# Genomic Data Sharing Standard Development with GA4GH and ELIXIR

## *Opportunities and Pitfalls in Federated Data Discovery*

Michael Baudis @ DMLS Lecture Series 2024-02-27

# Theoretical Cytogenetics and Oncogenomics

**Cancer Genomics | Data Resources | Methods & Standards for Genomics and Personalized Health**

Curators

**Data Parasites**

# Bioinformatics & Bioinformaticians are ...

| **Bio**informatician | Bio**informatician** |
|---|---|
| strong biological knowledge | sufficient biological background |
| provides hypothesis and / or dataset | provides statistical, analysis methods |
| **sufficient statistical** and **computational** expertise to correctly use bioinformatics tools & develop workflows (scripting ...) | **sufficient biological** or **medical** background to understand problems presented and identify pitfalls and hidden biases arising from data generation |
| expert **user** of informatics tools | **developer** of informatics tools |
| may get a Nobel | may get rich |

# Bioinformatics & Bioinformaticians are ...

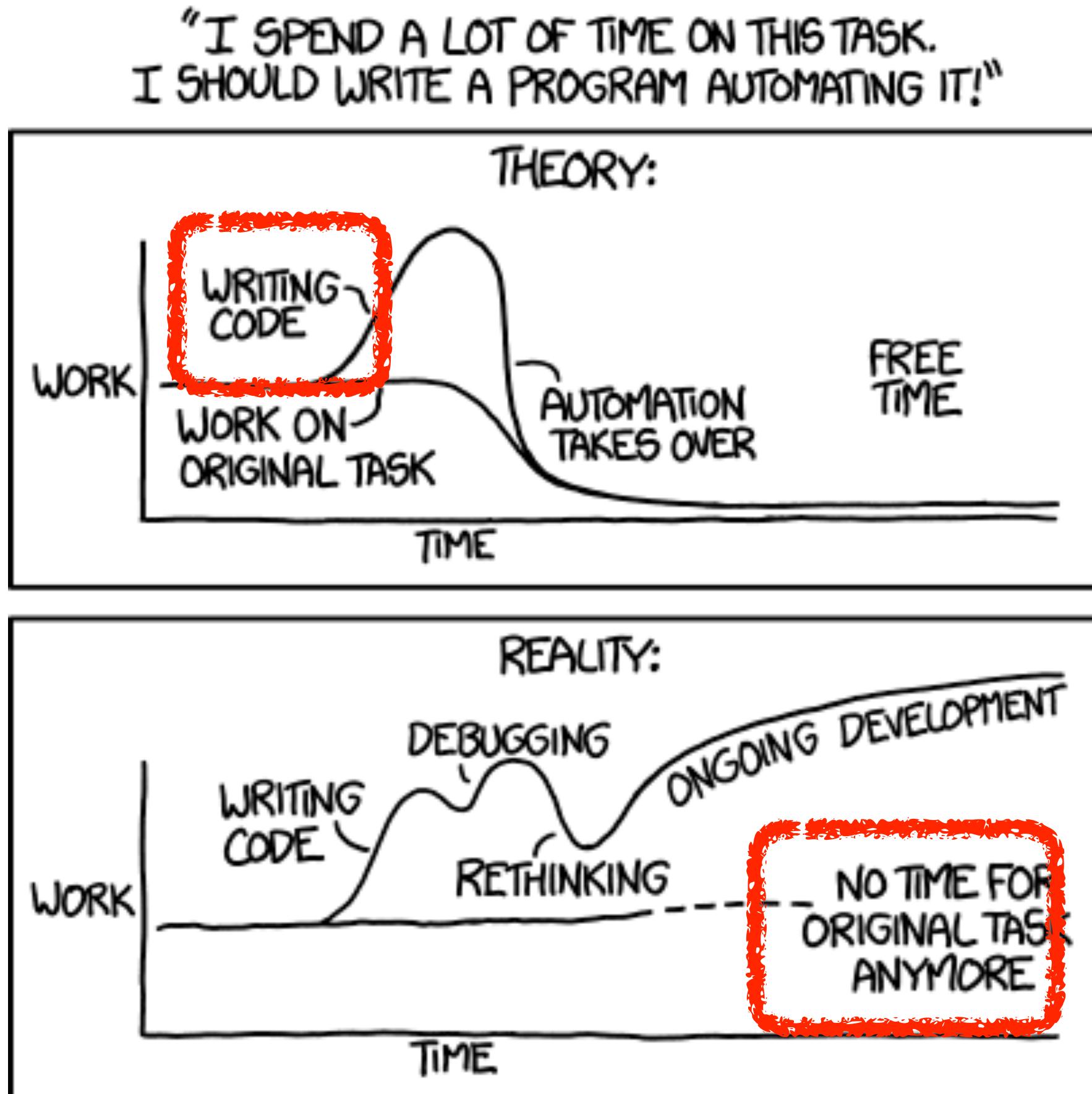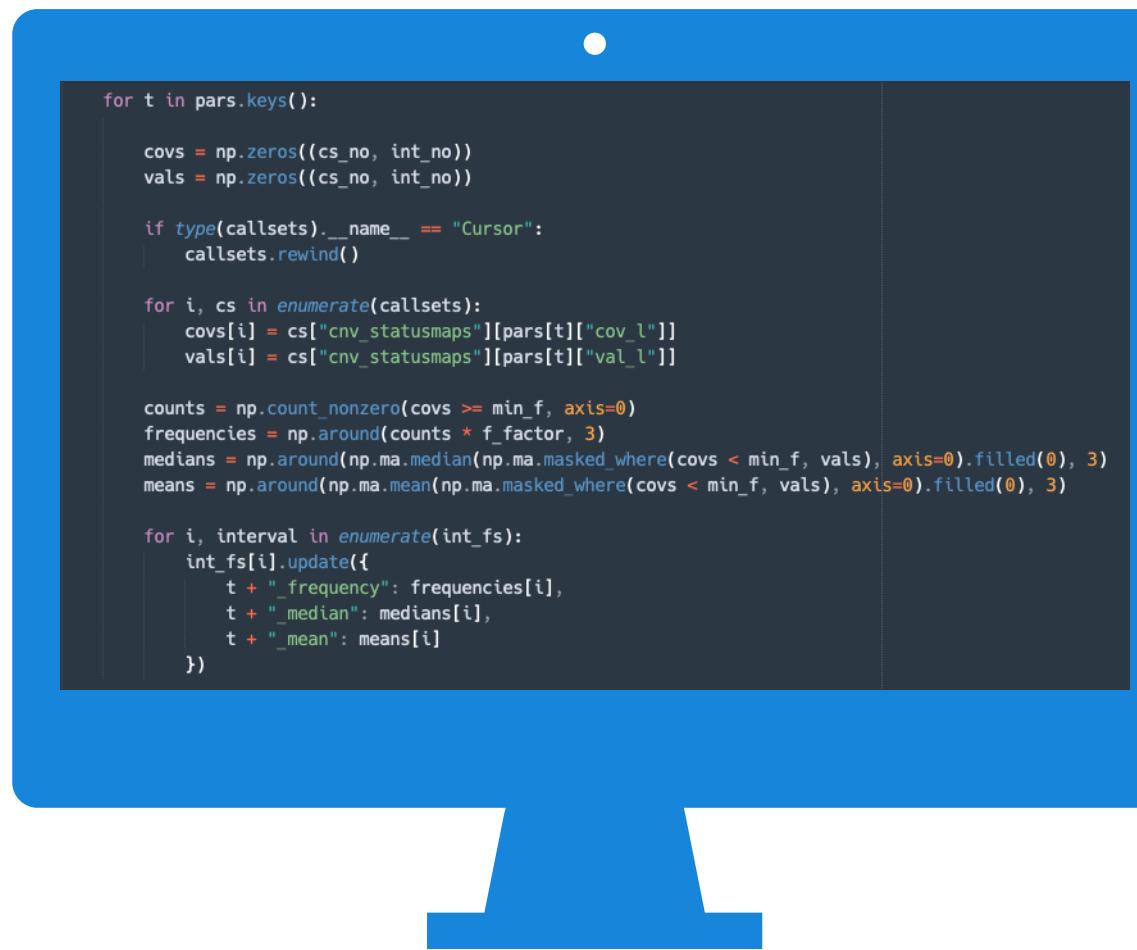| **Bio**informatician | Bio**informatician** |
|---|---|
| strong biological knowledge | sufficient biological background |
| provides hypothesis and / or dataset | provides statistical, analysis methods |
| **sufficient statistical** and **computational** expertise to correctly use bioinformatics tools & develop workflows (scripting ...) | **sufficient biological** or **medical** background to understand problems presented and identify pitfalls and hidden biases arising from data generation |
| expert **user** of informatics tools | **developer** of informatics tools |
| may get a Nobel | may get rich |

flux

# {BioInformaticsScience}
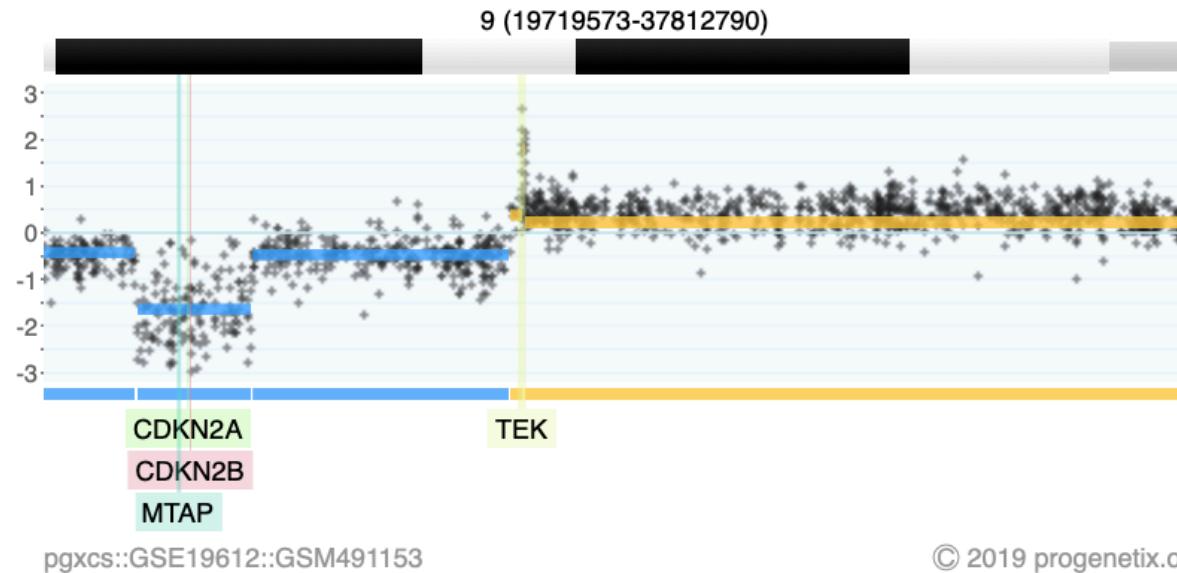


Randall Munroe - XKCD
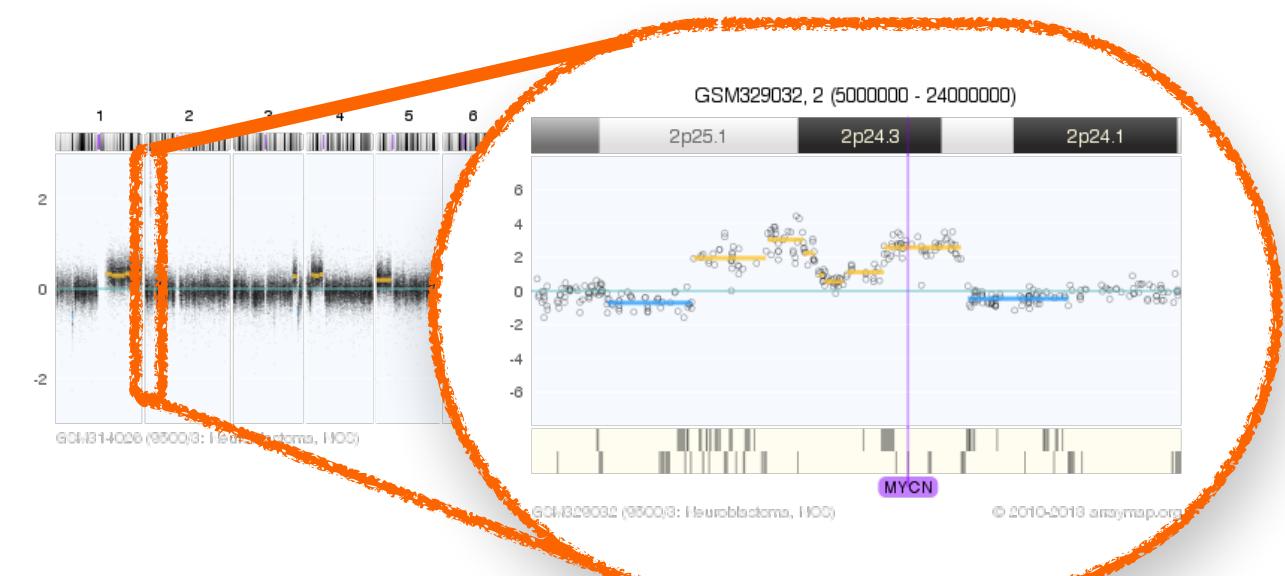https://xkcd.com/1319/

# Theoretical Cytogenetics and Oncogenomics Research | Methods | Standards

## Genomic Imbalances in Cancer - Copy Number Variations (CNV)
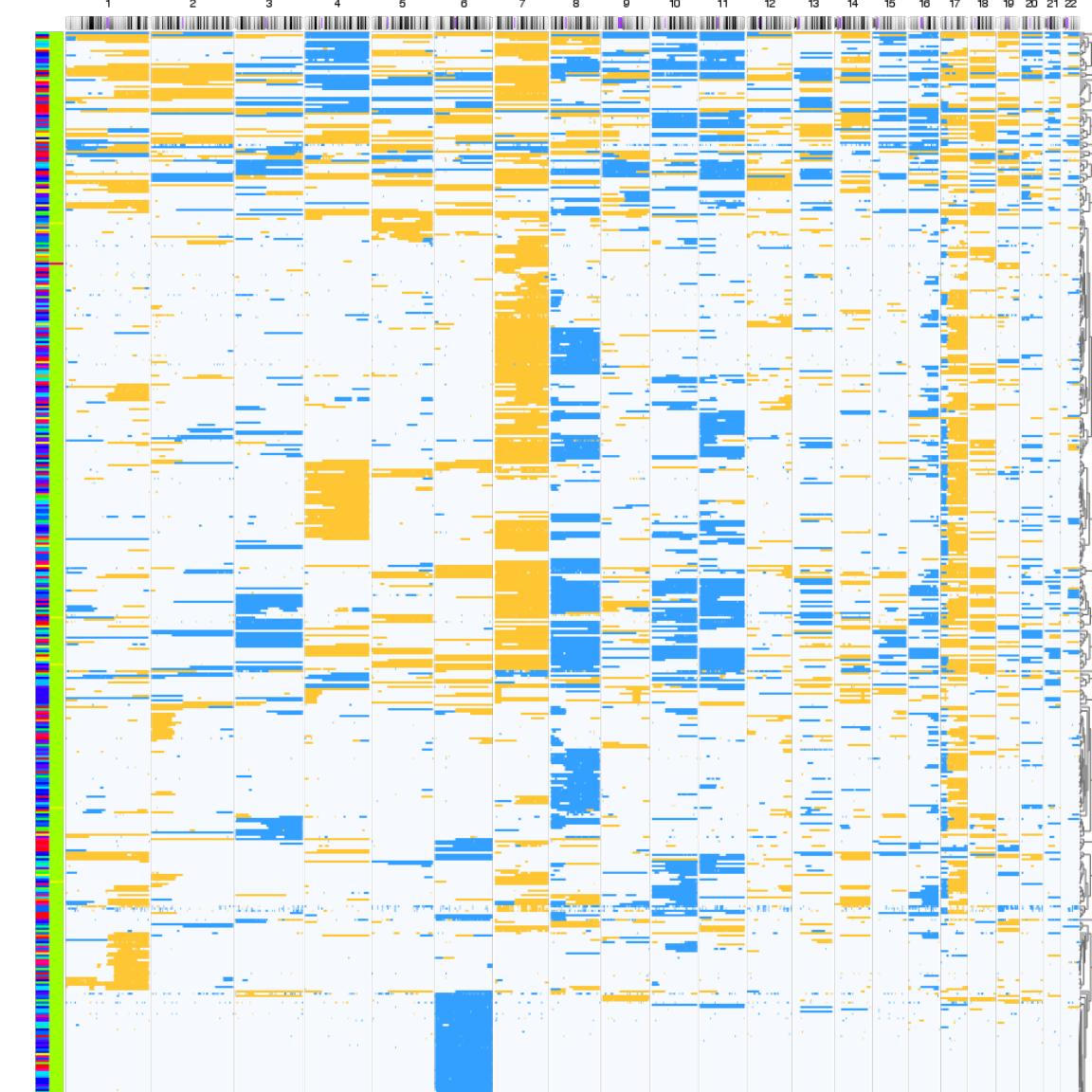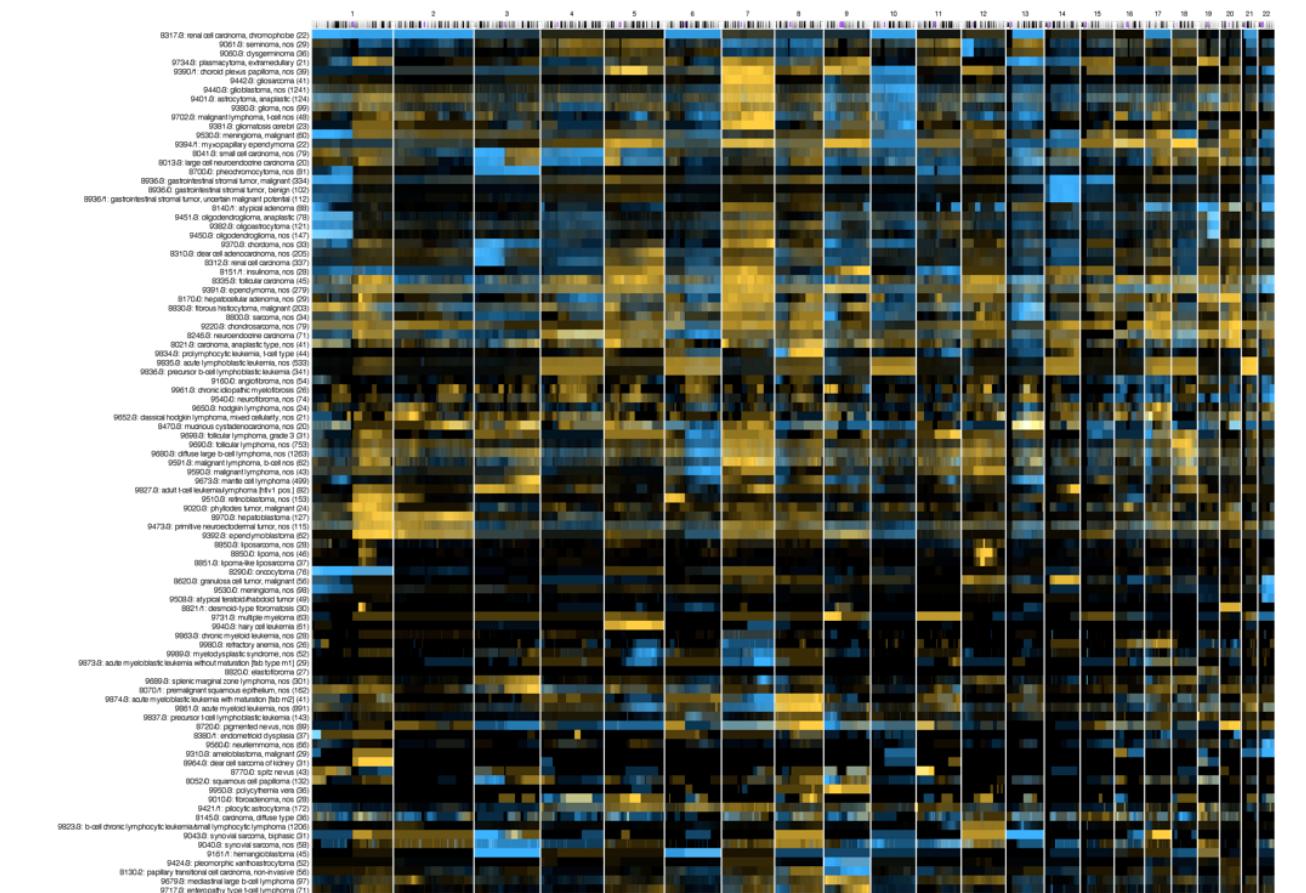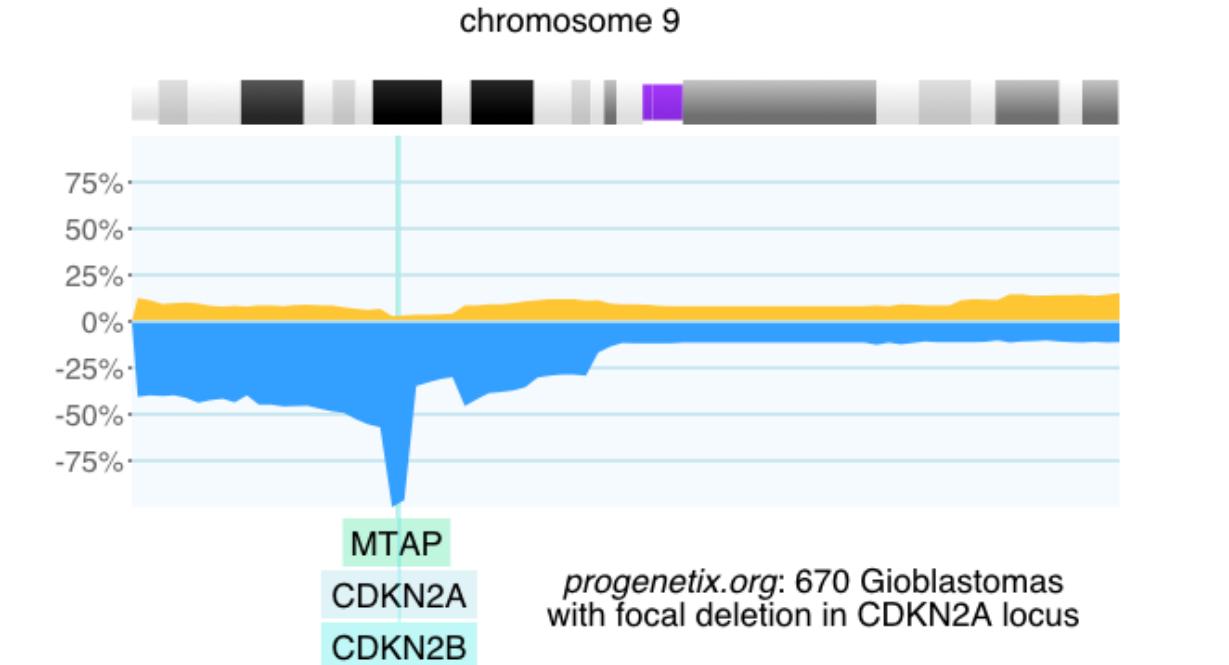
- Point mutations (insertions, deletions, substitutions)

- Chromosomal rearrangements

- **Regional Copy Number Alterations** (losses, gains)

- Epigenetic changes (e.g. DNA methylation abnormalities)



chromosome 9

*progenetix.org*: 670 Glioblastomas with focal deletion in CDKN2A locus







2-event, homozygous deletion in a Glioblastoma



MYCN amplification in neuroblastoma
(GSM314026, SJNB8_N cell line)

# progenetix.org

## 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

Universität Zürich UZH

progenetix

SIB Swiss Institute of Bioinformatics

---

progenetix

**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
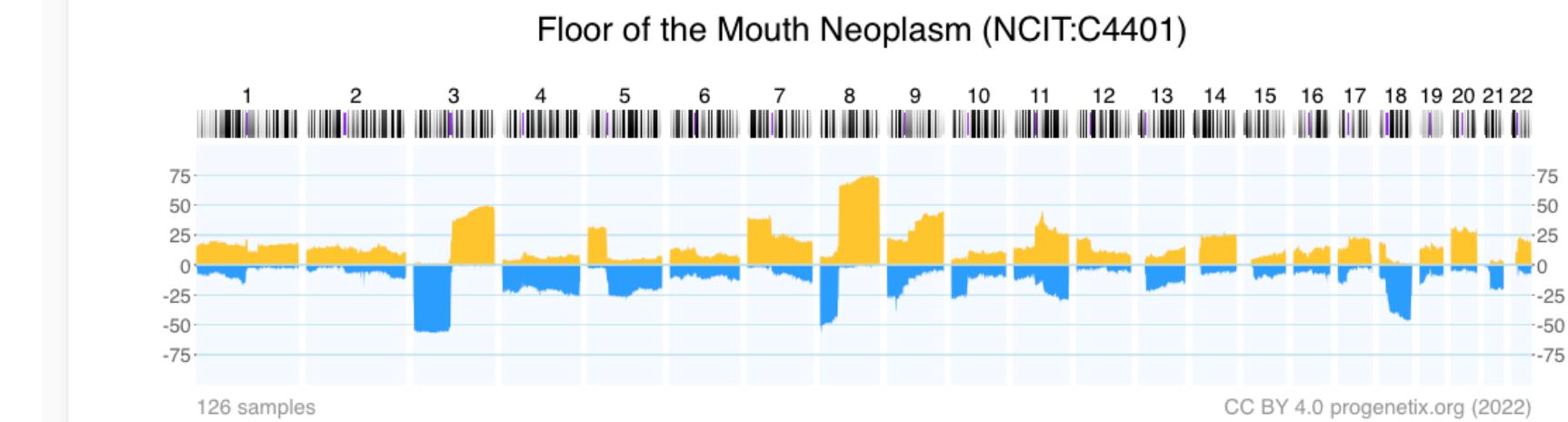
**Upload & Plot**

**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.
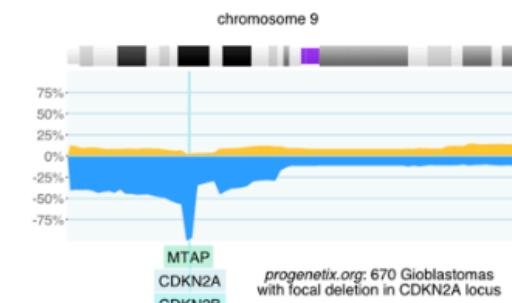
### Floor of the Mouth Neoplasm (NCIT:C4401)



126 samples                          CC BY 4.0 progenetix.org (2022)

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.

# progenetix.org

## 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



Universität Zürich UZH

progenetix

SIB Swiss Institute of Bioinformatics

# Cancer Cell Lines

## Cancer Genomics Reference Resource

- starting from >5000 cell line CNV profiles
  - ‣ 5754 samples | 2163 cell lines
  - ‣ 256 different NCIT codes
- genomic mapping of annotated variants and additional data from several resources (ClinVar, CCLE, Cellosaurus...)
  - ‣ 16178 cell lines
  - ‣ 400 different NCIT codes
- query and data delivery through Beacon v2 API

➡ integration in data federation approaches

cancercelllines.org

Lead: Rahel Paloots

---

Assembly: GRCh38  Chro: NC_000007.14  Start: 140713328  End: 140924929
Type: SNV

cellz

Matched Samples: 1058    UCSC region
Retrieved Samples: 1000  Variants in UCSC
Variants: 127            Dataset Responses (JSON)
Calls: 1444

Visualization options

Results | Biosamples | Variants | Annotated Variants

| Digest | Gene | Pathogenicity | Variant type | Variant Instances |
|---|---|---|---|---|
| 7:140834768-140834769:G>A | BRAF | | Missense variant | V: pgxvar-63ce6abca24c83054b B: pgxbs-3DfBeeAC |
| 7:140734714-140734715:G>A | BRAF | | Missense variant | V: pgxvar-63ce6acda24c83054b B: pgxbs-3fB2a14B |
| 7:140753334-140753339:T>TGTA | BRAF | Pathogenic | | V: pgxvar- |

---

cancercelllines

**Cancer Cell Lines by Cellosaurus ID**

The cancer cell lines in *cancercelllines.org* are labeled by th
hierarchially: Daughter cell lines are displayed below the pri
as a daughter cell line of **HeLa (CVCL_0030)** and so forth.

Sample selection follows a hierarchical system in which sam
response. This means that one can retrieve all instances and
for HeLa will also return the daughter lines by default - but

- Cancer Cell Lines°
- Search Cell Lines
- Cell Line Listing
- CNV Profiles by Cancer Type
- Documentation
  - News
- **Progenetix**
  - Progenetix Data
  - Progenetix Documentation
  - Publication DB

**Cell Lines (with parental/derived hierarchies**

Filter subsets e.g. by prefix    Hierarchy Depth

No Selection

- ☐ > cellosaurus:CVCL_0312: HOS (204 sa
- ☐ > cellosaurus:CVCL_1575: NCI-H650 (6
- ☐ > cellosaurus:CVCL_1783: UM-UC-3 (9
- ☐ ⌄ cellosaurus:CVCL_0004: K-562 (28 s
  - ☐ cellosaurus:CVCL_3827: K562/Ad
- ☐ > cellosaurus:CVCL_0589: Kasumi-1 (9

---

**Cell Line Details**

## HOS (cellosaurus:CVCL_0312)

**Subset Type**
- Cellosaurus - a knowledge resource on cell lines cellosaurus:CVCL_0312

**Sample Counts**
- 204 samples
- 57 direct *cellosaurus:CVCL_0312* code matches
- 21 CNV analyses

**Search Samples**
Select *cellosaurus:CVCL_0312* samples in the Search Form

**Raw Data (click to show/hide)**

HOS (cellosaurus:CVCL_0312)



21 CNV samples                                CC BY 4.0 progenetix.org (2023)

Download SVG | Go to cellosaurus:CVCL_0312 | Download CNV Frequencies

Gene Matches | Cytoband Matches | Variants

| ALK | . ABC-14 cells harbored no **ALK** mutations and were sensitive to ... crizotinib while also exhibiting MNNG **HOS** transforming gene ( MET ) | Rapid Acquisition of Alectinib Resistance in ALK-Positive Lung Cancer With High Tumor Mutation Burden (31374369) | ABSTRACT |
| AREG | crizotinib while also exhibiting MNNG **HOS** | Rapid Acquisition of Alectinib Resistance | ABSTRACT |

---

**Higher level of CNV coverage of the genomes of cancer cell lines compared to their origins**

Fold changes between genome CNV coverages of cell lines and tumors

Lead: Rahel Paloots

# Tumor subpopulations can be matched with highly similar cell lines

- Lung Small Cell Carcinoma Subpopulation

- Cell Lines:

  - CVCL_1140: COR-L279

  - CVCL_1455: NCI-H1105

  - CVCL_1527: NCI-H2107



Lead: Rahel Paloots

# Tumor subpopulations can be matched with highly similar cell lines?!



Lead: Rahel Paloots

# Somatic Mutations In Cancer: Patterns
## Making the case for genomic classifications
### Some related cancer entities show similar copy number profiles

9390/1: choroid plexus papilloma, nos (39)
9442/3: gliosarcoma (41)
9440/3: glioblastoma, nos (1241)
9401/3: astrocytoma, anaplastic (124)
9380/3: glioma, nos (99)
9702/3: malignant lymphoma, t-cell nos (48)
9381/3: gliomatosis cerebri (23)
9530/3: meningioma, malignant (60)
9394/1: myxopapillary ependymoma (22)

9451/3: oligodendroglioma, anaplastic (78)
9382/3: oligoastrocytoma (121)
9450/3: oligodendroglioma, nos (147)

9698/3: follicular lymphoma, grade 3 (31)
9690/3: follicular lymphoma, nos (753)
9680/3: diffuse large b-cell lymphoma, nos (1263)
9591/3: malignant lymphoma, b-cell nos (62)
9590/3: malignant lymphoma, nos (43)
9673/3: mantle cell lymphoma (499)

9984/3: refractory anemia with excess blasts in transformation [raebt] (24)
9983/3: refractory anemia with excess blasts [raeb] (38)
9867/3: acute myelomonocytic leukemia [fab type m4] (32)
9920/3: therapy-related acute myeloid leukemia, nos (32)
9891/3: acute monoblastic leukemia [fab m5] (23)

9051/3: desmoplastic mesothelioma (59)
9053/3: mesothelioma, biphasic, malignant (27)
9050/3: mesothelioma, nos (81)
9052/3: epithelioid mesothelioma, malignant (64)



arrayMap

# CNV profiles heterogeneity vs cancer classification
## Correspondance of genomic profiles to NCIT cancer hierarchy



Nerve Sheath Neoplasm (NCIT:C4972)

414 CNV samples

Chondrogenic Neoplasm (NCIT:C4755)

Neoplasm by Morphology

Lung Squamous Cell Carcinoma (NCIT:C3493)

1938 CNV samples

Lung Adenocarcinoma (NCIT:C3512)

4644 CNV samples

Lead: Ziying Yang

# Results
## Entity CNV heterogeneity: Glioblastoma



| group cluster | CNV features |
|---|---|
| 15968 | Dup 7 |
| 19069 | Del 10 |
| 19198 | Dup 7, Del 10 |
| 22279 | Dup 7, Del 10, Dup 19 |
| 22292 | Dup 7, Del 10, Del 13 |
| 28527 | Del 1p, Del 19q |
| 29242 | Dup 19 |
| 30914 | Dup 7, Del 10, Dup 19, Dup 20 |

Lead: Ziying Yang

# CNV Categorization

## different levels of CNV

( high amplitude )

Focal SCNAs

Arm-level SCNAs

( low amplitude e.g. single-copy changes )

Rameen et al 2010 Nature



GSM329032, 2 (5000000 - 24000000)

2p25.1     2p24.3     2p24.1

MYCN

MYCN amplification in neuroblastoma
(GSM314026,  SJNB8_N cell line)

Lead: Hangjia Zhao

## CopyNumberChange

*Copy Number Change* captures a categorization of copies of a molecule within a system, relative to a baseline. These types of Variation are common outputs from CNV callers, particularly in the somatic domain where integral CopyNumberCount are difficult to estimate and less useful in practice than relative statements. Somatic CNV callers typically express changes as relative statements, and many HGVS expressions submitted to express copy number variation are interpreted to be relative copy changes.

### Computational Definition

An assessment of the copy number of a Location or a Feature within a system (e.g. genome, cell, etc.) relative to a baseline ploidy.

### Information Model

Some CopyNumberChange attributes are inherited from Variation.

| Field | Type | Limits | Description |
|---|---|---|---|
| _id | CURIE | 0..1 | Variation Id. MUST be unique within document. |
| type | string | 1..1 | MUST be "CopyNumberChange" |
| subject | Location \| CURIE \| Feature | 1..1 | A location for which the number of systemic copies is described. |
| copy_change | string | 1..1 | MUST be one of "efo:0030069" (complete genomic loss), "efo:0020073" (high-level loss), "efo:0030068" (low-level loss), "efo:0030067" (loss), "efo:0030064" (regional base ploidy), "efo:0030070" (gain), "efo:0030071" (low-level gain), "efo:0030072" (high-level gain). |

# labelSeg

## segment annotation for tumor copy number variation profiles

Signal from probes in microarray or from reads in NGS

Segmentation

a step to split the chromosomes into regions of equal copy number that accounts for the noise in the data.



README.md

# labelSeg

This is an R package designed to identify and label different levels of Copy Number Alterations (CNA) in segmented profiles.

## Installation

To install the package, you can use the `devtools` package as follows:

```
install.packages("devtools")
devtools::install_github("baudisgroup/labelSeg")
```

Packages

No packages published

Languages

● R 100.0%

Lead: Hangjia Zhao

# Pipeline Development
## improve CNV calling in large numbers of heterogeneous cancer samples



Lead: Hangjia Zhao

# CNV Categorization

## different levels of CNV

## labelSeg: segment annotation for tumor copy number alteration profiles

Hangjia Zhao (iD) and Michael Baudis (iD)

Corresponding author: Michael Baudis, Department of Molecular Life Sciences, University of Zurich, Winterthurerstrasse 190, CH-8057 Zurich, Switzerland.
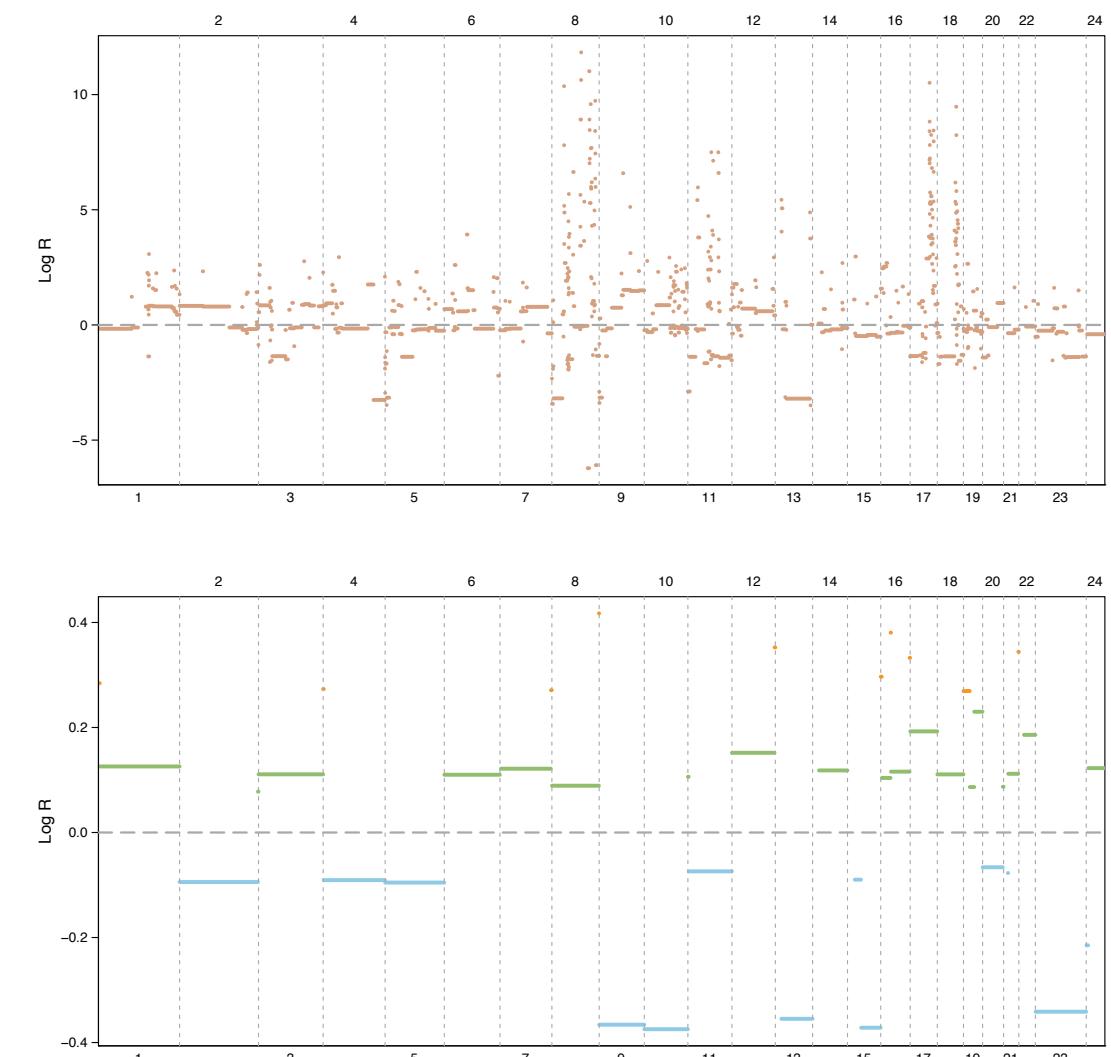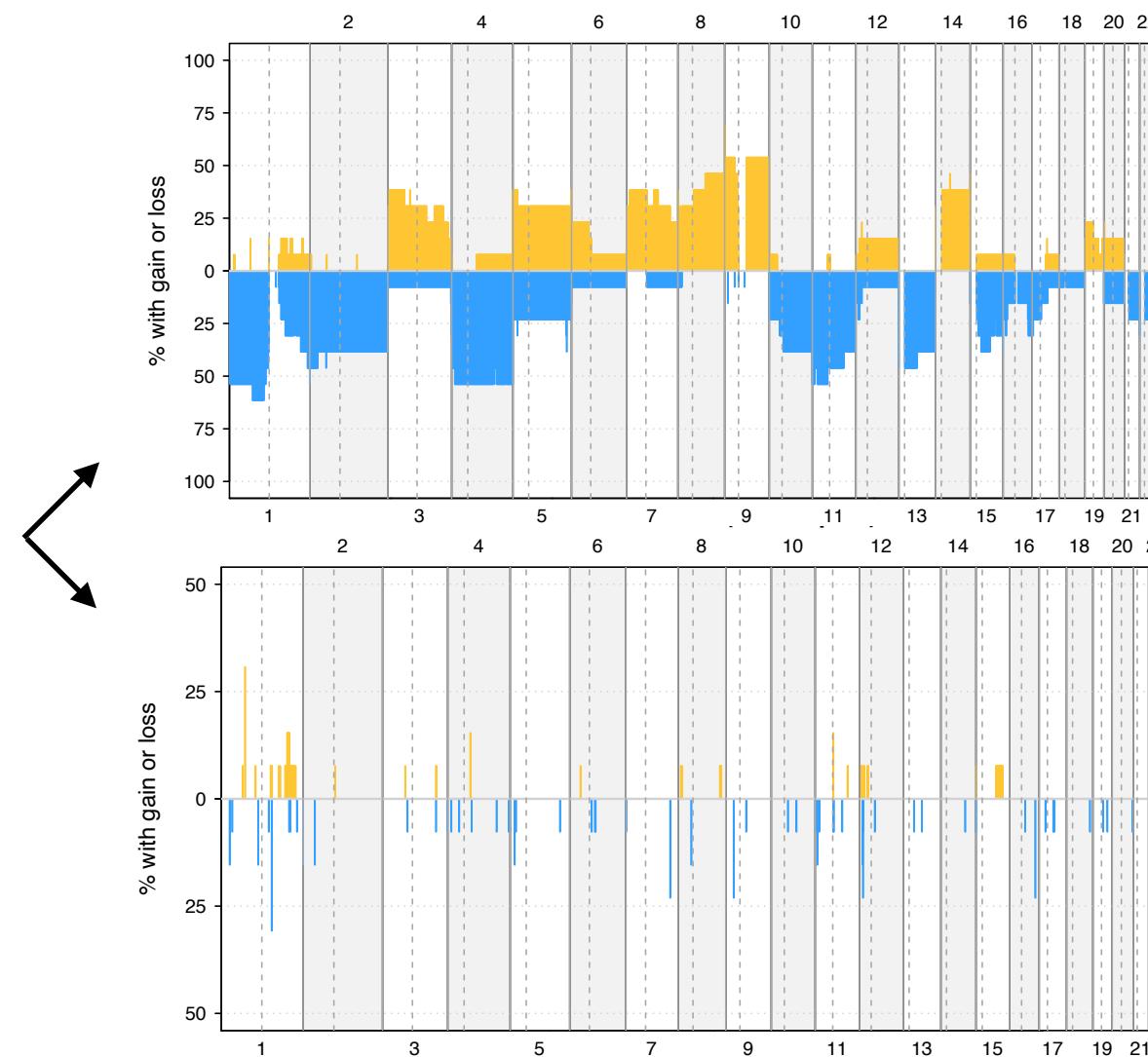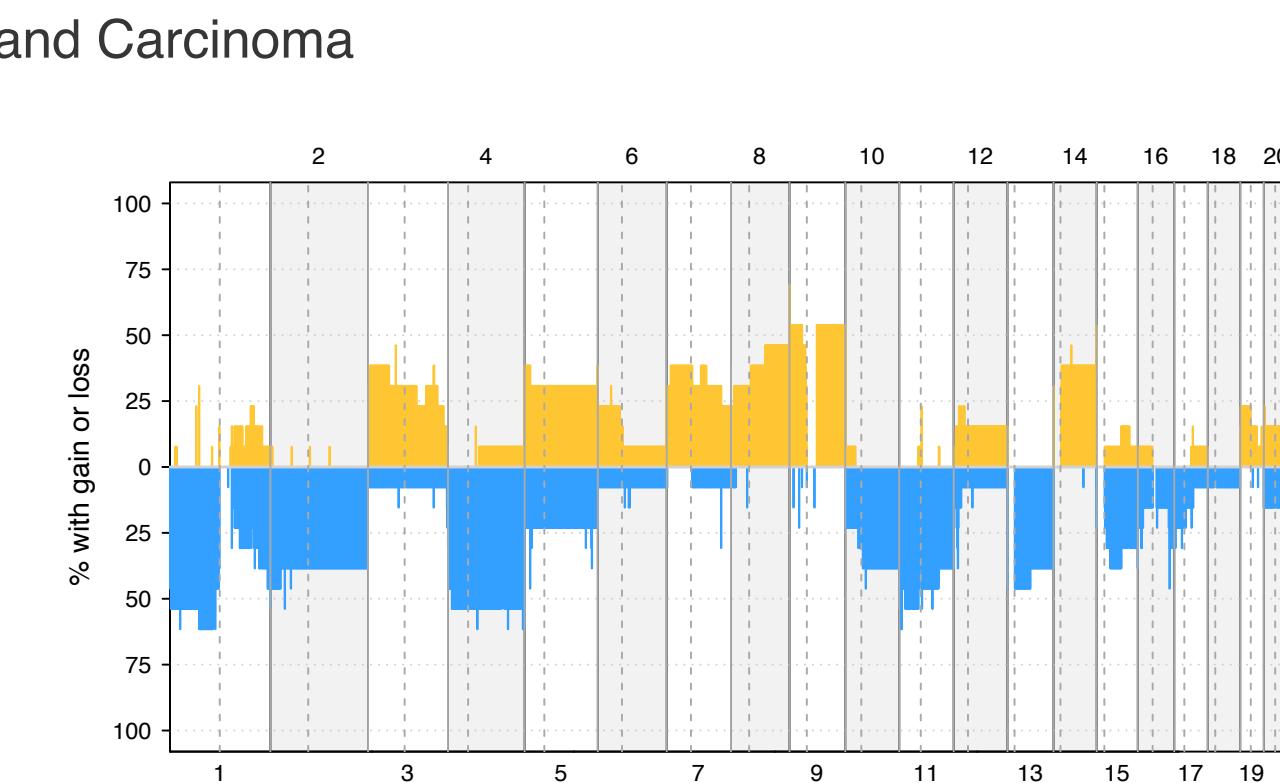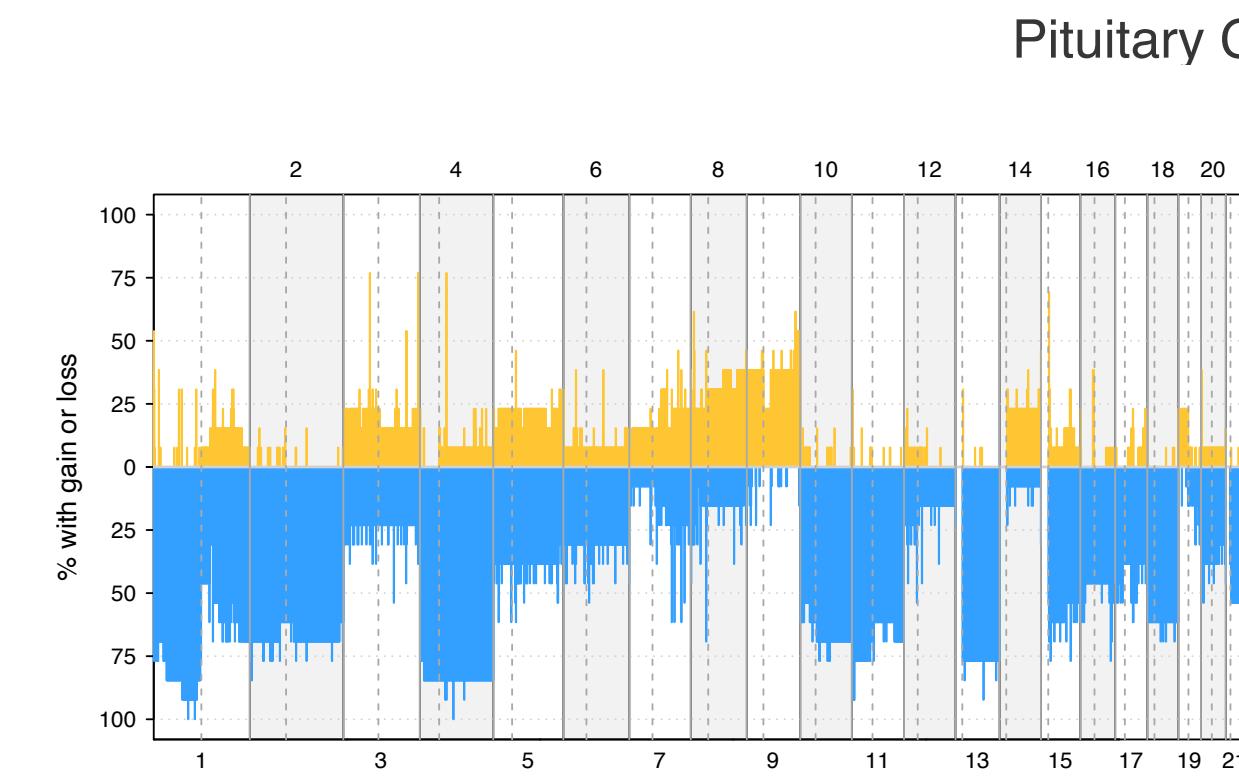
Tel.: (+41) 44 635 34 86; E-mail: michael.baudis@mls.uzh.ch

**CopyNumberChange**

Copy Number Change captures a categorization of copies of a molecule within a system, relative to a ... callers, particularly in the somatic ... and less useful in practice than ... is relative statements, and many ... interpreted to be relative copy ... a system (e.g. genome, cell,

MYCN amplification in neuroblastoma (GSM314026, SJNB8_N cell line)

| | | | |
|---|---|---|---|
| _id | CURIE | 0..1 | Variation id: MUST be unique within document. |
| type | string | 1..1 | MUST be "CopyNumberChange" |
| subject | Location \| CURIE \| Feature | 1..1 | A location for which the number of systemic copies is described. |
| copy_change | string | 1..1 | MUST be one of "efo:0030069" (complete genomic loss), "efo:0020073" (high-level loss), "efo:0030068" (low-level loss), "efo:0030067" (loss), "efo:0030064" (regional base ploidy), "efo:0030070" (gain), "efo:0030071" (low-level gain), "efo:0030072" (high-level gain). |

Lead: Hangjia Zhao

# TCGA BLCA project (pgx:TCGA.BLCA)

**TCGA BLCA project**
*Variant Distribution on Genome Coordinates*

CDKN2A

TP53

Beyond CNVs – Kay von Grünigen

# Where does Genomic Data Come From?
## Geographic bias in published cancer genome profiling studies

Articles

### Geographic assessment of cancer genome profiling studies

Paula Carrio-Cordo[1,2], Elise Acheson[3], Qingyao Huang[1,2] and Michael Baudis[1,*]

[1]Institute of Molecular Life Sciences, University of Zurich, Zurich, Switzerland [2]Swiss Institute of Bioinformatics, Zurich, Switzerland [3]Department of Geography, University of Zurich, Zurich, Switzerland

Map of the geographic distribution (by first author affiliation) of the 104'543 genomic array, 36'766 chromosomal CGH and 15'409 whole genome/exome based cancer genome datasets. The numbers are derived from the 3'240 publications registered in the Progenetix database.

# Global Alliance
## for Genomics & Health

## Collaborate. Innovate. Accelerate.

# A federated ecosystem for sharing genomic, clinical data

Silos of genome data collection are being transformed into seamlessly connected, independent systems

The Global Alliance for Genomics and Health*

# GA4GH Roadmap Development Process

**Global Alliance**
for Genomics & Health

| 2018 | 2019 | 2020 | 2021 | 2022 | 2023 | 2024 |
|------|------|------|------|------|------|------|

**V1 ROADMAP**

15 standards approved | 8 in development · · · · · · · · · · · · · · · · · · · · · · *Ongoing development work* · · · →

**GAP ANALYSIS** | **V2 ROADMAP**

Standards development · · · · · · · · · · · · · · *Ongoing development work* · · · →

**GAP ANALYSIS** | **V3 ROADMAP**

Standards development · · · →

ga4gh.org

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

ga4gh.org

# 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



**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**

# 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


European Genome-Phenome Archive

- EGA "owns" nothing; data controllers tell who is authorized to access *their* datasets

- EGA admins provide smooth "all or nothing" data sharing process



## # Files



- Array 444.037
- FASTQ 1.167.840
- VCF 904.852
- BAM-CRAM 1.449.676

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

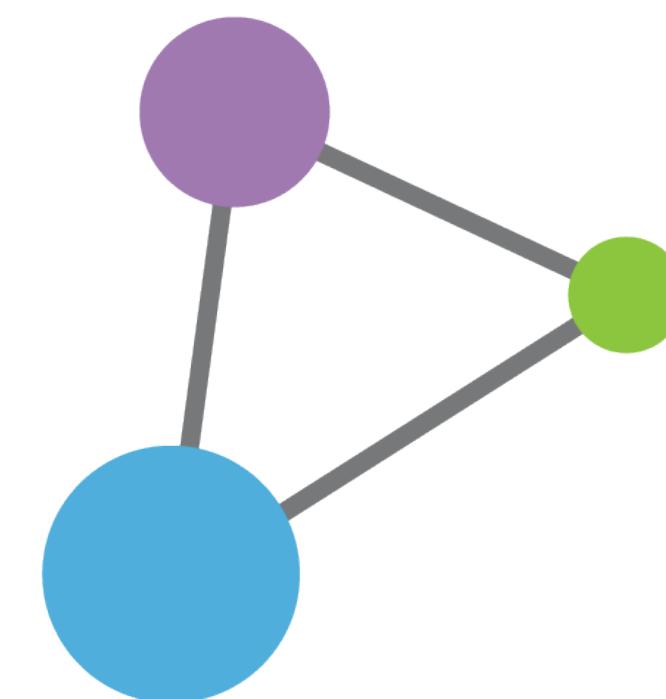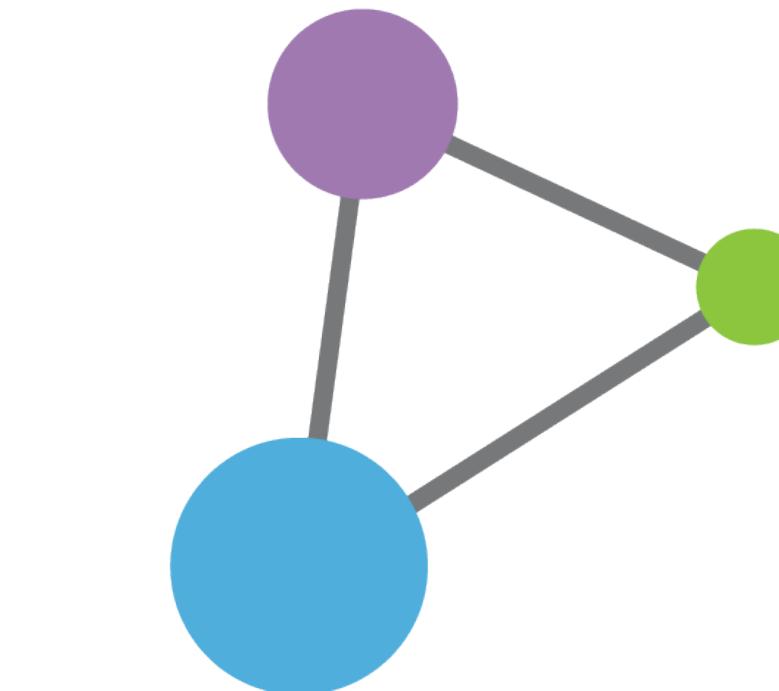# Different Approaches to Data Sharing



**Centralized Genomic Knowledge Bases**

**Data Commons**
Trusted, controlled repository of multiple datasets

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

**Linkage of distributed and disparate datasets**

# The Swiss Personalized Health Network

# Different Approaches to Data Sharing

**Centralized Genomic Knowledge Bases**

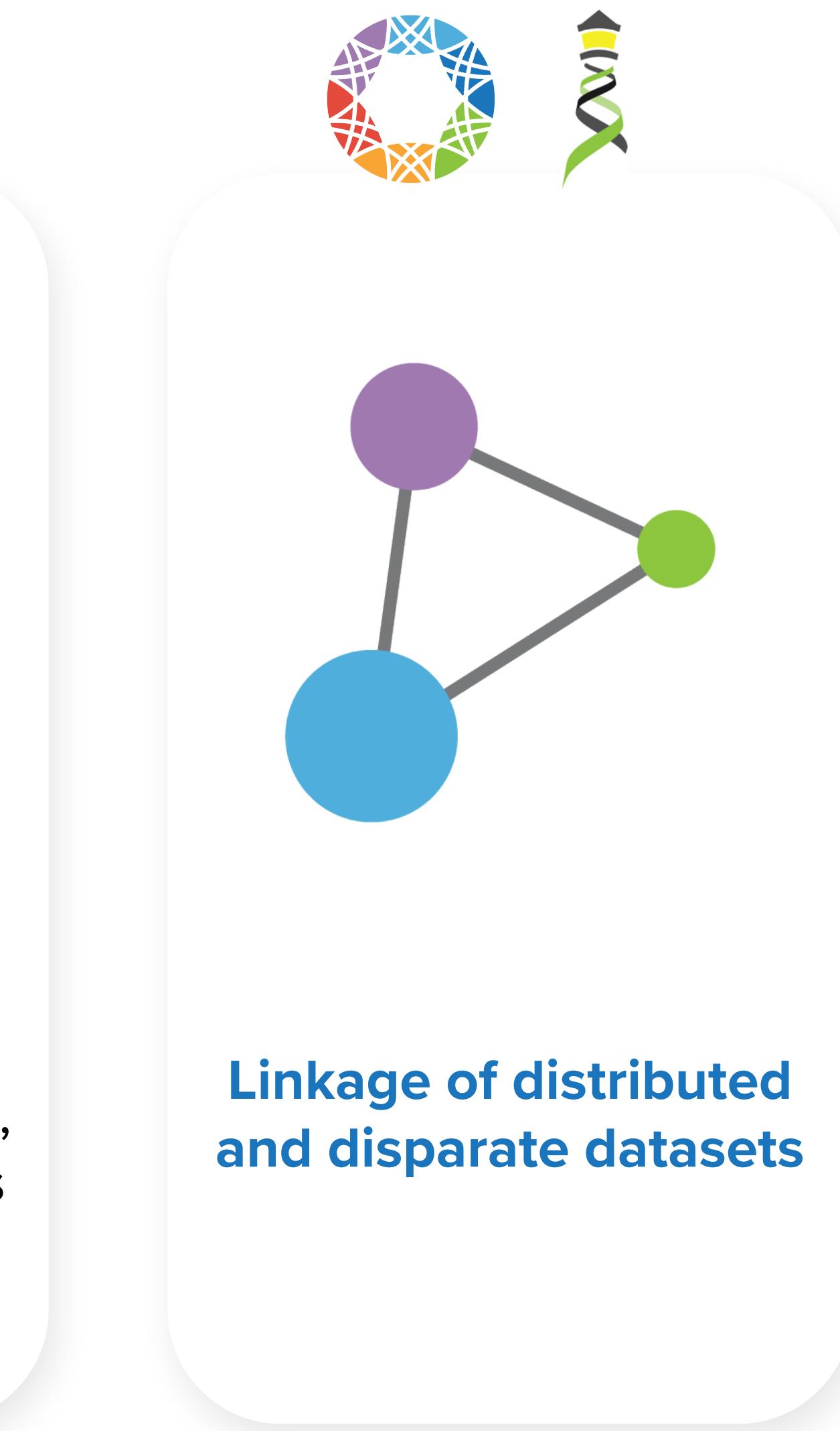**Data Commons**
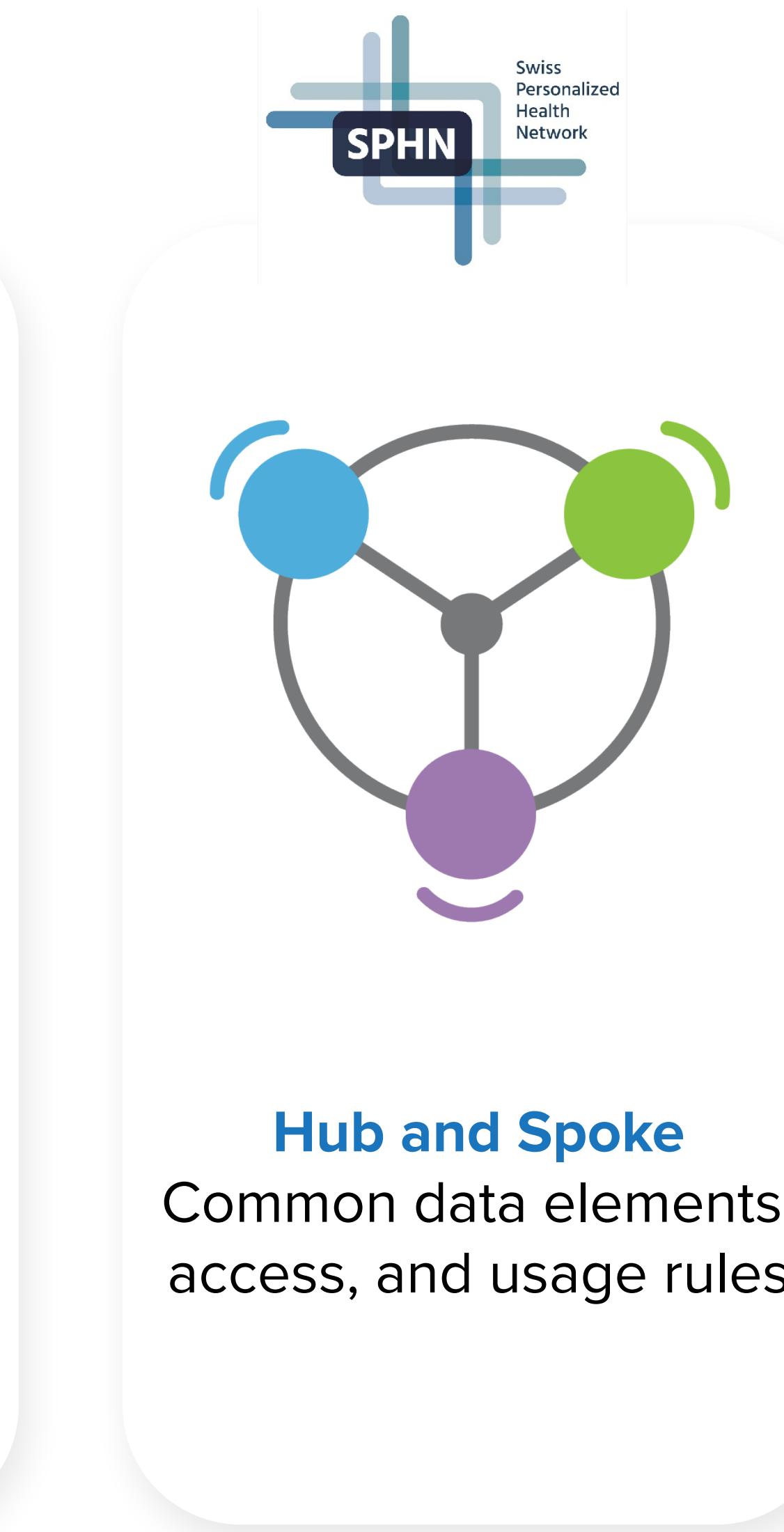Trusted, controlled repository of multiple datasets
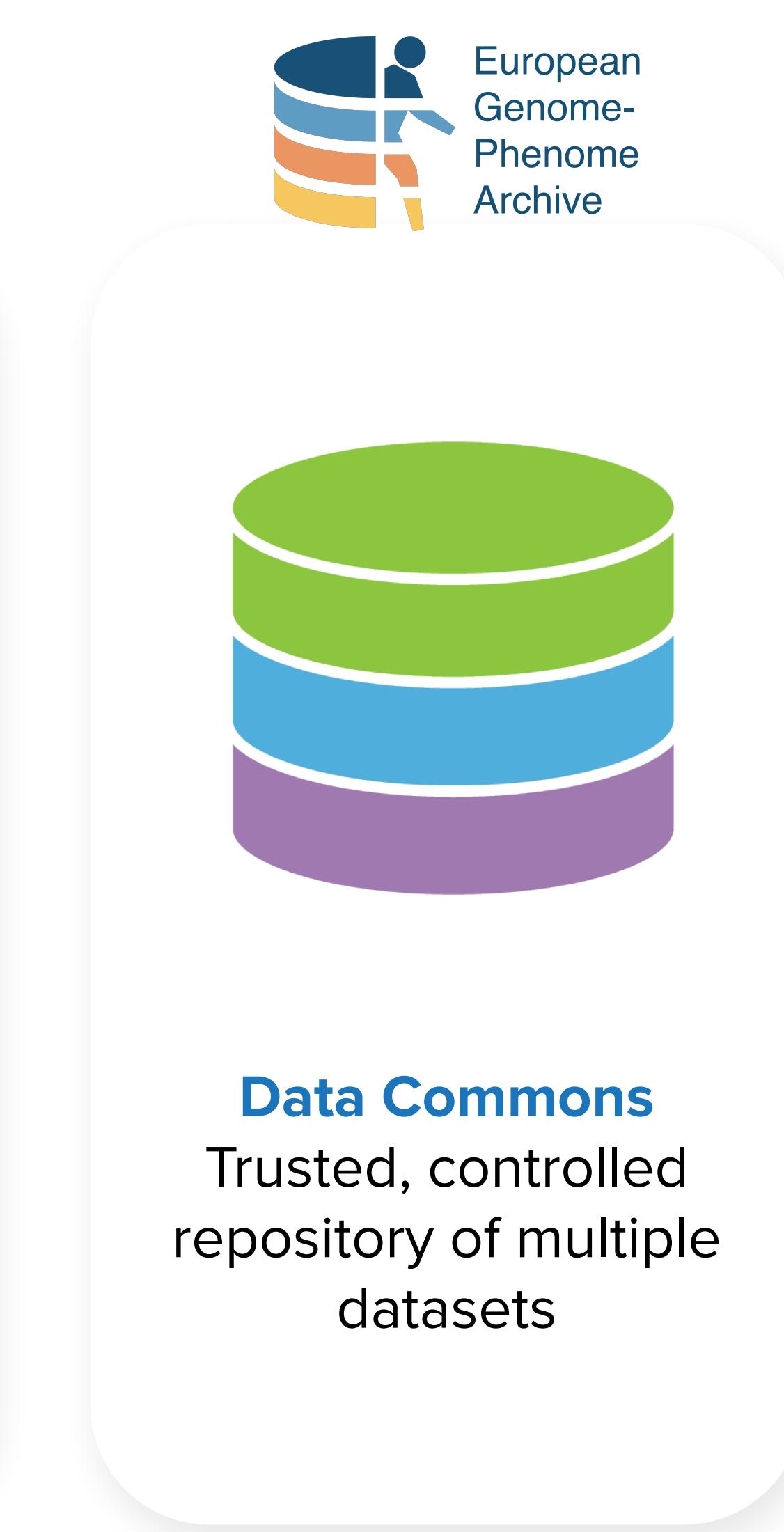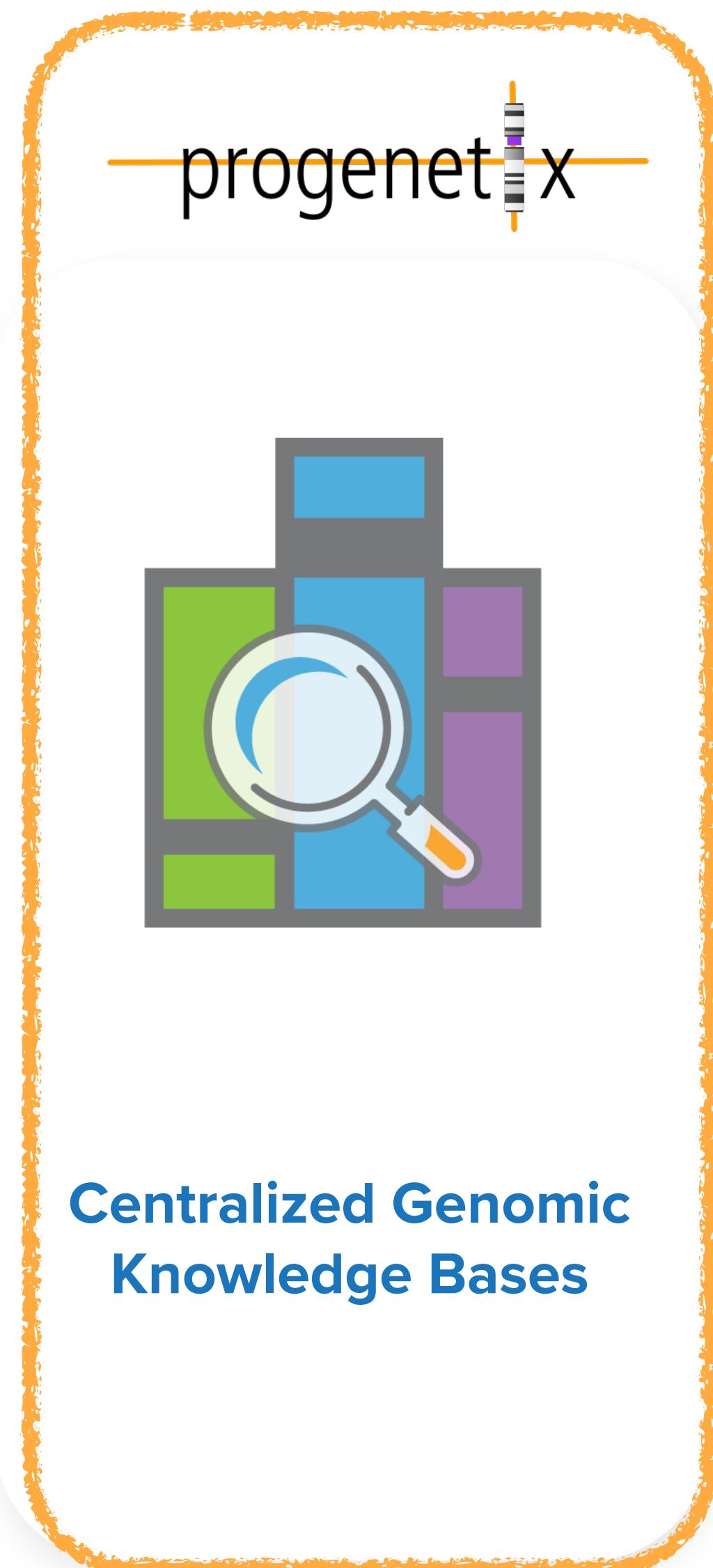
**Hub and Spoke**
Common data elements, access, and usage rules

**Linkage of distributed and disparate datasets**

**Federation**

ATTTATCTGCTCTCGTTG
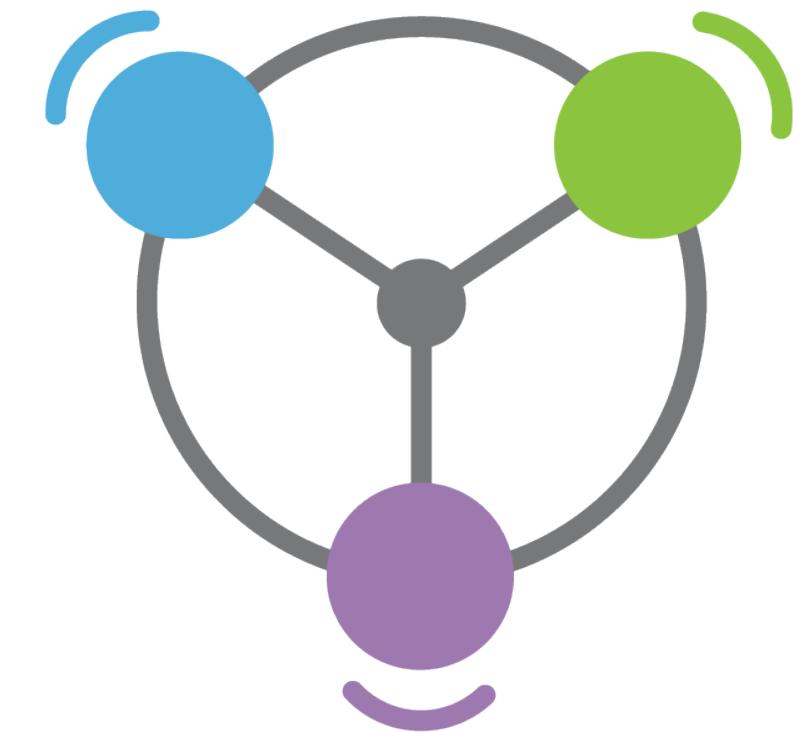GAAGTACAAAATTCATTAAT
GCTATGCACAAAATCTGTAG
CTAGTGTCCCATCTATTT

**A federated data ecosystem.** To share genomic data globally, this approach furthers medical research without requiring compatible data sets or compromising patient identity.

Genomics API

Framework for Responsible Sharing of Genomic and Health-Related Data

Privacy and Security Policy

| A | B |
| F | | C |
| E | D |

Beacon

Matchmaker Exchange

BRCA Challenge

Other International Data-Sharing Projects

Data are organized, secured, and made accessible through federated use of GA4GH tools

GENOMICS

# A federated ecosystem for sharing genomic, clinical data

Silos of genome data collection are being transformed into seamlessly connected, independent systems

Global Alliance for Genomics & Health

# A New Paradigm for Data Sharing

**Data Copying**

**Data Visiting**

# A New Paradigm for Data Sharing



FROM

TO

STANDARDS

Data Copying

Data Visiting

INFORMING HUMAN HEALTH & MEDICINE

**Governing Outputs**
- Return of Results Policy
- Standard Genomic Data Licenses & Agreements

**Find Datasets**
- Data Connect API
- Beacon API
- Data Use Ontology

**Discover Services**
- Service Registry
- Service Info

**Retrieve Datasets**
- refget
- htsget
- RNAget

**Analyze Datasets**
- TRS
- WES
- TES
- DRS

**Share Datasets**
- VRS
- VA
- Phenopackets
- Pedigree

**DATABASE**

**DATA DONOR**
- Your DNA Your Day

**Consents**
- Consent Policy
- Consent Clauses
- Machine-Readable Consent Guidance
- Data Use Ontology

**Genomic Sequencing**
- CRAM/BAM
- VCF
- Crypt4GH
- Data Privacy and Security Policy

Record

**RESEARCH ETHICS COMMITTEE**
- Ethics Review and Recognition Policy

Data transformation for database storage

**DATA STEWARD**
- Data Security Infrastructure Policy

**Apply for GA4GH Passport**
- GA4GH Passport

APPROVED

**Approval of Data Access Request**
- GA4GH Passport
- Data Use Ontology

**Request Access to Dataset**
- GA4GH Passport
- AAI

**RESEARCHER / CLINICIAN**

**DATA ACCESS COMMITTEE**
- Data Access Committee Review Standards

ga4gh.org

# Overview of GA4GH standards and frameworks

| | | | | | |
|---|---|---|---|---|---|
| **Clin/Pheno Data Capture** | Phenopackets | Pedigree Representation | Cohort Representation | | |
| **Cloud** | Workflow Execution Service | Tool Registry Service | Data Repository Service | Task Execution Service | Cloud Testbed Interoperability |
| **Discovery** | Beacon | Service Info | Service Registry | Data Connect | |
| **Data Security** | Authentication & Authorization Infrastructure | Data Security Infrastructure Policy | Risk Assessment | Bad Actors in Research Environments | Cloud Security & Privacy |
| **Data Use & Researcher Identity** | Data Use Ontology | GA4GH Passports | Machine Readable Consent Guidance | Data Access Committee Review Standards Toolkit | |
| **Genomic Knowledge Standards** | Variation Representation | Variation Annotation | Sequence Annotation | | |
| **Large Scale Genomics** | htsget API | refget API | SAM/BAM/CRAM | VCF | Crypt4GH | rnaget API | BED File Format |
| **Regulatory & Ethics** | Framework for responsible data sharing | Consent Toolkit | 20+ other policy tools/frameworks | Genetic Discrimination Toolkit | GDPR Forum | Public Attitudes for genomic policy |

Legend: Approved · Ongoing · In Development

ga4gh.org

## Phenopackets v2

Phenopackets is a standard schema for sharing phenotypic information.

**Approved:** June 24, 2021

## VCF/BCF

The Variant Call Format (VCF) specifies the format of a text file used in bioinformatics for storing gene sequence variations. The Binary Call Format (BCF) is the Binary equivalent, smaller and more efficient to process.

**Software Libraries:** htslib | htsjdk

**Tools:** Samtools | BCFtools

**Databases:** European Variation Archive (EVA) | dbGAP | dbSNP | 1000 Genomes Projects / IGSR

**Genome Browsers:** ENSEMBL | JBrowse | UCSC Genome Browser

**Example Users**

# CRAM

CRAM is a file format for storing compressed genomic data. To make files small and efficient, the algorithm compresses information by only storing the parts that are different from the reference human genome.



**1.5 million+** CRAM files store more than **4 petabytes** of compressed genomic data around the globe

*CRAM compresses data by only storing the difference.*

## Genomics England implements GA4GH API to provide secure access to genomic data for the NHS

Genomics England has implemented the standard GA4GH API hts
Genomes Program and the Genomic Medicine Service.

14 Feb 2024

## NIH and GA4GH commit to ongoing collaboration

**NIH and GA4GH strengthen their partnership to expand responsible data use for the benefit of human health through a Memorandum of Agreement.**



The United States National Institutes of Health (NIH) Office of Data Science Strategy (ODSS) and the Global Alliance for Genomics and Health (GA4GH) have announced a strategic collaboration in the form of a Memorandum of Agreement. This partnership aims to bolster the development of technology standards, tools, and policy frameworks to support responsible sharing of genomic and related health data on a global scale.

17 : 7577121 G > A

A ***Beacon*** answers a query for a specific genome variant against individual or aggregate genome collections
**YES** | **NO** | **\0**

17 : 7577121 G > A

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**.

# Global Alliance "Beacon" - Jim Ostell, NCBI, March 7, 2014

## Introduction

… I proposed a challenge application for all those wishing to seriously engage in ***international*** data sharing for human genomics. …

1.  Provide a public web service
2.  Which accepts a query of the form "Do you have any genomes with an "A" at position 100,735 on chromosome 3?"
3.  And responds with one of "Yes" or "No" …

"Beacon" because … people have been scanning the universe of human research for ***signs of willing participants in far reaching data sharing***, but … it has remained a dark and quiet place. The hope of this challenge is to 1) ***trigger the issues*** blocking groups … in way that isn't masked by the … complexities of the science, fully functional interfaces, and real issues of privacy, and to 2) in ***short order*** … see ***real beacons of measurable signal*** … from ***at least some sites*** … Once your "GABeacon" is shining, you can start to take the ***next steps to add functionality*** to it, and ***finding the other groups*** … following their GABeacons.

## Utility

Some have argued that this simple example is not "useful" so nobody would build it. Of course it is not the first priority for this application to be scientifically useful. …intended to provide a ***low bar for the first step of real*** … ***engagement***. … there is some utility in …locating a rare allele in your data, … not zero.

A number of more useful first versions have been suggested.

1.  Provide ***frequencies of all alleles*** at that point
2.  Ask for all alleles seen in a gene ***region*** (and more elaborate versions of this)
3.  Other more complicated queries

"I would personally recommend all those be held for **version 2**, when the beacon becomes a service."
Jim Ostell, 2014

## Implementation

1.  Specifying the chromosome … The interface needs to specify the ***accession.version*** of a chromosome, or ***build number***…
2.  Return values … right to ***refuse*** to answer without it being an error … DOS ***attack*** … or because …especially ***sensitive***…
3.  Real time response … Some sites suggest that it would be necessary to have a ***"phone home" response*** …

# Beacon v2

docs.genomebeacons.org

Biosamples

Individuals

Genomic Variations

Runs

Analyses

Datasets

Cohorts

*Beacon Model*

Beacon Framework (protocol)

# Beacon API v2

The Beacon API can be implemented as a web-accessible service that users may query for information about a specific allele.

**Approved:** April 21, 2022



PUBLIC — Accessible to users of the internet

REGISTERED — Accessible to users with an account
e.g. Bona fide researcher

CONTROLLED — Accessible to authorized users
e.g. Signed agreement, agree to data use conditions

9:18000000,21975098-
21967753,26000000:DEL
ncit:C3058
DUO:0000004
HP:0003621

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"**.

## Example Users

EMBL-EBI

Australian Genomics

SciLifeLab

elixir

UNIVERSITY OF CALIFORNIA SANTA CRUZ

BROAD INSTITUTE

EUROPEAN GENOME-PHENOME ARCHIVE

International Cancer Genome Consortium

CDKN2A:DEL
size<1Mb
granularity:record
ncit:C3058
DUO:0000004
HP:0003621

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

# Progenetix and GA4GH Beacon

## Implementation driven development of a GA4GH standard

| **Beacon v1 Development** | **Beacon v2 Development** | **Related ...** |
|---|---|---|

**2014**   GA4GH founding event; Jim Ostell proposes Beacon concept including "more features ... version 2"

**2015**
- beacon-network.org aggregator created by DNAstack

- ELIXIR starts Beacon project support

**2016**
- Beacon v0.3 release
  work on queries for structural variants (brackets for fuzzy start and end parameters...)

- Beacon\* concept implemented on progenetix.org
- concepts from GA4GH Metadata (ontologies...)
- entity-scoped query parameters ("individual.age")

**2017**
- OpenAPI implementation
- integrating CNV parameters (e.g. "startMin, statMax")

- Beacon\* demos "handover" concept

- GA4GH re-structuring (workstreams...)
- Beacon part of Discovery WS

**2018**
- Beacon v0.4 release in January; feature release for GA4GH approval process
- GA4GH Beacon v1 approved at Oct plenary

- new Beacon website (March)

**2019**
- ELIXIR Beacon Network

- Beacon hackathon Stockholm; settling on "filters"
- Barcelona goes Zurich developers meeting
- Beacon API v2 Kick off
- adopting "handover" concept

- Beacon publication at Nature Biotechnology

**2020**

- "Scouts" teams working on different aspects - filters, genomic variants, compliance ...
- discussions w/ clinical stakeholders

**2021**

- framework + models concept implemented
- range and bracket queries, variant length parameters
- starting of GA4GH review process

- Phenopackets v2 approved

**2022**

- further changes esp. in default model, aligning with Phenopackets and VRS
- unified beacon-v2 code & docs repository
- Beacon v2 approved at Apr GA4GH Connect

- *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 Queries

## Implementation of Current Options

- (so far) the Beacon model does not define explicit query types

- disambiguation of parameters is left to implementers

- implicit query types:

  - ➡ allele/sequence query

  - ➡ range query, w/ or w/o additional parameters

  - ➡ bracket query (e.g. sized CNVs)

  - ➡ aminoacid, HGVS, gene

beaconplus.progenetix.org

---

**Beacon+**   Progenetix   Help

**Beacon Query Types**

| Sequence / Allele | CNV (Bracket) | Genomic Range | Aminoacid | Gene ID | HGVS | Sam |

**Dataset**

Test Database - examplez ✕                                    ✕  |  ⌄

**Chromosome** ℹ                          **Variant Type** ℹ

Select...                      ⌄         Select...                      ⌄

**Start or Position** ℹ

19000001-21975098

**Reference Base(s)** ℹ                   **Alternate Base(s)** ℹ

N                                        A

**Select Filters** ℹ

Select...                                                        ⌄

Query Database

**Form Utilities**     ⚙ Gene Spans     ⚙ Cytoband(s)

**Query Examples**     CNV Example   SNV Example   Range Example   Gene Match

Aminoacid Example   Identifier - HeLa

# Beacon Queries

## Range ("anything goes") Request

- defined through the use of 1 start, 1 end

- any variant... but can be limited by type etc.



**Beacon Range Query**

Matching variants in a region



17'600'000          18'600'001 - 18'650'000          19'650'000

GENE01

C > TT

**Bold: Matched Variants**

**Shaded: Unmatched**

DEL (Copy Number Loss)     DUP (Copy Number Gain)     SNP / INDEL ...     Unknown Annotation

---

**Beacon Query Types**

| Sequence / Allele | CNV (Bracket) | Genomic Range | Aminoacid | Gene ID | HGVS | Sam |

**Dataset**

Test Database - examplez  ✕                                      ✕ | ⌄

**Chromosome** ⓘ                          **Variant Type** ⓘ

17 (NC_000017.11)            ⌄          SO:0001059 (any sequence alteration - S...   ⌄

**Start or Position** ⓘ                    **End (Range or Structural Var.)** ⓘ

7572826                                    7579005

**Reference Base(s)** ⓘ                    **Alternate Base(s)** ⓘ

N                                          A

**Select Filters** ⓘ

Select...                                                          | ⌄

**Chromosome 17** ⓘ

7572826



7579005

**Query Database**

**Form Utilities**        ⚙ Gene Spans        ⚙ Cytoband(s)

**Query Examples**    CNV Example    SNV Example    Range Example    Gene Match

Aminoacid Example    Identifier - HeLa

As in the standard SNV query, this example shows a Beacon query against mutations in the `EIF4A1` gene in the DIPG childhood brain tumor dataset. However, this range + wildcard query will return any variant with alternate bases (indicated through "N"). Since parameters will be interpreted using an "AND" paradigm, either Alternate Bases OR Variant Type should be specified. The exact variants which were being found can be retrieved through the variant handover [H—>O] link.

# Beacon Queries

## Bracket ("CNV") Query

- defined through the use of 2 start, 2 end
- any contiguous variant...



**Beacon Bracket Query**

**Example for complete regional match**

GENE01

17'600'000 – 18'600'000     18'600'001 – 18'650'000     18'650'001 – 19'6500'000

Bold: Matched Variants

Shaded: Unmatched

**DEL (Copy Number Loss)**     **DUP (Copy Number Gain)**

---

**Beacon Query Types**

Sequence / Allele | **CNV (Bracket)** | Genomic Range | Aminoacid | Gene ID | HGVS | Sam

**Dataset**

Test Database - examplez ✕ | ✕ | ⌄

**Chromosome** ⓘ | **Variant Type** ⓘ

9 (NC_000009.12) | ⌄ | EFO:0030067 (copy number deletion) | ⌄

**Start or Position** ⓘ | **End (Range or Structural Var.)** ⓘ

21000001-21975098 | 21967753-23000000

**Select Filters** ⓘ

NCIT:C3058: Glioblastoma (100) ✕ | ✕ | ⌄

**Chromosome 9** ⓘ

21000001 21975098

21967753 23000000

**Query Database**

**Form Utilities** | ⚙ Gene Spans | ⚙ Cytoband(s)

**Query Examples** | CNV Example | SNV Example | Range Example | Gene Match

Aminoacid Example | Identifier - HeLa

This example shows the query for CNV deletion variants overlapping the CDKN2A gene's coding region with at least a single base, but limited to "focal" hits (here i.e. <= ~2Mbp in size). The query is against the examplez collection and can be modified e.g. through changing the position parameters or data source.

# Progenetix Stack

- JavaScript front-end is populated for query results using asynchronous access to multiple handover objects
  - ‣ biosamples and variants tables, CNV histogram, UCSC .bed loader, .pgxseg variant downloads…
- the complete middleware / CGI stack is provided through the *bycon* package
  - ‣ schemas, query stack, data transformation (e.g. Phenopackets generation)…
- data collections mostly correspond to the main Beacon default model entities
  - ‣ no separate *runs* collection; integrated w/ analyses
  - ‣ *variants* are stored per observation instance

- *collations* contain pre-computed data (e.g. CNV frequencies, statistics) and information for all grouping entity instances and correspond to **filter values**
  - ‣ PMID:10027410, NCIT:C3222, pgx:cohort-TCGA, pgx:icdom-94703…
- *querybuffer* stores id values of all entities matched by a query and provides the corresponding access handle for **handover** generation



```
_id: ObjectId("6249bb654f8f8d67eb94953b"),
id: '0765ee26-5029-4f28-b01d-9759abf5bf14',
source_collection: 'variants',
source_db: 'progenetix',
source_key: '_id',
target_collection: 'variants',
target_count: 667,
target_key: '_id',
target_values: [
  ObjectId("5bab578b727983b2e00ca99e"),
  ObjectId("5bab578d727983b2e00cb505")
```

variants  analyses  biosamples  individuals

collations  geolocs  genespans  publications  qBuffer

**Entity collections**

**Utility collections**

## byconaut — progenetix / byconaut

Public

Edit Pins | Unwatch 2 | Fork 1 | Star 0

main | 2 branches

mbaudis get_plot_parameters

- bin
- docs
- exports
- imports
- local
- rsrc
- services
- tmp
- .gitignore
- LICENSE
- README.md
- __init__.py
- install.py
- install.yaml
- mkdocs.yaml

## beaconplus-web — progenetix / beaconplus-web

Public

forked from progenetix/progenetix-web

Edit Pins | Watch 0 | Fork 3 | Star 0

main | 1 branch | 0 tags

This branch is 44 commits ahead, 24 commits behind progenetix:main.

mbaudis code cleaning, no feature changes

| .github/workflows | cleanup |
| docs | still first implementation clean-up |
| extra | documentation |
| public | graphic refinement |
| src | code cleaning, no feature changes |
| .babelrc | Simplify query generation and add |
| .env.development | first working version |
| .env.local | first working version |
| .env.production | env |
| .env.staging | env |
| .eslintrc.json | BioSubsetsPage perf optimisations |

## bycon — progenetix / bycon

Public

Edit Pins | Unwatch 4 | Fork 6 | Starred 5

main | 4 branches | 25 tags

Go to file | Add file | Code

mbaudis 1.3.6 ... | be19a12 3 days ago | 852 commits

| .github/workflows | Create mk-bycon-docs.yaml | 8 months ago |
| bycon | 1.3.6 | 3 days ago |
| docs | 1.3.6 | 3 days ago |
| local | 1.3.5 preparation | 2 weeks ago |
| .gitignore | Update .gitignore | 3 months ago |
| LICENSE | Create LICENSE | 3 years ago |
| MANIFEST.in | major library & install disentanglement | 9 months ago |
| README.md | #### 2023-07-23 (v1.0.68) | 4 months ago |
| install.py | 1.3.6 | 3 days ago |
| install.yaml | v1.0.57 | 5 months ago |
| mkdocs.yaml | 1.1.6 | 3 months ago |
| requirements.txt | 1.3.6 | 3 days ago |
| setup.cfg | ... | 10 months ago |
| setup.py | 1.3.6 | 3 days ago |
| updev.sh | 1.3.6 | 3 days ago |

### About

Bycon - A Python Based Beacon API (beacon-project.io) implementation leveraging the Progenetix (progenetix.org) data model

- Readme
- CC0-1.0 license
- Activity
- 5 stars
- 4 watching
- 6 forks

Report repository

### Releases

25 tags

Create a new release

### Packages

No packages published
Publish your first package

# pgxRpi

## An interface API for analyzing Progenetix CNV data in R using the Beacon+ API

GitHub: https://github.com/progenetix/pgxRpi

Bioconductor

**README.md**

## pgxRpi

Welcome to our R wrapper package for Progenetix REST API that leverages the capabilities of Beacon v2 specification. Please note that a stable internet connection is required for the query functionality. This package is aimed to simplify the process of accessing oncogenomic data from Progenetix database.

You can install this package from GitHub using:

```
install.packages("devtools")
devtools::install_github("progenetix/pgxRpi")
```

For accessing metadata of biosamples/individuals, or learning more about filters, get started from the vignette Introduction_1_loadmetadata.

For accessing CNV variant data, get started from this vignette Introduction_2_loadvariants.

For accessing CNV frequency data, get started from this vignette Introduction_3_loadfrequency.

For processing local pgxseg files, get started from this vignette Introduction_4_process_pgxseg.

If you encounter problems, try to reinstall the latest version. If reinstallation doesn't help, please contact us.

## pgxRpi

| platforms | all | rank | 2218 / 2221 | support | 0 / 0 | in Bioc | devel only |
| build | ok | updated | < 1 month | dependencies | 144 | | |

DOI: 10.18129/B9.bioc.pgxRpi
This is the **development** version of pgxRpi; to use it, please install the devel version of Bioconductor.

## R wrapper for Progenetix

Bioconductor version: Development (3.19)

The package is an R wrapper for Progenetix REST API built upon the Beacon v2 protocol. Its purpose is to provide a seamless way for retrieving genomic data from Progenetix database—an open resource dedicated to curated oncogenomic profiles. Empowered by this package, users can effortlessly access and visualize data from Progenetix.

Author: Hangjia Zhao [aut, cre] (iD), Michael Baudis [aut] (iD)

Maintainer: Hangjia Zhao <hangjia.zhao at uzh.ch>

Citation (from within R, enter `citation("pgxRpi")`):

# 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*")

- **support** and/or get involved with international **data standards** efforts and projects

➡️ **Collaborate!**



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

```
CDKN2A:DEL
size<1Mb
granularity:record
ncit:C3058
DUO:0000004
HP:0003621
```

# 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*")

- **support** and/or get involved with international **data standards** efforts and projects

➡ **Collaborate!**



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

**"Internet of Genomics"**

```
CDKN2A:DEL
size<1Mb
granularity:record
ncit:C3058
DUO:0000004
HP:0003621
```

# The Beacon team through the ages

**European Genome-Phenome Archive**

**CRG** Centre for Genomic Regulation

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

**CSC**

Juha Törnroos
Teemu Kataja
Ilkka Lappalainen
Dylan Spalding

**University of Leicester**

**Cafe Variome Central**

**Tony Brookes**
**Tim Beck**
Colin Veal
Tom Shorter

**SPHN** Swiss Personalized Health Network

**University of Zurich** UZH

**Michael Baudis**
Rahel Paloots
Hangjia Zhao
Ziying Yang
Bo Gao
Qingyao Huang

**Genomics england**

**Augusto Rendon**
**Ignacio Medina**
Javier López
Jacobo Coll
Antonio Rueda

**cnag** centre nacional d'anàlisi genòmica / centro nacional de análisis genómico

**Sergi Beltran**
Carles Hernandez

**Inserm** Institut national de la santé et de la recherche médicale

David Salgado

**BSC Barcelona Supercomputing Center** Centro Nacional de Supercomputación

**Salvador Capella**
Dmitry Repchevski
JM Fernández

**DisGeNET**

**Laura Furlong**
Janet Piñero

**elixir** **B1MG**

**Serena Scollen**
Gary Saunders
Giselle Kerry
David Lloyd

**H3Africa** Human Heredity & Health in Africa

**Nicola Mulder**
Mamana
Mbiyavanga
Ziyaad Parker

**EU Can CAN.**
**David Torrents**

**AUTISM SPEAKS**

**Dean Hartley**

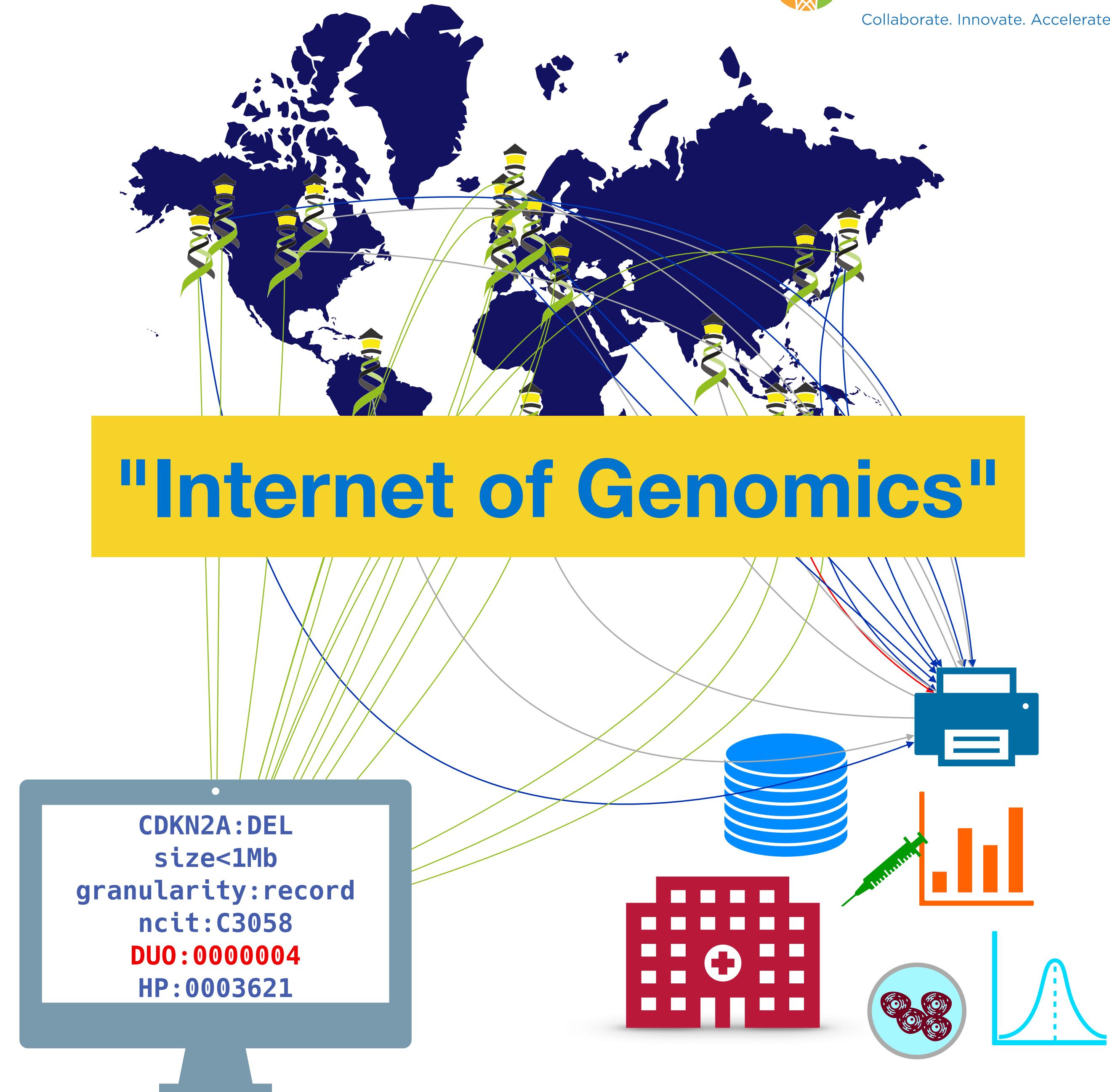**Junta de Andalucía** Fundación Progreso y Salud **CONSEJERÍA DE SALUD**

**Joaquin Dopazo**
Javier Pérez
J.L. Fernández
Gema Roldan

**CINECA**

**Thomas Keane**
Melanie Courtot
Jonathan Dursi

**Heidi Rehm**
Ben Hutton

Toshiaki
Katayama

**GEM Japan**

**Stephane Dyke**

**DNASTACK**

**Marc Fiume**
Miro Cupak

**BRCA EXCHANGE**

**Melissa Cline**

**ENA**

**EMBL-EBI**

Diana Lemos

**European Joint Programme RARE DISEASES**

**VICC** Variant Interpretation for Cancer Consortium

**GA4GH Phenopackets**
Peter Robinson
Jules Jacobsen

**GA4GH VRS**
Alex Wagner
Reece Hart

**Beacon PRC**
Alex Wagner
Jonathan Dursi
Mamana Mbiyavanga

Alice Mann
Neerjah Skantharajah

**elixir**

# ELIXIR hCNV Community

https://cnvar.org/

## ELIXIR Human Copy Number Variation community

**h-CNV Community**
Homepage & News
About ...
h-CNV Projects
CNV Annotation Standards
Databases & Resources
CNV References Project
Contacts
Genome Blog
h-CNV @ ELIXIR
Beacon Project

Among the different types of inherited and acquired genomic variants, regional genomic copy number variations (CNV) contribute - if measured by affected genomic sequences - contribute by far the largest amount of genomic changes, contributing both to many syndromic diseases as well as the vast majority of human cancers. The website of the *Human Copy Number Variation Community* (hCNV) is a resource originated in ELIXIR's h-CNV Community Implementation Study (2019-2021) with the aim to provide a resource hub and knowledge exchange space for scientists and practitioners working with - or being interested in - genomic copy number variations in health and diseases. However, the scope of the community extends beyond CNVs and includes definition of and work with other types of genomic variations with a focus on structural variants.

🇺🇦 Leaflet | Map data © OpenStreetMap

---

## CNV Annotation Formats

### CNV Term Use Comparison in Computational (File/Schema) Formats

This table is maintained in parallel with the Beacon v2 documentation.

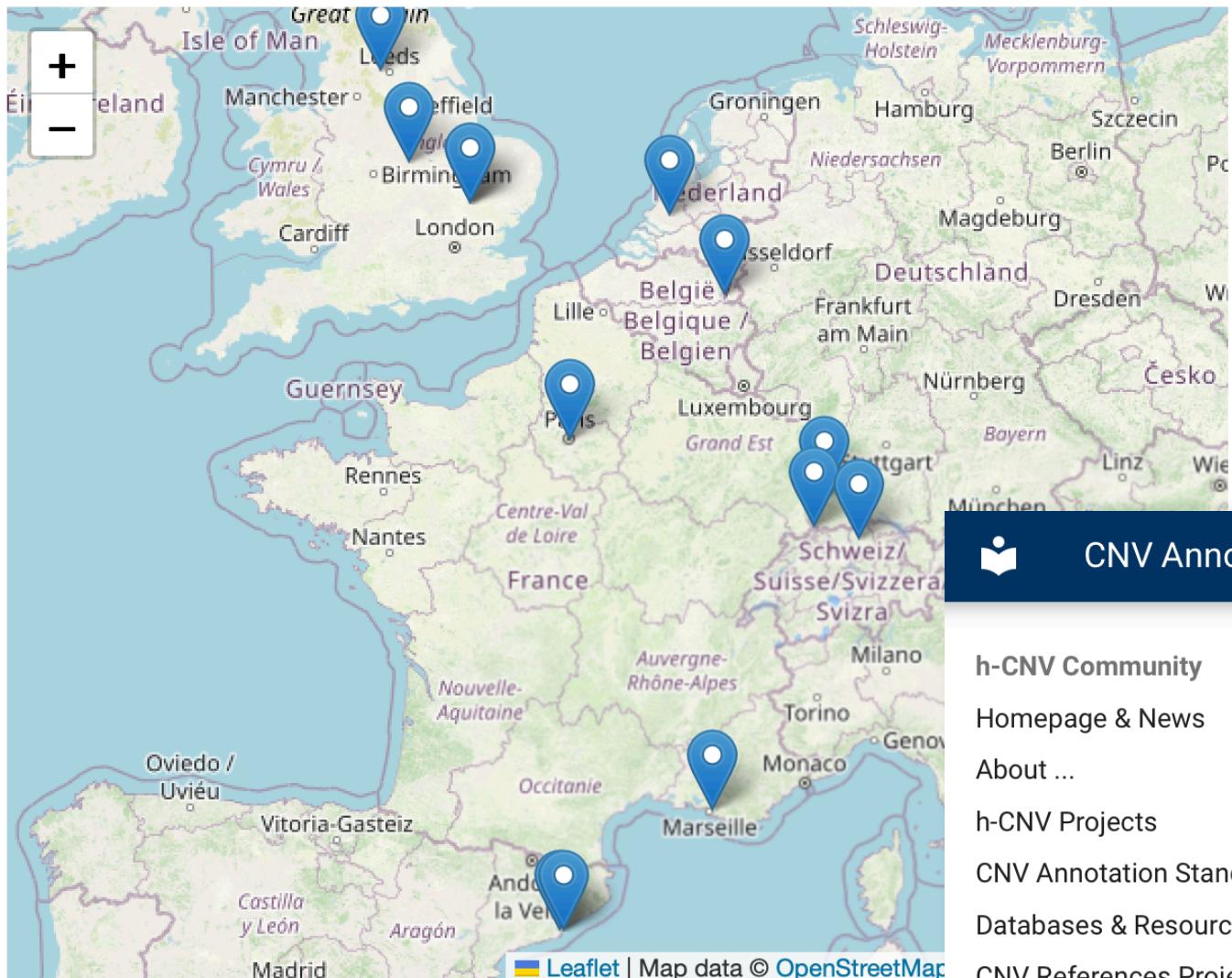| EFO | Beacon | VCF | SO | GA4GH VRS[1] | Notes |
|-----|--------|-----|-----|----------|-------|
| EFO:0030070 copy number gain | DUP[2] or EFO:0030070 | DUP SVCLAIM=D[3] | SO:0001742 copy_number_gain | EFO:0030070 gain | a sequence alteration whereby the copy number of a given genomic region is greater than the reference sequence |
| EFO:0030071 low-level copy number gain | DUP[2] or EFO:0030071 | DUP SVCLAIM=D[3] | SO:0001742 copy_number_gain | EFO:0030071 low-level gain | |
| EFO:0030072 high-level copy number gain | DUP[2] or EFO:0030072 | DUP SVCLAIM=D[3] | SO:0001742 copy_number_gain | EFO:0030072 high-level gain | commonly but not consistently used for >=5 copies on a bi-allelic genome region |
| EFO:0030073 focal genome amplification | DUP[2] or EFO:0030073 | DUP SVCLAIM=D[3] | SO:0001742 copy_number_gain | EFO:0030072 high-level gain[4] | commonly but not consistently used for >=5 copies on a bi-allelic genome region, of limited size (operationally max. 1-5Mb) |
| EFO:0030067 copy number loss | DEL[2] or EFO:0030067 | DEL SVCLAIM=D[3] | SO:0001743 copy_number_loss | EFO:0030067 loss | a sequence alteration whereby the copy number of a given genomic region is smaller than the reference sequence |
| EFO:0030068 low-level copy number loss | DEL[2] or EFO:0030068 | DEL SVCLAIM=D[3] | SO:0001743 copy_number_loss | EFO:0030068 low-level loss | |
| EFO:0020073 high-level copy number loss | DEL[2] or EFO:0020073 | DEL SVCLAIM=D[3] | SO:0001743 copy_number_loss | EFO:0020073 high-level loss | a loss of several copies; also used in cases where a complete genomic deletion cannot be asserted |

isplay a menu

ELIXIR hCNV Community

www.elixir-europe.org

# CNV Term Use Comparison

## in computational (file/schema) formats

| EFO | Beacon | VCF | SO | GA4GH VRS1.3 |
|---|---|---|---|---|
| **EFO:0030070**<br>copy number gain | DUP or **EFO:0030070** | DUP<br>SVCLAIM=D | SO:0001742<br>copy_number_gain | **EFO:0030070**<br>gain |
| **EFO:0030071**<br>low-level copy number gain | DUP or **EFO:0030071** | DUP<br>SVCLAIM=D | SO:0001742<br>copy_number_gain | **EFO:0030071**<br>low-level gain |
| **EFO:0030072**<br>high-level copy number gain | DUP or **EFO:0030072** | DUP<br>SVCLAIM=D | SO:0001742<br>copy_number_gain | **EFO:0030072**<br>high-level gain |
| EFO:0030073<br>focal genome amplification | DUP or EFO:0030073 | DUP<br>SVCLAIM=D | SO:0001742<br>copy_number_gain | **EFO:0030072**<br>high-level gain |
| **EFO:0030067**<br>copy number loss | DEL or **EFO:0030067** | DEL<br>SVCLAIM=D | SO:0001743<br>copy_number_loss | **EFO:0030067**<br>loss |
| **EFO:0030068**<br>low-level copy number loss | DEL or **EFO:0030068** | DEL<br>SVCLAIM=D | SO:0001743<br>copy_number_loss | **EFO:0030068**<br>low-level loss |
| **EFO:0020073**<br>high-level copy number loss | DEL or **EFO:0020073** | DEL<br>SVCLAIM=D | SO:0001743<br>copy_number_loss | **EFO:0020073**<br>high-level loss |
| **EFO:0030069**<br>complete genomic deletion | DEL or **EFO:0030069** | DEL<br>SVCLAIM=D | SO:0001743<br>copy_number_loss | **EFO:0030069**<br>complete genomic loss |

# GA4GH Ascona Connect

**REGISTER FOR THIS EVENT** ⧉

This hybrid working meeting aims to support GA4GH contributors in advancing product development and gathering feedback on needs.



**Image summary:** Join us for GA4GH Connect from 21 to 24 April 2024.

*https://www.ga4gh.org/event/ga4gh-connect-2/*

# Genomic Data & Privacy

**Risks & opportunities**

# Gattaca (1997)

**A genetically inferior man assumes the identity of a superior one in order to pursue his lifelong dream of space travel.**

- genetic determinism

  ‣ main character has been determined to be unsuitable for complex jobs based on genetic analysis

- genetic identification

  ‣ the use of genetic sampling for personal identification is daily routine

With information from https://www.imdb.com/title/tt0119177/

**17 : 7577121 G > A**

## Beacon

A ***Beacon*** answers a query for a specific genome variant against individual or aggregate genome collections
**YES** | **NO** | **\0**

# Genome *Beacons* Compromise Security?

Querying for thousands of specific SNV occurrences in a genomic data pool can identify individuals in an anonymized genomic data collection

## Stanford researchers identify potential security hole in genomic data-sharing network

Hackers with access to a person's genome might find out if that genome is in an international network of disease databases.

OCT 29 2015

Sharing genomic information among researchers is critical to the advance of biomedical research. Yet genomic data contains identifiable information and, in the wrong hands, poses a risk to individual privacy. If someone had access to your genome sequence — either directly from your saliva or other tissues, or from a popular genomic information service — they could check to see if you appear in a database of people with certain medical conditions, such as heart disease, lung cancer or autism.

Work by a pair of researchers at the Stanford University School of Medicine makes that genomic data more secure. Suyash Shringarpure, PhD, a postdoctoral scholar in genetics, and Carlos Bustamante, PhD, a professor of genetics, have demonstrated a technique for hacking a network of global genomic databases and how to prevent it. They are working with investigators from the Global Alliance for Genomics and Health on implementing preventive measures.

The work, published Oct. 29 in *The American Journal of Human Genetics,* also bears importantly on the larger question of how to analyze mixtures of genomes, such as those from different people at a crime scene.
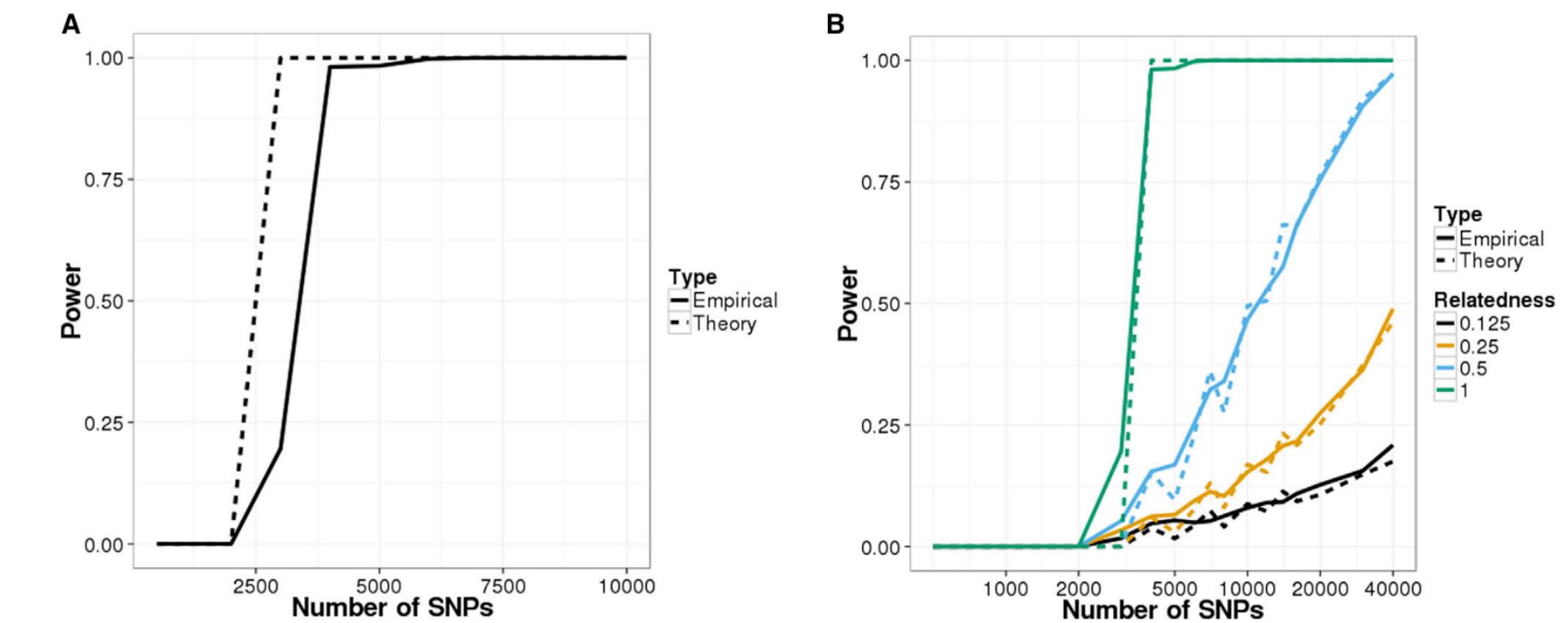
Stanford researchers are working with the Global Alliance for Genomics and Health to make genomic information in the Beacon Project more secure.
*Science photo/Shutterstock*

# IDENTIFICATION OF INDIVIDUALS FROM MIXED COLLECTIONS USING RARE ALLELES

## Privacy Risks from Genomic Data-Sharing Beacons

Suyash S. Shringarpure[1,*] and Carlos D. Bustamante[1,*]

The human genetics community needs robust protocols that enable secure sharing of genomic data from participants in genetic research. Beacons are web servers that answer allele-presence queries—such as "Do you have a genome that has a specific nucleotide (e.g., A) at a specific genomic position (e.g., position 11,272 on chromosome 1)?"—with either "yes" or "no." Here, we show that individuals in a beacon are susceptible to re-identification even if the only data shared include presence or absence information about alleles in a beacon. Specifically, we propose a likelihood-ratio test of whether a given individual is present in a given genetic beacon. Our test is not dependent on allele frequencies and is the most powerful test for a specified false-positive rate. Through simulations, we showed that in a beacon with 1,000 individuals, re-identification is possible with just 5,000 queries. Relatives can also be identified in the beacon. Re-identification is possible even in the presence of sequencing errors and variant-calling differences. In a beacon constructed with 65 European individuals from the 1000 Genomes Project, we demonstrated that it is possible to detect membership in the beacon with just 250 SNPs. With just 1,000 SNP queries, we were able to detect the presence of an individual genome from the Personal Genome Project in an existing beacon. Our results show that beacons can disclose membership and implied phenotypic information about participants and do not protect privacy a priori. We discuss risk mitigation through policies and standards such as not allowing anonymous pings of genetic beacons and requiring minimum beacon sizes.



**Figure 1. Power of Re-identification Attacks on Beacons Constructed with Simulated Data**
Power curves for the likelihood-ratio test (LRT) on (A) a simulated beacon with 1,000 individuals and (B) detecting relatives in the simulated beacon. The false-positive rate was set to 0.05 for all scenarios.

▸ rare allelic variants can be used to identify an individual (or her relatives) in a genome collection without having access to individual datasets

▸ however, such an approach requires previous knowledge about the individual's SNPs

"We're an information economy. They teach you that in school. What they don't tell you is that it's impossible to move, to live, to operate at any level without leaving traces, bits, seemingly meaningless fragments of personal information. Fragments that can be retrieved, amplified . . ."

**–William Gibson in "Johnny Mnemonic" (1986)**

# Phenotyping from DNA

## From DNA to "Wanted" Posters?



**Paragon Nanolabs Inc.**
**The Snapshot DNA Phenotyping Service**

- association of genomic variants with phenotypic data collection

- while hair, eye color are easy targets not useful for relevant phenotypic features especially if large environmental component

- huge biases based on input/collection data

- Belgium and Germany do not allow forensic DNA phenotyping

- Switzerland: Bundesrat decision on 2020-12-04 to allow phenotyping for law enforcement purposes

DNA → Genotype of Unknown Contributor

| $S_1$ | $S_2$ | $S_3$ | $S_4$ | $S_5$ | $S_6$ | $S_7$ | ... | $S_{900K}$ |
|---|---|---|---|---|---|---|---|---|
| AA | AT | CC | GG | CG | AA | TT | | CC |

**Model #1: Skin Color**
$(2.4) \cdot S_2 + (-1.7) \cdot S_5 + (0.6) \cdot S_{12}$

**Model #2: Eye Color**
$(5.3) \cdot S_{16} + (3.6) \cdot S_{21} + (-7.1) \cdot S_{35}$

**Model #3: Hair Color**
$(7.4) \cdot S_{12} + (4.3) \cdot S_5 + (1.4) \cdot S_{16}$

Snapshot Models

"When the New York Times ran an informal test of the Parabon system with one of its reporters, it failed badly." (ACLU.org)

https://snapshot.parabon-nanolabs.com/phenotyping

# Long-Range Familial Searches



Hi Michael,

Good news! We've discovered new DNA Matches for you.

- Commercial, "Direct to Customer" DNA analyses are provided through independent sites and such affiliated to genealogy services (MyHeritage, Ancestry.com, 23andMe...)
- Genealogy sites identify individuals with matching haplotype blocks & provide a prediction about degree of genetic relation
- Law enforcement agencies (and who else?!) can send individual SNP profiles (e.g. recovered from evidence many years after a crime) using a *Jane Doe* identity, to identify relatives of the suspect - **long range familial search**



Devaughn had never been a suspect until genetic genealogy put police on his trail several months ago. Earlier this year, police sent the DNA profile to Parabon, a private genetics company, to compare the suspect's DNA sample to a public genealogy DNA database looking for people with similar DNA profiles who might be kin to the suspect. That eventually led authorities to look at Devaughn.

## Rienzi man charged with 1990 Starkville murder

By William Moore Daily Journal   15 hrs ago   Comments

© Copyright 2018 Daily Journal, 1242 S Green St Tupelo, MS



### The New York Times

## *How a Genealogy Site Led to the Front Door of the Golden State Killer Suspect*

Investigators used DNA from crime scenes that had been stored all these years and plugged the genetic profile of the suspected assailant into an online genealogy database. One such service, GEDmatch, said in a statement on Friday that law enforcement officials had used its database to crack the case. Officers found distant relatives of Mr. DeAngelo's and, despite his years of eluding the authorities, traced their DNA to his front door.

The New York Times, April 26, 2018

But genotyping itself is for professional labs, right?

# Rapid re-identification of human samples

...

We developed a rapid, inexpensive, and portable strategy to re-identify human DNA using the MinION. Our strategy requires only ~60 min preparation and 5-30 minutes of MinION sequencing, works with low input DNA, and enables familial searches using Direct-to-Consumer genomic reference datasets. This method can be implemented in a variety of fields:

## Forensics

Identification of abandoned meterial using DNA fingerprinting is a common practice. The main challange currently being: time. Our method allows rapid sample preparation at the crime scene (see movie). We envision that the method can be adopted in the field for rapid checks, after a mass disaster, and can be adopted in border control to fight human traffacking.

## Clinic

Clinics procces many samples, either for analysis or, for example, organ donations. These samples are DNA fingerprinted to prevent sample mix-up mistakes. Our method can be implemnted in the clinic for rapid sanitiy-check of all incoming samples.

## Cell line identification

Cross contamination of cell lines in science is a major problem. It results in unreproducible data, and clinical trails based on inaccurate findings. This problem costs billions of dollars per year. We envision labs can adopt our identification method to ensure the purity of the cell line, and detect contamination.

**The MinION** (Oxford Nanopore)
Source: Sophie Zaaijer

# DEMOCRATIZING DNA FINGERPRINTING

**Sophie Zaaijer, Assaf Gordon, Robert Piccone, Daniel Speyer, Yaniv Erlich, 2016**
*ddf.teamerlich.org*



**MinION by Oxford Nanopore Technologies**

The MinION is the smallest DNA sequencer currently around. Its the size of a Mars bar, and can be simply plugged into a laptop with a USB3.0 port.

For more information about the MinION please click:

Oxford Nanopore Technologies

**Bento Lab**

The Bento lab is a miniature lab with a centrifuge, thermocycler and a electrophoresis compartment.

For more information about the Bento-lab please click:

Bento Lab

Data can be loaded into the person ID pipeline
matches inferred between 3-30 minutes

DNA sequencing for identification/fingerprinting soon "commodity" technology (in contrast with technological/data challenges in "precision medicine")

Generalkonsent

BENEFIT

BLOCKCHAIN

HEALTH

PRIVACY

SECURITY

CONSENT

ACCESS

Right to Research

HACKERS

LAWS

**G**enetic
**I**nformation
**N**ondiscrimination
**A**ct

**H**ealth
**I**nsurance
**P**ortability and
**A**ccountability
**A**ct

SAFETY

CRYPTOGRAPHY

# Share *YOUR* Genome data?

- The Beacon concept - balanced approach for accessing genome variant data from internationally distributed resources

- However: Genome data has the inherent "risk" of being identified and linked to a person

**Solutions from Technology or Society? Discourse!**

TECHNOLOGY

## How can a DNA firm lose half its users' data to 'Jew-hating' hackers?

Dark-web criminals cited the head of 23andMe's faith after a raid on the details of 6.9 million people — including her Google-founding ex. Now the lawsuits are coming

FAMILY MATTERS

## Hackers stole ancestry data of 6.9 million users, 23andMe finally confirmed

Majority of impacted users are now being notified

ASHLEY BELANGER · 12/4/2023, 11:48 PM

# Find out what your DNA says about you and your family.

- See how your DNA breaks out across 31 populations worldwide
- Discover DNA relatives from around the world
- Share reports with family and friends

**order now**  USD$99

It has now been confirmed that an additional **6.9 million 23andMe users had ancestry data stolen** after hackers accessed thousands of accounts by likely reusing previously leaked passwords.

... Wired estimated that "at least a million data points from 23andMe accounts" that were "exclusively about Ashkenazi Jews" and data points from "hundreds of thousands of users of Chinese descent" seemed to be exposed.

... a spokesperson to confirm that two groups of ... opted into the **DNA Relatives feature** had their ... a stolen.

... scribes the DNA Relatives feature as ... u to find and connect with genetic relatives and ... about your family." By **opting in**, users hope to ... ly members by **willingly** giving others access to ... like their birth year, current location, and ... ames and birth locations. Users can opt out at ...

... about 5.5 million users, was hacked after opting in to automatically sharing information with DNA Relatives, including their "**name, birth year, relationship labels, the percentage of DNA shared with relatives, ancestry reports, and self-reported location**," TechCrunch reported.

... about 1.4 million users, shared "Family Tree profile information" ... including display names, relationship labels, birth year, and self-reported location, TechCrunch reported.

TECHNOLOGY

# How can a DNA firm lose half its users' data to 'Jew-hating' hackers?

Dark-web criminals cited the head of 23andMe's faith after a raid on the details of 6.9 million people — including her Google-founding ex. Now the lawsuits are coming

FAMILY MATTERS —

# Hackers stole ancestry data of 6.9 million users, 23andMe finally confirmed

Majority of impacted users are now being notified.

ASHLEY BELANGER - 12/4/2023, 11:48 PM

ars TECHNICA


Bloomberg / Contributor | Bloomberg

It has now been confirmed that an additional **6.9 million 23andMe users had ancestry data stolen** after hackers accessed thousands of accounts by likely reusing previously leaked passwords.

... Wired estimated that "at least a million data points from 23andMe accounts" that were "exclusively about Ashkenazi Jews" and data points from "hundreds of thousands of users of Chinese descent" seemed to be exposed.

... prompting a spokesperson to confirm that two groups of users who opted into the **DNA Relatives feature** had their personal data stolen.

23andMe describes the DNA Relatives feature as ... "allowing you to find and connect with genetic relatives and learn more about your family." By **opting in**, users hope to find lost family members by **willingly** giving others access to information like their birth year, current location, and ancestors' names and birth locations. Users can opt out at any time ...

... about 5.5 million users, was hacked after opting in to automatically sharing information with DNA Relatives, including their "**name, birth year, relationship labels, the percentage of DNA shared with relatives, ancestry reports, and self-reported location**," TechCrunch reported.

... about 1.4 million users, shared "Family Tree profile information" ... including display names, relationship labels, birth year, and self-reported location, TechCrunch reported.

TECHNOLOGY

## How can a DNA firm ~~hand~~ ~~sell~~ it~~s~~
users' data to 'Jew-ha~~ters~~'

Dark-web criminals cited the head of 23a~~ndMe~~
the details of 6.9 million people — includ~~ing~~
Now the lawsuits are ~~...~~

FAMILY MATTERS

## Hackers sto~~le~~
users, 23an~~dMe~~

Majority of impacted user~~s~~

ASHLEY BELANGER · 12/4/2023, 11:48 ~~...~~

# THE WALL STREET JOURNAL.

SIGN IN    SUBSCRIBE



# 23andMe's Fall From $6 Billion to Nearly $0

From celebrity 'spit parties' to a drop in the bucket: The once-hot DNA-testing company is struggling to profit

Anne Wojcicki of 23andMe, center, remotely rang the Nasdaq opening bell the day the company went public in 2021. PETER DASILVA/REUTERS

By *Rolfe Winkler* [Follow]
Jan. 31, 2024 at 5:30 am ET

It has now been confirmed that an additional **6.9 million**
~~ance~~**ry data stolen** after hackers
~~acco~~unts by likely reusing previously

~~at le~~ast a million data points from
~~were~~ "exclusively about Ashkenazi
~~and~~ hundreds of thousands of users
~~also~~ to be exposed.

~~went on~~ to confirm that two groups of
~~the DN~~A Relatives feature had their

~~DNA~~ Relatives feature as ...
~~to con~~nect with genetic relatives and
~~share.~~" By **opting in**, users hope to
~~are~~ **willingly** giving others access to
~~birth ye~~ar, current location, and
~~past lo~~cations. Users can opt out at

~~wa~~s hacked after opting in to
~~inform~~ation with DNA Relatives,
~~birth y~~ear, **relationship labels, the**
~~shared~~ **with relatives, ancestry**
~~rted l~~**ocation,**" TechCrunch reported.
~~also sh~~ared "Family Tree profile
~~informa~~tion" — including display names, relationship labels,
birth year, and self-reported location, TechCrunch reported.

# Human Rights Foundation

**Global Alliance**
for Genomics & Health

## Universal Declaration of Human Rights (1948)

### 27(1) "The Right to Science"

"Everyone has the right freely to participate in the cultural life of the community, to enjoy the arts and **to share in scientific advancement and its benefits.**"

### 27(2) "The Right to Recognition"

"Everyone has the right to the protection of the moral and material interests resulting from any scientific, literary or artistic production of which he is the author."