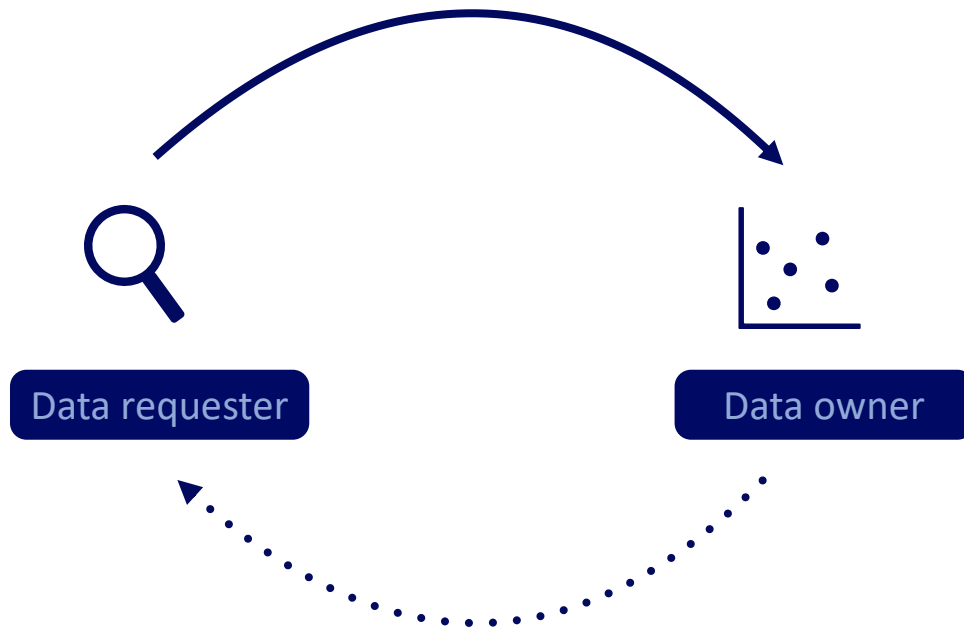# DATA EXCHANGE DEMO

Share data while retaining control and confidentiality of your data

**SURF**

# Gains and difficulties of sharing confidential data

**+** Access to non-public data.

**+** Potential new research and collaborations.

**−** More work to manage confidential data.

Data requester

Data owner

**∼** Possible to gain new insights.

**−** Risks on privacy and security.

**−** Additional work without direct return on investments (ROI).

⚠ Gain is usually with the data requester, burden is with the data provider

SURF

# Willingness to share data

**ROI** + **Trust**

Return on Investment (ROI) is determined by the balance between effort it takes to share data, and the gain received by sharing data

Trust is determined by the balance between the risks (due to privacy or competition), and the control (due to verification and security) of sharing and usage of data

Gains         Effort

Control         Risk

Return on investment

Trust

SURF

# Type of Data Owners

| Data aggregators | Hospitals + medical institutions | Onderzoekers + universiteiten | Bedrijven |
|---|---|---|---|
| Health care (Palga, NZa) Social-economic (CBS, municipalities) | Hospital (AMC, vuMC, St. Antonius) Insurance companies (Zilveren Kruis) | Universities (Twente, Wageningen, Groningen) Researchers | Friesland-Campina, Elsevier |

**Privacy sensitive**

**Competitive data**

SURF

# Example: Find the average income



**Run #1**

- 21 people

- Algorithm verified

- Outcome guaranteed not to be traceable to individual people

**Run #2**

- 22 people (same 21 and 1 other)

- Algorithm verified

- Outcome guaranteed not to be traceable to individual people

*Even if individual runs are fine, combining two runs may reveal confidential data*

# Different Methods to Ease Data Sharing

### Agreements

- Stipulation of what can/cannot be done
- Signing of contract or NDA
- Dispute resolution process

### Registration

- Authentication
- Verification of credential
- Reputation score
- Policy framework
- Audit trails

### Pseudonymization

- Filtering (on records)
- Pruning (on properties)
- Aggregation (combine records)
- Make coarse grained buckets
- Slight alteration of data
- One-way hashing
- One-time identifiers
- Synthetic data (mix records)

### Data Vault

- Data source retains control
- Delegate permissions
- No central data lake
- Data marketplace

### Secure Containers

- Bring algorithm to data
- At Trusted third party or at data provider
- Share output instead of data

### Secure Computing

- Secure multi-party computation
- Homomorphic encryption
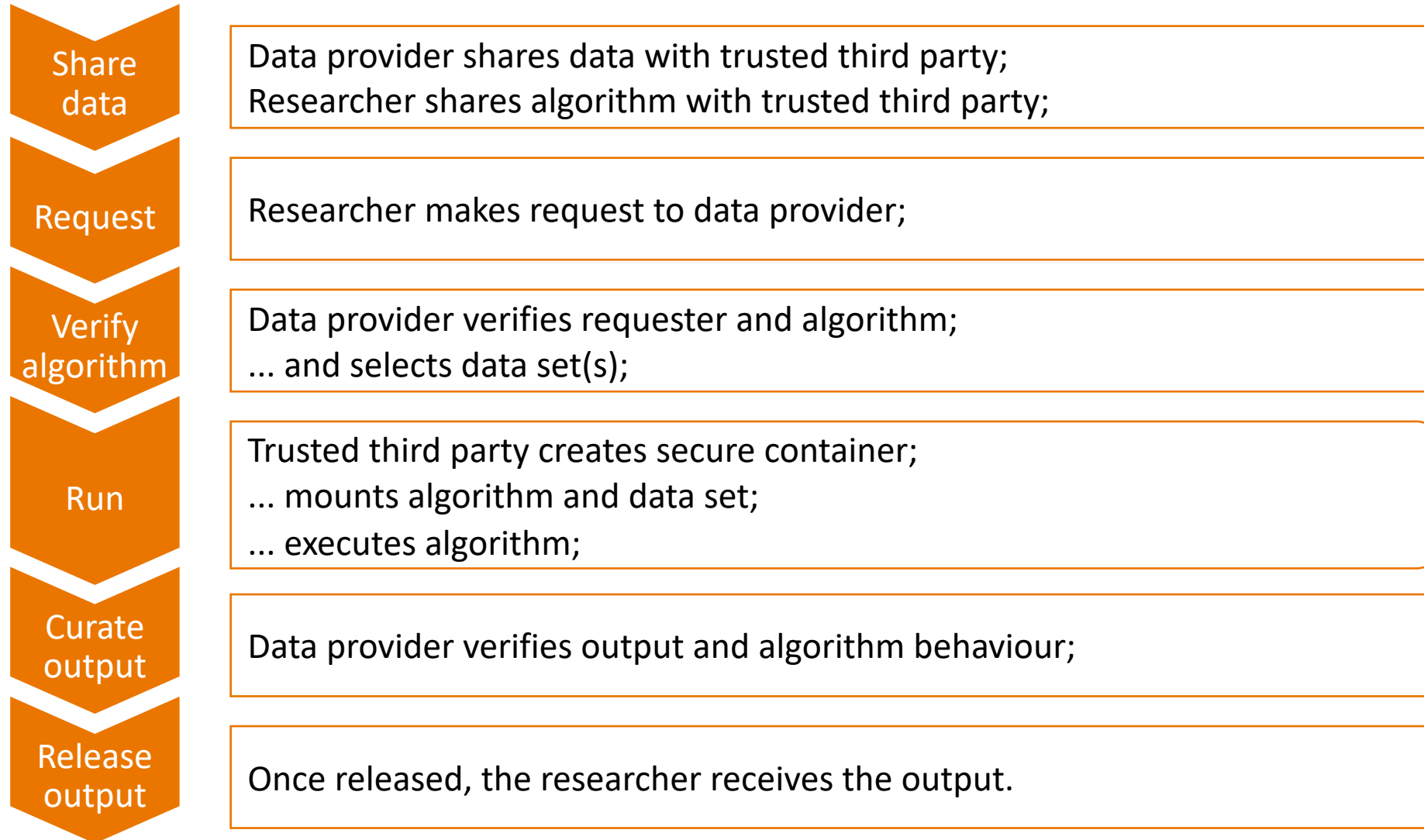- Garbled Circuits
- Zero-knowledge proof

SURF

# Data Exchange

| | | | |
|---|---|---|---|
| **VISION** | \multicolumn{3}{c}{**Realize a platform where data can easily be shared, while retaining control and confidentiality of the data**} | | |
| **TARGET GROUP** | **NEEDS** | **PRODUCT** | **BUSINESS GOALS** |
| Data providers with confidential data. E.g.<br>• Companies;<br>• Academic hospitals.<br><br>Researchers who like to use data from other organizations for a specific purpose. | Data providers like to share data, while<br>• retain control who can use the data for what purpose;<br>• adhere to legal limitations of processing data.<br><br>Data consumers (researchers) don't want to be limited to public datasets. | Proof of concept (demonstration).<br><br>Secure environment at trusted third party.<br><br>Performs calculations on data on behalf of a researcher, with explicit consent from the data provider. | Facilitate open science<br><br>Researchers make more use of data sources.<br><br>Provide a easy-to-use and trusted solution for both parties, data providers and researchers |

SURF

# Collaborating without direct Sharing Data



Data Provider

Data

Trusted Third Party

Secure container

Result

Curation of result

Code +Data

Researcher (Algorithm Provider)

Result

SURF

# Workflow

**Share data**
Data provider shares data with trusted third party;
Researcher shares algorithm with trusted third party;

**Request**
Researcher makes request to data provider;

**Verify algorithm**
Data provider verifies requester and algorithm;
... and selects data set(s);

**Run**
Trusted third party creates secure container;
... mounts algorithm and data set;
... executes algorithm;

**Curate output**
Data provider verifies output and algorithm behaviour;

**Release output**
Once released, the researcher receives the output.

SURF

# Permission Models

| One-off permission | Trust a researcher | Run on a data stream |
| --- | --- | --- |
| The data provider permits a researcher to **run** a specific algorithm **once** on a specific dataset. | The data provider permits a researcher to **run any algorithm on a specific dataset**.<br><br>The permission can be revoked at any time.<br>Example use cases:<br>• the data provider trust the researcher to always write benevolent code<br>• the researchers wants to tweak the algorithm, and run it on a sample dataset every time. | The data provider permits a researcher to **run a specific algorithm on any data set in a selected folder**. Every time a new dataset is added to the folder, the algorithm is automatically run.<br><br>The permission can be revoked at any time, but is also automatically revoked as soon as a change to the shared algorithm is detected. |

Currently supported permission models

SURF

# Implementation (Proof of Concept)

- Working prototype

- Non-production (not scalable nor fast, not rigorously tested)

- Data stored at ResearchDrive (OwnCloud implementation at SURF for researchers)

- Data sharing: https://dataexchange.surfsara.nl/
  (simple password to emphasis it is a demonstration only: demo / dex)

- **Goal is to understand user requirements**

Mike Kotsur    Rienk Koenders    Sijmen Schoon    Tijs Teulings    Sander van Wickeren    Axel Berg    Hylke Koers    Gerben van Malenstein    Freek Dijkstra

bit

# Technical Implementation of the prototype



**External integrations**

**Internal Components**

# Risks and Mitigations

| Risk | Mitigation |
|------|------------|
| Data is leaked to outside world | Researcher can never view the raw data, only the result |
| Data is used in other ways than intended | Data provider can review algorithm |
| Algorithm is leaked to outside world | Algorithm is not reviewed by data provider, researcher is trusted to write benevolent code only * |
| Output contains confidential information | Data provider curates output before releasing it to researcher |
| Malicious algorithm tries to copy data to remote server | No network access is allowed in secure container |
| Malicious algorithm tries to embed data in output | Data provider can review algorithm |
| Algorithm is altered after it is shared | Permissions involving this algorithm are automatic revoked |
| Researcher can no longer be trusted | Permission can be revoked by data provider at any time |
| Trusted third party can no longer be trusted | Sharing of data to trusted third party can be revoked at any time |
| Data is corrupt or data provider can no longer be trusted | Researcher should look for other data sources |
| Data can't leave premises, not even to a trusted third party | Secure container can be run at premises of data provider * |

\* Not yet implemented in the prototype

SURF

**Data is shared with the Data exchange**

14

**Algorithm is shared with the Data exchange by researcher**

**Researcher makes a request to the data provider**

**Data provider reviews request and selects dataset**

**Trusted Third Party runs algorithm on dataset**

**Data provider reviews output**

**Researcher can see released output**

**Data provider can at any time withdraw permissions**
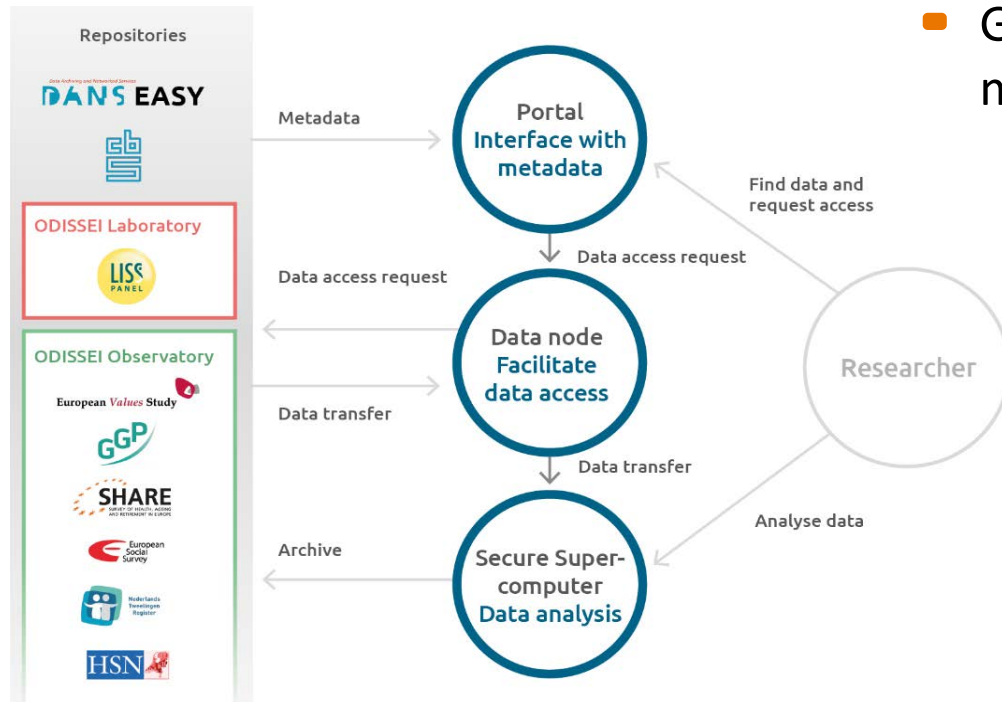
# Related Projects

## ODISSEI Secure Supercomputer (OSSC)

- In production

- Processes CBS micro-data on Cartesius

- Does pseudonymization as well



## AMdEX

- Collaboration of interested parties

- Initiated by Amsterdam Economic Board

- Goal is to build an infrastructure for multiple Data Marketplaces

# Partnership Questions

- Who may benefit from a data exchange?

    - Are there researchers that want to use confidential data?

    - Who are the data providers in this case?

    - Under what conditions would these data providers release their data?

- What should the role of SURF?

    - Service provider; software developer; community manager; …

- Should SURF turn this prototype into a pilot?

- Are there other projects we should collaborate with?

SURF

# Technical Questions

- Is a trusted third party the right approach?

- What is the trust relation?

  - Does the data provider trust the researcher?

  - Does the data provider trust the algorithm?

- More advances user scenarios (e.g. with 3 parties):

  - Patient trust a hospital with their data

  - Hospital trust a researcher with the patient data

  - What are the implications for the current demo with 2-part user-scenario?

  - Who gives what permissions, and is that a continuous permission? How to withdraw permissions?

SURF

# COLLABORATION WITHOUT SHARING DATA

**Freek Dijkstra**

**Freek.Dijkstra@surfsara.nl**

**www.surf.nl**

This presentation is available under the creative commons attribution 4.0 license

## Driving innovation together

SURF