

DeepSeek

DeepSeek-V3 is a Mixture-of-Experts (MoE) language model with 671 billion total parameters, of which only 37 billion are active per token. It improves training efficiency and inference speed while maintaining state-of-the-art performance in natural language understanding, coding, and mathematical reasoning, including:

DeepSeek LLM (General-purpose language model)

DeepSeek Coder (AI for code generation)

DeepSeek Vision (Upcoming multimodal model)

Key Innovations:

- Multi-Head Latent Attention (MLA) for memory-efficient attention.
- DeepSeekMoE with auxiliary-loss-free load balancing.
- Multi-Token Prediction (MTP) for faster inference and training.
- FP8 mixed-precision training for efficient computation.
- Optimized cross-node communication for large-scale distributed training.

1. Basic Architecture

Transformer-Based Model with Mixture-of-Experts

DeepSeek-V3 follows the Transformer architecture, but instead of a dense model, it uses a Mixture-of-Experts approach, where only a subset of parameters are active at any time. This enables training large-scale models with lower computational costs.

- Total Parameters: 671 billion
- Active Parameters per Token: 37 billion
- Mixture-of-Experts (MoE) layers distribute computations across multiple expert subnetworks.

Multi-Head Latent Attention

MLA is a memory-efficient attention mechanism that reduces the Key-Value cache size without losing accuracy.

- Compresses attention keys and values to save memory.
- Uses Rotary Positional Embedding (RoPE) for long-context learning.
- Maintains performance similar to Multi-Head Attention (MHA) while reducing computational cost.
- MLA significantly reduces memory usage during inference, making DeepSeek-V3 faster and more efficient.

DeepSeekMoE: Optimized Mixture-of-Experts Training

DeepSeek-V3 uses DeepSeekMoE, an advanced MoE technique that improves load balancing and reduces training inefficiencies.

- Shared & Routed Experts: Some layers are shared across all tokens, while others are routed based on token characteristics.
- Auxiliary-Loss-Free Strategy: Unlike previous MoE models that use additional loss functions for balancing experts, DeepSeek-V3 dynamically adjusts expert loads without introducing performance penalties.
- Node-Limited Routing: Limits communication costs by restricting expert allocation across compute nodes.
- MoE allows DeepSeek-V3 to scale up massively while keeping training and inference computationally feasible.

Multi-Token Prediction

DeepSeek-V3 introduces Multi-Token Prediction, a novel approach where the model predicts multiple future tokens in parallel, instead of just the next token.

- Speeds up training by generating multiple tokens at once.
- Improves inference by enabling speculative decoding.
- Maintains complete causal relationships across token predictions.
- MTP reduces training costs and allows DeepSeek-V3 to generate text more quickly than traditional autoregressive models.

2. Infrastructure

DeepSeek-V3 was trained on a highly optimized computing infrastructure with NVIDIA H800 GPUs.

Compute Cluster

- 2048 NVIDIA H800 GPUs used for training.
- 8 GPUs per node connected via NVLink (fast intra-node communication).
- InfiniBand interconnects for communication between nodes.
- This setup enables high-speed distributed training while minimizing communication overhead.

Training Framework

- DeepSeek-V3 was trained using **HAI-LLM**, a custom training framework optimized for large-scale AI models.
 - Key Optimization
 - DualPipe Parallelism – A novel pipeline parallelism method that overlaps computation and communication, improving efficiency.
 - Efficient Cross-Node Communication – Uses custom all-to-all communication kernels that maximize InfiniBand and NVLink bandwidth.
 - Memory Optimization – Recomputed normalization & low-precision storage reduce memory usage during training.
 - DeepSeek-V3 can train on extremely large datasets while keeping computational costs low.

3. Pre-Training

Dataset & Training Objective

- Trained on 14.8 trillion diverse high-quality tokens.
- Supports long-context training up to 128K tokens.
- Uses FP8 mixed precision for more efficient training.

FP8 Mixed Precision Training

DeepSeek-V3 is one of the first large-scale models to use FP8 training, which reduces memory usage and improves computational efficiency.

- Reduces memory consumption without loss in accuracy.
- Speeds up training by reducing precision bottlenecks.
- FP8 training allows DeepSeek-V3 to be trained using fewer resources while maintaining high accuracy.

4. Post-Training

Supervised Fine-Tuning (SFT)

- Trained on human-annotated datasets.
- Improves instruction-following ability for chatbot applications.

Reinforcement Learning with Human Feedback (RLHF-like)

- Uses a Reward Model and Group Relative Policy Optimization (GRPO).
- Fine-tunes the model to align with human preferences.

Feature	DeepSeek-V3	GPT-4	LLaMA 3	Mistral
Total Parameters	671B	~1T	405B	7B, 13B
Active Parameters per Token	37B	~200B	405B	7B, 13B
Open-Source	Yes	No	Yes	Yes
Efficient MoE	Yes	No	No	No
Pre-Training Tokens	14.8T	Unknown	Unknown	Unknown

DeepSeek is a Chinese AI company specializing in foundation models for natural language processing (NLP), code generation, and multimodal AI. Their main goal is to create highly capable AI models. Including

DeepSeek LLM

DeepSeek LLM is a transformer-based AI model designed for various NLP tasks, similar to GPT-4, LLaMA, and Mistral. It can be used for:

- **Text generation**
- **Summarization**
- **Translation**
- **Question answering**
- **Sentiment analysis**
- **Conversational AI**

Key Features of DeepSeek LLM

- Large-scale dataset: Trained on a massive dataset covering web data, books, research papers, and code repositories.
- Multilingual support: Works across multiple languages, making it useful for global applications.
- Open-source availability: Some versions are freely available for research and development.

- Improved factual accuracy: Designed to reduce hallucinations and improve reliability.
-
- ☐ Available on **Hugging Face** and **GitHub** for **free use and fine-tuning**.
 - ☐ API access for **developers and businesses**.

DeepSeek Coder

DeepSeek Coder is an AI coding assistant similar to GitHub Copilot and Code Llama. It helps developers by suggesting, completing, and debugging code in real-time.

Key Features of DeepSeek Coder

- Supports multiple programming languages: Python, Java, JavaScript, C++, ABAP, and more.
 - Code completion & generation: Suggests full functions, classes, and even entire scripts.
 - Intelligent debugging: Detects and fixes code errors.
 - Optimized for software engineering: Can handle complex programming patterns.
-
- ☐ Open-source models are available on **Hugging Face**.
 - ☐ Can be **integrated into IDEs** like VS Code.

DeepSeek Vision (Upcoming Model)

DeepSeek Vision is expected to be a multimodal AI model capable of processing text + images. This would make it an alternative to GPT-4V and Gemini.

Expected Features of DeepSeek Vision

- Image recognition
- Visual question answering (VQA)
- Image captioning
- Text + image reasoning

Availability & Access

DeepSeek provides open-source models as well as API-based access for businesses.

MODEL	AVAILABILITY	USAGE
DeepSeek LLM	Open-Source	Text-based AI, chatbot development, NLP tasks
DeepSeek Coder	Open-Source	Code generation, debugging, automation
DeepSeek Vision	Not Released Yet	Multimodal AI (Text + Image)
API Access	Available	Developers & businesses

- Hugging Face: Free access to model weights.
- GitHub: Open-source versions available.

- API Access: For enterprises and large-scale applications.

DeepSeek vs. Other AI Models

DeepSeek LLM vs. GPT-4, LLaMA, Mistral

Feature	DeepSeek LLM	GPT-4	LLaMA 2	Mistral
Open Source	(Partially)	(Closed)	(Open)	(Open)
Multimodal(Text + Image)	(Not Yet)	Yes	No	No
Performance	Fast	High	Moderate	Moderate
Fine-Tuning	Yes	No	Yes	Yes

DeepSeek Coder vs. Code Llama, GitHub Copilot

Feature	DeepSeek Coder	Code Llama	GitHub Copilot
Open-source	Yes	Yes	No
IDE Integration	Yes	Yes	Yes
Free Usage	Yes	Yes	No(paid)

Key Differences

- DeepSeek is partially open-source, while OpenAI's GPT-4 is closed.

- DeepSeek is optimized for coding, whereas GPT-4 is a general AI model.
- DeepSeek provides free access, while GitHub Copilot and OpenAI APIs are paid.

Usage Of DeepSeek

For NLP Applications

AI chatbots

summarization tools

Sentiment analysis

For Coding

AI-powered code assistant (like Copilot)

Debugging and error detection

Code completion for developers

For Research

Experimenting with LLM architectures

Fine-tuning AI models for specific tasks

Creating RAG Agent using DeepSeek -R1 and Ollama

Document Assisstant

In this I create end to end RAG application with the help of DeepSeek. This entire DeepSeek will be specifically installed in our local with the help of Ollama. In this I use Ollama embedding and show that How the vectors are created and how can we store that vector even in our local system

Accuracy

It's quite good when I chat with the RAG application. Performance wise also it's working really good.

First I download Ollama distilled Model and installed in my PC in My Command Prompt

This is the model I installed In my PC

```
ollama run deepseek-r1:1.5b
```

Models

DeepSeek-R1

```
ollama run deepseek-r1:671b
```

Distilled models

DeepSeek team has demonstrated that the reasoning patterns of larger models can be distilled into smaller models, resulting in better performance compared to the reasoning patterns discovered through RL on small models.

Below are the models created via fine-tuning against several dense models widely used in the research community using reasoning data generated by DeepSeek-R1. The evaluation results demonstrate that the distilled smaller dense models perform exceptionally well on benchmarks.

DeepSeek-R1-Distill-Qwen-1.5B

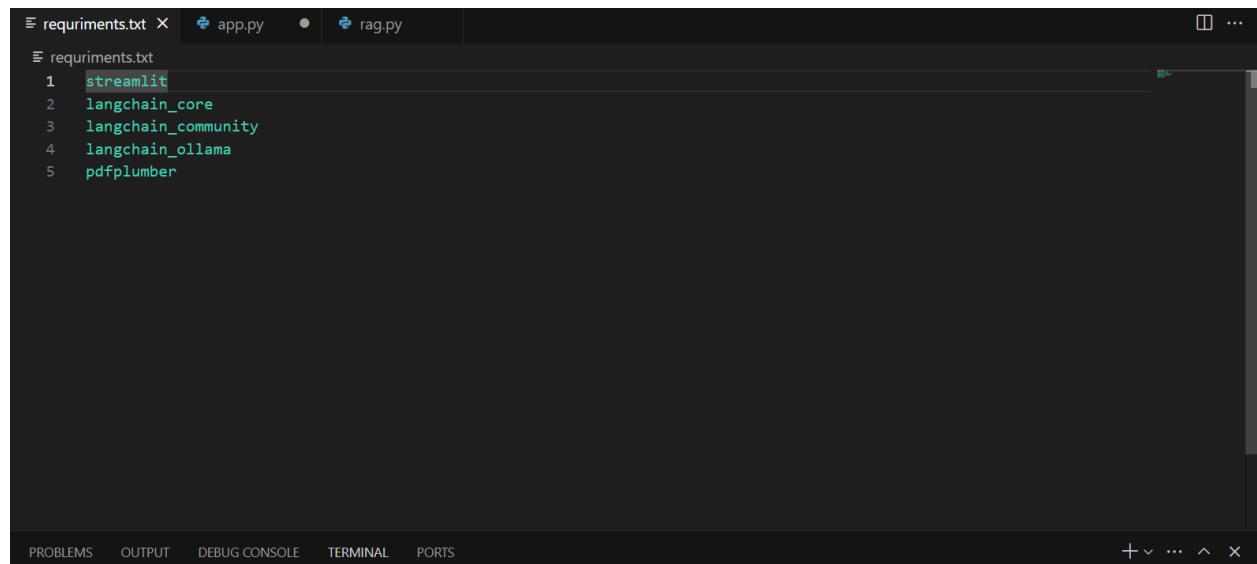
```
ollama run deepseek-r1:1.5b
```

DeepSeek-R1-Distill-Qwen-7B

```
ollama run deepseek-r1:7b
```

I created my RAG application in vscode and I push This in Github

These are the requirements first we pip installed before working on code

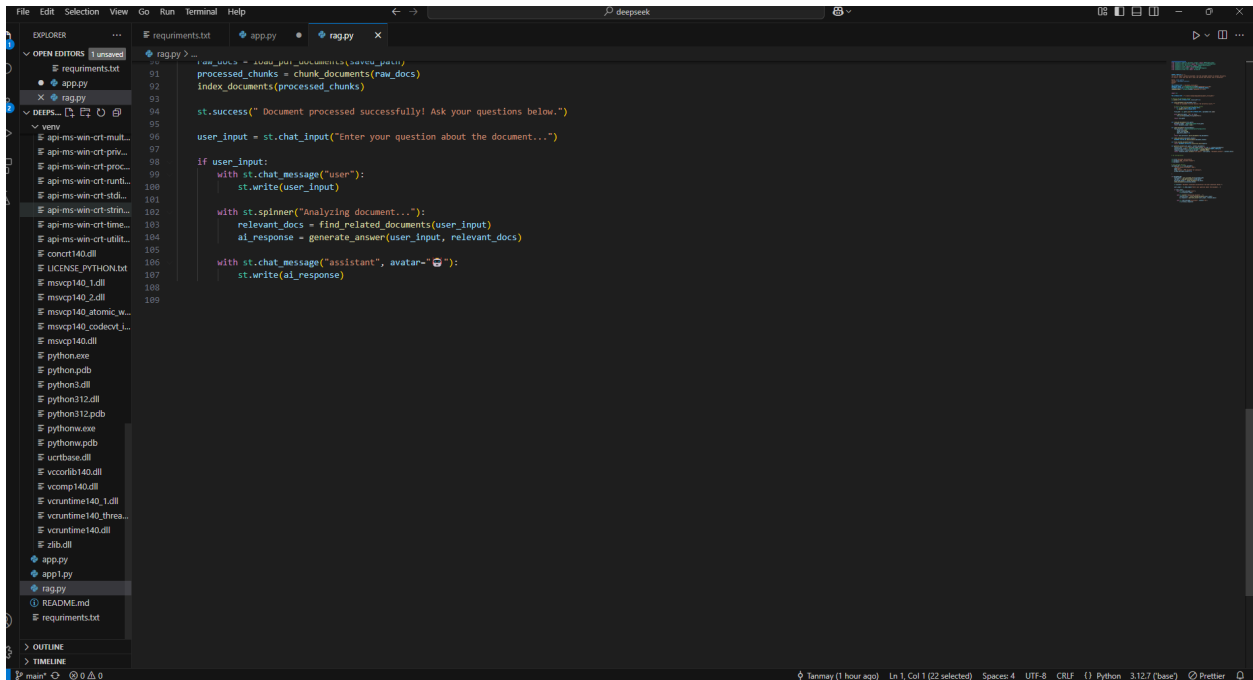
A screenshot of a code editor interface with a dark theme. The top bar shows three tabs: 'requirements.txt' (active), 'app.py', and 'rag.py'. The 'requirements.txt' tab is selected, and the file content is displayed in a large text area. The content consists of five lines of text, each preceded by a line number (1-5). The text is: 'streamlit', 'langchain_core', 'langchain_community', 'langchain_ollama', and 'pdfplumber'. The bottom of the editor has a panel with tabs for 'PROBLEMS', 'OUTPUT', 'DEBUG CONSOLE', 'TERMINAL' (which is active and underlined), and 'PORTS'. On the right side of the bottom panel, there are icons for window management: a plus sign, a downward arrow, a three-dot menu, an upward arrow, and a close 'X' button.

```
1 streamlit
2 langchain_core
3 langchain_community
4 langchain_ollama
5 pdfplumber
```

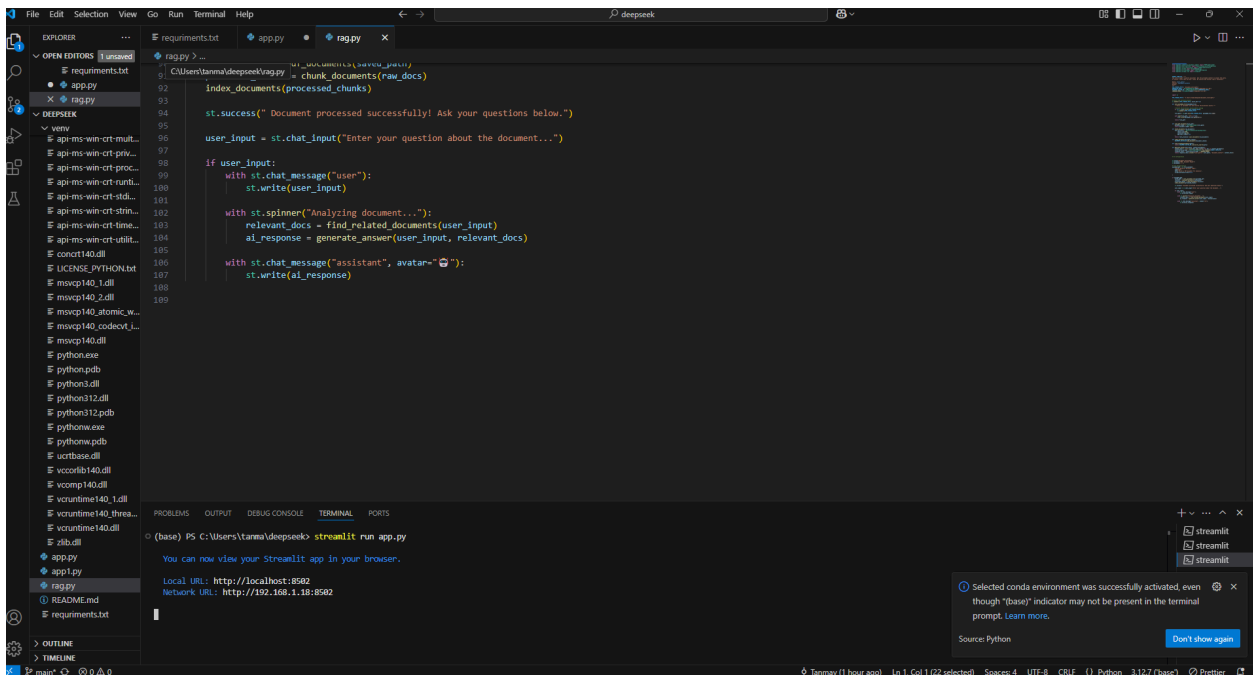
Then I created a python file name rag.py where I start doing My code this is my code

```
File Edit Selection View Go Run Terminal Help
requirements.txt app.py raggy x
EXPLORER
OPEN EDITORS
requirements.txt
app.py
raggy
DEEPSPEK
view
api-ms-win-crt-mult...
api-ms-win-crt-priv...
api-ms-win-crt-proc...
api-ms-win-crt-runli...
api-ms-win-crt-stdl...
api-ms-win-crt-strn...
api-ms-win-crt-time...
api-ms-win-crt-util...
conor140.dll
LICENSE_PYTHON.txt
msvcp140_1.dll
msvcp140_2.dll
msvcp140_atomic_w...
msvcp140_codecvt_l...
msvcp140.dll
python.exe
python.pdb
python3.dll
python312.pdb
python312.pdb
python.exe
python.pdb
ucrtbase.dll
vccorlib140.dll
vcomp140.dll
vcruntime140_1.dll
vcruntime140_thres...
vcruntime140.dll
zlib.dll
app.py
app.py
raggy
README.md
requirements.txt
OUTLINE
TIMELINE
main: C:\ 0 0 0 0
raggy.py
1 import streamlit as st
2 from langchain_community.document_loaders import PDFLoader
3 from langchain_text_splitters import RecursiveCharacterTextSplitter
4 from langchain_core.vectorstores import InMemoryVectorStore
5 from langchain_ollama import OllamaEmbeddings
6 from langchain_core.prompts import ChatPromptTemplate
7 from langchain_ollama.llms import OllamaLLM
11 PROMPT_TEMPLATE = """
12 You are an expert research assistant. Use the provided context to answer the query.
13 If unsure, state that you don't know. Be concise and factual (max 3 sentences).
14 Query: {user_query}
15 Context: {document_context}
16 Answer:
17 """
18 PDF_STORAGE_PATH = "document_store/pdfs/"
19 EMBEDDING_MODEL = OllamaEmbeddings(model="deepseek-r1:1.5b")
20 DOCUMENT_VECTOR_STORE = InMemoryVectorStore(EMBEDDING_MODEL)
21 LANGUAGE_MODEL = OllamaLLM(model="deepseek-r1:1.5b")
22
23 import os
24
25 PDF_STORAGE_PATH = "C:/Users/tanna/deepseek/document_store/pdfs/"
26
27 # Ensure the directory exists
28 os.makedirs(PDF_STORAGE_PATH, exist_ok=True)
29
30 def save_uploaded_file(uploaded_file):
31     """Saves an uploaded file and ensures the directory exists."""
32     # Ensure the directory exists before saving
33     if not os.path.exists(PDF_STORAGE_PATH):
34         os.makedirs(PDF_STORAGE_PATH)
35
36     file_path = os.path.join(PDF_STORAGE_PATH, uploaded_file.name)
37
38     with open(file_path, "wb") as file:
39         file.write(uploaded_file.getbuffer())
40
41     return file_path
42
43 def load_pdf_documents(file_path):
44     document_loader = PDFLoader(file_path)
45     raw_documents = document_loader.load()
46     text_processor = RecursiveCharacterTextSplitter(
47         chunk_size=1000,
48         chunk_overlap=200,
49         add_start_index=True
50     )
51     return text_processor.split_documents(raw_documents)
52
53 def index_documents(document_chunks):
54     DOCUMENT_VECTOR_STORE.add_documents(document_chunks)
55
56 def find_related_documents(query):
57     return DOCUMENT_VECTOR_STORE.similarity_search(query)
58
59 def generate_answer(user_query, context_documents):
60     context_text = "\n\n".join([doc.page_content for doc in context_documents])
61     conversation_prompt = ChatPromptTemplate.from_template(PROMPT_TEMPLATE)
62     response_chain = conversation_prompt | LANGUAGE_MODEL
63     return response_chain.invoke({"user_query": user_query, "document_context": context_text})
64
65 # UI Configuration
66 st.title("Document Assistant")
67 st.markdown("### Document Reader")
68 st.markdown("-----")
69
70 # File upload section
71 uploaded_pdf = st.file_uploader(
72     "Upload Research Document (PDF)",
73     type="pdf",
74     help="Select a PDF document for analysis",
75     accept_multiple_files=False
76 )
77
78 if uploaded_pdf:
79     saved_path = save_uploaded_file(uploaded_pdf)
80     raw_docs = load_pdf_documents(saved_path)
81     processed_chunks = chunk_documents(raw_docs)
82     index_documents(processed_chunks)
```

```
File Edit Selection View Go Run Terminal Help
requirements.txt app.py raggy x
EXPLORER
OPEN EDITORS
requirements.txt
app.py
raggy
DEEPSPEK
view
api-ms-win-crt-mult...
api-ms-win-crt-priv...
api-ms-win-crt-proc...
api-ms-win-crt-runli...
api-ms-win-crt-stdl...
api-ms-win-crt-strn...
api-ms-win-crt-time...
api-ms-win-crt-util...
conor140.dll
LICENSE_PYTHON.txt
msvcp140_1.dll
msvcp140_2.dll
msvcp140_atomic_w...
msvcp140_codecvt_l...
msvcp140.dll
python.exe
python.pdb
python3.dll
python312.pdb
python312.pdb
python.exe
python.pdb
ucrtbase.dll
vccorlib140.dll
vcomp140.dll
vcruntime140_1.dll
vcruntime140_thres...
vcruntime140.dll
zlib.dll
app.py
app.py
raggy
README.md
requirements.txt
OUTLINE
TIMELINE
main: C:\ 0 0 0 0
raggy.py
46
47 def load_pdf_documents(file_path):
48     document_loader = PDFLoader(file_path)
49     return document_loader.load()
50
51 def chunk_documents(raw_documents):
52     text_processor = RecursiveCharacterTextSplitter(
53         chunk_size=1000,
54         chunk_overlap=200,
55         add_start_index=True
56     )
57     return text_processor.split_documents(raw_documents)
58
59 def index_documents(document_chunks):
60     DOCUMENT_VECTOR_STORE.add_documents(document_chunks)
61
62 def find_related_documents(query):
63     return DOCUMENT_VECTOR_STORE.similarity_search(query)
64
65 def generate_answer(user_query, context_documents):
66     context_text = "\n\n".join([doc.page_content for doc in context_documents])
67     conversation_prompt = ChatPromptTemplate.from_template(PROMPT_TEMPLATE)
68     response_chain = conversation_prompt | LANGUAGE_MODEL
69     return response_chain.invoke({"user_query": user_query, "document_context": context_text})
70
71 # UI Configuration
72 st.title("Document Assistant")
73 st.markdown("### Document Reader")
74 st.markdown("-----")
75
76 # File upload section
77 uploaded_pdf = st.file_uploader(
78     "Upload Research Document (PDF)",
79     type="pdf",
80     help="Select a PDF document for analysis",
81     accept_multiple_files=False
82 )
83
84 if uploaded_pdf:
85     saved_path = save_uploaded_file(uploaded_pdf)
86     raw_docs = load_pdf_documents(saved_path)
87     processed_chunks = chunk_documents(raw_docs)
88     index_documents(processed_chunks)
```



Then I run the code and in terminal I run streamlit run rag.py
Streamlit app will open in browser



Streamlit app interface after uploading a file

Document Assisitent

Document Reader

Upload Research Document (PDF)

Drag and drop file here
Limit 200MB per file • PDF

Browse files

aily scrum.pdf 82.0KB

Document processed successfully! Ask your questions below.

tell me what we have done in 18th jan

<think> Okay, I need to figure out what was done on January 18th based on the context provided. The query asks specifically about "tell me what we have done in 18th jan," so I should focus on events that happened around that date. Looking at the context, there's a daily scrum entry from January 18, 2025, at 08:05 AM. Inside that entry, there are some details about updates and tasks. The task mentioned is analyzing Kaggle dataset data to predict if it's real or fake.

Enter your question about the document...

After uploading a file the question bar appeared where you can ask anything about the PDF. It will exactly answer like the way DeepSeek answer.

Upload Research Document (PDF)

Drag and drop file here
Limit 200MB per file • PDF

Browse files

aily scrum.pdf 82.0KB

Document processed successfully! Ask your questions below.

tell me what we have done in 18th jan

<think> Okay, I need to figure out what was done on January 18th based on the context provided. The query asks specifically about "tell me what we have done in 18th jan," so I should focus on events that happened around that date. Looking at the context, there's a daily scrum entry from January 18, 2025, at 08:05 AM. Inside that entry, there are some details about updates and tasks. The task mentioned is analyzing Kaggle dataset data to predict if it's real or fake.

There are two steps listed under NEXT STEPS:

1. Research attributes (function and employment type).
2. Plan preprocessing strategies for the dataset.

So, on January 18th, we performed an initial analysis of a Kaggle dataset related to job postings. The tasks involved understanding specific attributes and planning how to preprocess them. There are no other details mentioned about that day beyond this.

Enter your question about the document...

The best part is that it will work locally, all information will be stored locally no one will access your file and information, not even China can access this only you can access this.