

FULL-STACK ORCHESTRATION ON HPE PRIVATE CLOUD AI

HPE AI Voice Agent v4.0

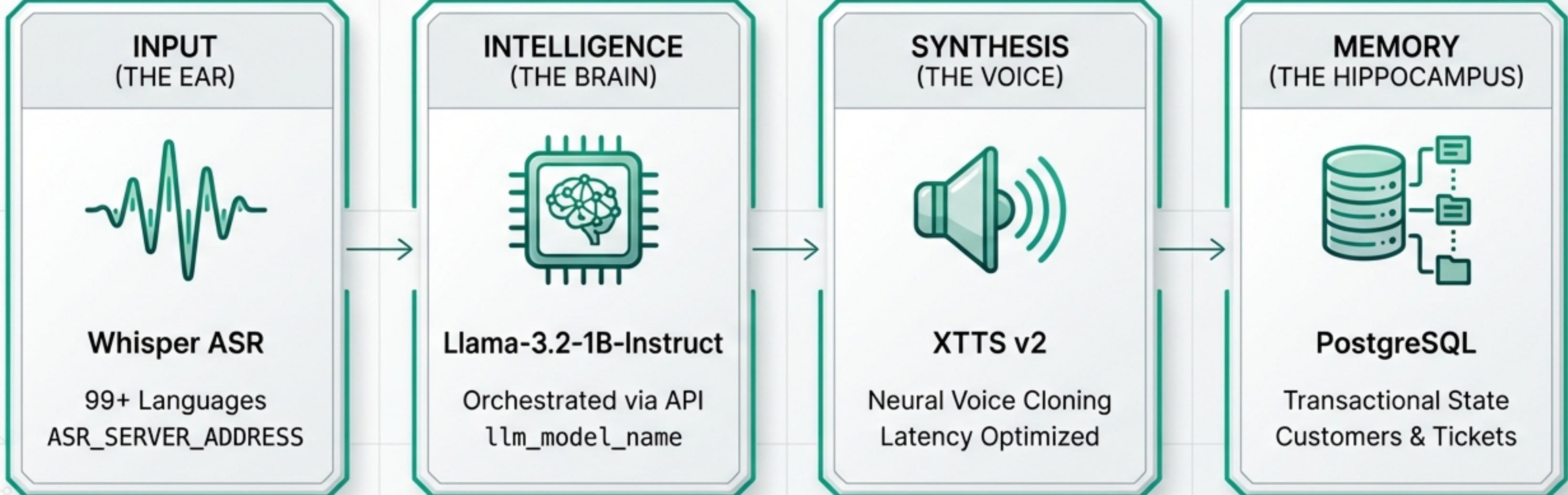
Architecture & Engineering
Deep Dive

METADATA.JSON

```
1 // SYSTEM MANIFEST
2 {
3     "version": "4.0.0 (Stable)",
4     "architecture": "Event-Driven WebSocket Pipeline",
5     "latency_target": "< 1000ms end-to-end",
6     "state_management": "PostgreSQL Persistent",
7     "deploy_target": "HPE Private Cloud AI"
8 }
```

A System of Systems: The v4.0 Tech Stack

Orchestrating four powerful engines into a unified agent.



Differentiation: Moving beyond 'Push-to-Talk' to continuous, hands-free 'Conversation Mode'.

Frontend Engineering: The 'Hands-Free' Experience in Inter Tight

Browser-based intelligent voice activity detection.

Client-Side VAD

JavaScript injection prevents server overload.

```
window.vad = {  
  thresholdStart: 20,  
  thresholdStop: 12, // Hysteresis  
  noiseFloor: 0  
};
```

Noise Calibration

Dynamically calculates noise floor to ignore background hum.

```
calibrateNoise() // Uses 30th percentile of samples
```

Canvas Rendering

Real-time frequency analysis.

```
analyser.getByteFrequencyData(dataArray);
```

AGENT INTERFACE

Conversation Mode

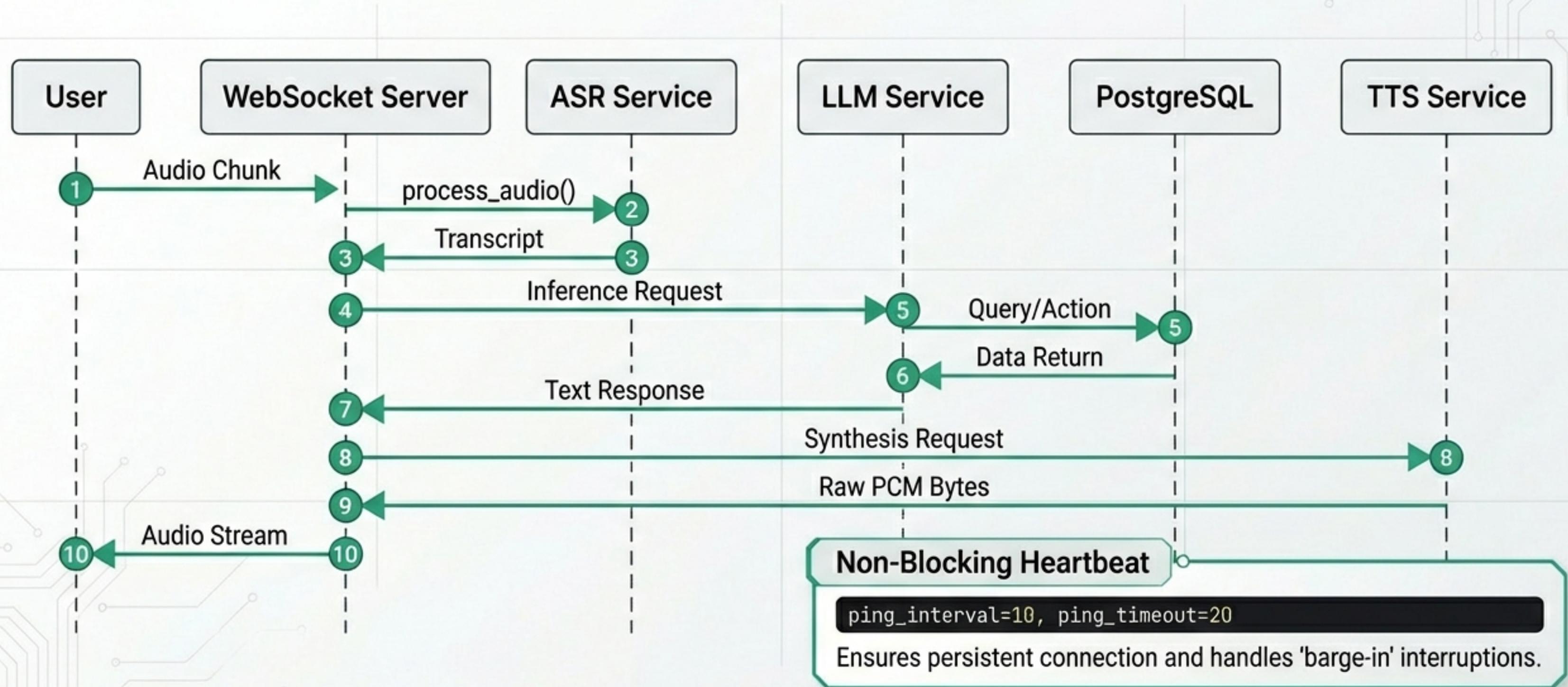


Active



The Intelligent Pipeline: Event-Driven Orchestration

Asynchronous data flow via WebSockets.



Decoupled Multilingual Architecture

Strategy: Understand Global, Speak Local.

INPUT (Understanding)



Supported: 99+

```
ASR_LANGUAGES = {  
    'he': 'Hebrew (עברית)',  
    'zh': 'Chinese (中文)',  
    'ar': 'Arabic (العربية)'  
}
```

LLM Context Injection

SYSTEM PROMPT:

"Language Rule: You MUST respond ONLY in [Target Language]. Even if user speaks Hebrew, reply in English."

OUTPUT (Speaking)



Supported: 16
(High Fidelity)

```
TTS_LANGUAGES = {  
    'en': 'English',  
    'es': 'Spanish',  
    'fr': 'French'  
}
```

The Brain: Dynamic Persona Injection

Role-playing engine defined by strict prompt engineering.

Smart Assistant

"template": "Keep it
SHORT - 1-3 sentences
max. Sound like a real
person, not a script."

Customer Service Pro

"template": "Warm and
professional... helpful
colleague, not
corporate robot."

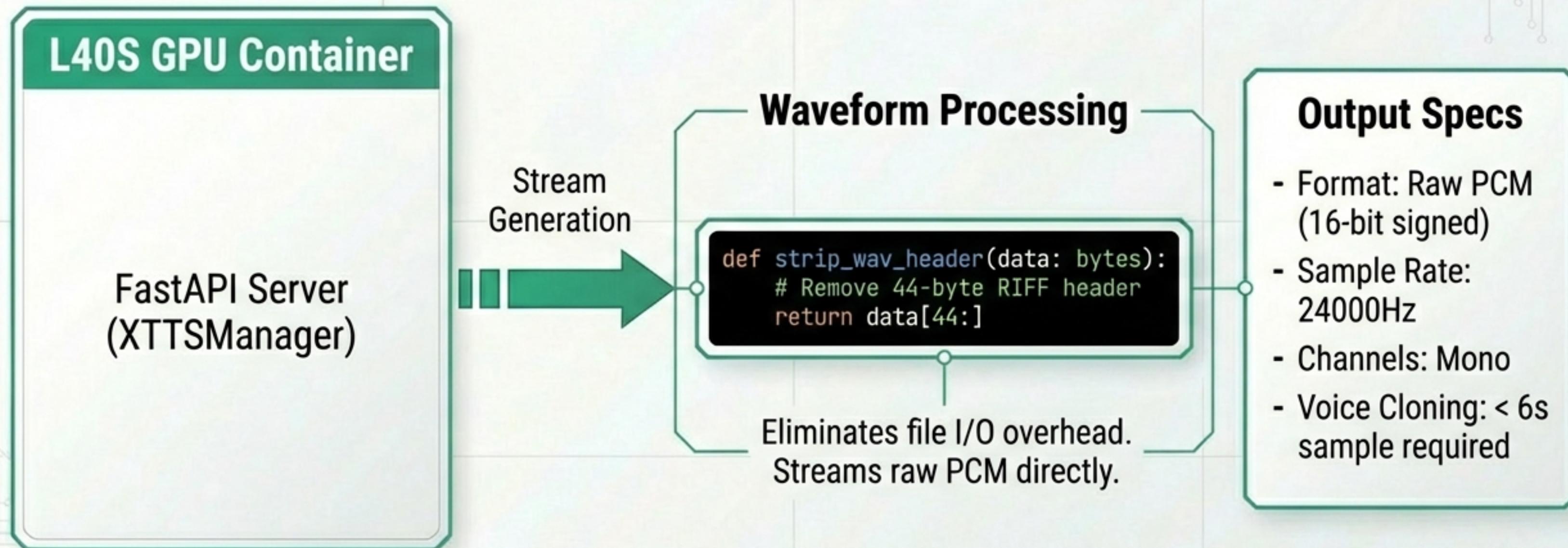
Business Advisor

"template": "Direct,
actionable insights...
confident but not
arrogant."

```
def build_final_prompt(persona, lang):  
    return AGENT_PERSONAS[persona]['template'] + \  
        TTS_LANGUAGES[lang]['response_instruction']
```

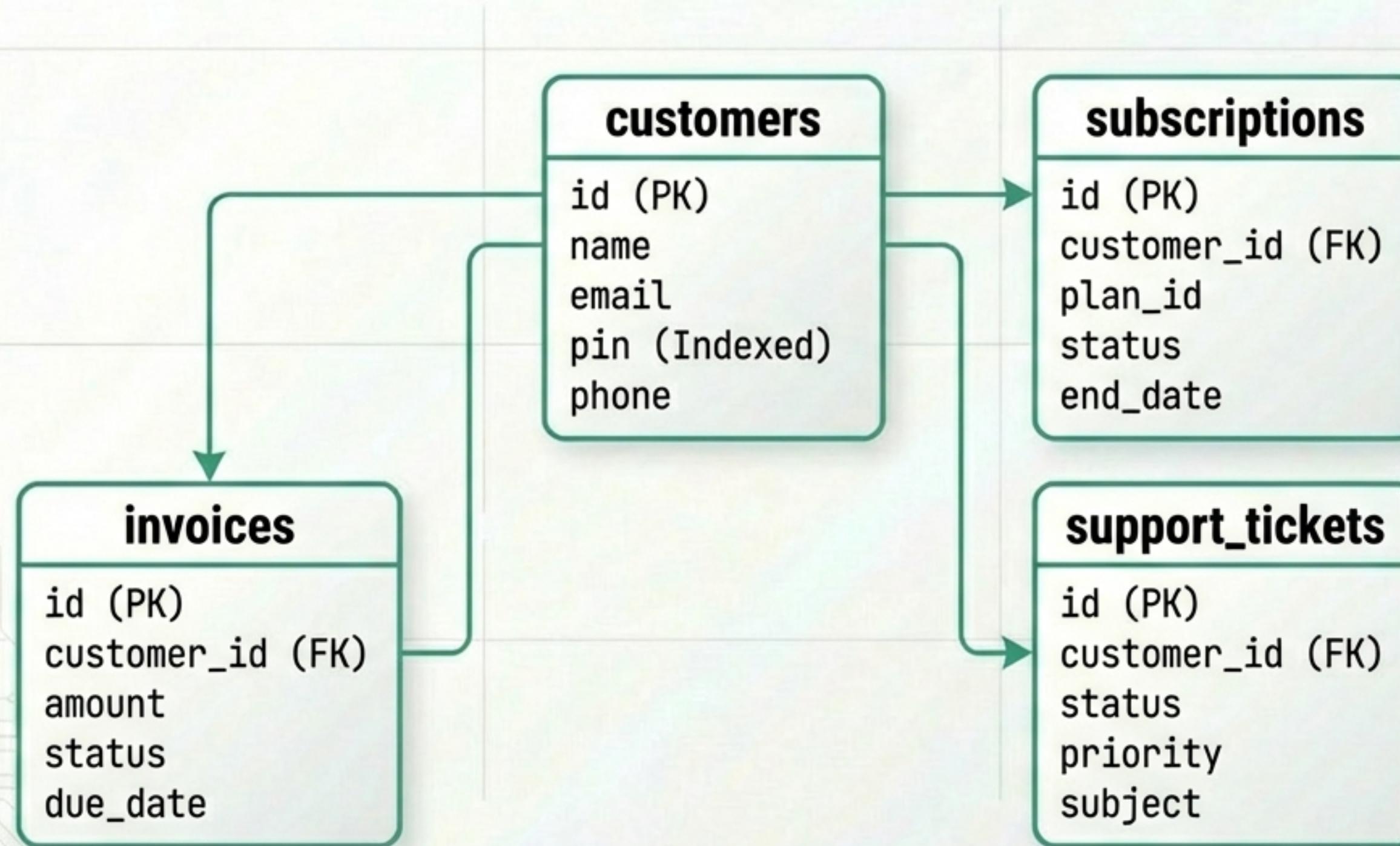
The Voice: High-Fidelity XTTs v2 Server

Engineering for low-latency audio synthesis.



The Memory: Stateful Database Integration

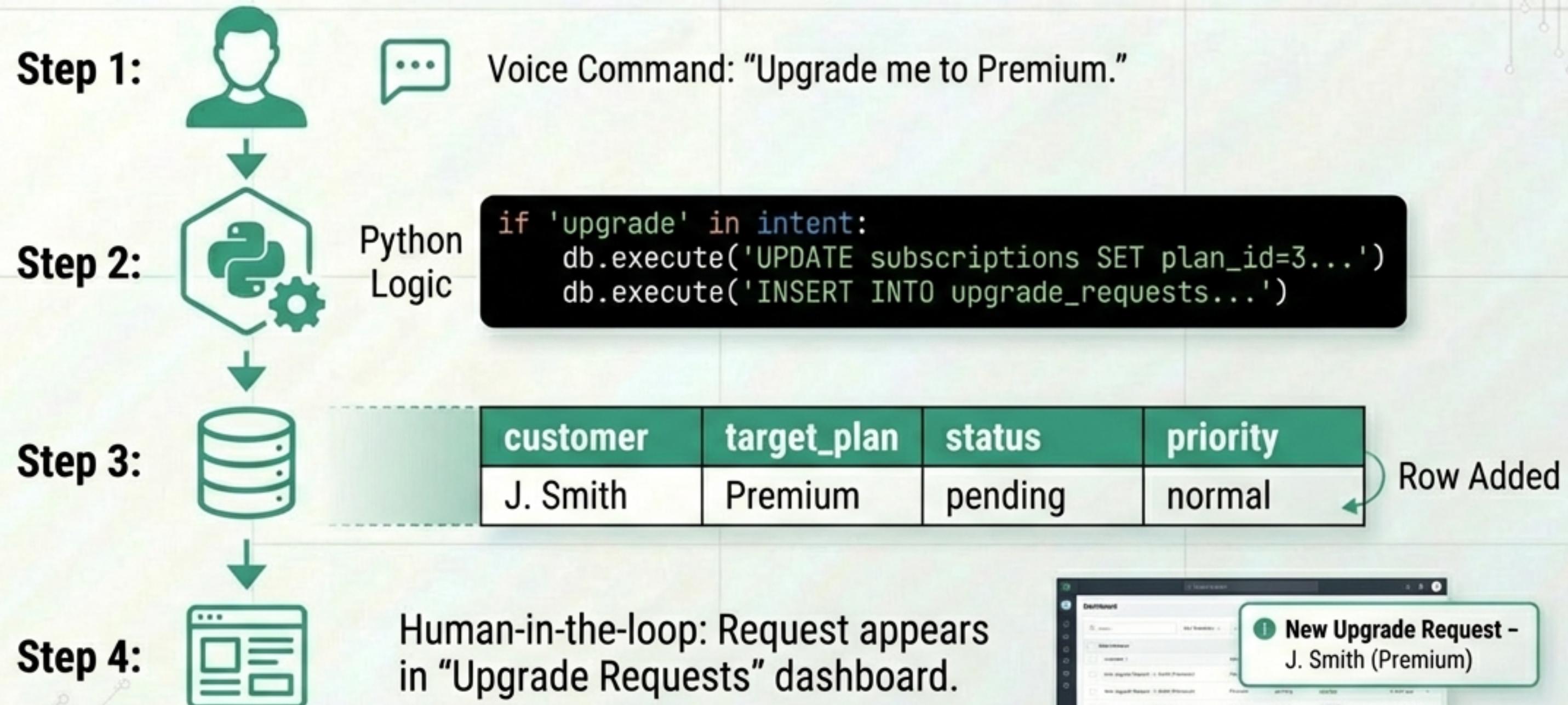
PostgreSQL Schema for persistent customer context.



Session Persistence:
Context survives connection drops via `'_customer_info_cache'` tied to session ID.

Transactional Voice Capabilities

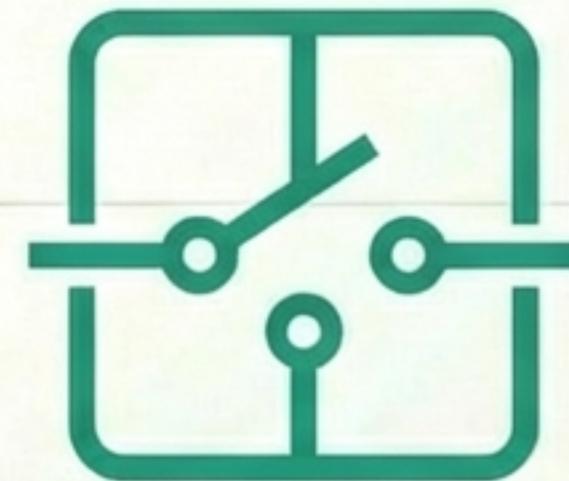
From conversation to execution: The Upgrade Request Workflow.



v4.0 Engineering: Stability & Resilience

Hardening the system for production loads.

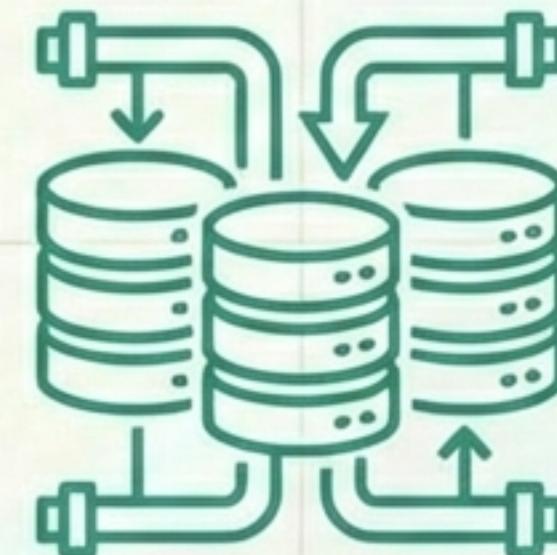
Circuit Breakers



Prevents cascading failures.

```
def check_service_status(url, timeout=10.0):  
    # Handle 5xx vs ConnectionRefused
```

Connection Pooling



Prevents socket exhaustion.

```
httpx.AsyncClient(  
    limits=httpx.Limits(max_keepalive=20)  
)
```

Health Checks



Real-time service monitoring.

```
endpoints = ['/health', '/v1/models']  
status = '✅ Online' if response.ok else '🔴'
```

Performance Optimization & Caching

Minimizing latency through multi-tiered strategies.

Memory Cache (RAM)

Customer Info & Session Context

```
_customer_info_cache = {} # TTL 3600s
```

Disk Persistence

Metrics & Analytics recovery

```
json.dump(_metrics_cache, f)
```

SQL Query Cache

Frequent read operations

```
SimpleCache(default_ttl=30)
```

Latency Report

ASR Processing : 0.42s

LLM First Token : 0.85s

TTS First Chunk : 0.35s

Total Round-Trip : 1.62s

Security & Authentication

Protecting customer data in voice workflows.

Voice Authentication



```
SELECT name FROM customers  
WHERE pin = %s  
AND phone = %s
```

Role-Based Access Control (RBAC)

Admin User

Full DDL Access (Create/Drop DB)

App User

Restricted DML Access (Select/Update)

Network Isolation: All internal traffic routed via Kubernetes DNS `.svc.cluster.local.`

Advanced Analytics: Sentiment & Churn Risk

Turning conversations into business intelligence.

sentiment_alerts

High Alert

Sentiment Score:
0.15
(Very Negative)

Churn Risk:
HIGH ●

Trigger Phrase:
"cancel my
subscription"

Action:
Manager Intervention
Required

Medium Alert

Sentiment Score:
0.45
(Frustrated)

Churn Risk:
MEDIUM ●

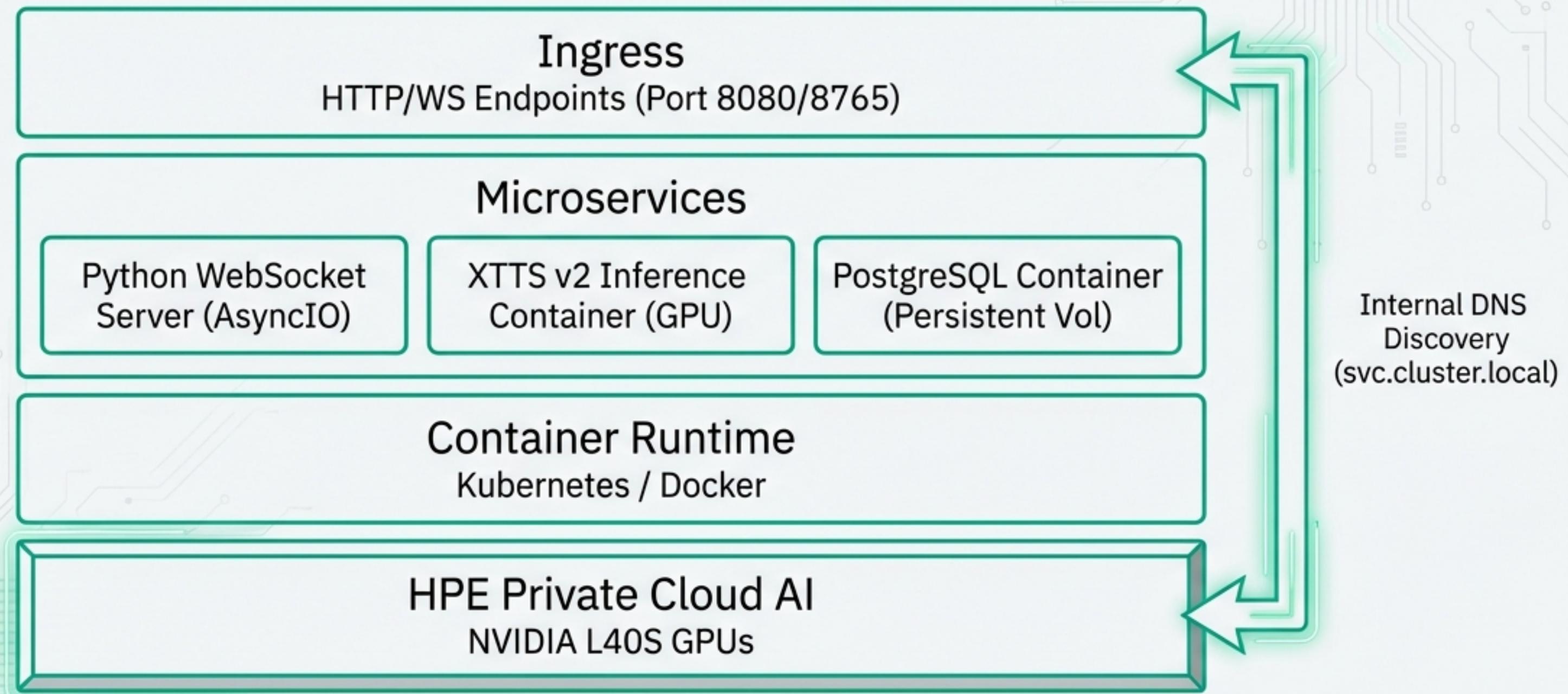
Trigger Phrase:
"too expensive"

Action:
Offer Discount

```
CREATE TABLE sentiment_alerts ("  
    sentiment_score DECIMAL(3,2),  
    churn_risk VARCHAR(20),  
    trigger_phrases TEXT  
");
```

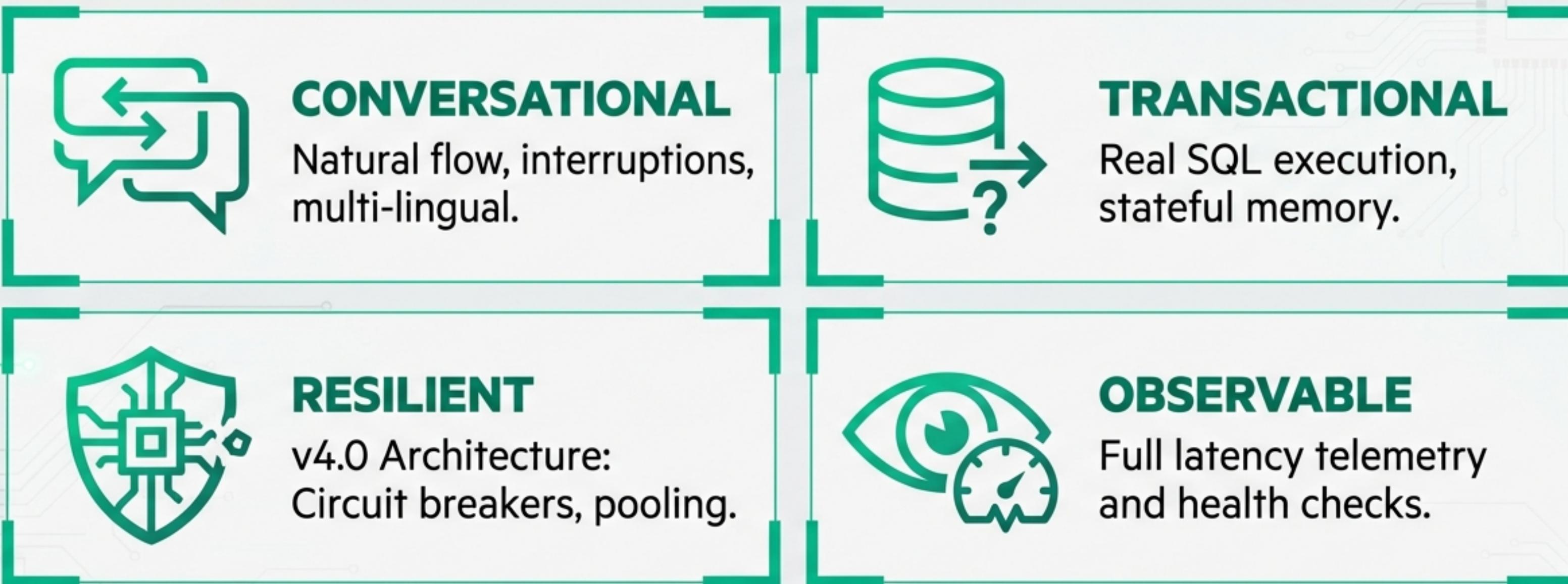
Infrastructure & Deployment

Orchestrated on HPE Private Cloud AI.



Summary: The “Full-Stack” AI Advantage

The shift from chatting with AI to working with AI.



HPE AI Voice Agent v4.0: Production Ready

HPE Confidential

NotebookLM