Introduction :

In the realm of natural language processing (NLP), document summarization and question answering on documents have emerged as vital components in the quest for efficient information retrieval and comprehension. This project stands at the intersection of these two fundamental tasks and leverages the power of Large Language Models (LLMs) to redefine how we approach document understanding.

In the past, endeavours in text summarization often relied on rule-based methods or extractive techniques that involved identifying and extracting key sentences or phrases from a document. Similarly, traditional question-answering systems relied on structured data sources or manually curated knowledge bases. While these methods have yielded valuable insights, they are inherently limited by their inability to capture the nuanced relationships and contextual intricacies present in unstructured text. Hence, the task of extracting valuable insights from vast textual datasets has been a longstanding challenge.

Recent advancements in NLP, particularly the advent of LLMs, have catalysed a paradigm shift in how we process and comprehend textual data. LLMs, such as GPT-3 and its successors, have demonstrated remarkable capabilities in understanding, generating, and summarising text [1]. Their innate capacity to grasp the semantics, context, and subtleties of language has opened new avenues for tackling complex NLP tasks.

The project introduces a novel approach that draws inspiration from earlier text summarization and question-answering systems but takes a giant leap forward by harnessing the immense potential of LLMs. Unlike traditional summarization techniques that often struggle with maintaining context or generating coherent summaries, the recursive summarization strategy employed harnesses LLMs' contextual understanding to produce highly informative and cohesive summaries. By recursively summarising document chunks, we preserve the essence of the source text while achieving concise abstractions.

Furthermore, the integration of retrieval augmented generation, a core component of this project, distinguishes it from its predecessors. By ingesting documents into a queryable format and coupling this with the retrieval augmented generation chain, we enable the model to fetch and incorporate relevant documents in real-time when posed with questions [2]. This dynamic and adaptive approach empowers to respond to a wide array of queries with up-to-date and contextually accurate answers.

In summary, while previous projects have made significant strides in text summarization and question answering, this approach stands out as a pioneering endeavour that capitalises on the capabilities of LLMs, thereby redefining the boundaries of what is achievable in document understanding. This project offers the potential to revolutionise information retrieval, enabling users to effortlessly distil insights from extensive textual resources while maintaining contextual richness and accuracy.

References:
[1] Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., ... & Amodei, D. (2020). Language models are few-shot learners.

[2] Dai, Z., Yang, Z., Yang, F., Cohen, W. W., Carbonell, J., Le, Q. V., & Salakhutdinov, R. (2020). Transformer-based large scale language models require a new training paradigm.