

Big data in life sciences

from theory to applications

Stefano Ceri

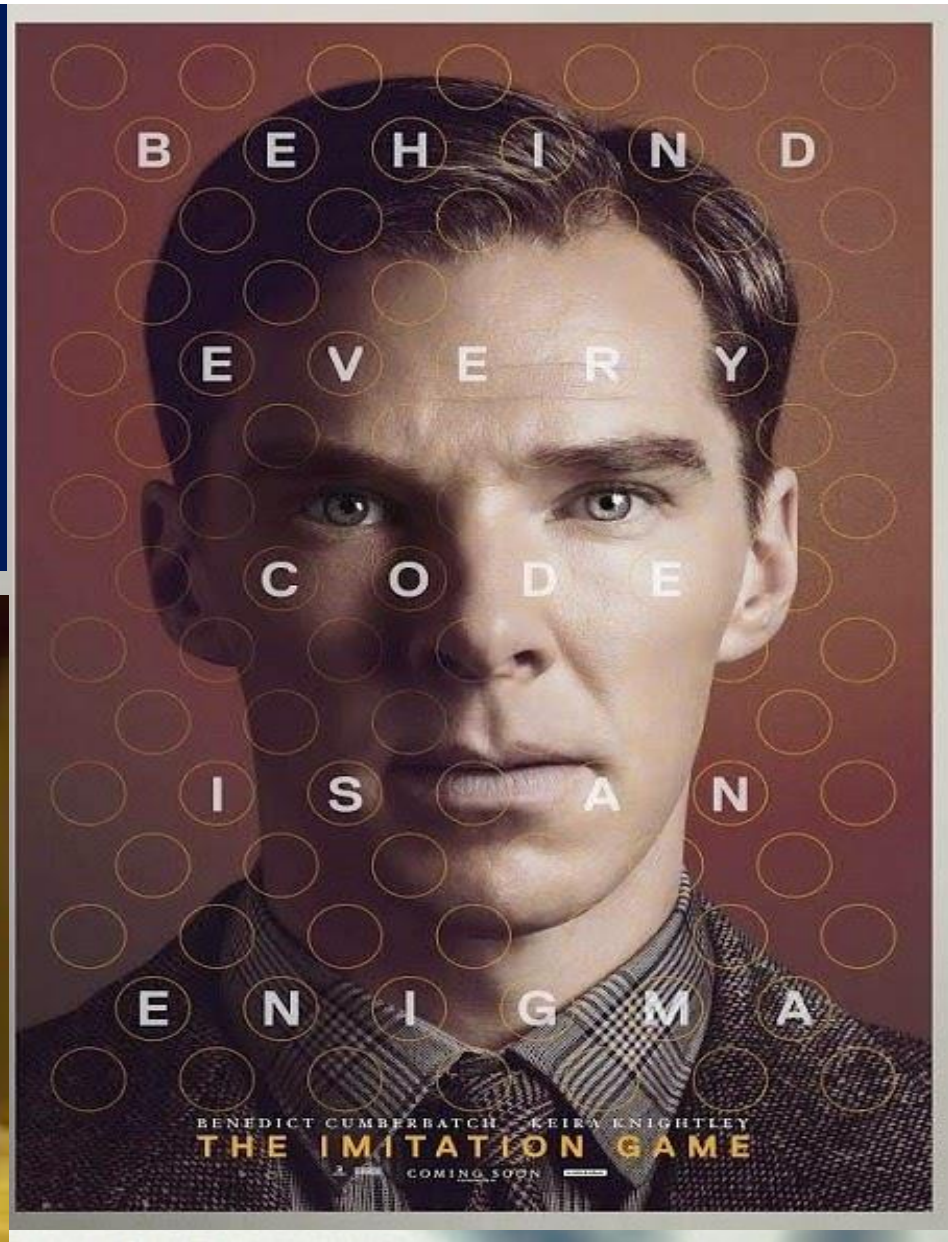
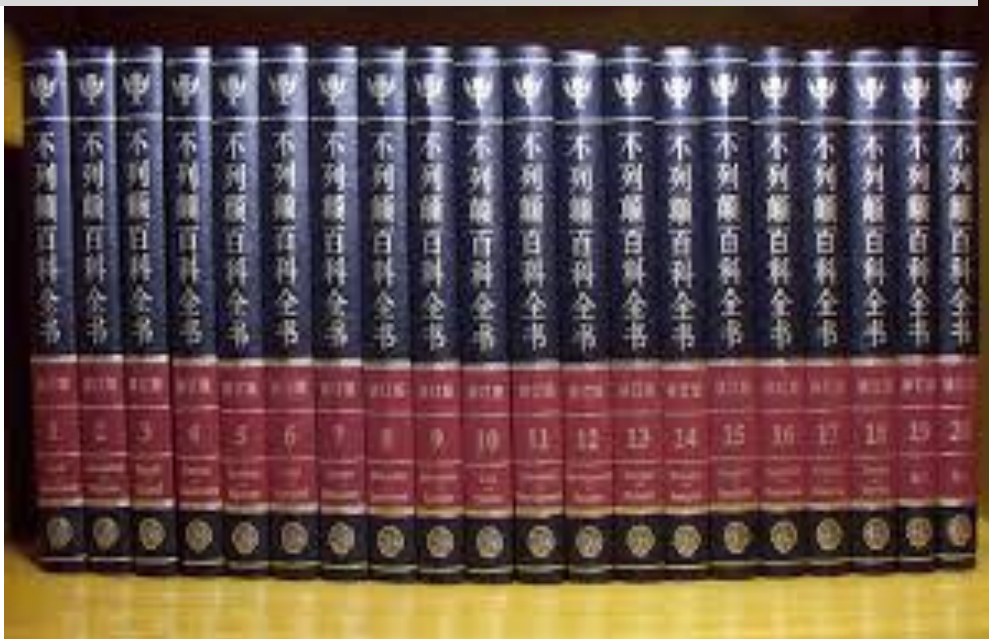
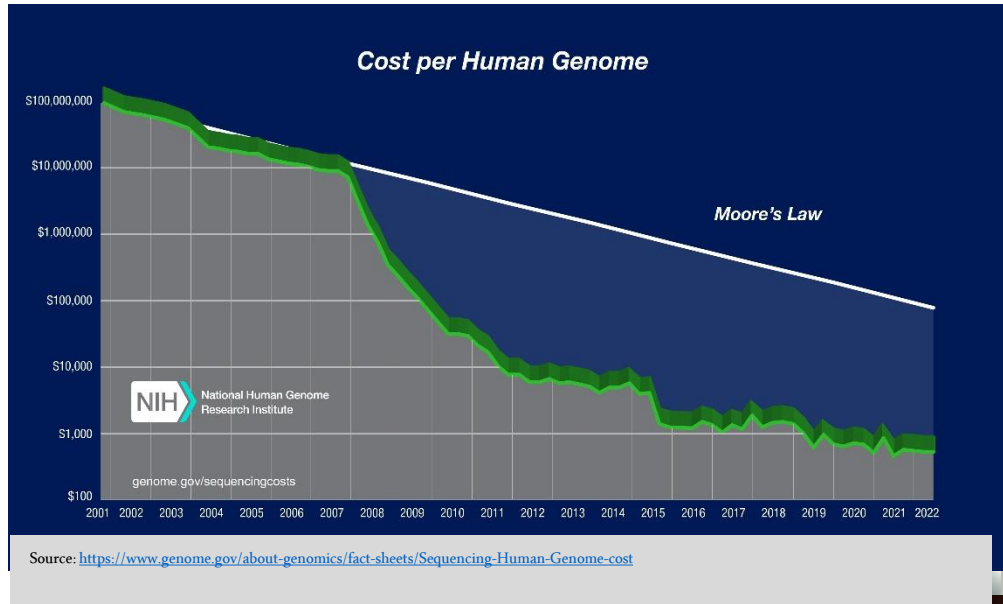
DEIB | Dipartimento di Elettronica,
Informazione e Bioingegneria

<https://ceri.faculty.polimi.it/>



POLITECNICO
MILANO 1863

The big data challenge --- human genomics



Data-Driven Genomic Computing (GeCo)

ERC Advanced Grant, 2016-2021

Focus: Data design, integration, extraction and analysis for genomic data

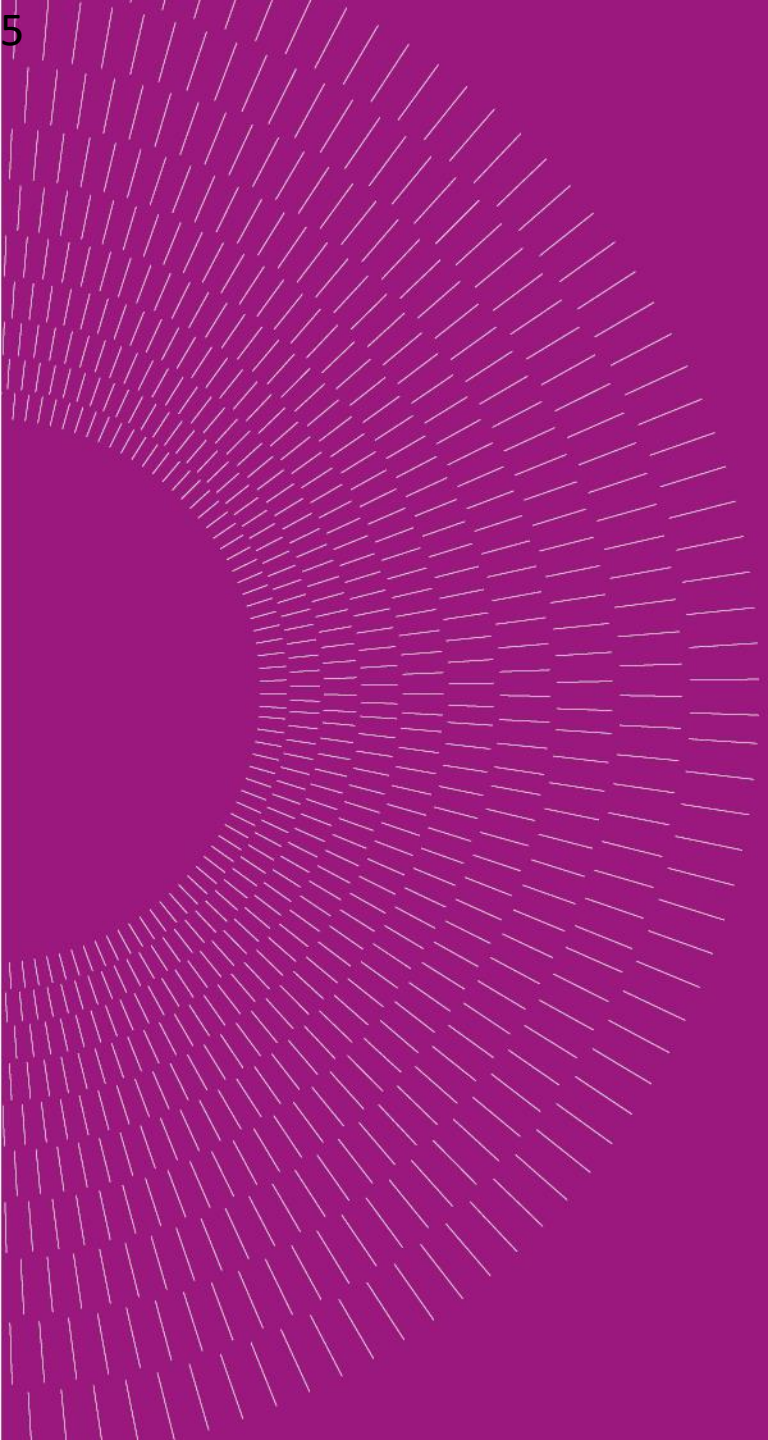
Approach: Radical change in data management abstractions for genomics (broader and simpler)

Objectives: Open-source systems for genomic data management

Results: Big data management systems + integrated repositories, both in human and viral genomics, demonstrated through biological and clinical research.

Human vs viral genomic research

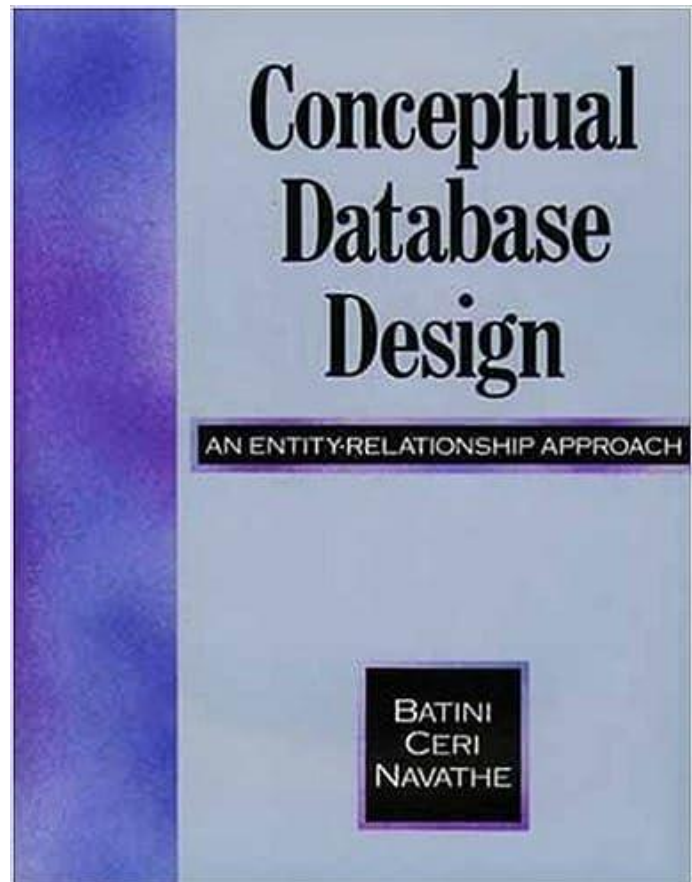
- **HUMAN GENOME: 3 Billion DNA base pairs, about 20K genes**
 - Data: mutations, expression, methylation, 3D contacts, ...
 - Biga data approach: design, integrate, search, analyze
 - Most relevant clinical applications: rare diseases and cancer – but genetics produces a huge number of human traits
 - Research highlights: basic science for genomics, disease prevention & treatment, personalization, drug repurposing.
- **SARS-CoV-2 VIRAL GENOME: 30K RNA nucleotides, 12 genes**
 - Data: viral sequences and their nucleotide mutations
 - Reuse our big-data approach, applied to a "smaller" problem
 - Research driven by pandemic emergency: Understand genome evolution (mutations, variants), their mechanisms (co-evolution, recombination – now reassortment for influenza) and threats to humans, e.g. effects on COVID-19 disease spreading and severity.

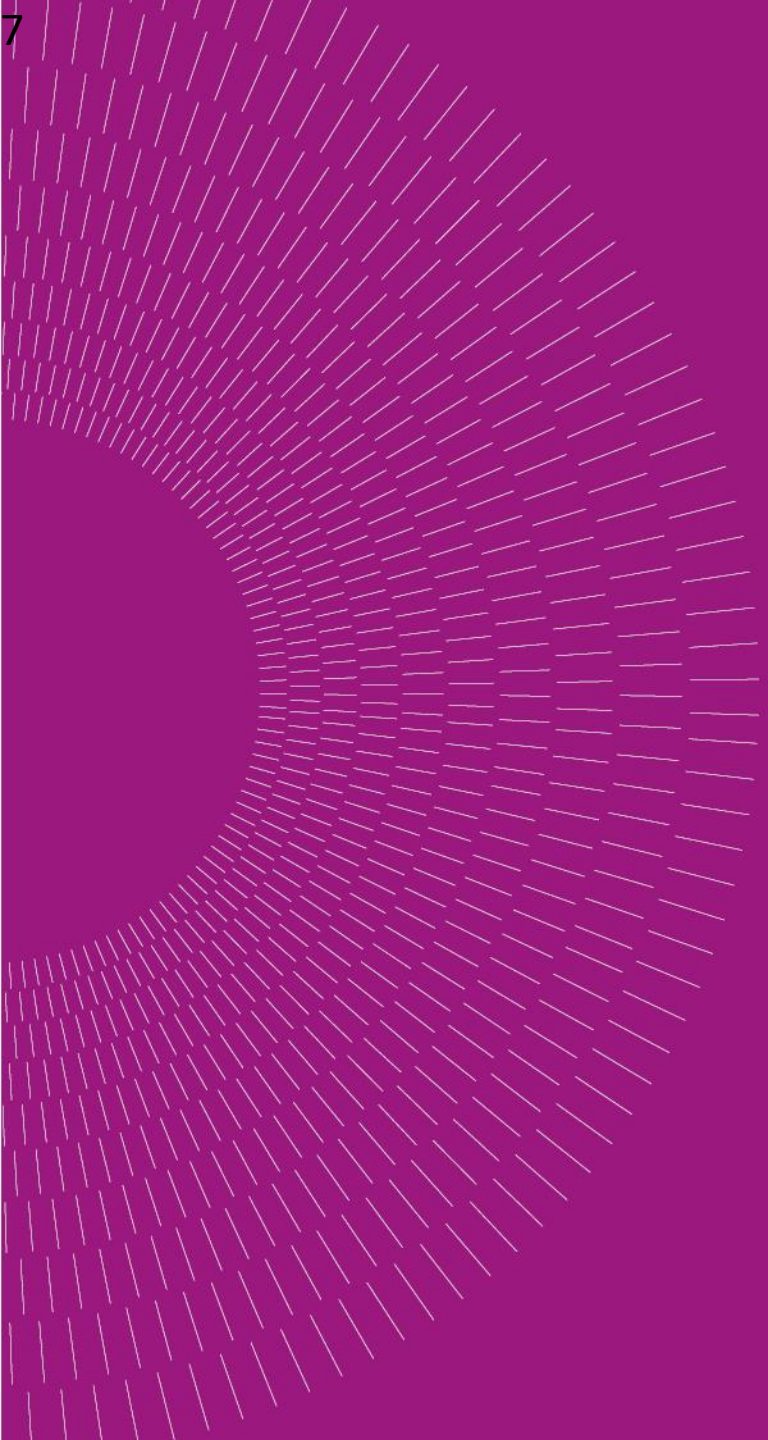


DATA DESIGN

Design objectives & challenges

- TOP-DOWN: Clean design, with nice abstractions
- BOTTOM-UP: Matching with existing data, providing content
- Model: Entity-Relationship (P. Chen)





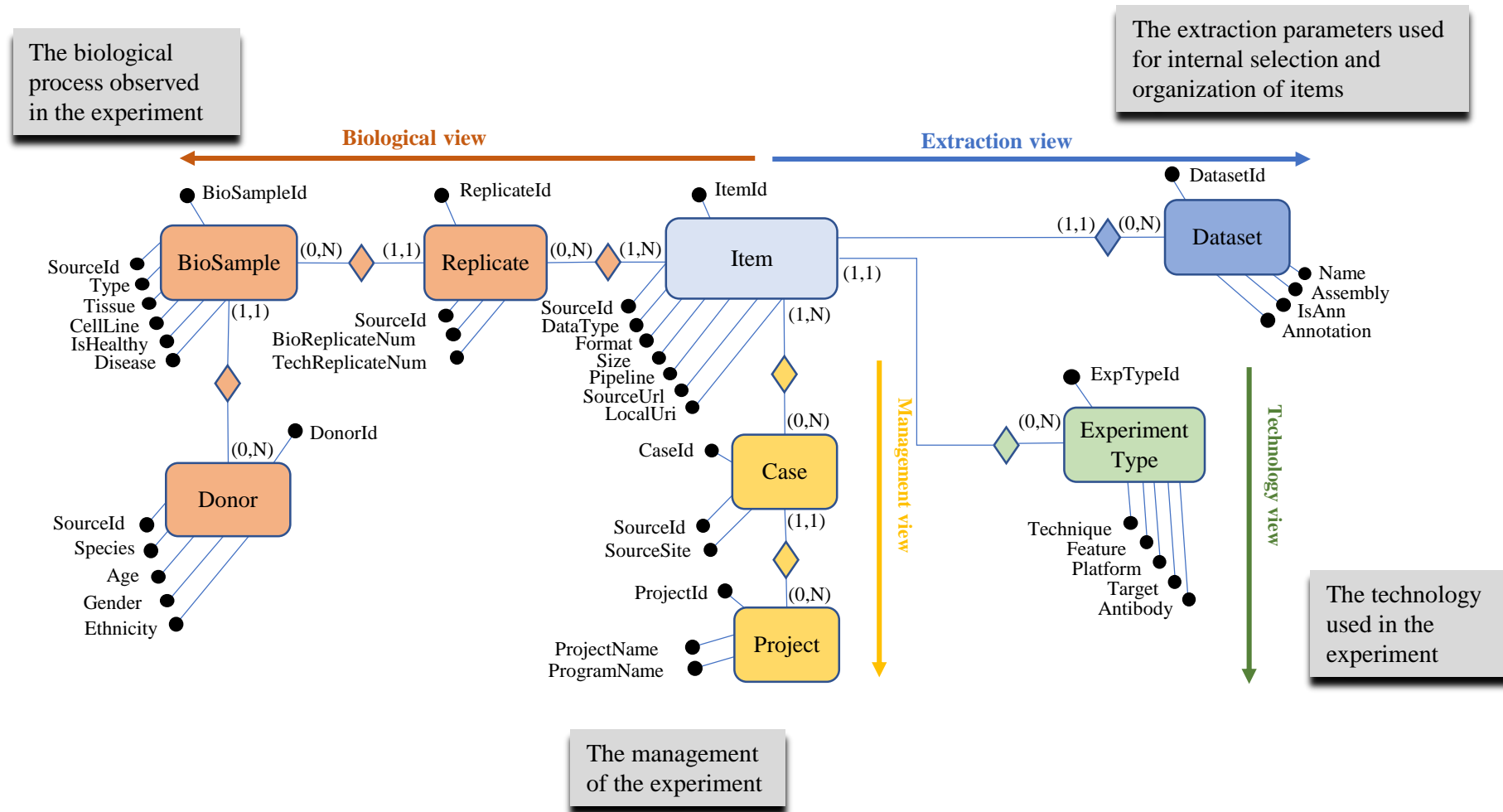
HUMAN GENOMICS:

DATA DESIGN

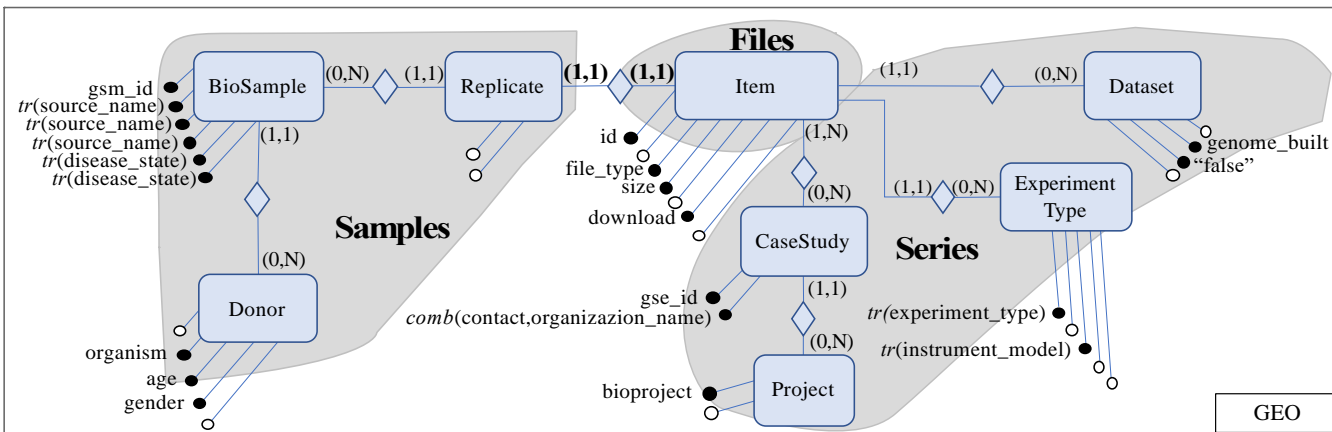
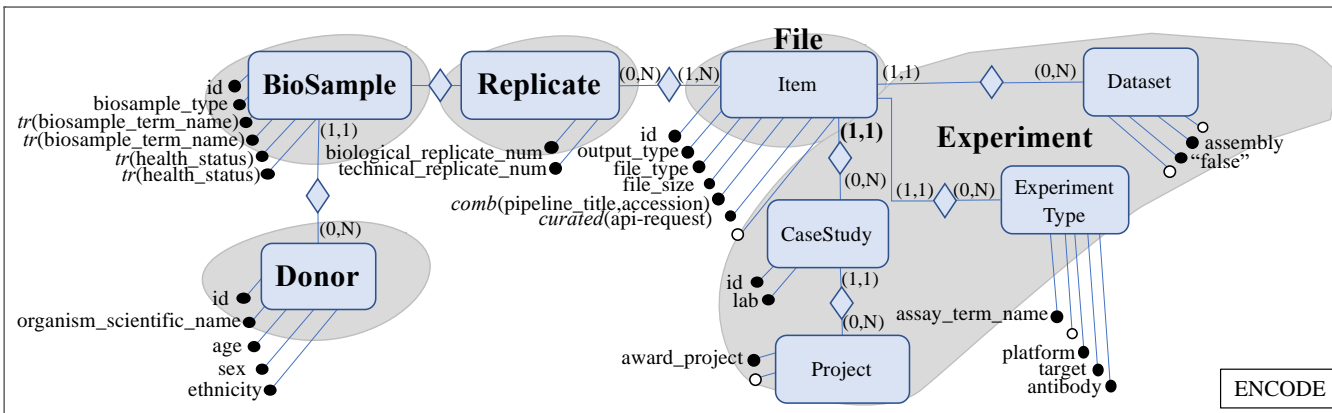
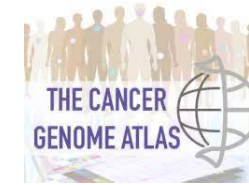
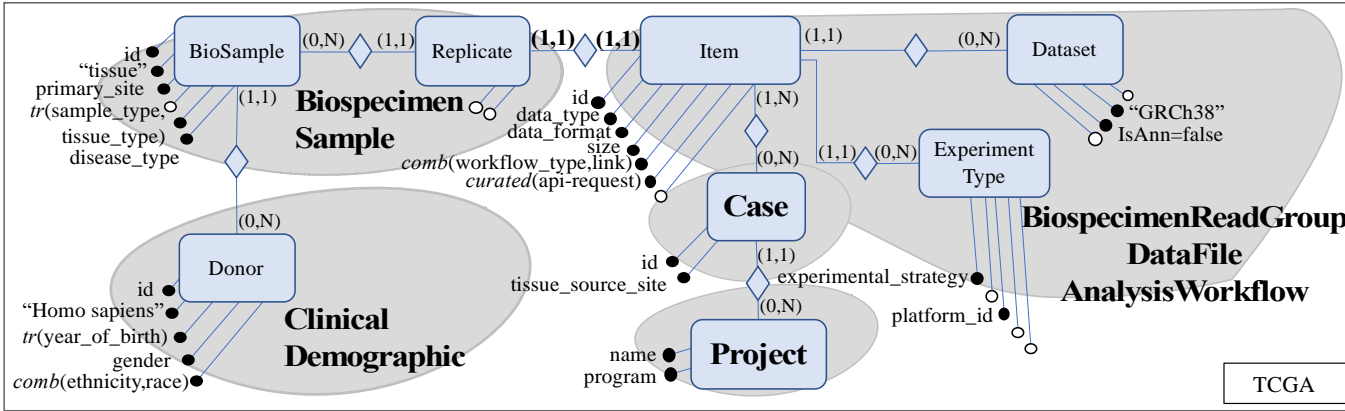
DATA INTEGRATION

TOOLS

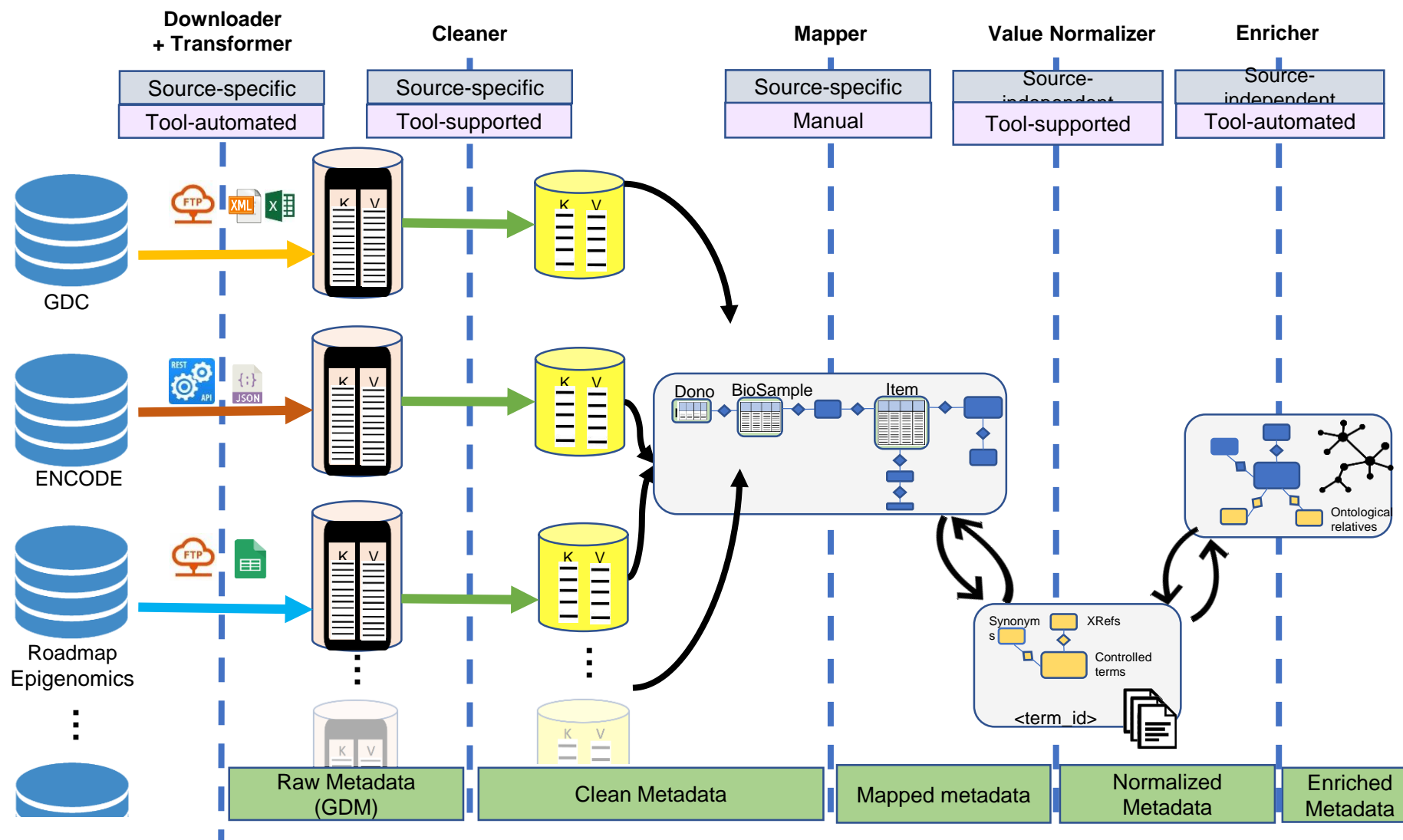
Modeling Genomic Metadata: The Genomic Conceptual Model



Sources elements are mapped into the global schema



META-BASE: phases of integration



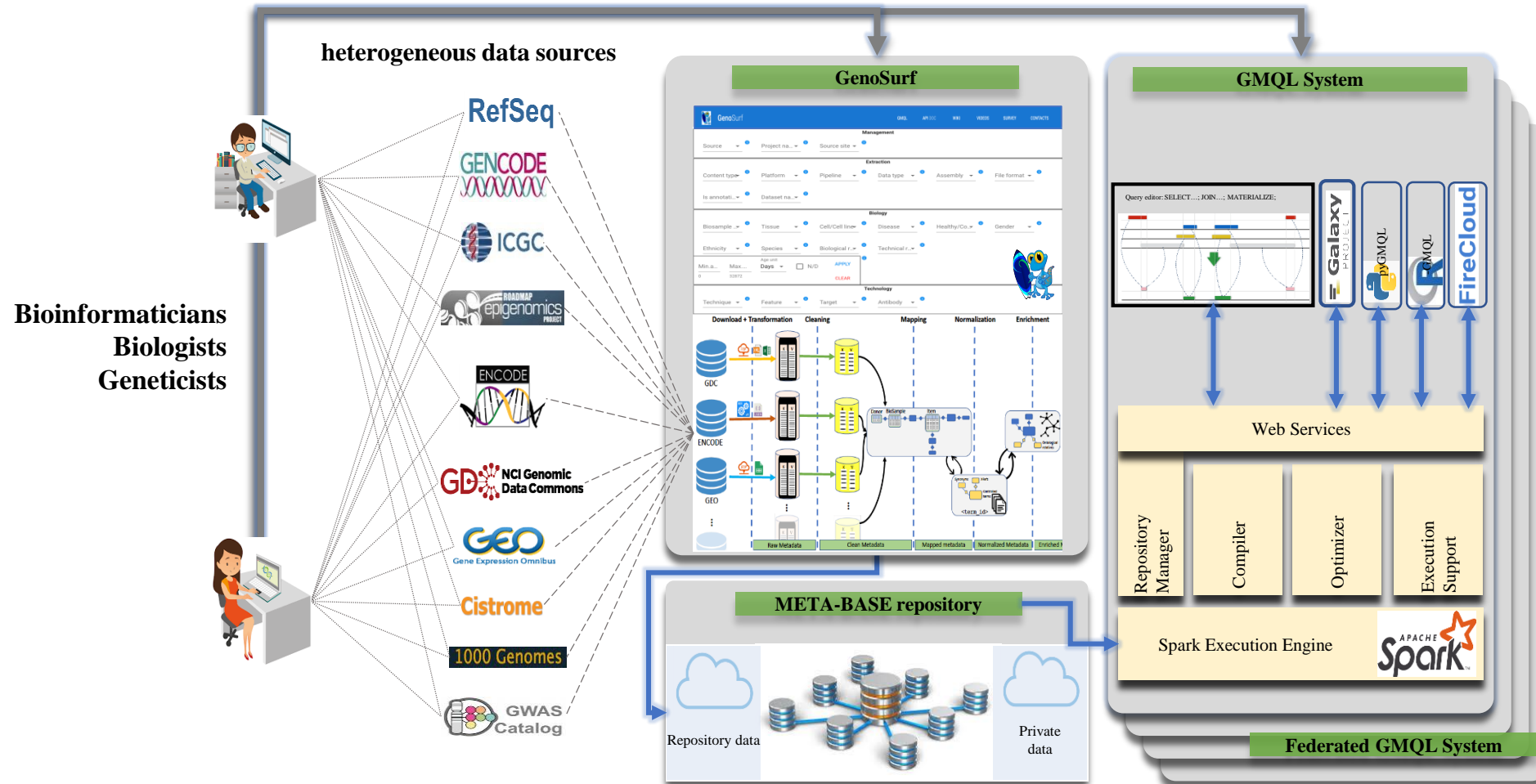
REPOSITORY

Consortium	Imported datasets	# of samples	File size (MB)
ENCODE	GRCh38_ENCODE_BROAD	850	6,869
	GRCh38_ENCODE_NARROW	11,573	128,316
	HG19_ENCODE_BROAD	844	18,382
	HG19_ENCODE_NARROW	10,342	111,925
ROADMAP EPIGENOMICS	HG19_ROADMAP_EPIGENOMICS_BED	156	968
	HG19_ROADMAP_EPIGENOMICS_BROAD	979	24,332
	HG19_ROADMAP_EPIGENOMICS_DMR	66	3,060
	HG19_ROADMAP_EPIGENOMICS_GAPPED	979	6,875
	HG19_ROADMAP_EPIGENOMICS_NARROW	1,032	11,788
	HG19_ROADMAP_EPIGENOMICS_RNA_expression	399	2,453
TCGA	HG19_TCGA_cnv	22,632	797
	HG19_TCGA_dnamethylation	12,860	247,742
	HG19_TCGA_dnaseq	6,914	286
	HG19_TCGA_mirnaseq_isoform	9,909	4,207
	HG19_TCGA_mirnaseq_mirna	9,909	746
	HG19_TCGA_rnaseq_exon	3,675	47,668
	HG19_TCGA_rnaseq_gene	3,675	5,327
	HG19_TCGA_rnaseq_spljxn	3,675	44,377
	HG19_TCGA_rnaseqv2_exon	9,825	124,343
	HG19_TCGA_rnaseqv2_gene	9,825	21,862
	HG19_TCGA_rnaseqv2_isoform	9,825	53,082
	HG19_TCGA_rnaseqv2_spljxn	9,825	115,088
GDC - TCGA	GRCh38_TCGA_copy_number	22,374	686
	GRCh38_TCGA_copy_number_masked	22,375	337
	GRCh38_TCGA_gene_expression	11,091	56,542
	GRCh38_TCGA_methylation	12,218	1,348,516
	GRCh38_TCGA_miRNA_expression	10,947	1,502
	GRCh38_TCGA_miRNA_isoform_expression	10,999	5,004
	GRCh38_TCGA_somatic_mutation_masked	10,188	2,280
GENCODE	GRCh38_ANNOTATION_GENCODE	24	1,798
	HG19_ANNOTATION_GENCODE	20	1,324
REFSEQ	GRCh38_ANNOTATION_REFSEQ	31	740
	HG19_ANNOTATION_REFSEQ	30	275
...
Grand total	67 datasets	564,422	9.23 TB

Experimental datasets and annotations from external databases

67 datasets
564k samples
9.23 TB

GECO Technology – bird eye's view



GMQL user Interface

The screenshot displays the GMQL user interface with the following components:

- Navigation Bar:** GMQL, GMQL-REST, Demo Video, Documentation, Example Queries, GeCo, Hello Demo User, Logout.
- Datasets:** A tree view showing a hierarchy of datasets under 'Public', including various GRCh38 annotations and TCGA data. Buttons for Add, Delete, Download, and UCSC are visible.
- Query editor:** A text area containing a GMQL query:


```

1 myExperiment = SELECT() UPLOADED;
2 myData = COVER(2,ANY) myExperiment;
3
4 genes = SELECT(annotation_type == 'gene'
5               AND provider == 'RefSeq' ) HG19_BED_ANNOTATION;
6 mutations = SELECT(type == "single_base_substitution") ICGC_REPOSITORY;
7
8 insideGene = JOIN(distance < 0;
9                  Output: right) genes myData ;
10
11 mutationCount = MAP() insideGene mutations;
12 mutationCountFilter = SELECT(region:count_insideGene_mutations > 0) mutationCount;
13
14 MATERIALIZE mutationCountFilter into result;
15

```

 Below the editor, there is a 'Query name' field (demo), an 'Output format' section (radio buttons for Tab delimited and GTF), and buttons for Show jobs, Compile, and Execute.
- Metadata browser:** A section for defining query conditions. It shows a query snippet:


```

DATA_SET_VAR = SELECT(annotation_type ==
"gene" AND provider == "RefSeq")
HG19_BED_ANNOTATION;

```

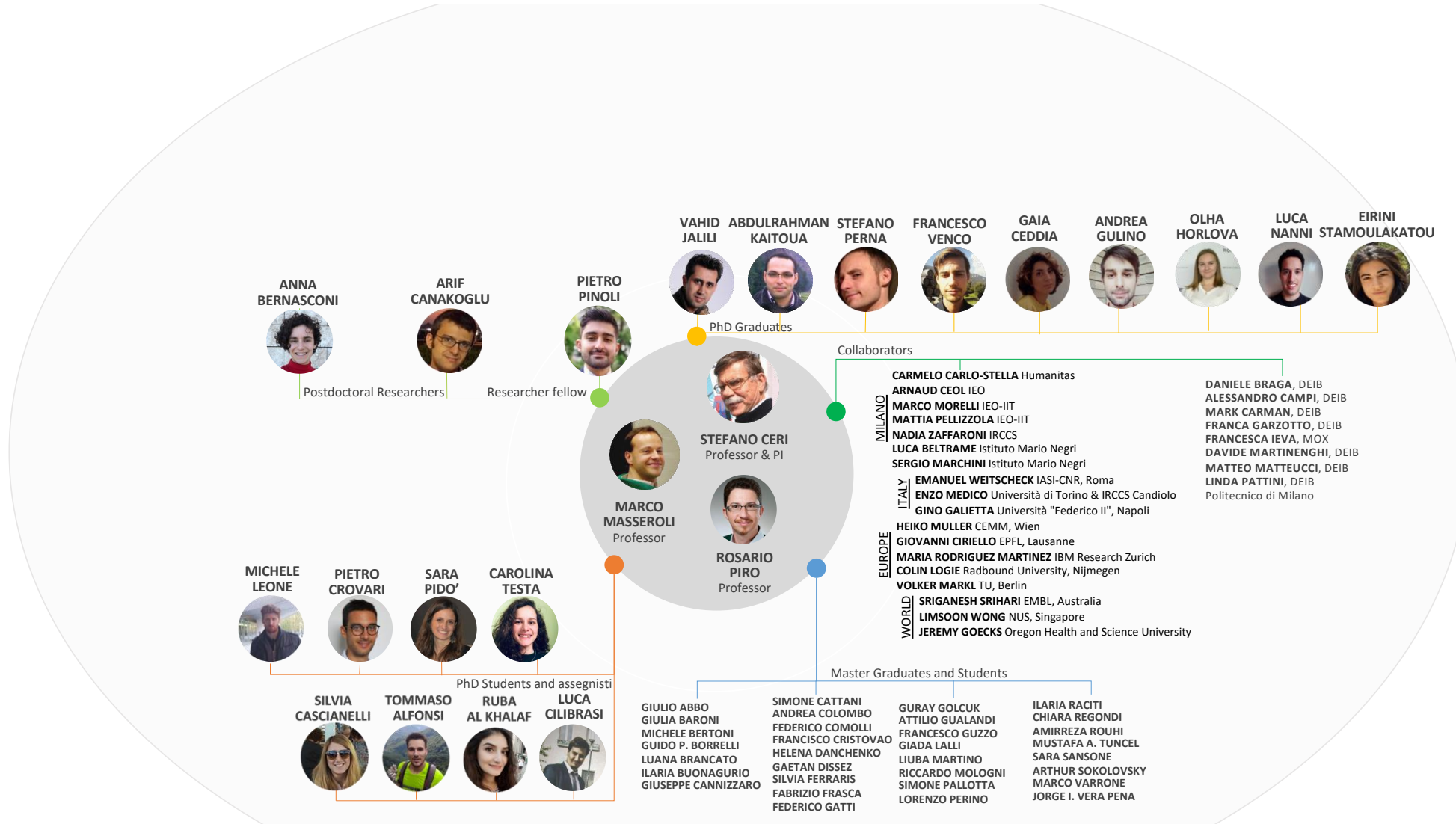
 Below this are buttons for '+ New condition' and 'Test', and two filter conditions: 'annotation_type (3)' set to 'gene' and 'provider (2 - 3)' set to 'RefSeq (1)'.
- Sample metadata:** A table showing attributes and their values:

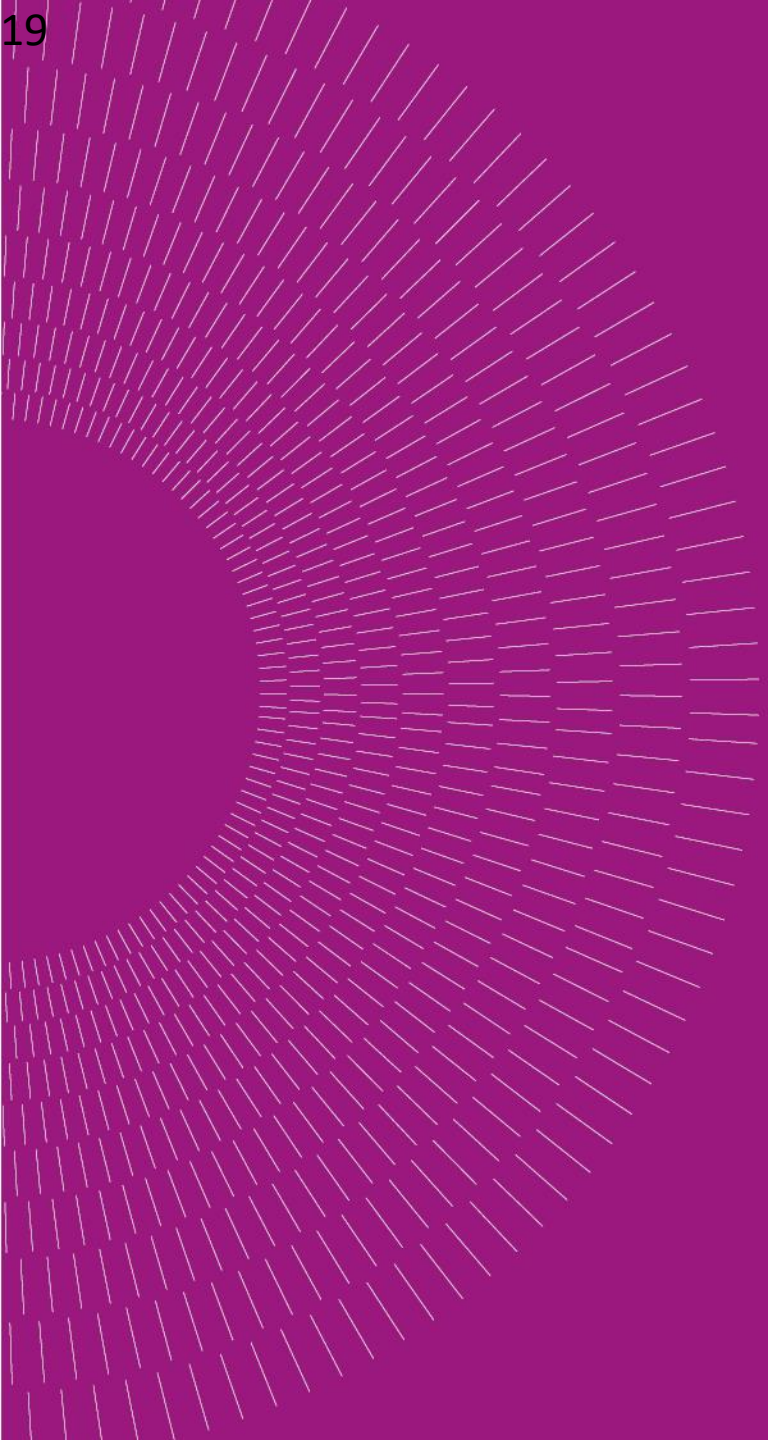
Attribute	Value
annotation_type	gene
assembly	hg19
name	RefSeqGenes
provider	RefSeq
- Schema:** A table showing the schema type 'tab' with columns for Field name, Field type, and Heat map:

Field name	Field type	Heat map
chr	STRING	
left	LONG	
right	LONG	
name	STRING	
score	DOUBLE	
strand	STRING	

Masseroli, M., Pinoli, P., Venco, F., Kaitoua, A., Jalili, V., Palluzzi, F., Muller, H. and Ceri, S., 2015. GenoMetric Query Language: a novel approach to large-scale genomic data management. *Bioinformatics*, 31(12), pp.1881-1888.

GECO Team --- 2016-2021





VIRAL GENOMICS:

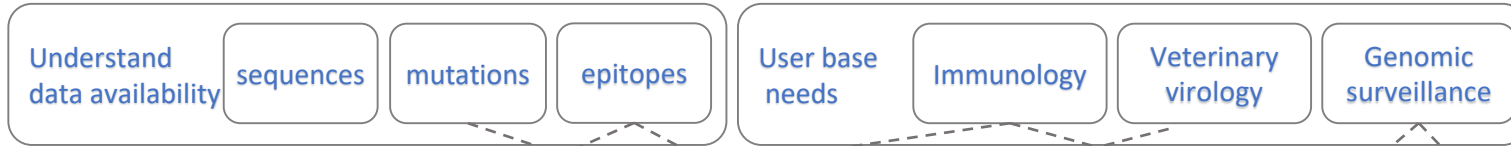
DATA DESIGN

DATA INTEGRATION

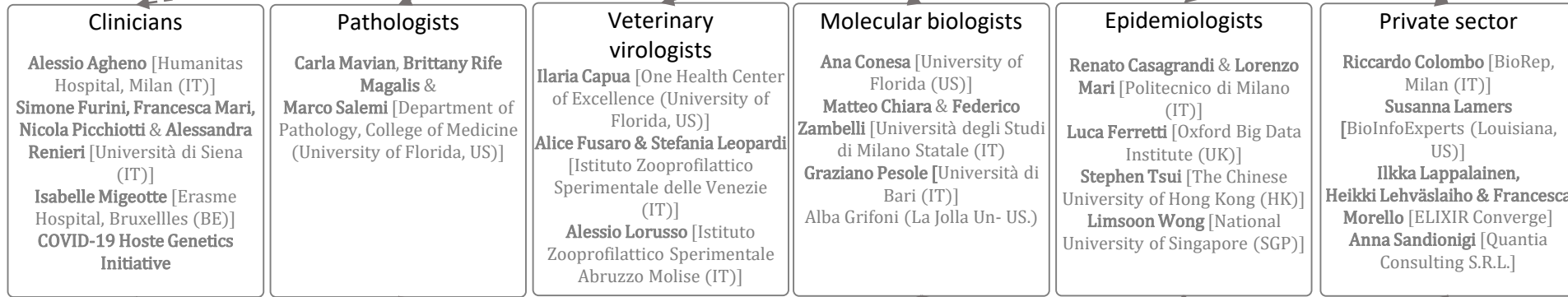
INTERVIEWS (2020)



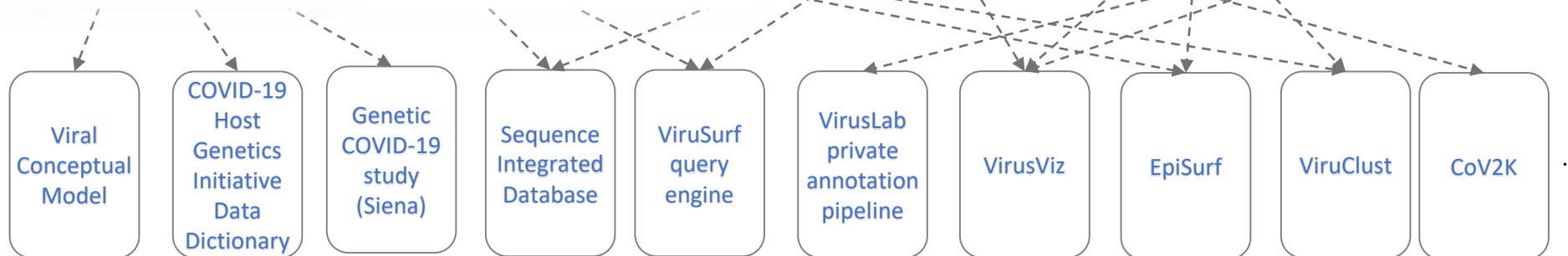
Identification of area of interest



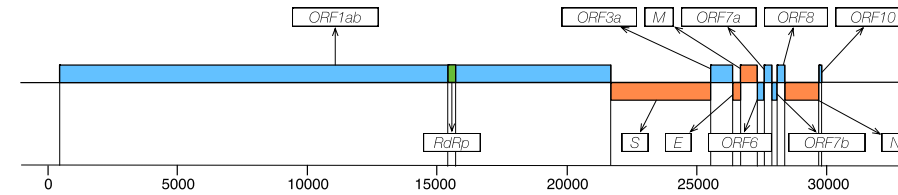
Requirements analysis



Systems and studies design

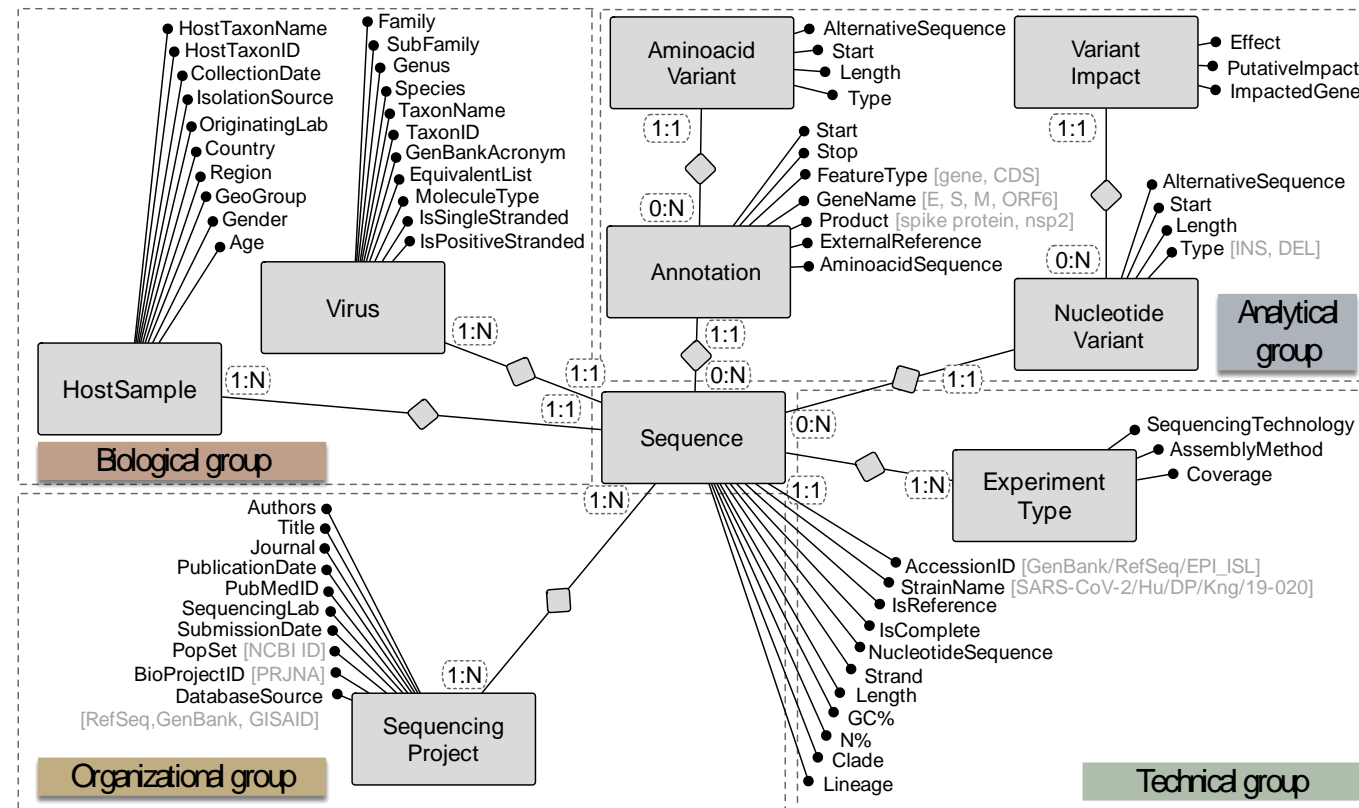


Viral Conceptual Model



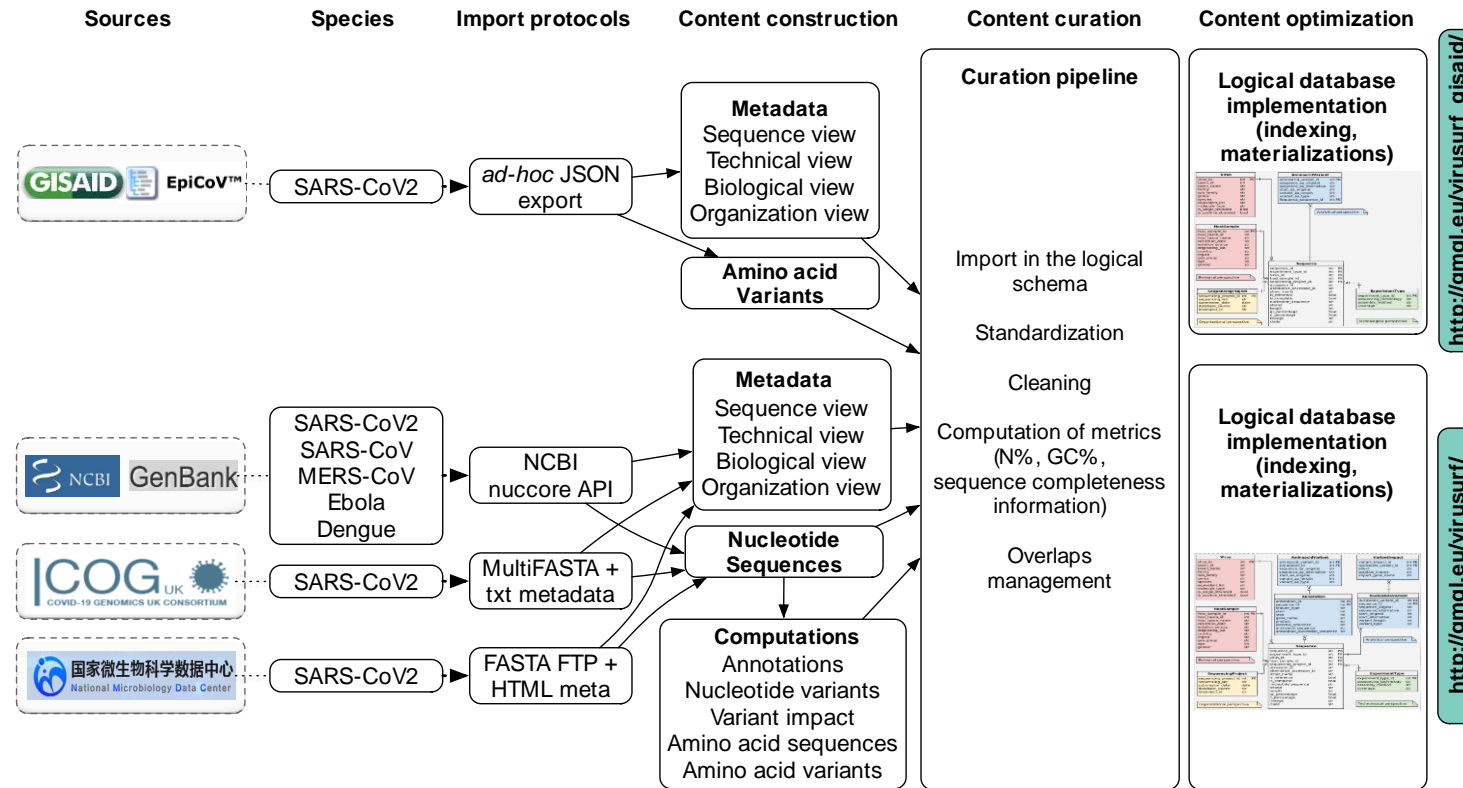
The **Viral Conceptual Model (VCM)**, centered on the virus **sequence** described from four perspectives:

- **biological perspective** (virus species and host environment)
- **technological perspective** (sequencing technology)
- **organizational perspective** (project responsible for producing the sequence)
- **analytical perspective** (properties of the sequence, such as known annotations and variants)



Bernasconi, Canakoglu, Pinoli, Ceri. Empowering virus sequence research through conceptual modeling. In *Proceedings of the International Conference on Conceptual Modeling (ER 2020)*

Data integration pipelines for viral genomes data



January 2022

Sources for ViruSurf-GISAID

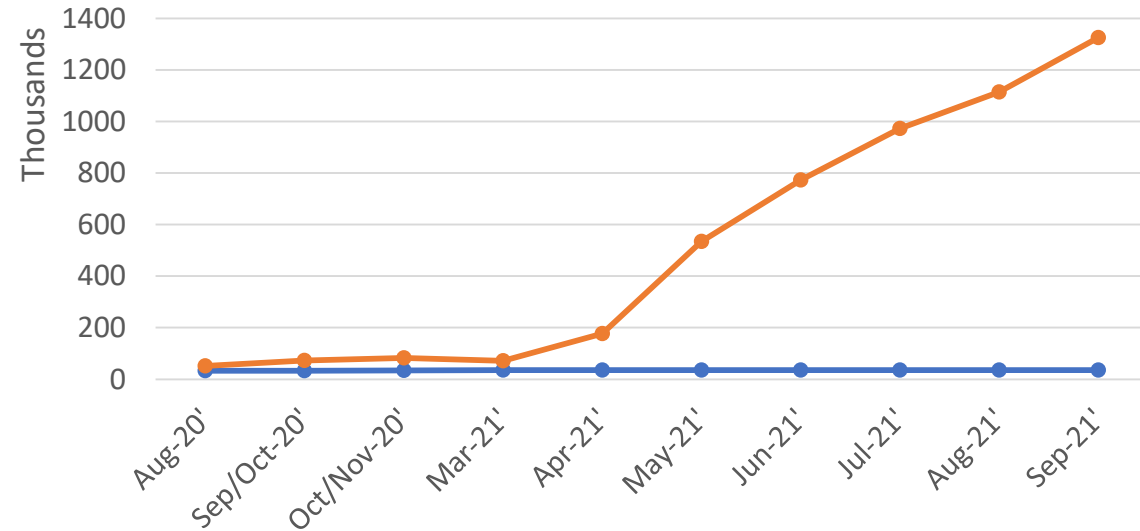
- GISAID EpiCoV™ db
~ 7M sequences

Sources for ViruSurf

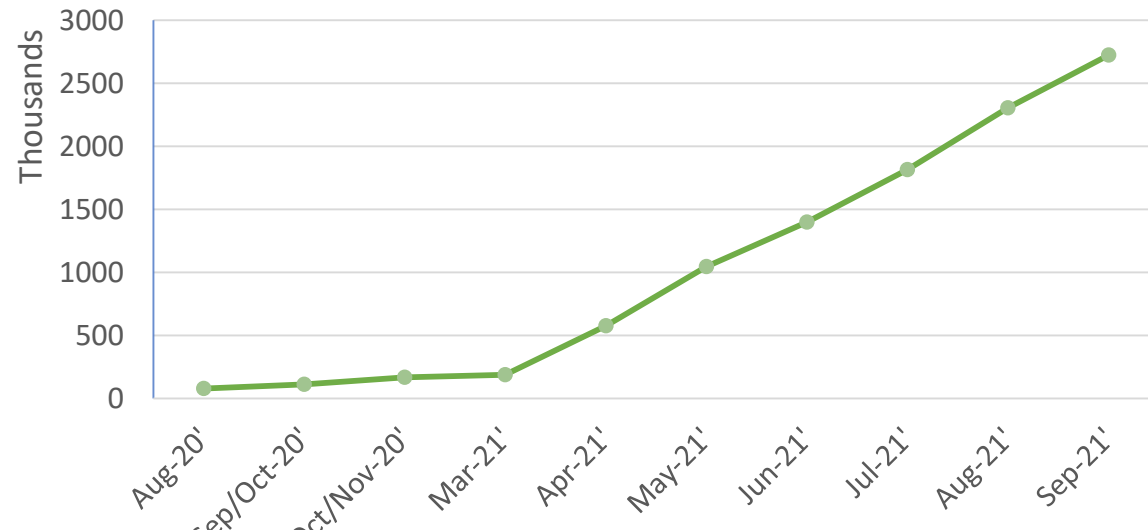
- GenBank ~ 3.2M sequences (SARS-CoV-2)
- GenBank ~ 35K sequences (other viruses)
- COG-UK ~ 800K sequences
- NMDC ~ 300 sequences

Numbers of the import process

Sequences imported from NCBI + COG-UK + NMDC



Sequences imported from GISAID



... ~16 million sequences July 2023!

Virusurf
search system

<http://gmql.eu/virusurf/>

1 Top bar

2 Metadata search

3 Variant search

4 Results visualization

Metadata search

taxon_name: ["severe acute respiratory syndrome coronavirus 2"], is_complete: {true}, n_percentage: {"min_val":null,"max_val":0.5,"is_null":false}

Variant search

Amino acid query:
gene_name: ["n"], sequence_aa_original: ["r"], sequence_aa_alternative: ["k"]
OR
gene_name: ["n"], sequence_aa_original: ["g"], sequence_aa_alternative: ["r"]

Nucleotide query:
n_gene_name: ["n"], start_original: {"min_val":28881,"max_val":28881}

Amino acid query:
product: ["spike (surface glycoprotein)"], start_aa_original: {"min_val":1000}

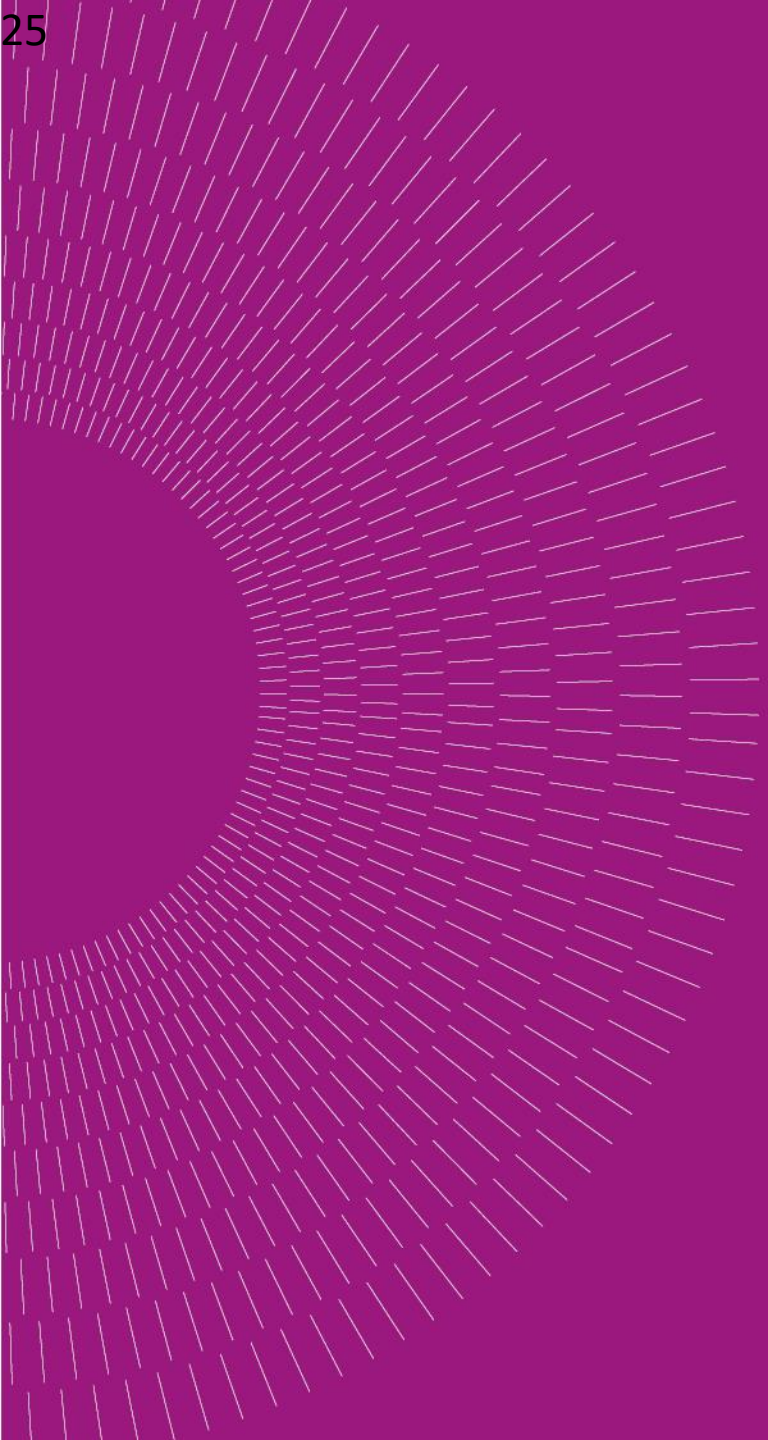
Product protein: spike (surface glycoprotein)
Change type: [dropdown]
Alternative sequence: [dropdown]
Position range: min: 1000

RESULT SEQUENCES

DOWNLOAD TABLE | DOWNLOAD SEQUENCE | Show control | Choose protein name to extract its sequence: FULL | SELECT/SORT FIELDS

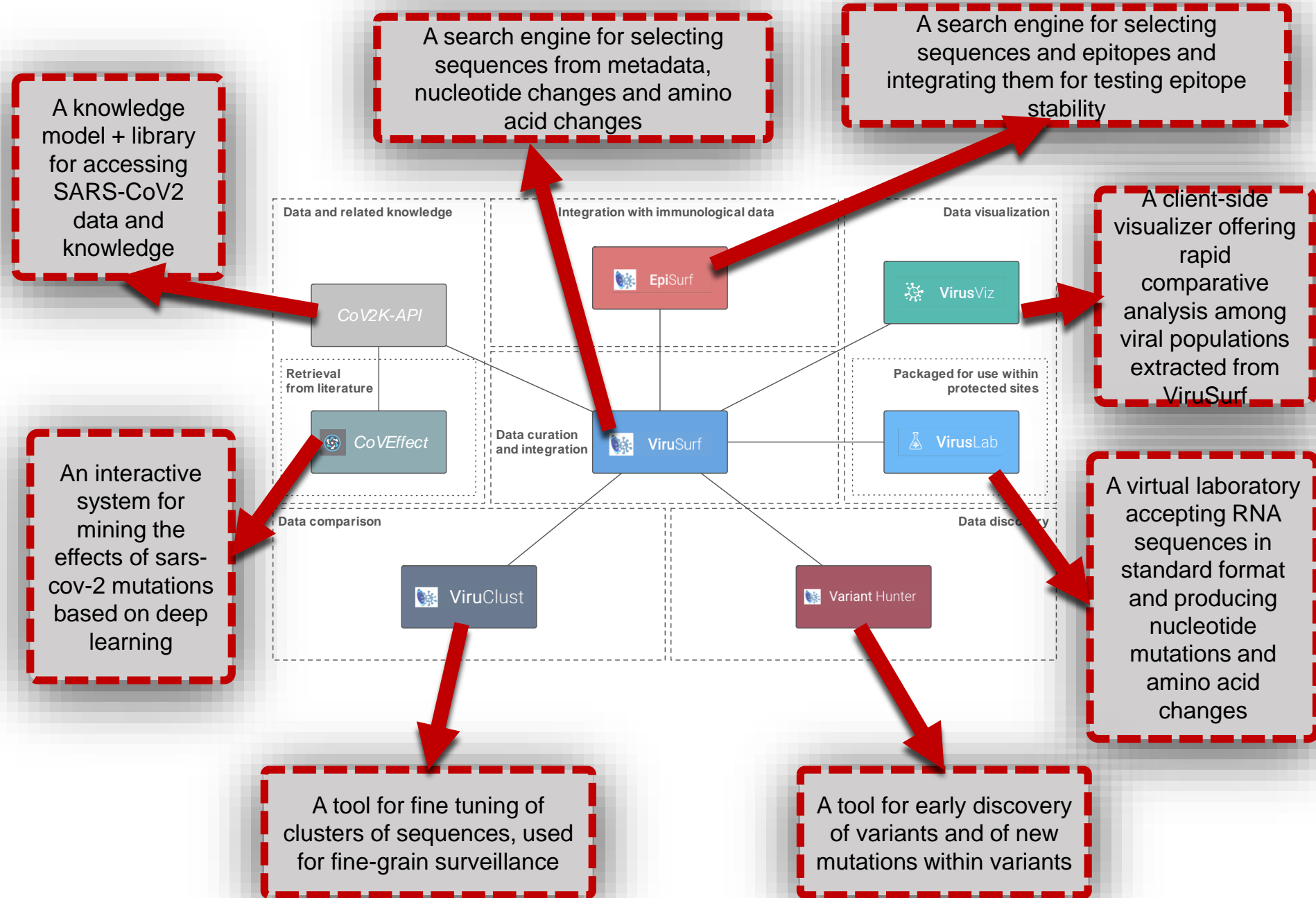
Source Page	Accession ID ↑	Strain name	Is reference	Is complete	Strand	Sequence Length	GC%	N%	Lineage (Clade)	Seq. Technology	Assembly Method	Coverage	Submission date
link	England/BIRM-5F8FB/2020	England/BIRM-5F8FB/2020	False	True	positive	29903	37.99	0.42	B.1.1.1 (N/D)	N/D	N/D	N/D	2020-03-04
link	England/BIRM-611BF/2020	England/BIRM-611BF/2020	False	True	positive	29903	37.98	0.41	B.1.1.1 (N/D)	N/D	N/D	N/D	2020-03-04
link	England/BRIS-18531B5/2020	England/BRIS-18531B5/2020	False	True	positive	29903	38.0	0.4	B.1.1 (N/D)	N/D	N/D	N/D	2020-03-04

Rows per page: 10 | 1-10 of 14 | 14 sequences found



TOOLS

Tools for viral genomics



<https://gmql.eu/epivirusurf>

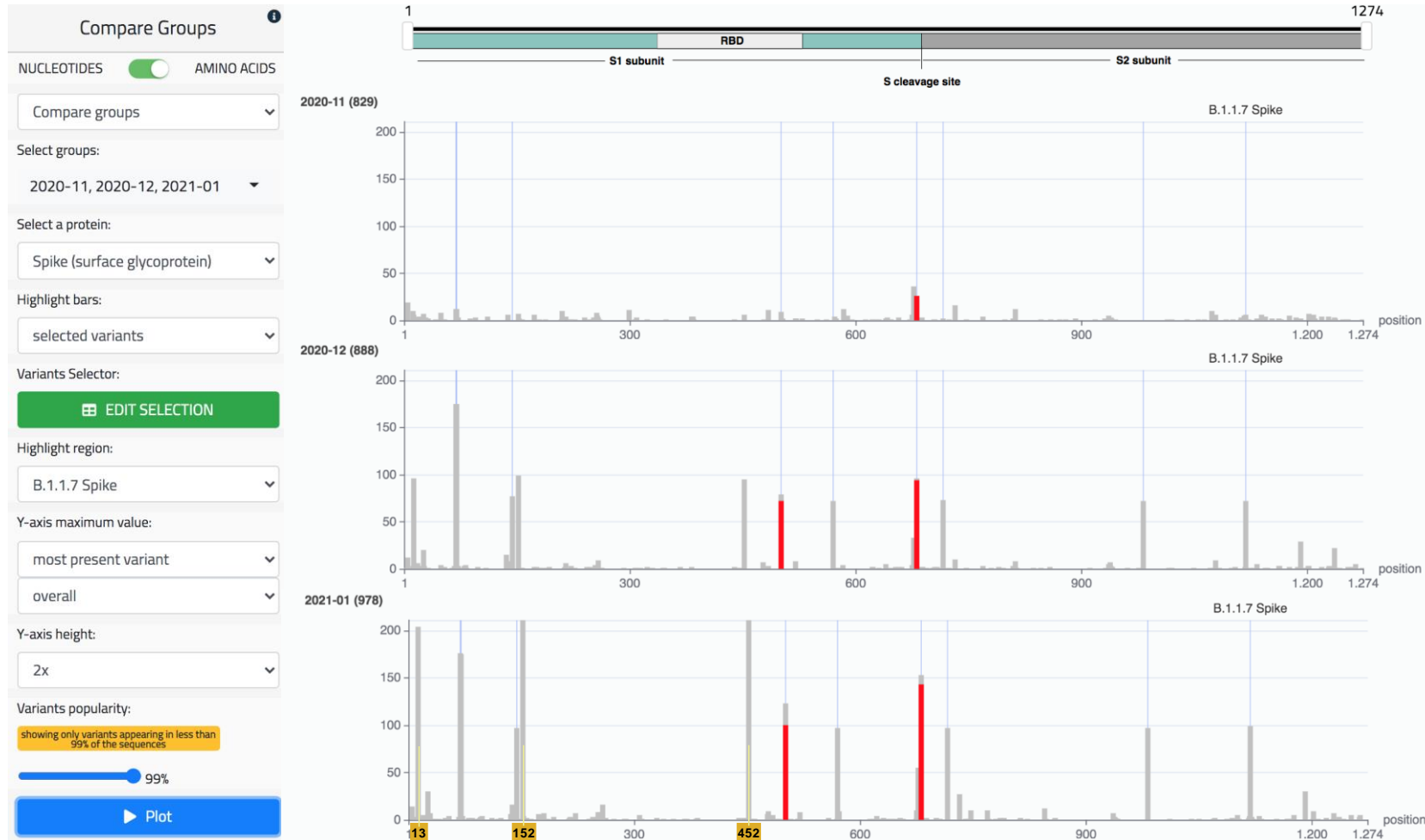
The screenshot shows the EpiSurf web application interface, divided into four main sections:

- Section 1 (Top bar):** Contains the EpiSurf logo, navigation links (EPISURF GISAID, VIRUSURF, VIRUSURF GISAID, GENOSURF, WIKI, ABOUT), a search bar with a "CLEAR YOUR QUERY" button, and a "Last update date: 2021-05-09" indicator.
- Section 2 (Sequence population search):** A form for searching sequences. It includes a "Virus" section with "Virus taxon name" (severe acute respiratory syndrome coronavirus 2) and "Virus species". A "Host Organism" section includes "Host taxon name" (homo sapiens), "Collection date", "Isolation source", "Continent", "Country", and "Region" (florida). A "Sequence properties and technology" section includes "Is reference", "Is complete", "Strand", "Sequence Length", "GC%", "N%", "Lineage", and "Sequencing technology". An "Organization" section includes "Database source".
- Section 3 (Epitope/Variant search):** A form for searching epitopes. It includes "Protein Name" (Spike (surface glycoprotein)), "Assay" (T cell), "HLA restriction" (HLA-A*02:01), "Is Linear", "Response Frequency", "Position Range" (min: 331, max: 524), and "Epitope IEDB ID". Buttons for "APPLY EPILOPE SEARCH" and "ADD CONDITION ON AMINO ACIDS" are present.
- Section 4 (Result visualization):** A table displaying search results. The table has columns for "EPILOPE IEDB ID", "REF PAGE", "VIRUSVIZ MUTATED SEQ", "VIRUSVIZ ALL POPULATION", "HLA RESTRICTION", "RESPONSE FREQUENCY", "EPILOPE SEQ", "POSITION RANGE", "NUM MUT SEQ", "TOT MUT", "MUT FREQ", and "MUT SEQ RATIO". The table contains several rows of data, including entries for HLA-A*02:01 with various epitope sequences and mutation statistics.

Top bar
Sequence population search
Epitope/Variant search
Result visualization

Bernasconi, A., Cilibrasi, L., Al Khalaf, R., Alfonsi, T., Ceri, S., Pinoli, P., & Canakoglu, A. (2021). EpiSurf: metadata-driven search server for analyzing amino acid changes within epitopes of SARS-CoV-2 and other viral species. Database, 2021, baab059.

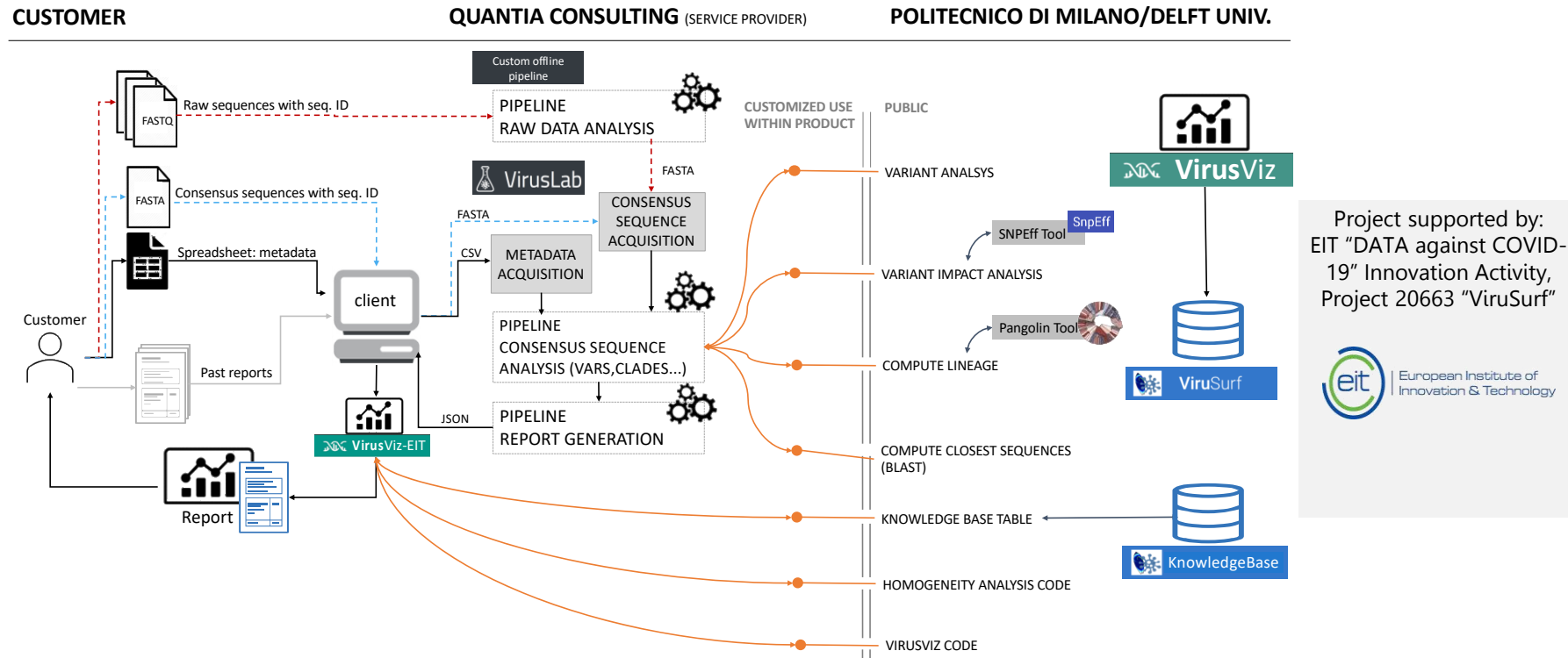
<https://gmql.eu/virusviz>



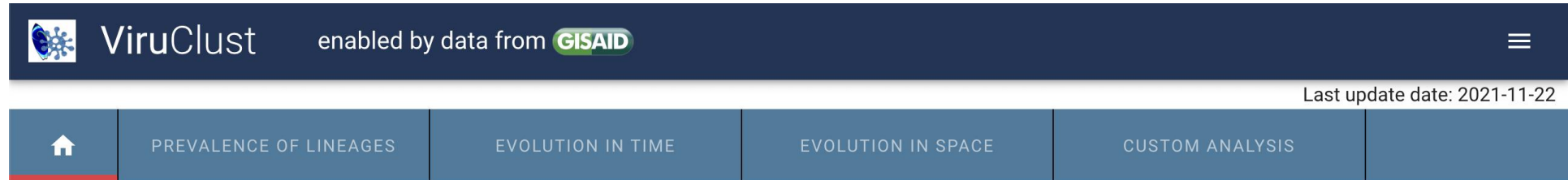
Bernasconi, A., Gulino, A., Alfonsi, T., Canakoglu, A., Pinoli, P., Sandionigi, A., & Ceri, S. (2021). VirusViz: comparative analysis and effective visualization of viral nucleotide and amino acid variants. *Nucleic Acids Research*, 49(15), e90-e90.

VirusLab: customizable services for use within protected sites

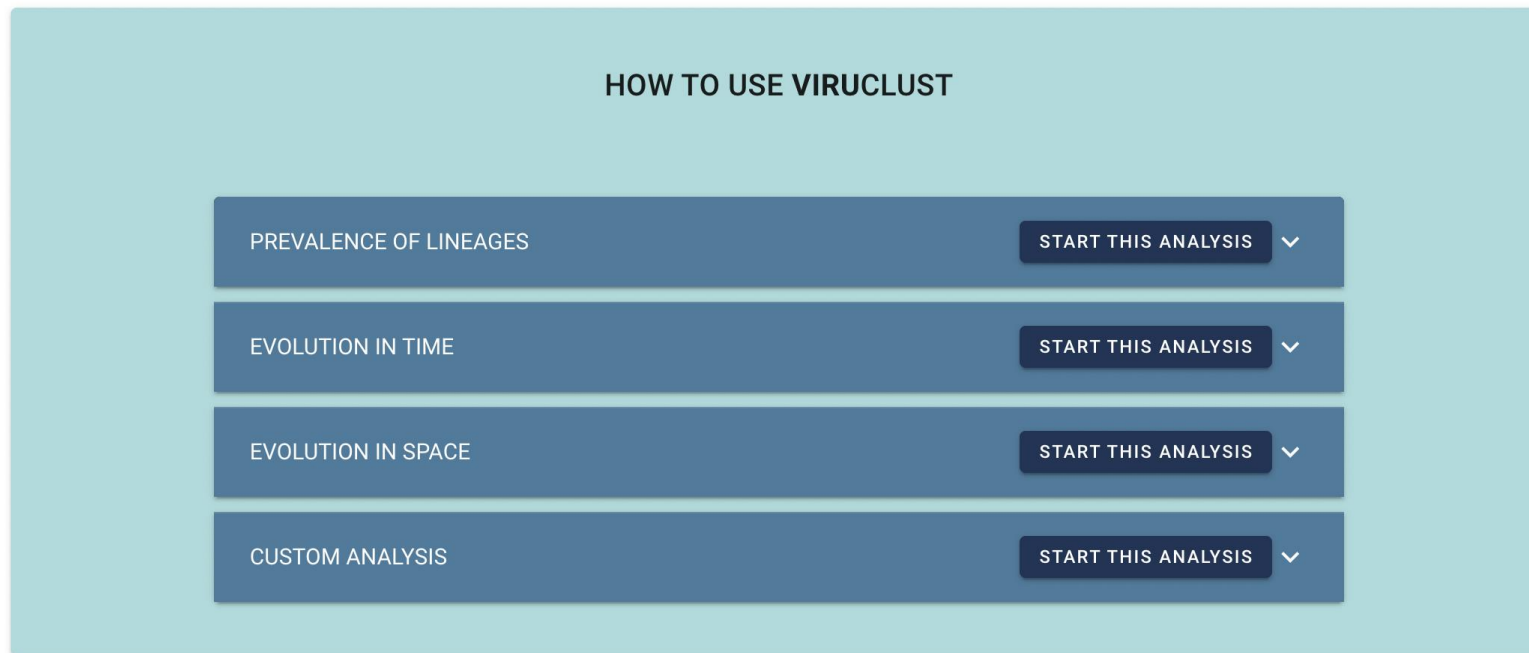
For supporting a laboratory wishing to perform secondary (row data) analysis and to add sensible metadata (e.g. clinical), still using our data and knowledge bases and visualization tools



<https://gmql.eu/virusclust>



The header navigation bar features the VirusClust logo on the left, followed by the text "enabled by data from GISAID". On the right side, there is a hamburger menu icon and the text "Last update date: 2021-11-22". Below this, a horizontal menu contains five items: a home icon, "PREVALENCE OF LINEAGES", "EVOLUTION IN TIME", "EVOLUTION IN SPACE", and "CUSTOM ANALYSIS".



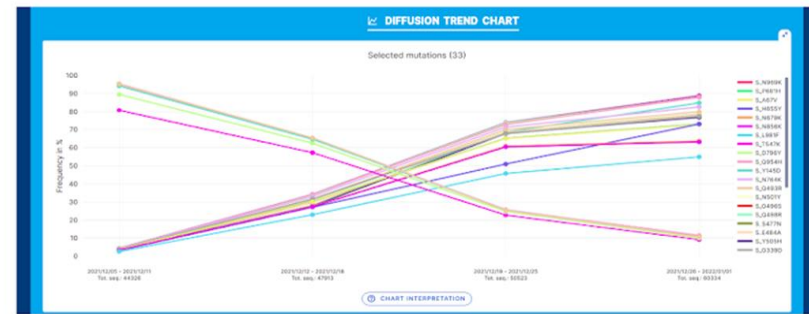
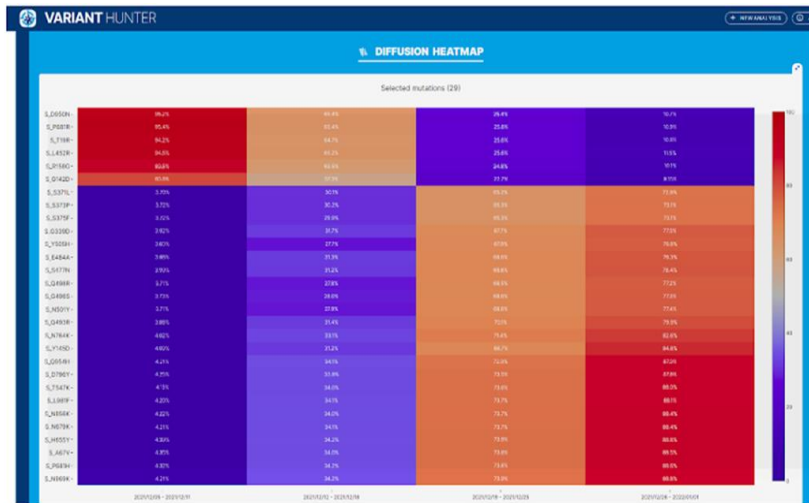
The "HOW TO USE VIRUCLUST" guide is presented as a list of four options, each with a corresponding "START THIS ANALYSIS" button and a dropdown arrow:

- PREVALENCE OF LINEAGES
- EVOLUTION IN TIME
- EVOLUTION IN SPACE
- CUSTOM ANALYSIS

Cilibrasi, L., Pinoli, P., Bernasconi, A., Canakoglu, A., Chiara, M., & Ceri, S. (2022). VirusClust: direct comparison of SARS-CoV-2 genomes and genetic variants in space and time. *Bioinformatics*, 38(7), 1988-1994.

https://gmql.eu/variant_hunter

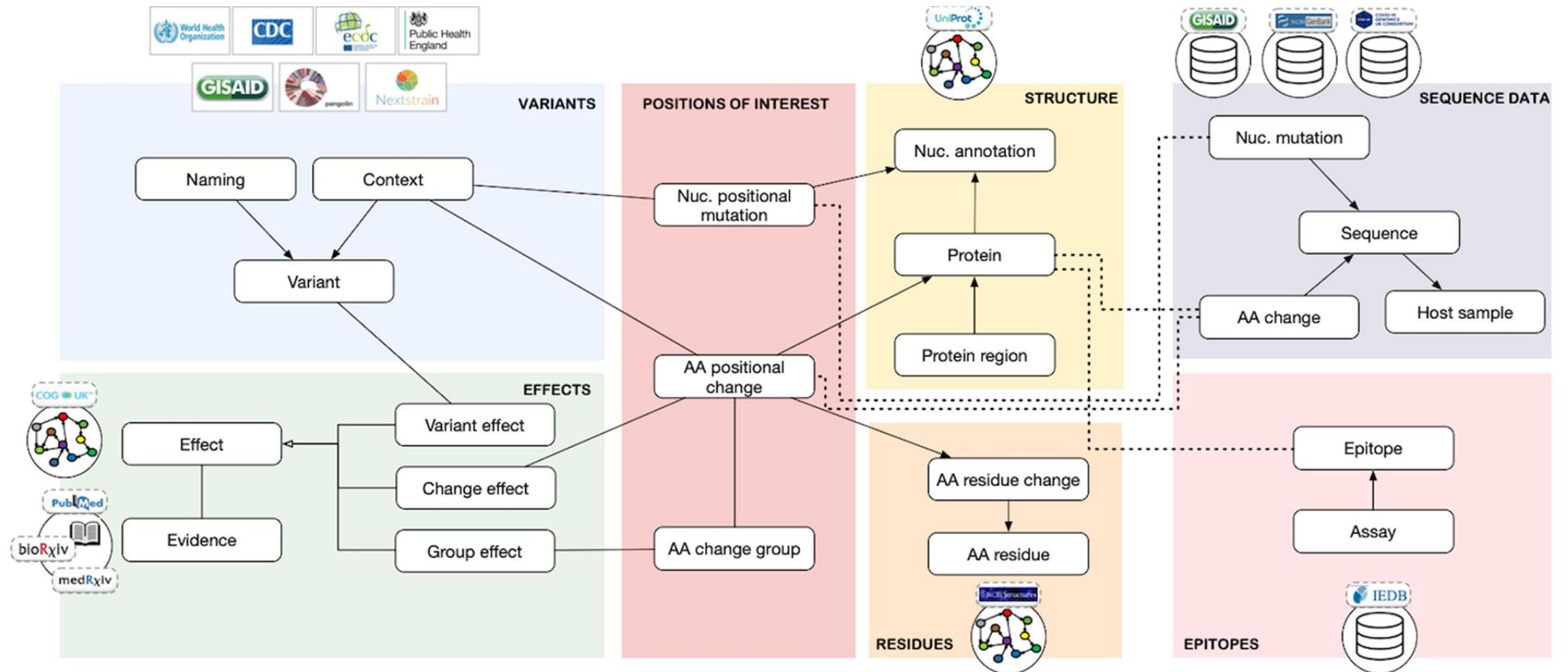
In the example: substitution of Delta variant by Omicron variant in Europe and North America



Pinoli, P., Canakoglu, A., Ceri, S., Chiara, M., Ferrandi, E., Minotti, L., & Bernasconi, A. (2023). VariantHunter: a method and tool for fast detection of emerging SARS-CoV-2 variants. Database, 2023, baad044.

CoV2K – Knowledge and data interplay

<http://gmql.eu/cov2k/api/>



<https://gmql.eu/coveffect>

ABOUT
CovEffect

Additional Search
New Session
Save Session
Load Session
Session Name: No Session Name

Annotated Papers
Selected Papers

Abstract

The **sensitivity** of SARS-CoV-2 variants of concern (VOCs) to neutralizing **antibodies** has largely been studied in the context of key receptor binding domain (RBD) mutations, including E484K and N501Y. Little is known about the epistatic effects of combined SARS-CoV-2 spike mutations. We now investigate the neutralization **sensitivity** of variants containing the non-RBD mutation Q677H, including B.1.525 (Nigerian isolate) and Bluebird (U.S. isolate) variants. The effect on neutralization of Q677H was determined in the context of the RBD mutations and in the background of major VOCs, including B.1.1.7 (United Kingdom, Alpha), B.1.351 (South Africa, Beta), and P1-501Y-V3 (Brazil, Gamma). We demonstrate that the Q677H mutation increases viral **infectivity** and syncytium formation, as well as **enhancing resistance** to **neutralization** for VOCs, including B.1.1.7 and P1-501Y-V3. Our work highlights the importance of epistatic interactions between SARS-CoV-2 spike mutations and the continued need to monitor Q677H-bearing VOCs.

Paper Info

Remove Paper

Annotated Papers:
0/2

ReadyToTrain Papers:
0

DOI:
[10.1128/mbio.02510-21](https://doi.org/10.1128/mbio.02510-21)

Title:
Neutralization of SARS-CoV-2 Variants of Concern Harboring Q677H

Authors:
Zeng, Cong; Evans, John P; Faraone, Julia N.; Qu, Panke; Zheng, Yi-Min; Saif, Linda; Oltz, Eugene M.; Lozanski, Gerard; Gumina, Richard J.; Liu, Shan-Lu

Year:
2021

Extracted Annotations

1:	Entity Type [100%] variant	Mutation(s)/Varia... B.1.1.7	Effect [56%] infectivity	Level [100%] higher		
2:	Entity Type [100%] variant	Mutation(s)/Varia... B.1.1.7	Effect [55%] sensitivity_to_a	Level [100%] lower		
3:	Entity Type [100%] variant	Mutation(s)/Varia... B.1.351	Effect [53%] infectivity	Level [100%] higher		
4:	Entity Type [100%] variant	Mutation(s)/Varia... B.1.351	Effect [54%] sensitivity_to_a	Level [100%] lower		
5:	Entity Type [100%] single	Mutation(s)/Varia... SPIKE_Q677H	Effect [100%] infectivity	Level [100%] higher		
6:	Entity Type [100%] single	Mutation(s)/Varia... SPIKE_Q677H	Effect [34%] sensitivity_to_a	Level [100%] lower		

SAVE

Editor

Selected Attribute:
Level

Confidence:
100%

Abstract doesn't contain this information

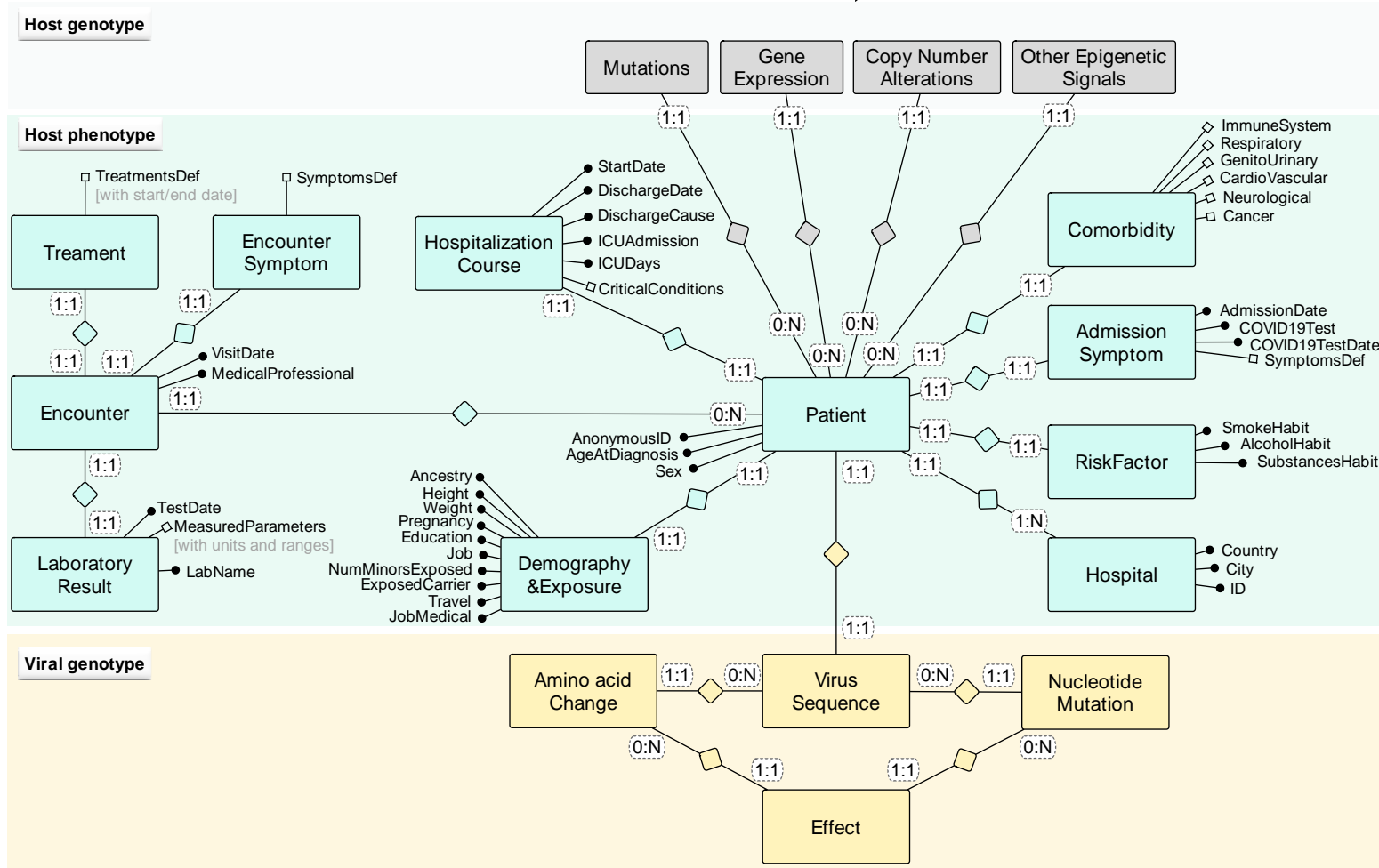
higher ▼



CLINICAL RESEARCH

COVID-19 Conceptual Model

Described by the Viral Conceptual Model



Described by the Genomic Conceptual Model

<http://gmql.eu/genosurf/>

GenoSurf

e.g., data dictionary of <https://www.covid19hg.org/>

PhenotypeDB

<http://gmql.eu/virusurf/>

Virusurf

Article | [Open Access](#) | [Published: 08 July 2021](#)

Mapping the human genetic architecture of COVID-19

[COVID-19 Host Genetics Initiative](#)

[Nature](#) **600**, 472–477 (2021) | [Cite this article](#)

284k Accesses | **370** Citations | **1900** Altmetric | [Metrics](#)

Article | [Open Access](#) | [Published: 07 March 2022](#)

Whole-genome sequencing reveals host factors underlying critical COVID-19

[Athanasios Kousathanas](#), [Erola Pairo-Castineira](#), [Konrad Rawlik](#), [Alex Stuckey](#), [Christopher A. Odhams](#), [Susan Walker](#), [Clark D. Russell](#), [Tomas Malinauskas](#), [Yang Wu](#), [Jonathan Millar](#), [Xia Shen](#), [Katherine S. Elliott](#), [Fiona Griffiths](#), [Wilna Oosthuyzen](#), [Kirstie Morrice](#), [Sean Keating](#), [Bo Wang](#), [Daniel Rhodes](#), [Lucija Klaric](#), [Marie Zechner](#), [Nick Parkinson](#), [Afshan Siddiq](#), [Peter Goddard](#), [Sally Donovan](#), [GenOMICC investigators](#), [23andMe investigators](#), [COVID-19 Human Genetics Initiative](#), ... [J. Kenneth Baillie](#) 

[+ Show authors](#)

[Nature](#) **607**, 97–103 (2022) | [Cite this article](#)

135k Accesses | **94** Citations | **1560** Altmetric | [Metrics](#)

Matters Arising | [Open Access](#) | [Published: 03 August 2022](#)

A first update on mapping the human genetic architecture of COVID-19

[COVID-19 Host Genetics Initiative](#)

[Nature](#) **608**, E1–E10 (2022) | [Cite this article](#)

30k Accesses | **57** Citations | **217** Altmetric | [Metrics](#)

Research empowered by availability of clinical + genetic data

Article | [Open Access](#) | Published: 17 January 2021

Employing a systematic approach to biobanking and analyzing clinical and genetic data for advancing COVID-19 research

Sergio Daga, Chiara Fallerini, Margherita Baldassarri, Francesca Fava, Floriana Valentino, Gabriella Doddato, Elisa Benetti, Simone Furini, Annarita Giliberti, Rossella Tita, Sara Amitrano, Mirella Bruttini, Ilaria Meloni, Anna Maria Pinto, Francesco Raimondi, Alessandra Stella, Filippo Biscarini, Nicola Picchiotti, Marco Gori, Pietro Pinoli, Stefano Ceri, Maurizio Sanarico, Francis P. Crawley, Giovanni Biolo, GEN-COVID Multicenter Study, ... [Elisa Frullanti](#) + Show authors

European Journal of Human Genetics 29, 745–759 (2021) | [Cite this article](#)

Cardiol Cardiovasc Med 2021; 5 (6): 673-694

DOI: 10.26502/jcm.92920232



Research Article

Post-Mendelian Genetic Model in COVID-19

Nicola Picchiotti^{1,2#}, Elisa Benetti^{3#}, Chiara Fallerini^{3,4#}, Sergio Daga^{3,4}, Margherita Baldassarri^{3,4}, Francesca Fava^{3,4,5}, Kristina Zguro³, Floriana Valentino^{3,4}, Gabriella Doddato^{3,4}, Annarita Giliberti^{3,4}, Rossella Tita⁵, Sara Amitrano⁵, Mirella Bruttini^{3,4,5}, Laura Di Sarno^{3,4}, Diana Alaverdian^{3,4}, Giada Beligni^{3,4}, Maria Palmieri^{3,4}, Susanna Croci^{3,4}, Mirjam Lista^{3,4}, Ilaria Meloni^{3,4}, Anna Maria Pinto⁵, Chiara Gabbi⁶, Stefano Ceri⁷, Antonio Esposito⁷, Pietro Pinoli⁷, Francis P. Crawley⁸, Elisa Frullanti^{3,4}, Francesca Mari^{3,4,5}, GEN-COVID Multicenter Study, Marco Gori^{1,9}, Alessandra Renieri^{3,4,5*}, Simone Furini^{3*}

¹University of Siena, DIISM-SAILAB, Siena, Italy

²Department of Mathematics, University of Pavia, Pavia, Italy

³Med Biotech Hub and Competence Center, Department of Medical Biotechnologies, University of Siena, Italy

⁴Medical Genetics, University of Siena, Italy

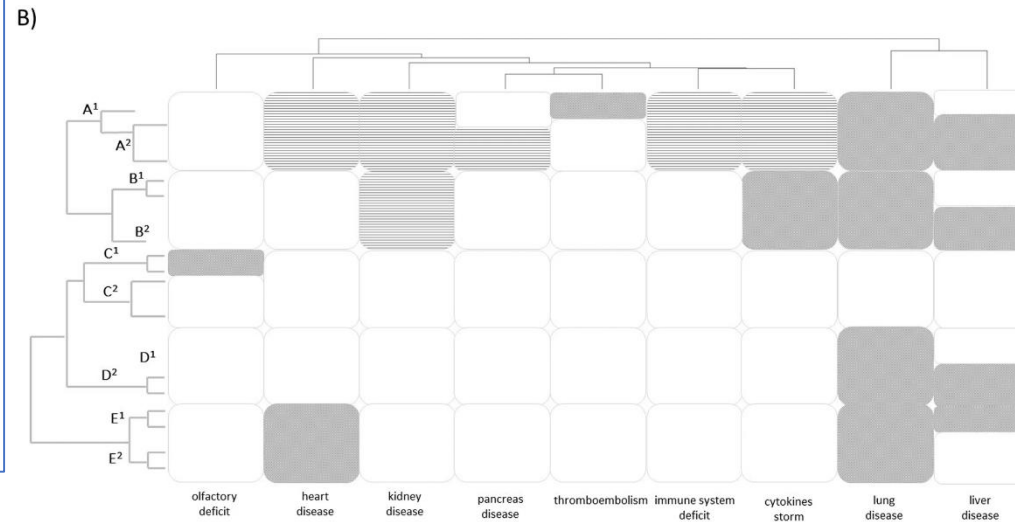
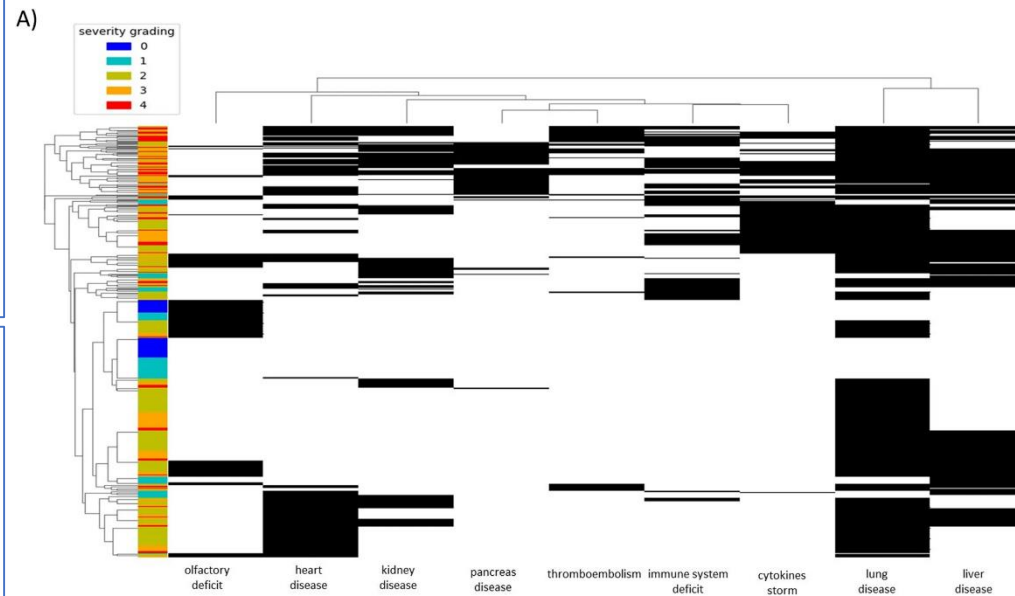
⁵Genetica Medica, Azienda Ospedaliero-Universitaria Senese, Italy

⁶Department of Biosciences and Nutrition, Karolinska Institutet, Stockholm, Sweden

⁷Politecnico di Milano, DEIB, Milano, Italy

⁸Good Clinical Practice Alliance-Europe (GCPA) and Strategic Initiative for Developing Capacity in Ethical Review-Europe (SIDCER), Leuven, Belgium.

⁹Universite Côte d'Azur, Inria, CNRS, I3S, Maasai



Association between SARS-CoV-2 viremia and COVID-19 mortality

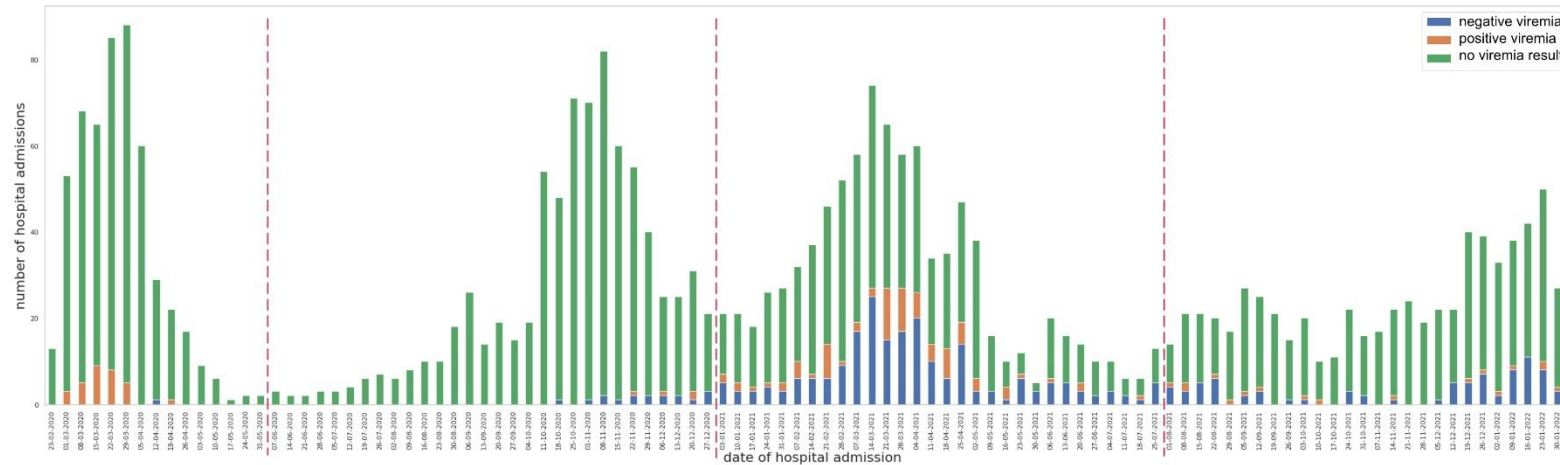
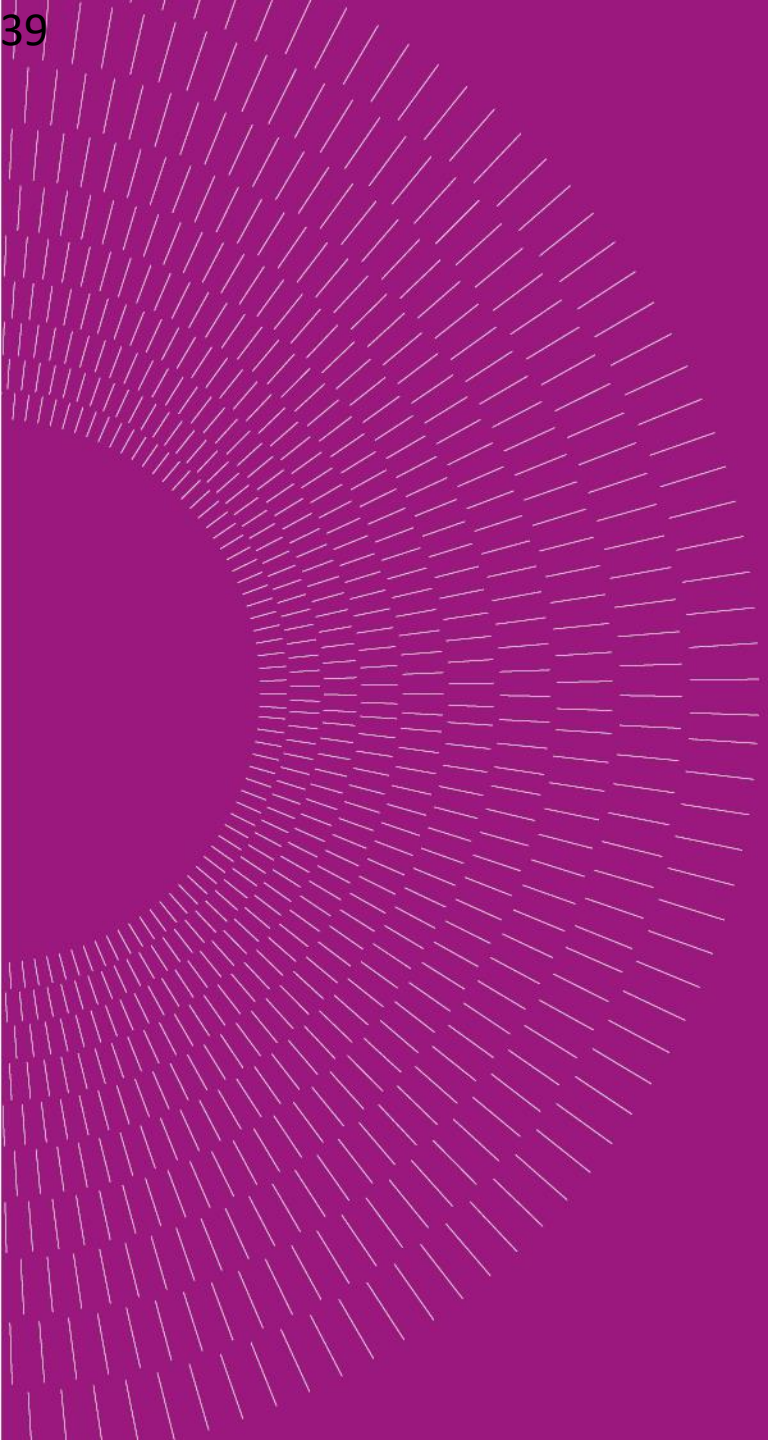


Table 2. Logistic regression analysis of factors associated with death.

Characteristics	OR	95%CI	aOR	95%CI
Females vs Males	0.57	0.35–0.94	0.63	0.38–1.03
Age (per 1 year more)	1.01	1.01–1.02	1.09	1.06–1.11
CCI (per one point more)	0.99	0.98–1.02	–	–
SARS-CoV-2 epidemic waves				
1 vs 2	1.71	0.54–5.41	–	–
1 vs 3	1.18	0.57–2.42	–	–
1 vs 4	0.72	0.31–1.67	–	–
Days from symptoms onset (per one day more)	1.00	1.00–1.01	–	–
Disease severity (Mild/moderate vs severe/critical)	1.85	1.12–3.06	1.55	0.95–2.54
Doses of vaccine				
0 vs 1	0.97	0.35–2.66	–	–
0 vs 2	1.02	0.45–2.30	–	–
0 vs 3	1.16	0.24–5.71	–	–
Positive SARS-CoV-2 viremia vs negative	5.16	3.15–8.45	6.48	4.05–10.55

List of abbreviations: OR, odds ratio; CI, confidence interval; aOR, adjusted odds ratio; CCI, Charlson comorbidity index.



METHODS FOR VIRAL GENOMICS

Time-series analysis of viral amino acid changes

Focus

Co-occurrence of mutations (\rightarrow amino acid changes) on the virus (\rightarrow variant)

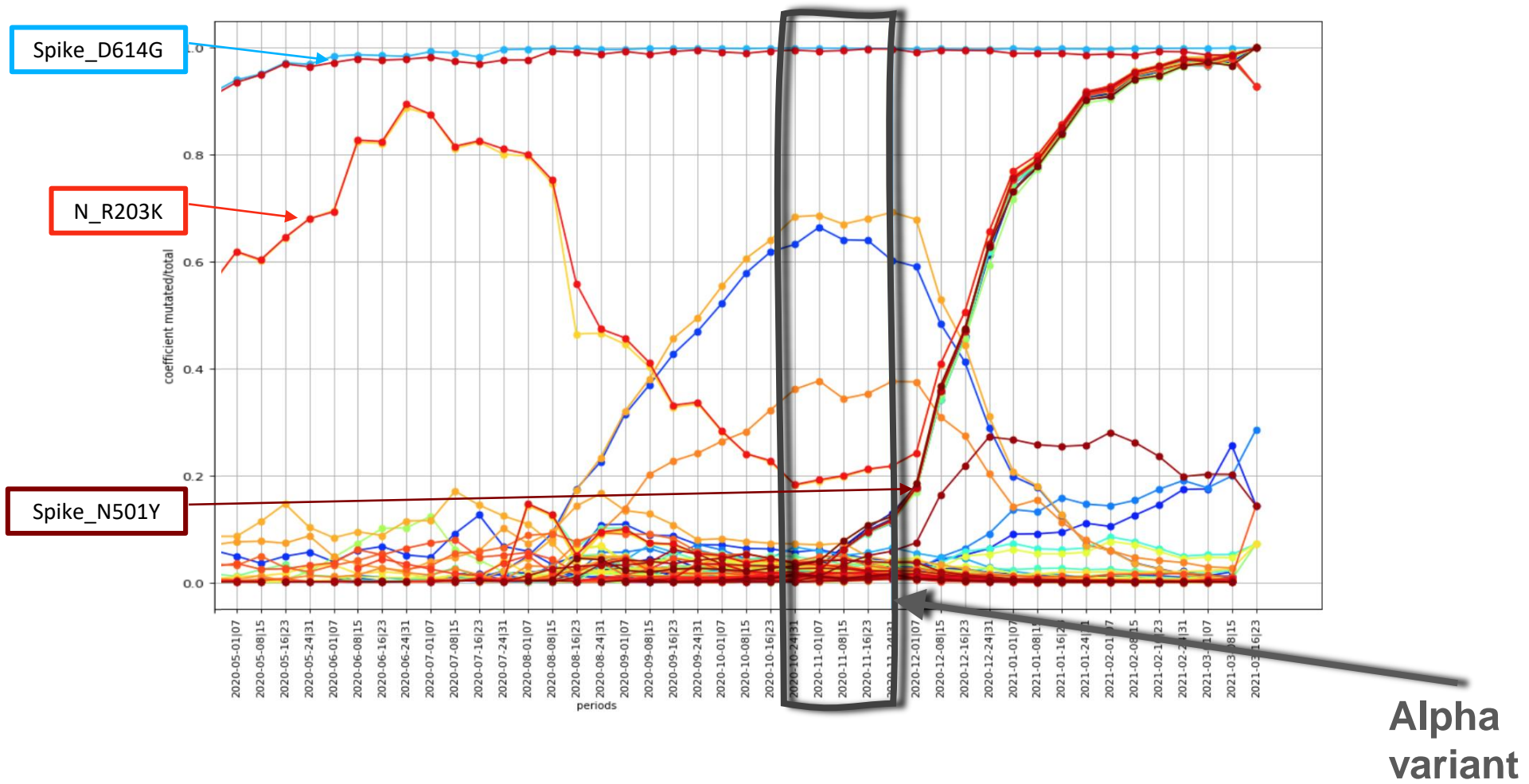
Idea

A variant can be identified by observing the dynamics of its amino acid changes (weekly prevalences in a geo-location)

Different changes could indicate the birth of a variant when:

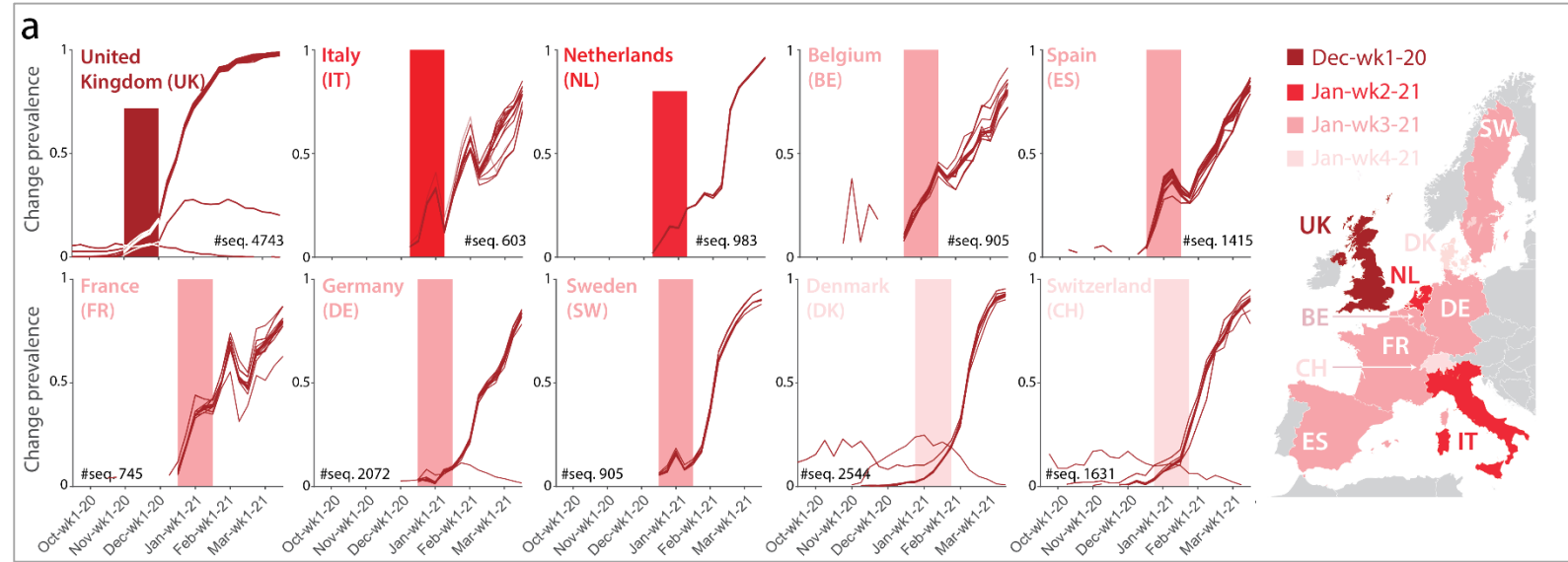
- ❖ their time-series are similar
- ❖ they are all growing

CoEvolution

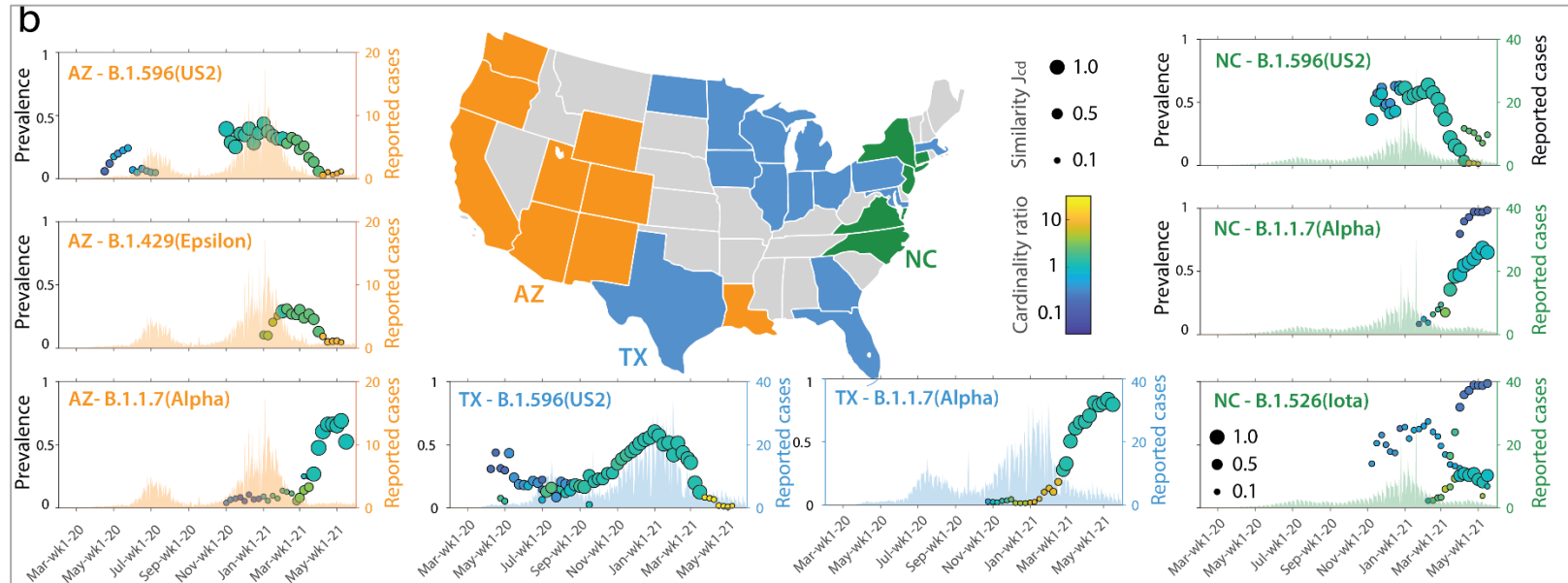


Europe and US variants

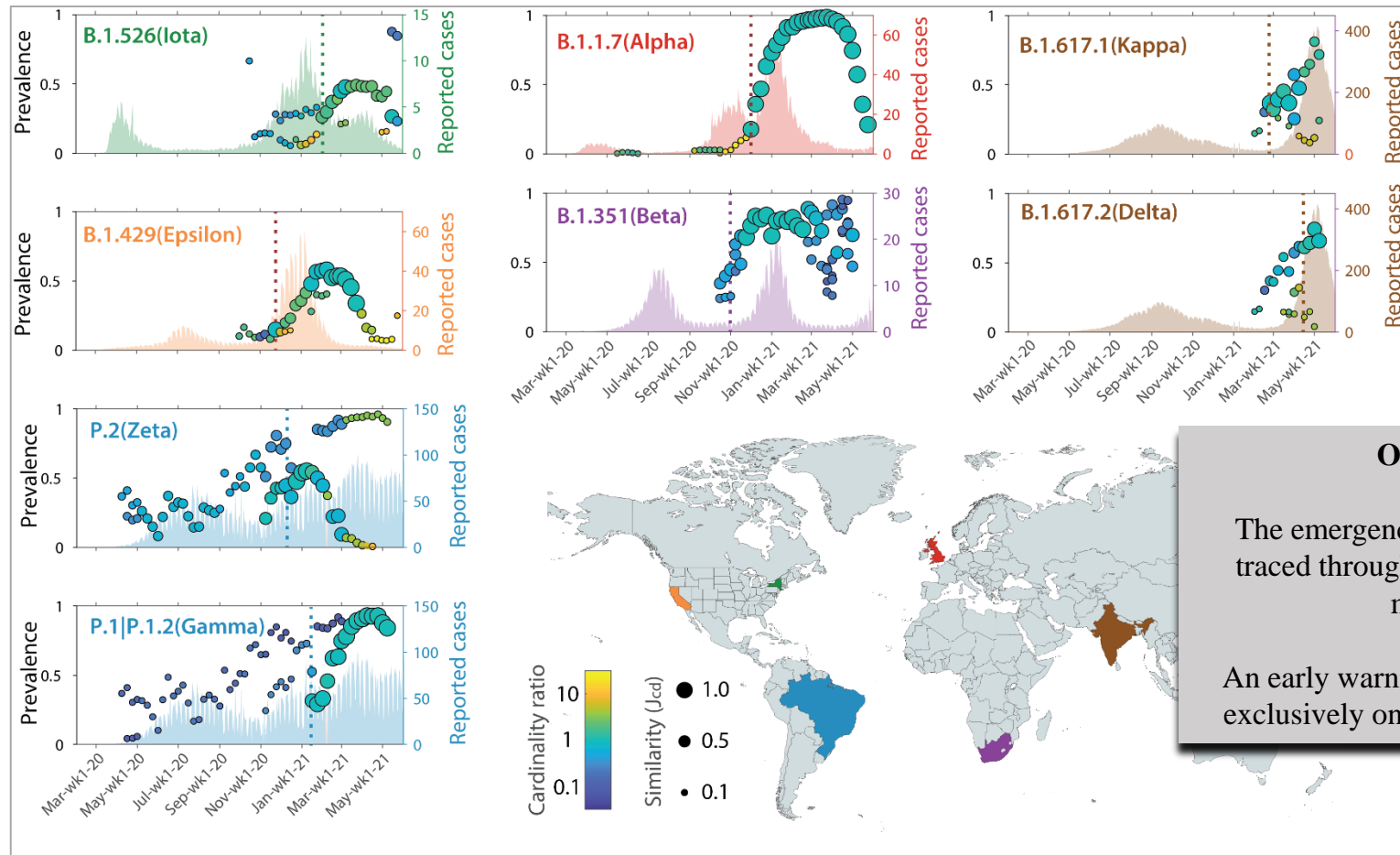
Temporal dynamics of the Alpha variant in European



Temporal dynamics of notable variants in the US



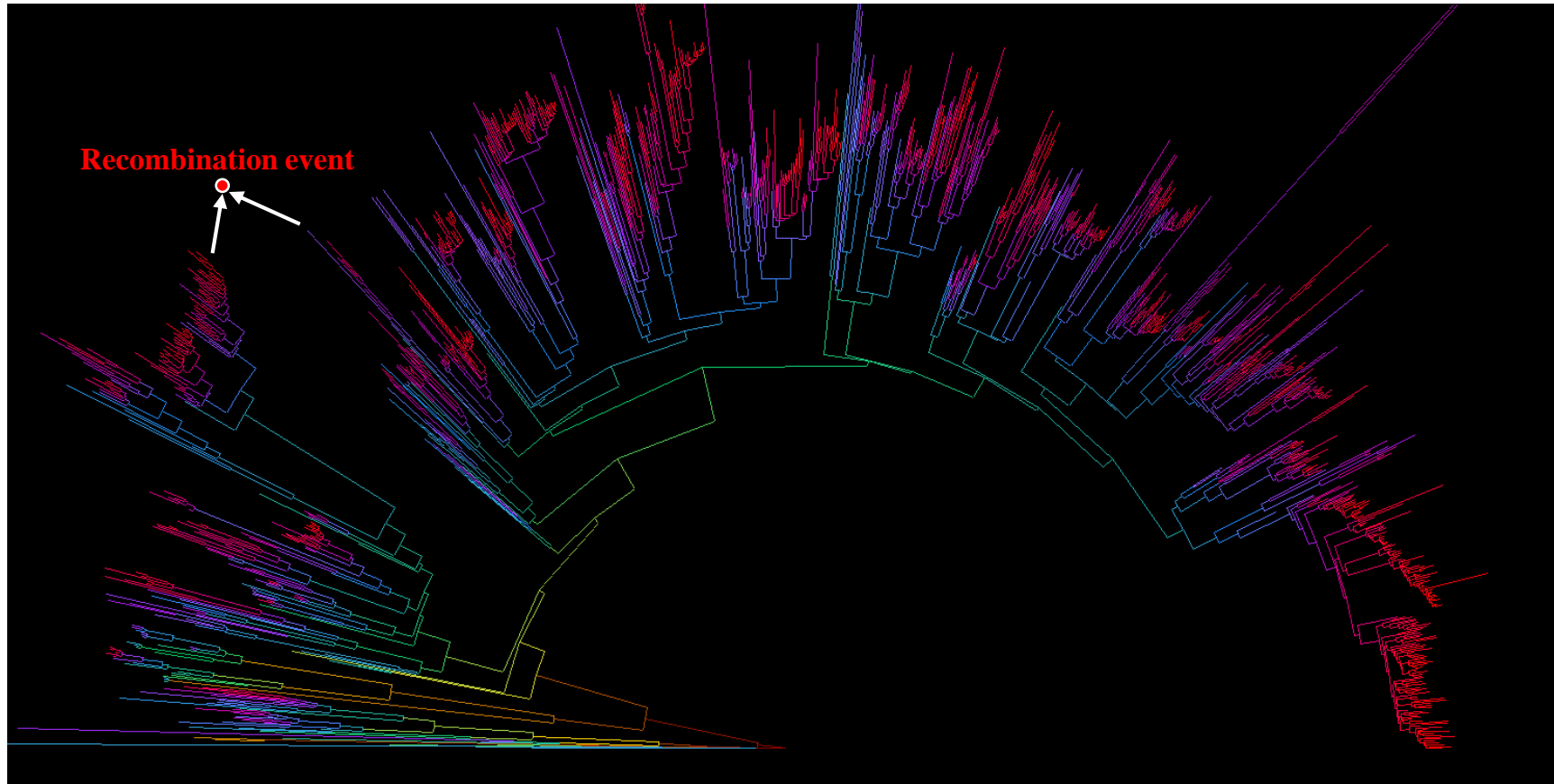
Assessment of the early warning system



Emergence of variants in their country of origin.

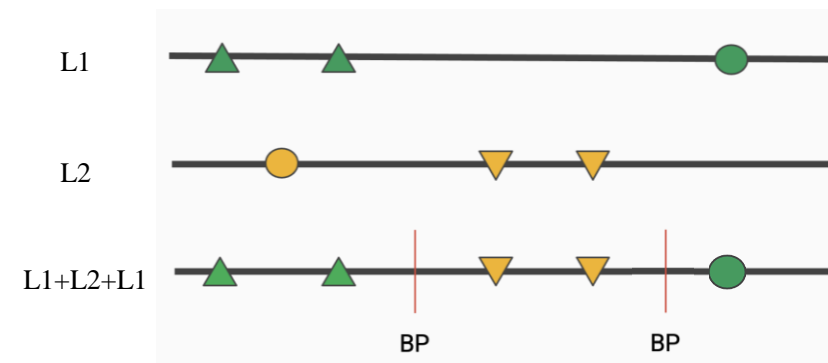
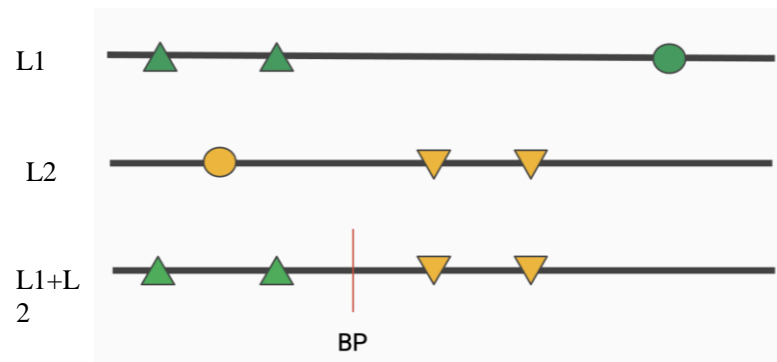
Bernasconi, A., Mari, L., Casagrandi, R., Ceri, S. **Data-driven analysis of amino acid change dynamics timely reveals SARS-CoV-2 variant emergence.** Scientific Reports 11, 21068 (2021). DOI: 10.1038/s41598

Viral recombination

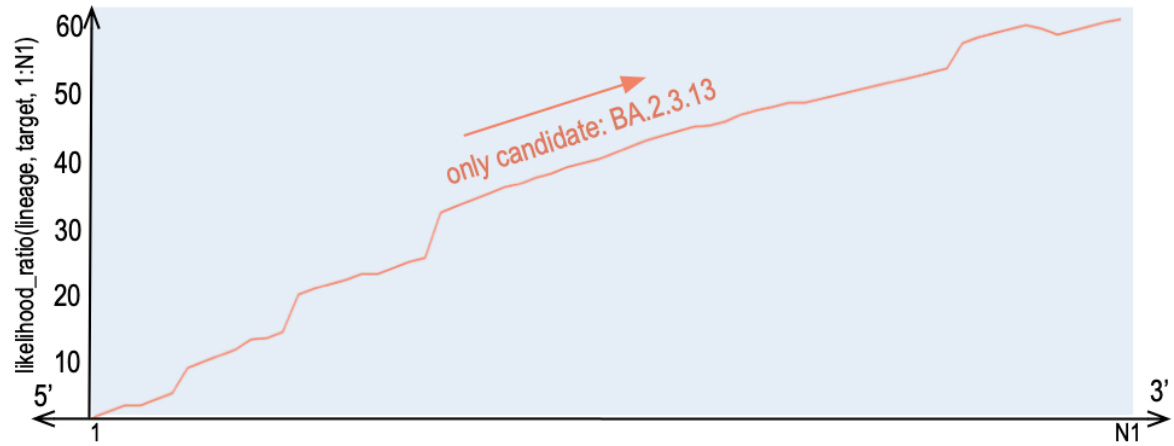


Recombination in SARS-CoV-2

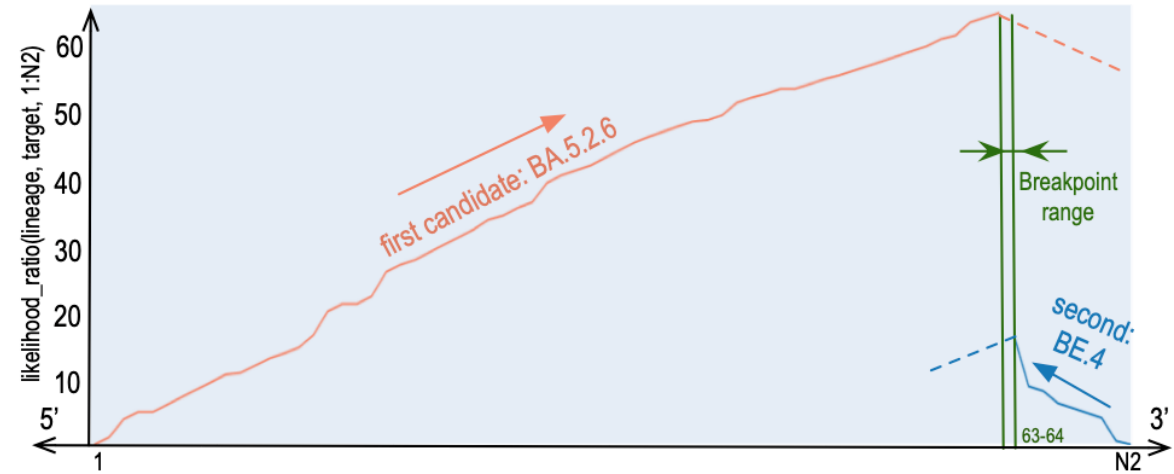
- 76 recombination events recognized in the SARS-CoV-2 virus evolution history (as of July 2023)
- Lineages XBB.1.5 and XBB.1.16 (descendants of the recombinant XBB) are considered **Variants Of Interest** (WHO, August 2023) leading to:
 - <https://www.who.int/activities/tracking-SARS-CoV-2-variants> tion
 - reduced efficacy of treatments
 - predicted increase in transmissibility or disease severity
- Typical recombination patterns:



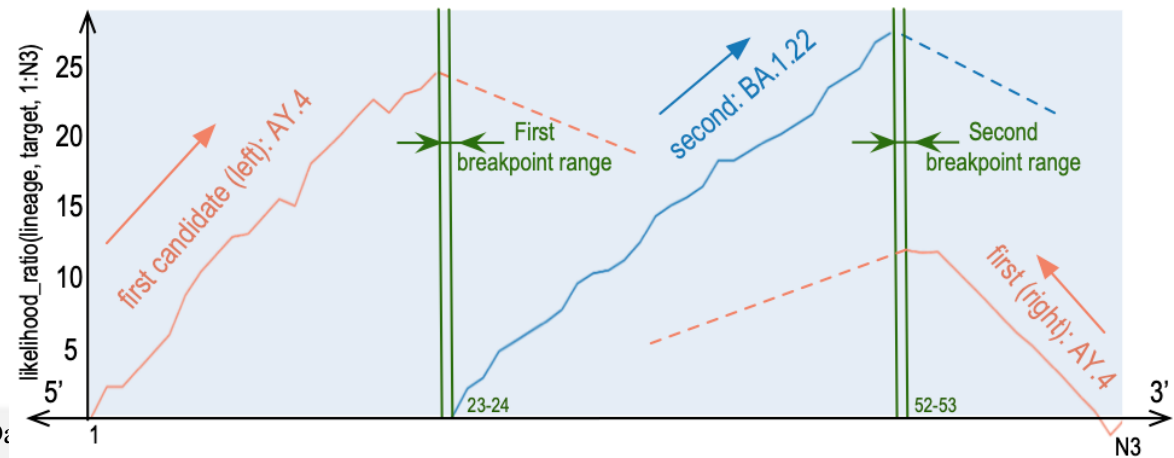
Non recombinant



1BP



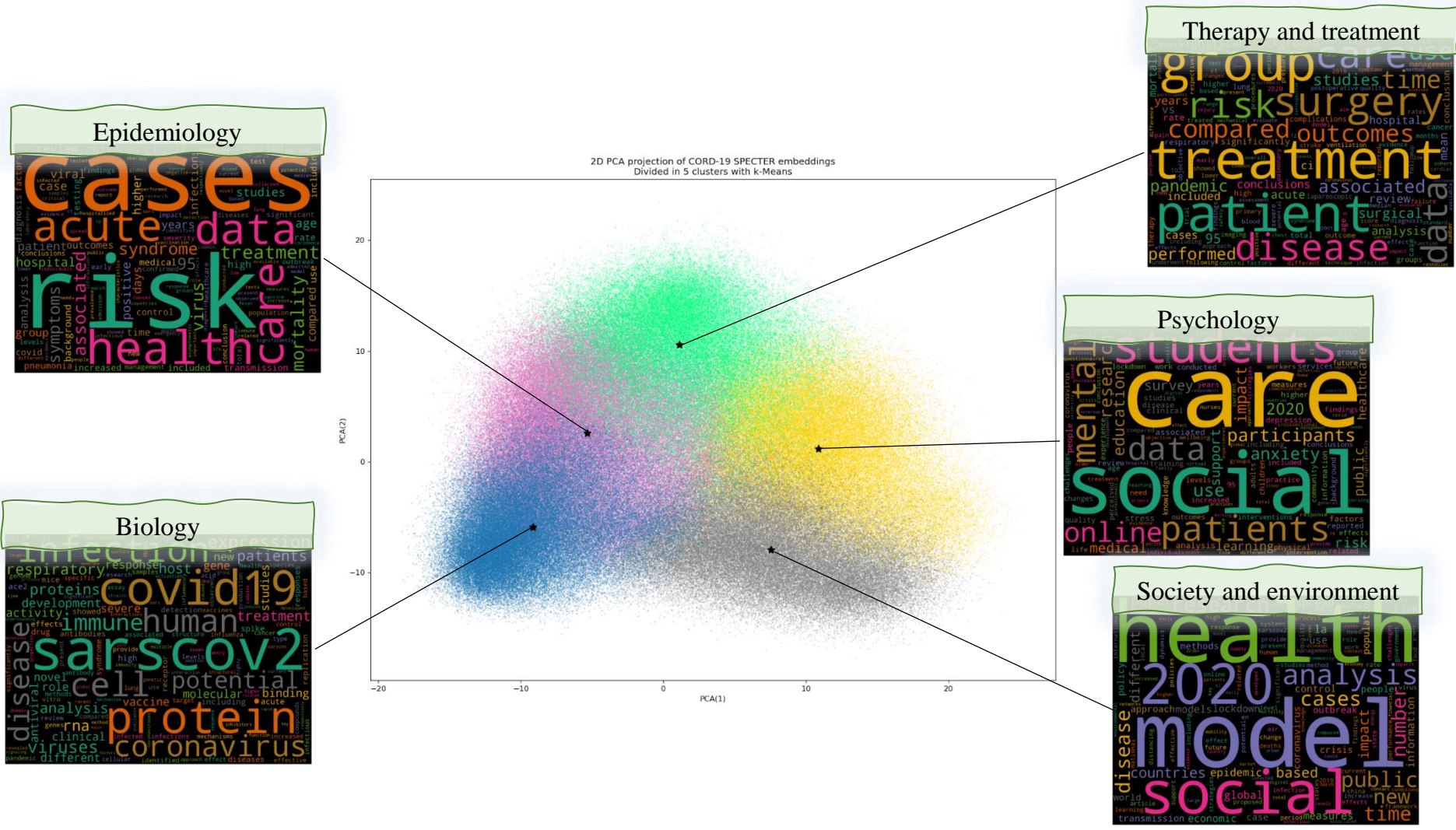
2BP





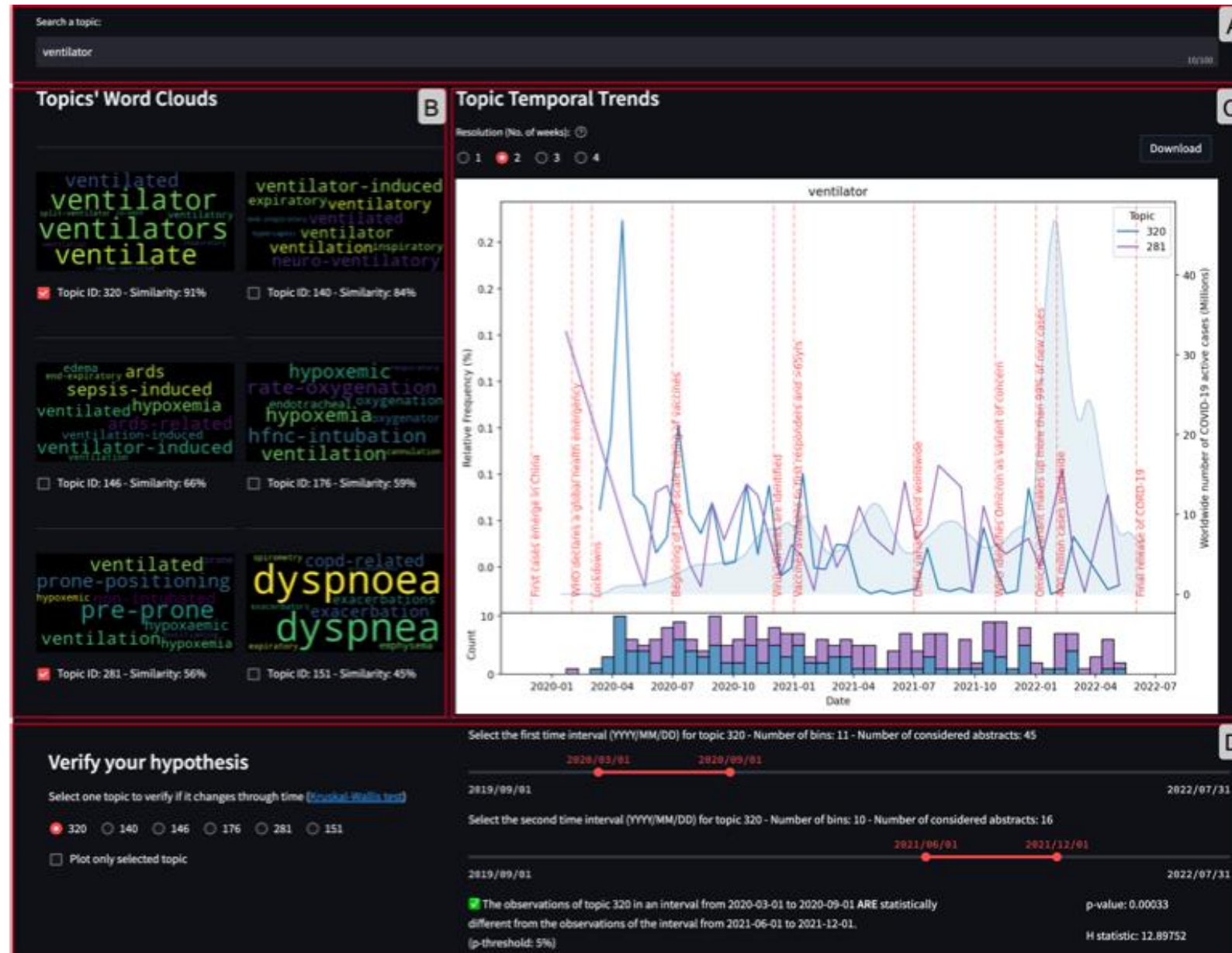
COVID-19 RELOADED

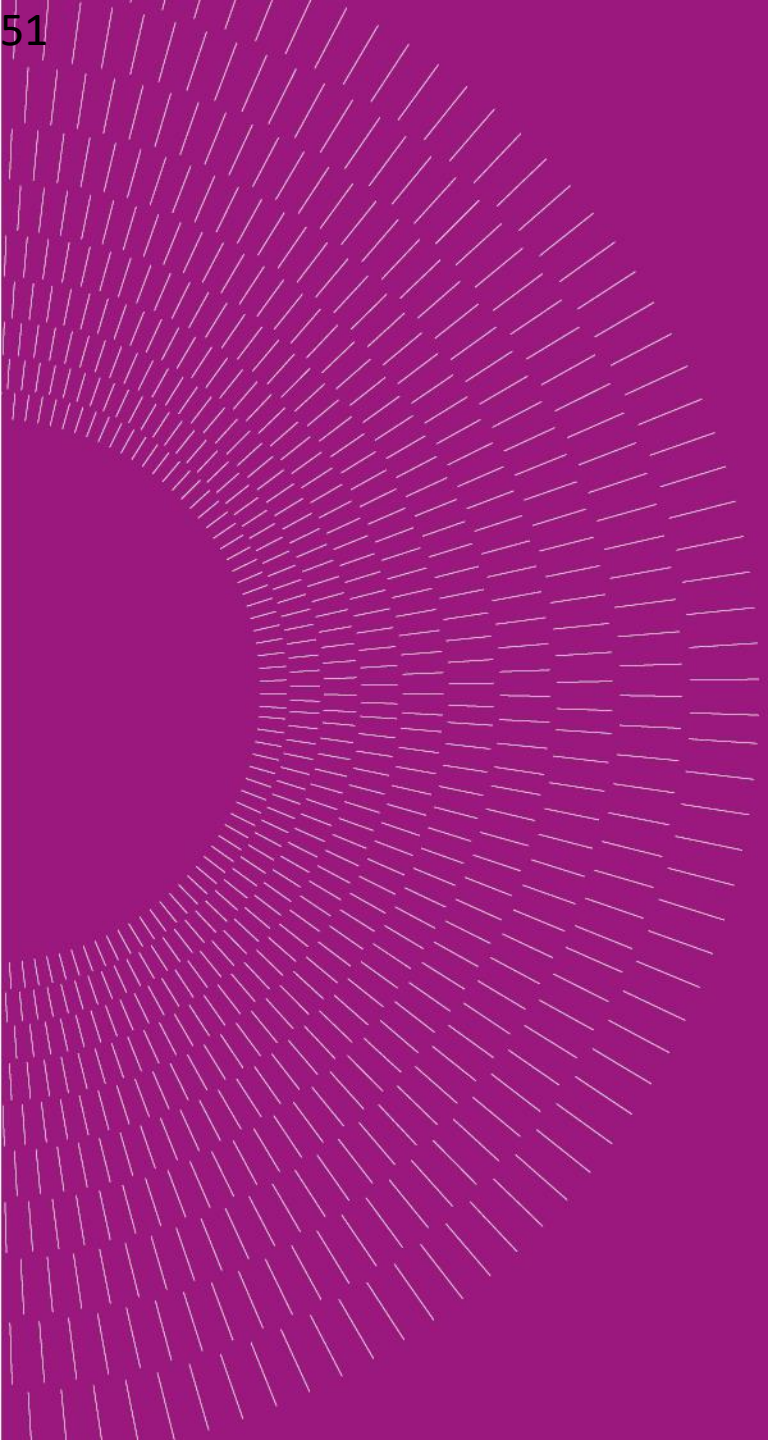
Five topic clusters in COVID-19



Our tool for reading the history of COVID-19: CORToViz dashboard

<http://gmql.eu/cortoviz>





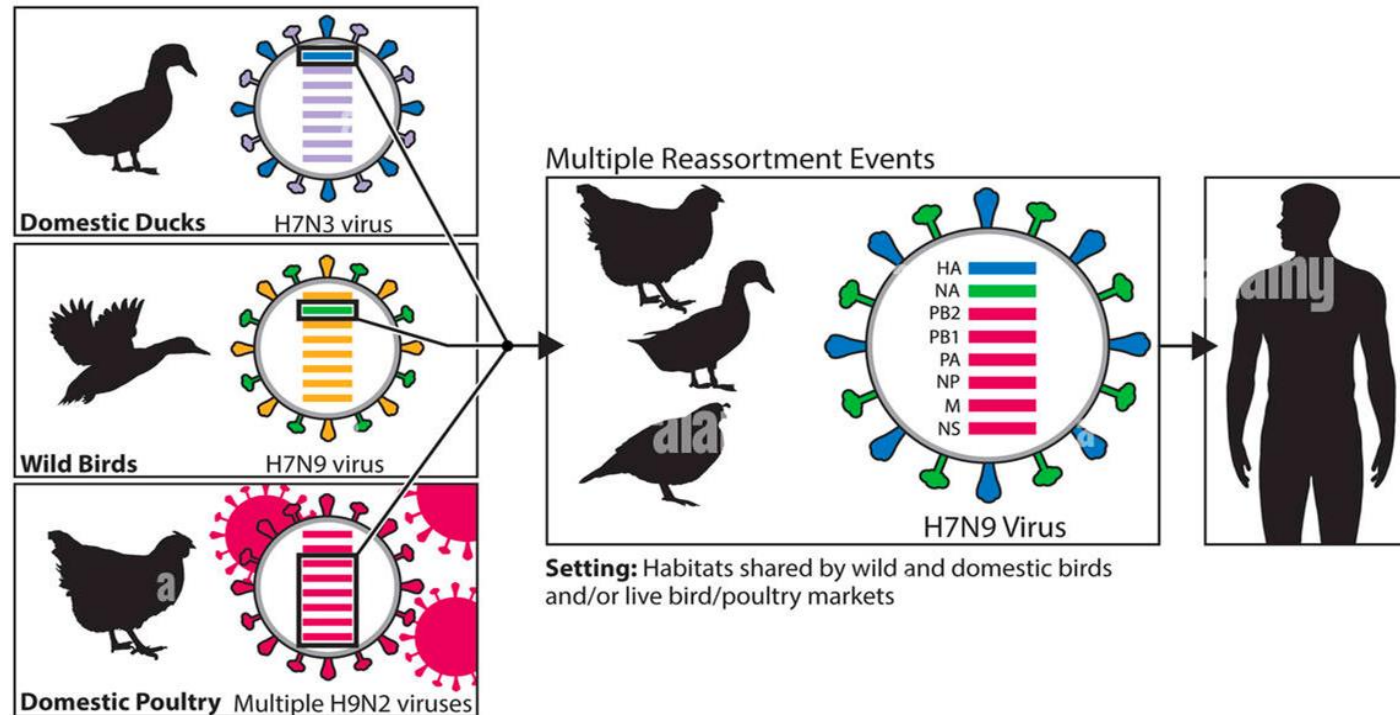
FUTURE RESEARCH HIGHLIGHTS

Recombination & reassortment in the Influenza virus

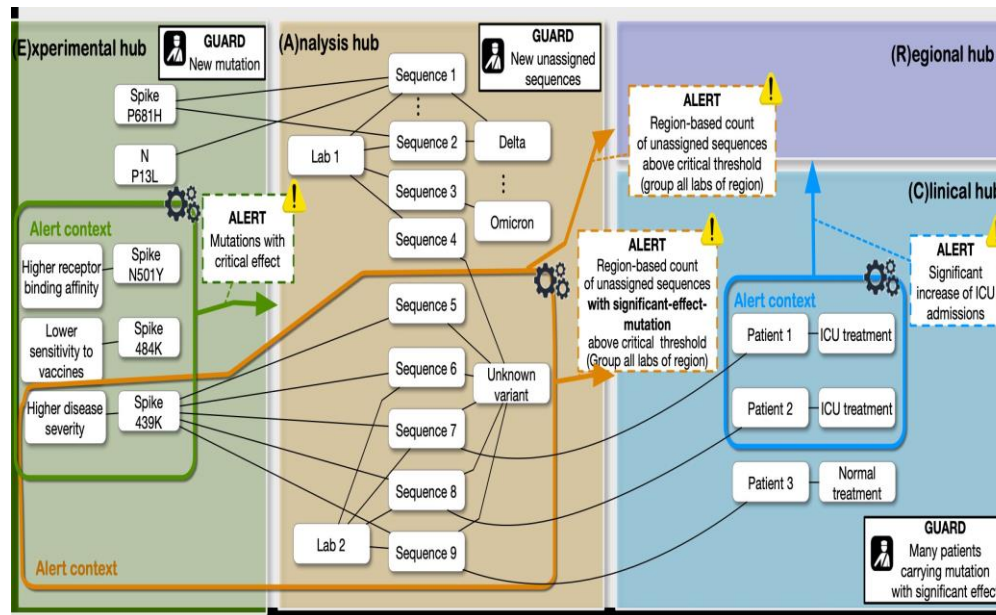
Source of variation in Influenza

- Reassortment
(exchange of genome segments)
- Intra-segment recombination
- Joint work with Ilaria Capua (virologist)

Genetic Evolution of H7N9 Virus in China, 2013



- Reactive Knowledge Management for COVID and Beyond

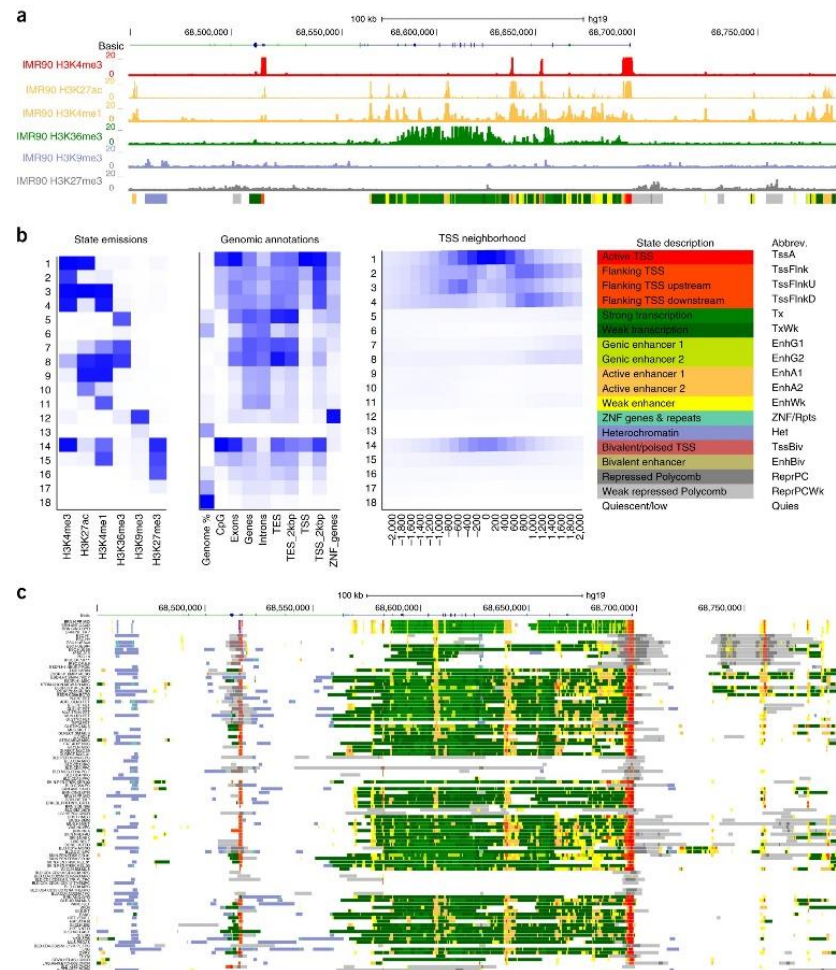


Highlights

- Partitioned knowledge graphs, partitions assigned to knowledge hub representing scientific communities/domains
- Reactive computations across domains make them aware of critical conditions
- Core work in knowledge management technology

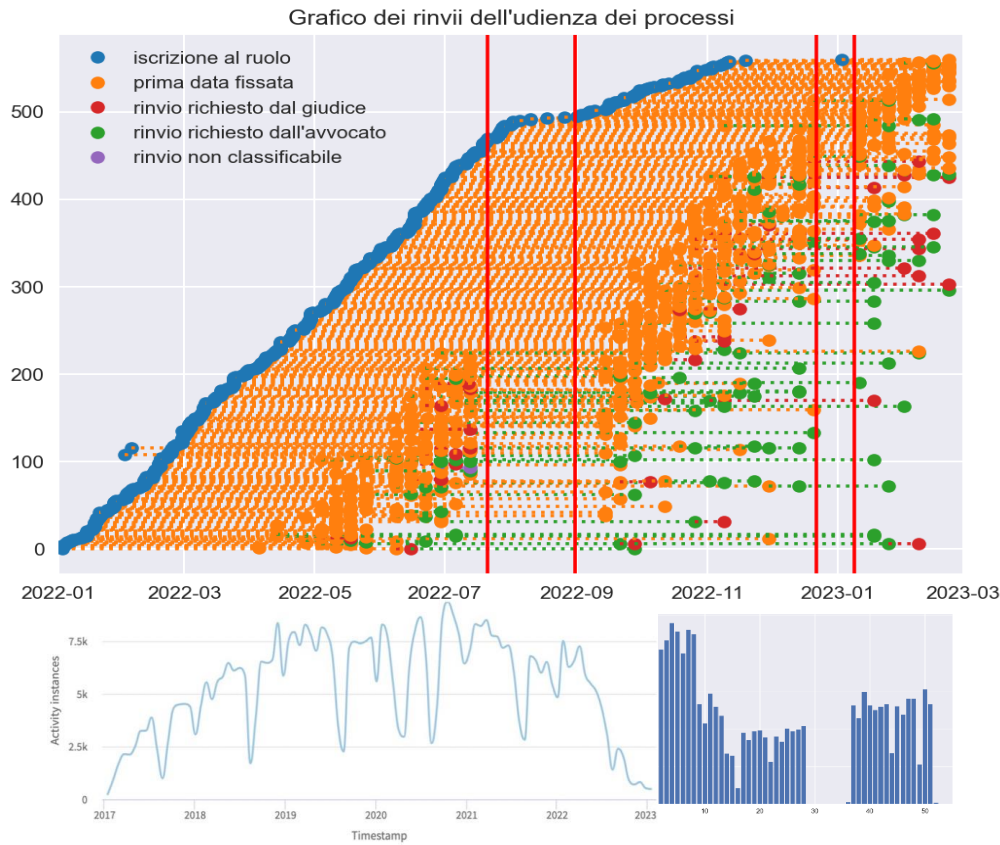
Sound of Genome

Highlights



- Regular structure of the genome (chromosomes, compartments, topological domains, loops)
- Already sonified (starting from the Integrated Genome Browser)
- Will the sound be useful, e.g. for recognizing cancer? Will it be nice to hear (e.g. Xenakis)?
- Project with: G. Haus (Musical informatics), A. Sarti (sound engineering), F. Avanzini, L. Nanni, P. Pinoli (genomics), F. Invernici, A. Bernasconi (Data Scientists) L. Francesconi (composer)

Time of Justice



Highlights

- Within Next Generation UPP, understand the times of justice (here: recent events in “Corte d’Appello”, Milano, courtesy of Barbara Pernici, after extremely complicated process analysis)
- We now have a large data warehouse upon which we will do data mart / data cube / machine learning for duration prediction
- Project with: B. Pernici, A. Campi, and several students
- No relationship to life science (but it is funded by NextGenEU and super-intriguing...)

GeCo – Viral Team

