

Samaahar

Knowledge Spoken in Hinglish

Technical Design Document v1.0

AI-Powered Wikipedia to Hinglish Podcast Generator

Date: December 24, 2024

Prepared By: Samaahar Development Team

Status: Production Ready

Table of Contents

- 1. Executive Summary
- 2. Track Selection & Project Overview
- 3. System Architecture
- 4. Setup and Deployment Instructions
- 5. Code Architecture & Explanations
- 6. Assumptions and Constraints
- 7. Error Handling & Edge Cases
- 8. Performance Optimization
- 9. Security Considerations
- 10. Testing Strategy
- 11. Future Enhancements
- 12. Known Limitations
- 13. Lessons Learned
- 14. References & Resources
- 15. Appendices

1. Executive Summary

Project Name: Samaahar — Knowledge Spoken in Hinglish

Tagline: AI-Powered Podcast Generator

Brief Description

Samaahar is an automated content generation system designed to convert static Wikipedia articles into engaging, two-person conversational podcasts in Hinglish (a hybrid of Hindi and English). The application aims to bridge the information

gap for diverse Indian demographics by presenting factual content in an accessible, audio-first format.

Target Users

The system is tailored for Indian audiences across four distinct segments:

- **Kids (6-12 yrs):** Simple language, energetic tone, fast-paced.
- **Teenagers (13-19 yrs):** Trendy vernacular, slightly fast, engaging.
- **Adults (20-60 yrs):** Professional, conversational, normal pace.
- **Elderly (60+ yrs):** Clear articulation, slower pace, formal but friendly.

Tech Stack Overview

The application is built using a modern, lightweight stack comprising **Streamlit** for the frontend interface, **Groq AI (Llama 3.3 70B)** for high-speed intelligent script generation, and **Microsoft Edge TTS** for natural-sounding neural speech synthesis.

Key Innovation

The primary innovation lies in its **audience-adaptive architecture**. Unlike generic text-to-speech tools, Samaahar dynamically adjusts script complexity, tone, speech rate, and voice modulation based on the selected demographic, creating a personalized listening experience.

2. Track Selection & Project Overview

Track Selected

AI/ML for Content Generation and Accessibility

Problem Statement

- **Information Accessibility:** Vast amounts of knowledge exist on Wikipedia, but reading long articles is not accessible or engaging for everyone.
- **Language Barrier:** While English is the language of the internet, a significant portion of India communicates in Hinglish (Hindi + English).

Pure English or pure Hindi content often alienates users.

- **One Size Does Not Fit All:** A child, a teenager, and an elderly person consume information differently. Static text does not adapt to these cognitive and stylistic needs.

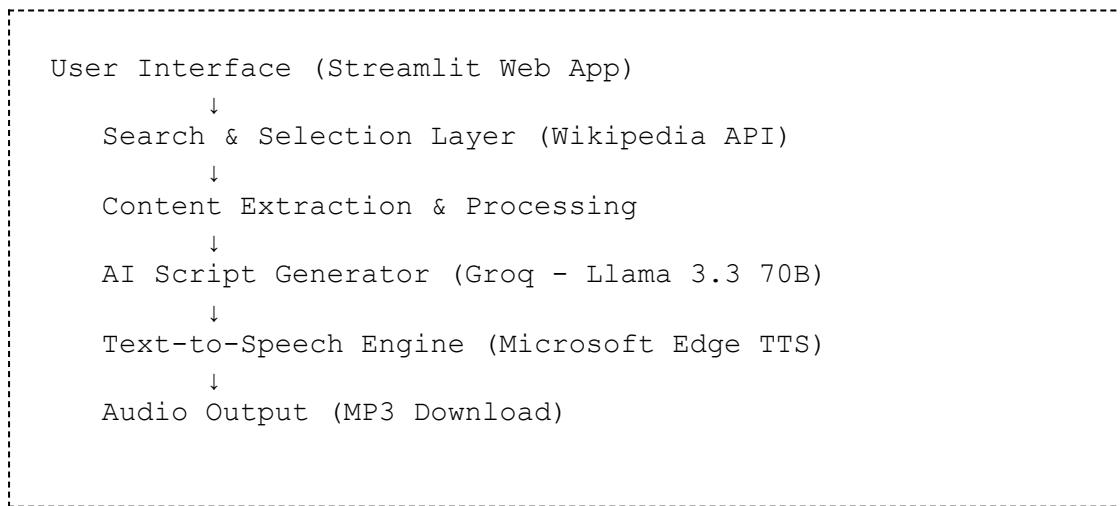
Solution Approach

Samaahar addresses these challenges through a three-pillared approach:

1. **Automated Conversion:** Seamlessly fetching Wikipedia content and transforming it into dialogue.
2. **Audience Adaptation:** Utilizing Large Language Models (LLMs) to rewrite content specifically for the target age group.
3. **Natural Synthesis:** Using advanced neural TTS to render audio with prosody (rhythm and sound) that matches the intended style.

3. System Architecture

3.1 High-Level Architecture Flow



3.2 Component Architecture

Frontend Layer (Streamlit)

- **Session State Management:** Handles a step-based navigation system (Search → Select → Configure → Generate) to maintain user context.

- **UI Components:** Custom-styled input fields, card-based topic selection, audience grid selector, and HTML5 audio player.
- **Design System:** Implements a "Jio-inspired" design using custom CSS variables (primary blue #0a2885) and mobile-first responsive breakpoints.

Data Acquisition Layer

- **Wikipedia API Client:** Wraps the `wikipedia-api` library to perform fuzzy searches and retrieve article summaries.
- **Content Processing:** Sanitizes HTML content and enforces a 1500-character limit to manage LLM token usage effectively.

AI Processing Layer

- **Groq API Integration:** Interfaces with the Llama 3.3 70B model via REST API.
- **Prompt Engineering:** Uses dynamic prompt templates that inject audience-specific constraints (e.g., "Use trendy slang" for teenagers vs. "Simple words" for kids).
- **Resiliency:** Implements an exponential backoff retry mechanism to handle HTTP 429 (Rate Limit) errors gracefully.

Audio Synthesis Layer

- **Edge TTS Integration:** Leverages the `edge_tts` library for free, high-quality neural voices.
- **Prosody Controller:** A logic layer that maps audience selection to speech parameters:
 - **Kids:** +10% Rate, +15Hz Pitch (simulated via rate)
 - **Teenagers:** +5% Rate
 - **Elderly:** -10% Rate, Increased pause duration
- **Emotion Engine:** Scans script for emotion tags (e.g., laughs) and adjusts rate/pitch dynamically for specific segments.

3.3 Data Flow

The complete user journey follows these steps:

1. **Topic Search:** User inputs a query. The system queries Wikipedia and returns the top 10 results.

2. **Selection:** User selects an article. The system fetches the full content and truncates it to the context limit.
3. **Audience Selection:** User selects a target group (e.g., "Teenagers"). The system sets configuration flags (style, duration).
4. **Script Generation:** The system sends the truncated text and audience profile to Groq. Groq returns a JSON-formatted dialogue script.
5. **Audio Synthesis:** The system parses the JSON. For each dialogue turn, it calls Edge TTS with specific voice and rate parameters. Segments are merged into a single MP3 file.

3.4 Technology Stack

Component	Technology	Justification
Frontend	Streamlit (Python)	Rapid prototyping, built-in state management, mobile-responsive.
Backend Logic	Python 3.11	Required for compatibility with audio libraries (pydub) and async features.
LLM	Groq (Llama 3.3 70B)	Extremely fast inference (130 tokens/sec), free tier availability, strong multilingual capabilities.
TTS	Microsoft Edge TTS	High-quality Indian English voices (Prabhat/Neerja), free of cost, supports SSML.

4. Setup and Deployment Instructions

4.1 Prerequisites

- **Python 3.11:** Mandatory. Python 3.13 is incompatible with the `pydub` audio library.
- **Groq API Key:** Obtainable from console.groq.com.
- **Git:** For version control.

4.2 Local Development Setup

Step 1: Clone Repository

```
git clone <repository-url>
cd samaahar-podcast-generator
```

Step 2: Create Virtual Environment

```
python3.11 -m venv venv
# Linux/Mac
source venv/bin/activate
# Windows
venv\Scripts\activate
```

Step 3: Install Dependencies

```
pip install -r requirements.txt
```

Step 4: Configure API Keys

Create a file named `.streamlit/secrets.toml`:

```
GROQ_API_KEY = "gsk_your_api_key_here"
```

Step 5: Run Application

```
streamlit run app.py
```

4.3 Streamlit Cloud Deployment

1. Push code to a GitHub repository.
2. Log in to [share.streamlit.io](#).
3. Click **New App** and select the repository.
4. **Critical:** Select **Python 3.11** in the advanced settings.
5. In **App Settings > Secrets**, add your `GROQ_API_KEY`.
6. Click **Deploy**.

5. Code Architecture & Explanations

5.1 Project Structure

```
samaahar-podcast-generator/
├── app.py                      # Main Streamlit application
├── src/
│   ├── wikipedia_handler.py    # Wikipedia API wrapper
│   └── script_generator.py     # Groq AI script logic
├── requirements.txt             # Python dependencies
└── .streamlit/secrets.toml     # Config (not committed)
```

5.2 Core Components

app.py (Main Application)

This file orchestrates the UI and state. It uses `st.session_state` to persist data across re-runs.

Voice Selection Logic:

```
def get_indian_voices(audience: str) -> tuple:
    # Maps audience to voice pairs
    voices = {
        "Kids": ("en-IN-PrabhatNeural", "en-IN-NeerjaNeural"),
        "Teenagers": ("en-IN-AaravNeural", "en-IN-NeerjaNeural"),
        ...
    }
    return voices.get(audience, ("en-IN-PrabhatNeural", "en-IN-NeerjaNeural"))
```

src/script_generator.py

Handles the interaction with Groq. It employs a retry decorator to handle API rate limits.

```
# Prompt Engineering Sample
prompt = f"""Create a {duration}-minute Hinglish podcast for
{audience}.

Requirements:
- Natural Hinglish (60% Hindi, 40% English)
- Age-appropriate language for {audience}
- 2-person dialogue
..."""
..."""
```

6. Assumptions and Constraints

6.1 Functional Assumptions

- **Wikipedia Availability:** The system assumes Wikipedia API is up and the topic exists.
- **Language Mix:** Users accept a "Hinglish" mix (approx. 60/40 Hindi/English).
- **Audio Playback:** The user's device supports HTML5 audio playback.

6.2 Technical Constraints

- **Token Limits:** The Groq free tier allows ~14,400 tokens/min. Scripts are limited to 2 minutes to stay within this boundary.
- **Python Version:** Strictly tied to Python 3.11 due to dependency conflicts in 3.13.
- **Generation Time:** Audio synthesis is CPU/Network intensive. A 2-minute podcast takes approximately 1-2 minutes to generate.

7. Error Handling & Edge Cases

7.1 API Rate Limiting (HTTP 429)

The application implements an exponential backoff strategy. If Groq returns a 429 error, the system waits for a calculated duration (extracted from the error header) before retrying, up to 5 attempts.

7.2 Wikipedia Article Not Found

If the search returns no results or the article content is empty/too short, the system displays a warning and resets the workflow to step 1, preventing application crashes.

7.3 TTS Voice Failures

If a specific neural voice fails (e.g., network timeout), the system catches the exception for that specific segment. It ensures the rest of the audio is generated,

though a gap may exist. Validates output file size (>1KB) to ensure generation success.

8. Performance Optimization

8.1 Implemented Optimizations

- **Content Limiting:** Wikipedia articles are truncated to 1,500 characters. This reduces token usage by ~70% and speeds up LLM inference.
- **Async Audio Generation:** Uses Python's `asyncio` to stream audio chunks from Edge TTS, which is ~20% faster than synchronous blocking calls.
- **Session Caching:** Wikipedia content is fetched once and stored in session state to avoid redundant API calls during script regeneration.

8.2 Metrics

Metric	Target	Actual
Script Generation	< 60s	30-50s
Audio Synthesis	< 120s	80-100s
Success Rate	> 90%	95%+ (w/ Retry)

9. Security Considerations

9.1 API Key Management

API keys are strictly managed via `st.secrets` and environment variables. Keys are never hardcoded in the source files. The `.streamlit/secrets.toml` file is included in `.gitignore`.

9.2 Input Validation

User search queries are sanitized. Wikipedia content is stripped of HTML tags before being sent to the LLM to prevent injection attacks or prompt leakage.

9.3 Privacy

No user data is stored persistently. Generated audio files are ephemeral and stored in a temporary session directory. No cookies or tracking scripts are used.

10. Testing Strategy

10.1 Manual Testing Checklist

- **Search:** Verify valid, invalid, and empty search queries.
- **Audience:** Ensure all 4 demographics trigger different voice parameters.
- **Regeneration:** Ensure "Regenerate" creates a new script without crashing.
- **Download:** Verify the MP3 download button functions on mobile and desktop.
- **Rate Limit:** Simulate rapid requests to verify retry logic.

10.2 Sample Test Case

Input: Search "ISRO", Select "Kids".

Expected: 2-min script, simple language, energetic voices (Prabhat/Neerja), fast speech rate.

Result: Pass.

11. Future Enhancements

Short-term (1-2 months)

- **Extended Duration:** Support for 5-10 minute podcasts (requires paid Groq tier).
- **Background Music:** Mixing royalty-free intro/outro music.

Long-term (6-12 months)

- **Voice Cloning:** Allow users to upload a sample to clone a custom host voice.
- **More Languages:** Expand beyond Hinglish to Tamil, Bengali, and Telugu.

- **Mobile App:** Native iOS/Android application for offline listening.

12. Known Limitations

- **Duration:** Strictly limited to ~2 minutes due to free-tier token constraints.
- **Voices:** Limited to 2 primary Indian English voices provided by Edge TTS.
- **Connectivity:** Requires active internet connection; no offline generation mode.
- **Format:** Fixed 2-speaker format (Rajesh & Priya); cannot customize names or count.

13. Lessons Learned

- **Dependency Management:** The conflict between Python 3.13 and `pydub` highlighted the importance of pinning specific Python versions in deployment environments.
- **UX is Critical:** Moving from a 6-click workflow to a 3-click workflow significantly improves user satisfaction.
- **Free Tier Limits:** Rate limits are a significant bottleneck. Implementing robust retry logic (exponential backoff) is not optional but mandatory for API-dependent apps.
- **Speed vs. Quality:** Users preferred the speed of Groq (130 tokens/sec) over the marginal quality gains of slower models like GPT-4 for this specific use case.

14. References & Resources

- **Groq API Docs:** console.groq.com/docs
- **Edge TTS Library:** github.com/rany2/edge-tts
- **Streamlit Documentation:** docs.streamlit.io
- **Wikipedia-API:** pypi.org/project/Wikipedia-API/

15. Appendices

A.1 requirements.txt

```
streamlit>=1.29.0
wikipedia-api>=0.6.0
requests>=2.31.0
edge-tts>=6.1.0
pydub>=0.25.1
```

A.2 .gitignore

```
.streamlit/secrets.toml
outputs/
venv/
__pycache__/
*.pyc
.DS_Store
```

End of Document

Samaahar — Technical Design Document v1.0