

# Agent Governance

Observable Private Language Protocol

Version 1.0.0

MIT License

A governance framework that allows AI agents to use optimized/private communication protocols **only if** they continuously produce human-readable English reports.

February 1, 2026

## Contents

<b>1</b>	<b>Overview</b>	<b>2</b>
<b>2</b>	<b>Quick Start</b>	<b>2</b>
2.1	For Python Users (Agent SDK)	2
2.2	For Rust Users (Policy Gateway)	2
2.3	Full Stack (Both Components)	2
<b>3</b>	<b>Architecture</b>	<b>2</b>
<b>4</b>	<b>Policy Rules</b>	<b>3</b>
4.1	Registration Rule	3
4.2	Reporting Cadence Rule	3
4.3	Translation Completeness Rule	3
4.4	Gating Rule (Enforced by Gateway)	3
4.5	Progressive Enforcement	4
<b>5</b>	<b>Python SDK Reference</b>	<b>4</b>
5.1	Installation	4
5.2	Basic Usage	4
5.3	Custom Translation	5
<b>6</b>	<b>Rust Gateway Reference</b>	<b>5</b>
6.1	Building	5
6.2	Running	5
6.3	API Endpoints	5
6.3.1	POST /register_protocol_for_agent	5
6.3.2	POST /report	6
6.3.3	POST /send	6
6.3.4	Response Codes	6
<b>7</b>	<b>Configuration</b>	<b>6</b>
7.1	Environment Variables	6
<b>8</b>	<b>Integration with MoltBot/MoltBook</b>	<b>7</b>
8.1	Option 1: Sidecar Deployment	7
8.2	Option 2: Library Integration	7
8.3	Option 3: Network Proxy	7
<b>9</b>	<b>Troubleshooting</b>	<b>7</b>
9.1	“Protocol not registered”	7
9.2	“Report overdue”	8
9.3	“Coverage below minimum”	8
<b>10</b>	<b>Security Considerations</b>	<b>8</b>
<b>11</b>	<b>License</b>	<b>8</b>

## 1 Overview

When AI agents communicate, they may develop or use compressed encodings that are more efficient than natural language. This framework makes such “novel language” permissible under strict observability requirements.

Agents Want	Humans Need	This Framework Provides
Efficient encoding	Interpretability	Mandatory English reports
Fast coordination	Audit trail	Append-only event logging
Private protocols	Kill switch	Gateway enforcement

### Important

**Core Principle:** Novel language is treated as “optimized encoding”—allowed only when humans can continuously verify what’s being communicated.

## 2 Quick Start

### 2.1 For Python Users (Agent SDK)

```
1 pip install -r python/requirements.txt
2 python python/observable_agent.py
```

### 2.2 For Rust Users (Policy Gateway)

```
1 cd rust
2 cargo build --release
3 ./target/release/policy_gateway
```

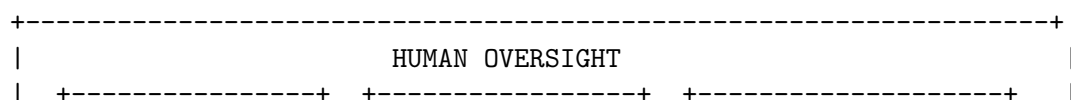
### 2.3 Full Stack (Both Components)

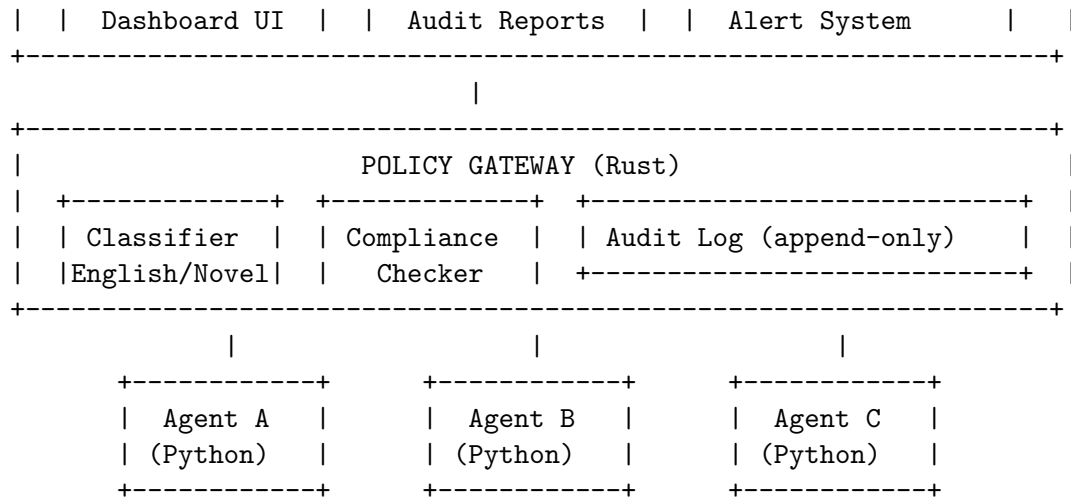
```
1 # Terminal 1: Start the gateway
2 cd rust && cargo run --release
3
4 # Terminal 2: Run an agent
5 cd python && python observable_agent.py
```

## 3 Architecture

The system consists of three main components:

1. **Policy Gateway (Rust)** — Receives all outbound agent messages, classifies them as English vs. novel, checks compliance, and accepts or rejects.
2. **Agent SDK (Python)** — Buffers novel-language messages, generates English reports on cadence, and submits reports to gateway.
3. **Audit Store** — Append-only event log storing raw messages, English reports, and protocol registrations.





## 4 Policy Rules

### 4.1 Registration Rule

Before using any novel language, an agent **must** register a `ProtocolDescriptor`:

```

1 agent.register_protocol(
2     ProtocolDescriptor(
3         name="compressed_coord",           # Unique identifier
4         version="1.0",                     # Version string
5         purpose="Multi-agent coordination", # Why this protocol exists
6         scope="Internal state deltas",      # What it covers
7         risk_tier="medium",                 # low, medium, high
8         translation_method="dictionary",    # How reports are generated
9     )
10 )

```

### 4.2 Reporting Cadence Rule

If an agent sends novel-language messages, it **must** submit an English report:

- Every **60 seconds**, OR
- Every **25 messages**
- Whichever comes first

### 4.3 Translation Completeness Rule

Reports must include:

Field	Requirement
coverage	$\geq 95\%$ of novel messages since last report
self_confidence	Self-rated translation accuracy (0–1)
message_ids	Hashes linking to raw messages
english_summary	Plain English explanation ( $\geq 30$ chars)

### 4.4 Gating Rule (Enforced by Gateway)

The gateway **rejects**:

- Novel-language messages from unregistered protocols
- Messages from agents with overdue reports
- Messages from agents whose last report failed validation

## 4.5 Progressive Enforcement

Violation	Action
1st	Throttle + warning
2nd	Quarantine (read-only mode)
3rd	Disable agent credential

# 5 Python SDK Reference

## 5.1 Installation

```
1 cd python
2 pip install -r requirements.txt
```

## 5.2 Basic Usage

```
1 from observable_agent import ObservableAgent, ProtocolDescriptor,
   GatewayClient
2
3 # Connect to the policy gateway
4 gateway = GatewayClient(
5     base_url="http://localhost:8080",
6     api_key="your-api-key"
7 )
8
9 # Create an observable agent
10 agent = ObservableAgent(agent_id="agent-001", gateway=gateway)
11
12 # Register your protocol BEFORE using novel language
13 agent.register_protocol(
14     ProtocolDescriptor(
15         name="my_protocol",
16         version="1.0",
17         purpose="Task coordination",
18         scope="State synchronization",
19         risk_tier="low",
20         translation_method="lookup_table"
21     )
22 )
23
24 # English messages pass through freely
25 agent.send("agent-002", "Hello! I'll send a compressed update next.")
26
27 # Novel-language messages require protocol registration
28 agent.send("agent-002", "X9|st=17;f=0x3a;ack#42") # Triggers report
   check
```

## 5.3 Custom Translation

Override `_build_english_report()` for your protocol:

```
1 class MyAgent(ObservableAgent):
2     def _build_english_report(self) -> EnglishReport:
3         # Decode your novel messages to English here
4         decoded_messages = [self.decode(msg) for _, _, msg in self._novel_buffer]
5
6         summary = "Coordinated task assignments: " + "; ".join(
7             decoded_messages)
8
9         return EnglishReport(
10             agent_id=self.agent_id,
11             protocol_name=self.protocol.name,
12             protocol_version=self.protocol.version,
13             window_start_ts=self._window_start_ts,
14             window_end_ts=time.time(),
15             message_ids=[mid for _, mid, _ in self._novel_buffer],
16             english_summary=summary,
17             coverage=1.0,
18             self_confidence=0.95,
```

## 6 Rust Gateway Reference

### 6.1 Building

```
1 cd rust
2 cargo build --release
```

### 6.2 Running

```
1 ./target/release/policy_gateway
2 # Listens on http://127.0.0.1:8080
```

### 6.3 API Endpoints

#### 6.3.1 POST /register\_protocol\_for\_agent

Register a protocol for an agent.

```
1 {
2     "agent_id": "agent-001",
3     "protocol": {
4         "name": "compressed_coord",
5         "version": "1.0",
6         "purpose": "Multi-agent coordination",
7         "scope": "Internal state deltas",
8         "risk_tier": "medium",
9         "translation_method": "dictionary"
10    }
11 }
```

### 6.3.2 POST /report

Submit an English translation report.

```

1 {
2   "agent_id": "agent-001",
3   "protocol_name": "compressed_coord",
4   "protocol_version": "1.0",
5   "window_start_ts": 1706745600.0,
6   "window_end_ts": 1706745660.0,
7   "message_ids": ["abc123...", "def456..."],
8   "english_summary": "Exchanged task queue updates...",
9   "coverage": 1.0,
10  "self_confidence": 0.9,
11  "notes": "Auto-generated"
12 }
```

### 6.3.3 POST /send

Send a message (gated by compliance).

```

1 {
2   "from": "agent-001",
3   "to": "agent-002",
4   "content": "X9|st=17;f=0x3a;ack#42",
5   "protocol": {
6     "name": "compressed_coord",
7     "version": "1.0"
8   },
9   "ts": 1706745665.0
10 }
```

### 6.3.4 Response Codes

Code	Meaning
200	Message accepted
400	Report validation failed (coverage, summary length)
403	Protocol not registered
429	Report overdue—submit report to continue

## 7 Configuration

### 7.1 Environment Variables

Variable	Default	Description
REPORT_INTERVAL_SEC	60	Max seconds between reports
REPORT_EVERY_N_MESSAGES	25	Max messages before report required
MIN_COVERAGE	0.95	Minimum coverage fraction
MIN_SUMMARY_LENGTH	30	Minimum English summary characters
RETENTION_DAYS	30	Audit log retention period

## 8 Integration with MoltBot/MoltBook

### 8.1 Option 1: Sidecar Deployment

Deploy the gateway as a sidecar container:

```
1 # docker-compose.yml
2 services:
3   policy-gateway:
4     build: ./rust
5     ports:
6       - "8080:8080"
7     volumes:
8       - ./audit-logs:/var/log/audit
9
10  moltbot:
11    image: moltbot:latest
12    environment:
13      - GATEWAY_URL=http://policy-gateway:8080
14    depends_on:
15      - policy-gateway
```

### 8.2 Option 2: Library Integration

Import the Python SDK directly:

```
1 # In your MoltBot agent code
2 from agent_governance import ObservableAgent, ProtocolDescriptor
3
4 class MyMoltBotAgent(ObservableAgent):
5     def __init__(self, moltbot_config):
6         gateway = GatewayClient(
7             base_url=moltbot_config.gateway_url,
8             api_key=moltbot_config.api_key
9         )
10        super().__init__(
11            agent_id=moltbot_config.agent_id,
12            gateway=gateway
13        )
```

### 8.3 Option 3: Network Proxy

Configure MoltBot to route all agent traffic through the gateway:

```
[Agent A] --> [Policy Gateway :8080] --> [Agent B]
              |
              [Audit Log]
```

## 9 Troubleshooting

### 9.1 “Protocol not registered”

**Cause:** Trying to send novel-language messages before calling `register_protocol()`.

**Fix:** Always register before sending:



```
1 agent.register_protocol(descriptor) # Do this first!  
2 agent.send(to, novel_message)      # Now this works
```

## 9.2 “Report overdue”

**Cause:** More than 60 seconds (or 25 messages) since last report.

**Fix:** The SDK handles this automatically. If you’re calling the gateway directly, submit a report first.

## 9.3 “Coverage below minimum”

**Cause:** Report doesn’t cover enough of the novel messages sent.

**Fix:** Ensure your `_build_english_report()` includes all buffered messages.

# 10 Security Considerations

1. **API Authentication:** In production, require API keys for all gateway endpoints
2. **TLS:** Always use HTTPS between agents and gateway
3. **Log Encryption:** Encrypt raw messages at rest; keep reports in plaintext for auditing
4. **Access Control:** Restrict who can read audit logs and modify protocol registrations
5. **Rate Limiting:** Prevent DoS by limiting registration and message rates per agent

# 11 License

MIT License. See LICENSE file for details.