# PASC22 Conference

## Enabling Industrialized Analysis of Textual Documents in Data Lakes

Pegdwendé Nicolas SAWADOGO

*pegdwende.sawadogo@univ-lyon2.fr*

June 27, 2022

# Outline

# Why Data Lakes?
### Welcome to the big data era

Tremendous growth of produced and available data

### Big data opportunities

# Why Data Lakes?
## Welcome to the big data era

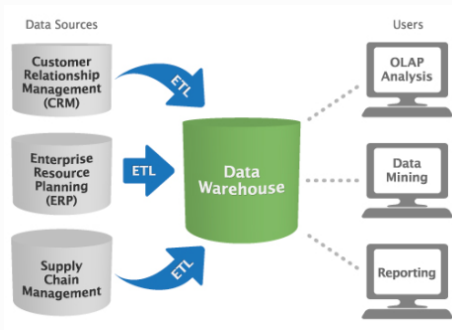Tremendous growth of produced and available data

Big data opportunities

Big data challenges



*https://www.ibm.com*

## Why Data Lakes?
From data warehouses to data lakes

- Data warehouses do achieve insights from big data.
- Distributed technologies to tackle **Volume**
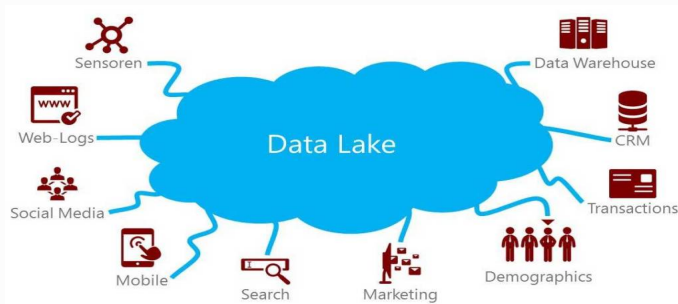- ...but **Variety** and **Velocity** pose great challenges.



*https://intellipaat.com*

## What is a Data Lake?
Definition

### James Dixon (2010)

A data lake is a large **repository** of **raw** and **heterogeneous data**, fed by external sources and allowing users to **explore**, sample and **analyze** the data.



*http://dwbimaster.com*

# What is a Data Lake?
Definition

### Sawadogo et al. (2019)

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.

## Data Lake Issues
Avoiding the data swamp

- **Schema-on-read** approach
- Efficient metadata system *required* for data access and querying
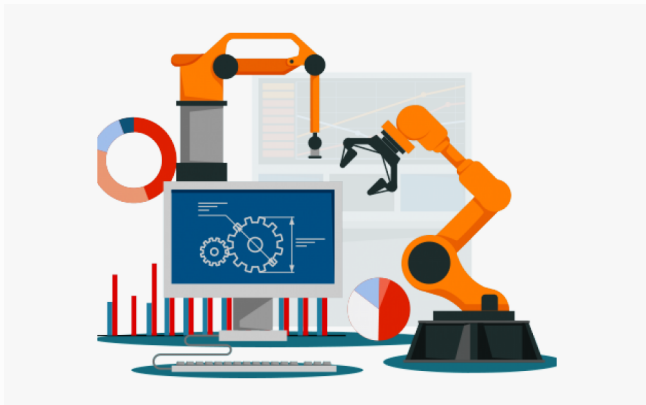- ...but how to design such a system?



*https://timoelliott.com*

# Data Lake Issues
## Enabling industrialized analyses

- Open the data lake to business users
- Make easier advanced analyses
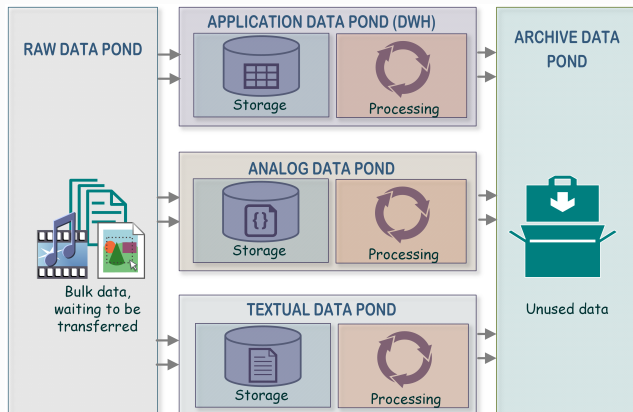- Automate metadata management

# Outline

Introduction
○○○○○○

DL and DWH
●○○

AUDAL implementation
○○○○○○○○

Textual analyses in AUDAL
○○○○○○○

Conclusion
○○

# DWH in a DL

- Approach proposed by Inmon
- Induces a data siloing

Introduction
oooooo

DL and DWH
o●o

AUDAL implementation
oooooooo

Textual analyses in AUDAL
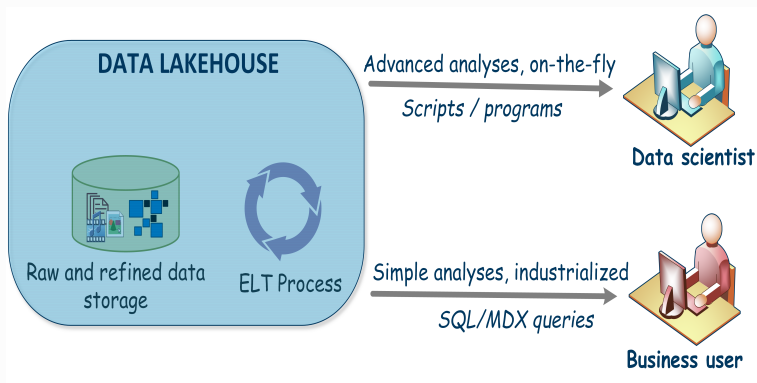ooooooo

Conclusion
oo

# DL ahead of a DWH

- Based on a functional distinction
- DL for ponctual analyses, DWH for industrialized analyses

# DL merged with a DWH

- Most recent approach (still maturing)
- Industrialization of analyses in the DL

# Outline

# AURA-PMI Project

*AURA-PMI, a multi-disciplinary project*

- Research project in management
- Analysis of the digitization of small enterprises
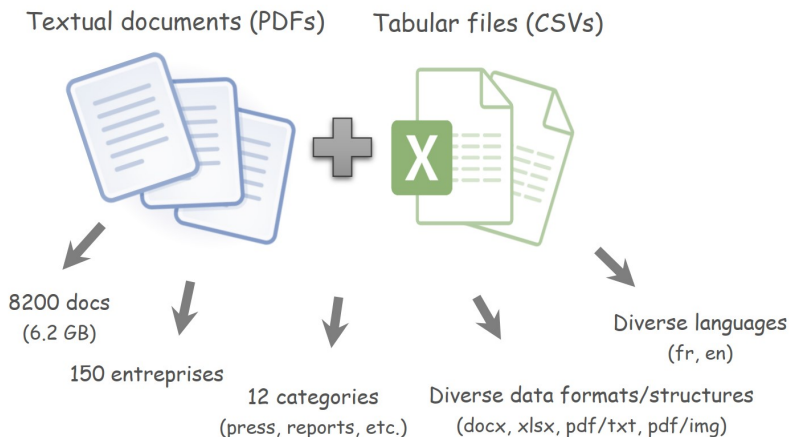- Comparison of digitization policies across categories of small enterprises

# AURA-PMI Project
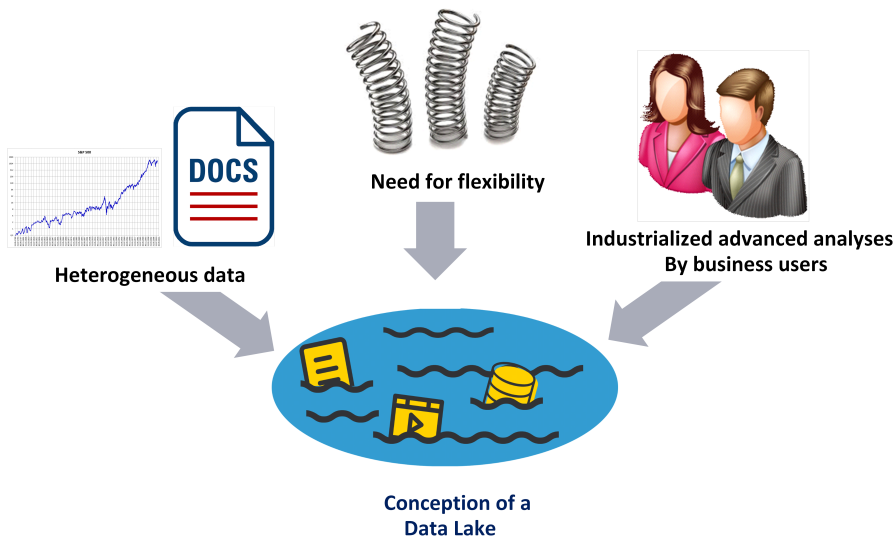*Extract insights from data*

- Enterprises' characteristics (region, nb. employees, domain, etc.)
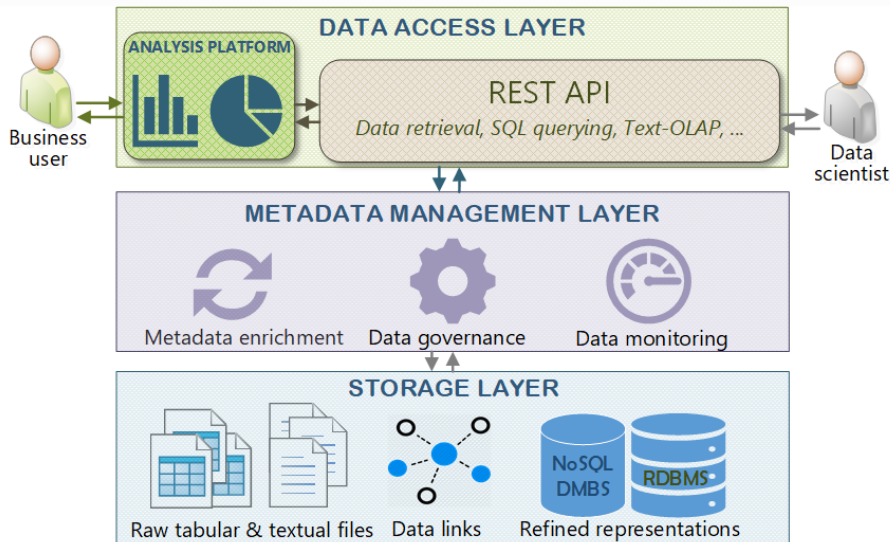- Annual reports, financial reports, etc.



Textual documents (PDFs)    Tabular files (CSVs)

8200 docs
(6.2 GB)

150 entreprises

12 categories
(press, reports, etc.)

Diverse data formats/structures
(docx, xlsx, pdf/txt, pdf/img)

Diverse languages
(fr, en)

Introduction
○○○○○○

DL and DWH
○○○

AUDAL implementation
○○●○○○○○

Textual analyses in AUDAL
○○○○○○○

Conclusion
○○

# AURA-PMI Project

*Need for a data lake*



**Need for flexibility**

**Industrialized advanced analyses
By business users**

**Heterogeneous data**

**Conception of a
Data Lake**

# Architecture of AUDAL
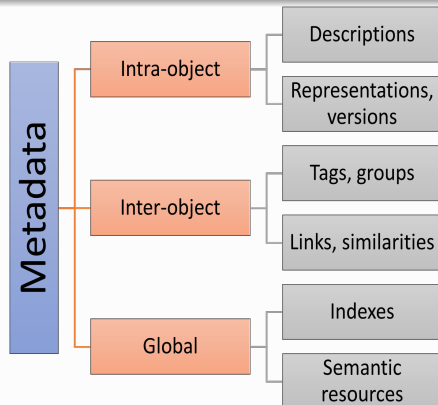
# Metadata management

An extensive vision of metadata

### Visengeriyeva (2020)

Metadata are **structured information** that **describes**, **explains**, **locates**, or otherwise makes it easier to **retrieve**, **use**, or **manage** information resources.
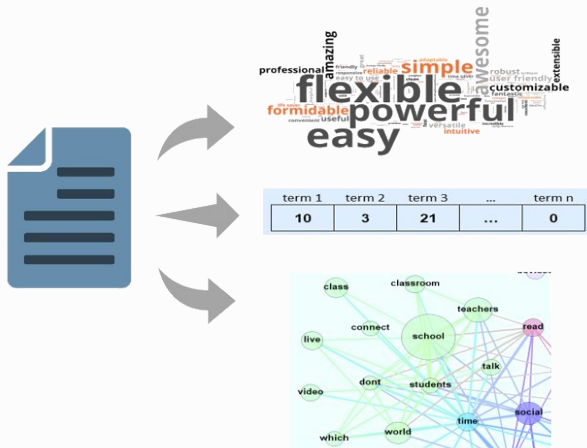
$$DL = \langle \mathscr{D}, \mathscr{M} \rangle$$
$$\mathscr{M} = \langle \mathscr{M}_{intra}, \mathscr{M}_{inter}, \mathscr{M}_{glob} \rangle$$
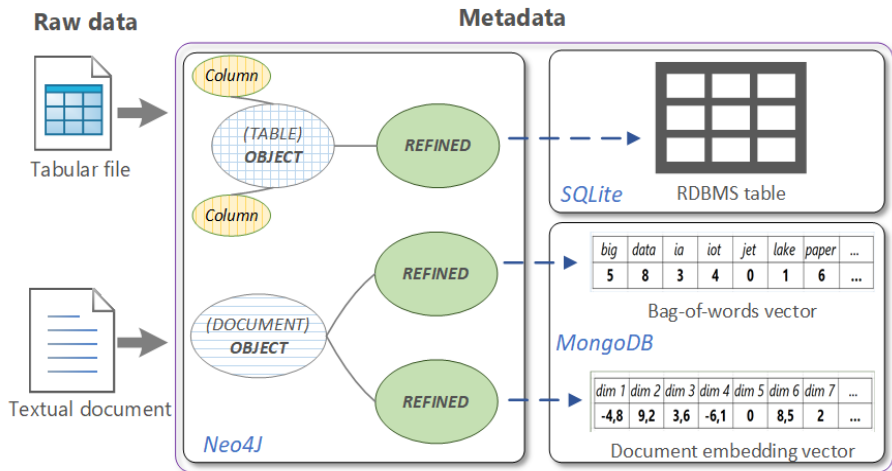
# Metadata management
Data polymorphism

- ▶ Simultaneously manage different representations of data
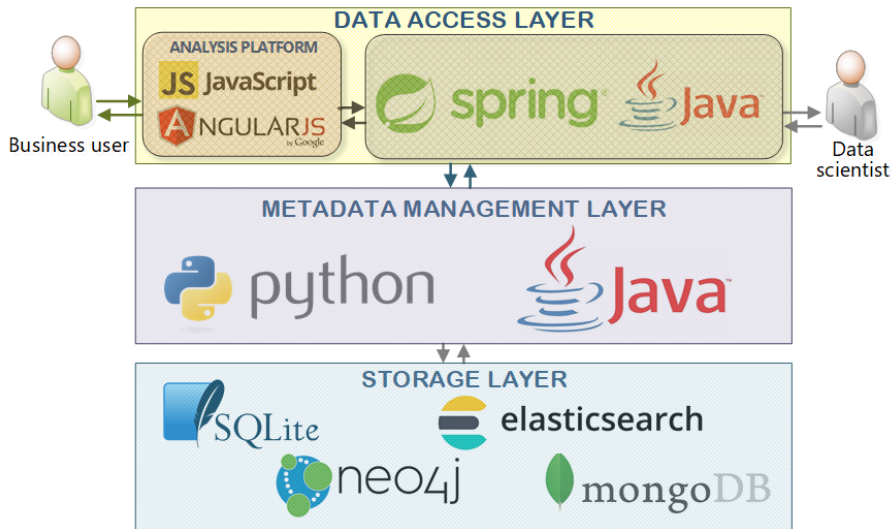- ▶ Such representations are viewed as metadata of raw data

# Metadata management in AUDAL (3/3)

## Intra-object metadata

# Technologies used in AUDAL
AUDAL = AUra-pmi DAta Lake

# Outline

Introduction
○○○○○○

DL and DWH
○○○

AUDAL implementation
○○○○○○○○

**Textual analyses in AUDAL**
●○○○○○○

Conclusion
○○

# Analyses dans AUDAL

# Analyses in AUDAL
*Data retrieval*

## Keyword-based filtering

🔍 Terms filtering     +

*Terms*

+ matching    [ ... ] +

        *relation*   *client*

− matching    [ ... ] +

*Parameters*

Strictness    **Any** | All

Fuzzy search    **Yes** | No

Terms extension    − | None ⌄ | +

## Category-based filtering

🏭 Groupings     −

*Groups*

▶ 1- cible

▶ 2- digitalNativity

▶ 3- docCategory

▶ 4- enterprise

▶ 5- language

▶ 6- mimeType

▶ 7- month

▼ 8- region

☑ *ALL*

☑ *Auvergne-Rhone-Alpes*    [1342]

☑ *Bourgogne-Franche-Comte*    [404]

☑ *Bretagne*    [45]

# Analyses in AUDAL
*Textual data aggregation*

**Highlighting documents' content with a wordcloud**

Introduction
○○○○○○
DL and DWH
○○○
AUDAL implementation
○○○○○○○○
**Textual analyses in AUDAL**
○○○●○○○
Conclusion
○○

# Analyses dans AUDAL
*Textual data aggregations*

## Highlighting documents' content with a Concordance

2 - Comptes sociaux semestriels_2015_ANEVIA.pdf[10]

La valorisation des technologies acquises est amortie sur 8 ans et celle

d'utilisation retenues sont les suivantes : Nature Durée retenues ns Te
***Relations***

d'entreprises, nécessaires pour la mise en œuvre de partenariats OSE
***relations***

***clients*** correspondent aux portefeuilles ***client***, stables et pérennes (cf
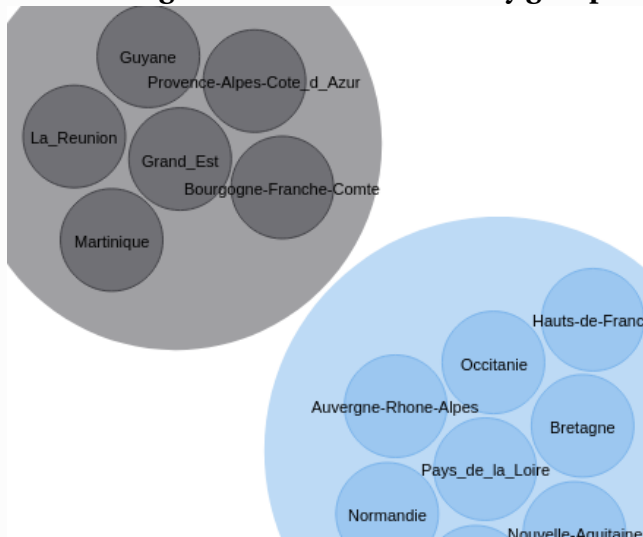
***Relation*** Total acquises brevets clientèle Valeurs brutes a au 31.12.15

La juste valeur de la ***relation*** clientèle est évaluée selon la méthode du

## Analyses in AUDAL
*Textual data aggregation*
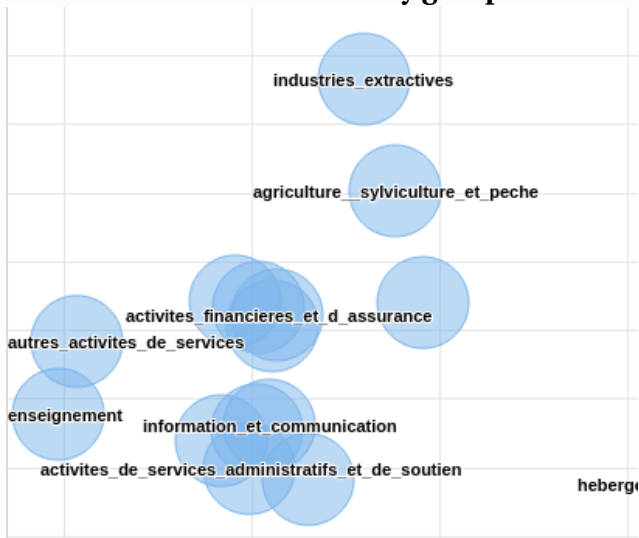
**Clustering KMeans of documents by groups**

Introduction
○○○○○○

DL and DWH
○○○

AUDAL implementation
○○○○○○○○

Textual analyses in AUDAL
○○○○○●○

Conclusion
○○

# Analyses in AUDAL
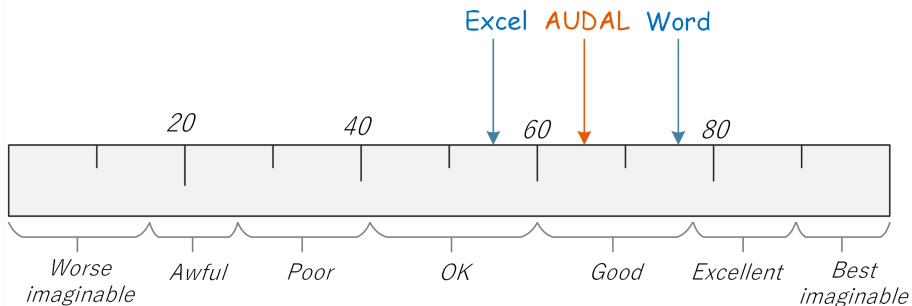*Textual data aggregation*

**PCA of documents by groups**

# User experience with AUDAL
*Usability*

- ▶ Method *System Usability Scale* (SUS)
- ▶ Protocol based on a feedback from 6 users via a questionnary

# Outline

## Conclusion

**We presented...**

★ **How DSS architectures evolved with the DL wave**

- DWH in the DL
- DL ahead of the DWH
- DL merged with the DWH

★ **How to activate industrialized analyses for textual documents in a DL**

- Using an extensive vision of metadata
- Thanks to a principle of data polymorphism
- With a combination of storage technologies

## Conclusion

**We presented...**

★ **How DSS architectures evolved with the DL wave**

- DWH in the DL
- DL ahead of the DWH
- DL merged with the DWH

★ **How to activate industrialized analyses for textual documents in a DL**

- Using an extensive vision of metadata
- Thanks to a principle of data polymorphism
- With a combination of storage technologies

**What's next ?**

▶ **Activate deeper textual data analyses**

- Sentiment analysis

▶ **Industrialized analyses for more unstructured data in DLs**

- Images
- Videos

## PASC22 Conference

## Enabling Industrialized Analysis of Textual Documents in Data Lakes

Pegdwendé Nicolas SAWADOGO

*pegdwende.sawadogo@univ-lyon2.fr*

June 27, 2022