

Spider HT platform: Elixir - course

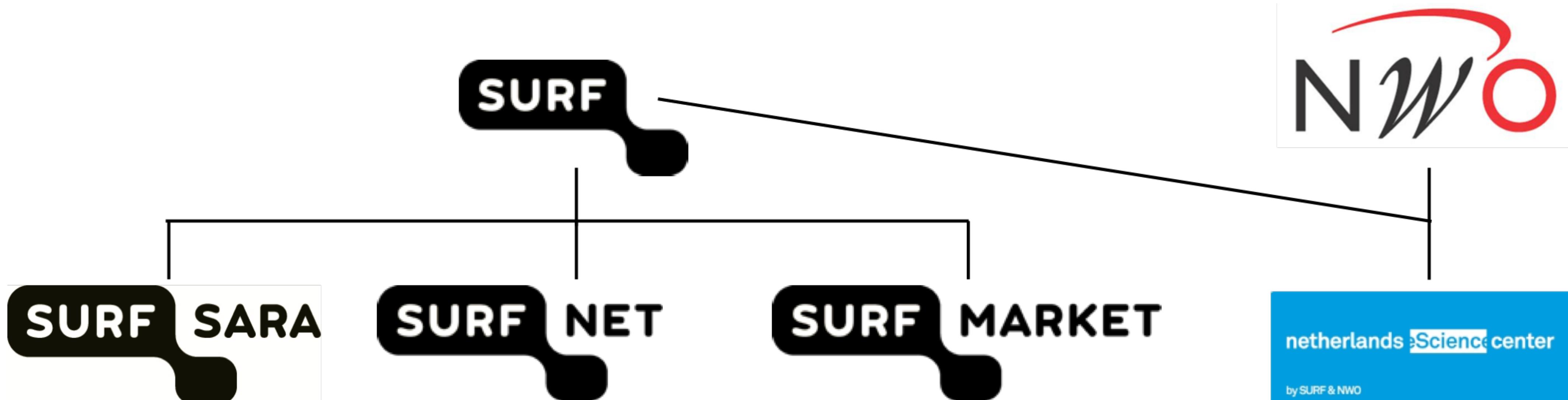
Utrecht, 28 Aug. 2019

Spider team @ SURFsara

SURF SARA

About SURF

ICT services for a strong and sustainable knowledge economy



Academic Medical Centres



100+ and growing

Universities for applied sciences

(+ Senior secondary vocational education institutions)



Universities



KNAW NWO & other research institutes



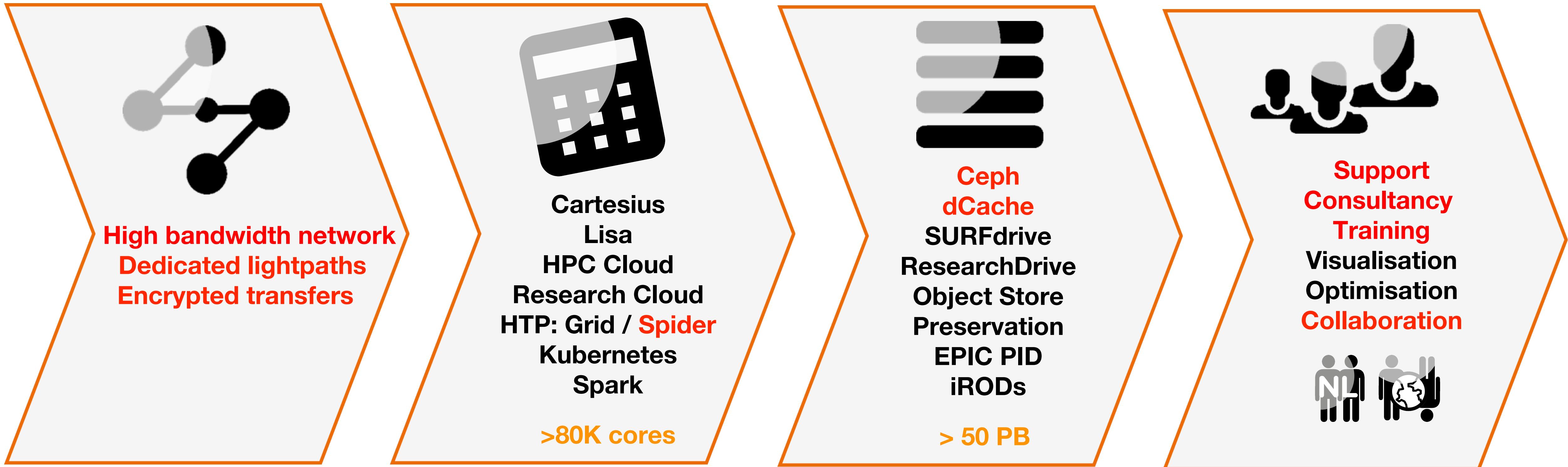


SURF Mission

Supporting Research & Education as:

- Service & infrastructure provider
- Innovator
- Partner in developing new services
- Service broker
- Expertise center

SURF services: (access via surf.nl)



- Data center: 100% renewable energy
- High information security: ISO 27001





**Digital Realty
Data Tower
(SURFsara Data
Center)**

SURFsara

1971: SARA founded by
CWI, UvA, VU
1984: 1st supercomputer
2013: joined SURF
cooperative

NLeSC

High throughput processing

- SURFsara is about making scientific discoveries happen – via high-end IT

High throughput processing

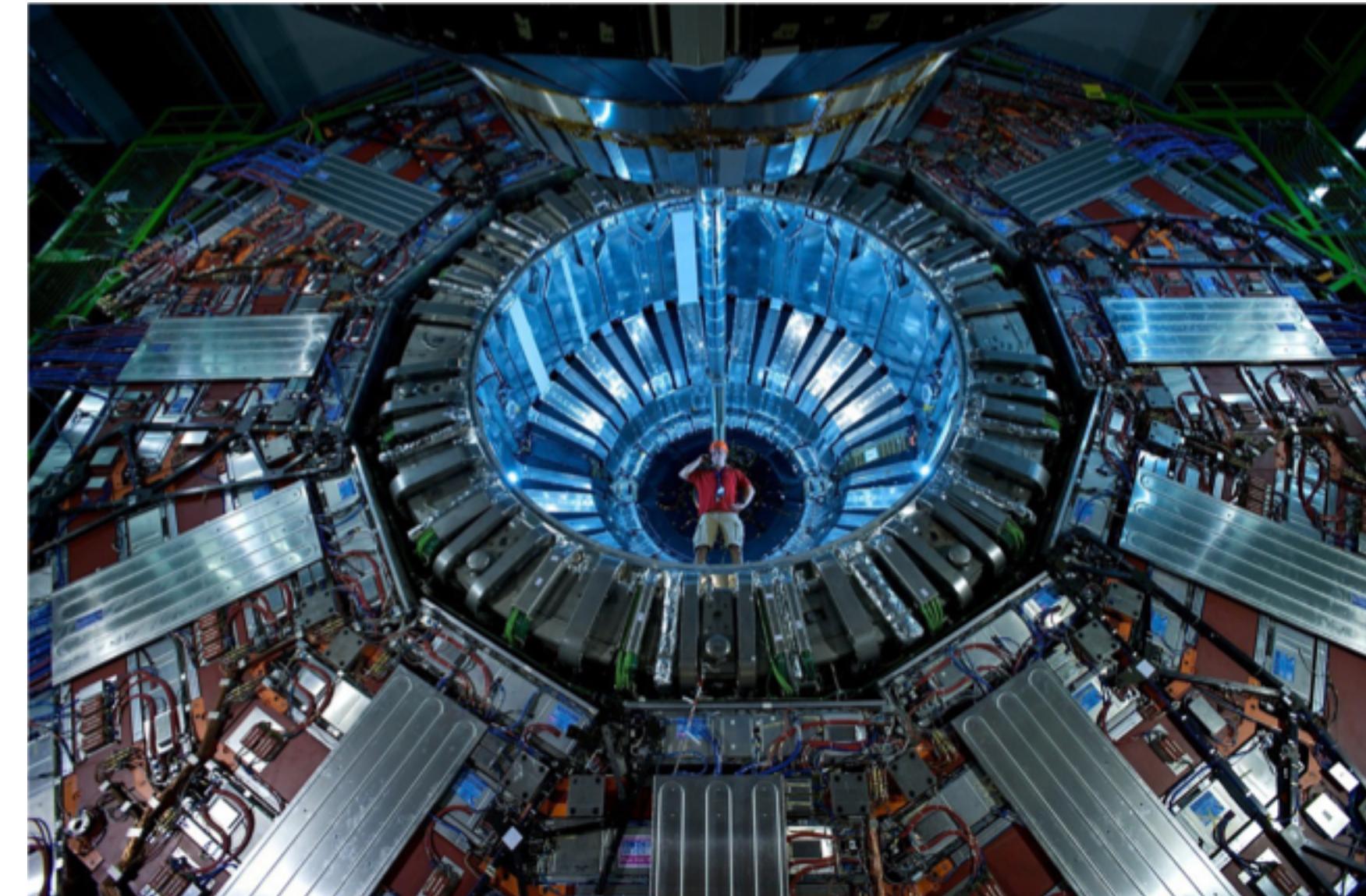
- HT processing is data-driven research => “ from noise to scientific discovery ”

High throughput processing

- HT processing is data-driven research => “ from noise to scientific discovery ”
- SURFsara Grid storage in total grows with **5-10 PB/yr** and has a throughput of ca. **400 PB/yr**
- Experimental data generation grows exponentially **>10x** (*2020-2030: 1 EB / yr per experiment*)

High throughput processing

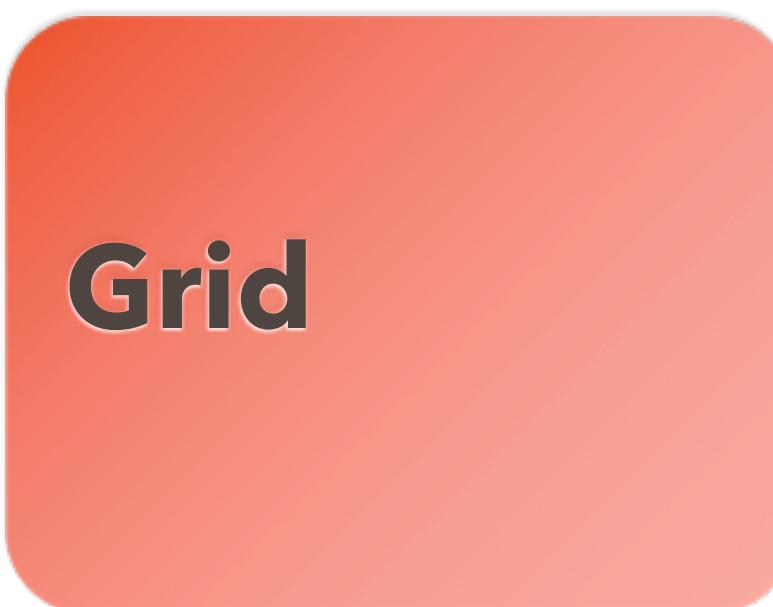
- HT processing is data-driven research => “ from noise to scientific discovery ”
- SURFsara Grid storage in total grows with **5-10 PB/yr** and has a throughput of ca. **400 PB/yr**
- Experimental data generation is increasing by **>10x** (*2020-2030: 1 EB / yr per experiment*)
- HTC: focused on Grid solution
 - WLCG (60 PB/yr , NL 10%)
 - International, distributed processing
 - Long-lived skilled IT teams
 - Simple, independent jobs



(CERN - CMS)



(Big experiments)



- International
- Distributed
- Inhomogeneous
- Certificates
- > Limited functionality



Applications

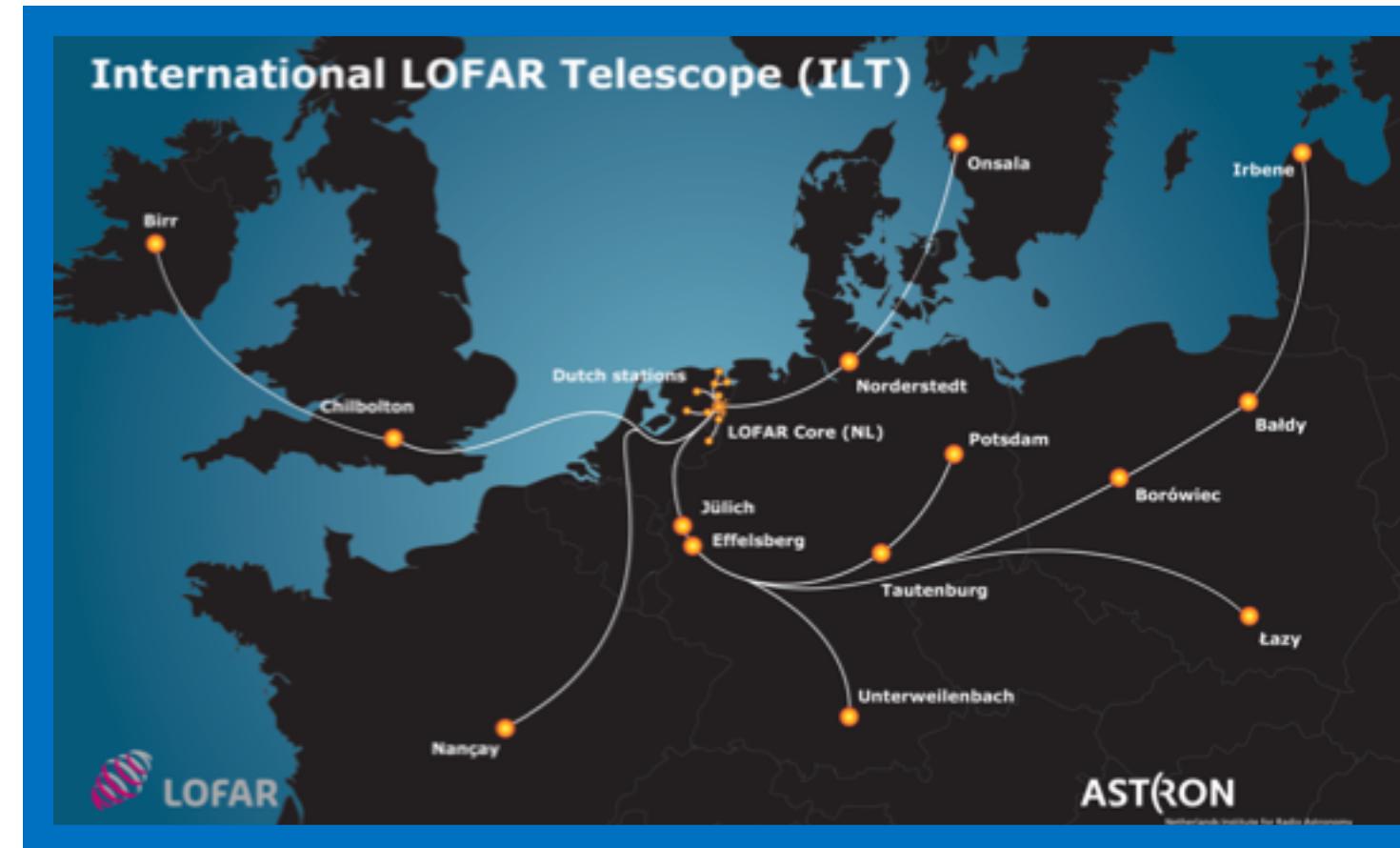
Funnel “rigidity” (Beck+)

Infra / HW



High throughput processing

- **Astronomy , Engineering / Climate , Life sciences : 2020 – 2030 contribute ca. 50% of the data**



LOFAR / SKA



Tropomi (Sentinel 5p)



e.g. Project MinE

- Grid solution insufficient: (a) *Complex workflows*, (b) *Fast moving targets*, (c) *Non-permanent staff*
- Need better solutions to accelerate scientific discovery in (emerging) data-driven research domains

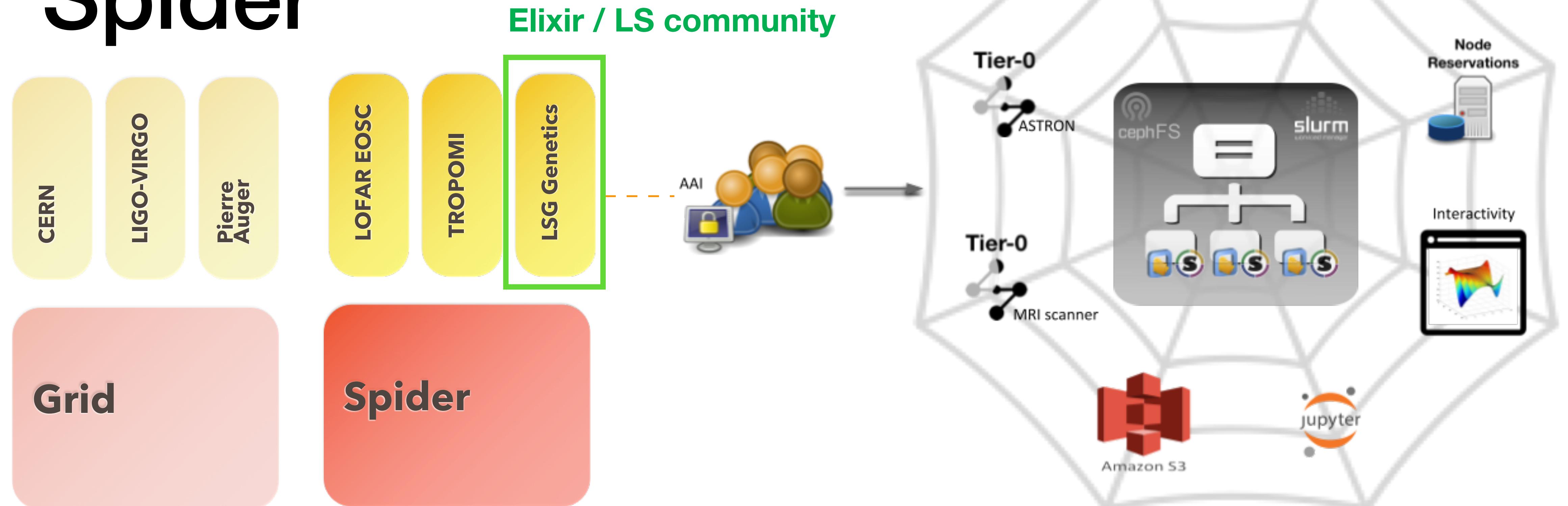
Infrastructure – “ cloudification ”

- DPS services move hardware **from physical to Virtual Data Center (VDC)**
- Flexible deployment & shift of resources for the data processing services
- VDC backbone:
 - Openstack: tuned for data processing (*managed Cloud*)
 - Powerful data processing nodes (*>256 GB Ram, >3 TB SSD*)
 - Ceph: shared, highly scalable storage (*3x Redundancy*)
 - High Bandwidth Network (*internal & external*)

} **Highly scalable !!**



Spider



Spider – aims & options for LS

* HT platform



- analyse large (>>1TB) data collections
high memory, SSD disks, shared filesystem & storage
Slurm scheduler + long running jobs

Spider – aims & options for LS

* HT platform



- analyse large (>>1TB) data collections
high memory, SSD disks, shared filesystem & storage
Slurm scheduler + long running jobs

* Software



- **reproducible workflows – containers & CVMFS**
support e.g., Chipster, Snakemake, Galaxy, DIRAC, Rucio, PiCas
std. protocols (e.g., CURL, rest-API's for data management)

Spider – aims & options for LS

* HT platform



- analyse large (>>1TB) data collections
high memory, SSD disks, shared filesystem & storage
Slurm scheduler + long running jobs

* Software



- reproducible workflows – containers & CVMFS
support e.g., *Chipster*, *Snakemake*, *Galaxy*, *DIRAC*, *Rucio*, *PiCas*
std. protocols (e.g., CURL, rest-API's for data management)

* Collaboration



- project roles & associated rights
project space (*/Data* , */Software*, */Shared*, */Public*)
scientific catalogs (e.g. UK Biobank)

Spider – aims & options for LS

* HT platform



- analyse large (>>1TB) data collections
high memory, SSD disks, shared filesystem & storage
Slurm scheduler + long running jobs

* Software



- reproducible workflows – containers & CVMFS
support e.g., *Chipster*, *Snakemake*, *Galaxy*, *DIRAC*, *Rucio*, *PiCas*
std. protocols (e.g., CURL, rest-API's for data management)

* Collaboration



- project roles & associated rights
project space (*/Data* , */Software*, */Shared*, */Public*)
scientific catalogs (e.g. UK Biobank)

* Interactive



- **Jupyter NB, public webviews, interactive Slurm jobs (w. graphics)**

Spider – aims & options for LS

* HT platform



- analyse large (>>1TB) data collections
high memory, SSD disks, shared filesystem & storage
Slurm scheduler + long running jobs

* Software



- reproducible workflows – containers & CVMFS
support e.g., *Chipster*, *Snakemake*, *Galaxy*, *DIRAC*, *Rucio*, *PiCas*
std. protocols (e.g., CURL, rest-API's for data management)

* Collaboration



- project roles & associated rights
project space (*/Data* , */Software*, */Shared*, */Public*)
scientific catalogs (e.g. UK Biobank)

* Interactive



- Jupyter NB, public webviews, interactive Slurm jobs (w. graphics)
- encrypted data transfers, secure light-paths, ISO 27001 certified

* Security



Spider – aims & options for LS

* HT platform



- analyse large (>>1TB) data collections
high memory, SSD disks, shared filesystem & storage
Slurm scheduler + long running jobs

* Software



- reproducible workflows – containers & CVMFS
support e.g., *Chipster*, *Snakemake*, *Galaxy*, *DIRAC*, *Rucio*, *PiCas*
std. protocols (e.g., CURL, rest-API's for data management)

* Collaboration



- project roles & associated rights
project space (/Data , /Software, /Shared, /Public)
scientific catalogs (e.g. UK Biobank)

* Interactive



- Jupyter NB, public webviews, interactive Slurm jobs (w. graphics)

* Security



- encrypted data transfers, secure light-paths, ISO 27001 certified

* Customize



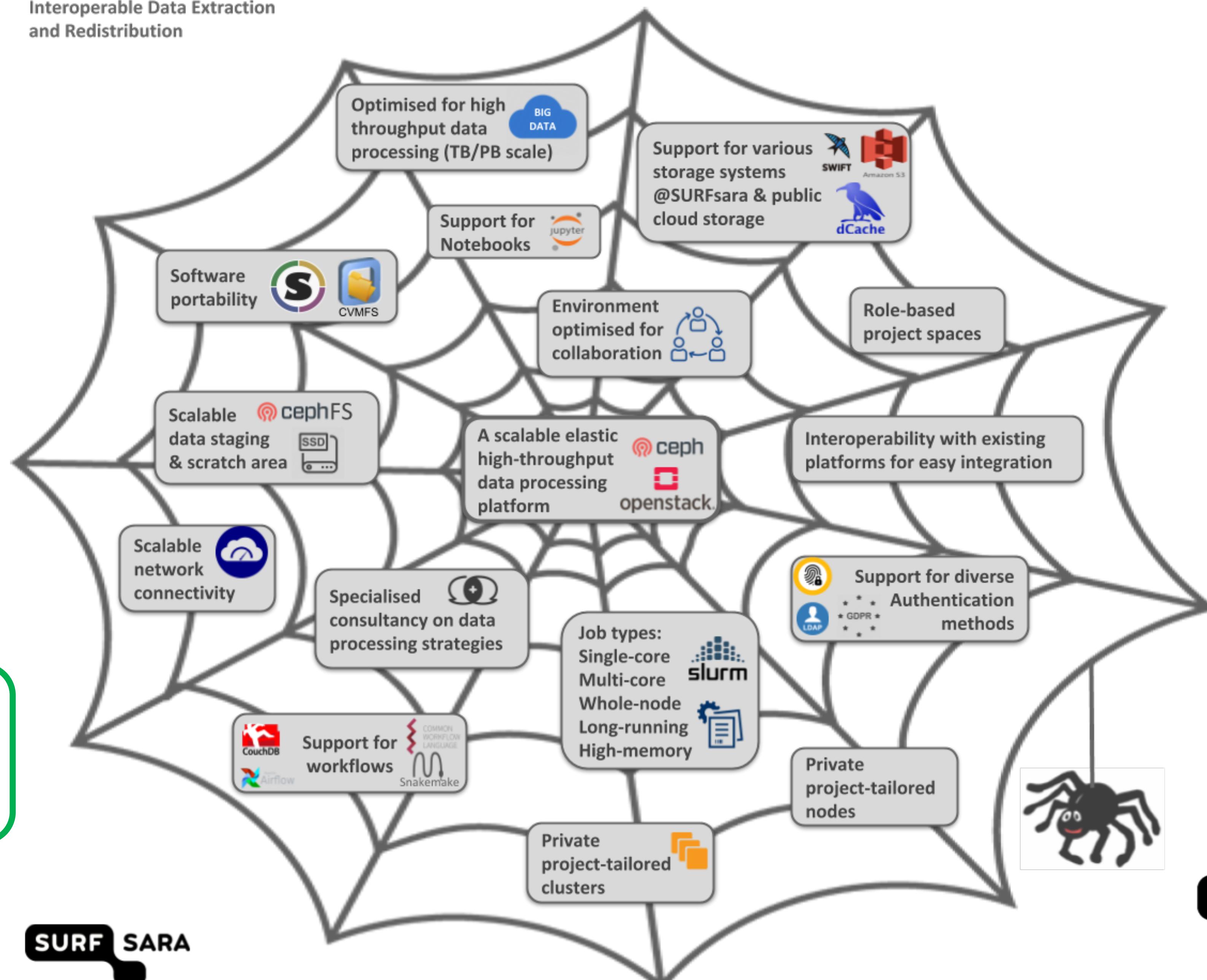
- private nodes (e.g. host private DB) / dedicated clusters

Phase 1

The ‘SPIDER’ project

Symbiotic Platform(s) for
Interoperable Data Extraction
and Redistribution

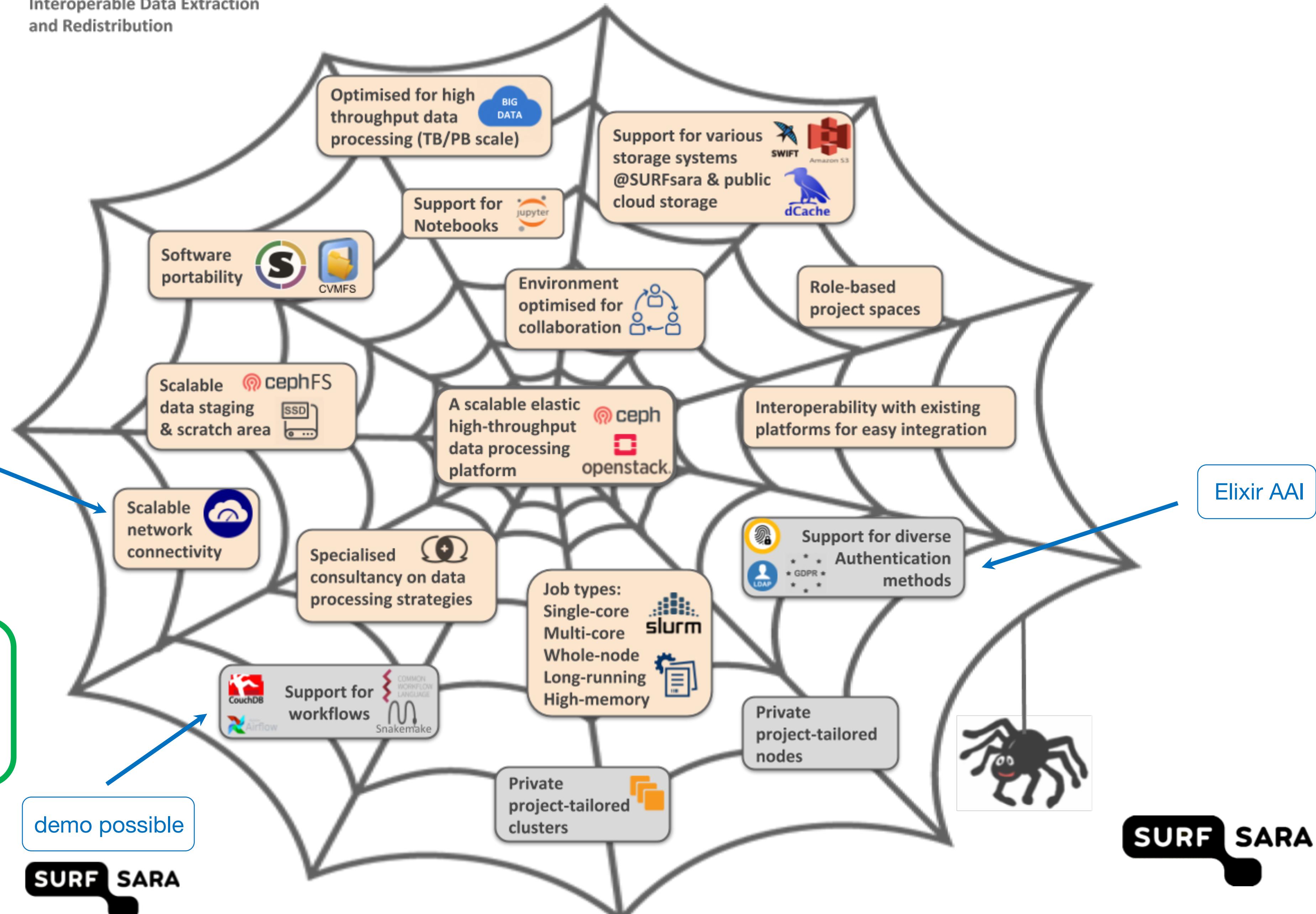
“Interoperable”

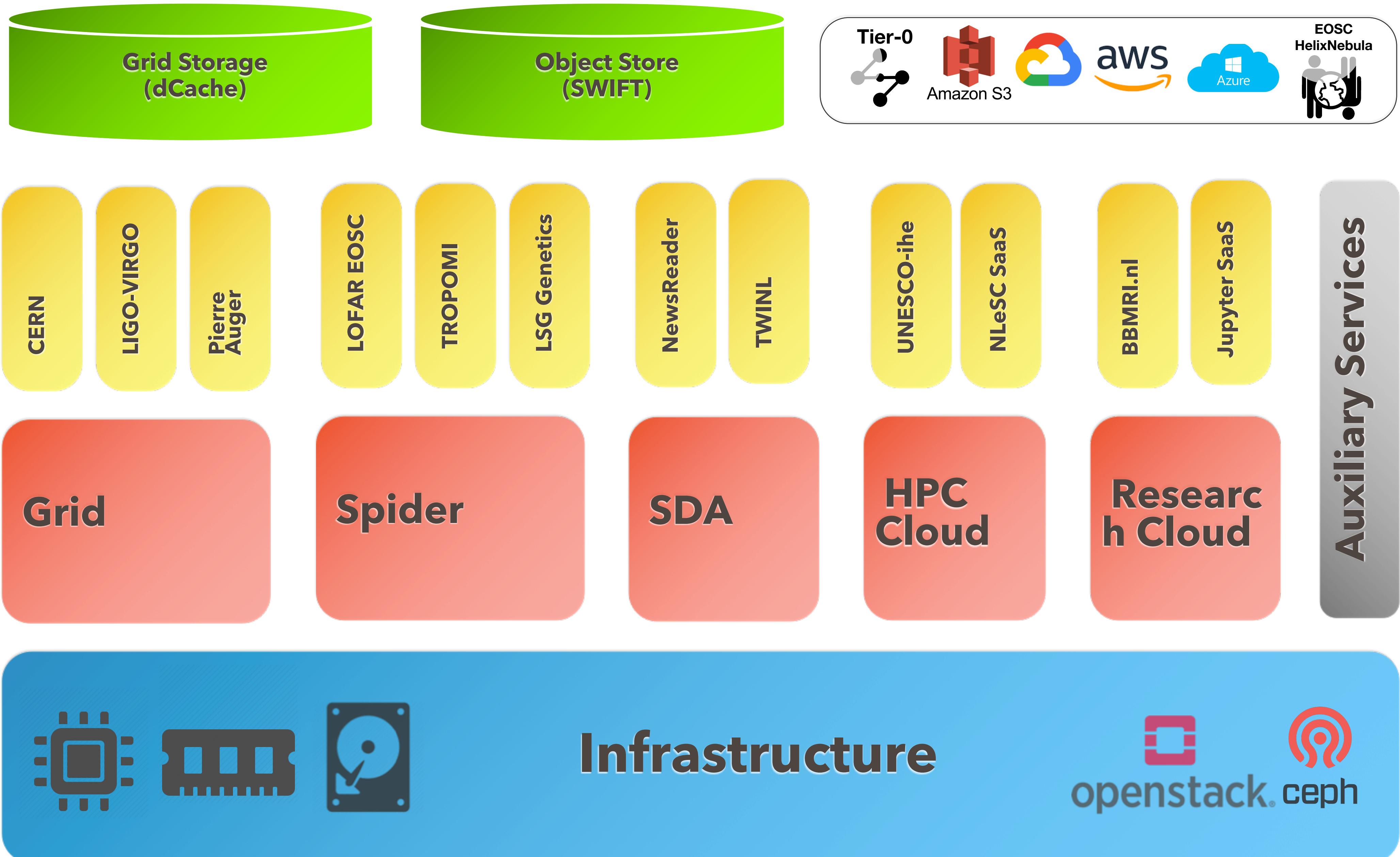


Current (beta)

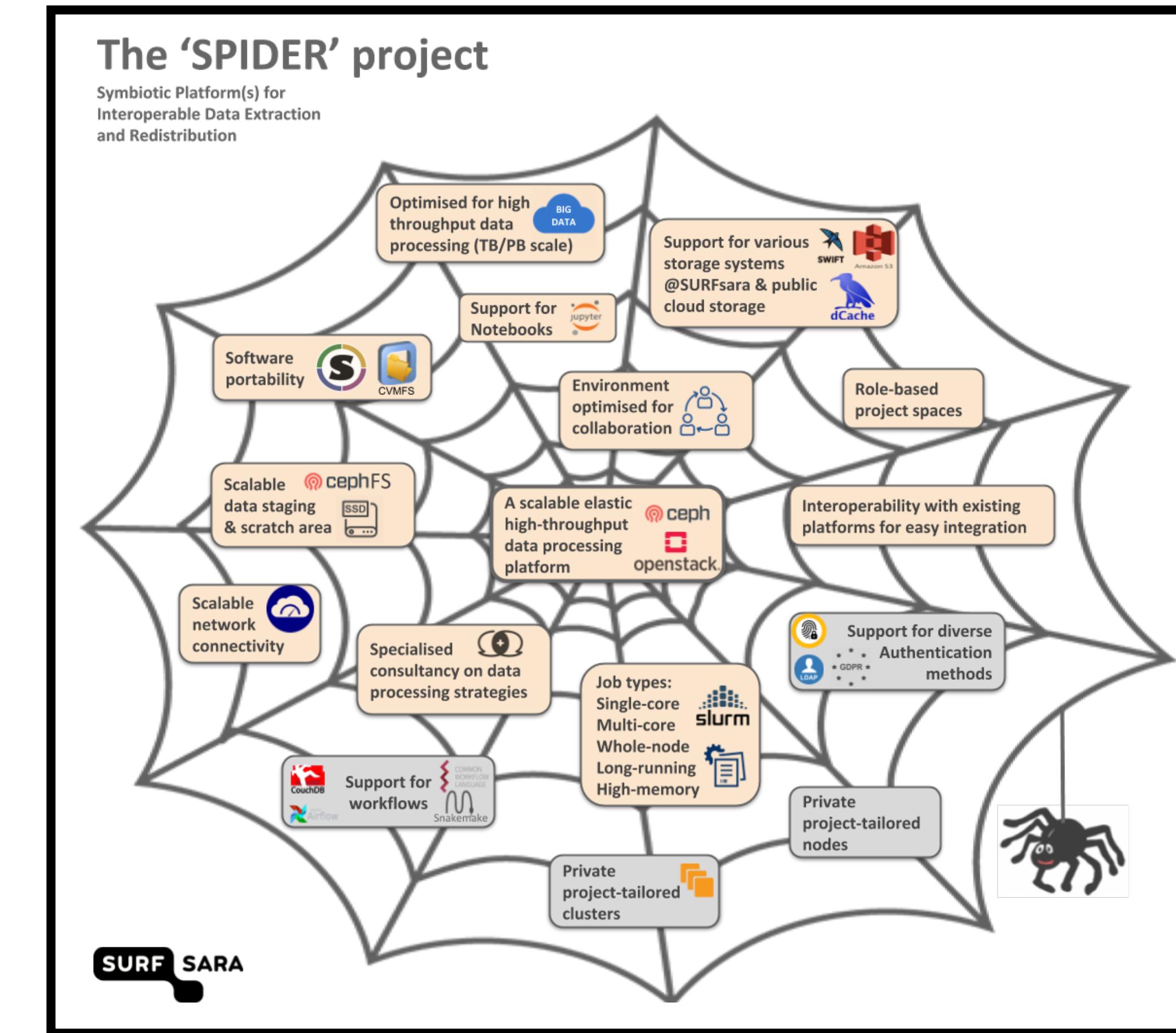
The 'SPIDER' project

Symbiotic Platform(s) for
Interoperable Data Extraction
and Redistribution





Please go to:



<https://github.com/sara-nl/spidercourse/>

Demo

Loui: Control center

dCache: Data storage

Gina: Data processing

PiCas: Data/Job metadata



: Softdrive SW distribution



: Container image for SW portability



: Todo token



: Done token



: Workload submission



: Input/output file

LOFAR Surveys KSP

