

# Utilizing Hadoop Ecosystem to Implement RAG for Chain of Thought Reasoning in Language Models

DBMS Capstone Project by  
MDA: Master in Database Administration





# Group Members

Aviral Pandey

21CS10089

Harsh Sharma

21CS30023

Ishan Raj

21CS10032

Prasanna Paithanakar

21CS30065

Utsav Dhanuka

21CS10087





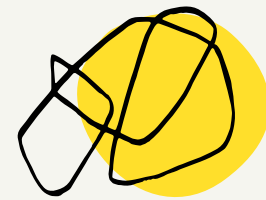
# Introduction

## Part 01





# The Problem

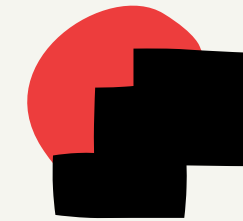


## What we want to solve

RAG is one of the latest developments in the fields of LLM's and Databases play an important role in implementation of those.

Furthermore, CoT reasoning is needed for improved generation.

We propose and implement a project to accomplish the above.



## Plan

To use MapReduce to find related data to prompt and utilise that in the RAG framework. The MapReduce would be implemented on a Solr distributed database.



# Solution abstract

This project proposes integrating Retrieval Augmented Generation (RAG) into the Hadoop ecosystem to enhance coherence and context retention in large language models' sequential reasoning. Leveraging Hadoop's scalable, distributed computing capabilities, the approach aims to enhance the efficiency and accuracy of natural language understanding and reasoning systems in queries that involve Chain of Thought Reasoning.



# Stack & Architecture

## Part 02





# Tech Stack

## Apache Zookeeper

Open-source server for distributed coordination, offering a hierarchical key-value store and essential services for large-scale distributed systems.

## Apache Solr

Highly reliable, scalable, fault-tolerant search engine powering major internet sites with distributed indexing, replication, and load-balanced querying.

## Langchain

A framework simplifying application development with large language models, supporting tasks such as document analysis, chatbots, and code analysis.

## Grafana

Multi-platform open-source analytics and visualization tool generating charts, graphs, and alerts from supported data sources.

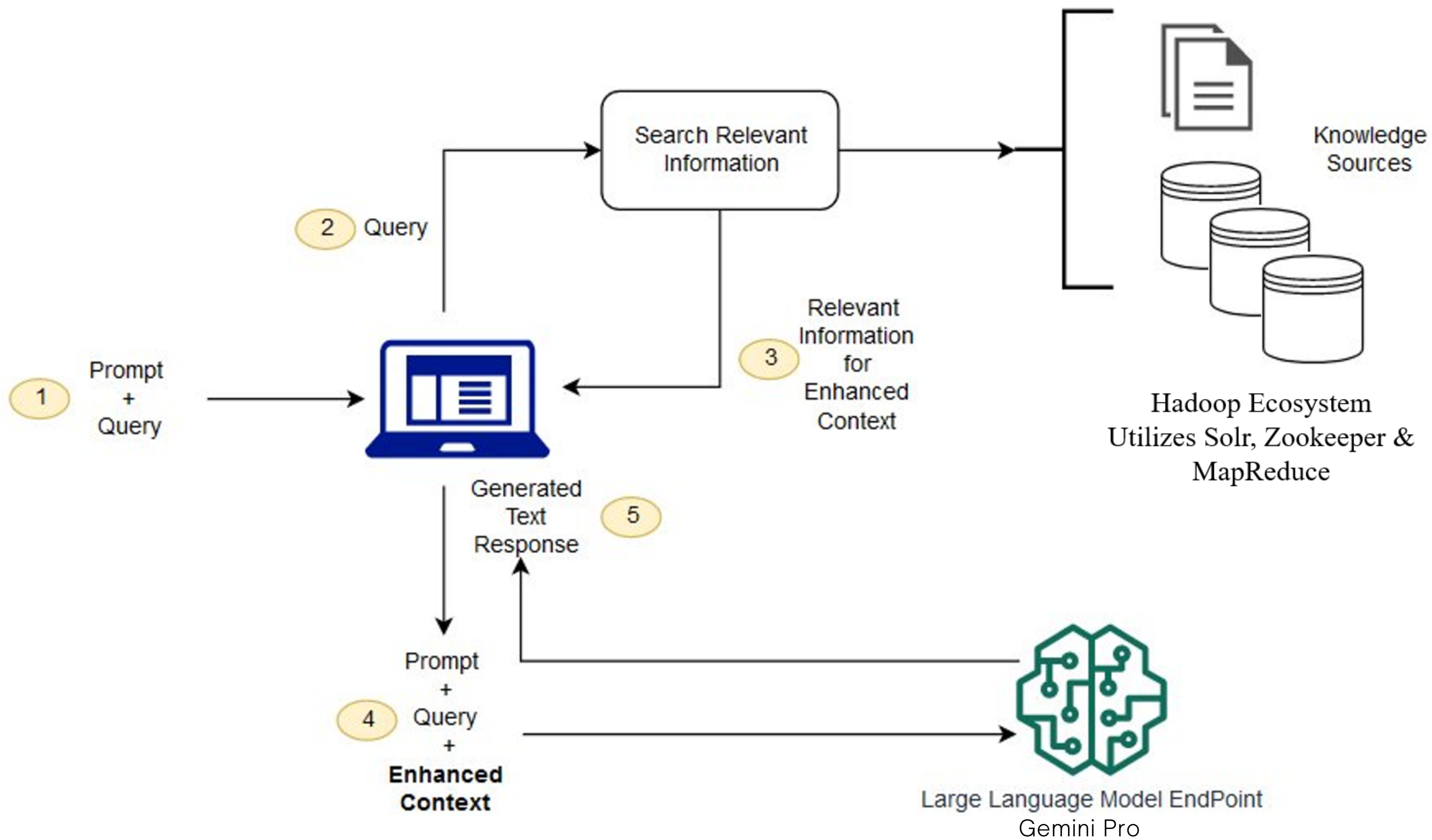
## Prometheus

Free event monitoring and alerting software with a time series database, flexible queries, and real-time alerting.

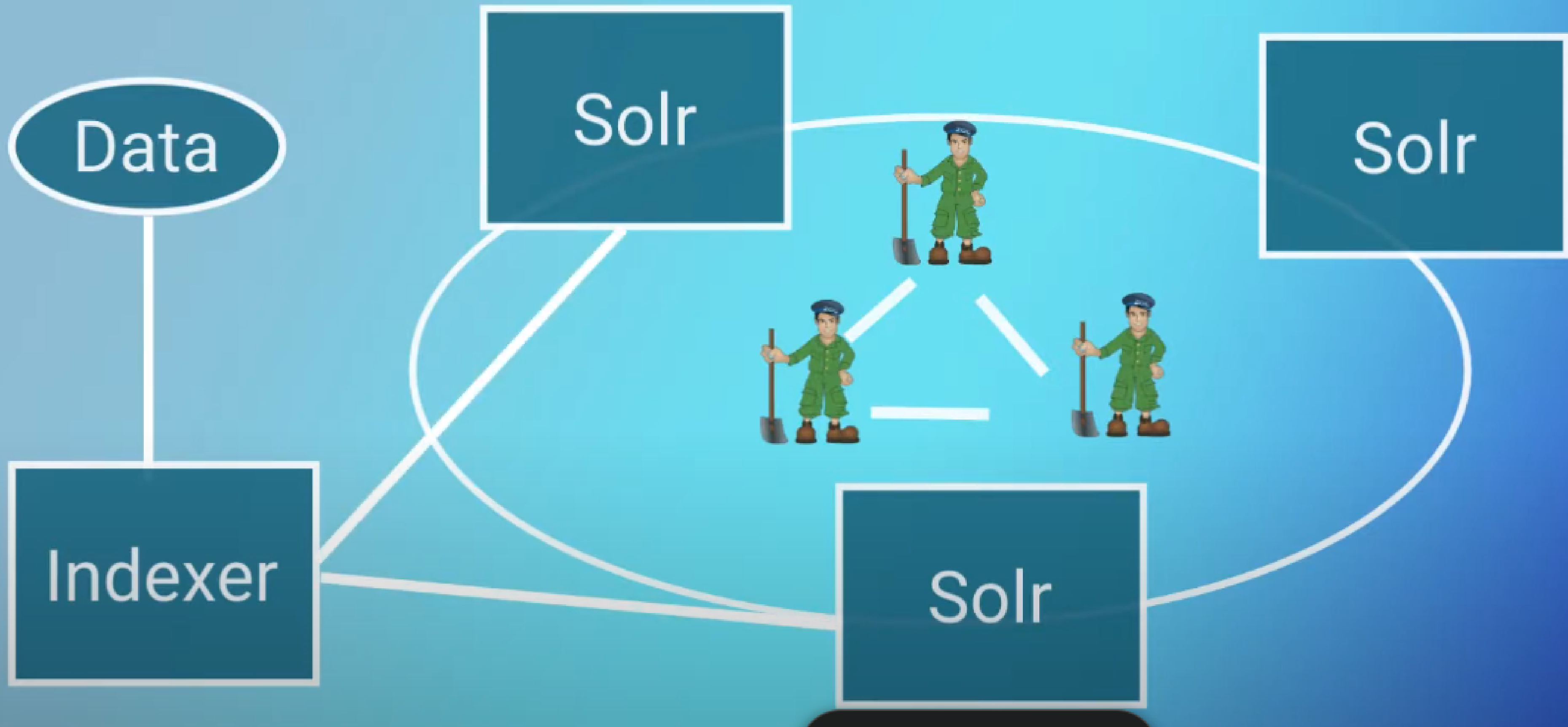
## Flask API

Lightweight Python web framework for effortlessly building and deploying RESTful APIs.











Solr-Zookeeper framework is used with 3 nodes of each connected to each other providing data security and load balancing.

---

DISTRIBUTED DATABASE

**Project highlights**

The documents stored as embedding in the database provide the Gemini Pro LLM an optimised prompt to generate better responses on complicated user prompts

---

CHAIN OF THOUGHT  
REASONING

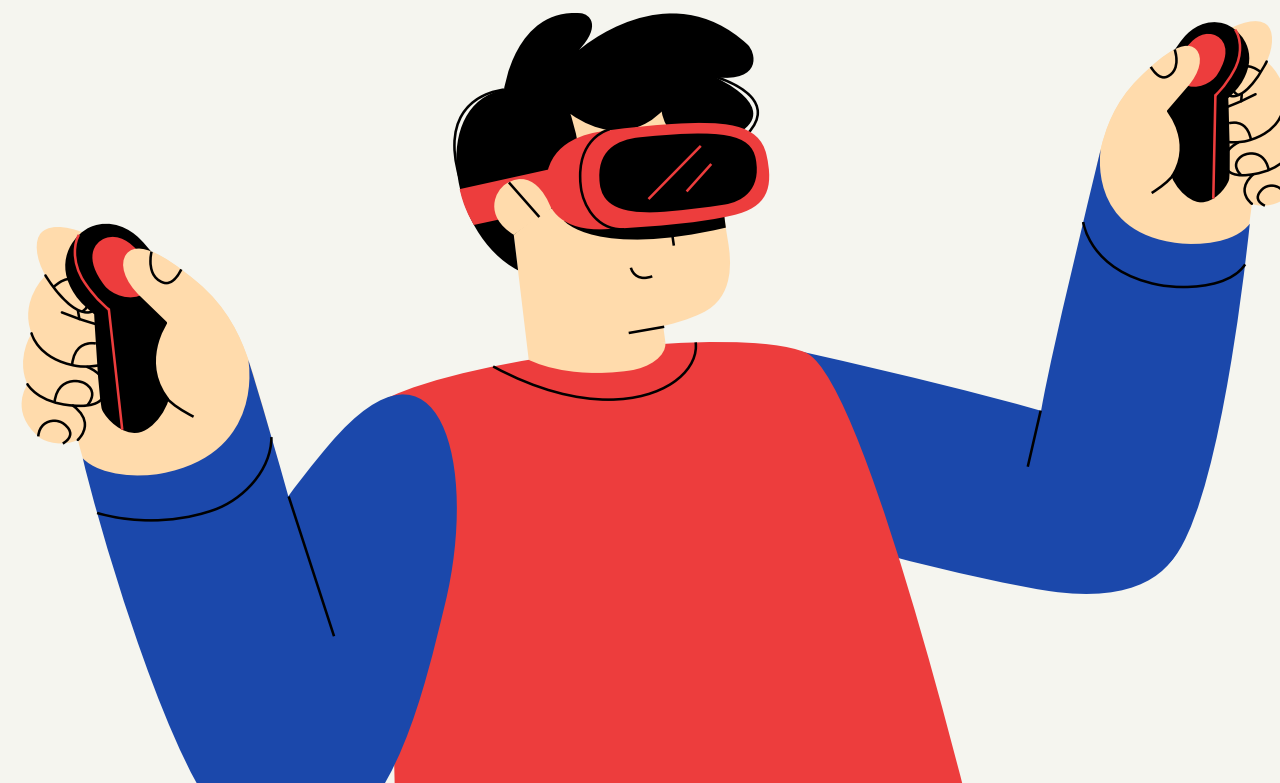
**Significant improvements**

We have used Prometheus and Grafana to monitor our nodes which makes the application scalable to industrial levels.

---

SCALABILITY

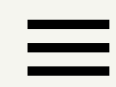
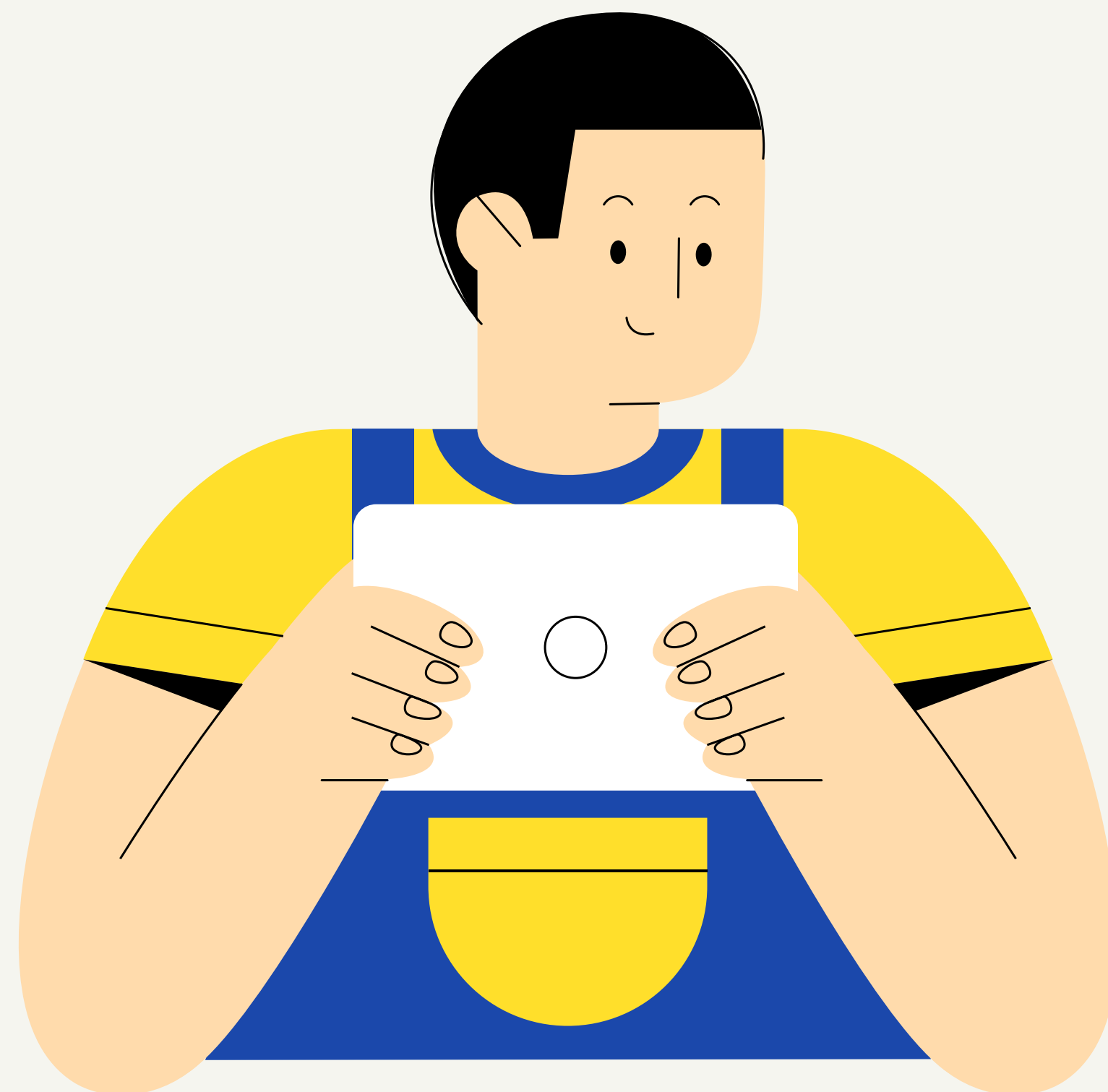
**Speciality**





# Demo

## Part 03



### Question

tell me about professor Pabitra Mitra

### Answer

Pabitra Mitra is a professor of Computer Science and Engineering at Indian Institute of Technology Kharagpur. He did his PhD from Indian Statistical Institute Calcutta and B. Tech from Indian Institute of Technology Kharagpur. He has been an Assistant Professor at IIT Kanpur and Scientist at Centre for AI and Robotics Bangalore. He received the INAE Young Engineer Award, IBM and Yahoo Faculty Awards. He has co-authored a book and about 100 research papers in pattern recognition and machine learning. Thanks for asking!

### Question

research areas of Pabitra Mitra

### Answer

Pabitra Mitra's research interests are in the fields of data mining, knowledge discovery, medical expert systems, and learning. Thanks for asking!

### Question

data mining related paper of pabitra mitra

### Answer

Pabitra Mitra has a paper titled "Data Mining: A Knowledge Discovery Approach" published in the book "Data Mining and Knowledge Discovery for Real-Life Applications" edited by A. Goswami, P. Mitra, R. Bhatt, and T.R. Rao. You can find this book on Amazon. Thanks for asking!

### Question

Tell me about hybrid architecture for learning to detect retinal vessel in fundus images

### Answer

Our proposed approach for learning to detect retinal vessels in fundus images is a hybrid architecture that combines a deep neural network (DNN) with a random forest (RF). The DNN is used for unsupervised learning of vesselness dictionaries using sparse trained denoising auto-encoders (DAE), followed by supervised learning of the DNN response using a random forest for detecting vessels in color fundus images. The DNN is a stacked denoising autoencoder (DAE) with three hidden layers. The first hidden layer has 1000 neurons, the second hidden layer has 500 neurons, and the third hidden layer has 250 neurons. The output layer has one

### Question

who is pabitra mitra give answer in 200 words

### Answer

Pabitra Mitra is an author who has contributed to the research paper titled "Ensemble of Deep Convolutional Neural Networks for Learning to Detect Retinal Vessels in Fundus Images." This paper focuses on developing a computational imaging framework using deep and ensemble learning for reliable detection of blood vessels in fundus color images. He is affiliated with the Indian Institute of Technology Kharagpur, Kharagpur, WB 721302, India. Thanks for asking!



Thank You

