# Data services

Christine Staiger, SURFsara

# Agenda

**1pm - 2.10pm**

      - Introduction (Christine)

      - GUI to archive, Cyberduck and FileZilla (Narges)
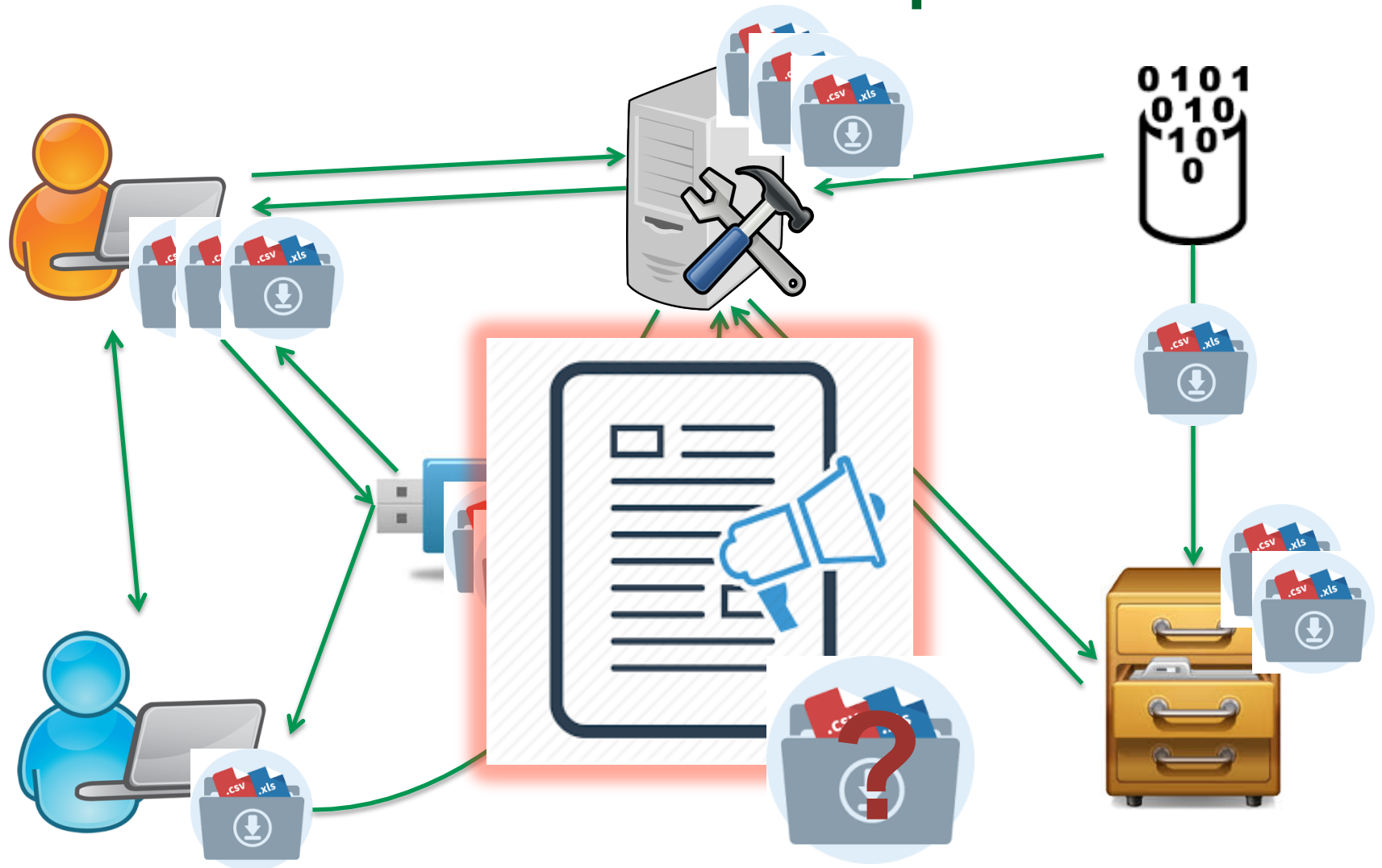
      - **Coffee break**

**2.10pm - 4.10pm**

      - scp + rsync (Jeroen)

      - Tarring data with tar and dmftar (Hans)

      - **Coffee break**

**4.10pm - 5pm**

      - Wrap-up (Christine)

      - Final Challenge:

            Putting it all together in a python workflow

# Data – where is the problem?

# Data services in the Netherlands – there are a lot of solutions

# Data services at SURFsara

# The data Life cycle



**CREATING DATA:** designing, planning consent, collection and management, capturing and creating metadata

**PROCESSING DATA:** entering, transcribing, checking & validating, anonymising and describing

**RE-USING DATA:** for follow-ups, new research, research reviews, scrutinising, teaching & learning

**ANALYSING DATA:** interpreting, deriving, producing outputs & publishing, preparing for sharing

**ACCESS TO DATA:** distributing, sharing, controlling access, promoting

**PRESERVING DATA:** migrating, backing-up, storing, creating metadata and documentation, archiving

CREATING DATA

PROCESSING DATA

RE-USING DATA

TRUST

GIVING ACCESS TO DATA

ANALYSING DATA

PRESERVING DATA
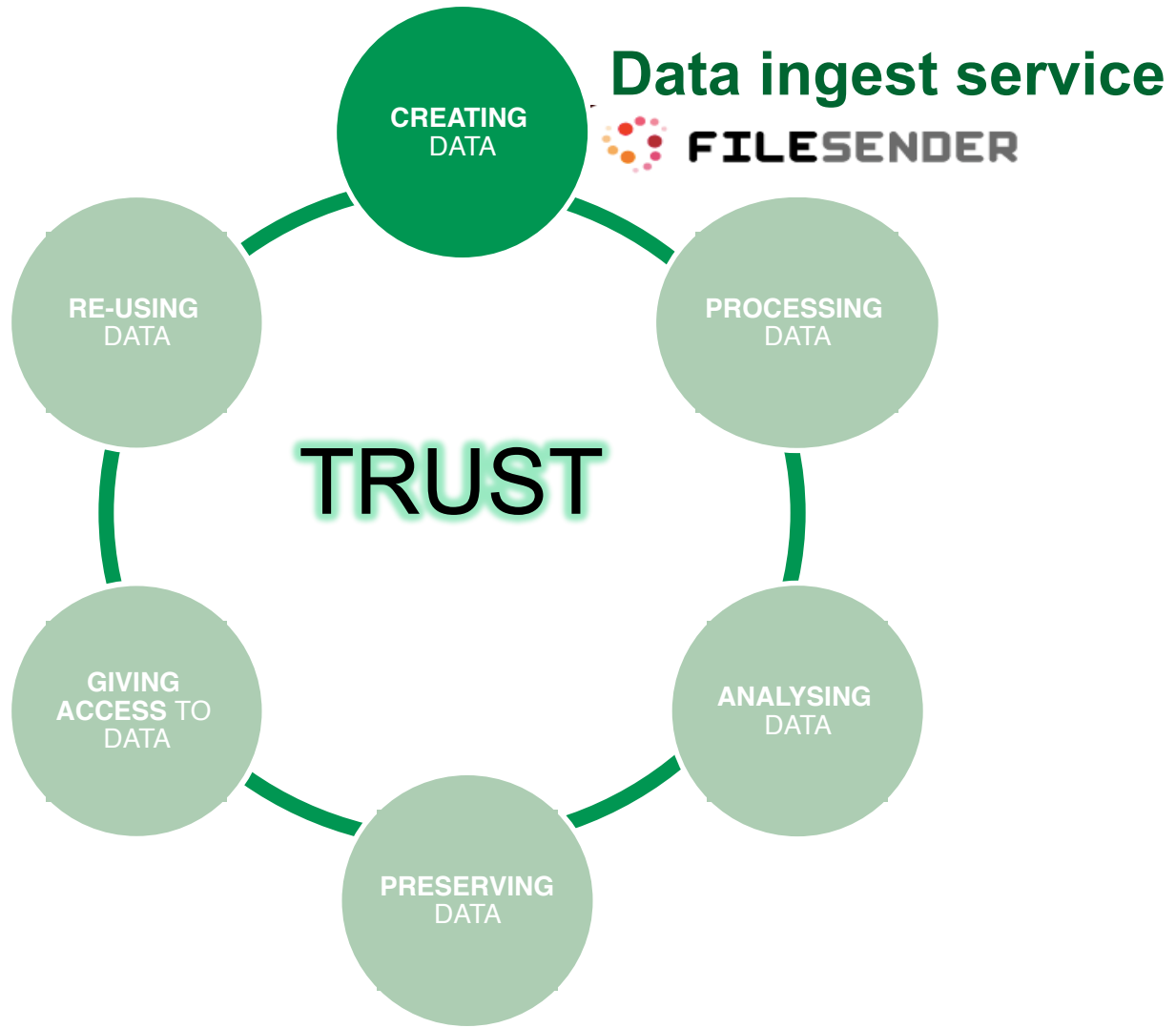
# The researchers' needs

- **Store** data during research
- **Share** data during and after research
- **Archive** data
- **Synchronise** data across different locations

- **Link** publication to processed and raw data
- **Publish** data
- **Find** data and **make data findable** by others

- Data **transfers**
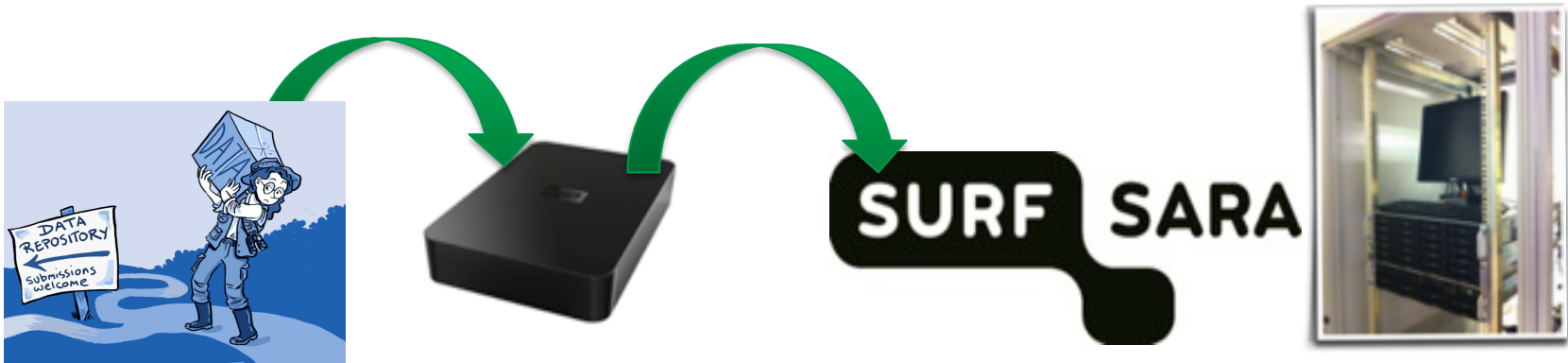- Data **provenance**: what happened with the data

- …

# The data Life cycle



TRUST

**CREATING DATA**

**Data ingest service**
FILESENDER

**PROCESSING DATA**

BeeHub
SURF DRIVE

**ANALYSING DATA**

**PRESERVING DATA**

**Archive**
SURF SARA

**GIVING ACCESS TO DATA**

BeeHub
SURF DRIVE
FILESENDER

**RE-USING DATA**

ePIC
Persistent Identifiers for eResearch
Handle.Net®

SURF SARA

# The data Life cycle

# Data ingest service

- Data often resides on external storage media, USB sticks, external hard drives
- Slow or no internet connection

- Easy way to upload large data from disk to SURFsara facilities
- Upload data from 45 disks in parallel

# FILESENDER

- Trusted community service
- Transferring BIG files from person to person
- File Transport service not File Storage (!)
- Simple interface
- Up to 200GB per file
- Option to apply end-to-end encryption (250MB browser limit)
- Vouchers for guest usage

# The data Life cycle

# Data services – what goes where?

# SURFdrive

- Trusted community cloud for personal storage
- Sharing smaller data files

- Collaboration between SURFsara, SURFnet and Dutch universities
- Specifications and service determined by end-users (universities)

- 100 GB storage capacity per user
- Based on ownCloud, synchronises with local storage



**SURF DRIVE**

# Beehub



- Sharing large data
- Mountable via webdav (no synchronisation)
- Capacity: 80 TB

- Temporary storage
- 100GB/user
- You can apply for more via SURF e-infrastructure grant

# The data Life cycle
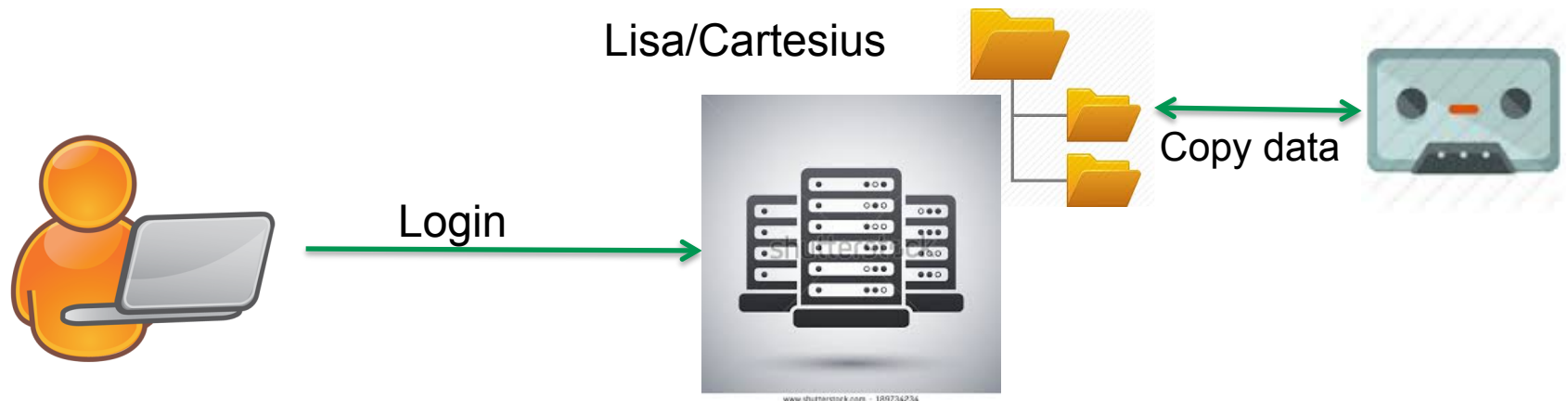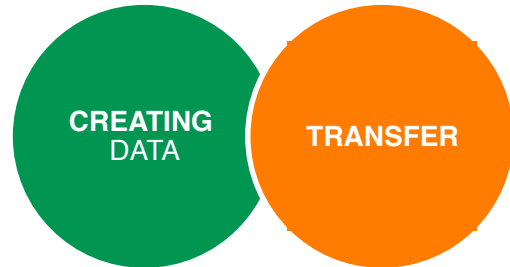
# Data services – what goes where?

# The Archive – Long-term storage

- Long-term storage of big data
- Storage medium: Tape
  → high latency

- Powerful transfer protocols:
  - gridfTP
  - rsync
  - scp

Lisa/Cartesius

Login

Copy data

- Easy access from HPC services lisa and cartesius via NFS mounts → use archive as yet another directory

- Access: NWO grant or SURF e-infrastructure grant

# Data services – what goes where?
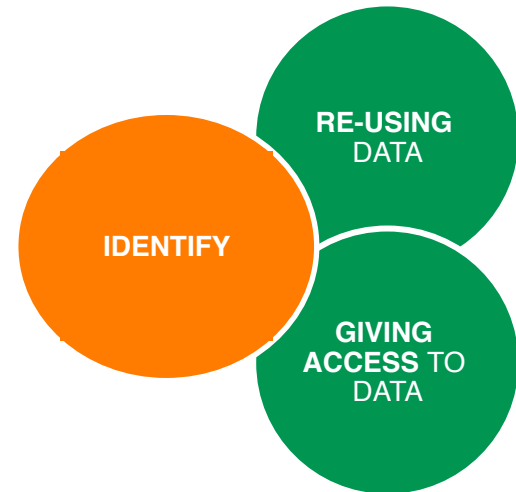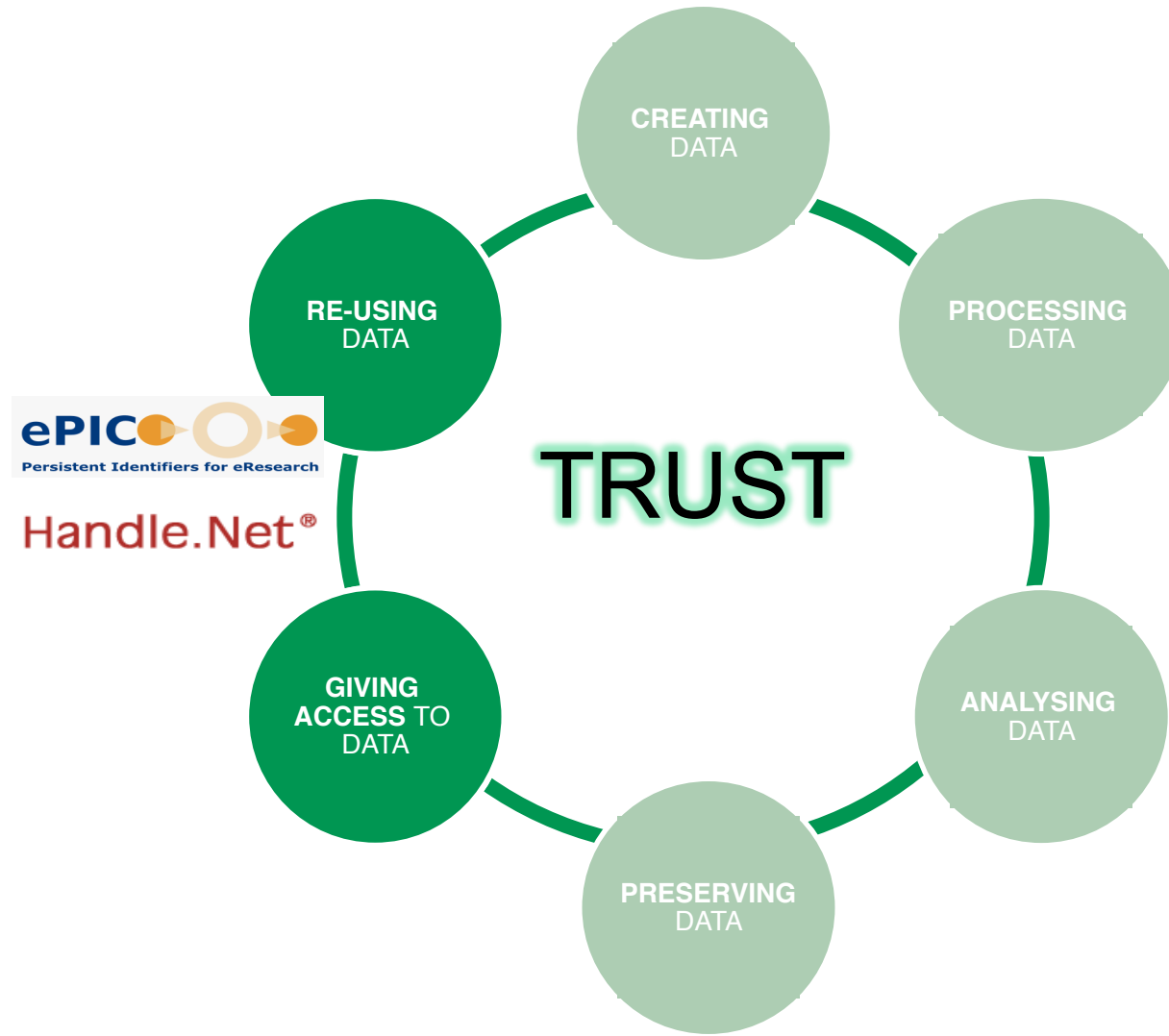


**CREATING DATA**

**TRANSFER**

**Data ingest service**

FILESENDER

doi

ePIC
**Persistent Identifiers for eResearch**

Handle.Net®

**IDENTIFY**

**RE-USING DATA**

**GIVING ACCESS TO DATA**

SURF SARA

# The data Life cycle

# PIDs

- PIDs (Persistent Identifiers) are
  - Pointers to resources like files, folders, webpages, real world objects
  - Globally unique
  - Resolvable via http

- Example resolvers: https://dx.doi.org/ and http://hdl.handle.net/

10.2307/748467

# PIDs – Handle, EPIC and DOIs

- Handle
  - Technology to create, store and update PIDs
  - Infrastructure and technology to resolve PIDs
  - Example: https://dx.doi.org/ and http://hdl.handle.net/

- EPIC (European Persistent Identifier Consortium)
  - Maintaining reliable joint PID service for storing
  - Mirroring service, distributed resolving
  - Employing Handle technology

- **http://www.ncdd.nl/pid-wijzer/**

# The data Life cycle