# **FAIR data in trustworthy repositories:**
# **How do I organise and preserve my research data?**

Workshop Digital Humanities – the perspective of Africa
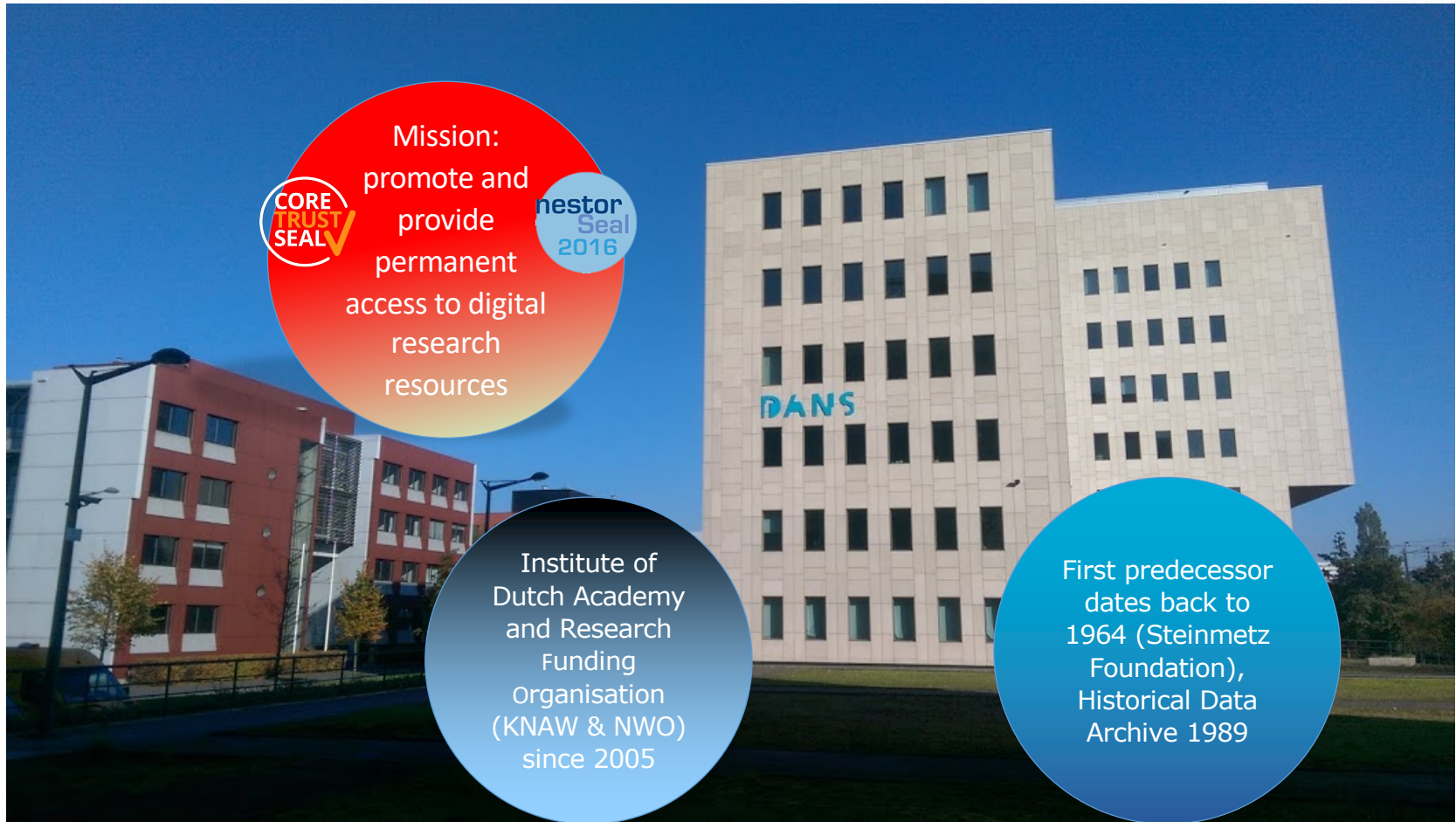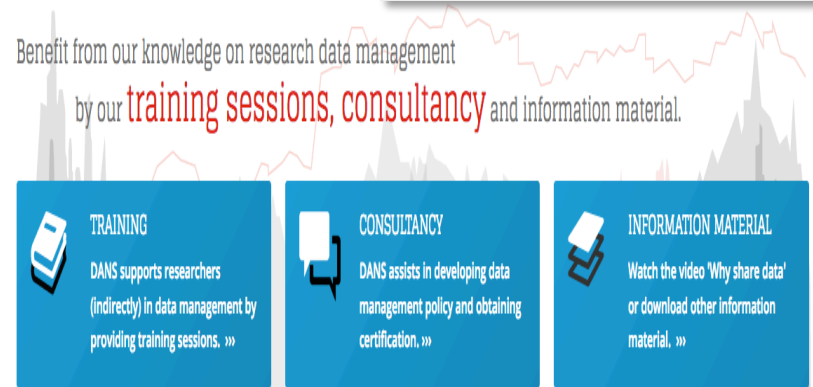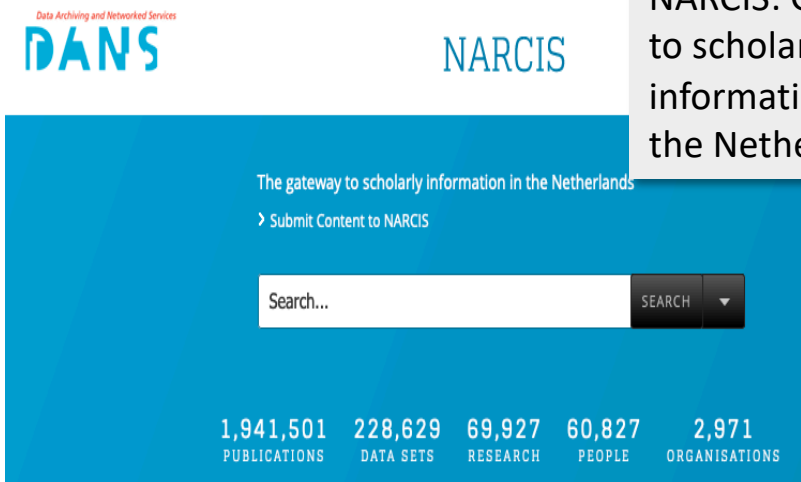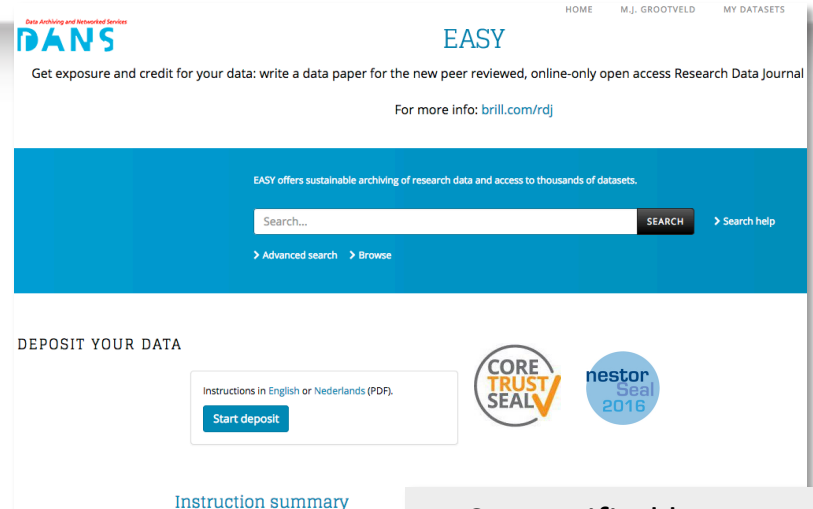
Marjan Grootveld

July 1st, 2019

1. What is DANS?
2. What is data?
3. Exercise 1: organise your data
4. Trustworthy data repositories
5. Exercise 2: find a relevant repository for your data
6. Take-aways from this session

https://www.pexels.com/photo/purple-blue-green-pink-orange-and-yellow-highlighter-159659/

# DANS   https://dans.knaw.nl/nl

Mission: promote and provide permanent access to digital research resources

CORE TRUST SEAL ✓

nestor Seal 2016

DANS

Institute of Dutch Academy and Research Funding Organisation (KNAW & NWO) since 2005

First predecessor dates back to 1964 (Steinmetz Foundation), Historical Data Archive 1989

DANS

# Core DANS services



DataverseNL for short- and mid-term data storage

EASY: certified long-term Electronic Archiving System for self-deposit

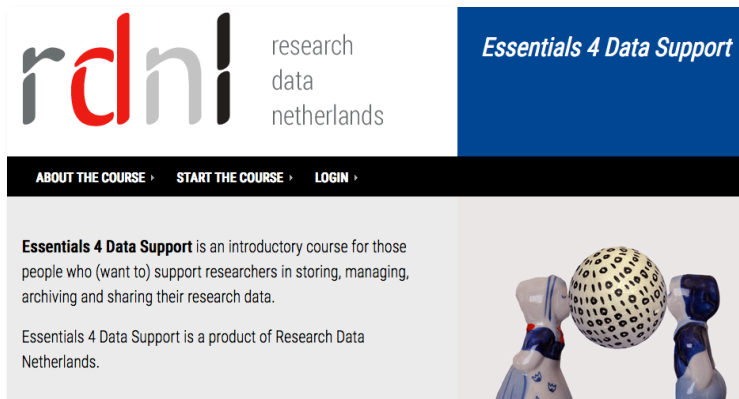NARCIS: Gateway to scholarly information in the Netherlands

# Data management training and consultancy

Partner in (inter)national projects:

FAIR and Open data

Data management planning

Trustworthy digital repositories



https://datasupport.researchdata.nl/en/



https://eudat.eu/
https://eoscpilot.eu/
https://www.eosc-hub.eu/
https://www.openaire.eu/
https://www.fairsfair.eu/

# What is research data?



An introduction to the basics of research data
https://www.youtube.com/watch?v=q2aiDJzJPuw

DANS

# From "real life" to research data: CLIWOC - climatological database for the world's oceans



Every yellow dot represents a ship report.
Image copied from https://www.knmi.nl/kennis-en-datacentrum/achtergrond/cliwoc
Project web site: http://pendientedemigracion.ucm.es/info/cliwoc/

DANS

# FAIR data principles

1. Findable – Easy to find by both humans and computer systems and based on mandatory description of the metadata that allow the discovery of interesting datasets;
2. Accessible – Stored for long term such that they can be easily accessed and/or downloaded with well-defined license and access conditions (Open Access *when possible*), whether at the level of metadata, or at the level of the actual data content;
3. Interoperable – Ready to be combined with other datasets by humans as well as computer systems;
4. Re-usable – Ready to be used for future research and to be processed further using computational methods.

http://www.nature.com/articles/sdata201618
www.force11.org/group/fairgroup/fairprinciples

https://librarycarpentry.org/Top-10-FAIR/  Top 10 FAIR Data and Software Things

DANS

# Simplified research data lifecycle

**CREATING DATA**: designing research, DMPs, planning consent, locate existing data, data collection and management, capturing and creating metadata

**PROCESSING DATA**: entering, transcribing, checking, validating and cleaning data, anonymising data, describing data, manage, store, back-up data

**RE-USING DATA**: follow-up research, new research, undertake research reviews, scrutinising findings, teaching & learning

**ANALYSING DATA**: interpreting, & deriving data, producing outputs, authoring publications, preparing for sharing

**ACCESS TO DATA**: distributing data, sharing data, controlling access, establishing copyright, promoting data

**PRESERVING DATA**: data archiving, migrating to best format & medium for long term, creating metadata and documentation

CREATING DATA

PROCESSING DATA

RE-USING DATA

GIVING ACCESS TO DATA

ANALYSING DATA

PRESERVING DATA

Based on UK Data Archive lifecycle: https://www.ukdataservice.ac.uk/manage-data/lifecycle
Used in OpenAIRE RDM briefing paper: https://www.openaire.eu/briefpaper-rdm-infonoads

DANS

# Exercise 1: Data organisation – 15 minutes

- Form a group of 4-5 people
- Read the *Veteran tapes* description
- Design a data organisation for this project:
  1. Folder structure
  2. File-naming convention
- Don't drown yourself in the details

0

15

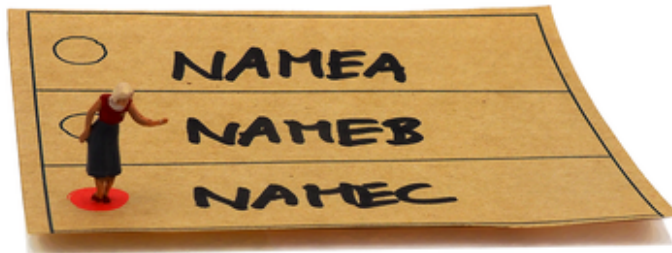# Folders and files: CESSDA DM Expert Guide



## Data Management Expert Guide

This guide is designed by European experts to help social science researchers make their research data Findable, Accessible, Interoperable and Reusable (**FAIR**).

You will be guided by different European experts who are - on a daily basis - busy ensuring long-term access to valuable social science datasets, available for discovery and reuse at one of the CESSDA social science data archives.

http://cessda.eu/dmeg   Chapter 2

**Data Management Expert Guide** ˅

1. Plan ›
2. Organise & Document ˅
   Designing a data file structure
    Organisation of variables
   **File naming and folder structure**
   Documentation and metadata
   Adapt your DMP: part 2
   Sources and further reading
3. Process ›
4. Store ›
5. Protect ›
6. Archive & Publish ›
7. Discover ›



NAMEA
NAMEB
NAMEC

**DANS**

# Plan data management for the full cycle

A Data Management Plan is a brief plan to define:

- how the data will be created
- how it will be documented
- who can access it
- where it will be stored
- **whether the data will be shared or "published"**
- **where it will be preserved**

Well, in a serious repository of course ;-)

DANS

# How to select a repository?



For giving (i.e. archiving & sharing) and for taking (i.e. reusing) data:

- Certification as a 'Trustworthy Data Repository' with an explicit ambition to keep the data available for the long term
- Matches your particular data needs:
  - e.g. file formats accepted;
  - mixture of open and restricted access;
  - usage licences
- Gives your submitted dataset a persistent and globally unique identifier for sustainable citations and to link back to particular researchers and grants
- Provides guidance on how to cite the deposited data

DANS

# Standards of trust in repositories

Formal

Extended

Core



http://www.iso16363.org/

DIN 31644
http://www.dnb.de/Subsites/nestor/
EN/Siegel/siegel.html

https://www.coretrustseal.org

DANS

# CoreTrustSeal Data Repository Certification

## 16 Requirements:

- Context (R0)
- Organisational infrastructure (R1-6)
- Digital object management (R8-14)
- Technology (R15-16)

https://doi.org/10.17026/dans-22n-gk35



25/08/2015                          Common Requirements/V2.1

**DSA–WDS Partnership
Working Group
Catalogue of Common Requirements**

**Introduction**

**Importance of Certification**

National and international funders are increasingly likely to mandate open data and data management policies that call for the long-term storage and accessibility of data.

If we want to be able to share data, we need to store them in a trustworthy digital repository. Data created and used by scientists should be managed, curated, and archived in such a way to preserve the initial investment in collecting them. Researchers must be certain that data held in archives remain useful and meaningful into the future. Funding authorities increasingly require continued access to data produced by the projects they fund, and have made this an important element in Data Management Plans. Indeed, some funders now stipulate that the data they fund must be deposited in a trustworthy repository.

Sustainability of repositories raises a number of challenging issues in different areas: organizational, technical, financial, legal, etc. Certification can be an important contribution to ensuring the reliability and durability of digital repositories and hence the potential for sharing data over a long period of time. By becoming certified, repositories can demonstrate to both their users and their funders that an independent authority has evaluated them and endorsed their trustworthiness.

**Basic Certification and its Benefits**

Nowadays certification standards are available at different levels, from a basic level to extended and formal levels. Even at the basic level, certification offers many benefits to a repository and its stakeholders.

DANS

# Main CoreTrustSeal requirements

R2. The repository maintains all applicable licenses covering data access and use and monitors compliance.

R3. The repository has a continuity plan to ensure ongoing access to and preservation of its holdings.

R4. The repository ensures, to the extent possible, that data are created, curated, accessed, and used in compliance with disciplinary and ethical norms.

R7. The repository guarantees the integrity and authenticity of the data.

R8. The repository accepts data and metadata based on defined criteria to ensure relevance and understandability for data users.

R10. The repository assumes responsibility for long-term preservation and manages this function in a planned and documented way.

R11. The repository has appropriate expertise to address technical data and metadata quality and ensures that sufficient information is available for end users to make quality-related evaluations.

R13. The repository enables users to discover the data and refer to them in a persistent way through proper citation.

R14. The repository enables reuse of the data over time, ensuring that appropriate metadata are available to support the understanding and use of the data.

https://www.coretrustseal.org/

DANS

# And other repositories?

Repositories without a trustworthy, long-term ambition may have a simpler process for depositing and preserving data:

- typically, they don't ask for preferred file formats – because they won't convert or migrate the data to new formats in future (mere "bit preservation");

- they may be less demanding (or helpful!) regarding metadata, and

- they won't remind data producers to add documentation – which probably diminishes the interpretability and reusability of the data;

- they may not have long-term budget, qualified staff, appropriate technical infrastructure nor a continuity plan, should the organisation or the budget fail.

# Exercise 2: Use re3data to find a repository

http://www.re3data.org/

Read the *Veteran tapes* project brief and identify what should be kept for the long term (3 mins)

Search re3data.org for repositories (10 mins), considering:
1. Data type(s)
2. Discipline
3. Repository features

0

15

DANS

# It's all about trust



- All data needs to be properly managed.

- Decisions made early affect what you can do later. For instance:
  - Folder structure with authorisations
  - File-naming conventions
  - Domain metadata
  - Accompanying documentation

- Strong resemblance CoreTrustSeal requirements and FAIR principles: ongoing access, explicitness & clarity, metadata, persistent references, documentation, data discovery, understandability, reuse, …

- Depositing data in a certified repository makes life easier for researchers and keeps FAIR data FAIR.

DANS

## All's FAIR that ends FAIR – any questions?



Acknowledgements:

https://eoscpilot.eu/

https://eudat.eu/

https://www.fairsfair.eu/

https://www.eosc-hub.eu/

https://www.openaire.eu/

marjan.grootveld@dans.knaw.nl

https://dans.knaw.nl/en/projects/