# Custom Search Engine On AWS

## Project Report

Aranya Singh Chauhan (CSE, SRM IST)
Eshaan Mathakari (ECE, SRM IST)

## Problem Statement

Search as a capability is an important feature that is required by almost all medium and large enterprises as search helps filter relevant and required information in the world of big data. Search helps find relevant information quickly and saves time to go through vast information.

When given a large dataset, it becomes extremely difficult for an individual or an organization to fetch a piece of information that they seek. In addition to this, even with traditional searching methods (searching for keywords in a document), it can be tough to locate that particular piece of information if one isn't sure about what keywords to look for in the dataset.
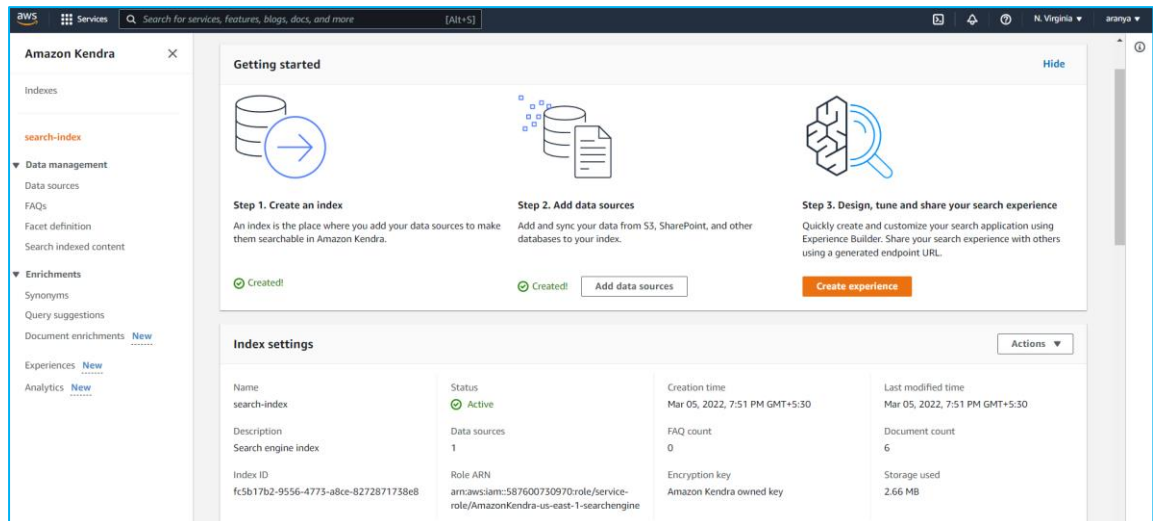
## Project Description

- The project aims to deploy a search engine on AWS, that not only eliminates the problems mentioned in the section above but also provides features that make the overall system keep on improving over time with the help of feedback provided by the users.

- The project makes use of **AWS Kendra**, which is an intelligent search service provided by machine learning. With AWS Kendra, we build a search engine that can use a dataset stored within an **AWS Simple Storage Service (Amazon S3) Bucket**.

- The search not only looks for keywords in the dataset, but it can also look for synonyms to the keywords that are being searched. This gives the user the ability to search for relevant content even if they are not aware of the exact keywords that they need to find. The users can even make use of natural language as their queries.

- Most commonly asked questions/queries related to the dataset can be added to the search engine as a configuration as well. This helps the search engine display the results extremely quickly if the query that the user is entering contains an FAQ, thus reducing the engine's time complexity.
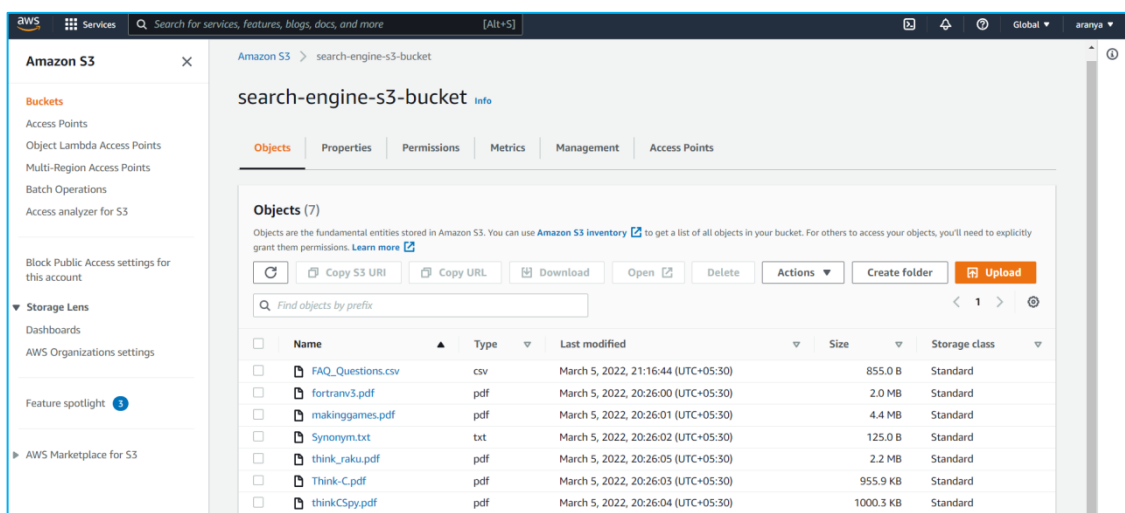
## Approach Used

1. **Index Creation on AWS Kendra:**
   An *index* serves the purpose of acting as a dataset source that is to be used by the search engine. Connectors are used to connect the dataset with AWS Kendra. Here, we have used an S3 bucket as the data source for the connector.
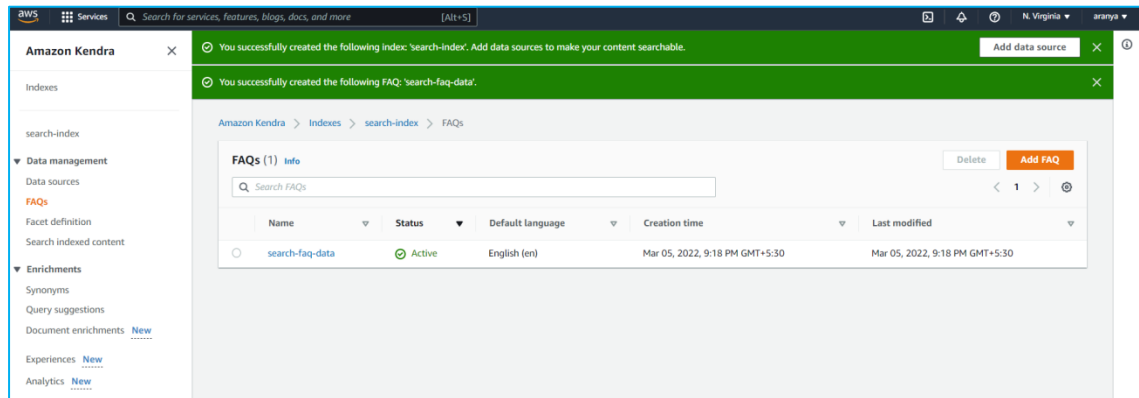


**AWS Kendra Index**

2. **Uploading Database:** A *data source* is a location, such as an Amazon Simple Storage Service (Amazon S3) bucket, where you store the documents for indexing. One can automatically synchronize data sources with an Amazon Kendra index so that new, updated, or deleted documents in the data source are also added, updated, or deleted in the index for searching on.
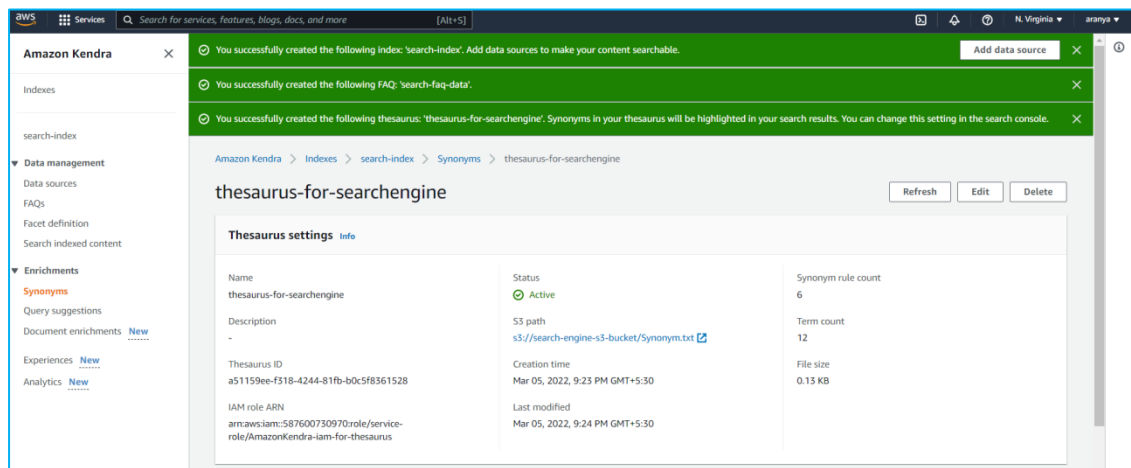


**Dataset for Search Engine using S3**

3. **Adding Search Functionality to the Engine:** To get answers, we need to query an index. With Kendra, users can use natural language in their queries. The response contains information, such as the title, a text excerpt, and the location of documents in the index that provide the best answer.

4. **Adding FAQs and Synonyms to optimize search results:**



Data for FAQ Search Function

In order for the search engine to be able to display the search results efficiently and quickly, we add documents containing FAQs related to the dataset along with synonyms for commonly searched terms. This not only gets the engine to display results faster, but also makes it easier for the user to search without having to worry about the keywords.



Thesaurus Data for Synonyms Function
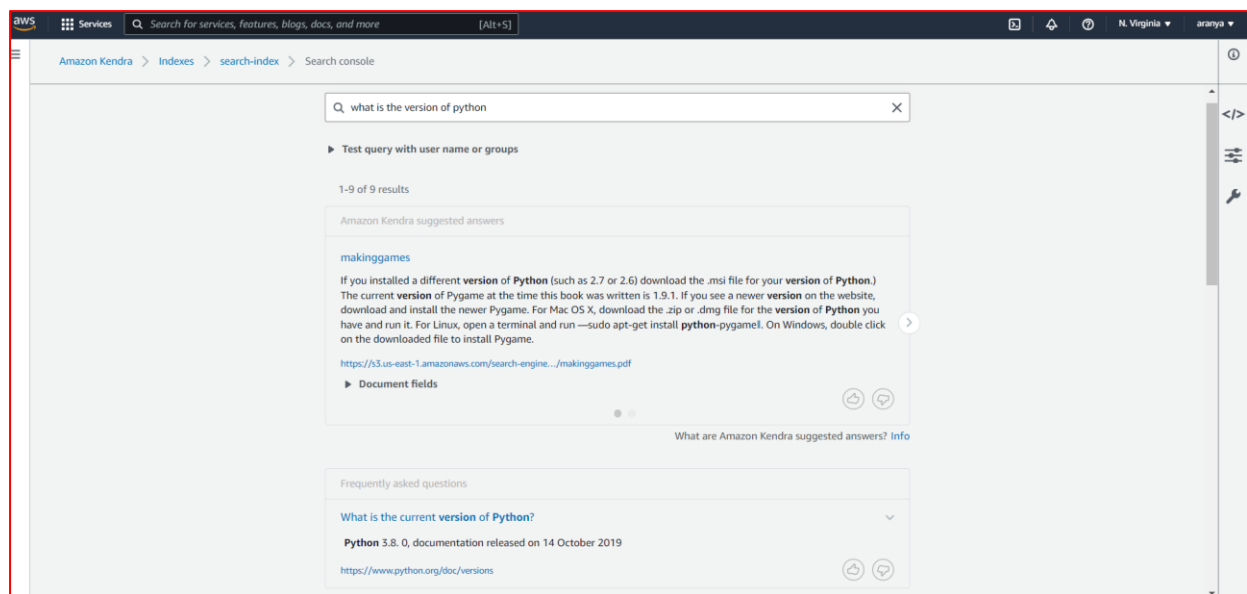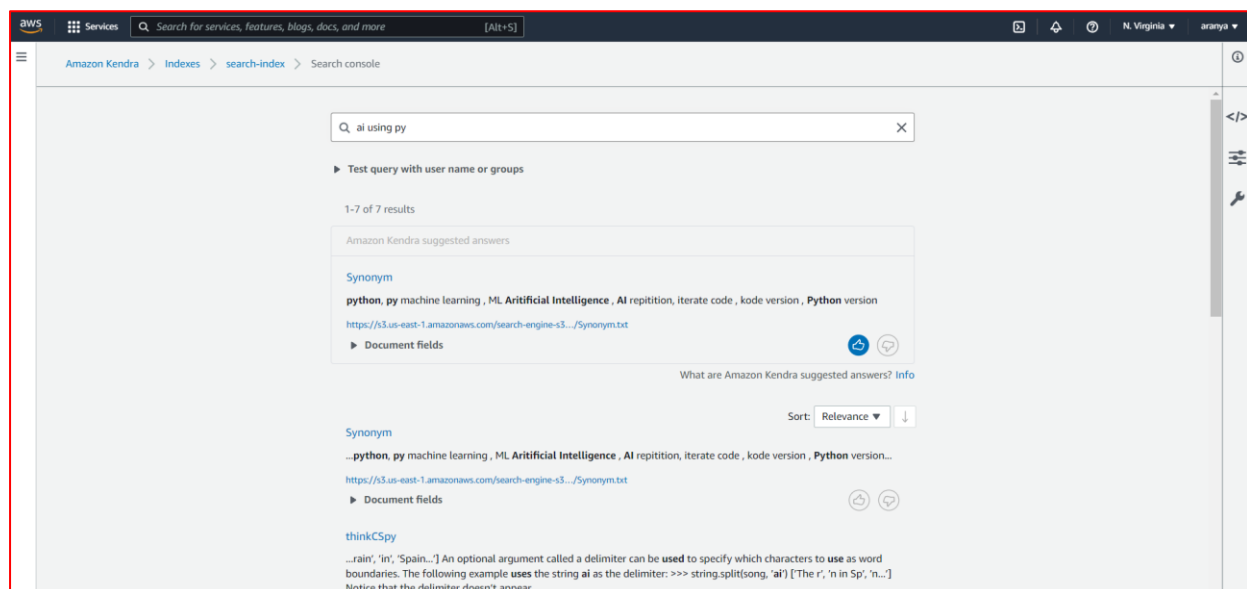
**Working of the Project:**



Figure 1FAQ Search Search Results for FAQs



Search with Synonyms, FAQ functions implemented.

**Conclusion**

The project was successfully deployed on AWS, using Kendra as the primary search tool and S3 for storage. Users can successfully search for content within the dataset without having to type-in the exact keywords. They can also provide feedback in regards to the results that are displayed, so that the tool can use its Machine Learning capabilities and improve when it comes to displaying accurate results.