Preserving domain knowledge in a data lake

Master Thesis Master of IT, Software Construction

> Mads Staberg Thomsen msth21@student.aau.dk

Supervisor: Hamdi Ben Hamadou, PhD Aalborg University, Denmark

Department of Computer Science Aalborg University Selma Lagerlöfs Vej 300, 9220 Aalborg, Denmark

2022-12-02



Table of Content

List of figures and tables	4
Abstract	6
Chapter 1 - Motivation	7
Chapter 2 - Problem Statement	8
2.1 Project Scope	8
Chapter 3 - Terminology	9
3.1 Data discovery	9
3.2 Data lake	9
3.3 Data lineage	9
3.4 Domain knowledge	9
3.5 Metadata	10
Chapter 4 - Methods	11
4.1 Metadata concept	11
4.2 Sample data	12
4.3 Quality Attribute Scenarios	16
Chapter 5 - Sandbox environment	17
5.1 Data catalog	17
5.2 Data lake	19
Chapter 6 - Analysis Model	20
6.1 Usability	20
6.1.1 Automatic infer	20
6.1.2 Manual infer	21
6.1.3 Semi-automatic infer	22
6.2 Performance	22
6.2.1 Implicit infer	23
Chapter 7 - Analyses and Results	24
7.1 Automatic infer	24
7.1.1 Logbook	24
7.2 Manual infer	27
7.2.1 Logbook	27
7.3 Semi-automatic infer	29
7.3.1 Logbook	30
7.4 Implicit infer	32

7.4.1 Logbook	33
7.5 Discussion	36
7.5.1 Usability	36
7.5.2 Performance	37
Chapter 8 - Conclusion	39
8.1 Reflections	40
Bibliography	41
Appendix	43
A. Production environment Size	43
A1 Landing zone	43
A2. Gold zone	44
B. Sandbox environment DevOps	45
B1. Repository content	45
C. Automatic infer Screenshots	46
C1. Initial scan	46
C2. Data discovery Browse assets	46
C3. Data discovery Asset overview	47
C4. Schema registry after change	47
D. Manual infer Screenshots	48
D1. Business glossary terms	48
D2. Data discovery Filters	48
D3. Data discovery Asset schema	49
D4. Data discovery Domain expert knowledge	49
E. Semi-automatic infer Screenshots	50
E1. Business glossary terms	50
E2. Classification rules	50
F. Implicit infer Screenshots	51
F1. Unstructured landing zone Test #1	51
F2. Unstructured landing zone Test #2	52

List of figures and tables

Figure 1: Distribution Network	7
Figure 2: Data discovery Self-service	10
Figure 3: Metadata concept Entities	11
Figure 4: Metadata concept Logical model	12
Figure 5: Sample data Measurement.csv	13
Figure 6: Sample data Sensor.csv	13
Figure 7: Sample data Sensorconfig.csv	13
Figure 8: Data model Sample data	14
Figure 9: Data model Sample data with metadata	14
Figure 10: Metadata concept Logical model with sample data	15
Figure 11: Quality Attribute Scenario Template	16
Figure 12: Data catalog Exploring data	17
Figure 13: Data catalog Metadata concept	18
Figure 14: Data catalog Data map	18
Figure 15: Data lake Zones	19
Figure 16: Data lake Sandbox storage accounts	19
Figure 17: Data lake Structured directory	23
Figure 18: User actions Automatic infer	24
Figure 19: User actions Manual infer	27
Figure 20: User actions Semi-automatic infer	30
Figure 21: User actions Semi-automatic infer	32
Table 1: Metadata concept Properties	12
Table 2: Metadata Domain expert knowledge	15

Table 3: Quality Attribute Scenario Usability Automatic infer	21
Table 4: Quality Attribute Scenario Usability Manual infer	21
Table 5: Quality Attribute Scenario Usability Semi-automatic infer	22
Table 6: Quality Attribute Scenario Performance Implicit infer	23
Table 7: Logbook Usability Automatic infer	26
Table 8: Logbook Usability Manual infer	29
Table 9: Logbook Usability Semi-automatic infer	32
Table 10: Logbook Performance Implicit infer	35
Table 11: Results Usability	36
Table 12: Results Performance	38

Abstract

This thesis is based on real-world challenges for a data-intensive application in a data lake context, e.g. difficulties with understanding the domain and the meaning of the data. It explores how to optimize data discovery by preserving domain knowledge in a time-saving way.

In order to transform domain knowledge into metadata, the thesis contributes with a definition of novel metadata concept. The metadata concept makes the task of implementing a data catalog and building a metadata fundament manageable and accessible. The concept guides the user to navigate through the implementation of metadata in the data catalog and classify the data. The concept is both tool, platform and domain independent.

Sample data and the corresponding metadata is introduced. A data catalog is implemented, which is a collection of metadata, combined with data management and search tools, that helps users to find the data that they need. The catalog serves as an inventory of available data, and provides information of the data quality. A small sandbox environment is set up with some sample data.

A custom analysis model is defined to learn how to handle the challenges. To ensure the analyses are conducted in a reproducible manner, the scenarios are expressed in a Quality Attribute Scenario framework. The scenario template outlines response measures, which are measurable.

Usability and performance qualities are in focus. The usability scenarios cover automatic, manual and semi-automatic infer of metadata into the data catalog. A performance scenario explores how implicit infer of metadata using a structured directory affects the fetching of data from the data lake.

The reader is guided through each analysis with a process map based on the metadata concept and a corresponding logbook. The user actions are marked and explained. The development logbook documents the steps in the analyses. Issues, solutions and learnings are elaborated. A result section at the end of each logbook clarifies the partial results.

All results are combined and grouped by usability and performance scenarios. The results are discussed based on the metadata concept and the problem statement. Learnings from the analyses are highlighted. Based on the usability and performance scenarios, it can be concluded that it is possible to preserve domain knowledge in a time-saving way and thus optimize data discovery. The best usability results are obtained with semi-automatic infer and storage of data in a structured directory provides the best performance.

However, it should be noted that automatic infer in some cases can be a good choice, as it is possible to create searchable technical metadata very quickly. This is a much better starting point than having no information in the data catalog.

Chapter 1 - Motivation

Aarhus Vand is a large water utility company, which handles drinking water in most of Aarhus municipality. The distribution network is divided into zones and monitored by a number of sensors. The sensors are distributed in the network based on domain-specific knowledge. Each sensor measures the water's temperature, flow and pressure every minute, and data is continuously transferred wirelessly to the backend Scada control system. A part of the southern distribution network is illustrated in Figure 1.



Figure 1: Distribution Network

One measurement every minute, 24 hours per day, 365 days per year creates approx 500.000 data rows. Aarhus Vand has more than 4.000 sensors. Every year, this results in billions of data rows being created and must be handled in a big data context. A small part of this data is uploaded to an Azure Data lake, where they are consumed by a third party analysis tool. Every hour approx. 8.100 measurements are uploaded to the data lake. The number of records increases by approx. 6 million measurements every month. The current files in the data lake contain approx. 52 million records. Further details are outlined in Appendix A.

Data is loaded into the data lake in their raw / non-conform state, so it is quick to add new data. Since the sensor measurements are being stored without any corresponding metadata, users find it challenging to inspect, understand and consume the data. Thus, the deep domain knowledge about the distribution network and the measurement's internal relationship is lost and use of the data may require domain expert knowledge. If data was stored in a data warehouse, then the domain experts will interpret and combine data according to the current need in a conforming manner. But it will be a heavy process that will limit future use of data.

Chapter 2 - Problem Statement

The purpose of this project is to investigate how domain knowledge can be preserved in a data lake to improve data discovery. Managing domain knowledge can be a time-consuming process and seems like a difficult task. Therefore, the project intends to answer the following problem statement:

Can domain knowledge be preserved in a time-saving way into a data lake in order to optimize the data discovery?

The criterias for domain knowledge, metadata and data discovery are determined by a custom analysis model. The problem statement will be elaborated through the hypothesis:

By constructing a custom analysis model, where automatic, manual, semi-automatic and implicit ingesting of metadata are explored, it can be evaluated whether domain knowledge in the data lake optimizes the data discovery or not.

The current setup of the data lake at Aarhus Vand presents some real-world challenges for the data consumers, eg.:

- Difficult to understand the domain and meaning of the data
- Impossible to know which sensors make sense to compare to get useful information
- Time-consuming to fetch measurements from a specific period or a certain sensor because all data are stored in one file

A metadata concept will help navigating through the exploration. Some metadata and sample data implemented in a sandbox environment will be used as a basis for the analyses, where the challenges are investigated and tested in different scenarios. By performing a series of analyses the aim is to learn more about how the problem can be handled.

2.1 Project Scope

Although the hypothesis, analysis model and analyses will be relevant to many domains, the project will be limited to use only a subset of the data from the drinking water domain at Aarhus Vand. Handling of access control and General Data Protection Regulation (GDPR) issues are not in scope.

Chapter 3 - Terminology

3.1 Data discovery

Data Discovery involves the collection and evaluation of data from various sources and is often used to understand trends and patterns in the data. It requires a progression of steps that organizations can use to understand their data. Data discovery can be treated as a synonym for Knowledge Discovery in Data (KDD) [Han2012]. The KDD process includes collecting data from multiple data sources, cleaning, integration, task-relevant selection, transformation, and performing analysis to gain knowledge and insights into business processes [Nwagu2017].

In this project scope, data discovery includes the workflow of 'find and understand' data, which includes integration of different sources and selection of task-relevant data.

3.2 Data lake

A data lake stores current and historical data for one or more systems for the purpose of analyzing the data. Data may include raw copies of source system data and transformed data used for tasks such as reporting, visualization, advanced analytics and machine learning. Files of different file formats generated from different sources can be stored in a data lake [Hamadou2020]. A data lake can be established on premises within an organization's data centers or be cloud based using cloud services.

The data lake will typically be structured in three zones [Thomsen2022]:

- Landing zone where the raw data enters in different formats including structured, unstructured and binary data
- Work zone where work-in-progress data is stored
- Gold zone with data that has been cleansed and processed

3.3 Data lineage

Data lineage uncovers the life cycle of data and aims to show the complete data flow, from start to finish. Data lineage is the process of understanding, recording, and visualizing data as it flows from data sources to consumption. This includes all transformations the data underwent along the way; where data came from and how the data was transformed [Gorelik2019].

3.4 Domain knowledge

Domain knowledge is knowledge of a specific, specialized discipline or field. People with domain knowledge are often regarded as specialists or experts in their field. Data does not "speak on its own" and attention is needed to also convey information about the problem domain as well [Videla2021]. In this project scope, domain knowledge will be expressed through metadata.

3.5 Metadata

Metadata is data about the data and provides information about one or more aspects of data, e.g. description, data quality, value ranges, dates for creation and update [Hamadou2020]. Metadata will give the user confidence in the data. Metadata unlocks the value of data by improving that data's usability and findability.

There are different types of metadata [Atlas2022]:

- Technical: schemas, data types, models etc.
- Business: data tags, classifications, mapping to business
- Operational: process output, performance, Extract-Transform-Load jobs

This project explores technical and business metadata. When data is tagged with metadata, it will help the user to become self-serviced in the 'find and understand' data discovery process. The self-service data discovery idea is illustrated in Figure 2.



Figure 2: Data discovery | Self-service

Chapter 4 - Methods

This chapter outlines the method for exploring the hypothesis. A metadata concept introduces a conceptual way of interpreting domain knowledge. Some metadata and sample data implemented will be used as a basis for the analyses. The scenarios will be described using a Quality Attribute Scenarios framework.

4.1 Metadata concept

To transform domain knowledge into metadata, this project contributes with the definition of a metadata concept. The concept offers a conceptual way of interpreting domain knowledge into different entities, including metadata types.

- **Asset**: a piece of data information about the data but not the data itself
- Technical metadata: column with data type and relation between assets
- **Business metadata**: tags based on classification and business terms
- **Data store**: source containing the asset, e.g. file or database

The entities are included as they precisely describe domain knowledge. The metadata concept is visualized in Figure 3.



Figure 3: Metadata concept | Entities

Each entity in the metadata concept can be extended with properties that define characteristics of the elements. The reason for the characteristics being chosen is explained in the Table 1.

Туре	Entities	Properties Description
Asset	Data asset	Includes a description, a timestamp for update of the schema and a domain expert validation flag.
Technical	Column	Name, data type and description. Collection of columns makes up the schema. Changes to the schema can be tracked in a schema registry.

1	ſechnical	Lineage	Shows the relationship among the data assets.
E	Business	Classification	Describes patterns to look for as input for tagging.
E	Business	Business term	Explains domain specific knowledge.
E	Business	Contact	Person who has ownership of data and an associated role.
E	Business	Collection	Shows logical boundaries in the organization.
Ι	Data store	Source	Contains data in several different formats.

Table 1: Metadata concept | Properties

The metadata concept with entities is outlined as a class UML diagram at the logical level in Figure 4. The metadata concept is tool independent and can be implemented in any data catalog tool. Furthermore, the concept is domain independent.



Figure 4: Metadata concept | Logical model

4.2 Sample data

The basis for the analyses are some sample data, which is a subset of the raw data with total measurements as outlined in Appendix A. The subset covers measurements per minute from the period 06.06 - 12.06.2022. Data is divided into three different files:

- Measurement.csv
- SensorConfig.csv

- Sensor.csv

The three files with sample data are uploaded in the project's code repository placed at Azure DevOps as described in Appendix B.

Measurement

The measurement from each sensor with a timestamp. The size of the measurement.csv file is 125 MB and contains approx. 1.450.000 records. The header and the first few records are shown in Figure 5.

Datetime 🔽	Sensorid	Measure	🕶 Updatetime 🔤	FileId 💌
2022-06-06 00:00:00.000	ABO.B10_E_W3_BF1-M_FLOW	2.369791746	2022-06-06 01:29:38.000	3
2022-06-06 00:00:00.000	BED.B12_E_B1-M_TRYK	6.62808609	2022-06-06 00:29:53.000	9
2022-06-06 00:00:00.000	BED.B12_E_B2-M_FLOW	219.816482544	2022-06-06 01:29:38.000	3
2022-06-06 00:00:00.000	BED.B12_E_G3_B1-M_TRYK	4.288676739	2022-06-06 00:29:10.000	10
2022-06-06 00:00:00.000	BED.B12_E_G3_BF1-M_FLOW	27.178453445	2022-06-06 11:29:57.000	4
2022-06-06 00:00:00.000	BUS.C04_E_C1_B1-M_NIVEAU	3.723090172	2022-06-06 00:31:42.000	8
2022-06-06 00:00:00.000	BUS.C04_E_G1_BF1-M_FLOW	18.900817871	2022-06-06 11:29:57.000	4
2022-06-06 00:00:00.000	BUS.C04_E_G1_BF1-M_FLOW	18.900817871	2022-06-06 01:30:15.000	2
2022-06-06 00:00:00.000	BUS.C04_E_G1_BP1-M_TRYK	4.951388836	2022-06-06 00:31:42.000	8

Figure 5: Sample data | Measurement.csv

Sensor

Metadata about each sensor including parameter, unit and zone according to Figure 6. File size is 10 KB and contains 93 records.

SensorId	Description	Parameter	Unit 💌	Media 💌	Time se 💌	Zone 💌	Location 🔹
OST.B09_E_G1_BF1-M_FLOW	Flow Østerby til Z80S	Q	m3/h	Drikkevand	Sensor	Z80Syd	Østerbyværket
BUS.C04_E_G1_BF1-M_FLOW	Flow fra Bushøj til Z125 Syd	Q	m3/h	Drikkevand	Sensor	Z80Syd	Beholder Bushøj
BUS.C04_E_G2_BF2-M_FLOW	Flow fra Bushøj til Z105 Syd	Q	m3/h	Drikkevand	Sensor	Z80Syd	Beholder Bushøj
BUS.C04_E_W1_BF1-M_FLOW	Flow fra Bushøj til Z80S	Q	m3/h	Drikkevand	Sensor	Z80Syd	Beholder Bushøj
VIV.B08_E_G1_BF1-M_FLOW	Flow udveksling Z50M til Z80S (viv)	Q	m3/h	Drikkevand	Sensor	Z80Syd	Beholder Vibyværket
OBS.C11_E_G1_BF1-M_FLOW	Flow udveksling Z50M til Z80S (obs)	Q	m3/h	Drikkevand	Sensor	Z80Syd	Beholder Observatoriet
DIS.D01_H_B84_BF1-M_FLOW	DB84 Flow S80 Grøndalsvej	Q	m3/h	Drikkevand	Sensor	Z80Syd	DB84: Skanderborgvej 226

Figure 6: Sample data | Sensor.csv

SensorConfig

A sensor configuration file containing a few metadata about the different zones and measurement types as displayed in Figure 7. File size is 1 KB and contains 11 records.

FileId 🔽 FileName	🔽 DecimalSep	arator 💌 ColumnDelimiter 💌
2 Flows_Z80Syd.csv	· · · · · ·	;
3 Flows_Z105Syd.csv	,	;
4 Flows_Z125Syd_Z80	Beder.csv ,	;
5 Temperatur_Z80Syd	.csv ,	;
6 Temperatur_Z105Sy	d.csv ,	;
7 Temperatur_Z125Sy	d_Z80Beder.csv ,	;
8 Tryk_Z80Syd.csv	,	;
9 Tryk_Z105Syd.csv	,	;
10 Tryk_Z125Syd_Z80Be	eder.csv ,	;
11 FlowTrykTemperatu	ır_Z125N.csv ,	;

Figure 7: Sample data | Sensorconfig.csv

The internal relationship between the files can be described at logical level in a class diagram with Unified Modeling Language (UML) notation. The sample data are .csv format thus without data type, but the types are stated in the class diagram in Figure 8.



Figure 8: Data model | Sample data

An interview with the domain expert unlocked knowledge hidden in the SensorId, as this uses a unique template. Business specific information can be extracted from the SensorId field. This is shown in Figure 9 where the extracted metadata based on the template is indicated below the dashed line in the 'Sensor' class.



Figure 9: Data model | Sample data with metadata

The interview with the domain expert also revealed some aspects of the network and measurements. This extended domain knowledge is reproduced in Table 2.

Subject	Domain expert knowledge	
Flow	Measurements with false 'IsMid' values have more noise as it is a snapshot of a given minute	
Pressure	Values can be used to detect pressure shocks and extend lifetime	
Temperature	Depending on the distance to the surface	
Context The fluctuations over the day have morning and evening pea		
	Increased flow creates increased pressure	
	Increased flow does not immediately result in lower temperature	

Table 2: Metadata | Domain expert knowledge

To get an idea of how to transform the domain knowledge and the metadata from the sample data, the classes in the metadata concept logical model can be filled with an example as shown in Figure 10.



Figure 10: Metadata concept | Logical model with sample data

4.3 Quality Attribute Scenarios

To ensure the analyses are conducted in a reproducible manner, the scenarios are expressed in Quality Attribute Scenarios [Bass2012, p. 196].

Quality Attribute Scenarios are a way of setting a testable frame around a specific quality the system should have, and it provides a structured approach for the evaluation of requirements.

To measure software we need to define the qualities we measure and some kind of metric. The framework handles quality attributes, metrics, measurements and uses a measurement template. The template captures requirements and measurements in a common format. The composition of a Quality Attribute Scenario (QAS) is illustrated in Figure 11.



Figure 11: Quality Attribute Scenario | Template

Key points of the QAS template [Christensen2020]:

- a *source* generates some *stimulus*
- arrives as some *artifact* in an *environment*
- must be dealt with a *response* with a satisfactory *response measure*

The custom analysis model constructed in chapter 6 will elaborate Usability and Performance quality attributes in different scenarios.

Chapter 5 - Sandbox environment

The following sections describe the implementation of a data catalog and data lake in the sandbox environment. The setup is documented in the project's code repository located on Azure DevOps and available in Appendix B.

5.1 Data catalog

A data catalog is a collection of metadata, combined with data management and search tools, that helps users to find the data that they need. The catalog serves as an inventory of available data, and provides information of the data quality. Searching a data catalog is a great tool for data discovery, especially for the workflow of 'find and understand' data as illustrated in Figure 12.



Figure 12: Data catalog | Exploring data

The data catalog in the sandbox is using the Azure Cloud platform, since Aarhus Vand already uses Azure Cloud. However, it is important to highlight that the metadata concept is tool independent and can be implemented in any data catalog tool. There are several cloud providers, e.g. Amazon Web Services, Google Cloud Services and Microsoft Azure Cloud. Exploration of the problem is independent of the cloud provider.

In Azure the data catalog is deployed as a part of Microsoft Purview, which is a unified data governance service. Purview Data Catalog provides a system for registration of enterprise data sources and a discovery service for searching data assets. Purview is based on Apache Atlas, an open source metadata management and data governance tool [Purview2022].

A collection of entities from the metadata concept introduced in the last chapter can be presented in a data catalog as shown in Figure 13.



Figure 13: Data catalog | Metadata concept

The data catalog is accessible through a web viewer and the data assets can be found on a data map as shown in Figure 14, where relevant domains are added as collections. The sample data belongs to the collection hierarchy Thesis Purview / SRO data / Drinking water.



Figure 14: Data catalog | Data map

Data assets are ingested during scans and populated in the data catalog. Some of the processes are performed automatically and others processes require manual

interaction to enrich the metadata. Both types of processes are covered in the analysis model as outlined in the next chapter.

5.2 Data lake

To make the architecture reproducible, the data lake is implemented via a cloud service rather than in an on premises data center. The sample data described are ingested. Focus is on the landing zone in the data lake architecture as shown in Figure 15.



Figure 15: Data lake | Zones

In Azure a data lake is implemented as a standard storage account resource with file blob storage. Content of structured and unstructured data is organized in blob containers, where each container represents a zone [Azure2022].

A data lake storage account is implemented in the sandbox environment as displayed in Figure 16. The data lake acts as a baseline with a landing zone including sample data from the distribution network.



Figure 16: Data lake | Sandbox storage accounts

Chapter 6 - Analysis Model

The analysis model used to evaluate the hypothesis is described in this chapter. The goal is to test if domain knowledge can be preserved in a time-saving way into a data lake in order to optimize the user's data discovery. The model consists of a number of Quality Attribute Scenarios (QAS), which each has a description of why it is relevant and an important architectural consideration.

6.1 Usability

Usability is important when the user searches the data catalog, as it should be easy to understand the data context, relationships and track the data lineage. The response to an usability quality attribute is intended to give the user appropriate feedback or assistance. In Software Architecture in Practice, Bass [Bass2012, p. 175] describes usability as:

Usability is concerned with how easy it is for the user to accomplish a desired task and the kind of user support the system provides. Over the years, a focus on usability has shown itself to be one of the cheapest and easiest ways to improve a system's quality.

Metadata can be inferred automatically into the data catalog during a scan. The scanning process establishes a connection to the data source and captures data assets with technical metadata like names, file size and columns. The scan can also apply classification tags on data assets.

However, only a subset of domain knowledge will be captured and ingested. The remaining part of domain knowledge must be added manually to obtain a full set of metadata. Business terms are never inferred automatically, as they are unique to the organization. The following usability scenarios deal with automatic and manual ingestion of metadata.

6.1.1 Automatic infer

This scenario explores how large a subset of domain knowledge can be derived automatically. The data catalog will be used out-of-the box with a minimum of human intervention. The process of automatic infer metadata by running a scan in the data catalog tool is explored in the first usability Quality Attribute Scenario as described in Table 3.

Scenario:		As a data manager, I want to be able to ingest metadata automatically into the data catalog by running a scan.
Quality Attributes:		Usability
Source:		Data manager
s	Stimulus:	Running a scan from the data map.
Part	Artifact:	Baseline data lake and the data catalog.
rio]	Environment:	Runtime
ena	Response:	The scan completes and data assets are ingested.
Sc	Response Measure:	All three files with sample data are ingested as data assets.

Table 3: Quality Attribute Scenario | Usability | Automatic infer

6.1.2 Manual infer

Not all metadata can be inferred automatically in the data catalog. Manually creating metadata allows for more detail and it is possible to add more domain knowledge manually into the data catalog. Table 4 outlines how to do this in data catalog tool and is the second usability Quality Attribute Scenario. The effort has been time-limited to reflect a real-world situation with a business demand.

Scenario:		As a data manager, I want to manually enrich metadata with domain knowledge including business terms and classification tags on assets.
Quality Attributes:		Usability
	Source:	Data manager
rts	Stimulus:	Provide information about business terms, classification and relationship between data assets.
o Pa	Artifact:	Baseline data lake and the data catalog.
lari	Environment:	Runtime
Scer	Response:	Data assets are updated with data types, classification, business terms and lineage. Domain knowledge is searchable in the data catalog.
Response Measure:		The work is completed within 90 minutes of work.

Table 4: Quality Attribute Scenario | Usability | Manual infer

6.1.3 Semi-automatic infer

Obtaining a full set of domain knowledge captured into the data catalog can be very time-consuming. Finding a time-saving way to populate metadata will be explored in the third Quality Attribute Scenario as shown in Table 5. The effort has been time-limited to reflect a real-world situation with a business demand.

Scenario:		As a data manager, I want to import business terms in a time-saving way and tag with custom classification automatically.		
Quality Attributes:		Usability		
	Source:	Data manager		
rts	Stimulus:	Import business terms via interface and define custom classification rules. Run a scan and edit.		
o Pa	Artifact:	Baseline data lake and the data catalog.		
ıari	Environment:	Runtime		
Scer	Response:	Business terms are imported and the scan will tag data assets with custom clarification.		
	Response Measure:	The work is completed within 60 minutes of work.		

Table 5: Quality Attribute Scenario | Usability | Semi-automatic infer

6.2 Performance

Performance is the ability to meet timing requirements for an event. The response to a performance quality attribute is intended to be generated within a time constraint as described by Bass [Bass2012, p. 135]:

The goal of performance tactics is to generate a response to an event arriving at the system within some time-based constraint. The event can be single or a stream and is the trigger to perform computation.

The sample data contains time series data and one of the biggest challenges is that data grows over time. One measurement every minute, 24 hours per day, 365 days per year creates approx 500.000 data rows! Big data is a challenge for relational database management systems (RDBMS). Because of the relations between tables, a selection of data may require a time-consuming table join.

Big data is often generated by the time dimension, but the relation model ignores the order of rows in a table. It is easier to get data in than out of a RDBMS, as queries do not scale linearly with data size[Jacobs2009]. In the baseline data lake, measurements are stored in one big file. Even though there are no relations to other files, extracting data will become increasingly time-consuming. Structuring a directory in the landing zone might improve performance. The directory can be divided into a hierarchy of subfolders with year / month / day.

6.2.1 Implicit infer

The performance scenario investigates how long time it takes to fetch data from the structured directory compared to the landing zone as shown in Table 6. The scenario contains two test periods, both of which include a specific sensorId and a time span. Each test period is run five times and the average execution time is calculated as a response measure, The measure is compared between structured directory and unstructured landing zone. 25% faster is chosen as a realistic goal.

Scenario:		A user fetches data from the landing zone and the structured directory in these periods: #1: 09.06 06:15 - 06:29 #2: 08.06 06:15 - 06:29 + 09.06 06:15 - 06:29 Only sensorId 'DIS.D01_H_B26_BP1-M_TRYK'		
Quality Attributes:		Performance		
	Source:	Data analyst		
rts	Stimulus:	The user fetches the data using Python language interaction and Apache Spark notebook.		
o Pa	Artifact:	Baseline data lake with a structured directory		
lari	Environment:	Normal mode		
Scer	Response:	All data from the period #1 and #2 and the sensorId are retrieved.		
	Response Measure:	Average response time is at least 25% faster from the structured directory compared to the landing zone in both periods.		

Table 6: Quality Attribute Scenario | Performance | Implicit infer

The structuring strategy based on time is about metadata, as the naming of folders implicitly becomes meaningful as displayed in Figure 17.



Figure 17: Data lake | Structured directory

Chapter 7 - Analyses and Results

Based on the sample data and scenarios outlined in the previous chapters, the analyses in this chapter will explore how to deal with the challenges of preserving domain knowledge in a data lake. The analyses will focus on whether or not it is possible to optimize the data discovery.

To document the process, a simple template is set up in each section with a development logbook that contains steps, issues, solutions, learnings and results for each scenario. The purpose of the logbooks is to document each step in the Quality Attribute Scenarios and make them reproducible. As a guide for navigating through the scenarios, each analysis section contains a process map showing the various steps in the infer process. The maps are based on the metadata concept and user actions are marked as yellow events.

7.1 Automatic infer

This analysis explores the scenario described in the Quality Attribute Scenario | Automatic infer. The baseline data lake in the sandbox environment is the foundation for the work. The goal is to have metadata automatically registered as data assets, enabling data discovery by searching the catalog. As a starting point, the data catalog is empty and contains no data assets. The automatic infer user actions are illustrated in Figure 18 and includes one event (yellow box):

1. Scan the baseline data lake



Figure 18: User actions | Automatic infer

7.1.1 Logbook



Date	Subject		
2022-11-05	Register baseline data lake as a source to 'Drinking Water' collection. Setup scan connector, use default scan rule set and run the scan. Three assets candidates for classification and schema are extracted. Search the data catalog for 'sensorId' and get two results.		
		Usability Issues	
Issue		Solutions and learnings	
Schema insufficient		The schema inferred automatically is not sufficient. Only column names, but not descriptive tags, are derived. All data types are set to strings, since the source files are csv. This must be handled manually.	
Metadata types		The automatic out-of-the-box scanning only discovers technical metadata. There is a lack of business metadata, which has to be provided manually.	
Classification missing		The default scan rule set includes all supported file types for schema extraction and classification. The system classifications rules look for passport and credit card numbers, date of birth, email addresses and passwords and so on in order to tag the data assets. However, the sample data in the baseline data lake contains domain specific data, hence none of the data sets are classified. Custom classifications rules need to be created.	
Schema registry		The schema is automatically inferred. If the schema is changed before the next scan, the column names of the data asset will be overridden without any warning (see Appendix C). It is actually possible to be alerted during ingesting via Azure event hub and messaging services to be aware about the change [AzureEvent2022]. To be able to use messaging in the data catalog, event hubs have to be enabled in the settings. This is not investigated further as it is out of project scope.	
Relation not clear		The three data assets are marked as related, but it is not stated which properties link the files together. The data manager must discover the lineage himself.	

Screenshots				
Additional screenshots are added in Appendix C. Starting point:				
🕞 No sources 🖩 No assets 🛄 No glossary terms				
After initial scan:				
🕞 1 source 🖩 6 assets) 🛄 No glossary terms				
Data discovery by searching the data catalog for 'sensorId'.				
Data catalog >				
Search results for sensorid				
Source type : all 🛛 😵 Clear all filters				
Tilter by keyword Showing 1-2 out of 2 results				
Object Type Measurement.csv O Dashboards Azure Data Lake Storage Gen2 File D Data pipelines https://mthbaselinedls.dfs.core.windows.net/landing/Measurement.csv Files Sensor.csv				
Glossary terms Azure Data Lake Storage Gen2 File https://mthbaselinedls.dfs.core.windows.net/landing/Sensor.csv				
Image: Image				
Data ratalon 1. Search results "censorid" 1.				
Measurement.csv Azure Data Lake Storage Gen2 File				
	er Bl Desktop			
Overview Properties Schema Lineage Contacts Related Updated on November 6, 2022 1:21 PM b	y automated scan			
▼ Filter by name				
Showing 5 of 5 items				
Column name Classifications Sensitivity label Glossary terms	Data type			
DateTime	string			
Sensorid	string			
Measure	string			
UpdateTime	string			
Filed	string			
Results				
The automatic scan is capable of ingesting all data assets. But the inferred schemas are not efficient and no business metadata tags are obtained. The subset of domain knowledge preserved in the data lake is small. However the user action is conducted fast and the domain knowledge can be discovered in the data catalog.				

Table 7: Logbook | Usability | Automatic infer

7.2 Manual infer

The previous Quality Attribute Scenario pinpointed that the automatic out-of-the-box scanning only creates a small subset of domain knowledge with some missing elements such as:

- column data types not specified
- no classification tags of patterns
- business metadata missing
- lineage with relations between files not explicit

This analysis explores how to ingest the missing elements as described in the Quality Attribute Scenario | Manual infer. The baseline data lake in the sandbox environment is the foundation for the work. As a starting point, the data catalog is populated with three data assets from the automatic scan.

The user actions are illustrated in Figure 19 and includes the following events (yellow boxes):

- 1. Create business terms in the business glossary.
- 2. Create custom classification with description.
- 3. Edit data assets to update descriptions, tags and certified labels.



Figure 19: User actions | Manual infer

7.2.1 Logbook

QAS usability Manual infer				
The data manager enriches metadata manually with domain knowledge and updates data assets with data types, classification, business terms and lineage. Domain knowledge is searchable in the data catalog. The work should be completed within 90 minutes of work.				
Development Logbook				
Date	Subject			
2022-11-07	Create business glossary terms in the data catalog.			

	Create custom classifications for columns in data assets. Relate business glossary to assets, change data type, add lineage and finally label enriched data assets as 'Certified'. Search the data catalog for 'sensorId' and get two results. Time consumption: 125 minutes			
		Usability Issues		
Issue		Solutions and learnings		
Glossary taxonomy		Defining business terms cannot be inferred from data assets. It is time consuming to enter this information. Metadata extraction from Figure 9 and domain expert knowledge in Table 2 is captured as 'Sensor ID template' with related terms (see Appendix D).		
Edit data assets		Tagging fields is done easily by picking glossary terms from a dropdown menu. Data type is a free text field, which may result in incorrect values. By adding lineage the relation between files becomes explicit.		
Classification		After creating custom classification they are manually added to data asset fields in the same way as glossary terms.		
Enable trust for the user		When data asset information is validated, it can be labeled with 'Certified' and regarded as reliable.		
Data discove	ry	Business glossary terms and custom classification can be filtered by when searching the data catalog.		
		Screenshots		
Additional screenshots are added in Appendix D. After creating business glossary terms: → 1 source → 1 source				
Lineage: Sensor.csv Path Measurement.csv Sensorconfig.csv				

Custom classification:					
Display name	Formal name	Description			
Aarhus Vand	Aarhus Vand	Assest owned by Aarhus Vand			
Configuration	Configuration	Metadata for consuming a time serie file			
Delimiter	Delimiter	Column delimiter			
Description	Description	Explanation			
File name	File name	Name of file			
FileID	FileID	Foreign key to file			
Location	Location	Geographical directions			
Measure	Measure	Value from time serie recorded on created datetime			
Media	Media	Utility type: Drinking water Rain water Resource water			
Parameter	Parameter	Measure parameter			
SensorID	SensorID	Unique sensor template: ID Node Segment Number Measure type Sensor type IsMid.			
Sensors	Sensors	Wireless Internet of Things (IOT) sensor			
Separator	Separator	Separator between measurement values			
Time serie	Time serie	Time serie data from wireless Internet of Things (IOT) sensor			
Time serie type	Time serie type	Source of time serie			
Timestamp	Timestamp	Timestamp with date and time			
Unit	Unit	Mathematical unit			
Zone	Zone	Section in the distribution network			
		Results			

All metadata are added by manual processes, including data types, tagging, business metadata and relation between files. The metadata from the data model of the sample data is also applied and even the domain expert knowledge in a searchable way. In fact, all domain knowledge is preserved in the data lake in a searchable way. But the process is very time consuming and infeasible in a real-world scenario with several data assets. The time limit of 90 minutes of work was not met, as the scenario took 125 minutes to conduct.

Table 8: Logbook | Usability | Manual infer

7.3 Semi-automatic infer

This analysis explores the scenario described in the Quality Attribute Scenario | Semi-automatic infer. The analysis will explore how to properly define business terms, so it can be imported in a time-saving manner.

The user actions are illustrated in Figure 20 and includes the following events (yellow boxes):

- 1. Import business terms in the business glossary using a template
- 2. Create custom classification with a pattern defined as regular expression.
- 3. Add scan rule using the created custom classifications and run scan
- 4. Edit data assets to update description, check tags and set certified labels.



Figure 20: User actions | Semi-automatic infer

7.3.1 Logbook

QAS usability Semi-automatic infer				
The data manager imports business terms in a time-saving way. After defining custom classification rules, the tagging is done automatically by running a scan. The work should be completed within 60 minutes of work.				
		Development Logbook		
Date	Date Subject			
2022-11-15	Fill out template and bulk import business terms in the catalog. Create custom classification rules and define data patterns with a regular expression or column patterns by name. Add scan rule set using the new custom classification rules. New scan with custom scan rule set. Edit data assets to update description, terms and certified labels Time consumption: 85 minutes			
		Usability Issues		
Issue	Issue Solutions and learnings			
Business terms import		Import can be done using either REST API endpoints for create, read, update or delete actions or a term template in .csv format [PurviewAPI2022]. In this scenario, the .csv template is used. After a few trials with validation errors, it was very quick to import the total glossary. The import template is added to the code repository placed at Azure DevOps as described in Appendix B.		

Classification rules		Defining data patterns with a regular expression to match the data stored in a data field takes a little practice. Using <u>https://regex101.com</u> (online regex builder), a regex for the unique sensor ID template is created. But for the untrained it takes time to build a solid regex pattern. Classification rules can also label sensitive data in order to restrict access. However, this is not in scope of this project.					
Scan with new rule set		A new scan with the custom scan rule set attached should tag the data assets with classifications according to the classification rule - the first few times without result! After refactoring the classification rules and using the 'test classification rule' option, the assets were tagged. This is much more time-saving than re-running a scan.					
Edit		Need to do taggin	g of imp	orted b	ousiness term	is ma	nually.
		Screensh	ots				
Additional screenshots are added in Appendix E. After creating business glossary terms:							
Template for importin	ng bus	<u>iness terms:</u>					
Name	Definitio	n	*	Status 💌	Parent Term Name 🔽	Term 1	[Attribute
_Created date time	Created o	late time		Draft		Thesis;	Import
_Metro maps	Function	map showing sensor locations		Draft		Thesis;	Import
_IOT sensor measurement	Measurer	nent from IOT sensor		Approved		Thesis;	Import
_Sensor ID template_Segment	Subset of	Scada network		Approved	_Sensor ID template	Thesis;	Import
_IUT sensor	Consocut	internet of Things (IOT) sensor	ific comont	Approved	Sonsor ID tomplate	Thesis;	Import
Sensor value	Sensor m	easurement	ine segment	Draft	_sensor to template	Thesis:	Import
File reference	Relation	to file		Draft		Thesis;	Import
Sensor ID template_Node	Scada net	vwork		Approved	_Sensor ID template	Thesis;	Import
_Utility type	Types: 	r≫ul≫li >Drinking water	Rain water <td>> Draft</td> <td></td> <td>Thesis;</td> <td>Import</td>	> Draft		Thesis;	Import
_Config file	Settings f	or using a measurement		Approved	-	Thesis;	Import
Sensor ID template_Measure type	Posible v	alues: M: measurement<!--</td--><td>li>Z: calcula</td><td>Approved</td><td>_Sensor ID template</td><td>Thesis;</td><td>Import</td>	li> Z: calcula	Approved	_Sensor ID template	Thesis;	Import
_sensor ID template_Isivita	Timestan	aiue: <ui>ii>irue: minute avera</ui>	ge of approx. 2	Draft	_sensor iD template	Thesis;	Import
Sensor ID template Sensor type	Posible v	alues: FlowPressu	ure Tem	Approved	Sensor ID template	Thesis:	Import
Sensor ID template	Unique se	ensor ID using template with info	rmation about	Approved	_	Thesis;	Import
_Sensor ID template_ID	Only ID o	f the Seonsor ID template		Approved	_Sensor ID template	Thesis;	Import
Import validation: Column Guidance Issue to resolve [Attribute][Thesis]Source* The field is required. Select a value from: I This required column is missing from the template, please mo and try again.			dify the file				

Custom classification rule Pattern of SensorId e.g. OST.B09_E_G1_BF1-M_FLOW:				
Data pattern 🛈				
(^[a-zA-Z]{3}\.[a-z	zA-Z][0-9]{2}_[a-:	zA-Z]_[a-zA-Z][0-9]+_[a-zA-Z]{2}[0-9]-[a-zA-Z]_[a-zA-Z]{4}\$ +		
Minimum match th	nreshold 🛈	0% 100% 70% 🗘		
<u>Classification tag</u>	g added autor	matically by scan:		
	AAV.SensorID			
Schema classifi	Applied by	Scan on November 15, 2022 12:56 AM		
(AAV.SensorID)	Туре	Custom		
	Sample count	11526		
Fully qualified r	Distinct count	134		
		Results		
The scan results in tagging the columns with classifications. Importing the business terms makes it fast to tag data assets with business knowledge. The time limit of 60 minutes of work was not met, as the scenario took 85 minutes to conduct. The time-consuming part was creating regular expressions. One way to make the process faster is to use column patterns with names, but it is not as robust as it does not look at the actual data content. Over time, the regular expressions will be reusable for the users, thus skipping the overhead.				

Table 9: Logbook | Usability | Semi-automatic infer

7.4 Implicit infer

This analysis explores the scenario described in the Quality Attribute Scenario | Implicit infer. A structured directory divided into a hierarchy of subfolders with year / month / day is added to the baseline data lake. The user fetches the data using Python programming language into an Apache Spark notebook running on Azure. The cluster attached to the notebook has 14 GB memory and 4 cores. The user actions are illustrated in Figure 21 and includes one event (yellow box):



1. Fetch data from source into an Apache Spark notebook

Figure 21: User actions | Semi-automatic infer

7.4.1 Logbook

QAS performance | Implicit infer

The data analyst fetches data from a structured directory and the landing zone in the period #1: 09.06 06:15 - 06:29 and #2: 08.06 06:15 - 06:29 + 09.06 06:15 - 06:29. Each test period is run five times and the average execution time is calculated. The goal for the structured directory is to be 25% faster than the landing zone.

Development Logbook

Date	Subject				
2022-11-22	 Period #1: 09.06 06:15 - 06:29 Fetch data from period #1 into a dataframe using 'loc' which selects data by columns. Average faster execution time: 58,6% 				
2022-11-23	 Periode #2: 08.06 06:15 - 06:29 + 09.06 06:15 - 06:29 Two files from the structured directory need to be combined in order to find the proper results. Average faster execution time: 69,8% 				
Performance Issues					
Issue	Solutions and learnings				
Combine files	After some failed attempts, the Python 'glob' module solves the challenge [GlobModule2022]. Glob patterns specify sets of filenames with wildcard characters similar to regular expressions [GlobPattern2022]. This enables iterating over the structured directory and loading only files of interest. This glob pattern find filenames from all month in year 2022 on days 08+09: .csv('/mnt/landing/structured (year=2022/month=*/day={0[8-9]})')				
SQL-like search	A Python panda dataframe loads all data to the notebook before selecting the result. To avoid this, the PySpark library processes and combines data before they are loaded to the notebook [PySpark2022]. The PySpark library applies spark.sql() to run arbitrary SQL-like queries on a dataframe saved as a table. Each row in the dataframe gets a filepath added: filepath dbfs:/mnt/landing/structured/year=2022/month=06/day=09. Measurement.csv				





Results						
Test period #1	Structured (ms)	Unstructured (ms)	Difference (%)			
Run #1	110	350	68,6			
Run #2	160	300	46,7			
Run #3	110	290	62,1			
Run #4	120	300	60,0			
Run #5	150	340	55,9			
Average	130	316	58,6			
Test period #2	Structured (ms)	Unstructured (ms)	Difference (%)			
Run #1	200	600	66,7			
Run #2	200	560	64,3			
Run #3	140	610	77,0			
Run #4	160	560	71,4			
Run #5	180	590	69,5			
Average	176	584	69,8			

Table 10: Logbook | Performance | Implicit infer

7.5 Discussion

To get an overview, all the results found in the above analyses based on the Quality Attribute Scenarios formulated in the analysis model, are gathered together and discussed in this section. The purpose of the analyses is to test if *domain knowledge* can be preserved in a *time-saving* way into a data lake in order to optimize the user's *data discovery*. Each analysis has its starting point in the *metadata concept*. Therefore, the overview is divided according to the highlighted words in bold and placed on a scale, where 'High' is best.

7.5.1 Usability

The usability quality is investigated through three Quality Attribute Scenarios with different user actions. The response measure for the scenarios indicates how easy it is for the user to accomplish the scenarios as described in Table 11.

Metadata o	concept				
Entities	Туре	Automatic	Manual	Semi	Legend
Data asset	Asset				Full
Column	Technical				Partly
Lineage	Technical				None
Business term	Business				
Classification	Business				
	Total scale	Low	High	Medium	
Domain knowledge					
Time consumption					
Data discovery					
Total scale		Low	Medium	High	

Table 11: Results | Usability

Looking at the metadata concept, automatic infer is at the low end of the scale, while manual infer is at the high end and in the middle is semi-automatic infer. In terms of optimizing the data discovery, the picture is different: at the high end is semi-automatic infer, in the middle is manual infer while automatic infer still is at the low end. The above result table is based on the result section from the logbooks:

Automatic infer

The automatic scan is capable of ingesting all data assets. However, the inferred schemas are not efficient and no business metadata tags are

obtained. The subset of domain knowledge preserved in the data lake is small. However the user action is conducted fast and the technical metadata is searchable and can be discovered in the data catalog.

Manual infer

All metadata are added by manual processes, including data types, tagging, business metadata and relation between files. The metadata from the data model of the sample data is also applied and even the domain expert knowledge in a searchable way. In fact, all domain knowledge is preserved in the data lake in a searchable way. But the process is very time consuming and infeasible in a real-world scenario with several data assets. The time limit of 90 minutes of work was not met, as the scenario took 125 minutes to conduct.

Semi-automatic infer

The scan results in tagging the columns with classifications. Importing the business terms makes it fast to tag data assets with business knowledge. The time limit of 60 minutes of work was not met, as the scenario took 85 minutes to conduct. The time-consuming part was creating regular expressions. One way to make the process faster is to use column patterns with names, but it is not as robust as it does not look at the actual data content. Over time, the regular expressions will be reusable for the users, thus skipping the overhead.

7.5.2 Performance

Performance quality is explored through one Quality Attribute Scenario with two test periods. The results are evaluated by comparing the average execution time in milliseconds from the structured directory versus the unstructured landing zone as outlined in Table 12.



Metadata concept				
Source	Structured	Unstructured	Difference	
Period #1	130 ms	316 ms	58,6 %	
Period #2	176 ms	584 ms	69,8 %	
Total scale	High	Low		Legend
Domain knowledge				Full
Time consumption				Partly
Data discovery				None
Total scale	High	Low		

Table 12: Results | Performance

Looking at the metadata concept, fetching data from the structured directory is at the high end of the scale, unstructured landing zone is at the low end. In terms of optimizing the data discovery, the picture is the same: at the high end is structured directory while unstructured landing zone still is at the low end. The above result table is based on the logbook:

Structured directory

Selection from the structuring directory based on time provides the user with useful knowledge about when the data is created. Combined with glob pattern, the fetching process only loads files of interest and serves the proper file paths. This approach is extremely efficient, very scalable and offers the user a quick data discovery with finding data. In addition, the user benefits from the SQL-like syntax. It is also possible to introduce additional parameters in the glob pattern, e.g. sensorId.

Unstructured landing zone

No implicit knowledge is offered from the unstructured landing zone. When storing all the measurements in one large file in the landing zone, the user must load the entire file to fetch the desired data. This is inefficient and the performance degrades significantly over time, as time series data grows over time. Thus the data discovery will be worse. Querying data from the dataframe columns is more time consuming to write than the Pyspark SQL-like syntax.

Chapter 8 - Conclusion

The purpose of this project is to investigate how domain knowledge can be preserved in a data lake to improve the user data discovery. Managing domain knowledge can be a time-consuming process and seems like a difficult task. Therefore, the project intends to answer the following problem statement:

Can domain knowledge be preserved in a time-saving way into a data lake in order to optimize the user's data discovery?

In order to answer the problem statement, a hypothesis is set up, which involves constructing a metadata concept and an analysis model with scenarios. The concept provides the user with an overview of the different metadata entities. The hypothesis is evaluated using a small sandbox environment with sample data.

By constructing a custom analysis model, where automatic, manual, semi-automatic and implicit ingesting of metadata are explored, it can be evaluated whether domain knowledge in the data lake optimizes the data discovery or not.

The hypothesis is verified, as it is possible to evaluate the extent to which *domain knowledge* can be preserved in a *time-saving* way into a data lake in order to optimize the user's *data discovery*.

Usability

The result of the usability part of the analysis model is visualized using a scale. At the high end is semi-automatic infer, in the middle is manual infer while automatic infer is at the low end. For each infer type, the domain knowledge and time consumption is analyzed to determine whether data discovery has improved.

The best optimization of data discovery is obtained by manual and semi-automatic infer. The most time-saving method is semi-automatic.

However, it should be noted that automatic infer in some cases can be a good choice, as it is possible to create searchable technical metadata very quickly. This is a much better starting point than having no information in the data catalog.

Performance

For the performance part of the analysis model, the result is very significant. Fetching data from a structured directory versus the unstructured landing zone is more than twice as fast. Looking at the metadata concept, fetching data from the structured directory is at the high end of the scale, unstructured landing zone is at the low end. Combined with the glob pattern, the structured directory performs great. The faster data can be retrieved, the faster discovery.

The structured directory with implicit domain knowledge is extremely efficient, very scalable and offers the user a quick data discovery.

Depending on data use, it can be considered to introduce additional parameters in the glob pattern, e.g. sensorId.

Answer to Problem Statement

Based on the results from the usability and performance part of the analysis model, the following answer to the problem statement is given:

Based on the usability and performance scenarios, it can be concluded that it is possible to preserve domain knowledge in a time-saving way and thus optimize the data discovery. The best results are obtained with semi-automatic infer and storage of data in a structured directory.

8.1 Reflections

The work to investigate the problem statement has been a learning journey in how to handle some real-world challenges for a data-intensive application in a data lake. The journey has not been straightforward, thus it is valuable to highlight some learnings.

Key takeaways:

- One of the most important learnings has been the importance of developing a concept that is both tool, platform and domain independent.
- Another learning is the value of enabling the user to become self-serviced.
- The metadata concept makes the task of implementing a data catalog and building a metadata fundament manageable and accessible.
- The concept guides the user to navigate through the implementation of metadata in the data catalog, classify the data and add business knowledge.
- Users can expand the metadata concept according to their experiences.
- Costs of running a data catalog in Azure are billed per capacity hour unlike other cloud resources.

Bibliography

[Atlas2022] Apache Atlas. *Metadata Management and Data Governance*. https://atlan.com/what-is-apache-atlas/. Accessed 25.10.2022.

[Azure2022] Microsoft Learn. *Azure Blob Storage documentation*. https://learn.microsoft.com/en-us/azure/storage/blobs/. Accessed 08.10.2022.

[AzureEvent2022] Microsoft Learn. *Azure Schema Registry in Azure Event Hubs*. https://learn.microsoft.com/en-us/azure/event-hubs/schema-registry-overview. Accessed 05.11.2022.

[Bass2012] Bass, Len, Paul Clements, and Rick Kazman. *Software architecture in practice*. SEI Series in Software Engineering, 2012.

[Christensen2020] Christensen, Henrik Bærbak. *Quality Attribute Scenarios* (*Slide deck*). Software Architecture in Practice (Master of IT course), 2020.

[GlobModule2022] Towards Data Science. *The Python Glob Module*. <u>https://towardsdatascience.com/the-python-glob-module-47d82f4cbd2d</u>. Accessed 22.11.2022.

[GlobPattern2022] GeeksforGeeks. *How to use Glob function to find files.* <u>https://www.geeksforgeeks.org/how-to-use-glob-function-to-find-files-recursiv</u> <u>ely-in-python/</u>. Accessed 23.11.2022.

[Gorelik2019] Gorelik, Alex. *The enterprise big data lake: Delivering the promise of big data and data science*. O'Reilly Media, 2019.

[Hamadou2020] Hamadou, Hamdi Ben, Torben Bach Pedersen, and Christian Thomsen. "The Danish National Energy Data Lake." *2020 IEEE International Conference on Big Data (Big Data)*. IEEE, 2020.

[Han2012] Han, Jiawei, Jian Pei, and Hanghang Tong. *Data mining: concepts and techniques*. Morgan Kaufmann, 2012.

[Jacobs2009] Jacobs, Adam. "The pathologies of big data." *Communications of the ACM* 52.8 (2009): 36-44.

[Nwagu2017] Nwagu et al. "Knowledge Discovery in Databases (KDD): an overview." *Int J Comput Sci Inf Secur (IJCSIS)* 15.12 (2017): 13-16.

[Purview2022] Microsoft Learn. *Microsoft Purview governance documentation*. https://learn.microsoft.com/en-us/azure/purview/. Accessed 08.10.2022. [PurviewAPI2022] Microsoft Learn. *Microsoft Purview Resource Rest API*. https://learn.microsoft.com/en-us/rest/api/purview/. Accessed 08.10.2022.

[PySpark2022] Apache Spark. *PySpark Documentation*. https://spark.apache.org/docs/latest/api/python/. Accessed 23.11.2022.

[Thomsen2022] Thomsen, Christian. *Data warehousing & OLAP (Slide deck)*. Business Intelligence - Analysis of large databases (Master of IT course), 2022.

[Videla2021] Videla, Alvaro. "Meaning and Context in Computer Programs: Sharing domain knowledge." *Queue* 19.5 (2021): 60-68.

Appendix

A. Production environment | Size

The data lake was started up in Azure in September 2021 where a small amount of sensor measurements were uploaded to the landing zone. Since November 2021 and until the summer of 2022, measurements from 45 sensors have been uploaded. It corresponds to approx. 6 million measurements every month (60 minutes x 24 hours x 31 days x 45 sensors x 3 types).

Every hour measurements per minute are uploaded for the last six hours. The size of files placed in the landing zone are 0,7 MB and contain 8.100 records.

The uploaded files are handled in a data factory, where an Extract - Transform - Load (ETL) process appends delta measurements to the files in the publish zone with raw data. The current size of files placed in the publish zone are 3,2 GB and contain approx. 52.000.000 records.

A1 Landing zone

	Microsoft Azure	٩	Search resources, service	s, and docs (G+/)						
Home	aavshareddlsprod Containers >									
	landingzone … ^{Container}									
>>	↑ Upload + Add Directory ひ Refresh	🖓 Rename 📋 De	elete 🛱 Change tier	တ် Acquire lease တိ	g Break lease					
(Authentication method: Access key (Switch to Azure AD User Account) Location: landingzone / SRO Search blobs by prefix (case-sensitive)									
	Name	Modified	Access tier	A. Blob type	Size	Lease state				
	🔲 🛅 []									
	Flows_Z105Syd.csv	6/29/2022, 2:26:42 PI	M Hot (Inferred)	Block blob	224.61 KiB	Available				
	Flows_Z125Syd_Z80Beder.csv	6/29/2022, 2:26:32 PI	M Hot (Inferred)	Block blob	94.5 KiB	Available				
	Flows_Z80Syd.csv	6/29/2022, 2:26:51 PI	M Hot (Inferred)	Block blob	181.61 KiB	Available				
	🗌 📄 FlowTrykTemperatur_Z125N.csv	6/29/2022, 2:26:13 PI	M Hot (Inferred)	Block blob	145.89 KiB	Available				
	🗌 📄 Temperatur_Z105Syd.csv	6/29/2022, 2:25:44 PI	M Hot (Inferred)	Block blob	160.62 KiB	Available				
	Temperatur_Z125Syd_Z80Beder.csv	6/29/2022, 2:26:54 PI	M Hot (Inferred)	Block blob	70.82 KiB	Available				
	🗌 📄 Temperatur_Z80Syd.csv	6/29/2022, 2:26:17 PI	M Hot (Inferred)	Block blob	116.13 KiB	Available				
	Tryk_Z105Syd.csv	6/29/2022, 2:25:32 PI	M Hot (Inferred)	Block blob	192.47 KiB	Available				
	Tryk_Z125Syd_Z80Beder.csv	6/29/2022, 2:25:49 PI	M Hot (Inferred)	Block blob	83.41 KiB	Available				
	Tryk_Z80Syd.csv	6/29/2022, 2:25:57 PI	M Hot (Inferred)	Block blob	133.33 KiB	Available				
					`					

A2. Gold zone

	Microsoft Azure	∠ Search reso	ources, services, and	docs (G+/)						
Hom	Home > aavshareddlsprod Containers >									
	En publish … Container									
>>	🕂 Upload 🕂 Add Directory 💍 Refresh │ 🤇 Rename 🛍 Delete 🔁 Change tier 🔗 Acquire lease 🖉 Break lease									
C	Authentication method: Access key (Switch to Azure AD User Account) Location: publish / Raw / SRO									
	Search blobs by prefix (case-sensitive)									
	Name	Modified	Access tier	A. Blob type	Size	Lease state				
	📄 📙 []									
	Staging					-				
	📄 📄 sensorconfig_old.csv	10/15/2021, 12:00:26 PM	Hot (Inferred)	Block blob	177 B	Available				
	📄 📄 sensorconfig.csv	10/15/2021, 12:01:04 PM	Hot (Inferred)	Block blob	387 B	Available				
	TotalSRO_Flows_Z105Syd.csv	6/29/2022, 2:40:07 PM	Hot (Inferred)	Block blob	575.62 MiB	Available				
	TotalSRO_Flows_Z125Syd_Z80Beder.csv	6/29/2022, 2:35:43 PM	Hot (Inferred)	Block blob	196.77 MiB	Available				
	TotalSRO_Flows_Z80Syd.csv	6/29/2022, 2:38:46 PM	Hot (Inferred)	Block blob	455.02 MiB	Available				
	🗌 📄 TotalSRO_FlowTrykTemperatur_Z125N.csv	6/29/2022, 2:36:10 PM	Hot (Inferred)	Block blob	329.69 MiB	Available				
	🗌 📄 TotalSRO_Temperatur_Z105Syd.csv	6/29/2022, 2:35:43 PM	Hot (Inferred)	Block blob	349.72 MiB	Available				
	🗌 📄 TotalSRO_Temperatur_Z125Syd_Z80Beder.csv	6/29/2022, 2:33:44 PM	Hot (Inferred)	Block blob	126.71 MiB	Available				
	🗌 📄 TotalSRO_Temperatur_Z80Syd.csv	6/29/2022, 2:34:32 PM	Hot (Inferred)	Block blob	233.88 MiB	Available				
	Distalseries Totalseries Total	6/29/2022, 2:36:30 PM	Hot (Inferred)	Block blob	431.98 MiB	Available				
	TotalSRO_Tryk_Z125Syd_Z80Beder.csv	6/29/2022, 2:33:09 PM	Hot (Inferred)	Block blob	166.33 MiB	Available				
	TotalSRO_Tryk_Z80Syd.csv	6/29/2022, 2:36:08 PM	Hot (Inferred)	Block blob	293.02 MiB	Available				

B. Sandbox environment | DevOps

The infrastructure behind the sandbox environment can be described as code using the Azure Resource Manager templates and is available in the project's code repository placed at Azure DevOps. Sample data, python notebooks and the business terms template are also uploaded to the repository.

Until the exam at 16th December 2022 there is granted anonymous access to the DevOps: <u>https://dev.azure.com/mth-thesis/Data-lake-sandbox/_git/Sandbox</u>

B1. Repository content

	zure DevOps mth-thesis / Data-la	ake-sandbo	x / Repos / Files / 🚸 Sandbox 🕥						
D	Sandbox	:	😢 main 🗸 🗈 / Type to find a file or folder						
+	> 🖿 4.2 Sample data		Files						
	🗧 🖿 5.1 Data catalog		Contents History						
	> 🖿 5.2 Sandbox data lake								
	✓ ■ 6.2 Structered directory		Name 1	Last change					
8	∨ 🖿 year=2022		🖿 4.2 Sample data	lørdag					
đ	∨ 🖿 month=06		🖿 5.1 Data catalog	lørdag					
¢	> 🖿 day=06		🖿 5.2 Sandbox data lake	lørdag					
ഭ	> 🖿 day=07		6.2 Structered directory	14h ago					
દુષ્	> 🖿 day=08		7.3 Business terms template	lørdag					
0	> 🖿 day=09		7.4 Apache spark	lørdag					
11	> 🖿 day=10		M↓ README.md	lørdag					
	> 🖿 day=11								
	> 🖿 day=12		Folder content						
-	🗧 🖿 7.3 Business terms template		 Section 4.2 Sample data: .csv files Section 5.1 Data catalog: ARM template 						
	🗸 🖿 7.4 Apache spark		Section 5.2 Sandbox data lake: ARM template						
	> 🖿 Databricks		 Section 6.2 Structured directory: .csv files Section 7.3 Business terms template: import .csv 	file					
	> 🖿 Python notebooks		 Section 7.4 Apache spark Databricks: ARM template 						
	MI README.md		Python notebooks: python files						

C. Automatic infer | Screenshots

C1. Initial scan

Azure Data Lake Storage Gen2				
🎖 New scan 🧷 Edit source 🧯] Delete source 🜔 Ref	sch		
Overview Scans		An asset where Microsoft Purview extracts schema and applies classifications during an automated scan. The scan rule set determines which assets get classified. If the asset is considered candidate for classification and no classifications are applied with the state of the scalar scal	da .	
Source ID: https://mthbaselinedls.	dfs.core.windows.net/	during scan time, an asset is still considered a classified asset.		Registered on
Scans	Discovered assets ①	Classified assets ①		11/06/2022, 1:13:25 PM
C	5	3		
'	5	3		Collection path
				Thesis Purview
Recent scans				🗗 SRO data
Scan name	Last run statu	s Scan rule set	Last scan time	Drinking water
Scan-initial	📀 Complete	d	11/06/2022, 1:18 PM	
See all applied scans				Source hierarchy
				Azure OneWater AaV Dev Subscription
Recent failed scans				
-	Ctatus	Suggestions		aav-thesis-rg-dev

C2. Data discovery | Browse assets

Browse assets	
🕐 Refresh	
By collection By source type	
< View collection tree	Source type : all
Drinking water Thesis Purview > SRO data >	Showing 1-6 out of 6 results
Sub collection(s)	anding
Related	🖷 Azure Data Lake Storage Gen2 File System
SBO data (Parent)	https://mthbaselinedls.dfs.core.windows.net/landing
Resource water	
	Measurement.csv
Narrow results by:	Azure Data Lake Storage Gen2 File
▼ Object Type	https://mthbaselinedls.dfs.core.windows.net/landing/Measurement.csv
O Dashboards OD Data pipelines Diffes Files Folders	Azure Storage Account https://mthbaselinedls.core.windows.net
Reports	🚔 . mthbacalinadic
Stored procedures	Bi Azure Data Lake Storage Gen2 Service
Tables	https://mthbaselinedls.dfs.core.windows.net
▼ Classification ···	Azure Data Lake Storage Gen2 File https://mthbaselinedis.dfs.core.windows.net/landing/Sensor.csv
* Contact	Sensorconfig.csv Azure Data Lake Storage Gen2 File

C3. Data discovery | Asset overview

Microsoft Purview	العند العندي ا	🧈 🛚 🕫 🧼 ?
~	Data catalog > Browse assets >	
📑 Data catalog	Measurement.csv	
🚯 Data map		
💕 🖁 Data share (preview)	Overview Properties Schema Lineage Contacts Related	Updated on Novembe
Data estate insights	Asset description	Collection path
Data policy (preview)	No description for this asset.	Thesis Purview
📑 Management	Managed attributes (preview)	🗗 SRO data
-	No Attributes for this asset.	Drinking water
	Classifications $^{\odot}$	Hierarchy
	No classifications for this asset.	Azure Storage Account
	Schema classifications $^{\odot}$	Azure Data Lake Storage Gen2 Service
	No classifications for this asset.	landing
	Fully qualified name $^{\odot}$	Azure Data Lake Storage Gen2 File System
	https://mthbaselinedis.dfs.core.windows.net/landing/Measurement.csv	Azure Data Lake Storage Gen2 File

C4. Schema registry after change

Data map > Sources > AzureDataLakeS	torage-vBL >					Showing 4 of 4 items Full scan		Showing 5 of 5 items Incremental scan
Scan-initial run history					ĺ	Column name		Column name
Thesis Purview > SRO data > Drinking w	iter					FileId		FileId
Run scan now V Bedit scan	Delete scan 🚫 🤇	Cancel runs 💍	Refresh ≡≡ R	Edit columns		FileName		FileName
Incremental scan	can assets that have beer run.	n modified or creat	ted since the			DecimalSeparator	\longrightarrow	Type New
Full scan ^① Statu	Run type	Scan type	Assets discov	ered Assets ingested		ColumnDelimiter		DecimalSeparator
ce49e1cb-bd6a-47e7-9cfc-bfa 🛛 🛇 C	ompleted Incremental	Manual	5	З	l			Separator Changed
bde26dd2-458a-418d-b63c-7. 🥑 C	ompleted Full scan	Manual	5	5			l	Changed

D. Manual infer | Screenshots

D1. Business glossary terms

Glossary terms										
+ New term 😤 Manage term templates 🖉 Edit ↔ Import terms 🗁 Export terms 🛍 Delete 🖒 Refresh 🛛 🚺										
▼ Filt	Term template : All Status : All Contact : All Term template : All									
Showing	17 terms Collapse all									
_										
	Name	Term template	Status	Definition						
	Config file	System default	Approved	Settings for using a measurement						
	Created date time	System default	🗋 Draft	Created date time						
	File reference	System default	🗋 Draft	Relation to file						
	IOT sensor	System default	Approved	Wireless Internet of Things (IOT) sensor						
	IOT sensor measurement	System default	Approved	Measurement from IOT sensor						
	Metro maps	System default	🗋 Draft							
	∨ Sensor ID template	System default	Approved	Unique sensor ID using template with information abo						
	ID	System default	Approved	Only ID of the Seonsor ID template						
	IsMid	System default	Approved	Boolean value: True: minute average of approx. 240 me						
	Measure type	System default	Approved	Posible values:M: measurementZ: calculated						
	Node	System default	Approved	Scada netvwork						
	Number	System default	🗋 Draft	Consecutive numbering of sensors in specific segment						
	Segment	System default	Approved	Subset of Scada network						
	Sensor type	System default	Approved	Posible valuesFlowPressureTemperatureDomain expert						
	Sensor value	System default	🗅 Draft	Sensor measurement						
	Uploaded date time	System default	🗅 Draft	Timestamp vwith date and time						
	Utility type	System default	🗅 Draft	Types: Drinking water Rain water Resource water						

D2. Data discovery | Filters

Object Type Data store	Classification Technical n	netadata	Assigned term Business met	adata-	Contact Business metadata
🗌 🞯 Dashboards	Aarhus Vand	3	File reference	2	Mads Thomsen (mads.thomsen 🛽
000 Data pipelines	FileID	2	Sensor ID template	2	_
Files	SensorID	2	Config file	1	
Folders	Configuration	1	Created date time	1	
Glossary terms	Delimiter	0	IOT sensor	0	
Reports	Description	1	IOT sensor measurement	1	
Stored procedures	File name	0	Metro maps	0	
🗌 🆩 Tables	Location	1	Sensor value	1	
	Measure	0	Uploaded date time	0	
Collection Business metadata	Media	0	Utility type	1	
Drinking water 6	See all				

D3. Data discovery | Asset schema

Microsoft Purview	urview	\mathcal{P} certified	×		<28 Q @ `	? R MTH@aarhusvand.dk ONEWATER
"	Data catalog > Search results "certified" 3	>				
📮 Data catalog	Measurement.csv	rtified				1 of 3 ↑ 🥠
🚸 Data map –	🖉 Edit 🕀 Select for bulk edit 🖵 R	lequest access 🕐 Refrest 📋 🛙	Delete 📰 E dit-solu mns 🕞 Share 🗸			Open in Power BI Desktop
Data share (preview)	Overview Properties Schema	Lineage Contacts Relate	ed	_	Updated on Nover	nber 8, 2022 9:46 PM by automated sca
🔮 Data estate insights	∀ Filter by name					
Data policy (preview)	Showing 5 of 5 items	Ļ				
🚔 Management	Column name	Classifications	Sensitivity label	Glossary terms	Data type	Asset description
	DateTime	Timestamp		Created date time	datetime	
	Sensorid	SensorID		Sensor ID template	string	Related to Sensor.SensorId
	Measure	Measure		Sensor value	double	
	UpdateTime	Timestamp		Uploaded date time	datetime	
	FileId	FileID		File reference	integer	Related to Sensorconfig.FileId

D4. Data discovery | Domain expert knowledge

~~	Data catalog >
📔 Data catalog	Search results for noise
🚸 Data map	Source type : all 🛛 😨 Clear all filters
💕 Data share (preview)	▼ Filter by keyword « Showing 1-1 out of 1 results
🖤 Data estate insights	Object Type Sensor ID template / Sensor type
Data policy (preview)	Organization Posible values Flow Pressure Temperature Domain expert knowledge Flow: Non mid values has more noise
🚔 Management	Files Folders
Data catalog > Se	arch results "noise" >
Sensor type	ault
🖉 Edit 🛅 De	lete 💍 Refresh 🕂 New child term
Overview Re	elated Contacts
Formal name Sensor ID temple	te_Sensor type
Definition	
Posible values • Flow • Pressure • Temperatur	re v
Domain expert kn Flow: Pressure: Temperatur The fluctua Increased fl Increased fl	Non mid values has more noise is it is a snapshot of a given minute Values can be used to detect pressure shocks and extend lifetime re: Depending on the distance to the surface tions over the day typically have morning and afternoon/evening peaks low creates increased pressure low does not immediately result in lower temperature

E. Semi-automatic infer | Screenshots

E1. Business glossary terms

ilossary terms						
🕂 New term 🎂 Manage term templa	ates 🖉 Edit 🛏	Import terms 🛛 🗁 Expo	rt terms 🛍 Delete 💍 Refresh 🛛 List view			
∀ Filter by keyword	Term te	mplate : All Status	: All Contact : All 🔀 Add filter			
Showing 34 terms Collapse all						
Name	Term template	Status	Definition			
Config file	Thesis	Approved	Settings for using a measurement			
_Created date time	Thesis	🗋 Draft	Created date time			
File reference	Thesis	🗋 Draft	Relation to file			
_IOT sensor	Thesis	Approved	Wireless Internet of Things (IOT) sensor			
_IOT sensor measurement	Thesis	Approved	Measurement from IOT sensor			
Metro maps	Thesis	🗋 Draft	Function map showing sensor locations			
Sensor ID template	Thesis	Approved	Unique sensor ID using template with information ab			
ID	Thesis	Approved	Only ID of the Seonsor ID template			
IsMid	Thesis	Approved	Boolean value: True: minute average of approx. 240 m			
Measure type	Thesis	Approved	Posible values: M: measurement Z: calculated			
Node	Thesis	Approved	Scada netvwork			
Number	Thesis	🗅 Draft	Consecutive numbering of sensors in specific segment			
Segment	Thesis	Approved	Subset of Scada network			
Sensor type	Thesis	Approved	Posible values Flow Pressure Temperature Domain ex			
Sensor value	Thesis	🗅 Draft	Sensor measurement			
Uploaded date time	Thesis	🗋 Draft	Timestamp with date and time			
Utility type	Thesis	🗋 Draft	Types: Drinking water Rain water Resource water			

E2. Classification rules

REGULAR EXPRESSION	5 matches (94 steps, 0.1ms)	MATCH INFORMATION	I
: [] [] []]]]]]]]]]]]]]])]-[a-zA- / gm	Match 1 0-23 0	DST.B09_E_G1_BF1-M_FLOW
TEST STRING		Match 2 24-47	BUS.C04_E_W1_BF1-M_FLOW
OST.B09_E_G1_BF1-M_FLOW* BUS.C04_E_W1_BF1-M_FLOW*		Match 3 48-71	VIV.B08_E_G1_BF1-M_FLOW
VIV.808_E_G1_BF1-M_FLOW- 085.C11_E_G1_BF1-M_FLOW-		Match 4 72-95	OBS.C11_E_G1_BF1-M_FLOW
DIS.D01_H_B70_BF1-M_FLOW_MID		Match 5 96-120	DIS.D01_H_B70_BF1-M_FLOW
EXPLANATION			
9]+_[a=zA=Z] (2) [0=9] - [a=zA=Z] [a=zA=Z] [4] a ssserts position at start of a line ①			
 Match a single character present in the list below [a=zA=Z] [3] matches the previous token exactly 3 times 			
a-z matches a single character in the range between a (index 97) and z (index 122) (case sensitive)			
A-Z matches a single character in the range between A (index 65) and Z (index 90) (case sensitive)			
), matches the character , with index $46_{10}(2E_{16}\text{or}56_8)$ literally (case sensitive)			

F. Implicit infer | Screenshots

F1. Unstructured landing zone | Test #1

Unstruct File Edit Cmd 1 1 # D 2 imp 3 4 # U	ured Python V View Run Help ataframe manageme ort pandas as pd nstructured land	Last edit was 9 minutes ago G ent	ive feedbad	:k		
5 dfM	easurement = pd.m	read_csv('/dbfs/mnt/landir	ng/Measur	ement.csv')		
Command	took 1.99 seconds	- by mads.thomsen@aarhusvand.d	k at 24.11	.2022 15.57.07 on mads	.thomsen@a	arhusvand.dk's Cluster
Cmd 2						
1 # F 2 dfR 3 4	ind measurements esult = dfMeasure	in period 2022-06-09 06:1 ement.loc[(dfMeasurement[(dfMeasurement[(dfMeasurement[15 – 2022 DateTime DateTime SensorId	-06-09 06:29 for se '] >= '2022-06-09 ('] <= '2022-06-09 ('] == 'DIS.D01_H_B2	ensorId D 06:15:00' 06:29:59' 26_BP1-M_	IS.D01_H_B26_BP1-M_TRY) &) & TRYK')]
Command	took 0.35 seconds	- by mads.thomsen@aarhusvand.d	k at 24.11	.2022 15.57.15 on mads	.thomsen@a	arhusvand.dk's Cluster
Cmd 3						
1 dfR	esult.head()					
	DateTin	ne Sensorid	Measure	UpdateTime	FileId	
679142	2022-06-09 06:15:00.0(00 DIS.D01_H_B26_BP1-M_TRYK	4.277757	2022-06-09 07:31:23.000	10	
679150	2022-06-09 06:16:00.00	00 DIS.D01_H_B26_BP1-M_TRYK	4.318294	2022-06-09 07:31:23.000	10	
670166	2022-06-09 06:17:00.00	00 DIS.DUI_H_B20_BF1-W_TR1K	4.109030	2022-06-09 07:31:23.000	10	
679174	2022-06-09 06:19:00 00	00 DIS.D01_H_B20_BP1-M_TRYK	4.210232	2022-06-09 07:31:23.000	10	
0.0111			1.201011	2022 00 00 01:01:20:000	10	
Command	took 0.09 seconds	- by mads.thomsen@aarbusvand.d	k at 24.11	.2022 15.57.27 on mads	.thomsen@a	arhusvand.dk's Cluster
Cmd 4		by moust chomochedar naovanara			renoiloenga	
1 # C 2 pri 3 pri 4 pri 5 pri	heck expected num nt('Number of mea nt('- Expected: nt('- Actual: nt('- Total:	mbers: 1 x 15 minutes asurements') ' + str(1 * 15)) ' + str(len(dfResult))) ' + str(len(dfMeasurement))))			
Number - Expec - Actua - Total	of measurements ted: 15 l: 15 : 1445943					

F2. Unstructured landing zone | Test #2

```
Cmd 5
      # Dataframe management
 1
 2
     import pandas as pd
 3
    # Unstructured landing zone
 4
     dfMeasurement = pd.read_csv('/dbfs/mnt/landing/Measurement.csv')
 5
  Command took 1.92 seconds -- by mads, thomsen@aarhusvand.dk at 24.11.2022 16.02.02 on mads.thomsen@aarhusvand.dk's Cluster
Cmd 6
 1
      # Find measurements in period 2022-06-08 06:15 - 2022-06-08 06:29 + 2022-06-09 06:15 - 2022-06-09 06:29
 2
      # for sensorId DIS.D01_H_B26_BP1-M_TRYK
 3
     dfResult = dfMeasurement.loc[(dfMeasurement['DateTime'] >= '2022-06-08 06:15:00') &
 4
                                     (dfMeasurement['DateTime'] <= '2022-06-08 06:29:59') &</pre>
 5
                                     (dfMeasurement['SensorId'] == 'DIS.D01_H_B26_BP1-M_TRYK')]
 6
     dfResult_09_06 = dfMeasurement.loc[(dfMeasurement['DateTime'] >= '2022-06-09 06:15:00') &
 7
                                     (dfMeasurement['DateTime'] <= '2022-06-09 06:29:59') &</pre>
                                     (dfMeasurement['SensorId'] == 'DIS.D01_H_B26_BP1-M_TRYK')]
 8
 9
      dfResult = dfResult.append(dfResult_09_06)
 Command took 0.60 seconds -- by mads.thorsen@aarhusvand.dk at 24.11.2022 16.02.14 on mads.thomsen@aarhusvand.dk's Cluster
```

```
Cmd 7
      dfResult.head()
 1
                       DateTime
                                                  Sensorld Measure
                                                                               UpdateTime FileId
  469922 2022-06-08 06:15:00.000 DIS.D01_H_B26_BP1-M_TRYK 4.202184 2022-06-08 07:30:51.000
                                                                                              10
  469930
          2022-06-08 06:16:00.000 DIS.D01_H_B26_BP1-M_TRYK 4.273755 2022-06-08 07:30:51.000
                                                                                              10
  469938
          2022-06-08 06:17:00.000 DIS.D01_H_B26_BP1-M_TRYK 4.295518 2022-06-08 07:30:51.000
                                                                                              10
  469946
          2022-06-08 06:18:00.000 DIS.D01_H_B26_BP1-M_TRYK 4.331699 2022-06-08 07:30:51.000
                                                                                              10
   469954
          2022-06-08 06:19:00.000 DIS.D01_H_B26_BP1-M_TRYK 4.194826 2022-06-08 07:30:51.000
                                                                                              10
  Command took 0.04 seconds -- by mads.thomsen@aarhusvand.dk at 24.11.2022 16.02.26 on mads.thomsen@aarhusvand.dk's Cluster
Cmd 8
  1
    # Check expected numbers: 2 * 15 minutes
     print('Number of measurements')
 2
 3
     print('- Exprected: ' + str(2 * 15))
     print('- Actual: ' + str(len(dfResult)))
print('- Total: ' + str(len(dfMeasurement)))
 4
  5
 Number of measurements
  - Exprected: 30
  - Actual:
                30
    Total:
                1445943
```