



University of  
Zurich<sup>UZH</sup>

Department of Molecular Life Sciences



# Data Mining in Genomics

**Resources | Standards | Protocols | Tools | Discourse  
for Genomic Research and Personalised Health Strategies**

Prof. Dr. Michael Baudis

Department of Molecular Life Sciences

University of Zurich

**SIB** | Swiss Institute of Bioinformatics

2020-03-19



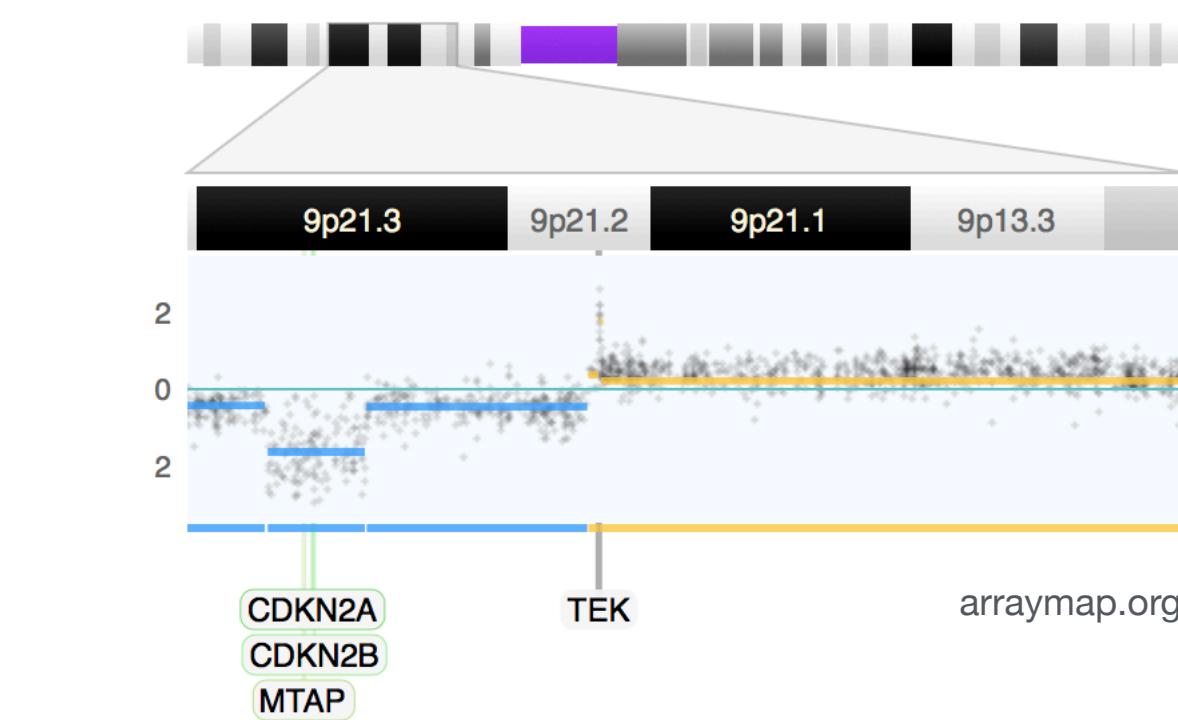
**Global Alliance**  
for Genomics & Health





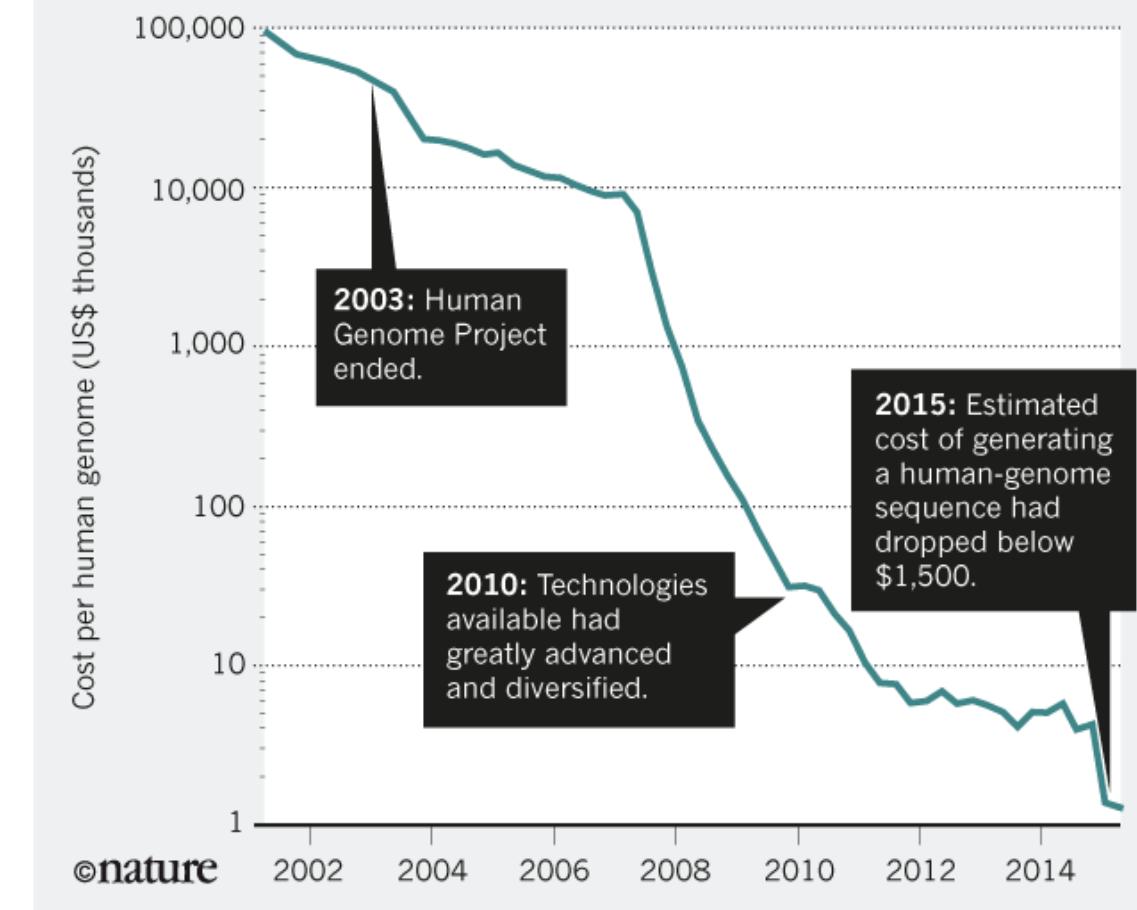
## Genome screening at the core of “Personalised Health”

- ▶ **Genome analyses** (including transcriptome, metagenomics) are core technologies for Personalised Health™ applications
- ▶ The unexpectedly large amount of **sequence variants** in human genomes - germline and somatic/cancer - requires huge analysis efforts and creation of **reference repositories**
- ▶ **Standardized data formats** and **exchange protocols** are needed to connect these resources throughout the world, for reciprocal, international **data sharing** and **biocuration** efforts
- ▶ Our work @ UZH:
  - ▶ **cancer genome repositories**
  - ▶ **biocuration**
  - ▶ **protocols & formats**

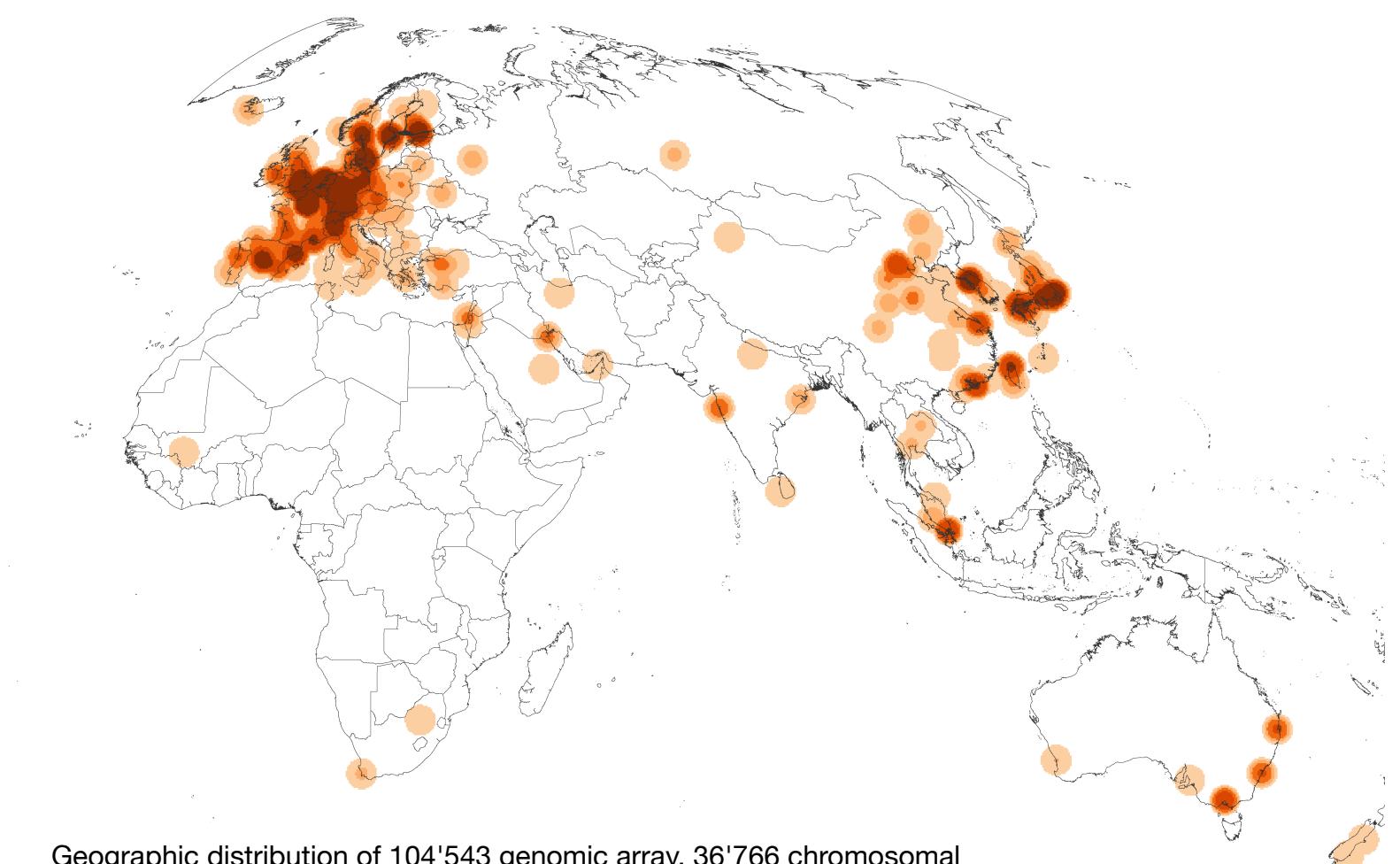
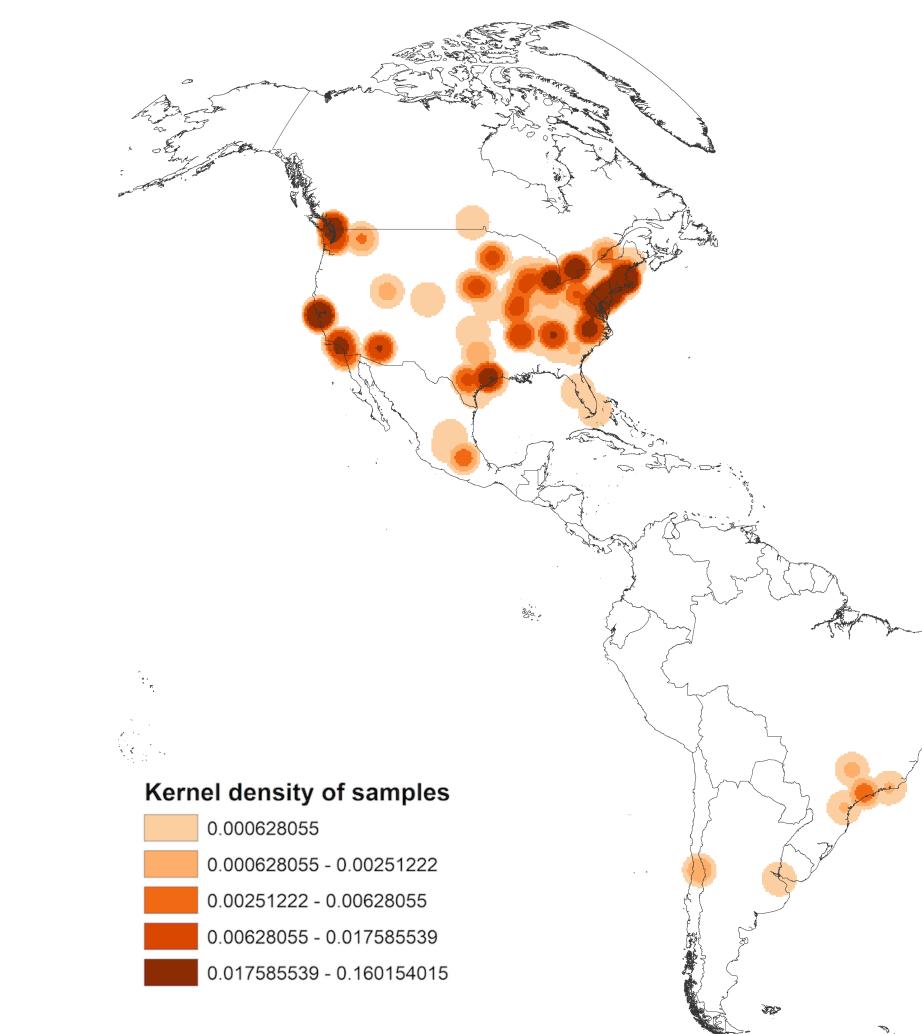


### BETTER, CHEAPER, FASTER

The cost of DNA sequencing has dropped dramatically over the past decade, enabling many more applications.



The future of DNA sequencing. Eric D. Green, Edward M. Rubin & Maynard V. Olson. Nature; 11 October 2017 (News & Views)



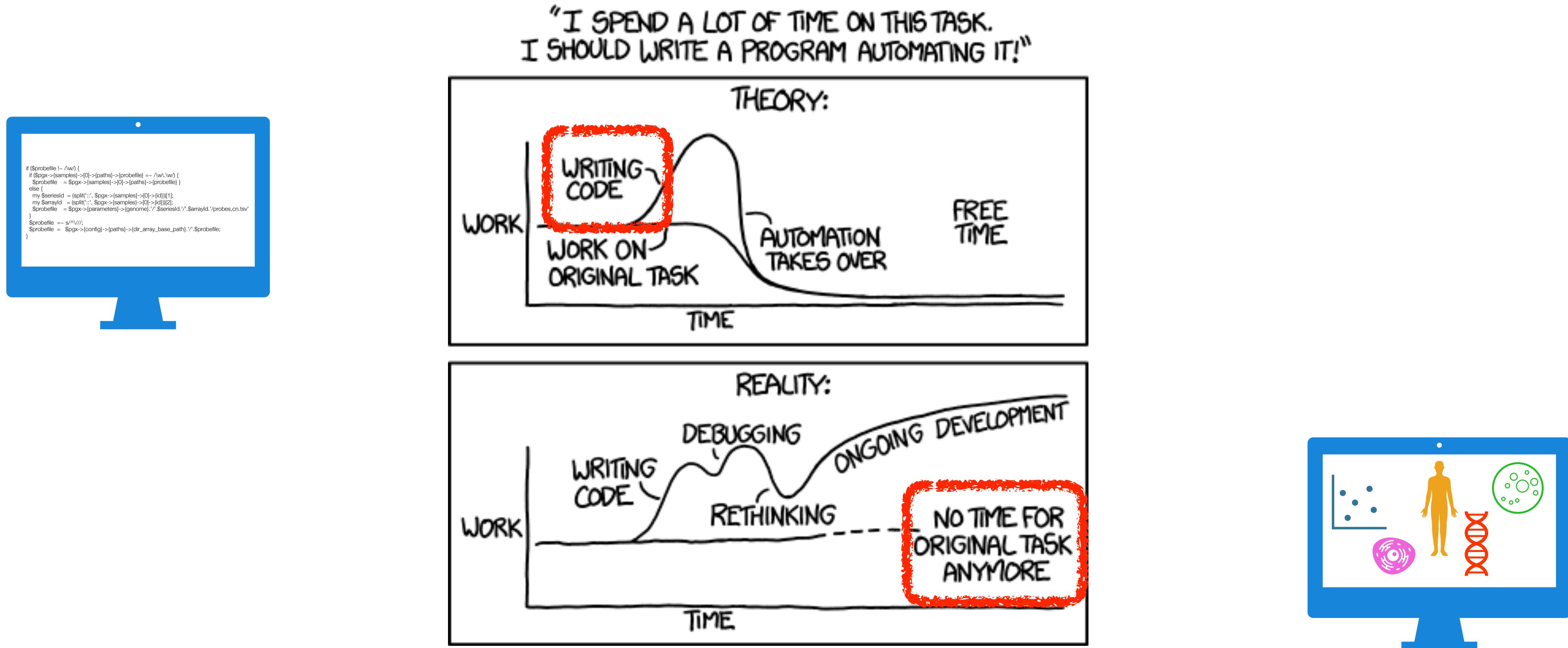
Geographic distribution of 104'543 genomic array, 36'766 chromosomal CGH and 15'409 whole genome/exome based cancer genome datasets

# {bio\_informatics\_science}

---



# {bio\_informatics\_science}





## Our contributions I: Cancer genome knowledge resources and research



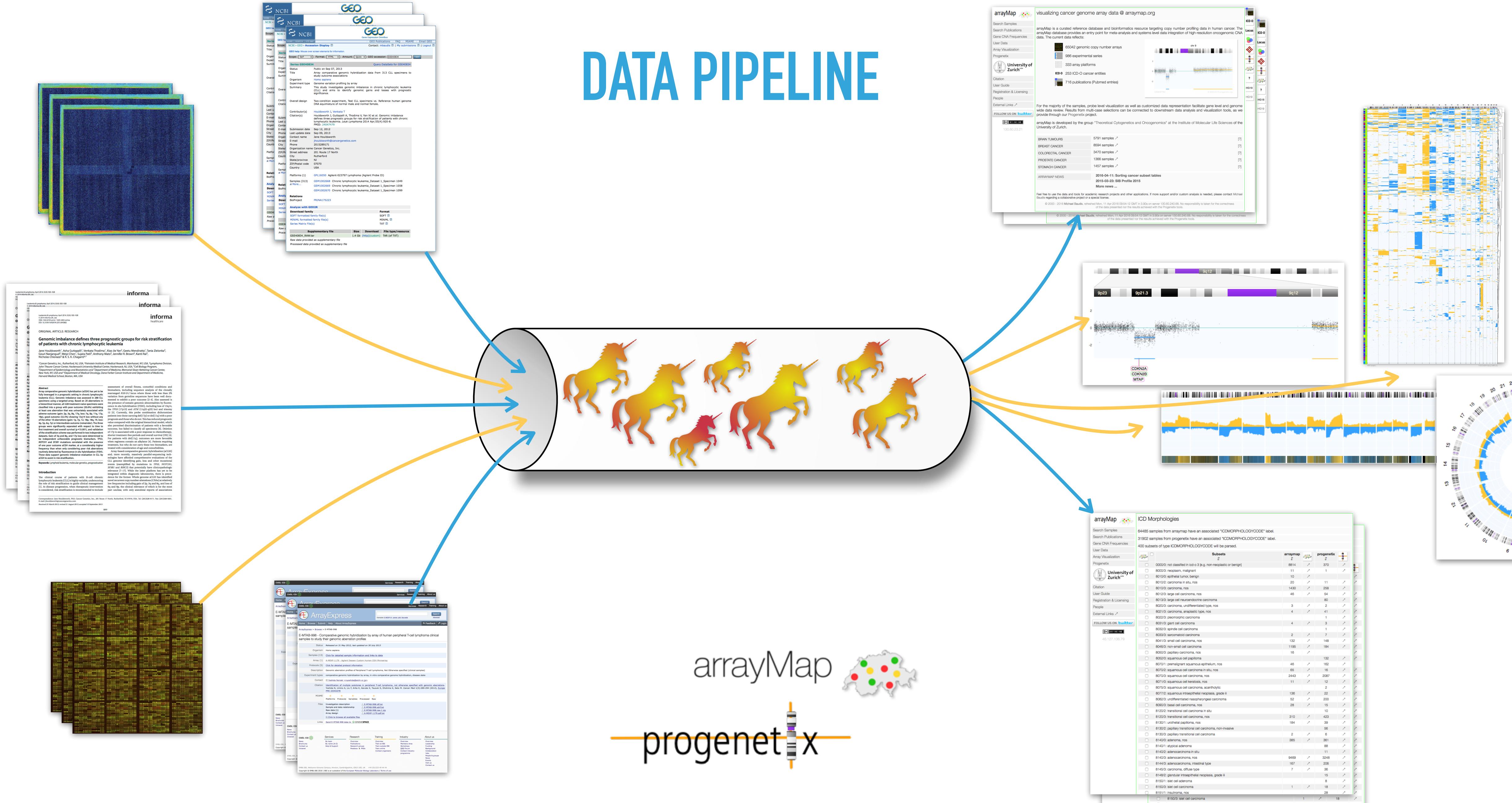
- ▶ curated cancer genome publication resource (more than 3200 manually curated articles)
- ▶ article metadata
- ▶ annotated genome profiles of >100'000 samples
- ▶ ontology mappings; clinical data where available
- ▶ epistemology: geographic and histologic sampling biases



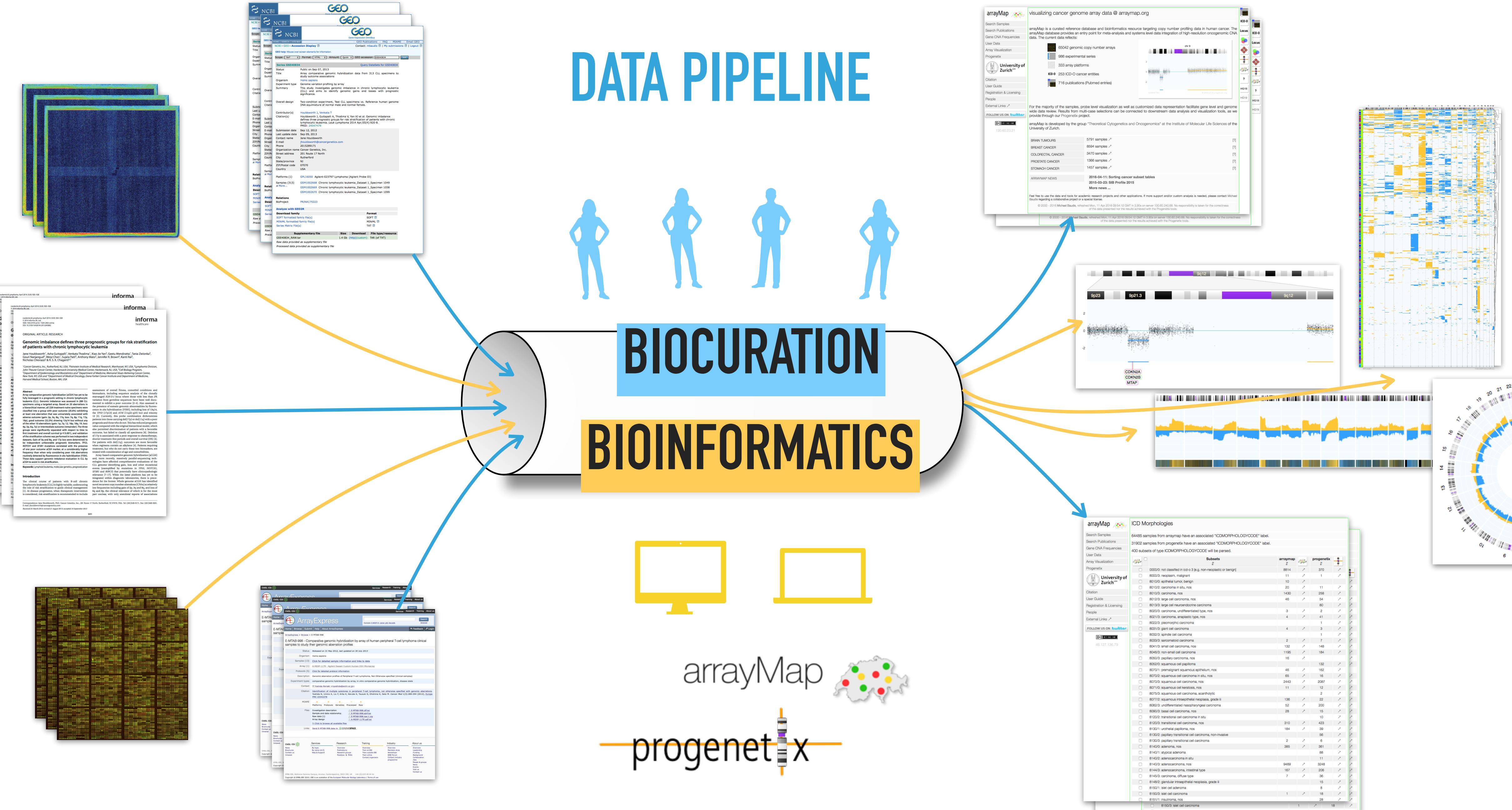
- ▶ more than 70'000 array based genome profiles
- ▶ probe level, copy number, metadata
- ▶ completely open data access through web interface, downloads and API calls
- ▶ re-annotated metadata (diagnostic coding, basic clinical) for all samples



# DATA PIPELINE



# DATA PIPELINE

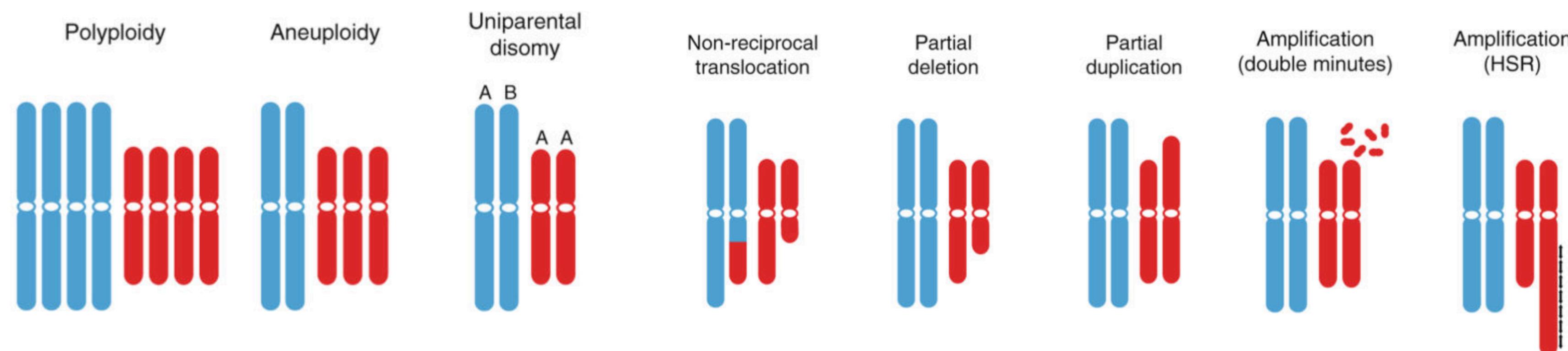


# Introduction

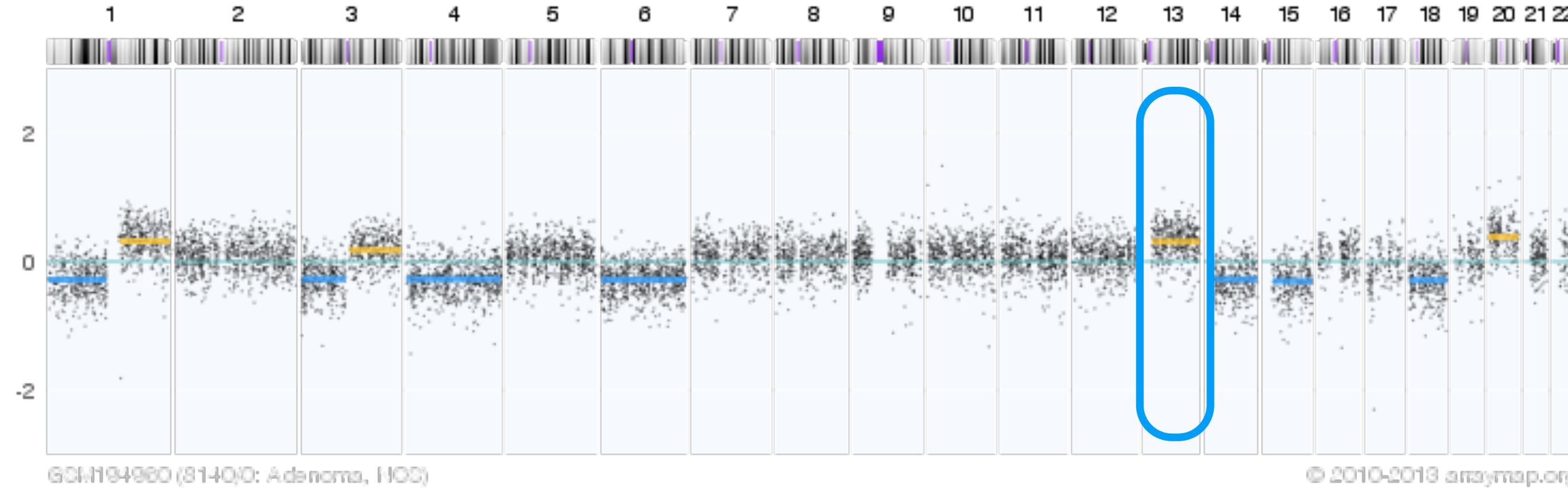
## Types of genomic alterations in Cancer

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

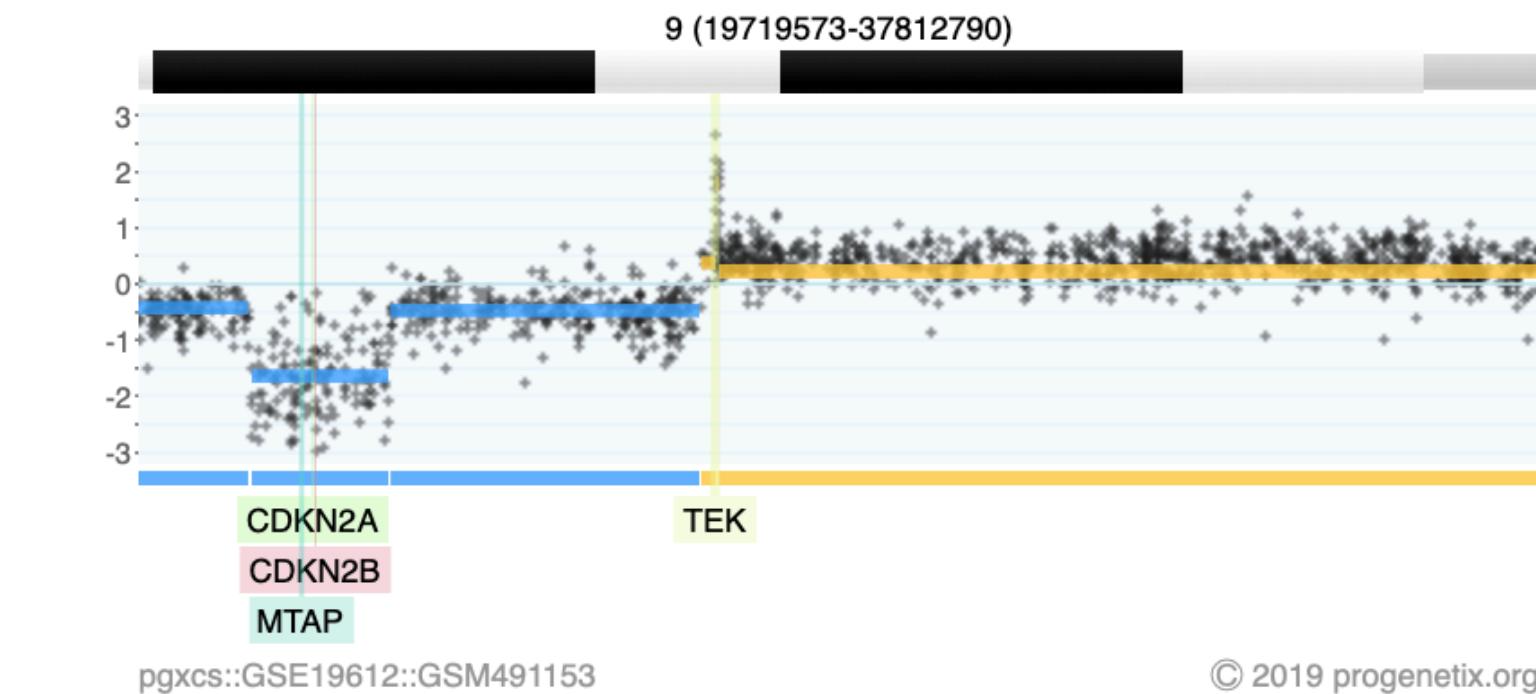
## Imbalanced Chromosomal Changes: CNV



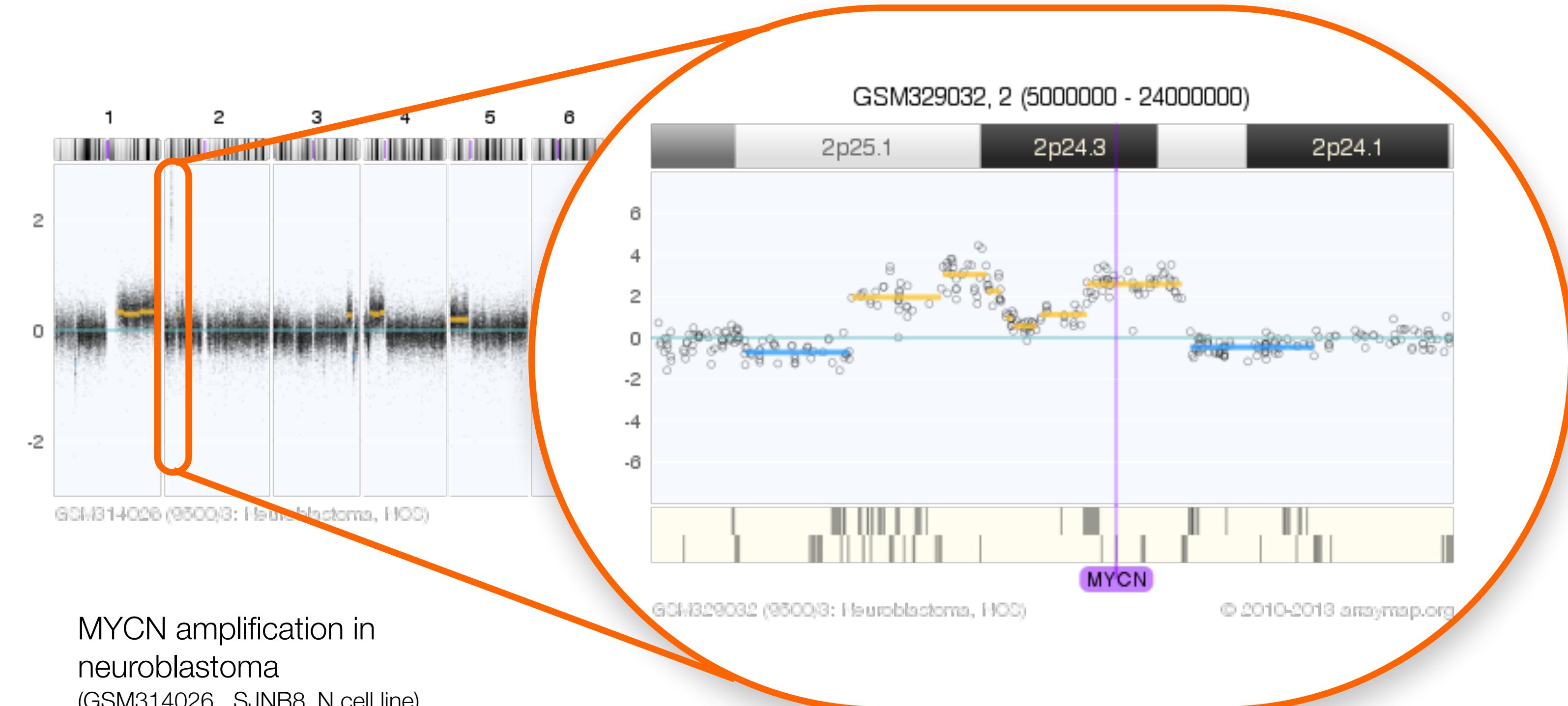
# Somatic Copy Number Variations



Gain of chromosome arm 13q in colorectal carcinoma



2-event, homozygous deletion in a Glioblastoma



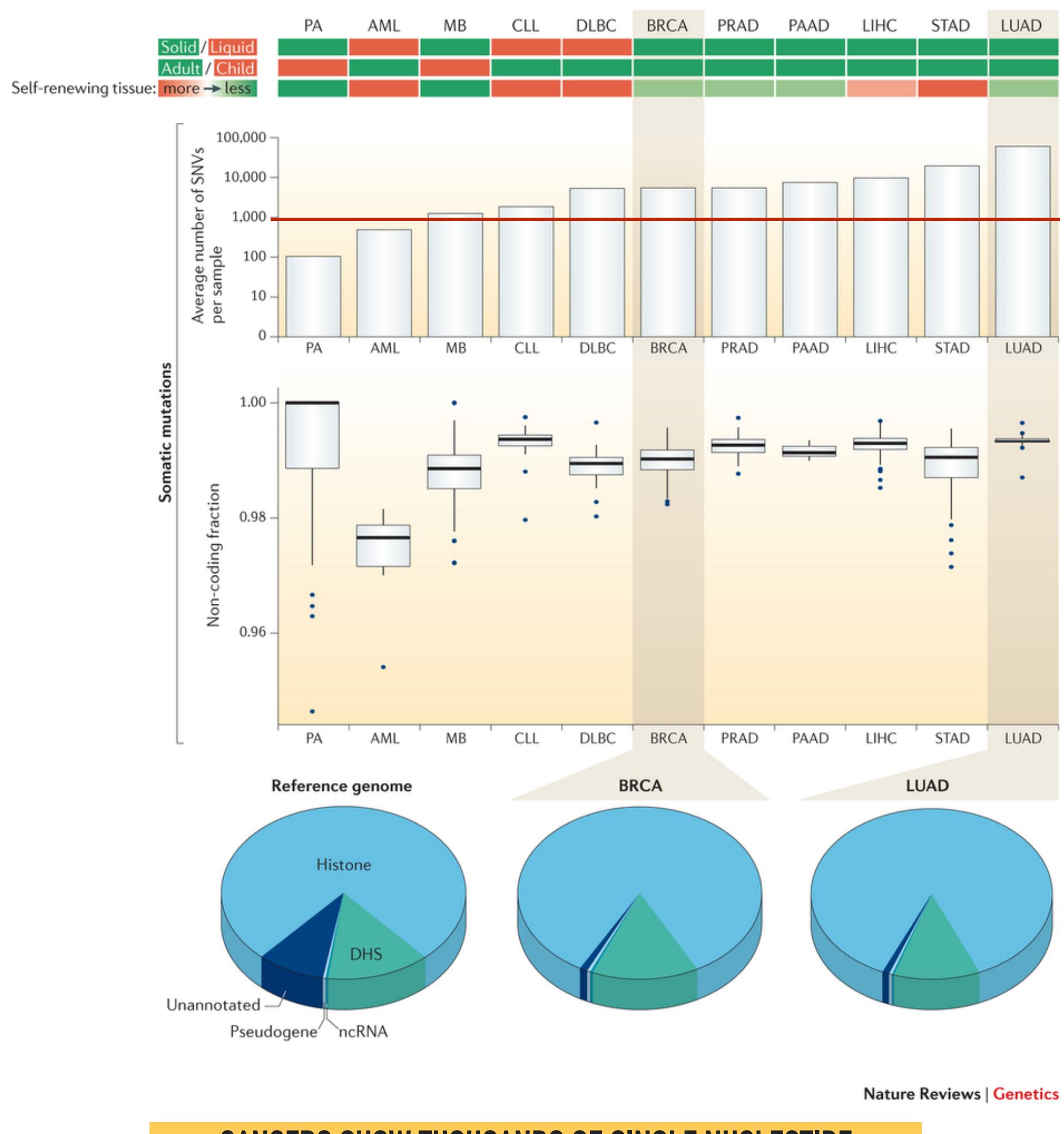
MYCN amplification in neuroblastoma  
(GSM314026, SJNB8\_N cell line)

**low level/high level** copy number alterations (CNAs)

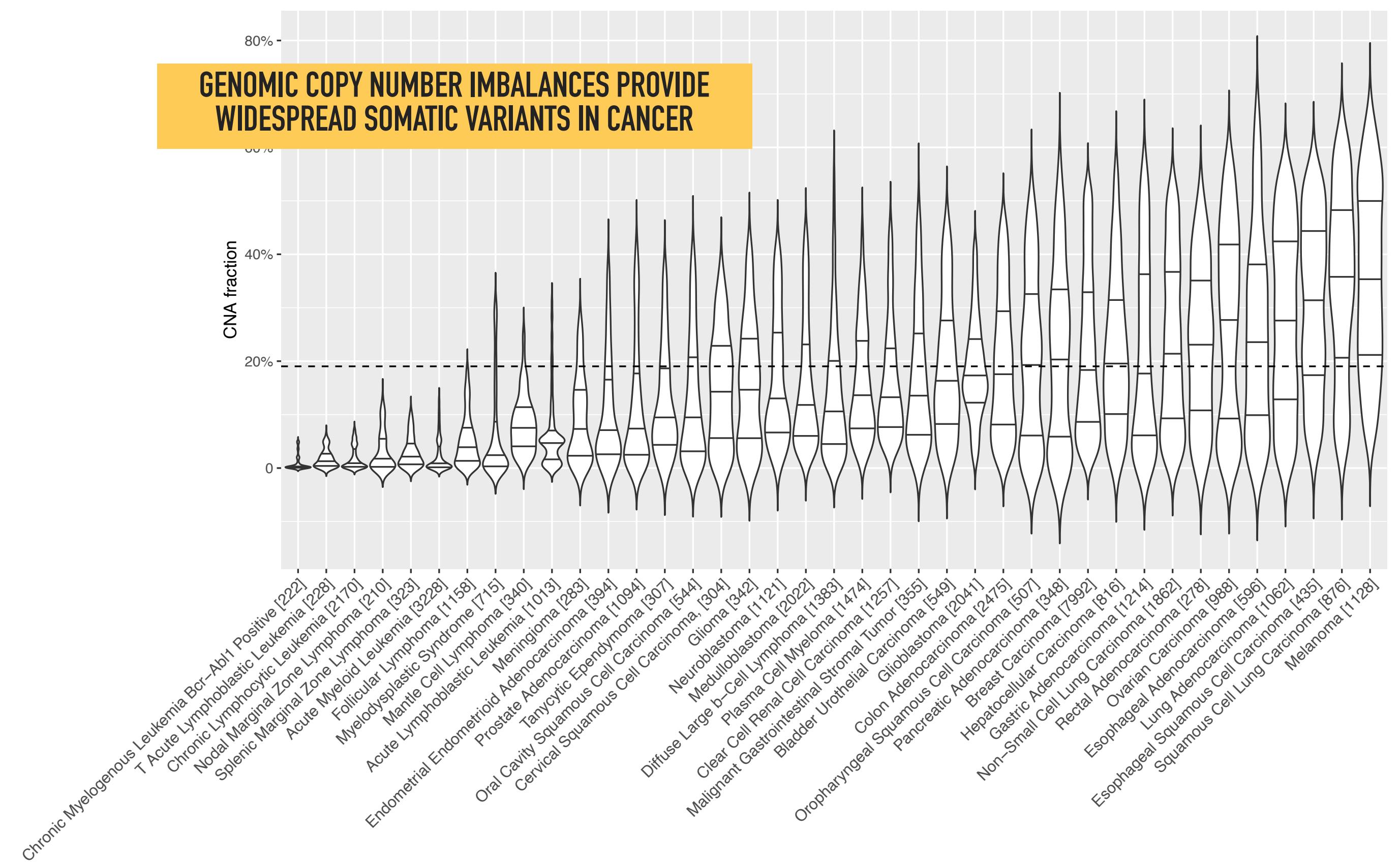
arrayMap



# 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))

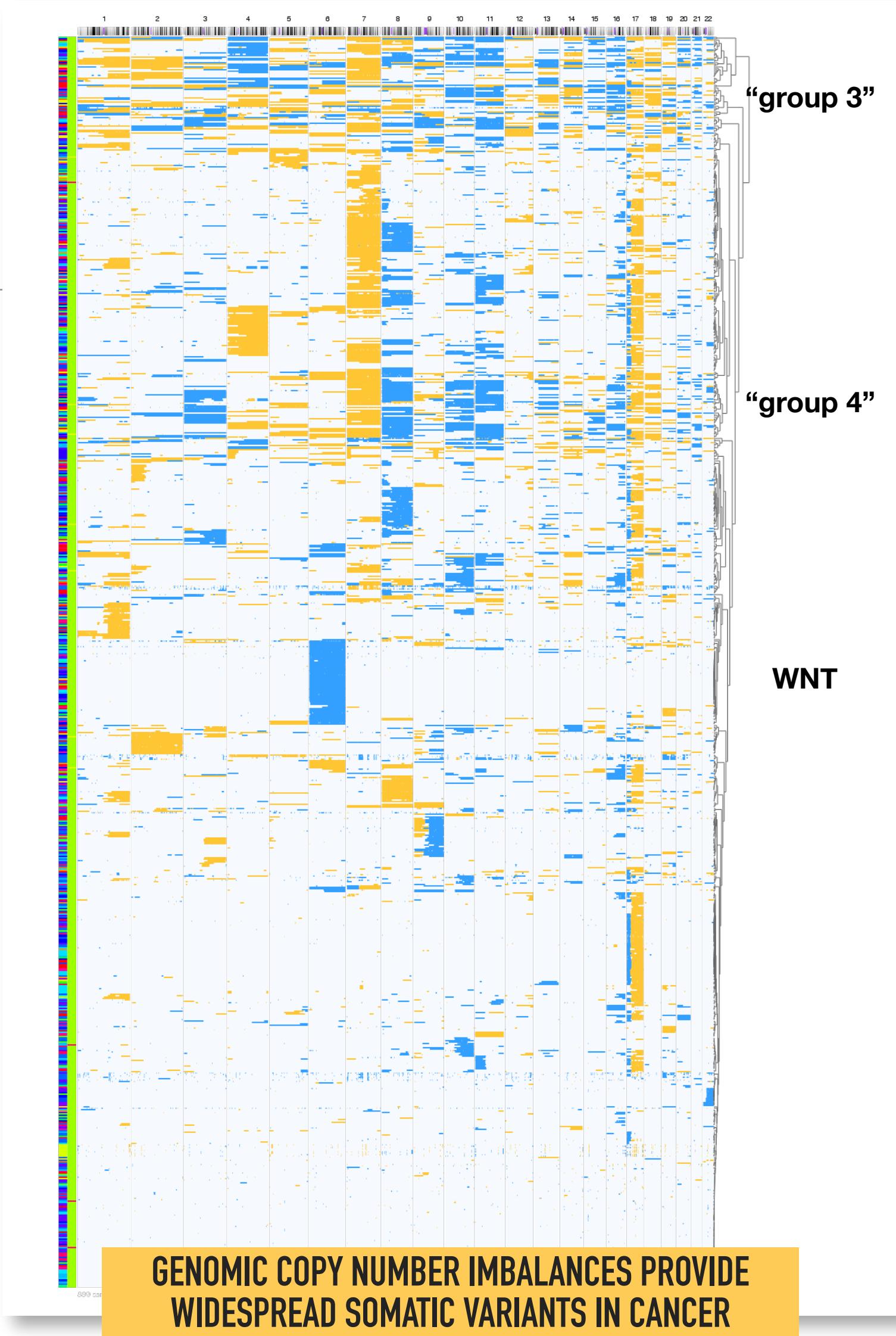
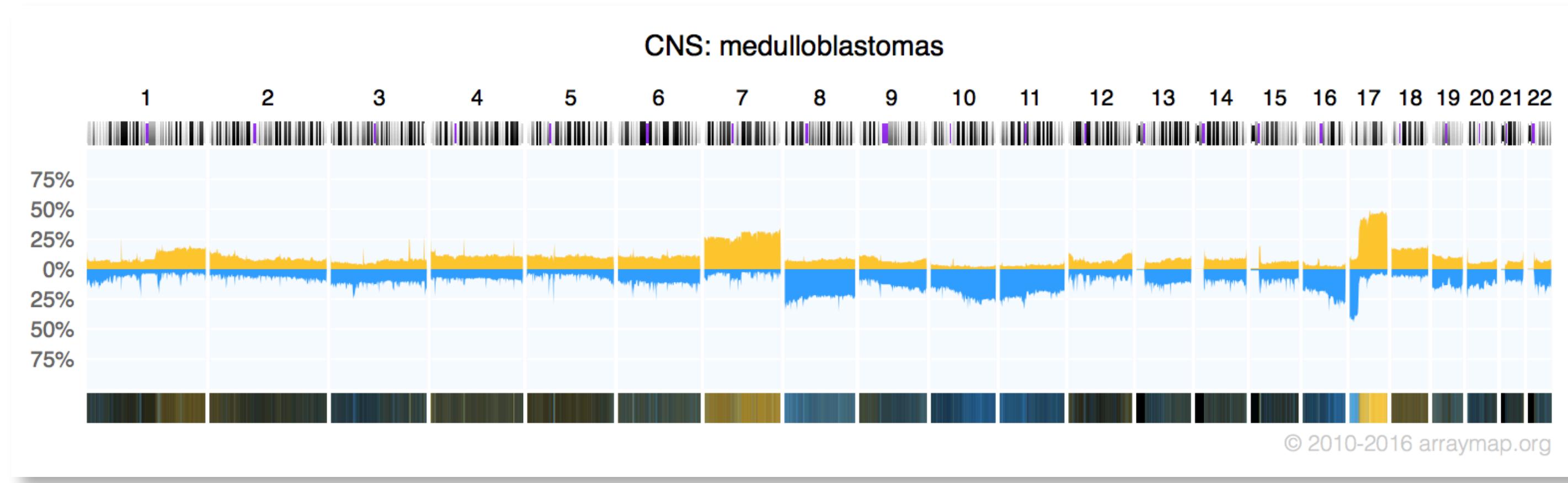


On average ~19% of a cancer genome are in an imbalanced state (more/less than 2 alleles); Original data based on 43654 cancer genomes from [progenetix.org](http://progenetix.org)

# Somatic CNVs In Cancer: Patterns

Many tumor types express **recurrent mutation patterns**

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



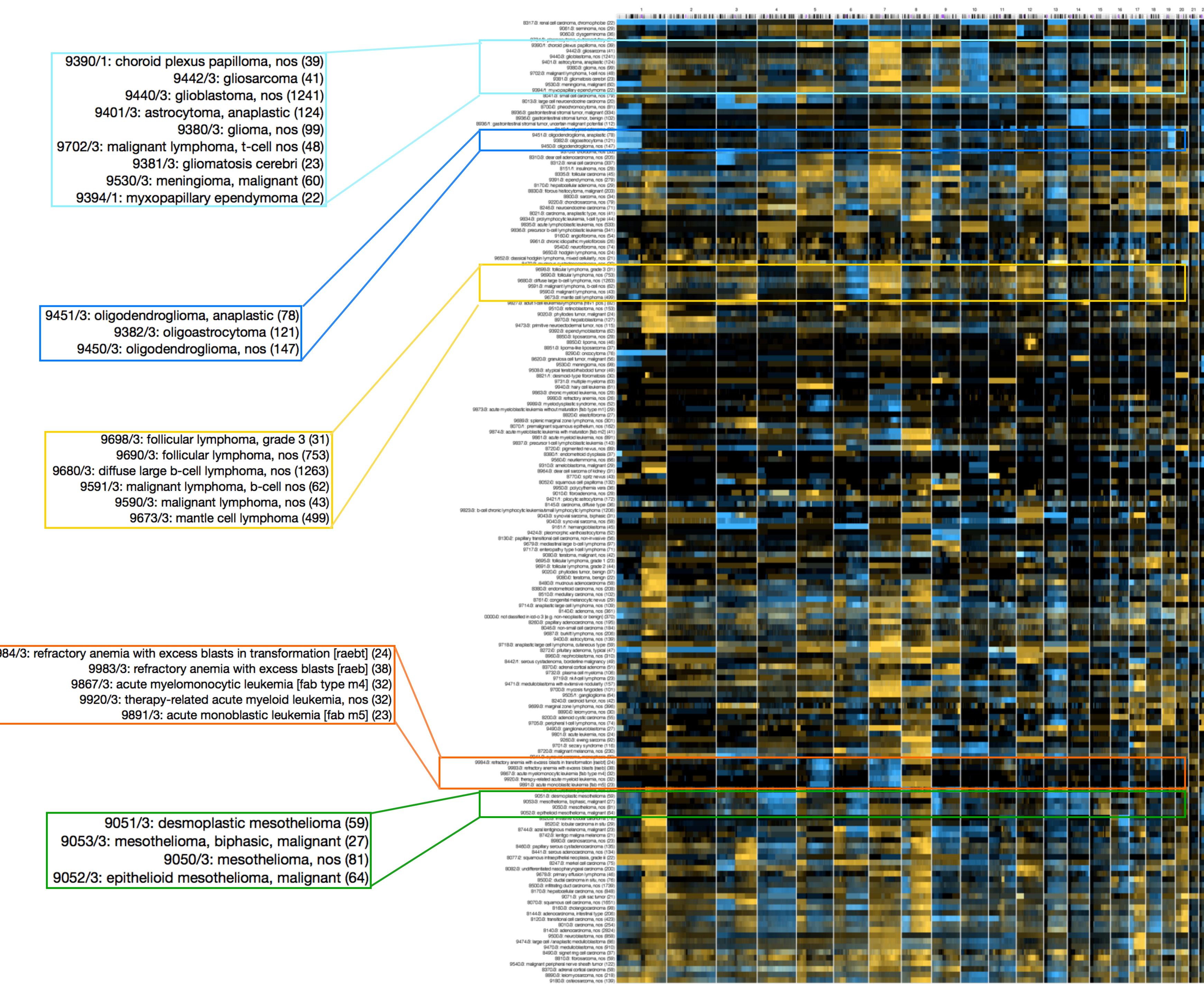
A genomic copy number histogram for malignant medulloblastomas, the most frequent type of pediatric brain tumors, displaying regions of genomic duplications and deletions. These can be decomposed into individual tumor profiles which segregate into several clusters of related mutation patterns with functional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



# Somatic Mutations In Cancer: Patterns

## Making the case for genomic classifications

Some related cancer entities show similar copy number profiles



# 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



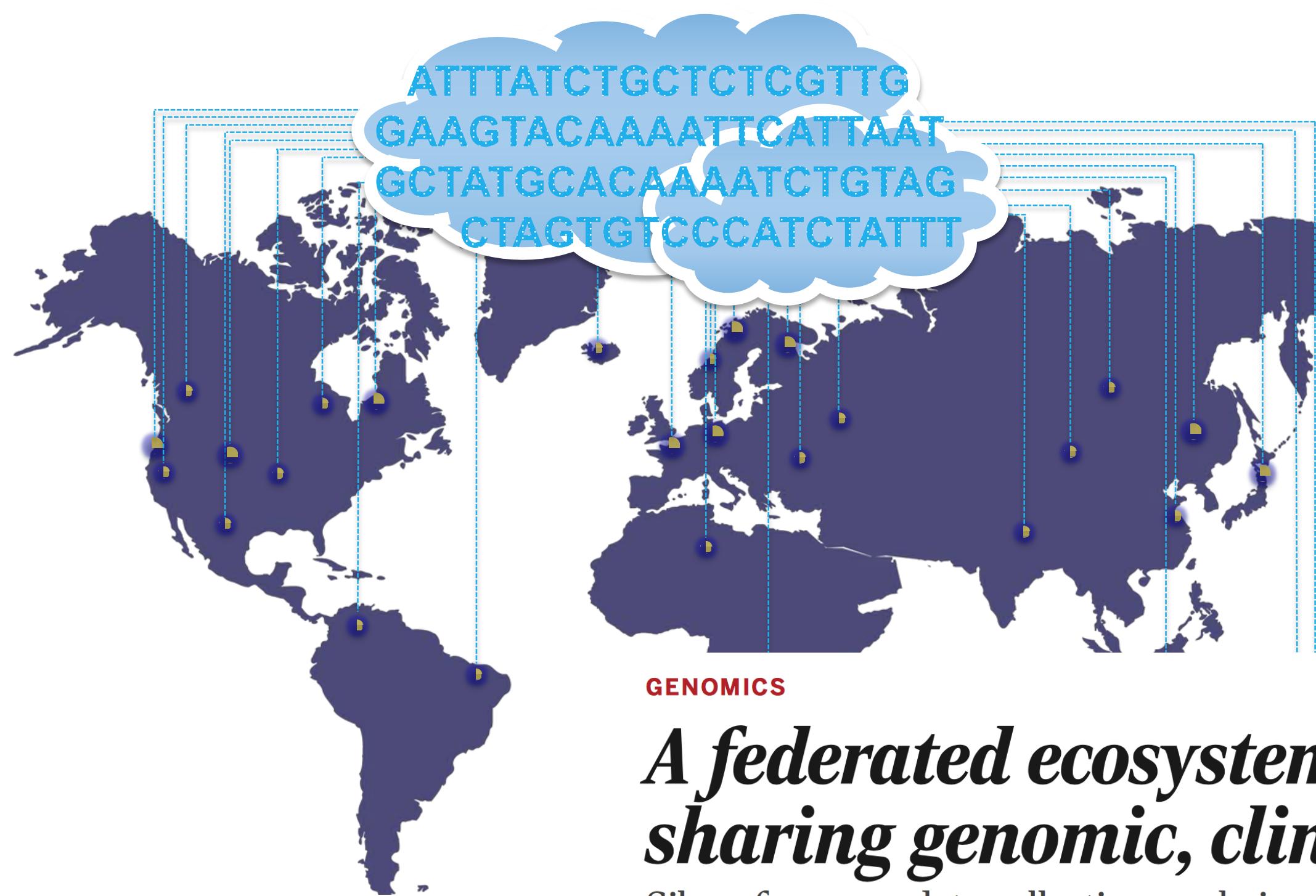
**Figure S1** The fraction or contribution of theoretical ancestors ( $k=9$ ) in reference individuals from 1000 Genomes Project with regard to nine SNP array platforms. The x-axis are individual samples, grouped by their respective population. Groups belonging to the same continent/superpopulation are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



University of  
Zurich<sup>UZH</sup>

Department of Molecular Life Sciences

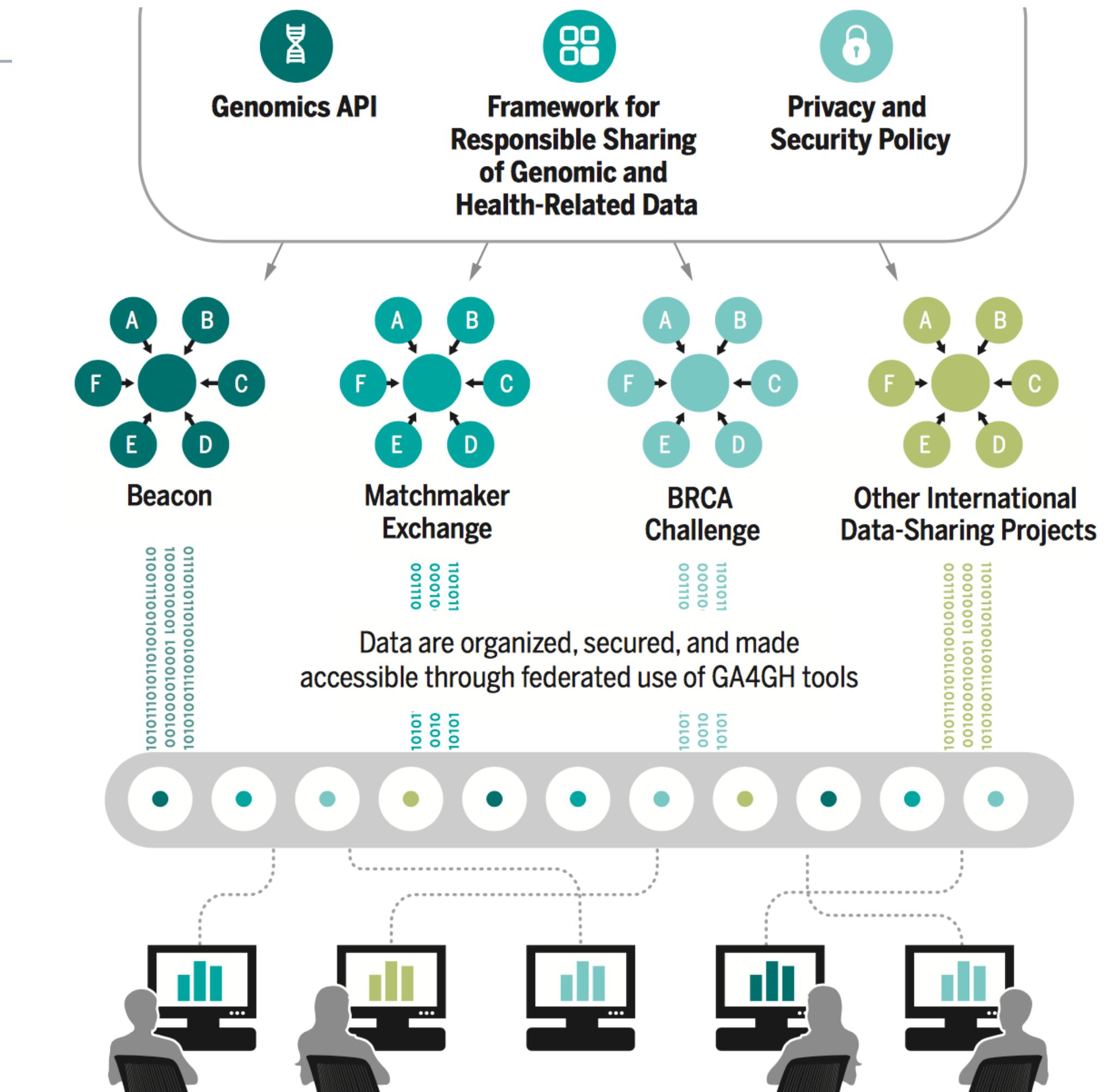
## Genome Data Sharing: The Global Alliance for Genomics and Health (GA4GH)



*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.



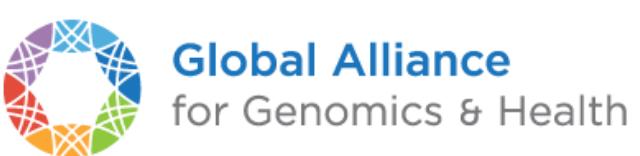


# University of Zurich<sup>UZH</sup>

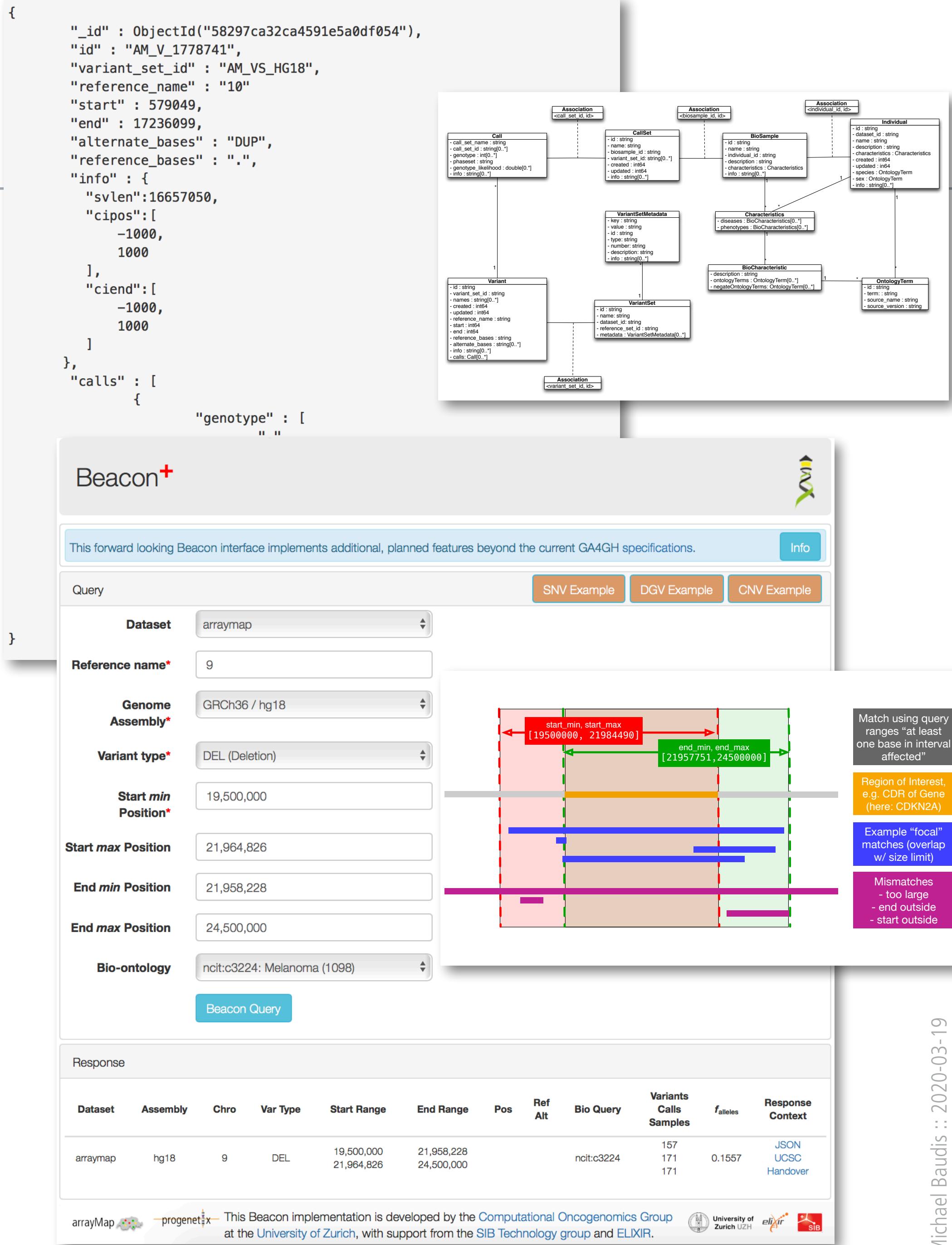
## Department of Molecular Life Sciences

### Our contributions II: Developing Genomic Knowledge Standards & Implementing Driver Projects

- ▶ Progenetix for data driven GA4GH development
  - **Metadata** schema development through implementation of Progenetix and arrayMap data (former ELIXIR project)
  - developing standards for structural genome variants as part of the **ELIXIR h-CNV** and **Beacon** initiatives
  - **SchemaBlocks** objects(w/ EBI, BBOP @ LBNL ...)
- ▶ **Beacon<sup>+</sup>**
  - Beacons - a GA4GH driver project - are **web services** for single-entry genome variant queries over >200 genome resources worldwide.
  - We are developing Beacon<sup>+</sup> protocols for **structural variants**, **metadata** and "handoff" retrieval protocols (co-lead ELIXIR Beacon).
  - Enabling **Implementation** driven development using our and external datasets (arrayMap, TCGA, DIPG ...)



arrayMap



# GA4GH {S}[B] SchemaBlocks

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

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

[schemablocks.org](https://schemablocks.org)



### GA4GH :: SchemaBlocks

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

[About {S}\[B\]](#)  
[News](#)  
[Participants](#)  
[Standards](#)  
[Schemas](#)  
[Examples, Guides & FAQ](#)  
[Meeting minutes](#)  
[Contacts](#)

**Related Sites**

[GA4GH](#)  
[GA4GH::Discovery](#)  
[Beacon Project](#)  
[Phenopackets](#)  
[GA4GH::CLP](#)  
[GA4GH::GKS](#)  
[Beacon+](#)

**Github Projects**

[SchemaBlocks](#)  
[ELIXIR Beacon](#)

**Tags**

Beacon CP Discovery FAQ GA4GH  
GKS MME admins code contacts  
contributors core dates developers  
documentation howto identifiers  
implemented issues leads news  
phenopackets playground press  
proposed sb-phenopackets tools  
website



## GA4GH SchemaBlocks Home

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

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

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

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

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

## SchemaBlocks Repositories

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

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

## SchemaBlocks “Status” Levels

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

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

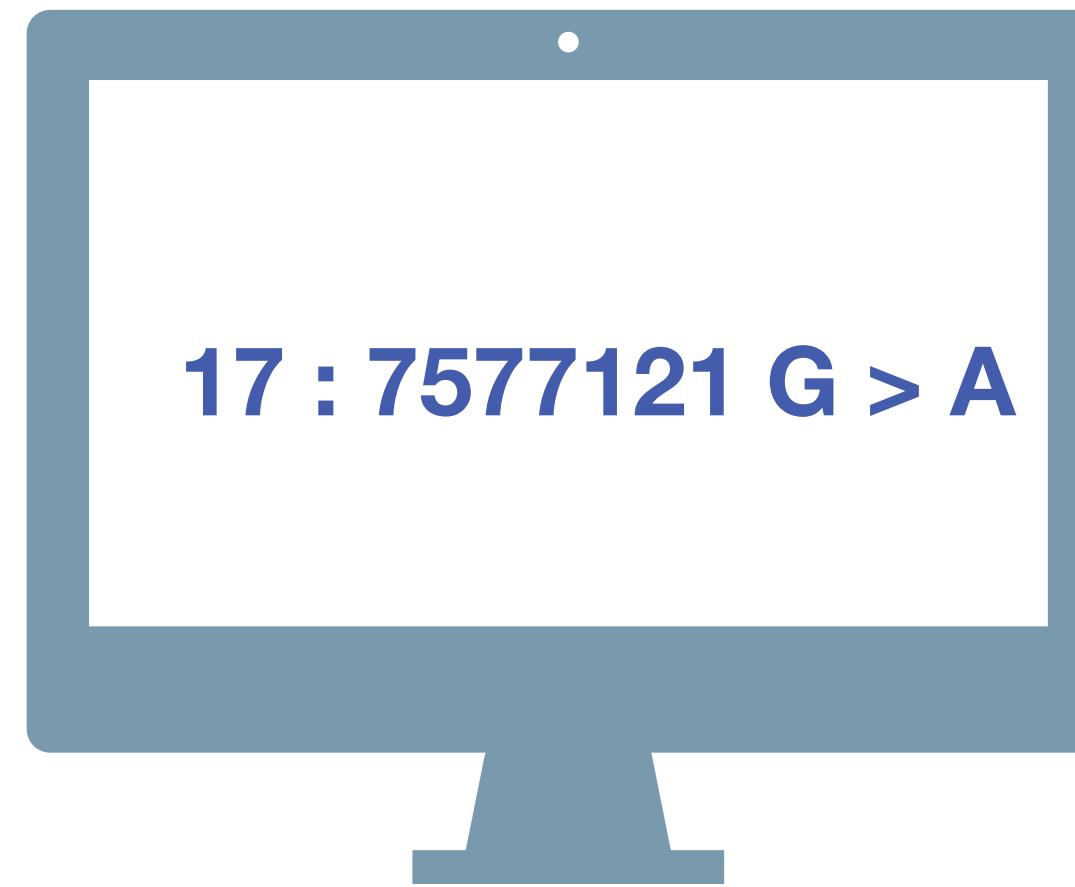
## SchemaBlocks {S}[B] Mission Statement

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

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

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

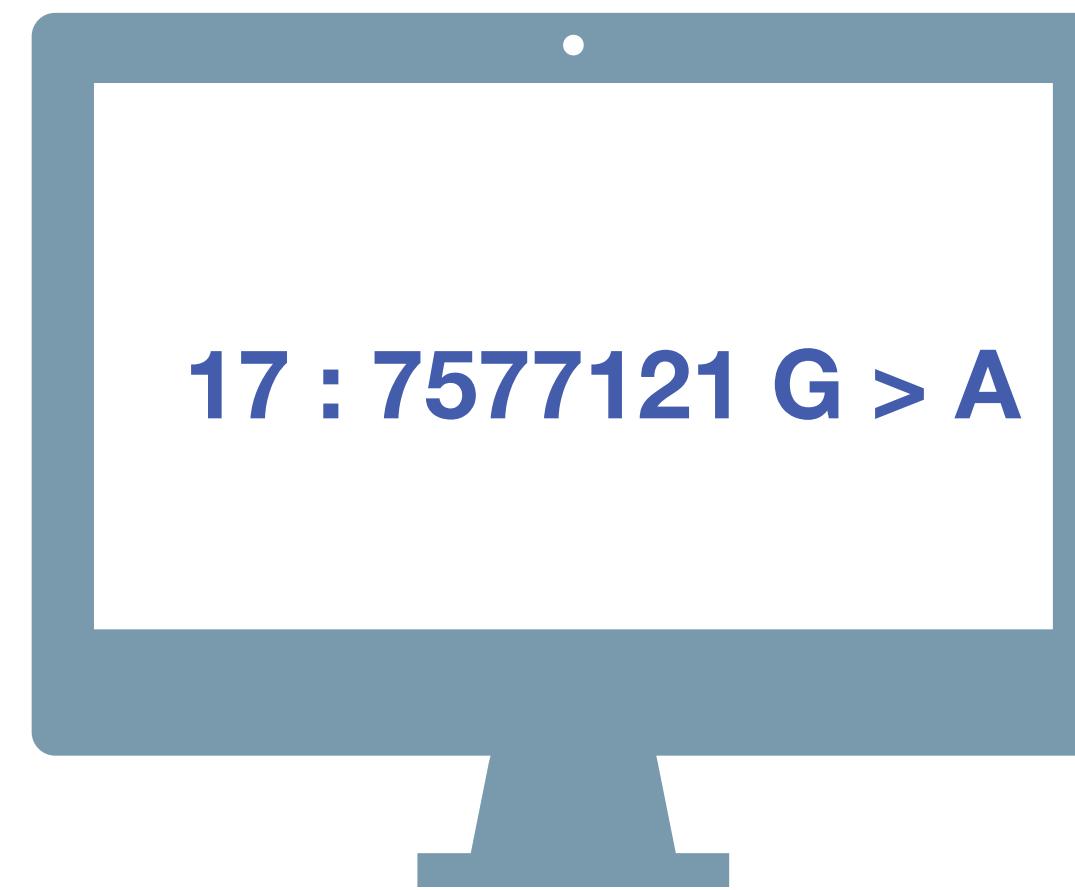




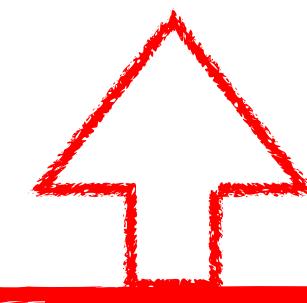
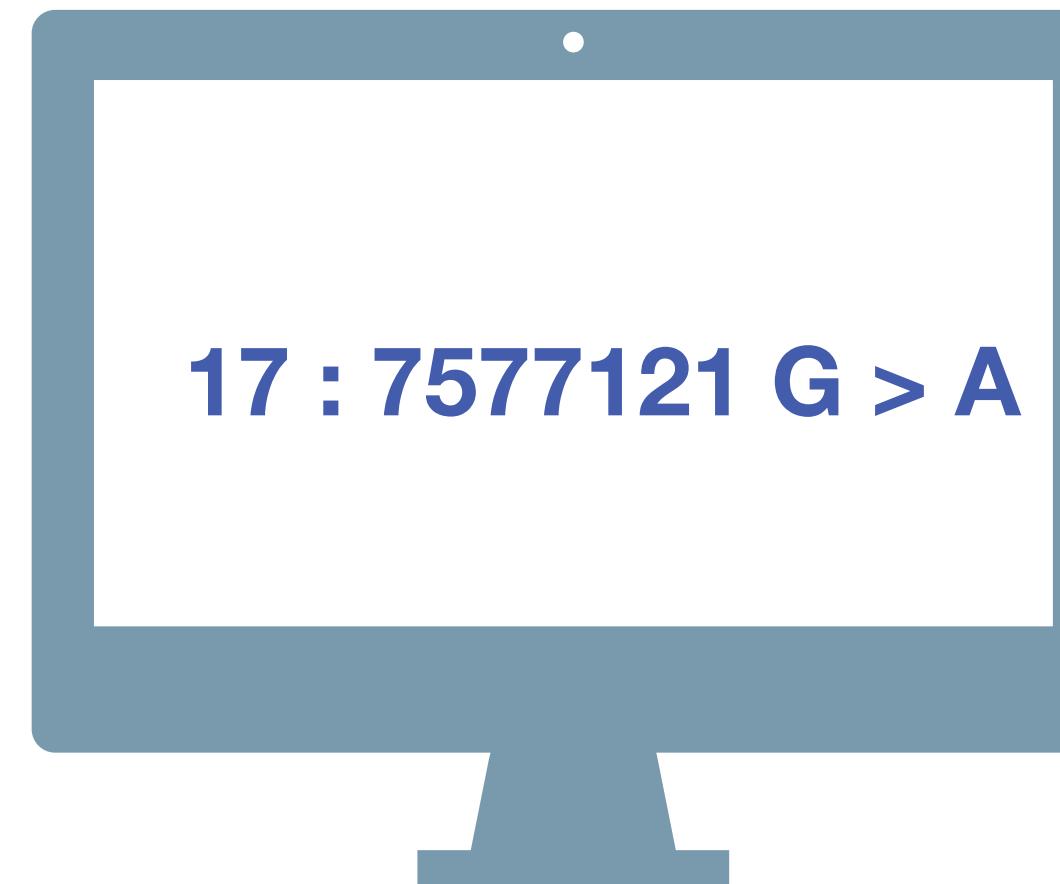
# Beacon

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

**YES | NO | \0**



A Beacon network federates  
*genome variant queries*  
across databases that  
support the ***Beacon API***

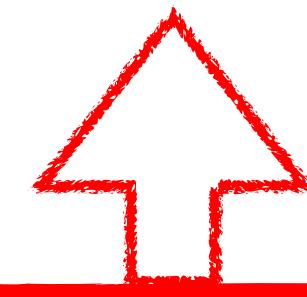
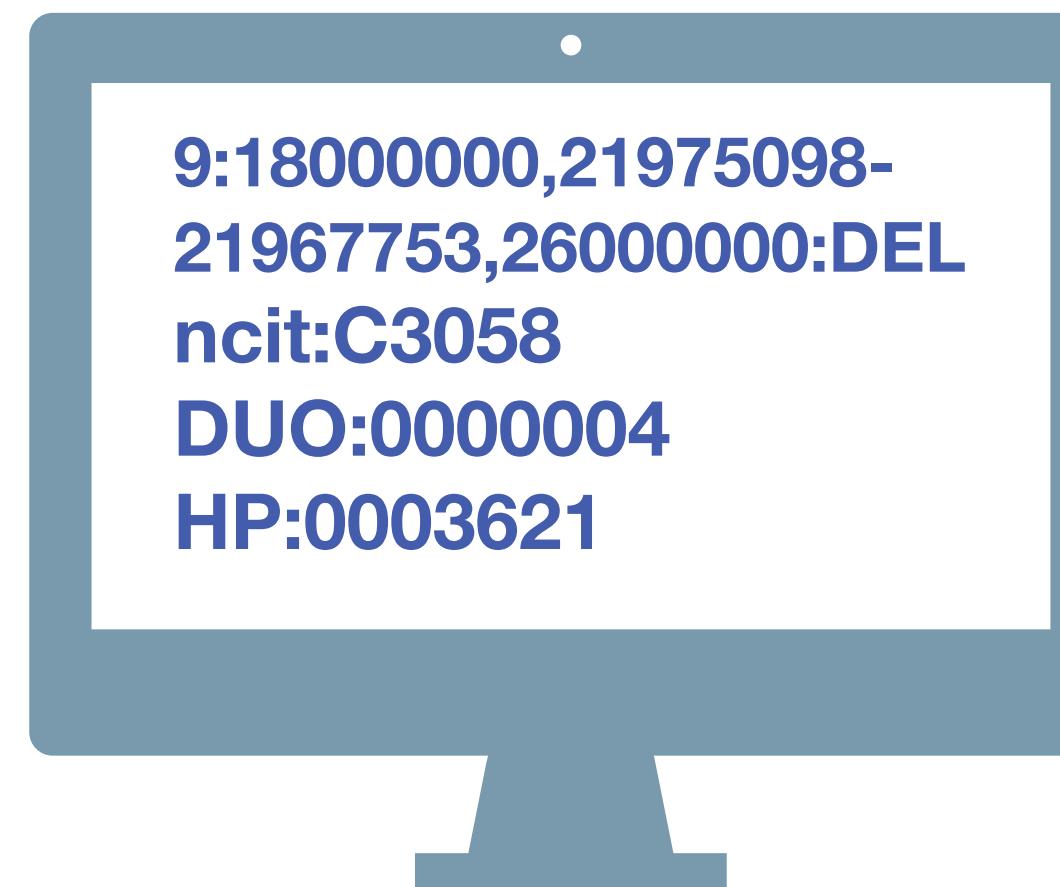


Have you seen this variant?  
It came up in my patient  
and we don't know if this is  
a common SNP or worth  
following up.

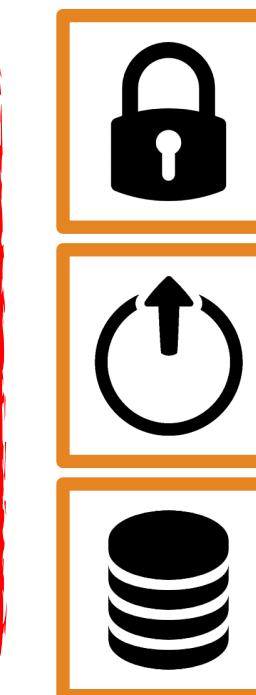


A Beacon network federates  
genome variant queries  
across databases that  
support the **Beacon API**

Here: The variant has  
been found in **few**  
resources, and those  
are from **disease**  
specific **collections**.

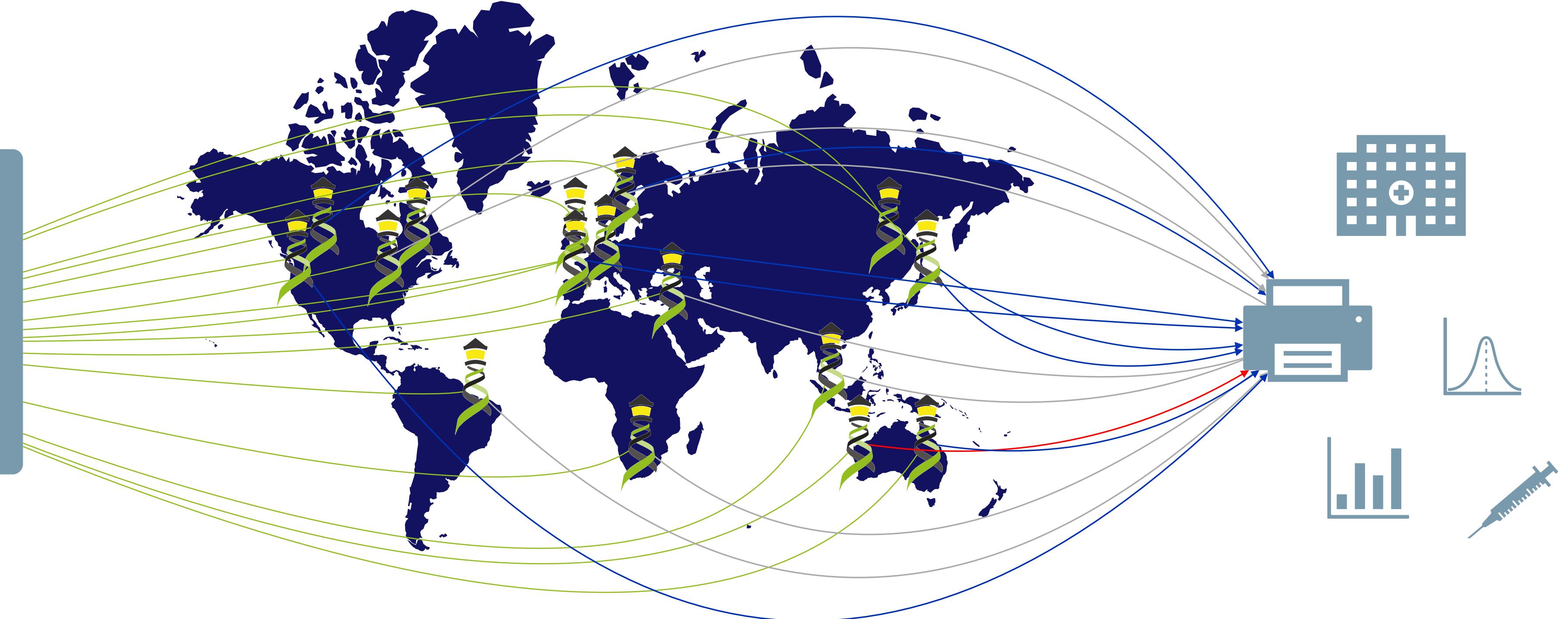


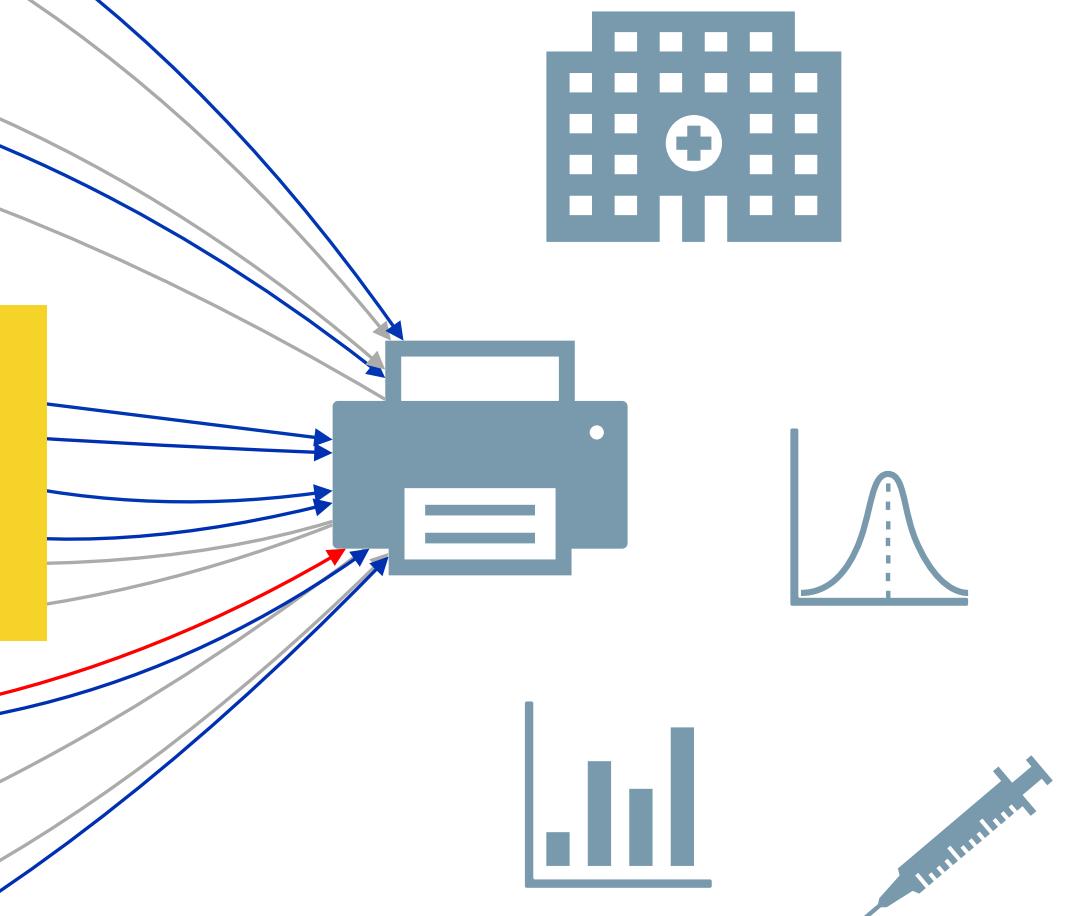
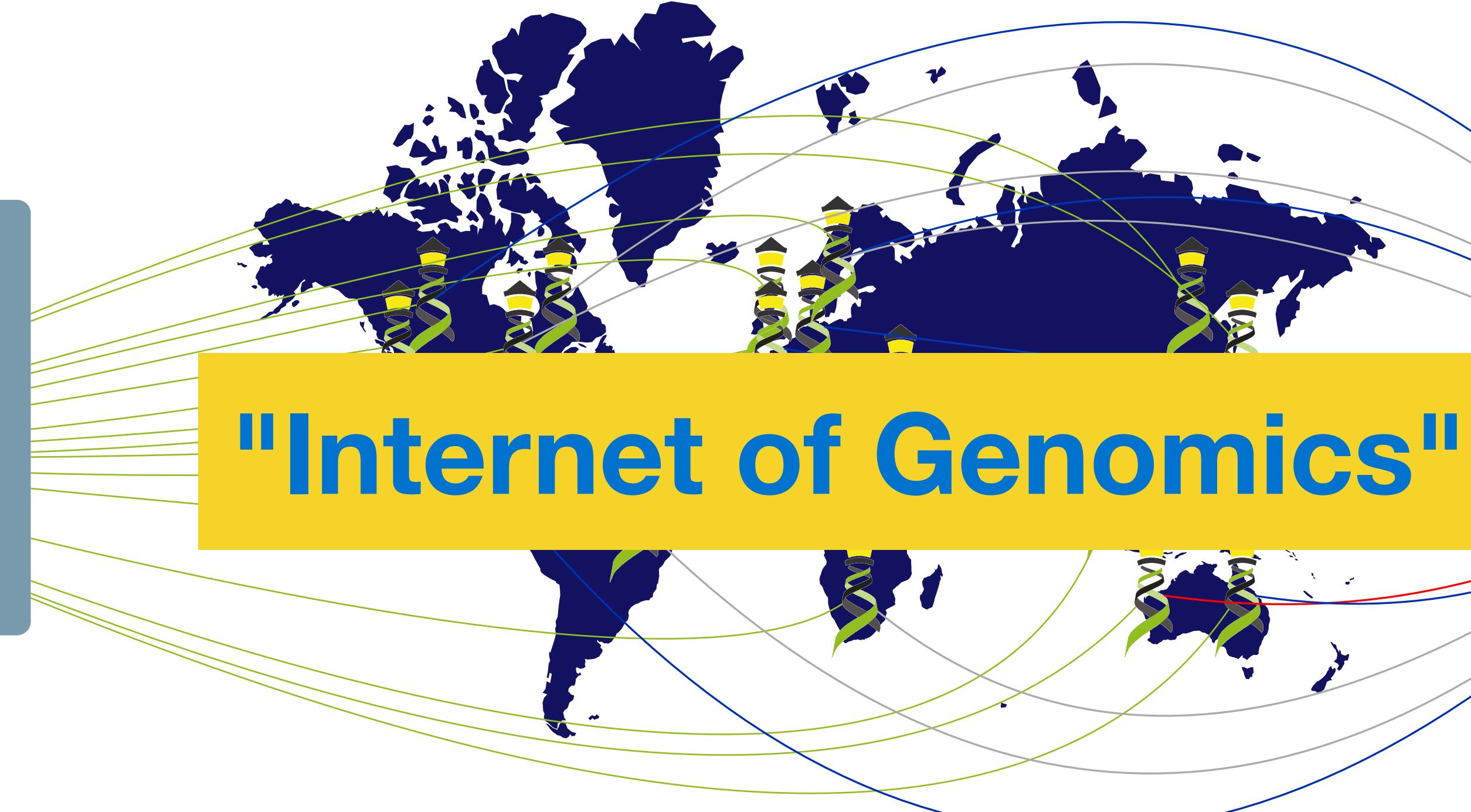
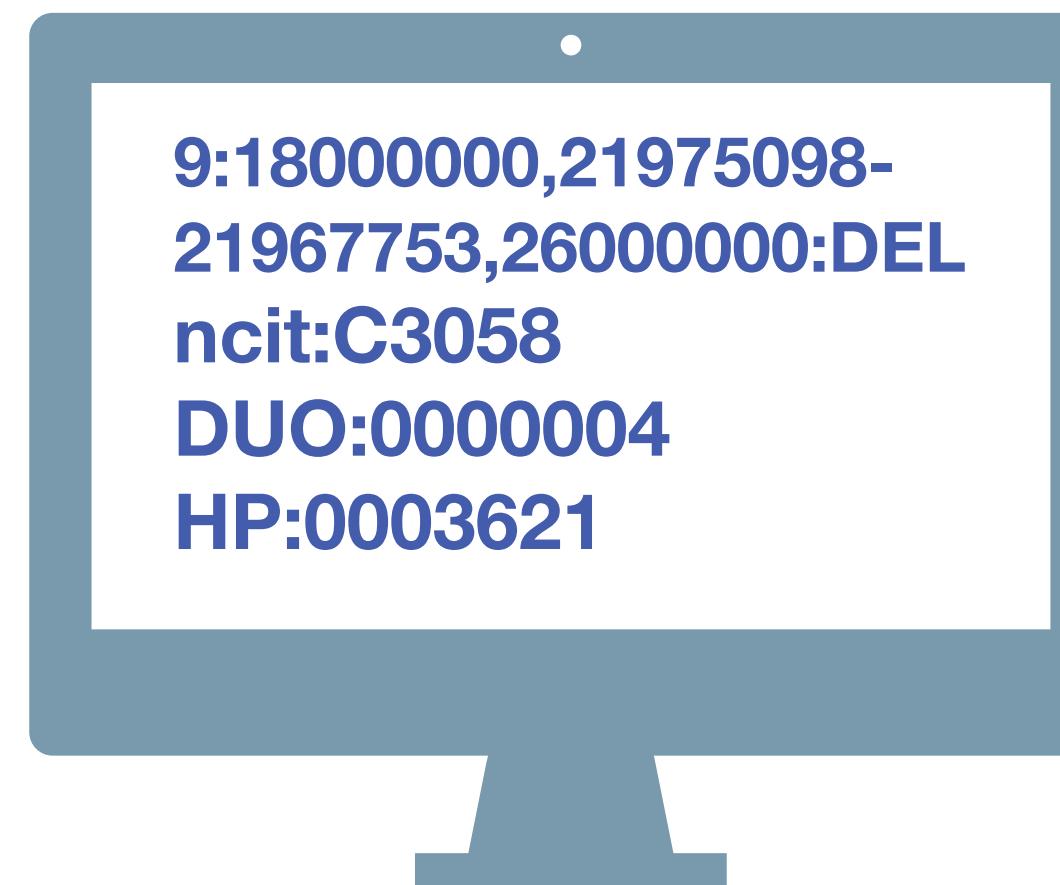
Have you seen deletions in this region on chromosome 9 in Glioblastomas from a juvenile patient, in a dataset with unrestricted access?



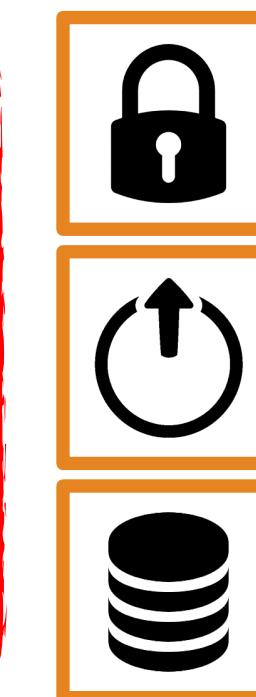
## Beacon v2 API

The Beacon API v2 proposal opens the way for the design of a simple but powerful "genomics API".





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".

## 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

## @WORLD

MATTHIAS ALTMAYER  
THOMAS EGGERMANN  
ROSA NOGUERA  
REINER SIEBERT  
CAIUS SOLOVAN



University of  
Zurich<sup>UZH</sup>



Global Alliance  
for Genomics & Health

## BBOP GROUP @ LNBL

CHRIS MUNGALL  
JUSTIN REESE  
DEEPAK UNNI

## GA4GH

MELANIE COURTOT  
MELISSA HAENDEL  
HELEN PARKINSON

## ELIXIR & CRG

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

## H-CNV

CHRISTOPHE BÉROUD  
DAVID SALGADO



University of  
Zurich<sup>UZH</sup>

Department of Molecular Life Sciences



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

[arraymap.org](http://arraymap.org)  
[progenetix.org](http://progenetix.org)  
[info.baudisgroup.org](mailto: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)



Global Alliance  
for Genomics & Health

