

# Progenetix Genomics Resource

## From Experiments to APIs

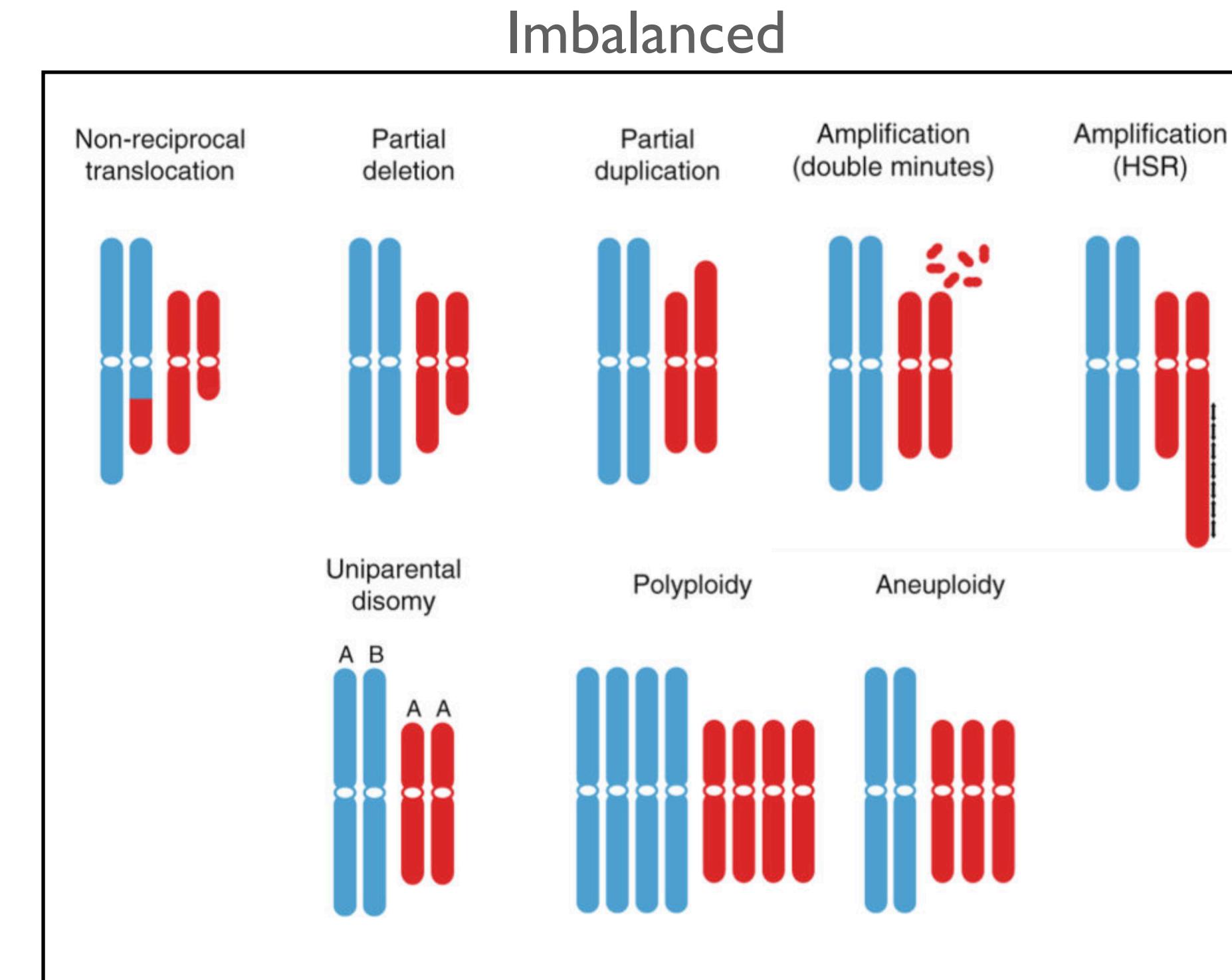
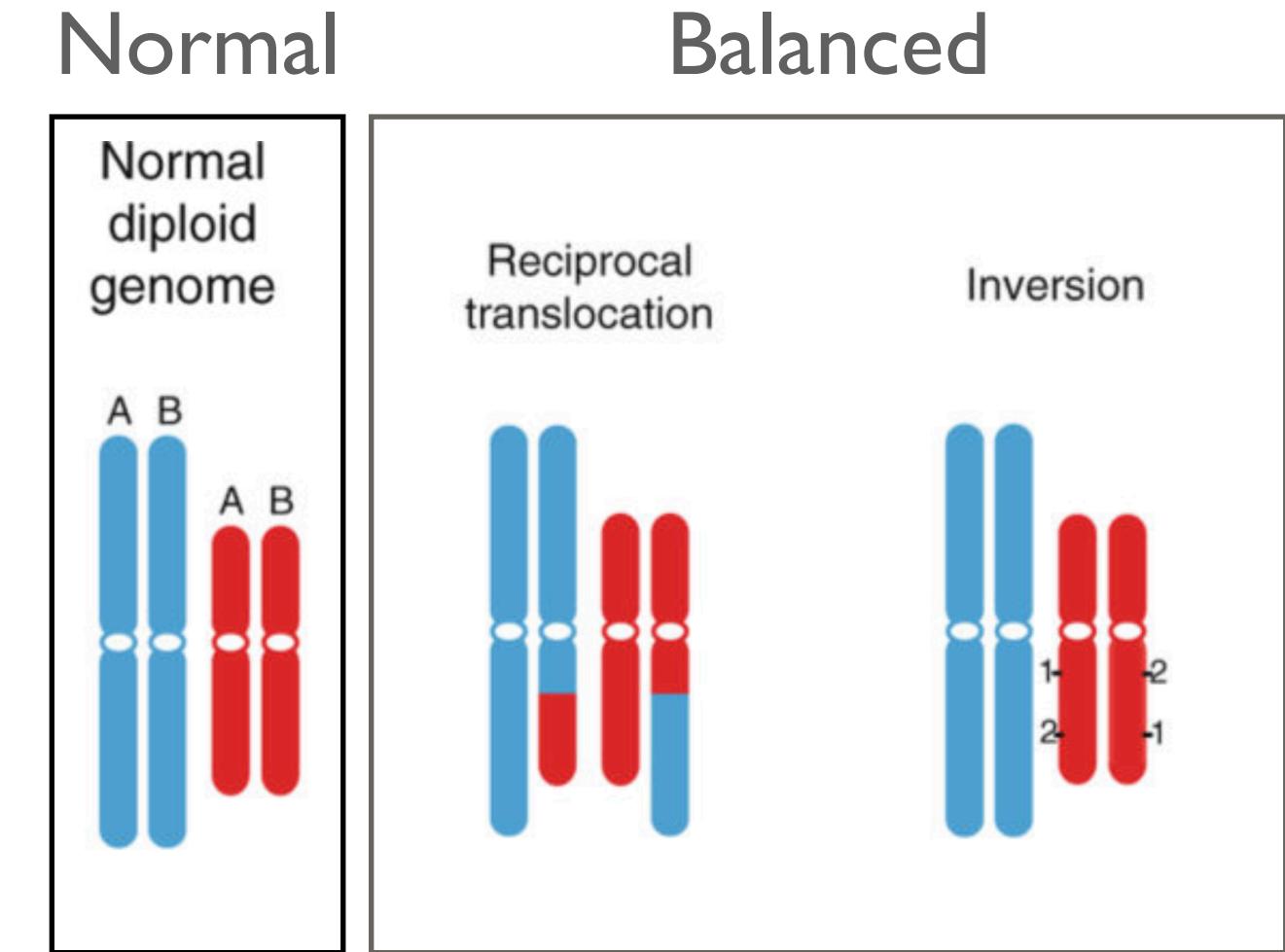
Michael Baudis | 2022



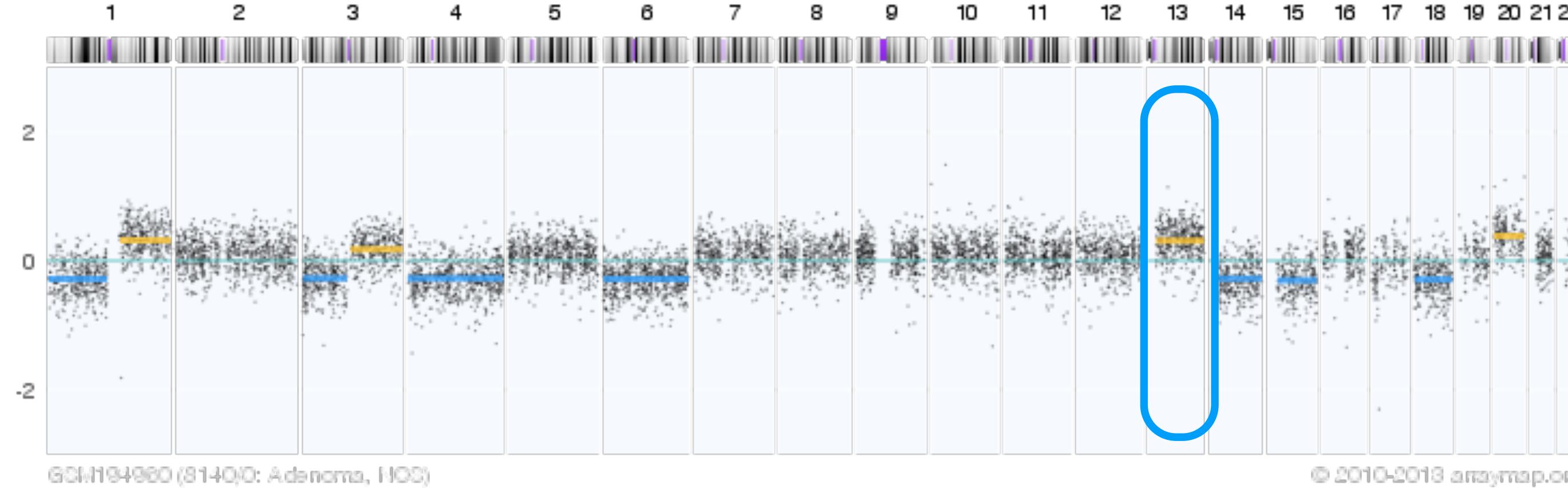
# Types of genomic alterations in Cancer

## Imbalanced Chromosomal Changes: CNV

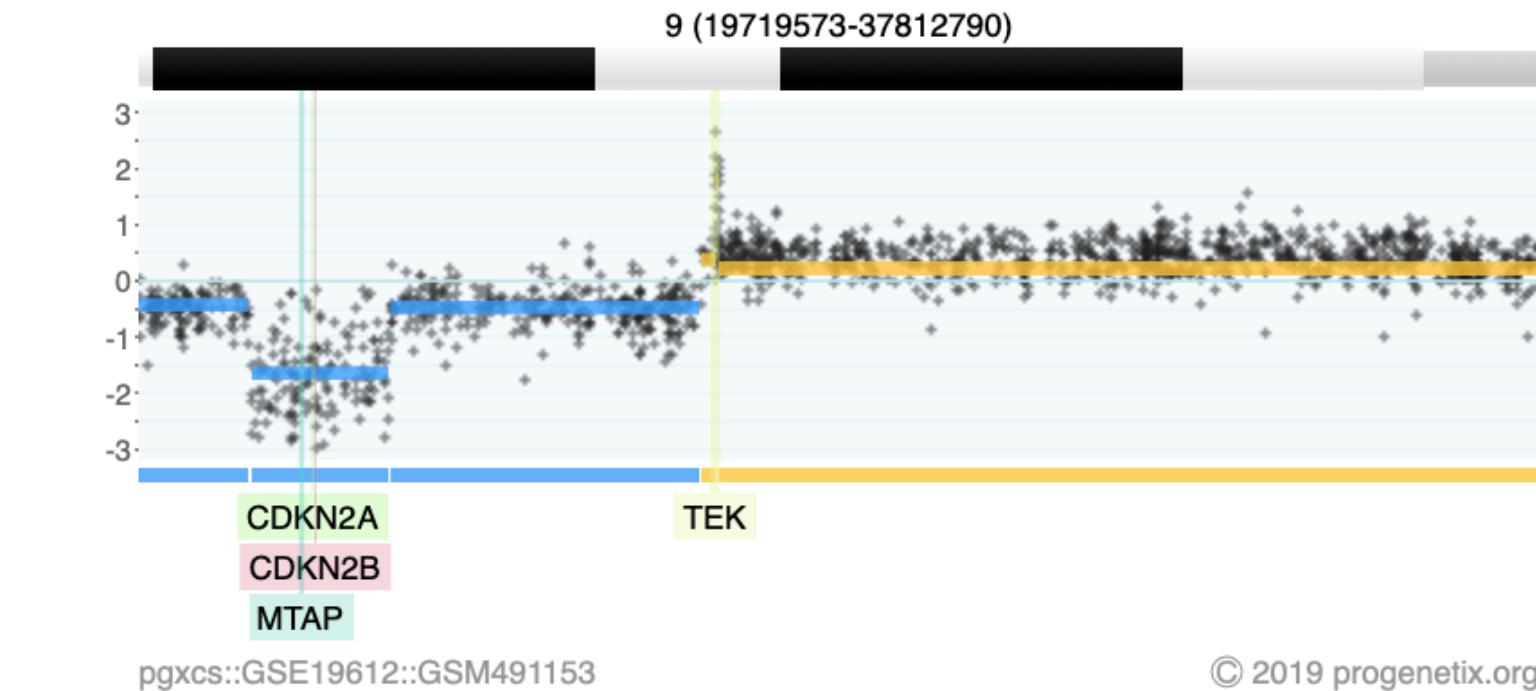
- Point mutations (insertions, deletions, substitutions)
- Chromosomal rearrangements
- Structural chromosomal Aberrations
- **Regional Copy Number Alterations**  
(losses, gains)
- Epigenetic changes (e.g. DNA methylation abnormalities)



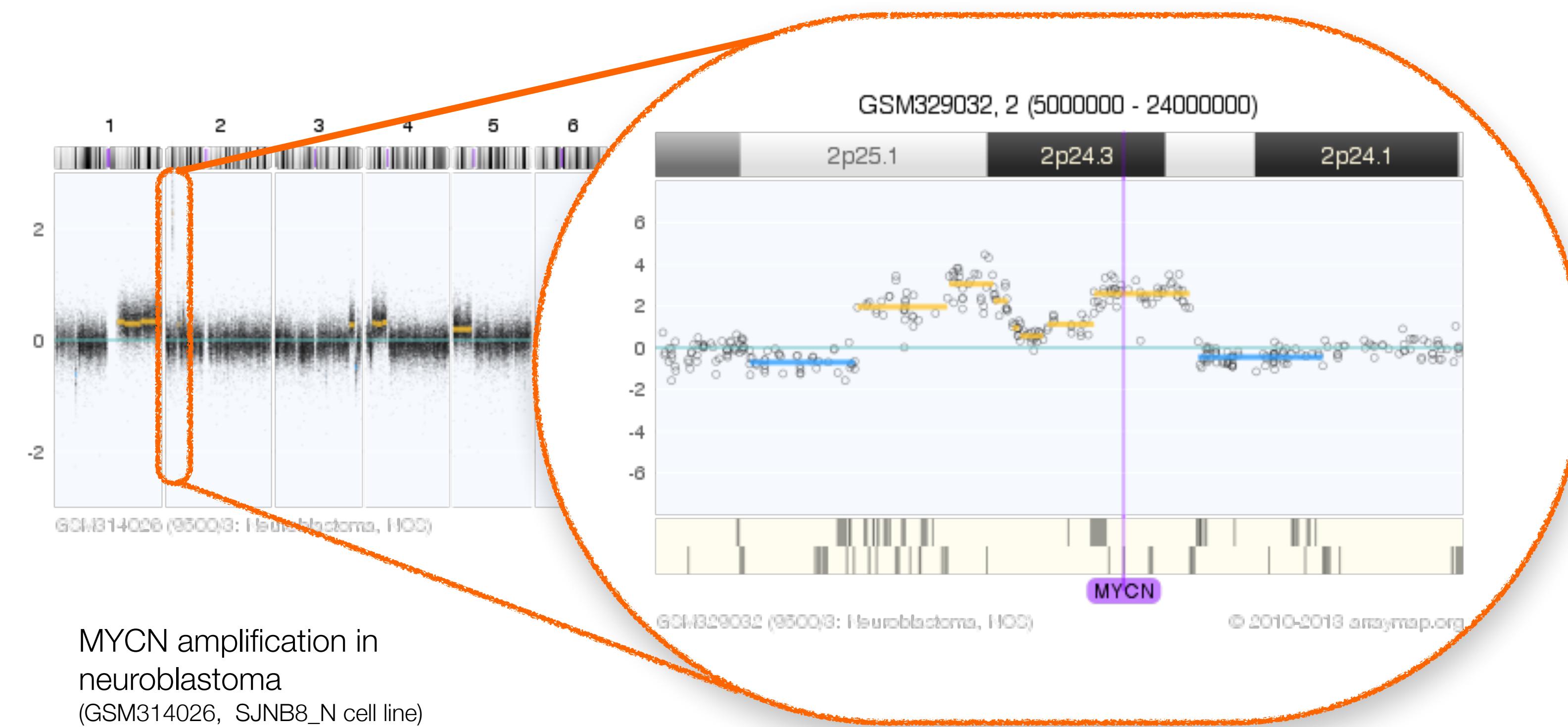
# Somatic Copy Number Variations



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)

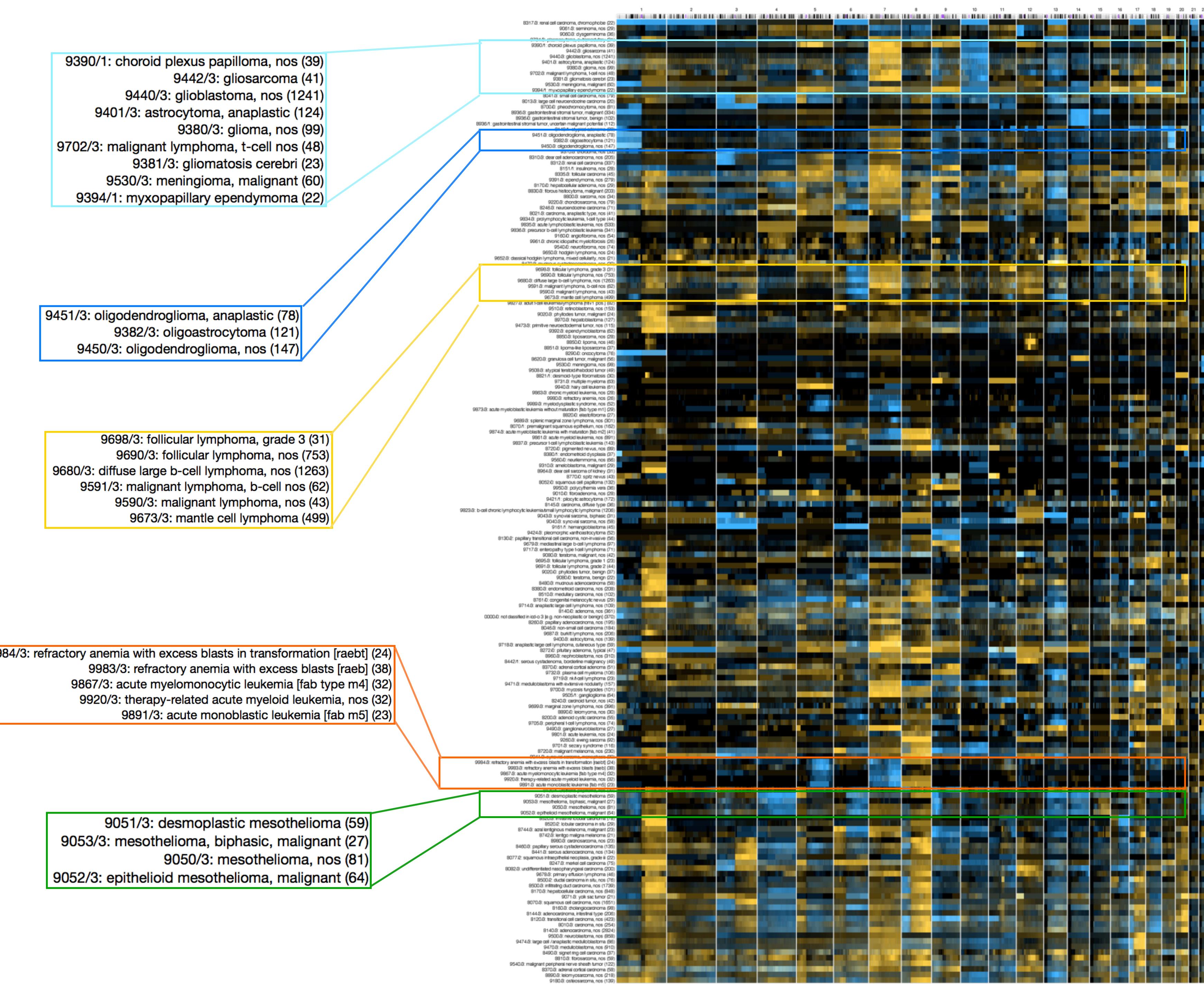
arrayMap



# Somatic Mutations In Cancer: Patterns

## Making the case for genomic classifications

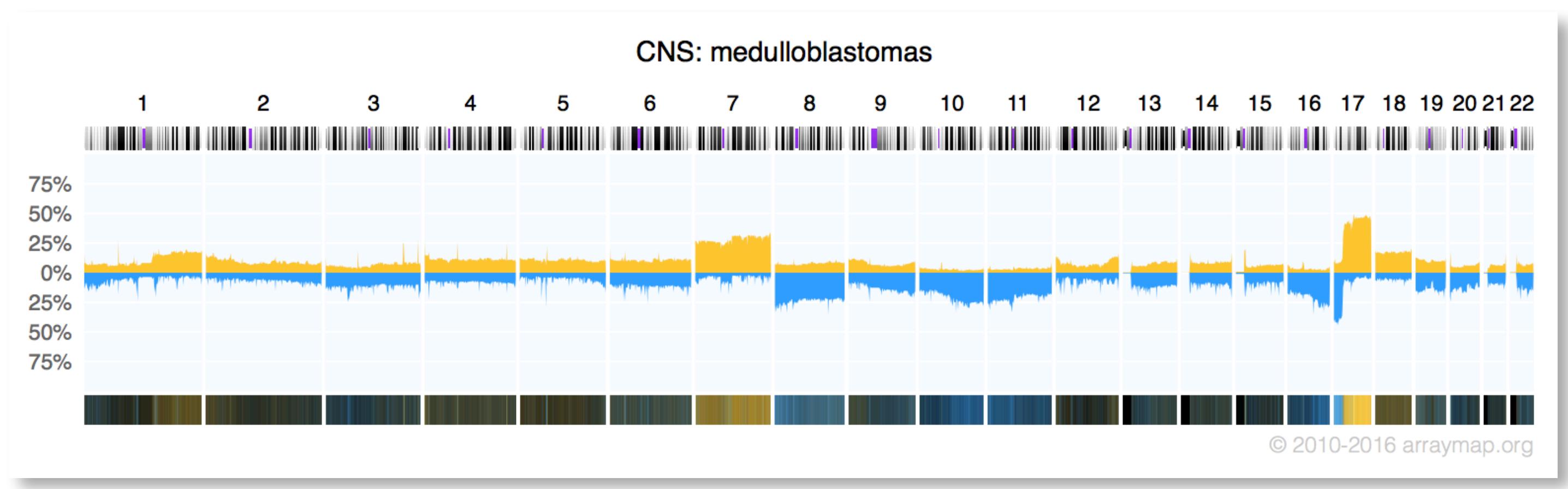
Some related cancer entities show similar copy number profiles



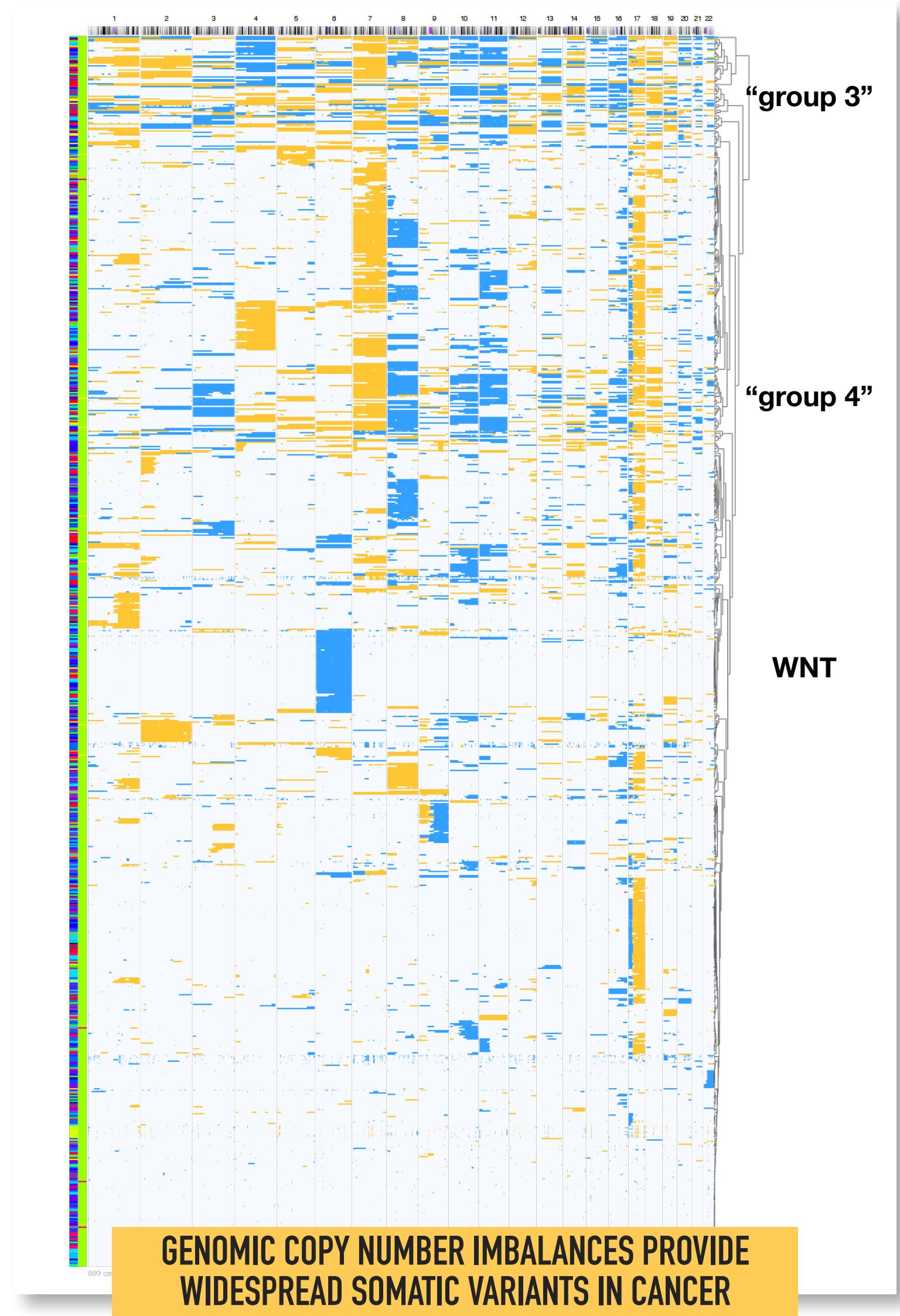
# Somatic CNVs In Cancer

## Recurrent mutation patterns

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

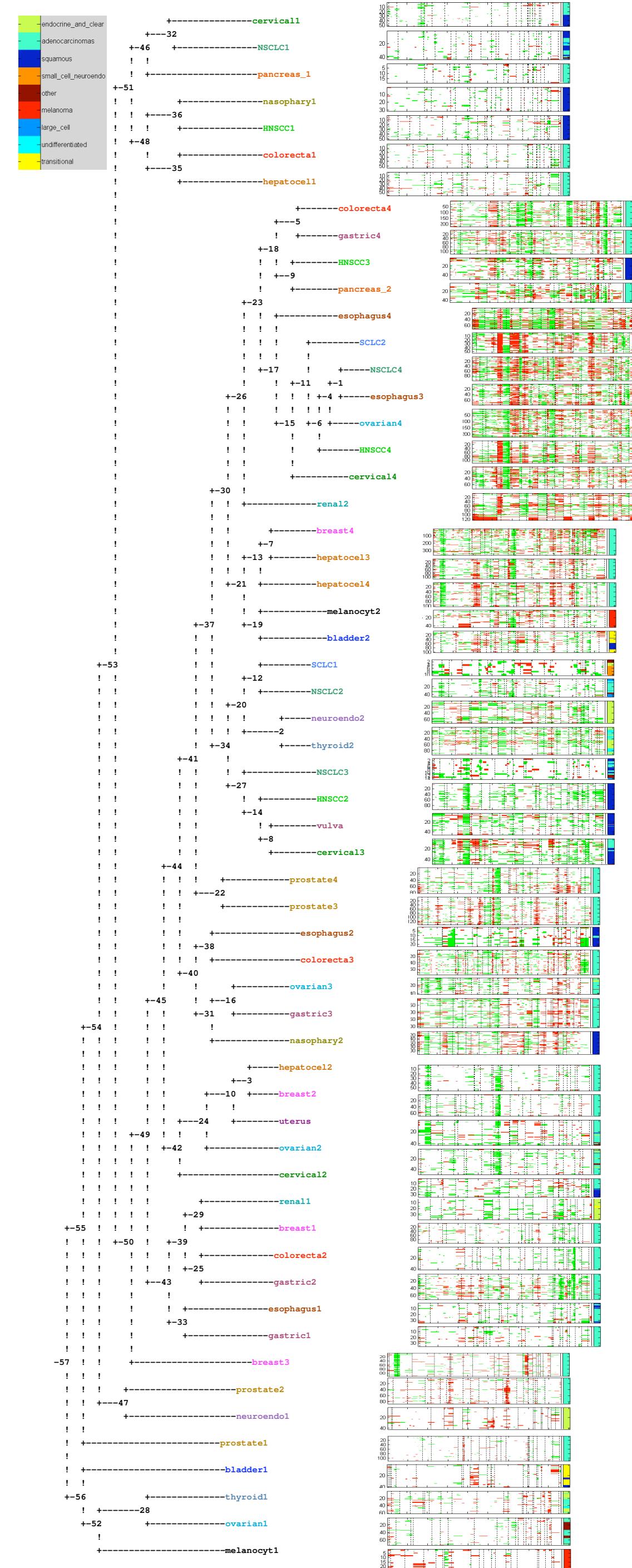
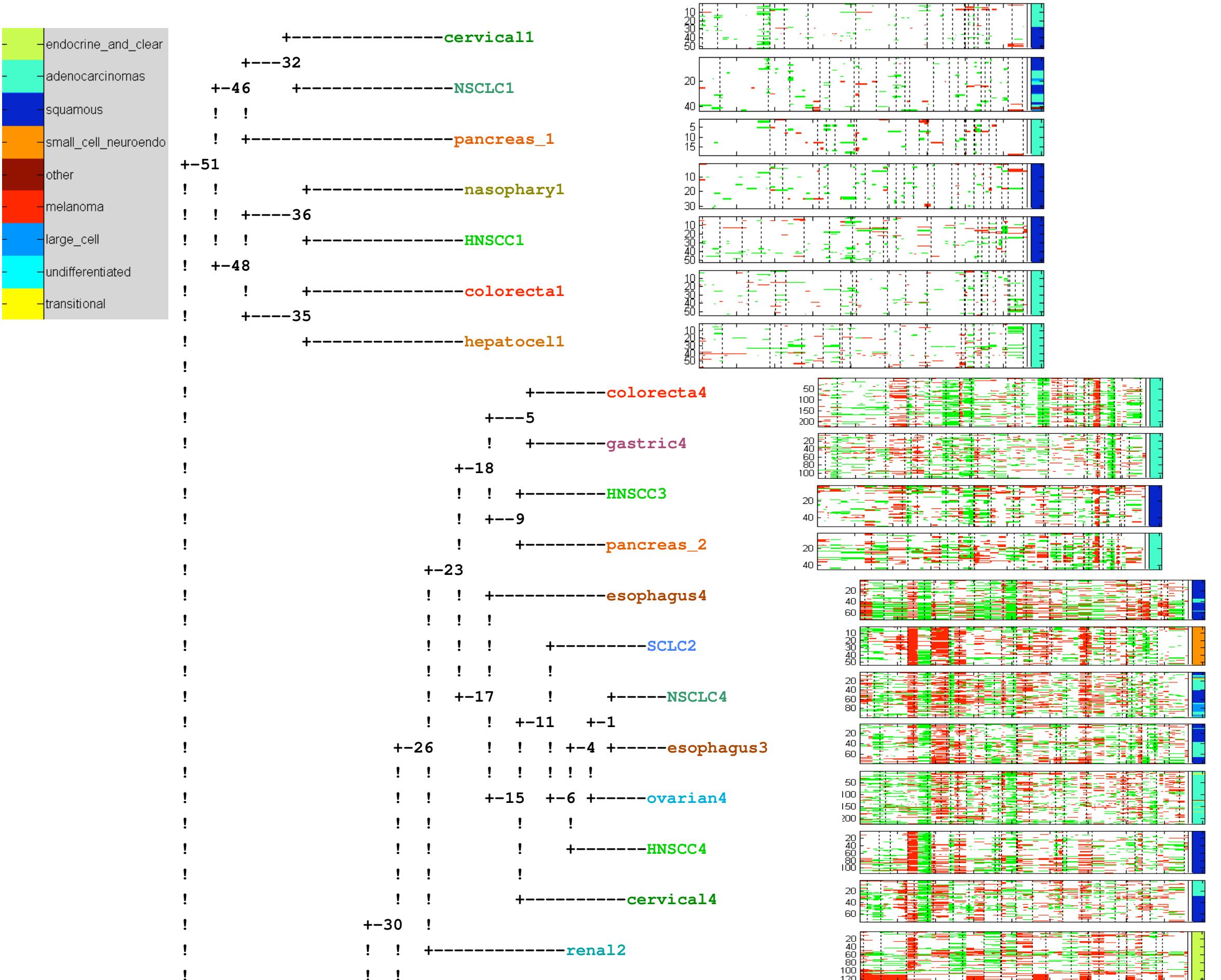


A genomic copy number histogram for malignant medulloblastomas, the most frequent type of pediatric brain tumors, displaying regions of genomic duplications and deletions. These can be decomposed into individual tumor profiles which segregate into several clusters of related mutation patterns with functional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



## Gene expression

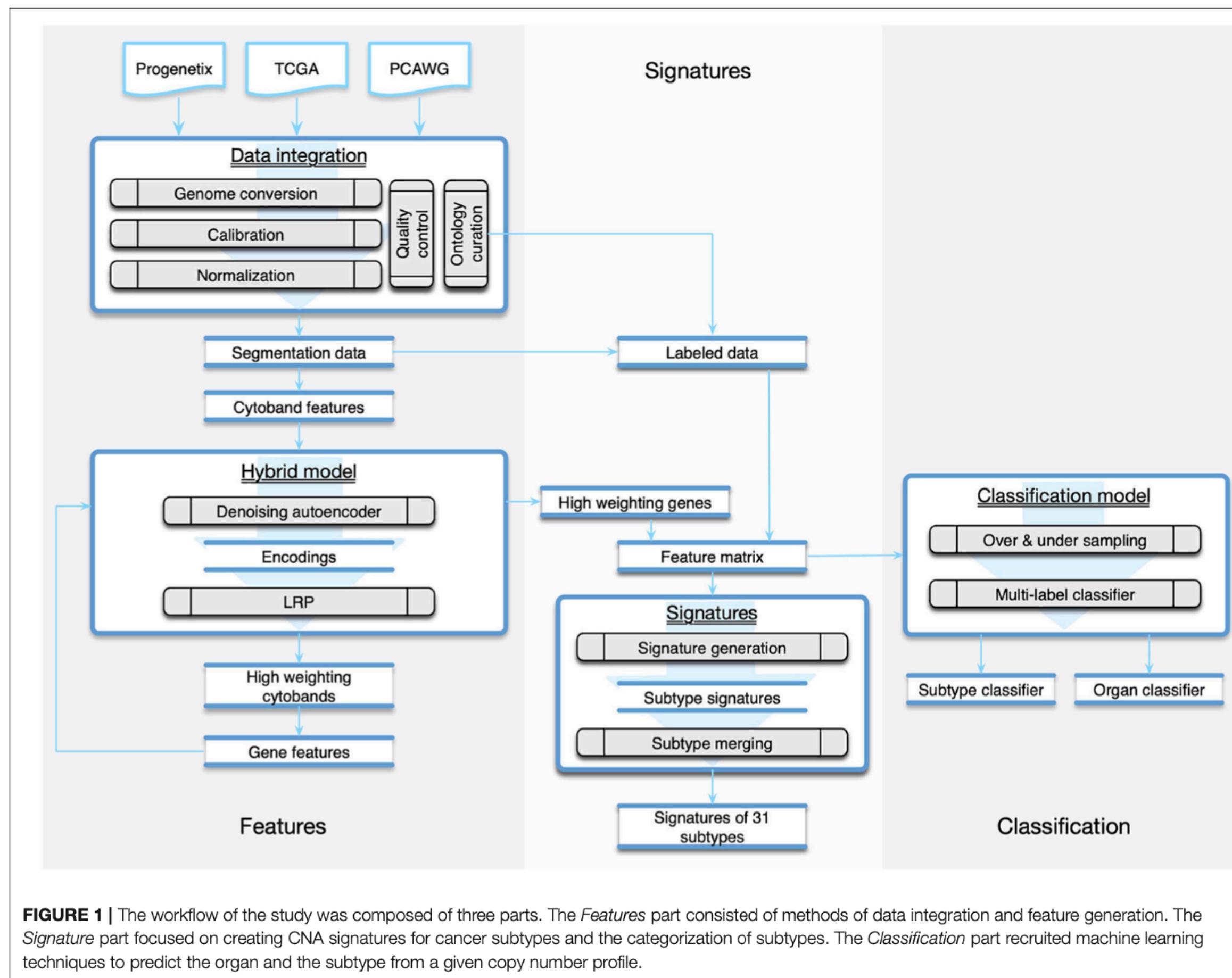
## Inferring progression models for CGH data

Jun Liu<sup>1</sup>, Nirmalya Bandyopadhyay<sup>1,\*</sup>, Sanjay Ranka<sup>1</sup>, M. Baudis<sup>2</sup> and Tamer Kahveci<sup>1,\*</sup><sup>1</sup>Computer and Information Science and Engineering, University of Florida, Gainesville, FL, USA and <sup>2</sup>Institute for Molecular Biology, University of Zurich, Zurich, Switzerland

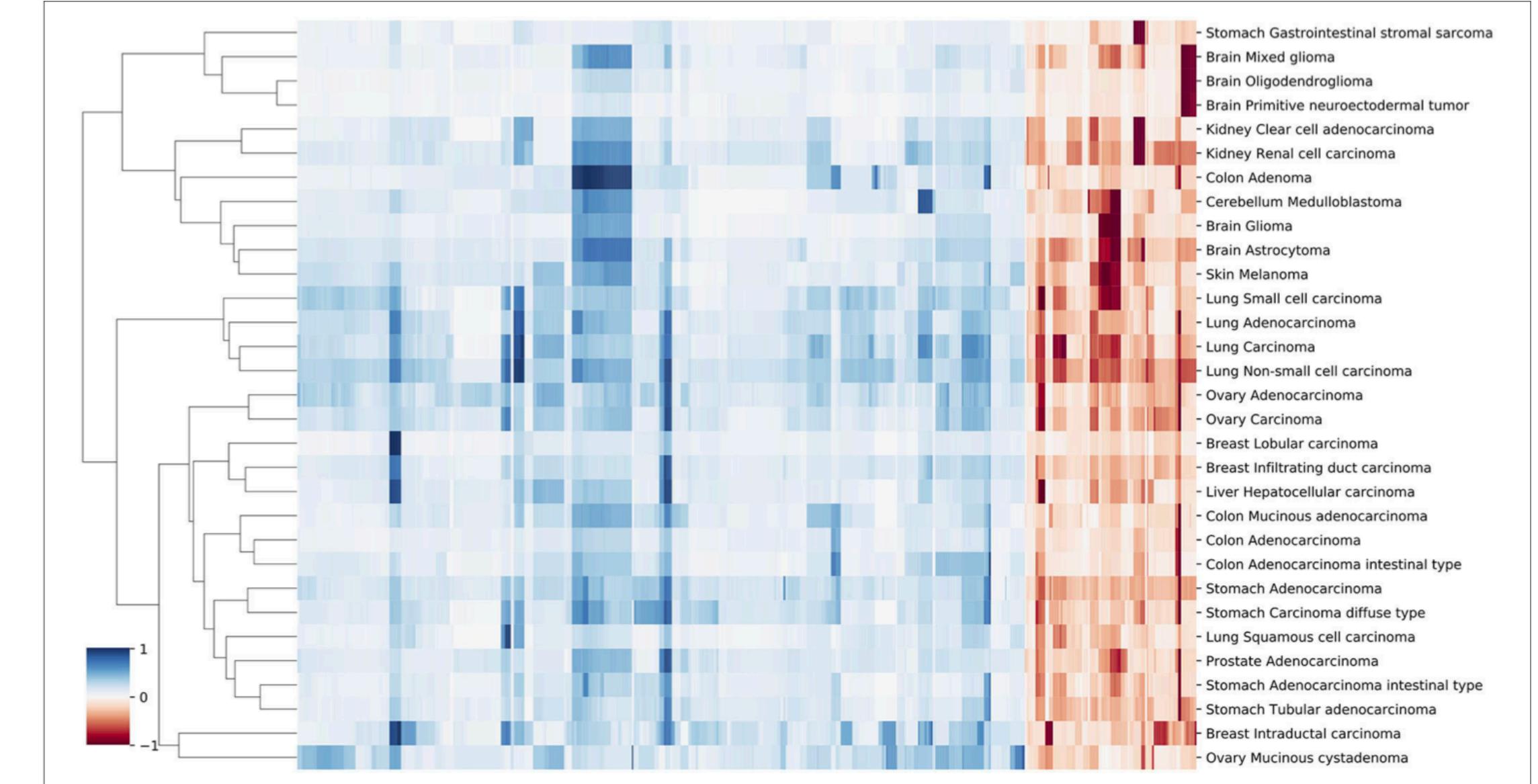


# Signatures of Discriminative Copy Number Aberrations in 31 Cancer Subtypes

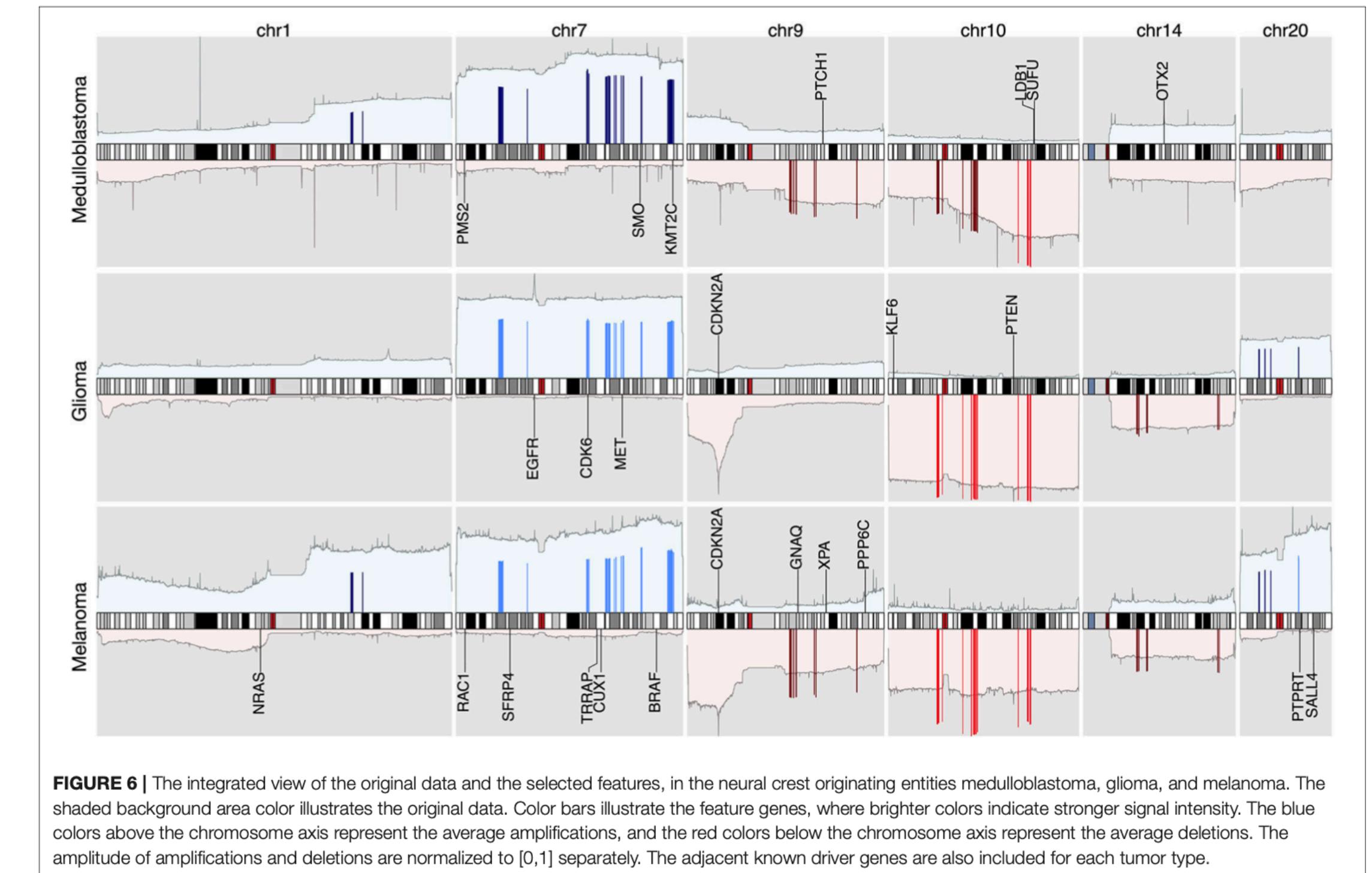
Bo Gao<sup>1,2</sup> and Michael Baudis<sup>1,2\*</sup>



**FIGURE 1 |** The workflow of the study was composed of three parts. The *Features* part consisted of methods of data integration and feature generation. The *Signature* part focused on creating CNA signatures for cancer subtypes and the categorization of subtypes. The *Classification* part recruited machine learning techniques to predict the organ and the subtype from a given copy number profile.



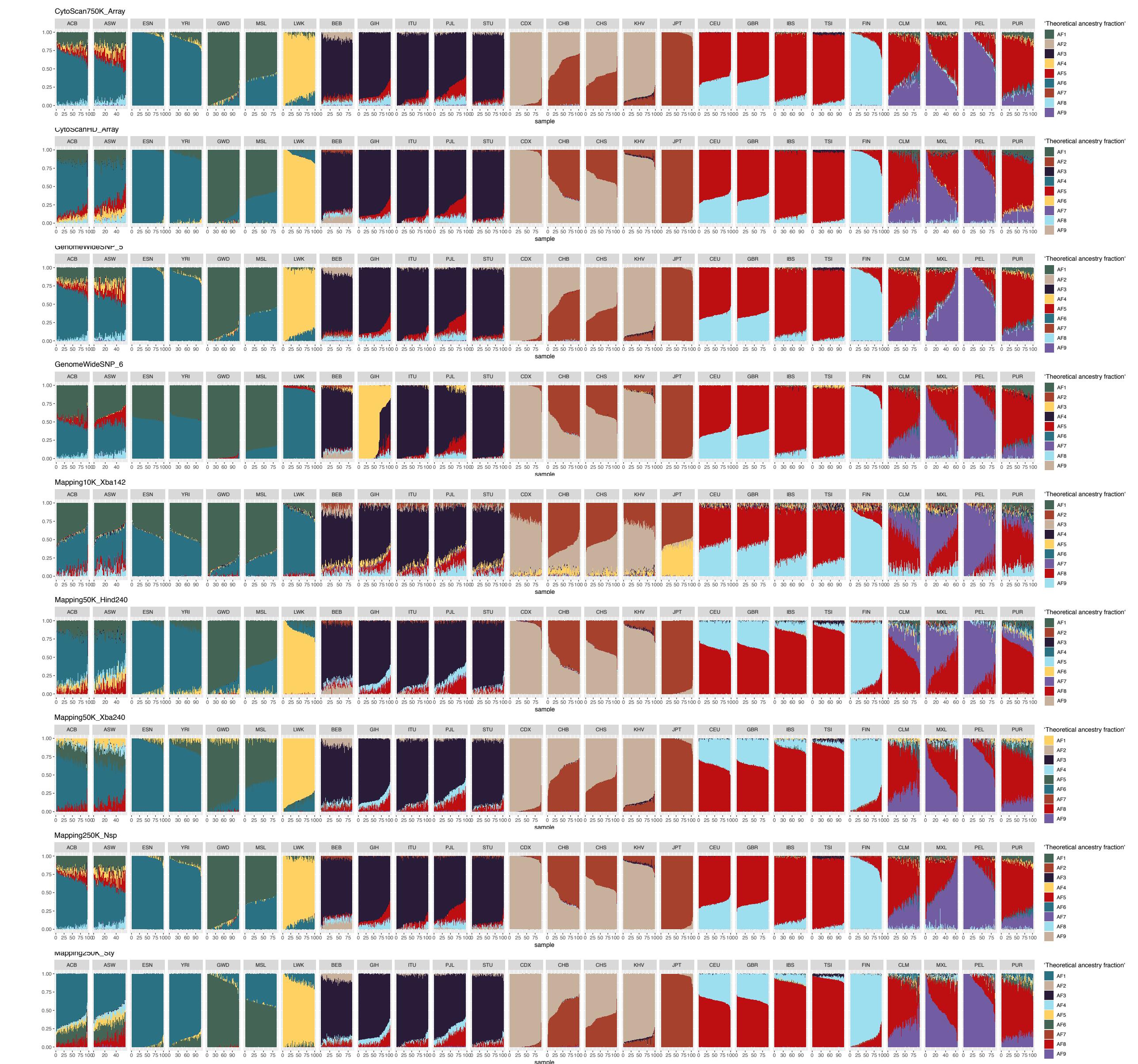
**FIGURE 5 |** A clustering heatmap of features in 31 signatures. Columns are normalized average CNV intensities of feature genes, where the blue colors are duplication features and red colors are deletion features. Duplication and deletion frequencies are normalized separately.



**FIGURE 6 |** The integrated view of the original data and the selected features, in the neural crest originating entities medulloblastoma, glioma, and melanoma. The shaded background area color illustrates the original data. Color bars illustrate the feature genes, where brighter colors indicate stronger signal intensity. The blue colors above the chromosome axis represent the average amplifications, and the red colors below the chromosome axis represent the average deletions. The amplitude of amplifications and deletions are normalized to [0,1] separately. The adjacent known driver genes are also included for each tumor type.

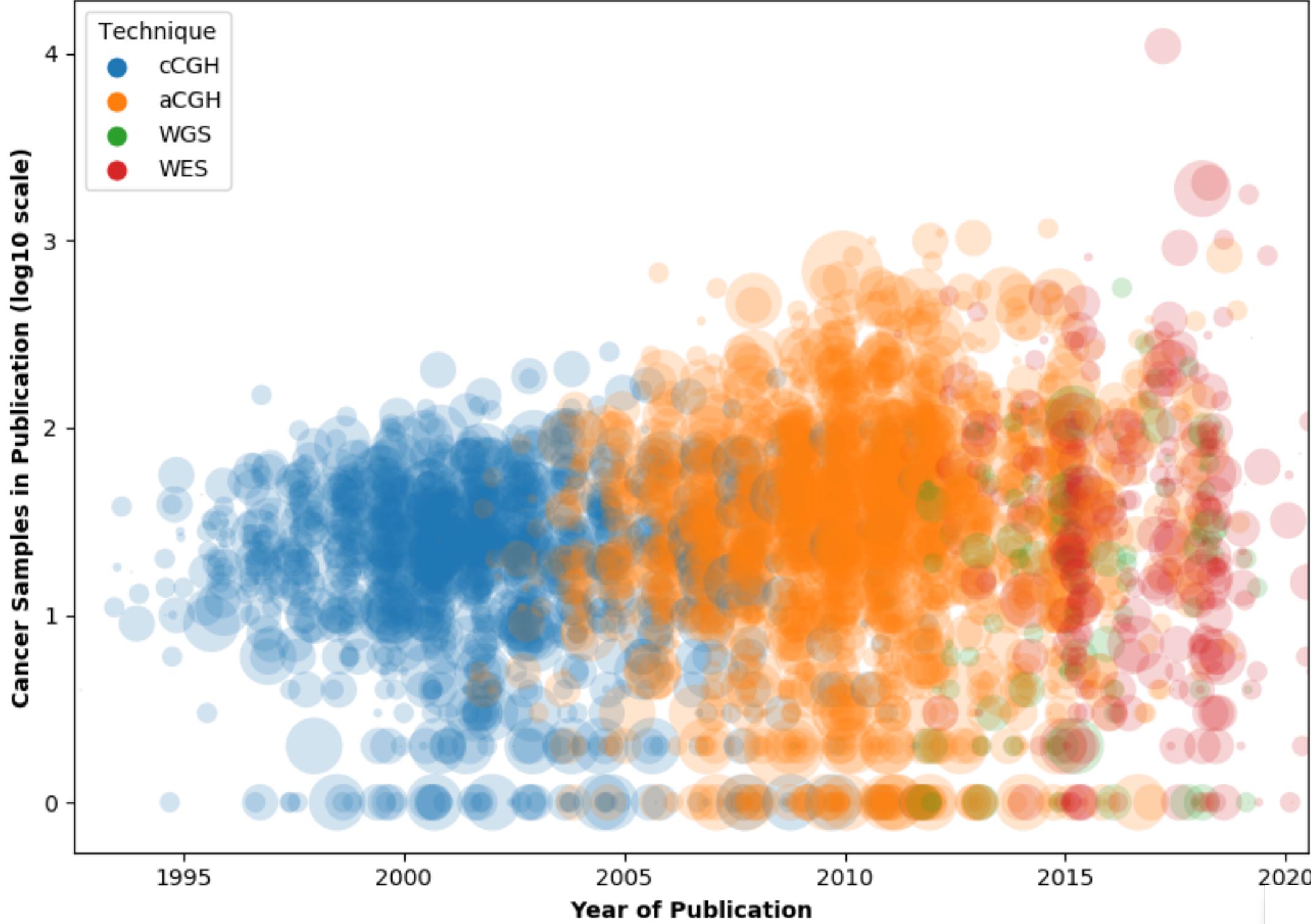
# Population stratification in cancer samples based on SNP array data

- 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



**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).

## Number of tumor samples for each publication across the years



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.



## Progenetix Publication Collection

The current page lists articles describing whole genome screening (WGS, WES, aCGH, cCGH) experiments in cancer, registered in the Progenetix publication collection. For each publication the table indicates the numbers of samples analysed with a given technology and if sample profiles are available in Progenetix.

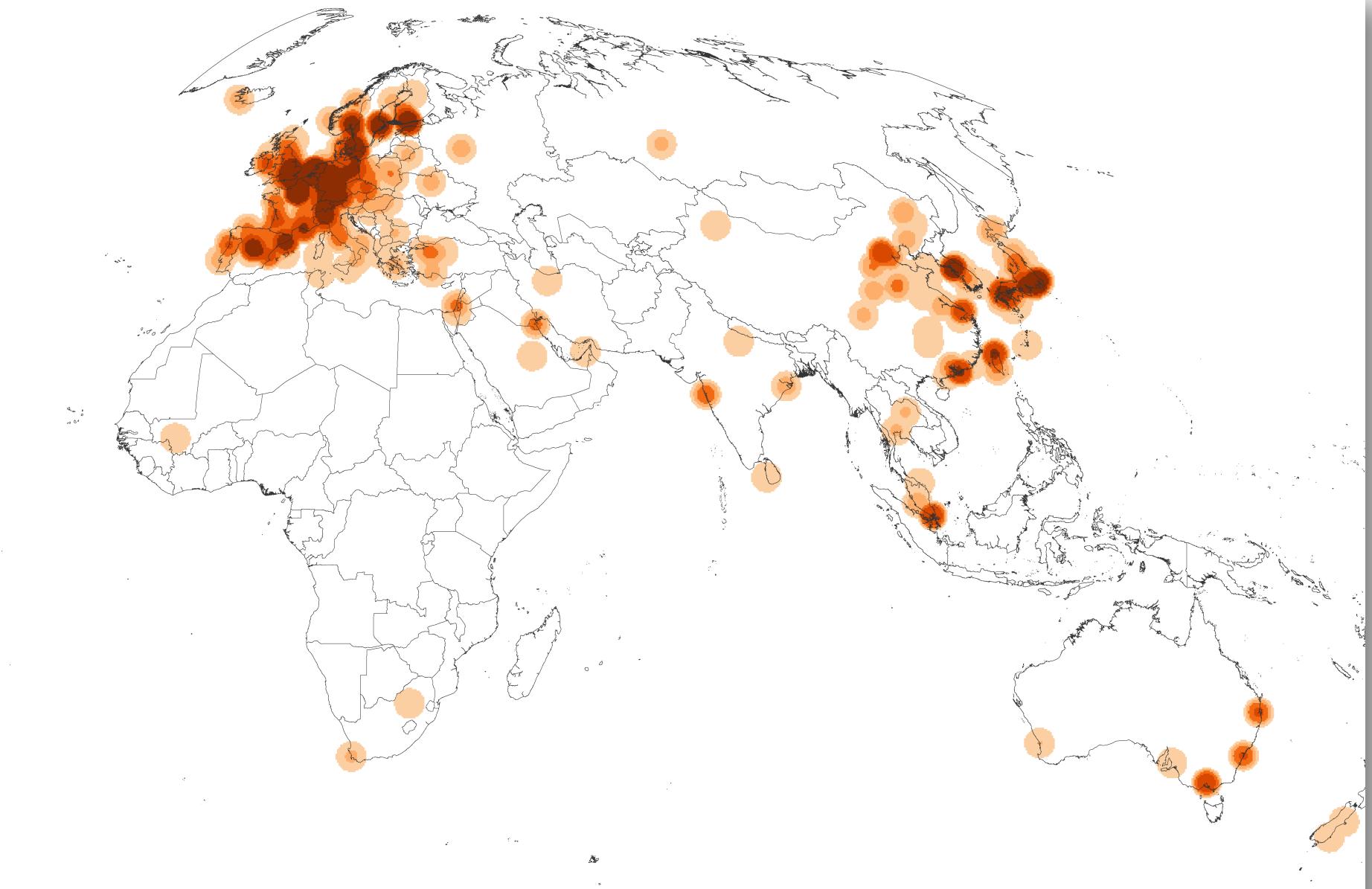
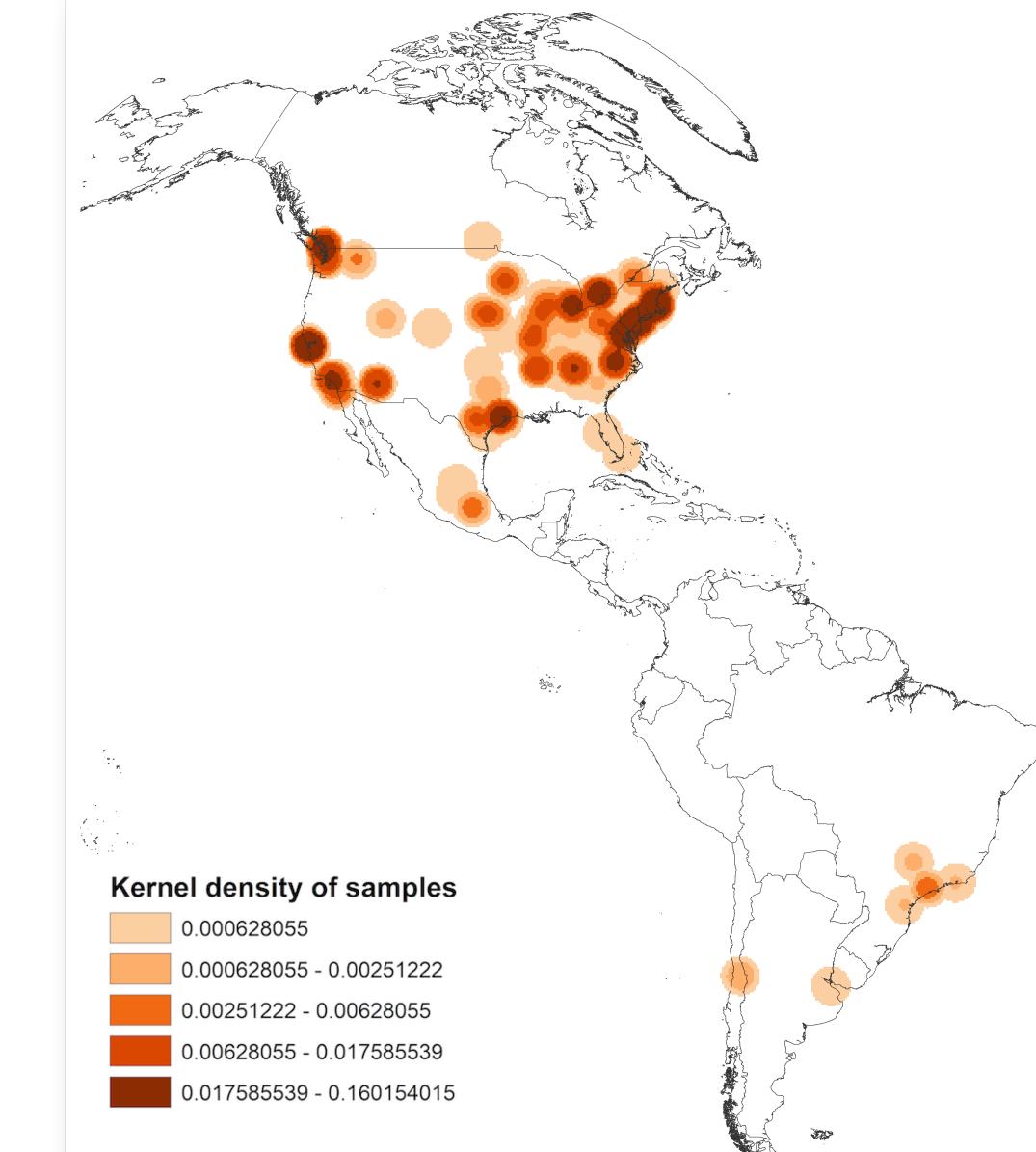
Please [contact us](#) to alert us about additional articles you are aware of. The inclusion criteria are described in the documentation [🔗](#).

Filter [i](#) City [i](#)

  Type to search... [▼](#)

### Publications (3324)

id <a href="#">i</a> ▾	Publication	cCGH	aCGH	WES	WGS	pgx
PMID:34103027	Peng G, Chai H, Ji W, Lu Y, Wu S et al. (2021) Correlating genomic copy number alterations with clinicopathologic findings in 75 cases of ... <i>BMC Med Genomics</i>	0	79	0	0	0
PMID:34059130	Tsui DWY, Cheng ML, Shady M, Yang JL et al. (2021) Tumor fraction-guided cell-free DNA profiling in metastatic solid tumor patients. ... <i>bioRxiv</i>	0	0	5	113	0



# Recent Publications

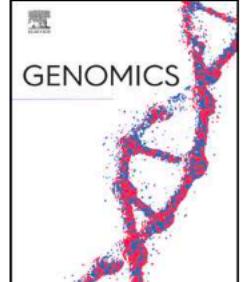
## CNV Data Analysis & Methods

- collaborative projects utilizing the Progenetix data for multi-omics analyses
  - data and bioinformatics analysis support for e.g. translational studies w/o "omics" focus



# **Signatures of Discriminative Copy Number Aberrations in 31 Cancer Subtypes**

Bo Gao<sup>1,2</sup> and Michael Baudis<sup>1,2\*</sup>



Cai et al. BMC Genomics 2012, 13:322

ELSEVIE

RESEARCH A

Bo Gao<sup>a,b</sup>, Michael Baudis<sup>a,b</sup>

# Chromothripsis-like patterns are recurring but heterogeneously distributed features in a survey of 22,347 cancer genome screens

Haoyang Cai<sup>1,2</sup>, Nitin Kumar<sup>1,2</sup>, Homayoun C Bagheri<sup>3</sup>, Christian von Mering<sup>1,2</sup>, Mark D Robinson<sup>1,2</sup>  
and Michael Baudis<sup>1,2\*</sup>

SOFTWARE TOOL ARTICLE

**REVISED** **segment\_liftover : a Python tool to convert segments between genome assemblies [version 2; peer review: 2 approved]**

Bo Gao<sup>1,2</sup>, Qingyao Huang<sup>1,2</sup>, Michael Baudis<sup>1,2</sup>

Ai et al. BMC Genomics (2016) 17:79  
DOI 10.1186/s12864-016-3074-7

ORIGINAL PAPER

# CNARA: reliability assessment for genomic copy number profiles

Ni Ai<sup>1\*</sup>, Haoyang Cai<sup>2</sup>, Caius Solovan<sup>3</sup> and Michael Baudis<sup>1\*</sup>

# Progenetix in 2021

## Cancer Genomics Reference Resource

- largest open resource for curated cancer genome profiling data, with focus on copy number variations (CNV)
- >116'000 cancer CNV profiles, mapped to >800 NCI<sup>t</sup> codes
- majority of data from genomic arrays with ~50% overall from SNP platforms with original data re-processing
- structured diagnostic encodings for NCI<sup>t</sup>, ICD-O 3, UBERON
- identifier mapping for PMID, GEO, Cellosaurus where appropriate
- core biosample and technical metadata annotations where accessible (TNM, genotypic sex, survival ...)
- publication database and code mapping services



[Cancer CNV Profiles](#)

[Search Samples](#)

[Studies & Cohorts](#)

[arrayMap](#)

[TCGA Samples](#)

[DIPG Samples](#)

[Gao & Baudis, 2021](#)

[Cancer Cell Lines](#)

[Publication DB](#)

[Services](#)

[NCIt Mappings](#)

[UBERON Mappings](#)

[Upload & Plot](#)

[Download Data](#)

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

[Progenetix Info](#)

[About Progenetix](#)

[Use Cases](#)

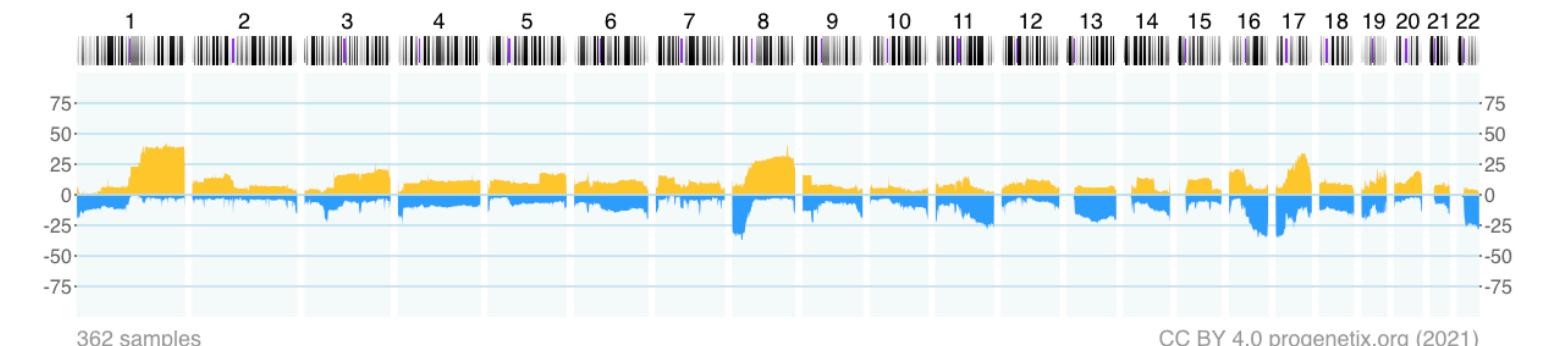
[Documentation](#)

[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* from currently **139448** samples.

Breast Cancer by AJCC v6 Stage (NCIT:C90513)

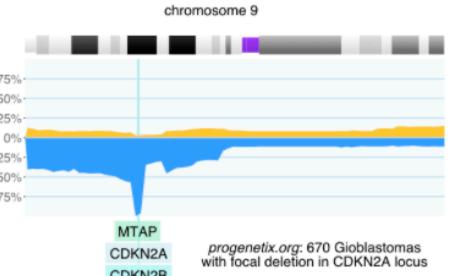


[Download SVG](#) | [Go to NCIT:C90513](#) | [Download CNV Frequencies](#)

Example for aggregated CNV data in 362 samples in Breast Cancer by AJCC v6 Stage. Here the frequency of regional **copy number gains** and **losses** are displayed for all 22 autosomes.

[Progenetix Use Cases](#)

[Local CNV Frequencies](#)



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

[Cancer CNV Profiles](#)

The progenetix resource contains data of **810** different cancer types (NCIt neoplasm classification), mapped to a variety of biological and technical categories. Frequency profiles of regional genomic gains and losses for all categories (diagnostic entity, publication, cohort ...) can be accessed through the [ [Cancer Types](#) ] page with direct visualization and options for sample retrieval and plotting options.

[Cancer Genomics Publications](#)

Through the [ [Publications](#) ] page Progenetix provides **4025** annotated references to research articles from cancer genome screening experiments (WGS, WES, aCGH, cCGH). The numbers of analyzed samples and possible availability in the Progenetix sample collection are indicated.

# Progenetix in 2021

## Cancer Genomics Reference Resource

- largest open resource for curated cancer genome profiling data, with focus on copy number variations (CNV)
- >116'000 cancer CNV profiles, mapped to >800 NCI<sup>t</sup> codes
- majority of data from genomic arrays with ~50% overall from SNP platforms with original data re-processing
- structured diagnostic encodings for NCI<sup>t</sup>, ICD-O 3, UBERON
- identifier mapping for PMID, GEO, Cellosaurus where appropriate
- core biosample and technical metadata annotations where accessible (TNM, genotypic sex, survival ...)
- publication database and code mapping services

progenetix

Search Samples Modify Query

Assembly: GRCh38 Chro: 9 Start: 21500001-21975098 End: 21967753-22500000  
Type: DEL Filters: NCIT:C3058

progenetix

Samples: 668 Variants: 286 Calls: 675

Found Variants (.pgxseg) All Sample Variants (.json) All Sample Variants (.pgxseg) Show Variants in UCSC

UCSC region JSON Response Visualization options

Results Biosamples Biosamples Map Variants

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22

75% 50% 25% 0% -25% -50% -75%

-75% -50% -25% 0% 25% 50% 75%

progenetix: 670 samples

CC BY 4.0 progenetix.org (2021)

Matched Subset Codes	Subset Samples	Matched Samples	Subset Match Frequencies
UBERON:0002021	4	1	0.250
icdot-C71.4	4	1	0.250
icdom-94403	4291	664	0.155
NCIT:C3058	4375	664	0.152
UBERON:0016525	14	2	0.143
icdot-C71.1	14	2	0.143
UBERON:0000955	7068	651	0.092
icdot-C71.9	7066	651	0.092
icdom-94423	84	4	0.048
NCIT:C3796	84	4	0.048
UBERON:0001869	1712	14	0.008
icdot-C71.0	1712	14	0.008

# Progenetix in 2021

## Cancer Genomics Reference Resource

- largest open resource for curated cancer genome profiling data, with focus on copy number variations (CNV)
- >116'000 cancer CNV profiles, mapped to >800 NCI<sup>t</sup> codes
- majority of data from genomic arrays with ~50% overall from SNP platforms with original data re-processing
- structured diagnostic encodings for NCI<sup>t</sup>, ICD-O 3, UBERON
- identifier mapping for PMID, GEO, Cellosaurus where appropriate
- core biosample and technical metadata annotations where accessible (TNM, genotypic sex, survival ...)
- publication database and code mapping services



[Cancer CNV Profiles](#)

[Search Samples](#)

[Studies & Cohorts](#)

[arrayMap](#)

[TCGA Samples](#)

[DIPG Samples](#)

[Gao & Baudis, 2021](#)

[Cancer Cell Lines](#)

[Publication DB](#)

[Genome Profiling](#)

[Progenetix Use](#)

[Services](#)

[NCI<sup>t</sup> Mappings](#)

[UBERON Mappings](#)

[Upload & Plot](#)

[Download Data](#)

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

[Progenetix Info](#)

[About Progenetix](#)

[Use Cases](#)

[Documentation](#)

[Baudisgroup @ UZH](#)

## Progenetix Publication Collection

The current page lists articles describing whole genome screening (WGS, WES, aCGH, cCGH) experiments in cancer, registered in the Progenetix publication collection. For each publication the table indicates the numbers of samples analysed with a given technology and if sample profiles are available in Progenetix.

Please [contact us](#) to alert us about additional articles you are aware of. The inclusion criteria are described in the documentation [↗](#).

**New Oct 2021** You can now directly submit suggestions for matching publications to the [oncopubs](#) repository on [Github ↗](#).

**Filter** [i](#)

**City** [i](#)

 Type to search... | [▼](#)

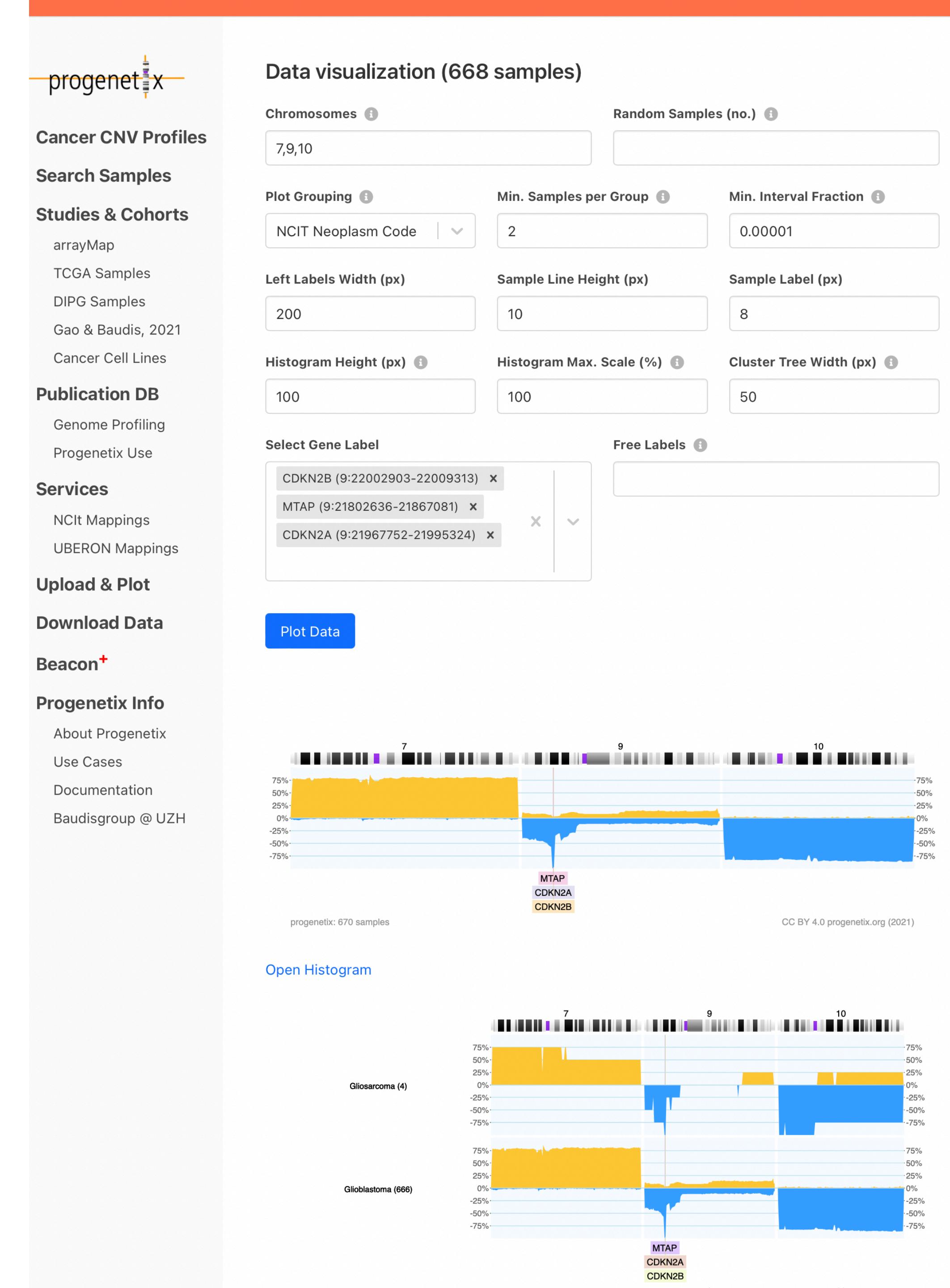
### Publications (3349)

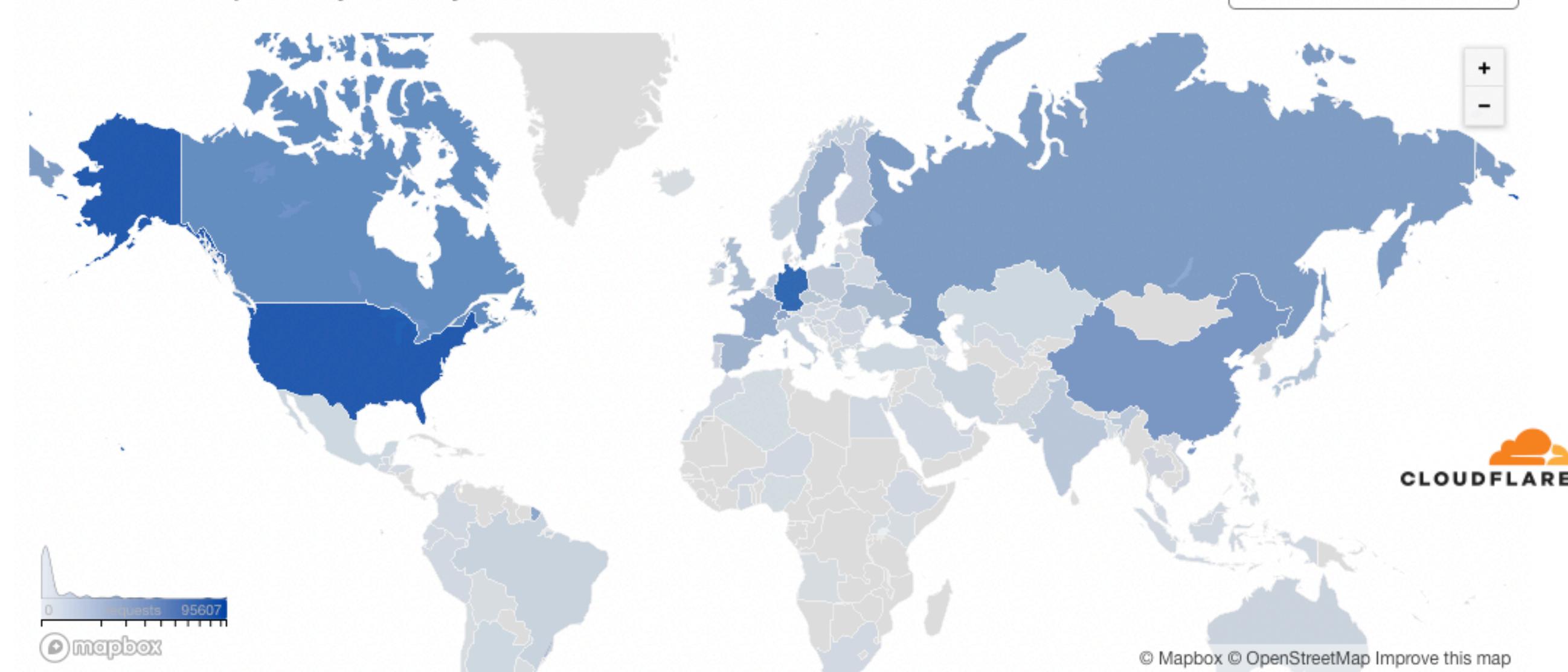
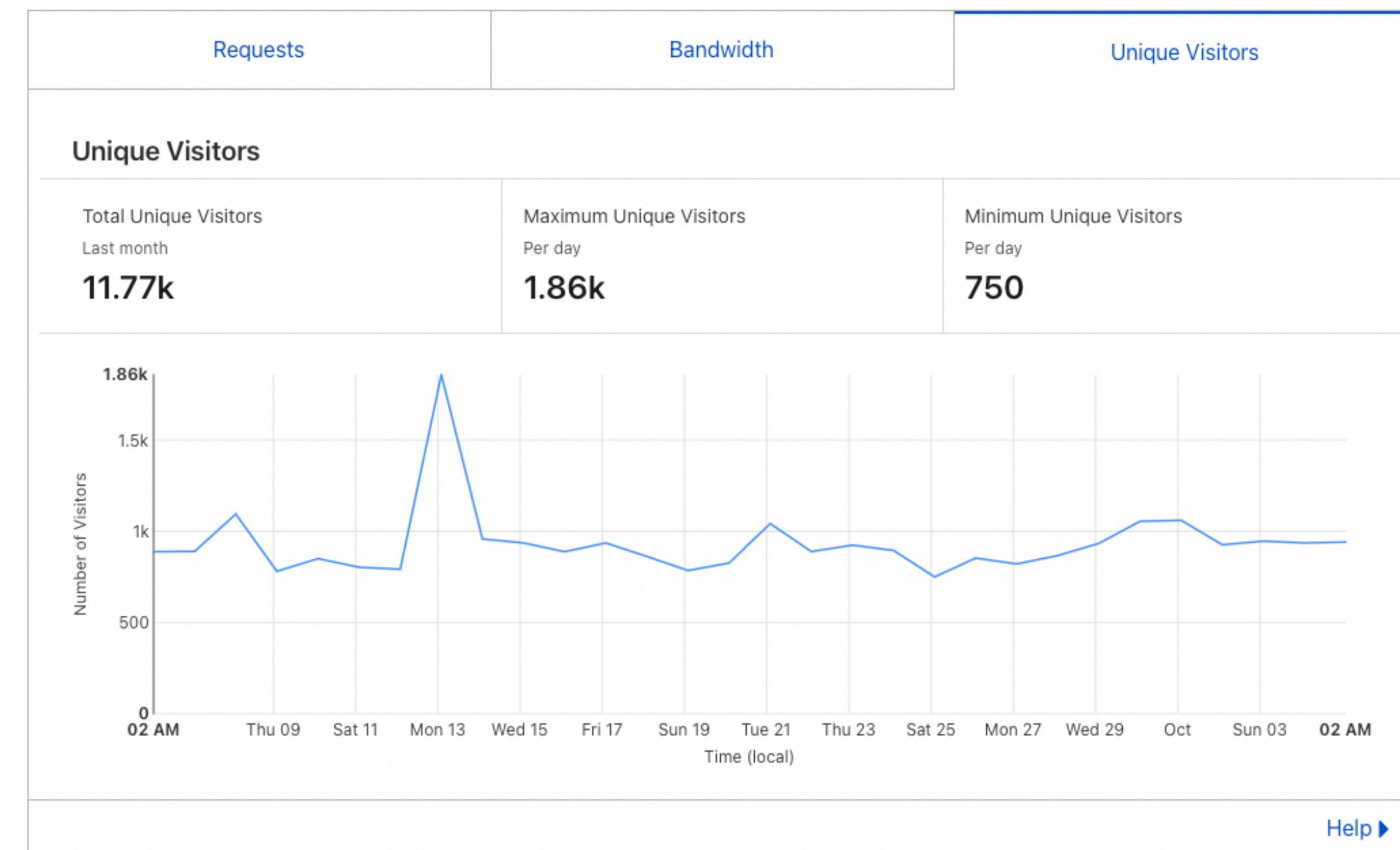
id <a href="#">i</a> ▾	Publication	Samples				
		cCGH	aCGH	WES	WGS	pgx
<a href="#">PMID:34604048</a>	Dai J, Jiang M, He K, Wang H, Chen P et al. (2021) DNA Damage Response and Repair Gene Alterations Increase Tumor Mutational Burden and ... <i>Front Oncol</i>	0	0	122	0	0
<a href="#">PMID:34573430</a>	Juhari WKW, Ahmad Amin Noordin KB et al. (2021) Whole-Genome Profiles of Malay Colorectal Cancer Patients with Intact MMR Proteins. ... <i>Genes (Basel)</i>	0	0	0	7	0
<a href="#">PMID:34307137</a>	Xu S, Li X, Zhang H, Zu L, Yang L et al. (2021) Frequent Genetic Alterations and Their Clinical Significance in Patients With Thymic Epithelial ... <i>Front Oncol</i>	0	0	0	123	0
<a href="#">PMID:34285259</a>	Erkizan HV, Sukhadia S, Natarajan TG et al. (2021) Exome sequencing identifies novel somatic variants in African American esophageal squamous cell ... <i>Sci Rep</i>	0	0	20	0	0
<a href="#">PMID:34205964</a>	Gross C, Engleitner T, Lange S, Weber J et al. (2021) Whole Exome Sequencing of Biliary Tubulopapillary Neoplasms Reveals Common Mutations in Chromatin ... <i>Cancers (Basel)</i>	0	0	17	0	0
<a href="#">PMID:34203905</a>	Chicano M, Carbonell D, Suárez-González J et al. (2021) Next Generation Cytogenetics in Myeloid Hematological Neoplasms: Detection of CNVs and Translocations. ... <i>Cancers (Basel)</i>	0	0	0	135	0
<a href="#">PMID:34103027</a>	Peng G, Chai H, Ji W, Lu Y, Wu S et al. (2021) Correlating genomic copy number alterations with clinicopathologic findings in 75 cases of ... <i>BMC Med Genomics</i>	0	79	0	0	0

# Progenetix in 2021

## Cancer Genomics Reference Resource

- largest open resource for curated cancer genome profiling data, with focus on copy number variations (CNV)
- >116'000 cancer CNV profiles, mapped to >800 NCI codes
- majority of data from genomic arrays with ~50% overall from SNP platforms with original data re-processing
- structured diagnostic encodings for NCI, ICD-O 3, UBERON
- identifier mapping for PMID, GEO, Cellosaurus where appropriate
- core biosample and technical metadata annotations where accessible (TNM, genotypic sex, survival ...)
- publication database and code mapping services





# Progenetix in 2021

## Cancer Genomics Reference Resource

- largest open resource for curated cancer genome profiling data, with focus on copy number variations (CNV)
- >116'000 cancer CNV profiles, mapped to >800 NCI codes
- majority of data from genomic arrays with ~50% overall from SNP platforms with original data re-processing
- structured diagnostic encodings for NCI, ICD-O 3, UBERON
- identifier mapping for PMID, GEO, Cellosaurus where appropriate
- core biosample and technical metadata annotations where accessible (TNM, genotypic sex, survival ...)
- publication database and code mapping services

# The Progenetix oncogenomic resource in 2021

Qingyao Huang<sup>1,2</sup>, Paula Carrio-Cordo<sup>1,2</sup>, Bo Gao<sup>1,2</sup>, Rahel Paloots<sup>1,2</sup> and Michael Baudis<sup>1,2,\*</sup>

<sup>1</sup>Department of Molecular Life Sciences, University of Zurich, Winterthurerstrasse 190, Zurich 8057, Switzerland

<sup>2</sup>Swiss Institute of Bioinformatics, Winterthurerstrasse 190, Zurich 8057, Switzerland

\*Corresponding author: Tel: +41 44 635 34 86; Email: [michael.baudis@mls.uzh.ch](mailto:michael.baudis@mls.uzh.ch)

Citation details: Huang, Q., Carrio-Cordo, P., Gao, B. et al. The Progenetix oncogenomic resource in 2021. *Database* (2021) Vol. 2021: article ID baab043; DOI: <https://doi.org/10.1093/database/baab043>

## Abstract

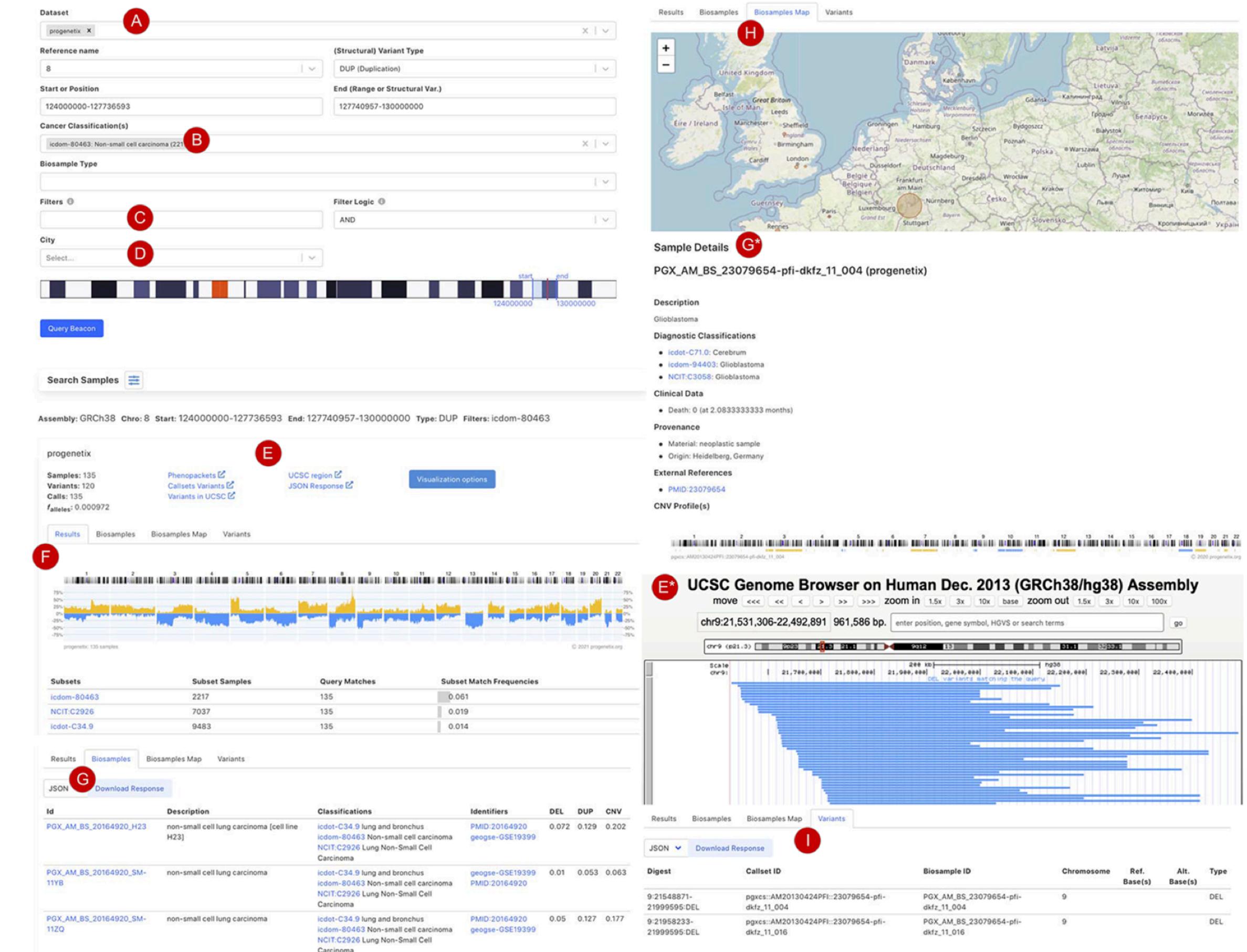
In cancer, copy number aberrations (CNAs) represent a type of nearly ubiquitous and frequently extensive structural genome variations. To disentangle the molecular mechanisms underlying tumorigenesis as well as identify and characterize molecular subtypes, the comparative and meta-analysis of large genomic variant collections can be of immense importance. Over the last decades, cancer genomic profiling projects have resulted in a large amount of somatic genome variation profiles, however segregated in a multitude of individual studies and datasets. The Progenetix project, initiated in 2001, curates individual cancer CNA profiles and associated metadata from published oncogenomic studies and data repositories with the aim to empower integrative analyses spanning all different cancer biologies. During the last few years, the fields of genomics and cancer research have seen significant advancement in terms of molecular genetics technology, disease concepts, data standard harmonization as well as data availability, in an increasingly structured and systematic manner. For the Progenetix resource, continuous data integration, curation and maintenance have resulted in the most comprehensive representation of cancer genome CNA profiling data with 138 663 (including 115 357 tumor) copy number variation (CNV) profiles. In this article, we report a 4.5-fold increase in sample number since 2013, improvements in data quality, ontology representation with a CNV landscape summary over 51 distinctive National Cancer Institute Thesaurus cancer terms as well as updates in database schemas, and data access including new web front-end and programmatic data access.

Database URL: [progenetix.org](http://progenetix.org)

**Table 1.** Statistics of samples from various data resources

Data source	GEO	ArrayExpress	cBioPortal	TCGA	Total
No. of studies	898	51	38	33	1939
No. of samples	63 568	4351	19 712	22 142	138 663
Tumor	52 090	3887	19 712	11 090	115 357
Normal	11 478	464	0	11 052	23 306
Classifications					
ICD-O (Topography)	100	54	88	157	209
ICD-O (Morphology)	246	908	265	140	491
NCIt	346	148	422	182	788
Collections					
Individuals	63 568	4351	19 712	10 995	127 549
Biosamples	63 568	4351	19 712	22 142	138 663
Callsets <sup>a</sup>	63 568	4351	19 712	22 376	138 930
Variants	5 514 126	118 4170	1 778 096	2 654 065	10 716 093

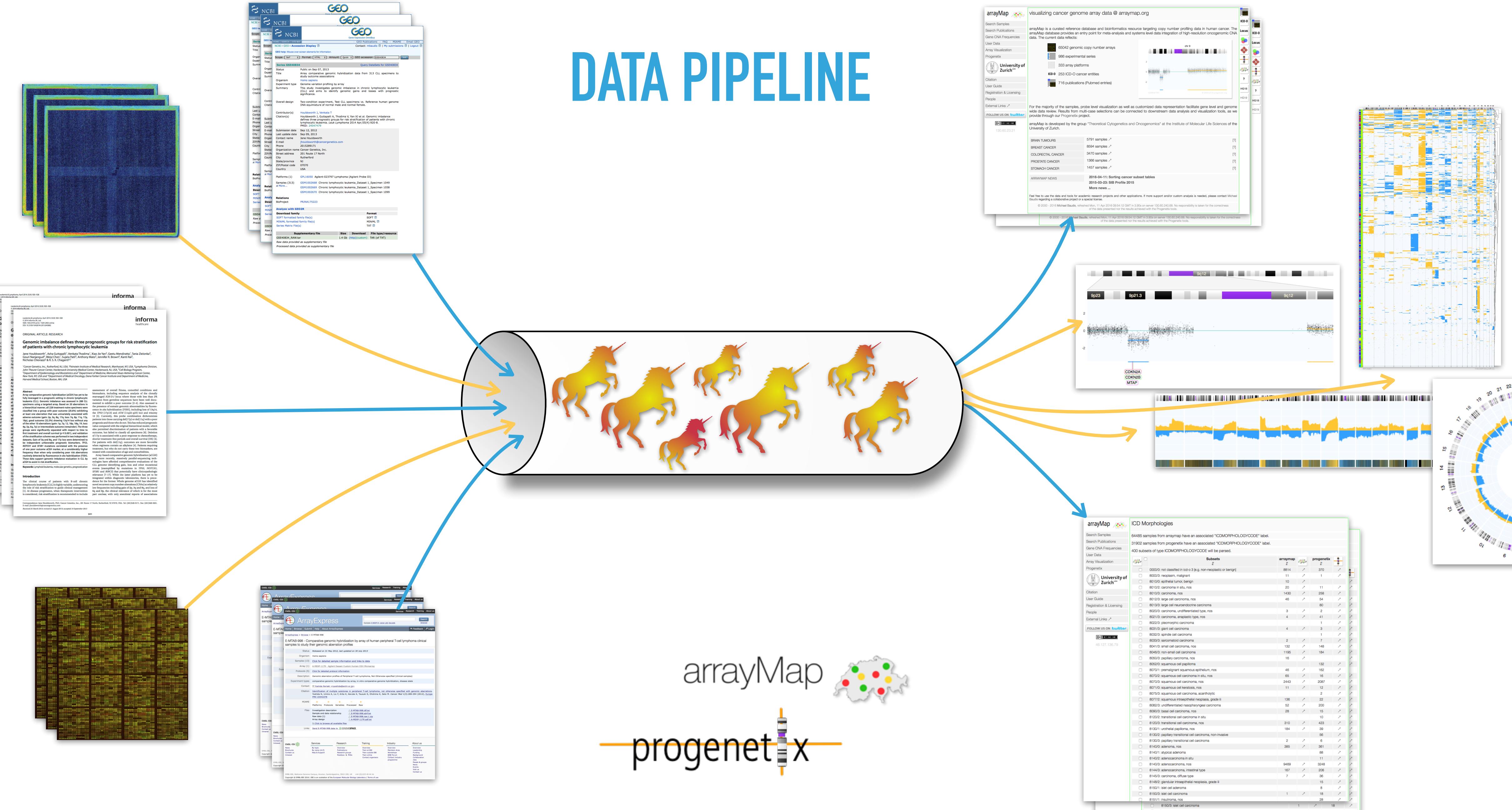
<sup>a</sup>set of variants from one genotyping experiment; ICD-O, International Classification of Diseases for Oncology; NCIt, National Cancer Institute Thesaurus.



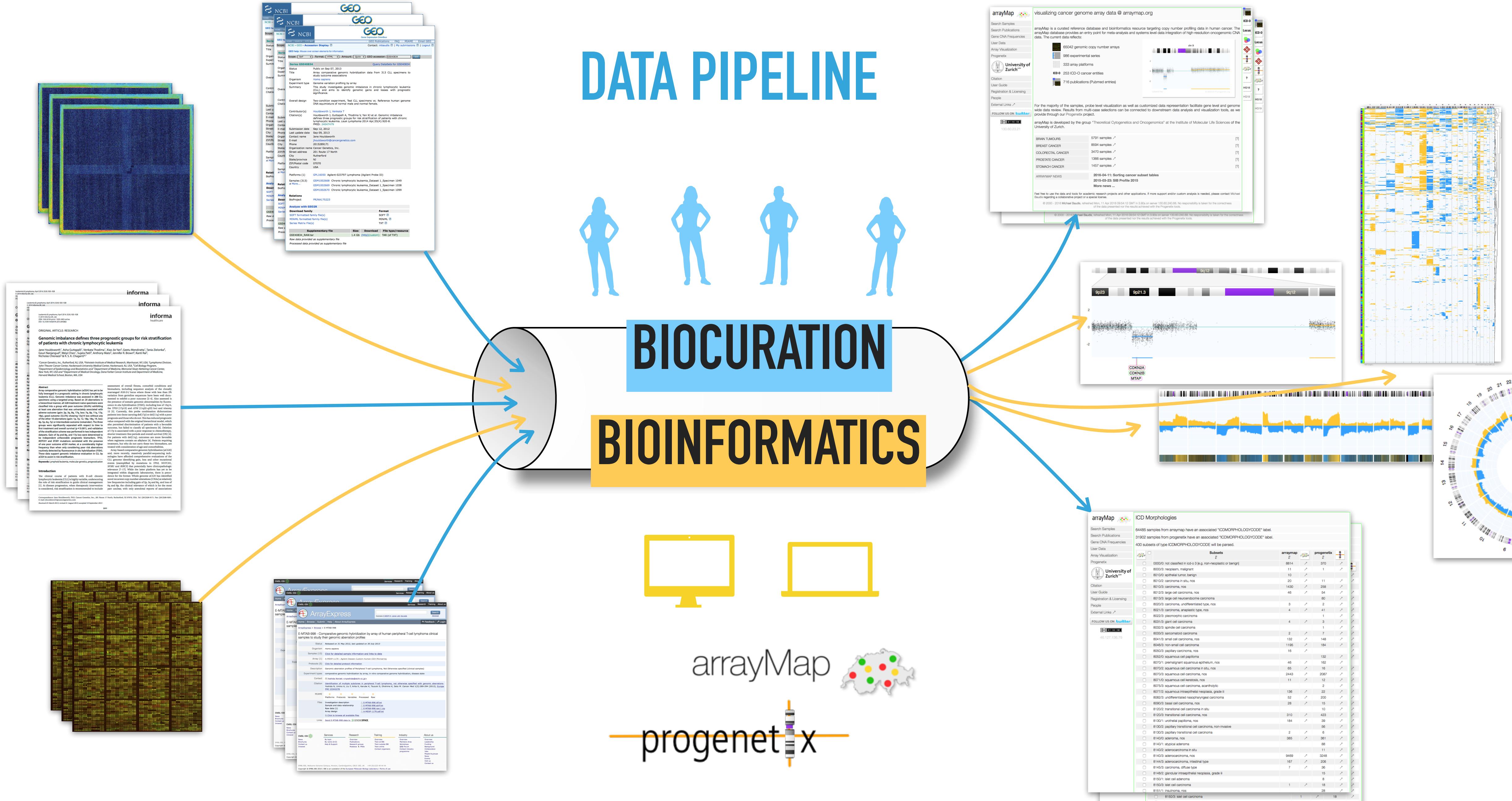
**Figure 3.** Beacon-style query using fuzzy ranges to identify biosamples with variants matching the CNA range. This example queries for a continuous, focal duplication covering the complete MYC gene's coding region with  $\leq 6$  Mb in size. A: Filter for dataset; B: filter for cancer classification (NCIt and ICD-O-3 ontology terms available); C: additional filter, e.g. Cellosaurus; D: additional filter for geographic location; E: external link to UCSC browser to view the alignment of matched variants; F: cancer type classification sorted by frequency of the matched biosamples present in the subset; G: list of matched variants with reference to biosamples can be downloaded in json or csv format.

# Data "Pipelines"

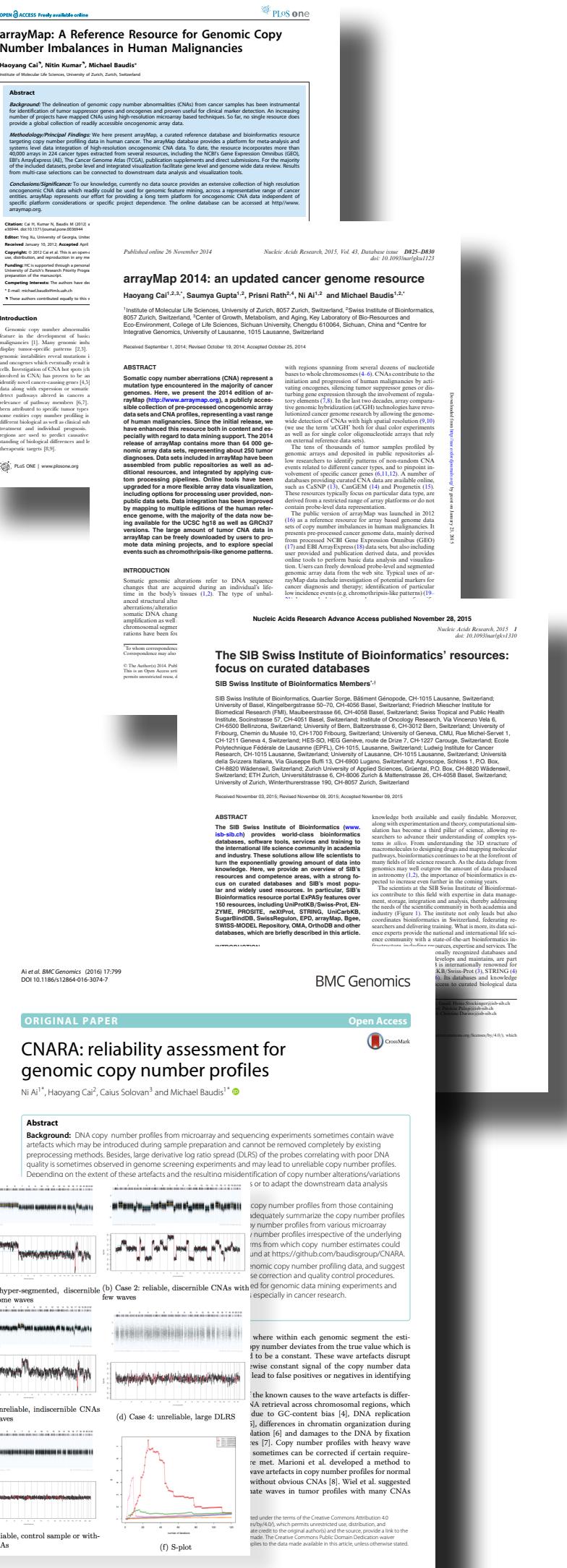
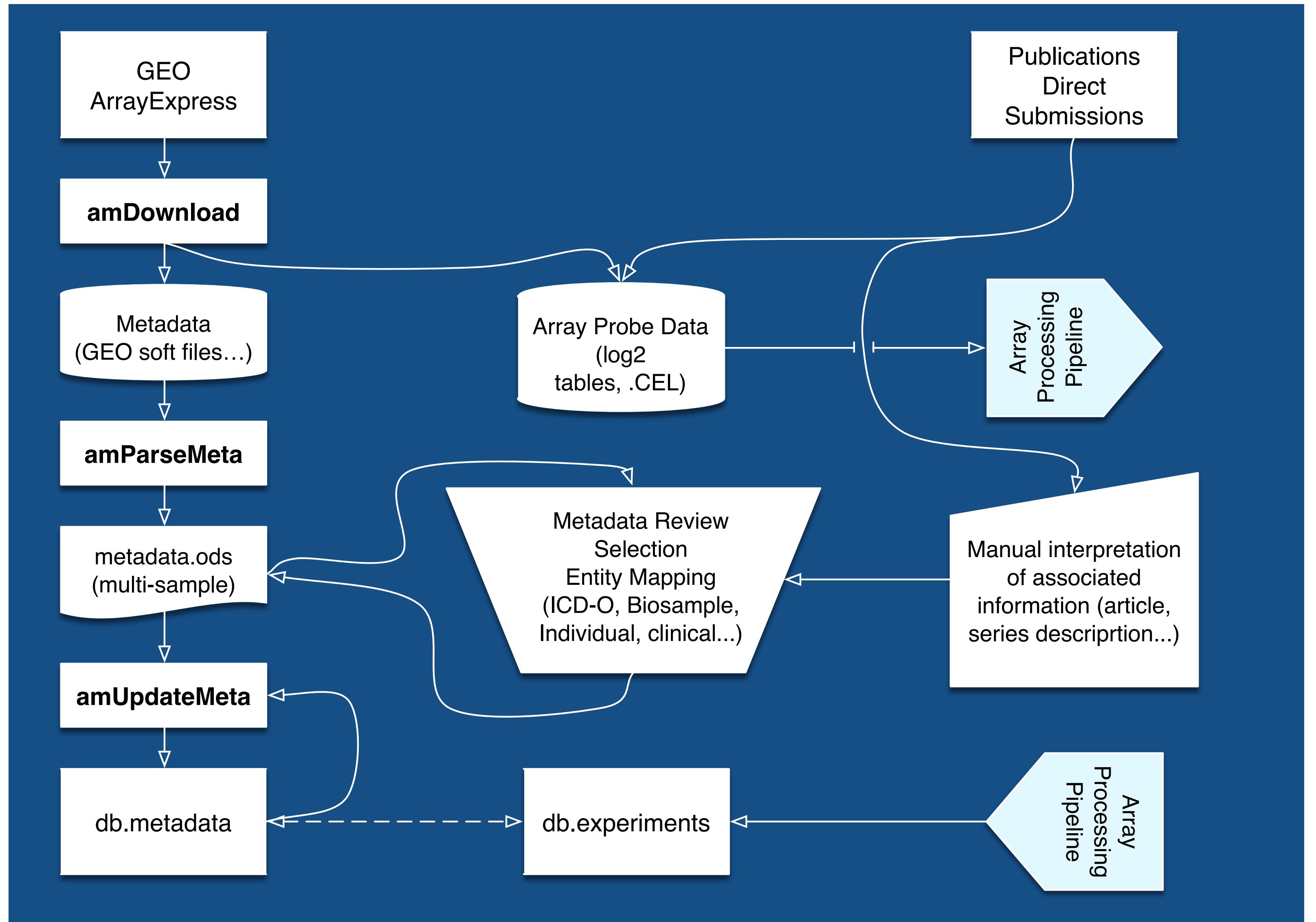
# DATA PIPELINE



# DATA PIPELINE



# Bioinformatics & Data Curation - arrayMap data “Pipeline”



# Progenetix & arrayMap: Data Scopes

## Biomedical and procedural "Meta"data types

- Diagnostic classification
  - mapping text-based cancer diagnoses to standard classification systems
- Provenance data
  - store identifier-based pointers
  - geographic attribution (individual, biosample, experiment)
- Clinical information
  - **core set** of typical cancer study values:
    - ➡ stage, grade, followup time, survival status, genomic sex, age at diagnosis
  - balance between annotation effort and expected usability



# Data Curation - Happy RegExing!

## Extracting clinical and technical metadata from GEO SOFT file

```
^SAMPLE = GSM174832
!Sample_title = 9194
!Sample_geo_accession = GSM174832
!Sample_status = Public on May 01 2007
!Sample_submission_date = Mar 13 2007
!Sample_last_update_date = Mar 13 2007
!Sample_type = genomic
!Sample_channel_count = 1
!Sample_source_name_ch1 = Bone marrow with 96% blasts
!Sample_organism_ch1 = Homo sapiens
!Sample_taxid_ch1 = 9606
!Sample_characteristics_ch1 = Immunotype: common ALL; Age: 9.2 yrs; Gender: F
!Sample_molecule_ch1 = genomic DNA
!Sample_extract_protocol_ch1 = QiaAmp purification kit (Qiagen)
!Sample_label_ch1 = biotin
!Sample_label_protocol_ch1 = Biotinylated DNA was prepared according to the standard
manual 701930 Rev.3 or 100k assay manual 701684 Rev.3, Affymetrix).
!Sample_hyb_protocol = Hybridizations were performed according to the standard
701930 Rev.3 or 100k assay manual 701684 Rev.3, Affymetrix) using an Affymetrix
!Sample_scan_protocol = Scanning performed according to the standard Affymetrix
or 100k assay manual 701684 Rev.3, Affymetrix) using an Affymetrix scanner 3000
!Sample_description = primary ALL diagnosis sample
!Sample_data_processing = copy number detection using CNAG2.0 software (http://
!Sample_platform_id = GPL3718
!Sample_contact_name = Roland,P.,Kuiper
!Sample_contact_email = r.kuiper@antrg.umcn.nl, e.verwiel@antrg.umcn.nl
!Sample_contact_phone = +31243610868
!Sample_contact_fax = +31243668752
!Sample_contact_department = Human Genetics
!Sample_contact_institute = Radboud University Nijmegen Medical Centre
!Sample_contact_address = Geert Grootplein 10
!Sample_contact_city = Nijmegen
!Sample_contact_zip/postal_code = 6525GA
!Sample_contact_country = Netherlands
!Sample_supplementary_file = ftp://ftp.ncbi.nlm.nih.gov/geo/samples/GSM174nnn/GSM174832/suppl/GSM174832.CEL.gz
```

```
foreach (grep { ! /characteristics_ch\d/ } @in) {
    my ($key, $value) = split(' = ', $_);
    $key =~ s/[\w]/_/g;
    if ($key =~ /submission_date/i) {
        $sample->{ YEAR } = $value;
        $sample->{ YEAR } =~ s/^.*?(\d\d\d\d)$/\1/;
    }
}
```

```
$mkey->{ samplekey } = 'AGE';
$mkey->{ matches } = [ qw( age )];

( $mkey->{ retv }, $mkey->{ retk } ) = _grepmeta( $mkey, $meta );

if ( $mkey->{ retv } =~ /^(.+?)$/ ) {
    if ( $mkey->{ retv } =~ /month/i ) {
        $mkey->{ retk } .= '_months';
        $mkey->{ retv } =~ s/[^d\.\.]/ /g;
    }

    $sample->{ $mkey->{ samplekey } } = _normNumber($mkey->{ retv });
    if ( $mkey->{ retk } =~ /month/i ) { $sample->{ $mkey->{ samplekey } } /= 12 }
    if ( $sample->{ $mkey->{ samplekey } } == 0 ) { $sample->{ $mkey->{ samplekey } } = 'NA' }
    $sample->{ $mkey->{ samplekey } } = sprintf "% .2f", $sample->{ $mkey->{ samplekey } };
}
```

# Data Curation - Happy RegExing!

## Extracting clinical and technical metadata from GEO SOFT file

```
^SAMPLE = GSM286922
!Sample_title = 481 - mAbID:75320
!Sample_geo_accession = GSM286922
!Sample_status = Public on Sep 04 2008
!Sample_submission_date = May 06 2008
!Sample_last_update_date = Nov 26 2008
!Sample_type = genomic
!Sample_channel_count = 2
!Sample_source_name_ch1 = Normal Lymphocytes
!Sample_organism_ch1 = Homo sapiens
!Sample_taxid_ch1 = 9606
!Sample_characteristics_ch1 = Tissue: lymphocytes
!Sample_molecule_ch1 = genomic DNA
!Sample_extract_protocol_ch1 = Sample DNA Extraction Protocol
!Sample_extract_protocol_ch1 = Other: The DNA was isolated by Qiagen DNe
!Sample_label_ch1 = cy5
!Sample_label_protocol_ch1 = NimbleGen Cy5 Sample Labeling Protocol
!Sample_label_protocol_ch1 = Other: Proprietary protocol information available at http://www.nimblegen.com/technology/index.html
!Sample_source_name_ch2 = 481
!Sample_organism_ch2 = Homo sapiens
!Sample_taxid_ch2 = 9606
!Sample_characteristics_ch2 = Gender: male
!Sample_characteristics_ch2 = Age: 49
!Sample_characteristics_ch2 = Tissue: lymph node
!Sample_characteristics_ch2 = Disease state: Lymphoma
!Sample_characteristics_ch2 = Individual: 481
!Sample_characteristics_ch2 = Clinical info: Submitting diagnosis: DLBCL
!Sample_characteristics_ch2 = Clinical info: Final microarray diagnosis: ABC DLBCL
!Sample_characteristics_ch2 = Clinical info: Follow up status: ALIVE
!Sample_characteristics_ch2 = Clinical info: Follow up years: 10.75
!Sample_characteristics_ch2 = Clinical info: Chemotherapy: CHOP-Like Regimen
!Sample_characteristics_ch2 = Clinical info: ECOG performance status: 2
!Sample_characteristics_ch2 = Clinical info: Stage: 4
!Sample_characteristics_ch2 = Clinical info: LDH ratio: 0.82
!Sample_characteristics_ch2 = Clinical info: Number of extranodal sites: 1
```

Channel 1 is normal -> Cave value swap!

Gender or "chromosomal sex"?

context indicates years, but if it would be a medulloblastoma...

Not yet registered way to express "alive"!

```
$mkey->{ samplekey } = 'DEATH';
$mkey->{ matches } = [
    'death',
    'dead ',
    'vital_status',
    'dead_alive',
    'alive_dead',
];
(
    $mkey->{ retv }, $mkey->{ retk } ) = _grepmeta( $mkey, $meta );
if ( $mkey->{ retv } =~ /^(.+?)$/ ) {
    $sample->{ $mkey->{ samplekey } } = _normDeath($mkey->{ retv }) }
```

# Data Curation

## Happy RegExing!



```
19 extraction_scopes:
20   description: >-
21     Detection and processing of clinical scopes goes through several stages:
22     1. line cleanup - so far run for the input before processing the individual
23       scopes
24     2. line match using some general pattern expected in all lines containing
25       data for the current scope (`filter` pattern)
26     3. finding and extracting the relevant data by looping over a list of
27       specific patterns with memorized matches (`find`)
28     4. post-processing using empirical cleanup replacements (`cleanup`)
29     5. checking the correct structure (`final_check` - a global pattern can be
30       used if other post-processing is performed)
31
32
33 survival_status:
34   filter: '(?i).*?(?:(:deaf?:d|th))|alive|surviv|outcome|status'
35   preclean:
36     - m: '(?i)days to death or last seen alive[^w]+?\d+?(?:[^w\.]|$)'
37     s: ''
38     - m: '[^w]+?NA(?:[^w\.]|$)'
39     s: ''
40     - m: 'remission status past double induction .cr. complete remission. RD. refractory disease. ED. early death[^w]+?ED'
41     s: 'survival: dead'
42     - m: 'remission status past double induction .cr. complete remission. RD. refractory disease. ED. early death[^w]+?NA'
43     s: ''
44     - m: 'remission status past double induction .cr. complete remission. RD. refractory disease. ED. early death[^w]+?CR'
45     s: 'survival: alive'
46     - m: 'remission status past double induction .cr. complete remission. RD. refractory disease. ED. early death[^w]+?RD'
47     s: '# alive but not responding to therapy so removed?'
48     - m: 'Event Free Survival[^w]+?no event'
49     s: 'recurrence: no'
50     - m: 'Event Free Survival.event'
51     s: 'recurrence: yes'
52     - m: 'Outcome[^w]+?no event'
53     s: 'survival: alive'
54     - m: 'Outcome[^w]+?event'
55     s: 'survival: dead'
56     - m: 'survival status[^w]+?0'
57     s: 'survival: dead'
58     - m: 'survival status[^w]+?1'
59     s: 'survival: alive'
60     - m: 'overall[^w]+?survival[^w]+?days[^w]+?NA'
61     s: ''
62     - m: 'survival(?: time|from diagnosis)?[^w]+?(days|months|years?)[^w]+?(\\d\\d?\\d?\\d?\\.?\\d?\\d?)'
63     s: 'survival: \\2\\1'
```

# Disease annotations in Progenetix

## From some text, somewhere, to ontology classes

- **diagnostic categories** are the **most important** labels to associate with genomic observations
- original data almost *never* uses **modern, hierarchical** classification systems but provides circumstantial ("breast cancer in pre-menopausal...") or domain-specific ("CLL Binet B", "colorectal carcinoma Dukes C") information
- clinical classifications (ICD-10 ...) have very limited relation to tumor biology
- concepts change over time ...
- for cancer, the "International Classification of Diseases in Oncology" (**ICD-O 3**) by IARC / WHO traditionally has been a good compromise to map to - but with non-hierarchical structure and is used by international reference projects

# From Classification to Hierarchical Ontology: ICD-O -> NCI

example_dx	ICDMORPHOLOGY	ICDOM	ICDTOPOGRAPHY	ICDOT	NCIT:CODE
malignant melanoma [metastatic cell line MaMel19]	Malignant melanoma NOS	8720/3	skin	C44	C3224
malignant melanoma [vagina]	Malignant melanoma NOS	8720/3	vagina and labia	C510	C3224
malignant melanoma [uvea metastasized]	Malignant melanoma NOS	8720/3	retina	C692	C3224
meningioma	Meningioma NOS	9530/0	meninges cerebral spinal	C700	C3230
mesothelioma	Mesothelioma NOS	9050/3	lung and bronchus	C34	C3234
pleural mesothelioma	Mesothelioma NOS	9050/3	pleura	C384	C3234
mesothelioma	Mesothelioma NOS	9050/3	connective and soft tissue NOS	C499	C3234
multiple myeloma	Plasma cell myeloma	9732/3	hematopoietic and reticuloendothelial system	C42	C3242
Mycosis fungoides	Mycosis fungoides	9700/3	skin	C44	C3246
Myelodysplastic syndrome	Myelodysplastic syndrome NOS	9989/3	hematopoietic and reticuloendothelial system	C42	C3247
Acute myeloblastic leukemia with maturation [FAB M2]	Acute myeloblastic leukemia with maturation [FAB M2]	9874/3	hematopoietic and reticuloendothelial system	C42	C3250
neuroblastoma	Neuroblastoma NOS	9500/3	peripheral nerves incl. autonomous	C47	C3270
Cerebral neuroblastoma [cerebral region midline frontal lobe]	Neuroblastoma NOS	9500/3	cerebrum	C710	C3270
neuroblastoma [adrenal gland cell line]	Neuroblastoma NOS	9500/3	adrenal gland	C76	C3270
Cutaneous neurofibroma	Neurofibroma NOS	9540/0	skin	C44	C3272
Plexiform neurofibroma	Neurofibroma NOS	9540/0	Nervous system NOS	C729	C3272
Oligodendrogioma [Supratentorial Frontal Lobe]	Oligodendrogioma NOS	9450/3	cerebrum	C710	C3288
oligodendrogioma	Oligodendrogioma NOS	9450/3	Brain NOS	C719	C3288
oligodendrogioma	Oligodendrogioma NOS	9450/3	brain nos	c719	C3288
Paraganglioma	Paraganglioma NOS	8680/1	Nervous system NOS	C729	C3308
paraganglioma	paraganglioma NOS	8680/1	adrenal cortex	C740	C3308

- since its beginning Progenetix samples have been classified using the 2 arms of the ICD-O system (morphology ~ histology/biology + topography ~ organ/tissue)
- over the last years we have established mappings between ICD-O code pairs and the NCIt "neoplasm" part of the NCI metathesaurus, thereby empowering hierarchical data structures for search and analysis

# Ontologies and Classifications



## Services: Ontologymaps (NCIt)

The **ontologymaps** service provides equivalency mapping between ICD-O and other classification systems, notably NCIt and UBERON. It makes use of the sample-level mappings for NCIT and ICD-O 3 codes developed for the individual samples in the Progenetix collection.

### NCIT and ICD-O 3

While NCIT treats diseases as **histologic** and **topographic** described entities (e.g. [NCIT:C7700: Ovarian adenocarcinoma](#)), these two components are represented separately in ICD-O, through the **Morphology** and **Topography** coding arms (e.g. here [8140/3 + C56.9](#)).

More documentation with focus on the API functionality can be found on the [documentation pages](#).

The data of all mappings can be retrieved trough this API call: [{JSON ↗}](#)

### Code Selection ⓘ

NCIT:C4337: Mantle Cell Lymphoma X | ▾

Optional: Limit with second selection | ▾

### Matching Code Mappings [{JSON ↗}](#)

NCIT:C4337: Mantle Cell Lymphoma	<a href="#">pgx:icdom-96733: Mantle cell lymphoma</a>	<a href="#">pgx:icdot-C77.9: Lymph nodes, NOS</a>
NCIT:C4337: Mantle Cell Lymphoma	<a href="#">pgx:icdom-96733: Mantle cell lymphoma</a>	<a href="#">pgx:icdot-C18.9: large intestine, excl. rectum and rectosigmoid junction</a>
NCIT:C4337: Mantle Cell Lymphoma	<a href="#">pgx:icdom-96733: Mantle cell lymphoma</a>	<a href="#">pgx:icdot-C42.2: Spleen</a>

More than one code groups means that either mappings need refinements (e.g. additional specific NCIT classes for ICD-O T topographies) or you started out with an unspecific ICD-O M class and need to add a second selection.

In Progenetix all cancer diagnoses are coded to both NCIt neoplasm codes and ICD-O 3 Morphology + Topography combinations. The matched mappings are provided as lookup-service since neither an official ICD-O ontology nor such a "disease defined by ICD-O M+T" concept is codified anywhere.

List of filters recognized by different query endpoints

Public Ontologies with CURIE-based syntax

CURIE prefix	Code/Ontology	Examples
NCIT	NCIt Neoplasm <sup>1</sup>	NCIT:C27676
HP	HPO <sup>2</sup>	HP:0012209
PMID	NCBI Pubmed ID	PMID:18810378
geo	NCBI Gene Expression Omnibus <sup>3</sup>	geo:GPL6801, geo:GSE19399, geo:GSM491153
arrayexpress	EBI ArrayExpress <sup>4</sup>	arrayexpress:E-MEXP-1008
cellosaurus	Cellosaurus - a knowledge resource on cell lines <sup>5</sup>	cellosaurus:CVCL_1650
UBERON	Uberon Anatomical Ontology <sup>6</sup>	UBERON:0000992
cBioPortal	cBioPortal <sup>9</sup>	cBioPortal:msk_impact_2017

### Private filters

Since some classifications cannot directly be referenced, and in accordance with the upcoming Beacon v2 concept of "private filters", Progenetix uses additionally a set of structured non-CURIE identifiers.

For terms with a `pgx` prefix, the [identifiers.org resolver](#) will

Filter prefix / local part	Code/Ontology	Example
pgx:icdom...	ICD-O 3 <sup>7</sup> Morphologies (Progenetix)	pgx:icdom-81703
pgx:icdot...	ICD-O 3 <sup>7</sup> Topographies(Progenetix)	pgx:icdot-C04.9
TCGA	The Cancer Genome Atlas (Progenetix) <sup>8</sup>	TCGA-000002fc-53a0-420e-b2aa-a40a358bba37
pgx:pgxcohort...	Progenetix cohorts <sup>10</sup>	pgx:pgxcohort-arraymap

# 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

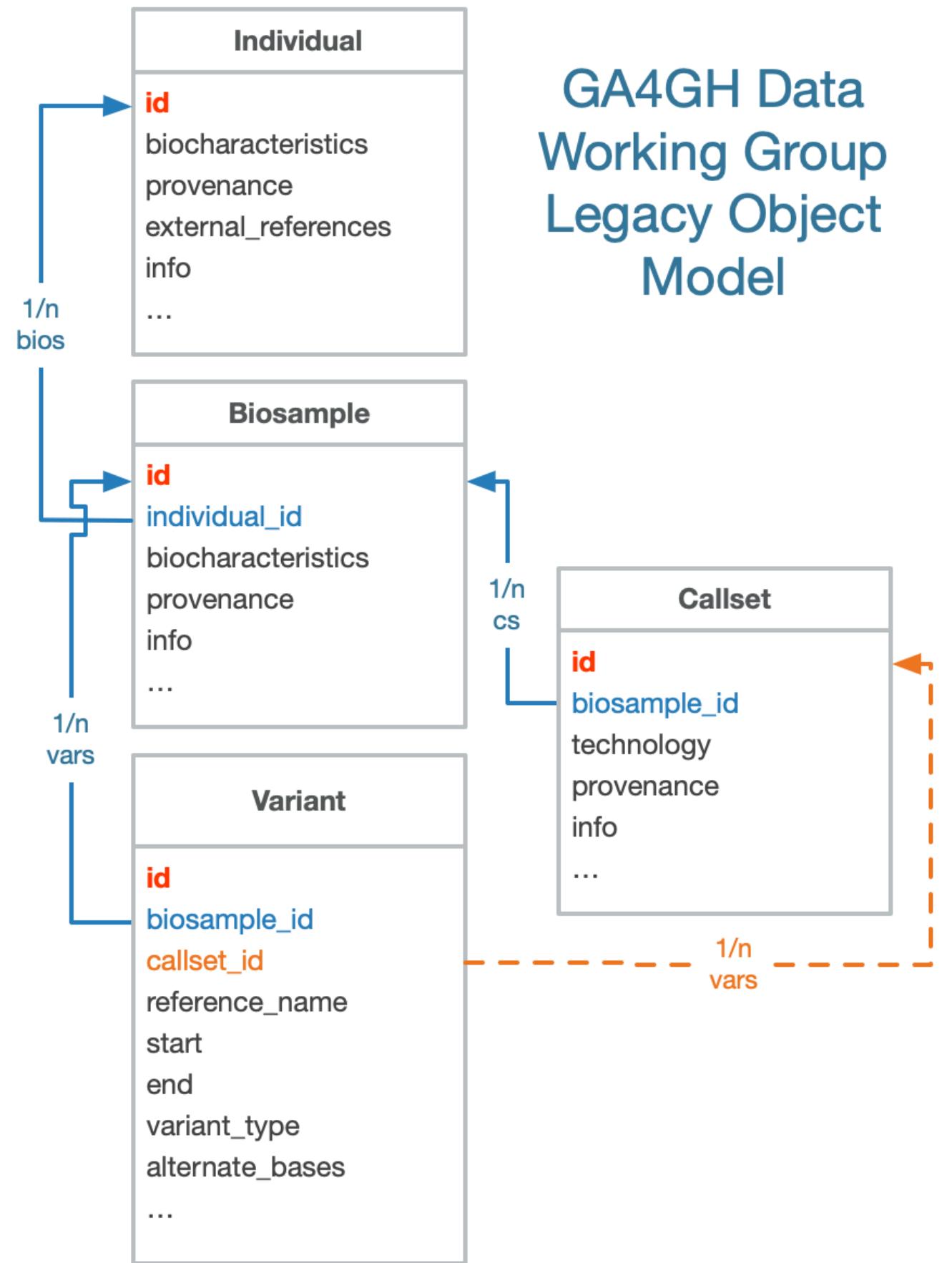


NCIt Neoplasm Core coded display (excerpt) for samples in the Progenetix cancer genome data resource allows sample selection on multiple hierarchy levels →

	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

# Database Structure

## From flat database to hierarchical object storage



- collections in Progenetix MongoDB database reflect a consensus domain model for genomic data repositories
- flexible linking and object structure facilitates rapid change-overs
- BSON/JSON format in DB

- equals data in JavaScript
- "equals" objects in Python, Perl

→ **rapid prototyping and implementation**

2021

- non-SQL document database (MongoDB)
- different object domains connected through identifiers
- data-driven website with JavaScript based frontend and data population through API calls

```
{  
  "id" : "pgxind-kftx394x",  
  "biocharacteristics" : [  
    {  
      "description" : "female",  
      "type" : {  
        "id" : "PATO:0020002",  
        "label" : "female genotypic sex"  
      }  
    },  
    {  
      "description" : null,  
      "type" : {  
        "id" : "NCBITaxon:9606",  
        "label" : "Homo sapiens"  
      }  
    }  
,  
  "data_use_conditions" : {  
    "label" : "no restriction",  
    "id" : "DUO:0000004"  
  },  
  "geo_provenance" : {  
    "label" : "Salamanca, Spain",  
    "precision" : "city",  
    "city" : "Salamanca",  
    "country" : "Spain",  
    "latitude" : 40.43,  
    "longitude" : -3.68  
  },  
  "info" : {  
    "legacy_id" : "PGX_IND_SMZL01"  
  },  
  "updated" : ISODate("2018-09-26T09:51:39.775Z")  
}  
  
  "digest" : "7:107200000-158821424:DEL",  
  "reference_name" : "7",  
  "variant_type" : "DEL",  
  "start" : 107200000,  
  "end" : 158821424,  
  "info" : {  
    "cnv_value" : null,  
    "cnv_length" : 51621424  
  },  
  "updated" : "2018-09-26 09:51:39.775397"  
}
```

```
{  
  "type" : {  
    "id" : "UBERON:0002106",  
    "label" : "spleen"  
  }  
,  
  {  
    "type" : {  
      "id" : "icdot-C42.2",  
      "label" : "Spleen"  
    }  
,  
    {  
      "type" : {  
        "id" : "icdom-96893",  
        "label" : "Splenic marginal zone B-cell lymphoma"  
      }  
,  
      {  
        "type" : {  
          "id" : "NCIT:C4663",  
          "label" : "Splenic Marginal Zone Lymphoma"  
        }  
      }  
,  
      "individual_id" : "pgxind-kftx394x",  
      "individual_age_at_collection" : "P67Y",  
      "info" : {  
        "death" : "0",  
        "followup_months" : 53,  
        "callset_ids" : [  
          "pgxcs-kftvv618"  
        ],  
        "legacy_id" : "PGX_AM_BS_SMZL01"  
      },  
      "external_references" : [  
        {  
          "type" : {  
            "id" : "PMID:11337382"  
          }  
        }  
      ],  
      "provenance" : {  
        "material" : {  
          "type" : {  
            "id" : "EFO:0009656",  
            "label" : "neoplastic sample"  
          }  
        }  
      },  
      "geo" : {  
        "label" : "Salamanca, Spain",  
        "precision" : "city",  
        "city" : "Salamanca",  
        "country" : "Spain",  
        "geojson" : {  
          "type" : "Point",  
          "coordinates" : [  
            -3.68,  
            40.43  
          ]  
        },  
        "ISO-3166-alpha3" : "ESP"  
      }  
    }  
},  
{  
  "type" : {  
    "id" : "UBERON:0002106",  
    "label" : "spleen"  
  }  
,  
  {  
    "type" : {  
      "id" : "icdot-C42.2",  
      "label" : "Spleen"  
    }  
,  
    {  
      "type" : {  
        "id" : "icdom-96893",  
        "label" : "Splenic marginal zone B-cell lymphoma"  
      }  
,  
      {  
        "type" : {  
          "id" : "NCIT:C4663",  
          "label" : "Splenic Marginal Zone Lymphoma"  
        }  
      }  
,  
      "individual_id" : "pgxind-kftx394x",  
      "individual_age_at_collection" : "P67Y",  
      "info" : {  
        "death" : "0",  
        "followup_months" : 53,  
        "callset_ids" : [  
          "pgxcs-kftvv618"  
        ],  
        "legacy_id" : "PGX_AM_BS_SMZL01"  
      },  
      "external_references" : [  
        {  
          "type" : {  
            "id" : "PMID:11337382"  
          }  
        }  
      ],  
      "provenance" : {  
        "material" : {  
          "type" : {  
            "id" : "EFO:0009656",  
            "label" : "neoplastic sample"  
          }  
        }  
      },  
      "geo" : {  
        "label" : "Salamanca, Spain",  
        "precision" : "city",  
        "city" : "Salamanca",  
        "country" : "Spain",  
        "geojson" : {  
          "type" : "Point",  
          "coordinates" : [  
            -3.68,  
            40.43  
          ]  
        },  
        "ISO-3166-alpha3" : "ESP"  
      }  
    }  
},  
{  
  "type" : {  
    "id" : "UBERON:0002106",  
    "label" : "spleen"  
  }  
,  
  {  
    "type" : {  
      "id" : "icdot-C42.2",  
      "label" : "Spleen"  
    }  
,  
    {  
      "type" : {  
        "id" : "icdom-96893",  
        "label" : "Splenic marginal zone B-cell lymphoma"  
      }  
,  
      {  
        "type" : {  
          "id" : "NCIT:C4663",  
          "label" : "Splenic Marginal Zone Lymphoma"  
        }  
      }  
,  
      "individual_id" : "pgxind-kftx394x",  
      "individual_age_at_collection" : "P67Y",  
      "info" : {  
        "death" : "0",  
        "followup_months" : 53,  
        "callset_ids" : [  
          "pgxcs-kftvv618"  
        ],  
        "legacy_id" : "PGX_AM_BS_SMZL01"  
      },  
      "external_references" : [  
        {  
          "type" : {  
            "id" : "PMID:11337382"  
          }  
        }  
      ],  
      "provenance" : {  
        "material" : {  
          "type" : {  
            "id" : "EFO:0009656",  
            "label" : "neoplastic sample"  
          }  
        }  
      },  
      "geo" : {  
        "label" : "Salamanca, Spain",  
        "precision" : "city",  
        "city" : "Salamanca",  
        "country" : "Spain",  
        "geojson" : {  
          "type" : "Point",  
          "coordinates" : [  
            -3.68,  
            40.43  
          ]  
        },  
        "ISO-3166-alpha3" : "ESP"  
      }  
    }  
},  
{  
  "type" : {  
    "id" : "UBERON:0002106",  
    "label" : "spleen"  
  }  
,  
  {  
    "type" : {  
      "id" : "icdot-C42.2",  
      "label" : "Spleen"  
    }  
,  
    {  
      "type" : {  
        "id" : "icdom-96893",  
        "label" : "Splenic marginal zone B-cell lymphoma"  
      }  
,  
      {  
        "type" : {  
          "id" : "NCIT:C4663",  
          "label" : "Splenic Marginal Zone Lymphoma"  
        }  
      }  
,  
      "individual_id" : "pgxind-kftx394x",  
      "individual_age_at_collection" : "P67Y",  
      "info" : {  
        "death" : "0",  
        "followup_months" : 53,  
        "callset_ids" : [  
          "pgxcs-kftvv618"  
        ],  
        "legacy_id" : "PGX_AM_BS_SMZL01"  
      },  
      "external_references" : [  
        {  
          "type" : {  
            "id" : "PMID:11337382"  
          }  
        }  
      ],  
      "provenance" : {  
        "material" : {  
          "type" : {  
            "id" : "EFO:0009656",  
            "label" : "neoplastic sample"  
          }  
        }  
      },  
      "geo" : {  
        "label" : "Salamanca, Spain",  
        "precision" : "city",  
        "city" : "Salamanca",  
        "country" : "Spain",  
        "geojson" : {  
          "type" : "Point",  
          "coordinates" : [  
            -3.68,  
            40.43  
          ]  
        },  
        "ISO-3166-alpha3" : "ESP"  
      }  
    }  
},  
{  
  "type" : {  
    "id" : "UBERON:0002106",  
    "label" : "spleen"  
  }  
,  
  {  
    "type" : {  
      "id" : "icdot-C42.2",  
      "label" : "Spleen"  
    }  
,  
    {  
      "type" : {  
        "id" : "icdom-96893",  
        "label" : "Splenic marginal zone B-cell lymphoma"  
      }  
,  
      {  
        "type" : {  
          "id" : "NCIT:C4663",  
          "label" : "Splenic Marginal Zone Lymphoma"  
        }  
      }  
,  
      "individual_id" : "pgxind-kftx394x",  
      "individual_age_at_collection" : "P67Y",  
      "info" : {  
        "death" : "0",  
        "followup_months" : 53,  
        "callset_ids" : [  
          "pgxcs-kftvv618"  
        ],  
        "legacy_id" : "PGX_AM_BS_SMZL01"  
      },  
      "external_references" : [  
        {  
          "type" : {  
            "id" : "PMID:11337382"  
          }  
        }  
      ],  
      "provenance" : {  
        "material" : {  
          "type" : {  
            "id" : "EFO:0009656",  
            "label" : "neoplastic sample"  
          }  
        }  
      },  
      "geo" : {  
        "label" : "Salamanca, Spain",  
        "precision" : "city",  
        "city" : "Salamanca",  
        "country" : "Spain",  
        "geojson" : {  
          "type" : "Point",  
          "coordinates" : [  
            -3.68,  
            40.43  
          ]  
        },  
        "ISO-3166-alpha3" : "ESP"  
      }  
    }  
},  
{  
  "type" : {  
    "id" : "UBERON:0002106",  
    "label" : "spleen"  
  }  
,  
  {  
    "type" : {  
      "id" : "icdot-C42.2",  
      "label" : "Spleen"  
    }  
,  
    {  
      "type" : {  
        "id" : "icdom-96893",  
        "label" : "Splenic marginal zone B-cell lymphoma"  
      }  
,  
      {  
        "type" : {  
          "id" : "NCIT:C4663",  
          "label" : "Splenic Marginal Zone Lymphoma"  
        }  
      }  
,  
      "individual_id" : "pgxind-kftx394x",  
      "individual_age_at_collection" : "P67Y",  
      "info" : {  
        "death" : "0",  
        "followup_months" : 53,  
        "callset_ids" : [  
          "pgxcs-kftvv618"  
        ],  
        "legacy_id" : "PGX_AM_BS_SMZL01"  
      },  
      "external_references" : [  
        {  
          "type" : {  
            "id" : "PMID:11337382"  
          }  
        }  
      ],  
      "provenance" : {  
        "material" : {  
          "type" : {  
            "id" : "EFO:0009656",  
            "label" : "neoplastic sample"  
          }  
        }  
      },  
      "geo" : {  
        "label" : "Salamanca, Spain",  
        "precision" : "city",  
        "city" : "Salamanca",  
        "country" : "Spain",  
        "geojson" : {  
          "type" : "Point",  
          "coordinates" : [  
            -3.68,  
            40.43  
          ]  
        },  
        "ISO-3166-alpha3" : "ESP"  
      }  
    }  
},  
{  
  "type" : {  
    "id" : "UBERON:0002106",  
    "label" : "spleen"  
  }  
,  
  {  
    "type" : {  
      "id" : "icdot-C42.2",  
      "label" : "Spleen"  
    }  
,  
    {  
      "type" : {  
        "id" : "icdom-96893",  
        "label" : "Splenic marginal zone B-cell lymphoma"  
      }  
,  
      {  
        "type" : {  
          "id" : "NCIT:C4663",  
          "label" : "Splenic Marginal Zone Lymphoma"  
        }  
      }  
,  
      "individual_id" : "pgxind-kftx394x",  
      "individual_age_at_collection" : "P67Y",  
      "info" : {  
        "death" : "0",  
        "followup_months" : 53,  
        "callset_ids" : [  
          "pgxcs-kftvv618"  
        ],  
        "legacy_id" : "PGX_AM_BS_SMZL01"  
      },  
      "external_references" : [  
        {  
          "type" : {  
            "id" : "PMID:11337382"  
          }  
        }  
      ],  
      "provenance" : {  
        "material" : {  
          "type" : {  
            "id" : "EFO:0009656",  
            "label" : "neoplastic sample"  
          }  
        }  
      },  
      "geo" : {  
        "label" : "Salamanca, Spain",  
        "precision" : "city",  
        "city" : "Salamanca",  
        "country" : "Spain",  
        "geojson" : {  
          "type" : "Point",  
          "coordinates" : [  
            -3.68,  
            40.43  
          ]  
        },  
        "ISO-3166-alpha3" : "ESP"  
      }  
    }  
},  
{  
  "type" : {  
    "id" : "UBERON:0002106",  
    "label" : "spleen"  
  }  
,  
  {  
    "type" : {  
      "id" : "icdot-C42.2",  
      "label" : "Spleen"  
    }  
,  
    {  
      "type" : {  
        "id" : "icdom-96893",  
        "label" : "Splenic marginal zone B-cell lymphoma"  
      }  
,  
      {  
        "type" : {  
          "id" : "NCIT:C4663",  
          "label" : "Splenic Marginal Zone Lymphoma"  
        }  
      }  
,  
      "individual_id" : "pgxind-kftx394x",  
      "individual_age_at_collection" : "P67Y",  
      "info" : {  
        "death" : "0",  
        "followup_months" : 53,  
        "callset_ids" : [  
          "pgxcs-kftvv618"  
        ],  
        "legacy_id" : "PGX_AM_BS_SMZL01"  
      },  
      "external_references" : [  
        {  
          "type" : {  
            "id" : "PMID:11337382"  
          }  
        }  
      ],  
      "provenance" : {  
        "material" : {  
          "type" : {  
            "id" : "EFO:0009656",  
            "label" : "neoplastic sample"  
          }  
        }  
      },  
      "geo" : {  
        "label" : "Salamanca, Spain",  
        "precision" : "city",  
        "city" : "Salamanca",  
        "country" : "Spain",  
        "geojson" : {  
          "type" : "Point",  
          "coordinates" : [  
            -3.68,  
            40.43  
          ]  
        },  
        "ISO-3166-alpha3" : "ESP"  
      }  
    }  
},  
{  
  "type" : {  
    "id" : "UBERON:0002106",  
    "label" : "spleen"  
  }  
,  
  {  
    "type" : {  
      "id" : "icdot-C42.2",  
      "label" : "Spleen"  
    }  
,  
    {  
      "type" : {  
        "id" : "icdom-96893",  
        "label" : "Splenic marginal zone B-cell lymphoma"  
      }  
,  
      {  
        "type" : {  
          "id" : "NCIT:C4663",  
          "label" : "Splenic Marginal Zone Lymphoma"  
        }  
      }  
,  
      "individual_id" : "pgxind-kftx394x",  
      "individual_age_at_collection" : "P67Y",  
      "info" : {  
        "death" : "0",  
        "followup_months" : 53,  
        "callset_ids" : [  
          "pgxcs-kftvv618"  
        ],  
        "legacy_id" : "PGX_AM_BS_SMZL01"  
      },  
      "external_references" : [  
        {  
          "type" : {  
            "id" : "PMID:11337382"  
          }  
        }  
      ],  
      "provenance" : {  
        "material" : {  
          "type" : {  
            "id" : "EFO:0009656",  
            "label" : "neoplastic sample"  
          }  
        }  
      },  
      "geo" : {  
        "label" : "Salamanca, Spain",  
        "precision" : "city",  
        "city" : "Salamanca",  
        "country" : "Spain",  
        "geojson" : {  
          "type" : "Point",  
          "coordinates" : [  
            -3.68,  
            40.43  
          ]  
        },  
        "ISO-3166-alpha3" : "ESP"  
      }  
    }  
},  
{  
  "type" : {  
    "id" : "UBERON:0002106",  
    "label" : "spleen"  
  }  
,  
  {  
    "type" : {  
      "id" : "icdot-C42.2",  
      "label" : "Spleen"  
    }  
,  
    {  
      "type" : {  
        "id" : "icdom-96893",  
        "label" : "Splenic marginal zone B-cell lymphoma"  
      }  
,  
      {  
        "type" : {  
          "id" : "NCIT:C4663",  
          "label" : "Splenic Marginal Zone Lymphoma"  
        }  
      }  
,  
      "individual_id" : "pgxind-kftx394x",  
      "individual_age_at_collection" : "P67Y",  
      "info" : {  
        "death" : "0",  
        "followup_months" : 53,  
        "callset_ids" : [  
          "pgxcs-kftvv618"  
        ],  
        "legacy_id" : "PGX_AM_BS_SMZL01"  
      },  
      "external_references" : [  
        {  
          "type" : {  
            "id" : "PMID:11337382"  
          }  
        }  
      ],  
      "provenance" : {  
        "material" : {  
          "type" : {  
            "id" : "EFO:0009656",  
            "label" : "neoplastic sample"  
          }  
        }  
      },  
      "geo" : {  
        "label" : "Salamanca, Spain",  
        "precision" : "city",  
        "city" : "Salamanca",  
        "country" : "Spain",  
        "geojson" : {  
          "type" : "Point",  
          "coordinates" : [  
            -3.68,  
            40.43  
          ]  
        },  
        "ISO-3166-alpha3" : "ESP"  
      }  
    }  
},  
{  
  "type" : {  
    "id" : "UBERON:0002106",  
    "label" : "spleen"  
  }  
,  
  {  
    "type" : {  
      "id" : "icdot-C42.2",  
      "label" : "Spleen"  
    }  
,  
    {  
      "type" : {  
        "id" : "icdom-96893",  
        "label
```

# Progenetix Documentation

## docs.progenetix.org

- information about concepts and practices
- API access
- Use cases
- Ontologies

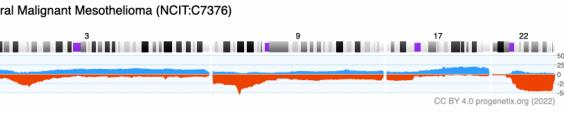
**Download or Plot CNV Frequencies**

**Collation plots**

The Progenetix resource provides pre-computed CNV frequencies for all its "collations" such as

- cancer types by e.g. NCIt, ICD-O morphology and topography codes
- experimental series, e.g. all samples from a given publication
- custom cohorts, e.g. all samples used in a Progenetix meta-analysis or external project such as TCGA

This data can be accessed through the Progenetix API in data and image format.



**Query and export segment copy number variant data**

You can download the copy number variant data of individual biosamples queried by pgxRpi or by Progenetix website. The variant data export formats, more information see vignettes.

```
variants <- pgxLoader(type="variant", biosample_id = c("pgxb5-kftva6du", "pgxb6-1kqjw5u"))
```

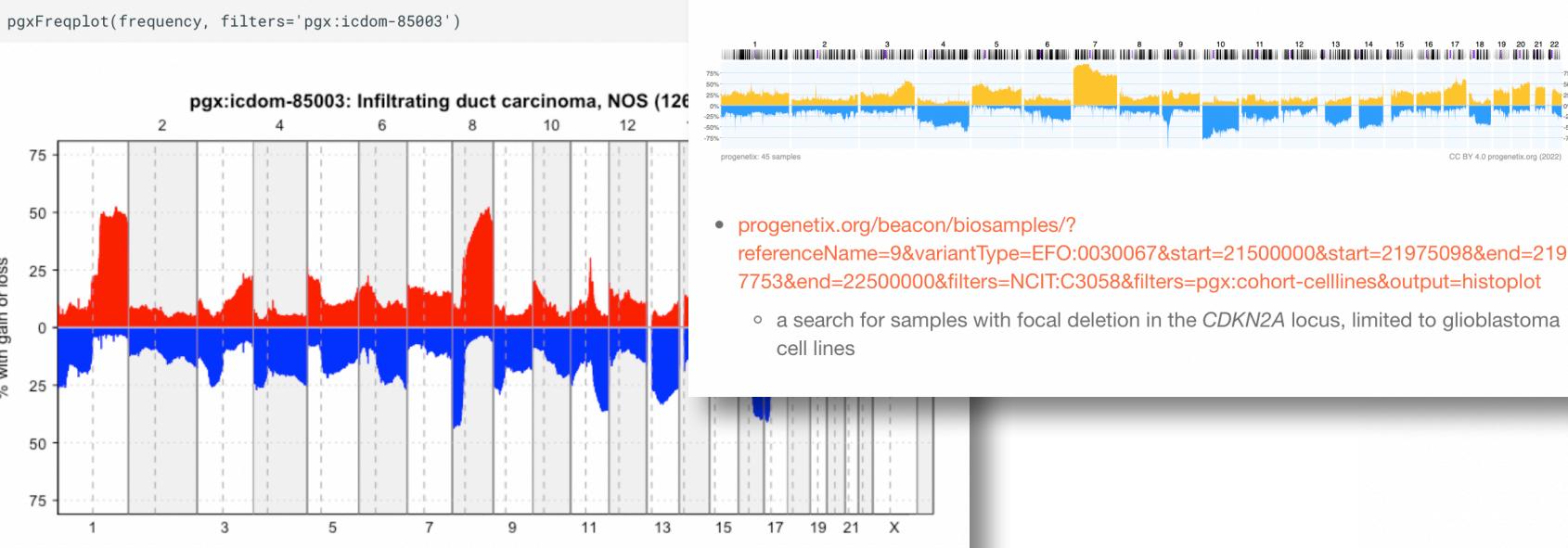
**Query and visualize CNV frequencies**

You can query the CNV frequency of specific filters, namely specific clinical filters available data formats. One is `.pgxseg`, good for visualization. Another for analysis.

```
frequency <- pgxLoader(type="frequency", output = "pgxseg",
                        filters=c("NCIT:C4638", "pgx:icdom-85003"),
                        codematches = TRUE)
```

The data visualization requires the input data with `.pgxseg` format. You can genome, by chromosomes, or plot like circos.

```
pgxFreqplot(frequency, filters='pgx:icdom-85003')
```



**Interval frequencies are per default stored in a 1Mb binned format. More information about the API use can be found in the [IntervalFrequencies API documentation](#).**

**Query-based histograms**

The Progenetix Beacon responses - depending on their type - usually contain a handover URL to retrieve CNV histogram and/or sample plots of the samples matched by the query. The `bycon` API now offers a direct access to the histograms without the need to deparse JSON response first. The switch to the histogram is initiated by adding `&output=histoplot` to the Beacon query URL. Then, the API will first query the samples and then perform a handover to the plotting API. Please be aware that this procedure is best suited for limited queries and may lead to a time-out.

**Examples:**



- [progenetix.org/beacon/biosamples/?referenceName=9&variantType=EFO:0030067&start=21500000&start=21975098&end=21967753&end=22500000&filters=NCIT:C3058&filters=pgx:cohort-celllines&output=histoplot](https://progenetix.org/beacon/biosamples/?referenceName=9&variantType=EFO:0030067&start=21500000&start=21975098&end=21967753&end=22500000&filters=NCIT:C3058&filters=pgx:cohort-celllines&output=histoplot)
  - a search for samples with focal deletion in the CDKN2A locus, limited to glioblastoma cell lines

Progenetix Documentation

[Documentation Home](#)

[Progenetix News](#)

[Use Case Examples](#)

[Services & API](#)

[Classifications, Ontologies & Standards](#)

[Publication Collection](#)

[Changelog](#)

[Data Review](#)

[Beacon+ & bycon](#)

[Technical Notes](#)

[Progenetix Data](#)

[Baudisgroup @ UZH](#)

## Progenetix Cancer Genomics Resource Documentation

The Progenetix database and cancer genomic information resource contains genome profiles of more than 100'000 individual cancer genome screening experiments. The genomic profiling data was derived from genomic arrays and chromosomal Comparative Genomic Hybridization (CGH) 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.

### Citation

Huang Q, Carrio-Cordo P, Gao B, Paloots R, Baudis M. (2021) **The Progenetix oncogenomic resource in 2021**. *Database (Oxford)*. 2021 Jul 17  
progenetix.org: **Progenetix oncogenomic online resource** (2022)

### Additional Articles & Citation Options

### Registration & Licenses

The **Progenetix** database and cancer genomic information resource was publicly launched in 2001, announced through an article in *Bioinformatics*. The database & software are developed by the [group of Michael Baudis](#) at the [University of Zurich](#) and the Swiss Institute of Bioinformatics ([SIB](#)).

Additional information - e.g. about contacts or related publications - is available through the [group page](#) of the Baudis group at the University of Zürich. For a list of publications by the Baudis group you can go to the [group's website](#), [EuropePMC](#) or any of the links here.

## Progenetix Source Code

With exception of some utility scripts and external dependencies (e.g. [MongoDB](#)) the following projects provide the vast majority of the software (from database interaction to website) behind Progenetix and [Beacon+](#).

### bycon

- Python based service based on the [GA4GH Beacon protocol](#)
- software powering the Progenetix resource
- [Beacon+](#) implementation(s) use the same code base

### progenetix-web

- website for Progenetix and its [Beacon+](#) implementations
- provides Beacon interfaces for the [bycon](#) server, as well as other Progenetix services (e.g. the [publications](#) repository)
- implemented as [React / Next.js](#) project

### PGX

- a Perl library providing utility functions for Progenetix CNV data
- used for data transformation, e.g. binning of segmental CNV data
- main purpose now in providing the various plots (CNV histograms, clustered CNV profiles, array plots)

## Table of contents

[Progenetix Source Code](#)

[bycon](#)

[progenetix-web](#)

[PGX](#)

[Additional Projects](#)

# **Summary: Offers & Needs**

# Progenetix Needs & Offers

## What we have ...

- ✓ collection of >4000 articles assessed for scope
  - training set for NLP & search engine generation
- ✓ cancer specific ontologies with cross-mappings (ICD-O vs. NCIt) based on >100k samples
  - existing service API
- ✓ metadata ontology mappings for some 10k samples, with varying coverage for grade / stage / survival / ...
- ✓ CNV profiles for >110k samples, >700 entities with disease codes and metadata
- ✓ cell line CNV profiles together with mapped variants with clinical evidences

## What we'd like...

- (semi-)automated detection of additional articles for scope (genome screening technologies, cancer samples, geographies)
- generation of a complete ICD-O terminology tree with NCIT (?) correspondence
  - improved service API & publication
- improved annotations using smarter source (article, annotation files) pre-/processing
- correlation between individual profiles, profile heterogeneity and external parameters
- relation between cell lines and native tumor types, with consideration of non-CNV parameters

# Progenetix Needs & Offers

## What we have ...

- ✓ collection of >4000 articles assessed for scope
  - training set for NLP & search engine generation
- ✓ cancer specific ontologies with cross-mappings (ICD-O → NCIT → TCGA → ...)
  - exist
- ✓ metadata ontology mappings for some 10k samples, with varying coverage for grade / stage / survival / ...
- ✓ CNV profiles for >110k samples, >700 entities with disease codes and metadata
- ✓ cell line CNV profiles together with mapped variants with clinical evidences

## Open Access API

## What we'd like...

- (semi-)automated detection of additional articles for scope (genome screening technologies, cancer samples, geographies)
- generation of a complete ICD-O terminology tree with NCIT (?) correspondence
  - improved service API & publication
- improved annotations using smarter source (article, annotation files) pre-/processing
- correlation between individual profiles, profile heterogeneity and external parameters
- relation between cell lines and native tumor types, with consideration of non-CNV parameters

# Progenetix Needs & Offers

## Publication resource

- collection of >4000 articles assessed for scope
- training set for NLP & search engine generation
- identification of content-similar articles for expansion of publication resource
- text-analysis for resource identifiers (e.g. GEO accessions...), geographic term, clinical annotations
  - for new data and in comparison to existing datasets in Progenetix (e.g. hints for usable metadata not yet in DB)
- **Opening up Publication resource for wider/ more general use?**

## What we'd like...

- (semi-)automated detection of additional articles for scope (genome screening technologies, cancer samples, geographies)
- generation of a complete ICD-O terminology tree with NCIT (?) correspondence
  - improved service API & publication
- improved annotations using smarter source (article, annotation files) pre-/processing
- correlation between individual profiles, profile heterogeneity and external parameters
- relation between cell lines and native tumor types, with consideration of non-CNV parameters

# NCIt <=> ICDO: Complete Mapping w/ Resource

## NCIT

Neoplasm by Site (R_1)
Breast Neoplasm (C2910)
Breast Adenoma (C40382)
Breast Apocrine Adenoma (C40383)
Breast Ductal Adenoma (C40384)
Breast Pleomorphic Adenoma (C40408)
Breast Tubular Adenoma (C62210)
Lactating Adenoma (C9473)
Breast Carcinoma (C4872)
Adenoid Cystic Breast Carcinoma (C5130)
Bilateral Breast Carcinoma (C8287)
Breast Adenocarcinoma (C5214)
Ductal Breast Carcinoma (C4017)
Acinic Cell Breast Carcinoma (C40367)
Ductal Breast Carcinoma In Situ (C2924)
Invasive Ductal Carcinoma, Not Otherwise Specified (C4194)
Medullary Breast Carcinoma (C9119)
Mucinous Breast Carcinoma (C9131)
Papillary Breast Carcinoma (C9134)
Intraductal Papillary Breast Carcinoma (C4190)
Invasive Papillary Breast Carcinoma (C36085)
Secretory Breast Carcinoma (C4189)
Tubular Breast Carcinoma (C9135)
Inflammatory Breast Carcinoma (C4001)
Lobular Breast Carcinoma (C3771)
Invasive Lobular Breast Carcinoma (C7950)
Lobular Breast Carcinoma In Situ (C4018)
Mixed Lobular and Ductal Breast Carcinoma (C5160)
Paget Disease of the Breast (C47857)
Breast Carcinoma by Gene Expression Profile (C53553)
Basal-Like Breast Carcinoma (C53558)
HER2 Positive Breast Carcinoma (C53556)
Luminal A Breast Carcinoma (C53554)
Luminal B Breast Carcinoma (C53555)
Normal Breast-Like Subtype of Breast Carcinoma (C53557)
Triple-Negative Breast Carcinoma (C71732)
Breast Small Cell Carcinoma (C6760)
Female Breast Carcinoma (C2918)
Hereditary Breast Carcinoma (C4503)
Male Breast Carcinoma (C3862)
Nipple Carcinoma (C28432)
Sporadic Breast Carcinoma (C7566)
Squamous Cell Breast Carcinoma (C5177)
Unilateral Breast Carcinoma (C46073)
Breast Fibroadenoma (C3744)
Breast Complex Fibroadenoma (C5194)
Breast Giant Fibroadenoma (C4273)
Breast Intracanalicular Fibroadenoma (C4271)
Breast Juvenile Fibroadenoma (C4276)
Breast Pericanalicular Fibroadenoma (C4272)



## ICDO (T/M)

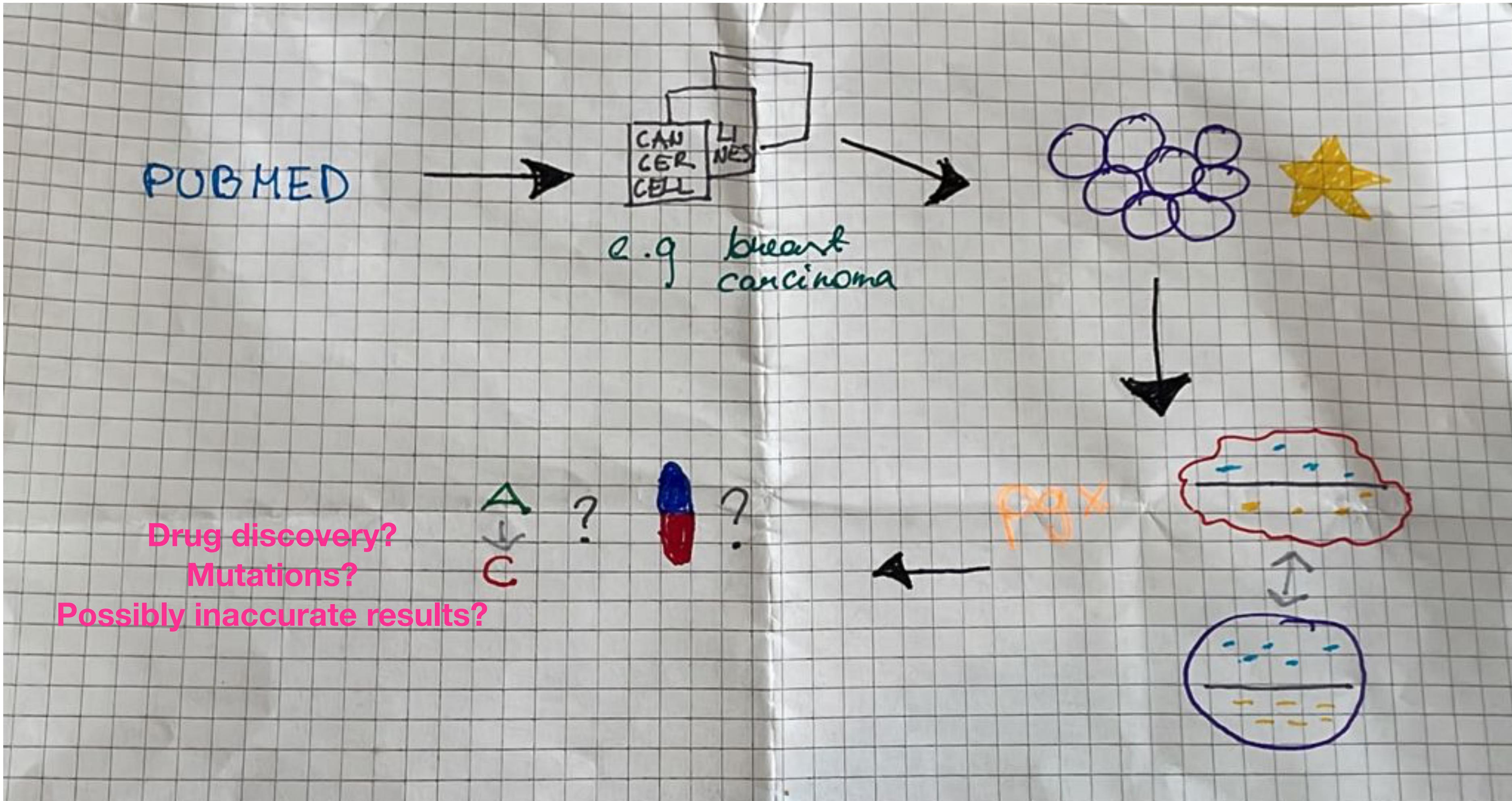
Abdomen	
C76.2	NOS
C47.4	autonomic nervous system
C49.4	connective tissue
C49.4	muscle
C47.4	peripheral nerve
C44.5	skin
C49.4	subcutaneous tissue
Abdominal	
C49.4	aorta
C15.2	esophagus
C77.2	lymph node
C49.4	vena cava
8822/1	Abdominal desmoid
8822/1	Abdominal fibromatosis
Abdominal wall	
C76.2	NOS
C44.5	NOS (carcinoma, melanoma, nevus)
C49.4	NOS (sarcoma, lipoma)
C49.4	adipose tissue

**Mature B-cell Neoplasms**

9737/3	ALK positive large B-cell lymphoma
9680/3	B-cell lymphoma, unclassifiable, with features intermediate between diffuse large B-cell lymphoma and Burkitt lymphoma
9596/3	B-cell lymphoma, unclassifiable, with features intermediate between diffuse large B-cell lymphoma and classical Hodgkin lymphoma
9833/3	B-cell prolymphocytic leukemia
9687/3	Burkitt lymphoma
9823/3	Chronic lymphocytic leukemia/small lymphocytic lymphoma
9680/3	Diffuse large B-cell lymphoma (DLBCL), NOS
9699/3	Extranodal marginal zone lymphoma of mucosa-associated lymphoid tissue (MALT lymphoma)
9734/3	Extraosseous plasmacytoma
9690/3	Follicular lymphoma
9762/3	Heavy chain diseases (alpha, gamma, mu)
9712/3	Intravascular large B-cell lymphoma
9738/3	Large B-cell lymphoma arising in HHV8-associated multicentric Castleman disease
9766/1	Lymphomatoid granulomatosis
9671/3	Lymphoplasmacytic lymphoma
9673/3	Mantle cell lymphoma
9699/3	Nodal marginal zone lymphoma
9591/3	Non-Hodgkin lymphoma, NOS; Splenic B-cell lymphoma/leukemia, unclassifiable

- expand our data-driven mapping to a complete coverage of all cancer codes
- integrate with Monarch initiative

# Cell lines



# **Progenetix and GA4GH Beacon**

## **Implementation driven development of a GA4GH standard**





GENOMICS

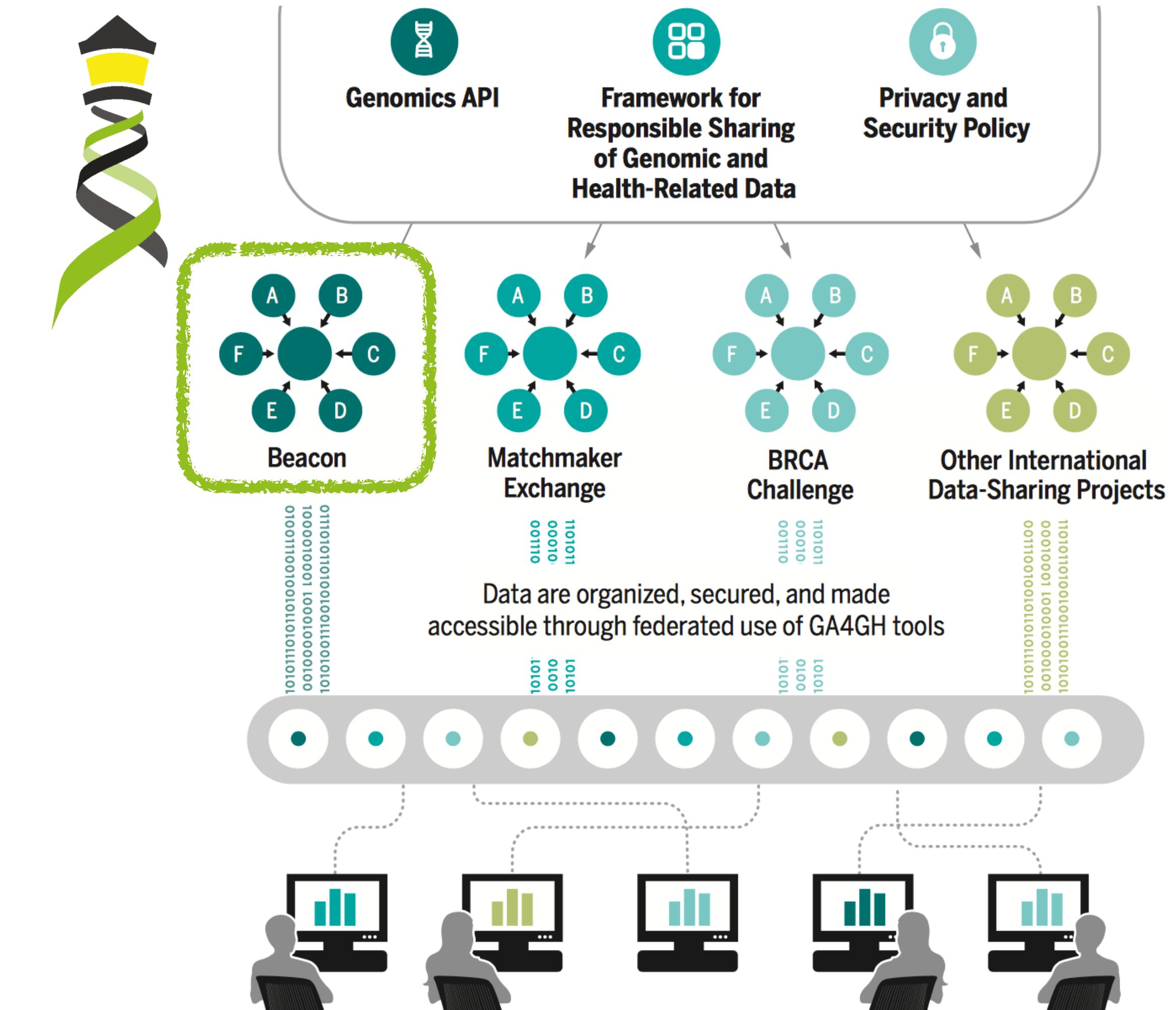
# *A federated ecosystem for sharing genomic, clinical data*

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 62

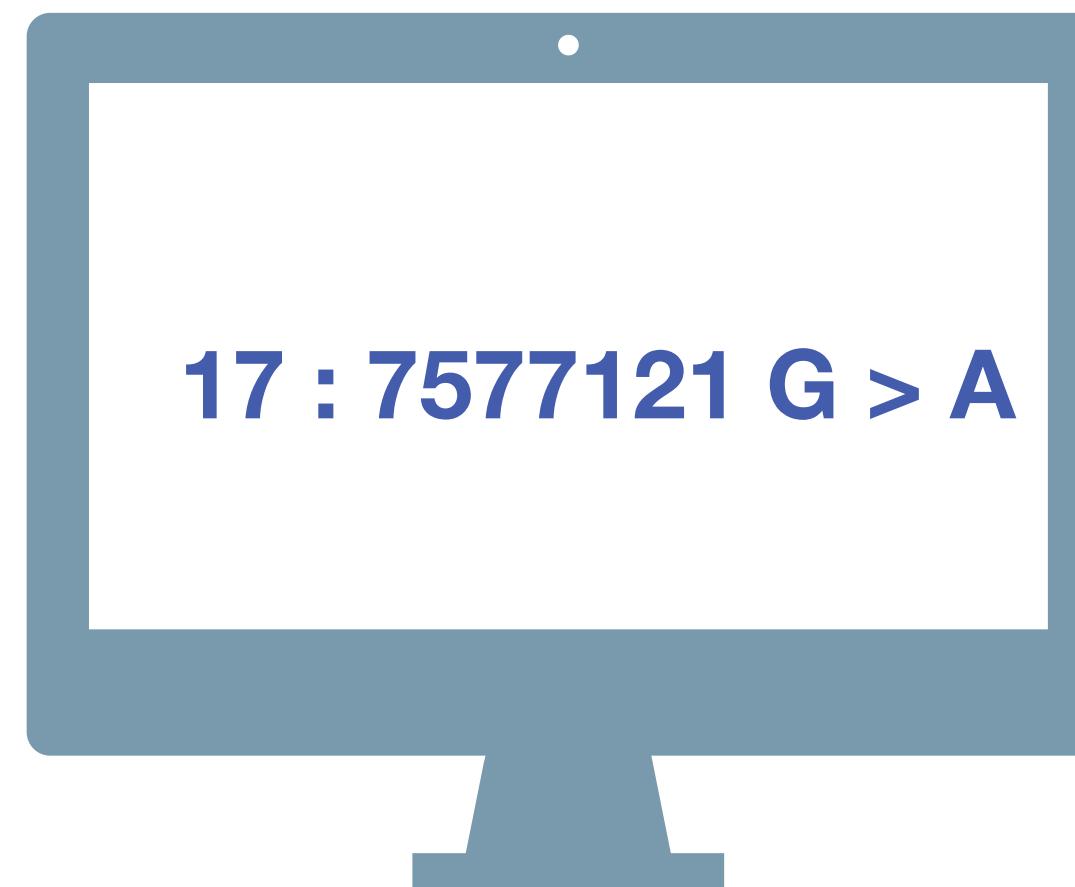
**A federated data ecosystem.** To share genomic data globally, this approach furthers medical research without requiring compatible data sets or compromising patient identity.



# DNASTACK



Global Alliance  
for Genomics & Health

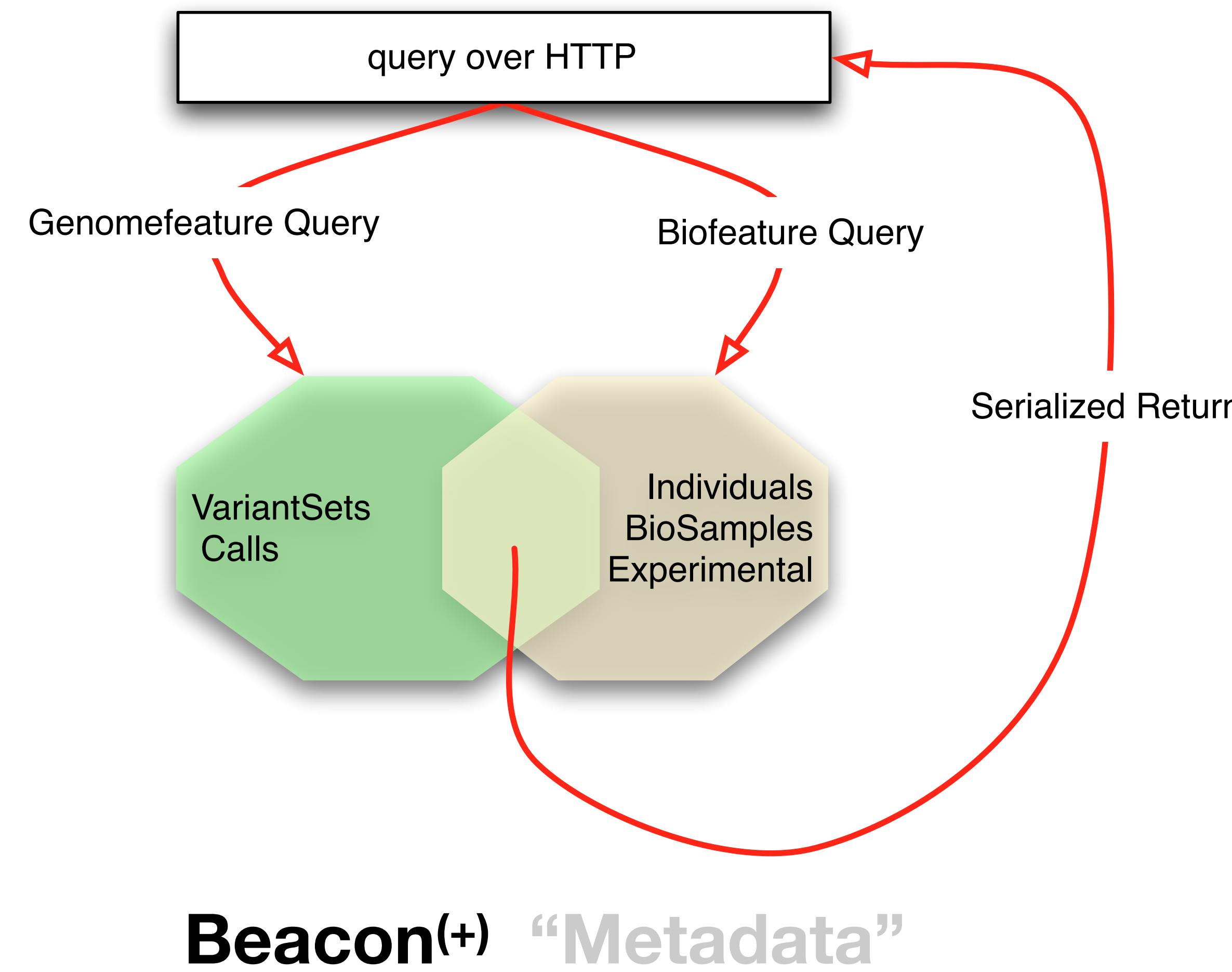


# Beacon

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

**YES | NO | \0**

# Minimal GA4GH query API structure



# Beacon v2 Filters

# **Example: Use of hierarchical classification systems (here NCI ICD neoplasm core)**

- 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



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

↓

progenetix

Variants: 0    falleles: 0    Callsets Variants ↗    UCSC region ↗    Legacy Interface ↗     [Show JSON Response](#)

Results    **Biosamples**

<b>Id</b>	<b>Description</b>	<b>Classifications</b>	<b>Identifiers</b>	<b>DEL</b>	<b>DUP</b>	<b>CNV</b>
PGX_AM_BS_MCC01	Merkel cell carcinoma	<a href="#">icdot-C44.9</a> Skin, NOS <a href="#">icdom-82473</a> Merkel cell carcinoma <a href="#">NCIT:C9231</a> Merkel Cell Carcinoma	<a href="#">PMID:9537255</a>	0.116	0.104	0.22
PGX_AM_BS_MCC02	Merkel cell carcinoma	<a href="#">icdot-C44.9</a> Skin, NOS <a href="#">icdom-82473</a> Merkel cell carcinoma <a href="#">NCIT:C9231</a> Merkel Cell Carcinoma	<a href="#">PMID:9537255</a>	0.154	0.056	0.21
PGX_AM_BS_MCC03	Merkel cell carcinoma	<a href="#">icdot-C44.9</a> Skin, NOS <a href="#">icdom-82473</a> Merkel cell carcinoma <a href="#">NCIT:C9231</a> Merkel Cell Carcinoma	<a href="#">PMID:9537255</a>	0.137	0.21	0.347
PGX_AM_BS_MCC04	Merkel cell carcinoma	<a href="#">icdot-C44.9</a> Skin, NOS <a href="#">icdom-82473</a> Merkel cell carcinoma <a href="#">NCIT:C9231</a> Merkel Cell Carcinoma	<a href="#">PMID:9537255</a>	0.158	0.056	0.214
PGX_AM_BS_MCC05	Merkel cell carcinoma	<a href="#">icdot-C44.9</a> Skin, NOS <a href="#">icdom-82473</a> Merkel cell carcinoma <a href="#">NCIT:C9231</a> Merkel Cell Carcinoma	<a href="#">PMID:9537255</a>	0.107	0.327	0.434

« < > »

Page 1 of 105

# Beacon v2 Requests

## POSTing Queries

- Beacon v2 supports a mix of dedicated endpoints with REST paths
- POST requests using JSON query documents
- final syntax for core parameters still in testing stages

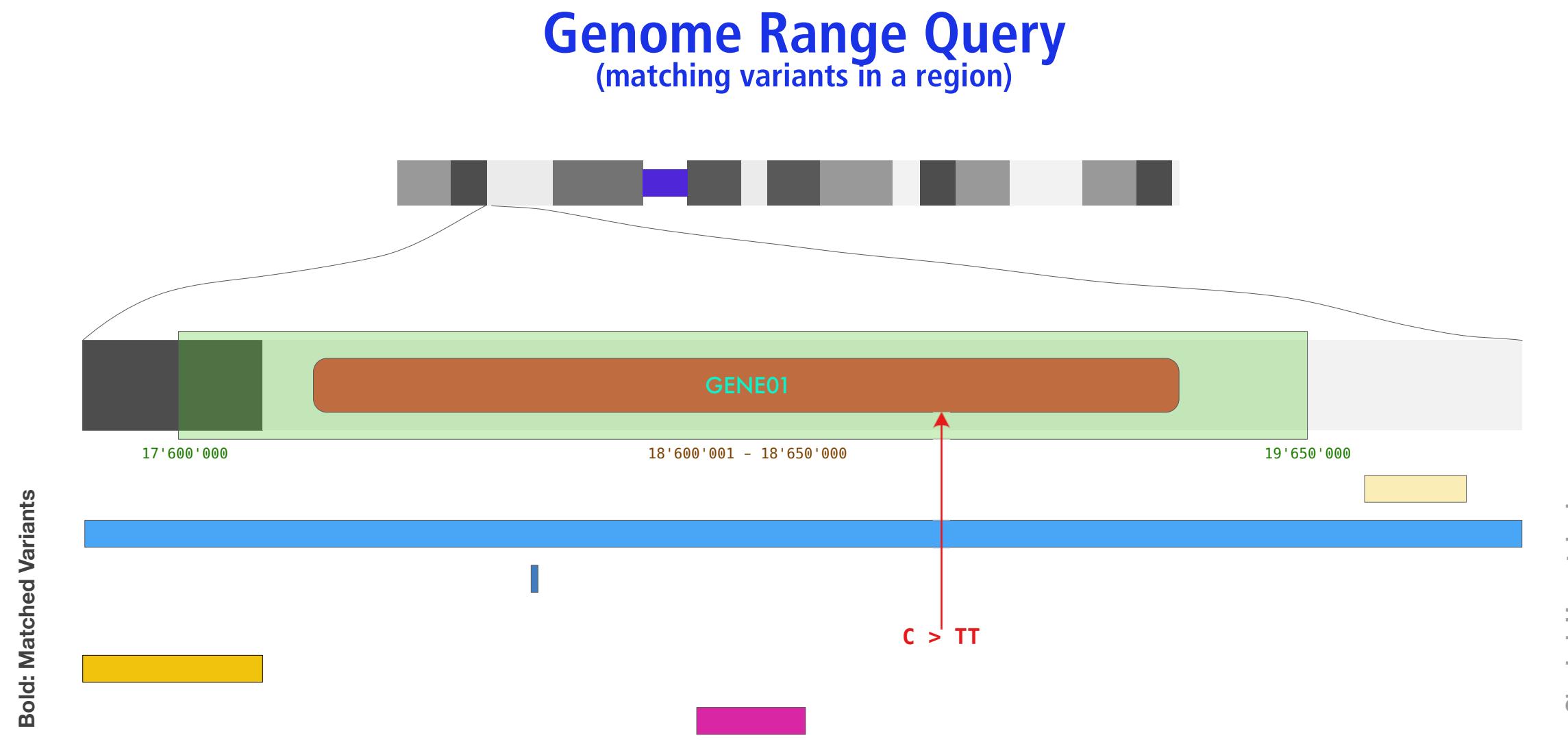
```
{  
  "$schema": "beaconRequestBody.json",  
  "meta": {  
    "apiVersion": "2.0",  
    "requestedSchemas": [  
      {  
        "entityType": "individual",  
        "schema": "https://progenetix.org/services/schemas/Phenopacket/"  
      }  
    ],  
    "query": {  
      "requestParameters": {  
        "datasets": {  
          "datasetIds": ["progenetix"]  
        }  
      },  
      "filterLogic": "OR"  
    },  
    "pagination": {  
      "skip": 0,  
      "limit": 10  
    },  
    "filters": [  
      { "id": "NCIT:C4536" },  
      { "id": "NCIT:C95597" },  
      { "id": "NCIT:C7712" }  
    ]  
  }  
}
```



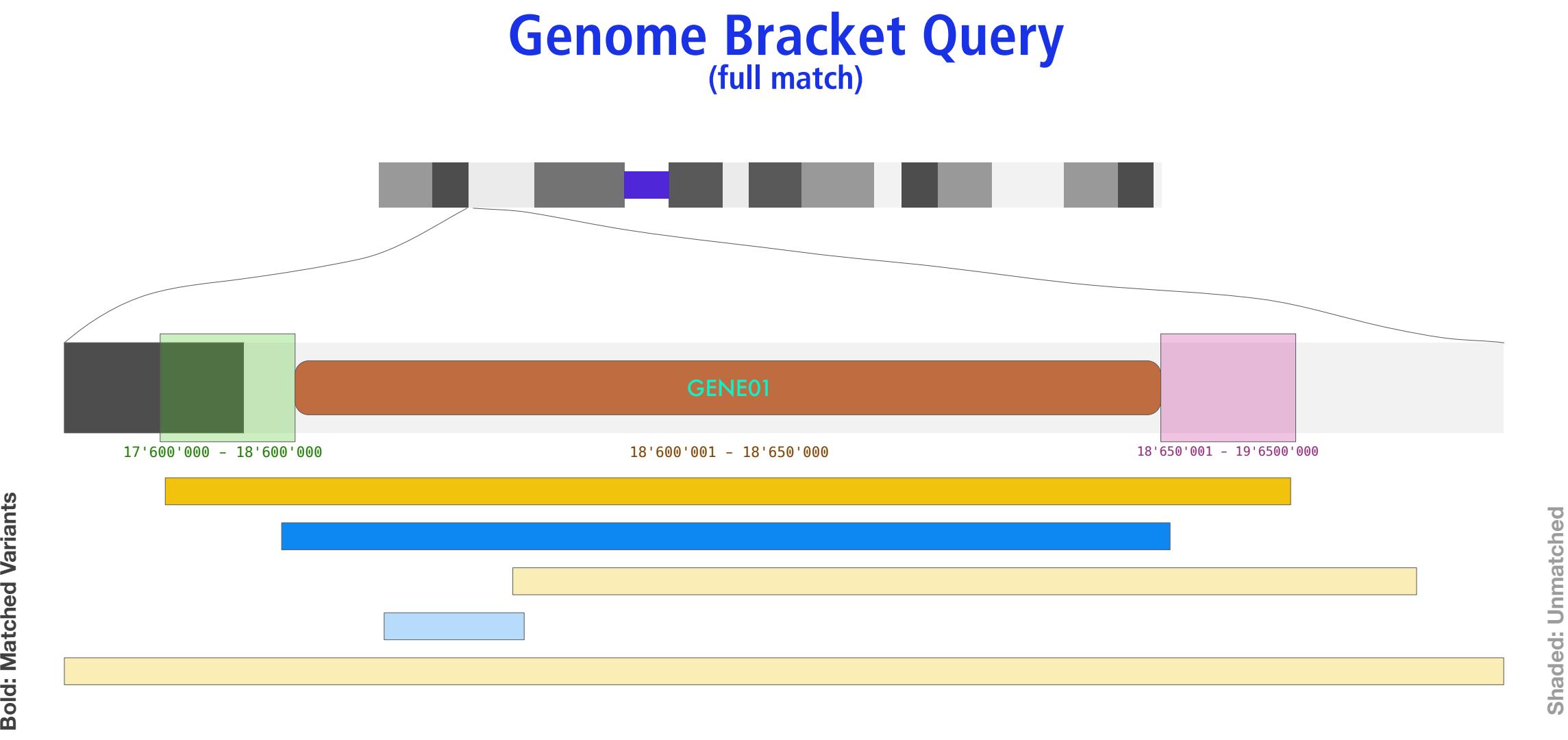


# Beacon v2: Extended Variant Queries

Range and Bracket queries enable positional wildcards and fuzziness



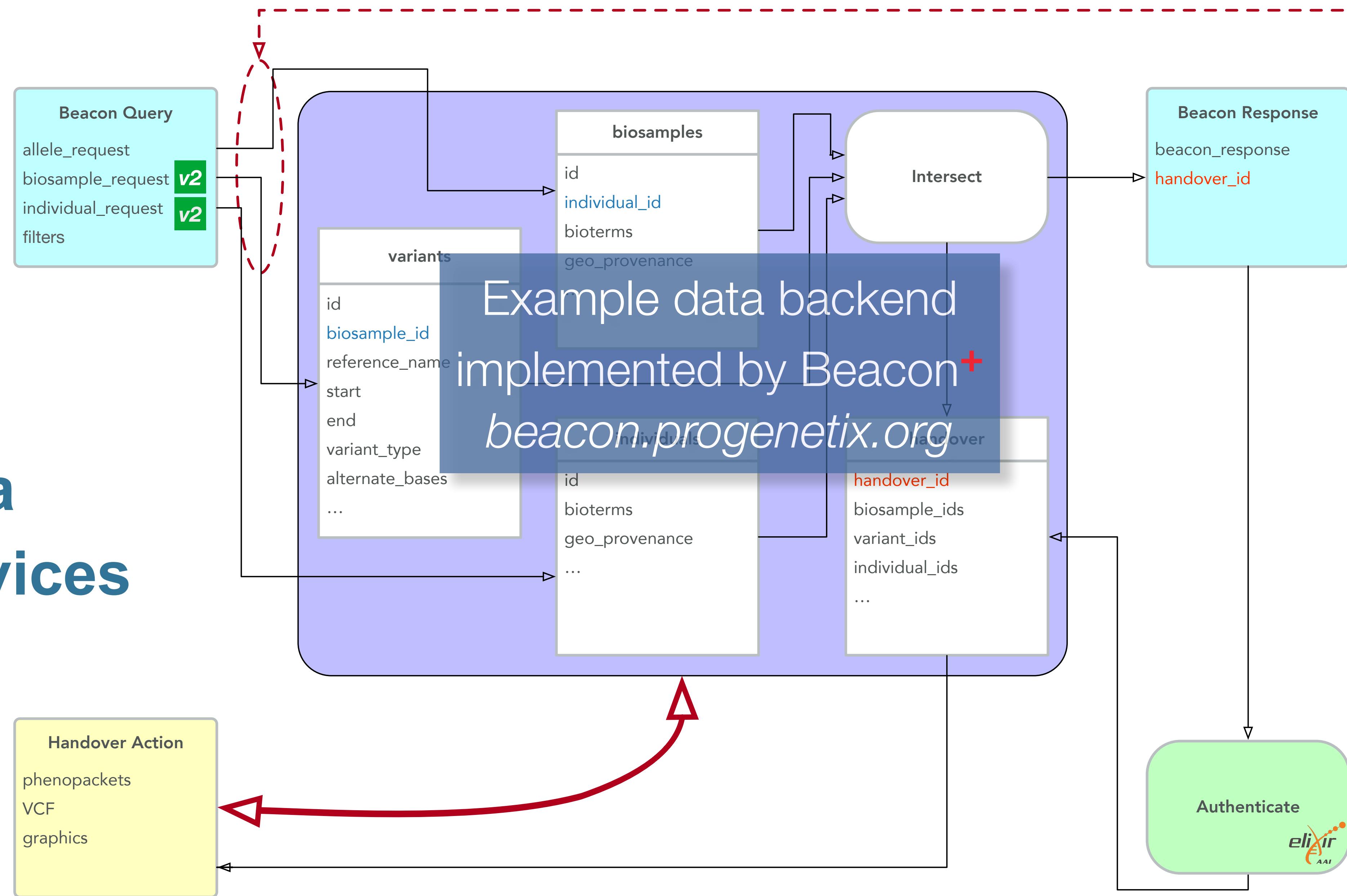
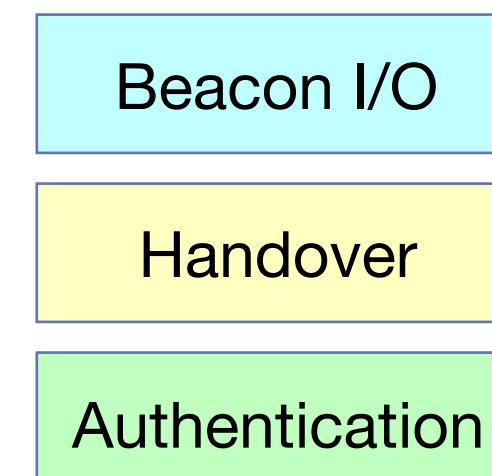
- Genome Range Queries provide a way to "fish" for variants overlapping an indicated region, e.g. the CDR of a gene of interest
- Additional parameters (e.g. variant type, reference or alternate bases) limit the scope of the responses
- new Beacon v2 size parameters to limit structural variants (e.g. "focal" CNVs)



- Genome Bracket Queries allow to search for structural variants with start and end positions falling into defined sequence ranges
- allows to query any contiguous genomic variant (and in principle also can step in for range queries)
- typical use case is e.g. the query for variants such as duplications covering the whole CDR of a gene, while limiting the allowed start or end regions

# Beacon & Handover

Beacons v1.1  
supports data  
delivery services



ga4gh-beacon / beacon-framework-v2 Public

Code Issues 18 Pull requests 2 Discussions Actions Wiki Security Insights Settings

main 7 branches 0 tags Go to file Add file Code About

jrambla Merge pull request #51 from ga4gh-beacon/configuration-typos-fixes ...

common	de-lining \n
configuration	speling in configuration -> filteringTermsSchema
requests	de-lining \n
responses	de-lining \n
.gitignore	Initial commit
LICENSE	Initial commit
README.md	Adding naming conventions to readme
endpoints.json	de-lining \n

README.md

## beacon-framework-v2

Beacon Framework version 2

### Introduction

The GA4GH Beacon specification is composed by two parts:

- the Beacon Framework (in *this* repo)
- the Beacon Model (in the [Models repo](#))

The Beacon Framework is the part that describes the overall structure of the API

progenetix / bycon Public

Code Issues Pull requests 1 Actions Projects Wiki Security Insights Settings

master 3 branches 0 tags Go to file Add file Code About

mbaudis Update README.md 5064e89 11 seconds ago 519 commits

beaconServer	datatables, genesRefresher	6 days ago
byconeer	datatables, genesRefresher	6 days ago
config	datatables, genesRefresher	6 days ago
lib	intervalFrequencies service & some library shuffling	5 months ago
schemas	datatables, genesRefresher	6 days ago
services	genespan method for gene request size reduction	2 days ago
remnants	biocharacteristics removal; shuffling of beaconsv2 references...	21 days ago
.gitignore	biocharacteristics removal; shuffling of beaconsv2 references...	21 days ago
LICENSE	Create LICENSE	12 months ago
README.md	Update README.md	11 seconds ago
__init__.py	intervalFrequencies service & some library shuffling	5 months ago
requirements.txt	add non-interactive mode	16 months ago

README.md

License CC0 1.0

### Bycon - a Python-based environment for the Beacon v2 genomics API

The `bycon` project - at least at its current stage - is a mix of Progenetix (i.e. GA4GH object model derived, MongoDB implemented) - data management, and the implementation of middleware & server for the Beacon API.

More information about the current status of the package can be found in the inline documentation which is also presented in an accessible format on the [Progenetix website](#).

About

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

Readme

CC0-1.0 License

Releases

No releases published [Create a new release](#)

Packages

No packages published [Publish your first package](#)

Contributors 4

mbaudis Michael Baudis
sofiapfund Sofia
qingyao
KyleGao Bo Gao

Languages

Python 99.9% Shell 0.1%

# Onboarding Demonstrating Compliance



- Progenetix Beacon+ has served as implementation driver since 2016
- Beacon v2 as service with protocol-driven registries for federation
- GA4GH approval process in the Spring 2022 session

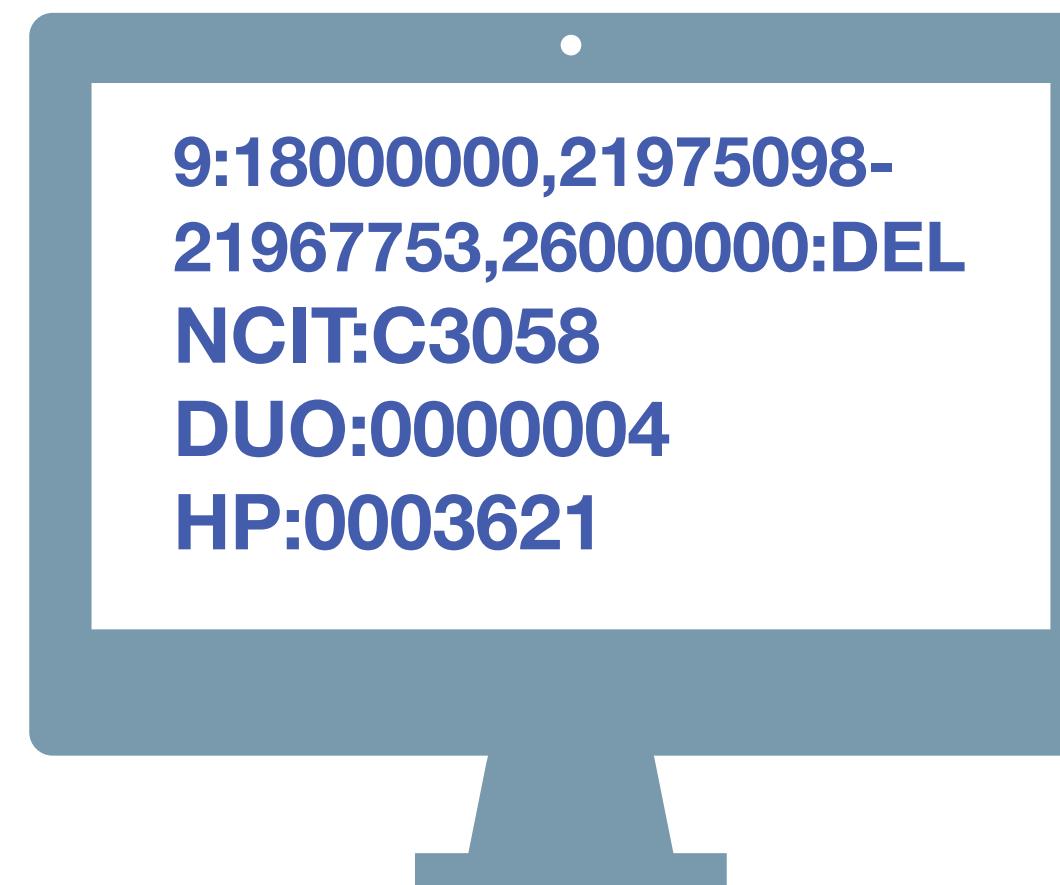
Beacon v2 GA4GH Approval Registry

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

The screenshot displays the Beacon v2 GA4GH Approval Registry interface. It shows five separate Beacon entries: European Genome-Phenome Archive (EGA), progenetix, baudisgroup at UZH and SIB, cnag, and Centre Nacional Analisis Genomica (CNAG-CRG). Each entry includes a summary card with links to visit the site, use the Beacon API, or contact them. Below each card is a detailed table of data types (BeaconMap, Bioinformatics analysis, Biological Sample, Cohort, Configuration, Dataset, EntryTypes, Genomic Variants, Individual, Info, Sequencing run) with colored bars indicating compliance status: green for matches the spec, red for not matching, and white for not implemented.

Beacon	BeaconMap	Bioinformatics analysis	Biological Sample	Cohort	Configuration	Dataset	EntryTypes	Genomic Variants	Individual	Info	Sequencing run
European Genome-Phenome Archive (EGA)	Matched	Matched	Matched	Matched	Matched	Matched	Matched	Matched	Matched	Matched	Matched
progenetix	Matched	Matched	Matched	Matched	Matched	Matched	Matched	Matched	Matched	Matched	Matched
baudisgroup at UZH and SIB	Matched	Matched	Matched	Matched	Matched	Matched	Matched	Matched	Matched	Matched	Matched
cnag	Matched	Matched	Matched	Matched	Matched	Matched	Matched	Matched	Matched	Matched	Matched
Centre Nacional Analisis Genomica (CNAG-CRG)	Matched	Matched	Matched	Matched	Matched	Matched	Matched	Matched	Matched	Matched	Matched
University of Leicester	Matched	Matched	Matched	Matched	Matched	Matched	Matched	Matched	Matched	Matched	Matched
Clinical Bioinformatics Area	Matched	Matched	Matched	Matched	Matched	Matched	Matched	Matched	Matched	Matched	Matched

Legend: ✓ Matches the Spec, ✘ Not Match the Spec, ⓘ Not Implemented



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".