# What is OMICS?

**CINECA**

CINECA has built OMICS, an **user-friendly web interface** designed for seamless **uploading** and **processing of genomic** data through NVIDIA Clara Parabricks, a software suite of GPU-accelerated genome analysis tools for genomic **read alignment**, **variant calling**, and GVCF filtering and post-processing.

## Main Features

✓ Data security and privacy compliance with EU GDPR regulations are maintained through a secure encrypted cloud Virtual Machine (VM)
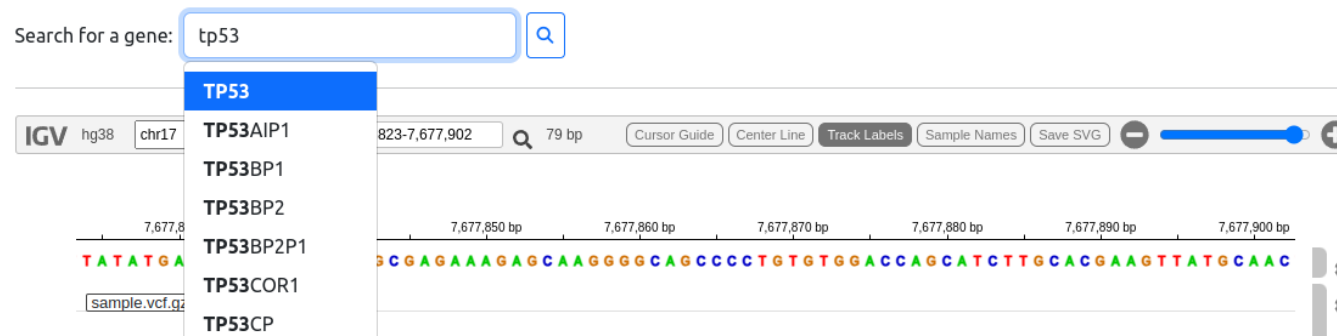
✓ VM is equipped with GPUs to make the analyses increasingly scalable

Finanziato
dall'Unione europea
NextGenerationEU

Ministero
dell'Università
e della Ricerca

Italiadomani
PIANO NAZIONALE
DI RIPRESA E RESILIENZA

ICSC
Centro Nazionale di Ricerca in HPC,
Big Data and Quantum Computing

# What is OMICS?

CINECA

- OMICS main dashboard allows users to **effortlessly upload their raw data** and **customize advanced options** for three containerized pipelines
- Capability to swiftly upload and download large datasets (e.g. raw sequence data)



- OMICS aims to streamline WGS data analysis **while prioritizing user experience**, data **privacy**, and **accelerated** processing through **GPU** technology, aiming at supporting activities of Spoke 0 and Spoke 8 of ICSC

# Il gruppo di lavoro CINECA

- Giuseppe Melfi
- Silvia Gioiosa
- Alessandro Grottesi
- Juan Mata Naranjo
- Xhulio Dhori

# Under the hood

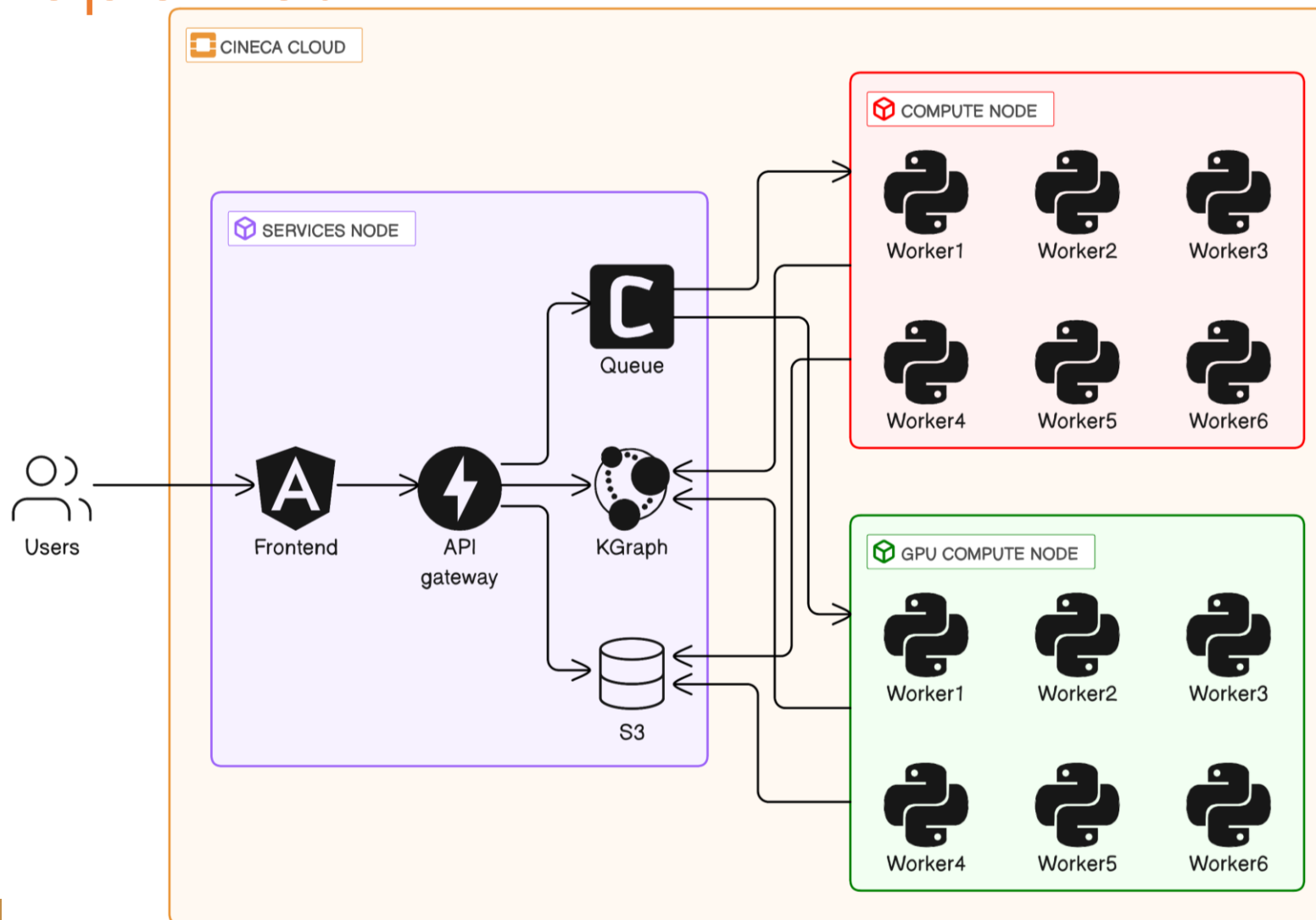Let's look at what's going on behind the scenes

# High level architecture

- Simplified/unified access
- Unified data layers
  - **APIs**
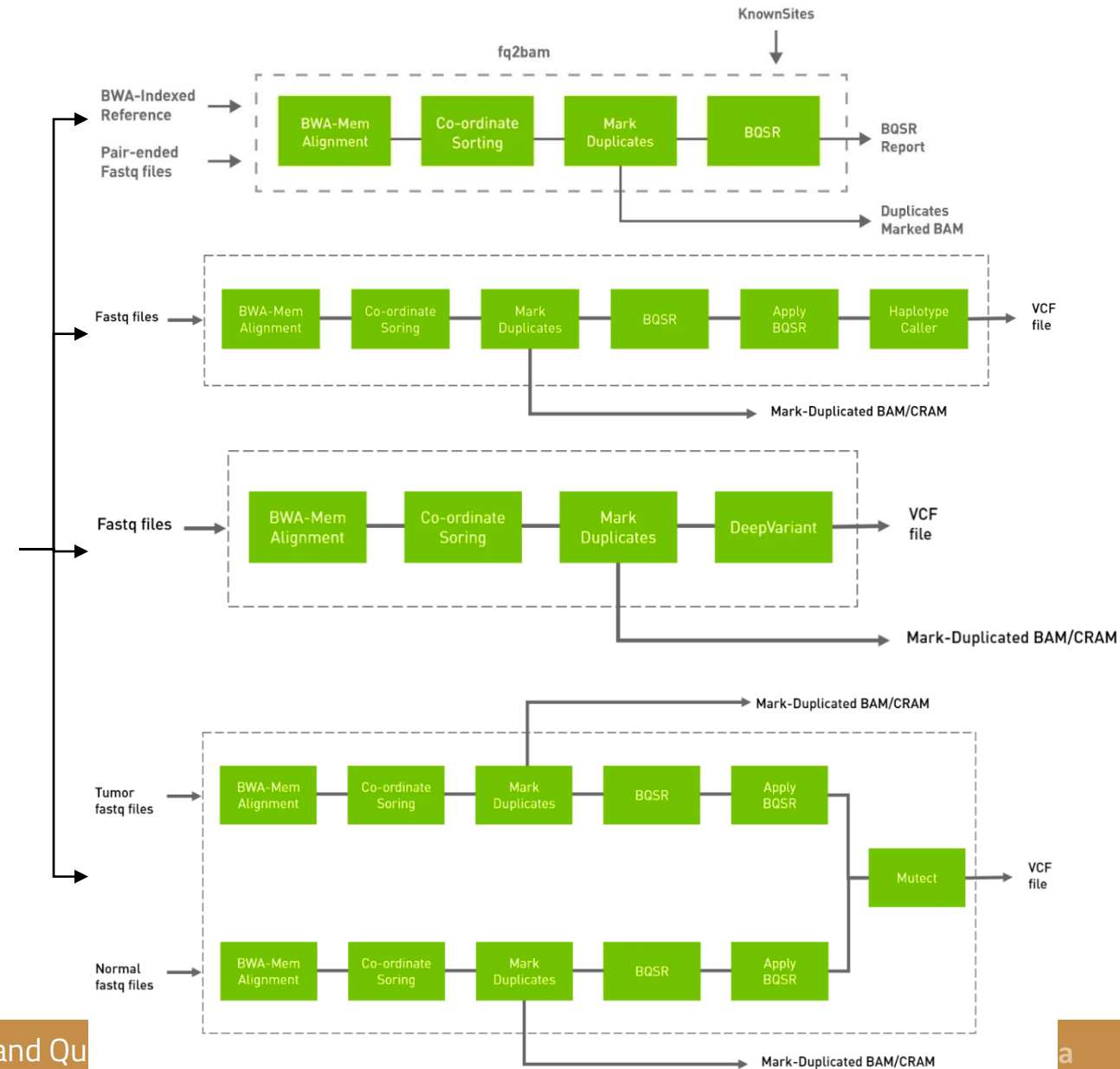  - **Knowledge Graph**
  - **Storage**
- Vertical custom tools



Frontend UI

REST API

Knowledge Graph

Service #1 Tools | Service #2 Tools | Service #3 Tools | Service #4 Tools

Storage

# Design properties

- **Open Source** techs
- Microservices oriented
- Cloud ready
- Easy portability
- Easy scalability
- **Private Cloud** resources

# Current Features

- **Multichunk Upload** (User can upload large fastq files directly from the interface)

- **Quality Check** (Fastqc and Multiqc report)

- **Ready to use pipelines** on GPUs:
    - Fq2bam
    - Germline
    - DeepGermline
    - Somatic

- **Annotation** tools (VEP)

- **Visualization** (IGV Browser)

- **User Guide**

Finanziato
dall'Unione europea
NextGenerationEU

Ministero
dell'Università
e della Ricerca

Italiadomani
PIANO NAZIONALE
DI RIPRESA E RESILIENZA

ICSC
Centro Nazionale di Ricerca in HPC,
Big Data and Quantum Computing

# BENCHMARKS

# NEXT STEPS

- Beta tester **feedbacks**
- Partners **feedback** on **data management** and permissions, **sharing** results within own organization
- **Batch analysis** on multiple samples
- **Re-run** analysis with the same parameters
- **Upload** with **command line**
- **Preprocess** your data (trimming…)
- Cloud **upgrade** (more GPUs available)
- **UX/UI** improvements

# *Demo live*

# DASHBOARD

## AND LOGIN

# UPLOAD FILES

# START GPU ACCELERATED ANALYSIS

# START GPU ACCELERATED ANALYSIS

# CONFIGURE ADVANCED OPTIONS

# CONFIGURE ADVANCED OPTIONS

# VEP ANNOTATION

**Variant classes**

Legend:
- sequence_alteration
- insertion
- deletion
- SNV
- Other

| Variant class | Count |
|---|---:|
| indel | 23 |
| sequence_alteration | 345 |
| insertion | 7,777 |
| deletion | 9,236 |
| SNV | 108,562 |

**Consequences (most severe)**



Legend:
- missense_variant
- splice_region_variant
- splice_polypyrimidine_trac…
- synonymous_variant
- 5_prime_UTR_variant
- 3_prime_UTR_variant
- non_coding_transcript_ex…
- intron_variant
- upstream_gene_variant
- downstream_gene_variant
- intergenic_variant
- Other

CINECA

# IGV VISUALIZATION FOR VCF AND BAM FILES

# FASTQC AND MULTIQC REPORT

# FASTQC AND MULTIQC REPORT