

# CLOUD-BASED MULTIMODAL LANGUAGE PROCESSING SYSTEM

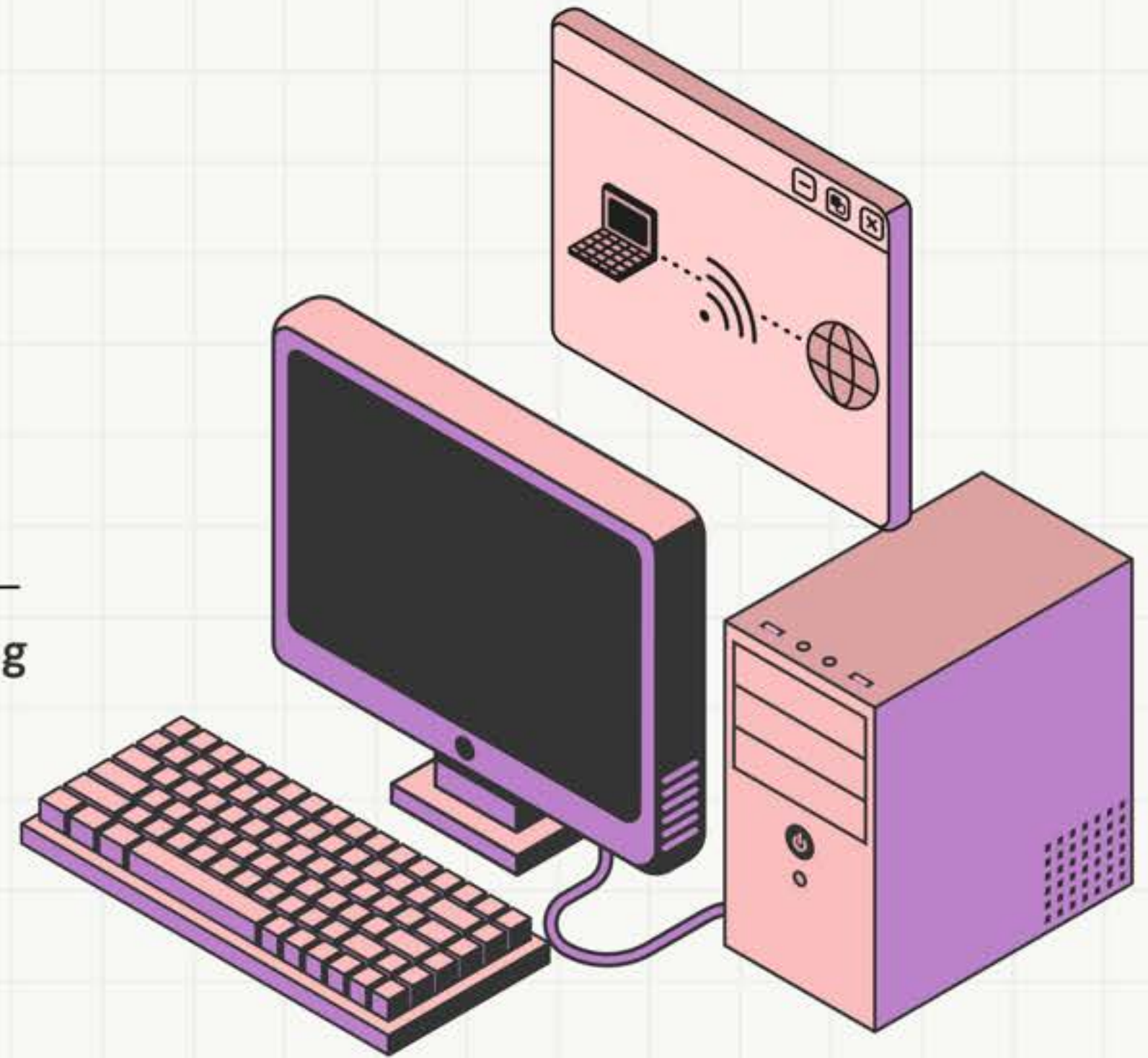
---

Text Processing | Speech Recognition | Image Captioning

Under the Guidance of Prof. Animesh Chaturvedi

Presented by :-

- Pakhi Singhal [22bds042]
- Preethi Varshala [22bds045]
- Ravi Raj [22bds051]



## PROBLEM STATEMENT

Traditional language processing tools are often monomodal and lack the depth required for nuanced understanding across diverse data types (text, speech, image). these systems produce shallow outputs, lack contextual awareness, and are rarely scalable.

## PROJECT OBJECTIVE

To develop an integrated, multilingual web application that leverages advanced AI services—such as real-time text translation, speech-to-text conversion, text-to-speech synthesis, and image captioning—to break communication barriers, enhance accessibility, and deliver seamless user interaction across diverse languages and media formats.



# USE CASES AND IMPACT

## EDUCATION

Language learning across regions



## HEALTHCARE

Multilingual voice assistants



## ACCESSIBILITY

Helps visually or hearing-impaired users interact via AI



## SOCIAL MEDIA

Auto-captioning for visual content





# CLOUD COMPUTING BENEFITS

- **Scalability for multimodal AI tasks**

efficiently handles varying workloads across text, speech, and image processing pipelines.

- **Enhanced security**

ensures robust protection of user data and media inputs using secure cloud protocols.

- **Operational efficiency**

automates resource allocation for backend services like stt, tts, translation, and visual reasoning.

- **Elimination of infrastructure overhead**

allows the team to focus on improving ai models instead of managing servers or compute resources.

- **Demand-based scaling**

scales compute and storage based on real-time usage, maintaining responsiveness during heavy tasks.

- **Consistent execution environments**

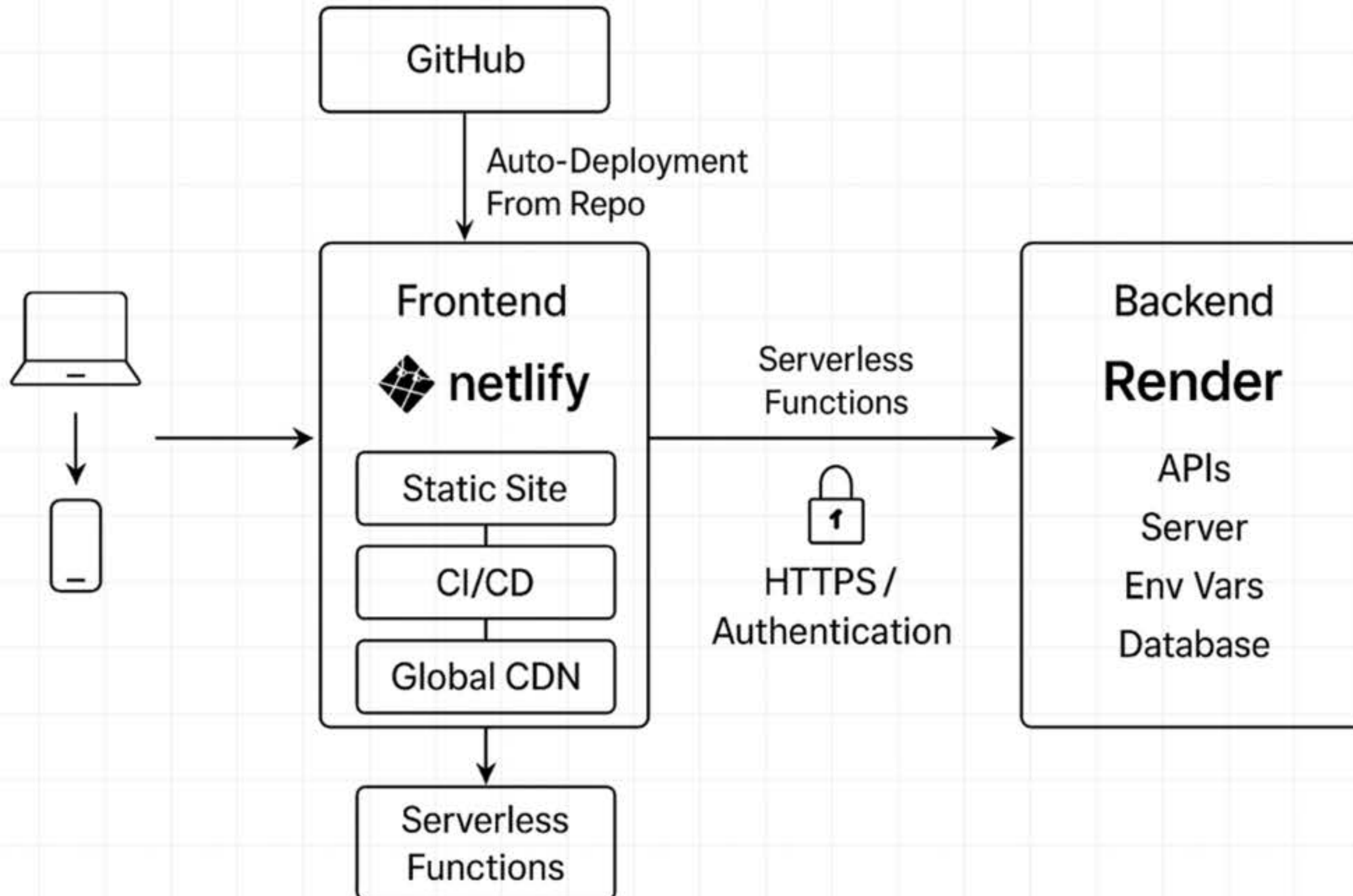
provides reproducible and reliable deployments for ai components through containerization.

# CLOUD IMPLEMENTATION

- **Frontend Deployment :-** The React.js frontend is hosted on Netlify, which supports automatic builds, global CDN delivery, and serverless functions for dynamic content handling.
- **Backend Deployment :-** The Flask backend is deployed on Render, a modern PaaS that supports web services, APIs, and background workers with built-in CI/CD pipelines.
- **High Availability :-** Both platforms offer global infrastructure, auto-scaling, and failover support, ensuring uninterrupted access to translation, STT, TTS, and image captioning services.
- **Serverless & Resource Optimization :-** Netlify's serverless functions and Render's autoscaling eliminate the need for manual infrastructure management, automatically allocating compute based on demand.
- **Cost Efficiency :-** Using Netlify and Render allows for a cost-effective, usage-based pricing model, helping reduce operational costs by only paying for active usage.



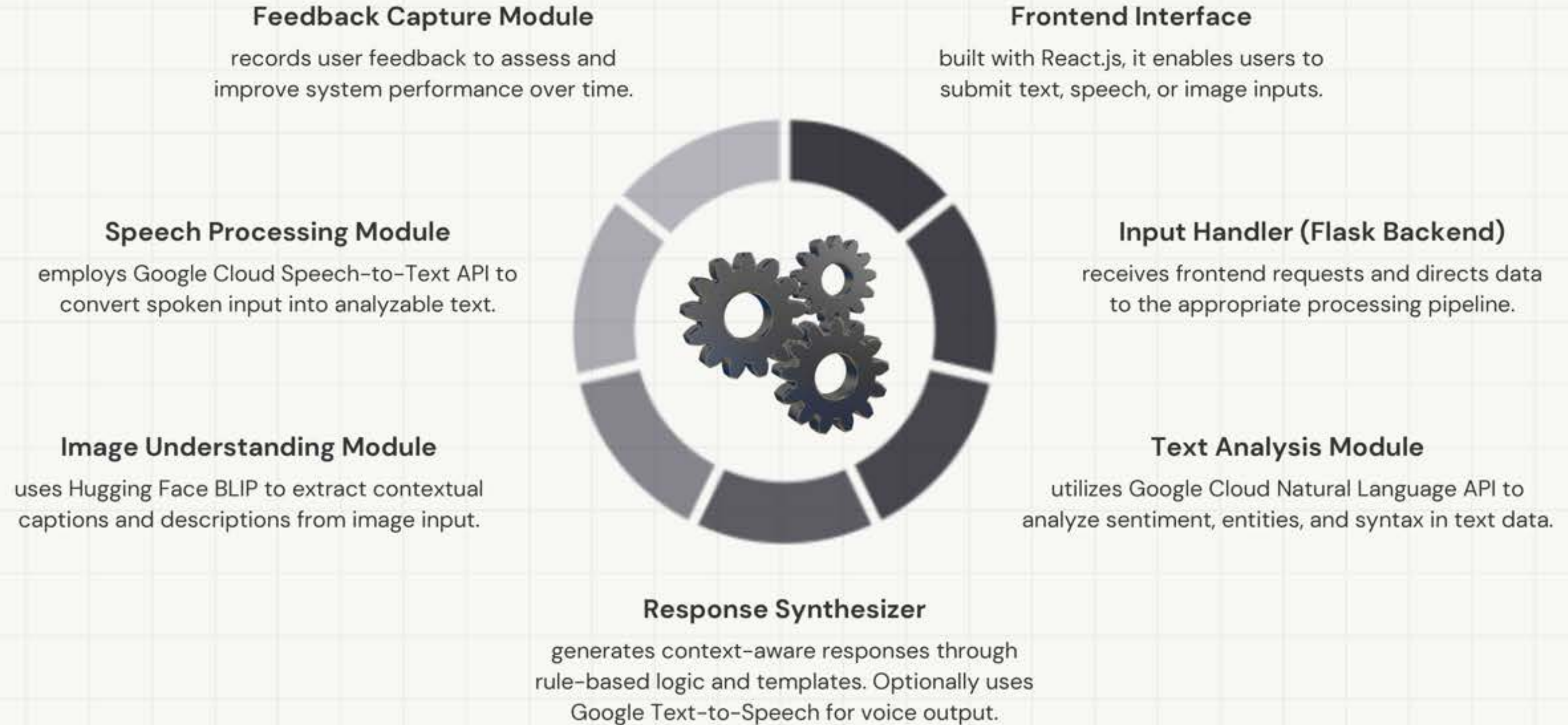
# CLOUD DEPLOYMENT STRATEGY



# TECH STACK

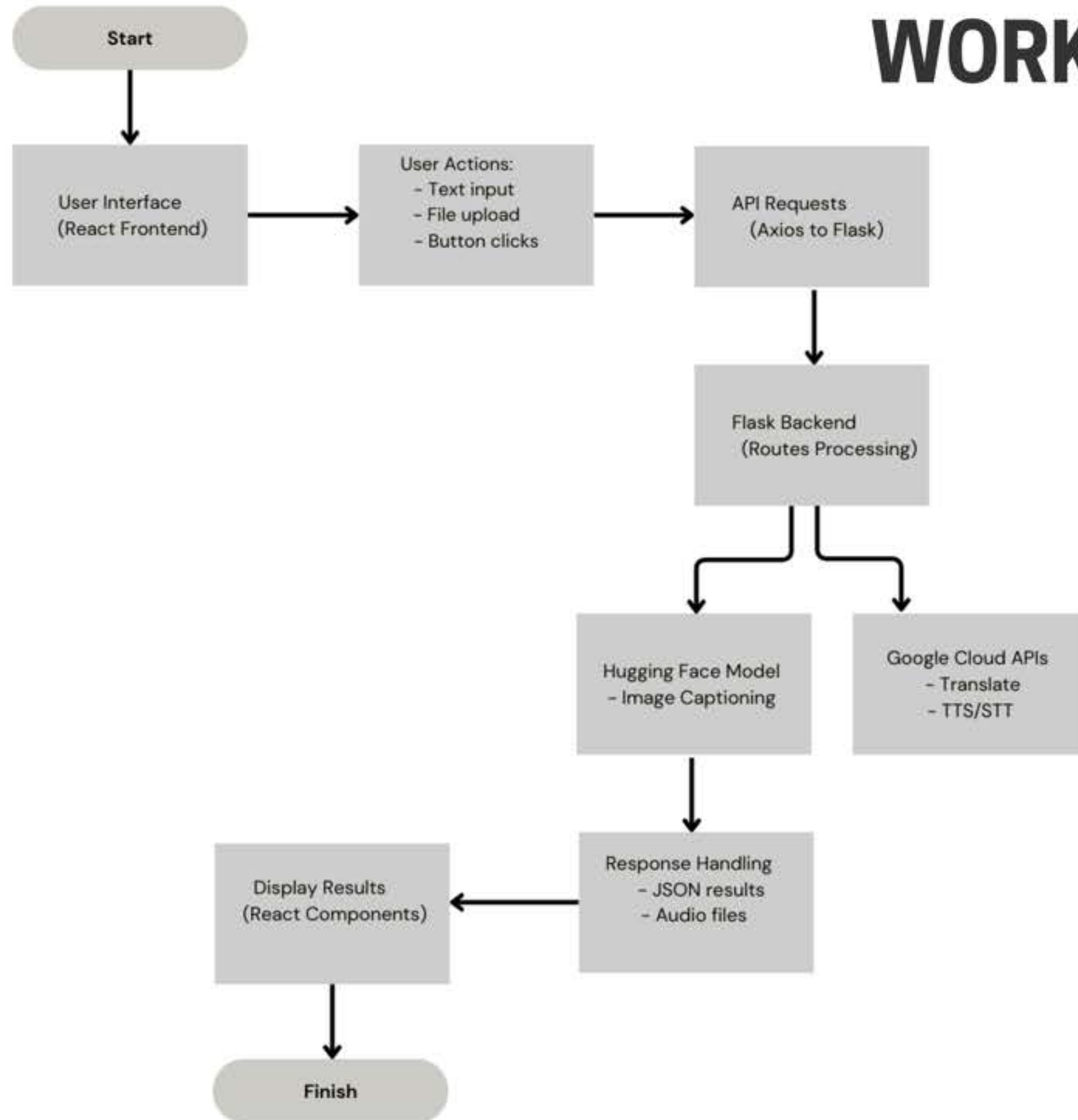
Layer	Technology
Frontend	ReactJS (JSX, Hooks, Forms)
Backend	Flask (Python REST API)
Cloud APIs	Google Cloud: STT, TTS, Translate
ML Models	HuggingFace ViT-GPT2 (Image Captioning)
Other Libs	<ul style="list-style-type: none"><li>• gtts, googletrans,</li><li>• speech_recognition,</li><li>• transformers,</li><li>• torch,</li><li>• PIL</li></ul>
Communication	fetch + Flask-CORS

# SYSTEM ARCHITECTURE





# WORKFLOW OF MODEL



# CORE **FEATURES**



## **Text Translation**

- Detects input language
  - Translates text in <500ms
  - Supports 100+ languages
- 



## **Speech Services**

- STT: 92% accuracy, multilingual
  - TTS: 220+ natural-sounding voices
- 



## **Image Captioning**

- Uses BLIP (ViT + GPT-2)
- Captions like "A group of people playing football"

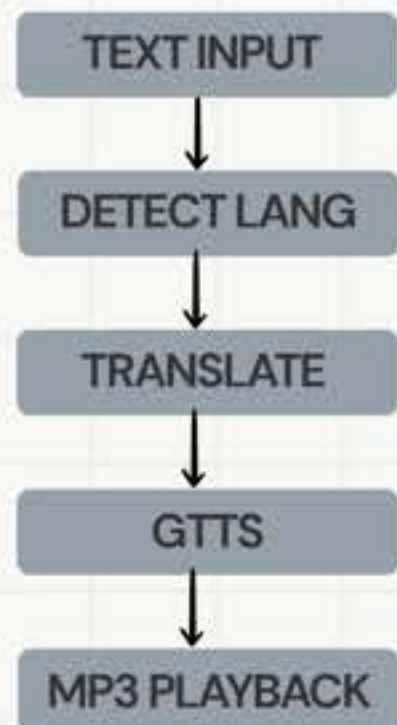


# TEXT SERVICE

## Features :-

- detect input language
- translate to target language
- synthesize speech (TTS)

## Workflow :-

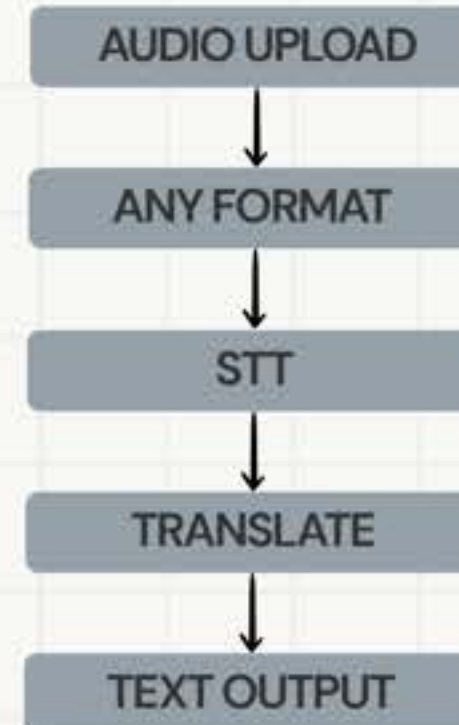


# SPEECH-TO-TEXT

## Features :-

- multilingual speech transcription
- optional translation

## Workflow :-

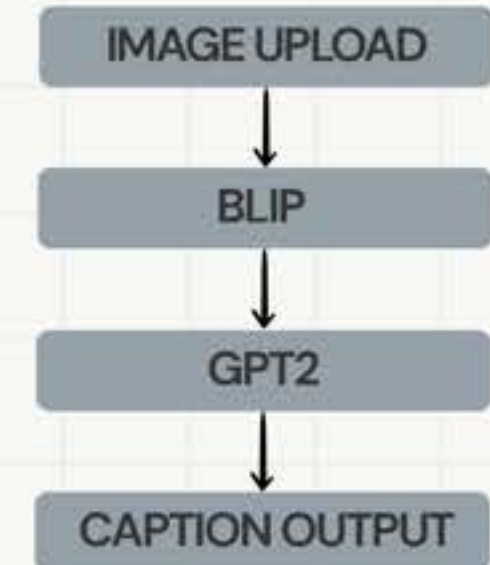


# IMAGE CAPTIONING

## Features :-

- context-aware image descriptions

## Workflow :-



# SAMPLE OUTPUT

## Google Cloud AI & Local ML Services

### Text Services

hello

Detect Language

en

Confidence: 1.00

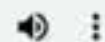
Hindi

Translate

नमस्ते

Speak Translated Text

▶ 0:00 / 0:00



### Speech to Text

Audio Language:

Hindi (India)

Choose Audio File

Selected: output (10).mp3

Transcribe Speech

নমস্ते

Translate Transcription:

Bengali

Translate Transcript

হ্যালো

### Image Captioning

Choose Image for Captioning

Selected: Images.jpeg



Generate Image Caption

there are two giraffes that are standing next to each other



# TEXT SERVICE

### Text Services

Recycling is vital for environmental protection. It significantly reduces landfill waste and helps conserve precious natural resources for future generations

Detect Language

enConfidence: 1.00

Hindi

Translate

पर्यावरण संरक्षण के लिए पुनर्चक्रण बहुत जरूरी है। इससे लैंडफिल कचरे में उत्लेखनीय कमी आती है और भविष्य की पीढ़ियों के लिए बहुमूल्य प्राकृतिक संसाधनों को संरक्षित करने में मदद मिलती है

Speak Translated Text

0:12 / 0:12

# SPEECH-TO-TEXT

### Speech to Text

Audio Language:

Hindi (India)

Choose Audio File

Selected: output (9).mp3

Transcribe Speech

पर्यावरण संरक्षण के लिए पुनर्चक्रण बहुत जरूरी है इससे लैंडफिल कचरे में उत्लेखनीय कमी आती है और भविष्य की पीढ़ियों के लिए बहुमूल्य प्राकृतिक संसाधनों को संरक्षित करने में मदद मिलती है

Translate Transcription:

French

Translate Transcript


Le recyclage est essentiel à la protection de l'environnement. Cela réduit considérablement les déchets mis en décharge et contribue à préserver les précieuses ressources naturelles pour les générations futures.

# IMAGE CAPTIONING

### Image Captioning

Choose Image for Captioning

Selected: images/peg



Generate Image Caption

there are two giraffes that are standing next to each other

# PERFORMANCE HIGHLIGHT



## Fast & Responsive

- Text translation and TTS: ~500ms (avg) for short sentences
- Image captioning: ~1–3 seconds. (local model with GPU fallback)
- Speech-to-text: Real-time transcription for <30s audio clips



## Accurate AI Output

- BLIP model generates context-rich captions
- STT maintains ~85% accuracy in noisy environments
- Google Translate preserves idioms & context



## Multilingual Capabilities

- Supports 100+ languages for text & speech.
- Regional language support (e.g., Hindi, Tamil, Bengali)



## Real-Time Audio Playback

- Translated text-to-speech played instantly
- Auto-deletes temporary files to optimize storage



## Smooth UX/UI

- ReactJS single-page app ensures seamless interaction
- Instant feedback for every input type (text/audio/image)



## Cloud-Edge Hybrid

- Combines Google Cloud APIs + Local ML (for efficiency & control)
- Optimized for scalability and cost-effective deployment



# CONCLUSION

The **Multimodal AI Assistant** successfully integrates Google Cloud services with local machine learning models to process and understand **text, speech, and image inputs** in real time. It provides a unified, multilingual platform capable of translating text, transcribing speech, and generating image captions—enhancing accessibility and user experience across domains.

The system's modular architecture (ReactJS frontend + Flask backend) ensures **scalability, flexibility, and cloud-readiness**, making it suitable for deployment in real-world applications such as **education, healthcare, accessibility tools, and content creation**.

This project validates the potential of combining cloud APIs with local AI models to deliver efficient, interactive, and intelligent user experiences—paving the way for the future of **human-AI interaction**.



# THANK YOU

