



This research note is restricted to the personal use of Saurabh Gupta (saurabh.gupta@ge.com).

Solution Path for Implementing a Comprehensive Architecture for Data and Analytics Strategies

Published 28 June 2018 - ID G00351281 - 77 min read

By Analysts [Carlton E. Sapp](#),

Supporting Key Initiative is [Data Management Strategies](#)

New analytic requirements need a flexible architecture for interoperating with diverse data sources and complementary analytic solutions. Technical professionals can use this step-by-step methodology to create a comprehensive architecture to support a data and analytics strategy.

Overview

Key Findings

- New analytics use cases are driving the need for greater accessibility in architectural capabilities and technologies that reduce the movement of data.
- Rising demands for analytic services are driving the need for cloud-based analytic services. Such services help support distributed data stores, near-real-time data pipelines, and AI technologies, such as machine learning, language processing and computer vision.
- Data lakes are taking new forms, such as cloud data processing systems and object stores that bridge the gap between raw data ingestions landing and a broader scope of end users.
- Architectures that support deploying diverse analytics into production drive user adoption. Automated deployment architectures that provision analytics into operations or embed them into business processes are emerging as a top priority for technical professionals.

Recommendations

Technical professionals focused on modernizing their data and analytics architecture and strategy:

- Move analytics to the data by leveraging architectural patterns that expose analytic functions natively on the data source or as a part of database functions.

- Develop a transition plan for analytics leveraging multivendor, hybrid cloud-based services. Cloud-enabled data and analytics achieve superior scalability and offer the high-performance computing needed by machine learning and other AI technologies.
- Repurpose and restructure the data lake to help support macro and micro architectural patterns that empower data scientists, citizen data scientists and remaining knowledge workers.
- Build in repeatable, analytic deployment architectures that support complementary analytic solutions by leveraging integrated analytic engines and model exchange formats.

Problem Statement

How can I build a holistic data management architecture to support current and future business intelligence, advanced analytics and machine learning?

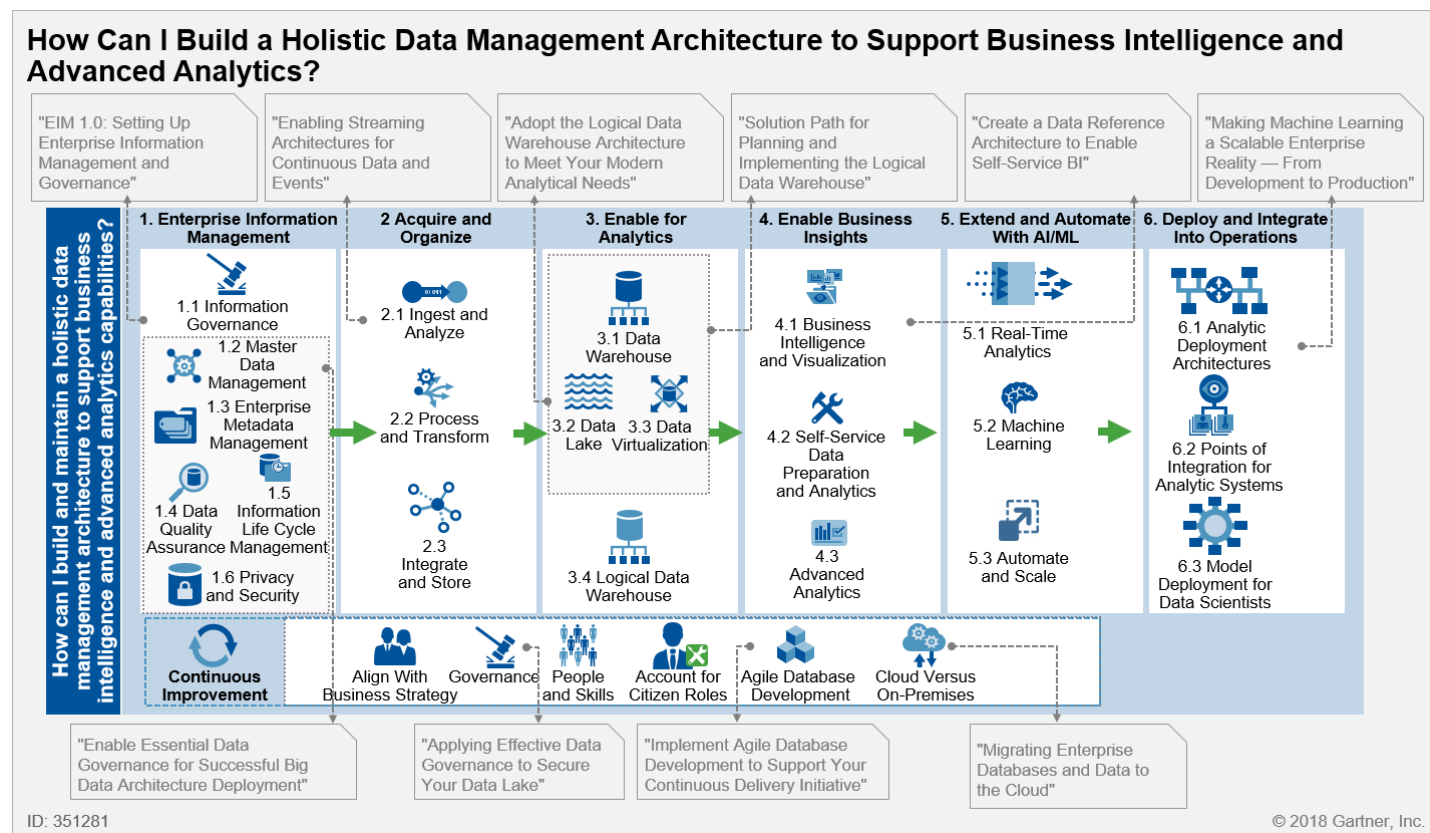
Solution Path Diagram

Many organizations struggle with how to implement a data management and analytics strategy. They ask Gartner about holistic data management architectures. Additionally, trends in big data drive technical professionals to constantly evolve their existing strategy.

Gartner recommends that organizations adopt a build-for-change mindset to establish a highly modular data architecture that delivers and interoperates with complementary analytic solutions. Each new capability should be implemented using an iterative execution model based on a framework that ensures consistency across key technical considerations.

This Solution Path outlines the steps in the framework (see Figure 1) and provides links to related Gartner research.

Figure 1. How Can I Build a Holistic Data Management Architecture to Support Business Intelligence and Advanced Analytics?



Source: Gartner (June 2018)

Solution Path

Digital business transformation demands a modern data and analytics architecture, but organizations are still struggling to maintain existing processes and coping with outdated technologies and skills. The old methods are no longer working.

Organizations need an iterative framework (as shown in Figure 1) that builds on the knowledge acquired at each stage of the process and integrates new and emerging technologies to address changing business demands. Moreover, organizations must constantly evolve their data and analytics architectures to support automated decision making and action based on an ever-changing flow of new information.

Using a framework also allows businesses to more quickly add new functionality and scale while ensuring that they retain governance and control.

At each stage of the framework, technical professionals must revisit common elements within the framework that ensure the consistent application of key decisions. These elements are shown across the bottom of the framework in Figure 1 and represent areas for continuous improvement.

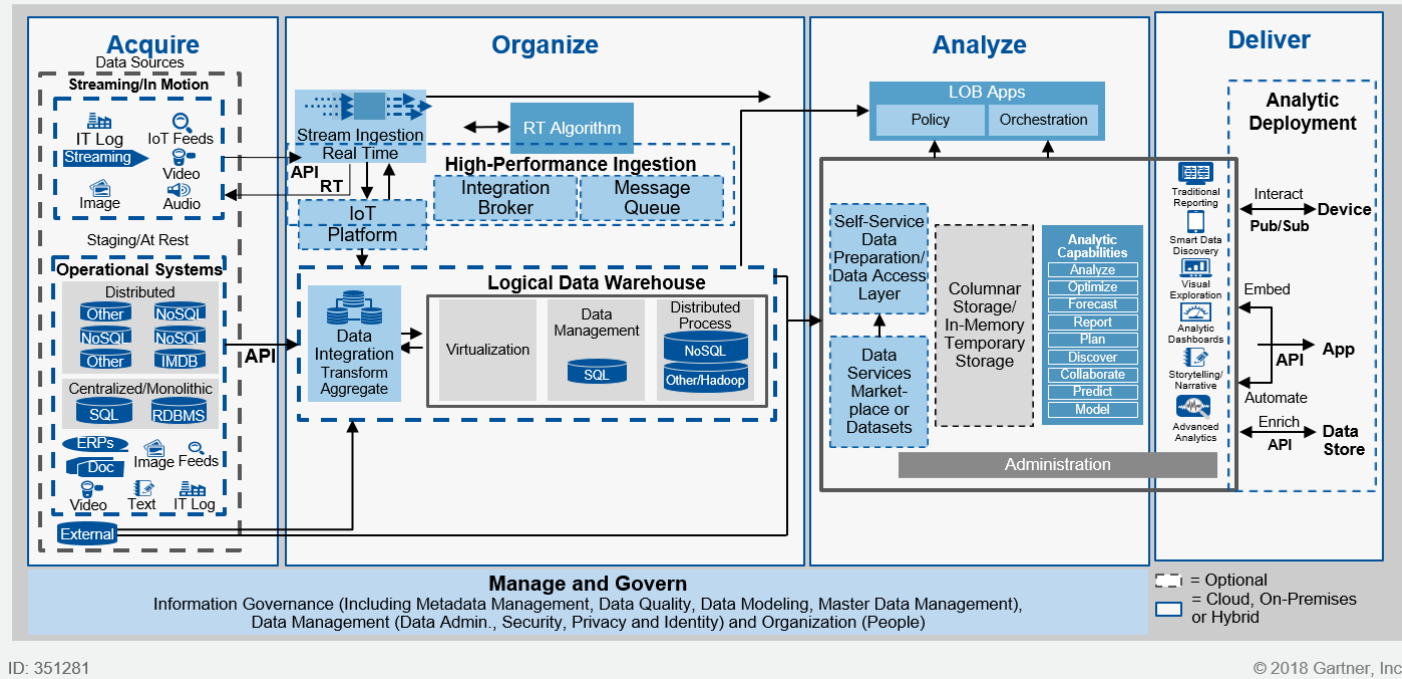
This Solution Path for data and analytics will help technical professionals plan their strategy by addressing key questions, such as:

- How can we ensure that information is properly classified and protected according to business requirements without sacrificing agility?
- How do we manage the acquisition and integration of data from disparate sources, including cloud and on-premises?
- What patterns and technologies are needed to organize and store data in support of operational and analytical projects?
- How do we extract value from data to support ongoing business decisions and create competitive advantage?
- How can we adopt an agile approach to analytics, delivering insights earlier and more frequently while learning from each initiative and evolving our approach?
- How can advances in machine learning and artificial intelligence help us automate actions based on these insights?
- How can we deploy/implement analytics in operations, or embed them into products or services?

The program described by this Solution Path helps data management professionals structure the business and technical projects that will lead to an efficient and agile data management strategy with advanced architectural capabilities supporting data through creation, capture, distribution and consumption. The comprehensive, end-to-end data and analytics architecture (as shown in Figure 2 reference architecture) supports the full spectrum of projects outlined in this research. Figure 2 provides a logical view of an end-to-end architecture meant to guide the technical professional toward building a comprehensive data architecture.

Figure 2. An End-to-End Data Reference Architecture

An End-to-End Data Reference Architecture



Source: Gartner (June 2018)

Step 1: Enterprise Information Management

The first planning step in this Solution Path, enterprise information management (EIM), builds on a business-focused foundation. Gartner defines enterprise information management as an integrative discipline for structuring, describing and governing information assets, regardless of organizational and technological boundaries, and to improve operational efficiency, promote transparency and enable business insight. EIM creates alignment with business at every stage of the architecture by promoting enterprisewide information governance, extending the business data dictionary, and encouraging data quality and context.

Review your information governance, metadata and business data dictionary. These underpin the architecture and should be a first step in developing a modern data and analytics strategy.

The EIM program is typically owned and championed by the CIO, or in the CIO's place, the chief data officer (CDO). It is the program that helps coordinate and organize all information initiatives to achieve alignment with business objectives. Although EIM is owned and led by the CDO, the technology components supporting the EIM program are normally managed by IT with support from and collaboration with business stakeholders.

EIM synchronizes decisions between strategic, operational and technical stakeholders, coordinating efforts to improve the organization's information capabilities.

A data management and analytics strategy benefits from an EIM program with technology, processes and tools supporting the key areas surrounding information governance:

- Metadata management
- Master data management
- Data quality management
- Information life cycle management
- Privacy and security

Each of these areas promotes the goals of information governance, as shown in Figure 3.

Figure 3. EIM Components and Their Relationship to Information Governance

EIM Components and Their Relationship to Information Governance

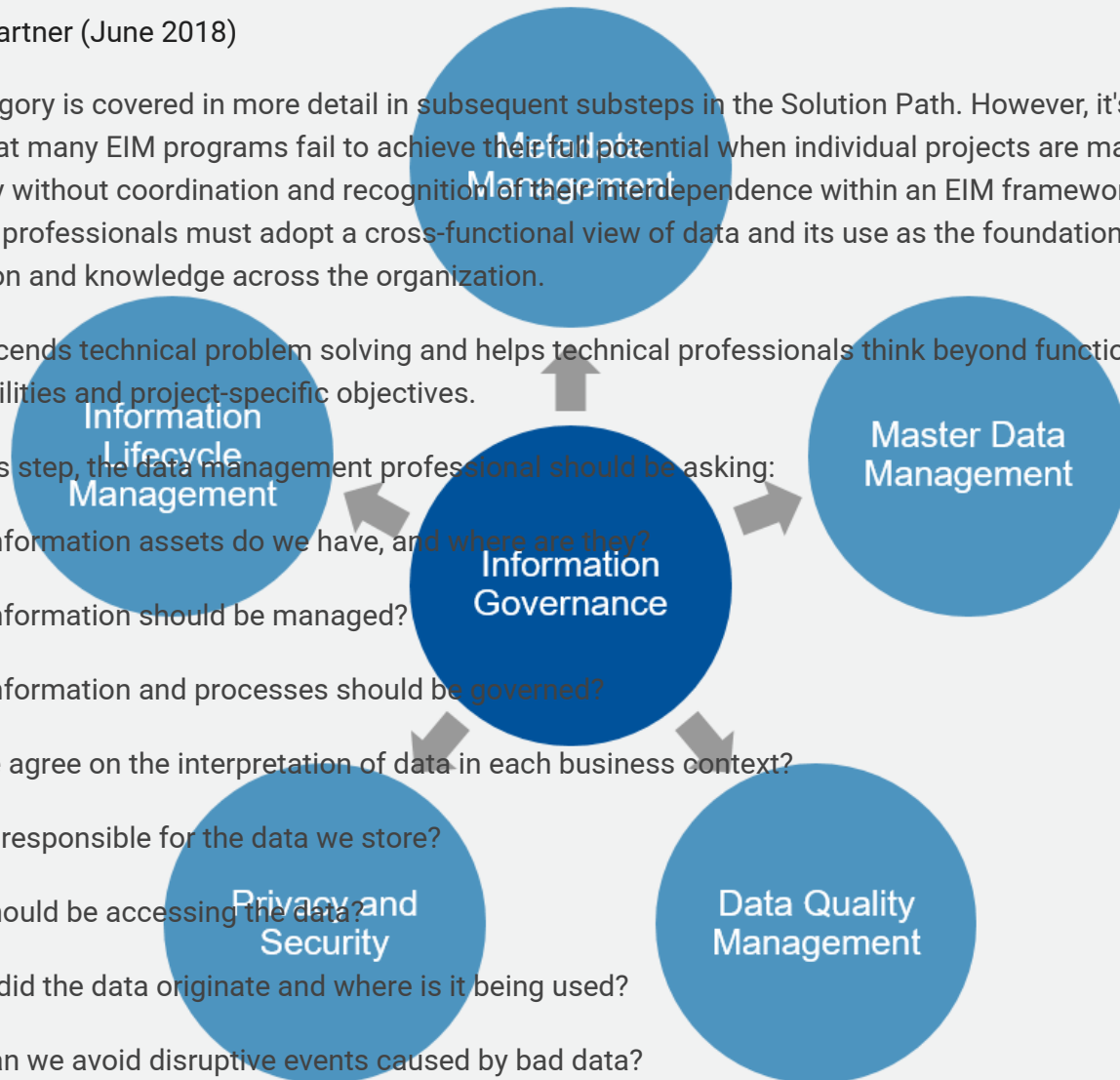
Source: Gartner (June 2018)

Each category is covered in more detail in subsequent substeps in the Solution Path. However, it's important to note that many EIM programs fail to achieve their full potential when individual projects are managed separately without coordination and recognition of their interdependence within an EIM framework. Technical professionals must adopt a cross-functional view of data and its use as the foundation for information and knowledge across the organization.

EIM transcends technical problem solving and helps technical professionals think beyond functional job responsibilities and project-specific objectives.

During this step, the data management professional should be asking:

- What information assets do we have, and where are they?
- What information should be managed?
- What information and processes should be governed?
- Can we agree on the interpretation of data in each business context?
- Who is responsible for the data we store?
- Who should be accessing the data?
- Where did the data originate and where is it being used?
- How can we avoid disruptive events caused by bad data?



ID: 351281

© 2018 Gartner, Inc.

Gartner has published a decision framework to help organizations through this process (see the Gartner research "[EIM 1.0: Setting Up Enterprise Information Management and Governance](https://www.gartner.com/document/code/342309?ref=grbody&refval=3880568)" (<https://www.gartner.com/document/code/342309?ref=grbody&refval=3880568>)). Gartner's method for designing and implementing EIM and data governance helps information management professionals to create EIM programs that are easy to understand and adaptable to change.

1.1 Information Governance

Information governance is a decision-making framework for assigning rights, responsibilities and authorities to ensure that an enterprise, its regulators and its shareholders receive reliable, authentic, accurate and timely information.

Although information governance is presented as an early step in the Solution Path, information governance is not a project, but an ongoing program integral to each new data management initiative. Organizations should consider this step as both a foundation and an ongoing program with regular checkpoints and considerations for new information sources (see the Iterate for Continuous Improvement section below).

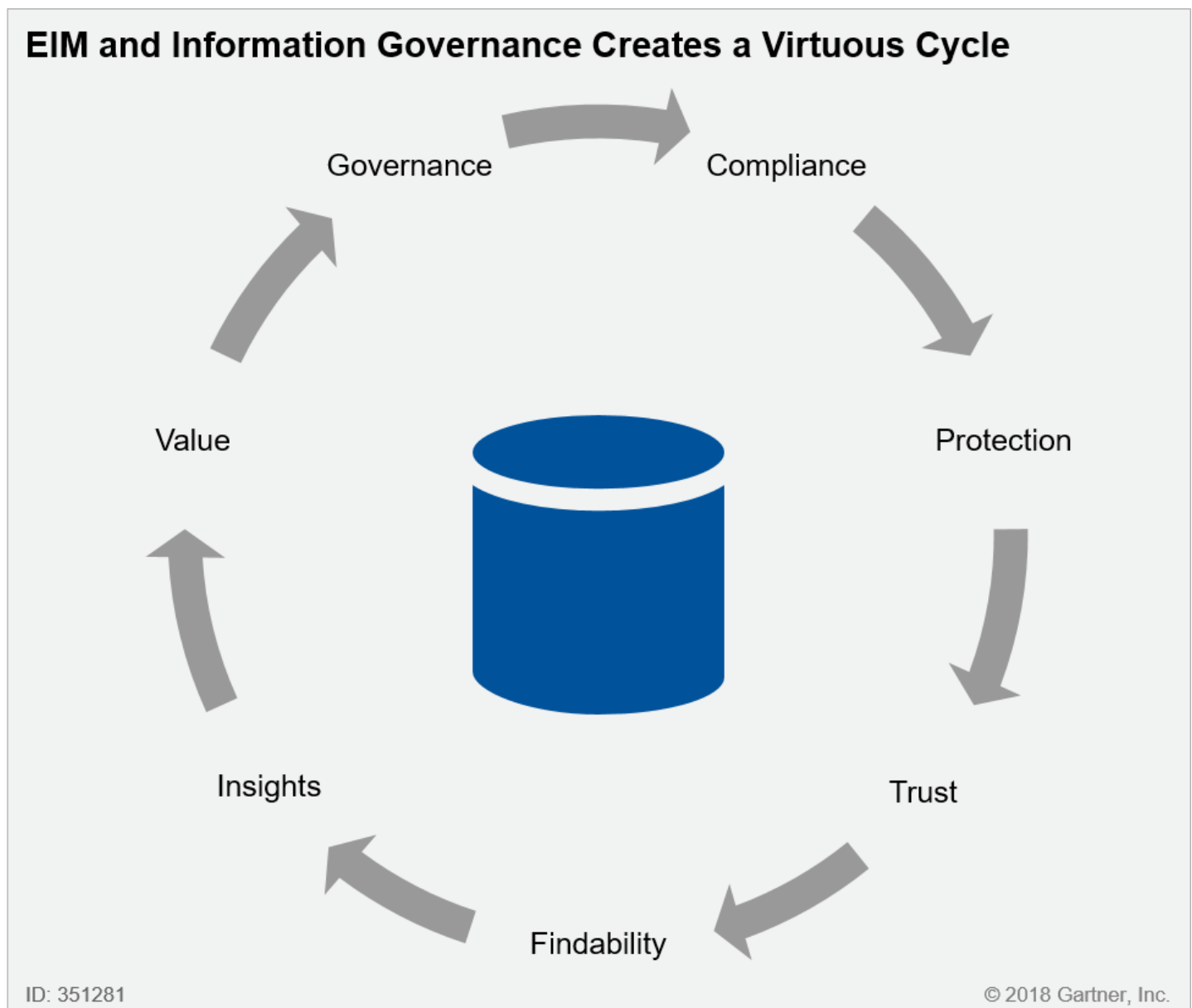
The foundational information governance program will establish processes and controls that will ensure that information is accurate and nonredundant. This requires the introduction of formal management controls in the form of systems, processes and data stewards who serve as custodians of the data.

Each new information management initiative must operate within the constraints of the information governance program to create ongoing consistency and quality. This means that each new data source must undergo data cleansing to prevent ingestion of corrupted, inaccurate, duplicated or nonessential data (see Step 1.4, Data Quality Assurance).

Information management and governance is more than just protecting sensitive data by restricting access and enforcing compliance (see Step 1.6, Privacy and Security). Successful information governance protects information while simultaneously promoting awareness and understanding of the business's data assets to help drive desired business outcomes.

This interrelationship between EIM components creates a virtuous cycle, as shown in Figure 4.

Figure 4. EIM and Information Governance Creates a Virtuous Cycle



Source: Gartner (June 2018)

1.2 Master Data Management

The proliferation of enterprise applications that store information about customers, products and other information assets has made it difficult to create, maintain and enable a single, trusted, shareable version of master data across business domains. Organizational changes, such as mergers and acquisitions, introduce additional challenges because new information assets must be reconciled to support business processes and decision making.

The contention between the mastering of shared data objects leads to inconsistencies in data quality and classification, requiring reconciliation through complex and error-prone data transformation processes. Data inconsistencies result in low confidence in data and inhibit business decisions that rely upon that data. Without enterprisewide agreement on commonly reused master data domains, entities and attributes, organizations cannot be totally effective or efficient in the execution of many business and IT programs. To address these inconsistencies, Gartner recommends master data management (MDM) to help data practitioners create a single version of truth for specific information assets.

MDM is a technology-enabled business discipline in which business and IT work together to ensure the uniformity, accuracy, stewardship, governance, semantic consistency and accountability of an enterprise's official shared master data assets. Master data is the consistent and uniform set of identifiers and extended attributes that describe the core entities of the enterprise.

Each organization should assess the cause of data inconsistencies to decide where it will focus its MDM efforts. There are six characteristics that will drive the MDM implementation's focus:

- Where the master data is authored
- Where the master data is verified
- The latency of master data movement
- The degree to which a physical "golden record" is instantiated
- The usage of the master data
- Search complexity

Gartner has published a guide to four MDM implementation styles to help organizations choose a strategy aligned with their unique requirements (see "[A Comparison of Master Data Management Implementation Styles](https://www.gartner.com/document/code/342604?ref=grbody&refval=3880568)" (<https://www.gartner.com/document/code/342604?ref=grbody&refval=3880568>)).

1.3 Enterprise Metadata Management

Enterprise data management can be a powerful enabler for multiple EIM functions, like MDM and data governance.

The next step is to establish enterprise metadata management (EMM). Because of the growing variety and volume of data, data lakes are becoming an essential part of the architecture. The data lake pattern employs technology that supports data from a variety of different sources, such as files, clickstreams, Internet of Things (IoT) sensors and social networks. The unconstrained storage model of the data lake offers extreme flexibility but lacks the assurances and context enforced within an enterprise data warehouse. EMM offers an effective solution to this problem.

An effective metadata management approach enables the following information capabilities:

- **Describe:** To collect knowledge about information assets
- **Organize:** To align and structure information assets so they can be readily found and easily consumed by other capabilities of the platform
- **Share:** To make data available to consumption points
- **Govern:** To provide for control, levels of consistency, protection, quality assurance, risk assessment and compliance

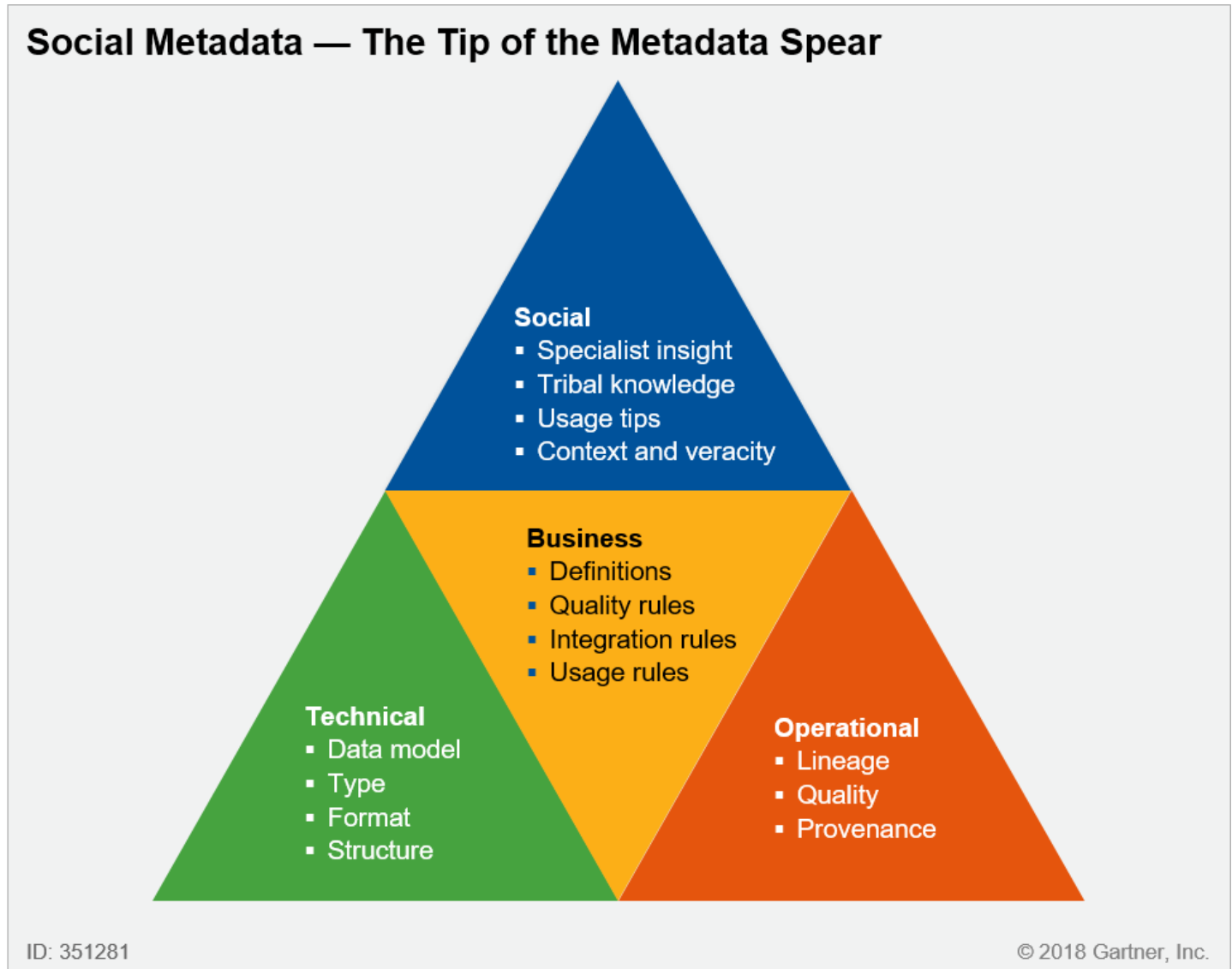
To understand the benefits of EMM, it's useful to review four types of metadata:

- **Technical metadata:** Describes the technical attributes including the form, type and structure of each dataset.
- **Operational metadata:** Captures the lineage and provenance of the data, along with audit details about the success or failure of job runs and update frequency.
- **Business metadata:** Captures specific business attributes to help promote context, findability and meaning. Business metadata gives end users a way to find and understand information assets based on business semantics.
- **Social metadata:** Gives end users the ability to tag and describe data with information meaningful within their specific business context. This crowdsourced metadata can offer tremendous collaborative benefits to citizen data scientists working with new and exotic data sources.

Social metadata promotes discovery, but more importantly, provides meaningful insights from specialist and citizen data scientists that evolve new insights about the data on each new project.

Social collaboration plays a critical new role in metadata strategy by providing context from users that work closely with the data (see Figure 5).

Figure 5. Social Metadata — The Tip of the Metadata Spear



Organizations can publish a centrally managed EMM data catalog/repository and business glossary to provide a consistent set of definitions and descriptions for each data element. Such a strategy adds coherence to business metadata and ensures reliability in data semantics.

By combining different types of metadata, the data catalog/repository gives insight into:

- The movement of data within jobs or between systems
- What data is available within the organization and how it can be used
- Impact analysis to determine the dependencies between information assets
- Common business vocabulary and accountability for its terms and definitions
- Audit trail for compliance

EMM technologies offer automated metadata discovery to inspect data elements for semantics and take specific actions based on workflow rules: for example, inspecting datasets for operational metadata combined with business metadata to determine the information classification of incoming data. This could be used to ensure policy/regulatory compliance or for the dynamic application of data masking or obfuscation techniques.

As previously mentioned, EMM is an important capability for information governance. To ensure consistency with information governance, metadata management processes must support the following:

- Reliable, sustainable, governed processes based on defined standards.
- Storage in a common model that enforces standards and is managed in an integrated repository.
- Data consumers need to be able to access the data from one central place. They must be able to provide feedback about the metadata to improve it over time.

Gartner has published a guide to deploying effective metadata management solutions to help organizations choose a strategy aligned with their unique requirements (see "[Deploying Effective Metadata Management Solutions](https://www.gartner.com/document/code/347645?ref=grbody&refval=3880568)" (<https://www.gartner.com/document/code/347645?ref=grbody&refval=3880568>)).

1.4 Data Quality Assurance

As mentioned earlier, information governance promotes awareness and understanding of the business's data assets to help drive desired business outcomes. In this step, a data quality assurance strategy is adopted or adapted to improve reliability and usability of information by ensuring that data is fit-for-purpose in downstream business processes.

These processes range from those used in core operations to those required by analytics and for decision making, regulatory compliance, and engagement and interaction with external entities. As a discipline, data quality assurance covers much more than technology. It also includes:

- Roles and organizational structures
- Processes for monitoring, measuring, reporting and remediating data quality issues
- Links to broader information governance activities via data-quality-specific policies

Successful data quality assurance depends on a strong partnership between IT and business stakeholders, along with clear roles and responsibilities. For example, organizations often rely on information stewardship for the enforcement of information governance policies and rules. Data steward is a role that is established within a line of business (LOB) — not IT. Data stewards work on behalf of the business stakeholders, enacting the policies created and working to ensure data conforms to expectations. They actively monitor the quality of the data (via data quality metrics and visualization techniques) and take corrective action when data in certified sources does not align with policy. In addition, data stewards influence the people around them to rely on certified information sources.

Technical professionals support the techniques and technologies used by data stewards for discovering and investigating data quality issues, such as duplication, lack of consistency, and lack of accuracy and completeness. This is accomplished by analyzing one or multiple data sources and collecting metadata that shows the condition of the data and enabling the further investigation into the origin of data errors.

Data quality assurance should be applied at all stages of the data and analytics architecture. However, efforts to improve data quality early in the data life cycle will reduce friction and improve efficiency.

Organizations should apply data quality assurance across the enterprise and data life cycle:

- Adopt comprehensive data quality assurance processes to ensure consistency and avoid duplication of efforts.
- If data quality is not recognized as a priority, enlist support from business sponsors to create a business case for data quality investments.
- Codify expectations of data usability by cataloging requirements for data completeness, accessibility, format, presentation, timeliness, relevance and accuracy.
- Use metadata to assess and classify data quality characteristics.
- Implement data profiling during ingestion and integration, but offer self-service data preparation to extend quality assurance capabilities to end users who are closest to the data.

1.5 Information Life Cycle Management

Information life cycle management (ILM) is an approach to data and storage management that recognizes that the value of information changes over time and that it must be managed accordingly. ILM seeks to classify data according to its business value and establish policies to migrate and store data on the appropriate storage tier and, ultimately, remove it altogether. ILM has evolved to include upfront initiatives like master data management and compliance.

Data management professionals should consider ILM classification strategies, such as:

- Categorizing data according to business life cycle rules
- Identifying frequently and infrequently accessed data and classifying as hot, warm, cold
- Identifying data by type for different storage tiers
- Considering data retention, compliance and security requirements mandated by regulations

Business continuity strategies strive to keep businesses operating in the event of a disaster. Data protection is a key component of business continuity.

Particularly with historical or analytical data use cases, organizations can take advantage of tiered storage models to move data from nodes with faster storage (and higher compute capability) to nodes with lower-cost storage models. This strategy allows for cost-effective and scalable storage platforms for data that would otherwise require higher-cost resources. Cloud service providers (CSPs) offer tiered storage to take advantage of different ILM classifications.

Information governance may alter or preclude ILM handling procedures. You should:

- Require encrypted storage for certain classifications
- Prevent storage or replication across regional boundaries
- Require archival or purging of data based on information retention policies

1.6 Privacy and Security

Data privacy and security form a significant part of the information governance program. Privacy and related issues are complex and will require both business and legal support to determine how best to address the risk of personal data processing. Data security provides tools and processes to protect the data from a variety of threats. Both dimensions require the orchestration of policies across disparate data stores and processes acting on data.

Any successful program needs continuous insight as to where the relevant data exists. Without such an understanding, any technology control will be at best partially successful and at worst useless. Data discovery products can help identify sensitive data and trigger appropriate data protection measures to:

- Apply metadata or labeling

- Apply data tokenization, masking and/or redaction
- Automated access rights remediation

See ["Improving Data Security Governance Using Classification Tools"](https://www.gartner.com/document/code/337209?ref=grbody&refval=3880568)

(<https://www.gartner.com/document/code/337209?ref=grbody&refval=3880568>) for a discussion on how and where such technologies can be found and used.

Note that while this may suffice for internal security requirements, it is unlikely to provide full compliance with privacy or similar regulations. Understanding your data "map" allows you to determine where security technologies and processes may be required to protect data. But for compliance, it's necessary to understand how the data is being used, why, by whom, under what legal and contractual basis, and under what limitations.

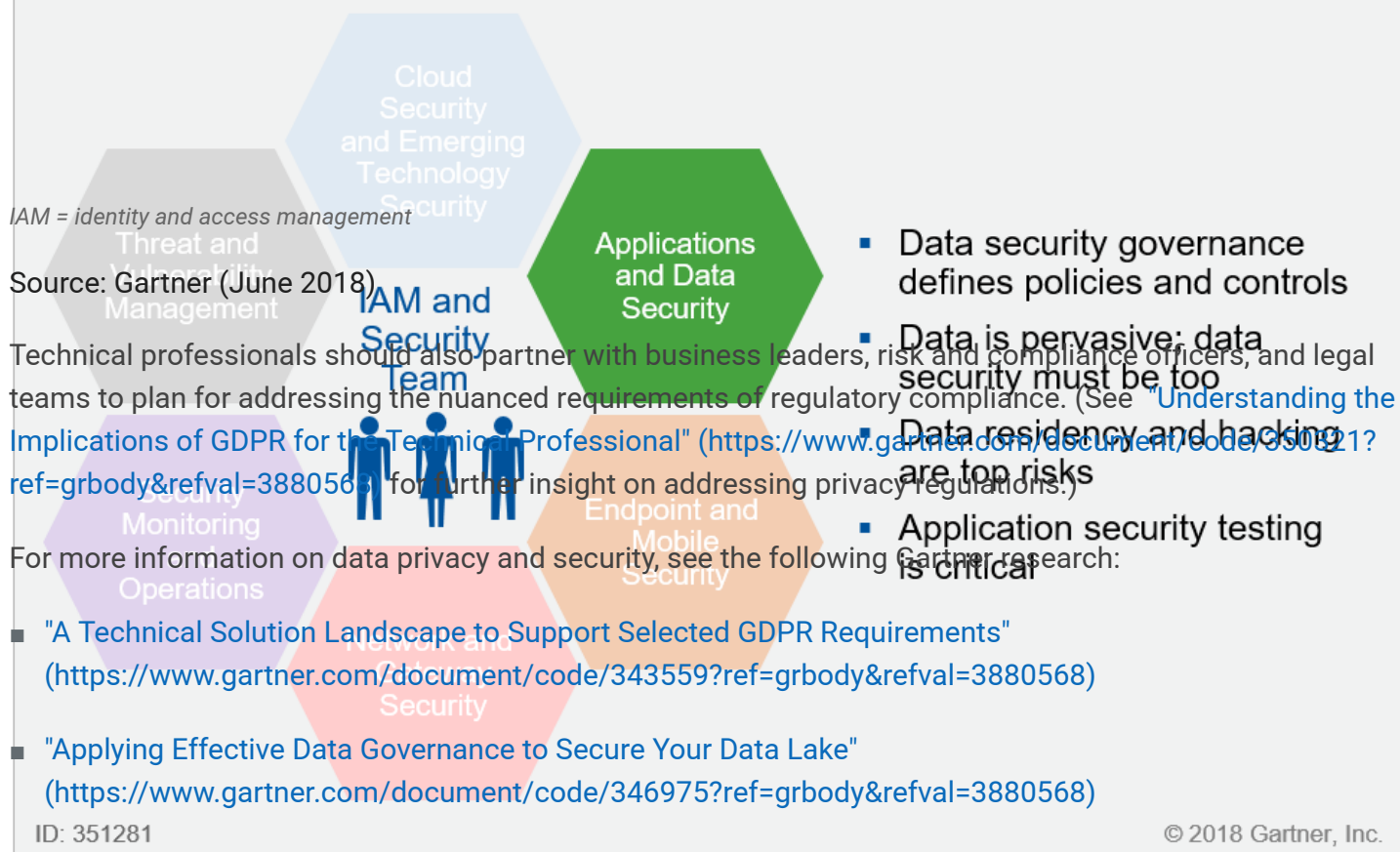
As a starting point, technical professionals should work with security teams to identify where and how to implement technologies to enforce information governance policies, such as:

- Requiring encryption at the application or data layer, including storage and at-rest
- Implementing role-based access controls to protect data elements from unauthorized consumption or alteration
- Segmenting data for different audiences based on information classifications

Data privacy and security are part of a much broader information security agenda (see Figure 6). Data management professionals should work closely with information security professionals to ensure appropriate security measures are applied.

Figure 6. Application and Data Security Overview

Application and Data Security Overview



- "Securing the Big Data and Advanced Analytics Pipeline" (<https://www.gartner.com/document/code/352648?ref=grbody&refval=3880568>)
- "Four Steps to Secure Modern Databases" (<https://www.gartner.com/document/code/290961?ref=grbody&refval=3880568>)
- "Protecting Big Data in Hadoop" (<https://www.gartner.com/document/code/271209?ref=grbody&refval=3880568>)

What to Consider Before Moving to the Next Step

Gartner recommends that most organizations should consider the following before moving on to the next step:

- Consolidate EIM initiatives to ensure functional collaboration, and reduce overlap and inconsistencies arising from fragmented projects.
- Establish a framework for applying data governance to new data management and analytics programs. For more information, see ["Enabling Essential Data Governance for Successful Big Data Architecture Deployment."](https://www.gartner.com/document/code/327532?ref=grbody&refval=3880568) (<https://www.gartner.com/document/code/327532?ref=grbody&refval=3880568>)

Step 2: Acquire and Organize

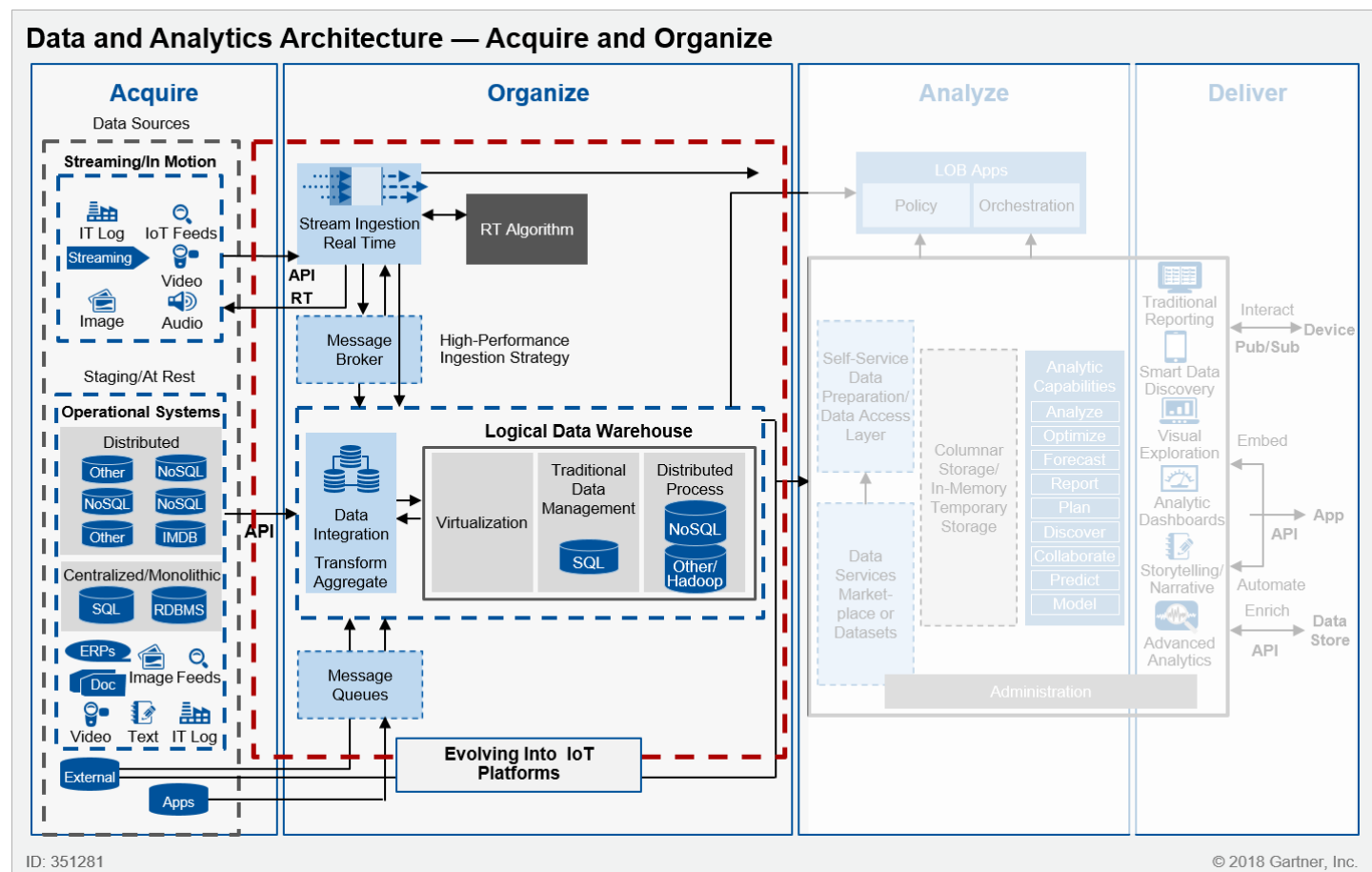
The next step of the Solution Path, Acquire and Organize, prepares the architecture to collect and assimilate information, regardless of frequency, structure or origin. An effective data and analytics strategy allows for discovery, ingestion and integration of all available data from the source, as fast as it is produced, in any format and quality.

This stage involves three related subtasks:

- **Ingest:** Ingestion describes the process of acquiring, importing, transferring or loading data for transformation, processing or storage.
- **Process and transform:** This is the stage where data is prepared, cleaned, modified or enriched before use or storage.
- **Integrate:** Integration describes the process by which data is retrieved from disparate sources and combined to provide a unified view of the data.

The architectural components that support ingestion, transformation and integration are closely related and are commonly used in concert throughout this stage (see Figure 7).

Figure 7. Data and Analytics Architecture — Acquire and Organize



Source: Gartner (June 2018)

Data practitioners should conceptualize this part of the architecture as a managed data pipeline composed of a series of stages through which the data flows.

The following patterns support different use cases for ingesting, transforming and integrating data:

- **Data replication and synchronization** is the process of copying data from one data store to another. This is a well-established integration pattern and is used in traditional extraction, transformation and loading (ETL) processes to manage the acquisition and integration of information from heterogeneous or compatible systems.
- **Batch processing** is an asynchronous method for ingesting and storing an immutable and constantly growing dataset in bounded intervals – that is, the incoming data has a beginning and an end. This method is extremely effective for processing high volumes of data for specified intervals and allows for the computation of arbitrary functions on the dataset, which in turn, are output in batches.
- **Stream processing** is an asynchronous method of ingesting an infinite, constantly growing data stream. Such streams often include unbounded, unordered datasets produced at high velocity. Examples include web logs, mobile events and sensor data, but can represent any unbounded dataset. Discrete portions of the stream, known as windows, can be captured through some characteristic of the data, such as event time stamps.
- **Data virtualization** is a form of data integration based on data abstraction and provides a consistent interface to data distributed across multiple, disparate data sources and repositories. Modern data virtualization tools provide both read and write access to a host of popular data types and sources (such as relational, Hadoop, NoSQL, flat files and cloud data stores).
- **Message-oriented** patterns are often associated with service-oriented architecture (SOA), messaging captures, transforms and delivers data through message brokers. Data is often delivered in near real time.

These patterns are not mutually exclusive and are often used together to provide a broad range of options within the architecture. When implementing each subtask, revisit these patterns to select the most appropriate style for each use case.

2.1 Ingest and Analyze

The ingest step ensures that the data and analytics architecture is capable of collecting, replicating and storing data from many different types of data sources with support for high volumes and low latency.

The acquisition of data requires different strategies depending on the data source, the velocity of information and the type of information to be consumed. Organizations should adopt an architecture capable of supporting multiple styles of data ingestion:

- Batch data ingestion and ingestion of data at rest
- Real-time data stream ingestion and stream processing with low latency

- Reliable and durable message ingestion

Attributes of these styles include:

- Scalability, to support high throughput of incoming data
- Decoupling of data sources from data subscribers
- Support for some transformation before handing over the output stream to consumers
- Ability to act as both input source and output sink to other services

To build a scalable ingestion process with a repeatable data pipeline:

- Embrace event-driven architectural patterns to deliver a stream-based data management strategy.
- Capture and add metadata during ingestion to understand where data originated (lineage) and identify where it was processed or stored.
- Implement ingestion workflow to halt processing, reroute data or take actions based on metadata, quality or other characteristics of the data.
- Implement support for a variety of endpoints and data types with low impact to data sources.
- Avoid custom code, and automate ingestion with managed, automated and repeatable pipeline definitions.
- Implement change data capture to optimize ingestion workloads.
- Regulate the ingestion and read rates of producers and consumers to make streaming applications production-ready.

2.2 Process and Transform

Once data is created or ingested, it usually requires some processing and transformation before it can be used. In this step, the data is prepared, cleaned, modified or enriched so that it will be fit for the business use case.

A traditional approach used an ETL process that involved extracting data from the source, transforming it to fit operational needs and loading it into staging tables or a data warehouse.

Another approach relies on an extraction, loading, and then transformation (ELT) process, where data is extracted from a source, then loaded directly into a staging table or data lake in its raw, unmodified state. Transformation occurs in the staging area before being loaded into a target database or data warehouse.

As with data ingestion, data transformation can benefit from a managed data pipeline:

- Monitor the quality of data as it flows through the pipeline, and implement predefined processing rules to apply transformation, enrichment or structure when needed.
- Apply data privacy and security policies by segmenting, masking or tokenizing data before it gets published for consumption.
- Automate scheduling and orchestration of data movement between heterogeneous storage environments.
- Apply data life cycle policies to move data across tiered storage or cloud and on-premises based on hot, warm or cold classification.

Gartner has published research to help organizations through this process (see ["Enabling Streaming Architectures for Continuous Data and Events With Kafka"](https://www.gartner.com/document/code/353112?ref=grbody&refval=3880568) (<https://www.gartner.com/document/code/353112?ref=grbody&refval=3880568>)). Gartner's method for deploying event-based architectures helps information management professionals create a stream-based data management paradigm to help drive a data ingestion strategy.

2.3 Integrate and Store

Data Integration

Data integration consists of the practices, architectural techniques and tools for collecting data from disparate sources and combining that data into a unified view to meet the data consumption requirements of all applications and business processes.

Unfortunately, the process of integrating data has become more complicated, due to several factors:

- Data is no longer centrally stored and managed in corporate data centers. Integration strategies must account for data sourced from diverse locations, cloud infrastructure and external parties.
- Continuous data streams demand different patterns and technologies with strategies for real-time integration and analytics.
- Today's integration techniques can't rely on predictable data structure and schema, and must be capable of dealing with unstructured data in a variety of formats.
- Business users are demanding self-service integration capabilities (see the Self-Service Data Preparation and Analytics section below).

These challenges necessitate an agile and flexible data integration strategy that can support different styles of data integration, each at the appropriate stage in the architecture. However, take care to avoid data integration strategies that rely on hard-coded, unmanaged interfaces. Without proper governance, point-to-point integrations can become unmanageable and result in redundant work and higher maintenance costs.

Data management professionals should be familiar with the following integration styles:

- **Embedded integration:** Typically delivered through commercial applications, embedded integration is a component of broader IT solutions. It includes three subtypes: embedded in databases, embedded in operational software and embedded in analytics software.
- **Stand-alone integration:** Independent of databases and applications, stand-alone integration is often delivered through separate integration middleware. Stand-alone integration is business-friendly and includes three subtypes:
 - Integration platform as a service (iPaaS)
 - Integration software as a service (iSaaS)
 - Data federation/virtualization (which is mentioned in more detail later in the next section)
- **Data preparation:** Data preparation provides business-friendly and data-centric capabilities, such as data access, data discovery, data cleansing, data transformation, data enrichment and data collaboration. Data preparation functionality can be offered either as stand-alone tools or as embedded capabilities in an analytics platform. More detail on data preparation can be found later in this document.

For more information about the pros and cons of data integration styles, see ["Use Data Integration Patterns to Build Optimal Architecture."](https://www.gartner.com/document/code/270543?ref=grbody&refval=3880568) (<https://www.gartner.com/document/code/270543?ref=grbody&refval=3880568>)

For additional information, please refer to the following research:

- ["Comparing Three Self-Service Integration Architectures"](https://www.gartner.com/document/code/297311?ref=grbody&refval=3880568) (<https://www.gartner.com/document/code/297311?ref=grbody&refval=3880568>)
- ["Comparing Four iPaaS-Based Architectures for Data and App Integration in Public Cloud"](https://www.gartner.com/document/code/289852?ref=grbody&refval=3880568) (<https://www.gartner.com/document/code/289852?ref=grbody&refval=3880568>)
- ["Deploying Effective iPaaS Solutions for Data Integration"](https://www.gartner.com/document/code/324279?ref=grbody&refval=3880568) (<https://www.gartner.com/document/code/324279?ref=grbody&refval=3880568>)

Data Stores

Data stores range from the simple, such as text files, to the complex, such as distributed processing and storage frameworks like Hadoop. There are implementation choices, such as persistent disk storage versus in-memory processing, and infrastructure considerations, such as cloud versus on-premises. Technical professionals have a multitude of options when it comes to data stores and database management systems (DBMSs), and often ask Gartner, "What data store is most appropriate to address my particular use case?"

Answering this question depends on several factors and characteristics of the data and the consumption model, including:

- Transactional guarantees
- Consistency
- Availability
- Schema flexibility
- Scalability and concurrency
- Cost performance
- Workload manageability
- Recoverability and partition tolerance

There are also architectural and business considerations, such as:

- General versus use-specific
- Community support versus vendor support
- Technical maturity
- On-premises versus cloud
- Managed versus unmanaged
- Single vendor versus multiple vendors

Choosing the components that make up the data storage and management architecture is a critical step and will have long-lasting implications for the maintainability, performance, availability and efficacy of the overall architecture.

Data architects will need to evaluate the requirements and constraints for each type of DBMS to decide which option is most appropriate for a particular use case. To assist with this process, refer to the following research:

- "Decision Point for Selecting a DBMS Architecture" (<https://www.gartner.com/document/code/274013?ref=grbody&refval=3880568>)
- "Identifying and Selecting the Optimal Persistent Data Store for Big Data Initiatives" (<https://www.gartner.com/document/code/322578?ref=grbody&refval=3880568>)

Today's cloud-based data services represent a strong alternative to on-premises deployments. Cloud-based operational and analytical database services offer deployment speed, flexibility and scalability that are difficult, or even impossible, for enterprises to achieve with on-premises systems.

- For a guide to evaluating cloud-based database services, see "[Evaluating Microsoft Azure's Cloud Database Services](https://www.gartner.com/document/code/311029?ref=grbody&refval=3880568)" (<https://www.gartner.com/document/code/311029?ref=grbody&refval=3880568>) and "[Evaluating the Cloud Databases From Amazon Web Services](https://www.gartner.com/document/code/303836?ref=grbody&refval=3880568)." (<https://www.gartner.com/document/code/303836?ref=grbody&refval=3880568>)
- For a comparison of cloud-based managed big data services, see "[Assessing Cloud-Based Big Data Services: Amazon EMR vs. Microsoft Azure HDInsight](https://www.gartner.com/document/code/317100?ref=grbody&refval=3880568)." (<https://www.gartner.com/document/code/317100?ref=grbody&refval=3880568>)

What to Consider Before Moving to the Next Step

Gartner recommends that most organizations consider the following before moving on to the next step:

- Review known and anticipated requirements necessitating new architectural capabilities, such as streaming ingestion, real-time data streams and integration between heterogeneous systems.
- Evaluate cloud services and hybrid cloud strategies to support data migration requirements.

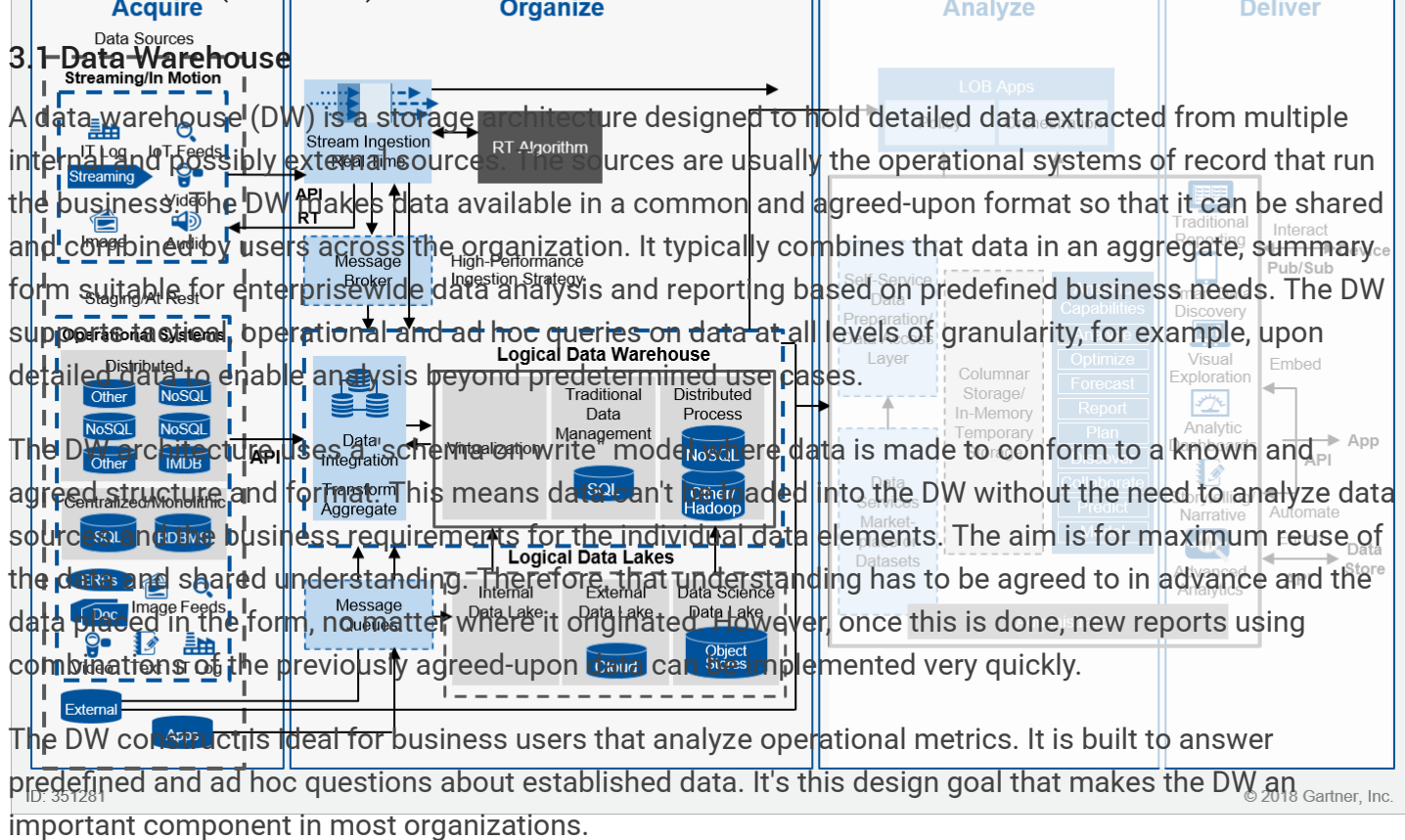
Step 3: Enable Data for Analytics

In this step, the data is prepared for consumption by end users, applications and business processes in support of analytics. The data could be stored in memory, on disk, and/or distributed across physical or logical data stores or data lakes (see Figure 8).

Figure 8. Data and Analytics Architecture — Store and Enable

Data and Analytics Architecture — Store and Enable

Source: Gartner (June 2018)



However, the DW model has certain shortcomings in modern digital business. Due to the need to structure and predefine data in advance, the DW finds it difficult to accommodate nontraditional data sources, such as web server logs, sensor data and social network activity.

Loading new data sources incur planning and expense. If it's not clear how to answer business questions, it will usually be excluded from the data warehouse. This makes analysis that requires access to these data sources very difficult, unless the DW is used in conjunction with a data lake, or as part of a logical data warehouse (see Section 3.4 below).

Note that some of the traditional issues with the data warehouse are being mitigated. For example, automatic discovery and profiling of data makes the consumption of data much simpler. Likewise, many relational database management systems (RDBMSs) have been extended to allow for semistructured data, such as JavaScript Object Notation (JSON). This, together with the ability to place this data on cost-effective storage, makes mixing structured and semistructured data much simpler. However, that said, the traditional data warehouse is mainly aimed at that part of the workload using structured and curated data and now works in conjunction with the other major data stores.

Implement a data warehouse when you need:

- A system that can integrate, aggregate and analyze data from transactional and operational systems
- A prestructured data model designed to support enterprise reporting and a common understanding to enable maximum reusability in the enterprise

- An efficient analytical system optimized for reporting and analysis of common and agreed-upon business subjects with high throughput, high concurrency and low latency

Gartner recommends that every data warehouse, whether new or old, be built on or be evolving toward a logical data warehouse (LDW) model to effectively support business intelligence and analytics. The LDW is covered in more detail in Section 3.4.

For more information on data warehouse strategies and platforms, refer to the following research:

- ["Comparing Cloud Data Warehouses: Amazon Redshift and Microsoft Azure SQL Data Warehouse"](https://www.gartner.com/document/code/309107?ref=grbody&refval=3880568) (<https://www.gartner.com/document/code/309107?ref=grbody&refval=3880568>)
- ["Solution Path for Planning and Implementing the Logical Data Warehouse"](https://www.gartner.com/document/code/320563?ref=grbody&refval=3880568) (<https://www.gartner.com/document/code/320563?ref=grbody&refval=3880568>)

3.2 Data Lake

A data lake represents a collection of storage instances of various data assets, which are often stored in a near-exact, or even exact, copy of the source format. Unlike the data warehouse, it is up to the users of the lake to interpret the data and to determine the best data applicability for the identified use cases.

Data lakes don't replace primary systems of record or data warehouses. Instead, they complement existing efforts and support the discovery of new questions. A data lake contains unrefined data where the data structure may not be known in advance, or when organizations want to increase analytics and operational agility by complementing their systems of record with systems of insight. Unlike the data warehouse's "schema on write" model, data lakes use a "schema on read" model.

Use a data lake for:

- Consolidating data in its raw, unrefined state from a variety of different data sources – structured, semistructured and unstructured
- A general-purpose staging area
- Collecting all data and attributes with the goal of deriving new insights from analyzing larger volumes of data or new types of data
- An active archive of historical data

The data lake is a useful construct for managing data in different zones or layers, with each layer optimized for different styles of consumption. These patterns are optional and can easily coexist within a lake, each layer serving its own use cases:

1. **Transient layer:** This acts as a landing area supporting the ingestion of many disparate data sources. This layer holds all raw, unrestricted data, including sensitive information (such as personally identifiable information [PII] or protected health information [PHI]) in its original, unaltered form. The layer is deemed

transient because it serves merely as a staging area to fill other layers. Access to this zone is highly restricted and generally available to admins only due to presence of sensitive information. The data in this zone should be encrypted at rest.

2. **Discovery layer.** This holds all the data that complies with corporate governance and compliance policies to promote exploration and discovery in a governed manner. For example, sensitive data is tokenized or masked to prevent unauthorized consumption. This layer can be thought of as an unrefined and unaltered sandbox for exploration and advanced analytics.
3. **Refined layer.** In this layer, enrichments and transformations are processed to create new datasets, which are made available to downstream applications and processes. Data is integrated into a common format, with data validation and cleansing techniques applied.
4. **Trusted layer.** The trusted layer is a step beyond the refined layer, with reference data reconciled to ensure consistency with master data policies. The format of the data complies with the end-user analytical tools. Many organizations have multiple trusted layers each with data prepared for separate business units, such as marketing and finance, because their needs and governance differ.

For a deeper analysis of data lake architectural styles, see ["Use Design Patterns to Increase the Value of Your Data Lake."](https://www.gartner.com/document/code/342255?ref=grbody&refval=3880568) (<https://www.gartner.com/document/code/342255?ref=grbody&refval=3880568>)

Along with the data warehouse, the data lake is an important component of the LDW.

Design Pattern for Data Lakes to Unleash Advanced Analytics and Machine Learning

Governed data lakes continue to emerge as a viable solution to support advanced analytics and data science initiatives. However, technical professionals should systematically tackle the design of the data lake using a hierarchical framework that includes the following to support these initiatives.

- A macro-level architecture that considers differing conceptual data lake architectural styles, such as and inflow data lake, outflow data lake and data science lab lake.
- A medium-level architecture that determines how the macro-level architectural styles can be broken down into manageable distinct units, also known as logical zones.
- A micro-level architectural style that defines and administers the logical zones and controls movement between them.

The data science lab style lake is best for enabling innovation through discovery and exploration for machine learning and advanced analytics. However, leveraging a data lake architecture guidance framework will ensure greater support for future capabilities.

For more information, see ["Use Design Patterns to Increase the Value of Your Data Lake."](https://www.gartner.com/document/code/342255?ref=grbody&refval=3880568) (<https://www.gartner.com/document/code/342255?ref=grbody&refval=3880568>)

Governing Data Lakes

Growth in the demand for advanced analytics and machine learning requires a more aggressive approach to governing data lakes, which are often the source for data science programs. A more automated data governance framework ensures data science projects using data lakes are properly governed and meet compliance requirements. The following elements should be considered to govern a data lake to support modern advanced analytics and machine learning initiatives:

- Autodiscovery of sensitive data and/or ability to tag attributes with sensitive classification
- Early data discovery
- Enabling operational analytics on data usage

Governing data lakes is also required to operationalize data lakes. In the early iterations of data lakes, the only consumers were data scientists. However, organizations are now expanding data lake workloads and adding business and data analysts as its consumers. This additional exposure necessitates that the data lake be governed and meet compliance requirements.

For more information on governing the data lake, see:

- ["Enabling Essential Data Governance for Successful Big Data Architecture Deployment"](https://www.gartner.com/document/code/327532?ref=grbody&refval=3880568)
(<https://www.gartner.com/document/code/327532?ref=grbody&refval=3880568>)
- ["Applying Effective Data Governance to Secure Your Data Lake"](https://www.gartner.com/document/code/346975?ref=grbody&refval=3880568)
(<https://www.gartner.com/document/code/346975?ref=grbody&refval=3880568>)

3.3 Data Virtualization

Data virtualization can provide a uniform interface to multiple data stores, allowing users easy access to all the organization's data (subject to security controls).

When combined with the data warehouse and data lake, data virtualization (DV) becomes an integral component of the LDW (see the next section). DV can provide a uniform interface to multiple data stores, allowing users easy access to all the organization's data (subject to security controls). DV can also be used by business intelligence (BI) and reporting tools to enable analytics on data residing in different repositories.

To understand the concept of data virtualization, it's useful to look at the different styles/types:

- **Embedded virtualization:** The virtualization technology is functionally embedded in a BI tool, allowing the BI software to make multiple calls to back-end databases to provide a consolidated view for analytics or reporting.
- **Physical virtualization:** Data is retrieved from disparate data sources and consolidated into new physical data structures for consumption from a unified source.
- **Dynamic virtualization:** In this model, a virtualization engine acts as a query orchestration manager, accepting queries and decomposing those queries into subqueries to be run against multiple data sources. The virtualization engine is capable of delegating subqueries by pushing them down to the back-end source system for independent processing. It can also cache data from multiple sources so that it can perform processing tasks on its own. Advanced DV tools will perform distributed, cost-based optimization of the queries.

DV is not a panacea. It is a very useful tool to have in your data and analytics toolkit, but it does not remove the need to physically consolidate collections of data onto servers. For large, complex processing tasks, physical consolidation of data onto servers and processing it there remains the most efficient and cost-effective method. However, for combining the data, the results provided by those servers' DV is extremely useful. Therefore, rather than expect DV to eliminate large servers, expect to see it enabling the integration a small number of large servers, plus enabling access to a range of smaller ones.

For a detailed analysis of data virtualization, see: "[Solution Path for Planning and Implementing the Logical Data Warehouse](https://www.gartner.com/document/code/320563?ref=grbody&refval=3880568)." (<https://www.gartner.com/document/code/320563?ref=grbody&refval=3880568>)

3.4 Logical Data Warehouse

The LDW is Gartner's recommended data management architecture for analytics, combining the strengths of traditional data warehouses with alternative data management and access strategies such as data lakes.

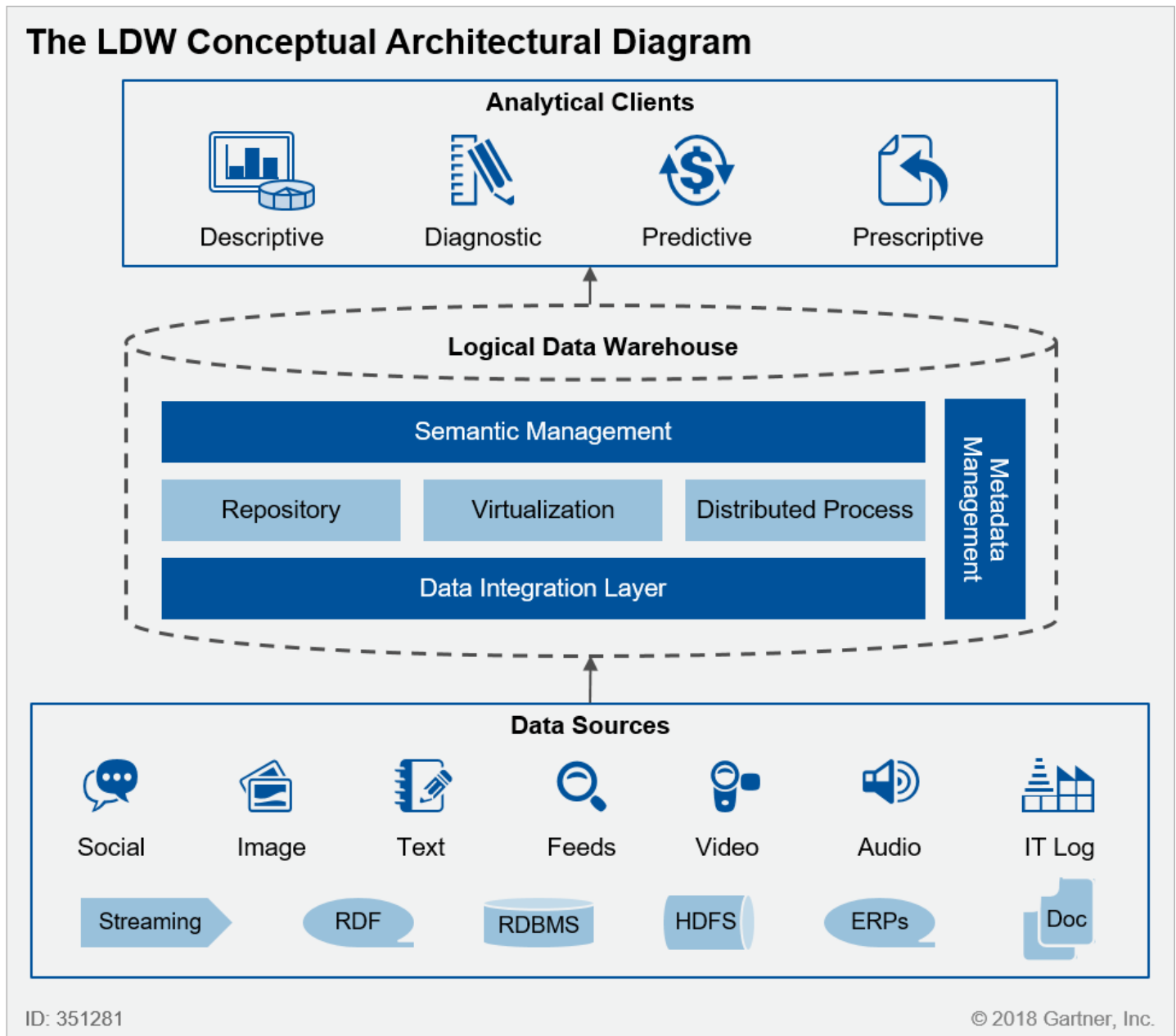
Figure 9 depicts the LDW as a conceptual layer that unifies a collection of architectural components into a connected logical view. This layer provides the logical definitions, processes and repositories that integrate the storage and persistence architecture underneath which LDW data resides.

The original goal of the classic data warehouse was to make all the data in an organization available for analysis. The data warehouse did this by physically copying and transforming all the data into a single special-purpose server, usually an RDBMS. Data today varies in size and nature too much to allow a single-server solution; therefore, the architecture must use multiple engines. It is still desirable to have a unified view of all the data. Therefore, the system must integrate the multiple engines. However, the data is now logically rather physically integrated.

The goal of the traditional data warehouse was to make available all of the data in an organization to enable analysis of the past, present and future.

The mission remains the same but the implementation has changed. The integration is now logical rather than physical.

Figure 9. The LDW Conceptual Architectural Diagram



Source: Gartner (June 2018)

Data architects can implement and evolve the LDW using a combination of three complementary architectural approaches:

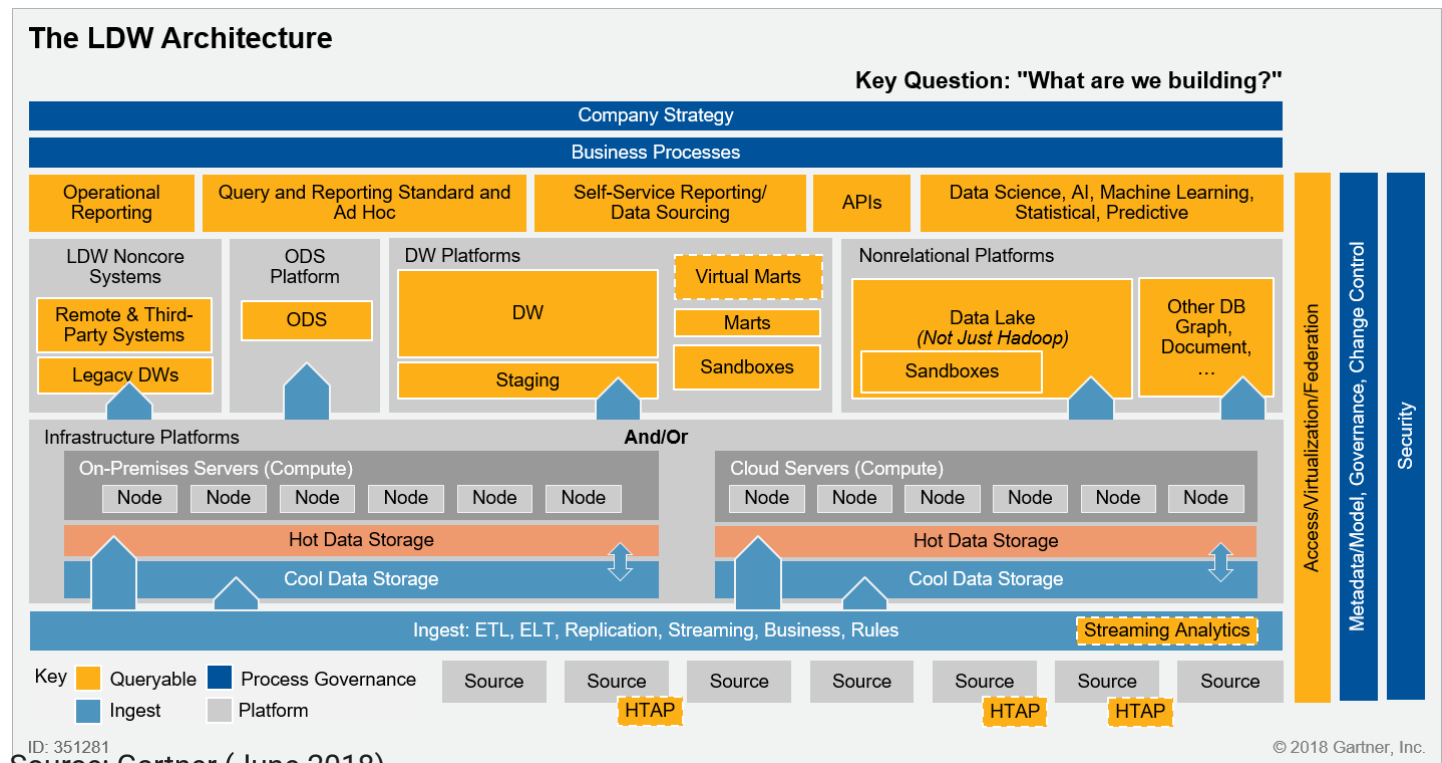
- **Repository, the classic data warehouse:** Build the traditional repository component to deliver high-performance, predictable, prequalified and (possibly) summarized data.
- **Distributed process, the data lake:** Enable working with very large scale and/or unstructured data with source data as near to native format and structure as possible.
- **Data virtualization:** Use DV to combine data from physical or logical data marts, the traditional data warehouse repository, data lakes and other sources. It also provides a uniform way of controlling access to disparate data sources and monitoring their usage.

The LDW is a concept based on the combination of architectural components and is not a commercial, off-the-shelf solution. The components of the LDW can be implemented by a range of software from different vendors and custom-built components. In this way, the LDW strategy offers a separation of concerns, with different parts of the architecture servicing different needs based on required capabilities (see Figure 10).

Modern software components increasingly include features that make it easier to integrate them into an LDW, for example, data virtualization, external database links, open interfaces and metadata interchange.

It is useful for the LDW architect to view these major components, repository (enterprise data warehouse [EDW], marts, operational data store [ODS]), distributed process (data lake) and data virtualization as complementary engines. This is more useful than seeing them as competing solutions.

Figure 10. The LDW Architecture



For example, relational DBMSs can be used for the DW and ODS, while Hadoop Distributed File System (HDFS) can be used for the data lake. Other database technologies, such as graph and document databases, can also be used to service special requirements.

Over time, the convergence of analytic and operational platforms will become more common. In-memory computing (IMC) provides performance that removes the need for competing operational and analytical strategies for accessing data. An emerging technology supporting such convergence is hybrid transaction/analytical processing (HTAP). By using IMC, HTAP systems are capable of performing both online transaction processing (OLTP) and online analytical processing (OLAP) for real-time processing and analytics. IMC is becoming a standard part of many mainstream offerings. This is because of the performance and price performance advantages it confers.

For assistance with LDW implementations and decisions, see:

- "Solution Path for Planning and Implementing the Logical Data Warehouse," (<https://www.gartner.com/document/code/320563?ref=grbody&refval=3880568>) which provides a step-by-step guide to planning and implementing the LDW, along with details about LDW architectures and styles.
- "Embrace Sound Design Principles to Architect a Successful Logical Data Warehouse," (<https://www.gartner.com/document/code/268536?ref=grbody&refval=3880568>) which reviews the LDW's strengths and weaknesses, along with a case study of a successful LDW.
- "Adopt the Logical Data Warehouse Architecture to Meet Your Modern Analytical Needs," (<https://www.gartner.com/document/code/342254?ref=grbody&refval=3880568>) which reviews the logical data warehouse components, architectural patterns, and how to use the LDW for modern analytical needs.

What to Consider Before Moving to the Next Step

Gartner recommends that most organizations should consider the following before moving on to the next step.

- Its flexibility makes the LDW a critical capability for modern data management and analytics. Data architects should align their user SLA requirements with use cases and then determine how to architect their LDW. Populate this design incrementally.
- Evaluate managed cloud-based services, including iPaaS and database platform as a service (dbPaaS) for "greenfield" deployments or where IT skills or resources are lacking.

Step 4: Enable Business Insights

Today's competitive business environment demands effective analytics and business intelligence to enable fast, reliable insights and decisions. Achieving this goal depends on the tasks outlined in this step.

During this stage, data practitioners should be asking the following questions:

- What are the business objectives, and what analytics capabilities are needed throughout the architecture to deliver business value?
- What roles, skills and tools will enable technical professionals to do their own jobs while empowering business users to perform their own data preparation and analysis when appropriate?

The objective should be to deliver analytics capabilities to the decision maker at the point where the decision will have the greatest impact. To accomplish this, the organization needs an architecture nimble enough to respond quickly to changing demands driven by new data sources (internal and external), and new types of data with increasing volume and velocity. In many cases, timely access to insights means pushing data services outside traditional IT boundaries in the form of self-service functionality.

4.1 Business Intelligence and Visualization

Many organizations rely on IT-centric, enterprise-reporting-based platforms for large-scale systems of record reporting. While these platforms serve an essential function, they fail to meet business expectations due to a continued focus on IT-centric analytic strategies. Such traditional BI platforms typically emphasize reports and dashboards based exclusively on OLAP technologies, managed as an isolated solution. In these environments, end users submit requests for reports or data and wait for IT to deliver what they need. In some cases, users have ad hoc access to data through web interfaces or SQL for limited data mining use cases.

However, the BI and analytics market has moved from traditional, descriptive tools purchased and managed by IT to predictive, prescriptive, self-service tools and applications bought and used by lines of business. While enterprise reporting remains an important requirement for organizations, technical professionals need to anticipate the shift to self-service BI and analytics and plan accordingly.

Technical professionals should consider the following trends:

- The need for analytic agility and business user autonomy outweighs the requirement for centrally provisioned, highly governed and scalable system-of-record reporting. As a result, The demand for relevant, easily used and on-demand decision support has eclipsed the capabilities of traditional BI tools that are focused on simply query and reporting.
- Platform buying decisions have shifted more heavily to the business with a de-emphasis on IT's role.
- Today's BI and analytics platforms must provide self-service capabilities, interactive visual exploration, analytic dashboards, and sharing and collaboration.

Basic BI, reporting and visualization capabilities are essential to any organization, and the implementation of these tools is a necessary stage in the data and analytics strategy. However, business modernization depends on empowering business users with the latest tools and the necessary data to explore new analytic opportunities. When evaluating BI platforms and tools, technical professionals should consider the BI platform's ability to support self-service exploration.

The demand for faster insights and proven value has put enormous pressure on IT to evolve existing BI environments. As a result, technical professionals must also identify and evaluate the capabilities of modern analytics and BI platforms to support diverse analytic requirements. Modern BI platforms offer pragmatic, ready-to-use capabilities to give the line-of-business users the insights they need, when they need them. However, aggressive changes in BI platform implementation strategies and vendor lock-in create challenges for technical professionals to decide which platform is fit-for-purpose.

Gartner recommends technical professionals using a multidimensional decision-making approach to decide which platform to align to their analytic workload. The decision made will typically depend on the analytics use cases the IT organization is looking to support, which may include:

- **Traditional enterprise reporting:** Creating trusted, sanctioned and controlled reports and dashboards, which are automatically distributed to a large number of users or embedded in applications.
- **Centralized agile BI/analytics content creation:** This is done by a central analytics and BI team, but with interactive analysis capabilities for report/dashboard/visualization consumers.
- **Decentralized analytics:** Analytics content created by distributed users, with assistance from augmented data preparation and analytics features of BI and data science platforms.
- **Governed data discovery:** IT-governed decentralized analytics, enabling distributed business teams of citizen developers to create their analytical content, but with IT governance, reusability and sharing features.
- **OEM or embedded BI:** The decision to embed analytical capabilities in a business process or an application/portal.
- **Extranet deployment:** Providing external customer access to centrally delivered and managed analytics content.

Many organizations, facing many different analytics and BI products, buy either to consolidate many capabilities into a single suite or to deliver high-value, new capabilities via a best-of-breed product. Moreover, depending on these use cases, organizations may choose a centralized, decentralized or federated support model for analytics, detailed in ["Create a Data Reference Architecture to Enable Self-Service BI."](https://www.gartner.com/document/code/333398?ref=grbody&refval=3880568) (<https://www.gartner.com/document/code/333398?ref=grbody&refval=3880568>)

Refer to the following research for more information about BI:

- ["The Evolving Capabilities of Analytics and Business Intelligence Platforms"](https://www.gartner.com/document/code/353081?ref=grbody&refval=3880568) (<https://www.gartner.com/document/code/353081?ref=grbody&refval=3880568>)
- ["Evaluation Criteria for Business Intelligence and Analytics Platforms"](https://www.gartner.com/document/code/297598?ref=grbody&refval=3880568) (<https://www.gartner.com/document/code/297598?ref=grbody&refval=3880568>)
- ["How to Build Data Visualization Capabilities as Part of a Modern Business Intelligence Platform"](https://www.gartner.com/document/code/297545?ref=grbody&refval=3880568) (<https://www.gartner.com/document/code/297545?ref=grbody&refval=3880568>)

4.2 Self-Service Data Preparation and Analytics

Modern business demands faster and deeper insights from a wider range of data sources than ever before. While IT oversees the architecture and technologies supporting data management, analytics happens throughout the enterprise.

Many use cases described in the previous section, such as traditional enterprise reporting and centralized, agile BI, are not new, but technical professionals are tasked with supporting decentralized analytics and governed data discovery with increasing frequency. As self-service and advanced analytics become more pervasive, technical professionals have an opportunity to shift and expand their role from being BI content creators and data access controllers to user and data enablers. This step is a critical stage of growth in the evolution of an organization's data and analytics strategy.

Technical professionals can enable business analytics by conceptualizing IT's role as providing "data as a service" and enabling analytics throughout the enterprise.

This strategy does not require a lessening of data governance, security and privacy requirements. In fact, data classification and handling becomes even more critical. Leverage the architectural elements mentioned earlier to support data governance, but embrace self-service capabilities to enable business transformation.

At this stage, technical professionals should endeavor to:

- Empower business users to prototype and model new strategies and concepts.
- Establish a culture of innovation with a cross-functional business analytics foundation.
- Allow small, opportunistic projects with limited life spans that serve the needs of one person or a few people to bypass formal architecture disciplines, even for the underlying information model and software technology layers.

Federated Reference Architecture for BI to Enable Self-Service

A federated reference, or data virtualized architecture for BI, enables self-service in organizations with multiple BI implementations. It should be a part of the enterprise architecture to ensure the interoperability and information sharing between semiautonomous noncentrally organized LOBs, IT systems and applications.

For more information on creating a reference architecture for multitool BI environments, refer to "[Create a Data Reference Architecture to Enable Self-Service BI.](https://www.gartner.com/document/code/333398?ref=grbody&refval=3880568)" (<https://www.gartner.com/document/code/333398?ref=grbody&refval=3880568>)

Specialist and Citizen Users

The data scientist now plays an indispensable role in many organizations. Unfortunately, the proliferation of data science use cases creates anxiety among veteran data management professionals, who ask:

- Who are these users?
- Why do they need access to so much data?
- How can we safeguard the data once it's out of our control?

Compounding the problem, business users and self-styled analysts are demanding unprecedented access to data and analytical tools, often under the banner of data science. These "citizen" users may have legitimate use cases, but technical professionals are dubious.

The technical professionals fear risks such as data leakage, excessive query load on production databases, misuse of data, or the combination of datasets that turn safe, governed data into noncompliant, risky data. In order to support these use cases, technical professionals must take advantage of machine intelligence and new technologies to support self-service workflow.

Support Self-Service Workflow

Self-service capabilities accelerate time to insight by giving business users a way to find, access, clean and prepare data for analytics. Technical professionals can support self-service workflow in the following ways:

Self-Service Data Preparation

Self-service data preparation technologies can be stand-alone or embedded in modern BI and advanced analytics platforms. They give users the tools to combine, prepare and manipulate data as well as collaborate with others by providing descriptive metadata. Self-service data preparation tools can access

and combine data from on-premises and cloud repositories and enable blending of enterprise data with data acquired from partners and third parties, such as data management platforms.

Provide Data Services

In lieu of, or along with, self-service data preparation, IT can help speed the creation of curated, trusted data for a range of distributed analytics content authors. This function is being addressed by the emerging role of data engineer: a centralized data specialist within IT.

Data engineers have the technical expertise and familiarity with IT processes, technologies and requirements to expedite curated datasets on behalf of business and specialist users. They typically possess software development experience and are capable of writing complex queries. They work closely with data scientists and analysts to understand requirements and design systems and processes that support their users.

4.3 Advanced Analytics

Once the traditional analytical tools that comprise basic BI and reporting are in place, organizations should turn their focus to advanced analytics capabilities that assist business decision making. Advanced analytics empower business stakeholders to conduct "what if" analysis to envision the outcome of events based on today's decisions. For example, manufacturers can analyze buying patterns, forecast future trends and optimize inventory accordingly.

From Explanatory to Exploratory

While traditional data reporting and analysis tend to look in the rearview mirror, advanced analytics peer into the future.

- **Predictive analytics:** Answers the question of "What is likely to happen?" This category relies on techniques such as predictive modeling, regression analysis, forecasting, multivariate statistics and pattern matching.
- **Prescriptive analytics:** Addresses the question of "What should be done?" or "What can we do to make X happen?" This category relies on techniques such as graph analysis, simulation, complex-event processing, recommendation engines, heuristics and, increasingly, neural networks and machine learning.

Technical professionals should evaluate advanced analytics through two different lenses: IT-related and business-related.

IT can benefit from advanced analytics by looking at its own internal processes, systems and metrics to identify opportunities to improve operational effectiveness.

In addition to using advanced analytics to improve IT, technical professionals play a critical role in supporting business analytics. Technical professionals provide the architecture, tools and in some cases, prepare the data supporting business-related analytics. This requires planning for new data types, data sources and real-time use cases.

What to Consider Before Moving to the Next Step

Gartner recommends that most organizations should assess capabilities before moving on to the next step. Data architects should ask themselves the following questions:

- Do we have the foundational BI and reporting capabilities?
- Can we support self-service capabilities for our BI platforms?
- Can we support specialist and citizen users with the data they need to accomplish their goals?
- Can we provide this data while applying appropriate safeguards within established data governance policies?

Step 5: Extend and Automate With Artificial Intelligence and Machine Learning

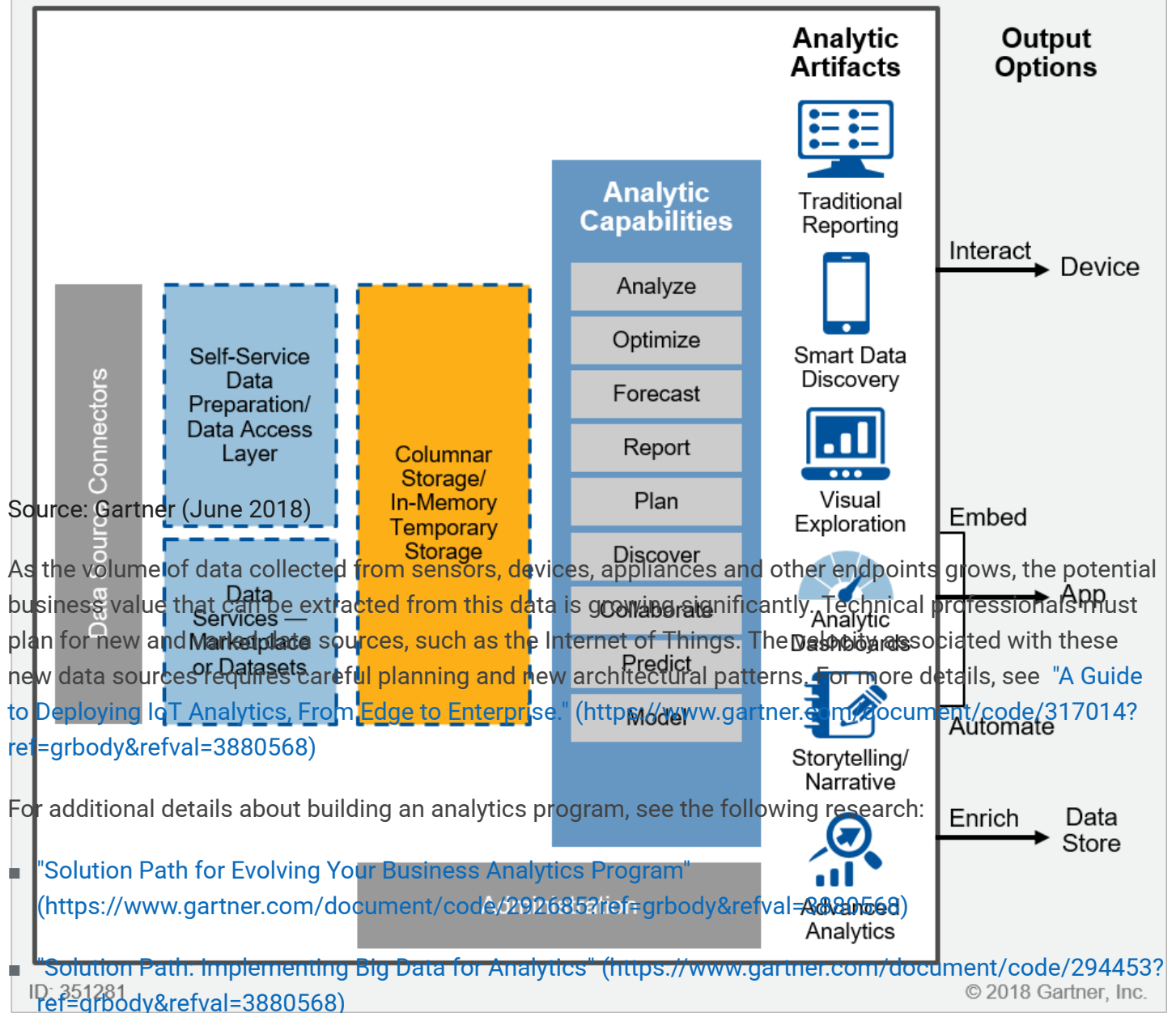
This step expands the strategy to enable the full spectrum of analytic capabilities by automating processes based on real-time data analysis and adopting machine learning to add intelligence to automation, process improvement and decision making.

5.1 Real-Time Analytics

Supporting the full range of analytic capabilities requires an architecture that can support ingestion of streaming data, perform real-time analysis, take automated actions and deliver business insights. Figure 11 provides a reference architecture for an overall model of the key components and analytic capabilities.

Figure 11. Business Analytic Reference Architecture — From Data to Insight to Action

Business Analytic Reference Architecture — From Data to Insight to Action



- "Hyperscaling Analytics: Comparing Streaming Analytics in the Cloud With AWS, Microsoft Azure and IBM" (<https://www.gartner.com/document/code/310019?ref=grbody&refval=3880568>)

5.2 Machine Learning

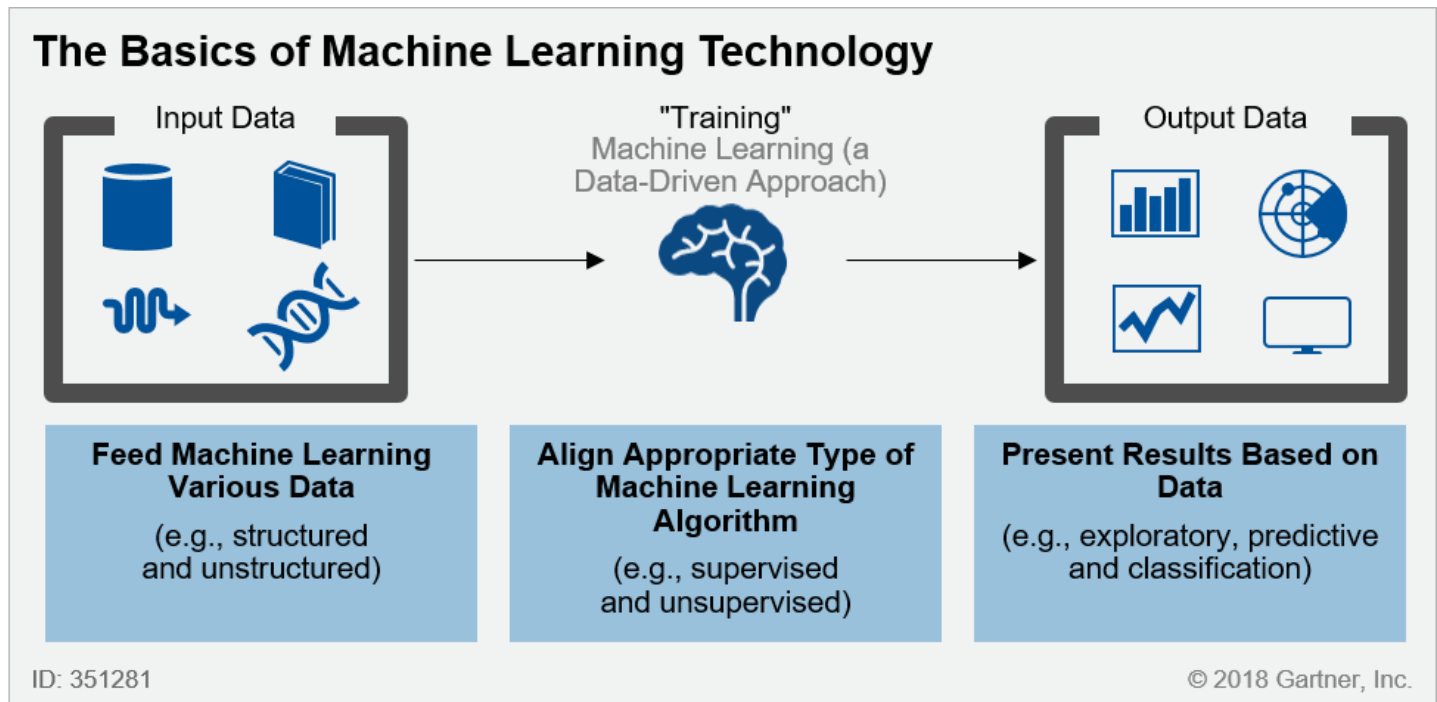
Although data mining and manual analysis can be powerful, many organizations don't have the resources to broadly apply data science. Machine learning (ML) is an attractive option for evolving and automating analysis.

ML, a subset of artificial intelligence (AI), is more than a technique for analyzing data. It's a system that uses data through algorithms to produce desired output — prediction and data-driven decisions — directly

from data. ML enables advanced systems that can appear to be trained, predict and even operate autonomously, rather than being programmed for only a finite set of actions.

ML can leverage data and perform computations to find insights in large datasets. Using input data, as shown in Figure 12, ML recognizes patterns that can be used to make a prediction or classify object and/or observations.

Figure 12. The Basics of Machine Learning Technology



Source: Gartner (June 2018)

Architect for Machine Learning

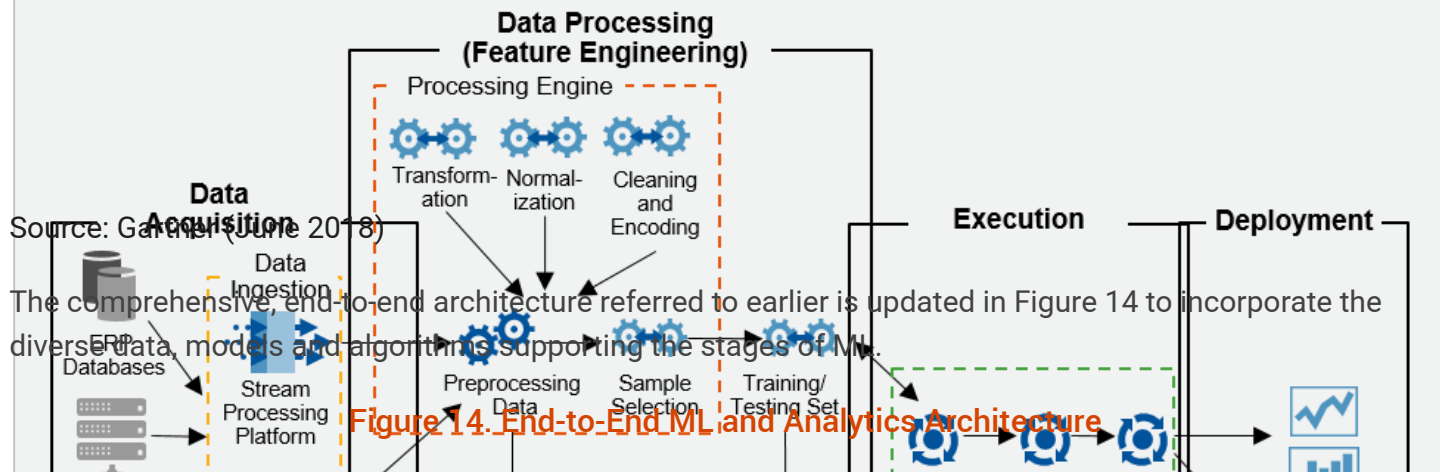
ML demands a flexible architecture capable of supporting elastic learning patterns, and consuming large and varying volumes of data. These design requirements often necessitate considerable processing power and storage versatility. With its elastic capabilities and scaling algorithms, cloud infrastructure is an

excellent proving ground for new ML initiatives. Figure 13 shows Gartner's suggested reference architecture for ML covering the functional areas required for the ML process:

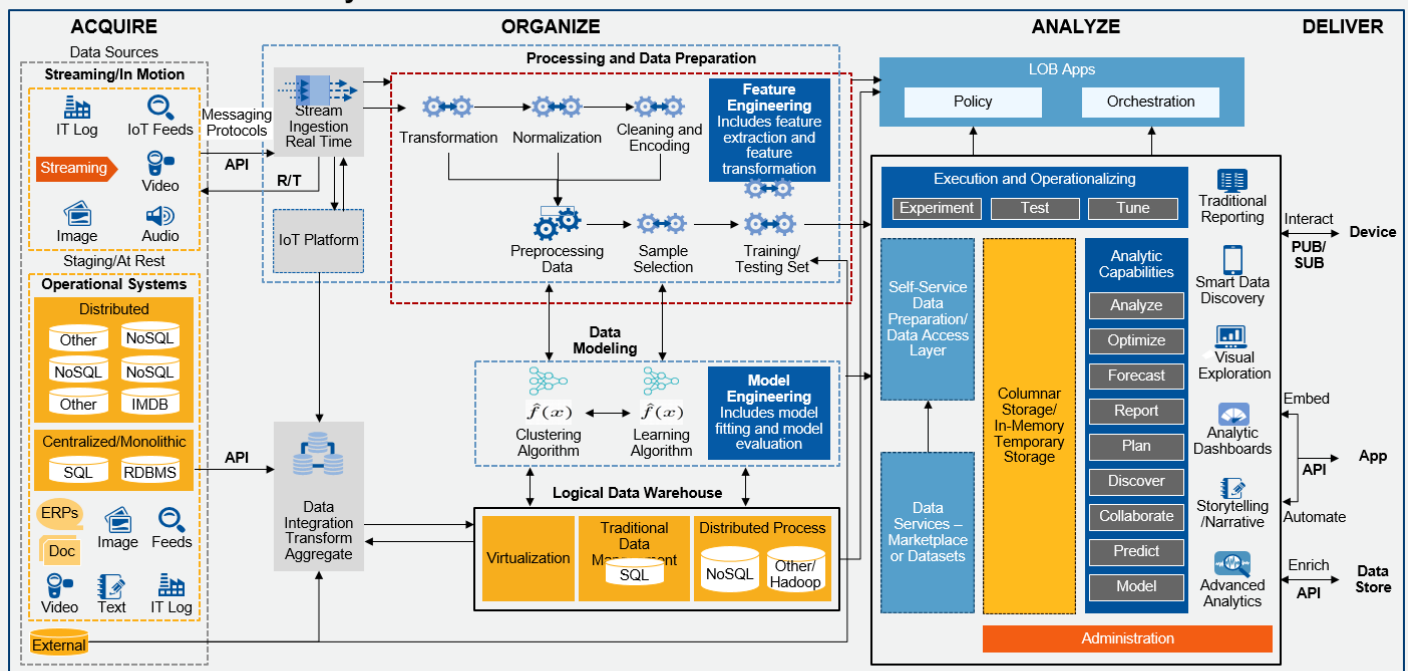
- **Data acquisition:** Where data is collected, prepared and forwarded for processing.
- **Data processing:** Where steps such as preprocessing, sample selection and the training of datasets take place, in preparation for execution of the ML routines:
 - Feature analysis or feature engineering (a subset of the data processing component), where features that describe the structures inherent in your data are analyzed and selected
- **Data modeling or model engineering:** Includes the data model designs and machine algorithms used in ML data processing (including clustering and training algorithms):
 - Model fitting, where a set of training data is assigned to a model to make reliable predictions on new or untrained data
 - Model evaluation, where models are evaluated based on performance and efficacy
- **Execution,** the environment where the processed and trained data is forwarded for use in the execution of ML routines (such as experimentation, testing and tuning).
- **Deployment,** where business-usable results of the ML process — such as models or insights — are deployed to enterprise applications, systems or data stores (for example, for reporting).

Figure 13. Machine Learning Architecture

Machine Learning Architecture



End-to-End ML and Analytics Architecture



Manage and Govern
Information Governance (including Metadata Management, Data Quality, Data Modeling, Master Data Management), Data Management (Data Admin, Security, Privacy and Identity), Organization (People)

□ = Optional

□ = Cloud, On-Premises or Hybrid

ID: 351281

© 2018 Gartner, Inc.

Source: Gartner (June 2018)

Steps to get started with machine learning:

- Learn about and experiment with ML concepts and technology.
- Work closely with data science teams and business users to identify a use case.
- Build a use case in the cloud.
- Iteratively expand your ML platform and services over time.

For more details, see the following Gartner research:

- "Preparing and Architecting for Machine Learning" (<https://www.gartner.com/document/code/317328?ref=grbody&refval=3880568>)
- "How to Create a Data Strategy for Machine Learning-Powered Artificial Intelligence" (<https://www.gartner.com/document/code/324342?ref=grbody&refval=3880568>)
- "Making Machine Learning a Scalable Enterprise Reality – From Development to Production" (<https://www.gartner.com/document/code/343614?ref=grbody&refval=3880568>)

Gartner recommends that technical professionals adopt ML techniques as part of their personal "tradescraft." This will improve their ability to support digital business efforts, as well as tackle data management and operations challenges that arise within IT.

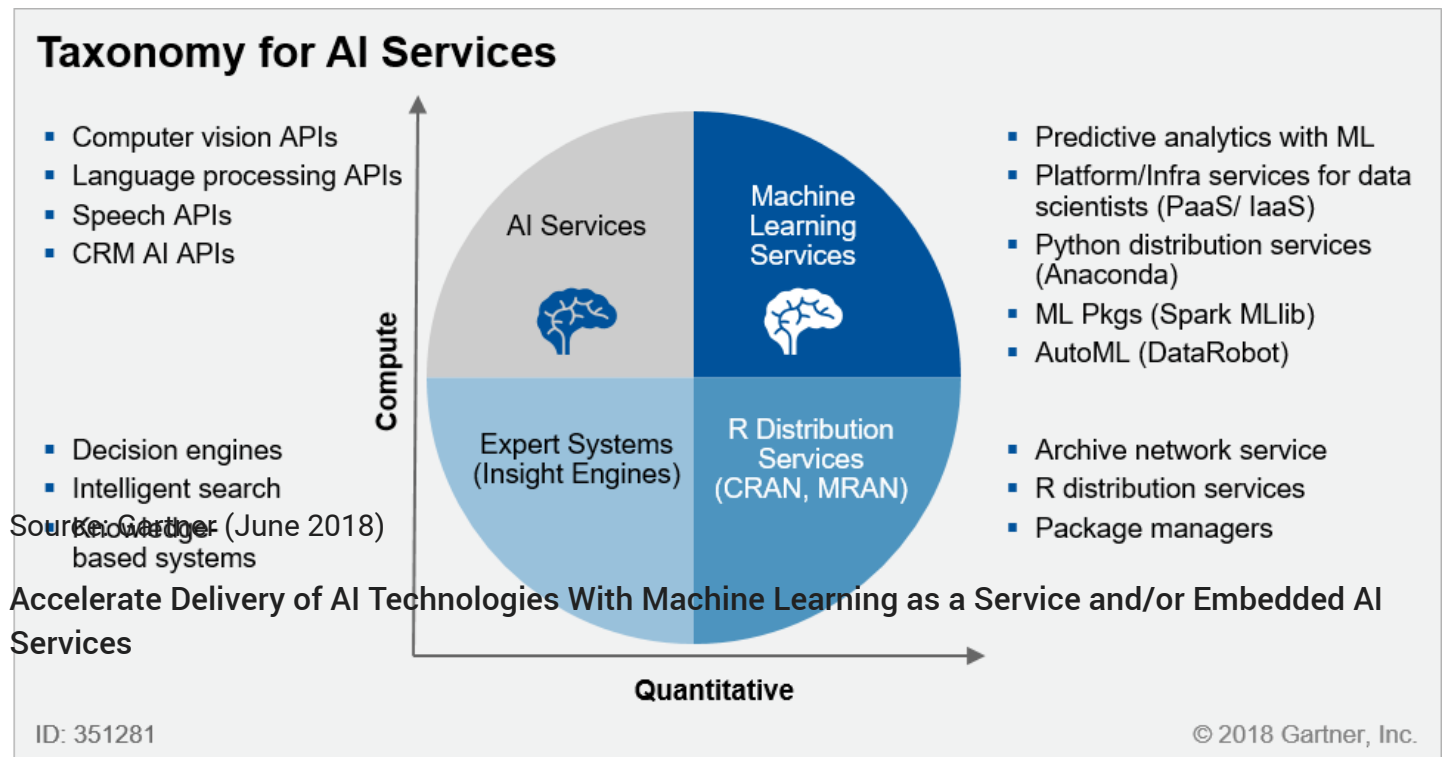
Leveraging Artificial Intelligence and Machine Learning API Services to Complement Your Data and Analytics Architecture

Gartner defines AI as a technology or system that can emulate partial human cognitive performance. AI leverages a group of technologies that include machine learning, natural language processing (NLP), deep learning and expert systems.

As AI and ML technologies become more pervasive, technical professionals have an opportunity to leverage services that seamlessly integrate into their existing architecture. Cloud-based AI and ML API services are commonly being used to leverage AI and ML solutions. These services must also seamlessly integrate into new and existing business operations.

Figure 15 provides a basic taxonomy of AI services commonly available for integration into new and existing business operations. These services are outlined based on their computational requirements and quantitative methods.

Figure 15. Taxonomy for AI Services



Simplify solutions by connecting to data assets for ML, rather than collecting data.

Machine learning as a service (MLaaS) is the application of specialized AI competencies through quantifiable measurement processes, which makes the development, environments, libraries and algorithms more accessible to multiple stakeholders. It can be used to accelerate the delivery of AI solutions, and is often used to reduce the burden on the technical professional or organizations that lack the skill and experience needed to deploy AI solutions.

MLaaS accelerates the delivery of AI solutions by offering:

- Predefined components that eliminate the need to build components from scratch
- Reusable components that are shared across the enterprise to improve collaboration
- Common libraries and toolkits available and shared across the enterprise to accelerate development

There are three critical capabilities that aid the technical professional in delivery AI or ML solutions:

- **Build capabilities:** Services can be used to build and train ML models across the enterprise
- **Deployment capabilities:** Services can be used to consistently deploy to a centralized or decentralized environment

- **Operational capabilities:** Services can be used to operate, maintain, monitor and manage environments on an ongoing basis.

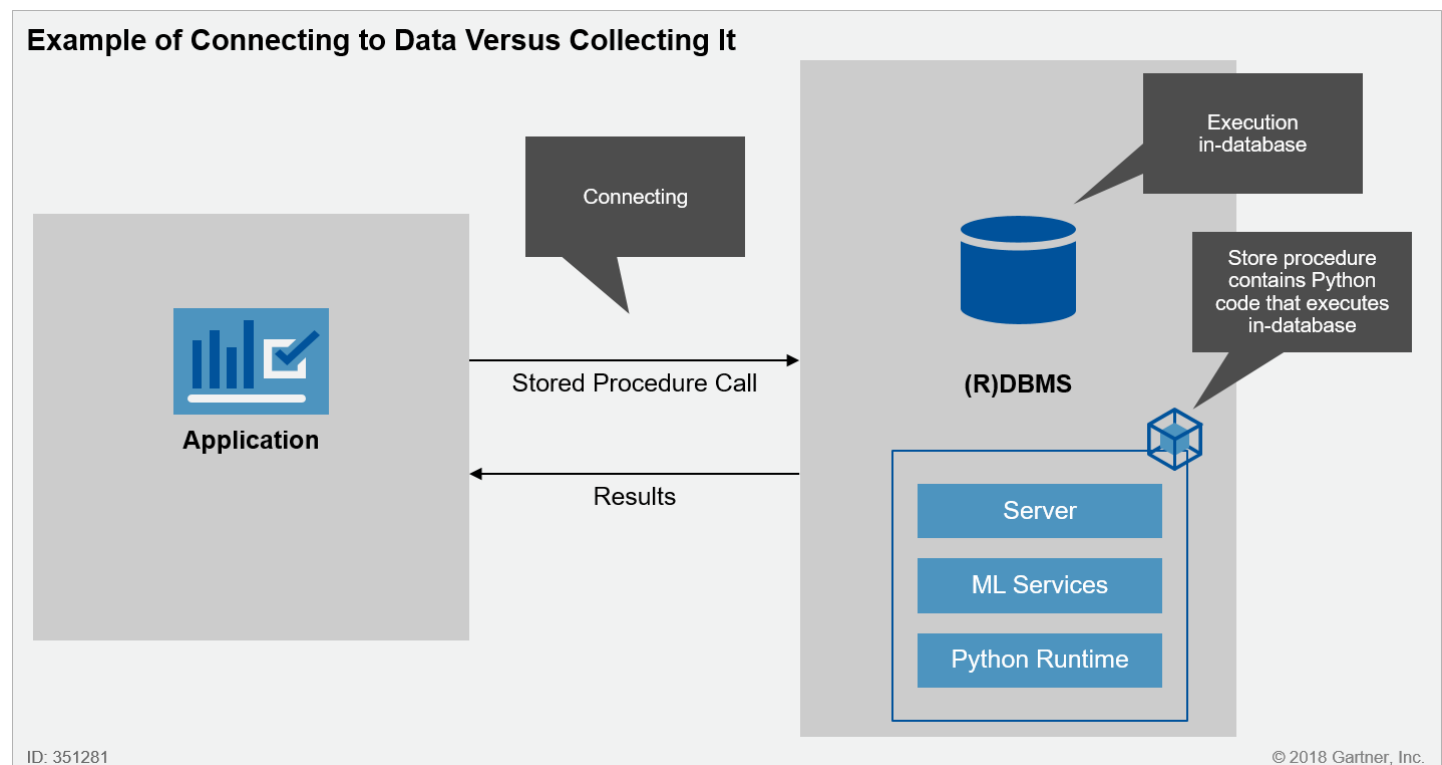
Cloud-based MLaaS and AI services offer a promising opportunity to augment the data and analytics architecture, as many AI and ML solutions remain computationally and process-intensive. Gartner recommends technical professionals leverage MLaaS and embedded AI services to augment the data and analytics architecture to support existing and future capabilities.

For example, Amazon SageMaker can be a useful, complementary component that enables the capability to build, train and deploy machine learning models at scale, while integrating with existing data pipelines. Technical professionals should consider adding similar components to their data and analytics strategies to empower knowledge workers to develop ML solutions.

Another method for accelerating the delivery of ML solutions is to leverage architectural patterns that allow you to connect to data, rather than collecting it. For example, many database management systems embed algorithms (pretrained and/or optimized) as user-defined functions (UDFs). This allows users to execute functions directly on the data, rather than having to collect and export the data to a different environment, as shown in Figure 16.

Microsoft SQL Server 2017 is an example of an RDBMS that allows users to use Python code and UDFs within T-SQL to execute ML as a function without having to transport the data to a different environment. This is an emerging pattern that should be a part of a comprehensive data strategy.

Figure 16. Example of Connecting to Data Versus Collecting It



Source: Gartner (June 2018)

For more details, see the following Gartner research:

<https://www.gartner.com/document/3880568?ref=feed>

- "Driving an Effective AI Strategy" (<https://www.gartner.com/document/code/356807?ref=grbody&refval=3880568>)
- "Comparison of Amazon, Google, IBM and Microsoft AI Cloud Services" (<https://www.gartner.com/document/code/332643?ref=grbody&refval=3880568>)

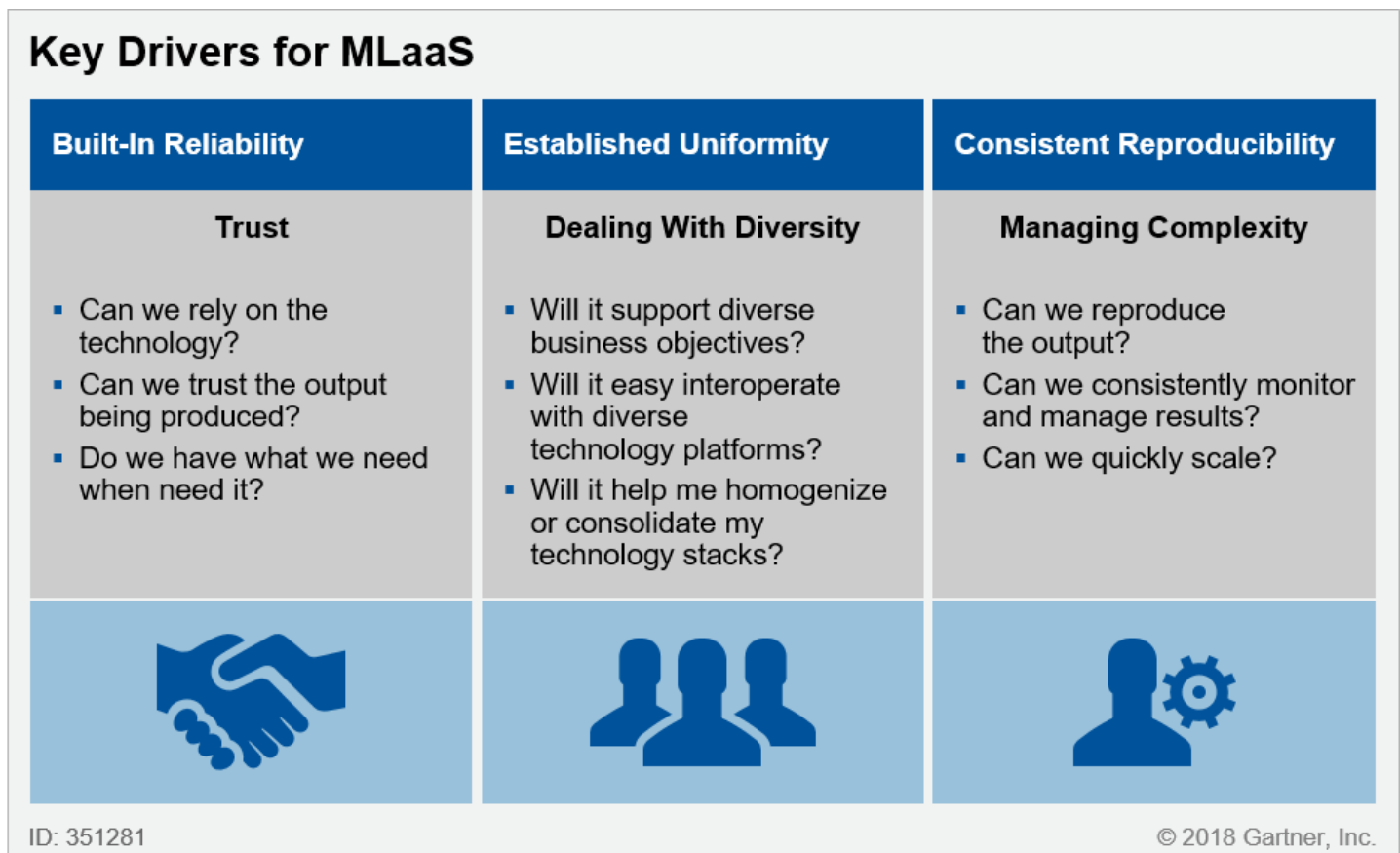
Get Started With Accelerating Delivery of AI Technologies by Focusing on These Drivers

Modern data and analytics strategies will most likely need to incorporate AI and ML technologies. However, before leveraging these technologies, Gartner recommends technical professionals focus on key drivers that increase the probability of success and trust within your data environments.

The following key drivers have a major impact on the performance and outcome of the analytic solutions, as outlined in Figure 17. Gartner recommends focusing attention on the drivers throughout the enterprise data strategy to ensure that analytic solutions are aligned with the business.

- Built-in reliability with adequate data quality
- Established uniformity with crowdsourced semantics
- Consistent reproducibility with automated analytic deployment

Figure 17. Key Drivers for MLaaS



Source: Gartner (June 2018)

5.3 Automate and Scale

This step focuses on supporting the architectural capabilities outlined in the previous steps with process automation and scalability.

Data Science at Scale

Analysts and data scientists leverage specialist, stand-alone tooling when data integration requirements exceed the capabilities embedded in analytics platforms. Many data scientists also supplement embedded capabilities with multiple open-source tools, such as Python, R and Scala, or tools that come with Hadoop distributions for different parts of the data pipeline process. Using these stand-alone tools, they extract small amounts of data from the data lake and perform analysis on their local systems.

There are several problems with this stand-alone strategy:

- Users may base models on incomplete data or datasets that are too small
- Storing data on personal laptops leads to data sprawl
- Data scientists spend time acquiring and preparing data rather than on high-value tasks
- Processing occurs on stand-alone systems instead of distributed, scale-out platforms
- Models and algorithms may require refactoring for production rollout

By implementing self-service data preparation tools, combined with support for data science utilities and platforms, technical professionals can accelerate and streamline the entire process, resulting in more efficient and productive data science teams.

This approach also requires a process and guidelines for how business analysts and data scientists can "promote" or operationalize findings and models into trusted, recurring analysis, either themselves or via formal data integration practices. See Step 6.3, Model Deployment for Data Scientists, for an example of processes and guidelines for promoting analytic models.

Decision Support

Many organizations claim that their business decisions are data-driven. But they often use the term "data-driven" to mean reporting key performance metrics based on historical data and using analysis of these metrics to support and justify business decisions that will, hopefully, lead to desired business outcomes. While this a good start, it is no longer enough. ML and predictive analytics can take decision support to new levels.

Using ML to Support Decisions

Forecasting is nothing new, and organizations have long relied upon a combination of intuition, expertise and past results to anticipate future state. For decades, organizations have applied data analysis based on statistical models to distinguish previously unknown relationships and patterns within data. However, the confluence of ML techniques and advances in computing and processing capabilities have increased the

speed and efficiency of predictive analytics — especially in cases where designing and programming explicit algorithms is either impractical or unfeasible.

ML is now being used in a variety of business and IT contexts to support decision making, as shown in Tables 1 and 2.

Table 1: ML Decision Support Examples — Business Context

Business Category ↓	Decision Context ↓
Supply Chain	Analyze buyer behavior and purchasing trends to identify changes in the marketplace and forecast demand.
Sales	Analyze CRM data to identify actions and events conducive to desirable sales outcomes.
Manufacturing	Use sensor data to anticipate failure based on events, temperatures, usage, etc.
Customer Service	Analyze churn probability to prioritize service and sales efforts to maximize retention.
Marketing	Identifying negative brand sentiment through mentions in social media

Source: Gartner (June 2018)

Table 2: ML Decision Support Examples — IT Context

IT Category ↓	Decision Context ↓
IT Operations Management	Anticipating equipment failure and implementing proactive maintenance to prevent disruption.
Information Security	Identifying and resolving IT and security incidents by automatically detecting anomalies and patterns in data

Source: Gartner (June 2018)

Automated Actions

In this step, data and analytics become the brain of the enterprise — becoming proactive as well as reactive, and coordinating a host of decisions, interactions and processes in support of business and IT outcomes.

Infuse predictive intelligence directly into your systems to combine ML and advanced analytics with workflows and automated processes:

- Shape and mold external and internal customer experiences, based on predicted preferences for how each individual and group wants to interact with the organization.
- Drive business processes, not only by recommending the next-best action but also by triggering those actions automatically.

As organizations become more adept at using data to drive decisions, the next logical stage is using data to take action. The ability to automate decisions and actions based on information is a powerful differentiator in modern business. Technical professionals can expand on the examples listed earlier to take automated actions (see Tables 3 and 4).

Table 3: ML Automation Support Examples – Business Context

Business Category ↓	Automated Actions ↓
Supply Chain	Trigger automatic inventory replenishment, implement discounts and promotions to reduce inventory.
Sales	Classify, prioritize and channel leads based on scoring models. Optimize models as data evolves and route leads automatically.
Recommendations	Make real-time product recommendations, ad placement, content filtering, and ranking.
Fraud Detection	Block transactions or suspend accounts when fraudulent activity is detected.
Failure Prediction	Ship replacement parts and schedule service.

Source: Gartner (June 2018)

Table 4: ML Automation Support Examples – IT Context

IT Category ↓	Description ↓
Security	Security actions can be based on data thresholds, patterns or trends. For example, security anomalies can trigger account lockouts or temporarily suspend transactions.
IT Operations	Engineers can have new nodes added automatically to a compute cluster when metrics exceed preset thresholds.
Service and Support	Algorithmic chatbots using ML can become more effective as they learn from each interaction and draw on contextual information about their users.

Source: Gartner (June 2018)

See "2018 Planning Guide for Data and Analytics." (<https://www.gartner.com/document/code/331851?ref=grbody&refval=3880568>)

Step 6: Deploy and Integrate Analytics Into Operations

Modern analytics requires a set of procedures and techniques concerned with deciding how best to deploy and operate analytic systems, usually under conditions requiring the allocation of scarce resources.

The last phase of the data strategy is the implementation of an analytic solution that has been tested and approved into operations. Delivering analytics for operations requires a deployment strategy. This step focuses on deploying and integrating the analytics developed in the previous steps into operations to drive user adoption and improve usability. Analytics may also be embedded into existing products or services. Successful use of analytics requires a process for integrating analytics within the organization and business domains. This process should focus on four elements:

- **Accessibility:** The analytic output should be accessible to everyone, but gated by security and identity access management rules.
- **Collaboration:** The capability to share analytic output, perspectives and insights to align stakeholders in solving business problems.
- **General-purpose utility:** The analytic output can be used across different business domains, regardless of geographic location.
- **Universal architecture:** There is a common language and architectural framework used to access the analytic output.

Both IT and business experts have an interest in governing analytics deploying into operations. As a result, Gartner views an optimal deployment strategy as a two-step process: planning and performance. Though there are many processes for deploying analytics into operations, most include:

- IT planning for transition activities, such as system integration and testing.
- Business performance intervention to determine how the analytics will be used, how it will be executed, and how it will drive the business strategy

6.1 Analytic Deployment Architectures

Look for solutions for deploying your analytic workloads anywhere – especially in existing business processes.

Analysts and data scientists develop and deploy analytic results that must be collaborative and made available across the enterprise and in growing diverse environments. Analytic deployment architectures focus on aspects of the analytic solution that are important after the solution/system has been tested and is ready to go into live production. The deployment architectures may reside on top of a centralized

infrastructure, such as an on-premises Hadoop distribution, a distributed one, such as a cloud provider (e.g., AWS, Google, etc.) or a hybrid cloud solution (e.g., Microsoft, IBM). Gartner recommends starting with a centralized infrastructure and evolving to a distributed infrastructure once the organization has adopted your enterprise data strategy.

Analytic workloads are often tightly coupled with platform environments and infrastructure, making deploying analytics stagnant. For example, much of the analytics of today are deployed to a business intelligence environment or vendor tool suite. Defining multiple analytic deployment architectures provides users with a more flexible delivery model – one that is based on pushing analytics to multiple delivery endpoints.

There are three common types of analytic deployment architectures:

- Open-source stack analytic deployment architectures for moving analytic workloads to a production-ready platform or service (e.g., Kubernetes/Kubeflow)
- Container-centric deployment architectures for preconfiguring virtual machines for analytic development and deployment (e.g., Microsoft Azure Data Science Virtual Machines)
- Proprietary architectures that enable distributed, yet collaborative, analytics at scale (e.g., Datameer)

There are several advantages of analytic deployment architectures:

- Users are able to integrate the results of their analyses into existing systems and business processes.
- Organizations enable faster decision making by pushing insights to the users, rather than pulling them via a request, query or inquiry.
- Users are able to share insights beyond BI reporting tools.
- Improves user adoption as insights are seamlessly delivered to the user.

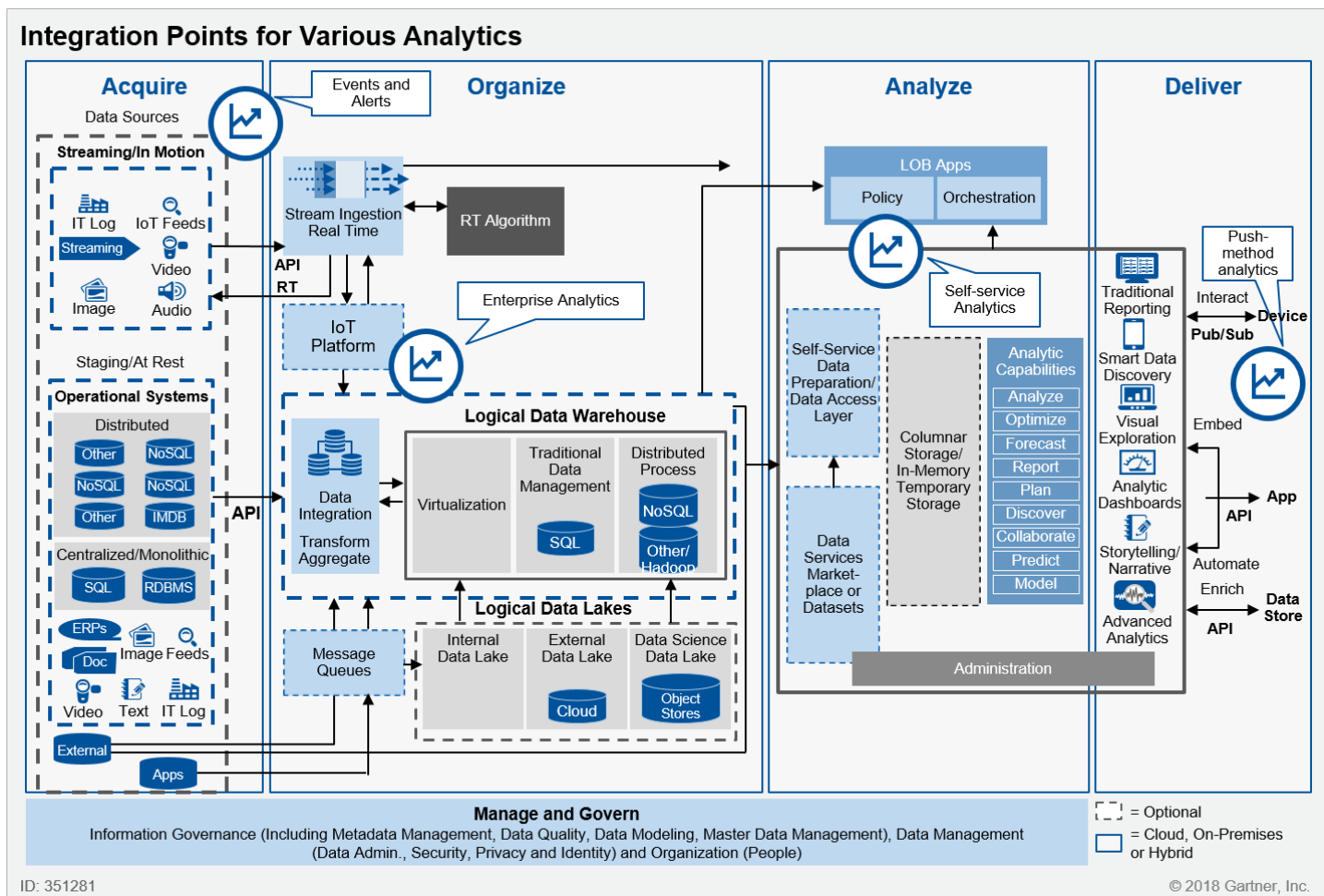
6.2 Points of Integration for Existing Analytic Systems

The data and analytic architecture and the applications that are built upon it are not meant to be a silo. To be actionable, the insights derived from each component of the data and analytics architecture must be shared, and the architecture that's the most capable of integration at various points of the architecture is the one that delivers more value. Gartner's end-to-end data and analytics architecture has numerous integration points for analytical applications, enterprise applications, repositories and plug-ins, making it comprehensive in terms of functionality and integration capabilities.

To summarize, each component of the data and analytics architecture should be available for analysis to the enterprise. In addition, points of integration should allow for interoperability and continuous deployment of new analytics to support operations and the LOBs.

Figure 18 illustrates Gartner's end-to-end data and analytics architecture with multiple points of integration for deploying different types of analytics.

Figure 18. Integration Points for Various Analytics



Source: Gartner (June 2018)

There are four main points of integration in the end-to-end data and analytics architecture:

- Event and alert analytics that support real-time analysis and decision making
- Enterprise analytics that support the enterprise reporting solutions
- Self-service analytics that support the democratization of analytics throughout the enterprise

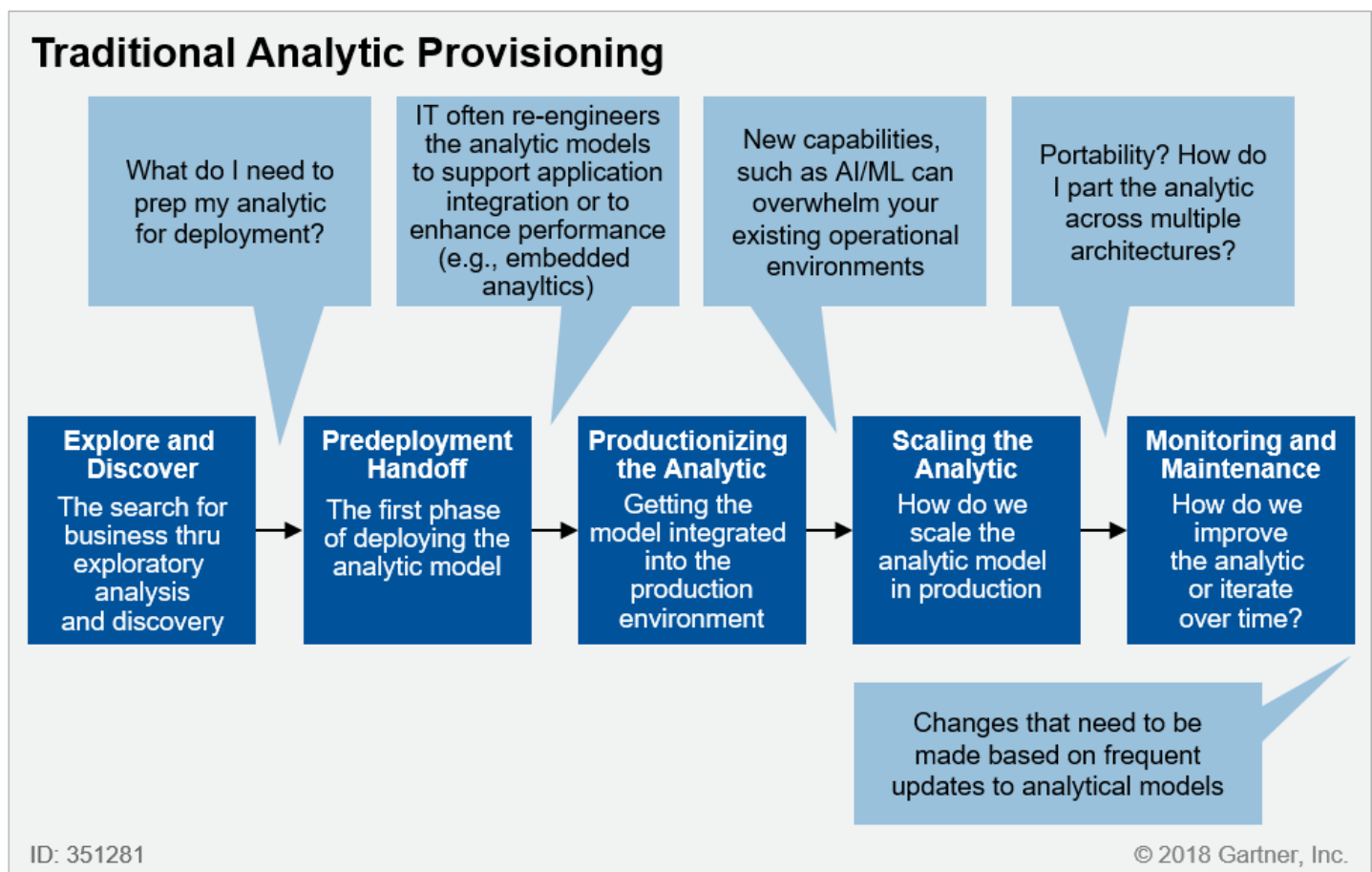
- Push-method analytics that support the delivery of analytic models and outputs to consumers, without having to query or develop logic to retrieve information for the purpose of analysis

6.3 Model Deployment for Data Scientists

There are five core competencies of analytic deployment. Although many analytic deployment strategies will be similar despite the complexity of the analytics, model deployment for data scientists should pay closer attention to these competencies to ensure successful delivery.

Figure 19 illustrates a typical analytic provisioning model and competencies around model deployment for data scientists.

Figure 19. Traditional Analytic Provisioning Model



Source: Gartner (June 2018)

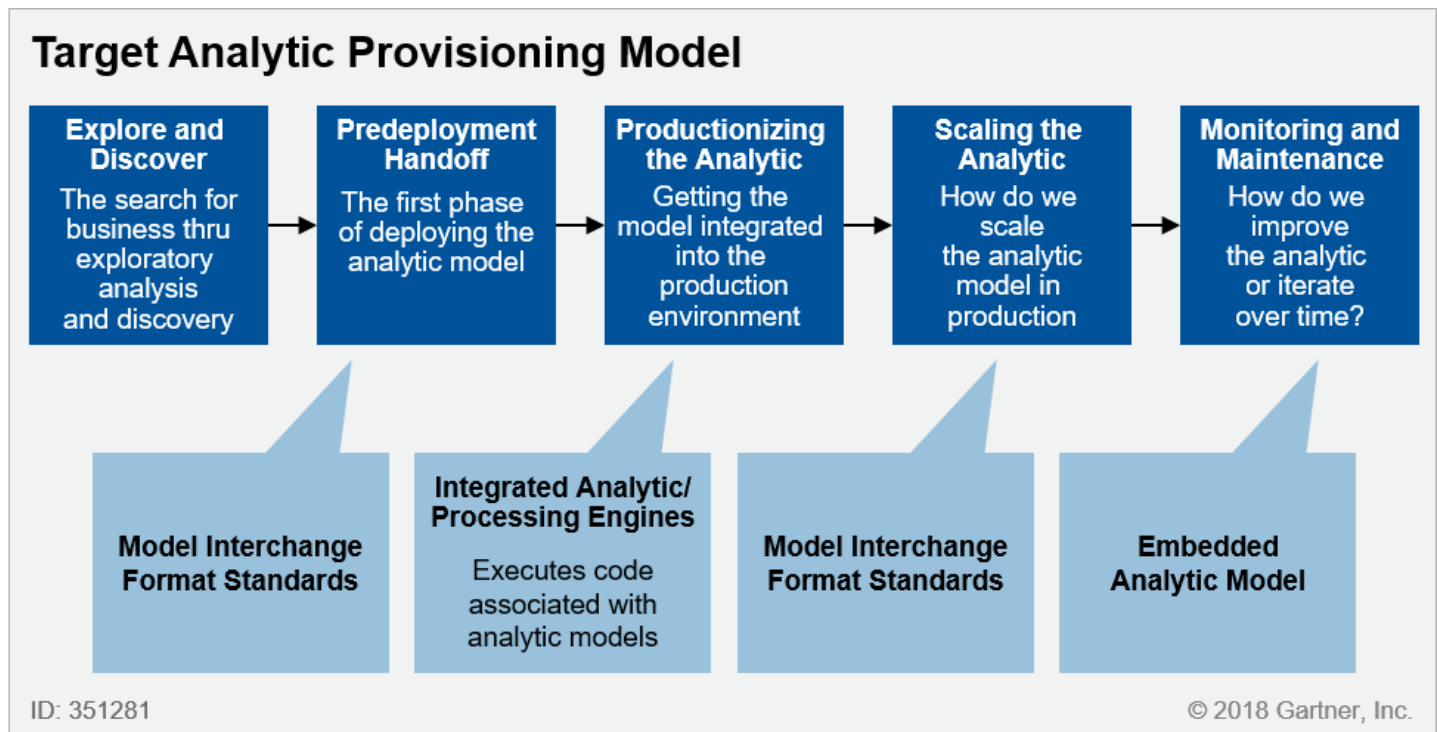
To answer the questions outlined in Figure 19, Gartner recommends leveraging a target analytic provisioning model based on three technical components:

- Model interchange format standards, such as the Predictive Model Markup Language (PMML) and Portable Format for Analytics (PFA), enable a common interface for model exchange.

- Integrated analytic processing engines allow for the execution of multiple types of machine learning models, such as Spark, TensorFlow and PyTorch.
- Embedded analytic model output allows for the integration of analytic output into existing systems.

Figure 20 illustrates a target provisioning framework that should be used by both data scientists and technical professionals.

Figure 20. Target Analytic Provisioning Model



Source: Gartner (June 2018)

Iterate for Continuous Improvement

Technical professionals should revisit the following steps when adjusting the architecture, adding new data stores, or implementing new technologies. Reviewing these steps ensures alignment with business goals, promotes comprehensive data governance and creates technical consistency in the selection of tools and services. Organizations should use the steps in this research to iteratively improve their data management and analytics architecture as they expand into new use cases.

Align With Business Strategy

At the beginning of each project, senior managers and business stakeholders should explain the vision and strategy of the company. This can occur through initial meetings, interviews and company documentation,

such as the annual report. This step ensures that the team understands how information needs relate to high-priority items in company strategy. It also allows prioritization to balance the longer-term strategic goals that may have very high value, and which are essential to company growth and survival, alongside other immediate and urgent requirements.

Review business drivers and business outcomes. Start with the problem you are trying to solve. Examples:

- How can we identify and reduce/prevent fraud?
- How can we optimize purchasing across divisions?

If possible, reframe business challenges within the context of analytical questions to be answered.

Identify the datasets needed to solve those problems: What are the use cases, and what changes are needed in the architecture?

Figure 21, from "Key Recommendations for Implementing Enterprise Metadata Management Across the Organization," (<https://www.gartner.com/document/code/315412?ref=grbody&refval=3880568>) shows the relationship between business and technical metrics.

Figure 21. Aligning Business Strategy With Development

Aligning Business Strategy With Development

Source: Gartner (June 2018)

Governance

Revisit the components of information governance outlined earlier. Take steps to ensure consistent application with each new project or additional technology.

It can be useful to adopt information governance stages. This will ensure that new initiatives are properly socialized within the organization, that there is agreement on policies and procedures, and that various aspects of information management are considered during the project (see Figure 22).

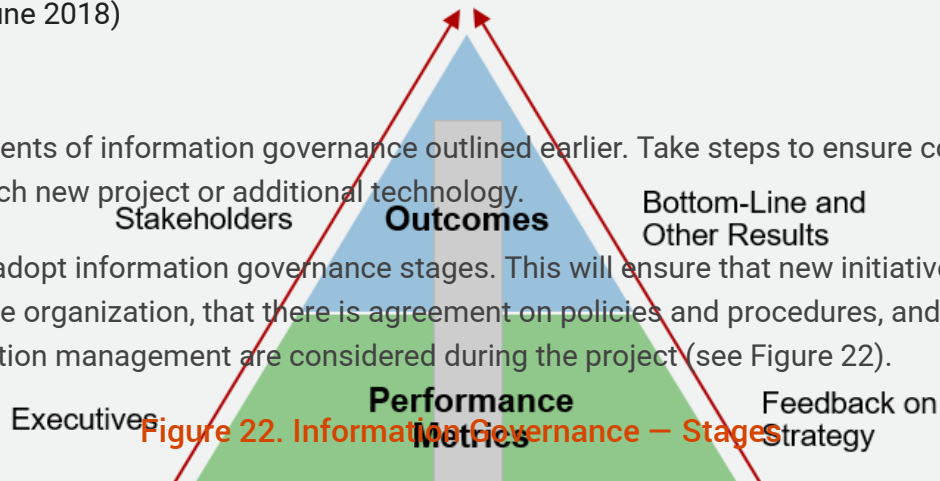


Figure 22. Information Governance – Stages

Information Governance — Stages



Awareness

- Promote visibility
- Training
- Auditing and reporting



Compliance

- Regulatory compliance
- Corporate policies and procedures
- IT policies and procedures



Information Profiling

- Data inventory and analysis
- Data quality, veracity
- Lineage and provenance



Privacy and Security

- Information classification
- Privacy management
- Role-based access controls



Information Life Cycle

- Retention policy
- Hot, warm, cold – tiered storage
- Archival

Source: Gartner (June 2018)

People and Skills

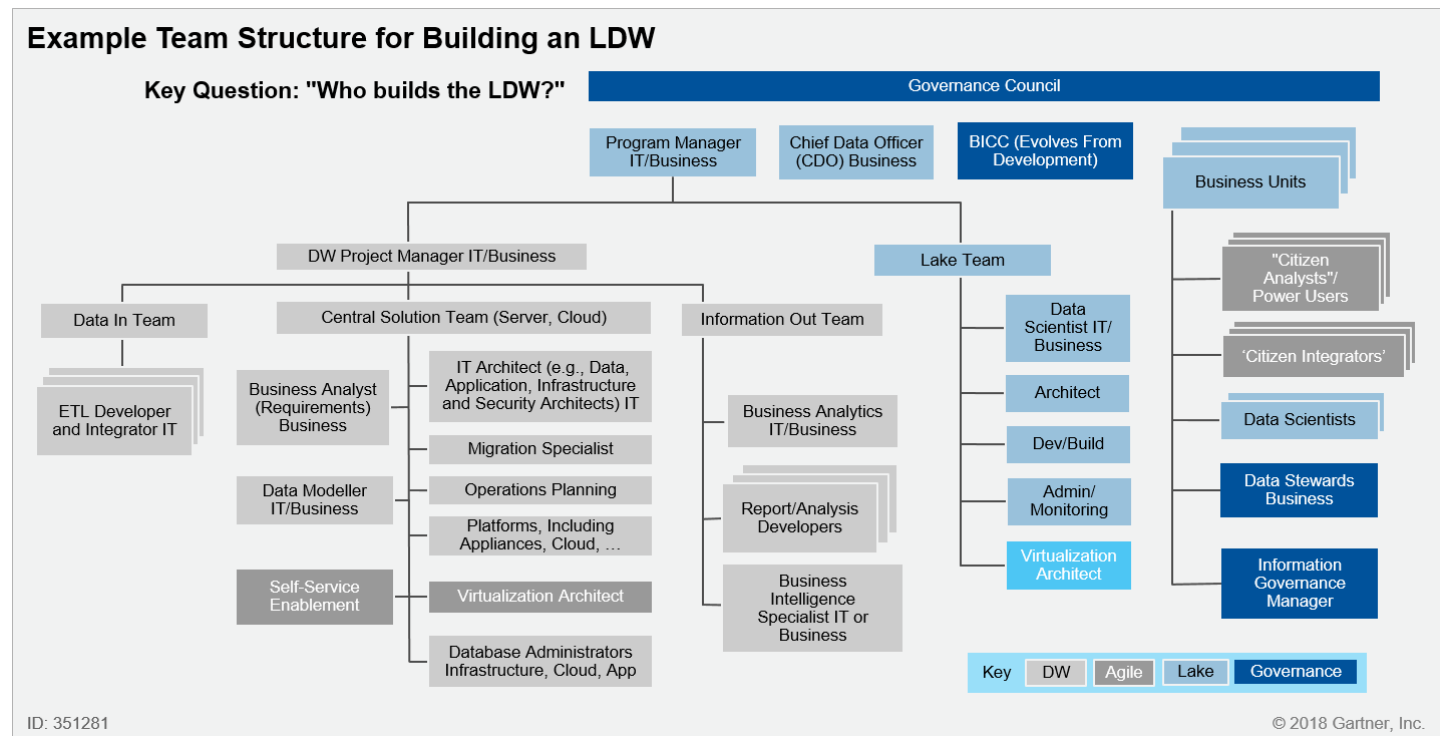
Identify roles and responsibilities along with skills needed for each new initiative. The following roles can be helpful in managing a modern data and analytics architecture.

- **Business user:** Someone who understands the domain area and usually benefits from the results. This person can consult and advise the project team on the context of the project, the value of the results and how the outputs will be operationalized. Usually, a business analyst, line manager or deep subject matter expert in the project domain fulfills this role.
- **Project sponsor:** Responsible for the genesis of the project. Provides the impetus and requirements for the project and defines the core business problem. Generally provides the funding and gauges the degree of value from the final outputs of the working team. This person sets the priorities for the project and clarifies the desired outputs.
- **Project manager:** Ensures that key milestones and objectives are met on time and at the expected quality.
- **Business intelligence analyst:** Provides business domain expertise based on a deep understanding of the data, key performance indicators (KPIs), key metrics and business intelligence from a reporting perspective. Business intelligence analysts generally create dashboards and reports, and have knowledge of the data feeds and sources.
- **Database administrator (DBA):** Provisions and configures the database environment to support the analytics needs of the working team. These responsibilities may include providing access to key databases or tables and ensuring the appropriate security levels are in place related to the data repositories.
- **Data engineer:** Leverages deep technical skills to assist with tuning SQL queries for data management and data extraction, and provides support for data ingestion into the analytic sandbox. Whereas the DBA sets up and configures the databases to be used, the data engineer executes the actual data extractions and performs substantial data manipulation to facilitate the analytics. The data engineer works closely with the data scientist to help shape data in the right ways for analyses, using a variety of tools and languages (e.g., Python, Java, JavaScript, ETL, etc.)
- **Data scientist:** Provides subject matter expertise for analytical techniques, data modeling and applying valid analytical techniques to given business problems. Ensures overall analytics objectives are met. Designs and executes analytical methods and approaches with the data available to the project.

A very useful foundation for a data and analytics system is the LDW. It provides sources of structured and unstructured data for analysis. It has associated with it a variety of tools to transform and quality-assure data. It also can contain simple mechanisms like the operational data store (ODS) to publish results to transactional systems. Figure 23 shows a sample team structure for building an LDW. Documenting the

roles and responsibilities in this way can be a useful reference when designing a data and analytics strategy.

Figure 23. Example Team Structure for Building an LDW



Source: Gartner (June 2018)

Account for Citizen Roles

Account for citizen roles performed by business users. These new, hybrid roles often span functions and departments, blending business and IT functions and blurring traditional organizational boundaries. Organizations must plan for the emergence of these hybrid roles and determine how the information management and analytics strategy will support these new functions.

With each new initiative, it's important to consider cross-functional boundaries and identify where IT's responsibilities begin and end. Citizen roles can't be ignored, so IT should provide tools, processes and support to enable a strong partnership with business stakeholders.

Agile Database Development

Agile and DevOps require application and data developers to adopt iterative and incremental design and implementation processes. To achieve continuous delivery, technical professionals focused on agile development must apply these same processes to their application database changes.

To effectively support agile development, the methods used to design, implement and evolve databases must follow the same incremental, evolutionary and automated approaches used to develop the application code. Techniques such as schema on read and the dynamic transformation of data enabled by in memory processing make data structures much more malleable than they have been in the past.

Agile developers, under the guidance and support of the DBA, can remove or overcome constraints on velocity by developing the knowledge and skills to:

- Create, maintain and incrementally evolve the data model as the project advances.
- Write, review, modify, optimize and validate changes to database interface code, SQL scripts, database schemas and functional objects, using agile technical practices.
- Manage, maintain and version the test data used for validation.
- Implement automation as part of the continuous integration and DevOps pipelines to minimize the amount of manual effort involved in supporting the tasks listed above.

Gartner has published guidance to help organizations through this process (see ["Implement Agile Database Development to Support Your Continuous Delivery Initiative"](https://www.gartner.com/document/code/310316?ref=grbody&refval=3880568) (<https://www.gartner.com/document/code/310316?ref=grbody&refval=3880568>)).

Cloud Versus On-Premises

With each new initiative, technical professionals should consider important technology placement and deployment considerations, such as:

- On-premises versus cloud
- Cloud service providers
- Multicloud and hybrid strategies
- Platform as a service (PaaS) and infrastructure as a service (IaaS)

Key factors in the decision might include:

- **Reducing or controlling costs:** To optimize business and reduce costs (such as costs associated with capacity upgrades or infrastructure operations), organizations opt for external IT sourcing alternatives, such as colocation, managed hosting and cloud.
- **Improving speed of delivery:** Organizations leverage technologies, such as virtualization and orchestration, to improve the speed of delivery and to adapt more easily to the changing business

demands.

- **Focusing on core competencies:** Organizations want to refocus on core competencies or capabilities that provide competitive differentiation to the business, and leave more commoditized services to external providers.
- **Overcoming skills shortages:** Building and operating a data center requires a complete range of skills — including people, processes and technology — that many organizations find difficult to maintain. External hosting models require a subset of the skills required by the data center; for example, facilities management skills are no longer required. In the case of public cloud, required skills include vendor management and business strategy.
- **Improving service levels:** Many corporate-owned data centers are overcrowded and lack availability features, such as backup generators or multiple internet connections. Moving applications from these data centers to more robust facilities can improve service levels and availability (see Figure 24).

Figure 24. Cloud Deployment Options for Data

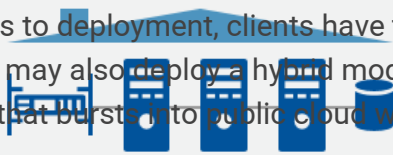
Cloud Deployment Options for Data

On-Premises

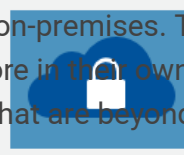
Client's Own Data Center

Source: Gartner (June 2018)

When it comes to deployment, clients have two basic choices: cloud or on-premises. Technical professionals may also deploy a hybrid model consisting of the data store in their own data center or in a private cloud that bursts into public cloud when it needs to meet loads that are beyond its design.



Private Cloud



As technical professionals migrate on-premises systems to the public cloud, they need to know which databases to migrate, which tools and techniques to use, and how to run a successful migration project. This is the traditional option, but it can tie businesses down to processes and maintaining infrastructure, and paying the associated costs.

Gartner has published guidance to help organizations through this process (see [Migrating Enterprise Databases and Data to the Cloud](https://www.gartner.com/document/code/317167?ref=ggrec&refval=3880568) (<https://www.gartner.com/document/code/317167?ref=ggrec&refval=3880568>)). Unlike public cloud, the infrastructure is dedicated to a single organization. This can help meet certain needs such as strong security control. Private clouds can also be hosted at an external provider. These are best-suited for mission-critical applications that have very high security and uptime requirements.

Document Revision History

Solution Path for Planning and Implementing a Data and Analytics Architecture - 16 June 2017

(<https://www.gartner.com/document/code/324344?ref=ddrec>)

Recommended by the Author

EIM 1.0: Setting Up Enterprise Information Management and Governance

(<https://www.gartner.com/document/code/342309?ref=ggrec&refval=3880568>)

Enabling Streaming Architectures for Continuous Data and Events With Kafka

(<https://www.gartner.com/document/code/353112?ref=ggrec&refval=3880568>)

Enabling Essential Data Governance for Successful Big Data Architecture Deployment

(<https://www.gartner.com/document/code/327532?ref=ggrec&refval=3880568>)

Adopt the Logical Data Warehouse Architecture to Meet Your Modern Analytical Needs

(<https://www.gartner.com/document/code/342254?ref=ggrec&refval=3880568>)

Applying Effective Data Governance to Secure Your Data Lake

(<https://www.gartner.com/document/code/346975?ref=ggrec&refval=3880568>)

Solution Path for Planning and Implementing the Logical Data Warehouse

(<https://www.gartner.com/document/code/320563?ref=ggrec&refval=3880568>)

Implement Agile Database Development to Support Your Continuous Delivery Initiative

(<https://www.gartner.com/document/code/310316?ref=ggrec&refval=3880568>)

Create a Data Reference Architecture to Enable Self-Service BI

(<https://www.gartner.com/document/code/333398?ref=ggrec&refval=3880568>)

Migrating Enterprise Databases and Data to the Cloud (<https://www.gartner.com/document/code/317167?ref=ggrec&refval=3880568>)

Making Machine Learning a Scalable Enterprise Reality – From Development to Production

(<https://www.gartner.com/document/code/343614?ref=ggrec&refval=3880568>)

Recommended For You

Use Design Patterns to Increase the Value of Your Data Lake

(<https://www.gartner.com/document/3876783?ref=ddrec&refval=3880568>)

© 2018 Gartner, Inc. and/or its affiliates. All rights reserved. Gartner is a registered trademark of Gartner, Inc. or its affiliates. This publication may not be reproduced or distributed in any form without Gartner's prior written permission. If you are authorized to access this publication, your use of it is subject to the [Gartner Usage Policy](#) posted on gartner.com. The information contained in this publication has been obtained from sources believed to be reliable. Gartner disclaims all warranties as to the accuracy, completeness or adequacy of such information and shall have no liability for errors, omissions or inadequacies in such information. This publication consists of the opinions of Gartner's research organization and should not be construed as statements of fact. The opinions expressed herein are subject to change without notice. Although Gartner research may include a discussion of related legal issues, Gartner does not provide legal advice or services and its research should not be construed or used as such. Gartner is a public company, and its shareholders may include firms and funds that have financial interests in entities covered in Gartner research. Gartner's Board of Directors may include senior managers of these firms or funds. Gartner research is produced independently by its research organization without input or influence from these firms, funds or their managers. For further information on the independence and integrity of Gartner research, see "[Guiding Principles on Independence and Objectivity](#)."