

Techniques I used for Dataset Development and Refinement involves–

Data Collection:

- **Development Phase:**
 - Curate a diverse set of questions related to your business domain. Ensure that questions cover various aspects and nuances of the topic.
 - Collect relevant documents or context passages that could potentially answer the questions. This might involve extracting information from manuals, articles, or other authoritative sources.
- **Refinement Phase:**
 - Clean the collected data to remove noise, irrelevant information, or biased content.
 - Annotate the data with labels indicating the relevance of each document to the corresponding question. Use domain experts for accurate annotations
- **Development Phase:**
 - Ensure a balanced distribution of question types and topics to avoid bias in the model.
 - Incorporate a mix of long and short questions to cover a range of user queries.
- **Refinement Phase:**
 - Review the distribution of annotated labels to ensure a balanced representation of relevant and non-relevant documents.
 - Adjust the balance if necessary, especially if the dataset shows skewness toward one class.

Data Augmentation:

- **Development Phase:**
 - Introduce variations in questions and context passages to simulate real-world scenarios.
 - Add paraphrased versions of questions to diversify language patterns.
- **Refinement Phase:**
 - Apply data augmentation techniques such as back-translation or synonym replacement to increase the dataset's diversity.
 - Ensure that augmented data retains semantic integrity.

Brief Comparison of Language Model Fine-Tuning

Approaches:

Traditional Supervised Fine-Tuning:

- **Description:**
 - Fine-tune the language model on a labeled dataset with paired input-output examples.
- **Pros:**
 - Straightforward and easy to implement.
 - Effective for tasks with clear labels.
- **Cons:**
 - May struggle with capturing nuances and subtleties in natural language.

Transfer Learning with Pre-trained Models:

- **Description:**
 - Utilize a pre-trained language model (e.g., BERT, GPT) and fine-tune it on a task-specific dataset.
- **Pros:**
 - Leverages pre-existing knowledge from large corpora.
 - Effective for various NLP tasks.
- **Cons:**
 - Requires a substantial amount of data for fine-tuning.

Reinforcement Learning Fine-Tuning:

- **Description:**
 - Apply reinforcement learning to fine-tune the model based on rewards or feedback.
- **Pros:**
 - Can optimize for specific objectives.
 - Effective for tasks where quality is subjective.
- **Cons:**
 - Requires careful reward design and might be sensitive to the reward structure.