# Healthcare research: HPC and cloud.

Arnaud Ceol, arnaud.ceol@ieo.it,

WORKSHOP GARR 2019, 8-10 ottobre 2019, Rome

IEO
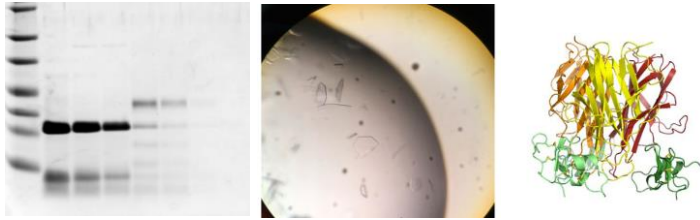Istituto Europeo
di Oncologia

IEO25

# IEO Hospital

# Department of Experimental Oncology



Clinical and basic research, Integration within IEO Clinical Programs, PhD program
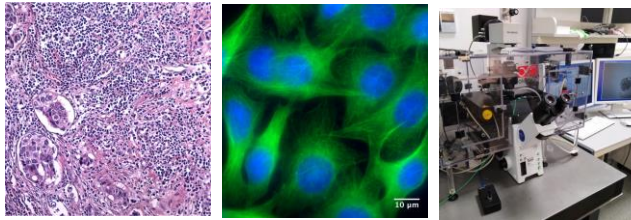› 20 bioinformaticians

IEO
Istituto Europeo
di Oncologia    IEO25

# Technological units: raw & processed data

## Crystallography



- Few users
- 1 TB/year
- Analyses: Memory and GPUs

## imaging
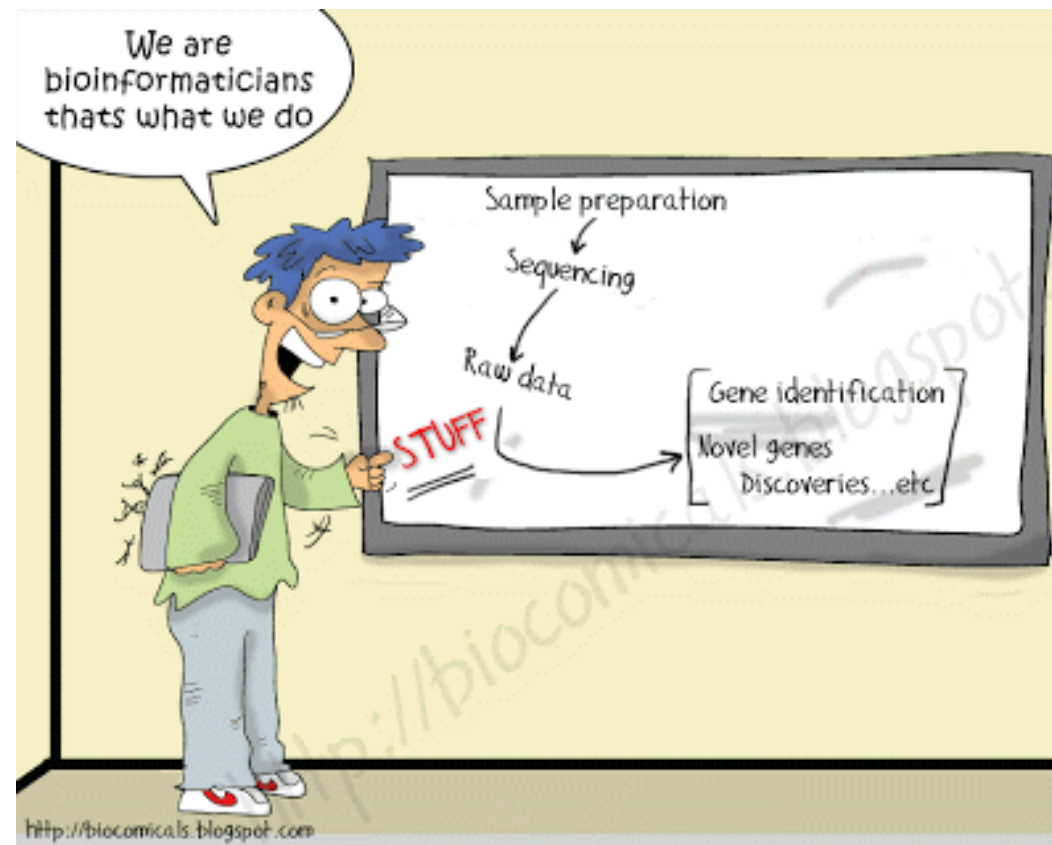


- ~50 TB/year
- ~10 users
- Big images visualization

## genomics



- > 20 TB/year
- cpus and memory for data preparation and analyses
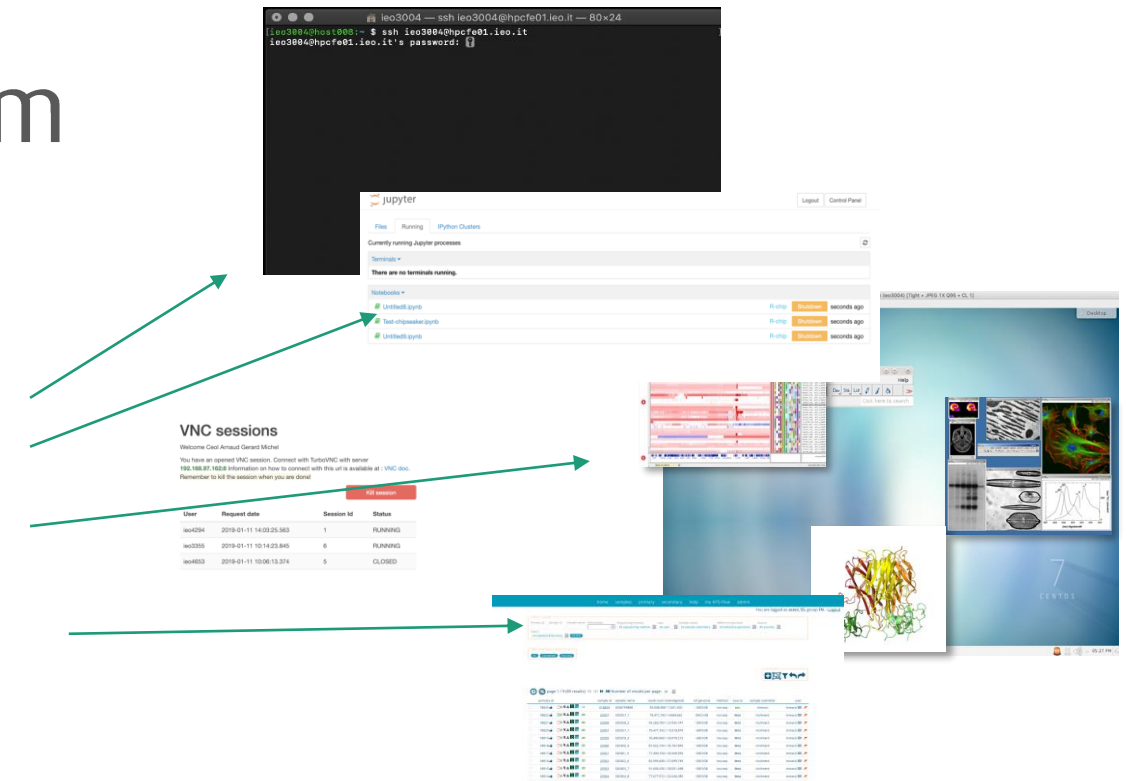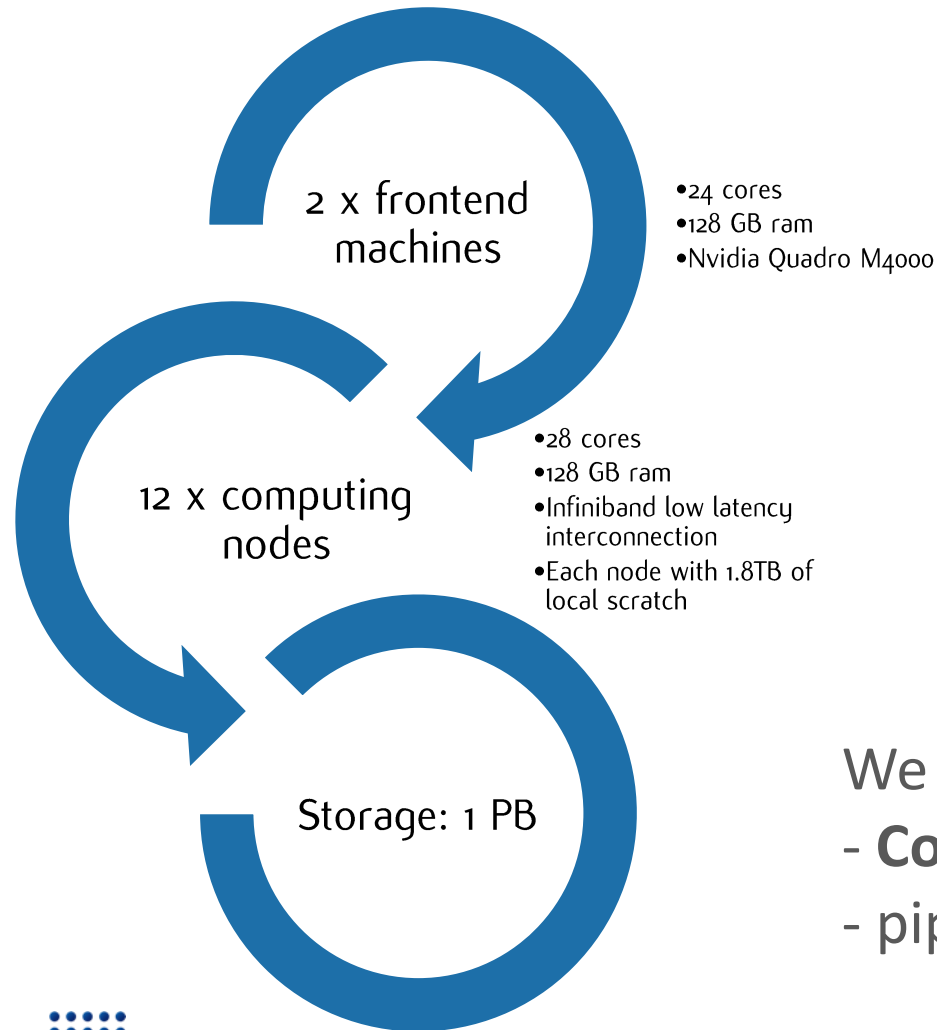
## Proteomics, Radiomics, …



- > 100 TB/year ?
- CPUs, GPUs

IEO
Istituto Europeo
di Oncologia    IEO25

Dry-lab Researchers (bioinformaticians): analyze the data

# The bioinfo (HPC) platform



**2 x frontend machines**

- 24 cores
- 128 GB ram
- Nvidia Quadro M4000

**12 x computing nodes**

- 28 cores
- 128 GB ram
- Infiniband low latency interconnection
- Each node with 1.8TB of local scratch

Storage: 1 PB

We encourage the usage of:
- **Containers**: 100% reproducible results
- pipeline managers (Nextflow, Snakemake)

# nextflow

```nextflow
1   #!/usr/bin/env nextflow
2
3   /*…
12  params.bam
13  params.parameters
14  params.regionBed
15  params.refGenome
16  params.outputDir
17
18
19  process tvc {
20
21      cpus 4
22      time '4h'
23
24      script:
25      """
26      mkdir -p ${params.outputDir}
27      variant_caller_pipeline.py --num-threads 4 \
28          --input-bam ${params.bam} \
29          --parameters-file ${params.parameters} \
30          --reference-fasta ${params.refGenome} \
31          --region-bed ${params.regionBed} \
32          --output-dir ${params.outputDir}
33
34      """
35  }
36
```

```
1   singularity.enabled = true
2   process.container = '/hpcnfs/techunits/bioinformatics/singularity/smith.simg'
3   singularity.runOptions = ' --bind /hpcnfs/ '
4   process.executor = 'pbs'
```
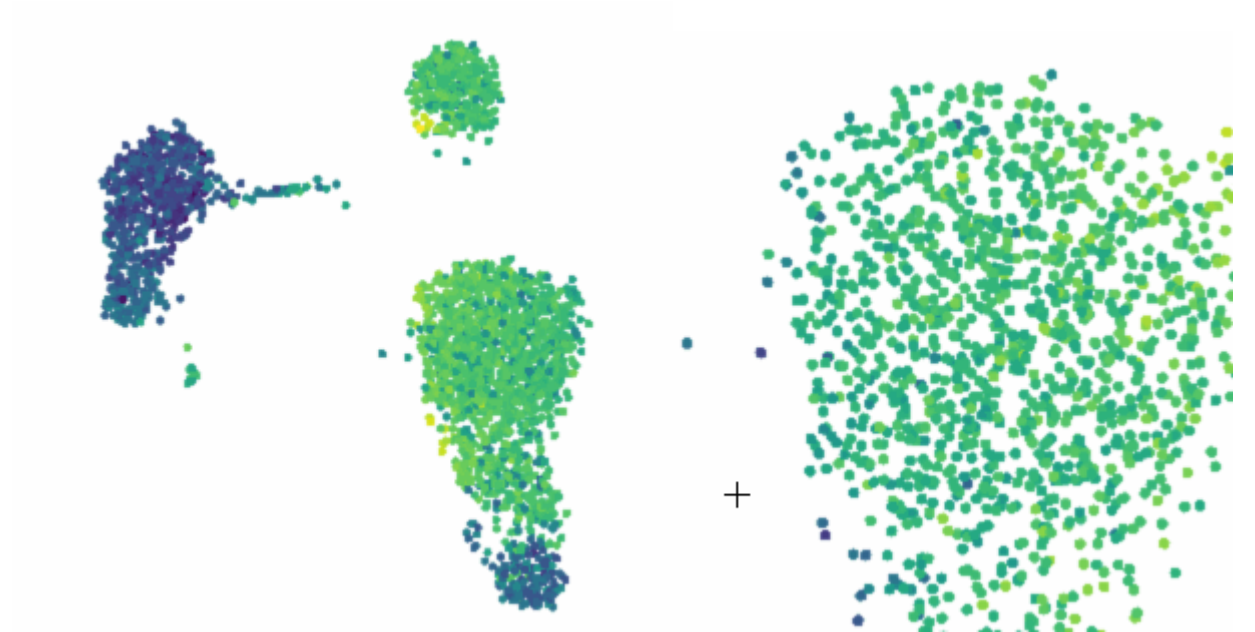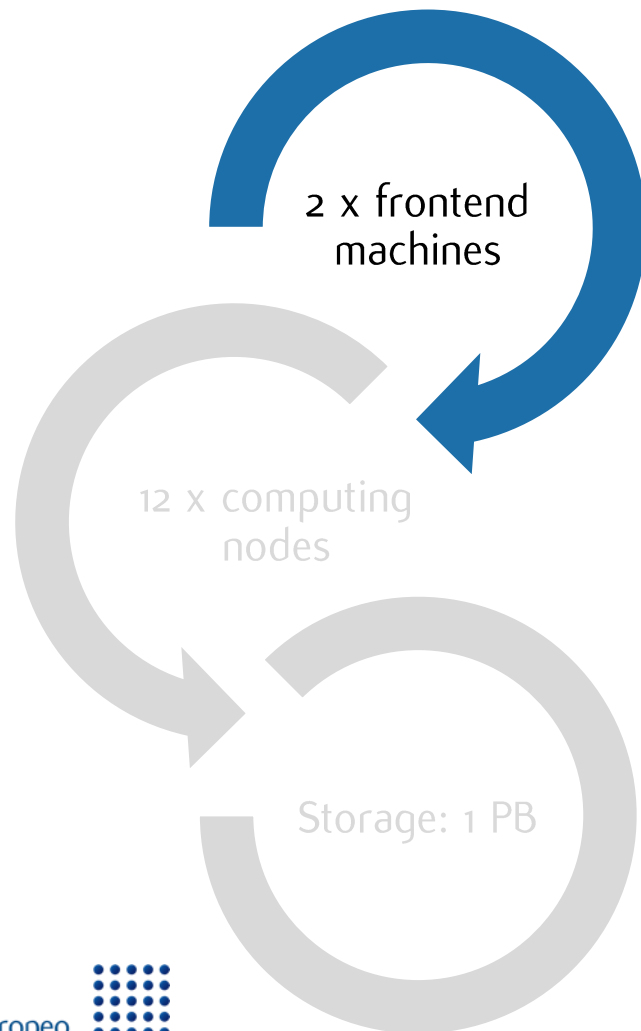
**⊟ Executors**

- Local
- SGE
- LSF
- SLURM
- PBS/Torque
- PBS Pro
- Moab
- NQSII
- HTCondor
- Ignite
- → Kubernetes
- → AWS Batch
- → Google Pipelines
- GA4GH TES

IEO Istituto Europeo di Oncologia   IEO25

# HPC: limitations, and how the cloud can help

# Limitation(s) 1: no admin, no Docker, no web-serving, no outside access



2 x frontend machines

12 x computing nodes

Storage: 1 PB

Example: testing **docker-based web application** for the visualization and analyses of single cell sequencing

IEO
Istituto Europeo
di Oncologia    IEO25

cell×gene

# Limitation 1: no admin and web serving



2 x frontend machines

12 x computing nodes

Storage: 1 PB

Cloud solution

Virtual Machine 1: Managed by the user: Docker + floating IP

Virtual Machine 2: Provide to all users of the institute a "HPC-like" image.

openstack™

openstack™

JAAS

The cloud allows to create testing environment & publish web servers/services

IEO
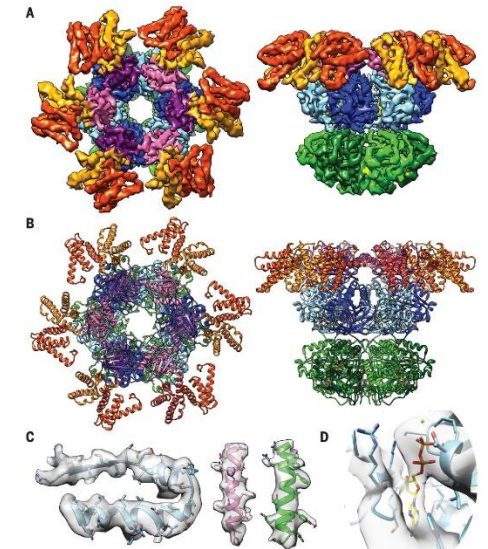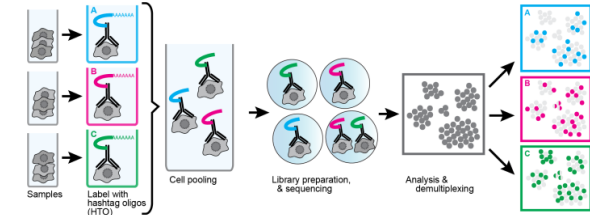Istituto Europeo
di Oncologia
IEO25

# Limitation 2: cores, GPU, RAM

**Example 1:** CRISPR screening ortoghonal interference of 1M cells with multiplexing-enabling. **5TB** raw data to be analized **100000 cpu/hour** and for the integration and analysis we need a **512GB ram** machine

**Example 2: DeepVariant** is an analysis pipeline that uses a deep neural network to call genetic variants from next-generation DNA sequencing data: needs VM with **64 CPUs 240 GB RAM**
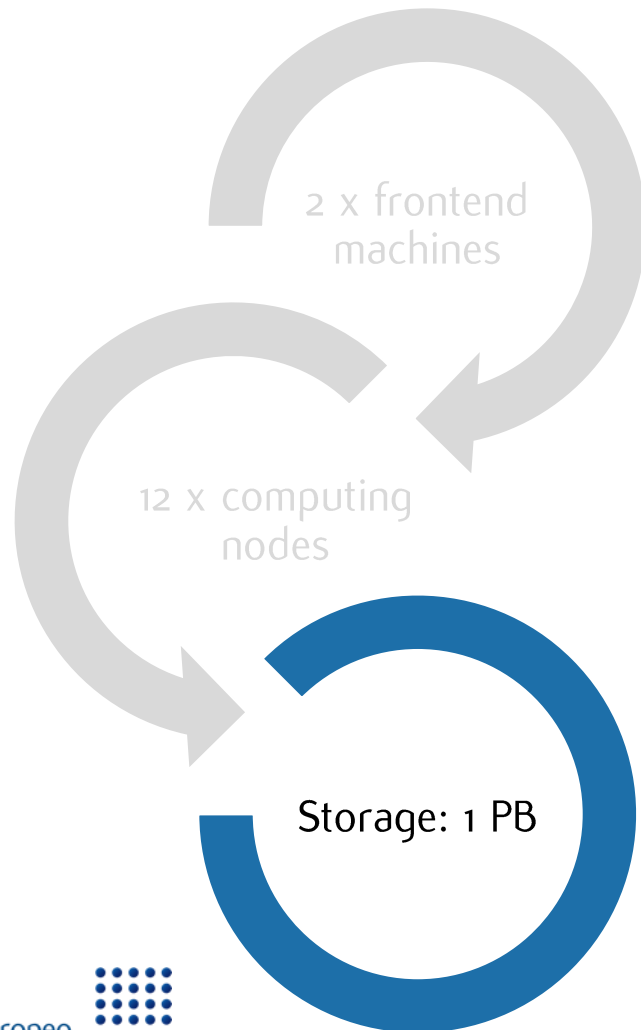
**Example 3: Cryo-EM** analyze large, complex and flexible structures. 1 project = **3 weeks to 6 months** with **2 GPUs**

2 x frontend machines

12 x computing nodes

Storage: 1 PB

IEO
Istituto Europeo di Oncologia    IEO25

Cloud solution

openstack

# Limitation 3: storage of archives

2 x frontend machines

12 x computing nodes

Storage: 1 PB

Cloud solution

Storage: Isilon (1PB)

Growth rate: ~ 100 TB/year (only raw data), but will increase.
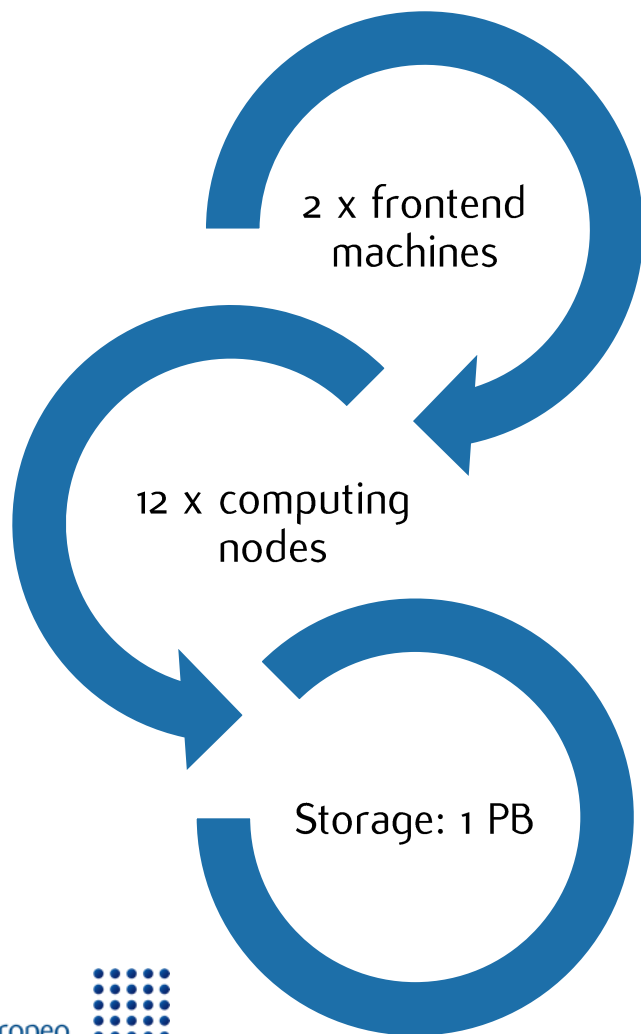
Cold data: > 300 TB
- 90 % will not be used anymore
- cannot be deleted (publication, patents, etc.)

Cloud storage:
-> admin-based management
-> user-based management

# Limitation 1: exchange and collaboration

2 x frontend machines

12 x computing nodes

Storage: 1 PB

Cloud solution

HPC and storage accessible only from the IEO network (intranet).

Consortium GARR

FILESENDER

RCLONE

ceph

IEO
Istituto Europeo di Oncologia    IEO25

# Summary and further considerations

Cloud for: **freedom**, **BIG stuff**, **long term**


Virtual data centers
for collaborative projects

Clinical/biomedical data    Solutions & guidance are welcome

**1994-2019**
25 anni di ricerca e innovazione
per la lotta al cancro,
25 anni di Istituto Europeo di Oncologia



Arnaud Ceol, arnaud.ceol@ieo.it