

Full Cost Breakdown: 15-Minute Emotionally Fluent Voice AI Conversation

Conversation Assumptions

Parameter	Value	Notes
Total duration	15 minutes	Similar to your Maya Conversation #3
Human speaking time	~7.5 minutes	50% of conversation
AI speaking time	~7.5 minutes	50% of conversation
Speaking rate	150 words/minute	Average conversational pace
Words spoken (each party)	~1,125 words	7.5 min × 150 words
Characters (AI output)	~5,625 characters	~5 chars/word
Exchanges	~18-20 turns	Back-and-forth dialogue

Component 1: Speech-to-Text (ASR)

Transcribing human speech so the AI can "hear" you

Provider	Rate	7.5 min cost
Deepgram Nova-3 (streaming)	\$0.0077/min	\$0.058
Deepgram Nova-3 (batch)	\$0.0043/min	\$0.032
OpenAI Whisper API	\$0.006/min	\$0.045
Google Cloud STT	\$0.016/min	\$0.120

Selected estimate: \$0.058 (Deepgram streaming for real-time)

Component 2: LLM Inference ("The Brain")

Generating intelligent, contextual responses

Token Estimation

- System prompt + personality: ~500 tokens (one-time)
- Conversation context accumulates: ~2,000 tokens average
- Per-turn input: ~150 tokens (user message + context window management)
- Per-turn output: ~100 tokens (AI response)
- Total across 20 turns:
 - **Input tokens: ~8,000** (includes growing context)
 - **Output tokens: ~2,000**

Model	Input Rate	Output Rate	Total Cost
Claude 3.5 Sonnet	\$3/1M	\$15/1M	\$0.054
Claude Haiku 3.5	\$0.80/1M	\$4/1M	\$0.015
GPT-4o	\$5/1M	\$20/1M	\$0.080
GPT-4o-mini	\$0.15/1M	\$0.60/1M	\$0.002
Gemini 2.5 Flash	\$0.15/1M	\$0.60/1M	\$0.002

Selected estimate: \$0.054 (Claude Sonnet for quality conversational AI)

Component 3: Text-to-Speech (TTS)

Converting AI text responses to natural speech

Provider	Rate	5,625 chars cost
Qwen3-TTS-Flash	\$0.011/1K chars	\$0.062
Hume Octave 2	\$0.063/1K chars	\$0.354
OpenAI TTS-1	\$0.016/1K chars	\$0.090
OpenAI TTS-1 HD	\$0.032/1K chars	\$0.180
ElevenLabs v3	\$0.189/1K chars	\$1.064
Deepgram Aura 2	\$0.032/1K chars	\$0.180

Selected estimate: \$0.062 (Qwen3-TTS for efficiency breakthrough referenced in article)

Component 4: Emotional Intelligence Layer (Optional)

What makes it "feel" human — the Hume-style layer

This is where Maya's magic lives. Two approaches:

Option A: Expression Measurement (analyzing user emotion)

Hume Feature	Rate	7.5 min cost
Audio-only analysis	\$0.0639/min	\$0.479
With video	\$0.0828/min	\$0.621

Option B: Emotionally Intelligent Voice Generation

Using Hume's EVI (Empathic Voice Interface) instead of standard TTS:

- Replaces Component 3
- Adds emotional responsiveness
- Estimated: \$0.30-0.50/minute = **\$2.25-\$3.75** for 7.5 min

Option C: Integrated Speech-to-Speech (emerging models)

Models like Nvidia PersonaPlex or FlashLabs Chroma handle listening + speaking in one model:

- Eliminates separate ASR + TTS
- Estimated: **\$0.10-0.20/minute = \$1.50-\$3.00** for full conversation

Selected estimate (basic emotional layer): \$0.479 (Hume audio analysis)

Total Cost Scenarios

Scenario 1: Budget Stack (functional but not emotionally fluent)

Component	Provider	Cost
ASR	Deepgram Nova-3 batch	\$0.032

Component	Provider	Cost
LLM	GPT-4o-mini	\$0.002
TTS	Qwen3-TTS-Flash	\$0.062
Emotion	None	\$0.00
TOTAL		\$0.096

~10 cents for 15 minutes. Functional chatbot, not Maya-level.

Scenario 2: Quality Stack (what powers a good voice assistant)

Component	Provider	Cost
ASR	Deepgram Nova-3 streaming	\$0.058
LLM	Claude 3.5 Sonnet	\$0.054
TTS	OpenAI TTS-1 HD	\$0.180
Emotion	None	\$0.00
TOTAL		\$0.29

~30 cents. Good voice assistant, but no emotional adaptation.

Scenario 3: Maya-Level Experience (emotionally fluent)

Component	Provider	Cost
ASR	Deepgram Nova-3 streaming	\$0.058
LLM	Claude 3.5 Sonnet	\$0.054
TTS	Hume Octave 2 (emotional)	\$0.354
Emotion Analysis	Hume Expression (audio)	\$0.479
TOTAL		\$0.95

~\$1 for 15 minutes of emotionally fluent conversation.

Scenario 4: Enterprise Premium (full Hume EVI stack)

Component	Provider	Cost
Speech-to-Speech	Hume EVI (integrated)	\$2.50-4.00
LLM enhancement	Claude Sonnet (optional)	\$0.054
TOTAL		\$2.50-\$4.00

Maximum emotional intelligence, still under the price of a latte.

Cost Per Hour (At Scale)

Scenario	15-min cost	Hourly cost	Daily (8 hrs)	Monthly (160 hrs)
Budget	\$0.10	\$0.40	\$3.20	\$64
Quality	\$0.29	\$1.16	\$9.28	\$186
Maya-Level	\$0.95	\$3.80	\$30.40	\$608
Enterprise	\$3.50	\$14.00	\$112.00	\$2,240

The Ubiquity Implication

At Maya-level quality (~\$1 per 15-minute conversation):

- A company could provide 1,000 emotionally fluent conversations/day for **\$1,000/day**
- That's **\$30,000/month** for 30,000 Maya-quality interactions
- Compare to: One human customer service rep costs ~\$4,000-6,000/month for ~160 hours

At budget quality (~\$0.10 per 15-minute conversation):

- 1,000 conversations/day = **\$100/day = \$3,000/month**
- Essentially free at enterprise scale

What The VentureBeat Breakthroughs Change

Problem Solved	Impact on Cost
Qwen3-TTS 12Hz tokenizer	Reduces TTS bandwidth/compute by 80%+
Inworld 120ms latency	Enables cheaper real-time without buffering overhead
Nvidia PersonaPlex (open)	Eliminates licensing fees for full-duplex
Hume enterprise deals	Volume pricing below published rates

Net effect: The ~\$1 Maya-level experience will likely drop to **\$0.25-0.50** within 12 months as these technologies mature.

The Human Cost Comparison

Experience Type	Cost	Duration
15-min Maya conversation	\$0.95	15 min
Cup of gas station coffee	\$1.50	5 min
Therapy session	\$150-300	50 min
Phone call with a friend	Free	Variable
Building real human trust	Immeasurable	Years

"Building trust without earning it" now costs less than a dollar.