

The Rig Veda API's

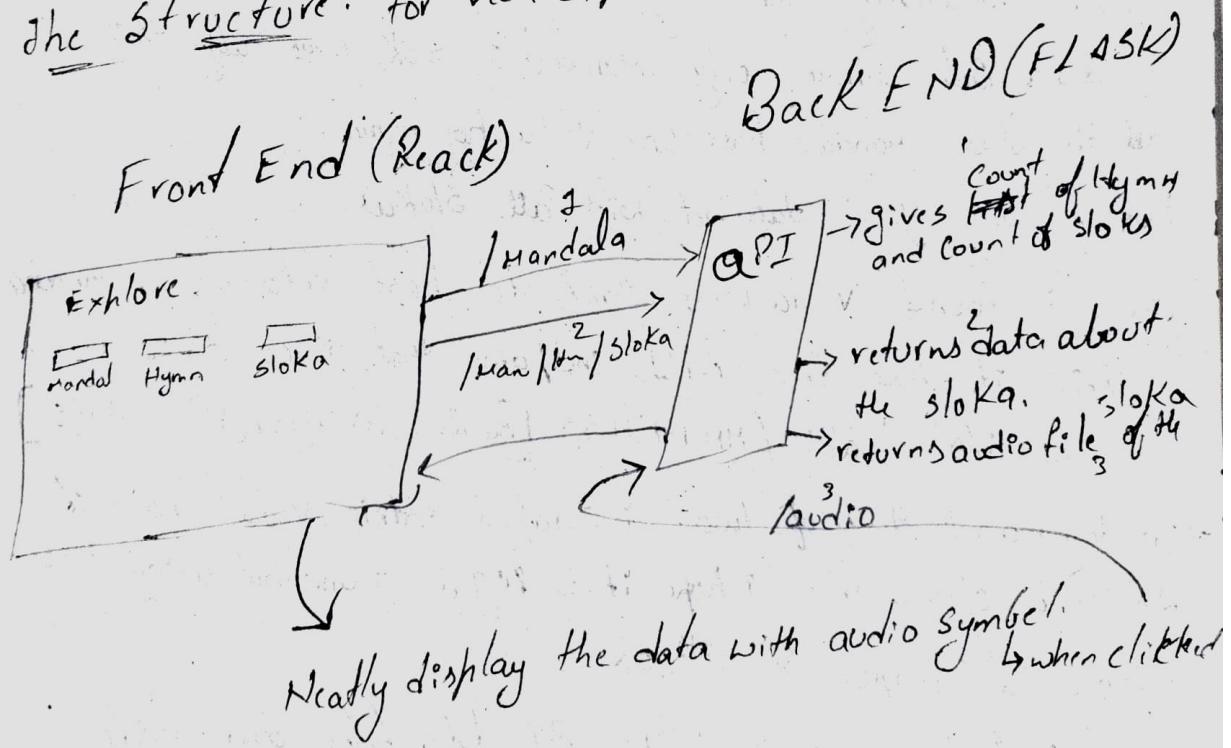
Initial idea was to get the dataset and make an explorer that anyone can explore the Rig Veda and get any sloka with audio and translation.

The Development plan is being made, but then I got an idea on adding chat GPT like bot to it, for only answering veda questions.

I got idea like training LLM with the dataset, but I don't have resources to deploy it, so, I got an idea like

- Based on user question, get relevant sloka
- feed it to get natural answers from LLM like GPT 3.5/ernie like
- Display sloka and answer

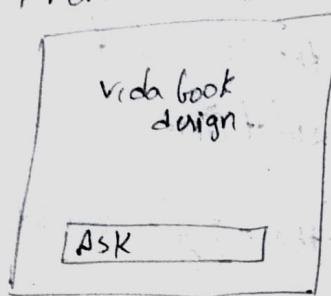
The Structure for Veda Explorer. (might change through development)



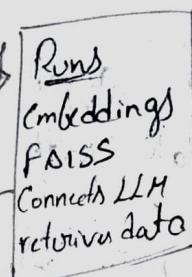
This used JSON file as RigVeda dataset not MongoDB
As it requires paid sub for high traffic

Structure for chat Goti

Frontend (React)



Backend (Flask)



1 Gives the User Query

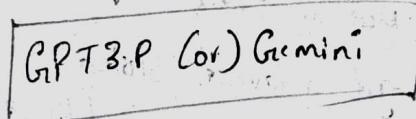
2 Dynamically gives the answer with slokas.

3 return natural answers
for the query with
slokas

4 Runs
embeddings
FAISS
Connects LLM
returns data

1 uses Prebuilt
embedding and retrieving
relevant sloka.

2 gives the query and returned Sloka.



Development until now:

- While searching for Rigveda dataset mostly i found a data set where only mandal 1 slokas were fully there all the other mandals have one sloka per hymn.
But i want good dataset with all slokas.
- Then i found vedaweb, thanks to host reference at announcement
→ their it has free API endpoint that gives sloka and its full details "api / MM HHHSSS [Mandala(M), Hymn(H), Sloka(S)]"
- On day 1 and Half of two, i scraped entire veda from it. As it is open source, i hope it is legal, i am not using it for commercial purpose so it's OK.
- Then i parsed all the RAW API data files and made initial per Mandala.json file and merged them for final dataset.
- Later i got audio files dataset from RV-audio-Rico from github, but it has all slokas combined into one MP3 per hymn.

→ I designed a formula to split the Entire Hymn .MP3 file into single per sloka MP3 files using slokas count and length of each sloka, i am Middly satisfied with the result but it is good for now.

→ Then on evening of Day 2 i designed and implemented React + flask communication.

→ At first i used MongoDB to store them dataset but it is more than 16MB So i split it for Mandal

→ Later when i search sloka it became difficult and traffic increasing so i made it local JSON files and backend directly access it.

Data set structure (Per Mandal)

```
{ "Mandala": 1,  
  "hymns": [  
    { "hymn-number": 1,  
      "stanzas": [  
        { "stanza-number": 1,  
          "location": "01-001-01",  
          "Sankrit": "...",  
          "Transliteration": "...",  
          ".... - scheme": "IAST",  
          "translations": {  
            "griffith": "...",  
            "Macdonell": "...",  
            "oldenburg": "...",  
            "Padav": { ".... 33  
          }  
        }  
      ]  
    }  
  ]  
}
```

→ Finally at last of second day i got working model for "IVE".

Day 3

Today i focused on creating Embeddings for the dataset to create the chat. Got

→ All the Morning i spent Coding Embedding creation

using gemini-embedding v1 but it didn't reached limits before embedding all. (it took half day to realize)

→ So i switched to local model (Bge-base-en-v1.5)

→ So i switched to local model (Bge-base-en-v1.5)
it took 3 hours to create all embeddings for each sloka

→ Yes i embedded each sloka (10524) rather than combining all slokas in Hymn into one embedding.

→ Then i created merged all embeddings to create final FAISS index and Slokas-mapping.

→ Then i tested the ~~FAISS~~ index Semantic Search using Py Script.

→ I have given a search query

→ It encodes it with the Model(embedding model) → creates vector

→ Compare/Search the vector with FAISS index and give top 3 index's

, Returns Text using slokas-mapping (returns the slokas)

→ It performs very good, giving slokas relevant to the search term.

Journal continued / Documentation of development

→ Today : only focuses on chatbot creation (Backend)
only.

→ I changed Backend file structure to modular approach
⇒ from one file to two folders/modules

↳ chat-bot

↳ sloka-explorer

each has its routes that handles the API and all of them
integrated with APP.py and its routes.

* The chatbot Processing flow is as follows.

→ user gives the query to chatbot.

* The driver function (get-answer) will handle all the processing.

* first what the user wants will be extended using
extract-intents-gemini

* Based on the intent the below functions will perform.

↳ Fetch Sloka by location

↳ Semantic Search

* Answering Question:

all the results will be the context for LLM (gemini),
along with query and sloka will be retrieved and summarized.
answer will be returned by LLM and it will be returned
to frontend.

* Now After testing the Backend with Postman via Cloudflame API call ; move on to frontend.

→ Before that i moved all the data to Backend /data folder.

* Frontend :
Created frontend Prototype. It took me like a week.
Because i got sick for few days.

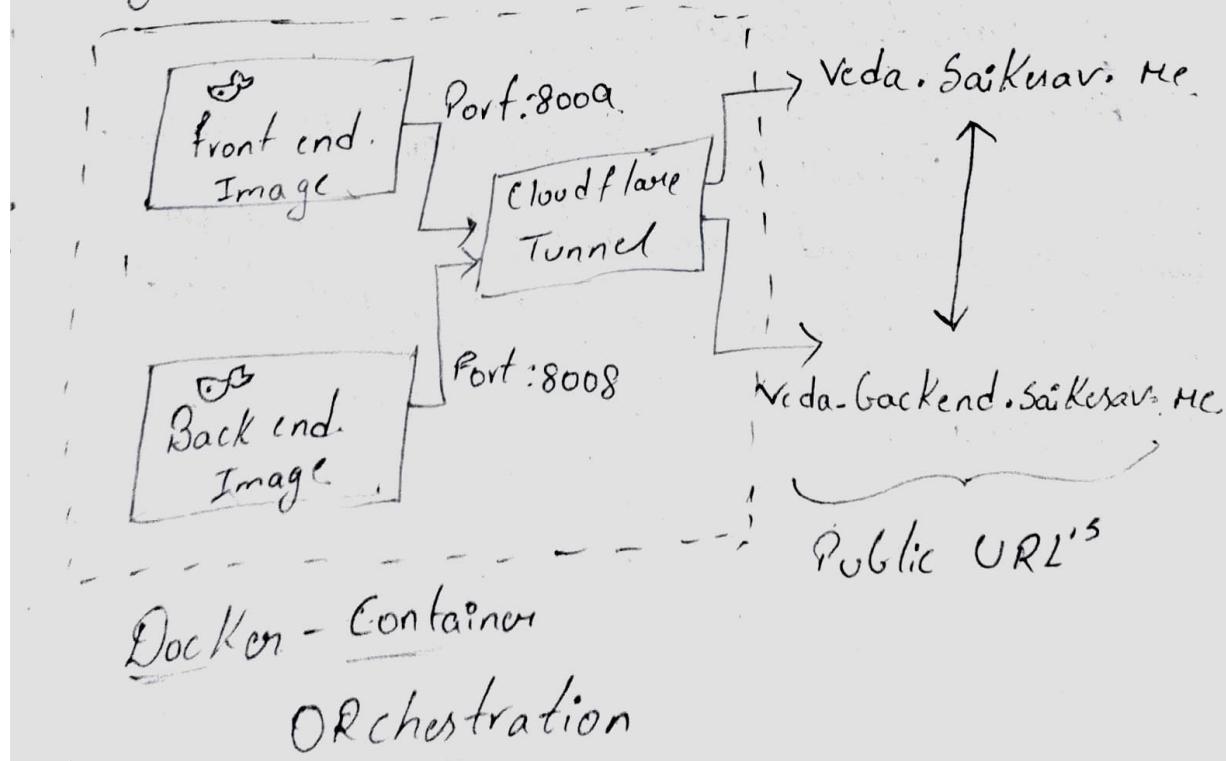
→ By 9th and 10th the initial version of the frontend was built. The design is very simple and easy to understand later ; will change / re-design it.

→ 11 and 12th and continues.

Tested Both components individually and integrated them. I used individual docker files.

→ After testing i want to deploy it in my server.

so i changed docker files and created this below deployment structure.



Deployment Method

I used Docker Orchestration technique to make a way of integrating multiple containers to work/communicate under single Network.

Approach:

I used Docker to containerized my front end & Backend.

As I already have code to make individual Image's now I tried to combine them into one image. But the overall size has become too much to be a single

image. This is a Problem (Centralized Image was 8GB+)

This is too much and creates a overload in the system and also it is not an efficient way of using Docker.

So I change the structure. I created two separate

Docker Images for frontend and Backend then I

↓ ↓
1.2GB 6+GB Including audio files

used a separate cloudflare Tunnel Image As shown in the side Drawing and created a Docker-compose file to common Bridge network.

Then ~~Deployed~~ the build both Images and uploaded to Docker Hub and pulled into the server and built the composed Image and hosted via the URL veda-saikumar.me

From now on I started using Git-version control as (v1.0) has developed and cleaned ~~for~~ without any sensitive data like APIs in them. So I used .gitignore & .dockerignore files to tranfer sensitive data from the code. If still you found any API keys don't use them please. Else you may

* Applying Threading to Backend

From 15 to 21

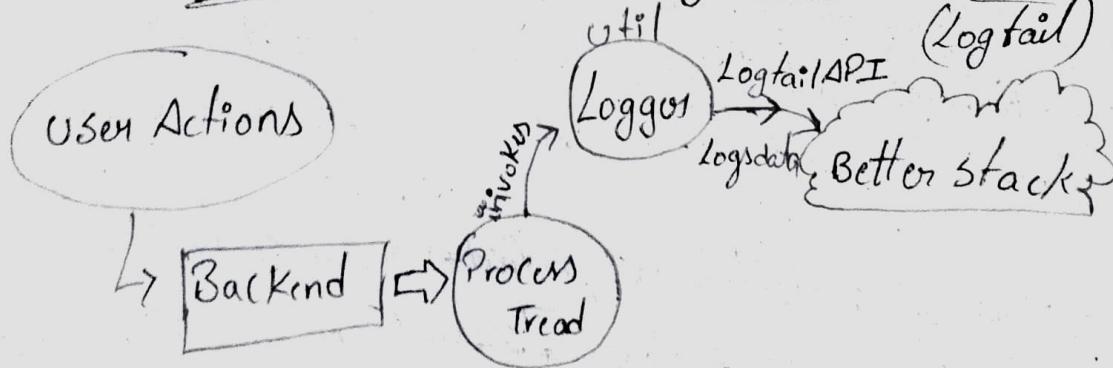
* Concurrent User Handling

- * As Prototype version is Ready. And also on 13-14 : implemented Session control and logging features.
After deployment : tested the website and found when multiple users logged in at once and used the AI the answers are mixing up and results in error from Backend.
→ So used Python Threading ^{and} to create multiple threads of the Backend instances and modularized resource sharing between instances thus achieving multi-thread performance.

when user access the api for like chatbot or something then Backend : separate thread will be created for that process and handle the request.

- * By the Threading and Resource Sharing is achieved.
- User Isolation
- Efficient Memory Usage.
- Resource sharing.

* Centralized cloud Logging with Better Stack*



After Deployment : want to know how Semantic Search, chat BOT responded to user, is the results are upto the user required. For that i developed cloud logging using Gitter Stack.

→ used logtail API token and created a source point

→ used logtail API token and created a source point
→ used logtail API token and created a source point
→ used logtail API token and created a source point

→ used logtail API token and created a source point
→ used logtail API token and created a source point

→ used logtail API token and created a source point
→ used logtail API token and created a source point

→ Anonymized IP address

→ Total time to handle the request.

→ User interaction data

→ System responses

for more indept details see `backend/utils/logging utility`.

After developing these features i deployed them

used Previously mention Orchestration Method.

For now initial version of the Project is done

From now on refinement and redesign/optimizing the features is Pending.

Phase 1 development Done.

Version 1.0

- * From now on every code change is recorded.
- * You can use Commit History to view the changes.

* Phase 2 : UI/UX Refinements (Oct 15-17, 2025)

- > Redesigned HomePage and added quotes
- > Updated veda Explore UI and added Progress Bar
- > Homepage animation is not up to mark so removed it, improved chatbot and Semantic Search UI.
- > Reduced LOC (Line of Code) in Frontend.

* Phase 3 : Feature Enchancement (Oct 17-18, 2025).

- > Now i noticed the audio (stotra audio) files it's size is 1GB, it is increasing Backend size. So i uploaded them int GitHub. (Reduce Backend size 8 → 7GB)
- > I add audio file retrieval ^{Module} from Github into My Backend and no frontend change as same ~~way~~ working.
- > Today while scrolling reels i saw a video when someone moving an object with hand gestures, then seeing that i got an idea why do i add an object in the webpage that user can interact with so i search google and found a 3D Model of OM. i put that homepage.
- > Then added new images in the webpage as icons.
- > Till now my Backend is about 6+GB and frontend is about 1GB.

→ The 6GB is because for Semantic Search model i am using "biggel-base-en-v1.5" it's about 480MB and it requires multiple variants of model to run in Docker Image in the next phase i want to optimize this and optimize the docker Image.

Phase 4 Production Readiness: (18-19, 2025)

- Today entire Semantic Search Model is changed old model alone takes upto 480x3 MB + other dependencies.
- To reduce that i used all-minilm-L6-v2 a 90MB Model.
- Smaller then base-en ~~but~~ and performs lower but i tested smaller then dataset using the Mini if perform upto the standard ; satisfied with the performance in testing so implemented it for entire dataset.
- First i created embeddings for each sloka using newer model and created new FAISS index and transferred the db index and embeddings to backend and changed backend code to use newer model. The performance is good. but semantic search formula giving errors so i changed it to adapt to the newer model.

* Reducing Backend Docker Image Size:

- After using lower Model for Semantic Search Docker image Docker Image size reduce 6 → 3.5GB, But still i want smaller size so i rewrote Dockerfile and reduced Image to 2.4GB. For more info see Commit History (dbfc68, 9964174)

* Frontend : Added Credits & Privacy Policy Pages.
if see those page for in details info about those website.

- Reduced 3D model texture file size. 6 MB → 2 MB so that every time page loads faster.
- changed Title from Rig Veda to Eternal Veda.

* Feature Completion

→ added navigate button for sloka Explorer to goto next sloka.

→ changed Frontend Color Palette.

for Backend : implemented API Authentication.

→ Till now my Backend is exposed to internet anyone can access LLM, Semantic model and entire dataset this is risky as my data and system resources can be breached.

so i implemented API Auth, i used Decorator and auth module to check every API request.

→ set up an API key in frontend and same in backend, on every request Backend checks API key is matching with its own key and returned the result(error) if no else.

gives access to the request resources.

For now i think this is the Final Phase, so i Merged both Backend & Frontend Git Head's into one and Published the Repo into GitHub.

→ Don't worry every API key in commit History has been stripped and it will not work :)

23 Oct 2025

- > Added Previous button in veda explorer just like next button.

Failover Protection using cloudflare + cloud deploy

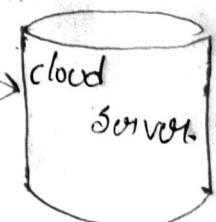
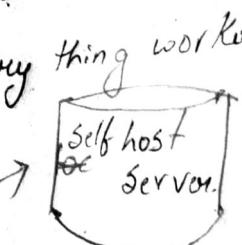
As I am using local server I had a serious risk of fail as it depends on regular household power & internet so now I want have a redundant way of deploying the website I want to use cloud service to act as backup server to the website.

-> I deployed frontend in vercel and till evening I tried deploying backend on services like rendev, fly etc but it didn't work the 512GB RAM of free tier is not enough for my ML + flask backend. I am looking for alternatives.

-> I handled user traffic between Primary (local server) & Backup (cloud deploy) using cloudflare workers.

-> I created a worker that checks if Primary is working then fetches data from there else fetches data from backup and writes appropriate console messages.

Deployed the worker and tested. Every thing works fine.



Traffic Handling cloudflare worker

24-26 Oct 2025

Improved audio files using New Extraction Method.

Previously I used audio length and text length attribute to create a formula to extract sloka audio from Hymn but, this method didn't give 100% accuracy. Some sloka's have non-uniform text length and audio length this cause merging and overlapping of audio's.

→ To solve this, first I tried using openAI's audio local model whisper - medium but Sanskrit is not detecting clearly and System Resources needed for model running is very high.

→ Then I used a new approach called "Silence-Based" approach.

→ Here I used Pydub's AudioSegment & detect_nonsilent method, to detect pauses between sloka's recitation → Then I took those segments and applied Sanskrit text to detect accurate audio length / segment's count.

• Based on the segments and text's metadata is extracted from sloka audio.

→ Identified all audio files and it's count with the text dataset to maintain accuracy.

→ Uploaded the audio to my github repo rig veda audio

Finally all the changes are made and deployed.

Thanks For Reading This Documentation.

The One Month Development Journey.

This development journey is a very special one in my life, it teached me new things and I learned a lot about Sanatan Dharma while building this website. Even if I not win, still I will not have any regret, because, I think this Hackathon is not about winning, it's about giving Eternal Knowledge to others. I am very thankfull to be apart of it. I also thank "Indra in Pixels by Ashritis" @indrainPixels for inspiring me to create this.

Still if you have any doubts or if you want to know more ^{about} this project or want to collaborate with me in future or want to be Friends - you can contact me via my website <https://www.saikesav.me>.

:)

Sathish