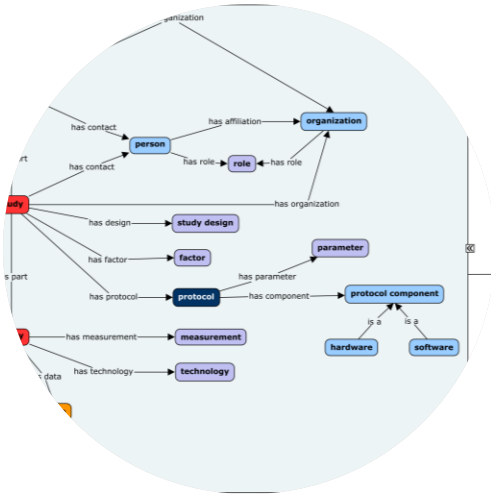


# FAIR Data & ISA standard



Using the ISA standard for collecting and sharing data

2024, Sven Warris & Rick van de Zedde



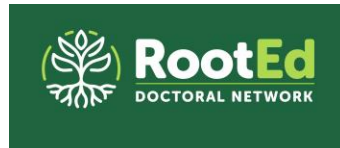
# Data challenges in genomics & phenomics

- Collecting and pre-processing terabytes of data
  - Single and combined experiments
  - Many different sensors, technology platforms
  - Many different pre-processing steps
- Sharing and publishing data
- Metadata is as important as the data itself
  - Organisms, treatment, samples, etc
  - Sensor type, settings, etc

# Project partners



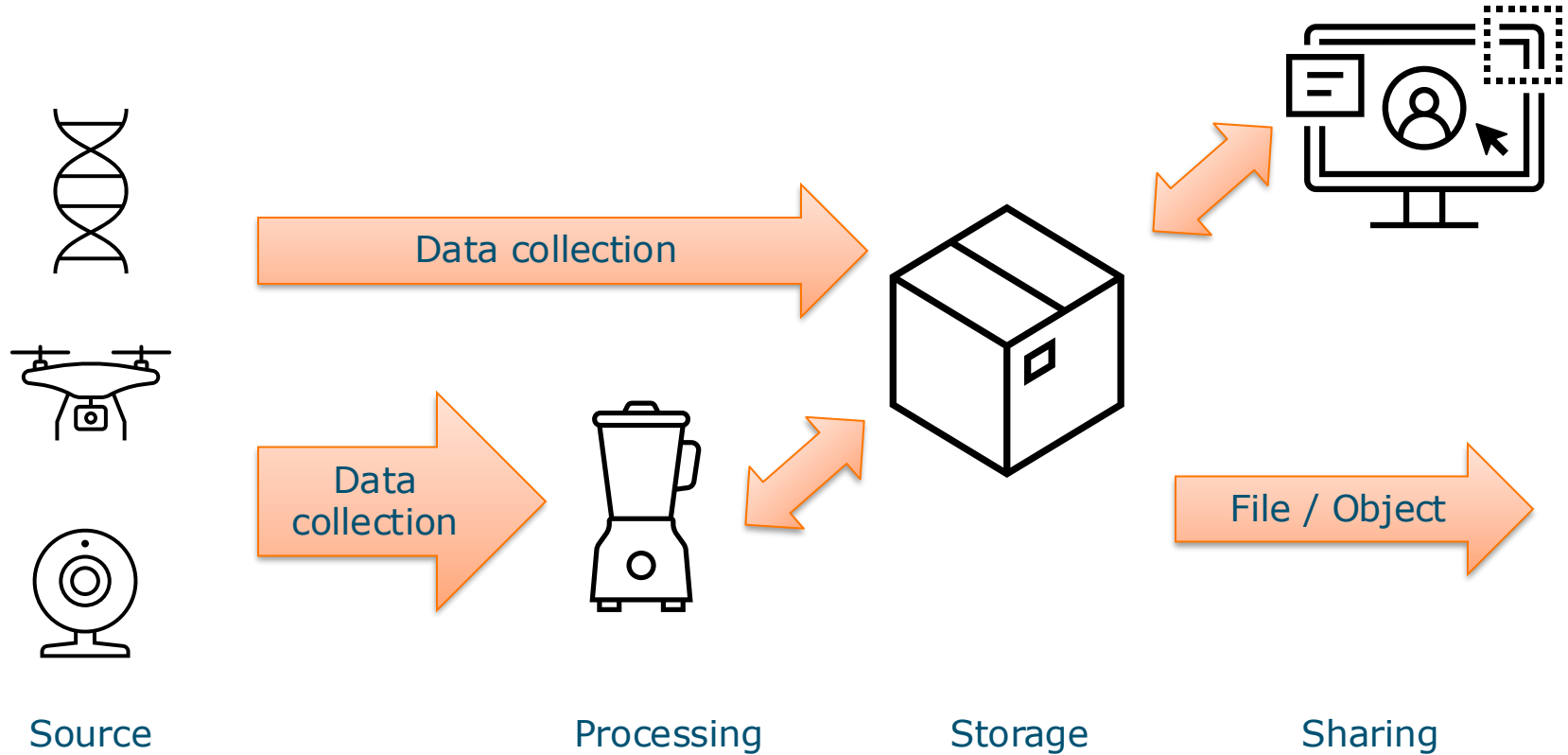
- Bioscience
- Biointeractions
- Biometrics



FB-Information Technology

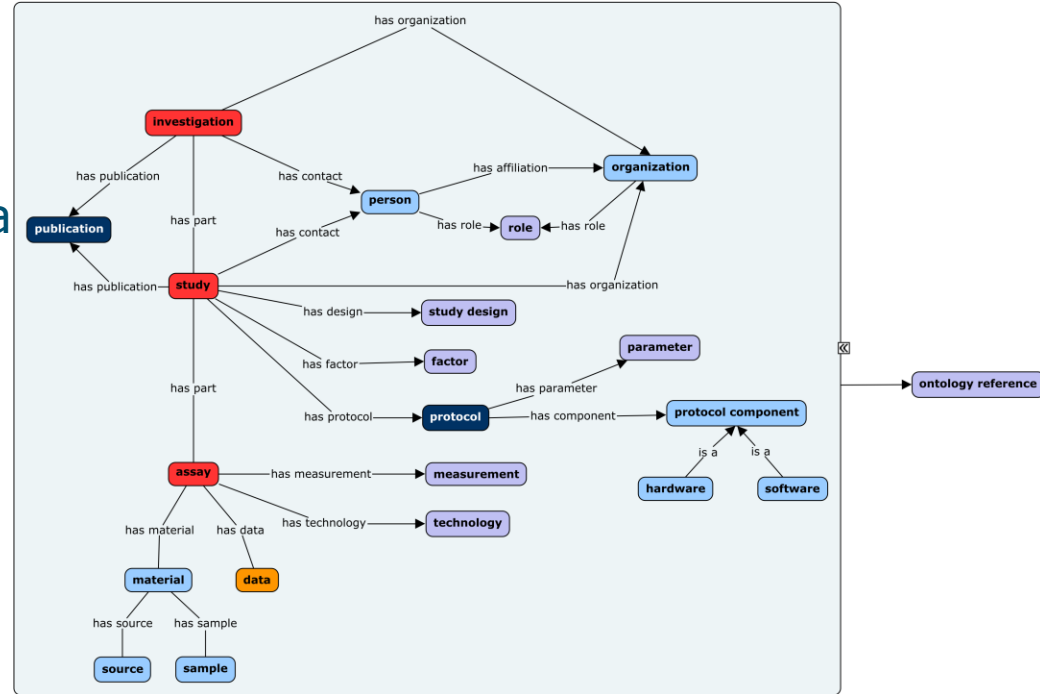


# Data flow



# Investigation / study / assay (ISA)

- Standardized way of structuring project metadata
- Data files, documents, etc
- Ontology-based
- Important entities:
  - Person, Assay, Data



- <https://isa-specs.readthedocs.io/en/latest/isamodel.html>

# Investigation / study / assay

- Structure is the same in, for example:
  - eLabJournal
  - FAIRDOM-seek
- Python support through isatools:
  - ISAJSON / ISATAB
- R support:
  - ISATAB
  - ISAJSON currently being developed

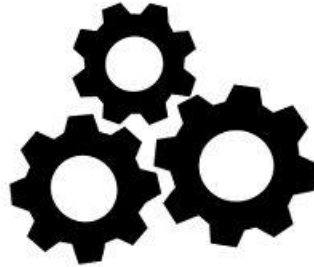
F  
indable



A  
ccessible



I  
nteroperable



R  
eusable



# FAIRDOM-seek

Open-source platform designed for the cataloguing and sharing of diverse scientific research data, including datasets, models, simulations, processes, and outcomes.

- **Preserves Associations:** It maintains relationships between various research components along with information about the involved people and organizations
- **ISA Infrastructure:** FAIRDOM-SEEK structures how experiments are part of broader studies and investigations
- **Configurable Structure:** structure is adaptable, making it suitable for various scientific fields



- **Sharing Permissions:** flexible and detailed sharing permissions, supporting collaboration at different research stages, from initial collaboration to publishing final results
- **DOI Generation:** Digital Object Identifiers for individual items or entire collections packaged as Research Objects
- **Semantic Technology:** Advanced queries over its content
- **Metadata Collection:** standard Excel tools and processes (RightField)
- **MIAPPE** support: metadata from *Minimum Information About Plant Phenotyping Experiments*

# Metadata implementation (MIAPPE)

```
20220111-WUR_test24-metadata.json x
1 {
2   "Name": "Minimum Information About Plant Phenotyping Experiment",
3   "Models": [
4     {
5       "Name": "Environment",
6       "Definition": "Environment",
7       "Attributes": [
8         {
9           "Name": "Altitude",
10          "Line": "ENV-19",
11          "Definition": "Altitude",
12          "Value": "14.2",
13          "Format": "Degrees in the decimal degrees",
14          "Marker": "EnvironmentIdentity"
15        },
16        {
17          "Name": "Organism",
18          "Line": "ENV-20",
19          "Definition": "Organism",
20          "Value": "14.2",
21          "Format": "Degrees in the decimal degrees",
22          "Marker": "EnvironmentIdentity"
23        },
24        {
25          "Name": "Light",
26          "Line": "ENV-21",
27          "Definition": "Light",
28          "Value": "9.8",
29          "Format": "m (metre)",
30          "Marker": "EnvironmentIdentity"
31        }
32      ],
33      "Name": "Plot",
34      "Definition": "Plot position and size",
35      "Attributes": [
36        {
37          "Name": "Geographic location (DM-19)",
38          "Line": "DM-19",
39          "Definition": "Coordinate point",
40          "Value": "51.988212150",
41          "Format": "Degrees in the decimal degrees",
42          "Marker": "PlotIdentity"
43        },
44        {
45          "Name": "Geographic location (DM-20)",
46          "Line": "DM-20",
47          "Definition": "Coordinate point",
48          "Value": "5.661708225",
49          "Format": "Degrees in the decimal degrees",
50          "Marker": "PlotIdentity"
51        },
52        {
53          "Name": "Geographic location (DM-21)",
54          "Line": "DM-21",
55          "Definition": "INS height of plot",
56          "Value": "9.8",
57          "Format": "m (metre)",
58          "Marker": "PlotIdentity"
59        }
60      ],
61      "Name": "Observation unit ID",
62      "Line": "DM-70",
63      "Definition": "Plotname",
64      "Value": "WUR_test24",
65      "Format": "Text",
66      "Marker": "PlotIdentity"
67    },
68    {
69      "Name": "Plot coordinate A (Longitude)",
70      "Line": "",
71      "Definition": "",
72      "Value": "5.661030201",
73      "Format": "Degrees in the decimal degrees",
74      "Marker": "ImageContext"
75    },
76    {
77      "Name": "Plot coordinate B (Latitude)",
78      "Line": "",
79      "Definition": "",
80      "Value": "51.989074658",
81      "Format": "Degrees in the decimal degrees",
82      "Marker": "ImageContext"
83    },
84    {
85      "Name": "ImagingSession",
86      "Definition": "Imaging session parameters",
87      "Attributes": [
88        {
89          "Name": "Event accession",
90          "Line": "DM-66",
91          "Definition": "Foldernam",
92          "Value": "WUR_test24-20220111-09:58",
93          "Format": "Text",
94          "Marker": "ImageIdentity"
95        },
96        {
97          "Name": "Event date",
98          "Line": "DM-68",
99          "Definition": "Date and time",
100         "Value": "20220111-09:58",
101         "Format": "Date/Time (ISO 8601)",
102         "Marker": "ImageContext"
103       },
104       {
105         "Name": "Plot entry coordinate A",
106         "Line": "",
107         "Definition": "",
108         "Value": "5.661021788",
109         "Format": "Degrees in the decimal degrees",
110         "Marker": "ImageContext"
111       },
112       {
113         "Name": "Plot entry coordinate B",
114         "Line": "",
115         "Definition": "",
116         "Value": "51.989074658",
117         "Format": "Degrees in the decimal degrees",
118         "Marker": "ImageContext"
119       }
120     ],
121     {
122       "Name": "LIDAR",
123       "Definition": "Settings and specifications of LIDAR",
124       "Attributes": [
125         {
126           "Name": "LIDAR used",
127           "Line": "",
128           "Definition": "",
129           "Value": "Yes",
130           "Format": "Text",
131           "Marker": "DeviceSetting"
132         },
133         {
134           "Name": "LIDAR type",
135           "Line": "",
136           "Definition": "",
137           "Value": "Sick LMS400-2000",
138           "Format": "Text",
139           "Marker": "DeviceProperty"
140         },
141         {
142           "Name": "LIDAR product number",
143           "Line": "",
144           "Definition": "Manufacturers product number",
145           "Value": "1041725",
146           "Format": "Text",
147           "Marker": "DeviceIdentity"
148         },
149         {
150           "Name": "LIDAR serial number",
151           "Line": "",
152           "Definition": "Serial number of sensor",
153           "Value": "14420197",
154           "Format": "Text",
155           "Marker": "DeviceIdentity"
156         }
157       ]
158     }
159   ]
160 }
20220111-WUR_test24-metadata.json x
20220111-WUR_test24-metadata.json x
20220111-WUR_test24-metadata.json x
```

# Extended ISA directory structure

- With thousands of measurements per experiment, the ISA structure is not suited for file storage
  - Operating / File systems and file browsers cannot deal with that many files in a single directory
- Extended ISA structure:  
Experiment / exp[id] / [Sample type] / [Pot ID] / [Assay timestamp] / Imaging / PlantEye / [data type]
- Makes data browsable (aka Findable) again

# Extended ISA directory structure

- ▼ Ryegrass\_Experiment21\_Gantry
  - ▼ Experiment21
    - animations
    - derived
  - ▼ Pot
    - ▼ NPEC54.20220822TW.CK1.Bar52.Drought.1
      - ▼ 20220822T140847
        - ▼ Imaging
          - ▼ PlantEye
            - derived
            - pointcloud
    - > 20220822T230620
    - > 20220823T060618

Name

- f00067\_20220822T140847\_full\_sx000\_sy000.ply.gz
- f00067\_20220822T140847\_full\_sx000\_sy000.ply.gz.ndvi.PNG
- f00067\_20220822T140847\_full\_sx000\_sy000.ply.gz.png
- f00067\_20220822T140847\_mg\_sx000\_sy000.ply.gz
- f00067\_20220822T140847\_mr\_sx000\_sy000.ply.gz
- f00067\_20220822T140847\_sl\_sx000\_sy000.ply.gz

# ISA-JSON

- JSON file containing:
  - Project metadata
  - Samples & organisms
  - Sensor technologies
  - Location of the data files
  - And the type for each file (raw, derived)
- Computer (and little bit human) readable format
- Readily in and out Python & R isatools

## Assay (part)

```
{
  "@id": "#sample/a721c29a-5f86-4921-a5ce-024c39833d04",
  "characteristics": [
    {
      "category": {
        "@id": "#characteristic_category/040b5fbf-5d50-4631-ad8f-8d509ddcb0a5"
      },
      "comments": [],
      "value": {
        "@id": "#ontology_annotation/cb10d77d-dac2-4d6f-a218-e571e730208d",
        "annotationValue": "Plant",
        "comments": [],
        "termAccession": "http://purl.bioontology.org/ontology/NCBITAXON/33090",
        "termSource": "NCBITaxon"
      }
    }
  ],
  "comments": [],
  "derivesFrom": [
    {
      "@id": "#source/1cf2584b-e509-46ff"
    }
  ],
  "factorValues": [],
  "name": "Fl_6"
},
```

```
{
  "@id": "#data_file/0f65eadf-f045-4c0c-af40-bf91970627a9",
  "comments": [
    {
      "name": "fullPath",
      "value": "running_exp28_Gantry4/exp28/Pot/F2_5/20230725T111409/Imaging/PlantEye/derived/20230725T111409_psri.csv"
    }
  ],
  "name": "20230725T111409_psri.csv",
  "type": "Derived Data File"
},
],
"filename": "20230725T111409",
"materials": {
  "otherMaterials": [],
  "samples": [
    {
      "@id": "#sample/a575a1be-6b2b-4b41-af26-c82fcabb6aa8"
    }
  ]
},
"measurementType": {
  "@id": "#ontology_annotation/2279abf0-7b79-478d-a11f-7da48bf11403",
  "annotationValue": "",
  "comments": [],
  "termAccession": "",
  "termSource": ""
},
"processSequence": [],
"technologyPlatform": "PlantEye",
"technologyType": {
  "@id": "#ontology_annotation/4dbb5726-90f0-4df2-bb43-eb16af6a3b6a",
  "annotationValue": "Imaging",
  "comments": [],
  "termAccession": "http://jermontology.org/ontology/JERMontology#Imaging",
  "termSource": ""
},
},
```

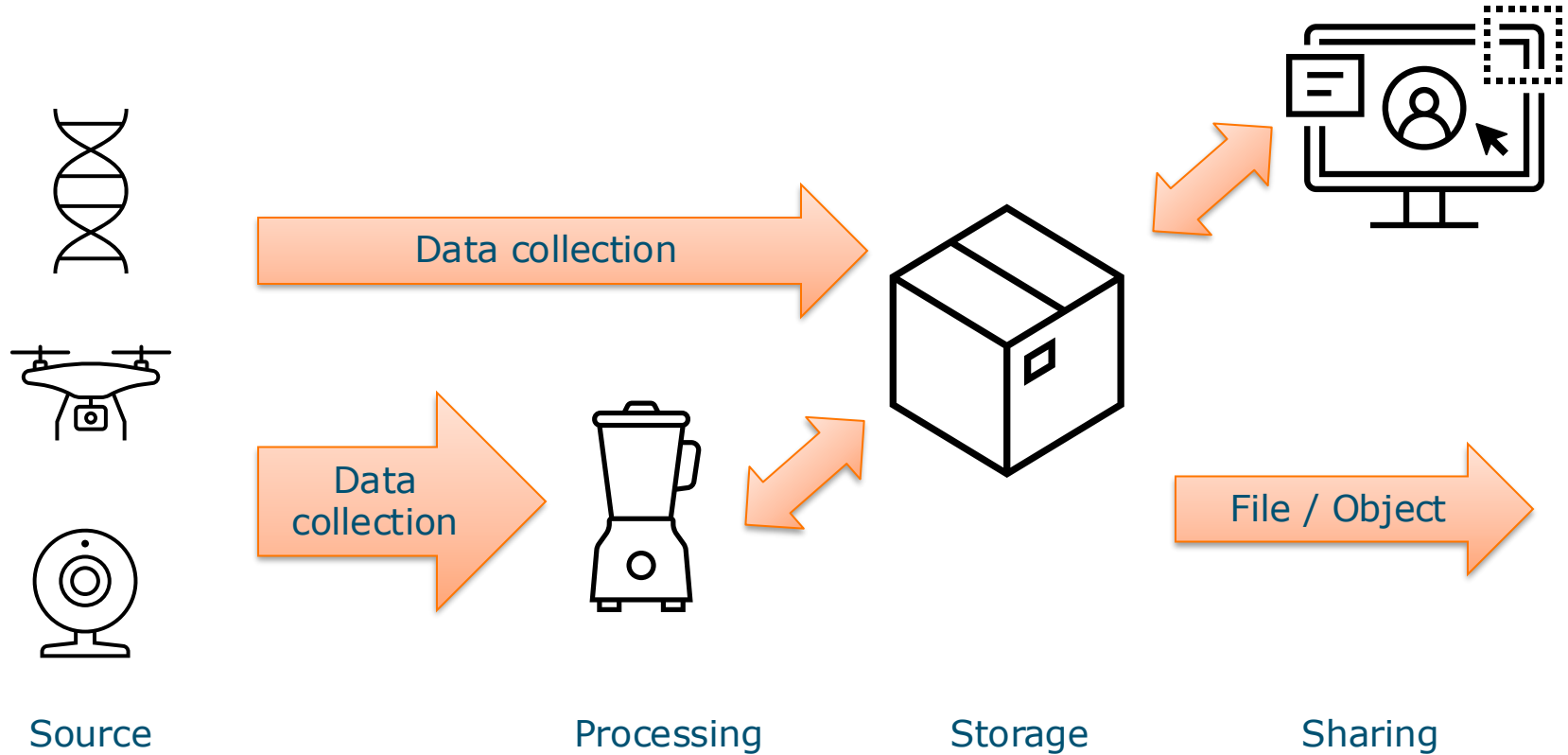
## Sample

# Use of ISA and ISA-JSON

- Adding pre-processing steps:
  - Get the relevant files
  - Process and store the location also in JSON
- For researchers:
  - Direct access to the relevant data files
  - Sharing makes the data more FAIR
- IT:
  - Use the metadata to store, archive or provide access to data



# Data flow implementation





# Data flow implementation



NPEC



PHENET

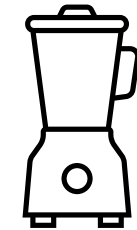
PHENOTYPING & ENVIROTYPING  
SOLUTIONS FOR AGROECOLOGY

Source

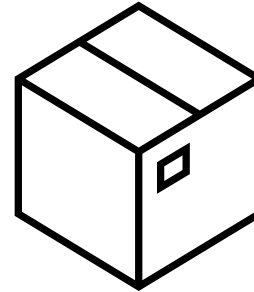


Data collection

Data  
collection



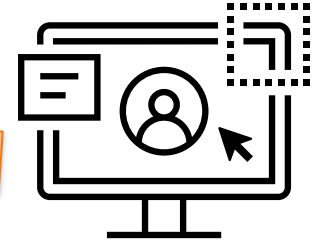
Processing



Storage

File / Object

Sharing



# Data flow implementation



NPEC



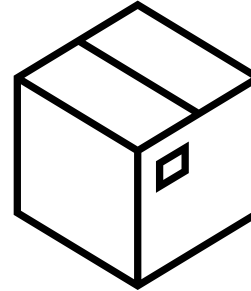
PHENET

PHENOTYPING & ENVIROTYPING  
SOLUTIONS FOR AGROECOLOGY

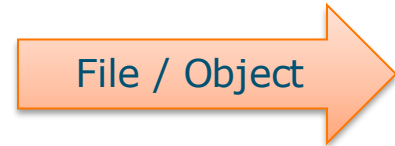
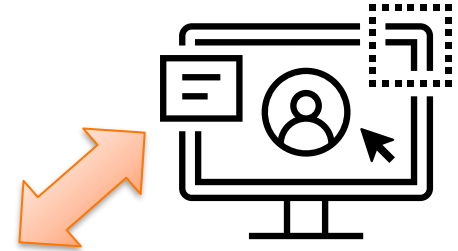
Source



Processing



Storage



Sharing

# Data flow implementation



NPEC



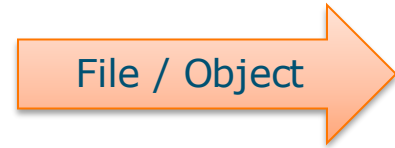
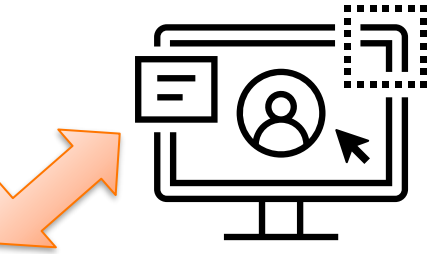
PHENET

PHENOTYPING & ENVIROTYPING  
SOLUTIONS FOR AGROECOLOGY



lustre®

iRODS®



Source

Processing

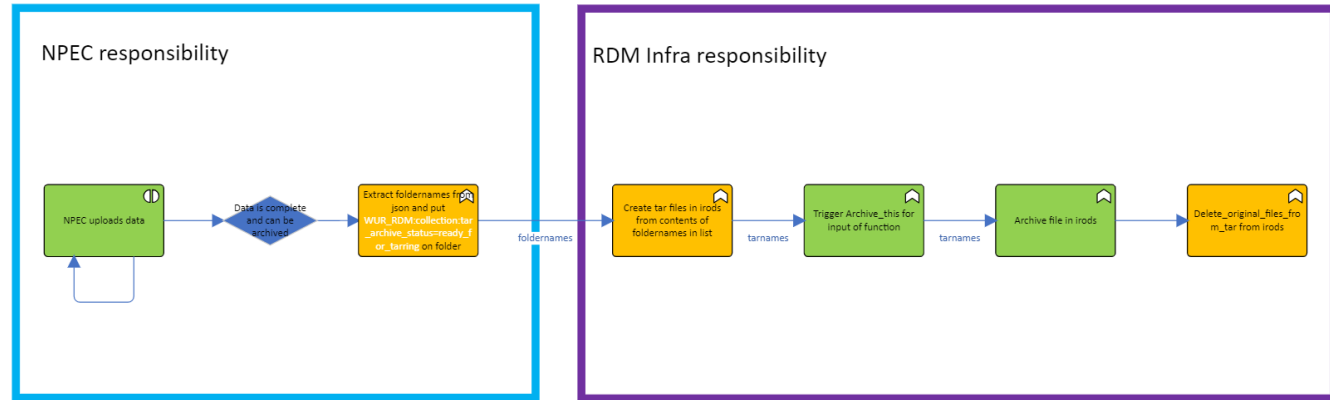
Storage

Sharing



# Archiving data file – ISA-JSON - iRODS

- Cluster data files (TAR) on assay level
- Archive these TAR-files to tape
- Add metadata (original full file path of original collection name) to these TAR-files



# Data flow implementation



Source



Data collection

Data collection



Processing

lustre

iRODS

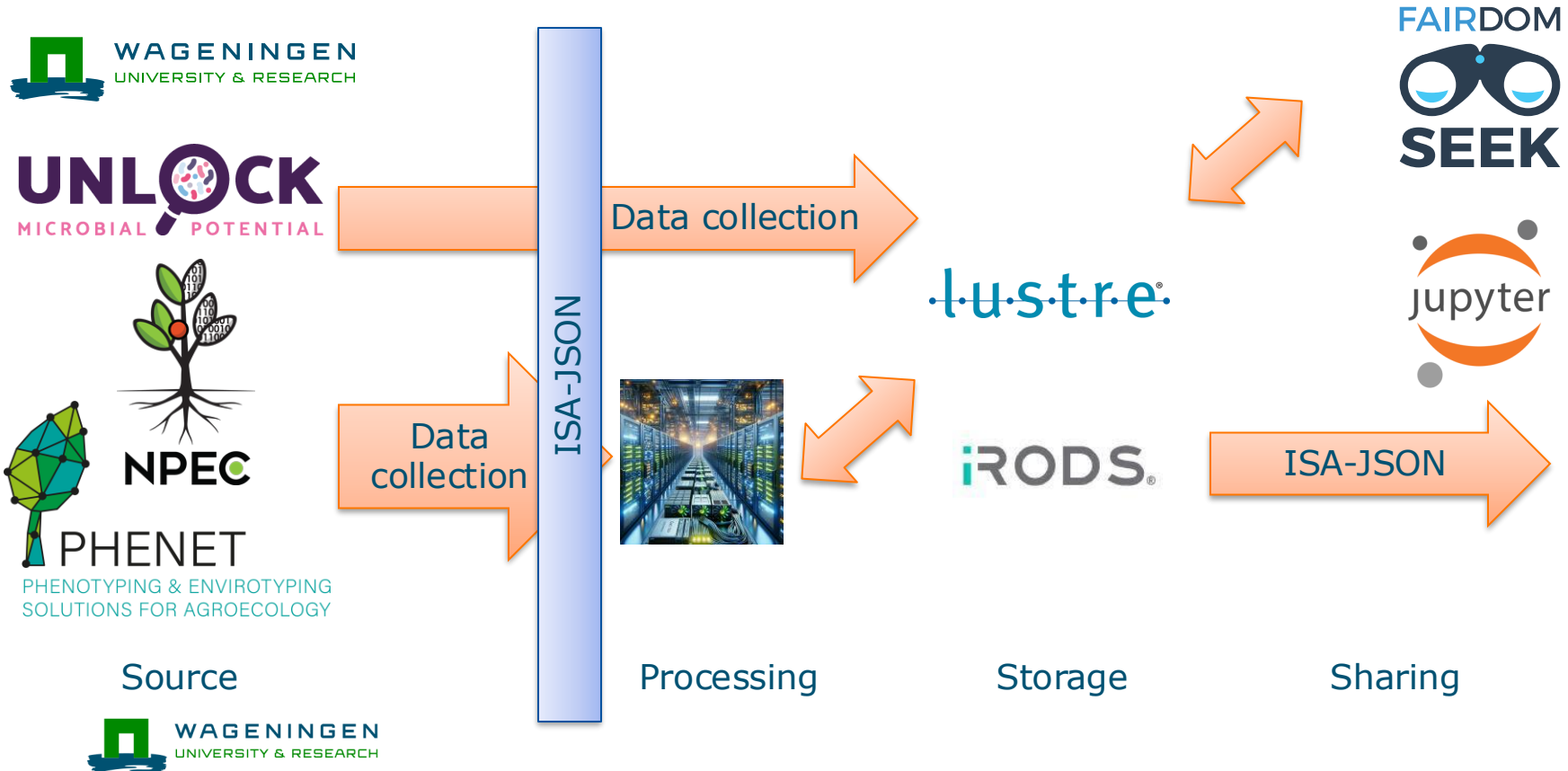
Storage



ISA-JSON

Sharing

# Data flow implementation



# Ongoing implementations

- Data in iRODS
  - URL of file location in iRODS
- extISA structure for many projects
  - Automatic archiving on tape
- Create experiment related ISA-JSON files
- Processing ISA-JSON & data in iRODS using Python notebooks
  - Configurable upload to FAIRDOM-seek
- Upload datasets for data publications incl DOI

# Glossary

- **MIAPPE:** Minimal Information About a Plant Phenotyping Experiment.  
*This is an open, community driven, data standard designed to harmonize data from plant phenotyping experiments. MIAPPE provides a specification including a checklist and a data model of metadata required to adequately describe plant phenotyping experiments. [www.miappe.org](http://www.miappe.org)*
- **ISA-Tab:** The Investigation/Study/Assay (ISA) tab-delimited (TAB) format.  
*This is a known and general purpose framework with which to collect and communicate complex metadata (i.e. sample characteristics, technologies used, type of measurements made) from 'omics-based' experiments employing a combination of technologies. Created by core developers from the University of Oxford, ISA-TAB v1.0 was released in November 2008.*
- **JSON:** JavaScript Object Notation.  
*This is an open standard file format and data interchange format that uses human-readable text to store and transmit data objects consisting of attribute–value pairs. JSON is a language-independent data format. It was derived from JavaScript, but many modern programming languages include code to generate and parse JSON-format data. JSON filenames use the extension .json.*



# More info:

- <https://www.npec.nl/phenotyping-modules/module-7-data/>
- [https://seek4science.org/about\\_us.html](https://seek4science.org/about_us.html)
- <https://www.miappe.org/>

# Question?

Balazs Brankovics

Bart-Jan van Rossum

Rick van de Zedde

Tim van Daalen

Sven Warris

Many others

