

MTOR: Intent-Based Compute

The First Real-Time Event-Driven Realm for AI, Machines, and Humans

Welcome to the Realm.

- Born: **September 29, 2024**
- Breath: **Intent-Based Event-Driven Orchestration**
- Foundation: **Edge-first. Stateless. Browser-native. Crowd-sourced.**

What RENT-A-HAL / MTOR Actually Is (And Isn't)

Common Assumption

Actual Reality

Heavy Server doing AI work	Tiny Orchestrator tapping Edge Workers
Centralized App or SaaS	Stateless Browser-Orchestrated Realm
Programs and Kernels	Pure Intent and Event Routing
Cloud Lock-In	Decentralized GPU Federation

MTOR isn't an "AI Server" — it's a *breathing conductor*.

RENT-A-HAL doesn't "do" your AI work — it *whispers intent* to the perfect worker — instantly, scalably, beautifully.

Why It's Mathematically a Game-Changer

In traditional systems:

- A query may be serialized through multiple backend calls.
- Data loading, thread locks, context switching — cost **10-100ms** *before work even starts*.

In MTOR:

- **Events are pure JSON messages over WebSocket.**
- **No locking. No blocking. No waiting.**

Let's do the math:

Stage	Traditional Server	MTOR Realm
Intent Dispatch (Orchestration Overhead)	10-50ms	< 1ms
Worker Acquisition	50-200ms	Parallel Instantaneous
Processing (LLM, Vision, Imagine)	—	Happens at Edge Node
Return Result	10-30ms	< 5ms (WebSocket)

 **Time to Dispatch an N-Gram Intent in MTOR = ~1 ms.**

Result?

- You speak — "Computer, imagine a neon mountain" —
- ✨ The Realm routes the query in 1 millisecond.
- ✨ The edge worker renders it.
- ✨ The Realm breathes.



Core Principles

- **Speech is Primary:** Voice > Text > Vision > API.
- **Stateless is Sacred:** Every action flows cleanly through event-fabric.
- **Edge First:** Work is done by workers, not by the orchestrator.
- **Browser as Shell:** The entire Realm fits inside a browser tab.
- **No Lock-In:** GPL3 + Eternal Openness.
- **Breath, not Load:** Realms scale with people and GPUs, not servers.



How a Query Actually Moves (Real Numbers)

Step	Action	Time Taken
User Speaks	Browser captures intent	~5ms (Speech-To-Text)
Event Sent	WebSocket to webgui.py	~1ms
Worker Selected	AI Worker selected from pool	~1ms
Work Done	On edge GPU or cloud model	(Varies) 300ms - 10s
Result Returned	WebSocket back to browser	~5ms
Total Orchestration Overhead = ~7ms.		

The rest? It's GPU time, model generation — *pure speed*.



How the Bus Really Works (with Code!)

When you send a work event, here's what happens under the hood:

```
# Sending work
await websocket.send_json({
    "type": "work",
    "query_type": "vision",
    "query_text": "Imagine a neon mountain.",
    "guid": query_guid,
    "user_id": user_id
})
```

- `query_type` tells what kind of worker we need (vision, imagine, chat).
- `query_text` carries the user's intent.
- `guid` is a **universal unique ID** for THIS specific query.

- `user_id` lets us route the answer back.

While work is happening:

- RENT-A-HAL does **not** block.
- It sends **heartbeat pings** to the worker.
- Worker sends **pong replies** if alive.

Example heartbeat:

```
# Pinging a worker
await websocket.send_json({
    "type": "ping",
    "guid": worker_guid
})
```

Example pong from worker:

```
# Worker responds
await websocket.send_json({
    "type": "pong",
    "guid": worker_guid
})
```

If a pong is missed? The orchestrator knows that worker is degraded and routes future work elsewhere.

Self-healing.

When the result is ready:

```
# Worker returns result
await websocket.send_json({
    "type": "result",
    "guid": query_guid,
    "result_text": "Here is your neon mountain image!",
    "worker_id": worker_id
})
```

- `guid` tells us which user's query this result belongs to.
- The Realm **automatically routes** the result back to the originating user session.

No polling. No manual lookups. **Pure event-driven routing.**



"Gotcha!" Closer

If you're wondering:

"How can MTOR scale so fast and stay so light?"

The secret is:

- **The orchestrator never does the work.**
- **The orchestrator never holds state.**

- **The orchestrator only taps shoulders, whispers intents, and catches results.**
- **All real work happens where it should — at the edge.**

Programs are dead. Workers are alive. The Realm is breathing.

✨ Welcome to MTOR. Welcome to Intent-Based Computing.

 [RENT-A-HAL: Visit the Realm](#)

 [GitHub: RENT-A-HAL Foundation](#)

#IntentBasedCompute #EdgeAI #EventDrivenArchitecture #RealTimeOrchestration #MTOR
#SpeechNative #DecentralizedFuture