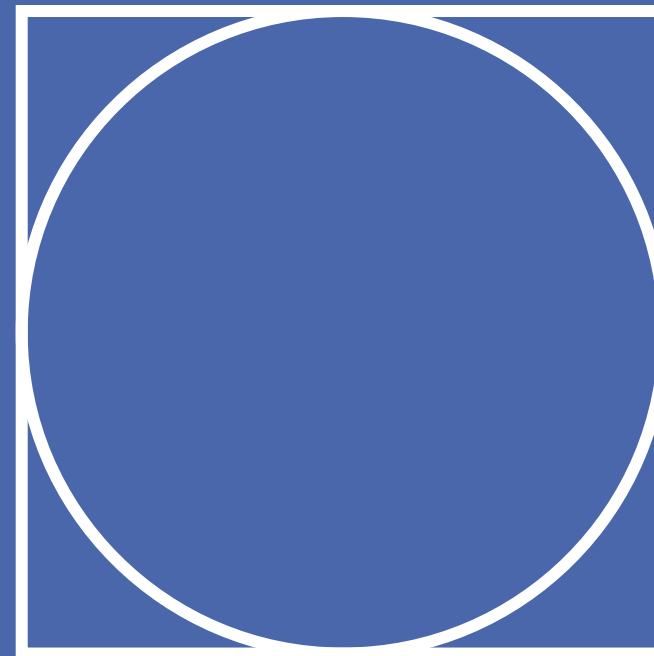# FAIR DATA

Mareike Christina Harms Buss (mabu.lib@cbs.dk)
RDM Senior Adviser, Copenhagen Business School

# Funders require FAIR data

FAIR and responsible data management is a central funder requirement:

## Open science: research data management

The beneficiaries must manage the digital research data generated in the action ('data') responsibly, in line with the FAIR principles and by taking all of the following actions:

- establish a data management plan ('DMP') (and regularly update it)

- as soon as possible and within the deadlines set out in the DMP, deposit the data in a trusted repository; if required in the call conditions, this repository must be federated in the EOSC in compliance with EOSC requirements

- as soon as possible and within the deadlines set out in the DMP, ensure open access — via the repository — to the deposited data, under the latest available version of the Creative Commons Attribution International Public License (CC BY) or Creative Commons Public Domain Dedication (CC0) or a licence with equivalent rights, following the principle 'as open as possible as closed as necessary', unless providing open access would in particular:

  - be against the beneficiary's legitimate interests, including regarding commercial exploitation, or

  - be contrary to any other constraints, in particular the EU competitive interests or the beneficiary's obligations under this Agreement; if open access is not provided (to some or all data), this must be justified in the DMP

Horizon Europe Annotated Model Grant Agreement

# Journals require reproducible or replicable data

Data reproducibility or replicability is a central journal requirement:

It is the policy of the American Economic Association to publish papers only if the data and code used in the analysis are clearly and precisely documented and access to the data and code is nonexclusive to the authors.

Authors of accepted papers that contain empirical work, simulations, or experimental work must provide, prior to acceptance, information about the data, programs, and other details of the computations sufficient to permit replication, as well as information about access to data and programs.

AEA Data and Code Availability Policy

The intent of *Journal of Marketing*'s Research Transparency policy is to (1) ensure the availability of the material necessary to evaluate and, as appropriate, replicate findings reported in the *Journal* as part of a robust review process, and (2) ensure that papers published in the *Journal* contribute to the development of cumulative, reliable, and applicable knowledge. Closing transparency gaps and ensuring safe data retention will bolster confidence not only in individual articles but also in the larger body of knowledge offered by the *Journal*.

JoM's policy of research transparency

CBS

# What do researchers want to do with their *primary data* after the end of the project?
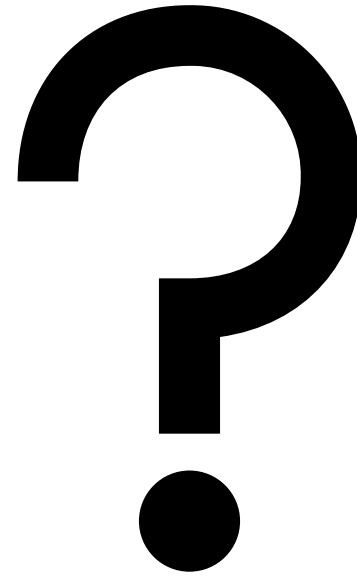
- *Publish the data for reuse*, e.g. as a data paper or in a data repository:

  - ➢ FAIR data package

- *Publish the data for replication*, e.g. as supplementary material to an article:

  - ➢ Replication package

  - ➢ (Executable) research compendium

- *Not publish the data*, but curate them for own data reuse

**CBS**

**But what about secondary data?**

**And, more specifically, confidential secondary data like register data?**

?

# HOW TO CREATE
# FAIR DATA PACKAGES

# The FAIR principles

| 地区 | 以工作日为计的供货 |
|---|---|
| 澳大利亚 | 5-8 |
| 奥地利 | 7-10 |
| 比利时 | |
| 巴西 | |
| 加拿大 | 10-14 |
| 中国 | 7-10 |
| 捷克共和国 | 7-10 |
| 丹麦 | 7-10 |
| 芬兰 | 7-10 |
| 法国 | 7-10 |
| 德国 | 7-10 |
| 香港 | 8-10 |
| 匈牙利 | 7-10 |
| 爱尔兰 | 7-10 |

Page not found

http://www.lego.com/404

Slide by Falco Hüser, Royal Library

CBS

# The FAIR principles



Illustration by Patrick Hochstenbach published in the Open Science Training Handbook under CC0

# FAIR PRINCIPLES

## Findable:

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

F2. data are described with rich metadata;

F3. metadata clearly and explicitly include the identifier of the data it describes;

F4. (meta)data are registered or indexed in a searchable resource;

## 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;

## Interoperable:

I1. (meta)data use a formal, accessible, shared, and broadly applicable language for knowledge representation.

I2. (meta)data use vocabularies that follow FAIR principles;

I3. (meta)data include qualified references to other (meta)data;

## Reusable:

R1. (meta)data are richly described with 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 detailed provenance;

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

# Making your data FINDABLE

Publish your data in a **searchable ressource** like Dataverse, Zenodo or Figshare

# Making your data FINDABLE

Assign **persistent identifiers**, e.g. ORCID, DOI, ROR

# Making your data FINDABLE

## Provide rich metadata
Example: Zenodo



Data set title

Author name & ORCID

Data set description

File names

Keywords

Persistent identifier (DOI)

License

# Making your data ACCESSIBLE

- Data are **retrievable by their identifier** using a standard protocol
- **Metadata are accessible**, even if the data are closed
  - ➢ Open Access
  - ➢ Embargoed Access
  - ➢ Restricted Access
  - ➢ No Access

# Making your data INTEROPERABLE

- Use a **formal machine-readable language** for the metadata

*Examples*
RDF
JSON

```
<rdf:Description>
  <dc:creator>Peter Noeller</dc:creator>
  <dc:title>Algebra</dc:title>
  <dc:subject>mathematics</dc:subject>
  <dc:date>2008-04-23</dc:date>
  <dc:language>EN</dc:language>
  <dc:description>
     An Introduction to Algebra
  </dc:description>
</rdf:Description>
```
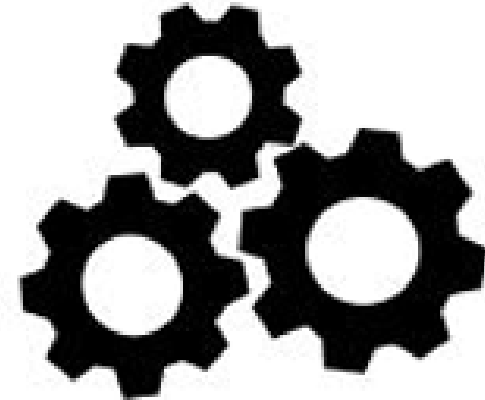
Interoperable

# Making your data INTEROPERABLE

- Use a **formal machine-readable language** for the metadata

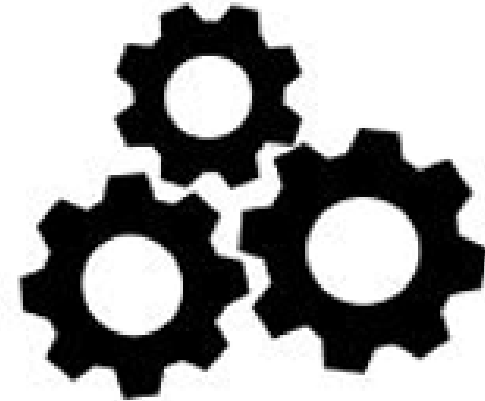- Use **open file formats** (if possible)

    *Examples*
    CSV for tabular data
    RTF for textual data
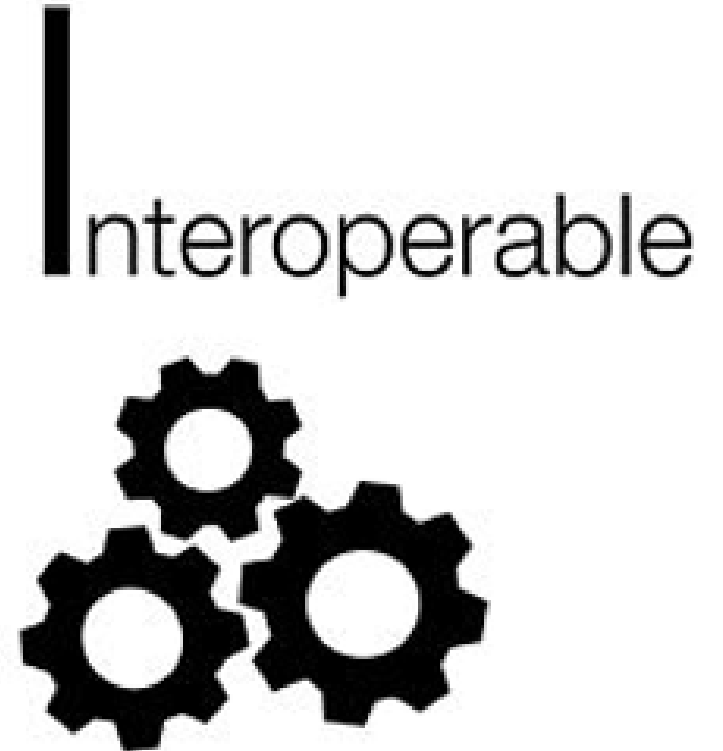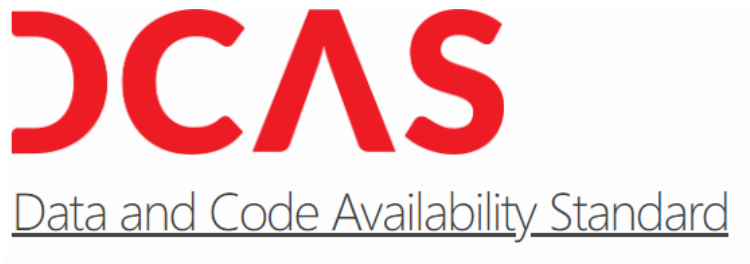    TIFF for images

# Making your data INTEROPERABLE

- Use a **formal machine-readable language** for the metadata
- Use **open file formats** (if possible)
- Use **common standards** (discipline)

Examples
➢ **Data and Code Availability Standard v1.0**

# Making your data REUSABLE

- Provide **rich accurate metadata**

- Give **detailed provenance information**

- Refer to **community standards**



https://social-science-data-editors.github.io/template_README/template-README.html

**CBS**

# Making your data REUSABLE

- Provide **rich accurate metadata**
- Give **detailed provenance information**
- Refer to **community standards**
- Use **clear licenses** for reuse, e.g. creative commons licenses

# Shades of FAIR

# Repositories support data FAIRification

## Findable:

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

✓ F2. data are described with rich metadata;

✓ F3. metadata clearly and explicitly include the identifier of the data it describes;

✓ F4. (meta)data are registered or indexed in a searchable resource;

## Interoperable:

✓ I1. (meta)data use a formal, accessible, shared, and broadly applicable language for knowledge representation.

✓ I2. (meta)data use vocabularies that follow FAIR principles;
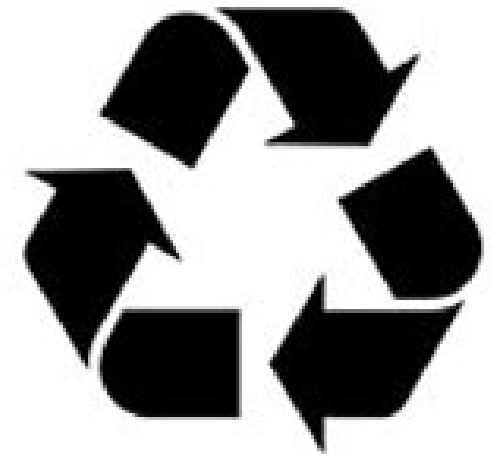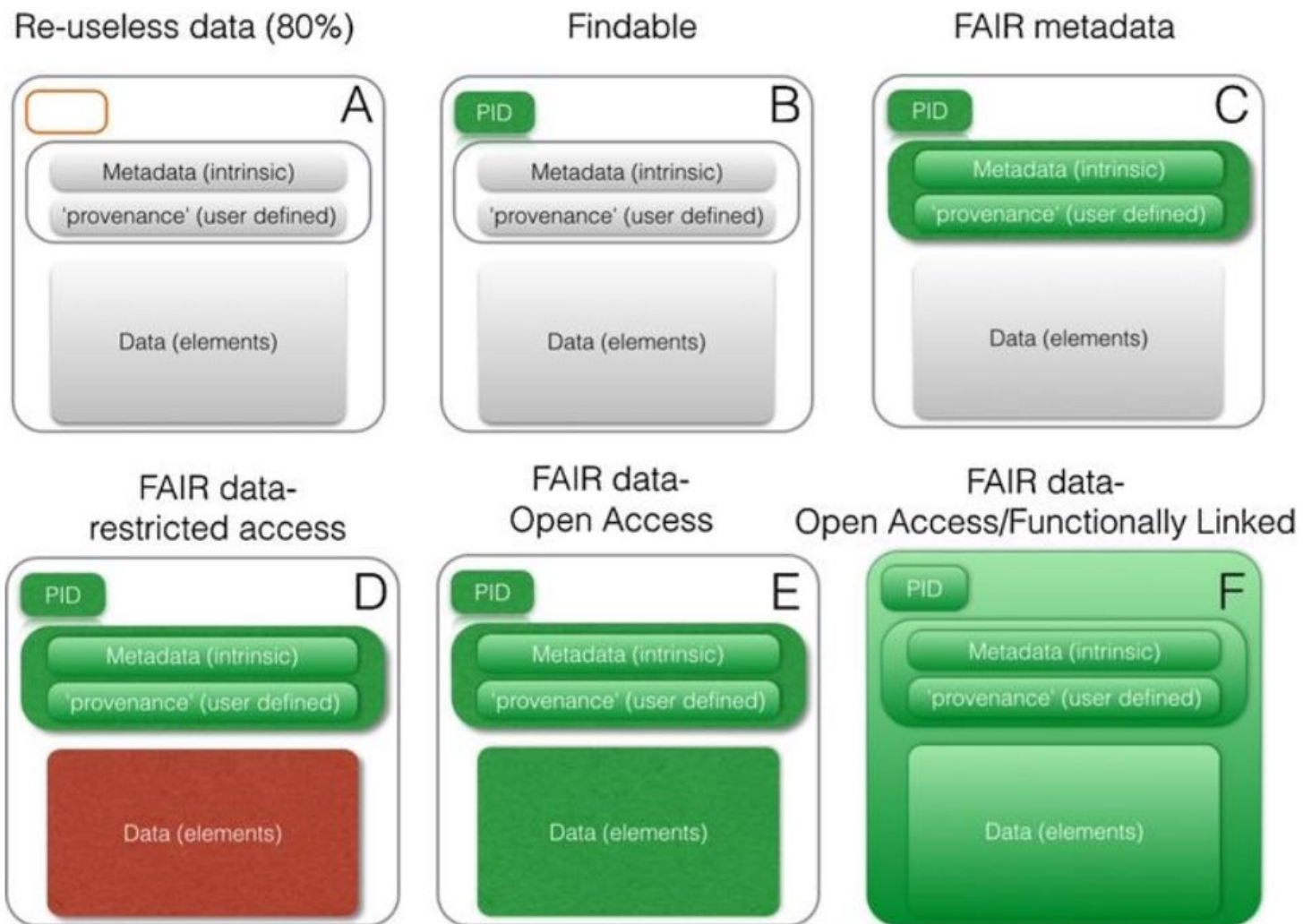
✓ I3. (meta)data include qualified references to other (meta)data;

## 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;

## Reusable:

✓ R1. (meta)data are richly described with 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 detailed provenance;

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

CBS

# Register data are FAIR by default

**Findable**

- Rich, standardized metadata
- Variable lists

**Accessible**

- Accessible (under an authorization)

**Access to data**

Access to data under the Research Scheme

**Interoperable**

- With other register data
- With external data, if linked by key variable

**Reusable**

- Reusability is at the core of DST

**STATISTICS DENMARK**

Burmeister, N. B. (2022). Arbejdsgruppe A - Use Case – FAIR Access to Register Data. Zenodo. https://doi.org/10.5281/zenodo.7437332

# Register data are FAIR by default – with limitations

**F**indable

- Easily findable - for experts only
- You don't always find what you need

**A**ccessible

- Only Danish institutions can be authorized
- (Language barrier)

**I**nteroperable

- Limited interoperability with other register data from the Nordics

**R**eusable

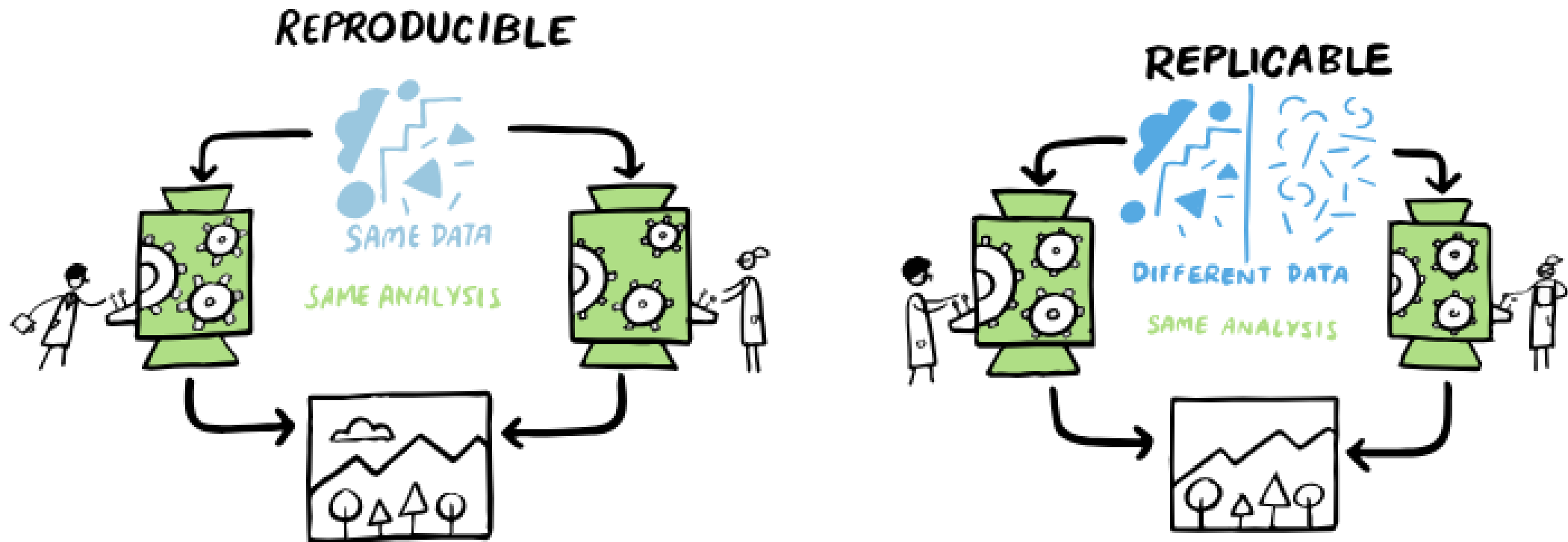- Changing variables (provenance!)
- Difficult access to previous projects

Burmeister, N. B. (2022). Arbejdsgruppe A - Use Case – FAIR Access to Register Data. Zenodo. https://doi.org/10.5281/zenodo.7437332

# HOW TO CREATE REPLICABLE DATA PACKAGES

# Reproducibility vs. replicability

# DCAS – the gold standard for replication packages

Data and Code Availability Standard (DCAS) requires the following elements:

1. Data availability statement
2. Raw data
3. Analysis data
4. Format
5. Metadata
6. Citation

CBS

# DCAS – the gold standard for replication packages

Data and Code Availability Standard (DCAS) requires the following elements:

7. Data transformation

8. Analysis

9. Format

CBS

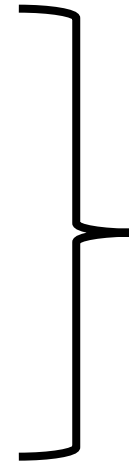# DCAS – the gold standard for replication packages

Data and Code Availability Standard (DCAS) requires the following elements:

10. Instruments

11. Ethics

12. Pre-registration

13. Documentation

14. Location

15. License

16. Omissions

CBS

# How to comply with DCAS when working with register data?

Create a package consisting of:

> ➢ Data citations
>
> ➢ Data access description
>
> ➢ Code
>
> ➢ Supplementary material

**FAIR protocol**

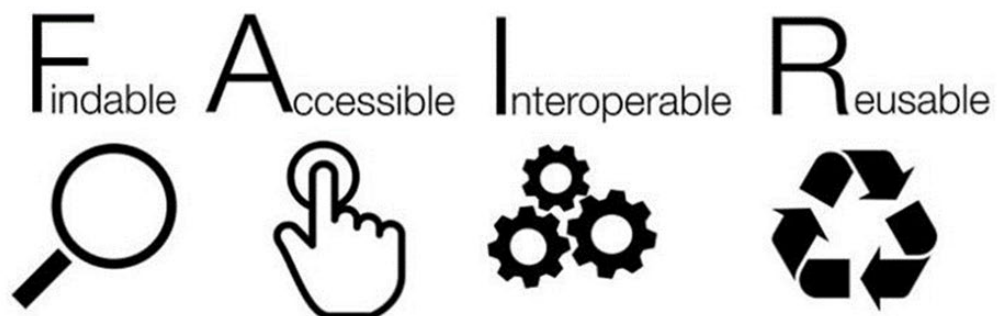# Squaring the circle: FAIR and replicable register data packages

**CAS**
Data and Code Availability Standard

CBS