

# Digital data collections: characteristics and properties

---

**Geneveva Vargas-Solar**  
**French Council of Scientific Research, LIG**  
**[geneveva.vargas@imag.fr](mailto:geneveva.vargas@imag.fr)**

<http://vargas-solar.com/data-centric-smart-everything/>

<https://classroom.google.com/c/MTQ4MzcwMjY1MDEz?cjc=5bz2tk6>

Slack channel: [https://join.slack.com/t/colenationale-5jr8199/shared\\_invite/zt-hhf9euv7-bmp7Kn9LL68RyzdhJnbKxA](https://join.slack.com/t/colenationale-5jr8199/shared_invite/zt-hhf9euv7-bmp7Kn9LL68RyzdhJnbKxA)





# DATA

# BIG DATA DEFINITION

- Data collections with characteristics difficult to process on single machines or traditional databases
- A new generation of tools, methods and technologies to collect, process and analyse massive data collections
  - Tools imposing the use of parallel processing and distributed storage



## 5v: Value

Which is the real value of data?



**VOLUME**  
DATA SIZE



**VELOCITY**  
SPEED OF CHANGE



**VARIETY**  
DIFFERENT FORMS  
OF DATA SOURCES



**VERACITY**  
UNCERTAINTY OF  
DATA

# BIG DATA PROPERTIES



3V

4V

5V

...

10V

- **Volume** (size)
- **Velocity** (production rate)
- **Variety** (data types & format)
- **Variability** (inconsistencies by constant meaning changes)
- **Veracity** (truth and consistency)
- **Value** (how much information)

V's models [Jagadish 2014]

“Big Data can really be very small and not all large datasets are big!”

- Mike 2.0 [Hillard 2012]

# HOW BIG IS YOUR DATA ?

*For Starters..*

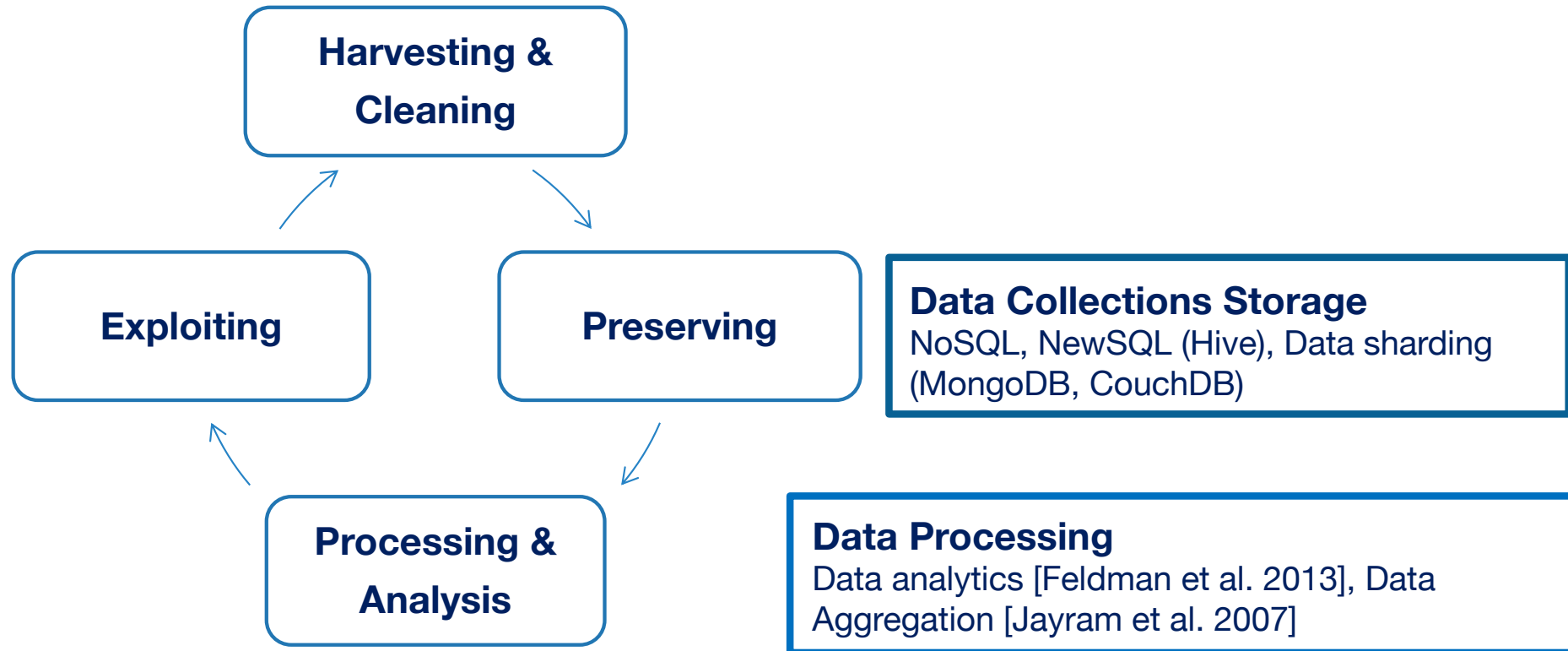
BIT	=	A BINARY DIGIT SET TO EITHER A 1 OR 0
BYTE	=	8 BITS
KB	=	1,000 BYTES
MB	=	1,000,000 BYTES
GB	=	1,000,000,000 BYTES

**Helluva lot of data !!**

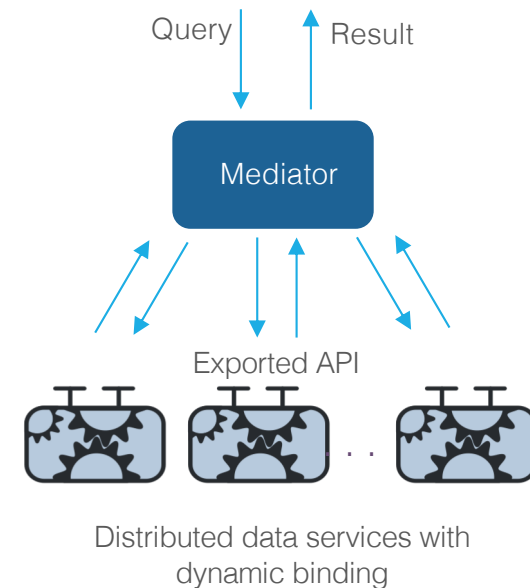
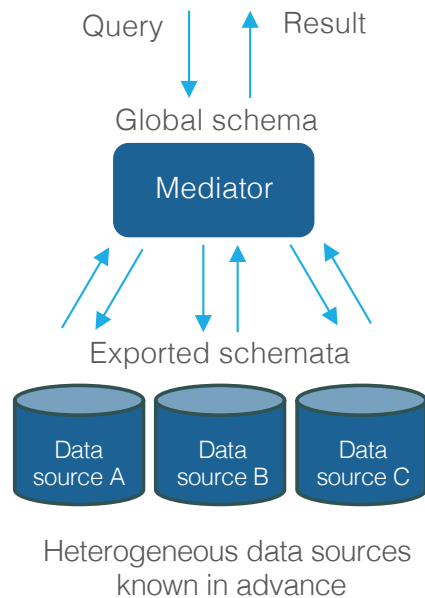
<http://spectrum.ieee.org/computing/software/beyond-just-big-data>

1 Brontobyte	1 000 Yottabytes
1 Geopbyte	1 000 Brontobytes

# BIG DATA LIFE CYCLE



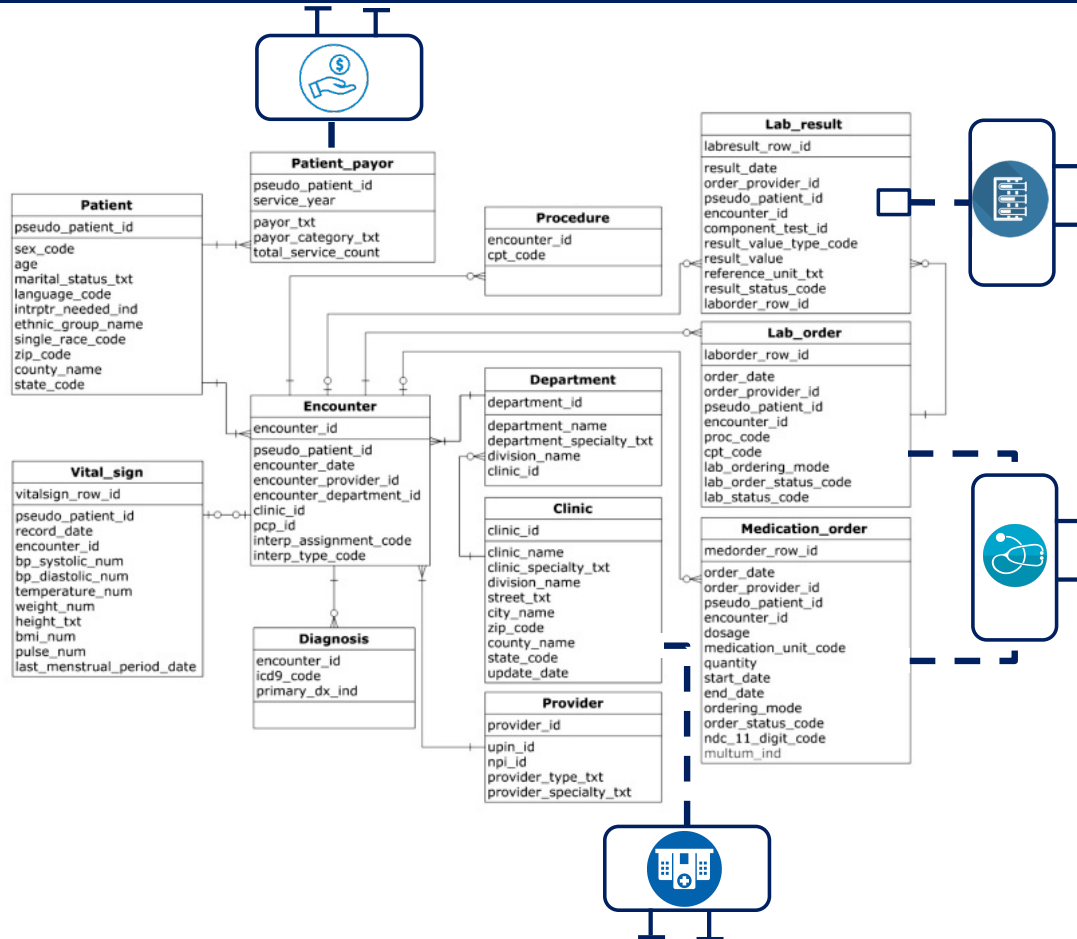
# DATA DESIGN & CONSUMPTION PHILOSOPHIES



- **Data:** model of a mini-world, it is a set of facts structured according to some data model
- **Query:** precisely stated it can include terms, operators (and/or/negation, relational, aggregation), and constraints
- **Result:** collection of items that completely or partially correspond to consumers requirements (precision & recall)



# ASKING FACTS ABOUT DIABETES

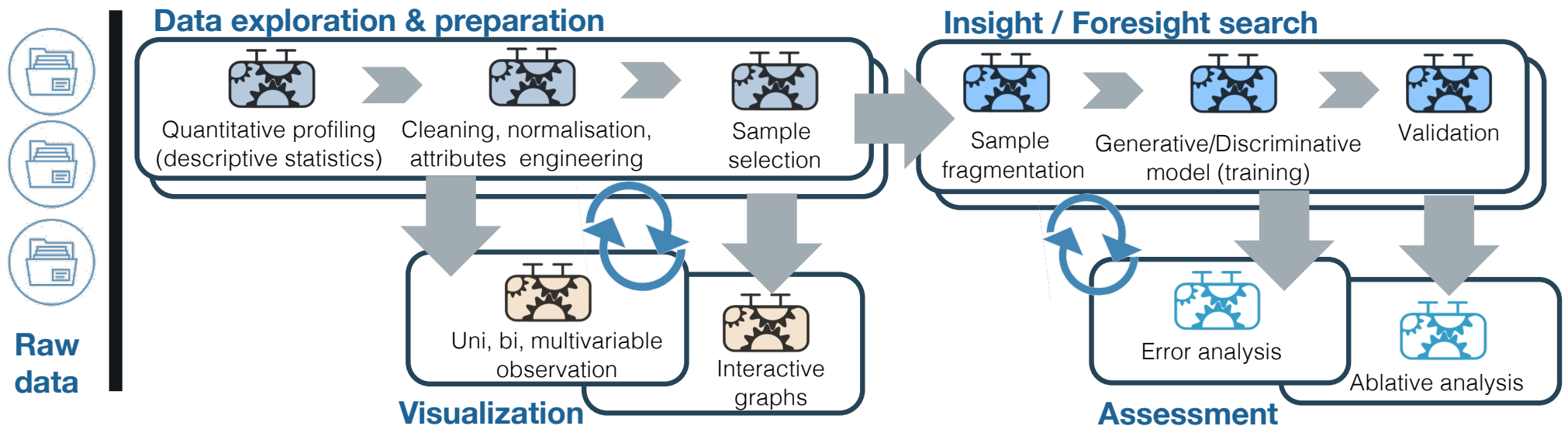


- Average of lab & medication orders per patient and physician in a given clinic
- Vital signs and lab results used to emit a diagnosis for a given patient
- Number of patients with diabetes followed per clinic



# DATA CONSUMPTION PHILOSOPHIES

## Diogenes approach



- **Data:** observations of phenomena often described as series of features/attributes
- **Query:** analytics objective (looks for insights or foresights) expressed as a pipeline of operations guided by the conditions and characteristics of the data
- **Result:** a model or prediction with associated assessment indexes, not definitive accepted with an associated error margin, accepted by comparison

<https://www.kaggle.com/hamzael1/hospital-beds-by-country>



# DIGITAL DATA COLLECTIONS

## NEITHER MANAGEABLE NOR EXPLOITABLE AS SUCH

### RAW DATA

- Heterogeneous (**variety**)
- Huge (**volume**)
- Incomplete, unprecise, missing, contradictory (**veracity**)
- Continuous releases produced at different rates (**velocity**)
- Proprietary, critical, private (**value**)

# THE RIGHT DATA FOR THE RIGHT ANALYTICS

*Different sizes, evolution in structure, completeness, production conditions & content, access policies modification ...*

*Is the data clean? If not what has to be done to make it clean?*

*What information is available in these collections?*

*Are there relations between data collections which could be exploited for better prediction?*

*What are the update rates and in what way does this affect the collection?*



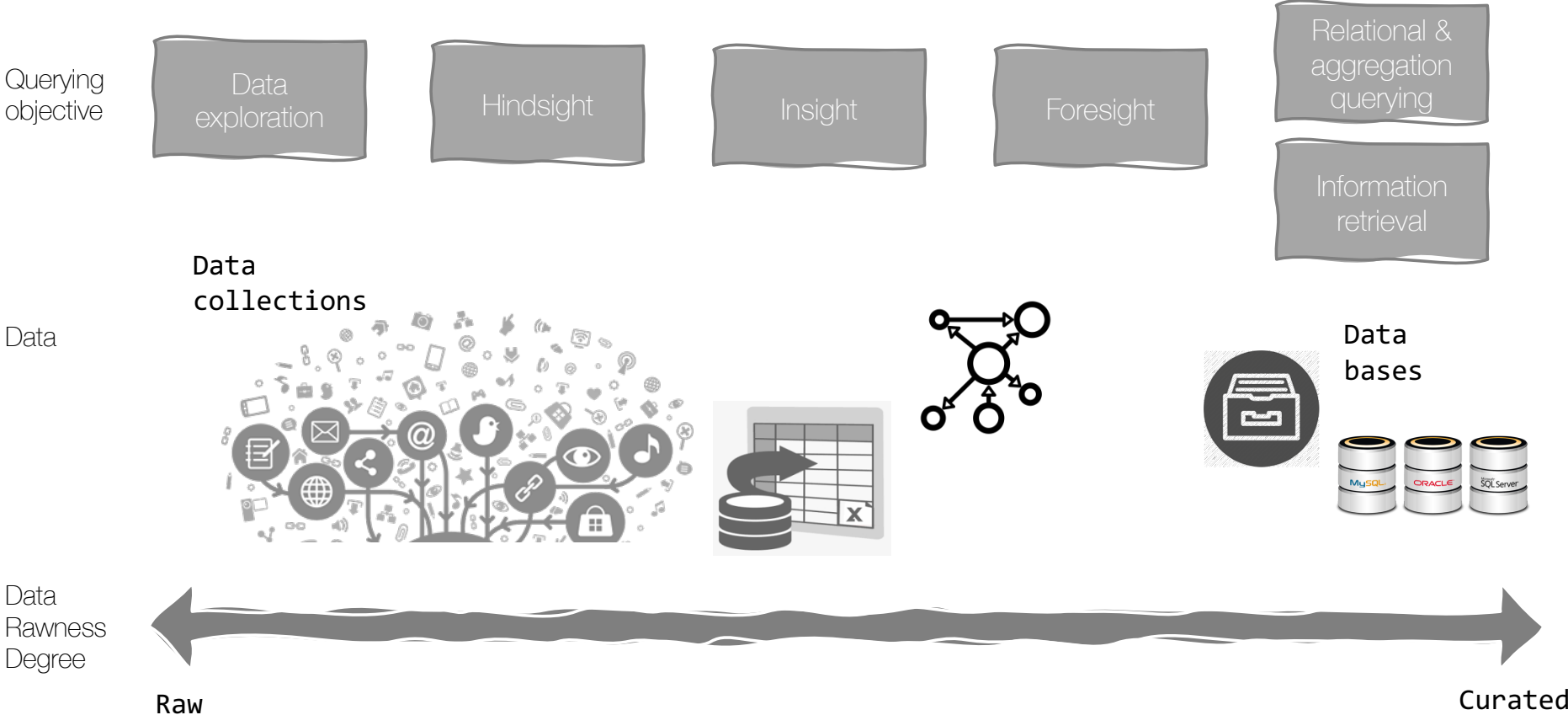
# DATA CURATION

## Data management through its **lifecycle of interest & usefulness**

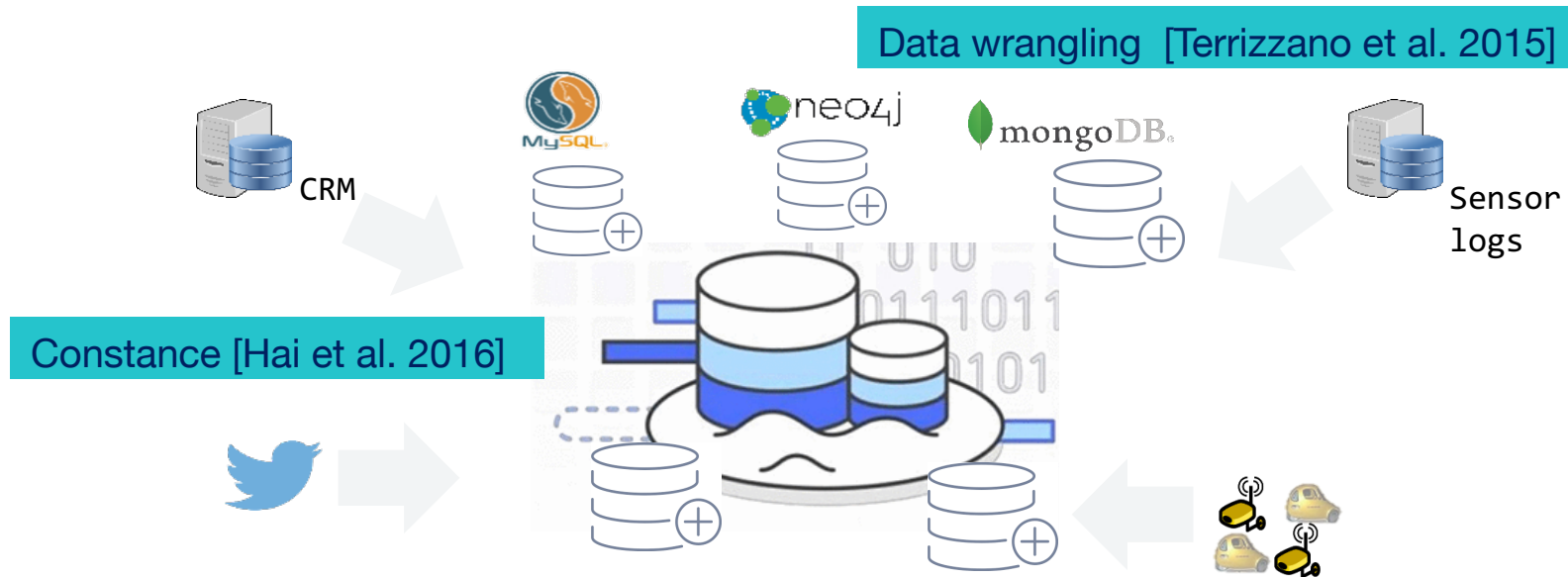
- Enable data **discovery & retrieval**
- Maintain data quality
- Add value
- Provide for **re-use over time**



# DATA SPECTRUM



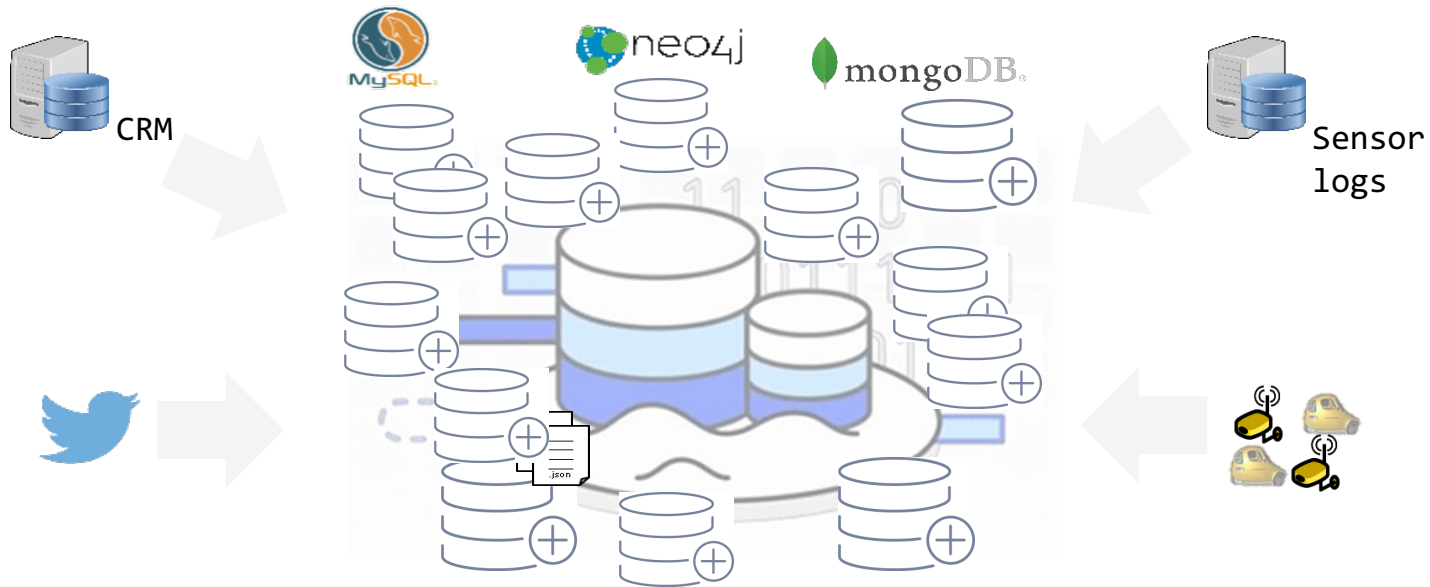
# DATA LAKE



Centralized repository containing virtually inexhaustible amounts of raw data to be analysed

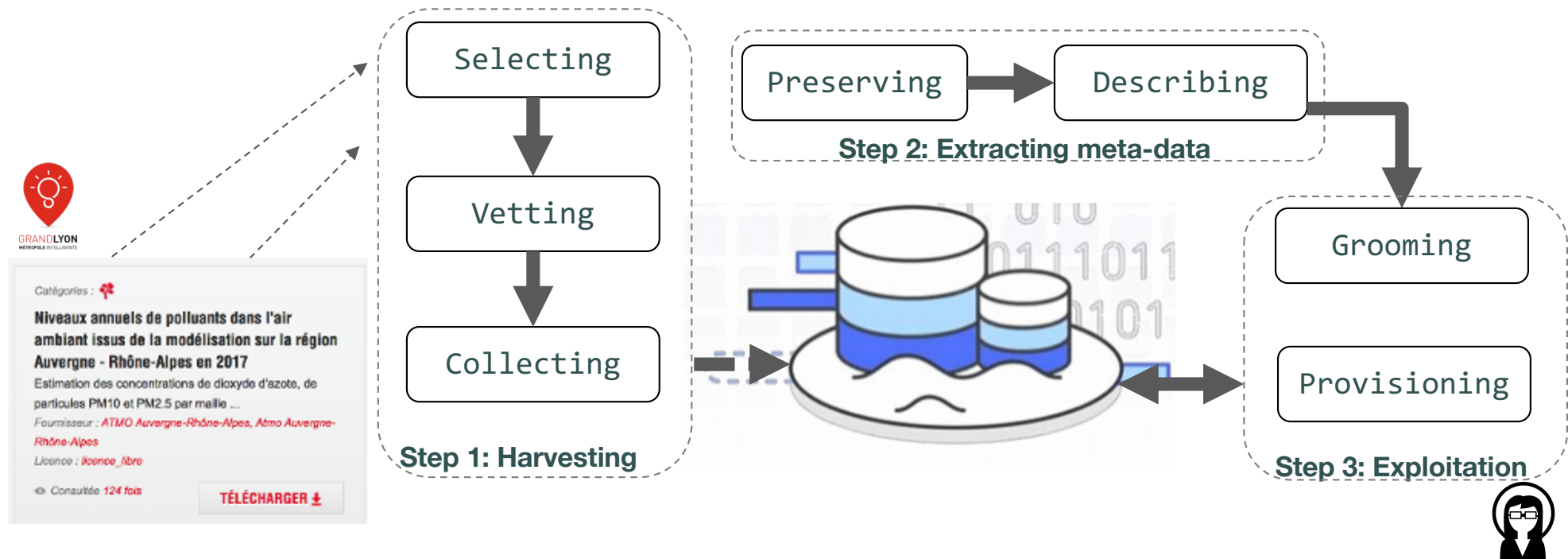


# DATA SWAMP



Repositories grow ever bigger and complex to the point that a lake becomes a **swamp**

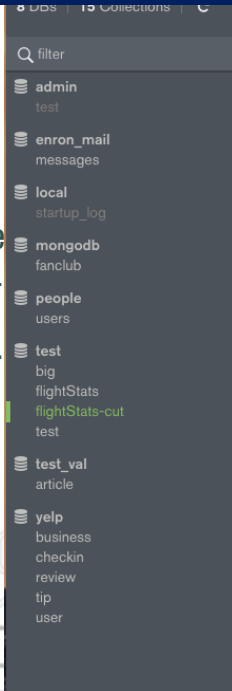
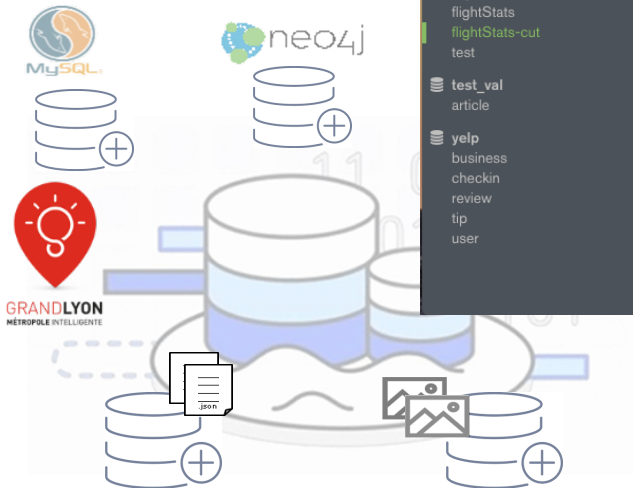
# DATA CURATION WORKFLOW



Data wrangling [Terrizzano et al. 2015]

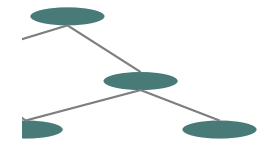
# EXTRACTING META-DATA

Structural meta-data  
(e.g. Schemata, DTD, XML)



Quantitative meta-data  
(e.g. Compass MongoDB)

data approaches)

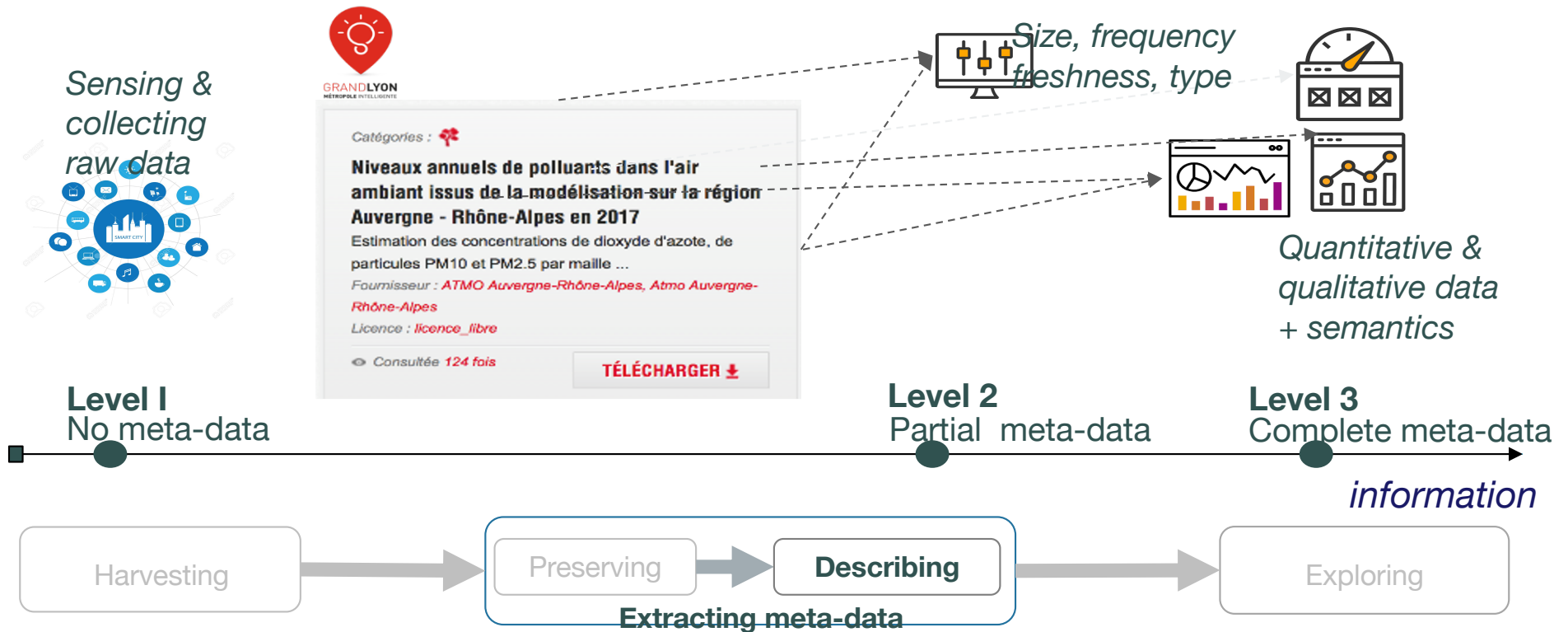


ology

ONS

- Topic
- Geographic location
- Temporal releases

# EXTRACTING META-DATA



GRAND LYON  
MÉTROPÔLE INTELLIGENTE

Catégories :

**Niveaux annuels de polluants dans l'air ambiant issus de la modélisation sur la région Auvergne - Rhône-Alpes en 2017**

Estimation des concentrations de dioxyde d'azote, de particules PM10 et PM2.5 par maille ...

Fournisseur : *ATMO Auvergne-Rhône-Alpes, Atmo Auvergne-Rhône-Alpes*

Licence : *licence\_libre*

Consultée 124 fois

TÉLÉCHARGER ↓

[Stonebraker et al. 2015]

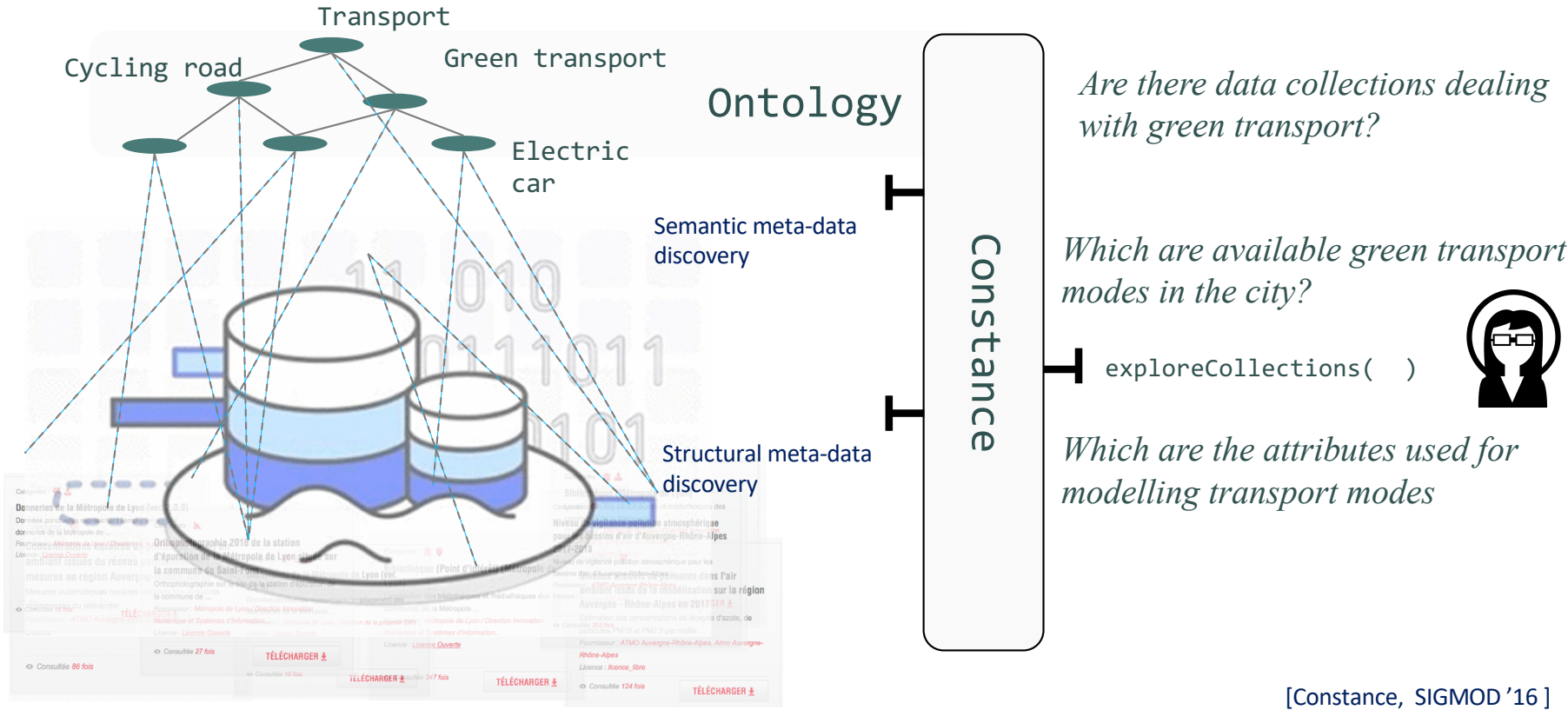
# APPROACHES FOR EXTRACTING META-DATA

	Approach	Principle	Pros	Cons
<b>Level 1:</b> simple tags as data is harvested	<b>Curation at Source</b> [Curry et al. 2010]	Simple <b>tagging &amp; information extraction</b> directly from the sensor	Pre-processes data according to <b>samples</b>	No awareness of the whole data collection
<b>Level 2:</b> pivot vocabulary & semantics	<b>Master Data Management</b> [Weatherspoon et al. 2013]	Provide a <b>standard vocabulary</b> at the company scale	Standardizes the language used in data collections	Does not improve data structures
	<b>Semantic linking</b> <b>Constance</b> [Hai et al. 2016]	Identify <b>attributes</b> referring to similar topics in different data sets	Explores data collections	Does not improve data structures
	<b>Data set clustering</b> <b>Goods</b> [Halevy 2016]	Discover <b>similar data sets</b> using clustering	Creates an <b>synthesized representation</b> of several data collections	Difficult to define a <b>similarity criterion</b> to cluster data with <b>low quality</b> (missing & null values, types)
<b>Level 3:</b> manual and collaborative	<b>Crowdsourcing and Collaboration spaces</b> [Doan et al. 2005]	<b>Communities</b> produce, maintain & tag data (crowdsourcing)	Improves the <b>quality</b> of raw data	Manual & a lot of human resources

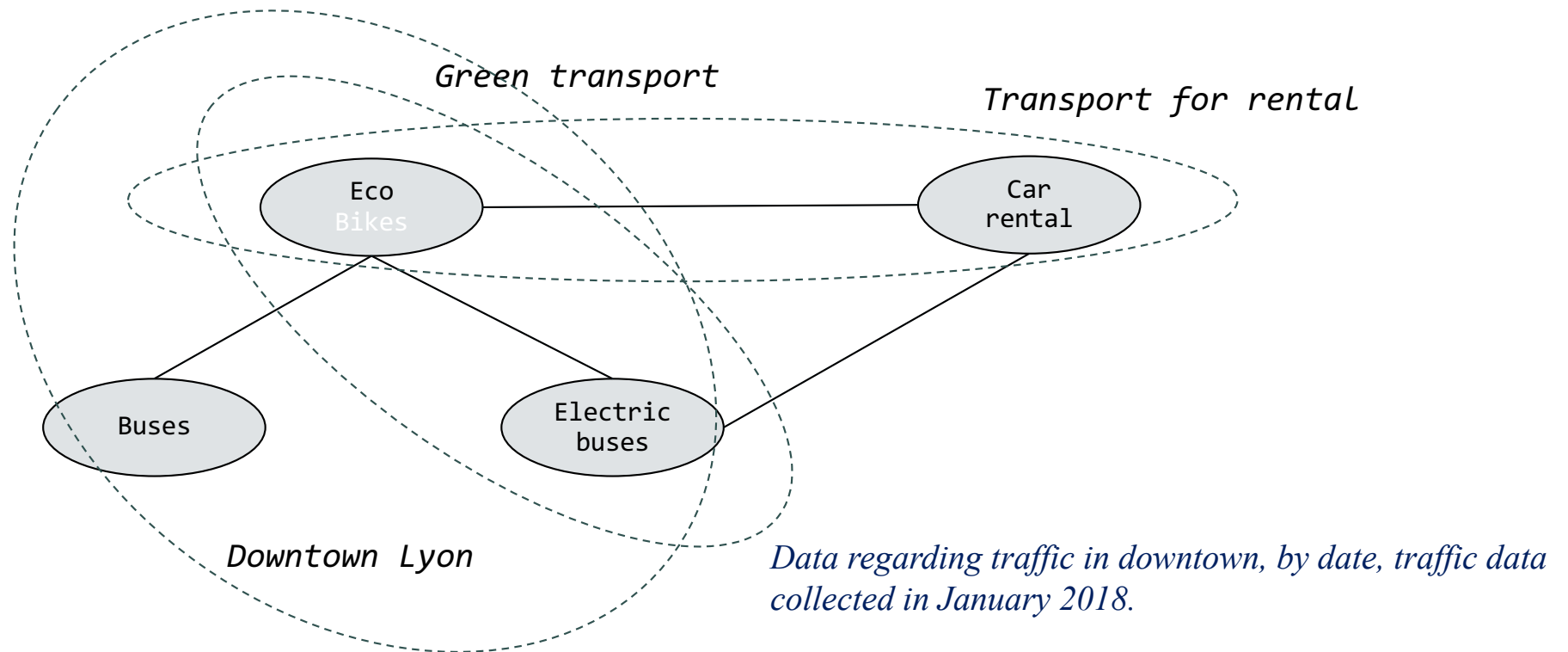
# DATA CURATION ON THE LAKES

Approach	Principle	Pros	Cons
<b>Meta data</b> : Constance [Hai et al 2016], [Stonebraker et al 2013] Data wrangling [Terrizzano et al 2015]	Extraction of semantic meta-data or mapping items with concepts Query using SPARQL, regular expression	Possibility to express declarative queries; difficult to automate completely	No statistics, iterative data querying for exploring data
<b>Data content description</b> : descriptive statistics, processing according to data types, schema extraction	Explore data structure for extracting the schema Compute descriptive statistics functions for every element	Aggregated view of the content despite data types. Simple to visualize and scale if important volume.	Adapted for (semi)structured data, manual tagging for multimedia content
<b>Curation</b> : [CoreKG, Curry 2016, QoS MOS 2018, Tacit knowledge management 2017]	API with methods for preserving data collections. Querying and data fusion operations	Exploit Compass tool from MongoDB for providing a quantitative vision of data content	Data transformation from CSV to JSON. No semantic knowledge (terms, functional dependencies)

# FISHING DATA IN THE LAKE

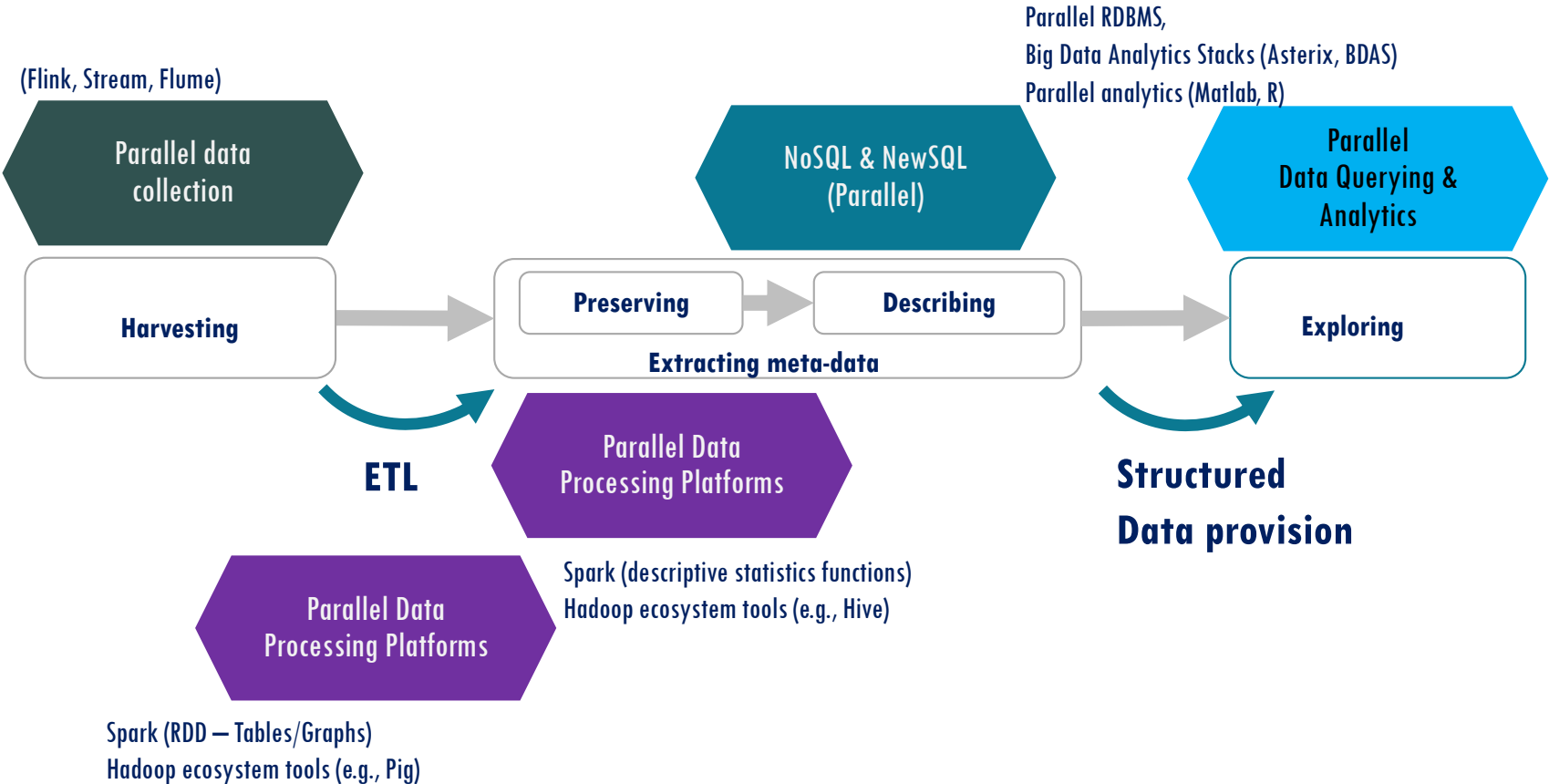


# EXTRACTING META-DATA





# PREPARING DATA COLLECTIONS



# CLOUD COMPUTING AND BIG DATA

Ready to use environments for  
**Storing Big Data & Running greedy processing tasks**

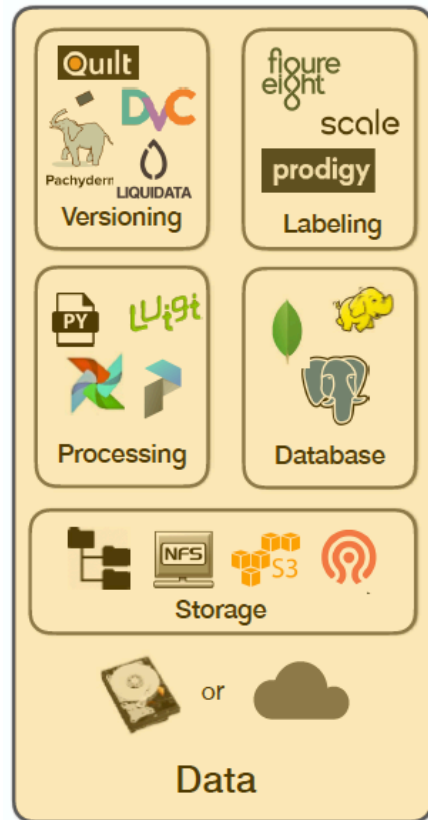
Platform as a Service (**PaaS**)  
Database systems, frameworks

Infrastructure as a Service (**IaaS**)  
CPU, RAM, Disk



Software as a Service (**SaaS**)  
Full functional software

# DATA LABS SERVICES



-  [kaggle.com](https://www.kaggle.com)
-  Google Colab
-  Azure Notebooks



