

# **Core DS Knowledge Model**

Organization: DSW Global

Created by: Gerald W (e0125536@student.tuwien.ac.at)
Based on: Common ELIXIR Knowledge Model, 1.0.0

Project Phase: Before Submitting the Proposal

Created at: 22.04.2019

# I. Design of experiment

Before you decide to embark on any new study, it is nowadays good practice to consider all options to keep the data generation part of your study as limited as possible. It is not because we can generate massive amounts of data that we always need to do so. Creating data with public money is bringing with it the responsibility to treat those data well and (if potentially useful) make them available for re-use by others.

### Report

#### Indications

Answered	0
Unanswered	0

#### **Metrics**

Metric	Score
Findability	0
Accessibility	0
Interoperability	0
Reusability	0
Good DMP Practice	0
Openness	0

1

#### Is there any pre-existing data?

Are there any data sets available in the world that are relevant to your planned research?

Data Stewardship for Open Science: <u>atq</u>

☑ External Links:

# b. Yes

1.b.1

## Will you be using any pre-existing data (including other people's data)?

Will you be referring to any earlier measured data, reference data, or data that should be mined from existing literature? Your own data as well as data from others?

■ Data Stewardship for Open Science: ezi

External Links:

b. Yes

1.b.1.b.1

#### What reference data will you use?

Much of todays data is used in comparison with reference data. A genome for instance is compared with a reference genome to identify genomic variants. If you use reference data, there are several other issues that you should consider. What are the reference data sets that you will use?

- Data Stewardship for Open Science: <u>quc</u>
- ☑ External Links:

#### **Answers**

Item 1.b.1.b.1.a Fremont Bridge Hourly Bicycle Counts by Month October 2012 to present

1.b.1.b.1.a.1

#### Do you know where and how is it available?

Do you know where the reference data is available, what the conditions for use are, and how to reference it?

- Data Stewardship for Open Science: <u>ckt</u>
- ☑ External Links:

b. Yes

1.b.1.b.1.a.2

### Do you know in what format the reference data is available?

Do you know the data format of the reference data? Is this suitable for your work? Does it need to be converted?

- Data Stewardship for Open Science: <u>jxb</u>
- □ External Links:

a. I can directly use it

1.b.1.b.1.a.3

# Is the reference data resource versioned?

Many reference data sets evolve over time. If the reference data set changes, this may affect your results. If different versions of a reference data set exist, you need to establish your "version policy".

- Data Stewardship for Open Science: <u>rgy</u>
- External Links:

a. No

1.b.1.b.1.a.4

### How will you make sure the same reference data will be available to reproduce your results?

Will the reference data in the version you use be available to others?

a. I will keep a copy and make it available with my results

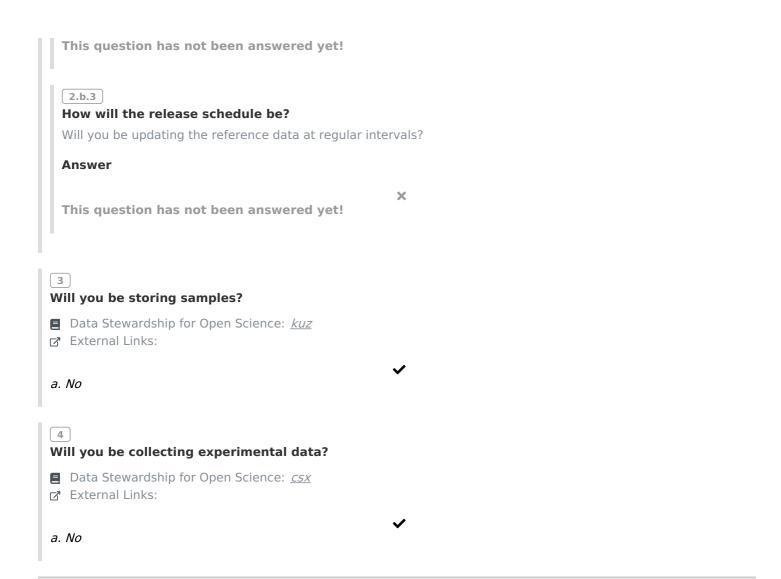
1.b.1.b.2

#### What existing non-reference data sets will you use?

Even if you will be producing your own data, you often will also be relying on existing data sets (e.g. from earlier . You may need to integrate your new data with an existing data set or retrieve additional information from related data bases. Will you be doing such things?

- Data Stewardship for Open Science: wya
- External Links:

# **Answer** × This question has not been answered yet! 1.b.2 Do you need to harmonize different sources of existing data? If you are combining data from different sources, harmonization may be required. You may need to re-analyse some original data. Data Stewardship for Open Science: wht ☑ External Links: a. No 1.b.3 Will you be using data that needs to be (re-)made computer readable first? Some old data may need to be recovered, e.g. from tables in scientific papers or may be punch cards. ■ Data Stewardship for Open Science: <u>pth</u> ☑ External Links: a. No 2 Will reference data be created? Will any of the data that you will be creating form a reference data set for future research (by others)? ■ Data Stewardship for Open Science: *rbz* ☑ External Links: b. Yes 2.b.1 What will the Intellectual Property be like? Who will own the rights to the reference data set? Who will be able to use it? ■ Data Stewardship for Open Science: <u>hct</u> ☑ External Links: **Answer** This question has not been answered yet! 2.b.2 How will you maintain it? How will maintenance be paid for in the long run? Will you host it yourself or deposit it with a repository? How will you deal with requests for help? And with requests for adding data? ■ Data Stewardship for Open Science: <u>usx</u> External Links: **Answer**



# II. Data design and planning

In the data design and planning phase, we will make sure that we know what data comes when, that we have enough storage space and compute power to deal with it, and that all the responsibilities have been taken care of.

### **Report**

#### **Indications**

Answered	0
Unanswered	0

#### **Metrics**

Metric	Score
Findability	0.5
Accessibility	0.5
Interoperability	1
Reusability	0.5
Good DMP Practice	0.5
Openness	0

# What data formats/types will you be using? Have you identified types of data that you will use that are used by others too? Some types of data (e.g. genetic variants in the life sciences) are used by many different projects. For such data, often common standards exist that help to make these data reusable. Are you using such common data formats? Data Stewardship for Open Science: <u>njy</u> ☑ External Links: **Answers** Item 1.a CSV 1.a.1 Is this a standard data format used by others too? b. Yes 1.a.2 Does this data format enable sharing and long term archiving? Complicated (binary) file formats tend to change over time, and software may not stay compatible with older versions. Also, some formats hamper long term usability by making use of patents or being hampered by restrictive licensing b. Yes 2 Will you be using new types of data? Sometimes the type of data you collect can not be stored in a commonly used data format. In such cases you may need to make your own, keeping interoperability as high as possible. Data Stewardship for Open Science: <u>ikk</u>

- ☑ External Links:
- a. No, all of my data will fit in common formats

3

# How will you be storing metadata?

For the re-usability of your data by yourself or others at a later stage, a lot of information about the data, how it was collected and how it can be used should be stored with the data. Such data about the data is called metadata, and this set of questions are about this metadata

- Data Stewardship for Open Science: <u>rhm</u>
- rオ External Links:
- a. Explore

3.a.1

Do suitable 'Minimal Metadata About ...' (MIA...) standards exist for your experiments?

a. No

Did you really check a service like fairsharing.org to verify this?

3.a.1.a.1

Do you have a good idea of what metadata is needed to make it possible for others to read and interpret your data in the future?

3.a.2 Do you know how and when you will be collecting the necessary metadata? Often it is easiest to make sure you collect the metadata as early as possible. a. No 3.a.3 Will you consider re-usability of your data beyond your original purpose? Adding more than the strict minimum metadata about your experiment will possibly allow more wide re-use of your data, with associated higher data citation rates. Please note that it is not easy for yourself to see all other ways in which others could be reusing your data. a. No, I will just document the bare minimum 3.a.4 Did you consider how to monitor data integrity? Working with large amounts of heterogenous data in a larger research group has implications for the data integrity. How do you make sure every step of the workflow is done with the right version of the data? How do you handle the situation when a mistake is uncovered? Will you be able to redo the strict minimum data handling? ■ Data Stewardship for Open Science: <u>spg</u> ☐ External Links: a. Explore 3.a.4.a.1 Will you be keeping a master list with checksums of certified/correct/canonical/verified data? Data corruption or mistakes can happen with large amounts of files or large files. Keeping a master list with data checksums can be helpful to prevent expensive mistakes. It can also be helpful to keep the sample list under version control forcing that all changes are well documented. b. Yes Will you define a way to detect file or sample swaps, e.g. by measuring something independently?

a. No

#### Do all datasets you work with have a license?

It is not always clear to everyone in the project (ad outside) what can and can not be done with a data set. It is helpful to associate each data set with a license as early as possible in the project. A data license should ideally be as free as possible: any restriction like 'only for non-commercial use' or 'attribution required' may reduce the reusability and thereby the number of citations. If possible, use a computer-readable and computer actionable license.

b. Yes

3.a.5.b.1 Will you store the licenses with the data at all time? It is very likely that data will be moved and copied. At some point people may lose track of the origins. It can be helpful to have the licenses (of coarse as open as possible) stored in close association with the data. ■ Data Stewardship for Open Science: <u>atw</u> ☑ External Links: b. Yes 3.a.6 How will you keep provenance? To make your experiments reproducible, all steps in the data processing must be documented in detail. The software you used, including version number, all options and parameters. This information together for every step of the analysis is part of the so-called data provenance. There are more questions regarding this in the chapter on data processing and curation. a. All steps will be documented in an (electronic) lab notebook 3.a.7 How will you do file naming and file organization? Putting some thoughts into file naming can save a lot of trouble later. a. Explore 3.a.7.a.1 Did you make a SOP (Standard Operating Procedure) for file naming? It can help if everyone in the project uses the same naming scheme. **Answer** × This question has not been answered yet! 3.a.7.a.2 Will you be keeping the relationships between data clear in the file names? Advice: Use the same identifiers for sample IDs etc throughout the entire project. a. No 3.a.7.a.3

### Will all the metadata in the file names also be available in the proper metadata?

The file names are very useful as metadata for people involved in the project, but to computers they are just identifiers. To prevent accidents with e.g. renamed files metadata information should always also be available elsewhere and not only through the file name.

b. Yes, all metadata is also explicitly available elsewhere

4

Much of the raw data you have will need to be archived for your own later use somewhere. This is often done off-line on tape, not on the disks of the compute facility. Please note that this does not refer to the data publication.
■ Data Stewardship for Open Science: <u>kjp</u> ☑ External Links:
a. No
4.a.1 Can the original data be regenerated?
<ul><li>■ Data Stewardship for Open Science: <u>ixr</u></li><li>☑ External Links:</li></ul>
a. No
4.a.2 When is the raw data archived?
• All at an accuit the the manufact the and of the manifest
c. All at once with the results at the end of the project
Will you need a shared working space to work with your data?
<b>✓</b>
a. No
5.a.1
Are data all project members store adequately backed up and traceable?
Are data all project members store adequately backed up and traceable?  •  a. No
<b>✓</b>
<b>✓</b>
a. No  6 Is the risk of information loss, leaks and vandalism acceptably low?
a. No
a. No  6 Is the risk of information loss, leaks and vandalism acceptably low?  There are many factors that can contribute to the risk of information loss or information leaks. They are often part of the behavior of the people that are involved in the project, but can also be steered by properly planned
a. No  6 Is the risk of information loss, leaks and vandalism acceptably low?  There are many factors that can contribute to the risk of information loss or information leaks. They are often part of the behavior of the people that are involved in the project, but can also be steered by properly planned infrastructure.
a. No  Is the risk of information loss, leaks and vandalism acceptably low?  There are many factors that can contribute to the risk of information loss or information leaks. They are often part of the behavior of the people that are involved in the project, but can also be steered by properly planned infrastructure.  a. Explore
a. No  6 Is the risk of information loss, leaks and vandalism acceptably low?  There are many factors that can contribute to the risk of information loss or information leaks. They are often part of the behavior of the people that are involved in the project, but can also be steered by properly planned infrastructure.  a. Explore  6.a.1  Do project members store data or software on computers in the lab or external hard drives
a. No  6 Is the risk of information loss, leaks and vandalism acceptably low?  There are many factors that can contribute to the risk of information loss or information leaks. They are often part of the behavior of the people that are involved in the project, but can also be steered by properly planned infrastructure.  a. Explore  6.a.1  Do project members store data or software on computers in the lab or external hard drives connected to those computers?  When assessing the risk, take into account who has access to the lab, who has (physical) access to the computer
a. No  6 Is the risk of information loss, leaks and vandalism acceptably low?  There are many factors that can contribute to the risk of information loss or information leaks. They are often part of the behavior of the people that are involved in the project, but can also be steered by properly planned infrastructure.  a. Explore  6.a.1  Do project members store data or software on computers in the lab or external hard drives connected to those computers?  When assessing the risk, take into account who has access to the lab, who has (physical) access to the computer hardware itself. Also consider whether data on those systems is properly backed up
a. No  6 Is the risk of information loss, leaks and vandalism acceptably low? There are many factors that can contribute to the risk of information loss or information leaks. They are often part of the behavior of the people that are involved in the project, but can also be steered by properly planned infrastructure.  a. Explore  6.a.1  Do project members store data or software on computers in the lab or external hard drives connected to those computers?  When assessing the risk, take into account who has access to the lab, who has (physical) access to the computer hardware itself. Also consider whether data on those systems is properly backed up  a. No
a. No  6 Is the risk of information loss, leaks and vandalism acceptably low?  There are many factors that can contribute to the risk of information loss or information leaks. They are often part of the behavior of the people that are involved in the project, but can also be steered by properly planned infrastructure.  a. Explore  6.a.1 Do project members store data or software on computers in the lab or external hard drives connected to those computers?  When assessing the risk, take into account who has access to the lab, who has (physical) access to the computer hardware itself. Also consider whether data on those systems is properly backed up  a. No  6.a.2

	6.a.3
	Do project members store project data in cloud accounts?
	Think about services like Dropbox, but also about Google Drive, Apple iCloud accounts, or Microsoft's Office365
	a. No
	6.a.4 Do project members send project data or reports per e-mail or other messaging services?
	✓
	a. No
	6.a.5  Do all data centers where project data is stored carry sufficient certifications?
	<b>✓</b>
	b. Yes
ì	
	Are all project web corvices addressed via secure bttp (bttps://)?
	Are all project web services addressed via secure http (https://)?
	<b>✓</b>
	b. Yes
ı	
Ì	6.a.7
	Have project members been instructed about the risks (generic and specific to the project)?
	Project members may need to know about passwords (not sharing accounts, using different passwords for each
	service, and two factor authentication), about security for data they carry (encryption, backups), data stored in their own labs and in personal cloud accounts, and about the use of open WiFi and https
	<b>✓</b>
	a. No
ì	6.a.8
	Did you consider the possible impact to the project or organization if information is lost?
	bia you consider the possible impact to the project of organization is instru
	b. Yes; the effect is small
	D. Tes, the effect is sinall
i	
	6.a.9
	Did you consider the possible impact to the project or organization if information leaks?
	<b>✓</b>
	b. Yes; the effect is small
ĺ	
	6.a.10
	Did you consider the possible impact to the project or organization if information is vandalized?
	<b>✓</b>
	a. No

**7** 

# Do you need to do compute capacity planning?

If you require substantial amounts of compute power, amounts that are not trivially absorbed in what you usually have abailable, some planning is necessary. Do you think you need to do compute capacity planning?



# Report

#### **Indications**

Answered	0
Unanswered	0

#### **Metrics**

Metric	Score
Findability	0
Accessibility	0
Interoperability	1
Reusability	0
Good DMP Practice	0
Openness	0

1

Please specify what data sets you will acquire using measurement equipment

### **Answers**

Item 1.a Bicycle count automatically via induction loops by Seattle Department of Transportation

1.a.1

# Who will do the measurements? And where?

Are there easily accessible specialized service providers for data capture?

c. External party

consider making them partner in the project

1.a.1.c.1

Has formal ownership of the data been established?

a. The party measuring the data owns it

1.a.1.c.2

Has responsibility for long term safe keeping of the raw data been established? Who will deal with data publication?

d. We have made other arrangements

1.a.1.c.2.d.1

What other arrangements?

# IV. Data processing and curation

### Report

**Indications** 

Answered 0

#### **Metrics**

Metric	Score
Findability	0
Accessibility	0
Interoperability	0
Reusability	0.21052631578947367
Good DMP Practice	0
Openness	0



#### **Workflow development**

It is likely that you will be developing or modifying the workflow for data processing. There are a lot of aspects of this workflow that can play a role in your data management, such as the use of an existing work flow engine, the use of existing software vs development of new components, and whether every run needs human intervention or whether all data processing can be run in bulk once the work flow has been defined.

a. This has been arranged



How will you make sure to know what exactly has been run?

**~** 

a. Explore

2.a.1

Will you keep results together with all processing scripts or workflows including documentation of the versions of the tools that have been run?

a. No

2.a.2

# Will you make use of the metadata fields in your output data files to register how the data was obtained?

File formats like VCF (for genetics) and TIFF (for images) have possibilities to document metadata in the file header. It is a good idea to use work flow tools that use these fields to document what was done to obtain the data.

a. No

**/** 

2.a.3

## Will you use a central repository for all tools and their versions as used in your project?

Especially if analysis and processing of data in the project is done on multiple different computers by different people, it is a good idea to have your own repository of tools and their blessed versions.

■ Data Stewardship for Open Science: <u>pzq</u>

r External Links:

a. No

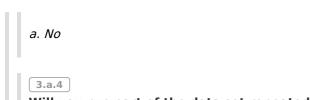
**✓** 

2.a.4 Will you use a central repository for reference data used in your project? Especially if analysis and processing of data in the project is done on multiple different computers by different people, it is a good idea to have your own repository of reference data versions. ■ Data Stewardship for Open Science: pzg r₹ External Links: b. Yes Will you make use of standard workflow engines and automatic work flows for all data analysis in the project? It is much easier to guarantee consistency and reproducibility if all data processing is done using automated work flows, especially if the workflow engine automatically keeps adequate provenance data. a. No 2.a.6 Are all software tools in the work flow professionally maintained, with version control? Will you be able to find and reproduce exactly which version was used for any analysis? Not only for the major tools in the workflows, but also for all 'glue' code and small tools you created especially for the project? b. Yes 3 How will you validate the integrity of the results? a. Explore 3.a.1 Will you run a subset of your jobs several times across the different compute infrastructures you are There are surprisingly many complications that can cause (slight) inconsistencies between results when workflows are run on different compute infrastructures. A good way to make sure this does not bite you is to run a subset of all jobs on all different infrastructure to check the consistency. a. No 3.a.2 Will you be instrumenting the tools into pipelines and workflows using automated tools? Surrounding all tools in your data processing and analysis workflows with the 'boilerplate' code necessary on the computer system you are using is tedious and error prone. Especially if you are using the same tools in multiple different work flows and/or on multiple different computer architectures. Automated instrumentation, e.g. by using a workflow management system, can prevent many mistakes.

a No

# Will you use independently developed duplicate tools or workflows for critical steps to reduce or eliminate human errors?

Validation of results without a golden standard is very hard. One way of doing it is to develop two solutions for a problem (two independent workflows or two independently developed tools) to check whether the results are identical or comparable.



Will you run part of the data set repeatedly to catch unexpected changes in results?

Running a small subset of the data repeatedly can be useful to catch unexpected problems that would otherwise be very hard to detect.

- Data Stewardship for Open Science: <u>egv</u>
- ☑ External Links:

a. No



# Do you have a contingency plan?

What will you do if the compute facility is down?

a. We will wait until the problem is fixed

# V. Data integration

# Report

# Indications

Answered	0
Unanswered	0

#### **Metrics**

Metric	Score
Findability	0
Accessibility	0
Interoperability	0.38095238095238093
Reusability	0
Good DMP Practice	0
Openness	0

1

How will you be doing the integration of different data sources?

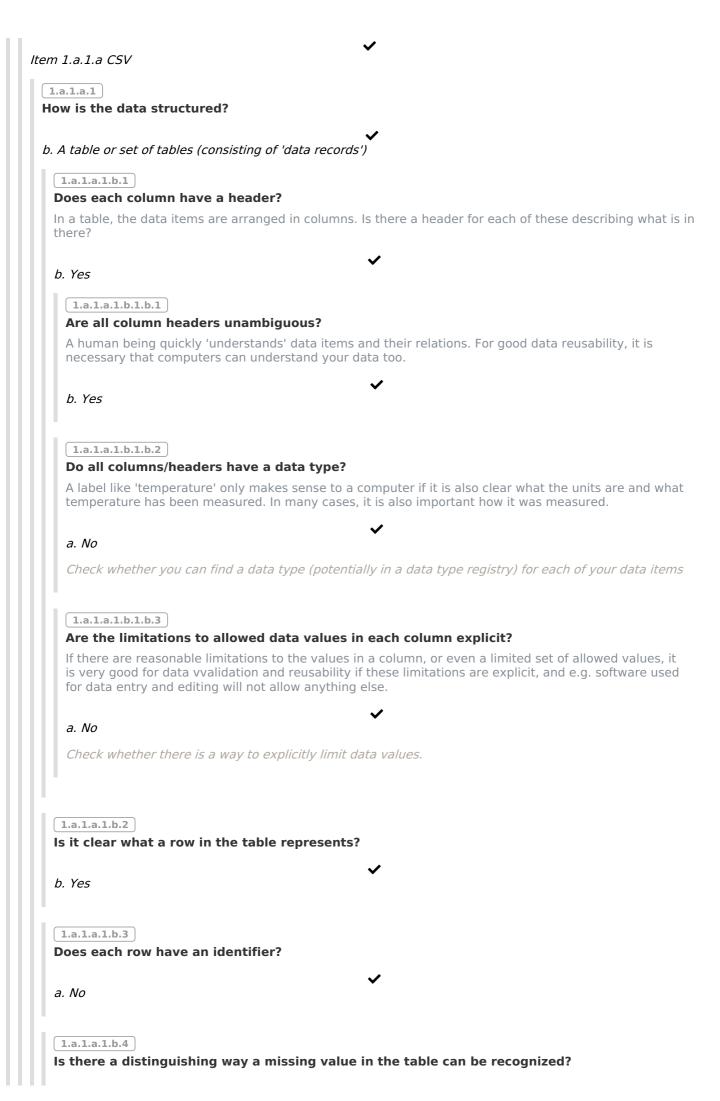
a. Explore

1.a.1

List the data formats you will be using for data integration

Answer some questions for each

Answers



Sometimes, and empty field or a zero is indicating a missing value. But is that really unique? Could there be valid empty or zero fields? Has the convention for missing values been made explicit somewhere? b. Yes 1.a.1.a.1.b.5 Is the relation between each of the columns and the record identifier clear? It may appear that in a table with 'patients' as rows, a column labeled 'disease' coupled to an ontology has a clear meaning. But that is not always explicit enough! A 'disease' could e.g. be the disease that the patient is suffering from, but it could also be an earlier diagnose, a suspected diagnose, or the disease a family member recently died of. a. No 1.a.1.a.1.b.6 Are all the relations between the column headers explicit? For a good understanding of tabular data, you need to make the relationship between each pair of columns explicit. E.g. if one column is 'disease' and another is 'treatment', you want to make sure that this is the chosen treatment that this person is undergoing for the given disease. a. No 1.a.2 Will you be using a workflow for data integration, e.g. with tools for database access or conversion? Data Stewardship for Open Science: <u>gab</u> ☑ External Links: a. No 1.a.3 Will you use a 'linked data' approach? a. No Will you be using common or exchangeable units? a. No Will you be using common ontologies? a. No Will there be potential issues with statistical normalization?

3

4

Will you be integrating different data sources to get more samples or more data points?
a. No
6 Will you be integrating different data sources in order to get more information for each sample or data point?
a. No
7 Do you have all tools to couple the necessary data types?

# VI. Data interpretation

# **Report**

a. No

a. No

#### **Indications**

Answered	0
Unanswered	0

### **Metrics**

2

Metric	Score
Findability	0
Accessibility	0
Interoperability	0
Reusability	1
Good DMP Practice	0
Openness	0

Will data interpretation and modeling require significant compute infrastructure capacity?a. No

How will you be making sure there is good provenance of the data analysis?

Data analysis is normally done manually on a step-by-step basis. It is essential to make sure all steps are properly documented, otherwise results will not be reproducible.

3

Will you be doing (automated) knowledge discovery?

- Data Stewardship for Open Science: <u>bzu</u>
- ☑ External Links:

a. No

# VII. Information and insight

# Report

#### Indications

Answered	0
Unanswered	0

#### **Metrics**

Metric	Score
Findability	0.4
Accessibility	0
Interoperability	0
Reusability	1
Good DMP Practice	0
Openness	0.666666666666666

Will you be working with the philosophy 'as open as possible' for your data?

□ Data Stewardship for Open Science: jvm
□ External Links:

b. Yes

Can all of your data become completely open immediately? *▶ b. Yes* 

3

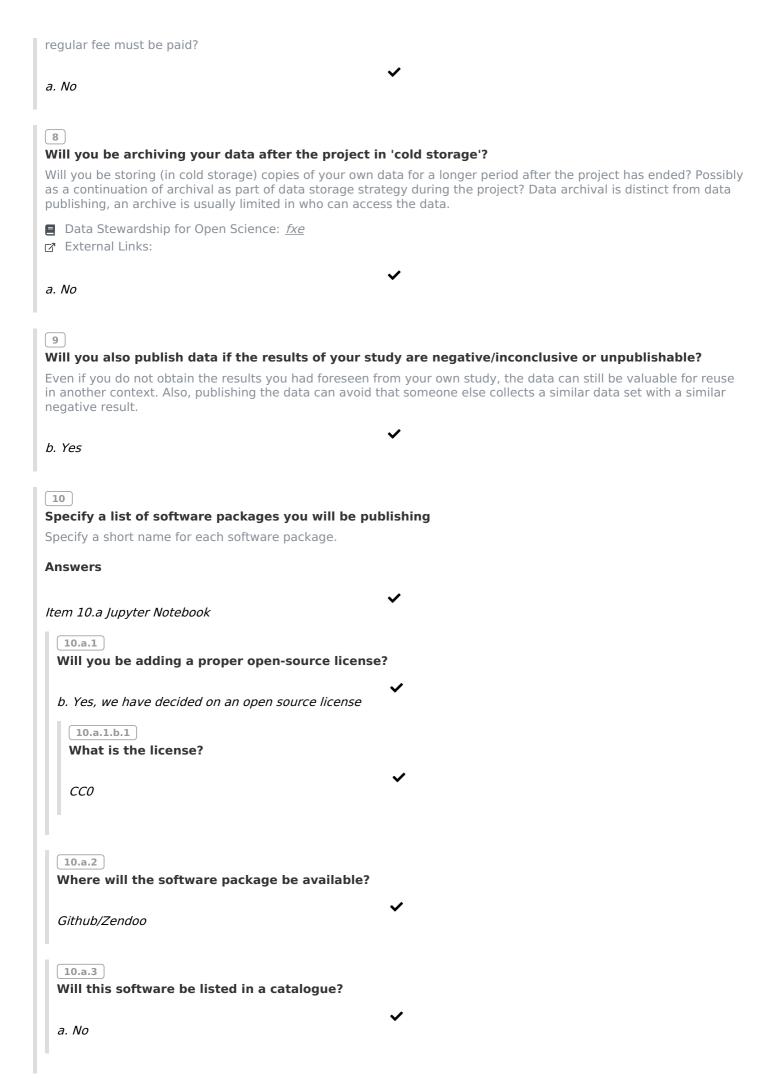
# Specify a list of data sets you will be publishing

Specify a short name for each data set, sufficient for yourself to know what data it is about. It is useful to think about a data set as some collection of data that will be ending up in the same place.

#### **Answers**

Are there any recurring fees to keep data or documents available?

Are you using any commercially licensed products to keep data, software or documents available, for which a



11

Will there be planning of valorization or translational returns?

**~** 

a. No

Data Managament Plan generated by Data Stewardship Wizard < https://ds-wizard.org>