

# Governance-Ready Data Sharing Infrastructure

Policy-Safe Marketplaces Powered by Privacy-Tech and Verifiable Contracts

**Author:** *Mehtab A. Rosul*

Director of R&D, EncryptArx — Senior Technical Researcher & AI-ML Engineer

**Date:** August 28, 2025

## Abstract

Data markets promise huge value by enabling firms, researchers, and platforms to purchase, sell, and license data. To date, however, privacy, regulatory constraints, and commercial incentives have limited broad adoption. Privacy-preserving data marketplaces-coupled with synthetic-contract frameworks that integrate legal contracts, cryptographic proofs, and programmable enforcement-offer a producible path forward. This article synthesizes technical building blocks, proposes layered product architecture for marketplaces, describes a practical synthetic-contract model for transactions and provenance, and outlines regulatory, governance, and business models that make such systems commercially viable. We conclude with an implementation roadmap, risk register, and concrete recommendations for product leaders, policy makers, and buyers/sellers of data.

## 1. Why these matters - problem statement and opportunity

High-quality data is the fuel of modern AI. Yet most valuable data is sensitive-personal health records, financial transactions, enterprise telemetry-and tightly regulated. Two persistent problems block efficient data exchange:

Privacy & legal risk: Transfer of raw data often violates privacy laws or contractual obligations. Organizations rightly refuse to share because of compliance and reputational risk.

Misaligned incentives: Data sellers seek fair compensation and control over downstream uses. Buyers want rich, verifiable, well-documented datasets with predictable utility.

A privacy-preserving data marketplace is a place that enables commerce in data while making sure that no party, including the operator of the marketplace, gets inappropriate access to identifiable raw data. The synthetic-contract framework ties together the economics and legal enforceability of transactions to cryptographic and auditable technical primitives: data provenance, usage policies, verifiable computation, and on-chain or off-chain settlement.

If done correctly, these systems are a win-win: sellers monetize assets with reduced risk of leakage; buyers get high-utility datasets or insights under auditable controls; and regulators gain the transparency and verifiability they need.

## **2. Core technical primitives (brief, essential primer)**

A production system depends on composable privacy technologies. Each of the primitives solves a different problem, and together they create useful tradeoffs.

**Differential Privacy.** A mathematically rigorous framework that provides quantifiable bounds on disclosure risk when releasing aggregate statistics or trained models. DP represents the most practical tool today for publishing useful outputs with provable privacy guarantees. Useful for analytics, model training, and sanitization of synthetic data.

**Synthetic Data Generation:** Machine-learning methods that generate realistic but non-identifiable data replicas, including GANs, VAEs, diffusion models, and probabilistic simulators. Synthetic data reduces leakage risk while preserving downstream model utility when validated properly.

**Secure MPC:** A protocol that allows mutually distrustful parties to compute together a function over private inputs, without revealing these inputs, while being suitable for collaborative analytics or training where the raw data cannot leave its owner.

**Homomorphic encryption:** It enables direct computation on encrypted data and is suitable for outsourced analytics for which the operator must not see plaintext. Practical for constrained classes of operations and with continued performance progress.

**TEEs:** hardware enclaves—such as SGX-style—running code in isolated and attested environments so sensitive data can be processed safely when parties trust enclave attestation.

**Verifiable Computation & Zero-Knowledge Proofs:** As the names suggest, these are cryptographic techniques to prove that a computation has been executed correctly on given inputs without revealing those inputs. Useful for auditability of marketplace computations.

**Smart Contracts & Programmable Enforcement:** Blockchain or ledger primitives that enable escrow, automated payments, and immutable recording of agreements and provenance for dispute resolution and incentives.

Each of these primitives has trade-offs-performance, trust assumptions, cost, among others. A practical marketplace stitches the primitives together to support many transaction patterns, rather than forcing one-size-fits-all approaches.

### **3. Marketplace product architecture — layered, pragmatic design**

A market aimed at enterprise and regulated domains must be modular. Following is an operational architecture that provides a balance between privacy, latency, and usability.

#### **Layer 1 — Registry & Provenance Layer**

An artifact registry storing metadata: dataset descriptors, schema, lineage, owner assertions, policy templates, and cryptographic hashes for artifact verification.

Provenance ledger (append-only): records dataset creation, transformations, adapter deployments, and contract bindings. This ledger may be implemented by a permission blockchain or a signed audit log.

#### **Layer 2 — Policy & Contract Layer (Synthetic-Contract)**

Policy templates: machine-readable usage policies (who, purpose, duration, retention, allowed analyses). Standardized languages are to be used for automation, such as ODRL-style or JSON-schema.

Synthetic-contract engine: The policy gets bound to a marketplace transaction and then enforces it, combining cryptographic guarantees with legalese. Example contracts include payout terms, escrow conditions, audit clauses, and revocation rules.

Enforcement primitives: smart contracts for escrow + off-chain enforcement agents that validate compliance by means of attested computation or verifiable proofs.

### **Layer 3 — Privacy Execution Layer**

Support multiple privacy execution modes. A seller and a buyer choose the appropriate pattern for each transaction.

**Mode A - Synthetic Data Delivery.** Vendor supplies DP-sanitized synthetic dataset (with measured utility/DP  $\epsilon$ ). Purchaser pays and gets dataset plus evaluation harness held in registry showing expected model performance on held utility tests.

**Mode B - Query-as-Service (QaaS).** Buyers execute preapproved analytic queries on seller's data via an MPC/TEE pipeline; only aggregate, DP-sanitized results are returned to buyer.

**Mode C - Model-Training as Service.** Buyer submits training logic or model spec; model trains locally or via MPC on sellers' encrypted data; resulting model is either returned - with DP guarantees - or hosted behind an API with access controls.

**Mode D- Verifiable Insights.** Vendor executes analytics, generates a ZKP that attests to the correctness of the outputs and policy compliance; the client receives the attested results and pays automatically in case of verification success.

### **Layer 4 — Marketplace UI & Risk Scoring**

Instrumented UI for sellers to profile datasets: data quality, sensitivity score, and compliance flags; buyers to evaluate expected utility and compliance fit; and legal teams to preview contracts.

Risk-scoring engine combining automated sensitivity classification, regulatory-context rules, and historical incident data to recommend privacy modes and pricing.

### **Layer 5 — Settlement & Escrow**

Payment channels - fiat on-ramp, stablecoins, or hybrid - with escrow. Synthetic contract conditions released upon verifiable proof of correct delivery or performance, such as the model meeting advertised metrics on an independent test harness.

This layered approach allows the onboarding of a wide range of sellers and buyers incrementally because it starts with Mode A/B and gradually incorporates heavier cryptographic modes as demand and performance permit.

## **4. Synthetic-Contract frameworks: design and enforcement**

A synthetic contract is both a legal agreement and a technical contract: it encodes obligations and enforcement paths that use verifiable computation and programmable settlement. Key components:

Canonical contract schema. Human-readable legal clauses (licenses, liability, termination) mapped to machine-readable policy primitives: allowed queries, retention windows, access scopes.

Proof obligations. Clauses that require provable artifacts (e.g., "*Seller must provide a DP proof that  $\epsilon \leq 1.0$  and supply seed/hash of generator model*"; "*Buyer may request audit evidence for at most two transactions per quarter*").

Escrow & penalty rules: Financial holdbacks for large deals released upon passing independent verification or arbitrable conditions.

Audit & dispute protocols: Upon dispute, the parties agree to provide an independent arbiter with sealed artifacts under NDA. Attested logs and signed manifests give the arbiter evidence necessary to rule disputes.

Revocation & updates: dynamic policies, which may be revoked or updated, with the consent duly recorded, and lead to automatic reversion of access or obligations.

### **Enforcement modes map to technical primitives:**

The Marketplace enforces "no raw data export" by using MPC/TEE/QaaS modes.

To provide "DP guarantees," the seller should provide proof of DP application, auditable pipeline, and optionally a ZKP of correct noise addition

To prove "utility claims," the seller can bind an evaluation harness held in registry that runs buyer's black-box model on a small held test set to verify expected performance, releasing only summary metrics.

This hybrid approach couples legal deterrence with technical verifiability, substantially lowering the possibility of successful misuse.

## **5. Business models & commercial hooks**

Privacy-preserving marketplaces are not a charity; they must align incentives. Possible viable monetization and go-to-market hooks include:

QaaS: per-query pricing. The buyer only pays per approved query execution. The seller gets recurring revenue and retains the raw data. This model suits analytics and regulated sectors.

Subscription & API access; hosting the trained models behind APIs-model leasing, where buyers will pay for accessing model predictions instead of the data. Models can be periodically retrained under fresh DP constraints, with revenue shared

Outcome-based pricing: Payments contingent upon model performance on independently verified benchmarks. Useful in verticals such as drug discovery. Synthetic-contract escrow releases payments based on attainment of pre-agreed metrics.

Data cooperatives & revenue sharing: Communities (hospitals, banks) pool data under governance rules and share revenue from aggregated, privacy-preserving offerings. Cooperative governance increases supply and social license.

Compliance as a Service: Provide packaged compliance artifacts like signed manifests, audit reports, and DP proofs for the enterprise buyer, who would use them as evidence bundles for regulators.

Marketplace tooling & consulting: sell integration, redaction and model-validation services to bridge the legacy systems into privacy modes.

A realistic commercial rollout will blend a few of these models, starting with the low-friction per-query and subscription services, then adding outcome-based deals as trust and verification tooling mature.

## **6. Policy and regulatory considerations**

A legitimate marketplace operates in harmony with data protection and consumer regulations.

Data protection law GDPR & equivalents. Even synthetic and aggregated outputs constitute personal data depending on the level of re-identification risk. Marketplaces need to document a Data Protection Impact Assessment, implement Data Protection where possible, and provide robust deletion and access controls.

Contractual liability and insurance: Synthetic contracts should clearly indicate the liability of each party, such as the seller for misrepresentation and the buyer for misuse. Such contracts may also need indemnity and cyber insurance clauses.

Competition & antitrust: mitigate risks of market power concentration-design marketplace to avoid vendor lock-in, interoperable standards and enable data portability.

Regulatory acceptability of synthetic data: Policymakers increasingly ask for standards quantifying the privacy/utility tradeoff of synthetic datasets. Marketplaces should produce standard KPIs: reidentification risk metrics, utility scores.

Cross-border data flows - leverage geofenced execution modes and local TEEs to satisfy sovereignty regulations.

Regulators are likely to look upon systems favorably that: provide auditable attestations of privacy guarantees, furnish machine-readable evidence packs to enable supervision, and support rapid incident response.

## **7. Risk analysis & mitigation checklist**

### **Key risks and mitigations:**

Re-identification attacks on synthetic data. Mitigation: conservative DP parameters, membership inference testing, external red teaming, and certification of synthetic pipelines.

Malicious buyers attempting to reconstruct data through adaptive queries. Mitigation: query budgets, DP noise on query results, pre-approved query templates, and anomaly detection on query patterns.

MPC/HE cryptographic modes: performance/latency. Mitigation: Hybrid placement-use MPC for high-value low-frequency analytics and use TEEs when medium latency is sufficient; optimize using batching and approximate/quantized computations.

Trust vulnerability: side-channel attacks from TEEs. Mitigation: multi-mode redundancy, MPC fallback, enclave updates, and layered logging/attestation.

Legal ambiguity around synthetic data status: Mitigation may involve comprehensive provenance records, legal opinions, and adherence to conservative privacy budgets.

A formal risk register mapped to SLA obligations and insurance premiums should form a part of the marketplace operating manual.

## **8. Implementation roadmap (12–24 months, phased)**

### **Phase 0 — Discovery & pilot (0–3 months)**

Identify cohorts of sellers-healthcare networks and financial institutions-and pilot use cases with buyers.

Build artifact registry and simple provenance ledger.

Pilot Mode A - Synthetic data delivery with clear DP reporting.

**Phase 1 — QaaS & contracts (3–9 months)**

Launch QaaS with TEE-backed execution and per-query DP sinks.

Implement synthetic-contract engine with escrow and basic dispute flow.

Add marketplace UI and risk scoring.

### **Phase 2 — Advanced cryptography & verification (9–18 months)**

Integrate MPC for collaborative analytics and HE prototypes for constrained workloads.

Add ZKP attestation for selected computations with automatic release of escrowed funds on verified proofs.

### **Phase 3 — Scale & governance (18–36 months)**

Expand supply via data cooperatives; add outcome-based contracting; and formalize external audit partnerships.

Publish compliance artifacts, formalize DPIA processes, and engage regulators for pilot approvals.

This roadmap balances early revenue in synthetic data and QaaS with progressively harder engineering bets in MPC/HE and ZK proofs as customer trust grows.

## **9. Recommendations - practical guidance for product leaders**

Start with Synthetic Delivery + DP: It's the quickest way to market, and it informs pricing and utility metrics.

Provide several privacy modes. Allow sellers to choose risk profiles, i.e., higher prices for tighter guarantees.

Standardize machine-readable policies. Interoperability speeds up enterprise adoption.

Invest in independent audits & red teams. Third-party certification engenders buyer confidence and regulatory goodwill.

Establish clear dispute and remediation pathways. Business continuity and legal clarity reduce friction for large deals.

Design for portability: use open formats for manifests and proofs, avoiding lock-in and facilitating regulator review.

## 10. Conclusion

Privacy-preserving data marketplaces coupled with synthetic contract frameworks are a pragmatic, economically attractive, and policy-aligned approach that can help unlock sensitive data for AI. They are not silver bullets: careful technical choices should be made, conservative privacy budgets adopted, robust verification performed, and clear legal instruments established. However, by putting together well-understood primitives—differential privacy, TEEs, MPC-pragmatic delivery modes—synthetic datasets, QaaS-and enforceable synthetic contracts, product teams can create markets that pay data owners, verifiably provide utility to the buyers, and satisfy regulators.

The successful marketplace will be incremental: immediate value delivered with synthetic datasets, verifiable compliance demonstrated, trust reinvested into higher-assurance cryptographic modes, and expansion into outcome-based contracts. For the leaders building such platforms, the time to act is now high-quality, compliant data access is in growing demand, and well-engineered marketplaces will capture long-term value while preserving fundamental privacy and social license.

## References to look for

- C. Dwork & A. Roth (2014). The Algorithmic Foundations of Differential Privacy. Foundational text about differential privacy. I. Goodfellow et al. (2014).
- Generative Adversarial Nets. Foundational paper on GANs and synthetic data generation. G. Ateniese et al. (2000).
- Provable Data Possession at Untrusted Stores. (On cryptographic proofs and data integrity — relevant to verifiable computations and provenance.)
- Craig Gentry (2009). A Fully Homomorphic Encryption Scheme. (Seminal work on HE — practical advances build on this.) Yao, A. C. (1982).
- Protocols for secure computations (Yao's Millionaires' Problem).
- (Foundational MPC work.) Vitalik Buterin (2016). On-chain vs Off-chain Contracts & Smart Contract Design.
- (Practical perspective on programmable enforcement and escrow.) European Data Protection Board / GDPR guidance - selected technical opinions. Background on data protection obligations relevant to marketplaces.