Realtime Switch: Unified Voice AI API

Author: Apurav Chauhan, Indu Chauhan

Date: August 2025

Version: 1.0

Abstract

While exploring voice AI APIs for a real-time application, we encountered a frustrating reality: every provider had completely different event structures, connection protocols, and data formats. OpenAI's Realtime API worked one way, Google's Gemini Live worked entirely differently, and each new provider meant starting integration from scratch.

When my primary provider started showing performance issues and users began experiencing poor audio quality, I faced a dilemma that many developers know too well: Do I spend weeks integrating another provider from the ground up? What happens tomorrow when a cheaper, better provider emerges – do I go through this painful process all over again? Or do I stay locked into one provider and watch customers churn due to poor performance?

That frustration sparked a question: What if we had a unified API that abstracted away all these provider differences?

This is how the Realtime Switch Architecture was born – a comprehensive solution that enables seamless switching between voice AI providers without breaking user conversations or requiring client code changes. Developers can use any existing API style (OpenAI Realtime API or Gemini Live API specifications) with our universal connection, ensuring out-of-the-box compatibility with their current implementations. This document presents the architectural framework that transforms the chaotic landscape of voice AI integration into a simple, unified experience, solving the critical problem of vendor lock-in while maintaining optimal performance and reliability.

Table of Contents

- 1. System Overview & Problem Statement
- 2. High-Level System Architecture
- 3. Core Event Architecture
- 4. Event Extraction Pattern
- 5. Provider Event Management
- 6. Bidirectional Event Transformation
- 7. Intelligent Provider Switching
- 8. Conversation Persistence & State Management
- 9. Pipeline Orchestration
- 10. Performance Test Results
- 11. Cost Analysis and Comparison
- 12. Current Implementation Status and Future Roadmap

1. System Overview & Problem Statement

The Challenge: Voice AI Vendor Lock-in

Modern voice AI applications face critical vendor lock-in challenges when integrating with different providers:

- OpenAl Realtime API offers excellent quality but has performance inconsistencies
- Google Gemini Live provides superior latency but uses incompatible event formats
- Provider switching traditionally breaks active conversations and loses context

Our Solution: Universal Abstraction Layer

The system provides a single WebSocket API that:

- 1. Abstracts all providers behind a unified interface
- 2. Preserves existing client event formats (zero code changes required)
- 3. Switches providers dynamically based on performance metrics
- 4. Maintains conversation context during switches
- 5. Transforms events bidirectionally between provider formats

Core Value Proposition

- Developer Experience: Write once, work with any voice AI provider
- Performance Optimization: Automatic switching to best-performing provider
- Future-Proofing: Easy integration of new providers without client changes
- Context Preservation: Seamless conversations regardless of underlying provider

2. High-Level System Architecture

2.1 System Components Overview

The system consists of six primary architectural components that work together to provide seamless voice AI provider abstraction:

Core Components

- EventManager: Base class providing pub/sub architecture for all system communication
- Pipeline: Central orchestrator that coordinates all components and manages the event flow
- Switch: Intelligent provider switching system that monitors performance and triggers provider changes
- SessionManager: Captures and manages session configuration state for replay across provider switches
- Checkpoint: Conversation persistence system that maintains dialogue history across sessions
- AuthService: Centralized authentication service for secure WebSocket connections

Processing Components

- Extractors: Parse provider-specific events into standardized callbacks
- Transformers: Convert events bidirectionally between different provider formats
- Providers: EventManager implementations that handle WebSocket connections to external voice AI services

2.2 Component Interaction Flow

Component Relationships

- Pipeline acts as the central coordinator, instantiating and connecting all other components
- SessionManager sits at the entry point, capturing configuration while passing events through
- 3. Switch monitors Provider performance and triggers Pipeline reconfiguration
- 4. Extractors are embedded within Transformers to parse provider-specific event formats
- 5. Checkpoint receives transformed events to maintain persistent conversation history
- 6. Providers communicate with external services while reporting performance to Switch

2.3 System Boundaries and Interfaces

External Interfaces

- Client WebSocket: Receives events from client applications in their expected API format
- Provider WebSockets: Communicates with external voice AI services (OpenAI, Gemini)

Internal Communication

- ProvidersEvent: Universal event format used throughout the system
- Callback Registration: Clean pattern for component communication
- Performance Stats: Metrics flow from Providers to Switch for decision making

2.4 Operational Modes

Static Mode (No Switching)

 Pipeline configured with matching API style and provider (e.g., OpenAI client → OpenAI provider)

- Transformers use NoOp pattern for direct passthrough
- Switch remains inactive, no performance monitoring

Dynamic Mode (Cross-Provider)

- Pipeline configured with different API style and provider (e.g., OpenAI client → Gemini provider)
- Transformers perform bidirectional event conversion
- Switch actively monitors performance and can trigger provider changes

Hot-Swap Mode (Runtime Switching)

- Switch detects performance degradation and triggers provider change
- Pipeline reconfigures all components without service interruption
- SessionManager replays configuration to new provider
- Checkpoint continues logging in original API format

3. Core Event Architecture

3.1 Client Events Interface

The ClientEvents interface defines what clients can send TO the system:

Core Client Capabilities

- userAudio: Raw audio data from client microphone
- sessionUpdate: Configuration changes (instructions, tools, voice settings)
- toolResponse: Function call responses from client applications

Design Principles

- Minimal Interface: Only three core event types for maximum simplicity
- Provider Agnostic: Client doesn't need to know which provider is active
- Stateless: Each event is self-contained with all necessary context

3.2 Server Events Interface

The ServerEvents interface defines what the system emits TO clients:

Core Server Capabilities

- userTranscript: Real-time transcription of user speech
- responseTranscript: Al response text transcription
- responseAudio: Streamed audio response data
- toolCall: Function call requests to client
- turn: Conversation turn completion signals

Event Flow Characteristics

- Real-time Streaming: All events support delta-based streaming
- Turn Management: Explicit turn boundaries for conversation state
- Function Calling: Bidirectional tool integration support
- Error Handling: Comprehensive error and interruption signaling

4. Event Extraction Pattern

4.1 The Extraction Challenge

Each voice AI provider uses fundamentally different event structures:

- OpenAl: Flat event hierarchy with explicit type fields
- Gemini: Nested serverContent structure with implicit semantics

4.2 Callback-Based Extractor Architecture

The system uses a callback registration pattern to connect extractors with transformers. Transformers register individual callback functions with extractors for each event type, and extractors invoke these callbacks when parsing relevant events. This approach avoids circular object references and enables proper cleanup through explicit callback

deregistration.

ServerEventsExtractor Interface

Defines the contract for parsing provider events:

- onUserTranscript: Callback registration for user speech events
- onResponseTranscript: Callback for AI response transcription
- onResponseAudio: Callback for audio response data
- onToolCall: Callback for function call events
- onTurn: Callback for turn completion events
- · cleanup: Explicit memory cleanup method

4.3 Provider-Specific Extractor Implementations

OAIServerEventsExtractor

Handles OpenAI's flat event hierarchy:

- Direct Type Mapping: Uses event.payload.type for routing
- Function Call Detection: Identifies response.output_item.done with function_call type
- Turn Completion: Processes response.done events

GeminiServerEventsExtractor

Manages Gemini's complex nested structure:

- serverContent Parsing: Extracts events from nested serverContent object
- Setup Protocol: Handles Gemini-specific connection initialization
- Turn Logic: Distinguishes between generationComplete, interrupted, and turnComplete
- Model Dependencies: Accounts for model-specific event generation patterns

4.4 Client-Side Extraction

Similar pattern applies to client events:

GeminiClientEventsExtractor: Handles realtimeInput, setup, toolResponse

 OAIClientEventsExtractor: Processes input_audio_buffer, session.update, conversation.item.create

5. Provider Event Management

5.1 EventManager Base Architecture

Observer Pattern Implementation

The EventManager base class provides:

- Subscriber Management: Dynamic subscription to event streams
- Event Broadcasting: Efficient one-to-many event distribution
- Cleanup Coordination: Proper resource management
- Chain of Responsibility: Sequential event processing pipeline

Core EventManager Capabilities

- addSubscribers: Register multiple event consumers
- emitEvent: Broadcast events to all subscribers
- receiveEvent: Abstract method for event processing
- cleanup: Resource cleanup and subscriber management

5.2 Provider-Specific Implementations

OAIEventManager: OpenAl Integration

Connection Management:

- Direct WebSocket to wss://api.openai.com/v1/realtime
- Model: gpt-4o-realtime-preview-2024-12-17
- Authorization via Bearer token with OpenAl-Beta headers

Performance Monitoring:

5-second ping intervals with timestamp correlation

- Latency calculation via pong response timing
- Stats reporting to switching system

Event Processing:

- JSON parsing of WebSocket messages
- Direct event forwarding to subscribers
- Error handling and reconnection logic

GeminiEventManager: Google Gemini Live Integration

Connection Protocol:

- WebSocket to wss://generativelanguage.googleapis.com
- Critical Setup Requirement: First message MUST be setup with models/gemini-2.0-flash-live-001
- Model Dependency: Only gemini-2.0-flash-live-001 generates proper generationComplete events

Setup Protocol Enforcement:

- Single setup message per connection (multiple setup messages cause errors)
- Wait for setupComplete before sending other messages
- Connection-scoped configuration (cannot be updated like OpenAI sessions)

Performance Integration:

- Ping/pong latency monitoring
- Stats callback integration for switching decisions
- Connection state management

5.3 Provider Connection Management Interface

ProviderConManager Contract

Enables performance monitoring and intelligent switching:

- onConnected: Callback for connection establishment
- onConnectionStats: Performance metrics reporting
- sendPing: Active latency measurement

Integration with Switching System

- Automatic registration of performance callbacks
- Real-time latency data collection
- Connection state change notifications

6. Bidirectional Event Transformation

6.1 Transformation Architecture Overview

The system implements bidirectional event transformation to enable seamless provider switching while maintaining client compatibility.

Two-Layer Transformation Model

- 1. Client Transformations: Convert client events TO provider format
- 2. Server Transformations: Convert provider events BACK TO client format

6.2 Server Transformation Logic

Gemini → OpenAl Server Transformation

Design Philosophy: Direct 1:1 event mapping without pseudo-event generation

Key Transformations:

- User Transcription: serverContent.inputTranscription.text → conversation.item.input_audio_transcription.delta
- Response Audio: serverContent.modelTurn.parts[].inlineData.data → response.audio.delta
- Response Transcript: serverContent.outputTranscription.text → response.audio_transcript.delta
- Function Calls: toolCall.functionCalls[] → Multiple response.output_item.done events
- Turn Completion: generationComplete/interrupted → response.done

Function Call Transformation (1:N): Single Gemini event with multiple function calls

transforms to multiple OpenAI events:

Gemini: {toolCall: {functionCalls: [func1, func2]}}

→ OpenAl: [{type: "response.output_item.done", item: func1},

{type: "response.output_item.done", item: func2}]

OpenAl → Gemini Server Transformation

Design Philosophy: Event filtering with function call accumulation

Key Transformations:

- User Transcription: conversation.item.input_audio_transcription.delta → serverContent.inputTranscription
- Response Audio: response.audio.delta → serverContent.modelTurn.parts[].inlineData
- Response Transcript: response.audio_transcript.delta → serverContent.outputTranscription
- Function Calls: Multiple response.output_item.done → Single toolCall.functionCalls[]
- Turn Completion: response.done → generationComplete/interrupted + turnComplete

Function Call Accumulation (N:1): The transformer immediately processes each function call without accumulation, aligning with observed Gemini behavior:

```
OpenAl: [{type: "response.output_item.done", item: func1}]

→ Gemini: {toolCall: {functionCalls: [func1]}}

OpenAl: [{type: "response.output_item.done", item: func2}]

→ Gemini: {toolCall: {functionCalls: [func2]}}
```

6.3 Client Transformation Logic

Session Configuration Transformation

Tool Definition Conversion:

- Type Case Transformation: "string" ↔ "STRING", "object" ↔ "OBJECT"
- Structural Reorganization: session.tools[] ↔ setup.tools[].functionDeclarations[]

Audio Data Transformation

Format Standardization:

- OpenAl: {type: "input_audio_buffer.append", audio: "base64"}
- Gemini: {realtimeInput: {audio: {data: "base64", mimeType: "audio/pcm;rate=24000"}}}

Tool Response Transformation

Response Format Conversion:

- Data Format: JSON string ↔ Object
- ID Field Mapping: call_id ↔ id
- Event Structure: conversation.item.create ↔ toolResponse

6.4 NoOp Transformations

When client API style matches provider (e.g., OpenAI client with OpenAI provider):

- Direct Passthrough: No transformation applied
- Performance Optimization: Minimal processing overhead
- Consistency: Same extractor pattern for uniform behavior

7. Intelligent Provider Switching

7.1 Performance Monitoring System

Real-Time Metrics Collection

The system currently monitors provider performance through:

Ping/Pong Latency: 5-second intervals with timestamp correlation

Additional performance metrics can be added in the future, such as:

- Connection Stability: WebSocket state monitoring and reconnection tracking
- Response Quality: Analysis of response times and completion rates
- Error Rate Monitoring: Track provider-specific error patterns

7.2 Switch Class Implementation

Decision Algorithm

Configurable Threshold Approach:

- Latency Threshold: Configurable via SWITCH_LATENCY_THRESHOLD_MS environment variable (default: 500ms)
- Failure Count: Configurable via SWITCH_FAILURE_COUNT environment variable (default: 3 consecutive violations)
- Reset Logic: Statistics reset after switching to prevent oscillation

Switch State Management

- Current Provider Tracking: Maintains active provider state
- Latency History: Per-provider performance arrays
- Callback System: Notifies Pipeline of switching decisions

Switching Logic Flow

- 1. Stats Collection: Receive performance data from active provider
- 2. Threshold Evaluation: Check if latency exceeds 500ms limit
- 3. Failure Counting: Increment consecutive failure counter
- 4. Switch Decision: Trigger switch after 3 consecutive failures
- 5. Provider Update: Notify Pipeline to reconfigure components
- 6. Stats Reset: Clear failure history for fresh evaluation

7.3 Pipeline Integration

Automatic Provider Registration

- Connection Callbacks: Register performance monitoring on provider connection
- Stats Forwarding: Route provider performance data to Switch instance
- Dynamic Reconfiguration: Handle provider changes without service interruption

Session State Preservation

- Configuration Replay: SessionManager maintains and replays session config
- Context Continuity: Conversation history preserved across switches
- Checkpoint Integration: Persistent conversation state management

8. Conversation Persistence & State Management

8.1 SessionManager Architecture

Configuration State Management

Session Configuration Capture:

- Extractor Integration: Uses ClientEventsExtractor to identify session updates
- Merging Strategy: Intelligent merging of session configuration changes
- Replay Capability: Automatic configuration replay after provider switches

Pipeline Integration

- Event Passthrough: Maintains event flow while capturing session state
- Provider Coordination: Ensures session config reaches new providers
- State Isolation: Session state independent of provider switching

9. Pipeline Orchestration

9.1 Pipeline Architecture

Component Orchestration

The Pipeline class coordinates all system components:

- Provider Management: Dynamic provider instantiation and switching
- Transformation Chain: Client → Provider → Server transformation flow
- Session Management: Configuration capture and replay
- Performance Monitoring: Switch integration and callback registration

Event Flow Pipeline

Client → SessionManager → ClientTransformer → ProviderEventManager → ServerTransformer → [SocketEventManager, Checkpoint]

11. Performance Test Results

Test Environment

- Server: MacBook Pro 16GB, 12-core M2, macOS 15.6.1
- Client: MacBook Air 8GB, 8-core M2, macOS 14.7.7
- Network: Home network with IP aliasing for multi-connection testing
- Test Duration: ~2.5 minutes

Test Methodology and Bandwidth Analysis

Sample AI Response Text Used: "That's a great question! I think the best approach would be to start with the basics and then build from there. What do you think about that idea?"

Real-World Speaking Test:

- Text characteristics: 28 words, 36 tokens per response
- Speaking speed test: Maximum 8 responses per minute when speaking continuously

- Realistic conversation pattern: 2 user inputs + 6 Al responses per minute
- Audio encoding: WAV to Base64, 972,848 bytes per response

Bandwidth Calculation:

- Al responses per minute: 6
- Bandwidth per response: 972,848 bytes
- Total bandwidth per connection: 6 × 972,848 = 5,837,088 bytes = 5.57 MB/minute

WebSocket Connection Scale Test Results

Peak Performance Metrics (Before Server's resource limits reached and system crashed and auto restarted):

- Maximum Connections: 505,069 concurrent WebSockets
- Peak Theoretical Bandwidth: 2.78 TB/minute at full capacity
- Per-Connection Load: 5.57 MB/minute (realistic voice Al audio stream)
- Connection Rate: ~3,367 connections/second sustained
- Test Duration: 2.5 minutes before server crash and restart

System Limitations Identified:

- Server Crash Point: 505K connections (M2 MacBook Pro hardware limit)
- Client Bottleneck: CPU usage reached 100% eventually becoming bottleneck

HMAC Authentication Performance:

 Authentication Processing: HMAC was performed for every call which brought CPU usage to ~60% for 1 core on server

This demonstrates the architecture can handle enterprise-scale deployments with hundreds of thousands of concurrent voice AI conversations while maintaining security and performance. The theoretical bandwidth capacity of 2.78 TB/minute positions this system for large-scale production environments with realistic voice AI conversation patterns. The only bottleneck could be the underlying Voice providers like OpenAI, Gemini etc and they could trigger rate limits after a point.

12. Cost Analysis and Comparison

Background: Total Cost Components

Complete Voice AI Deployment Costs: Voice AI applications require two primary cost components that must be considered together:

- 1. Infrastructure Cost: Hosting the application that manages voice Al connections
 - VPS federation (Eg: Hetzner/OVH etc)
 - Cloud hosting (AWS/GCP)
 - Network bandwidth and server capacity
- 2. Al Inference Cost: Payment to voice Al providers for processing
 - OpenAl Realtime API (token-based pricing)
 - Google Gemini Live (token-based pricing)
 - ElevenLabs (subscription + usage-based pricing)

Test Scenario: 500K concurrent users × 1-minute conversations = 500K total connection-minutes

Network Infrastructure & Pricing Models

Network Speed vs Bandwidth Allowance: Voice Al applications face dual constraints - network interface speed (how fast data can flow) and bandwidth allowances (how much total data is permitted monthly). Network speed determines maximum concurrent connections, while bandwidth allowances affect operational costs.

Provider Infrastructure Comparison:

Provider Type	Network Speed	Bandwidth Model	Pricing Philosophy
Budget VPS	200 Mbps - 2.5 Gbps	Fixed monthly allowances	Unlimited within limits
Cloud Providers	1-100+ Gbps	Pay-per-GB egress	Scalable but expensive
Dedicated Hosting	10-100 Gbps	Various hybrid models	Premium performance

Key Providers Analysis:

- VPS Providers: 2.5 Gbps network, 2TB daily allowance, €5.50/month (representative high-performance VPS)
- Cloud Providers (AWS/GCP): Variable network (1-100+ Gbps), \$0.01-0.09 per GB egress

Case Studies: 500K Concurrent Voice Al Connections

Baseline Requirements:

- 500,000 concurrent connections for 1-minute duration
- 5.57 MB per connection per minute (realistic voice AI audio stream)
- Peak bandwidth needed: 500K × 5.57 MB = 2.78 TB/minute
- Network speed conversion: 2.78 TB/minute ÷ 60 seconds = 46.33 GB/second × 8 bits = 370.67 Gbps sustained
- Total data transfer: 500K × 5.57 MB = 2.785 GB total

Case 1: VPS Federation Strategy

Technical Specifications:

- Network speed: 2.5 Gbps per server
- Servers required: 370.67 Gbps ÷ 2.5 Gbps = 148 servers
- Bandwidth allowance: 2TB/day per server

Bandwidth Allowance Verification:

- Bandwidth per server: 370.67 ÷ 148 = 2.5 Gbps = 0.31 GB/s
- Daily usage per server: 0.31 × 86,400 seconds = 27 GB/day
- Allowance check: 27 GB/day << 2TB/day (73x under limit)

Cost Analysis:

- Monthly server cost: 148 × €5.50 = €814 (\$884)
- Annual cost: \$10,608
- Cost per connection (annual): \$0.021

Case 2: AWS Cloud Deployment

Instance Requirements:

Instance type: r8g.large (12.5 Gbps network, \$0.11782/hour)

• Instances needed: 370.67 Gbps ÷ 12.5 Gbps = 30 instances

Monthly compute cost: 30 × \$84.84 = \$2,545

Bandwidth Requirements:

Total data transfer: 2.785 GB

Bandwidth cost: 2.785 GB × \$0.01 = \$28

Total AWS Cost: \$2,545 + \$28 = \$2,573

Case 3: GCP Cloud Deployment

Instance Requirements:

• Instance type: c4-standard-2 (10 Gbps network, \$0.096866/hour)

• Instances needed: 370.67 Gbps ÷ 10 Gbps = 37 instances

Monthly compute cost: 37 × \$69.74 = \$2,580

Bandwidth Requirements:

• Total data transfer: 2.785 GB

Bandwidth cost: 2.785 GB × \$0.01 = \$28

Total GCP Cost: \$2,580 + \$28 = \$2,608

Case 4: ElevenLabs Enterprise Deployment

ElevenLabs Pricing Options:

Basic plan: \$11/month (250 minutes included, \$0.12/minute overage)

- Enterprise plan: \$1,320/month (13,750 minutes included, \$0.096/minute overage)
- Test scenario: 500K users × 1-minute conversations = 500K minutes needed

Basic Plan Analysis (Not suitable for enterprise):

• Base cost: \$11/month

• Overage minutes: 500,000 - 250 = 499,750 minutes

Overage cost: 499,750 × \$0.12 = \$59,970

• Total: \$59,981/month

Enterprise Plan Analysis:

• Base cost: \$1,320/month

Overage minutes: 500,000 - 13,750 = 486,250 minutes

Overage cost: 486,250 × \$0.096 = \$46,680

Total: \$48,000/monthAnnual cost: \$576,000

• Cost per connection (annual): \$1.15

Summary Analysis

Network Speed Impact: Network interface speed emerges as the primary technical constraint. High-performance VPS providers with fast network speeds (2.5 Gbps) offer significantly better server efficiency compared to lower-tier alternatives.

Infrastructure Cost Comparison (Monthly fixed costs):

1. VPS Federation: \$884/month

2. AWS Cloud: \$2,573/month

3. GCP Cloud: \$2,608/month

4. ElevenLabs Enterprise: \$0/month (SaaS model)

Infrastructure Efficiency Rankings:

Most Efficient: High-performance VPS (baseline)

Moderate: AWS/GCP Cloud (2.9x more expensive)

• Premium: ElevenLabs SaaS (no infrastructure management required)

Bandwidth Allowance Analysis: High-performance VPS providers typically offer generous bandwidth allowances:

- Usage: Only 1.4% of typical daily allowances (27GB of 2TB)
- Conclusion: Network speed becomes the primary constraint, not bandwidth allowances

Voice Al Provider Costs Analysis

Based on the realistic conversation pattern (2 user inputs + 6 Al responses per minute):

Token Usage Calculation (500K users × 1-minute each):

- Input tokens per user: 2 responses × 36 tokens = 72 tokens
- Output tokens per user: 6 responses × 36 tokens = 216 tokens
- Total usage: 500K × 72 = 36M input tokens, 500K × 216 = 108M output tokens

Al Provider Cost Comparison:

Provider	Input Cost	Output Cost	Total Cost
Gemini Live	36M × \$3/1M = \$108	108M × \$12/1M = \$1,296	\$1,404
OpenAl Realtime	36M × \$40/1M = \$1,440	108M × \$80/1M = \$8,640	\$10,080
ElevenLabs Enterprise	Flat cost	500K minutes × overage	\$48,000

ElevenLabs Enterprise Analysis (500K × 1-minute):

- Basic plan: \$11 + (499,750 × \$0.12) = \$59,981 (not viable)
- Enterprise plan: \$1,320 + (486,250 × \$0.096) = \$48,000 (appropriate tier)
- Business model: SaaS pricing with no infrastructure management required

Complete Deployment Costs (Infrastructure + Al Provider)

Optimal Configuration: Gemini Live + VPS

Al Provider: \$42,120/monthInfrastructure: \$884/month

• Total: \$43,004/month

Cost per user: \$0.086/month

Cloud Alternative: Gemini Live + AWS

• Al Provider: \$42,120/month

• Infrastructure: \$1,273,000/month

Total: \$1,315,120/monthCost per user: \$2.63/month

Complete Cost Analysis: Infrastructure + Al Inference

For 500K concurrent users × 1-minute conversations, here are the total deployment costs:

Infrastructure Costs (Monthly fixed costs for peak capacity):

• VPS Federation: \$884 (148 servers × €5.50)

• AWS Cloud: \$2,573 (30 r8g.large instances + bandwidth)

• GCP Cloud: \$2,608 (37 c4-standard-2 instances + bandwidth)

Al Inference Costs (Usage-based for 500K × 1-minute):

• Gemini Live: \$1,404 (36M input + 108M output tokens)

• OpenAl Realtime: \$10,080 (36M input + 108M output tokens)

• ElevenLabs Enterprise: \$48,000 (500K minutes with enterprise plan)

Final Total Cost Comparison

Configuration	Infrastructure	Al Inference	Total Cost	Cost per User
VPS + Gemini	\$884	\$1,404	\$2,288	\$0.0046
VPS + OpenAI	\$884	\$10,080	\$10,964	\$0.022
AWS + Gemini	\$2,573	\$1,404	\$3,977	\$0.008

GCP + Gemini	\$2,608	\$1,404	\$4,012	\$0.008
AWS + OpenAl	\$2,573	\$10,080	\$12,653	\$0.025
ElevenLabs Enterprise	\$0	\$48,000	\$48,000	\$0.096

Strategic Implications

The optimal configuration (VPS + Gemini) provides the most cost-effective solution at \$2,288 for $500K \times 1$ -minute conversations.

Key Insights:

- VPS efficiency: 1.7x cheaper infrastructure than cloud providers
- Al provider impact: OpenAl costs 7x more than Gemini for identical functionality
- Infrastructure vs inference: Al costs dominate for short conversations (61% vs 39%)
- ElevenLabs premium: 21x more expensive than optimal self-hosted solution
- Scalability consideration: Fixed infrastructure costs amortize better with higher usage

Practical Business Scenario Analysis

Beyond the peak capacity test, here's a realistic business usage pattern analysis:

Test Case: 100 Calls per Hour Business Operation

Business Pattern:

- 100 calls/hour × 30 minutes each = 3,000 call-minutes/hour
- Operating schedule: 8 hours/day × 5 days/week × 4 weeks/month = 160 hours/month
- Total monthly usage: 480,000 call-minutes

Peak Concurrent Load: 100 calls \times 30 minutes = 3,000 concurrent connections Peak Bandwidth: 3,000 \times 5.57 MB/min = 2.23 Gbps sustained

Token Usage:

• Input tokens: 480,000 × 72 = 34.56M tokens

• Output tokens: 480,000 × 216 = 103.68M tokens

Infrastructure Requirements:

• VPS: 1 server (2.5 Gbps capacity) = \$6/month

• AWS: 1 × r8g.large (12.5 Gbps network) = \$85/month + minimal egress

• GCP: 1 × c4-standard-2 (10 Gbps network) = \$70/month + minimal egress

Al Inference Costs:

• Gemini Live: \$104 + \$1,244 = \$1,348/month

• OpenAl Realtime: \$1,382 + \$8,294 = \$9,676/month

• ElevenLabs Enterprise: $$1,320 + (466,250 \times $0.096) = $46,080/month$

Complete Business Scenario Cost Comparison

Configuration	Infrastructure	Al Inference	Total Monthly	Cost per Call
VPS + Gemini	\$6	\$1,348	\$1,354	\$0.084
VPS + OpenAl	\$6	\$9,676	\$9,682	\$0.605
AWS + Gemini	\$85	\$1,348	\$1,433	\$0.090
GCP + Gemini	\$70	\$1,348	\$1,418	\$0.089
ElevenLabs Enterprise	\$0	\$46,080	\$46,080	\$2.88

Business Insights:

- Infrastructure becomes negligible: 0.4-6% of total costs for sustained usage
- Al provider hierarchy: Gemini (\$0.084/call) < OpenAl (\$0.605/call) < ElevenLabs

(\$2.88/call)

- VPS maintains cost advantage: \$0.084 vs \$0.089-0.090 per call for cloud
- Provider switching impact:
 - Gemini vs OpenAI: \$8,328/month savings
 - Gemini vs ElevenLabs: \$44,726/month savings (34x cost difference)
- ElevenLabs premium: 34x more expensive than optimal solution for business operations

13. Current Implementation Status and Future Roadmap

Current Implementation Status

Fully Implemented & Tested

- Core Architecture: All design patterns and base classes complete
- Provider Integration: OpenAI and Gemini fully operational
- Bidirectional Transformations: All transformation combinations working
- Function Call Support: Complete function calling with response handling
- Conversation Persistence: Checkpoint system with file-based storage
- Performance Monitoring: Real-time latency tracking and switching

Comprehensive Test Coverage

- Baseline Tests: Individual provider behavior validation
- Pipeline Tests: Cross-provider transformation verification
- Function Call Tests: Multi-function scenarios with response handling
- Audio Processing: WAV/PCM conversion and streaming

Future Roadmap

Open Source Initiative

 Opensourcing the product: Making the voice AI abstraction layer available to the broader developer community

Provider Ecosystem Expansion

 Add more providers: Integration with additional voice AI services (Anthropic Claude Voice, Azure Speech, AWS Transcribe, etc.)

Enhanced Intelligence

 Advanced switch rules: Machine learning-based switching algorithms, cost optimization strategies, and quality-aware provider selection

Conclusion

This unified abstraction layer eliminates voice AI integration hell by providing a single, consistent interface across multiple providers. Developers no longer need to manage complex provider-specific implementations, reducing development time and increasing system reliability.

The architecture is built in a scalable and extensible way, enabling seamless provider switching without conversation interruption while maintaining optimal performance and user experience. This approach significantly eases development complexity and provides the flexibility to adapt to the rapidly evolving voice AI landscape.