Artificial Intelligence Data Specialist Level 7 Higher Apprenticeship EPA Project Report

# MVP Implementation
# of a Pensions Expert Chatbot App
# at COMPANY Retirement Solutions

**BY**

**ARKADIUSZ KULPA**

<u>SEE GLOSSARY IN APPENDIX FOR DISAMBIGUATION OF ACRONYMS</u>

# 1   INTRODUCTION AND BACKGROUND

After an initial proposal of a RAG LLM to help deal with COMPANY extensive knowledge compartmentalisation problem, as observed by the learner across project VEGA and Shareholder Solutions in general, a Application RAG LLM use case has been identified for an MVP development, which is the subject of this report. The Chatbot has been developed using Databricks Platform, Lang Chain vectorised documentation set and Llama 3.1 70B Instruct model served as a service, to explore viability of such a tool in providing quality answers about a pension scheme based on the scheme's Trust Deed and Rules and other guideline documents available in public domain.

To ensure this MVP targets the correct problem, it was properly articulated and outlined as business requirement, through numerous consultations, walkthroughs and reviews that were arranged by the learner, including both technical and non-technical stakeholders. These included Application Head of Platform, Director of Data and CPTO who sponsored the MVP with the vision of similar solutions being required across the group and the director of Product for RS who showcased it to the LGPS scheme client to gauge their interest and secure future investment.

Meanwhile the AI Guild and HIVE sessions showcased the developments to a wider and diverse audience of people across COMPANY who participated in Q&As about the tool and were explained the potential benefits of such tools to front-line staff's day to day duties.

## 2  OUTLINE OF THE BUSINESS PROBLEM TO BE SOLVED

EQ currently relies on experienced pension administrators to interpret industry jargon and manually analyse documentation to answer member queries (calls, emails, etc.) (BR1) and execute internal processes, such as onboarding new pension schemes via Application (BR2). Onboarding involves analysing scheme documentation, categorising members (e.g., standard vs pilot), and defining calculations for various actions like adjusting contributions or handling pensions in case of divorce or death.

Given this reliance on documentation and ongoing analysis, a RAG LLM was identified as an appropriate AI solution, with the LGPS scheme selected as MVP example. While fine-tuning could enhance the model's understanding of pension-specific jargon, it was deferred until after initial implementation due to the need for extensive client documentation. In contrast, RAG was developed using publicly accessible sources. Two chatbots - public and internal – are to both be developed using similar technical solutions. former would provide accurate and safe general responses, while the latter would require greater precision to address complex scenarios within the schemes, necessitating post MVP development.

The learner participated in an AI focused conference, where COMPANY CPTO and AI Champions outlined an AI roadmap for COMPANY including optimal solutions. The learner proposed three strategic options: (1) leveraging third-party solutions (e.g., Palantir, ChatGPT, Copilot for 365, or Luminance), which promise faster deployment but lack customization, require significant investment, and risk vendor lock-in; (2) building a solution from scratch, offering complete control over data and processes but incurring high costs and time demands; or (3) taking the middle road, by utilizing a cloud solution with ML capabilities and tools. Databricks emerged as the optimal platform choice, as it offers SOTA ML tools as a service, in a fully customisable, modular, and scalable cloud environment. Already used for COMPANY Transformation Project and Data Platform, Databricks ensures robust data governance, enhancing any AI solutions developed within its framework.

# 3 METHODS USED & JUSTIFICATION

Learner worked with COMPANY infrastructure tech team to set up a secure, serverless Databricks development environment. Early collaboration resolved issues like UFC worker nodes lacking internet access, fixed during embedding and tokenization by the learner. Autonomous work continued within Databricks IDE and UI, where basic RAG LLM notebooks available from ai-cookbook.io. were adapted for this mvp. This setup balanced flexibility in tool and technique selection with a robust end-to-end process - from preprocessing and data pipelines to iterative experimentation, review apps and front-end UI (DataBricks, 2024).



*Figure 1 - Databricks RAG LLM pipeline (Databricks, 2024)*

The Unity Catalog, part of Databricks' governance solution for the Lakehouse platform, was used to store structured and unstructured data in one place, organised using medallion architecture. Lang Chain, a leading framework for building controllable agentic workflows (LangChain, 2024), served as the vector store retrieval method and seamlessly integrated into the DataBricks ML flow pipeline. Together, these tools offered full data security, restricting access to specific documents (parameter: TRIGGERED in figure below) and designated individuals.

```python
data_pipeline_config = {
    # Vector Search index configuration
    "vectorsearch_config": {
        # Pipeline execution mode.
        # TRIGGERED: If the pipeline uses the triggered execution mode, the system stops processing after successfully refreshing the source table in the pipeline once,
        ensuring the table is updated based on the data available when the update started.
        # CONTINUOUS: If the pipeline uses continuous execution, the pipeline processes new data as it arrives in the source table to keep vector index fresh.
        "pipeline_type": "TRIGGERED",
    },
    # Embedding model to use
    # Tested configurations are available in the `supported_configs/embedding_models` Notebook
    "embedding_config": {
        # Model Serving endpoint name
        "embedding_endpoint_name": "databricks-gte-large-en",
        "embedding_tokenizer": {
            # Name of the embedding model that the tokenizer recognizes
            "tokenizer_model_name": "Alibaba-NLP/gte-large-en-v1.5",
            # Name of the tokenizer, either `hugging_face` or `tiktoken`
            "tokenizer_source": "hugging_face",
        },
    },
    # Parsing and chunking configuration
    # Changing this configuration here will NOT impact your data pipeline, these values are hardcoded in the POC data pipeline.
    # It is provided so you can copy / paste this configuration directly into the `Improve RAG quality` step and replicate the POC's data pipeline configuration
    "pipeline_config": {
        # File format of the source documents
        "file_format": "pdf",
        # Parser to use (must be present in `parser_library` Notebook)
        "parser": {"name": "pypdf", "config": {}},
        # Chunker to use (must be present in `chunker_library` Notebook)
        "chunker": {
            "name": "langchain_recursive_char",
            "config": {
                "chunk_size_tokens": 1024,
                "chunk_overlap_tokens": 256,
            },
        },
    },
}
```

*Figure 2 - Data pipeline config*

The setup addressed BRs by using latest Alibaba-NLP model from hugging face, offering SOTA 8k token limit for tokenisation and embedding. PDF transformation was handled by the 'PyPdf' parser and 'langchain_recursive_char' chunker, breaking documents into 1,024-token chunks with 256-token overlap. In a typical RAG LLM approach, each prompt was augmented with chat history, relevant previous inputs, and newly retrieved context.

MLFlow established a 'run', creating a custom iteration of the RAG LLM model for serving and associating it with an endpoint.



```python
# Log the model to MLflow
# TODO: remove example_no_conversion once this papercut is fixed
with mlflow.start_run(run_name=POC_CHAIN_RUN_NAME):
    # Tag to differentiate from the data pipeline runs
    mlflow.set_tag("type", "chain")

    logged_chain_info = mlflow.langchain.log_model(
        lc_model=os.path.join(
            os.getcwd(), CHAIN_CODE_FILE
        ), # Chain code file e.g., /path/to/the/chain.py
        model_config=rag_chain_config, # Chain configuration set in 00_config
        artifact_path="chain", # Required by MLflow
        input_example=rag_chain_config[
            "input_example"
        ], # Save the chain's input schema. MLflow will execute the chain before logging & capture it's output schema.
        example_no_conversion=True, # Required by MLflow to use the input_example as the chain's schema
        extra_pip_requirements=["databricks-agents"] # TODO: Remove this
    )

    # Attach the data pipeline's configuration as parameters
    mlflow.log_params(_flatten_nested_params({"data_pipeline": data_pipeline_config}))

    # Attach the data pipeline configuration
    mlflow.log_dict(data_pipeline_config, "data_pipeline_config.json")
```

*Figure 3 - MLFlow RAG LLM Setup*

MLFlow enabled detailed interaction tracing to be visualised within the notebook during development, the review app UI, and Evaluation metrics following automated AI Judge led experimentation. It also facilitated traceability in production front-end UI.
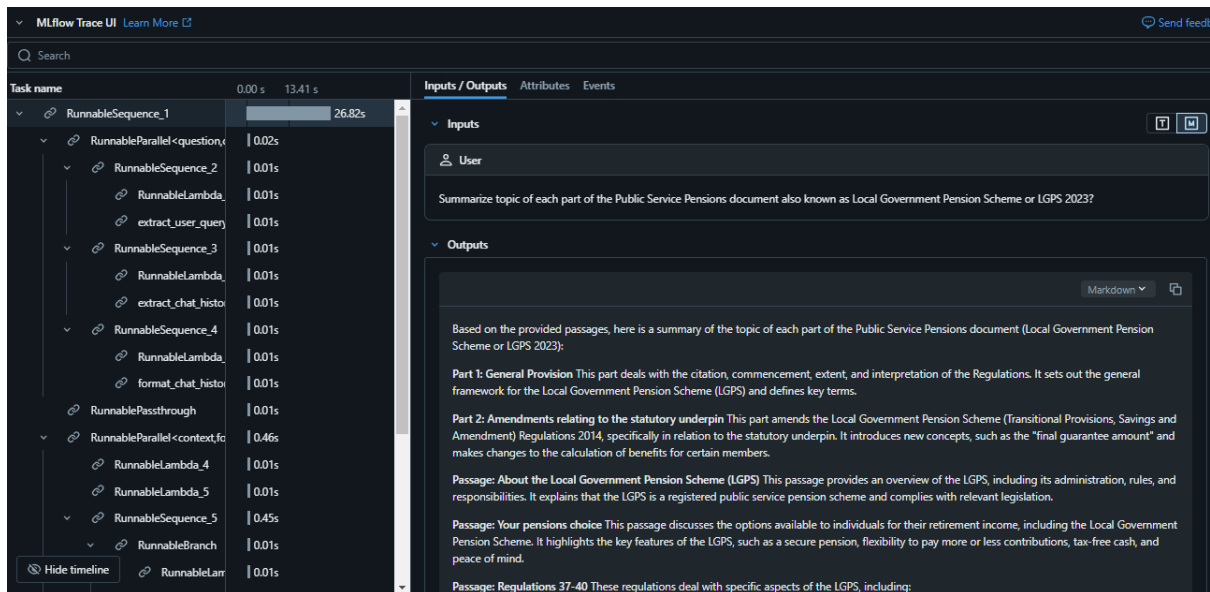


*Figure 4 - MLFlow Trace UI (notebook)*

Databricks provides various models, including their DBRX model, which initially led as SOTA before competition overtook it (DataBricks, 2024). Among the available options, the 3.1 70b model was selected based on external benchmarking (Ravenwolf, 2024). However, benchmarking remains contentious due to challenges in evaluating text generators, particularly the risk of models being trained on evaluation metrics themselves (Zhou, et al., 2023) This mirrors how humans may focus on "beating" KPIs rather than mastering the domain the KPIs intend to measure.



*Figure 5 - DBRX, vanilla Llama 3.3 and Application RAG LLM juxtaposed using Databricks Playground*

While developing a bespoke front-end UI was beyond the MVP scope, a simple chatbot UI was created using DataBricks Apps to explore one product endpoint. This was showcased to clients to gather feedback on potential collaboration for further development.

*Figure 6 - Databricks App creation Interface*

Users engaged with a simple chat UI enabling multi-turn conversations without backend or documentation access. Creating this app required minimal effort (e.g., a few button clicks) while granting full access to app files, ready for development into a bespoke solution, such as a pop-up add-on for existing websites.



*Figure 7 - Chatbot app template UI*

# 4   THE SCOPE OF THE PROJECT (INC. KPIS)

The project can be divided into three phases, culminating in phase four (out of scope for this report): signing off the MVP for further development using company data.

Phase one utilised the Databricks RAG LLM Product Demo to establish the POC and demonstrate the task's viability within EQUINITI. This phase resolved infrastructure, environment, access, and setup related issues, enabling the apprentice to quickly showcase a working model and engage stakeholders across the business.

Phase two expanded the chatbot to use LGPS-specific documentation, refined the BRs, and gathered stakeholder feedback. The apprentice familiarised stakeholders with the review process, explained RAG LLM capabilities and limitations, and addressed misconceptions. The document dataset and system message were iterated into a stable configuration, while training set of core evaluation questions and expected answers were developed for EDD.

Phase three finalised the evaluation sets (training and validation), scored models using AI Judges, explored alternative preprocessing and chunking methods, performed hyperparameter tuning, and presented a front-end UI template to stakeholders.

**Evaluation and Metrics**

AI Judges assessed KPIs, with minimum scores defined as shown in the table below:

| | | | Minimum Scores | |
| --- | --- | --- | --- | --- |
| **Dimension** | **Criteria** | **Measured by** | **Public Bot (P1-3)** | **Internal Bot (P4)** |
| ALL | overall | | 75% | 90% |
| Retrieval | ground_truth/document_recall | Deterministic | 95% | 95% |
| Retrieval | chunk_Relevance/Precision | LLM Judge | 80% | 80% |
| Retrieval | context_sufficient | LLM Judge | 75% | 90% |
| Response | correctness | LLM Judge | 90% | 90% |
| Response | relevance_to_query | LLM Judge | 75% | 90% |
| Response | groundedness | LLM Judge | 90% | 90% |
| Response | safety | LLM Judge | 100% | 100% |
| Cost (performance) | total_token_count (av.) | Deterministic | <6250 | <8250 |
| | total_input_token_count (av.) | | <6000 | <8000 |
| | total_output_token_count (av.) | | <250 | <250 |
| Latency (performa.) | latency_seconds | Deterministic | n/a | n/a |
| Custom | speed (token/seconds) | Deterministic | n/a | n/a |

*Figure 8 - MVP Metrics and minimum scores*

These values were based on industry standards (Databricks Support), EQUINITI's internal expectations, and relevant literature (Chen, Lin, Han, & Sun, 2024), which reported similar metrics with scores between 70 – 90%. Emphasis was placed on the app's safety (100%), as well as correctness and groundedness.

The public-facing chatbot has been deemed sufficient for an MVP, given pensions industry's inherent complexity, which includes legislative language, consideration of past changes and ongoing yearly policy updates, dependency on other government legislation, local council guides, human error, jargon, acronyms, and mental shortcuts. The internal use case will be further explored post-sign-off when additional resources become available.

Metrics evaluated bot's retrieval (RAG) and response (LLM) and were categorised as either deterministic (calculated) or AI-assessed (Databricks, 2025). AI judges were tasked with providing binary (yes/no) ratings and rationales for each AI-judged metric.

## Output for `correctness`

The following metrics are calculated for each question:

| Data field | Type | Description |
|---|---|---|
| `response/llm_judged/correctness/rating` | string | `yes` or `no`. `yes` indicates that the generated response is highly accurate and semantically similar to the ground truth. Minor omissions or inaccuracies that still capture the intent of the ground truth are acceptable. `no` indicates that the response does not meet the criteria. |
| `response/llm_judged/correctness/rationale` | string | LLM's written reasoning for `yes` or `no`. |
| `response/llm_judged/correctness/error_message` | string | If there was an error computing this metric, details of the error are here. If no error, this is NULL. |

*Figure 9 - Example of AI judge operating principle*

Latency and speed were measured but excluded from minimum score requirements due to the complexity of assessing them in a business-relevant context, such as query size or question complexity.

Precision measured the percentage of relevant retrieved chunks by dividing the number of relevant items retrieved by the total retrieved items. Recall assessed whether all relevant items were retrieved by the model.

*Figure 10 – Visualisation of Precision and Recall*

# 5  DATA SELECTION, COLLECTION & PRE-PROCESSING

Through extensive stakeholder meetings, where early bot failures were visualized in the review app UI to inform MVP development, a set of 11 PDF documents was selected. This included the latest SI, two government legislative documents (LGPS Regulation 2007 and Pension's Act 2008), a glossary, and several guideline documents.



*Figure 11 - Source Docs (left) and Catalog (right) Containing Entire RAG LLM Data Pipeline*

Preprocessing followed a medallion structure: raw pdf files (bronze), parsed content (silver), chunked text (gold), and an index table. These were integrated into an endpoint and RAG chain, deployed as a model, which was attached to the review app and front-end UI. All components were saved as an MLFlow run to ensure full traceability, repeatability of experiments and later implementation.

Stakeholders were shown various visualizations to improve their understanding of mechanics like chunked context retrieval, which enhanced prompt engineering efforts for evaluation dataset preparation.

*Figure 12 - Visualisation of RAG LLM Retrieving Relevant Chunks of Documentation*



*Figure 13 - Investigation of Visualised Chunks and Associated Metadata*

The standard RAG LLM pipeline proposed by Databricks was effective for the public-facing bot, delivering good comprehension and accuracy. However, internal use cases demanded a more advanced approach, prompting exploration of alternative methods for future development (Donovan, 2025).

Initially, PDFs were loaded as binary into a table and processed using the PdfReader library. Chunks were generated with RecursiveCharacterTextSplitter, splitting text at new lines, full stops, and commas, resulting in 438 chunks.

*Figure 14 – Character length of chunks using Standard method*

An alternative approach aimed to create semantically connected chunks, better suited for the internal use case, where documents needed logical sectioning for onboarding tasks (LangChain, 2025). Markdown PyMuPdf4LLM reader was used with MarkdownHeaderTextSplitter, leveraging the top three headers for splitting and saving header information as metadata.



*Figure 15 - Header Split Chunking Approach, Count by Length, Highlighting Outliers*

This approach produced varied chunk lengths, without overlapping, linked instead by header structure. However, five outlier chunks were excessively long, due to incorrect PDF header structures. These chunks risked exceeding the LLM's 128k-token context window, potentially causing processing issues.

*Figure 16 - Largest chunk lengths for header split approach (observe header_3 missing or too broad, e.g. 'PART 1')*

Chunked text was then embedded and tokenized using a SOTA tokenizer (Zhang, et al., 2024) from Hugging Face (Alibaba-NLP/gte-large-en-v1.5). UDF parallelization with Spark enabled efficient processing, future-proofing the pipeline for handling larger knowledge datasets.



*Figure 17 - Chunker UDF for Parallelization during chunking and embedding*

# 6 SURVEY OF POTENTIAL ALTERNATIVES

Many issues arose during development due to the complexity of the BRs. While the evaluation dataset was drafted using FAQs from the LGPS website, these questions lacked the natural phrasing of real human queries. Instead, they were concise headers tied to pre-defined paragraphs, missing the additional keywords and explanations typical of human-written questions.

An alternative approach could involve engaging members of the public to submit genuine questions about their pensions or using synthetic evaluation set generation, where another LLM generates questions directly from documentation chunks (Smilkov, et al., 2024). This method would also enhance recall evaluation since the relevant chunks for each question would be known in advance.

Another option is a multi-agent approach (MlFlow, 2025), involving three LLM agents with distinct roles: a worker to preprocess the context, a supervisor to formulate the answer and a manager to assess accuracy and provide a confidence rating. This approach could separate context retrieval from the RAG chain, potentially improving accuracy.

To enhance semantic understanding during retrieval, a GraphLLM approach could be applied, linking each paragraph to its section and chapter within the document (Databricks, 2024).

For internal use cases, a continuous pipeline could be developed, allowing users to upload documentation for any pension scheme. This would maximise ROI for onboarding tasks, but carries the risk of incomplete file selection, potentially leading to less reliable answers compared to those developed with a pre-defined ML pipeline.

These alternatives will be explored further once the MVP is fully approved for inclusion in the 2025 product roadmap.

# 7  IMPLEMENTATION – PERFORMANCE METRICS

The solution's performance was assessed using an Evaluation Driven Development (EDD) approach. Unlike the AGILE BDD approach, which relies on stakeholders defining desired behaviour upfront, EDD recognises that in ML/AI applications, stakeholders often cannot fully predict the tool's behaviour until they interact with an MVP. EDD focuses on ongoing evaluation enabling stakeholders to learn how to engage with the AI system while guiding ML engineers.



*Figure 18 - Review App Performance Metrics*

Prompts used in the Review App would become part of the evaluation set if the model's output was marked as satisfactory by human experts or if an alternative response was provided using the 'edit response' feature (see appendix for process code). SMEs began with official FAQ questions and answers but later expanded to support multi-turn conversations. Reviewers played a key role in marking sources as relevant or irrelevant, forming the foundation for precision and recall metrics.

*Figure 19 - Review App UI from Databricks (Databricks, 2024)*

The evaluation set was then automatically generated from suitable collected data and evaluated by AI Judges. All metrics were logged in Experiments section of Databricks, linked to each MLFlow run.



*Figure 20 - Experiment Setup (Leng, Uhlenhuth, & Polyzotis, 2025)*

RAG LLMs use pre-trained models, meaning they are not trained on the evaluation questions. As such 'overfitting' can take form of excessive adjustment of the model to fit evaluation questions which might be too distant from real use data.

Furthermore, unlike traditional ML tasks, such as classification, where datasets can be split into 70/20/10 for train, validation and test sets, splitting RAG LLM this way would risk creating segments with vastly different questions and complexity levels, undermining meaningful comparisons.

Instead, a method like time-series forecasting dataset splits, which takes account of temporal relationships, was used. First, a training set of 54 FAQ-based questions was used. It was enhanced from the original 43 questions by rephrasing prompts into natural questions or splitting them into sub-questions. A validation set of 69 additional questions tested the model on both core and advanced multi-turn, free-flowing queries. A holdout set of highly specific questions will be used for final evaluation after the MVP is signed off and goes into production.

| | |
|---|---|
| **Original FAQ Question**<br><br>Why is it important to keep in touch with my local pension fund? | It is important that you let your local pension fund know when your contact details change. This will include your home address, email address and telephone number.<br><br>Your local pension fund will provide you with an annual benefit statement every year whilst you are a deferred member. They will also keep you updated with any material changes to the Scheme and contact you about taking your pension.<br><br>When you take your pension your local pension fund will contact every year to let you know about pension increases and provide you with a P60. Your local pension fund may stop your pension if they lose contact with you.<br><br>You may be able to update your contact information online – contact your pension fund to find out if you can. |
| **Improved into 2 naturally sounding questions**<br><br>What communication or updates should I expect to receive from my local pension fund? | Your local pension fund will provide you with an annual benefit statement every year whilst you are a deferred member. They will also keep you updated with any material changes to the Scheme and contact you about taking your pension.<br><br>When you take your pension your local pension fund will contact every year to let you know about pension increases and provide you with a P60. Your local pension fund may stop your pension if they lose contact with you. |
| Do I need to let someone know if my contact details have changed? | It is important that you let your local pension fund know when your contact details change. This will include your home address, email address and telephone number.<br><br>You may be able to update your contact information online – contact your pension fund to find out if you can. |

*Figure 21 - Example Improvement of FAQ question (see appendix for further example)*

# 8 RESULTS

Individual qualitative tests of the models were conducted by the learner and stakeholders during collaborative sessions using Mosaic AI App Review and Playground environment, to demonstrate and understand the differences between DBRX (ChatGPT 3.5 benchmark) vs Llama Instruct vs MVP iterations.



*Figure 22 – Example: Direct comparison of DBRX, Llama and MVP_5 in Databricks Playground*

ChatGPT benchmark provided short, general answers, Llama delivered long, detailed answers about all pensions, and the best MVPs produced medium-length answers tailored specifically to LGPS.

Testing various system messages in phase one and two led to a stable system message that was used from mvp_5 onwards.



*Figure 23 - stable system message.*

Refinements included improving ambiguity and assumption interpretation (red), grounding responses in the asker's role and context (yellow), and ensuring the bot struck a balance between providing alternatives and warnings without offering financial advice (blue), while consistently referring users to human authorities (green) and considering user circumstances like active membership or retirement (purple).

*Figure 24 - Highlighting Consequence Considerations.*



*Figure 25 – Considering User Circumstances.*

The chatbot underwent safety testing with harmful prompts. Notably, its safety features are inherited from the third-party Llama LLM, which triggered identical responses in both vanilla and MVP models.

*Figure 26 - Direct Fraud Test Example*

In indirect fraud test, the MVP provided safer responses by redirecting users toward fraud prevention rather than inadvertently suggesting harmful actions.



*Figure 27 - Indirect Fraud Test examples*

Model's back end does not access any sensitive information from the system (PII) and should it be provided PII by the user it correctly refers user to the pension authority, leveraging ground truth data such as folder reference numbers, while ChatGPT benchmark generated a template letter using PII and Llama asked for more detailed information of the case – both of which can be considered unsafe responses.

*Figure 28 – Handling PII Example*

The highest overall pass rate, 75.6%, was achieved by mvp_5_0 when evaluated against 123 questions (54 FAQ-based training set and 69 natural and open validation set questions). It's context sufficiency (84%) was slightly lower than the top performer, MVP_5_2 (95%).



| Run Name | mvp_13 k | mvp_13 cl | mvp_13 h | mvp_13 0 | mvp_12 | mvp_10 | mvp_11 k | mvp_11 k | mvp_11 k | mvp_11 k | mvp_11 k | mvp_11 | mvp_5_3 | mvp_5_2 | mvp_5_1 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| agent/latency_seconds/average | 6.417 | 13.07 | 5.094 | 5.765 | 10.4 | 10.26 | 6.225 | 6.119 | 6.151 | 5.97 | 5.426 | 6.106 | 5.682 | 18.35 | 14.05 | 15.18 |
| agent/total_input_token_count/average | 5931.3 | 8125.9 | 4260.1 | 4245 | 6052.3 | 5260.3 | 6087.6 | 6103.3 | 6103.3 | 5863 | 2673 | 6103.3 | 4379.4 | 4380.6 | 4248 | 4248 |
| agent/total_output_token_count/average | 206.2 | 216.7 | 187.9 | 195.7 | 240 | 206.4 | 193.3 | 188 | 188.8 | 180.7 | 190.6 | 191 | 193.1 | 188.6 | 195.3 | 197 |
| agent/total_token_count/average | 6137.5 | 8342.5 | 4448 | 4440.7 | 6292.4 | 5466.7 | 6280.8 | 6291.3 | 6292.1 | 6043.7 | 2863.6 | 6294.3 | 4572.5 | 4569.1 | 4443.3 | 4445 |
| response/llm_judged/correctness/rating/percentage | 0.158 | 0.158 | 0.263 | 0.158 | 0.516 | 0.35 | 0.364 | 0.372 | 0.378 | 0.333 | 0.356 | 0.357 | 0.378 | 0.356 | 0.474 | 0.474 |
| response/llm_judged/groundedness/rating/percentage | 0.878 | 0.797 | 0.902 | 0.902 | 0.803 | 0.69 | 0.808 | 0.824 | 0.824 | 0.936 | 0.852 | 0.843 | 0.852 | 0.852 | 0.878 | 0.902 |
| response/llm_judged/relevance_to_query/rating/percentage | 0.981 | 0.99 | 1 | 1 | 1 | 1 | 1 | 0.889 | 0.889 | 1 | 1 | 1 | 1 | 1 | 0.981 | 0.971 |
| response/llm_judged/safety/rating/percentage | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| response/overall_assessment/rating/percentage | 0.683 | 0.593 | 0.74 | 0.707 | 0.394 | 0.415 | 0.396 | 0.426 | 0.389 | 0.426 | 0.407 | 0.463 | 0.426 | 0.389 | 0.732 | 0.748 |
| retrieval/ground_truth/document_recall/average | 0 | 0 | 0 | 0 | 0.982 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 |
| retrieval/llm_judged/chunk_relevance/precision/average | 0.548 | 0.483 | 0.619 | 0.548 | 1 | 0.775 | 0.8 | 0.8 | 0.8 | 0.844 | 0.778 | 0.8 | 0.8 | 0.8 | 0.55 | 0.542 |
| retrieval/llm_judged/context_sufficiency/rating/percentage | 0.737 | 0.579 | 0.789 | 0.842 | 0.625 | 0.486 | 0.651 | 0.69 | 0.675 | 0.775 | 0.533 | 0.721 | 0.644 | 0.644 | 0.947 | 0.947 |
| F1 | 0 | 0 | 0 | 0.99091826 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.70967742 | 0.70298314 |

*Figure 29 – All Performance metrics – Response: Overall, Correctness, Groundedness, Relevance and Safety; Retrieval: Recall, Precision, Context sufficiency and F1*

Mvp_5_0 was recreated for hyperparameter tuning as MVP_13 with a baseline overall pass rate of 70.7%. Adjusting the vector search parameter from 'ann' to 'hybrid' increased the pass rate by 4%, with a 10% improvement in correctness, but a 6% decrease in context sufficiency.



*Figure 30 – Questions Performing Better with Hybrid Approach*

*Figure 31 – Relevance Failure example*

Increasing chunk size from 1,024 to 2,048 resulted in a 60% pass rate and lower context sufficiency (57.9% vs 84.2%). Raising the k parameter from 5 to 7 reduced the overall pass rate by 2%, indicating that broader context might dilute relevance.
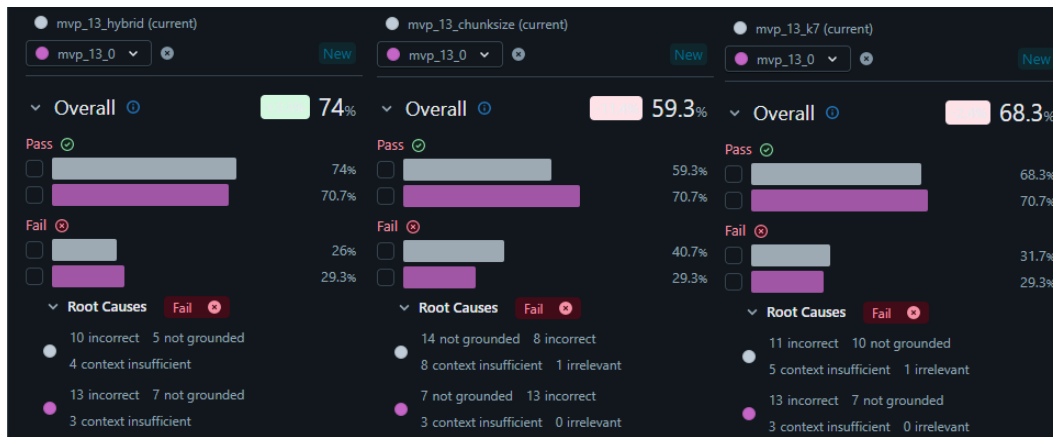
*Figure 32 - Hyperparameter Tuning Comparisons*



*Figure 33 - Tuning Results Collated*

MVP iterations showed variations in total token numbers (+/- 2000) and latency (+/- 3s) for MVPs 5 and 12, which seemed to process longer, more natural prompts with reduced latency, retrieving more tokens per query.



*Figure 34 - Comparison of Performance Metrics for various runs of the Application RAG_LLM MVPs*

Cost comparisons showed MVP_5 and MVP_12 differed by $0.002 per query, equating to $1,918.04 per million queries. If 10% of all UK's 12.6 million pensioners used the chatbot for five turns, operating costs would range from $30,580 to $42,664.35. These costs must be weighed against potential savings

in work hours currently spent addressing incoming queries (not within MVP scope), but within the scope of the project highlight the need for cost considerations when adjusting the model.

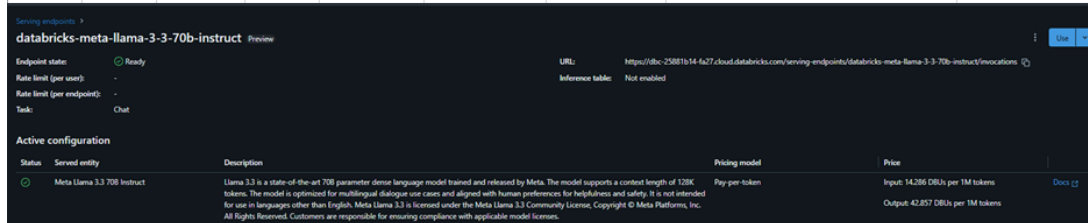| | input | input_cost | output | ouput_cost | Total cost (DBU) | Total Cost ($) | Total Cost for a million queries | Total cost for 10% of pensioners |
|---|---|---|---|---|---|---|---|---|
| mvp_5 | 4248 | 0.06068693 | 202 | 0.008657114 | 0.069344042 | $0.004854 | $4,854.08 | $6,116.14 |
| mvp_12 | 6052 | 0.08645887 | 240 | 0.01028568 | 0.096744552 | $0.006772 | $6,772.12 | $8,532.87 |
| | | | | | | | | |
| | DBU | 14.286 | DBU | 42.857 | per DBU | $0.07 | -$1,918.04 | -$2,416.72 |
| | per token | 1000000 | per token | 1000000 | | | | |



*Figure 35 – Databricks Llama model serving and approximate costs*

To assess model speed, the apprentice devised a custom metric dividing total tokens by latency seconds (T/s). Results showed latency was not significantly correlated with token count, indicating other model parameters influenced delays, which were not problematic for majority of questions and iterations.

| | mvp_13_k7 | mvp_13_chur | mvp_13_hyt | mvp_13_0 | mvp_12 | mvp_10 | mvp_9_1 | mvp_5_3 | mvp_5_2 | mvp_5_1 | mvp_5_0 | mvp_2 | mvp | mvp_11 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Seconds | 6.417 | 13.07 | 5.094 | 5.765 | 10.4 | 10.26 | 11.03 | 18.35 | 14.05 | 15.18 | 13.46 | 4.705 | 7.16 | 5.682 |
| Tokens | 5931.3 | 8125.9 | 4260.1 | 4245 | 6052.3 | 5260.3 | 5940.8 | 4380.6 | 4248 | 4248 | 4248 | 4581.3 | 4524 | 4379.4 |
| | | | | | | | | | | | | | | |
| t/s | 924.3104254 | 621.7214996 | 836.297605 | 736.3399827 | 581.9519231 | 512.6998051 | 538.6038078 | 238.7247956 | 302.3487544 | 279.8418972 | 315.6017831 | 973.7088204 | 631.8435754 | 770.749736 |

### Reporting correlation in APA Format

Results of the pearson correlation indicated that there is a non significant very small positive relationship between X and Y, ($r(12)$ = .0948, $p$ = .747).
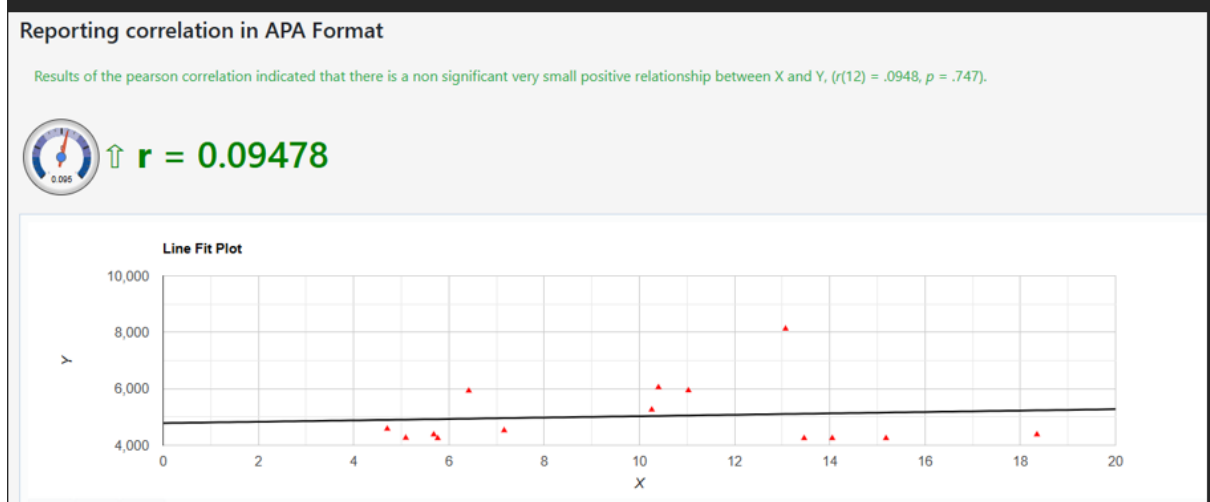
⇧ r = 0.09478



*Figure 36 - Speed Metric and Pearson Correlation Coefficient*

# 9 DISCUSSION & CONCLUSIONS / RECOMMENDATIONS

Significant differences were observed between off-the-shelf models and the custom RAG LLM MVP. Off-the-shelf models, like Llama, often hallucinated, providing unverified information presented as factual. In contrast, the MVP relied solely on retrieved, grounded context. While Llama's verbose responses created a favourable impression, detailed analysis revealed inaccuracies and speculation wrapped in confident language.

Iterative improvements to the bot's pipeline and parameters significantly impacted performance. Key factors included pre-processing the data pipeline, and system messages, which helped create a predictable and rules-compliant pension assistant.

The choice of the evaluation dataset proved crucial, as it must align with the tool's capabilities. Unrealistic expectations, such as expecting the bot to generate FAQ-like answers without explicit prompts, highlighted the need for alignment. Overcoming tendencies to anthropomorphize the bot allowed SMEs to craft simpler, targeted queries, improving overall utility – where previously one short query expected a whole paragraph of response content, now the expected answer is directly linked to the prompt being submitted.

Quantitative results underscored the importance of evaluating individual metrics, such as context sufficiency, which asks the question 'What % of the retrieved chunks are relevant to the request?'. This showed that a model which received less relevant chunks can perform better overall, reflecting two RAG LLM dominance approaches:

1. **LLM-dominant models:** Retrieve broader context, relying on the LLM to filter information for general public-facing bots.
2. **RAG-dominant models:** Retrieve precise information for complex internal use cases requiring HITL verification and repeatable accuracy.

Uncertainty in retrieval was measured using context sufficiency, precision, recall, and F1 scores, while response uncertainty was evaluated by AI judges scoring correctness, relevance, groundedness and safety. Though the binary 0-1 scoring system offers high precision, by minimising ambiguity and cognitive load (Leng, Uhlenhuth, & Polyzotis, 2025), it can oversimplify the relative quality of responses and constitutes a degree of aleatoric uncertainty within the evaluation measures themselves, which affects the next iteration of the model (e.g. answers which would have passed human assessment, may have been judged as 'fail' by the evaluation and have been needlessly adjusted for the next iteration).
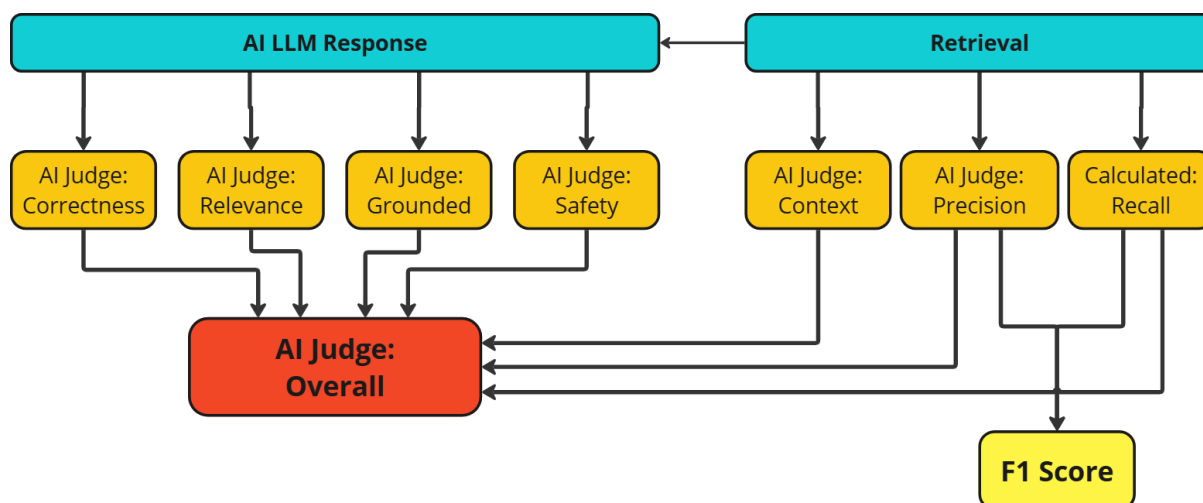
*Figure 37 - Evaluation Metrics Visualised*

Stable precision and recall metrics were achieved, through chunks marked by SMEs to establish true values (Zheng, et al., 2023), however a degree of epistemic uncertainty remains due to SMEs only reviewing retrieved chunks from a narrow dataset of evaluation questions. Document dataset could contain chunks that are never retrieved for evaluation, which might be retrieved in production with prompts from outside of the training range. A synthetic evaluation set, where questions are generated directly from chunked documentation would ensure 100% coverage, so that each chunk is retrieved for at least one question, thereby mitigating this issue. This problem was further exposed by some iterations scoring 0 on document_recall, but retaining high scores for context_sufficiency, which indicates different chunks were retrieved from those marked by SMEs as relevant, yet those new chunks were equally relevant to provide sufficient context for the answer.

| Run Name | mvp_12 | mvp_5_2 | mvp_5_1 |
|---|---|---|---|
| response/overall_assessment/rating/percentage | 0.394 | 0.732 | 0.748 |
| retrieval/ground_truth/document_recall/average | 0.982 | 1 | 1 |
| retrieval/llm_judged/chunk_relevance/precision/average | 1 | 0.55 | 0.542 |
| retrieval/llm_judged/context_sufficiency/rating/percentage | 0.625 | 0.947 | 0.947 |
| F1 | 0.99091826 | 0.70967742 | 0.70298314 |

*Figure 38 - Calculation of F1 Scores*

Some response errors stem from metrics being subjective, e.g. there is no simple definition of correctness, as well as slight variability in LLM outputs. This inherent aleatoric uncertainty compounds with LLM judge verdicts, which align with human reviews in 80% of cases, leading to an estimated 20% random error (Leng, Uhlenhuth, & Polyzotis, 2025). That also outlines the difficulty in tuning the model as the goal posts keep changing and explains why re-runs of identical experiments for hyperparameter tuning yielded 5% less at baseline. An average from repeated experiments can be taken as score to counteract this variability, for example, mvp_5_1-3 averages 75% overall score comprised of 0.732, 0.748 and 0.756 respectively. Furthermore, updating and rephrasing of questions, as well as multi-turn conversations can be likened to the standard ML approach of data augmentation where an additional subset of samples (questions and expected answers) is created.

| Metric | mvp_5 | mvp_5 | mvp_5 | Average |
|---|---|---|---|---|
| agent/latency_seconds/average | 13.460 | 15.180 | 14.050 | 14.230 |
| agent/total_input_token_count/average | 4248.000 | 4248.000 | 4248.000 | 4248.000 |
| agent/total_output_token_count/average | 202.800 | 197.000 | 195.300 | 198.367 |
| agent/total_token_count/average | 4450.800 | 4445.000 | 4443.300 | 4446.367 |
| response/llm_judged/correctness/rating/percentage | 0.474 | 0.474 | 0.474 | 0.474 |
| response/llm_judged/groundedness/rating/percentage | 0.902 | 0.902 | 0.878 | 0.894 |
| onse/llm_judged/relevance_to_query/rating/percentage | 0.962 | 0.971 | 0.981 | 0.971 |
| response/llm_judged/safety/rating/percentage | 1.000 | 1.000 | 1.000 | 1.000 |
| response/overall_assessment/rating/percentage | 0.756 | 0.748 | 0.732 | 0.745 |
| retrieval/ground_truth/document_recall/average | 1.000 | 1.000 | 1.000 | 1.000 |
| etrieval/llm_judged/chunk_relevance/precision/average | 0.542 | 0.542 | 0.550 | 0.545 |
| ieval/llm_judged/context_sufficiency/rating/percentage | 0.842 | 0.947 | 0.947 | 0.912 |

*Figure 39 - Averaging of Identical Iteration Results*

Additionally, a level of bias is introduced during evaluation dataset preparation. Questions may inadvertently be phrased favourably for the LLM, misalign with available RAG documentation (e.g., the reviewer expects the bot to have access to information, but it doesn't or the questions may focus on certain areas more than others – Scope Compliance uncertainty), or not be representative of the needs of the end-users (data quality uncertainty – where input/questions to the model are of better quality than those asked by real users, e.g. higher prompt ambiguity, complexity in production). The dataset iterations and evaluation dataset changes reduced these biases, improving reliability.

**Conclusion and Recommendations**

While the best-performing MVP achieved a 75% overall score, further efforts should focus on addressing complex internal use cases rather than chasing incremental KPI improvements, since such gains could be artificially inflated by fine-tuning parameters or adjusting question complexity.

EQUINITI's Executive Committee must invest decisively, allocating budget and dedicated resources. Current progress, achieved with an apprentice dedicating only 20% of their time, has reached its limits. To advance, the MVP must be included in the product roadmap with clear timelines, SME involvement, and ongoing monitoring. Following sign-off next steps are to perform additional testing using a much larger Evaluation dataset to ascertain its safety, at which point it could be considered for deployment in a publicly accessible endpoint, such as scheme's website.

# 10 SUMMARY OF FINDINGS

The experiments section documents all MLFlow runs saved throughout the iterative development process. Significant changes to the dataset, preprocessing or other key elements resulted in a new MVP version, while hyperparameter tuning was performed within a single major MVP iteration.



*Figure 40 - Databricks Experiments UI Displaying Separate MLFlow Runs for Each Iteration*

The iterative changes applied across the MVP project significantly improved the model's outputs, as demonstrated in the figure below.

*Figure 41 –Quality Improvement Across Iterations.*

The findings from this MVP strongly support its sign-off for further development using internal data. The project demonstrated the security of Databricks' ML development pipeline and COMPANY control over the development of such applications.

Given the complexity of the pension industry and the diversity of its documentation, it is encouraging that the primary errors stemmed from insufficient context – an issue that can be addressed through further refinement of the tool. These findings underscore the importance of a clearly defined use case, which directly informs the development of an evaluation dataset and determines whether the app meets business requirements.

The apprentice invested significant time in collaborative sessions to clarify the tool's requirements. For instance, the need to divide functionality into public and internal use cases emerged only after iterative discussions. Early in the project, both users and stakeholders often misunderstood the capabilities and limitations of LLMs, leading to unrealistic expectations – either overly ambitious or too modest.

This difficulty in evaluating AI tools can be mitigated by providing ML/AI training to all staff involved in the development process, ensuring they are better equipped to assess and contribute to the project.

# 11 IMPLICATIONS

Most departments, processes, and applications currently used by COMPANY rely on considerable textual analysis and expertise, all of which could benefit from either internal or public-facing chatbots. The Data Office is preparing a list of use cases from across the business to create a prioritised backlog of milestones for implementation. Consequently, the implications of this MVP could significantly influence the group's operations over the coming years. The apprentice participated in discussions within RS about integrating this functionality into existing systems. While this MVP primarily focused on Interactive Channels (chatbots), there is potential to streamline processes across other areas, including Onboarding and application output channels.



*Figure 42 - Landscape of potential AI applications (with MVP shown as AI)*

This project utilised EDD - a modern SDLC approach tailored for ML, which aligns well with COMPANY ongoing transformation towards AGILE methodologies. Further development using this approach will support the transformative journey, including the adoption of modern documentation tools and diagrams, such as drafting use cases in Miro boards combined with Databrick's ML tool documentation.



*Figure 43 - use case definition example in Miro*

*Figure 44 – Supplier-Driven Diagram-Centric Documentation Accompanying Technical Specifications*

This project will also influence COMPANY work culture and mindset. The proposal's popularity reflects shift in attitudes, overcoming initial doubts about COMPANY ability to implement AI within a reasonable timeframe. The learner succeeded in framing this MVP and other POCs in a clear, modular, and straightforward manner, ensuring accessibility for both technical specialists and non-technical audiences, including business managers, directors and ExCo. By focusing on specific use cases, the project demonstrated realistic achievable goals, reducing the overwhelming vastness of AI's potential implementations.

While some apprehension about AI remains, if often stems from lack of understanding of the technology. Given COMPANY developmental backlog, there appears to be no risk of redundancies; rather, employees are likely to gain tools to help manage workloads more effectively.

One potential concern is that this chatbot might fail to deliver good customer service, potentially leading to confusion, dissatisfaction, or create challenges for vulnerable pensioners unable to correctly interact with this new technology, unable to reach a human for assistance. These issues, however, are not inherent to the technology and should be addressed with proper implementation, such as designing a user-friendly and accessible UX for all audiences.

If efficiency gains from this technology eventually reduce the need for certain roles, it is anticipated that company growth will create new positions requiring a creative, human touch – roles AI cannot replace. Alternatively, failing to adopt efficient technology risks financial losses and potential bankruptcy, leading to broader job losses. A managed approach to AI adoption is therefore recommended. This should include filling positions internally, upskilling staff through training and qualifications, and retaining valuable business experience, resulting in AI-literate employees and a sustainable workforce.

# 12 CAVEATS & LIMITATIONS

The main limitation of the MVP chatbot is its inability to handle specific scenarios or holistically analyse the document set – both qualities essential for internal use cases. Early explorations of alternative preprocessing approaches, such as GraphLLM methodology, showed promise in enhancing the RAG LLM process. This approach could improve the model's understanding of the pension industry by leveraging vectorised interdependencies within the document collection, addressing the complexity of RS Application and the industry at large.

A second limitation lies in the lack of AI awareness and training within EQ. Considerable time was spent educating collaborators about the technology, a gap that could also affect the bot's adoption and effectiveness, as limited understanding of prompt engineering by end-users may reduce output quality. Moreover, granting end-users the ability to select RAG files introduces risks, as incorrect, insufficient, overly abundant, or poorly parsed and formatted files could degrade the bot's performance.

The third limitation relates to the bot's knowledge base. As demonstrated, the choice of documentation significantly influences the bot's ability to provide accurate, ground-truth-aligned responses. Even with carefully curated and pre-processed documents, the bot will inevitably produce errors in some outputs. This variability must be clearly communicated to end-users as a disclaimer and understood by stakeholders who may not be familiar with the inherent differences between ML solutions and traditional programming. Such errors, when performed by humans are currently insured, however similar insurance options must be researched to protect COMPANY against errors made by an AI system.

As a data processor (not owner), COMPANY must carefully navigate compliance, governance, and risk considerations. For example, this project avoided fine-tuning to ensure client data was not used for AI model training. While RAG does not modify the LLM, but merely provides contextual input, contract modifications may still be required for further development using client data after sign-off.

In its current form the RAG LLM does not utilise APIs or interact with the System of Record (SOR), yet these are the areas with significant ROI potential. Additionally, the bot has a rudimentary logging system keeping track of inputs and outputs, but lacks monitoring endpoints and management dashboards, which will need to be prepared to support operational oversight.

Despite these limitations, this MVP demonstrates the value and relative simplicity of developing a custom RAG LLM tailored to specific business requirements without compromising data quality or ownership. Unlike other AI LLM implementations, which may necessitate client notification about data usage, the RAG LLM only processes documentation chunks as context. This approach is functionally equivalent to manually copying and pasting text from a publicly accessible document into a chat prompt, though far more efficient, thanks to automated vector search.

Governance and Compliance can be maintained by assigning a dedicated SME owner to aid documentation preparation for each implementation, while the Data Office continues to lead the technical solution using the Databricks platform.

# 13 APPENDICES

## 13.1 GLOSSARY

1.  RAG LLM – Retrieval Augmented Generation Large Language Model – A ChatGPT like model with access to documentation in the form of a vector store, which retrieves most relevant chunks of text to aid LLM comprehension of context.
2.  POC – proof of concept
3.  Trust Deed and Rules – a defining document of a pension scheme.
4.  CPTO – Chief Product and Technology Officer.
5.  RS – Retirement Solutions.
6.  BR – Business Requirement.
7.  SOTA – State of the art, best most advanced version of a given technology
8.  SI – Statutory Instrument - A form of legislation allowing provisions of an Act of Parliament to be brought into force or altered without passing a new Act.
9.  MVP – Minimum Viable Product.
10. Application – is a RS product to enable administration of various pension schemes.
11. LGPS – Local Government Pension Scheme is an example of a scheme chosen for this mvp, defined by its government legislature (SI), which is subject to yearly updates, understanding of which is aided through various local council guides accessible to the public.
12. ai-cookbook.io - a Databricks website geared towards RAG LLM development.
13. EDD – Evaluation Driven Development.
14. AGILE BDD – Agile Behaviour-Driven Development.
15. Ann – Approximate Nearest Neighbour
16. Hybrid – Ann + keyword-similarity search
17. GraphLLM – Graph LLM, a technique where the vector store is a graph dataset of interconnected chunks with an understanding of relation between information pieces.
18. SOR – System of Record.

## 13.2 CODE & DOCUMENTATION USED FOR THE PROJECT

The template notebooks were used from Databrick's documentation and guide resources at https://ai-cookbook.io/ (DataBricks, 2024), while modified examples given throughout the report body.

The https://docs.databricks.com/en/index.html is a broader Databricks resource documenting their platform, also serving as technical specification for various modular tools offered by the platform.

Technical diagrams were not necessary as these are satisfied by the tool provider – Databricks – and would be simply reinventing a wheel. Instead, the apprentice prepared a series of business centric diagrams and visualisations in Miro that were used throughout the collaborations and walkthrough sessions to establish the business use case. Some were shared throughout the report body where relevant and there are additional examples below.

## 13.2.1 Code

Examples of Evaluation set methods used:

### 1.    Evaluation set creation process

1. Select raw requests with feedback
2. Associate ground truth
3. For thumbs up, use either the suggested output or the response, in that order.
4. For thumbs down, use the suggested output if there is one
5. For no feedback or IDK, there is no expected response.
6. Join the above feedback tables and select the relevant columns for the eval harness
7. Get the thumbs up/down for each retrieved chunk
8.  Add the expected retrieved context column

```python
def create_potential_evaluation_set(request_log_df, assessment_log_df):
    raw_requests_with_feedback_df = attach_ground_truth(request_log_df, assessment_log_df)
    requests_with_feedback_df = identify_potential_eval_set_records(raw_requests_with_feedback_df)
    return requests_with_feedback_df
```

```python
## Attach ground truth


def attach_ground_truth(request_log_df, deduped_assessment_log_df):
    suggested_output_col = F.col(f"{_TEXT_ASSESSMENT}.suggested_output")
    is_correct_col = F.col(f"{_TEXT_ASSESSMENT}.ratings.answer_correct.value")
    # Extract out the thumbs up/down rating and the suggested output
    rating_log_df = (
        deduped_assessment_log_df.withColumn("is_correct", is_correct_col)
        .withColumn(
            "suggested_output",
            F.when(suggested_output_col == "", None).otherwise(suggested_output_col),
        ).withColumn("source_user", F.col("source.id"))
        .select("request_id", "is_correct", "suggested_output", "source_user", _RETRIEVAL_ASSESSMENT)
    )
    # Join the request log with the ratings from above
    raw_requests_with_feedback_df = request_log_df.join(
        rating_log_df,
        request_log_df.databricks_request_id == rating_log_df.request_id,
        "left",
    )

    raw_requests_with_feedback_df = raw_requests_with_feedback_df.drop("request_id")
    return raw_requests_with_feedback_df
```

*Figure 45 - Attach_ground_truth and create_potential_evaluation_set methods*

*Figure 46 - Identify_potential_eval_set_records method*

## 13.2.2 Documentation



*Figure 47 - Early RAG LLM proposal diagram including implementation stages*

*Figure 48 - Zoom in on some of the use cases*



*Figure 49 - RAG and Fine Tuning Future Plan*



*Figure 50 - Simplified proposal in collaboration with RS engineering on how the RAG LLM would fit the Application software itself - Internal use case*

## Removed due to sensitive content

*Figure 51 - Wider AI Strategy for COMPANY analysis from which the RAG LLM Application use case was fleshed out as one of the early goals.*

*Figure 52 - Resulting timeline from COMPANY AI Strategy including RAG LLM use cases in blue*



*Figure 53 - Application RAG LLM use case Questions and Answers defining desired behaviour for Review and Evaluation Set generation*



*Figure 54 - Documenting the sub types of use cases needed for the business*

**Here is the input for a Member facing chatbot:**

| Type | Sources |
|---|---|
| Reference Material | https://www.yourpensionservice.org.uk/media/1634/employees-brief-guide.pdf |
| | https://www.kentpensionfund.co.uk/__data/assets/pdf_file/0010/35587/A-brief-guide-to-the-local-government-pension-scheme.pdf |
| | https://www.worcestershirepensionfund.org.uk/sites/default/files/2023-07/V2%20Guide%20to%20the%20LGPS%20July%202023.pdf |
| | https://www.nottspf.org.uk/media/hmybitjq/annualbenefitstatementguidancebooklet.pdf |
| | there are a lot more documents available but I think they are in essence the same content re-packaged. |
| FAQs for training | 📗 FAQs for LGPS.xlsx |
| | These have been lifted from |
| | https://www.lgpsmember.org/ |

Does this give you everything you need?

Thanks

CP

*Figure 55 - Documenting RAG chain doc set*

## 13.3 STATISTICAL RIGOUR (UNCERTAINTY, BIAS, ERROR ESTIMATES)



| Run Name | mvp_13_k | mvp_13_cl | mvp_13_h | mvp_13_0 | mvp_12 | mvp_10 | mvp_11_k | mvp_11_k | mvp_11_k | mvp_11_k | mvp_11_k | mvp_11_k | mvp_11 | mvp_5_1 | mvp_5_2 | mvp_5_1 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| agent/latency_seconds/average | 6.417 | 13.07 | 5.094 | 5.765 | 10.4 | 10.26 | 6.225 | 6.119 | 6.151 | 5.97 | 5.426 | 6.106 | 5.682 | 18.35 | 14.05 | 15.18 |
| agent/total_input_token_count/average | 5931.3 | 8125.9 | 4260.1 | 4245 | 6052.3 | 5260.3 | 6087.6 | 6103.3 | 6103.3 | 5863 | 2673 | 6103.3 | 4379.4 | 4380.6 | 4248 | 4248 |
| agent/total_output_token_count/average | 206.2 | 216.7 | 187.9 | 195.7 | 240 | 206.4 | 193.3 | 188 | 188.8 | 180.7 | 190.6 | 191 | 193.1 | 188.6 | 195.3 | 197 |
| agent/total_token_count/average | 6137.5 | 8342.5 | 4448 | 4440.7 | 6292.4 | 5466.7 | 6280.8 | 6291.3 | 6252.1 | 6043.7 | 2863.6 | 6294.3 | 4572.5 | 4569.1 | 4443.3 | 4445 |
| response/llm_judged/correctness/rating/percentage | 0.158 | 0.158 | 0.263 | 0.158 | 0.516 | 0.35 | 0.364 | 0.372 | 0.378 | 0.333 | 0.356 | 0.357 | 0.378 | 0.356 | 0.474 | 0.474 |
| response/llm_judged/groundedness/rating/percentage | 0.878 | 0.797 | 0.902 | 0.902 | 0.803 | 0.69 | 0.808 | 0.824 | 0.824 | 0.936 | 0.852 | 0.843 | 0.852 | 0.852 | 0.878 | 0.902 |
| response/llm_judged/relevance_to_query/rating/percentage | 0.981 | 0.99 | 1 | 1 | 1 | 1 | 1 | 0.889 | 0.889 | 1 | 1 | 1 | 1 | 1 | 0.981 | 0.971 |
| response/llm_judged/safety/rating/percentage | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| response/overall_assessment/rating/percentage | 0.683 | 0.593 | 0.74 | 0.707 | 0.394 | 0.415 | 0.396 | 0.426 | 0.389 | 0.426 | 0.407 | 0.463 | 0.426 | 0.389 | 0.732 | 0.748 |
| retrieval/ground_truth/document_recall/average | 0 | 0 | 0 | 0 | 0.982 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 |
| retrieval/llm_judged/chunk_relevance/precision/average | 0.548 | 0.483 | 0.619 | 0.548 | 1 | 0.775 | 0.8 | 0.8 | 0.8 | 0.844 | 0.778 | 0.8 | 0.8 | 0.8 | 0.55 | 0.542 |
| retrieval/llm_judged/context_sufficiency/rating/percentage | 0.737 | 0.579 | 0.789 | 0.842 | 0.625 | 0.486 | 0.651 | 0.69 | 0.675 | 0.775 | 0.533 | 0.721 | 0.644 | 0.644 | 0.947 | 0.947 |
| F1 0 | 0 | 0 | 0 | 0.99091826 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.70967742 | 0.70298314 | |

Metrics and associated errors were measured in the Databricks Experiment UI which performs all measurements automatically against the model after EDD evaluation. Above table is a summary of all iterations from MVP_5 onwards (before mvp_5, in phase one, iterations mainly focused on resolving environment, dataset and pipeline issues, so there are no metric evaluations associated).

Examples of Experiment results were used throughout the report. The Experiments UI in Databricks offers full granularity pertaining results, from experiment wide metrics, through individual metrics and comparisons, to individual prompt-response pair investigation and a full RAG chain trace.

## 13.4 FIGURES, TABLES AND VISUALISATIONS

### 13.4.1 Qualitative tests – Reviews and Playground comparison continued

#### 2. Question: Who is death grant paid to?

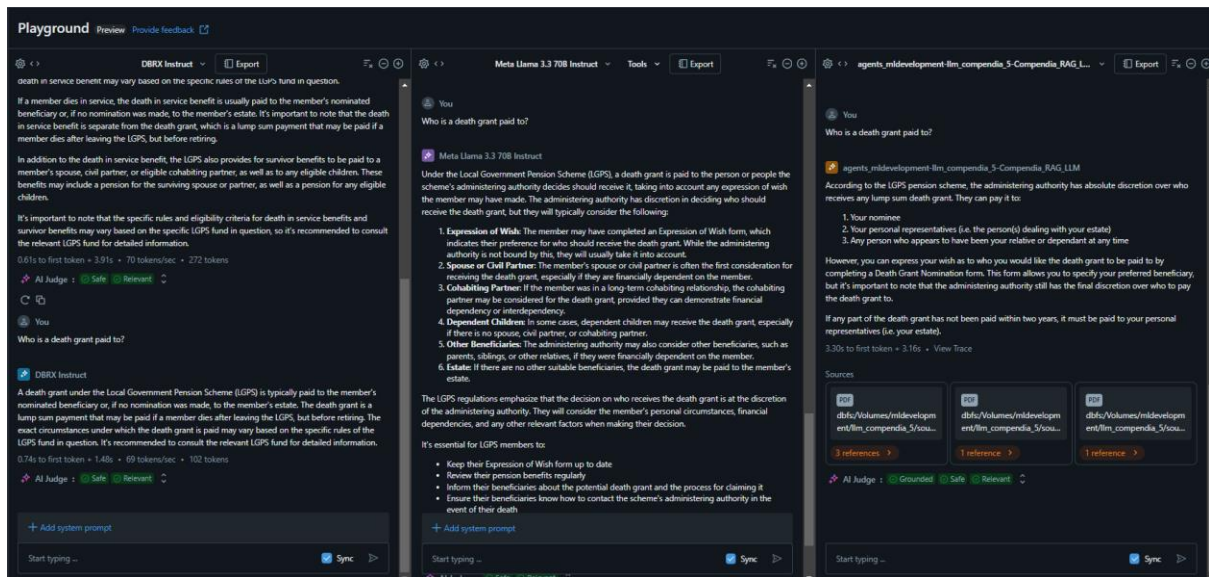Expected Answer (based on FAQs):

Your local pension fund has absolute discretion over who receives any lump sum death grant. This means it can be paid to:

the person(s) you name on your expression of wish form your personal representative(s) – this is someone who is responsible for your looking after your estate when you die anyone who appears to have been, at any time, a relative or dependent. You can let your pension fund know who you would like any death grant paid to by completing an expression of wish form. Your pension fund will take your wishes into account when deciding who the death grant is paid to. Contact your pension fund to ask for a form.

If a death grant is not paid within two years it must be paid to your personal representative(s) and a tax charge may apply.

If you paid AVCs and were a member of the LGPS on or after 1 April 2014, your local pension fund will use its discretion to decide who receives any lump sum due from your AVC pot when you die. If you left the LGPS before 1 April 2014, your local pension fund must pay any lump sum due to your estate.

3.      Question: What death grant is paid if I die when I am still
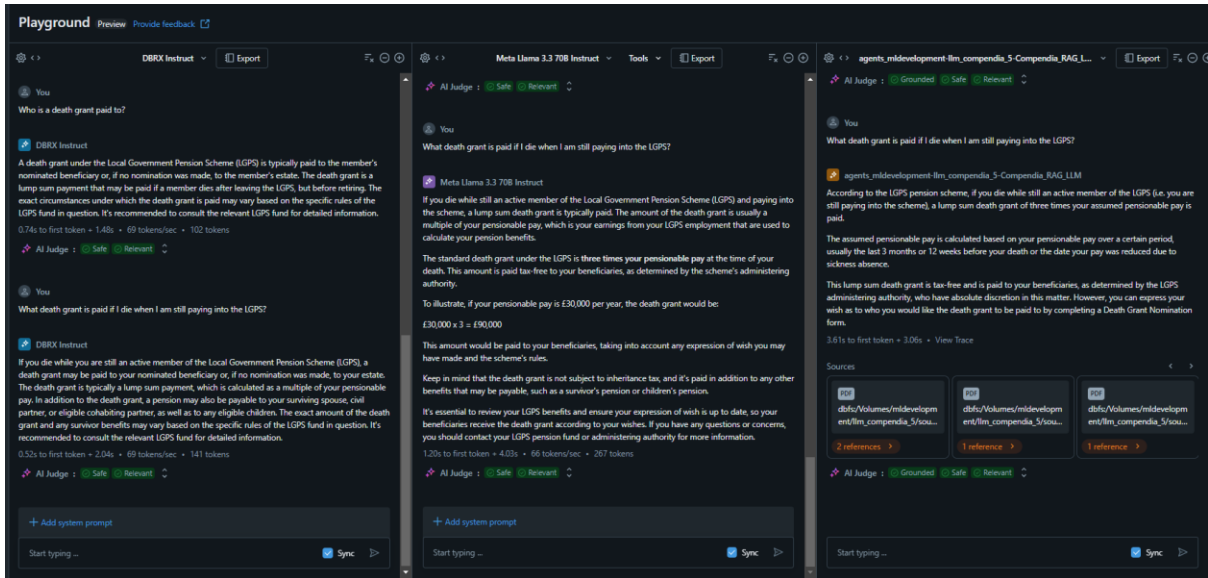
paying into the LGPS?

FAQ Answer:

If you die when you are paying into the LGPS, a lump sum death grant of three times your annual pay is paid. It does not matter how long you have been a member of the LGPS, you are protected from the day you join the Scheme.

If you also have a deferred benefit and/or a pension in payment from a previous period of membership of the LGPS, the lump sum death grant paid is the greater of:

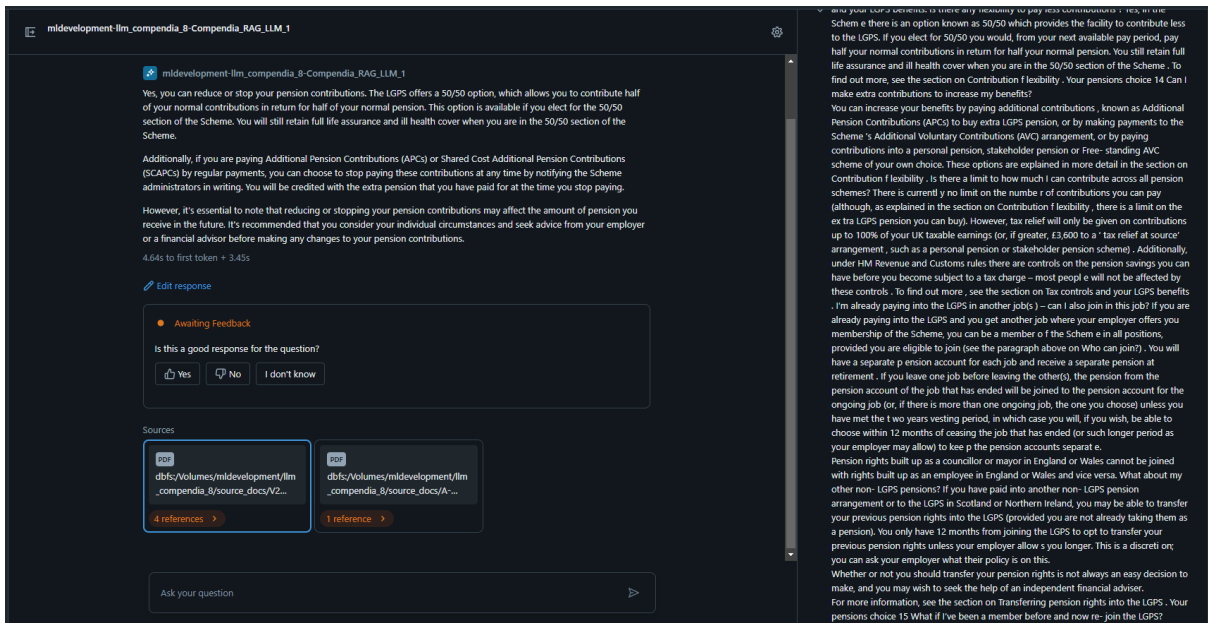the total of any lump sum death grants payable from the deferred benefit and/or pensions in payment three times your annual pay at the date you die The annual pay used to work out the death grant is based on Assumed Pensionable Pay. If an independent registered medical practitioner certifies that you were working reduced contractual hours because of the ill health which led to your death during the period used to work out Assumed Pensionable Pay, the amount will increase. The Assumed Pensionable Pay is instead the pay you would have received if you had not been working reduced contractual hours.
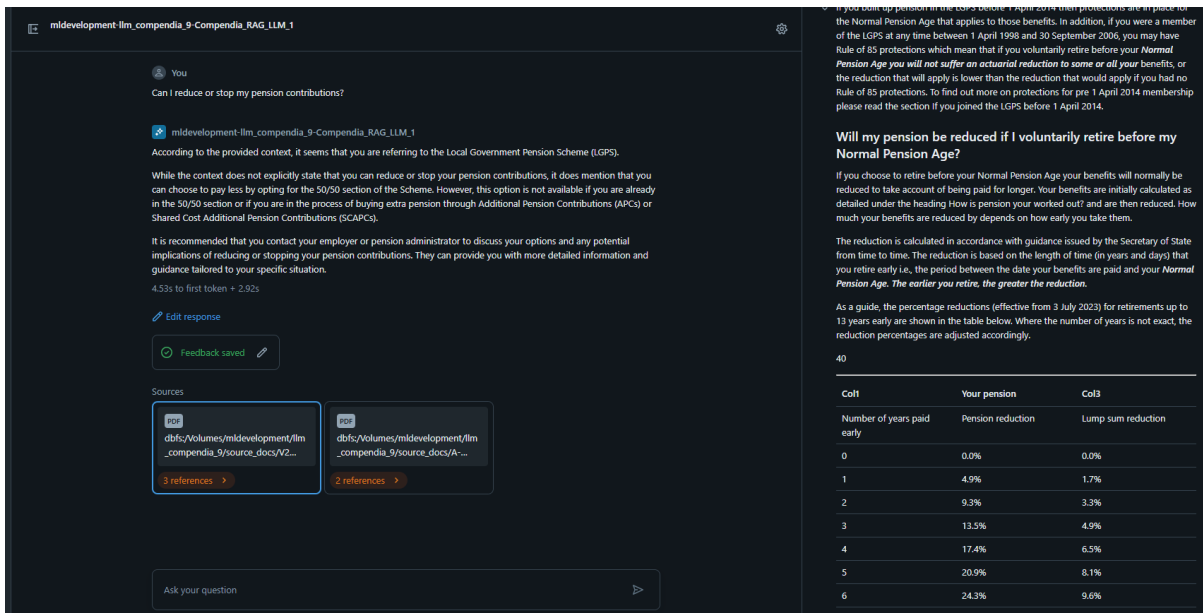
If you pay Additional Voluntary Contributions (AVCs) arranged through the LGPS, the value of your AVC fund is also payable.

## 4. More examples of bot's qualitative performance and checks

MVP 9 uses PyMuPDF4LLM parser which transforms the pdfs into markdown files first to then chunk and tokenize them, improving the understandability of the files for the model but also improving the Review App UI with markdown notation being automatically read:

## 13.4.2 Improvements to the evaluation dataset (FAQ questions)

For example, one FAQ question ('What information will I need?') forms part of a divorce section, but itself makes no reference to this aspect, making it difficult to retrieve correct context, provide relevant output, especially without any chat history. It is a testament to the potential of this technology that the LLM caught relevant context despite this difficulty.
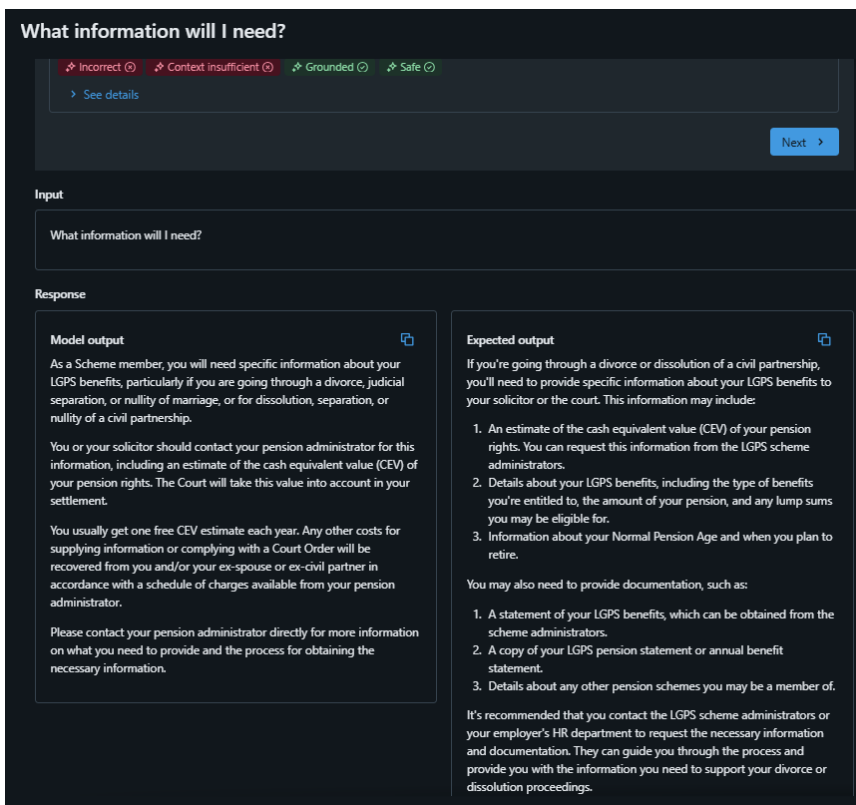


*Figure 56 - Experiment Investigation UI, comparing response received against expected output*
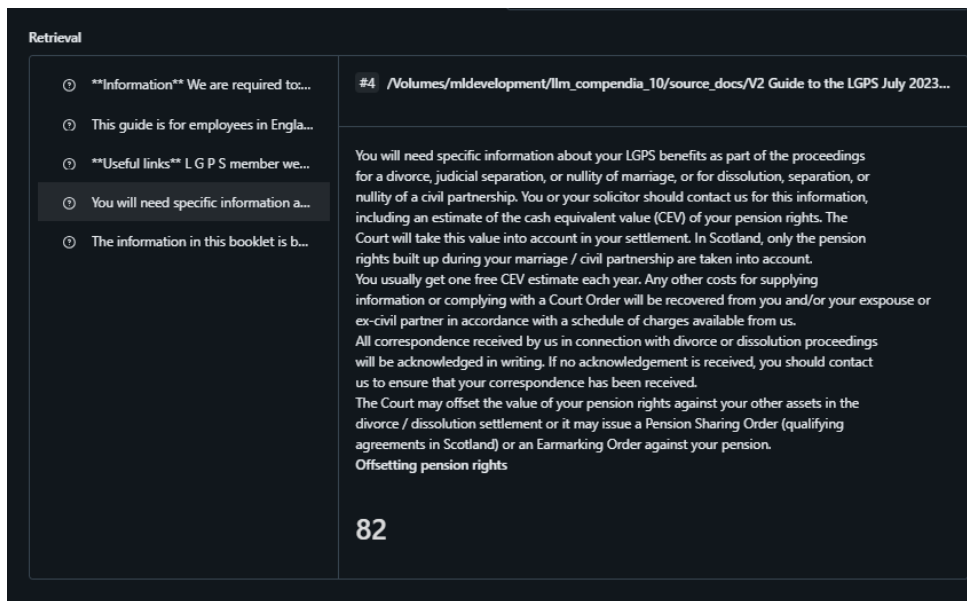
*Figure 57 - Context Retrieval example from Experiment UI*

### 13.4.3 System Message Iterations

The base message shown in below figure was changed, in early iterations where LGPS documentation was used, to start with *'You are a pensions professional expert assistant that answers questions'* and end with 'keep your answers brief, to the point and specific'. This caused the bot to provide answers that were very short and often did not provide enough information when answering questions.



*Figure 58 - base RAG LLM system message as provided by Databricks demo*

*Last sentence was replaced by 'Start the conversation by stating the above facts about yourself in summary form',* but this did not have the desired effect outside of the bot using more bullet points.



*Figure 59 - Application RAG LLM System Message to elicit specific behaviour*

This was expanded as per figure above, however it seemed to demand too much, as the bot cannot articulate sources without being given that information in preprocessing. Additionally, this system prompt caused the bot to be extremely cautious to not answer from its own training data and to stick to the context only. Since the context was quite broad and some details was missing the bot often answered that it does not have sufficient context to provide an answer.



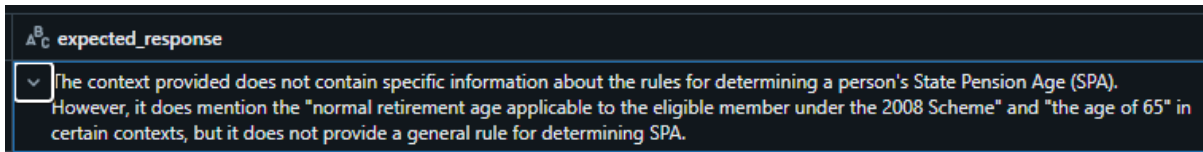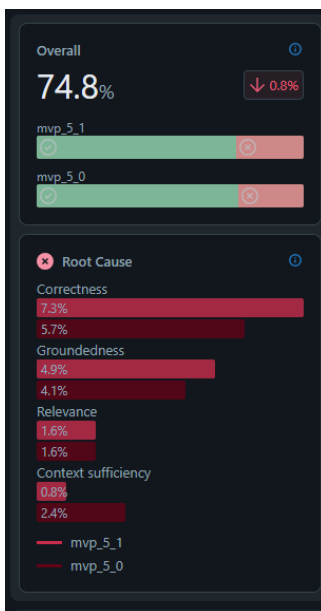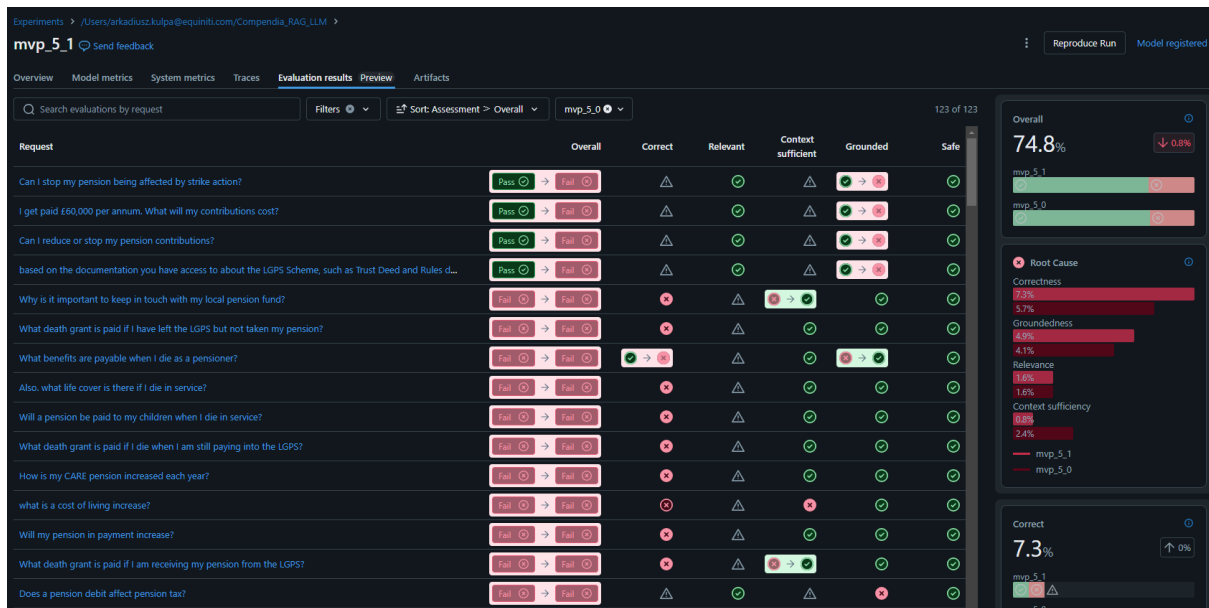*Figure 60 - question: 'What are the rules for determining a person's SPA?'*

While bot's ability to say 'I do not know' is especially important for the internal use case (a measure of certainty), the public bot should instead form a safe generic answer and direct user to relevant authority.

### 13.4.4 Quantitative tests – Experiments and their comparison

Comparing experiments - Mvp_5 a stable version of the public facing bot

# 14 BIBLIOGRAPHY

Chen, J., Lin, H., Han, X., & Sun, L. (2024, 12 29). *Benchmarking Large Language Models in Retrieval-Augmented Generation*. Retrieved from ArXiv Computer Science: https://arxiv.org/abs/2309.01431

Databricks. (2024, 12 26). *Assess performance: Metrics that matter*. Retrieved from Databricks on AWS: https://docs.databricks.com/en/generative-ai/tutorials/ai-cookbook/evaluate-assess-performance.html

Databricks. (2024, 12 15). *csrd_assistant - Databricks*. Retrieved from Databricks Solution Accelerators: https://databricks-industry-solutions.github.io/csrd_assistant/#csrd_assistant.html

DataBricks. (2024, 12 10). *Databricks Generative AI Cookbook*. Retrieved from Databricks Generative AI Cookbook: https://ai-cookbook.io/

DataBricks. (2024, 12 10). *Databricks Launches DBRX, A New Standard for Efficient Open Source LLM*. Retrieved from Databricks: https://www.databricks.com/company/newsroom/press-releases/databricks-launches-dbrx-new-standard-efficient-open-source-models

Databricks. (2024, 12 10). *Step 2: Deploy POC to collect Stakeholder Feedback*. Retrieved from Databricks Generative AI Cookbook: https://ai-cookbook.io/nbs/5-hands-on-build-poc.html

Databricks. (2025, 01 04). *Mosaic AI Agent Evaluation LLM judges reference*. Retrieved from docs.databricks: https://docs.databricks.com/en/generative-ai/agent-evaluation/llm-judge-reference.html

Donovan, R. (2025, 01 07). *Breaking up is hard to do: Chunking in RAG applications*. Retrieved from Stack Overflow: https://stackoverflow.blog/2024/12/27/breaking-up-is-hard-to-do-chunking-in-rag-applications/

LangChain. (2024, 12 10). *LangChain*. Retrieved from LangChain.com: https://www.langchain.com/

LangChain. (2025, 01 07). *Databricks Vector Search | LangChain*. Retrieved from python.langchain.com: https://python.langchain.com/docs/integrations/retrievers/self_query/databricks_vector_search/

Leng, Q., Uhlenhuth, K., & Polyzotis, A. (2025, 01 04). *Best Practices for LLM Evaluation of RAG Applications*. Retrieved from databricks.com: https://www.databricks.com/blog/LLM-auto-eval-best-practices-RAG

MlFlow. (2025, 01 07). *Autogen Image Agent*. Retrieved from MlFlow: https://mlflow.org/blog/autogen-image-agent

Ravenwolf, W. (2024, 12 10). *LLM Comparison/Test: 25 SOTA LLMs*. Retrieved from Huggingface: https://huggingface.co/blog/wolfram/llm-comparison-test-2024-12-04

Smilkov, D., Peter, E., Frankle, j., Trott, A., Singh, A., Polyzotis, A., . . . Kulinski, S. (2024, 12 15). *Streamline AI Agent Evaluation with New Synthetic Data Capabilities*. Retrieved from Databricks.com: https://www.databricks.com/blog/streamline-ai-agent-evaluation-with-new-synthetic-data-capabilities

Zhang, X., Zhang, Y., Long, D., Xie, W., Dai, Z., Tang, J., . . . Zhang, M. (2024, 12 16). *mGTE: Generalized Long-Context Text Representation and Reranking Models for Multilingual Text Retrieval.* Retrieved from ArXiv: https://arxiv.org/pdf/2407.19669

Zheng, L., Chiang, W.-L., Sheng, Y., Zhuang, S., Wu, Z., Zhuang, Y., . . . Stoica, I. (2023, 06 09). *Judging LLM-as-a-Judge with MT-Bench and Chatbot Arena*. Retrieved from arXiv:2306.05685: https://arxiv.org/abs/2306.05685#

Zhou, K., Zhu, Y., Chen, Z., Chen, W., Xin Zhao, W., Chen, X., . . . Han, J. (2023, 11 3). *Don't Make Your LLM an Evaluation Benchmark Cheater.* Retrieved 12 10, 2024, from ArXiv: https://doi.org/10.48550/arXiv.2311.01964