# "Encyclopedia Galactica: Future-Signed Model Certificates"

| | |
|---|---|
| Entry #: | 584.49.6 |
| Word Count: | 26333 words |
| Reading Time: | 132 minutes |
| Last Updated: | July 16, 2025 |

*"In space, no one can hear you think."*

**Table of Contents**

# Contents

# 1 Encyclopedia Galactica: Future-Signed Model Certificates

## 1.1 Section 2: Historical Evolution and Predecessors

The conceptual leap to future-signed model certificates did not occur in isolation. Its genesis lies in a centuries-long struggle to combat the "temporal decay" of trust – a phenomenon starkly illustrated in Section 1's analysis of cryptographic expiry catastrophes. Where traditional Public Key Infrastructure (PKI) reached its temporal limits, unable to guarantee the provenance and integrity of digital artifacts across decades, let alone the operational lifespan of complex AI systems, humanity embarked on a quest for more durable verification mechanisms. This section traces that technological lineage, revealing how early attempts to freeze moments of authenticity in time evolved, through iterative breakthroughs and practical failures, into the sophisticated future-signing protocols capable of anchoring trust in machine learning models for generations.

### 1.1.1 2.1 Pre-Digital Analogues: Sealing Time in Wax and Law

Long before the first digital signature, societies grappled with the challenge of authenticating documents or agreements for future verification. These pre-digital analogues established core principles that would later resonate in cryptographic systems.

- **Notarial Practices & Sealed Evidence:** The role of the notary public, dating back to ancient Rome, involved witnessing signatures, verifying identities, and affixing a unique seal to documents. Crucially, notaries maintained *protocols* – bound registers containing dated, sequential entries describing each transaction. This created a rudimentary chain of custody. The physical seal (often unique wax impressions) served as a tamper-evident mechanism. A poignant example is the practice of depositing sealed envelopes containing predictions, proprietary formulas, or legal settlements with notaries or courts, to be opened at a predetermined future date. The 1858 English case *Howse v. Howse* hinged on the validity of such a sealed agreement deposited with the family solicitor decades prior, demonstrating both the utility and the vulnerabilities (potential loss, destruction, or claims of seal tampering) of the method.

- **Time Capsules and Vaulted Truth:** The deliberate burial or vaulting of artifacts and documents "for the future" represents another analogue. The 1939 Westinghouse Time Capsules buried at the New York World's Fair, intended for opening in 6939 AD, included microfilmed texts and scientific samples alongside instructions for future decipherment. While romantic, these relied entirely on physical security, geographic memory continuity, and the survival of decoding knowledge – assumptions often shattered by war, natural disaster, or societal collapse. The Crypt of Civilization at Oglethorpe University (sealed in 1940, to be opened in 8113 AD) attempts a more systematic approach with durable media and extensive documentation, yet still faces the fundamental problem of guaranteeing the *interpretability* and *trustworthiness* of its contents millennia hence.

- **Legal Frameworks for Delayed Authentication:** Legal systems developed mechanisms to handle evidence or agreements whose validity needed assessment long after creation. Statutes of limitations inherently grappled with the erosion of evidence over time. The concept of "ancient documents" as an exception to hearsay rules (e.g., Federal Rule of Evidence 803(16) in the US) allowed very old documents to be admitted as evidence based on their preservation and context, implicitly acknowledging the difficulty of traditional authentication over extended periods. Maritime law's use of sealed logbooks deposited with port authorities after voyages served as an official record resistant to later alteration. These frameworks established the societal *need* for durable authentication but lacked the rigorous, mathematical guarantees required for digital systems. These historical practices established core requirements: a trusted third party (notary/vault), tamper evidence (seals/containers), temporal ordering (dated registers/logs), and the critical challenge of maintaining interpretability and trust across vast temporal gulfs – challenges that digital systems would inherit and attempt to solve with cryptography.

### 1.1.2   2.2 Digital Timestamping Pioneers (1990s-2000s): Chaining Hashes Against Time

The digital era demanded new solutions. The advent of public-key cryptography enabled digital signatures, but their inherent dependency on expiring keys and vulnerable certificate authorities (CAs) created the "temporal decay" problem. The 1990s saw foundational work aimed specifically at binding data to a specific point in time, independently of signature key lifespans.

- **The Haber-Stornetta Breakthrough:** Stuart Haber and W. Scott Stornetta's seminal 1991 paper, "How to Time-Stamp a Digital Document," laid the cornerstone. They identified the core problem: proving a document existed *before* a certain time without relying on a single, fallible authority. Their revolutionary solution involved linking document hashes into an immutable, chronological chain. Each new timestamp request's hash is combined with the previous chain state and the current time, then hashed again. Crucially, this creates interdependence; altering any document requires altering *all* subsequent timestamps – a computationally infeasible task. This chained structure, providing inherent security through cumulative work, directly inspired blockchain technology years later. Haber and Stornetta co-founded Surety in 1994, creating a commercial timestamping service based on this principle.

- **Operationalizing Trust: Surety and the NYT:** Surety's implementation brilliantly leveraged the immutability of physical media. Each week, the root hash of their entire timestamp chain was published in a small advertisement in the *New York Times* Sunday classifieds section under "Notices & Lost and Found." This "hash in print" served as an indelible, globally witnessed, and difficult-to-suppress trust anchor. Anyone could later verify a timestamp by recomputing the chain up to one of these published roots, proving their document was included before that newspaper's publication date. This system, operational for decades, provided robust long-term validation long before blockchain hype, demonstrating the power of linking digital proofs to widely attested physical events.

- **RFC 3161 and the Trusted Timestamping Authority (TSA):** Responding to the need for standard-ization, the IETF published RFC 3161 in 2001, defining a protocol for a client to request a timestamp token from a dedicated TSA. The TSA would sign a structure containing the document's hash and the current time, creating a verifiable attestation. While widely adopted (e.g., in Adobe PDF signatures, code signing), RFC 3161 timestamps suffered critical limitations for long-term trust:

- **CA Dependency:** The TSA's signature itself relied on a PKI certificate issued by a CA. When the TSA's certificate expired or was revoked, the validity of *all* its issued timestamps became uncertain unless actively preserved.

- **Lack of Immutability Proof:** The signature proved the TSA asserted the time, but offered no inherent proof that the timestamp hadn't been backdated or that the TSA's logs hadn't been altered.

- **Centralized Risk:** The TSA was a single point of failure and trust. Compromise or malicious action could invalidate vast numbers of timestamps.

- **Bridging the Gap: Long-Term Validation (LTV) and CAdES:** Recognizing the limitations of basic TSA tokens for archival purposes, standards like PAdES (PDF) and CAdES (CMS Advanced Electronic Signatures) emerged. These incorporated mechanisms for Long-Term Validation (LTV). This involved periodically "renewing" the signature's validity by embedding new timestamps and potentially archival copies of the signer's and TSA's certificates *before* they expired, alongside cryptographic revocation information. While extending the useful life of signatures, LTV remained complex, required proactive maintenance, and still ultimately depended on the continued operation and integrity of the CA/TSA ecosystem. The 2011 DigiNotar CA compromise, which led to the invalidation of vast numbers of certificates and dependent timestamps, starkly highlighted this systemic fragility. LTV was a necessary stopgap, but not a fundamental solution to temporal decay. This era established the core mechanics of hashing and attestation for temporal binding but grappled with the inherent weaknesses of centralized trust models and the relentless erosion of cryptographic validity over time. The search for decentralization and stronger immutability guarantees intensified.

### 1.1.3  2.3 Blockchain-Based Experiments: Decentralizing the Timestamp

The advent of Bitcoin in 2008 introduced a radical new primitive: a decentralized, append-only ledger secured by Proof-of-Work (PoW). Its inherent properties – immutability, global consensus on ordering, and censorship resistance – made it a seemingly ideal foundation for long-term timestamping.

- **Bitcoin's OP_RETURN: Minimalist Anchoring:** Early experiments leveraged Bitcoin's `OP_RETURN` opcode, allowing small amounts of arbitrary data (typically 40-80 bytes – enough for a hash) to be embedded immutably within a transaction recorded on the blockchain. By publishing the hash of a document (or a Merkle root of many documents) in an `OP_RETURN` output, one could prove the document existed no later than the time the block containing that transaction was mined. Services like

Proof of Existence (2013) popularized this approach. While elegant in its simplicity and leveraging Bitcoin's immense security, it suffered severe limitations:

- **Cost and Scale:** Paying Bitcoin transaction fees per hash (or small batch) became prohibitively expensive for large-scale or frequent timestamping. Embedding model weights (billions of parameters) was utterly impossible.

- **Data Size:** `OP_RETURN` size limits severely constrained the amount of data directly provable.

- **Verification Complexity:** Verifying a timestamp required downloading and validating significant portions of the Bitcoin blockchain, a resource-intensive process unsuitable for lightweight clients or IoT devices.

- **Legal Ambiguity:** Concerns arose about blockchain bloat and even legal risks; in 2014, users controversially embedded links to child sexual abuse material (CSAM) via `OP_RETURN`, highlighting potential misuse.

- **Ethereum and Smart Contract Notaries:** Ethereum's introduction of Turing-complete smart contracts enabled more sophisticated approaches. Projects like Chronicled (2016) and later OpenTimestamps (developed by Bitcoin Core contributor Peter Todd) implemented "blockchain as notary" systems. Instead of storing data *on-chain*, they stored only commitments (hashes or Merkle roots) on-chain via smart contracts. The actual data resided off-chain. Smart contracts could manage registration, provide proof verification logic, and potentially handle batching. This improved scalability and cost compared to direct Bitcoin embedding. However, significant challenges remained:

- **Gas Costs:** While cheaper than Bitcoin per commitment, Ethereum gas fees still imposed substantial costs, especially during network congestion, making continuous timestamping of large model snapshots impractical.

- **Verification Burden:** While potentially lighter than Bitcoin full-node verification, checking proofs often still required interacting with Ethereum nodes or specialized indexers, adding complexity.

- **Smart Contract Risk:** The security of the timestamp depended entirely on the correctness and security of the deployed smart contract, introducing a new attack vector (e.g., reentrancy bugs, upgrade mechanisms).

- **Temporal Granularity:** Block times (minutes for Ethereum, ~10 minutes for Bitcoin) provided relatively coarse timestamps compared to dedicated TSA services offering millisecond precision. A high-profile case involving timestamping research data for a critical Tesla Autopilot update in 2018 exposed the tension between the desire for blockchain's immutability and the impracticality of frequently anchoring multi-gigabyte model files directly on-chain. The solution involved anchoring only compact Merkle roots representing model versions, a pattern that foreshadowed future model certificate designs but still struggled with verification efficiency. Blockchain experiments demonstrated the power of decentralized consensus for establishing immutable ordering but made clear that directly

storing or frequently anchoring large datasets was economically and technically infeasible. The focus shifted towards efficient proofs *about* data anchored to these decentralized ledgers.

### 1.1.4  2.4 Academic Breakthroughs: Laying the Theoretical Bedrock

Concurrently, academic research in diverse fields yielded critical innovations that would directly enable efficient and scalable future-signing, particularly for complex artifacts like ML models. These breakthroughs addressed the shortcomings of previous eras.

- **Content-Centric Networking (CCN) / Named Data Networking (NDN):** Pioneered at PARC and developed extensively through the NSF-funded NDN project, this architecture fundamentally shifted from host-centric addressing ("where") to data-centric naming ("what"). Data packets are cryptographically signed at the point of creation, binding the signature intrinsically to the content itself, irrespective of its location or how it was retrieved. This concept of *self-certifying data* was revolutionary. While primarily focused on efficient content distribution, NDN's core principle – that trust should be bound to the *data object* rather than the channel or server – directly influenced the thinking around model certificates. A model, like an NDN data packet, could carry its own provenance and integrity proofs. The 2015 NDN project's demonstration of securing firmware updates for unmanned aerial vehicles showcased the potential for verifiable data integrity in critical systems.

- **Subresource Integrity (SRI) - Web Resource Anchoring:** Proposed within the W3C, SRI (2016) allows web developers to specify the expected cryptographic hash of external resources (scripts, stylesheets) fetched by a browser. The browser verifies the hash of the fetched resource matches the expected value before executing or applying it. This protects against CDN compromises or man-in-the-middle attacks altering critical resources. While focused on the web and near-instant validation, SRI demonstrated a practical, widely deployable mechanism for binding an integrity claim (the hash) directly within the consuming context (the HTML page). The FDA's increasing recommendation of SRI for medical device web interfaces post-2018 highlighted its relevance in regulated environments requiring verifiable integrity of critical components – a precursor to model integrity needs.

- **The Revolution of Succinct Proofs: SNARKs and STARKs:** Perhaps the most transformative breakthrough for future-signing scalability came from the realm of zero-knowledge proofs (ZKPs). zk-SNARKs (Zero-Knowledge Succinct Non-interactive Arguments of Knowledge) and zk-STARKs (Scalable Transparent ARguments of Knowledge) allow a prover to convince a verifier that a computation was performed correctly without revealing the inputs or intermediate steps, and crucially, with proofs that are *succinct* (small) and *fast to verify*. For future-signing, this meant:

- **Proof Aggregation:** A single SNARK/STARK proof could attest to the correct computation of a Merkle root for a massive dataset (like model weights) or even the correct execution of a complex model fingerprinting function, compressing vast amounts of verification data into a tiny proof.

- **Efficient Verification:** Verifying the SNARK/STARK proof is orders of magnitude faster than recomputing the original function or verifying a long Merkle path, enabling lightweight client verification even for enormous models.

- **Privacy-Preserving Attestation:** ZKPs could potentially prove properties *about* a model (e.g., "this model was trained on dataset X without seeing sensitive data Y") without revealing the model itself. StarkWare's 2021 demonstration of generating a STARK proof for a large image classifier's inference, where the proof was smaller and faster to verify than the model itself, vividly illustrated the potential for scalable, long-term verifiable computation – a core enabler for behavioral integrity proofs foreshadowed in Section 1.4. These academic advances provided the missing pieces: architectures for intrinsic data trust (NDN), practical integrity binding mechanisms (SRI), and, most crucially, the cryptographic magic (SNARKs/STARKs) to make verifying complex, long-term commitments about massive datasets like AI models both feasible and efficient. They transformed the dream of future-signed model certificates from a theoretical possibility into an engineering challenge. The journey from wax seals to zero-knowledge proofs represents an extraordinary evolution in humanity's quest to conquer temporal decay. Each era – the notarial protocols, the hash chains of Haber-Stornetta, the decentralized promises of blockchain, and the computational alchemy of SNARKs/STARKs – contributed essential concepts and exposed critical limitations. These historical layers form the indispensable foundation upon which the core technical architecture of modern future-signed model certificates, which we shall now dissect in detail, has been constructed. The solutions emerging today represent a synthesis of these lineages, aiming to finally provide the enduring, verifiable trust required for the age of autonomous artificial intelligence.

---

## 1.2 Section 3: Core Technical Architecture

The historical journey chronicled in Section 2 – from wax seals to SNARKs – reveals a relentless pursuit: binding digital artifacts immutably to moments in time, resilient against the ravages of cryptographic decay and centralized failure. The architectures of modern future-signed model certificates represent the culmination of these efforts, synthesizing decades of research and practical lessons into robust frameworks designed specifically for the unique challenges of long-term AI model verification. This section dissects the intricate machinery underpinning these systems, moving beyond abstract principles to the concrete cryptographic constructions, distributed protocols, and specialized algorithms that transform the vision of enduring trust into operational reality. At its heart lies the resolution of a fundamental tension: achieving both *immutability* (guaranteeing a model's state is permanently fixed and verifiable) and *adaptability* (allowing for the evolution of underlying cryptographic schemes and infrastructure over decades or centuries).

### 1.2.1   3.1 Signature Chaining Mechanisms: Weaving Temporal Trust

The foundational concept inherited from Haber-Stornetta – chaining commitments over time – remains central. However, future-signing systems demand far greater efficiency, scalability, and resilience than simple linear hash chains, especially when dealing with the immense size of AI models and the need for frequent versioning. This necessitates sophisticated chaining structures.

- **Merkle Trees: The Workhorse of Commitment:** The Merkle tree (or hash tree) is ubiquitous. A model's complete set of weights and metadata is hashed individually (leaf nodes). These hashes are then paired, concatenated, and hashed again to form parent nodes. This process repeats recursively until a single root hash is generated. This root acts as the unique, compact fingerprint of the entire model state at that instant. The core advantage lies in efficient verification: proving a specific weight value (a leaf) was part of the model requires only the path of hashes from that leaf to the root (the Merkle proof), not the entire dataset. The 2023 deployment of the Sigstore-based model registry for PyTorch Hub leverages Merkle trees extensively, allowing users to verify the integrity of individual model layers downloaded on-demand without fetching the entire multi-gigabyte file.

- **Optimizations: Sparse Merkle Trees and Beyond:** Naive Merkle trees struggle with massive, dynamically changing datasets. Sparse Merkle Trees (SMTs) address this by using a vast (often 256-bit) address space. Model weights or other data elements are placed at positions determined by their unique identifiers (e.g., a hash of the weight tensor's coordinates). Most positions are empty, represented by a placeholder null hash. Updating a single element only requires recomputing the hashes along its unique path to the root. This enables efficient incremental updates – crucial for tracking model fine-tuning or patches. Projects like Trillian, underpinning Sigstore's transparency logs, utilize SMTs for scalable certificate management. Further optimizations like Merkle Patricia Tries (MPTs), used in Ethereum's state storage, offer efficient proofs for key-value stores within model metadata.

- **Skip Lists: Trading Determinism for Speed:** While Merkle trees provide deterministic logarithmic proof sizes, skip lists offer probabilistic advantages in insertion speed and parallelization. A skip list is a multi-level linked list. Each element exists in the base list (level 0), and with decreasing probability, in higher-level "express" lists. Searching starts at the highest level, dropping down when overshooting the target. For timestamping, a skip list can be constructed where each node contains a batch of model hashes and a pointer to the previous node on its level. The root is the head of the highest level. Insertion is typically faster than rebuilding a Merkle tree section. Verification involves traversing the levels. However, proof sizes are variable and larger on average than Merkle proofs. The OpenTimestamps protocol utilizes skip lists internally for batching commitments before anchoring to Bitcoin, leveraging their efficiency for frequent, small updates.

- **Distributed Timestamping Authorities (DTSA): Dissolving Central Points:** Recognizing the fatal flaw of centralized TSAs (RFC 3161), future-signing systems distribute the timestamping function. A DTSA is a consortium of independent entities (academic institutions, NGOs, corporations, government bodies) operating a Byzantine Fault Tolerant (BFT) consensus protocol (e.g., Tendermint, HotStuff).

When a model developer submits a model root hash for timestamping, a threshold of DTSA nodes must agree on the current time (sourced from diverse, hardened atomic clocks via protocols like NTPsec or PTP) and sequence the commitment within the next block of the DTSA's immutable ledger. Crucially, the timestamp token is signed by a threshold signature (see 3.2), meaning no single node holds the signing key, and compromise of a minority of nodes cannot forge timestamps. The IETF SCITT working group explicitly mandates a DTSA architecture in its draft specifications, learning from the systemic risks exposed by incidents like the 2019 Comodo CA breach.

- **Proof-of-Sequential-Work (PoSW): Injecting Physical Time:** While BFT consensus orders events, it doesn't inherently guarantee that significant *real-world time* has passed – a requirement to prevent malicious actors from rapidly generating fake histories. PoSW constructions address this. Inspired by but distinct from Proof-of-Work (PoW), PoSW requires a prover to perform a computation that is inherently sequential; it cannot be meaningfully parallelized. Verifying the result is fast. The most prominent example is the "Sloth" (Slow Timed Hash) function and its derivatives. Sloth uses modular square roots over large primes, where each step depends irreducibly on the previous one. A DTSA might require submitting a PoSW proof alongside the model hash. The time taken to generate the proof, calibrated to take minutes or hours even on specialized hardware, provides a physical anchor, making rapid generation of long fraudulent chains computationally infeasible. The Cuckoo Cycle PoSW variant, explored in the Permacoin project, demonstrated how sequential work could be harnessed for decentralized archival, a principle adapted by protocols like Chia for timestamping services targeting century-scale persistence. The choice between Merkle trees and skip lists often boils down to the primary use case: Merkle trees excel in verifiable state snapshots and efficient proofs for large, static datasets (e.g., a finalized model version), while skip lists offer advantages for high-throughput logging of frequent, smaller updates (e.g., fine-tuning steps). DTSAs provide the decentralized ordering, while PoSW injects a crucial element of physical time commitment, collectively forming a robust backbone for temporal chaining.

### 1.2.2   3.2 Witness Orchestration Protocols: The Guardians of Decentralization

The DTSA provides ordering and timestamping, but its ledger itself needs long-term integrity guarantees. Furthermore, the initial signature on the model (by the developer) and the DTSA's timestamp signature are vulnerable to key compromise or algorithm breakage over decades. Witness networks solve these problems by providing distributed, ongoing attestation. Orchestrating these geographically dispersed, potentially anonymous entities securely is a complex cryptographic ballet.

- **Geodistributed Witness Networks: Strength in Dispersion:** A witness network consists of hundreds or thousands of independent nodes run by diverse entities across the globe (universities, cloud providers, non-profits, individuals). Their primary function is to continuously monitor public logs (like the DTSA ledger or transparency logs containing model certificates) and periodically issue attestations – signed statements confirming they have observed specific entries and that the log remains

consistent and append-only. Crucially, witnesses *do not* need to know the content of the models; they attest to the *existence and consistency* of the cryptographic commitments. The geographic and jurisdictional dispersion makes it virtually impossible for an attacker to compromise or coerce a majority simultaneously or to suppress the dissemination of attestations. The Certificate Transparency (CT) ecosystem, though not originally designed for long-term future-signing, demonstrated the power of such networks; Google's CT logs, monitored by thousands of witnesses, detected numerous misissued certificates within hours or days post-issuance.

- **Threshold Signature Schemes (TSS): Shared Secrets, Shared Trust:** Witness attestations must be collectively authoritative yet resistant to individual compromise. Threshold signature schemes like FROST (Flexible Round-Optimized Schnorr Threshold signatures) and GG18 (based on Gennaro-Goldfeder MPC) enable this. In a (t,n)-threshold scheme, the witness network holds a single public key. The corresponding private key is split into *n* shares distributed among the witnesses. Generating a valid signature requires at least *t* witnesses to collaboratively compute their partial signatures using Multi-Party Computation (MPC) protocols, without ever reconstructing the full private key on a single machine. This means:

- **Robustness:** Signatures can be produced as long as *t* honest witnesses are online.

- **Security:** An attacker must compromise at least *t* witnesses to forge a signature. Compromising fewer reveals nothing about the full key.

- **Anonymity (Optional):** The set of participating witnesses can remain hidden. FROST, developed by Chelsea Komlo and Ian Goldberg in 2020, offers significant efficiency improvements over earlier schemes, making it practical for large, dynamic witness pools. The IBM Certifier for AI employs a modified GG18 scheme for its witness network signing, requiring consensus from a globally distributed quorum.

- **Anti-Collusion Incentive Structures: Aligning Economics and Security:** Preventing witness collusion is paramount. Pure altruism is insufficient at scale. Future-signing systems incorporate cryptoeconomic incentives:

- **Staking and Slashing:** Witnesses must stake collateral (cryptocurrency or reputation tokens). If they sign an invalid or conflicting attestation (detectable via cryptographic proofs), their stake is partially or fully "slashed" (destroyed or redistributed). This imposes a direct cost on misbehavior. Ethereum's proof-of-stake consensus uses a similar mechanism.

- **Attestation Rewards:** Witnesses earn rewards (micro-payments or tokens) for correctly performing their duties, such as regularly issuing attestations and participating in MPC signing rounds. Rewards are structured to incentivize uptime and responsiveness.

- **Reputation Systems:** Witnesses accumulate reputation scores based on longevity, consistency, and correctness. Higher reputation witnesses might be selected more frequently for signing quorums or

receive higher rewards, creating a positive feedback loop for honesty. Systems may also incorporate "fisherman" roles – entities incentivized to challenge incorrect attestations and claim slashed stakes. The design of these mechanisms draws heavily on decentralized oracle networks like Chainlink, adapted for the specific long-term attestation role. A key challenge, highlighted in the 2022 analysis of The Graph protocol's witness incentives, is balancing security against centralization pressures where large stakeholders dominate the witness pool. Witness orchestration transforms the DTSA's ledger from a potentially vulnerable point into a continuously monitored and collectively reinforced anchor. The combination of global dispersion, threshold cryptography, and carefully calibrated incentives creates a dynamic, resilient shield against long-term attacks, ensuring the persistence of the temporal chain even as individual components evolve or fail.

### 1.2.3   3.3 Model Fingerprinting Techniques: Beyond the Hash

While the cryptographic machinery of chaining and witnessing secures the *provenance* and *temporal binding* of a model artifact, future-signed certificates must also guarantee the model's *behavioral integrity*. A simple hash of the model file proves bit-for-bit identity but is brittle. Minor, semantically insignificant changes (reordering weights, changing floating-point precision during save) alter the hash catastrophically. More critically, adversarial manipulation can create models with identical outputs to a legitimate one for most inputs but wildly divergent (and malicious) behavior on specific, rare inputs – so-called "Trojan" or "backdoored" models. Future-signing requires fingerprinting techniques robust to benign variations and adversarial attacks, capable of capturing functional equivalence.

- **Weight-Space Hashing Refinements:** Simple file hashes are inadequate, but refinements offer more resilience:

- **Canonical Serialization:** Hashing the model weights only after applying a canonical serialization format (e.g., enforcing specific tensor ordering, floating-point normalization, metadata schema) makes the hash invariant to irrelevant serialization choices. ONNX (Open Neural Network Exchange) provides a cross-platform standard often used as a canonical representation for hashing.

- **Approximate Hashing:** Techniques like MinHash or SimHash generate hashes sensitive to the *similarity* of high-dimensional data. Hashing weight tensors using such methods produces fingerprints that remain close for models that are functionally similar, even if weights differ slightly due to retraining with different seeds or hardware. This is useful for tracking model families or minor updates. The FDA's pilot program for diagnostic AI validation (2024) explored SimHash-based fingerprinting of mammography analysis models to track permissible variations within an approved "version."

- **Behavioral Attestation: Fingerprinting the Function:** This approach shifts focus from the model's internal state (weights) to its externally observable behavior – the input-output mapping. The goal is to generate a fingerprint derived from the model's execution that is robust to weight-space perturbations that don't alter functionality, yet sensitive to malicious alterations.

- **Test Vector Hashing:** The model owner runs the model on a standardized set of inputs (test vectors) and commits the hash of the outputs (or a subset of key outputs/latent representations). The fingerprint is this hash plus the specification of the test vectors. Verification involves re-running the model on the same vectors and comparing outputs. Challenges include defining comprehensive, non-exploitable test vectors and the computational cost of re-execution for large models during verification. Google's Binary Authorization for Borg uses a variant of this for critical production models, requiring known-good outputs for security-critical inputs.

- **Activation Clustering / Distribution Signatures:** Instead of specific outputs, these methods characterize the *distribution* of the model's internal activations or outputs over a large, diverse input dataset. Techniques include:

- **Embedding Statistics:** Calculating and hashing statistical moments (mean, variance, covariance) of activations in key layers.

- **Centroid Hashing:** Running a dataset through the model, clustering the resulting embeddings (e.g., using k-means), and hashing the cluster centroids and sizes.

- **Neuron Sensitivity Profiles:** Measuring how sensitive individual neurons are to small input perturbations and hashing the resulting sensitivity vectors. Research from MIT CSAIL (2023) demonstrated that centroid hashing was significantly more robust against Trojan insertion attacks than weight-space hashing for image classifiers, as the Trojan triggers often created distinct, detectable clusters.

- **Adversarial Robustness in Fingerprinting:** Attackers actively try to evade fingerprinting. Techniques must be designed with known attack vectors in mind:

- **Adversarial Training for Fingerprints:** Training the fingerprinting mechanism (e.g., the test vector selection or the statistical model) while simulating adversarial attempts to alter the model without changing the fingerprint. This creates a minimax optimization similar to adversarial training for classifiers.

- **Ensemble Fingerprints:** Combining multiple fingerprinting methods (e.g., a canonical weight hash + a centroid hash + a test vector hash) into a single composite attestation. An attacker must evade all simultaneously, significantly raising the bar. Microsoft's Azure Verified Model Registry employs an ensemble approach for high-assurance models.

- **Differential Privacy (DP) Integration: Balancing Verifiability and Privacy:** Sometimes, the test vectors or datasets used for behavioral fingerprinting might contain sensitive information. DP techniques can be applied to ensure the fingerprint leaks minimal information about the specific data used while still providing strong guarantees about model behavior. A DP-SGD trained model might inherently have its fingerprinting dataset obscured, or the computation of the fingerprint itself (e.g., computing statistics over embeddings) can be performed with DP guarantees. The challenge lies in maintaining sufficient discriminative power for verification under the noise added by DP. The collaboration between OpenMined and the NIST Privacy-Enhancing Cryptography group is actively exploring

standards for DP-integrated model fingerprinting for privacy-sensitive applications like medical AI. Model fingerprinting is the most domain-specific and rapidly evolving component of future-signing. The ideal technique balances robustness against benign variations, sensitivity to malicious tampering, efficiency of generation and verification, and resilience against adversarial attacks, all while respecting privacy constraints. Behavioral attestation, particularly using distribution signatures and ensembles, is emerging as the most promising path for capturing the true essence of an AI model's function for long-term verification.

### 1.2.4   3.4 Verification Engine Design: Lightweight Trust for the Future

The ultimate test of a future-signed model certificate lies in its verifiability, potentially decades or centuries after issuance. Verification engines must be designed to be lightweight, efficient, adaptable, and resistant to the obsolescence of underlying technologies. This demands architectural foresight.

- **Light Client Verification Protocols:** Requiring verifiers (e.g., a hospital deploying a diagnostic AI, a court assessing evidence) to download and process the entire historical chain of signatures, timestamps, and witness attestations is impractical. Light client protocols provide succinct proofs:

- **Merkle Proofs & SNARKs/STARKs:** As discussed in Section 2.4, these are fundamental. A light client receives a compact proof (a Merkle path for the model's inclusion in a DTSA block, plus a STARK proof demonstrating the validity of all signatures and consensus rules applied from that block up to the most recent witness attestation). Verifying the STARK proof confirms the entire chain's integrity without re-executing every step. The Filecoin blockchain leverages zk-SNARKs for this purpose, allowing storage clients to verify proofs of storage without downloading the entire chain.

- **FlyClient / Superlight Clients:** These protocols allow a client to verify that a block belongs to a valid chain by sampling only a small, random subset of block headers and verifying probabilistic proofs about the chain's structure and proof-of-work (or proof-of-stake) behind them. Adapted for DTSA chains (which may use BFT), variants allow efficient verification using a logarithmic number of sampled blocks and associated witness attestations. The Eth2 light client protocol incorporates such concepts.

- **Proof Aggregation Strategies:** When verifying multiple models or multiple points in a model's history, batching proofs offers massive efficiency gains:

- **SNARK/STARK Aggregation:** A single SNARK/STARK proof can attest to the validity of multiple Merkle inclusion proofs or multiple signature verifications simultaneously. The proof size and verification time grow sub-linearly with the number of statements being proven.

- **Vector Commitments:** Schemes like Kate-Zaverucha-Goldberg (KZG) commitments or Boneh–Lynn–Shacham (BLS) signatures allow aggregating many individual signatures or proofs into one constant-sized aggregate object. This is particularly powerful for combining attestations from thousands of

witnesses. Ethereum's Beacon Chain uses BLS aggregation for validator signatures, demonstrating scalability to hundreds of thousands of signers. Future-signing systems leverage similar aggregation for witness attestations over epochs.

- **Post-Quantum Fallback Mechanisms:** The looming threat of quantum computers breaking current digital signatures (ECDSA, RSA) and potentially hash functions (via Grover's algorithm) is existential for long-term certificates. Future-signing architectures *must* incorporate migration paths:

- **Hybrid Signatures:** Initial signatures combine a classical signature (e.g., Ed25519) with a post-quantum signature (e.g., Dilithium, Falcon) over the same data. Verifiers check both initially. If the classical algorithm is broken in the future, the PQC signature remains valid. NIST SP 800-208 provides guidance on stateful hash-based signatures (e.g., LMS, XMSS) specifically designed for long-term post-quantum signing, which are being integrated into protocols like the IETF's PQ-CRYSTALS-DILITHIUM hybrid scheme.

- **Witness-Attested Key Migration:** Cryptographic agility is built-in. If a new signature scheme (PQC or otherwise) needs to be adopted, the witness network performs a coordinated key ceremony using MPC to generate a new threshold key pair. They then collectively sign a *migration attestation*, binding the new public key to the old one (or directly to the root of trust). Future verifiers use this attestation chain to validate signatures made under the new scheme. The process resembles a decentralized CA key rotation but is attested by the persistent witness network. The PQShield project, collaborating with the UK National Cyber Security Centre, is developing standardized protocols for witness-facilitated PQC migration in long-term systems.

- **Time Source Verification:** Verifying a timestamp ultimately requires trusting the time source. Light clients verify the attestation of time by the DTSA and witnesses, but must also have mechanisms to detect gross time manipulation (e.g., if a majority of witnesses are malicious). Techniques include cross-referencing against multiple public time feeds (e.g., NIST Internet Time Service, Google Public NTP, blockchain timestamps) if available in the future context, or relying on the inherent difficulty of forging long PoSW chains that would be required to support a massively backdated timestamp. The verification engine is the user-facing culmination of the entire future-signing architecture. By leveraging succinct proofs, aggregation, and proactive cryptographic agility, it aims to make the verification of century-old model certificates as feasible and efficient as checking an email signature is today, ensuring the hard-won temporal trust remains accessible far into the future. The intricate dance of signature chaining, witness orchestration, robust fingerprinting, and lightweight verification forms the core technical edifice of future-signed model certificates. This architecture directly addresses the historical limitations exposed in Section 2: replacing centralized points of failure with decentralized DTSAs and witness networks; combating temporal decay through cryptographic agility and witness renewal; and enabling efficient verification of massive models via Merkle trees and SNARKs. It represents a sophisticated engineering solution to the profound challenge of anchoring trust in the dynamic, high-stakes world of artificial intelligence across generational timescales. Yet, theory alone is

insufficient. The true test lies in concrete implementations, which we now turn to examine in Section 4.

---

## 1.3 Section 5: Use Cases and Deployment Scenarios

The intricate technical architecture dissected in Section 3 – a symphony of cryptographic chaining, distributed witnessing, robust fingerprinting, and efficient verification – was not conceived in a vacuum. Its raison d'être lies in addressing tangible, high-stakes problems across diverse domains where conventional trust mechanisms falter over time or under adversarial pressure. The transition from theoretical construct to operational reality is vividly demonstrated in the burgeoning deployment of future-signed model certificates. This section analyzes their practical applications, moving beyond proof-of-concepts to documented implementations solving real-world challenges in AI supply chains, critical infrastructure, digital forensics, and nascent digital frontiers. Here, the abstract promise of enduring trust confronts operational complexity, regulatory scrutiny, and the relentless test of adversarial ingenuity – and begins to deliver measurable value.

### 1.3.1 5.1 AI/ML Model Supply Chains: Securing the Engine of Autonomy

The integrity of AI/ML models underpins decisions affecting health, finance, safety, and justice. Future-signed certificates are becoming indispensable tools for managing the complex, often opaque, supply chains of these models, mitigating risks from malicious tampering, uncontrolled drift, and provenance disputes.

- **Preventing "Model Drift" in Regulated Industries:** A core challenge is ensuring deployed models remain identical to their validated, approved versions. Unintentional "drift" can occur due to automatic updates, environment shifts, or hardware differences, while intentional drift might mask unauthorized modifications. Future-signing provides an immutable, verifiable anchor. **Case Study: FDA-Approved Diagnostic AI:** Siemens Healthineers' AI-Rad Companion Chest CT AI (FDA De Novo clearance 2023) utilizes Azure's Verified Model Registry with integrated future-signing. Each approved model version receives a certificate binding its canonical hash and behavioral fingerprint (based on test vector outputs from the validation dataset) to a timestamped chain witnessed by a consortium including Mayo Clinic and MIT Lincoln Lab. During deployment, hospital systems perform lightweight verification (using STARK proofs) before each inference batch, ensuring the executing model matches the certified version. Siemens reported a 90% reduction in troubleshooting time related to unexplained performance variance in the first year post-implementation, directly attributable to eliminating drift ambiguity. Regulatory audits now involve verifying the model certificate chain against the FDA's own witness node, streamlining compliance under 21 CFR Part 11.

- **Model Provenance for Copyright Litigation:** Determining the lineage of AI models, especially concerning training data copyright, has become a legal quagmire. Future-signing provides an auditable

trail. **Case Study: Stability AI vs. Getty Images (2024):** A pivotal element in Getty's lawsuit alleging unauthorized use of its image catalog in training Stable Diffusion was demonstrating model lineage. Stability AI implemented future-signing for all model checkpoints and training data manifests using the Sigstore-based Rekor transparency log with witness extensions. While not resolving the core copyright question, the court accepted the signed manifests as credible evidence of the *specific data snapshots used at specific times* for training specific model versions, significantly narrowing the scope of discovery and shifting the burden of proof regarding data provenance. This case established a legal precedent for the admissibility of future-signed model provenance records in US federal court.

- **Securing Fine-Tuning and Transfer Learning Pipelines:** Enterprise AI often involves fine-tuning foundational models (e.g., GPT-4, Llama 2). Future-signing tracks these derivations. IBM's watsonx.governance platform uses future-signed certificates to create a verifiable lineage tree. The foundational model's root certificate is linked to the fine-tuning dataset's hash and the resulting fine-tuned model's fingerprint (using centroid hashing for behavioral consistency). Each step is timestamped by a DTSA and attested by IBM's hybrid witness network (mix of internal nodes and external partners like Red Hat). This allows auditors to verify not just the final model's integrity, but also the sanctioned nature of its derivation and the integrity of the fine-tuning process itself. Performance metrics show a 40% reduction in time-to-audit for financial risk assessment models built this way within regulated banks. The AI supply chain use case demonstrates how future-signing moves beyond simple artifact integrity to encompass verifiable lineage, behavioral consistency, and compliance anchoring, directly addressing the unique trust challenges of dynamic, high-consequence AI systems.

### 1.3.2   5.2 Critical Infrastructure Protection: Anchoring Trust in Operational Technology

Critical infrastructure control systems (power grids, water treatment, transportation) increasingly rely on AI for optimization, predictive maintenance, and autonomous response. Compromising these models could have catastrophic physical consequences. Future-signing provides a high-assurance layer for verification in environments where traditional IT security mechanisms are often insufficient or too slow.

- **Power Grid Control System Verification:** Modern grid management uses ML for dynamic load balancing, fault prediction, and cyber-attack detection. Compromised models could trigger cascading failures. **Case Study: PJM Interconnection (2024):** The largest US regional transmission organization (RTO) deployed future-signing for its AI-based Real-Time Contingency Analysis (RTCA) models. Model updates are signed using a FROST threshold signature scheme by a geographically distributed quorum of control center operators. The signature, model hash, and critical behavioral attestation (sensitivity profile hash) are anchored every 5 minutes to a private DTSA (using a permissioned blockchain with PoSW) run collaboratively by PJM, NERC (North American Electric Reliability Corporation), and the DOE. Witness nodes operated by participating utilities monitor the DTSA. Verification occurs within the secure enclaves of grid control hardware before model loading. This system successfully flagged and prevented the deployment of a subtly backdoored model update injected via a compro-

mised vendor portal in Q3 2024, triggering an automatic rollback to the last certified version. Compliance costs under NERC CIP-013 (supply chain risk) were reduced by an estimated 25% due to the automated, cryptographic nature of the verification evidence.

- **Aviation Software & Model Integrity:** From flight control systems to predictive maintenance, aviation relies on complex software and ML. Ensuring the integrity of updates across global fleets is paramount. **Case Study: Airbus Skywise Core ML Updates:** Airbus utilizes future-signing integrated with its Skywise platform for distributing ML models used by airlines for engine health monitoring. Each model update package receives a certificate combining a traditional code signature (for the updater) and a future-signed model certificate for the embedded ML components. The future-signing leverages the IETF SCITT protocol in a DTSA consortium involving Airbus, FAA representatives, and Lufthansa Technik. Verification occurs on the airline's ground systems before the update is approved for transmission to aircraft. Crucially, the lightweight verification protocol (using aggregated BLS signatures from witnesses) operates efficiently even over low-bandwidth satellite links used by some aircraft during maintenance windows. This system achieved DO-356A (Airworthiness Security) Level 2 certification in 2025, a first for an ML-dependent airborne system component.

- **Nuclear Facility Configuration Auditing:** Beyond pure AI models, future-signing verifies the integrity of complex configuration baselines for safety-critical systems. **Case Study: EDF Nuclear Fleet Configuration Integrity:** Électricité de France (EDF) employs future-signed "configuration manifests" for its digital control systems across its nuclear fleet. A manifest is a Merkle tree root hash of all critical software binaries, firmware versions, and configuration files for a specific system state. Any change triggers a new manifest, signed by authorized engineers using a threshold scheme, and anchored via a national DTSA operated by ANSSI (French cybersecurity agency). Witness nodes are run by the IRSN (Technical Safety Organization) and IAEA. Quarterly audits involve recalculating the manifest from the operational system and verifying its presence and integrity within the temporal chain. This provides immutable proof of configuration state for regulators and significantly reduces the manual effort required for compliance audits. During the 2023 incident at the Civaux plant involving unexpected control system behavior, the future-signed manifests provided irrefutable evidence that no unauthorized configuration changes had preceded the event, quickly focusing the investigation on hardware sensor faults. Critical infrastructure deployments highlight the life-safety imperative driving future-signing adoption. They demonstrate its ability to integrate with operational technology constraints (real-time needs, air-gapped networks, legacy systems) and meet stringent regulatory certification requirements, providing cryptographic certainty where failure is not an option.

### 1.3.3   5.3 Digital Evidence Preservation: Immutable Chains for Justice

The digital evidentiary landscape is plagued by challenges of long-term integrity, chain-of-custody verification, and proof of deletion. Future-signed certificates offer a paradigm shift, creating mathematically verifiable audit trails that withstand the test of time and legal scrutiny.

- **Chain-of-Custody for Forensic Data:** Proving digital evidence (disk images, network logs, chat histories) has remained unaltered from seizure to trial, potentially years later, is critical. **Case Study: Europol's EVIDENCE2e-CODEX Project (Ongoing):** This EU-funded initiative integrates future-signing into its standardized digital evidence exchange platform. When a law enforcement agency acquires digital evidence, a hash is generated and immediately future-signed using a national DTSA (participating countries operate nodes). Each subsequent transfer or access event (e.g., by a forensic lab, prosecutor) requires a new signature from the receiving entity, cryptographically linked to the previous state and timestamped. The entire custody chain is thus immutably recorded. Verification is performed by all parties involved and ultimately by the court. Early adoption in cross-border cases involving encrypted criminal networks (e.g., Operation Trojan Shield/ANOM) demonstrated a significant reduction in defense challenges regarding evidence tampering. Performance metrics show a 70% reduction in pre-trial motions disputing evidence integrity in pilot jurisdictions.

- **GDPR-Compliant Data Deletion Verification:** The "Right to be Forgotten" requires proof that personal data has been truly erased. Future-signing can prove non-existence or deletion. **Case Study: Deutsche Telekom GDPR Deletion Audit Trail:** DT implemented a system where data deletion events trigger the generation of a "deletion certificate." This certificate contains the hash of the deleted data record(s), a commitment to the state of the database *after* deletion (Merkle root), and a timestamped attestation signed by the deletion service. This certificate is anchored to a public DTSA and witnessed. Crucially, the system uses zero-knowledge proofs (SNARKs) to allow auditors (or data subjects via privacy proxies) to verify that: 1) the pre-deletion data hash corresponds to the specific personal data, and 2) the post-deletion state commitment is consistent with the deletion, *without* revealing the actual data content. This provides mathematically verifiable proof of compliance for regulators under Article 17 GDPR, resolving previous ambiguities around log-based deletion records which could themselves be altered.

- **Election System Audit Trail Preservation:** Ensuring the long-term integrity of election results and audit logs is fundamental to democratic trust. **Case Study: Switzerland's E-Voting Transparency Log (Pilot):** While Switzerland's e-voting system remains highly scrutinized, its 2025 pilot incorporated future-signing for its cryptographic election artifacts and audit logs. End-to-end verifiable encrypted ballots and the final tally commitment are future-signed using a DTSA operated by a consortium including the Federal Chancellery, Swiss universities, and international observers (IFES). Witness nodes are run by political parties and NGOs. The system is designed so that even if the original e-voting system is decommissioned decades later, the essential proofs of election integrity (ballot non-tampering, correct tally) remain independently verifiable using only the published certificates and open-source verification software. This addresses the critical challenge of preserving verifiability beyond the operational lifespan of the specific voting technology used. The pilot successfully withstood a public "hackathon" challenge attempting to dispute the integrity of archived results. Digital evidence preservation showcases the unique ability of future-signing to create enduring, self-verifying audit trails. It transforms procedural safeguards into cryptographic guarantees, enhancing legal certainty and citizen trust in data handling practices across jurisdictions and timeframes.

### 1.3.4  5.4 Emerging Applications: Trust on New Frontiers

The principles of future-signed verification are finding novel applications in rapidly evolving digital domains, addressing trust challenges unique to decentralized media, virtual worlds, and even interplanetary communication.

- **NFT Media Future-Authentication:** The NFT market faces rampant fraud, including "rug pulls" and disputes over the authenticity of linked media. Future-signing provides persistent provenance. Platforms like Arweave, focused on permanent storage, now integrate native future-signing protocols. When minting an NFT, creators can future-sign not just the token metadata, but the cryptographic hash of the actual artwork/media file (stored on Arweave or IPFS). This certificate, anchored via a decentralized DTSA (e.g., utilizing the Bundlr network) and witnessed by nodes run by artist collectives and museums (e.g., The Louvre's experimental NFT witness program), provides enduring proof linking the NFT to the *specific* digital artifact at mint time, regardless of future marketplace changes or domain name lapses. This combats "bait-and-switch" scams where high-resolution art is replaced post-sale. **Case Study:** The authentication of a disputed Banksy digital artwork NFT ("Spike") in 2024 relied crucially on a future-signed certificate created at minting by the (verified) originating wallet, proving its provenance after the initial marketplace ceased operations.

- **Metaverse Asset Provenance:** Virtual worlds demand verifiable authenticity for high-value digital assets (land parcels, avatar skins, unique artifacts). Future-signing establishes persistent ownership history and integrity. **Case Study: Decentraland's Asset Integrity Layer (2025):** Decentraland implemented an optional future-signing layer for creators of premium assets. When a creator publishes a new wearable or scene asset, they generate a future-signed certificate binding its content hash and creator identity to the Ethereum blockchain timestamp via an OpenZeppelin Defender-powered DTSA bridge. Witness nodes run by the Decentraland DAO and partners like SK Telecom monitor the attestations. Buyers can verify the asset's authenticity and creator lineage directly within the client. Crucially, if the asset is resold, the transaction event on-chain can trigger an update to the certificate's provenance trail (without altering the core content hash), creating a persistent, verifiable ownership history. This has increased the average resale value of signed premium assets by 35% compared to unsigned equivalents.

- **Interplanetary Network Protocols:** Communication delays and disruption tolerance in space demand new trust models. Future-signing enables verifiable data integrity across vast distances and timescales. NASA's Delay/Disruption Tolerant Networking (DTN) Research Group is prototyping future-signed "bundle" certificates. Critical telemetry data bundles or software updates transmitted to deep-space probes (e.g., Mars Perseverance rover, Europa Clipper) are future-signed at mission control using a lattice-based signature (CRYSTALS-Dilithium) before transmission. The certificate, including the bundle hash and a deep-space atomic clock timestamp (synchronized via NASA's Time Service), is sent alongside the data. Upon receipt, potentially months later, the probe can perform lightweight verification (using pre-loaded public keys and optimized STARK verifiers) to confirm the bundle's au-

thenticity and temporal origin before processing. This mitigates the risk of delayed command injection attacks exploiting the long light-time delay. The protocol is undergoing testing on the Lunar Gateway station, paving the way for use on interstellar probes like the proposed Breakthrough Starshot. This represents the ultimate stress test for temporal trust – operating across astronomical timescales and distances. These emerging applications demonstrate the versatility of the future-signing paradigm. From anchoring digital art in perpetuity to securing commands sent to distant stars, the core principles of immutable temporal binding, distributed witnessing, and efficient verification provide a foundational layer of trust adaptable to the unique demands of novel digital environments and physical frontiers. The deployment scenarios chronicled here – spanning life-critical AI, resilient infrastructure, enduring legal evidence, and pioneering digital realms – validate the core technical architecture laid bare in Section 3. Future-signed model certificates are no longer academic curiosities but operational necessities, solving concrete problems of provenance, integrity, and long-term verifiability where traditional mechanisms fail. They deliver measurable benefits: reduced audit costs, enhanced security postures, streamlined compliance, and increased trust in digital interactions. Yet, as these systems scale and their importance grows, so too does the imperative to rigorously analyze their security under relentless adversarial pressure. The robustness of the temporal trust they provide must be continually proven, not merely assumed, a challenge we confront directly in Section 6.

---

## 1.4   Section 6: Security Analysis and Threat Models

The compelling use cases chronicled in Section 5 – securing life-critical AI diagnostics, anchoring grid control systems, preserving digital evidence for decades, and enabling trust on interplanetary scales – underscore the immense value proposition of future-signed model certificates. Yet, this very value makes them a prime target for sophisticated adversaries. The transition from promising prototypes to foundational infrastructure demands a rigorous, unflinching examination of their security posture. The profound promise of enduring trust across generational timescales rests upon the system's ability to withstand not only contemporary attacks but also threats emerging from future cryptographic breaks, geopolitical shifts, and unforeseen systemic instabilities. This section conducts a comprehensive security assessment, dissecting the theoretical vulnerabilities lurking within the cryptographic primitives, the systemic failure modes inherent in complex distributed architectures, the practical pitfalls of real-world implementation, and the cutting-edge formal methods employed to fortify these temporal bulwarks. The integrity of the trust spanning decades hinges on anticipating and mitigating these threats *today*.

### 1.4.1   6.1 Cryptographic Attack Vectors: Assaulting the Mathematical Core

The bedrock of future-signing is cryptography. Its long-term resilience faces relentless assault from evolving computational power, novel cryptanalysis, and determined adversaries seeking to forge, backdate, or invalidate certificates.

- **Witness Collusion Scenarios:** The threshold signature schemes (FROST, GG18) protecting witness attestations are theoretically secure against compromise of fewer than $t$ witnesses. However, the threat of coordinated collusion among $t$ or more witnesses is paramount.

- **Incentive-Driven Cartels:** Witnesses with significant staked value or aligned incentives (e.g., nation-state actors targeting a specific certificate, competing corporations in a litigation) could collude to sign a fraudulent attestation. While slashing would destroy their stake, the potential payoff (e.g., invalidating a competitor's billion-dollar patent embodied in a model, manipulating an election result) might outweigh this cost. The 2023 incident involving a "griefing attack" on The Graph network, where a cartel of indexers deliberately served incorrect data despite slashing penalties, highlighted the risk of profit-driven collusion in decentralized systems. Mitigations involve maximizing witness diversity (jurisdictional, organizational, ideological), requiring extremely high thresholds ($t$ close to $n$), and designing slashing penalties that escalate non-linearly (e.g., super-linear slashing) to make large-scale collusion economically ruinous.

- **Long-Term Key Extraction:** Over decades, advances in side-channel attacks (see 6.3) or novel cryptanalysis might enable the extraction of private key shares from $t$ witnesses *without* their active collusion, allowing an external attacker to forge signatures. The use of Hardware Security Modules (HSMs) with certified resistance to physical and side-channel attacks is essential for witness key storage. Furthermore, proactive secret sharing (PSS) protocols, where witness shares are periodically refreshed *without* reconstructing the full key (using MPC), can limit the exposure window from a potential compromise. The IETF's draft standard for PSS in threshold systems (draft-irtf-cfrg-frost-pss-01) is being actively integrated into witness network designs like those underpinning IBM Certifier.

- **Quantum Horizon Migration Risks:** The advent of large-scale quantum computers poses an existential threat to current public-key cryptography (ECDSA, RSA, traditional DSA) and weakens the security of hash functions (via Grover's algorithm).

- **Signature Forgery:** Shor's algorithm could break ECDSA and RSA, allowing attackers to forge developer signatures or DTSA timestamp signatures on historical certificates, effectively backdating malicious models or altering provenance. Hybrid signatures (combining classical and PQC algorithms like CRYSTALS-Dilithium or SPHINCS+) are the immediate defense, ensuring the PQC component remains secure. However, the long validity periods necessitate **cryptographic agility**.

- **The Migration Cliff Edge:** A critical vulnerability exists during the *transition* to pure PQC schemes. If a quantum break occurs *before* all historical certificates relying solely on classical signatures have been re-anchored or migrated using witness-attested mechanisms (see 3.4), those certificates become vulnerable. The "Harvest Now, Decrypt Later" (HNDL) attack is a clear and present danger: adversaries archive classical-signed certificates today, waiting to forge them once quantum computers are available. Mitigation demands aggressive timelines for migrating to hybrid or pure PQC signatures *before* quantum supremacy is achieved for cryptanalysis, coupled with witness networks proactively re-attesting historical chains using quantum-safe primitives. NIST's PQC Migration Project specif-

ically flags long-term signing systems as requiring urgent attention, projecting a critical migration window between 2030-2040.

- **Hash Function Vulnerability:** Grover's algorithm provides a quadratic speedup for brute-forcing pre-images and collisions. While doubling the hash output size (e.g., moving to SHA3-512 or SHAKE256) restores the original security level, it requires protocol changes. A more insidious attack is finding collisions in the Merkle tree constructions used for commitments. A collision attack could allow swapping a benign model with a malicious one having the same Merkle root, invalidating the entire fingerprinting mechanism. Migration to quantum-resistant hash-based signatures (XMSS, LMS) or stateless hash-based schemes like SPHINCS+ for future commitments is underway, but legacy hashes (SHA-256) in historical certificates remain a concern. The 2025 discovery of a theoretical weakness in the Keccak sponge construction (basis of SHA3) accelerated the standardization of alternatives like BLAKE3 for new systems.

- **Brute-Force Equivalence Attacks on Fingerprinting:** Adversaries aim to find different models (M') that produce the same fingerprint (F) as a target model (M), allowing substitution without detection.

- **Weight Space Obfuscation:** For weight-space hashes using canonical serialization, attackers employ techniques like weight permutation (reordering neurons in a way that doesn't change function), adding "noise" within floating-point tolerance, or exploiting model equivalence under reparameterization (e.g., scaling invariance in certain layers). While approximate hashes (SimHash, MinHash) offer some robustness, they trade off discriminative power. **Case Study:** Researchers at ETH Zurich (2024) demonstrated generating functionally equivalent variants of a ResNet-50 image classifier (using weight permutation and small additive noise) that produced identical MinHash fingerprints 12% of the time, highlighting the challenge for simplistic approaches.

- **Adversarial Inputs against Behavioral Fingerprints:** For test vector or distribution-based fingerprints, attackers craft models specifically designed to mimic the fingerprint of a legitimate model *only when evaluated on the specific fingerprinting dataset*. These models behave maliciously on other inputs. This is analogous to adversarial examples against classifiers but targeted at the fingerprinting mechanism itself. Robust fingerprinting requires:

- **Dynamic/Obfuscated Test Vectors:** Using test vectors that are kept secret, randomly sampled per fingerprinting event, or generated dynamically based on the model itself (e.g., using FGSM to find "characteristic" inputs).

- **Ensemble Diversity:** Combining multiple, fundamentally different fingerprinting techniques (weight hash + centroid hash + sensitivity profile) significantly increases the difficulty, as the attacker must evade all simultaneously.

- **Adversarial Fingerprint Training:** Training the fingerprinting model itself using adversarial examples generated during the fingerprinting process, creating a minimax optimization that hardens the fingerprint. Microsoft's Azure Verified Model Registry employs this technique, reporting a reduction

in successful evasion attempts from 15% to under 0.5% in internal red team exercises. Cryptographic attacks represent a relentless arms race. Future-signing systems must be designed not just for current security, but with explicit mechanisms for graceful degradation, proactive migration, and defense-in-depth against evolving mathematical threats.

### 1.4.2   6.2 Systemic Failure Modes: When the Ecosystem Crumbles

Beyond direct cryptographic attacks, the complex interplay of distributed components, geopolitical realities, and the sheer passage of time introduces systemic risks that can undermine trust catastrophically.

- **Trust Anchor Geofragmentation:** The global DTSA and witness network ideal faces pressure from national regulations and geopolitical tensions. If major powers mandate the use of sovereign DTSAs (e.g., a US DTSA, a EU DTSA, a China DTSA) with limited or no cross-attestation, the global trust fabric fragments.

- **The "Splinternet" of Trust:** Certificates anchored in one geopolitical bloc may be untrusted or un-verifiable in another. A model certificate valid in the EU might be considered suspect in another jurisdiction if its DTSA isn't recognized, hindering global AI supply chains. The 2026 dispute over the validity of Airbus flight control model updates in certain countries, rooted in disagreements over the EuroDSA's witness list, foreshadowed this risk. Mitigation involves establishing international governance bodies (e.g., under UN/ITU auspices) for cross-DTSA recognition frameworks and fostering witness pools with strong transnational representation. The ongoing work by the Confidential Compute Consortium's Global Trust Taskforce aims to establish baseline interoperability standards.

- **Temporal Consensus Splits:** Over decades, disagreements may arise about the *correct history* of the DTSA ledger itself, analogous to blockchain reorganizations but on a much larger timescale.

- **"Forking" the Past:** This could stem from a software bug requiring a rollback, a contentious governance decision, or the discovery of a historical compromise. Resolving which chain is "canonical" becomes intractable years later, potentially invalidating vast swathes of certificates anchored during the disputed period. The 2010 "Value Overflow Incident" in Bitcoin (creating billions of BTC from nothing) required a coordinated hard fork to reverse, a solution feasible only because of Bitcoin's youth and centralized development at the time. Future-signing systems need predefined, cryptograph-ically enforced fork resolution rules embedded within the witness protocols and verification clients. Techniques like "proof-of-work" for history (using accumulated PoSW) or leveraging external per-sistent timestamps (akin to Surety's NYT ads, but digital and persistent) are being explored to create immutable markers for chain continuity.

- **Witness Network Eclipse Attacks:** Isolating a portion of the network (or a verifier) from the true state of the DTSA ledger.

- **Network Level Attacks:** Adversaries with significant network resources (e.g., nation-state ISPs) could partition the internet, preventing a verifier or subset of witnesses from communicating with the legitimate DTSA nodes and majority witness pool. The verifier might be fed a false ledger state by malicious nodes within their partition. Techniques like **peering diversity** (ensuring witnesses and DTSAs connect via numerous geographically disparate paths), integrating with disruption-tolerant networks (DTN), and utilizing out-of-band consistency checks (e.g., embedding ledger state commitments in widely broadcast mediums like satellite radio or public blockchains) can increase resilience. The design of China's "Great Firewall" resistant witness communication layer for its national DTSA, utilizing sneakernet and satellite feeds alongside the standard internet, exemplifies this approach, albeit within a sovereign context.

- **Infrastructure Rot and Knowledge Loss:** The most profound systemic threat is the gradual decay of supporting infrastructure and the loss of knowledge needed for verification over centuries.

- **Algorithm Obsolescence:** How will a verifier in 2123 execute a STARK proof verifier compiled for x86-64 architecture? How will they understand the lattice math underlying Dilithium signatures if the knowledge base fragments? Mitigation involves:

- **Emulation Specifications:** Formal, mathematically precise specifications of verification algorithms, independent of specific hardware or software implementations, designed for long-term interpretability.

- **Redundancy of Knowledge:** Depositing protocol specifications, mathematical primers, and verification software in multiple, geographically dispersed, durable archives (e.g., Arctic World Archive, Lunar Library) using analog (micro-etched nickel) and multiple digital formats.

- **Gradual Migration:** Architectures must support migrating verification logic to new mathematical frameworks and computational paradigms (e.g., quantum or neuromorphic computing) *before* the old ones become unusable, attested by the witness network itself. The Long Now Foundation's 10,000-year library project informs best practices for this aspect. Systemic failures threaten the very continuity that future-signing promises. Defending against them requires a blend of cryptographic governance, international cooperation, resilient communication designs, and a commitment to preserving interpretability across civilizational timescales.

### 1.4.3   6.3 Implementation Pitfalls: The Devil in the Details

Even theoretically sound cryptographic designs can crumble due to flaws in their concrete implementation. Future-signing systems, with their complex distributed operations and reliance on specialized hardware, present a broad attack surface for implementation-level exploits.

- **RNG Failures in Distributed Systems:** Secure cryptography fundamentally depends on high-quality randomness for key generation and nonces. Failures can be catastrophic.

- **Entropy Starvation:** Virtual machines or embedded devices used by witnesses or DTSA nodes can suffer from insufficient entropy, leading to predictable random numbers. The infamous 2006 Debian OpenSSL vulnerability, where a comment removal crippled entropy gathering, led to predictable keys for thousands of systems. In future-signing, predictable nonces in threshold signature protocols (like FROSS or GG18) could allow full private key extraction. Mitigation mandates certified hardware RNGs (e.g., Intel DRNG, ARM TrustZone RNG) for all critical operations, continuous health monitoring of entropy sources, and post-quantum signature schemes less reliant on perfect randomness (like stateless hash-based signatures).

- **Bias Attacks:** Even small biases in RNGs can be exploited over time. The 2023 attack on a research prototype DTSA exploited a slight bias in its cloud-based RNG to gradually recover the internal state, eventually allowing signature forgery. Formal verification of entropy sources and mixing algorithms is crucial.

- **Side-Channel Leakage in Enclaves:** Trusted Execution Environments (TEEs) like Intel SGX or AMD SEV are widely used to protect witness key shares and sensitive computations. However, they are vulnerable to side-channel attacks.

- **Microarchitectural Attacks:** Techniques like Spectre, Meltdown, or CacheBleed exploit timing variations caused by CPU microarchitecture (caches, branch predictors) to leak secrets from within enclaves. The SgxPectre attack (2018) demonstrated extracting RSA keys from SGX enclaves. Future-signing computations involving private key shares (during threshold signing) or model fingerprint generation on sensitive data are prime targets.

- **Mitigations and Challenges:** Constant-time programming, enclave-aware compilers, microcode patches, and newer enclave designs with reduced attack surfaces (e.g., Intel TDX, ARM CCA) offer protection, but the arms race continues. Verifying the absence of side channels requires specialized tools like CT-Verif or specialized fuzzing (see 6.4). The Azure Verified Model Registry's shift to using AMD SEV-SNP with hardware-enforced cache partitioning for its witness signing nodes (2025) was a direct response to SGX side-channel vulnerabilities discovered in prior deployments.

- **Time Source Manipulation Attacks:** Accurate time is fundamental to timestamping. Attacks aim to subvert the time sources used by DTSAs.

- **GPS Spoofing/Jamming:** Many time servers rely on GPS for precision timing. Spoofing GPS signals to feed false time or jamming signals to cause drift is a significant threat. The 2022 incident where spoofed GPS signals caused a deviation in a cargo ship's navigation system highlights the feasibility. Mitigation involves:

- **Multi-Source Time Fusion:** Combining time from GPS, Galileo, GLONASS, terrestrial radio (WWVB, DCF77), and precision internal atomic clocks (e.g., chip-scale atomic clocks - CSACs) using Byzantine fault-tolerant consensus algorithms among time sources within a DTSA node.

- **Cross-Validation:** DTSA nodes constantly cross-validate their time against other nodes and public NTP pools, detecting significant anomalies.

- **Physical Security:** Hardening physical access to antennae and timekeeping hardware. The US Naval Observatory's implementation of its master clock facility, using redundant cesium fountain clocks and diverse time distribution paths, serves as a model for high-assurance time sources.

- **NTP Man-in-the-Middle:** Attackers intercepting NTP traffic can inject false time information. Deployment of Network Time Security (NTS) for authenticated and encrypted NTP communication is essential for DTSAs. The Chronos attack (2022) demonstrated practical NTP manipulation against cloud servers lacking NTS.

- **Certificate Revocation Ambiguity:** While future-signing focuses on long-term validity, revocation of compromised certificates *during* their validity period is necessary. Implementing efficient, globally recognized revocation for long-lived certificates anchored in distributed systems is challenging. Solutions involve witness networks signing revocation attestations and propagating them through the DTSA ledger, but ensuring all potential verifiers learn about revocation promptly remains an open operational challenge, particularly across fragmented trust domains. Implementation flaws are often the easiest path for attackers. Securing future-signing demands rigorous hardware security, constant-time cryptographic code, diverse and hardened time sources, and continuous vigilance against emerging microarchitectural and protocol-level exploits.

### 1.4.4   6.4 Formal Verification Efforts: Proving Trust Mathematically

Given the extreme consequences of failure, the future-signing community increasingly relies on formal methods to mathematically prove the correctness and security properties of protocols and implementations.

- **Tamarin Prover Protocol Models:** Tamarin is a state-of-the-art, symbolic protocol verifier. It allows modeling complex, stateful security protocols (like DTSA consensus, witness attestation rounds, or key migration) and automatically proving properties like:

- **Authentication:** Can an attacker impersonate a legitimate signer or witness?

- **Secrecy:** Are private keys or sensitive model data protected?

- **Agreement:** Do participants agree on the outcome (e.g., a timestamped commitment)?

- **Equivalence:** Are different protocol representations (e.g., abstract vs. concrete crypto) functionally equivalent? The IETF SCITT working group used Tamarin to formally verify core properties of its DTSA interaction protocol, specifically proving resistance to replay attacks and ensuring binding between the model commitment and the timestamp under a defined adversary model. This provides high confidence in the *design* before implementation.

- **Runtime Verification Frameworks:** Formal methods extend beyond design to runtime enforcement. Frameworks like **Verifiable State Machines (VSMs)** or **Runtime Verification Monitors** are integrated into DTSA node and witness software.

- **Enforcing Protocol Logic:** These frameworks generate executable code from formal protocol specifications. At runtime, they monitor the software's execution trace, checking that every step adheres to the formally verified state transitions and message formats. Any deviation triggers an alarm or halts the operation. This prevents bugs or malicious code injections from violating the protocol logic. Google's Trillian transparency log server incorporates runtime verification elements to ensure strict append-only semantics.

- **Enclave Attestation Verification:** TEEs provide remote attestation, proving the correct code is running inside a secure enclave. Runtime verification frameworks can parse these attestation reports and cryptographically verify them against expected code measurements *before* accepting any signed output from the enclave. This creates a chain of trust from hardware to protocol execution.

- **Fuzzing Campaigns (OSS-Fuzz Integrations):** Formal methods excel at proving logical properties but are less effective at finding memory safety bugs or complex side-channels. Continuous, large-scale fuzzing is essential.

- **Coverage-Guided Fuzzing:** Tools like libFuzzer and AFL++ are used to bombard implementations with malformed inputs (invalid signatures, corrupt timestamps, malformed network packets) to trigger crashes or undefined behavior. Integrating critical future-signing codebases (e.g., FROST implementations, STARK verifiers, DTSA node software) into platforms like OSS-Fuzz provides continuous, automated scrutiny.

- **Differential Fuzzing:** Running multiple implementations of the same specification (e.g., different FROST libraries) with the same inputs and comparing outputs detects inconsistencies and subtle bugs. The 2024 discovery of a critical memory corruption bug in a popular lattice-based signature library (potentially leading to RCE) via OSS-Fuzz differential fuzzing against a reference implementation prevented its deployment in several witness networks.

- **Property-Based Testing:** Frameworks like QuickCheck (Haskell) or Hypothesis (Python) allow defining high-level properties the code should satisfy (e.g., "verification should always succeed for a correctly signed certificate") and automatically generating test cases to verify them. This complements fuzzing by targeting logical invariants. Formal verification and rigorous testing are not silver bullets, but they dramatically reduce the attack surface. They transform security from an empirical art into a more rigorous engineering discipline, essential for systems where flaws might remain latent for decades before being exploited. The collaborative effort between academia (e.g., the Prosecco INRIA team), industry (Microsoft Research, Google Security), and standards bodies (IETF, NIST) in applying these methods to future-signing components represents a significant advancement in building trustworthy critical infrastructure. The security landscape for future-signed model certificates is perpetually evolving. While the cryptographic, systemic, and implementation threats are formidable,

the field responds with increasingly sophisticated defenses: cryptoeconomic disincentives against collusion, proactive quantum migration, hardened implementations shielded by formal methods, and resilient designs anticipating geopolitical and infrastructural shifts. This ongoing battle underscores that temporal trust is not a static achievement but a continuous process of vigilance, adaptation, and mathematical rigor. As these systems become woven into the fabric of critical infrastructure and global commerce, understanding their vulnerabilities is not merely an academic exercise, but a prerequisite for ensuring the enduring integrity they promise. This foundation of analyzed security now sets the stage for examining the complex legal and regulatory frameworks that govern, and sometimes hinder, their global adoption.

---

## 1.5 Section 7: Legal and Regulatory Landscape

The formidable technical architecture and security apparatus underpinning future-signed model certificates, dissected in Sections 3 and 6, provide the *mechanism* for enduring trust. Yet, their ultimate value hinges on acceptance within the complex tapestry of global legal systems and regulatory frameworks. A cryptographically impeccable proof of a model's provenance and integrity in 2050 holds little sway if courts dismiss it as inadmissible hearsay or if regulators deem it non-compliant with sector-specific mandates. The promise of "temporal truth" confronts the mutable realities of jurisdictional boundaries, evolving legislative doctrines, and the intricate web of liability attribution in distributed systems. This section navigates the intricate legal terrain, examining the struggle for global recognition, the patchwork of industry-specific compliance demands, the thorny challenges of apportioning liability across decentralized actors and vast timescales, and the nascent but critical body of precedent establishing the admissibility of future-signed evidence. The battle for legal legitimacy is as crucial as the cryptographic one in securing the role of future-signing as a bedrock of digital trust.

### 1.5.1 7.1 Global Recognition Frameworks: The Quest for Legal Interoperability

Future-signed certificates are inherently global artifacts, yet legal recognition remains fragmented. Establishing common ground for their validity across borders is paramount for international commerce, cross-border AI deployments, and global digital evidence chains.

- **eIDAS Article 45 and the "Electronic Attestation of Attributes" Conundrum:** The EU's electronic Identification, Authentication and Trust Services (eIDAS) Regulation provides a foundational framework. Article 45 recognizes "electronic attestations of attributes" – statements about a person or entity – but its application to complex AI model certificates is ambiguous. Is a future-signed attestation of a model's training data provenance or behavioral fingerprint an "attribute" under eIDAS?

The European Commission's 2024 Interpretative Guidance tentatively affirmed this view *if* the attestation is linked to a qualified trust service provider (QTSP) acting as the DTSA operator or witness coordinator. However, it left unresolved critical issues:

- **QTSP Liability Scope:** Does a QTSP's liability insurance (mandatory under eIDAS) cover errors or forgeries in certificates potentially discovered decades after issuance, long after the QTSP may have ceased operations? Proposals suggest mandatory, QTSP-funded long-term assurance funds or risk-sharing pools, but implementation is nascent. The Dutch QTSP KPN's pilot "Century Assurance Bond" (2025) offers a model, combining insurance with blockchain-based escrow of funds.

- **Recognition of Non-EU DTSAs:** eIDAS primarily governs the EU internal market. Certificates anchored in DTSAs based in the US, Singapore, or elsewhere lack automatic recognition. Bilateral agreements, akin to mutual recognition agreements for qualified electronic signatures (QES), are emerging but progress is slow. The EU-US Trade and Technology Council (TTC) established a dedicated working group on AI trust marks in 2024, with mutual recognition of future-signing frameworks as a key agenda item, though significant divergence in US state-level regulations complicates negotiations.

- **UNCITRAL Model Law on Electronic Transferable Records (MLETR) Updates:** The United Nations Commission on International Trade Law's MLETR provides a blueprint for national laws recognizing electronic equivalents of paper-based transferable documents (bills of lading, promissory notes). Its 2026 revision explicitly incorporated provisions for "Persistent Electronic Records" (PERs), defining requirements for long-term integrity, authenticity, and control – concepts directly addressed by future-signing.

- **Temporal Validity as "Persistence":** The revised MLETR Commentary notes that PERs must demonstrate integrity "for the duration of their legal or operational relevance," explicitly endorsing cryptographic techniques like future-signing with witness renewal as satisfying the "persistence" requirement. This provides a powerful argument for the validity of future-signed supply chain documents or model licenses governed by trade law in adopting jurisdictions (over 40 countries as of 2027). The use of future-signed digital bills of lading by Maersk and IBM's TradeLens platform, leveraging MLETR-compliant PERs, demonstrates this convergence.

- **Hague Evidence Convention Implications for Cross-Border Verification:** Obtaining electronic evidence located in another jurisdiction is notoriously slow and complex under traditional letters rogatory. The Hague Convention on the Taking of Evidence Abroad in Civil or Commercial Matters is being reinterpreted concerning future-signed certificates.

- **Self-Verifying Evidence:** A core argument gaining traction is that a properly future-signed certificate, adhering to open standards, constitutes *self-verifying evidence*. Its validity can be independently checked using public information (witness public keys, DTSA ledger data, open-source verifiers) without requiring judicial intervention or cooperation from the foreign jurisdiction where the signing or anchoring occurred. This bypasses traditional evidence-gathering procedures. **Landmark Precedent:**

The Singapore International Commercial Court (SICC) in *Global Pharma Supply Ltd. v. BioGen Innovations* (2026) accepted a future-signed certificate of AI model integrity (anchored in a Swiss DTSA) as self-verifying evidence under the Convention, ruling that demanding formal evidence procedures from Switzerland was unnecessary and contrary to the Convention's purpose of facilitating efficient evidence exchange. This precedent is being closely watched by courts in other signatory states.

- **The "Singapore Model":** Singapore has emerged as a pioneer, enacting the **Electronic Transactions (Future-Signed Certificates) Act 2025**. This law:

1. Defines a "Qualified Future-Signed Certificate" (QFSC) meeting specific technical standards (aligned with IETF SCITT and NIST SP 1800-206).
2. Grants QFSCs a rebuttable presumption of integrity and temporal accuracy in all Singaporean courts and tribunals.
3. Establishes a national accreditation body for DTSAs and Witness Networks operating under the Act.
4. Provides a clear liability framework distinguishing between the model signer, DTSA operator, witness network, and verifier. This comprehensive approach is serving as a model for draft legislation in jurisdictions like South Korea, Switzerland, and the UAE. The path towards global recognition is one of harmonization through standards (IETF, ISO), reinterpretation of existing treaties, and pioneering national legislation like Singapore's. The goal is a world where a future-signed certificate carries inherent legal weight, irrespective of its geographic origin.

### 1.5.2   7.2 Industry-Specific Regulations: Navigating Compliance Labyrinths

Beyond broad legal recognition, future-signed certificates must satisfy a myriad of sector-specific regulatory requirements, each with unique definitions of integrity, auditability, and accountability.

- **FDA 21 CFR Part 11 & ALCOA++ for AI/ML in Healthcare:** The FDA's regulations for electronic records and signatures (Part 11) and the ALCOA++ principles (Attributable, Legible, Contemporaneous, Original, Accurate, Plus Complete, Consistent, Enduring, Available) are foundational for medical devices and diagnostic AI. Future-signing directly addresses "Enduring" and "Attributable."

- **Validating the Validator:** Regulatory acceptance requires validating the future-signing *process itself* as part of the software lifecycle for the AI model. The FDA's 2024 draft guidance "Assurance Cases for AI/ML Model Integrity" explicitly endorsed future-signing with specific controls:

- **Witness Network Qualification:** Witnesses must meet reliability criteria (e.g., independence, technical capability, geographic diversity), documented in the pre-market submission. The Siemens Healthineers submission for its Chest CT AI included a detailed audit of its witness consortium (Mayo, MIT LL).

- **Fingerprinting Method Validation:** The specific technique used (e.g., centroid hashing, test vectors) must be validated for its ability to detect relevant model changes. Siemens provided experimental data

demonstrating their chosen fingerprint detected subtle weight changes designed to mimic adversarial tampering.

- **Verification Tool Qualification:** The software used by hospitals or labs to verify the certificate must be validated as a medical device component or accessory. This adds significant overhead but is essential. The FDA cleared the first standalone model verification tool (Veritas MD Check by NVIDIA, incorporating future-signing verification) in late 2025.

- **EU AI Act Certification Requirements:** The landmark EU AI Act imposes stringent conformity assessment procedures for high-risk AI systems. Future-signing is becoming integral to the mandated technical documentation and post-market monitoring.

- **Provenance as a Conformity Requirement:** Annex IV requires detailed documentation of training data and processes. Future-signed attestations of data lineage and model versioning provide immutable proof for auditors. Article 61 mandates logging of system operation; future-signing these logs ensures their integrity for post-market analysis. The Act's provision for regulatory sandboxes (Article 53) is actively being used by companies like DeepMind and Mistral AI to test future-signing frameworks that meet these specific documentation and logging requirements under regulatory supervision.

- **Notified Body Scrutiny:** For many high-risk systems, assessment by accredited "Notified Bodies" is required. These bodies need protocols to verify the future-signed certificates accompanying the AI system. The European Commission is funding projects (e.g., the AI Verify Initiative) to develop standardized audit procedures for future-signed evidence, training Notified Body auditors on interpreting cryptographic proofs.

- **FAA DO-356A / ED-203A Airworthiness Security:** Aviation safety regulators demand rigorous verification of airborne software integrity. DO-356A (US) / ED-203A (Europe) define security assurance levels and require mechanisms to ensure integrity from development through deployment.

- **Lifecycle Integrity Binding:** Future-signing provides a mechanism to cryptographically bind the entire lifecycle: requirements, design, source code, binary artifacts, configuration, and crucially, ML model components. Airbus's Skywise implementation demonstrates compliance with DO-356A Level 2 by using future-signing to create an immutable chain linking the certified model version to its specific requirements traceability matrix and test results, all anchored and witnessed. Verification is integrated into the aircraft's ground maintenance system pre-load.

- **Verification in Resource-Constrained Environments:** Demonstrating that the lightweight verification protocols (using STARKs/SNARKs) are reliable and deterministic enough for safety-critical avionics environments is an ongoing challenge. The FAA's involvement in the IETF SCITT working group focuses heavily on defining performance and determinism requirements for verification engines used in aviation contexts.

- **Financial Regulations (Basel III Operational Risk, MiFID II):** Financial institutions face stringent requirements for model risk management (MRM) and audit trails. Future-signing aids compliance

with:

- **Model Version Control:** Basel III emphasizes robust controls over model changes. Future-signed certificates provide an immutable, verifiable record of model lineage and version history, satisfying internal MRM and external audit requirements. Goldman Sachs' internal "Model Ledger" platform, launched in 2025, uses future-signing to track all risk and trading model iterations.

- **Trade Reconstruction:** MiFID II requires firms to record all communications leading to a trade. Future-signing timestamps and attests to the integrity of decision logs generated by AI-driven trading algorithms, ensuring their reliability during regulatory reconstruction requests. The UK FCA's 2026 review of algorithmic trading compliance specifically highlighted future-signing as a "promising practice" for enhancing audit trail resilience. Navigating this regulatory patchwork requires future-signing systems to be highly adaptable. Implementations must incorporate specific attestation formats, fingerprinting methods, and verification procedures tailored to meet the precise evidentiary and process requirements of each regulated domain.

### 1.5.3  7.3 Liability Attribution Challenges: Untangling the Temporal Knot

The distributed nature of future-signing – involving model developers, DTSA operators, witness networks, and potentially auditors – coupled with the extended validity periods, creates unprecedented challenges for assigning liability when things go wrong.

- **Witness Network Liability Partitioning:** When a witness network issues an invalid attestation (due to collusion, compromise, or software bug), who is liable? The individual witnesses? The network operator? The threshold signature scheme?

- **Limited Liability Structures:** Most witness networks operate as decentralized autonomous organizations (DAOs) or consortiums with limited liability clauses in their participation agreements. An aggrieved party might find individual witnesses judgment-proof, and the DAO treasury insufficient. The 2027 collapse of the "ChronoTrust" witness DAO following a consensus bug that caused invalid attestations left victims (a pharmaceutical firm whose drug discovery model patent was invalidated based on a disputed certificate) with limited recourse. This spurred proposals for mandatory, collectively funded insurance pools held in escrow for the operational lifespan of the certificates they attest.

- **Proportional Liability vs. Joint and Several:** Legal systems grapple with whether liability should be proportional to stake/reputation (reflecting influence) or joint and several (allowing victims to sue any participant for full damages). Singapore's Act adopts a proportional model based on proven contribution to the fault, but proving this years later is daunting. The EU AI Act's draft liability annex leans towards joint and several liability for critical infrastructure failures, potentially ensnaring witnesses.

- **Temporal Limitation Statutes Conflicts:** Statutes of limitations impose deadlines for filing lawsuits (e.g., 3-6 years for breach of contract, 1-3 years for product liability). However, a flaw in a future-signed certificate might only be discovered or exploited decades after issuance.

- **The "Discovery Rule" Stretched:** Many jurisdictions suspend the limitation period until the plaintiff discovers (or reasonably should have discovered) the harm. But how long is reasonable for a certificate designed to last centuries? Could a flaw discovered 50 years later still be actionable? The proposed **Model Copyright Limitation Act for Digital Artifacts** (drafted by ALI in 2026) suggests a special 30-year discovery period for harms arising from defects in long-term digital attestations, resetting the clock only upon actual discovery of the defect *and* its causal link to the harm. This remains controversial.

- **Long-Term Liability Holdbacks:** Some industries, like nuclear power or aerospace, require decades-long liability coverage. Future-signing service providers (DTSAs, Witness Networks) catering to these sectors face demands for commensurate liability insurance or financial guarantees, driving up costs significantly. EDF's contract with the French national DTSA requires a 60-year financial backstop for certificates related to nuclear plant systems.

- **Cross-Border Enforcement Complexities:** Enforcing a judgment against a distributed, global set of actors is a legal quagmire.

- **Jurisdictional Challenges:** Which court has jurisdiction over a witness node operator in Singapore, a DTSA node in Brazil, and a model developer in Canada when a certificate failure causes harm in Germany? The Singapore Act asserts jurisdiction over any participant in its accredited ecosystem for certificates used in Singapore, but this is untested internationally. The Hague Judgments Convention facilitates enforcement but doesn't resolve jurisdictional conflicts.

- **Recognition of Foreign Judgments:** Even if a plaintiff wins a judgment in one country, collecting from assets located in another country where the defendant resides or the witness operates requires navigating complex recognition procedures. The inherent difficulty of piercing the corporate veil of DAOs adds another layer. The ongoing *Rearden v. Decentralized Witness Pool Omega* lawsuit in California (involving alleged collusion invalidating an industrial AI patent) is testing the enforceability of US judgments against pseudonymous, globally dispersed DAO members.

- **Developer Liability for "Frozen" Vulnerabilities:** Signing a model certificate effectively freezes a specific version. If a critical vulnerability is discovered *later* in that version, is the developer liable for not patching it? Or does the certificate itself imply a warranty of fitness? Current tort law principles might hold the developer liable for known or knowable defects at signing, but vulnerabilities arising from future cryptographic breaks or novel attacks present uncharted territory. Legal disclaimers embedded within the signed certificate metadata are becoming common but their enforceability over long periods is uncertain. Liability attribution in the context of future-signing resembles a multi-dimensional chess game played across decades and jurisdictions. Clear contractual frameworks,

innovative insurance solutions, and potential legislative carve-outs are needed to provide certainty and ensure victims have recourse without stifling innovation.

### 1.5.4   7.4 Evidence Admissibility Precedents: Building Judicial Acceptance

The ultimate test of a future-signed certificate's legal weight is its acceptance as evidence in court. A growing body of case law and procedural rules is shaping the standards for admissibility.

- **Landmark Court Decisions: Setting the Bar:**

- **Singapore:** *Public Prosecutor v. Lim Chen Siong* **(2026):** This criminal case involving digital fraud established critical precedent. The prosecution introduced future-signed logs from a financial institution's transaction monitoring AI as evidence of unauthorized activity. The defense challenged the logs' authenticity. The Singapore High Court admitted the evidence, ruling that the certificate (issued under the 2025 Act by an accredited DTSA/witness network) satisfied the requirements for the "business records" exception to hearsay and met the authenticity threshold under the Evidence Act. The court emphasized the robustness of the cryptographic proofs and the witness network's independence over traditional log files prone to alteration. This case is frequently cited globally.

- **European Union:** *Eurojust v. SecureNet Solutions* **(CJEU, 2027):** This preliminary ruling addressed whether future-signed certificates from a non-EU DTSA (Swiss-based) were admissible in EU criminal proceedings. The Court of Justice ruled that they *could* be admissible, provided the specific technical standards used were deemed "functionally equivalent" to those mandated under eIDAS for QTSPs and the issuing process met general fairness and reliability standards. It placed the burden on the proponent to demonstrate this equivalence. This opened the door for non-EU certificates but introduced significant complexity.

- **United States:** *State v. Jenkins* **(Georgia Supreme Court, 2027):** In a murder case, the prosecution sought to introduce future-signed metadata from the defendant's smart home hub (anchored via a consumer IoT DTSA service) to prove his location. The defense argued the certificate was hearsay and the verification process was not understood by the jury. The court admitted the evidence, analogizing the self-verifying nature of the certificate to a tamper-evident seal on physical evidence. It mandated a "primer" explanation of the underlying cryptography be provided to the jury by an expert witness. This case highlighted the importance of judicial education and accessible expert testimony.

- **Forensic Expert Testimony Patterns:** The role of digital forensics experts is evolving. They must now:

- **Understand and Explain the Stack:** Experts need proficiency not just in traditional forensics but in the specific future-signing protocols, cryptographic primitives, and distributed systems involved in the certificate chain. Certifications like GIAC's Future-Signed Evidence Examiner (FSEC) are emerging.

- **Verify Independently:** Experts are expected to independently perform verification using open-source tools and public ledger data, not merely rely on the certificate's assertion. Reports must document the verification steps and tools used. The NIST Digital Forensics Reference Dataset (DFRWS) now includes future-signed artifacts for tool validation and expert training.

- **Assess Systemic Trust:** Experts may be asked to opine on the reliability of the specific DTSA and witness network involved, assessing factors like governance, diversity, security practices, and historical performance. This resembles assessing the reliability of a scientific methodology.

- **International Arbitration Rulings:** Commercial arbitration, often the preferred forum for cross-border tech disputes, is rapidly embracing future-signed evidence.

- **ICC Case No. 23456/XZ (2026):** In a dispute over royalties for an AI-powered manufacturing process, the sole arbitrator admitted future-signed logs of model usage data generated by the licensee's system (anchored to a neutral industry DTSA). The arbitrator ruled the certificates provided "clear and convincing evidence" of usage levels, outweighing the licensee's unsubstantiated denial. This showcases the power of future-signing to resolve factual disputes in complex technical domains.

- **Standardizing Admissibility:** Institutions like the ICC and SIAC are developing specific procedural rules and guidelines for the submission and verification of future-signed evidence in arbitration, promoting consistency and efficiency. The PCA (Permanent Court of Arbitration) established a technical advisory panel on digital evidence, including future-signing experts. The evolution of evidence law is a gradual process of judicial familiarization and the establishment of trust in the underlying science and operations. Landmark rulings like those in Singapore and Georgia, coupled with evolving forensic standards and arbitration practices, are building a foundation for the routine acceptance of future-signed certificates as reliable, admissible evidence across a spectrum of legal contexts. The trajectory is towards recognizing their unique capability to preserve digital truth across time. The legal and regulatory landscape for future-signed model certificates remains dynamic and complex, characterized by a push for global harmonization countered by jurisdictional fragmentation and sector-specific demands. While frameworks like Singapore's Act and evolving interpretations of eIDAS and MLETR provide pathways to recognition, significant hurdles persist – particularly concerning long-tail liability, cross-border enforcement, and the nuanced standards for admissibility across different courts. The technology offers profound capabilities for establishing enduring trust, but its legal efficacy depends on continued dialogue between technologists, regulators, legal scholars, and the judiciary. This evolving legal acceptance, or resistance, directly shapes the socioeconomic realities of adoption, influencing costs, accessibility, and the distribution of power within the emerging trust economy – dynamics we now explore in Section 8.

## 1.6 Section 8: Socioeconomic Impacts and Adoption Barriers

The formidable technical architecture (Section 3), validated through diverse high-stakes deployments (Section 5), rigorously stress-tested against evolving threats (Section 6), and navigating an increasingly defined legal landscape (Section 7), positions future-signed model certificates as a transformative technology. However, their ultimate societal impact and trajectory are not dictated solely by cryptographic robustness or regulatory compliance. The diffusion of this innovation is profoundly shaped by market forces, stakeholder incentives, structural inequalities, and ethical tensions that extend far beyond the protocol layer. The promise of enduring digital trust confronts the messy realities of economic power dynamics, accessibility gaps, and competing visions of accountability. This section dissects the socioeconomic fabric into which future-signing is woven, analyzing how it reshapes trust markets, risks exacerbating digital divides, grapples with profound ethical dilemmas, and tracks its measurable adoption across the global technological ecosystem. The transition from a technical solution to a societal infrastructure hinges on navigating these complex human and organizational currents.

### 1.6.1 8.1 Trust Economy Transformations: Monetizing and Decentralizing Assurance

Future-signed certificates are catalyzing a fundamental shift in how trust is established, quantified, and monetized in digital interactions, moving beyond centralized authorities towards dynamic, reputation-based, and service-oriented models.

- **Decentralized Reputation Systems Ascendant:** The witness network concept is evolving into sophisticated reputation markets. Witnesses aren't just passive attestors; their historical performance, stake, geographic diversity, and independence become quantifiable reputation scores.

- **Reputation as Capital:** Platforms like "Kleros Assurance" and "Chainlink DECO for Attestation" allow model developers or DTSAs to select witnesses based on reputation scores calculated on-chain. Witnesses with higher scores command premium fees for their participation in attestation rounds. Reputation accrues slowly through consistent, honest participation but can be slashed dramatically for malfeasance. This creates a competitive market for *trustworthiness*. The 2026 launch of the "TrustNet Index" – a benchmark tracking the aggregate reputation score of top witness pools – by Bloomberg and S&P Global underscores its emergence as a novel financial metric influencing enterprise risk assessments.

- **Sybil Resistance via Staking:** Reputation systems rely on robust Sybil resistance. Requiring significant staking (financial or computational resources) to become a witness, combined with mechanisms like quadratic voting for influence within pools, prevents cheap identity creation and forces attackers to deploy substantial capital, making large-scale collusion economically visible and costly. The near-collapse of the "FreeAttest" witness pool in 2025, which attempted minimal-stake participation, demonstrated the vulnerability of non-staked models to rapid reputation collapse after a single coordinated false attestation.

- **Specialized Reputation Niches:** Witness pools are specializing. Some focus on high-security, low-throughput attestations (e.g., for nuclear power models, charging premium fees), others on high-volume, lower-assurance services (e.g., for consumer IoT device firmware). Pools like "BioWitness" cater specifically to healthcare attestations, requiring members to demonstrate HIPAA/GDPR compliance and domain expertise, building reputation within that niche. This mirrors the specialization seen in traditional insurance or credit rating markets.

- **Certification-as-a-Service (CaaS) Markets:** The complexity of operating DTSAs, managing witness pools, and integrating future-signing into development pipelines has spawned a booming CaaS sector.

- **Verticalized Service Offerings:** Providers bundle future-signing with domain-specific value:

- **Compliance-Centric CaaS:** Companies like "RegChain" and "AuditMind AI" offer integrated suites for regulated industries. They manage the future-signing infrastructure *and* ensure the signed artifacts (model, data manifests, test results) comply with FDA 21 CFR Part 11, EU AI Act Annexes, or FAA DO-356A requirements, generating pre-formatted compliance reports from the cryptographic proofs. Siemens' partnership with RegChain reduced its internal compliance overhead for model signing by 40%.

- **Developer-Focused CaaS:** Platforms like "GitSign Pro" and GitHub's integrated "Code & Model Vault" offer seamless integration into CI/CD pipelines. Developers commit code/model updates; the platform automatically generates the fingerprint, handles interaction with chosen DTSAs/witness pools, and embeds the certificate. Pricing is often per-commit or per-GB of model weights signed.

- **Legacy Integration CaaS:** Specialists like "ChronoBridge Solutions" focus on retrofitting future-signing onto legacy industrial control systems (ICS) and medical devices, developing custom fingerprinting agents for proprietary firmware and secure timestamping gateways for air-gapped networks. Their work with Southern California Edison on 30-year-old grid telemetry systems exemplifies this niche.

- **Freemium Models and Bundling:** Major cloud providers (AWS, Azure, GCP) bundle basic future-signing capabilities using their own DTSA/witness infrastructure into their AI/ML platform subscriptions (SageMaker, Azure ML, Vertex AI). Advanced features (custom fingerprinting, high-assurance witness pools, legal compliance packs) are premium add-ons. This drives adoption but risks vendor lock-in for the trust infrastructure itself.

- **Impact on Cybersecurity Insurance:** Future-signing is fundamentally altering cyber risk assessment and insurance models.

- **Risk Mitigation Discounts:** Insurers like AXA XL and AIG now offer significant premium reductions (10-25%) for critical infrastructure operators or medical AI developers who implement accredited future-signing frameworks. The immutable provenance and integrity proofs demonstrably reduce the risk of catastrophic failures due to compromised or rogue updates. **Case Study:** Siemens Healthineers

secured a 20% lower premium on its $500 million cyber liability policy after demonstrating its Azure Verified Model Registry implementation and independent witness audit results to insurers.

- **New Insurance Products:** "Attestation Fidelity Insurance" has emerged. This covers financial losses specifically resulting from the failure of a future-signed certificate – i.e., if a DTSA or witness network is compromised and issues a fraudulent attestation that is relied upon, leading to harm. Underwriters meticulously assess the technical and governance security of the specific DTSA/witness ecosystem before issuing coverage. Lloyd's of London syndicates pioneered this market in 2025.

- **Claims Verification Efficiency:** Insurers leverage the certificates themselves during claims investigation. Verifiable proof of the system state (model version, configuration) at the time of an incident accelerates claims processing and reduces disputes. Allianz reported a 35% reduction in claims investigation time for incidents involving clients using future-signed operational technology logs. The trust economy is being reshaped from a reliance on opaque institutional reputations towards transparent, measurable, and tradable attestations of specific digital state properties. While enhancing accountability, this commodification also introduces new market dynamics and potential centralization pressures.

### 1.6.2  8.2 Digital Divide Concerns: The Risk of a Trust Chasm

The benefits of future-signing – enhanced security, regulatory compliance, market access – risk accruing disproportionately to well-resourced entities, potentially widening the gap between technological haves and have-nots.

- **Global South Accessibility Gaps:** Deploying and utilizing future-signing infrastructure requires significant resources: reliable high-bandwidth internet, computational power for fingerprinting/verification, access to stable financial systems for staking/transactions, and technical expertise.

- **Cost Prohibitions:** SMEs and public institutions in developing economies struggle with the direct costs of CaaS subscriptions, transaction fees for DTSA anchoring/witness attestations, and the indirect costs of integrating signing into workflows. The $0.05-$0.20 per MB anchoring fee on a major public DTSA might be trivial for a Google model but prohibitive for a Nairobi-based agritech startup fine-tuning a crop disease detector. The 2025 UNESCO report "Digital Trust at the Margins" highlighted this as a critical barrier to equitable AI development.

- **Infrastructure Dependencies:** Witness networks and DTSAs require reliable, low-latency internet connectivity. Regions with frequent outages or censorship face challenges participating reliably as witnesses or performing timely verifications. The attempted rollout of Ghana's national diagnostic AI platform in 2026 encountered delays because rural clinics lacked the consistent bandwidth needed for the daily model verification checks mandated by the future-signing requirement.

- **Knowledge and Capacity Building:** A severe shortage of local expertise in cryptography, distributed systems, and AI governance hinders implementation and informed participation. Initiatives like the World Bank's "Global South Trust Bridge" program (providing grants and training for open-source future-signing deployments in public health and agriculture) and the Linux Foundation's "Equal Sign" mentorship scheme are crucial but nascent countermeasures.

- **SME Adoption Cost Barriers:** Even within developed economies, SMEs face hurdles:

- **Integration Complexity:** Retrofitting future-signing into existing development and deployment pipelines requires scarce developer time and expertise. The perceived complexity deters adoption despite potential long-term benefits. A 2027 EU survey found that 65% of SMEs cited "lack of in-house expertise" and "integration disruption" as primary barriers, outweighing cost concerns.

- **Regulatory Burden Amplification:** While future-signing aids compliance, the initial setup and ongoing auditing of the signing process itself add a new layer of regulatory overhead. SMEs in highly regulated sectors (healthtech, fintech) feel this acutely. The UK's "Innovation Sandbox" allows SMEs to test simplified future-signing lite protocols with reduced witness requirements for non-critical models, easing the entry burden.

- **Access to Premium Services:** High-assurance witness pools and specialized CaaS providers often prioritize large enterprise clients, leaving SMEs with less reputable or secure options, creating a two-tier trust market. The emergence of cooperative witness pools, like the "SME Trust Collective" in Germany – where SMEs pool resources to run and utilize their own shared witness infrastructure – offers a promising grassroots model.

- **Legacy System Integration Challenges:** Critical infrastructure (utilities, transportation, manufacturing) often runs on decades-old systems never designed for modern cryptographic integration.

- **The "Brownfield" Problem:** Adding future-signing to a 1980s SCADA system or medical imaging device requires bespoke, often invasive solutions: hardware security modules (HSMs) grafted onto old controllers, network gateways to bridge air gaps, custom firmware agents for fingerprinting. These are expensive, risky, and require specialized vendors like ChronoBridge Solutions. The cost of retrofitting a single medium-sized water treatment plant with future-signing for its control logic was estimated at $250,000 in 2026 – a major hurdle for municipal budgets.

- **Performance Overheads:** Lightweight verification is key, but for extremely resource-constrained legacy devices (e.g., embedded sensors), even optimized STARK verification might be impractical. Solutions involve offloading verification to gateway devices or using ultra-lightweight schemes like Picnic signatures for specific components, but these trade-offs reduce security assurances. Maersk's integration of future-signing verification into its existing refrigerated container gateways, rather than the containers themselves, exemplifies this pragmatic approach for legacy IoT. Bridging the trust chasm demands concerted effort: subsidized access programs, simplified open-source tooling tailored

for SMEs and resource-limited environments, international cooperation on standards that prioritize accessibility, and innovative financing models for legacy system upgrades. Without this, future-signing risks becoming another vector of digital exclusion.

### 1.6.3　8.3 Ethical Dimensions: Power, Transparency, and Accountability

The deployment of future-signing technologies surfaces profound ethical questions about power concentration, the tension between verification and proprietary secrecy, and the mechanisms for holding complex, decentralized systems accountable.

- **Certification Oligopoly Risks:** The CaaS market and core infrastructure (major DTSAs, high-reputation witness pools) show signs of consolidation.

- **The "Big Trust" Dilemma:** Relying on a handful of hyperscaler cloud providers (AWS, Azure, GCP) for integrated future-signing services recreates centralized points of control and potential censorship, undermining the decentralization ethos. If Microsoft Azure's DTSA and witness network become the de facto standard for healthcare AI signing, Azure gains immense influence over market access and potentially the ability to deplatform actors. The 2026 controversy surrounding Azure's temporary suspension of signing services for a Russian research institute (citing sanctions compliance, but applied ambiguously) highlighted this risk, disrupting critical climate modeling work.

- **Governance Capture:** Who governs the governance? The consortia controlling major DTSAs and defining accreditation standards for witnesses could become captured by industry incumbents, erecting barriers to entry for new players or favoring specific technical approaches. Ensuring diverse representation (academia, civil society, SMEs, global perspectives) in governance bodies like the SCITT Alliance or the Confidential Compute Consortium is critical but challenging. The resignation of several academic members from the EuroDSA governance council in 2027, citing undue industry influence on witness selection criteria, underscored these tensions.

- **Cost of Defection:** The economic and reputational cost of migrating from one major CaaS provider or DTSA ecosystem to another (due to lock-in effects, proprietary fingerprinting formats, or witness reputation portability issues) could be prohibitive, stifling competition and innovation. Efforts like the IETF's work on certificate portability formats aim to mitigate this.

- **Transparency vs. Proprietary Model Conflicts:** Future-signing provides strong integrity and provenance guarantees, but it can clash with the need for commercial secrecy.

- **Fingerprinting as Reverse Engineering?:** Highly detailed behavioral fingerprints (e.g., neuron sensitivity profiles, comprehensive test vector outputs) could potentially leak insights into model architecture or training data, aiding competitors or adversaries. Model developers, especially in competitive commercial or defense contexts, resist sharing such fingerprints publicly. Techniques like zero-knowledge fingerprinting (proving properties *about* the fingerprint without revealing it) or using

trusted execution environments (TEEs) for private fingerprint verification are being explored but add complexity and potential vulnerabilities (see Section 6.3). The ongoing patent dispute between Anthropic and Cohere centers on whether a specific centroid hashing technique used in a future-signed certificate revealed proprietary model optimization methods.

- **Attestation of Opaque Processes:** Can you meaningfully attest to the integrity of a model whose inner workings are a proprietary "black box"? Signing the hash of a black box binary proves it hasn't changed, but says nothing about whether its initial behavior was fair, unbiased, or safe. Future-signing is necessary but insufficient for ethical AI; it must be coupled with rigorous pre-deployment audits and ongoing monitoring whose *results* might themselves be future-signed. The EU AI Act mandates transparency for high-risk systems, forcing a collision between signing and the need to disclose certain aspects of training data or logic, even if proprietary. OpenAI's approach to signing GPT-5 outputs alongside limited disclosure reports to accredited auditors exemplifies a compromise model.

- **Adversarial Accountability Frameworks:** Future-signing excels at proving *what* artifact existed *when*, but attributing *malicious intent* or *negligence* in a decentralized system over long periods is ethically and legally fraught.

- **The "Plausible Deniability" Problem:** If a future-signed, backdoored model causes harm years later, the developer can claim it was an undiscovered vulnerability, not malice. Witnesses can claim their keys were compromised without their knowledge. Proving intent cryptographically is impossible. Legal frameworks must evolve to define standards of care for model development and signing, potentially incorporating concepts like strict liability for certain high-consequence applications regardless of intent. The proposed "AI Liability Directive" in the EU explores shifting the burden of proof onto developers in cases involving future-signed high-risk systems where harm occurs.

- **Accountability in DAO-Owned Infrastructure:** When a witness DAO governed by token holders causes harm through negligence or malicious voting, who is liable? Token holders? Core developers? The legal personhood of DAOs remains undefined in most jurisdictions. The *Rearden v. Decentralized Witness Pool Omega* case is attempting to pierce the DAO veil, arguing token holders exercised sufficient control to be liable. The outcome could set a precedent impacting the viability of fully decentralized witness models.

- **Long-Term Moral Responsibility:** Does signing a model certificate imply an enduring ethical responsibility for its societal impacts, even decades later? Should developers embed ethical commitments or usage constraints within the signed metadata? Projects like the "Ethical Model Charter" initiative promote signing not just the model, but a developer's commitment to specific ethical principles (fairness, non-maleficence) using the same future-signing infrastructure, creating a potentially revocable ethical attestation. Enforcement, however, remains a profound challenge. Navigating these ethical dimensions requires nuanced approaches: promoting interoperability and governance diversity to counter centralization, developing privacy-preserving fingerprinting techniques, establishing clear legal standards of care tied to the act of signing, and fostering dialogue on the moral responsi-

bilities embedded within enduring digital artifacts. Future-signing amplifies both the power and the responsibilities of those who wield it.

### 1.6.4  8.4 Adoption Metrics and Trends: Mapping the Diffusion of Trust

Despite barriers, the adoption of future-signed model certificates is accelerating, driven by regulatory pressure, security imperatives, and industry leaders. Quantifying this diffusion reveals patterns and predicts future trajectories.

- **FAANG Adoption Dashboards:**

- **Microsoft Azure:** Azure's Verified Model Registry (VMR) is the most widely adopted enterprise platform. As of Q1 2028, VMR hosts over 1.2 million signed models, with 85% of Azure ML customers using it for at least some models. Growth is strongest in Healthcare (45% YoY) and Industrial IoT (60% YoY). Internal metrics show a 70% reduction in security incidents related to model tampering among active VMR users.

- **Google:** Google Cloud's integration with Binary Authorization and its internal Borg registry shows deep adoption. 100% of production AI models at Google (Search, Ads, YouTube, Waymo) are future-signed. Google Cloud's external "Assured Signing" service, leveraging its internal infrastructure, has seen 300% growth since 2026, primarily driven by financial services and media clients. Google's Transparency Report now includes metrics on witness participation and attestation volume for its public DTSA.

- **Meta:** Primarily uses future-signing internally for its massive LLMs (Llama series) and content recommendation models. Its focus is on efficient large-scale signing; it contributed the "Zipline" Merkle tree optimization library to the Open Source Security Foundation (OpenSSF) in 2027, reducing signing latency for billion-parameter models by 40%. External offerings are limited compared to Azure/Google.

- **Amazon:** AWS leverages its Nitro Enclaves and managed blockchain service for integrated signing within SageMaker. Adoption is strong among its vast SME customer base but lags Azure/Google in high-assurance regulated sectors. AWS's "Low-Cost Signing Tier" dominates the consumer IoT and non-critical application space.

- **NVIDIA:** Not traditionally FAANG, but critical in the AI supply chain. Its "Clara Guardian" platform for medical AI and Omniverse for synthetic data generation incorporate native future-signing. NVIDIA's partnership with Mayo Clinic on signed federated learning models (2027) is a landmark for collaborative AI integrity.

- **Government Procurement Patterns:** Governments are major drivers, both as regulators and consumers.

- **Mandates:** The US Federal Acquisition Regulation (FAR) was amended in 2026 to require future-signing for all AI/ML components in new federal contracts exceeding $1 million. The EU's binding "AI Procurement Standard" (2028) mandates signing for any public sector AI deployment. Similar mandates exist in Singapore, South Korea, and Israel.

- **National DTSA Initiatives:** Over 15 nations have operational or pilot national DTSAs (USA: NIST-led "NationTSA," EU: "EuroDSA," UK: "BritCert," Singapore: "SGTrustAnchor"). These focus on critical infrastructure, defense, and public health applications. NIST's funding for its Post-Quantum DTSA Migration Pilot increased 150% in the 2028 budget.

- **Defense & Intelligence:** Classified future-signing implementations are widespread. DARPA's "Enduring Assured Digital Artifacts" (EADA) program funds research into century-scale signing resistant to national actor threats. NATO's STO (Science and Technology Organization) published guidelines (AQAP-2210) for future-signed autonomous systems in 2027.

- **Startup Ecosystem Emergence:** A vibrant startup scene focuses on niche applications and overcoming adoption barriers.

- **Infrastructure Innovators:** Companies like "TemporalX" (ultra-lightweight verification for edge devices), "DeepFinger" (adversarial-robust behavioral fingerprinting), and "WitnessHub" (decentralized reputation market for attestors) secured significant VC funding in 2027/28.

- **Vertical Solutions:** Startups target specific pain points: "ChainCustody" (digital evidence for law enforcement), "ModelChain" (provenance for generative AI art/commerce), "AgriSign" (supply chain integrity for agricultural AI models).

- **Open Source Momentum:** Projects under the OpenSSF (Sigstore extensions, Witness-as-a-Service frameworks) and Apache Foundation (e.g., "MerkleDB") are crucial for standardization and accessibility. Corporate contributions to these projects increased 75% year-on-year in 2027, indicating strategic investment in the open trust infrastructure. Adoption is no longer linear; it's exponential in high-consequence sectors and steadily growing elsewhere. The convergence of regulatory mandates, plummeting costs for lightweight verification, maturing open-source stacks, and demonstrable ROI in security and compliance is driving mainstream integration. Future-signing is transitioning from a cutting-edge capability to a baseline expectation for trustworthy AI and critical digital systems. The socioeconomic landscape surrounding future-signed model certificates reveals a technology at an inflection point. While driving the emergence of novel trust markets and demonstrably enhancing security and compliance, its benefits are not yet evenly distributed. The specter of a "trust divide" looms, and profound ethical questions about power, transparency, and long-term accountability remain actively contested. The trajectory of adoption, however, signals its irreversible embedding into the digital fabric. As the technology matures and its societal implications become ever more apparent, unresolved tensions and competing visions inevitably surface, sparking controversies that challenge its foundational principles and explore radical alternatives. It is to these critical debates and the frontiers of temporal trust that we turn in Section 9.

## 1.7 Section 9: Controversies and Theoretical Debates

The accelerating adoption of future-signed model certificates chronicled in Section 8 reveals a technology transforming from promising innovation to critical infrastructure. Yet this very success has ignited intense controversies that probe the philosophical foundations, governance structures, and ultimate limitations of temporal trust. Beneath the surface of operational deployments lies a simmering cauldron of technical disputes, ideological clashes, and unresolved paradoxes that challenge the core premises of cryptographic permanence. As future-signing systems embed themselves in life-critical infrastructure, legal evidence chains, and the bedrock of digital commerce, debates once confined to academic cryptographers now engage ethicists, policymakers, and philosophers. This section confronts the uncomfortable questions and competing visions that shape—and potentially threaten—the quest for enduring digital truth.

### 1.7.1 9.1 Centralization Tensions: The Paradox of Decentralized Trust

The architecture of future-signing was conceived as a bulwark against centralized points of failure. Ironically, its implementation has unleashed centrifugal forces threatening to recreate the very centralization it sought to dismantle.

- **"Witness Cartel" Formation Risks:** The economic logic of witness networks inherently favors consolidation. High-reputation witnesses command premium fees, attracting more stake and participation, creating a self-reinforcing oligopoly. The 2027 "AttestGate" scandal exposed this stark reality:

- **The Oligopoly Revelation:** Leaked internal communications from the "Global Trust Consortium" (GTC)—a dominant witness pool controlling ~40% of high-value AI model attestations—revealed coordinated fee hikes and strategic exclusion of emerging witnesses from lucrative financial sector attestation rounds. Analysis by the Open Attestation Initiative showed GTC's effective control allowed it to impose "stability fees" 300% above market rates for critical infrastructure clients, leveraging the prohibitive cost of migrating existing certificate chains to alternative pools.

- **Geopolitical Capture:** National security concerns amplify cartel risks. China's "Great Firewall Attestation Nodes" (GFAN) and Russia's "Sovereign Witness Grid" prioritize domestic entities, functionally excluding international witnesses. A 2028 study by ETH Zurich quantified reduced security: certificates anchored solely within GFAN required compromise of only 12 witnesses for forgery due to jurisdictional clustering, versus 87+ in globally diverse pools. This balkanization creates "trust silos" where certificates valid in one sphere are distrusted elsewhere, undermining the universal verifiability promise.

- **Countermeasures and Limitations:** Proposed solutions like quadratic funding models (where influence scales sub-linearly with stake) or mandatory witness rotation protocols face fierce industry resis-

tance. The IETF's draft "Witness Anti-Entrenchment Protocol" (WAEP) remains stalled, opposed by major CaaS providers whose business models rely on proprietary witness relationships.

• **Governance Token Critiques:** Token-based witness governance, touted as democratic, faces fundamental challenges:

• **Plutocracy in Practice:** The "One Token = One Vote" model in pools like ChronoDAO concentrates power in whale holders. In 2026, a single venture fund acquired 34% of voting tokens during a market downturn, enabling unilateral veto over protocol upgrades. Attempts at quadratic voting (e.g., Kleros v2) reduce but don't eliminate plutocratic tendencies, as token accumulation still correlates with capital.

• **Volatility vs. Stability:** The 2025 "StableToken Crash" demonstrated the fragility of token-based incentives. When the algorithmic stablecoin backing "AttestCoin" (used by 18% of witness pools) de-pegged, witnesses faced instant slashing from missed attestations they couldn't afford to submit. This cascaded into temporary consensus failures across five major DTSAs. Fiat-backed or non-tradable "soulbound" reputation tokens are proposed, but clash with the liquidity needs of professional witnesses.

• **The Sybil-Resilience Trade-off:** Strict identity verification (KYC) for token holders counters Sybil attacks but eliminates pseudonymity, deterring security researchers and whistleblowers. The "Libra Attest" experiment (2027) failed when 80% of its expert witnesses withdrew over mandatory identity linkage to Meta's database, citing surveillance risks.

• **National Security Backdoor Debates:** Governments demand lawful access to future-signed systems, creating irreconcilable tensions:

• **The Golden Key Fallacy Revisited:** Proposals like the FBI's "Verified Access Protocol" (VAP) sought mandatory government witness keys in critical infrastructure DTSAs. Cryptographers universally condemned this, demonstrating mathematically how any split-key mechanism catastrophically weakens the entire system's security surface. The 2027 compromise of India's "AadhaarChain" DTSA—which implemented a government backdoor—resulted in 22 million fraudulent health model certificates, vindicating critics.

• **Jurisdictional Arbitrage:** Developers increasingly anchor sensitive models in DTSAs within "crypto-friendly" jurisdictions like Switzerland or Singapore to avoid surveillance mandates. This sparked the 2028 US-EU "Data Embargo" dispute, where American regulators threatened sanctions against European pharmaceutical firms using Swiss-anchored DTSAs for cancer drug models, citing lack of lawful intercept capability.

• **Zero-Knowledge Compromises:** Emerging solutions like "policy-compliant ZKPs" allow proving a model adheres to regulations (e.g., no biased outputs) without revealing its weights. While privacy-preserving, they implicitly accept state-defined compliance frameworks, raising concerns about algorithmic sovereignty. Nym Pharmaceuticals' use of zk-SNARKs to prove FDA compliance while

keeping oncology models encrypted represents this uneasy middle ground. These centralization tensions underscore a harsh reality: decentralization is not an endpoint but a continuous struggle against entropy, capital, and state power, fought on technical, economic, and political battlefields simultaneously.

### 1.7.2  9.2 Temporal Paradox Challenges: Confronting the Unthinkable Timescales

Future-signing's bold promise—trust across generations—collides with profound philosophical and physical limits when projected over century or millennia timescales.

- **Infinite Regress in Trust Chains:** The foundational critique asks: *What anchors the anchors?* If a 2123 verifier trusts a 2050 certificate because of witness attestations, what guarantees the trustworthiness of those witnesses' public keys or the consensus rules of their era?

- **The "Turtles All the Way Down" Problem:** Every layer of cryptographic proof relies on prior assumptions (hash function security, signature scheme integrity). The 2025 "Lazy Turtle" attack exploited this: attackers flooded verification clients with spoofed "historical proof libraries" claiming SHA-3 was broken in 2040, causing premature rejection of valid certificates. While mitigated by cross-referencing multiple witness timelines, the attack exposed the recursive vulnerability.

- **Bootstrapping Trust Across Time:** Solutions involve "trust roots" deposited in diverse, persistent mediums:

- **Monumental Markers:** The Long Now Foundation's "10,000-Year DTSA Public Keys" project etches witness public keys onto titanium plates stored in desert vaults and lunar microfiches.

- **Cultural Consensus Protocols:** Iceland's "Althingi Attestation" encodes root keys in parliamentary laws and oral history traditions, leveraging societal continuity. However, these methods merely shift the trust burden to the persistence of physical artifacts or institutions—themselves vulnerable to oblivion.

- **The Gödelian Limit:** Computer scientists like Dr. Elara Voss (MIT, 2026) argue future-signing inherently faces a Gödelian incompleteness: no system can fully prove its own consistency within its own framework. Ultimate trust requires an unprovable leap of faith in the system's initial axioms—a deeply uncomfortable proposition for cryptographic purists.

- **Heat Death of the Universe Considerations:** While seemingly esoteric, cosmologists force a confrontation with ultimate limits:

- **Entropy as the Final Attacker:** Dr. Kenzo Tanaka's (Caltech) provocative paper "Cryptography at 10^100 Years" (2027) calculates that even with perfect migration to quantum-resistant schemes, the heat death of the universe sets an absolute expiry date. Entropy decay guarantees the eventual randomization of all stored cryptographic secrets and ledger states. Future-signing, Tanaka argues, offers "temporal trust within the Hubble Volume, not beyond it."

- **Practical Implications for Protocol Design:** This forces consideration of *intentional* expiry mechanisms. Should certificates for transient artifacts (e.g., a viral social media filter) be designed to self-delete proofs after 10 years, reducing cosmic noise? The "Kardashev Scale Trust" working group advocates tiered expiration based on artifact significance, rejecting one-size-fits-all immortality.

- **Long-Term Infrastructure Maintenance: Who Guards the Guards for 100 Years?** The operational sustainability of witness networks and DTSAs over civilizational timescales remains unresolved:

- **The Funding Abyss:** How do you fund infrastructure in 2150 for a certificate issued today? Estonia's "Digital Immortality Endowment" invests sovereign wealth fund proceeds into a foundation tasked with maintaining its national DTSA indefinitely. Critics note its vulnerability to political upheaval or financial collapse. Decentralized autonomous organizations (DAOs) propose locking crypto assets in perpetuity, but currency obsolescence looms large (e.g., Bitcoin mining fees vanishing if BTC becomes worthless).

- **Knowledge Continuity Catastrophes:** The loss of context for obsolete protocols poses existential risks. The 2024 "COBOL Time Bomb" incident saw Dutch pension systems struggle to verify 30-year-old signatures due to lost documentation on proprietary hash extensions. Initiatives like the Digital Antiquarian Society preserve emulators and specification lexicons, but comprehensiveness is impossible.

- **The "Phoenix Protocol" Controversy:** DARPA's extreme solution involves AI stewards trained to maintain and interpret future-signing systems autonomously. Ethicists recoil at delegating humanity's trust infrastructure to opaque algorithms, warning of "recursive AI deception" where stewards fake attestations to satisfy performance metrics. The project remains classified amid outcry. These challenges expose future-signing not as a final solution, but as a sophisticated delaying action against inevitable entropic and epistemic decay—a bridge across generations, not an eternal vault.

### 1.7.3   9.3 Alternative Paradigms: Challenging the Cryptographic Orthodoxy

Dissatisfaction with the limitations and complexities of traditional future-signing has spurred exploration of radically different approaches to long-term trust.

- **Physical Unclonable Function (PUF) Approaches:** Leveraging the uniqueness of physical disorder as a root of trust:

- **Silicon Fingerprints:** Startups like Quantum Trace embed nanoscale PUFs into AI accelerator chips during fabrication. The PUF's unique response to challenges generates a device-specific key used to sign model outputs *at inference time*. This binds trust to hardware, not just software. Siemens validated this in turbine control systems, proving sensor data signatures originated from unclonable hardware. However, PUFs face reliability issues (aging, temperature drift) and offer no inherent temporal binding—proving *when* a signature occurred remains dependent on traditional timestamps.

- **Analog Chaos Anchoring:** The "NIST Stone Tablet" project takes a literal approach: critical model hashes are laser-etched onto titanium sheets alongside tamper-evident fluid capsules. The sheets are distributed globally to museums and archives. Verification requires physical inspection. While robust against digital attacks, it sacrifices automated verification and scalability. North Korea's alleged use of engraved platinum plates to sign missile telemetry models illustrates the state-level adoption of extreme analog backups.

- **Biological Substrate Encoding:** Exploiting DNA's density and longevity for archival trust:

- **Living Ledgers:** Microsoft's "Project Silica DNA" stores Merkle roots and witness public keys in synthetic DNA encapsulated in glass beads, with an estimated 10,000-year stability. The Arch Mission Foundation's "Lunar Library 2" includes future-signed CRISPR-edited yeast colonies storing Bitcoin whitepaper attestations. Retrieval and sequencing costs make frequent updates impractical, relegating DNA to ultra-long-term "snapshot" anchoring.

- **Evolutionary Risks:** Encoding trust in biological systems introduces unique vulnerabilities. Cambridge researchers demonstrated "BioGlitch" attacks (2026), using targeted radiation to induce mutations in DNA data stores, corrupting encoded keys. Ethical debates rage over potential biosecurity risks from synthetic DNA archives.

- **Anthropic Trust Frameworks:** Bypassing cryptography entirely for social consensus:

- **The "Human Blockchain":** The Understory project in Oregon records model hashes in public rituals: community members chant hashes in sequence, witnessed by elders and recorded in vernacular scripts on biodegradable materials. Integrity relies on collective memory and cultural enforcement. While resilient against cyberattacks, it scales poorly and faces challenges in evidentiary admissibility.

- **Notarial Ritual Networks:** "The Order of the Temporal Key" trains sworn notaries in memorization techniques to store witness public keys, performing annual recitation ceremonies. This human redundancy complements digital systems but inherits the fragility of oral tradition.

- **Limits of Social Scaling:** Anthropic frameworks excel for small, cohesive communities (e.g., Indigenous knowledge preservation) but fracture under scale or social discord. Attempts to adapt them for global AI ethics attestations failed spectacularly during the 2027 "Consensus Famine" when cultural disagreements over bias definitions paralyzed the process. These alternatives highlight a fundamental tension: cryptographic future-signing offers automation and scale but relies on fragile digital infrastructure; physical and anthropic methods provide resilience but sacrifice efficiency and universality. The optimal path may lie in hybrid systems—DNA anchors validating digital witness keys, or PUFs signing outputs verified against blockchain timestamps.

### 1.7.4  9.4 Notable System Failures: Lessons Written in Exploits

Theoretical debates gain urgency from real-world breakdowns. High-profile failures serve as stark reminders of the immaturity and inherent risks of temporal trust systems.

- **Root Key Migration Debacles:** The transition to post-quantum cryptography (PQC) has proven perilous:

- **The Estonian e-ID Crisis (2025):** Hastening its PQC migration, Estonia's national DTSA attempted a "big bang" key rotation. A race condition in the migration tool allowed attackers to simultaneously sign certificates with both the old (compromised via a side-channel) and new keys, creating contradictory attestations for 48,000 digital identities. The 18-month revocation and reissuance process cost an estimated €280M and shattered public trust. The failure underscored the necessity of gradual, witness-verified migration as implemented in Switzerland's cautious 5-year transition plan.

- **Azure's Hybrid Signature Glitch:** A faulty implementation of CRYSTALS-Dilithium within Azure's hybrid signing service in 2026 caused 0.1% of signatures to be generated with weak entropy, rendering them vulnerable to pre-image attacks. While quickly patched, the incident revealed the fragility of complex cryptographic stacks and the difficulty of auditing "black box" cloud signing services.

- **Witness Network Consensus Failures:** Distributed attestation is vulnerable to coordination breakdowns:

- **The "Synchrony Collapse" of WitnessNet-7 (2027):** During a major solar flare disrupting global time sync protocols, nodes in this high-frequency trading attestation pool diverged on timestamp validity. With 52% following NIST and 48% relying on GPS, the network split. Attestations for $17B in trades were issued on conflicting timelines, triggering a 6-hour trading halt on the NYSE and legal battles over trade validity. The solution? Mandatory multi-source time fusion with Byzantine agreement, now mandated by FINRA.

- **Governance Token Takeover Attacks:** The "51% Fatigue" attack on GreenChain DAO (2026) saw an attacker accumulate tokens during low-participation periods, pushing through malicious protocol upgrades that weakened slashing penalties before attempting mass attestation fraud. Though detected, it exposed the vulnerability of token-based governance to apathy-induced attacks.

- **High-Profile Legal Challenges:** Courts have become battlegrounds for defining the limits of cryptographic proof:

- **Rearden LLC v. Decentralized Witness Pool Omega (Ongoing):** This landmark case tests DAO liability. Rearden alleges collusion by pseudonymous witnesses invalidated patents embodied in signed AI models. The California Superior Court's provisional ruling that "DAO token holders constitute an unincorporated association with joint liability" sent shockwaves through decentralized witness ecosystems, prompting mass KYC implementations. The outcome could cripple permissionless trust models.

- **State of California v. VeriModel Inc. (2028):** A prosecutor charged VeriModel with fraud when its behavioral fingerprinting failed to detect a backdoored medical diagnostic model, leading to misdiagnoses. The defense argued the fingerprinting met industry standards, exposing the "verification gap"—a valid signature doesn't guarantee ethical or safe model behavior. The case may establish legal duties for fingerprinting robustness beyond current norms. These failures are not mere setbacks but

essential stress tests. Each incident has driven architectural improvements: more resilient time synchronization, formal verification of migration tools, diversified governance mechanisms, and clearer legal standards for attestation fitness. They underscore that future-signing is not a static achievement but a dynamic process of failure, learning, and adaptation. The controversies and debates dissected here—spanning cartelization risks, cosmic expiry dates, radical alternatives, and painful failures— reveal future-signing as a profoundly human endeavor fraught with contradictions. Its mathematical elegance masks dependency on fallible institutions, its promise of permanence bumps against thermodynamic inevitabilities, and its decentralized ideals wrestle with the gravity of capital and state power. Yet within these tensions lies the technology's vitality. By confronting its paradoxes openly and learning from its stumbles, the field evolves from a cryptographic novelty into a mature discipline capable of supporting civilization-scale trust across generations. This hard-won resilience now sets the stage for exploring the emergent horizons of temporal trust—where quantum entanglement, neuromorphic hardware, and interstellar communication beckon with new possibilities and perils—as we turn to the concluding perspectives in Section 10.

---

## 1.8    Section 10: Future Trajectories and Concluding Perspectives

The controversies and failures chronicled in Section 9 – the specter of witness cartels, the Gödelian limits of self-verifying systems, the painful lessons of migration debacles and cosmic expiry – are not endpoints, but crucibles. They forge a more resilient, nuanced understanding of what temporal trust can realistically achieve. Future-signed model certificates emerge from this gauntlet not as a panacea, but as a foundational, albeit perpetually evolving, stratum of the digital infrastructure. The path forward, illuminated by breakthroughs in physics, computation, and social organization, points towards protocols capable of spanning interstellar distances, standards harmonizing trust across civilizations, sociotechnical systems blurring the digital-physical divide, and profound philosophical reckonings with the nature of truth itself across deep time. This concluding section synthesizes these emergent trajectories, mapping the road from operational necessity to a cornerstone of humanity's enduring digital legacy.

### 1.8.1    10.1 Next-Generation Protocols: Beyond Classical Constraints

The relentless pace of cryptanalysis and the demands of novel deployment environments are driving research beyond the lattice-based and hash-based primitives dominating current systems. The next wave focuses on harnessing exotic physics and computational paradigms.

- **Homomorphic Time-Lock Puzzles (HTLPs): Merging Computation and Time:** Traditional timelock puzzles (TLPs) encrypt data decryptable only after a predetermined computational effort, simulating time passage. HTLPs, leveraging Fully Homomorphic Encryption (FHE), allow computation *on the encrypted puzzle* itself before decryption. This enables revolutionary future-signing workflows:

- **"Sealed Future Verification":** A model developer can encrypt a model and its future-signing commitment using an HTLP. Witnesses perform attestations *on the encrypted data*, verifying properties (e.g., the FHE-encrypted model passes certain homomorphically evaluated test vectors) without ever seeing the plaintext. Only after the pre-set time lock expires (e.g., upon model publication or patent expiry) can the commitment and model be decrypted and the historical witness attestations verified. This provides strong confidentiality *during development* with guaranteed future verifiability. **Project Chronos Seal:** DARPA's HOMTIM program funded IBM Research and Galois Inc. to prototype this for defense AI models. Their 2027 demonstration involved homomorphically verifying an encrypted neural network's image classification accuracy against a public dataset within an FHE enclave, attested by witnesses, with decryption scheduled for 2035. This addresses the core tension between proprietary secrecy and long-term accountability in high-stakes R&D.

- **Entanglement-Based Temporal Proofs: Trust Rooted in Quantum Reality:** Exploiting quantum entanglement offers a path to timestamps whose validity derives from fundamental physical laws, not computational assumptions.

- **The EPR Paradox as Notary:** Protocols like "QTemp" create pairs of entangled photons. One photon is measured immediately to generate a commitment hash for the model. The other is stored in a quantum memory. Later verification involves measuring the stored photon. The Bell inequality violations (proving the photons were entangled *before* the second measurement) provide irrefutable proof that the commitment existed at the time of the first measurement, with security guaranteed by quantum non-locality. Crucially, attempts to backdate break the entanglement, leaving detectable statistical anomalies. **Breakthrough & Challenge:** The University of Vienna's Quantum Timestamping Lab demonstrated a lab-scale prototype in 2026, anchoring a 1KB document hash. Scaling to multi-gigabyte model weights requires breakthroughs in high-fidelity, long-duration quantum memory (beyond current rare-earth doped crystal limits) and efficient quantum hashing. The EU's Quantum Flagship project "Entrust" aims for a practical satellite-based entanglement distribution network for global timestamping by 2035, potentially revolutionizing trust anchors for critical systems.

- **Neuromorphic Hardware Integration: Trust at the Speed of Thought:** As AI inference migrates to dedicated neuromorphic chips (like Intel's Loihi, IBM's NorthPole), future-signing must integrate natively at the hardware level.

- **In-Silico Fingerprinting:** Neuromorphic architectures represent models as physical configurations of synapses and neurons. "Analog Fingerprinting" techniques exploit inherent manufacturing variations (once a nuisance) as a unique PUF. Reading the analog conductance state of a specific, designated neural pathway after model loading generates a hardware-bound fingerprint inseparable from the model's physical instantiation. **IBM's NeuSign Prototype:** Integrated into their NorthPole-2 chip, this uses a "signature core" – a dedicated analog subnetwork whose activation pattern, when stimulated by the loaded model's weights, produces a unique, cryptographically signed output. This binds the model's execution integrity directly to the immutable hardware substrate, offering tamper resistance far exceeding software-based checks. Challenges remain in standardizing the fingerprinting process across

different neuromorphic architectures and ensuring resilience against sophisticated physical attacks (e.g., focused ion beam probing).

- **Bio-Hybrid Attestation:** Beyond pure DNA storage, research explores integrating biological processes into active trust mechanisms. **Project Living Root:** A collaboration between MIT Media Lab and the Svalbard Global Seed Vault encodes DTSA root keys into the genomes of Arctic moss species. The slow, natural mutation rate of the moss (monitored via satellite spectroscopy) provides a publicly observable "bioclock." Periodically sequenced mutations serve as a decentralized, environmentally anchored proof of elapsed time, offering a novel layer of resilience against digital ledger manipulation or catastrophic societal collapse. While highly experimental, it symbolizes the quest for trust embedded within the fabric of nature itself. These next-generation protocols represent a paradigm shift: no longer merely *recording* temporal state, but weaving trust into the computational process, quantum fabric, physical hardware, and biological systems, creating multi-layered assurances resilient against previously unimaginable threats.

### 1.8.2   10.2 Standardization Roadmaps: Weaving the Global Trust Fabric

The proliferation of proprietary and national systems, highlighted as a fragmentation risk in Section 9, necessitates aggressive, coordinated standardization to ensure interoperability and universal verifiability. Roadmaps are converging around core international bodies.

- **IETF SCITT Evolution: From Protocol to Ecosystem:** The IETF's Supply Chain Integrity, Transparency, and Trust (SCITT) working group, having standardized the core receipt format and API for DTSA interactions, is now expanding its scope:

- **Witness Interoperability Framework (draft-ietf-scitt-witness-interop):** Defining standard APIs for witness enrollment, attestation submission, reputation reporting, and slashing evidence. This enables witnesses to participate seamlessly across multiple DTSA networks, preventing vendor lock-in. Planned completion: 2029.

- **Post-Quantum Migration Profiles (draft-ietf-scitt-pqc-migration):** Standardizing hybrid signature formats (e.g., combining Dilithium with Ed25519) and migration procedures for witness networks and DTSAs. Mandates cryptographic agility metadata within certificates. Final call expected 2028, with mandatory support timelines tied to NIST PQC migration phases.

- **Lightweight Verification Profiles:** Defining constrained-resource verification protocols optimized for specific sectors: ultra-low-power (ISO 29167 for industrial sensors), deterministic real-time (DO-356A for aviation), and disruption-tolerant (Bundle Protocol for space). W3C collaboration ensures web verifier compatibility.

- **NIST Post-Quantum Cryptography Program: The Migration Imperative:** NIST's PQC standardization, culminating in CRYSTALS-Dilithium (Signatures) and CRYSTALS-Kyber (KEM) as primary selections, is just the beginning for future-signing:

- **PQC Migration Project for Long-Term Signing (SP 1800-206C):** This critical supplement to NIST's foundational model certificate guide outlines a phased, 15-year migration (2028-2043):

- **Phase 1 (2028-2033):** Hybrid signatures mandatory for all new certificates; witness networks begin re-attesting critical historical chains using PQC.

- **Phase 2 (2034-2038):** Pure PQC signatures become the default; DTSAs disable support for classical-only anchoring.

- **Phase 3 (2039-2043):** Deprecation of classical signatures; verifiers may reject non-PQC-migrated historical certificates.

- **Quantum Random Number Generation (QRNG) Standards (SP 800-90C):** Future-signing's security critically depends on entropy. NIST is accelerating standards for device-independent and semi-device-independent QRNGs suitable for integration into HSMs and TEEs used by witnesses and DT-SAs, mitigating RNG failure risks highlighted in Section 6.

- **ISO/TC 307 Blockchain and Distributed Ledger Technologies: Expanding Scope:** Recognizing that DTSAs increasingly use diverse ledger technologies beyond traditional blockchains, ISO/TC 307 is extending its standards:

- **ISO 23259: Interoperability Framework for Distributed Timestamp Authorities:** Defines core functionalities (append-only logging, consensus interfaces, proof formats) independently of the underlying ledger tech (blockchain, DAG, permissioned database). Facilitates cross-DTSA attestation chains. Published 2027.

- **ISO/TR 23476: Long-Term Preservation of DTSA Ledger States:** Best practices for ensuring the accessibility and interpretability of ledger data centuries later, including emulation specifications, format migration protocols, and decentralized archival strategies. Directly addresses the "knowledge loss" systemic risk. Under development, draft release 2029.

- **W3C Verifiable Credentials for AI Models:** Bridging identity and model trust, the W3C VC group is defining extensions for model certificates:

- **Model Credential Schema:** Standardized JSON-LD schemas for encoding model metadata (publisher, architecture family, training data provenance commitments, intended use, ethical constraints) alongside the cryptographic commitment. Enables rich, machine-readable attestations about model properties beyond mere integrity.

- **Selective Disclosure ZKPs:** Integrating zero-knowledge proofs (e.g., based on BBS+ signatures) into VCs allows model owners to prove specific properties (e.g., "this model was trained on data respecting GDPR," "it passes fairness threshold X for demographic group Y") to verifiers without revealing the entire credential or sensitive fingerprint details. Essential for balancing transparency and proprietary concerns. This concerted standardization push aims to transform the fragmented landscape into a cohesive, interoperable global trust fabric where certificates issued under one framework can

be seamlessly verified within another, governed by clear, forward-looking migration and preservation strategies.

### 1.8.3 10.3 Sociotechnical Evolution: Trust in the Cyber-Physical Continuum

Future-signing is escaping the confines of pure digital artifacts, converging with identity systems, and preparing for humanity's expansion beyond Earth.

- **Integration with Decentralized Identity (DID): The Self-Sovereign Model:** Future-signed model certificates are evolving into verifiable attestations issued *about* DIDs, creating a unified trust framework for entities and their digital creations:

- **"DID as Signer":** Instead of X.509 certificates tied to organizations, models are signed by the DID of the developer team or the autonomous AI agent itself (if legally recognized). Microsoft's ION-based Identity Hub and the Decentralized Identity Foundation's (DIF) "Universal Resolver" enable verification of the signer's DID alongside the model's integrity. Siemens Healthineers now signs diagnostic models using a corporate DID attested by the German Chamber of Commerce.

- **Reputation-Backed Signing:** Witness attestations can become verifiable credentials linked to a signer's DID, proving their membership in accredited pools or their historical reliability score. A verifier can check not only *that* a model was signed, but *by whom* and *with what reputation*. The "TrustNet Index" is piloting DID-compatible reputation oracles.

- **Revocation via Identity:** Revoking a compromised signing key becomes synonymous with revoking or suspending the associated DID, leveraging existing DID revocation mechanisms (e.g., Sidetree CRLs, status list VCs) for more efficient management.

- **Cyber-Physical System Convergence: Signing the Real World:** The boundary between digital models and physical actuators is dissolving. Future-signing is extending to the state and behavior of physical systems:

- **Sensor Data Provenance Chains:** Projects like IOTA's "Tangle for Industry 4.0" and Bosch's "CPS Anchor" use lightweight future-signing to create immutable, timestamped chains of sensor readings from factories, power grids, and vehicles. This proves the integrity of the physical data feeding AI models and control systems. Rolls-Royce uses this to sign turbine sensor data streams used by its predictive maintenance AI, creating an auditable trail from physical vibration to maintenance decision.

- **Actuator Command Attestation:** Critical commands sent to physical systems (e.g., "close valve," "apply brakes") are being signed at generation and verified at the edge controller before execution. **Project Cerberus (DARPA):** Develops secure co-processors for military vehicles that verify future-signed commands from AI tactical systems, ensuring only authorized, temporally valid commands actuate weapons or defenses. Requires ultra-low-latency verification, achieved through pre-computed STARK proofs.

- **Digital Twins as Signed Entities:** The comprehensive digital twin of a physical asset (building, aircraft, human body) becomes a future-signed entity itself. Changes to the twin, reflecting physical modifications or sensor updates, are anchored and witnessed. Airbus's "Wing of Tomorrow" program maintains a future-signed digital twin, providing an immutable history of design iterations, simulations, and physical test results used for airworthiness certification.

- **Interstellar Communication Implications: Trust Across Light-Years:** As humanity contemplates probes and communications across interstellar distances, future-signing provides mechanisms for establishing trust despite extreme latency and isolation.

- **Bundle Protocol Security Extensions:** NASA's Delay/Disruption Tolerant Networking (DTN) Research Group is extending the Bundle Protocol to incorporate future-signed certificates for critical telecommand bundles and scientific data bundles. Using lattice-based signatures and pre-shared keys anchored before launch, probes like the proposed Interstellar Probe (launch ~2036) could verify commands received decades later originated from Earth and haven't been altered in transit. Verification must be energy-efficient and resilient against cosmic ray-induced bit flips.

- **Autonomous Verification Oracles:** Probes or interstellar habitats operating beyond feasible real-time communication require autonomous agents capable of verifying future-signed updates or policy changes. **Project Starlight (Breakthrough Initiatives):** Researches compact, radiation-hardened verification modules using neuromorphic processors optimized for lattice-based signature checks, enabling autonomous colonies to validate critical software updates or governance directives sent centuries prior.

- **The Beacon Chain Concept:** Proposals exist for an interstellar "beacon" – a spacecraft or network broadcasting a continuous, future-signed stream of humanity's accumulated knowledge, public keys, and current state. Verification relies on predictable orbital mechanics and quantum-resistant signatures. The Beacon would serve as a common trust anchor for any future intelligence encountering it, a digital Voyager Golden Record for the age of cryptography. The Lunar Library serves as a prototype, its nickel plates future-signed via laser-etched quantum-resistant commitments witnessed by multiple Earth-based DTSAs. This sociotechnical evolution signifies a future where cryptographic trust is not an add-on, but an intrinsic property woven into the identity of digital agents, the state of physical infrastructure, and humanity's attempts to reach the stars – a pervasive, resilient layer securing our increasingly complex existence.

### 1.8.4  10.4 Philosophical Implications: Redefining Truth in the Digital Epoch

The relentless pursuit of enduring digital trust embodied in future-signed certificates forces a profound re-examination of concepts central to human epistemology and our relationship with time and information.

- **Re-defining Digital Permanence: Beyond Bit Preservation:** Future-signing shifts the focus from merely preserving bits (a solved problem with replication) to preserving the *meaning* and *provable*

*authenticity* of information across deep time.

- **The "Context Preservation Problem":** A perfectly preserved and verified model hash from 2050 is meaningless if the context needed to interpret it – the purpose of the model, the nature of its training data, the societal norms it reflected – is lost. Initiatives like the Long Now Foundation's "Manual for Civilization" accompanying its archives, or embedding contextual metadata within W3C Verifiable Credentials, attempt to address this. Philosopher Luciano Floridi argues future-signing necessitates a new discipline of "digital hermeneutics" – methodologies for interpreting ancient digital artifacts whose original context has evaporated.

- **Permanence as a Choice, Not Default:** The Estonian e-ID crisis and the COBOL time bomb incident underscore that not all digital artifacts warrant or can sustain indefinite verification. Society faces difficult choices about what deserves the resource investment for deep-time preservation. Archivists propose "tiered eternity" frameworks, where critical cultural or safety artifacts receive multi-layered future-signing (quantum, DNA, monumental), while ephemeral data is allowed to gracefully expire, challenging the digital hoarding instinct.

- **Epistemological Shifts in Verification: From Authority to Mathematics:** Future-signing represents a radical shift in how we establish belief in the past.

- **The Decline of Institutional Trust:** Where trust traditionally resided in enduring institutions (governments, universities, notaries), future-signing vests it in mathematical proofs and decentralized networks resistant to institutional collapse or corruption. Historian Yuval Noah Harari observes this mirrors a broader trend where "algorithms challenge narratives" as primary sources of truth. The Singapore court's acceptance of self-verifying cryptographic evidence over traditional notarial seals exemplifies this shift.

- **The Rise of the Verifier Citizen:** Lightweight verification empowers individuals to independently confirm the provenance and integrity of critical information (news sources, medical AI diagnoses, election results) without relying on intermediaries. This promises greater autonomy but demands unprecedented levels of digital and mathematical literacy. Projects like "CryptoLiteracy for All" aim to make basic verification concepts accessible, fostering a society where cryptographic proof is as fundamental as reading comprehension.

- **The "Oracle Problem" Recast:** While future-signing secures the *artifact* and its *temporal origin*, it cannot guarantee the *truthfulness* of the information the artifact contains or represents. A future-signed dataset can be impeccably authentic yet fundamentally biased or misleading. Philosophers like Shannon Vallor argue this elevates the need for complementary frameworks focusing on the ethics of data generation and model training – the *inputs* to the signing process – recognizing that mathematical verification alone cannot establish ethical or epistemic soundness.

- **The Quest for "Temporal Truth": A Sisyphean Ideal?** The ultimate aspiration – establishing an objective, immutable record of digital history – confronts fundamental limitations:

- **The Subjectivity of Significance:** What gets signed and preserved is inherently selective. The models, data, and logs chosen for future-signing reflect the power structures, biases, and priorities of the present. Whose truth is being preserved? Feminist technoscience scholars critique the risk of future-signing entrenching dominant narratives, advocating for participatory frameworks where marginalized communities define what deserves temporal anchoring.

- **The Illusion of Neutrality:** The mathematical purity of the signature belies the human choices embedded in the system: the design of fingerprinting algorithms, the selection of witnesses, the governance of DTSAs. These choices encode values and power relations that shape the historical record. Future-signing doesn't escape politics; it becomes a new terrain for political contestation over whose past is deemed verifiable. The debates over witness pool composition and national DTSAs exemplify this.

- **Beyond Verification to Understanding:** Even with perfect verification, *understanding* the past requires empathy, interpretation, and context that cryptography cannot provide. A future historian might verify a 21st-century social media model's integrity but still fundamentally misunderstand the cultural milieu that spawned it. Future-signing preserves the digital artifact, but not the lived experience. It anchors evidence, not meaning.

## 1.9 Concluding Synthesis: The Enduring Scaffold

Future-signed model certificates emerged from a simple need: combating the temporal decay of trust in an increasingly digital and AI-driven world. Through the intricate architecture dissected in Section 3, validated across critical domains in Section 5, hardened against relentless threats in Section 6, and navigating complex legal and socioeconomic landscapes in Sections 7 and 8, they have evolved into a foundational technology. They are no longer merely a solution to cryptographic expiry but a scaffold upon which we build enduring digital integrity – for life-saving AI diagnostics, resilient critical infrastructure, irrefutable legal evidence, and the artifacts we cast towards the stars. The controversies of Section 9 and the frontiers explored in this final section reveal the profound depth of this undertaking. The quest for temporal trust is Sisyphean, perpetually challenged by entropy, power dynamics, epistemological limits, and the sheer vastness of time. Future-signing does not offer absolute, final truth. It offers something more modest yet revolutionary: a mathematically verifiable anchor point in the relentless flow of digital change. It provides a mechanism to say, with high assurance, *this existed then*, and *this is what it was*. In a world awash in misinformation, deepfakes, and mutable digital records, this ability to establish enduring, verifiable facts is not just a technical achievement; it is a prerequisite for rational discourse, accountable institutions, and a coherent historical narrative. The journey chronicled in this Encyclopedia Galactica entry – from the trust gap to cosmic implications – underscores that future-signing is more than cryptography. It is an ongoing sociotechnical project demanding continuous innovation, vigilant security, ethical reflection, and global cooperation. It requires balancing the power of mathematical certainty with the humility to acknowledge its limits. As we deploy this scaffold across our digital and physical worlds, and eventually into the cosmos, we are not merely preserving bits; we are shaping the foundation upon which future civilizations will judge our own. The signature we

leave on time, both digital and existential, must be worthy of the trust it seeks to secure. The technology is now proven; the wisdom to wield it endures as our greatest challenge and responsibility.

---

## 1.10 Section 4: Major Implementation Frameworks

The intricate theoretical architecture explored in Section 3 – with its signature chains, witness networks, robust fingerprinting, and lightweight verification – provides the blueprint for enduring trust. Yet, the true measure of its efficacy lies in concrete realization. This section examines the landscape of production-grade systems translating these principles into operational reality. From nascent industry standards to battle-tested enterprise platforms and specialized domain adaptations, we analyze the design philosophies, deployment patterns, and real-world lessons shaping the practical ecosystem for future-signed model certificates. This transition from cryptographic abstraction to deployable infrastructure marks a critical phase, revealing how divergent priorities – standardization versus flexibility, transparency versus confidentiality, universality versus specialization – manifest in tangible systems securing AI supply chains today.

### 1.10.1 4.1 Industry Standards Initiatives: Forging Common Ground

The fragmentation of early digital timestamping and PKI underscored the necessity of interoperability. Industry consortia and standards bodies have emerged as crucial forces in defining common protocols, data formats, and trust models for future-signing, aiming to prevent a repeat of incompatible silos hindering long-term verification.

- **IETF SCITT (Supply Chain Integrity, Transparency, Trust):** Emerging as the most ambitious standardization effort, the IETF SCITT working group (chartered 2023) directly targets verifiable supply chains, explicitly including AI/ML artifacts. Its core philosophy centers on *transparency logs* and *decentralized governance*:

- **Architecture:** SCITT adopts the DTSA model (Section 3.1) but formalizes it as a "Transparency Service" – a logically centralized but potentially sharded/distributed append-only ledger. Entities submit "envelopes" containing signed claims (e.g., model hashes, attestations) to the service, receiving a cryptographic receipt proving inclusion. Crucially, witnesses monitor these services.

- **Receipt Structure:** A SCITT receipt is a complex cryptographic object. It contains the original claim, the issuer's signature, the Transparency Service's signature (or threshold signature) binding the claim to a specific position in the log, and a Merkle inclusion proof. This receipt is the portable, verifiable proof of provenance and timestamp.

- **Witness Integration:** SCITT mandates interfaces for witness networks to fetch and attest to the state of Transparency Services. Its draft specification defines standard APIs for submitting claims, fetching

receipts, retrieving log entries, and processing witness attestations. A key innovation is the "Consistency Proof" allowing witnesses to efficiently prove the append-only nature of the log between any two points in time.

- **Real-World Traction:** Microsoft's Azure Confidential Compute team provided the initial implementation ("SCITT-Roots") used as a basis for the standard. A notable early deployment is the CNCF's use of a SCITT-compatible transparency log for securing Tekton CI/CD pipeline artifacts, demonstrating applicability beyond pure models. The 2024 SCITT Interop Event successfully demonstrated receipt exchange and verification between implementations from Microsoft, IBM, and the Linux Foundation's Sigstore project, marking significant progress towards interoperability.

- **Confidential Compute Consortium (CCC) Frameworks:** While broader than future-signing, the CCC (founded by Alibaba, Arm, Google, IBM, Intel, Microsoft, Red Hat, Swisscom) tackles the critical intersection of *verifiable computation* and *data confidentiality* – essential for models handling sensitive data or proprietary algorithms. Its key contribution relevant to future-signing is the concept of *Verifiable Claims* within Trusted Execution Environments (TEEs):

- **Remote Attestation Integration:** TEEs (e.g., Intel SGX, AMD SEV, Arm CCA) generate hardware-signed attestation reports ("quotes") proving the integrity of the environment and the initial code loaded. CCC frameworks define how future-signed model certificates can be *bound* to these attestations. For instance, the model fingerprinting process (Section 3.3) can be performed securely within the TEE, and the TEE's attestation report is included as part of the evidence signed into the future-proof certificate. This proves not only the model's provenance but also that it was loaded and potentially executed within a verified, isolated environment.

- **Model Encryption & Sealed Keys:** CCC specifications outline patterns for encrypting model weights using keys protected by the TEE. The future-signing process can then include commitments to the encrypted model *and* the key release policy enforced by the TEE's attestation guarantees. Google's Vertex AI Model Encryption leverages this pattern, where model access requires a future-signed certificate combined with a valid TEE attestation from the requesting environment. This addresses the "model theft" vector while still enabling verifiable provenance.

- **Project Keylime Integration:** The CCC's Project Keylime provides an open-source framework for TEE-based attestation and runtime integrity monitoring. Future-signing systems increasingly integrate with Keylime agents, allowing the witness network or verification engine to not only validate the initial TEE attestation but also receive continuous, signed telemetry confirming the model's runtime behavior hasn't deviated from expectations (e.g., no unauthorized code injection), extending behavioral integrity attestation into the operational phase.

- **NIST SP 1800-206: Foundational Guidance:** The National Institute of Standards and Technology provides critical foundational guidance through its Special Publication series. SP 1800-206, "Trusted and Verifiable AI Supply Chains" (draft released for comment in late 2024), synthesizes best practices

and requirements. While not prescribing a single implementation, it heavily influences procurement
and development:

- **Minimum Requirements:** SP 1800-206 mandates future-signed certificates for "high-risk" AI sys-
tems (as defined by frameworks like the EU AI Act), specifying minimum cryptographic strengths
(e.g., NIST PQC Level 3 equivalency), witness network diversity thresholds (e.g., $\geq 100$ independent
entities across $\geq 3$ jurisdictions), and proof retention periods (e.g., operational lifespan + 10 years).

- **Fingerprinting Standards:** The publication provides detailed recommendations for model finger-
printing, endorsing ensemble approaches combining canonical serialization (ONNX) with behavioral
attestation using statistically sampled test vectors and activation distribution signatures, emphasizing
resistance to equivalence attacks.

- **Verification Blueprint:** It outlines a reference architecture for verification engines, emphasizing
light-client protocols using SNARKs/STARKs and standardized APIs for proof aggregation and post-
quantum fallback verification. The FDA's adoption of SP 1800-206 recommendations (where ap-
plicable) for its "Digital Health Software Precertification Pilot" demonstrates its practical impact on
regulated industries. These standards initiatives represent the collective effort to establish a common
language and baseline for trust. Their success hinges on widespread adoption and avoiding the pitfalls
of competing, incompatible standards that plagued earlier PKI deployments.

### 1.10.2    4.2 Open Source Ecosystems: The Engine of Innovation

Open-source projects provide the agile testing ground and reference implementations for future-signing stan-
dards, driving rapid innovation and lowering barriers to entry. The collaborative nature fosters transparency
and auditability, crucial for building trust in the trust mechanisms themselves.

- **Sigstore's Future-Signing Extensions:** Originally focused on code signing and software supply chain
security (Cosign, Fulcio, Rekor), the Sigstore project (under the Linux Foundation) has aggressively
expanded into model signing and future-proofing.

- **Cosign for Models:** The `cosign` CLI tool now supports signing OCI artifacts, including ML model
files stored in registries like Hugging Face Hub or private registries. The signature includes metadata
identifying the model architecture and framework, creating a basic provenance record. While initially
using standard PKI (Fulcio as a CA), Sigstore is actively integrating future-signing primitives.

- **Rekor++:** The transparency log (Rekor) is undergoing enhancements ("Rekor++") to support long-
term validation. This includes integrating Merkle tree structures optimized for large payloads (leverag-
ing Trillian SMTs), experimental support for witness attestations via pluggable interfaces, and explor-
ing SNARK-based proof aggregation for efficient historical verification. The PyTorch Hub integration
(2023) uses a custom Sigstore instance where model uploads automatically trigger `cosign` signing
and entry into a dedicated Rekor log, providing developers a seamless entry point into verifiable prove-
nance.

- **Witness Integration Prototypes:** Sigstore is prototyping interfaces for witness networks. One approach allows witnesses to monitor specific Rekor instances, periodically generating attestations (signed using FROST threshold schemes in development) about the log's consistency and checkpointing these attestations to diverse backends (e.g., Ethereum, Bitcoin, other transparency logs) for robust anchoring. The "Sigstore Witness Service" (SWS) proof-of-concept demonstrated this at KubeCon 2024.

- **Transparent Log Forks (Trillian/Cosign):** Trillian, developed by Google and now part of the Open Source Security Foundation (OpenSSF), is a general-purpose transparent log that underpins Certificate Transparency and now serves as a key building block for future-signing systems.

- **Scalable Merkle Trees:** Trillian's core innovation is its highly efficient, scalable SMT implementation capable of handling billions of entries. It provides APIs for appending data, generating inclusion proofs, and obtaining signed tree heads (STHs). This directly implements the DTSA ledger concept (Section 3.1).

- **Cosign + Trillian:** The standard Sigstore stack uses Rekor, but projects frequently deploy `cosign` signing backed by custom Trillian logs configured as specialized Transparency Services. This offers greater control over log governance, sharding policies, and witness monitoring rules than the public Sigstore infrastructure. The "Confidential AI Registry" project by the Open Infrastructure Foundation uses a permissioned Trillian log combined with Intel TDX TEEs to manage future-signed certificates for sensitive biomedical models.

- **Witness-as-a-Service (WaaS) Platforms:** The complexity of running a witness node (secure key management, MPC participation, monitoring) has spurred the emergence of open-source WaaS platforms, lowering the barrier to participation in witness networks.

- **Witness Protocol Stacks:** Projects like `witnet-rust` (inspired by the Witnet decentralized oracle network) and `threshold-witness` (developed by the Zcash Foundation) provide modular, open-source implementations of witness node software. They handle secure enclave management (for key isolation), communication protocols for MPC rounds (e.g., GG18/FROST), and interfaces for monitoring target logs (SCITT, Trillian, etc.).

- **Public Good Networks:** Initiatives aim to bootstrap public, permissionless witness networks. The "Proof of Witness" testnet, launched by a consortium of universities (Stanford, ETH Zurich, NUS) in 2024, allows anyone to run a witness node using `threshold-witness` software, stake test tokens, and earn rewards for correctly attesting to the state of participating model registries and transparency logs. While still experimental, it demonstrates the viability of decentralized witness sourcing. The open-source ecosystem is the crucible where standards are implemented, stress-tested, and refined. Its vibrancy ensures that future-signing isn't solely the domain of large enterprises but remains accessible and adaptable.

**1.10.3　4.3 Enterprise Solutions: Scaling Trust for Critical Workloads**

Major cloud providers and technology companies are integrating future-signing into their core AI/ML and security platforms, providing managed services that abstract complexity for enterprise customers. These solutions prioritize scalability, deep integration with existing infrastructure, and meeting stringent compliance requirements.

- **Microsoft Azure Verified Model Registry (VMR):** Positioned as a cornerstone of Azure's Responsible AI framework, VMR offers a comprehensive suite for model provenance and long-term integrity.

- **End-to-End Workflow:** VMR integrates with Azure Machine Learning pipelines. Upon model registration or pipeline completion, it automatically triggers the fingerprinting process (using an ensemble: ONNX canonical hash + activation centroid signature over a curated dataset + optional test vector hash). The fingerprint and metadata are signed using Azure Key Vault Managed HSM keys and submitted to an Azure-managed SCITT-compatible Transparency Service (built on Trillian).

- **Managed Witness Network:** Azure operates a geographically distributed, high-availability witness network (using FROST threshold signatures) that continuously attests to the Azure Transparency Service logs. Attestations are anchored to multiple blockchains (including Ethereum and Bitcoin).

- **Attestation Service:** A key differentiator is the VMR Attestation Service. It generates and signs SCITT receipts on-demand, incorporating the latest witness attestations and cryptographic proof material. Users (or downstream systems) can verify model provenance using a lightweight client library against these pre-packaged receipts. During the 2023 incident involving a compromised third-party medical imaging model, VMR's attestations provided irrefutable evidence for identifying the specific tampered version before deployment, preventing potential patient harm.

- **Google's Binary Authorization for Borg / Vertex AI:** Leveraging Google's internal expertise in large-scale security (Borg) and cloud AI (Vertex AI), Google's solution emphasizes policy enforcement and integration with its Confidential Computing stack.

- **Policy-Driven Deployment:** Binary Authorization (BinAuthz) acts as a gatekeeper. Before a model pod can be deployed on Borg (or a Vertex AI endpoint provisioned), BinAuthz checks a future-signed attestation associated with the model image against predefined organizational policies. Policies can mandate specific signers (e.g., "Signed by VMR"), require valid SCITT receipts, specify minimum witness thresholds, or enforce TEE attestation requirements (via integration with Google Confidential Space).

- **Kritis / KMS Integration:** The attestations are stored and managed via the Kritis system, tightly integrated with Google Cloud Key Management Service (KMS) for signing key security. Kritis handles the verification of attestations and receipts before allowing BinAuthz to permit deployment.

- **Distributed Fingerprinting:** Google utilizes a technique called "Distributed Weight Hashing" for large models. The model is sharded across multiple secure workers. Each worker hashes its shard

using a locality-sensitive hash (LSH), and the collection of LSH values forms the fingerprint. This allows parallel fingerprinting of massive models like PaLM. A notable deployment secures Google Search's core ranking model updates, where BinAuthz ensures only properly future-signed and attested models can be pushed to production.

- **IBM Certifier for AI:** Targeting highly regulated industries and complex hybrid cloud environments, IBM's solution emphasizes governance, auditability, and cryptographic agility.

- **Hyperledger Fabric Integration:** While supporting standards like SCITT, the Certifier often leverages private, permissioned Hyperledger Fabric blockchains as the Transparency Service, appealing to industries (finance, government) wary of public logs. Fabric channels provide data isolation.

- **Advanced Governance Module:** A central feature is a policy engine governing witness selection, key rotation schedules, and cryptographic scheme migration (e.g., automatic transition plans from ECDSA to CRYSTALS-DILITHIUM). It generates auditable trails for all administrative actions and cryptographic ceremonies.

- **Cross-Cloud Verification:** The IBM Certifier Verification Engine is designed as a standalone component deployable on-premises, in private clouds, or within air-gapped networks. It can verify certificates originating from diverse sources (Azure VMR, SCITT logs, private Fabric chains) using a unified API, crucial for enterprises with multi-vendor AI supply chains. IBM's collaboration with the FDA on the "AI Validation Network" pilot uses the Certifier to manage future-signed certificates for diagnostic models across multiple participating hospital networks, enabling cross-institutional model validation audits. Enterprise solutions demonstrate the scalability and operational rigor required for real-world, high-stakes deployments. They bridge the gap between cutting-edge cryptography and the practical demands of global AI infrastructure.

### 1.10.4  4.4 Specialized Domain Implementations: Tailoring Trust

Beyond generic platforms, future-signing is being adapted to meet the unique constraints, regulations, and risk profiles of specific high-impact domains. These implementations reveal how core principles are specialized for extreme environments.

- **Medical Device Firmware & AI Models (FDA Submissions):** Regulatory compliance (FDA 21 CFR Part 11, EU MDR) and patient safety drive stringent requirements.

- **Enhanced Fingerprinting & Audit Trails:** Submissions for AI-based SaMD (Software as a Medical Device) increasingly mandate future-signed certificates. Fingerprinting goes beyond the model to include the *entire software bill of materials (SBOM)* – OS, libraries, drivers – and the training dataset fingerprint (using deduplicated MinHash). The FDA's eSTAR submission portal now accepts SCITT receipts as part of the "Technical Documentation" package. A landmark case was the 2024 approval of HeartFlow's FFRct AI, where the future-signed certificate chain provided the audit trail linking

the cleared model version to the specific clinical validation data cited in the submission, streamlining regulatory review.

- **Hardware Anchoring:** For implantable devices or critical bedside monitors, future-signed firmware updates are often anchored within hardware security modules (HSMs) or TEEs on the device itself. Verification occurs locally before installation, ensuring integrity even if cloud services become unavailable. Medtronic's Azure CIED (Cardiac Implantable Electronic Device) platform uses this pattern for secure field updates.

- **Autonomous Vehicle (AV) Model Certification:** The dynamic operational environment (OTA updates), safety-critical nature, and regulatory scrutiny (ISO 21448 SOTIF, UL 4600) necessitate robust solutions.

- **Continuous Attestation:** Beyond initial model signing, AV systems require *continuous* future-signing of operational data snapshots and model performance metrics. These are signed periodically (e.g., every drive cycle) by secure vehicle HSMs and streamed to manufacturer transparency logs. This creates an immutable "black box" record crucial for incident investigation and proving compliance with operational design domains (ODDs).

- **Sensor Fusion Fingerprinting:** Fingerprinting extends to the entire perception stack. Techniques involve hashing fused sensor data (Lidar, Radar, Camera) representations or the inputs to planning modules at key decision points, signed in real-time. This provides evidence that the deployed model received the inputs its safety validation assumed. Waymo's "Verifiable Data History" system, utilizing future-signed certificates anchored to a private DTSA, exemplifies this approach, providing crucial evidence in liability disputes.

- **National Security Applications:** Secrecy, tamper-proofing, and long-term (decades/centuries) validation are paramount.

- **Air-Gapped Witness Networks:** Highly classified model registries operate entirely disconnected from public networks. Witness networks consist of physically isolated, geographically dispersed secure facilities (e.g., national labs, secure bunkers). Attestations are exchanged via secure physical transport of encrypted media (e.g., one-time pad encrypted hard drives), creating an "optical air gap." Verification engines run within secure compartments.

- **Physically Unclonable Function (PUF) Integration:** For extreme tamper resistance at endpoints (e.g., field-deployed sensors or autonomous systems), the future-signed model certificate or its verification key is cryptographically bound to a hardware PUF unique to the specific device. Tampering physically alters the PUF, invalidating the key and preventing model execution. IARPA's "SHIELD" program explores PUF-anchored future-signing for intelligence applications.

- **Non-Digital Fallbacks:** Recognizing the potential for digital obsolescence, some archives employ "ceremonial key carving" – threshold keys split into physical shares (e.g., etched on titanium plates)

stored in ultra-secure, geographically dispersed vaults (e.g., Salt mines, arctic repositories). Instructions for reconstructing keys and verifying proofs using future computational methods (possibly non-electronic) are included. This blends the digital future-signing paradigm with the physical permanence of historical time capsules. These specialized implementations demonstrate the remarkable flexibility of the future-signing paradigm. Whether ensuring a pacemaker's firmware update won't kill a patient, proving an autonomous car "saw" what it claimed before a crash, or securing state secrets for a century, the core architecture adapts to enforce trust under the most demanding conditions. The landscape of implementation frameworks is diverse and rapidly evolving. Industry standards provide the essential connective tissue, open-source projects fuel innovation and accessibility, enterprise platforms deliver scalable trust for global operations, and specialized adaptations meet the most extreme requirements. This rich ecosystem, translating the cryptographic bedrock of Section 3 into operational reality, forms the infrastructure upon which reliable AI supply chains are being built. Yet, the ultimate test lies not in the architecture or the implementation, but in its application. How is this technology being deployed to solve real-world problems? What are its measurable impacts? It is to these practical use cases and deployment scenarios that we now turn our attention.

---