# Human Technopole
from startup phase to a large-scale research infrastructure

Albino Zamboni

Human Technopole

WORK SHOP GARR 2022

NET MAKERS

# HT Fundation

- The Human Technopole Foundation was established by financial law n. 232, 11 December 2016. The founding members of the Foundation are the **Ministry of Economy and Finance, the Ministry of Health and the Ministry of Education, University and Research** which are responsible for supervising the Foundation.

- The purpose of the Foundation, as indicated in art.1, c. 116 of the above mentioned law, is the **creation of a multidisciplinary scientific and research infrastructure of national interest**, integrated in the fields of health, genomics, nutrition, data and decision science and in the implementation of the Human Technopole scientific and research project ("HT Project").
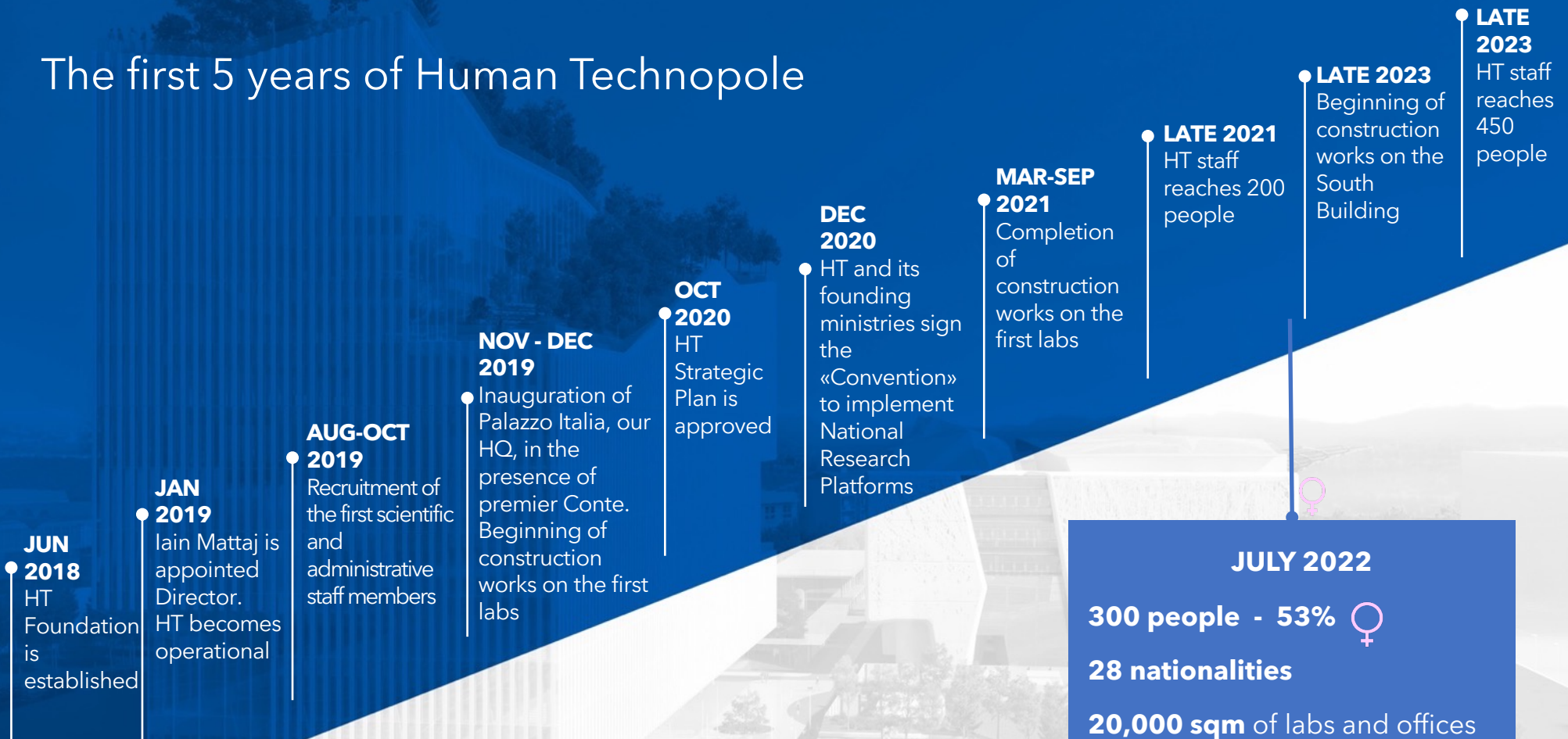
WORK SHOP GARR 2022

NET MAKERS

# Our Mission

- Improve human health and well-being, including healthy ageing.

- Carry out frontier research to improve people's health and well-being.

- Set up and operate a large-scale research infrastructure with interdisciplinary laboratories for the development of precision medicine.

- Act as an open hub to support the growth of the Italian life science research community.

- Engage in industrial cooperation and technology transfer support activities.

- Employ 1,000 scientists including biologists, bioinformatics, chemists, engineers, mathematicians and computer scientists.

# The first 5 Years of HT

## The first 5 years of Human Technopole

**JUN 2018**
HT Foundation is established

**JAN 2019**
Iain Mattaj is appointed Director. HT becomes operational

**AUG-OCT 2019**
Recruitment of the first scientific and administrative staff members

**NOV - DEC 2019**
Inauguration of Palazzo Italia, our HQ, in the presence of premier Conte. Beginning of construction works on the first labs

**OCT 2020**
HT Strategic Plan is approved

**DEC 2020**
HT and its founding ministries sign the «Convention» to implement National Research Platforms

**MAR-SEP 2021**
Completion of construction works on the first labs

**LATE 2021**
HT staff reaches 200 people

**LATE 2023**
Beginning of construction works on the South Building

**LATE 2023**
HT staff reaches 450 people

### JULY 2022

**300 people - 53%** ♀

**28 nationalities**

**20,000 sqm** of labs and offices

# HT Today

- **63%** of the scientific team from **international institutions**.

- **70 Italians back from abroad**.

- **Fast scientific recruitment** (1 scientist per week for the next 3 years).

- **8.5-million-euro grant** (Jan. '21- Mar. '22).

- **20,000 sqm** of labs and offices.

*Cryo-Electron Microscopy*
*Light Imaging & Image Analysis*
*Genomics*
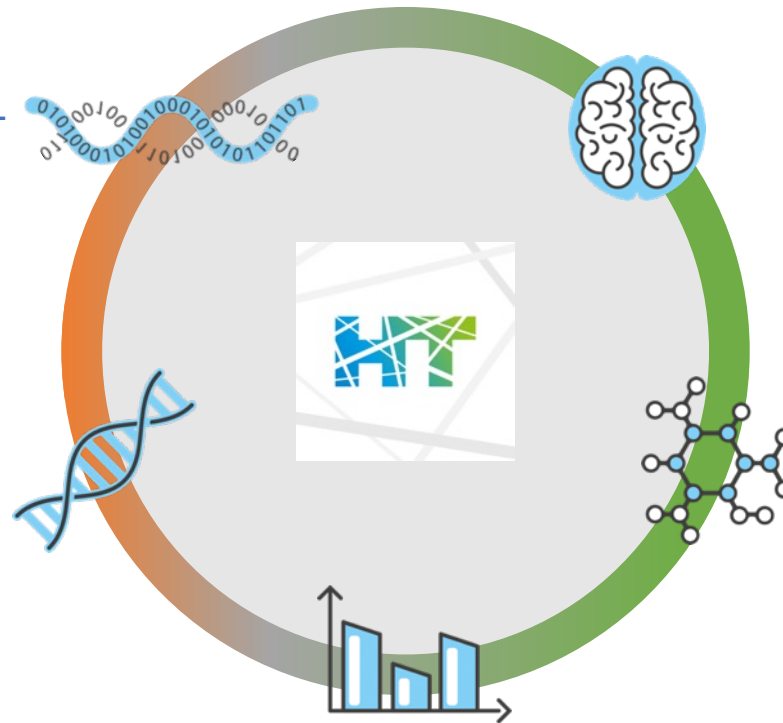*Data Centre*

# Our Research Centers



OUR LINES OF RESEARCH

**GENOMICS**
FUNCTIONAL
POPOULATION & MEDICAL

**NEUROGENOMICS**

**COMPUTATIONAL BIOLOGY**

**STRUCTURAL BIOLOGY**

**HEALTH DATA SCIENCE**

WORK SHOP GARR 2022

NET MAKERS

# Our Research Centers

Our biomedical research aims at developing predictive and  personalised medicine to treat cancer, cardiovascular and neurodegenerative diseases.



**Genomics**

The Centre studies genomics characteristics and traits to identify how hereditable genetic information is shared in view of identifying more personalised treatments



**Neurogenomics**

The Centre studies neuropsychiatric and neurological diseases to probe the structure, function and development of the nervous system.



**Structural biology**

The Centre aims at gaining precise knowledge of the structure of macromolecules, a fundamental step in understanding the function of cells.

# Our Research Centers

Our biomedical research aims at developing predictive and  personalised medicine to treat cancer, cardiovascular and neurodegenerative diseases.



**Computational biology**

The Centre develops solutions for the analysis, management and integration of data produced by other Centres, making it available to the wider scientific community.
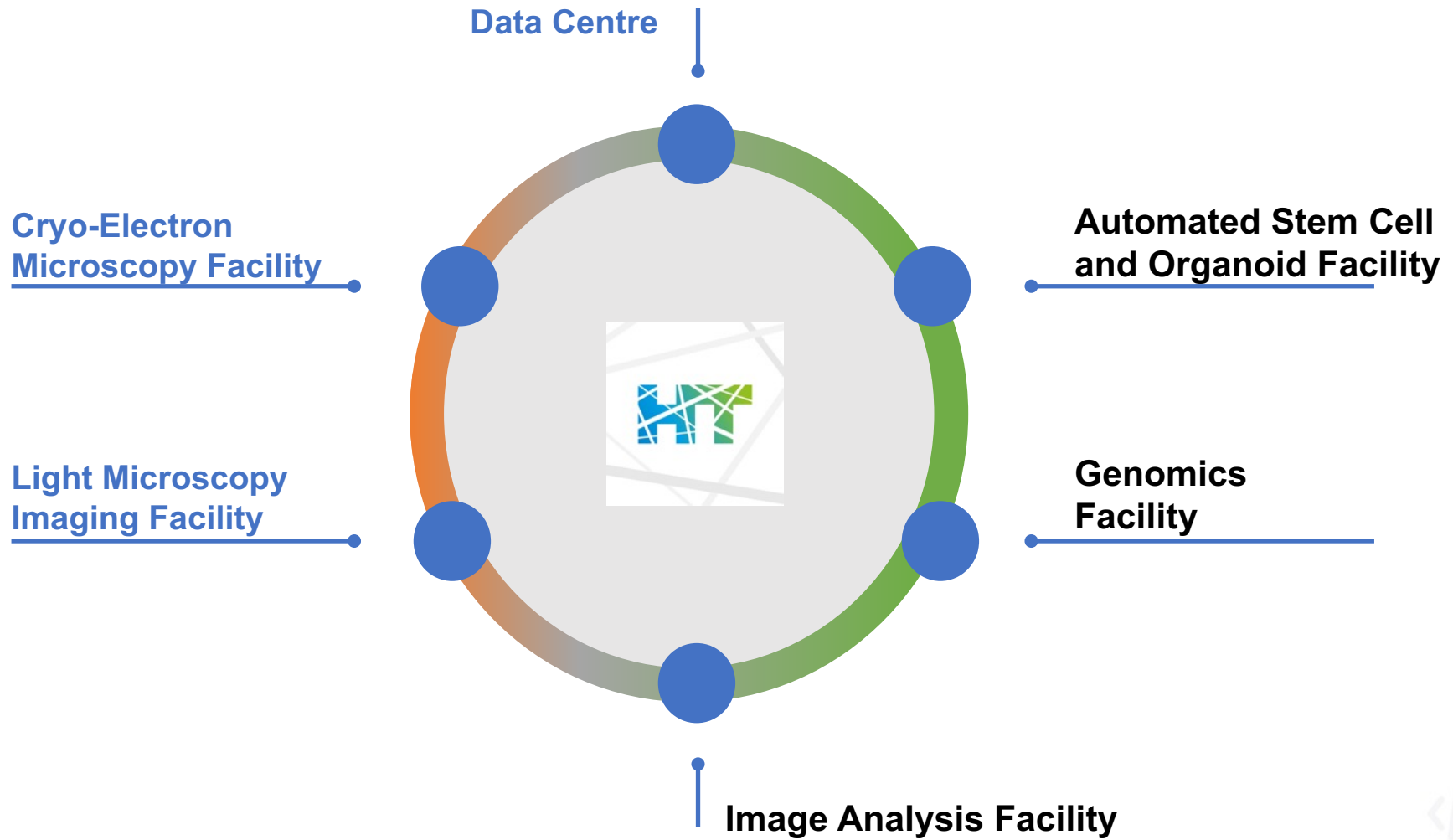


**Health Data Science**

The Centre analyses clinical and socio-economic data to provide advice to different stakeholders, in particular policymakers, mainly to the national health system.

# Our Facilities



Data Centre

Cryo-Electron Microscopy Facility

Light Microscopy Imaging Facility

Automated Stem Cell and Organoid Facility

Genomics Facility

Image Analysis Facility

# Our Facilities

HT is a national hub and centre of reference for life science research.
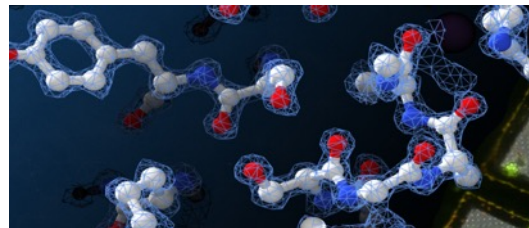
Our facilities are available to HT scientists and researchers as well as to the external scientific community who will access them through open selection procedures based on merit.
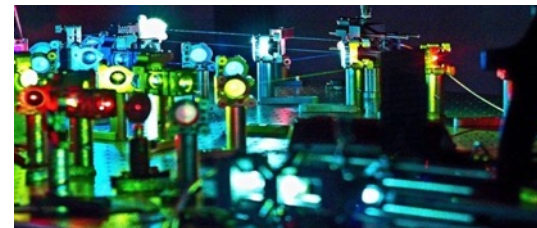
**Six research facilities:**



**Genomics**
Large-scale DNA/RNA sequencing infrastructure to conduct population studies and support national screening initiatives.



**Cryo-Electron Microscopy**
Italy's most comprehensive CryoEm infrastructure: five state of the art microscopes to freeze molecules and observe them at atomic level.
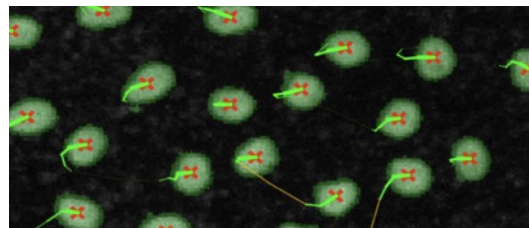


**Light Imaging**
With a focus on 3D imaging it will photograph rare, dynamic and constantly evolving processes.
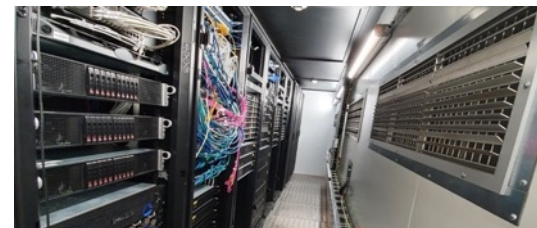


**Automated Stem Cell and Organoid Facility**
It will engage in cell re-programming, genome editing and organoid culture.



**Image Analysis**
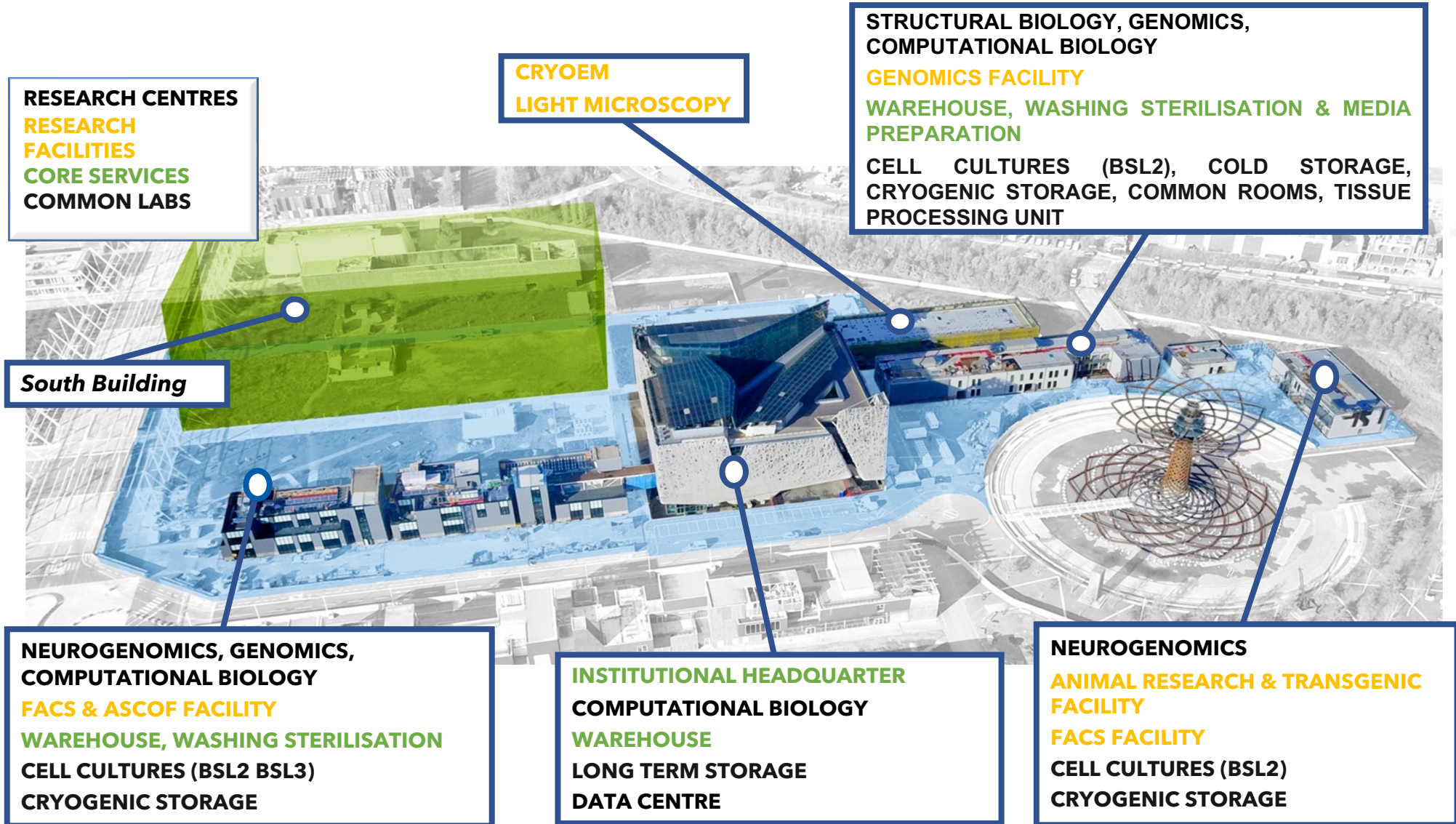Solutions for for image restoration, real-time image analysis, big data management and visualization.



**Data Centre**
High storage and computing capacity to support researchers in the storage and analysis of huge amounts of data.

# Our Campus

Ongoing works to build a large scale research infrastructure



**RESEARCH CENTRES**
**RESEARCH FACILITIES**
**CORE SERVICES**
**COMMON LABS**

**CRYOEM**
**LIGHT MICROSCOPY**

**STRUCTURAL BIOLOGY, GENOMICS, COMPUTATIONAL BIOLOGY**
**GENOMICS FACILITY**
**WAREHOUSE, WASHING STERILISATION & MEDIA PREPARATION**
**CELL CULTURES (BSL2), COLD STORAGE, CRYOGENIC STORAGE, COMMON ROOMS, TISSUE PROCESSING UNIT**

*South Building*

**NEUROGENOMICS, GENOMICS, COMPUTATIONAL BIOLOGY**
**FACS & ASCOF FACILITY**
**WAREHOUSE, WASHING STERILISATION**
**CELL CULTURES (BSL2 BSL3)**
**CRYOGENIC STORAGE**

**INSTITUTIONAL HEADQUARTER**
**COMPUTATIONAL BIOLOGY**
**WAREHOUSE**
**LONG TERM STORAGE**
**DATA CENTRE**

**NEUROGENOMICS**
**ANIMAL RESEARCH & TRANSGENIC FACILITY**
**FACS FACILITY**
**CELL CULTURES (BSL2)**
**CRYOGENIC STORAGE**

# South Building



**Once completed, it will host labs for 800 scientists as well as offices, event spaces, workshops, and training courses.**

The winning project has been awarded in April 2020, and the building, which is to be constructed, will be ready in 2027.
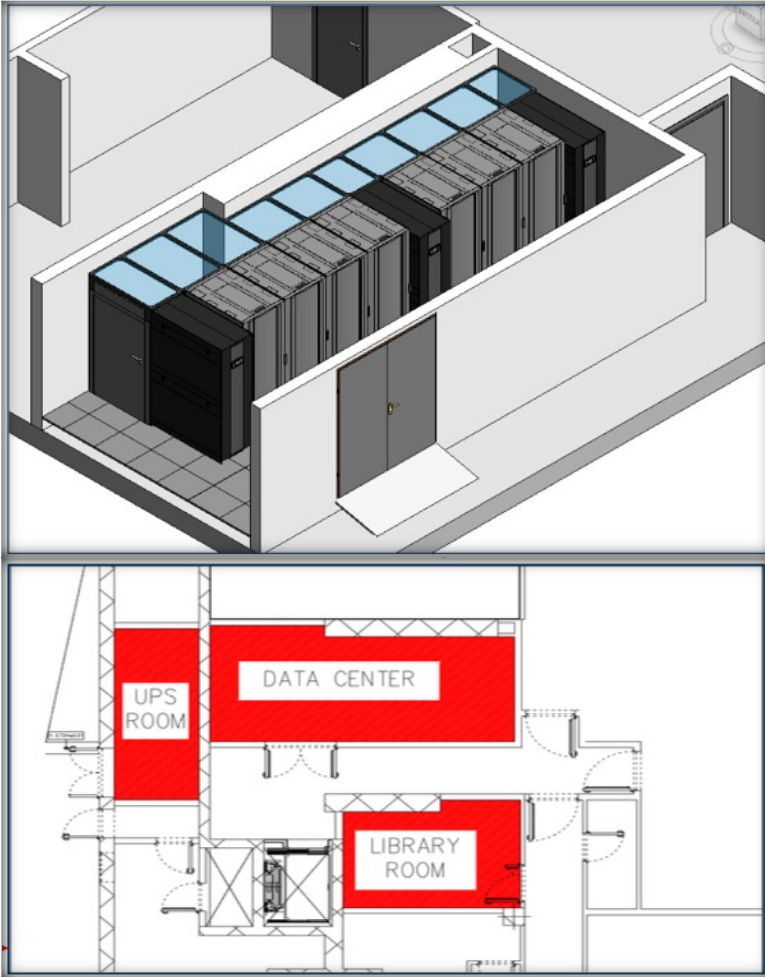
# Our Campus

## IT Data Center Container

# Our Campus
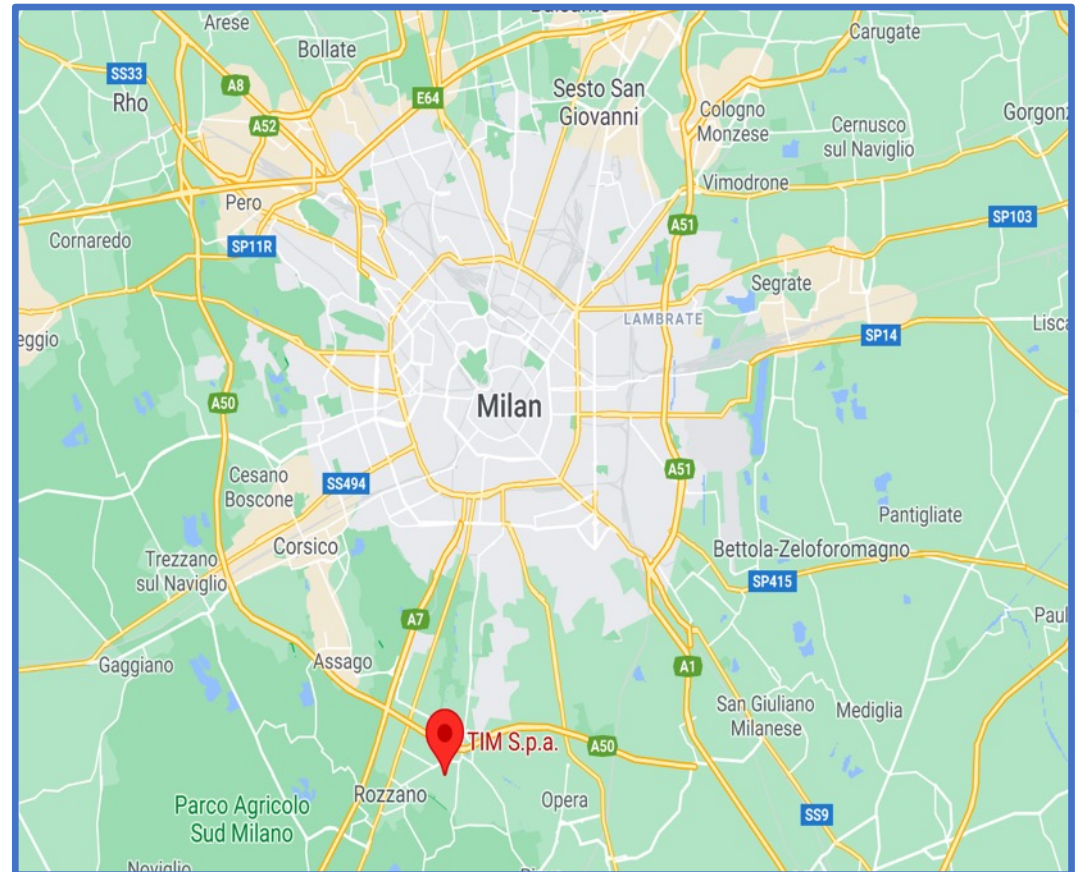
Extending resources of External shelter

# Our Colocation

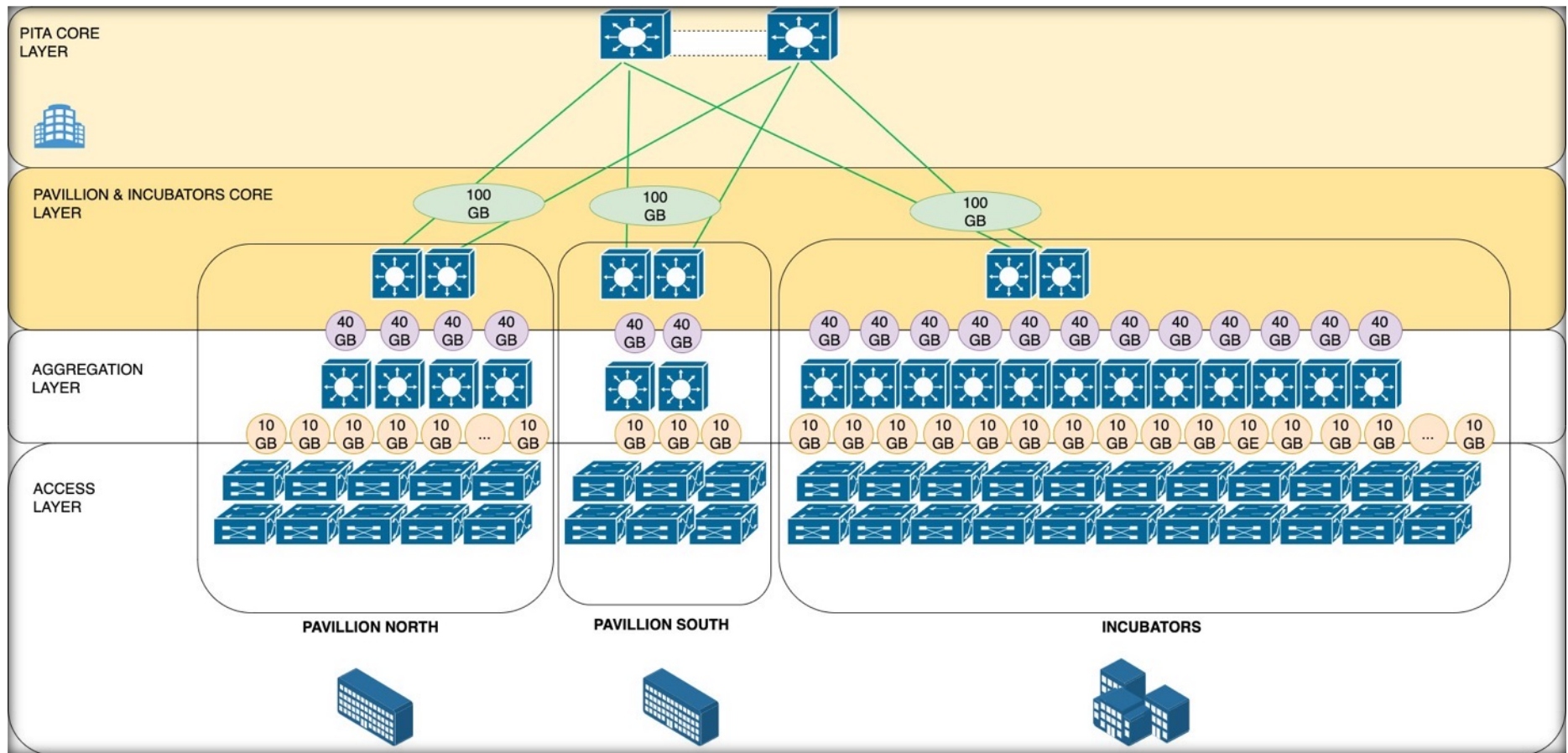## Offsite backup and data processing

**Redundant services:**

- **HPC**
  - 12 compute nodes
  - 1 PB attached storage
- **Virtualisation**
  - 3 nodes
  - 30 TB Fiber Channel storage
- **Storage:**
  - 1PB central scientific storage
- **Connection:**
  - 10Gbps DWDM redundant

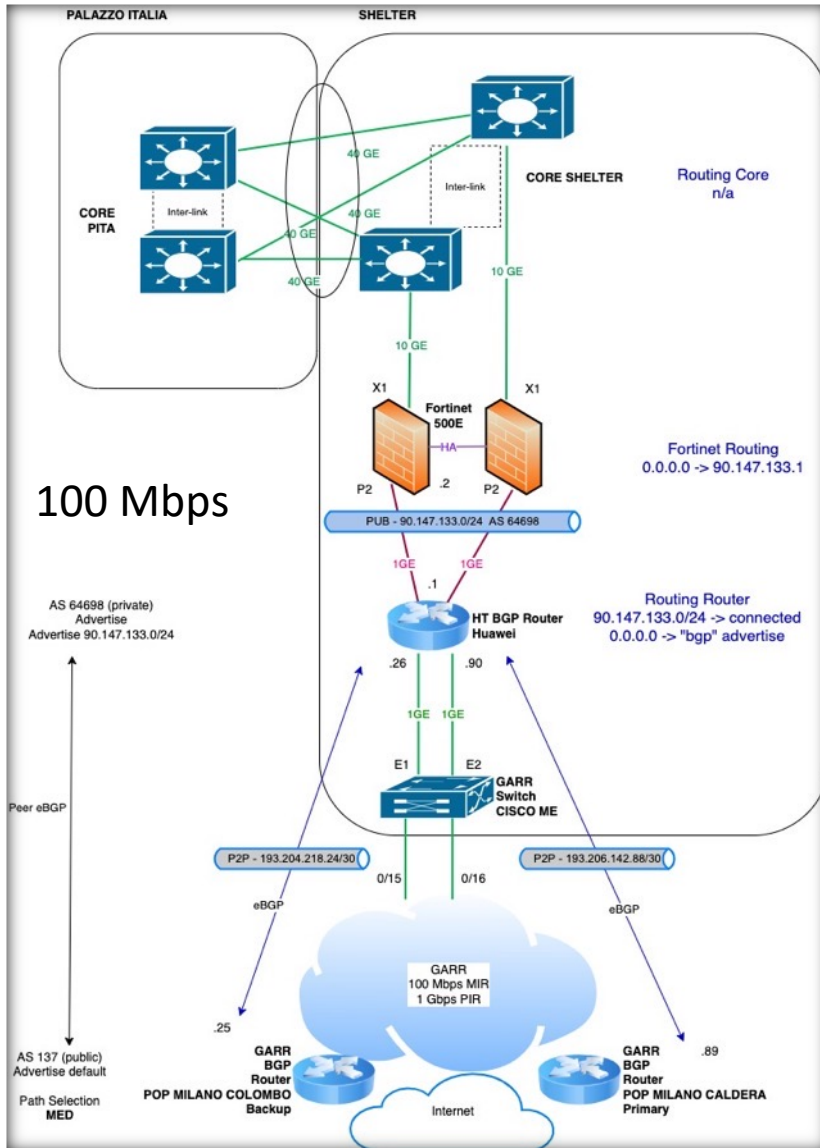**Low latency network connection**

# Our Campus LAN connections
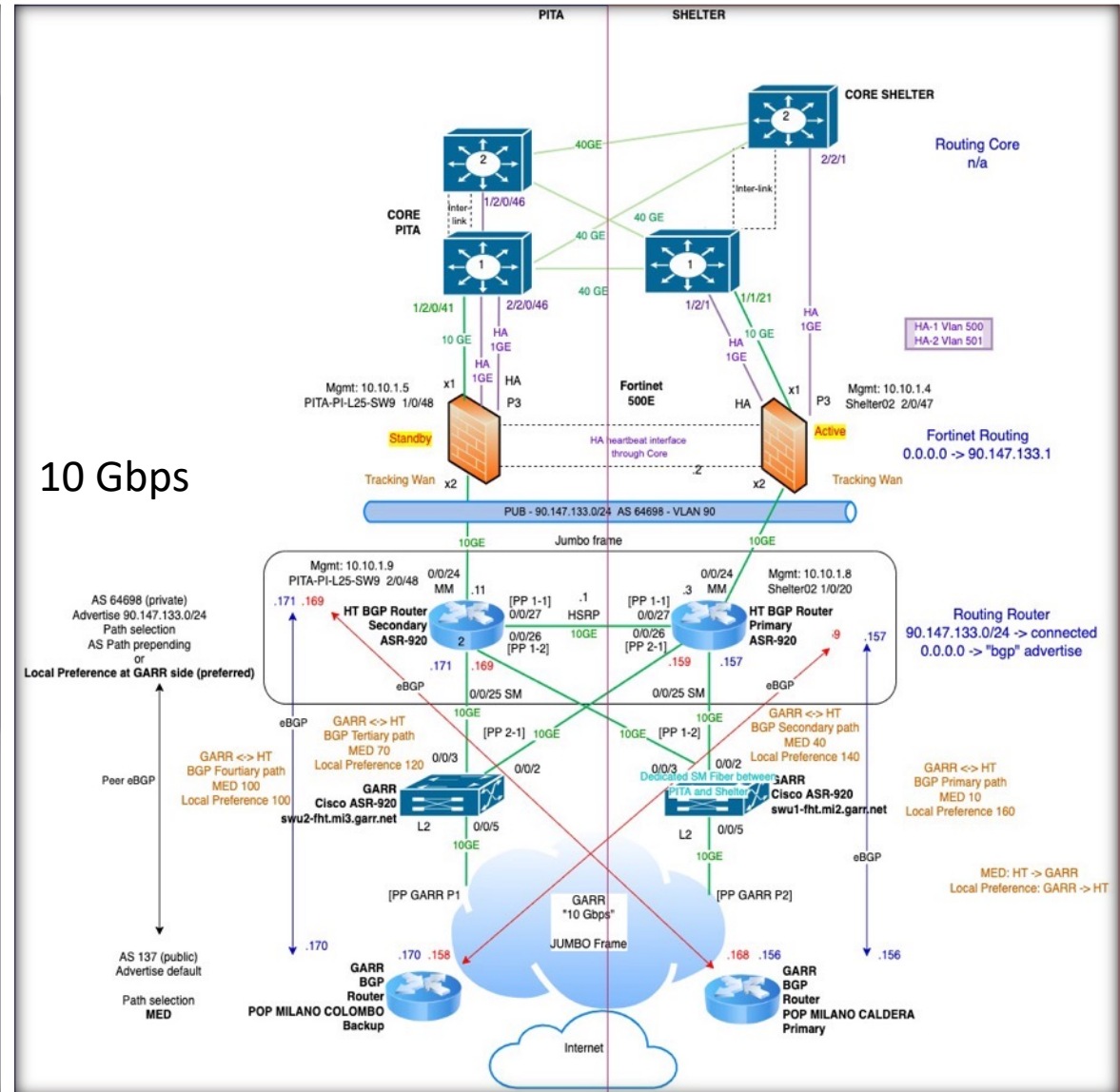
x 100Gbps inter-buildings connections

# Our Campus WAN connections

2021

2022

100 Mbps

10 Gbps

# Our Network Topology

Spine-Leaf approach

Virtual Port-Channels

Low latency switches

# High Performance Computing

**DELL HPC**

**Usage:** Computational processing

**Operating System:** CentOS 8

**Scheduler:** SLURM

**Cluster Interconnect :** 100 Gbps

| | **25 CPU Nodes** | **10 GPU Nodes** | **5 FAT Nodes** |
|---|---|---|---|
| CPU GHz | 2.9 | 2.2 | 2.2 |
| N° Core | 32 | 36 | 36 |
| RAM | 576 GB | 576 GB | 1.1 TB |
| GPU | | 4 x V100 | 4 x V100 |
| Total CPU Core Count | | | 1340 |
| Total RAM | | | 26 TB |
| **Calculated performance** | | | **150 TFLOPS** |

**Available Queues:**
- **CPU Standard** - CPU nodes (CPU Intensive)
- **GPU Standard -** GPU nodes (GPU intensive)
- **High Intensive -** Specially allocated nodes (Memory or GPU Intensive)

**Access:** 2 x Login Nodes

# Storage

## HPC access

**High Performance storage**

**Usage:** High Performance IO Operations

**Filesystem:** Parallel BeeGFS with HA

**Capacity:** 2.1 PB

**Throughput:** 18 GB/s with **sequential** R/W

**Cluster interconnect:** 100 Gbps Infiniband HDR

# Virtualisation

## Traditional Virtualisation + High Performance GPU-Driven VDI

| Core Services | |
|---|---|
| Physical Servers | 5 |
| Cores | 560 |
| Memory (GB) | 3840 |
| Storage (TB) | 30 |

| VDI - Double Precision GPUs - Intel Based | |
|---|---|
| Physical Servers | 5 |
| Cores | 320 |
| GPUs (NVDIA V100s) | 20 |
| Memory (GB) | 3840 |
| Storage (TB) | 30 |

Q1 2023

| VDI - Single Precision GPUs - AMD Based | |
|---|---|
| Physical Servers | 2 |
| Processors | 64 |
| GPUs (NVDIA A40) | 12 |
| Memory (GB) | 4096 |
| Storage (TB) | 30 |

| VDI - Single Precision GPUs - AMD Based | |
|---|---|
| Physical Servers | 12 |
| Processors | 384 |
| GPUs (NVDIA A40) | 72 |
| Memory (GB) | 24576 |
| Storage (TB) | 180 |

## Mellanox low latency switches

- 300nsec for 100GbE port-to-port          - Flat latency across L2 and L3 forwarding

# Storage

## Central Scientific Storage

**Usage:** Tier 1 Storage for Group Shares and User Homes

**Protocols:** NFSv4, SMB 3.1, and S3

**Group Share (Default) Size:** 10TB

**Home Share (Default) Size:** 200GB

**Backup Policy:**
- Hourly Snapshots – 1 Day Retention
- Daily Snapshots – 1 Week Retention
- Weekly Snapshots – 1 Month Retention

**RAW** Central Storage Total Capacity

# 9PB

Available on Campus

# 7PB approx

8/52 nodes moved to colocation

# > 2PB

# Backup

## Second line of defence for recovery

**Backup System:** Bacula Enterprise

**Capacity:** 2.5 PB
**Technology:** Disk Backup
**Retention:**

- **Group Shares:** 1 year
- **VMware VMs:** 3 months
- **Laptops:** 1 month
- **Office 365:** 1 month



**Upgrade Plans (2022 - 2023):** Long Term Archive (10 years) on Tape
    **Initial deployment plan capacity:** 60 PB based on LTO8
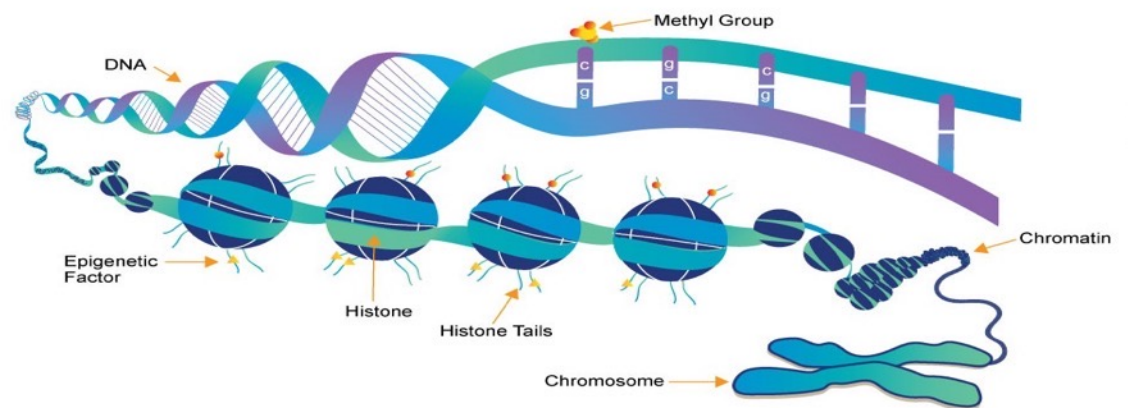    **Scale as needed strategies:**
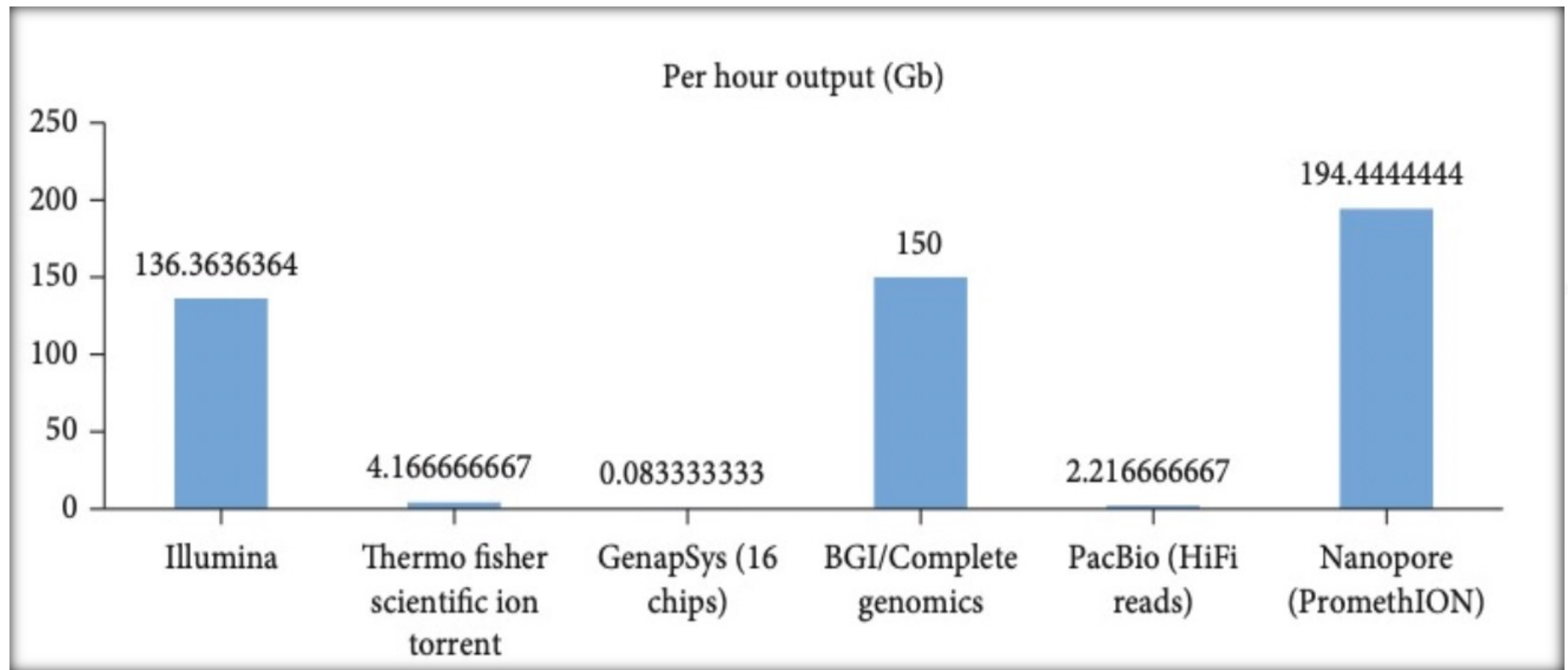- New Expansion shelf
- Tape Generation

# Genomic-Sequencers

New sequencers generation:

- High Throughput - Gbps

- Long run experiment - hours

- High cost run – tens of k €

- Cloud AI monitoring

# Genomic-Sequencers

Sequencers can produce up tp 200 Gbyte data/h

# Genomics-Sequencers

Experiments can take up to 72 run hours and produce 14 Tbyte data

| Manufacturer | Read length | Data output | Max. run time (hours) | Chemistry | Key applications** |
|---|---|---|---|---|---|
| Illumina (NovaSeq 6000) | 300 PE | 6 Tb (6000 Gb) | 44 | Sequencing by synthesis | SS-WGS and TGS, TGEP, 16sMGS, WES, SCP, LS-WGS, CA, MS, MGP, CFS, LBA |
| Thermo Fisher Scientific Ion Torrent (Ion GeneStudio S5 Prime) | 600 SE | 50 Gb | 12 | Sequencing by synthesis | WGS, WES, TGS |
| GenapSys (16 chips) | 150 SE | 2 Gb | 24 | Sequencing by synthesis | TS, SS-WGS, GEV, 16S rRNA sequencing, sRNA sequencing, TSCAS |
| QIAGEN (GeneReader) | 100 SE | Not available | Not available | Sequencing by synthesis | Cancer research and identifying mutations |
| BGI/Complete Genomics | 400 SE | 6 Tb (6000 Gb) | 40 | DNA nanoball | Small and large WGS, WES and TGS |
| PacBio (HiFi Reads) | 25 Kb | 66.5 Gb | 30 | Real-time sequencing | DN sequencing, FT, identifying ASI, mutations, and FPM |
| Nanopore (PromethION) | 4 Mb | 14 Tb (14000 Gb) | 72 | Real-time sequencing | SV, GS, phasing, DNA and RNA base modifications, FT, and isoform detection |

# Genomics-Sequencers

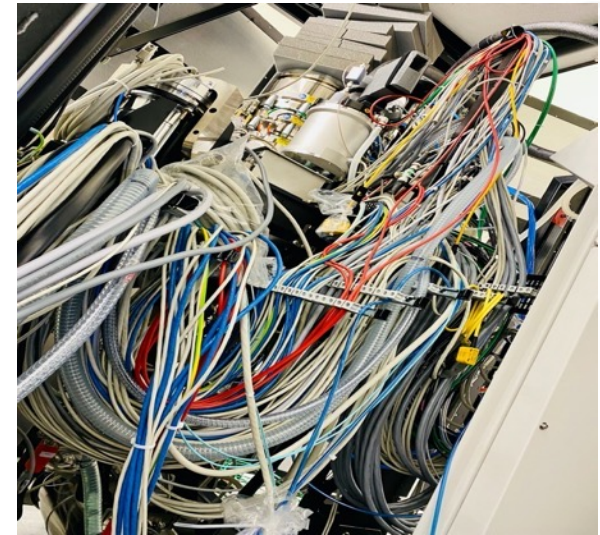A single sequencer can produce up to 676 Tbyte data/year

| Instrument | WEEK | MONTH | YEAR |
|---|---|---|---|
| NovaSeq 1 | 14 Tb | 56 Tb | 672 Tb |
| NovaSeq 2 | 14 Tb | 56 Tb | 672 Tb |
| MiSeq | 30 Gb | 120 Gb | 1.4 Tb |
| NextSeq | 840 Gb | 3.3 Tb | 40.3 Tb |
| PromethION | 1.1 Tb | 4.5 Tb | 56 Tb |

# Cryo Electron Microscopes

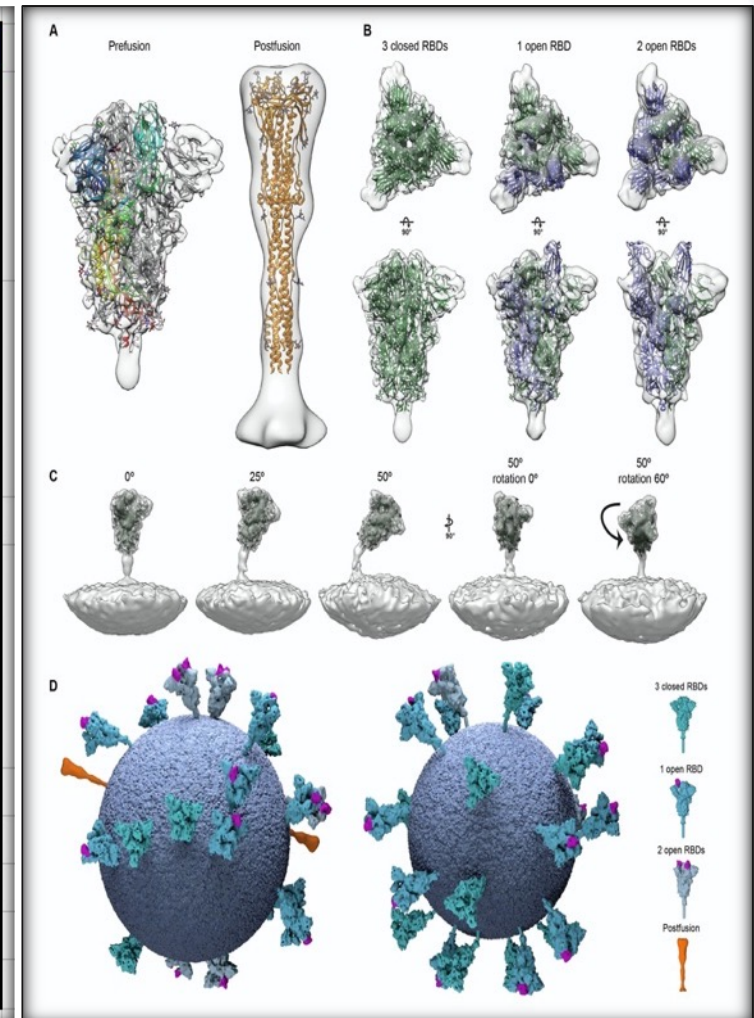To get high-resolution images we need to freeze and keep the sample at -175 °C

- **Cryo-electron microscopy** allows to photograph the inner section of a frozen cell and individual molecules at high resolution

- We **freeze** the sample (cells, enzymes, DNA, viruses, etc.) to preserve them and get better images

- It'is a completely new level of **observation**, this is why the technique was awarderd the 2017 Nobel Prize for Chemistry

# Cryo Electron Microscopes

Microscopes can produce up to 15 Gbyte image data/minute

| Instrument | Maximum Throughput (MB/s) |
|---|---|
| Krios | Min: 1 image (500MB to 1.5GB) per minute.<br>Max: 10 image (500MB to 1.5GB) per minute.<br>NO intermediate files.<br>Usual dataset size: 1-4 TB. |
| Glacios | Min: 1 image (500MB to 1.5GB) per minute.<br>Max: 10 image (500MB to 1.5GB) per minute.<br>NO intermediate files.<br>Usual dataset size: 1-4 TB. |
| Spectra | Max: 1 image per second (64 MB/s) |
| Tundra | Min: 1 image (500MB to 1.5GB) per minute.<br>Max: 10 image (500MB to 1.5GB) per minute.<br>NO intermediate files.<br>Usual dataset size: 1-4 TB. |
| Talos | Max: 1 image per second (64 MB/s) |
| Leica Stellaris | Max: 1 image per second (32 MB/s) |
| Leica Thunder | Max: 1 image per second (32 MB/s) |
| Aquilos 2 | Max: 1 image per second (32 MB/s) |
| Arctis | Max: 1 image per second (64 MB/s) |

# Cryo Electron Microscopes

A single microscope can produce up to 2 PB image data/year

| Instrument | TB/year |
|---|---|
| Krios | 770 |
| Glacios | 770 |
| Krios | 770 |
| Spectra | 1.925 |
| Tundra | 770 |
| Talos | 1.925 |
| Leica Stellaris | 962 |
| Leica Thunder | 962 |
| Aquilos 2 | 1.925 |
| Arctis | 1.925 |

WORK SHOP GARR 2022 NET MAKERS

# Inter-Institutes connections

**Institutes/research centers:**

•**IEO**

Istituto Europeo di Oncologia (ieo.it)

•**IRCCS** (research hospitals in Italy)

https://www.hsantalucia.it/en/irccs

•**CRG**

Centre for Genomic Regulation Website (crg.eu)

•**MRC**

MRC Laboratory of Molecular Biology (cam.ac.uk)

•**NCBI**

National Center for Biotechnology Information (nih.gov)

•**EBI**

The European Bioinformatics Institute < EMBL-EBI

•**EMBL**

Heidelberg | EMBL.org

**Protocols:**

•Globus

•sftp (ssh)

•Ncftp

•ncftpput

Whole Genome Sequencing (**WGS**) files size to transfer:

**1.8 PB**

17 days at 10Gbps

# What's next!

## More resources

We have to be ready to host:

- Tens of Genomics Sequencers

- Tens of Cryo-EM Microscopes

- More WGS and other scientific data to transfer

- More scientists and administrative employees (1000 – 2000)

**Upgrade** →

- More computing: HPC, GPU VDI

- More storage: Tens PetaByte

- More WAN bandwith: 40 / 100 Gbps

- More Security and Data Governance

# Fundamental deployment

**It has to be rock solid!**

**Stability**

– Scientist work 24/7

and if they don't, many of their long running jobs do!

–Some researchers run irreproducible experiments

and some are expensive!

**Speed**

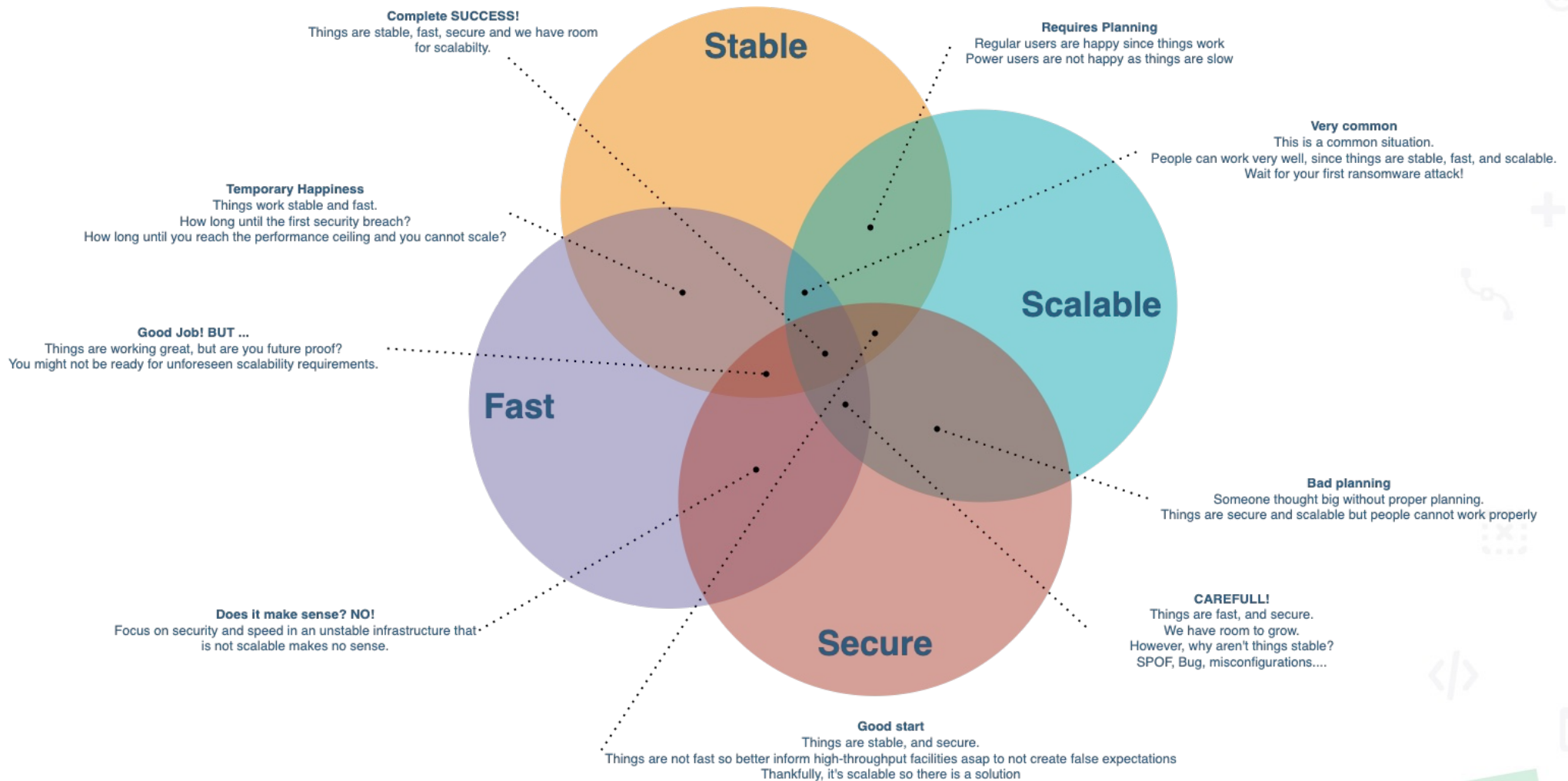– Some facilities require on top of stability throughput

**Scalability**

– Prepare for what's coming

– Are we future-proof? Is there room for improvement / growth?

Stable

Fast

Temporary happiness

Success

Requires planning

instability is unacceptable

Scalable

# Let's add a bit of complexity

## What about Security?

# CyberSecurity

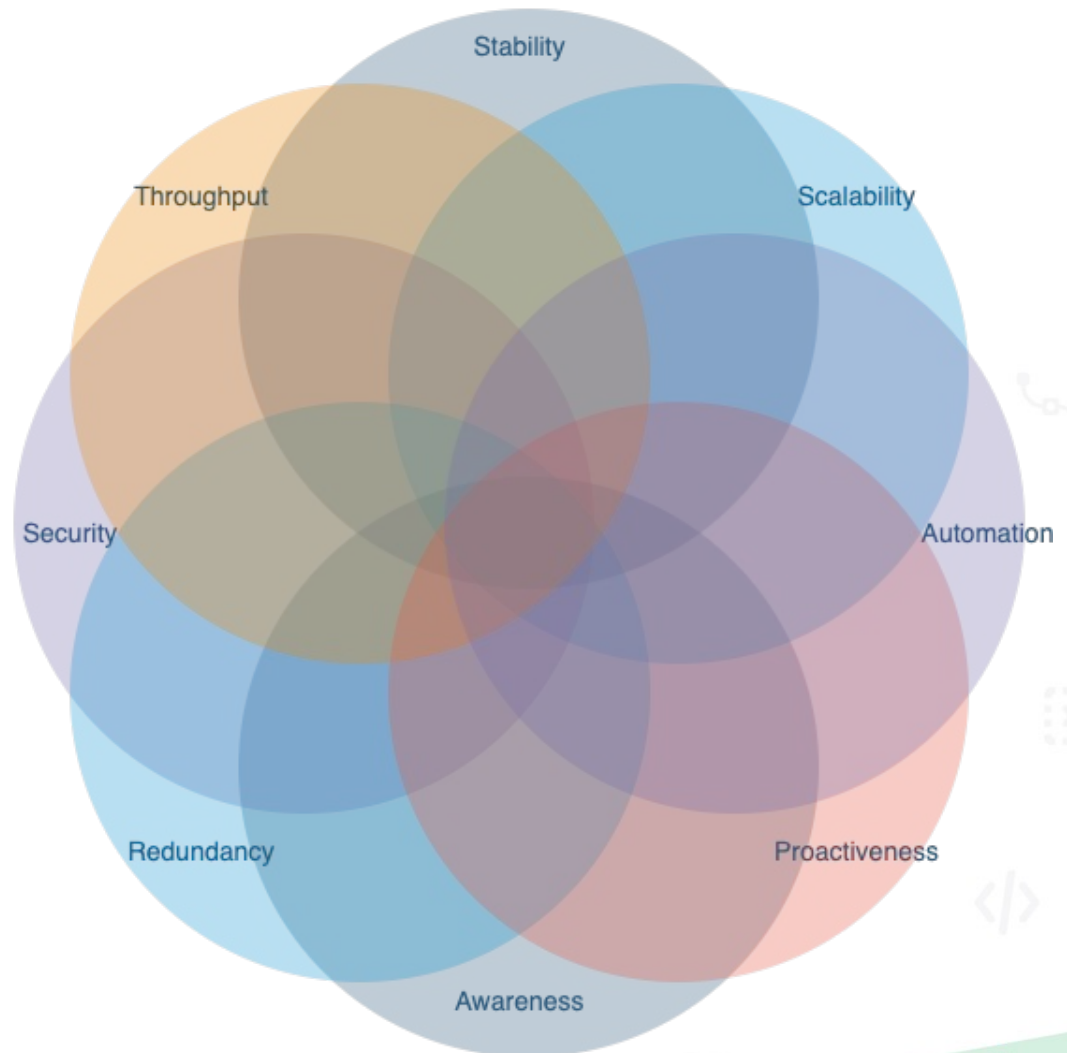| | | |
|---|---|---|
| **Network Segmentation** | **NAC** | **Multi Factor Authentication** |
| **Web Filtering** | **EDR-XDR** | **MDM** |
| **Vulnerability Assessment** | **Anti phishing Campaign** | **Patch Management** |

# Data Governance

✓ **STRUCTURED APPROACH TO AQUIRE EXTERNAL DATA**

   **(HEALTHCARE DATASETS FROM INSTITUTIONAL DATA PROVIDERS)**

✓ **GDPR  COMPLIANCE - PROCESSING OF PERSONAL DATA**

✓ **STRUCTURED PROCESS TO MONITOR AND MANAGE DATA AND ITS**

   **QUALITY**

✓ **STRUCTURED DATA ACHITECTURE ENABLING FASTER AND BETTER**

   **ANALYSIS**

✓ **REDUCTION OF THE RISKS ASSOCIATED WITH THE MANAGEMENT OF**

   **TREATMENTS CONTRARY TO THE LAW**

# What about more complexity

Redundancy, monitoring, alerting, automation… you name it.
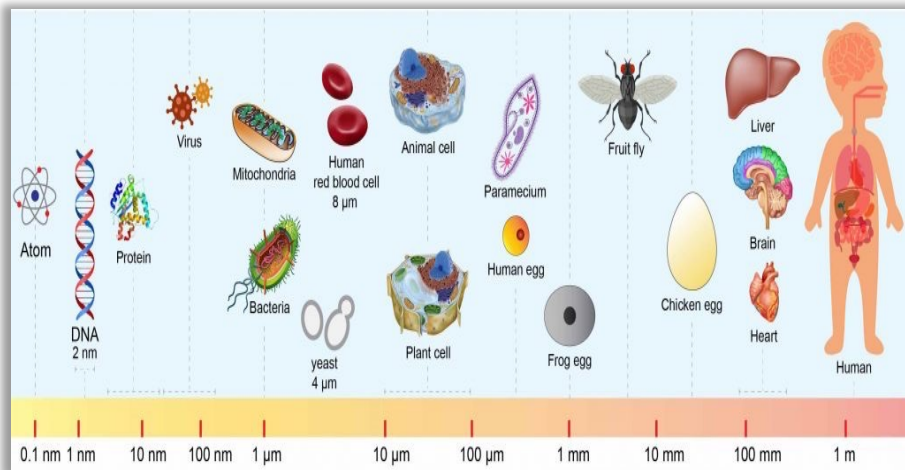
$$\sum_{i=2}^{n} \frac{n!}{i!\,(n-i)!} = 2^n - n - 1$$

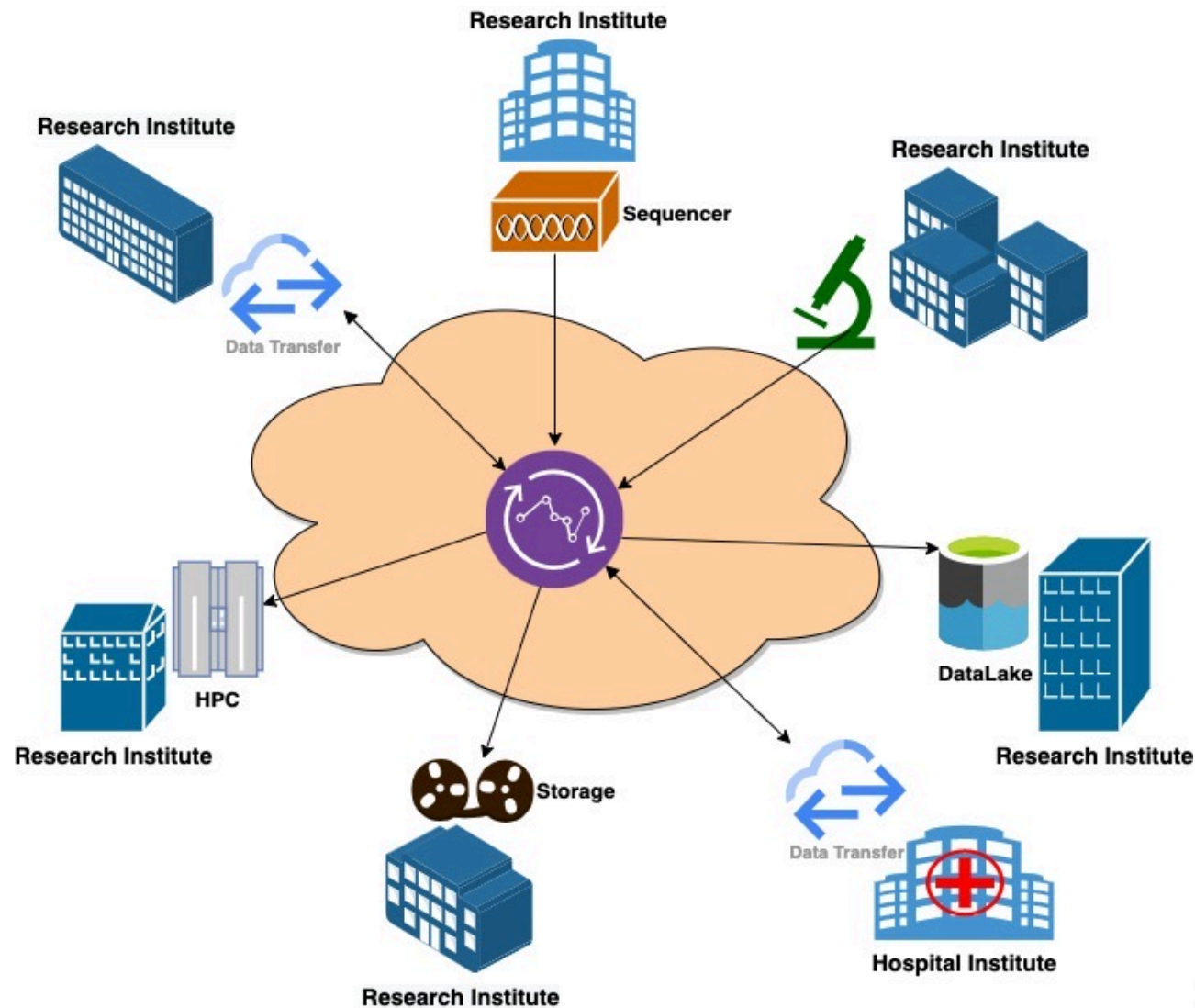| Variables | Overlapping regions |
|:---:|:---:|
| 3 | 4 |
| 4 | 11 |
| 5 | 26 |
| 6 | 57 |
| 7 | 120 |
| 8 | 247 |

# Research Institutes WAN networks

Research Institutes WAN networks are used for:

- Normal business operations including email, web browsing, O365, VPN SSL, SaaS among others. The network must also be built with security features.

- The scientific research process as scientists depend on this infrastructure to share, store, and analyse research data from many different external sources.

- Networks optimized for business operations are neither designed for nor capable of supporting the data movement requirements of data intensive science.

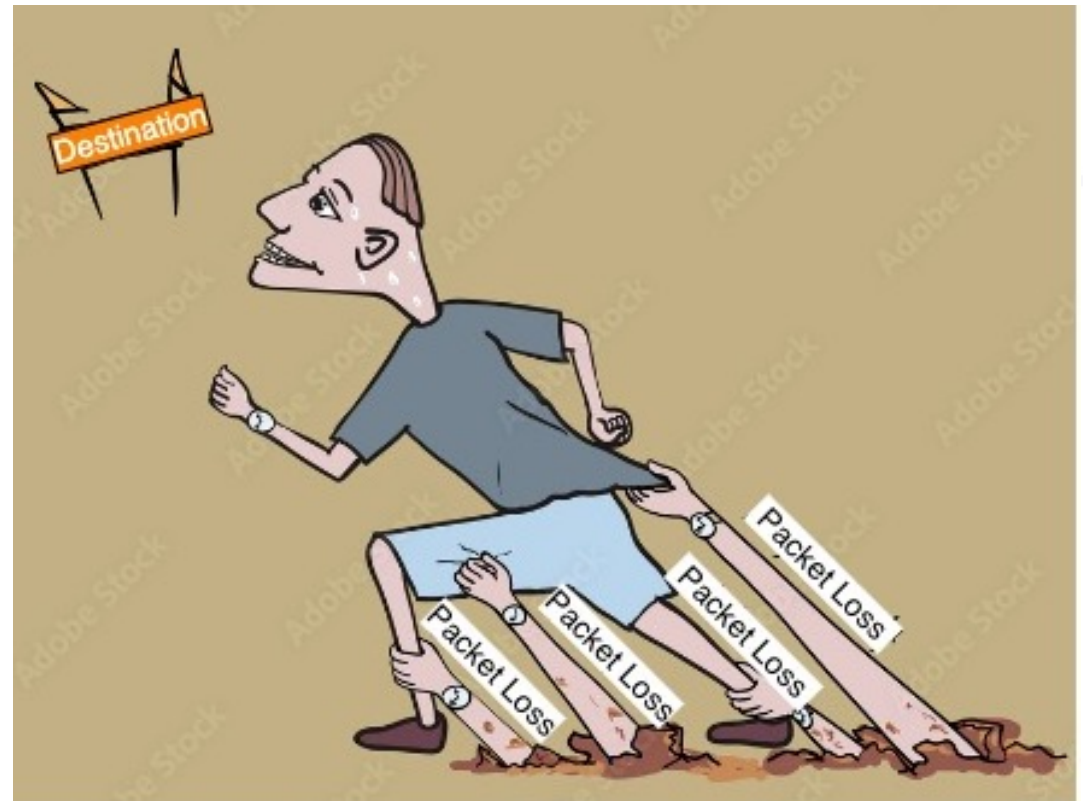# Inter-Institutes connections

These connections could be possible

# Research Institutes WAN networks

TCP performance issues:

- While most science applications that need reliable data delivery use TCP-based tools for data movement, TCP's interpretation of packet loss can cause performance issues..

- TCP interprets packet loss as network congestion, and so when loss is encountered TCP dramatically reduces its sending rate.

- The rate slowly ramps up again.

- This becomes more dramatic as the distance between communicating hosts is increased and with MTU path not in jumbo frames.
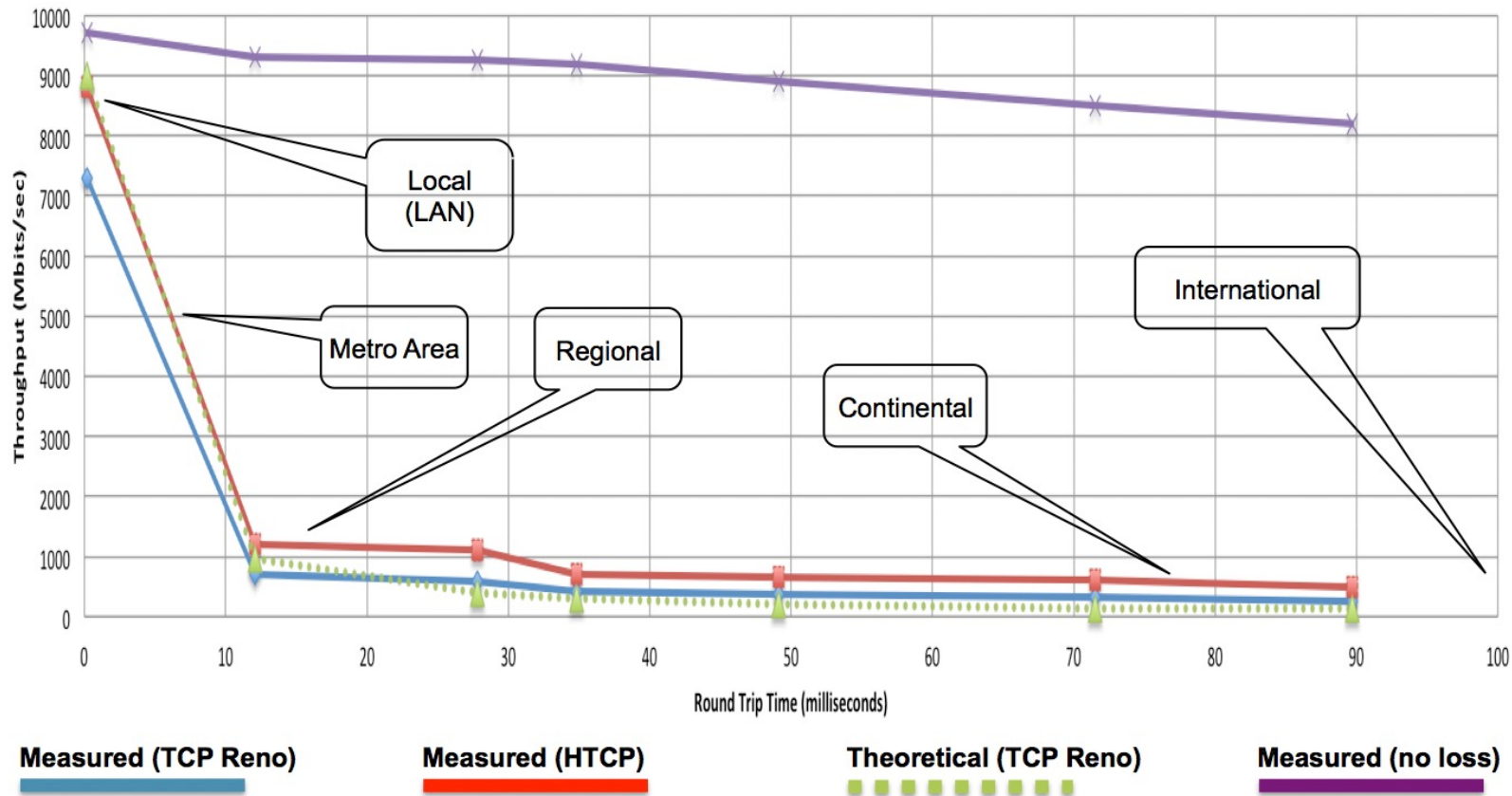
## TCP Marathon

# Research Institutes WAN networks
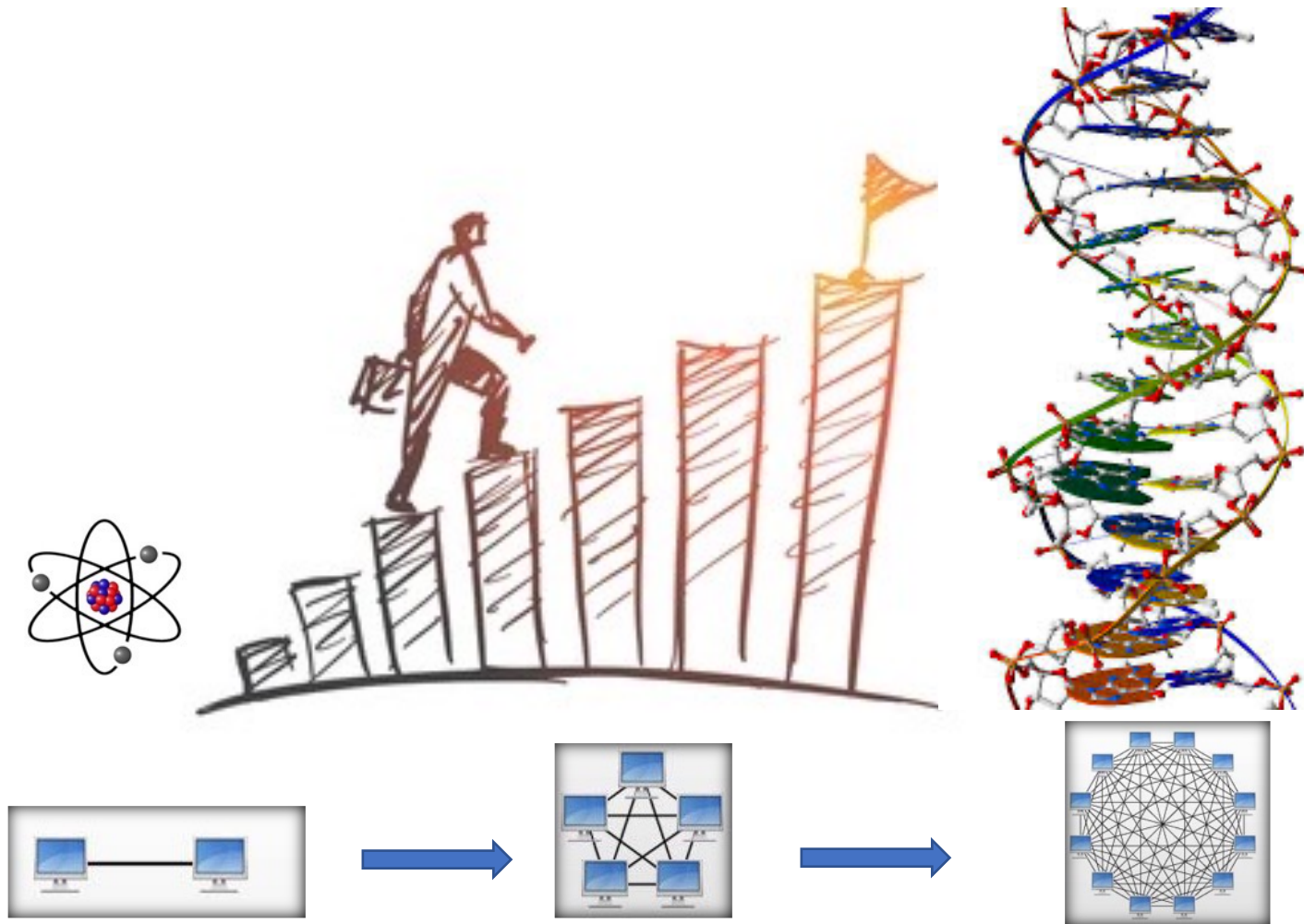
## TCP performance issues:

- A tiny amount of loss (much less than 1%) is enough to reduce TCP performance by over a factor of 50.

**Throughput vs. increasing latency on a 10Gb/s link with _0.0046%_ packet loss**



Legend labels on chart: Local (LAN), Metro Area, Regional, Continental, International

X-axis: Round Trip Time (milliseconds)
Y-axis: Throughput (Mbits/sec)

Legend: Measured (TCP Reno) · Measured (HTCP) · Theoretical (TCP Reno) · Measured (no loss)

# Research Institutes WAN networks

The challenge for a future DMZ Network
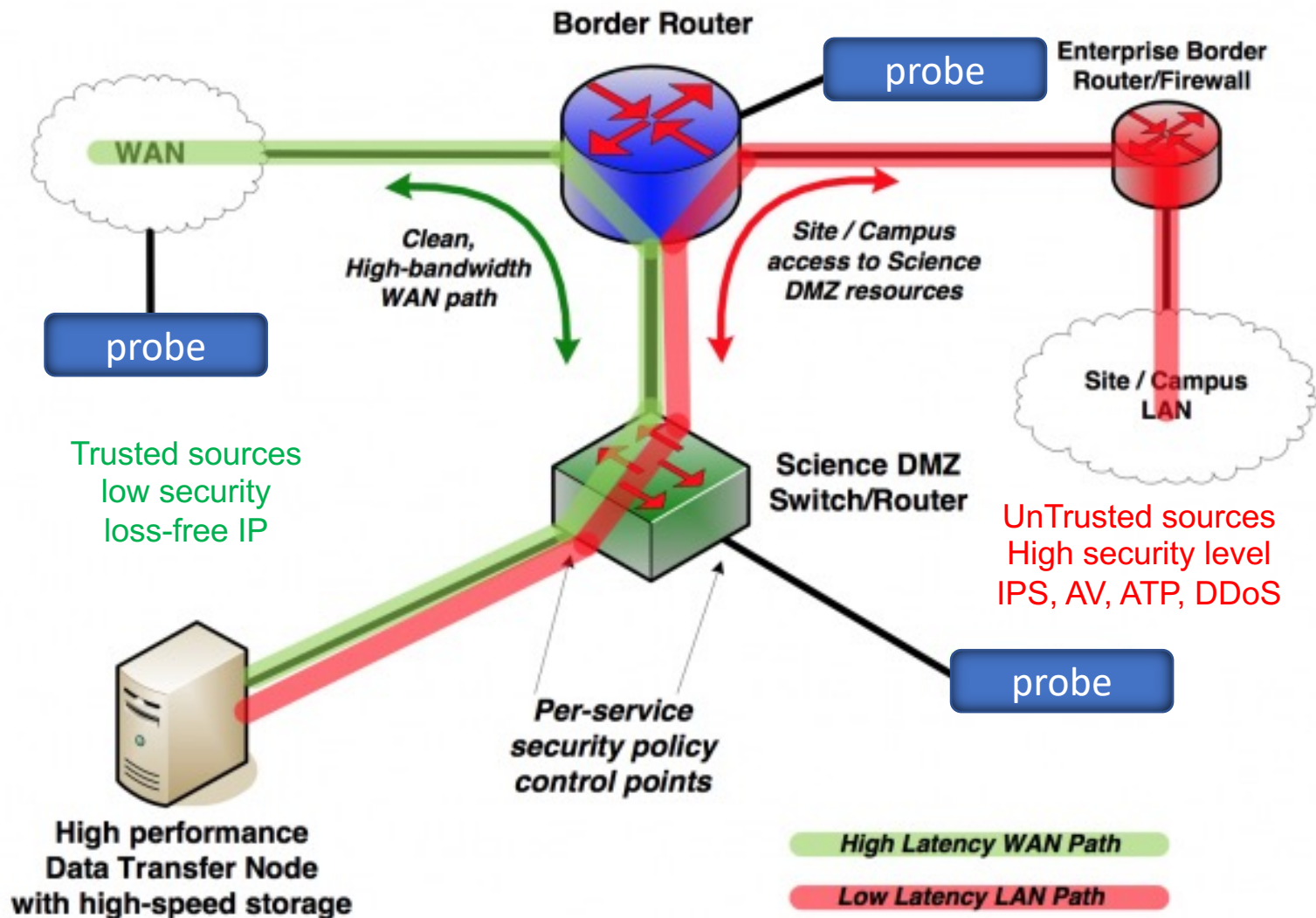
# Research Institutes WAN networks

## Science DMZ

- The **Science DMZ** Model accomplishes this by explicitly creating a portion of the network that is specifically engineered for science applications and does not include support for general-purpose use.

- By separating the high-performance science network (the Science DMZ) from the general-purpose network, each can be optimized without interfering with the other.

- The Science DMZ model allows a laboratory, campus, or scientific facility to build a special-purpose infrastructure that can provide the necessary services to allow high-performance applications to be successful.
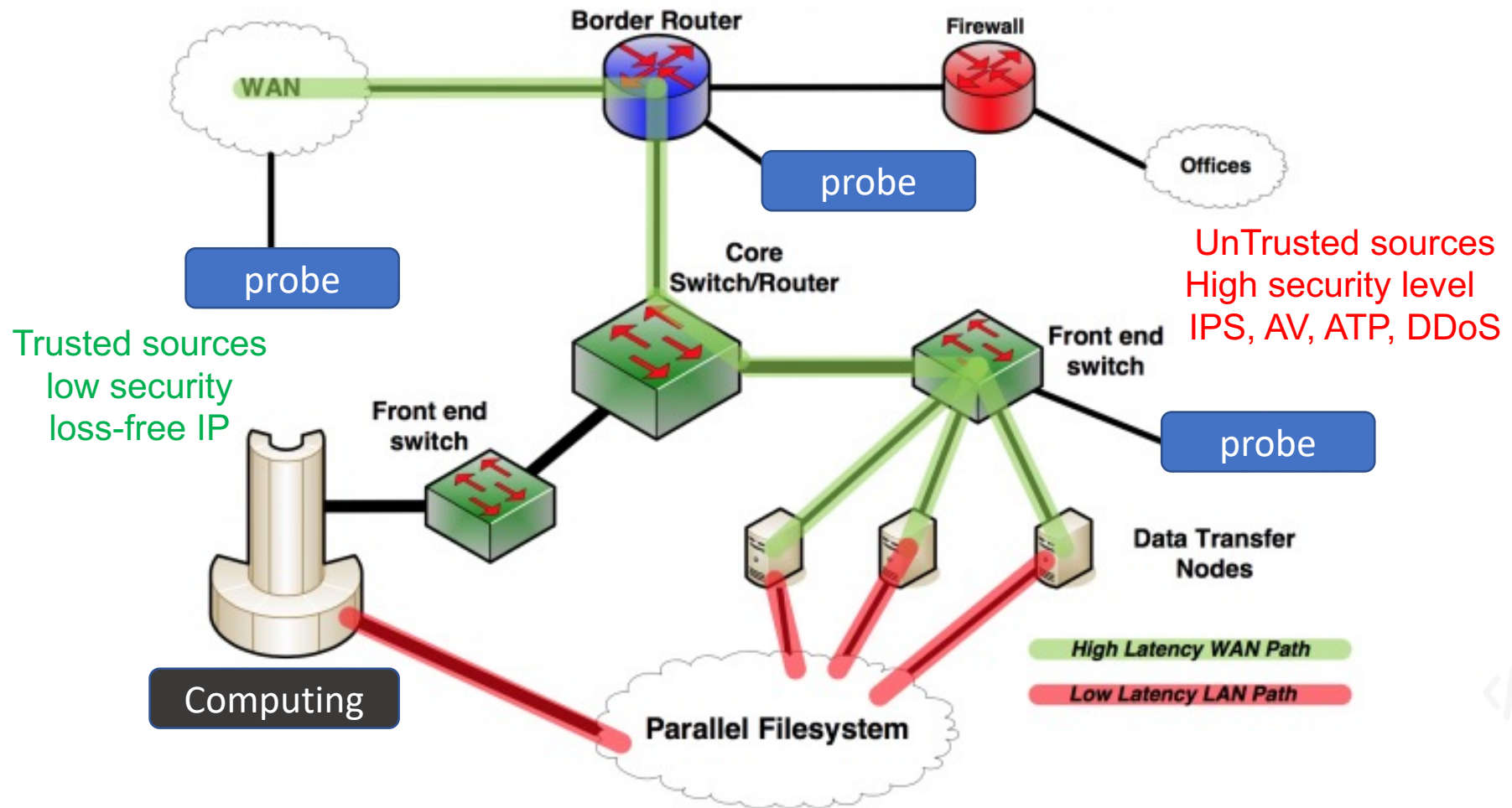
# Simple Science DMZ architecture

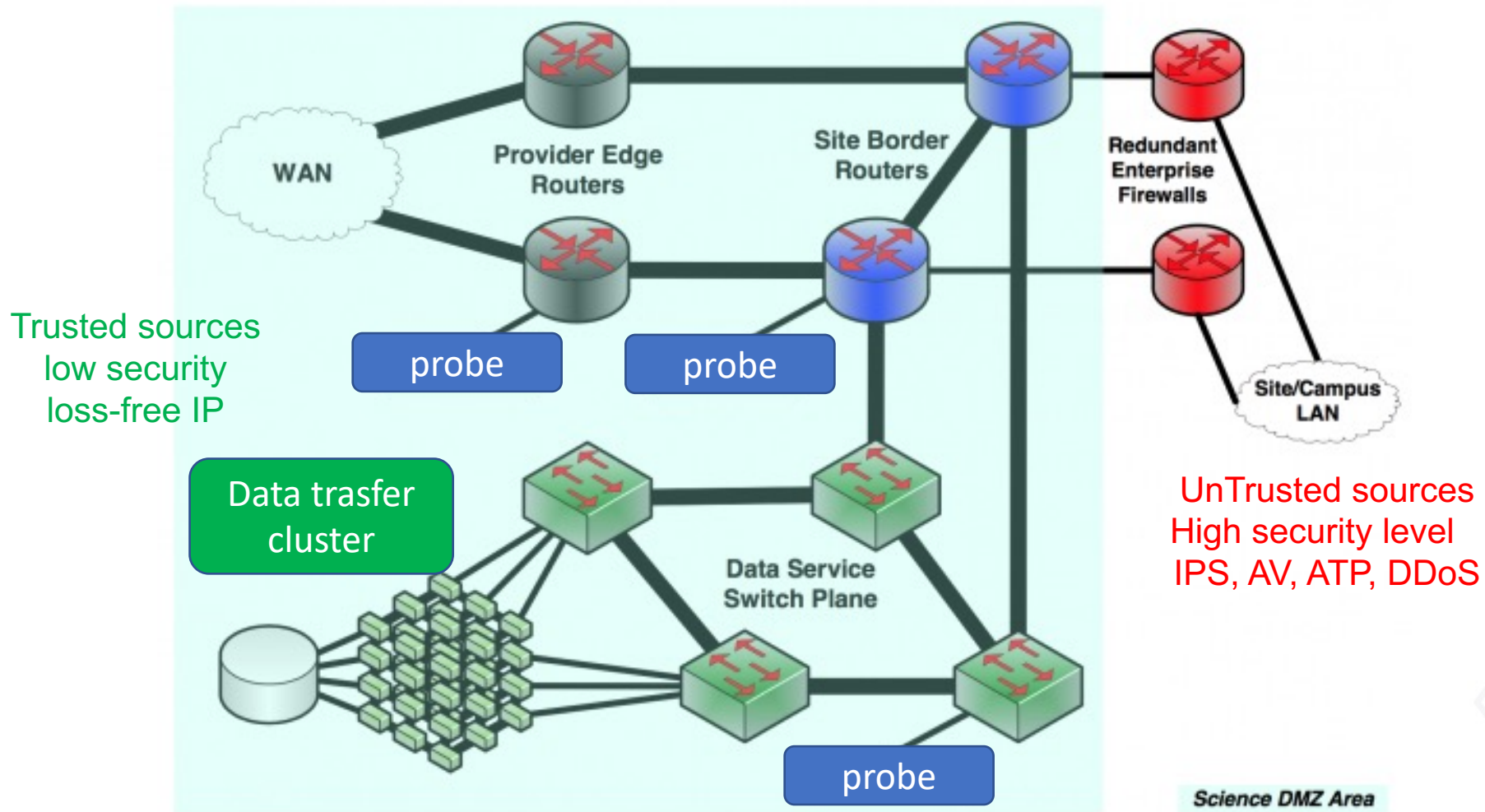High-volume bulk data transfer, remote experiment control

# Science DMZ architecture

High-volume data processing transfer, remote experiment control

# Science DMZ architecture

High-volume parallel bulk data transfer, remote experiment control



Trusted sources
low security
loss-free IP

UnTrusted sources
High security level
IPS, AV, ATP, DDoS

# Thank you