



# Cancer Data

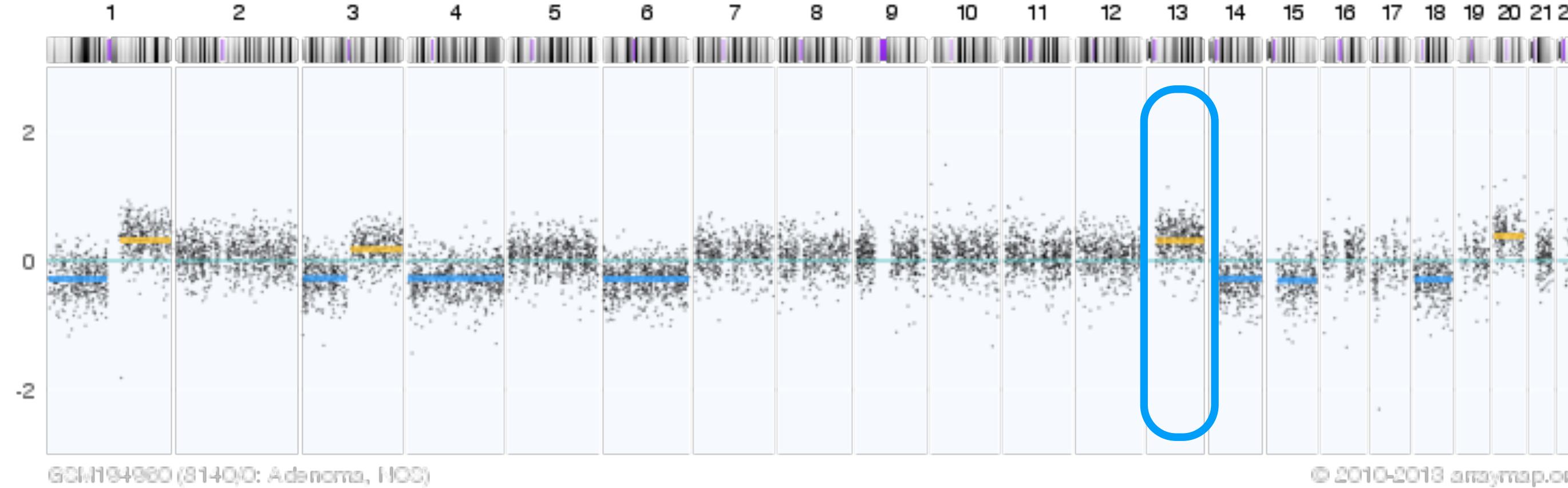
**ELIXIR::GA4GH: Advancing genomics resources through standards and ontologies**

Michael Baudis | ECCB vBarcelona | 2020-09-02

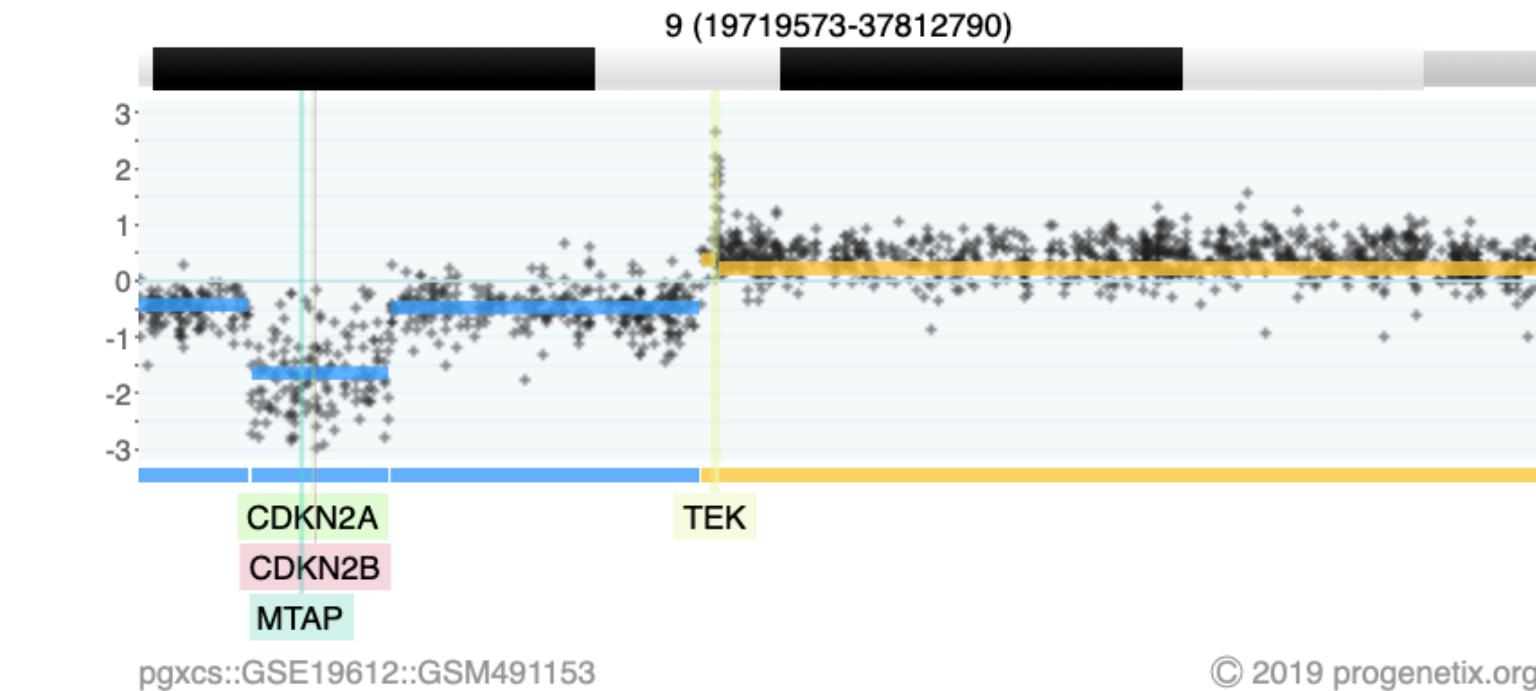


Global Alliance  
for Genomics & Health

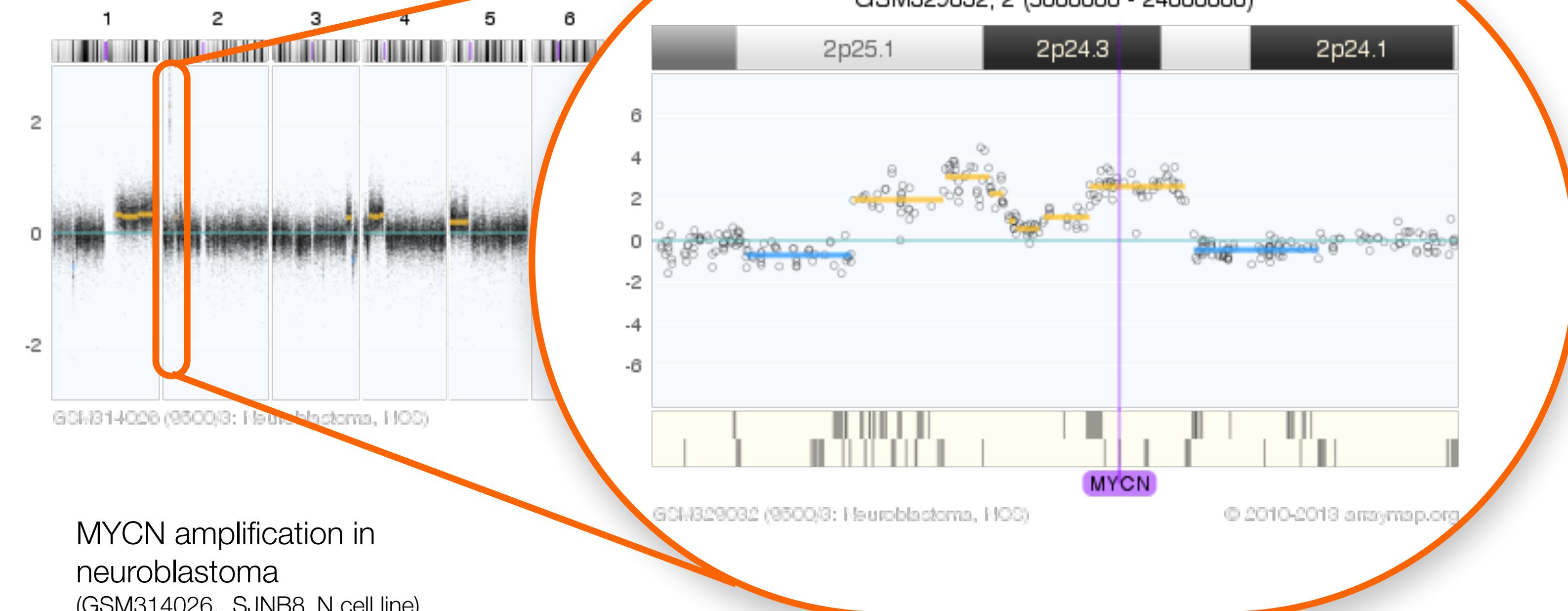
# Somatic Copy Number Variations



Gain of chromosome arm 13q in colorectal carcinoma



2-event, homozygous deletion in a Glioblastoma



MYCN amplification in neuroblastoma  
(GSM314026, SJNB8\_N cell line)

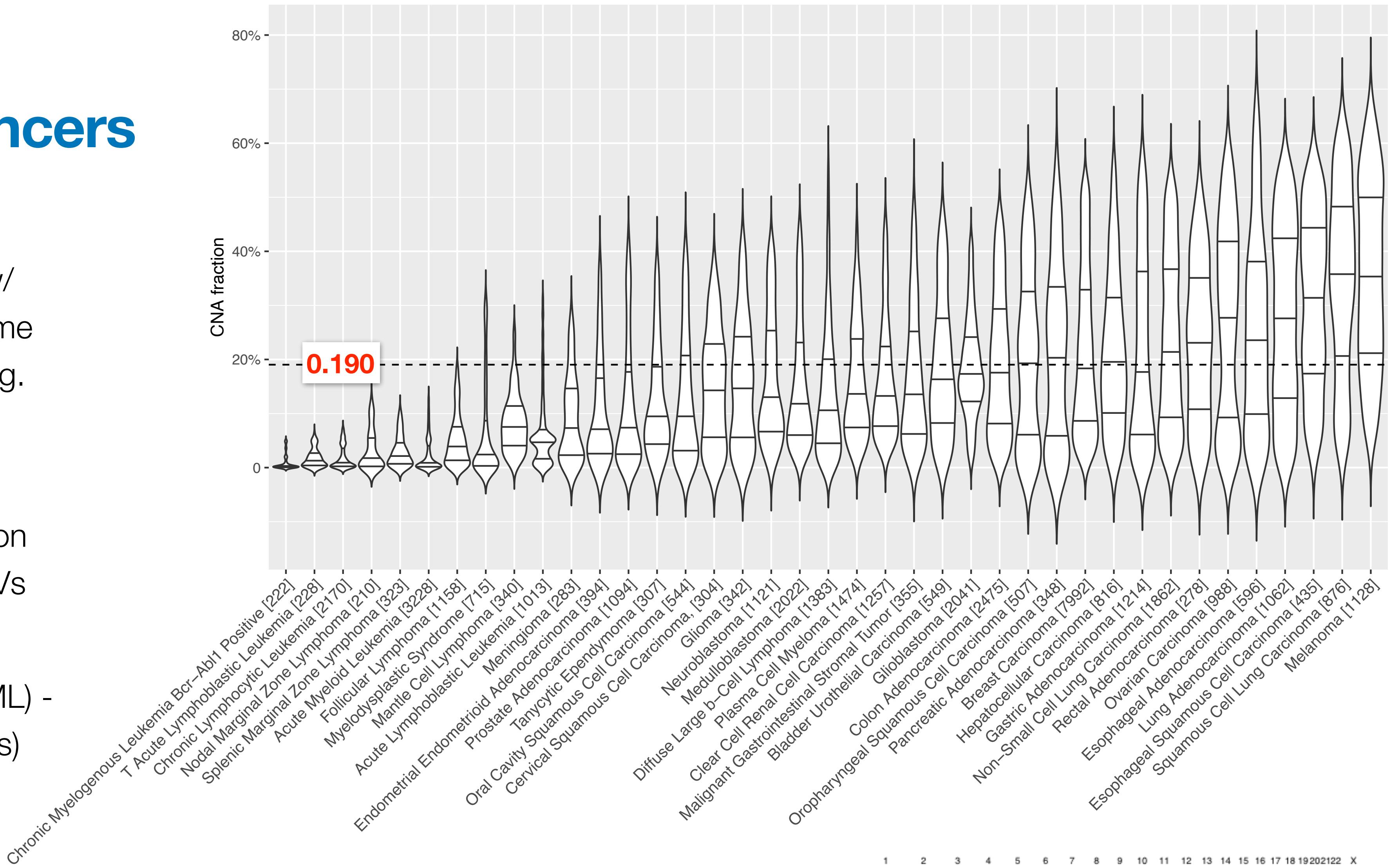
low level/high level copy number alterations (CNAs)

arrayMap

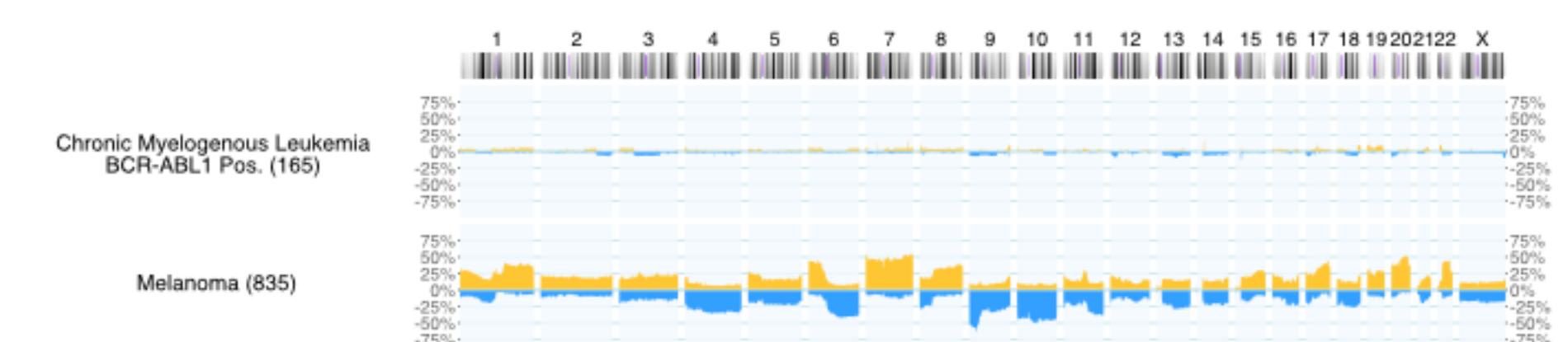


# Genome CNV coverage in Cancers

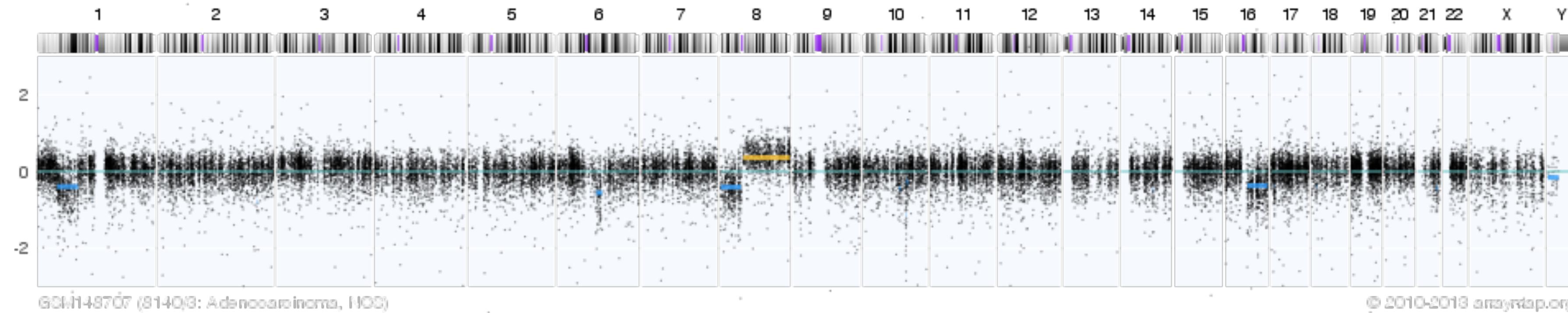
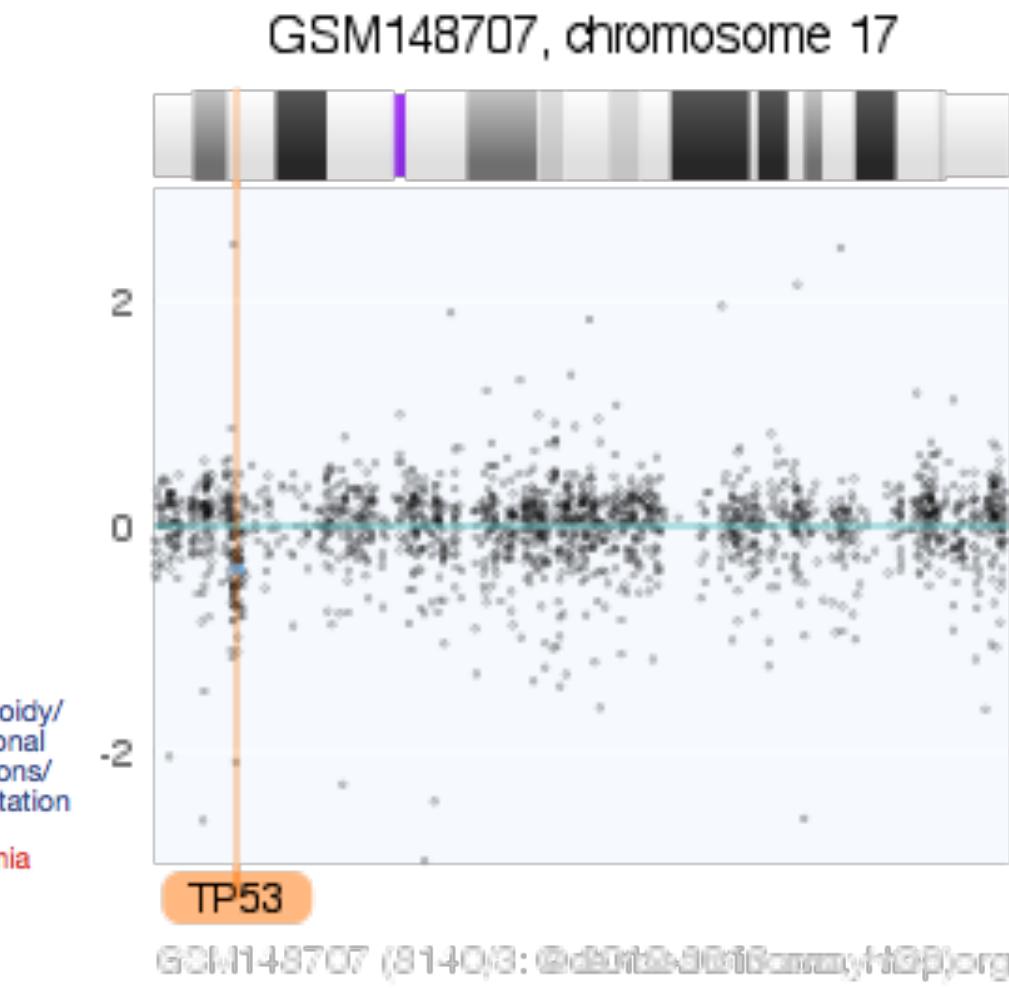
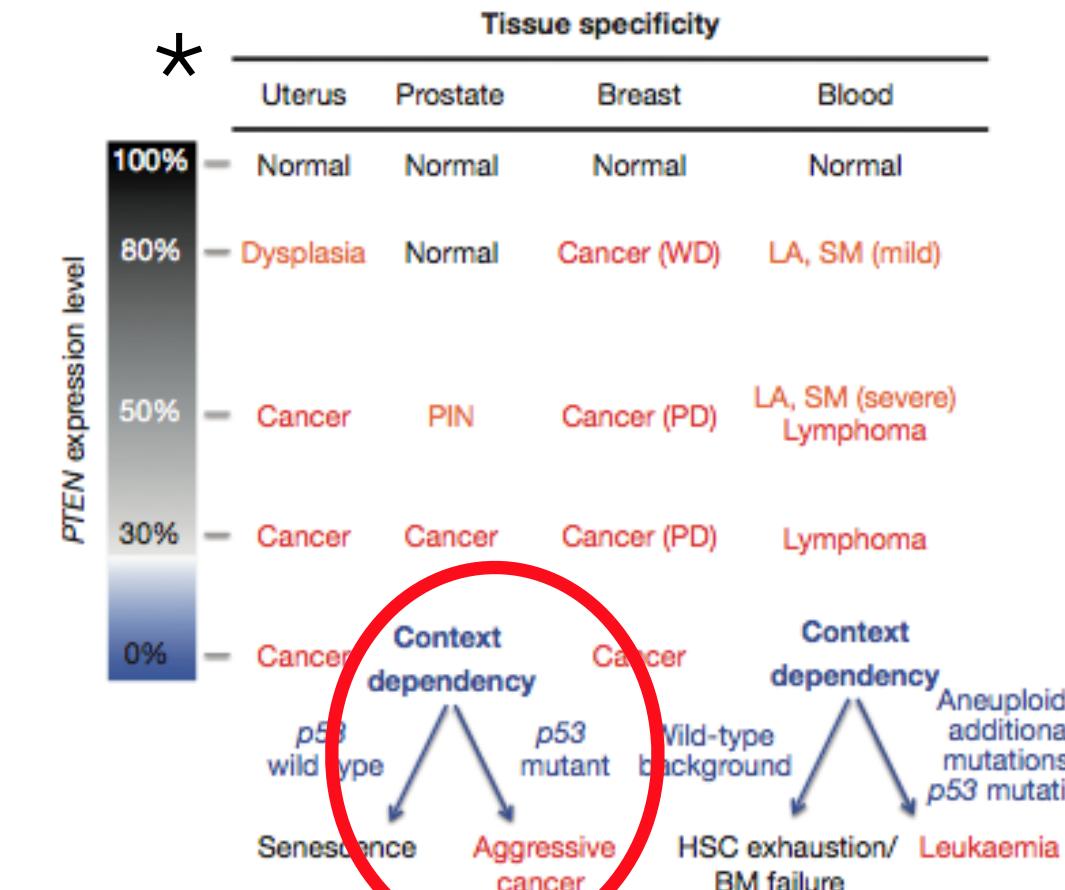
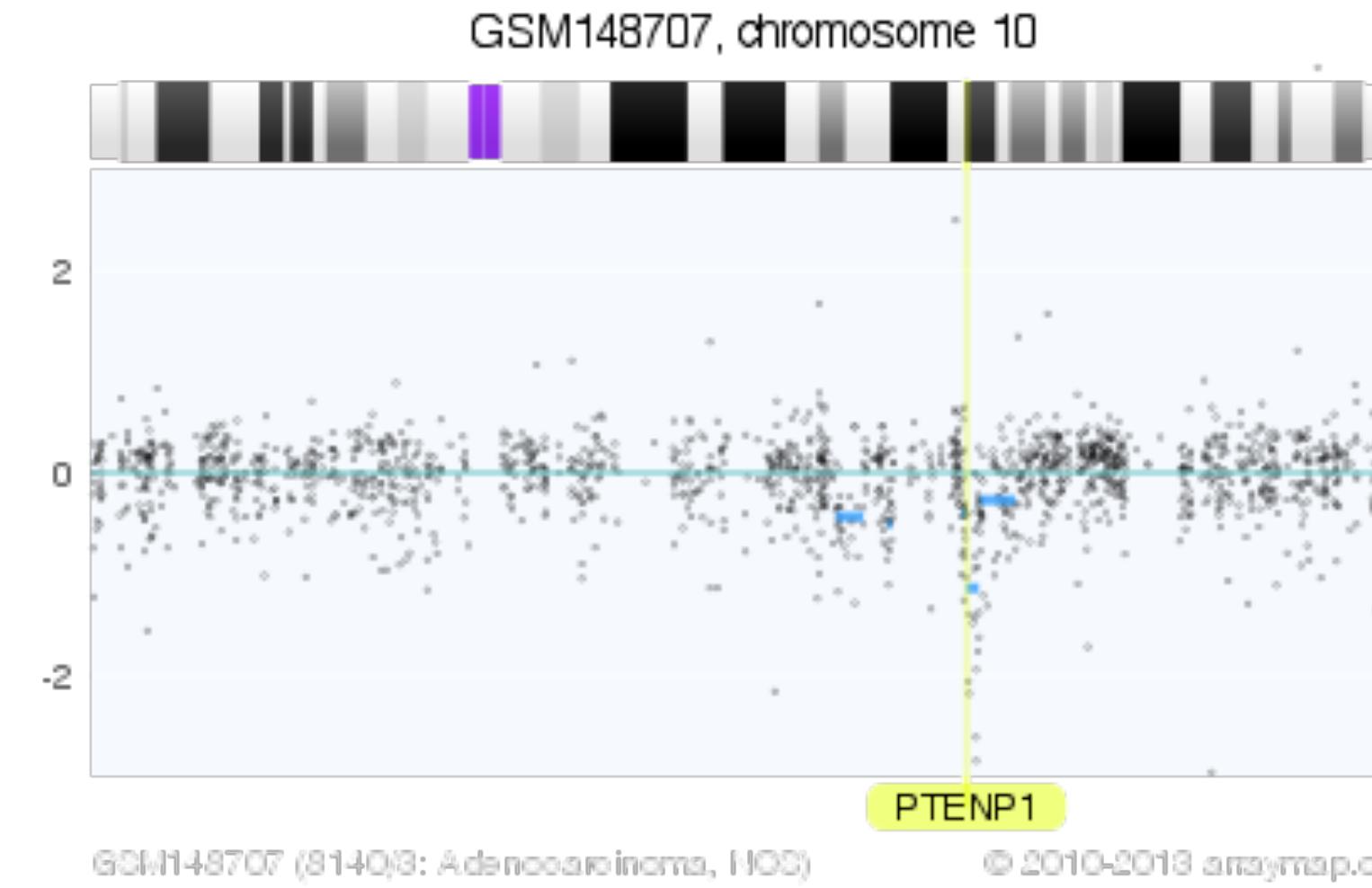
- 43654 out of 93640 CNV profiles; filtered for entities w/ >200 samples (removed some entities w/ high CNV rate, e.g. sarcoma subtypes)
- Single-sample CNV profiles were assessed for the fraction of the genome showing CNVs (relative gains, losses)
- range of medians 0.001 (CML) - 0.358 (malignant melanomas)



Lowest / Highest CNV fractions =>



# Gene dosage phenomena beyond simple on/off effects



Combined heterozygous deletions involving *PTEN* and *TP53* loci in a case of prostate adenocarcinoma  
(GSM148707, PMID 17875689, Lapointe et al., CancRes 2007)

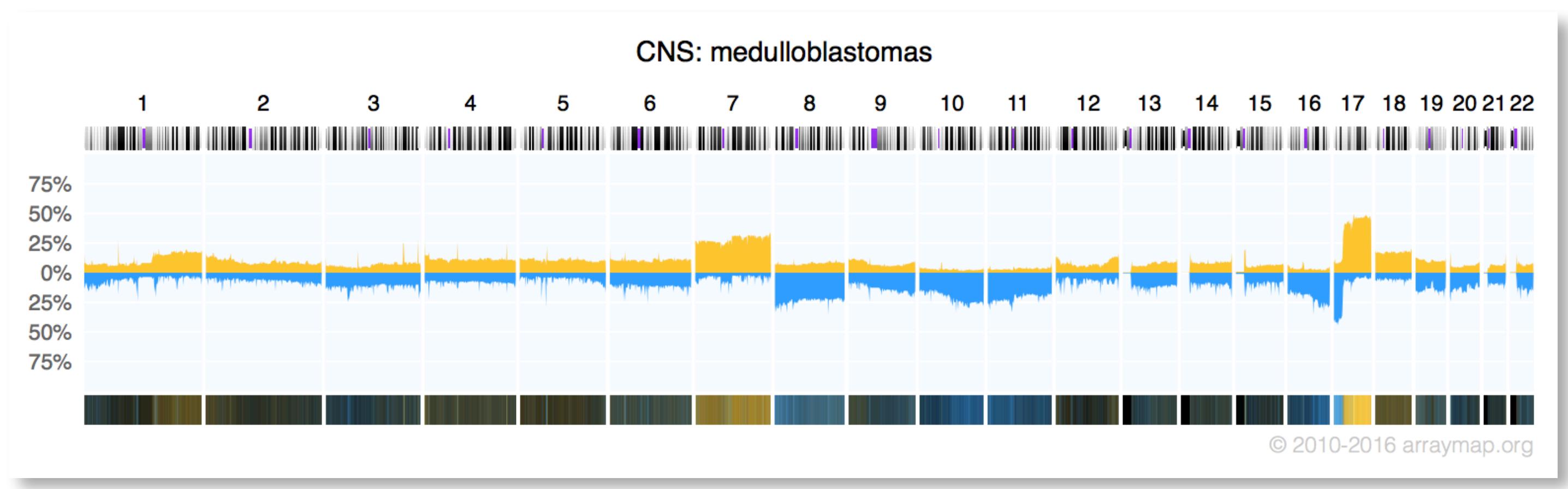
\* A. H. Berger, A. G. Knudson, and P. P. Pandolfi, "A continuum model for tumour suppression," *Nature*, vol. 476, no. 7359, pp. 163–169, Aug. 2011.



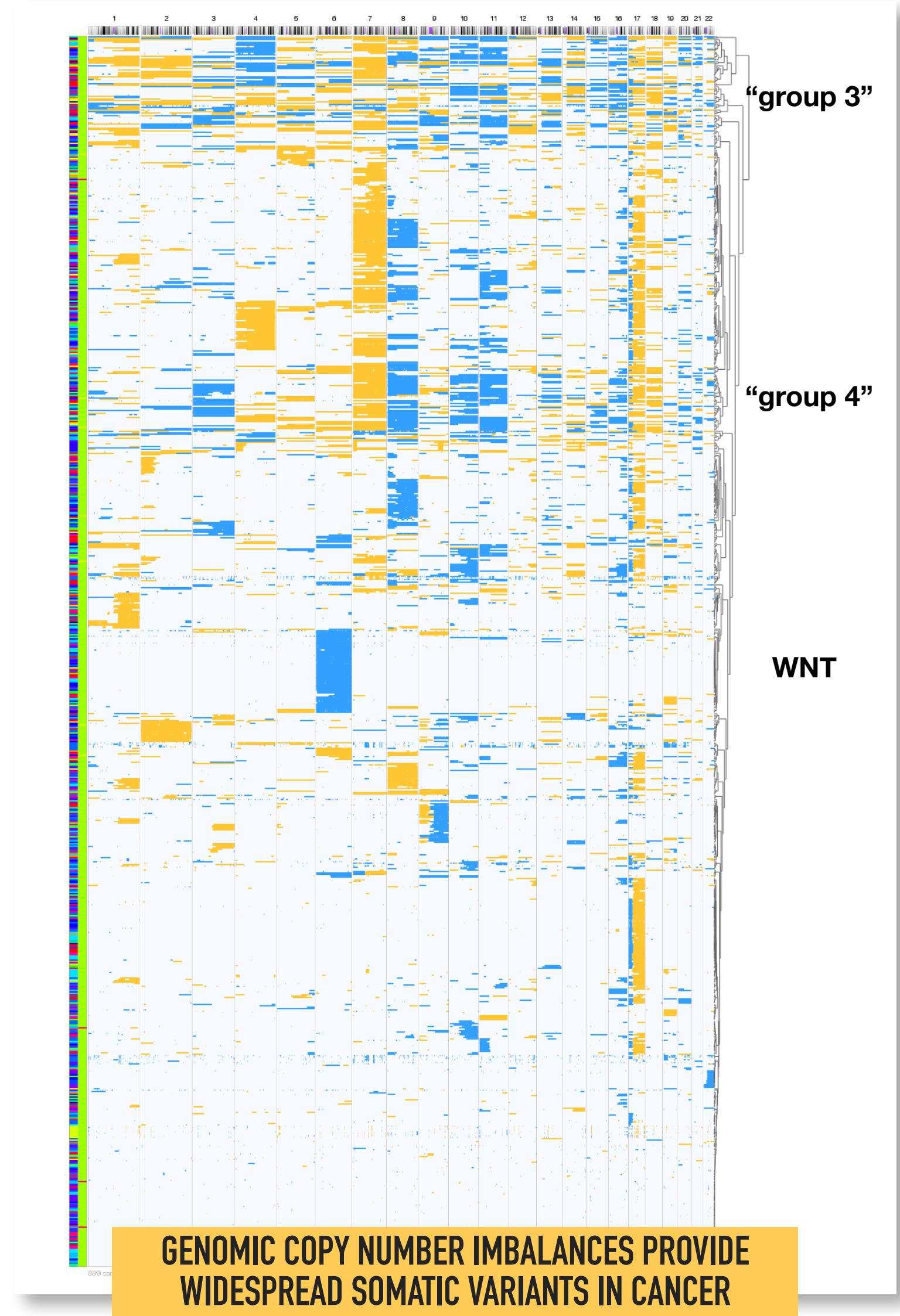
# Somatic CNVs In Cancer

## Recurrent mutation patterns

How can those patterns be used for classification and determination of biological mechanisms?



A genomic copy number histogram for malignant medulloblastomas, the most frequent type of pediatric brain tumors, displaying regions of genomic duplications and deletions. These can be decomposed into individual tumor profiles which segregate into several clusters of related mutation patterns with functional relevance and clinical correlation.

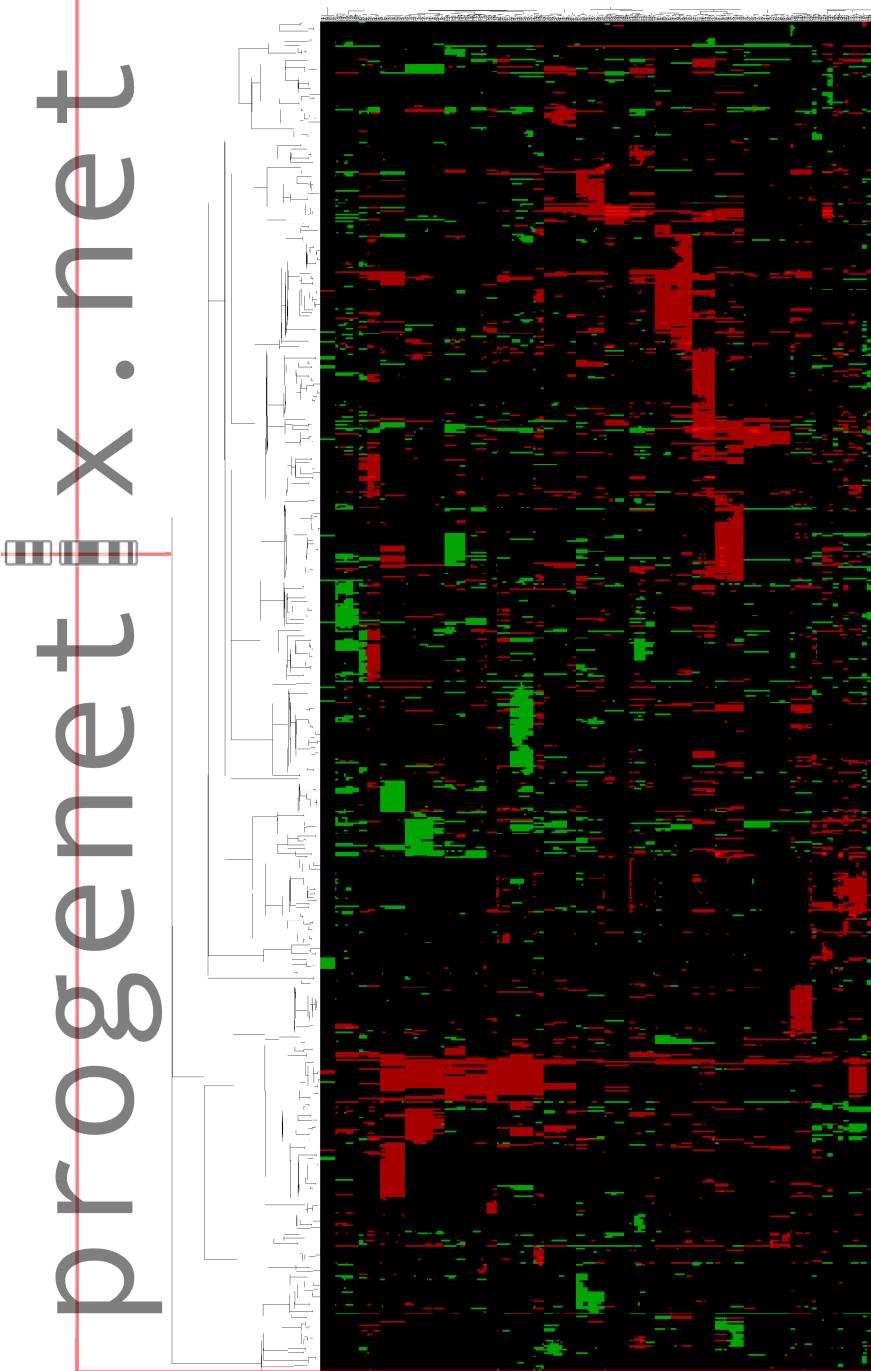


Over the last decade, techniques for the genome wide scanning for genomic imbalances in malignant neoplasia have been developed, e.g. Comparative Genomic Hybridization (CGH).

Currently, no comprehensive online source for CGH data with a standardized format suitable for data mining procedures has been made available for public access. Such a data repository could be valuable in identifying genetic aberration patterns with linkage to specific disease entities, and provide additional information for validating data from large scale expression array experiments.

A case and band specific aberration matrix was selected as most suitable format for the mining of CGH data. The [progenetix.net] data repository was developed to provide the according data to the research community for a growing number of human malignancies.

In the current implementation, two main purposes are being served. First, access to the band specific pattern of chromosomal imbalances allows the instantaneous identification of genomic "hotspots". Second, the band specific aberration matrices can be included in data mining efforts. As an example, the clustering off all informative cases from the current (September 2001) dataset is shown here (online source under [www.progenetix.net/bcats/clustered.png](http://www.progenetix.net/bcats/clustered.png)).



#### Data selection

PubMed is searched for publications applying CGH to the analysis of malignant tumors. Articles are selected according to their online availability and the description of genomic imbalances on a per case basis.

#### Transformation of input data

Chromosomal aberration data is transformed via customized parsing commands to a common format adherent to ISCN 1995 recommendations. In some cases, aberration data was transcribed from graphical representations or provided by the authors.

#### Data storage

Currently, the primary data is stored in a dedicated "off-line" database. Besides case identifier and ISCN adapted chromosomal imbalance data, tumor classification and source information including the PubMed identifier is recorded. Disease entities are reclassified to ICD-O-3 codes.

#### Text parsing and generation of aberration matrix

For the generation of the case and band specific aberration matrix, a dedicated text pattern comparison model was developed using Perl. Briefly, for each chromosomal band, the aberration field of each case is searched for a variety of patterns containing aberration information applying to that band. A matrix with currently 324 band resolution is generated, annotating chromosomal gains with "+" and losses with "-"; localized high-level gains are designated "2".

#### Website generation

For graphical representation of chromosomal imbalances, HTML pages containing different views of the underlying aberration matrices are generated using Perl. Graphics are implemented using HTML syntax. Besides band specific, whole genomic overviews, chromosome specific pages with links to all involved cases are generated for each ICD-O-3 entity as well as for each registered project. Additionally, those representations are available for several subsets combining related data (e.g. all lymphoid neoplasias, breast carcinoma cases). For each of the groups, the according aberration matrix is linked for download.

Hierarchical clustering of band specific chromosomal imbalances from 999 human neoplasias, contained in the [progenetix.net] collection. Cases without aberrations were excluded.



## Progenetix.net: an online repository for molecular cytogenetic aberration data

Michael Baudis <sup>1, 2,\*</sup> and Michael L. Cleary<sup>2</sup>

<sup>1</sup>Medizinische Klinik und Poliklinik V der Universität Heidelberg, Germany

<sup>2</sup>Department of Pathology, Stanford University Medical Center, Stanford, CA 94305, USA

Received on July 5, 2001; revised on July 9, 2001; accepted on July 16, 2001

#### ABSTRACT

**Summary:** Through sequencing projects and, more recently, array-based expression analysis experiments, a wealth of genetic data has become accessible via online resources. In contrast, few of the (molecular-) cytogenetic aberration data collected in the last decades are available in a format suitable for data mining procedures. [www.progenetix.net](http://www.progenetix.net) is a new online repository for previously published chromosomal aberration data, allowing the addition of band-specific information about chromosomal imbalances to oncologic data analysis efforts.

**Availability:** <http://www.progenetix.net>  
**Contact:** mbaudis@stanford.edu

Neoplastic transformation and progression is the result of genetic defects arising in normal cells and giving rise to a malignant clone. During the process of oncogenesis, some of the usually multiple steps required for acquisition of the full neoplastic phenotype may represent themselves as numerical or structural abnormalities in the chromosomes of the transformed cells.

Over the last decades, the analysis of chromosomal abnormalities in malignant cells has gained importance in oncologic research as well as in clinical practice. A vast number of genetic abnormalities has been identified in the virtually complete range of human neoplasias. Several attempts have been undertaken for collection and classification of those abnormalities, the most widely recognized being the catalog by Mitelman and co-workers (Mitelman, 1994; online access through <http://cgap.nci.nih.gov/Chromosomes/Mitelman>).

In addition to metaphase analysis of short-term cultivated tumor cells or tumor cell lines, molecular cytogenetic techniques have recently been applied to the analysis of chromosomal abnormalities in primary tumor tissues. One of the more widely used screening techniques is Comparative Genomic Hybridization (CGH; Kallion-

iem et al., 1992; du Manoir et al., 1993). Briefly, this method is based on the competitive *in-situ* hybridization of differentially labeled tumor versus normal genomic DNA to normal human metaphase spreads. The calculation of the intensity ratios of the two fluorochromes gives an overview about relative gains and losses of DNA in the tumor genome with mapping to the respective chromosomal bands. The identification of frequently imbalanced regions in tumor entities may point towards tumor suppressor gene or proto-oncogenes mapping to the respective chromosomal bands. Usually, the result of those experiments is communicated either in text format according to the International System for Cytogenetic Nomenclature (Mitelman, 1995) or graphically, with aberration bars next to chromosomal ideograms for the representation of chromosomal gains and losses.

Because in each experiment CGH analysis covers the whole number of chromosomes, the comparison of data sets from related malignancies could lead to the delineation of common as well as divergent genetic pathways defining the respective malignant phenotypes. Although an extremely large number of malignant tumors has been analyzed using this technique, no comprehensive CGH database with band-specific chromosomal aberration information is publicly available<sup>†</sup>.

A minimal requirement for such a database would be the conversion of the text or graphical information used in publications to data tables, representing the information about the aberration status of single chromosomal bands for each case. For the site discussed here, this process includes: (1) the transformation of the published results in a format adapted from the ISCN, and (2) the automatic generation of the band specific aberration table.

Due to format variations of the published data, step 1 consists of the manual conversion of the text data or evaluation and conversion of the graphical representations, respectively. Due to the (in computational terms) odd

<sup>†</sup>Links to a number of online CGH resources with different scopes can be found at [www.progenetix.net](http://www.progenetix.net).

\*To whom correspondence should be addressed.



[Cancer CNV Profiles](#)

[Search Samples](#)

[Publication DB](#)

[Info](#)

[Beacon<sup>+</sup>](#)

[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.

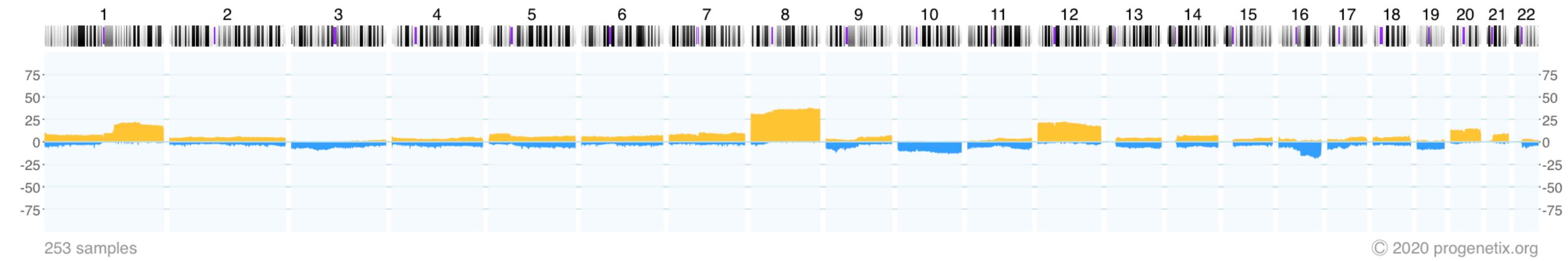
### Data Content and Provenance

The resource currently contains genome profiles of **113322** individual experiments. The genomic profiling data was derived from array and chromosomal [Comparative Genomic Hybridization \(CGH\)](#) experiments as well as Whole Genome or Whole Exome Sequencing (WGS, WES) studies. Genomic profiles are either processed from various raw data formats or are extracted from published experimental results.

Besides genomic profiling data, Progenetix contains sample specific biological, technical and provenance information which so far has been curated from **1600** articles.

Original diagnoses are mapped to (hierarchical) classification systems and represents **420** and **542** different cancer types, according to the International classification of Diseases in Oncology (ICD-O) and NCIt "neoplasm" classification, respectively.

Progenetix: Ewing sarcoma (icdom-92603)



For exploration of the resource it is suggested to either start with [Cancer Types](#) or by [searching](#) for CNVs genes of interest.

Additionally to genome profiles and associated metadata, the website present information about publications (currently **3962** articles) referring to cancer genome profiling experiments.

### Access, Maintenance and Contributions

The content of the progenetix resource is freely accessible for research and commercial purposes, with attribution.

The database & software are developed by the [group of Michael Baudis](#) at the [University of Zurich](#) and the Swiss Institute of Bioinformatics [SIB](#).

Many previous members and external collaborators have contributed to data content and resource features. Participation (features, data, comments) by volunteers are welcome.

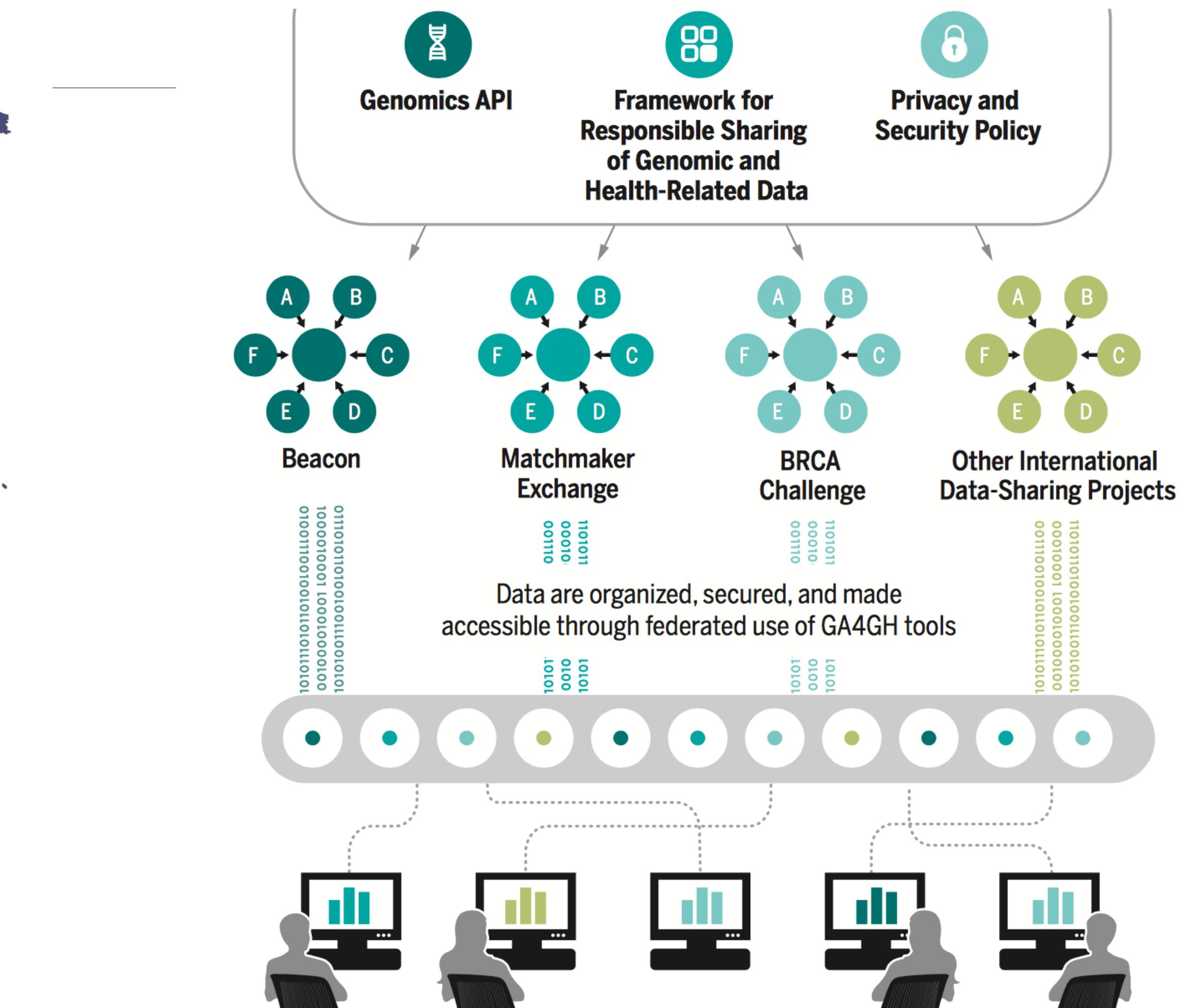


## GENOMICS

# *A federated ecosystem for sharing genomic, clinical data*

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

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



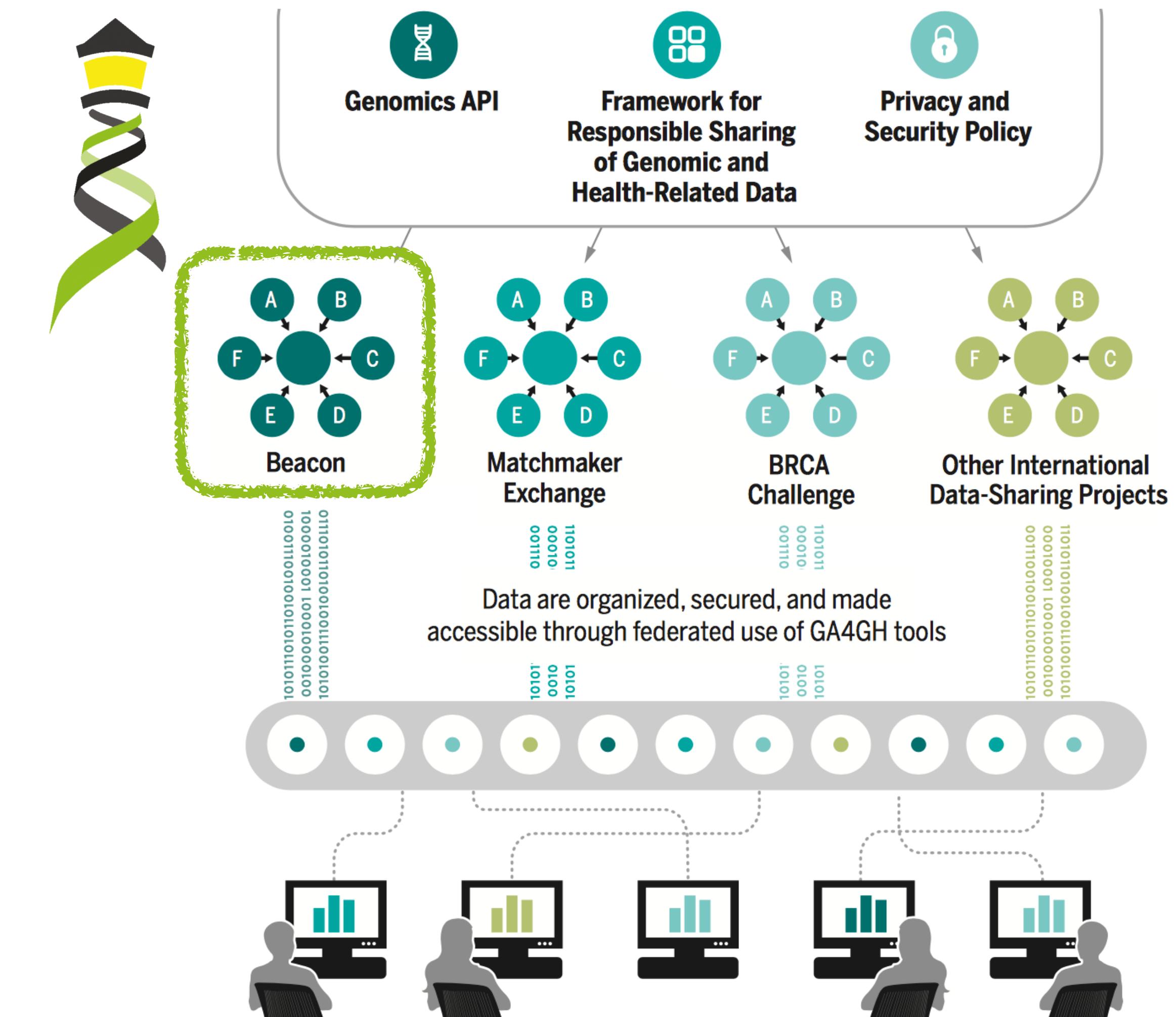


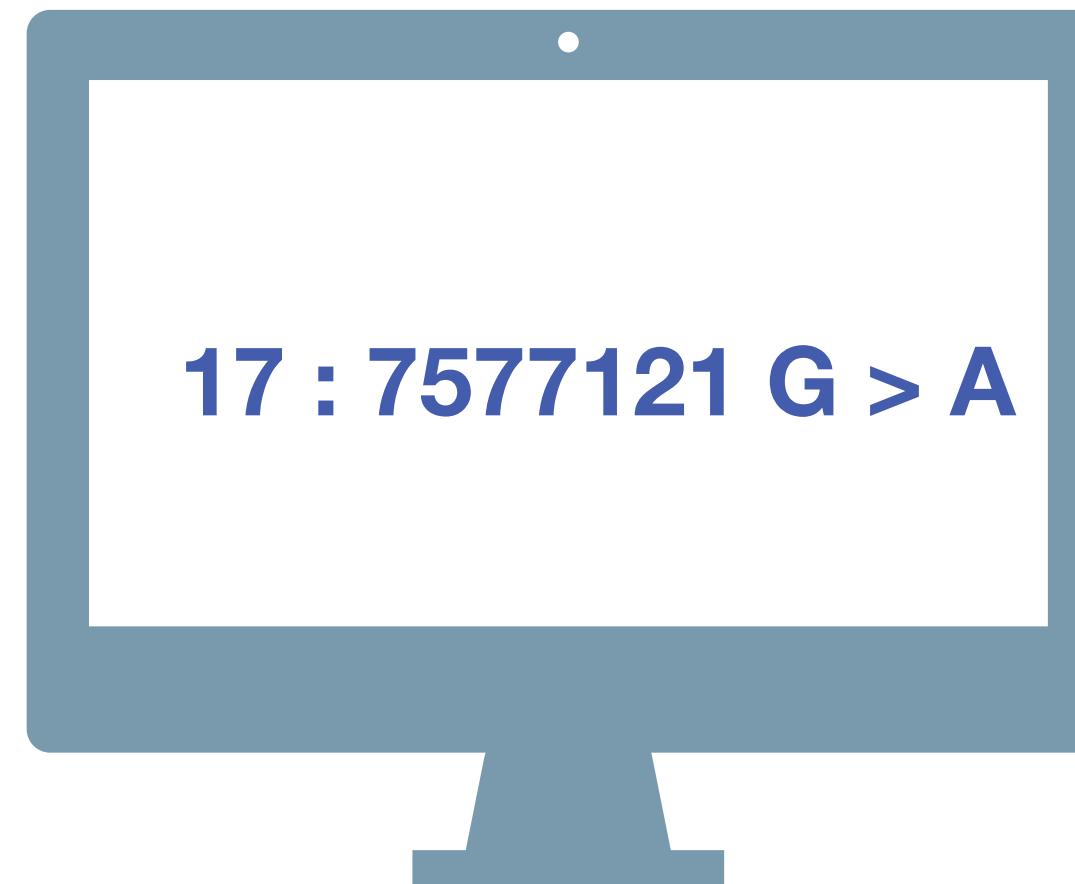
## GENOMICS

# *A federated ecosystem for sharing genomic, clinical data*

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

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





# Beacon

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

**YES | NO | \0**

# ELIXIR Beacon Project

- Driver project on GA4GH roadmap
- aligns with Discovery Work Stream
- strong impact on GA4GH developments as a concrete, funded project

The screenshot shows two cards. The left card is titled 'Driver Projects' and contains text about real-world genomic data initiatives. The right card is titled 'ELIXIR Beacon' and provides links to its implementation studies and champions.

**Driver Projects**  
GA4GH Driver Projects are real-world genomic data initiatives that help guide our development efforts and pilot our tools. Stakeholders around the globe advocate, mandate, implement, and use our frameworks and standards in their local contexts.

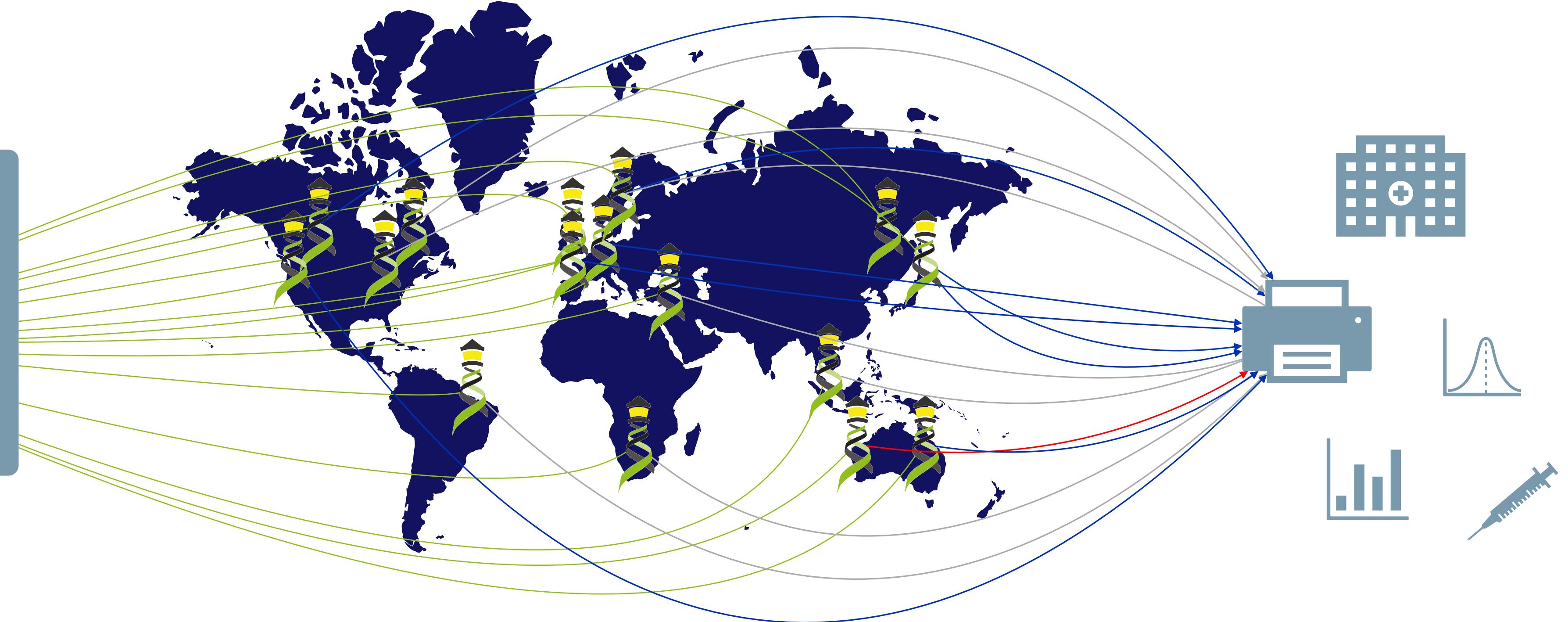
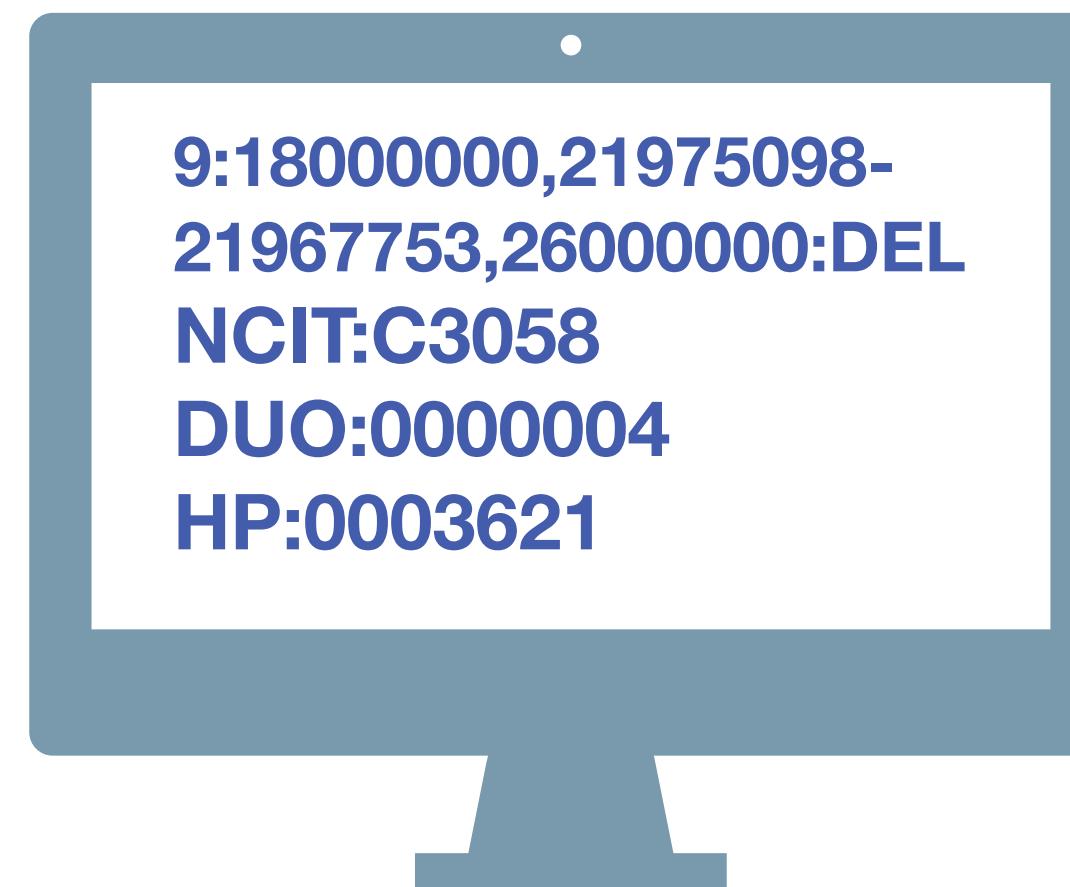
**ELIXIR Beacon**  
<https://www.elixir-europe.org/about/implementation-studies/beacons>

Europe  
**Champions:** Jordi Rambla, Juha Tornroos, Gary Saunders

## v1.1 and roadmap

- structural variations** (DUP, DEL) in addition to SNV
- ... more structural queries (translocations/fusions...)
- Beacon queries as entry for **data handover** (outside Beacon protocol)
- layered authentication system using **ELIXIR AAI**
- v2** **filters** for phenotypic & technical metadata
- v2** Extended quantitative responses
- Ubiquitous **deployment** (e.g. throughout ELIXIR network)





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".



# Features and Possibilities of the current Beacon Specification

Beyond "testing the willingness for data sharing"...

- precise variant queries (chr17: 7673767 C>T)
- range queries ("any variant from here to there")
- variant frequencies
- structural genome variants, e.g. CNVs ("any deletion overlapping CDKN2A CDR coordinates")
- delivery of any kind of data matching a given query (variants, sample information, patient data ...) utilising "**handover**" objects (anonymous links to external services with their own security / privacy implementations)
- networking of v1.n Beacons with AAI integration as demonstrated by the ELIXIR Beacon Network

# DX Ontologies

## Hierarchical NCIt Neoplasm Core replaces heterogeneous primary annotations

- heterogeneous and inconsistent diagnostic annotations are common in clinical reports and research studies ("text", ICD-10, ICD-O 3, OncoTree, domain-specific classifications)
- highly variable granularity of annotations is a major road block for comparative analyses and large scale data integration
  - ▶ "Colorectal Cancer" or "Rectal Mucinous Adenoca."
- initiatives and services such as Phenopackets, MONDO, OXO ... rely on and/or provide mappings to hierarchical ontologies



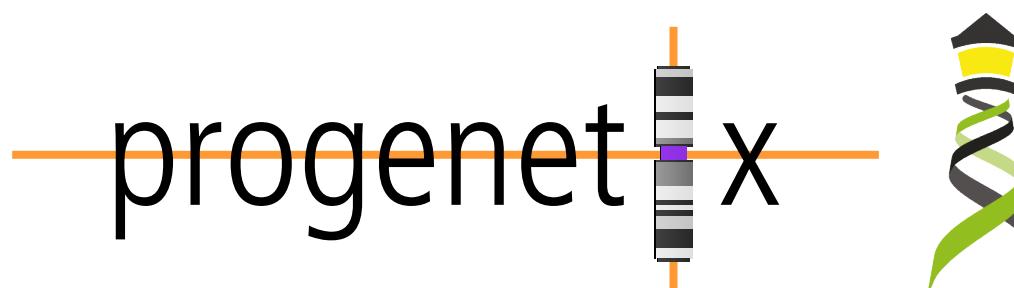
NCIt Neoplasm Core coded display (excerpt) for samples in the Progenetix cancer genome data resource allows sample selection on multiple hierarchy levels →

	Subsets	Samples
□	▼ NCIT:C3262: Neoplasm	88844
□	▼ NCIT:C3263: Neoplasm by Site	84747
□	▼ NCIT:C156482: Genitourinary System Neoplasm	11616
□	▼ NCIT:C156483: Benign Genitourinary System Neoplasm	219
□	▼ NCIT:C4893: Benign Urinary System Neoplasm	90
□	▼ NCIT:C4778: Benign Kidney Neoplasm	90
□	NCIT:C159209: Kidney Leiomyoma	1
□	NCIT:C4526: Kidney Oncocytoma	82
□	NCIT:C8383: Kidney Adenoma	7
□	▼ NCIT:C7617: Benign Reproductive System Neoplasm	129
□	▼ NCIT:C4934: Benign Female Reproductive System Neoplasm	129
□	▼ NCIT:C2895: Benign Ovarian Neoplasm	58
□	▼ NCIT:C4510: Benign Ovarian Epithelial Tumor	58
□	▼ NCIT:C40039: Benign Ovarian Mucinous Tumor	58
□	NCIT:C4512: Ovarian Mucinous Cystadenoma	58
□	▼ NCIT:C4060: Ovarian Cystadenoma	58
□	NCIT:C4512: Ovarian Mucinous Cystadenoma	58
□	▼ NCIT:C3609: Benign Uterine Neoplasm	71
□	▼ NCIT:C3608: Benign Uterine Corpus Neoplasm	71
□	NCIT:C3434: Uterine Corpus Leiomyoma	71
□	▼ NCIT:C156484: Malignant Genitourinary System Neoplasm	11171
□	▼ NCIT:C157774: Metastatic Malignant Genitourinary System Neoplasm	2
□	▼ NCIT:C146893: Metastatic Genitourinary System Carcinoma	2
□	NCIT:C8946: Metastatic Prostate Carcinoma	2
□	▼ NCIT:C164141: Genitourinary System Carcinoma	10561
□	▼ NCIT:C146893: Metastatic Genitourinary System Carcinoma	2
□	NCIT:C8946: Metastatic Prostate Carcinoma	2
□	▼ NCIT:C3867: Fallopian Tube Carcinoma	19

# NCIt Neoplasm Core

# Beacon v2 filters make use of hierarchical classification systems

- Beacon v2 "filters" assumes inclusion of child terms when using hierarchical classifications
    - ➡ implicit *OR* with otherwise assumed *AND*
  - implementation of hierarchical annotations overcomes some limitations of "fuzzy" disease annotations
  - data *handover* (Beacon v1.1+) enables further data exploration and export scenarios



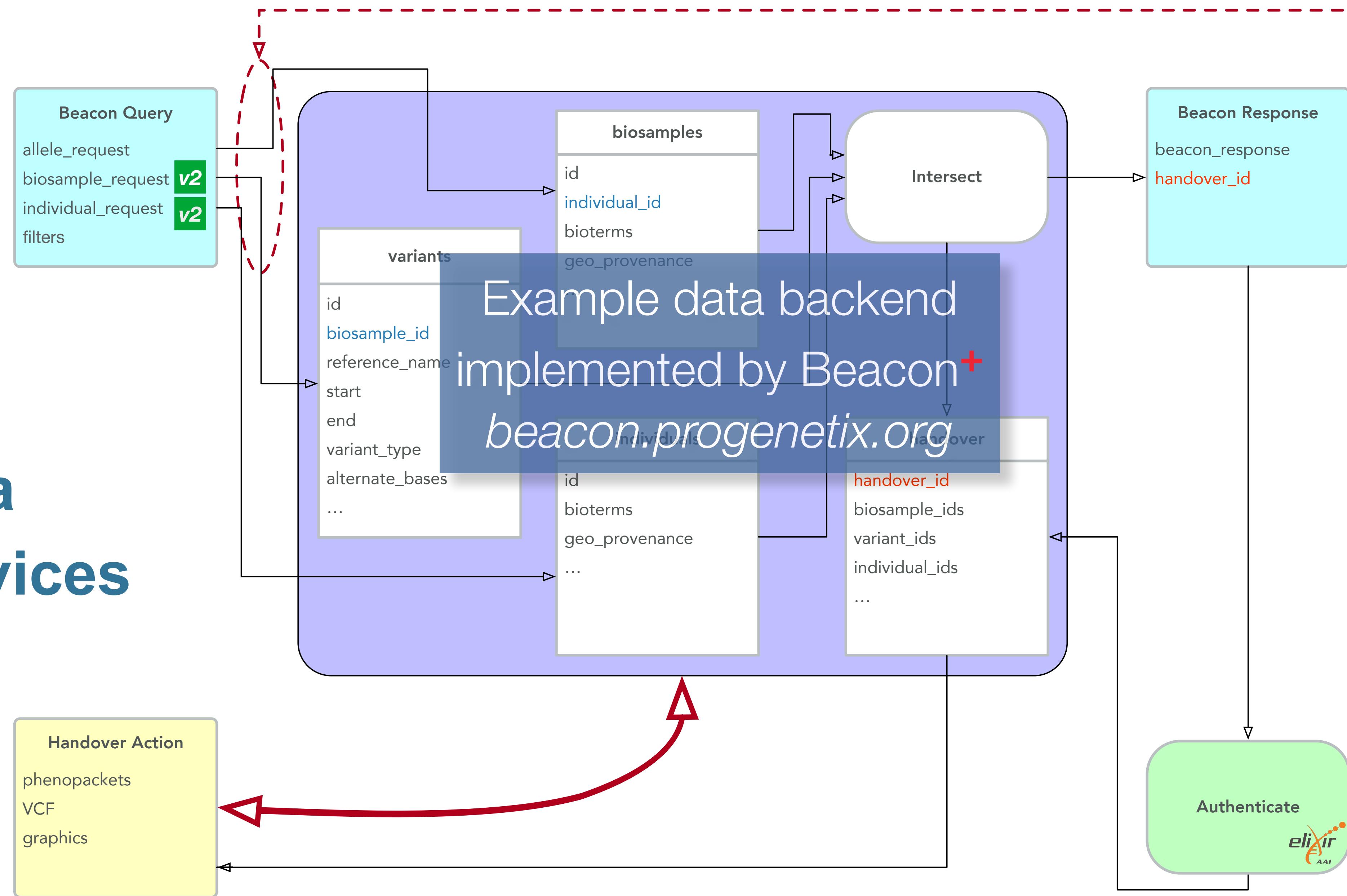
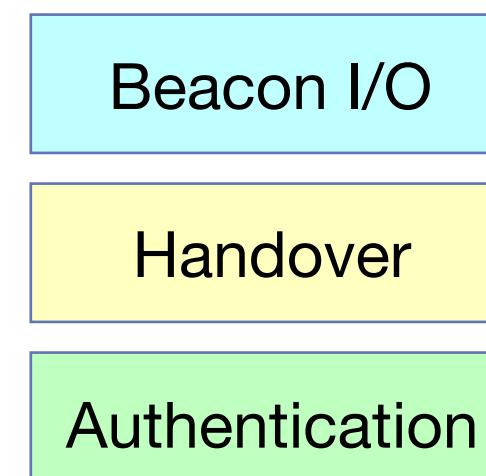
Beacon+ specific: Multiple term selection with OR logic

<input checked="" type="checkbox"/>	> <a href="#">NCIT:C4914: Skin Carcinoma</a>	213
<input type="checkbox"/>	> <a href="#">NCIT:C4475: Dermal Neoplasm</a>	109
<input checked="" type="checkbox"/>	> <a href="#">NCIT:C45240: Cutaneous Hematopoietic and Lymphoid Cell Neoplasm</a>	310

**Filters:** NCIT:C4914, NCIT:C4819, NCIT:C9231, NCIT:C2921, NCIT:C45240, NCIT:C6858, NCIT:C3467, NCIT:C45340, NCIT:C7195, NCIT:C3246, NCIT:C7217

# Beacon & Handover

Beacons v1.1  
supports data  
delivery services



# Beaconized Progenetix

## From Beacon Query to Explorative Analyses of CNV Patterns

- The consistent use of hierarchical diagnostic codes allows the use of Beacon "filters" for histopathological/clinically scoped queries
- Beacon's handover protocols can be utilized for data retrieval and, well, handing over to additional services, e.g.
  - downloads
  - visualization
  - use of external services (UCSC browser display...)



**Search**

**Dataset**: arraymap x | ▾

**Genome Assembly**: GRCh38 / hg38 | ▾

**Reference name**: 9 | ▾

**(Structural) Variant**: DEL (Deletion) | ▾

**Start**: 21000001-21975098 End Position: 21967753-23000000 | ▾

**Bio-ontology**: NCIT:C3058: Glioblastoma (2119) x | ▾

**Beacon Query**

Assembly: GRCh38 Chro: 9 Start: 21000001-21975098 End: 21967753-23000000 Type: DEL Ref. Base(s): N  
Filters: NCIT:C3058

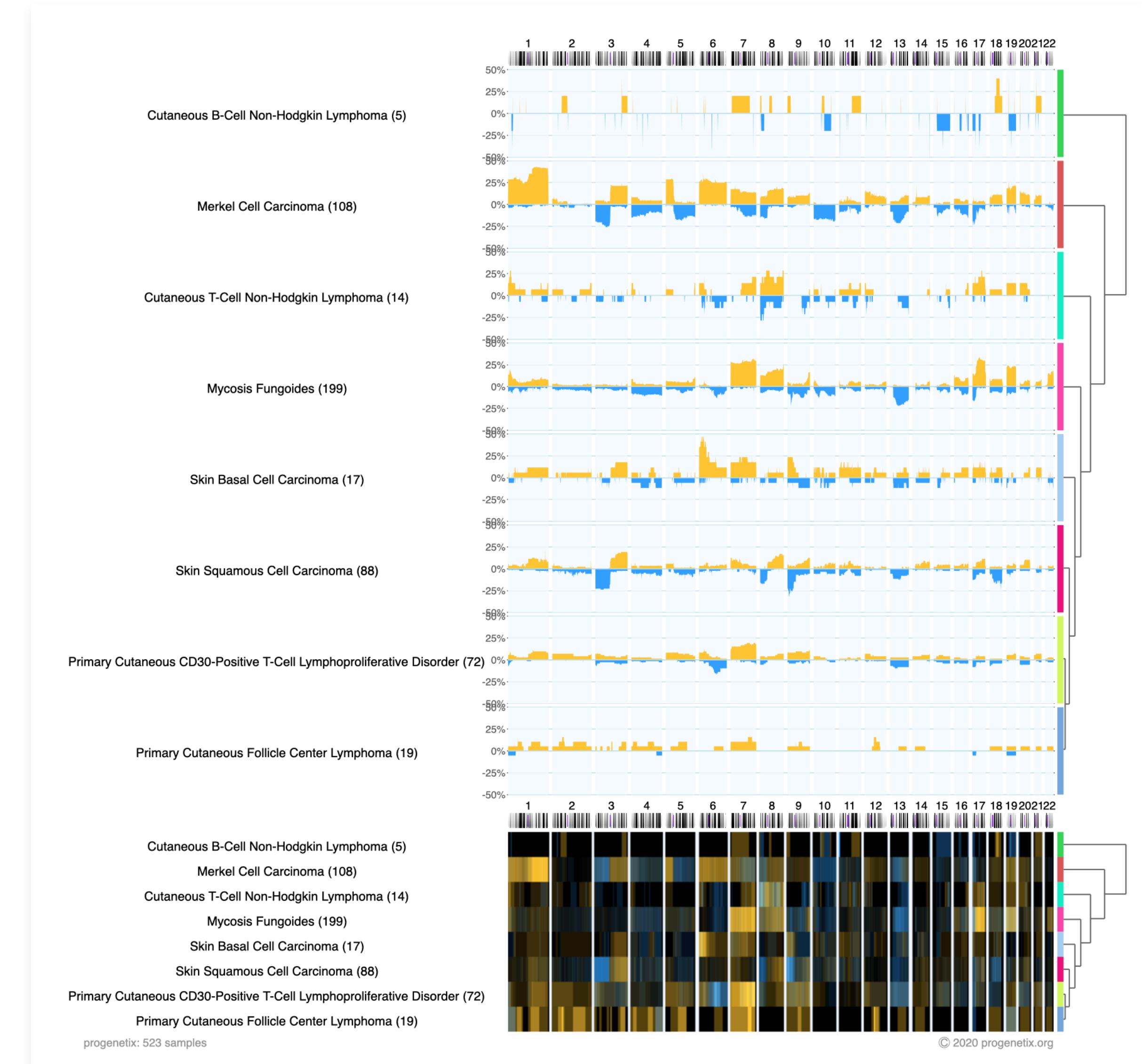
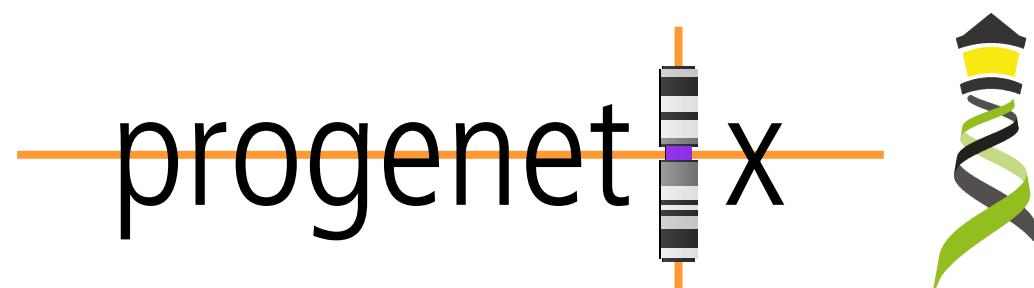


Subsets	Subset Samples	Query Matches	Subset Match Frequencies
icdom-94403	3251	298	0.092
NCIT:C3058	3303	298	0.090
icdot-C71.9	4760	297	0.062
icdot-C71.0	1712	1	0.001

# Beaconized Progenetix

## From Beacon Query to Explorative Analyses of CNV Patterns

- The consistent use of hierarchical diagnostic codes allows the use of Beacon "filters" for histopathological/clinically scoped queries
- Beacon's handover protocols can be utilized for data retrieval and, well, handing over to additional services, e.g.
  - ▶ downloads
  - ▶ visualization
  - ▶ use of external services (UCSC browser display...)



# Standardized Data

**Data re-use depends on standardized, machine-readable metadata**

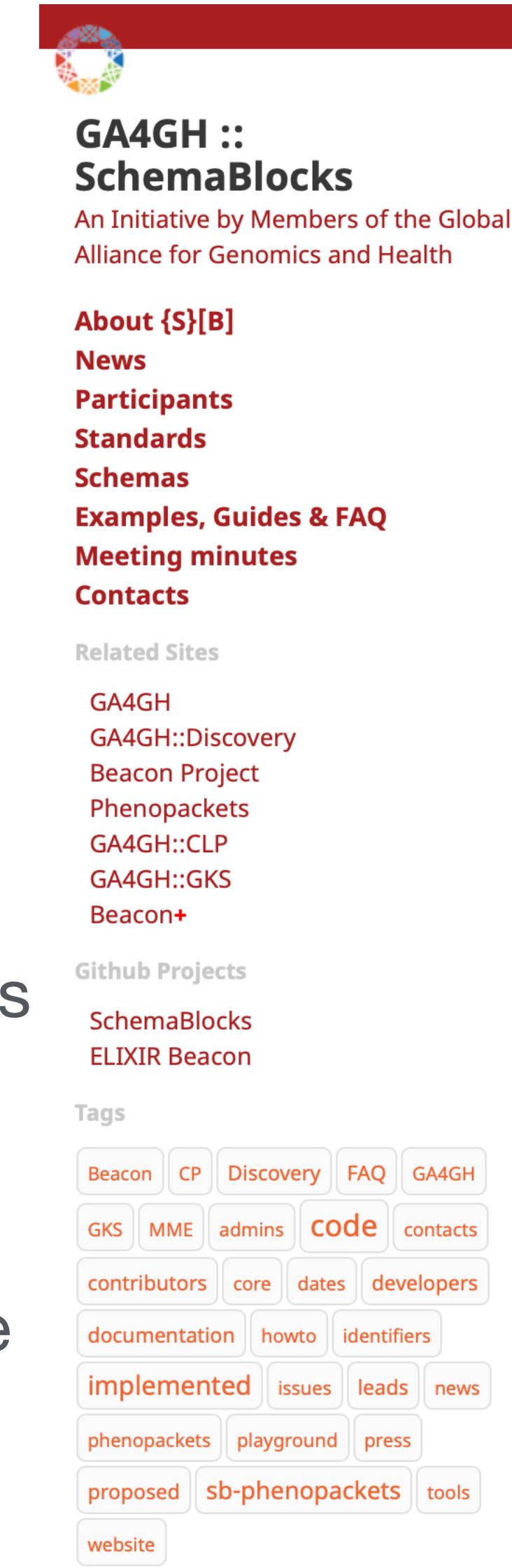
- Multiple international initiatives (ELIXIR, GA4GH, MONARCH...) and resource providers (EBI, NCBI ...) work on the generation and implementation of data annotation standards
- emerging / established principles are the use of hierarchical coding systems where individual codes are represented as CURIEs
- other formats for non-categorical annotations based on international standards, e.g.
  - ISO (ISO 8601 time & period, ISO 3166 country codes ...)
  - IETF (GeoJSON ...)
  - W3C (CURIE ...)
- these standards become pervasive throughout GA4GH's ecosystem (e.g. Phenopackets ...)

```
"data_use_conditions" : {  
    "label" : "no restriction",  
    "id" : "DUO:0000004"  
},  
  
"provenance" : {  
    "material" : {  
        "type" : {  
            "id" : "EFO:0009656",  
            "label" : "neoplastic sample"  
        }  
    },  
    "geo" : {  
        "label" : "Zurich, Switzerland",  
        "precision" : "city",  
        "city" : "Zurich",  
        "country" : "Switzerland",  
        "latitude" : 47.37,  
        "longitude" : 8.55,  
        "geojson" : {  
            "type" : "Point",  
            "coordinates" : [  
                8.55,  
                47.37  
            ]  
        },  
        "IS0-3166-alpha3" : "CHE"  
    },  
    {  
        "age": "P25Y3M2D"  
    }  
}
```

# GA4GH {S}[B] SchemaBlocks

Standardized formats and data schemas for developing an "Internet of Genomics"

- “cross-workstreams, cross-drivers” initiative to document GA4GH object **standards** and **prototypes**
- launched in December 2018
- documentation and implementation examples provided by GA4GH members
- not a rigid, complete data schema
- object **vocabulary** and **semantics** for a large range of developments
- ▶ **Beacon** as contributor and user



The screenshot shows the GA4GH SchemaBlocks homepage. At the top left is the GA4GH logo, followed by the text "GA4GH :: SchemaBlocks" and "An Initiative by Members of the Global Alliance for Genomics and Health". A sidebar on the left contains links for "About {S}[B]", "News", "Participants", "Standards", "Schemas", "Examples, Guides & FAQ", "Meeting minutes", and "Contacts". Below this are sections for "Related Sites" (GA4GH, GA4GH::Discovery, Beacon Project, Phenopackets, GA4GH::CLP, GA4GH::GKS, Beacon+) and "Github Projects" (SchemaBlocks, ELIXIR Beacon). A "Tags" section at the bottom lists various terms like Beacon, CP, Discovery, FAQ, GA4GH, GKS, MME, admins, code, contacts, contributors, core, dates, developers, documentation, howto, identifiers, implemented, issues, leads, news, phenopackets, playground, press, proposed, sb-phenopackets, tools, and website.

## GA4GH SchemaBlocks Home

SchemaBlocks is a “cross-workstreams, cross-drivers” initiative to document GA4GH object standards and prototypes, as well as common data formats and semantics.

Launched in December 2018, this project is still to be considered a “community initiative”, with developing participation, leadership and governance structures. At its current stage, the documents can **not** be considered “authoritative GA4GH recommendations” but rather represent documentation and implementation examples provided by GA4GH members.

While future products and implementations may be completely based on *SchemaBlocks* components, this project does not attempt to develop a rigid, complete data schema but rather to provide the object vocabulary and semantics for a large range of developments.

The SchemaBlocks site can be accessed through the permanent link [schemablocks.org](https://schemablocks.org). More information about the different products & formats can be found on the workstream sites. For reference, some of the original information about recommended formats and object hierarchies is kept in the [GA4GH Metadata repositories](#).

For more information on GA4GH, please visit the [GA4GH Website](#).

## SchemaBlocks Repositories

The SchemaBlocks Github organisation contains several specifically scoped repositories. Please use the relevant *Github Issues* to and/or GH pull requests comment and contribute there.

@mbaudis 2019-11-19: [more ...](#)

## SchemaBlocks “Status” Levels

SchemaBlocks schemas (“blocks”) provide recommended blueprints for schema parts to be re-used for the development of code based “products” throughout the GA4GH ecosystem. We propose a labeling system for those schemas, to provide transparency about the level of support those schemas have from {S}[B] participants and observers.

@mbaudis 2019-07-17: [more ...](#)

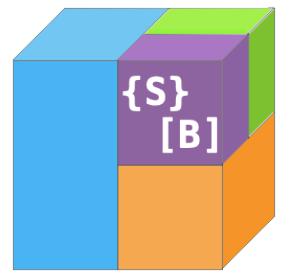
## SchemaBlocks {S}[B] Mission Statement

SchemaBlocks aims to translate the work of the workstreams into data models that:

- Are usable by other internal GA4GH deliverables, such as the Search API.
- Are usable by Driver Projects as an exchange format.
- Aid in aligning the work streams across GA4GH.
- Do not create a hindrance in development work by other work streams.

@mbaudis 2019-03-27: [more ...](#)





## BeaconAlleleRequest beacon ↗

{S}[B] Status [i]	implemented
Provenance	◦ Beacon API
Used by	◦ Beacon ◦ Progenetix database schema (Beacon+ backend)
Contributors	◦ Marc Fiume ◦ Michael Baudis ◦ Sabela de la Torre Pernas ◦ Jordi Rambla ◦ Beacon developers...
Source (v1.1.0)	◦ raw source [JSON] ◦ Github

### Attributes

Type: object

Description: Allele request as interpreted by the beacon.

### Properties

Property	Type
alternateBases	string
assemblyId	string
datasetIds	array of string
end	integer
endMax	integer
endMin	integer
mateName	<a href="https://schemablocks.org/schemas/beacon/v1.1.0/Chromosome">https://schemablocks.org/schemas/beacon/v1.1.0/Chromosome</a> [HTML]
referenceBases	string
referenceName	<a href="https://schemablocks.org/schemas/beacon/v1.1.0/Chromosome">https://schemablocks.org/schemas/beacon/v1.1.0/Chromosome</a> [HTML]
start	integer (int64)
startMax	integer
startMin	integer
variantType	string

### alternateBases

- type: string

The bases that appear instead of the reference bases. Accepted values: [ACGTN]\*. N is a wildcard, that denotes the position of any base, and can be used as a standalone base of any type or within a partially known sequence. For example a sequence where the first and last bases are known, but the middle portion can exhibit countless variations of [ACGT], or the bases are unknown: ANNT the Ns can take any form of [ACGT], which makes both ACCT and ATGT (or any other combination) viable sequences.

Symbolic ALT alleles (DEL, INS, DUP, INV, CNV, DUP:TANDEM, DEL:ME, INS:ME) will be represented in variantType.

Optional: either alternateBases or variantType is required.

### alternateBases Value Example

#### assemblyId

- type: string

Assembly identifier (GRC notation, e.g. GRCh37).

### assemblyId Value Example

## Curie sb-vr-spec ↗

{S}[B] Status [i]	implemented
Provenance	◦ vr-spec
Used by	◦ vr-spec
Contributors	◦ Reece Hart ◦ Michael Baudis

### Attributes

Type: string

Pattern: ^\w[^:]+:\$

Description: A string that refers to an object uniquely. This is a standard CURIE.

VR does not impose any constraints on strings used as identifiers, the VR Specification RECOMMENDS that implementers use CURIEs as identifiers. String CURIEs are represented as prefix:reference (W3C namespace:accession or namespace:local id colloquially).

The VR specification also RECOMMENDS that prefix be omitted. The reference component is an unconstrained string.

FHIR mapping: Specimen.

A CURIE is a URI. URIs may locate objects (i.e., specify what they point to). VR uses CURIEs primarily as a naming mechanism.

Implementations MAY provide CURIE resolution mechanisms.

Using internal IDs in public messages is strongly discouraged.

### Curie Value Examples

"ga4gh:GA\_01234abcde"

"DUO:0000004"

"orcid:0000-0003-3463-0775"

"PMID:15254584"

## Biosample sb-phenopackets ↗

{S}[B] Status [i]	implemented
Provenance	◦ Phenopackets
Used by	◦ Phenopackets
Contributors	◦ GA4GH Data Working Group ◦ Jules Jacobsen ◦ Peter Robinson ◦ Michael Baudis ◦ Melanie Courtot ◦ Isuru Liyanage

### Attributes

Type: object

Description: A Biosample refers to a unit of biological material from which the substrate molecules (e.g. genomic DNA, RNA, proteins) for molecular analyses (e.g. sequencing, array hybridisation, mass spectrometry) are extracted.

Examples would be a tissue biopsy, a single cell from a culture for single cell genome sequencing, or a fraction from a gradient centrifugation.

Several instances (e.g. technical replicates) or types of experiments (e.g. genomic array as well as experiments) may refer to the same Biosample.

FHIR mapping: Specimen.

### Properties

Property	Type
ageOfIndividualAtCollection	<a href="https://schemablocks.org/schemas/sb-phenopackets/v1.0.0/Age.json">https://schemablocks.org/schemas/sb-phenopackets/v1.0.0/Age.json</a> [SRC] [HTML]
ageRangeOfIndividualAtCollection	<a href="https://schemablocks.org/schemas/sb-phenopackets/v1.0.0/AgeRange.json">https://schemablocks.org/schemas/sb-phenopackets/v1.0.0/AgeRange.json</a> [SRC] [HTML]
description	string
diagnosticMarkers	array of <a href="https://schemablocks.org/schemas/sb-phenopackets/v1.0.0/OntologyClass.json">https://schemablocks.org/schemas/sb-phenopackets/v1.0.0/OntologyClass.json</a> [SRC] [HTML]
histologicalDiagnosis	<a href="https://schemablocks.org/schemas/sb-phenopackets/v1.0.0/OntologyClass.json">https://schemablocks.org/schemas/sb-phenopackets/v1.0.0/OntologyClass.json</a> [SRC] [HTML]
htsFiles	array of <a href="https://schemablocks.org/schemas/sb-phenopackets/v1.0.0/HtsFile.json">https://schemablocks.org/schemas/sb-phenopackets/v1.0.0/HtsFile.json</a> [SRC] [HTML]
id	string
individualId	string
isControlSample	boolean
phenotypicFeature	array of <a href="https://schemablocks.org/schemas/sb-phenopackets/v1.0.0/PhenotypicFeature.json">https://schemablocks.org/schemas/sb-phenopackets/v1.0.0/PhenotypicFeature.json</a> [SRC] [HTML]
procedure	<a href="https://schemablocks.org/schemas/sb-phenopackets/v1.0.0/Procedure.json">https://schemablocks.org/schemas/sb-phenopackets/v1.0.0/Procedure.json</a> [SRC] [HTML]
sampledTissue	<a href="https://schemablocks.org/schemas/sb-phenopackets/v1.0.0/Tissue.json">https://schemablocks.org/schemas/sb-phenopackets/v1.0.0/Tissue.json</a> [SRC] [HTML]

## Checksum sb-checksum ↗

{S}[B] Status [i]	proposed
Provenance	◦ GA4GH DRS (`develop` branch)
Used by	◦ GA4GH DRS ◦ GA4GH TRS
Contributors	◦ Susheel Varma

### Attributes

Type: object

Description: Checksum

### Properties

Property	Type
checksum	string
type	string

### checksum

- type: string

The hexadecimal encoded (Base16) checksum for the data.

### checksum Value Example

"77af4d6b9913e693e8d0b4b294fa62ade6054e6b2f1ffb617ac955dd63fb0182"

### type

- type: string

The digest method used to create the checksum. The value (e.g. sha-256) SHOULD be listed as Hash Name String in the GA4GH Hash Algorithm Registry. Other values MAY be used, as long as implementors are aware of the issues discussed in RFC6920.

GA4GH may provide more explicit guidance for use of non-IANA-registered algorithms in the future.

### type Value Example

"sha-256"

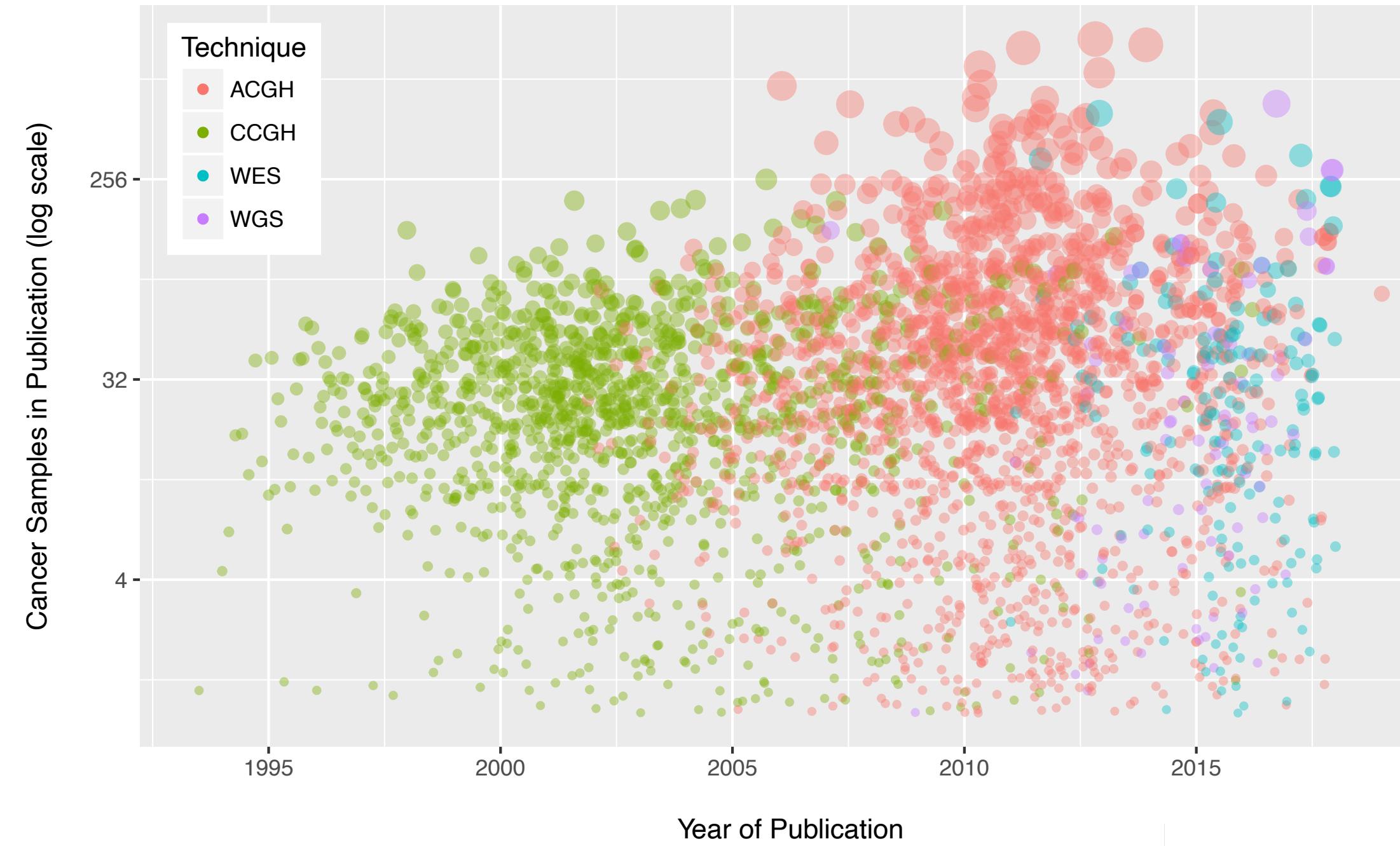
# Cancer (Genomics) Use Case

## Lessons from Implementing a Cancer CNV Resource w/ GA4GH & ELIXIR Standards

- Beacon v1.1 -> 2.alpha is already very capable for data search & delivery, at least for the "minimal expected subset" of associated data
  - some additions (query logic, data aggregation...) helpful in increasing utility w/o breaking the standard
- the Beacon v2 "filters" paradigm coupled with consistent use of CURIEs has proven very reliable & should drive adoption of standard ontologies throughout resources
  - may require remapping or translation service integration by resource providers
- GA4GH standards such as DUO for and Passport together w/ ELIXIR AAI (&open implementations) will enable opening of controlled access to yet closed resources
- several ELIXIR initiatives, platforms and projects work on integrating and driving GA4GH standards, partially through use of standard documentation initiatives like **SchemaBlocks**
- consistent use of metadata formats allows new insights into "accidentally collected data"



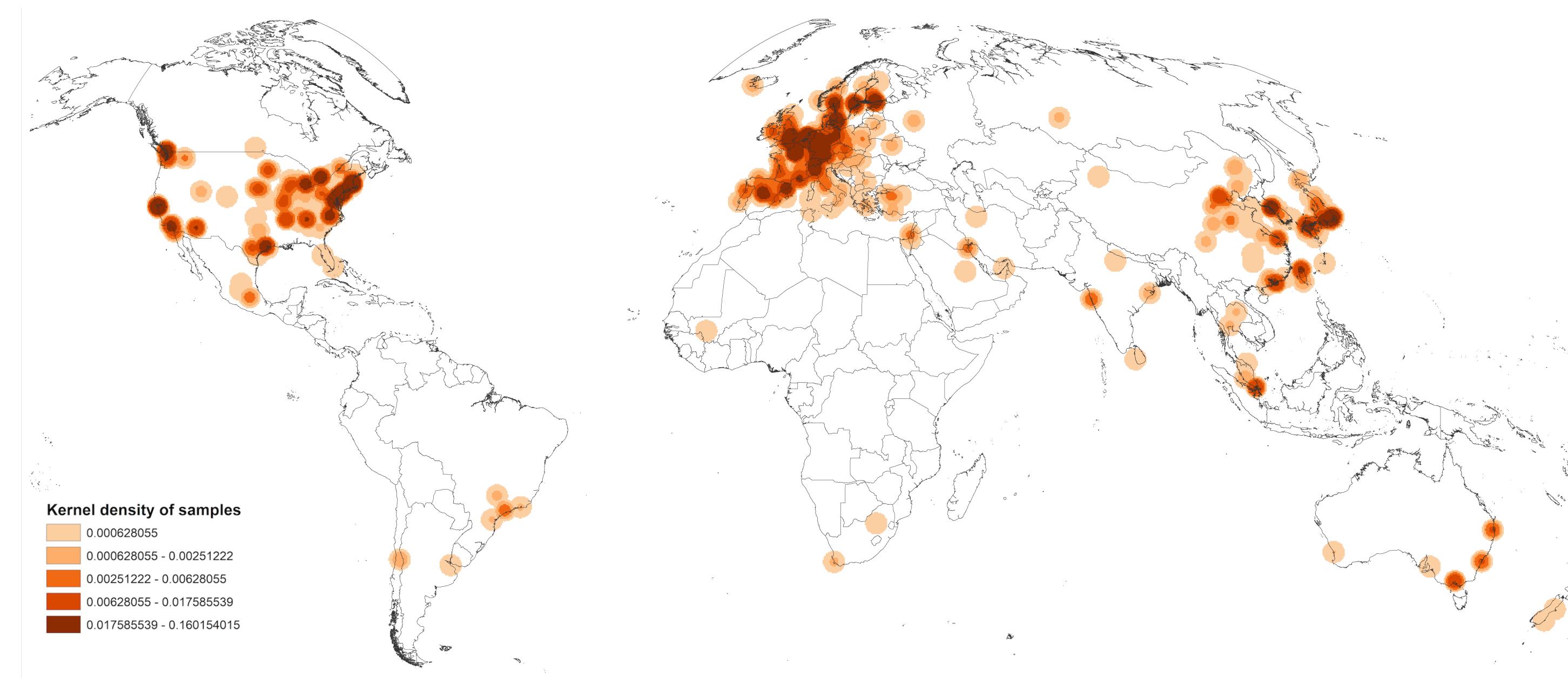
Map of the geographic distribution (by first author affiliation) of the 104'543 genomic array, 36'766 chromosomal CGH and 15'409 whole genome/exome based cancer genome datasets.  
The numbers are derived from the 3'240 publications registered in the Progenetix database.



## Publication Landscape of Cancer CNV Profiling

Publication statistics for cancer genome screening studies. The graphic shows our assessment of publications reporting whole-genome screening of cancer samples, using molecular detection methods (chromosomal CGH, genomic array technologies, whole exome and genome sequencing).

For the years 1993-2018, we found 3'229 publications reporting 174'530 individual samples in single series from 1 to more than 1000 samples. Y-axis and size of the dots correspond to the sample number; the color codes indicate the technology used.



# ELIXIR Beacon Network



- developed under lead from ELIXIR Finland
- **authenticated access** w/ ELIXIR AAI
- **incremental extension**, starting with ELIXIR Beacon resources adhering to the **latest specification** (contrast to legacy networks)
- service details provided by individual Beacons, using **GA4GH service-info**
- **registration service**
  - integrator** throughout ELIXIR Human Data
  - starting point for "**beyond ELIXIR**" **feature rich** federated Beacon services

GRCh38 ▾ 17 : 7577121 G > A

[Example variant query](#)

[Advanced Search](#)

baudisgroup at UZH and SIB  
Progenetix Cancer Genomics Beacon+

Beacon+ provides a forward looking implementation of the Beacon API, with focus on structural variants and metadata based on the cancer and reference genome profiling data represented in the Progenetix oncogenomic data resource (<https://progenetix.org>).

[Visit Us](#) · [Beacon API](#) · [Contact Us](#)

National Bioinformatics Infrastructure Sweden  
SweFreq Beacon

Beacon API Web Server based on the GA4GH Beacon API

[Visit Us](#) · [Beacon API](#) · [Contact Us](#)

LCSB at University of Luxembourg  
ELIXIR.LU Beacon

ELIXIR.LU Beacon

[Visit Us](#) · [Beacon API](#) · [Contact Us](#)

Research Programme on Biomedical Informatics  
DisGeNET Beacon

Variant-Disease associations collected from curated resources and the literature

[Visit Us](#) · [Beacon API](#) · [Contact Us](#)

European Genome-Phenome Archive (EGA)  
EGA Beacon

This [Beacon](https://beacon-project.io/) is based on the GA4GH Beacon [v1.1.0](https://github.com/ga4gh/beacon/specification/blob/develop/beacon.yaml)

[Visit Us](#) · [Beacon API](#) · [Contact Us](#)

University of Tartu Institute of Genomics, Estonia  
Beacon at the University of Tartu, Estonia

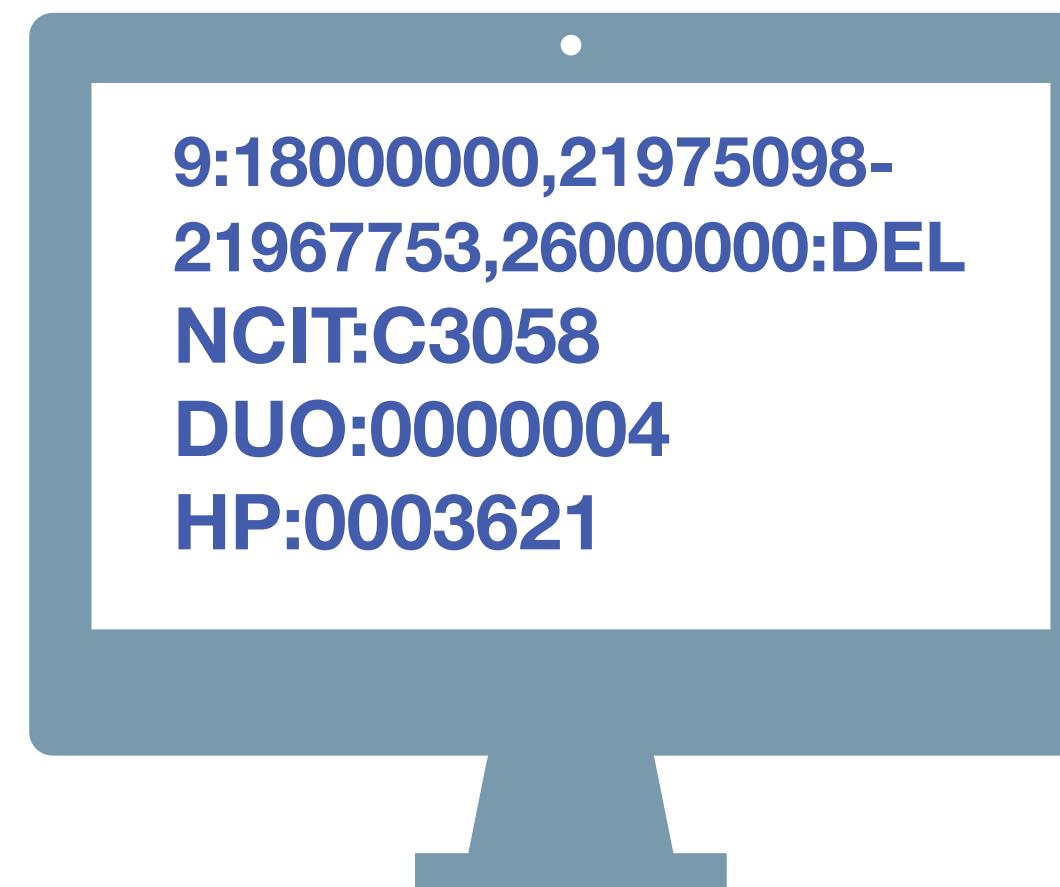
Beacon API Web Server based on the GA4GH Beacon API

[Visit Us](#) · [Beacon API](#) · [Contact Us](#)

CSC - IT Center for Science Production Beacon

Beacon API Web Server based on the GA4GH Beacon API

[Visit Us](#) · [Beacon API](#) · [Contact Us](#)



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".





## ELIXIR Cancer Data Focus Group



[www.elixir-europe.org](http://www.elixir-europe.org)

# hCNV Community

- Community officially approved in February 2019
- Leadership



Christophe Bérroud  
(ELIXIR France)



David Salgado  
(ELIXIR France)



Gary Saunders  
(Human Data Coordinator,  
ELIXIR Hub)



Michael Baudis  
(ELIXIR Switzerland)

## Mission statement

Despite the fact that **Copy Number Variations** are the **most prevalent genetic mutation type**, identifying and interpreting them is still a major challenge. The ELIXIR human Copy Number Variation (hCNV) Community aims to implement processes to make the **detection, annotation** and **interpretation** of these variations easier





Cancer CNV Profiles

Search Samples

Publication DB

Info

Beacon+

Baudisgroup @ UZH

**Cancer genome data @ progenetix.org**

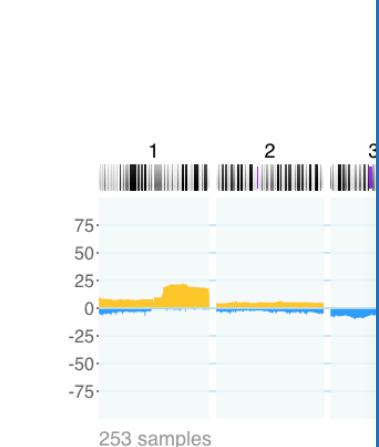
The Progenetix database contains cancer genome abnormalities (CNV / CNA)

**Data Content and Processing**

The resource currently contains cancer genome data derived from array and chip sequencing. Genomes are in Whole Exome or Whole Genome data formats or are extracted from other sources.

Besides genomic profiling, clinical information which so far has been collected.

Original diagnoses are mainly based on cancer types, according to the International Classification of Diseases for Oncology, respectively.



For exploration of the resources of interest.

Additionally to genome profiles, there are also publications (currently 3962 articles) available.

**Access, Maintenance and Contributions**

The content of the progenetix resource is freely accessible for research and educational purposes.

The database & software are developed by the [group of Michael Baudis](#) at the Institute of Bioinformatics SIB.

Many previous members and external collaborators have contributed to data collection and analysis. Participation (features, data, comments) by volunteers are welcome.

[progenetix.org](http://progenetix.org)  
[github.com/progenetix](https://github.com/progenetix)

**Search**

Allele Request Range Query All

**Baudisgroup @ UZH**

(Ni Ai)

Michael Baudis

(Haoyang Cai)

Paula Carrio Cordo

Bo Gao

Qingyao Huang

(Saumya Gupta)

(Nitin Kumar)

Rahel Paloots

Pierre-Henri Toussaint

**Start**

19000001-21975098

**End Position**

21967753-24000000

**Bio-ontology**

Select...

**Beacon Query****{S}[B] and GA4GH**

Melanie Courtot

Helen Parkinson

many more ...

**Beacon API Leads**

Jordi Rambla

Anthony Brooks

Juha Törnroos

**Discovery WS**

Michael Baudis (Beacon)

Marc Fiume (Networks)

**ELIXIR**

Gary Saunders

David Lloyd

Serena Scollen

The Beacon protocol defines an open standard for genomics data discovery, developed by members of the Global Alliance for Genomics & Health. It provides a framework for public web services responding to queries against genomic data collections, for instance from population based or disease specific genome repositories.

This repository contains the specification for the v2 major version upgrade of the Beacon API. It is now (2020) under active development and has *not* seen a stable code release.

For further information, please follow the work here and consult the [Beacon Project website](#).

Sabela de la Torre Pernas

Push your first package

Contributors 3

 sdelatorrep sdelatorrep

 mbaudis mbaudis

 blankdots blankdots

[next.progenetix.org/beacon-plus](http://next.progenetix.org/beacon-plus)

[beacon-project.io](http://beacon-project.io)  
[github.com/ga4gh-beacon/](https://github.com/ga4gh-beacon/)

