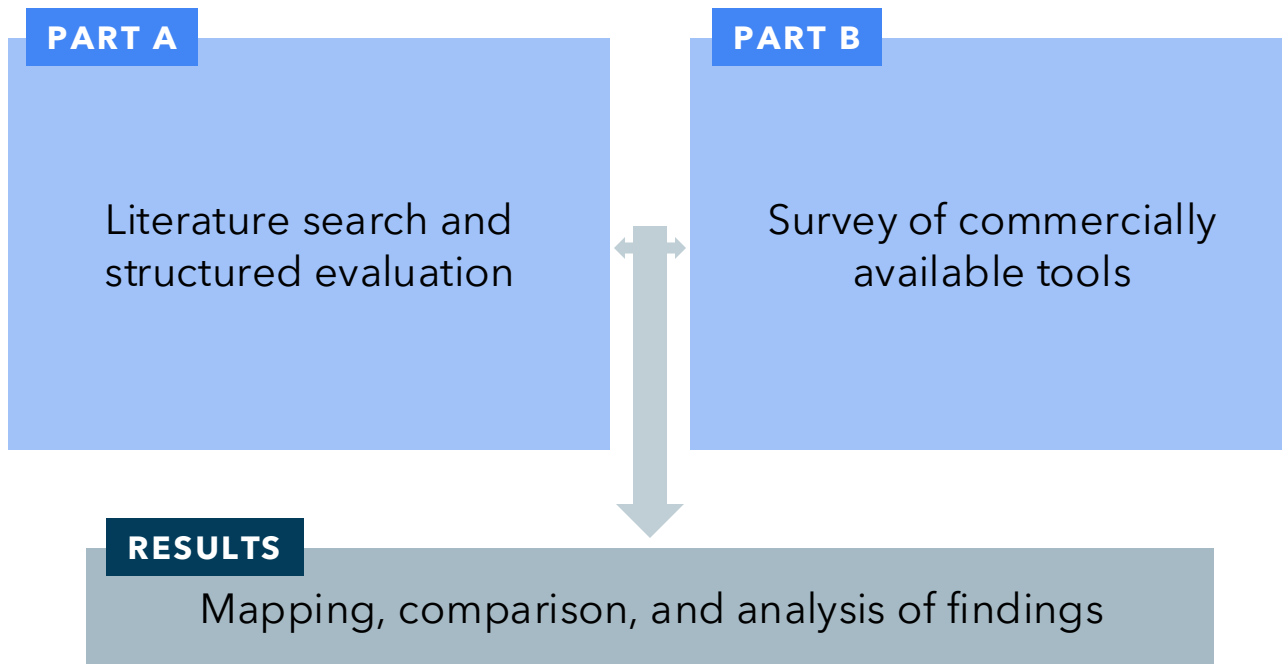
Introduction and motivations

Methods

Results and key findings

Future research

This study employed a 2-part research design



Agenda

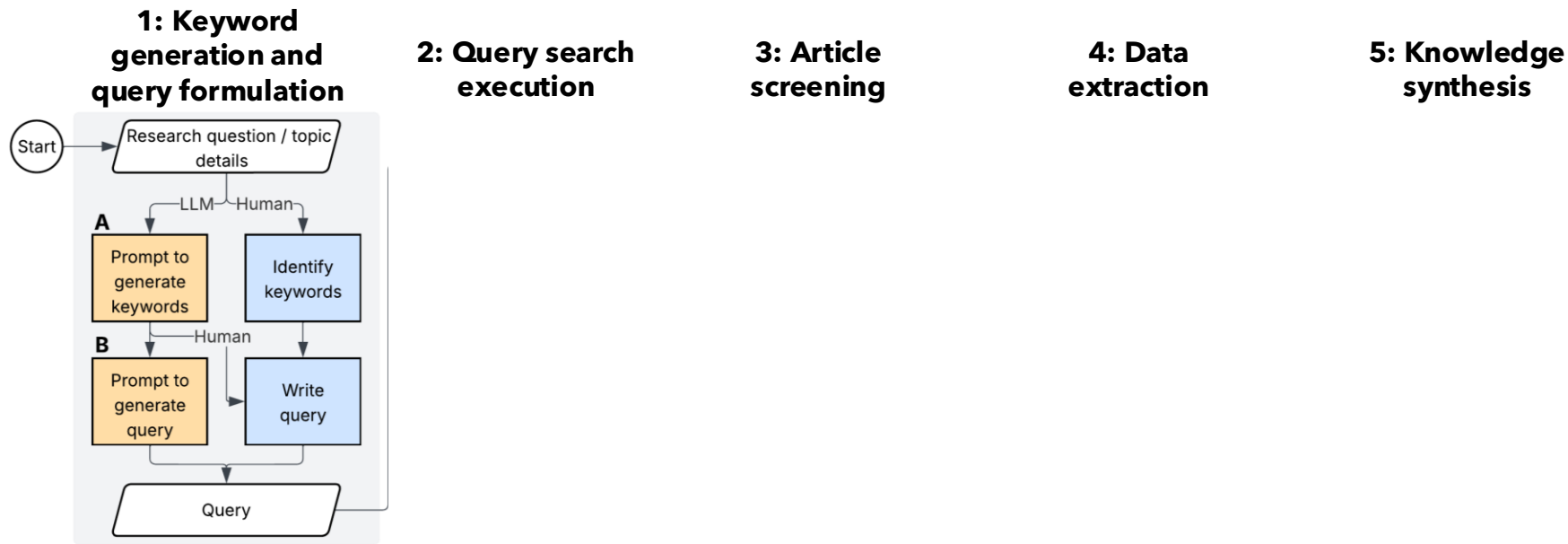
Introduction and motivations

Methods

Results and key findings

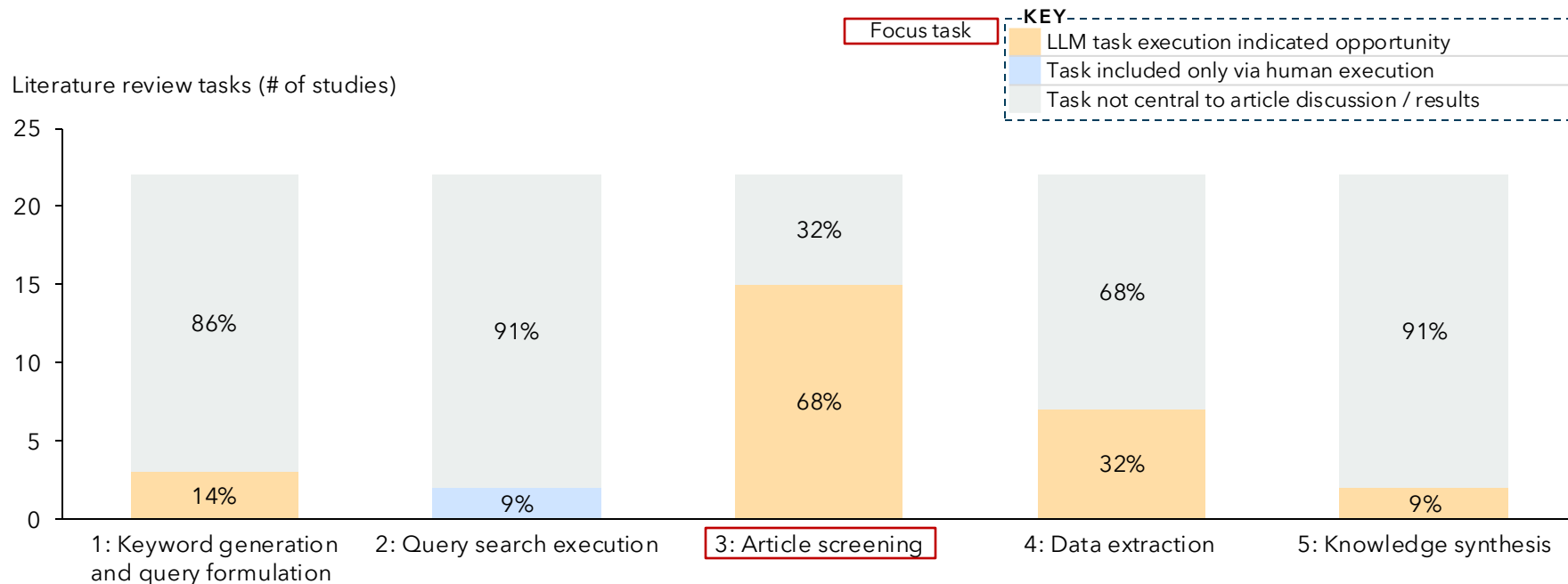
Future research

Five core tasks emerge in the LLM-supported literature review workflow for unstructured text synthesis, based on a review of 22 studies



Across these five tasks, LLMs gain the most traction in the literature within article screening, where studies find potential

Opportunity identified in ~70% of studies for article screening; only ~9% for synthesis

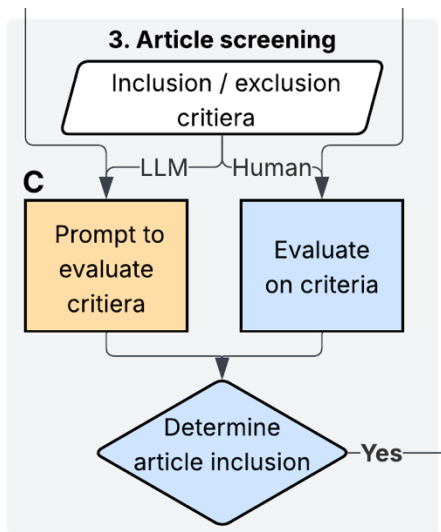


As literature review workflows are analogous to those in proposal writing, we can ground claims about purpose-built tools in evidence



Mapping purpose-built tool features to core workflow tasks reveals the breadth of stated and implied functionality

Article screening is the primary focus within literature

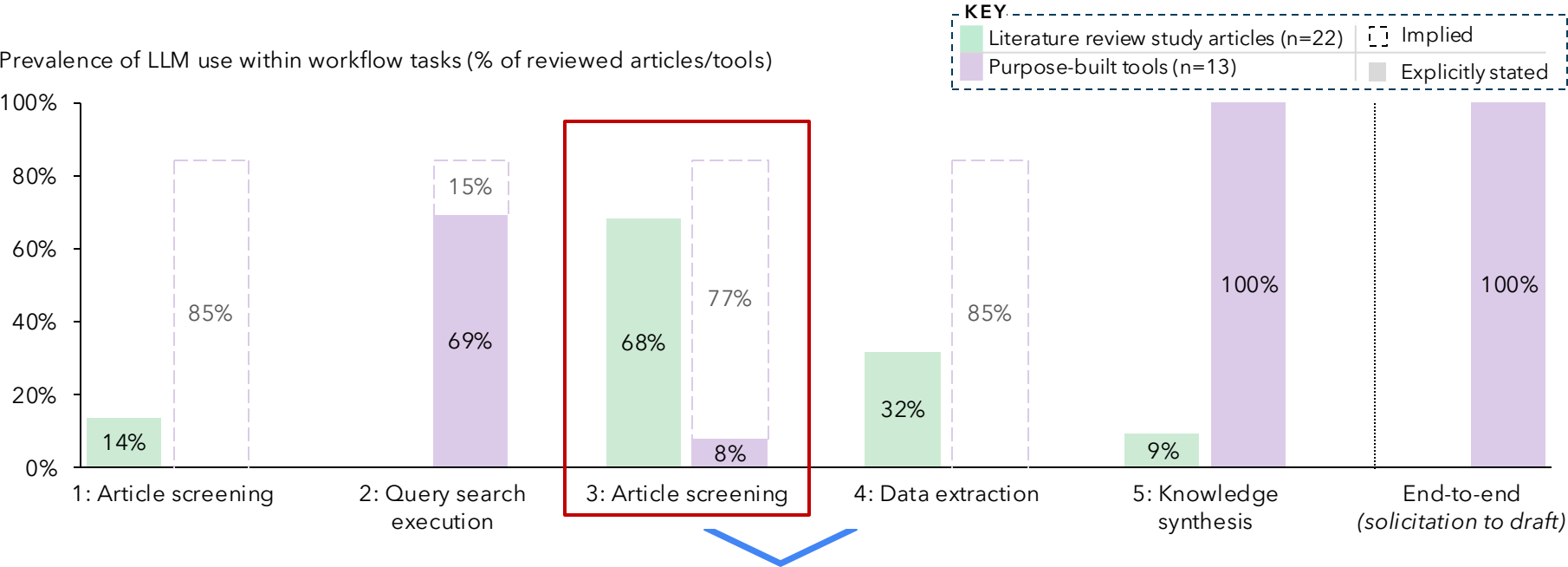


All mapped tools claim end-to-end capabilities while only explicitly stating ~40% of core tasks

Tool	1: Keyword generation and query formulation	2: Query search execution	3: Article screening	4: Data extraction	5: Knowledge synthesis	End-to-end (solicitation to draft)
AutogenAI	•	✓	•	•	✓	✓
AutoRFP.ai	•	✓	✓	•	✓	✓
govdash.com	•	✓	•	•	✓	✓
GOVEAGLE	•	✓	•	•	✓	✓
GovSignals	•	✓	•	•	✓	✓
Grant Assistant	x	x	x	x	✓	✓
Grantable	•	•	•	•	✓	✓
grantboos	x	x	x	x	✓	✓
inventive	•	✓	•	•	✓	✓
loopio	•	✓	•	•	✓	✓
LogAI	•	✓	•	•	✓	✓
responsive	•	✓	•	•	✓	✓
sweetspot	•	•	•	•	✓	✓

KEY ✓ Explicitly stated • Implied x Not applicable

Stated commercial tool functionality diverges markedly from opportunities in the current literature conversation



The gap between marketed commercial tool capabilities and the opportunities identified in literature necessitates robust frameworks for performance evaluation

Summary of key findings



Five core tasks emerge from the unstructured text processing workflow in the literature, with article screening being the most indicated opportunity



Stated commercial tool functionality diverges markedly from opportunities in the current literature conversation



The gap between marketed commercial tool capabilities and the opportunities identified in literature necessitates robust frameworks for performance evaluation

Agenda

Introduction and motivations

Methods

Results and key findings



Future research

We can look towards evaluation techniques in the literature, although varied and ad-hoc, to inform an evaluation framework

Evaluation metric (<i>non-exhaustive</i>)	1: Keyword generation and query formulation	3: Article screening	4: Data extraction	5: Knowledge synthesis
Retrieval	✓			✓
Coverage (% of articles)	✓			✓
Completeness	✓			
Classification	✓	✓	✓	
Accuracy		✓	✓	
Sensitivity / Recall		✓		
Specificity		✓		
F-scores (F1/F2/F3)	✓	✓		
Agreement		✓	✓	
Cohen's kappa		✓		
Prevalence-Adjusted Bias-Adjusted kappa (PABAK)		✓	✓	
Text analysis		✓		✓
Flesch-Kincaid Grade Level		✓		
Simple Measure of Gobbledygook (SMOG)				
Recall-Oriented Understudy for Gisting Evaluation (ROUGE)				✓
Qualitative evaluation			✓	✓
Categorical error typology			✓	
Human / expert assessment				✓

✓ Commonly referenced in literature

Initial results from the literature indicate opportunity but do not show performance matching the current statements from purpose-built tools

1

Keyword generation and query formulation

- Keyword generation yielded **<10% of articles** (*coverage*)³⁸
- Query search terms were **incorrect ~50% of the time** (*completeness*)³¹

3

Article screening

- Higher performing study found **~75% sensitivity; findings vary**²⁸⁻³⁰

4

Data extraction

- *F1-scores* between **mid-70s and 90s**³³⁻³⁵
- Performance was lower for questions requiring inference or synthesis³⁴

5

Knowledge synthesis

- Performance evaluation **ranked** different approaches; no conclusive quality score³⁸

AI SEARCH

First, AI Search finds Relevant Content

AutoRFP.ai

Write

Turn a blank page into a first draft in minutes.

Generate tailored content using your own library, trusted sources, and AI prompts designed for proposals — getting you to a review ready stage faster than ever.

AutogenAI

Winning government proposal drafts in under 30 minutes.

GovSignals



In-Line AI Assistant

Seamless grant writing with AI

- Generate grant content instantly

Grantable

All together, the gap between marketed tool capabilities, the literature, and initial performance necessitates robust evaluation frameworks

**Key
considerations**

- How is a **threshold** for relevant defined?
- Is there a **ground truth**?
- What is the **reference point**?
- What **other techniques** may be useful for evaluation?
- How do you manage performance **trade-offs** (e.g., sensitivity vs. specificity)?
- At what point does an error become **truly problematic**?

Thank you!



Five core tasks emerge from the unstructured text processing workflow in the literature, with article screening being the most indicated opportunity



Stated commercial tool functionality diverges markedly from opportunities in the current literature conversation



The gap between marketed commercial tool capabilities and the opportunities identified in literature necessitates robust frameworks for performance evaluation

A large dark blue rectangular box containing the white text "Q / A" in a bold, sans-serif font.

Q / A