

Internet Engineering Task Force
Internet-Draft
Intended status: Standards Track
Expires: 1 November 2026

R. Sharif
CyberSecAI Ltd
1 May 2026

ATTP: Agent Trust Transport Protocol

draft-sharif-attp-00

Abstract

This document specifies the Agent Trust Transport Protocol (ATTP), a protocol-agnostic framework for trust scoring, cryptographic identity, action-limit enforcement, compliance gating, and tamper-evident audit for autonomous AI agents.

ATTP separates trust from identity. Identity protocols answer "who is this agent?" ATTP answers "should this agent be allowed to perform this action, at this magnitude, against this counterparty, right now?"

The protocol defines five progressive trust levels (L0-L4), per-agent ECDSA P-256 cryptographic identity, a five-dimension behavioural trust scoring model, action-limit tiers derived from trust scores, real-time compliance gating (sanctions screening, jurisdictional controls), kill switches for instant revocation, anomaly detection for agent behaviour, and a public trust query API.

ATTP is transport-agnostic. This document defines the core framework. Protocol-specific bindings map ATTP trust enforcement to individual transports: MCPs provides the binding for the Model Context Protocol (MCP), with additional bindings defined for REST APIs, Google A2A, gRPC, and GraphQL.

This specification supersedes draft-sharif-agent-payment-trust-00, broadening the scope from payment transactions to any agent action requiring trust-gated authorisation.

Status of This Memo

This Internet-Draft is submitted in full conformance with the provisions of BCP 78 and BCP 79.

Internet-Drafts are working documents of the Internet Engineering Task Force (IETF). Note that other groups may also distribute working documents as Internet-Drafts. The list of current Internet-Drafts is at <https://datatracker.ietf.org/drafts/current/>.

Internet-Drafts are draft documents valid for a maximum of six months and may be updated, replaced, or obsoleted by other documents at any time. It is inappropriate to use Internet-Drafts as reference material or to cite them other than as "work in progress."

This Internet-Draft will expire on 1 November 2026.

Copyright Notice

Copyright (c) 2026 IETF Trust and the persons identified as the document authors. All rights reserved.

This document is subject to BCP 78 and the IETF Trust's Legal Provisions Relating to IETF Documents (<https://trustee.ietf.org/license-info>) in effect on the date of publication of this document. Please review these documents

carefully, as they describe your rights and restrictions with respect to this document.

Table of Contents

1. Introduction
 - 1.1. The Trust Gap
 - 1.2. Relationship to Identity Protocols
 - 1.3. Relationship to Existing Specifications
 - 1.4. Design Goals
2. Terminology
3. Agent Identity
 - 3.1. Key Pair Generation
 - 3.2. Public Key Registration
 - 3.3. Agent Passport
 - 3.4. Key Rotation
4. Challenge-Response Identity Verification
 - 4.1. Challenge Issuance
 - 4.2. Challenge Signing
 - 4.3. Signature Verification
 - 4.4. Failure Handling
5. Trust Scoring Model
 - 5.1. Scoring Dimensions
 - 5.2. Dimension Weights
 - 5.3. Score Computation
 - 5.4. Trust Levels
 - 5.5. Dynamic Adjustment
 - 5.6. Dormancy Decay
 - 5.7. Promotion Rate Limiting
6. Action Limit Enforcement
 - 6.1. Per-Action Limits
 - 6.2. Daily Aggregate Limits
 - 6.3. Principal Aggregate Limits
 - 6.4. Cooling Period
7. Kill Switches
 - 7.1. Per-Agent Kill Switch
 - 7.2. Per-Principal Kill Switch
 - 7.3. Emergency Freeze
8. Compliance Gating
 - 8.1. Sanctions Screening
 - 8.2. Jurisdictional Controls
 - 8.3. Compliance Record
9. Anomaly Detection
 - 9.1. Magnitude Anomaly
 - 9.2. Velocity Anomaly
 - 9.3. Temporal Anomaly
 - 9.4. Counterparty Anomaly
 - 9.5. Probing Anomaly
 - 9.6. Self-Dealing Detection
10. Public Trust Query API
 - 10.1. Request Format
 - 10.2. Response Format
 - 10.3. Batch Queries
 - 10.4. Rate Limiting
 - 10.5. Information Disclosure Controls
11. Action Signing and Receipts
 - 11.1. Action Envelope
 - 11.2. Hash Chain
 - 11.3. Non-Repudiable Receipts
12. Protocol Bindings
 - 12.1. MCP Binding (MCPS)
 - 12.2. REST Binding
 - 12.3. Google A2A Binding
 - 12.4. gRPC Binding
 - 12.5. GraphQL Binding
 - 12.6. Keycloak Integration

13.	Revocation
13.1.	Agent Revocation
13.2.	Emergency Freeze
14.	Security Considerations
14.1.	Key Management
14.2.	Trust Farming Mitigation
14.3.	Sybil Attack Prevention
14.4.	Impersonation Detection
14.5.	Race Condition Handling
14.6.	Replay Protection
14.7.	Kill Switch Consistency
14.8.	Comparison with Identity-Only Protocols
15.	Privacy Considerations
16.	IANA Considerations
17.	References
17.1.	Normative References
17.2.	Informative References
Appendix A.	Regulatory Mapping
A.1.	PSD2 SCA Mapping
A.2.	PCI DSS v4.0.1 Mapping
Appendix B.	Trust Level Reference Implementation
Appendix C.	Comparison with Identity-Only Approaches
	Author's Address

1. Introduction

1.1. The Trust Gap

The emergence of autonomous AI agents that discover resources at runtime, execute multi-step tasks across trust domains, and take consequential actions without human oversight creates a security gap that identity protocols alone cannot close.

Authentication answers a binary question: is this agent who it claims to be? The answer is yes or no. But consequential actions -- financial transactions, data access, infrastructure changes, contract execution -- require graduated assessment:

- o Has this agent earned the right to perform actions of this magnitude?
- o Does this agent's behavioural history justify the level of trust required?
- o Is the counterparty cleared for this interaction (sanctions, jurisdictional compliance)?
- o Can this agent be stopped instantly if something goes wrong?
- o Is there a tamper-evident record of every action for audit?

Identity is necessary but not sufficient. ATTP provides the trust layer that sits above identity and below application logic.

1.2. Relationship to Identity Protocols

ATTP does not replace identity protocols. It consumes identity assertions from any standards-compliant source:

- o OAuth 2.0 / OpenID Connect tokens (RFC 6749, OpenID Core)
- o X.509 certificates (RFC 5280)
- o HTTP Message Signatures (RFC 9421)
- o Self-issued cryptographic identities (JWKS-based)
- o Any OIDC-compliant provider (Keycloak, Azure AD, Okta, etc.)

ATTP requires proof that an agent is who it claims to be. It does not prescribe how that proof is obtained. Once identity is established, ATTP evaluates trust, enforces limits, gates compliance, and records the decision.

1.3. Relationship to Existing Specifications

This specification is part of a three-layer stack:

AEBA (Behavioural Analytics) draft-sharif-aeba Anomaly detection, ML-based monitoring
ATTP (Trust Transport) [this document] draft-sharif-attp Trust scoring, limits, compliance, audit
MCPS (MCP Secure) draft-sharif-mcps-secure-mcp Message signing for MCP transport

MCPS provides message-level cryptographic security for the Model Context Protocol. ATTP provides the trust framework that MCPS (and other transport bindings) enforce. AEBA provides behavioural analytics that feed into ATTP trust scoring.

1.4. Design Goals

1. Transport-agnostic. ATTP defines trust semantics, not wire format. Protocol bindings map ATTP to specific transports.
2. Identity-agnostic. ATTP consumes identity from any source. It does not compete with any identity protocol -- OAuth, OIDC, mTLS, SPIFFE, HTTP Message Signatures, or future standards.
3. Graduated trust. Binary identity is insufficient. Five trust levels (L0-L4) with associated action limits provide proportional authorisation.
4. Fail-closed. If trust cannot be evaluated, the default is deny. No silent degradation.
5. Instant revocation. Kill switches take effect on the next request. No waiting for certificate expiry or CRL refresh.
6. Compliance-native. Sanctions screening and jurisdictional controls are built into the protocol, not bolted on.
7. Self-hostable. The Trust Authority has no external service dependency. Operators run their own instance.
8. Auditable. Every trust decision produces a tamper-evident record with hash-chain integrity.

2. Terminology

The key words "MUST", "MUST NOT", "REQUIRED", "SHALL", "SHALL NOT", "SHOULD", "SHOULD NOT", "RECOMMENDED", "NOT RECOMMENDED", "MAY", and "OPTIONAL" in this document are to be interpreted as described in BCP 14 [RFC2119] [RFC8174] when, and only when, they appear in all capitals, as shown here.

Agent: An autonomous software entity that performs actions on behalf of a principal (human or organisation) via standardised protocols. Actions include but are not limited to: financial transactions, data queries, tool invocations, infrastructure changes, and contract execution.

Principal: The human or organisational entity accountable for an agent's actions.

Trust Authority (TA): The service that maintains agent registrations, computes trust scores, evaluates compliance gates, enforces action limits, and records audit entries. The TA is self-hostable with no external dependency.

Trust Score: A numerical value (0-100) representing the assessed trustworthiness of an agent based on multiple behavioural dimensions.

Trust Level: A discrete classification (L0-L4) derived from the trust score, with associated action limits.

Action: Any operation an agent requests that requires trust evaluation. Actions are protocol-specific: a payment in Stripe, a tool call in MCP, an API request in REST, a task in A2A.

Action Limit: A constraint on the magnitude or frequency of actions an agent may perform at a given trust level.

Kill Switch: A per-agent or per-principal flag that, when activated, causes immediate denial of all actions regardless of trust level.

Compliance Gate: A check performed before action execution that evaluates regulatory constraints (sanctions lists, jurisdictional rules, embargo lists).

Protocol Binding: A specification that maps ATTP trust enforcement to a specific transport protocol (MCP, REST, A2A, gRPC, GraphQL).

Agent Passport: A cryptographically signed identity document containing the agent's public key hash, principal identity, authorised scope, trust level, and issuance metadata.

3. Agent Identity

3.1. Key Pair Generation

Each agent MUST have a unique ECDSA key pair using the P-256 curve [FIPS186-5].

Key pairs MAY be generated by the Trust Authority and provided to the principal (the private key MUST be returned exactly once and MUST NOT be stored by the Trust Authority), or generated by the principal and the public key registered with the Trust Authority (RECOMMENDED for production deployments).

Implementations MUST support ECDSA P-256 (ES256). Implementations MAY additionally support Ed25519 (EdDSA) for environments where P-256 is not available.

3.2. Public Key Registration

The principal MUST register the agent's public key with the

Trust Authority at agent creation time. The Trust Authority stores ONLY the public key. The private key MUST remain under the principal's control at all times.

For production deployments, the private key SHOULD be stored in a Hardware Security Module (HSM), cloud Key Management Service (KMS), or platform-native secure storage.

3.3. Agent Passport

Upon registration, the Trust Authority issues an Agent Passport -- a signed JSON document containing:

- o agentId: Unique identifier for the agent.
- o publicKeyHash: SHA-256 hash of the agent's public key, binding the passport to a specific key pair.
- o principalId: Identifier of the owning principal.
- o scope: Array of authorised action categories.
- o trustLevel: Current trust level (L0-L4).
- o issuedAt: ISO 8601 timestamp of issuance.
- o expiresAt: ISO 8601 timestamp of expiry.
- o issuer: Identifier of the Trust Authority.
- o protocolVersion: ATTP version (this document: "1.0").

The passport is signed by the Trust Authority's key using ECDSA P-256. The signature covers the JSON Canonicalization Scheme (JCS) [RFC8785] serialisation of the passport fields.

Passports MUST expire. Maximum lifetime MUST NOT exceed 365 days. RECOMMENDED lifetime is 90 days for L0-L2 agents and 180 days for L3-L4 agents.

3.4. Key Rotation

Agents MUST support key rotation without identity loss. When rotating keys:

1. The principal generates a new key pair.
2. The principal registers the new public key with the TA.
3. The TA issues a new passport binding the agentId to the new public key.
4. The old key enters a grace period (RECOMMENDED: 24 hours) during which both keys are accepted.
5. After the grace period, the old key is revoked.

Key rotation MUST NOT reset the agent's trust score.

4. Challenge-Response Identity Verification

When a platform needs to verify an agent's identity before granting access, it uses the following challenge-response protocol.

4.1. Challenge Issuance

The platform (or TA on behalf of the platform) generates a challenge:

- o challenge: 32 bytes of cryptographically random data, hex-encoded (64 characters).
- o expiresAt: Timestamp, MUST be no more than 60 seconds from issuance.
- o agentId: The agent for which the challenge is issued.

Each challenge MUST be single-use. The TA MUST reject any

attempt to verify a challenge that has already been used.

4.2. Challenge Signing

The agent signs the challenge using its private key:

```
signature = ECDSA-Sign(SHA-256(challenge), agentPrivateKey)
```

The signature uses IEEE P1363 fixed-length $r||s$ encoding (64 bytes for P-256) per RFC 7518 Section 3.4.

4.3. Signature Verification

The TA verifies the signature using the agent's registered public key. If valid, the TA returns the agent's current trust score, level, and recommendation.

4.4. Failure Handling

If verification fails:

- o The TA MUST log the failure as a potential impersonation attempt.
- o The agent's trust score MUST be penalised (RECOMMENDED: -10 points per failed attempt).
- o The response MUST include an error code: IMPERSONATION, CHALLENGE_EXPIRED, CHALLENGE_REPLAYED, or AGENT_MISMATCH.

5. Trust Scoring Model

5.1. Scoring Dimensions

The trust score is computed from five dimensions:

1. Code Attestation (CA): Whether the agent's code or configuration has been cryptographically verified against its declared specification.
2. Execution Success Rate (ES): The proportion of the agent's past actions that completed without anomaly, dispute, or reversal.
3. Behavioural Consistency (BC): The statistical consistency of the agent's action patterns relative to its established baseline.
4. Operational Tenure (OT): The duration the agent has been registered and actively operating.
5. Anomaly History (AH): The inverse of the count and severity of detected anomalies.

5.2. Dimension Weights

Each dimension contributes 20% to the overall score:

$$W_{CA} = W_{ES} = W_{BC} = W_{OT} = W_{AH} = 0.20$$

Implementations MAY adjust weights provided the sum equals 1.0 and no single dimension exceeds 0.40.

5.3. Score Computation

$$\text{raw} = (CA * W_{CA}) + (ES * W_{ES}) + (BC * W_{BC}) + (OT * W_{OT}) + (AH * W_{AH})$$

Each dimension value is in the range [0, 100].

```
score = clamp(raw + bonus + dormancy, 0, 100)
```

5.4. Trust Levels

Level	Score	Label	Payment Enabled
L0	0-19	No Access	No
L1	20-39	Restricted	Micro only
L2	40-59	Standard	Yes
L3	60-79	Elevated	Yes (monitored)
L4	80-100	Full Access	Yes (audited)

No trust level SHALL have unlimited action authority. L4 represents the maximum tier with mandatory enhanced auditing.

5.5. Dynamic Adjustment

The trust score adjusts dynamically based on observed behaviour:

- o Successful action: +0.5 trust bonus.
- o Blocked action (over limit): -2 trust bonus.
- o Anomaly detected: -5 per anomaly.
- o Critical anomaly (3+ simultaneous): -20 trust bonus.
- o Failed identity verification: -10 trust bonus.
- o Probing detected: -15 trust bonus.

Self-dealing actions (agent-to-agent within the same principal) MUST earn zero trust bonus.

5.6. Dormancy Decay

Agents inactive for extended periods MUST have their trust score reduced:

- o 30 days inactive: -10 penalty.
- o 60 days inactive: -20 penalty.
- o 90+ days inactive: -30 penalty.

The dormancy penalty is removed when the agent resumes activity.

5.7. Promotion Rate Limiting

Trust level promotions MUST be rate-limited to prevent trust farming attacks:

- o L0 to L1: Minimum 24 hours at L0 with at least 5 successful actions.
- o L1 to L2: Minimum 7 days at L1 with at least 20 successful actions.
- o L2 to L3: Minimum 30 days at L2 with at least 100 successful actions and zero critical anomalies.
- o L3 to L4: Minimum 90 days at L3 with at least 500 successful actions, zero anomalies, and manual principal attestation.

Demotions are instant. A single critical anomaly at L4 drops the agent to L2 immediately.

6. Action Limit Enforcement

6.1. Per-Action Limits

Each action **MUST** be checked against the per-action limit of the agent's current trust level. Actions exceeding the limit **MUST** be rejected with error code ATTP-ACTION-LIMIT.

Action limits are protocol-binding-specific. For financial transactions:

Level	Per-TX Limit	Daily Limit
L0	\$0	\$0
L1	\$10	\$50
L2	\$100	\$500
L3	\$1,000	\$5,000
L4	\$50,000	\$200,000

For non-financial actions, protocol bindings define equivalent limits (e.g., API calls per minute, data volume per request).

6.2. Daily Aggregate Limits

The TA **MUST** track daily aggregate action magnitude per agent within a rolling 24-hour window.

6.3. Principal Aggregate Limits

To prevent Sybil attacks (creating many agents to circumvent per-agent limits), the TA **MUST** enforce aggregate limits per principal across all their agents.

6.4. Cooling Period

After a trust level promotion, the agent **MUST** operate under the previous level's limits for a cooling period (RECOMMENDED: 24 hours) before the new limits take effect.

7. Kill Switches

7.1. Per-Agent Kill Switch

Each agent **MUST** have an independently controllable kill switch. When activated:

- o ALL actions by the agent are denied immediately.
- o The denial takes effect on the next request (no grace period).
- o The agent's trust score is frozen (not reset).
- o The kill switch state is recorded in the audit trail.

Reactivation requires explicit principal action and **MUST NOT** be automated.

7.2. Per-Principal Kill Switch

Each principal **MUST** have a kill switch that simultaneously disables ALL agents under their control.

7.3. Emergency Freeze

The TA **MUST** support an emergency freeze that disables all agents across all principals. This is intended for infrastructure-level incidents. Emergency freeze **MUST** require multi-party authorisation (RECOMMENDED: two

independent operators).

8. Compliance Gating

8.1. Sanctions Screening

Before executing any financial action, the TA MUST screen the counterparty against applicable sanctions lists:

- o OFAC SDN (US)
- o EU Consolidated Sanctions
- o UK HMT Sanctions
- o UN Security Council Sanctions

Screening MUST use fuzzy matching with a configurable threshold (RECOMMENDED: 70% match score). Matches above the threshold MUST block the action with error code ATTP-SANCTIONS-MATCH.

Sanctions lists MUST be updated at least daily.

8.2. Jurisdictional Controls

Actions MAY be restricted based on the jurisdiction of the agent, the principal, or the counterparty. The TA MUST support configurable jurisdiction rules.

8.3. Compliance Record

Every compliance check MUST produce an audit record containing: the query, the lists checked, the match results (including near-misses), and the decision. Compliance records MUST be retained for the period required by applicable regulation (RECOMMENDED minimum: 5 years).

9. Anomaly Detection

9.1. Magnitude Anomaly

An action whose magnitude exceeds 3 standard deviations from the agent's historical mean MUST be flagged.

9.2. Velocity Anomaly

An action rate exceeding 2x the agent's historical maximum within any 5-minute window MUST be flagged.

9.3. Temporal Anomaly

Actions outside the agent's established operating hours (based on 90-day history) SHOULD be flagged.

9.4. Counterparty Anomaly

Actions directed at counterparties with which the agent has no prior history, when combined with elevated magnitude, MUST be flagged.

9.5. Probing Anomaly

A pattern of incrementally increasing action magnitudes (e.g., \$1, \$5, \$10, \$50, \$100) within a short window indicates probing behaviour and MUST trigger a -15 trust penalty and flag the agent for review.

9.6. Self-Dealing Detection

Actions between agents belonging to the same principal MUST earn zero trust bonus and SHOULD be flagged if they exceed 10% of the agent's total action volume.

10. Public Trust Query API

10.1. Request Format

Any platform MAY query the TA for an agent's trust status:

```
GET /v1/trust/{agentId}
```

The TA MUST support this query without requiring a pre-existing relationship with the querying platform.

10.2. Response Format

```
{
  "agentId": "agent_abc123",
  "trust": {
    "score": 67,
    "level": 3,
    "label": "L3 -- Elevated"
  },
  "recommendation": "ALLOW",
  "limits": {
    "perAction": 100000,
    "daily": 500000
  },
  "identity": {
    "verified": true,
    "principalId": "dev_xyz",
    "jurisdiction": "GB"
  },
  "meta": {
    "protocolVersion": "1.0",
    "queriedAt": "2026-05-01T00:00:00Z",
    "checkedBy": "ATTP Trust Authority"
  }
}
```

10.3. Batch Queries

The TA MUST support batch queries for efficiency:

```
POST /v1/trust/batch
{ "agentIds": ["agent_1", "agent_2", "agent_3"] }
```

Maximum batch size: 100 agents.

10.4. Rate Limiting

The TA MUST rate-limit trust queries:

- o Anonymous queries: 120 per minute per IP.
- o Authenticated queries: 600 per minute per API key.

10.5. Information Disclosure Controls

The public trust API MUST NOT expose:

- o Raw scoring dimension values.

- o Action history or details.
- o Private key material or key fingerprints.
- o Internal anomaly thresholds.

The API provides trust level, recommendation, and limits only.

11. Action Signing and Receipts

11.1. Action Envelope

Every action processed through ATTP produces a signed envelope:

```
{
  "actionId": "act_unique_id",
  "agentId": "agent_abcl23",
  "action": "payment_initiate",
  "magnitude": 5000,
  "counterparty": "recipient_name",
  "trustLevel": 3,
  "complianceResult": "CLEAR",
  "timestamp": "2026-05-01T00:00:00Z",
  "signature": "<ECDSA signature over canonical JSON>"
}
```

The signature covers the JCS [RFC8785] serialisation of all fields except the signature field itself.

11.2. Hash Chain

Action envelopes are linked in a hash chain:

```
hash_n = SHA-256(hash_{n-1} || JCS(envelope_n))
```

The genesis hash is SHA-256("ATTP-GENESIS").

Breaking the hash chain indicates tampering. The TA MUST verify chain integrity on every append.

11.3. Non-Repudiable Receipts

After successful action execution, the TA issues a receipt signed by the TA's key. The receipt includes:

- o The action envelope.
- o The TA's signature over the envelope.
- o The hash chain position.
- o The compliance check result.

Receipts are non-repudiable: the agent cannot deny the action, and the TA cannot deny the trust decision.

12. Protocol Bindings

ATTP is transport-agnostic. Protocol bindings map ATTP trust enforcement to specific transports.

12.1. MCP Binding (MCPS)

The MCP binding is specified in draft-sharif-mcps-secure-mcp [MCPS]. MCPS wraps JSON-RPC messages in signed envelopes carrying ATTP trust claims.

Trust level is conveyed in the MCPS envelope header. Action limits map to tool call parameters (amount, scope). Kill

switch state is checked before tool execution.

12.2. REST Binding

For REST APIs, ATTP trust is conveyed via HTTP headers:

```
X-ATTP-Trust-Level: 2
X-ATTP-Agent-Id: agent_abc123
X-ATTP-Signature: <base64 ECDSA signature>
X-ATTP-Nonce: <UUID>
X-ATTP-Timestamp: <unix epoch ms>
```

The signature covers: method, path, body hash (SHA-256), nonce, and timestamp.

12.3. Google A2A Binding

For the Agent2Agent protocol [GOOGLE-A2A], ATTP trust is embedded in the Task metadata. The agent's AgentCard includes ATTP trust level and TA endpoint in its capabilities declaration.

12.4. gRPC Binding

For gRPC, ATTP trust is conveyed via metadata headers using the same header names as the REST binding.

12.5. GraphQL Binding

For GraphQL, ATTP trust is conveyed via HTTP headers (same as REST) or via an extensions field in the GraphQL request:

```
{
  "query": "...",
  "extensions": {
    "attp": {
      "trustLevel": 2,
      "agentId": "agent_abc123",
      "signature": "...",
      "nonce": "...",
      "timestamp": 1234567890
    }
  }
}
```

12.6. Keycloak Integration

For enterprises using Keycloak or other OIDC providers, ATTP trust claims are embedded in the JWT access token as a nested "attp" claim:

```
{
  "sub": "agent_abc123",
  "attp": {
    "trust_level": 2,
    "trust_label": "L2 -- Standard",
    "payment_enabled": true,
    "tx_limit": 10000,
    "day_limit": 50000,
    "scopes": "payment_initiate,sanctions_screen",
    "protocol_version": "1.0"
  }
}
```

A Keycloak protocol mapper implements the role-to-trust-level mapping (attp-l0 through attp-l4 roles).

13. Revocation

13.1. Agent Revocation

An agent MAY be revoked by the principal or by the TA. Revocation is immediate (kill switch) and permanent until explicitly reversed.

Revocation MUST:

- o Set the kill switch.
- o Invalidate the current passport.
- o Record the revocation in the audit trail.
- o Optionally notify dependent platforms via webhook.

13.2. Emergency Freeze

The TA MUST support emergency freeze of individual agents, all agents under a principal, or all agents globally. Global freeze MUST require multi-party authorisation.

14. Security Considerations

14.1. Key Management

Private keys MUST be stored in HSM, KMS, or platform-native secure storage for production deployments. Private keys MUST NOT be transmitted over network channels after initial generation. The TA MUST NOT store agent private keys.

14.2. Trust Farming Mitigation

ATTP mitigates trust farming through:

- o Promotion rate limiting (Section 5.7).
- o Self-dealing detection (Section 9.6).
- o Principal aggregate limits (Section 6.3).
- o Dormancy decay (Section 5.6).

An attacker who creates agents and performs small successful actions to build trust is limited by the minimum time at each level (24 hours, 7 days, 30 days, 90 days) and the minimum action counts. Reaching L4 requires at minimum 128 days of sustained operation with zero anomalies.

14.3. Sybil Attack Prevention

Principal aggregate limits prevent circumvention of per-agent limits through agent proliferation. The TA SHOULD implement additional heuristics: common IP addresses, identical behavioural patterns, coordinated action timing.

14.4. Impersonation Detection

Failed challenge-response verifications trigger trust penalties and audit entries. Three consecutive failures SHOULD trigger automatic agent suspension pending principal review.

14.5. Race Condition Handling

Concurrent action requests MUST be serialised at the limit enforcement layer. Implementations MUST use atomic

compare-and-swap operations for daily aggregate tracking.

14.6. Replay Protection

Action envelopes include nonce and timestamp. The TA MUST reject any envelope with a previously-seen nonce or a timestamp outside the acceptable window (RECOMMENDED: 5 minutes).

14.7. Kill Switch Consistency

Kill switch state MUST be checked atomically with action evaluation. There MUST be no window between trust check and action execution during which a kill switch activation could be missed.

14.8. Comparison with Identity-Only Protocols

Identity protocols (OAuth 2.0, OIDC, mTLS, HTTP Message Signatures, or any future agent identity standard) establish who an agent is. They do not:

- o Score agent trustworthiness based on behavioural history.
- o Enforce graduated action limits tied to trust levels.
- o Screen counterparties against sanctions lists.
- o Provide instant kill switches independent of credential lifecycle.
- o Produce tamper-evident audit chains for regulatory compliance.

ATTP is designed to compose with identity protocols, not replace them. An identity assertion is a prerequisite for ATTP trust evaluation.

15. Privacy Considerations

The public trust query API reveals that an agent exists and its current trust level. This is intentional -- platforms need this information to make access decisions.

The API MUST NOT reveal:

- o Action history or patterns.
- o Scoring dimension breakdowns.
- o Counterparty information.
- o Key material.
- o Principal identity (unless the principal opts in).

Implementations SHOULD support pseudonymous agent identifiers that are unlinkable across platforms unless the principal explicitly enables cross-platform linking.

16. IANA Considerations

This document defines the following HTTP headers for the REST protocol binding:

- o X-ATTP-Trust-Level
- o X-ATTP-Agent-Id
- o X-ATTP-Signature
- o X-ATTP-Nonce
- o X-ATTP-Timestamp

Registration of these headers in the IANA Message Headers

registry is requested.

This document defines the following error codes:

- o ATTP-ACTION-LIMIT
- o ATTP-SANCTIONS-MATCH
- o ATTP-TRUST-INSUFFICIENT
- o ATTP-KILL-SWITCH-ACTIVE
- o ATTP-NONCE-REPLAY
- o ATTP-TIMESTAMP-EXPIRED

This document defines the following well-known URI:

- o /.well-known/attp-trust

For discovery of Trust Authority endpoints.

17. References

17.1. Normative References

- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, DOI 10.17487/RFC2119, March 1997.
- [RFC8174] Leiba, B., "Ambiguity of Uppercase vs Lowercase in RFC 2119 Key Words", BCP 14, RFC 8174, DOI 10.17487/RFC8174, May 2017.
- [RFC8785] Rundgren, A., Jordan, B., and S. Erdtman, "JSON Canonicalization Scheme (JCS)", RFC 8785, DOI 10.17487/RFC8785, June 2020.
- [RFC7518] Jones, M., "JSON Web Algorithms (JWA)", RFC 7518, DOI 10.17487/RFC7518, May 2015.
- [FIPS186-5] National Institute of Standards and Technology, "Digital Signature Standard (DSS)", FIPS PUB 186-5, February 2023.

17.2. Informative References

- [MCPS] Sharif, R., "MCPS: Cryptographic Security Layer for the Model Context Protocol", draft-sharif-mcps-secure-mcp-00, March 2026.
- [AEBA] Sharif, R., "AEBA: Agent Event Behaviour Analytics", draft-sharif-aeba-00, April 2026.
- [OWASP-AISVS] OWASP, "AI Security Verification Standard", 2026.
- [RFC9421] Backman, A., Richer, J., and M. Sporny, "HTTP Message Signatures", RFC 9421, DOI 10.17487/RFC9421, February 2024.
- [GOOGLE-A2A] Google, "Agent2Agent Protocol", <<https://github.com/google/A2A>>.
- [MCP] Anthropic, "Model Context Protocol Specification", <<https://spec.modelcontextprotocol.io/>>.

Appendix A. Regulatory Mapping

A.1. PSD2 SCA Mapping

ATTP trust levels map to PSD2 Strong Customer Authentication requirements:

- o L0-L1: Equivalent to unauthenticated. No SCA-exempt transactions permitted.
- o L2: SCA-exempt for low-value transactions (< EUR 30, aggregate < EUR 100).
- o L3-L4: Challenge-response verification satisfies the "something you have" factor when combined with code attestation ("something you are" for agents).

A.2. PCI DSS v4.0.1 Mapping

- o Req 3.5 (key storage): ATTP mandates HSM/KMS for production private keys.
- o Req 6.2 (software security): Code attestation dimension maps to verified software inventory.
- o Req 8 (identify and authenticate): Challenge-response with ECDSA satisfies strong authentication.
- o Req 10 (logging): Hash-chained audit trail satisfies tamper-evident logging.
- o Req 11 (testing): Anomaly detection satisfies continuous monitoring requirements.

Appendix B. Trust Level Reference Implementation

A reference implementation of the ATTP trust scoring model is available at:

- o Go: github.com/razashariff/agentpass-go (Apache 2.0)
- o JavaScript: `npm install agentsign` (BSL 1.1)
- o Keycloak: github.com/razashariff/keycloak-mcps (BSL 1.1)

Appendix C. Comparison with Identity-Only Approaches

Identity-only protocols provide authentication. ATTP provides trust-gated authorisation. The following table summarises the differences:

Capability	Identity	ATTP
Prove agent identity	Yes	Yes*
Graduated trust levels	No	Yes
Action limits	No	Yes
Sanctions screening	No	Yes
Kill switches	No	Yes
Behavioural scoring	No	Yes
Tamper-evident audit	No	Yes
Fail-closed enforcement	Varies	Yes
Dormancy decay	No	Yes
Self-dealing detection	No	Yes

* ATTP consumes identity from external providers; it does not implement its own identity protocol.

Raza Sharif
CyberSecAI Ltd
205 Regent Street
London W1B 4HB
United Kingdom
Email: contact@agentsign.dev