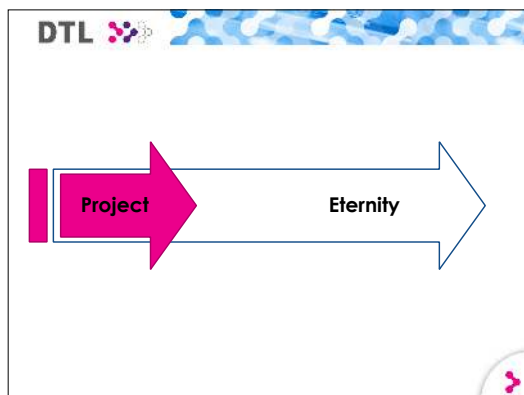


I think in many cases for projects in life sciences the budget for dealing with the data (time as well as infrastructure) is much larger than the world wide research average of 5%. The fraction is so significant that:

- * It will actually help if we can save part of that budget by planning properly
- * We can no longer afford having “forgot” to budget for part of the data management.
- * We can no longer afford to “play until it works”
- * We can no longer afford to do something really wrong and having to re-do



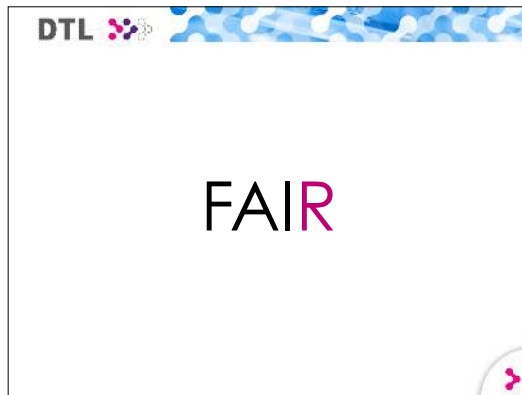
We need to plan how we work with the data, starting thoughts before the project. And, to be able to deal with changes in the project, the data stewardship plan should be updated all along with the project. That way it will also function as part of the documentation of the data after the project is finished.

Big project? Hire special data expert. They should know a lot, and also know what they do not know, and know where to find the experts.



We want to make our data Reusable.

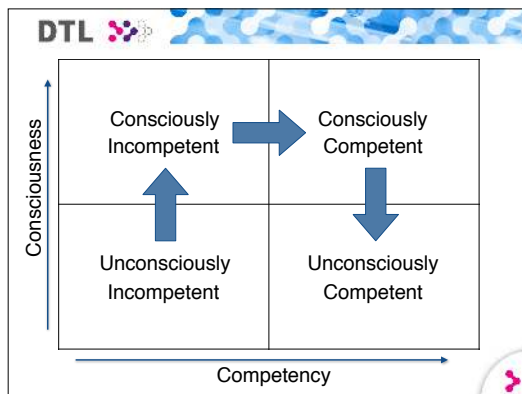
- * Reusable for others
- * Usable and Reusable for ourselves
 - * Now: because it is high volume, complex, and done with others
 - * Later: because you won't remember, and because you struggle to reproduce



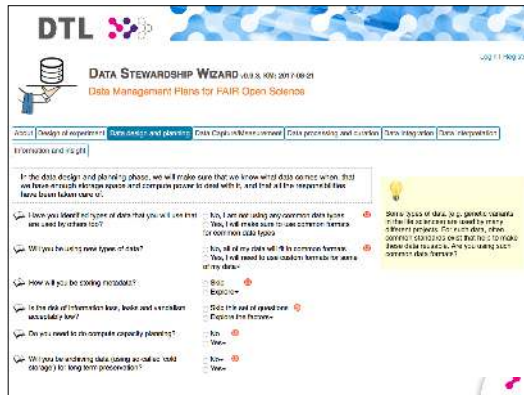
To make our data re-usable, we need to make it F, A, I.

For good data management, at every step of the way, you need to be asking yourself: what do I need to do to make sure this data will be FAIR.

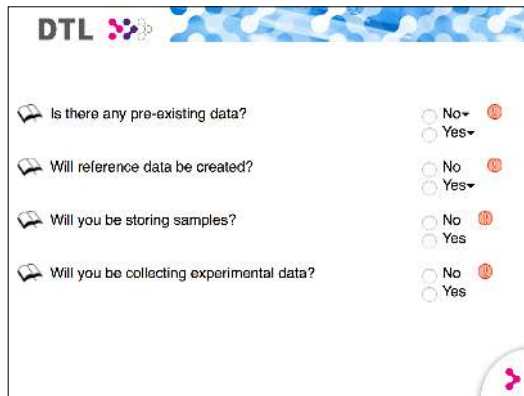
We wanted to make a tool to help with these choices. With broad target audience.



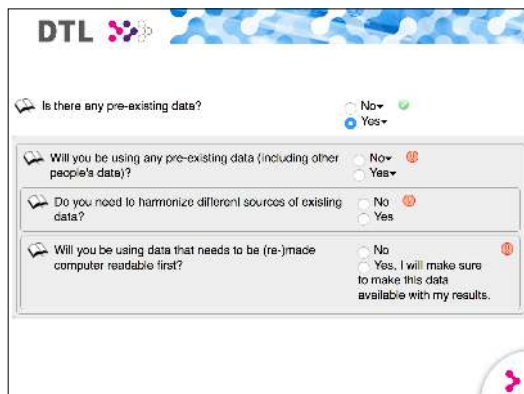
There are four stages of “knowing” things. Like car drivers of 12yo, 17yo, 19yo, 40yo. See: https://en.wikipedia.org/wiki/Four_stages_of_competence
We want the researchers to become at least consciously incompetent. They need help for their projects from people that are at least consciously competent.



Together with ELIXIR partners we are working on a tool that is not based on Funder questions, but on our own insights, collected over the last 4 years.



It presents questions in a hierarchical fashion: answers with a triangle open more questions.



This helps by only exposing issues that are applicable to a project. There are also blocks of questions that can be repeated for different resources, like different data sources or different data sets produced.

DTL

Will you be using any pre-existing data (including other people's data)? No Yes

What reference data will you use? List all the items below.

| Item | Do you know where and how is it available? | Do you know in what format the reference data is available? | Is the reference data resource versioned? | How will you make sure the same reference data will be... |
|------|---|--|---|---|
| | <input type="radio"/> No <input checked="" type="radio"/> Yes | <input type="radio"/> I can directly use it <input checked="" type="radio"/> I need to convert it before using | <input type="radio"/> No <input checked="" type="radio"/> Yes | <input type="radio"/> I will ensure... |

For Data Stewards, this system works as a checklist. For researchers that want to "do it themselves", it is a discovery system for traps and pitfalls. But we do not leave them in danger, but added guidance wherever possible.

DTL

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

Next to every question we have primary guidance popping up.

DTL

Data Stewardship for Discovery : Chapter 1.4

WHERE IS IT AVAILABLE?

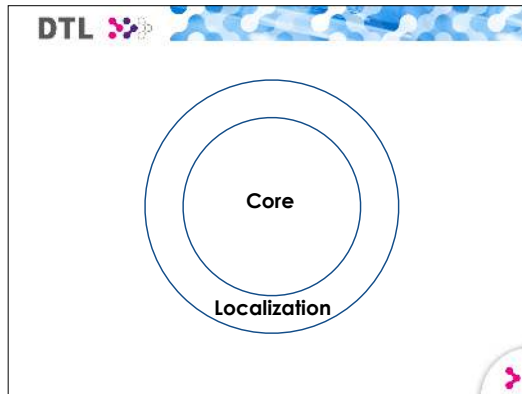
What's up?

Data sets (including reference data) may be available at different locations (in practice) and with different service level agreements attached. Before use, users should check policy and a sustainability plan to ensure that the data will be available and properly accessed for a long period of time. Some data sets need to be downloaded from external servers several times a day. It is especially important to use data sets from a stable source in immediate context. The data sets are usually available for use. It is wise to use those resources that are most stable and available and sustainable. Data that are available in a repository that is not approved or certified and/or those data for your experiment and they are no longer available later or because for instance the repository went off line or got closed off, you may get into trouble for follow up experiments or reproducibility and review issues.

Do

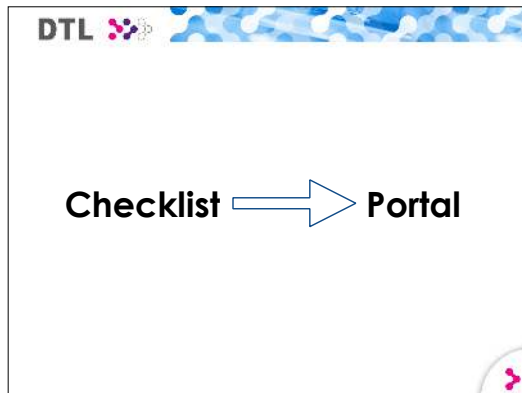
- Make sure you use data online only when there is a time horizon that it will remain available (under the same conditions) indefinitely or long time.

Further guidance: pointers to training materials and experts.



Data in the tool is highly configurable.

It is separated in Generic and Local: So we can separate life science specific (like “human health data”), ELIXIR specific (like names of experts), and even for local institutes we can adapt it.

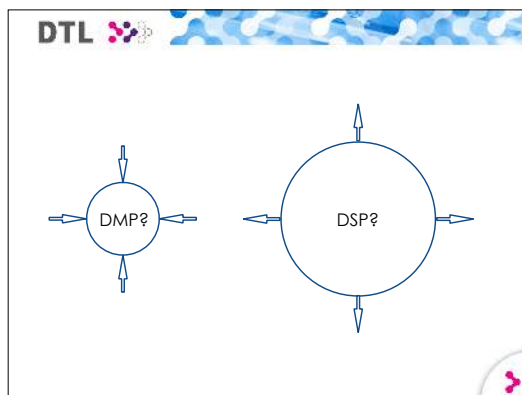


Current version is already able to help as checklist with over 300 pieces of guidance.

In the future, we want to extend the functionality:

- * Make, edit, review, publish Data Management Plans

For this, we are taking part in RDA ADMP and its working groups, and in contact with DCC DMPonline to see what we can learn and/or reuse from each other



Tendency for DMP ("data management plans") requirements is shrinking. They are often perceived as an obligation and funders don't have time to judge detailed plans.

We want to grow DSP ("data stewardship plans"). They benefit science. Carrot instead of Stick.

Three groups that can benefit are:

- * It is a checklist, like for pilots flying the data plane.