

## A) Introduction

This document outlines the findings from a comprehensive analysis of technologies for the AI-Powered Talking Avatar MVP. The primary goal was to identify a toolchain that delivers a high-quality user experience while minimizing engineering complexity and development time.

After evaluating multiple options, the recommended stack is **ElevenLabs for voice cloning** and **HeyGen for avatar generation**.

## B) Recommended Stack & Justification

### 1. Voice Synthesis: ElevenLabs

- **State-of-the-Art Quality:** ElevenLabs is widely recognized for producing the most realistic and emotionally expressive AI-generated speech.
- **Instant Voice Cloning:** With just a few minutes of audio, we can rapidly prototype the desired voice.
- **Developer-First API:** The platform's robust API allows seamless integration into our backend, minimizing engineering overhead compared to hosting open-source models.

### 2. Avatar Generation: HeyGen

- **API-Driven Workflow:** HeyGen's API can programmatically generate lip-synced videos from an audio file and an image.
- **High-Quality 2D Avatars:** Produces realistic avatars with natural head movements and expressions.
- **MVP-Friendly:** Provides a polished result without requiring 3D rendering or GPU-heavy infrastructure.

## C) Analysis of Rejected Alternatives

### Why Not Synthesia?

- Optimized primarily for **pre-recorded corporate/training videos** rather than live conversational use.

- API is less developer-focused compared to HeyGen.
- While capable, it introduces more friction for rapid prototyping.

### Why Not MetaHuman?

- MetaHuman is a **3D character creation tool**, not an end-to-end avatar service.
- Requires:
  - Unreal Engine integration
  - Real-time lip-sync pipelines (e.g., NVIDIA Audio2Face)
  - GPU servers and streaming infrastructure
- This results in **massive engineering overhead and cost**, unsuitable for a lean MVP.

### D) Core Technical Decisions & Strategy

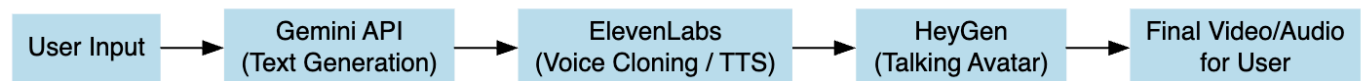
- **Architecture:** Adopting a “*Generate-and-Play*” request-response model instead of real-time streaming to simplify implementation.
- **Latency:** Accepting ~5–10 seconds per interaction:
  - Gemini API: ~0.5–1s
  - ElevenLabs: ~1–2s
  - HeyGen: ~3–7s
- **Backend Orchestration:** Smooth chaining of the pipeline: (*User Input* → *Gemini* → *ElevenLabs* → *HeyGen* → *User*).
- **Minimalist Frontend:** Standard HTML/CSS/JS, keeping focus on backend and AI integration.

### E) Ethical, Cost, and Scalability Considerations

- **Ethics:** Voice cloning requires explicit permission. Licensed/stock assets will be used for demos.
- **Cost:** Both HeyGen and ElevenLabs offer free tiers adequate for MVP testing. Scaling beyond MVP requires paid plans.

- **Scalability:** The chosen stack is MVP-friendly but can later evolve into more advanced pipelines (e.g., MetaHuman + Unreal for ultra-realism).

Tool	Strengths	Weaknesses	MVP Suitability
HeyGen	Easy API, realistic avatars, fast rendering	2D only, small latency	✓ Best choice
ElevenLabs	High-quality TTS, instant cloning, streaming	Paid beyond free tier	✓ Best choice
Synthesia	Stable, polished, widely trusted	Pre-recording focus, higher latency	Δ Less suitable
MetaHuman	Ultra-realism, high control	GPU heavy, costly, complex pipeline	✗ Not suitable



## F) Free Tier vs Paid Tier Strategy

- For the **MVP demonstration**, we will leverage the **free tiers** of HeyGen and ElevenLabs.
  - This is sufficient to test functionality end-to-end, though the free tier introduces **higher latency (up to ~60 seconds per response cycle)** due to queueing and rate limits.
  - While this latency is not ideal for live interaction, it is acceptable at the MVP stage since the goal is to validate the pipeline and user experience flow.
- Once the MVP proves successful (ignoring the latency factor), we can upgrade to **paid plans**.
  - Paid tiers support **real-time or near real-time streaming** with significantly reduced latency (~1–3 seconds).
  - This will allow us to evolve from a proof-of-concept demo into a **production-ready, live conversational avatar**

## G) Conclusion

The proposed stack of **ElevenLabs + HeyGen** is the result of rigorous evaluation and offers the most direct, efficient, and effective path to launching a high-quality MVP.

Synthesia was excluded due to higher latency and pre-recording focus, while MetaHuman was deemed impractical due to engineering and infrastructure overhead.

This approach balances **realism, speed, and ease of implementation**, while remaining extensible for future research into real-time and 3D pipelines.