# PDMR CHAT BOT

**Team Members:** Nathan Steurer, Daniel Harradine, Hruthin Muddasani, Yu-Jen Chen, John Graham Chiles, Tejomay Marathe, Vachan Arora, Ken Thai Nguyen, Madhavan Prasanna, Sharan Suthaharan

Turfgrass Science at Purdue
Department of Horticulture & Landscape Architecture
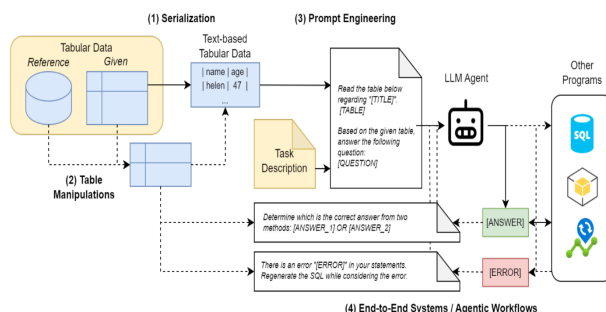
MathWorks

## The Purpose of this Project
- This project aims to create a learned language model (LLM) to optimize the search for turf grass disease treatments using 26 years of published field trials stored in Plant Disease Management Reports (PDMRs).
- Currently, the turfgrass treatment resides in the PDMR's that must be manually searched – a cumbersome process.
- By creating an LLM with the ability to search through PDMR's and create a database of the information, we improve the accessibility of the information for stakeholders.

## Our Plan
To increase the ability for an LLM to properly read the PDMR's so it can give accurate responses we split the team into three different groups.
- Parsing – First, for the LLM to more easily read the PDMR's the files needed to be more standardized.
- Database – Next, the PDMR's need to be stored in a data so the LLM could easily search for relevant PDMR's needed for the response.
- LLM – The LLM itself needed to be engineered in a way that it would look for not only the correct treatment but the most recent one.

## Research Into The Project

### Parsing of the PDMRs
Throughout this semester the parsing team looked at three different ways of parsing the PDMR's:
- Custom Computer Vision (CV) solution, splitting the work into five stages: Fetching, Parsing, Data Cleaning 1, Data Cleaning 2, and Output generation.
- Matlab's Text Analytics Toolbox
- ChatGPT to parse the PDMRs and transform tables into CSV files.

### Storing The PDMR's
The best form for storing the CSV files was on a vector database. Which would allow the LLM to more easily find relevant files to any question. This was done through a myriad of tools.
- Docker
- VS Code.
- Milvus

The first step was to code a vector database. Then to encode the data into the database. The final step was to create indicator's allowing for files to be queried faster.

### Training the LLM
After much deliberation between open-source or closed-source, we eventually choose an open-source because of the reduced cost.
- The LLM that we settled with was lama 3.2. We were able to more easily prompt engineer the LLM by using Ollama to run it without many limits.
- To decrease the amount of time it takes the LLM to respond we had Purdue's Anvil host the Ollama where we could import our desired LLM.

## Parsing Team: ①
- After handling text extraction from PDMRs with Computer Vision, we clean it up and separate the paragraphs
- Trained the learning algorithm of ChatGPT to build accuracy
- Used a pipeline into ChatGPT from Python to handle different table formats in PDMRs

## Database Team: ②
- To improve LLM performance, we created a vector database for retrieval augmented-generation (RAG)
- Hosted the vector and SQL database of PDMR's on Docker and Milvus as backend
- Training LLM using unsupervised learning on vector database
- Establishing public hosting on Anvil for easy connection/ setup

*(1) Serialization — (3) Prompt Engineering*

Tabular Data / Text-based Tabular Data / LLM Agent / Other Programs

*(2) Table Manipulations*

*(4) End-to-End Systems / Agentic Workflows*

## User Interface
- Created user interface (UI) using pre-existing tooling in order to achieve professional results akin to the DeepSeek or ChatGPT interface.
- Utilized tools such as AI-SDK along with Vercel in order to simplify and speed up deployment process, while also enabling easy switching of LLMs if necessary.

## LLM Team: ③
- Used database to query the LLM using RAG.
- Fine-tuning and prompt engineering to ensure correct output.
- Researched state of the art methods to handle tabular data as well as textual data.
- Conducted exhaustive testing with a set question bank in order to verify outputs

## Conclusion
- We were able to increase much in the two semesters that we worked on this project. Finding the many different challenges working with LLM's.
- There are many different solutions to the problems that we faced. Leading multiple branched methods of parsing, storing, and reading through the LLM.
- However, in the end we were able to successfully parse the PDMR's into readable CSV files. Also being able to import them into a vector datable using Milvus. With the final out from the LLM at least being more capable compared to the previous version.

## Future Goals
- The most pressing goal would be continuing to improve the way the PDMR's can be parsed in a more efficient method, handling edge cases such as tables spanning multiple pages, or non-standardized formats.
- Prompt engineer the LLM to give better responses and reduce the chance that the LLM gives a false response or hallucinates.
- Generalize the solution to account for other types of PDMRs, not just turfgrass.
- Fine-tune response times by resolving bottlenecks in processing and inference.
- Enable LLM to help perform data analysis and other tasks on the PDMRs as directed by the user.

| Treatment and rate/1000 ft² | 'U-3' Phytotoxicity[Z] | Quality[Y] | 'Riviera' Phytotoxicity | Quality | 'Tifway 419' Phytotoxicity | Quality |
|---|---|---|---|---|---|---|
| Non-fungicide Treated Control[X] | 0.88 D | 8.57 A | 1.37 | 8.25 | 1.50 | 8.25 |
| Trinity 1.67 SC 0.5 fl oz | 1.13 CD | 8.44 AB | 1.63 | 8.13 | 1.50 | 8.25 |
| Trinity 1.67 SC 0.75 fl oz | 1.13 CD | 8.44 AB | 1.63 | 8.13 | 1.75 | 8.13 |
| Trinity 1.67 SC 1.0 fl oz | 1.00 D | 8.50 A | 1.63 | 8.00 | 1.75 | 8.13 |
| Trinity 1.67 SC 1.5 fl oz | 1.38 BC | 8.31 BC | 1.50 | 8.31 | 2.00 | 8.00 |
| Trinity 1.67 SC 2.0 fl oz | 1.13 CD | 8.44 AB | 1.63 | 8.13 | 1.75 | 8.13 |
| Banner MAXX 1.24 MEC 2.0 fl oz | 2.13 A | 7.81 D | 1.63 | 8.13 | 2.00 | 8.00 |
| Bayleton 50 WG 2.0 oz | 2.13 A | 7.88 D | 1.63 | 8.13 | 2.00 | 8.00 |
| Eagle 40 WP 1.2 oz | 1.63 B | 8.19 C | 1.50 | 8.13 | 2.00 | 8.00 |

[Z] Phytotoxicity was assessed on 7-Aug and based on a scale of 0 - 10 where 0 = no change, 1 = slight change in turfgrass appearance, 5 = distinct change in turfgrass appearance, 10 = very abnormal turfgrass. Means followed by the same letter are not significantly different according to Fisher's test of protected least significant difference where; LSD=0.51; $P$-value=0.02.
[Y] Turfgrass quality was assessed on 7-Aug and based on a scale of 1 - 9 where 1 = no turf present, 5 = unacceptable turfgrass, 7 = acceptable turf, 9 = dense, dark color, thick stand of turfgrass. Means followed by the same letter are not significantly different according to Fisher's test of protected least significant difference where; LSD=0.23; $P$-value=0.01.
[X] Fungicide application was the sub-plot effect of the split-block experimental design.

## Acknowledgements
We would like to thank Dr. Miller, Gen Sasaki, Carmen Lee, and Sathvik Hegde for all the help on this project. We would also like to thank the staff at Data Mine team with all the support they have gave us.

## References
Ollama. (n.d.). *Ollama Documentation.* Retrieved from https://github.com/ollama/ollama/blob/main/docs/README.md
Milvus. (n.d.). *Milvus Documentation.* Retrieved from https://milvus.io/docs
All the various turfgrass related PDMRs.
Sui, Y., Zhou, M., Zhou, M., Han, S., & Zhang, D. (2024, March). Table Meets LLM: Can Large Language Models Understand Structured Table Data? A Benchmark and Empirical Study. *The 17th ACM International Conference on Web Search and Data Mining (WSDM '24).*

**The Data Mine Corporate Partners Symposium 2025**