

Original Article

LLM-Native Data Architecture: How Foundation Models Change Data Storage and Access

Dr. Ramesh Babu¹, Dr. Deepika Arora²

¹Department of Information Technology, National Institute of Public Administration, India

²School of Cybersecurity and Digital Governance, Government Technology University, India

Abstract: *The rapid adoption of large language models (LLMs) and other foundation models is fundamentally reshaping how modern information systems store, retrieve, and reason over data. Traditional data architectures—designed primarily for deterministic query execution, rigid schemas, and structured workloads—are increasingly misaligned with the probabilistic, semantic, and context-driven access patterns introduced by LLM-powered applications. This paper introduces and formalizes the concept of LLM-native data architecture, an architectural paradigm in which foundation models and their learned representations are treated as first-class components of the data layer rather than peripheral consumers of data services. LLM-native data architecture is motivated by the observation that contemporary AI systems no longer interact with data through narrowly defined SQL queries or static APIs alone. Instead, they rely on semantic similarity search, retrieval-augmented generation (RAG), multimodal fusion, and iterative reasoning over large corpora of heterogeneous data. These workloads impose new requirements on data systems, including support for dense vector representations, hybrid symbolic-semantic indexing, dynamic schema evolution, provenance-aware retrieval, and cost-aware orchestration across model backends. As a result, existing architectures—centered on relational databases, data lakes, and keyword-based search—must be extended or reimaged to meet these demands. In this paper, we propose a comprehensive architectural framework for LLM-native data systems. We define core design principles such as representation-first storage, hybrid indexing, contextual chunking, model-aware query planning, and semantic caching. Building on these principles, we describe key architectural components, including a semantic persistence layer for managing raw data and derived model artifacts, vector index engines optimized for large-scale similarity search, hybrid query planners that account for model context limits and cost constraints, and governance layers that ensure privacy, security, and auditability. We emphasize the importance of provenance tracking and verifiable retrieval to mitigate hallucinations and establish trust in model-generated outputs.*

The paper further examines data modeling strategies that support multi-granularity storage, versioned embeddings, and incremental reprocessing as models evolve. We analyze indexing and retrieval techniques that combine approximate nearest neighbor search with predicate filtering and cross-encoder reranking to balance accuracy, latency, and cost. Operational considerations—including consistency models, caching strategies, scalability, observability, and privacy-preserving retrieval—are discussed in detail, highlighting the trade-offs inherent in deploying LLM-native systems at enterprise scale. Finally, we outline evaluation methodologies and open research challenges, such as embedding space alignment across model upgrades, efficient private vector search, trustworthy attribution mechanisms, and standardization for interoperability. By articulating both practical system designs and a forward-looking research agenda, this paper aims to serve as a foundational reference for researchers and practitioners seeking to build robust, scalable, and trustworthy data architectures in the era of foundation models.

Keywords: *LLM-Native Data Architecture; Foundation Models; Semantic Retrieval; Vector Databases; Hybrid Indexing; Retrieval-Augmented Generation; Data Governance; Embeddings; AI Data Systems; Model-Aware Query Planning*

I. INTRODUCTION

The emergence of large language models (LLMs) and other foundation models represents one of the most significant paradigm shifts in the history of data-intensive computing. Trained on massive corpora and capable of general-purpose reasoning, summarization, and generation, these models are rapidly becoming the primary interface through which users and applications interact with data. From enterprise knowledge assistants and intelligent search systems to autonomous agents and decision-support platforms, LLM-powered applications are redefining expectations around how data should be accessed, interpreted, and synthesized. Traditional data architectures were designed under a fundamentally different set of assumptions. Relational databases and data warehouses optimized for structured schemas, deterministic queries, and transactional

consistency. Data lakes and lakehouses expanded scalability and flexibility but still relied on explicit schemas, batch-oriented processing, and query engines rooted in symbolic logic. Search engines introduced keyword-based retrieval, but their indexing and ranking models were largely orthogonal to generative reasoning. In contrast, LLM-centric workloads rely on semantic similarity, probabilistic inference, and contextual reasoning across unstructured and semi-structured data. This mismatch has exposed architectural limitations that cannot be fully addressed through incremental extensions to existing systems.

A defining characteristic of LLM-based applications is their reliance on learned representations—dense vector embeddings that encode semantic meaning rather than explicit structure. Instead of querying for exact matches, applications retrieve information based on conceptual similarity, relevance, and contextual alignment with a natural language prompt. Moreover, LLMs operate under practical constraints such as finite context windows, non-deterministic outputs, and significant computational cost. These constraints directly influence how data must be stored, indexed, retrieved, and assembled into prompts. Consequently, data architecture decisions increasingly depend on model behavior, inference cost, and retrieval quality rather than purely on storage efficiency or query throughput. In response to these shifts, a growing ecosystem of vector databases, embedding stores, and retrieval-augmented generation (RAG) pipelines has emerged. While these tools provide critical building blocks, they are often integrated as add-ons rather than as cohesive, end-to-end architectures. Many current systems treat LLMs as external services layered on top of conventional data platforms, leading to fragmented pipelines, duplicated storage, weak governance, and limited observability. As deployments scale, these shortcomings manifest as rising costs, inconsistent behavior across model upgrades, privacy risks, and difficulties in verifying or auditing model-generated outputs.

This paper argues that addressing these challenges requires a more fundamental rethinking of the data layer itself. We introduce the concept of LLM-native data architecture, in which foundation models and their artifacts—such as embeddings, summaries, extracted entities, and provenance traces—are treated as first-class data assets. In an LLM-native architecture, data storage and access mechanisms are explicitly designed around semantic retrieval, contextual assembly, and model-aware optimization. Rather than forcing LLM workloads to adapt to legacy data abstractions, the architecture adapts to the operational realities of foundation models. An LLM-native approach has several important implications. First, data is stored at multiple semantic granularities, enabling efficient retrieval of passages, summaries, or structured facts depending on the task and available context budget. Second, indexing strategies combine vector similarity search with symbolic filters, allowing precise control over scope, access rights, and relevance. Third, query planning becomes inherently model-aware, accounting for token limits, retrieval cost, and the trade-offs between additional context and in-model reasoning. Fourth, governance mechanisms must extend beyond raw data to encompass derived representations and generated outputs, ensuring privacy, compliance, and traceability.

Equally important is the need for trust and verifiability. LLMs are known to hallucinate or produce confident but incorrect statements when retrieval is incomplete or ambiguous. Without explicit provenance linking model outputs to authoritative sources, organizations risk deploying systems that are opaque and difficult to audit. An LLM-native data architecture therefore embeds provenance tracking, citation mechanisms, and validation workflows directly into the data access layer, enabling evidence-backed generation and post hoc verification. The motivation for this work is both practical and forward-looking. Practically, enterprises across domains—including finance, healthcare, manufacturing, and research—are struggling to operationalize LLMs at scale while maintaining cost control, data quality, and regulatory compliance. Forward-looking, the trajectory of AI systems points toward increasingly autonomous agents that plan, retrieve, and act over data with minimal human intervention. Such systems will place even greater demands on data architectures, requiring tighter integration between reasoning, retrieval, and governance. This paper makes three primary contributions. First, it provides a clear conceptual definition of LLM-native data architecture and articulates its core design principles. Second, it presents a detailed architectural framework that integrates semantic persistence, hybrid indexing, retrieval augmentation, model orchestration, and governance into a cohesive system. Third, it identifies key evaluation metrics and open research challenges that must be addressed to advance the state of the art.

By synthesizing ideas from databases, information retrieval, and machine learning systems, this work aims to bridge disciplinary boundaries and provide a common reference point for architects, researchers, and practitioners. As foundation models continue to reshape how humans and machines interact with information, LLM-native data architecture offers a principled path toward building scalable, trustworthy, and future-proof data systems.

II. MOTIVATION AND PROBLEM STATEMENT

Traditional data architectures were conceived and optimized for an era in which data access patterns were largely deterministic, well-structured, and explicitly defined by application logic. Online Transaction Processing (OLTP) systems prioritize strong ACID guarantees, normalized schemas, and predictable access paths to ensure correctness and consistency under concurrent updates. Online Analytical Processing (OLAP) platforms and data warehouses extend these principles to large-scale analytics, favoring columnar storage, batch-oriented processing, and carefully curated schemas that support aggregations and reporting. Even more recent paradigms such as data lakes and lakehouses, while more flexible in terms of schema enforcement and data formats, still fundamentally assume that queries are symbolic, explicit, and executed by humans or predefined applications. These assumptions no longer hold when large language models become primary consumers and producers of data. LLM-driven workloads introduce a qualitatively different mode of interaction with data. Rather than issuing precise queries against predefined schemas, models operate on semantic intent expressed in natural language. The goal is no longer to retrieve exact matches, but to identify information that is conceptually relevant, contextually appropriate, and useful for downstream reasoning or generation. This shift from symbolic retrieval to semantic retrieval fundamentally challenges existing storage and indexing mechanisms. Keyword-based search, primary keys, and relational joins are insufficient for capturing the latent meaning encoded in text, images, audio, and other unstructured sources. Instead, dense vector representations derived from neural models become the primary access path, requiring data systems to efficiently manage similarity search at scale.

In addition to semantic retrieval, LLM workloads rely heavily on contextual aggregation. Most production deployments today use retrieval-augmented generation, in which relevant pieces of external data are dynamically selected and injected into the model's prompt to ground its responses. This introduces new architectural requirements that traditional systems were never designed to address. Data must be stored and indexed in a way that supports flexible chunking, ranking, and recomposition under strict context window constraints. Long documents cannot be treated as monolithic records; they must be segmented, summarized, and reassembled in multiple forms depending on the task, the query, and the model's token budget. As a result, data access is no longer a single query-response interaction, but a multi-stage pipeline that interleaves retrieval, transformation, and reasoning. The growth of model context windows further complicates this picture. While modern LLMs can process increasingly large inputs, context remains finite and expensive. This forces data systems to make intelligent decisions about what information to retrieve, at what granularity, and in what order. Such decisions are inherently model-aware: they depend on token limits, attention behavior, and the trade-offs between providing more raw context versus compressed summaries. Traditional query planners optimize join order and index usage based on cost models tied to CPU, I/O, and memory. In contrast, LLM-centric systems must also reason about inference cost, latency, and the diminishing returns of additional context. Existing data architectures lack abstractions for expressing and optimizing these concerns.

Another defining characteristic of LLM workloads is their inherently multimodal nature. Modern foundation models increasingly operate across text, images, audio, video, and structured signals within a single reasoning process. Data architectures that evolved around homogeneous data types struggle to support such workflows. Relational databases excel at structured records, search engines focus on text, and media stores manage binary assets, but LLM applications require unified access across modalities. Semantic similarity must be computed not only within a single modality but across them, such as retrieving an image based on a textual description or linking an audio transcript to related documents and structured metadata. This requires consistent representation, indexing, and retrieval strategies that cut across traditional system boundaries. Cost and performance considerations further expose the limitations of existing approaches. Generating embeddings, summaries, and other derived artifacts is computationally expensive, particularly at scale. LLM-native workloads therefore rely on amortization strategies, caching, and reuse of model outputs. Unlike traditional query result caches, these artifacts are tightly coupled to specific model versions, prompts, and preprocessing pipelines. When models are upgraded or prompts change, cached representations may become stale or incompatible, necessitating selective reprocessing. Conventional data systems provide limited support for managing such model-dependent derived data, leading to inefficiencies, duplicated storage, and operational complexity.

Schema management presents another fundamental challenge. Traditional systems assume that schemas are defined a priori by humans and evolve slowly through controlled migrations. In LLM-driven environments, however, schemas can emerge dynamically as models extract entities, infer relationships, or synthesize new concepts from unstructured data. These emergent structures may not fit neatly into predefined tables or ontologies, yet they are valuable for downstream reasoning and retrieval. Supporting such fluid, model-driven schema evolution without sacrificing queryability, governance, or backward compatibility is

beyond the scope of most existing architectures. Trust and verifiability introduce perhaps the most critical gap. LLMs are powerful but imperfect; they can hallucinate facts, conflate sources, or generate plausible but incorrect statements when retrieval is incomplete or ambiguous. In traditional data systems, results can usually be traced back to specific records and queries. In contrast, LLM outputs are probabilistic and synthesized, making it difficult to understand how a particular answer was produced. Without explicit support for provenance tracking, citation, and auditability, organizations risk deploying systems that are opaque, untrustworthy, and difficult to regulate. This is especially problematic in domains such as healthcare, finance, law, and scientific research, where accountability and explainability are essential.

These challenges collectively point to a fundamental architectural problem rather than a collection of isolated feature gaps. Simply adding vector search capabilities or plugging an LLM into an existing data platform does not address the deeper misalignment between traditional data assumptions and model-centric workloads. Storage systems, indexes, and access patterns must be redesigned to natively support semantic representations, contextual retrieval, multimodal fusion, and model-aware optimization. They must also provide robust mechanisms for caching, versioning, governance, and provenance that reflect the realities of foundation model inference. The central problem, therefore, is how to design data systems that can deliver low-latency and cost-effective access to semantically relevant information while remaining robust, auditable, and adaptable in the face of rapidly evolving models. Such systems must bridge the gap between symbolic data management and probabilistic reasoning, enabling LLMs to interact with large-scale data in a controlled and trustworthy manner. Addressing this problem requires rethinking not only individual components such as storage engines or indexes, but the overall architecture of the data layer itself, placing foundation models and their representations at the core of data access and management.

III. DESIGN PRINCIPLES OF LLM-NATIVE ARCHITECTURES

LLM-native data architectures are guided by a distinct set of design principles that emerge directly from the interaction between foundation models and data systems. These principles are not incremental optimizations of existing database design practices; rather, they reflect a shift in perspective in which learned representations, probabilistic reasoning, and model constraints shape how data is stored, indexed, accessed, and governed. Together, they provide a conceptual foundation for building data platforms that can reliably support LLM-centric workloads at scale. At the core of an LLM-native architecture is a representation-first mindset. In traditional systems, raw data is the primary asset, and all derived structures—indexes, materialized views, or caches—are considered secondary optimizations. In contrast, LLM-driven applications depend fundamentally on model-derived representations such as embeddings, summaries, extracted entities, and provenance traces. These representations are not merely performance accelerators; they define how the system understands and retrieves information. Treating them as first-class data requires persistent storage, explicit versioning, metadata management, and lifecycle controls comparable to those applied to raw source data. This shift enables consistent reuse of expensive model outputs, controlled reprocessing when models change, and systematic governance over both original and derived artifacts.

Effective retrieval in LLM-native systems also depends on hybrid indexing. Semantic similarity search using dense vector embeddings is essential for capturing conceptual relevance, but it is insufficient on its own. Real-world applications impose constraints related to tenant isolation, security policies, temporal scope, and structured attributes that cannot be expressed through vector similarity alone. Hybrid indexing strategies integrate vector indexes with traditional inverted indexes and relational or key-based indexes, allowing systems to combine semantic relevance with precise symbolic filtering. This fusion enables high-precision retrieval while maintaining control over access boundaries and query intent, bridging the gap between neural representations and classical data management techniques. Contextual chunking is another foundational principle. Large documents and multimodal assets cannot be consumed directly by LLMs without preprocessing, as models operate on finite context windows and attend unevenly to long inputs. LLM-native architectures therefore store data at multiple semantic and structural granularities, such as tokens, sentences, paragraphs, full documents, and generated abstracts. Maintaining explicit mappings between these granularities allows the system to retrieve and assemble context dynamically based on task requirements. Fine-grained chunks support precise attribution and targeted retrieval, while coarser representations enable efficient summarization and ranking. This multi-level storage model is essential for balancing retrieval accuracy, coherence, and token efficiency.

Query planning in LLM-native architectures must be explicitly model-aware. Traditional query optimizers focus on minimizing I/O, CPU usage, or execution time under deterministic cost models. In contrast, LLM-centric workloads introduce additional dimensions, including token budgets, inference latency, and the marginal utility of additional context. A model-aware query planner reasons about whether a question should be answered through additional retrieval, in-prompt reasoning, or

cached responses, and how much context to supply at each stage. Such planners must account for the characteristics of specific models, including context window size, attention behavior, and cost per token, making query optimization an inherently AI-aware process. Caching and memoization play a central role in making LLM-native systems economically viable. Generating embeddings, reranking candidates, and producing high-quality responses are computationally expensive operations. Unlike traditional result caches, however, cached artifacts in LLM systems are tightly coupled to model versions, prompts, and preprocessing pipelines. Effective memoization therefore requires versioned caches keyed by model identifiers and prompt fingerprints, enabling safe reuse while avoiding inconsistencies when models or prompts change. By elevating caching to a first-class architectural concern, LLM-native systems can significantly reduce latency and cost without sacrificing correctness.

Provenance and auditability are essential principles for trust and operational reliability. Because LLM outputs are synthesized and probabilistic, it is often unclear how a specific response was produced unless the system explicitly records retrieval sources, prompt history, and model versions. LLM-native architectures embed provenance tracking directly into the data access pipeline, linking generated outputs to the underlying chunks, summaries, and representations that informed them. This enables verification, debugging, and compliance auditing, and supports evidence-backed generation in which responses can be traced to authoritative sources. Privacy and access control requirements are amplified in LLM-driven systems. Learned representations can inadvertently encode sensitive information, and semantic retrieval can surface data that would not be accessible through explicit queries. LLM-native architectures must therefore support fine-grained access policies that apply not only to raw data but also to embeddings and derived artifacts. Techniques such as encryption, tenant-aware indexing, differential privacy, and private retrieval mechanisms become integral components of the data layer, ensuring that semantic power does not come at the expense of confidentiality or regulatory compliance.

Cost-aware routing is another defining principle. Modern deployments often involve a heterogeneous mix of local models, specialized small models, and remote foundation model APIs, each with different cost, latency, and accuracy characteristics. An LLM-native architecture dynamically routes requests based on these factors, selecting cached answers, local inference, or external calls according to policy and service-level objectives. This adaptive routing capability allows systems to scale economically while maintaining acceptable performance. Schema evolution in LLM-native systems must be flexible and continuous. As models extract new entities, infer relationships, or generate structured interpretations from unstructured data, schemas may evolve organically rather than through centralized design. Supporting such evolution while preserving backward compatibility and queryability requires abstractions that decouple logical structure from physical storage. LLM-native architectures therefore embrace schema-on-read approaches enriched with model-generated metadata, enabling emergent structure without breaking existing applications. Finally, operational observability underpins all other principles. Semantic systems are subject to drift as data distributions change, models are retrained, and embedding spaces evolve. LLM-native architectures must continuously monitor retrieval quality, embedding distributions, and model behavior to detect degradation early. Observability extends beyond traditional metrics such as latency and throughput to include semantic measures, such as relevance scores and citation coverage, providing operators with the insight needed to maintain reliable and trustworthy systems over time.

Together, these design principles define an architectural blueprint that aligns data systems with the capabilities and constraints of foundation models. By embedding semantic understanding, model awareness, and governance directly into the data layer, LLM-native architectures enable a new class of intelligent, scalable, and accountable data-driven applications.

IV. ARCHITECTURAL COMPONENTS

An LLM-native data architecture is composed of tightly integrated layers that collectively support semantic storage, efficient retrieval, model-aware planning, and trustworthy governance. Rather than treating data management, retrieval, and model serving as loosely coupled subsystems, this architecture unifies them into a coherent stack optimized for foundation-model-driven workloads. Conceptually, the architecture can be organized into three major subsystems: the Semantic Data and Indexing Layer, the Intelligent Retrieval and Model Execution Layer, and the Governance, Security, and Control Layer. Each subsystem encapsulates multiple responsibilities while maintaining clear interfaces that enable scalability, observability, and evolution over time.

A. Semantic Data and Indexing Layer

The Semantic Data and Indexing Layer forms the foundation of an LLM-native architecture by managing both raw data and the rich set of derived representations produced by models. At its core lies the Semantic Persistence Layer (SPL), which stores original documents alongside multiple semantic derivatives, including embeddings at various granularities, canonical

summaries, extracted entities, vector fingerprints, and detailed provenance metadata. Unlike traditional storage layers that treat derived artifacts as ephemeral or disposable, the SPL treats these representations as durable, versioned assets. This approach enables consistent reuse, controlled reprocessing when models evolve, and reliable governance over the entire semantic lifecycle of data.

Data within the SPL is organized around explicit relationships between documents, chunks, embeddings, and metadata. Long-form content is decomposed into semantically meaningful chunks—such as sentences or paragraphs—each of which can be independently embedded, summarized, and indexed. Persistent links between these layers allow the system to trace which specific chunks contributed to a given retrieval result or generated response. This traceability is essential for auditability and debugging, particularly in regulated environments. Closely coupled with the SPL is the Vector Index Engine, which provides high-performance semantic retrieval over dense representations. The index engine supports approximate nearest neighbor search at scale, enabling efficient similarity queries across millions or billions of vectors. To accommodate real-world workloads, the index is designed for dynamic updates, incremental ingestion, and per-tenant sharding. Hybrid indexing capabilities allow vector similarity search to be combined with scalar filters derived from structured metadata, ensuring that semantic relevance does not override access constraints or contextual boundaries.

Performance optimization is a key concern at this layer. Frequently accessed vectors are maintained in memory-resident structures, such as hierarchical navigable small world graphs, while colder data is stored in disk-based indexes using compressed or quantized representations. GPU acceleration may be employed to further reduce latency for dense retrieval workloads. By integrating semantic storage and indexing into a unified layer, LLM-native architectures enable fast, scalable, and semantically rich access to data that traditional storage systems cannot provide.

B. Intelligent Retrieval and Model Execution Layer

Sitting above the semantic data foundation is the Intelligent Retrieval and Model Execution Layer, which orchestrates how queries are interpreted, how context is assembled, and how models are invoked. This layer is responsible for translating user intent—expressed as natural language prompts, structured filters, or programmatic queries—into optimized retrieval and generation workflows that respect both system constraints and model capabilities. At the heart of this layer is the Hybrid Query Planner. Unlike conventional query optimizers, the planner operates under a model-aware cost framework. It considers not only retrieval latency and compute cost, but also token budgets, model inference pricing, and service-level objectives. Given a query, the planner determines which chunks to retrieve, at what granularity, and in what sequence, balancing the trade-off between retrieving more external context and relying on in-model reasoning. This planning process is adaptive and may vary depending on the target model, workload criticality, and available caches.

The Retrieval Augmentation Engine executes the planner's decisions through a multi-stage retrieval pipeline. Initial candidate selection typically uses fast, coarse-grained vector similarity search to identify a broad set of potentially relevant contexts. These candidates are then refined through deduplication, ranking, and optional reranking using lightweight cross-encoder models that provide higher precision at increased cost. The final selection is constrained by the model's context window and may include dynamically generated summaries to compress less critical information. This staged approach ensures both efficiency and relevance in retrieval-augmented generation workflows. Model execution is abstracted through a unified Model Interface and Serving Layer. This abstraction decouples the architecture from specific model implementations, enabling seamless integration of local open-source models, cloud-based foundation model APIs, and specialized task-specific models. The serving layer supports orchestration patterns such as tool calling, chained reasoning, and function invocation, while preserving end-to-end provenance. By consolidating retrieval and model execution into a cohesive layer, the architecture enables flexible, cost-aware, and extensible AI-driven data access.

C. Governance, Security, and Control Layer

The Governance, Security, and Control Layer ensures that the semantic power of LLM-native architectures is deployed in a responsible, auditable, and compliant manner. As models gain the ability to infer and synthesize information across large corpora, traditional access control mechanisms become insufficient. Governance must extend beyond raw data to encompass embeddings, derived artifacts, and generated outputs. Access control policies in this layer are enforced consistently across storage, indexing, retrieval, and generation. Fine-grained permissions determine not only which documents a user can access, but also which embeddings can be queried and which contexts can be injected into model prompts. This prevents semantic retrieval from inadvertently bypassing security boundaries. Sensitive data may be protected using client-side encryption, tenant-aware indexing, or execution within secure enclaves, ensuring confidentiality even during model inference.

Provenance tracking is a central function of this layer. Every query, retrieval operation, and model invocation is logged with references to the underlying data, model version, and prompt configuration. These audit trails enable post hoc analysis, regulatory compliance, and systematic debugging of incorrect or unexpected outputs. In environments where trust is critical, such as healthcare or finance, provenance data supports evidence-backed generation, allowing responses to be traced to authoritative sources. Privacy-preserving analytics and monitoring are also addressed within this layer. Aggregate metrics on retrieval quality, model usage, and system performance can be computed using differential privacy techniques to avoid leaking sensitive information. Together, these governance mechanisms provide the operational safeguards required to deploy LLM-native data architectures at enterprise scale without sacrificing transparency, accountability, or trust.

Table 1: Core Components of an LLM-Native Data Architecture

Architectural Layer	Key Components	Primary Responsibilities	Data Artifacts Managed
Semantic Data and Indexing Layer	Semantic Persistence Layer, Vector Index Engine	Store raw and derived data, manage embeddings, enable semantic and hybrid retrieval	Documents, chunks, embeddings, summaries, metadata
Intelligent Retrieval and Model Execution Layer	Hybrid Query Planner, Retrieval Augmentation Engine, Model Serving Interface	Optimize retrieval, assemble context, orchestrate model inference	Retrieval sets, prompts, intermediate rankings, generated outputs
Governance, Security, and Control Layer	Access Control, Provenance Tracking, Privacy Mechanisms	Enforce security, auditability, compliance, and observability	Access policies, audit logs, provenance traces, privacy metrics

V. DATA MODELING AND STORAGE FORMATS

Data modeling in LLM-native architectures departs significantly from traditional schema-centric approaches. Because foundation models consume and produce information at multiple semantic levels, the data layer must support flexible granularity, heterogeneous representations, and continuous evolution as models and tasks change. This section describes how LLM-native systems structure data for semantic access, how different representation stores coexist within a unified architecture, and how versioning and updatability are handled without disrupting system stability.

A. Multi-Granularity Semantic Data Modeling

A defining requirement of LLM-native data architectures is the ability to represent the same source information at multiple semantic and structural granularities. Unlike conventional databases, where records are typically stored at a single level of abstraction, LLM workloads require data to be accessible at fine-grained, intermediate, and coarse-grained levels depending on the retrieval and reasoning task. This necessity arises from the constraints of model context windows, attention mechanisms, and the need for precise attribution in generated outputs. At the finest level, data is represented as tokens or spans, enabling precise extraction, citation, and attribution. Token-level representations are particularly important for tasks such as entity extraction, factual verification, and evidence-backed generation, where the system must identify exactly which fragment of text supports a given claim. Maintaining offsets and alignment with the original source ensures that generated responses can be traced back to authoritative locations within documents.

At an intermediate level, sentences and paragraphs form the primary unit of semantic retrieval. These chunks balance contextual coherence with retrieval precision, making them well suited for similarity search using embeddings. Paragraph-level chunks often serve as the default retrieval unit in retrieval-augmented generation pipelines, as they provide sufficient context for reasoning while remaining compact enough to fit within model input limits. Importantly, these chunks are not isolated artifacts; they are linked to both finer-grained spans and higher-level document representations through explicit parent-child relationships. At the coarsest level, entire documents, sections, or generated abstracts provide a high-level semantic overview. These representations are valuable for initial ranking, corpus-level navigation, and tasks that require broad understanding rather

than detailed evidence. Abstracts and document-level embeddings allow systems to quickly narrow the search space before performing more expensive fine-grained retrieval. Maintaining consistent mappings across all granularities allows the architecture to dynamically traverse levels of abstraction, selecting the most appropriate representation based on query intent, cost constraints, and model capacity.

B. Representation Stores for Semantic and Symbolic Access

LLM-native architectures rely on multiple specialized representation stores, each optimized for a different access pattern but integrated through shared identifiers and metadata. The most prominent of these is the embedding store, which manages dense vector representations used for semantic similarity search. Embeddings are stored primarily within vector indexes to support low-latency retrieval, but are also persisted in durable storage systems such as blob stores or column-oriented databases. This redundancy enables reindexing, auditing, and large-scale reprocessing when models or embedding strategies change. Each embedding is accompanied by rich metadata, including the model identifier, generation timestamp, dimensionality, and normalization scheme, ensuring that semantic compatibility can be assessed during retrieval. Complementing the embedding store is the symbolic representation store, which manages structured information extracted or inferred by models. Entities, relationships, and attributes are stored in relational tables or graph databases to support deterministic queries, joins, and rule-based reasoning. This symbolic layer plays a critical role in grounding LLM outputs, enabling the system to validate generated claims against authoritative structured data. By combining symbolic and semantic representations, the architecture supports both probabilistic retrieval and exact lookups within a single unified system.

In addition to embeddings and symbolic data, LLM-native architectures maintain stores for summaries and canonical forms. These artifacts provide compressed representations of larger content and standardized references for entities or concepts that may appear under multiple surface forms. Summaries are particularly valuable for reducing token usage during retrieval and for presenting concise information to users. Canonicalization helps mitigate ambiguity and duplication, improving both retrieval quality and downstream reasoning. Importantly, these representations are treated as durable, queryable data rather than ephemeral outputs, reinforcing the representation-first philosophy of the architecture.

C. Versioning, Evolution, and Updatability of Model-Derived Data

Continuous model evolution is an inherent characteristic of LLM-based systems, and data architectures must accommodate this dynamism without sacrificing stability or correctness. In LLM-native architectures, all model-derived artifacts—including embeddings, summaries, and extracted structures—are explicitly versioned. Each version is associated with the specific model configuration and preprocessing pipeline that produced it, enabling precise tracking of semantic lineage over time. When embedding models are upgraded, the resulting changes in embedding space can invalidate direct similarity comparisons with older representations. Rather than requiring disruptive full reprocessing, LLM-native architectures support progressive migration strategies. Re-embedding jobs can be scheduled incrementally, prioritizing frequently accessed data, while compatibility mappings or alignment transforms are applied to enable temporary coexistence of old and new embeddings. This approach minimizes downtime and allows systems to evolve continuously.

Table 2: Data Modeling and Storage Components in LLM-Native Architectures

Data Layer Component	Representation Type	Storage Format	Primary Purpose
Multi-granularity Store	Tokens, chunks, documents, abstracts	Object store + metadata index	Flexible retrieval and attribution
Embedding Store	Dense vectors with metadata	Vector index + blob/column store	Semantic similarity search
Symbolic Store	Entities, relations, attributes	Relational or graph database	Deterministic lookup and validation
Summary & Canonical Store	Compressed text, normalized entities	Document store	Token-efficient retrieval and display
Versioning Metadata	Model and artifact versions	Metadata catalog	Evolution, auditability, compatibility

Atomic updatability is another critical feature. Derived artifacts such as embeddings or canonical summaries can be updated independently of the original source document. By decoupling raw data from its representations, the architecture avoids unnecessary rewrites and preserves the integrity of historical records. Persistent identifiers and links ensure that updated representations are seamlessly incorporated into retrieval and generation workflows. Together, these versioning and update mechanisms provide the flexibility required to support rapidly advancing foundation models while maintaining a consistent and reliable data layer.

VI. INDEXING AND RETRIEVAL STRATEGIES

Efficient indexing and retrieval are central to the performance and accuracy of LLM-native architectures. Unlike traditional systems that rely exclusively on symbolic queries or full-text search, LLM-centric workloads demand hybrid strategies capable of combining semantic understanding, contextual reasoning, and structured filtering. The architecture must address challenges including high-dimensional vector search, token-limited model inputs, multi-stage retrieval, and cost-effective storage. This section describes three complementary strategies: hybrid filtering, two-stage retrieval with reranking, and context management with compression and quantization.

A. Hybrid Semantic and Predicate Filtering

High-precision retrieval in LLM-native systems requires the integration of semantic and structured search paradigms. Semantic retrieval relies on dense embeddings to identify content conceptually relevant to a query, while structured or symbolic constraints—such as tenant identifiers, temporal ranges, tags, or categorical attributes—remain essential for governance, scope limitation, and filtering. Hybrid filtering fuses these approaches by first restricting the candidate set according to predicates and metadata filters and subsequently ranking the reduced set using vector similarity measures. By combining structured constraints with semantic distance, the system can eliminate irrelevant or unauthorized candidates before performing computationally intensive similarity computations, thereby improving both retrieval precision and query efficiency. Hybrid filtering also supports dynamic weighting between structured and semantic criteria, allowing operators to adjust prioritization depending on the application context, regulatory requirements, or query intent.

B. Multi-Stage Retrieval and Reranking

A critical strategy for balancing speed, relevance, and cost is multi-stage retrieval. In the initial stage, candidate generation identifies a broad set of potentially relevant chunks using approximate nearest neighbor (ANN) search over embeddings. These embeddings may represent entire documents, paragraphs, or coarsened semantic representations, depending on query requirements. The candidate set is intentionally permissive to avoid missing relevant content while keeping computation manageable. The second stage employs reranking using a smaller contextual transformer model, often referred to as a cross-encoder, which scores each candidate relative to the query and other retrieved contexts. The cross-encoder is capable of modeling fine-grained interactions between query tokens and candidate content, producing a ranking that reflects not only semantic similarity but also contextual relevance, recency, and inferred trustworthiness. This staged approach separates computationally cheap retrieval from expensive scoring, allowing systems to scale to large corpora while ensuring that the final top-k candidates fed into the model are of high quality. Multi-stage retrieval also facilitates flexible adaptation to varying latency or token budget constraints, enabling the architecture to optimize trade-offs between speed, accuracy, and cost.

C. Context Assembly, Compression, and Quantization

Even after top-ranked candidates are identified, practical limitations arise due to model context windows and token budgets. Effective context management strategies include sliding windows over overlapping chunks, prioritization of the most relevant and trustworthy content, and on-the-fly summarization of less critical material. Sliding windows ensure coverage of important content without truncating key information, while overlapping chunks maintain continuity across segment boundaries. Chunk prioritization can incorporate multiple signals, including semantic relevance, recency, provenance confidence, and trust scores. If the token budget is insufficient to include all relevant content, compressed summaries of lower-priority chunks are used to retain essential information while conserving space.

Storage and retrieval efficiency are further enhanced through vector quantization and compression. Product quantization (PQ), optimized PQ (OPQ), and other retrieval-friendly compression schemes reduce storage footprint while preserving approximate nearest neighbor search quality. Summarized or distilled representations, which retain core facts but eliminate extraneous detail, allow rapid retrieval when exact quoting is unnecessary. Together, these strategies ensure that retrieved context is both semantically rich and computationally tractable, enabling LLMs to generate coherent and accurate outputs even under strict model constraints.

Table 3: Indexing and Retrieval Strategies in LLM-Native Architectures

Strategy	Key Components	Purpose	Implementation Notes
Hybrid Filtering	Predicate filters, vector similarity	Combine semantic relevance with structured constraints	Filter candidates by metadata first, then rank by embeddings
Multi-Stage Retrieval	ANN candidate generation, cross-encoder reranking	Balance retrieval speed and accuracy	Coarse retrieval first, fine-grained reranking second
Context Assembly & Compression	Sliding windows, prioritization, summarization, quantization	Fit context within token budgets while preserving information	Use overlapping chunks, compressed summaries, PQ/OPQ vector storage

VII. CONCLUSION

As foundation models continue to redefine the landscape of intelligent systems, the underlying data architectures that support them must evolve in tandem. Traditional data systems were designed around deterministic queries on structured schemas supported by symbolic indexes and relational storage. While those systems excelled at transactional integrity and analytical workloads, they were not built to serve the semantic, generative, and context-driven access patterns that characterize large language model (LLM)-centric applications. This paper has introduced the concept of LLM-native data architecture as a principled response to these emerging demands, articulating its foundational motivations, core design principles, architectural components, data models, indexing strategies, and operational considerations. In a world where LLMs increasingly interact with and reason over large, heterogeneous datasets, several underlying themes emerge. First, representation matters: raw data alone is insufficient for practical retrieval and reasoning. Instead, derived artifacts such as dense embeddings, summaries, and structured extractions become first-class citizens of the data layer. By integrating these representations directly into storage and indexing structures, LLM-native architectures unlock semantic retrieval capabilities that are both efficient and robust. Second, semantic understanding must be coupled with traditional governance and access controls. Semantic similarity search, if left unchecked, can inadvertently bypass security boundaries; hybrid filtering mechanisms that combine predicate logic with vector search provide the precision needed for controlled access.

At the heart of the architecture is a seamless interplay among three major subsystems: the Semantic Data and Indexing Layer, the Intelligent Retrieval and Model Execution Layer, and the Governance, Security, and Control Layer. The Semantic Data and Indexing Layer unifies raw content with model-derived representations such as multi-granularity chunks and vector indexes, enabling high-throughput neural search at scale. The Intelligent Retrieval Layer orchestrates context-aware query planning, multi-stage retrieval, and cost-sensitive model invocation to fit within finite context windows and performance constraints. Finally, the Governance Layer enforces access controls, tracks provenance, and preserves auditability—qualities that are critically important for trustworthy and compliant use of AI systems in regulated domains. The paper also outlined the importance of data modeling that spans multiple granularities, from tokens and spans up to full documents and abstracts. This flexible data modeling enables the system to adapt retrieval granularity to the task at hand, balancing precision, coherence, and model cost. Representation stores such as embedding databases, symbolic stores for structured facts, and canonical summary repositories collectively support hybrid reasoning that leverages both semantic similarity and exact symbolic queries. Through explicit versioning and update strategies, the architecture also accommodates the continuous evolution of models—a reality of modern AI environments—without sacrificing data hygiene or system performance.

A central theme throughout the exposition is the necessity of model-aware design across all layers. Query planners, retrieval engines, and context assemblers cannot be agnostic to the cost, capacity, and semantics of the models they serve. Optimizing retrieval under token budget constraints, batching embeddings, prioritizing high-confidence contexts, and managing model-specific caches are all elements of this broader paradigm. Such model-awareness extends to operational metrics as well: observability systems must track semantic drift, embedding distribution changes, retrieval quality, and inference cost patterns to ensure that the architecture continues to meet quality-of-service expectations over time. As LLMs are increasingly deployed in mission-critical applications across healthcare, finance, law, and scientific research, trust and verifiability become paramount. The architecture's emphasis on provenance, audit trails, and explainable retrieval serves not only operational needs but also

ethical imperatives. Systems built on LLM-native principles can provide evidence-backed generation and accountability, mitigating risks of hallucination and enhancing confidence in automated outputs.

The envisioned architectural blueprint is not merely a conceptual exercise, but a practical foundation for real-world systems that must scale, adapt, and comply with organizational and regulatory requirements. While this paper has established core design principles and architectural components, numerous avenues for future research remain. These include aligning embedding spaces across evolving model versions, developing efficient private retrieval techniques, standardizing representation metadata formats, and constructing comprehensive benchmarks for semantic retrieval and generation quality. In conclusion, foundation models are not merely another class of applications; they represent a fundamental shift in how data is accessed, interpreted, and generated. Meeting the needs of these systems requires rethinking data architecture from the ground up—placing learned representations, semantic access patterns, model constraints, and governance at the core of design. LLM-native data architectures embody this shift, offering a cohesive framework for building intelligent, scalable, trustworthy, and future-proof data systems in the era of foundation models.

VIII. REFERENCES

- [1] Vaswani, A., Shazeer, N., Parmar, N., et al. (2017). *Attention Is All You Need*. NeurIPS.
- [2] Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2019). *BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding*. NAACL.
- [3] Brown, T., Mann, B., Ryder, N., et al. (2020). *Language Models Are Few-Shot Learners*. NeurIPS.
- [4] Lewis, M., Perez, E., Piktus, A., et al. (2020). *Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks*. NeurIPS.
- [5] Khandelwal, U., Levy, O., Jurafsky, D., & Zettlemoyer, L. (2020). *Generalization Through Memorization: Nearest Neighbor Language Models*. ICLR.
- [6] Johnson, J., Douze, M., & Jégou, H. (2019). *Billion-Scale Similarity Search with GPUs*. arXiv.
- [7] Malkov, Y. A., & Yashunin, D. A. (2018). *Efficient and Robust Approximate Nearest Neighbor Search Using Hierarchical Navigable Small World Graphs*. IEEE TPAMI.
- [8] Reimers, N., & Gurevych, I. (2019). *Sentence-BERT: Sentence Embeddings Using Siamese BERT-Networks*. EMNLP.
- [9] Karpukhin, V., Oguz, B., Min, S., et al. (2020). *Dense Passage Retrieval for Open-Domain Question Answering*. EMNLP.
- [10] Guu, K., Lee, K., Tung, Z., Pasupat, P., & Chang, M. (2020). *REALM: Retrieval-Augmented Language Model Pre-Training*. ICML.
- [11] Radford, A., Wu, J., Child, R., et al. (2019). *Language Models Are Unsupervised Multitask Learners*. OpenAI Blog.
- [12] Lee, K., Piktus, A., Lewis, M., et al. (2021). *Beyond Goldfish Memory: Long-Context Retrieval for LLMs*. arXiv.
- [13] Shokri, R., & Shmatikov, V. (2015). *Privacy-Preserving Deep Learning*. CCS.
- [14] Dwork, C., & Roth, A. (2014). *The Algorithmic Foundations of Differential Privacy*. Foundations and Trends in Theoretical Computer Science.
- [15] Kairouz, P., Bonawitz, K., et al. (2019). *Advances and Open Problems in Federated Learning*. arXiv.
- [16] Patterson, D., Gonzalez, J., Le, Q., et al. (2021). *Carbon Emissions and Large Neural Network Training*. arXiv.
- [17] Chen, M., Tworek, J., Jun, H., et al. (2021). *Evaluating Large Language Models Trained on Code*. arXiv.
- [18] Ramesh, A., Dhariwal, P., Nichol, A., et al. (2021). *Zero-Shot Text-to-Image Generation (DALL·E)*. arXiv.
- [19] Croft, W. B., Metzler, D., & Strohman, T. (2009). *Search Engines: Information Retrieval in Practice*. Addison-Wesley.
- [20] Manning, C. D., Raghavan, P., & Schütze, H. (2008). *Introduction to Information Retrieval*. Cambridge University Press.