

WHITE PAPER

# HOW DESIGN-DRIVEN SYNTHETIC DATA ENABLES THE CONVERGENCE OF TEST DATA MANAGEMENT WITH AI

### **Executive Summary**

Two data-centric disciplines—Test Data Management (TDM) and AI/ML training data generation—have historically operated in parallel, each serving distinct purposes with separate tools, strategies, and success metrics. TDM has been focused on accelerating test coverage, ensuring compliance, and improving software quality across DevOps pipelines. AI/ML data provisioning, on the other hand, has centered around statistical distributions, large-scale data generation, and eliminating bias in training sets to support intelligent systems.

While the business goals and implementation details differ, both arenas share a deep reliance on one critical resource: **high-quality, secure, context-aware data**. As AI becomes embedded into enterprise software applications, QA processes, and CI/CD workflows, these once-separate domains are converging—both in practice and in purpose.

This white paper explores the dynamics of that convergence and argues that **GenRocket's Design-Driven Data platform** is the only synthetic data solution uniquely qualified to address both markets simultaneously. At the center of this convergence is the shared requirement for data that is not only **realistic and compliant** but also **intentional and fit for purpose**. In software testing, quality data means fewer defects and more reliable deployments. In AI/ML, it means higher prediction accuracy, stronger generalization to unseen data, and reduced bias or unintended behavior in models.

### Market Overview: Two Segments on a Collision Course

#### Test Data Management (TDM)

According to the latest market data published by <u>Verified Market Research</u>, the TDM market is valued at **\$1.54 billion in 2024** and projected to reach **\$2.5 billion by 2030**, growing at a **CAGR of 11.2%**. This steady growth is fueled by:

- The rise of agile and DevOps methodologies
- Increased software release velocity
- Regulatory pressure around data privacy (HIPAA, GDPR, CCPA)
- The complexity of modern enterprise systems

In this space, synthetic test data plays a crucial role. Unlike masked production data, synthetic test data provides safer, faster, and more flexible means to test new features, automate regression testing, and validate integrations.

#### AI/ML Training Data

In contrast, the AI training dataset market is experiencing exponential growth. It its latest research, <u>Grandview Research</u> estimated the market at **\$2.6 billion in 2024** and it's expected to reach **\$8.6 billion by 2030**, with a CAGR of **21.9%**.

#### This surge is driven by:

- 1. Enterprise-wide AI adoption, fueling demand for domain-specific training data at scale
- 2. Data privacy regulations (e.g., GDPR, HIPAA) that restrict use of real-world data, accelerating the shift to synthetic datasets
- **3.** The rise of large-scale models and LLMs, which require massive, diverse, and highquality inputs
- 4. A growing focus on fairness, auditability, and bias mitigation, increasing the need for controlled, rule-based datasets

As AI becomes embedded in software functionality, the data used to train models is just as critical as the algorithms themselves.

## **Comparison: Distinct Needs, Common Foundation**

Attribute	Test Data for QA	Training Data for AI/ML
Primary Goal	Improve test coverage, reduce defects	Improve model accuracy and decision- making
Volume & Scale	Controlled, deterministic	High-volume, statistically accurate
Structure	Schema-driven, often relational	Structured or unstructured data
Generation Logic	Business rules, test cases	Statistical profiles and rule-based
Deployment	CI/CD pipelines, functional testing	Training pipelines, retraining, validation
Quality Definition	Pass/fail logic, defect exposure	Precision, robustness, fairness, predictability
Compliance Needs	Data masking, referential integrity	Bias mitigation, privacy, hallucination avoidance

Despite these differences, both disciplines require synthetic data that is **accurate, secure, realistic, and tailored to the task at hand**.

### Why Convergence Is Already Happening

Convergence between TDM and AI/ML data provisioning isn't theoretical—it's visible across multiple dimensions:

- Al is now embedded in software: Predictive recommendations, fraud alerts, and anomaly detection models are part of the product.
- **Testing is becoming intelligent**: Model-assisted test generation is growing. Al is being used to predict test gaps and identify defects.
- **Shared governance**: Enterprise privacy, compliance, and data governance policies now span both QA and data science.
- Merged DevOps & MLOps: CI/CD pipelines increasingly include training steps, model validation, and feedback loops.

This convergence demands a unified approach to synthetic data: one that supports both precision testing and intelligent training at scale.

## Introducing Design-Driven Data™

**GenRocket's Design-Driven Data**<sup>™</sup> is a synthetic data paradigm that allows organizations to create data based on **intention and control**—not just random data subsets or synthetic replication of production data values.

#### Key Features:

- **Rule-Based Data Generation**: QA teams can define test cases, edge conditions, and schema rules to generate structured test data aligned with business logic.
- **Statistical Conditioning**: Data scientists can define feature distributions, class balance, or rare-event amplification to simulate real-world behavior.
- **Massive Scale**: Supports both record-level precision for testing and bulk volume for generating millions of training data records in minutes.
- Schema & Metadata Awareness: Data generation is schema-driven and supports referential integrity—crucial for relational systems and tabular models.
- **Policy-Driven Compliance**: Enforces data protection rules such as synthetic-only environments, GDPR safeguards, and HIPAA workflows.

By combining these capabilities, GenRocket empowers QA, DevOps, and ML teams to collaborate on a **single synthetic data platform** that speaks to both accuracy and agility.

## How GenRocket Supports Both Use Cases

Feature	Software Testing Benefit	AI/ML Training Benefit
Design-Driven Workflows	Accurate test coverage	Precise model scenario construction
Data Conditioning	Edge-case simulation, negative testing	Bias reduction, class balancing
High-scale Generation	Parallelized regression and performance tests	Bulk data creation for model training
Schema & Referential Integrity	Valid relational test datasets	Structured, consistent inputs for tabular models
Policy Enforcement	Regulatory compliance	Ethical AI training and auditability
API/DevOps Integration	CI/CD automation	MLOps orchestration

## **Real-World Use Cases**

### Healthcare - Claims Testing & NLP

- QA engineers generate HL7 and X12 messages for integration testing.
- ML teams build models to predict claim denials or auto-classify diagnoses.
- GenRocket provisions structured, compliant data for both pipelines.

### Fraud Detection - Banking

- QA tests simulate suspicious transaction patterns (e.g., location velocity).
- ML teams use amplified signals to train anomaly detection models.
- Synthetic data fills rare-pattern gaps in production logs.

### Loan Origination - Financial Services

- Synthetic personas test edge cases like income thresholds or risk flags.
- Al models use the same features to assess creditworthiness or defaults.
- Shared data definitions improve consistency across engineering teams.

### **Retail - Recommendation Engines**

- Product inventory and user session events are generated synthetically for QA.
- Training datasets simulate diverse shopper behaviors for personalization.
- Shared logic prevents test/model divergence.

### Data Quality: The Unifying Thread

Data quality may look different depending on the audience:

- In QA: It means every scenario is covered, every test passes, and defects are minimized.
- In ML: It means accurate, non-biased predictions and well-generalized models.

But the **underlying requirement is the same**: data must be purpose-built and trustworthy. GenRocket enables teams to design the data they need to meet those outcomes.

## **Converging Synthetic Data Markets**



#### Security, Privacy & Compliance Built-In

Both QA and ML workflows are heavily scrutinized for security risks:

- QA teams must never use production data in lower environments.
- ML models must not learn from sensitive or biased historical data.

GenRocket's platform ensures:

- In-place masking when required
- Policy enforcement across all generation modules
- Synthetic-only workflows with zero risk of data breaches

This means that every stakeholder—developer, QA, data scientist, compliance officer can trust the data.

#### Toward a Unified Synthetic Data Market

What was once two separate solutions is becoming one strategic imperative: a **unified synthetic data platform** that powers both test automation and intelligent systems.

With AI being woven into every stage of the software lifecycle, and software teams becoming increasingly data-driven, the enterprise needs a platform that:

- Provisions data for testing, training, validation, and monitoring
- Integrates into both **DevOps** and **MLOps**
- Supports **compliance and security** without compromising agility

## Strategic Implications for the Enterprise

As digital transformation accelerates, synthetic data is no longer a QA tool—it's a strategic asset.

- **DevOps and MLOps are converging**: Organizations must deliver quality code and intelligent functionality together.
- **Release frequency is rising**: Synthetic data speeds up validation, enabling faster deployment.
- Security posture must be unified: Test and training data must meet the same compliance standards.
- **Data provisioning is a bottleneck**: Design-Driven Data removes friction by giving teams what they need—on demand.

Enterprises adopting a unified synthetic data platform can reduce data preparation time, improve model ROI, increase test coverage, and maintain a strong compliance posture—all from a **single data provisioning platform**.

### Conclusion: Why GenRocket Leads the Convergence

With the synthetic data market projected to exceed \$11 billion by 2030, the winners will not be tools that serve one niche, but platforms that unify functions.

GenRocket's Design-Driven Data™ is that platform.

- It meets the precision of TDM and the scale of AI training data.
- It integrates with CI/CD and MLOps pipelines.
- It is governed, secure, and designed to meet enterprise policy standards.
- It empowers teams to generate **exactly the data they need** for test coverage, anomaly detection, bias mitigation, or AI observability.

Where other vendors offer either speed or control, GenRocket delivers both. Where others serve QA or ML, GenRocket serves **both—and brings them together.** 

It's not just synthetic data. It's synthetic data by design.