UTS Group 24

# NLP ASSIGNMENT 3

*Organising a text corpus using K-means clustering and applying retrieval-augmented generation (RAG) to a chatbot*

Warit Boonmasiri - 25399522
Mohammad Alhajjeh - 24832800
Jada Gamis - 24823787
Yingrong Zhang - 25428842
Aiden Blishen Cuneo - 24971160

# Overview

## Problem

Our project addresses the problem of the inability to keep up with an influx of information across a range of scientific fields, using natural language processing (NLP) techniques introduced in this course to visualise and organise scholarly literature.

## Aim

Clustering a text corpus into related groups, then use RAG to allow a chatbot to use one of these clusters as context. The goal is for the user to be able to swap between these clusters very quickly, allowing them to decide which specific context they want to give to the chatbot.

## Solution

- Use K-means clustering on the dataset
- Allow the user to swap between clusters to choose the context for RAG
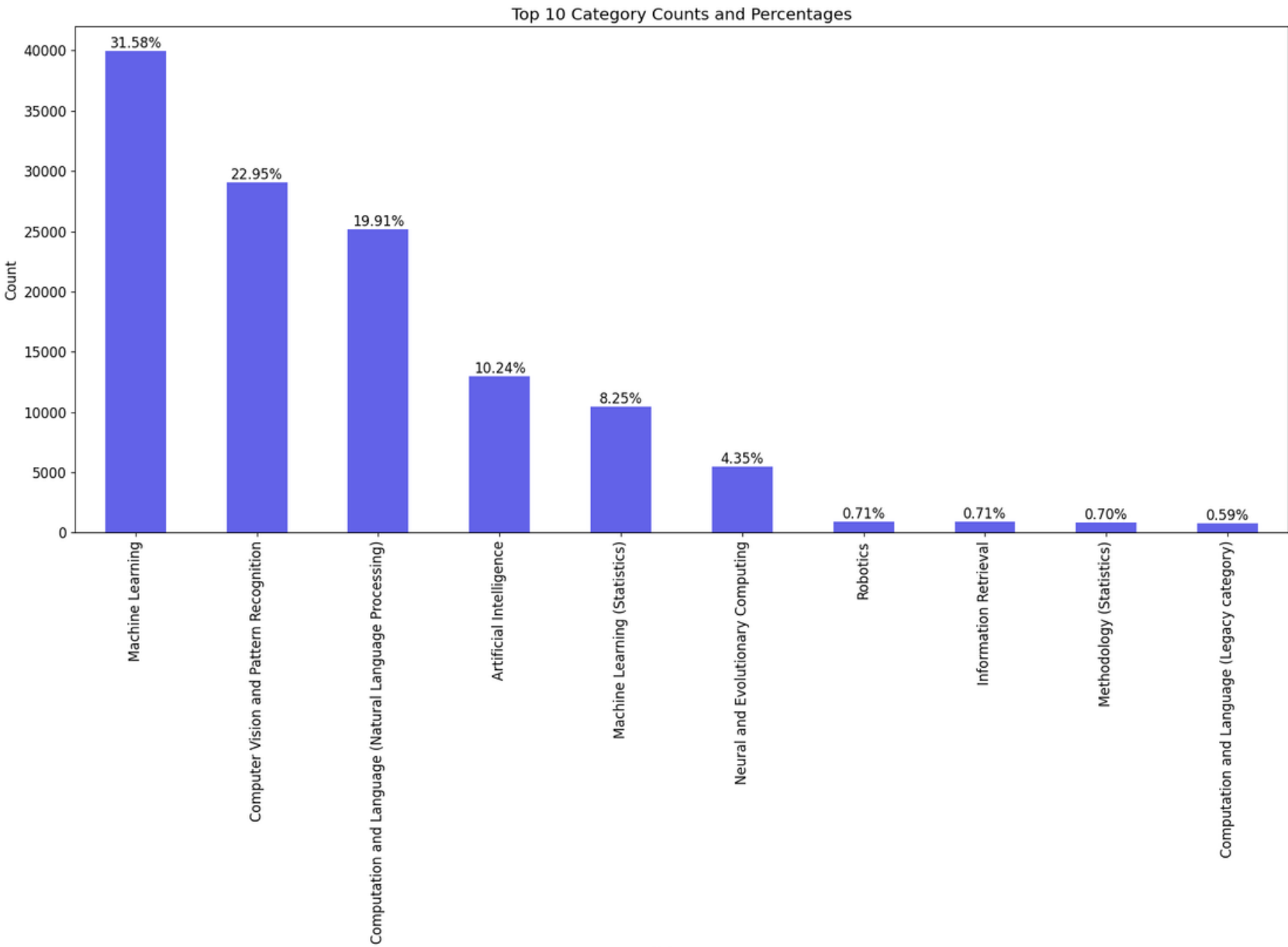- The model generates a response based on this context

# Dataset

Figure 1: Bar graph showing article categories and their frequencies



Figure 2: Sample rows from the chosen dataset

- arXiv Scientific Research Papers Dataset
- Contains over 100,000 articles
- Topics such as Artificial Intelligence, Machine Learning, Computer Science, etc.

# Clustering

**01** **How clustering was performed**

- Convert text summaries to vector representation
- Apply PCA to reduce dimensionality
- Use K means clustering
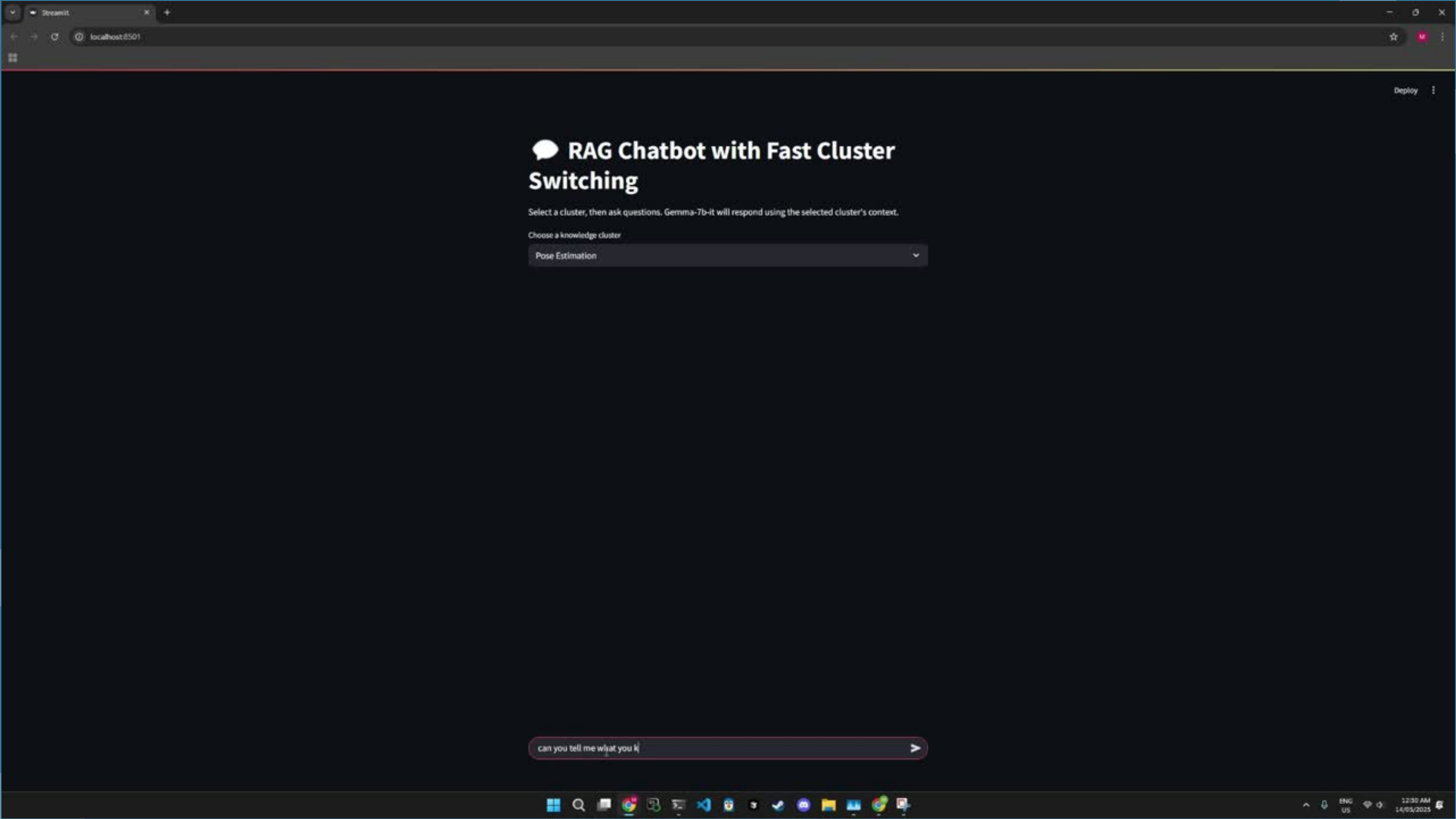- Do final dimension reduction to visualise in 2D

**02** **Why cluster the data?**

- Clustering splits the summaries into more specific areas
- RAG can be used on these topics
- Performance boost for a smaller set of summaries



Title: ZK-GanDef: A GAN based Zero Knowledge Adversarial Training Defense for Neural Networks
Author(s): ['Guanxiong Liu', 'Issa Khalil', 'Abdallah Khreishah']
Title: Towards Speeding up Adversarial Training in Latent Spaces
Author(s): ['Yaguan Qian', 'Qiqi Shao', 'Tengteng Yao', 'Bin Wang', 'Shouling Ji', 'Shaoning Zeng', 'Zhaoquan Gu', 'Wassim Swaileh']
Title: Effective and Robust Detection of Adversarial Examples via Benford-Fourier Coefficients
Author(s): ['Chengcheng Ma', 'Baoyuan Wu', 'Shibiao Xu', 'Yanbo Fan', 'Yong Zhang', 'Xiaopeng Zhang', 'Zhifeng Li']

# Demo of UI

# THANK YOU!