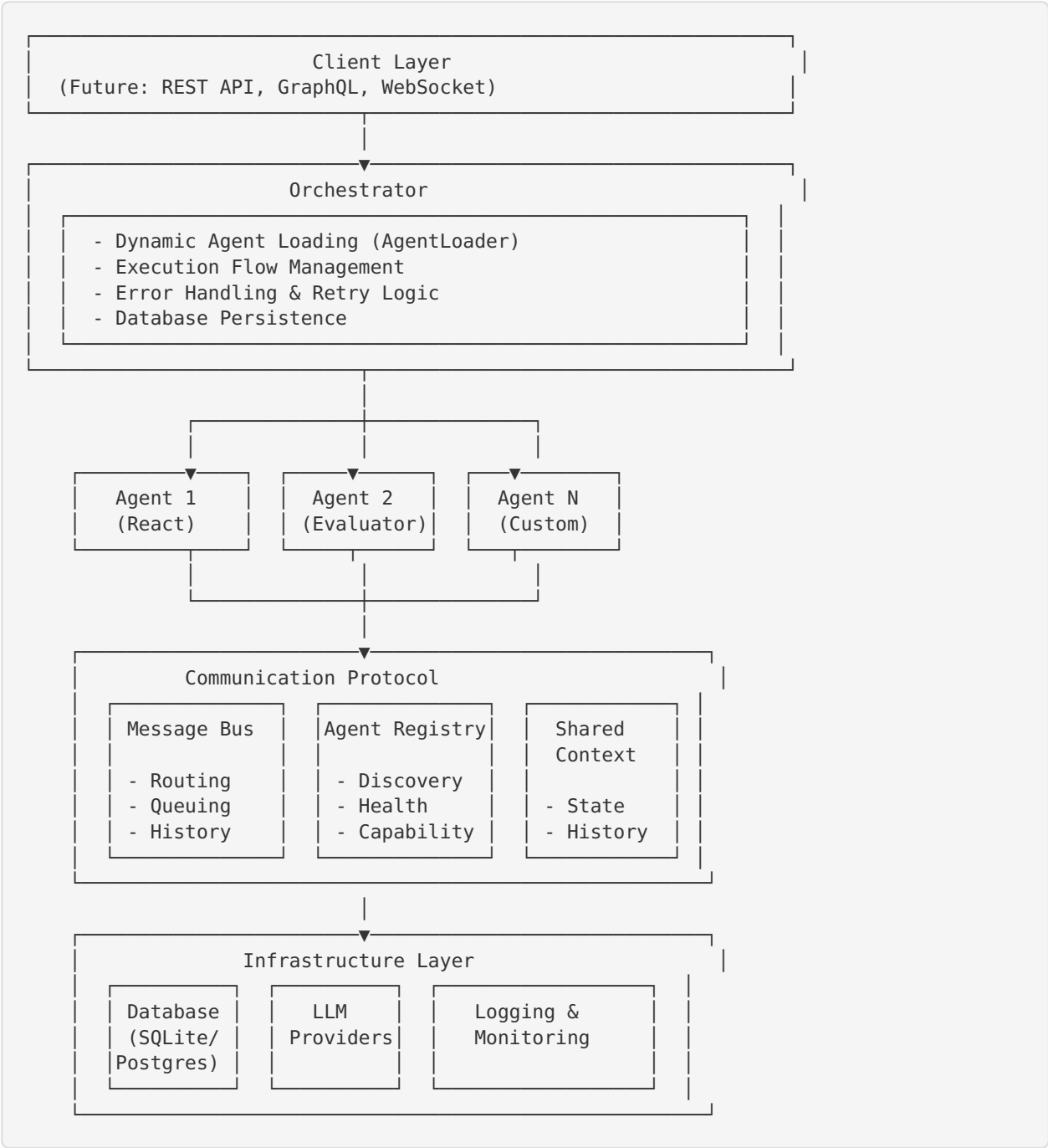


Architecture Documentation

System Architecture

High-Level Overview



Core Components

1. Communication Protocol (★ CRITICAL)

The communication protocol is the foundation that enables the modular business model.

Message Bus

Responsibility: Route messages between agents

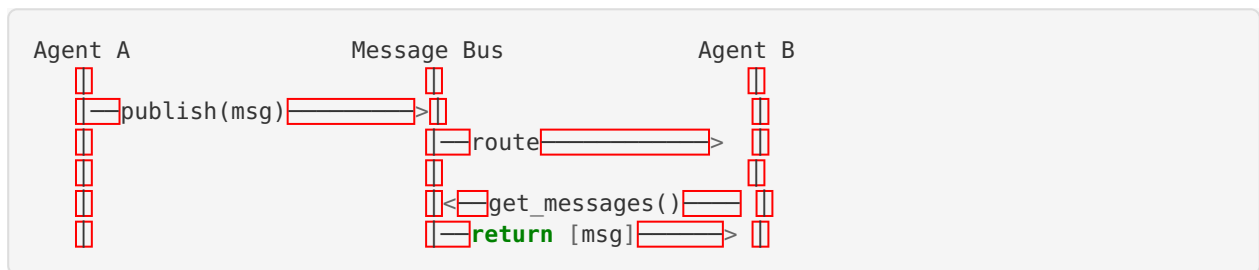
Key Features:

- Direct messaging (point-to-point)
- Broadcast messaging (one-to-many)
- Message queuing with priority
- Subscription-based routing
- Message history for audit

Implementation Details:

- Thread-safe with locks
- In-memory queues (deque) per agent
- Configurable queue size and retention
- Handler registration for callbacks

Message Flow:



Agent Registry

Responsibility: Service discovery and health monitoring

Key Features:

- Agent registration with capabilities
- Capability-based search
- Health monitoring via heartbeats
- Load balancing support (least busy agent)

Data Structure:

```

{
  "agent_name": AgentInfo(
    name="agent_name",
    agent_type="react",
    capabilities=["reasoning", "planning"],
    status="active",
    last_heartbeat=datetime.now(),
    message_count=42
  )
}
  
```

Shared Context

Responsibility: Shared state management

Key Features:

- Global state (accessible to all)

- Namespaced state (per-agent private)
- State history tracking
- Watch mechanism for notifications
- Thread-safe operations

State Organization:

```
Global State:
- task: "Analyze data"
- run_id: "uuid-123"
- status: "running"

Agent Namespaces:
agent_1:
- memory: {...}
- internal_state: {...}
agent_2:
- cache: {...}
```

2. Orchestrator

Responsibility: Manage multi-agent execution

Execution Strategies:

1. **Sequential:** Agents run one after another

```
Agent1 → Agent2 → Agent3 → Result
```

2. **Parallel:** Agents run simultaneously

```
Agent1 ↘
Agent2 → Aggregator → Result
Agent3 ↗
```

3. **Adaptive:** Dynamic execution based on outputs

```
Agent1 → Decision → Agent2 or Agent3 → Result
```

Error Handling:

- Retry logic with exponential backoff
- Graceful degradation (continue with other agents)
- Comprehensive error logging
- Database persistence of failures

3. Agent Loader

Responsibility: Dynamic agent discovery and instantiation

Discovery Process:

1. Scan `agents/` directory
2. Import each Python module
3. Find classes inheriting from `BaseAgent`
4. Register in internal registry
5. Instantiate on demand

Benefits:

- Zero configuration for new agents
- Hot-reload capability (development)

- Validation of agent types
- Metadata extraction (capabilities, docs)

4. LLM Abstraction Layer

Responsibility: Unified interface to multiple LLM providers

Provider Architecture:

```

BaseLLMProvider (Abstract)
├── OpenAIProvider
│   └── GPT-4, GPT-3.5-turbo
├── AnthropicProvider
│   └── Claude 3 (Opus, Sonnet, Haiku)
└── CustomProvider (extensible)
  
```

Unified Interface:

- `invoke()` : Synchronous generation
- `invoke_streaming()` : Streaming generation
- `count_tokens()` : Token counting
- Automatic retry logic
- Error normalization

5. Database Layer

Models:

```

















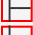


Tenant (Multi-tenancy)
├── Project (Grouping)
│   └── Run (Execution)
│       ├── AgentRun (Per-agent execution)
│       └── AgentMessage (Communication log)
  
```

Key Design Decisions:


- UUID primary keys for distributed systems
- Timestamps for audit trail
- JSON columns for flexible metadata
- Enum types for status fields
- Cascade deletes for data integrity

Data Flow


Typical Execution Flow

1. Client Request
 -  `Orchestrator.run(task, project_id)`
2. Initialization
 -  `Create Run` record in DB
 -  `Load` agents via AgentLoader
 -  Register agents with Protocol
 -  Initialize shared `context`
3. Agent Execution (Sequential)
 - For** each agent:
 -  Create AgentRun record
 -  Update agent status **to** "busy"
 -  Execute agent.execute()
 -  Agent may invoke LLM
 -  Agent may send messages
 -  Agent may access shared state
 -  Agent **returns** output
 -  Update AgentRun with output
 -  Update agent status **to** "idle"
 -  Add output **to** context
4. Finalization
 -  Update `Run` status **to** "completed"
 -  Persist all messages **to** DB
 -  **Return** aggregated results
 -  Cleanup
5. Communication (Parallel)

Agent A



Agent B



```

sequenceDiagram
    participant A as Agent A
    participant B as Agent B
    A->>B: send_message()
    B-->A: respond()
  
```

Scalability Considerations

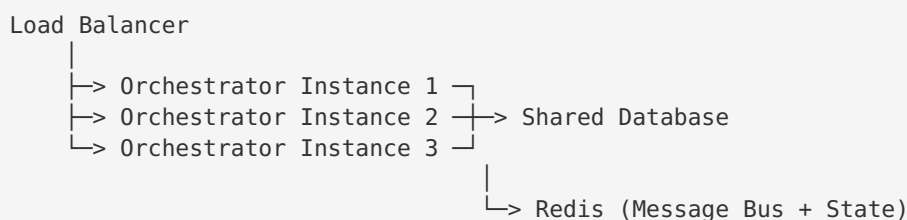
Current Implementation (Single Process)

- In-memory message bus
- Thread-safe operations
- SQLite database
- Suitable for: Development, demos, small deployments

Production Scaling Path

1. **Database:** PostgreSQL with connection pooling
2. **Message Bus:** Redis pub/sub or RabbitMQ
3. **Orchestrator:** Multiple instances with shared DB
4. **Agents:** Distributed workers (Celery, Temporal)
5. **State:** Redis for shared context
6. **Monitoring:** Prometheus + Grafana

Horizontal Scaling



Security Architecture

Multi-Tenancy

- Tenant ID in all database queries
- Row-level security (future: PostgreSQL RLS)
- API key per tenant (future)

Data Protection

- Environment variables for secrets
- No hardcoded credentials
- Encrypted connections (TLS)
- Audit logging of all operations

Input Validation

- Pydantic models for all inputs
- Type checking at runtime
- Sanitization of user inputs
- Rate limiting (future)

Extension Points

Adding New Agents

1. Create `agents/my_agent.py`
2. Inherit from `BaseAgent`
3. Implement `execute()` method
4. Define `CAPABILITIES`
5. Done! (Auto-discovered)

Adding New LLM Providers

1. Create provider class inheriting `BaseLLMProvider`
2. Implement required methods
3. Register with `LLMFactory`
4. Configure in settings

Adding New Communication Patterns

1. Extend `MessageType` enum
2. Add handler in `MessageBus`
3. Update `CommunicationProtocol` if needed

Performance Characteristics

Message Bus

- **Throughput:** ~10,000 messages/sec (in-memory)
- **Latency:** <1ms for local delivery
- **Memory:** $O(n)$ where $n = \text{queue size} \times \text{agents}$

Agent Execution

- **Overhead:** ~10-50ms per agent (without LLM)
- **LLM Latency:** 1-10 seconds (depends on provider)
- **Retry Overhead:** Exponential backoff (2s, 4s, 8s)

Database

- **SQLite:** Suitable for <100 concurrent runs
- **PostgreSQL:** Suitable for 1000+ concurrent runs
- **Write Latency:** ~5-10ms per record

Monitoring & Observability

Metrics to Track

1. **Agent Metrics**
 - Execution time per agent
 - Success/failure rate
 - Token usage (LLM)
2. **Communication Metrics**
 - Messages sent/received per agent
 - Message queue depth
 - Message delivery latency
3. **System Metrics**
 - Active agents
 - Concurrent runs
 - Database query time

Logging Strategy

- **Structured Logging:** JSON format for parsing
- **Log Levels:** DEBUG, INFO, WARNING, ERROR
- **Context:** Run ID, Agent name in all logs
- **Retention:** Configurable (default: 30 days)

Future Enhancements

1. **Real-time Communication:** WebSocket support
2. **Async Execution:** Full async/await support
3. **Distributed Tracing:** OpenTelemetry integration
4. **Advanced Routing:** Content-based routing
5. **Agent Marketplace:** Plugin system for third-party agents

6. **Visual Workflow Builder:** Drag-and-drop agent composition
7. **A/B Testing:** Compare different agent configurations
8. **Cost Optimization:** LLM call caching and batching

Design Principles

1. **Modularity:** Each component is independent and replaceable
 2. **Extensibility:** Easy to add new agents, providers, features
 3. **Reliability:** Comprehensive error handling and retries
 4. **Observability:** Extensive logging and metrics
 5. **Simplicity:** Clear abstractions, minimal magic
 6. **Production-Ready:** Multi-tenancy, security, scalability
-

This architecture enables the **modular B2B business model** where clients can select any combination of agents, and they will communicate seamlessly through the robust communication protocol.