# Towards AI Evaluation in Domain-Specific RAG Systems: The AgriHubi Case Study

Md. Toufique Hasan*, Ayman Asad Khan*, Mika Saari*, Vaishnavi Bankhele*, Pekka Abrahamsson*

*Faculty of Information Technology and Communication Sciences, Tampere University

Tampere, Finland

{mdtoufique.hasan, ayman.khan, mika.saari, vaishnavi.bankhele, pekka.abrahamsson}@tuni.fi

*Abstract*—**Large language models show promise for knowledge-intensive domains, yet their use in agriculture is constrained by weak grounding, English-centric training data, and limited real-world evaluation. These issues are amplified for low-resource languages, where high-quality domain documentation exists but remains difficult to access through general-purpose models. This paper presents *AgriHubi*, a domain-adapted retrieval-augmented generation (RAG) system for Finnish-language agricultural decision support. AgriHubi integrates Finnish agricultural documents with open PORO family models and combines explicit source grounding with user feedback to support iterative refinement. Developed over eight iterations and evaluated through two user studies, the system shows clear gains in answer completeness, linguistic accuracy, and perceived reliability. The results also reveal practical trade-offs between response quality and latency when deploying larger models. This study provides empirical guidance for designing and evaluating domain-specific RAG systems in low-resource language settings.**

*Keywords*—*Empirical Software Engineering, Generative AI, RAG, LLMs, System Design, System Implementation, Agricultural Decision Support, Human-Centered Evaluation*

## I. INTRODUCTION

Agricultural decision making increasingly relies on large and fragmented sources of information, including research reports, regulatory documents, advisory guidelines, and environmental records. Although much of this material is publicly available, practitioners often struggle to access, interpret, and apply it effectively in practice. At the same time, recent advances in generative artificial intelligence (GenAI) have created new opportunities for retrieving and presenting domain knowledge through conversational interfaces powered by large language models (LLMs) [1].

General-purpose LLMs are poorly suited for specialized domains. They often produce confident but weakly grounded answers because they fail to connect domain facts consistently [2]. In agriculture, where advice depends on complex climate–soil relationships, this results in frequent hallucinations due to limited domain-specific training data.

Retrieval-augmented generation (RAG) has emerged as a promising approach to address these shortcomings by combining document retrieval with generative models. Prior work has shown that grounding responses in curated agricultural sources improves both relevance and reliability. Samuel et al. [3] demonstrated that retrieving context from specialized agricultural literature yields more accurate and useful answers than standalone generation. Fanuel et al. [4] further showed that prioritizing region-specific documents during retrieval significantly improves trustworthiness, underscoring the importance of localized knowledge for agricultural applications.

However, a critical gap remains. Most existing agricultural RAG systems focus on English-language data, and there is limited empirical evidence on domain-specific RAG systems designed for low-resource languages and evaluated in real-world settings. Finnish-language agricultural decision support, in particular, has received little attention, despite the availability of high-quality national documentation and the emergence of Finnish-capable open LLMs such as the PORO family.

This paper addresses this gap by presenting *AgriHubi*, a domain-adapted RAG system that integrates Finnish agricultural documents with open LLMs from the PORO family. AgriHubi is designed for real-world use, with explicit source grounding and a built-in feedback mechanism that supports iterative refinement. The work is aligned with national AI and sustainability initiatives in Finland, including *Hiilestä kiinni* and *Ruokavirasto*, which emphasize responsible data use and practical decision support.

The study is guided by the following research questions:

- **RQ1:** How can RAG architecture be adapted to effectively support Finnish-language agricultural data using open-source LLMs?
- **RQ2:** How does model choice within the PORO family affect response quality, latency, and linguistic accuracy in the Finnish agricultural context?
- **RQ3:** How does iterative refinement based on user feedback influence the reliability and usability of a deployed RAG system over time?

To answer these questions, we have designed, implemented, and evaluated AgriHubi through multiple development iterations and two structured user evaluation rounds.

## II. BACKGROUND

This section examines the prior work on multilingual language models, retrieval-augmented generation, and evaluation practices in agricultural AI, highlighting persistent gaps in low-resource language support and empirically validated domain-specific deployments.

The dominance of English in foundational model development has created a prominent performance disparity for smaller languages. To address this English bias, Lai et al. [5] introduced xLLMs-100, demonstrating that scaling multilingual instruction tuning across 100 languages substantially enhances cross-lingual alignment. However, for languages with limited training corpora, language-specific optimization remains vital. Luukkonen et al. [6] showed that multilingual training can benefit low-resource languages, with their Poro 34B model outperforming previous Finnish-focused models despite data scarcity. This aligns with the survey by Ali and Pyysalo [7], which suggests that, while multilingual models are improving, language-specific models typically yield superior performance when sufficient high-quality data is available.

Beyond linguistic challenges, general-purpose models often struggle with the specialized knowledge required for agricultural decision support. A comprehensive review by Yin et al. [8] highlights that adapting foundation models to specific farming contexts, particularly through retrieval-augmented and multimodal approaches, is essential for tasks such as crop disease detection and yield prediction. Data quality plays a pivotal role in this adaptation; Jiang et al. [9] found that constructing knowledge bases from professional literature yields significantly better performance than relying on general internet data. However, localized deployment faces hurdles and, as Owiti and Kipkebut [10] noted, the lack of high-quality annotated datasets in native languages remains a primary bottleneck for training models that genuinely understand local farming terminology.

To mitigate these limitations, recent research has moved toward more sophisticated architectures to enhance reliability. Zhang et al. [11] demonstrated with BeefBot that combining fine-tuning with retrieval-augmented generation (RAG) significantly reduces hallucinations compared to general-purpose models. Further architectural innovations include the "Tri-RAG" framework proposed by Yang et al. [12], which integrates dense retrieval, sparse retrieval, and knowledge graphs to improve reasoning on complex tasks. Moving toward agentic workflows, Yang et al. [13] developed ShizishanGPT, confirming that integrating external tools for specialized tasks, such as phenotype prediction, considerably outperforms standard model capabilities.

Despite these architectural advancements, gaps remain in evaluation and practical deployment. Brown et al. [14] note that evaluation practices are inconsistent, often relying on basic metrics such as Exact Match rather than assessing retrieval quality or grounding effectiveness. Furthermore, Hasan et al. [15] argue that, while theoretical RAG research is abundant, there is a lack of empirical studies on domain-specific systems deployed in real-world settings. Addressing practical deployment challenges, Suvitie et al. [16] showed that hybrid human-in-the-loop workflows can greatly reduce data creation costs. Their approach achieved a cost of about $0.25 per 100 pages. This makes it possible to build high-quality domain datasets without large annotation teams.

These limitations motivated the present work, which focuses on the design, deployment, and evaluation of a domain-specific, Finnish-language RAG system for real-world agricultural decision support.

## III. SYSTEM DESIGN AND ARCHITECTURE

This section describes the design of AgriHubi and how its components work together. The system is organized around data ingestion, retrieval, model inference, and user interaction, forming a complete RAG pipeline that transforms agricultural documents into context-aware answers. The following subsections outline the overall architecture, key components, and underlying technologies.

### A. Overall System Architecture

AgriHubi follows a RAG pipeline with four core parts: document store, retriever, generative model, and user interface. When a user submits a query, the system embeds it and searches a FAISS index built from preprocessed agricultural PDFs. The highest-scoring text chunks are grouped into a context window and sent along with the query to the selected model (`Llama 3.2`, `PORO-34B`, `PORO-2-8B`, or `PORO-2-70B`) through the `Farmi`/`GPT-Lab` APIs.

The Streamlit interface streams the model's response back to the user and logs the query, retrieved passages, and ratings in a local SQLite database. Figure 1 shows the full workflow used in deployment.
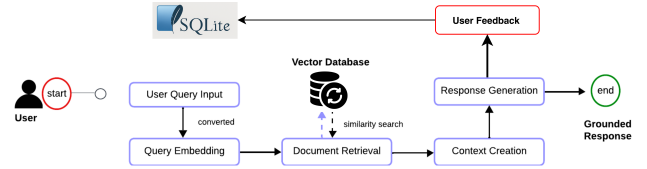


Fig. 1. AgriHubi RAG workflow from query embedding to retrieval, context construction, generation, and feedback storage.

### B. Components

AgriHubi is built from modular components that handle ingestion, retrieval, generation, and user interaction. The UML Class Diagram in Figure 2 shows how these pieces fit together.

*1) Data Ingestion and Preprocessing:* This component converts raw agricultural PDF documents into structured, searchable text. The pipeline includes:

- PDF text extraction with OCR fallback for scanned documents,
- segmentation of extracted text into coherent chunks,
- metadata tagging to preserve document identity and support traceability.

The preprocessing scripts `pdf_to_txt.py` and `txt_to_embedding.py` implement these steps and are part of the internal AgriHubi codebase.
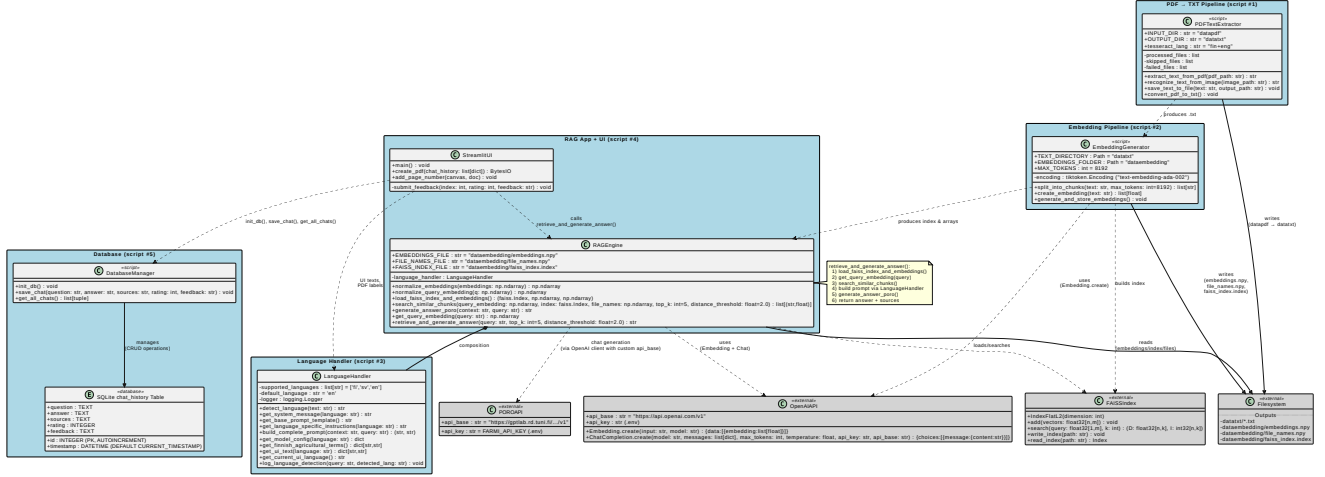
Fig. 2. UML Class Diagram showing the main modules, data flows, and interactions within the AgriHubi system.

*2) Embedding and Indexing:* The system embeds text chunks using OpenAI's `text-embedding-ada-002` model and stores them in a FAISS index. The setup uses L2-normalized vectors, cosine similarity scoring, and then persistent files in `/dataembedding/`. This enables fast and reliable similarity search.

*3) Generative Model Integration:* AgriHubi supports multiple Finnish-capable LLMs across its development timeline:

- `Llama 3.2` for early testing,
- `PORO-34B` for better Finnish output,
- `PORO-2-8B` for multilingual robustness,
- `PORO-2-70B` for the strongest final performance.

The scripts `main_poro34b.py`, `main_poro2_8b.py`, and `main_poro2_70b.py` combine retrieved context with prompt templates and call the models through the OpenAI/Farmi APIs.

*4) User Interface:* The Streamlit-based user interface supports Finnish, Swedish, and English, with automatic translation of interface text handled by a centralized language handler. User interactions follow a chat-style format with streamed responses from the model. The interface also allows users to export conversation histories as PDF files and provides a built-in five-point rating mechanism for collecting structured feedback on answer quality.

*5) Database and Feedback Loop:* The SQLite database (`chat_history.db`) stores each query, the retrieved document chunks, the model's response, the user rating and any written feedback, along with timestamps and model identifiers. This record of interactions formed the core of the feedback loop that guided improvements throughout development.

### C. Technical Stack

AgriHubi uses a Python-based architecture built for clarity, modularity, and reproducibility.

*1) Core Technologies:*

- `Python 3.11.11` for the main application,
- FAISS for vector retrieval,
- OpenAI embedding for text representation,
- Streamlit for the user interface,
- SQLite for persistent evaluation data,
- `dotenv` for secure API key handling.

*2) Model Hosting:* The PORO models (`PORO-2-70B`, `PORO-2-8B`, `PORO-34B`) run on `GPT-Lab/Farmi` servers, so no local GPUs are required.

*3) Deployment, and Supporting Utilities:* AgriHubi can be deployed locally or on laboratory servers using the dependencies listed in `requirements.txt`. Its lightweight, modular design supports containerized setups where frontend and backend components run independently. The system also includes utilities for text preprocessing, database management, and prompt and language handling, with architectural diagrams maintained under `/diagram/`. These components support maintainability, iterative development, and future system extensions.

## IV. IMPLEMENTATION AND ITERATIONS

AgriHubi was developed through eight iterations between January and August 2025. Each iteration introduced targeted improvements based on testing and user feedback. Over time, the system evolved from an initial RAG prototype into a stable, Finnish-optimized platform with improved context handling, terminology, and retrieval reliability.

### A. Iterations 1–3: Establishing the Core System

The first three iterations focused on building the core RAG pipeline. We added PDF extraction, OCR for scanned documents, text chunking, embedding generation, and a FAISS index for semantic retrieval. The early system relied on `Llama 3.2`. It worked for initial testing but struggled with Finnish agricultural terminology. Iteration 2 introduced `PORO-34B`, which immediately improved fluency

and domain coverage. We also released the first Streamlit interface so internal testers could try the system directly. In Iteration 3, the document collection was expanded and Finnish terminology refined based on tester feedback. These stages form the foundation of the repository's `main` branch.

### B. Iterations 4–6: Feedback and Model Refinement

Iteration 4 added a built-in feedback mechanism with a five-point rating scale stored in SQLite. This gave us clear insight into answer quality. The collected data showed issues with accuracy, incomplete reasoning, and occasional ambiguity. Iteration 5 included the first structured external evaluation (67 responses). Almost half of the answers were rated 1–2, which led to changes in the retrieval strategy, chunking logic, and similarity thresholds. During this phase, we also tested the `Viking-33B` model, but its responses were short and unstable, so it was removed. Iteration 6 transitioned the system to `PORO-2-8B`, which produced more stable Finnish and Swedish output and handled terminology more consistently. Preprocessing and embedding generation were also improved. These updates are captured in the `v6.0` branch.

### C. Iterations 7–8: Finalizing the System

Iteration 7 introduced `PORO-2-70B`, the strongest model in terms of accuracy and completeness. Its answers were longer, clearer, and more grounded in retrieved context. Supporting it required backend upgrades, including larger model calls, improved streaming, and updated timeout logic. These appear in the `v7.0` branch in the internal AgriHubi codebase. Iteration 8 concentrated on usability and system stability. We replaced filename-based retrieval with full-chunk retrieval, increased the maximum answer length from 700 to 2000 tokens, improved error handling, and fully localized the interface into Finnish. These updates addressed repeated user requests for more complete answers. The `v8.0` changes led to a clear jump in evaluation scores in the second testing round (47 responses), especially regarding completeness and clarity.

User feedback guided most design decisions. Testers reported problems with terminology, missing context, unclear answers, and slow responses. Their ratings led to changes in retrieval settings, model selection, interface design, and server setup. As a result, the system was improved through repeated testing and refinement rather than a fixed development process.

## V. EVALUATION AND RESULTS

This section explains how we evaluated AgriHubi and what the results show. Two user studies and several internal tests reveal how the system's accuracy, fluency, and responsiveness changed over time. Together, these findings highlight where the system has been improved and where limitations remain.

### A. Methodology

AgriHubi was evaluated in two rounds in April and August 2025. In both rounds, users asked real agricultural questions and rated the answers on a five-point Likert scale (1 = poor, 5 = excellent). The system also logged each query, retrieved document segments, model responses, and timestamps in a local SQLite database, enabling consistent comparison across iterations. The April evaluation tested the `PORO-34B` version and collected 67 ratings. The August evaluation assessed the updated system using `PORO-2-70B`, and in some cases `PORO-2-8B`, with 47 ratings. Question topics were kept similar to support direct comparison between the two rounds.

Across both phases, four primary evaluation metrics were used:

- user-perceived quality via Likert ratings,
- end-to-end response latency,
- factual accuracy and completeness of generated answers,
- fluency of Finnish agricultural terminology.

These were complemented by internal system measurements and qualitative feedback from users.

### B. Quantitative Results

**Round 1 (April 2025, `PORO-34B`).** Table I summarizes the 67 collected ratings. The distribution was as follows: 15 responses (22%) with rating 1, 16 responses (24%) with rating 2, 17 responses (25%) with rating 3, 17 responses (25%) with rating 4, and 2 responses (3%) with rating 5. Almost half of the answers (46%) therefore fell into the lowest two categories, while only 3% achieved the top score. Testers most often reported missing details, weak factual grounding, and inconsistent Finnish terminology.

**Round 2 (August 2025, `PORO-2-70B`).** Table I also reports the 47 ratings collected with the updated system. Here, 9 answers (19%) received rating 1, 9 answers (19%) rating 2, 15 answers (32%) rating 3, 4 answers (9%) rating 4, and 10 answers (21%) rating 5. Compared to Round 1, the share of low ratings (1–2) dropped from 46% to 38%, while the proportion of top-rated answers (score 5) increased from 3% to 21%. Ratings of 3 also increased from 25% to 32%, suggesting more stable and complete responses.

TABLE I
COMPARISON OF LIKERT EVALUATION RESULTS FOR APRIL AND AUGUST 2025

| Rating | April 2025 | | August 2025 | |
|---|---|---|---|---|
| | Responses | % | Responses | % |
| 1 | 15 | 22% | 9 | 19% |
| 2 | 16 | 24% | 9 | 19% |
| 3 | 17 | 25% | 15 | 32% |
| 4 | 17 | 25% | 4 | 9% |
| 5 | 2 | 3% | 10 | 21% |

**Comparison.** User ratings improved markedly between April and August. Low ratings (1–2) decreased from 46%

to 38%, while top ratings (5) increased from 3% to 21%. Mid-range ratings (3) rose from 25% to 32%, and ratings of 4 declined from 25% to 9%.

### C. Internal Performance Evaluation

Internal evaluations were conducted across all major model backends: `LLaMA`, `Viking-33B`, `PORO-34B`, `PORO-2-8B`, and `PORO-2-70B`. They were assessed along four criteria:

- **User feedback:** average rating stability and distribution,
- **Qualitative errors:** factual accuracy and evidence grounding,
- **Latency:** end-to-end response time,
- **Domain coverage:** robustness of Finnish agricultural terminology.

`PORO-34B` delivered the first stable Finnish output but often produced incomplete answers. `PORO-2-8B` improved fluency with moderate latency, while `PORO-2-70B` achieved the highest accuracy at the cost of increased memory use and latency. The `Viking-33B` model was discontinued due to unstable, error-prone outputs, and `LLaMA` was used only as an early baseline. These internal results confirm a predictable trade-off: larger models yield higher answer quality but at the cost of increased latency and resource consumption.

### D. Qualitative Feedback

Qualitative feedback helped explain the rating trends. In the first round, users often described answers as incomplete, overly general, inconsistent in Finnish terminology, or slow to generate. These comments match the low scores observed in April. In the second round, feedback shifted noticeably. Users described the system as clearer and more reliable, noting improved terminology, more natural Finnish, and better use of retrieved sources. Some issues remained, including occasional errors in highly technical questions and weaker performance for Swedish queries. Overall, the feedback shows a clear transition from a prototype to a usable decision-support tool, with steady gains in answer quality, completeness, and linguistic accuracy across iterations.

## VI. Discussion

The evaluation shows that AgriHubi improved steadily across the development iterations. The strongest gains came from system-level refinements rather than model scaling alone. Improvements in retrieval quality, preprocessing, and context construction had a greater impact on answer quality than switching to larger models. This confirms that stable system design is critical for domain-specific RAG systems.

### A. Technical Observations

Most technical improvements resulted from better retrieval and preprocessing. Early versions suffered from OCR noise, uneven chunking, and missing metadata, which prevented models from grounding answers correctly. Once these issues had been addressed, the PORO models used retrieved content more effectively and produced more complete responses.

A clear trade-off emerged between accuracy and latency. `PORO-2-70B` generated the most accurate and detailed answers but responded more slowly, especially under shared GPU load. Testers often associated slower responses with lower trust, showing that latency affects both usability and perceived reliability.

### B. Domain Impact

AgriHubi showed clear progress in handling Finnish agricultural terminology. Testers reported more natural phrasing, more accurate use of domain concepts, and stronger alignment between answers and retrieved evidence. These gains resulted from both improved model capability and refinements in prompts and language handling.

Support for Swedish remained limited. Sparse training data and inconsistent document coverage led to unstable responses, reducing trust among bilingual users. This highlights the need for more balanced multilingual data and retrieval support in future versions.

### C. Limitations

Several limitations influenced the evaluation. Hardware constraints limited the throughput of `PORO-2-70B`, and variable response times affected user ratings. The sample sizes of 67 and 47 responses provided useful insights but were insufficient for fine-grained statistical analysis. Privacy concerns also required care, as retrieved content sometimes included verbatim text from internal documents. In addition, differences in question types between evaluation rounds may have introduced minor scoring bias.

### D. Lessons Learned

Three lessons stand out. First, incremental improvements driven by user feedback were more effective than large architectural changes. Second, human review was essential for identifying subtle errors in meaning, terminology, and context that automated signals could not reliably detect, especially in a low-resource language setting. Third, the lack of a consistent evaluation process across iterations made structured comparison difficult, pointing to the need for more systematic evaluation practices in future work.

Overall, the results show that AgriHubi evolved from an early prototype into a more reliable decision-support system. While gains in answer quality and user trust are clear, challenges related to latency, multilingual consistency, and production readiness remain open.

## VII. Future Work

Future work will focus on expanding the system's data ecosystem, strengthening contextual reasoning, and extending AgriHubi toward action-oriented decision support. A

natural next step is deeper integration with external data sources, including sensor-based and API-driven inputs such as weather services, soil measurements, and environmental monitoring systems. Incorporating such real-time signals would allow the retrieval layer to reason not only over static documents but also over current conditions, enabling more situationally aware responses. Another important direction is the introduction of longer-term system memory. Persisting interaction history, retrieved context, and past recommendations would allow AgriHubi to maintain continuity across sessions and support cumulative reasoning over time. This form of contextual memory can improve consistency, reduce repeated explanations, and enable more informed follow-up guidance, especially in seasonal or longitudinal agricultural scenarios. Building on these capabilities, AgriHubi can gradually evolve from an information-centric RAG system into a broader decision-support platform. Future extensions could support lightweight planning tasks, scenario exploration, and recommendation workflows that connect retrieved knowledge with concrete actions, such as crop management suggestions or policy-compliant guidance. Achieving this vision will require careful data governance, controlled update cycles, and explicit handling of uncertainty when multiple data sources provide conflicting signals. Together, these directions will move AgriHubi beyond document question answering toward a grounded, context-aware agricultural AI system that supports real-world decision making.

## VIII. Summary

This paper introduced AgriHubi, a domain-specific retrieval-augmented generation system developed to support Finnish-language agricultural decision making. By integrating Finnish agricultural documents with open large language models from the PORO family, AgriHubi demonstrates how localized retrieval and language-aware generation can improve the reliability, clarity, and practical usefulness of AI-assisted decision support. The system was developed through eight iterative design cycles and evaluated in two structured user studies. Over time, AgriHubi progressed from an early prototype into a more stable and dependable system. Quantitative ratings and qualitative feedback show clear improvements in answer completeness, linguistic accuracy, and user trust, particularly after refinements to the retrieval pipeline and the transition to larger PORO models. The evaluation also revealed tradeoffs between response quality and latency, reinforcing the central role of retrieval quality in grounding model outputs.

Beyond the Finnish agricultural context, this study offers practical insights into how domain-adapted RAG systems can be engineered, deployed, and evaluated in applied settings. The results highlight the importance of iterative refinement, user-centered evaluation, and transparent grounding when applying large language models in domain-specific and high-stakes environments. Taken together, the findings position AgriHubi both as a func-

tional decision-support system and as a reference case for future research on multilingual, domain-specific retrieval-augmented generation.

## References

[1] A. Nguyen-Duc, B. Cabrero-Daniel, A. Przybylek, C. Arora, D. Khanna, T. Herda, U. Rafiq, J. Melegati, E. Guerra, K.-K. Kemell *et al.*, "Generative artificial intelligence for software engineering-a research agenda," *Software: Practice and Experience*, vol. 55, no. 11, pp. 1806–1843, 2025.

[2] B. Jiang, Y. Wang, Y. Luo, D. He, P. Cheng, and L. Gao, "Reasoning on efficient knowledge paths: knowledge graph guides large language model for domain question answering," in *2024 IEEE International Conference on Knowledge Graph (ICKG)*. IEEE, 2024, pp. 142–149.

[3] D. J. Samuel, I. Skarga-Bandurova, D. Sikolia, and M. Awais, "Agrollm: Connecting farmers and agricultural practices through large language models for enhanced knowledge transfer and practical application," *arXiv preprint arXiv:2503.04788*, 2025.

[4] M. Fanuel, M. N. Mahmoud, C. C. Marshal, V. Lakhotia, B. Dari, K. Roy, and S. Zhang, "Agriregion: Region-aware retrieval for high-fidelity agricultural advice," *arXiv preprint arXiv:2512.10114*, 2025.

[5] W. Lai, M. Mesgar, and A. Fraser, "Llms beyond english: Scaling the multilingual capability of llms with cross-lingual feedback," *arXiv preprint arXiv:2406.01771*, 2024.

[6] R. Luukkonen, J. Burdge, E. Zosa, A. Talman, V. Komulainen, V. Hatanpää, P. Sarlin, and S. Pyysalo, "Poro 34b and the blessing of multilinguality," in *Proceedings of the joint 25th nordic conference on computational linguistics and 11th baltic conference on human language technologies (nodalida/baltic-hlt 2025)*, 2025, pp. 367–382.

[7] W. Ali and S. Pyysalo, "A survey of large language models for european languages," *arXiv preprint arXiv:2408.15040*, 2024.

[8] S. Yin, Y. Xi, X. Zhang, C. Sun, and Q. Mao, "Foundation models in agriculture: A comprehensive review," *Agriculture*, vol. 15, no. 8, p. 847, 2025.

[9] J. Jiang, L. Yan, H. Liu, Z. Xia, H. Wang, Y. Yang, and Y. Guan, "Knowledge assimilation: Implementing knowledge-guided agricultural large language model," *Knowledge-based systems*, vol. 314, p. 113197, 2025.

[10] T. L. Owiti and A. K. Kipkebut, "Enhancing ai-driven farming advisory in kenya with efficient rag agents via quantized fine-tuned language models," in *Proceedings of the Sixth Workshop on African Natural Language Processing (AfricaNLP 2025)*, 2025, pp. 24–30.

[11] Z. Zhang, C.-A. Wilson, R. Hay, Y. Everingham, and U. Naseem, "Beefbot: Harnessing advanced llm and rag techniques for providing scientific and technology solutions to beef producers," in *Proceedings of the 31st International Conference on Computational Linguistics: System Demonstrations*, 2025, pp. 54–62.

[12] B. Yang, Y. Zhang, L. Feng, Y. Chen, J. Zhang, X. Xu, N. Aierken, Y. Li, Y. Chen, G. Yang *et al.*, "Agrigpt: A large language model ecosystem for agriculture," *arXiv preprint arXiv:2508.08632*, 2025.

[13] S. Yang, Z. Liu, W. Mayer, N. Ding, Y. Wang, Y. Huang, P. Wu, W. Li, L. Li, H.-Y. Zhang *et al.*, "Shizishangpt: An agricultural large language model integrating tools and resources," in *International Conference on Web Information Systems Engineering*. Springer, 2024, pp. 284–298.

[14] A. Brown, M. Roman, and B. Devereux, "A systematic literature review of retrieval-augmented generation: Techniques, metrics, and challenges," *arXiv preprint arXiv:2508.06401*, 2025.

[15] M. T. Hasan, M. Waseem, K.-K. Kemell, A. A. Khan, M. Saari, and P. Abrahamsson, "Engineering rag systems for real-world applications: Design, development, and evaluation," in *Software Engineering and Advanced Applications*, D. Taibi and D. Smite, Eds. Cham: Springer Nature Switzerland, 2026, pp. 143–158.

[16] N. Suvitie, M. Saari, and P. Abrahamsson, "From PDF to Dataset: Semi-Automatic Extraction of Fine-Tuning Data," in *49th ICT and Electronics Convention – MIPRO 2026*. IEEE, 2026, submitted.