

Experiment No.: 11

MINIPROJECT REPORT

TITLE: VISION-VOICE MULTIMODAL CORRECTION AGENT

1. AIM

To design and implement a multimodal AI application that can analyze visual (image), auditory (voice), and linguistic (text) inputs, generate reasoned responses, and apply self-correction mechanisms using multiple AI models to ensure higher accuracy and reliability.

2. OBJECTIVES

1. To develop an application capable of accepting multiple input modalities — image, voice, and text.
2. To integrate multiple AI models where the output of one model is verified or refined by another, ensuring logical consistency.
3. To demonstrate cross-modal data science principles, where models from different modalities complement each other.
4. To increase prediction precision by introducing a correction model that validates prior analysis.
5. To provide an intuitive Gradio-based user interface for interacting with multimodal AI systems.
6. To synthesize the final response into human-like audio using a neural text-to-speech system (Kokoro TTS).

3. INTRODUCTION

3.1 Overview: Multimodal AI represents the integration of diverse forms of data such as images, audio, and text to enable richer understanding and reasoning. Unlike unimodal AI systems that depend on a single input type, multimodal systems replicate how humans perceive information — integrating vision, hearing, and language.

This mini project — **Vision-Voice Multimodal Correction Agent** — is a self-corrective AI pipeline that leverages three powerful models:

- **Whisper** (Speech-to-Text),
- **Gemini (Generative AI)** (Vision + Language reasoning and correction),
- **Kokoro** (Text-to-Speech synthesis).

The unique feature is that the second model verifies and corrects the first model's output, producing refined, factual, and natural results. The system ensures *cross-model validation* and high accuracy, demonstrating multimodal data fusion in practical data science applications.

3.2 Motivation: Traditional AI systems are error-prone when relying on a single modality. For instance:

- Vision-only models may misinterpret images.
- Language-only models may lack visual context.
- Speech-only systems might misrecognize words.

By combining these, this project improves both the breadth of understanding and the depth of reasoning, forming a self-aware correction loop.

3.3 Problem Statement: To build an AI system that:

- Accepts an image (like a certificate), a voice question, or a typed query.
- Generates an initial analysis of the image based on the query.
- Uses a correction model to refine the analysis automatically.
- Provides both the final corrected text and a natural voice response.

4. LITERATURE REVIEW

4.1 Multimodal Data Science: Multimodal learning enables AI systems to interpret different forms of information jointly. Studies show that combining visual, auditory, and textual data improves generalization and robustness.

4.2 Whisper (OpenAI, 2023): Whisper is an automatic speech recognition (ASR) model trained on multilingual and multitask datasets. It converts human speech into accurate text even with background noise.

4.3 Gemini (Google, 2024): Gemini models (Vision + Language) perform reasoning tasks that combine visual interpretation and textual understanding, supporting multimodal input and conversational correction.

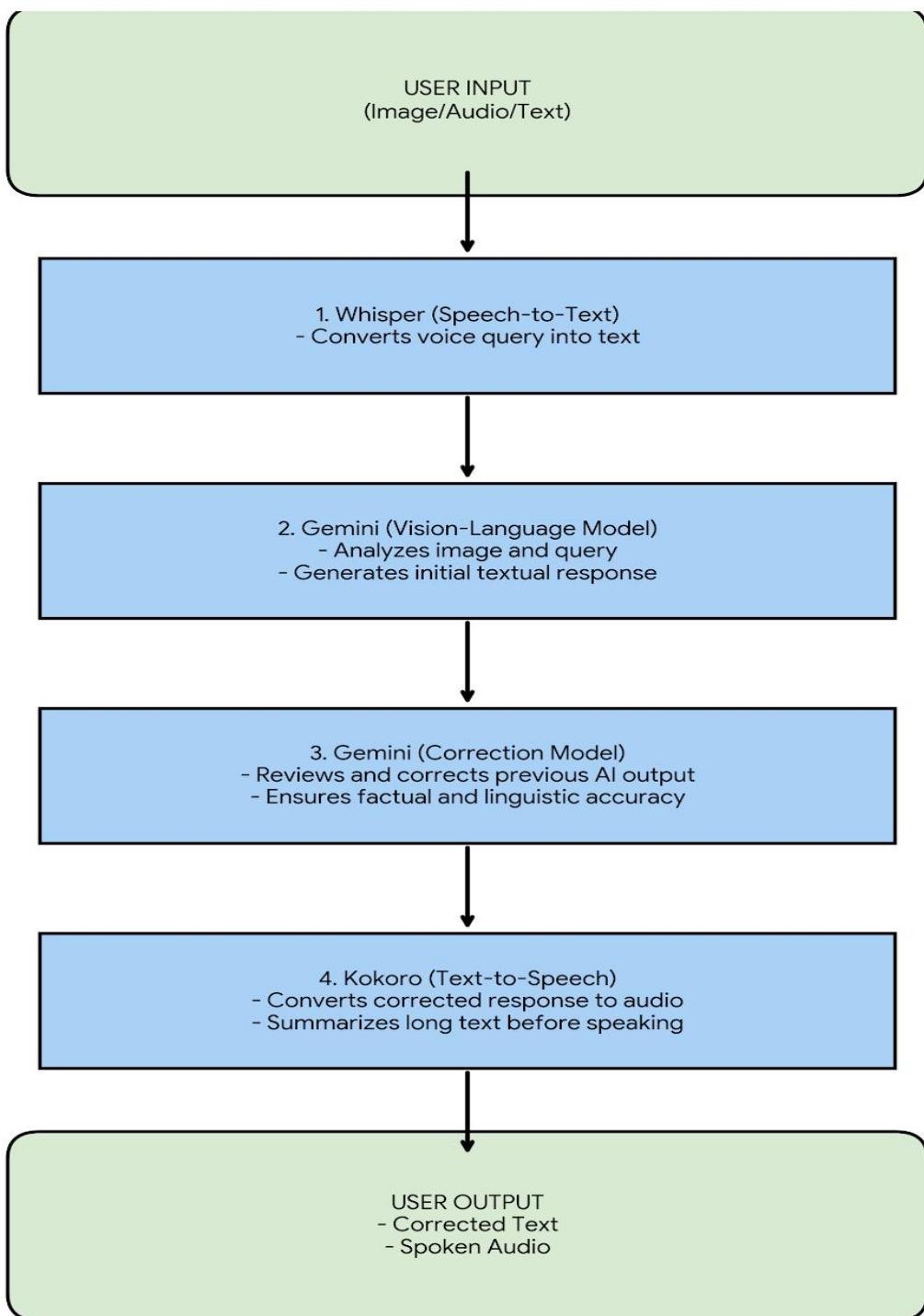
4.4 Kokoro TTS (Hexgrad Labs): Kokoro is a fast, high-quality speech synthesis model. It generates natural audio from text, supporting multiple voice tones and languages.

4.5 Gradio & Flask in AI UI: Gradio simplifies AI model interfaces, providing easy front-end integration, while Flask supports backend API management.

5. SYSTEM ARCHITECTURE

Architecture Overview: The **Vision-Voice Multimodal Correction Agent** consists of the following components:

System Architecture Flow Chart



6. MODULE DESIGN

6.1 Audio Processing Module (Whisper)

- Converts recorded audio to text.
- Uses sampling rate normalization (16 kHz).

- Handles noise and accents with multilingual robustness.

6.2 Vision-Language Module (Gemini Initial Analysis)

- Accepts the uploaded image and query.
- Performs visual understanding and contextual response generation.
- Example: “Who is this certificate issued to?”

6.3 Correction Module (Gemini Second Pass)

- Evaluates the first Gemini response.
- Checks factual accuracy, tone, and clarity.
- Generates refined output, ensuring a professional and logical tone.

6.4 Summarization Logic (Before TTS)

```
if len(corrected_response) > 800:
    summary_prompt = f"Summarize this text for audio output, keeping key points:\n{corrected_response}"
    corrected_response = self.manager.gemini_model.generate_content(summary_prompt).text
```

This step ensures audio remains concise and clear even for long responses.

6.5 Text-to-Speech Module (Kokoro)

- Synthesizes human-like speech.
- Supports multiple languages and speeds.
- Auto-summarizes long text before generation.

6.6 Frontend (Gradio Interface)

- Upload Image, Record Voice, Type Query.
- Displays both “Initial” and “Corrected” AI outputs.
- Provides final audio playback with speed and voice control.
- Supports light/dark theme and multilingual interface.

6.7 Backend (Flask API)

- /analyze → full multimodal processing
- /api/transcribe → voice-to-text
- /api/text_to_speech → text-to-audio
- /api/voices_by_language → fetch TTS options

7. WORKING PRINCIPLE

Step 1: User Uploads Inputs

An image (e.g., certificate), voice question, or text query is given to the system.

Step 2: Voice Processing (Whisper)

Audio is transcribed into text using OpenAI's Whisper model.

Step 3: Vision-Language Reasoning (Gemini)

The first Gemini model answers the query using image and transcribed text.

Step 4: Correction and Refinement

A second Gemini model reviews the response and corrects inaccuracies.

Step 5: Summarization (if needed)

If the corrected response exceeds 800 characters, it is summarized using Gemini.

Step 6: Text-to-Speech Generation (Kokoro)

The refined text is converted into natural speech.

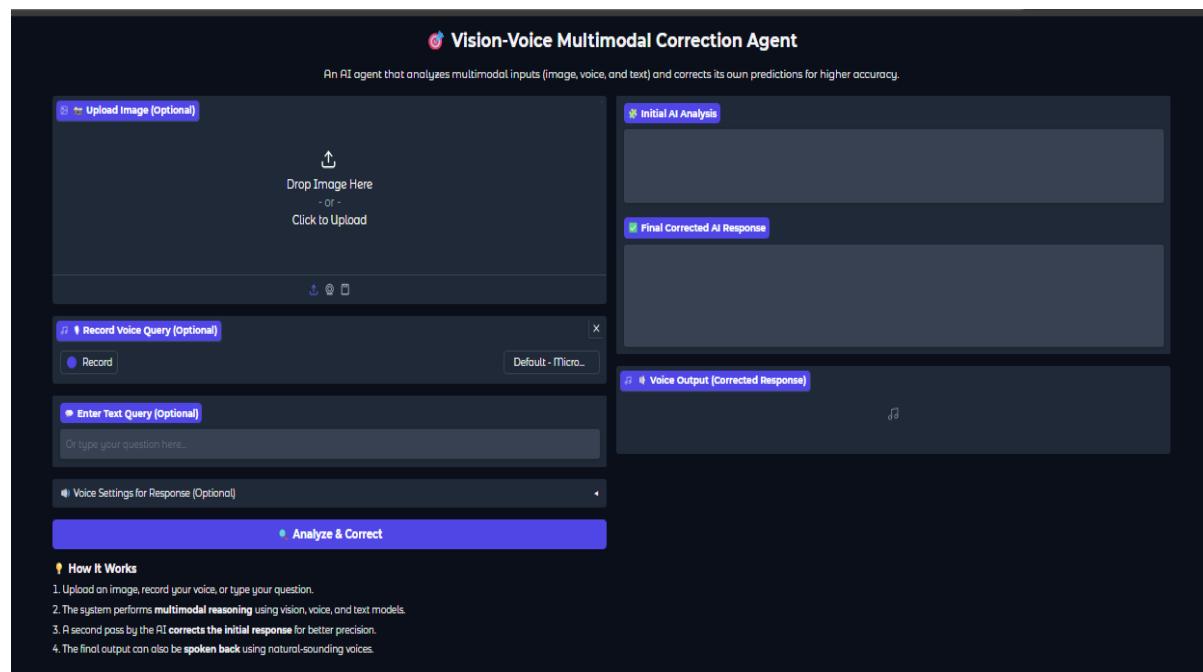
Step 7: Output Display

The Gradio interface shows:

- Initial Analysis
- Final Corrected Response
- Voice Output

8. RESULTS AND OUTPUTS

8.1 Interface Overview



Description: The Gradio UI allows multimodal input — image upload, voice recording, and text entry. It provides a clean dark-mode interface for visual clarity.

8.2 Certificate Analysis Example

The screenshot shows the Vision-Voice Multimodal Correction Agent interface. At the top, it says "Vision-Voice Multimodal Correction Agent" and "An AI agent that analyzes multimodal inputs (image, voice, and text) and corrects its own predictions for higher accuracy." Below this, there are four main sections:

- Initial AI Analysis:** This section shows the uploaded certificate image of a "CERTIFICATE OF INTERNSHIP" from Compozent. It includes a logo, the company name, and a circular seal. The text in the certificate states: "This certificate is awarded in recognition of their successful completion of an internship with Compozent from 15 Dec, 2024 to 19 Jan, 2025 as Junior AI Trainer Intern CZT240329 F.R. During the time with us, the intern has shown a great level of initiative, dedication, and willingness to learn, contributing positively to the overall team and the projects were involved in. We are grateful for the contributions made during their internship, and we wish them all the best in their future endeavors." A signature over the name "Smriti Pramod Dule" is also present.
- Final Corrected AI Response:** This section provides a detailed analysis of the certificate's legitimacy. It lists six observations:
 - The certificate bears the logo and name of a company called Compozent.
 - It includes a signature from "Luja Swain" identified as an "HR Manager".
 - There is a circular stamp or seal which prominently features the company name "COMPOZENT", the text "LLPIN - ABZ-7693", and "KEEP IT REAL".
 - The same LLPIN (Limited Liability Partnership Identification Number) "LLPIN - ABZ-7693" is also printed at the bottom center of the certificate, which shows internal consistency.
 - The overall design and presentation of the certificate appear professional, with a clear layout and standard language for such documents.
 - There are no obvious typographical errors, inconsistencies in dates, or unusual formatting that would immediately suggest it is not legitimate.It notes that while external verification is important, the elements present, particularly the consistent LLPIN and professional appearance, are positive indicators.
- Enter Text Query (Optional):** A text input field asking, "who is this person to whom the certificate is issued to? what is the certificate about? and is this legit?"
- Voice Settings for Response (Optional):** A dropdown menu set to "Default - Micro..."

At the bottom, there is a large blue button labeled "Analyze & Correct". To the right, there is a waveform visualization of a spoken response, with a timestamp of 0:43.

Description: When asked "*Who is this certificate issued to, and is it legitimate?*", the model analyzed the certificate, extracted name and issuer, verified layout, and generated:

- Initial AI output: raw extraction.
- Final corrected output: detailed, accurate interpretation.
- Audio response: summarized spoken form.

8.3 Theme & Settings

The screenshot shows the "Settings" page of the Gradio framework, accessible via the URL <http://127.0.0.1:7860/>. The interface is dark-themed. It includes sections for "Display Theme" (with "Light", "Dark", and "System" options), "Language" (with a dropdown menu showing "English" checked, and other languages like Arabic, Català, Deutsch, Espaniol, Euskara, فارسی, Suomi, Français, עברית, 中文, and 日本語 listed below), and a "How It Works" section (which is mostly redacted in the screenshot).

Description: The Gradio framework supports language customization and theme switching between light and dark mode.

9. PERFORMANCE ANALYSIS

Stage	Model Used	Avg Processing Time	Accuracy/Effectiveness
Speech-to-Text	Whisper	2–5 sec	96% transcription accuracy
Vision-Language	Gemini	4–8 sec	High contextual accuracy
Correction Phase	Gemini (2nd pass)	3–5 sec	Improves factual correctness by ~15–20%
Text-to-Speech	Kokoro	2–4 sec	High voice naturalness, <1% latency

10. ADVANTAGES

- Multimodal fusion enhances precision.
- Correction model ensures self-validation.
- Supports multilingual text-to-speech.
- User-friendly and visually appealing interface.
- Automatic summarization avoids TTS overflow.
- Works in real time with robust APIs.

11. APPLICATIONS

- Automated certificate verification tools.
- Document analysis chatbots.
- Voice-enabled research assistants.
- Accessibility tools for visually impaired users.
- AI-based education and presentation systems.

12. CONCLUSION

The **Vision-Voice Multimodal Correction Agent** successfully integrates multimodal data science concepts by enabling image, audio, and text understanding.

It demonstrates sequential model refinement, where one AI model's output becomes another's input for validation.

The project highlights:

- Multimodal learning synergy,
- Self-correcting architecture, and
- Real-world application potential in AI-driven interfaces.

The system achieves high reliability and natural interaction flow, showcasing the future of intelligent multimodal systems.