# The European COVID-19 Data Platform

Nadim Rahman
European Nucleotide Archive (ENA), EMBL-EBI

rahman@ebi.ac.uk

# The question is not *if* the next pandemic will happen but *when*

TIME magazine, May 2017

Recent outbreaks:

- 2002-2004 SARS (Severe Acute Respiratory Syndrome)
- 2009 Swine Flu
- 2011 Germany E. coli O104:H4
- 2013-2016 Western African Ebola virus
- 2015-2016 Zika virus

BY-COVID is developing the infrastructure to deal with the next pathogen.

**Funded under HORIZON-INFRA-2021-EMERGENCY-01:** FAIR and open data sharing in support to European preparedness for COVID-19 and other infectious diseases



TIME

WARNING:
WE ARE NOT READY FOR THE NEXT PANDEMIC

SCIENCE KNOWS HOW TO FIGHT AN OUTBREAK— BUT POLICY STILL GETS IN THE WAY
BY BRYAN WALSH

HOW TO KEEP THE WORLD SAFE
BY BILL GATES

# BY-COVID in a nutshell
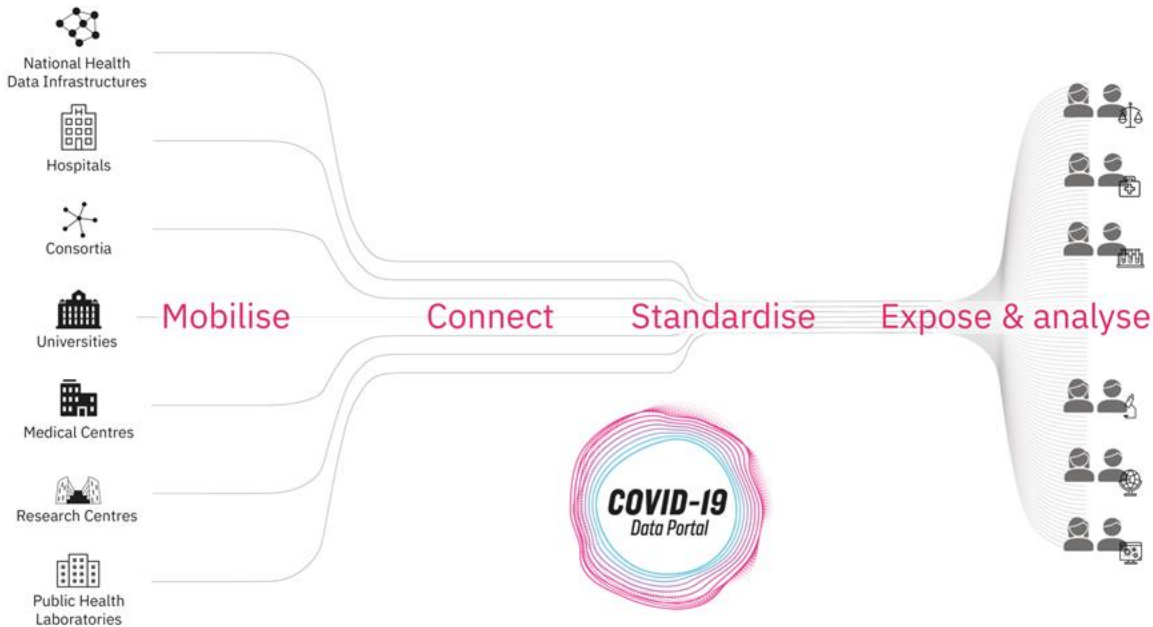
# European COVID-19 Data Platform

*Obj 3: Establish a sustainable and federated infrastructure enabling open sharing of scientific results*

BY-COVID is one of the Horizon Europe projects that supports the operation of the **European COVID-19 Data Platform**

The **COVID-19 Data Portal contains >14 million viral sequences and >1million open scientific articles related to COVID-19.** Since its launch in April 2020, the Portal received 8.8M visits from 500k unique visitors by April 2023.

Raw data from **105 countries and regions, 216 unique centres.**

# Background

- Sit on top of existing infrastructure at EMBL-EBI
- Includes 3 main components:

## COVID-19 Data Portal

Interface for COVID-19 life sciences data

**Publication DOI**: 10.1093/nar/gkab417



## SARS-CoV-2 Data Hubs

Tools to support submission, analysis, visualisation and presentation of COVID-19 sequence data in the COVID-19 Data Portal*



## Federated European Genome-phenome Archive (FEGA)

Controlled-access sharing of human COVID-19 biomolecular and phenotypic data, to present in the COVID-19 Data Portal

# European COVID-19 Data Platform

**Access**

**Data Hubs**

**Human data**



- 25M records across molecular, literature, imaging and social science, backed up by a network of 11 national Data Portals

- Mobilisation (>6M) and systematic analysis (>4M) data sets from SARS-CoV-2 isolates from 121 countries

- 15 VEO public health reports

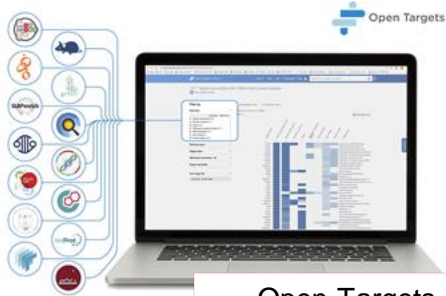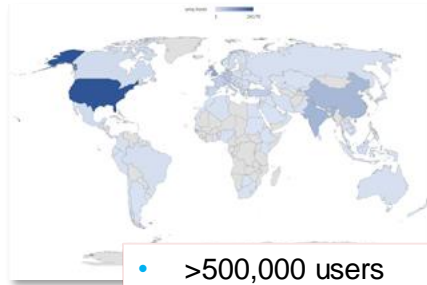- Demonstrated protocols to link between sensitive research subject and pathogen data, leveraging (federated) European infrastructure
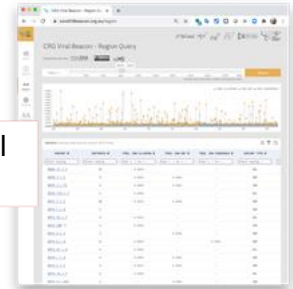
# Usage



Open Targets

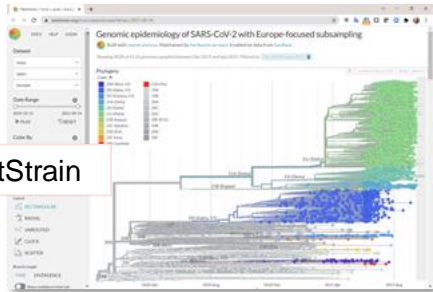- >500,000 users
- 8.8 million requests

Cloud workflows

CRG Viral Beacon

NextStrain

VEO Public Health reporting

National data integration

# Sustainability by design

- Molecular biology components built on existing ELIXIR data resources
  - Resources with institutional commitment and history of sustained activity
  - Typically globally connected through data exchange collaborations
  - Long-term sustainability: Global Biodata Coalition
- Open standards and software
- High level of distribution of expertise and effort

# Infrastructure already shared across projects

**Common indexing framework**



- Common indexing framework, supporting
  - Metadata indexing with three tiers of granularity
  - Domain-level classification system to define data "partitions", e.g. "blue domain"
  - Coverage spans ELIXIR Deposition and Core Data Resources and broadly connected via FAIRSharing
  - Coverage goes beyond ELIXIR including for image and social science data with further data resources being added
- Continued development work under both BY-COVID and Blue-Cloud26
- Supports metadata feeds to future initiatives in both domains

# A broader Pathogens Platform

- Scope
  - All pathogen, all disease approach
  - Hosts, vectors and pathogens
  - Antimicrobial resistance
- Preparedness and Outbreaks
  - Pathogens Portal
  - Pathogen Data Hubs - responding to outbreaks as they arise
- Applicability
  - Climate change → disease outbreaks
  - Cancer associated viral/bacterial infection
- Future roles
  - Food security, e.g. plant pathogens
  - Biodiversity loss

https://www.pathogensportal.org/

# EarlyCause - early life stress

- Data portal, search functionalities
- Biodata to support investigations into lifelong effects of early-life stress
- https://portal.earlycause.eu/
- Data on organism level or data type (Mouse, Rat, Human, Cell Lines, Literature, Cohorts)
- Reusable infrastructure and framework to bring forward biodata (e.g. soil biodata)

# Thanks!

**EMBL-EBI**

Ahmad Zyoud
Alexey Sokolov
Amonida Zadissa
Andrew Parton
Andrii Ludin
Andy Yates
Carla Cummins
Claire O'Donovan
Claire Rye
Colman O'Cathail
Craig Russell
Dipayan Gupta
Dylan Spalding
Eloise Stapleton
Gabi Rinck
Galabina Yordanova
Geetika Malhotra
Giselle Kerry
Guy Cochrane
Helen Parkinson
Henning Hermjakob
Jeena Rajan
Jeff Knaggs
Joseph Rosetto

Josephine Burgin
Karoly Erdos
Laura Harris
Mallory Freeberg
Manish Kumar
Marianna Ventouratou
Matt Pearce
Melanie Courtot
Mihai Glont
Milena Mansurova
Nadim Rahman
Nicola Buso
Oana Stroe
Ossama Edbali
Pablo Moreno
Peter Harrison
Peter Walter
Raheela Aslam
Rasko Leinonen
Rodica Petrusevschi
Rodrigo Lopez
Rolf Apweiler
Ross Thorne
Sam Holt
Sandeep Kadam
Sandeep Selvakumar

Sarah Hunt
Senthil Vijavaraja
Simon Kay
Stefan Gutnick Allen
Suran Jayathilaka
Thomas Keane
Timothee Cezard
Tony Burdett
Tracey Mahoney
Vishnu Kadhirvelu
Youngmi Park
Zahra Waheed
Zamin Iqbal

**Independent**
Robert Petryszak

**ELIXIR**
Katharina Lauer
Niklas Blomberg

**DTU**
Jose Luis Bellod Cisneros
Martin Christen Frølund Thomsen
Johanne Ahrenfeldt
Rolf Sommer Kaas
Lukasz Dariusz Dynowski
Frank Aarestrup
Jeffrey Skiby
Judit Szarvas
Camilla Hundahl Johnsen
Rene S. Hendriksen
Martin Koliba
Philip Clausen
Máté Gulyás

**ELTE**
János Márk Szalai-Gindl
Balint Pataki
Jozsef Steger
Dávid Visontai
Krisztian Papp
Istvan Csabai
Ágnes Becsei
Ákos Gellért
Anikó Mentes
Orsolva Pipek

**EMC**
Marion Koopmans
Clara Amid
David van de Vijver
Mariolein Poen
Miranda de Graaf
Maarten Hoek
David Nieuwenhuijse
Divyae Prasad
Marie-Catherine
Bouquieaux

**FLI**
Dirk Hoeper
Ariane Belka
Maria Jenckel
Claudia Wylezich
Martin Beer
Anne Pohlmann

**RIVM**
Dennis Schmitz
Florian Zwagemaker
Annelies Kroneman

# Thanks for listening!
# Any questions?

**Nadim Rahman**
rahman@ebi.ac.uk