

Theoretical Cytogenetics and Oncogenomics

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

Michael Baudis

Professor of Bioinformatics

University of Zürich

Swiss Institute of Bioinformatics **SIB**

Member GA4GH Strategic Leadership Committee

Co-lead ELIXIR Beacon API Development

Co-lead ELIXIR hCNV Community



Universität
Zürich^{UZH}



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



SIB
Swiss Institute of
Bioinformatics





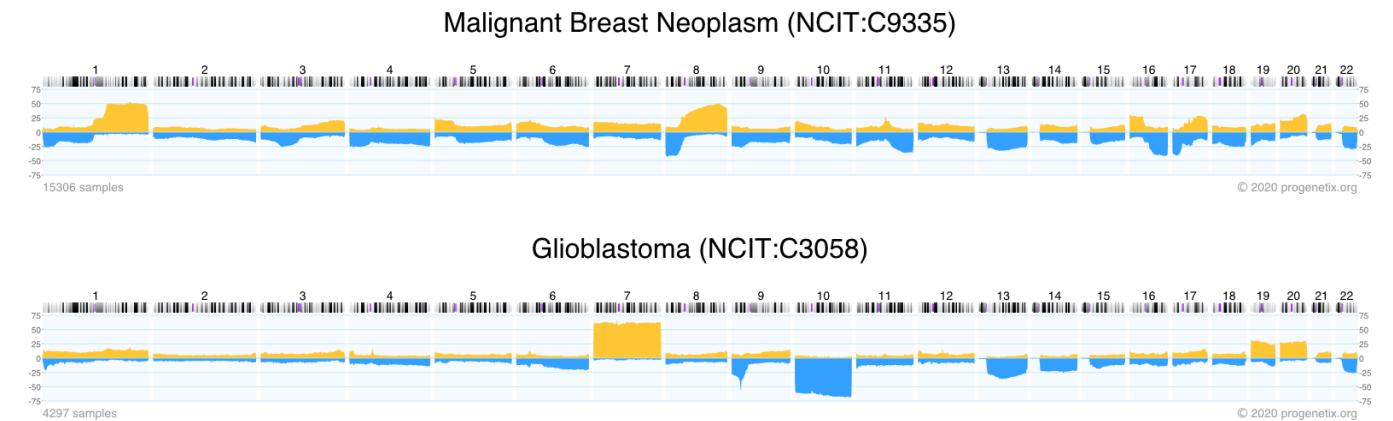
**University of
Zurich^{UZH}**
Department of Molecular Life Sciences



Theoretical Cytogenetics and Oncogenomics

... but what does this entail @baudisgroup?

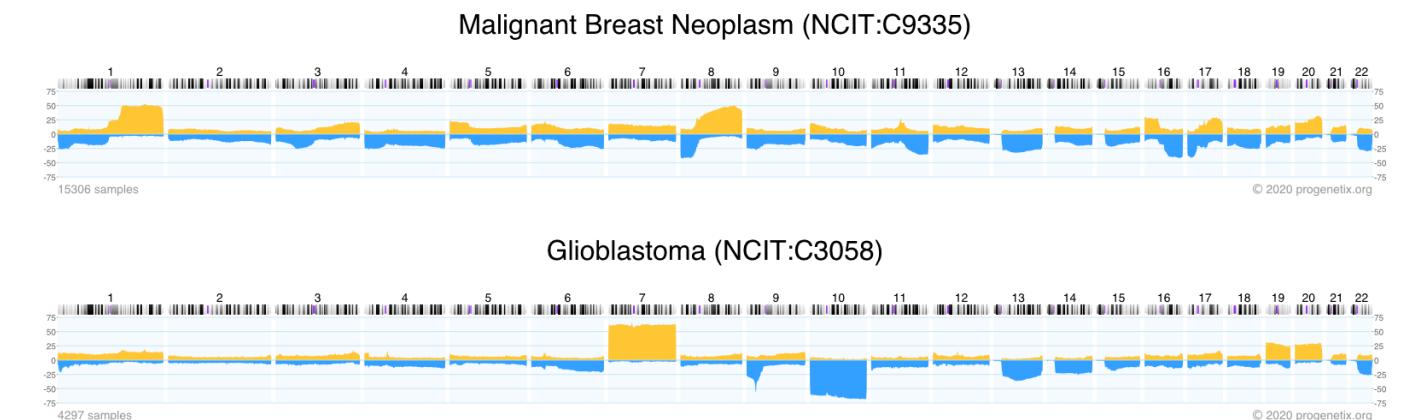
- **patterns & markers** in cancer genomics, especially somatic structural genome variants
- bioinformatics support in **collaborative** studies
- **reference resources** for curated cancer genome variations
- bioinformatics **tools & methods**
- **standards** and **reference implementations** for data sharing in genomics and personalized health
- **open research data** "ambassadoring"



Theoretical Cytogenetics and Oncogenomics

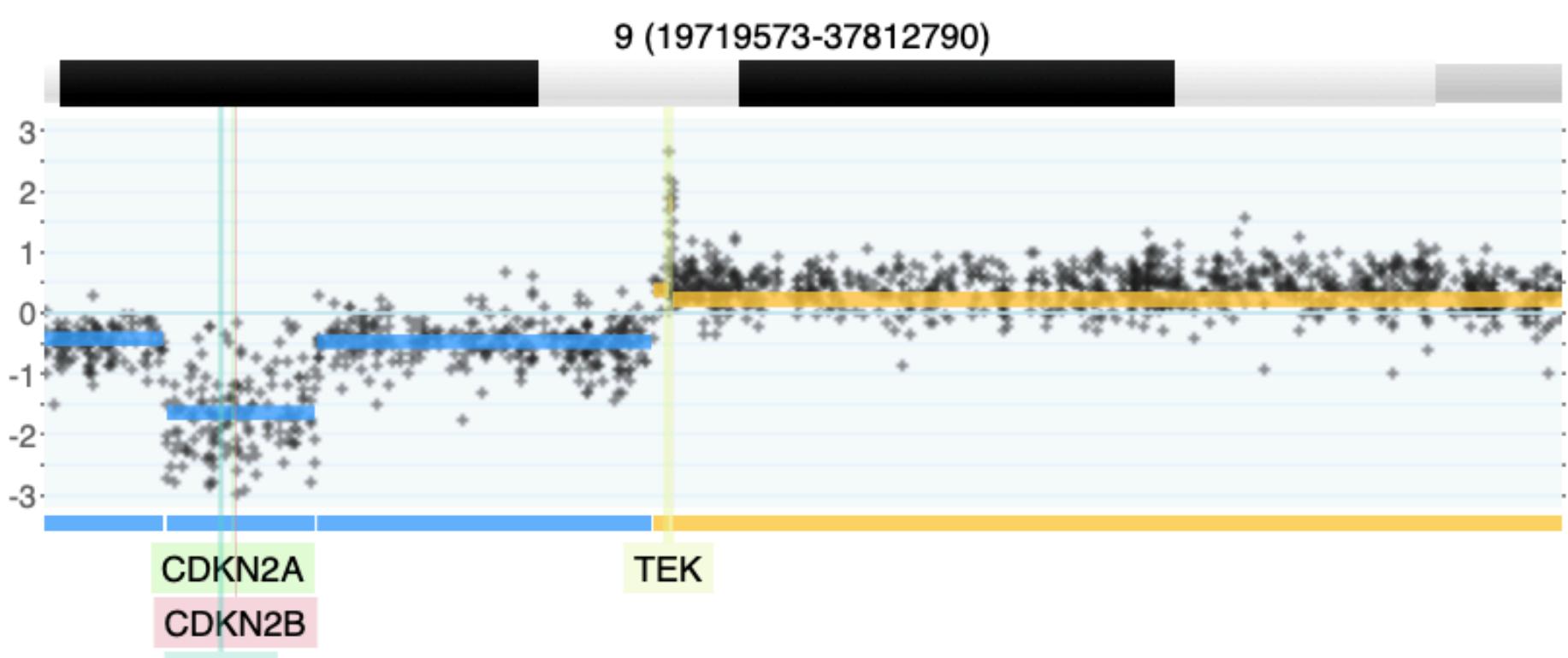
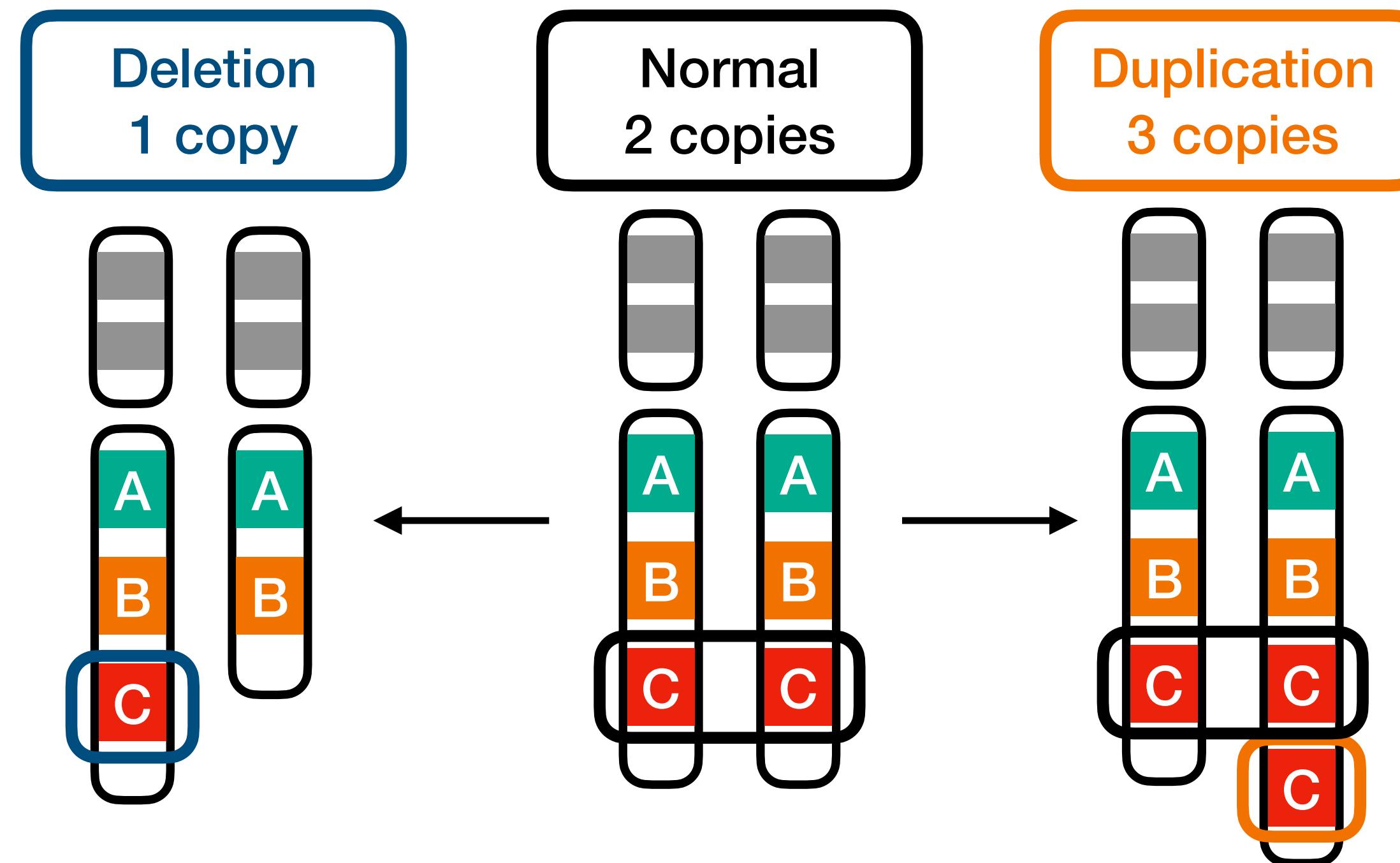
... but what does this entail @baudisgroup?

- **patterns & markers** in cancer genomics, especially somatic structural genome variants
- bioinformatics support in **collaborative** studies
- **reference resources** for curated cancer genome variations
- bioinformatics **tools & methods**
- **standards** and **reference implementations** for data sharing in genomics and personalized health
- **open research data** "ambassadoring"



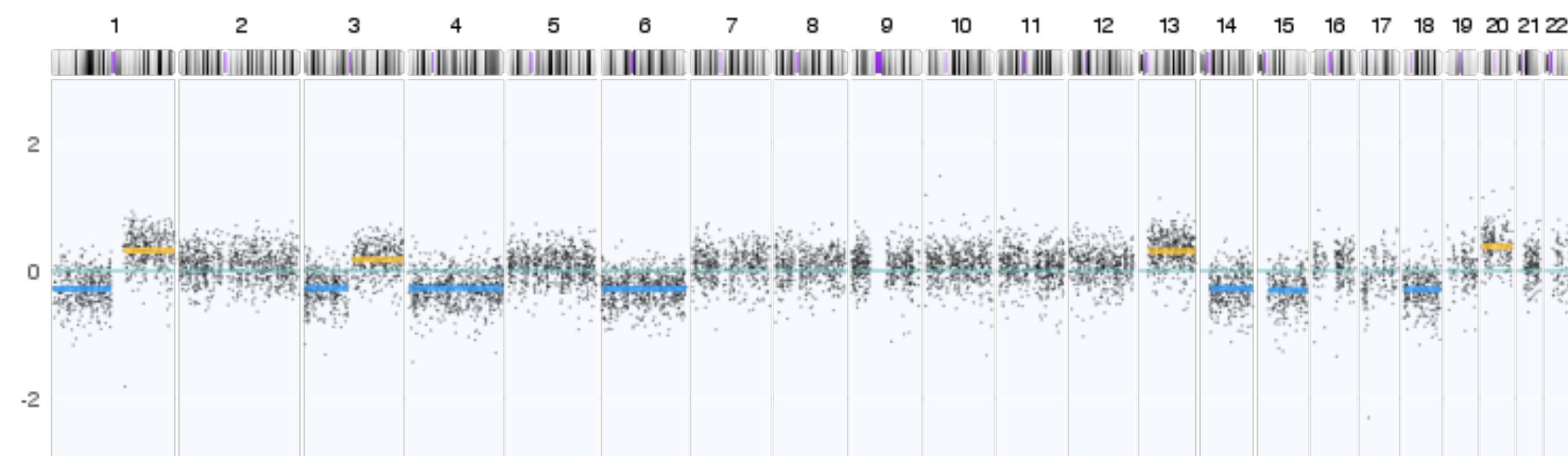
Curators
~~Data Parasites~~

Copy Number Variant (CNV)



2-event, homozygous deletion in a Glioblastoma

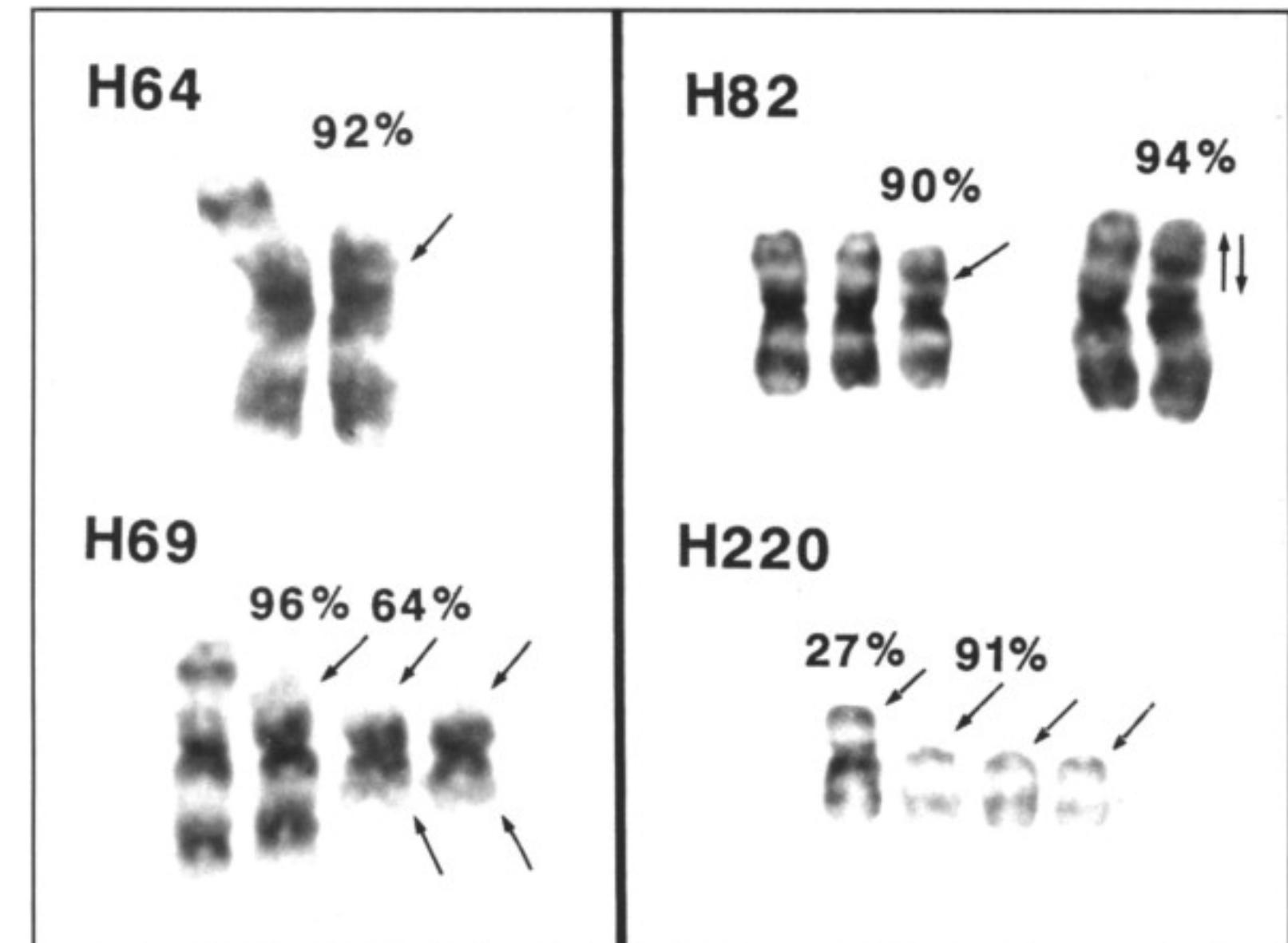
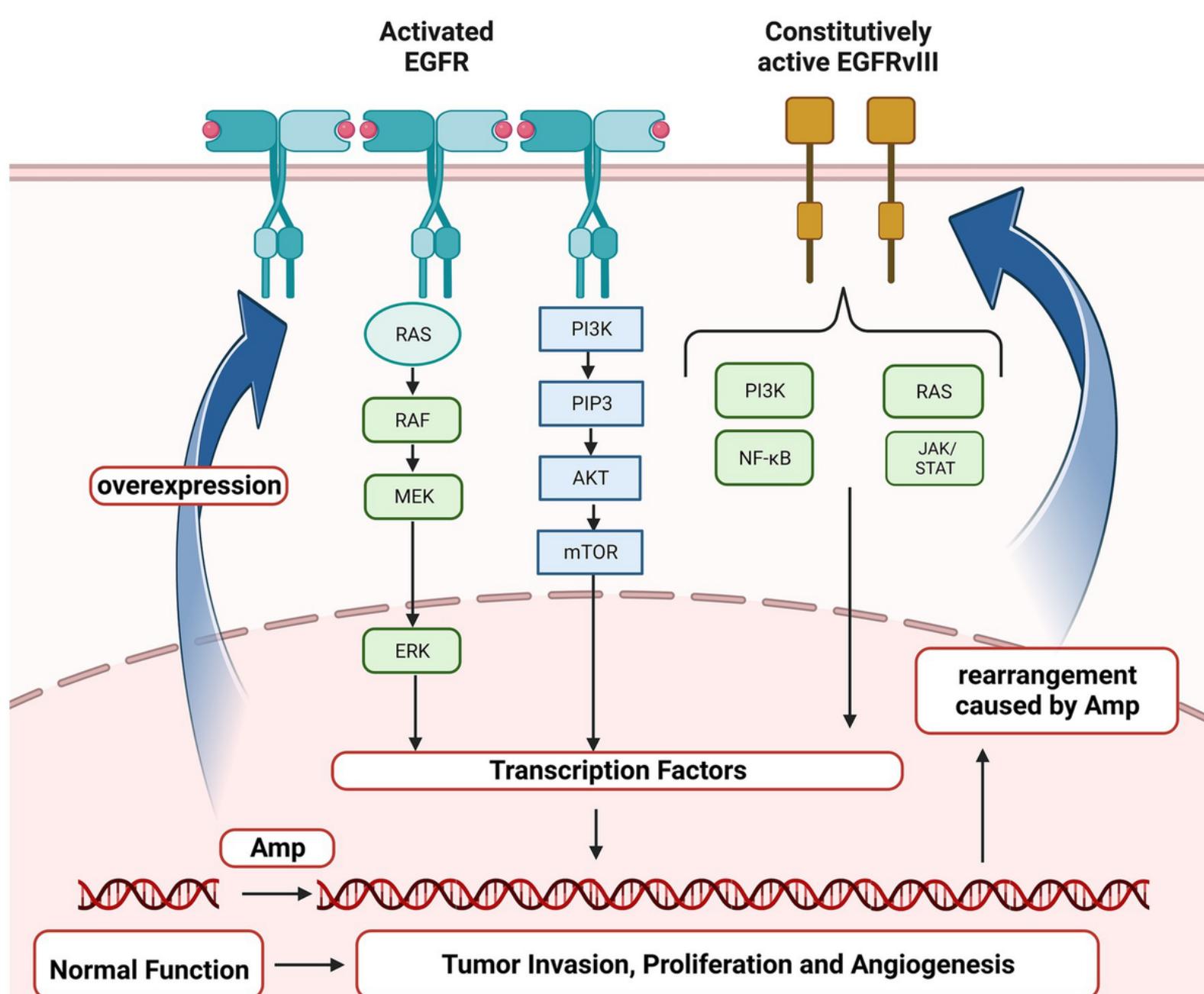
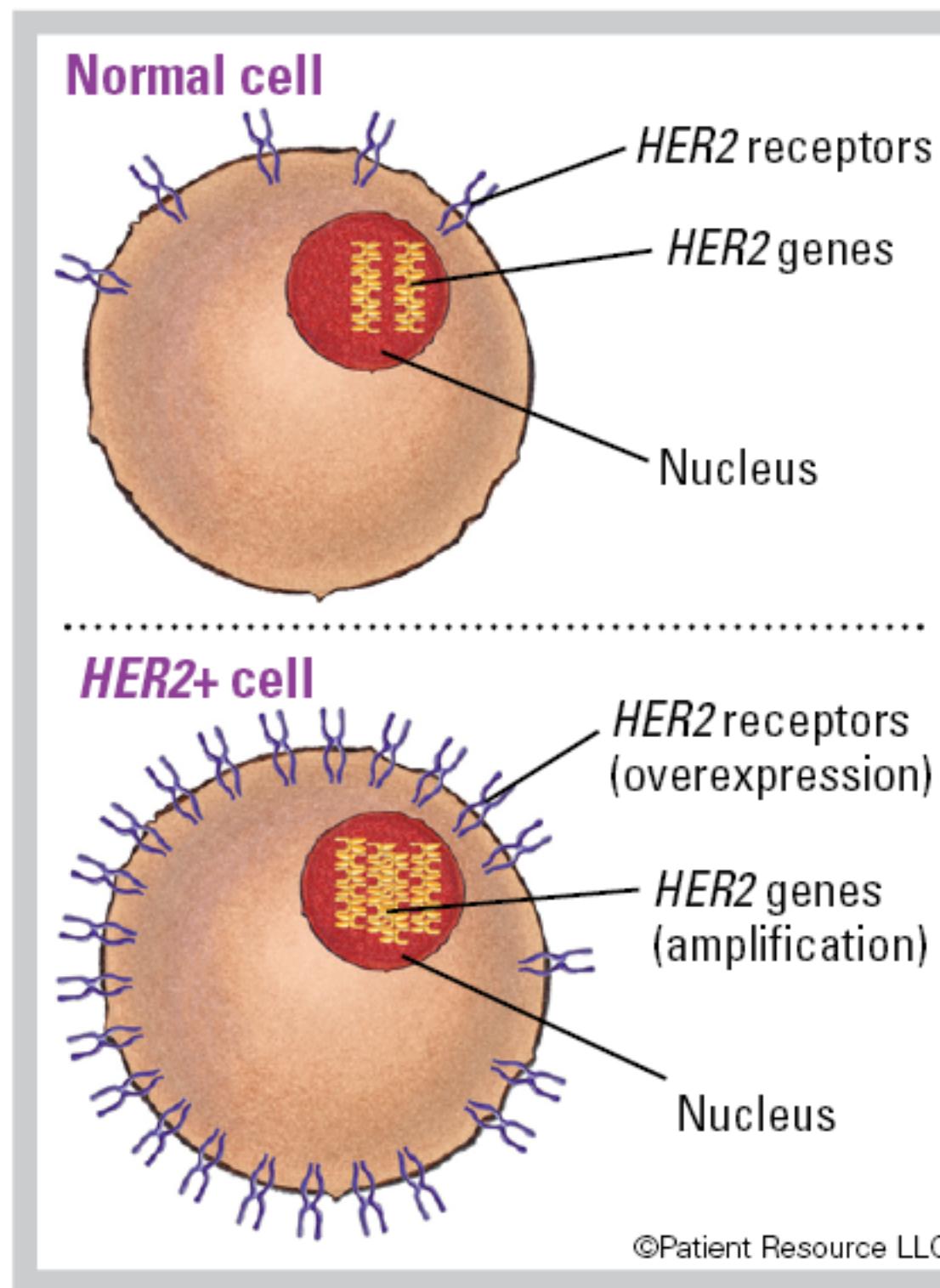
- Intermediate-scale genetic change
- Size: 1kb to multiple megabase
- Additional copies of sequence (**duplications**) and losses of genetic material (**deletions**)



Gain of chromosome arm 13q in colorectal carcinoma

Somatic CNVs related to cancer

▲ FIGURE 1
BREAST CELLS



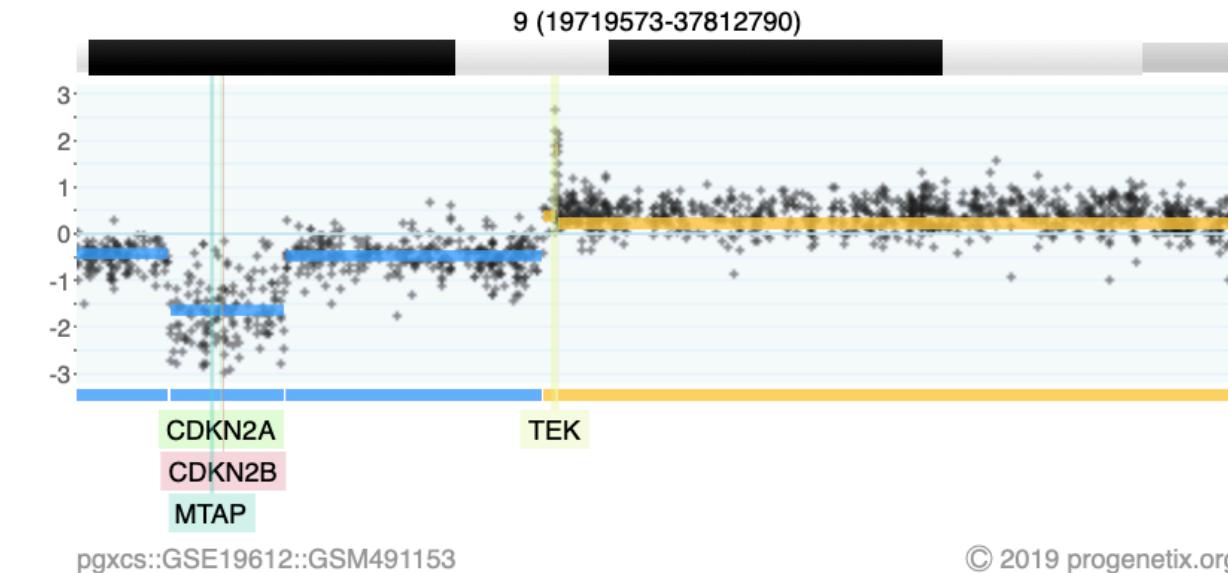
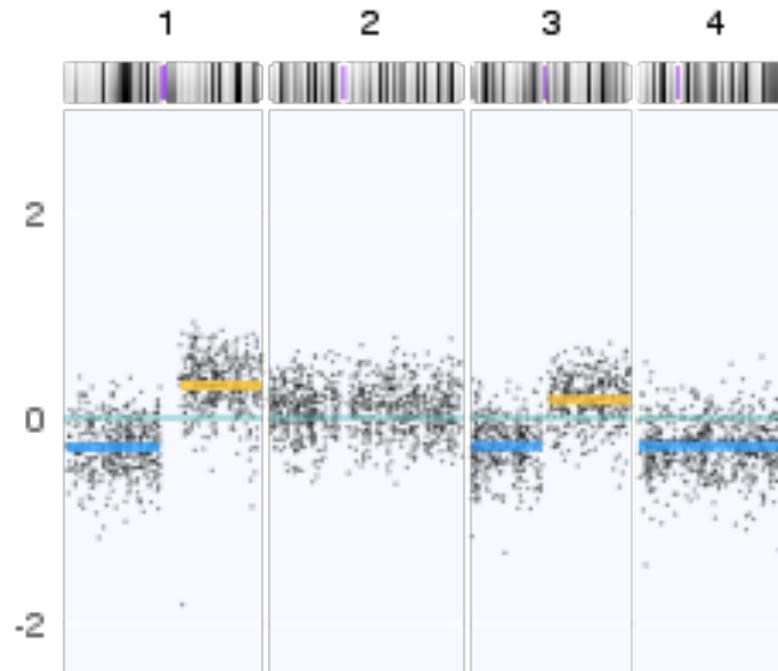
- Somatic CNVs in cancers:
 - HER2 Amplification in Breast Cancer
 - EGFR Amplification in Glioblastoma
 - Chromosome 3p Deletion in Lung Cancer

Theoretical Cytogenetics and Oncogenomics

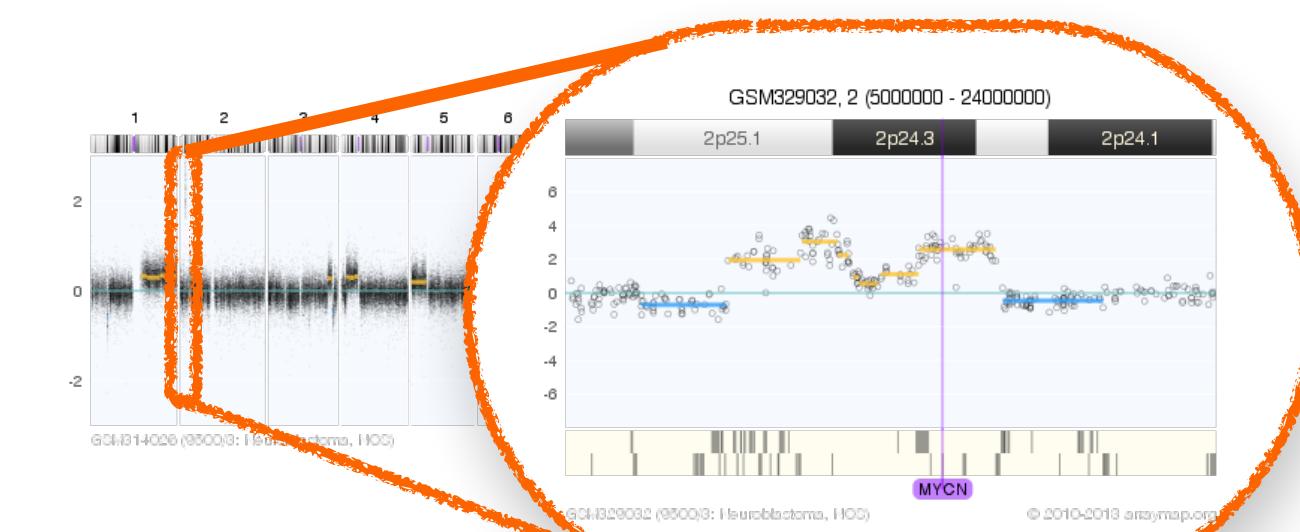
Research | Methods | Standards

Genomic Imbalances in Cancer - Copy Number Variations (CNV)

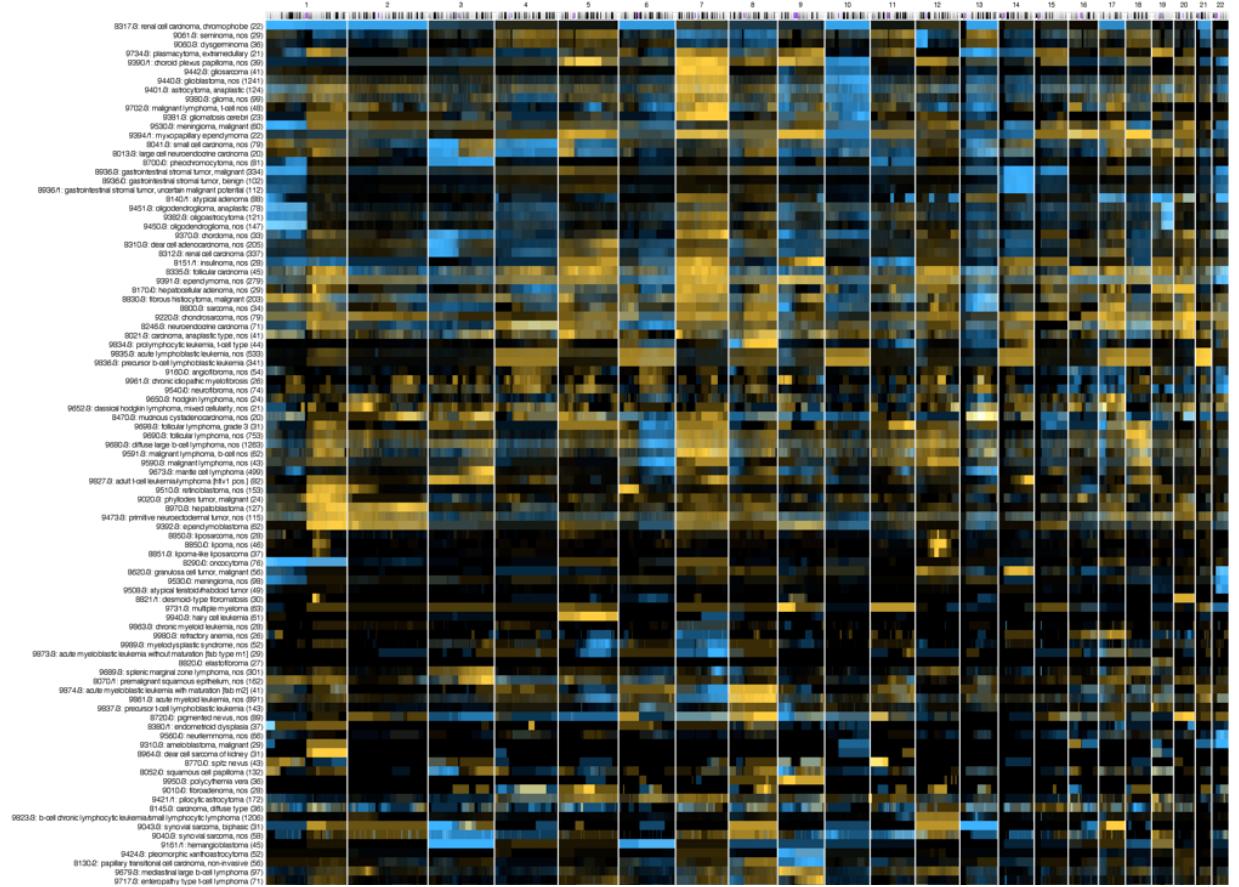
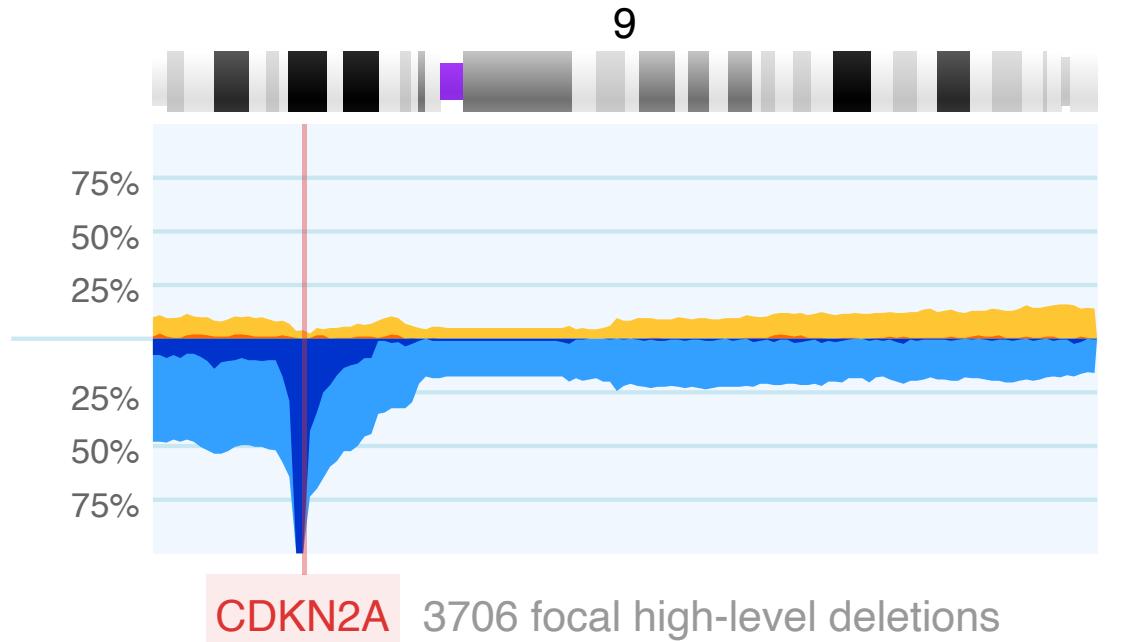
- Point mutations (insertions, deletions, substitutions)
- Chromosomal rearrangements
- **Regional Copy Number Alterations (losses, gains)**
- Epigenetic changes (e.g. DNA methylation abnormalities)



homozygous deletion in a Glioblastoma



MYCN amplification in neuroblastoma
(GSM314026, SJNB8_N cell line)



Cancer Genomics Reference Resource

- **open** resource for oncogenomic profiles
- over **240'000 cancer CNV profiles**
- more than 1'100 diagnostic types
- inclusion of reference sets (e.g. TCGA, GENIE...)
- standardized encodings (e.g. NCIIt, 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



CNV Profiles

- ... by NCIT
- ... by ICD-O Morphology
- ... by ICD-O Site
- ... by TNM & Grade

Search Samples

arrayMap

- TCGA Data
- cBioPortal Studies

Publication DB

Progenetix Use

NCIT - ICD-O Mappings

UBERON Mappings

Upload & Plot

OpenAPI Paths and Examples

Cancer Cell Lines

Beacon+

Documentation

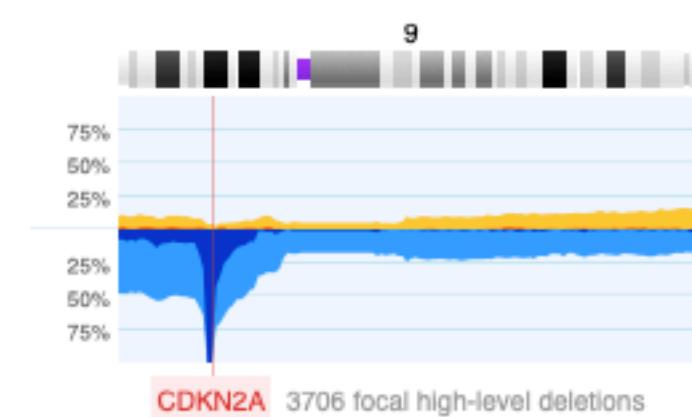
Baudisgroup @ UZH

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* of currently **240600** samples from **1126** different cancer types (NCIt neoplasm classification)

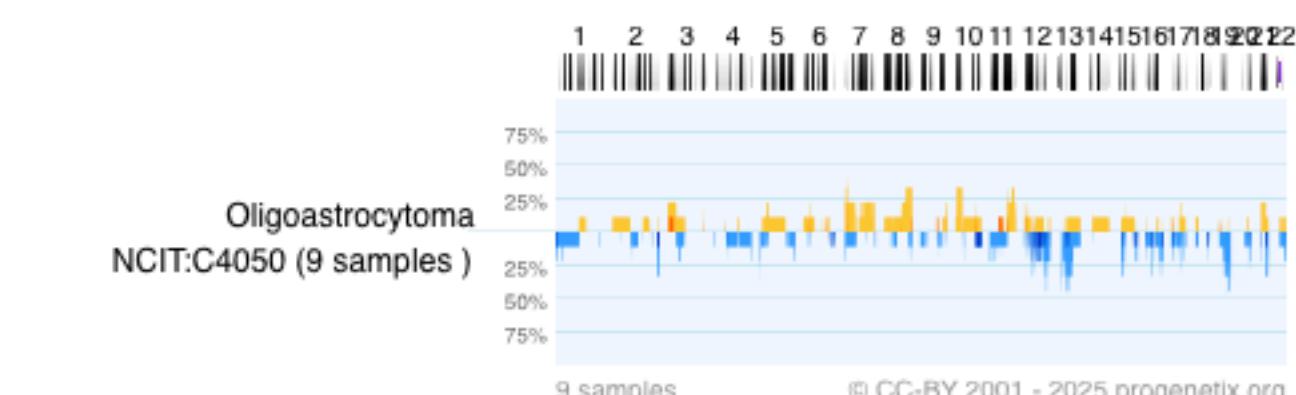
Local CNV Frequencies

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

Frequency profiles of regional genomic gains and losses for all categories (diagnostic entity, publication, cohort ...) can be accessed through the respective Cancer Types pages (e.g. [NCIT Neoplasia Codes](#) ) and compared through the [Compare CNV Profiles](#)  option. Below is an example of aggregated CNV data in 11 samples in Oligoastrocytoma with the frequency of regional **copy number gains (high level)** and **losses (high level)** displayed for the 22 autosomes.



© CC-BY 2001 - 2025 progenetix.org

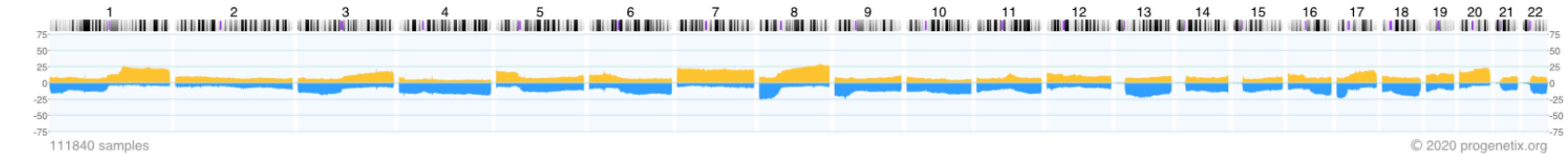
[Download SVG](#) | [Go to NCIT:C4050](#) | [Download CNV Frequencies](#)

Cancer Genomics Publications

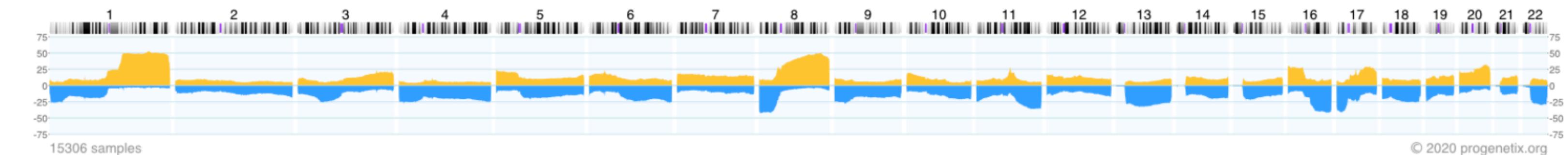
Through the [[Publications](#)] page Progenetix provides annotated references to research articles from cancer genome screening experiments (WGS, WES, aCGH, cCGH). The numbers of analyzed samples and possible availability in the Progenetix sample collection are indicated.

Somatic CNV in Cancer Example Frequency Profiles

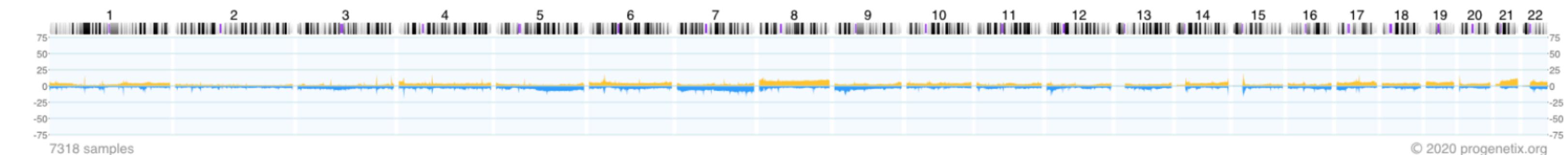
Progenetix: Regional CNV Frequencies in 111'840 Neoplasm (NCIT:C3262)



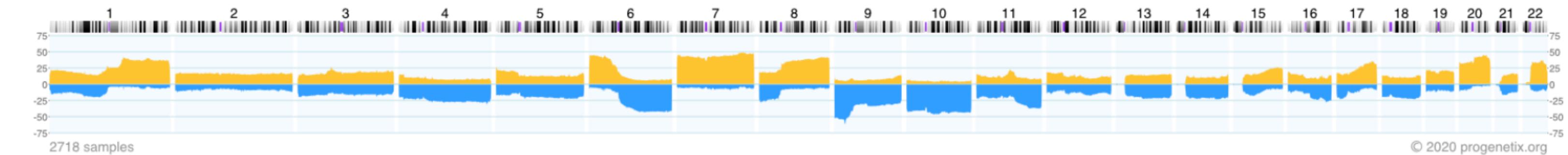
Malignant Breast Neoplasm (NCIT:C9335)



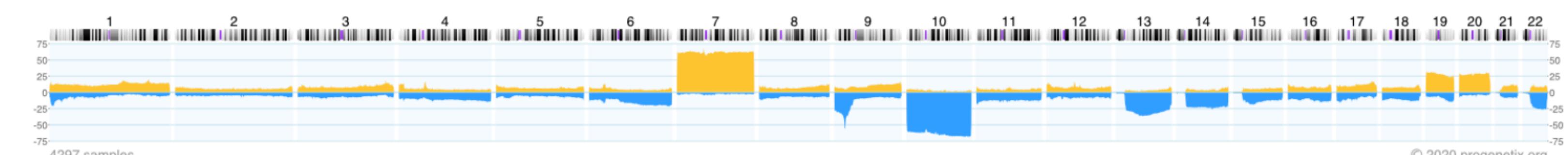
Acute Leukemia (NCIT:C9300)



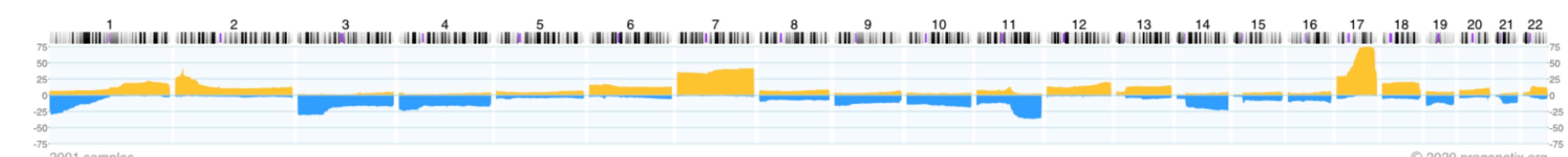
Melanoma (NCIT:C3224)



Glioblastoma (NCIT:C3058)



Neuroblastoma (NCIT:C3270)



2001 samples

© 2020 progenetix.org

Cancer Genomics Reference Resource

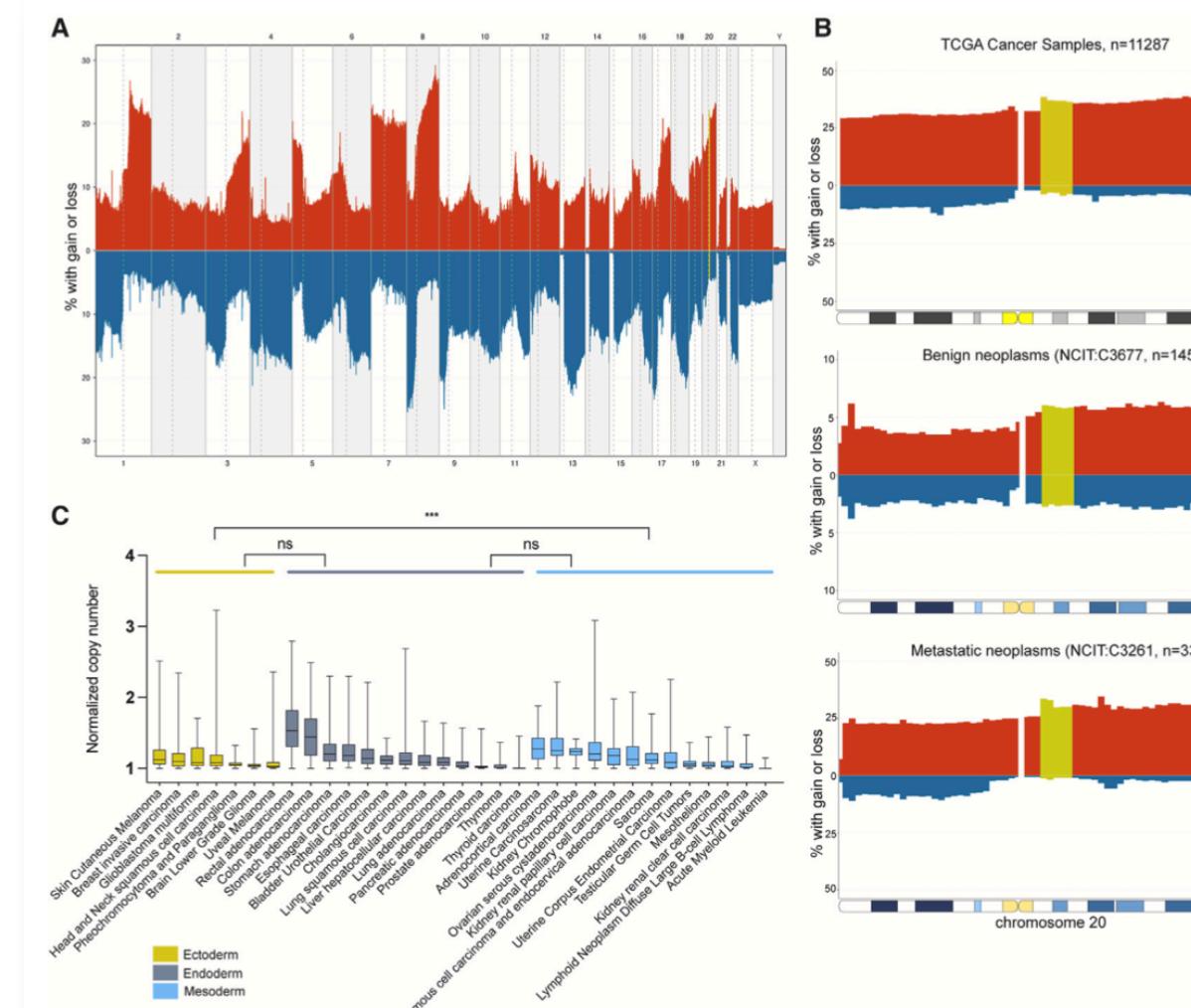
- **open** resource for oncogenomic profiles
- over **240'000 cancer CNV profiles**
- SNV data for some series (e.g. TCGA)
- more than **1100 diagnostic types**
- inclusion of reference datasets (e.g. TCGA, GENIE, cBioPortal)
- standardized encodings (e.g. NCIIt, ICD-O 3)
- identifier mapping for PMID, GEO, Cellosaurus, TCGA, cBioPortal where appropriate
- core clinical data (TNM, sex, survival ...)
- data mapping services

The screenshot shows the Progenetix interface with a sidebar on the left containing links like 'CNV Profiles', 'Search Samples', 'arrayMap', 'Publication DB', 'NCIT - ICD-O Mappings', 'Upload & Plot', 'OpenAPI Paths and Examples', 'Cancer Cell Lines', 'Beacon+', 'Documentation', and 'Baudisgroup @ UZH'. The main area is titled 'Cancer Types by National Cancer Institute NCI Code' and contains a detailed hierarchical tree of cancer types. A search bar at the top says 'Filter subsets e.g. by prefix ...' and a dropdown says 'Hierarchy Depth: 5 levels'. The tree starts with 'Neoplasm' (138552 samples, 144862 CNV profiles), which branches into 'Neoplasm by Site' (133029 samples, 139114 CNV profiles), then into 'Genitourinary System Neoplasm' (21582 samples, 23171 CNV profiles), 'Benign Genitourinary System Neoplasm' (243 samples, 243 CNV profiles), 'Benign Urinary System Neoplasm' (98 samples, 98 CNV profiles), 'Benign Urinary Tract Neoplasm' (3 samples, 3 CNV profiles), 'Benign Kidney Neoplasm' (95 samples, 95 CNV profiles), 'Benign Reproductive System Neoplasm' (145 samples, 145 CNV profiles), 'Benign Female Reproductive System Neoplasm' (145 samples, 145 CNV profiles), 'Malignant Genitourinary System Neoplasms' (20567 samples, 22154 CNV profiles), 'Metastatic Malignant Genitourinary System Neoplasms' (2 samples, 2 CNV profiles), 'Metastatic Genitourinary System Carcinoma' (2 samples, 2 CNV profiles), 'Genitourinary System Carcinoma' (19462 samples, 20921 CNV profiles), 'Metastatic Genitourinary System Carcinoma' (2 samples, 2 CNV profiles), 'Female Reproductive System Carcinoma' (5746 samples, 5974 CNV profiles), 'Male Reproductive System Carcinoma' (7022 samples, 7808 CNV profiles), 'Urinary System Carcinoma' (6694 samples, 7139 CNV profiles), and 'Recurrent Malignant Genitourinary System Neoplasms' (3 samples, 3 CNV profiles).

Progenetix Use

- CNV data is used e.g. as reference data in cancer genomics studies
- diagnosis specific CNV profiles serve as "fast look-up" in clinical genomics laboratories
- we loosely track publications in our literature database but there is no systematic check-back mechanism...

Example: 2025 article using Progenetix' *pgxRpi* Beacon/R interface to retrieve & visualize 117'587 cancer CNV profiles for a study into pluripotent stem cells' genomics



Progenetix References

Articles Citing - or Using - Progenetix

This page lists articles which we found to have made use of, or referred to, the Progenetix resource ecosystem. These articles may not necessarily contain original case profiles themselves. Please [contact us](#) to alert us about additional articles you are aware of. Also, you can now directly submit suggestions for matching publications to the [oncopubs repository on Github](#).

Filter

Publications (121)	Samples		
id	Publication	Genomes	pgx
PMID:38157850	Krivec N, Ghosh MS et al. (2024) Gains of 20q11.21 in human pluripotent stem cells: Insights from cancer research. ... Stem Cell Reports	0	0
PMID:37627037	Austin BK, Firooz A, Valafar H et al. (2023) An Updated Overview of Existing Cancer Databases and Identified Needs. Biology (Basel)	0	0
PMID:37393410	Liu SC, Wang CI, Liu TT, Tsang NM et al. (2023) A 3-gene signature comprising CDH4, STAT4 and EBV-encoded LMP1 for early diagnosis ... Discov Oncol	0	0

Stem Cell Reports Review



OPEN ACCESS

Gains of 20q11.21 in human pluripotent stem cells: Insights from cancer research

Nuša Krivec,^{1,2} Manjusha S. Ghosh,^{1,2} and Claudia Spits^{1,2,*}

¹Research Group Reproduction and Genetics, Faculty of Medicine and Pharmacy, Vrije Universiteit Brussel, Brussels, Laarbeeklaan 103, 1090 Brussels, Belgium

²These authors contributed equally.

*Correspondence: claudia.spits@vub.be
<https://doi.org/10.1016/j.stemcr.2023.11.013>

Figure 2. Copy-number alterations of human chromosome 20q11.21 in cancers

(A) Aggregated copy-number variation (CNV) data of 117,587 neoplasms (NCIT: C3262) from the Progenetix database ([Huang et al., 2021](#)) were plotted using R library *pgxRpi*. The percentage of samples with aberrations (red, gain; blue, loss) for the whole chromosome are indicated on the y axis. Chromosomal regions are depicted on the x axis; the minimal region of interest at chr20:31216079-35871578 is marked in moss green. NCIT, National Cancer Institute Thesaurus.

(B) Top to bottom: Aggregated CNV data of 11,287 TCGA cancer samples, 336 metastatic neoplasms (NCIT: C3261), and 1,455 benign neoplasms (NCIT: C3677) from the Progenetix database ([Huang et al., 2021](#)), respectively, were plotted using R library *pgxRpi*. The percentage of samples with aberrations (red, gain; blue, loss) for the whole chromosome are indicated on the y axis. Chromosomal regions are depicted on the x axis; the minimal region of interest at chr20:31216079-35871578 is marked in moss green.

{Bio|informatics}Science}

```
for t in pars.keys():

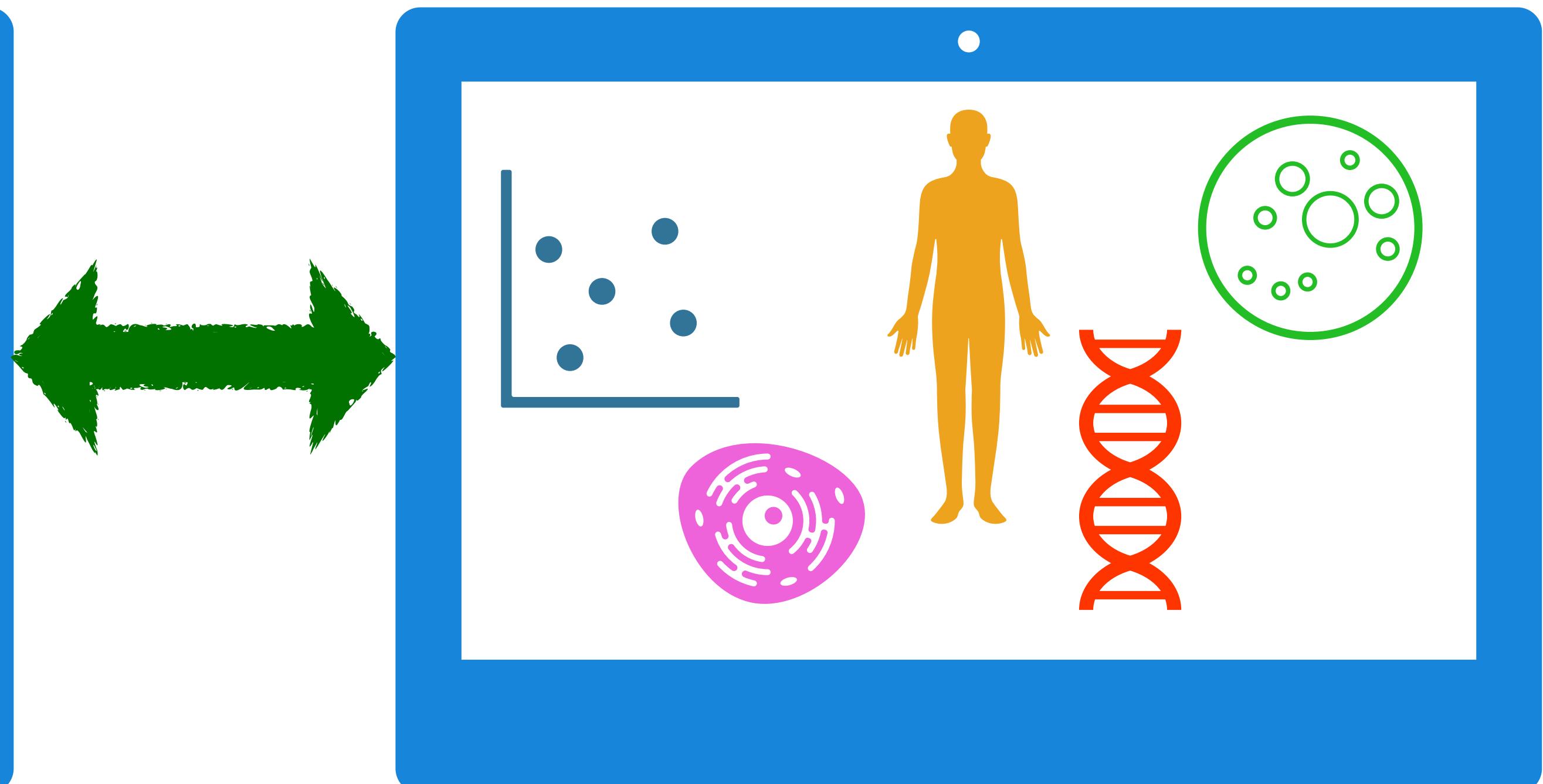
    covs = np.zeros((cs_no, int_no))
    vals = np.zeros((cs_no, int_no))

    if type(callsets).__name__ == "Cursor":
        callsets.rewind()

    for i, cs in enumerate(callsets):
        covs[i] = cs["cnv_statusmaps"][pars[t]["cov_l"]]
        vals[i] = cs["cnv_statusmaps"][pars[t]["val_l"]]

    counts = np.count_nonzero(covs >= min_f, axis=0)
    frequencies = np.around(counts * f_factor, 3)
    medians = np.around(np.ma.median(np.ma.masked_where(covs < min_f, vals), axis=0).filled(0), 3)
    means = np.around(np.ma.mean(np.ma.masked_where(covs < min_f, vals), axis=0).filled(0), 3)

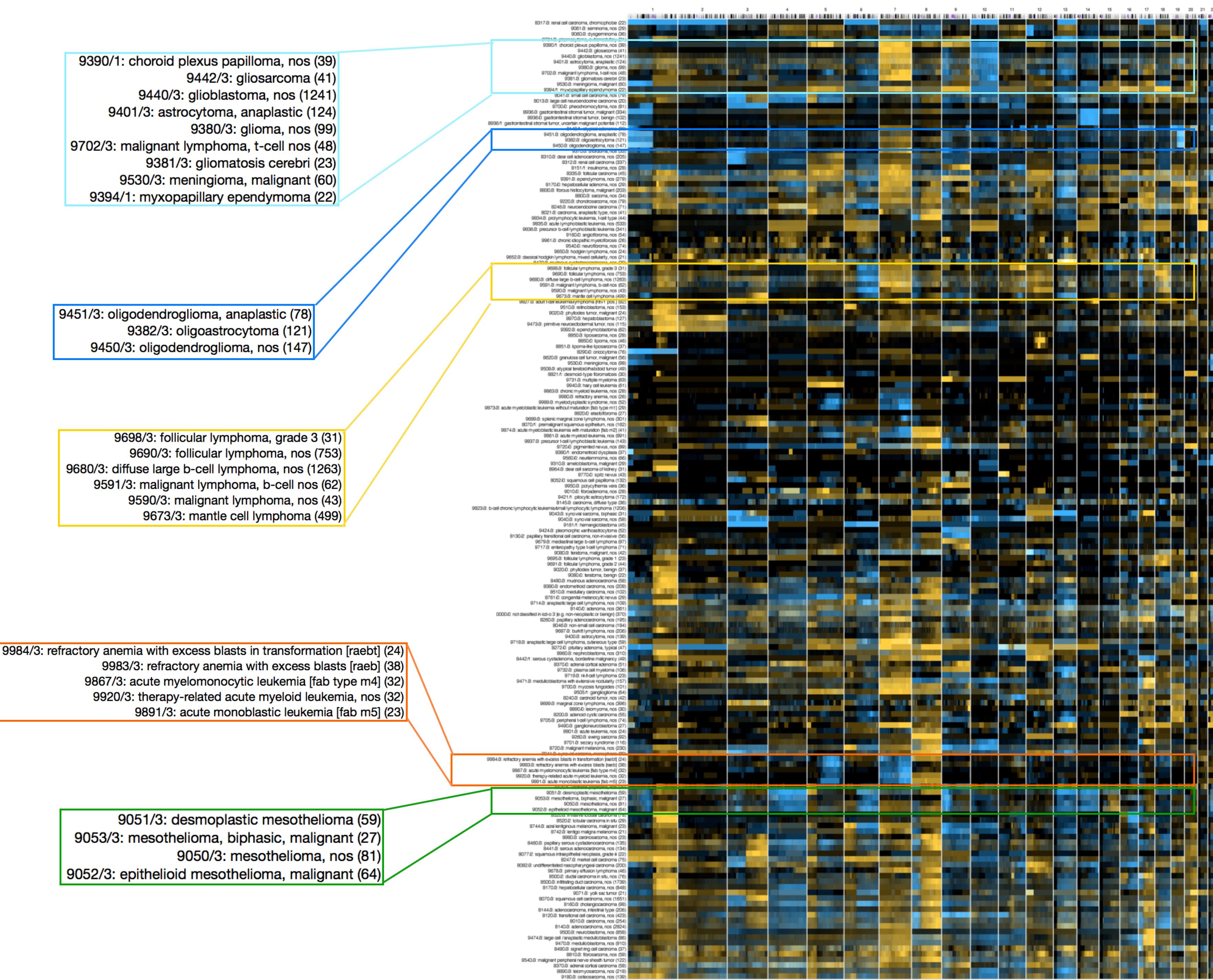
    for i, interval in enumerate(int_fs):
        int_fs[i].update({
            t + "_frequency": frequencies[i],
            t + "_median": medians[i],
            t + "_mean": means[i]
        })
```



Somatic Mutations In Cancer: Patterns

Making the case for genomic classifications

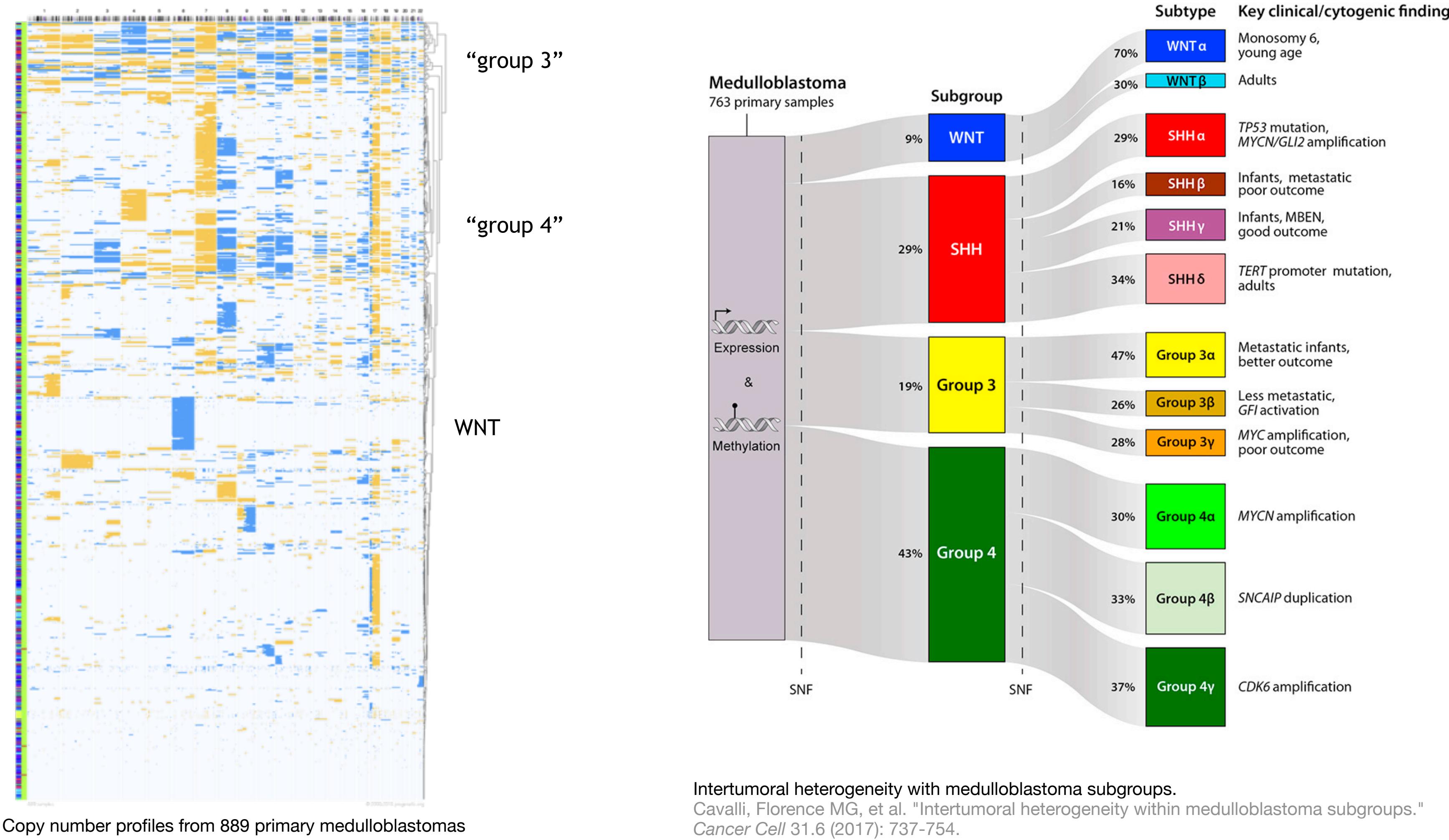
Some related cancer entities show similar copy number profiles



CNA & Cancer heterogeneity

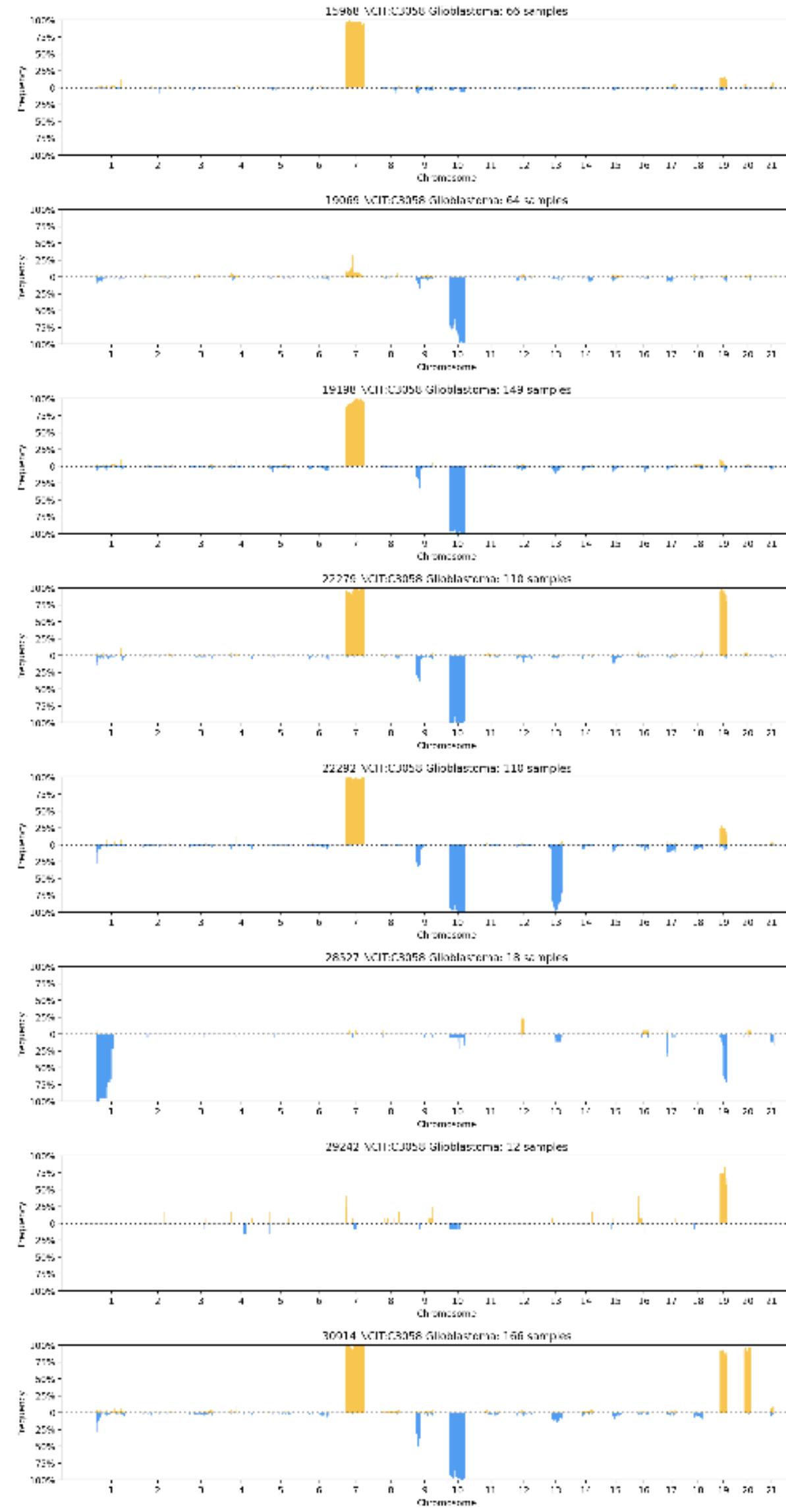
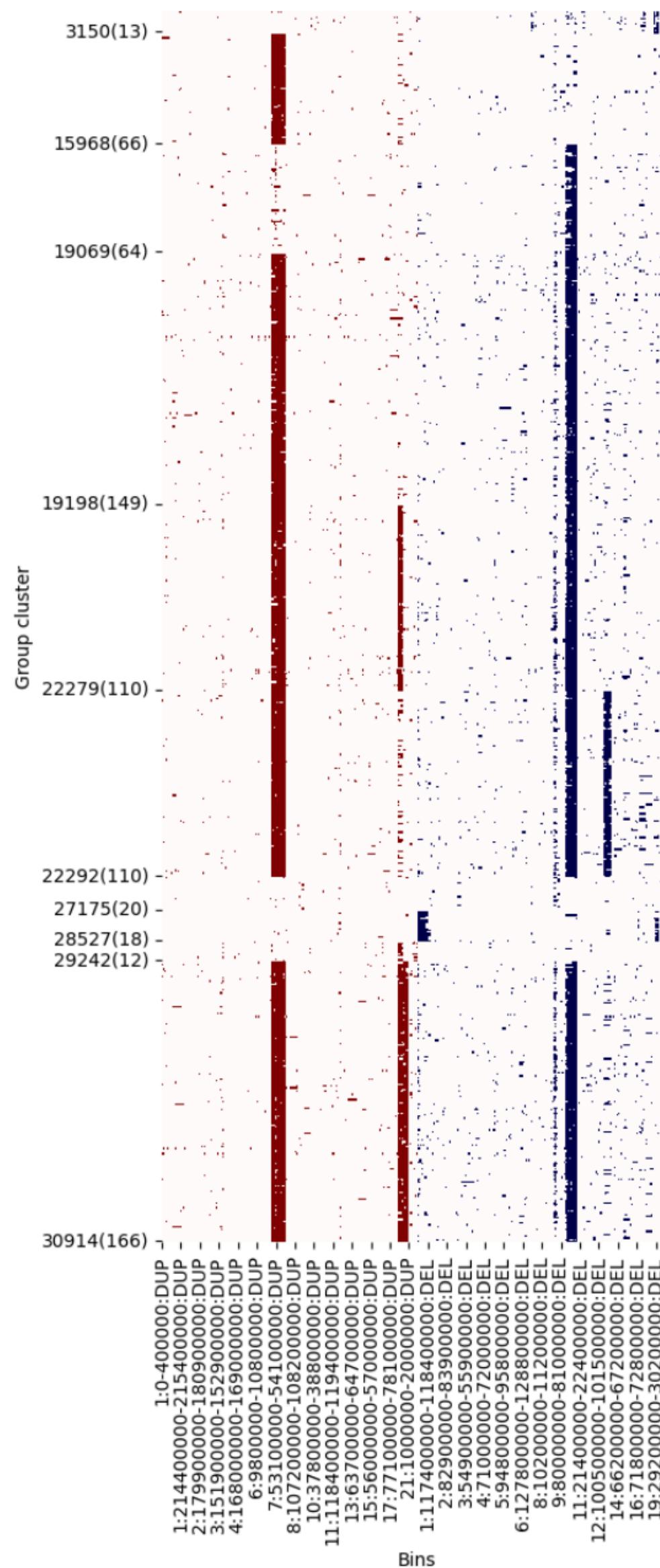
Cancer type definitions can be improved by the addition of molecular parameters as subtype markers or even complete re-evaluation of entity definitions from molecular subtypes with distinct functional mechanisms and clinical trajectories.

Lead: Ziyang Yang



Results

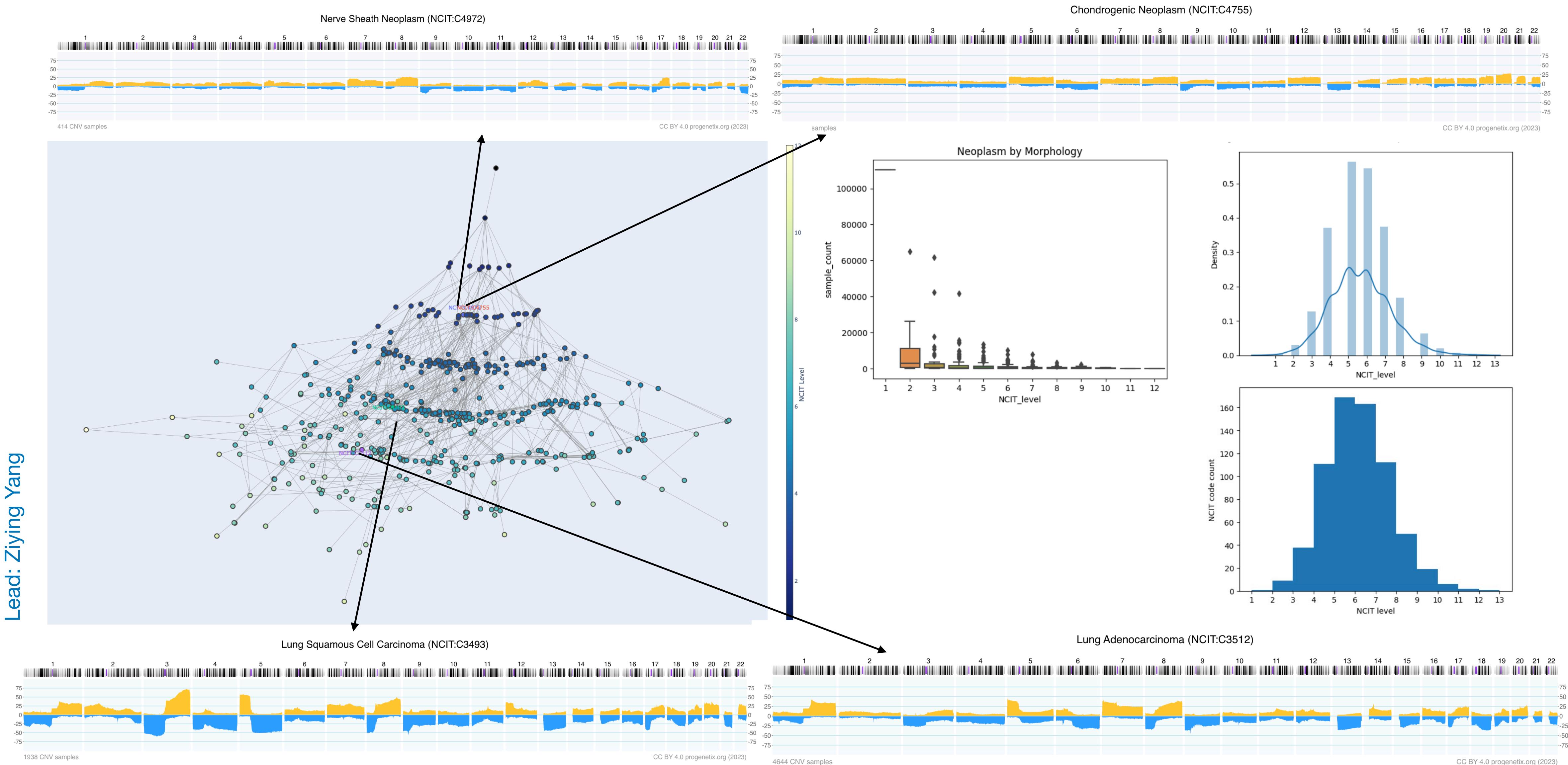
Entity CNV heterogeneity: Glioblastoma



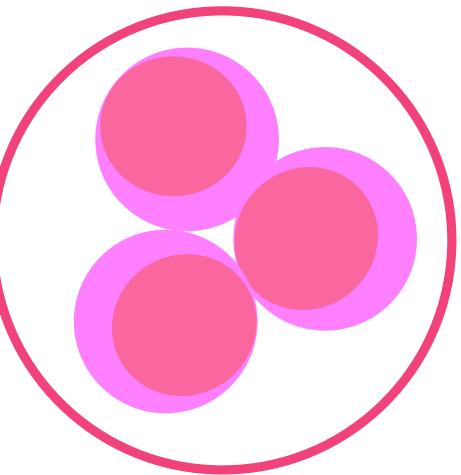
group cluster	CNV features
15968	Dup 7
19069	Del 10
19198	Dup 7, Del 10
22279	Dup 7, Del 10, Dup 19
22292	Dup 7, Del 10, Del 13
27175	Del 1p, Del 19q
28527	Del 1p, Del 19q
29242	Dup 19
30914	Dup 7, Del 10, Dup 19, Dup 20

CNV profiles heterogeneity vs cancer classification

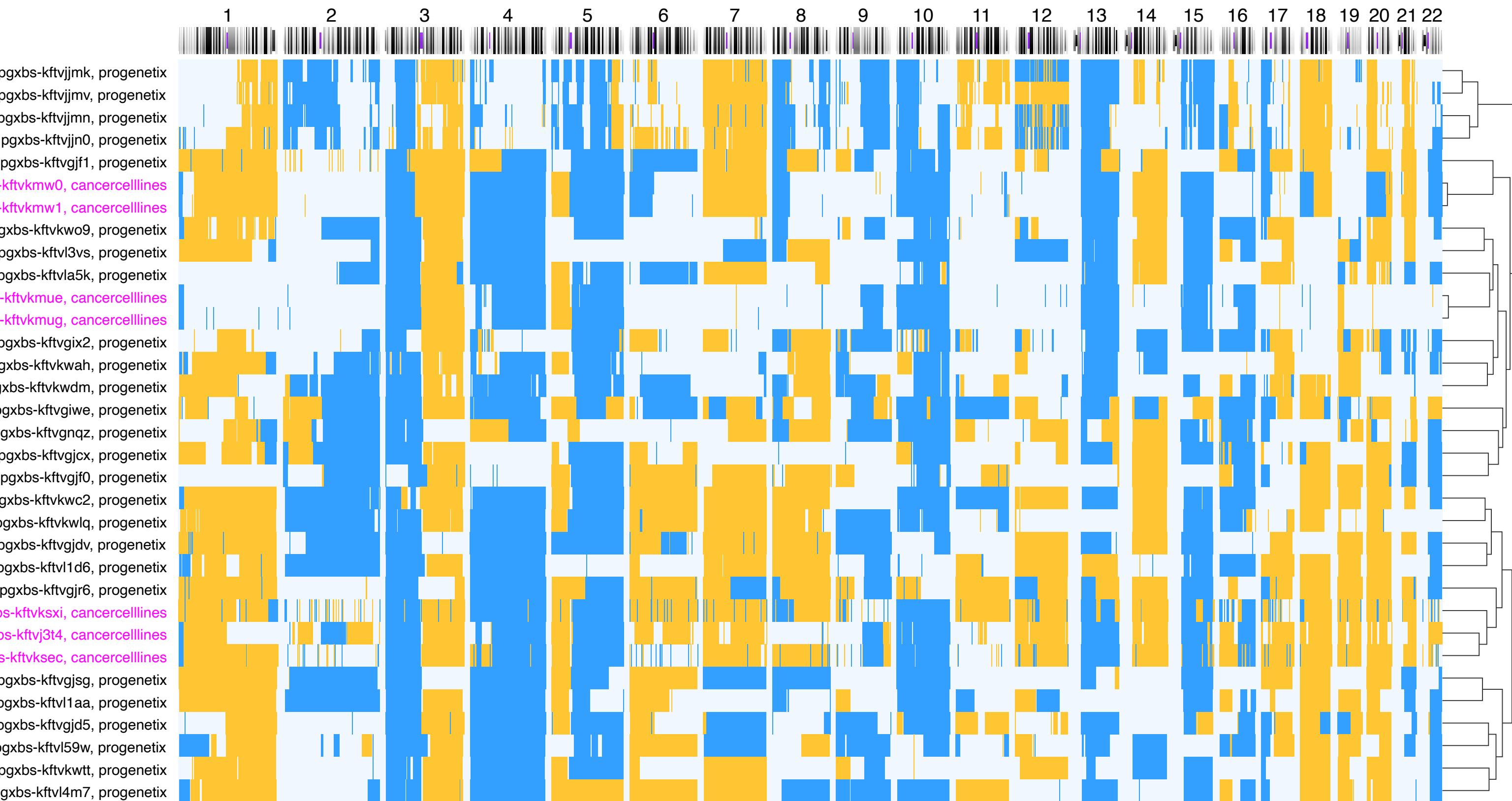
Correspondance of genomic profiles to NCIT cancer hierarchy



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



Population stratification in cancer samples based on SNP array data

- Despite extensive somatic mutations of cancer profiling data, consistency between germline and cancer samples reached 97% and 92% for 5 and 26 populations
- Comparison of our benchmarked results with self-reported meta-data estimated a matching rate between 88 % to 92%.
- Ethnicity labels indicated in meta-data are vague compared to the standardized output from our tool



Lead: Qingyao Huang

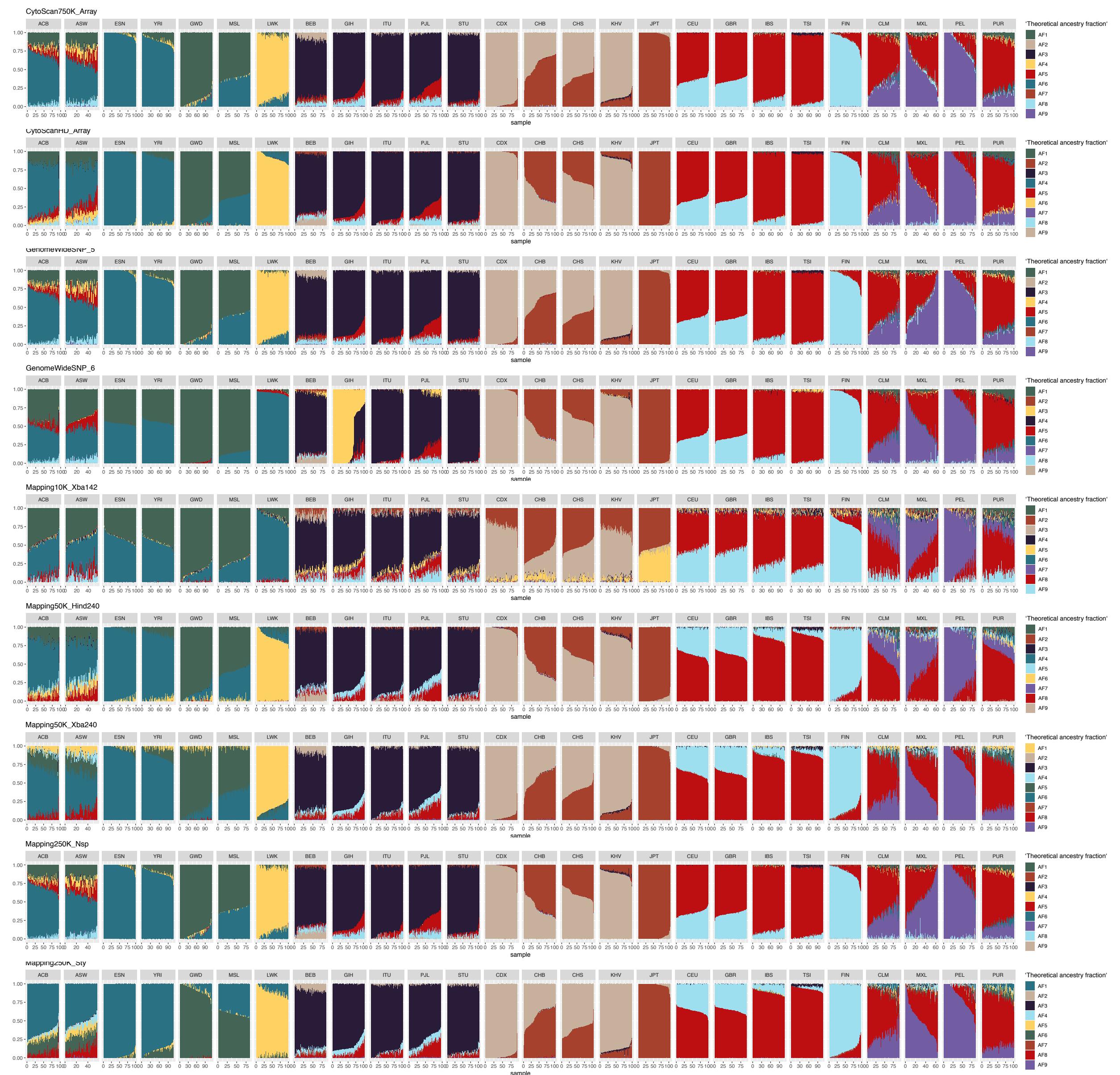


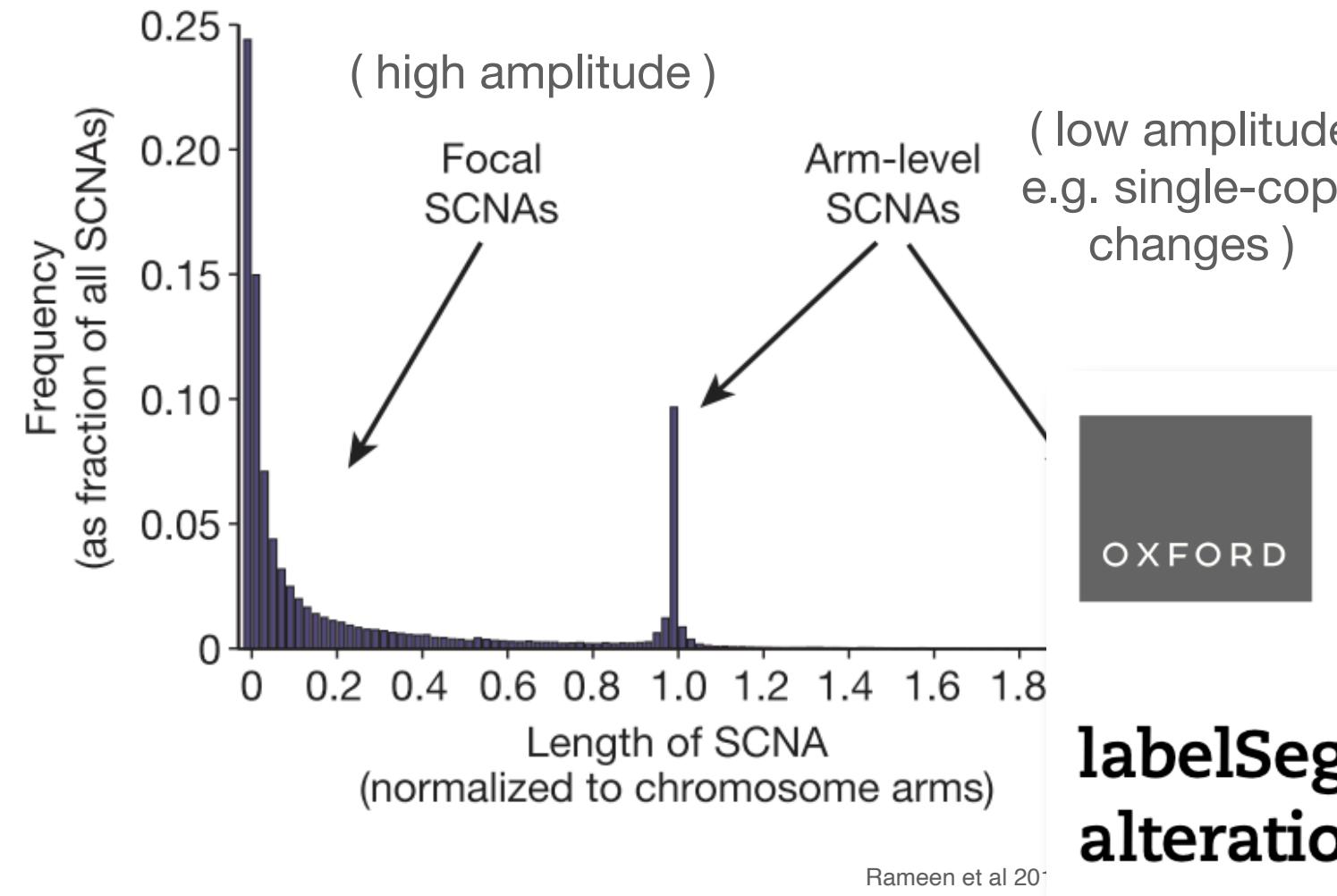
Figure S1 The fraction or contribution of theoretical ancestors ($k=9$) in reference individuals from 1000 Genomes Project with regard to nine SNP array platforms. The x-axis are individual samples, grouped by their respective population. Groups belonging to the same continent/superpopulation are placed neighboring to each other: AFR (1-7), SAS (8-12), EAS (13-17), EUR (18-22), AMR (23-26).

CNV Categorization Method Supporting GA4GH Standards

GA4GH Variation Representation Specification



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

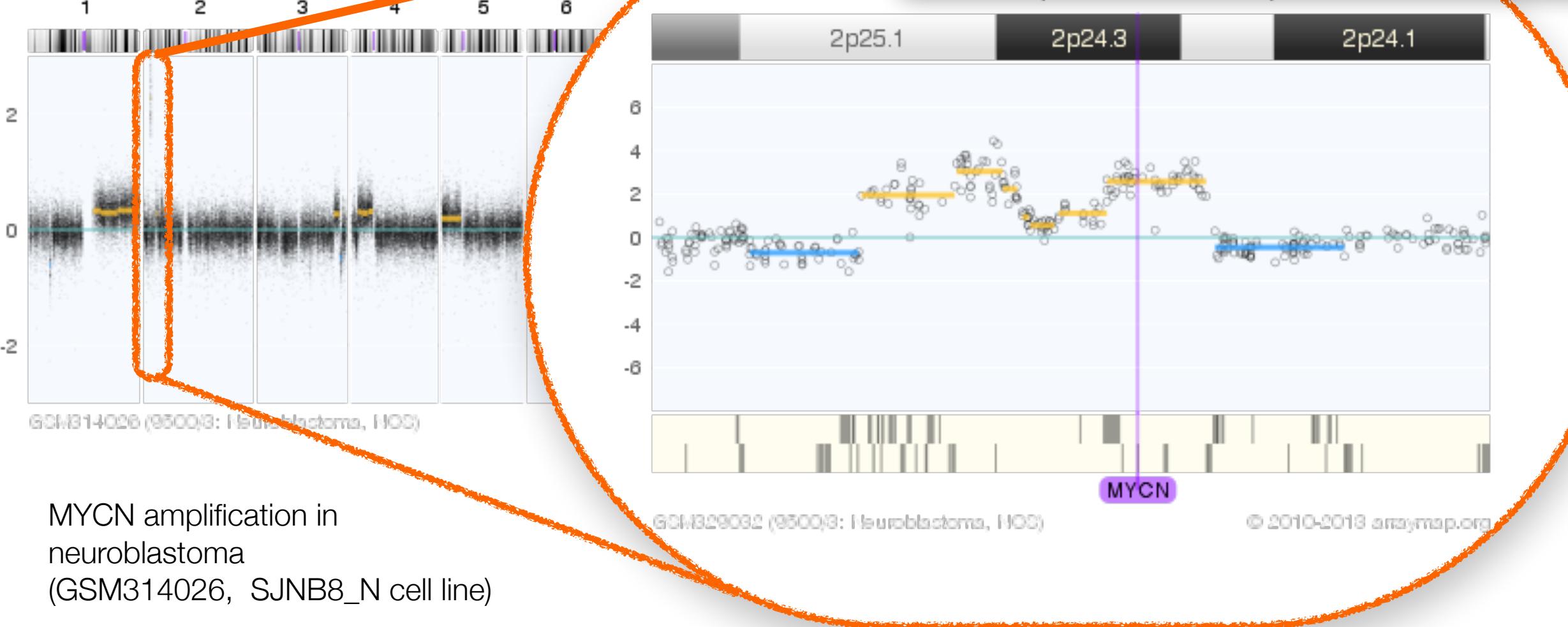


labelSeg: segment annotation for tumor copy number alteration profiles

Hangjia Zhao and Michael Baudis

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

GSM329032, 2 (5000000 - 24000000)



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 express copy number variation are interpreted to be relative copy

Briefings in Bioinformatics, 2024, 25(2), 1–12

<https://doi.org/10.1093/bib/bbad541>

Problem Solving Protocol

number of a [Location](#) or a [Feature](#) within a system (e.g. genome, cell,

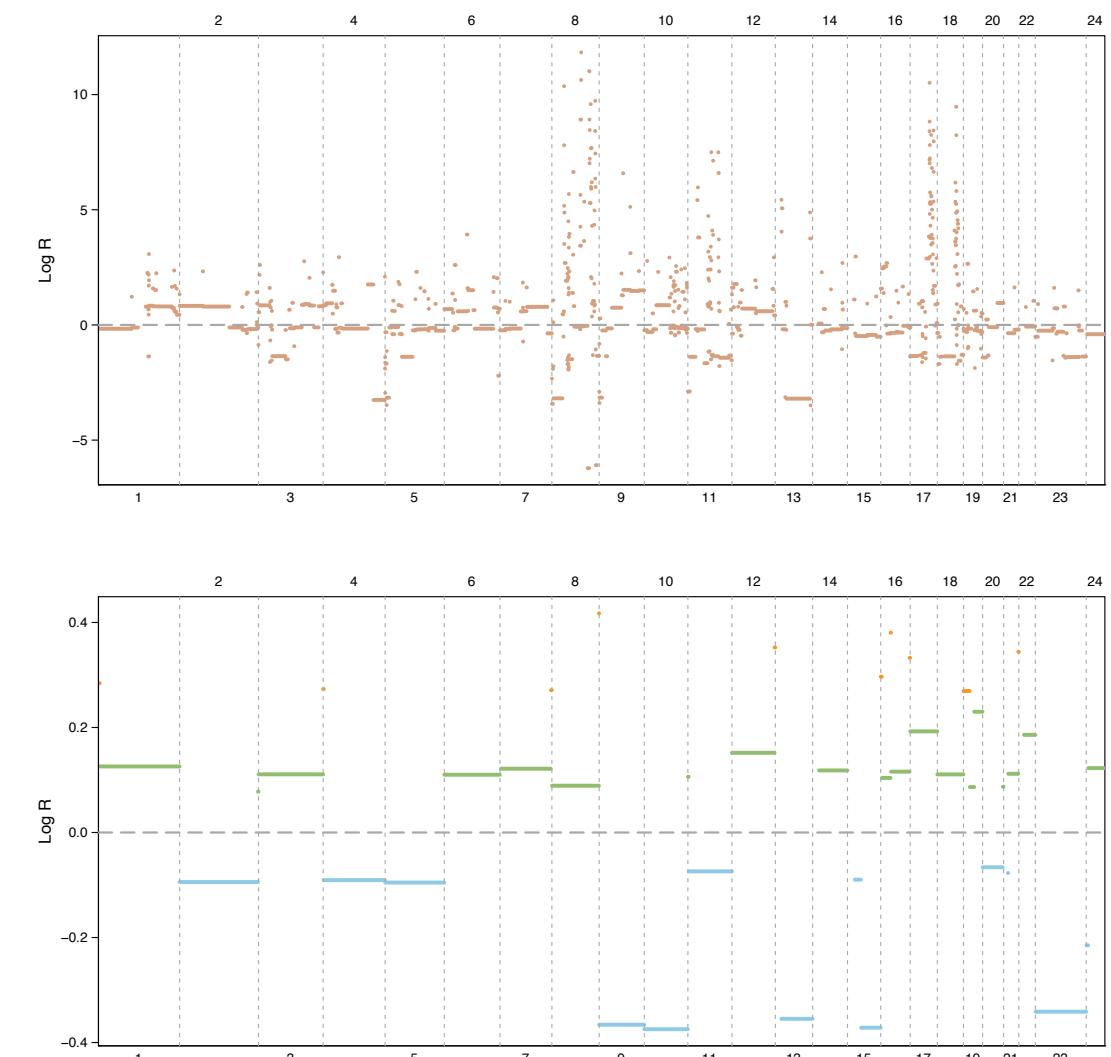
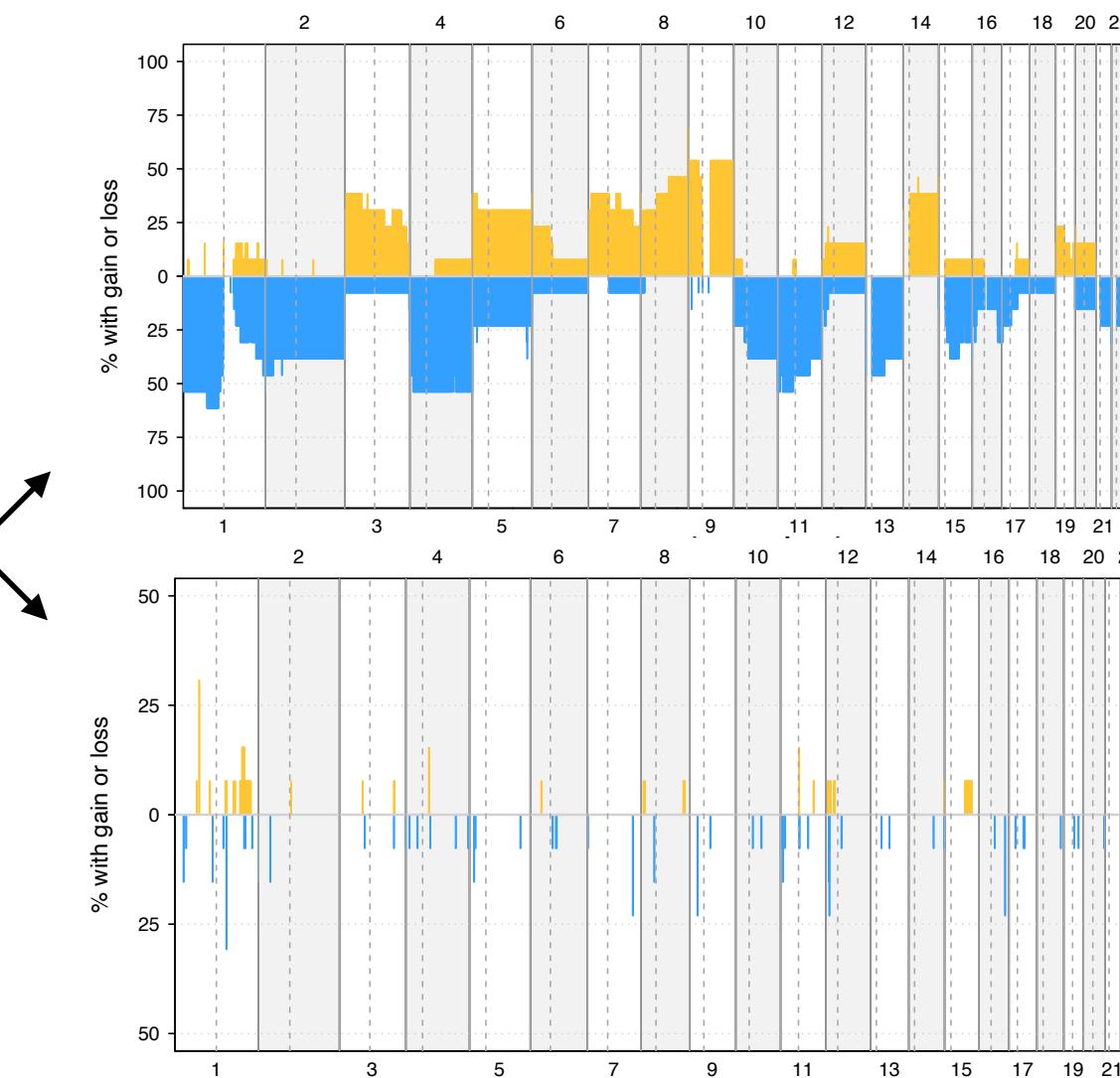
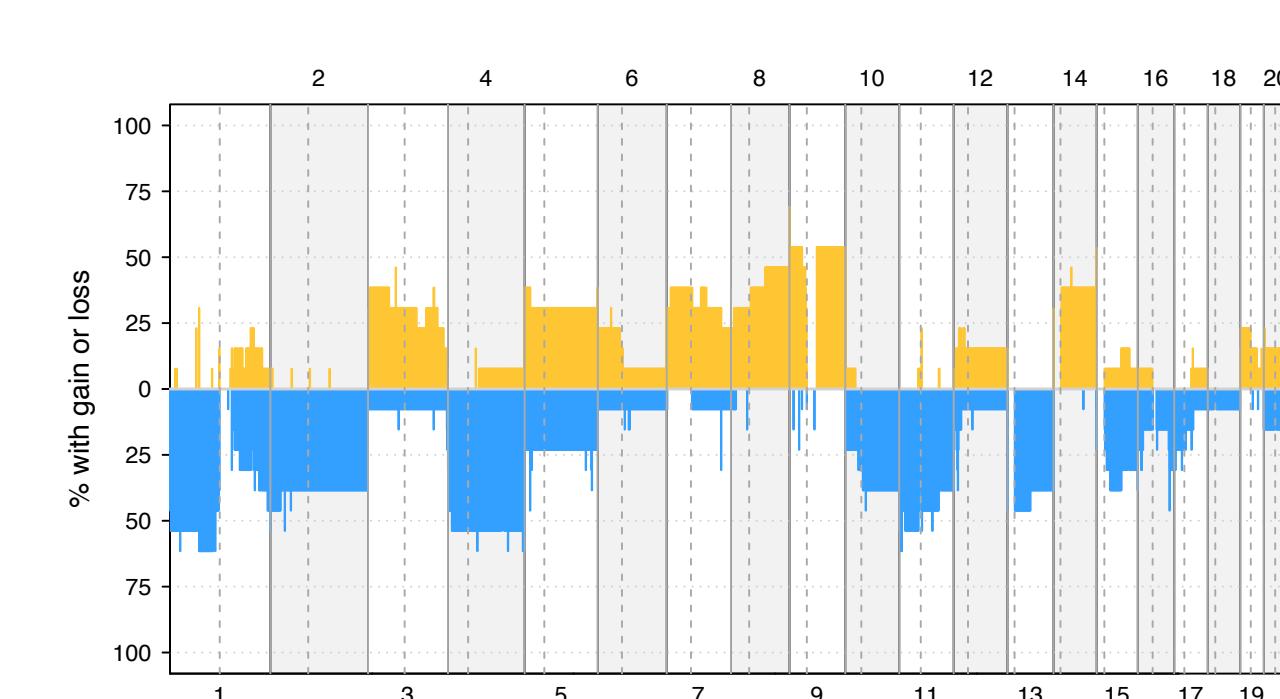
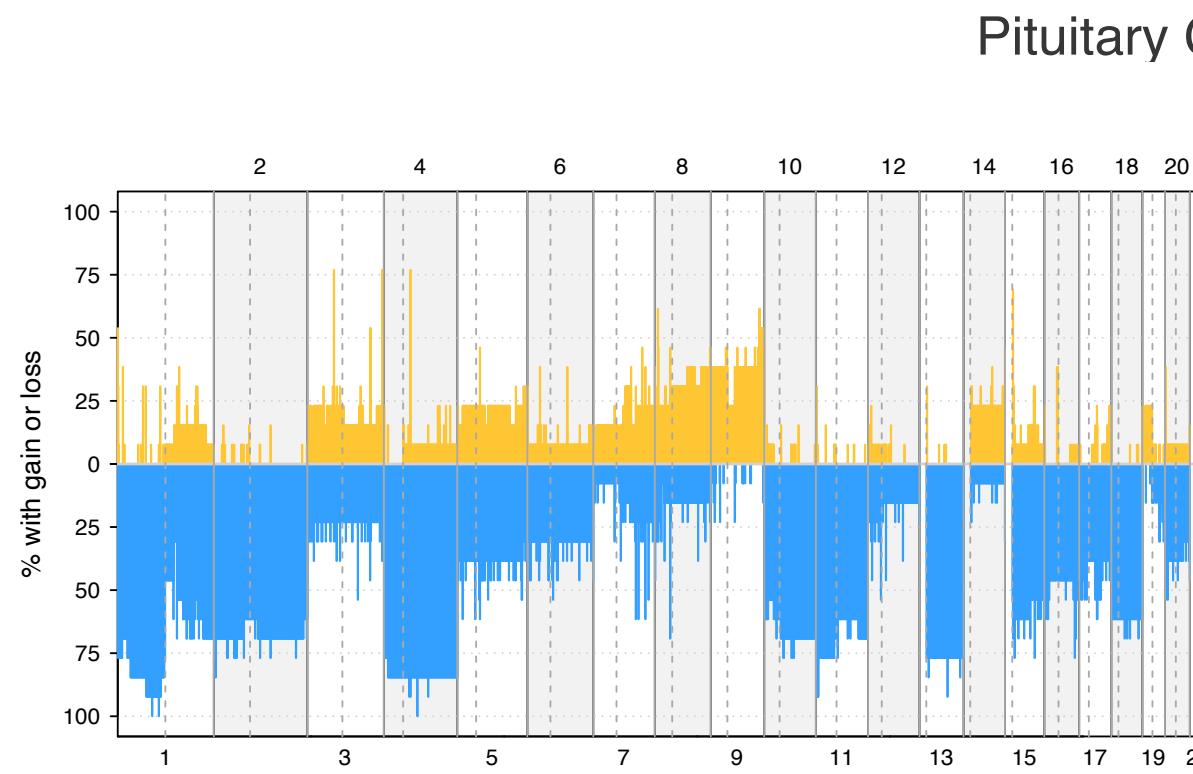
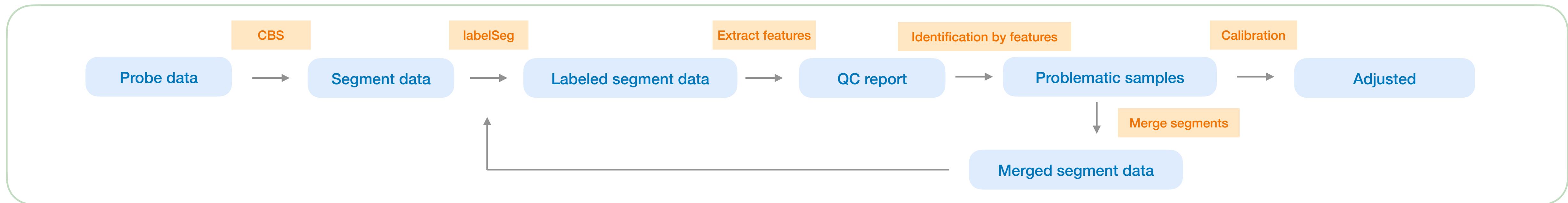
organism, tissue, cell type, individual, population, etc.). Copy number variations are inherited from [Variation](#).

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).

Pipeline Development

improve CNV calling in large numbers of heterogeneous cancer samples

nextflow

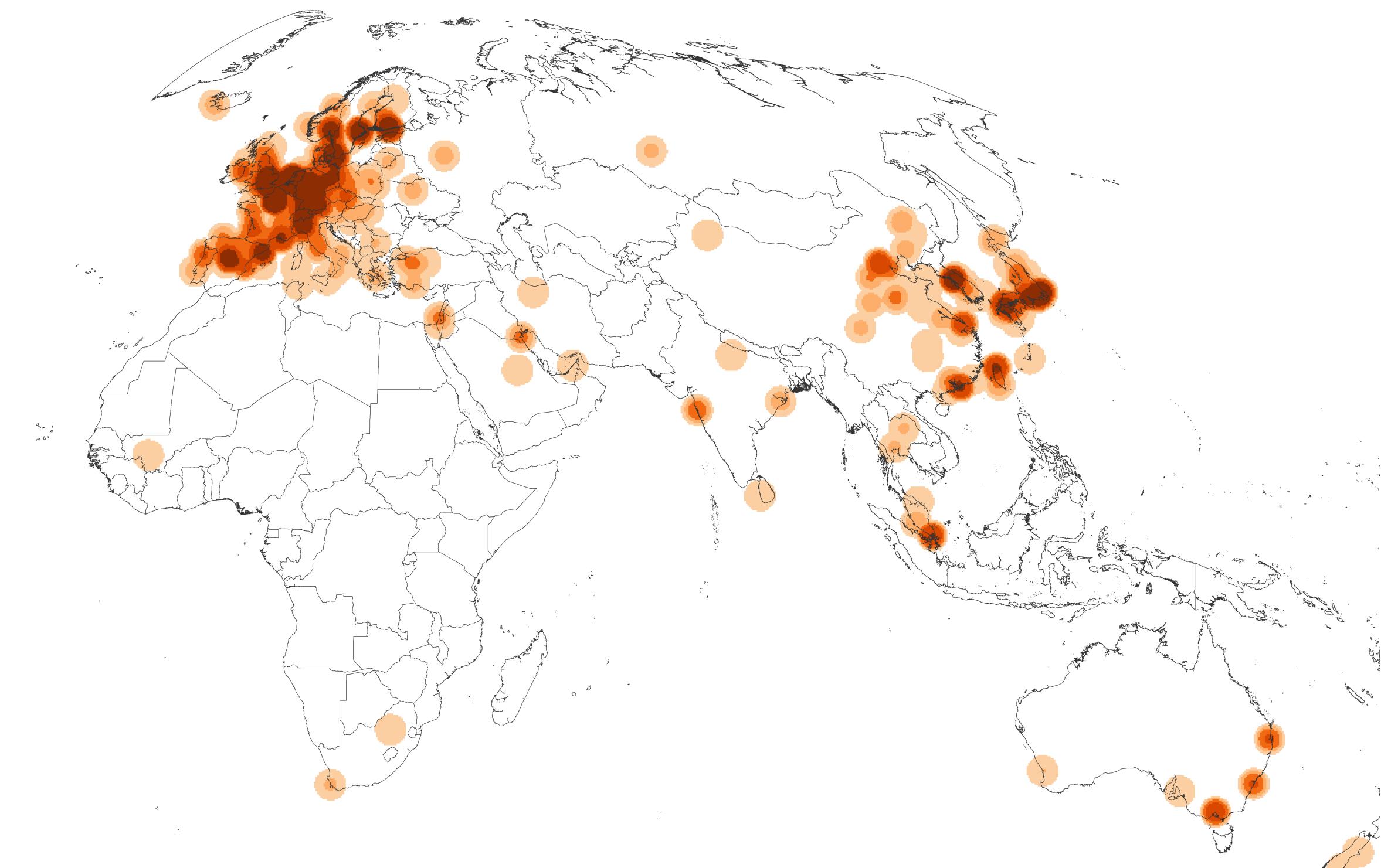
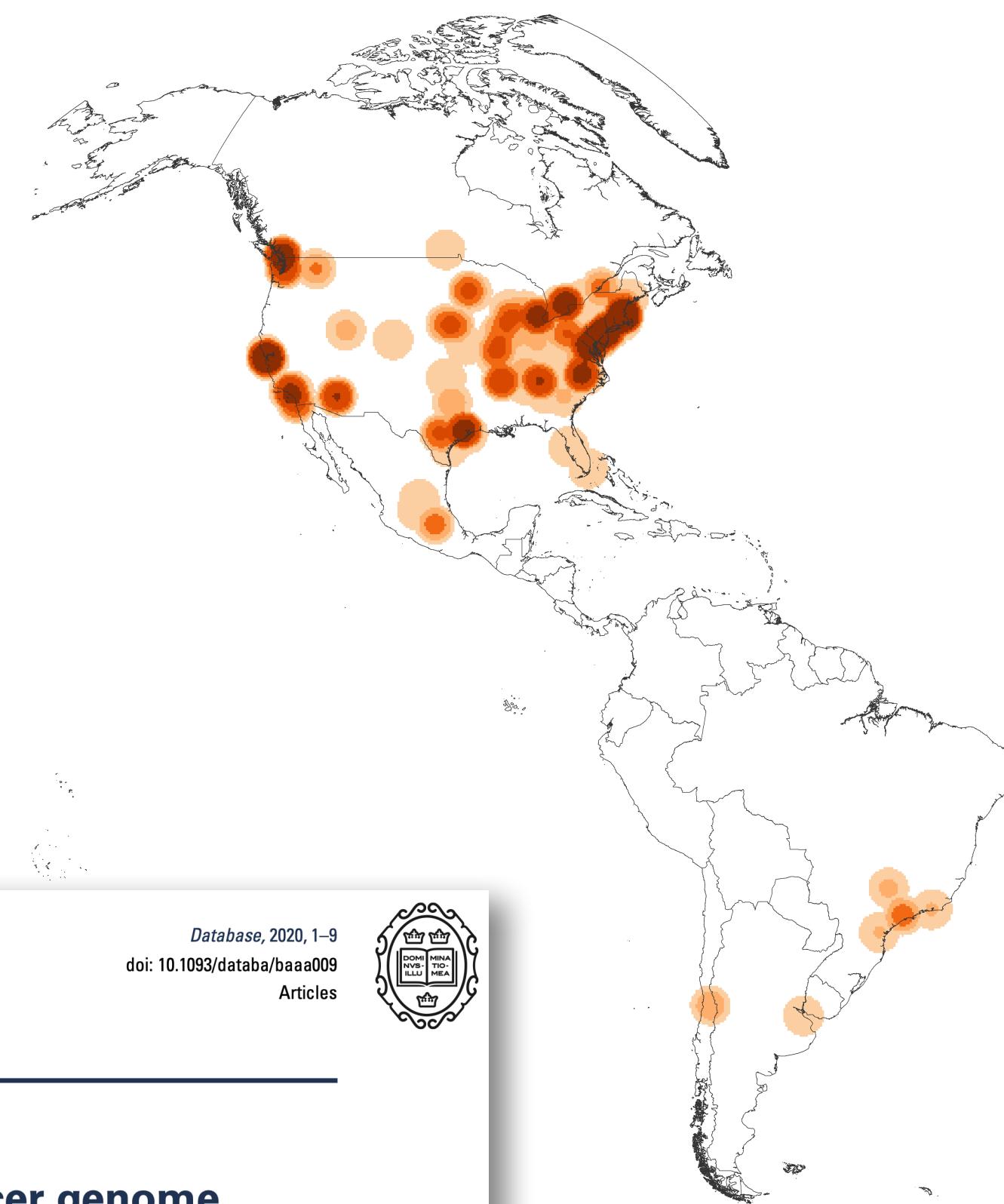
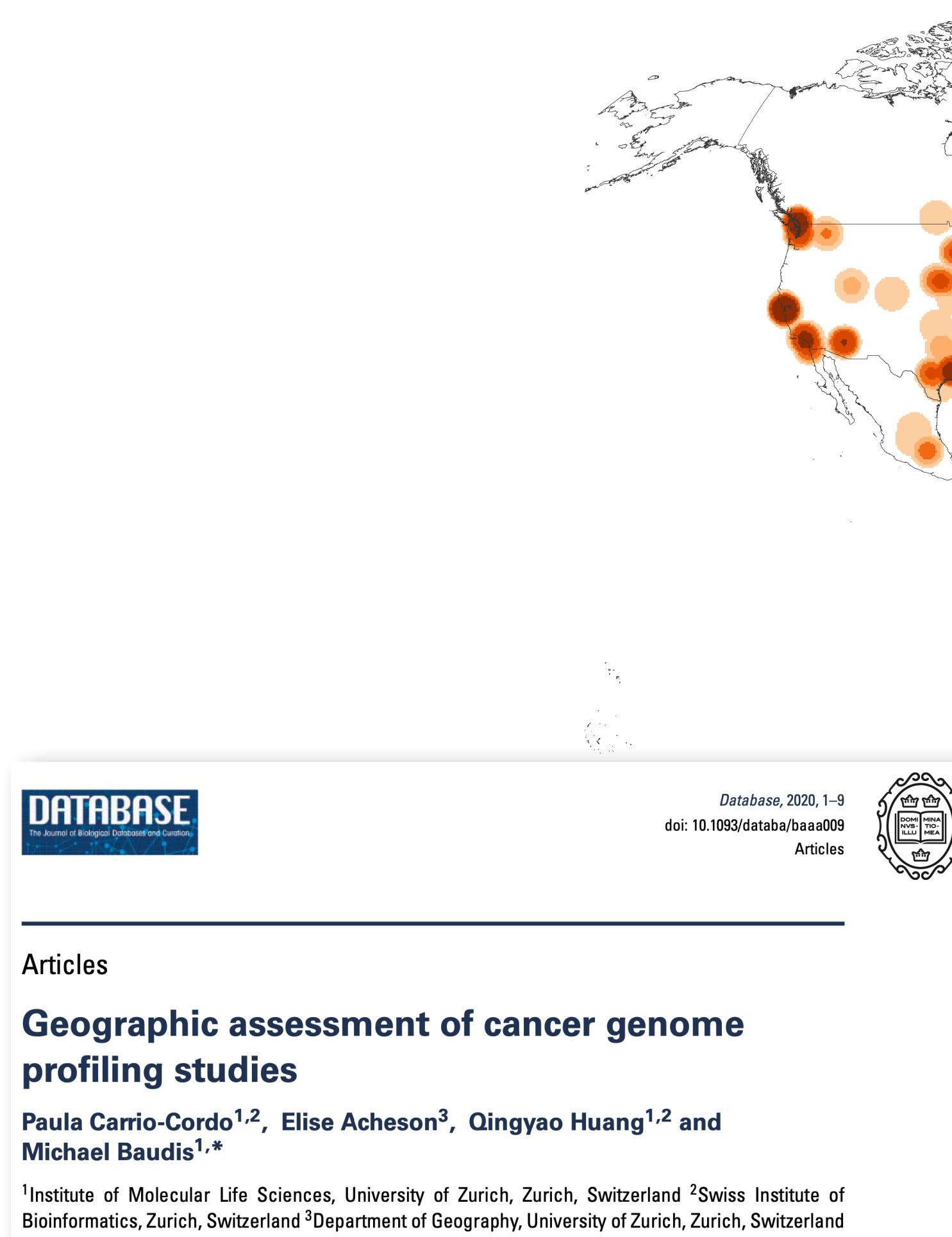


TCGA BLCA project (pgx:TCGA.BLCA)



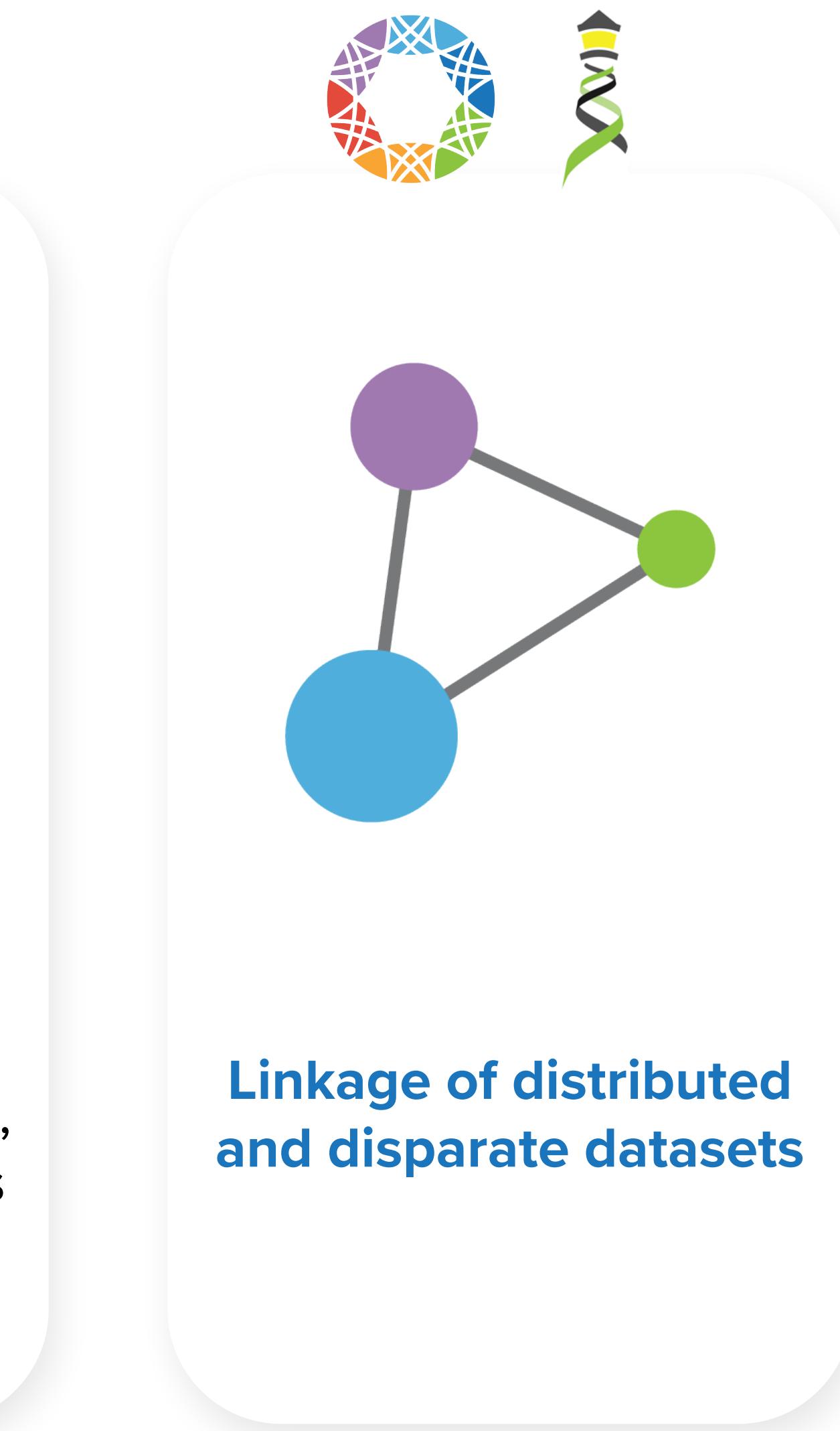
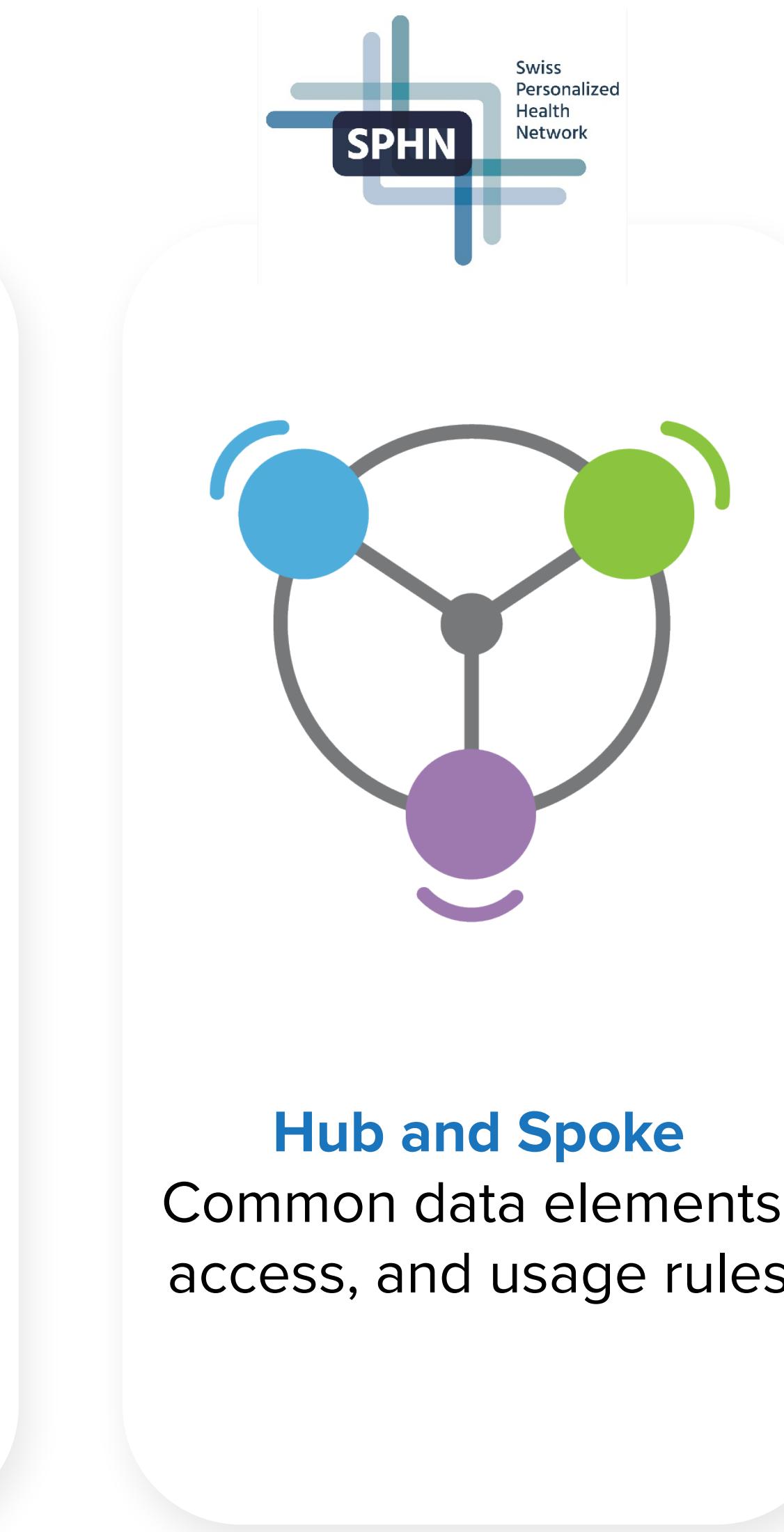
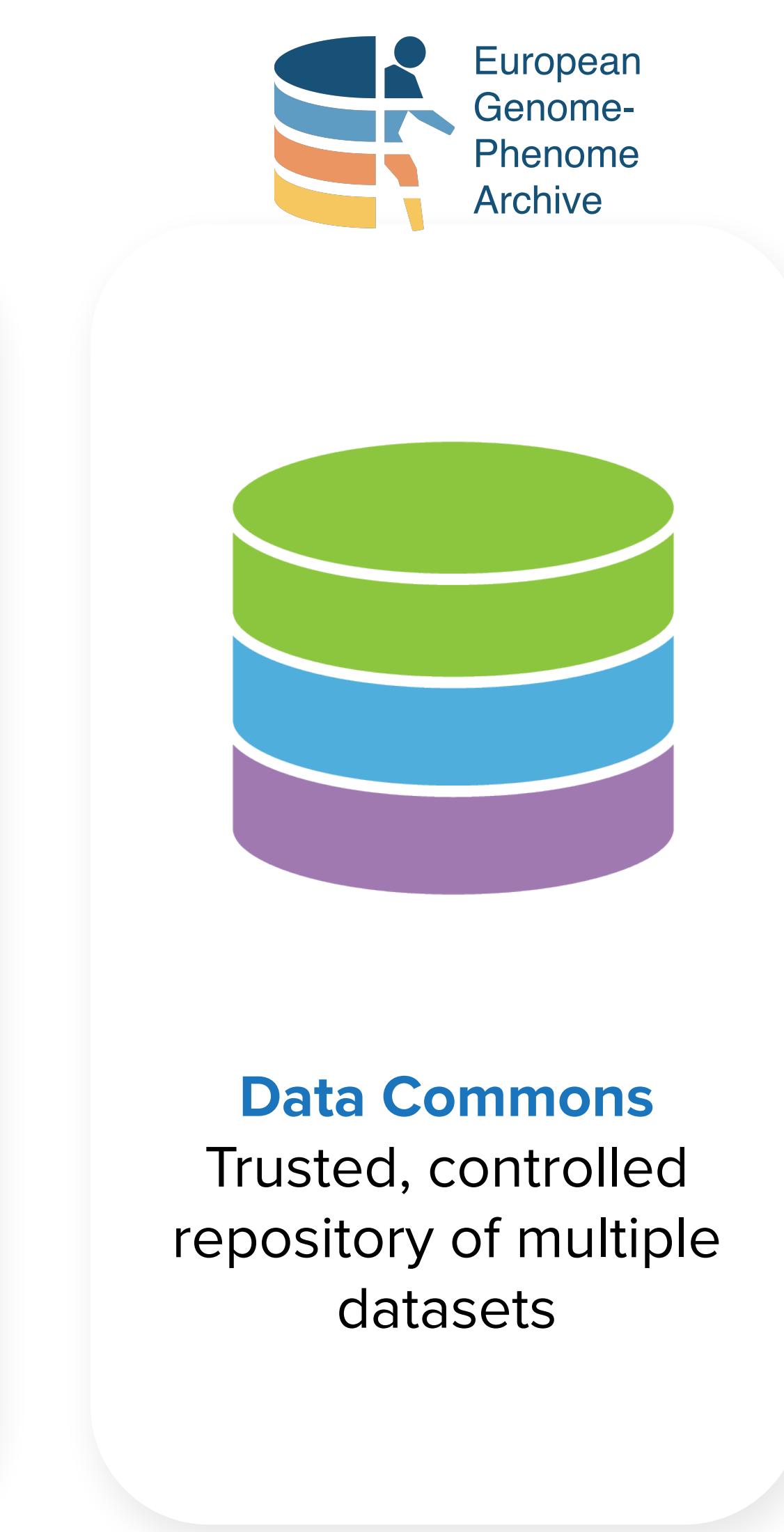
Where does Genomic Data Come From?

Geographic bias in published cancer genome profiling studies



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.

Different Approaches to Data Sharing



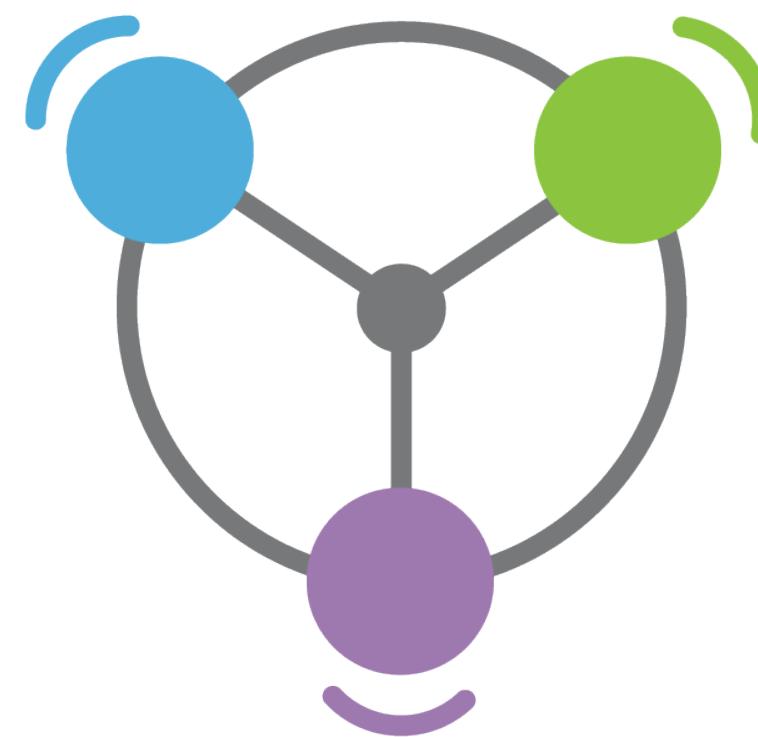
Different Approaches to Data Sharing



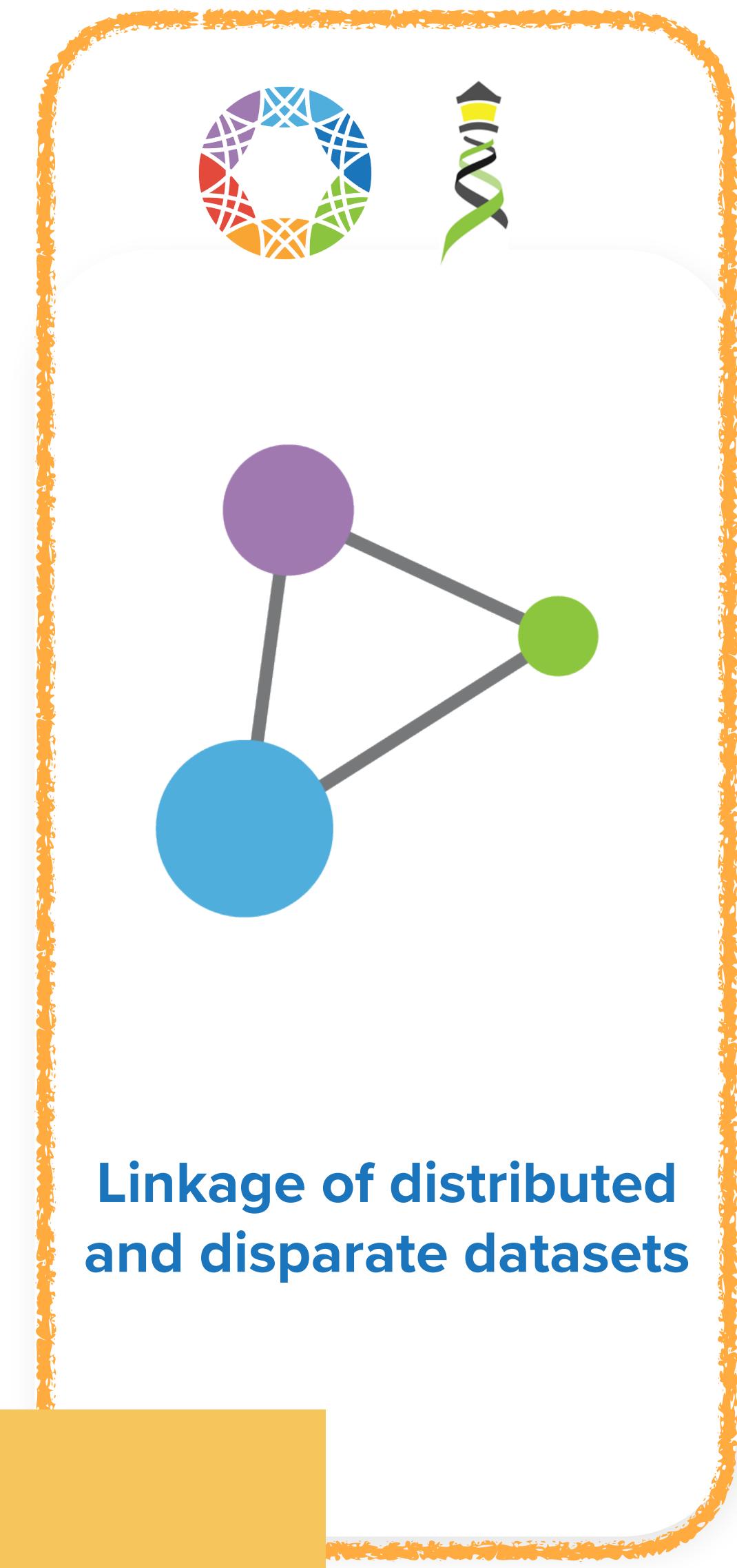
Centralized Genomic Knowledge Bases



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

INFORMATICS

Beacon v2 and Beacon networks: federated data discovery in biome

Commentary

International federation of genomic medicine databases using GA4GH standards

Adrian Thorogood,^{1,2,*} Heidi L. Rehm,^{3,4} Peter Goodhand,^{5,6} Angela J.H. Page,^{4,5} Yann Joly,² Michael Baudis,⁷ Jordi Rambla,^{8,9} Arcadi Navarro,^{8,10,11,12} Tommi H. Nyronen,^{13,14} Mikael Linden,^{13,14} Edward S. Dove,¹⁵ Marc Fiume,¹⁶ Michael Brudno,¹⁷ Melissa S. Cline,¹⁸ and Ewan Birney¹⁹

Jordi Rambla^{1,2} | Michael Baudis³ | Roberto Ariosa¹ | Tim Beck⁴ |
 Lauren A. Fromont¹ | Arcadi Navarro^{1,5,6,7} | Rahel Paloots³ |
 Manuel Rueda¹ | Gary Saunders⁸ | Babita Singh¹ | John D. Spalding⁹ |
 Juha Törnroos⁹ | Claudia Vasallo¹ | Colin D. Veal⁴ | Anthony J. Brookes⁴

Cell Genomics

Technology

The GA4GH Variation Representation Specification A computational framework for variation representation and federated identification

Alex H. Wagner,^{1,2,25,*} Lawrence Babb,^{3,*} Gil Alterovitz,^{4,5} Michael Baudis,⁶ Matthew Brush,⁷ Daniel L. Cameron,^{8,9} Melissa Cline,¹⁰ Malachi Griffith,¹¹ Obi L. Griffith,¹¹ Sarah E. Hunt,¹² David Kreda,¹³ Jennifer M. Lee,¹⁴ Stephanie Li,¹⁵ Javier Lopez,¹⁶ Eric Moyer,¹⁷ Tristan Nelson,¹⁸ Ronak Y. Patel,¹⁹ Kevin Riehle,¹⁹ Peter N. Robinson,²⁰ Shawn Rynearson,²¹ Helen Schuilenburg,¹² Kirill Tsukanov,¹² Brian Walsh,⁷ Melissa Konopko,¹⁵ Heidi L. Rehm,^{3,22} Andrew D. Yates,¹² Robert R. Freimuth,²³ and Reece K. Hart^{3,24,*}

Cell Genomics

Perspective

GA4GH: International policies and standards for data sharing across genomic research and healthcare

Heidi L. Rehm,^{1,2,47} Angela J.H. Page,^{1,3,*} Lindsay Smith,^{3,4} Jeremy B. Adams,^{3,4} Gil Alterovitz,^{5,47} Lawrence J. Babb,¹ Maxmillian P. Barkley,⁶ Michael Baudis,^{7,8} Michael J.S. Beauvais,^{3,9} Tim Beck,¹⁰ Jacques S. Beckmann,¹¹ Sergi Beltran,^{12,13,14} David Bernick,¹ Alexander Bernier,⁹ James K. Bonfield,¹⁵ Tiffany F. Boughtwood,^{16,17} Guillaume Bourque,^{9,18} Sarion R. Bowers,¹⁵ Anthony J. Brookes,¹⁰ Michael Brudno,^{18,19,20,21,38} Matthew H. Brush,²² David Bujold,^{9,18,38} Tony Burdett,²³ Orion J. Buske,²⁴ Moran N. Cabili,¹ Daniel L. Cameron,^{25,26} Robert J. Carroll,²⁷ Esmeralda Casas-Silva,¹²³ Debyani Chakravarty,²⁹ Bimal P. Chaudhari,^{30,31} Shu Hui Chen,³² J. Michael Cherry,³³ Justina Chung,^{3,4} Melissa Cline,³⁴ Hayley L. Clissold,¹⁵ Robert M. Cook-Deegan,³⁵ Mélanie Courtot,²³ Fiona Cunningham,²³ Miro Cupak,⁶ Robert M. Davies,¹⁵ Danielle Denisko,¹⁹ Megan J. Doerr,³⁶ Lena I. Dolman,¹⁹

(Author list continued on next page)



Global Alliance for Genomics & Health

GENOMICS

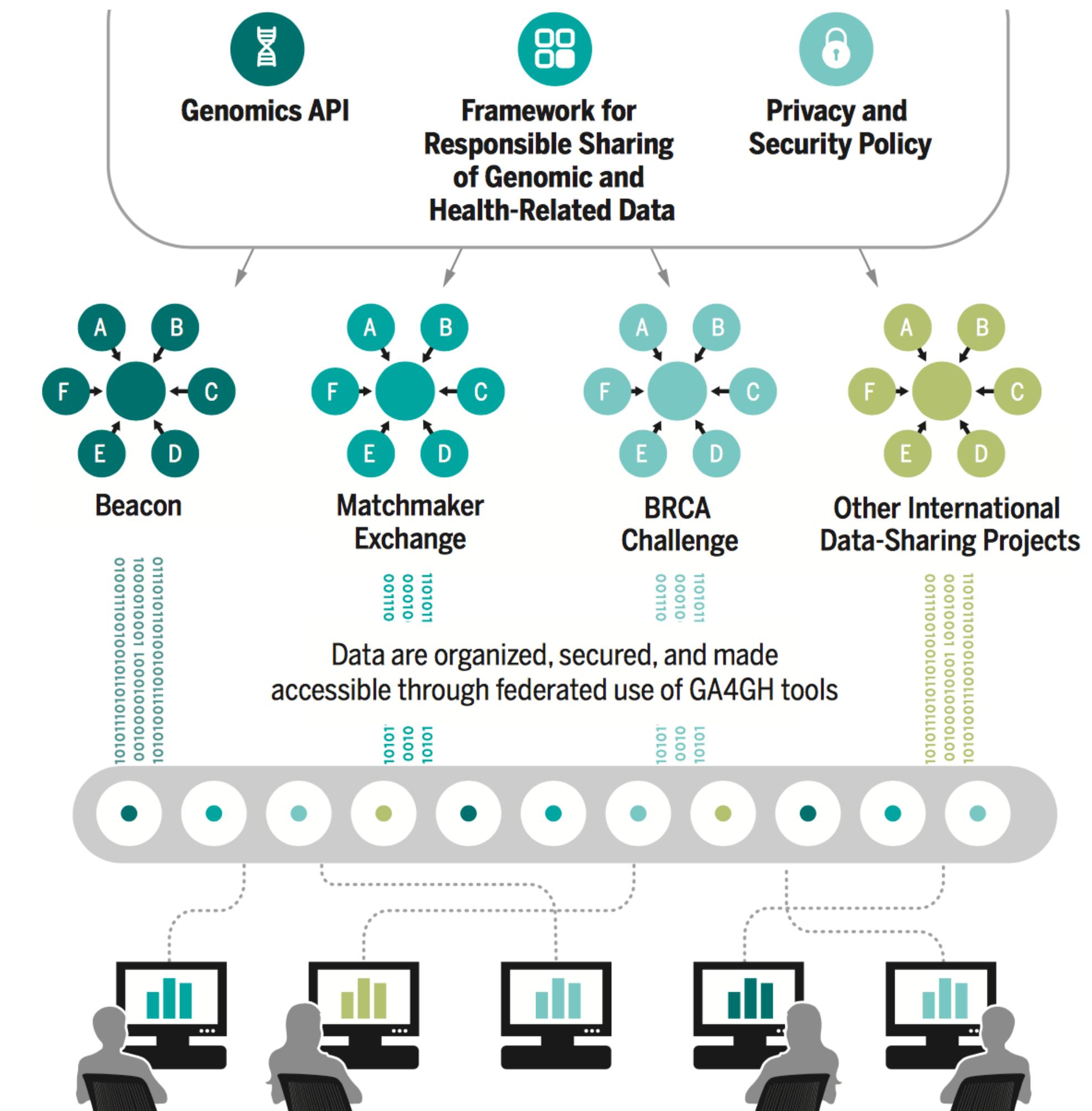
A federated ecosystem for sharing genomic, clinical data

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

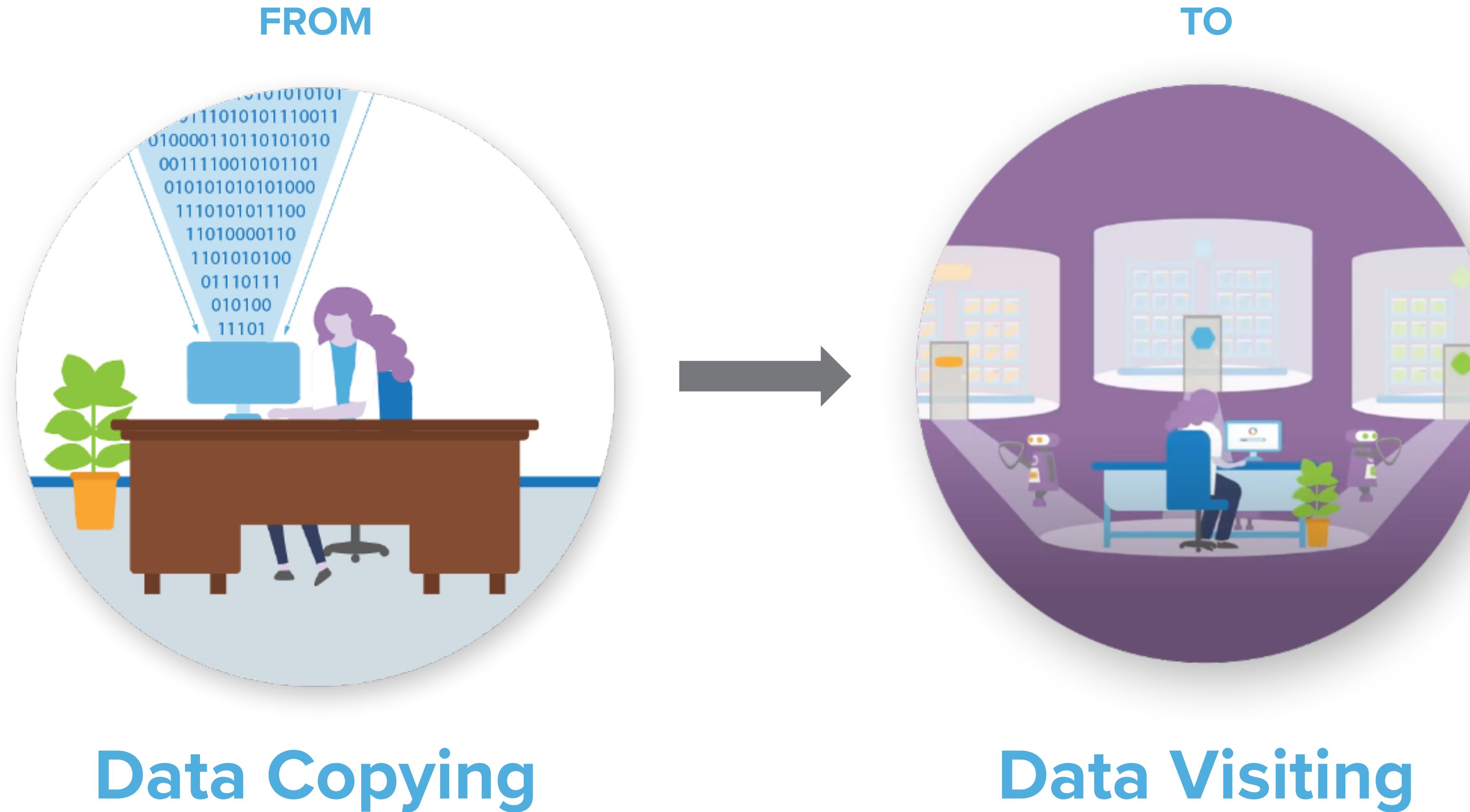
The Global Alliance for Genomics and Health*

SCIENCE 10 JUNE 2016 • VOL 352 ISSUE 6299

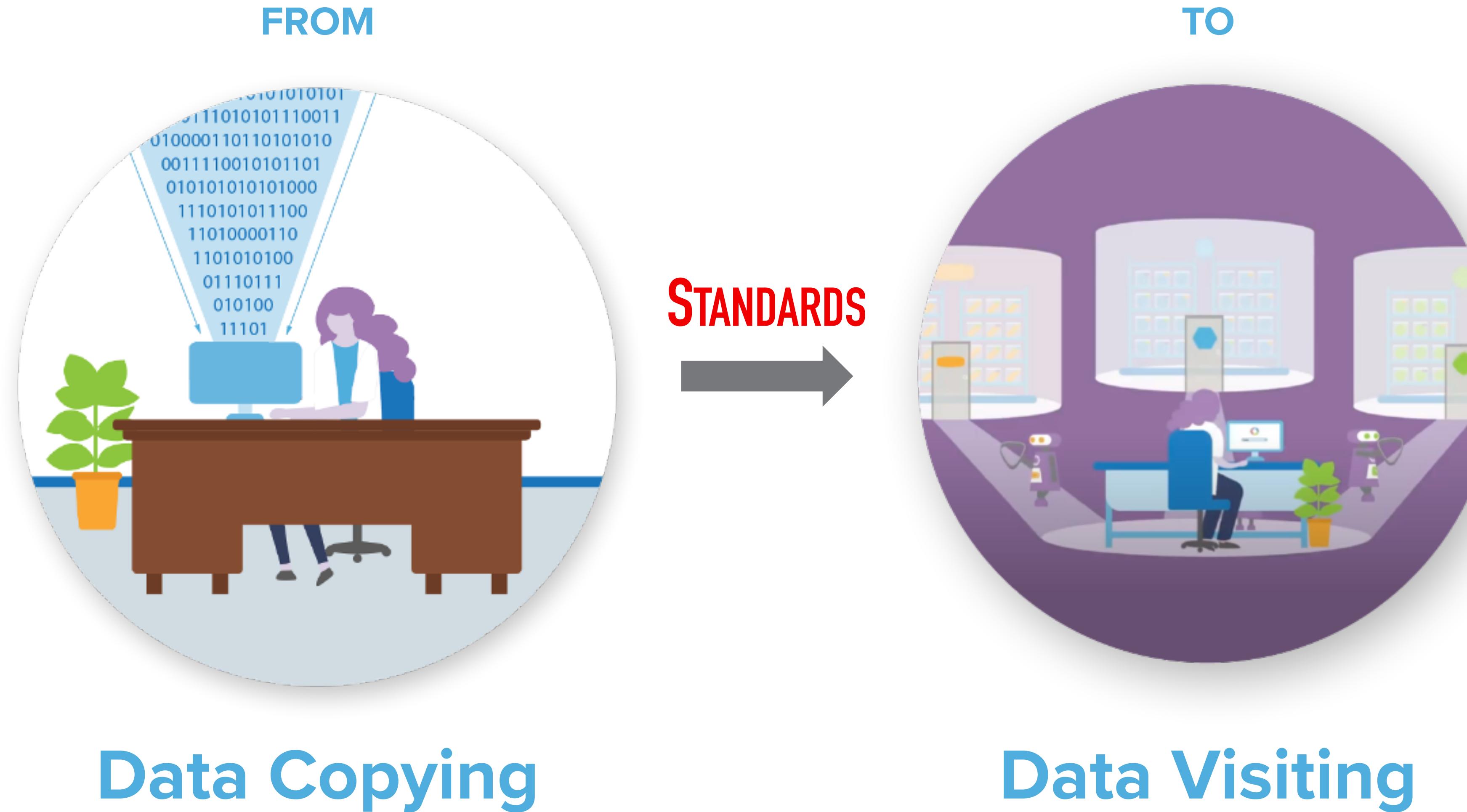
A federated data ecosystem. To share genomic data globally, this approach furthers medical research without requiring compatible data sets or compromising patient identity.



A New Paradigm for Data Sharing



A New Paradigm for Data Sharing



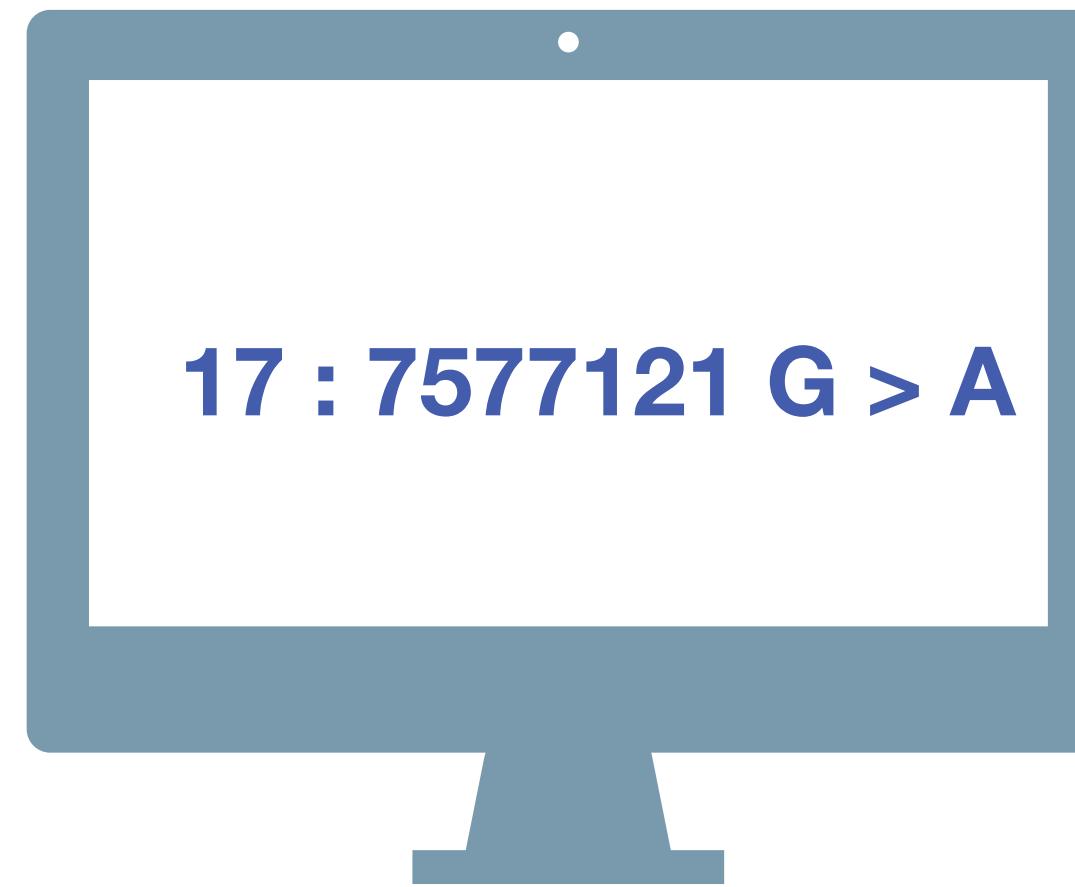


Global Alliance
for Genomics & Health
Collaborate. Innovate. Accelerate.



The GA4GH Beacon Protocol

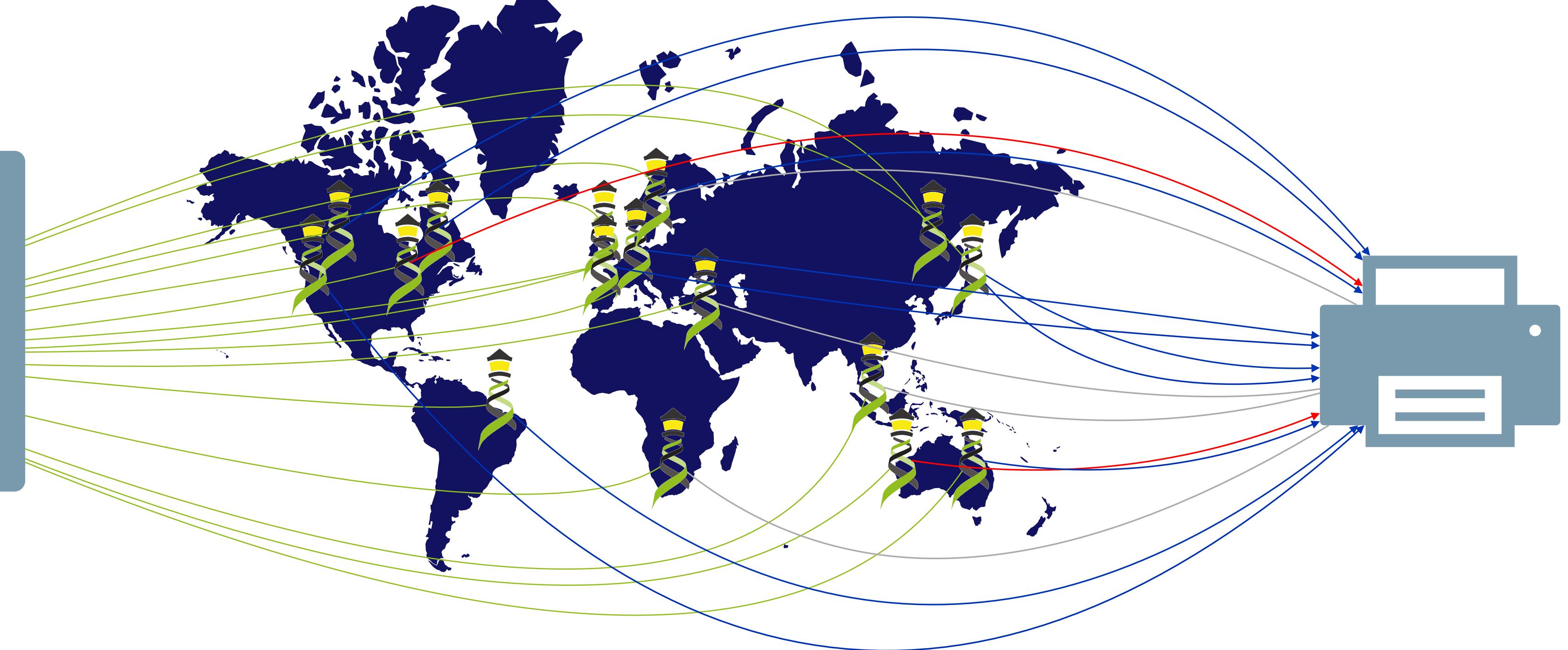
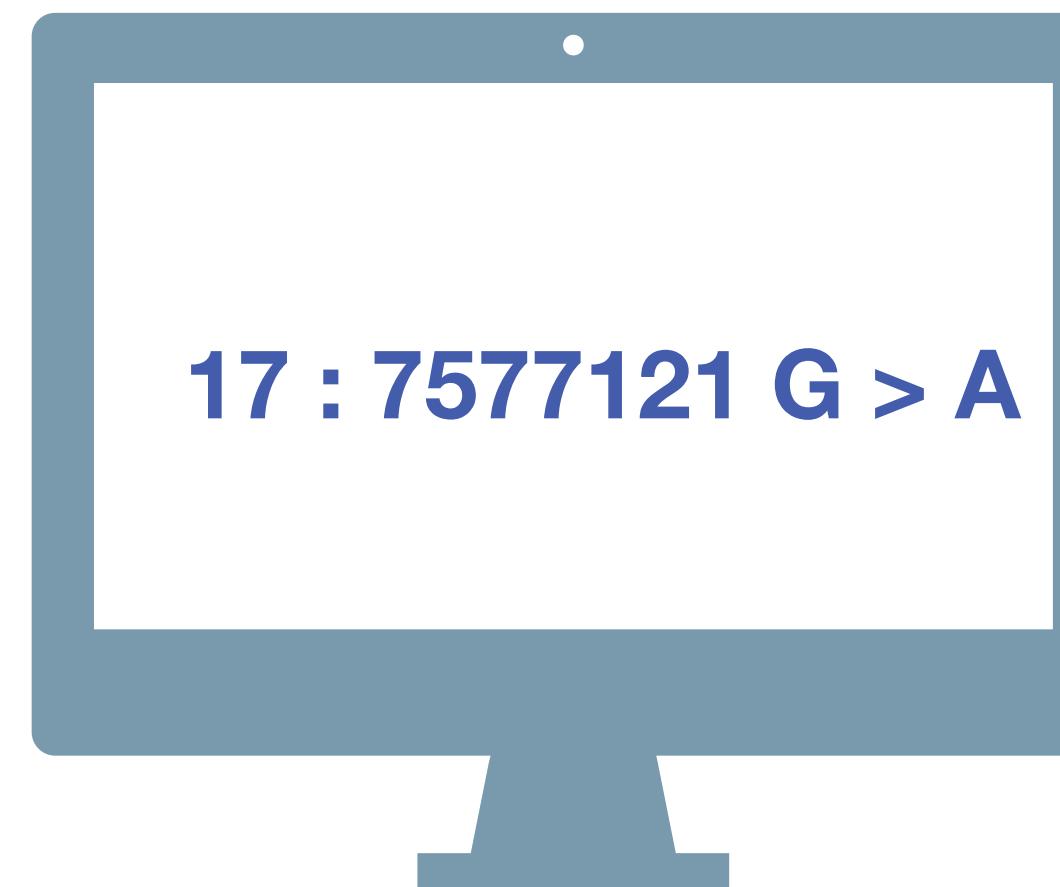
Federating Genomic Discoveries



Beacon

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

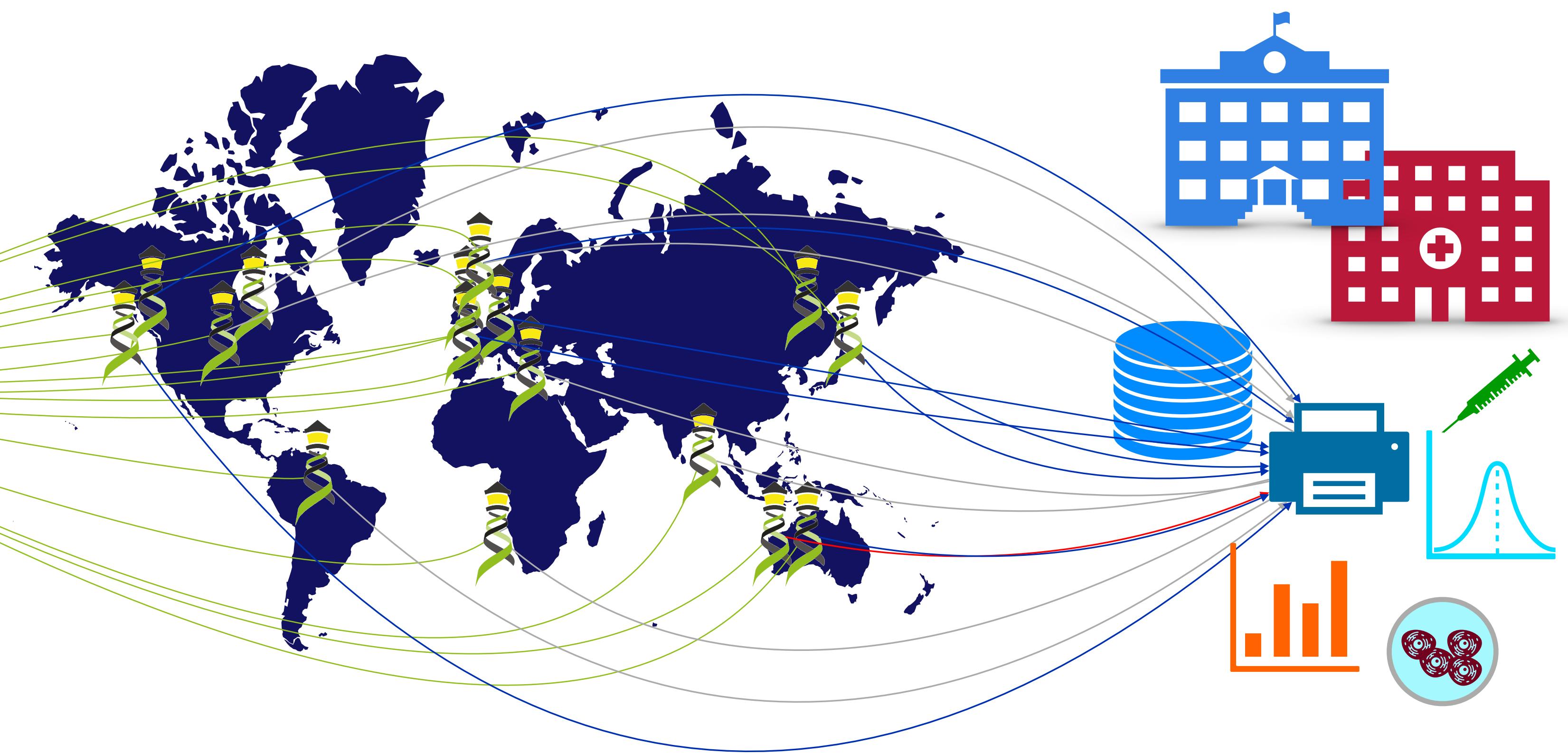
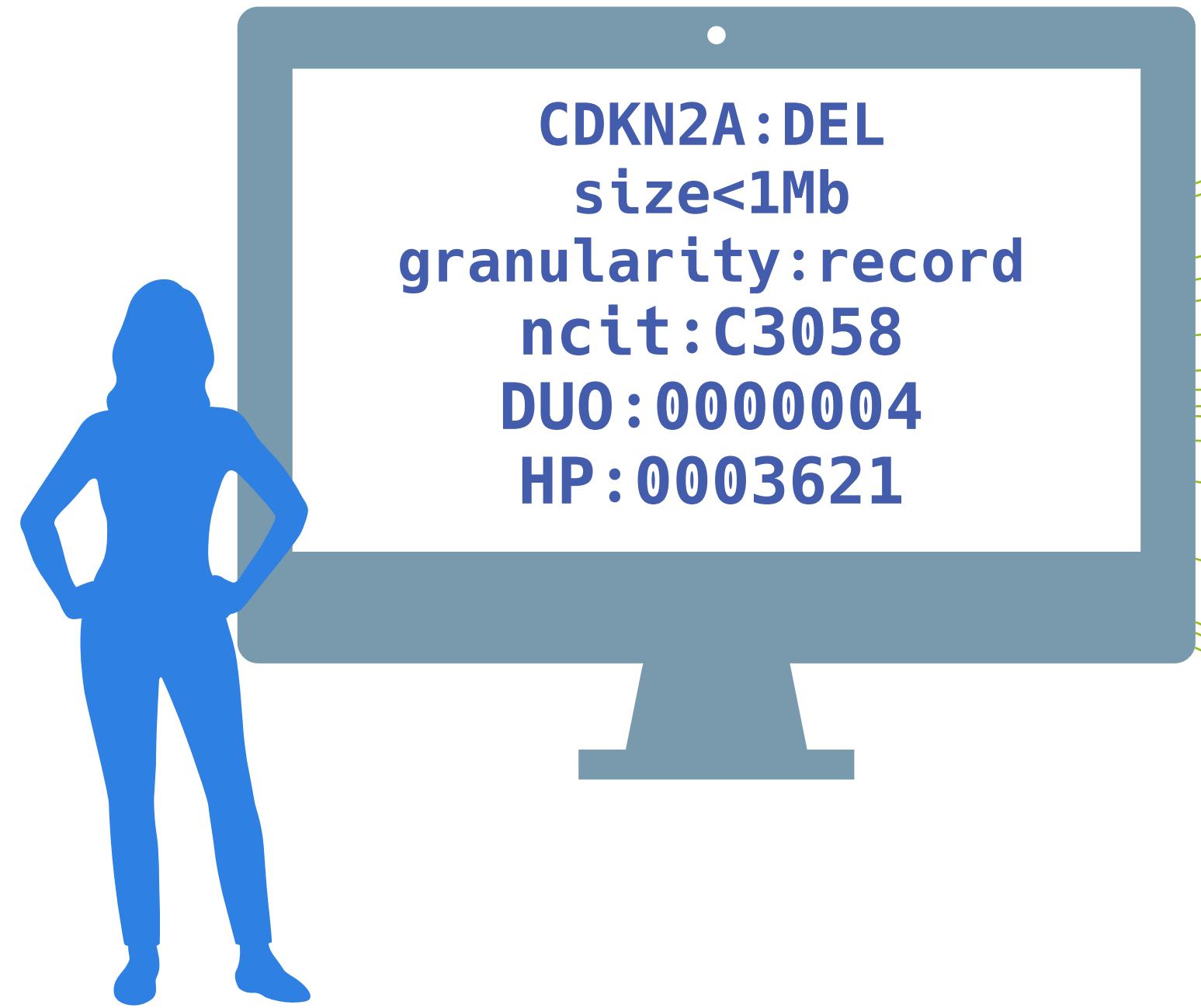
YES | NO | \0



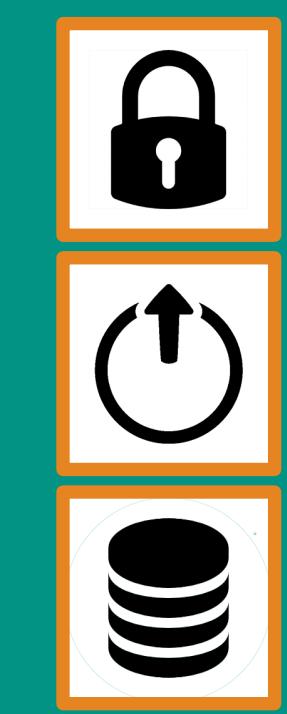
Have you seen this variant?
It came up in my patient
and we don't know if this is
a common SNP or worth
following up.

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**.



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



Beacon API

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

bycon Beacon+

Implementation driven standards development

- Progenetix' Beacon+ has served as implementation driver since 2016
- the *bycon* package is used to prototype advanced Beacon features such as
 - structural variant queries
 - data handovers
 - Phenopackets integration
 - variant co-occurrences
 - ...

Beacon protocol response verifier at time of GA4GH approval Spring 2022

Beacon v2 GA4GH Approval Registry

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

European Genome-Phenome Archive (EGA)

GA4GH Approval Beacon Test

This [Beacon](#) is based on the GA4GH Beacon [v2.0](#)

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

progenetix

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	Green
Beacon UI	Green
Beacon API	Green
Contact us	Green
BeaconMap	Green
Bioinformatics analysis	Green
Biological Sample	Green
Cohort	Green
Configuration	Green
Dataset	Green
EntryTypes	Green
Genomic Variants	Green
Individual	Green
Info	Green
Sequencing run	Green

cnag

Centre Nacional Analisis Genomica (CNAG-CRG)

Beacon @ RD-Connect

This [Beacon](#) is based on the GA4GH Beacon [v2.0](#)

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

University of Leicester

Cafe Variome Beacon v2

This [Beacon](#) is based on the GA4GH Beacon [v2.0](#)

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

✓ Matches the Spec ✗ Not Match the Spec ● Not Implemented



Website populated by asynchronous retrieval of Beacon query results using handovers

 [Edit Query](#)

CNV Profiles
... by NCIT
... by ICD-O Morphology
... by ICD-O Site
... by TNM & Grade

Search Samples

arrayMap
TCGA Data
cBioPortal Studies

Publication DB
Progenetix Use

NCIT - ICD-O Mappings
UBERON Mappings

Upload & Plot

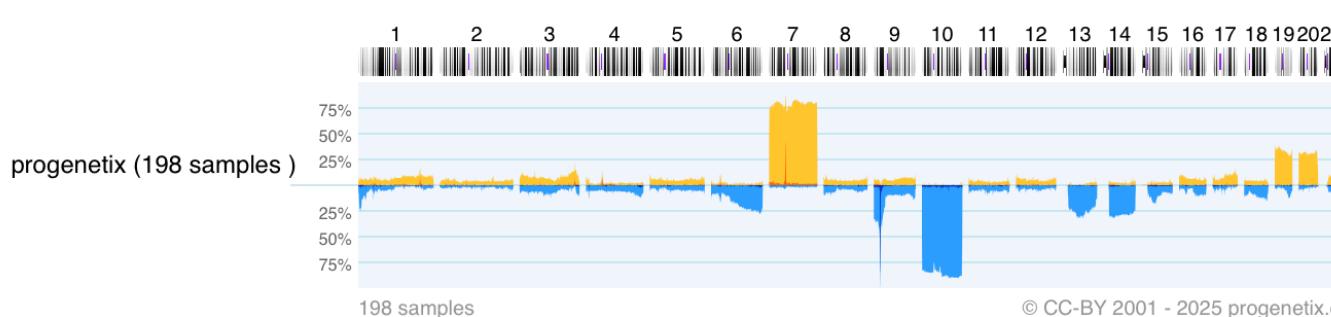
OpenAPI Paths and Examples

Cancer Cell Lines

progenetix

Matched Samples: 969 UCSC region 
Retrieved Samples: 200 Geographic Map 
Variants: 984 Variants in UCSC 
Calls: 976 Dataset Responses (JSON) 

[Results](#) [Biosamples](#) [Variants](#)



[Reload histogram in new window](#) 

Matched Subset Codes	Subset Samples	Matched Samples	Subset Match Frequencies
pgx:icdot-C71.4	4	1	0.250
pgx:icdot-C71.1	14	1	0.071
pgx:icdom-94403	4816	200	0.042
NCIT:C3058	4900	200	0.041
pgx:icdot-C71.9	13758	192	0.014
pgx:icdot-C71.0	1714	6	0.004

progenetix Data Downloads

Download Sample Data (TSV)
Part1  Part2  Part3  Part4  Part5 

Download Sample Data (JSON)
Part1  Part2  Part3  Part4  Part5 

Download Variants (Beacon VRS)
Part1  Part2  Part3  Part4  Part5 

Download Variants (VCF)

Results			
Biosample	Dx Classifications	Identifiers	Variants
pgxbs-kftvl1hz	pgx:icdom-94403 Glioblastoma, NOS pgx:icdot-C71.9 Brain, NOS NCIT:C3058 Glioblastoma	pubmed:28481359 Zehir A, Benayed R et al. (2017): Mutational landscape of metastatic cancer revealed... cbiportal:msk_impact_2017	
pgxbs-kftvl7f4	pgx:icdom-94403 Glioblastoma, NOS pgx:icdot-C71.9 Brain, NOS NCIT:C3058 Glioblastoma	pubmed:28481359 Zehir A, Benayed R et al. (2017): Mutational landscape of metastatic cancer revealed... cbiportal:msk_impact_2017	
pgxbs-kftvhm6s	pgx:icdom-94403 Glioblastoma, NOS pgx:icdot-C71.9 Brain, NOS NCIT:C3058 Glioblastoma	pgx:TCGA-GBM Glioblastoma Multiforme 18772890 Cancer Genome Atlas Research Network. (2008): Comprehensive genomic characterization defines human glioblastoma...	
Biosamples			
Digest	Gene	Pathogenicity	Variant type
9:21626201- 21981584:EFO_0030068			CopyNumberChange V: pgxvar- 665749ab2d6be9a260e55de8 A: pgxcs-kftwnmzs B: pgxbs-kftvjywz I: pgxind-kftx5yjj
9:21846286- 22201587:EFO_0030068			CopyNumberChange V: pgxvar- 6656fc5fbe3f6845a3555b82 A: pgxcs-kftw53z6 B: pgxbs-kftvi872 I: pgxind-kftx3t9l
9:21949762- 22004847:EFO_0020073			CopyNumberChange V: pgxvar- 6657226e8f6b96158261aa6 A: pgxcs-kftw3vh5 B: pgxbs-kftvi4e1 I: pgxind-kftx3vx4
9:21164528- 21990552:EFO_0020073			CopyNumberChange V: pgxvar- 66572dfe2d6be9a260e3d189 A: pgxcs-kftw95rl B: pgxbs-kftvilhz I: pgxind-kftx4au8



front-end showcases query strategies

Beacon Search Demonstrator

This search form shows parameter combinations and examples for different Beacon search patterns. Please be aware that search types and examples are *independent* of each other, so not all combinations are automatically adjusted.

Additionally, the search options here might extend the latest stable version of the Beacon API in a sense of "implementation driven development" but are supported through this version of the [bycon](#) library.

Search Samples

Compare CNV Profiles

CNV Profiles by Cancer Type

NCIT Neoplasia Codes
ICD-O Morphologies
ICD-O Organ Sites
TNM & Grade

OpenAPI Paths and Examples

Documentation

Progenetix

Baudisgroup @ UZH

Search Samples

CNV (Bracket) Range Gene ID Sequence Genomic Fusion Sample Data

Dataset(s) Test Database - examplez

Variant Type EFO:0030067 (copy number deletion)

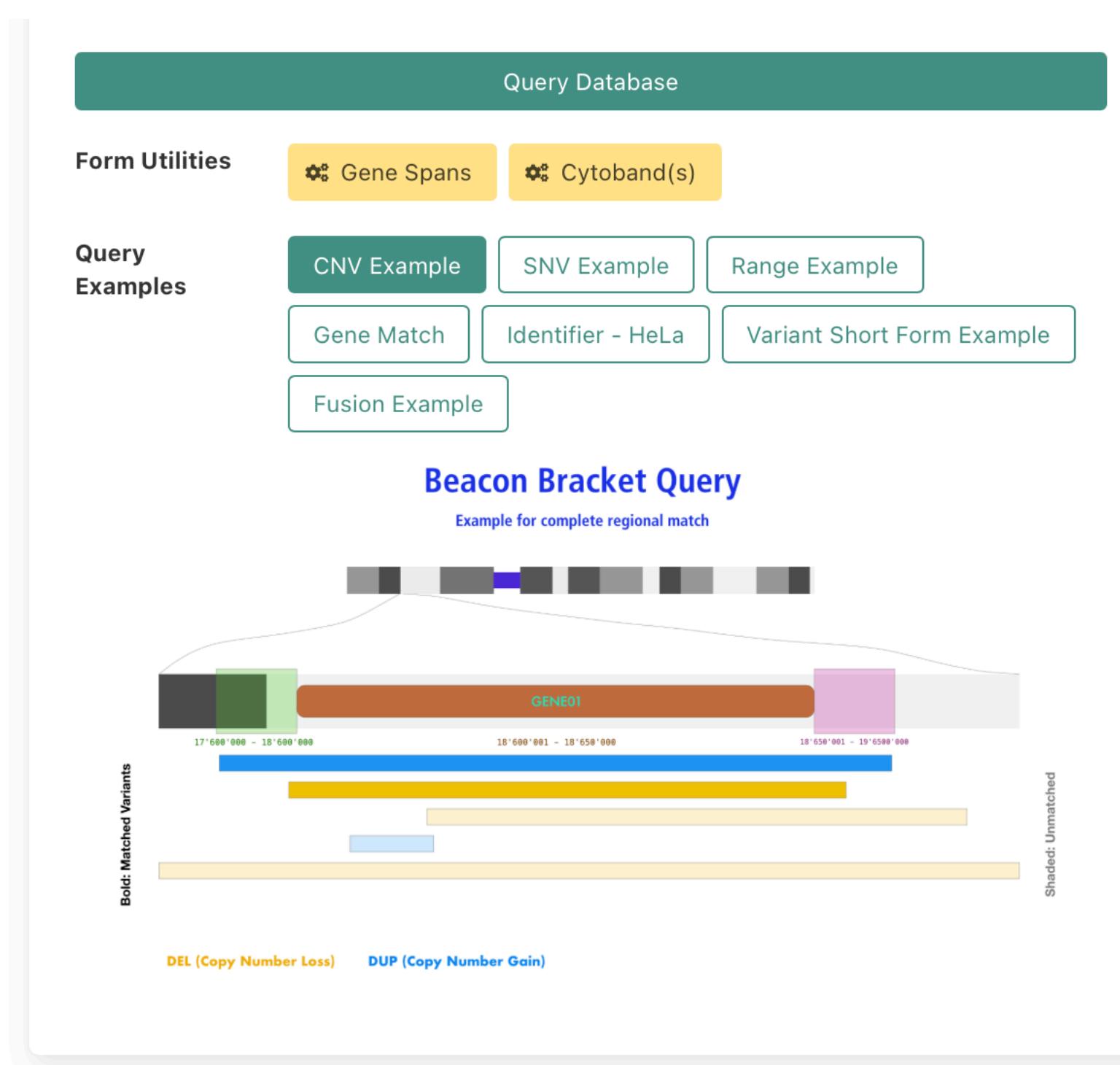
Chromosome 9 (NC_000009.12) **Start or Position** 21000001,21975098 **End (Range or Structural Var.)** 21967753,23000000

Cancer Classification(s) Select...

Genotypic Sex Select...

Various Subsets NCIT:C3058: Glioblastoma (28)

Chromosome 9 21000001,21975098 21967753,23000000

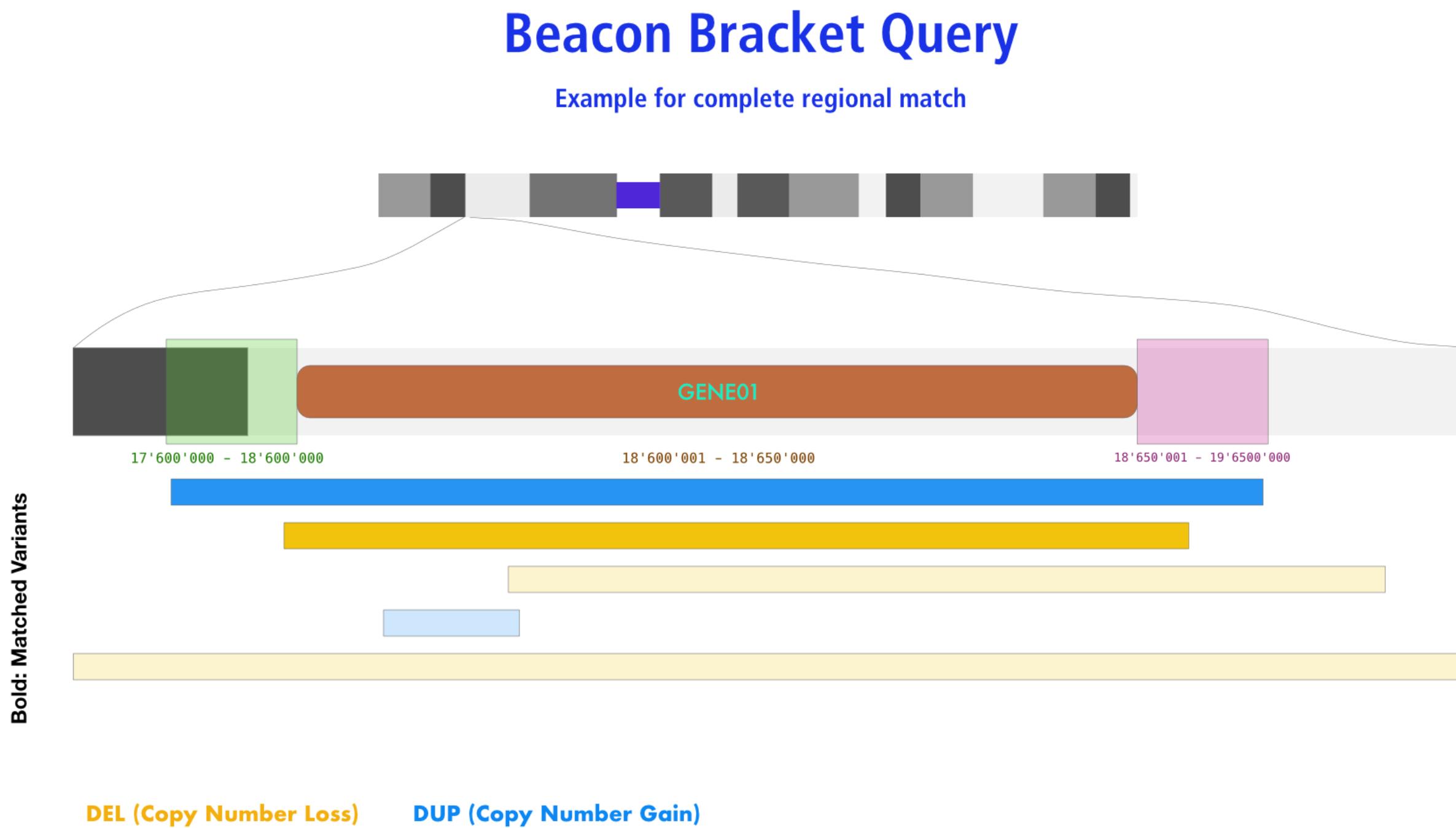


- React based website included in *bycon*
- exposes Beacon query options
- package provided examples can be extended by local ones

Beacon Queries

Bracket ("CNV") Query

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



Beacon Query Types

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

Dataset

Test Database - examplez x | v

Chromosome

9 (NC_000009.12) | v

Variant Type

EFO:0030067 (copy number deletion) | v

Start or Position

21000001-21975098

End (Range or Structural Var.)

21967753-23000000

Select Filters

NCIT:C3058: Glioblastoma (100) x | v

Chromosome 9

21000001-21975098



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.

Standards Development & Implementation: CNV Terms

in computational (file/schema) formats

- EFO:0030064
- EFO:0030067
 - | - EFO:0030068
 - \ - EFO:0020073
 - \ - EFO:0030069
- EFO:0030070
 - | - EFO:0030071
 - \ - EFO:0030072

GA4GH VRS1.3+	Beacon v2	VCF v4.4	SO
EFO:0030070 gain	DUP or EFO:0030070	DUP SVCLAIM=D	SO:0001742 copy_number_gain
EFO:0030071 low-level gain	DUP or EFO:0030071	DUP SVCLAIM=D	SO:0001742 copy_number_gain
EFO:0030072 high-level gain	DUP or EFO:0030072	DUP SVCLAIM=D	SO:0001742 copy_number_gain
EFO:0030072 high-level gain	DUP or EFO:0030073	DUP SVCLAIM=D	SO:0001742 copy_number_gain
EFO:0030067 loss	DEL or EFO:0030067	DEL SVCLAIM=D	SO:0001743 copy_number_loss
EFO:0030068 low-level loss	DEL or EFO:0030068	DEL SVCLAIM=D	SO:0001743 copy_number_loss
EFO:0020073 high-level loss	DEL or EFO:0020073	DEL SVCLAIM=D	SO:0001743 copy_number_loss
EFO:0030069 complete genomic loss	DEL or EFO:0030069	DEL SVCLAIM=D	SO:0001743 copy_number_loss

Website populated by asynchronous retrieval of Beacon query results using handovers

progenetix

[Edit Query](#)

CNV Profiles
... by NCIT
... by ICD-O Morphology
... by ICD-O Site
... by TNM & Grade

Search Samples

arrayMap
TCGA Data
cBioPortal Studies

Publication DB
Progenetix Use

NCIT - ICD-O Mappings
UBERON Mappings

Upload & Plot

OpenAPI Paths and Examples

Cancer Cell Lines

Chro: refseq:NC_000009.12 **Start:** 21000001,21975098 **End:** 21967753,23000000 **Type:** EFO:0030067
Filters: NCIT:C3058

progenetix

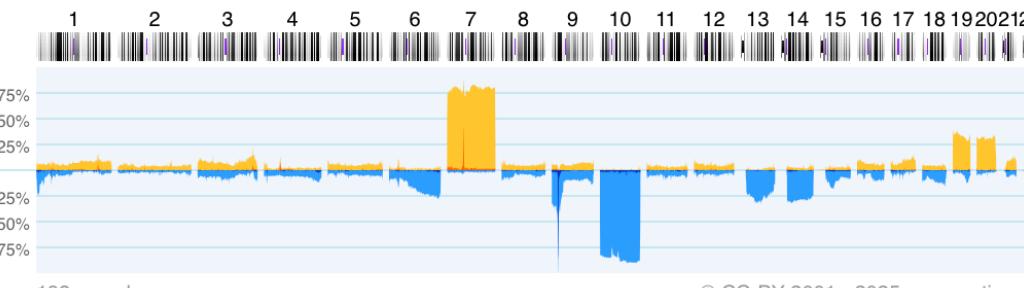
Matched Samples: 969 Retrieved Samples: 200 Variants: 984 Calls: 976

UCSC region Geographic Map Variants in UCSC Dataset Responses (JSON)

[Visualization options](#)

[Results](#) [Biosamples](#) [Variants](#)

progenetix (198 samples)



© CC-BY 2001 - 2025 progenetix.org

[Reload histogram in new window](#)

Matched Subset Codes	Subset Samples	Matched Samples	Subset Match Frequencies
pgx:icdot-C71.4	4	1	0.250
pgx:icdot-C71.1	14	1	0.071
pgx:icdom-94403	4816	200	0.042
NCIT:C3058	4900	200	0.041
pgx:icdot-C71.9	13758	192	0.014
pgx:icdot-C71.0	1714	6	0.004

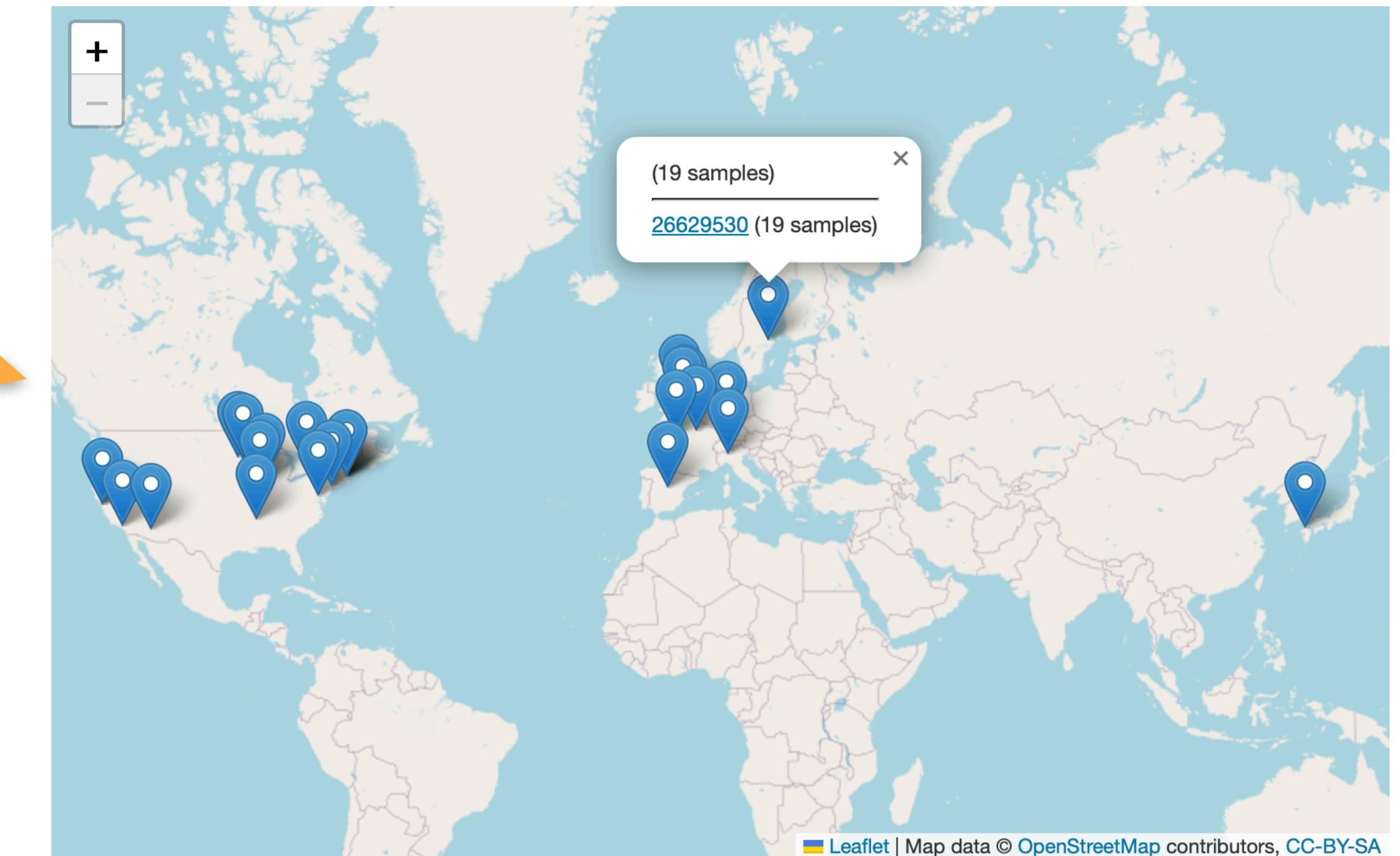
progenetix Data Downloads

[Download Sample Data \(TSV\)](#)
Part1 Part2 Part3 Part4 Part5

[Download Sample Data \(JSON\)](#)
Part1 Part2 Part3 Part4 Part5

[Download Variants \(Beacon VRS\)](#)
Part1 Part2 Part3 Part4 Part5

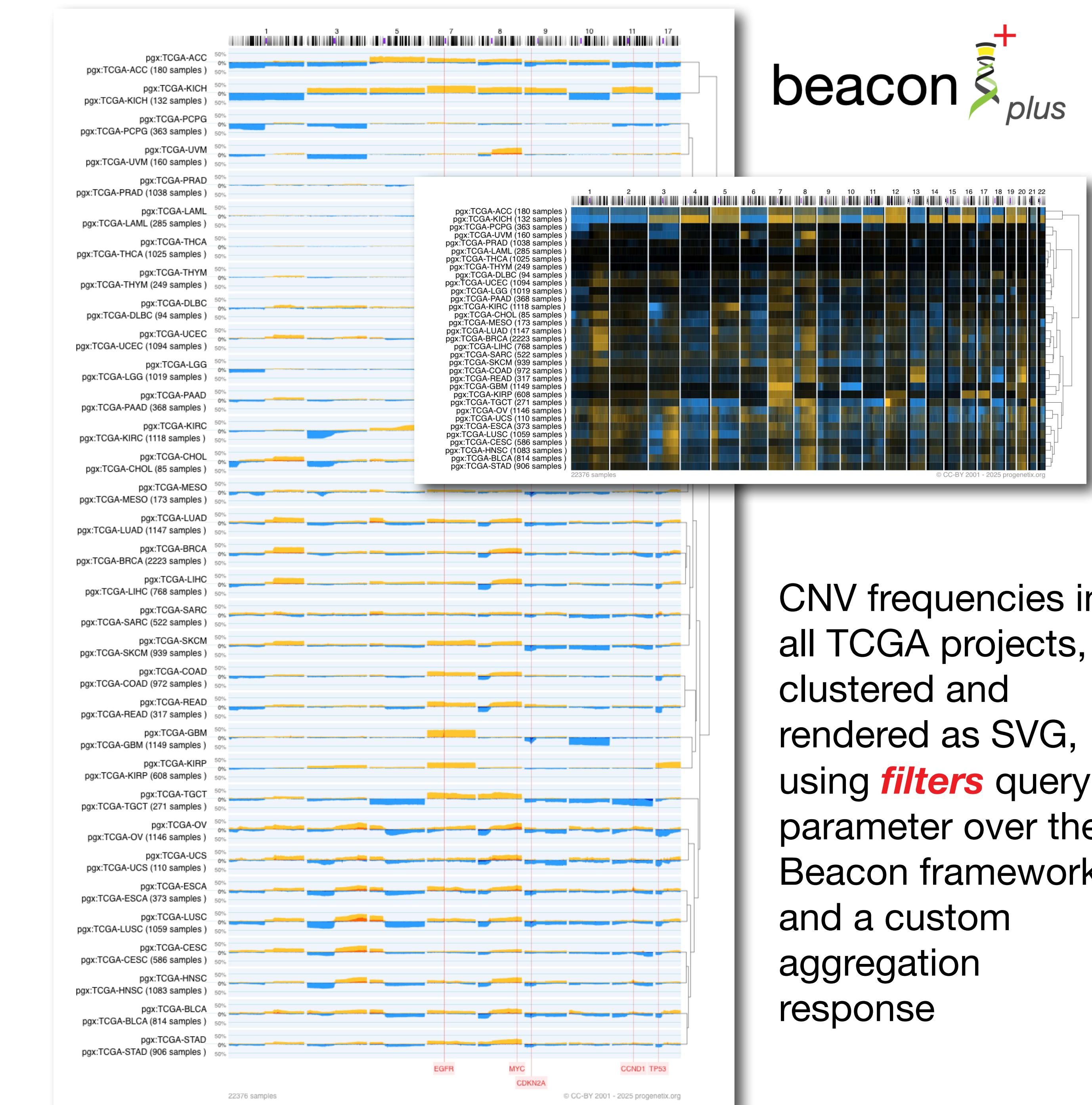
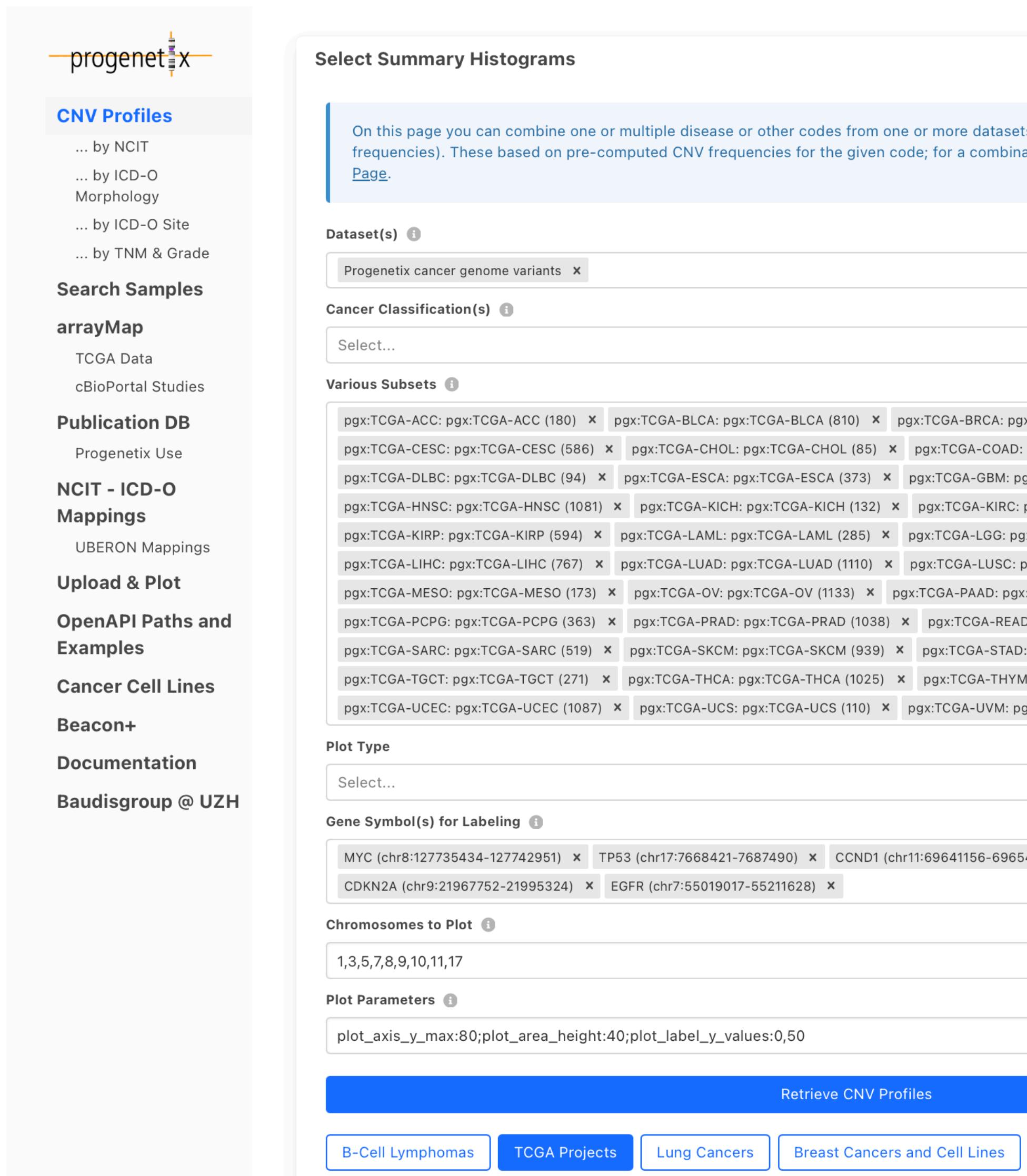
[Download Variants \(VCF\)](#)
Part1 Part2 Part3 Part4 Part5



Pushing the standard: Biosamples in Progenetix have geographic attribution in the form of GeoJSON objects, for query & display...

Pushing the envelope...

Custom Beacon aggregation response for displaying CNV frequencies

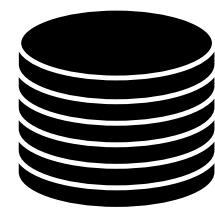


CNV frequencies in all TCGA projects, clustered and rendered as SVG, using ***filters*** query parameter over the Beacon framework and a custom aggregation response

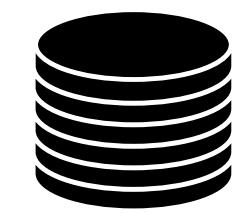
bycon based Beacon+ Stack

progenetix

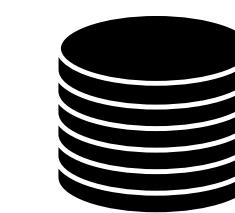
- *collations* contain pre-computed data (e.g. CNV frequencies, statistics) and information for all grouping entity instances and correspond to **filter values**
 - ▶ [pubmed:10027410](#), [NCIT:C3222](#), [pgx:cohort-TCGA](#), [pgx:icdom-94703...](#)
 - ▶ precomputed frequencies per collection informative e.g. in form autfills
- *querybuffer* stores id values of all entities matched by a query and provides the corresponding **accessid** for **handover** generation
- complete query aggregation; i.e. individual queries are run against the corresponding entities and ids are intersected
 - retrieval of any entity, e.g. all individuals which have queried variants analyzed on a given platform
 - allows multi-variant queries, i.e. all bio samples or individuals which had matches of all of the individual variant queries



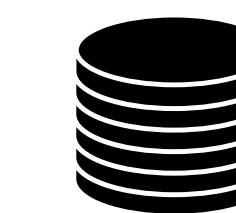
variants



analyses



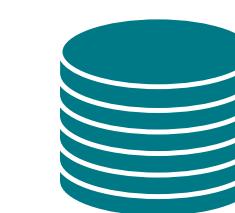
biosamples



individuals



collations



geolocs



genespans

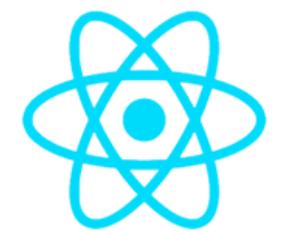


qBuffer

Entity collections

Utility collections

github.com/progenetix/bycon



React



ga4gh-beacon / **beacon-v2**

Type / to search

Code Issues Pull requests Discussions Actions Projects Security Insights Settings

beacon-v2 Public

Edit Pins Unwatch 10 Fork 22 Starred 32

add-aggregation-resp... 37 Branches 5 Tags Go to file Add file Code About

This branch is 25 commits ahead of main.

#259

mbaudis re-adding distributions 9682bed · 2 days ago 717

.github/workflows adding github actions demo file

bin fixes for aggregation PR

docs fixes for aggregation PR

framework re-adding distributions

models measures => measurements re-fix (this branch)

.gitattributes re-structuring intro pages

.gitignore Fix file naming conflict error in schemas-md on macOS A...

CHANGELOG.md Merge branch 'main' into schema-urgent-fixes

LICENSE Initial commit

README.md Merge branch 'develop' into develop_changelog

mkdocs.yaml some v2 naming/version use cleanup

requirements.txt switch to mermaid2 plugin

README License Security

Unified repository for Beacon Code & Documentation

progenetix / **bycon**

Type / to search

Code Issues Pull requests Actions Projects Wiki Security Insights Settings

bycon Public

Edit Pins Unwatch 3 Fork 10 Starred 8

main 10 Branches 49 Tags Go to file Add file Code About

mbaudis Merge pull request #44 from mbaudis/main 68ee58b · last month 931 Commits

.github/workflows docs & formatting last month

beaconServer 2.4.3 "Bologna" 4 months ago

beaconplusWeb vrsifier and vrs format last month

bycon VCF sequence fix; some clean-up last month

byconServices going VRSv2 alpha last month

docs VCF sequence fix; some clean-up last month

housekeepers going VRSv2 alpha last month

importers refactor importers 2 months ago

local vrsifier and vrs format last month

rsrc going VRSv2 alpha last month

tests going VRSv2 alpha last month

.gitignore 2.1.2 9 months ago

LICENSE Create LICENSE 5 years ago

README.md export tables last month

install.py 2.4.9 2 months ago

markdowner.py v2.4.7 "Thessaloniki" 3 months 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

Custom properties

8 stars

3 watching

10 forks

Report repository

Releases 15

v2.5.0 "Forked" Latest on Jul 30

+ 14 releases

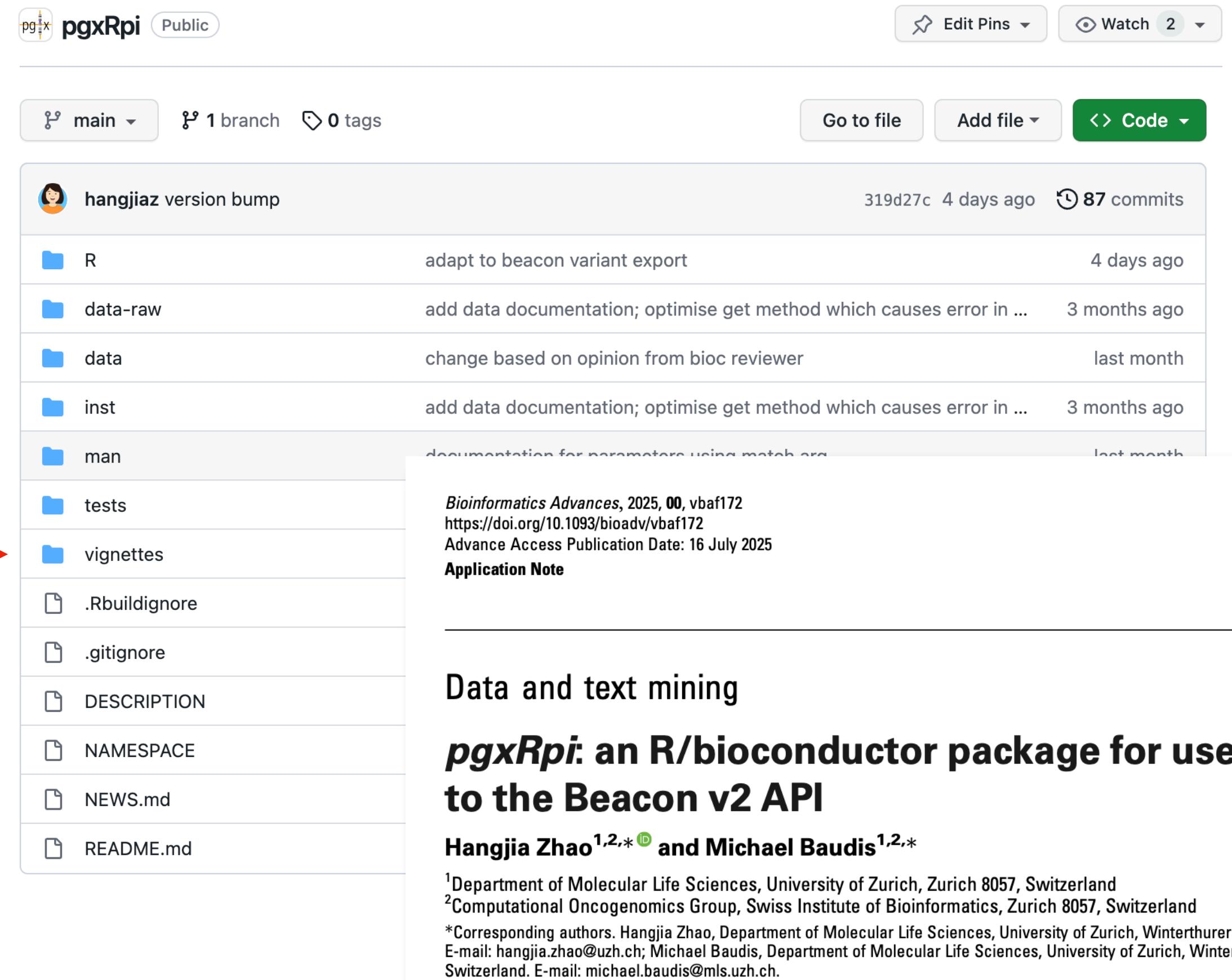
Packages

No packages published Publish your first package

Contributors 6

pgxRpi

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



pgxRpi Public

Edit Pins Watch 2

main 1 branch 0 tags

Go to file Add file Code

hangjiaz version bump 319d27c 4 days ago 87 commits

R adapt to beacon variant export 4 days ago

data-raw add data documentation; optimise get method which causes error in ... 3 months ago

data change based on opinion from bioc reviewer last month

inst add data documentation; optimise get method which causes error in ... 3 months ago

man documentation for parameters using match_... last month

Bioinformatics Advances, 2025, 00, vba172
https://doi.org/10.1093/bioadv/vba172
Advance Access Publication Date: 16 July 2025

Application Note

OXFORD

Data and text mining

pgxRpi: an R/bioconductor package for user-friendly access to the Beacon v2 API

Hangjia Zhao^{1,2,*} and Michael Baudis^{1,2,*}

¹Department of Molecular Life Sciences, University of Zurich, Zurich 8057, Switzerland

²Computational Oncogenomics Group, Swiss Institute of Bioinformatics, Zurich 8057, Switzerland

*Corresponding authors. Hangjia Zhao, Department of Molecular Life Sciences, University of Zurich, Winterthurerstrasse 190, Zurich 8057, Switzerland. E-mail: hangjia.zhao@uzh.ch; Michael Baudis, Department of Molecular Life Sciences, University of Zurich, Winterthurerstrasse 190, Zurich 8057, Switzerland. E-mail: michael.baudis@mls.uzh.ch.

2 Retrieve metadata of samples

2.1 Relevant parameters

type, filters, filterLogic, individual_id, biosample_id, codematches, limit, skip

2.2 Search by filters

Filters are a significant enhancement to the [Beacon](#) query API, providing a mechanism for specifying rules to select records based on their field values. To learn more about how to utilize filters in Progenetix, please refer to the [documentation](#).

The `pgxFilter` function helps access available filters used in Progenetix. Here is the example use:

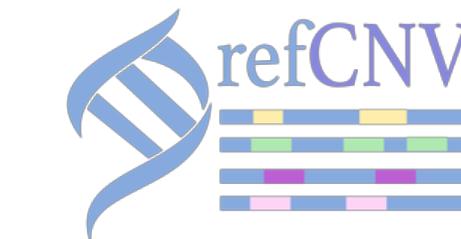
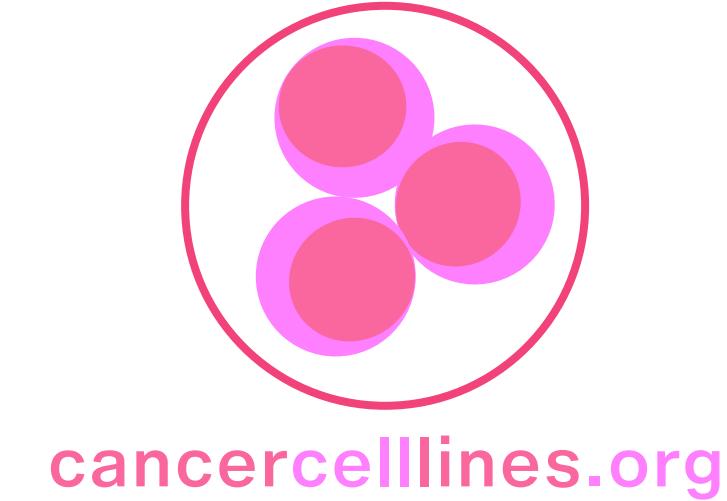
```
# access all filters
all_filters <- pgxFilter()
# get all prefix
all_prefix <- pgxFilter(return_all_prefix = TRUE)
# access specific filters based on prefix
ncit_filters <- pgxFilter(prefix="NCIT")
head(ncit_filters)
#> [1] "NCIT:C28076" "NCIT:C18000" "NCIT:C14158" "NCIT:C14161" "NCIT:C28077"
#> [6] "NCIT:C28078"
```

The following query is designed to retrieve metadata in Progenetix related to all samples of lung adenocarcinoma, utilizing a specific type of filter based on an NCIt code as an ontology identifier.

```
biosamples <- pgxLoader(type="biosample", filters = "NCIT:C3512")
# data looks like this
biosamples[c(1700:1705),]
#>   biosample_id group_id group_label individual_id callset_ids
#> 1700 pgxbs-kftvjjhx NA NA pgxind-kftx5fyd pgxcs-kftwjevi
#> 1701 pgxbs-kftvjjhz NA NA pgxind-kftx5fyf pgxcs-kftwjew0
#> 1702 pgxbs-kftviji1 NA NA pgxind-kftx5fyh pgxcs-kftwjewi
#> 1703 pgxbs-kftvjjn2 NA NA pgxind-kftx5g4r pgxcs-kftwjg5r
#> 1704 pgxbs-kftvjjn4 NA NA pgxind-kftx5g4t pgxcs-kftwjg6q
#> 1705 pgxbs-kftvjjn5 NA NA pgxind-kftx5g4v pgxcs-kftwjg78
```

baudisgroup @ UZH

Project Opportunities



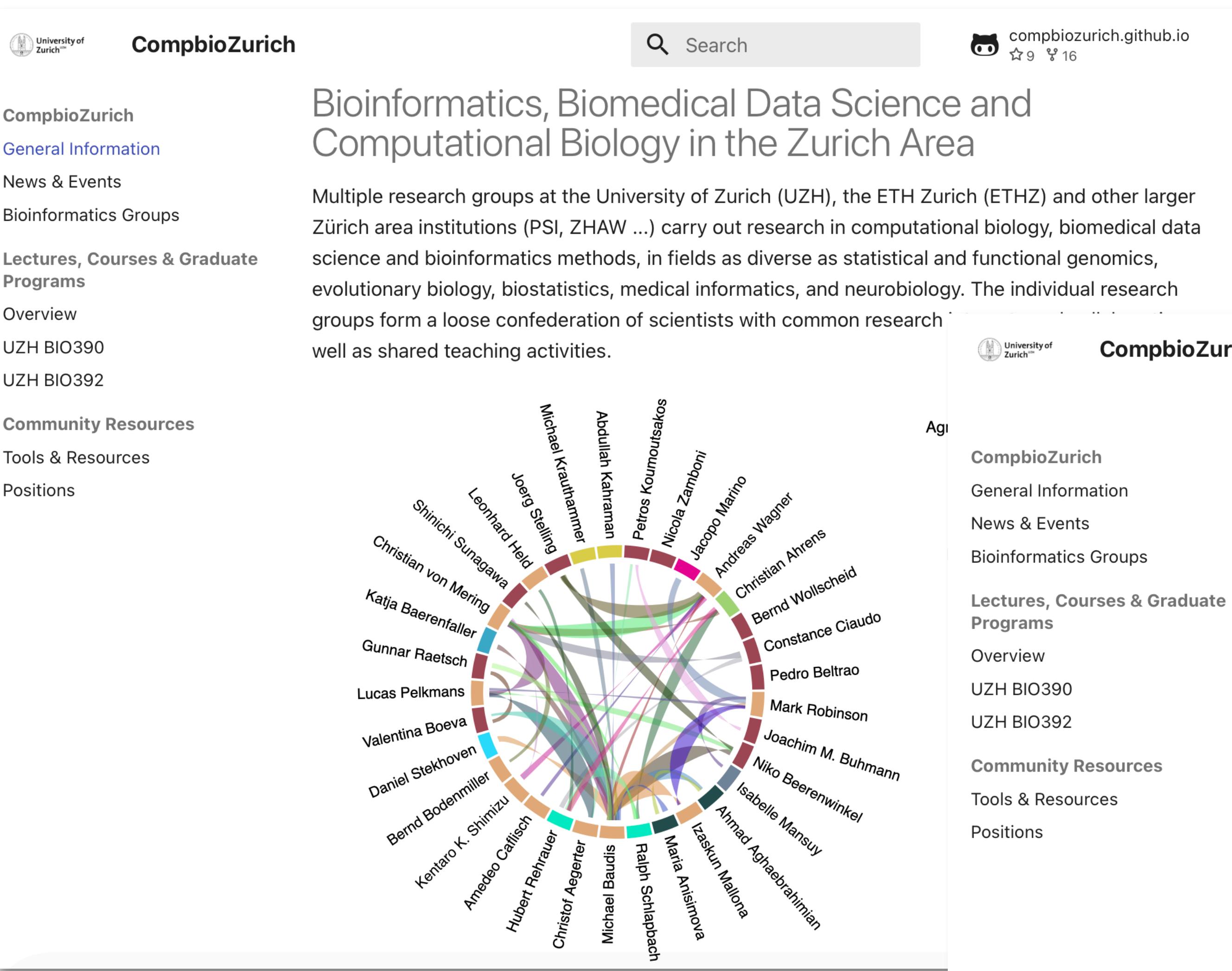
- data & metadata analysis projects
 - ▶ disease relations, trees, trajectories
 - ▶ geographies, disease statistics and clinical explorations
- data curation and resource generation
 - ▶ identification, retrieval and conversion of CNV and metadata from publications and repositories
 - ▶ ontology and classification mapping
 - ▶ clinical and metadata enrichment
 - ▶ beyond CNV mutation data & double hit events

baudisgroup @ UZH

Project Opportunities: Beacon Ecosystem



- schema development and implementation
 - ▶ adapting and contributing to **GA4GH** standards (Beacon but also VRS and core schemas)
- the **bycon** software stack
 - ▶ **Python** & **MongoDB** code development and optimization
 - ▶ packaging and distribution
 - ▶ security & authentication
 - ▶ aggregator...
- front-end, visualization and web deployment
 - ▶ **Java Script / React** development
 - ▶ queries, maps, genome plots ...

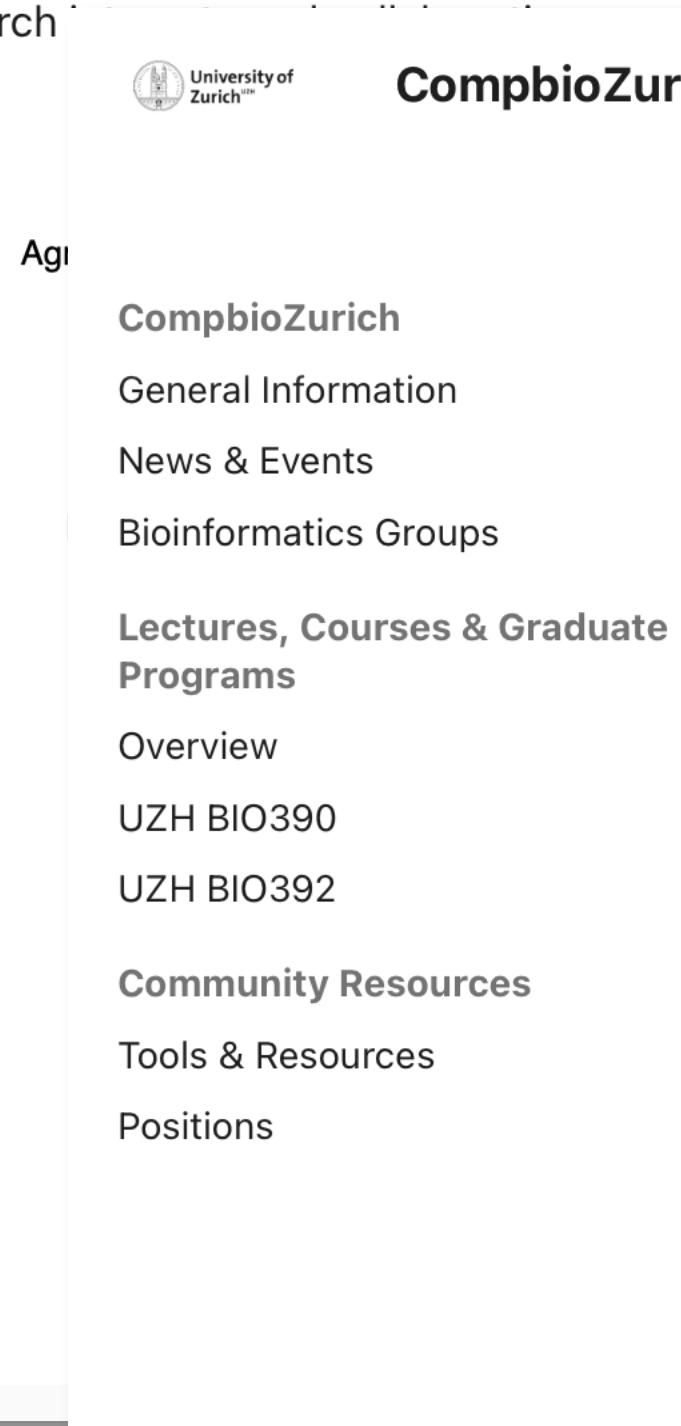


compbiozurich.org

Zürich Seminars in Bioinformatics

Thursdays 12:15-13:00

UZH Irchel Campyus & Zoom



News, Events, Seminars ...

Associating cancer-specific variants with methylation changes using nanopore genome data

Halimat Atanda (PhD student, University of Queensland)

ZURICH SEMINARS IN BIOINFORMATICS

 September 11, 2025

- 10:30 UZH Irchel Y17-H-05

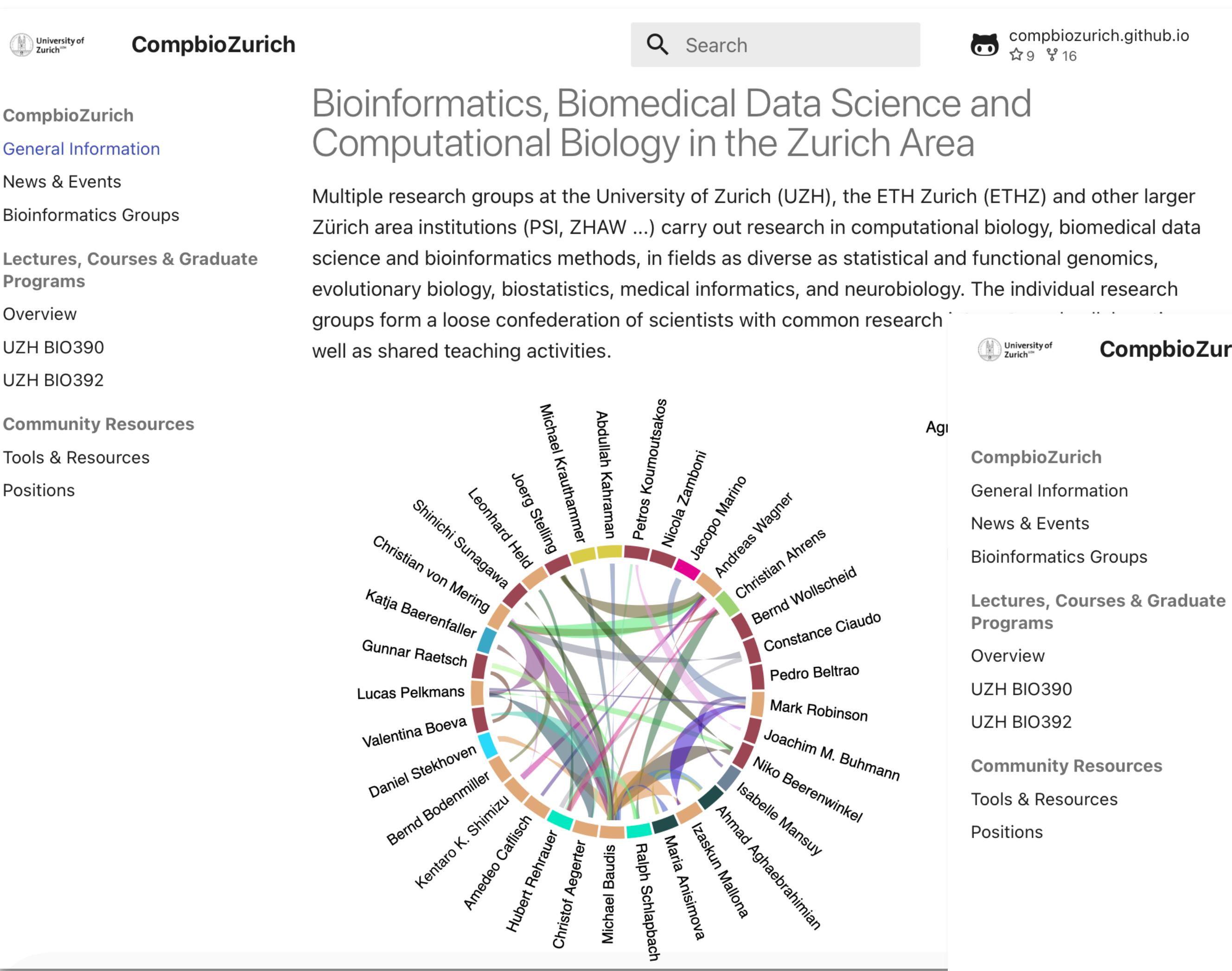
Abstract Advancement in high throughput sequencing technologies has eased the extraction of the whole genome from biological samples and the identification of genomic variants that define molecular phenotypes. This has resulted in a relative abundance of genomic data, especially variants, with limited understanding of their functional or regulatory impacts. → [Continue reading](#)

Benchmarking Cell Segmentation Approaches for High-Resolution Spatial Transcriptomic

Raphael De Gottardi (CBB M.Sc. Student, FGCZ)

ZURICH SEMINARS IN BIOINFORMATICS

  August 28, 2025

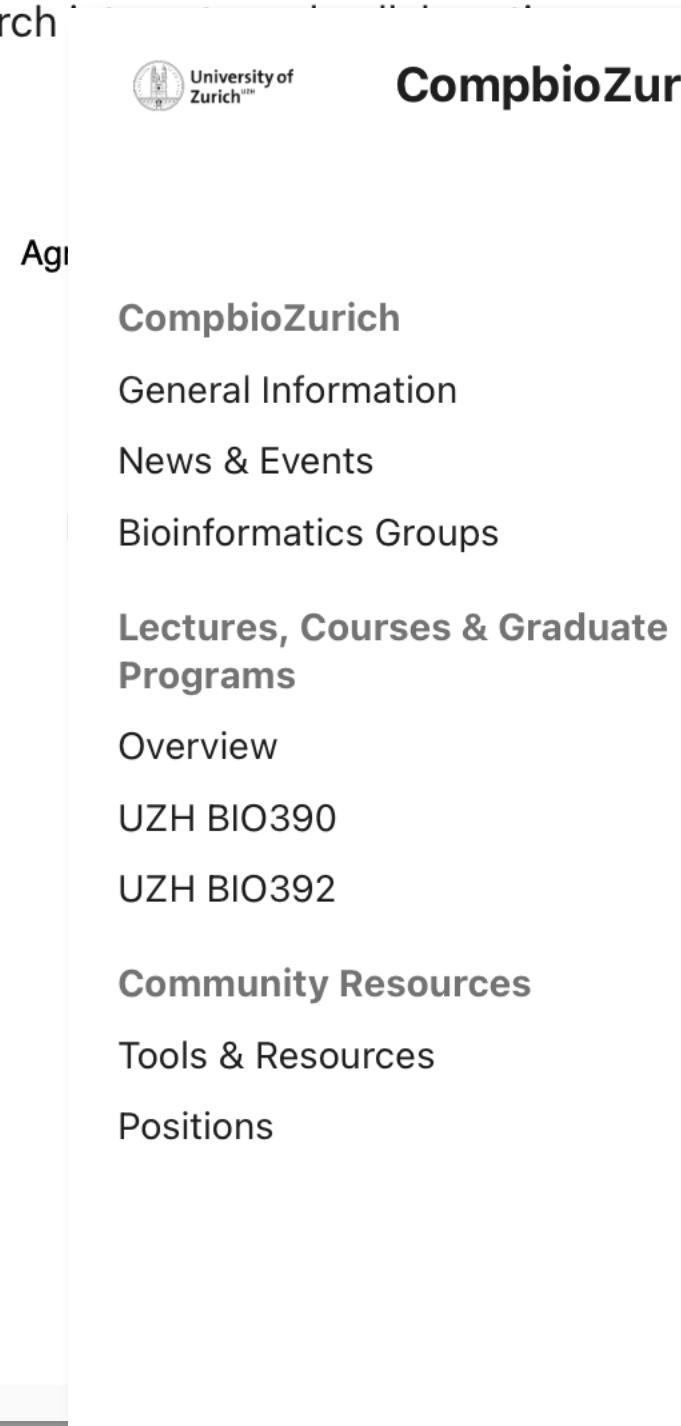


compbiozurich.org

Zürich Seminars in Bioinformatics

Thursdays 12:15-13:00

UZH Irchel Campyus & Zoom



News, Events, Seminars ...

Associating cancer-specific variants with methylation changes using nanopore genome data

Halimat Atanda (PhD student, University of Queensland)

ZURICH SEMINARS IN BIOINFORMATICS

 September 11, 2025

- 10:30 UZH Irchel Y17-H-05

Abstract Advancement in high throughput sequencing technologies has eased the extraction of the whole genome from biological samples and the identification of genomic variants that define molecular phenotypes. This has resulted in a relative abundance of genomic data, especially variants, with limited understanding of their functional or regulatory impacts. → [Continue reading](#)

Benchmarking Cell Segmentation Approaches for High-Resolution Spatial Transcriptomic

Raphael De Gottardi (CBB M.Sc. Student, FGCZ)

ZURICH SEMINARS IN BIOINFORMATICS

  August 28, 2025



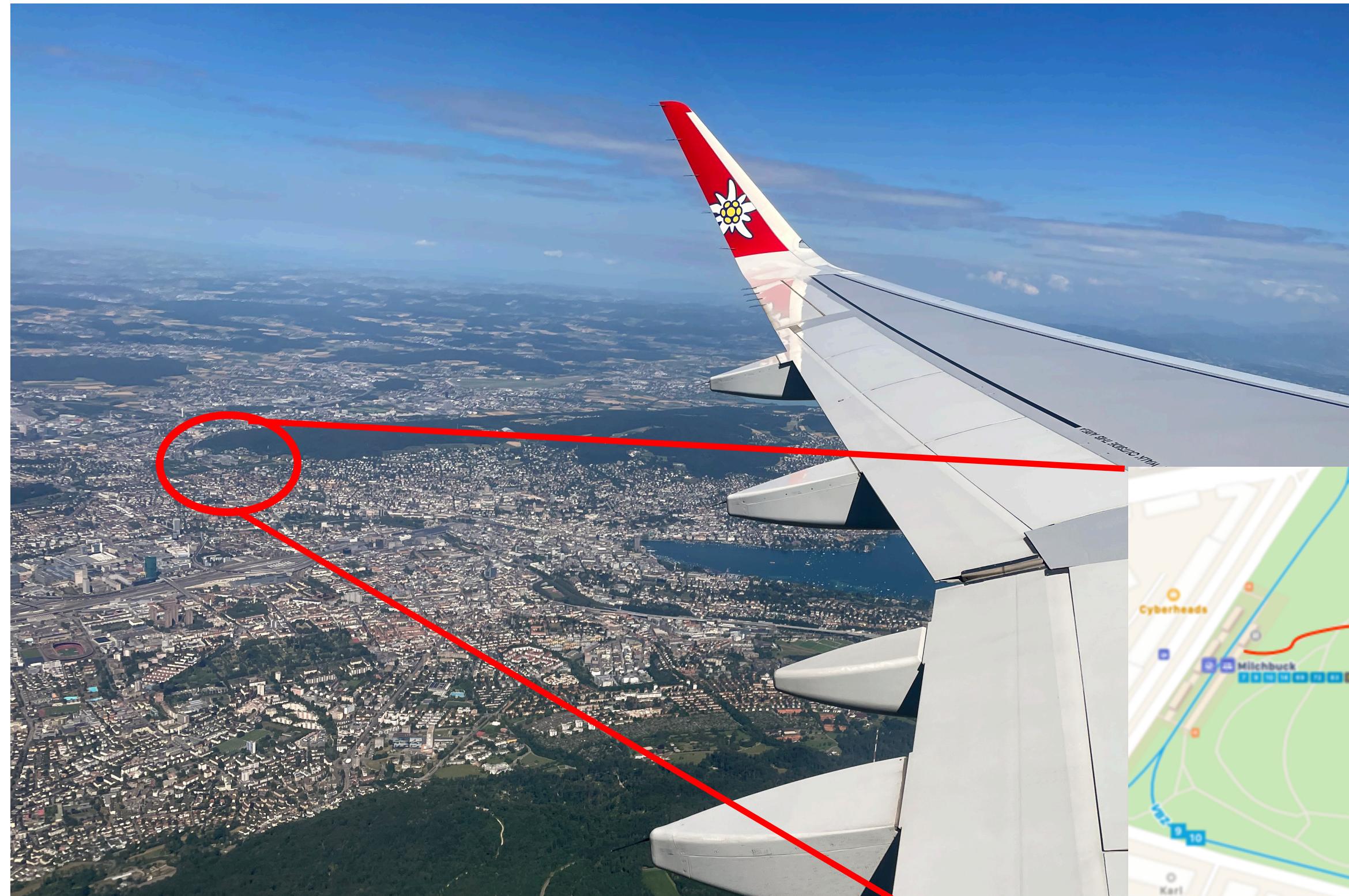
University of
Zurich^{UZH}
Department of Molecular Life Sciences



Global Alliance
for Genomics & Health



Swiss Institute of
Bioinformatics



Prof. Dr. Michael Baudis
Institute of Molecular Life Sciences
University of Zurich
SIB | Swiss Institute of Bioinformatics
Winterthurerstrasse 190
CH-8057 Zurich
Switzerland