# The End of CPaaS: Why SignalWire's Al-Integrated Call Fabric Redefines Real-Time Communication

# A Technical Blueprint for the Stack Everyone Else is Still Trying to Build

# 1. Executive Summary

The telecom and CPaaS landscape has remained stagnant for over a decade, a landscape plagued by bolt-on APIs, vendor lock-in, and fragmented media handling. As the market rushes to retrofit AI into traditional voice stacks, latency balloons, complexity surges, and developers are left struggling to integrate components that were never designed to work together.

SignalWire redefines this reality with a radical new paradigm: **Call Fabric**, a programmable, addressable network layer that treats every participant, room, and AI as a first-class, routable resource. It's a unified, programmable, and composable media infrastructure where **AI is embedded directly into the real-time communication stack**. With **sub-500ms latency**, seamless support for **voice**, **video**, **messaging**, **PSTN**, **SIP**, **and WebRTC**, and the powerful declarative scripting language SWML, SignalWire delivers the future of communications today.

Others are still raising capital to build what SignalWire has already deployed globally.

# 2. The Problem with CPaaS, Vertical AI, and "Real-Time" Infrastructure

#### 2.1 CPaaS is Not Programmable Enough

Traditional CPaaS exposes APIs, but developers are still limited to rigid voice and messaging workflows, such as predefined IVRs or hardcoded call trees with limited conditional logic. Customization often requires external logic hosted elsewhere, leading to disjointed experiences. You can't embed intelligence into the stack, you wrap it around the stack, and the seams always show.

#### 2.2 Voice AI Is Still a Bolt-On

Platforms like Twilio's Voice AI and Vonage AI Connect stream raw audio through WebSockets to external LLMs. This adds latency (1.5 to 3 seconds), increases costs, and fractures context management. These systems cannot maintain fluid, lifelike conversations.

#### 2.3 Channels Are Siloed

SIP, PSTN, WebRTC, video, and messaging are typically handled by distinct services or even entirely separate vendors. Developers have to stitch together identity, routing, and AI logic across tools that were never meant to interoperate. This makes omnichannel AI nearly impossible.

#### 2.4 Fragmentation Adds Latency and Cost

Every additional service adds hops. Each hop adds delay. Developers face escalating costs for bandwidth, compute, and coordination, while users endure awkward pauses and brittle handoffs between systems.

# 3. Introducing SignalWire's Call Fabric

**Call Fabric** is SignalWire's next-generation programmable communications layer. It treats every element, media stream, user, AI, room, or script, as a composable, addressable resource.

#### 3.1 Composable Resources

- Subscribers: SIP endpoints, mobile apps, or authenticated user identities
- Rooms: Multiparty audio/video bridges
- SWML Scripts: JSON/YAML logic documents that describe real-time call flows
- Al Agents: Digital employees that interact with users across modalities
- Queues: Intelligent routing and dispatch primitives

#### 3.2 Addressed Like the Web

Every resource has an address (e.g., /support/agent\_bot) that can be referenced, dialed, and interacted with, just like web URLs. This means developers can route calls or messages between resources using declarative logic, with full AI orchestration in between.

#### 3.3 Execute Across All Channels

SignalWire supports PSTN, SIP, WebRTC, SMS, IP messaging, and video conferencing as native, first-class channels. Developers can initiate, bridge, and route interactions across any modality in real time, enabling dynamic workflows that preserve context, memory, and state across the entire conversation lifecycle, even as the channel changes.

#### 3.4 Real-World Composition

For example, consider an inbound PSTN call answered by an AI Agent. The agent greets the caller, identifies their intent using natural language processing, queries a business hours API, and determines whether to route the call immediately or offer voicemail. If escalation is required, it dynamically places the caller into a Room with an available representative. After the session ends, the agent sends a personalized summary and follow-up via SMS, all executed seamlessly within a single, declarative workflow.

# 4. The AI Voice Stack: Embedded, Composable, and Live

#### 4.1 AI Embedded in the Media Plane

Unlike competitors, SignalWire's AI kernel lives **inside** the media and signaling stack. This allows real-time transcription, inference, and synthesis with **an average 500ms round-trip latency**. No WebSockets. No proxy servers. No detours.

#### 4.2 Full Runtime with STT, TTS, and LLMs

Speech-to-text, text-to-speech, and natural language reasoning are all tightly coupled, which ensures faster response cycles, reduced failure points, and more context-aware interactions for the user. The AI understands the media context, maintains conversation state, and can trigger functions, route calls, or escalate to humans, all declaratively.

#### 4.3 AI Agents as Resources

Agents can be invoked by address, reused in scripts, transferred into, or even embedded inside Rooms. You can treat AI like a SIP user, but with intelligence.

#### 4.4 Serverless Tools via SWAIG

**SWAIG** lets developers define actions the AI can perform using only JSON. These actions can be templated serverless workflows or proxied HTTP requests with runtime variable expansion. This makes it trivial to plug AI into CRMs, knowledge bases, payment APIs, and more, with no glue code.

# 5. SWML: A Language to Program Communication

#### 5.1 Declarative Control

**SWML (SignalWire Markup Language)** is a document-based, JSON/YAML-friendly language designed for real-time control of communication and AI logic.

#### 5.2 One-Liners for Everything

Want an agent to answer and collect info? One line.

```
Unset
- ai:
    prompt:
    text: "You're a receptionist. Greet the caller and collect
their name and reason for calling."
```

#### 5.3 Full Flow Programming

With execute, goto, cond, and switch, you can branch flows, collect input, loop, and embed reusable logic. Unlike low-code IVR builders, SWML scales to real-world, multi-step workflows.

#### 5.4 Post-Call Intelligence

Use post\_prompt or post\_prompt\_url to summarize, log, or escalate sessions. The AI can convert natural conversation into JSON, ready for back-end systems.

#### 6. LiveKit Reality Check: They're Building What We Already Built

In April 2024, LiveKit raised \$22 million to become the programmable media infrastructure for AI. While their ambitions mirror SignalWire's, their product does not.

LiveKit provides WebRTC infrastructure. It has no SIP, no PSTN, and no AI runtime. According to their own April 2024 Series B announcement, developers are responsible for stitching together LLMs, memory systems, and transcription services externally, requiring custom logic, orchestration middleware, and integration layers for even basic agent interactions.

#### SignalWire already delivers all of this today.



Composable Telecom PrimitivesImage: Rooms, Subs, Queues, Al<br/>AgentsImage: Media Rooms Only<br/>Media Rooms Only<br/>Media Rooms Only<br/>Media Rooms Only<br/>Image: Image: ProductStage of ProductImage: DeployedImage: Image: Image:

# 7. Displacing the CPaaS + AI Stack

Feature	SignalWire	Twilio	Vonage	PolyAl	Dialogflow	LiveKit
Al Built Into Media Stack	$\checkmark$	×	×	×	×	×
Fully Addressable Resource Model		×	×	×	×	×
~ 500ms Round Trip Latency	$\checkmark$	×	×	×	×	X
SIP + WebRTC + Video + PSTN	$\checkmark$		$\checkmark$	×	×	X
Declarative AI + Call Logic		×	×	×	1 Prompt only	×
No Glue Code or Middleware	$\checkmark$	×	×	×	×	×

# 8. Developer Experience: One Line to Start, Infinite Depth

SignalWire was designed with the developer in mind, from the first line of SWML to full-scale applications that can power global call centers. The developer experience is built on three principles: simplicity, composability, and control.

Getting started with AI-driven communications takes a single line of SWML. Developers can spin up an AI Agent with a purpose-specific prompt, route a call, and deploy it to a phone number or SIP endpoint in moments. There's no infrastructure to manage, no API keys to juggle across providers, and no need to manually connect speech recognition to a language model.

As needs grow, so does the depth of customization. Developers can integrate retrieval-augmented generation (RAG) systems to feed the AI live data, add memory to support stateful interactions, and define tools using SWAIG to allow agents to take meaningful actions like checking order status, submitting forms, or scheduling appointments. All of this is expressed declaratively in a single SWML document, with optional external APIs or webhooks for deeper integration.

Whether routing based on user intent, escalating to human agents, sending SMS follow-ups, or transcribing calls in real time, developers remain in full control of the experience. SignalWire removes the burden of orchestration and instead delivers a model where the platform handles real-time execution, while developers focus on outcomes.

SignalWire doesn't just simplify development, it elevates it. Developers get a fully programmable, AI-native comms stack that responds in milliseconds, scales automatically, and speaks the language of modern software.

# 9. Performance Engineering at Global Scale

SignalWire's architecture is performance-first, engineered from the ground up to handle real-time audio, video, and AI processing with near-zero latency. Unlike platforms that pass media through cloud microservices or stream audio to third-party processors, SignalWire keeps everything local to its own global edge network. This minimizes round-trip delay and ensures a seamless, responsive user experience.

The platform is built on FreeSWITCH, the open-source engine trusted by Tier 1 carriers and VoIP providers worldwide, and has been battle-tested in some of the most demanding telecom environments on the planet. SignalWire extends FreeSWITCH's low-level power into a modern cloud-native fabric, allowing any developer to tap into its capabilities without needing to manage servers or SIP infrastructure.

Al processing is embedded within the media path. This design eliminates the typical latency spikes that occur when raw audio is shipped to distant LLMs or transcription services. Instead, speech-to-text, natural language processing, and synthesis all happen inside the same runtime. The result: average latencies below 500 milliseconds, even for complex interactions involving memory and external API calls.

SignalWire's infrastructure is geographically distributed, automatically routing traffic through the closest media edge location. This ensures that users across North America, Europe, and Asia experience consistent performance and quality. The system supports dynamic scaling, horizontal load distribution, and automatic failover, all essential features for enterprise-grade deployment.

Security and compliance are also core to the performance model. Audio data is encrypted end-to-end, sensitive information can be handled outside the AI's memory space, and metadata can be tokenized to protect user identity while preserving context. Whether it's for healthcare, finance, or government, SignalWire is built to meet modern compliance standards without sacrificing speed.

# **10. The Strategic Future**

SignalWire's architecture isn't just solving today's telecom problems, it's laying the foundation for a new era of programmable communication. In this new model, every conversation is orchestrated dynamically in real time, with media, memory, logic, and modality all managed within a single, coherent system. As more organizations realize the limitations of retrofitting legacy systems with AI, the shift toward native, embedded intelligence will accelerate. SignalWire is already leading that transition with infrastructure built for scale, security, and real-time action.

The concept of **Call Fabric** mirrors the early internet, where every resource is addressable, composable, and built to interoperate. Just as URLs transformed how we access content, SignalWire's resource model allows communication endpoints, AI Agents, Rooms, Subscribers, to act like building blocks for new applications. This architecture positions SignalWire not just as a telecom company, but as the communications layer of the modern internet.

In the coming years, we will see AI Agents evolve from assistants to collaborators. They won't just answer questions, they'll complete tasks, escalate conversations intelligently, and learn over time. SignalWire enables this by making AI a first-class citizen of the media stack, not an external service tacked on with middleware and latency.

As digital identity becomes more portable and interaction channels multiply, SignalWire's composable infrastructure ensures continuity across modalities. A conversation that starts on the phone can shift to video, then to messaging, with context and state preserved. This fluid, multimodal approach unlocks entirely new forms of customer experience, support, and automation.

The future of communications isn't a product, it's a platform. And SignalWire's Call Fabric is already powering it.

# 11. Conclusion

The communications industry is standing at a pivotal inflection point. Legacy CPaaS vendors and next-gen media startups alike are racing to stitch AI into telecom, but they're burdened by outdated paradigms, fragmented tooling, and an architectural gap between conversation and action.

SignalWire changes that. It doesn't offer another API, chatbot, or voice add-on. Instead, it delivers a platform where AI, media, logic, and infrastructure are one and the same. Call Fabric unifies programmable communication across voice, video, and messaging, while embedding real-time AI agents into the stack itself. The result: faster development, natural conversations, and an execution model that scales.

While others define roadmaps and raise rounds to catch up, SignalWire is already there. Live in production. Deployed globally. Built by the same minds who revolutionized telecom with FreeSWITCH, now doing it again with programmable AI.

The shift is no longer coming. It's here, and it's composable, real-time, and built on SignalWire.

### Appendices

- A. Full AI Agent SWML Example
- B. LiveKit PR & Competitive Notes
- C. Benchmarks: Latency, Call Setup, Bandwidth
- D. Call Fabric Resource Diagram
- E. Glossary: PUC, SWAIG, SWML, etc.