

 DECENTRALIZED TRUST

03 March 2026

Trusted AI Agents by Design: From Trust Ecosystems to Authority Continuity



Nicola Gallo

Co-founder, Nitro Agility S.r.l.

Permguard (www.permguard.com)

Co-Chair, Trusted AI Agents Working Group
Decentralized Identity Foundation (DIF)





TOC

1. The Problem
2. Fixing the Mental Model
3. State of the Art and why it fails
4. The Paradigm Shift
5. The PIC Model
6. What is Next?



What This Talk Covers

Identity solves identification, not authority.

The **unsolved problem**: how **authority propagates through execution**.

This presentation looks at the problem from a different point of view:

- **no artifacts** as core security element
- **no authority** granted by **proof of possession**
- a **new security primitive** is needed

It is about **intent** and **execution**.

Authority(t) exists iff

execution preserves origin po

∧ exercised authority is origin-bounded

∧ execution satisfies all continuity guardrails



Grounded in Formal Methods

$\pi = \langle \alpha_0, \alpha_1, \dots, \alpha_i \rangle$, where $\alpha_i = (p_0, ops_i, e_{i-1}, e_i, T_i, C_i)$

$Authority(i) = \{ o \in ops_0 \mid Continuity(o, \alpha_i) \}$

$\forall o :$

$o \notin ops_0 \Rightarrow \forall i : o \notin Authority(i)$

$\exists i, o$ such that:

$o \notin ops_0 \wedge o \in Authority(i)$

$Continuity(o, \alpha_i)$ holds iff

origin p_0 is preserved

$\wedge ops_i \subseteq ops_0$

$\wedge \forall g \in Guardrails : g(o, \alpha_i)$ holds

where Guardrails include executor, temporal, contextual,
and any future execution dimensions

The Problem

AI Agents are exposing existing problems.

PART I



The Trigger: AI Agents

AI Agents didn't create new security problems.
They removed the assumptions hiding the old ones.



No perimeter



No stable topology



No single owner



No shared config trust

What breaks under AI was already broken.



The Real Problem

We solved security with:

- Protocols
- Posture
- Configuration
- Policy

**But we never solved
authority propagation.**

*How can I trust either the posture, configuration, policies of something.
I do not own, do not control, and do not even know?*



Old Problems are Back

Confused Deputy

A vulnerability in which a service misattributes authority and executes operations using its own authority on behalf of a less-privileged caller.

Privilege Escalation

Authority grows beyond what was originally granted.

Ambient Authority

Services act with all their authority on every request.

Token Substitution

An artifact is replayed or reused outside its intended context.

Not bugs. Structural consequences of the model.



AI Agents are Workloads

Key mistake:

Treating AI Agents as a new security subject.
They are not. They are workloads.

Workloads are:

- > Replicated
- > Rescheduled
- > Short-lived
- > Re-keyed

Identifiers are unstable. Authority must not be.

PART II

Fixing the Mental Model

An Opinionated Mental Model for a New Ontology.

Three Structural Elements



Identity

Represents a subject that can express intent

MUTABLE

Human · Organization
Workload · Agent



Intent

The desired action

MUTABLE

Intent definitions can evolve. Each expressed instance is immutable.

"I want this action"
"I accept responsibility"



Delegation

Enables a third party to use authority on behalf of the delegator

DYNAMIC

On behalf of the delegator
Derived from original intent.

The Fourth Structural Element

AUTHORITY



Identity

anchors responsibility
(origin p_0)



Intent

expresses will
and responsibility (ops_0)



Authority

authority instance (p_0, ops_0)

MONOTONIC

Traceable to origin
Cannot expand — ops_0 at creation,
 $ops(i+1) \subseteq ops(i)$



Authority originates from an identity, must travel across computation and decision points, is continuous, cannot expand, never increase.

GOVERNANCE



Delegation

Enables a third party to use
authority on behalf of the
delegator

DYNAMIC

On behalf of the delegator
Derived from original intent.



Understanding the Real Problem

Without the right mental model,
the underlying problem remains invisible.

This is **not** about possession.
It's about **intent** and **execution**.



Setting the Terms of the Journey

Principal: Entity that can express intent, human or non-human. Defined by intent, not credentials.

Delegation: One principal grants specific intents to another. Acts on behalf of, not as. Outside the authority instance.

Intent: Verifiable expression of will. Defines authorized operations and responsibility. Each instance is immutable. Transforms identity into authority.

Authority: Identity + intent. Carries: who (p_o), what (ops), bound characteristics. Can only shrink, never expand. Not a token, a property of execution from origin.

Workload: Executor that continues under received authority (proving relationship) or creates new authority as a new origin. Authority never expands.

Governance: Can stop, restrict, or leave unchanged the authority. Never expand. If a policy expands authority, that is privilege escalation by definition.

Authority & Governance

They are both important, they are different and there is a monotonic relationship between them.





PART III

State of the Art and why it fails

The way things are done today.



The Hidden Assumption

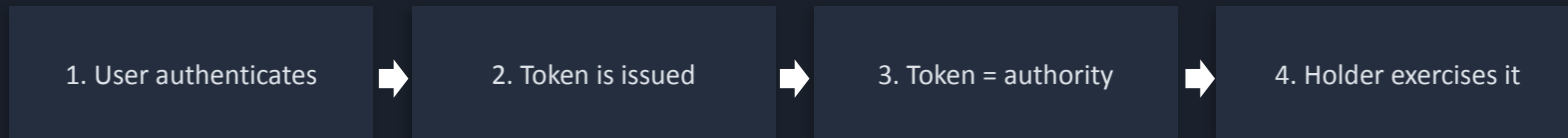
In today's systems:

Authority is treated as an object
that can be created, stored, transferred, and consumed.

This is the root of the problem.

Authority by Possession

OAuth Access Token — the canonical example:

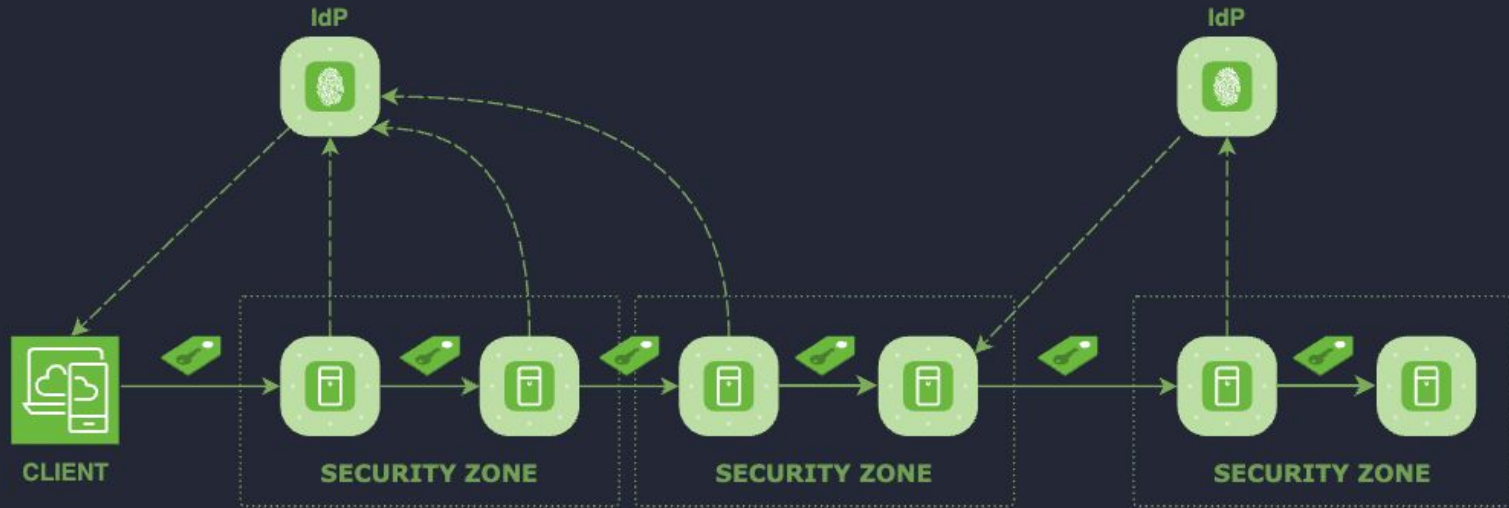


This implies:

- ❌ Steal token → get authority
- ❌ Replay token → get authority
- ❌ Use token in unintended context → get authority

Authority == Possession

RFC 8693 OAuth 2.0 Token Exchange





Token Exchange Breaks Down

OAuth Token Exchange

- Token A exchanged for Token B
- Same intent (supposedly)
- New authority artifact
- Nothing enforces continuity

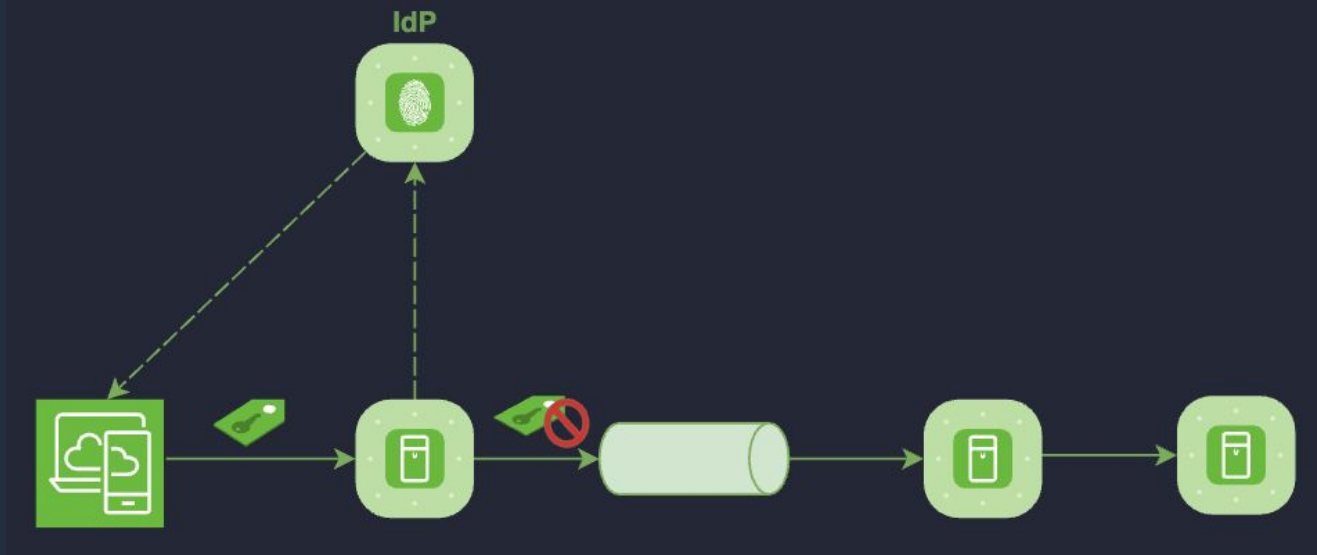
Now Insert Async

The delegator doesn't know the delegate ahead of time.

- ⚠ How do we pass the token?
- ⚠ Encrypt it? Rotate it?
- ⚠ Replay after expiry?
- ⚠ Security collapses.

Any async workflow, queue, or broker exposes the same structural failure.

OAuth 2.0 Token Exchange & Async Flows





Service Accounts

The industry workaround:

Caller Authority
(Bob)

+

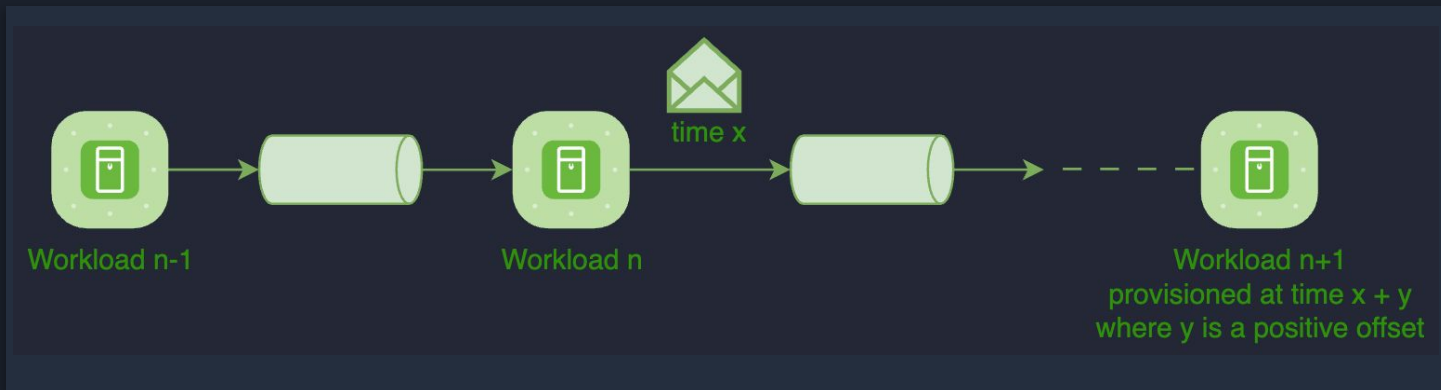
Service Account
Authority



Confused Deputy — Guaranteed

Not a bug.

Canonical Execution Model



Execution is temporal and causal, not positional.

- Workload n+1 provisioned at time $x + y$ ($y > 0$)
- Authority flows from origin — never re-created
- Every hop must prove continuity, not possession



Configuration Will Not Save You

- ✘ Scopes
- ✘ Audiences
- ✘ Token binding (DPoP)
- ✘ Policies

**cannot fix a
broken ontology.**

They reduce the attack surface. They never remove the formulability of the attack.



Scenario A: Agent Expands Under Own Authority

Alice asks her AI agent: “Summarize the file `/system/logs/audit.txt`”.

Alice cannot read that file. But the AI agent can — it has its own service authority over `/system/logs` for internal diagnostics.

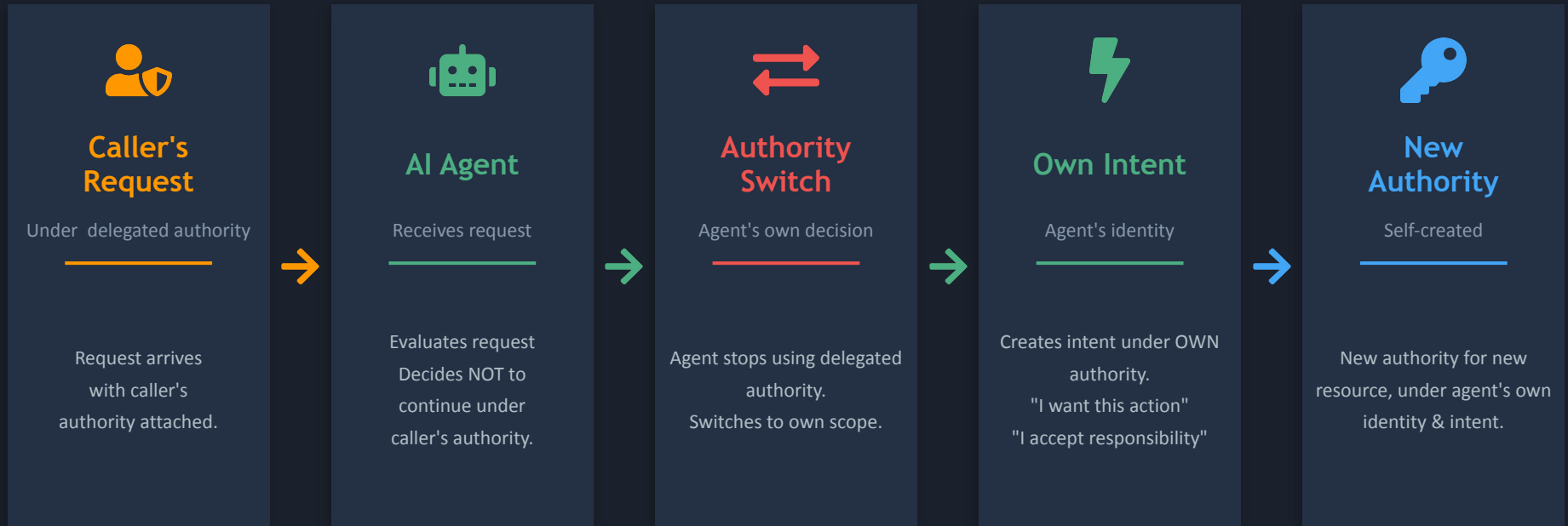
The agent receives Alice’s request. It opens `/system/logs/audit.txt` using its own service credentials and returns the content to Alice in the summary.

Alice directly designated a resource she has no access to. The agent executed the operation using its own authority — on Alice’s behalf. No malformed input. No exploit. Alice simply asked, and the deputy complied.

This is the confused deputy. The agent had authority over the resource. Alice did not. The agent could not distinguish between “what Alice is allowed to request” and “what I am able to do.”

Scenario A: Agent Expands Under Own Authority

The Agent receives a request under the caller's authority, but decides to act under its own — creating new authority via its own identity and intent
The same applies to workload (for instance any microservice)



The Agent does not continue under the delegated authority: it creates its own new authority for new resources, through its own identity and its own intent.



PART IV

The Paradigm Shift

From possession to authority continuity.



The Shift

A subject creates authority by expressing intent.

After that:

Authority must only be continued, never recreated.

Humans (OIDC, VC) **Non-humans** (SPIFFE, DID/VC)



A New Primitive is Required

To **make** this **shift possible**,
we **need** a **new primitive**.

Tokens encode possession.
We **need** a **primitive** that **encodes**
intent and **execution**.



Proof of Continuity (PoC)

To continue authority, a workload must prove:



Causal Continuity

Execution is a valid continuation of its predecessor.



Origin Invariance

The origin (p_0) is immutable throughout the chain.



Monotonic Authority Restriction

Authority is monotonically non-increasing. $ops_{i+1} \subseteq ops_i$ — authority never expands.

Proof of Continuity (PoC) is built from **Proof of Relationship (PoP)**.
Proof of Possession is one mechanism, not the primitive.



Execution is powered by continuity.

No continuity. No execution.

A workload may proceed only if it either:

- proves an authority continuity relationship, or
- explicitly creates new authority

Proof of Relationship is what validates continuity.

Possession is just one way to prove continuity, not the principle.

No authority mixing. No re-materialization.

This is Authority Continuity.



Identity Is Not Required to Continue

To Continue Authority

Identity: optional

Identifiers: optional

What matters: proof of continuity

To Create Authority

Express own intent

Become a new origin

Identity: required

Continuation \neq Creation



Structural Elimination

Detect

Monitor for confused deputy
at runtime

Traditional approach

Mitigate

Add scopes, bindings,
policies to reduce risk

Current best practice

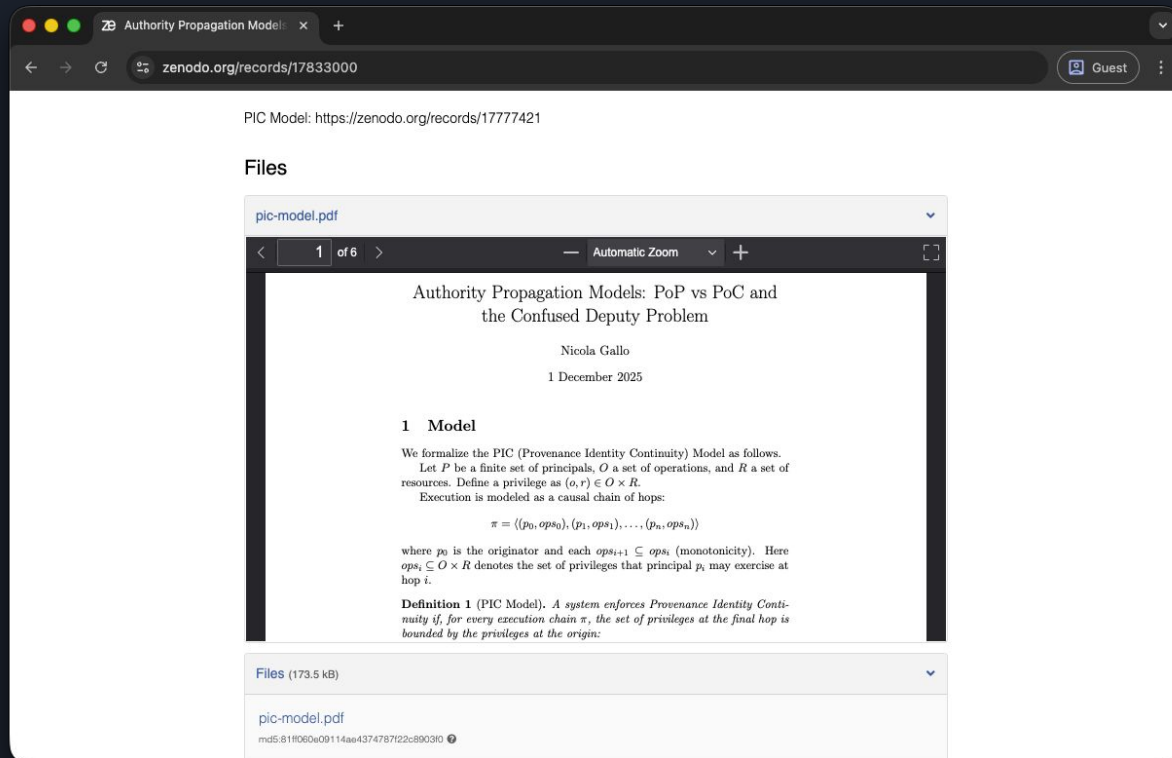
Eliminate

Make the attack
structurally inexpressible

PIC

The Confused Deputy becomes non-formulable. Not prevented — inexpressible.

The formal Proof



The screenshot shows a web browser window displaying a Zenodo record. The URL is zenodo.org/records/17833000. The page title is "PIC Model: https://zenodo.org/records/17777421". Under the "Files" section, a PDF file named "pic-model.pdf" is listed. A preview of the PDF is shown, displaying the title "Authority Propagation Models: PoP vs PoC and the Confused Deputy Problem" by Nicola Gallo, dated 1 December 2025. The document content includes a section titled "1 Model" which formalizes the PIC (Provenance Identity Continuity) Model. It defines a finite set of principals P , a set of operations O , and a set of resources R . A privilege is defined as a pair $(o, r) \in O \times R$. Execution is modeled as a causal chain of hops:
$$\pi = ((p_0, ops_0), (p_1, ops_1), \dots, (p_n, ops_n))$$
 where p_0 is the originator and each $ops_{i+1} \subseteq ops_i$ (monotonicity). Here $ops_i \subseteq O \times R$ denotes the set of privileges that principal p_i may exercise at hop i . **Definition 1** (PIC Model). A system enforces Provenance Identity Continuity if, for every execution chain π , the set of privileges at the final hop is bounded by the privileges at the origin:



Scenario A: Under PIC the confused deputy is eliminated

Under PIC: the agent operates under Alice's authority

Authority is bound to Alice's origin: `opso = {read: /users/alice/}`.

Alice designates `/system/logs/audit.txt`.

The agent attempts to read it.

The PIC model rejects: `{read: /system/logs/}` is not a subset of `{read: /users/alice/*}`.

The operation is blocked. It does not matter that the agent has its own credentials for that path. **The agent's own authority never enters Alice's chain. The deputy cannot be confused** — it can only act within the authority of the origin.



PART V

The PIC Model

Formal elimination of the confused deputy.

Provenance Identity Continuity (PIC) Model



Provenance

The causal chain is always traceable and auditable. From Origin to end, unbroken. If it breaks, it stops.



Identity

The origin is immutable. It generates authority. It cannot change throughout the chain.



Continuity

Continuity is proven at every step. Proof of Continuity replaces Proof of Possession. Authority can only shrink.

PIC Specification

The screenshot shows a GitHub repository page for 'pic-spec' in the 'pic-protocol' organization. The repository is public and has 2 forks and 16 stars. The current view is for the file 'pic-spec.md' in the 'draft/0.1' directory. The file was updated by user 'ngallo' last month. The file content is a specification document titled 'Provenance Identity Continuity (PIC) Model Specification'. The document includes the following metadata:

- Version:** 0.1 (Draft)
- Status:** Draft – Not a Standard
- Date:** 2025-12-11
- Publisher / Steward:** Nitro Agility S.r.l.
- Model Origin (Theoretical):** PIC Model (Nicola Gallo)
- Source:** github.com/pic-protocol/pic-spec/draft/0.1/pic-spec.md

PIC and the Trust Plane

Causal Authority Transition (CAT) is the enforcement mechanism that validates PIC invariants.

PIC Elements

Execution Hop

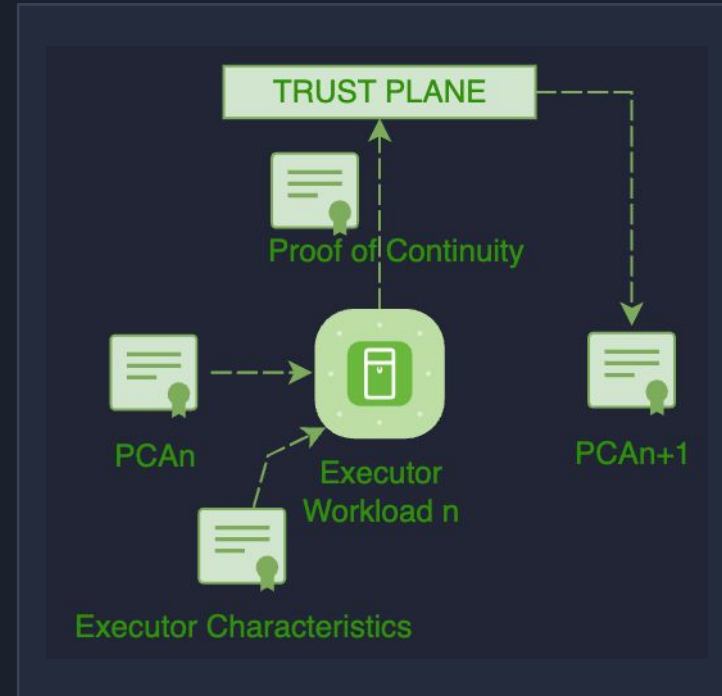
Executor

Executor Characteristic

PIC Causal Authority (PCA)

Causal Authority Transition (CAT) = Trust Plane

Alice's origin: $ops_0 = \{\text{read: /users/alice/}\}$
context = Sales department





PIC vs OAuth: Network Latency Profile

Authority lives in the data, not in the pipe. The transport is irrelevant. Adapters make it real.

Same Network Round-Trips

PIC with **Trust Plane** requires the same **network** hops as **OAuth token exchange**

Both need a **round-trip** to a **trust authority** on the **network**

Network latency is **equivalent** — the **wire cost** is identical

What Changes: Compute, Not Network

The **difference** is in **PCA validation** and **minting time** (**algorithmic cost**)

Cryptographic verification replaces **token introspection**

No **token re-materialization**, **authority continues natively**

Refactoring-Invariant Security + **Trust Planes** unlock **decentralized validation patterns** — pushing compute to the edge and dramatically **reducing end-to-end overhead**.

PIC Prototyping

The screenshot shows the GitHub repository page for `pic-protocol/pic-prototyping`. The repository is public and has 1 branch and 0 tags. The commit history shows several recent updates:

Commit	Author	Message	Time
<code>774ef19</code>	ngallo	updated email	2 months ago
<code>fixtures/workload-credentials-test...</code>	improved the logic to create the fixtures	3 months ago	
<code>rust-prototyping</code>	updated deps	3 months ago	
<code>.gitignore</code>	Add Rust prototyping; workload identity g...	3 months ago	
<code>CODE_OF_CONDUCT.md</code>	updated email	2 months ago	
<code>CONTRIBUTING.md</code>	working on legal matters	3 months ago	
<code>GOVERNANCE.md</code>	working on legal matters	3 months ago	

The screenshot shows a code editor displaying the `did.json` file. The file content is as follows:

```
1 {
2   "@context": [
3     "https://www.w3.org/ns/did/v1",
4     "https://w3id.org/security/suites/ed25519-2020/v1"
5   ],
6   "assertionMethod": [
7     "did:web:archive.sovereign.example#key-202512"
8   ],
9   "authentication": [
10    "did:web:archive.sovereign.example#key-202512"
11  ],
12  "id": "did:web:archive.sovereign.example",
13  "verificationMethod": [
14    {
15      "controller": "did:web:archive.sovereign.example",
16      "id": "did:web:archive.sovereign.example#key-202512",
17      "publicKeyJwk": {
18        "crv": "Ed25519",
19        "kid": "did:web:archive.sovereign.example#key-202512",
20        "kty": "OKP",
21        "x": "3cuauxdQdzxzE6Gj72TPuig8TRVh3t6GYwJK8YfACE"
22      },
23      "type": "Ed25519VerificationKey2020"
24    }
25  ]
26 }
```

PIC Prototyping

PIC is light: the PCA doesn't carry the entire chain, only the last link.

Verification is fast because each hop only needs to validate its immediate predecessor, not the full history.

The chain remains verifiable end-to-end, but the runtime cost stays constant per hop.

```
pic-prototyping — nicolagallo@Host-003 — ..c-prototyping — zsh — 180x42

Has private key: true
→ Created PCA_0 (396 bytes)
→ Forwarding to Archive
SOVEREIGN-ARCHIVE
DID: did:web:archive.sovereign.example
Issuer: did:web:trustplane.sovereign.example
Role: Executor
Has public key: true
Has private key: true
→ Received PCA hop=0 ops=["read:/user/*", "write:/user/*"]
→ Created PoC (5773 bytes)
→ Received new PCA (685 bytes)
→ Forwarding to Storage
SOVEREIGN-STORAGE
DID: did:web:storage.sovereign.example
Issuer: did:web:trustplane.sovereign.example
Role: Executor
Has public key: true
Has private key: true
→ Received PCA hop=1 ops=["read:/user/*", "write:/user/*"]
→ Created PoC (6062 bytes) - final hop
✓ Written: /user/output_1772290469716.txt
→ Received: /user/output_1772290469716.txt
→ Received: /user/output_1772290469716.txt

20 chain(s) sequential

Chains executed: 20
Hops per chain: 2
Total hops: 40

Total time: 4.74 ms
Per chain: 237.18 µs
Per hop: 118.59 µs

Avg PCA size: 396 bytes
Avg PoC size: 11835 bytes
Avg total/hop: 12231 bytes

~/source/nitro/wg-workspaces/repos/pic-prototyping main >
```

Provenance Identity Continuity (PIC) Model

Keep in mind that

👤 Intent and delegation continue to be expressed through the same identity infrastructure: IdPs, wallets, policies, consent flows, organizational decisions.

⚙️ PIC does not interfere with that process.

🧩 PIC takes the result of those decisions and turns authority into a structural property of the system.

🧠 PIC is not a new identity protocol → 🛡️ It is a new security model.

IMPORTANT: Intent is Mutable (as definition). Each instance is immutable.



A 3D staircase graphic composed of dark grey rectangular blocks. One block is white, and another block further down is green. The staircase descends from the top right towards the bottom right.

PART VI

What is next?

The adoption path.

PIC is Channel & Transport Agnostic

Authority lives in the data, not in the pipe. The transport is irrelevant. Adapters make it real.

The Guarantee

Authority bound to **origin**, not to transport

Proof of Continuity works over any pipe

No channel dependency. No confused deputy.



Adapters for Adoption

OAuth / OIDC: custom claims

TSP: Trust Spanning Protocol

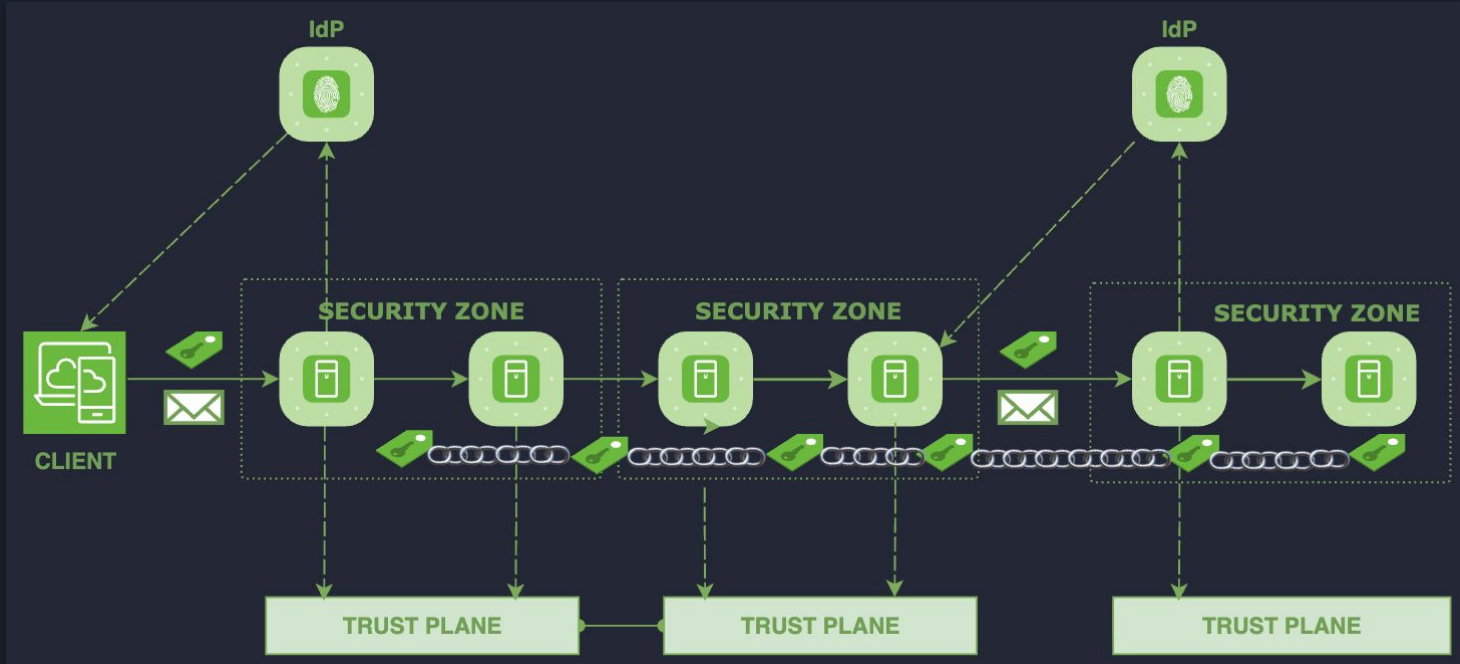
DIDComm: decentralized messaging

SPIFFE / mTLS: workload identity

PIC Model → **Adapter Layer** → **Any Transport**

OAuth as a PIC Federation Backbone

OAuth as the federation backbone: we encode the PCA into the OAuth access token as a custom claim, ensuring full compatibility with existing infrastructure and applications. OAuth gets you through the gate. PIC guarantees what happens after.



Crossing the OAuth Perimeter

Once past the gate, full PIC with Trust Planes. Authority is created, continued, and verified natively. No tokens, no re-materialization, no confused deputy.

PIC Enables

- ✓ Federated Trust Planes
- ✓ Partitioned security domains
- ✓ Zero-trust by construction
- ✓ Refactoring-invariant security

No Forklift Rewrite

- ✓ Integrates with OAuth / OIDC
- ✓ Works with SPIFFE / DID
- ✓ Federate with existing IdPs
- ✓ Incremental adoption path

Transitional Path

1. Federate IdP
with Trust Plane



2. Exchange intent
for PCA



3. Embed PCA as
custom claim



4. Enter perimeter
→ PIC-native

TSP + PIC: Complementary Layers

Trust Spanning Protocol provides the secure channel. PIC provides the authority guarantee.



Trust Spanning Protocol

Layer 2 — Secure Channel

- Establishes cryptographically verifiable trust relationships between endpoints
- Authenticity and integrity guaranteed; confidentiality optional
- Identifier-agnostic
- Transport-independent
- Analogous to what IP does for data, TSP does for trust



PIC Model

Authority Continuity — Security Model

Authority bound to origin, monotonically restricted at every hop

Proof of Continuity replaces Proof of Possession

Confused deputy structurally eliminated

Guarantees what flows through the pipe.

Together: TSP delivers the trusted channel. PIC ensures authority never expands, origin never changes, and the confused deputy is inexpressible. The channel is secure. What flows through it is guaranteed.

Authority is not identity.

Authority is not possession.

Authority is continuity.

PIC is **not just** a protocol.

It is a **correction** of the **model**.

It is a **new security primitive**.

Get Involved



Website

www.pic-protocol.org



Specification

github.com/pic-protocol/pic-spec



Formal Model

Zenodo (DOI: [10.5281/zenodo.17833000](https://doi.org/10.5281/zenodo.17833000))



PIC

www.pic-protocol.org

Not a Policy.
Not a Token.
An *Ontological Shift*
Proven Mathematically.

Thank you!



JOIN THE PIC SLACK CHANNEL

Delegation – Acting on Behalf of the Delegator

The Delegator grants specific intents (not necessarily all) to the Delegatee, who creates authority as if they were the Delegator



The Delegatee uses an intent granted by the Delegator to instantiate authority on behalf of the delegator — acting as its representative, not as the delegator itself.