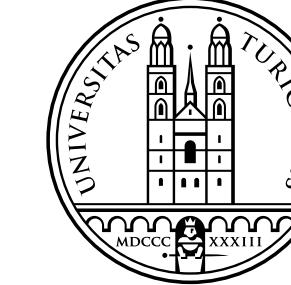




Global Alliance
for Genomics & Health
Collaborate. Innovate. Accelerate.



Swiss Institute of
Bioinformatics



Universität
Zürich^{UZH}

Beaconize this: Federated Data Discovery in Biomedical Genomics



Michael Baudis

Professor of Bioinformatics
University of Zürich
Swiss Institute of Bioinformatics **SIB**
GA4GH Workstream Co-lead *DISCOVERY*
Co-lead ELIXIR Beacon API Development



Republic of Korea
Switzerland. 60th
anniversary
1963–2023

200+ Genomic Data Initiatives Globally

Clinical/Genomic
Medicine



Research



National



Cohorts



How Many Genomes?

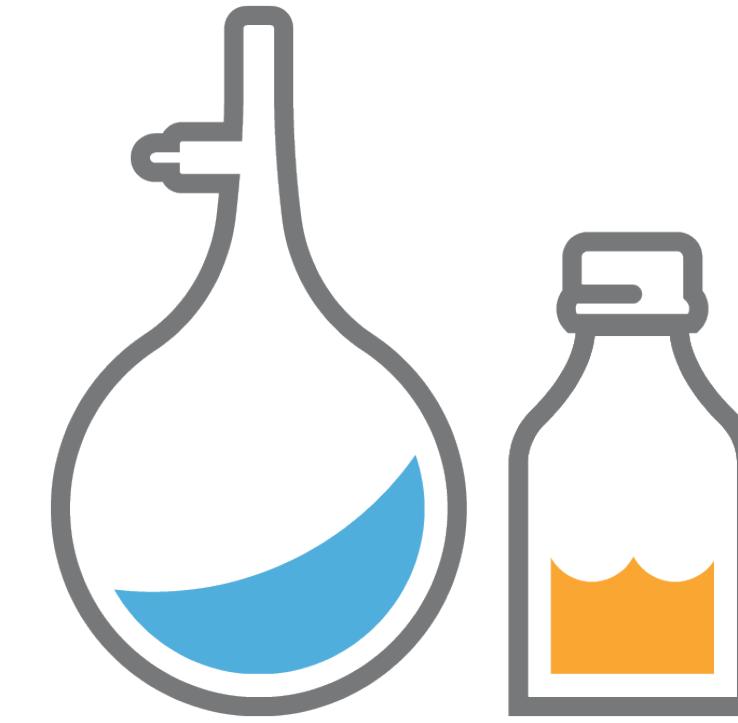


RESEARCH



HEALTHCARE

60M individuals
132.5M sequences



CLINICAL TRIALS

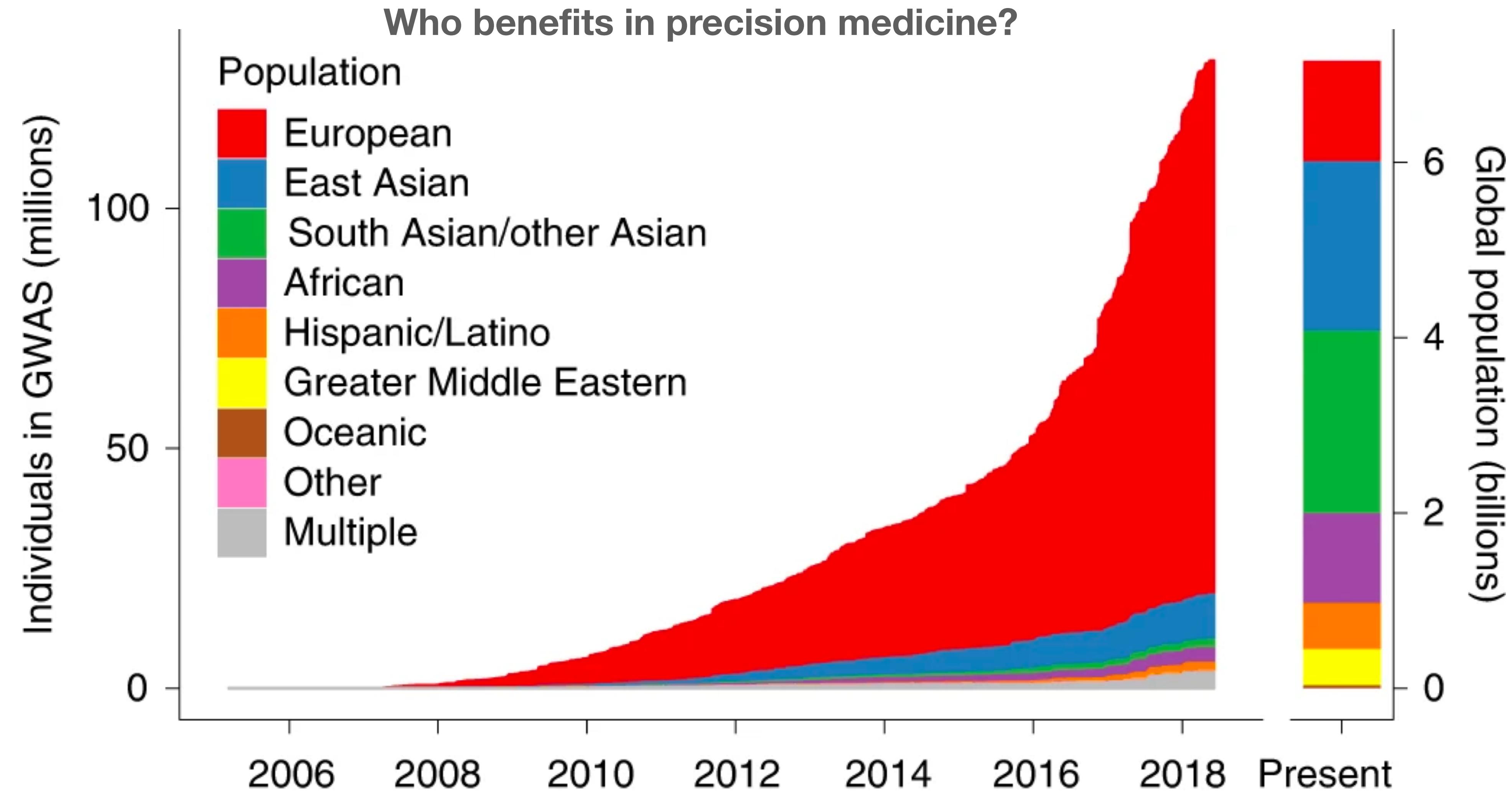
2.7-3M individuals



COHORTS

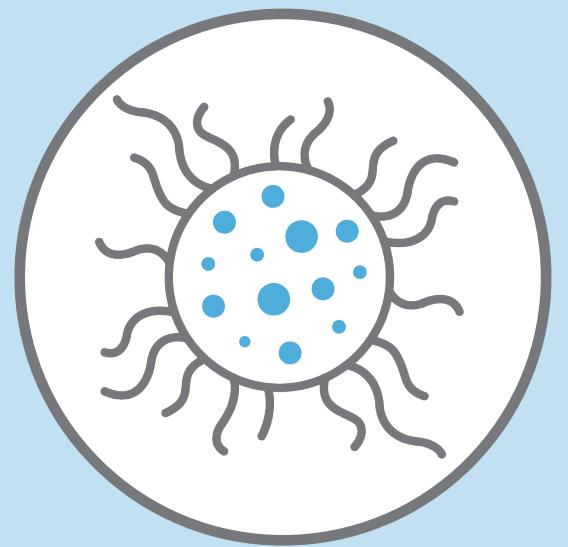
140M individuals

Genomic research has long-standing problems with diversity

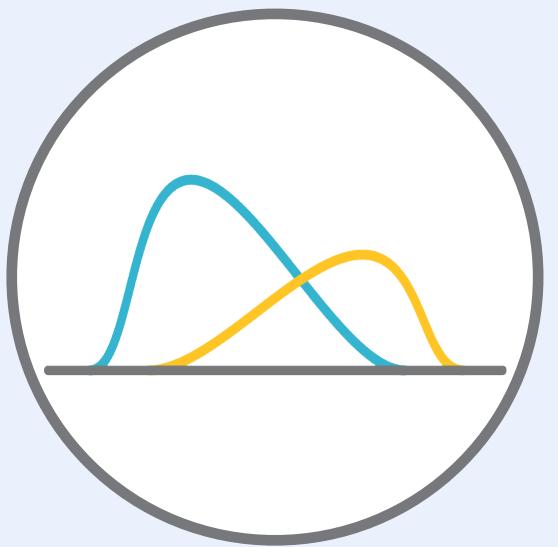




Global Genomic Data Sharing Can...



Demonstrate
patterns in health
& disease



Increase statistical
significance of
analyses



Lead to
“stronger” variant
interpretations



Increase
accurate
diagnosis



Advance
precision
medicine

Different Approaches to Data Sharing



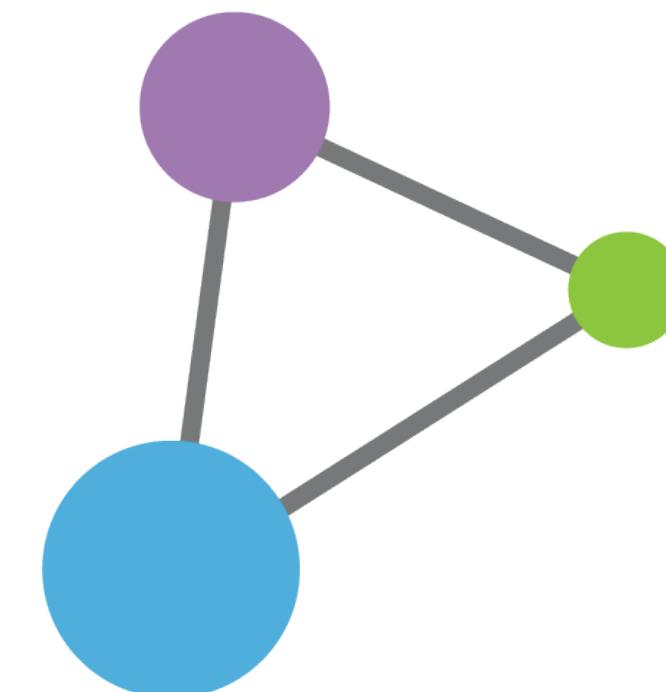
Centralized Genomic Knowledge Bases



Data Commons
Trusted, controlled repository of multiple datasets



Hub and Spoke
Common data elements, access, and usage rules



Linkage of distributed and disparate datasets

Different Approaches to Data Sharing



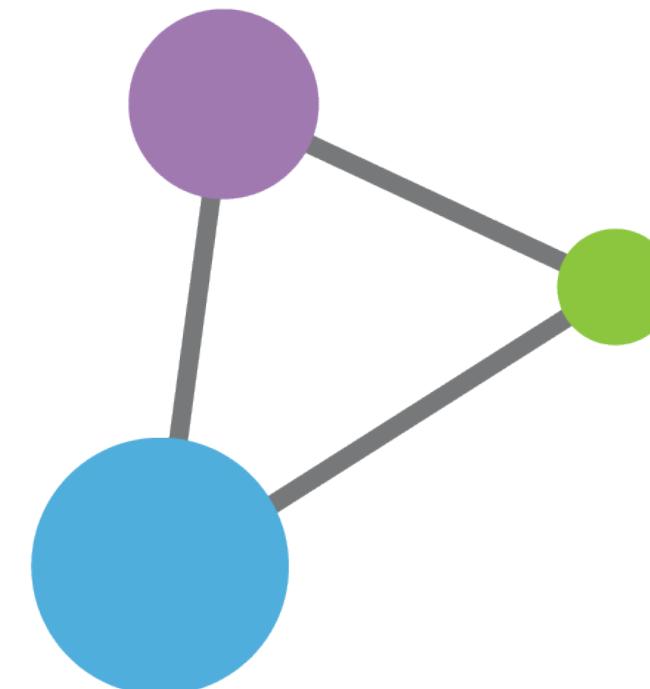
Centralized Genomic Knowledge Bases



Data Commons
Trusted, controlled repository of multiple datasets

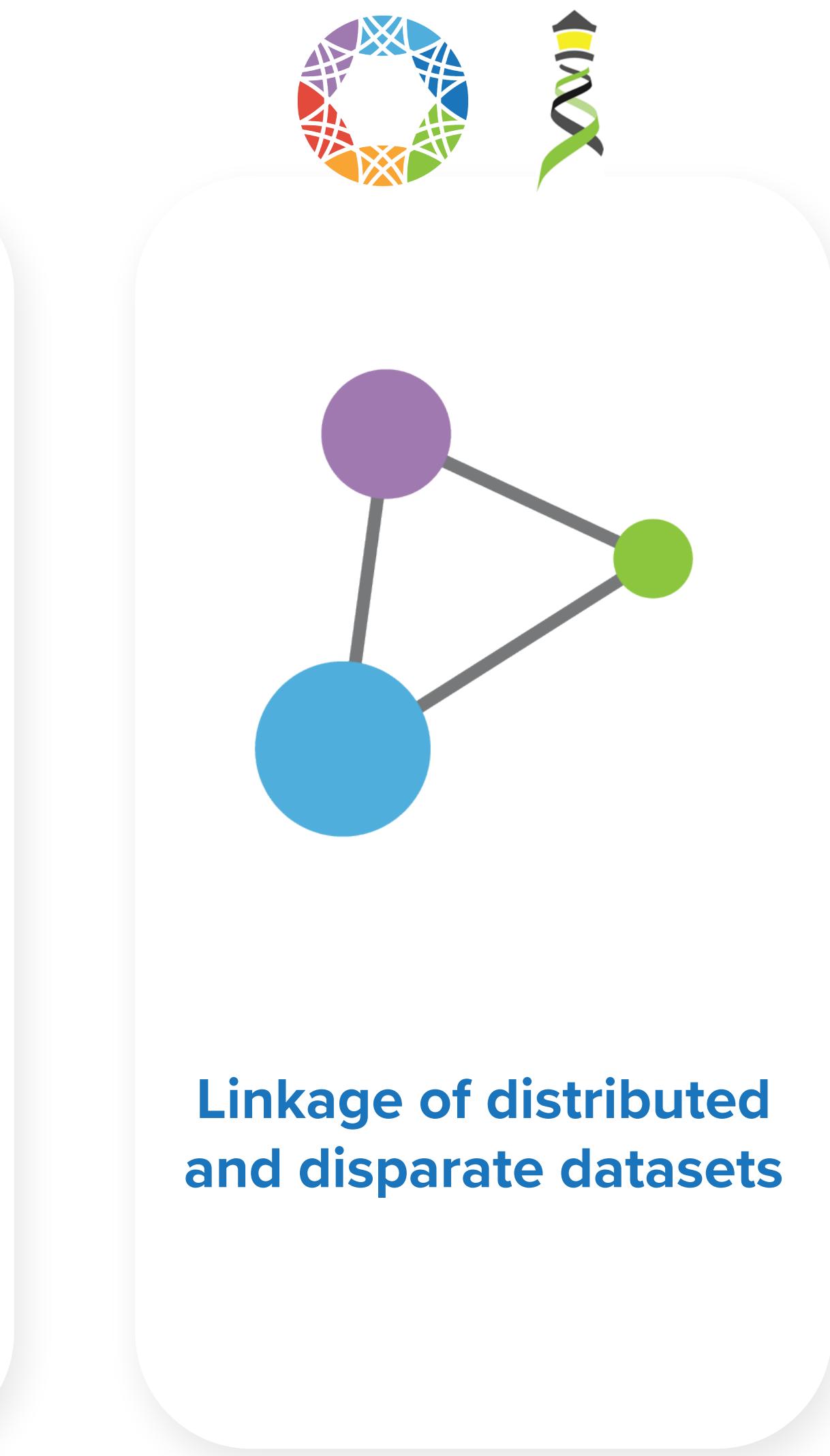
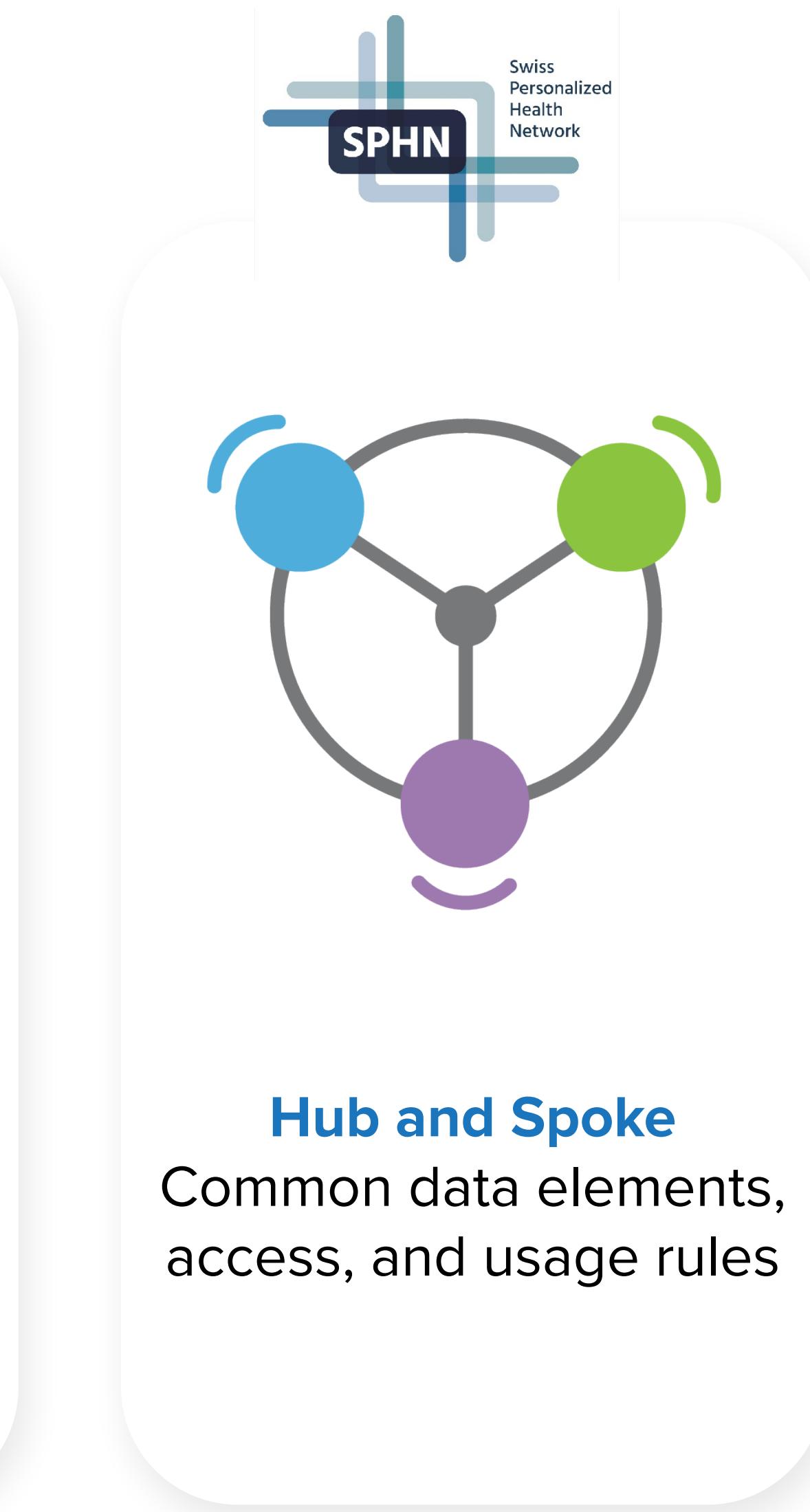
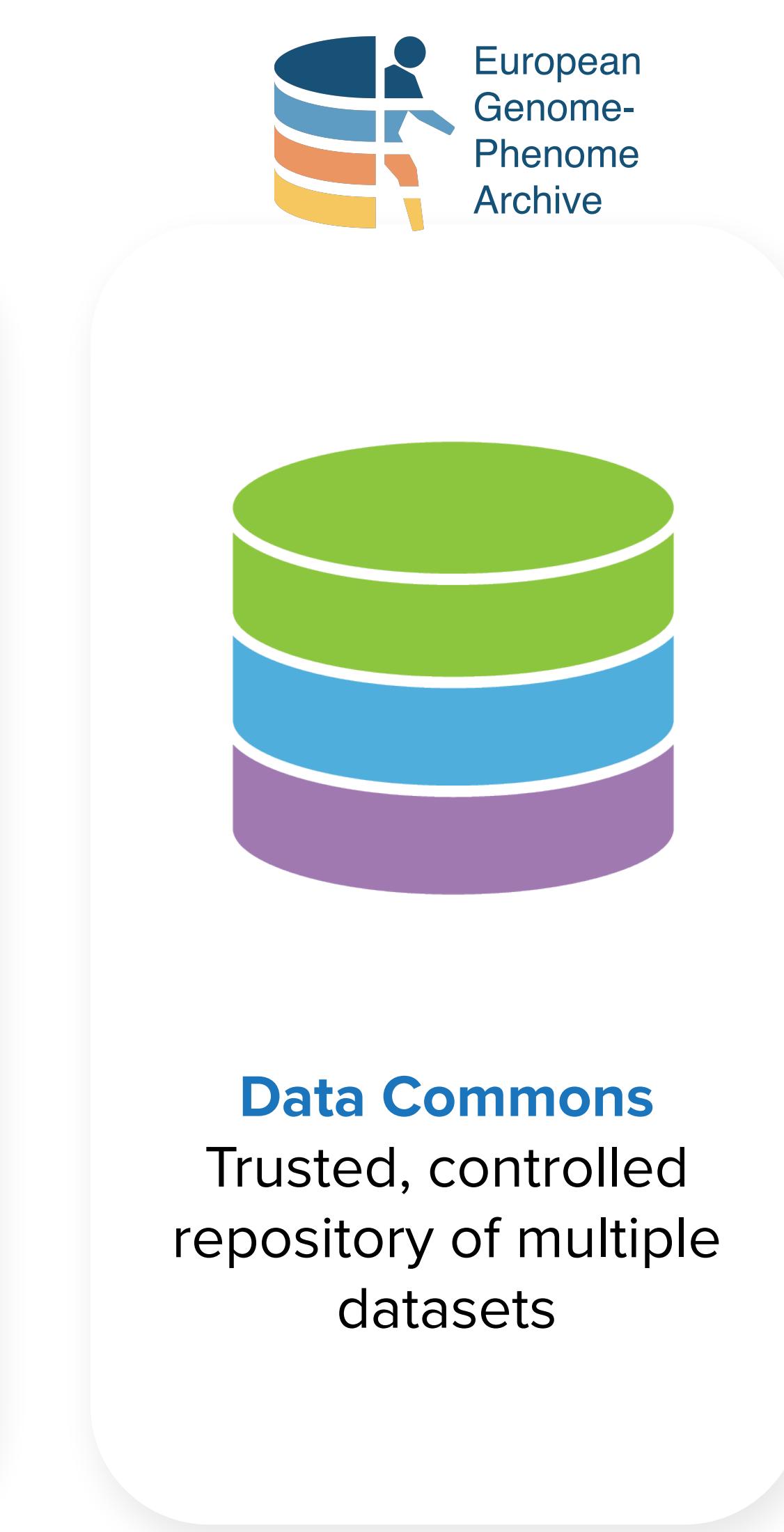


Hub and Spoke
Common data elements, access, and usage rules



Linkage of distributed and disparate datasets

Different Approaches to Data Sharing



Cancer Genomics Reference Resource

- **open** resource for oncogenomic profiles
- over **116'000 cancer CNV profiles**
- more than **800 diagnostic types**
- inclusion of reference datasets (e.g. TCGA)
- standardized encodings (e.g. NCIt, ICD-O 3)
- identifier mapping for PMID, GEO, Cellosaurus, TCGA, cBioPortal where appropriate
- core clinical data (TNM, sex, survival ...)
- data mapping services
- recent addition of SNV data for some series



Cancer CNV Profiles

ICD-O Morphologies
ICD-O Organ Sites
Cancer Cell Lines
Clinical Categories

Search Samples

arrayMap
TCGA Samples
1000 Genomes
Reference Samples
DIPG Samples
cBioPortal Studies
Gao & Baudis, 2021

Publication DB

Genome Profiling
Progenetix Use

Services

NCIt Mappings
UBERON Mappings

Upload & Plot

Beacon⁺

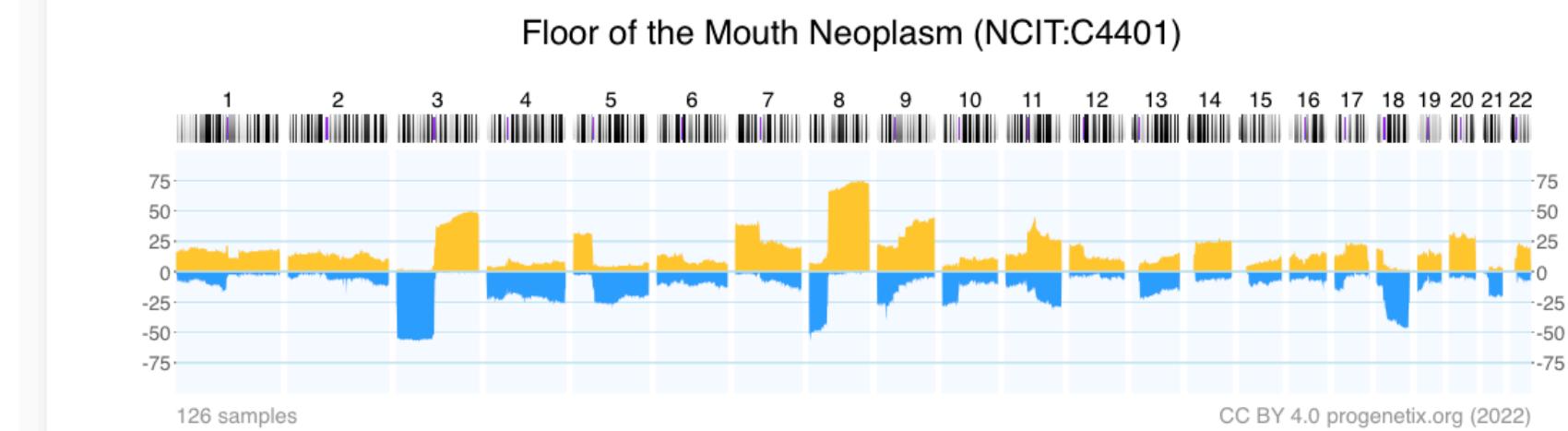
Documentation

News
Downloads & Use
Cases
Sevices & API

Baudisgroup @ UZH

Cancer genome data @ progenetix.org

The Progenetix database provides an overview of mutation data in cancer, with a focus on copy number abnormalities (CNV / CNA), for all types of human malignancies. The data is based on *individual sample data* from currently **142063** samples.



[Download SVG](#) | [Go to NCIT:C4401](#) | [Download CNV Frequencies](#)

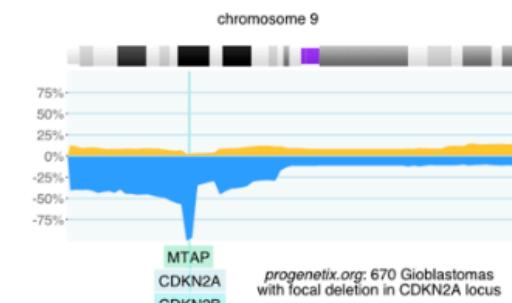
Example for aggregated CNV data in 126 samples in Floor of the Mouth Neoplasm.

Here the frequency of regional **copy number gains** and **losses** are displayed for all 22 autosomes.

Progenetix Use Cases

Local CNV Frequencies

A typical use case on Progenetix is the search for local copy number aberrations - e.g. involving a gene - and the exploration of cancer types with these CNVs. The [\[Search Page \]](#) provides example use cases for designing queries. Results contain basic statistics as well as visualization and download options.



Cancer CNV Profiles

The progenetix resource contains data of **834** different cancer types (NCIt neoplasm classification), mapped to a variety of biological and technical categories. Frequency profiles of regional genomic gains and losses for all categories (diagnostic entity, publication, cohort ...) can be accessed through the [\[Cancer Types \]](#) page with direct visualization and options for sample retrieval and plotting options.

Cancer Genomics Publications

Through the [\[Publications \]](#) page Progenetix provides **4164** annotated references to research articles from cancer genome screening experiments (WGS, WES, aCGH, cCGH). The numbers of analyzed samples and possible availability in the Progenetix sample collection are indicated.

Different Approaches to Data Sharing



Centralized Genomic Knowledge Bases



Hub and Spoke
Common data elements, access, and usage rules



Linkage of distributed and disparate datasets

The EGA



Long term secure archive for human biomedical research sensitive data, with focus on reuse of the data for further research (or “*broad and responsible use of genomic data*”)



The EGA



- EGA “owns” nothing; data controllers tell who is authorized to access ***their*** datasets
- EGA admins provide smooth “all or nothing” data sharing process

A screenshot of the EGA DAC interface. It shows two main sections: 'Requests' and 'History'.

My DACs - EGAC5000000005 - Requests

EuCanImage DAC
This is a DAC for EuCanImage data

Type something for filter the requests...

My DACs - EGAC5000000005 - History

EuCanImage DAC
This is a DAC for EuCanImage data

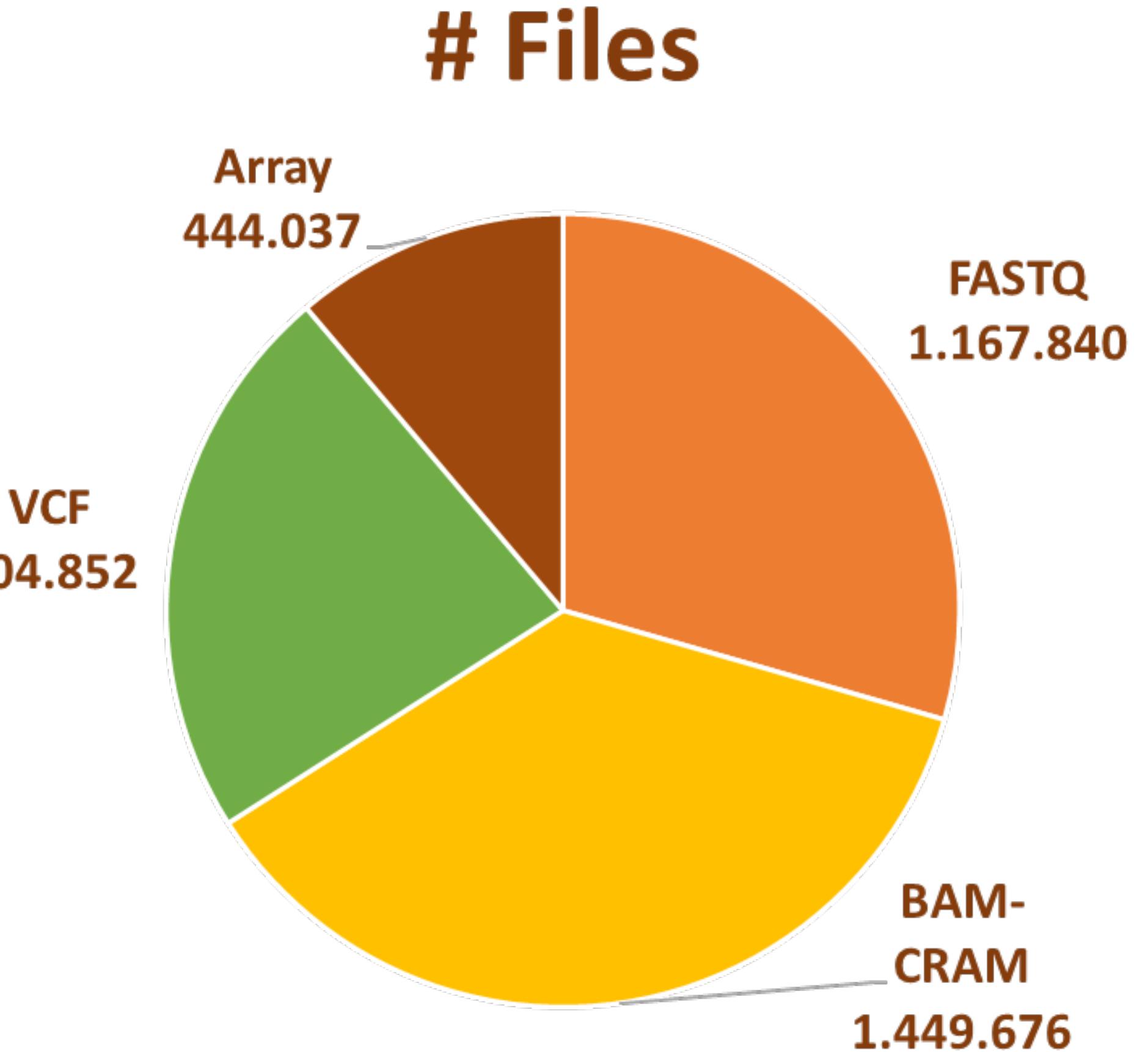
EGAD5000000032: EuCanImage

Dr Teresa Garcia Lezana teresa

Start typing user's name, e-mail or Start typing dataset ID or title... Select a date or a range... Filter

Date	Requester	Dataset	DAC Admin/Member
18 August 2022	gemma.milla@crg.eu	EGAD5000000032	Dr Lauren A Fromont
17 August 2022	Dr Teresa Garcia Lezana	EGAD5000000033	Dr Teresa Garcia Lezana
16 August 2022	Dr Teresa Garcia Lezana	EGAD5000000032	Dr Lauren A Fromont

revoke permission APPLY



4,328 Studies released
10,470 Datasets
2,309 Data Access Committees

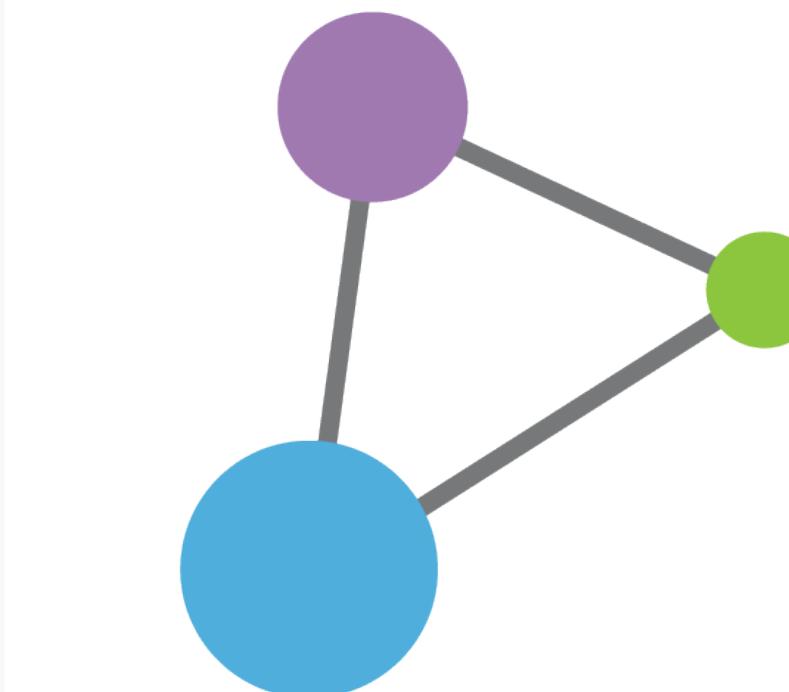
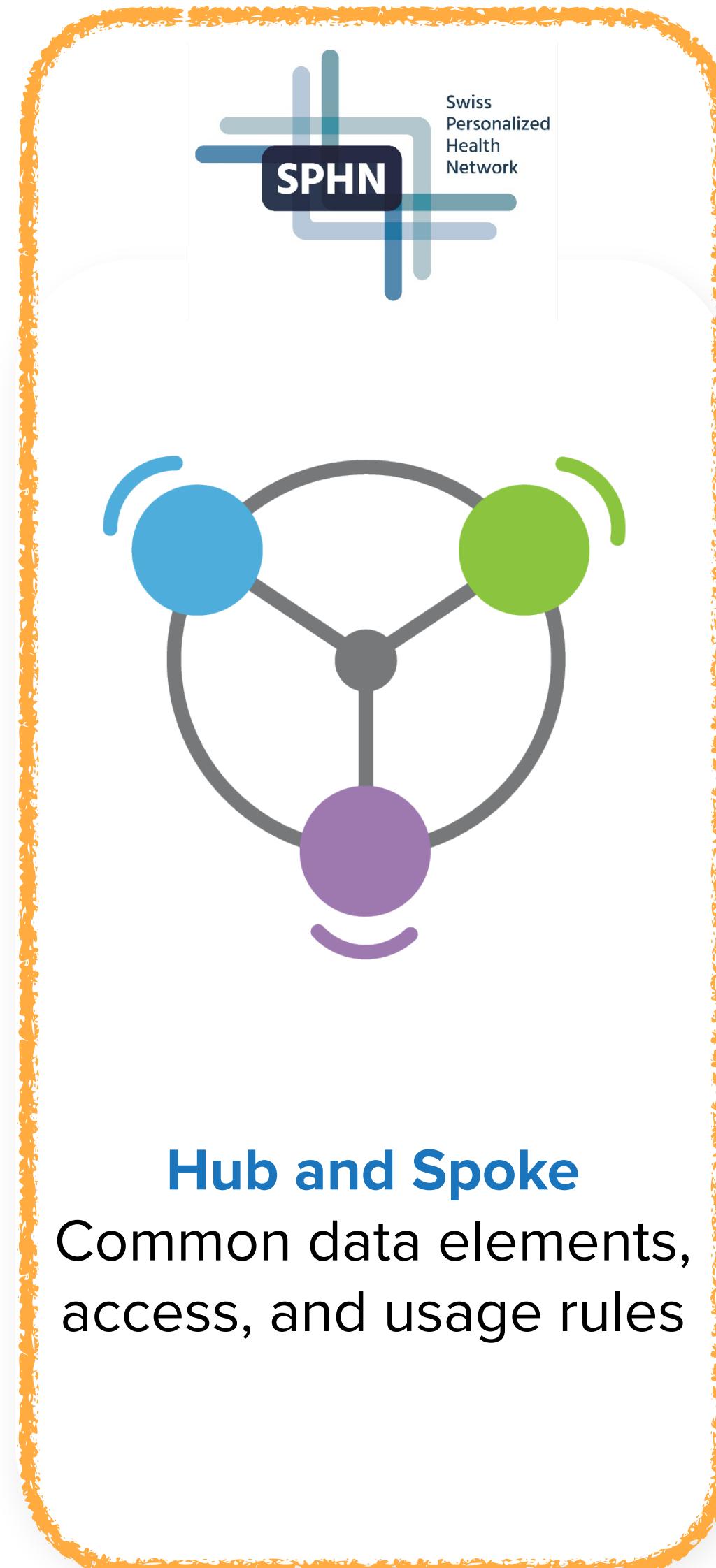
Different Approaches to Data Sharing



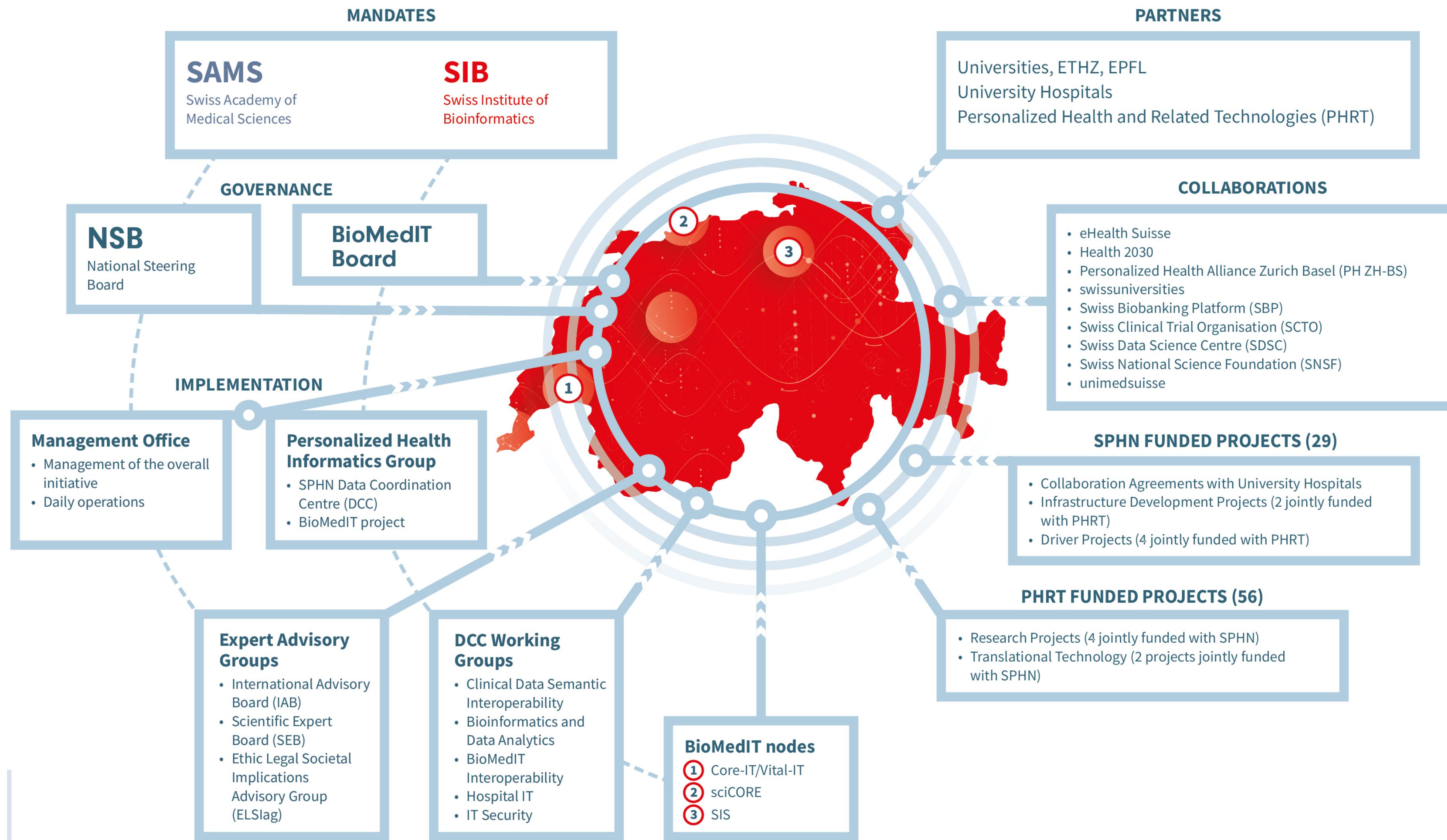
Centralized Genomic Knowledge Bases



Data Commons
Trusted, controlled repository of multiple datasets



SPHN: a federal mandate



Different Approaches to Data Sharing



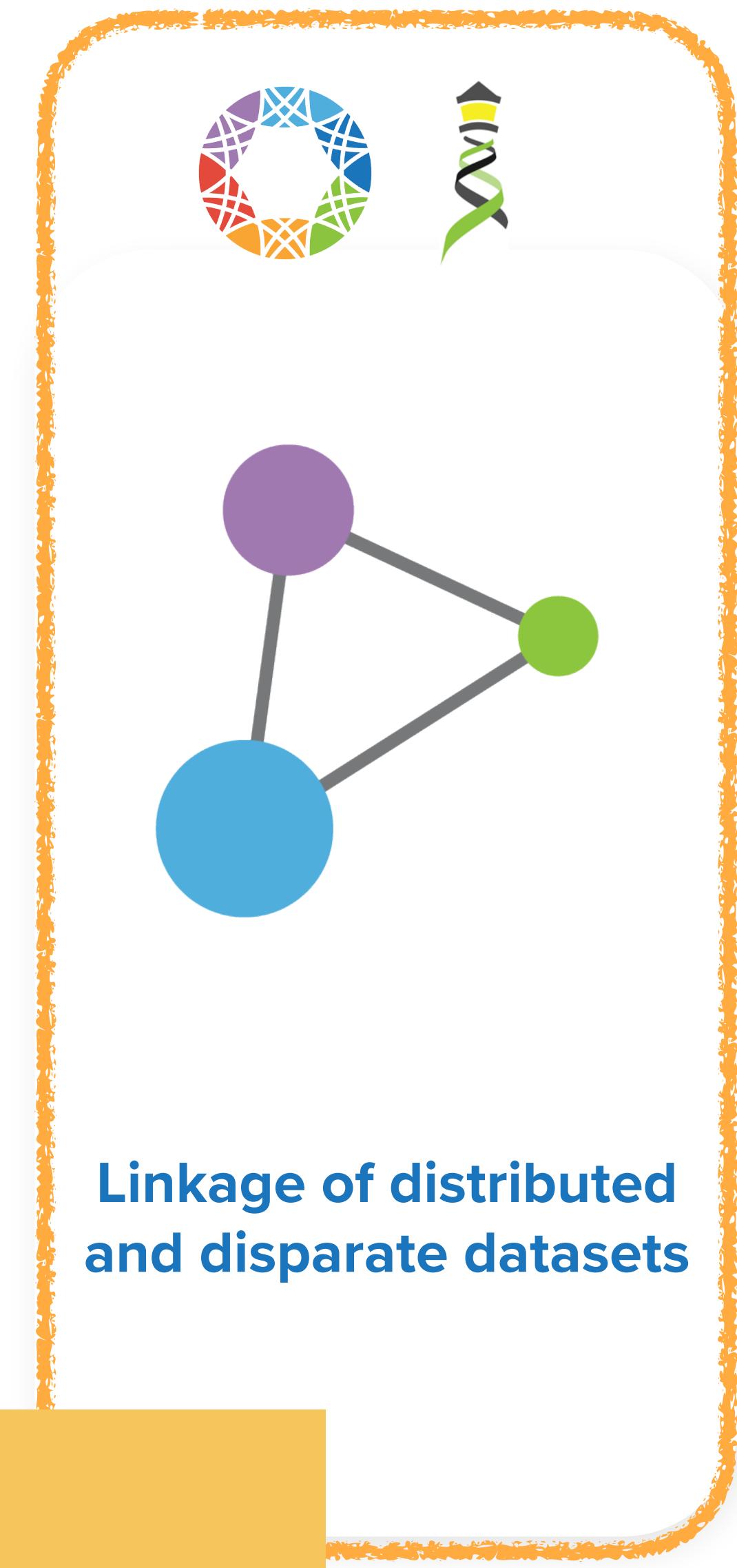
Centralized Genomic Knowledge Bases



Data Commons
Trusted, controlled repository of multiple datasets



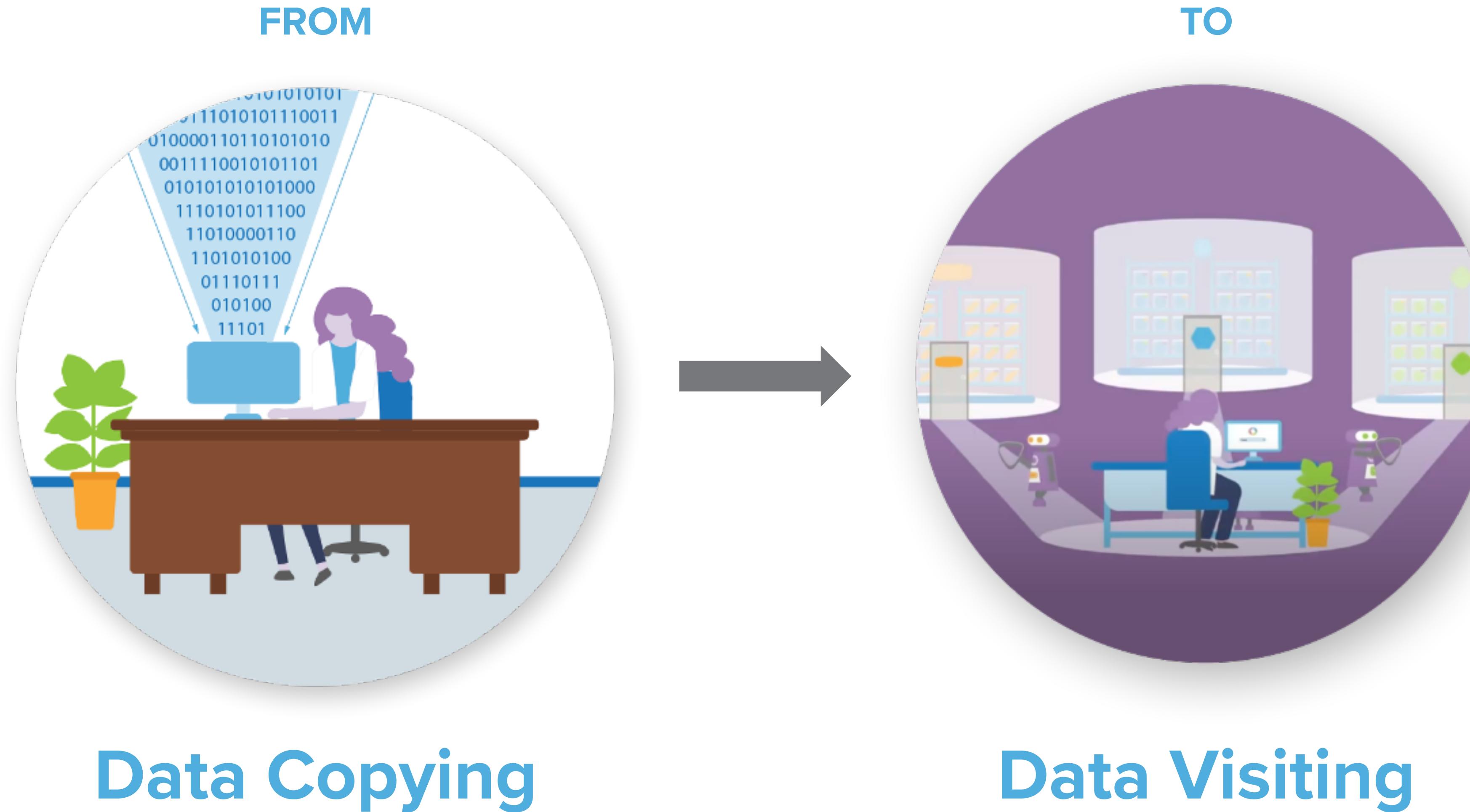
Hub and Spoke
Common data elements, access, and usage rules



Linkage of distributed and disparate datasets

Federation

A New Paradigm for Data Sharing



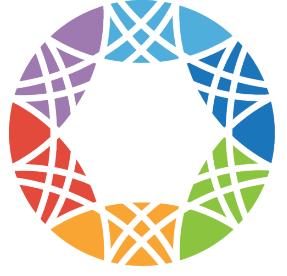
Federation



“A grouping of autonomous organizations and datasets with a centralised control”

Federation allows us to....

1. **Move analysis to data**, not aggregate data close to each researcher
2. **Have broad, reciprocal data access methods** which respect national processes and patient consent
3. Transfer methods and skills into the **healthcare sector**
4. **Leverage healthcare data** to make more discoveries on humans

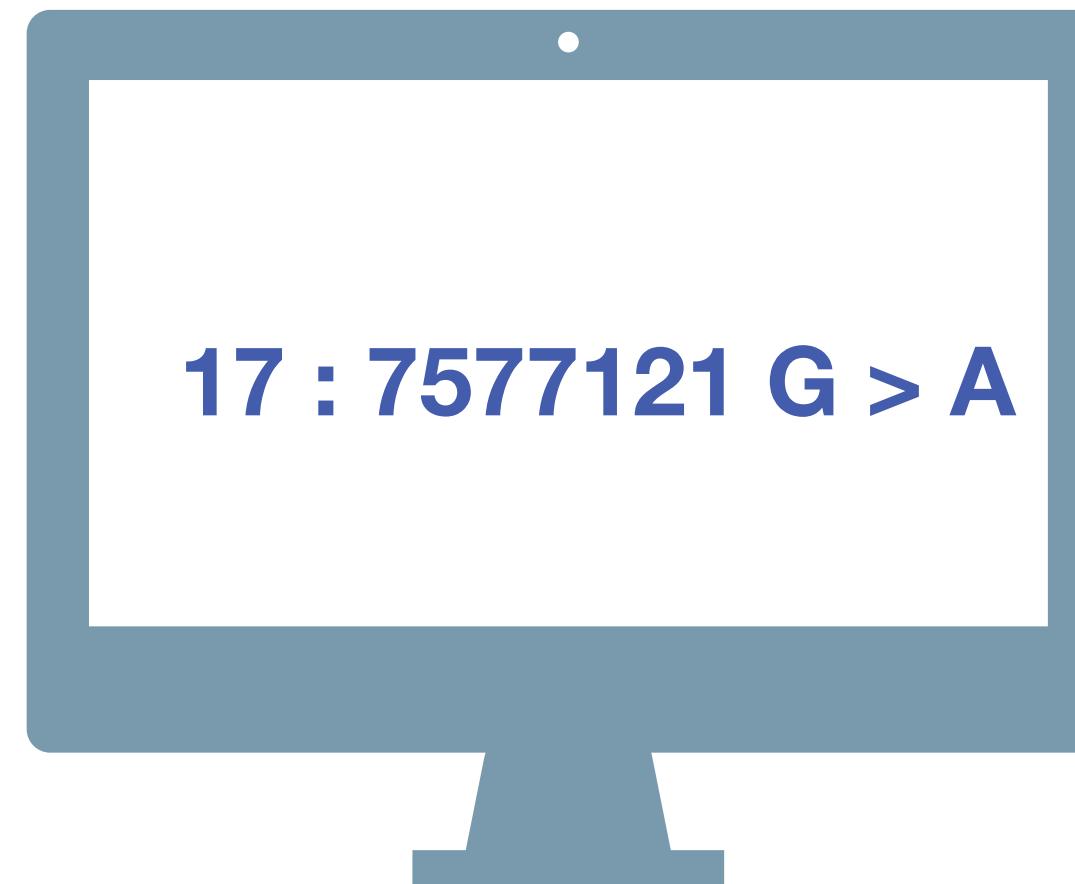


Global Alliance
for Genomics & Health
Collaborate. Innovate. Accelerate.



The GA4GH Beacon Protocol

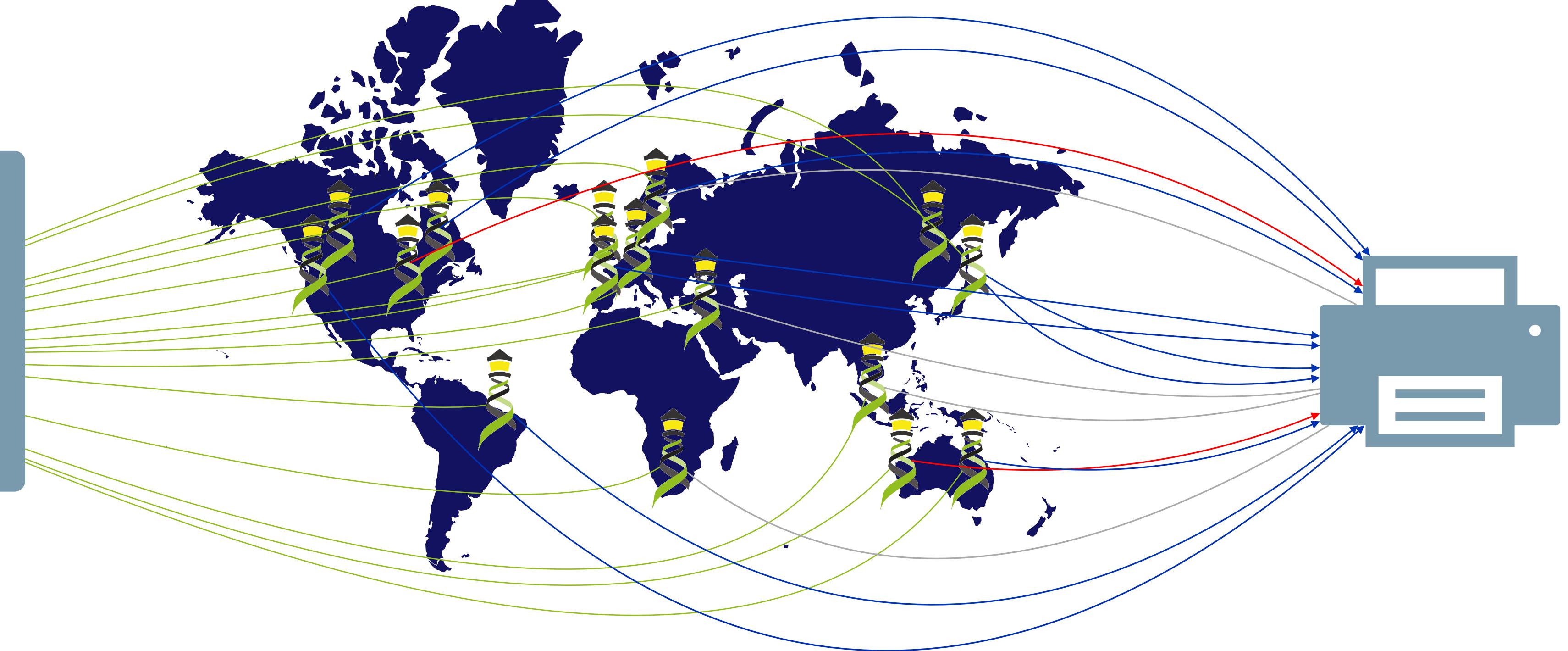
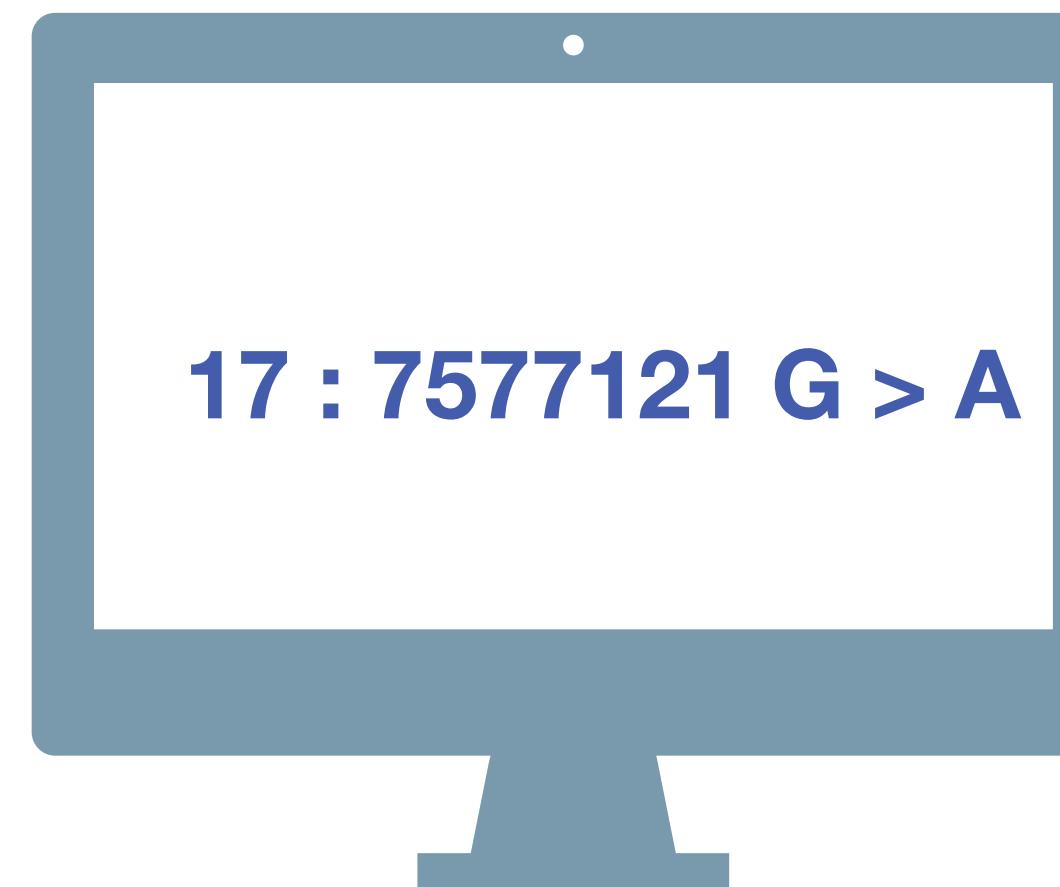
Federating Genomic Discoveries



Beacon

A **Beacon** answers a query for a specific genome variant against individual or aggregate genome collections

YES | NO | \0



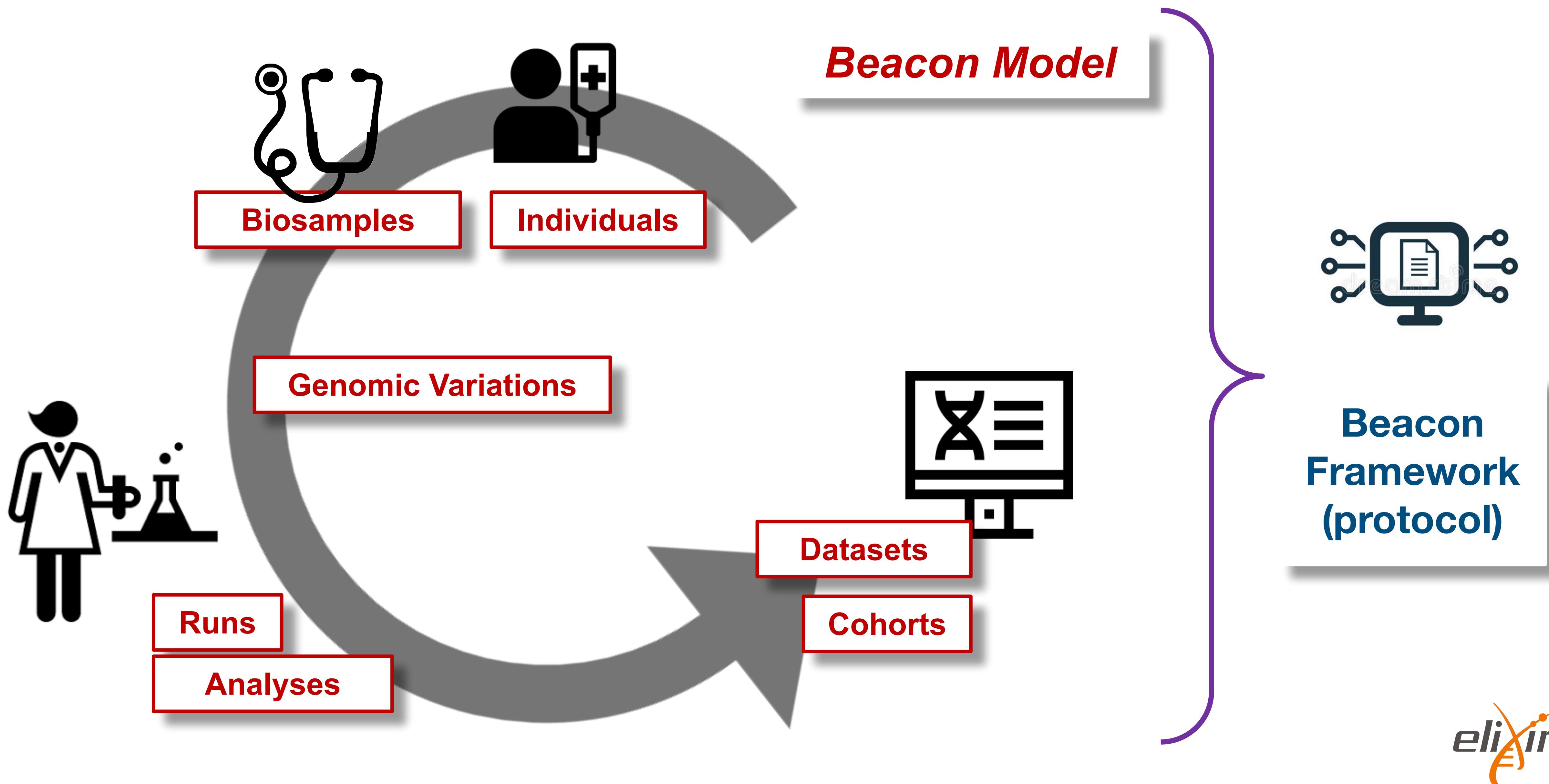
Have you seen this variant?
It came up in my patient
and we don't know if this is
a common SNP or worth
following up.

A Beacon network federates
genome variant queries
across databases that
support the **Beacon API**

Here: The variant has
been found in **few**
resources, and those
are from **disease**
specific **collections**.

Beacon v2

docs.genomebeacons.org



Progenetix & Beacon

Implementation driven standards development

- Progenetix Beacon+ has served as implementation driver since 2016
- prototyping of advanced Beacon features such as
 - structural variant queries
 - data handovers
 - Phenopackets integration

Beacon v2 GA4GH Approval Registry

Beacons: European Genome-Phenome Archive | progenetix | cnag | University of Leicester

Theoretical Cytogenetics and Oncogenomics group at UZH and SIB
Progenetix Cancer Genomics Beacon+ provides a forward looking implementation of the Beacon v2 API, with focus on structural genome variants and metadata based on the...

Visit us | Beacon UI | Beacon API | Contact us

European Genome-Phenome Archive (EGA)
GA4GH Approval Beacon Test
This Beacon is based on the GA4GH Beacon v2.0

BeaconMap
Bioinformatics analysis
Biological Sample
Cohort
Configuration
Dataset
EntryTypes
Genomic Variants
Individual
Info
Sequencing run

Visit us | Beacon UI | Beacon API | Contact us

BeaconMap
Bioinformatics analysis
Biological Sample
Cohort
Configuration
Dataset
EntryTypes
Genomic Variants
Individual
Info
Sequencing run

Visit us | Beacon UI | Beacon API | Contact us

Centre Nacional Analisis Genomica (CNAG-CRG)
Beacon @ RD-Connect
This Beacon is based on the GA4GH Beacon v2.0

BeaconMap
Bioinformatics analysis
Biological Sample
Cohort
Configuration
Dataset
EntryTypes
Genomic Variants
Individual
Info
Sequencing run

Visit us | Beacon UI | Beacon API | Contact us

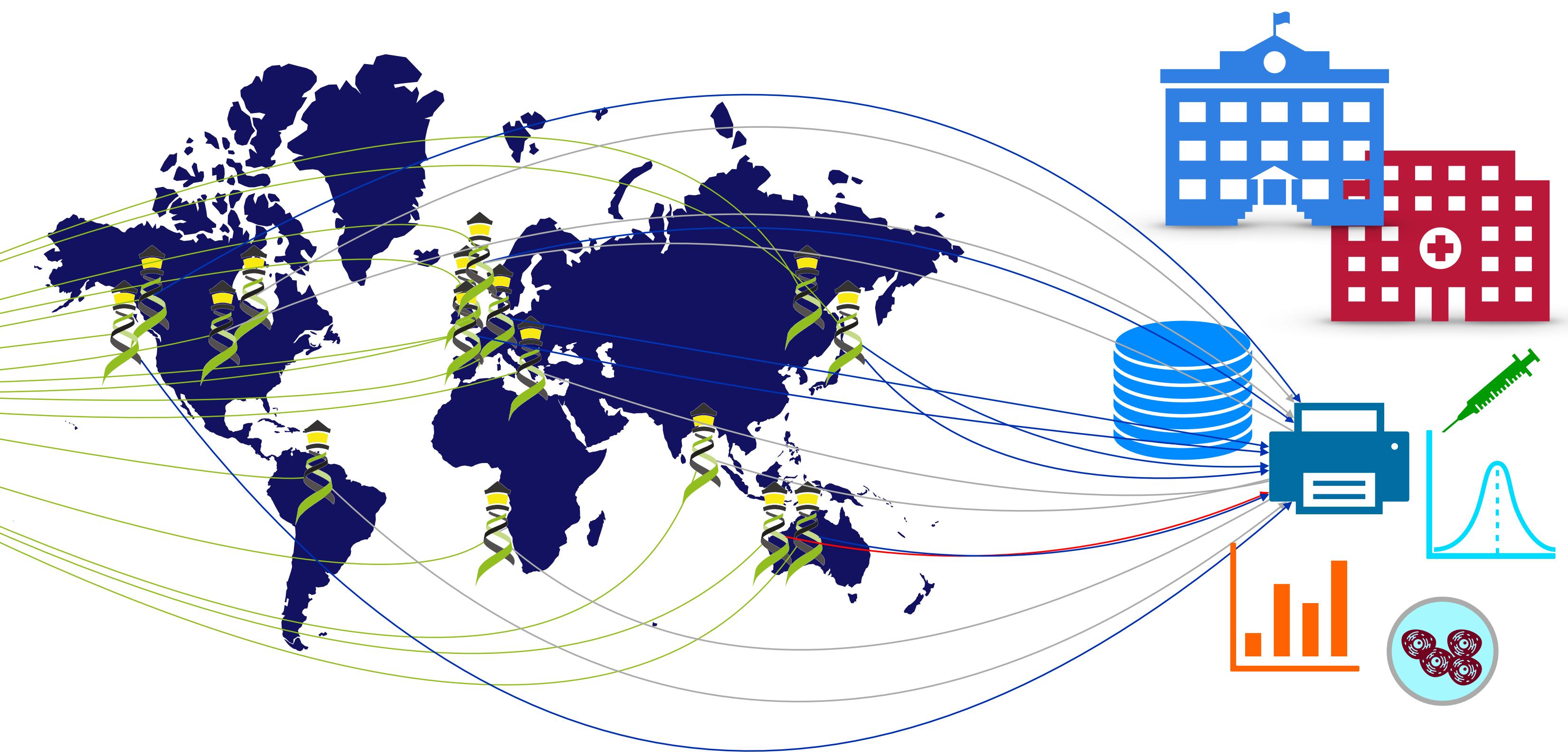
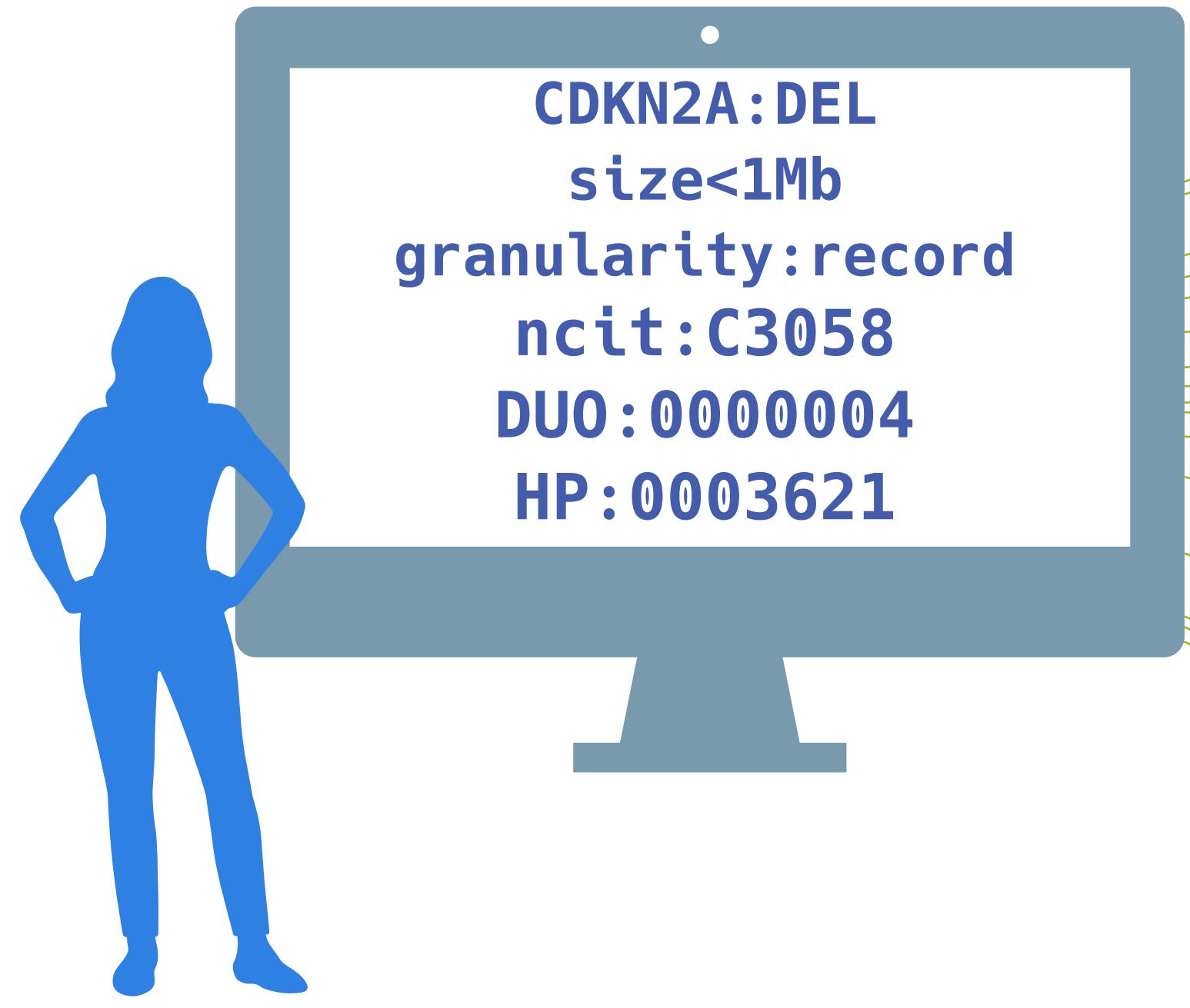
University of Leicester
Cafe Variome Beacon v2
This Beacon is based on the GA4GH Beacon v2.0

BeaconMap
Bioinformatics analysis
Biological Sample
Cohort
Configuration
Dataset
EntryTypes
Genomic Variants
Individual
Info
Sequencing run

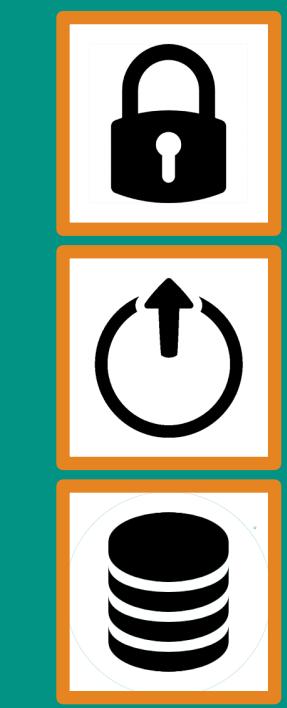
Visit us | Beacon UI | Beacon API | Contact us

Beacon protocol response verifier at time of GA4GH approval Spring 2022

Matches the Spec | Not Match the Spec | Not Implemented



Can you provide data about focal deletions in CDKN2A in Glioblastomas from juvenile patients with unrestricted access?

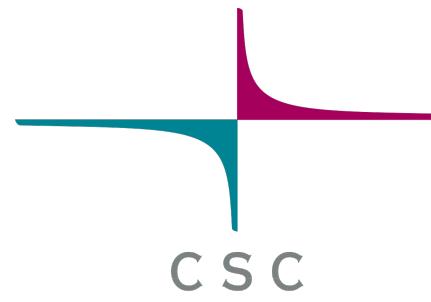


Beacon **v2** API

The Beacon API v2 represents a simple but powerful **genomics API** for **federated** data discovery and retrieval



Jordi Rambla
Arcadi Navarro
Roberto Ariosa
Manuel Rueda
Lauren Fromont
Mauricio Moldes
Claudia Vasallo
Babita Singh
Sabela de la Torre
Marta Ferri
Fred Haziza



Juha Törnroos
Teemu Kataja
Ikkka Lappalainen
Dylan Spalding



Tony Brookes

Tim Beck

Colin Veal

Tom Shorter



Michael Baudis

Rahel Paloots

Hangjia Zhao

Ziying Yang

Bo Gao



Augusto Rendon

Ignacio Medina

Javier López

Jacobo Coll

Antonio Rueda



centre nacional d'anàlisi genòmica
centro nacional de análisis genómico

Sergi Beltran

Carles Hernandez



Institut national
de la santé et de la recherche médicale

David Salgado



Salvador Capella

Dmitry Repchevski

JM Fernández



Laura Furlong

Janet Piñero



Serena Scollen

Gary Saunders

Giselle Kerry

David Lloyd



Nicola Mulder

Mamana

Mbiyavanga

Ziyaad Parker



David

Torrents



Dean Hartley



Fundación Progreso y Salud
CONSEJERÍA DE SALUD

Joaquin Dopazo

Javier Pérez

J.L. Fernández

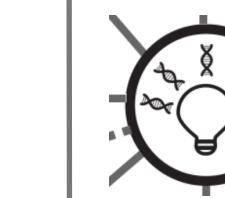
Gema Roldan



Thomas Keane

Melanie Courtot

Jonathan Dursi

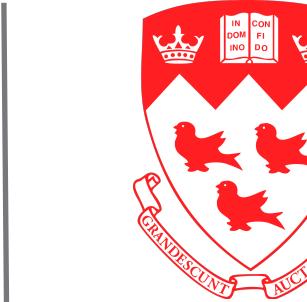


Heidi Rehm

Ben Hutton



Toshiaki Katayama



Stephane Dyke



Marc Fiume

Miro Cupak



Melissa Cline



Diana Lemos



GA4GH Phenopackets

Peter Robinson
Jules Jacobsen



GA4GH VRS

Alex Wagner
Reece Hart

Beacon PRC

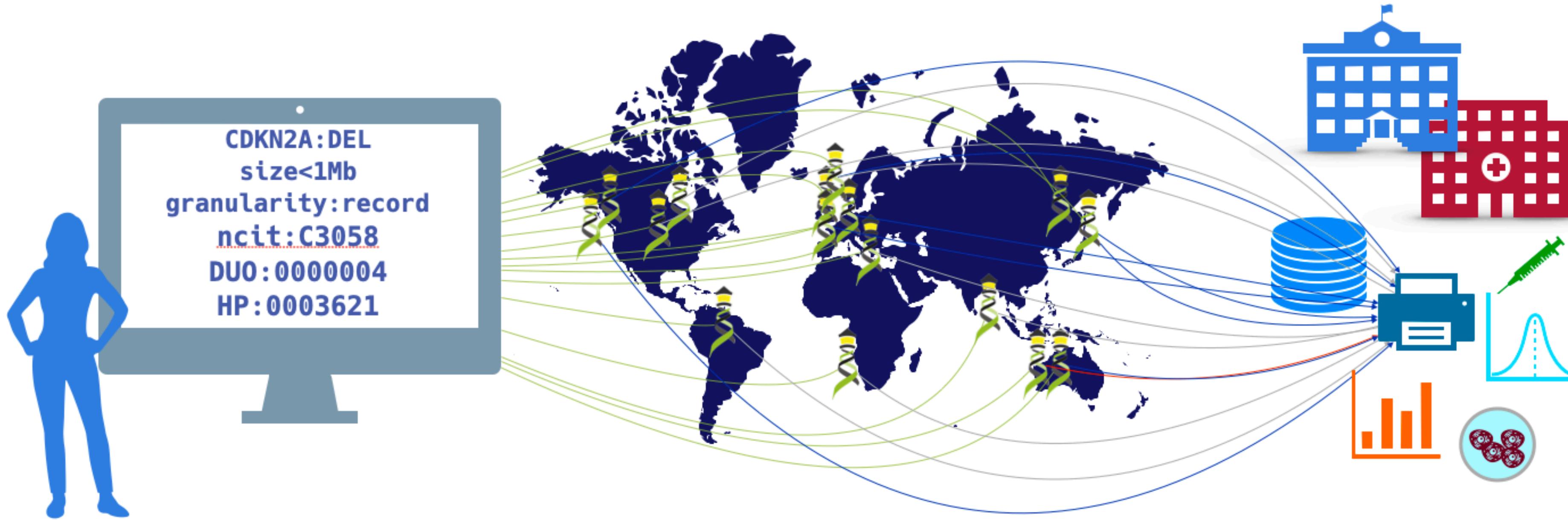
Alex Wagner
Jonathan Dursi
Mamana Mbiyavanga
Alice Mann
Neerjah Skantharajah



The Beacon team through the ages



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



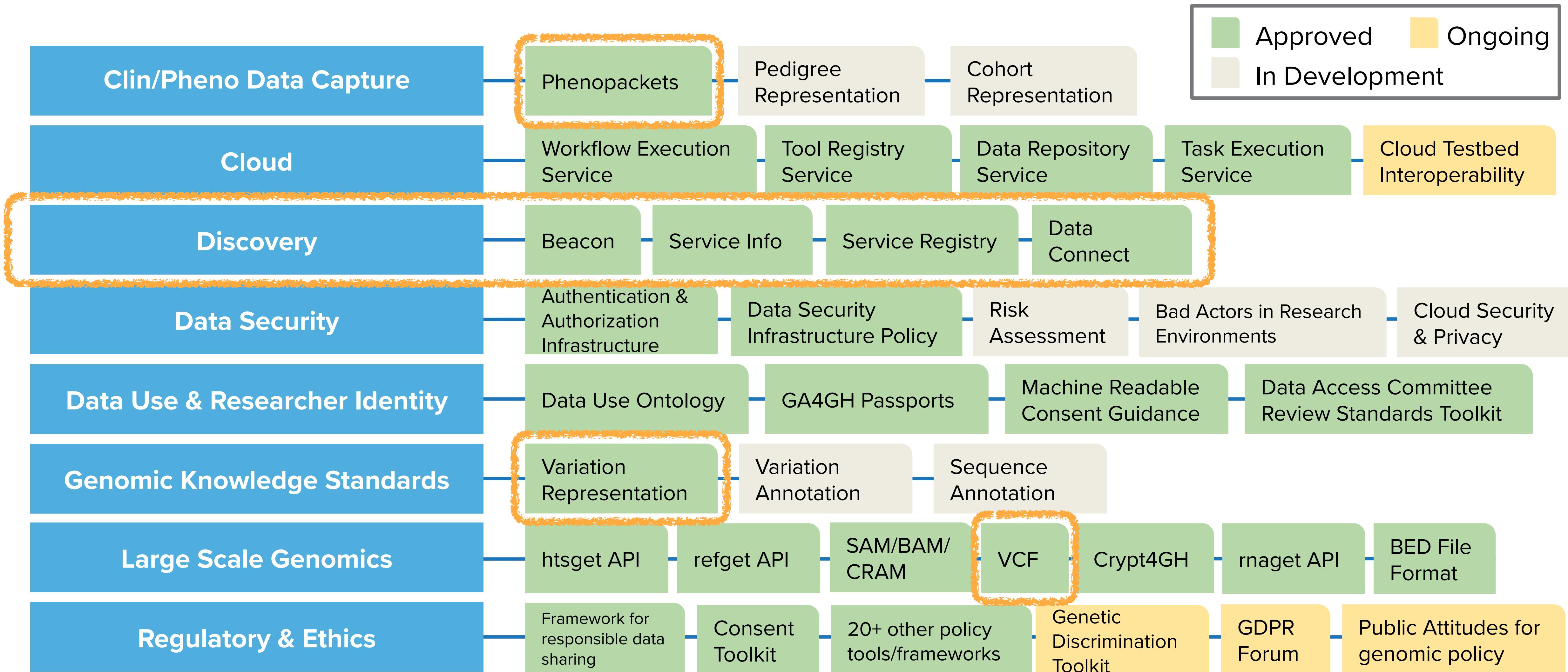
The GA4GH Beacon Protocol

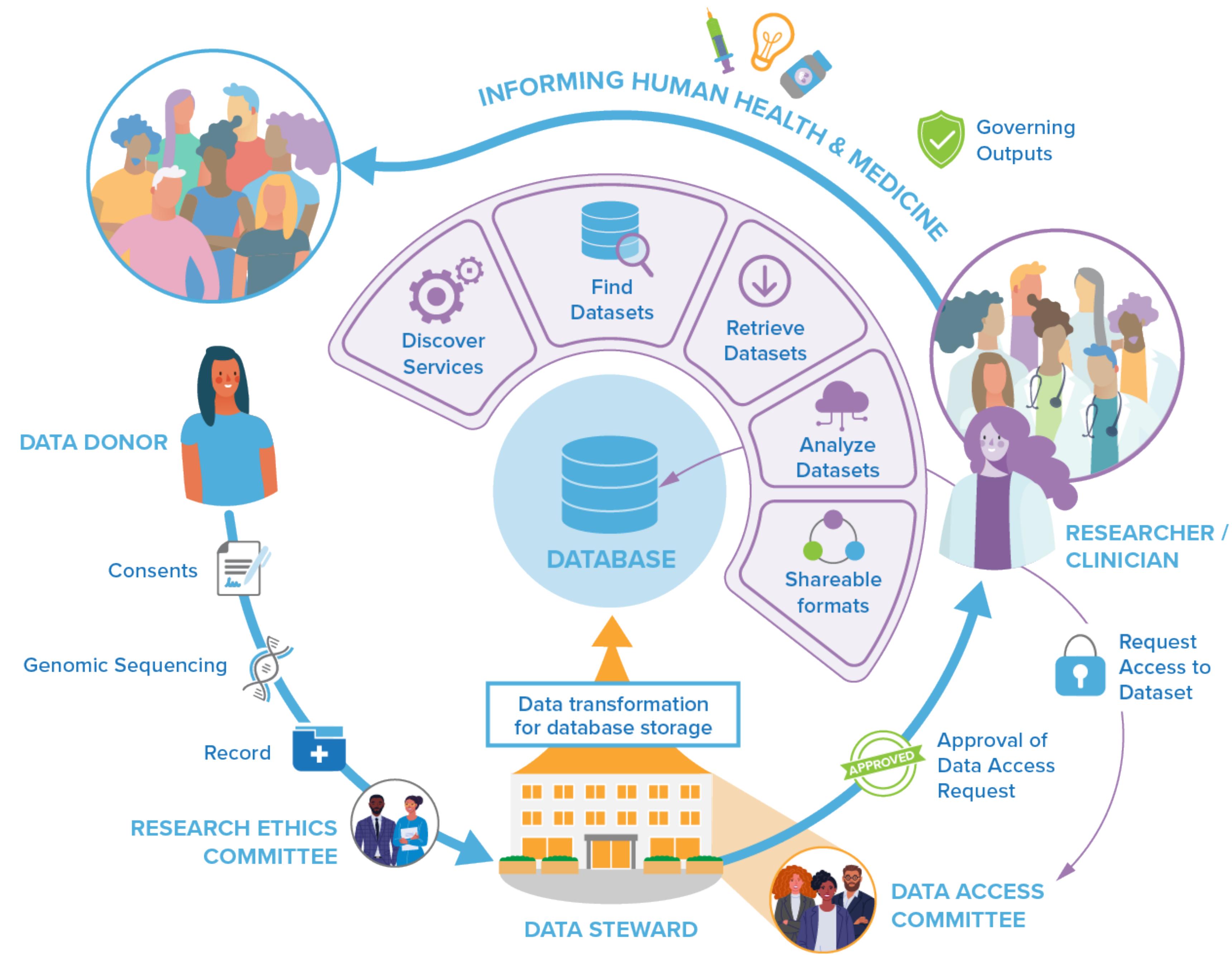
... is only one part of the GA4GH ecosystem

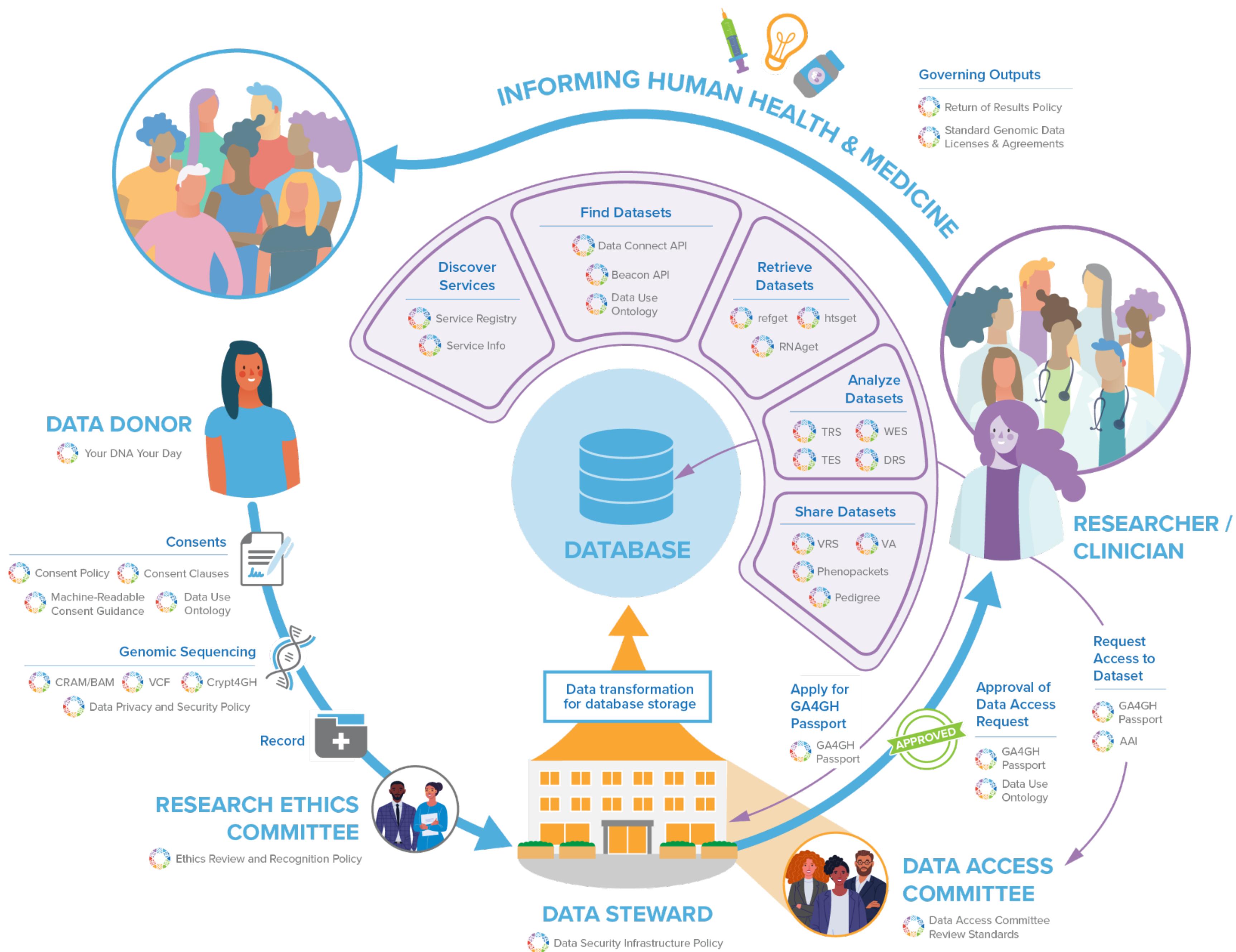
GA4GH ecosystem and outputs



Overview of GA4GH standards and frameworks







Global Collaborations



International HundredK+
Cohorts Consortium



Global Genomic
Medicine Consortium



Beyond One Million
Genomes Project



Maps to
Mechanisms to
Medicine



WORLD
ECONOMIC
FORUM



THE MEDICAL GENOME
INITIATIVE



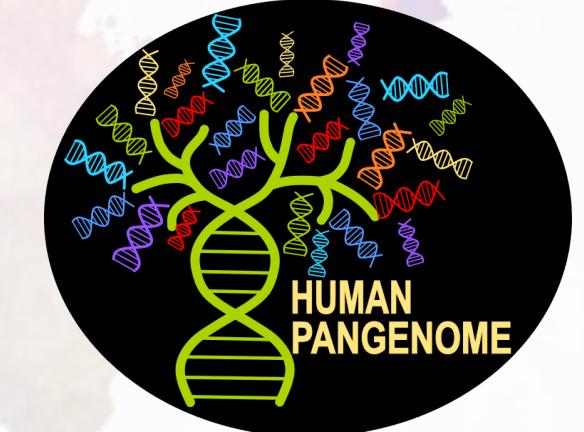
World Health
Organization
Science Council



Public Health Alliance for
Genomic Epidemiology



Global Alliance
for Genomics & Health



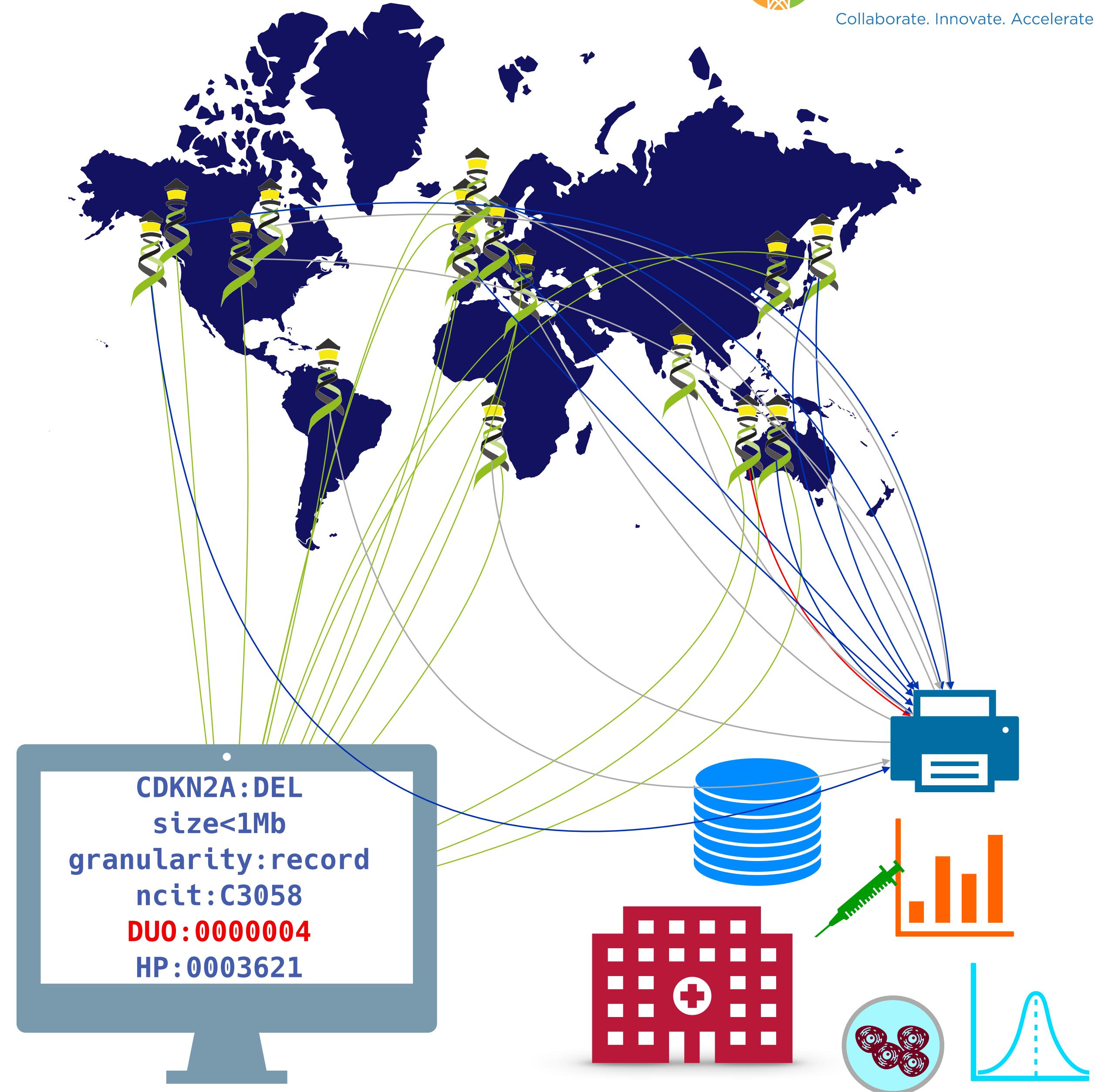
Human Pangenome
Project (HPP) / Human
Pangenome Reference
Consortium (HPRC)

What Can You Do?

- implement procedures and standards supporting **data discovery** (FAIR principles) and federation approaches
- forward looking consent and data protection models adhering to **ORD** principles ("as secure as necessary, as open as possible")
- get involved with international **data standards** efforts and projects



Collaborate!

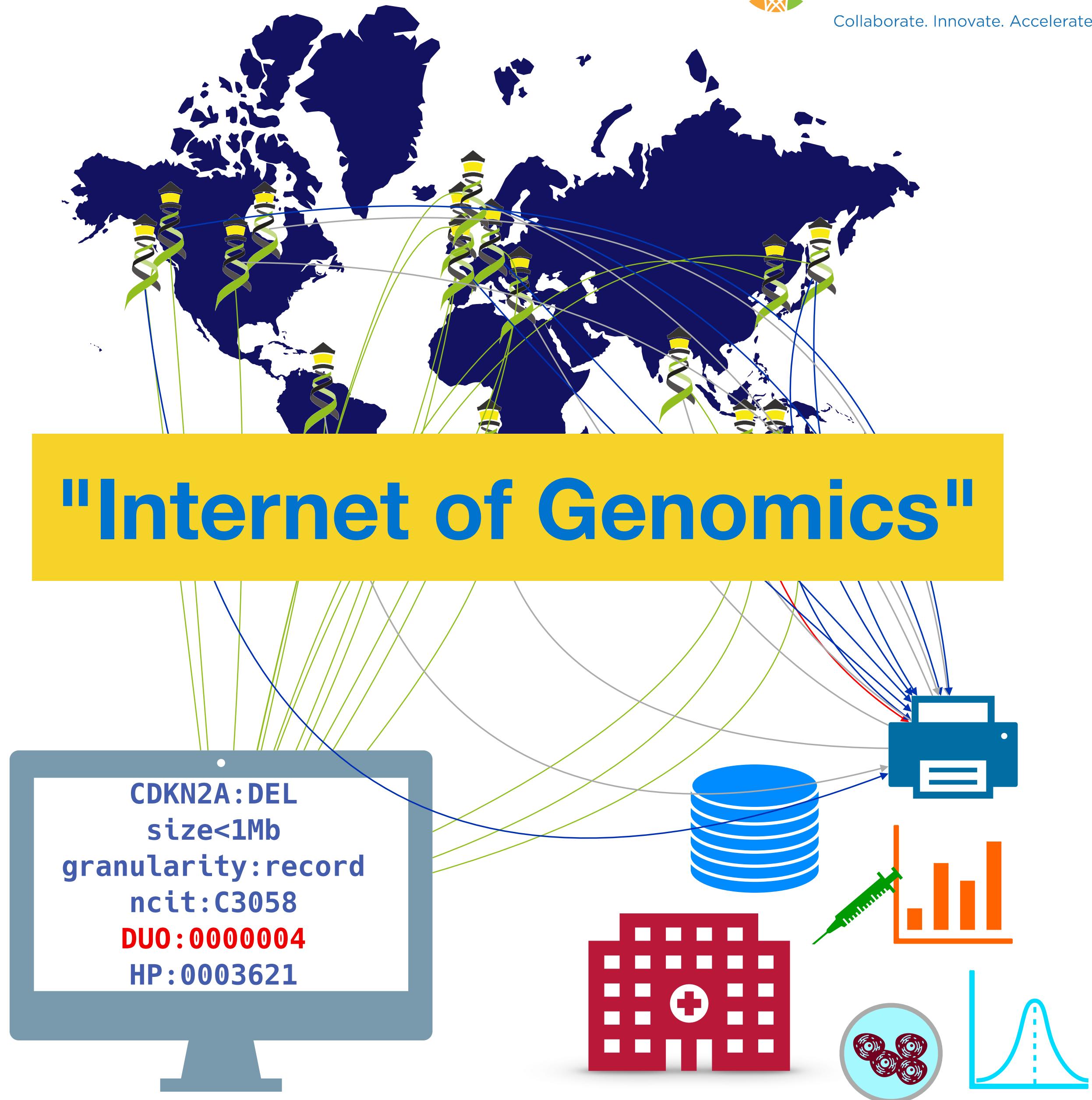


What Can You Do?

- implement procedures and standards supporting **data discovery** (FAIR principles) and federation approaches
- forward looking consent and data protection models adhering to **ORD** principles ("as secure as necessary, as open as possible")
- get involved with international **data standards** efforts and projects



Collaborate!





Universität
Zürich^{UZH}



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



Swiss Institute of
Bioinformatics



