# Usage of the "big" data

asov - paris - 2015.03.24



david.languignon@obspm.fr

# Astro data services: **Observations**



**Métadonnées**:
- positions
- instruments
- domaine spectral
- filtres
- temps d'exposition
- configuration instrument

~10 quantités

Instruments

Traitement

Meta données

Data

Database

VO

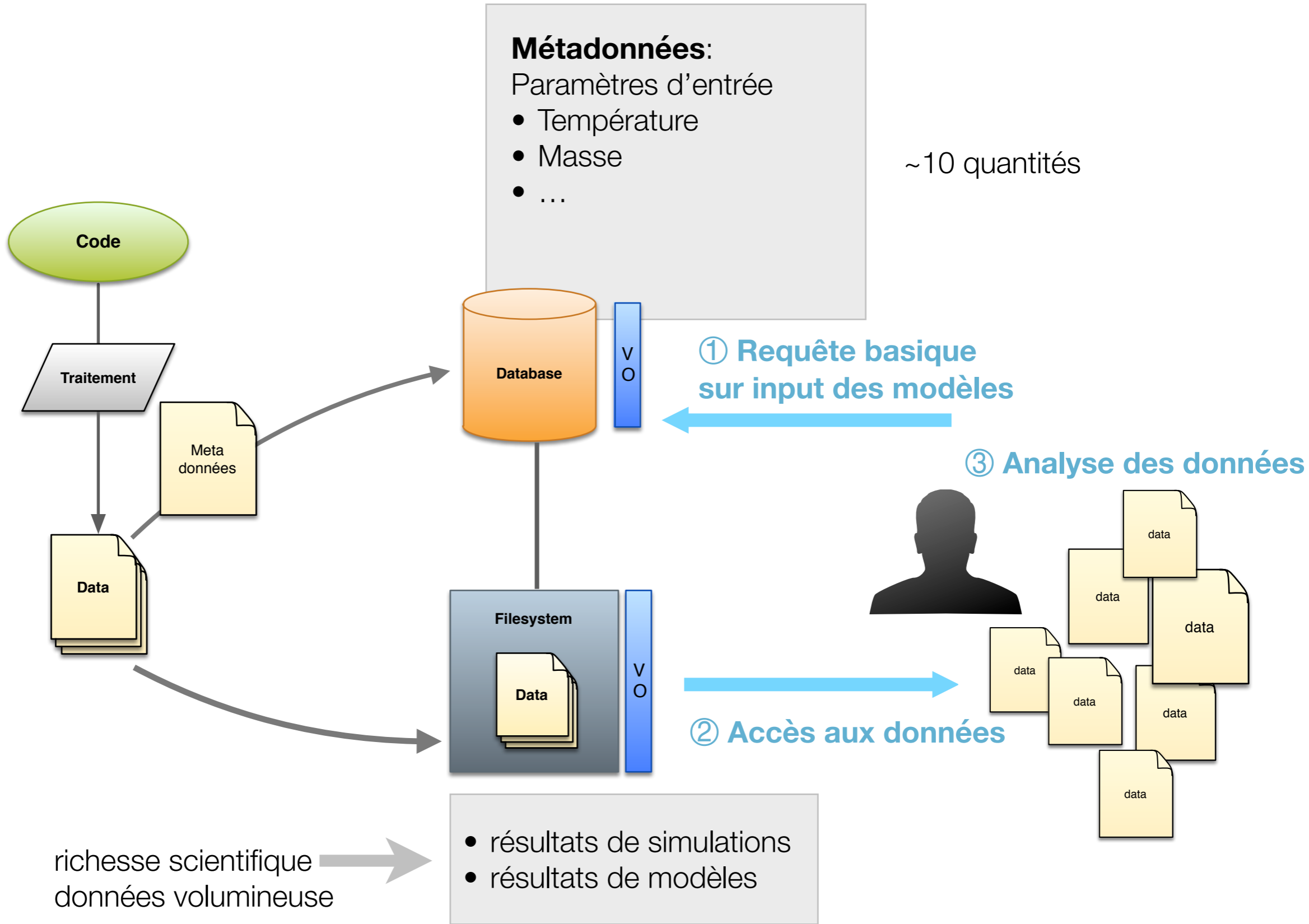① **Requête**
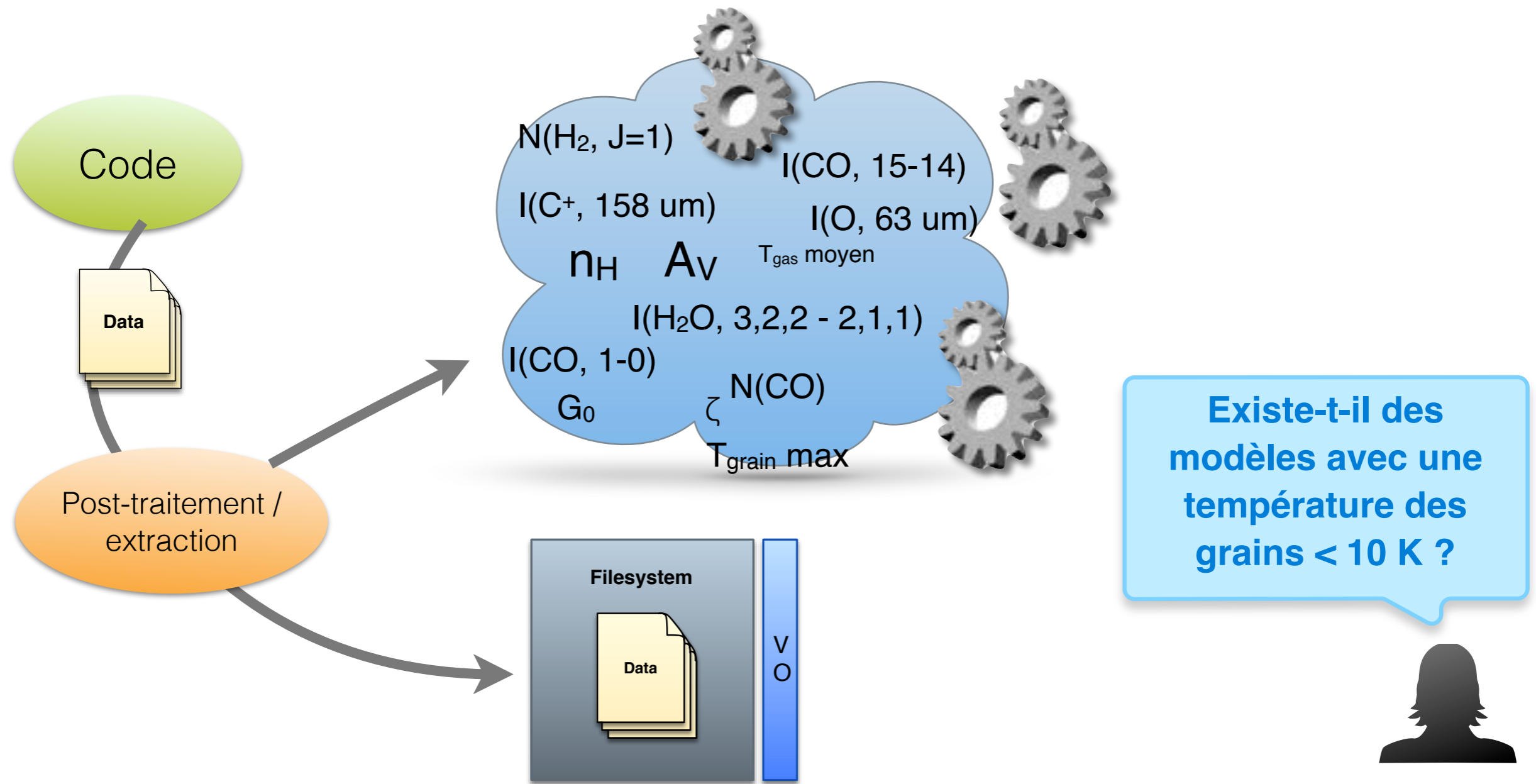
③ **Analyse des données**

data

Filesystem

Data

VO

② **Accès aux données**

richesse scientifique
données volumineuse

- spectres
- images
- cubes

# Astro data services: **Numerical simulations**



**Métadonnées**:
Paramètres d'entrée
- Température
- Masse
- …

~10 quantités

Code

Traitement

Meta données

Data

Database

V O

① **Requête basique sur input des modèles**

③ **Analyse des données**

Filesystem

Data

V O

② **Accès aux données**

data
data
data
data
data
data
data
data
data

richesse scientifique
données volumineuse

- résultats de simulations
- résultats de modèles

# What we would like



- 10 **code input quantities** can be queried
- thousands of models
- tens of thousands of data

→

- hundreds of thousands of **code output quantities** can be queried
- thousands of models
- millions of data

# Technical challenges

- Simulated data very heterogeneous:
  - **dimension nature** (mass, line intensity, x, y, z, …)
  - **number of dimensions** ($[10; 10^{5+}[$)
  - **number of objects** (dm halos nbr < particles nbr < …)

- Human-computer interaction while manipulating large meta-data amounts

# Technical challenges

**few common dimensions**

**huge number of dimensions**

**huge number of objects**

| obj | A | B | C | D | E | F | G | ... |
|---|---|---|---|---|---|---|---|---|
| 1 | | | 23 | | | | | |
| 2 | 2E+12 | | | 3.12e2 | | | | |
| 3 | | | | | 1.2 | 3E+12 | | |
| ... | | | | | | | | |
| $10^{15+}$ | | | | | | | | |

sparse matrix

# Technical challenges

**Unfortunately, no one tech solution to rule them all**

# Handle a lot of columns

The current services are built on top of relational db / SQL
Ex: One of the most used VO standard: Table Access Protocol (TAP)

**Problem**: the number of columns a relational database can handle is limited

|  | Table Size | Number of col. | Nbr of rows | Col. name size |
|---|---|---|---|---|
| **MySQL** | 64 Tb | 4096 | a lot | 64 |
| **Postgress** | 32 TB | 250 - 1600 | unlimited | 63 |
| **Oracle** | 4GB * block size | 1000 | unlimited | 30 |
| **Microsoft** | 524272 TB | 30 000 | limited by storage | 128 |

➡ The classical relational db approach doesn't fit the number of columns, neither the data heterogeneity very well

# How to manage a huge number of dimensions ?

Solutions ?

- noSQL & other new db designs
  - are the data more like documents than like table ?
  - do the db engine provide the management convenience we need ? (given the amount of data: clustering, memory setting etc…)
  - what logic must be moved to the application side when switching to schedules db ?

**Actually, the problem is just moved, not solved
The new problems are often harder to solve**

➡ **Be careful about new tech trends…
They are often intended to very specific cases, which are not yours**

# LISP (1958): association lists

(tag value)



v
20

B
5

type
C

shock

int(H2)
2e-7

cd(H2)
2e20

int(CO)
4e-8

Interstellar cloud

**Object = ( (tag1 value1) (tag2 value2) … )**

## Data tagging

- **Unlimited number of tags for an object**
  - Solve the high dimensionality challenge
- **Unlimited tag combinations to a given object**
  - Solve the dimensions heterogeneity challenge (sparse matrix)
- **Unlimited number of objects can be tagged**
  - Solve the large number of object challenge

**+**

- Abstract enough to be implemented on top of many technologies
  - RDBMS (EAV)
  - key/value engines
  - noSQL

**But at the cost of complex data query**

# Already used in biomed



- Used for a long time in Biomed
- Similar design used for years in RDF & Prolog

- Can benefit from robustness of RDBMS engine as implementation layer

There is **no silver bullet** here,
just choose the **right technology**
for the **right problem**

➡ For simulated data with few dimensions,
a classical relational schema is the best solution.

# What about huge number of objects ?

- Assume infinite collections
  - Use the **Stream abstraction** (infinite lazy list -> LISP)
- **Provide high level library** allowing easy handling of streams
  - ex: on top of a votable paginated api

# Human - machine interactions

How would a human being easily interact
with such a system ?

# Human - machine interactions

Interface of the VLA archive:



23 search parameters
classical "complexe" Interface

PDR services :
150 000+ parameters

Brian Glendenning, NRAO, InterOp Heidelberg 2013

# Human - machine interactions

2 steps

- What are the available dimensions I can query ?
- What query can I do against a dimension ?

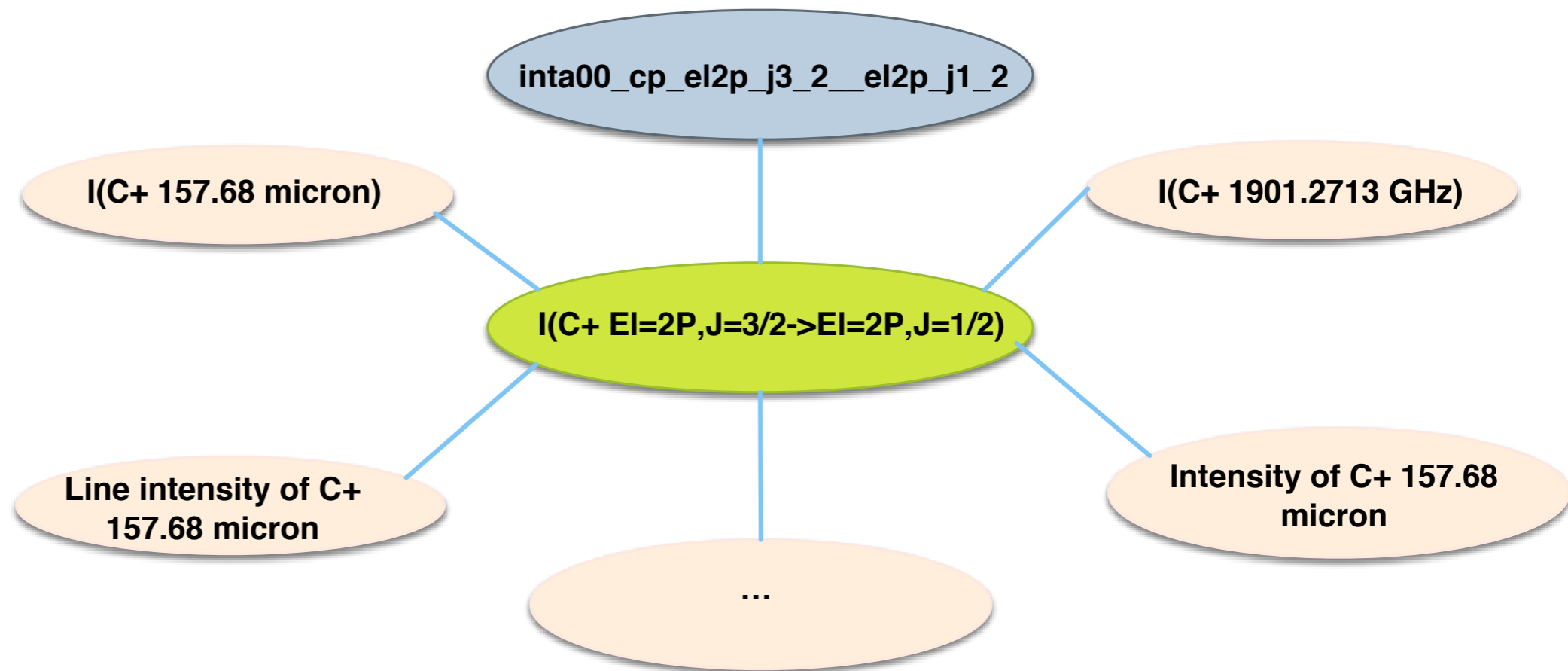Google France

Recherche Google     J'ai de la chance

+

Semantic web

# Dimension discovery

## Vocabulary for the PDR code

- The thousands of quantities handled by the PDR code are tagged
  - ID
  - human readable name
  - unit + description

- Creation of the synonyms list

- Currently : ~ 300 000 terms

# Dimension query

simple DSL: a tiny subset of SQL

**Axis constraints**

ex: N(Fe+) > 6e12

| Add | N(H2)| |

N(H2) > 8.0E20
N(H2) < 8.8E20
N(CO) > 1.0E13
N(CO) < 1.0E14
I(C+ El=2P,J=3/2->El=2P,J=1/2 angle 00 deg) >  3.6E-6

# Applications

- Observations analysis
- Statistical aggregated analysis on grid of models
- Cross grid (different codes) consolidation
- Machine learning (we have started that with Emetic Bron)

# In particular

- Today
  - **Manual features extraction**
  - Manage a lot of features

- Tomorrow
  - **Automatic features extraction** (+ scientist checking)
  - Pattern recognition / analysis
    (ex: generic quantities relationships)

# Demo

## Plot axis

x | nH (input parameter) | (cm-3) | ☑ log scale

y | ISRF scaling factor (back side) | (Mathis_unit) | ☑ log scale

## Fixed axis

AVmax | (mag) | 1

## Axis constraints

ex: N(Fe+) > 6e12

[ Add ]  N(H2)|

N(H2) > 8.0E20
N(H2) < 8.8E20
N(CO) > 1.0E13
N(CO) < 1.0E14
I(C+ EI=2P,J=3/2->EI=2P,J=1/2 angle 00 deg) >  3.6E-6
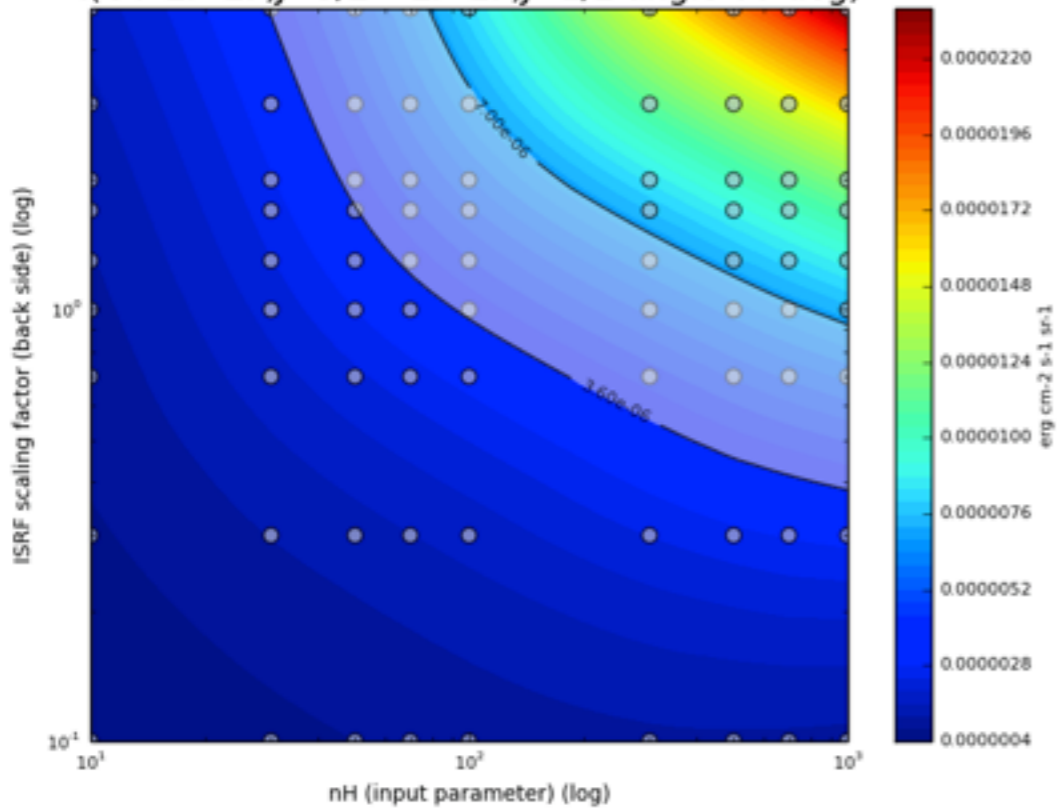
### N(H2)

**name**: N(H2)

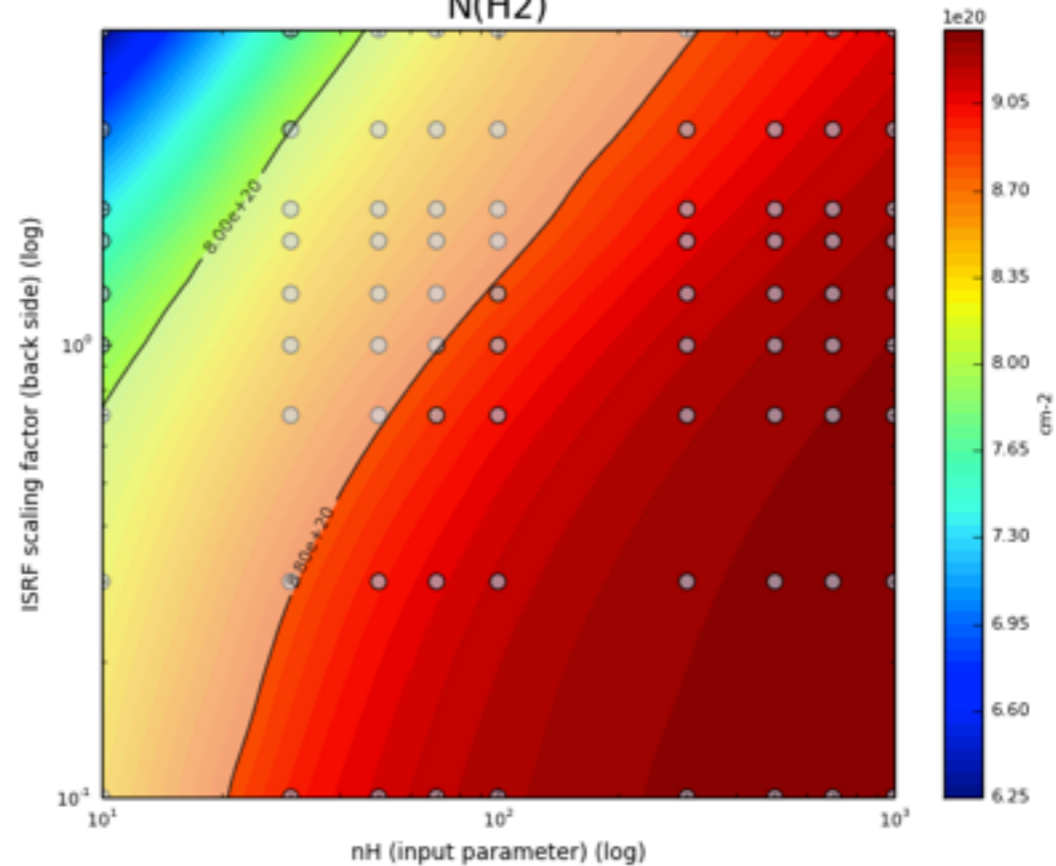**doc**: none

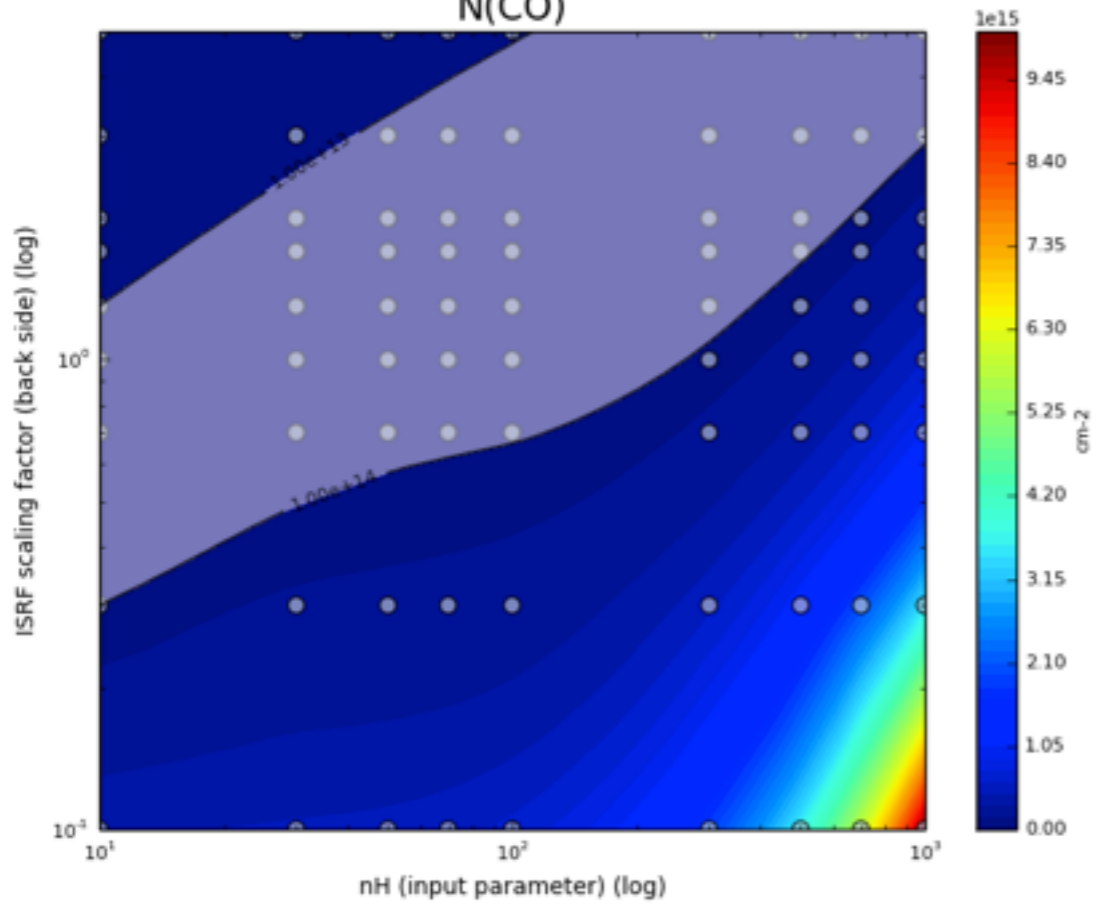**range**: [3.99e+19, 1.87e+21]

**unit**: cm-2

[ Plot ]

# Demo
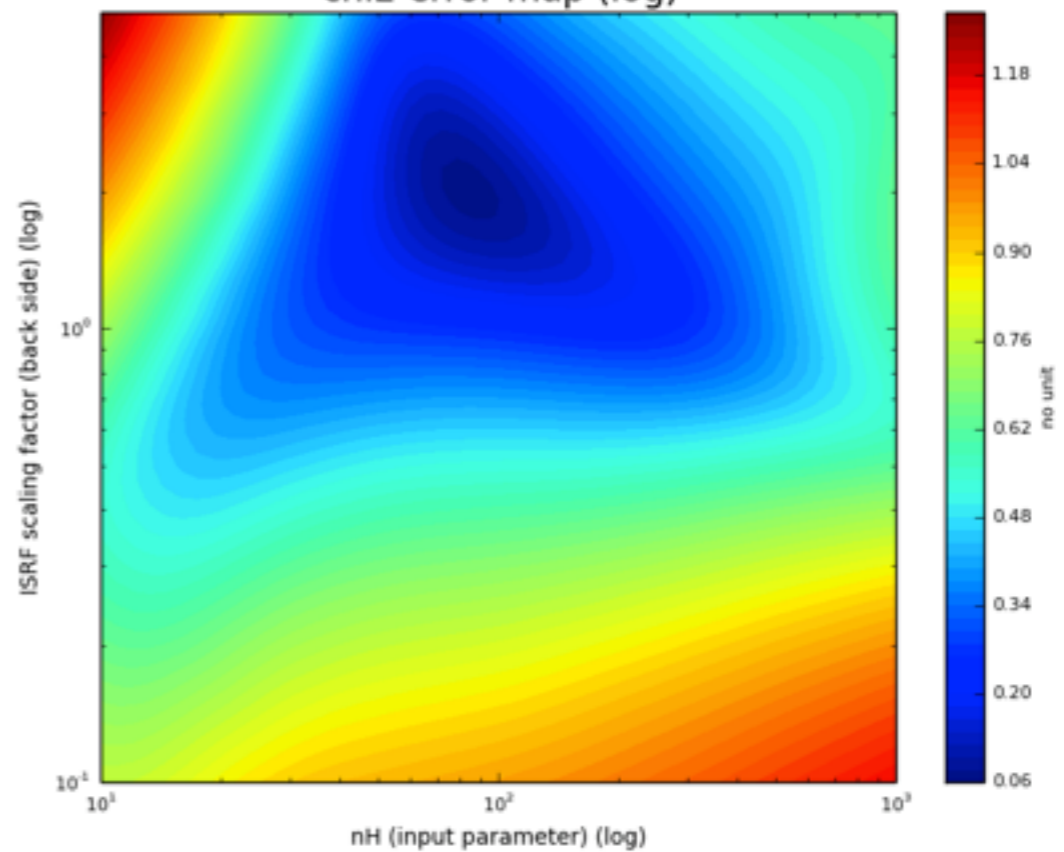
# Conclusion

- The "Big Data" is under-utilised.
- Use the machine to increase the ROI of processing time slots / big shared equipments.
- Help to identify the quantities where a real scientific work is required.
  - Machine Learning
- Consolidate the semantic dbs (vocabularies)
- **Document how the machine works / process (because there is no magic !)**

**Focus on the real scientific questions, let the machine do the dumb job**

**Summary:**

Quite heavy to set up for now
- Complexity shifted from the client to the server
- Exploratory programming: perimeter of the problem at hand not always well known/defined

Opens interesting opportunities