# Building a local AI server: Hardware, architecture, and implementation guide

The optimal local AI server for running 3 simultaneous LLM agents on a ~$5,000 budget centers on **an RTX 5090 (32GB) or dual RTX 3090 setup running 3× Llama 8B instances**—not 70B models, which require **35-42GB VRAM per instance** and simply won't fit on consumer hardware. For December 2024, the best strategy is waiting for the RTX 5090 launch (January 30, 2025) or purchasing a pre-built system with RTX 4090, as standalone 4090 cards are selling 40-75% above MSRP due to discontinued production. Starlink is unsuitable for trading bots due to **25-50ms latency and CGNAT restrictions,** ( Starlink ) though it works for remote access via Tailscale. iOS cannot support always-on listening—Apple's privacy architecture prohibits background microphone access—so dedicated hardware running local Whisper is required.

---

## RTX 50-series launches January 2025, reshaping the GPU landscape

The RTX 50-series was **announced at CES 2025 (January 6)** with retail availability beginning January 30, 2025. As of December 2024, these cards were not yet purchasable—only rumors existed.

| Model | VRAM | Memory Bandwidth | MSRP | TDP |
|---|---|---|---|---|
| **RTX 5090** | 32GB GDDR7 | 1,792 GB/s | $1,999 | 575W |
| **RTX 5080** | 16GB GDDR7 | 960 GB/s | $999 | 360W |
| **RTX 5070 Ti** | 16GB GDDR7 | 896 GB/s | $749 | 300W |
| **RTX 5070** | 12GB GDDR7 | 672 GB/s | $549 | 250W |

The RTX 5090's **77% memory bandwidth increase** over the RTX 4090 is transformative for LLM inference, which is memory-bound. Benchmarks show the 5090 achieving **213 tokens/second on Llama 3 8B** versus 127-128 tok/s for the 4090—a 67% improvement. The 32GB VRAM also enables running quantized 70B models that don't fit on the 4090's 24GB.

**Current RTX 4090 market reality** is challenging: prices range from $2,200-$3,500 (versus $1,599 MSRP) due to production discontinuation. Pre-built systems often offer better value—the Alienware Aurora R16 with i9-14900KF and RTX 4090 dropped to **$3,699 during holiday sales** (from $4,694). ( Tom's Hardware )
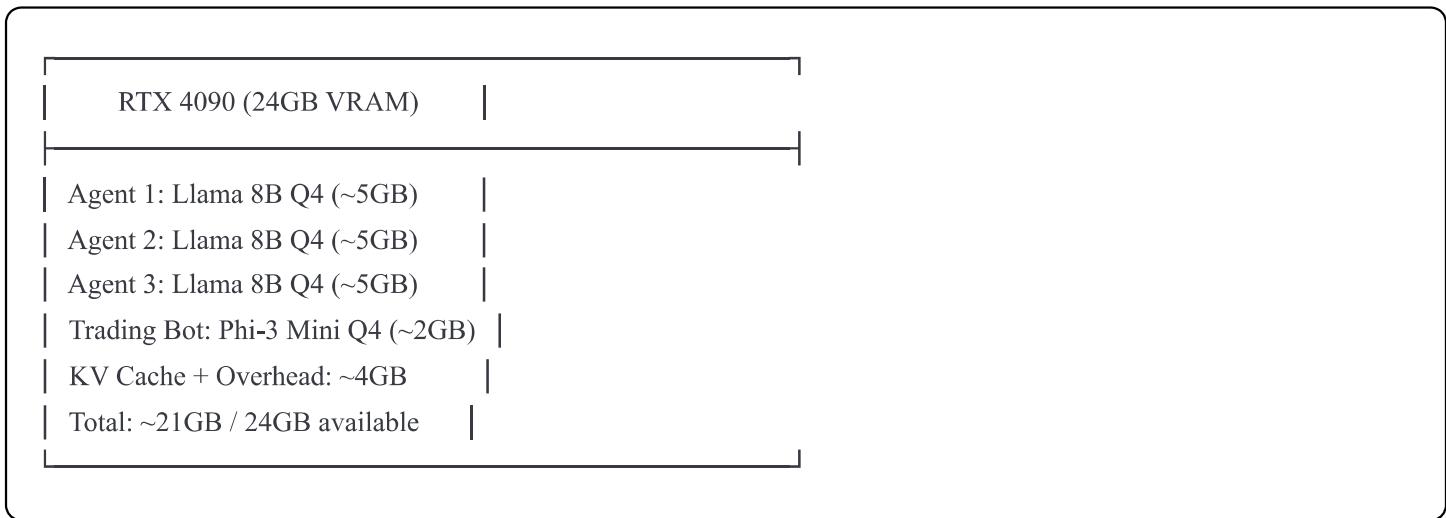
---

## Why 70B models won't work for your 3-agent use case

The mathematics of VRAM requirements eliminate running multiple 70B instances on consumer hardware. A single Llama 3.1 70B at 4-bit quantization requires **35-42GB VRAM** including KV cache overhead— ( Novita AI )exceeding even the RTX 5090's 32GB capacity.

| Model Size | FP16 | INT8 | 4-bit (Q4) |
|---|---|---|---|
| 8B | ~16 GB | ~8 GB | **~4-5 GB** |
| 13B | ~26 GB | ~13 GB | ~7-8 GB |
| 70B | ~140 GB | ~70 GB | **~35-42 GB** |

**The practical solution for 3 concurrent agents**: Run 3× Llama 3.1 8B Q4_K_M instances (~15GB total) plus a smaller model (Phi-3 Mini) for trading bots (~2GB), totaling approximately 21GB—comfortably within a 24GB RTX 4090 or 32GB RTX 5090.

Ollama 0.2+ supports concurrent requests via the OLLAMA_NUM_PARALLEL environment variable. (X) For production multi-agent deployments, **vLLM delivers 10-35× higher throughput** than Ollama through PagedAttention and dynamic batching. Red Hat benchmarks show vLLM achieving **793 tokens/second** versus Ollama's 41 tokens/second at scale. (Red Hat)

```
       RTX 4090 (24GB VRAM)      |

  Agent 1: Llama 8B Q4 (~5GB)    |
  Agent 2: Llama 8B Q4 (~5GB)    |
  Agent 3: Llama 8B Q4 (~5GB)    |
  Trading Bot: Phi-3 Mini Q4 (~2GB)  |
  KV Cache + Overhead: ~4GB      |
  Total: ~21GB / 24GB available  |
```

# Recommended $5,000 build configurations

**Option A: Wait for RTX 5090 (Best Performance, ~$4,000-4,500)**

| Component | Recommendation | Price |
|---|---|---|
| GPU | RTX 5090 32GB | $1,999 |
| CPU | AMD Ryzen 9 7950X | $400-450 |
| Motherboard | MSI X670E Carbon | $250-300 |
| RAM | 128GB DDR5-5600 (4×32GB) | $350-400 |
| Storage | Samsung 990 Pro 2TB NVMe | $180-200 |
| PSU | Corsair HX1500i (1500W Platinum) | $350-400 |
| Cooler | Corsair 360mm AIO | $150-200 |
| Case | Fractal Torrent | $150-200 |

**Option B: Dual RTX 3090 for Maximum VRAM (~$3,500-4,000)**

For 48GB total VRAM enabling larger models: two used RTX 3090s at $800-900 each provide 87% of RTX 4090 performance at a fraction of the cost. (The Ai Hardware) This configuration can run 70B quantized models that single-GPU setups cannot.

**Critical power considerations**: The RTX 5090's **575W TDP** requires a 1000-1200W PSU minimum. Running dual RTX 3090s (700W combined) demands 1500W+. Premium 80+ Platinum efficiency is essential for 24/7 operation.

---

# Network infrastructure must prioritize wired connections for trading

**Always use wired Ethernet for trading systems and AI servers.** WiFi introduces variable latency, jitter, and potential packet loss that disqualify it for latency-sensitive applications. Home trading setups should expect 10-50ms latency depending on location and ISP.

**Starlink is fundamentally unsuitable for hosting trading bots.** Current Starlink residential latency runs **25-50ms** with variable jitter—acceptable for browsing but problematic for competitive trading where professional co-located servers achieve sub-1ms latency. Additionally, Starlink's CGNAT prevents inbound connections without workarounds.

| Starlink Metric | Typical Range |
|---|---|
| Latency | 25-50ms |
| Download | 50-250 Mbps |
| Upload | **8-25 Mbps** (major limitation) |

**Recommended network architecture**:

- **Router**: Ubiquiti Cloud Gateway Fiber ($300-400) with three 10GbE ports and built-in VPN/IDS

- **Switch**: MikroTik CRS305-1G-4S+IN ($150) for 10Gbps backbone

- **Remote access**: Tailscale (free tier supports 100 devices) works through CGNAT with zero port forwarding

- **UPS**: CyberPower CP1500PFCLCD ($250) provides 10-15 minutes runtime for graceful shutdown

For trading bots requiring low latency, **host on a VPS near the exchange** (Chicago for CME, ~$50-200/month for sub-1ms latency) and access remotely via Starlink. This separates the latency-critical trading execution from your home network limitations.

---

## Always-on listening requires dedicated hardware—iOS won't cooperate

Apple's privacy architecture creates absolute barriers to always-on voice assistants on iOS. Third-party apps **cannot access Siri conversation history, phone call audio, or maintain background microphone access**. Even Apple's own iOS 18.1 call recording feature plays mandatory audible notifications to all parties and saves transcripts inaccessibly in Notes.

**What IS possible on iOS**: User-initiated foreground recordings, Apple's Speech Recognition API (with permission), and Shortcuts automation for specific bounded tasks.

**The solution is dedicated hardware running local Whisper**:

| Whisper Model | VRAM Required | Relative Speed |
|---|---|---|
| tiny | ~1 GB | ~10× faster |
| small | ~2 GB | ~4× faster |
| **turbo** | ~6 GB | ~8× faster |
| large-v3 | ~10 GB | 1× (baseline) |

The **turbo** model offers near-large accuracy with 8× faster inference—the optimal choice for real-time transcription alongside LLM workloads. Whisper.cpp runs efficiently on consumer GPUs, achieving **3,000+ words per minute on an RTX 4090**.

**Privacy-first architecture pattern**:

1. **Wake word detection** (Porcupine, ~1MB RAM) runs continuously—audio never leaves device

2. Upon wake word, buffer **30-60 seconds of audio**

3. **Local Whisper transcription** processes audio

4. Extract action items, **delete raw audio immediately**

5. Store only encrypted text; maintain audit log

Porcupine achieves **97%+ detection accuracy with <1 false alarm per 10 hours** (GitHub) and runs on Raspberry Pi through desktop GPUs.

---

## Multi-agent architecture patterns for shared infrastructure

Running personalized AI agents on shared compute requires **namespace-based isolation**, **priority scheduling**, and **per-agent memory management**.

**Recommended stack for 3 Business Partner agents**:

- **LLM backend**: vLLM (10× throughput vs Ollama) with quantized Llama 8B

- **Agent framework**: CrewAI for simplicity, (Helicone) LangGraph for production control (LangChain)

- **Memory**: Letta (formerly MemGPT) for per-agent memory blocks with shared organizational context

- **Vector DB**: ChromaDB self-hosted with namespace-per-agent pattern

- **Persistence**: PostgreSQL + pgvector for unified storage

**Per-agent context isolation** works through namespacing in vector databases. Each agent gets its own namespace for RAG retrieval, preventing context bleed between personalities:

```python
# Query retrieves only from agent's namespace
results = index.query(
    vector=query_embedding,
    namespace=f"agent_{agent_id}",
    top_k=5
)
```

**System prompt management** maintains distinct personalities: each agent loads its unique role definition, behavioral guidelines, tool access permissions, and persona blocks at session start. Memory blocks persist across sessions, enabling agents to "remember" preferences and prior interactions. (Letta)

**Resource scheduling** prevents any agent from monopolizing the GPU through priority queuing. Assign higher priority to time-sensitive agents (trading-related queries) and lower priority to research tasks. vLLM's continuous batching naturally handles concurrent requests, (House Of FOSS) while Ollama requires explicit (OLLAMA_NUM_PARALLEL) configuration. (GitHub)

---

# Conclusion: Practical implementation path

For a December 2024 start, the optimal path is purchasing a pre-built RTX 4090 system (Alienware Aurora at $3,699 offers best value) while waiting for RTX 5090 availability in late January 2025. Run 3× Llama 8B instances rather than attempting 70B models—the VRAM math doesn't support multiple large model instances on consumer hardware.

Key technical decisions that differentiate success from frustration:

- **vLLM over Ollama** for multi-agent deployments—the throughput difference is dramatic (Arsturn)

- **Dedicated VPS for trading bots** near exchange, accessed remotely—Starlink's latency disqualifies it from hosting time-sensitive operations

- **Whisper turbo model** balances accuracy and speed for transcription alongside LLM workloads

- **Tailscale for remote access** works through any network topology including CGNAT

- **Namespace-based vector DB isolation** prevents agent context bleed while sharing infrastructure

- **8B parameter models** running concurrently outperform a single 70B model for multi-agent scenarios on consumer GPUs

The $5,000 budget comfortably covers a production-capable system with headroom for UPS, networking, and potential GPU upgrade when RTX 5090 stock normalizes.