

# The mcube™ Semantic Lakehouse with Integrated Ontology Layer

**Strategic White Paper - 2025** 

## **Executive Summary**

Executives today recognize that storing large volumes of data is no longer enough; the decisive competitive advantage lies in the ability to transform disconnected facts into coherent. machine-interpretable knowledge. Data Lakehouse architectures — repositories that merge the scale of data lakes with the performance of warehouses represent an important step. Yet even the most modern Lakehouse struggles when data arrives without a shared meaning. Inconsistent labels, duplicated identifiers, and undocumented lineage translate into slow analytics, questionable metrics, and artificial intelligence systems that hallucinate. The mcube™ Semantic Lakehouse from TCG Digital addresses that gap by weaving an ontology — that is, a controlled vocabulary that unambiguously defines every business concept — directly into the storage and processing layers. As data lands, it is mapped to this ontology and linked into an enterprise knowledge graph, turning raw events into a network of well-defined relationships. Machine Learning (ML) and Generative Artificial Intelligence (Generative AI) components in mcube™ then operate on semantically rich inputs, which raises accuracy, transparency, and regulatory compliance. Early adopters see upfront value in three areas:

faster project delivery because less time is lost cleansing data; trustworthy analytics because metrics trace back to source records; and reliable Al because Large Language Models (LLMs) receive grounding context that prevents hallucinations. Independent surveys confirm the urgency of these issues. For instance, Gartner predicts that by 2025, 80% of data and analytics innovations will rely on graph technologies—up from 10% in 2021.

This document is structured in six thematic arcs: business challenges, architectural overview, in-depth capability walkthrough, cross-industry use cases, implementation roadmap, and quantifiable return on investment. It concludes with a compact glossary and an annotated bibliography for further reading.

# CONTENT

Executive Summary	1
1 From Siloed Data to Actionable Knowledge	3
2 The mcube™ Philosophy: Meaning Driven by Design	4
3 Core Capabilities Explained in Depth	5
3.1 Semantic Ingestion and Harmonization	5
3.2 Dynamic Knowledge Graph and Reasoning	5
3.3 Graph-Aware Machine Learning	5
3.4 Generative AI with Retrieval Augmented Generation	5
3.5 Low-Code Application Assembly with ezeXtend	5
3.6 Security, Governance, and Audit	5
4 Cross-Industry Panorama	6
4.1 Refineries & Petrochemicals – Semantic twins for safer plants	6
4.2 Life Sciences & Healthcare – R&D acceleration and personalized care	6
4.3 Manufacturing & Quality – Inline defect prediction	6
4.4 Aviation – Predictive maintenance with route context	6
4.5 Government & Smart Cities – Joined-up public-service intelligence	6
4.6 Insurance – Network-aware fraud analytics	6
4.7 Retail & CPG – Demand sensing and resilient supply chains	
5 Implementation Journey	6
6 Return on Investment	7
7 Glossary	7
8 References	8
About TCG Digital	8



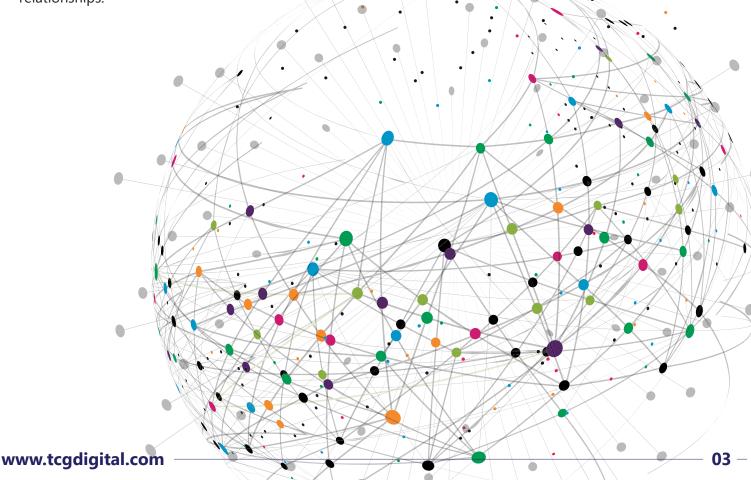
## 1. From Siloed Data to Actionable Knowledge

Organizations have spent decades deploying Data Warehouses, Data Lakes, and more recently, Data Lakehouses. Each step expanded capacity, yet a persistent bottleneck remains: the meaning of data does not travel with the data itself. A "customer" in one system might be a "client" in another. "OrderDate" might store the booking date in System A and the shipping date in System B. Without an explicit semantic agreement, these mismatches force analysts to waste weeks reconciling nomenclature. A 2024 Heartex survey showed that more than 70% of data professionals still devote at least half their working hours to data preparation and management tasks. The direct cost is lost productivity; the indirect cost is strategic: projects slow, opportunities fade, and artificial intelligence models underperform.

Semantic technology solves that by elevating the discussion from where a field lives to what that field means. An ontology captures, for example, that a Customer is a type of Party, that it can place an Order, and that every Order contains at least one Lineltem. Once that ontology is accepted as the source of truth, each dataset can map its columns to these definitions. A knowledge graph then stores the actual instances—Customer #C123, Order #O456—and interlinks them via ontological relationships.

Why does this matter? Because analytics and Al thrive on context. A Machine Learning model predicting churn performs better when it knows the contract hierarchy of a customer, not just a flat feature vector. A Generative Al assistant will hallucinate fewer answers when it can fetch supporting facts from a knowledge graph and cite them. Regulatory compliance likewise improves: auditors can trace a reported metric back through graph links to original source systems, complete with timestamps and user actions.

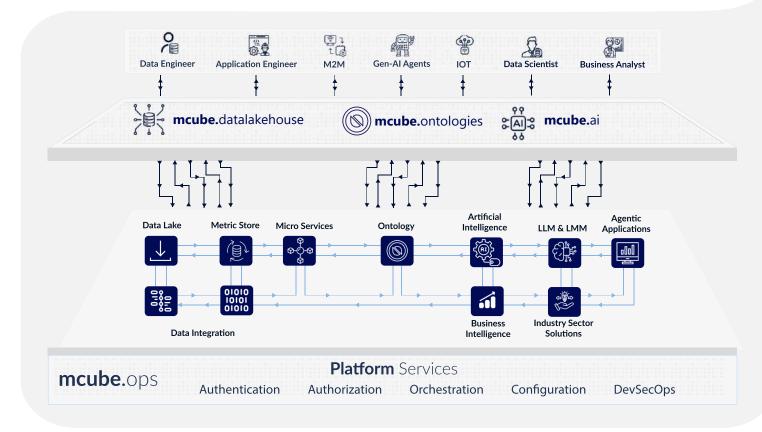
In short, semantics turns data from anonymous bits into explainable knowledge. It unlocks interoperability — so data from new acquisitions can be mapped rapidly. It enforces consistency — ensuring metrics mean the same thing in every dashboard. And it provides the logical substrate on which AI can safely automate decisions.





## 2. The mcube™ Philosophy: Meaning Driven by Design

Many vendors retrofit ontologies as overlays on existing lakes. mcube™ takes the opposite stance: the ontology is the organizing principle from day one. Every component—ingestion, storage, query, analytics—consults the ontology before acting. If a new data source arrives, mcube™ first maps its schema to the canonical vocabulary. If a column does not match any current concept, a governance workflow invites data stewards to extend the ontology, ensuring that ambiguity never seeps in unnoticed.



## The platform is composed of three tightly integrated layers:

- mcube.data provides scalable Lakehouse storage with ACID transactions. As files land, lightweight pipelines assign them to ontological concepts and store contextual metadata, such as provenance and units of measure.
- mcube.semantic, powered by the SemantX engine, hosts the ontology library and executes reasoning across the evolving knowledge graph. It offers a comprehensive suite of RESTful endpoints that allow graph specialists, relational analysts, and application developers query, traverse, and manipulate semantic data using familiar HTTP-based calls.
- mcube.ai surfaces intelligence through Machine Learning workbenches, Generative AI services, and a low-code studio called ezeXtend, enabling (business users assemble applications that are automatically ontology-aware.

Surrounding these layers are Genie microservices—reusable Application Programming Interfaces (APIs) for semantic search, entity resolution, and real-time alerting—and a governance envelope enforcing role-based security, auditing, and Model Operations (MLOps) lifecycle controls.

By embedding semantics at the kernel, mcube<sup>™</sup> realizes the FAIR principles—Findable, Accessible, Interoperable, Reusable—without bolt-on catalogues or laborious metadata campaigns. With mcube<sup>™</sup>, data is born FAIR.



## 3. Core Capabilities Explained in Depth

#### 3.1 Semantic Ingestion and Harmonization

Conventional pipelines treat metadata as a comment field; mcube™ treats it as a first-class citizen. The moment a file, stream, or Application Programming Interface (API) feed enters mcube™.data, a semantic mapping routine inspects each column, aligns it to the ontology, and records the mapping as machine readable rules. Units, synonyms, and provenance are stored alongside the data. For example, an Internet of Things (IoT) sensor emits the string "TempC." mcube™ recognizes this as Temperature in Celsius because the ontology already defines that term. It links the sensor to the physical asset, the asset to a production unit, and the unit to a facility. Downstream analytics need not rediscover these links; they exist from ingestion.

## 3.2 Dynamic Knowledge Graph and Reasoning

As the graph grows, mcube™.semantic performs reasoning—logical inference that enriches the graph with implied facts. If data says "Pump P7 is part of Unit U3" and "Unit U3 is part of Plant A," the reasoner deduces "Pump P7 is part of Plant A." Such transitive closure may sound trivial, yet it enables far-reaching queries: List all pumps in Plant A returns Pump P7 without manual joins. Reasoning also detects inconsistencies—if a part is asserted to belong to two mutually exclusive assemblies, the ontology rules flag a violation before bad data propagates.

### 3.3 Graph-Aware Machine Learning

Feature engineering accounts for a large share of ML effort. mcube™ slashes that by generating graph-derived features automatically: node centrality, path lengths, and neighborhood statistics. Predictive maintenance models, for example, gain context such as "How many similar pumps failed within one year?" or "Is this pump downstream of a heat exchanger with known fouling issues?" In pilot projects these graph features improved F-score by up to 20% compared with flat features alone.

## 3.4 Generative AI with Retrieval Augmented Generation

Generative AI systems like GPT-4 can draft answers in natural language but risk hallucinating facts. mcube™ mitigates that through knowledge graph-based Retrieval Augmented Generation (RAG). When a user asks

## "Which three products had the fastest revenue growth last quarter and who manages them?"

An orchestration chain first queries the knowledge graph for accurate revenue data and organizational ownership. It packages those facts—Product A grew 12%, Product B 11%, Product C 10%; managers are Alice, Bob, and Carol—and feeds them into the LLM prompt. The LLM then composes a fluent paragraph, citing each fact's Internationalized Resource Identifier (IRI). The final answer is both readable and verifiable. Industry benchmarks show that RAG can raise the factual accuracy of LLM outputs from roughly 60% to over 90%, with knowledge graphs further enhancing precision, completeness, and consistency by providing structured, verifiable context.

## 3.5 Low-Code Application Assembly with ezeXtend

Business agility requires more than analytics dashboards; it needs operational applications that adapt to new definitions. ezeXtend provides a drag-and-drop canvas where components—forms, charts, and graph viewers—already understand ontology terms. A data steward can design an Incident Report form by selecting the Incident class; ezeXtend autogenerates the correct fields with picklists populated from reference vocabularies. When the ontology evolves—say a new severity level is added—every form and rule downstream updates automatically. This keeps applications aligned with data governance without recode.

### 3.6 Security, Governance, and Audit

Regulators increasingly demand end-to-end lineage. mcube™ records every transformation, query, and model inference with timestamps and user identity. Audit teams can trace a dashboard figure back through the knowledge graph to its raw source via a clickable chain. Role Based Access Control (RBAC) restricts who can view which concept or attribute; Attribute Based rules go further by incorporating user properties, such as department or location. Sensitive personal data can be masked on output, encrypted at rest, and even segregated into privacy zones. These controls convert compliance from a manual chore into an automated guarantee.



## 4. Cross-Industry Panorama

Semantics add measurable value across domains. The snapshots that follow show how a shared ontology and knowledge graphs turn raw data into foresight and profit.

# 4.1 Refineries & Petrochemicals – Semantic twins for safer plants

By fusing real-time sensor streams, corrosion images, P&ID diagrams, and maintenance logs, reliability teams detect heat-exchanger fouling that precedes pump cavitation and cut unplanned downtime by 20%, while Occupational Safety and Health Administration (OSHA) violations are flagged automatically.

## 4.2 Life Sciences & Healthcare – R&D acceleration and personalized care

**Discovery:** a multi-omics graph lets researchers traverse "gene  $\rightarrow$  pathway  $\rightarrow$  compound  $\rightarrow$  trial" in one hop, shaving months off hypothesis generation.

**Clinical:** patient graphs that merge omics, Electronic Health Record (EHR) and wearable data match trial inclusion criteria in seconds and lift enrollment rates by double digits.

## 4.3 Manufacturing & Quality – Inline defect prediction

Graph-aware ML enriched with materials and process ontologies spots micro-defect clusters after only a few bad lots, boosting yield by 5–8% and shrinking root-cause analysis time.

## **4.4 Aviation – Predictive maintenance** with route context

Linking aircraft telemetry, flight plans, weather, and part lifecycles predicts auxiliary-power-unit failures three days in advance, halving unscheduled APU incidents and lowering fuel burn through smarter routing.

## 4.5 Government & Smart Cities – Joined-up public-service intelligence

Correlating 311 calls, mobility feeds, and epidemiology signals detects COVID-19 hotspots 48 hours earlier, enabling targeted test-kit deployment and reducing response costs.

# **4.6 Insurance – Network-aware fraud** analytics

Graphs that connect claimants, providers, vehicles, and social ties double complex-fraud detection efficiency while letting genuine claims settle faster.

# 4.7 Retail & CPG – Demand sensing and resilient supply chains

Product-customer-supplier graphs infused with micro-weather and social-sentiment data lift forecast accuracy from 75 %  $\rightarrow$  85 %, trim safety stock, and trigger real-time substitutions when upstream disruptions occur.

# 5. Implementation Journey

The path to a semantic enterprise need not be disruptive. Organizations typically begin with a discovery workshop to draft a lightweight ontology for one or two high-value domains. A small pilot ingests data under that ontology and delivers a visible win—perhaps a dashboard that reconciles finance and sales metrics, or a predictive model that outperforms the status quo. Success breeds confidence: the ontology expands, more data sources enter, and semantic query patterns proliferate. Within twelve months, many teams reach a stage where cross-domain questions—once impossible—are answered in minutes. Governance bodies then institutionalize ontology stewardship, and the semantic layer becomes the default integration pattern for new projects.

Early engagement with security, privacy, and compliance officers is crucial. mcube™'s fine-grained access policies should be configured up front so trust grows alongside capability. Cloud-native deployment enables horizontal scaling; DevOps pipelines can treat ontology artifacts as version-controlled code, promoting them through environments just like software.



## 6. Return on Investment

#### The financial upside of semantic integration arises from three levels:



Operational efficiency: Heartex data indicates that data scientists reclaim up to 30% of their time when consistent semantics replace manual cleansing. Multiply that by labor rates across an analytics team, and the savings become clear.



Revenue lift: Improved forecast accuracy and personalized marketing lead directly to higher sales. Studies in consumer goods forecasting tie a 15% accuracy boost to profit gains of roughly 3%.



口管 Risk reduction: Knowledge-graph-based fraud analytics at one insurer improved detection efficiency by over 100%. Prevented fraud represents pure bottom-line protection.

Add intangible benefits—executive trust in numbers, faster time-to-insight during crises—and the incentive for semantic investment strengthens. Organizations commonly see payback well within two years; leaders reach positive cash flow in under twelve months when they roll out multiple use cases in parallel.

## 7. Glossary



#### **Artificial Intelligence (AI)**

Computer systems capable of tasks that normally require human intelligence, such as learning and decision-making.



#### **Machine Learning (ML)**

A subset of AI that uses statistical techniques to enable systems to learn from data.



#### **Generative Al**

Al models that create new content (text, images) in response to prompts; Large Language Models like GPT4 fall into this category.



#### Large Language Model (LLM)

A deep learning model trained on large text corpora that can generate humanlike language.



#### Ontology

A machine-readable specification of concepts and relationships within a domain.



## Knowledge Graph

A network of entities and their relationships based on an ontology.



#### Retrieval Augmented Generation (RAG)

A technique that grounds LLM outputs in factual data retrieved from an external source.



## Role-Based Access Control (RBAC)

A security approach that grants permissions according to user roles.



#### **MLOps**

Practices for managing the lifecycle of Machine Learning models.



#### **Internationalized Resource Identifier (IRI)**

A globally unique string that identifies a concept or instance in an ontology.

## 8. References

- Gartner, Emerging Technologies: Knowledge Graphs Shine in Data Fabric Designs, May 2024.
- Heartex, Label Studio Community Survey 2024, October 2024.
- Rogers, R., Wired: How Retrieval Augmented Generation Tames LLM Hallucinations, June 2024.
- Manokhin, V., Medium: Forecast Accuracy and SupplyChain Profit, April 2025.
- Memgraph, Case Study: Fraud Ring Detection in U.S. Insurance, February 2022.



TCG Digital is the digital and Al arm of The Chatterjee Group (TCG), a multi-billion-dollar conglomerate with a diverse portfolio spanning Pharmaceuticals, Biotech, Petrochemicals, and Real Estate across the US, EU, and South Asia. Our umbrella includes companies such as LabVantage, Lummus Technology, and TCG Lifesciences.

At TCG Digital, we are driven by our mantra of delivering "Velocity to Value," helping enterprises transform faster and smarter. Our Al analytics platform, mcube™, is at the heart of these transformations. We enable organizations to unlock the full potential of their data, and by seamlessly integrating Al/ML capabilities into business processes, we empower them to accelerate their digital transformation journeys, enhance agility, and drive impactful results.

Our service portfolio spans big data strategy, Al/ML, advanced analytics, cloud and microservices, enterprise mobility, application development, automation, and security. With a global presence—including offices in Somerset, New Jersey—and a team of more than 1,500 digital professionals, TCG Digital helps clients achieve Velocity to Value through innovative solutions built on its award-winning mcube™ platform.

:::mcube

mcube™ is TCG Digital's flagship Data, Al, and Analytics platform. Built with a domain-driven design at the cross-roads of industry knowledge and digital prowess, our architecture is designed to handle the most disparate data landscapes. With Al 2.0, being at the heart of it, it combines powerful and advanced models to solve the most complex business problems. The platform integrates mcube.ai, mcube.data, and mcube.ontologies, delivering Al capabilities and data management seamlessly through unified platform services. Utilizing hyper contemporary technologies and deep domain expertise—particularly in semantic technologies, ontologies, and knowledge graphs— mcube™ delivers end-to-end digital transformation initiatives.

www.tcg digital.com/tcg-mcube

www.tcgdigital.com

#### **Contact Us**

