

# HPE AI Voice Agent v4.0

A Full-Stack Conversational Intelligence Solution Powered by HPE Private Cloud AI.



Microservices Architecture

Frontend: Gradio/React

Backend: WebSocket/FastAPI

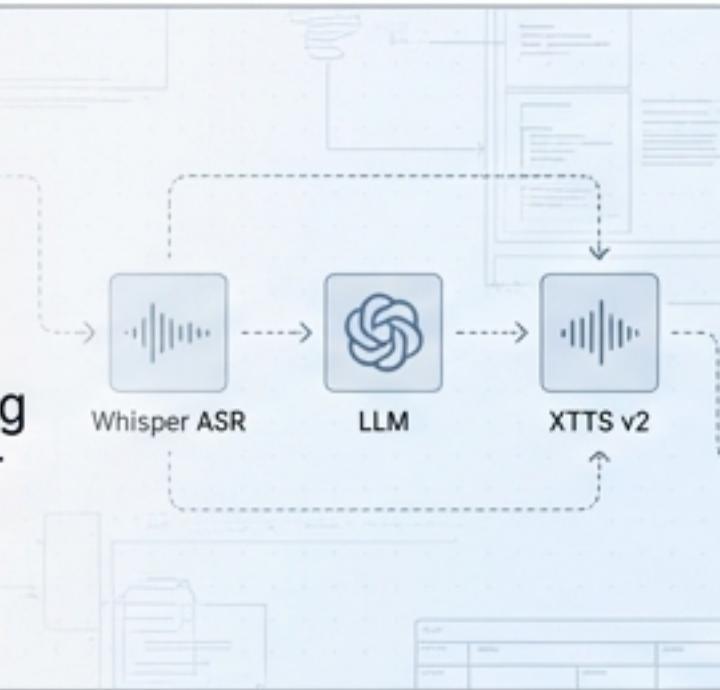
Storage: PostgreSQL

# Beyond Chatbots: A Transactional Voice Operating System



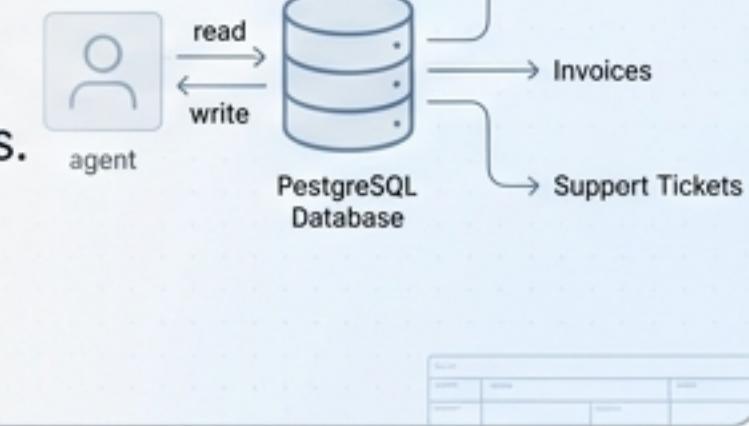
## Real-Time Orchestration

Low-latency event pipeline combining Whisper ASR for recognition, LLM for intelligence, and XTTs v2 for neural synthesis.



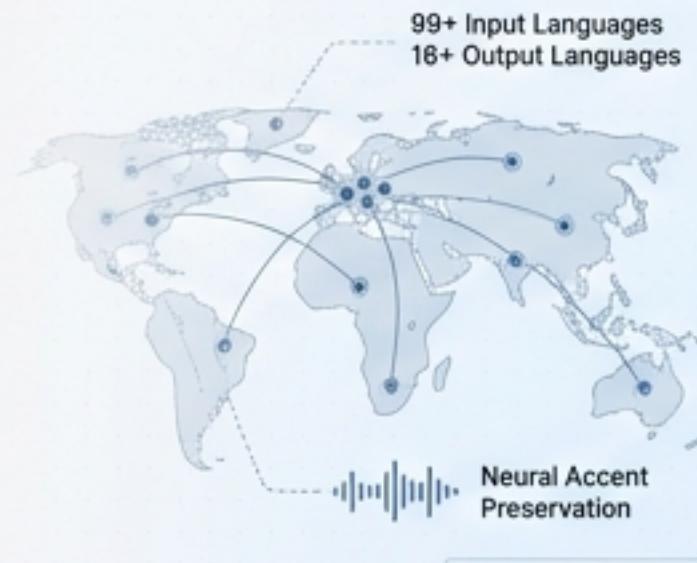
## Deep Business Integration

Direct PostgreSQL read/write access. The agent manages subscriptions, processes invoices, and updates support tickets in real-time.



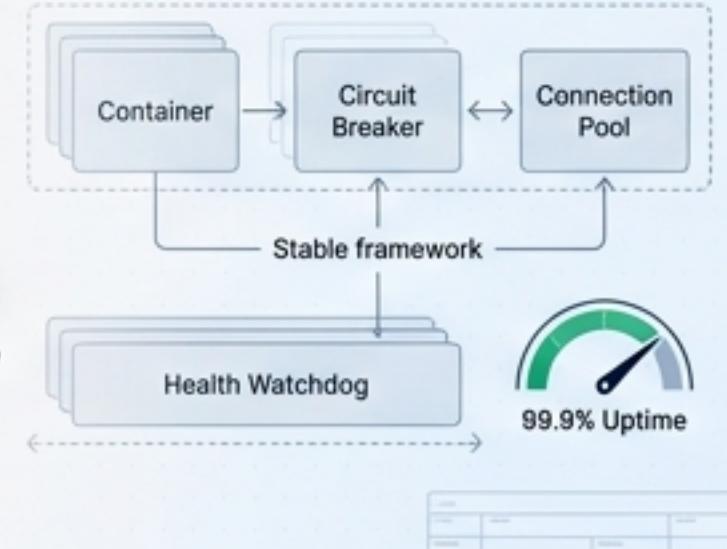
## Global Accessibility

Input support for 99+ languages with auto-detection. Output in 16+ languages (EN, ES, FR, DE, JP, etc.) with neural accent preservation.

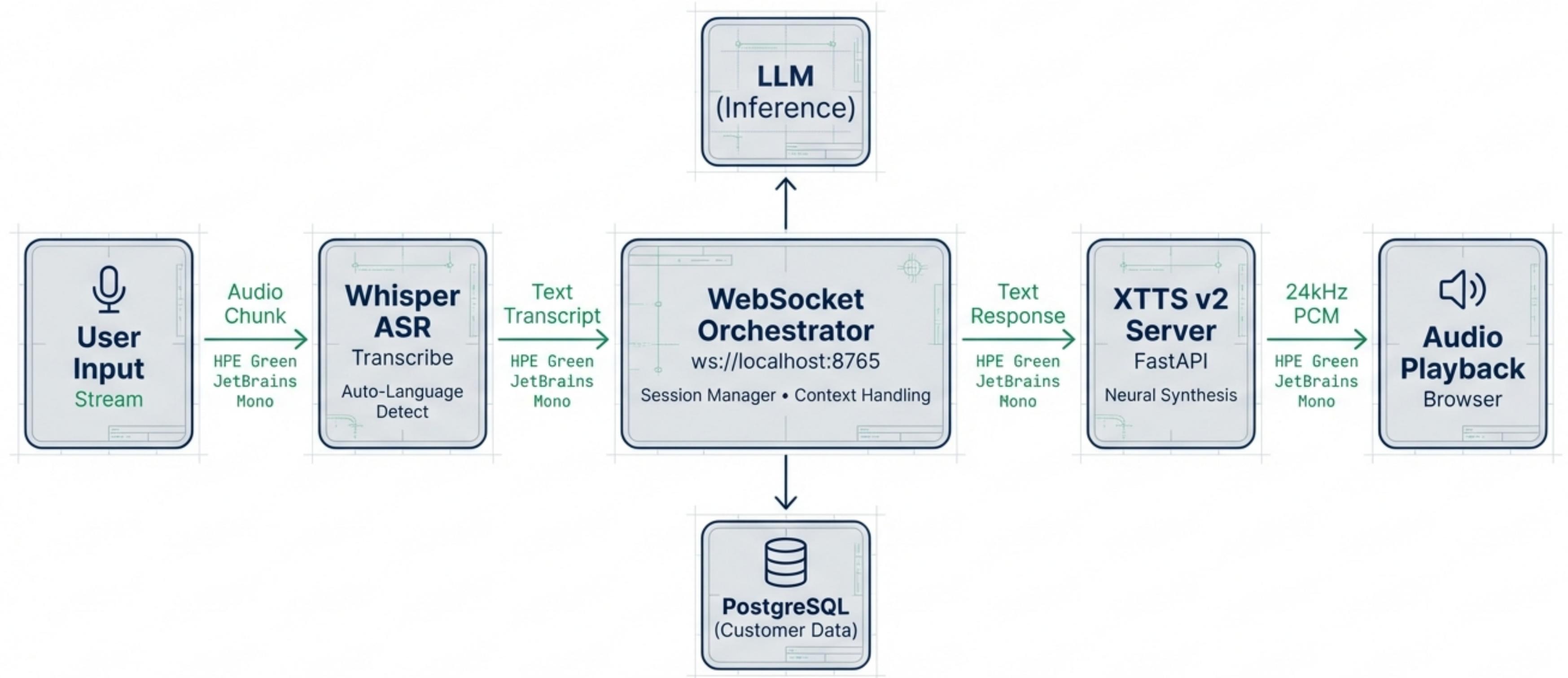


## Enterprise Reliability

Containerized architecture featuring circuit breakers, connection pooling, and automated health watchdogs for 99.9% uptime.



# End-to-End Event Pipeline



Orchestrator manages state via session\_id hashing and persistent WebSocket connections.

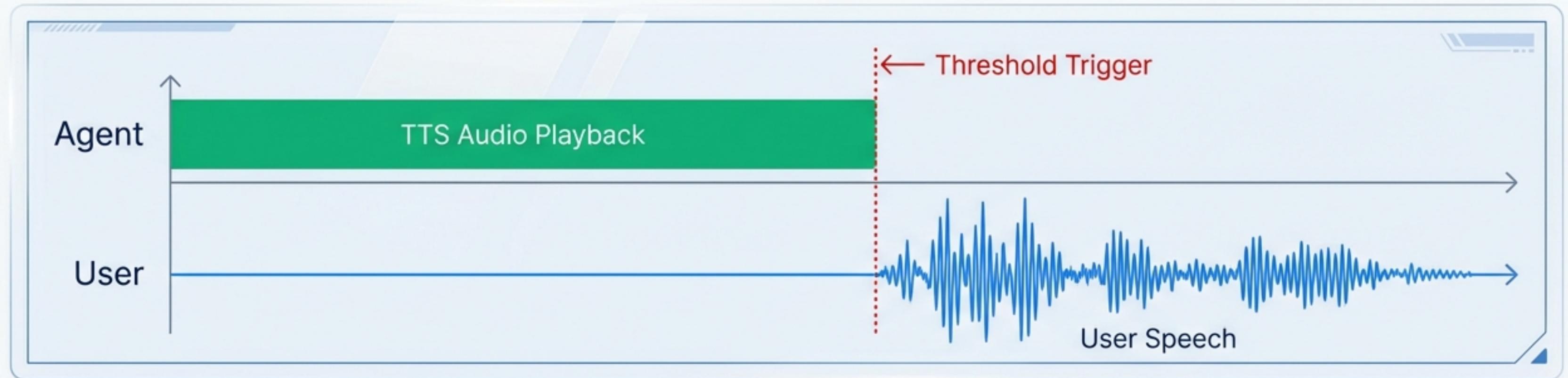
# Immersive User Experience (UI/UX)

The frontend leverages modern CSS glassmorphism and real-time JavaScript visualizers to create a responsive, high-end interface.

- **Glassmorphism:** translucent panels with `backdrop-filter: blur(10px)`.
- **Live Visualization:** 'Siri-style' multi-wave rendering.
- **Responsive Feedback:** Animated buttons and status badges.



# Conversational Mechanics: The "Barge-In" Logic



## Voice Activity Detection (VAD)

Tunable sensitivity slider (Range 5-50) with hysteresis logic (Start Threshold vs. Stop Threshold). Includes automated noise floor calibration to prevent false positives.

## Interruption Event

When user audio amplitude > threshold:

```
window.vad.bargeIn = true;  
audioPlayer.pause();  
audioPlayer.currentTime = 0;
```

# Global Input, Controlled Output

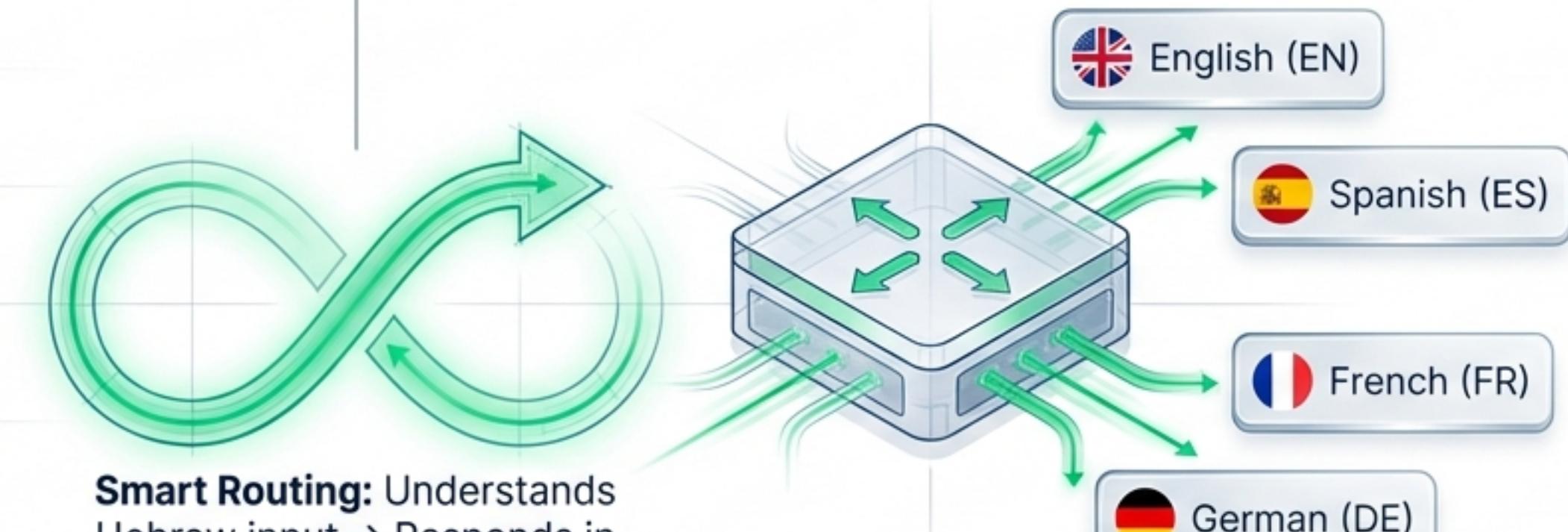
Decoupled logic for understanding vs. speaking.

## INPUT (Whisper ASR)



Auto-detects 99+ languages.  
Zero configuration required.

## OUTPUT (XTTS v2)

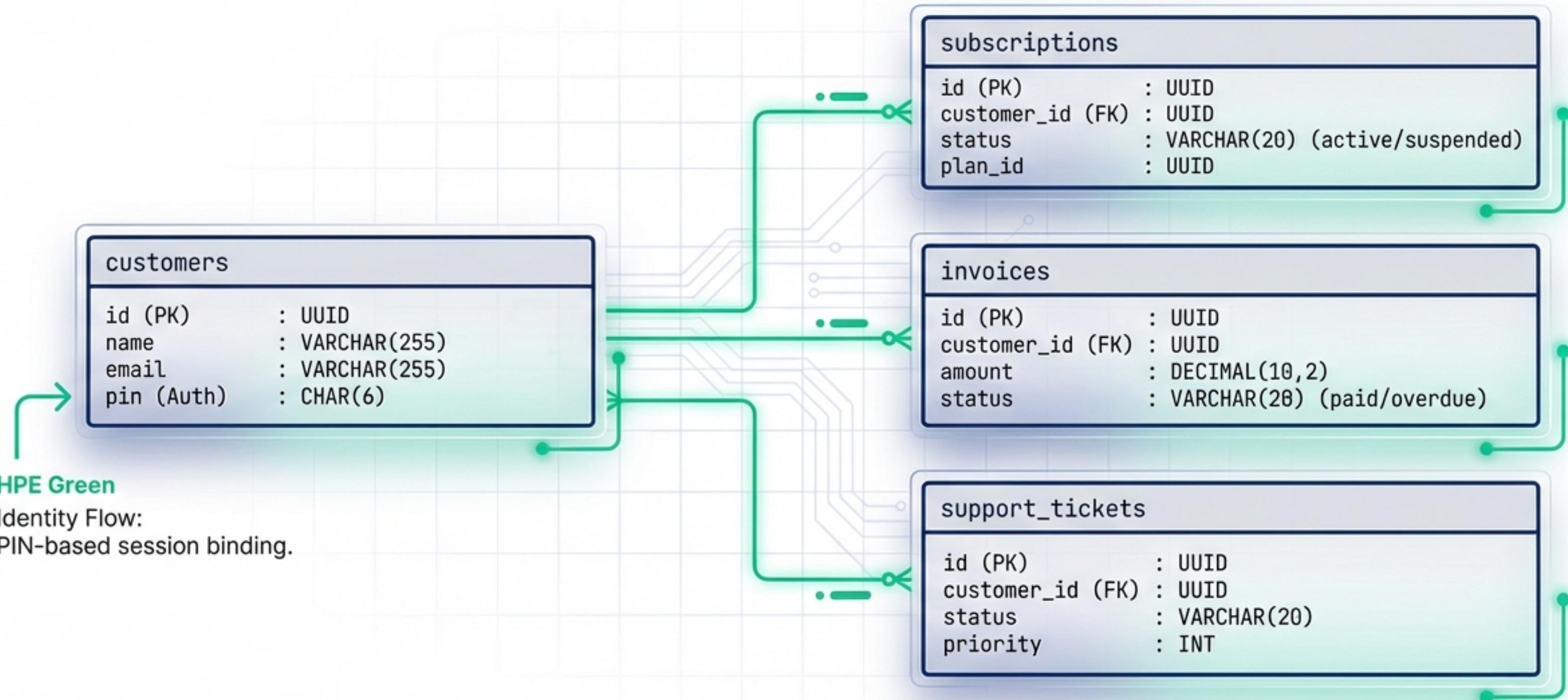


**Smart Routing:** Understands Hebrew input → Responds in English (if configured)

16+ Neural Voices. Backend enforces  
specific `response\_instruction` prompts.

# Persistent Memory & Identity Resolution

Schema definition for the `customer\_service` PostgreSQL database.



# The Agent's View: Dynamic Customer Context

Constructed in real-time via ``format_customer_info_html``, this is the exact data context provided to the AI.

 **John Smith** Since: Jan 2024

 john.smith@email.com

Open Tickets: 1 Overdue Invoices: 0

**Subscription**  
Premium Plan - \$49.99/mo  
**Active**

**Recent History**

- Cannot access premium features (High)
- Billing question (Low)

[Close Ticket](#) [Update Ticket](#) [Full History](#)

# Actionable Intelligence & Transactions

The agent performs CRUD operations, not just text generation.

## Ticket Management



**Action:**  
Create / Update / Delete

**Context:**  
Voice commands trigger  
SQL updates to  
'support\_tickets' table.

## Billing & Plans



**Action:**  
Query / Upgrade

**Context:**  
Process plan changes via  
'upgrade\_requests' table  
logic.

## Sentiment Analysis



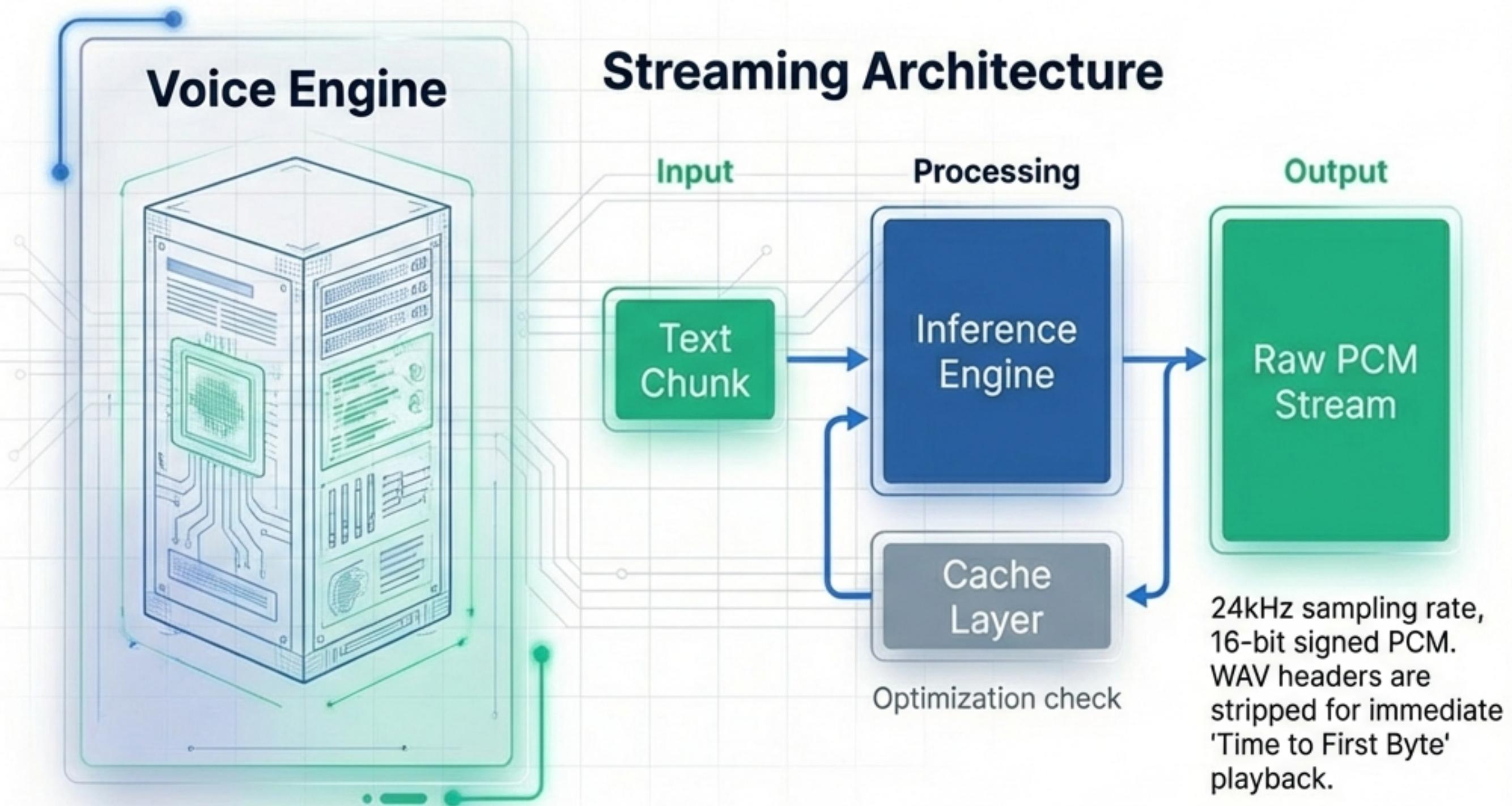
**Action:**  
Churn Detection

**Context:**  
Real-time logging to  
'sentiment\_alerts'.  
Triggers: '**Risk: High**',  
**'Intent: Cancellation'**.

# High-Fidelity Neural Synthesis (XTTS v2)

## Model Specs

- **Model:** XTTS v2 (Coqui)
- **Framework:** FastAPI / Python
- **Hardware:** Optimized for L40S GPU



# Voice Cloning & Persona Library

## The Speaker Roster



Daisy (US)



Adde (FR)



Suad (AR)



Kazuhiko (JP)

Cross-lingual capabilities: A cloned American voice can speak fluent Japanese without retraining.

## Cloning Technology



Reference  
Sample (<10s)



Endpoint: /speakers/upload



New Speaker  
Avatar

# Engineering for Reliability & Scale



# Observability & Performance Metrics

Real-time latency reporting and persistence.

## Total Round Trip Time (1.2s)

ASR Processing  
**0.3s**

LLM Token Gen  
**0.6s**

TTS First Chunk  
**0.3s**

```
{  
  "type": "latency_report",  
  "data": {  
    "asr_processing_time": 0.32,  
    "first_llm_token_time": 0.58,  
    "first_tts_chunk_time": 0.29  
  }  
}
```

Metrics are cached (\_metrics\_cache) and persisted to disk to track historical trends.

# Use Case: The ‘Unhappy Customer’ Flow



# The Evolution of Voice AI

From simple chatbots to context-aware,  
transactional, multilingual agents.

---

**Powered by HPE Private Cloud AI**