

# Vapi Integration & Analysis

Generated on 2025-10-30 08:49:42

## Executive Summary

This document analyses the suitability of Vapi (Vapi AI) as an orchestration layer for your TTS benchmarking and production voice-agent needs. It recommends when to use Vapi and when to integrate TTS providers directly.

## Key Recommendations

- Use Vapi when building end-to-end voice agents (STT -> LLM -> TTS) that require orchestration, phone/telephony integration, or advanced session handling.
- For pure TTS latency benchmarking, call TTS providers directly (Cartesia, ElevenLabs, Deepgram) to avoid extra orchestration overhead when measuring raw TTS latency.
- You can combine both approaches: benchmark TTS providers directly, then plug the selected providers into Vapi for production testing (end-to-end).

## Integration Patterns

1) Direct TTS (Benchmarking) — Call providers directly using REST or WebSocket and measure latency. 2) Vapi as Orchestration — Use Vapi to host the voice-session, then configure Vapi to call your chosen TTS endpoint for output. This measures end-to-end latency inside an agent flow. 3) Hybrid — Benchmark providers directly for raw performance, then validate selected provider(s) through Vapi for real-world E2E performance.

## Example Vapi Integration Flow

Basic flow for plugging a custom TTS into Vapi: 1. Vapi receives the text to speak from the agent. 2. Vapi calls your custom TTS endpoint (could be Cartesia/ElevenLabs through your proxy or direct) and waits for audio. 3. Vapi streams the audio back to the caller or plays it on the telephony channel. Key measurement points: request time to TTS, time-to-first-byte on audio stream, total audio generation time, and network RTT between Vapi and the TTS provider.

## Risks & Caveats

- Additional network hop: routing TTS through Vapi introduces additional network hop and potential serialization overhead.
- Vendor-specific streaming protocols: Vapi supports plugging custom endpoints, but streaming control frames may need adapter code.
- Cost & licensing: adding Vapi may increase costs and operational complexity; evaluate only if you need full agent features.

## Conclusion

Vapi is recommended as a production orchestration layer when you need voice-agent capabilities. For benchmarking and raw latency-focused evaluation, use direct integration with TTS providers. After you select the TTS providers, validate them end-to-end through Vapi to ensure real-world performance.