Open in app ↗

✦ Member-only story

# Building LLM Applications: Open-Source RAG (Part 7)

Vipra Singh · Follow

27 min read · Mar 16, 2024

👏 223        💬 1                              🔖        ▶        ⬆️        •••

*Learn Large Language Models (LLM) through the lens of a Retrieval Augmented Generation (RAG) Application.*

## Posts in this Series

9. *Serving LLMs*

10. *Advanced RAG*

# Table of Contents

# vishnu, read the best stories from industry leaders on Medium.

The author made this story available to Medium members only. Upgrade to instantly unlock this story plus other member-only benefits.

✦  Access all member-only stories on Medium

✦  Become an expert in your areas of interest

✦  Get in-depth answers to thousands of questions about technical

✦  Grow your career or build a new one

**Marc-André Giroux**
Sr. Software Developer
Netflix

**Carlos Arguelles**
Sr. Staff Engineer
Google

**Tony Yiu**
Director
Nasdaq

**Brandeis Marshall**
CEO
DataedX

**Cassie Kozyrkov**
Chief Decision Scientist
Google

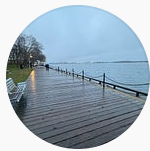**Memo Akten**
Asst. Professor
UCSD

**Vitali Zaidman**
Software Architect
Meta
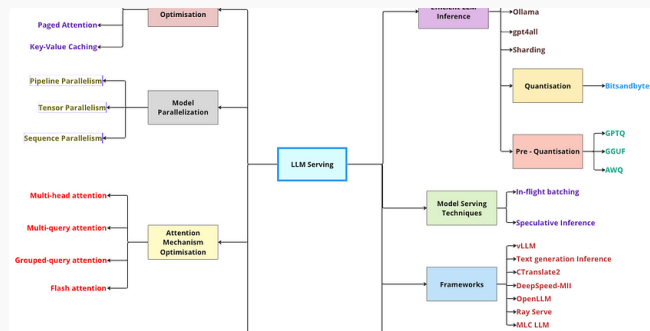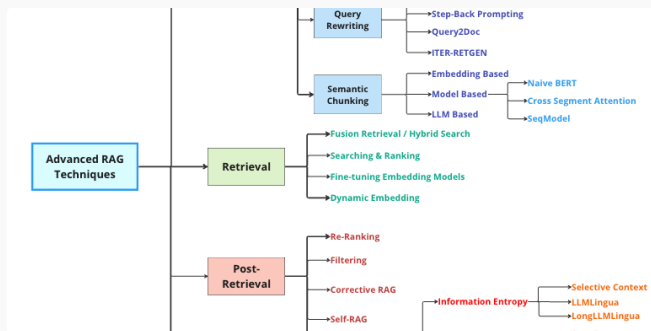
**Camille Fournier**
Head of Engineering
JPMorgan Chase

Upgrade

# Written by Vipra Singh

Follow

## 1K Followers

## More from Vipra Singh



Vipra Singh

### Building LLM Applications: Advanced RAG (Part 10)

Learn Large Language Models ( LLM ) through the lens of a Retrieval Augmented...

✦ · 48 min read · Apr 27, 2024

👏 384    💬 2                    🔖    •••



Vipra Singh

### Building LLM Applications: Serving LLMs (Part 9)

Learn Large Language Models ( LLM ) through the lens of a Retrieval Augmented...
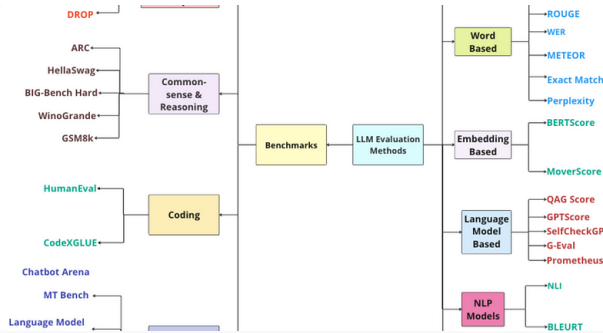
✦ · 50 min read · Apr 17, 2024

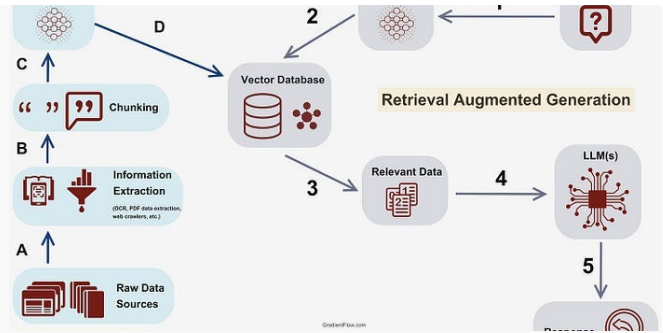👏 546    💬 3                    🔖    •••

Vipra Singh

## Building LLM Applications: Evaluation (Part 8)

Learn Large Language Models ( LLM ) through the lens of a Retrieval Augmented...

✦ · 48 min read · Apr 7, 2024

👏 276    💬 1                    🔖    ⋯



Vipra Singh

## Building LLM Applications: Introduction (Part 1)

Learn Large Language Models ( LLM ) through the lens of a Retrieval Augmented...
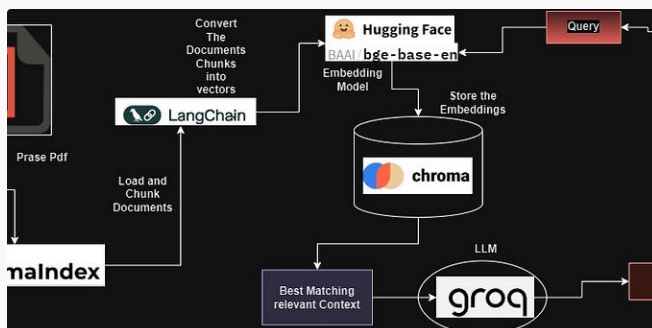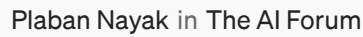
5 min read · Jan 8, 2024
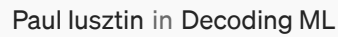
👏 905    💬 3                    🔖    ⋯

( See all from Vipra Singh )

# Recommended from Medium

Plaban Nayak in The AI Forum                Paul Iusztin in Decoding ML

## RAG on Complex PDF using LlamaParse, Langchain and Groq

## The 4 Advanced RAG Algorithms You Must Know to Implement

Retrieval-Augmented Generation (RAG) is a new approach that leverages Large Languag…

Implement from scratch 4 advanced RAG methods to optimize your retrieval and post-…

13 min read · Apr 7, 2024                     15 min read · May 4, 2024

678      9                                1.4K      10

## Lists

### Natural Language Processing
1494 stories · 1011 saves

### AI Regulation
6 stories · 473 saves

### ChatGPT prompts
47 stories · 1642 saves

### Predictive Modeling w/ Python
20 stories · 1254 saves





Barsha Rani Swain in GoPenAI            Dr. Leon Eversberg in Towards Data Science

## Advanced RAG: Corrective Retrieval Augmented Generation…

## How to Build a Local Open-Source LLM Chatbot With RAG

CRAG enhances the traditional RAG by introducing a retrieval evaluator to assess th…

Talking to PDF documents with Google's Gemma-2b-it, LangChain, and Streamlit

10 min read · Apr 23, 2024                ⭐ · 12 min read · Mar 31, 2024

👤 Fareed Khan in Level Up Coding          👤 Vishal Rajput 🔷 in AIGuys

## Building LLaMA 3 From Scratch with Python
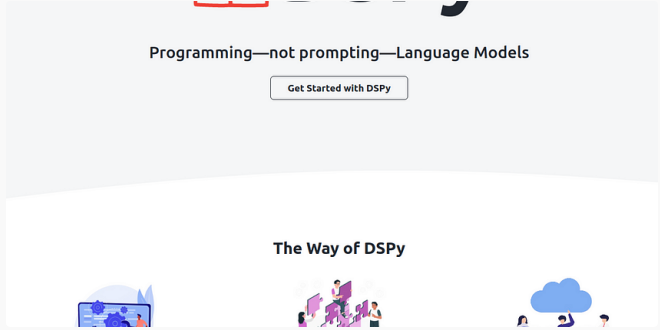
Code Your Own Billion Parameter LLM

29 min read  ·  May 28, 2024

## Prompt Engineering Is Dead: DSPy Is New Paradigm For Prompting

DSPy Paradigm: Let's program — not prompt — LLMs

✦  ·  11 min read  ·  May 29, 2024

( See more recommendations )