

Better infrastructure
Better science
Better society



ODISSEI

Open Data Infrastructure for Social Science and Economic Innovations

The ODISSEI Secure Supercomputer

Annette Langedijk, Lucas van der Meer • ICTeSSH • 30 June 2020



ODISSEI

Open Data Infrastructure for Social Science and Economic Innovations

Who are we?



Annette Langedijk

Community
manager SSH SURF
Management Board
ODISSEI



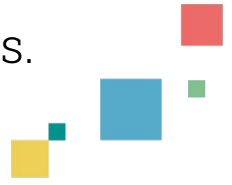
Lucas van der Meer

Operational manager
ODISSEI

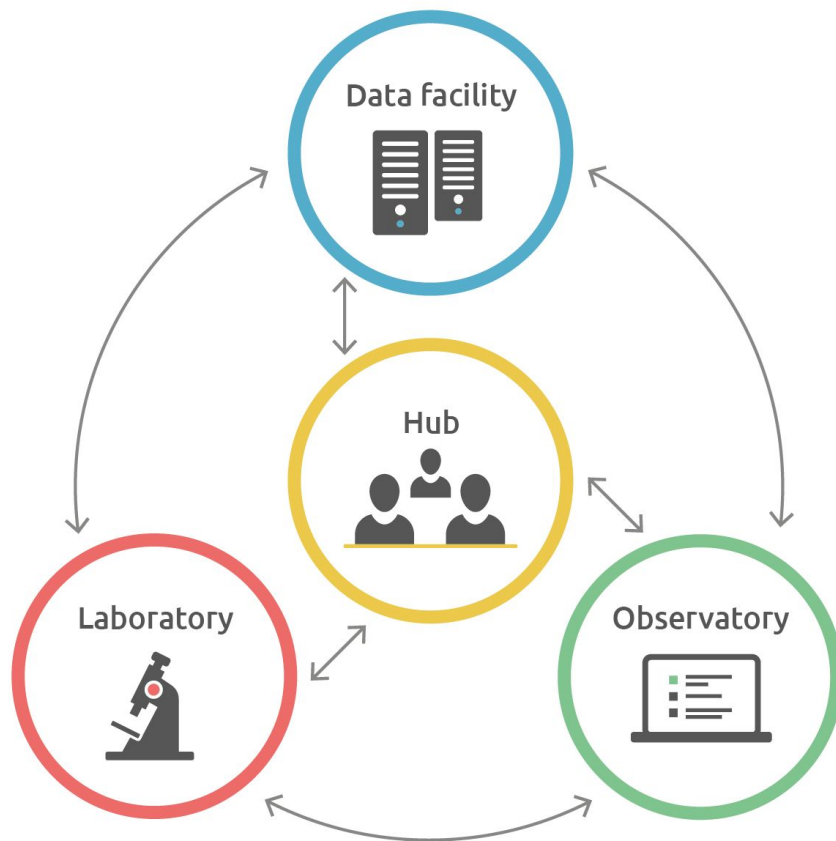


A unique non-profit collaborative consortium

ODISSEI creates a federated data infrastructure for the social and economic sciences in the Netherlands, on behalf of 40 member organisations.



The ODISSEI research infrastructure



From traditional 2D survey data...

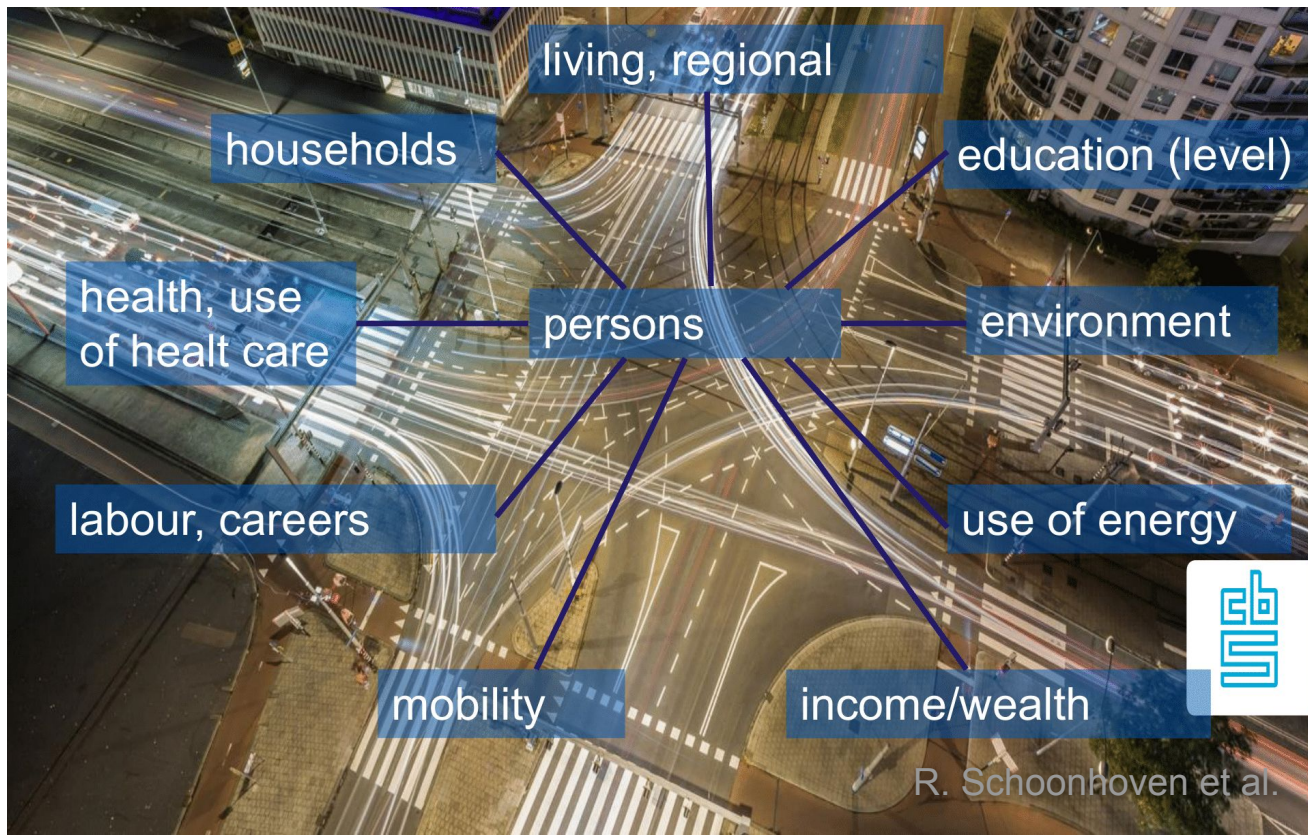
Person#	Sex	Occupation	...
454a87	M	Unemployed	
986c77	F	Plumber	
...			



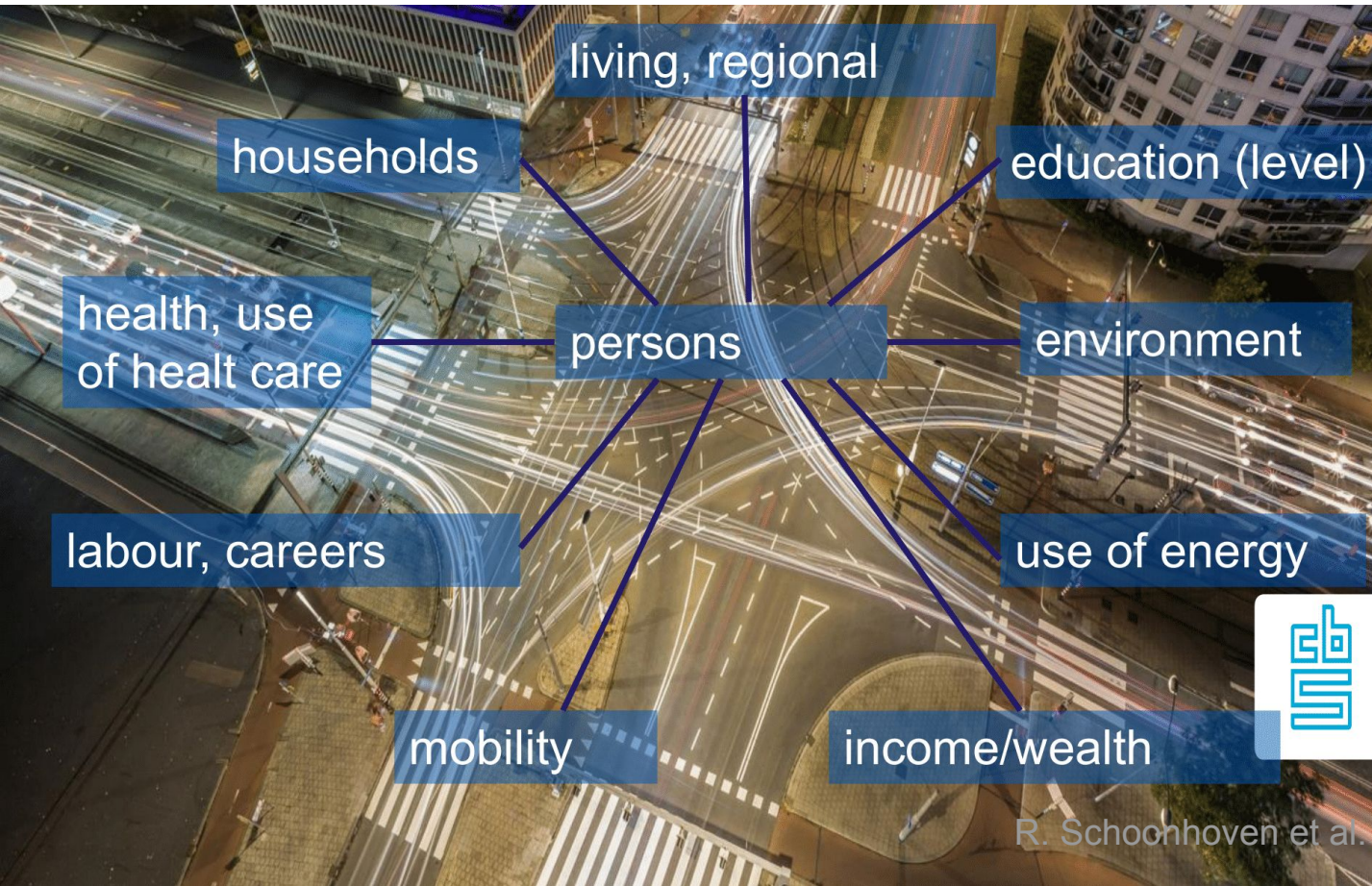
From traditional 2D survey data...

Person#	Sex	Occupation	...
454a87	M	Unemployed	
986c77	F	Plumber	
...			

...to rich multi-dimensional, linkable data



Microdata from Statistics Netherlands (CBS)



Microdata:

linkable data on individuals, companies, and addresses



Example Microdata research

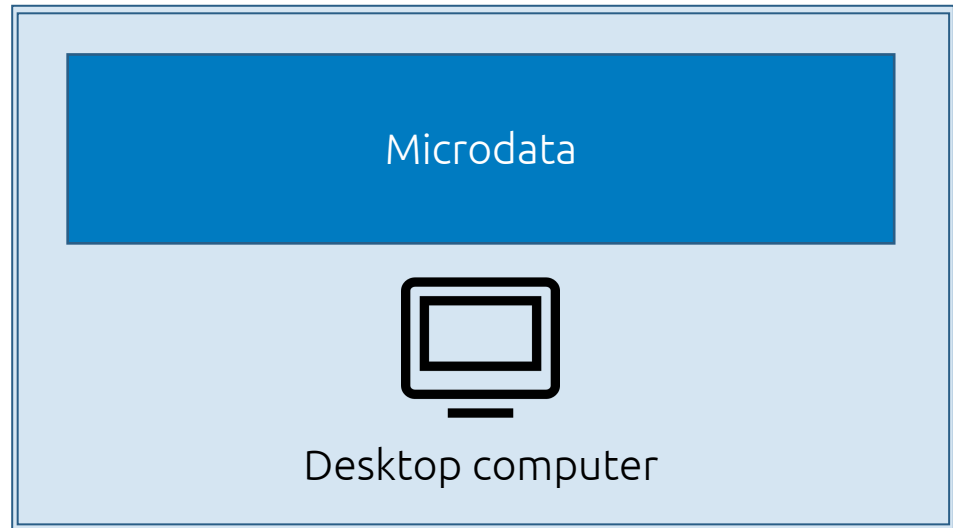
THE EFFECT OF EDUCATIONAL FAILURE ON MENTAL HEALTH

	2006	2007	2008	2009	2010	2011	2012	2013	2014	2015	2016	2017	2018
Costs				■	■	■	■	■	■	■	■	■	■
Medication	■	■	■	■	■	■	■	■	■	■	■	■	■
Diagnoses						■	■	■	■	■	■	■	■
Education			■	■	■	■	■	■	■	■	■	■	■

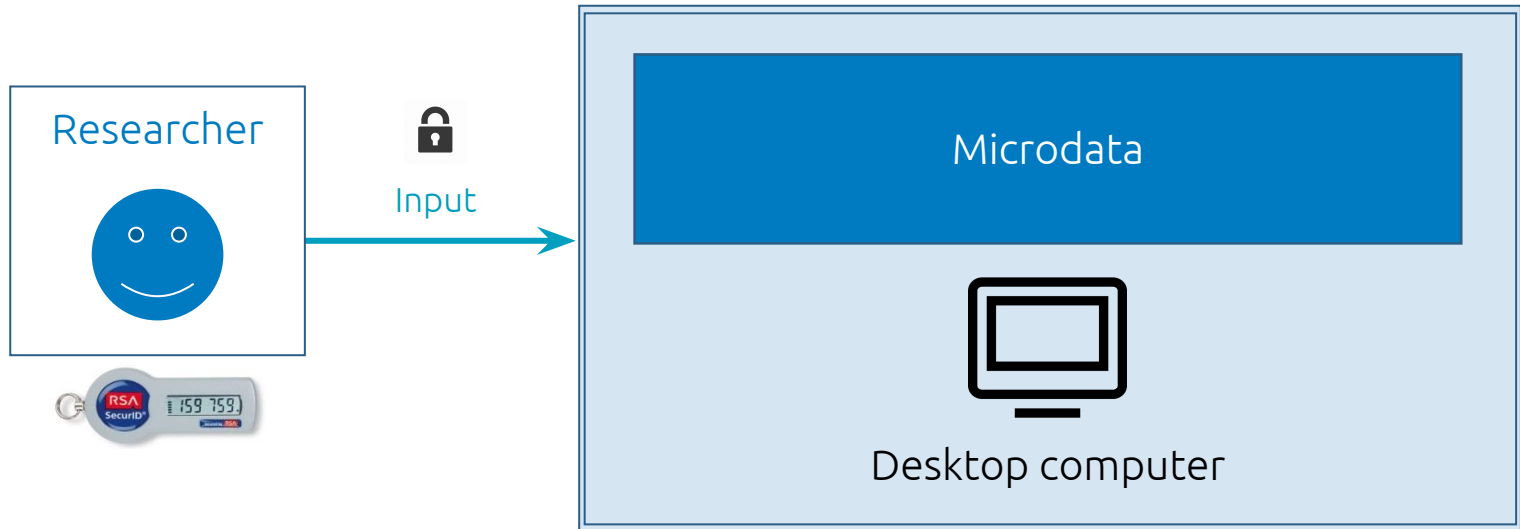
■ Adults + children ■ Only adults



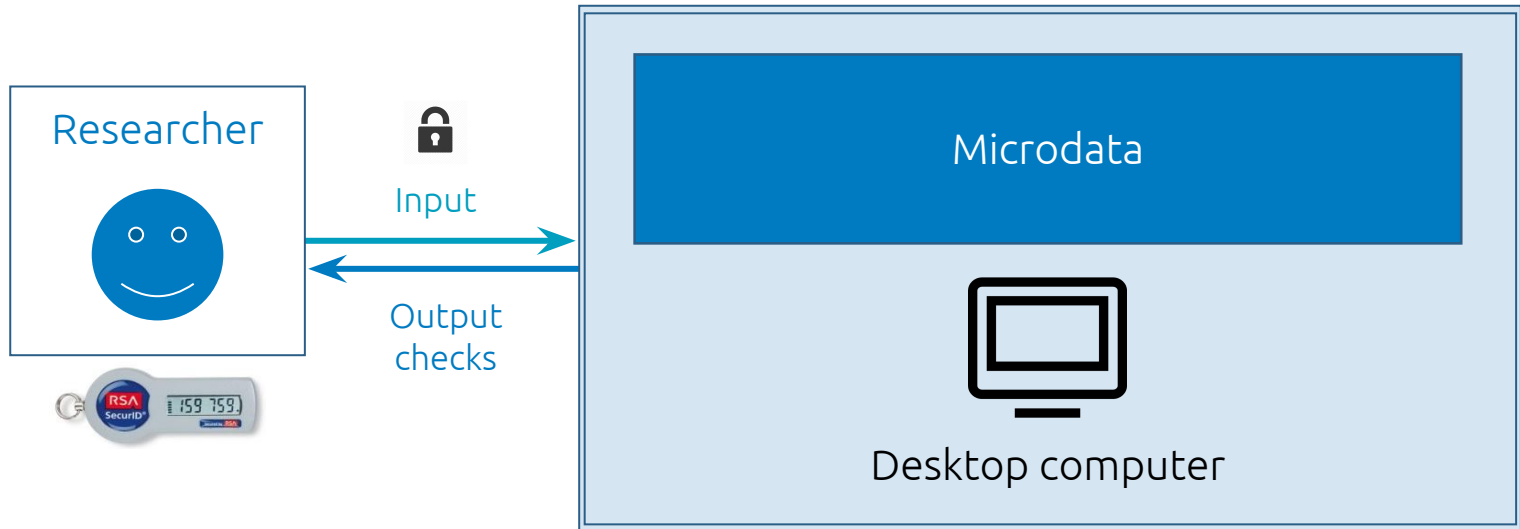
Statistics NL Remote Access environment



Statistics NL Remote Access environment



Statistics NL Remote Access environment



Challenge: billions of data points

585,000 populated cells
× 101 scales × 2 variables
× 15 years = 1.8 billion
data points

→ 4 months continuous
calculations

ODISSEI Secure Supercomputer



ODISSEI

Open Data Infrastructure for Social Science and Economic Innovations

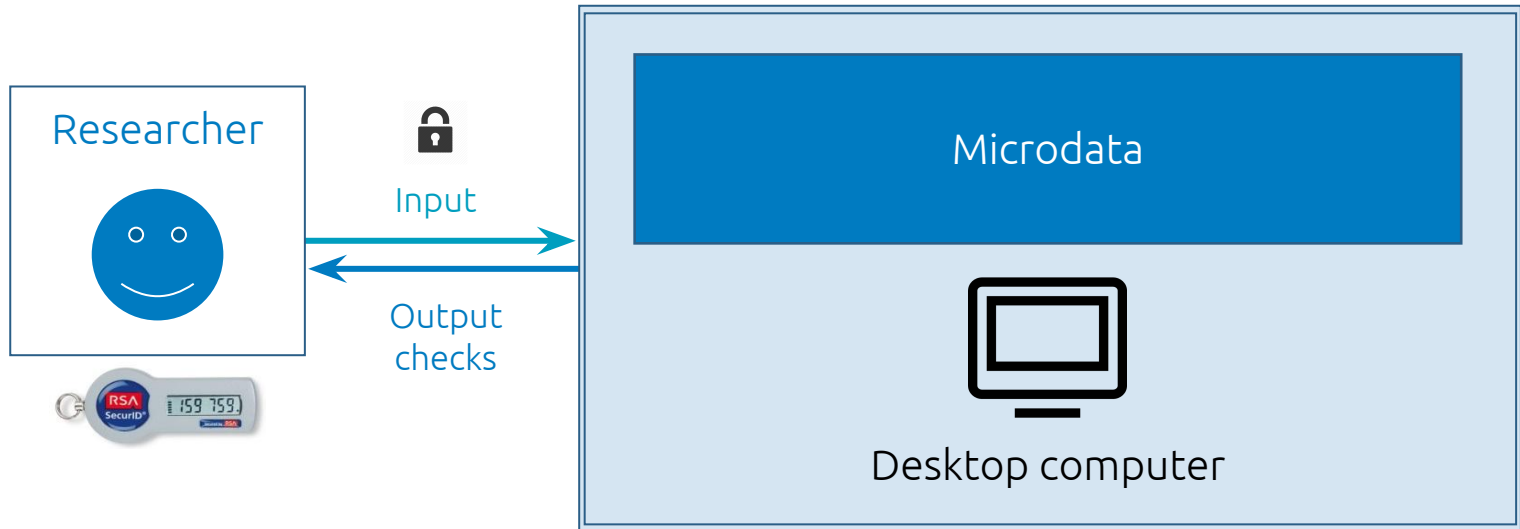
Cartesius, the Dutch National supercomputer

SURF

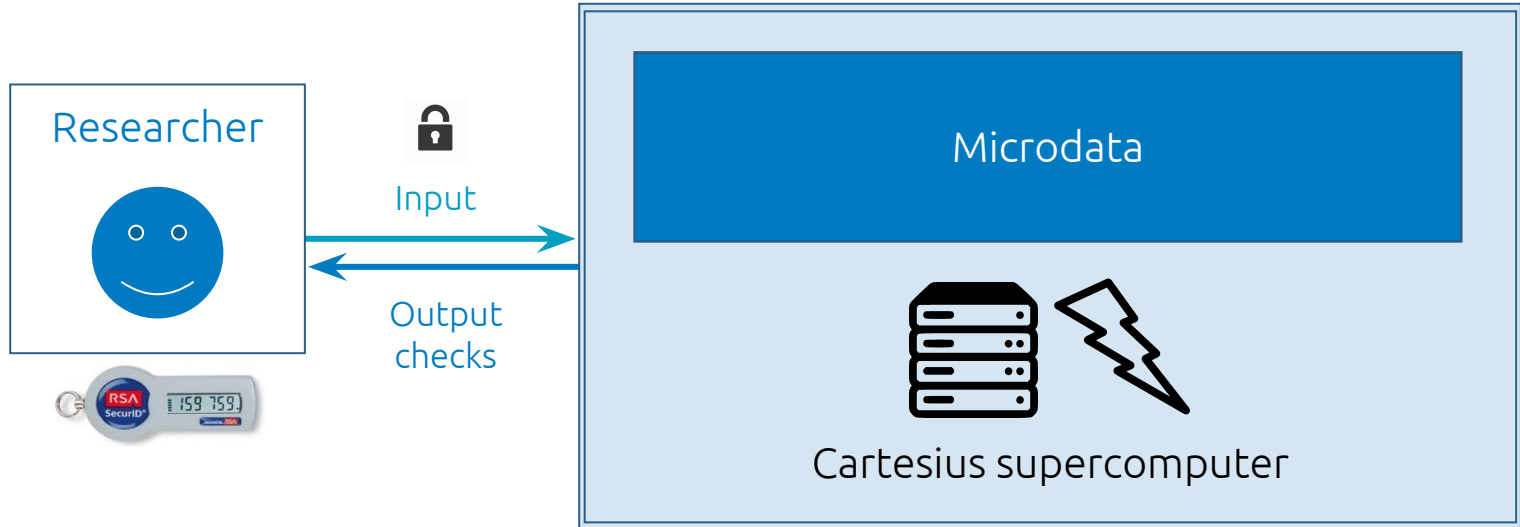
The collaborative
organisation for ICT in
Dutch education and
research



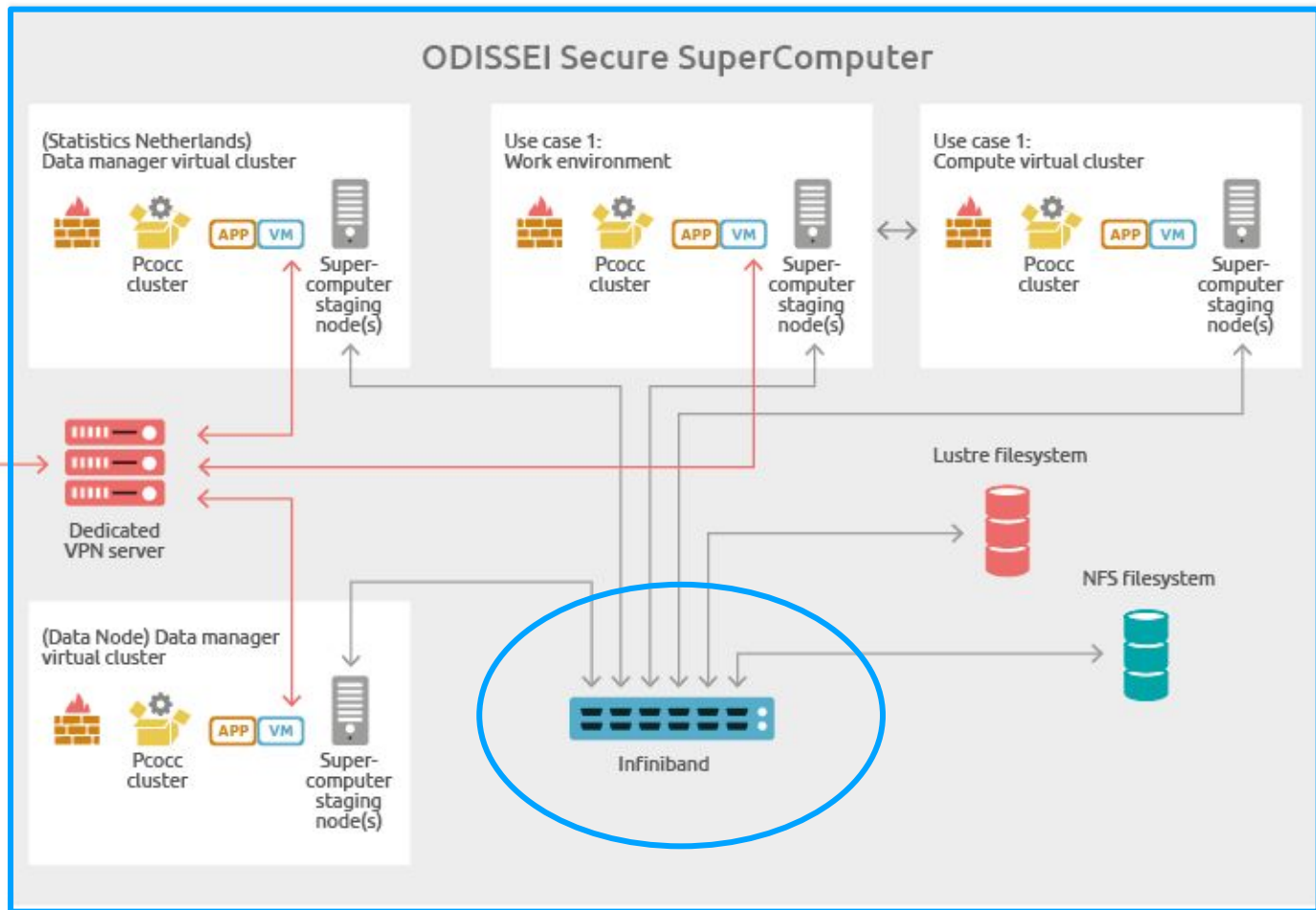
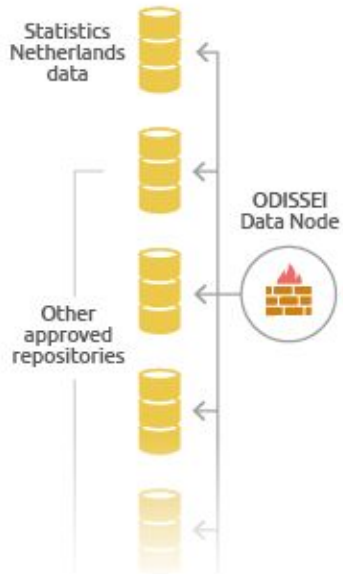
Statistics NL Remote Access environment



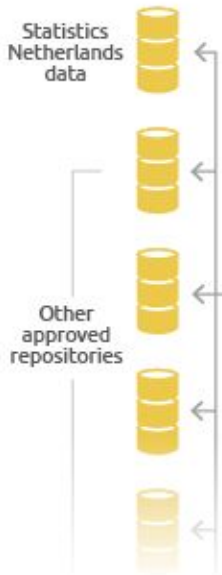
The ODISSEI Secure Supercomputer (OSSC)



Data sources



Data sources



ODISSEI Data Node

VPN via internet



(Statistics Netherlands) Data manager virtual cluster



Pcocc cluster



Super-computer staging node(s)

ODISSEI Secure SuperComputer

Use case 1: Work environment



Pcocc cluster



Super-computer staging node(s)

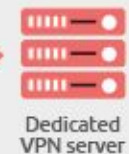
Use case 1: Compute virtual cluster



Pcocc cluster



Super-computer staging node(s)



Dedicated VPN server

(Data Node) Data manager virtual cluster



Pcocc cluster



Super-computer staging node(s)

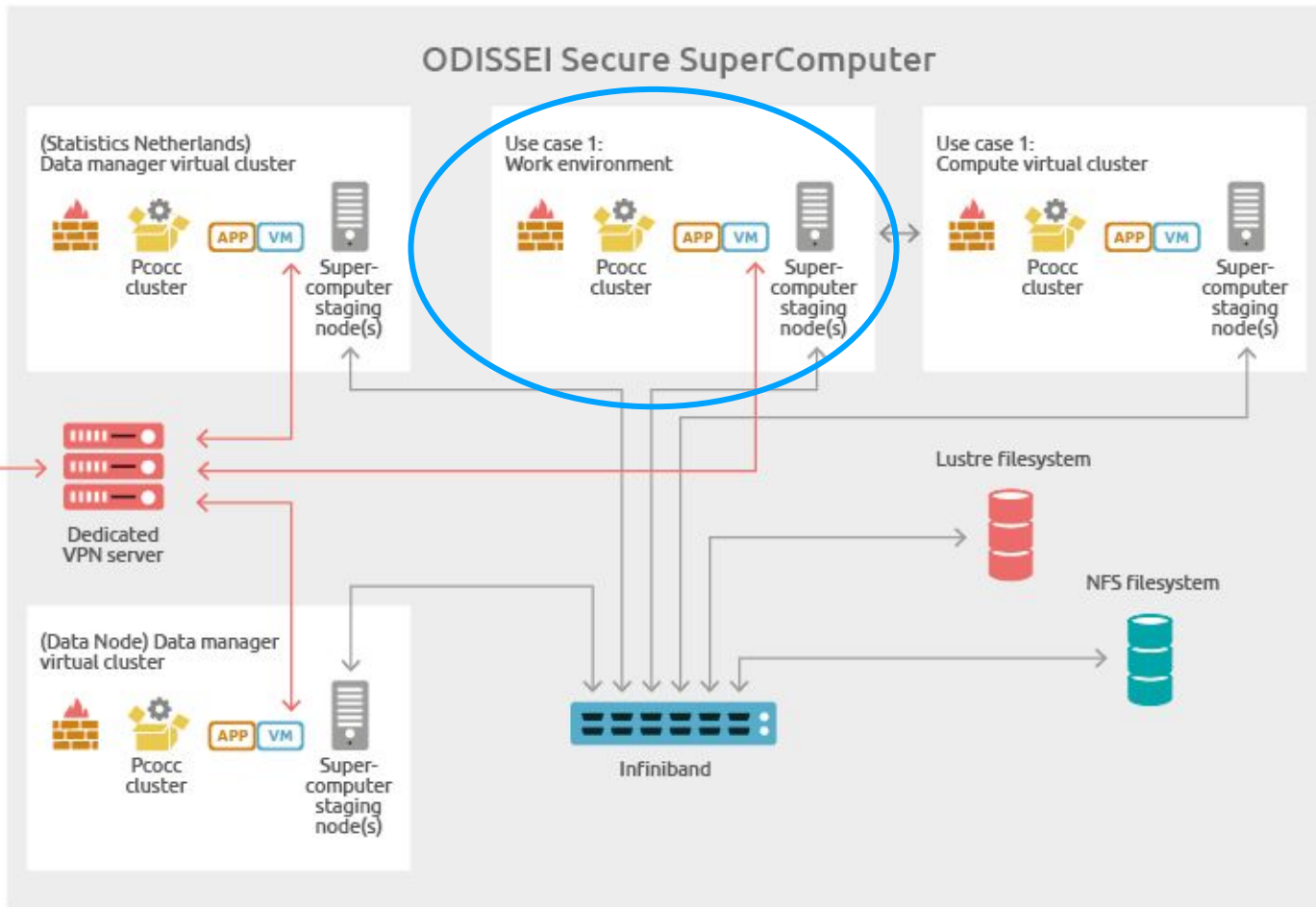
Lustre filesystem



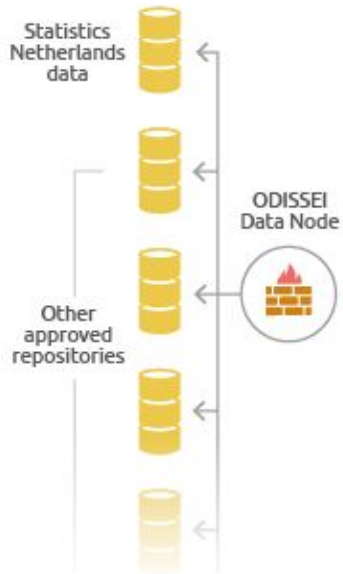
NFS filesystem



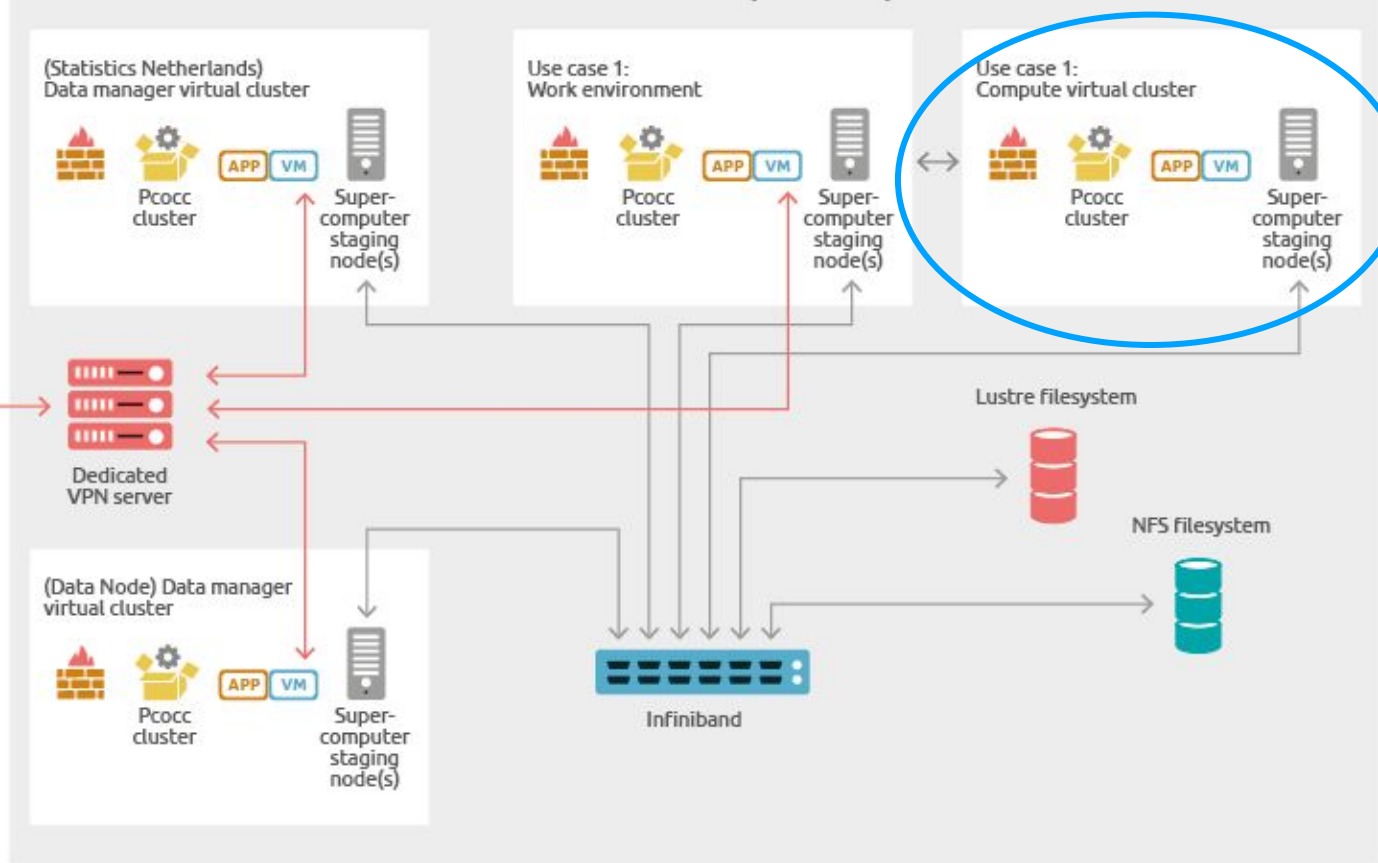
Infiniband



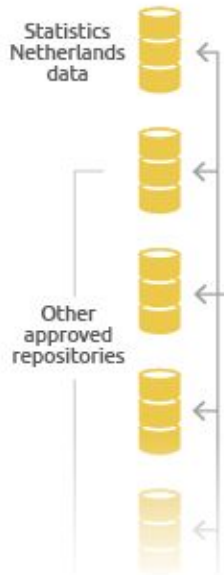
Data sources



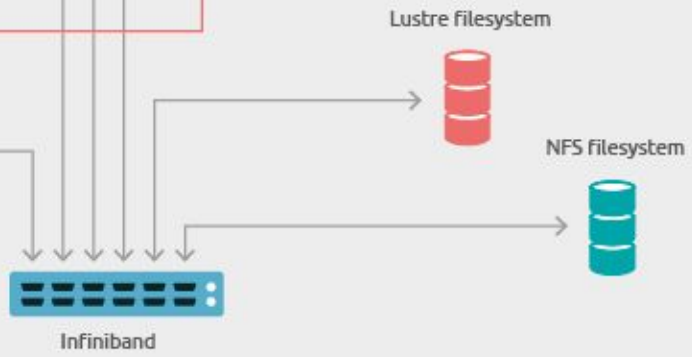
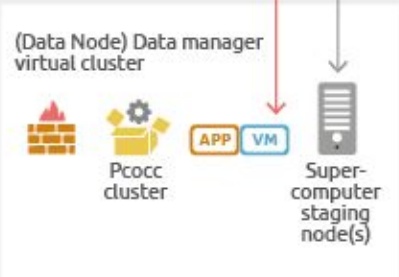
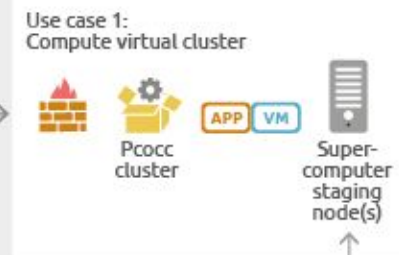
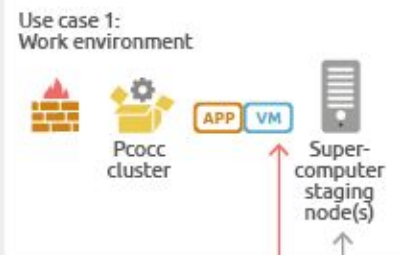
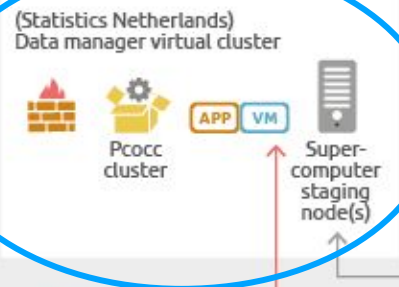
ODISSEI Secure SuperComputer



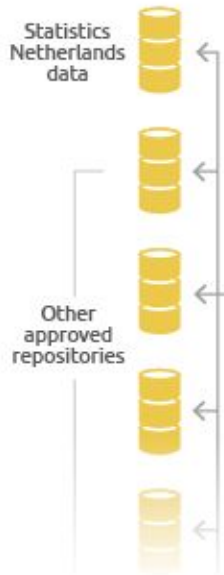
Data sources



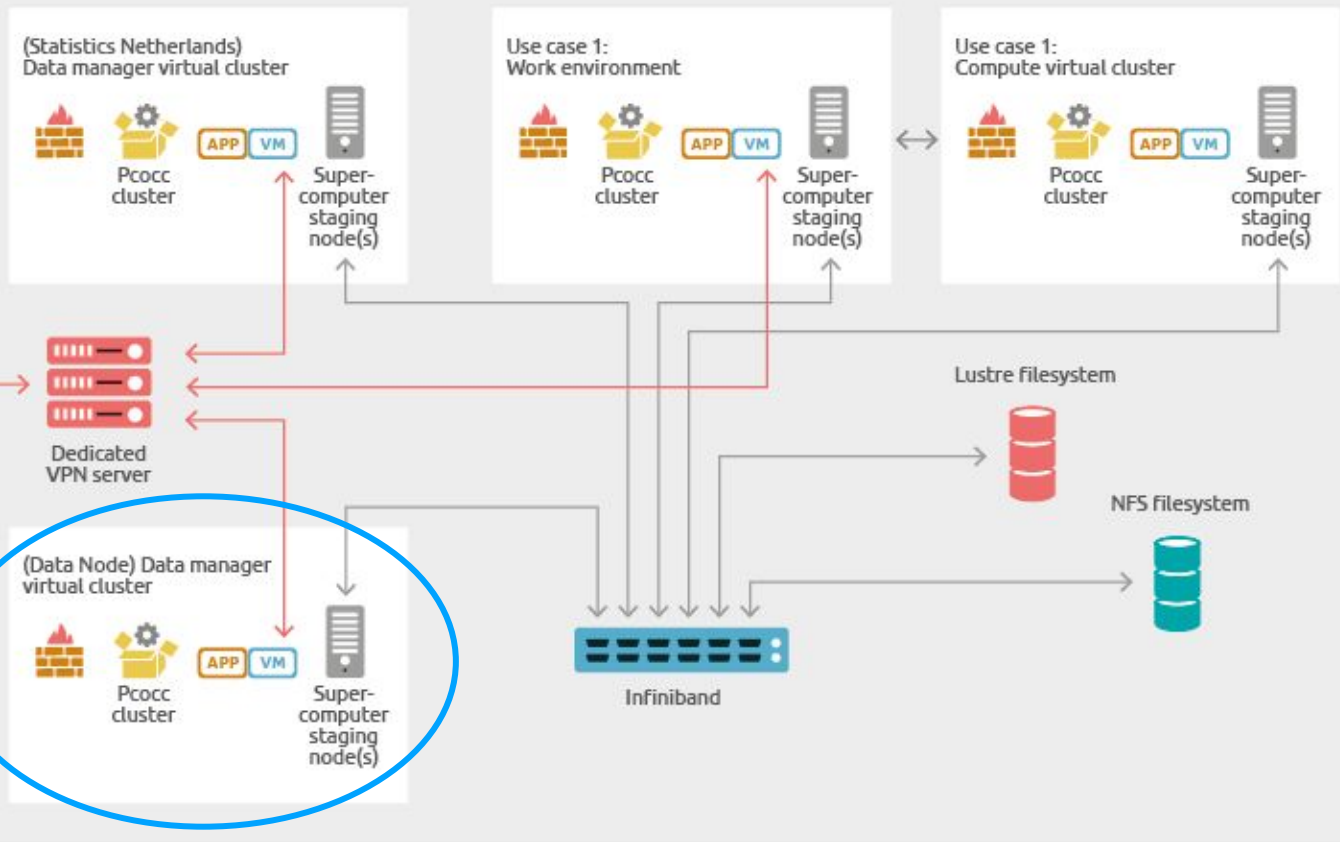
ODISSEI Secure SuperComputer



Data sources



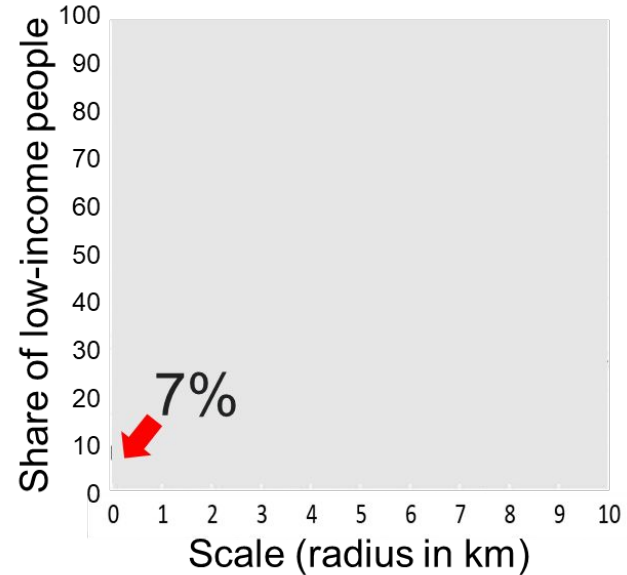
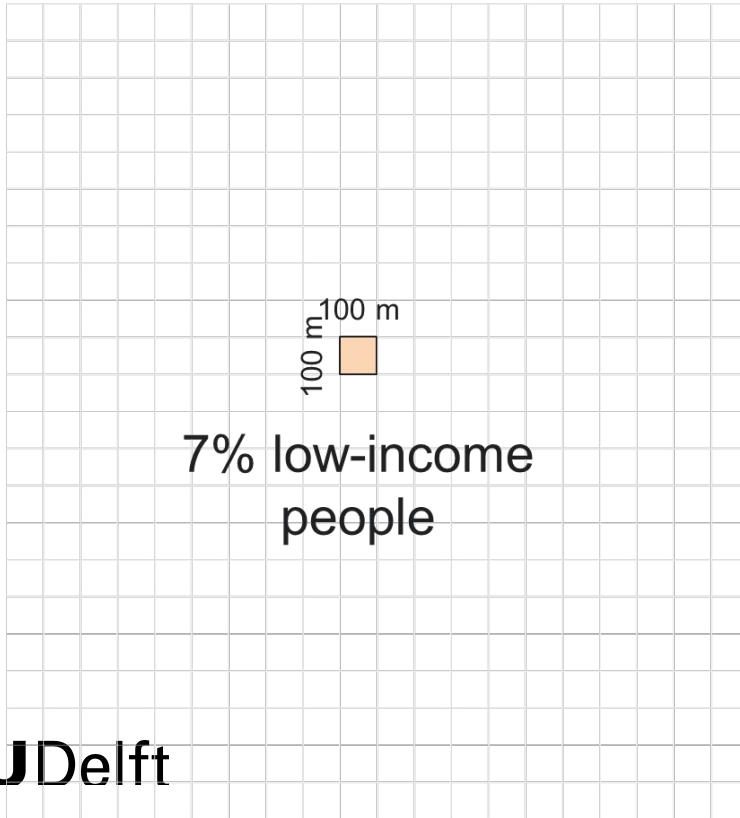
ODISSEI Secure SuperComputer





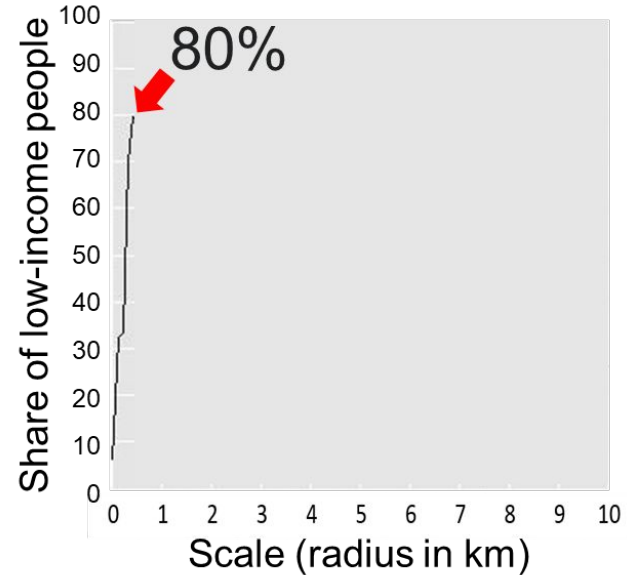
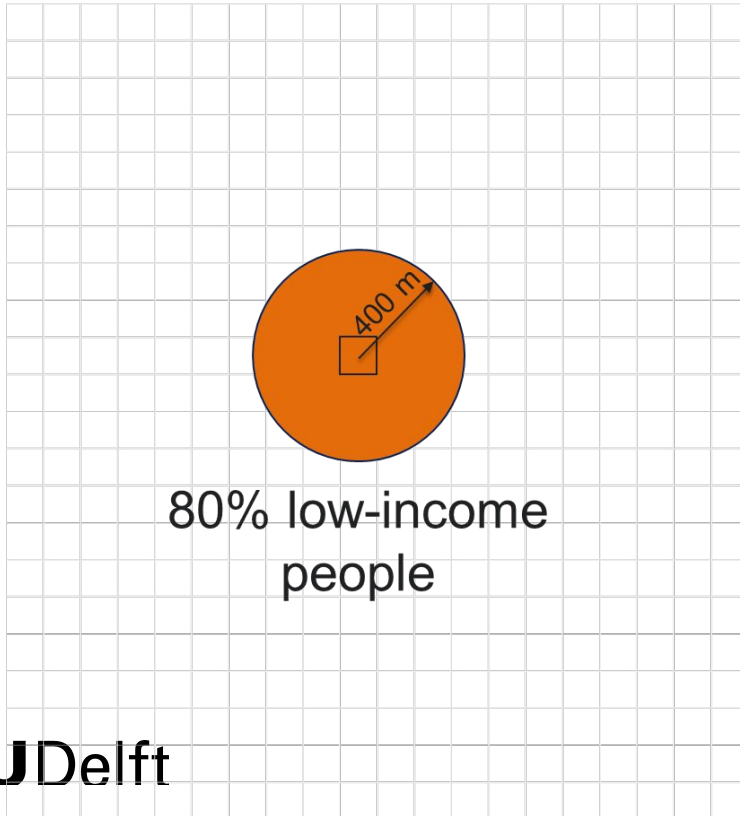
Use Case:
Contextual effects at multiple spatial scales
for the full population of the Netherlands

Exposure to spatial context at 101 scales



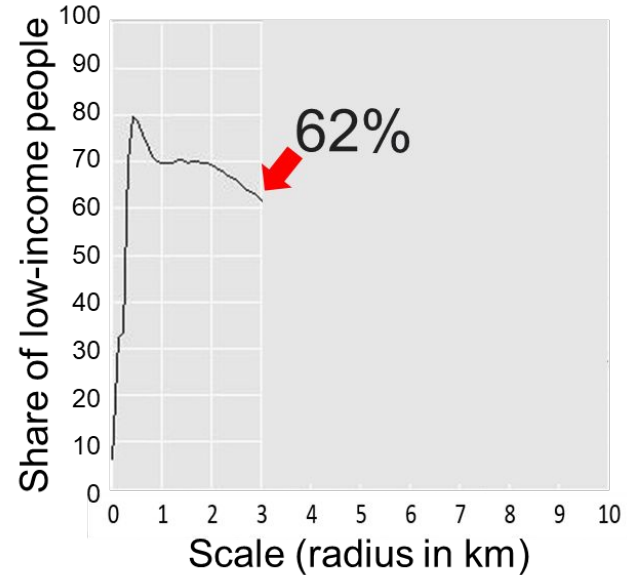
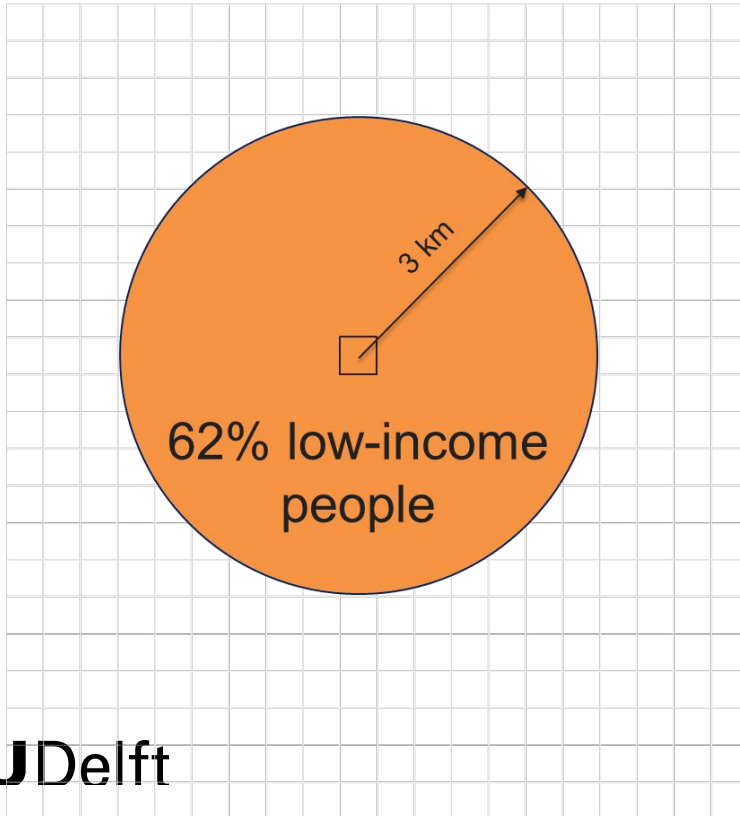
(see Petrović, van Ham, & Manley, 2018)

Exposure to spatial context at 101 scales



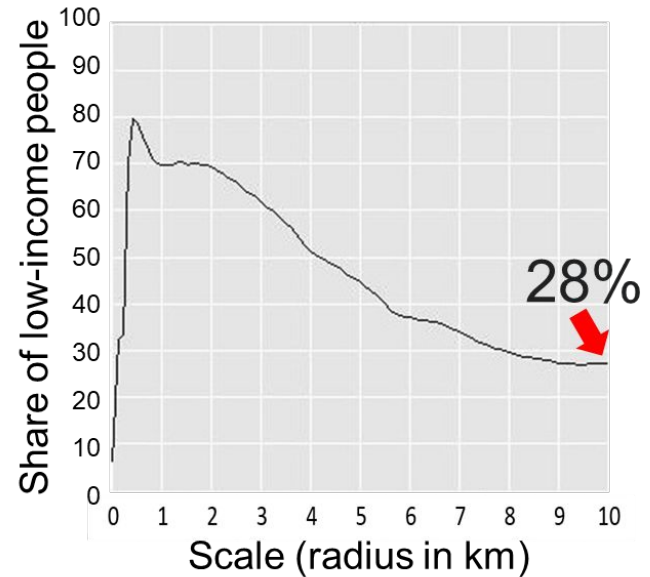
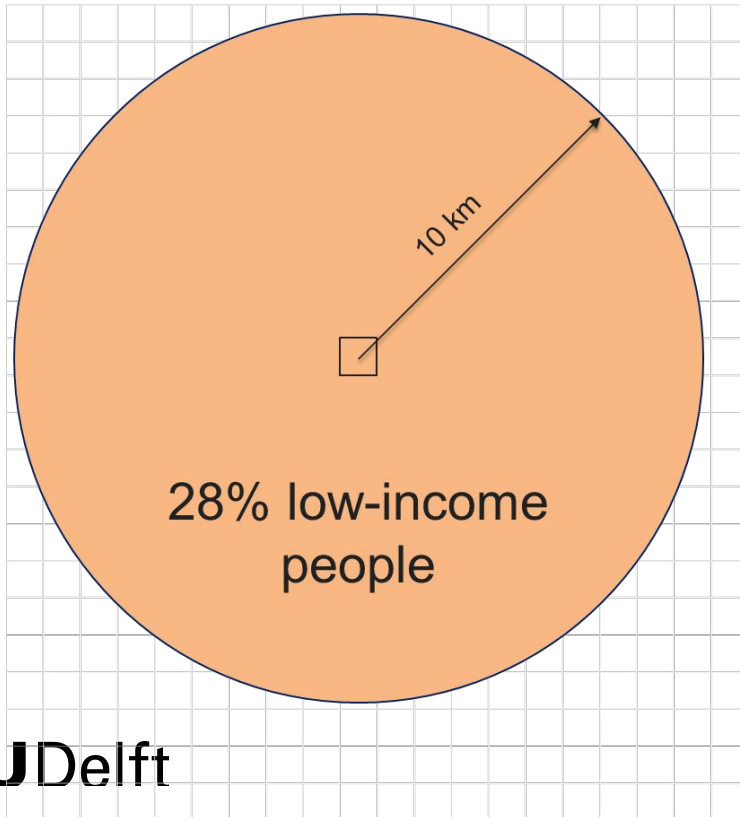
(see Petrović, van Ham, & Manley, 2018)

Exposure to spatial context at 101 scales



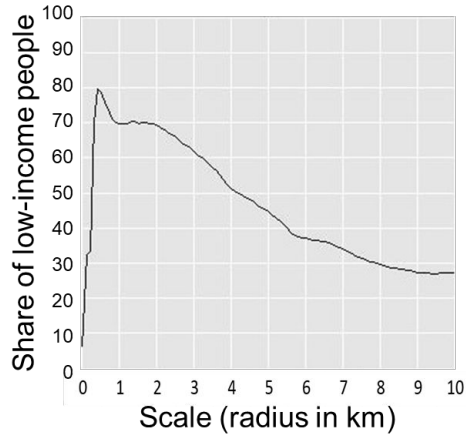
(see Petrović, van Ham, & Manley, 2018)

Exposure to spatial context at 101 scales

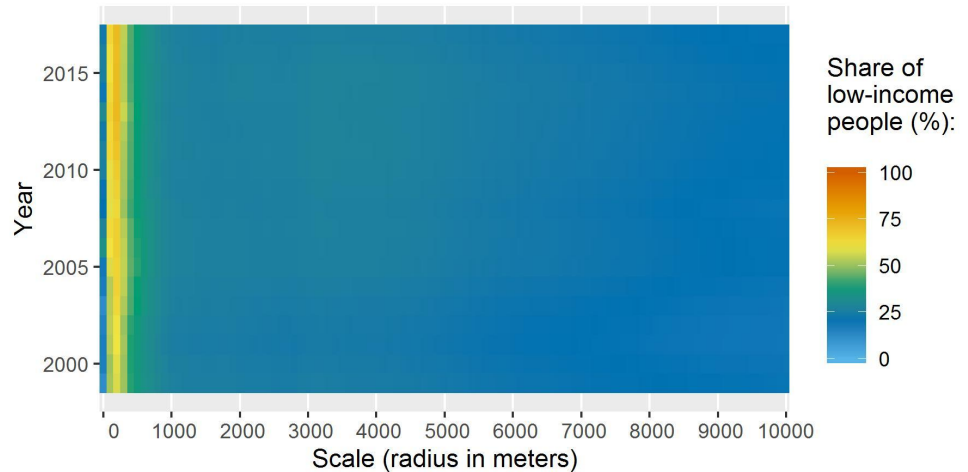


(see Petrović, van Ham, & Manley, 2018)

Share of low-income people in one residential location, measured at 101 scales in 19 years



× 19 =



Challenge: billions of data points

585,000 populated cells
× 101 scales × 2 variables
× 15 years = 1.8 billion
data points

→ 4 months continuous
calculations

Reduced computing time on ODISSEI Secure Supercomputer

585,000 populated cells
× 101 scales × 2 variables
× 15 years = 1.8 billion
data points

→ 4 months continuous
calculations

585,000 populated cells
× 101 scales × 3 variables
× 19 years = 3.4 billion
data points

→ 1 week on 24 nodes



The ODISSEI Secure Supercomputer (OSSC): Unique

Nowhere else is there an infrastructure that allows to analyse register and survey data on a high-performance computer, which is available to the entire community of social scientists.





The ODISSEI Secure Supercomputer (OSSC): Future

The OSSC has been extensively piloted and will be opened up to ODISSEI member organisations in September 2020.

The OSSC design can be adapted to any sensitive data provider, not just Statistics Netherlands. Health insurers, private companies...



Users SSH, builders ICT

- Take time to understand each other
- Plan for early user testing and allow for design adaptations
- Arrange support for researchers (Linux, parallelisation)
- Computational turn in SSH: show, don't tell



Thank you

info@odissei-data.nl

www.odissei-data.nl/oss



ODISSEI

Open Data Infrastructure for Social Science and Economic Innovations

The ODISSEI Secure Supercomputer

Annette Langedijk, Lucas van der Meer • ICTeSSH • 30 June 2020



ODISSEI

Open Data Infrastructure for Social Science and Economic Innovations