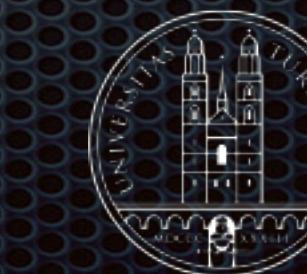
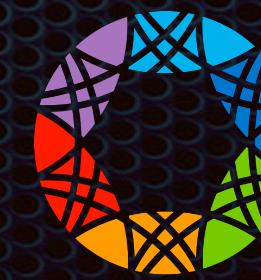


Connecting the silos

Genomic Data Standards, Resources and the
Global Alliance for Genomics and Health



University of
Zurich^{UZH}



Global Alliance
for Genomics & Health

1992



2001



2003



2006



2007



Heidelberg

Stanford

Gainesville

Aachen

Zürich

Student of medicine | doctoral thesis in molecular cytogenetics @ DKFZ (Peter Lichter) | resident in clinical hematology/oncology | data, clinical studies & cancer systematics

Post-doc in hemato-pathology (Michael Cleary) | molecular mechanisms of leukemogenesis | transgenic models | expression arrays | systematic cancer genome data collection | *Progenetix* website

Assistant professor in paediatric haematology | molecular mechanisms of leukemogenesis | focus on bioinformatics for cancer genome data analysis

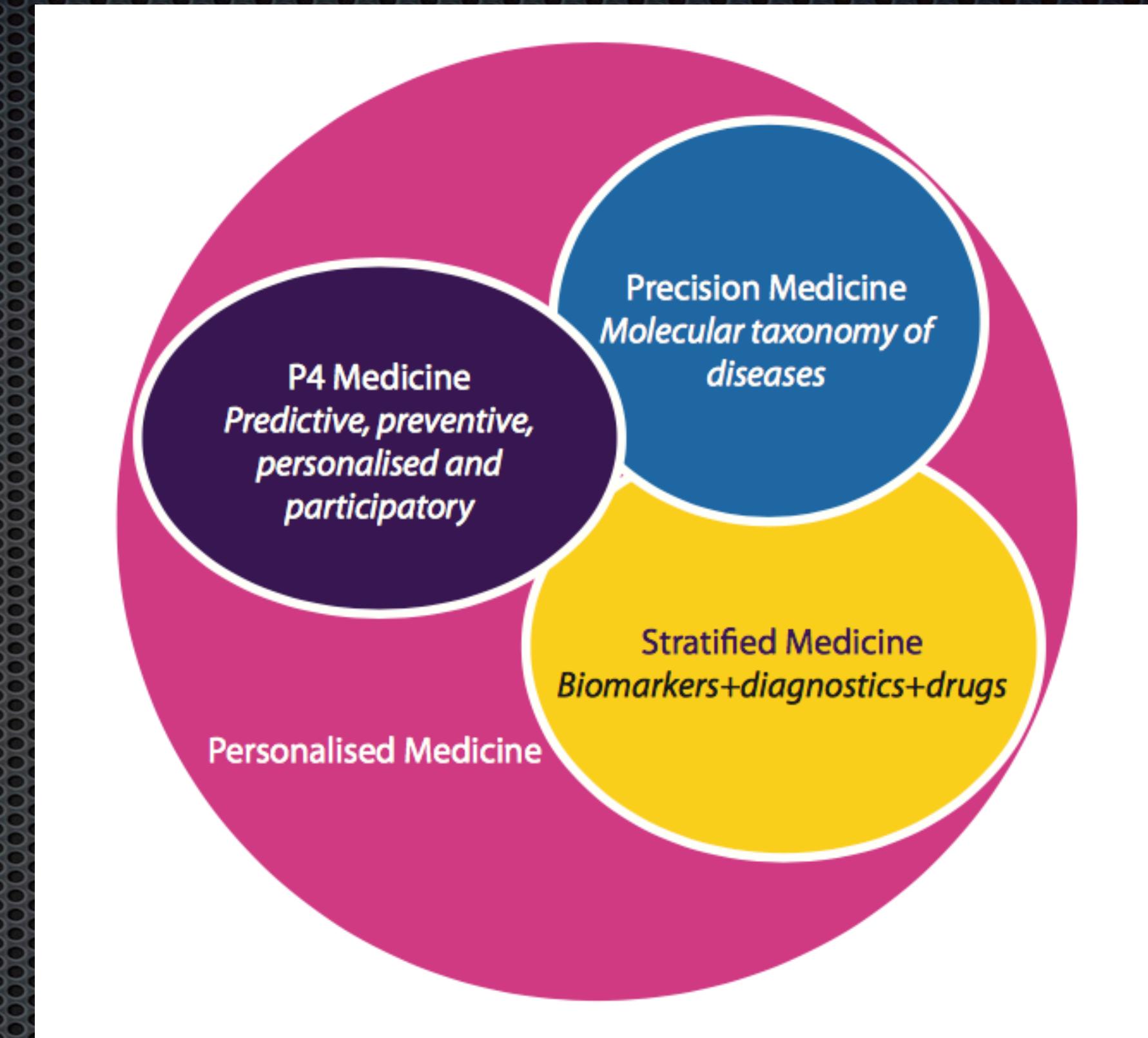
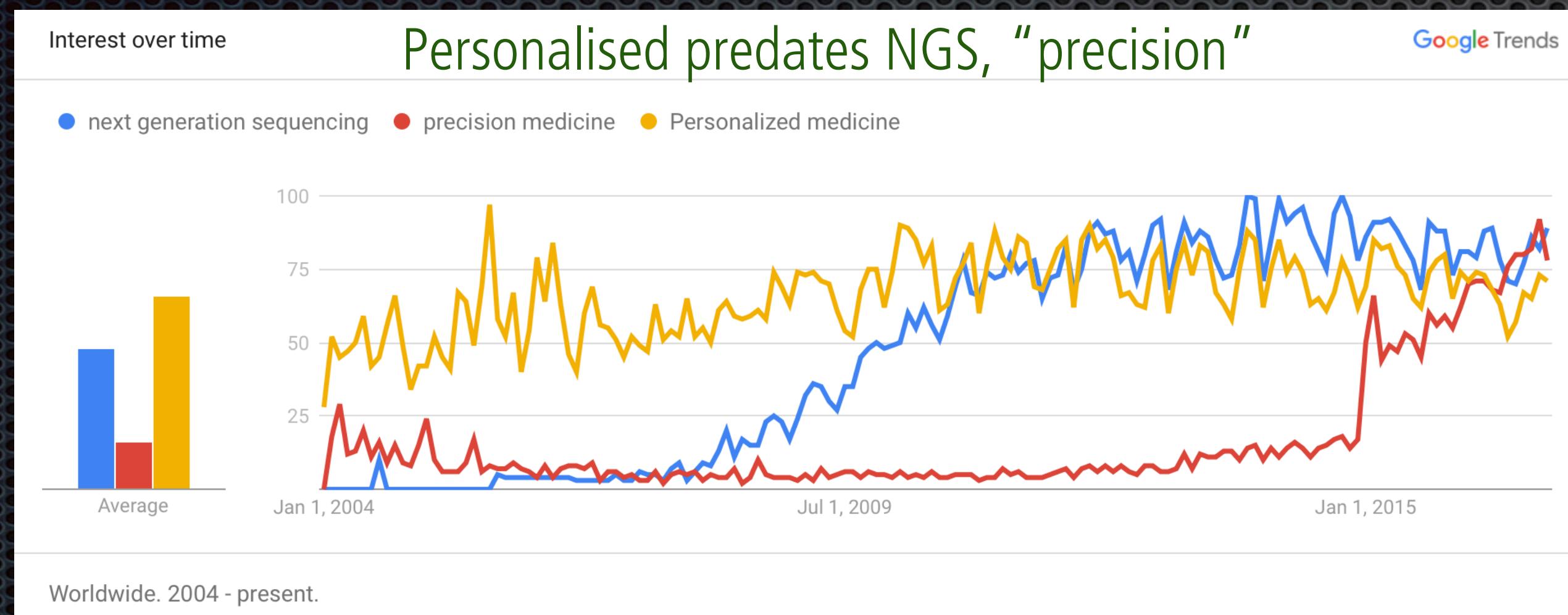
Research group leader in genetics | genomic array analysis for germline alterations | descriptive analysis of copy number aberration patterns in cancer entities

Professor of bioinformatics @ IMLS (2015) | systematic assembly of oncogenomic data | databases and software tools | patterns in cancer genomes | *arraymap* online resource | GA4GH | SPHN

Genomic Background + Disease Parameters

Personalised Medicine **Precision Medicine**

...

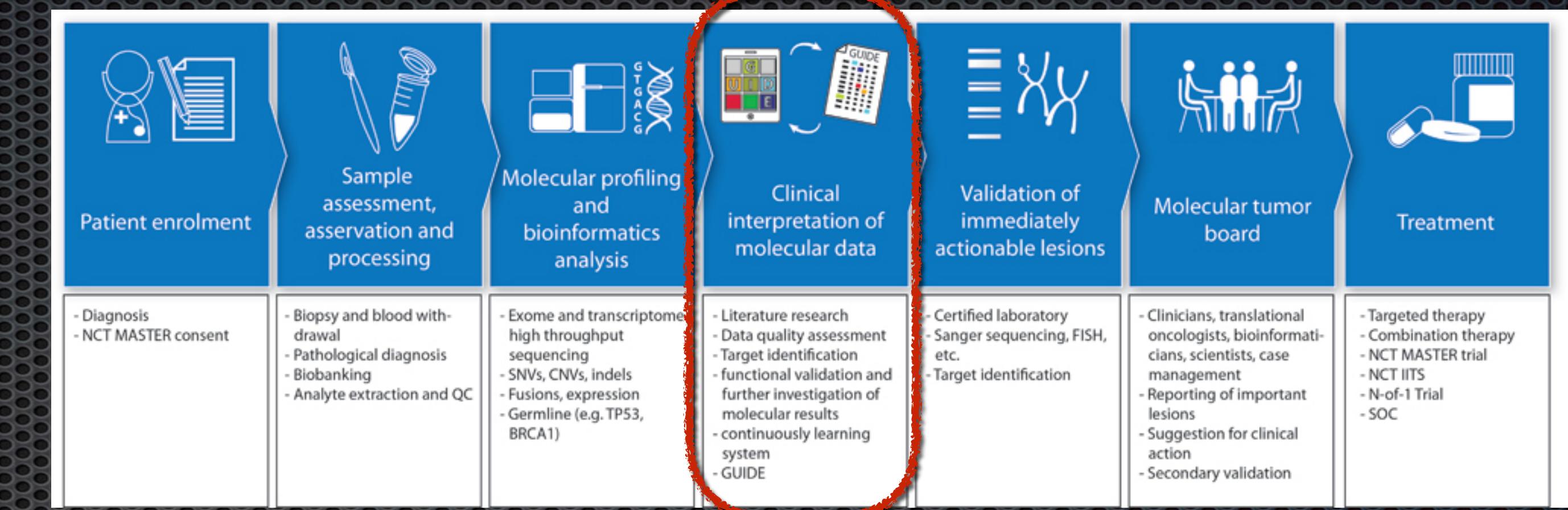


Source: PHG Foundation

While medicine has always been "personal" and "precise" in the given context of available knowledge and technologies, the concept of "**Personalised Medicine**" describes the use of **individual molecular analysis data** and individually **targeted therapies**.

Personalised Medicine in Cancer - A Genome Based Approach

- personalized cancer therapy uses information about the **individual genetic background** and **tumor sequence analysis** for the identification of somatic variants

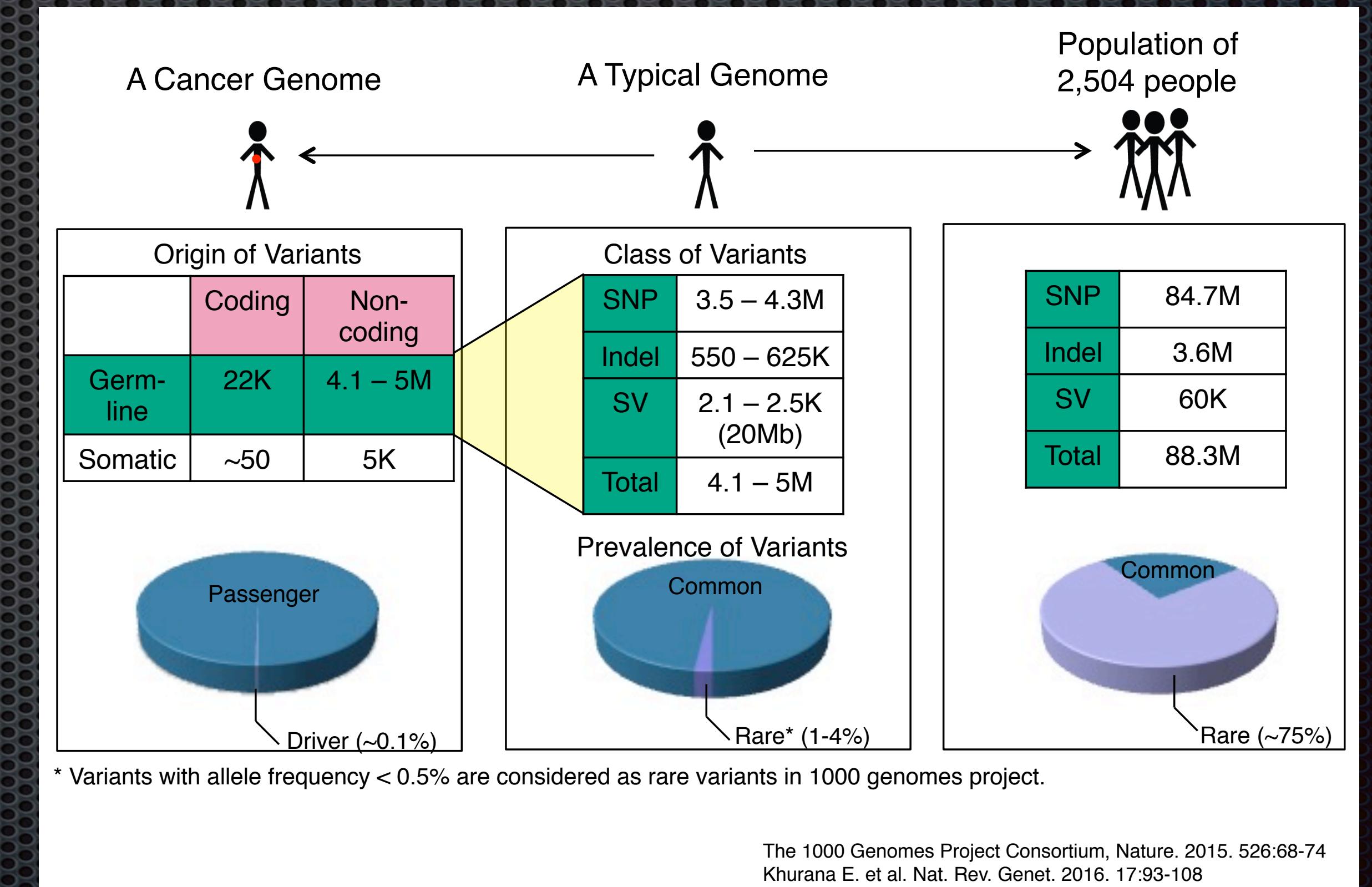


Workflow of a cancer treatment protocol based on "personalized" assessment of actionable genomic lesions (source: NCT Heidelberg).

- currently mostly use of **targeted / panel sequencing** for identification of tens - hundreds of most common "actionable" mutations
- knowledge resources and literature search for interpretation of non-standard variants

Finding Somatic Mutations In Cancer: Many Needles in a Large Haystack

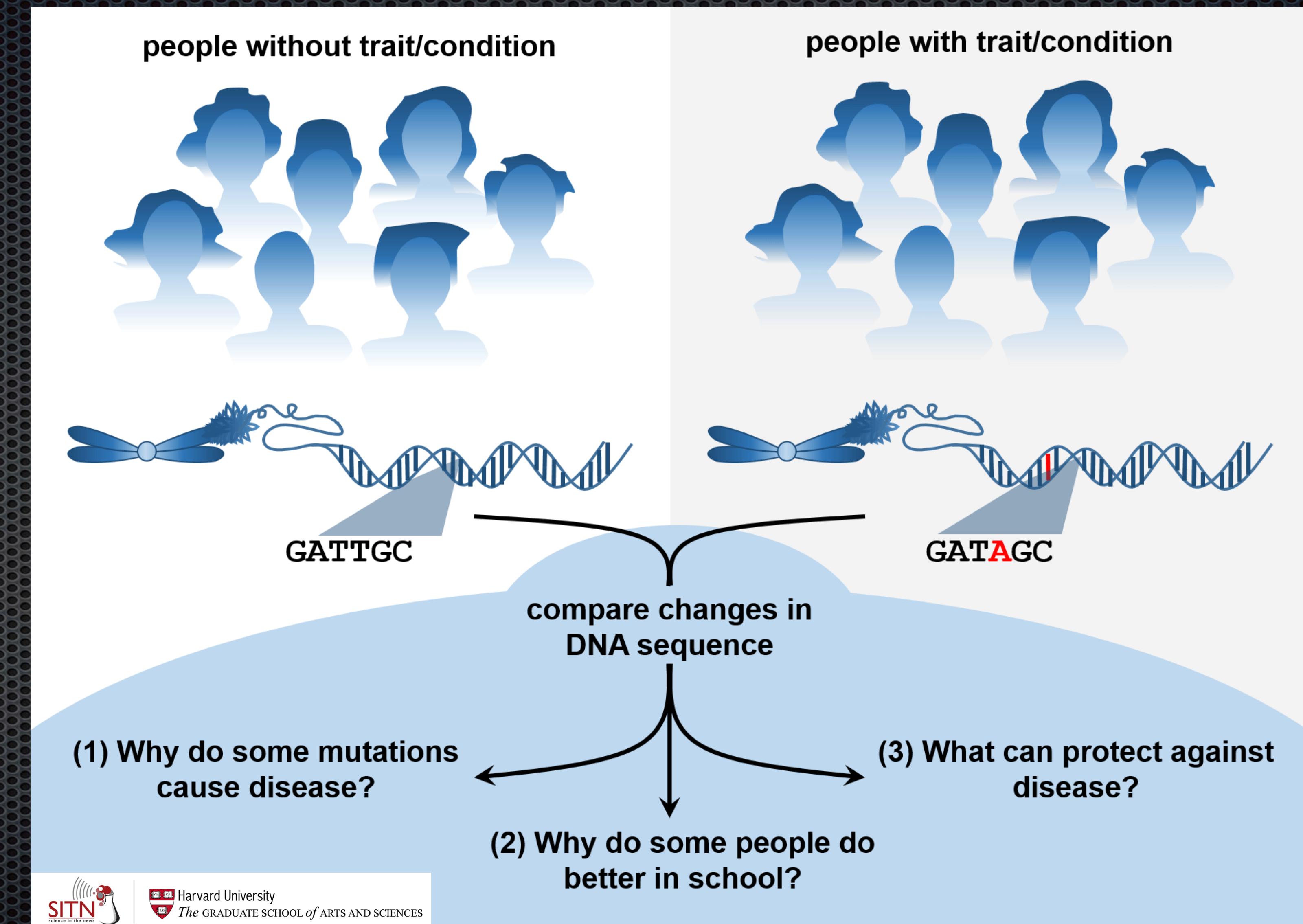
- a typical human genome (~3 billion base pairs) has ~5 million variants
- most of them are "rare"; i.e. can only be identified as recurring when sequencing thousands of people
- cancer cells accumulate additional variants, only **few** of which ("**drivers**") are relevant for the disease



Graphic adapted from Mark Gerstein (GersteinLab.org; @markgerstein)

Population Genomics - The Basic Concept

Studying the genomes of many thousands of people at once, known as population sequencing, is an exciting area of research and is producing countless new insights into human biology and medicine. Typically, a population sequencing experiment will involve sequencing the exomes or genomes of large groups of individuals with and without a trait, such as a disease like cancer, then comparing the differences in the genes of the two groups. Once researchers identify DNA changes associated with the trait of interest, they can then use that information to answer many important questions that help guide drug development, clinical practice, and therapeutic selection.





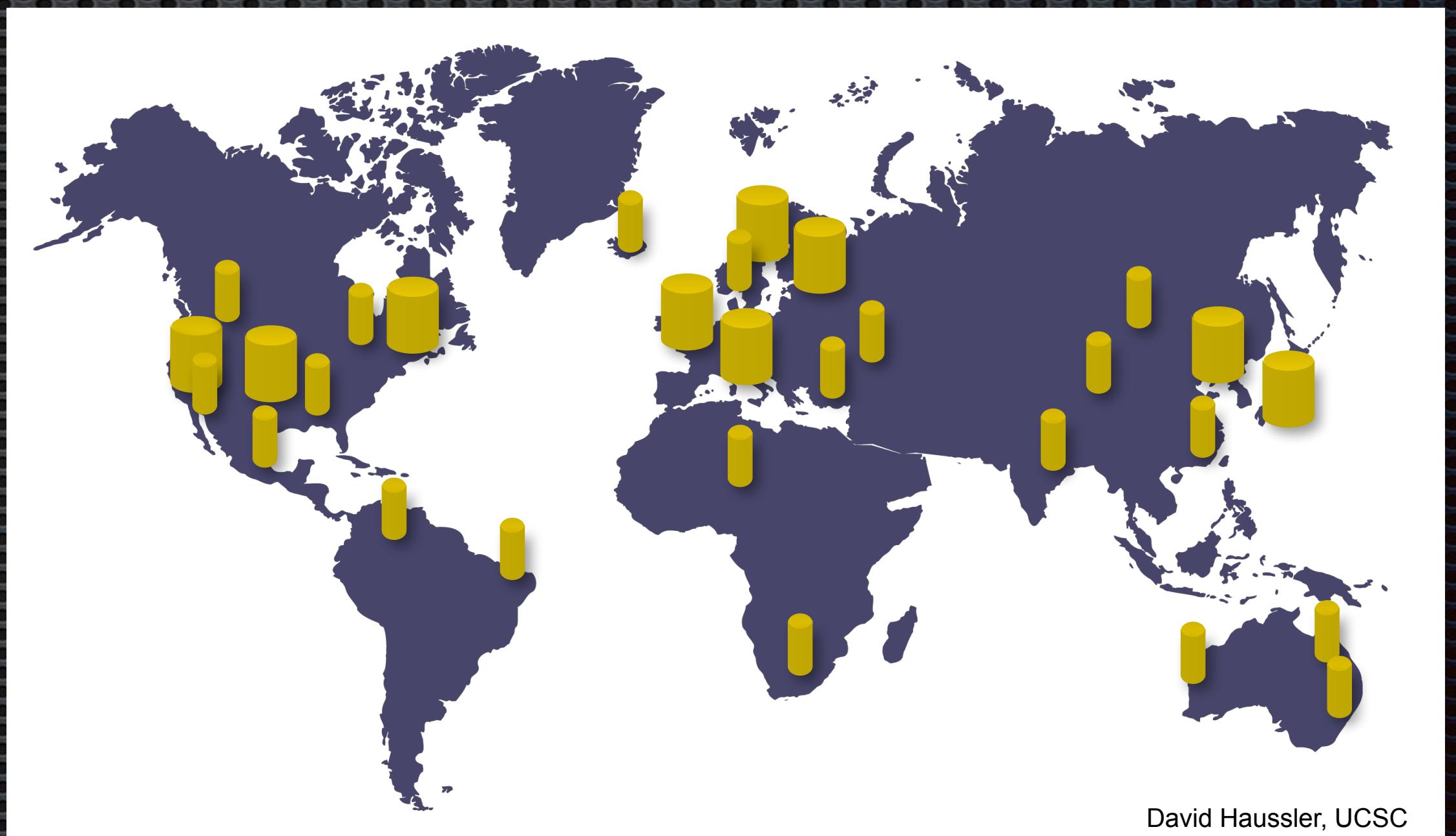
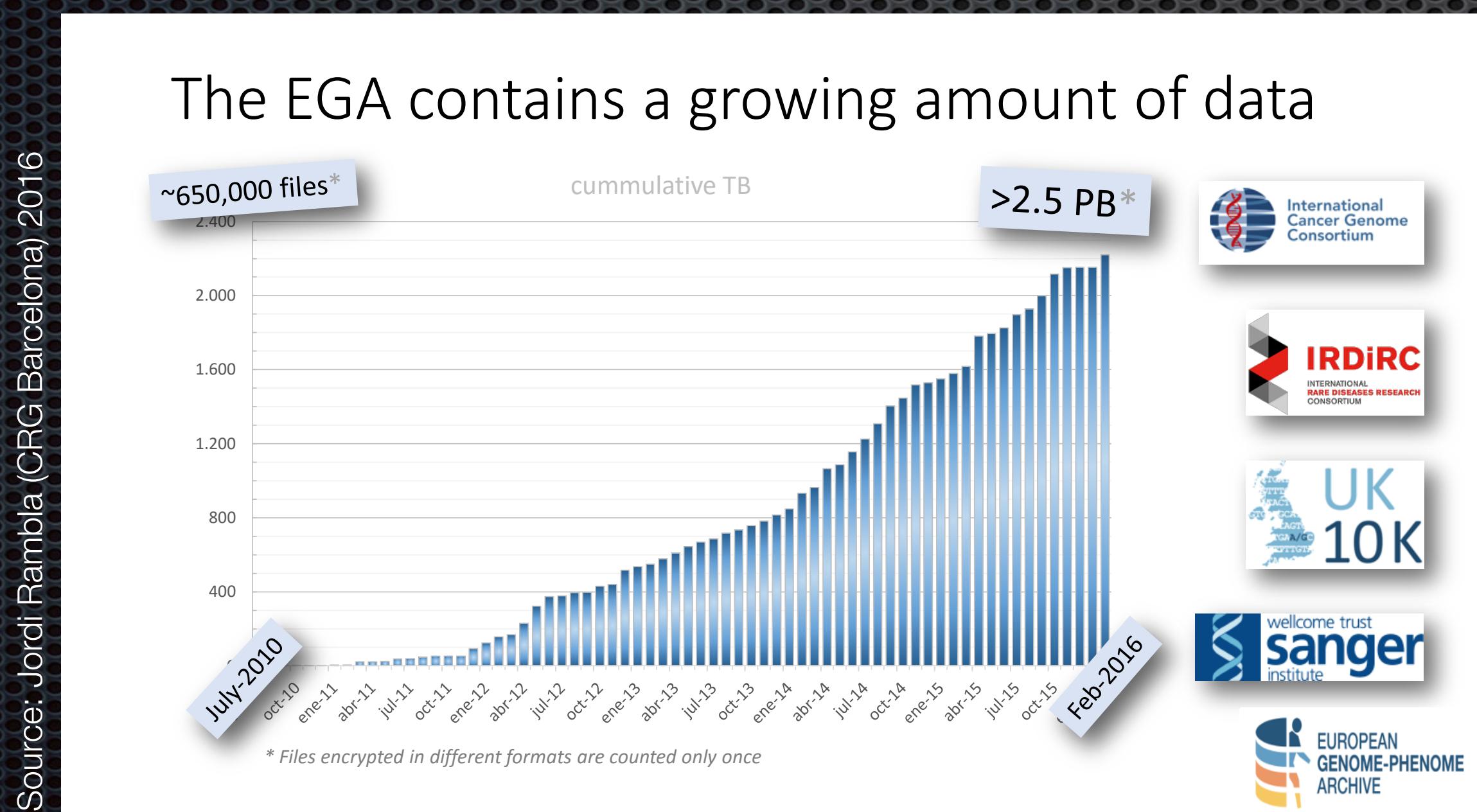
Genomes Everywhere

Large Genome Data Generation, Analysis & Sharing Initiatives

Organization / Initiative: Name	Organization / Initiative: Category	Cohort	
100K Wellness Project	Research Project	107 unaffected individuals (scaling up to 100,000)	>1'000'000
23andMe	Organization	>1 million customers (>80% consented to research)	
Actionable Cancer Genome Initiative (ACGI)	Data-Sharing Project	Goal: 100,000 individuals	
Ancestry.com	Organization	1.4 million customer DNA samples (what % consented to research?)	100'000
BioBank Japan	Repository	Specimens from >200,000 patients and unaffected controls	20'000+
Cancer Moonshot2020	Consortium	Phase 1: 20,000 cancer patients	
Children's Hospital of Philadelphia Biorepository	Repository	Capacity for 8.6 million samples	
China Kadoorie Biobank	Repository	>512,000 participants (general population, China). Genotyping data available for ~100,000.	
CIMBA	Consortium	>15,000 BRCA1 carriers, >8,000 BRCA2 carriers	
Clinical Sequencing Exploratory Research (CSER)	Consortium	~4,000 patients and healthy controls	
DECIPHER	Repository	19,014 patients (international)	
deCode Genetics	Organization	500,000 participants (international)	
East London Genes & Health	Research Project	100,000 unaffected individuals (East London, Pakistani or Bangladeshi heritage)	
Electronic Medical Records and Genomics (eMERGE) Network	Repository, Consortium, Research Project	55,028 patients	
European Network for Genetic and Genomic Epidemiology (ENGAGE)	Research Project	80,000 GWAS scans, and DNA and serum/plasma from >600,000 individuals	
Exome Aggregation Consortium (ExAC)	Consortium	60,706 individuals	
GENIE/AACR	Data-Sharing Project	>17,000 cancer patients (international)	17'000+
Genome Asia 100K	Consortium	Goal: 100,000 individuals (Asia)	
Genomics England	Organization	Goal: 100,000 genomes from 70,000 individuals (rare disease & cancer patients, and their relatives)	
GoT2D	Consortium, Data-Sharing Project	Multiple case-control cohorts	
International Cancer Genome Consortium (ICGC)	Consortium	PCAWG currently data from >8'000 genomes	8'000+
International Genomics of Alzheimer's Project (IGAP)	Consortium	40,000 patients with Alzheimer's disease	
International Multiple Sclerosis Genetics (IMSG) Consortium	Consortium	Goal: >50,000 patients with MS	
Kaiser Permanente: Genes, Environment, and Health (RPGEH)	Repository, Research Project	200,000 DNA samples (scaling up to 500,000)	>2-500'000
Leiden Open Variation Database (LOVD)	Repository	>170,000 individuals	
Million Veteran Program	Research Project	Goal: 1 million individuals; first 200,000 is complete.	
MyCode® Community Health Initiative	Repository, Research Project	Goal: >250,000 patients	
Precision Medicine Initiative	Research Project	Goal: >1 million participants, starting in 2016 (US)	>1'000'000
Psychiatric Genomics Consortium (PGC)	Consortium	>170,000 subjects	
Resilience Project	Research Project	589,306 individuals	
Saudi Human Genome Program	Research Project	Goal: ~100,000 patients and controls (Saudi Arabia)	
Scottish Genomes Partnership (SGP)	Research Project	>3,000 individuals (Scotland)	
T2D-GENES	Consortium, Data-Sharing Project	10,000 patients and controls (five ethnicities); 600 individuals (Mexican American)	
TBResist	Consortium	>2,600 samples	
UK Biobank	Repository, Consortium, Research Project	500,000 individuals (age 40-69 years; UK)	500'000
UK10K	Research Project	10,000 participants (6,000 patients and 4,000 controls)	
Vanderbilt's BioVU	Repository	>215,000 samples	

Genome Datasets: Rapid Growth, Limited Access

population based and cancer research studies produce a rapidly increasing amount of genome sequence data



genome data is stored in an increasing number of institutional and core repositories, with **incompatible data** structures and **access** policies

GA4GH to solve genome
data access....



Enabling genomic data sharing for the benefit of human health

The Global Alliance for Genomics and Health (GA4GH) is a policy-framing and technical standards-setting organization, seeking to enable responsible genomic data sharing within a **human rights framework**



**Genomic Data
Toolkit**



**Regulatory & Ethics
Toolkit**



**Data Security
Toolkit**



[VIEW OUR LEADERSHIP](#)

[MORE ABOUT US](#)

[BECOME A MEMBER](#)

GA4GH HISTORY & MILESTONES

- January 2013 - 50 participants from eight countries
- June 2013 - White Paper, over next year signed by 70 “founding” member institutions (e.g. SIB, UZH)
- March 2014 - Working group meeting in Hinxton & 1st plenary in London
- October 2014 - Plenary meeting, San Diego; interaction with ASHG meeting
- June 2015 - 3rd Plenary meeting, Leiden
- September 2015 - GA4GH at ASHG, Baltimore
- October 2015 - DWG / New York Genome Centre
- April 2016 - Global Workshop @ ICHG 2016, Kyoto
- October 2016 - 4th Plenary Meeting, Vancouver
- May 2017 - Strategy retreat, Hinxton
- October 2017 - 5th plenary, Orlando
- May 2018 - Vancouver
- October 2018 - 6th plenary, Basel

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 6291



Global Alliance
for Genomics & Health

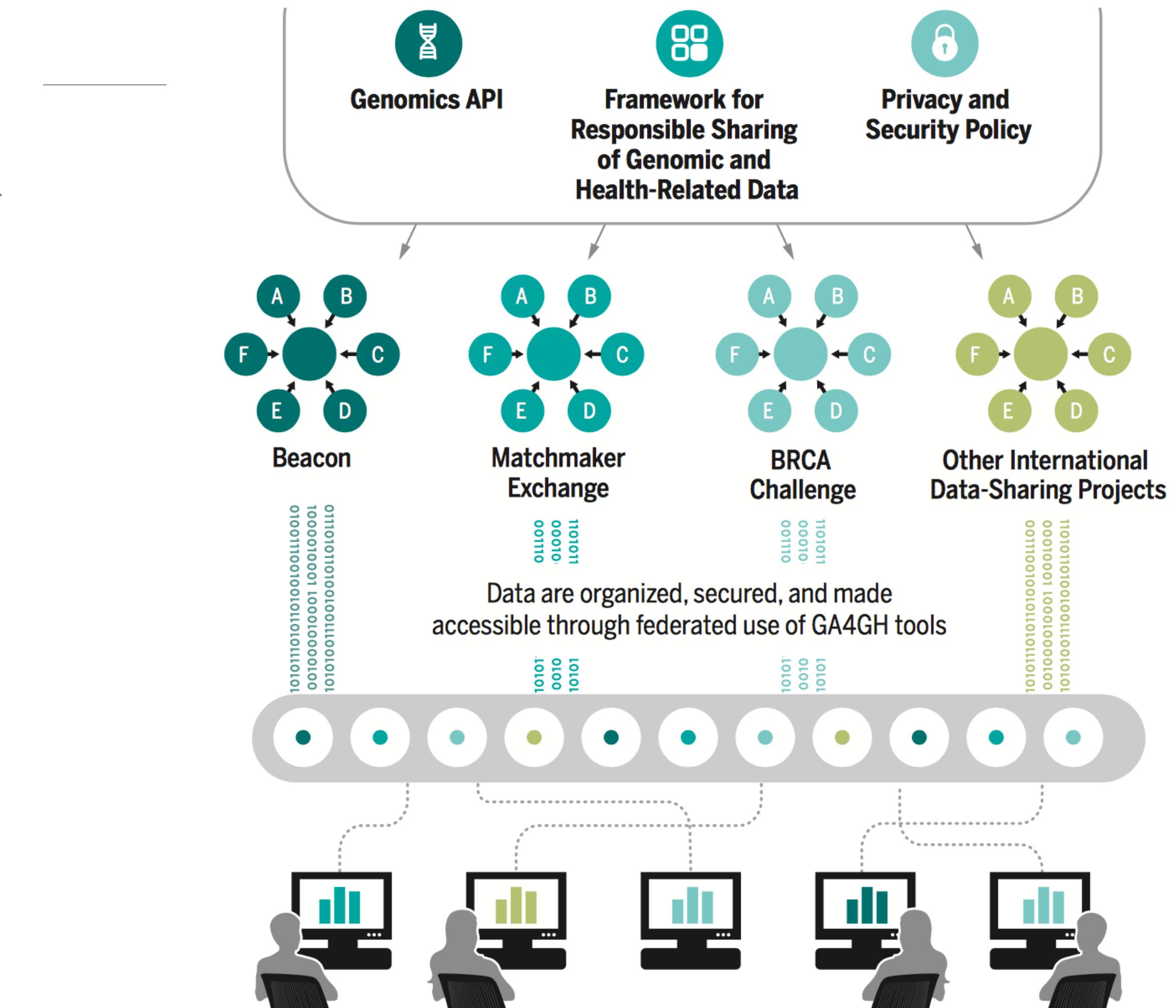
GA4GH API promotes sharing



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.





		Real-World Driver Projects								Partner Engagement
		Project 1	Project 2	Project 3	Project 4	Project 5	Project 6	Project 7	Project 8	
Technical Work Streams	Discovery	✓		✓		✓		✓		
	Large-Scale Genomics		✓		✓		✓		✓	
	Data Use & Researcher IDs	✓		✓		✓		✓		
	Cloud		✓	✓					✓	
	Genomic Knowledge Standards		✓					✓	✓	
	Clinical & Phenotypic Data Capture	✓			✓	✓		✓		
	Regulatory & Ethics									
Foundational Work Streams	Data Security									

GA4GH Driver Projects

BRCA Challenge

The BRCA Challenge aims to advance understanding of the genetic basis of breast and other cancers using data from around the world.



Beacon Project

Beacon Project is an open web service that tests the willingness of international sites to share genetic data. It is being implemented on the websites of the world's top genomic research organizations.



Matchmaker Exchange

Matchmaker Exchange is a federated network of databases whose goal is to find genetic causes of rare diseases by matching similar phenotypic and genotypic profiles.



Matchmaker
Exchange

Beacon Project

An open web service that tests the willingness of international sites to share genetic data.



Beacon Network

Search Beacons

Search all beacons for allele

GRCh37 ▾ 10:118969015 C / CT Search

Response All None
 Found 16
 Not Found 27
 Not Applicable 22

Organization All None
 AMPLab, UC Berkeley
 BGI
 BioReference Laborato...
 Brazilian Initiative on ...
 BRCA Exchange
 Broad Institute
 Centre for Genomic R...
 Centro Nacional de A...
 Curoverse
 EMBL European Bio...
 Global Alliance for G...
 Google
 Institute for Systems ...
 Instituto Nacional de ...

Response	All	None
<input checked="" type="checkbox"/> Found	16	
<input type="checkbox"/> Not Found	27	
<input type="checkbox"/> Not Applicable	22	

BioReference BioReference Hosted by BioReference Laboratories Found

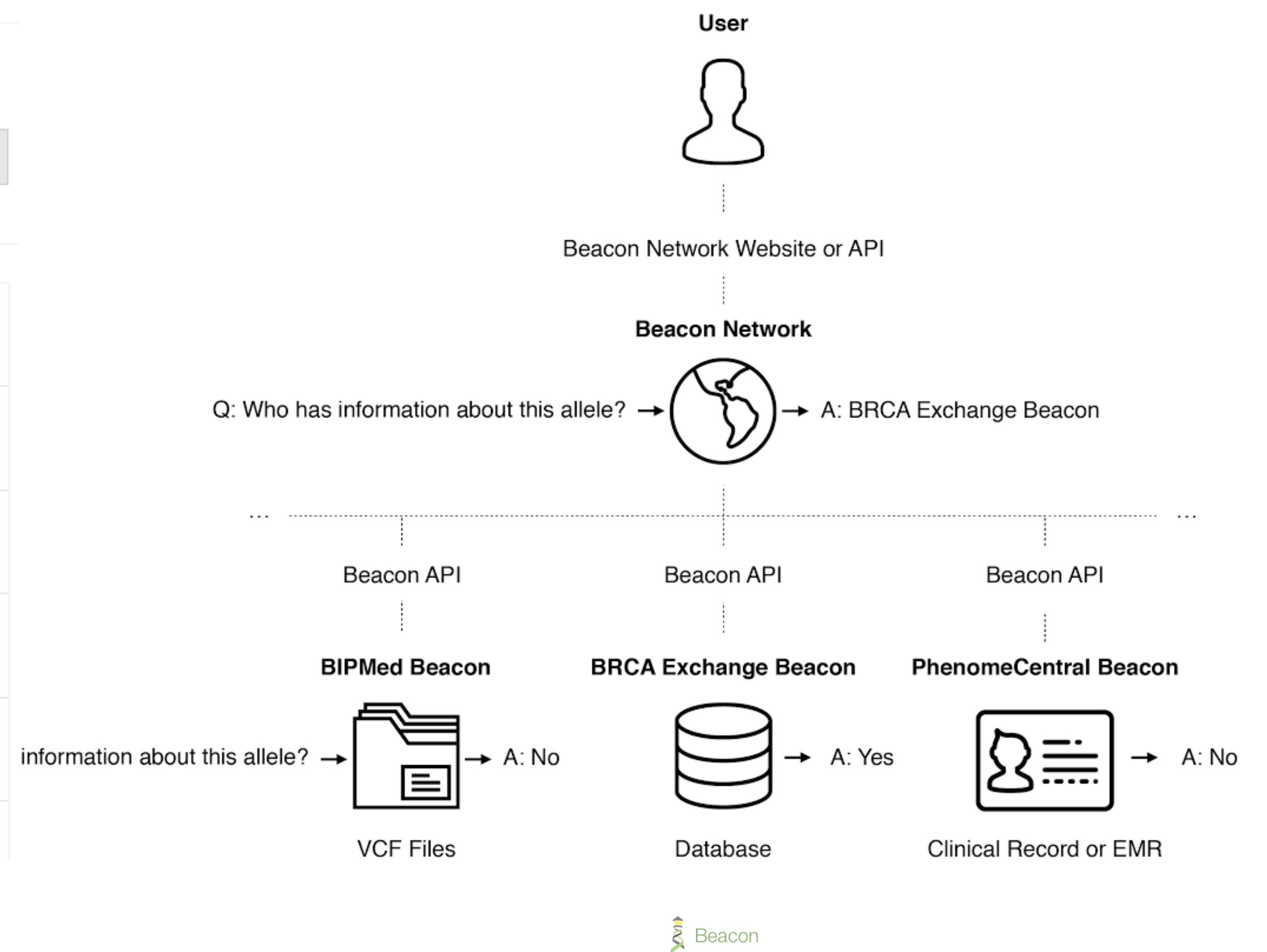
Catalogue of Somatic Mutations in Cancer Catalogue of Somatic Mutations in Cancer Hosted by Wellcome Trust Sanger Institute Found

Cell Lines Cell Lines Hosted by Wellcome Trust Sanger Institute Found

Conglomerate Conglomerate Hosted by Global Alliance for Genomics and Health Found

COSMIC COSMIC Hosted by Wellcome Trust Sanger Institute Found

dbGaP: Combined GRU Catalog and NHLBI Exome Seq... dbGaP: Combined GRU Catalog and NHLBI Exome Seq... Found



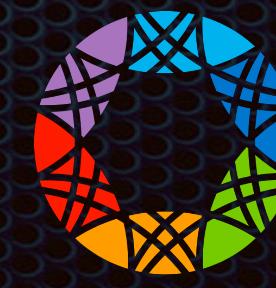
Date	Tag	Title
2018-01-24	v0.4.0	Beacon
2016-05-31	v0.3.0	Beacon

baudisgroup @ UZH & SIB

Computational Cancer Genomics | Data
Sharing | Populations Background | Tools



University of
Zurich UZH

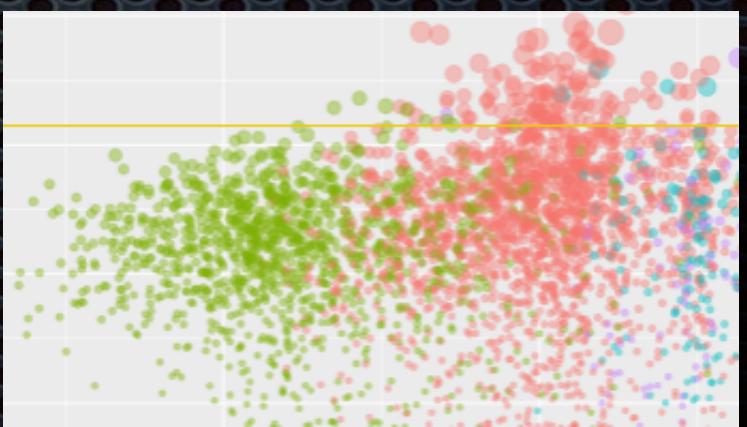
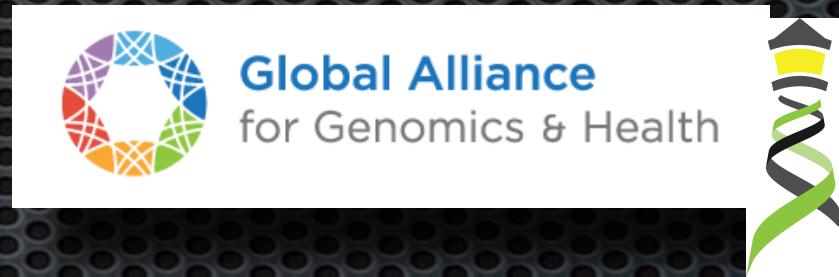
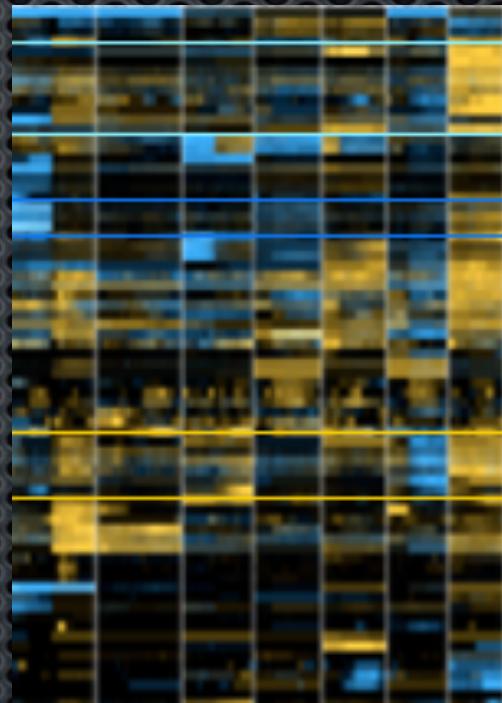
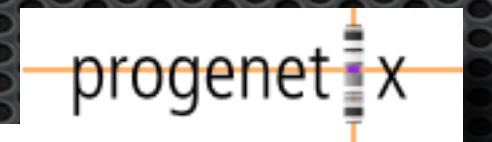


Global Alliance
for Genomics & Health

baudisgroup @ UZH & SIB

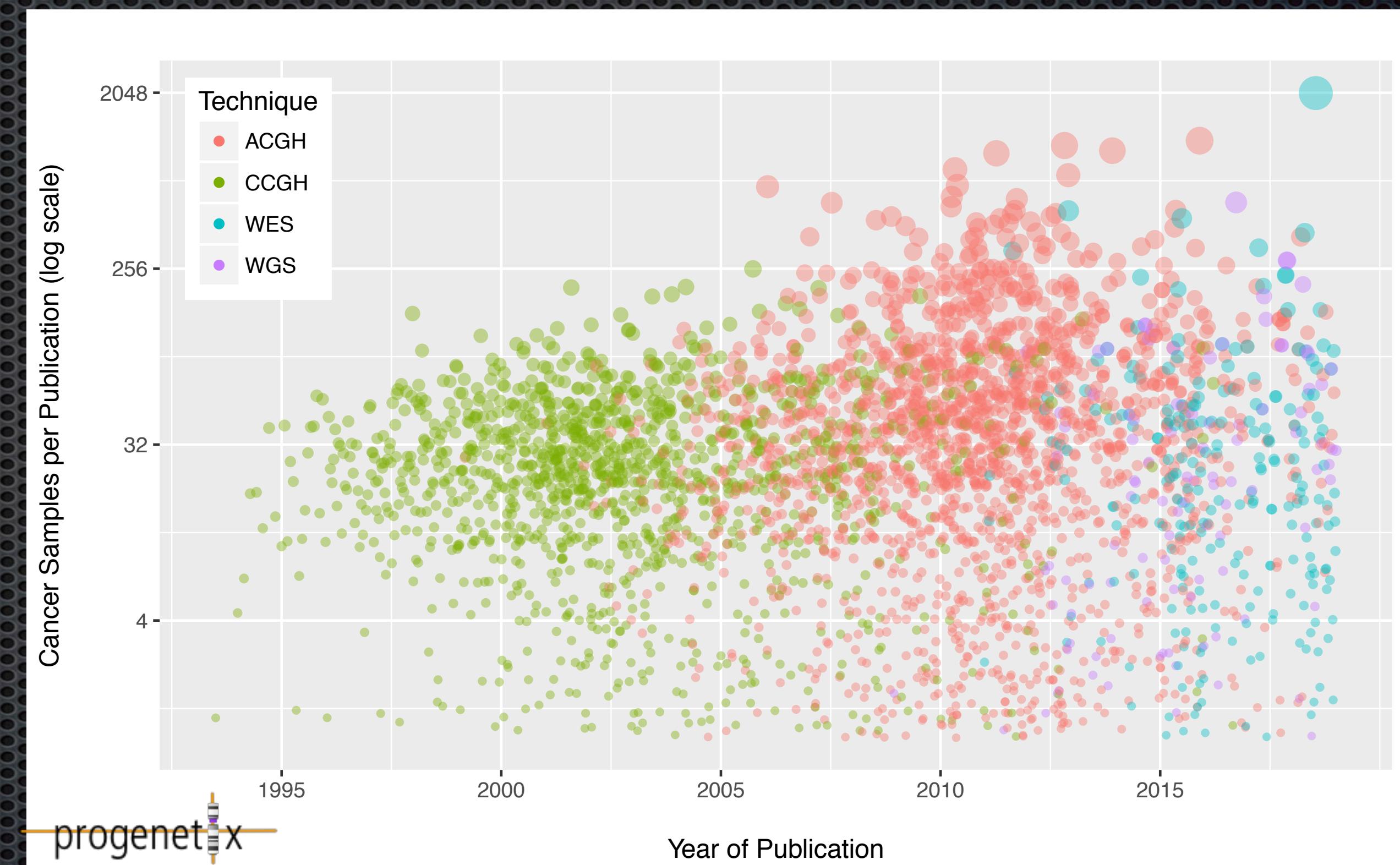
beyond cancer genomes

- Cancer genome data resources
- Software tools for data analysis & visualisation
- Parsing the cancer genome landscape: Patterns & targets
- GA4GH: **Standards & Beacons**
- Quantifying cancer genomic research
- The Swiss Personalized Health Network initiative
- Collaborations!



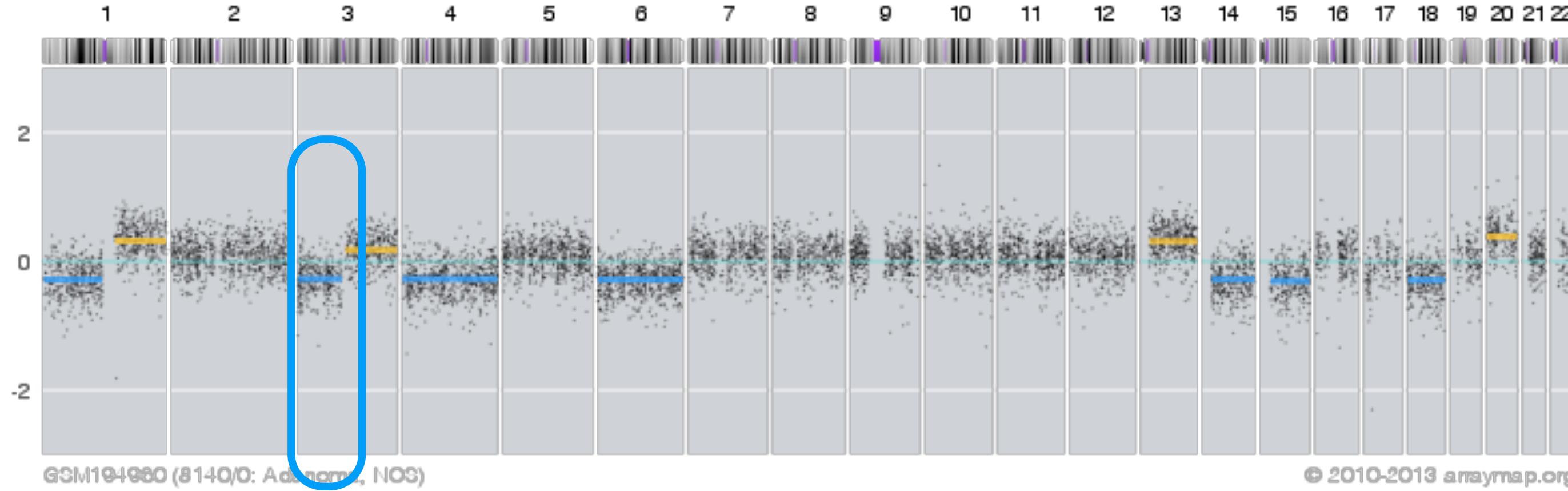
Molecular Cytogenetic & Sequencing Studies for **Whole Genome Profiling** in Cancer

- genome screening to identify mutations in cancer samples
- for diagnostic purposes and therapeutic target identification
 - karyotyping (~1968)
 - Comparative Genomic Hybridization (1992)
 - genome **microarrays** (aCGH, SNP arrays ...; 1997)
 - Whole Exome Sequencing** (2010)
 - Whole Genome Sequencing** (2011)

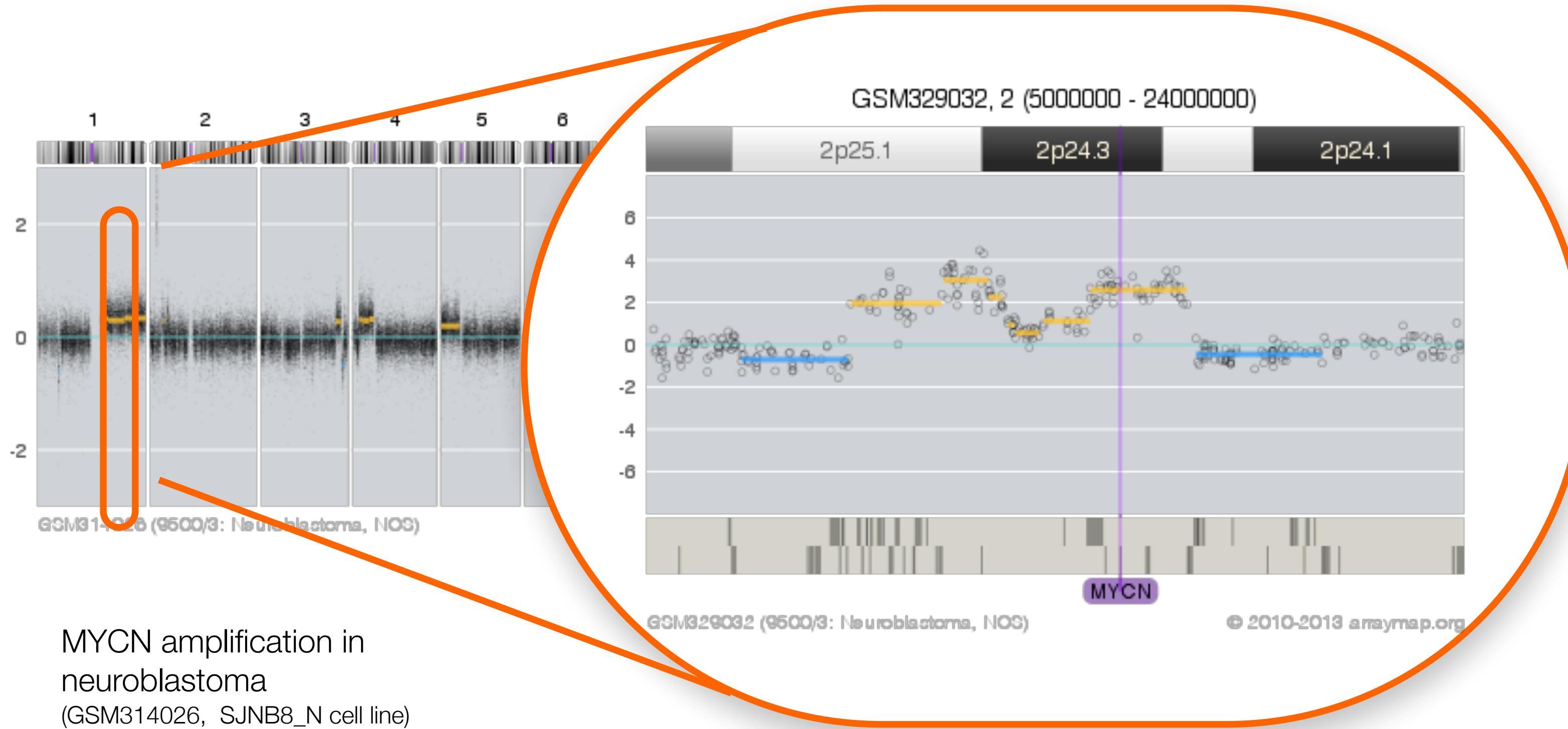


Overview of publications reporting whole-genome screening analysis of cancer samples, by molecular-cytogenetic or genome sequencing methods. The data represents 3357 articles assessed for the progenetix.org cancer genome data resource (M. Baudis, 2001-2019)

Copy number aberrations



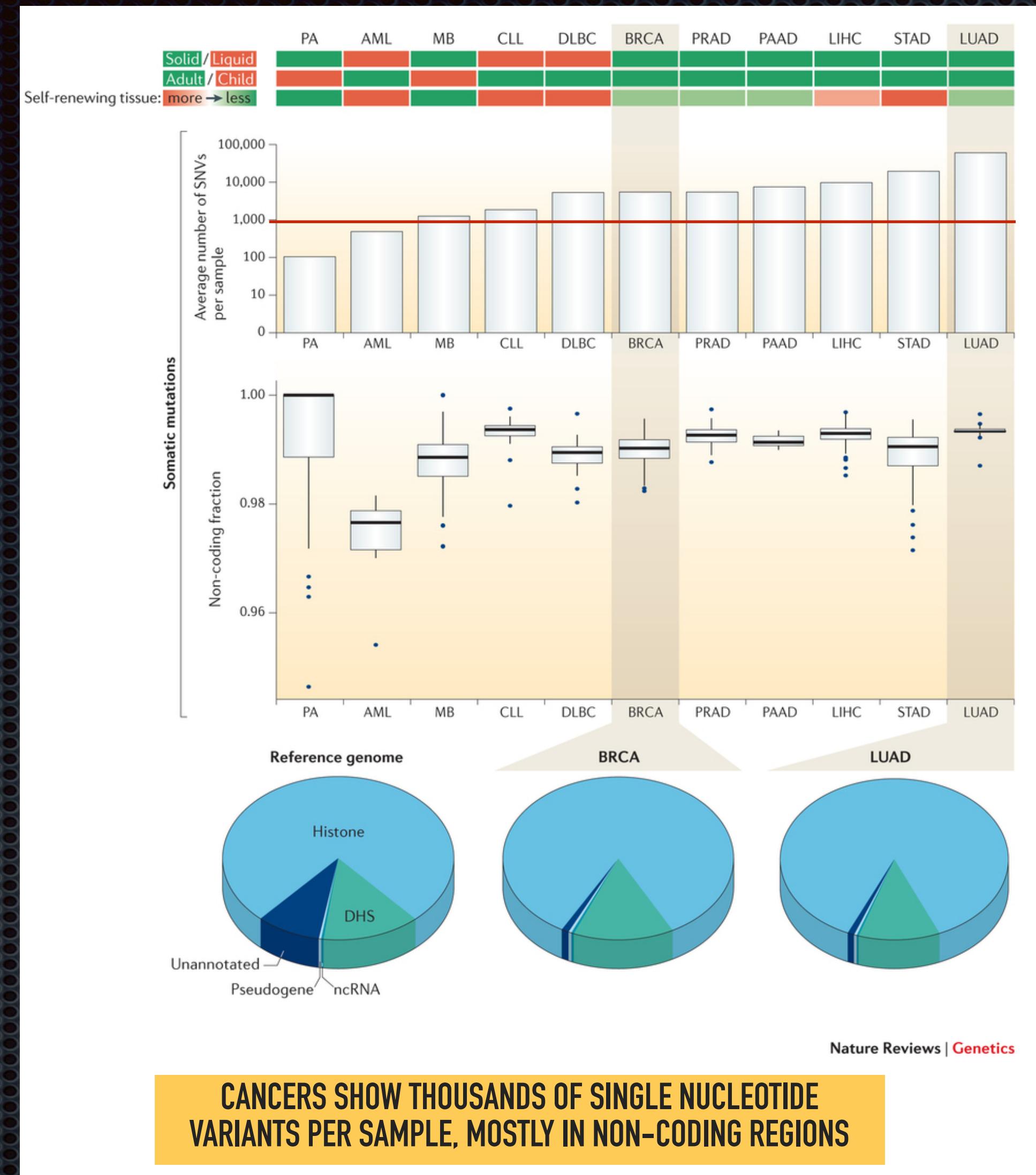
Gain of chromosome arm 3q in colorectal carcinoma



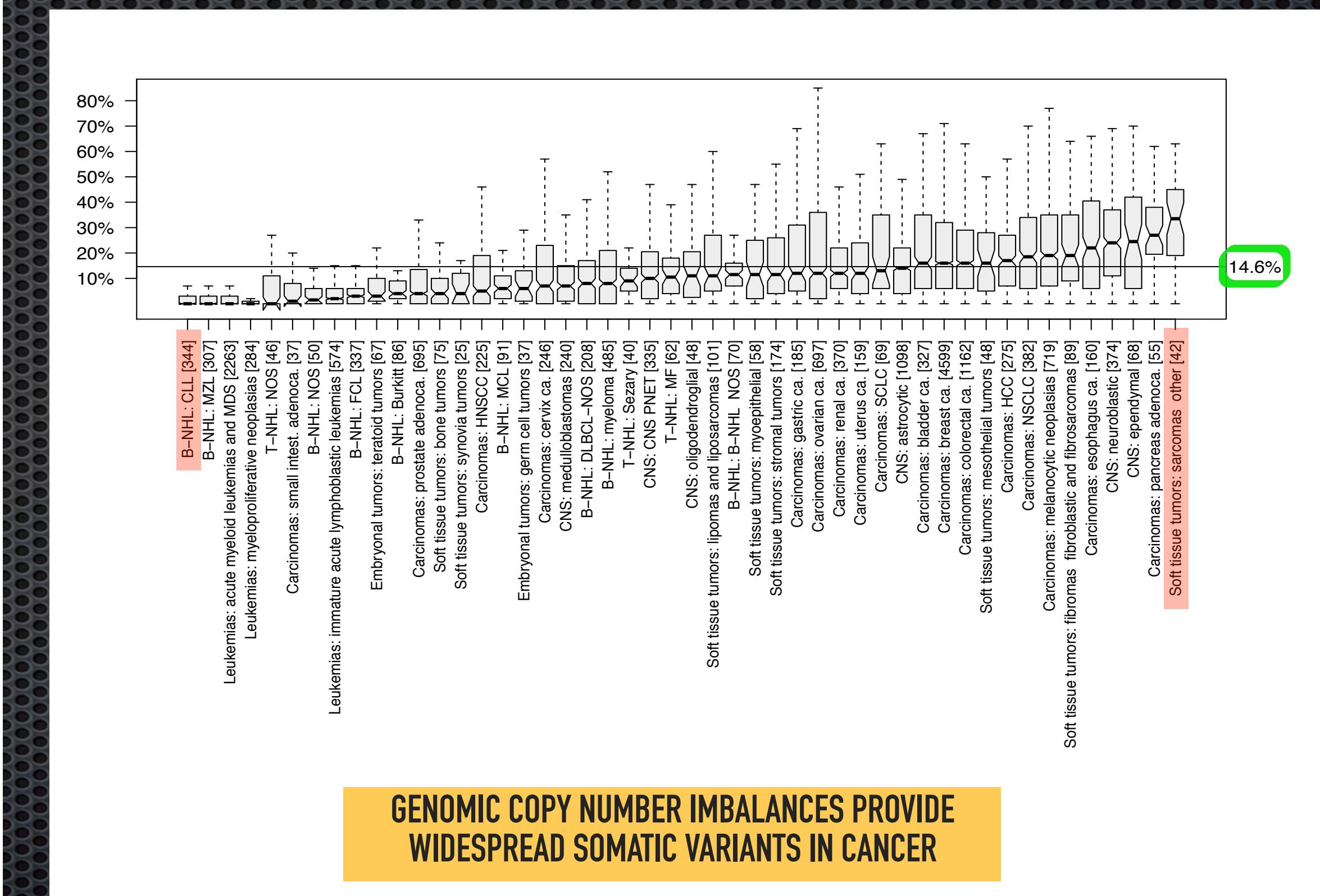
low level/high level copy number alterations (CNAs)

arrayMap





Quantifying Somatic Mutations In Cancer



Pan-Cancer Analysis of Whole Genomes (PCAWG) data show widespread mutations in non-coding regions of cancer genomes (Khurana et al., Nat. Rev. Genet. (2016))

On average ~15% of a cancer genome are in an imbalanced state (more/less than 2 alleles);
Original data based on >30'000 cancer genomes from arraymap.org

Reference Resources for Cancer Genome Profiling

- continuously updated reference resources for cancer genome profiling data and related information
- basis for own research activities, collaborative projects and external use
- structured information serves for implementing GA4GH concepts



arrayMap

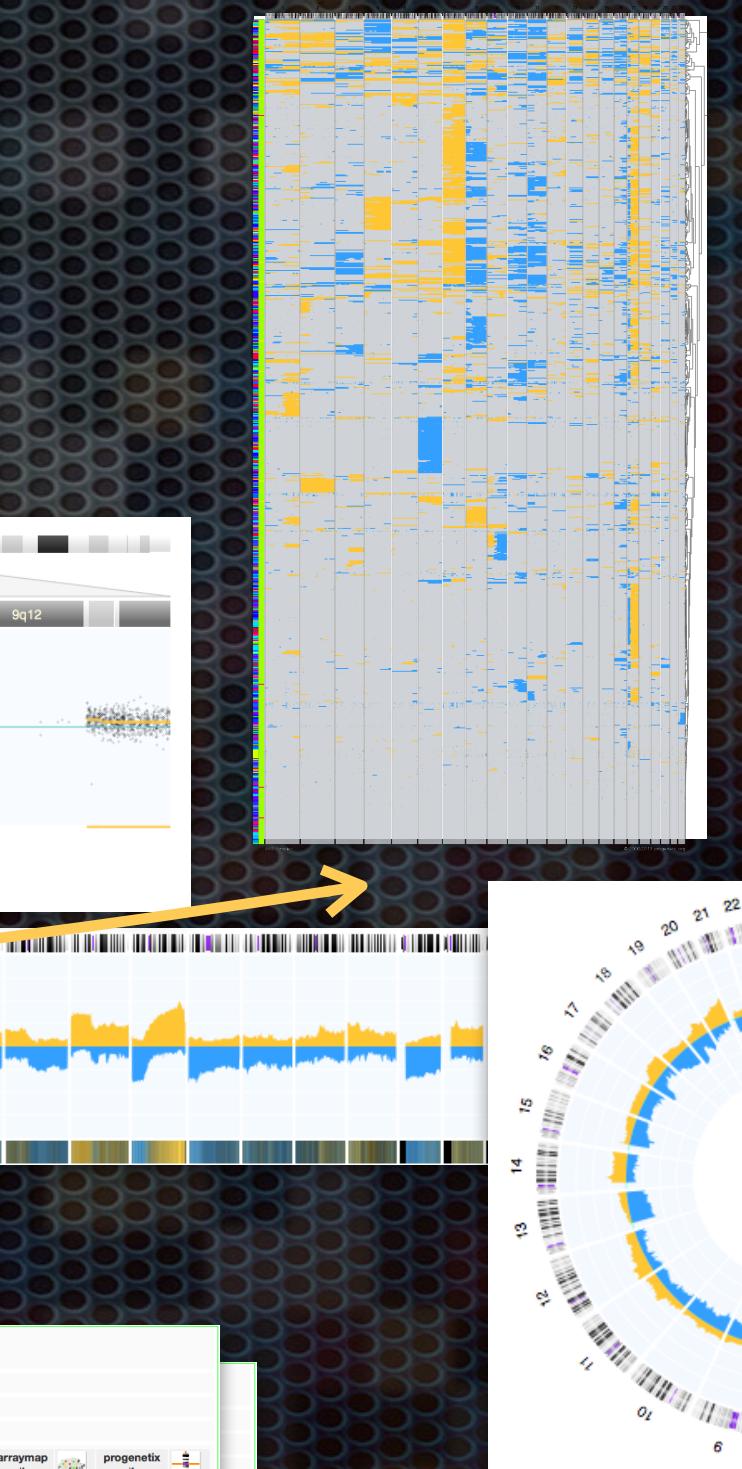
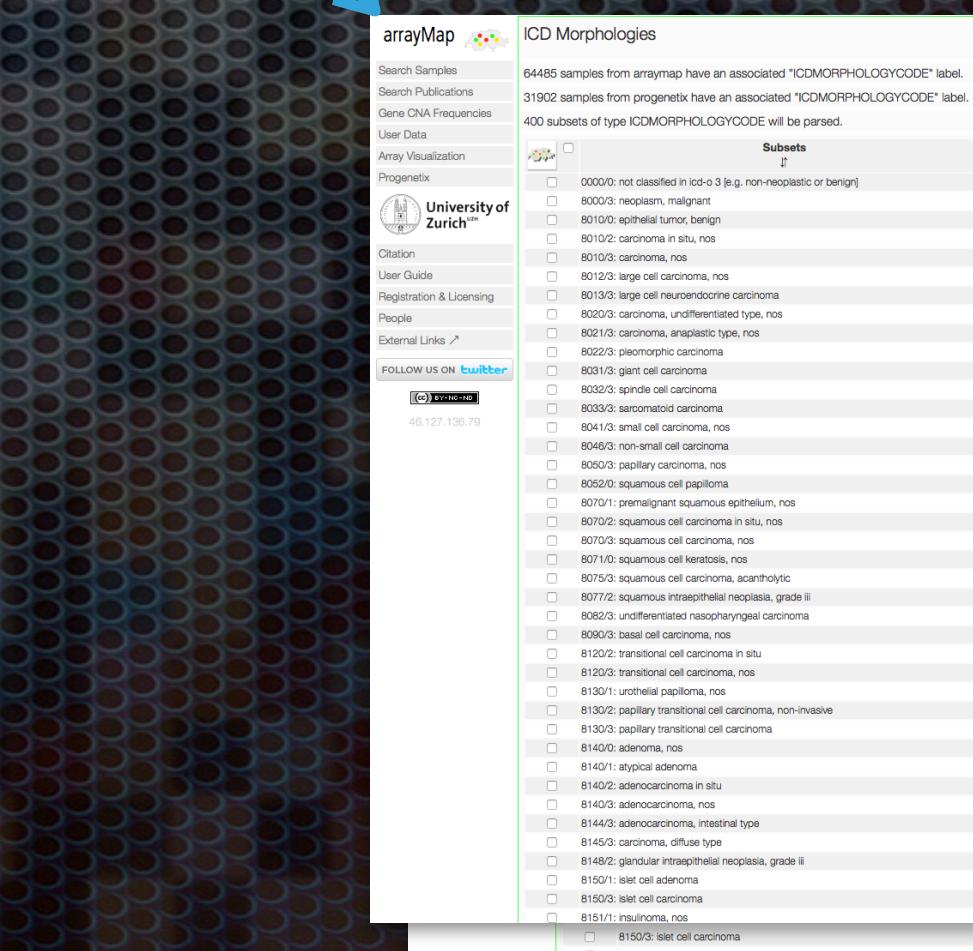
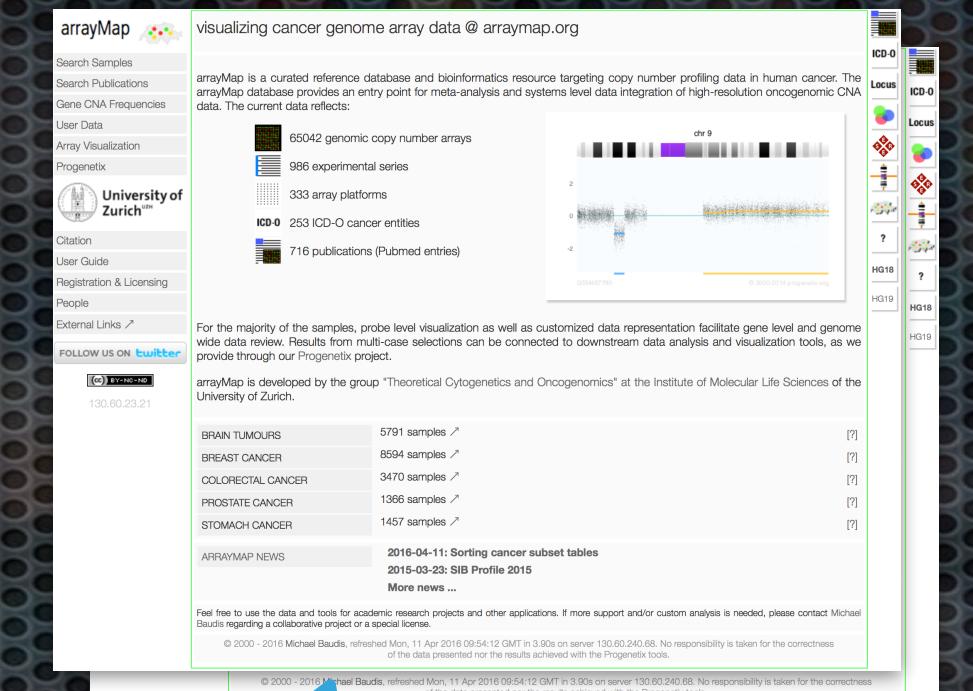
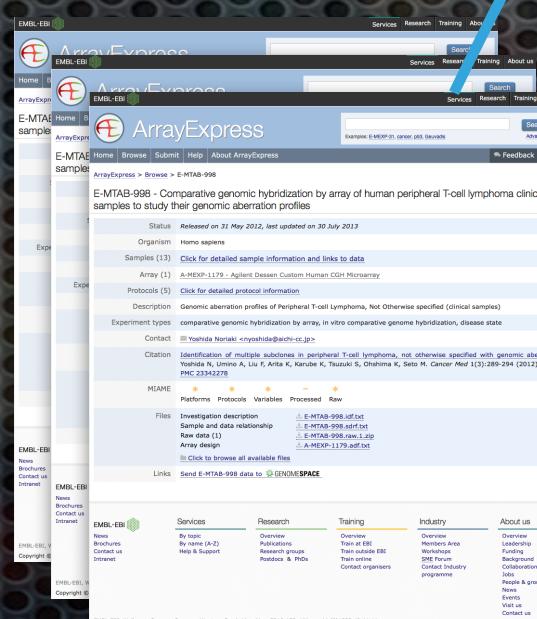
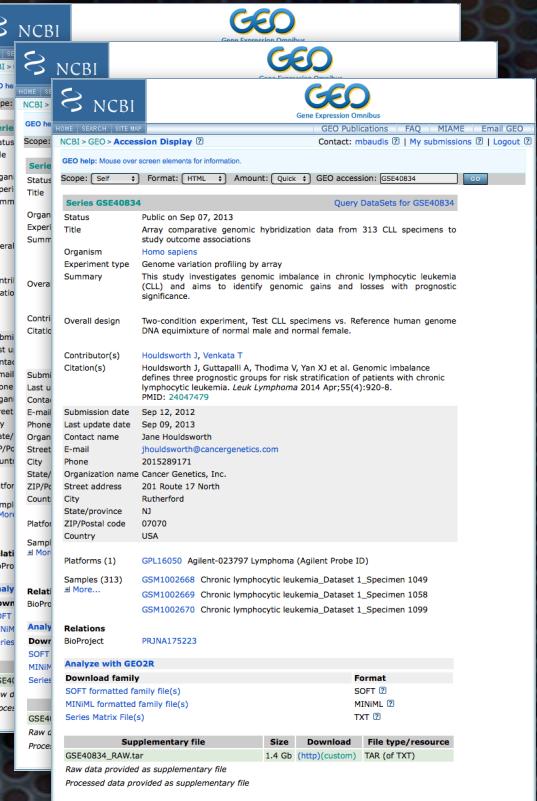
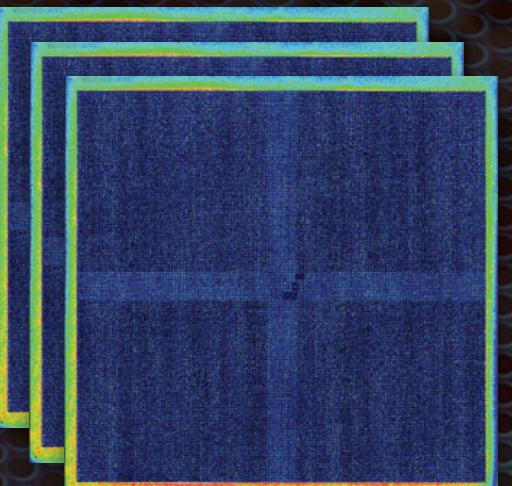


techniques	cCGH, aCGH, WES, WGS	aCGH (+?)
scope	sample (e.g. combination of several experiments)	experiment
content	>31000 samples	>60000 arrays
raw data presentation	no (link to sources if available)	yes (raw, log2, segmentation if available)
per sample re-analysis	no; supervised result (mostly as provided through publication)	yes (re-segmentation, thresholding, size filters ...)
final data	annotated/interpreted CN status for GP and cytogenetic regions	unsupervised CN status for GP and cytogenetic regions
main purposes	<ul style="list-style-type: none">• Distribution of CNA target regions in most tumor types (>350 ICD-O)• Cancer classification	<ul style="list-style-type: none">• Gene specific hits• Genome feature correlation (fragile sites ...)

DATA PIPELINE

Automagical Processing Engine™

arrayMap

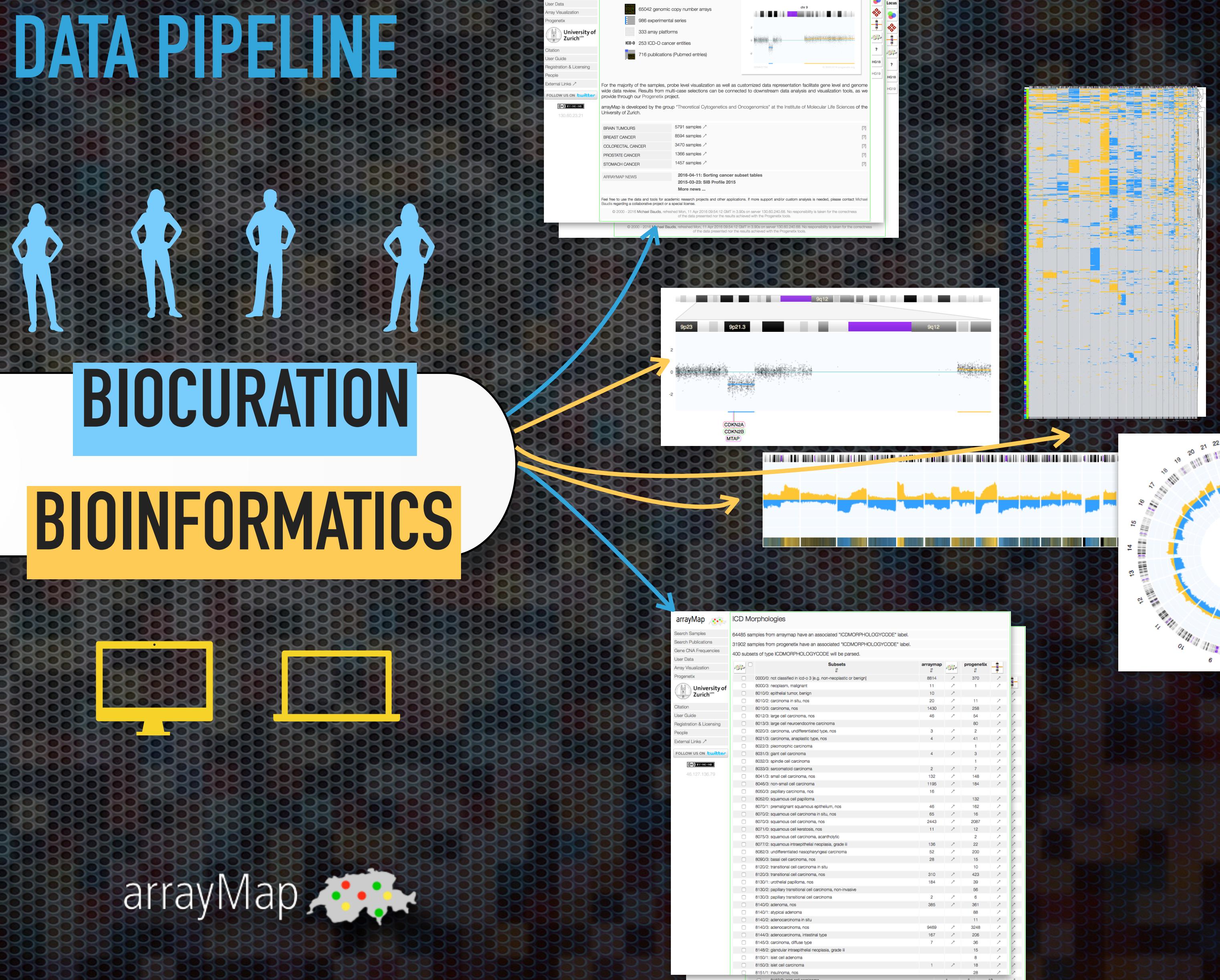
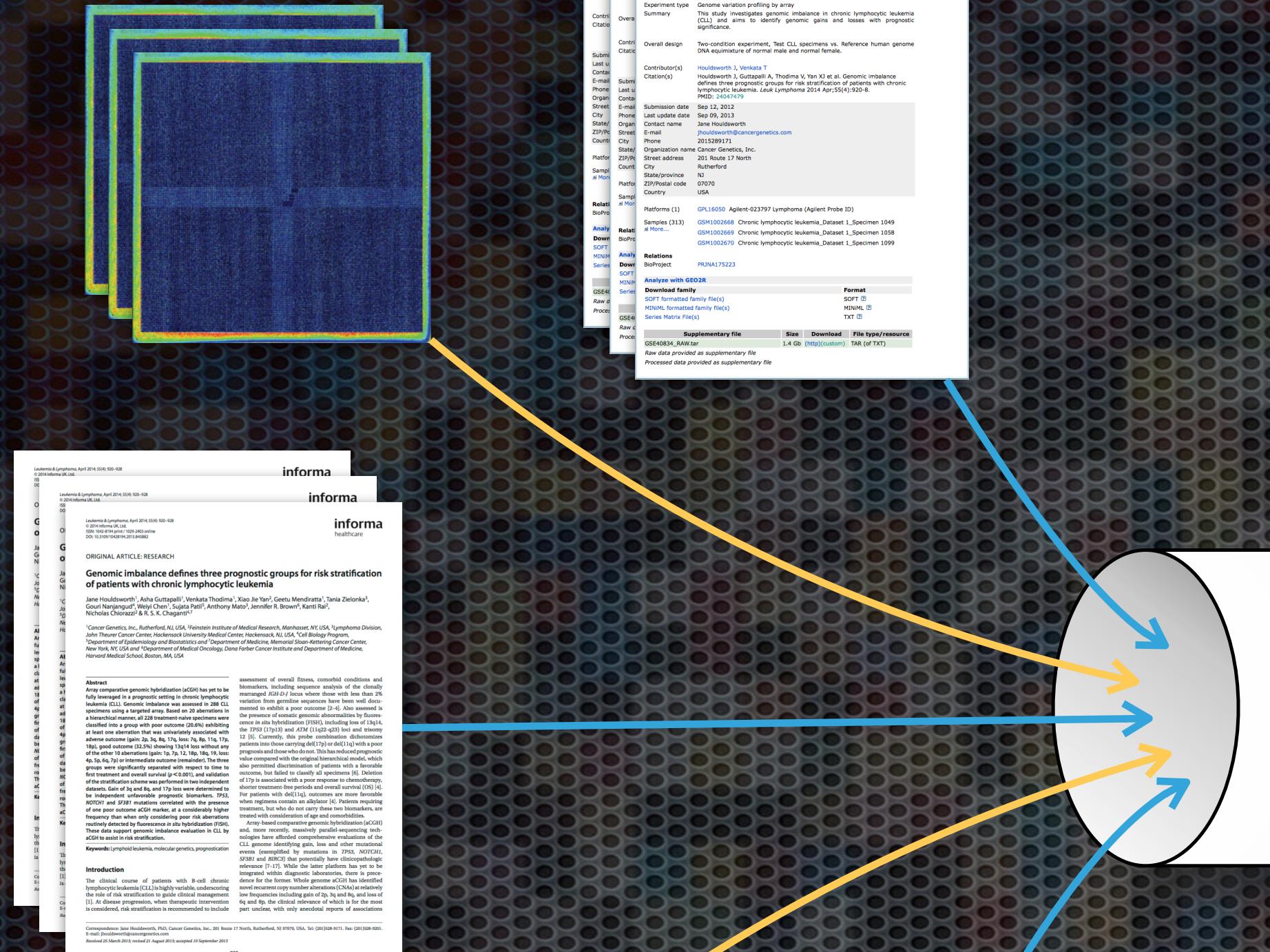


DATA PIPELINE

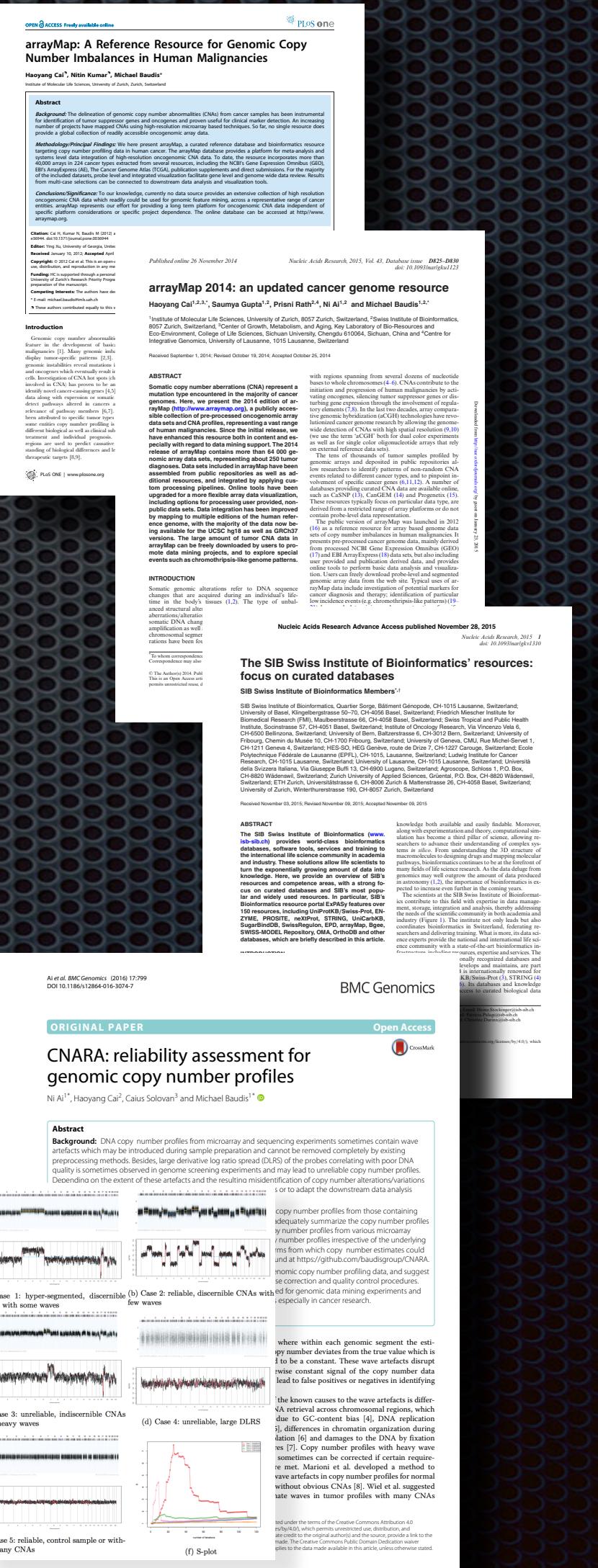
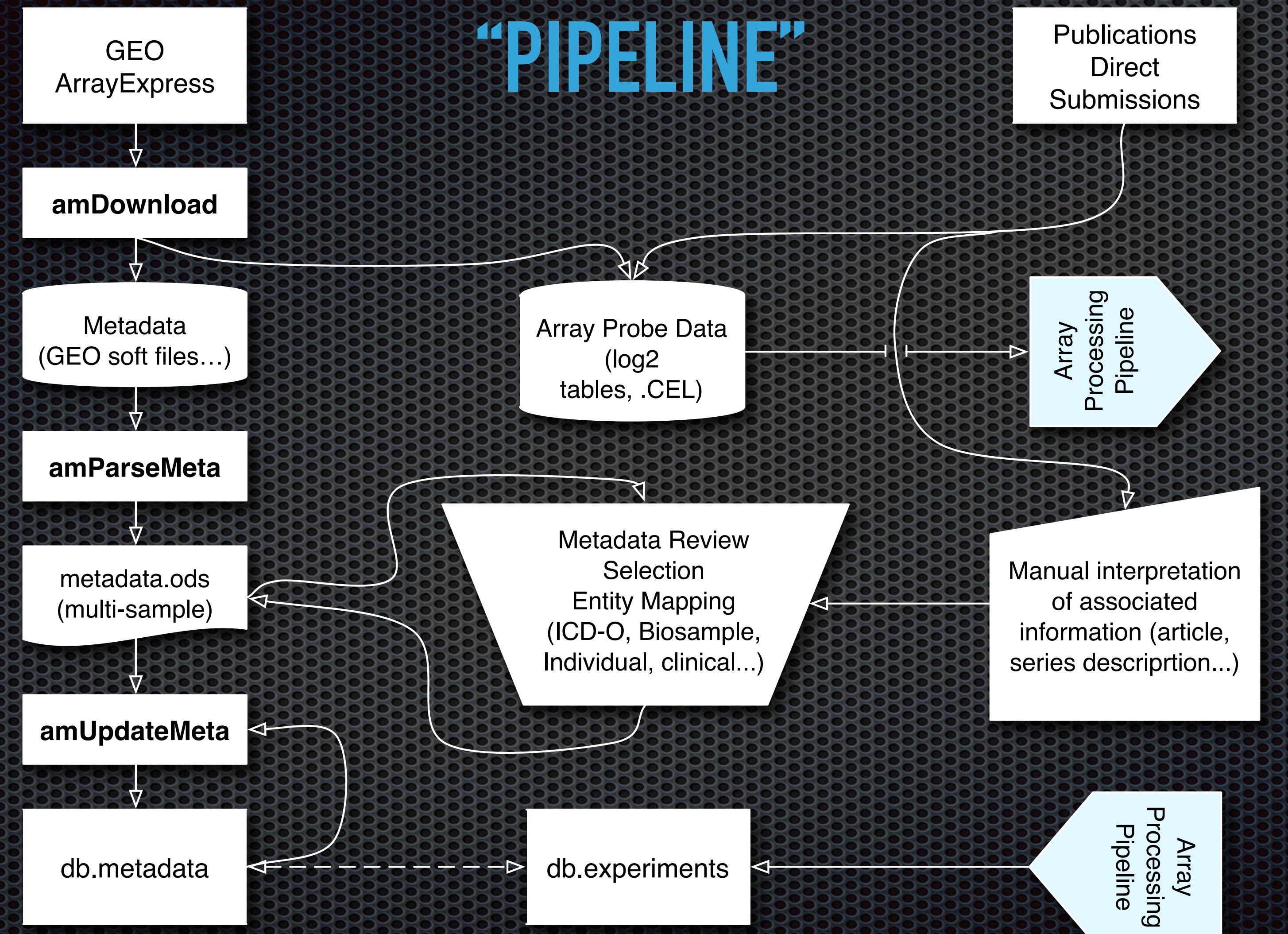
BIOCURATION

BIOINFORMATICS

arrayMap



ARRAYMAP DATA



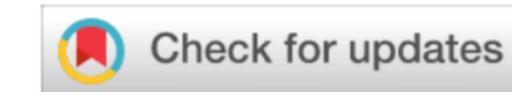
Batch Tool for Genome Liftover

Challenge:

1. Keep the integrity of copy number segments after Liftover.
2. 10% data lost from straight Liftover.
3. 1TB segment and probe data in over 2,000 nested directories

Solution:

1. Algorithm to lift segments.
2. Algorithm for fuzzy remapping.
3. Parallel processing and failure recovery mechanism



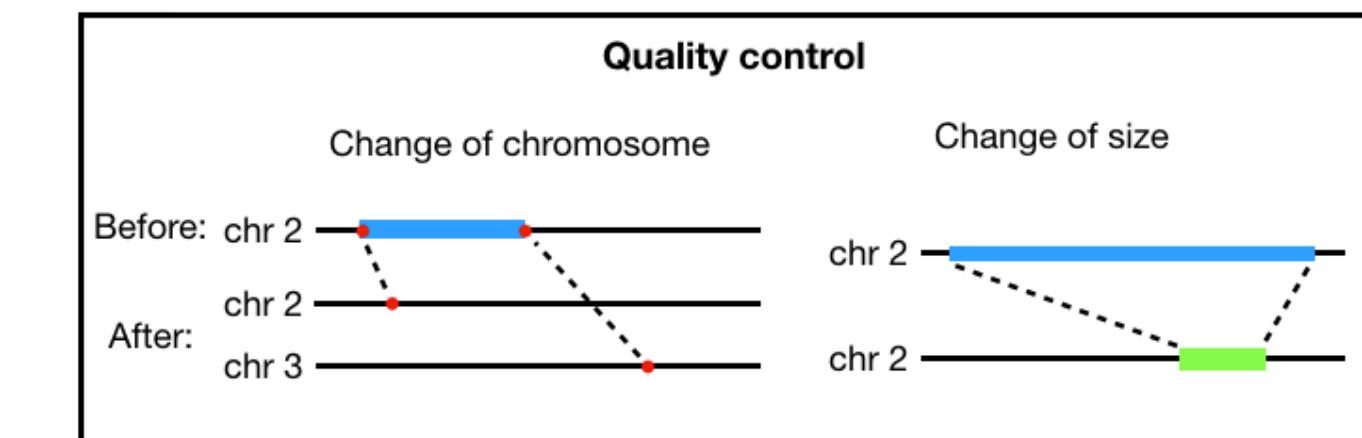
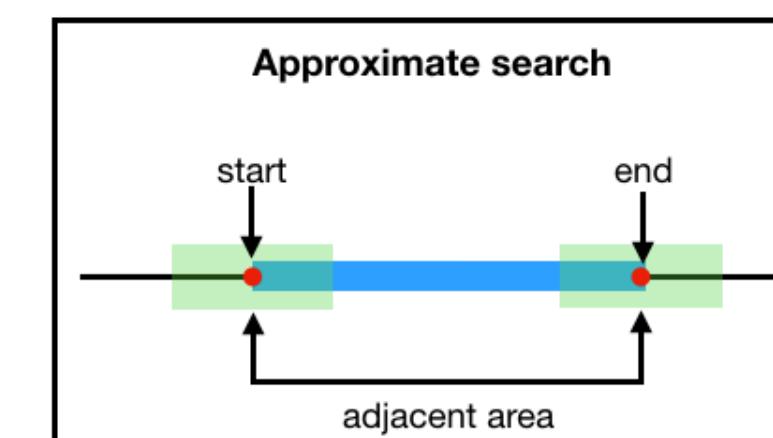
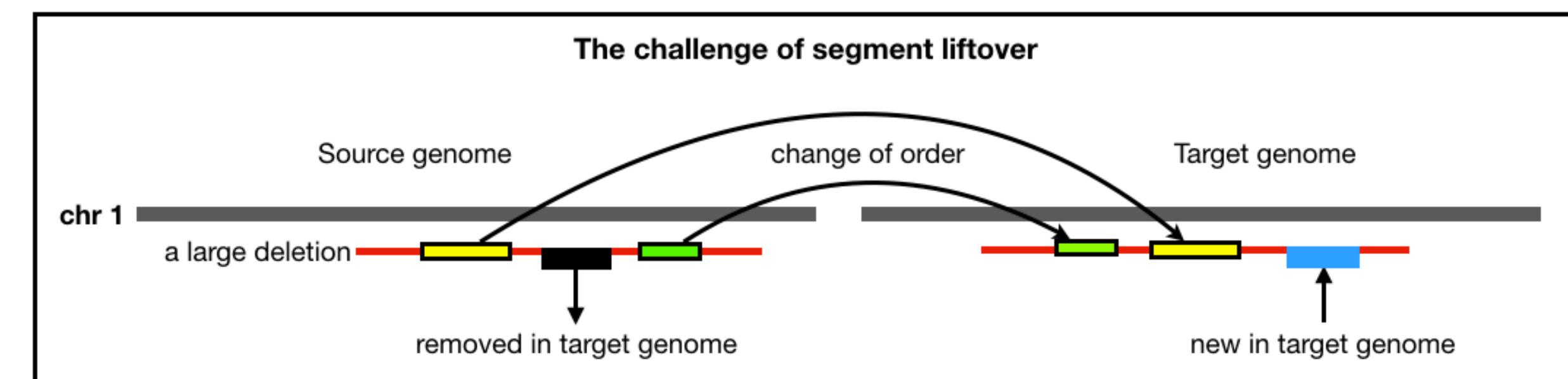
SOFTWARE TOOL ARTICLE

REVISED **segment_liftover : a Python tool to convert segments between genome assemblies [version 2; referees: 2 approved]**

Bo Gao 1,2, Qingyao Huang 1,2, Michael Baudis 1,2

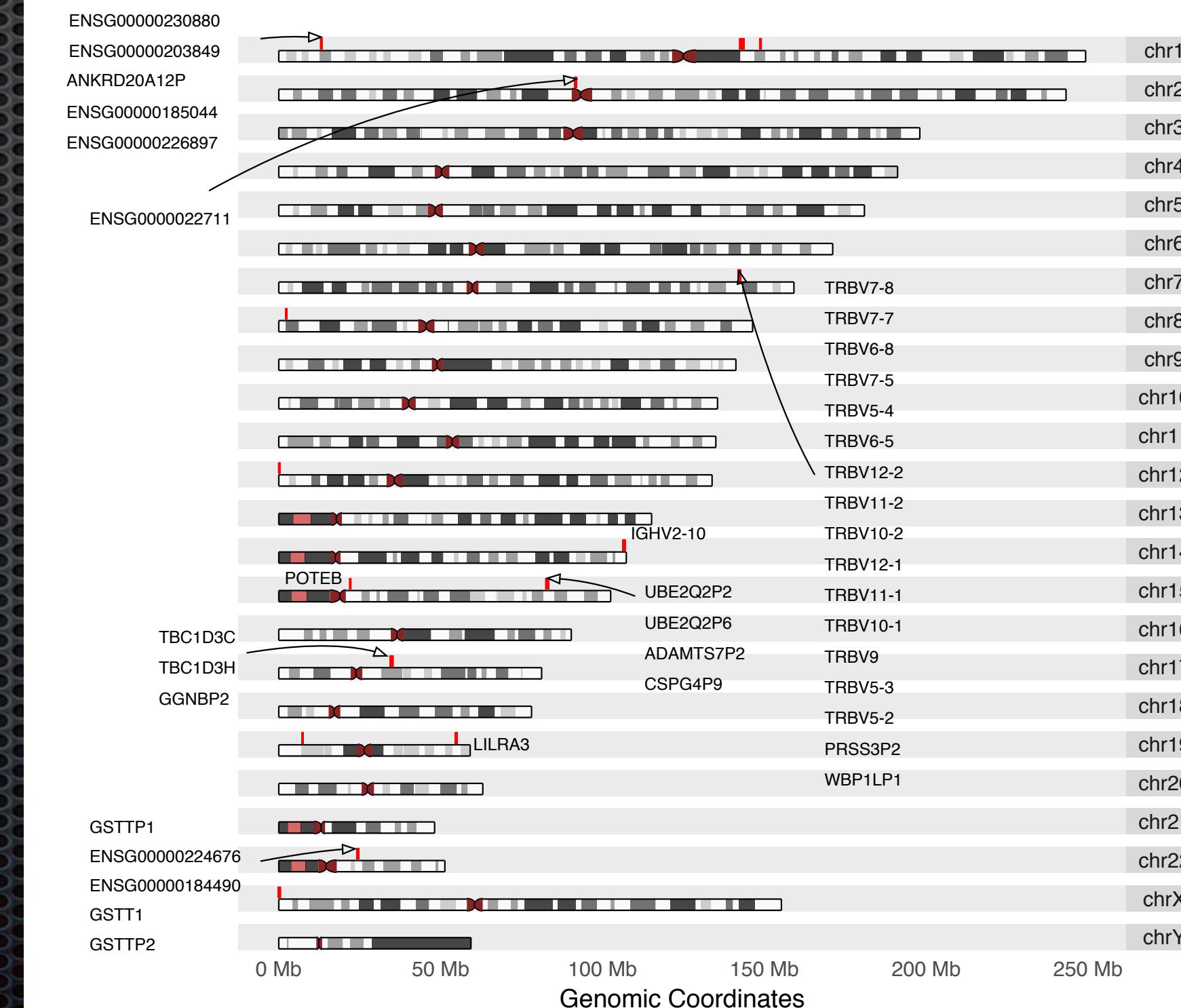
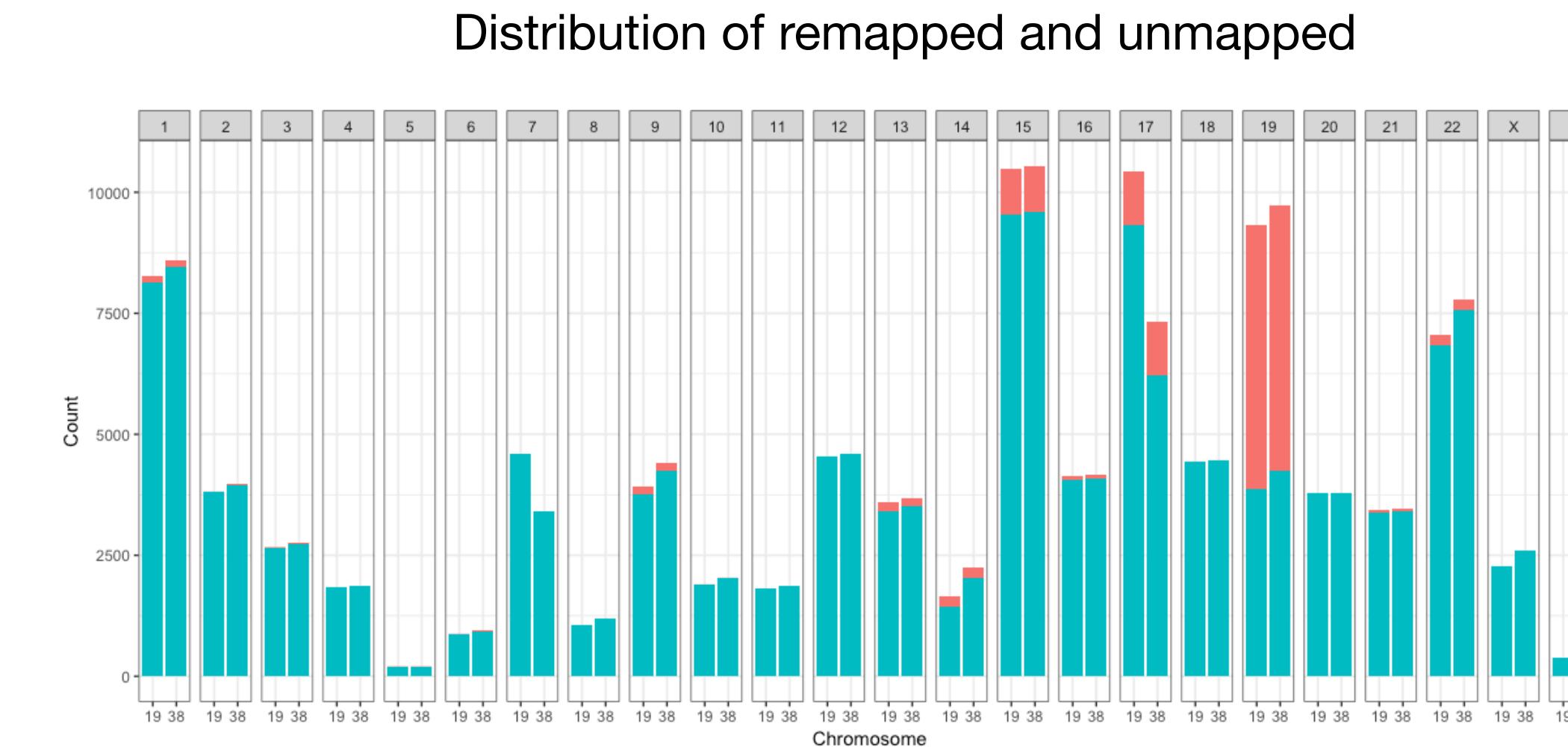
¹Institute of molecular Life Sciences, University of Zürich, Zürich, CH-8057, Switzerland

²Swiss Institute of Bioinformatics, University of Zürich, Zürich, CH-8057, Switzerland

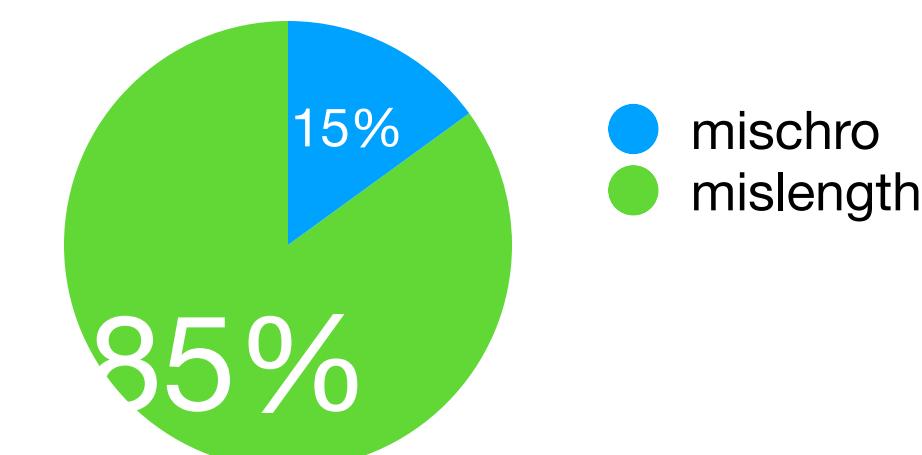


Results of Liftover

- Converted 44,632 probe files and 44,471 segment files from hg19 to hg38.
- The probe data were generated from nine Affymetrix genotyping array platforms, which currently only support annotations for hg19.
- It took 42 hours to convert 5.5 billion probes and 40 minutes to convert 4.8 million segments. (12-core, 128GB RAM machine with 8 parallel processes)
- Reduced information lost from 10% to less than 0.1%



Reason of unmapped segments



arrayMap

Resource for copy number variation data in cancer

arrayMap 

visualizing cancer genome array data @ arraymap.org

arrayMap is a curated reference database and bioinformatics resource targeting copy number profiling data in human cancer. The arrayMap database provides an entry point for meta-analysis and systems level data integration of high-resolution oncogenomic CNA data.

The current data reflects:

- 63060 genomic copy number arrays
- 763 experimental series
- 145 array platforms
- 141 ICD-O cancer entities
- 554 publications (Pubmed entries)

 University of Zurich ^{UZH}

Citation User Guide Registration & Licensing People External Links ↗ FOLLOW US ON [twitter](#)

 130.60.23.21

For the majority of the samples, probe level visualization as well as customized data representation facilitate gene level and genome wide data review. Results from multi-case selections can be connected to downstream data analysis and visualization tools, as we provide through our Progenetix project.

arrayMap is developed by the group "Theoretical Cytogenetics and Oncogenomics" at the Institute of Molecular Life Sciences of the University of Zurich.

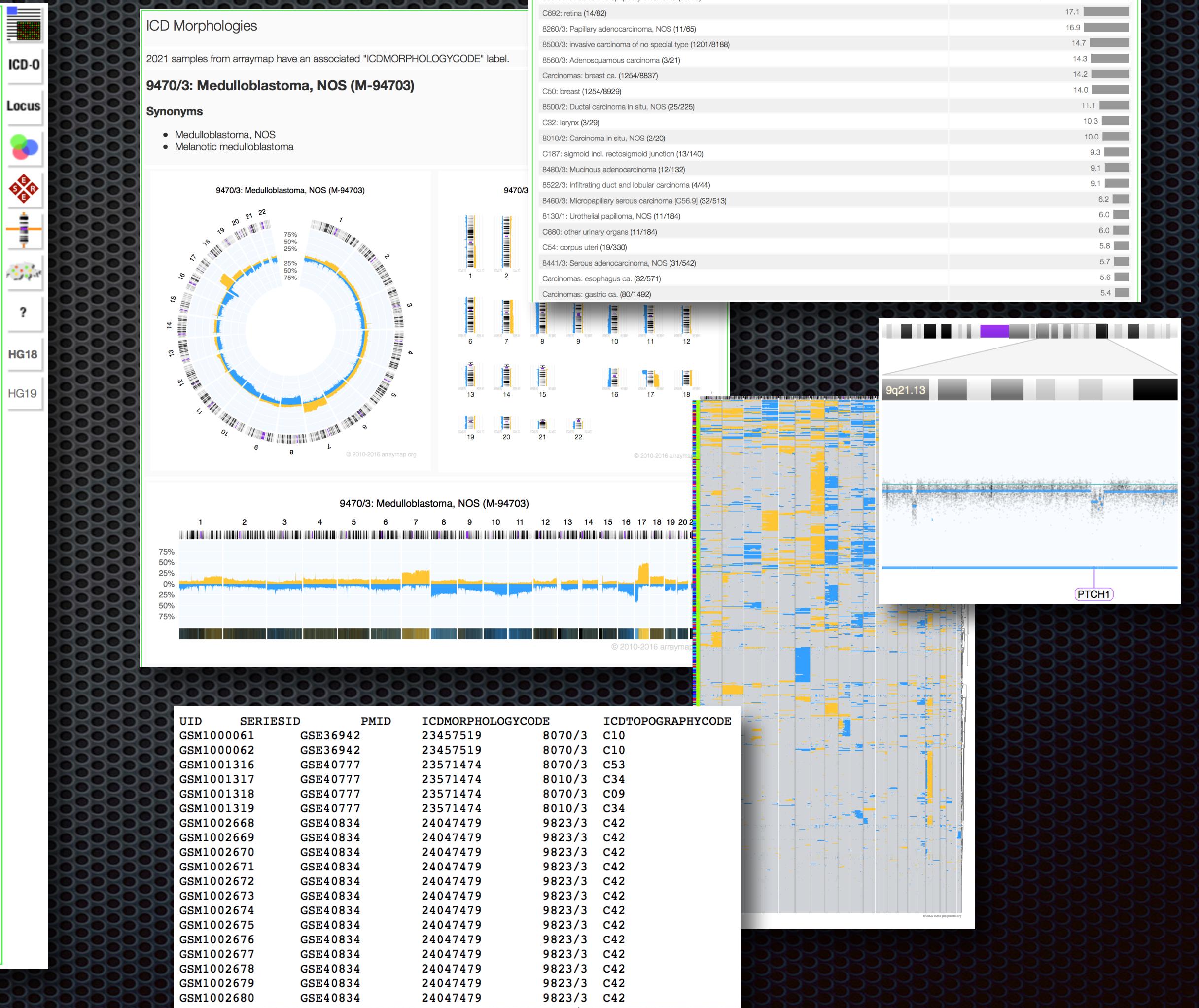
BRAIN TUMOURS	5653 samples ↗	[?]
BREAST CANCER	8329 samples ↗	[?]
COLORECTAL CANCER	3238 samples ↗	[?]
PROSTATE CANCER	991 samples ↗	[?]
STOMACH CANCER	1062 samples ↗	[?]

ARRAYMAP NEWS

- 2016-08-03: SVG graphics
- 2016-05-17: Transitioning to Europe PMC
- More news ...

Feel free to use the data and tools for academic research projects and other applications. If more support and/or custom analysis is needed, please contact Michael Baudis regarding a collaborative project or a special license.

© 2000 - 2016 Michael Baudis, refreshed Mon, 19 Sep 2016 10:20:09 GMT in 6.87s on server 130.60.240.68. No responsibility is taken for the correctness of the data presented nor the results achieved with the Progenetix tools.



arrayMap

Resource for copy number variation data in cancer

- Typical use: Gene CNA frequencies
- allows for (gene symbol supported) query for DUP/DEL overlapping a region
- size windowing ("focal" changes)
- no thorough statistics implemented yet (e.g. modeling overall CNA background)

arrayMap

Search Samples
Search Publications
Gene CNA Frequencies
User Data
Progenetix
 University of Zurich UZH
Citation
User Guide
Registration & Licensing
People


130.226.87.158

Search for single or combined imbalances

FIND CNAS BY GENE OR REGION
Gene Symbol: TP53 [CDKN2A] 9:21957750-21984490:DEL
Minimal size: minimum CNA segment
Maximum size: 5000

CNA REGION SIZE (KB)
SELECT CANCER TYPES (ICD-O 3)
Type here for selection ...

SELECT CANCER LOCI
Type here for selection ...

SELECT CANCER TYPES (NCIT NEOPLASMS)
Type here for selection ...

ARRAY SERIES IDS
CITY (AND DISTANCE FROM IT)

Query Database

4098 of 62104 cases matched the selection criteria.

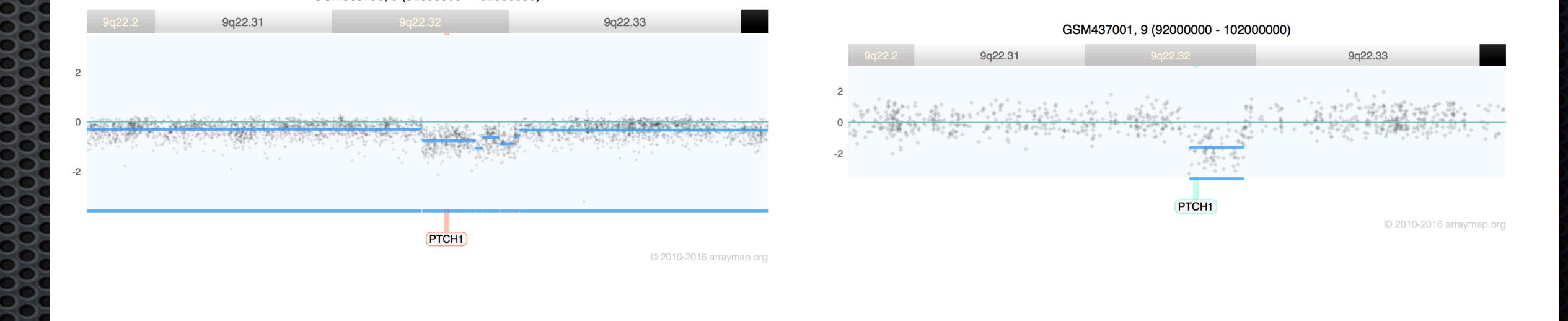
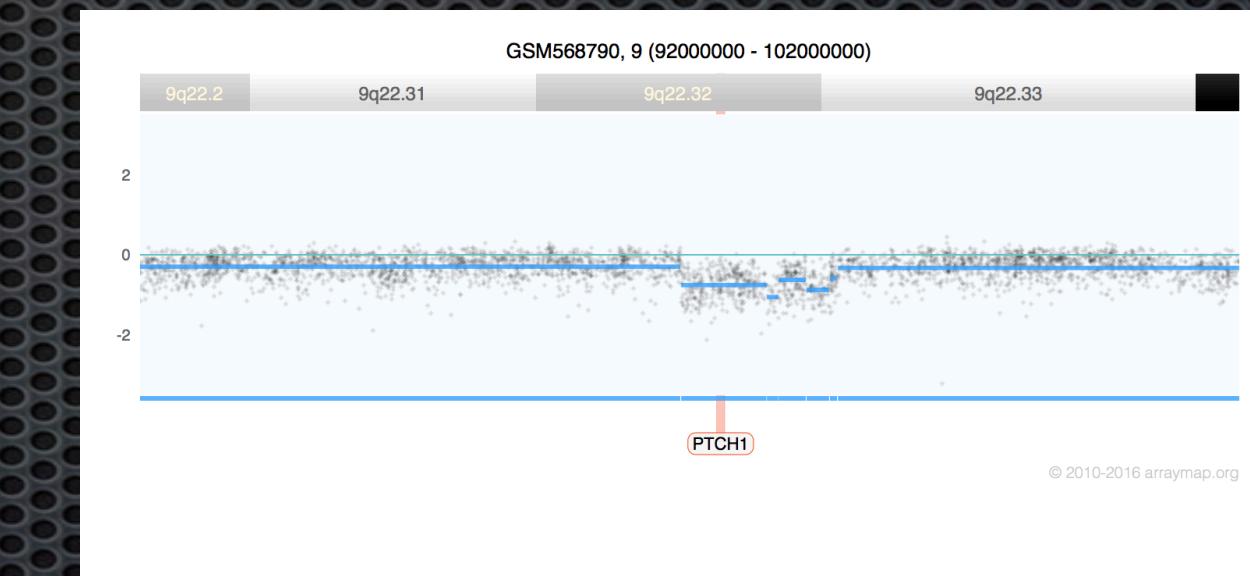
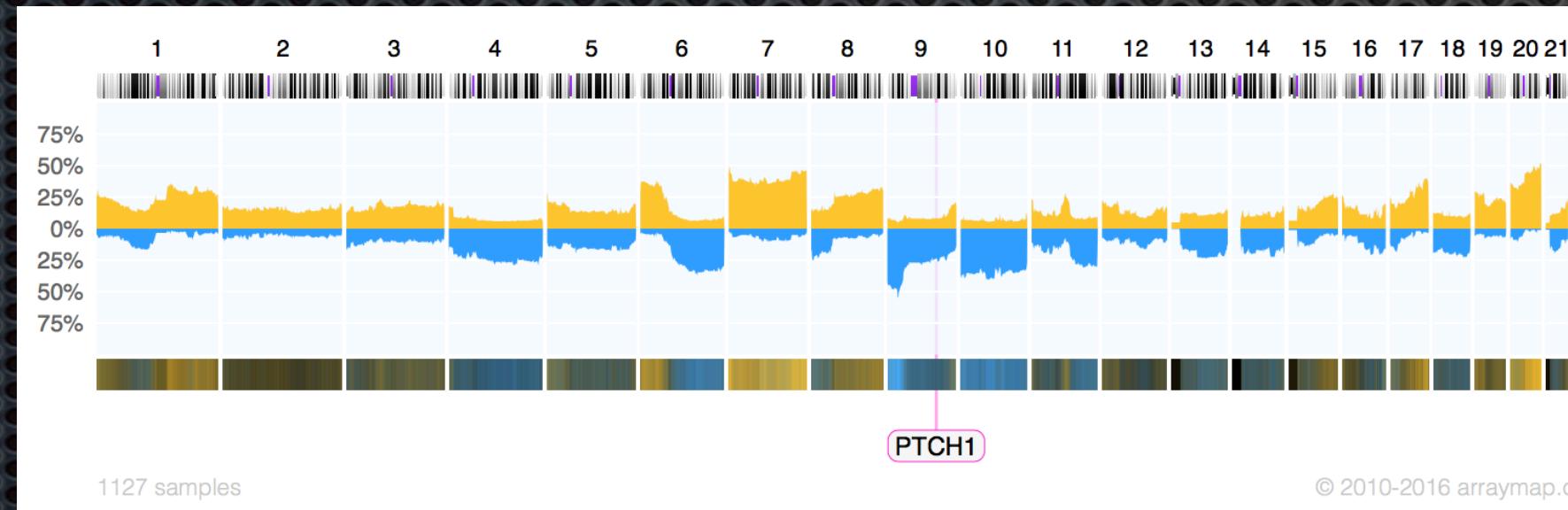
¶ Subsets	¶ Samples	¶ Observations	¶ Frequency
9590/3: 9590/3: Malignant lymphoma, NOS	6	5	0.833
9052/3: 9052/3: Epithelioid mesothelioma, malignant	22	14	0.636
9050/3: 9050/3: Mesothelioma, malignant	87	55	0.632
9729/3: 9729/3: Precursor T-cell lymphoblastic lymphoma	30	14	0.467
8130/1: 8130/1: Papillary transitional cell neoplasm of low malignant potential	38	12	0.316
8250/3: 8250/3: Bronchiolo-alveolar adenocarcinoma, NOS	23	7	0.304
9801/3: 9801/3: Acute leukemia, NOS	10	3	0.300
9440/3: 9440/3: Glioblastoma, NOS	2047	607	0.297
8560/3: 8560/3: Adenosquamous carcinoma	21	6	0.286
8072/3: 8072/3: Squamous cell carcinoma, large cell, nonkeratinizing, NOS	96	27	0.281
9827/3: 9827/3: Adult T-cell leukemia/lymphoma (HTLV-1 positive)	179	49	0.274
9837/3: 9837/3: Precursor T-cell lymphoblastic leukemia	233	59	0.253
8500/2: 8500/2: Intraductal carcinoma, noninfiltrating, NOS	147	36	0.245
8430/3: 8430/3: Mucoepidermoid carcinoma	21	5	0.238
8811/3: 8811/3: Fibromyxosarcoma	43	10	0.233
8160/3: 8160/3: Cholangiocarcinoma	39	9	0.231
8012/3: 8012/3: Large cell carcinoma, NOS	48	11	0.229

ICD-O
Locus
NCIt
S
ER
C
?

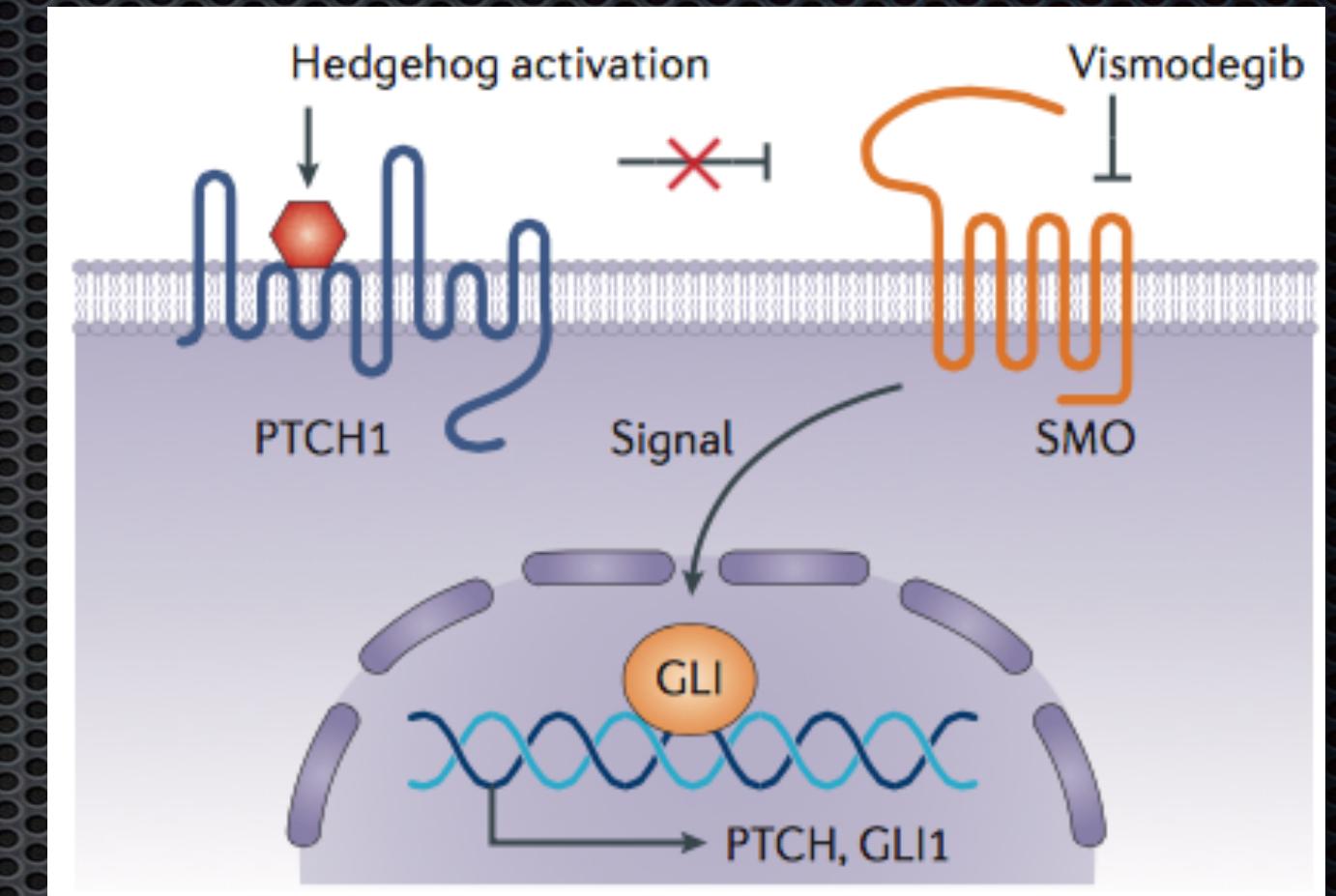
Rare Events & Hidden Therapeutic Options?

Example: PTCH1 deletions in malignant melanomas

- PTCH1 is a actionable tumor suppressor gene, which has been demonstrated in e.g. basalomas and medulloblastomas
- analysis of 1127 samples from 26 different publications could identify **focal** deletions in 4 samples
- a current project addresses the focal involvement of all mapped genes, in >50'000 cancer genome profiles



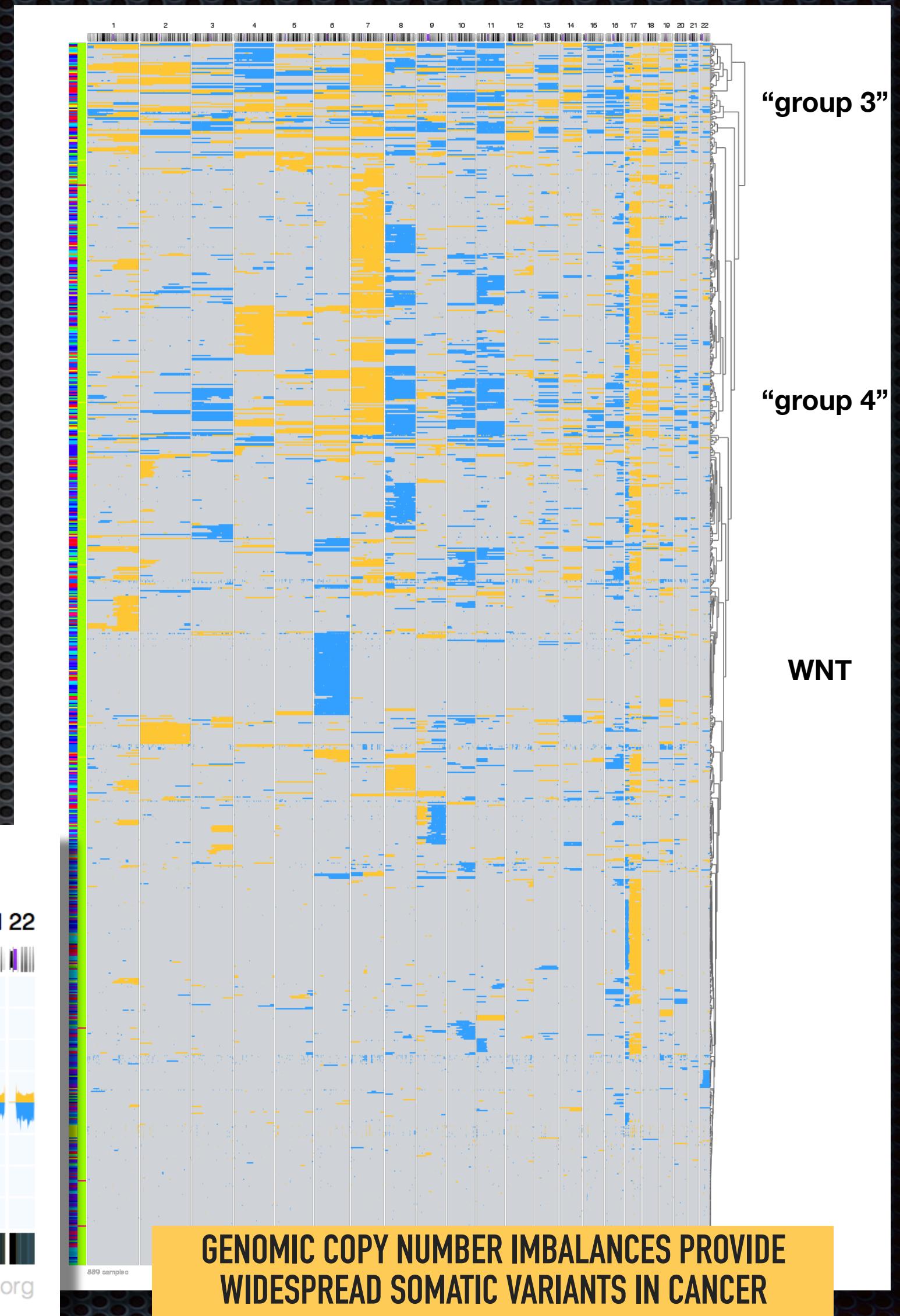
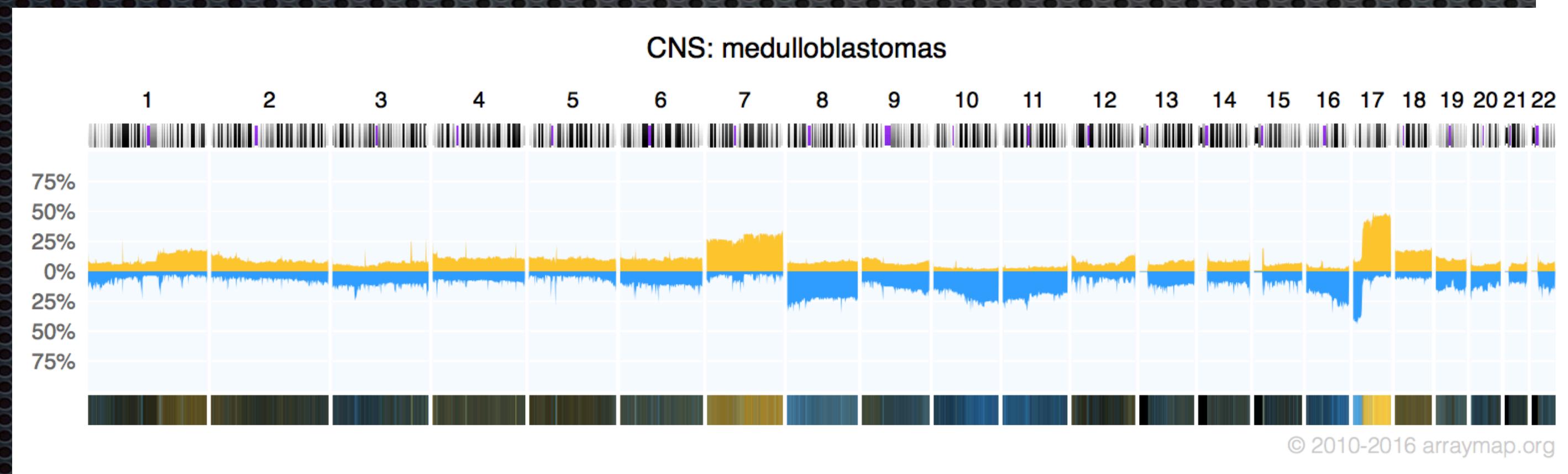
Examples of focal / homozygous PTCH1 deletions detected in the analysis of 1127 genomic array datasets. Focal somatic imbalance events are considered an indicator for oncogenic involvement of the affected target genes.



In its normal function, PTCH1 is a tumor suppressor gene in the sonic hedgehog pathway and inhibits SMO driven transcriptional activation. A loss of PTCH1 function (mutation, deletion) can be mitigated through drugs antagonistic to SMO activation.

Somatic Mutations In Cancer: Patterns

- many tumor types express **recurrent mutation patterns**
- How can** those patterns be used for classification and determination of biological mechanisms?

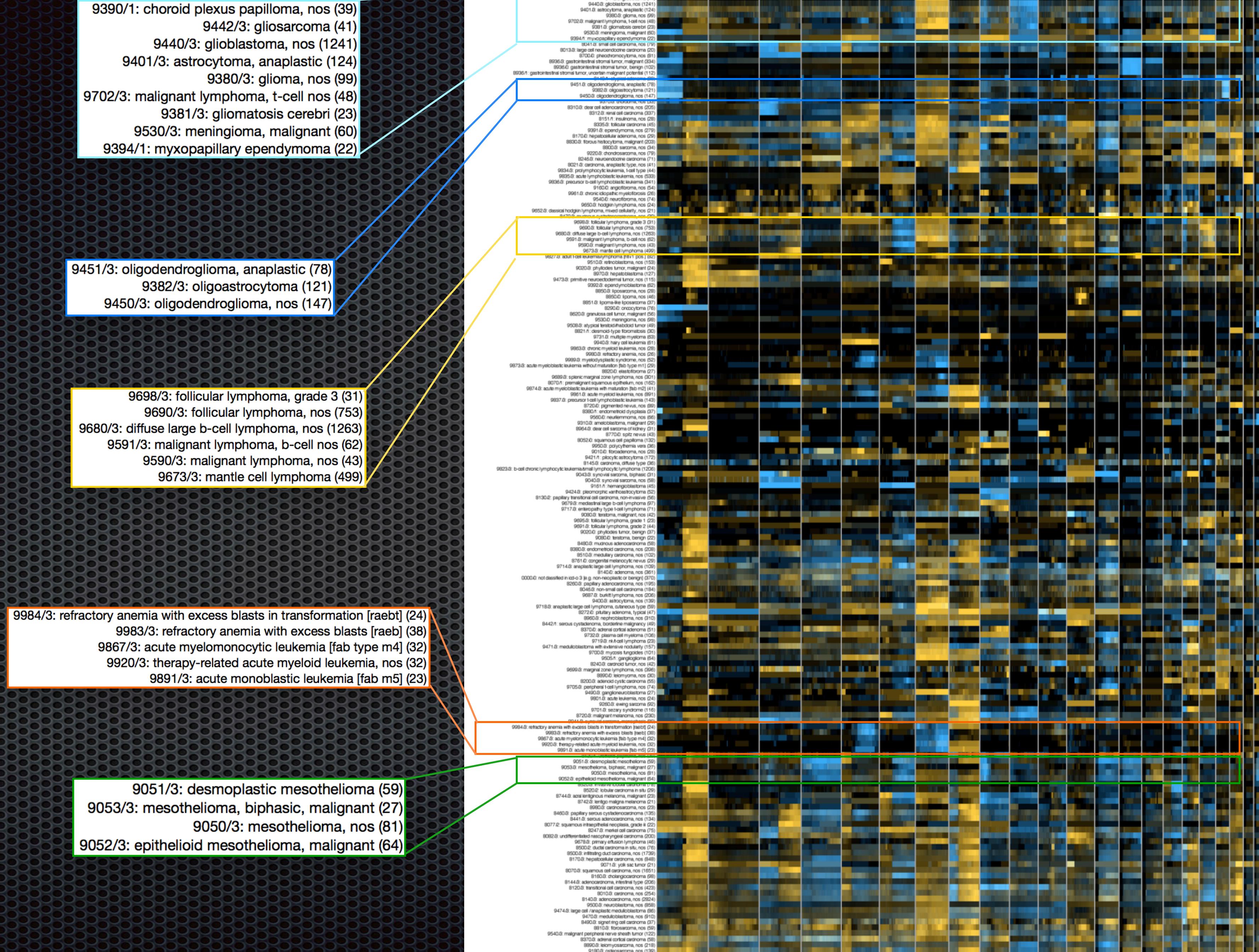


A genomic copy number histogram for malignant medulloblastomas, the most frequent type of pediatric brain tumors, displaying regions of genomic duplications and deletions. These can be decomposed into individual tumor profiles which segregate into several clusters of related mutation patterns with functional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

Somatic Mutations In Cancer: Patterns III

Making the case for genomic classifications

Some related cancer entities show similar copy number profiles



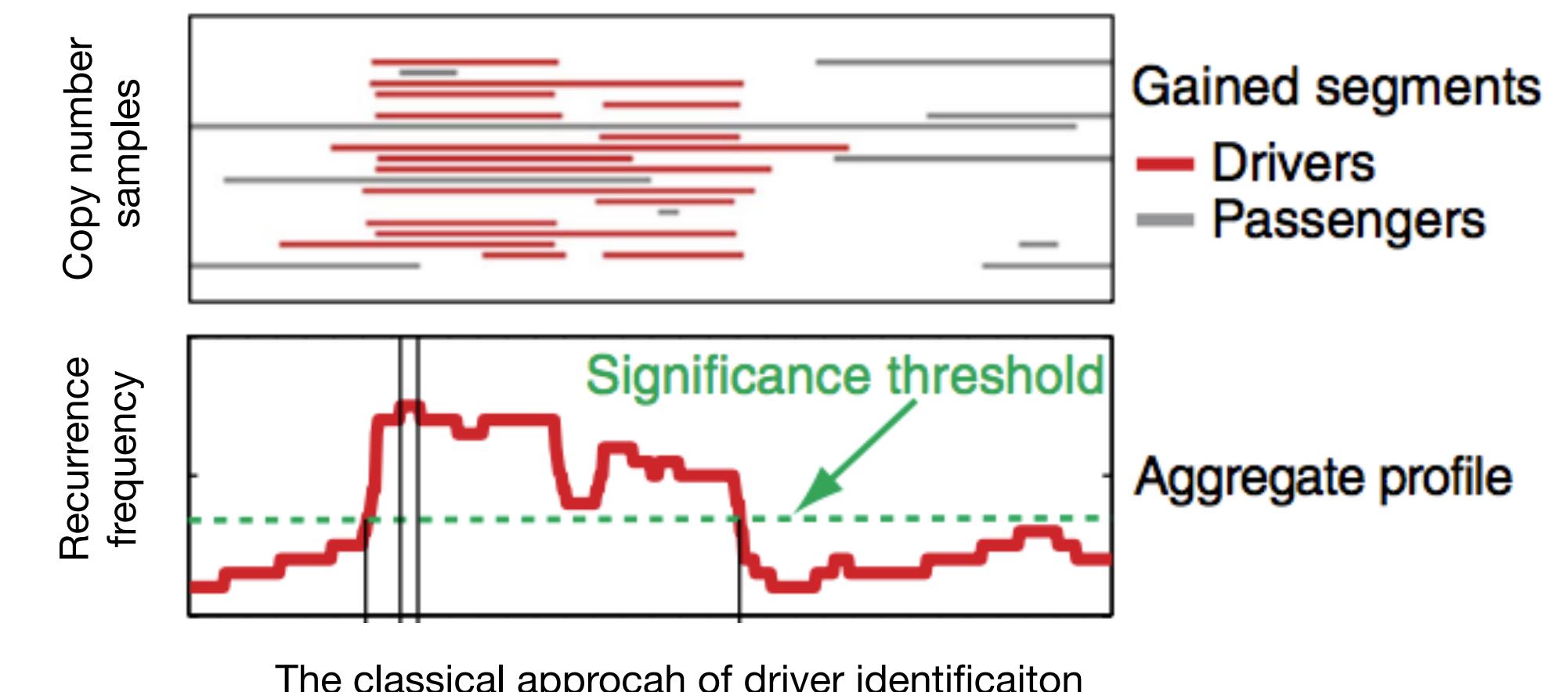
A machine learning approach to identify driver genes in cancer development

Traditional approach

- method: finding focal signal peaks
- data: point mutation or copy number variation

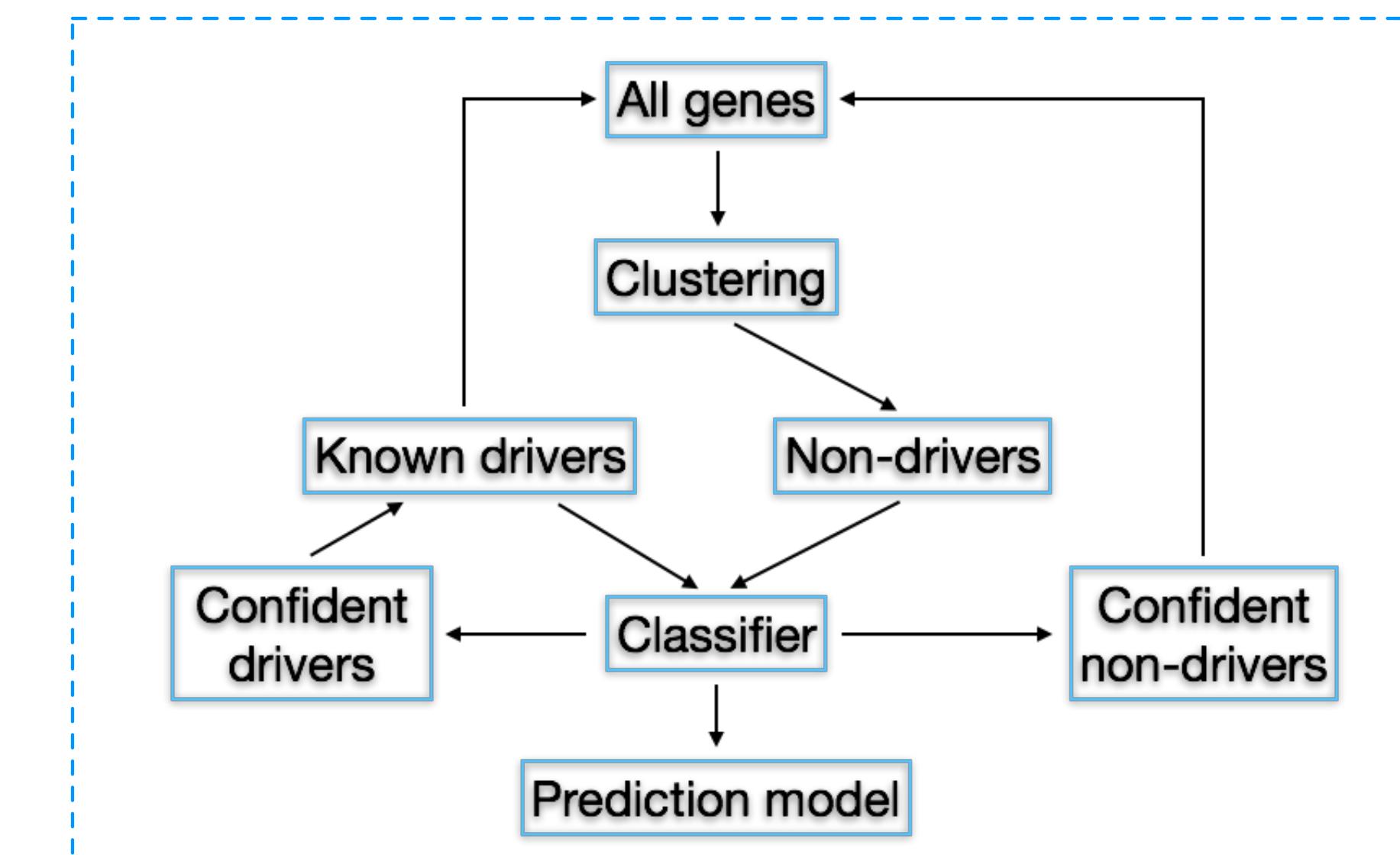
New approach

- method: unsupervised machine learning
- data: point mutation, copy number variation and methylation



The classical approach of driver identificaiton

Cited from: Van Dyk, Ewald, et al. *Nature communications* 7 (2016).



The recursive learning workflow

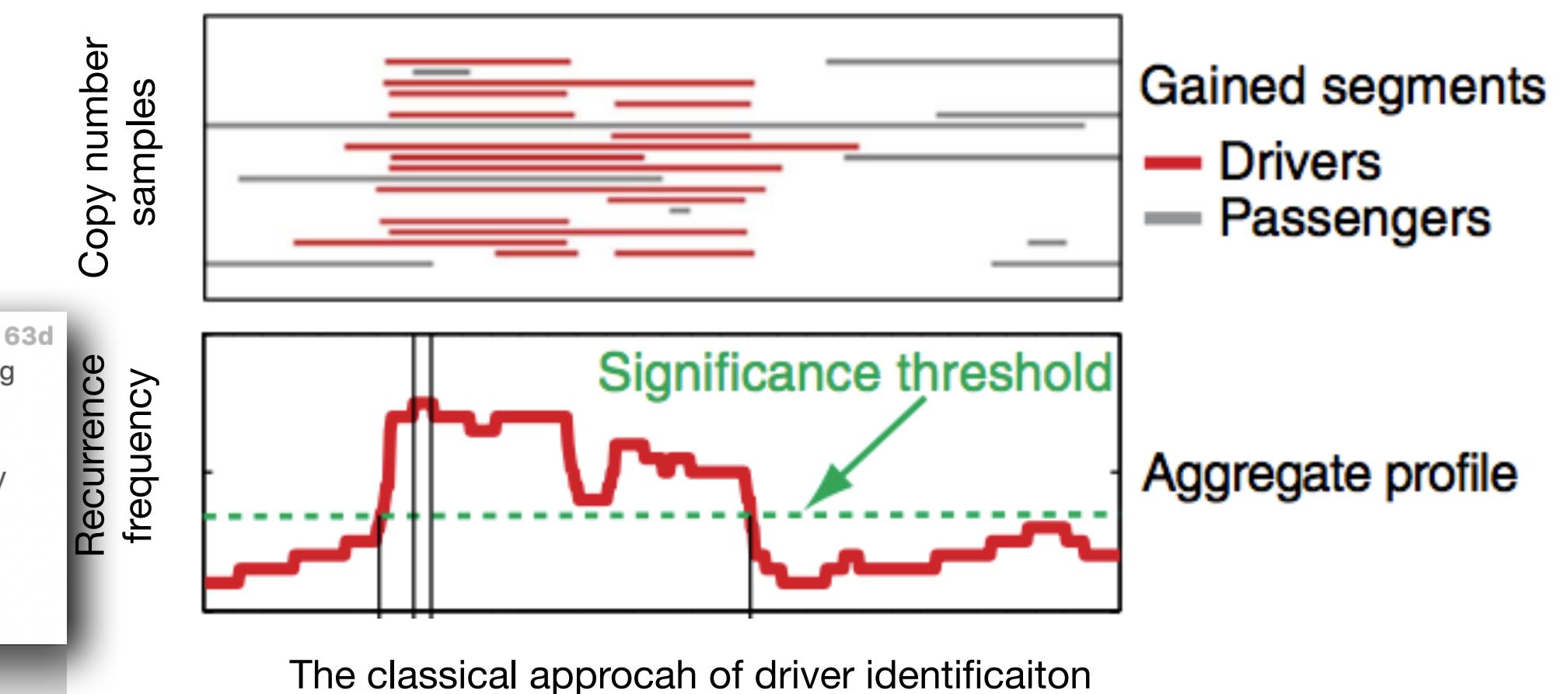
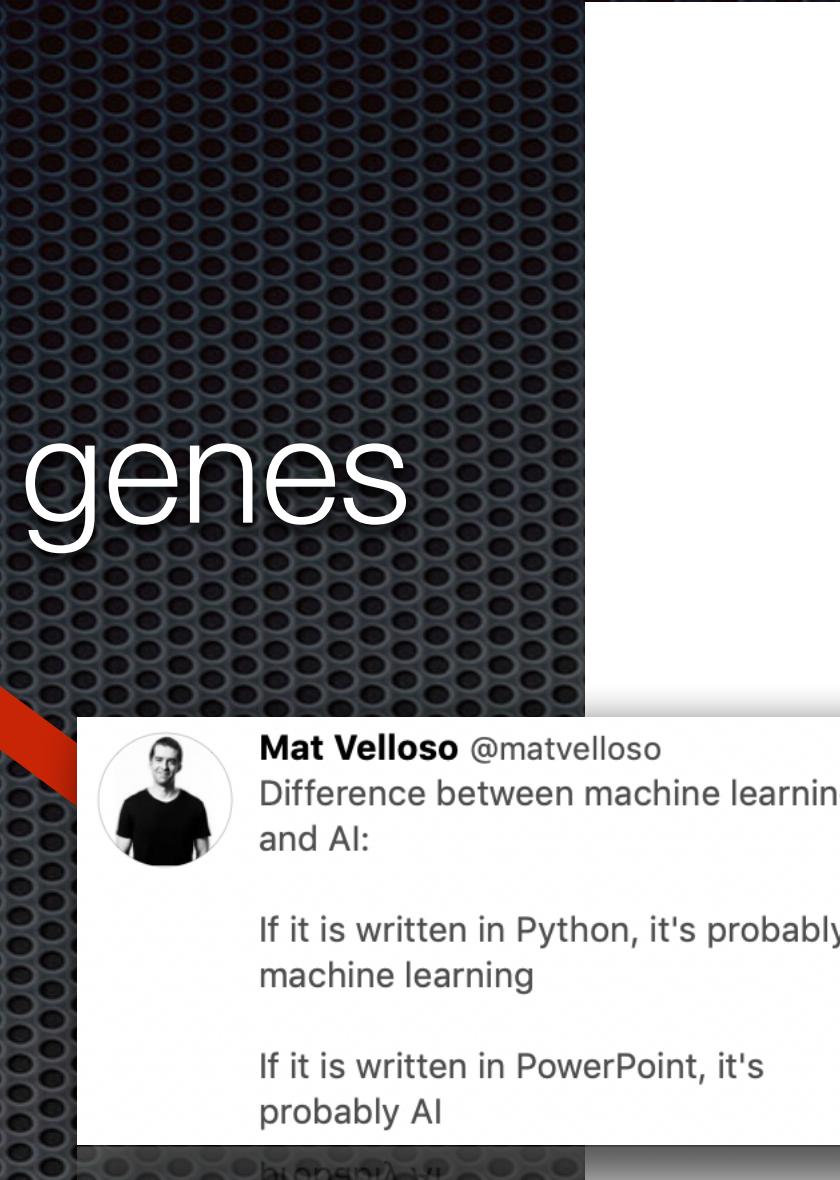
An Artificial Intelligence approach to identify driver genes in cancer development

Traditional approach

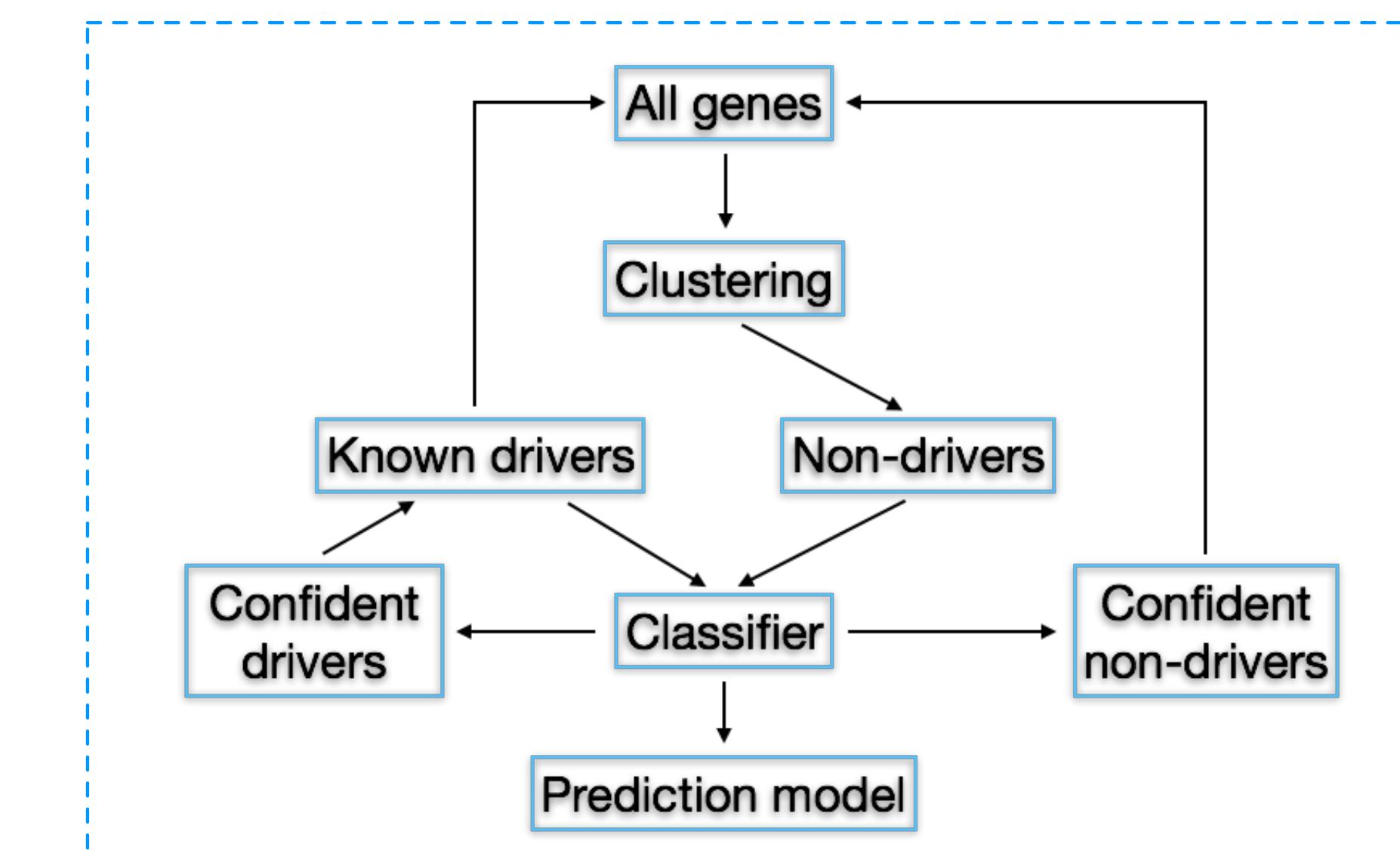
- method: finding focal signal peaks
- data: point mutation or copy number variation

New approach

- method: unsupervised machine learning
- data: point mutation, copy number variation and methylation

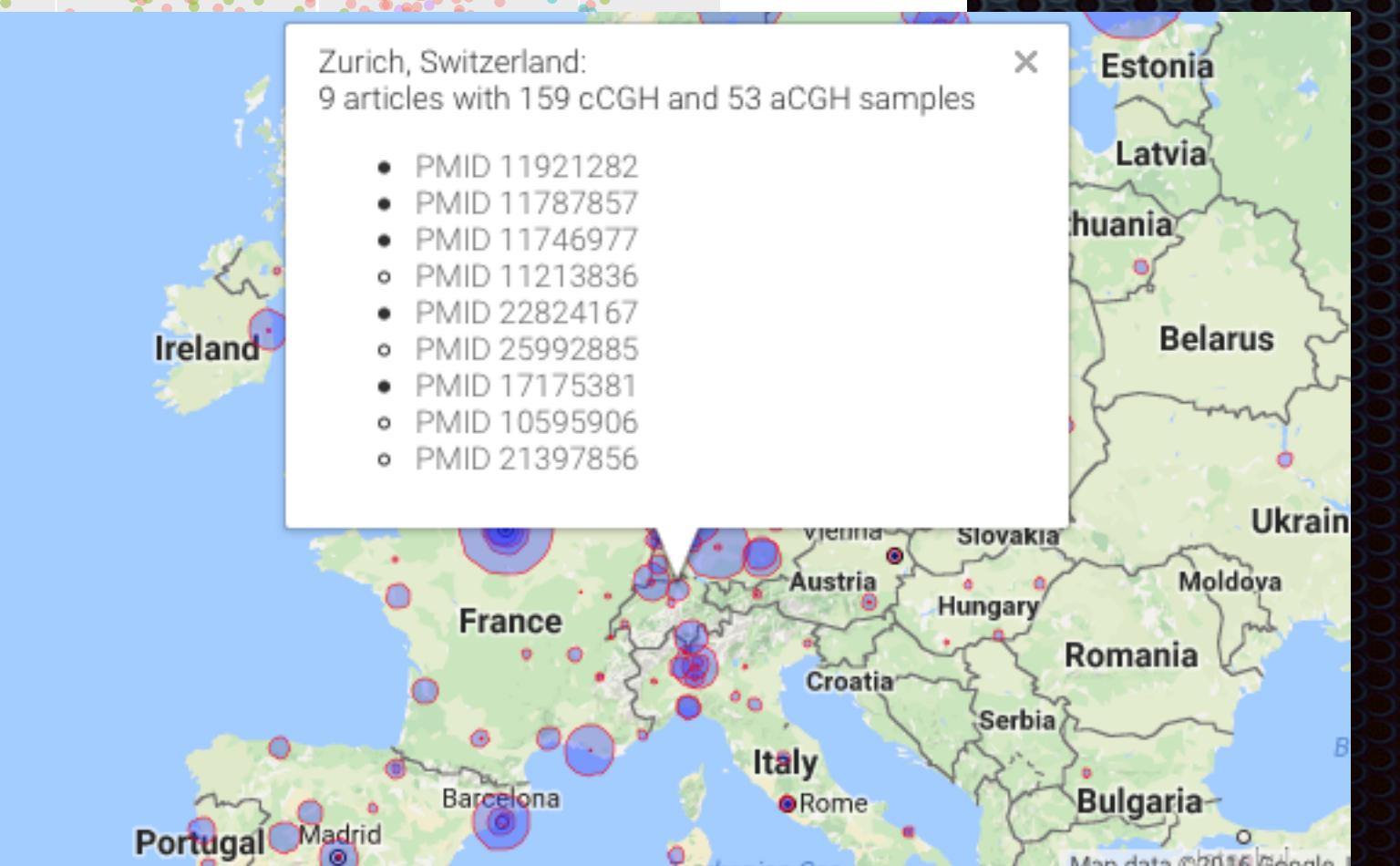
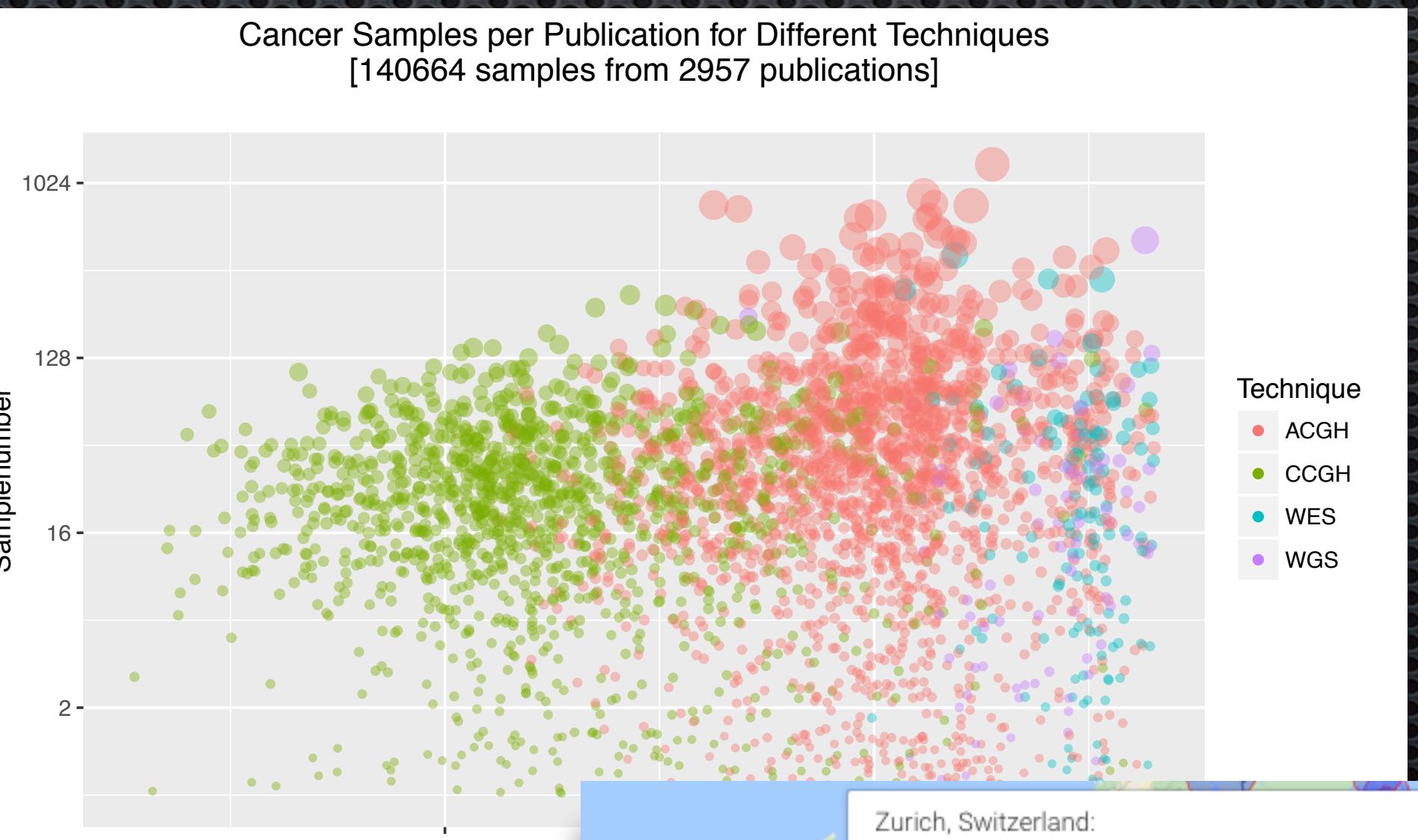
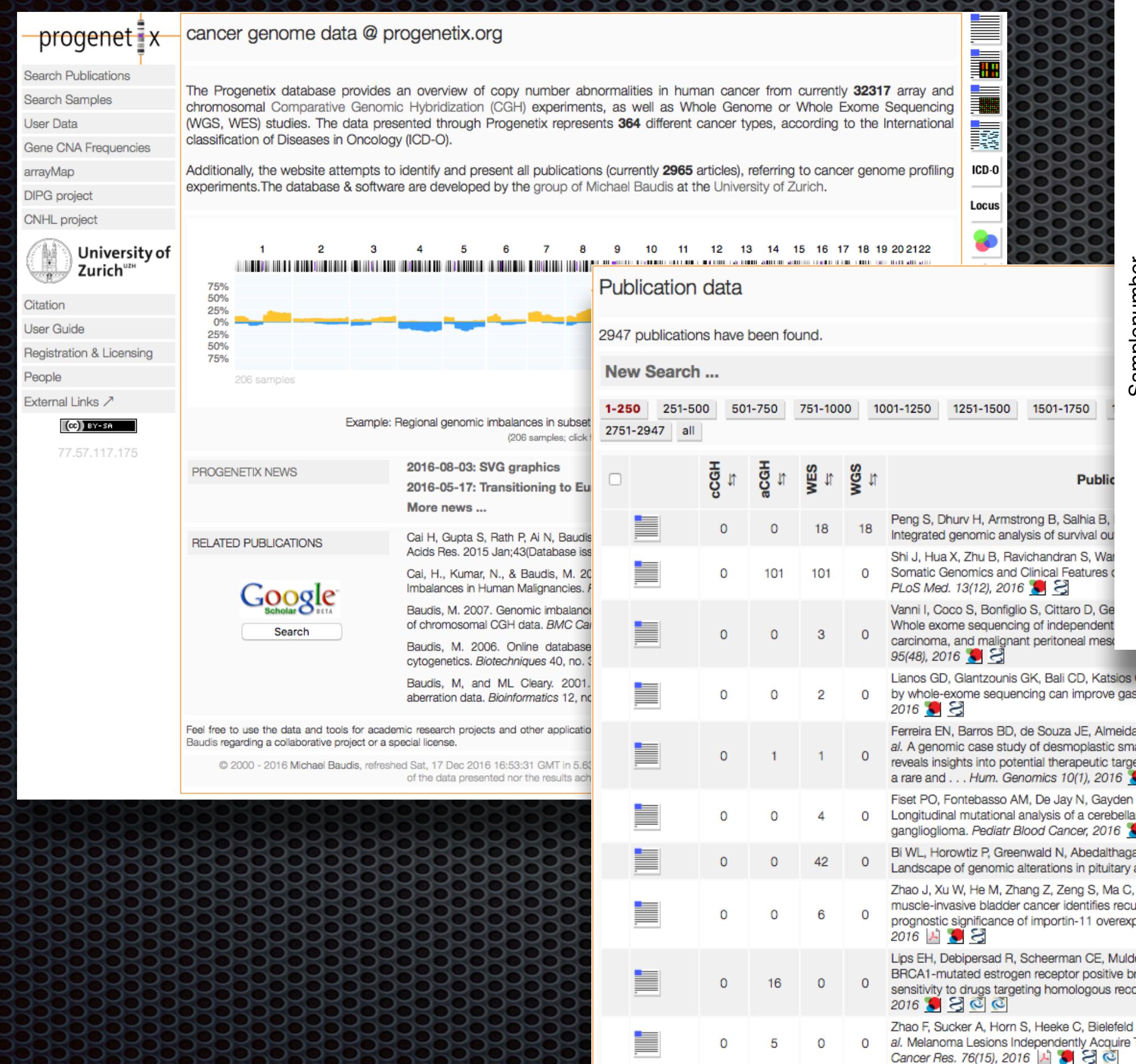


Cited from: Van Dyk, Ewald, et al. *Nature communications* 7 (2016).



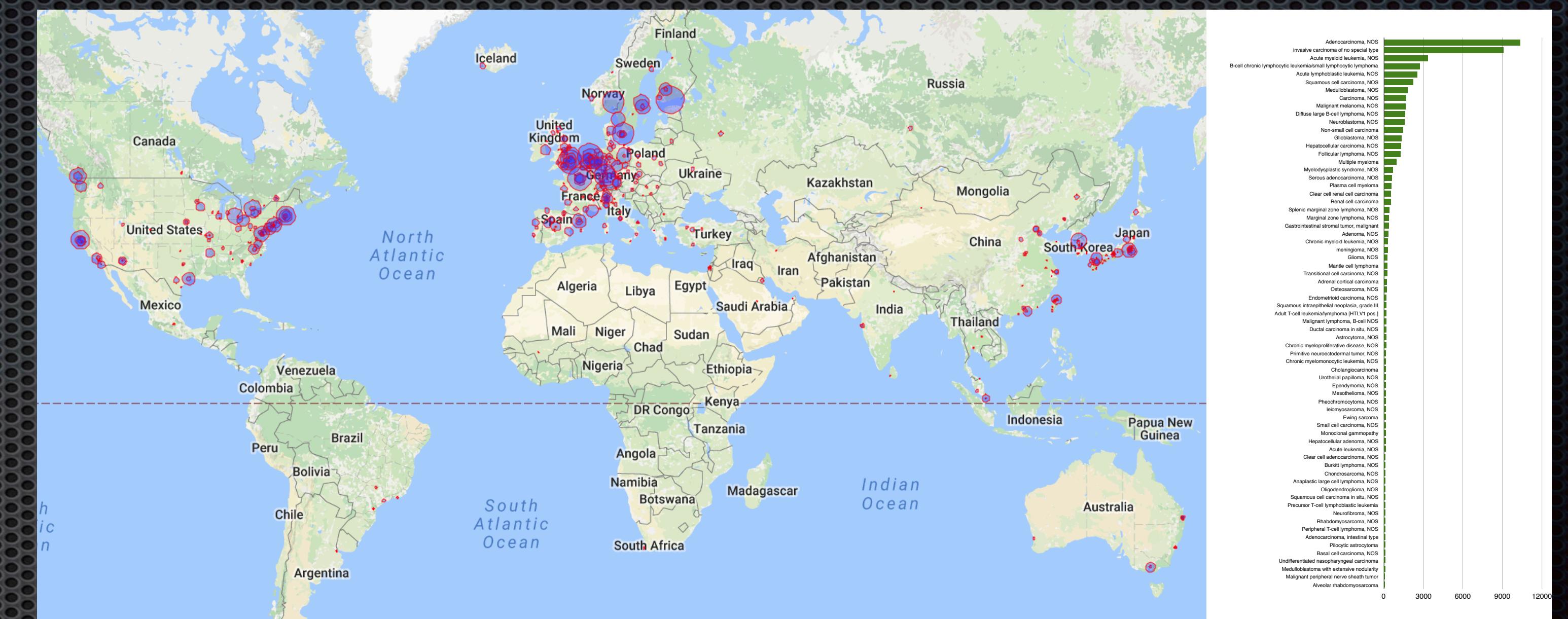
The recursive learning workflow

Progenetix: Cancer Genome Profiles, Article Metrics, Epistemology, Resource Hub



Bias in Ascertainment / Background / Environment in Cancer Genome Studies

- the frequency of many genome variants depends on the genetic background
- cancer incidence & type can correlate to environmental factors
- geographic analysis can support interpretation and point to knowledge gaps



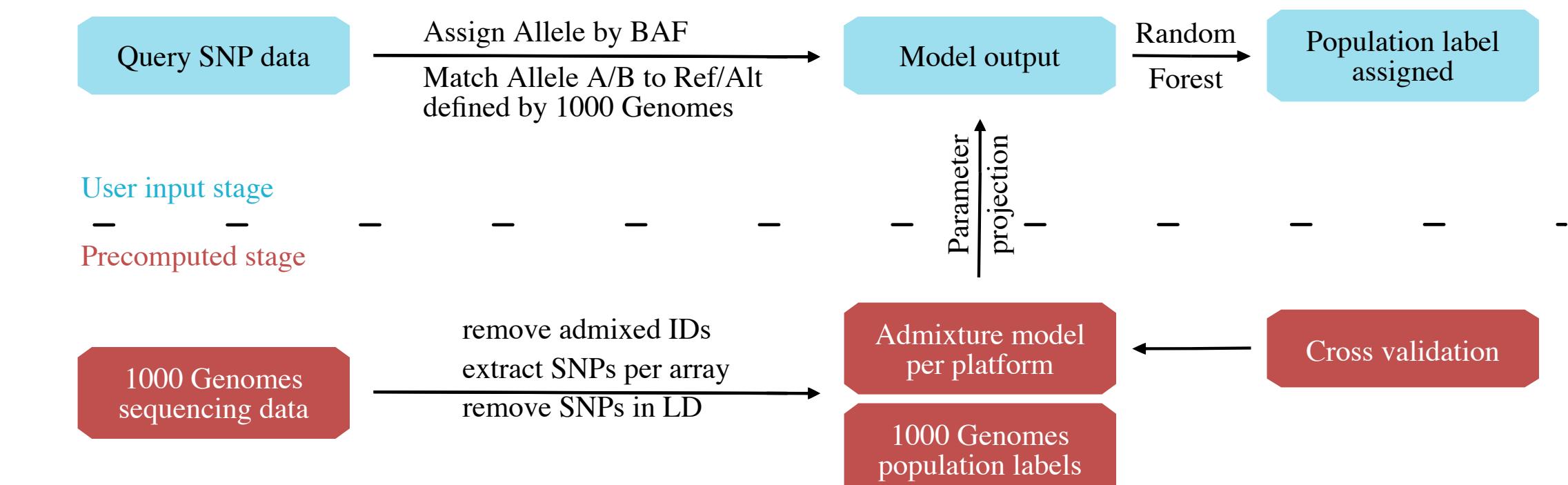
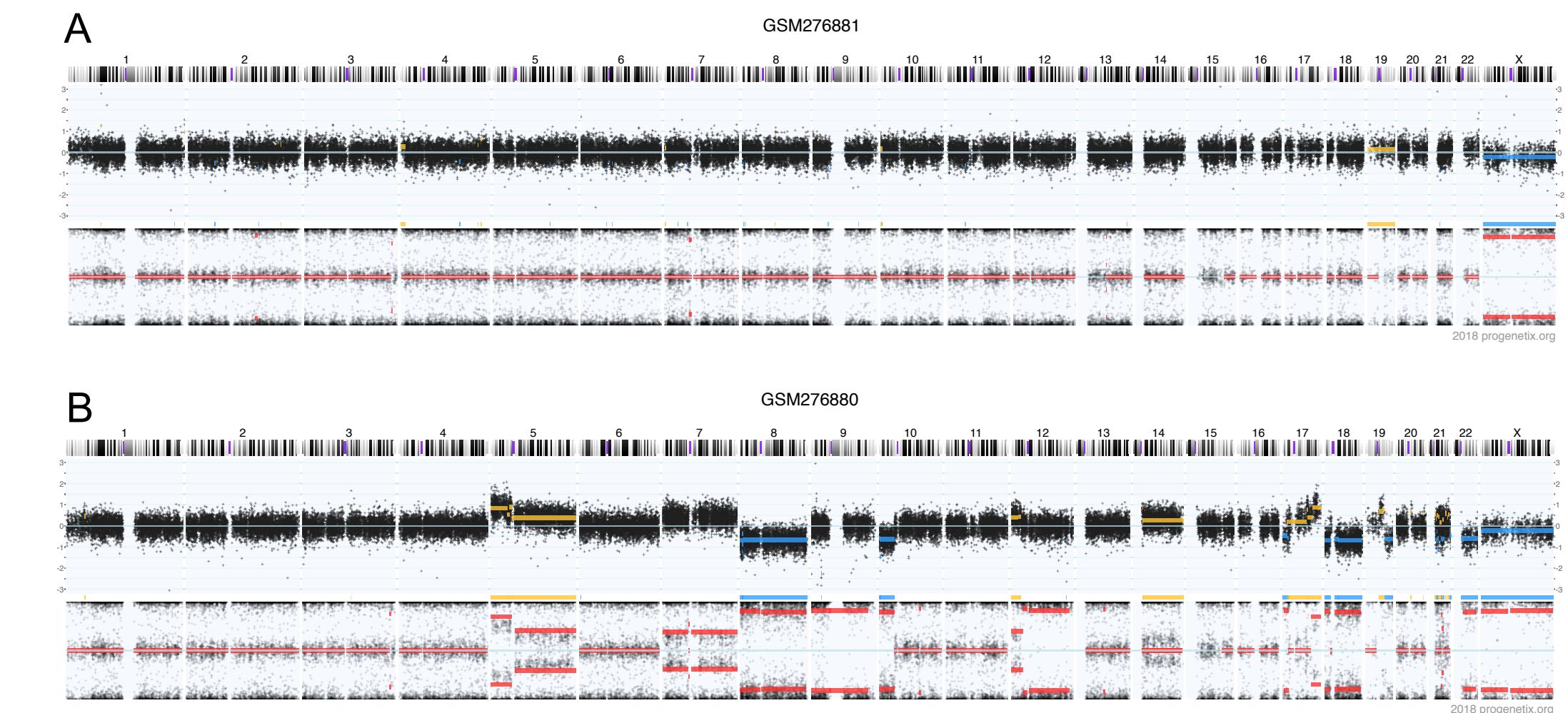
Geographic distribution of >140'000 cancer genome profiles reported in the literature. The numbers are derived from the 2947 publications registered in the Progenetix database.

Population stratification in cancer samples based on SNP array data

- 2504 genome profiles from 1000 Genome project phase 1 as reference
- 5 (or 26) superpopulations: South Asia, Europe, South America, East Asia and Africa.
- SNP positions used in 9 Affymetrix SNP arrays are extracted to train a population admixture model.

Enabling population assignment from cancer genomes with SNP2pop

Qingyao Huang^{1,2} and Michael Baudis^{1,2✉}



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 metadata estimated a matching rate between 88 % to 92%.
- Ethnicity labels indicated in metadata are vague compared to the standardized output from our tool

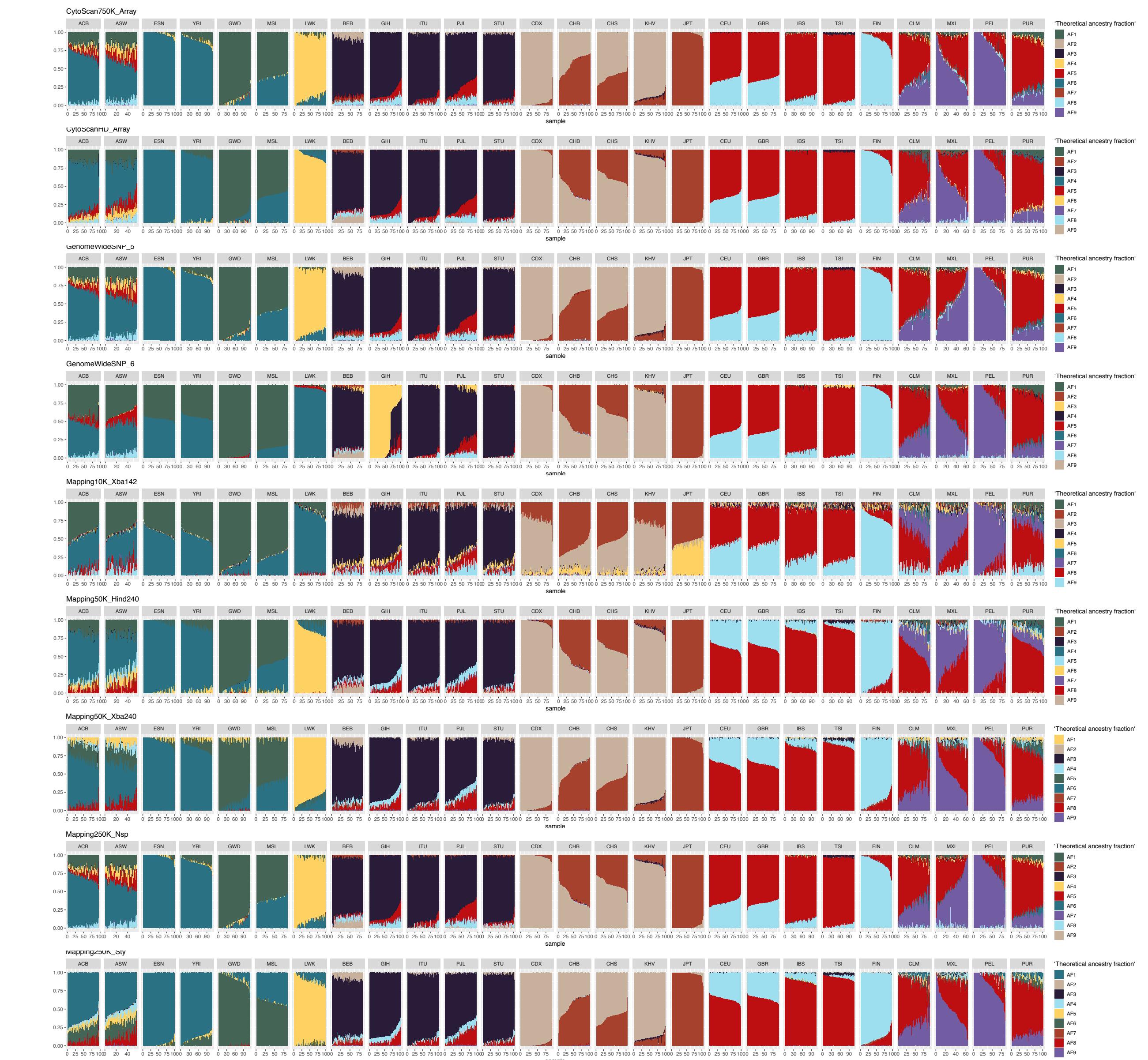


Figure S1 The fraction or contribution of theoretical ancestors ($k=9$) in reference individuals from 1000 Genomes Project with regard to nine SNP array platforms. The x-axis are individual samples, grouped by their respective population. Groups belonging to the same continent/superpopulation are placed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

Developing the GA4GH Schemas

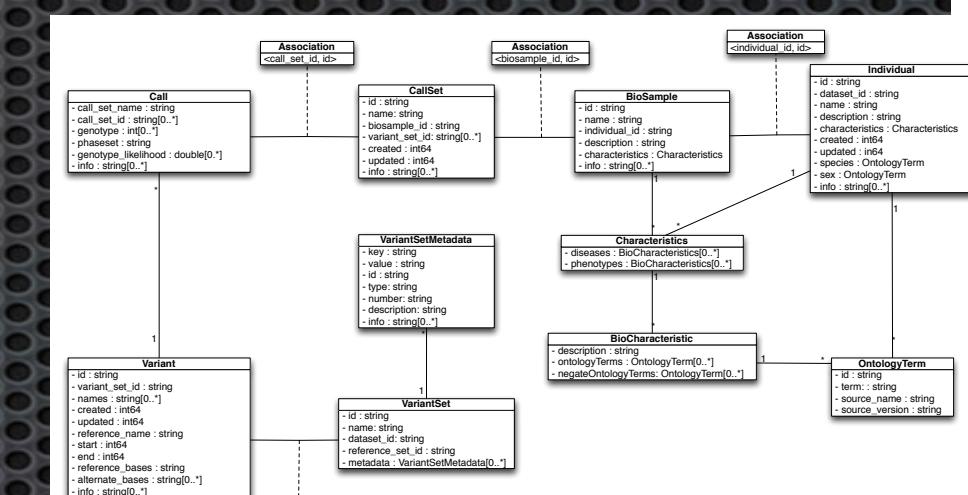
▶ arrayMap for GA4GH

- metadata schema development through implementation of arrayMap resource data
 - OntologyTerm objects for biodata
 - implementation w/ ontology services

Driving Beacon Development

Beacon+

- CNV/CNA as first type of structural variants
 - disease specific queries
 - quantitative reporting



```
{  
    "_id" : ObjectId("58297ca32ca4591e5a0df054"),  
    "id" : "AM_V_1778741",  
    "variant_set_id" : "AM_VS_HG18",  
    "reference_name" : "10"  
    "start" : 579049,  
    "end" : 17236099,  
    "alternate_bases" : "DUP",  
    "reference_bases" : ".",  
    "info" : {  
        "svlen":16657050,  
        "cipos":[  
            -1000,  
            1000  
        ],  
        "ciend": [  
            -1000,  
            1000  
        ]  
    },  
    "calls" : [  
        {  
            "genotype" : [  
                ".",".  
                ".  
            ],  
            "call_set_id" : "AM_CS_TCGA-61-1917-01A-01D-0648-01",  
            "info" : {  
                "segvalue" : 0.5491  
            }  
        }  
    ],  
    "created" : ISODate("2016-11-14T08:33:58.202Z"),  
    "updated" : ISODate("2016-11-14T08:33:58.202Z"),  
}
```

Baudisgroup: Wrangling GA4GH Schemas

- object model
Individual - Biosample - Callset
- referencing of ontologies instead of text descriptors
- **these** are no “real” open ontologies
- data standards use (e.g. ISO)
- fallback to generic object map for unassigned data; this should disappear over time

```
"id" : "PGX_AM_BS_GSM510730",
"individual_id" : "PGX_IND_GSM510730",
"name" : "PGX_AM_BS_GSM510730",
"description" : "breast carcinoma",
"bio_characteristics" : [
  {
    "description" : "breast carcinoma",
    "ontology_terms" : [
      {
        "term_id" : "ncit:C4017",
        "term_label" : "Ductal Breast Carcinoma"
      },
      {
        "term_id" : "pgx:icdom:8500_3",
        "term_label" : "invasive carcinoma of no special type"
      },
      {
        "term_id" : "pgx:icdot:C50",
        "term_label" : "breast"
      }
    ],
    "negated_ontology_terms" : [ ],
  },
  "individual_age_at_collection" : "P47Y",
  "attributes" : {
    "tnm" : {
      "values" : [
        {
          "string_value" : "T1N0M0"
        }
      ]
    },
    "location" : {
      "geo_label" : "Oslo, Norway",
      "latitude" : 59.91,
      "longitude" : 10.75,
      "geo_precision" : "city"
    },
    "external_identifiers" : [
      {
        "database" : "Pubmed",
        "identifier" : "20592421",
        "relation" : "part_of"
      },
    ],
    "updated" : ISODate("2017-03-20T08:37:07.771Z"),
  }
]
```

Ontologies need an Einstein to sort them out



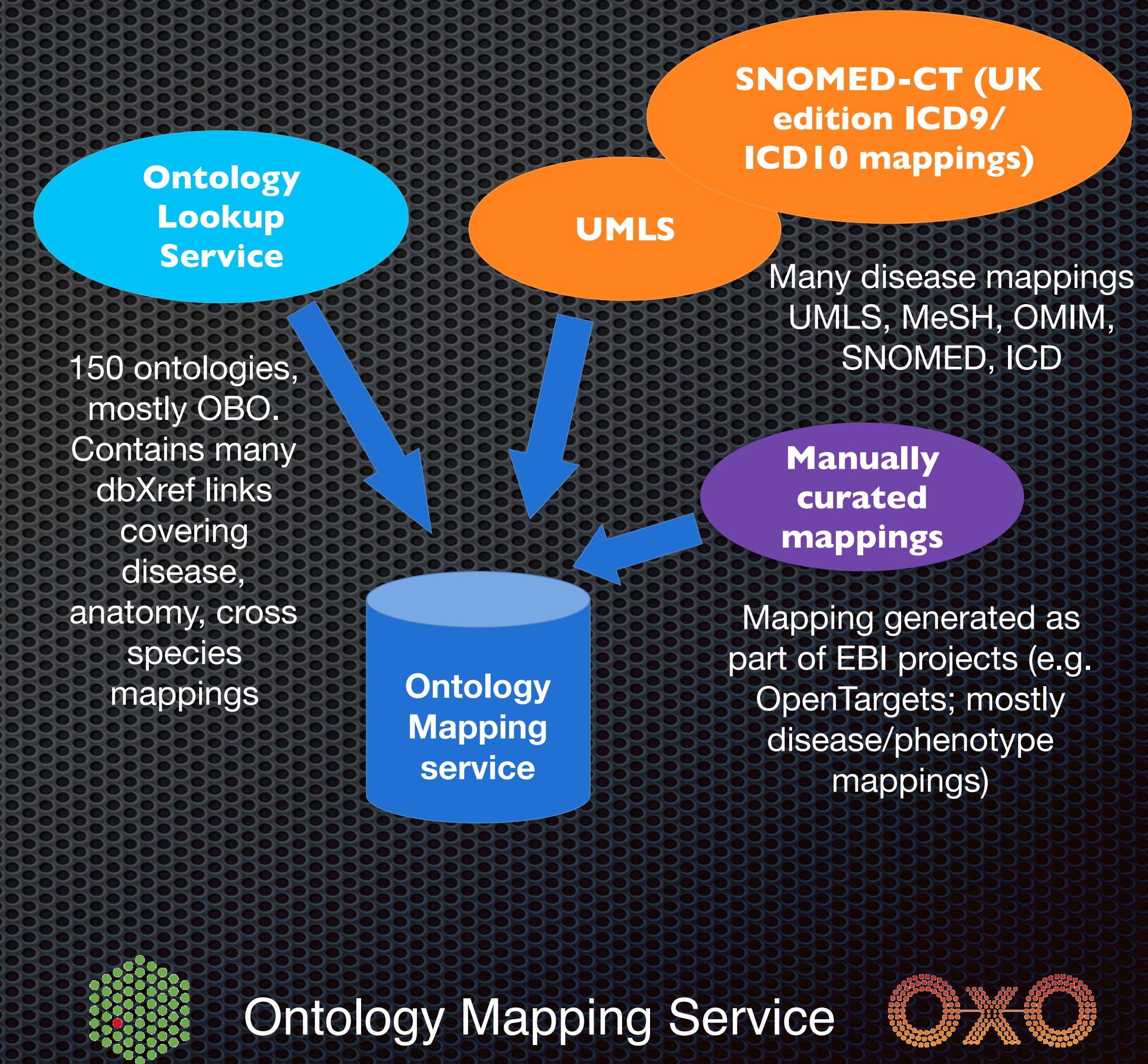
DRAGTS NCI:038 NCI:DRM10 MORTHOLOGY038 ICD10GR038
GSM393858 C2853 Acute Myeloid Leukemia Not Otherwise Specified9861/3 C42
GSM302285 C2852 Adenocarcinoma 8140/3 C34
GSM918983 C3222 Medulloblastoma 9480/3 C716
GSM551398 C4017 Ductal Breast Carcinoma 8500/3 C50
GSM412384 C3163 Chronic Lymphocytic Leukemia 9823/3 C42
GSM1218286 C4017 Ductal Breast Carcinoma 8500/3 C50
GSM714412 C2852 Adenocarcinoma 8140/3 C569
GSM1109923 C9306 Soft Tissue Sarcoma 8800/3 C499
GSM711848 C2852 Adenocarcinoma 8140/3 C25
GSM746294 C89426 8022/2 C53
GSM1981528 C4017 Ductal Breast Carcinoma 8500/3 C50
GSM281399 C8949 8500/2 C50
GSM533469 C9349 Plasmacytoma 9831/3 C42



Making Ontologies Work for GA4GH Implementation Studies

- biomedical "metadata" in different resources frequently follows incompatible classification systems
- medical coding systems are driven by different paradigms compared to biological ontologies (e.g. for cross-species comparisons)
- frequently used classifications (ICD, Snomed...) are either not "ontologised" or cannot be referenced in open resources

Federated queries across resources need **curated mappings** of classifications/ontologies





GA4GH::SchemaBlocks

An Initiative by Members of the Global Alliance for Genomics and Health

- [News](#)
- [Participants](#)
- [Data Formats](#)
- [Data Schemas](#)
- [Ontology terms](#)
- [Examples, Guides & FAQ](#)
- [Meeting minutes](#)
- [Contacts](#)
- [Related Sites](#)
 - GA4GH::Discovery
 - GA4GH::CLP
 - GA4GH::GKS
 - SchemaBlocks at Metadata
 - ELIXIR Beacon
 - Phenopackets
 - GA4GH
 - Beacon+
- [Tags](#)
 - Beacon
 - CP
 - admins
 - coordinates



Schemas

Schema elements previously developed as part of various GA4GH efforts had been assembled in the [SchemaBlocks demonstrator](#). Those schemas and documentation will be re-implemented in this space.

GA4GH::SchemaBlocks

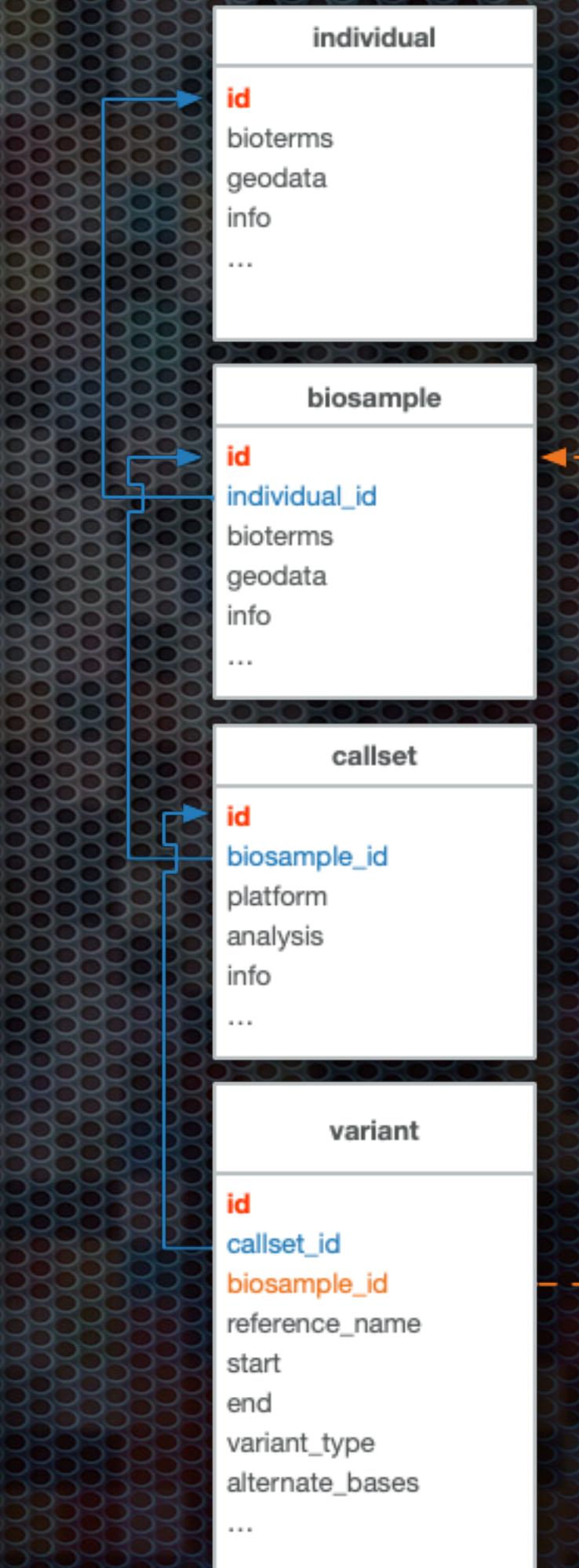
An Initiative by Members of the Global Alliance for Genomics and Health

- [News](#)
- [Participants](#)
- [Data Formats](#)
- [Identifiers and CURIES](#)
- [Genome Coordinates](#)
- [Dates & Times](#)
- [Data Schemas](#)
- [Examples, Guides & FAQ](#)
- [Meeting minutes](#)
- [Contacts](#)
- [Related Sites](#)
 - GA4GH::Discovery
 - GA4GH::CLP
 - GA4GH::GKS
 - SchemaBlocks at Metadata
 - ELIXIR Beacon
 - Phenopackets
 - GA4GH
 - Beacon+
- [Tags](#)
 - Beacon
 - CP
 - admins
 - coordinates
 - Discovery
 - GA4GH
 - GKS
 - MME
 - code
 - contributors
 - leads
 - press
 - times
 - contacts
 - dates
 - developers
 - identifiers



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

CODED BY: JONATHAN WOODWARD





ELIXIR - Towards Biomedical Beacons

Needs & Models Beyond Basic Variant Discovery

Michael Baudis, University of Zurich | **SIB**

Basel 2019



Global Alliance
for Genomics & Health

Beacon Project

An open web service that tests the willingness of international sites to share genetic data.



Beacon Network

Search Beacons

Search all beacons for allele

GRCh37 ▾ 10:118969015 C / CT Search

Response All None
 Found 16
 Not Found 27
 Not Applicable 22

Organization All None
 AMPLab, UC Berkeley
 BGI
 BioReference Laborato...
 Brazilian Initiative on ...
 BRCA Exchange
 Broad Institute
 Centre for Genomic R...
 Centro Nacional de A...
 Curoverse
 EMBL European Bio...
 Global Alliance for G...
 Google
 Institute for Systems ...
 Instituto Nacional de ...

Response	All	None
<input checked="" type="checkbox"/> Found	16	
<input type="checkbox"/> Not Found	27	
<input type="checkbox"/> Not Applicable	22	

BioReference BioReference Hosted by BioReference Laboratories Found

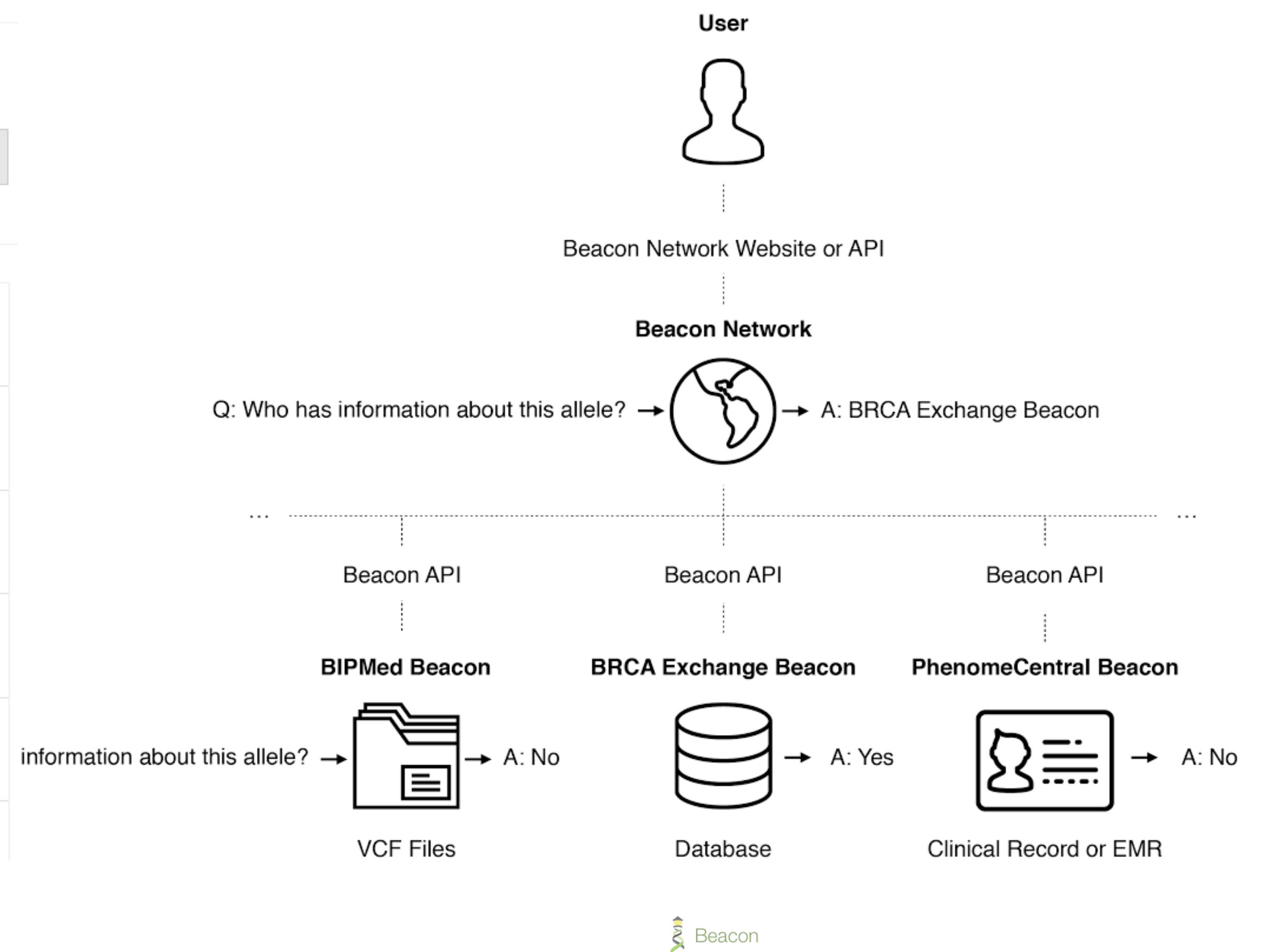
Catalogue of Somatic Mutations in Cancer Catalogue of Somatic Mutations in Cancer Hosted by Wellcome Trust Sanger Institute Found

Cell Lines Cell Lines Hosted by Wellcome Trust Sanger Institute Found

Conglomerate Conglomerate Hosted by Global Alliance for Genomics and Health Found

COSMIC COSMIC Hosted by Wellcome Trust Sanger Institute Found

dbGaP: Combined GRU Catalog and NHLBI Exome Seq... dbGaP: Combined GRU Catalog and NHLBI Exome Seq... Found



Date	Tag	Title
2018-01-24	v0.4.0	Beacon
2016-05-31	v0.3.0	Beacon

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 the 'Driver Projects' section of the GA4GH website. It features a red circular icon with a white rocket ship. Below it, the text 'Driver Projects' is displayed. A detailed description follows: '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 local contexts.' To the right, there is a box for the 'ELIXIR Beacon' project, which includes the ELIXIR logo, the text 'ELIXIR Beacon', the URL 'www.elixir-europe.org', the location 'Europe', and the 'Champions: Serena Scollen, Ilkka Lappalainen, Michael Baudis'.

Beacon forward



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

Beacon⁺

This forward looking Beacon interface implements additional, planned features.

Query

Dataset	tcga
Reference name*	9
Genome Assembly*	GRCh38 / hg38
Start min Position*	19,500,000
Start max Position	21,975,098
End min Position	21,967,753
End max Position	24,500,000
Alt. Base(s)*	DEL
Bio-ontology	icdot:c50.9: (4065)

Beacon Implementations

- implementing existing resources with Beacon protocol
- e.g. TCGA cancer variants (structural and SNV)

Beacon Response

- quantitative (counts for variants, callsets and samples)
- *Handover* to authentication system for data retrieval
- **no exposure** of data beyond standard Beacon response and additional pointer to matched data

Prototyping Query Extensions

- testing e.g. bio-metadata queries using ontology terms

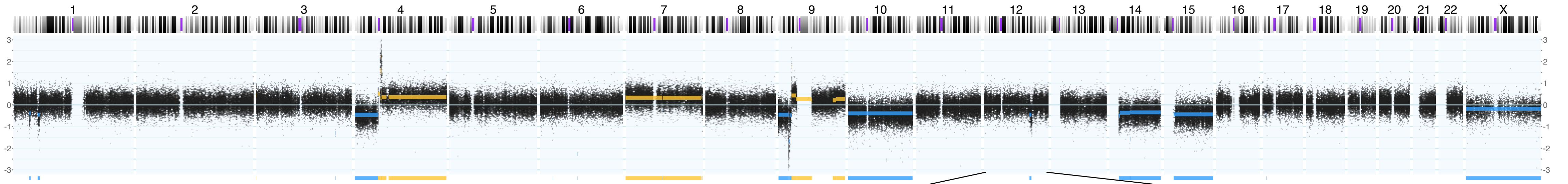
Dataset	Assembly	Chro	Start Range	End Range	Pos	Ref Alt	Bio Query	Variants Calls Samples	f _{alleles}	Response Context
tcga	hg38	9	19,500,000 21,975,098	21,967,753 24,500,000		DEL	icdot:c50.9	54 54 54	0.0243	JSON UCSC Handover

arrayMap progenetix This Beacon implementation is developed by the Computational Oncogenomics Group at the University of Zurich, with support from the SIB Technology group and ELIXIR.

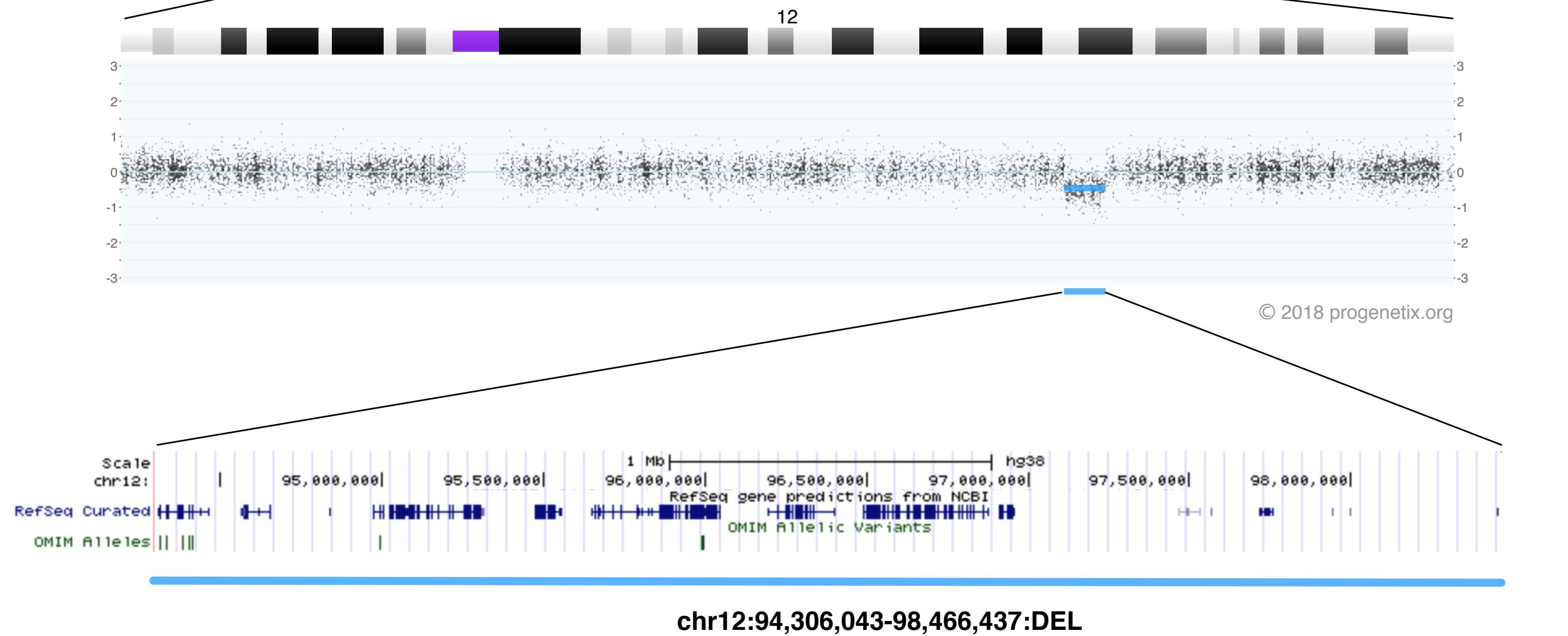
University of Zurich UZH ELIXIR SIB



GSM491153



- Beacon+ **range queries** allow the definition of a genome region of interest, containing a specified variant (or other mappable feature)
- “fuzzy” matching of region ends is essential for features without base specific positions
- current Beacon implementation addresses CNV (<DUP>,), as are specified in VCF && GA4GH variant schema



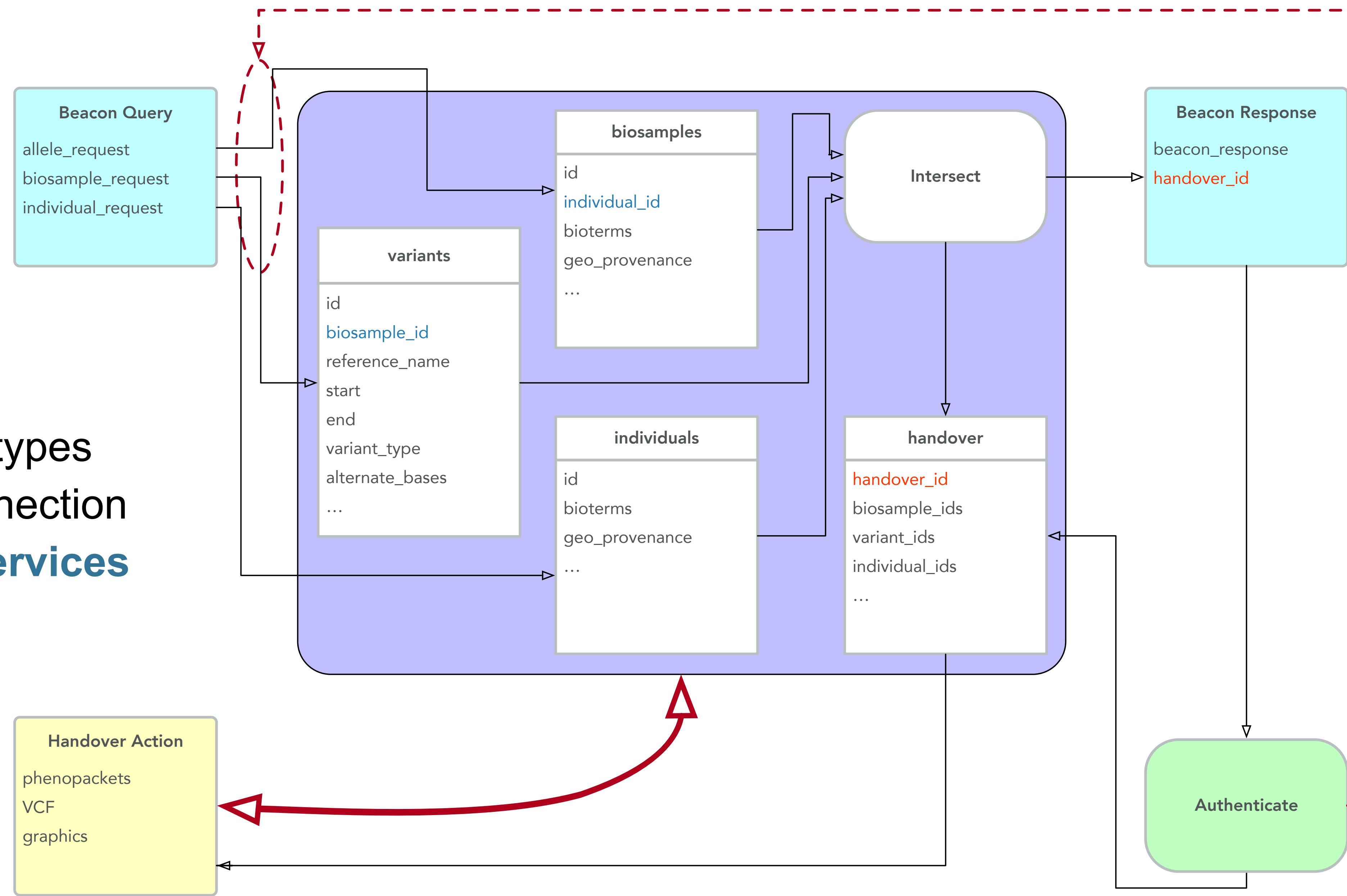
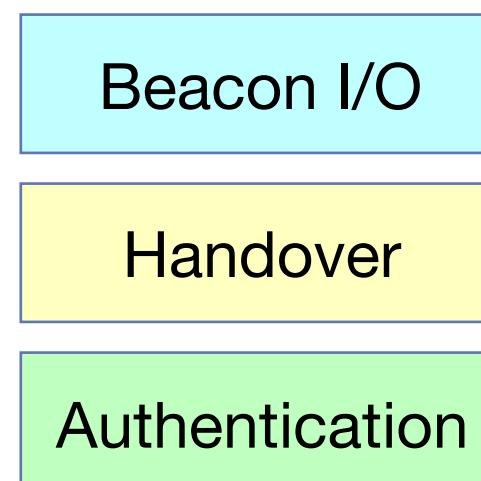
start_min: 94,000,000
start_max: 94,500,000
variant_type: "BND"

reference_name: "9"
variant_type: "DEL"

end_min: 98,200,000
end_max: 98,700,000
variant_type: "BND"

Beacon & Handover

Future Beacons to support advanced types of queries and connection of **data delivery services**



Beacon query => Handover Handle => Authentication => Data Retrieval



Beacon+ example implementation using public somatic variation data

Beacon+

This is an implementation of a Beacon "handover" concept, in which a Beacon query response additionally describes a representation of the query results (i.e. callsets, biosamples, metadata ...), which can then be accessed after ("yes"|"no") or quantitative ("n matches") Beacon response from a data delivery mechanism.

The current implementation exemplifies some possible scenarios:

- providing a histogram of regional gain/loss frequencies (DUP, DEL) for samples with structural variations
- returning data of the associated callsets which matched the Beacon query (this is for feature demonstration)
- returning the metadata (diagnoses etc.) of the biosamples from which the matching callsets were derived

This demonstrator does not implement authentication procedures yet; login & password fields can be left empty.

Handover Action

Plot DUP/DEL histogram

Export Callset Data

Export Biosample Data

Login

Password

••••••••••••••

Process Data

Global Alliance
for Genomics & Health

Beacon Handover Demonstrator

- only exposure of access handle to data stored in secure system
- here one-step authentication and selection of *handover* action; other scenarios possible / likely
- *handover response formats not managed by Beacon protocol - phenopackets, VCF, htsget ...*

Beacon Handover Demonstrator

- only exposure of access handle to data stored in secure system
 - here one-step authentication and selection of *handover* action; other scenarios possible / likely
 - *handover response outside of Beacon protocol / system*
- ```
{
 "datasetAlleleResponses": [
 {
 "callCount": 588,
 "variantCount": 588,
 "datasetHandover": [
 {
 "description": "retrieve data of the biosamples matched by the query",
 "handoverType": {
 "id": "pgx:handover:biosamplesdata",
 "label": "Biosamples"
 },
 "url": "https://beacon.progenetix.org/beaconplus-server/
beacondeliver.cgi?do=biosamplesdata&accessid=d50b0a14-1a2c-11e9-af9a-fd65cc50531f"
 },
 {
 "url": "https://beacon.progenetix.org/beaconplus-server/
beacondeliver.cgi?do=cnvhistogram&accessid=d513a3c3-1a2c-11e9-af9a-fa62516f26af",
 "description": "create a CNV histogram from matched callsets",
 "handoverType": {
 "id": "pgx:handover:cnvhistogram",
 "label": "CNV Histogram"
 }
 },
 ...
]
 }
]
}
```



# Check it Out!

- managed, participation driven projects living on Github:  
**ga4gh-schemablocks** & **ga4gh-beacon**
- test datasets & code available through our  
**progenetix** &  
**baudisgroup** repositories
  - test
  - comment
  - suggest
  - propose
  - complain ...

ga4gh / ga4gh-schemas      Unwatch 108 ★ Star 212 Fork 113

progenetix / beaconplus-server      Unwatch 3 ★ Star 0 Fork 0

progenetix / beaconplus-ui      Unwatch 5 ★ Star 0 Fork 0

ga4gh / beacon-team      Unwatch 36 ★ Star 15 Fork 15

progenetix / arraymap2ga4gh      Unwatch 5 ★ Star 2 Fork 1

No dependencies

GA4GH

Branches

This branch

Branch: master New pull request

85 commits 2 branches 0 releases 3 contributors

KyleGao Multi genome editions      Latest commit 5db07db 6 days ago

data Multi genome editions      6 days ago

examples update DIPG examples      2 months ago

tools remove the per scripts => beaconplus-server      2 months ago

README.md link      2 months ago

schema.pdf updated schema diagram      6 months ago

README.md

## Implementation of the GA4GH schema based on genome profiles and metadata from arrayMap

This repository will contain data and information regarding the arrayMap based implementation of a GA4GH schema structure. While it is not expected that GA4GH compliant resources mirror the schema in their internal structure, this project is aimed at showing the principle feasibility of such an approach, mainly to test & drive schema development.

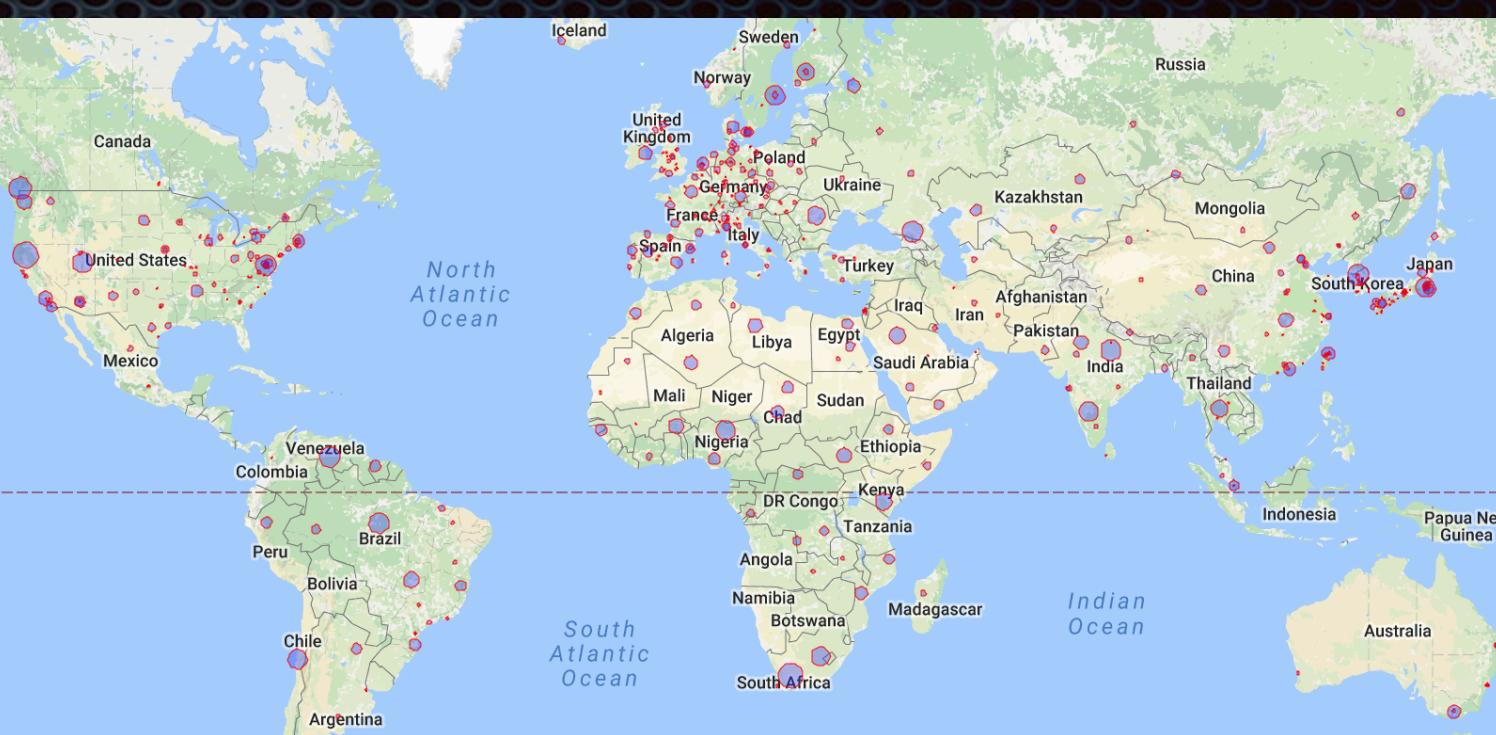
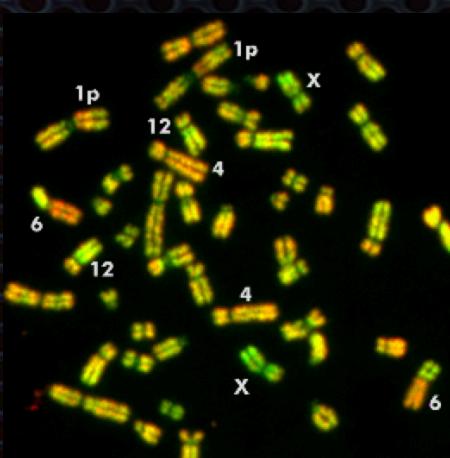
Data & schemas represented here are not kept in a stable/versioned status, but are updated together with or anticipating GA4GH schema changes.

Global Alliance for Genomics & Health



# Cancer Genome Data: Where Do We Need To Go

- balancing clinical medicine (**panel** sequencing of **actionable** cancer targets) and the necessary knowledge generation (**complete genomes**, multi"omes", genetic background)
- "genetic **awareness**" and regulatory **security**
- vastly increasing **data curation efforts** to make best use of existing data
- data **sharing frameworks** & **federated** analysis
- **everywhere in the world...**



## BAUDISGROUP @ UZH

NI AI  
MICHAEL BAUDIS  
(HAOYANG CAI)  
PAULA CARRIO CORDO  
BO GAO  
QINGYAO HUANG  
SAUMYA GUPTA  
(NITIN KUMAR)  
(RAHEL PALOOTS)

## SIB

AMOS BAIROCH  
HEINZ STOCKINGER  
DANIEL TEIXEIRA

THOMAS EGGERMANN  
ROSA NOGUERA  
REINER SIEBERT  
CAIUS SOLOVAN



University of  
Zurich<sup>UZH</sup>



Global Alliance  
for Genomics & Health



## GA4GH

LARRY BABB  
ANTHONY BROOKES  
MELANIE COURTOT  
MELISSA HAENDEL  
MICHAEL MILLER  
HELEN PARKINSON  
GUNNAR RÄTSCH  
ANDY YATES

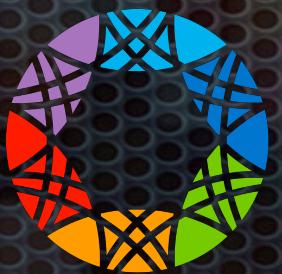
## ELIXIR & CRG

JORDI RAMBLA DE ARGILA  
GARY SAUNDERS  
ILKKA LAPPALAINEN  
S. DE LA TORRE PERNAS  
SERENA SCOLLEN  
JUHA TÖRNROOS



University of  
Zurich<sup>UZH</sup>

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



Global Alliance  
for Genomics & Health



[arraymap.org](http://arraymap.org)  
[progenetix.org](http://progenetix.org)  
[info.baudisgroup.org](http://info.baudisgroup.org)  
[sib.swiss/baudis-michael](http://sib.swiss/baudis-michael)  
[imls.uzh.ch/en/research/baudis](http://imls.uzh.ch/en/research/baudis)  
[beacon-project.io](http://beacon-project.io)  
[schemablocks.org](http://schemablocks.org)



