

Name: Jake Brulato

Date: 28 Feb, 2025

Project Title: Prompt Engineering and Framework Conversion for MRM

Mentor and Company: Anwesha Bhattacharyya, Wells Fargo

Period Covered: 01 Jan – 17 January 2025

Project Objective: Developing a centralized repository of tailored prompts for Large Language Models, specifically optimized for the DSAI team within Wells Fargo's MRM sector, incorporating a ReAct (Reason + Action) framework for open-source LLMs after prompt hub creation.

Progress Summary:

Work Completed:

- Task 4.2 Building ReAct and Testing with Public Data (Complete)
  - Summary:
    - Internal package conflicts are still occurring, and came to conclusion that the tools for the ReAct agent needed to be built in-house.
    - Created a RAG framework using Llama 2 on local machine with defined functions to process multiple document types using external data to test and chunk it.
    - Implemented a multi-step “thought”, “action, and “reflection” based on the user query.
  - Lessons Learned:
    - Llama 2 may not be strong enough to handle multi-document processing and will need to be changed to Llama 3.1 or at least Llama 3.
    - Different embedding models change the speed of processing chunked data, using a better model is recommended.
  - Risk & Mitigation:
    - Risk: File structure internally may differ from normal a normal pdf or csv in the future, possibly causing issues in the document processing/chunking of data.
    - Mitigation: Added to documentation example file usage and how it would work if used with internal documents.
- Task 5.2 Testing of Visual Capabilities and Text Generation (Complete)
  - Summary:
    - Completed testing of Deepseek-R1 distilled Qwen, SmolVLM, and Qwen2.5 (3B & 7B) with completed examples for use by the ATOM GenAI team.
    - Consolidated to a report showing the text capabilities, nuances, and image/video performance across all models with potential recommendations.
  - Lessons Learned:
    - Model performance isn’t entirely limited to the parameter amount used to train the model.

- Smaller models can have the potential to do great in one task compared to another that can do multiple tasks well.
- Risk & Mitigation:
  - Risk: Qwen 2.5 and Distilled Deepseek are still very experimental and found to be prone to some hallucinations when prompted in long context.
  - Mitigation: Wrote out specifics that a person should avoid when model is in use, provided proper files for image resizing and processing as well.

#### Work Scheduled Next Sprint

- Task 4.3 Add ReAct to Loan Harmonization Models
  - Refine the user query to completely follow the Loan Harmonization project guidelines based on example work from another division.
  - Use RAG for the text extraction of fed files through previous outputs. Make it think critically and adapt with each iteration.
- Task 6. Present Findings on Generative Models in AIA Seminar
  - Present in the weekly seminar that MRM hosts the proper guidelines and functionality of core models used in the prompt hub.
- Risk & Mitigation:
  - Risk: Metadata and guidelines may be too complex for the current model to process in a long document string. May cause hallucinations if not properly checked.
  - Mitigation: Pull similarity score from chunks that were used and train a model if need be. If scores are too low, request more resources or switch to another model with a better training method.

#### Work Package Updates:

	WP	Work Package Title	Scheduled Dates	Status
✓	<del>1.</del>	<del>Onboarding</del>	<del>28<sup>th</sup> Oct — Nov 22<sup>nd</sup></del>	<del>Complete</del>
✓	<del>2.</del>	<del>Testing performance of Large Language Models and Configuration</del>	<del>Nov 25<sup>th</sup> — Dec 13<sup>th</sup></del>	<del>Complete</del>
✓	<del>3.</del>	<del>Creation of Prompt Hub for DSAI and MRM</del>	<del>Dec 16<sup>th</sup> — Jan 30<sup>th</sup></del>	<del>In-Progress</del>
✓	<del>3.1.</del>	<del>Research Model Specifications and Prompting Techniques</del>	<del>Dec 16<sup>th</sup> — Dec 18<sup>th</sup></del>	<del>Complete</del>
✓	<del>3.2.</del>	<del>Testing Prompts around research</del>	<del>Dec 16<sup>th</sup> — Dec 20<sup>th</sup></del>	<del>Complete</del>
✓	<del>3.3.</del>	<del>Creation of Prompt Hub</del>	<del>Dec 18<sup>th</sup> — Dec 30<sup>th</sup> Jan 30<sup>th</sup></del>	<del>Backlog</del>
	4.	Implement ReAct Framework	Jan 1 <sup>st</sup> — Mar 28 <sup>th</sup>	In-Progress
✓	4.1.	Meet with Team to Discuss ReAct for Loan Harmonization	Jan 13 <sup>th</sup> — Jan 17 <sup>th</sup>	Complete
✓	4.2.	Building ReAct and Testing with Public Data	Jan 6 <sup>th</sup> — Jan 17 <sup>th</sup> Feb 16 <sup>th</sup> Feb 22 <sup>th</sup>	Complete

	4.3.	Add ReAct to Loan Harmonization Models	Jan 17 <sup>th</sup> – Mar 28 <sup>th</sup>	In-Progress
✓	5.	<del>More Generative Model Development</del>	<del>Jan 28<sup>th</sup> – Mar 28<sup>th</sup></del>	<del>Complete</del>
✓	5.1.	<del>Research Chat Template for Visual/Large Language Models</del>	<del>Jan 31<sup>st</sup> – Feb 5<sup>th</sup></del>	<del>Complete</del>
✓	5.2.	<del>Testing of Visual Capabilities and Text Generation</del>	<del>Feb 5<sup>th</sup> – Feb 14<sup>th</sup></del>	<del>Complete</del>
	6.	Present Findings on Generative Models in AIA Seminar	Feb 28 <sup>th</sup> – Mar 13 <sup>th</sup>	In-Progress