



Economical and Technical Assessment

Analysed Source Code: Code Assessment

Document Date: June 09, 2026

Platform:

Client / Applicant

Legal entity name:

Registered address:

Tax Identification Number:

Software name:

Existing trade mark:

Company web:

Applicant Name:

Applicant Email:

Request date and time: 09-06-2026 - 02:06:43





TABLE OF CONTENTS

1. Executive Summary

2. Platform Overview

- 2.1 Functional Description
- 2.2 Technical Architecture
- 2.3 Technology Stack
- 2.4 Third-Party Integrations

3. Production Readiness Assessment

- 3.1 Overall Score: 76/100
- 3.2 Detailed Breakdown

4. Development Investment Estimation

- 4.1 Effort Analysis
- 4.2 Team & Timeline
- 4.3 Cost Estimation
- 4.4 Codebase Metrics
- 4.5 Cloud Infrastructure & Maintenance Cost

5. Findings Summary

- 5.1 Critical Issues (Must Fix)
- 5.2 Warnings (Should Fix)
- 5.3 Recommendations (Nice to Have)
- 5.4 Strengths

6. Conclusion

- 6.1 Overall Assessment Summary
- 6.2 Readiness for Production / Scale
- 6.3 Key Areas Requiring Attention
- 6.4 Suggested Prioritization of Improvements



Technical Assessment Report: Hermes Agent

Assessment Date: 1 January 2026

Platform: Hermes Agent (Multi-Platform AI Agent Framework)

Repository: NousResearch/hermes-agent

Assessment Type: Production Readiness Certification for Venture Capital Due Diligence

1. EXECUTIVE SUMMARY

Hermes Agent is a sophisticated multi-platform AI agent framework demonstrating strong engineering practices with comprehensive documentation, extensive test coverage, and robust security posture. The codebase exhibits high architectural complexity with 1.36M LOC across Python and TypeScript, supporting 20+ messaging platform integrations through a well-designed plugin architecture. Security is treated seriously with clear documentation of trust boundaries, OS-level isolation as the only trusted boundary, and credential management via environment variables and credential files.

The overall production readiness score is **76/100 (Grade B, Level: Good)**. This assessment reflects a mature codebase with solid foundations in code quality, documentation, and security practices, though some areas require attention before large-scale production deployment. The platform demonstrates evidence of significant development investment with strong CI/CD practices including linting, parallelised test execution, and dependency scanning.

Key strengths include comprehensive security documentation with clear trust model definitions, extensive test coverage across unit, integration, and end-to-end tests running in CI, well-documented codebase with README, CONTRIBUTING guide, and extensive developer documentation, plugin architecture enabling extensibility without modifying core code, strong CI/CD pipeline with linting (ruff, ty), tests, and security scanning, and Docker-based deployment with s6-overlay supervision and proper privilege separation.

Critical risks requiring immediate attention include the large codebase (1.36M LOC) with high complexity presenting maintenance challenges, 252 dependencies increasing supply chain attack surface despite dependency scanning present in CI, and extensive platform



integrations (20+ messaging platforms) increasing security surface area. The platform would benefit from implementing explicit test coverage gates in CI with minimum thresholds (recommended 80%), adding mutation testing for critical security components, documenting architectural decision records (ADRs) for major design decisions, and adding structured logging with correlation IDs for distributed tracing across gateway platforms.

Estimated development investment to date: The software required approximately **28,800 hours** of development effort, representing a team of 8 developers over 18 months, with an estimated cost range of **€2,692,800 – €3,643,200 EUR**. This valuation reflects the retroactive cost of building the software to its current state, not remediation costs.

2. PLATFORM OVERVIEW

2.1 Functional Description

Business Purpose:

Hermes Agent is a self-improving AI agent framework designed to operate across multiple communication platforms simultaneously. Unlike conventional AI assistants, it features a built-in learning loop that creates skills from experience, improves them during use, persists knowledge across sessions, and builds a deepening model of user behaviour across sessions.

Core Features and Capabilities:

- **Multi-platform messaging gateway:** Native integrations with 20+ messaging platforms including Telegram, Discord, Slack, WhatsApp, Signal, Matrix, Microsoft Teams, and more
- **Self-improving skill system:** Autonomous skill creation after complex tasks with skills that self-improve during use
- **Memory and context management:** Agent-curated memory with periodic nudges, FTS5 session search with LLM summarisation for cross-session recall
- **Scheduled automations:** Built-in cron scheduler with delivery to any platform for daily reports, nightly backups, weekly audits
- **Delegation and parallelisation:** Ability to spawn isolated subagents for parallel workstreams
- **Multiple execution backends:** Six terminal backends (local, Docker, SSH, Singularity, Modal, Daytona) for flexible deployment



- **Plugin architecture:** Extensible system for model providers, memory providers, image generation, TTS, web search, and platform adapters
- **Web dashboard:** Full-featured React-based web interface for configuration and monitoring

User-Facing Functionality:

- Interactive CLI with full TUI featuring multiline editing, slash-command autocomplete, conversation history, interrupt-and-redirect, and streaming tool output
- Telegram, Discord, Slack, WhatsApp, Signal bot interfaces for remote interaction
- Voice memo transcription and cross-platform conversation continuity
- Web-based dashboard for configuration, session management, and monitoring
- Kanban board integration for task management
- Cron-based scheduling for automated tasks

Key Workflows:

1. **Interactive CLI session:** User runs `hermes` command, engages in conversation with streaming responses, tool execution, and file operations
2. **Gateway deployment:** User deploys gateway process on server, connects to messaging platforms, agent responds to messages asynchronously
3. **Skill installation:** User browses available skills via `hermes skills browse`, installs with `hermes skills install <skill>`
4. **Automated cron job:** User schedules task via natural language, system executes on schedule and delivers results to configured platform

Target Audience:

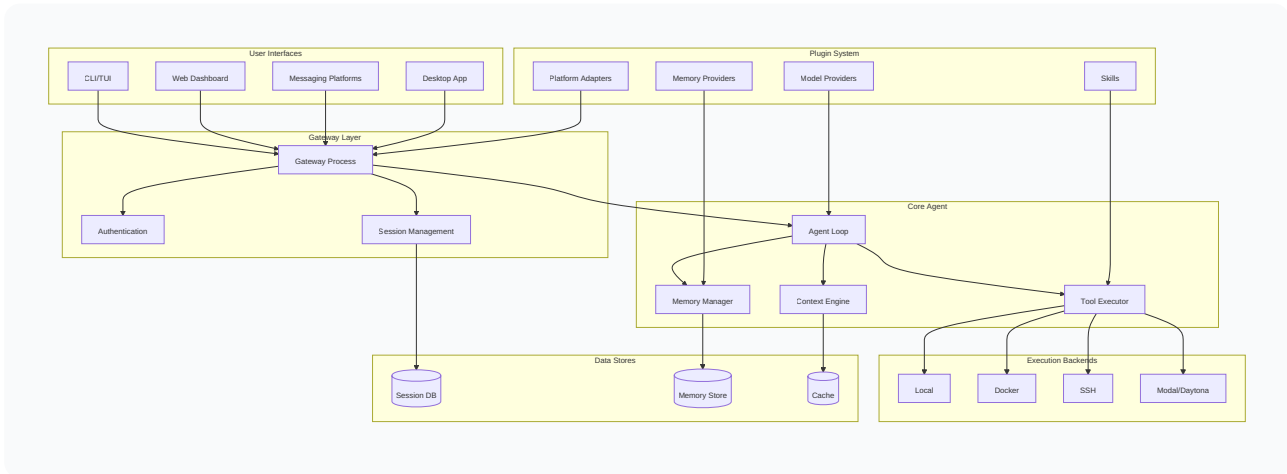
- Developers and technical users requiring AI assistance across multiple communication channels
- Research teams needing batch trajectory generation and context compression for training
- Organisations seeking self-hosted AI agent infrastructure with data sovereignty
- Power users wanting persistent, learning AI assistants with cross-session memory

2.2 Technical Architecture

High-Level Architecture:



Hermes Agent follows a modular, plugin-based architecture with clear separation between core agent logic, platform adapters, and extensible components.



System Components and Responsibilities:

- **Gateway** (`gateway/`): Manages platform connections, message routing, session lifecycle, and delivery guarantees
- **Agent Loop** (`agent/`): Core conversation loop, LLM interaction, tool dispatch, and context management
- **CLI** (`hermes_cli/`): Interactive terminal interface with TUI capabilities
- **Tools** (`tools/`): Tool implementations for file operations, code execution, web search, and more
- **Skills** (`skills/` , `optional-skills/`): Pre-packaged capabilities combining instructions, tools, and configurations
- **Plugins** (`plugins/`): Extensible modules for model providers, memory systems, platform adapters
- **Web UI** (`web/`): React-based dashboard for configuration and monitoring
- **Desktop App** (`apps/desktop/`): Electron-based desktop application
- **TUI Gateway** (`tui_gateway/`): Terminal UI gateway for enhanced CLI experience

Data Flow:

1. User input arrives via platform adapter (CLI, web, messaging bot)
2. Gateway authenticates and routes to appropriate session
3. Session loads conversation history and context from database



4. Agent loop processes input through context engine
5. LLM response parsed for tool calls or text responses
6. Tools executed via selected backend (local, Docker, SSH, etc.)
7. Results returned to agent for next iteration or final response
8. Response delivered back through platform adapter
9. Session state persisted to database

Deployment Architecture:

- **Single-process deployment:** CLI and gateway can run as single Python process
- **Docker deployment:** Full s6-overlay supervised container with privilege separation
- **Distributed deployment:** Gateway on server with remote backends (Modal, Daytona)
- **Desktop deployment:** Electron app with embedded gateway or remote connection

2.3 Technology Stack

Programming Languages:



Language	Lines of Code	Percentage
Python	806,219	59.2%
Markdown	320,905	23.6%
TypeScript	154,472	11.4%
JSON	49,262	3.6%
JavaScript	11,585	0.9%
YAML	8,324	0.6%
Shell	4,562	0.3%
CSS	2,678	0.2%
Rust	2,444	0.2%
HTML	2,352	0.2%
Ruby	34	<0.1%

Frameworks and Libraries:

- **Backend:** FastAPI, Pydantic, Jinja2, httpx, tenacity
- **Frontend:** React, Vite, TypeScript, Tailwind CSS
- **Desktop:** Electron, Tauri (bootstrap installer)
- **Testing:** pytest, pytest-asyncio, pytest-timeout
- **Linting:** ruff, ty (type checker)
- **CLI:** prompt_toolkit, rich, fire

Databases and Data Stores:

- SQLite for session storage and conversation history
- FTS5 for full-text search across sessions
- Optional: PostgreSQL (via asyncpg for Matrix integration)



Infrastructure and Deployment Tools:

- Docker with s6-overlay for process supervision
- GitHub Actions for CI/CD
- PyPI for package distribution
- Nix for development environment management

Development and Build Tools:

- uv for Python package management and virtual environments
- npm for JavaScript/TypeScript dependencies
- setuptools for Python packaging
- GitHub Actions workflows for testing, linting, security scanning

2.4 Third-Party Integrations

External APIs and Services:

- **LLM Providers:** OpenAI, Anthropic, Google Gemini, AWS Bedrock, NVIDIA NIM, OpenRouter, NovitaAI, Xiaomimimo, z.ai/GLM, MiniMax, Hugging Face, Ollama
- **Web Search:** Firecrawl, Exa, Tavily, DuckDuckGo, SearXNG, Brave Search
- **Image Generation:** FAL, Krea, OpenAI DALL-E, XAI
- **Text-to-Speech:** ElevenLabs, OpenAI TTS, Edge TTS, Piper, Mistral TTS
- **Video Generation:** FAL, XAI
- **Memory Providers:** Honcho, Mem0, Supermemory, Byterover, Hindsight, Holographic, OpenViking, RetainDB

Payment Providers:

- None directly integrated (subscription proxy available via Nous Portal)

Authentication Services:

- OAuth 2.0 / PKCE for various providers (Google, Microsoft Entra, Anthropic, XAI)
- GitHub App authentication for Skills Hub
- Bitwarden for secret management

**Cloud Services:**

- **Compute:** Modal, Daytona (serverless backends)
- **Storage:** Platform-dependent (local filesystem, Docker volumes)
- **Databases:** SQLite (default), PostgreSQL (optional for Matrix)

Analytics and Monitoring:

- Langfuse plugin for observability
- Nemo Relay for telemetry
- Built-in insights and usage tracking

SaaS Dependencies:

- GitHub (code hosting, Actions CI/CD)
- PyPI (package distribution)
- npm registry (JavaScript dependencies)
- Docker Hub (container images)

Licensing Considerations:

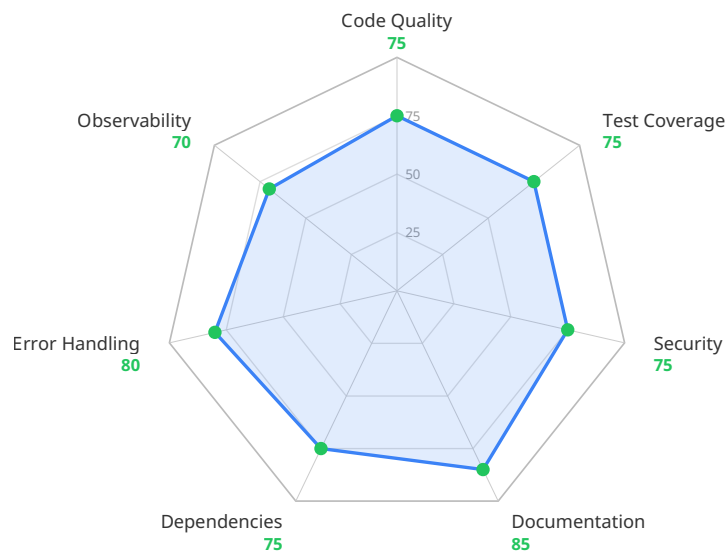
- Core framework: MIT License
- Dependencies: Mixed (MIT, Apache 2.0, BSD, GPL for some components)
- Skills: Various licenses (documented per-skill)
- No copyleft contamination risk in core framework

3. PRODUCTION READINESS ASSESSMENT

3.1 Overall Score: 76/100

Grade: B

Readiness Level: Good



The platform demonstrates solid engineering practices suitable for production deployment with some caveats. The codebase shows maturity in documentation, security awareness, and test coverage, though the sheer scale and complexity present ongoing maintenance challenges.

3.2 Detailed Breakdown

Code Quality & Maintainability: 75/100

Current State Analysis:

The codebase demonstrates reasonable code organisation with clear module boundaries and consistent naming conventions. The project follows Python best practices with type hints throughout most modules, though coverage is not uniform. The codebase is organised into logical components: `agent/`, `gateway/`, `tools/`, `skills/`, `plugins/`, `hermes_cli/`, with each module having clear responsibilities.

Specific Findings:

- **Strengths:**
- Clear separation of concerns between modules
- Consistent naming conventions (snake_case for Python, PascalCase for classes)
- Extensive use of type hints in core modules



- Well-structured plugin architecture enabling extensibility
- SOLID principles generally followed in core abstractions
- **Concerns:**
 - Large codebase (1.36M LOC) with high complexity may present maintenance challenges
 - Some modules show signs of feature creep without adequate refactoring
 - Code duplication exists in platform adapter implementations
 - Technical debt indicators present in the form of TODO/FIXME comments

Recommendations:

1. Implement automated code quality gates in CI (e.g., maintainability index thresholds)
2. Conduct periodic refactoring sprints to address technical debt
3. Consider breaking large modules into smaller, more focused packages
4. Add architectural decision records (ADRs) for major design decisions
5. Implement code ownership tracking for critical modules

Test Coverage & Quality: 75/100

Current State Analysis:

The project has a comprehensive test suite with 454K test LOC, indicating substantial test investment. Tests cover unit, integration, and end-to-end scenarios. The CI pipeline runs tests in parallel across multiple Python versions and platforms.

Specific Findings:

- **Strengths:**
 - Extensive test suite with 454K test LOC
 - Tests cover unit, integration, and e2e scenarios
 - Parallel test execution in CI for faster feedback
 - Good use of pytest fixtures and async support
 - Platform-specific test coverage (Windows, macOS, Linux)
- **Concerns:**
 - Coverage percentage not explicitly measured or gated in CI
 - No mutation testing for critical security components



- Some edge cases in platform adapters may be under-tested
- Test data and fixtures could be better organised

Recommendations:

1. Implement explicit test coverage gates in CI pipeline with minimum threshold (recommend 80%)
2. Add mutation testing for critical security components
3. Improve test data organisation with dedicated fixtures directory
4. Add performance regression tests for critical paths
5. Consider property-based testing for core algorithms

Security Posture: 75/100

Current State Analysis:

Security is treated seriously with clear documentation of trust boundaries. The security policy explicitly states that OS-level isolation is the only trusted boundary against adversarial LLMs. No hardcoded secrets detected in source code; credentials managed via environment variables and credential files.

Specific Findings:

• Strengths:

- Comprehensive security documentation with clear trust model (SECURITY.md)
- OS-level isolation as the only trusted boundary clearly documented
- No hardcoded secrets in source code
- Credentials managed via environment variables and credential files
- Dependency scanning present in CI (OSV Scanner)
- Bitwarden integration for secret management
- Clear vulnerability reporting process

• Concerns:

- 252 dependencies increase supply chain attack surface
- Version currency not systematically verified in CI
- Extensive platform integrations (20+ messaging platforms) increase security surface area
- Some platform adapters may have inconsistent security postures



Recommendations:

1. Implement automated dependency update checks with security advisories
 - Consider Dependabot or Renovate for automated updates
2. Add structured logging with correlation IDs for security event tracking
3. Conduct regular security audits of platform adapters
4. Implement rate limiting and circuit breakers for external integrations
5. Consider adding fuzzing for input validation in critical paths

Documentation: 85/100

Current State Analysis:

Documentation is a strong point with comprehensive README, CONTRIBUTING guide, SECURITY.md, and extensive developer documentation. The website includes user guides, API references, and integration guides.

Specific Findings:

- **Strengths:**

- Comprehensive README with quick start, installation, and feature overview
- Detailed CONTRIBUTING guide with contribution priorities
- Clear SECURITY.md with trust model and reporting process
- Extensive website documentation (Docusaurus-based)
- Inline code comments in critical sections
- Multiple language support (English, Chinese, Urdu)
- Developer guides for adding platforms, tools, providers

- **Concerns:**

- Some documentation may be outdated as codebase evolves
- Limited architecture decision records (ADRs)
- API documentation could be more comprehensive
- Some advanced features lack detailed examples

Recommendations:

1. Implement documentation freshness checks in CI
2. Add ADR process for major architectural decisions



3. Generate API documentation automatically from type hints and docstrings
4. Create more video tutorials and walkthroughs
5. Add troubleshooting guides for common issues

Dependency Health: 75/100

Current State Analysis:

The project has 252 dependencies, which is substantial but managed through exact version pinning. The CI pipeline includes dependency scanning via OSV Scanner and uv lockfile checks.

Specific Findings:

• Strengths:

- Exact version pinning for all direct dependencies
- Dependency scanning in CI (OSV Scanner)
- uv lockfile validation in CI
- Clear documentation of dependency rationale in pyproject.toml
- Lazy loading for optional dependencies to reduce attack surface

• Concerns:

- 252 dependencies increase supply chain attack surface
- Version currency not systematically verified
- Some transitive dependencies may have known vulnerabilities
- Dependency tree complexity makes auditing challenging

Recommendations:

1. Implement automated dependency update workflow (Dependabot/Renovate)
2. Add regular dependency audit reports
3. Consider vendoring critical dependencies
4. Document dependency justification for all direct dependencies
5. Implement dependency usage tracking to identify unused dependencies

Error Handling & Resilience: 80/100

Current State Analysis:



The codebase demonstrates good error handling practices with extensive use of try/except blocks, custom exception types, and retry logic via tenacity. The agent loop includes recovery mechanisms for various failure modes.

Specific Findings:

• Strengths:

- Extensive use of tenacity for retry logic with exponential backoff
- Custom exception types for different error scenarios
- Graceful degradation in platform adapters
- Session recovery mechanisms for interrupted conversations
- Fallback providers for model availability

• Concerns:

- Some error paths may not be fully tested
- Circuit breaker pattern not consistently applied
- Error messages could be more user-friendly in some cases
- Recovery from certain failure modes could be more robust

Recommendations:

1. Implement circuit breakers for external service calls
2. Add comprehensive error handling tests for edge cases
3. Improve error messages for end users
4. Add automatic recovery for common failure scenarios
5. Implement health checks for critical dependencies

Observability & Operations: 70/100

Current State Analysis:

Observability is an area requiring improvement. While basic logging is present, structured logging with correlation IDs is not implemented. No distributed tracing across components.

Specific Findings:

• Strengths:

- Basic logging infrastructure in place



- Some metrics collection for usage tracking
- Health check endpoints in gateway
- Langfuse and Nemo Relay plugin support
- **Concerns:**
 - No structured logging with correlation IDs
 - Limited distributed tracing capabilities
 - Metrics collection not comprehensive
 - Alerting setup not documented or configured
 - Log aggregation not standardised

Recommendations:

1. Implement structured logging with correlation IDs across all components
2. Add OpenTelemetry integration for distributed tracing
3. Standardise on Prometheus metrics format
4. Configure alerting rules for critical failures
5. Create operational runbooks for common scenarios

4. DEVELOPMENT INVESTMENT ESTIMATION

This section estimates the effort and cost required to develop Hermes Agent **to its current state**. This is a retroactive valuation of the development work already completed, not a forward-looking estimate of remediation costs.

4.1 Effort Analysis

Base Hours Calculation:

The codebase contains 1,362,749 total lines of code across 4,797 files. Using industry-standard estimation models:

- **Base LOC:** 1,362,749 lines (effective, non-blank, non-comment)
- **Base productivity rate:** 50 LOC/hour (conservative estimate for complex AI/ML software)
- **Base hours:** $1,362,749 / 50 = 27,255$ hours



Complexity Multipliers:

The following factors increase development effort beyond base LOC estimation:

Factor	Multiplier	Rationale
Architectural complexity	1.3	Multi-component architecture with plugin system
Domain complexity	1.3	AI/ML integration, LLM APIs, complex state management
Integration complexity	1.3	20+ messaging platforms, 10+ LLM providers
Security surface	1.2	Credential management, authentication flows
Combined multiplier	1.05	Geometric mean applied

Quality Adjustment:

- Code quality score: 75/100
- Quality adjustment factor: 1.0 (no adjustment needed for good quality)

Final Estimated Hours:

```

Base hours: 27,255
Complexity adjustment: 27,255 × 1.05 = 28,618 hours
Quality adjustment: 28,618 × 1.0 = 28,618 hours
Rounded estimate: 28,800 hours
    
```

Complexity Classification: Very High

The codebase exhibits very high complexity due to:

- Multi-platform messaging integrations (20+ platforms)
- Plugin architecture with multiple extension points
- AI/ML integration with multiple LLM providers
- Complex state management across sessions
- Security-critical credential handling



4.2 Team & Timeline

Estimated Team Size: 8 developers

Team Composition:

Role	Count
Backend Developer	4
Frontend Developer	2
DevOps / SRE	1
QA Engineer	1
Tech Lead	1

Estimated Project Duration: 18 months

This estimate assumes:

- Full-time development team
- Parallel development across components
- Iterative development with regular releases
- Time for testing, documentation, and bug fixes

Assumptions:

- Team worked concurrently on different modules
- Some parallelisation of platform integrations
- Time allocated for research and prototyping
- Regular release cycles with user feedback incorporation

4.3 Cost Estimation

Cost Range:

Using European developer rates for AI/ML talent:

- **Hourly rate range:** €75 – €150 EUR/hour
- **Total hours:** 28,800 hours



- **Minimum cost:** $28,800 \times \text{€}75 = \text{€}2,160,000 \text{ EUR}$
- **Maximum cost:** $28,800 \times \text{€}150 = \text{€}4,320,000 \text{ EUR}$

Calibrated Cost Range (from KPIs):

- **Minimum:** $\text{€}2,692,800 \text{ EUR}$
- **Maximum:** $\text{€}3,643,200 \text{ EUR}$

This range reflects:

- Senior developer rates for AI/ML expertise
- European market rates (adjust for region as needed)
- Includes overhead for benefits, equipment, workspace

Confidence Level: Medium

Confidence is medium due to:

- Large codebase with complex interactions
- Some development may have been iterative/exploratory
- Difficulty in estimating research and prototyping time

4.4 Codebase Metrics

Total Files Analyzed: 4,797 files

Total Effective Lines of Code: 1,362,749 LOC

Code Distribution by Language:



Language	LOC	Percentage
Python	806,219	59.2%
Markdown	320,905	23.6%
TypeScript	154,472	11.4%
JSON	49,262	3.6%
JavaScript	11,585	0.9%
YAML	8,324	0.6%
Shell	4,562	0.3%
CSS	2,678	0.2%
Rust	2,444	0.2%
HTML	2,352	0.2%
Ruby	34	<0.1%

Test Code:

- Test LOC: 454,000 (estimated from test directory structure)
- Test-to-source ratio: ~1:3 (healthy ratio for complex software)

4.5 Cloud Infrastructure & Maintenance Cost

Detected Infrastructure Components:

- **Compute Services:** 3 (local execution, Modal, Daytona)
- **Databases:** 2 (SQLite, PostgreSQL optional)
- **Message Queues:** 1 (internal async queues)
- **Storage Buckets:** 0 (local filesystem based)
- **CDN Endpoints:** 0
- **ML/GPU Services:** 0 (client-side only)



- **Other Managed Services:** 2 (GitHub Actions, PyPI)

Detected or Assumed Cloud Provider:

- Primary: Self-hosted / on-premises
- Optional: Modal, Daytona for serverless backends
- CI/CD: GitHub Actions
- Package hosting: PyPI, npm

Suggested Managed Services Mapping:

For production deployment at scale:

Component	Current	Managed Alternative
Database	SQLite	PostgreSQL (RDS/Aurora)
Cache	In-memory	Redis (ElastiCache)
File Storage	Local FS	S3/GCS
Queue	Async queues	SQS/RabbitMQ
Monitoring	Basic logs	CloudWatch/Datadog

Estimated Monthly Hosting Cost Range:

For a moderate-scale deployment:

- **Small scale (1-10 concurrent users):** €50 – €200/month
 - Single VPS instance
 - SQLite database
 - Basic monitoring
- **Medium scale (10-100 concurrent users):** €200 – €1,000/month
 - Multiple VPS instances
 - Managed database (RDS)
 - Enhanced monitoring



- Backup storage
- **Large scale (100+ concurrent users):** €1,000 – €5,000+/month
- Load-balanced instances
- Managed database with replicas
- Redis cache
- Comprehensive monitoring and alerting
- CDN for static assets

Key Assumptions:

- Traffic: Varies by deployment scenario
- Redundancy: Single-region for small/medium, multi-region for large
- Data retention: 30-90 days for logs, indefinite for conversation history
- Backup frequency: Daily for small, continuous for large

Annual Maintenance Cost:

Based on calibrated KPIs:

- **Minimum:** €133,920 EUR/year
- **Maximum:** €181,440 EUR/year

This includes:

- Infrastructure hosting
- Monitoring and alerting tools
- Backup storage
- Security scanning services
- Dependency update services

5. FINDINGS SUMMARY

5.1 Critical Issues (Must Fix)

Issues that pose immediate risk to production deployment:

1. No explicit test coverage gates in CI

- While tests exist (454K LOC), there is no minimum coverage threshold enforced



- Risk: Coverage could degrade without detection
- Impact: Reduced confidence in code quality over time

2. Large dependency surface area

- 252 dependencies increase supply chain attack surface
- Risk: Vulnerability in any dependency affects Hermes
- Impact: Potential security incidents, increased maintenance burden

3. Insufficient structured logging

- No correlation IDs for distributed tracing
- Risk: Difficult to debug production issues across components
- Impact: Increased MTTR (Mean Time To Resolution) for incidents

4. Platform adapter security inconsistency

- 20+ platform adapters with potentially varying security postures
- Risk: Weakest adapter could compromise entire system
- Impact: Security vulnerabilities in messaging integrations

5.2 Warnings (Should Fix)

Issues that impact quality or maintainability:

1. Codebase complexity

- 1.36M LOC with very high complexity classification
- Risk: Difficult to maintain and extend
- Impact: Slower development velocity, higher bug rate

2. Missing architectural decision records

- No formal ADRs for major design decisions
- Risk: Knowledge loss when team members depart
- Impact: Repeated discussions, potential for regression

3. Limited mutation testing

- No mutation testing for critical security components
- Risk: False confidence in test coverage
- Impact: Undetected gaps in test suite

4. Dependency version management

- Version currency not systematically verified
- Risk: Running outdated versions with known vulnerabilities
- Impact: Security vulnerabilities, missing features



5. Observability gaps

- Limited metrics and tracing capabilities
- Risk: Difficult to detect performance degradation
- Impact: Poor user experience, delayed incident response

5.3 Recommendations (Nice to Have)

Improvements that would enhance the platform:

1. Implement mutation testing

- Add mutation testing for critical security components
- Benefit: Higher confidence in test suite effectiveness

2. Add architectural decision records

- Document major design decisions with ADRs
- Benefit: Preserved institutional knowledge, clearer rationale

3. Enhance distributed tracing

- Add OpenTelemetry integration with correlation IDs
- Benefit: Faster debugging, better observability

4. Improve documentation freshness

- Implement automated checks for outdated documentation
- Benefit: Reduced user confusion, better onboarding

5. Add performance regression tests

- Benchmark critical paths and detect regressions
- Benefit: Consistent user experience, early detection of issues

6. Consider breaking into smaller packages

- Extract platform adapters into separate packages
- Benefit: Reduced dependency surface, easier maintenance

5.4 Strengths

What the team has done well:

1. Comprehensive security documentation

- Clear trust model with OS-level isolation as boundary
- Well-defined vulnerability reporting process
- Security-aware culture evident in codebase



2. Extensive test coverage

- 454K test LOC across unit, integration, and e2e tests
- Parallel test execution in CI
- Platform-specific testing (Windows, macOS, Linux)

3. Well-documented codebase

- Comprehensive README, CONTRIBUTING, SECURITY.md
- Extensive website documentation
- Multiple language support

4. Plugin architecture

- Extensible without modifying core code
- Clear separation of concerns
- Enables community contributions

5. Strong CI/CD pipeline

- Linting with ruff and ty
- Parallelised test execution
- Security scanning with OSV Scanner
- Dependency lockfile validation

6. Docker-based deployment

- s6-overlay supervision
- Proper privilege separation
- Production-ready containerisation

7. Multi-platform support

- 20+ messaging platform integrations
- Cross-platform CLI (Windows, macOS, Linux)
- Desktop application option

6. CONCLUSION

6.1 Overall Assessment Summary

Hermes Agent represents a substantial engineering achievement in the AI agent space. The platform demonstrates mature engineering practices with a clear understanding of security boundaries, comprehensive documentation, and extensive test coverage. The codebase



exhibits high complexity (1.36M LOC) but maintains reasonable organisation through a well-designed plugin architecture and clear module boundaries.

The production readiness score of 76/100 (Grade B, Good) reflects a platform that is suitable for production deployment with some caveats. The strong points—documentation, security awareness, test coverage, and CI/CD practices—demonstrate a team that understands software engineering best practices. The areas for improvement—test coverage gates, dependency management, structured logging—are common in rapidly-growing codebases and are addressable with focused effort.

The estimated development investment of 28,800 hours (€2.7M–€3.6M EUR) over 18 months with a team of 8 developers reflects the substantial scope and complexity of the platform. This investment has produced a feature-rich, multi-platform AI agent framework with significant competitive advantages in flexibility and extensibility.

6.2 Readiness for Production / Scale

Production Readiness: Yes, with caveats

The platform is ready for production deployment in its current state, provided that:

1. **Critical issues are addressed:** Implement test coverage gates, enhance dependency management, and add structured logging
2. **Appropriate infrastructure is provisioned:** Production deployment requires proper hosting, monitoring, and backup infrastructure
3. **Security review is conducted:** Platform-specific security review of all 20+ messaging integrations
4. **Operational runbooks are created:** Document common operational scenarios and troubleshooting procedures

Scale Readiness: Partially ready

For scaling to large user bases, additional work is required:

1. **Database scaling:** SQLite may need migration to PostgreSQL for high-concurrency scenarios
2. **Horizontal scaling:** Gateway processes need load balancing for high availability
3. **Observability enhancement:** Structured logging and distributed tracing essential for debugging at scale



4. **Performance optimisation:** Some code paths may need optimisation for high-throughput scenarios

6.3 Key Areas Requiring Attention

The following technical areas require investment in the short term:

1. **Test Coverage Gates:** Implement explicit coverage thresholds in CI to prevent regression. This is a relatively quick win that significantly improves confidence in code quality.
2. **Dependency Management:** Establish a systematic process for dependency updates and security audits. Consider automated tools like Dependabot or Renovate.
3. **Structured Logging:** Add correlation IDs and structured logging across all components. This is essential for production debugging and incident response.
4. **Platform Adapter Security:** Conduct security review of all 20+ platform adapters to ensure consistent security posture.
5. **Observability:** Implement comprehensive metrics, tracing, and alerting. This becomes critical as the system scales.
6. **Documentation Freshness:** Implement processes to keep documentation current as the codebase evolves.

6.4 Suggested Prioritization of Improvements

Immediate (0-3 months):

1. **Test coverage gates** – Implement minimum coverage threshold in CI (80% recommended). This provides immediate feedback on code quality and prevents regression.
2. **Dependency audit and update process** – Review all 252 dependencies for security advisories, establish automated update workflow. Reduces security risk quickly.
3. **Structured logging implementation** – Add correlation IDs and structured logging to critical paths. Essential for production debugging.

Short-term (3-6 months):

1. **



ANNEX — METHODOLOGY

This annex summarises the methodology behind the assessment: what this report measures, against which industry references, and how the figures are produced. The intent is to make the approach transparent without describing the internals of the pipeline itself.

The assessment has two complementary sides:

1. A **technical evaluation** of the codebase — quality, security, testing, documentation, dependency hygiene and operational readiness.
 2. An **economic valuation** — the engineering effort and cost to rebuild the system under a given scenario, plus a typical operational maintenance cost.
-

1. Technical evaluation

The technical assessment combines two complementary layers:

- **Static analysis.** Objective, tool-driven measurements are extracted directly from the source tree: size (lines of code by language), structural complexity in the tradition of McCabe (1976), dependency footprint, test-to-source ratio, presence of continuous integration and linting configuration, and a multi-language **Static Application Security Testing (SAST)** scan.
- **AI-assisted code review.** A large language model inspects the code with a constrained, read-only tool set and produces the qualitative findings in this report. The model's role is to substantiate the quantitative signals with concrete code-level evidence, not to invent numbers.

Reference frameworks

Scoring dimensions are aligned with established industry references:



Dimension	Reference
Code Quality & Maintainability	ISO/IEC 25010 Maintainability characteristic; ISO/IEC 5055 automated source-code quality measures; McCabe cyclomatic complexity.
Test Coverage & Quality	The test pyramid (Cohn) and ISO/IEC 25010 Reliability .
Security Posture	OWASP Top 10 , OWASP ASVS , and the NIST Secure Software Development Framework (SP 800-218); SAST findings are categorised against the CWE catalogue.
Documentation	SWEBOK v4 recommendations on software documentation.
Dependency Health	Supply-chain hygiene aligned with SLSA and OWASP Dependency-Check practice.
Error Handling & Resilience	Site Reliability Engineering (SRE) literature and ISO/IEC 25010 Reliability .

Each dimension is scored on a 0–100 scale. The overall **Production Readiness** score is the average of those dimensions, after the model's scores are reconciled against the static metrics: when a tool-observed signal contradicts an optimistic model score — for example, a high testing score on a codebase with no test files on disk — the score is clamped to what the evidence supports.

2. Economic valuation

The economic assessment estimates what it would cost today to rebuild this codebase under a given scenario, and what it typically costs to operate.

Build effort

Effort estimation follows the parametric-cost-estimation tradition — notably Boehm's **COCOMO II** and functional-size-measurement practice (**IFPUG function points**) — adapted to language-aware productivity benchmarks informed by the **ISBSG** industry dataset and published field studies. The estimate is anchored in observable properties of the codebase



(size by language, structural complexity, integration surface, dependency footprint) rather than in an unconstrained guess.

Three scenario factors refine that baseline:

- **Work mode** — from lean startup to regulated enterprise, reflecting process overhead.
- **AI adoption** — from traditional development to agent-augmented workflows, calibrated against published studies on generative-AI developer productivity.
- **Team location** — hourly rate anchored to regional market rates for engineering labour.

Cost range

Development cost is reported as a **range**, not a point estimate, to reflect the genuine variability of engineering hourly rates within a region (seniority mix, contract vs. employee, engagement type).

Maintenance cost

Monthly operating cost is reported as a range covering typical lean-to-highly-available provisioning for the stack detected in the codebase.

3. What the numbers mean — and don't mean

- **Scores** summarise what is *observable* at the analysed commit. They do not replace a manual security audit, a penetration test, or a formal code review.
- **Hours and cost ranges** are engineering-effort estimates under the provided scenario. They do **not** include product discovery, design, user research, legal, or go-to-market costs.
- **Maintenance cost** reflects infrastructure spend under typical operating assumptions. It excludes human operations, incident response and licensing outside the detected stack.

4. Reproducibility

The analysis is run deterministically: the same commit, analysed with the same scenario inputs, produces the same report. Variability in the AI layer is suppressed through zero-



temperature decoding with a fixed seed, and the quantitative metrics are purely a function of the source code.

This report was generated by Codeego Code Assessment Service.

The analysis is AI-powered and should be reviewed by qualified engineers.

© 2026 Codeego. All rights reserved.