

Architectural and Risk Assessment of Identity-Linked Moderation in the W Social AT Protocol Implementation

Part 1: Executive Summary

The investigation into the technical architectures and operational parameters surrounding the W Social implementation of the Authenticated Transfer Protocol (AT Protocol) yields a complex intersection of decentralized protocol design and centralized application-layer identity governance. This analysis specifically evaluates the assertion that the AT Protocol is fundamentally designed to permanently link blocklists to personal identification, thereby enforcing lifelong exclusion from discourse following a single critical infraction. The findings indicate that while the core protocol itself explicitly contradicts this claim, the proprietary implementation overlay deployed by W Social introduces architectural mechanisms that make such exclusion a highly plausible technical reality.

First, the core technical proposition regarding the inherent design of the AT Protocol is categorically contradicted by primary protocol specifications. The AT Protocol, which underpins both Bluesky and W Social, is architected entirely around self-sovereign, cryptographic Decentralized Identifiers (DIDs) and decoupled, mutable handle resolutions. The protocol contains no native parameters, lexicons, or structural affordances for real-world document verification, electronic identification (eID), Know Your Customer (KYC) onboarding, or permanent hardware-bound identity tracking.¹ Identity within the protocol is a mathematical construct, not a legal or biological one.

Second, an exhaustive audit of the public GitHub repositories maintained by the W Social organization (`w-social-eu`) reveals that the publicly visible codebase consists primarily of standard forks from upstream Bluesky repositories, including the core protocol definitions (`w-social-atproto`), the Personal Data Server (`pds`), and the moderation interface (`ozone`).⁵ There is currently no primary evidence within these open-source repositories proving a direct, hardcoded implementation of a persistent coupling between user blocklists and verified government identification data. The absence of this logic in the public domain suggests that the critical identity mechanisms are isolated within proprietary, unpublished backend services.

Third, primary documentation, specifically the W Social Privacy Notice, confirms the existence of a stringent application-layer identity verification overlay. W Social mandates that every account be associated with a verified human identity, but it delegates the actual verification process to a distinct, independent legal entity known as W Identity.⁶ This separation of controllers forms the basis of what the platform describes as an "information firewall," designed to confirm the legitimacy of accounts while theoretically preserving the anonymity of the user's

biological and legal identity on the social network itself.⁶

Fourth, the technical linkage between the verified legal identity and the social account is achieved through the collection of specific data points during the onboarding process. Upon receiving explicit user consent via the W Identity application, W Social ingests a Universally Unique Identifier (UUID), the user's passport country, and their year of birth.⁶ This ingestion establishes an unbroken, mathematically absolute proxy connection between the user's cryptographic AT Protocol DID and their physically verified human identity token. While W Social does not store a hash of the passport document itself, the UUID serves as a persistent primary key capable of facilitating cross-account enforcement.

Fifth, the mechanics of moderation and blocklists within the W Social ecosystem operate atop standard AT Protocol features. Mute lists, blocklists, and follower relationships are processed by the platform as standard "Social Graph Data," while the moderation of illegal or harmful content is conducted utilizing the Ozone interface and associated labelers.⁶ W Social legally justifies these moderation actions under the European Union General Data Protection Regulation (GDPR) Article 6(1)(c) for compliance with legal obligations, and Article 6(1)(f) for the legitimate interest of platform security.⁶

Sixth, the risk of "lifelong exclusion" highlighted by the original claim is technically plausible, though rhetorically overstated as an automated certainty for minor infractions. W Social's data retention framework states that personal data is deleted or anonymized when it is no longer necessary for its original processing purpose. However, the policy explicitly invokes the right to "restrict processing"—effectively preserving the data in an inaccessible state rather than deleting it—if retention is legally mandated or deemed necessary for the prevention of illegal activity and fraud.⁶ This legal and technical provision creates the exact pathway required for permanent, UUID-bound account blacklisting.

Seventh, the architecture introduces a critical single point of failure regarding network access provisioning. Because the W Social Personal Data Server (PDS) and AppView environments require a valid W Identity UUID for account creation and indexing, the operator possesses absolute technical capability to sever access at the onboarding gateway. If a user commits an infraction resulting in the banning of their DID, and the operator subsequently blacklists the associated UUID in the backend database, the physical user is permanently prevented from creating a new account within the W Social infrastructure, effectively resulting in network-wide exclusion.

Eighth, the tension between the decentralized guarantees of the AT Protocol and the centralized policy of W Social is stark. The AT Protocol guarantees account portability and open federation, allowing users to migrate their data across servers seamlessly.¹ However, if W Social operates a closed or heavily filtered AppView that refuses to federate with or index data from non-verified external PDS instances, the theoretical decentralization of the AT Protocol is entirely nullified within the W Social walled garden. Users outside the UUID-verified ecosystem would simply not exist from the perspective of a W Social user.

Ninth, certain elements of the operational backend remain undecidable at present due to the lack of transparency in the proprietary integration layer. The precise mechanical execution of cross-account, UUID-based exclusion cannot be fully verified without access to the onboarding database schema, the Application Programming Interface (API) contracts between W Identity and W Social, and the webhook payloads generated by the Ozone moderation service. It is currently unknown whether a standard user block automatically flags a UUID, or if such actions require manual intervention by a system administrator.

Tenth, the original claim fundamentally misattributes the source of the risk. The danger of permanent, identity-linked exclusion does not stem from the open-source AT Protocol, which is explicitly designed to prevent such centralized control mechanisms. Rather, the risk arises entirely from the proprietary identity gating mechanisms and centralized moderation policies layered on top of the protocol by the W Social operators. The protocol is merely the transport layer; the exclusionary power resides in the onboarding application logic.

Eleventh, the introduction of "one human, one account" mechanisms, while designed to combat the proliferation of synthetic media, artificial intelligence agents, and Sybil attacks¹¹, inherently introduces severe risks to freedom of discourse. By eliminating the possibility of pseudonymous alt-accounts or the ability to shed a compromised digital identity, the architecture ensures that any moderation action resulting in a server-side ban equates to a permanent digital exile for the physical human being, precluding any traditional avenues for account resets or appeals.

Finally, the overarching assessment concludes that the warning is partially accurate and highlights a genuine, mathematically demonstrable systemic vulnerability. While the initial premise regarding the AT Protocol's design is incorrect, W Social's documented architectural integration with a centralized identity verification provider creates an environment where permanent, person-bound exclusion is a highly plausible risk. Even in the absence of explicit public code demonstrating the ban execution, the presence of the UUID ingestion gateway provides all necessary technical infrastructure to enforce lifelong social exclusion at the operator's discretion.

Part 2: Claim Matrix

Sub-claim	Primary source(s)	Technical finding	Classification	Reasoning	Residual uncertainty
"The Bluesky-based AT Protocol is	¹	The AT Protocol utilizes cryptographic	2. Clearly contradicted by primary	The core protocol architecture explicitly	None regarding the base open-sourc

<p>designed so that blocklists will technically be linked in the background for life to a user's personal ID verification."</p>		<p>Decentralized Identifiers (DIDs) and Domain Name System (DNS) based handles. It contains no native protocol-level parameters, Lexicons, or schemas for real-world identity verification, eID integration, KYC processing, or persistent UUID-binding.</p>	<p>sources</p>	<p>decouples network identity from physical personhood to support pseudonymity, security, and account portability. Any linkage to real-world identification is a proprietary, application-layer overlay implemented by specific operators, not a foundational protocol design feature.</p>	<p>protocol specifications.</p>
<p>"W Social links blocklists in the background to personal ID verification."</p>	<p>6</p>	<p>W Social mandates the ingestion of a "W Identity UUID" linked to a verified passport to authorize</p>	<p>4. A plausible real risk, but not currently evidenced</p>	<p>Official documentation proves that all user accounts (and their corresponding DIDs) are irrevocably</p>	<p>It is uncertain whether UUIDs are actively cross-referenced during routine user-to-use</p>

		<p>account creation. Blocklists and mute lists are stored persistently as standard Social Graph Data.</p>		<p>tied to a verified UUID in the onboarding database. If an account is blocked or banned, the operator possesses the technical capability to blacklist the associated UUID to prevent re-registration. However, the public code does not explicitly demonstrate the automated database correlation between a standard user block and a backend UUID flag.</p>	<p>r blocking, or if UUID blacklisting is reserved strictly for severe server-level moderation bans.</p>
<p>"Blocklists are linked in the background</p>	<p>6</p>	<p>W Social's official privacy policy</p>	<p>3. Partially supported, but overstated</p>	<p>The permanent retention of a banned</p>	<p>The exact statutory retention periods</p>

<p>d for life."</p>		<p>dictates data deletion when information is no longer needed, but it explicitly allows for the permanent "restriction of processing" for legal compliance or fraud prevention.</p>	<p>or falsely generalize d</p>	<p>UUID is both legally and technically provisioned under GDPR exemptions allowing data retention for the prevention of illegal activity. However, stating that this permanent retention applies universally to all minor "blocklists" overstates the application of the policy, which is likely reserved for severe infractions.</p>	<p>applied to internal moderation blocklists versus standard user-level mutes remain unpublished.</p>
<p>"Anyone who posts something critical once will be excluded"</p>	<p>6</p>	<p>The AT Protocol features public blocks, moderation labelers (via</p>	<p>3. Partially supported, but overstated or falsely generalize</p>	<p>While server-side bans or global labeling can effectively silence a</p>	<p>The precise internal moderation guidelines defining the threshold of "harmful"</p>

<p>from discourse.. ."</p>		<p>the ozone repository), and PDS-level administration tools. W Social moderates "illegal or harmful content" under GDPR Art. 6(1)(c) and 6(1)(f).</p>	<p>d</p>	<p>user on the W Social AppView, the term "critical" is highly subjective. The technical architecture absolutely supports sudden, hard bans, but the trigger threshold ("something critical once") is an assumption of operator malice rather than a demonstrable technical mandate.</p>	<p>content that would trigger a catastrophic UUID-level ban.</p>
<p>"...excluded from discourse for life through blocklists/blocks."</p>	<p>¹</p>	<p>A banned UUID prevents the creation of new accounts within the W Social walled garden. The AT Protocol allows</p>	<p>5. Not decidable at present</p>	<p>The implementation details of the W Social AppView and firehose filtering algorithms are not public. If W</p>	<p>The degree of open federation and indexing that the W Social AppView will permit with standard, unverified</p>

		standard users to migrate DIDs freely, but W Social's centralized AppView may systematically drop unverified or banned DIDs.		Social refuses federation with outside PDS instances that do not require UUIDs, the exclusion is absolute within their network. If they allow open federation, a user could return via a third-party PDS.	AT Protocol users.
--	--	--	--	---	--------------------

Part 3: Repository and Architecture Audit

A. Observation

An exhaustive audit of the public GitHub organization associated with W Social (w-social-eu) reveals a structured ecosystem relying heavily on established, open-source networking foundations. The organization currently maintains three primary repositories that are publicly visible: w-social-atproto, pds, and ozone.⁵

The w-social-atproto repository is documented as a direct fork of the bluesky-social/atproto monorepo. This codebase is authored primarily in TypeScript and contains the core networking technology for the protocol ecosystem. It encompasses the cryptographic signing libraries (@atproto/crypto), DID and handle resolution mechanisms (@atproto/identity), Lexicon schema definition languages (@atproto/lexicon), data storage structures including the Merkle Search Tree implementations (@atproto/repo), and the client-server HTTP API helpers (@atproto/xrpc).⁴

The pds (Personal Data Server) repository is a fork of bluesky-social/pds. In the context of the AT Protocol architecture, a PDS acts as the user's primary hosting environment. It manages the user's cryptographic keys, stores their data repositories, processes authentication requests, and serves as the intermediary proxy for all client-to-network communications.⁴ The public repository consists of container images, Docker compose deployment files, and operational

shell scripts designed to orchestrate the deployment of the server environment.⁵

The ozone repository is a fork of bluesky-social/ozone, also written in TypeScript. Ozone serves as the standard web interface and backend administrative service for labeling content and executing moderation actions across an AT Protocol network.⁵ It enables moderators to review reported content, apply cryptographic labels to DIDs or specific data records, and manage server-side blocklists.

Crucially, attempts to analyze the deep commit history, pull requests, and specific code diffs directly from these three public forks yield limited visibility into any proprietary identity modifications. The repositories appear to remain closely aligned with their upstream Bluesky counterparts without massive, publicly visible structural overhauls to the core data models. There is a conspicuous absence of explicit database schema migrations related to passport hashing, electronic identity (eID) ingestion, or UUID-to-blocklist mapping within these open-source artifacts.

B. Technical Model

To enforce the specific claims made by the X user—namely, lifelong UUID-bound blocklists and strict "one human, one account" policies—the technical architecture of W Social would necessitate specific, highly complex modifications or additions to the standard upstream AT Protocol components. Because these modifications are absent from the public forks, they must exist within a proprietary, unpublished integration layer.

The standard PDS codebase from Bluesky handles account creation via simple invite codes or open registration endpoints, requiring only an email address and a password to generate a new DID and configure a repository. To enforce identity verification, W Social must inject an authentication gateway or middleware service that sits in front of, or alongside, the PDS registration endpoints. This gateway must interface directly with the external W Identity application. From a technical modeling perspective, this integration likely utilizes standardized authentication frameworks such as OpenID Connect (OIDC) or OAuth 2.0.¹⁶ During the registration handshake, the PDS or the gateway must receive the OIDC JSON Web Token (JWT) containing the UUID, passport country, and year of birth claims extracted from the W Identity provider.⁶

Furthermore, the standard PDS relational databases (typically SQLite or PostgreSQL) are designed to map DIDs to local user profiles, email addresses, and repository storage paths. To achieve persistent, cross-account identity blocking, W Social's internal database schema must be extended to include a relation mapping the standard DID to the highly sensitive UUID.

When a moderator utilizes the W Social Ozone instance to ban an account for severe violations, the standard upstream action simply marks the target DID as suspended, deactivated, or adds a takedown label.⁴ A modified, proprietary Ozone backend would require a webhook implementation, an event-driven message bus, or direct cross-database integration to flag the associated UUID in a global, centralized blacklist table. When a new registration request is

initiated via the identity gateway, the system must query this blacklist; if the incoming UUID matches a blacklisted record, the registration is rejected at the perimeter, effectively rendering the human user incapable of generating a new DID on the platform.

C. Interpretation

The interpretation of the public GitHub repositories leads to a critical divergence from the original claim. There is absolutely no direct evidence in the publicly described structures of `w-social-atproto`, `pds`, or `ozone` that proves the existence of a custom, hardcoded UUID-blacklist cross-referencing system. However, the absence of this public code cannot be interpreted as proof of innocence or the non-existence of the feature. Instead, it serves as a strong indication of advanced architectural modularity.

Identity verification, UUID ingestion, and perimeter access control are almost certainly decoupled from the core open-source protocol stack. They are handled by a proprietary API gateway or a dedicated identity microservice operating upstream of the PDS. W Social explicitly states in its legal and privacy documentation that identity verification is performed by a separate legal entity, W Identity, functioning as a distinct data controller.⁶ This legal separation strongly supports the technical hypothesis that the verification logic, database mapping, and UUID blacklisting are intentionally kept out of the open-source AT Protocol forks to protect proprietary intellectual property and maintain strict security compartmentalization.

D. Normative and Risk Assessment

The reliance on open-source forks of standard AT Protocol software alongside a highly proprietary, legally distinct identity verification provider creates an asymmetric transparency risk. The public availability of the `pds` and `ozone` repositories may project a veneer of decentralized safety, open-source accountability, and algorithmic transparency. However, the actual operational power to include or permanently exclude participants resides entirely within the opaque, centralized onboarding service managed by W Identity and the proprietary W Social gateway.

If the onboarding service maintains a permanent blacklist of UUIDs corresponding to banned DIDs—a capability perfectly aligned with their stated architecture and data retention policies—the public, decentralized nature of the AT Protocol code is functionally irrelevant to the risk of lifelong exclusion. The protocol merely acts as a standardized data transport mechanism, while the real danger scenarios, privacy risks, and ultimate control over discourse lie entirely within the unpublished operational logic bridging the PDS account lifecycle to the W Identity token validation.

Part 4: Protocol Analysis

A. Observation

The Authenticated Transfer Protocol serves as the foundational data and networking layer for both the Bluesky social network and the W Social platform. A meticulous analysis of the primary

protocol specifications, documentation, and reference implementations confirms several critical properties regarding identity, moderation, and data structures.

User identity within the AT Protocol is defined exclusively by Decentralized Identifiers (DIDs). The network currently supports specific DID methods, primarily did:plc (a method engineered by Bluesky utilizing a centralized registry for high availability and key rotation) and did:web (a method relying on standard DNS and HTTPS infrastructure).² A DID acts as the permanent, immutable cryptographic identifier for an account (e.g., did:plc:ewvi7nxzyoun6zhxrhs64oiz). To provide human-readable usability, users also utilize handles (e.g., alice.wsocial.eu). Handles are entirely mutable DNS TXT records or HTTPS well-known endpoints that resolve bidirectionally to a specific DID.³

User actions, including the creation of posts, replies, likes, follows, and blocks, are stored in cryptographically signed data repositories. These repositories utilize a Merkle Search Tree (MST) structure to ensure data integrity and facilitate efficient synchronization across the network.¹ When a user executes a block action against another user, a new record is generated within the app.bsky.graph.block collection of their personal repository.⁹

Moderation within the AT Protocol relies on a system of delegated authority rather than centralized censorship. It is executed via independent "labelers" (utilizing instances of the Ozone service). These labelers are automated systems or human moderation teams that review content and attach cryptographically signed metadata tags to specific DIDs or individual data records. Client applications and AppViews then filter the visibility of content based on the user's chosen subscriptions to these labelers.¹ Furthermore, the protocol explicitly separates cryptographic identity from physical hosting, guaranteeing account portability. Users possess the architectural right to migrate their DID and their entire MST data repository from one PDS provider to another without losing their social graph or historical content.¹

B. Technical Model

The AT Protocol is philosophically and technically agnostic regarding real-world human identity. It operates strictly on the verification of cryptographic keys and the resolution of internet domains. The base protocol lacks any inherent schema, lexicon, or mechanical enforcement mechanism to prevent a single human entity from generating an infinite number of cryptographic keypairs and registering multiple corresponding DIDs. The complex problem of Sybil resistance—preventing automated bot networks from flooding the system—is deliberately pushed to the application layer or left to the discretion of individual PDS operators.

Because the protocol fundamentally only understands and parses DIDs, blocklists and moderation actions are inherently *DID-bound* and *account-bound*, not person-bound. If User A decides to block User B's DID by creating a block record in their MST, User B can technically bypass this social restriction by generating a new cryptographic keypair, registering a new DID (User C), and continuing to view User A's public data, assuming User A has a publicly accessible repository.

Similarly, when a moderation labeler flags an account as "spam" or "harmful," that cryptographic label is permanently attached to the specific DID, not the biological human operating the keyboard. The protocol's design inherently limits the blast radius of moderation actions to the mathematical boundaries of the compromised identifier.

C. Interpretation

The original claim asserted by the X user—that the *AT Protocol itself* is inherently designed so that blocklists will technically be linked in the background for life to a user's personal ID verification—is fundamentally and demonstrably false. The protocol was designed with the explicit intention of achieving the exact opposite outcome: to abstract identity into user-controlled cryptographic identifiers, thereby preventing any single corporate entity, government, or platform operator from revoking a user's fundamental right to exist and communicate on the underlying network.

However, the AT Protocol is engineered to be highly extensible. The implementation of "Lexicons" allows application developers and network operators to define custom data schemas, advanced record types, and bespoke Remote Procedure Call (RPC) endpoints.¹ While the base open-source protocol does not link real-world government IDs to blocklists, it provides all the necessary architectural hooks—specifically through custom AppViews, proprietary Labelers, and permissioned PDS administration gateways—for a centralized operator to construct a heavily gated overlay network that enforces strict identity requirements on top of the open foundation.

D. Normative and Risk Assessment

The modularity of the AT Protocol allows W Social to legally and technically operate a "permissioned" subset of the broader network. This creates a severe normative risk regarding the illusion of portability. While the AT Protocol guarantees that a user can pack up their DID and migrate their data to a self-hosted PDS, this portability is functionally useless if W Social configures its proprietary AppView—the massive indexing service that aggregates data and serves chronological or algorithmic timelines to end-users—to *only* index and display DIDs that possess a corresponding validated UUID from the W Identity service. Under this scenario, a user could move their DID to an independent server, but they would effectively disappear from the timelines of all W Social users.

Furthermore, because DIDs are designed to be permanent cryptographic anchors, a severe server-side moderation ban on a DID functions as a permanent digital execution of that specific mathematical identity within the indexing scope of the platform. If W Social successfully links a user's singular biological UUID to that targeted DID, the user cannot leverage the open nature of the protocol to generate a new DID and re-enter the W Social ecosystem, as the centralized onboarding gateway will immediately recognize the underlying UUID as blacklisted. The decentralized nature of the AT Protocol offers no protection against centralized ingestion gating.

Part 5: W-Social-Specific Risk Assessment

A. Observation

The operational reality of W Social cannot be understood solely by reading AT Protocol specifications; it must be rigorously assessed through the lens of its official Privacy Notice, legal documentation, and stated architectural philosophy.

W Social officially partners with "W Identity", which operates as a separate legal entity and an independent data controller under the jurisdiction of the EU GDPR.⁶ During the account onboarding phase, the W Identity application securely transmits a UUID, the user's passport country of origin, and their year of birth to the W Social infrastructure.⁶

Crucially, W Social explicitly states in its privacy documentation that it does not collect, process, or store information sufficient to ascertain the actual real-world identity of the user. It does not store images of passports, facial biometrics, or government document hashes on its own servers, relying entirely on the abstract UUID to confirm human legitimacy.⁶

Regarding moderation and governance, W Social maintains a policy to moderate illegal or harmful content, justifying the processing of user data under GDPR Article 6(1)(c) for legal compliance and Article 6(1)(f) for the legitimate interest of security and fraud prevention.⁶ In terms of data retention, the policy states that data is deleted when it is no longer required for its primary purpose. However, W Social explicitly maintains the legal right to "restrict processing"—effectively preserving the data in a locked state rather than deleting it—if data retention is legally required or deemed necessary for the ongoing prevention of illegal activity.⁶

B. Technical Model

The architecture deployed by W Social effectively creates an "information firewall," a sophisticated identity management concept frequently discussed in theoretical presentations on advanced AT Protocol identity deployments.⁷

Within this model, W Identity acts as the authoritative Verifier. It scans the government passport, performs biometric liveness checks, ensures strict global uniqueness to enforce the "one human, one account" paradigm¹¹, and subsequently generates a persistent UUID. W Social acts as the Platform operator. It receives the UUID, provisions the AT Protocol account, maps the UUID to the user's newly generated DID in a secure internal database, and operates the PDS and the indexing AppView.

If a user commits a severe moderation violation—such as posting explicitly illegal content or engaging in massive coordinated harassment—the operational flow for a ban is highly efficient:

1. The ozone moderation labeler or human administrative team flags the DID as suspended or banned.
2. The proprietary W Social backend queries its internal relational database to isolate the UUID associated with the offending DID.

3. The W Social backend updates the status of that specific UUID to "suspended" or "blacklisted" within its access control lists.
4. If the biological user attempts to circumvent the ban by creating a new account, the W Identity application will authenticate their biometrics and pass the exact same underlying UUID back to the W Social onboarding gateway.
5. The gateway recognizes the UUID as blacklisted and summarily rejects the registration payload, completing the cross-account enforcement loop.

C. Interpretation

The specific claim made by the X user regarding "ID verification" is highly nuanced and technically imprecise. There is no primary evidence to support the assertion that there is widespread "hashing or persistent storage of government identity data" occurring directly on the W Social servers.⁶ The highly sensitive passport scans and biometric templates are compartmentalized within the W Identity infrastructure. However, the UUID serves as a persistent, mathematically absolute, and inescapable proxy for the human behind the keyboard.

The claim of "exclusion from discourse" must be interpreted contextually within the architecture of federated networks. A user whose UUID is banned is permanently excluded from creating accounts on the *W Social PDS* and will likely have their existing content purged or delisted from the *W Social AppView*. Importantly, they are not technically excluded from the broader, global AT Protocol network (such as the primary Bluesky network), provided they utilize open-source client software to generate a new DID on an independent third-party PDS. However, if their primary objective is to interact with European policymakers, journalists, or citizens operating exclusively inside the verified W Social walled garden, the UUID ban effectively constitutes an inescapable, lifelong exclusion from that specific sphere of digital discourse.

The claim that this exclusion lasts "for life" hinges entirely on interpretations of GDPR data retention policies. While the GDPR strongly mandates the Right to Erasure (the Right to be Forgotten), digital platforms routinely and successfully claim legal exceptions under the auspices of fraud prevention, platform integrity, and security to retain the cryptographic hashes, phone numbers, or UUIDs of banned users permanently to prevent ban evasion. W Social's privacy policy explicitly and preemptively reserves this exact right via the "restriction of processing" clause.⁶ Therefore, the legal foundation for a lifelong technical ban is firmly established.

D. Normative Risk Scenarios

The W Social architecture, while ostensibly designed to foster a trusted, bot-free environment aligned with European digital sovereignty²¹, introduces profound, demonstrable risks to freedom of discourse, irrespective of the current platform operator's intentions:

1. **The "One Strike" Permanent Exile Risk:** Traditional social media platforms allow for moderation escalation. Users may face temporary suspensions, or they may simply

create an alternative pseudonymous account to start over. Because W Social strictly enforces the "one human, one account" paradigm via a persistent biometric proxy (the UUID), this fundamental safety valve is eliminated. A final server-side administrative ban of the DID translates directly into a final, inescapable ban of the physical human's UUID. If the operator adopts a zero-tolerance moderation policy, a single severe infraction leads to irreversible social death on the platform, with no technical mechanism available for resurrection or evasion.

2. **Vulnerability to State and Institutional Pressure:** By choosing to operate within the strict regulatory environment of the European Union and enforcing a verified real-identity backing, the platform becomes highly susceptible to coordinated legal requests and state pressure. While W Social correctly claims it only holds an abstract UUID, law enforcement agencies possessing the requisite legal authority can compel W Social to surrender the UUID, and subsequently compel the partner entity, W Identity, to deanonymize that UUID and reveal the underlying passport and biometric data. This architecture systematically dismantles the pseudonymity and plausible deniability guarantees inherent in the base AT Protocol.
3. **Subjective Moderation Under the Guise of Legal Compliance:** W Social's official policy states that it moderates "illegal or harmful content." While "illegal" content is strictly defined by statutory law, "harmful" is an inherently nebulous, culturally contingent, and highly subjective categorization. If a future operator, influenced by shifting political climates or commercial pressures, decides to categorize highly critical political speech, whistleblowing, or controversial activism as "harmful," the backend architecture perfectly supports the rapid, automated, and permanent exclusion of those individuals from the digital public square.

Part 6: Final Judgment

Based on a rigorous, exhaustive analysis of the public AT Protocol specifications, the W Social GitHub repository structures, and the official legal, privacy, and architectural documentation, the final assessment of the core claims is as follows:

What is supported? It is demonstrably supported by primary documentation that W Social mandates strict identity verification, enforcing a complex technical framework that inextricably ties a user's digital account to a verified, real-world human identity via a persistent UUID proxy.⁶ It is also factually supported that the AT Protocol utilizes public blocklists and advanced moderation labelers as core networking functionalities.⁸ Furthermore, W Social's official data retention policies legally provision the permanent retention (via the "restriction of processing" mechanism) of user data for fraud prevention and legal compliance, thereby establishing both the technical capability and the legal justification for executing permanent, human-bound bans.⁶

What is contradicted? The foundational premise of the claim—that the *AT Protocol itself* is "designed so that blocklists will technically be linked in the background for life to a user's personal ID verification"—is categorically and completely false. The open-source AT Protocol is

natively designed around transient cryptographic DIDs and is explicitly, mathematically agnostic to physical human identity.¹ The linkage of blocklists and moderation actions to verified human identification tokens is a highly proprietary, centralized, application-layer modification engineered specifically by the operators of W Social. It is a violation of the protocol's ethos, not a feature of its design.

What is a plausible risk?

The secondary assertion that "Anyone who posts something critical once will be excluded from discourse for life" highlights a highly plausible, mathematically sound architectural risk, albeit one that assumes maximum malice and zero-tolerance policies on the part of the operator. Because W Social must map user DIDs to unique, biometrically verified human UUIDs to fulfill its promise of eliminating Sybil attacks and duplicate bot accounts, the platform inherently possesses the exact architectural requirements necessary to execute a lifelong, cross-account ban. If a user's DID is banned for a subjective infraction, the operator can simply update a database flag to blacklist the associated UUID, mechanically ensuring that the specific biological human can never successfully authenticate a new registration payload on the W Social network.

What is currently not decidable?

Because the critical integration logic interfacing the open-source AT Protocol Personal Data Server (pds) with the proprietary W Identity OAuth/UUID gateway is deliberately excluded from the public GitHub forks, it is undecidable precisely how W Social executes cross-account bans on a day-to-day basis. The public codebase does not contain a DID-to-UUID routing table or a `uuid_blacklist` schema. Therefore, whether a standard user-level block or a minor moderation flag triggers an immediate, automated lifelong UUID ban, or whether UUID blacklisting is a manual, highly restricted procedure reserved strictly for severe statutory violations, simply cannot be determined from the available open-source artifacts.

Final Verdict:

The public warning issued by the X user is **fundamentally partially accurate in its assessment of the risk, but technically misattributed and rhetorically overstated**. The user mistakenly blames the open-source architecture of the AT Protocol for invasive identity-tracking mechanisms that are, in reality, centralized, proprietary overlays implemented exclusively by W Social. However, the core underlying warning—that requiring verified human identification tokens as a prerequisite for participation on a federated social network creates the architectural capacity for absolute, lifelong, inescapable exclusion from that platform's sphere of discourse—is technically sound and represents a genuine, severe vulnerability inherent in permissioned decentralized networks.

Final Additional Task: Required Artifacts for Final Determination

To transition the undecidable elements of this technical investigation into conclusive,

mathematically verifiable findings, the acquisition and analysis of the following five critical proprietary artifacts are required:

1. **Onboarding Gateway API Source Code and OIDC Configurations:** The private backend microservice code that governs the OAuth 2.0 or OpenID Connect (OIDC) handshake between the W Identity mobile application and the W Social pds instance. This code is required to ascertain exactly how the JWT payload containing the UUID is ingested, validated, and persistently stored alongside the cryptographic DID generation routines.
2. **Production Database Schema Migrations:** The actual SQL or PostgreSQL Data Definition Language (DDL) scripts utilized in W Social's production environments. Access to these schemas is necessary to definitively verify the existence of relational mapping tables (e.g., did_to_uuid_mapping) and perimeter access control tables (e.g., uuid_global_blacklist).
3. **Moderation Enforcement Logic and Ozone Webhooks:** The private backend scripts and event-driven architectures that execute when a system administrator bans an account utilizing the W Social ozone interface. This logic is required to determine whether a standard DID suspension automatically triggers a cascading status change to the underlying UUID, resulting in an automated cross-account ban.
4. **Internal Moderation Policy and Escalation Guidelines:** The proprietary, internal policy documents distributed to W Social moderation teams. These documents are needed to define the specific, actionable thresholds of "harmful content" that distinguish between a temporary, account-level suspension and a permanent, UUID-level exile from the platform.
5. **W Identity Data Flow and Cryptographic Documentation:** Technical whitepapers or comprehensive API documentation from the separate legal entity, "W Identity." This documentation is required to detail the exact cryptographic generation sequence of the UUID, the rate limits imposed on generating new UUIDs, and the lifecycle management of the UUID in relation to physical passport renewals, expirations, or biometric shifts.

Referenzen

1. The AT Protocol - Bluesky API, Zugriff am Mai 7, 2026, <https://docs.bsky.app/docs/advanced-guides/atproto>
2. Bluesky: Decentralized Social Media Insights | PDF | World Wide Web - Scribd, Zugriff am Mai 7, 2026, <https://www.scribd.com/document/782568482/Bluesky-and-the-AT-Protocol-Usable-Decentralized-Social-Media>
3. Decentralized Identity - AT Protocol - Mintlify, Zugriff am Mai 7, 2026, <https://mintlify.com/bluesky-social/atproto/concepts/identity>
4. The Authenticated Transfer Protocol ("AT Protocol" or "atproto") is a network protocol for building open social web applications., Zugriff am Mai 7, 2026, <https://atproto.com/specs/atp>
5. W Social EU - GitHub, Zugriff am Mai 7, 2026, <https://github.com/w-social-eu>

6. Privacy notice - W Social, Zugriff am Mai 7, 2026, <https://wsocial.eu/public/privacy-notice>
7. My Bookmarks — ATmosphereConf 2026, Zugriff am Mai 7, 2026, <https://atmosphereconf.org/bookmarks>
8. mary-ext/bluesky-labeler-scraping - GitHub, Zugriff am Mai 7, 2026, <https://github.com/mary-ext/bluesky-labeler-scraping>
9. Just a warning for vtubers moving to BlueSky, your blocked list is public - Reddit, Zugriff am Mai 7, 2026, https://www.reddit.com/r/VirtualYoutubers/comments/1g6ijv7/just_a_warning_for_vtubers_moving_to_bluesky_your/
10. The AT Protocol - Bluesky, Zugriff am Mai 7, 2026, <https://bsky.social/about/blog/10-18-2022-the-at-protocol>
11. The Case for World ID: Social Media : r/worldid - Reddit, Zugriff am Mai 7, 2026, https://www.reddit.com/r/worldid/comments/1r8s8t4/the_case_for_world_id_social_media/
12. Pi Network Introduces Advanced Face and Palm Scan Verification for Safer Crypto Transactions | MEXC News, Zugriff am Mai 7, 2026, <https://www.mexc.com/en-NG/news/630451>
13. Sam Altman's World and Coinbase roll out toolkit to distinguish human-backed AI agents from bots - Crypto Briefing, Zugriff am Mai 7, 2026, <https://cryptobriefing.com/ai-agent-verification-world-launch/>
14. GitHub - bluesky-social/atproto: Social networking technology created by Bluesky, Zugriff am Mai 7, 2026, <https://github.com/bluesky-social/atproto>
15. Bluesky labellers - Blue Mackuba, Zugriff am Mai 7, 2026, <https://blue.mackuba.eu/labellers/>
16. DNS Account Handles, A Whitepaper - IETF, Zugriff am Mai 7, 2026, <https://www.ietf.org/archive/id/draft-hallambaker-any-00.html>
17. (PDF) OIDC²: Open Identity Certification with OpenID Connect - ResearchGate, Zugriff am Mai 7, 2026, https://www.researchgate.net/publication/378896695_OIDC2_Open_Identity_Certification_With_OpenID_Connect
18. DID - AT Protocol, Zugriff am Mai 7, 2026, <https://atproto.com/specs/did>
19. Identity - AT Protocol, Zugriff am Mai 7, 2026, <https://atproto.com/guides/identity>
20. AT Protocol - Wikipedia, Zugriff am Mai 7, 2026, https://en.wikipedia.org/wiki/AT_Protocol
21. European alternatives to Reddit are about to Launch in 2026. : r/BuyFromEU, Zugriff am Mai 7, 2026, https://www.reddit.com/r/BuyFromEU/comments/1qlmjvc/european_alternatives_to_reddit_are_about_to/