



AI-Powered Voice Booking Agent (ServiceTitan Integration)

This project implements an **AI voice agent** that answers customer calls and **automatically schedules service appointments** in ServiceTitan. It handles two call scenarios (new vs. existing customer), collects customer information (name, address, service request), and then uses ServiceTitan's APIs to create or retrieve the customer and **book a job** at an available time slot. The system is deployed as an MVP on AWS (using custom backend APIs) and leverages a telephony interface (e.g. Twilio) for call handling. For example, Retell.ai's integration page describes this exact use case – using a voice agent with ServiceTitan to “schedule appointments during customer calls” ¹ .

Client Pain Points

The client's customers often **missed calls** and manual booking was labor-intensive. Key pain points include:

- **Missed Opportunities:** Unanswered calls mean lost revenue. AI assistants can “pick up right away” and schedule bookings without human intervention ² . Industry studies note that every missed call is money lost, whereas a voice agent can capture those leads automatically ² .
- **Leads Falling Through the Cracks:** Without automation, callers may drop off or need callback. An AI receptionist ensures “every call” is answered even after hours, so no leads slip away ³ ² .
- **Scalability & Cost:** A 24/7 virtual agent handles high call volume without overtime. This is cost-effective – AI receptionists cost on the order of ~\$65-\$200/month, far less than a full-time human operator ⁴ .
- **Data Entry Overhead:** Manual scheduling requires duplicating information into the CRM. Automating booking via API integration eliminates redundant data entry and errors.

By automating call intake and booking, the voice agent directly addresses these issues – increasing answer rates and lead capture ² ³ .

Solution and Key Features

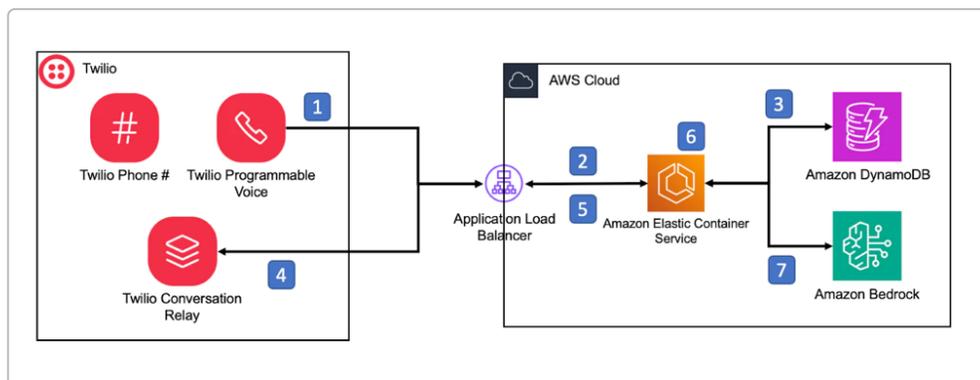
The voice agent provides the following capabilities:

- **New vs. Existing Customer Flow:** On each call, it determines if the caller exists in ServiceTitan. For new customers, it conversationally collects full name, address, and service reason; for existing ones, it looks up the customer record (e.g. via phone number).
- **Customer Creation:** If the customer is new, the agent uses the ServiceTitan API to create a customer and location. ServiceTitan requires both a **customer ID** and a **location ID** before booking a job ⁵ . The system automatically invokes the appropriate “CreateCustomer” and “CreateLocation” endpoints behind the scenes.

- **Appointment Scheduling:** Once the customer/location exist, the agent checks available times and asks the caller which slot works best. It then calls ServiceTitan’s “BookJob” API with the customer ID, location ID, chosen date/time, and other required fields to **create the appointment** ⁵ .
- **ServiceTitan Configuration:** Certain fields (business unit ID, job type ID, campaign ID, priority, etc.) are sent with the booking. In this MVP they are currently set as static placeholders, but in production they must match the client’s ServiceTitan setup. (For instance, if an account maps job types to business units, the supplied IDs must align ⁶ .)
- **Natural Dialogue:** The agent uses voice prompts to collect data (“Please provide your name and address,” “What date/time works best?”) and to confirm success (“Your appointment has been scheduled.”). It hands off to a live agent only if needed (e.g. unexpected input), preserving caller context.

By fully automating data collection and API calls, the solution **eliminates manual scheduling**. It ensures callers get an immediate response and booked appointment (or callback scheduling) without human intervention, 24/7.

Architecture and Implementation



Architecture: Twilio’s ConversationRelay ties the phone system to an AWS backend for processing. An incoming call to a configured number triggers a TwiML response that connects to the ConversationRelay WebSocket. Twilio then streams the caller’s speech-to-text to our application and plays back our responses as speech ⁷ ⁸ . The cloud backend (hosted on AWS) receives the transcribed text, runs the AI/logic (customer lookup, slot selection, API calls), and sends text back to be spoken. This design leverages proven cloud components to achieve scalable, low-latency voice interactions ⁸ ⁷ .

Key technical elements include: - **Telephony Interface:** We use Twilio (or a SIP trunk) to handle inbound calls. A TwiML endpoint initiates the ConversationRelay session, which transparently handles real-time speech recognition and text-to-speech ⁷ . - **Serverless Backend:** The conversation logic runs in AWS. For example, one implementation uses API Gateway and AWS Lambda (deployed via AWS SAM) to process Twilio events ⁹ . Incoming text messages (from Twilio) hit a WebSocket API (managed by API Gateway) whose Lambda function handles `connect`, `message`, and `disconnect` events ⁹ . - **State Management:** The system can use DynamoDB (or similar) to maintain session state (e.g. storing the ongoing conversation context or temporary data) during the call. This ensures information (like the name/address just collected) carries through to booking. - **AI/LLM Integration:** The agent’s dialogue and slot recommendation could leverage large language models (e.g. via AWS Bedrock) for natural response generation. In this MVP, most prompts are scripted or simple logic, but an LLM could enrich language understanding. AWS provides out-

of-the-box speech-to-text and text-to-speech integrations (via Twilio), and the application can call AI services as needed. - **ServiceTitan API Calls:** Custom code in the backend handles HTTP calls to ServiceTitan's REST API. It uses OAuth credentials (Client ID/Secret from ServiceTitan's Developer portal) to authenticate, then invokes the `CreateCustomer`, `CreateLocation`, and `BookJob` endpoints as needed ⁵. All API responses are processed to confirm success (e.g. retrieving the new customer ID). - **MVP Configuration:** The current demo uses placeholder values for optional fields (business unit, job type, etc.), which were logged for manual setup. In a full deployment, these IDs would be determined from the account's configuration or user input. According to ServiceTitan documentation, "you'll need a customer ID and location ID" to book a job ⁵, and any provided business-unit/job-type IDs must match the account's setup ⁶.

Technologies & Tools

- **Cloud:** Amazon Web Services (AWS) – e.g. API Gateway, Lambda (or ECS/Fargate), DynamoDB, IAM, etc. We used AWS for scalability and quick deployment. The solution can run in a serverless fashion (Lambda) as shown in reference architectures ⁹.
- **Telephony:** Twilio Programmable Voice with ConversationRelay for phone calls and real-time transcription ⁷. A Twilio phone number (or SIP connection) is purchased and connected to our webhooks.
- **Backend:** Node.js or Python (serverless functions) handling the conversation flow and API integration. The TwiML endpoint and WebSocket listener are implemented as separate Lambda functions (as in the AWS SAM example ⁹).
- **ServiceTitan API:** We use ServiceTitan's V2 REST API (with OAuth 2.0). The "Job Planning" endpoints are used as documented ⁵.
- **Version Control & CI/CD:** (Not detailed by client, but typically Git for source control and AWS CodePipeline or GitHub Actions for deployment.)
- **Monitoring:** Logs (CloudWatch for AWS, Twilio call logs) to trace calls and debug issues.

Benefits and Next Steps

By integrating the voice agent with ServiceTitan, the client gains **automated call answering and booking 24/7**, ensuring leads are captured even outside business hours. Every incoming request is turned into a ServiceTitan job without manual handoff. Industry benchmarks show such AI agents can dramatically increase answer rates (e.g. *"We've increased our answer rate by 175% with Voice Assist"*) and generate more booked jobs ². ServiceTitan itself notes that its built-in AI voice agents "never miss a call" and book jobs instantly from your existing data ¹⁰, mirroring our custom solution's goals.

Moving forward, the system could be refined by hooking up a real telephony number (the demo used a test environment), handling edge cases (e.g. multiple-choice logic), and replacing placeholder IDs with dynamic lookups. Once live, it will directly address the client's pain points: no more "calls going into the void," consistent lead capture, and seamless CRM integration ² ³. This project – fully deployed by the developer on AWS – showcases a modern solution to an age-old business problem, and can be highlighted in the portfolio as an example of **building a production-ready AI voice assistant** that solves real client needs.

Tech Stack: Twilio Voice, Retell AI platform (for call flows), AWS (API Gateway, Lambda, DynamoDB, Bedrock), ServiceTitan API ⁷ ⁵.

Result: A working MVP voice agent that creates customer profiles and books service appointments in ServiceTitan, eliminating missed calls and saving manual effort ² ¹ .

References: Industry sources on AI voice agents ² ³ ; ServiceTitan API docs ⁵ ⁶ ; Twilio architecture guides ⁷ ⁸ ; Retell AI integration info ¹ .

¹ Connect AI call agent to ServiceTitan

<https://www.retellai.com/integrations/servicetitan>

² ⁴ AI Receptionists in 2025: Turn Missed Calls into Revenue

<https://www.callrail.com/blog/ai-receptionist>

³ ¹⁰ AI Voice Agent for Contact Centers | ServiceTitan Pro

<https://www.servicetitan.com/features/pro/contact-center/voice-agents>

⁵ ⁶ 521. Docs: Job Planning - ServiceTitan API - developer portal

<https://developer.servicetitan.io/docs/api-resources-job-planning/>

⁷ ⁹ How to Build a Serverless Voice AI Assistant for Telephony in AWS using Twilio ConversationRelay – WebRTC.ventures

<https://webrtc.ventures/2025/07/how-to-build-a-serverless-voice-ai-assistant-for-telephony-in-aws-using-twilio-conversationrelay/>

⁸ ConversationRelay Architecture for Voice AI Applications Built on AWS using Fargate and Bedrock | Twilio

<https://www.twilio.com/en-us/blog/developers/tutorials/product/reference-architecture-aws-conversationrelay-voice-ai-app>