

On the EGA* Ecosystem

(*European Genome-phenome Archive)

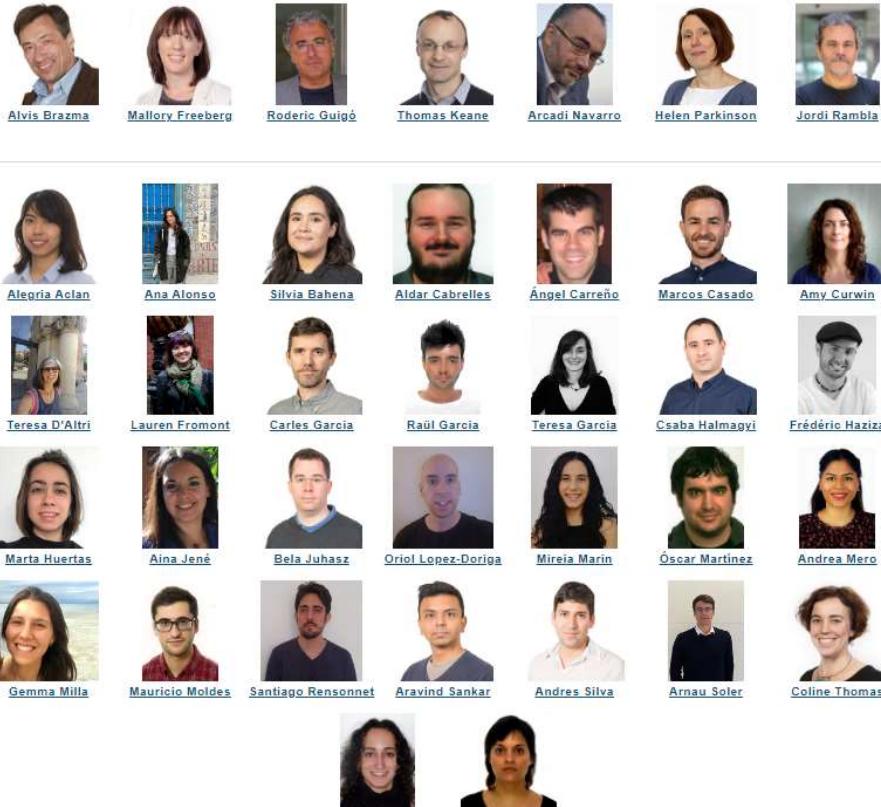
Arcadi Navarro

November the 13th, 2023



Thanks!

The Team



Core organizations:



Additional sources:



"Una manera de hacer Europa"



And critical support from the following sources:





Outline

1. Big Data & Genomic Big Data
2. Genomic Data
3. The sharing dilemma
4. What is the EGA for?
5. Beacon
6. The EGA Federation



Outline

1. Big Data & Genomic Big Data
2. Genomic Data
3. The sharing dilemma
4. What is the EGA for?
5. Beacon
6. The EGA Federation



what do all of these have in common?



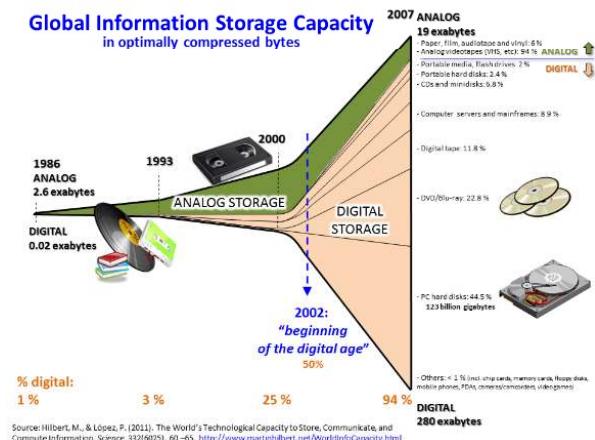
what is “BigData”?

Think for a sec... what does “Big Data” mean?

There are huge volumes of data in the world:

- From the beginning of recorded time until 2003, humans created 5 billion gigabytes of data.
- In 2011, the same amount was created every two days
- In 2013, every 10 minutes.
- The best estimation for 2020 was that 1.7MB of data was created **every second by every person on earth***. This is 13PB per second or >1 exabyte per minute.
- By 2025, there will be 175 zettabytes of data in the global datasphere.

Byte (B)	8 bits	(10^3)
Kilobyte (KB)	1,024 bytes	(10^6)
Megabyte (MB)	1,024 KB	(10^9)
Gigabyte (GB)	1,024 MB	(10^{12})
Terabyte (TB)	1,024 GB	(10^{15})
Petabyte (PB)	1,024 TB	(10^{18})
Exabyte (EB)	1,024 PB	(10^{21})
Zettabyte (ZB)	1,024 EB	(10^{24})
Yottabyte (YB)	1,024 ZB	(10^{27})



*<https://www.domo.com/solution/data-never-sleeps-6>

what is “BigData”?

Think for a sec... what does “Big Data” mean?

This is growing so much that 90% of the data in the world today has been created in less than the last year. This data comes from everywhere: sensors used to gather climate information, posts to social media sites, digital pictures and videos, medical records, purchase transaction records or cell phone GPS signals to name a few.

This data is **Big Data**

what is “BigData”?

...and a tiny fraction of these data belongs to the
fields of Medical Genomics

...but this tiny fraction is still a huge
amount of data!

where does genomic data come from?

The promise in the 90s: “In 10 years we will unravel the genetic bases of complex diseases!!” ... **It was the reason under the Human Genome Project**



20TH Anniversary of The Human Genome Project (2021)

Science aspects: High impact on science and healthcare

Data sharing aspects: Bermuda & Fort Lauderdale Principles > or “Good vs Evil”

Technological advances that these projects had generated:

- Genomic Browsers (and tracks)

- Haplotype viewers

- VCF format

- ...

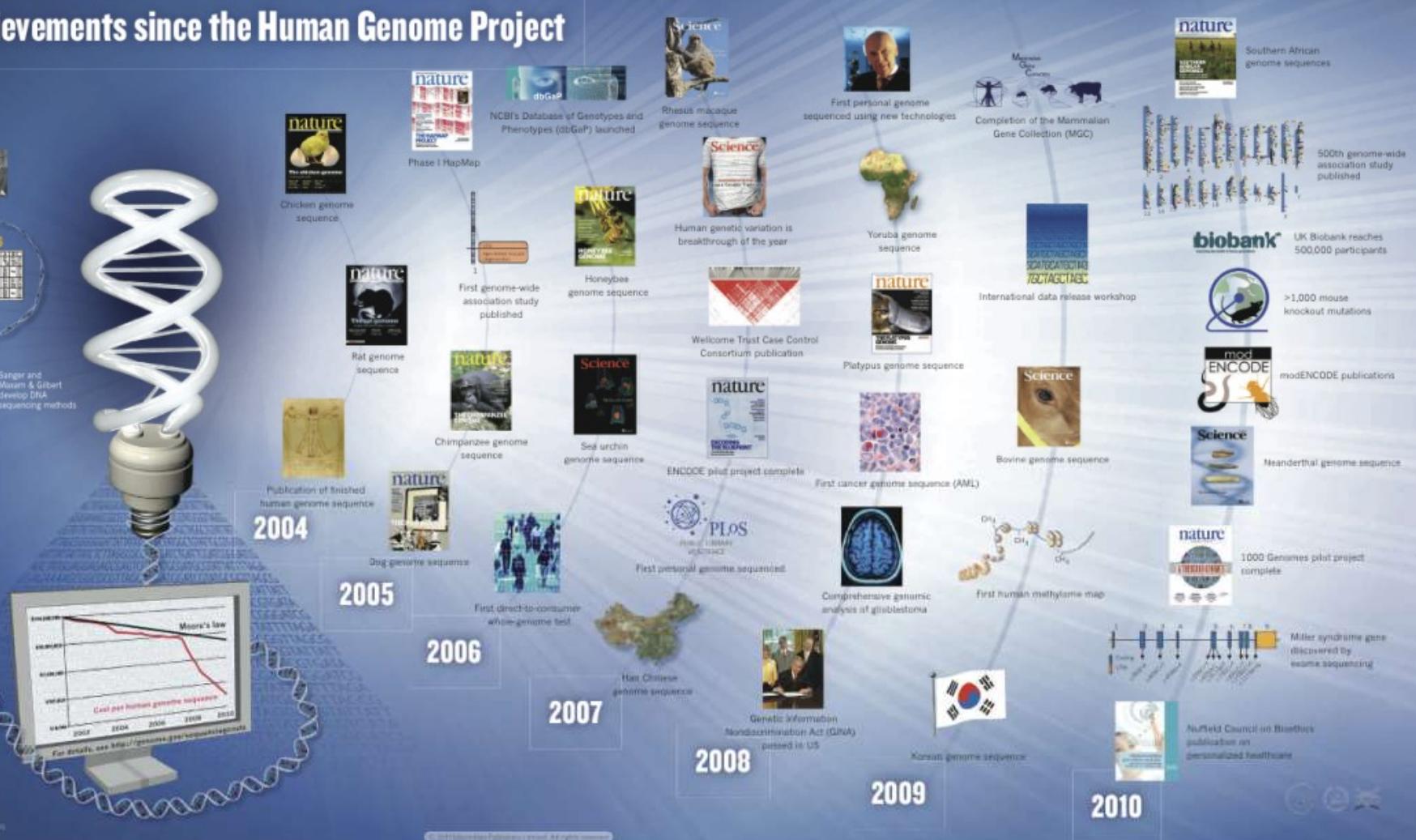
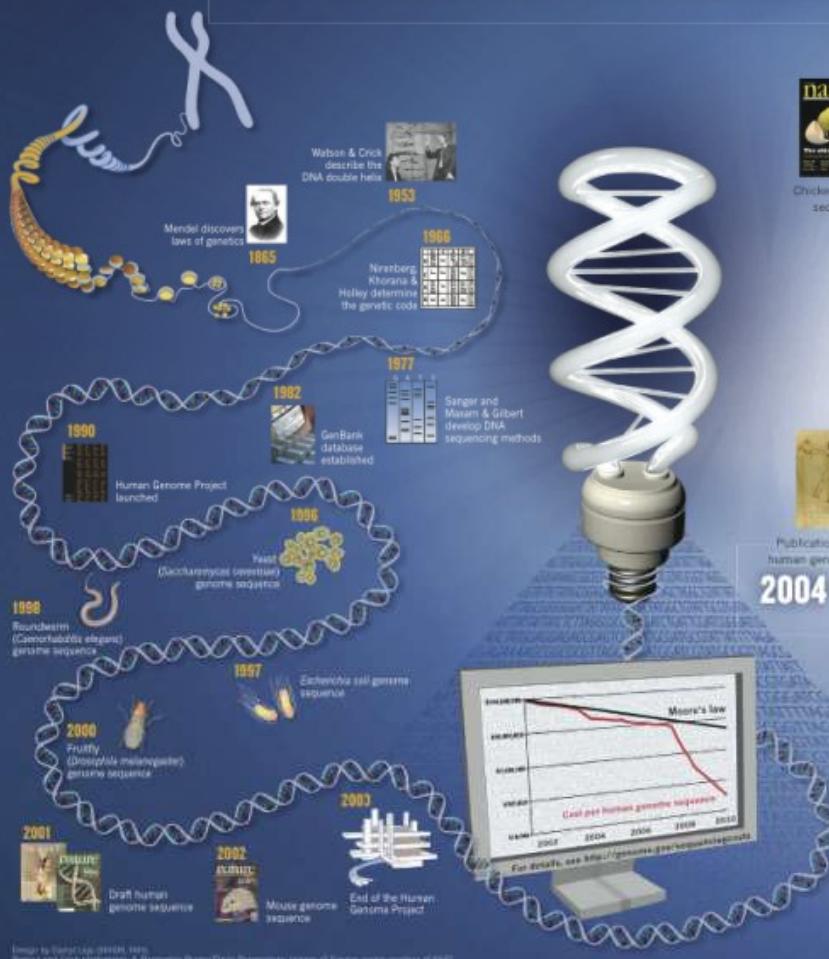
Issues:

- Lack of diversity

- Going backwards on data sharing?



Genomic achievements since the Human Genome Project



Outline

1. Big Data & Genomic Big Data
2. Genomic Data
3. The sharing dilemma
4. What is the EGA for?
5. Beacon
6. The EGA Federation



The enablers: The HapMap project



The International HapMap Consortium
A haplotype map of the human genome
Nature **437**, 1299–1320 (2005)
<https://doi.org/10.1038/nature04226>

The HapMap is a map of haplotype blocks and the specific SNPs that identify the haplotypes (tag SNPs).

The HapMap is valuable by reducing the number of SNPs required to examine the entire genome for association with a phenotype from the 10 million SNPs that exist to roughly 500,000 tag SNPs.

269 DNA samples from four “populations”

The project follows the data release principles of sharing information rapidly and without restriction on its use.

- Officially started October 27 to 29, 2002
- The complete data Phase I published 27 October 2005.
- Phase II dataset published October 2007.
- Phase III dataset released spring 2009 published in September 2010.

Data available at [Index of /hapmap \(nih.gov\)](#)

The enablers: 1000 Genomes Project

<https://www.internationalgenome.org>



Launched in January 2008. Since 2015 managed by the IGSR: *The International Genome Sample Resource*

Has had several phases: Pilot ([2010](#)), Phase I ([2012](#)): 1,092 samples, Phase II (technical), Phase III ([2015](#))

Uses a combination of low-coverage whole-genome sequencing, deep exome sequencing, and dense microarray genotyping

In early 2020, released high-coverage (30x) data for an additional 698 samples.

The samples are anonymous and have **no associated medical or phenotype data**. The project holds self-reported ethnicity and gender. All participants declared themselves to be healthy at the time the samples were collected.

Data was quickly made available to the scientific community through freely accessible public databases

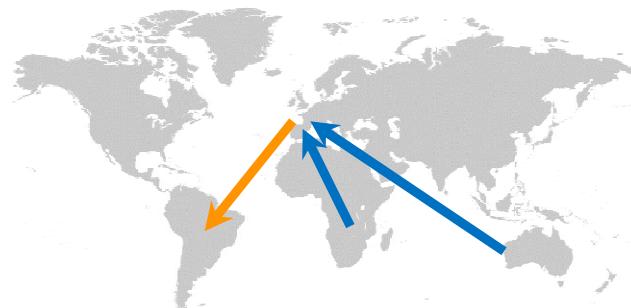
Files accessible from FTP sites hosted by EMBL-EBI and NCBI.

Now also hosted on Amazon Web Services (AWS) and AnVIL.

Thousands of papers have been published that use these datasets

Data access models

Fully Open research data



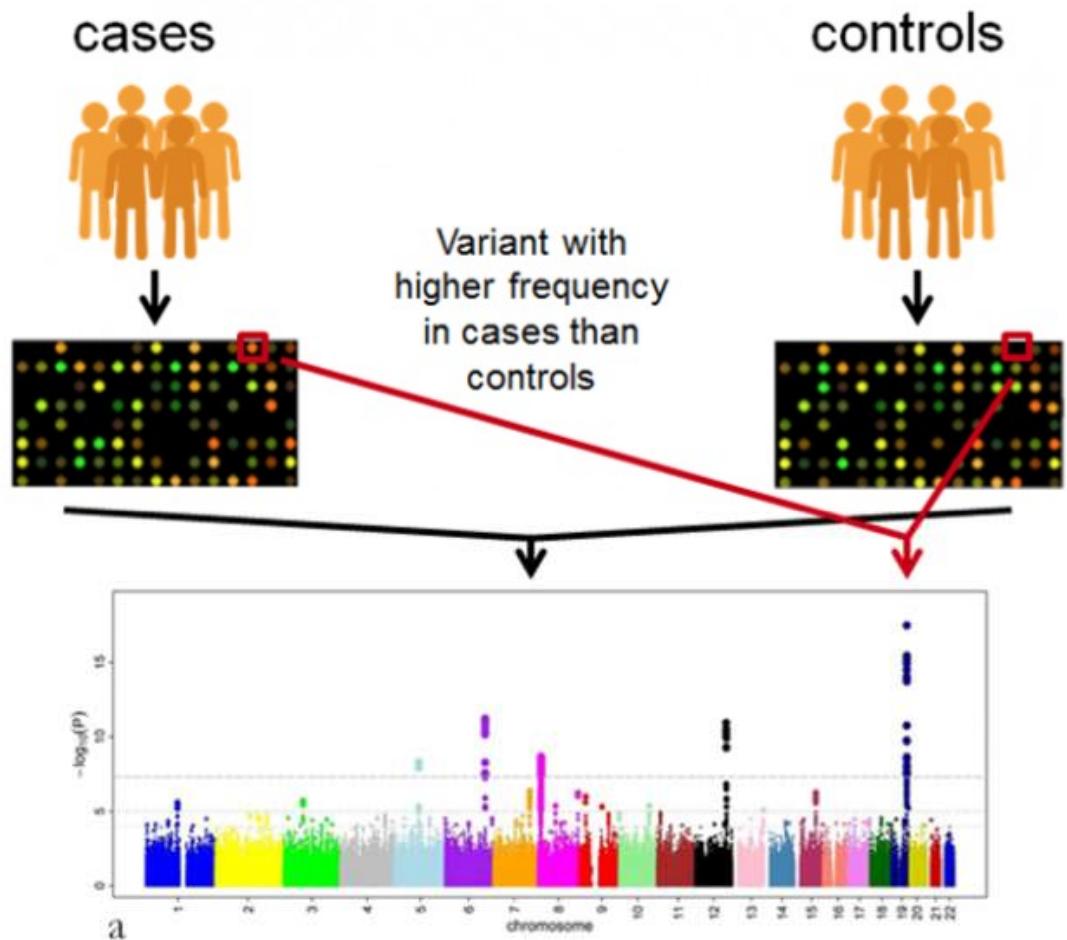
Aggregate data globally

Download, analyse locally

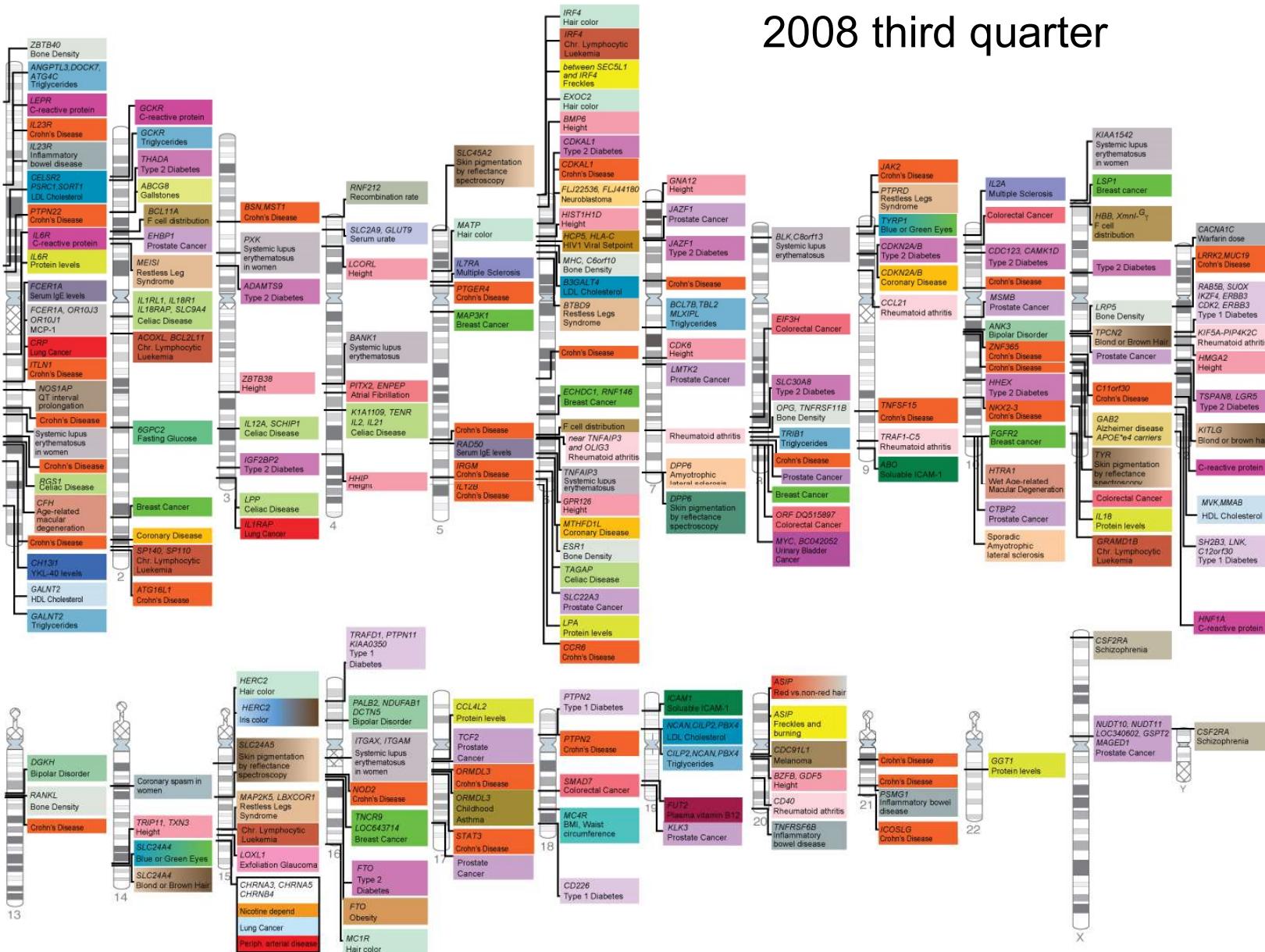
e.g., Plant and animal genomes, Population-based results, Some GWAS summary statistics

Accumulation of human data

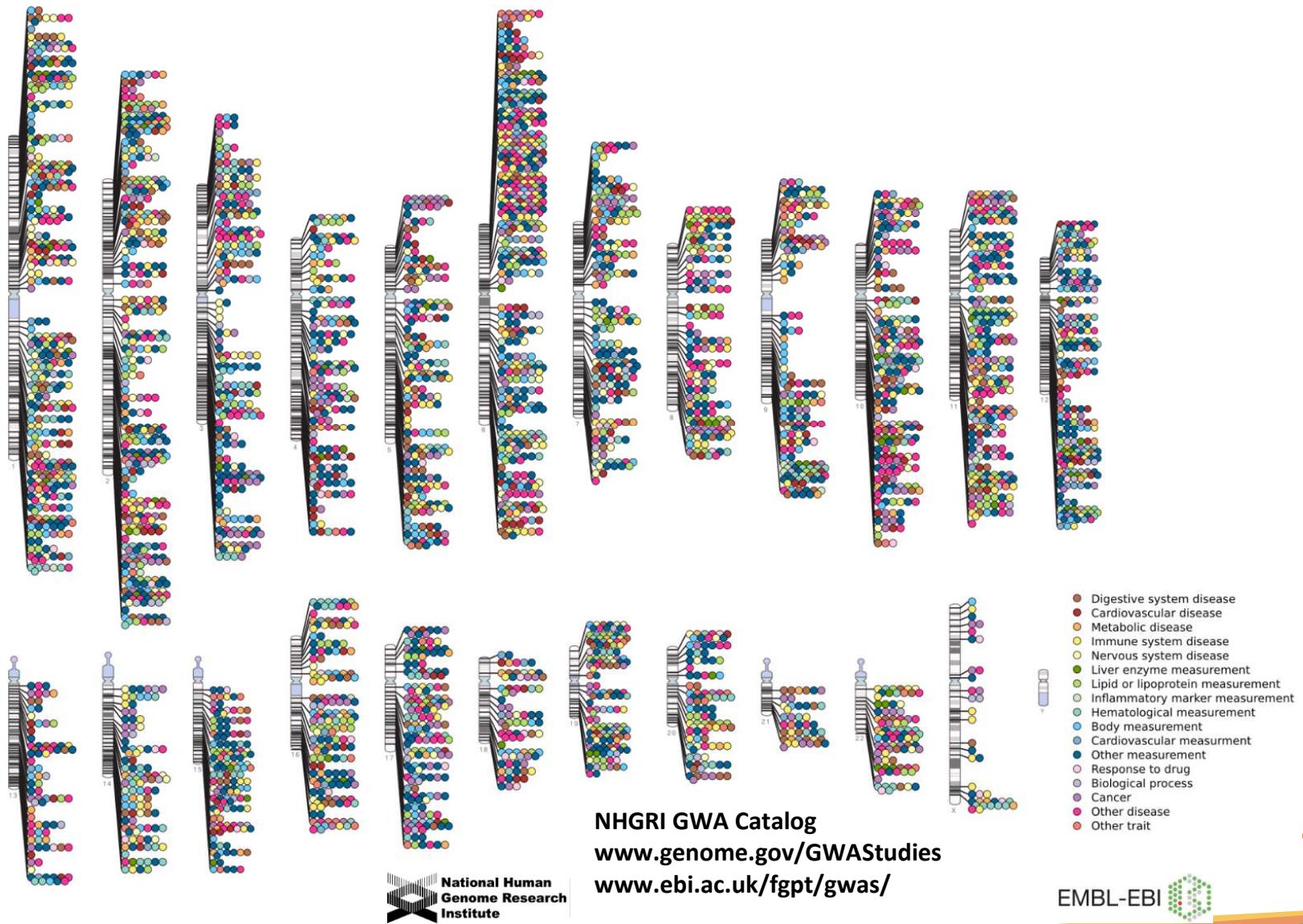
Using all these advances, many data relevant for human health has accumulated. For instance, take Genome Wide Association Studies (GWAS), which started in 2005-2007



2008 third quarter



Published Genome-Wide Associations through 12/2013
Published GWA at $p \leq 5 \times 10^{-8}$ for 17 trait categories



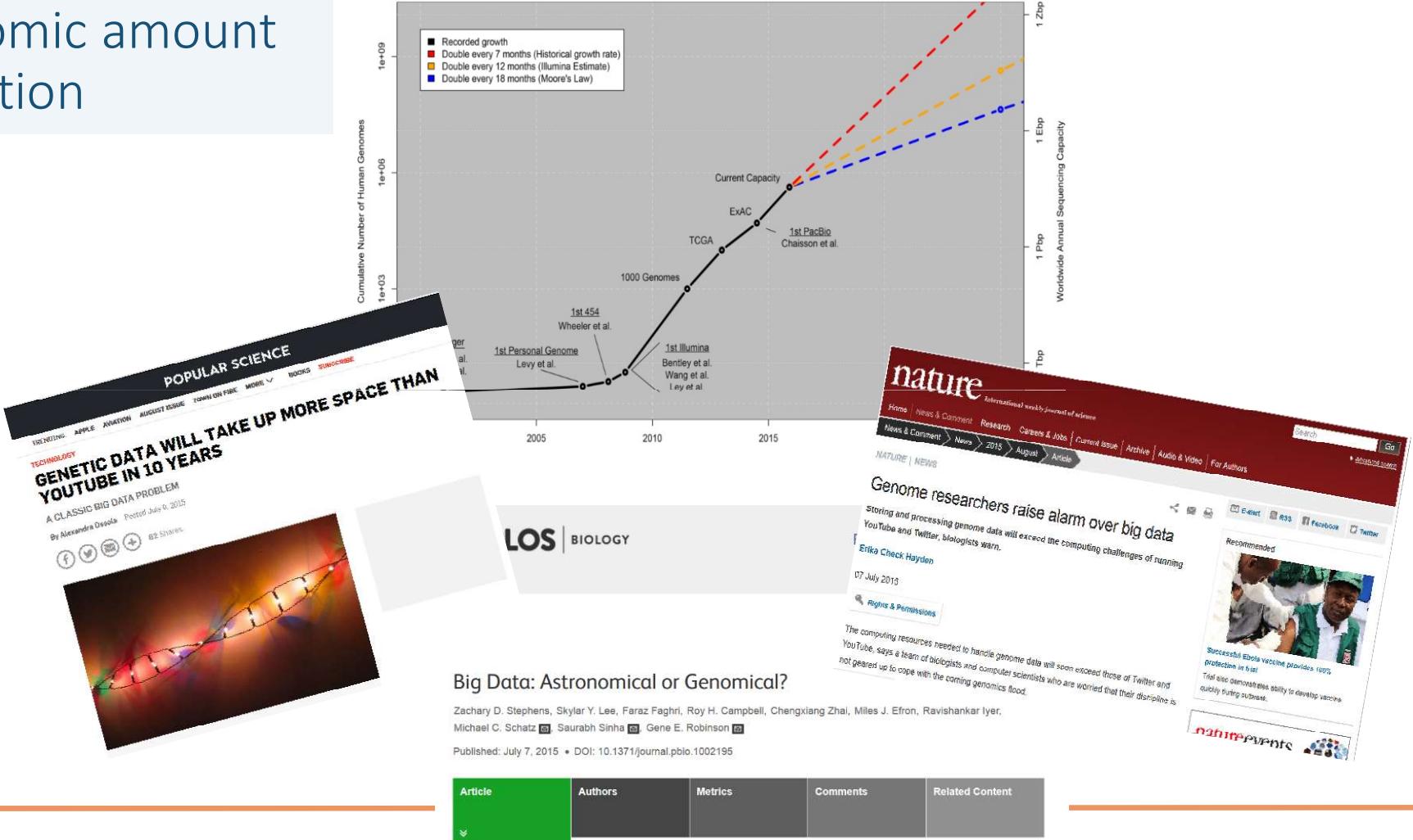
Yesterday's version





-

An astronomic amount of information



Amazing results!

nature reviews drug discovery

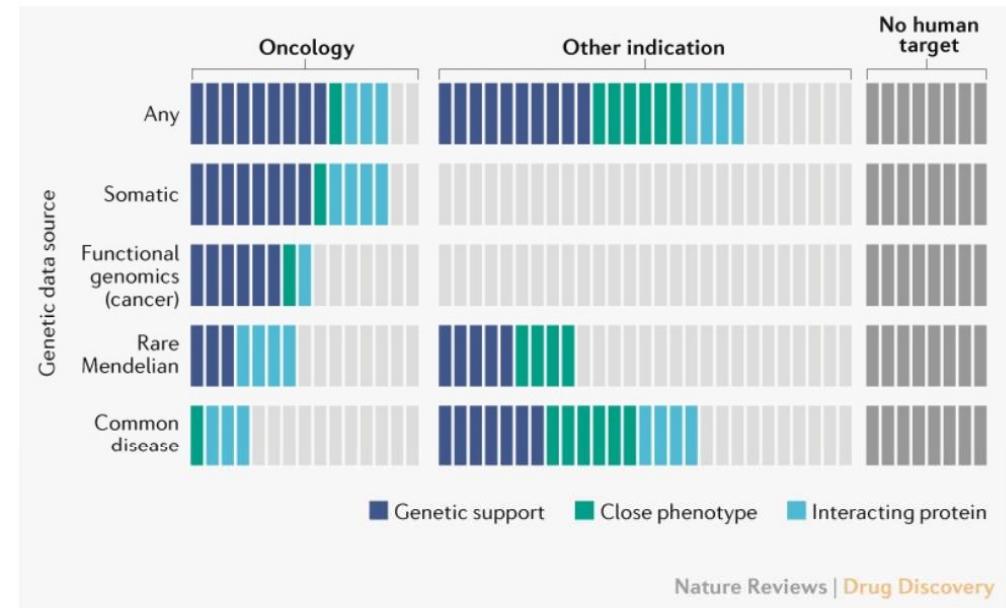
Explore content ▾ About the journal ▾ Publish with us ▾ Subscribe

[nature](#) > [nature reviews drug discovery](#) > [biobusiness briefs](#) > article

BIOBUSINESS BRIEFS | 08 July 2022

Human genetics evidence supports two-thirds of the 2021 FDA-approved drugs

David Ochoa✉, Mohd Karim, Maya Ghoussaini, David G. Hulcoop, Ellen M. McDonagh & Ian Dunham



Nature Reviews | Drug Discovery

Fig. 1 | Supporting human genetic evidence for new drugs approved by the FDA in 2021. Availability of genetic evidence implicating each of the 50 drug targets (columns) as likely causal genes for the disease for which the drug is indicated. Genetic data sources are stratified (rows) based on the predominant nature of the genetic information. Evidence is classified based on whether it directly supports the gene-indication pair or any of the closely related phenotypes or proteins. See Supplementary information for details and an expanded figure. Data source: [Open Targets Platform](#) (November 2021).

Outline

1. Big Data & Genomic Big Data
2. Genomic Data
3. The sharing dilemma
4. What is the EGA for?
5. Beacon
6. The EGA Federation



But... what did actually
happen in medical genomics?
(without many people noticing)

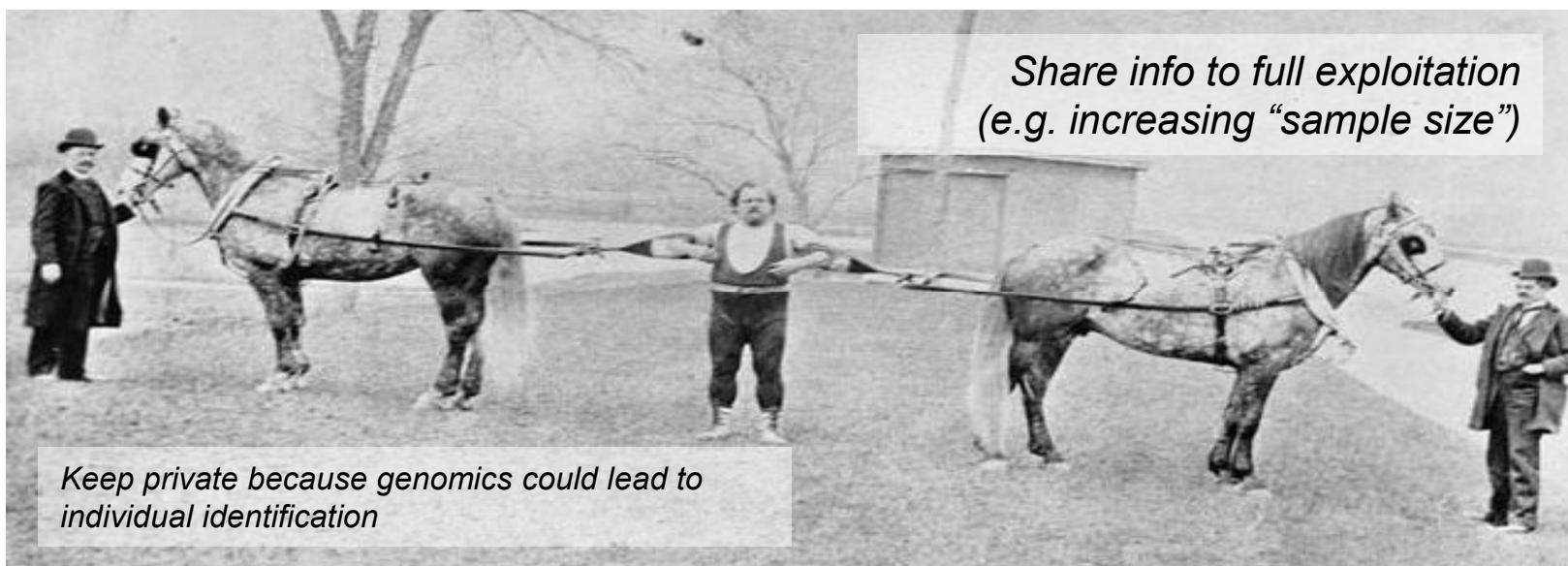
*“We will only share your data for
research purposes.”*

The power of a simple sentence... in a consent form

What is needed for you to guarantee that?

the “Sharing dilemma”

A human rights problem* with genomic information

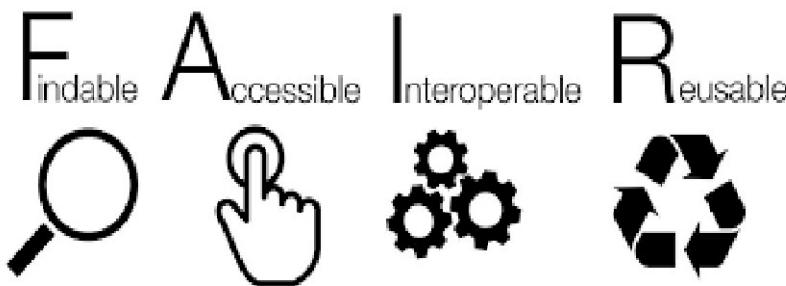


- Private Identifying Information – need to be kept confidential
- Generated with taxpayer’s money – need to be made public

*Arts. 12 and 27 of the Universal Declaration of Human rights

my own position in the dilemma:

Scientific data must be FAIR





 [Login](#) [Register](#) [Need Help?](#)

ABOUT DISCOVERY SUBMISSION ACCESS

 Search...

The European Genome-phenome Archive (EGA) is a service for permanent archiving and sharing of personally identifiable genetic, phenotypic, and clinical data generated for the purposes of biomedical research projects or in the context of research-focused healthcare systems.

I want to...



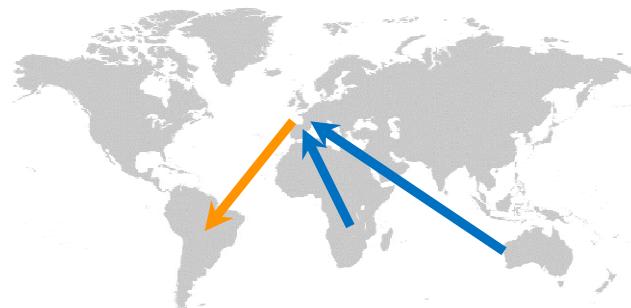
Events

16-17 NOVEMBER 2023 | ELIXIR COMMUNITIES | BARCELONA, SPAIN / VIRTUAL
ELIXIR Federated Human Data Community Day

Latest Blog Posts

Data access models

Fully Open research data



Aggregate data globally

Download, analyse locally

e.g., Plant and animal genomes, Population-based results, Some GWAS summary statistics

Managed access data



Different **controlled access** silos

Download, analyse locally

e.g., UK BioBank, dbGaP/EGA, GWAS consortia

Outline

1. Big Data & Genomic Big Data
2. Genomic Data
3. The sharing dilemma
4. What is the EGA for?
5. Beacon
6. The EGA Federation



The EU Solution: the European Genome-phenome Archive (EGA)

The EGA is an Archive for human genomic data that should have controlled access because it allows identifying individuals

- ✓ Data is provided by research centers and health care institutions
- ✓ Access is controlled by Data Access Committees
- ✓ Data requesters are researchers from other research or health care institutions

Our goal is to foster responsible data reuse for the good of science and all humankind

<https://ega-archive.org>



A screenshot of the European Genome-phenome Archive (EGA) website. The top navigation bar includes links for 'Login', 'Register', and 'Need Help?'. Below the navigation is a search bar. The main content area features a large image of a person holding a tablet, with text overlaying it that reads 'I want to...'. Three circular buttons below this are labeled 'DOWNLOAD DATA', 'DEPOSIT DATA', and 'MANAGE DATA'. To the right of this is a smaller window showing the EGA homepage with sections for 'REST API', 'SUBMISSION REST API IS LIVE', 'EGA WINTER SUMMIT 2016', '17TH ICDC SCIENTIFIC WORKSHOP', 'THE FINAL PHASE OF THE 1000 GENOMES PROJECT', and 'EGAS EXCELLENCE KICKS-OFF'. On the far right, there are 'EVENTS', 'IMPROVEMENTS', 'SCIENTIFIC NEWS', and 'TAGS' sections. At the bottom left, the text 'Events' is visible.

What is the EGA for?

- Science needs data:
 - Data you collect,
 - ... that journals & funders will ask you to share for reproducibility,
 - ...and other scientist will leverage in re-analysis and meta-analyses
- EGA supports Open Science for human sensitive data, in the frontier between pure and health-bound biomedical research

Who manages the EGA?

The EGA was created in 2008, to host the W

CRG joined in 2013

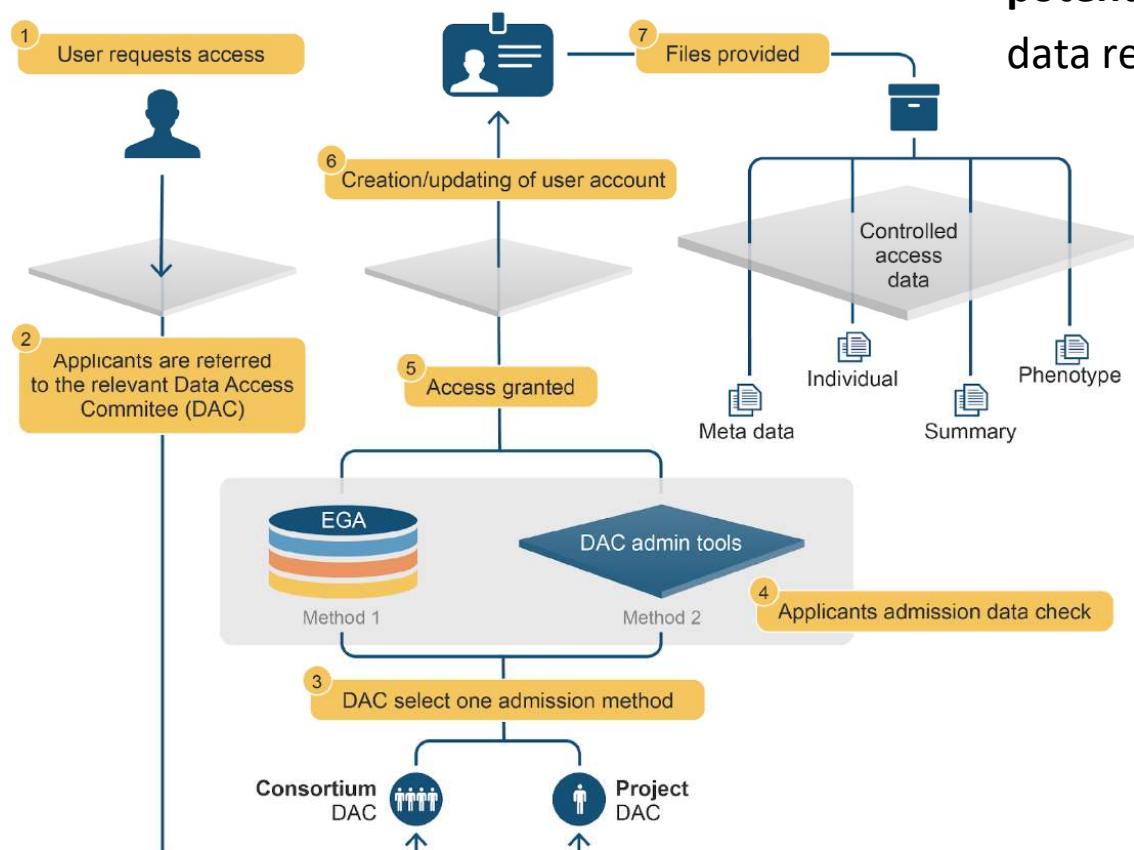


Project goal:

To transform the EGA to a joint



EGA's differential assets



The EGA is **free worldwide** resource for **permanent secure archiving** and sharing of all types of **potentially identifiable** biomolecular and *phenoclinic* data resulting from biomedical research projects.

Data is provided by research centers and health care institutions.

Control of access is kept by every study Data Access Committee (*data controller*).

Data requesters are researchers from other research or health care institutions.

The EGA contains a huge variety of data

Studies ~7,600

Datasets ~15,000

DACs ~3,500

Submitters ~2,250

Requesters ~32,500

Files ~3.5 Million

Archive size ~17,4 Petabytes

Updated Nov 11th

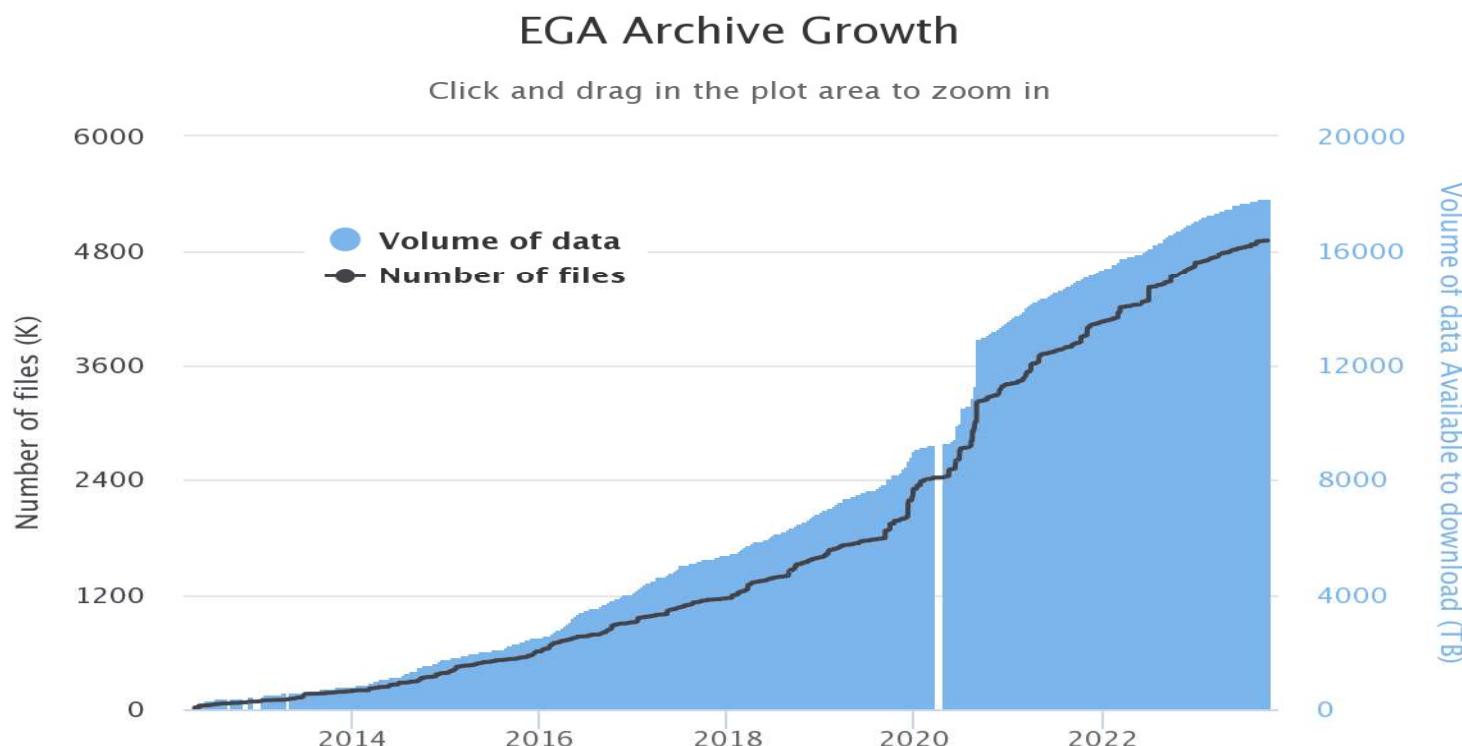


The screenshot shows the EGA homepage with the following key features:

- Header:** Includes the EGA logo, search bar, navigation menu (ABOUT, SUBMISSION, BROWSE, ACCESS, DOWNLOAD, METADATA), and user options (Helpdesk, Log in).
- Section: What is in the EGA?** A bar chart titled "Studies in the EGA by disease" showing the number of studies across different disease categories:

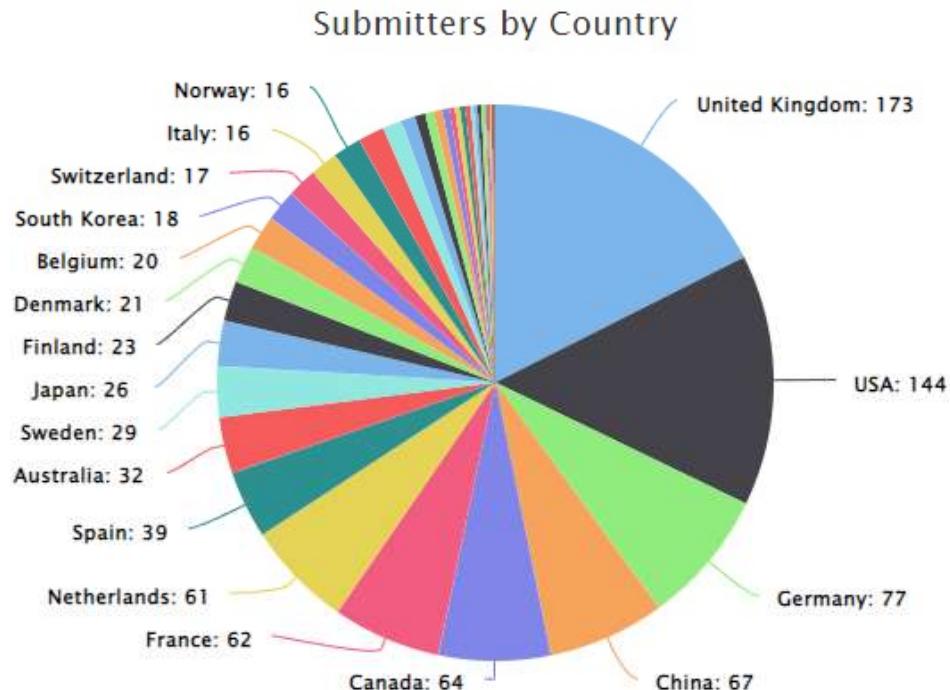
Disease Category	Number of Studies
Cancer	1,579
Cardiovascular	168
Infectious	44
Inflammatory	96
Neurological	78
Other	1,466
- Section: Latest studies** featuring a news item about SARS-CoV-2 receptor ACE2 and TMPRSS2 expression in bronchial transient secretory cells.
- Footer:** Includes links to ELIXIR, CRG, and EMBL-EBI, along with copyright information and legal notices.

EGA's Growth

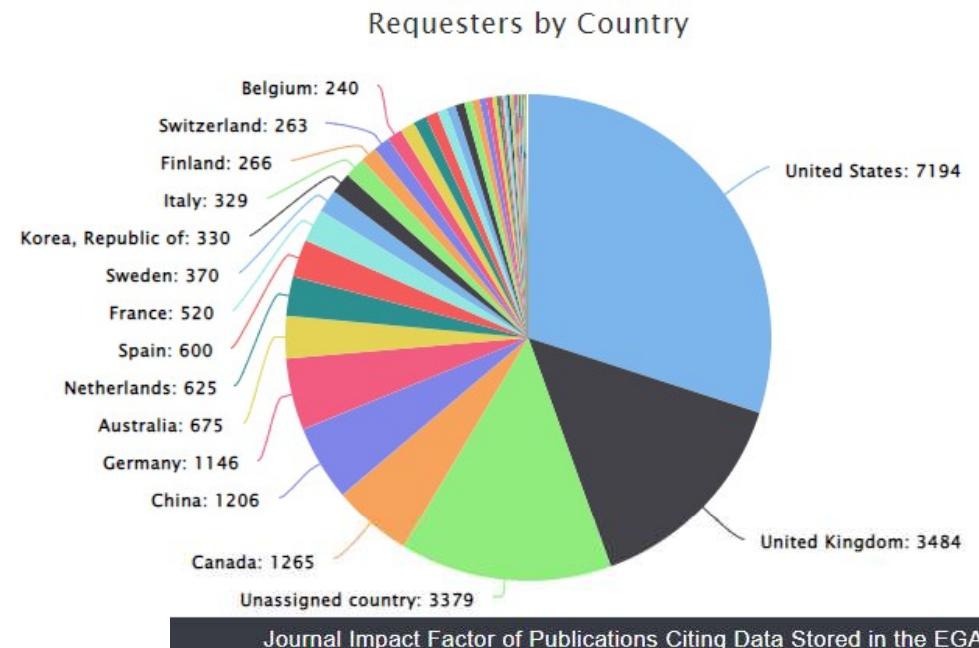


Who is using the EGA?

Submitter Accounts Created by Country



Requester Accounts Created by Country



Journal Impact Factor of Publications Citing Data Stored in the EGA



Outline

1. Big Data & Genomic Big Data
2. Genomic Data
3. The sharing dilemma
4. What is the EGA for?
5. Beacon
6. The EGA Federation



is all this enough?

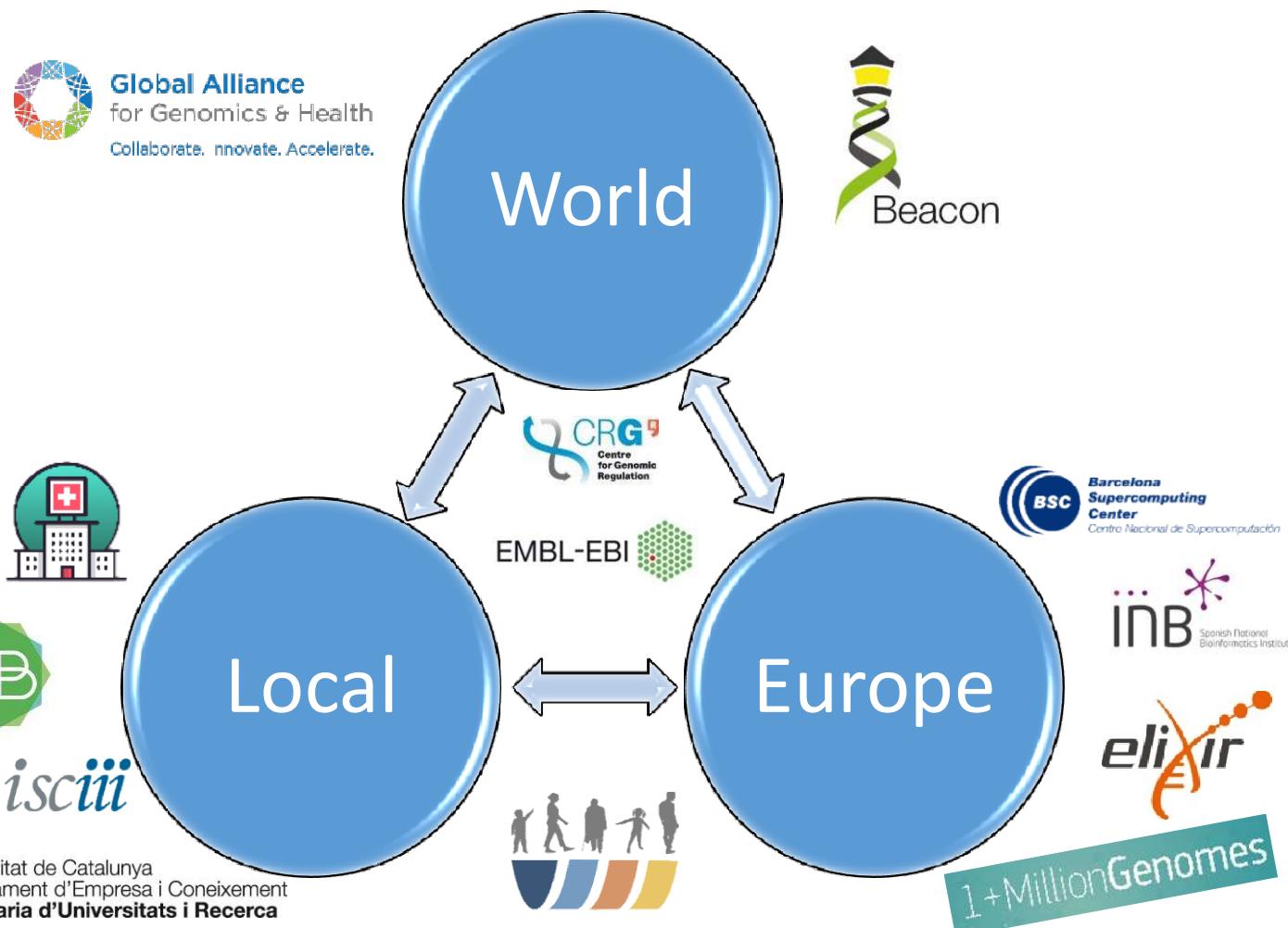
GDPR and its multifaceted interpretations, plus different jurisdictions and interests, plus practical and economic issues make the things difficult.

Fortunately...

*GDPR is the EC's General Data Protection Regulation

...the EGA is part of
a rich ecosystem

Pioneering GDPR-compliant data stewardship (comprehensive data discovery, secure data storage and legally-compliant data sharing)



The GA4GH the Global Alliance for Genomics and Health

<https://ga4gh.org>



[GA4GH: The Global Standards Organization for Genomics - YouTube](#)



Creating a Global Alliance to Enable Responsible Sharing of Genomic and Clinical Data

June 3, 2013

An initial draft of this White Paper was prepared for the January 2013 meeting and has since been revised substantially based on discussions at and since the meeting. A list of contributors and participants is provided at the end of this document.

Mission

The Global Alliance for Genomics and Health aims to accelerate progress in genomic research and human health by cultivating a common framework of standards and harmonized approaches for effective and responsible genomic and health-related data sharing.

Members

Experts in healthcare, research, patient advocacy, life science, and information technology.





GA4GH **Driver Projects** are real-world genomic data initiatives that provide input on the standards most needed for the international genomics community to share data. Driver Projects put finished standards and frameworks to immediate use.

GA4GH Driver Projects

National Initiatives



Pan-African



Pan-European

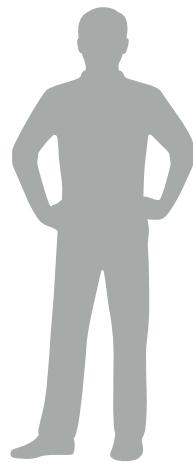


International



lighting beacons

we have made an enormous effort to distinguish discovery from sharing

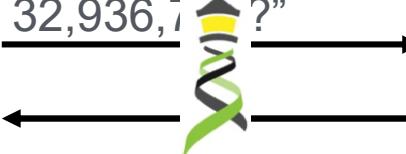


Geneticist



“Do you have a ‘C’
at chromosome 13
at position

32,936,7 ?”



“Yes” (or “no”)

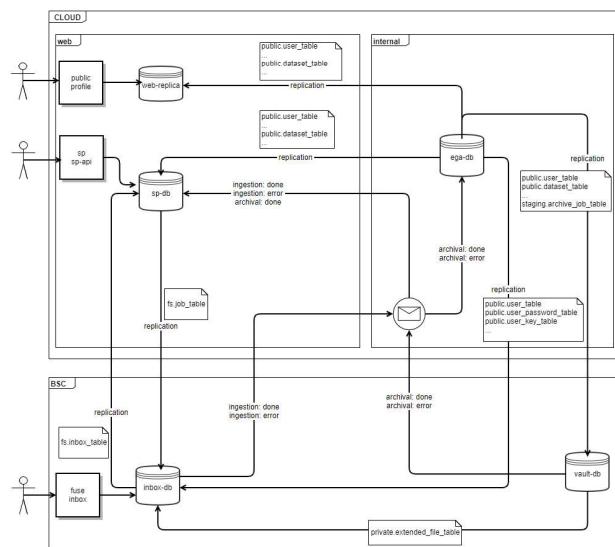
chr13
32,936,732

C

Data Holder

lighting beacons

Beacon allows for increasingly sophisticated queries



Have you seen deletions in this region on chromosome 9 in Glioblastomas from a juvenile patient, in a dataset with unrestricted access?

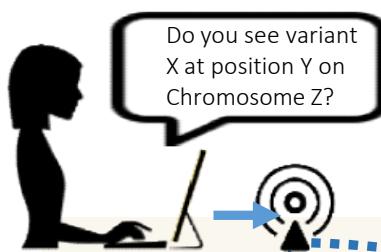
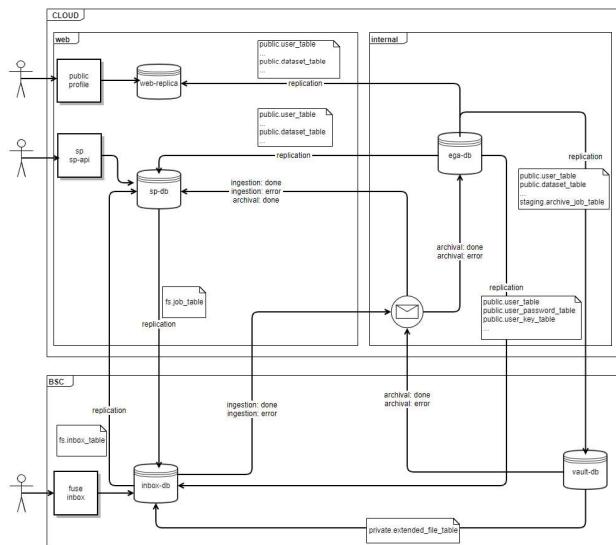


Beacon v2 API

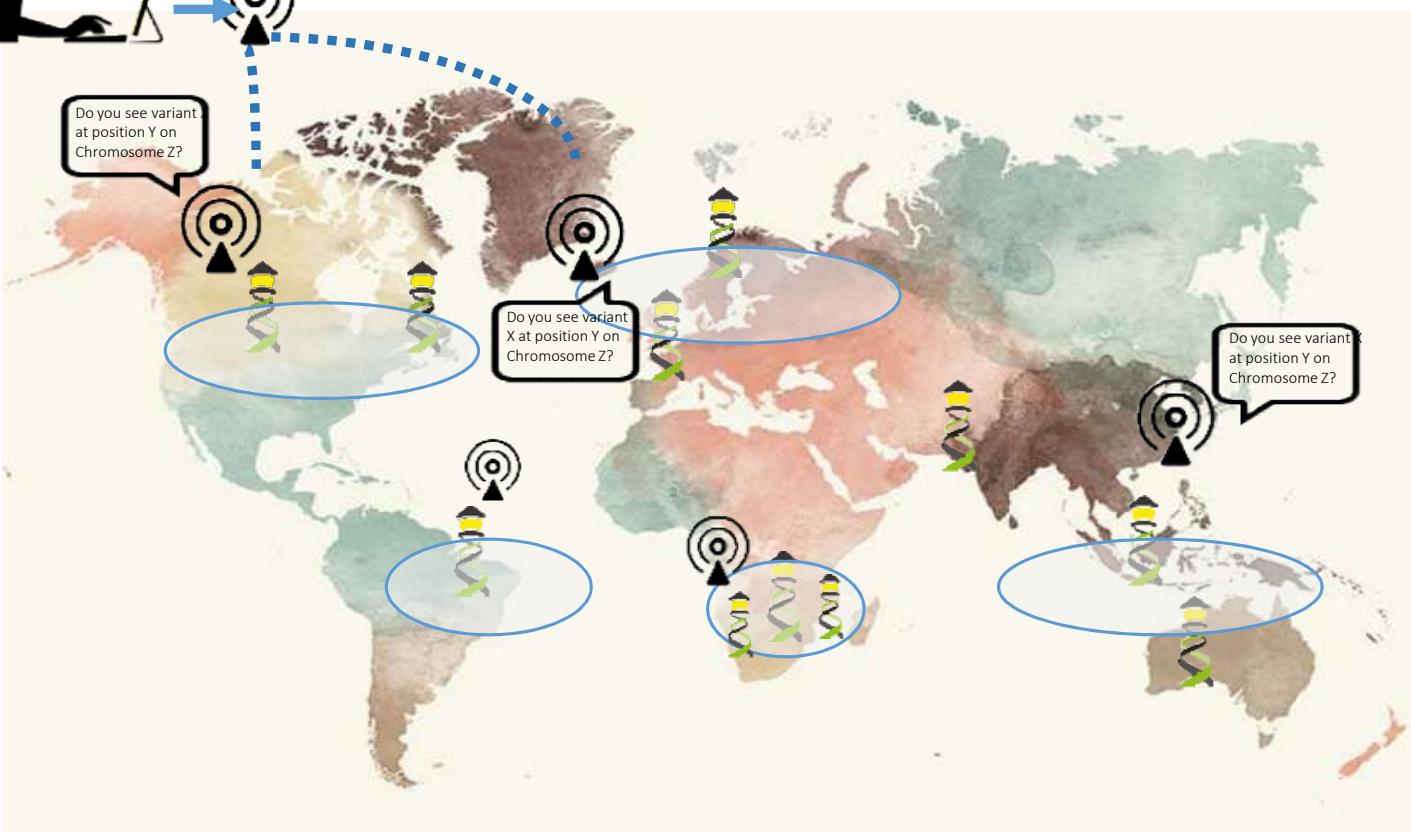
The Beacon API v2 proposal opens the way for the design of a simple but powerful **"genomics API"**.

lighting beacons

we have made an enormous effort to distinguish discovery from sharing

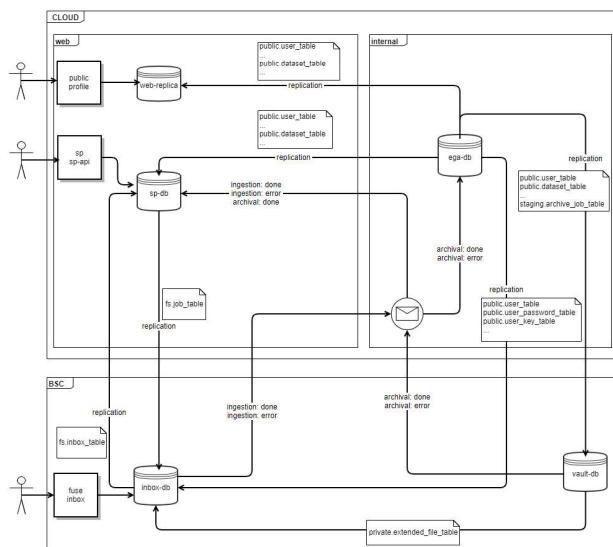


A growing number of Beacons

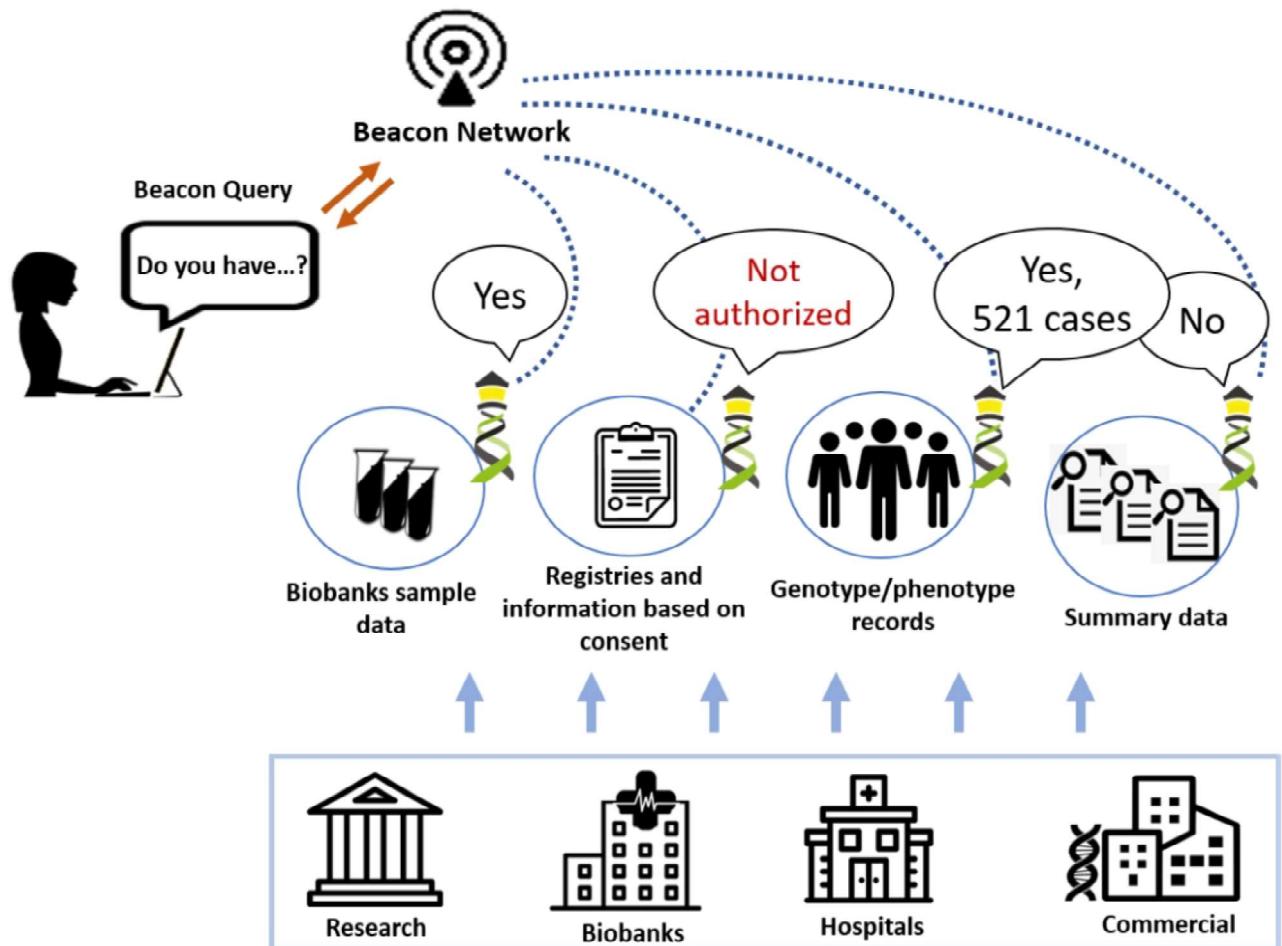


lighting beacons

And the Beacon Network allows for simple entry points

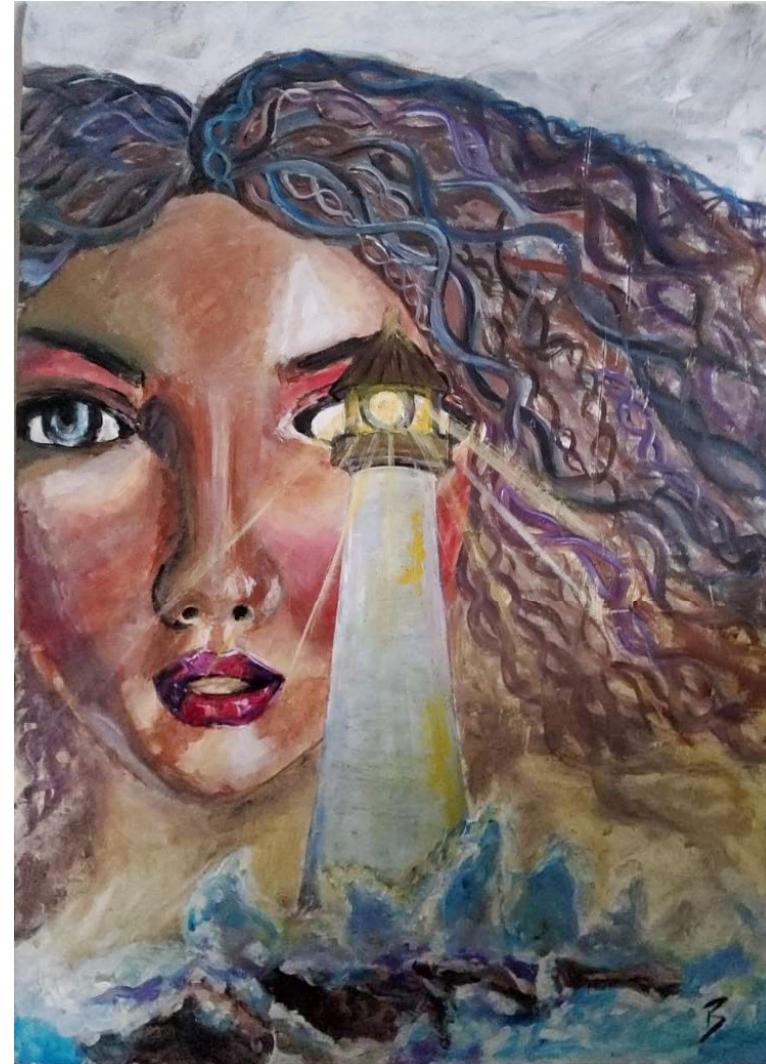


Beacon Network: Aggregator of Beacons



lighting beacons

And so, Beacon is becoming more and more useful, even beyond genomics



Towards the Internet of Genomics

It starts with a local Beacon network...



The map shows the outline of Barcelona with several orange dots representing Beacon locations. One dot is blue, located near the port area. The word "BARCELONA" is printed at the bottom right of the map.

La Marató 3
2019 Malalties minoritàries

"Xarxa interhospitalària catalana de variants genètiques"

- Hospital Sant Joan de Déu (centro coordinador)
IP: Dèlia Yubero
- Hospital Clínic de Barcelona
IP: Eva González
- Hospital Vall d'Hebron
IP: Elena García Arumí
- Hospital de la Santa Creu i Sant Pau
IP: Benjamín Rodríguez
- Hospital de Bellvitge
IP: Ariadna Padró
- Centre de Regulació Genòmica
IP: Babita Singh

Towards the Internet of Genomics

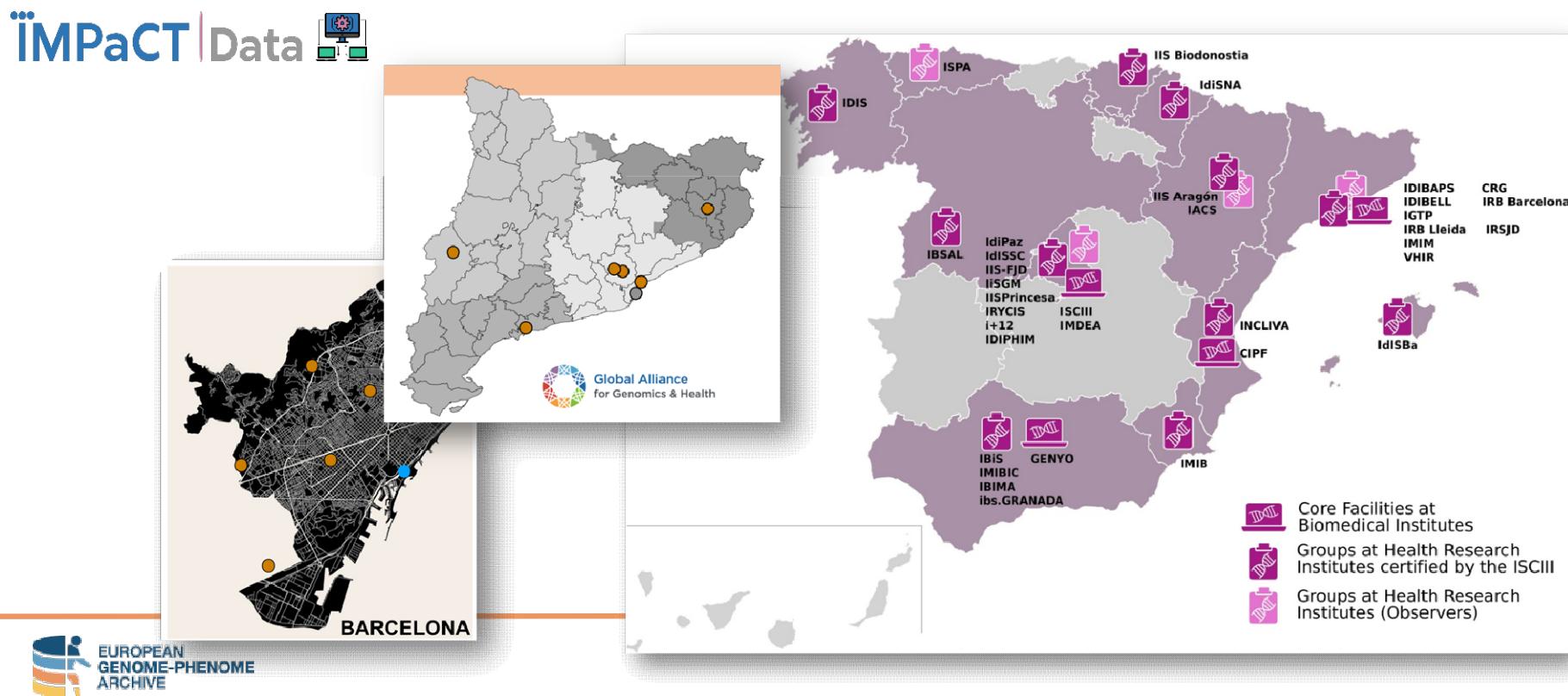
It starts with a local Beacon network...

The collage includes:

- A photograph of a meeting room where a presentation titled "Beacon v2 is a specification" is being given to a group of people seated around tables.
- A map of Catalonia showing the locations of various genomic hubs marked with orange dots.
- Logos for several organizations involved in the network:
 - CLÍNIC BARCELONA Hospital Universitari
 - HOSPITAL DE LA SANTA CREU I SANT PAU FUNDACIÓ DE GESTIÓ SANITÀRIA UNIVERSITAT AUTÒNOMA DE BARCELONA
 - Vall d'Hebron Barcelona Campus Hospitalari
 - CRG Centre for Genomic Regulation
 - Bellvitge Hospital Universitari
 - SJD Sant Joan de Déu Barcelona · Hospital
 - Global Alliance for Genomics & Health
 - IMIM Institut Hospital del Mar d'Investigacions Mèdiques
 - Parc Taulí Hospital Universitari
 - Germans Trias i Pujol Hospital
 - Fundació Puigvert
 - GCAT Genomes for Life
 - IGTP Instituto de Genética Molecular de Barcelona
 - CST Consorci Sanitari de Terrassa
 - Hospital Universitari Doctor Josep Trueta
 - HJ23 Hospital Joan XXIII
 - Hospital Universitari Arnau de Vilanova Lleida
 - La Marató 3

Towards the Internet of Genomics

That may be part of a state-wide network...



is all this enough?

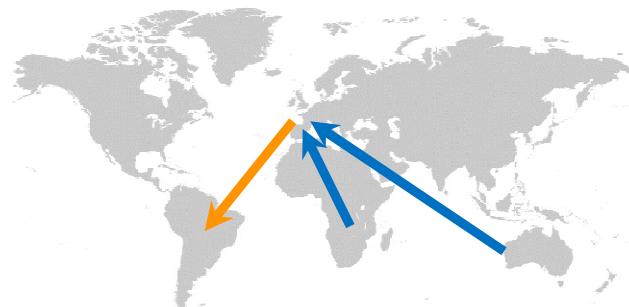
GDPR and its multifaceted* interpretations, plus different jurisdictions and interests, plus practical and economic issues make the paperwork cumbersome, plus the Centralized Stewardship model difficult to sustain

Oh, gosh!!

*and nefarious

Data access models

Fully Open research data



Aggregate data globally

Download, analyse locally

e.g., Plant and animal genomes, Population-based results, Some GWAS summary statistics

Managed access data



Different **controlled access** silos

Download, analyse locally

e.g., UK BioBank, dbGaP/EGA, GWAS consortia

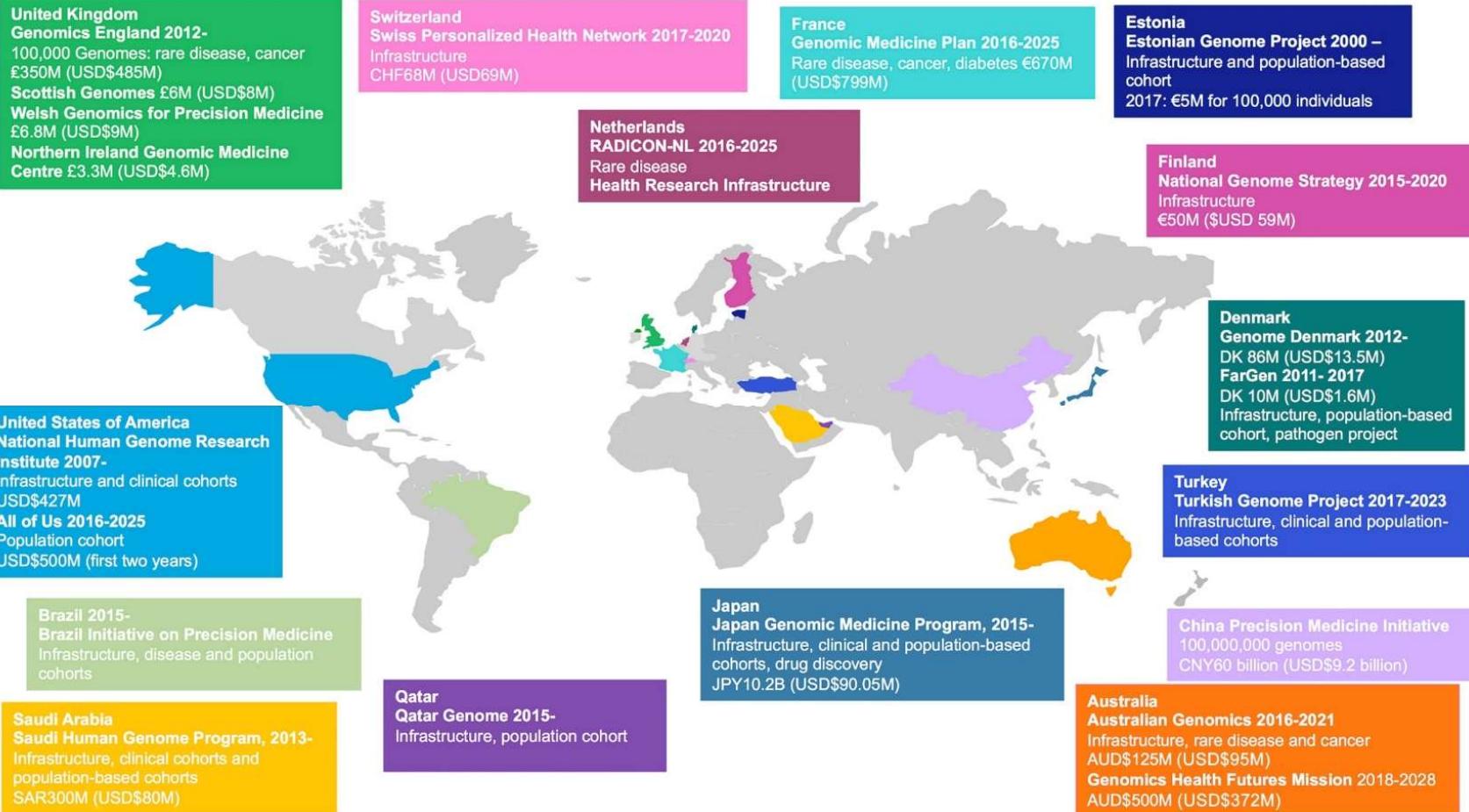
Integrating Genomics into Healthcare: A Global Responsibility

Zornitza Stark, Lena Dolman, Teri A. Manolio, Brad Ozenberger, Sue L. Hill, Mark J. Caulfield, Yves Levy, David Glazer, Julia Wilson, Mark Lawler, Tiffany Boughtwood, Jeffrey Braithwaite, Peter Goodhand, Ewan Birney, Kathryn N. North



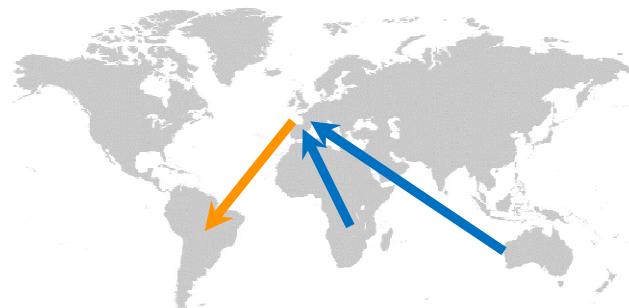
Global Alliance
for Genomics & Health

The American Journal of Human Genetics 2019 104:13-20 DOI: (10.1161/ajhg.2018.11.014)



Data access models

Fully Open research data

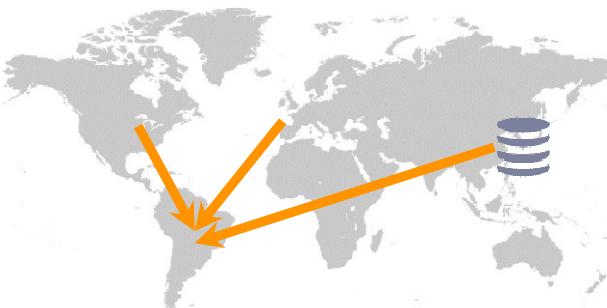


Aggregate data globally

Download, analyse locally

e.g., Plant and animal genomes, GWAS results

Managed access data

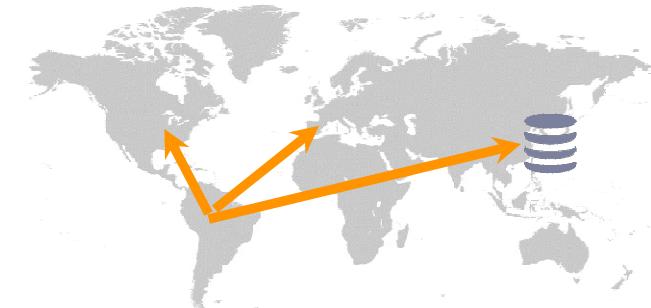


Different **controlled access** silos

Download, analyse locally

e.g., UK BioBank, dbGaP/EGA, GWAS consortia

Secondary Research use of healthcare data



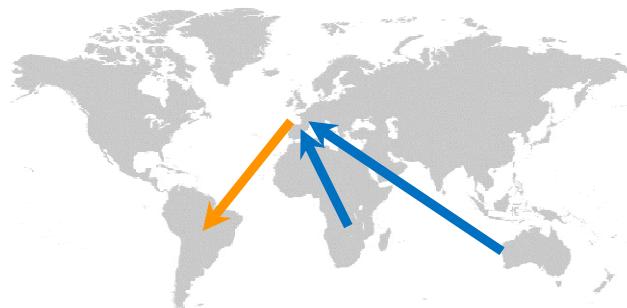
Different jurisdictionally restricted datasets

Send analysis containers into trusted research environments

e.g., Genomics England data, Danish Genome Cohort

Data access models

Fully Open research data



Aggregate data globally

Download, analyse locally

e.g., Plant and animal genomes, GWAS results

Managed access data

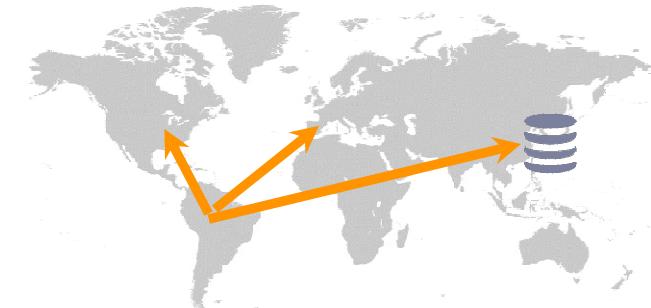


Different **controlled access** silos

Download, analyse locally

e.g., UK BioBank, dbGaP/EGA, GWAS consortia

Secondary Research use of healthcare data



Different jurisdictionally restricted datasets

Send analysis containers into trusted research environments

e.g., Genomics England data, Danish Genome Cohort

Outline

1. Big Data & Genomic Big Data
2. Genomic Data
3. The sharing dilemma
4. What is the EGA for?
5. Beacon
6. The EGA Federation



A potential solution

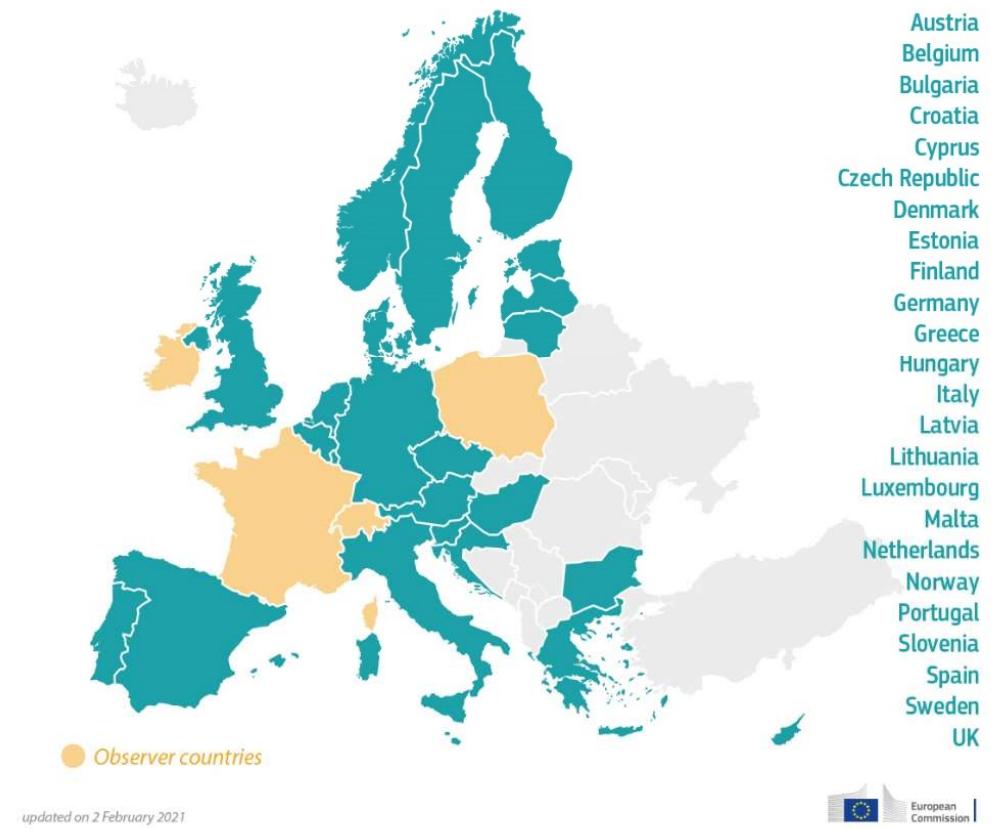
Federated EGA

Setting the path for actual federated sharing of sensitive human data

1+ Million Genomes



Countries that have signed the 1+MG Declaration since 2018



a complex European context

Endeavors such as the 1+Million Genomes Initiative & the Beyond 1MGenomes Project are made possible by the Federated EGA or the Federated EGA Technology



Mariya Gabriel, Commissioner for Digital Economy and Society welcomed this new step forward:

"EU cooperation on linking genomics and health data becomes ever more important to personalised health care, better research and innovation, and excellence in research. We are delighted to welcome Latvia joining the other 18 signatory States and look forward to enhancing our collaboration."

Countries that have signed the 1+MG Declaration since 2018



the European context

Endeavors such as the 1+Million Genomes Initiative & the Beyond 1MGenomes Project are made possible by the **Federated EGA or the Federated EGA Technology**

Partnerships and community formation

ELIXIR Human Data Communities

- Federated Human Data
- Rare Diseases
- human Copy Number Variation



Our role in the GDI* project

*European Genome Data Infrastructure

- CEGA and FEGA nodes are **key** partners of the European Genomic Data Infrastructure initiative
- FEGA software is one possible solution to become **operational** within GDI
- Joining FEGA is a possible way to **comply** with GDI mandate



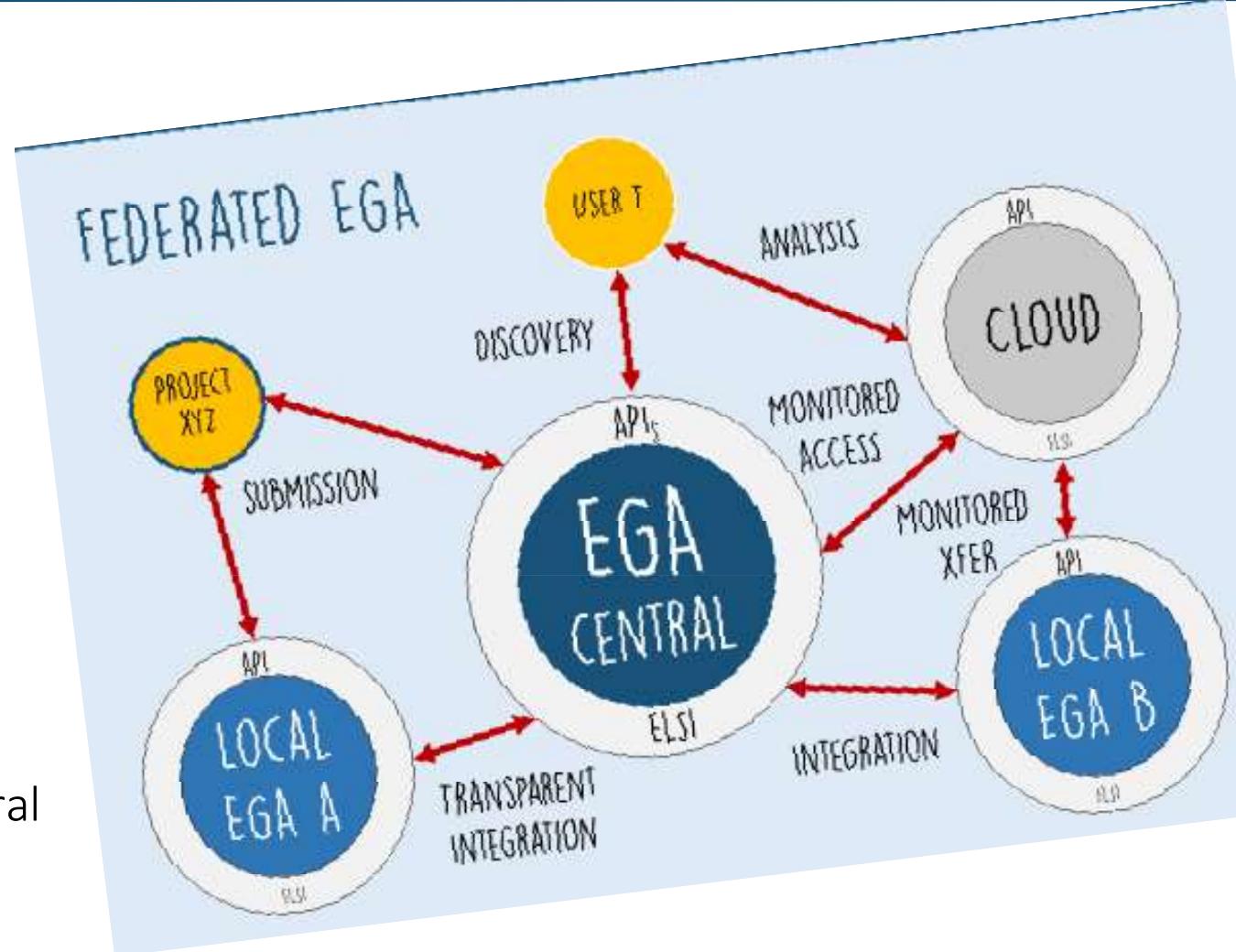
Federated EGA: The Vision

Federated EGA strives to support the discovery of and secure access to human data globally, while respecting national data protection regulations, with the goal of accelerating disease research and understanding and improving human health.

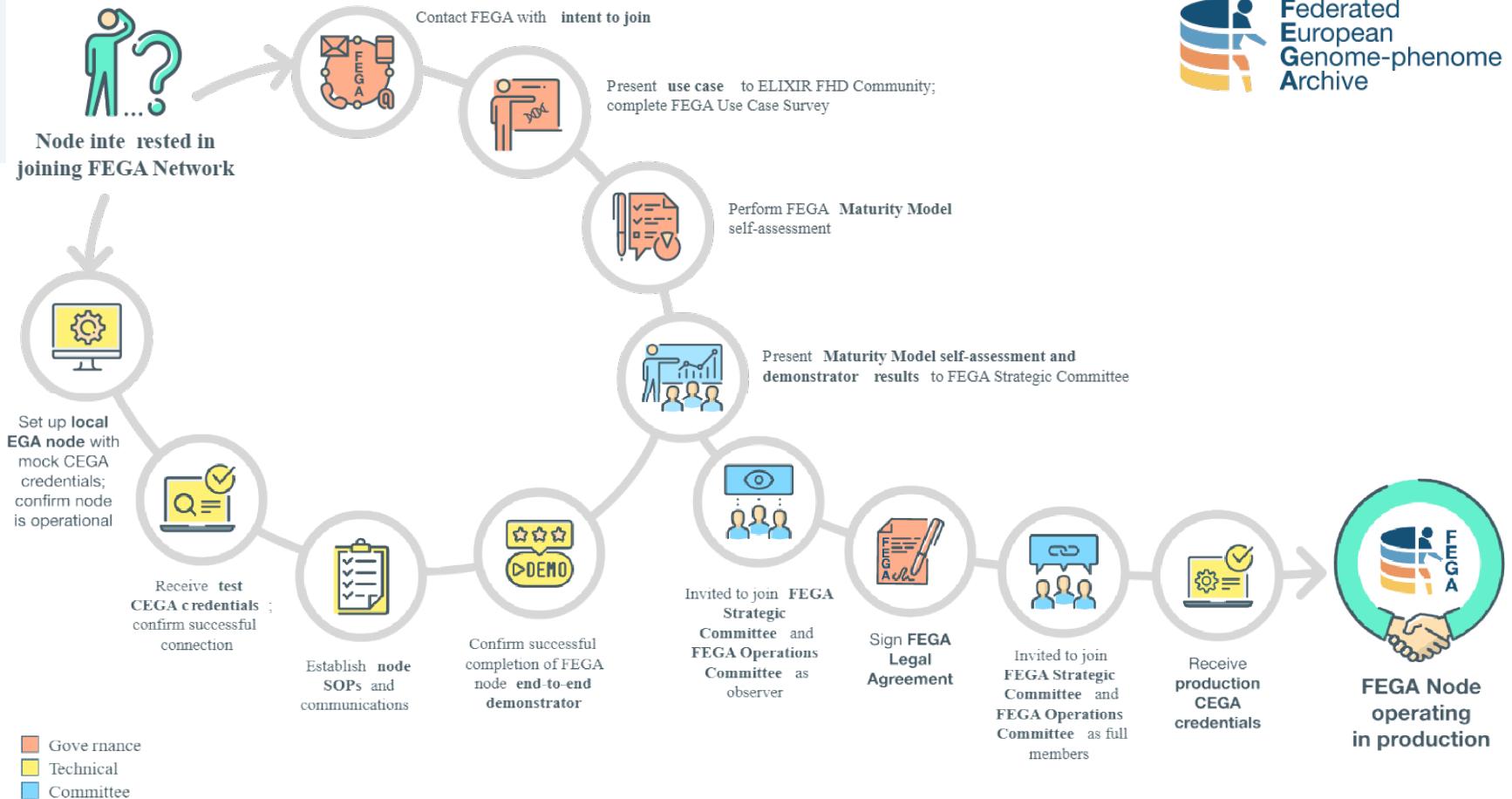


the European context

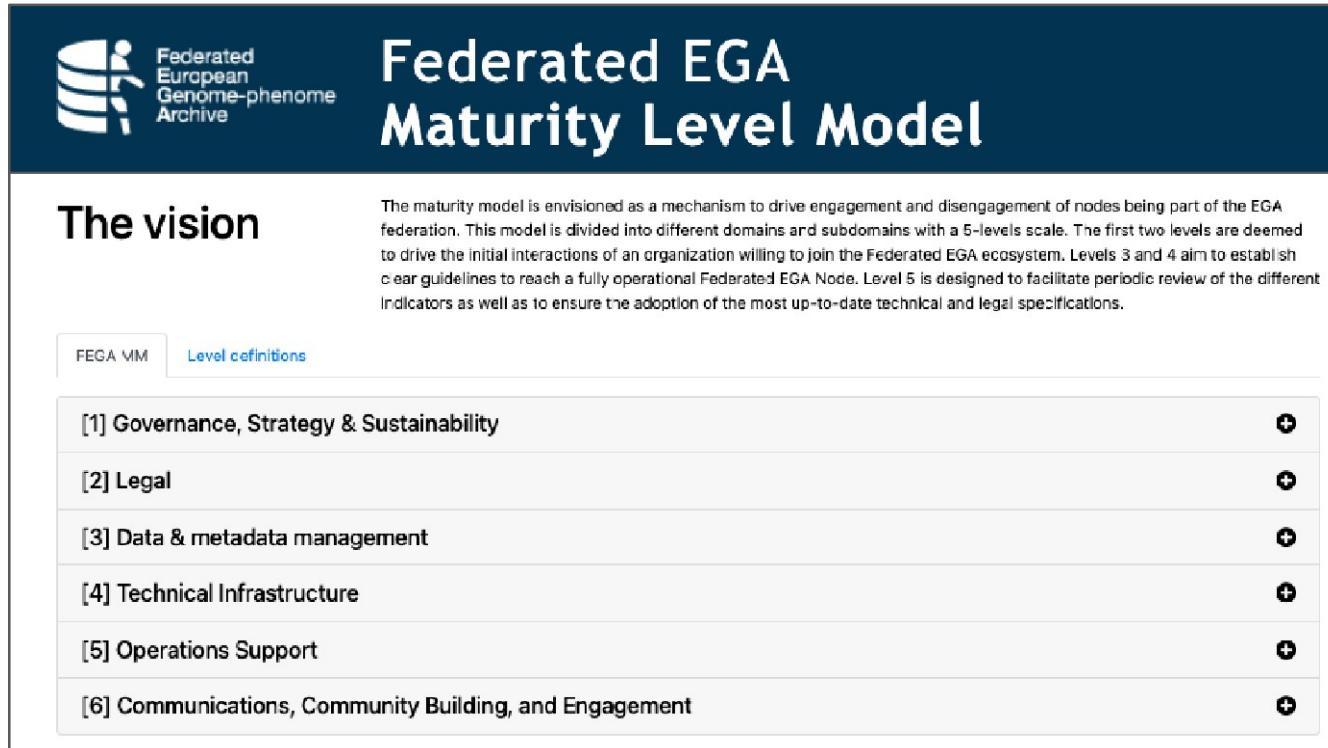
The EGA Federation is a way to achieve secure (and legal!) data discovery and sharing, including cloud features to allow for **federated analysis** (sending the algorithm to bits of data, rather than sending the data to a central computer with the algorithm)



How does the journey look like?



Started self-assessments of FEGA Maturity Model



The screenshot shows the 'Federated EGA Maturity Level Model' interface. At the top left is the FEGA logo. The main title 'Federated EGA Maturity Level Model' is centered above a section titled 'The vision'. Below this, there are two tabs: 'FEGA MM' (selected) and 'Level definitions'. A detailed description of the maturity model follows. A vertical list of six domains is shown, each with a corresponding circular icon to its right.

Domain	Icon
[1] Governance, Strategy & Sustainability	•
[2] Legal	•
[3] Data & metadata management	•
[4] Technical Infrastructure	•
[5] Operations Support	•
[6] Communications, Community Building, and Engagement	•

Align with other
Maturity Models



FAIRplus

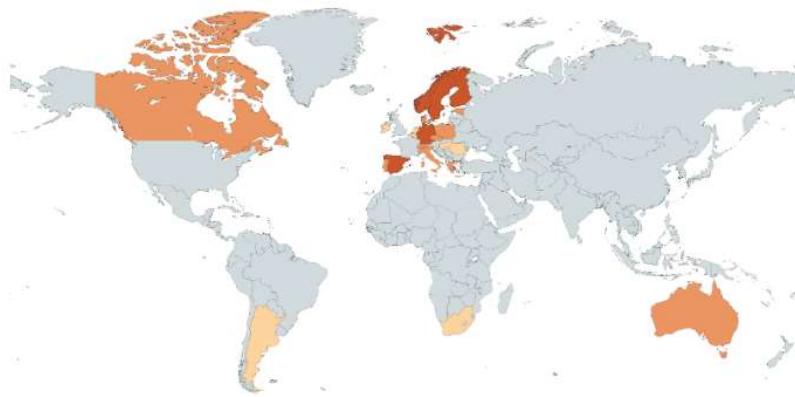
Global Alliance
for Genomics & Health
Collaborate. Innovate. Accelerate.



Albert Hornos, Amy Curwin, Kjell Petersen, Laura Portell-Silva, Mallory Freeberg, Melissa Konopko, Teresa D'Altri, Giselle Kerry, Salvador Capella-Gutierrez

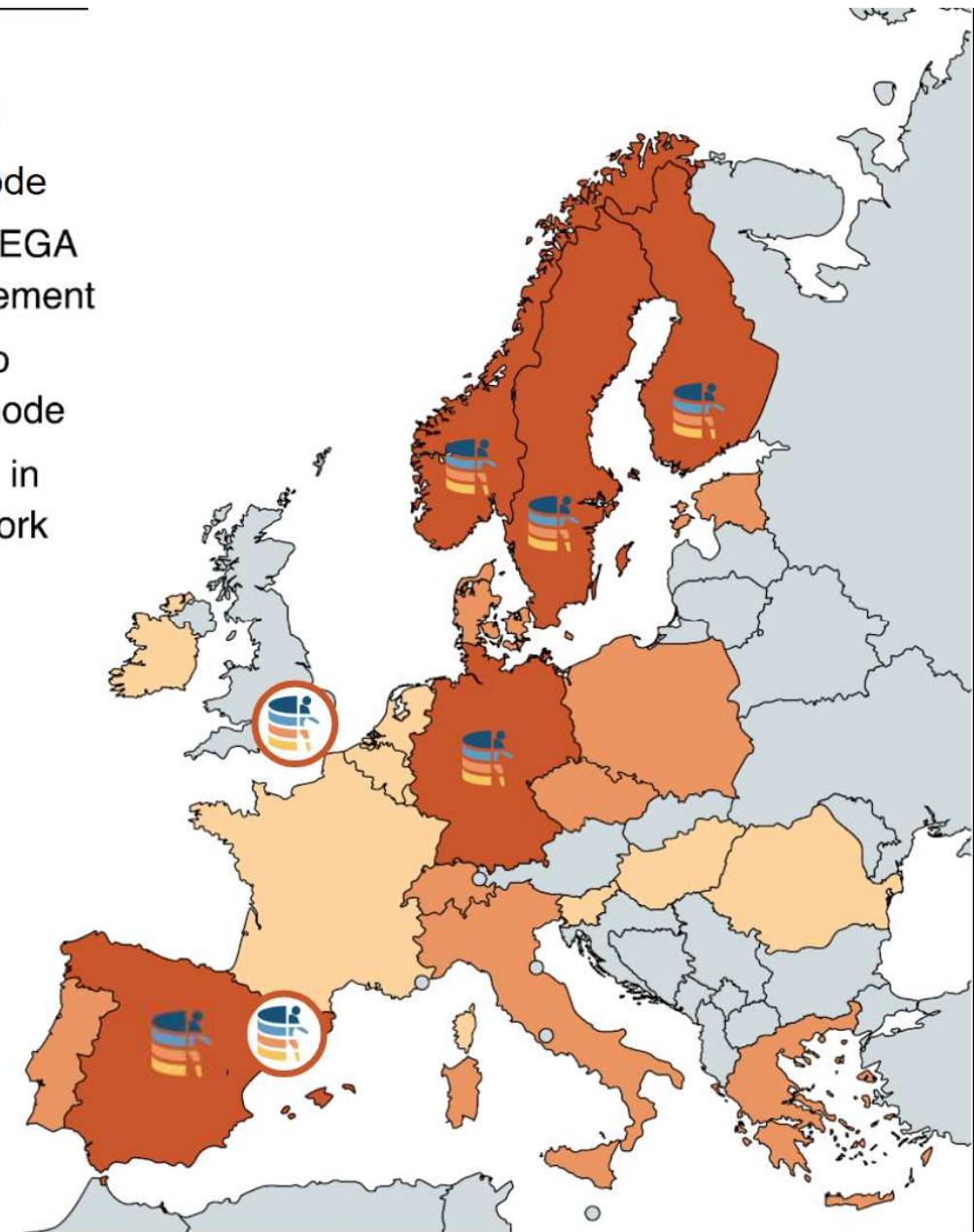
Status as to September 2022

>20 countries /
national initiatives
and counting!



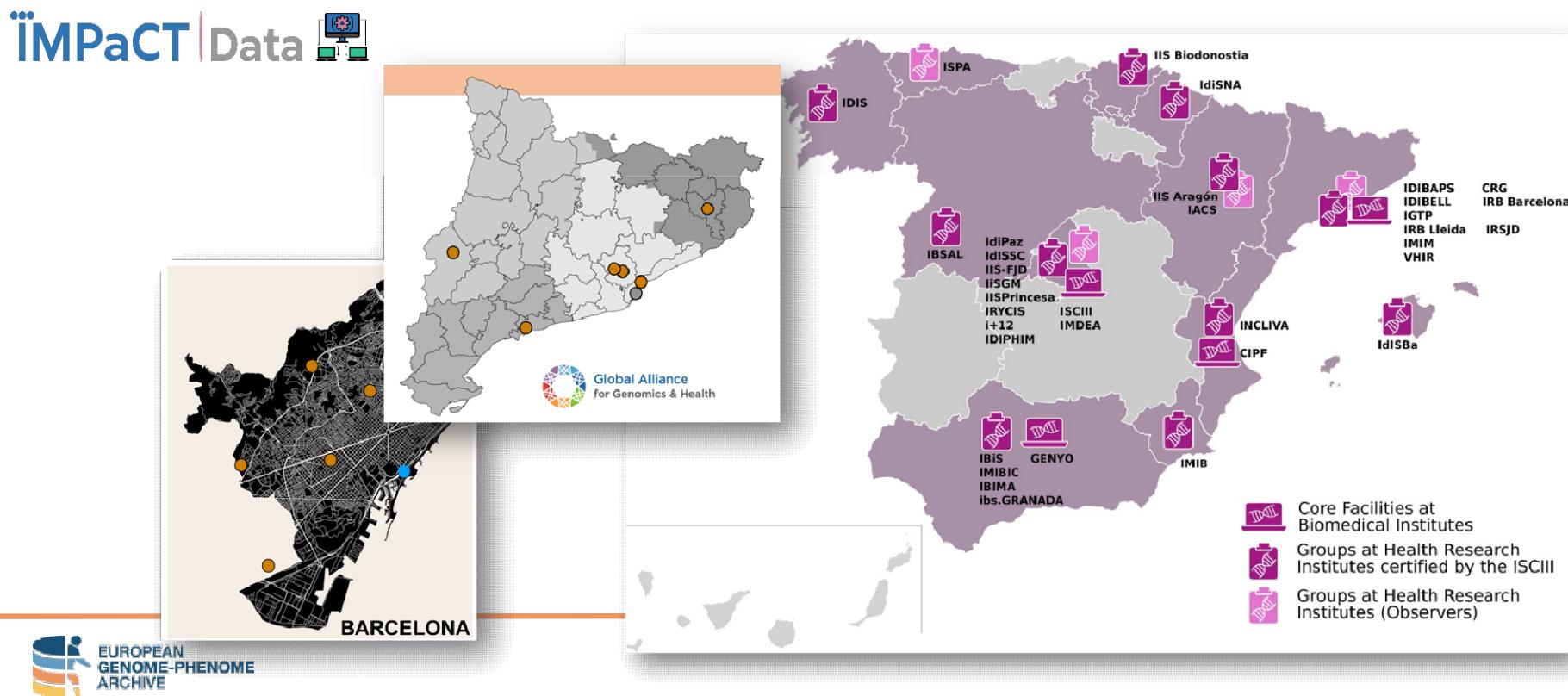
Also Canada, Australia, Argentina, &
H3Africa / South Africa

-
-  Central EGA node
 -  Federated EGA node
 -  Preparing to sign FEGA Collaboration Agreement
 -  Engaging in work to establish a FEGA node
 -  Expressing interest in joining FEGA Network



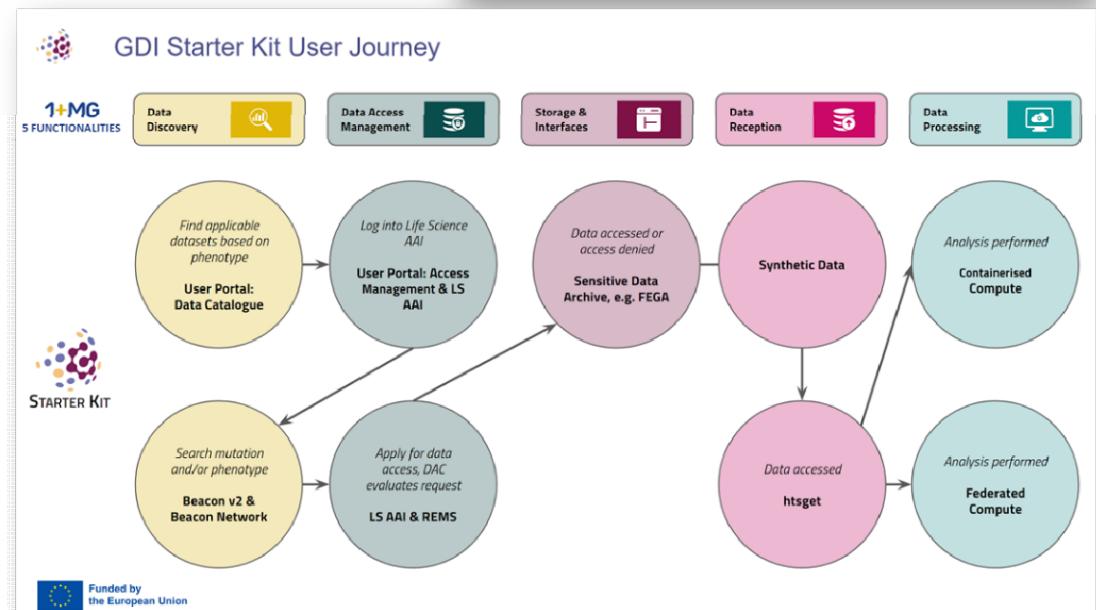
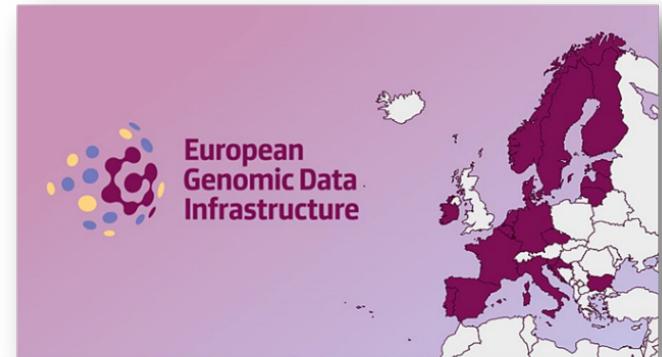
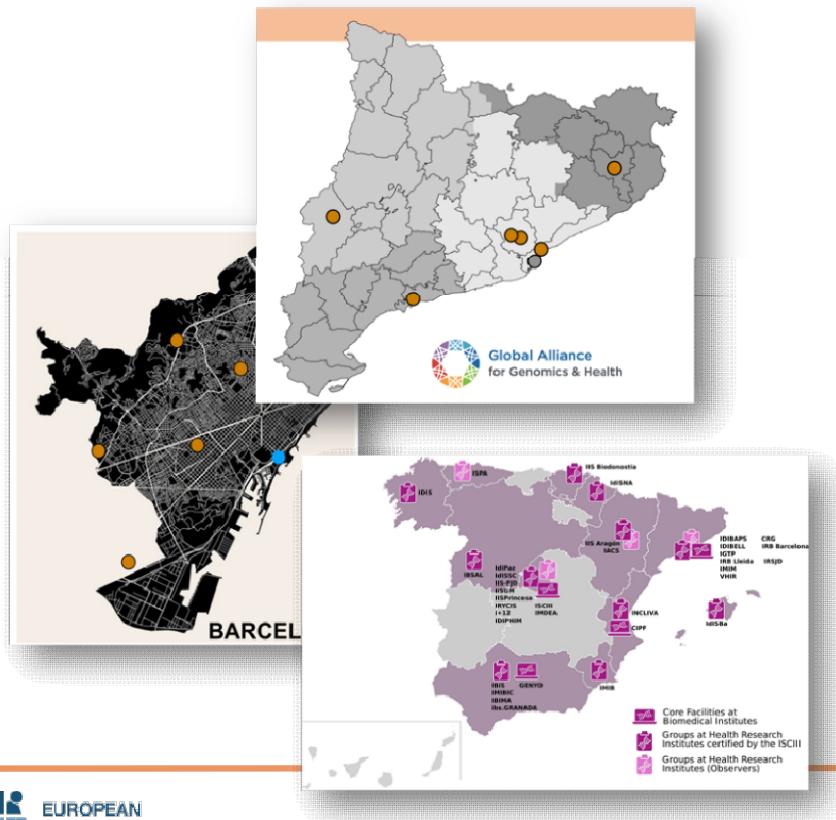
Towards the Internet of Genomics

That may be part of a state-wide network...



Towards the Internet of Genomics

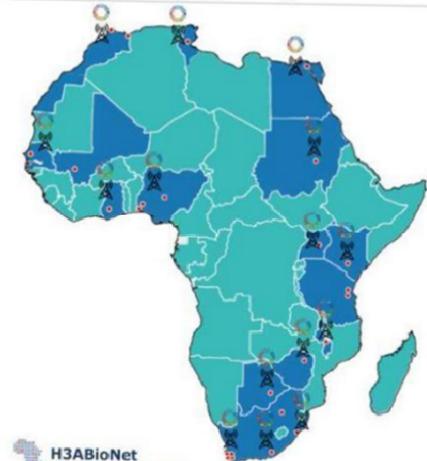
Which is part of a European Network



Towards the Internet of Genomics

And perhaps not just European

Pan-african Beacon Network



H3ABioNet
Pan-African Bioinformatics Network for H3Africa

NEWS | 26 February 2023

Plan for network of Genomics Centres of Excellence across Africa

Closing the gap on access to genomics technologies.

14th International Congress of Human Genetics
22 - 26 February 2023
COMING HOME | ICHG 2023 | AFRICA
Cape Town
www.ichg2023.com

EUROPEAN GENOME-PHENOME ARCHIVE

H3Africa
Human Heredity & Health in Africa

20th Meeting of the H3Africa Consortium
27th - 28th February 2023



Will we have a global network?

