
AI LEARNS TO TALK

THE \$47B VOICE REVOLUTION

From chatbots to human-like
voice agents in every industry

\$47.5B
by 2034

**Voice AI Agents Market
Growing at 35% CAGR**

From \$2.4B in 2024 → 20x in 10 years

Voice-First Computing



60% of Smartphone Users Use Voice

Up from 45% in 2023. Voice assistants now standard in phones, speakers, cars, and wearables. Hands-free is the new default.



90% of Hospitals Using AI Agents by 2025

Healthcare leads adoption. Voice AI handles scheduling, triage, patient follow-ups, and medication reminders at scale.



\$290B Conversational Commerce by 2025

Retail and e-commerce drive massive investment. Voice shopping, support bots, and AI assistants transform customer experience.

Voice-to-Voice AI Arrives



End-to-End Audio Models

No more text transcription. GPT-4o and EVI 2 process audio directly → audio output. Sub-500ms latency. Natural conversation finally possible.



Emotional Intelligence

AI detects tone, emotion, and intent from voice patterns. Responds with appropriate inflection. Empathic interfaces emerging.



Multilingual by Default

30+ languages supported. Real-time translation. Voice AI learns languages organically from training data.

Call Centers Go AI



87% Frustrated with Human Transfers

AI agents eliminate hold times and transfers. 24/7 availability. Consistent quality. Customers actually prefer it for simple tasks.



BFSI Leads with 32.9% Market Share

Banking, insurance, and finance deploy voice AI for fraud alerts, account inquiries, claims processing, and compliance calls.



\$0.07/min vs \$1+/min Human Agents

10-15x cost reduction. Voice AI platforms like Retell and Vapi handle thousands of concurrent calls. ROI in weeks.

||

Voice-to-voice models bring to fruition what Siri and Alexa had long promised — natural conversation with AI.

— Hume AI Research

Key Insight: OpenAI dropped Realtime API pricing 60-87% in Dec 2024. Cost barrier evaporating. Every app can now have a voice interface — the question is who builds the best experience.

When did you last use a voice assistant?



Siri, Alexa, ChatGPT Voice, or something else?

KEY NUMBERS

\$3.3B

ElevenLabs valuation — voice synthesis leader

1,000

Years of AI audio generated annually on ElevenLabs

75ms

ElevenLabs Flash model latency — faster than thought

\$3.3B

VALUATION

\$90M+

ARR (2024)

\$281M

TOTAL RAISED

60%

OF FORTUNE 500

Polish-founded voice AI powerhouse. Ultra-realistic text-to-speech in 30+ languages. Voice cloning from short samples. The "OpenAI for audio."

🔥 **Jan 2025:** \$180M Series C led by a16z. Voice Library marketplace pays creators \$2M+ in royalties. Disney Accelerator alum. ElevenLabs Reader app launched.

OA

OpenAI Voice

CHATGPT ADVANCED VOICE MODE

GPT-4o

MODEL

9

VOICE OPTIONS

Realtime

API (OCT 2024)

60-87%

PRICE DROP DEC '24

Native multimodal voice in ChatGPT. End-to-end audio-to-audio processing. Memory and custom instructions supported. The mainstream benchmark.

 **Dec 2024:** Realtime API pricing slashed — now accessible to all developers. Screen sharing and video coming. Free users get GPT-4o-mini voice daily.

HU

Hume AI

EMPATHIC VOICE INTERFACE

\$50M

SERIES B (2024)

EVI 2

LATEST MODEL

\$4.32/hr

API COST

Emotion

DETECTION

Founded by ex-Google DeepMind scientist Alan Cowen. Trained on cross-cultural emotional data. Detects and responds to user emotions from voice patterns.

🔥 **EVI 2:** First truly empathic voice AI. End-to-end audio model. 2x cheaper than OpenAI Realtime. Optimized for human well-being and mental health applications.

DG

Deepgram

SPEECH RECOGNITION LEADER

\$86M+
TOTAL RAISED

Nova-3
LATEST MODEL

200K+
DEVELOPERS

54%
LOWER WER

Enterprise-grade speech-to-text. 40x faster than competitors.
Processes 50,000+ years of audio annually. Sub-300ms latency for
real-time transcription.

 **Nova-3:** Trained on 47B tokens from 6M+ sources. Best-in-class accuracy for call centers, meetings, and voice analytics. Aura TTS for real-time conversations.

RE

Retell AI

ENTERPRISE VOICE AGENTS

\$0.07
PER MINUTE

31+
LANGUAGES

HIPAA
SOC2 / GDPR

<800ms
LATENCY

Developer-first platform for production voice agents. Full control over conversation logic. Built for healthcare, finance, and compliance-heavy industries.

 **Enterprise Ready:** Automatic PII redaction. Verified phone numbers reduce spam flags. Knowledge base integration for accurate answers. Cal.com scheduling built-in.

VA

Vapi
OPEN SOURCE VOICE SDK

\$0.05
PER MINUTE BASE

<500ms
LATENCY

Open
SOURCE

Multi
LLM SUPPORT

Open-source voice agent SDK for developers who want maximum customization. Supports GPT-4o, Claude, and custom models. Self-host or use their cloud.

 **Developer Favorite:** WebSocket streaming. Bring your own LLM. Thousands of configurations possible. Popular for rapid prototyping and custom telephony.

End-to

END INFRA

No-Code

BUILDER

Voice

CLONING

Enterprise

SCALE

Full-stack AI phone agent platform. Owns entire infrastructure for lowest latency. No-code builder plus API for developers. Built for enterprise scale.

 **Differentiated:** Context memory across calls. Built-in summarization and confidence scoring. Manages all CRM and telephony integrations. Enterprise privacy controls.

AG

Alexa + Google

CONSUMER VOICE GIANTS

500M+
ALEXA DEVICES

1B+
GOOGLE ASSISTANTS

Gemini
LIVE (GOOGLE)

LLM
ALEXA UPGRADE

The incumbents with massive installed bases. Both racing to integrate LLMs. Google launched Gemini Live. Amazon reportedly partnering with Anthropic for new Alexa.

🔥 **Playing Catch-up:** Legacy assistants struggled vs ChatGPT. Now investing heavily. Samsung Bixby also upgraded in 2024. Smart home dominance at stake.

PlayAI + SoundHound

VOICE INFRASTRUCTURE

Meta

ACQUIRED PLAYAI

\$177M

SOUNDHOUND '25 REV

\$1.2B

SH BOOKINGS

\$140B

TAM

PlayAI acquired by Meta to power voice for Meta AI and future products. SoundHound (public) sees massive automotive and restaurant demand.

 **Consolidation:** Big tech buying voice infrastructure. CB Insights flags ElevenLabs and Cartesia as top M&A targets. Voice is strategic, not commodity.

\$115M+
ASSEMBLYAI RAISED

Ultra
LOW LATENCY

Sonic
CARTESIA TTS

5-10x
FASTER TTS

AssemblyAI: Universal speech transcription and audio intelligence.
Cartesia: Ultra-low latency TTS optimized for real-time voice agents. Both strategic acquisition targets.

🔥 **Speed Race:** Cartesia's Sonic model 5-10x faster than competitors. AssemblyAI's Universal-1 handles any audio. Racing to sub-50ms response times.

FINAL THOUGHTS

The Voice-First Future



Every App Gets a Voice

APIs now affordable. Voice interfaces becoming table stakes. Expect voice-first experiences in banking, healthcare, retail, and productivity apps.



AI + Human Hybrid Models

Voice AI handles 80%+ of routine queries. Humans focus on complex, high-value interactions. Smith.ai and others blend both for best experience.



Emotional AI Emerges

Beyond words — understanding tone, sentiment, intent. Mental health, eldercare, and companionship apps on the rise. Voice as empathic interface.

 FOLLOW FOR MORE AI INSIGHTS

JJ Shay

bit.ly/jjshay