



HAL
open science

On data lake architectures and metadata management

Pegdwendé Nicolas Sawadogo, Jérôme Darmont

► **To cite this version:**

Pegdwendé Nicolas Sawadogo, Jérôme Darmont. On data lake architectures and metadata management. *Journal of Intelligent Information Systems*, 2021, 56 (1), pp.97-120. 10.1007/s10844-020-00608-7. hal-03114365

HAL Id: hal-03114365

<https://hal.science/hal-03114365>

Submitted on 22 Jul 2021

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

On Data Lake Architectures and Metadata Management

Pegdwendé Sawadogo · Jérôme Darmont

Received: date / Accepted: date

Abstract Over the past two decades, we have witnessed an exponential increase of data production in the world. So-called big data generally come from transactional systems, and even more so from the Internet of Things and social media. They are mainly characterized by volume, velocity, variety and veracity issues. Big data-related issues strongly challenge traditional data management and analysis systems. The concept of data lake was introduced to address them. A data lake is a large, raw data repository that stores and manages all company data bearing any format. However, the data lake concept remains ambiguous or fuzzy for many researchers and practitioners, who often confuse it with the Hadoop technology. Thus, we provide in this paper a comprehensive state of the art of the different approaches to data lake design. We particularly focus on data lake architectures and metadata management, which are key issues in successful data lakes. We also discuss the pros and cons of data lakes and their design alternatives.

Keywords Data lakes · Data lake architectures · Metadata management · Metadata modeling

Acknowledgement

The research accounted for in this paper was funded by the Université Lumière Lyon 2 and the Auvergne-Rhône-Alpes Region through the COREL and AURA-PMI projects, respectively. The authors also sincerely thank the anonymous reviewers of this paper for their constructive comments and suggestions.

Pegdwendé Sawadogo
Université de Lyon, Lyon 2, ERIC UR 3083
E-mail: pegdwende.sawadogo@univ-lyon2.fr

Jérôme Darmont
Université de Lyon, Lyon 2, ERIC UR 3083
E-mail: jerome.darmont@univ-lyon2.fr

1 Introduction

The 21st century is marked by an exponential growth of the amount of data produced in the world. This is notably induced by the fast development of the Internet of Things (IoT) and social media. Yet, while big data represent a tremendous opportunity for various organizations, they come in such volume, speed, heterogeneous sources and structures that they exceed the capabilities of traditional management systems for their collection, storage and processing in a reasonable time [37]. A time-tested solution for big data management and processing is data warehousing. A data warehouse is indeed an integrated and historical storage system that is specifically designed to analyze data. However, while data warehouses are still relevant and very powerful for structured data, semi-structured and unstructured data induce great challenges for data warehouses. Yet, the majority of big data is unstructured [37]. Thus, the concept of data lake was introduced to address big data issues, especially those induced by data variety.

A data lake is a very large data storage, management and analysis system that handles any data format. It is currently quite popular and trendy both in the industry and academia. Yet, the concept of data lake is not straightforward for everybody. A survey conducted in 2016 indeed revealed that 35% of the respondents considered data lakes as a simple marketing label for a preexisting technology, i.e., Apache Hadoop [17]. Knowledge about the concept of the data lake has since evolved, but some misconceptions still exist, presumably because most of data lakes design approaches are abstract sketches from the industry that provide few theoretical or implementation details [43]. Therefore, a survey can be useful to give researchers and practitioners a better comprehension of the data lake concept and its design alternatives.

To the best of our knowledge, the only literature reviews about data lakes are all quite brief and/or focused on a specific topic, e.g., data lake concepts and definitions [8, 33], the technologies used for implementing data lakes [35] or data lakes inherent issues [16, 43]. Admittedly, the report proposed in [47] is quite extensive, but it adopts a purely industrial view. Thus, we adopt in this paper a wider scope to propose a more comprehensive state of the art of the different approaches to design and exploit a data lake. We particularly focus on data lake architectures and metadata management, which lie at the base of any data lake project and are the most commonly cited issues in the literature (Figure 1).

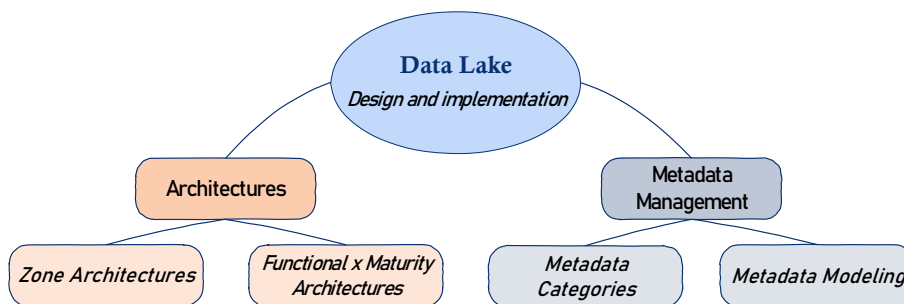


Fig. 1 Issues addressed in this paper

More precisely, we first review data lake definitions and complement the best existing one. Then, we investigate the architectures and technologies used for the implementation of data lakes, and propose a new typology of data lake architectures. Our second main focus is metadata management, which is a primordial issue to avoid turning a data lake into an inoperable, so-called data swamp. We notably classify data lake metadata and introduce the features that are necessary to achieve a full metadata system. We also discuss the pros and cons of data lakes.

Eventually, note that we do not review other important topics, such as data ingestion, data governance and security in data lakes, because they are currently little addressed in the literature, but could still presumably be the subject of another full survey.

The remainder of this paper is organized as follows. In Section 2, we define the data lake concept. In Section 3, we review data lake architectures and technologies to help users choose the right approach and tools. In Section 4, we extensively review and discuss metadata management. Eventually, we recapitulate the pros and cons of data lakes in Section 5 and conclude the paper in Section 6 with a mind map of the key concepts we introduce, as well as current open research issues.

2 Data Lake Definitions

2.1 Definitions from the Literature

The concept of data lake was introduced by Dixon as a solution to perceived shortcomings of datamarts, which are business-specific subdivisions of data warehouses that allow only subsets of questions to be answered [10]. In the literature, data lakes are also referred to as data reservoirs [7] and data hubs [15,29], although the terms data lake are the most frequent. Dixon envisions a data lake as a large storage system for raw, heterogeneous data, fed by multiple data sources, and that allows users to explore, extract and analyze the data.

Subsequently, part of the literature considered data lakes as an equivalent to the Hadoop technology [11,15,40]. According to this point of view, the concept of data lake refers to a methodology for using free or low-cost technologies, typically Hadoop, for storing, processing and exploring raw data within a company [11]. The systematic association of data lakes to low cost technologies is becoming minority in the literature, as the data lake concept is now also associated with proprietary cloud solutions such as Azure or IBM [33,51] and various data management systems such as NoSQL solutions and multistores. However, it can still be viewed as a data-driven design pattern for data management [47].

More consensually, a data lake may be viewed as a central repository where data of all formats are stored without a strict schema, for future analyses [8,26,35]. This definition is based on two key characteristics of data lakes: data variety and the schema-on-read approach, also known as late binding [11], which implies that schema and data requirements are not fixed until data querying [26,32,53]. This is the opposite to the schema-on-write approach used in data warehouses.

However, the variety/schema-on-read definition may be considered fuzzy because it gives little detail about the characteristics of a data lake. Thus, Madera and Laurent introduce a more complete definition where a data lake is a logi-

cal view of all data sources and datasets in their raw format, accessible by data scientists or statisticians for knowledge extraction [33].

More interestingly, this definition is complemented by a set of features that a data lake should include:

1. data quality is provided by a set of metadata;
2. the lake is controlled by data governance policy tools;
3. usage of the lake is limited to statisticians and data scientists;
4. the lake integrates data of all types and formats;
5. the data lake has a logical and physical organization.

2.2 Discussion and New Definition

Madera and Laurent’s definition of data lakes is presumably the most precise, as it defines the requirements that a data lake must meet (Section 2.1). However, some points in this definition are debatable. The authors indeed restrain the use of the lake to data specialists and, as a consequence, exclude business experts for security reasons. Yet, in our opinion, it is entirely possible to allow controlled access to this type of users through a navigation or analysis software layer.

Moreover, we do not share the vision of the data lake as a logical view over data sources, since some data sources may be external to an organization, and therefore to the data lake. Since Dixon explicitly states that lake data come from data sources [10], including data sources into the lake may therefore be considered contrary to the spirit of data lakes.

Finally, although quite complete, Madera and Laurent’s definition omits an essential property of data lakes: scalability [26,37]. Since a data lake is intended for big data storage and processing, it is indeed essential to address this issue. Thence, we amend Madera and Laurent’s definition to bring it in line with our vision and introduce scalability [49].

Definition 1 A data lake is a scalable storage and analysis system for data of any type, retained in their native format and used *mainly* by data specialists (statisticians, data scientists or analysts) for knowledge extraction. Its characteristics include:

1. a metadata catalog that enforces data quality;
2. data governance policies and tools;
3. accessibility to various kinds of users;
4. integration of any type of data;
5. a logical and physical organization;
6. scalability in terms of storage and processing.

3 Data Lake Architectures and Technologies

Existing reviews on data lake architectures commonly distinguish pond and zone architectures [16,45]. However, this categorization may sometimes be fuzzy. Thus, we introduce in Section 3.1 a new manner to classify data lakes architectures that we call Functional \times Maturity. In addition, we present in Section 3.2 a list

of possible technologies to implement a data lake. Eventually, we investigate in Section 3.3 how a data lake system can be associated with a data warehouse in an enterprise data architecture.

3.1 Data Lake Architectures

3.1.1 Zone Architectures

Pond architecture Inmon designs a data lake as a set of data ponds [23]. A data pond can be viewed as a subdivision of a data lake dealing with data of a specific type. According to Dixon's specifications, each data pond is associated with a specialized storage system, some specific data processing and conditioning (i.e., data transformation/preparation) and a relevant analysis service. More precisely, Inmon identifies five data ponds (Figure 2).

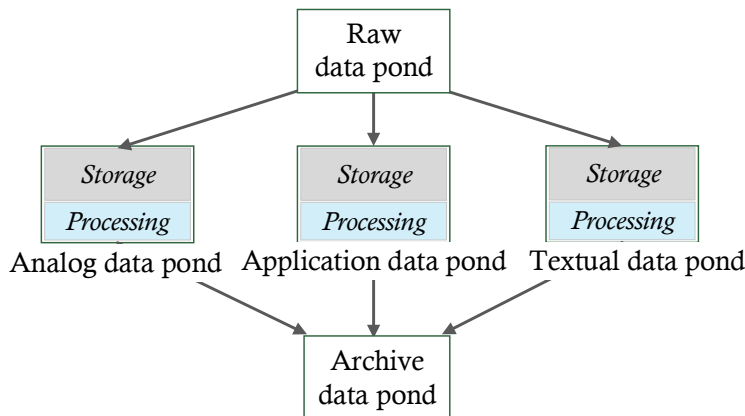


Fig. 2 Data flow in a pond architecture

1. The **raw data pond** deals with newly ingested, raw data. It is actually a transit zone, since data are then conditioned and transferred into another data pond, i.e., either the analog, application or textual data pond. The raw data pond, unlike the other ponds, is not associated with any metadata system.
2. Data stored in the **analog data pond** are characterized by a very high frequency of measurements, i.e., they come in with high velocity. Typically, semi-structured data from the IoT are processed in the analog data pond.
3. Data ingested in the **application data pond** come from software applications, and are thus generally structured data from relational Database Management Systems (DBMSs). Such data are integrated, transformed and prepared for analysis; and Inmon actually considers that the application data pond is a data warehouse.
4. The **textual data pond** manages unstructured, textual data. It features a textual disambiguation process to ease textual data analysis.

5. The purpose of the **archival data pond** is to save the data that are not actively used, but might still be needed in the future. Archived data may originate from the analog, application and textual data ponds.

Zone architectures So-called zone architectures assign data to a zone according to their degree of refinement [16]. For instance, Zaloni’s data lake [28] adopts a six-zone architecture (Figure 3).

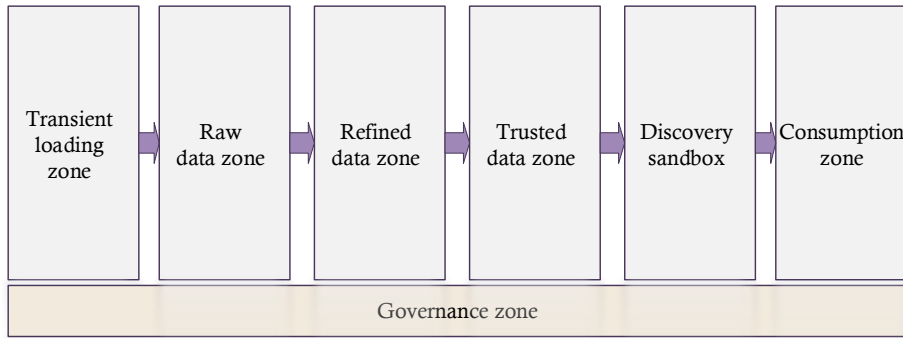


Fig. 3 Zaloni’s zone architecture [28]

1. The **transient loading zone** deals with data under ingestion. Here, basic data quality checks are performed.
2. The **raw data zone** handles data in near raw format coming from the transient zone.
3. The **trusted zone** is where data are transferred once standardized and cleansed.
4. From the trusted area, data move into the **discovery sandbox** where they can be accessed by data scientists through data wrangling or data discovery operations.
5. On top of the discovery sandbox, the **consumption zone** allows business users to run “what if” scenarios through dashboard tools.
6. The **governance zone** finally allows to manage, monitor and govern metadata, data quality, a data catalog and security.

However, this is but one of several variants of zone architectures. Such architectures indeed generally differ in the number and characteristics of zones [16], e.g., some architectures include a transient zone [28, 56, 59] while others do not [18, 45].

A particular zone architecture often mentioned in the data lake literature is the lambda architecture [24, 35]. It indeed stands out since it includes two data processing zones: a batch processing zone for bulk data and a real-time processing zone for fast data from the IoT [24]. These two zones help handling fast data as well as bulk data in an adapted and specialized way.

Discussion In both pond and zone architectures, data are pre-processed. Thus, analyses are quick and easy. However, this come at the cost of data loss in the pond architectures, since raw data are deleted when transferred to other ponds.

The drawbacks of the many zone architectures depend on each variant. For example, in Zaloni’s architecture [28], data flow across six areas, which may lead to multiple copies of the data and, therefore, difficulties in controlling data lineage. In the Lamda architecture [24], speed and batch processing components follow different paradigms. Thus, data scientists must handle two distinct logics for cross analyses [35], which makes data analysis harder, overall.

Moreover, the distinction of data lake architectures into pond and zone approaches is not so crisp in our opinion. The pond architecture may indeed be considered as a variant of zone architecture, since data location depends on the refinement level of data, as in zone architectures. In addition, some zone architectures include a global storage zone where raw and cleansed data are stored altogether [24, 43], which contradicts the definition of zone architectures, i.e., components depend on the degree of data refinement.

3.1.2 Functional \times Maturity Architectures

To overcome the contradictions of the pond/zone categorization, we propose an alternative way to group data lake architectures regarding the type of criteria used to define components. As a result, we distinguish functional architectures, data maturity-based architectures and hybrid architectures (Figure 4).

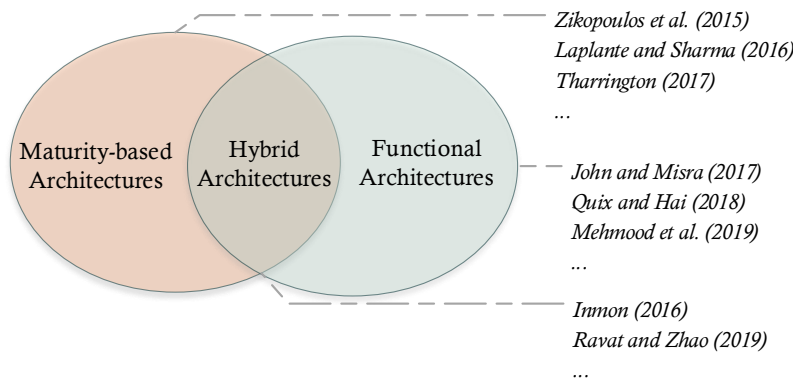


Fig. 4 Architecture typology proposal

Functional architectures follow some basic functions to define a lake’s components. Data lake basic functions typically include [28]:

1. a data ingestion function to connect with data sources;
2. a data storage function to persist raw as well as refined data;
3. a data processing function;
4. a data access function to allow raw and refined data querying.

Quix and Hai, as well as Mehmood et al., base their data lake architectures on these functions [36, 43]. Similarly, John and Misra’s lambda architecture [24] may be considered as a functional architecture, since its components represent data lake functions such as storage, processing and serving.

Data maturity-based architectures are data lake architectures where components are defined regarding data refinement level. In other words, it is constituted of most zone architectures. A good representative is Zaloni’s data lake architecture [28], where common basic zones are a transient zone, a raw data zone, a trusted data zone and a refined data zone [28, 56, 59].

Hybrid architectures are data lake architectures where the identified components depend on both data lake functions and data refinement. Inmon’s pond architecture is actually a hybrid architecture [23]. On one hand, it is a data maturity-based architecture, since raw data are managed in a special component, i.e., the raw data pond, while refined data are managed in other ponds, i.e., the textual, analog and application data ponds. But on the other hand, the pond architecture is also functional because Inmon’s specifications consider some storage and process components distributed across data ponds (Figure 2). Ravat and Zhao also propose such an hybrid data lake architecture (Figure 5 [45]).

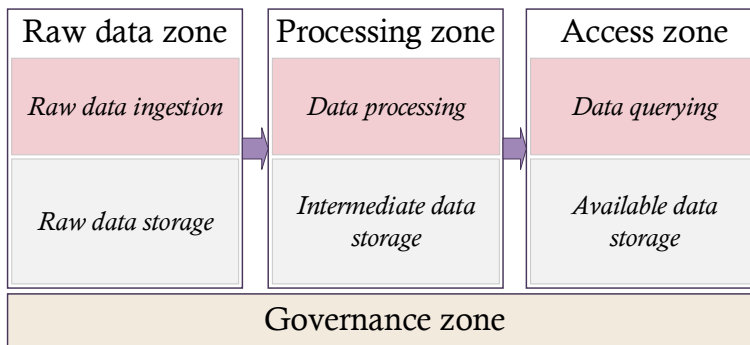


Fig. 5 Ravat and Zhao’s hybrid architecture [45]

Discussion Functional architectures have the advantage of clearly highlighting the functions to implement for a given data lake, which helps match easily with the required technologies. By contrast, data maturity-based architectures are useful to plan and organize the data lifecycle. Both approaches are thus limited, since they only focus on a unique point of view, while it is important in our opinion to take both functionality and data maturity into account when designing a data lake.

In consequence, we advocate for hybrid approaches. However, existing hybrid architecture can still be improved. For instance, in Inmon’s data pond approach, raw data are deleted once they are refined. This process may induce some data loss, which is contrary to the spirit of data lakes. In Ravat and Zhao’s proposal, data access seems only possible for refined data. Such limitations hint that a more complete hybrid data lake architecture is still needed nowadays.

3.2 Technologies for Data Lakes

Most data lake implementations are based on the Apache Hadoop ecosystem [8, 26]. Hadoop has indeed the advantage of providing both storage with the Hadoop

Distributed File System (HDFS) and data processing tools via MapReduce or Spark. However, Hadoop is not the only suitable technology to implement a data lake. In this section, we go beyond Hadoop to review usable tools to implement data lake basic functions.

3.2.1 Data Ingestion

Ingestion technologies help physically transfer data from data sources into a data lake. A first category of tools includes software that iteratively collects data through pre-designed and industrialized jobs. Most such tools are proposed by the Apache Foundation, and can also serve to aggregate, convert and clean data before ingestion. They include Flink and Samza (distributed stream processing frameworks), Flume (a Hadoop log transfer service), Kafka (a framework providing real time data pipelines and stream processing applications) and Sqoop (a framework for data integration from SQL and NoSQL DBMSs into Hadoop) [24,35,54].

A second category of data ingestion technologies is made of common data transfer tools and protocols (wget, rsync, FTP, HTTP, etc.), which are used by the data lake manager within data ingestion scripts. They have the key advantage to be readily available and widely understood [55]. In a similar way, some Application Programming Interfaces (APIs) are available for data retrieval and transfer from the Web into a data lake. For instance, CKAN and Socrata provide APIs to access a catalogue of open data and associated metadata [55].

3.2.2 Data Storage

We distinguish two main approaches to store data in data lakes. The first way consists in using classic databases for storage. Some data lakes indeed use relational DBMSs such as MySQL, PostgreSQL or Oracle to store structured data [3, 26]. However, relational DBMSs are ill-adapted to semi-structured, and even more so to unstructured data. Thus, NoSQL (Not only SQL) DBMSs are usually used instead [3,16,26]. Moreover, assuming that data variety is the norm in data lakes, a multi-paradigm storage system is particularly relevant [39]. Such so-called multistore systems manage multiple DBMSs, each matching a specific storage need.

The second main way to store data and the most used is HDFS storage (in about 75% of data lakes [47]). HDFS is a distributed storage system that offers a very high scalability and handles all types of data [24]. Thus, it is well suited for schema-free and bulk storage that are needed for unstructured data. Another advantage of this technology is the distribution of data that allows high fault-tolerance. However, HDFS alone is not sufficient to handle all data formats, especially structured data. Thus, it should ideally be combined with relational and/or NoSQL DBMSs.

3.2.3 Data Processing

In data lakes, data processing is very often performed with MapReduce [8,24, 26,35,53,54], a parallel data processing paradigm provided by Apache Hadoop. MapReduce is well-suited to very large data, but is less efficient for fast data because it works on disk [58]. Thus, alternative processing frameworks are used,

from which the most famous is Apache Spark. Spark works like MapReduce, but adopts a full in-memory approach instead of using the file system for storing intermediate results. Thence, Spark is particularly suitable for real-time processing. Similarly, Apache Flink and Apache Storm are also suitable for real-time data processing [24, 26, 35, 54, 58]. Nevertheless, these two approaches can be simultaneously implemented in a data lake, with MapReduce being dedicated to voluminous data and stream-processing engines to velocious data [24, 54].

3.2.4 Data Access

In data lakes, data may be accessed through classical query languages such as SQL for relational DBMSs, JSONiq for MongoDB, XQuery for XML DBMSs or SPARQL for RDF resources [12, 14, 18, 29, 42]. However, this does not allow simultaneously querying across heterogeneous databases, while data lakes do store heterogeneous data, and thus typically require heterogeneous storage systems.

One solution to this issue is to adopt the query techniques from multistores (Section 3.2.2) [39]. For example, Spark SQL and SQL++ may be used to query both relational DBMSs and semi-structured data in JSON format. In addition, the Scalable Query Rewriting Engine (SQRE) handles graph databases [19]. Finally, CloudMdsQL also helps simultaneously query multiple relational and NoSQL DBMSs [30]. Quite similarly, Apache Phoenix can be used to automatically convert SQL queries into a NoSQL query language, for example. Apache Drill allows joining data from multiple storage systems [3]. Data stored in HDFS can also be accessed using Apache Pig [24].

Eventually, business users, who require interactive and user-friendly tools for data reporting and visualization tasks, widely use dashboard services such as Microsoft Power BI and Tableau over data lakes [8, 47].

3.3 Combining Data Lakes and Data Warehouses

There are in the literature two main approaches to combine a data lake and a data warehouse in a global data management system. The first approach pertains to using a data lake as the data source of a data warehouse (Section 3.3.1). The second considers data warehouses as components of data lakes (Section 3.3.2).

3.3.1 Data Lake Sourcing a Data Warehouse

This approach aims to take advantage of the specific characteristics of both data lakes and data warehouses. Since data lakes allow an easier and cheaper storage of large amount of raw data, they can be considered as staging areas or Operational Data Stores (ODSs) [11, 47], i.e., intermediary data stores ahead of data warehouses that gather operational data from several sources before the ETL process takes place.

With a data lake sourcing a data warehouse, possibly with semi-structured data, industrialized OLAP analyses are possible over the lake's data, while on-demand, ad-hoc analyses are still possible directly from the data lake (Figure 6).

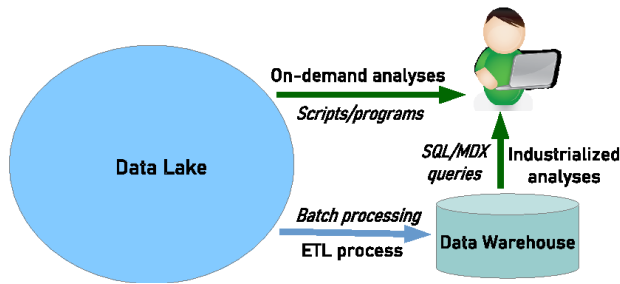


Fig. 6 Data lake and data warehouse architecture

3.3.2 Data Warehouse within a Data Lake

As detailed in Section 3.1.1, Inmon proposes an architecture based on a subdivision of data lakes into so-called data ponds [23]. For Inmon, structured data ponds sourced from operational applications are, plain and simple, data warehouses. Thus, this approach acts on a conception of data lakes as extensions of data warehouses.

3.3.3 Discussion

When a data lake sources a data warehouse (Section 3.3.1), there is a clear functional separation, as data warehouses and data lakes are specialized in industrialized and on-demand analyses, respectively. However, this comes with a data siloing issue.

By contrast, the data siloing syndrome can be reduced in Inmon’s approach (Section 3.3.2), as all data are managed and processed in a unique global platform. Hence, diverse data can easily be combined through cross-reference analyses, which would be impossible if data were managed separately. In addition, building a data warehouse inside a global data lake may improve data lifecycle control. That is, it should be easier to track, and thus to reproduce processes applied to the data that are ingested in the data warehouse, via the data lake’s tracking system.

4 Metadata Management in Data Lakes

Data ingested in data lakes bear no explicit schema [37], which can easily turn a data lake into a data swamp in the absence of an efficient metadata system [54]. Thence, metadata management plays an essential role in data lakes [29,26]. In this section, we detail the metadata management techniques used in data lakes. First, we identify the metadata that are relevant to data lakes. Then, we review how metadata can be organized. We also investigate metadata extraction tools and techniques. Finally, we provide an inventory of desirable features in metadata systems.

4.1 Metadata Categories

We identify in the literature two main typologies of metadata dedicated to data lakes. The first one distinguishes functional metadata, while the second classifies metadata with respect to structural metadata types.

4.1.1 Functional Metadata

Oram introduces a metadata classification in three categories, with respect to the way they are gathered [41].

1. **Business metadata** are defined as the set of descriptions that make the data more understandable and define business rules. More concretely, these are typically data field names and integrity constraints. Such metadata are usually defined by business users at the data ingestion stage.
2. **Operational metadata** are information automatically generated during data processing. They include descriptions of the source and target data, e.g., data location, file size, number of records, etc., as well as process information.
3. **Technical metadata** express how data are represented, including data format (e.g., raw text, JPEG image, JSON document, etc.), structure or schema. The data structure consists in characteristics such as names, types, lengths, etc. They are commonly obtained from a DBMS for structured data, or via custom techniques during the data maturation stage.

Diamantini et al. enhance this typology with a generic metadata model [9] and show that business, operational and technical metadata sometimes intersect. For instance, data fields relate both to business and technical metadata, since they are defined in data schemas by business users. Similarly, data formats may be considered as both technical and operational metadata, and so on (Figure 7).

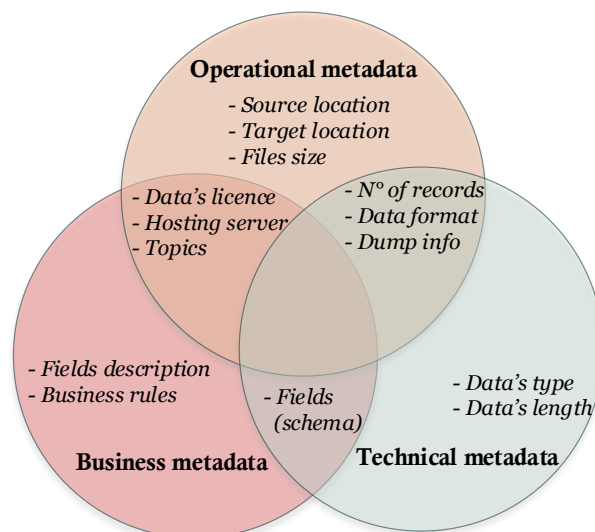


Fig. 7 Functional metadata model [9]

4.1.2 Structural Metadata

In this classification, Sawadogo et al. categorize metadata with respect to the “objects” they relate to [49]. The notion of object may be viewed as a generalization of the dataset concept [32], i.e., an object may be a relational or spreadsheet table in a structured data context, or a simple document (e.g., XML document, image file, video file, textual document, etc.) in a semi-structured or unstructured data context. Thence, we use the term “object” in the remainder of this paper.

Intra-object metadata belong to a set of characteristics associated with single objects in the lake. They are subdivided into four main subcategories.

1. **Properties** provide an object’s general description. They are generally retrieved from the filesystem as key-value pairs, e.g., file name and size, location, date of last modification, etc.
2. **Previsualization and summary metadata** aim to provide an overview of the content or structure of an object. For instance, metadata can be extracted data schemas for structured and semi-structured data, or wordclouds for textual data.
3. **Version and representation metadata** are made of altered data. When a new data object o' is generated from existing object o in the data lake, o' may be considered as metadata for o . Version metadata are obtained through data updates, while representation metadata come from data refining operations. For instance, a refining operation may consist of vectorizing a textual document into a bag-of-words for further automatic processing.
4. **Semantic metadata** involve annotations that describe the meaning of data in an object. They include such information as title, description, categorization, descriptive tags, etc. They often allow data linking. Semantic metadata can be either generated using semantic resources such as ontologies, or manually added by business users [18, 44].

Inter-object metadata represent links between two or more objects. They are subdivided into three categories.

1. **Object groupings** organize objects into collections. Any object may be associated with several collections. Such links can be automatically deduced from some intra-object metadata such as tags, data format, language, owner, etc.
2. **Similarity links** express the strength of likeness between objects. They are obtained via common or custom similarity measures. For instance, [32] define the affinity and joinability measures to express the similarity between semi-structured objects.
3. **Parenthood links** aim to save data lineage, i.e., when a new object is created from the combination of several others, these metadata record the process. Parenthood links are thus automatically generated during data joins.

Global metadata are data structures that provide a context layer to make data processing and analysis easier. Global metadata are not directly associated with any specific object, but potentially concern the entire lake. There are three subcategories of global metadata.

1. **Semantic resources** are knowledge bases such as ontologies, taxonomies, thesauri, etc., which notably help enhance analyses. For instance, an ontology can be used to automatically extend a term-based query with equivalent terms. Semantic resources are generally obtained from the Internet or manually built.
2. **Indexes** (including inverted indexes) enhance term-based or pattern-based data retrieval. They are automatically built and enriched by an indexing system.
3. **Logs** track user interactions with the data lake, which can be simple, e.g., user connection or disconnection, or more complex, e.g., a job running.

4.1.3 Discussion

Oram’s metadata classification is the most cited, especially in the industrial literature [9, 28, 46, 47], presumably because it is inspired from metadata categories from data warehouses [45]. Thus, its adoption seems easier and more natural for practitioners who are already working with it.

Yet, we favor the second metadata classification, because it includes most of the features defined by Oram’s. Business metadata are indeed comparable to semantic metadata. Operational metadata may be considered as logs and technical metadata are equivalent to previsualization metadata. Hence, the structural metadata categorization can be considered as an extension, as well as a generalization, of the functional metadata classification.

Moreover, Oram’s classification is quite fuzzy when applied in the context of data lakes. Diamantini et al. indeed show that functional metadata intersect (Section 4.1.1) [9]. Therefore, practitioners who do not know this typology may be confused when using it to identify and organize metadata in a data lake.

Table 1 summarizes commonalities and differences between the two metadata categorizations presented above. The comparison addresses the type of information both inventories provide.

Table 1 Comparison of Oram’s and Sawadogo et al’s metadata categories

| Type of information | Functional metadata | Structural metadata |
|--|---------------------|---------------------|
| Basic characteristics of data (size, format, etc.) | ✓ | ✓ |
| Data semantics (tags, descriptions, etc.) | ✓ | ✓ |
| Data history | ✓ | ✓ |
| Data linkage | | ✓ |
| User interactions | | ✓ |

4.2 Metadata Modeling

There are in the literature two main approaches to represent a data lake’s metadata system. The first, most common approach, adopts a graph view, while the second exploits data vault modeling.

4.2.1 Graph Models

Most models that manage data lake metadata systems are based on a graph approach. We identify three main subcategories of graph-based metadata models with respect to the main features they target.

Data provenance-centered graph models mostly manage metadata tracing, i.e., the information about activities, data objects and users who interact with a specific object [54]. In other words, they track the pedigree of data objects [20]. Provenance representations are usually built using a directed acyclic graph (DAG) where nodes represent entities such as users, roles or objects [3,21]. Edges are used to express and describe interactions between entities, e.g., through a simple timestamp, activity type (read, create, modify) [3], system status (CPU, RAM, bandwidth) [54] or even the script used [21]. For instance Figure 8-a shows a basic provenance model with nodes representing data objects and edges symbolizing operations. Data provenance tracking helps ensure the traceability and repeatability of processes in data lakes. Thus, provenance metadata can be used to understand, explain and repair inconsistencies in the data [3]. They may also serve to protect sensitive data, by detecting intrusions [54].

Similarity-centered graph models describe the metadata system as an undirected graph where nodes are data objects and edges express a similarity between objects. Such a similarity can be specified either through a weighted or unweighted edge. Weighted edges show the similarity strength, when a formal similarity measure is used, e.g., affinity and joinability [32] (Figure 8-b). In contrast, unweighted edges serve to simply detect whether two objects are connected [13]. Such a graph design allows network analyses over a data lake [13], e.g., discovering communities or calculating the centrality of nodes, and thus their importance in the lake. Another use of data similarity may be to automatically recommend to lake users some data related to the data they currently observe [32].

Composition-centered graph models help decompose each data object into several inherent elements. The lake is viewed as a DAG where nodes represent objects or attributes, e.g., columns, tags, etc., and edges from any node A to any node B express the constraint $B \subseteq A$ [9,20,38]. This organization helps users navigate through the data [38]. It can also be used as a basis to detect connections between objects. For instance, [9] used a simple string measure to detect links between heterogeneous objects by comparing their respective tags.

4.2.2 Data Vault

A data lake aims at ingesting new data possibly bearing various structures. Thus, its metadata system needs to be flexible to easily tackle new data schemas. Nogueira et al. propose the use of a data vault to address this issue [39]. Data vaults are indeed alternative logical models to data warehouse star schemas that, unlike star schemas, allow easy schema evolution [31]. Data vault modeling involves three types of entities [22].

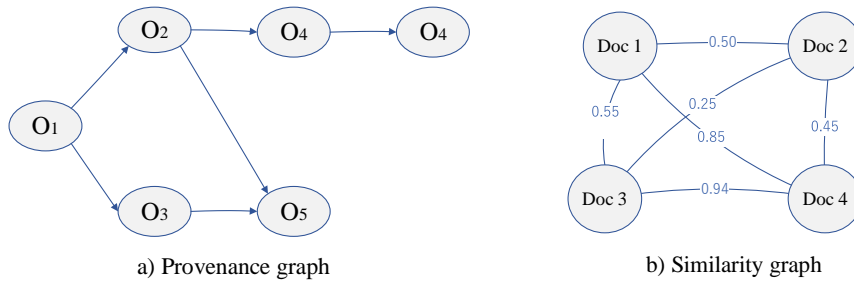


Fig. 8 Sample graph-based metadata models

1. A **hub** represents a business concept, e.g., customer, vendor, sale or product in a business decision system.
2. A **link** represents a relationship between two or more hubs.
3. **Satellites** contain descriptive information associated with a hub or a link. Each satellite is attached to a unique hub or link. In contrast, links or hubs may be associated with any number of satellites.

In Nogueira et al.'s proposal, metadata common to all objects, e.g., title, category, date and location, are stored in hubs; while metadata specific to some objects only, e.g., language for textual documents or publisher for books, are stored in satellites (Figure 9). Moreover, any new type of object would have its specific metadata stored in a new satellite.

4.2.3 Discussion

Data vault modeling is seldom associated with data lakes in the literature, presumably because it is primarily associated with data warehouses. Yet, this approach ensures metadata schema evolutivity, which is required to build an efficient data lake. Another advantage of data vault modeling is that, unlike graph models, it can be intuitively implemented in a relational DBMS. However, several adaptations are still needed for this model to deal with data linkage as in graph models.

Graph models, though requiring more specific storage systems such as RDF or graph DBMSs, are still advantageous because they allow to automatically enrich the lake with information that facilitate and enhance future analyses. Nevertheless, the three subcategories of graph models need to be all integrated together for this purpose. This remains an open issue because at most two of these graph approaches are simultaneously implemented in metadata systems from the literature [9,20]. The MEDAL metadata model does include all three subcategories of graph models [49], but is not implemented yet.

4.3 Metadata Generation

Most of the data ingestion tools from Section 3.2.1 can also serve to extract metadata. For instance, Suriarachchi and Plale use Apache Flume to retrieve data provenance in a data lake [54]. Similarly, properties and semantic metadata can be obtained through specialized protocols such as the Comprehensive Knowledge Archive Network (CKAN), an open data storage management system [55].

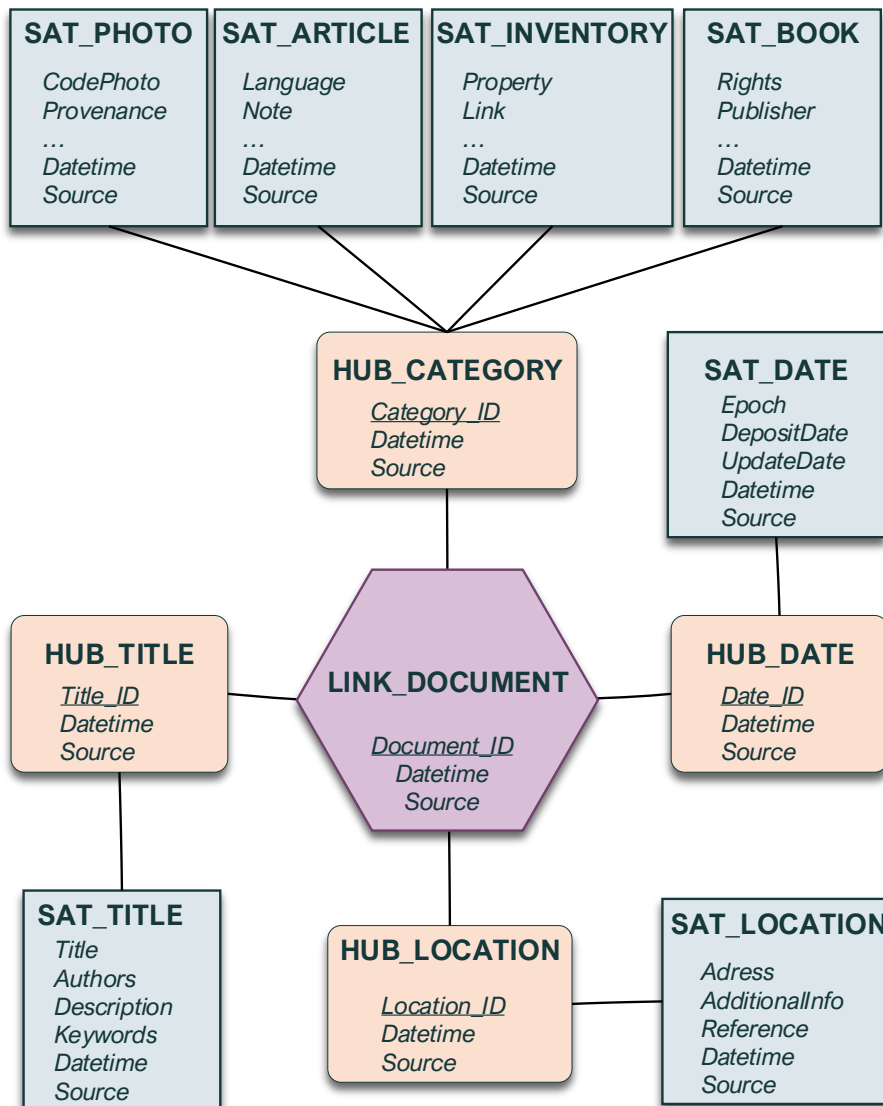


Fig. 9 Sample metadata vault model [39]

A second kind of technologies is more specific to metadata generation. For instance, Apache Tika helps detect the MIME type and language of objects [44]. Other tools such as Open Calais and IBM’s Alchemy API can also enrich data through inherent entity identification, relationship inference and event detection [12].

Ad-hoc algorithms can also generate metadata. For example, Singh et al. show that Bayesian models allow detecting links between data attributes [50]. Similarly, several authors propose algorithms to discover schemas or constraints in semi-structured data [3, 27, 44].

Last but not least, Apache Atlas [57], a widely used metadata framework [47], features advanced metadata generation methods through so-called hooks, which are native or custom scripts that rely on logs to generate metadata. Hooks notably help Atlas automatically extract lineage metadata and propagate tags on all derivations of tagged data.

4.4 Features of Data Lake Metadata Systems

A data lake made inoperable by lack of proper metadata management is called a data swamp [26], data dump [23, 54] or one-way data lake [23], with data swamp being the most common terms. In such a flawed data lake, data are ingested, but can never be extracted. Thus, a data swamp is unable to ensure any analysis. Yet, to the best of our knowledge, there is no objective way to measure or compare the efficiency of data lake metadata systems. Therefore, we first introduce in this section a list of expected features for a metadata system. Then, we present a comparison of eighteen data lake metadata systems with respect to these features.

4.4.1 Feature Identification

Sawadogo et al. identify six features that a data lake metadata system should ideally implement to be considered comprehensive [49].

1. **Semantic Enrichment (SE)** is also known as semantic annotation [18] or semantic profiling [2]. It involves adding information such as title, tags, description and more to make the data comprehensible [55]. This is commonly done using knowledge bases such as ontologies [2]. Semantic annotation plays a vital role in data lakes, since it makes the data meaningful by providing informative summaries [2]. In addition, semantic metadata could be the basis of link generation between data [44]. For instance, data objects with the same tags could be considered linked.
2. **Data Indexing (DI)** is commonly used in the information retrieval and database domains to quickly find a data object. Data indexing is done by building and enriching some data structure that enables efficient data retrieval from the lake. Indexing can serve for both simple keyword-based retrieval and more complex querying using patterns. All data, whether structured, semi-structured or unstructured, benefit from indexing [50].
3. **Link Generation (LG)** consists in identifying and integrating links between lake data. This can be done either by ingesting pre-existing links from data sources or by detecting new links. Link generation allows additional analyses. For instance, similarity links can serve to recommend to lake users data close to the data they currently use [32]. In the same line, data links can be used to automatically detect clusters of strongly linked data [13].
4. **Data Polymorphism (DP)** is the simultaneous management of several data representations in the lake. A data representation of, e.g., a textual document, may be a tag cloud or a vector of term frequencies. Semi-structured and unstructured data need to be at least partially transformed to be automatically processed [9]. Thus, data polymorphism is relevant as it allows to store and reuse transformed data. This makes analyses easier and faster by avoiding the repetition of certain processes [52].

5. **Data Versioning (DV)** expresses a metadata system’s ability to manage update operations, while retaining the previous data states. It is very relevant to data lakes, since it ensures process reproducibility and the detection and correction of inconsistencies [5]. Moreover, data versioning allows branching and concurrent data evolution [21].
6. **Usage Tracking (UT)** consists in managing information about user interactions with the lake. Such interactions are commonly creation, read and update operations. This allows to transparently follow the evolution of data objects. In addition, usage tracking can serve for data security, either by explaining data inconsistencies or through intrusion detection. Usage tracking and data versioning are related, since update interactions often induce new data versions. However, they are distinct features as they can be implemented independently [3,54].

4.4.2 Metadata System Comparison

We present in Table 2 a comparison of eighteen state-of-the-art metadata systems and models with respect to the features they implement [49]. We distinguish metadata models from implementations. Models are indeed quite theoretical and describe the conceptual organization of metadata. In contrast, implementations follow a more operational approach, but are usually little detailed, mainly focusing on a description of the resulting system instead of the applied methodology. This comparison also considers metadata systems that are not explicitly associated with the concept of data lake by their authors, but whose characteristics allow to be considered as such, e.g., the Ground metadata model [21].

The comparison shows that the most comprehensive metadata system with respect to the features we propose is MEDAL, with all features covered. However, it is not implemented yet. The next best systems are GOODS and CoreKG, with five out of six features implemented. However, they are black box metadata systems, with few details on metadata conceptual organization. Thus, the Ground metadata model may be preferred, since it is much more detailed and almost as complete (four out of six features).

Eventually, two of the six features defined in Section 4.4.1 may be considered advanced. Data polymorphism and data versioning are indeed mainly found in the most complete systems such as GOODS, CoreKG and Ground. Their absence from most of metadata systems can thus be attributed to implementation complexity.

5 Pros and Cons of Data Lakes

In this section, we account for the benefits of using a data lake instead of more traditional data management systems, but also identify the pitfalls that may correspond to these expected benefits.

An important motivating feature in data lakes is **cheap storage**. Data lakes are ten to one hundred times less expensive to deploy than traditional decision-oriented databases. This can be attributed to the usage of open-source technologies such as HDFS [26,53]. Another reason is that the cloud storage often used to build data lakes reduces the cost of storage technologies. That is, the data lake owner pays only for actually used resources. However, the use of HDFS may still fuel

| System | Type | SE | DI | LG | DP | DV | UT |
|----------------|------|----|----|----|----|----|----|
| SPAR [14] | ◆# | ✓ | ✓ | ✓ | | | ✓ |
| [1] | ◆ | ✓ | | ✓ | | | |
| [55] | ◆ | ✓ | ✓ | | | ✓ | ✓ |
| Constance [18] | ◆ | ✓ | ✓ | | | | |
| GEMMS [44] | ◇ | ✓ | | | | | |
| CLAMS [12] | ◆ | ✓ | | | | | |
| [54] | ◇ | | | | ✓ | | ✓ |
| [50] | ◆ | ✓ | ✓ | ✓ | ✓ | | |
| [13] | ◆ | | | ✓ | | | |
| GOODS [20] | ◆ | ✓ | ✓ | ✓ | | ✓ | ✓ |
| CoreDB [3] | ◆ | | ✓ | | | | ✓ |
| Ground [21] | ◇# | ✓ | ✓ | | | ✓ | ✓ |
| KAYAK [32] | ◆ | ✓ | ✓ | ✓ | | | |
| CoreKG [4] | ◆ | ✓ | ✓ | ✓ | ✓ | | ✓ |
| [9] | ◇ | ✓ | | ✓ | ✓ | | |
| [36] | ◆ | ✓ | ✓ | | | | |
| CODAL [48] | ◆ | ✓ | ✓ | ✓ | ✓ | | |
| MEDAL [49] | ◇ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |

◆ : Implementations ◇ : Metadata models
 # : Model or implementation akin to a data lake

Table 2 Comparison of data lake metadata systems [49]

misconceptions, with the concept of data lake remaining ambiguous for many potential users. It is indeed often considered either as a synonym or a marketing label closely related to the HDFS technology [1, 17].

Another feature that lies at the core of the data lake concept is **data fidelity**. Unlike in traditional decision-oriented databases, original data are indeed preserved in a data lake to avoid any data loss that could occur from data pre-processing and transformation operations [15, 53]. Yet, data fidelity induces a high risk of **data inconsistency** in data lakes, due to data integration from multiple, disparate sources without any transformation [40].

One of the main benefits of data lakes is that they allow exploiting and analyzing **unstructured data** [15, 29, 53]. This is a significant advantage when dealing with big data, which are predominantly unstructured [37]. Moreover, due to the schema-on-read approach, data lakes can comply with any data type and format [6, 15, 26, 33]. Thence, data lakes enable a wider range of analyses than traditional decision-oriented databases, i.e., data warehouses and datamarts, and thus show better **flexibility and agility**. However, although the concept of data lake dates back from 2010, it has only been put in practice in the mid-2010's. Thus, implementations vary, are still maturing and there is a **lack of methodological and technical standards**, which sustains confusions about data lakes. Finally, due to the absence of an explicit schema, **data access services** and APIs are essential to enable knowledge extraction in a data lake. In other words, a data access service is a must to successfully build a data lake [1, 23], while such a service is not always present.

Next, an acclaimed advantage of data lakes over data warehouses is **real-time data ingestion**. Data are indeed ingested in data lakes without any transformation, which avoids any time lag between data extraction from sources and their

ingestion in the data lake [15,29]. But as a consequence, a data lake requires an efficient **metadata system** for ensuring data access. However, the problem lies in the “how”, i.e., the use of inappropriate methods or technologies to build the metadata system can easily turn the data lake into an inoperable data swamp [1].

More technically, data lakes and related analyses are typically implemented using distributed technologies, e.g., HDFS, MapReduce, Apache Spark, Elasticsearch, etc. Such technologies usually provide a **high scalability** [11,37]. Furthermore, most technologies used in data lakes have replication mechanisms, e.g., Elasticsearch, HDFS, etc. Such technologies allow a high resilience to both hardware and software failure and enforce **fault tolerance** [24].

Eventually, data lakes are often viewed as sandboxes where analysts can “play”, i.e., access and prepare data so as to perform various, specific, **on-the-fly analyses** [47,53]. However, such a scenario requires **expertise**. Data lake users are indeed typically data scientists [26,33], which contrasts with traditional decision systems, where business users are able to operate the system. Thus, a data lake induces a greater need for specific, and therefore more expensive, profiles. Data scientists must indeed master a wide knowledge and panoply of technologies.

Moreover, with the integration in data lakes of structured, semi-structured and unstructured, expert data scientists can discover links and **correlations between heterogeneous data** [15]. Data lakes also allow to easily integrate data “as is” from external sources, e.g., the Web or social media. Such external data can then be associated with proprietary data to generate new knowledge through **cross-analyses** [29]. However, several statistical and Artificial Intelligence (AI) approaches are not originally designed for parallel operations, nor for streaming data, e.g., K-means or K-Nearest Neighbors. Therefore, it is necessary to **readjust classical statistical and AI approaches** to match the distributed environments often used in data lakes [40], which sometimes proves difficult.

6 Conclusion

In this survey paper, we establish a comprehensive state of the art of the different approaches to design, and conceptually build a data lake. First, we state the definitions of the data lake concept and complement the best existing one. Second, we investigate alternative architectures and technologies for data lakes, and propose a new typology of data lake architectures. Third, we review and discuss the metadata management techniques used in data lakes. We notably classify metadata and introduce the features that are necessary to achieve a full metadata system. Fourth, we discuss the pros and cons of data lakes. Fifth, we summarize by a mind map the key concepts introduced in this paper (Figure 10).

Eventually, in echo to the topics we chose *not* to address in this paper (Section 1), we would like to open the discussion on important current research issues in the field of data lakes.

Data integration and transformation have long been recurring issues. Though delayed, they are still present in data lakes and made even more challenging by big data volume, variety, velocity and lack of veracity. Moreover, when transforming such data, User-Defined Functions (UDFs) must very often be used (MapReduce tasks, typically). In ETL and ELT processes, UDFs are much more

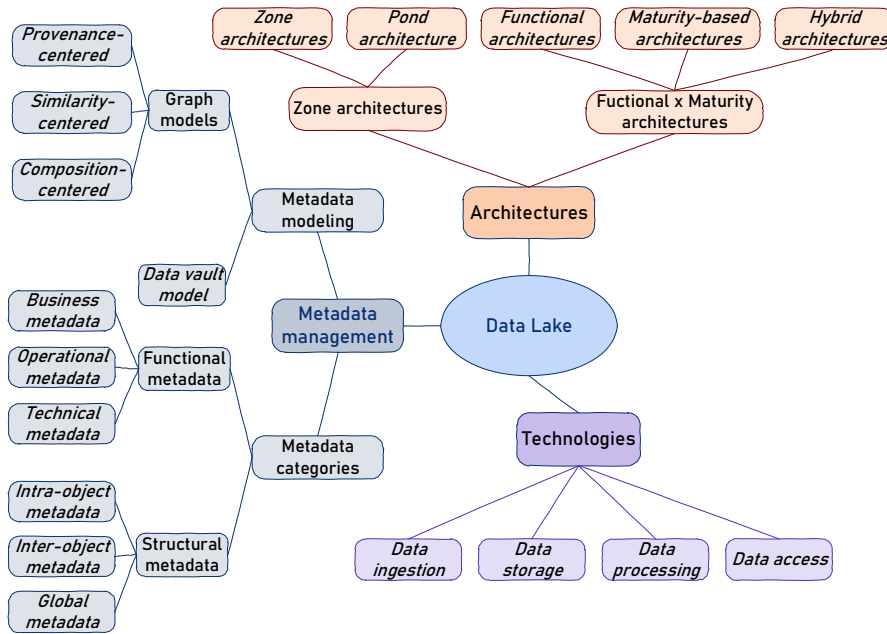


Fig. 10 Key concepts investigated in this survey

difficult to optimize than classical queries, an issue that is not addressed yet by the literature [52].

With data storage solutions currently going beyond HDFS in data lakes, **data interrogation** through metadata is still a challenge. Multistores and polystores indeed provide unified solutions for structured and semi-structured data, but do not address unstructured data, which are currently queried separately through index stores. Moreover, when considering data gravity [33], virtual data integration becomes a relevant solution. Yet, mediation approaches are likely to require new, big data-tailored query optimization and caching approaches [43,52].

Unstructured datasets, although unanimously acknowledged as ubiquitous and sources of crucial information, are very little specifically addressed in data lake-related literature. Index storage and text mining are usually mentioned, but there is no deep thinking about global querying or analysis solutions. Moreover, exploiting other types of unstructured data but text, e.g., images, sounds and videos, is not even envisaged as of today.

Again, although all actors in the data lake domain stress the importance of **data governance** to avoid a data lake turning into a data swamp, data quality, security, life cycle management and metadata lineage are viewed as risks rather than issues to address *a priori* in data lakes [33]. Data governance principles are indeed currently seldom turned into actual solutions.

Finally, **data security** is currently addressed from a technical point of view in data lakes, i.e., through access and privilege control, network isolation, e.g., with Docker tools [6], data encryption and secure search engines [34]. However, beyond these issues and those already addressed by data governance (integrity, consistency, availability) and/or related to the European General Data Protection Regulation

(GDPR), by storing and cross-analyzing large volumes of various data, data lakes allow mashups that potentially induce serious breaches of data privacy [25]. Such issues are still researched as of today.

References

1. Alrehamy, H., Walker, C.: Personal Data Lake With Data Gravity Pull. In: IEEE 5th International Conference on Big Data and Cloud Computing(BDCloud 2015), Dalian, china, *IEEE Computer Society Washington*, vol. 88, pp. 160–167 (2015). DOI 10.1109/BDCloud.2015.62
2. Ansari, J.W., Karim, N., Decker, S., Cochez, M., Beyan, O.: Extending Data Lake Metadata Management by Semantic Profiling. In: 2018 Extended Semantic Web Conference (ESWC 2018), Heraklion, Crete, Greece, pp. 1–15 (2018)
3. Beheshti, A., Benatallah, B., Nouri, R., Chhieng, V.M., Xiong, H., Zhao, X.: CoreDB: a Data Lake Service. In: 2017 ACM on Conference on Information and Knowledge Management (CIKM 2017), Singapore, Singapore, ACM, pp. 2451–2454 (2017). DOI 10.1145/3132847.3133171
4. Beheshti, A., Benatallah, B., Nouri, R., Tabebordbar, A.: CoreKG: A Knowledge Lake Service. *Proceedings of the VLDB Endowment* **11**(12), 1942–1945 (2018). DOI 10.14778/3229863.3236230
5. Bhattacharjee, S., Deshpande, A.: RStore: A Distributed Multi-Version Document Store. In: IEEE 34th International Conference on Data Engineering (ICDE), Paris, France, pp. 389–400 (2018). DOI 10.1109/ICDE.2018.00043
6. Cha, B., Park, S., Kim, J., Pan, S., Shin, J.: International Network Performance and Security Testing Based on Distributed Abyss Storage Cluster and Draft of Data Lake Framework. *Hindawi Security and Communication Networks* **2018**, 1–14 (2018). DOI 10.1155/2018/1746809
7. Chessell, M., Scheepers, F., Nguyen, N., van Kessel, R., van der Starre, R.: Governing and Managing Big Data for Analytics and Decision Makers. IBM (2014)
8. Couto, J., Borges, O., Ruiz, D., Marczak, S., Prikladnicki, R.: A Mapping Study about Data Lakes: An Improved Definition and Possible Architectures. In: 31st International Conference on Software Engineering and Knowledge Engineering (SEKE 2019), Lisbon, Portugal, pp. 453–458 (2019). DOI 10.18293/SEKE2019-129
9. Diamantini, C., Giudice, P.L., Musarella, L., Potena, D., Storti, E., Ursino, D.: A New Metadata Model to Uniformly Handle Heterogeneous Data Lake Sources. In: New Trends in Databases and Information Systems - ADBIS 2018 Short Papers and Workshop, Budapest, Hungary, pp. 165–177 (2018). DOI 10.1007/978-3-030-00063-9_17. URL https://doi.org/10.1007/978-3-030-00063-9_17
10. Dixon, J.: Pentaho, Hadoop, and Data Lakes (2010). URL <https://jamesdixon.wordpress.com/2010/10/14/pentaho-hadoop-and-data-lakes/>
11. Fang, H.: Managing Data Lakes in Big Data Era: What’s a data lake and why has it become popular in data management ecosystem. In: 5th Annual IEEE International Conference on Cyber Technology in Automation, Control and Intelligent Systems (CYBER 2015), Shenyang, China, IEEE, pp. 820–824 (2015). DOI 10.1109/CYBER.2015.7288049
12. Farid, M., Roatis, A., Ilyas, I.F., Hoffmann, H.F., Chu, X.: CLAMS: Bringing Quality to Data Lakes. In: 2016 International Conference on Management of Data (SIGMOD 2016), San Francisco, CA, USA, ACM, pp. 2089–2092 (2016). DOI 10.1145/2882903.2899391
13. Farrugia, A., Claxton, R., Thompson, S.: Towards Social Network Analytics for Understanding and Managing Enterprise Data Lakes. In: Advances in Social Networks Analysis and Mining (ASONAM 2016), San Francisco, CA, USA, IEEE, pp. 1213–1220 (2016). DOI 10.1109/ASONAM.2016.7752393
14. Fauduet, L., Peyrard, S.: A data-first preservation strategy: Data management in spar. In: 7th International Conference on Preservation of Digital Objects (iPRES 2010), Vienna, Austria, pp. 1–8 (2010). URL <http://www.ifs.tuwien.ac.at/dp/ipres2010/papers/fauduet-13.pdf>
15. Ganore, P.: Introduction To The Concept Of Data Lake And Its Benefits. <https://www.esds.co.in/blog/introduction-to-the-concept-of-data-lake-and-its-benefits> (2015)

16. Giebler, C., Gröger, C., Hoos, E., Schwarz, H., Mitschang, B.: Leveraging the Data Lake - Current State and Challenges. In: Proceedings of the 21st International Conference on Big Data Analytics and Knowledge Discovery (DaWaK 2019), Linz, Austria (2019)
17. Grosser, T., Bloeme, J., Mack, M., Vitsenko, J.: Hadoop and Data Lakes: Use Cases, Benefits and Limitations. Business Application Research Center – BARC GmbH (2016)
18. Hai, R., Geisler, S., Quix, C.: Constance: An Intelligent Data Lake System. In: International Conference on Management of Data (SIGMOD 2016), San Francisco, CA, USA, ACM Digital Library, pp. 2097–2100 (2016). DOI 10.1145/2882903.2899389
19. Hai, R., Quix, C., Zhou, C.: Query Rewriting for Heterogeneous Data Lakes. In: 22nd European Conference on Advances in Databases and Information Systems (ADBIS 2018), Budapest, Hungary, *LNC5*, vol. 11019, pp. 35–49. Springer (2018). DOI 10.1007/978-3-319-98398-1_3. URL https://doi.org/10.1007/978-3-319-98398-1_3
20. Halevy, A.Y., Korn, F., Noy, N.F., Olston, C., Polyzotis, N., Roy, S., Whang, S.E.: Goods: Organizing Google’s Datasets. In: Proceedings of the 2016 International Conference on Management of Data (SIGMOD 2016), San Francisco, CA, USA, pp. 795–806 (2016). DOI 10.1145/2882903.2903730
21. Hellerstein, J.M., Sreekanti, V., Gonzalez, J.E., Dalton, J., Dey, A., Nag, S., Ramachandran, K., Arora, S., Bhattacharyya, A., Das, S., Donsky, M., Fierro, G., She, C., Steinbach, C., Subramanian, V., Sun, E.: Ground: A Data Context Service. In: 8th Biennial Conference on Innovative Data Systems Research (CIDR 2017), Chaminade, CA, USA (2017). URL <http://cidrdb.org/cidr2017/papers/p111-hellerstein-cidr17.pdf>
22. Hultgren, H.: Data Vault modeling guide: Introductory Guide to Data Vault Modeling. Genessee Academy, USA (2016)
23. Inmon, B.: Data Lake Architecture: Designing the Data Lake and avoiding the garbage dump. Technics Publications (2016)
24. John, T., Misra, P.: Data Lake for Enterprises: Lambda Architecture for building enterprise data systems. Packt Publishing (2017)
25. Joss, A.: The Rise of the GDPR Data Lake. <https://blogs.informatica.com/2016/06/16/rise-gdpr-data-lake/> (2016)
26. Khine, P.P., Wang, Z.S.: Data Lake: A New Ideology in Big Data Era. In: 4th International Conference on Wireless Communication and Sensor Network (WCSN 2017), Wuhan, China, *ITM Web of Conferences*, vol. 17, pp. 1–6 (2017). DOI 10.1051/itmconf/2018170302
27. Klettke, M., Awolin, H., Stürl, U., Müller, D., Scherzinger, S.: Uncovering the Evolution History of Data Lakes. In: 2017 IEEE International Conference on Big Data (BIGDATA 2017), Boston, MA, USA, pp. 2462–2471 (2017). DOI 10.1109/BigData.2017.8258204
28. LaPlante, A., Sharma, B.: Architecting Data Lakes Data Management Architectures for Advanced Business Use Cases. O’Reilly Media, Inc. (2016)
29. Laskowski, N.: Data lake governance: A big data do or die. <https://searchcio.techtarget.com/feature/Data-lake-governance-A-big-data-do-or-die> (2016)
30. Leclercq, E., Savonnet, M.: A Tensor Based Data Model for Polystore: An Application to Social Networks Data. In: Proceedings of the 22nd International Database Engineering & Applications Symposium (IDEAS 2018), Villa San Giovanni, Italy, pp. 110–118 (2018). DOI 10.1145/3216122.3216152
31. Linstedt, D.: Super Charge your Data Warehouse: Invaluable Data Modeling Rules to Implement Your Data Vault. CreateSpace Independent Publishing (2011)
32. Maccioni, A., Torlone, R.: KAYAK: A Framework for Just-in-Time Data Preparation in a Data Lake. In: International Conference on Advanced Information Systems Engineering (CAiSE 2018), Tallin, Estonia, pp. 474–489 (2018). DOI 10.1007/978-3-319-91563-0_29
33. Madera, C., Laurent, A.: The next information architecture evolution: the data lake wave. In: 8th International Conference on Management of Digital EcoSystems (MEDES 2016), Biarritz, France, pp. 174–180 (2016). DOI 10.1145/3012071.3012077
34. Maroto, C.: Data Lake Security – Four Key Areas to Consider When Securing Your Data Lake. <https://www.searchtechnologies.com/blog/data-lake-security> (2018)
35. Mathis, C.: Data Lakes. *Datenbank-Spektrum* **17**(3), 289–293 (2017). DOI 10.1007/s13222-017-0272-7
36. Mehmood, H., Gilman, E., Cortes, M., Kostakos, P., Byrne, A., Valta, K., Tekes, S., Riekk, J.: Implementing big data lake for heterogeneous data sources. In: 2019 IEEE 35th International Conference on Data Engineering Workshops (ICDEW), pp. 37–44 (2019). DOI 10.1109/ICDEW.2019.00-37

37. Miloslavskaya, N., Tolstoy, A.: Big Data, Fast Data and Data Lake Concepts. In: 7th Annual International Conference on Biologically Inspired Cognitive Architectures (BICA 2016), NY, USA, *Procedia Computer Science*, vol. 88, pp. 1–6 (2016). DOI 10.1016/j.procs.2016.07.439
38. Nargesian, F., Pu, K.Q., Zhu, E., Bashardoost, B.G., Miller, R.J.: Optimizing Organizations for Navigating Data Lakes. <https://arxiv.org/abs/1812.07024> (2018)
39. Nogueira, I., Romdhane, M., Darmont, J.: Modeling Data Lake Metadata with a Data Vault. In: 22nd International Database Engineering and Applications Symposium (IDEAS 2018), Villa San Giovanni, Italia, pp. 253–261. ACM, New York (2018)
40. O’Leary, D.E.: Embedding AI and Crowdsourcing in the Big Data Lake. *IEEE Intelligent Systems* **29**(5), 70–73 (2014). DOI 10.1109/MIS.2014.82
41. Oram, A.: *Managing the Data Lake*. Zaloni (2015)
42. Pathirana, N.: *Modeling Industrial and Cultural Heritage Data*. Master’s thesis, Université Lumière Lyon 2, France (2015)
43. Quix, C., Hai, R.: *Data Lake*, pp. 1–8. Springer International Publishing (2018). DOI 10.1007/978-3-319-63962-8_7-1
44. Quix, C., Hai, R., Vatov, I.: Metadata Extraction and Management in Data Lakes With GEMMS. *Complex Systems Informatics and Modeling Quarterly* (9), 289–293 (2016). DOI 10.7250/csimq.2016-9.04
45. Ravat, F., Zhao, Y.: Data Lakes: Trends and Perspectives. In: 30th International Conference on Database and Expert Systems Applications (DEXA 2019), Linz, Austria (2019)
46. Ravat, F., Zhao, Y.: Metadata management for data lakes. In: 23rd European Conference on Advances in Databases and Information Systems (ADBIS 2019), Bled, Slovenia (2019)
47. Russom, P.: *Data Lakes Purposes, Practices, Patterns, and Platforms*. TDWI research (2017)
48. Sawadogo, P.N., Kibata, T., Darmont, J.: Metadata Management for Textual Documents in Data Lakes. In: 21st International Conference on Enterprise Information Systems (ICEIS 2019), Heraklion, Crete, Greece, pp. 72–83 (2019). DOI 10.5220/0007706300720083
49. Sawadogo, P.N., Scholly, E., Favre, C., Ferey, É., Loudcher, S., Darmont, J.: Metadata Systems for Data Lakes: Models and Features. In: BI and Big Data Applications - ADBIS 2019 Short Papers and Workshop, Bled, Slovenia (2019)
50. Singh, K., Paneri, K., Pandey, A., Gupta, G., Sharma, G., Agarwal, P., Shroff, G.: Visual Bayesian Fusion to Navigate a Data Lake. In: 19th International Conference on Information Fusion (FUSION 2016), Heidelberg, Germany, IEEE, pp. 987–994 (2016)
51. Sirosh, J.: *The Intelligent Data Lake* (2016). URL <https://azure.microsoft.com/fr-fr/blog/the-intelligent-data-lake/>
52. Stefanowski, J., Krawiec, K., Wrembel, R.: Exploring Complex and Big Data. *International Journal of Applied Mathematics and Computer Science* **27**(4), 669–679 (2017). DOI 10.1515/amcs-2017-0046
53. Stein, B., Morrison, A.: The enterprise data lake: Better integration and deeper analytics. PWC Technology Forecast (2014). URL <http://www.smalllake.kr/wp-content/uploads/2017/03/20170313.074222.pdf>
54. Suriarachchi, I., Plale, B.: Crossing Analytics Systems: A Case for Integrated Provenance in Data Lakes. In: 12th IEEE International Conference on e-Science (e-Science 2016), Baltimore, MD, USA, pp. 349–354 (2016). DOI 10.1109/eScience.2016.7870919
55. Terrizzano, I., Schwarz, P., Roth, M., Colino, J.E.: Data Wrangling: The Challenging Journey from the Wild to the Lake. In: 7th Biennial Conference on Innovative Data Systems Research (CIDR 2015), Asilomar, CA, USA, pp. 1–9 (2015). URL http://cidrdb.org/cidr2015/Papers/CIDR15_Paper2.pdf
56. Tharrington, M.: *The Dzone Guide to Big Data, Data Science & Advanced Analytics*. DZone (2017)
57. The Apache Software Foundation: *Apache Atlas – Data Governance and Metadata framework for Hadoop*. <https://atlas.apache.org/> (2019)
58. Tiao, S.: Object Storage for Big Data: What Is It? And Why Is It Better? <https://blogs.oracle.com/bigdata/what-is-object-storage> (2018)
59. Zikopoulos, P., deRoos, D., Bienko, C., Buglio, R., Andrews, M.: *Big Data Beyond the Hype*. McGraw-Hill Education (2015)