

# Multilingual Retrieval-Augmented Generation (RAG) System Report

## OCR to Create Text File:

- Used ocr\_pdf.py to extract text from hsc26\_bangla\_1st\_paper.pdf using pytesseract and pdf2image with Poppler.
- Saved the output as output.txt with a DPI of 300 for better quality.
- Faced formatting issues (e.g., broken sentences, mixed numbers like "109475৩") and fixed them with regex preprocessing.

## Initial Attempt with Gemini:

- Tried ChatGoogleGenerativeAI (Gemini) from langchain-google-genai with FAISS and HuggingFace embeddings.
- Encountered limitations, prompting a switch to another model.

## Switch to GPT with OpenAI API:

- I switched to ChatOpenAI (gpt-4o) using langchain\_openai.
- Installed openai and langchain-openai, then upgraded LangChain packages to resolve conflicts (e.g., langchain-core, langsmith).
- Fixed LLMChain validation errors by configuring ChatOpenAI with a temperature parameter.

## Troubleshooting and Iteration:

- Added a custom prompt to improve interpretive answers for the map\_reduce chain.
- Successfully answered factual queries like "কাকে অনুপমের ভাগ্য দেবতা বলে উল্লেখ করা হয়েছে?".

## Testing and Refinement:

- Tested Bangla queries, noting success with factual questions but struggles with interpretive ones (e.g., "অনুপমের ভাষায় সুপুরুষ কাকে বলা হয়েছে?").
- Explored solutions like better prompts and document context improvements.
- 

## English Queries:

### Bangla Queries

Query:

কাকে অনুপমের ভাগ্য দেবতা বলে উল্লেখ করা হয়েছে?

Answer:

অনুপম তার মামাকে ভাগ্য দেবতা এবং ভাগ্য দেবতার প্রধান এজেন্ট বলে উল্লেখ করেছে।

Groundedness Score: 0.65

Relevance Score: 0.78

Source Snippets:

- সেন এক জেড়া এয়ারিং এগিয়ে সেন সেকরার হাতে তা পরখ করে দেখার জন্য। কেননা সেটা ছিল অনুপমের মামার দেওয়া বিলাতি জিনিস, যাতে সেনার ভাগ আছে সামান্যই। গ. উদ্দীপকের সবিতা ও অপরিচিতা গল্পের কন্যাণী উভয়েই যৌ...
- চলেবে নিচের কোনটি সঠিক ক। খ। গ।।, 1 যা। 11, ৯৪। যরিশের বর্ণনায় মেয়ের বাবার পরিচয় সম্পর্কে জানা যায়... এককালে তাদের বংশে লক্ষ্মীর মঙ্গলঘট উপভূত করা ছিল দেশে বংশমর্যাদা রক্ষা করে চলা কঠিন বলে পশ...
- স্পষ্ট কথা বলার মতো সাহস তার নেই। নিজের সিদ্ধান্ত সে নিজে নিতে পারে না। সেকরা দিয়ে গহনা যাচাই যে কনেপকের অপমান তা অনুপম বুঝতে পারে না। এতে তার ব্যক্তিত্বহীনতার চরম প্রকাশ লক্ষ করা যায়। এসব বিক বিচ...
- আমি অন্নপূর্ণার কোলে গজাননের ছোটো ভাইটি। মতো তিনি আমাদের সমস্ত সংসারটাকে নিজের অন্তরের মধ্যে স্থাষিয়া লইয়াছেন। তাহাকে না খুঁড়িয়া এখানকার এক গণ্ডষও রস পাইবার জো নাই। এই কারণে কোনো-কিছুর জন্যই আমাকে ক...
- মধ্য দিয়ে গল্পের শেষাংশে কন্যাণীর স্তুতিস্তত্র আত্মপ্রকাশও ভবিষ্যতের নতুন নায়ের আগমনীর ইঙ্গিতে পরিসমাপ্ত। অপরিচিতা মনজাপে ভেঙেপড়া এক ব্যক্তিত্বহীন যুবকের স্বীকারোক্তির গল্প, তার পাপযালনের অকপট কথাম...
- মিল আছে ক হরিশের খ মামার গ শিক্ষকের ঘ বিনুর ৪। উক্ত চরিত্রে প্রাধান্য পেয়েছে - 1 দৌরাভ্য 1 যীনমনাতা 11 লোভ নিচের কোনটি ঠিক ক।। ও যা ও গ।।৬ ঘ।। 1 ও ৫ অনুপমের বয়স কত বছর ক পটিন্স খ ছাবিবিশ গ সা...
- বলাও তিনি আবশ্যক বোধ করিলেন না। কারণ, প্রমাণ হইয়া গেছে, আমি কেহই নই। তারপরে যা হইল সে আমি বলিতে ইচ্ছা করি না। বাড়লগঠন ভাঙিয়া-চুরিয়া, জিনিসপত্র লণ্ডভণ্ড করিয়া, ও কন্নট একসঙ্গে বাজিল না এবং অস্ত্র...
- হয়েছে ঠিকই কিন্তু গল্পের পরিণতিতে বৃত্তভাঙা ভিন্ন এক ব্যক্তি হিসেবে তাকে পাওয়া যায়। মন্তব্যটি যুক্তিসূক্ত। প্রতিষ্ঠিত হয় যখন সে পীমাবদ্ধতার গণ্ডি পেরিয়ে অসীমের সন্ধান পায়। তখন মানুষের চিত্ত হয়...
- কোন চরিত্রের ইঙ্গিতবহ ক হরিশ ঘ বিনুনা গ মামা ঘ শঙ্কনাথ উত্তর গ ৩৭। শাওনের কোন কোন বৈশিষ্ট্য অনুপমের মধ্যে থাকলে অনুপমের বিয়েটা ঠিকে যেত। সাহসিকতা ব্যক্তিষ্ট গভীর ভালোবাসা নিচের কোনটি সঠিক ক 7....
- সংকলন গল্পসম্প্রক-এ এবং পরে, গল্পগুচ্ছ তৃতীয় খণ্ডে ১৯২৭। অপরিচিতা গল্পে অপরিচিতা বিশেষণের আড়ালে যে বলিষ্ঠ ব্যক্তিত্বের অধিকারী নায়ীর কাহিনি বর্ণিত হয়েছে, তার নাম কন্যাণী। অমানবিক যৌতুক প্রচার নিম্নম...

[Ask another query](#)

Figure 1

Query:

অনুপমের বয়স কত বছর?

Answer:

অনুপমের বয়স ছাব্বিশ বছর।

Groundedness Score: 0.62

Relevance Score: 0.74

Source Snippets:

- কারণে গ গয়না নিয়ে মনোমালিন্যের কারণে ঘ কনের বাবার আত্মগরিমার কারণে ৭৫। আমার পুরোপুরি বয়সই হলো না কথাটি ঘুরা কী বোঝানো হয়েছে ক তরুণ বয়সী খ অপরিণত বয়সী গ অতি নির্ভরশীল ঘ চিত্রায় অপরিণত ৭৬। তামা...
- চলেবে নিচের কোনটি সঠিক ক। খ। গ।।, 1 যা। 11, ৯৪। যরিশের বর্ণনায় মেয়ের বাবার পরিচয় সম্পর্কে জানা যায়... এককালে তাদের বংশে লক্ষ্মীর মঙ্গলঘট উপভূত করা ছিল দেশে বংশমর্যাদা রক্ষা করে চলা কঠিন বলে পশ...
- সেনাপতি কার্তিকেশ্বর অপুর শোভা পায়। বড়ো হয়েও অনুপম কার্তিকের মতো মায়ের কাছাকাছি থেকে মাতৃআজ্ঞা পালনে ব্যস্ত থাকে। তাই পরিণত বয়সেও তার স্বাধীন ব্যক্তিত্বের বিকাশ ঘটে না। পাঠ্য গল্পের অনুপম পরিবারত...
- সেন এক জেড়া এয়ারিং এগিয়ে দেন সেকরার হাতে তা পরখ করে দেখার জন্য। কেননা সেটা ছিল অনুপমের মামার দেওয়া বিলাতি জিনিস, যাতে সেনার ভাগ আছে সামান্যই। গ. উদ্দীপকের সবিতা ও অপরিচিতা গল্পের কন্যাণী উভয়েই যৌ...
- তেইশ, এখন ইইয়াছে সাতশ। এখনা আশা ছাড়ি নাই, কিন্তু মাতুলকে ছাড়িয়াছি। নিতান্ত এক ছেলে বলিয়া মা আমাকে ছাড়িতে পারেন নাই। তোমরা মনে কর্তেছে, আমি বিবাহের আশা করি না, কোনো কালেই না। আমার মনে আছে, কেবল স...
- সংখ্যা ৬। অপরিচিতা গল্পে নায়কের বয়স কত বলা হয়েছে ক ২৮ বছর খ ২৬ বছর গ ২৭ বছর ঘ ২৫ বছর ৭ তরু হবার বিশেষ মূল্য আছে এখানে কীসের মূল্যের কথা বলা হয়েছে ক জীবনের খ মরণের গ কর্মের ঘ ধর্মের ৮। ছেলেবেলায় প...
- আমি অন্নপূর্ণার কোলে গজাননের ছোটো ভাইটি। মতো তিনি আমাদের সমস্ত সংসারটাকে নিজের অন্তরের মধ্যে স্থাষিয়া লইয়াছেন। তাহাকে না খুঁড়িয়া এখানকার এক গণ্ডষও রস পাইবার জো নাই। এই কারণে কোনো-কিছুর জন্যই আমাকে ক...
- যাঁহাও, কিন্তু মামাকে দমানো শক্ত। তিনি শঙ্কনাথবাবুর চুপচাপ ভাব দেখিয়া ভাবিলেন, লোকটা নিতান্ত নিজীব, একেবারে কোনো তেজ নাই। বেয়াই-সম্পদায়েের আর যাই থাক, তেজ থাকটা দোষের, অতএব মামা মনে মনে খুশি হইলেন।।...
- মিল আছে ক হরিশের খ মামার গ শিক্ষকের ঘ বিনুর ৪। উক্ত চরিত্রে প্রাধান্য পেয়েছে - 1 দৌরাভ্য 1 যীনমনাতা 11 লোভ নিচের কোনটি ঠিক ক।। ও যা ও গ।।৬ ঘ।। 1 ও ৫ অনুপমের বয়স কত বছর ক পটিন্স খ ছাবিবিশ গ সা...
- হাতে লইয়া দেখিলেন, এই এয়ারিং দিয়াই কন্যাকে তাহার আশীর্বাদ করিয়াছিলেন। মামার মুখ লাল হইয়া উঠিল। দরিদ্র তাহাকে ঠকাইতে চাখিব কিন্তু তিনি ঠকিবেন না এই আনন্দ-সমুদ্রে হইতে বঞ্চিত হইলেন এবং তাহার উপরেও ...

[Ask another query](#)

Figure 2

English Queries

#### Query:

Anupom er boyosh koto?

#### Answer:

অপরিচিতা গল্পে অনুপমের বয়স ২৭ বছর বলা হয়েছে।

**Groundedness Score: 0.58**

**Relevance Score: 0.62**

#### Source Snippets:

- সংখ্যা ৬। অপরিচিতা গল্পে নায়কের বয়স কত বলা হয়েছে ক ২৮ বছর খ ২৬ বছর গ ২৭ বছর ঘ ২৫ বছর ঙ তবু ইহার বিশেষ মূল্য আছে এখানে কীসের মূল্যের কথা বলা হয়েছে ক জীবনের খ মরণের গ কর্মের ঘ ধর্মের ঙ। ছেলেবেলায় প...
- ৫০। কাকে অনুপমের ভাগ্য দেবতা বলে উল্লেখ করা হয়েছে ক হরিশকে খ মামাকে গ বিনুদাকে ঙ। কার টাকার প্রতি আসক্তি বেশি ক শত্ৰুনাথের খ কল্যাণীর গ অনুপমের ঙ। কিছুদিন পূর্বে এমএ পাশ করিয়াছি। উজ্জ্বলি কার ক মামার...
- হয়েছে কোন চরিত্র ক আনোদি খ আসমা গ কল্যাণী ঘ মাদাম মোহিসেন উত্তর গ 29 30 কল আনলাইন ব্যাচ 11072 1 984159 ঙ। ছেলেবেলায় অনুপমের চেয়ারা নিয়ে বিদ্রপ করার সময় পণ্ডিতমশায় কোন দুটি ফুল ও ফলের সঙ্গে তুলন...
- ৩৫। গাটার সম্পর্কটাকে স্থায়ী করিবার ইচ্ছা আমার নাই। উজ্জ্বলি ক বিনুদাদার খ অনুপমের গ মামার ৩৬। অনুপম কাকে নিয়ে তীর্থযাত্রা শুরু করে ক কল্যাণীকে খ মাকে গ হরিশকে ৩৭। মা-পুত্রের তীর্থযাত্রার বাহন কী ছ...
- ঘ ৬। অপরিচিতা গল্পে রেলকর্মচারী কলটি টিকিট বেঞ্চে খুসি হয়েছিল খ বো ২২ ক একটি খ দুইটি গ তিনটি ঘ চারটি উত্তর খ ঙ। অপরিচিতা গল্পে কল্যাণী বিয়েতে কোন রঙের শাড়ি পরেছে বলে অনুপম কল্পনা করে কু বো ২২ ক খ...
- ক অনুপমের মামা খ অনুপমের মা গ শত্ৰুনাথ বারু ঘ হরিশ উত্তর গ ৩৩। শিবিনের সাথে কল্যাণীর মিল জোখায়, বাকার আঁজাবোই নিচ্ছে কোনটি সঠিক ক ১ খা, 11 গ, যা, 1 উত্তর ক নিচের উল্লিখিত পড়ে ৩৪ ও ৩৫ নম্বর প...
- মিল আছে ক হরিশের খ মামার গ শিবকের ঘ বিনুর ৪1 উক্ত চরিত্রে প্রধান্য পেয়েছে - 1 দৌরাখ্য 1 ইনন্দনাথ 11 লোভ নিচের কোনটি সঠিক ক। ও খ। ও গ।1 ও ঘ।1 1 ও ৫ অনুপমের বয়স কত বছর ক পঁচিশ খ ছাব্বিশ গ সা...
- মিল আছে ক হরিশের খ মামার গ শিবকের ঘ বিনুর উত্তর খ বাখ্যা দীপুর চাচা ও অপরিচিতা গল্পের অনুপমের মামার লোভ সীমাহীন। তারা উভয়েই ষোড়কলোড়ী। এই লোকী মানসিকতার দিক দিয়ে তাদের মধ্যে মিল রয়েছে। ৪। উক্ত...
- জন্মগ্রহণ করেন ক ১২৬১ খ ১২৬৮ গ ১২৭০ ঘ ১২৭২ ঙ ৫। অনুপম আহারে বসতে পারল না কেন ক তেমন ক্ষুধা ছিল না বলে খ আহার স্বাভাবিক ছিল না বলে গ মন কষাকষি হয়েছিল বলে ঘ মামার অনুমতি ছিল না বলে ঙ ৬। রবীন্দ্রনাথ ঠাকুরে...
- গল্পের অনুপমের মামার সাদৃশ্য ও বৈসাদৃশ্য দুটোই রয়েছে। 42 43 কল আনলাইন ব্যাচ 11072 1 984159 প্রশ্ন- ২ পড়াশুনা শেষ করে সবিতা এখন গ্রামের একটি সরকারি প্রাইমারি স্কুলে শিক্ষকতা করেন। বছর কয়েক আগে শহরে...

[Ask another query](#)

Figure 3

## Used Tools and Libraries

- OCR: pytesseract and pdf2image (in ocr\_pdf.py) with Poppler.
- FastAPI: Web app framework.
- LangChain: RAG implementation (langchain, langchain-openai, etc.).
- OpenAI: ChatOpenAI (e.g., gpt-4o).
- HuggingFace: sentence-transformers/paraphrase-multilingual-mpnet-base-v2 for embeddings.
- FAISS: Vector store.
- pdfplumber: PDF text extraction.
- NLTK: Text preprocessing.
- NumPy & scikit-learn: Similarity calculations.
- Pillow: Image handling for OCR.
- uvicorn: ASGI server.

## API Documentation

- **/ (GET):** Shows a query form.
- **/query (POST):**
  - Params: query (string), language (e.g., "bn", "en").
  - Response: HTML with answer, scores, and snippets.
- **/reinit (GET):** Re-initializes RAG chain, returns JSON {"message": "RAG chain re-initialized"}.

## Evaluation Matrix

- **Groundedness Score:** Support from documents (0-1, cosine similarity mean).
- **Relevance Score:** Best query-chunk match (0-1, cosine similarity max).

## Must Answer Questions

- **Text Extraction Method:**

- Used pytesseract and pdf2image in ocr\_pdf.py with Poppler for scanned PDFs.
- Faced formatting issues (e.g., broken text), fixed with regex preprocessing.
- **Chunking Strategy:**
  - RecursiveCharacterTextSplitter (size 1500, overlap 300).
  - Works well for semantic retrieval with fragmented text, though larger chunks could help.
- **Embedding Model:**
  - sentence-transformers/paraphrase-multilingual-mpnet-base-v2.
  - Chosen for Bangla support and semantic accuracy, captures meaning via contextual vectors.
- **Query Comparison:**
  - Cosine similarity with FAISS storage.
  - Chosen for robustness and speed in high-dimensional spaces.
- **Meaningful Comparison:**
  - Embeddings map query and chunks to a shared space; custom prompt aids inference.
  - Vague queries return "দুঃখিত, এই প্রশ্নের উত্তর নির্ধারণ করা সম্ভব নয়।"
- **Relevance of Results:**
  - Relevant for factual queries, not interpretive ones.
  - Improvements: Larger chunks, xlm-roberta-base embedding, or more document context.