# Managing dataflows in a research hospital

Managed Research Data Management with iRODS

NETHERLANDS
CANCER
INSTITUTE
ANTONI VAN LEEUWENHOEK

# The Netherlands Cancer Institute

The Netherlands Cancer Institute comprises an internationally acclaimed research institute as well as a dedicated cancer clinic. This combination ensures rapid translation of basic research into clinical applications: today's research for tomorrow's cure.
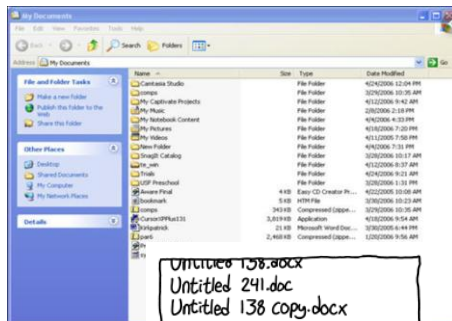
**Today's research for tomorrow's cure**

*We believe that cancer does not need to be a deadly disease. Our researchers and doctors are highly motivated to make this vision a reality by unraveling the biology of cancer and using this knowledge to improve the prospects of cancer patients. They conduct innovative and excellent fundamental, translational, and clinical research. Through close collaboration between lab and clinic, we create maximum impact for cancer patients.*

# The challenge

- At the Antoni van Leeuwenhoek hospital and Netherlands Cancer Institute (NKI-AVL) the demands for data storage and data management are rapidly evolving. **Departments in our institute increasingly integrate their data acquisition and analysis, sparking interdisciplinary research projects**. Furthermore, national and international regulations require researchers to **make their data FAIR** (Findable, Accessible, Interoperable, Reusable). Also, the development of the "-omic" techniques, such as genomic and proteomics, massively **increases the size of acquired data**. After analysis, this data should be **archived for longer periods of time** (>10y). Storing this data on rapid and available storage is **a waste of resources and money**. All these developments necessitate **meta-data driven data management**.
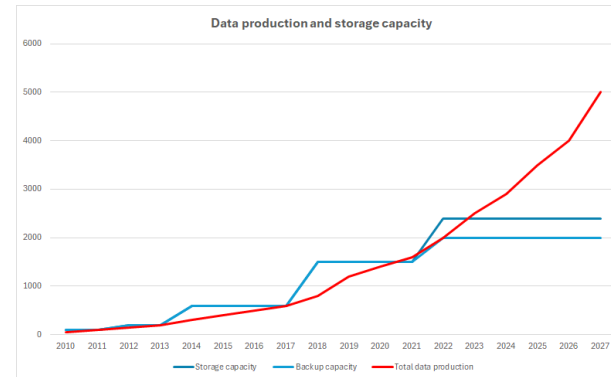
### RDM tools available



### Data production vs capacity



### (Data) Interaction between clinic and research



### Cost and budget



Hardware and financial investment.

3

# Research data (storage) challenges

- The digital era is well under way

- More and more data is being generated
  - Size and number of the data is increasing
  - Data is producing new data  (fueling itself)

- How to keep data findable and reusable data (FAIR)?

- Were do we put the data at what cost

- Is data safe and can we find the data?

- How do we share data (safe and controlled)



***Data production is increasing rapidly***

# Looking at the data; a research perspective

- RDM not a priority for researcher (just wants to store data)

  - (Manual) RDM distracts from research

- Data is not/poorly structured & described

  - Finding the data is hard/data gets lost

  - Do I have the correct data?

  - Poor to no insight in data linage

- Research outcome and reproducibility?

- Sharing data?

- Data automation is hard

NETHERLANDS
CANCER
INSTITUTE
ANTONI VAN LEEUWENHOEK
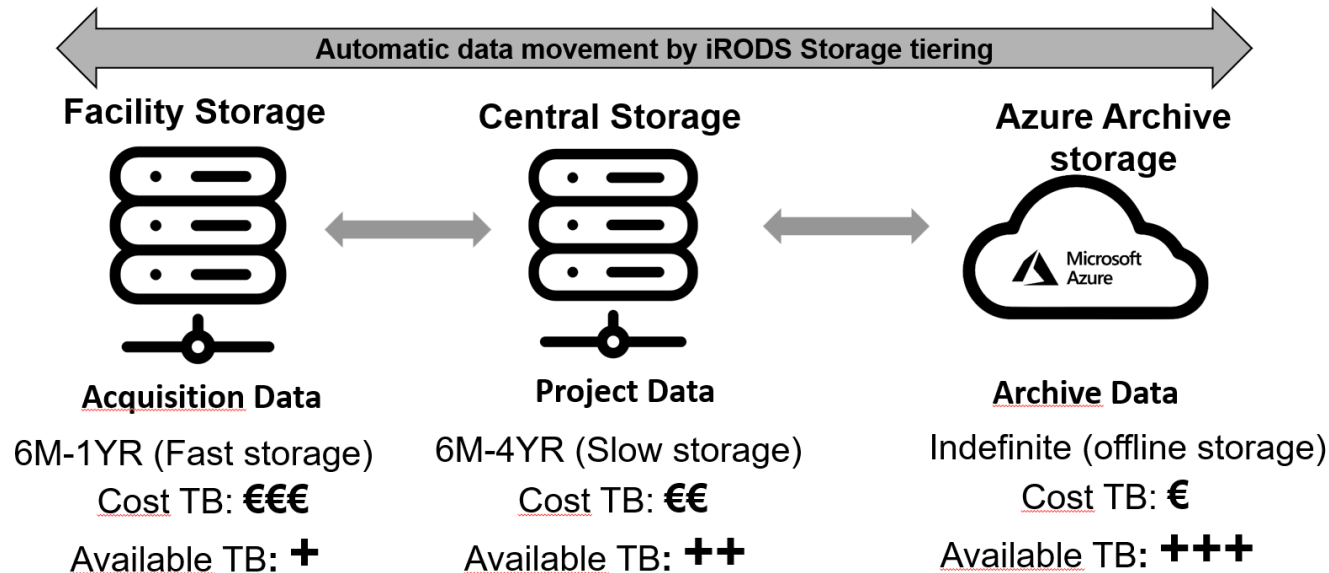
# Looking at facts of the data

- Data production and growth is (becoming) exponential

- Data is never deleted

- Big differences in data between research groups and facilities
  - Facilities are data producers, groups are data users

- 80 % of data is not being used any more (cold data)
  - Data becomes cold after ~6 months

- Why do we keep old/cold data on expensive storage (online)
  - and backup old/cold data in operational backup?
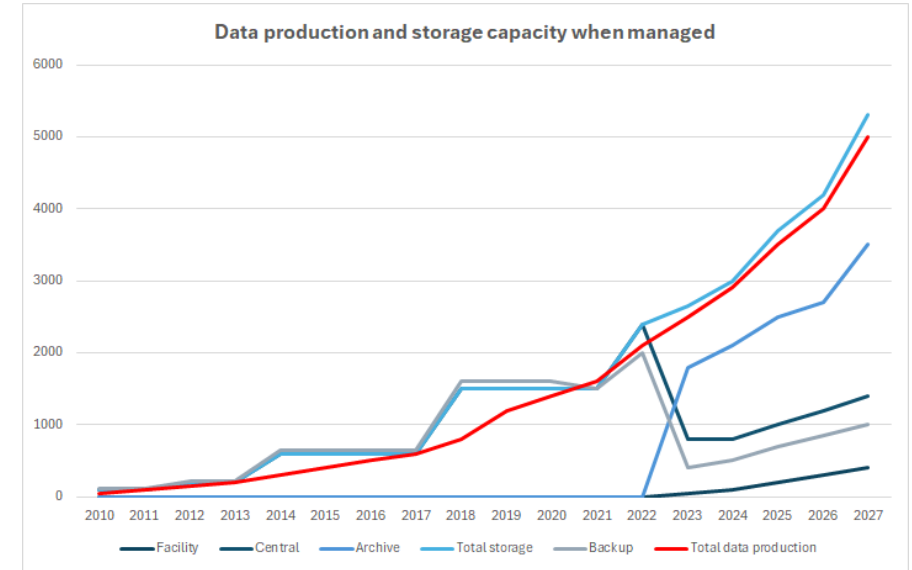
- Data is lost (not to be found)

NETHERLANDS
CANCER
INSTITUTE
ANTONI VAN LEEUWENHOEK

# iRODS solution for storage

**Data lifetime, storage and costs**

iRODS



Data production and storage capacity when managed

**Automatic data movement by iRODS Storage tiering**

| Facility Storage | Central Storage | Azure Archive storage |
|---|---|---|
| **Acquisition Data** | **Project Data** | **Archive Data** |
| 6M-1YR (Fast storage) | 6M-4YR (Slow storage) | Indefinite (offline storage) |
| Cost TB: €€€ | Cost TB: €€ | Cost TB: € |
| Available TB: **+** | Available TB: **++** | Available TB: **+++** |

*\* Data policies meta data driven and flexible*

**Managed data**
- Manage storage on lifetime and usage
- Implement storage at the correct location and store data where it's needed (onlive vs archive)
- Data archiving is automated, flexible and transparent for end user
- Annotate data with meta data so it can be found
- Multiple way's/API's to access data

NETHERLANDS
CANCER
INSTITUTE
ANTONI VAN LEEUWENHOEK

7

# Business case / financial savings

| | | |
|---|---|---:|
| Local storage cost / TB / Yr | € | 200,00 |
| Cloud archive storage cost / TB / ' € | | 30,00 |

| StorageUnit | CuttOfMonths | NrColdFiles | NrActiveFIles | NrModifiedFiles | NrTotalFiles | TBColdData | TBActive | TBModified | TBTotal |
|---|---|---|---|---|---|---|---|---|---|
| \\Facility Local | 6 | 1479588 | 660113 | 640601 | 2139701 | 29,6413 | 28,964 | 28,4546 | 58,6051 |
| \\Large Storage | 6 | 74632360 | 10912335 | 6671078 | 85544695 | 1120,3744 | 307,2595 | 165,7107 | 1427,8446 |
| \\General Research | 6 | 87615366 | 2800975 | 916289 | 90416341 | 266,297 | 36,1467 | 20,409 | 302,5035 |
| Total | | 163727314 | 14373423 | 8227968 | 178100737 | 1416,3 | 372,4 | 214,6 | 1789,0 |
| | | 92% | 8% | 5% | | 79% | 21% | 12% | |

| COST: | | | | | |
|---|---|---|---|---|---:|
| | No RDM active | | | € | 357.790,64 |
| | With RMD/iRODS | € 42.489 | € 74.474 | € | 116.963,42 |
| | **Yearly savings** | | | **€** | **240.827,22** |

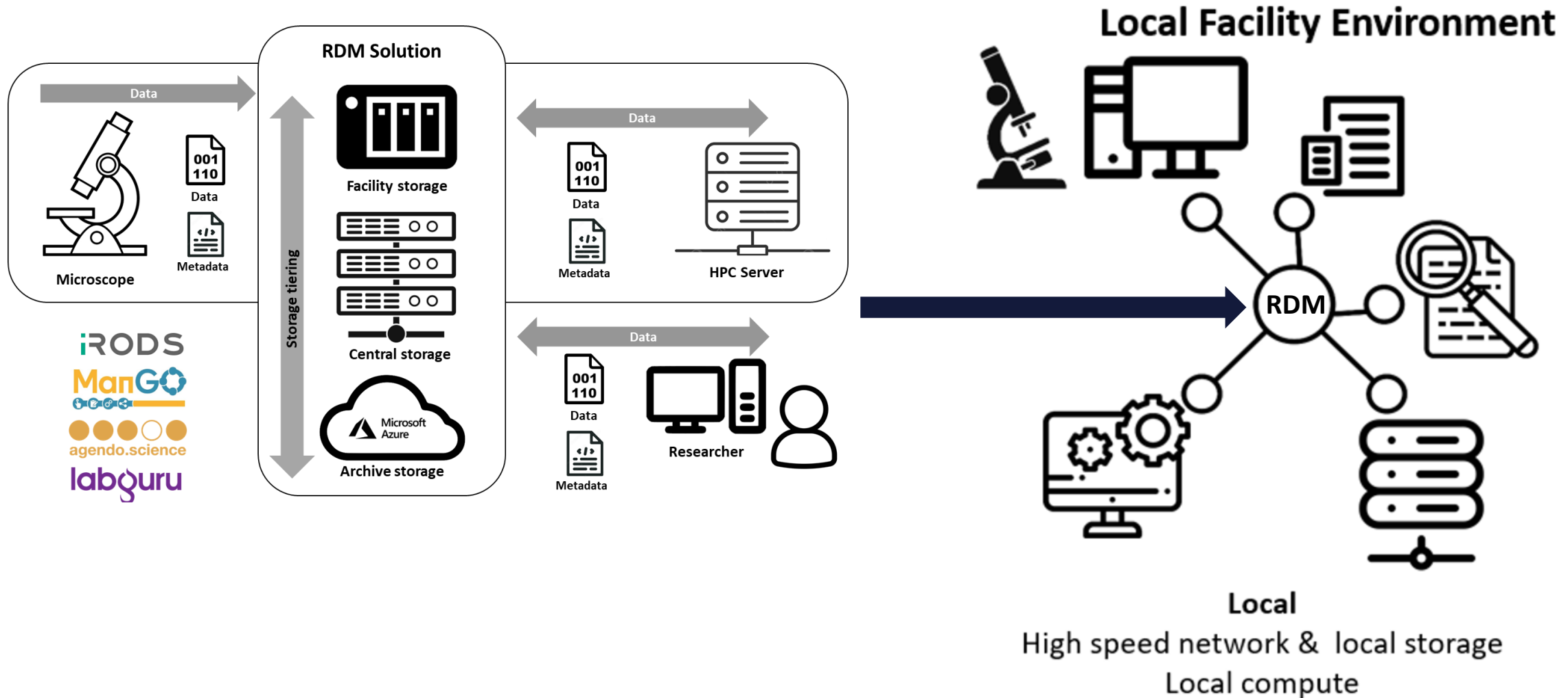| Old Data | | 800,00 TB | | |
|---|---|---|---|---:|
| | | Online costs | € | 160.000,00 |
| | | Archived costs | € | 24.000,00 |
| **Savings** | | | **€** | **136.000,00** |

**Total yearly savings**   **€   376.827,22**

**Business-case to support RDM/iRODS:**

*By correctly scaling and managing storage, financial investments in IT are reduced and research can be kept going*
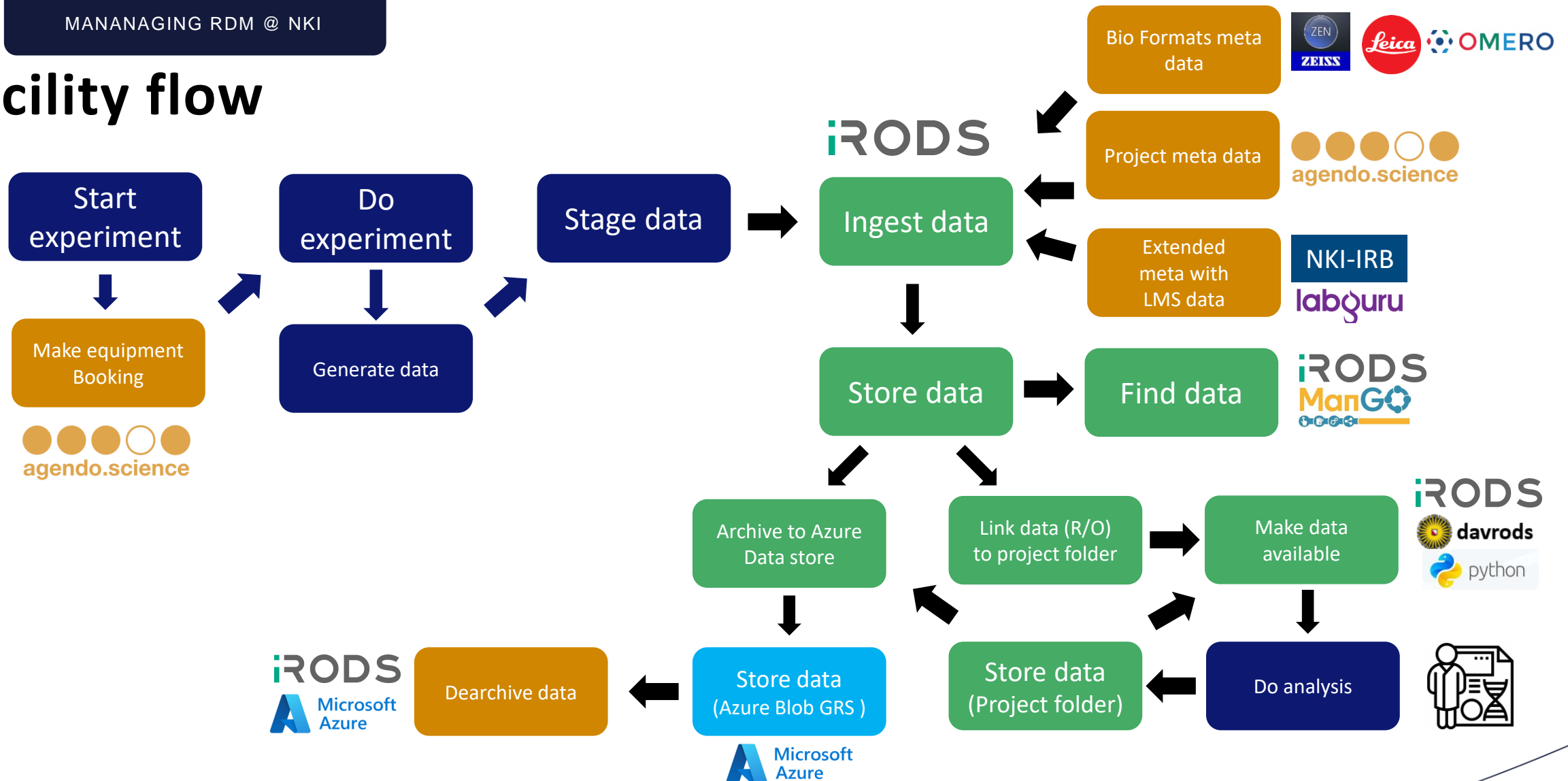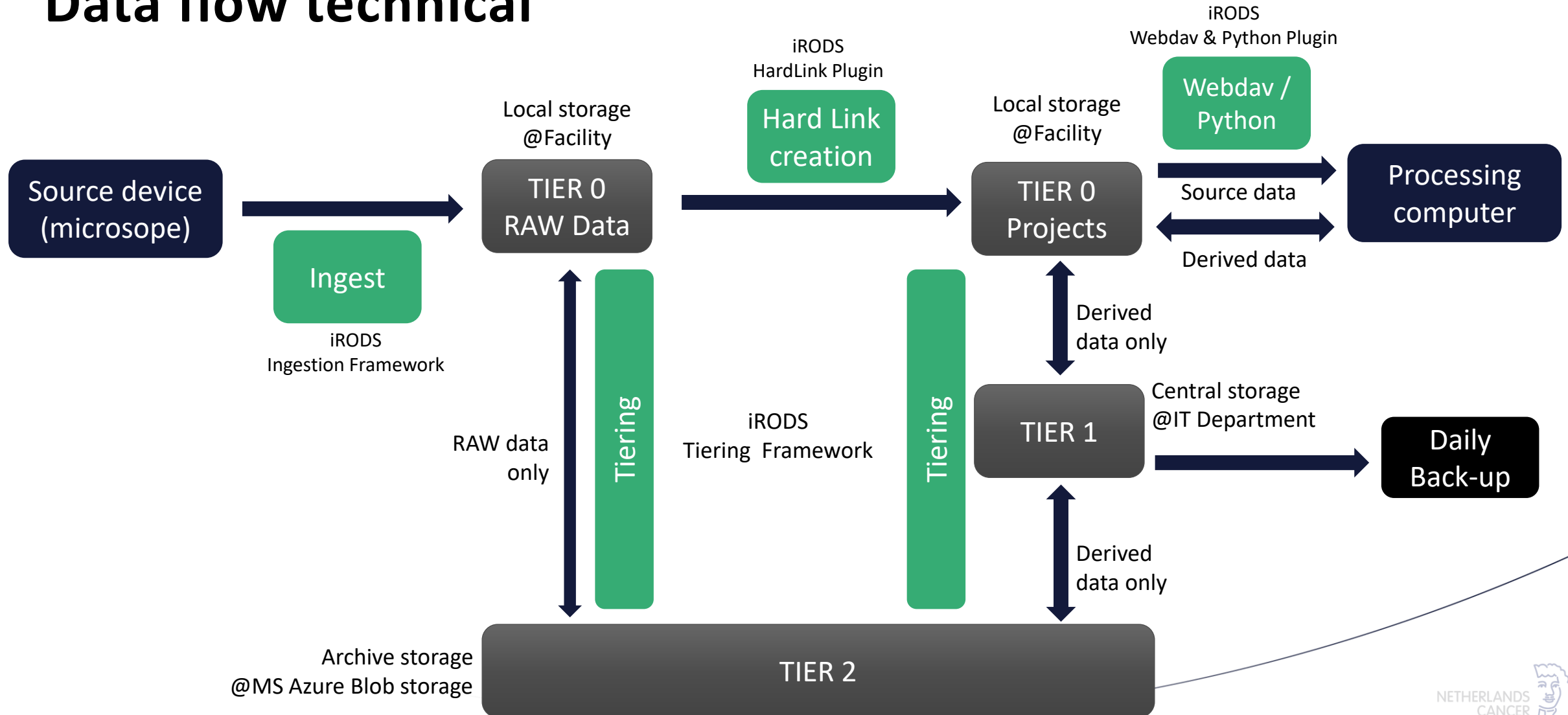
# RDM infrastructure for Bio Imaging facility

# RDM data flow

# Facility flow

# Data flow technical

# RDM Web User interfaces

# NKI Facilities RDM/iRODS Expansion

**Multi Facility**

**Single Facility**

# Clinical Pathology and research with AI/ML

## Current situation
- Pathology slides are scanned at Clinical Pathology department for diagnostic
- Pathology slides are scanned at research facility for research purposes
- Pathology slides are scanned and stored twice, no data carry over/data re-use
- Manual labor involved for older slides (>6month). Must be retrieved out of physical archive.
- No digital archive for clinical pathology department.
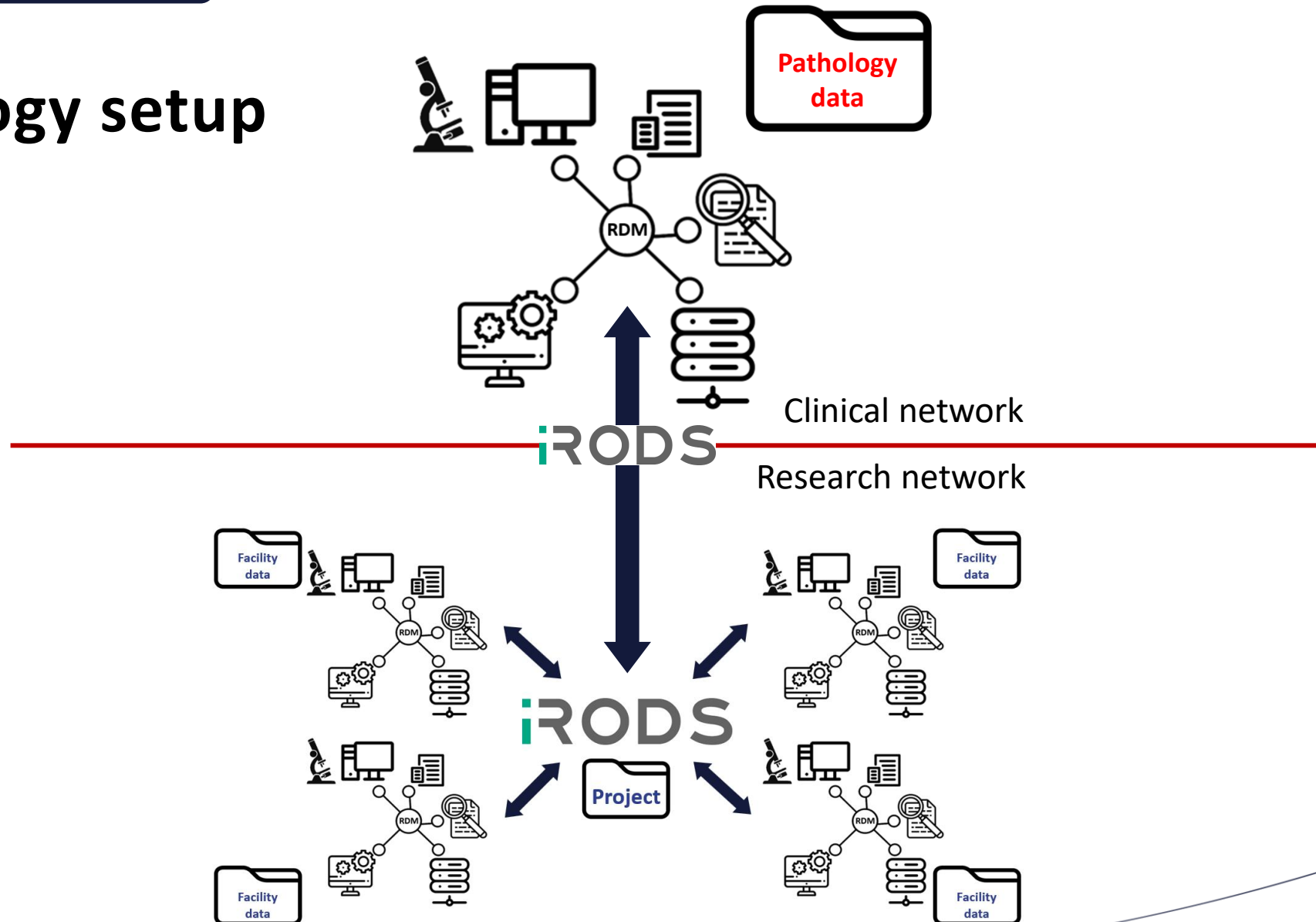- Complexity: clinical network is separated from research network.

## Desired situation
- Pathology slides are scanned at Clinical Pathology department for diagnostic and caried over to research for research purposes.
- Scanned slides are described and stored in long term archive for re-use for both clinic and research.
- Safe data, coordinated and auditable data transfer from clinic to research
- Data (at scale) available for AI/ML research workloads.

## Solution
- Use iRODS for safe and auditable data transport, meta data handling, data storage and managed data access

NETHERLANDS
CANCER
INSTITUTE
ANTONI VAN LEEUWENHOEK

# Pathology setup

# Pathology flow

PROSCIA

Process data in IMS → Diagnostic process → Diagnostic outcome

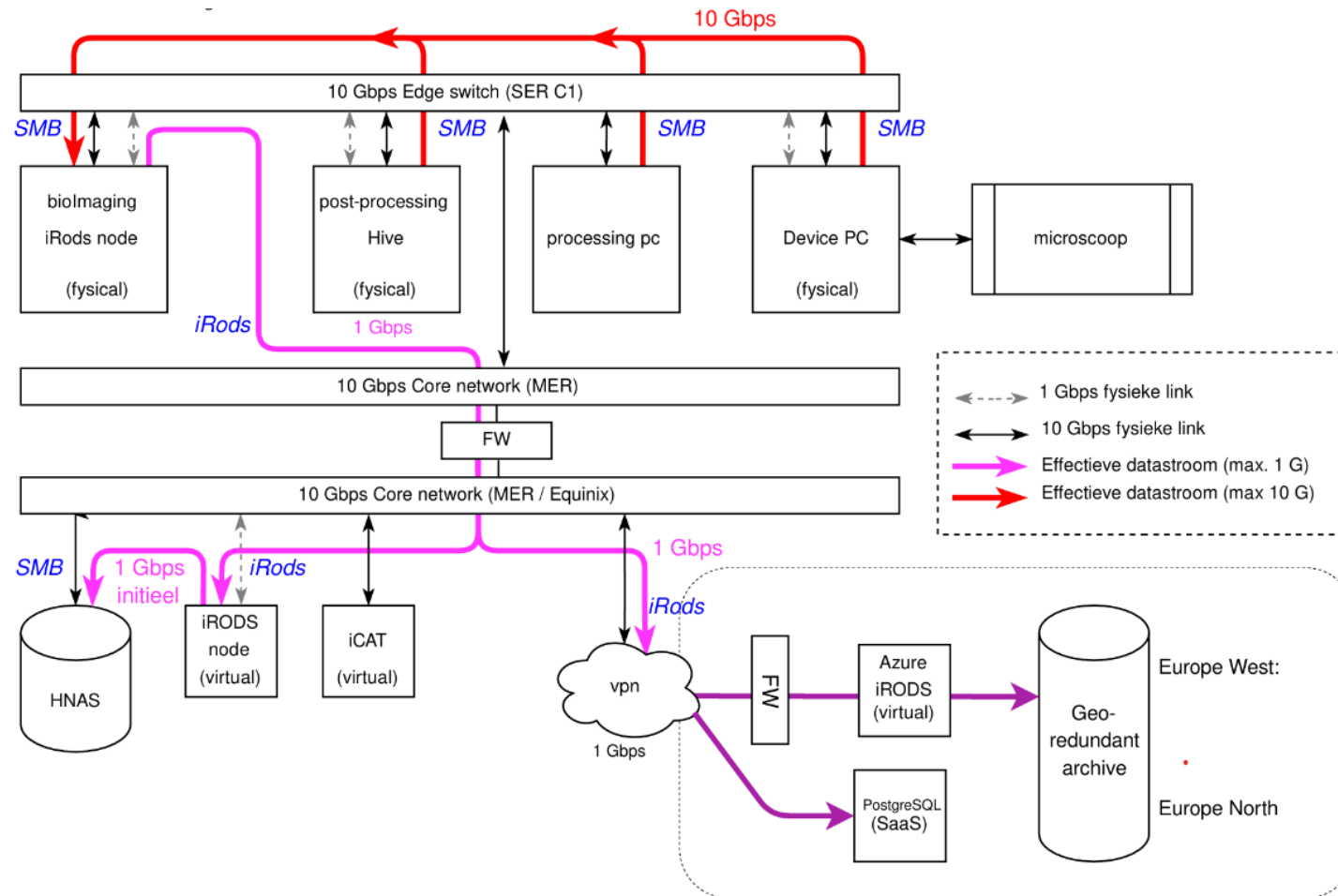Slide creation → Scan slide → Stage slide

Scan slide → Generate data

Stage slide → Data copy to staging areas for IMS and iRODS

Ingest data ← Slide scan meta data

Ingest data ← Extended meta and LMS data

NKI-IRB

Store data → Find data

iRODS ManGO python

Move to Azure Data Lake

Find data ↓

Store data (Azure Blob LRS ) → Get data

Dearchive data ← Store data (Azure Blob LRS )

iRODS Microsoft Azure

Microsoft Azure

Get data → Process data

Get data → Perform research

python

Microsoft Azure

AI

Daily:
~ 500 slides
400-600 GB of data

NETHERLANDS CANCER INSTITUTE
ANTONI VAN LEEUWENHOEK

17

# Basic iRODS network architecture

# Some numbers

- Total departments: 4
  - 2 operational departments
  - 2 for archive purposes
- Total files: 23.131.950
- Total TB: 270 TB
- Total TB in archive: 160 TB (and counting)
- Total unique meta data points: 52.317.828
- Database size: 150 GB

- servers: 7
  - 1 iCAT, 1 Database,
    4 storage nodes, 1 Web (Mango)
- API interfaces: 4
  - Agendo
  - NKI-IRB
  - LMS
  - Azure
- Data access API's: 3
  - iRODS
  - Python
  - WebDav
  - Web (Mango)

NETHERLANDS
CANCER
INSTITUTE
ANTONI VAN LEEUWENHOEK