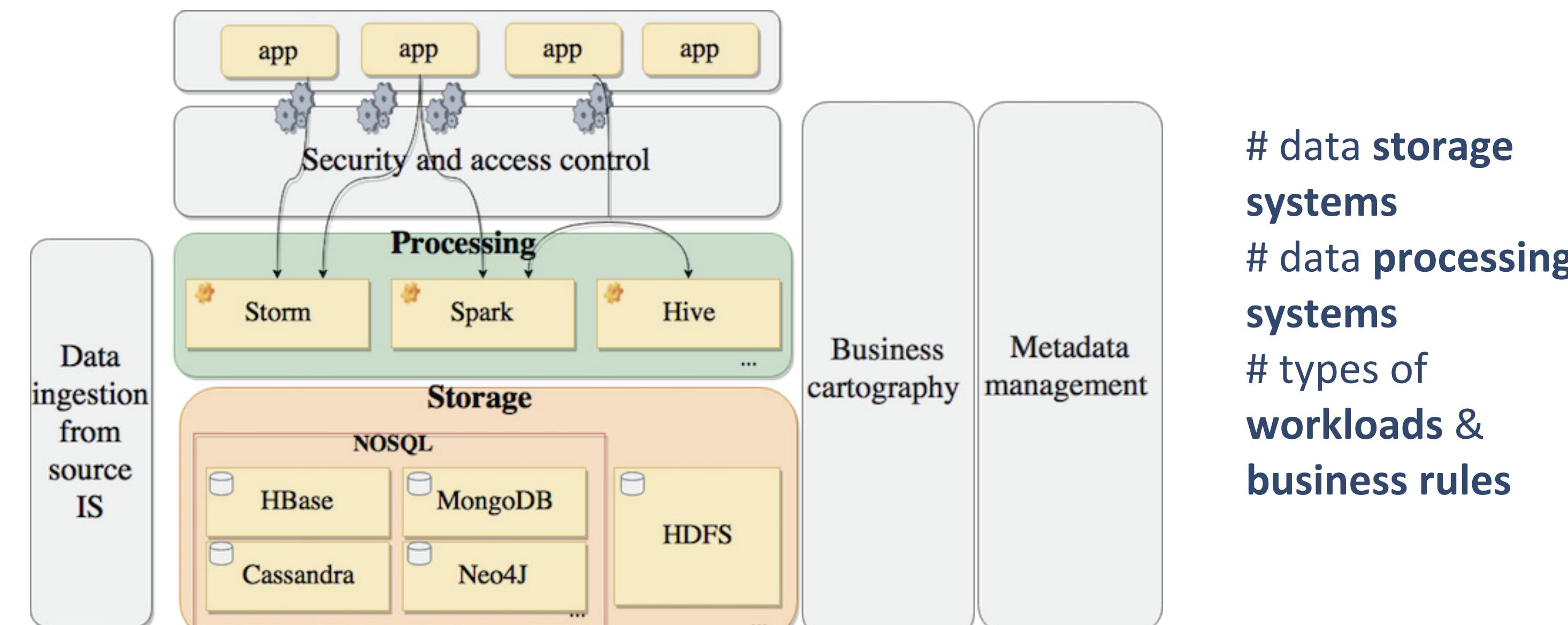
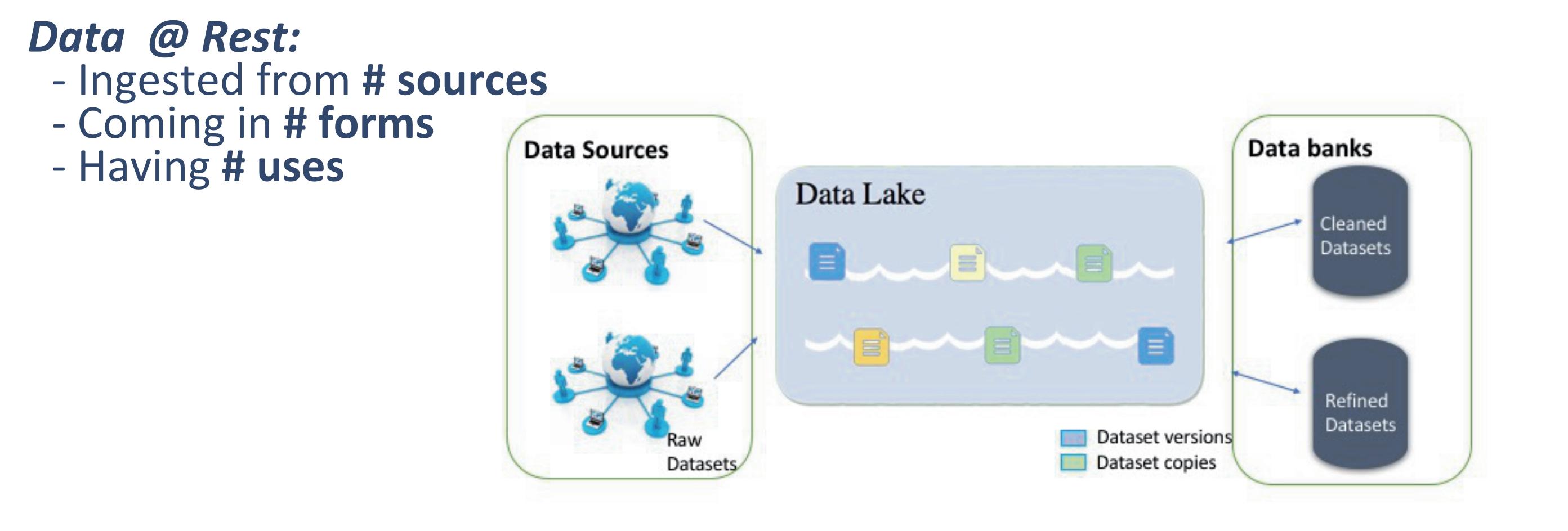
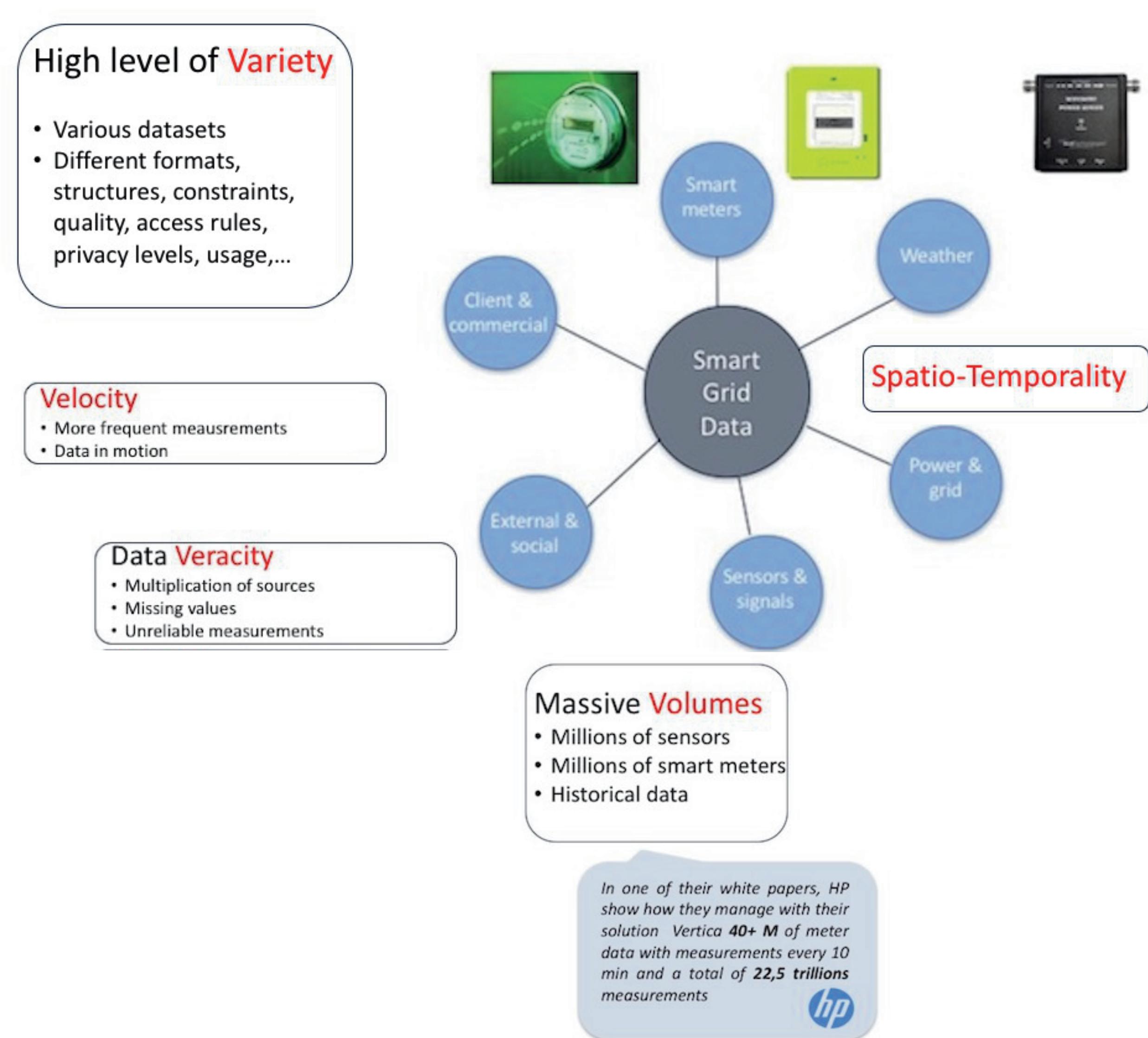




## Smart grid Big data ecosystem



**Complexity of the Data ecosystem!**

## Placement of datasets

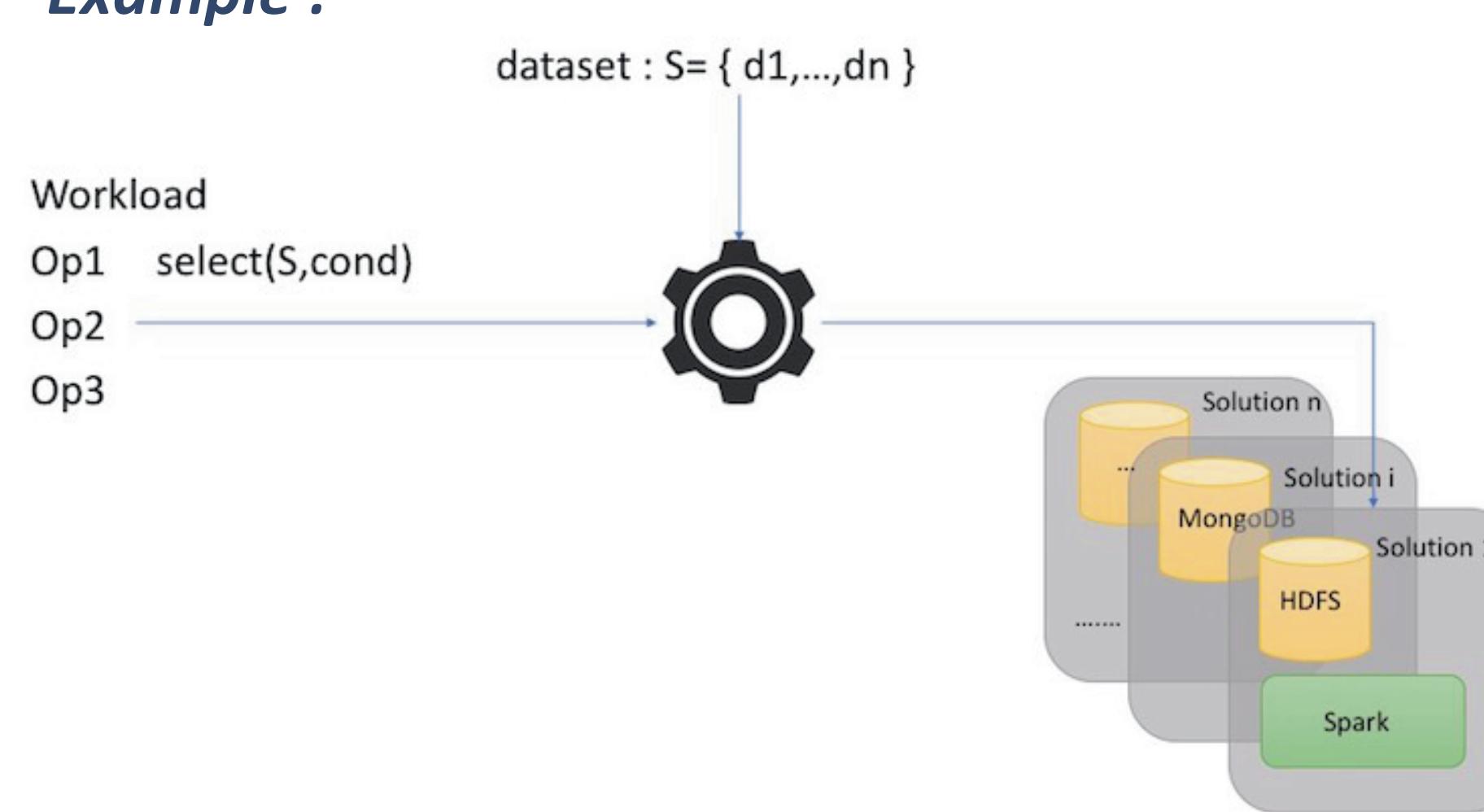
### Objective :

- Effective placement and processing of datasets

### Constraints :

- Efficient workload processing  
- Compliance with business environment & the data ecosystem

### Example :



### Where to place datasets ?

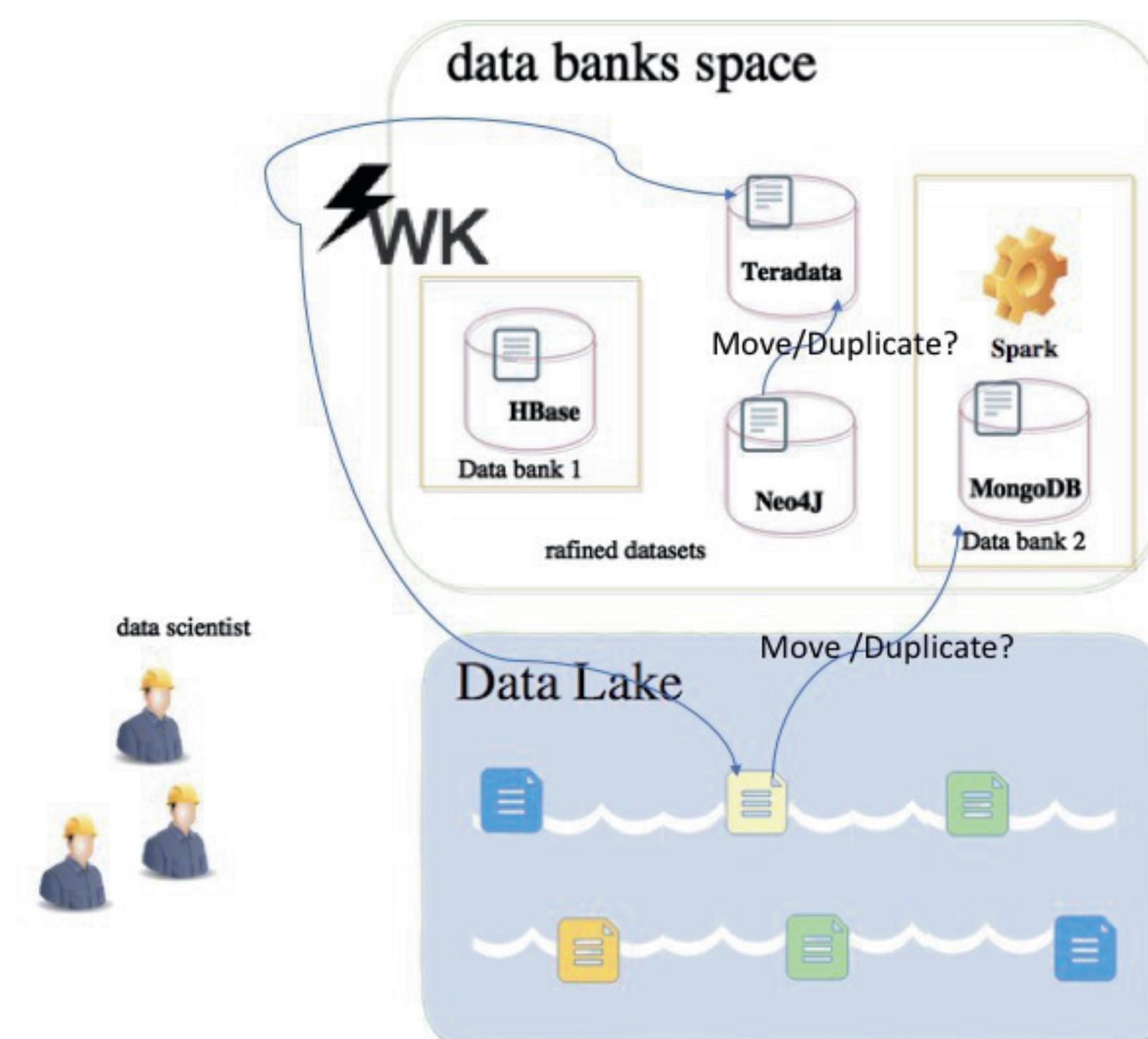
- in which databank ?
- in the lake ?

### How to place datasets ?

- duplicate / move the dataset

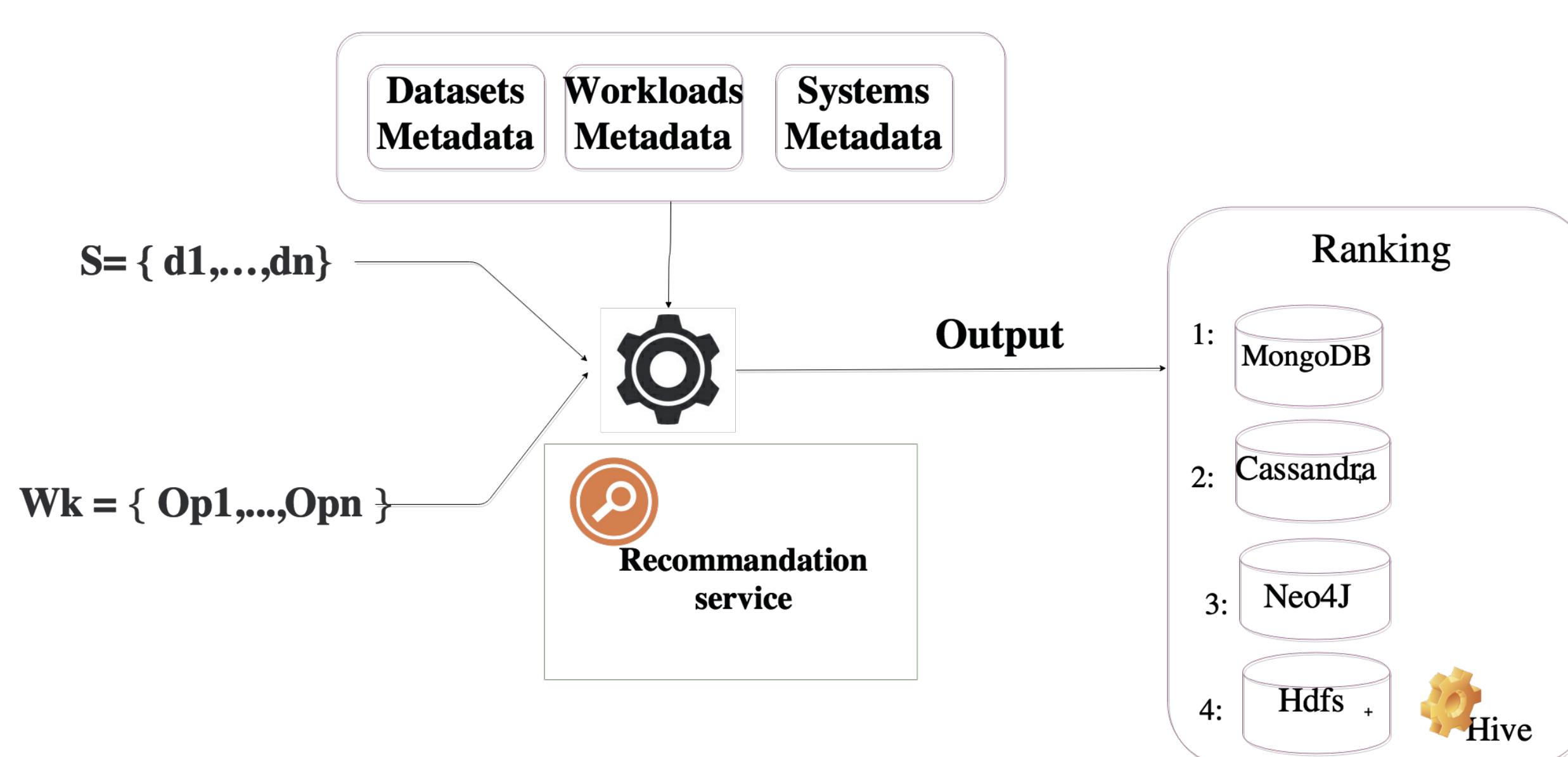
### When to place datasets ?

- before processing (on the input datasets)
- after processing (on the result, intermediate datasets)



## Workload based placement approach

### Recommendation for dataset placement



### Recommendations based on metadata about :

- Characteristics of the target systems
- Efficiency of the processing of the  $\langle wk, s \rangle$  in those systems

### Metadata

#### Workloads and applications

- Workloads metadata
  - Statistics
  - Plans
  - Patterns
  - Business rules
  - Annotations

#### Datasets

- Schema
- Statistics
- Lineage

#### Data Management Systems

- Data model
- Data partitioning
- Access API & Languages
- Data storage
- Systems configuration & Distribution

### Placement criteria

#### Feasibility

- Identification of the suitability of data storage & processing systems properties (data model, query model, partitioning, ...) with datasets and workloads metadata

#### Performance

- Response time :
  - Time needed for the execution of a workload  $wk$  on a data store DS
  - Based on cost estimation given workload, data stores and datasets characteristics
- How scalable is performance ?
- Is it crucial to the workload ?

#### Conformity

- With the **business rules** of Smart Grid Data Ecosystem

