

***Datamanagment, OpenData***

*In*

*Horizon 2020*



Björn Usadel

RWTH Aachen University and Forschungszentrum Jülich

Originally a trained biochemist, moved to Bioinformatics early  
coordinator of the German Plant Primary  
Database (>12 year project funding, then Helmholtz)

FW7 DROPS: WP Leader

FW7 EPPN: WP Leader

H2020 Goodberry: WP Leader

H2020 EPPN2020 (provisional): Co-Leader IT	} Co leading Germany +France
ESFRI EMPHASIS : Co -Leader IT	

*Why data management ?*

# Why data management ?

## Olivier Voinnet

From Wikipedia, the free encyclopedia

**Olivier Voinnet** (\*1973)<sup>[1]</sup> is a French biologist and (currently) professor of RNA biology at the [ETH Zurich](#).<sup>[2][3]</sup> Voinnet obtained his PhD in 2001 in England in the group of [David Baulcombe](#) and later obtained a position as independent group leader at the [CNRS](#) in Strasbourg where he was promoted to *Directeur de Recherche* in 2005. In 2010, he moved to [ETH Zurich](#) where he was appointed full professor of RNA Biology.<sup>[1][2]</sup>

### Contents [\[hide\]](#)

- 1 [Manipulation investigated](#)
- 2 [Bans, suspensions](#)
- 3 [Awards](#)
  - 3.1 [Taken back](#)
- 4 [References](#)
- 5 [Links](#)
  - 5.1 [Peer reviews retractions](#)
  - 5.2 [Inquiries](#)
    - 5.2.1 [Articles, discussions](#)
    - 5.2.2 [Press releases ETH, CHRS](#)
  - 5.3 [Early years \(CV early years\)](#)



**“The original lab book data provided to the editors of *Science* showed that these errors did not alter the data in any material way...”**

## Manipulation investigated [\[edit\]](#)

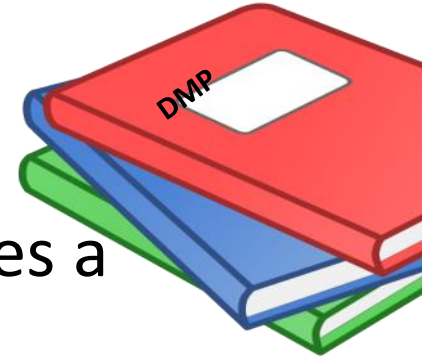
At 2015, his work was investigated for manipulation.<sup>[4]</sup> The investigation at ETH Zurich found that the scientist “breached his duty of care in the handling of figures as well as in his supervisory duties as a research director... and will receive an admonition in relation to his conduct” but also concluded that “this is not a case of scientific misconduct as defined in ETH Zurich’s Rules of Procedure”.<sup>[5]</sup> Another, independent, investigation by CNRS established “the existence of deliberate chart/diagram manipulations, in breach of the ethical standards applicable to the presentation of scientific results” and went on to say that such “inappropriate presentation of experimental data, however, does not amount to fabrication.”<sup>[6]</sup>

# *Why data management ?*

- Make best use of data and share data
- Make sure all partners have access to the data  
*(Usually EU projects feature large consortia and last a long time, one has to make sure that everyone can use and re-use data)*
- Documentation and archival
- Make sure data is not lost, and bring new lab members up to speed
- Good Laboratory Praxis !



# *EU H2020 The data management plan*



- The EU makes things easy, as it [usually] requires a **data management plan (DMP)**

“Horizon 2020 have produced [guidelines !\[\]\(eafc244b53721dd1ec133f0772f70fc7\_img.jpg\) Research Data Management \[PDF 151KB\]](#). All project proposals submitted to 'Research and Innovation actions' as well as 'Innovation actions' should include a section on research data management.”

- If properly written and thought through the plan will actually help you **structure** your thoughts and potentially **make the science** of your proposal **better**

# *EU H2020 The data management plan*

- What you need to do

**“When to write and revise your the Data Management Plan**

The first version of the DMP is expected to be delivered **within the first 6 months** of the project. More elaborated versions of the DMP can be delivered at later stages of the project.

**The DMP should be updated as a minimum in time with the periodic evaluation/assessment of the project. [...] the consortium can define a timetable for review in the DMP itself.**

New versions of the DMP should be created whenever important changes to the project occur due to inclusion of new data sets, changes in consortium policies or external factors.”

- What you probably should do

*Think through at least some points before grant submission!*

# Data Management plan at submission stage

However, **good research data management as such should be addressed under the impact criterion**, as relevant to the project. Your application should address the following issues:

- What standards will be applied?
- How will data be exploited and/or shared/made accessible for verification and reuse?  
If data cannot be made available, why?
- How will data be curated and preserved?
- Your policy should also:
  - reflect the current state of consortium agreements on data management

Participating in the ORD Pilot does **not necessarily mean opening up all your research data**. [...]

**be consistent with exploitation and Intellectual Property Rights (IPR) requirements**

**Use common sense!**



# *EU H2020 Writing the proposal*

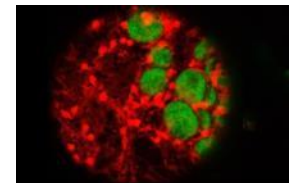
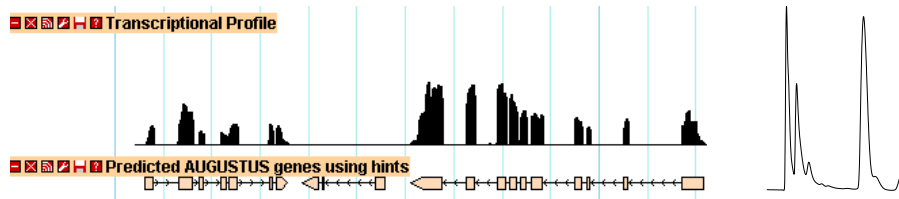
In section “**2.2 - Measures to maximise impact**” one should already mention how data is treated. This covers

- What types of data will the project generate/collect?
- What standards will be used?
- How will this data be exploited and/or shared/made accessible for verification and re-use. If data cannot be made available, explain why.
- How will this data be curated and preserved?

# Section 2.2. Measure to maximise impact

- **What types of data will the project generate/collect?**

Possible items can be e.g. qRT PCR results, microarray studies, chromotographics traces, marker enzymes, next generation sequencing data\*, health profiles\*, pictures and images....  
But also textual data such as patient interviews



\* Be careful about the ethical implications

# Section 2.2. Measure to maximise impact

- **What standards will be used?**

For many molecular biology derived datasets standards exist! For example you might have heard about the whole set of MIAXXXX (Minimal information about ....) Many standards can be found on biosharing.org

Search Standards

3 records in view

Registry	Name	Abbreviation	Type	Domain	Taxonomy	Related Database	Related Standard	Related Policy	In Collection/Recommendation	Status
<input checked="" type="checkbox"/>	Minimum Information for Publication of Quantitative Real-Time PCR Experiments	MIQE	Standard	<a href="#">Assay</a> <a href="#">DNA</a> <a href="#">Life Science</a> <a href="#">RNA</a> <a href="#">Real Time Polymerase Chain Reaction Assay (RT-PCR)</a> <a href="#">Plus 1 More...</a>	<input checked="" type="checkbox"/> All	None	RDML	Oxford University Press - Journal Of The National Cancer Institute - Manuscript Preparation	Minimum Information for Biological and Biomedical Investigations FAIRDOM Community Standards	

<https://biosharing.org/standards>

# Section 2.2. Measure to maximise impact

- What standards will be used?

Make sure the **standard is suitable** for your research! Say “standards like MIQE” “MIAPPE 1.0 or later” especially if the field is still evolving

HOW STANDARDS PROLIFERATE:  
(SEE: A/C CHARGERS, CHARACTER ENCODINGS, INSTANT MESSAGING, ETC)



# Section 2.2. Measure to maximise impact

- How will this data be exploited and/or shared/made accessible for verification and re-use. If data cannot be made available, explain why.

Likely you will share the data (→ OpenData), think of ways how you can make the data available to others. The best is **to make use of public, established** resources. Check [openaire.eu](https://www.openaire.eu)

The screenshot displays the OpenAIRE website interface. At the top left is the OpenAIRE logo. The navigation bar includes links for PARTICIPATE, SEARCH, MONITOR, SUPPORT, and OPEN ACCESS. Below the navigation bar is the heading "OR LOCATE DATA PROVIDER IN MAP".

On the left side, there are several filter categories with checkboxes:

- Publication Repository Aggregator
- Thematic Repository
- Publication Repository
- Journal Aggregator/Publisher
- Journal
- Registry
- Institutional Repository Aggregator

Next to these are filter options for OpenAIRE funding and compatibility:

- OpenAIRE 2.0 EC funding
- OpenAIRE Basic DRIVER OA
- collected from a compatible aggregator
- OpenAIRE 2.0 EC funding
- proprietary

On the right, there is a dropdown menu for country selection with options: Algeria, Argentina, Australia, Austria, and Belarus.

The main area shows a map of Europe with several blue location pins. A tooltip for "Juelich Shared Electronic Resources" is visible, providing details: type: Institutional Repository, compatibility: OpenAIRE 3.0 OA funding, organization: Forschungszentrums Jülich, country: Germany.

On the right side of the map, there is a promotional banner for Zenodo, stating "THE PLACE TO SHARE YOUR RESEARCH RESULTS" and "OpenAIRE's catch-all repository hosted by CERN". Below this, it says "OPENAIRE USES OpenDOAR" and "re3data.org" (with a note "LIST OF REPOSITORIES / PORTALS"). At the bottom right, it says "as authoritative repository registries".

The URL <https://www.openaire.eu> is displayed at the bottom right of the screenshot.

## Section 2.2. Measure to maximise impact

- **How will this data be exploited and/or shared/made accessible for verification and re-use. If data cannot be made available, explain why.**

Nature Scientific Data has a brief overview of some important repositories

urages authors to archive data to one of the above data-type specific repositories. If a data-type specific repository is not available, we recommend the following generalist repository for the storage of data. Generalist repositories may also be appropriate for archiving associated data, supplementing the primary data in a data-type specific repository.

Information on fees/costs	Size limits	Integrated with Scientific Data manuscript submission system
20 USD for first 20 GB, and \$50 USD for each additional 10 GB	None stated	Yes ✓
100 GB free per Scientific Data manuscript. Additional fees apply for larger datasets	1 TB per dataset	Yes ✓ - To qualify for the 100 GB free storage, data must be uploaded to figshare via our submission system. Download instructions.
Contact repository for datasets over 1 TB	2.5 GB per file, 10 GB per dataset	No

# Section 2.2. Measure to maximise impact

- How will this data be exploited and/or shared/made accessible for verification and re-use. If data cannot be made available, explain why.

Re3data allows you search and find repositories

The screenshot displays the Re3data search interface. On the left is a 'Filter' sidebar with various categories like Subjects, Content Types, Countries, etc. The main search area shows a search bar with 'pcr' entered, a search button, and a 'Toggle short help' link. Below the search bar are navigation controls: '← Previous', '1', and 'Next →'. The results section shows 'Found 4 result(s)'. The first result is 'database of Sequence Tagged Sites' (dbSTS) from the United States. The second result is 'ChIP-Seq Transcription Factor Data' from Canada. Each result card includes a title, a subject line, content type(s), country, and a brief description.

www.re3data.org/repository/3d100010649

# Zenodo a generalist repository

[Upload](#)[Communities](#)[Log in](#)[Sign up](#)

February 27, 2017

Dataset

Open Access

## Dataset for "Soil fluxes of carbonyl sulfide (COS), carbon monoxide, and carbon dioxide in a boreal forest in southern Finland"

Sun, W.; Kooijmans, L. M. J.; Maseyk, K.; Chen, H.; Mammarella, I.; Vesala, T.; Levula, J.; Keskinen, H.; Seibt, U.

This is the dataset (ver. 2017.02.13) for the manuscript "Soil fluxes of carbonyl sulfide (COS), carbon monoxide, and carbon dioxide in a boreal forest in southern Finland" submitted to the journal *Atmospheric Chemistry and Physics*.

Preview

hyy15\_chflux\_20170213.zip

hyy15_blank.csv	1.3 kB
hyy15_chflux_release.csv	1.3 MB
hyy15_moss.csv	2.9 kB
readme.md	4.4 kB
readme.pdf	462.1 kB

### Publication date:

February 27, 2017

### DOI:

DOI [10.5281/zenodo.322936](https://doi.org/10.5281/zenodo.322936)

### Keyword(s):

carbonyl sulfide

carbon monoxide

soil-atmosphere gas exchange

boreal forest

### Grants:

[European Commission:](#)

- INGOS - Integrated non-CO2 Greenhouse gas Observing System (284274)

### Related identifiers:

Cited by:

[10.5194/acp-2017-180](#)

### Communities:

[European Commission Funded Research \(OpenAIRE\)](#)

[Zenodo](#)

### License (for files):

[Creative Commons Attribution 4.0](#)

### Share

### Cite as

Sun, W., Kooijmans, L. M. J., Maseyk, K., Chen, H., Mammarella, I., Vesala, T., ... Seibt, U. (2017).



## *Section 2.2. Measure to maximise impact*

- **How will this data be curated and preserved?**

If you can use one of the previously mentioned providers → easy use them.

- Is my data safe with you / What will happen to my uploads in the unlikely event that Zenodo has to close?

Yes, your data is stored in CERN Data Center. Both data files and metadata are kept in multiple online replicas and independent replicas. CERN has considerable knowledge and experience in building and operating large scale digital repositories and a commitment to maintain this data centre to collect and store 100s of PBs of LHC data as it grows over the next 20 years. In the highly unlikely event that Zenodo will have to close operations, we guarantee that we will migrate all content to other suitable repositories, and since all uploads have DOIs, all citations and links to Zenodo resources (such as your data) will not be affected.

# EU H2020 Writing the proposal

Add (milestones and) deliverables for the DMP! By default you have to have a DMP.

If you have a workpackage about data management already, it makes a lot of sense to add it there and like this you can reiterate some points in the project summary

## Task 5.4 Data Management Plan

To safeguard open access and sustainable access to GoodBerry data, a data management plan will be elaborated. This will include SOPs to measure traits as detailed in D1.1 and D4.5 but will also deal with ontologies, general best practices, minimal requirements such as MIAMET, MINSEQe etc. data storage, backups and data accessibility (see also Task 5.1)(**D5.1** and **D5.5**).

*Partners involved: RWTH AACHEN*

*Duration: month 1 – month 48*

## Deliverables

N°	Brief description	Month of delivery
D5.1	Data Management Plan first version R1.0	6

# *EU H2020 Writing the proposal*

- Check your environment! Maybe your institution /library can help you! As they are already involved in data management schemes.
- Make sure you also **allocate necessary resources!**



# *Summary: Writing the proposal*

- Think about data **requirements**, plan what kind of data is going to be important and how much you will have
- Familiarize yourself with **standards** (probably you know them already anyway)
- See if there are already **data repositories**/standards tackling exactly what you need

# *Writing the DMP*

**H2020 Programme**

Guidelines on

FAIR Data Management in Horizon 2020

---

[http://ec.europa.eu/research/participants/data/ref/h2020/grants\\_manual/hi/oa\\_pilot/h2020-hi-oa-data-mgt\\_en.pdf](http://ec.europa.eu/research/participants/data/ref/h2020/grants_manual/hi/oa_pilot/h2020-hi-oa-data-mgt_en.pdf)

This also includes a „template“



# *DMP EU Template*

## 1. Data Summary

What is the **purpose** of the data collection/generation and its relation to the objectives of the project?

**What types and formats** of data will the project generate/collect?

Will you **re-use any existing data** and how?

What is the **origin** of the data?

What is the expected **size** of the data?

To whom might it be **useful** ('data utility')?

Most questions you can answer from the grant proposal.

For some sets to whom data might be useful: it might be just the consortium. If you explain this properly this is a just answer.

But be creative there is often many more uses



# *DMP EU Template*

## 2. The FAIR principle

### **The FAIR Guiding Principles**

- **To be Findable**
- **To be Accessible**
- **To be Interoperable**
- **To be Reusable**

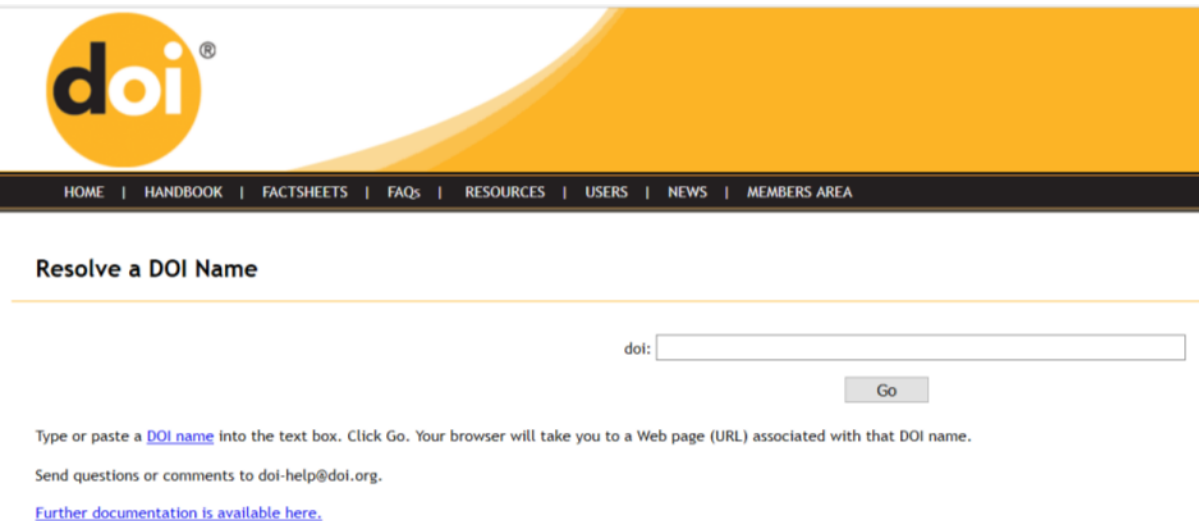
# To be Findable

F1. (meta)data are assigned a globally unique and eternally persistent identifier.

F2. data are described with rich metadata.

F3. (meta)data are registered or indexed in a searchable resource.

F4. metadata specify the data identifier.



The screenshot shows the top navigation bar of the DOI.org website with the DOI logo and links for HOME, HANDBOOK, FACTSHEETS, FAQs, RESOURCES, USERS, NEWS, and MEMBERS AREA. Below the navigation bar is a section titled "Resolve a DOI Name" which contains a text input field labeled "doi:" and a "Go" button. Below the input field, there is a small instruction: "Type or paste a DOI name into the text box. Click Go. Your browser will take you to a Web page (URL) associated with that DOI name." At the bottom of the section, there is a link to "Send questions or comments to doi-help@doi.org" and another link: "Further documentation is available here."

What is this?  
What was done?  
Who did it ?  
How was it done?  
Who submitted?

<https://guidelines.openaire.eu/>

<https://www.force11.org/group/fairgroup/fairprinciples>



# To be Accessible

A1 (meta)data are retrievable by their identifier using a standardized communications protocol.

A1.1 the protocol is open, free, and universally implementable.

A1.2 the protocol allows for an authentication and authorization procedure, where necessary.

A2 metadata are accessible, even when the data are no longer available.

# To be Interoperable

11. (meta)data use a formal, accessible, shared, and broadly applicable language for knowledge representation.
12. (meta)data use vocabularies that follow FAIR principles.
13. (meta)data include qualified references to other (meta)data.

“Metadata being **machine readable** is a *conditio sine qua non* for FAIRness.”

Use e.g. ontologies or controlled vocabularies

# To be Reusable

R1. meta(data) have a plurality of accurate and relevant attributes.

R1.1. (meta)data are released with a clear and accessible data usage license.

R1.2. (meta)data are associated with their provenance.

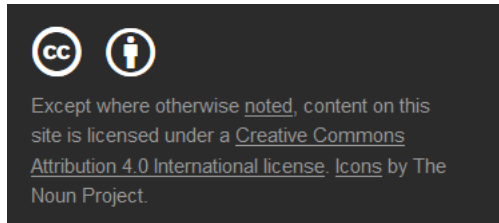
R1.3. (meta)data meet domain-relevant community standards.

[Try out our simple License Chooser.](#)

## Three “Layers” Of Licenses



From  
<https://creativecommons.org/licenses/>



# Selecting a license

Figure 1: Locating the license selector in the B2SHARE service.

## Usage

During usage, the tool will ask the user some questions which will narrow down the possibilities of license choice.

### Choose a License

Answer the questions or use the search to find the license you want


What do you want to deposit?

Search for a license...

---

**Public Domain Mark (PD)**



The work identified as being free of known restrictions under copyright law, including all related and neighboring rights.



---

**Public Domain Dedication (CC Zero)**



CC Zero enables scientists, educators, artists and other creators and owners of copyright- or database-protected content to waive those interests in their works and thereby place them as completely as possible in the public domain, so that others may freely build upon, enhance and reuse the works for any purposes without restriction under copyright or database law.

---

**Creative Commons Attribution (CC-BY)**

This is the standard creative commons license that gives others maximum freedom to do what they want with your work.

---

**Creative Commons Attribution-ShareAlike (CC-BY-SA)**

# FAIR: sounds unfairly difficult and technical?

Many repositories are FAIR already (or come close)  
Don't reinvent the wheel!

### FAIRDOMHub

The **FAIRDOMHub** is built upon the **SEEK** software suite, which is an open source web platform for sharing scientific research assets, processes and outcomes. For more information about SEEK please visit <http://seek4science.org>

**FAIRDOM** is an initiative to develop a community, and establish an internationally sustained Data and Model Management service to the European Systems Biology community. FAIRDOM is a joint action of ERA-Net EraSysAPP and European Research Infrastructure ISBE.  
For more information about FAIRDOM please visit <http://fair-dom.org>

If you are interested in using **FAIRDOMHub** within your own funding programme, or have any other questions related to FAIRDOM and SEEK, including feature requests or how to get involved, please contact us at [support@fair-dom.org](mailto:support@fair-dom.org).  
**Now available: self-management of your own programmes and projects**

#### News

[Progesterone signalling in broiler skeletal muscle is associated with divergent feed efficiency](#)  
Most Recent Articles: BMC Systems Biology - 4 ago

[Network topology of Nav1.7 mutations sodium channel-related painful disorder](#)  
Most Recent Articles: BMC Systems Biology - 4 ago

[Dynamic genome-scale metabolic model of the yeast \*Pichia pastoris\*](#)  
Most Recent Articles: BMC Systems Biology - 7 ago

#### Tags [show all]

[Bacillus subtilis](#) [Biochemistry](#) [Biochemist](#) [protein analysis](#) [Bioinformatics](#) [Computational](#) [theoretical biology](#) [Computational Systems I](#) [Data Management](#) [Databases](#) [Fermi](#) [Genetics](#) [Image analysis](#) [Mathematical](#) [matlab](#) [Metabolomics](#) [Microarray](#) [Microbiology](#) [Molecular](#) [Molecular biology techniques](#) [RNA/DNA](#) [parameter estimation](#) [Proteomics](#) [R](#) [SI](#) [Systems Biology](#) [Transcriptome](#)

Get started

Who is using it?

Get involved

Latest additions

Latest downloads

# An example plant phenotyping: Discussions started a long time ago



## Standardisation of Plant Phenotyping Experiment Description

[DOWNLOADS](#)[ISA-TAB FOR PHENOTYPING](#)[MIAPPE](#)[STANDARDISATION GROUP](#)

working group

### PAGES

[Downloads](#)[ISA-Tab for Phenotyping](#)[MIAPPE](#)[Standardisation group](#)

### RECENT POSTS

A new paper published  
November 24, 2016

Configuration change  
April 1, 2016

Configuration change  
September 1, 2015

Working on MIAPPE  
August 24, 2015

Opinion paper published  
June 15, 2015

### RECENT COMMENTS

#### MINIMUM INFORMATION

## OPINION PAPER PUBLISHED

🕒 JUNE 15, 2015   👤 HCWI

Following discussions of transPLANT partners and EPPN representatives, an opinion paper on the need of standardisation in the field of plant phenotyping has appeared:

Pawel Krajewski\*, Dijun Chen, Hanna Cwiek, Aalt D.J. van Dijk, Fabio Fiorani, Paul Kersey, Christian Klukas, Matthias Lange, Augustyn Markiewicz, Jan Peter Nap, Jan van Oeveren, Cyril Pommier, Uwe Scholz, Marco van Schriek, Bjoern Usadel, and Stephan Weise

### **Towards recommendations for metadata and data handling in plant phenotyping**


*Journal of Experimental Botany*, 2015. [doi:10.1093/jxb/erv271](https://doi.org/10.1093/jxb/erv271)

[◀ GOOD PRACTICE](#)[◀ MINIMUM INFORMATION](#)[◀ STANDARDISATION](#)

---

[PREVIOUS POST](#)

# Describing the file format is necessary



Standardisation of Plant Phenotyping Experiment Description

working group

Search ...

PAGES

- Downloads
- ISA-Tab for Phenotyping
- MIAPPE
- Standardisation group

RECENT POSTS

- A new paper published  
November 24, 2016
- Configuration change  
April 1, 2016
- Configuration change  
September 1, 2015
- Working on MIAPPE  
August 24, 2015
- Opinion paper published  
June 15, 2015

RECENT COMMENTS

DOWNLOADS   ISA-TAB FOR PHENOTYPING   MIAPPE   STANDARDISATION GROUP

FORMAT, MINIMUM INFORMATION

## A NEW PAPER PUBLISHED

NOVEMBER 24, 2016   HCWI

Following the opinions expressed in [Krajewski et al.](#), we are presenting a new paper summarising the solutions proposed for the improvement of phenotypic data description:

Hanna Ćwiek-Kupczyńska, Thomas Altmann, Daniel Arend, Elizabeth Arnaud, Dijun Chen, Guillaume Cornut, Fabio Fiorani, Wojciech Frohberg, Astrid Junker, Christian Klukas, Matthias Lange, Cezary Mazurek, Anahita Nafissi, Pascal Neveu, Jan van Oeveren, Cyril Pommier, Hendrik Poorter, Philippe Rocca-Serra, Susanna-Assunta Sansone, Uwe Scholz, Marco van Schriek, Ümit Seren, Björn Usadel, Stephan Weise, Paul Kersey and Paweł Krajewski

**Measures for interoperability of phenotypic data: minimum information requirements and formatting**  
*Plant Methods*, 2016. DOI [10.1186/s13007-016-0144-4](https://doi.org/10.1186/s13007-016-0144-4)

BEST PRACTICES   STANDARDISATION

Next steps new workshops, (re) defining ontologies and needs

### Minimum Information About a Plant Phenotyping Experiment (MIAPPE)

Attributes (concepts, subconcepts - in terms of ontology) marked by asterisk (\*) are essential for a description of experiment (e.g. by Poorter et al. [26]); the rest forms an extended description. For some attributes examples of possible values are listed.

Checklist section	Attributes	Source list / Biosharing ID / Reference	Recommended ontologies
General metadata	Unique identifier* Title* Description* Submission date Public release date Publications Laboratory address and contact details	Default ISA-Tab configuration [1]	OBI, Ontology for Biomedical Investigations [2]  CRO, Crop Research Ontology [3]
Timing and location	Timing: Start of experiment (date)* Duration (days/months/years)*  Experiment location: Geographic location* Latitude and longitude Altitude Inclination and aspect Habitat	Poorter et al. [4]  Morrison et al. [5]  CIMR [6]: Environmental Analysis Context [7]	OBI, Ontology for Biomedical Investigations [2]  GAZ, Gazetteer [9]
Biosource	Organism (taxon)* Intraspecific_name* Intraspecific_rank Common name	MiXs Plant-associated environmental package [10]  Yilmaz et al. [11]	UNIPROT Taxonomy [13]  NCBI Taxonomy [14]



INTRODUCTION

Welcome

Download & Source Code

APPLICATIONS

Running Repositories

Demos and Samples

SCIENTIFIC PUBLICATIONS

Paper - Talks - Poster

DOCUMENTATION

Development

Run

ABOUT

License and

## e!DAL-MetaData-API - store, cite and share primary data

e!DAL is a lightweight software framework for publishing and sharing research data. Its main features are version tracking, metadata management, information retrieval, journal and founding agency proven registration of persistent identifiers (DOI), an embedded HTTP(S) server for public data access, access as a network file system, and a scalable storage backend.

In any research that make use of the e!DAL components **please honour the author's work** and cite:

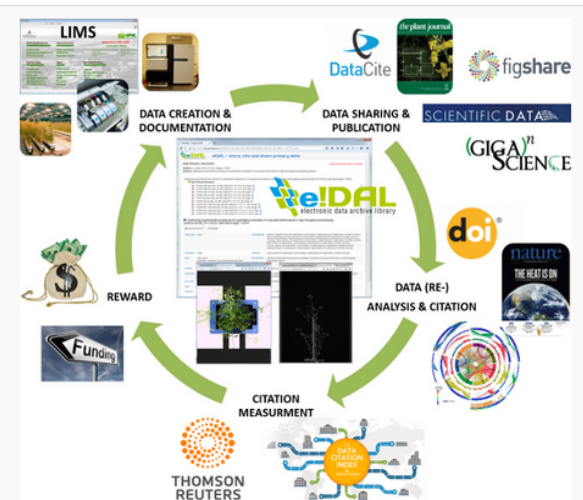
[Arend et al. e!DAL - a framework to store, share and publish research data](#)

### Public Data Repositories

e!DAL driven repositories publish already a high number of scientific citable and Data Cite registered research data. Registered users may submit data sets to the related repository.

### Personal Research Data Repository

Packaged as e!DAL server, all required API components are compiled as executable JAR archive. This can be executed at any platform to operate an own data publication infrastructure. Please follow the instructions to [set-up and execute an e!DAL server](#). Java projects may access a e!DAL Server using [remote client-API](#).



# You might still use your own database analysis tools to integrate data

Protein search Gene browser Plant

Genomes - overview of published *Angiospermae* genomes by taxonomic classification  
*ssica napus* (Rapeseed) data

protein name

Enter keyword or protein/gene name

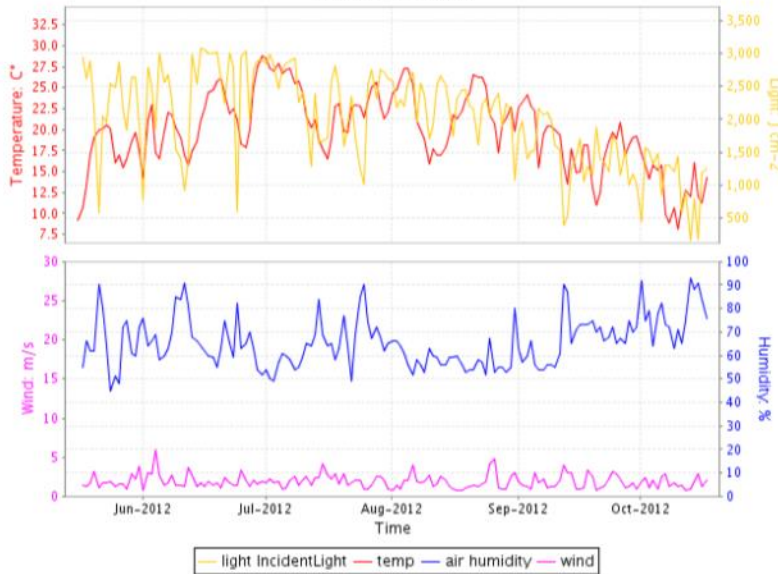
vine  
 seed  
 (Potato)  
 beet

enter keyword or name

reset search

	initial set							optional
sensor type	air temperature	air humidity	rain	radiation	wind speed	rate	optional	
sensor number	1	1	1	1	1	2		
sensor locations	200 m (sr) - 600 m	200 m (sr) - 600 m	on field - 200 m (sr)	on field - 200 m (sr)	field - 200 m (sr) - 600 m (drou)	** 30 km distance to the field		
recording frequency	60 minutes	60 minutes	60 minutes	60 minutes	60 minutes			
				Average method		rain gauge diameter:		
remarks							Use weather sensor status. Make program for logal data. Make program for the rate sensor of air sensor status. Make program for the rate sensor of air sensor status. Make program for the rate sensor of air sensor status. Make program for the rate sensor of air sensor status.	
TIME:								
dd/mm/yyyy hh:mm	continuous	alt continuous		global radiation	m/s			
	%	%	mm	W/m <sup>2</sup>	m/s	n/a	mm	
11 14/06/2013 17:00	19.4	77		140	1.1			
12 14/06/2013 18:00	19.76	75		82	1.1			
13 14/06/2013 19:00	18.89	81		20	0.7			
14 14/06/2013 20:00	17.88	88		0	0.3			
15 14/06/2013 21:00	17.54	92		0	0.6			
16 14/06/2013 22:00	17.05	96		0	0.2			
17 14/06/2013 23:00	16.72	98		0	0			
18 14/06/2013 00:00	16.37	99		0	0			
19 14/06/2013 01:00	16.33	99		0	0			
20 14/06/2013 02:00	16.15	99		0	0			
21 14/06/2013 03:00	16.02	97		0	0			
					0.2			
					1.4			
					1.6			
					1.3			
					1.3			
					21.2	1.3		
					27.0	1		
					40.9	1.7		
					50.6	1.5		
					57.1	1.2		
					57.0	1.9		
					64.7	1.7		
					64.0	1.7		
					100	1.5		
					105	1.9		
					99	0.8		
					25	0.2		
					0			

Weather Data

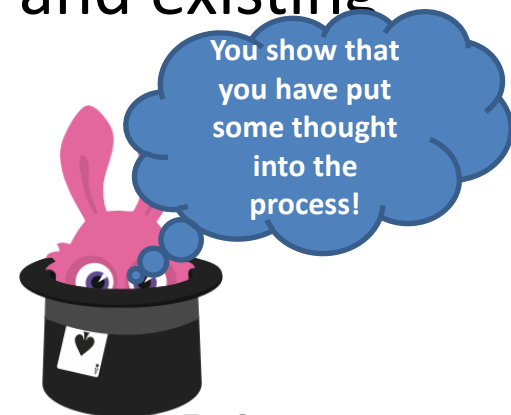


Plabipd already covers the genomic base, within PROGRESS additional data



# *Summary DMP: FAIR*

- Making data FAIR without any knowledge in the field can be difficult (but it is not impossible!)
- If you can, use acknowledged repositories and existing standards.
- Potentially liaise/explore with repositories



(Maybe EGI, EUDAT or the Research data alliance RDA can help)



# *DMP template*

## **3. Allocation of resources**

- What are the costs for making data FAIR in your project?
- How will these be covered? Note that costs related to open access to research data are eligible as part of the Horizon 2020 grant (if compliant with the Grant Agreement conditions).
- Who will be responsible for data management in your project?
- Are the resources for long term preservation discussed (costs and potential value, who decides and how what data will be kept and for how long)?

**Be aware you can claim costs, but data should be preserved long term also when the project is over! Make sure to be able to keep websites and domain names as well. Do this in your own interest. (In the future somebody might search your EU project website and find .... something unexpected)**



# \$\$\$ Who pays?

Plant genomes / Solanum pennellii (new cultivar)



WEB CONTENT DISPLAY

## Sequencing the gigabase plant genome of the wild tomato species *Solanum pennellii* using Oxford Nanopore single molecule sequencing

### Contributors

[Maximilian Schmid](#)<sup>1</sup>, [Alexander Vogel](#)<sup>1</sup>, Alexandra Wormit<sup>1</sup>, Alisandra Denton<sup>1</sup>, Anthony Bolger<sup>1</sup>, Henri van de Geest<sup>2</sup>, Benjamin Istace<sup>6</sup>, Marie E. Bolger<sup>3</sup>, Saleh Alseekh<sup>4</sup>, Janina Maß<sup>3</sup>, Christian Pfaff<sup>3</sup>, Ulrich Schurr<sup>3</sup>, Jean-Marc Aury<sup>6</sup>, Alisdair R. Fernie<sup>4</sup>, Dani Zamir<sup>5</sup>, [Björn Usadel](#)<sup>1,3</sup>.

<sup>1</sup>Institute for Botany and Molecular Genetics, BioEconomy Science Center, RWTH Aachen University, Aachen, Germany.

<sup>2</sup>Wageningen Plant Research, Droevendaalsesteeg 1, 6708 PB, Wageningen, The Netherlands

<sup>3</sup>Institute for Bio- and Geosciences (IBG-2: Plant Sciences), Forschungszentrum Jülich, Jülich, Germany.

<sup>4</sup>Department of Molecular Physiology, Max Planck Institute of Molecular Plant Physiology, Potsdam-Golm, Germany.

<sup>5</sup>Faculty of Agriculture, Hebrew University of Jerusalem, Rehovot, Israel.

<sup>6</sup>Genoscope (CEA) and UMR 8030 CNRS-Genoscope-Université d'Evry, 2 rue Gaston Crémieux, BP5706, 91057 Evry, France.

### Background

One dataset 6 weeks of lab work... more than 50TB of data.

By current measures and standards

**ALL THIS IS RAW DATA THAT NEEDS TO BE PRESERVED**

**In this case it is 3rd generation sequencing data, so it should be no problem**



# *DMP template*

## **4. Data security**

- What provisions are in place for data security (including data recovery as well as secure storage and transfer of sensitive data)?
  
- Is the data safely stored in certified repositories for long term preservation and curation?

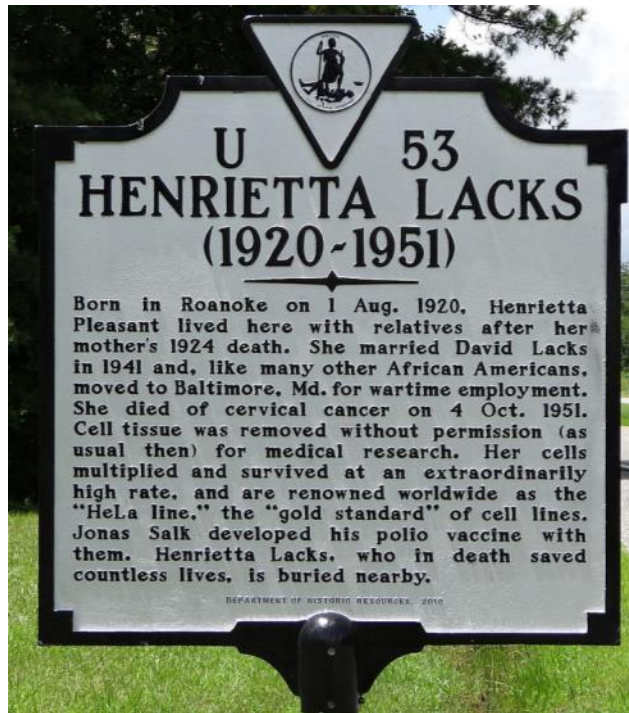




# DMP template

## 5. Ethical aspects

- Are there any ethical or legal issues that can have an impact on data sharing? These can also be discussed in the context of the ethics review. If relevant, include references to ethics deliverables and ethics chapter in the Description of the Action (DoA).
- Is informed consent for data sharing and long term preservation included in questionnaires dealing with personal data?



**Large scale DNA is never anonymous!**

HeLa cell line data published  
The family was not amused!



According to an interview with technologyreview  
Decode Genetics can infer almost all Icelander's  
Genomic make up . But can not warn due to Ethics rules



# DMP template

## 6. Other issues

Do you make use of other national/funder/sectorial/departmental procedures for data management? If yes, which ones?

The screenshot shows the website's navigation menu with 'Research Data Management' selected. The breadcrumb trail reads: Home > Research > Research Data Management > Research Data Management at RWTH Aachen. The left sidebar contains a menu with 'Research Data Management at RWTH Aachen' highlighted, followed by 'Leitlinie zum Forschungsdatenmanagement an der RWTH Aachen', and several question-based links: 'What Are Research Data?', 'Why Manage Data?', 'What Does My Funder or Publisher Require?', 'How Do I Get Help?', 'Contact', and 'Governance'. The main content area features the title 'Research Data Management at RWTH Aachen' and a paragraph explaining that digital data is both a product and an object of research, with researchers responsible for its protection and secure processing. Below this, it states that RWTH Aachen supports its researchers with standards, training, and infrastructure to ensure efficient and legally informed research work. A 'CONTACT' box on the right includes a photo of staff, the text 'ServiceDesk Research Data Management', and contact details: a phone icon with '+49 241 80 24680' and an email icon with 'Send Email'. A blue thought bubble on the right contains the text: 'This can actually be helpful and they might have templates services for you'.



This can actually be helpful and they might have templates services for you



# *DMP Examples and Resources*

## In this section

[Briefing Papers](#)
[How-to Guides & Checklists](#)
[Developing RDM Services](#)
[Curation Lifecycle Model](#)
[Curation Reference Manual](#)
[Policy and legal](#)

### Data Management Plans

[Checklist](#)
[DMPonline](#)
[FAQ on DMPonline](#)
[FAQ on Data Management Plans](#)
[Funders' requirements](#)
[Guidance and examples](#)
[Tools](#)
[Case studies](#)
[Repository audit and assessment](#)
[Standards](#)
[Publications and presentations](#)

## Checklist for a Data Management Plan

The DCC synthesises requirements for Data Management Plans and best practice within the wider community. This allows us to provide a Checklist that presents the main questions or themes that researchers may want to cover when writing a DMP.

In 2013 the DCC reviewed and shortened its Checklist. The current version is available to download below.

### [Checklist for a Data Management Plan \(v.4.0, 2014\)](#)

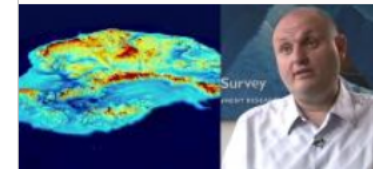
Also see the [DMP Checklist flyer](#), a handy foldout version of the Checklist. Hard copies are available if you would like some for events.

**\*\* This publication is available in print and can be ordered from our [online store](#) \*\***

### Earlier versions

The DCC first created a Content Checklist for a Data Management Plan in 2009. This was put out for consultation and over the years it developed into a comprehensive list of questions that researchers may wish to cover.

### Managing Research Data (video)



This short documentary, Digital Curation Centre: Managing Research Data, offers a unique insight into the importance of providing access to research data and the risks of not managing data effectively.

[Watch this video](#)

## In this section

Briefing Papers

How-to Guides & Checklists

Developing RDM Services

Curation Lifecycle Model

Curation Reference Manual

Policy and legal

### Data Management Plans

Checklist

**DMPonline**

FAQ on DMPonline

FAQ on Data Management Plans

Funders' requirements

Guidance and examples

Tools

Case studies

Repository audit and assessment

Standards

Publications and presentations

Roles

Curation journals

Informatics research

External resources

Online Store

## DMPonline

Research funders and organisations increasingly require data management plans, both during the bid-preparation stage and after funding has been secured.

**DMPonline** is the DCC's data management planning tool. It provides tailored guidance and examples to help researchers write data management plans.

The tool includes a number of templates for funders in the UK and overseas so researchers can write DMPs according to the specific requirements they need to meet. It can also be customised by institutions so they can add their own templates and guidance.

A [screencast](#) provides an overview of how the tool works.

Try the tool for yourself at <http://dmponline.dcc.ac.uk>

Anyone can use DMPonline. If your organisation is not listed, just select 'other organisation' or ask for it to be added.

If you would like to create a foreign language version of DMPonline, please contact us on [dmponline@dcc.ac.uk](mailto:dmponline@dcc.ac.uk)



To keep up with DMPonline news, [subscribe to the RSS feed](#) or [watch GitHub](#) for updates.

All our code is open and available for you to redistribute and/or modify it under the terms of the GNU Affero General Public License.

## Useful links

DMPonline

### IDCC



The 12th International Digital Curation Conference (IDCC) will take place at The Royal College of Surgeons of Edinburgh, UK

20 - 23 February 2017.

[Read more](#)

## Meaningful file names

Below are tips on meaningful and consistent file names. Read more in ['Naming files and folders'](#).<sup>(2)</sup>

- ❑ Make sure to use consistent file names. When you use a date in the file name, choose a notation (for instance, YYYYMMDD or yymmdd).
- ❑ Do not use strange characters like ?\!@\*%{[<> in the file name.
- ❑ Use traceable file names, such as Project\_Instrument\_locatie\_YYYYMMDD.ext.
- ❑ Make sure to only use each file once in the folder structure. If you store a file in more than one place, several versions of the same file can unwillingly be created.
- ❑ See also [version management](#).

It is good practice to note the file naming and its meaning in a readme.txt.

Even if a researcher is well underway with his project consistent file naming is still an option by using a [bulk file rename utility](#).<sup>(3)</sup> It is important, however, to check if this bulk renamer delivers on its promises.



white\_data\_20140708.csv



blue\_data\_20140708.docx



red\_data\_20140708.R



red\_data\_20140708\_v02.R

*File naming and version management*

# *The Open Data Pilot*

RESEARCH ARTICLE

Open Access

# A systematic review of barriers to data sharing in public health

Willem G van Panhuis<sup>1\*</sup>, Prama Paul<sup>1</sup>, Claudia Emerson<sup>2</sup>, John Grefenstette<sup>1</sup>, Richard Wilder<sup>3</sup>, Abraham J Herbst<sup>4,5</sup>, David Heymann<sup>6</sup> and Donald S Burke<sup>1</sup>

## Abstract

**Background:** In the current information age, the use of data has become essential for decision making in public health at the local, national, and global level. Despite a global commitment to the use and sharing of public health data, this can be challenging in reality. No systematic framework or global operational guidelines have been created for data sharing in public health. Barriers at different levels have limited data sharing but have only been anecdotally discussed or in the context of specific case studies. Incomplete systematic evidence on the scope and variety of these barriers has limited opportunities to maximize the value and use of public health data for science and policy.

**Methods:** We conducted a systematic literature review of potential barriers to public health data sharing. Documents that described barriers to sharing of routinely collected public health data were eligible for inclusion and reviewed independently by a team of experts. We grouped identified barriers in a tax international dialogue on solutions.

**Results:** Twenty potential barriers were identified and classified in six categories: technical, political, legal and ethical. The first three categories are deeply rooted in well-known challenge systems for which structural solutions have yet to be found; the last three have solutions that dialogue aimed at generating consensus on policies and instruments for data sharing.

**Table 1 Evidence for barriers to sharing of routinely collected public health data**

Category	Barrier	Peer-reviewed		Non peer-reviewed
		Empirical data	Non-empirical*	
Technical	1. Data not collected	[6,21,24,31]	[2,4,7,18,22,14,26-28,30]	[3,23,25]
	2. Data not preserved		[33]	[3,32,34,35]
	3. Data not found		[45]	[3,34]
	4. Language barrier			[36]
	5. Restrictive data format		[40]	[3,34,36-39,41]
	6. Technical solutions not available		[42]	[37]
	7. Lack of metadata and standards	[21,24,43]	[40,44,45]	[1,35-37,39,41,46]
Motivational	8. No incentives		[27,45,49]	[35]
	9. Opportunity cost	[51,52]	[13,33,50,53]	[35]
	10. Possible criticism		[33]	[32]
	11. Disagreement on data use	[21]	[49]	
Economic	12. Possible economic damage		[7,26,27,30]	[55]
	13. Lack of resources	[56,21]	[13,27,28,30,42,53,57]	[3,23,34-36,39,37]
Political	14. Lack of trust	[19,59,60]	[33,61]	[34-37]
	15. Restrictive policies		[30]	
	16. Lack of guidelines		[45,62,65]	[37,41,63,64]
Legal	17. Ownership and copyright		[62,65,66,69]	[37,63,64,67]
	18. Protection of privacy	[12,19,59,73,75]	[44,57,62,66,72,74]	[36,37,64,67,68,70,71]
Ethical	19. Lack of proportionality			[76]
	20. Lack of reciprocity	[51,52]	[50,77,78]	
Number of unique documents (% of total)		14 (21.5%)	30 (46.2%)	21 (32.3%)

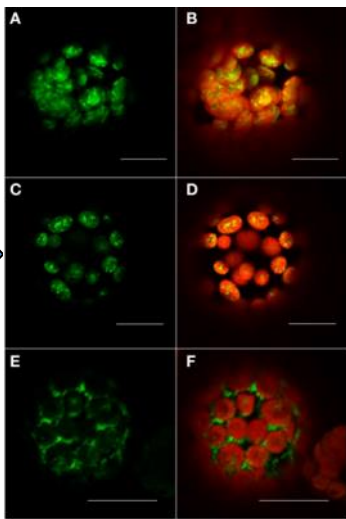
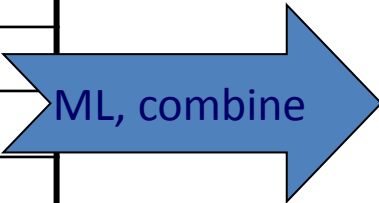
\*No or little original data presented.

# What is the Open Research Data Pilot?

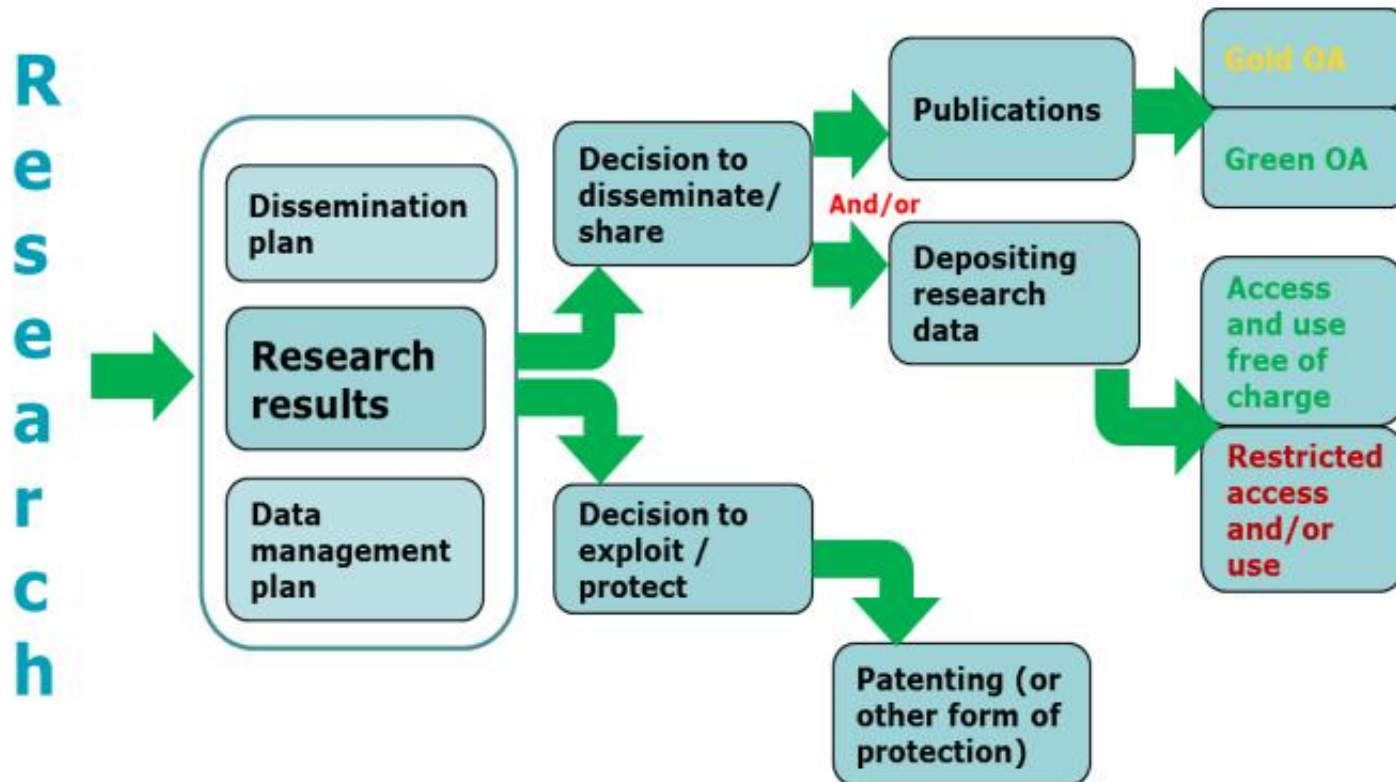
Updated on 15 November 2016

“Don’t panic - **you are not expected to share sensitive data or breach any IPR agreements** with industrial partners. You **do not need to deposit all the data** you generate during the project either – only that which underpins published research findings and/or has longer-term value. In addition to supporting your research’s integrity, openness has many other benefits. **Improved visibility** means your research will reach more people and have a greater impact – for science, society and your own career. “

	$XP_1$	$XP_2$	...	$XP_n$
Gene 1	1.1	1.2	23.1	22.1
Gene 2	1.5	5.3	12.1	5.6
Gene 3	9.1	4.2	9.2	4.3
..				
Gene m	8.1	4.3	1.3	1.5



Machine Learning, data mining

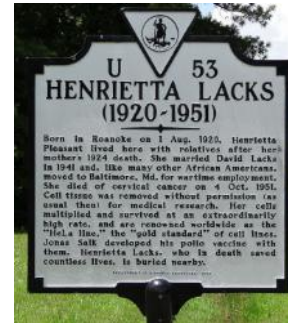


**Graph: Open access to scientific publication and research data in the wider context of dissemination and exploitation**



# Opting OUT

- participation is incompatible with the obligation to protect results that can
- **reasonably be expected to be commercially or industrially exploited**
- participation is incompatible with the need for confidentiality in connection with security issues
- **participation is incompatible with rules on protecting personal data**
- participation would mean that the project's main aim might not be achieved
- the project will not generate / collect any research data or
- there are other legitimate reasons (you can enter these in a free-text box at the proposal stage).



*Thank you*

*Questions?*