# DATA MANAGEMENT IN COMPBIOMED MOVING TOWARDS FAIR DATA

NARGES ZARRABI

**Machine Learning meets Modelling and Simulation Methods**

**17th March 2020**

SURF

# What is SURF?

SARA **(1971)** ➡ SURFsara (2013) ➡ **SURF (2020)**

**SURF is the collaborative organisation for ICT** in Dutch education and research

Driving innovation together

SURF

# Fields of work



**Education**

Flexible education

Diverse learning resources

Using study data



**Research**

Unlimited access

World-class facilities

Stimulating Open Science



**Cooperative facilities**

On campus

Security in the digital world

User-centred

SURF

# The Dutch National Supercomputer Cartesius

- Total cores:          47767 CPU + 132 GPU

- Total memory:     117 TB

- Peak performance:    1.843 Pflop/sec

- Disk space:         180 TB home,
  7.7 PB project/scratch

- Operating system:    BullX (GNU/Linux)

- Network:            Mellanox InfiniBand
  56 GBps bandwidth,
  3µs latency
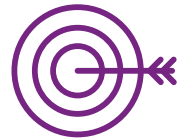
Top 500 largest supercomputers

- 2014: #45

- 2018: #360
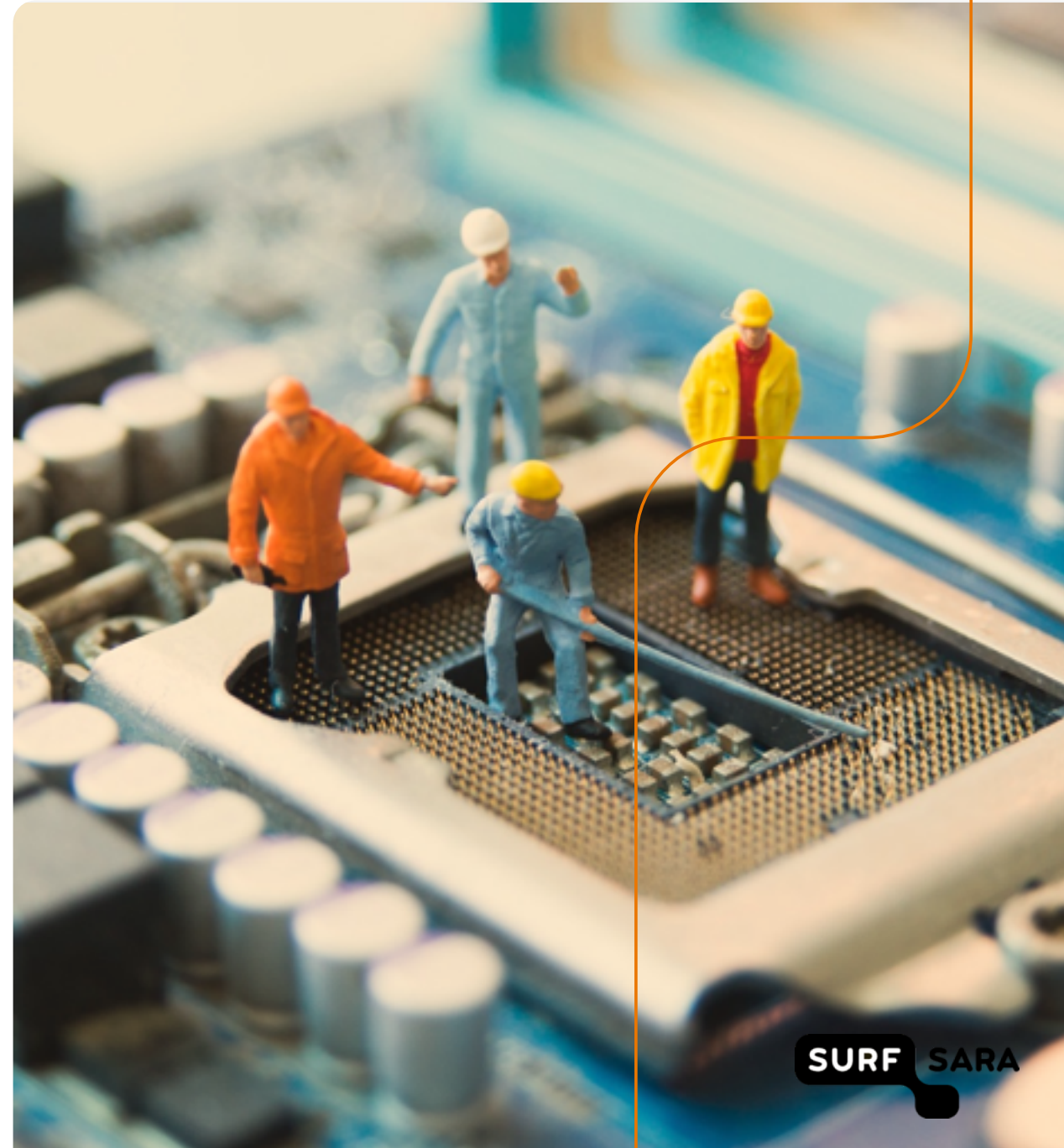
SURF

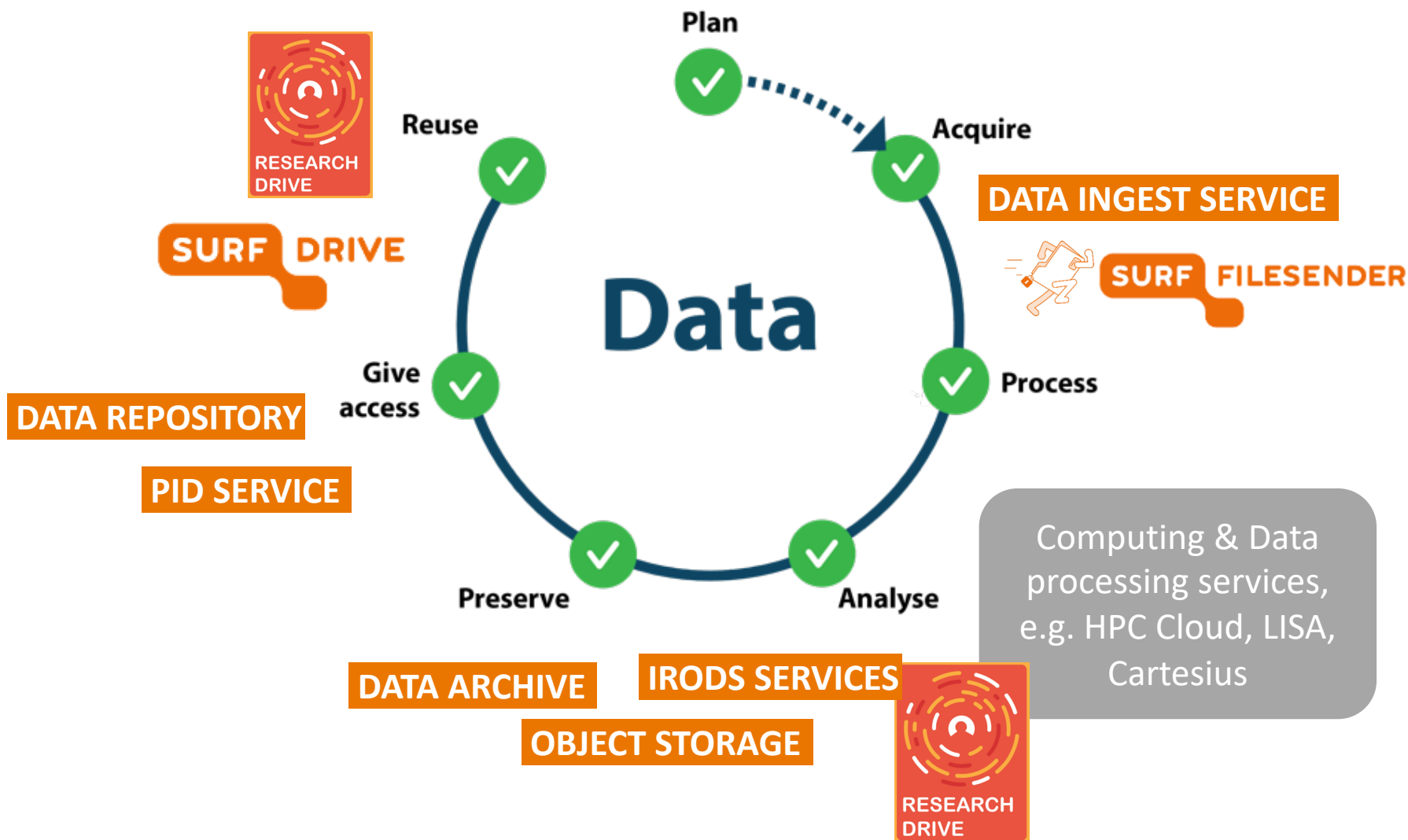# SURF is more..
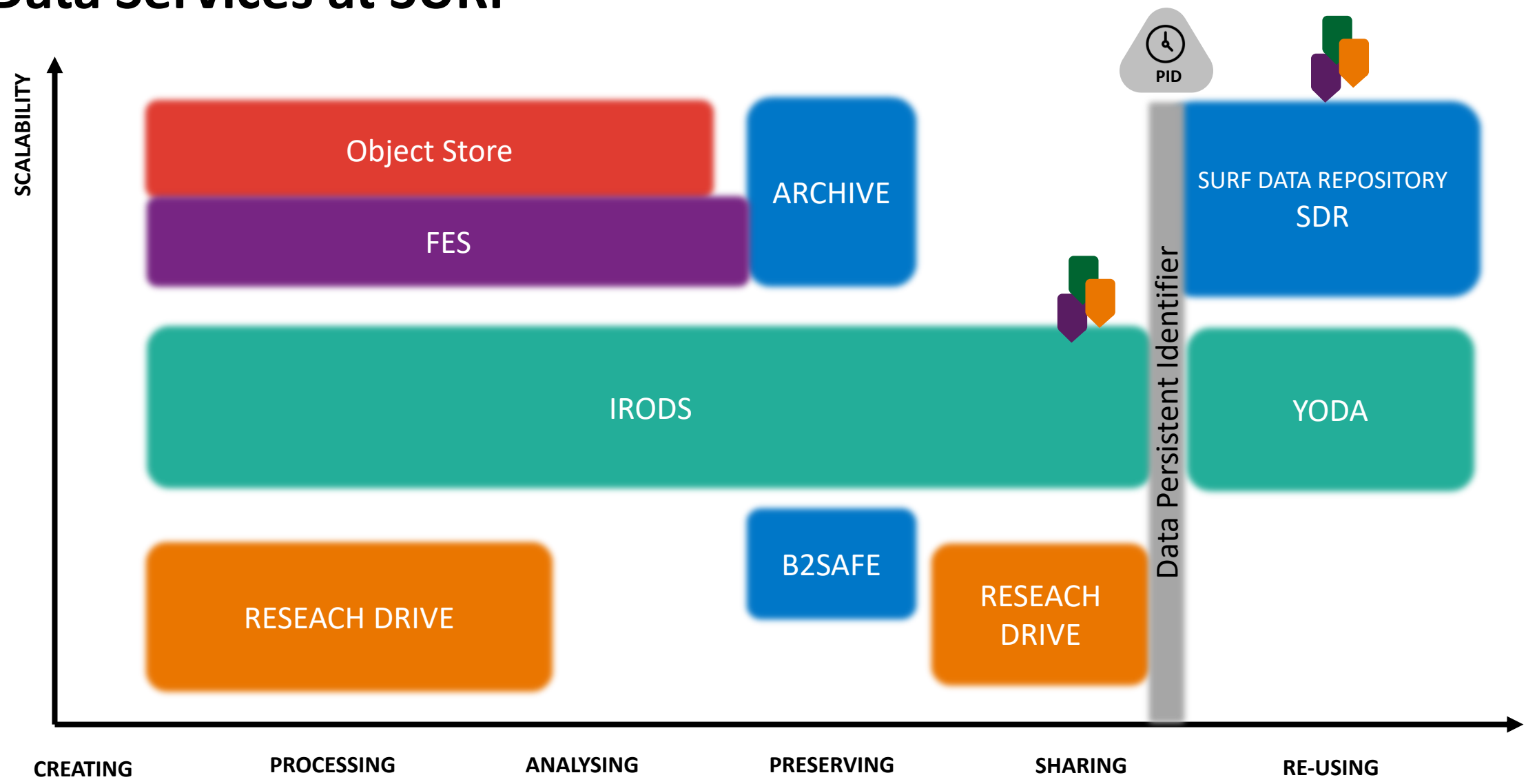# than just big systems

Consultancy

Training

Knowledge Exchange

SURF SARA

# Data Services at SURF



For more detail, please see https://www.surf.nl/diensten-en-producten/categorie/datadiensten

# Data Services at SURF

# Data Services Projects

# Data Requirements in Reseach Communities

- **More efficient data access, sharing and transfer**
    - *Intensive data-sharing and transfer*
    - *Restricted data-sharing and transfer*
- **Preserving research data**
    - *Storage, backup and archiving large data, synchronizing data over distributed places*
    - *data provenance*
- **Accessible research Data**
    - *Making data accessible to research communities, PIDs*
    - *Publishing data with domain specific metadata*
    - *Linking published data to processed and raw data*
- **Findable research data**
    - *A major challenges scientific communities is to discover data from research data collections and repositories*

SURF

# What is… FAIR ?

## 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;

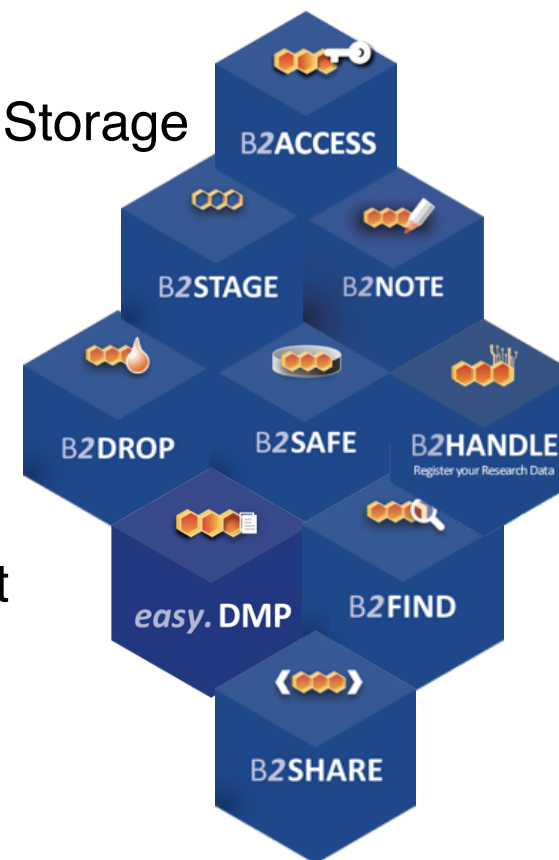F4. (meta)data are registered or indexed in a searchable resource;

## 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;

**EUDAT**

**SURF**

# EUDAT B2Service Suite

- **B2ACCESS** – Authentication and Authorisation
- **B2DROP** – Data Workspace
- **B2SAFE** – Distributed, Secure Policy Based Data Storage
- **B2SHARE** – Searchable Data Repository
- **B2STAGE** – High Performance Data Movement
- **B2FIND** – Searchable Metadata Aggregator
- **B2HANDLE** – Persistent Identifier Provider
- **B2NOTE** – Semantic Metadata Annotation
- **easy.DMP** – Data Management Planning Assistant
- **Gitlab** – Git repository and collaborative software development platform

# Use case: Workflow using Alya Application

- **Step 1: Data creation and transfer:** The raw data is collected at ESRF in France. The data is being stored locally on tapes. Currently, a copy of the data is transferred to BSC.

- **Step 2: Data pre-processing:** In BSC, researchers pre-process the data which includes manual and automated steps for image stitching, segmentation and meshing.

- **Step 3: Data replication:** The preprocessed data needs to be replicated from BSC to SURFsara and EPCC. The replicated data will then be used to run simulations on the supercomputers in these sites.

- **Who**
  - Community Data Managers
  - 'Sophisticated' Organizations

- **What**
  - Provide an abstraction layer which virtualizes large-scale data resources
  - Guard against data loss in long-term **archiving and preservation**
  - **Optimize access** for users from different **regions** and to **computing** resources
  - Data management on basis of **policies**

- **Why**
  - Performance
  - Replication between trusted sites
  - Data Preservation

# Data Replication Pilot (CompBioMed1)

**Data replication pilot**

- Safe data replication data preservation

- Allocation of PIDs to replicated data

- Facilitate large data transfer

- Bring data close to compute

- Scale-up compute power



EPCC

SURF

**Federated network of B2SAFE endpoints**

BSC

**HPC Centers:** BSC, SURFsara, EPCC
**Resources:** allocation of at least 24 TB storage at each of the HPC centers

# Future work pilot (CompBioMed2)

- Allocation of resources and replication of the real data (24 TB per Centre)

- Extend the network of B2SAFE endpoints including more HPC centres

- Try other replication scenarios

- Integration B2SAFE and B2SHARE

**EUDAT Data Repository for publishing data**

**Who**

- Small to Medium Teams

**What**

- **Store** data (incl. software) and add domain meta data
- **Share** registered research data worldwide
- **Preserve** (small-scale) research data for long-term

**Why**

- Register Data for Publications (FAIR)
- Make known to wider community

https://b2share.eudat.eu/

# CompBioMed Community in B2SHARE

# Thank you!

Thanks to:

- SURF Data Preservation Services team

- SURF Data Management Services team

- Projects: