

The idea

- Create a notes summarization website
- **Functionalities Overview**
 1. **Transcript Fetching:** Fetch the transcript from a YouTube link or from a user-uploaded lecture.
 2. **Lecture Recording:** Record audio through the browser and convert it into text (transcript).
 3. **Summarization and Note-Making:** Summarize the transcript and generate notes.
 4. **Save Notes:** Store the notes in a user-specific folder with options to edit and add images.
 5. **Quiz Generation:** Create quizzes based on the transcript content.
 6. **Flashcards Generation:** Generate flashcards for study purposes.
 7. **Q&A:** Implement a Q&A system that allows users to ask questions about the transcript.

The process

1. Frontend Development

UI/UX Design

- Design a user-friendly interface where users can:
 - Enter a YouTube link or upload a transcript.
 - Record their lecture.
 - View and edit summarized notes.
 - Access quizzes and flashcards.

2. Components

- **Transcript Fetching Form:** A form to input the YouTube link or upload files.
- **Quiz and Flashcard Viewers:** Components to display quizzes and flashcards.

3. Backend Development

Transcript Fetching

- **YouTube Transcript:** Use YouTube Data API to fetch the transcript if available.

<https://github.com/jdepoix/youtube-transcript-api>

- **Speech-to-Text:** Implement speech-to-text functionality using an API

SpeechRecognition library with Google Web Speech API (which has a free tier for small usage).

Summarization and Note-Making

Text summarizer <https://github.com/Mohit-Kundu/Text-Summarization-Web-Application>
<https://github.com/tkmanabat/Text-Summarization>
[https://github.com/nicknochnack/PegasusSummarization/blob/main/Pegasus Tutorial.ipynb](https://github.com/nicknochnack/PegasusSummarization/blob/main/Pegasus%20Tutorial.ipynb)
<https://www.youtube.com/watch?v=Yo5Hw8aV3vY>
https://colab.research.google.com/drive/1cQgDmiYZzQlu_dJPI_vqOJp5GhjunnyF4
<https://www.youtube.com/watch?v=3V-MJhJvRWg>

- **Summarization Model:**

In dono ke model se comparison hoga pegasus x sum ka

- **BART (Bidirectional and Auto-Regressive Transformers):**
 - **Strengths:** BART excels in summarization tasks, particularly for longer texts, making it suitable for transcripts.
 - **Advantages:** It can produce coherent and contextually relevant summaries, especially when fine-tuned on domain-specific data.
 - **Use Case:** Effective for both extractive and abstractive summarization of lecture transcripts.
- **T5 (Text-to-Text Transfer Transformer):**
 - **Strengths:** T5 is versatile and can handle various text-based tasks, including summarization. Its text-to-text framework allows it to work well with different formats.
 - **Advantages:** It performs well across diverse datasets and can generate concise summaries of transcripts.
 - **Use Case:** Good for converting long lecture transcripts into clear summaries.
- **Pegasus X sum: (this is the main one we are implementing)**
 - **Strengths:** Designed specifically for abstractive summarization, Pegasus can create high-quality summaries from extensive texts.
 - **Advantages:** It has shown state-of-the-art performance on summarization benchmarks and is effective with complex narratives, such as lecture content.
 - **Use Case:** Suitable for summarizing detailed lectures or presentations.

Pinecone

For your NoteAI project, Pinecone is an excellent choice for implementing a **retrieval-augmented generation (RAG)** architecture with Llama. To ensure that you can generate responses for all types of videos or articles, your Pinecone database should be designed with the following considerations:

1. Data Structure

Use a **vector database** schema optimized for semantic search and retrieval. Here's a structure you can adopt:

- **ID:** A unique identifier for each video/article (e.g., video ID or URL).
 - **Metadata:** Includes details like title, description, tags, and category (e.g., "science", "technology", "history").
 - **Text Segments:** Break down the transcript/article into smaller chunks (e.g., 256-512 tokens) for better retrieval performance.
 - **Embeddings:** Store the vector representation of each text segment, generated using a model like **OpenAI's text-embedding-ada-002** or **Hugging Face sentence-transformers**.
-

2. Embedding Model

To ensure high-quality responses across diverse content:

- Use a **general-purpose embedding model** for generating vector representations, such as:
 - [Hugging Face's sentence-transformers](#): Example models include all-MiniLM-L6-v2 or multi-qa-mpnet-base-dot-v1.
 - [OpenAI's embeddings API](#).
 - Fine-tune the embedding model (if necessary) for your dataset to improve domain-specific retrieval.
-

3. Index Configuration

Set up the Pinecone index with:

- **Metric:** cosine similarity or dot product for efficient nearest neighbor search.
 - **Index Size:** Ensure the index can scale with your dataset by choosing a Pinecone plan that supports large-scale embeddings.
 - **Sharding:** Use sharding for better performance across multiple queries if you anticipate high throughput.
-

4. Implementation for RAG with Llama

Here's how you can integrate Pinecone with RAG:

1. **Ingest Data:** Add the transcript/article segments and their embeddings into Pinecone.
2. **Query Pipeline:**
 - Tokenize the user's question.
 - Embed the question using the same embedding model.
 - Query Pinecone to fetch the top-N relevant segments (e.g., top 5).

- Concatenate the retrieved segments to form the context.
3. **Generate Answer:** Pass the context and question to **Llama** for answer generation.
-

5. Optimization Tips

- **Preprocessing:** Normalize text (e.g., lowercasing, removing special characters) to improve embedding quality.
 - **Metadata Filters:** Use metadata filters in Pinecone to restrict search by category, date, or language.
 - **Chunking Strategy:** Use overlapping sliding windows (e.g., 512 tokens with 128-token overlap) when chunking text to preserve context across boundaries.
-

6. Additional Features

- **Multimodal Support:** Extend your database schema to include video/audio embeddings (e.g., CLIP or Whisper embeddings) for richer context.
- **Dynamic Updates:** Implement real-time indexing in Pinecone to handle new transcripts/articles dynamically.

With this approach, Pinecone will enable fast, accurate retrieval for your RAG-based Q&A system, ensuring compatibility with diverse content types.

Q&A

Technologies:

https://github.com/CoderNitu/Question_Answering_Web_App?tab=readme-ov-file

To build a question-answering application using the open-source framework **Haystack** that incorporates **transcript fetching**, you can follow these detailed steps. Here's a thorough explanation of how the system would work, integrating transcript fetching as a core feature.

BERT/roBERTa for Question Answering (for comparison)

- **Description:** BERT and its variants (like roBERTa and DistilBERT) are designed for extractive question answering tasks. They identify spans of text from a given context that best answer the question.
- **Advantages:**
 - **Effective for Extractive QA:** These models are pre-trained on QA datasets (like SQuAD) and excel at finding exact answers within the provided context.
 - **Performance:** They can provide highly accurate answers when the answer is present in the text.
- **Usage Scenario:** Best suited for transcripts where the questions are likely to have direct answers in the text, such as factual information.

Example Code for BERT/roBERTa

python

Copy code

```
from transformers import pipeline
```

```
# Load the QA pipeline
```

```
qa_pipeline = pipeline("question-answering")
```

```
# Example context and question
```

```
context = "The Eiffel Tower is located in Paris, France."
```

```
question = "Where is the Eiffel Tower located?"
```

```
# Get the answer
```

```
result = qa_pipeline(question=question, context=context)
```

```
print(result['answer']) # Output: Paris
```

we are making one with RAG

A. Data Ingestion

- Text segmentation.
- Embedding generation using Hugging Face Sentence Transformer.
- Storing embeddings in Pinecone for vector search.

B. Query Processing

- Tokenizing user queries.
- Retrieving relevant text segments from Pinecone.

C. Answer Generation

- Contextual answer generation using Llama.

Quiz Generation

- **MCQ and Flashcards:** Automatically generate multiple-choice questions and flashcards based on key concepts identified in the transcript.
- <https://git.foosoft.net/alex/anki-connect> for flashcards

<https://github.com/PragatiVerma18/MLH-Quizzet?tab=readme-ov-file>

Here's a step-by-step **implementation plan** for your pipeline:

Step 1: Text Preprocessing

1. **Objective:** Clean the text to remove unnecessary characters and standardize the format.
2. **Actions:**
 - Remove non-alphanumeric characters using regex.
 - Convert text to lowercase for uniformity.
 - Optionally, tokenize the text into sentences or words.

Step 2: Named Entity Recognition (NER)

1. **Objective:** Extract key entities (e.g., names, dates, places, quantities) from the text.
2. **Actions:**
 - Use a pre-trained NER model (e.g., spaCy, Hugging Face) to identify entities.
 - Extract entity types such as PERSON, DATE, ORG, LOC, QUANTITY, etc.

Step 3: Masking Sentences

1. **Objective:** Replace identified entities with placeholders to create "fill-in-the-blank" questions.
2. **Actions:**
 - For each sentence containing an entity, replace the entity with a mask (e.g., _____).
 - Keep track of the original entity as the answer.

Step 4: Ranking Entities

1. **Objective:** Prioritize entities to determine which ones are the most relevant for quiz questions.
2. **Actions:**
 - Rank entities by frequency (most mentioned entities in the text).
 - Use the entity's position in the text (e.g., entities in the introduction or summary might be more important).
 - Optionally, assign custom weights to specific entity types (e.g., PERSON > DATE).

Step 5: Generating Alternatives

1. **Objective:** Create plausible distractors (incorrect options) for multiple-choice questions.

2. **Actions:**

- Use WordNet to generate synonyms, antonyms, or related terms.
 - Use a text generation model (e.g., GPT, T5) to create paraphrased alternatives.
 - Randomly select entities of the same type from other parts of the text as distractors.
-

Step 6: Creating Multiple-Choice Questions

1. **Objective:** Combine masked sentences with correct answers and distractors to form complete quiz questions.
 2. **Actions:**
 - For each masked sentence, list the correct answer and distractors as options.
 - Shuffle the options to randomize their order.
 - Format the question into a readable multiple-choice format.
-

Step 7: Compiling the Quiz

1. **Objective:** Organize all questions into a cohesive quiz format.
 2. **Actions:**
 - Number each question.
 - Include the correct answer and distractors for multiple-choice questions.
 - If required, generate a separate answer key for validation.
-

Example Workflow

1. Input: Raw text (e.g., *"Newton discovered gravity in 1687. The capital of France is Paris."*).
2. Cleaned Text: *"Newton discovered gravity in 1687 The capital of France is Paris"*
3. NER Output:
 - Entities:
 - Newton (PERSON)
 - 1687 (DATE)
 - Paris (LOC)
4. Masked Sentences:
 - *"_____ discovered gravity in 1687."*
 - *"The capital of France is _____."*

5. Ranked Entities:
 - Priority: Paris > Newton > 1687
 6. Generated Quiz:
 - Q1: *"Who discovered gravity in 1687?"*
 - (a) Newton (b) Einstein (c) Tesla (d) Darwin
 - Q2: *"The capital of France is ____."*
 - (a) Paris (b) Berlin (c) Rome (d) Madrid
-

Step 8: Optional Enhancements

1. **Flashcards:** Use Anki or other tools to create interactive flashcards for the quiz questions.
2. **Difficulty Levels:** Adjust questions based on complexity (e.g., fill-in-the-blank vs. MCQs with 4 options).
3. **Dynamic Quiz Selection:** Allow users to specify quiz length or focus on specific entity types.

This structured implementation ensures a systematic approach to creating an engaging and informative quiz.

AnkiConnect

Overview:

- **AnkiConnect** is a popular free plugin for Anki, a well-known flashcard application. It allows you to interact with Anki via an API, enabling you to programmatically create and manage flashcards.

Features:

- Create and update flashcards.
- Integrate with Anki for a powerful flashcard experience.

How to Use:

- Install Anki and the AnkiConnect plugin.
- Use the API to send requests to Anki from your application. For example, you can create flashcards by sending a JSON payload with details about the cards.

Example API Request:

json

Copy code

```
{  
  "action": "addNote",  
  "version": 6,
```



```

"params": {
  "note": {
    "deckName": "Default",
    "modelName": "Basic",
    "fields": {
      "Front": "What is the capital of France?",
      "Back": "Paris"
    },
    "tags": []
  }
}

```

Implement Automatic Question Generation (AQG)

NLP for Question Generation:

- Use NLP libraries or models to generate questions from the extracted text.

Python Libraries:

- **spaCy:** For NLP tasks and named entity recognition.
- **Transformers:** For advanced models like GPT or BERT.

Generate MCQs:

- You can use NLP models to generate MCQs. Although there's no direct pre-built model for MCQs, you can adapt text generation models to help create question-answer pairs with options.

Example Code Using Transformers:

python

Copy code

```
from transformers import pipeline
```

```
def generate_mcqs(text):
```

```
    question_generator = pipeline('text2text-generation', model='t5-small')
```

```
    # This might generate questions; you need to process these to MCQs
```

```
    questions = question_generator(text, max_length=100, num_beams=5, early_stopping=True)
```

```
mcqs = []

for question in questions:

    # Dummy options; you would need to implement logic to generate real options
    options = ["Option A", "Option B", "Option C", "Option D"]

    mcqs.append({

        'question': question['generated_text'],

        'options': options

    })

return mcqs
```

Database Management

- **User Authentication:** Implement user login/signup functionality using Django's built-in authentication or a third-party service like Firebase.
- **Data Storage:** Store transcripts, notes, quizzes, and flashcards in a structured format in the database.

4. Integrating Features

Recording and Transcription

- Send the audio to the backend for transcription using the chosen Speech-to-Text API.
- Ensure the transcription process is robust and handles different accents and noise conditions.

Summarization Workflow

- Once the transcript is ready, pass it through the summarization pipeline. The backend should handle this process asynchronously, allowing users to continue using the website while the summarization is happening.

Saving and Editing Notes

- Implement a saving mechanism where the summarized notes are stored in a user-specific directory or database. Allow editing and image embedding using the rich-text editor.

Quiz and Flashcard Management

- Design a backend process that generates quizzes and flashcards once the notes are saved. Store these in the database for later retrieval.

5. Deployment

- **Server Setup:** Deploy the backend on a cloud service like AWS or Heroku. Deploy the frontend on Vercel or Netlify.

- **CDN and Caching:** Use a CDN (Content Delivery Network) like Cloudflare to improve the loading times of your website.
- **Monitoring:** Implement monitoring tools like Sentry for error tracking and performance monitoring.

6. Additional Considerations

- **Scalability:** Consider how your application will scale with multiple users. Implement caching for repeated summarization tasks.
- **Security:** Implement HTTPS, secure API keys, and protect against common vulnerabilities like SQL injection and XSS.
- **User Experience:** Ensure that the website is responsive, accessible, and provides real-time feedback during processes like recording and summarization.

7. Testing

- **Unit Testing:** Test individual components, like the summarization model and the transcript-fetching functionality.
- **Integration Testing:** Ensure that the frontend and backend interact seamlessly, especially during asynchronous operations.
- **User Testing:** Conduct user testing sessions to identify usability issues and gather feedback for improvements.

8. Launch and Maintenance

- **Beta Launch:** Start with a beta version to gather feedback and fix bugs.
- **Continuous Updates:** Regularly update the models and improve functionality based on user feedback.
- **Documentation:** Provide clear documentation for users, especially if they need to interact with advanced features like uploading custom transcripts or recording lectures.

Abstract

A concise overview of the system, its objectives, methodologies, and outcomes.

Keywords

Automated Summarization, Question Answering, Quiz Generation, Flashcards, YouTube Transcripts, Wikipedia Articles, Natural Language Processing.

I. Introduction

- Overview of the problem and its significance.
- Motivation for automating summarization, Q&A, and learning content generation.

- Objectives of the project.
-

II. Input and Data Acquisition

A. Transcript and Article Input

- Handling YouTube videos: YouTube Transcript API.
- Handling Wikipedia articles: Web scraping using BeautifulSoup.

B. Preprocessing

- Text normalization.
 - Removing special characters and unwanted formatting.
-

III. Summarization System

A. Model Selection

- Pegasus-XSum: Features and advantages.

B. Summarization Workflow

- Tokenization of input.
 - Passing tokenized text through the summarization model.
-

IV. Question Answering System

A. Data Ingestion

- Text segmentation.
- Embedding generation using Hugging Face Sentence Transformer.
- Storing embeddings in Pinecone for vector search.

B. Query Processing

- Tokenizing user queries.
- Retrieving relevant text segments from Pinecone.

C. Answer Generation

- Contextual answer generation using Llama.
-

V. Quiz Generation

A. Preprocessing

- Text cleaning and tokenization.

B. Entity Recognition and Masking

- Named Entity Recognition (NER) using spaCy or Hugging Face.
- Creating fill-in-the-blank questions.

C. Distractor Generation

- Generating plausible distractors using WordNet or language models (e.g., GPT, T5).

D. Question Formatting

- Shuffling options.
 - Compiling multiple-choice questions with answer keys.
-

VI. Flashcard Generation

A. Methodology

- Integration with AnkiConnect API.

B. Flashcard Creation

- Front: Masked sentences or key terms.
 - Back: Correct answers or definitions.
-

VII. Challenges and Limitations

- Scalability with large datasets.
 - Accuracy of NER and distractor generation.
 - Latency in Q&A and summarization pipelines.
-

VIII. Future Enhancements

- Expanding to support more input types.
 - Integrating additional summarization and Q&A models.
 - Incorporating adaptive learning features.
-

IX. Conclusion

- Summary of achievements.
 - Contributions to learning and content comprehension.
-

X. References

- Cite all models, APIs, and libraries used (e.g., Pegasus-XSum, Llama, Hugging Face, Pinecone).
- Papers and resources referenced in building the system.

This structure balances technical detail with the flow of a standard research paper and ensures every component of your project is well-documented.