# Path to Open Science @ AGATA
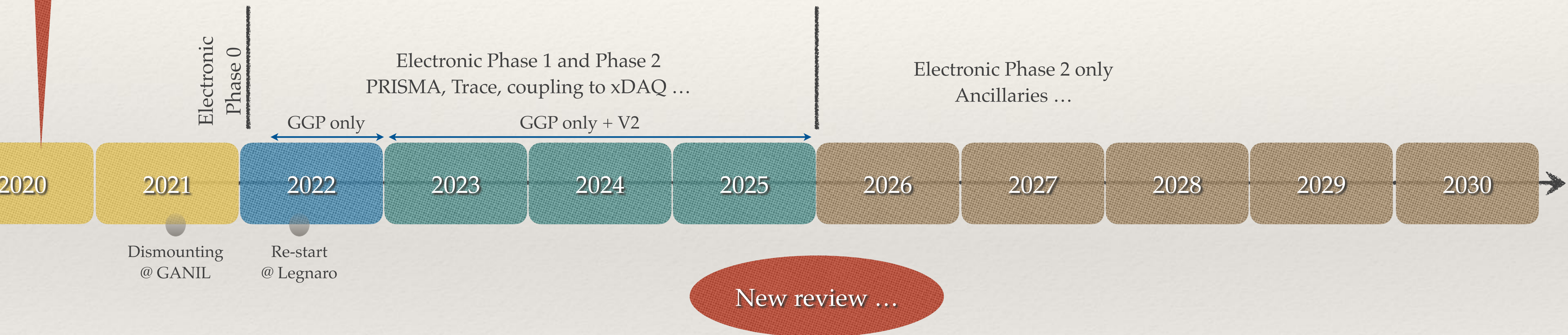
O.Stézowski,
as member of the AGATA Collaboration

Many thanks to the AGATA Data Processing Team

G.Baulieu, N.Dosme, J.Dudouet, S. Elloumi, Ph. Gauron, A. Goasduff, M.Gulmini,
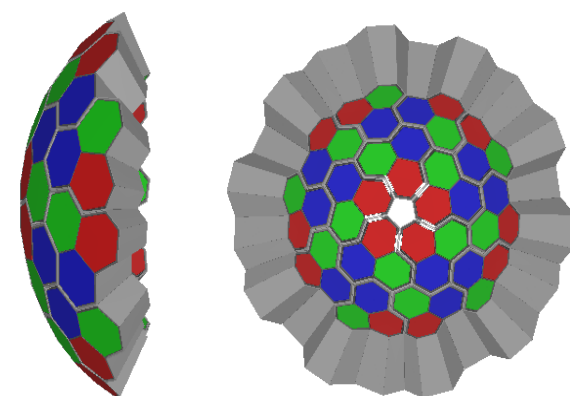A. Korichi,
J. Jacob, V. Lafage, E. Legay, P. Lejeannic, J. Ljungvall, G.Philippon, R.Molina, M. Roetto, M. Tauriga-Quere, N.Toniolo
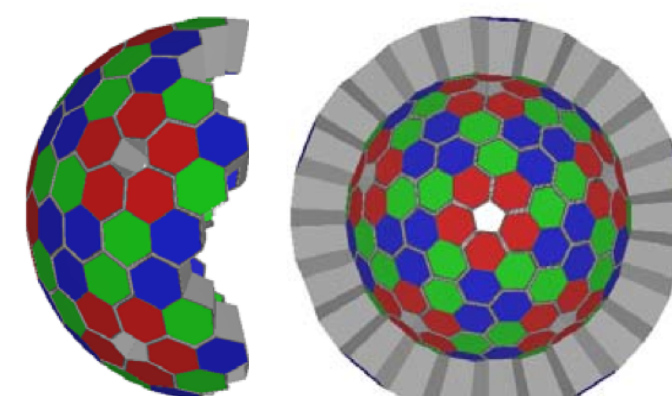
# The Data Processing « Phase 2 » Time Line

International Review
KDP2 IN2P3 } ⇒ New MoU AGATA Phase 2 (signed now)

➡ Data Management Plan <u>required</u>

Electronic
Phase 0

Electronic Phase 1 and Phase 2
PRISMA, Trace, coupling to xDAQ …

Electronic Phase 2 only
Ancillaries …

GGP only

GGP only + V2

| 2020 | 2021 | 2022 | 2023 | 2024 | 2025 | 2026 | 2027 | 2028 | 2029 | 2030 |

Dismounting
@ GANIL

Re-start
@ Legnaro

New review …



4PI/3
60 Ge
20 TC

2PI
90 Ge
30 TC

3PI
135 Ge
45 TC

# The Data Processing « Phase 2 » Time Line

International Review
KDP2 IN2P3  } ⇒ New MoU AGATA Phase 2 (signed now)
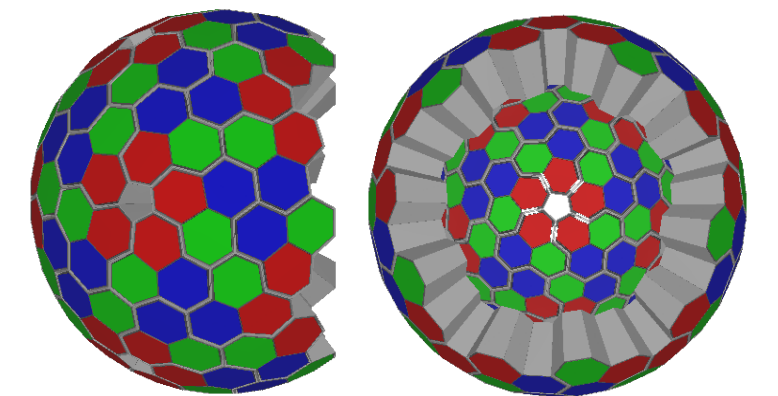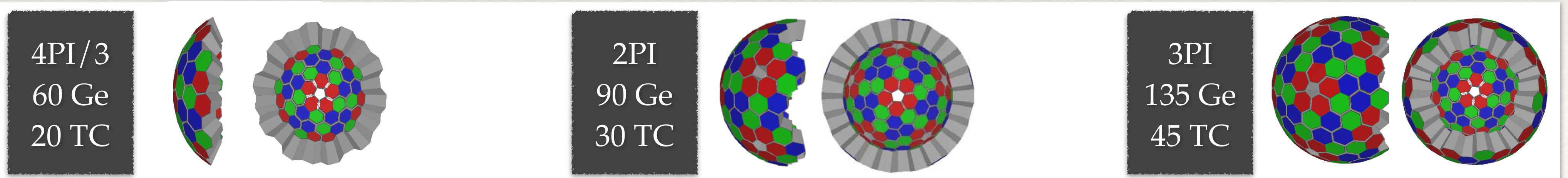
➡ Data Management Plan <u>required</u>

Electronic Phase 0

Electronic Phase 1 and Phase 2
PRISMA, Trace, coupling to xDAQ …

Electronic Phase 2 only
Ancillaries …

GGP only

GGP only + V2

| 2020 | 2021 | 2022 | 2023 | 2024 | 2025 | 2026 | 2027 | 2028 | 2029 | 2030 |

Dismounting
@ GANIL

Re-start
@ Legnaro

DMP ?

4 meetings

… a path to …



4PI/3
60 Ge
20 TC

2PI
90 Ge
30 TC

3PI
135 Ge
45 TC

# Step #1 : what is a DMP ?

It means we should establish a plan to deal with the cycle of life of the data produced by AGATA …

# Step #1 : what is a DMP ?

It means we should establish a plan to deal with the cycle of life of the data produced by AGATA …

Creation

Acquisition

Plan
- DMP -

From prior phase

Re-use
Discovery

Collect

Data life cycle

**Phase1**

AGATA
ADVANCED GAMMA
TRACKING ARRAY

So far we have worked here !
and not too badly
(From a DMP point of view)

To next phase

Share
Easy access

Processing
Analysis

Preserve
- Archive Destroy -

# Step #1 : what is a DMP ?

It means we should establish a plan to deal with the cycle of life of the data produced by AGATA …

New to be added for the Phase 2 [?]

Creation

Plan
DMP -

Acquisition

From prior phase

Discovery

Collect

Phase1

we worked here !
not too badly
(DMP point of view)

To next phase

Data life cycle

Share
Easy access

Processing
Analyse

Preserve
- Archive Destroy -

To be enhanced for the Phase 2 [?]

# Step #1 : what should we do ?

## How to fulfil the cycle ?



Cartoon from this presentation

# What AGATA produced as Data ?

Schematic view of the Data Production

# What AGATA produced as Data ?

Schematic view of the Data Production

Online Data stored (for distribution) on the grid
+ part of the meta-data

Data (meta-data) produced offline (@home)
**Lost or almost lost ... !**

Good, with have :
- Backups, access through VO,
- A policy to use data

...
**but** :
- No catalog
- Poor meta-data

...

Ancillaires

XXX capsules

```
1
0   1      0        0
  1   1      1      0   1
0   0  1    1    1   0   0
1   0  0  0  0      0   0
0   0  0  1  0      0   0
    1      1        1
1   1    1  1  1  1    1
1   1  .  .  .  1      1
.   .      .        .
.   .      .        .
.          .        .
```

```
1
0   1
0   0   1
1   0   0   0
0   1   1   1
1   1   1   1
1
```

**ONLINE**

A G A P R O

**OFFLINE**

P S A

T R A C K I N G

Are
the simulated data set,
the produced analysis,
the software,
the meta data
...
saved ? ... **NO** !!

Hardware / software

Hardware / software

meta-data electronic + ? Configurations

**Lost**

Most of the Software in svn, git or Hg ✓
... but are the version used to produce a particular data set saved ?

# Phase 1 : our workflow, more details

Producer
Intermediary
Consumer

AD ≡ AGATA Data

AD_Level_0
Traces

AD_Level_1
Hits

AD_Level_0
Ancillary Data

Ancillary

PrePSA    PSA    PostPSA

PrePSA    PSA    PostPSA

PrePSA    PSA    PostPSA

PCI

PCI

Production online @ different levels
Level 0 requires running again an experiment !
In principle from level 0
one can reproduce offline the other levels
ONLY IF we do have all the meta-data required

Builder    Merger    Tracking

AD_Level_2
Correlated Hits

AD_Level_3
Correlated Hits
+
Ancillary

AD_Level_4
Tracked Gamma
+    Ancillary
+ Correlated Hits

◄———— Local Level Processing ————►    ◄———— Global Level Processing ————►

# Phase 1 : our workflow, more details



Producer

Intermediary

Consumer

**AD ≡ AGATA Data**

**AD_Level_0**
**Ancillary Data**

Ancillary

**AD_Level_0**
**Traces**

**AD_Level_1**
**Hits**

Meta data saved so far are :
Online Configuration files (algo. parameters)
Online AGATA calibration files
*They are saved at the same place that the data*

PrePSA    PSA    PostPSA

PCI    PrePSA    PSA    PostPSA

PCI    PrePSA    PSA    PostPSA

Builder    Merger    Tracking

**AD_Level_4**
**Tracked Gamma**
+    **Ancillary**
+    **Correlated Hits**

**AD_Level_2**
**Correlated Hits**

**AD_Level_3**
**Correlated Hits**
+
**Ancillary**

←————— Local Level Processing —————→    ←————— Global Level Processing —————→

# First step toward a DMP for the Phase 2

## FAIRification Process !!!



Cleaning, documentation … ok … but what means FAIR ???

FAIR means **Findability**, **Accessibility**, **Interoperability**, **Reusability**

➥ FAIRification process, make sure the data (+meta) produced are FAIR
   ➥ **likely to have an impact on the way we produce, store etc … our data !**
➥ There are guidelines for that (see for instance https://www.go-fair.org/fair-principles/)
➥ let's have a look at some recommandations to be FAIR

### Findable

F1. Data (and meta-data) are assigned a globally unique and persistent identifier (PID)*
F2. **Data are described with rich meta data**
F3. Meta data clearly and explicitly include the identifier of the data they described
F4. (Meta)data are registered or indexed in a searchable resource

\* an example of PID is DOI. PID ≡ web page stored in a repository
   (See for instance zenodo)

F1. Obviously not the case
➥ *it might be good to start with at least a standard name for AGATA experiments*

F2. We have only a minimal amount of meta data **and** only for online data

F3. Our meta data are stored inside the data …
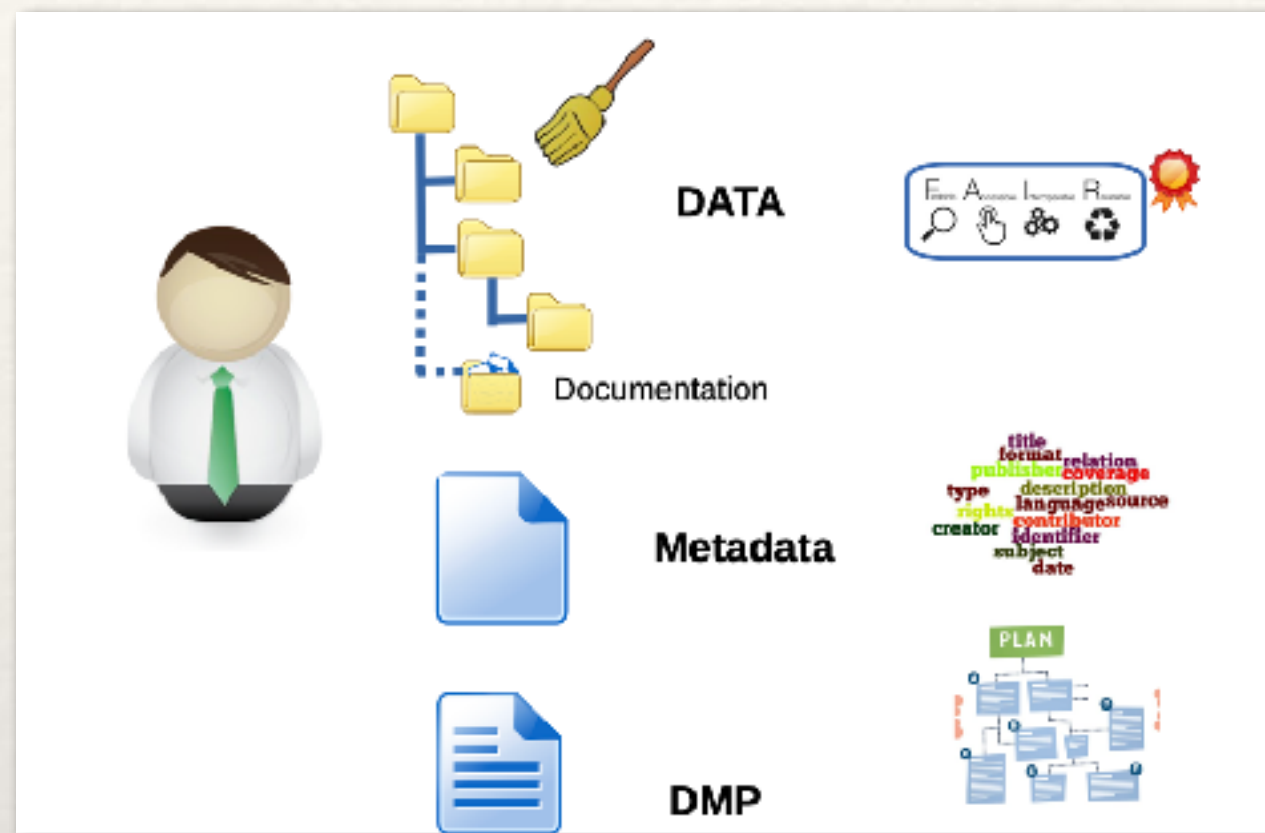➥ *metadata and data should be separated, see also A2*

F4 . Obviously not the case (see again zenodo)
➥ *searchable by humans and computers !*
➥ *Ex of a search: try and find all the data set produced at GANIL with NEDA ?*

# First step toward a DMP for the Phase 2

FAIRification of the data (see also the document called ADP-Part2, to be sent to the ACC)



Cleaning, documentation … ok … but what means FAIR ???

**FAIR** means **F**indability, **A**ccessibility, **I**nteroperability, **R**eusable

➡ FAIRification process, make sure the data (+meta) produced are FAIR

➥ **likely to have an impact on the way we produce, store etc … our data !**

➡ There are guidelines for that (see for instance https://www.go-fair.org/fair-principles/)

➡ let's have a look at some recommandations to be FAIR

Phase1

### Accessible

A1. (Meta)data are retrievable by their identifier using a standardised communication protocol
A1.1.The protocol is open, free and universally implementable
A1.2.The protocol, where necessary, allows for an authentification & authorisation procedure
A2. **Metadata are accessible, even when the data are no longer available**

A1. Grid access
➡ *difficult to retrieve a particular data set without browsing all*

A1.1 Grid access
➡ *Not completely universal, heavy for the collaboration, time to simplify if possible*
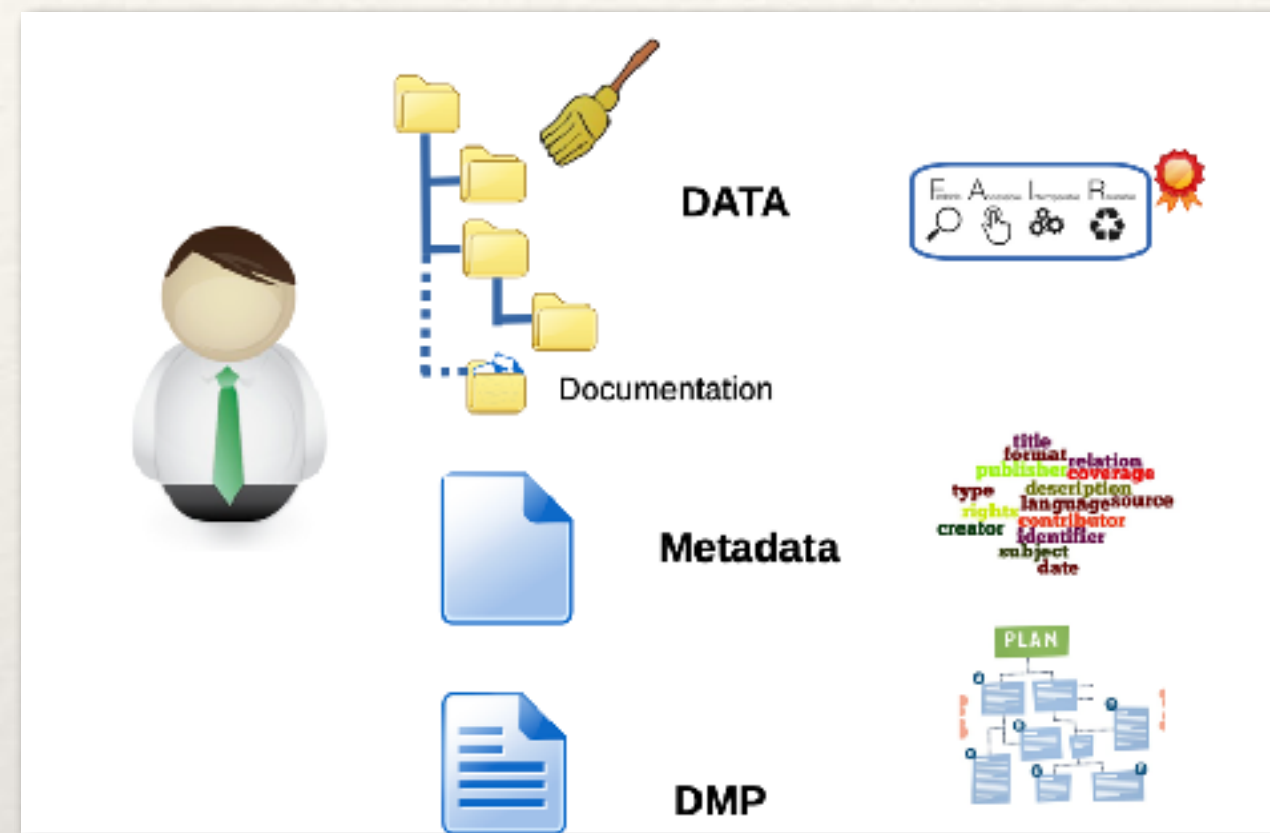
A1.2 AGATA Virtual Organisation
➡ *is this enough ?*

A2. Obviously not the case
➡ *AGAIN meta data should be separated from data*

# First step toward a DMP for the Phase 2

FAIRification of the data (see also the document called ADP-Part2, to be sent to the ACC)



Cleaning, documentation … ok … but what means FAIR ???

**FAIR** means **Findability**, **Accessibility**, **Interoperability**, **Reusable**

➡ FAIRification process, make sure the data (+meta) produced are FAIR

    ✍ **likely to have an impact on the way we produce, store etc … our data !**

➡ There are guidelines for that (see for instance https://www.go-fair.org/fair-principles/)

➡ let's have a look at some recommandations to be FAIR

### Re Usable

R1. (Meta)data are **richly described with a plurality of accurate and relevante attributes**
R1.1. (meta)data are released with a **clear and accessible usage licence**
R1.2. (meta)data are **associated with detailed provenance**
R1.3. (meta)data meet **domain-relevant community standards**

This is 'others' to play with AGATA data …
**Almost nothing done so far to help in that path …**

### Interoperable

I1.(meta)data use a normal, accessible, shared and broadly applicable language for knowledge representation
I2. **(meta)data use vocabularies that follow FAIR principles**
I3. Meta-data qualified references to other (meta)data

This is for integration of AGATA data with other data …
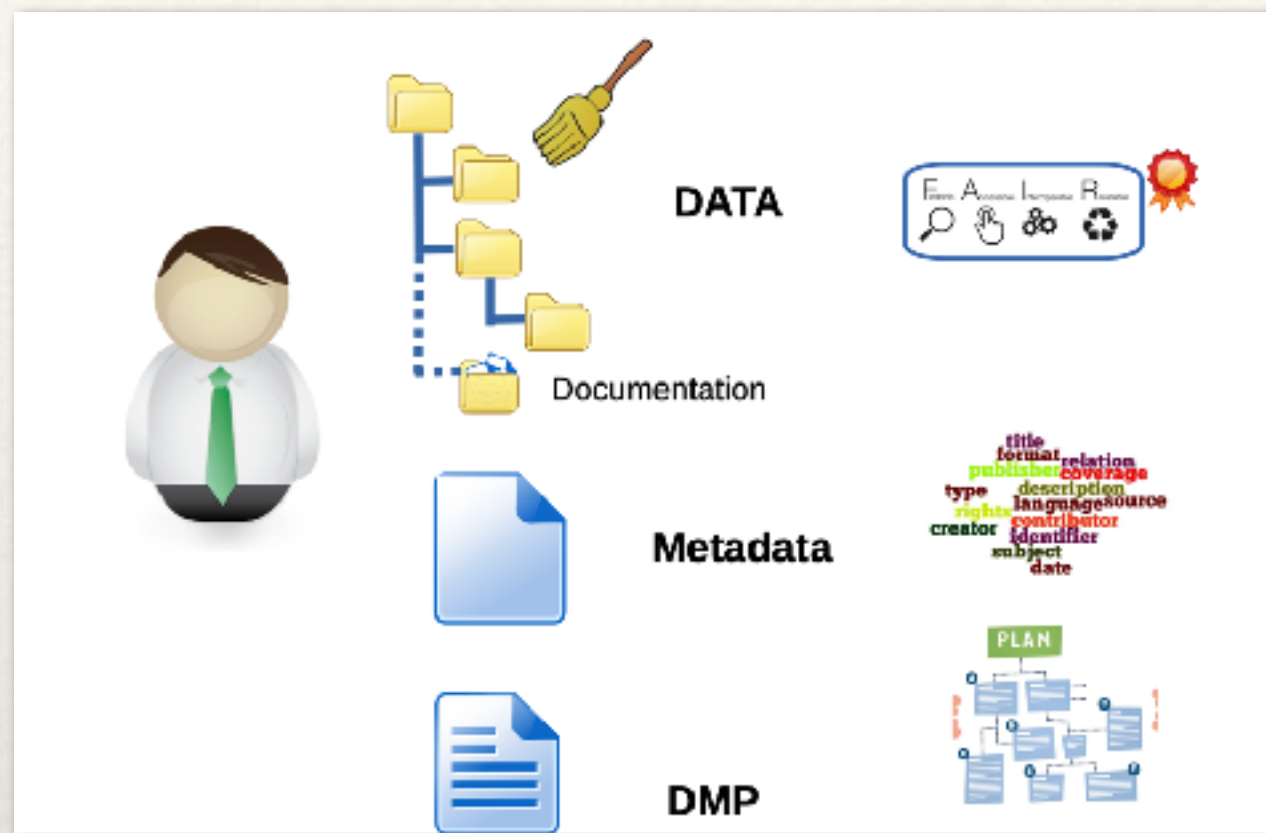**Almost nothing done so far to help in that path …**

# We do have a first DMP written ...



Many guidelines to really write a DMP
Here from the French ANR agency
  All are quite similar
  Moving from one to another one should be easy

1. **Data description and collection, re-use of existing data**

A. How will new data be collected or produced and/or **how will existing data be re-used?**

B. What data (for example the kinds, formats, and volumes) will be collected or produced?

2. **Documentation and data quality**

A. **What metadata** and documentation (for example the methodology of data collection and way of organizing data) will accompany data?

B. What data quality control measures will be used?

3. **Storage and backup during research process**

A. **How will data and metadata be stored and backed up during the research process?**

B. How will data security and protection of sensitive data be taken care of during the research?

4. **Legal and ethical requirements, codes of conduct**

A. if personal data are processed, how will compliance with legislation on personal data and on data security be ensured?

B. How will other legal issues, such as **intellectual property rights and ownership, be managed?** What legislation is applicable?

C. How will possible ethical issues be taken into account, and codes of conduct followed?

5. **Data sharing and long-term preservation**

A. **How and when will data be shared? Are there possible restrictions to data sharing or embargo reasons?**

B. **How will data for preservation be selected, and where will data be preserved long-term (for example a data repository or archive)?**

C. **What methods or software tools will be needed to access and use the data?**

D. How will the application of a unique and persistent identifier (such as a Digital Object Identifier (DOI)) to each data set be ensured?
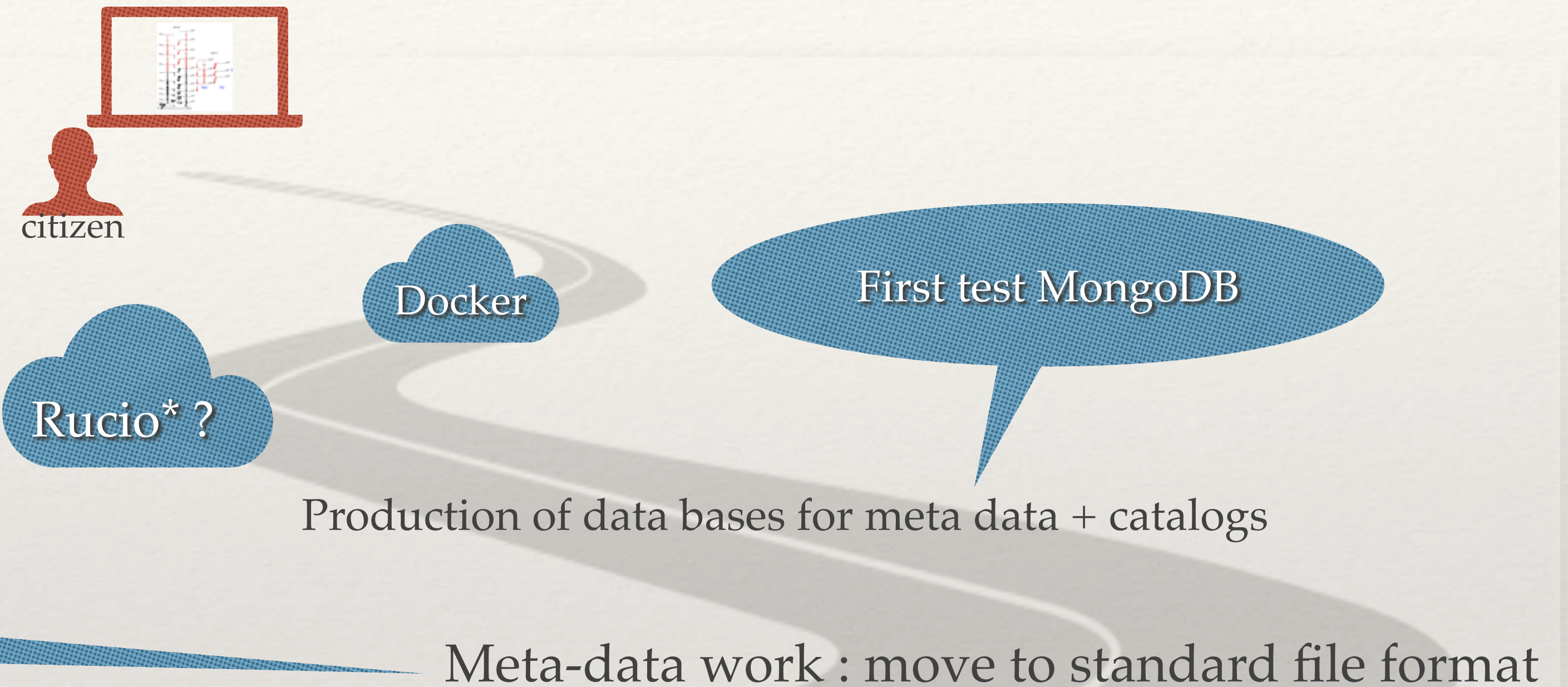
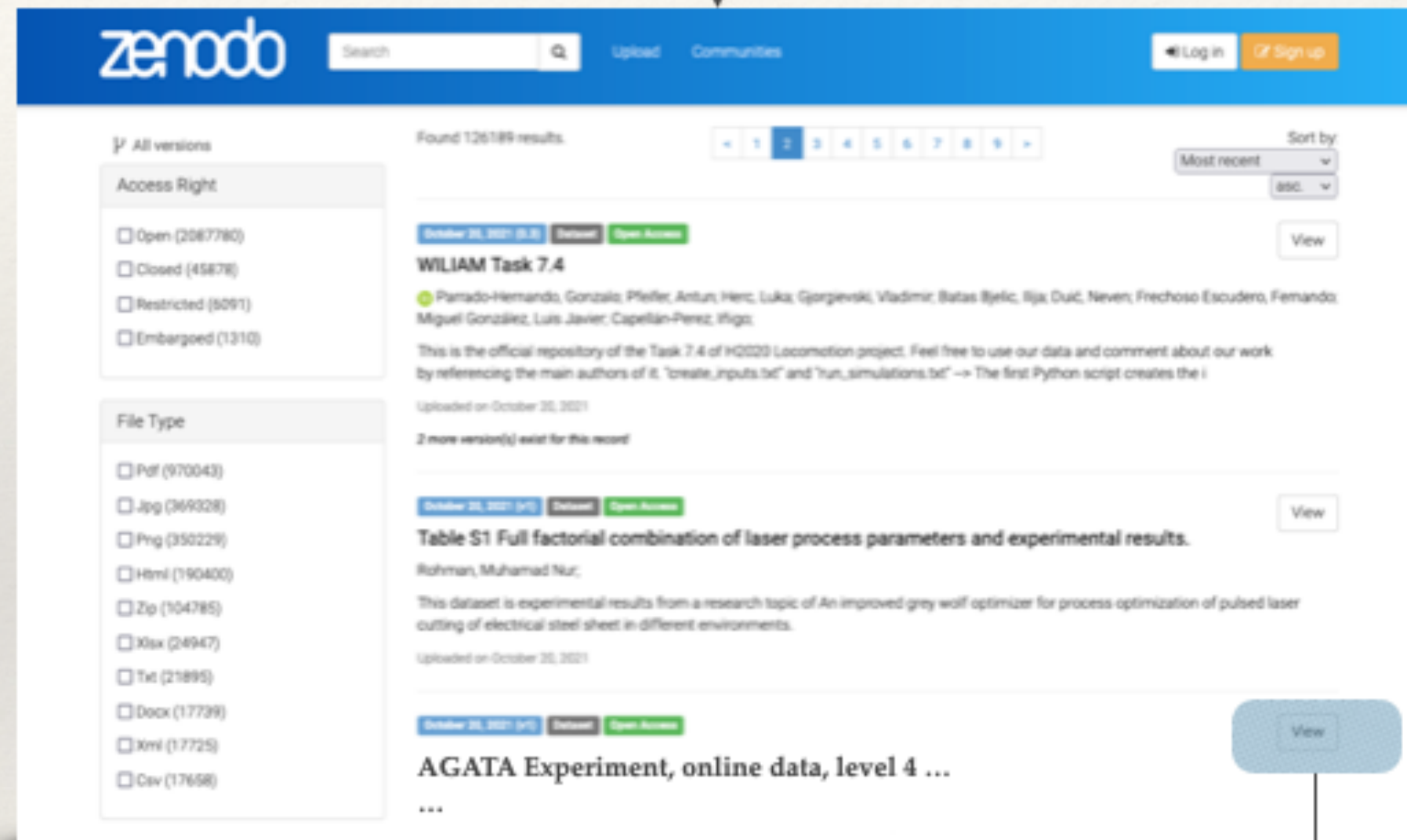6. **Data management responsibilities and resources**

A. Who (for example role, position, and institution) will be **responsible for data management** (i.e. the **data steward**)?

B. What resources (for example financial and time) will be dedicated to data management and ensuring that data will be FAIR (Findable, Accessible, Interoperable, Re-usable)?

AGATA
Data Policy
for
Publications

# Conclusions



citizen

Docker

First test MongoDB

Rucio* ?

Production of data bases for meta data + catalogs

Json instead of user's format

Meta-data work : move to standard file format

❷ We need to draw it using the DMP (A DMP is regularly modified)
We do not start from 0, we have identified good / bad practices
☞ progressive modification (FAIRification) of our way to work

❶ We are at the beginning of the 'open' path

* https://rucio.cern.ch/

# Thank you for listening

## Questions ?

# Phase 1 : our practices so far …

+ femul on grid

PC

Zone
Partagée

V2

V1

V0

Cluster

R&D PSA - Tracking          Ancillaires          DMP

Emulator [femul] - Offline 'infrastucture' - NO REAL TIME

PSA

Tracking

+

Ancillaire

Analyses de Physique

DAQ BOX [DCOD] - online 'infrastructure' - REAL TIME

Electonique AGATA          Host laboratory network          PSA / Tracking          Many Ancillaries