

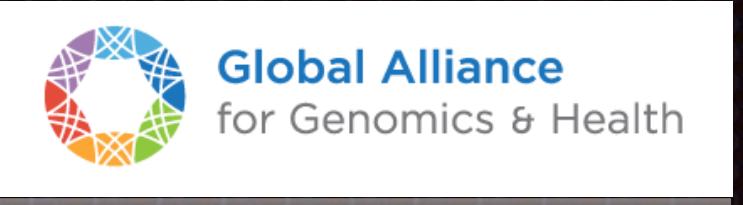
Genome data access

Application to personalized health and cancer research

Michael Baudis
Computational Oncogenomics



University of
Zurich^{UZH}





Janet Rowley (1972/73)

Chromosomal translocations in cancer

- Recurrent chromosomal translocations in leukemias and lymphomas
- "Philadelphia chromosome" in CML (Nowell & Hungerford, 1960) represents a reciprocal translocation between chromosomes 9 and 22
- 1972: t(8;21) ALL manuscript rejected by NEJM
- 1973: t(9;22) manuscript rejected by *Nature* "with some reasonable comments and some truly wrong"
- Clinical implications: **Tyrosine Kinase inhibitors** as standard first-line therapy in CML
 - first trials in 1998 (STI-571; Imatinib/Gleevec)
 - cf. Druker BJ, Lydon NB (2000). Lessons learned from the development of an Abl tyrosine kinase inhibitor... *J Clin Invest* 2000;105:3-7

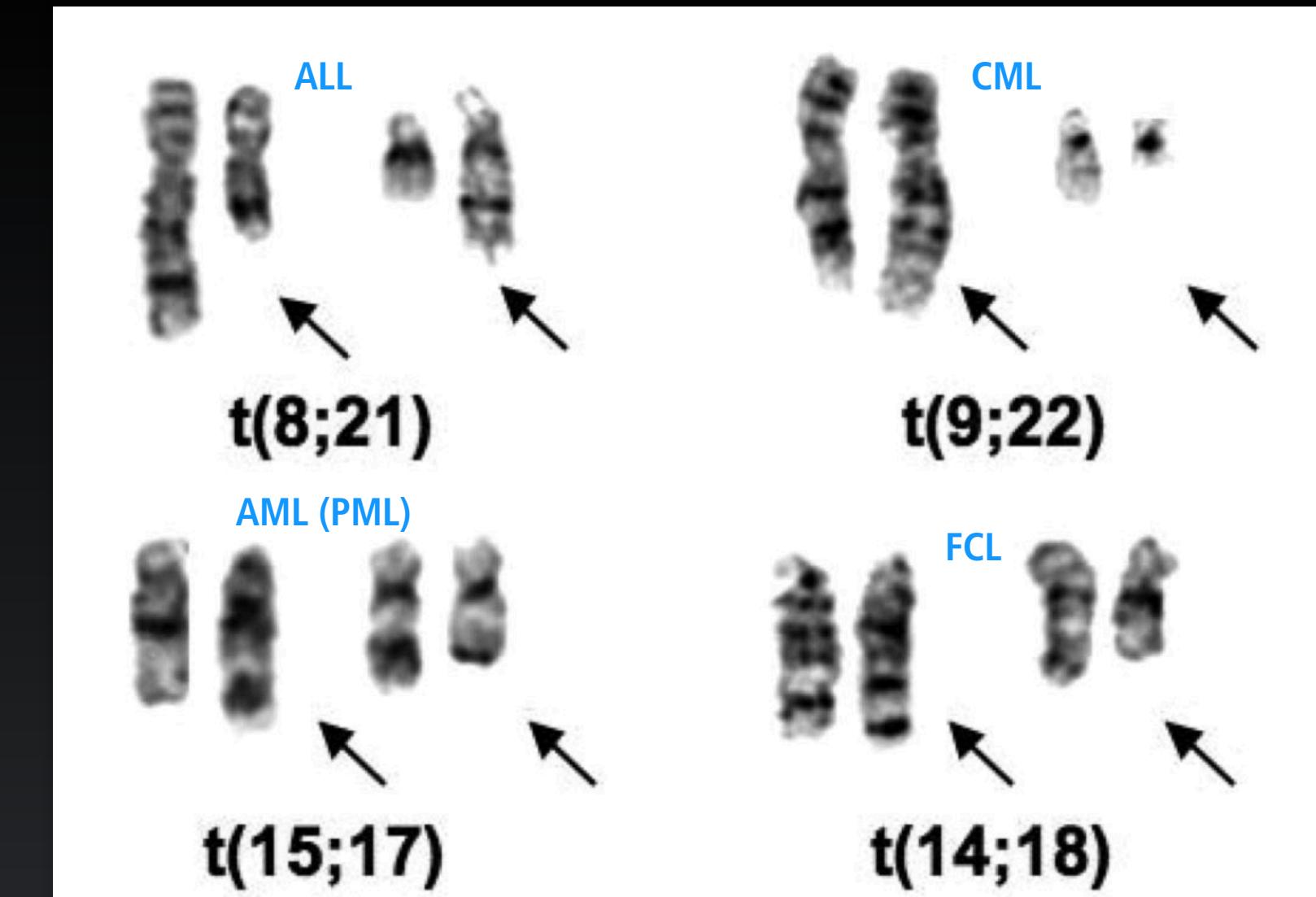
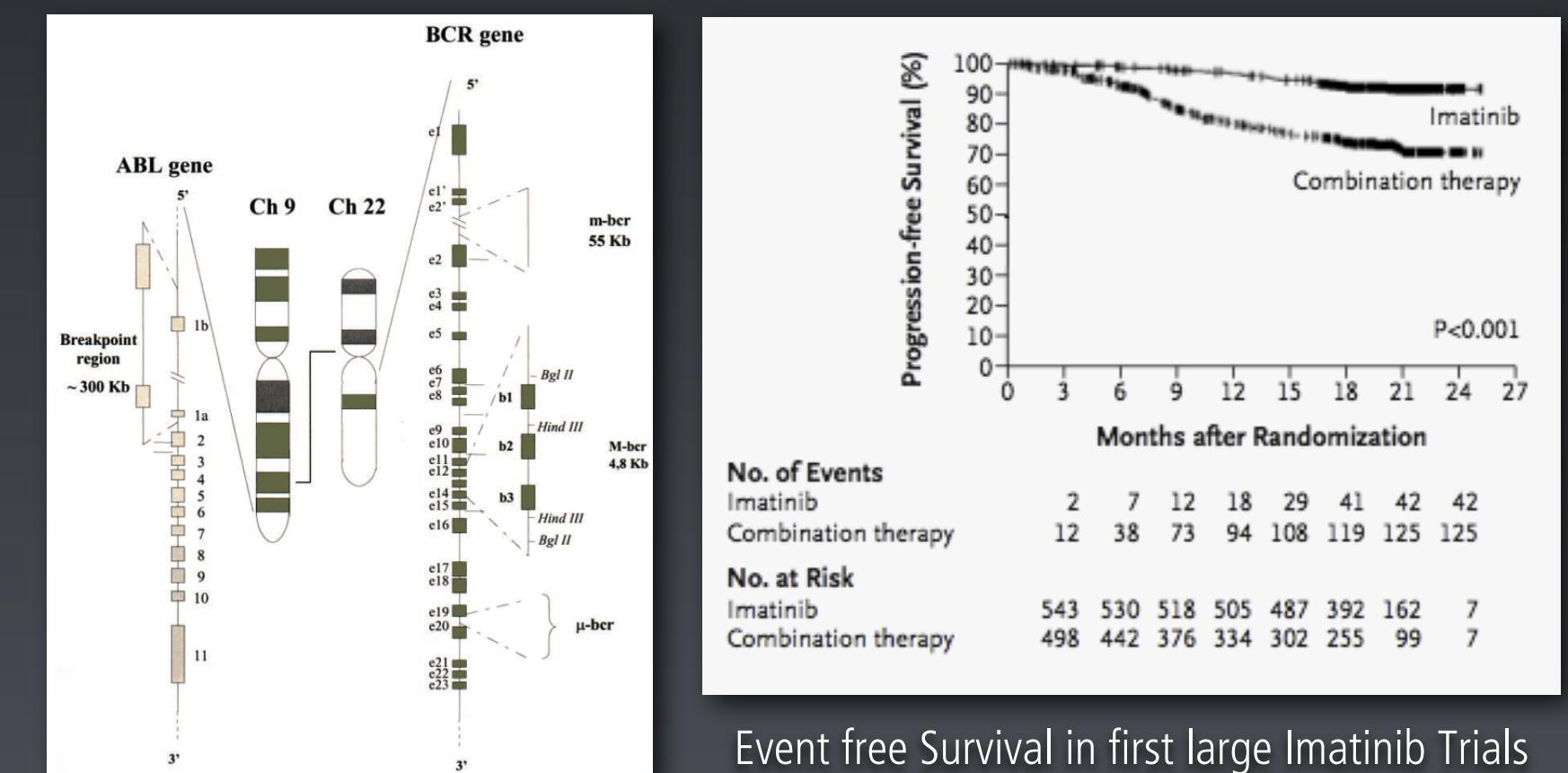


Figure 1. Partial karyotypes of common translocations discovered by Rowley.
The translocations appear in the order in which they were discovered.

Janet D Rowley. Chromosomal translocations: revisited yet again *Blood* (2008), 112(6)



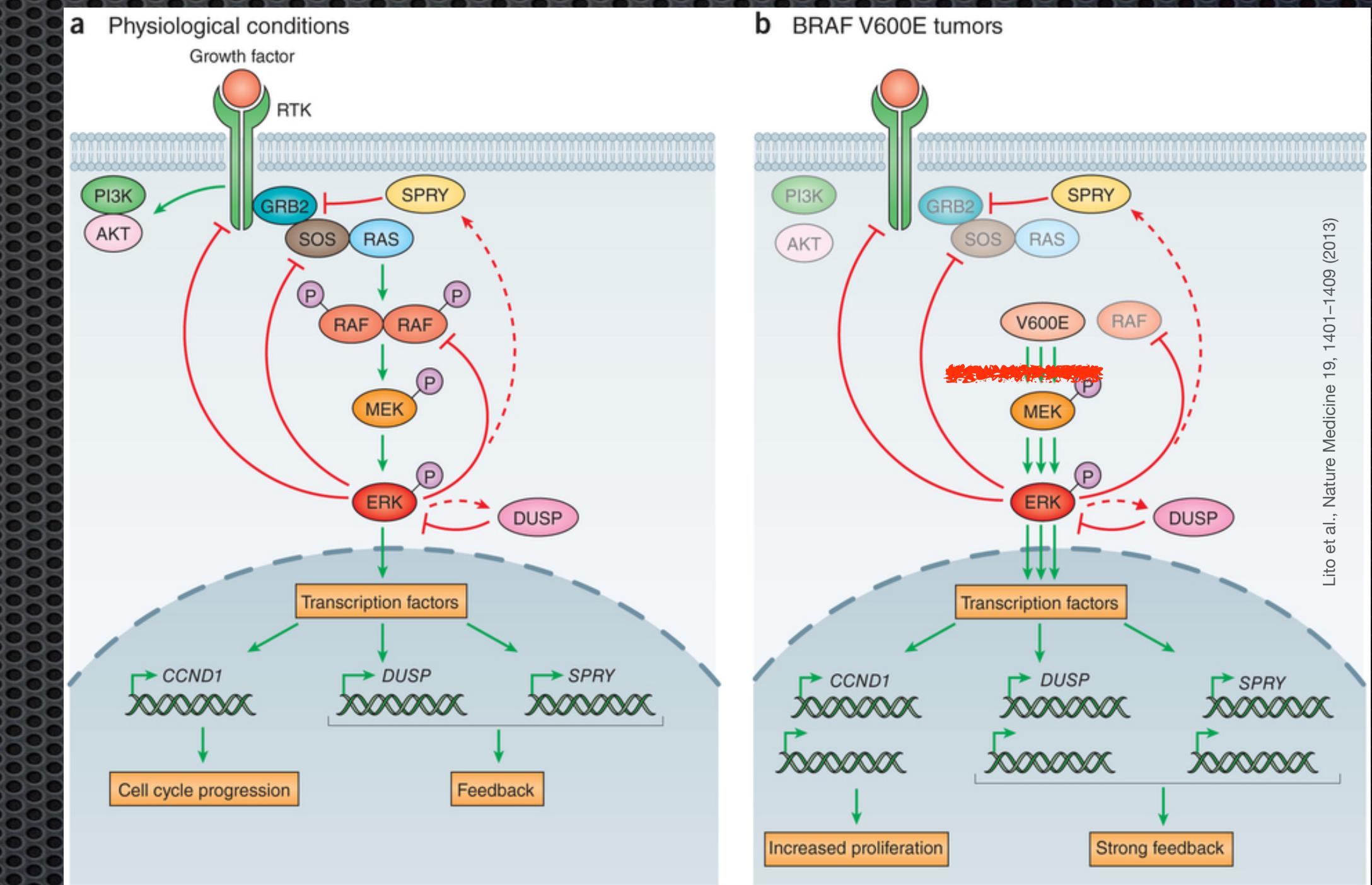
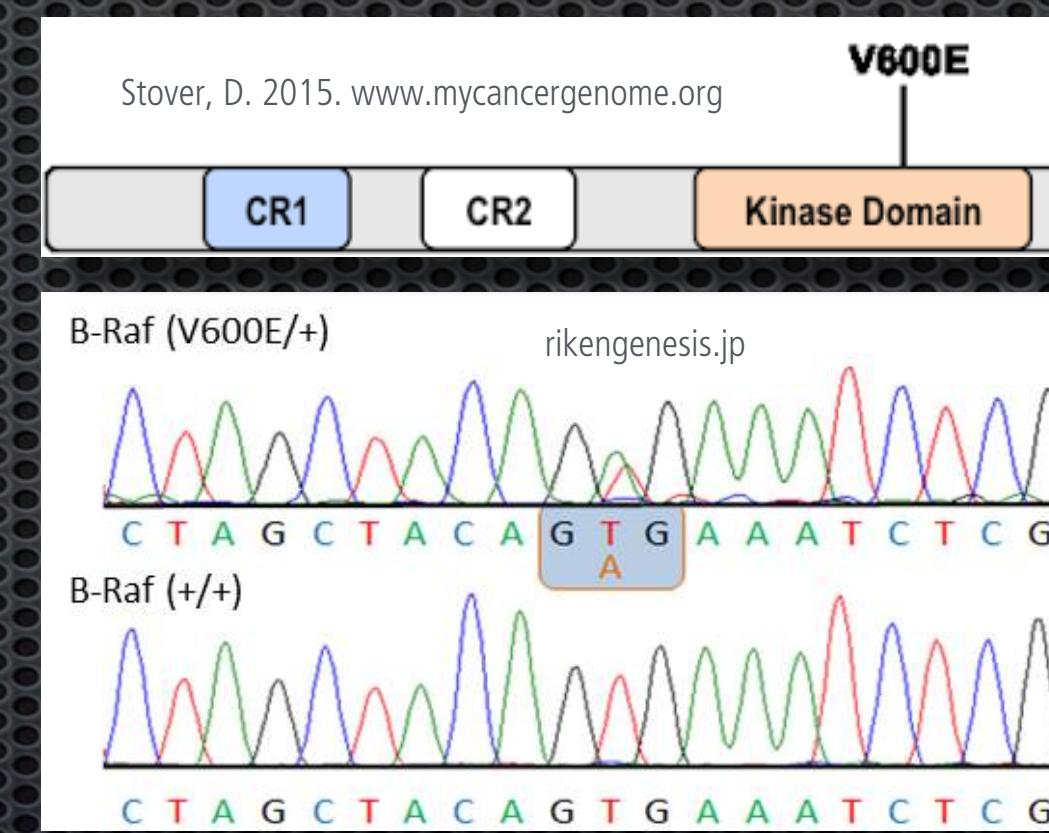
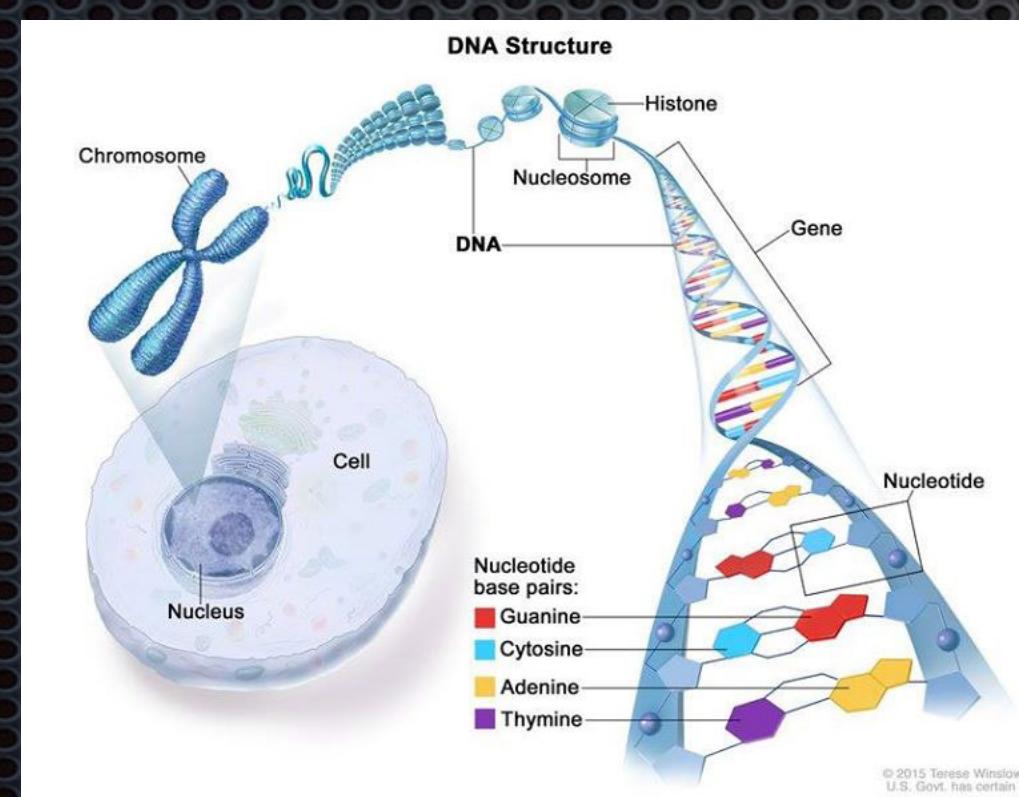
Event free Survival in first large Imatinib Trials

Pane et al. BCR/ABL genes
Oncogene (2002), 21 (56)

O'Brien et al. Imatinib compared with interferon and low-dose cytarabine... *NEJM* (2003) vol. 348 (11)

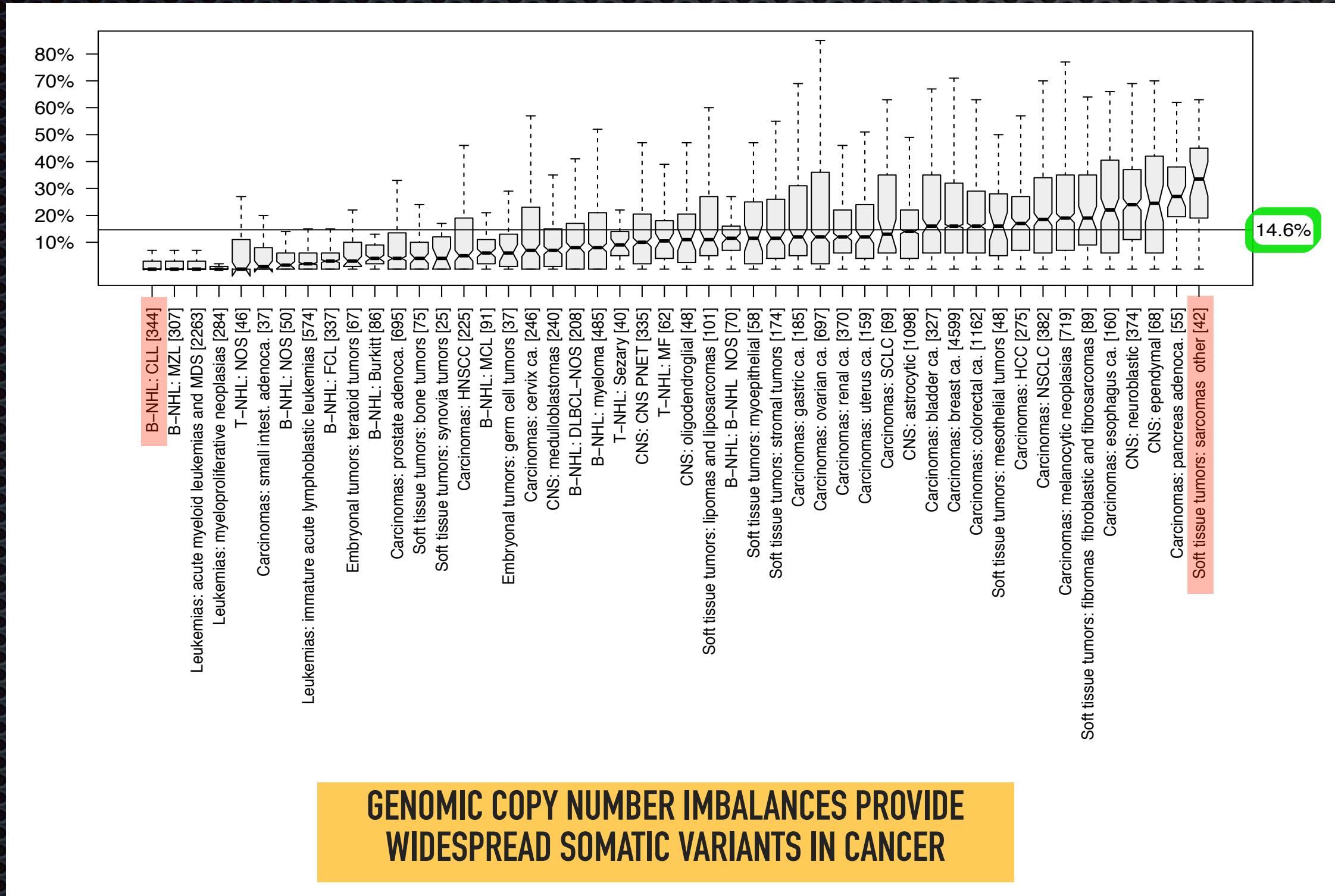
BRAF V600E (c.1799T>A) Mutation Oncogene Activation by Single Nucleotide Alteration

- a single nucleotide exchange Thymidine > Adenine leads to continuous RAF based activation of the MEK-ERK pathway
- BRAF V600E is a frequent mutation in >50% of malignant melanomas, but also CRC, lung ADC ...
- pharmacologic block of B-Raf (e.g. through **Vemurafenib**)

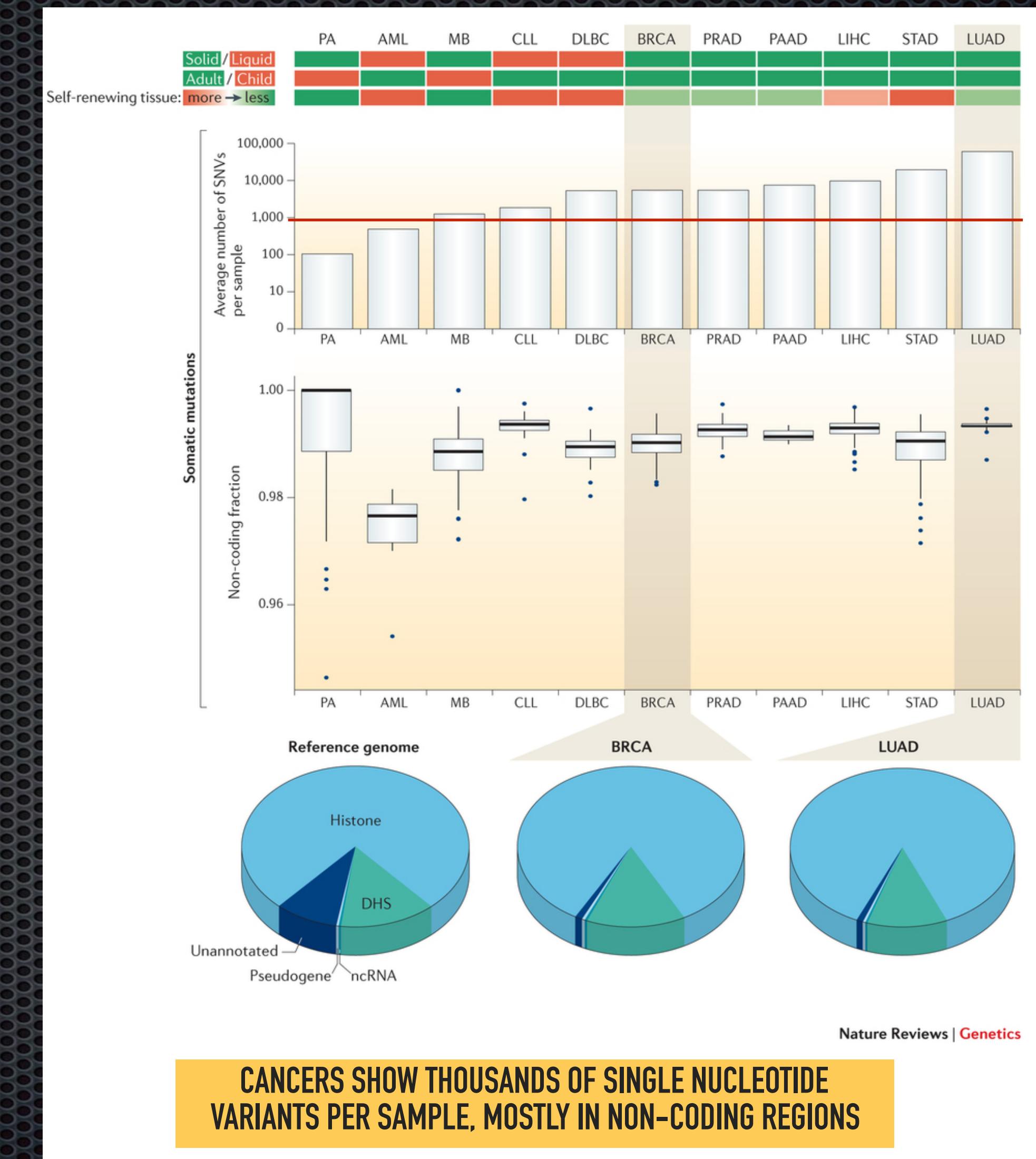


The BRAF V600E mutation leads to continuous phosphorylation of MEK, without the need for receptor based activation of the upstream pathway and loss of inhibitory feedback control.

Quantifying Somatic Mutations In Cancer

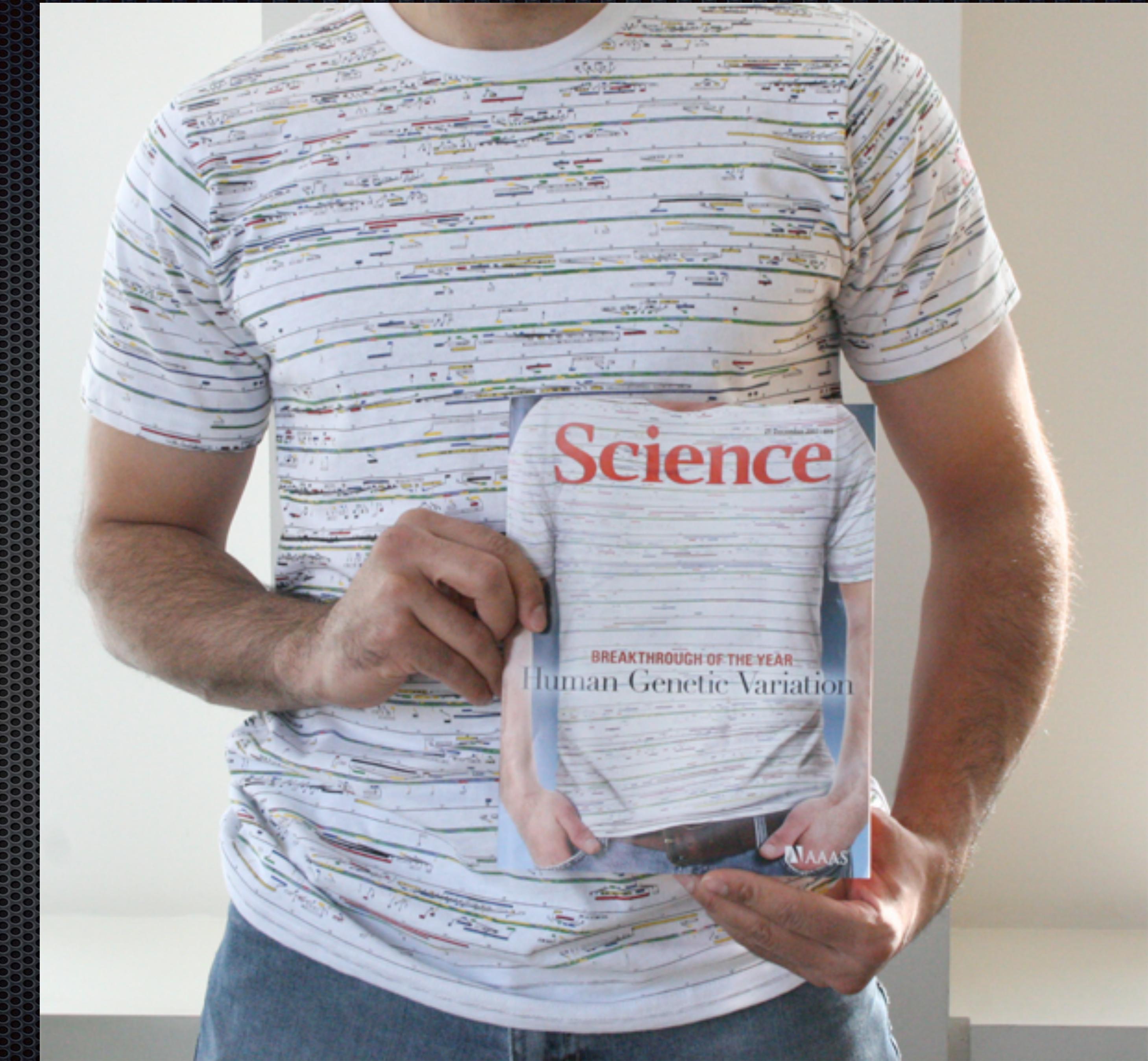


On average ~15% of a cancer genome are in an imbalanced state (more/less than 2 alleles);
Original data based on >30'000 cancer genomes from arraymap.org



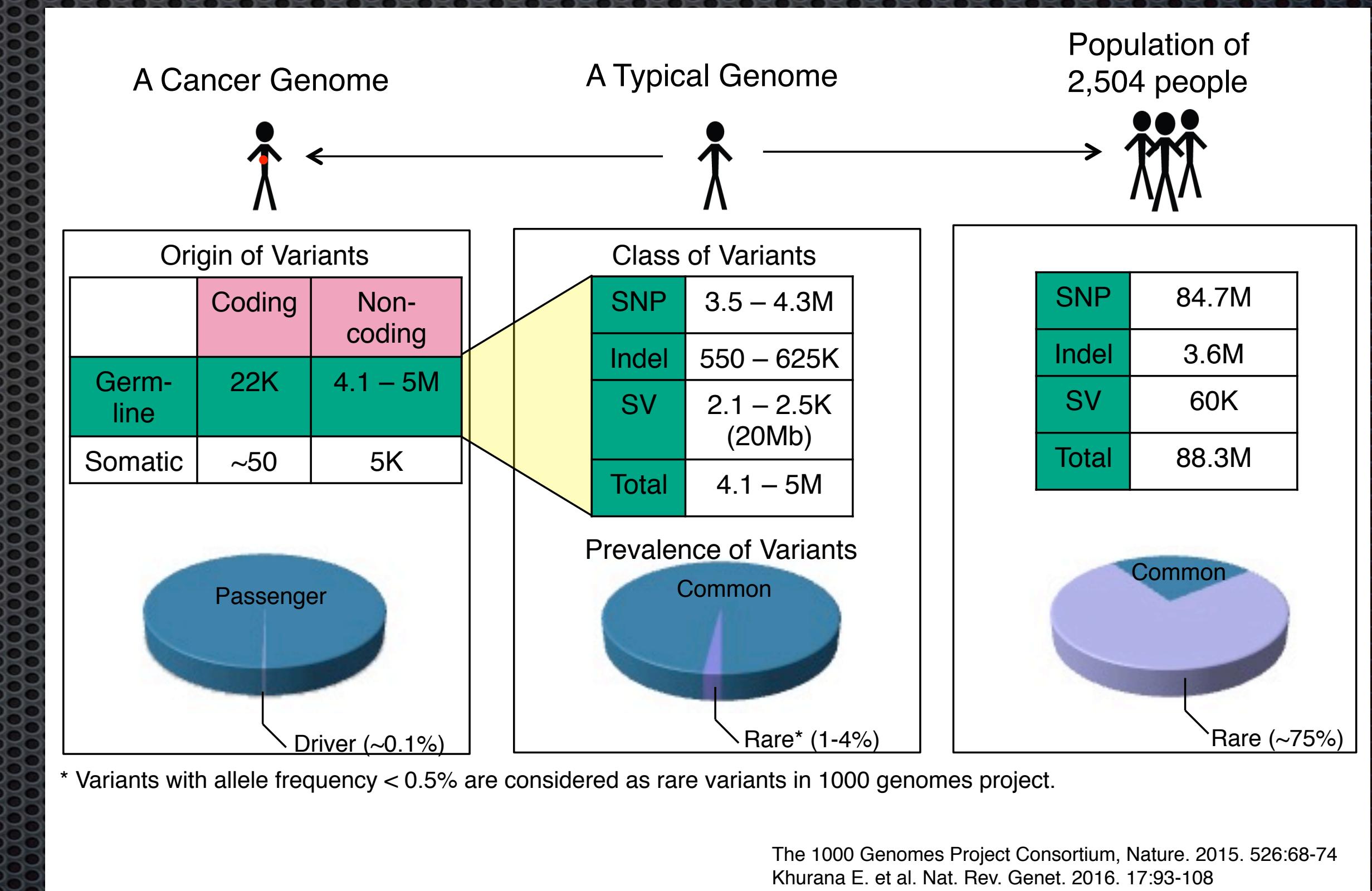
Pan-Cancer Analysis of Whole Genomes (PCAWG) data show widespread mutations in non-coding regions of cancer genomes (Khurana et al., Nat. Rev. Genet. (2016))

The trouble with human genetic variation



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

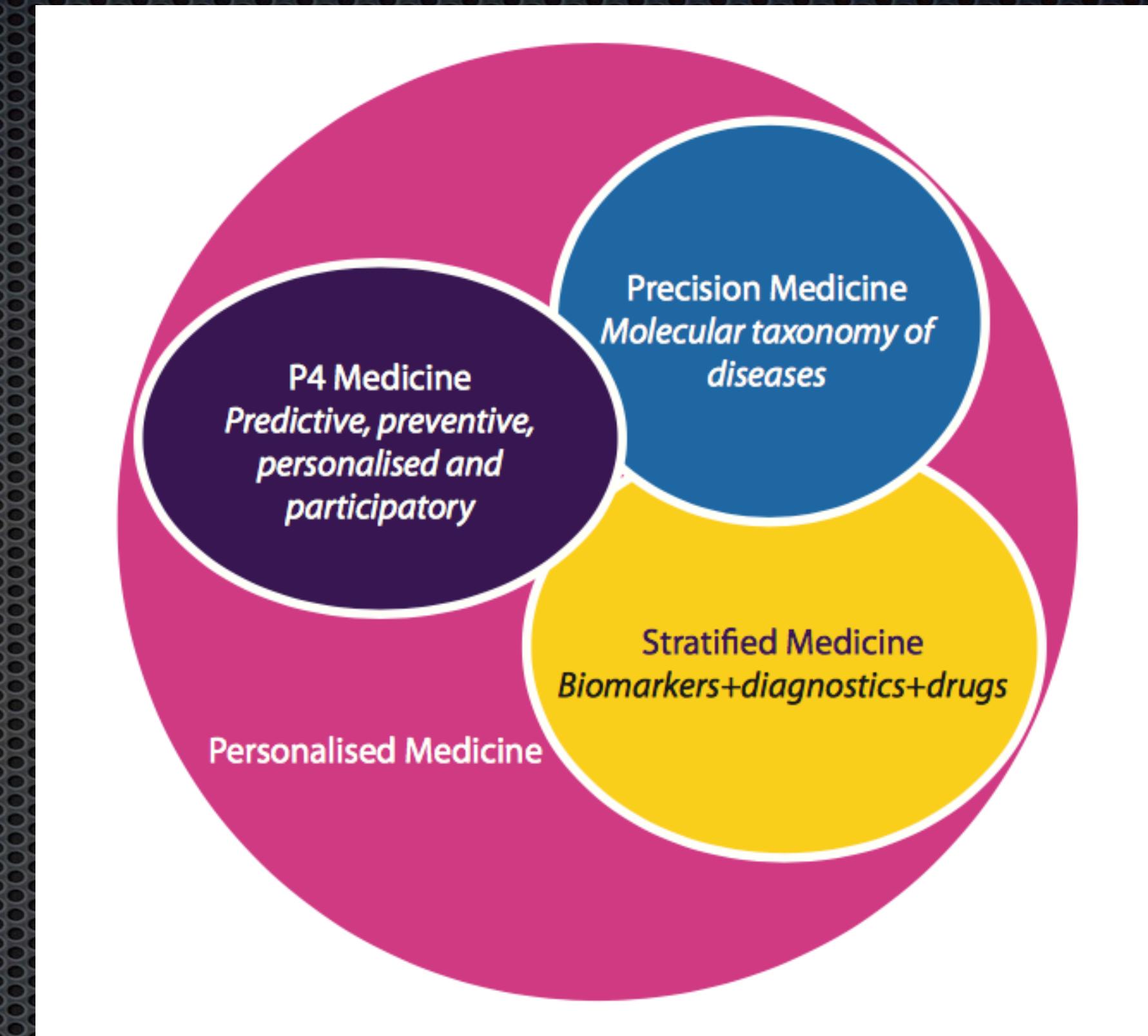
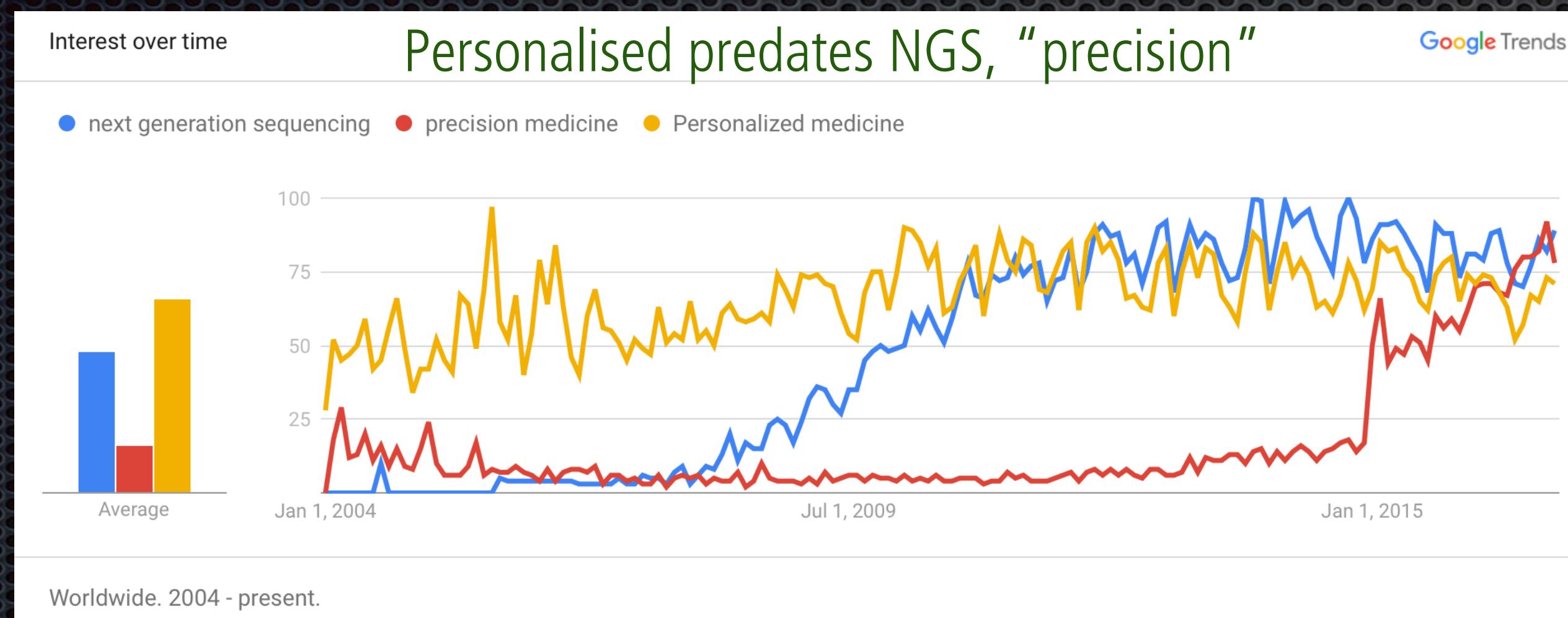


Graphic adapted from Mark Gerstein (GersteinLab.org; @markgerstein)

Genomic Background + Disease Parameters

Personalised Medicine **Precision Medicine**

...

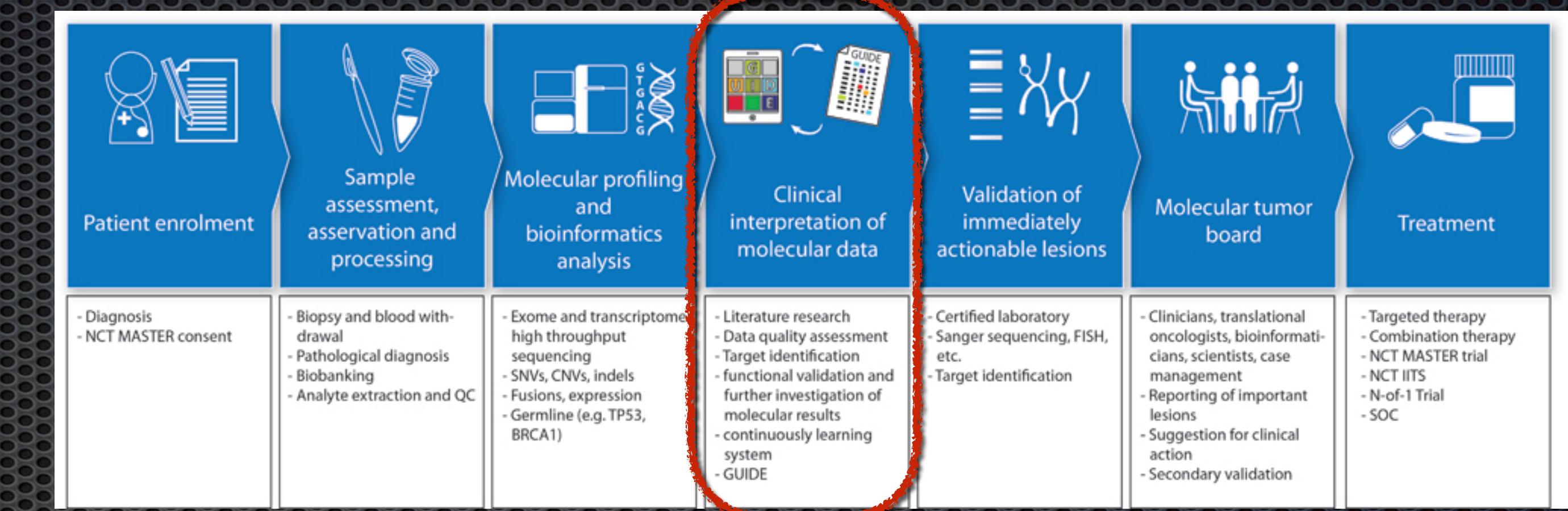


Source: PHG Foundation

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 genome information 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

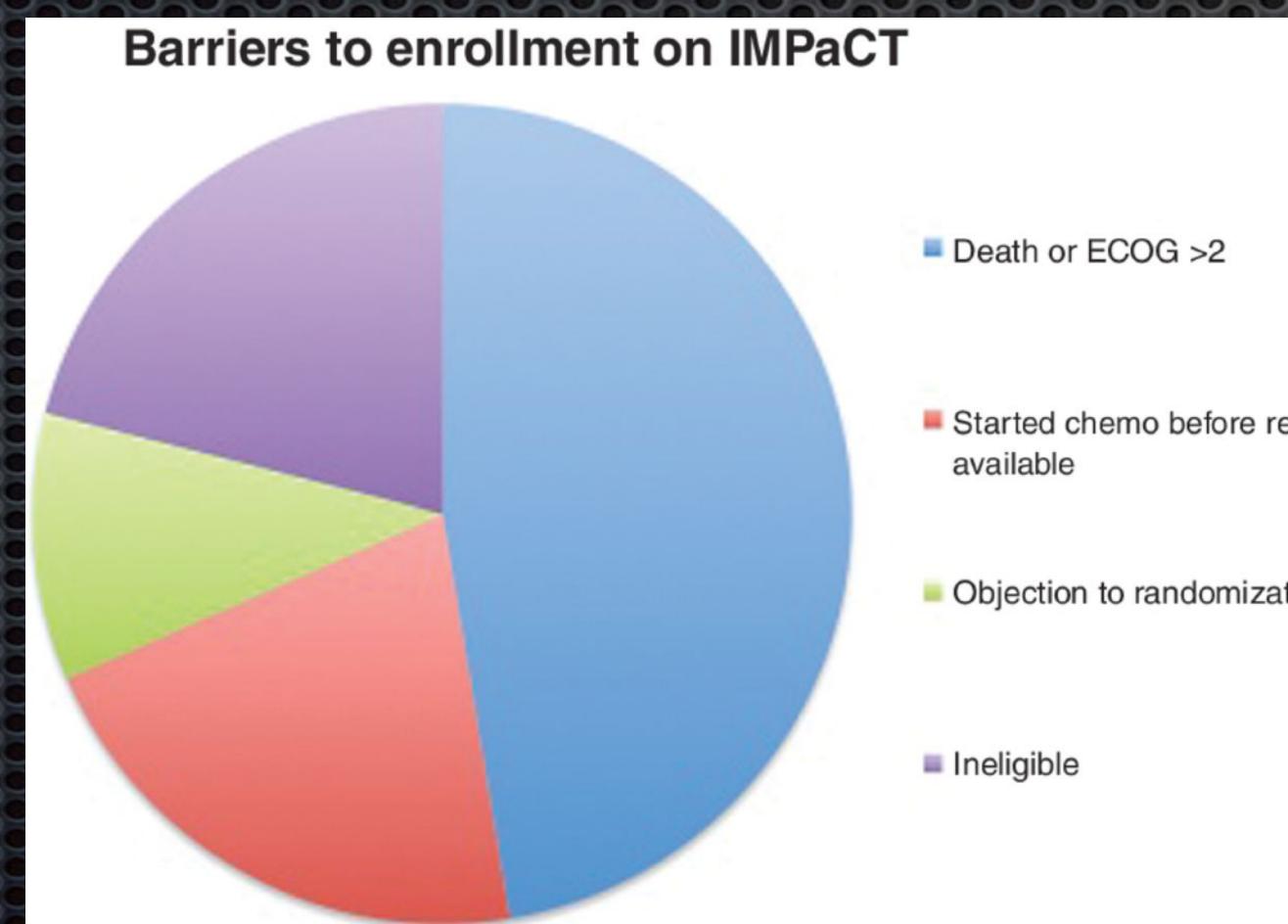


Workflow of a cancer treatment protocol based on "personalized" assessment of actionable genomic lesions (source: NCT Heidelberg).

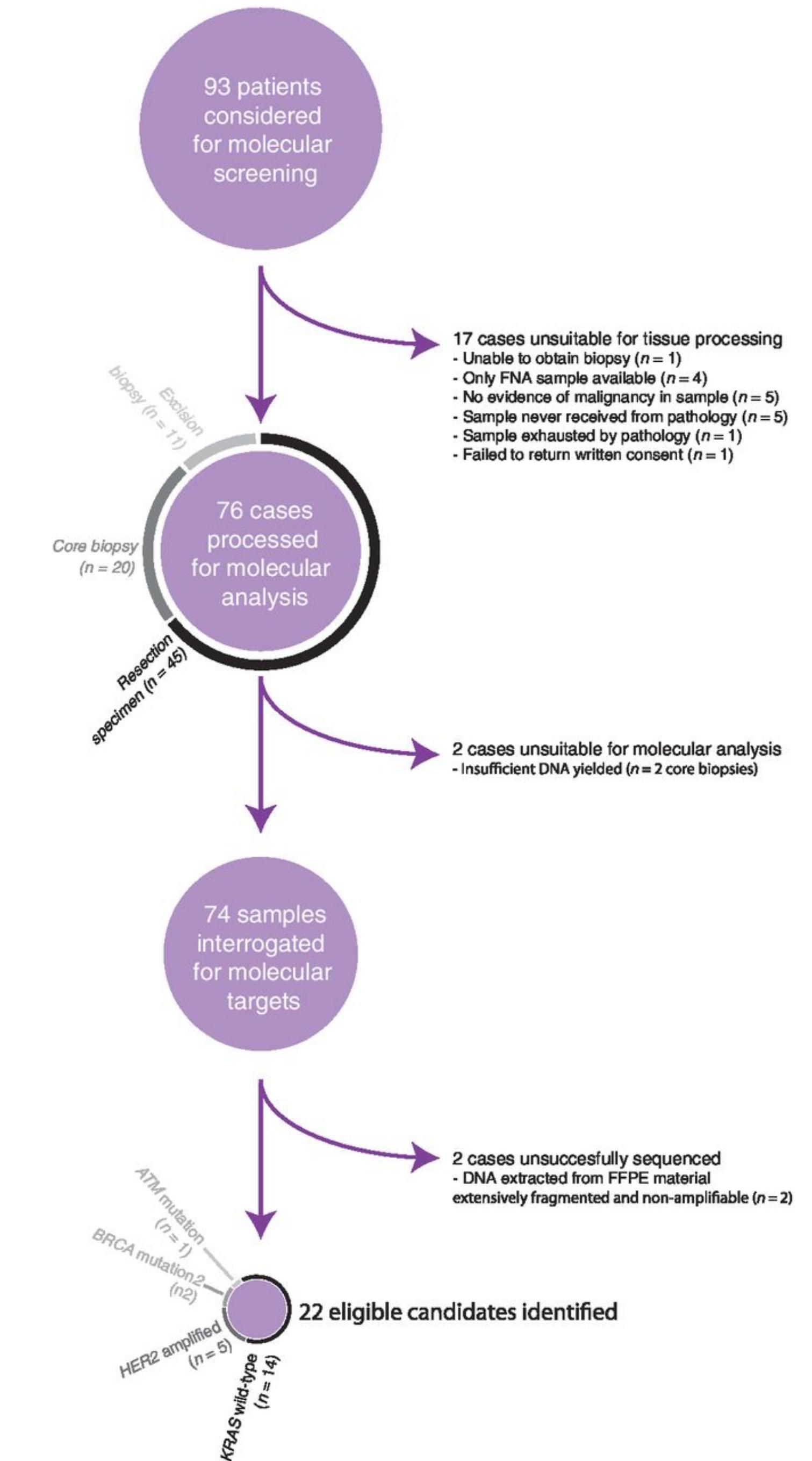
- currently mostly use of **targeted / panel sequencing** for identification of tens - hundreds of most common "actionable" mutations
- knowledge resources and literature search for interpretation of non-standard variants

Personalised Cancer Therapy Needs More Targets

- with the current knowledge, targeted molecular analysis will not lead to the identification of actionable interventions in a majority of cancer cases



Precision Medicine for Advanced Pancreas Cancer: The Individualized Molecular Pancreatic Cancer Therapy (IMPaCT) Trial.
Chantrill *et al.*, Clin Cancer Res. 2015 May 1;21(9):2029-37



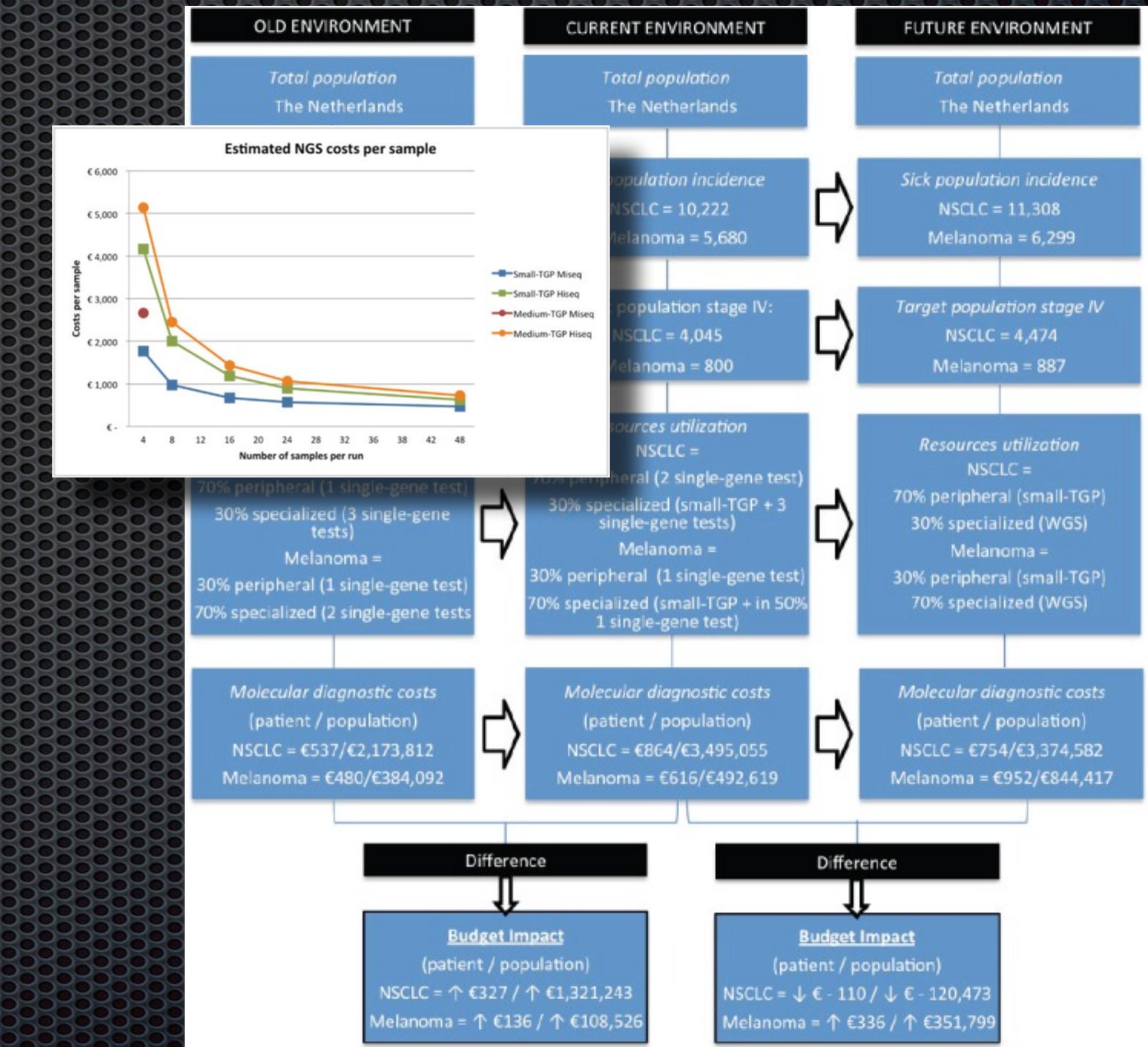
Next-generation sequencing strategies in malignancies

	Targeted panels	Whole exome	Whole genome
Pro	<ul style="list-style-type: none"> ▶ High depth of coverage ▶ Readily standardisable ▶ Rapid interpretation for clinical use ▶ Low costs ▶ Easy clinical implementation 	<ul style="list-style-type: none"> ▶ Detection of unknown variants ▶ Detection of CNVs ▶ Research applications ▶ Feasible in clinical routine ▶ Low price/performance ratio 	<ul style="list-style-type: none"> ▶ Comprehensive assessment of cancer genomes ▶ Highest resolution of genomic alterations ▶ SNVs in enhancer/promoter and ncRNA regions ▶ Decreasing costs ▶ Subject to future studies
Contra	<ul style="list-style-type: none"> ▶ Limited, 'peephole' observations ▶ Limited value for research ▶ Limited assessment of complex aberrations 	<ul style="list-style-type: none"> ▶ Not fully comprehensive ▶ Lower CNV resolution ▶ Amplification or exon capture necessary ▶ High bioinformatic effort ▶ Demanding clinical interpretation ▶ Time-consuming workflow 	

CNV, copy number variant; ncRNA, non-coding RNA; SNV, single nucleotide variant.

Horak P, et al. ESMO Open 2016;1:e000094.

WGS use is expected to reduce the cost of additional tests



Estimated costs per sample for small- and medium TGP and WGS; van Amerongen et al., Ecancermedicalscience. 2016

Curated Variant Data Resources as Backbone of Personalised Cancer Therapy

- cancer variant interpretation resources apply manual **data curation** and **bioinformatics** methods to provide information about putative targets and possible interventions

Database	Institute	Organized by
TARGET	BROAD	Gene
PCT	MD Anderson	Gene
cBioPortal / OncoKB	MSK	TCGA diseases
COSMIC	Sanger	Gene
IntOGen	University Pompeu Fabra	Gene
My Cancer Genome	Vanderbilt	Disease
CIViC	Washington University	Variant
DGIdb	Washington University	Drug/gene interaction

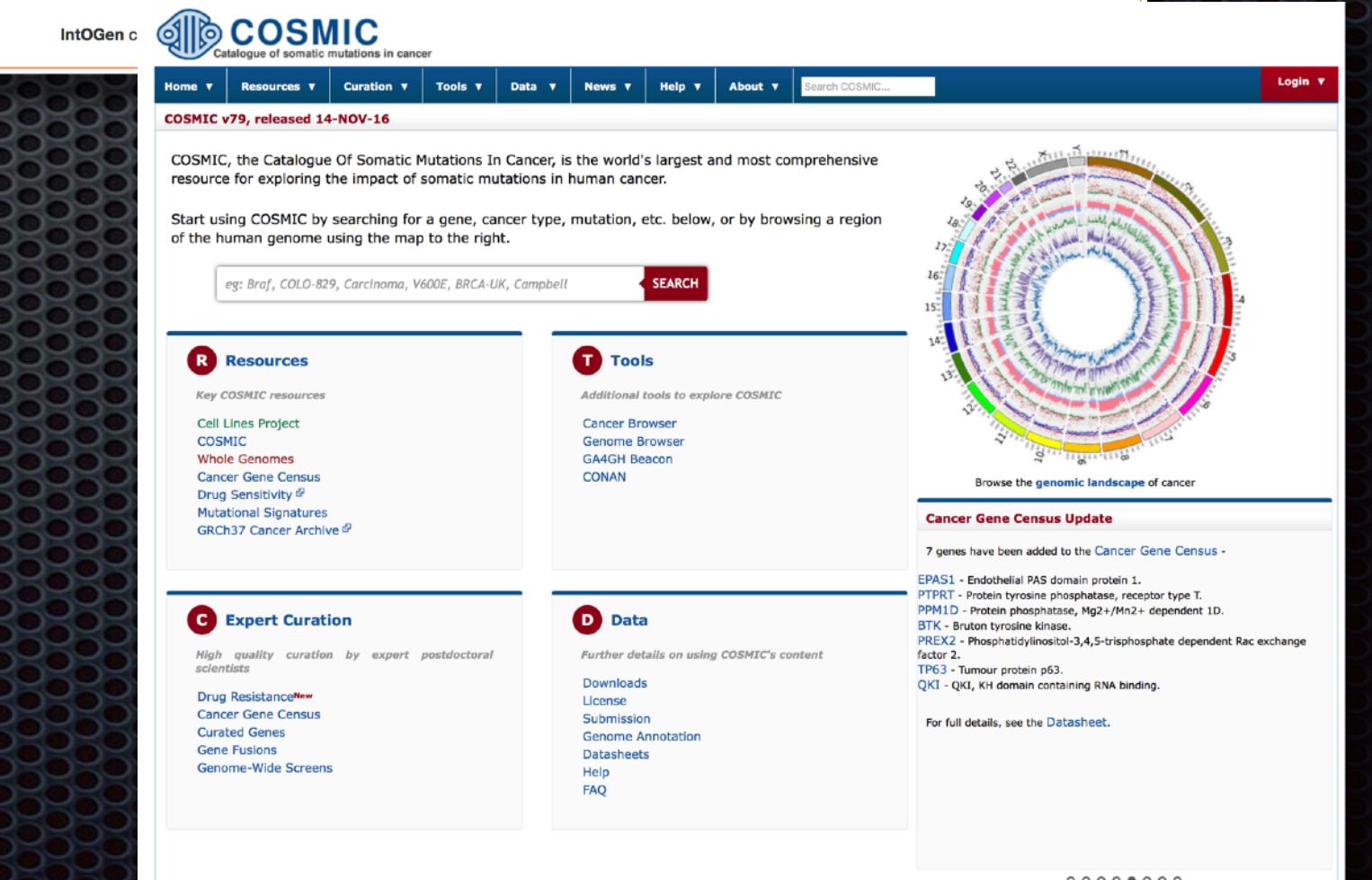
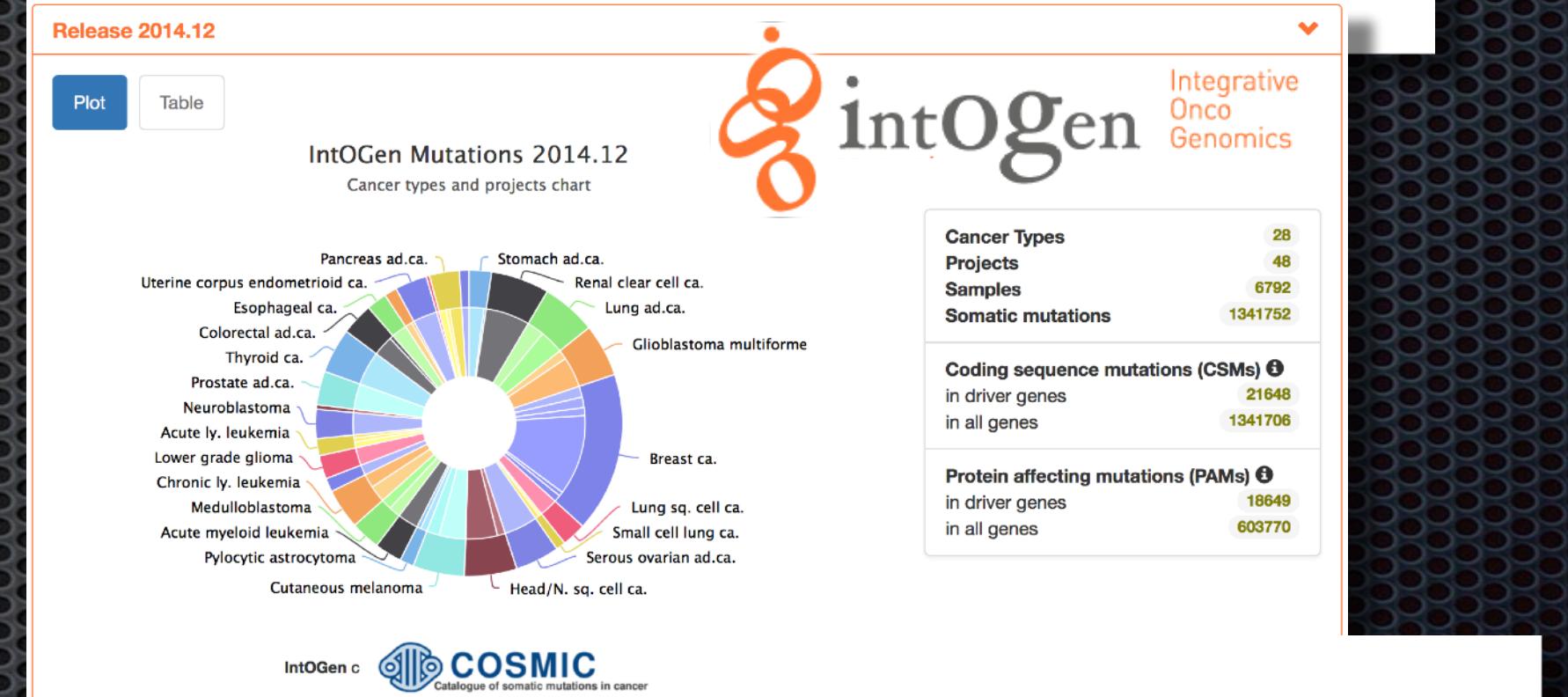
OncoKB Home About Team Levels of Evidence Actionable Genes Data Access News

OncoKB
Precision Oncology Knowledge Base
Annotation of Somatic Mutations in Cancer

418 Genes 3332 Variants 50 Tumor Types 71 Drugs

Search Gene

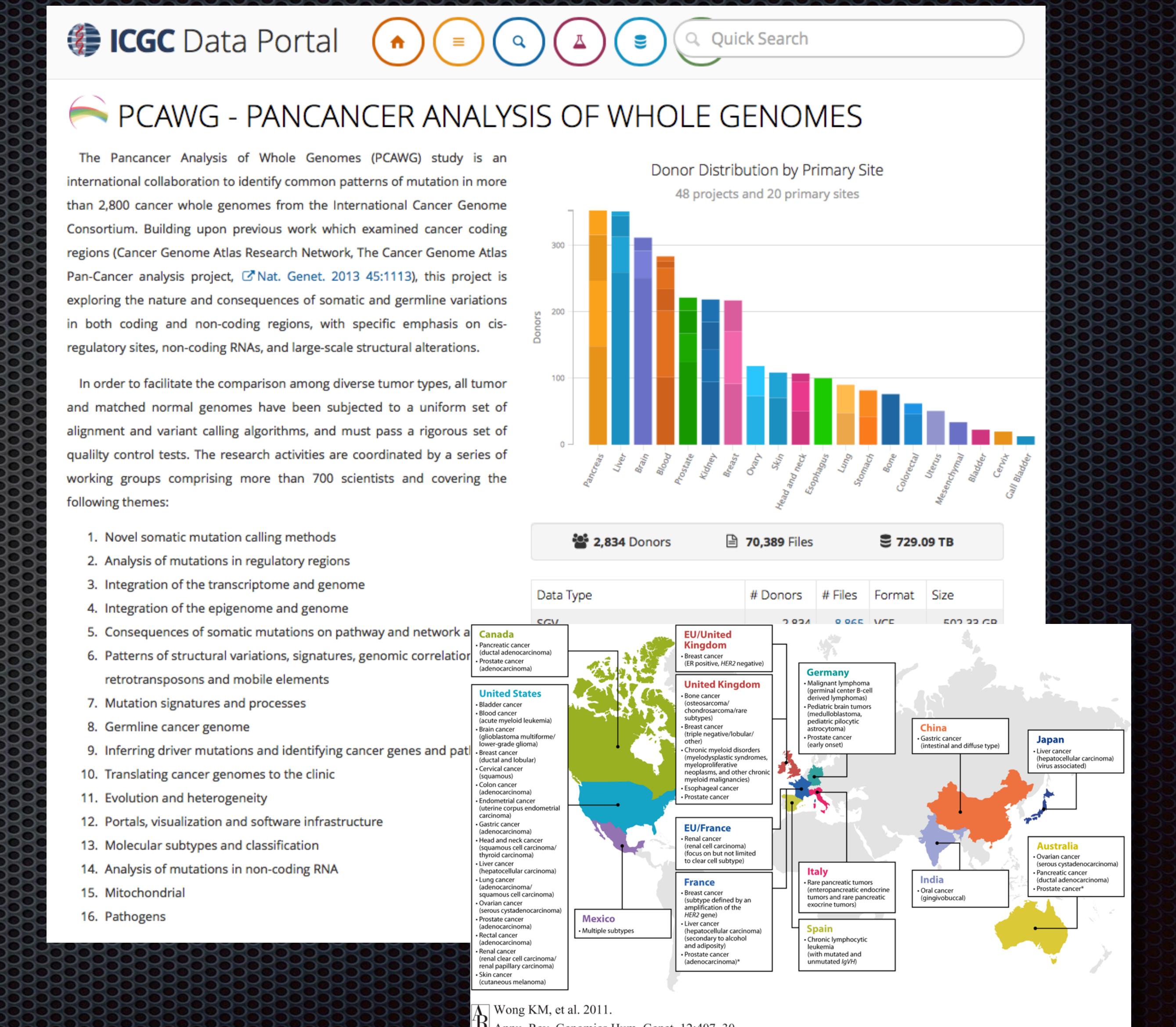
Level 1 FDA-approved Level 2 Standard-of-care Level 3 Clinical evidence



Genome Data Access in Cancer Mining for **New** Knowledge

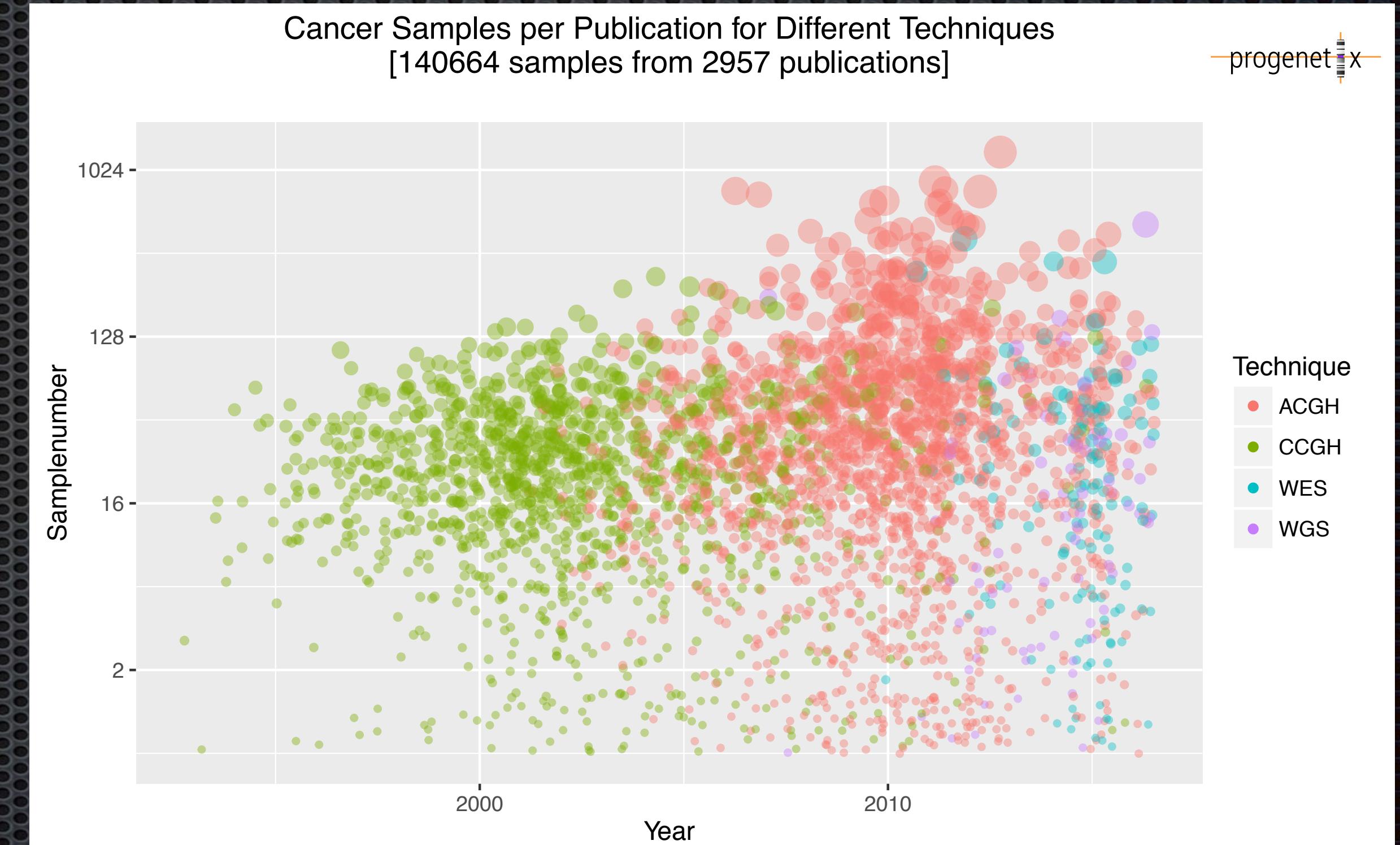
Genome-wide multi"omics" data generation for understanding tumor biology

- the International Cancer Genome Consortium (ICGC) as leading example of deep analysis of multiple cancer entities
- international collaboration of leading research centers for each of ~20 tumor types
- limitations:
 - focus on prominent cancer types w/ limited representation of rare entities
 - data access policies influenced by national regulations and legal frameworks
 - technical heterogeneity



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

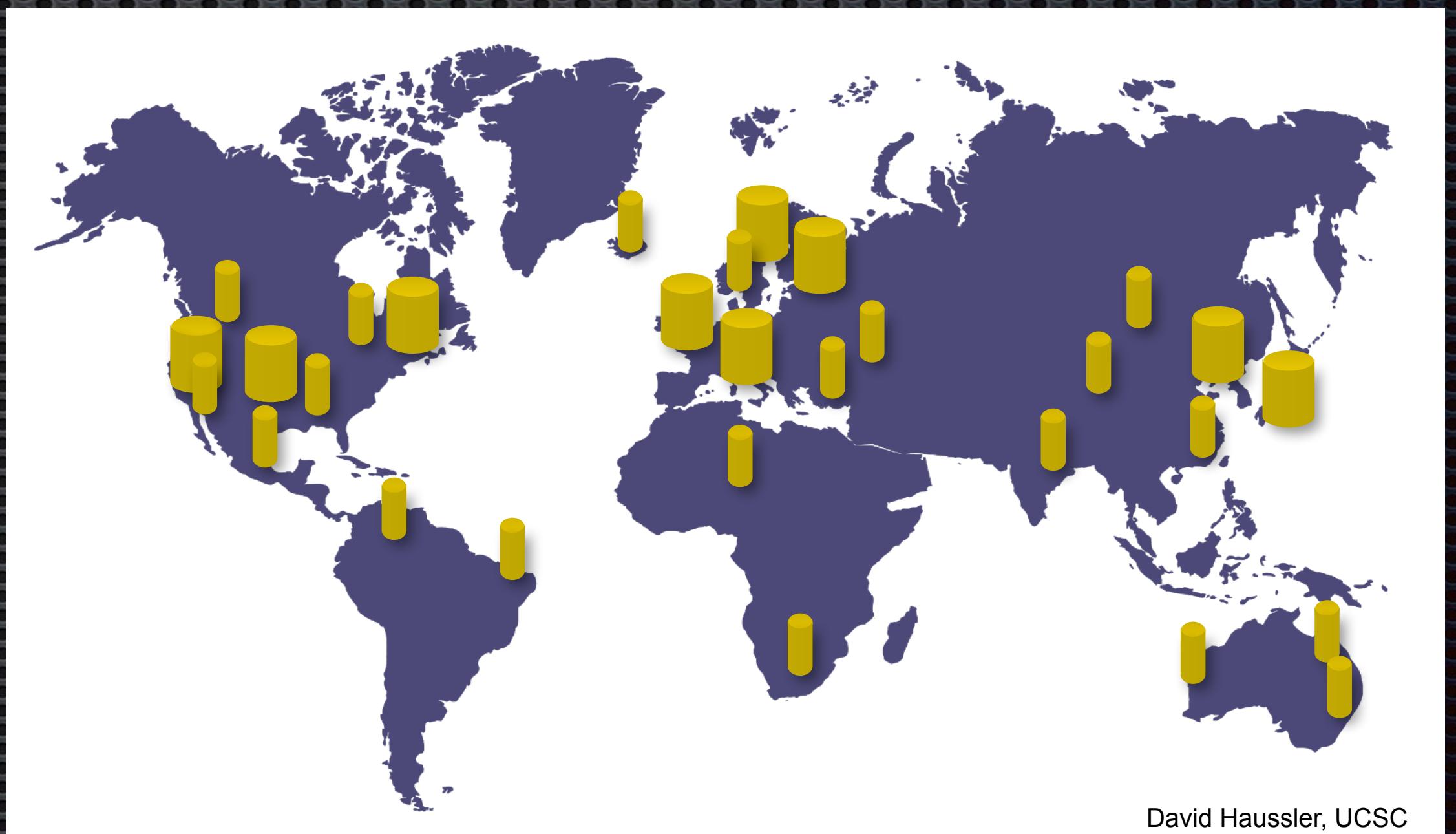
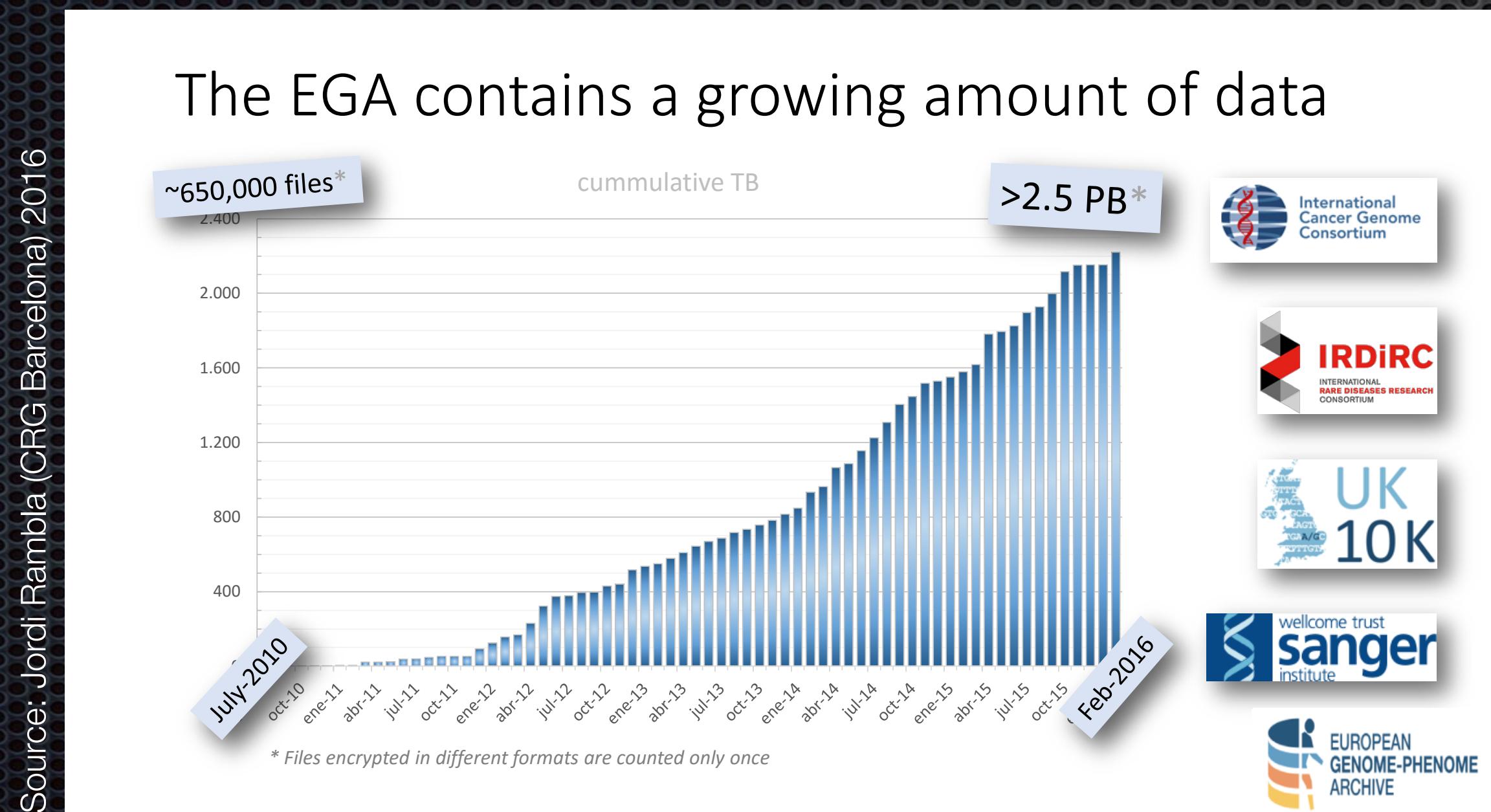
- genome screening to identify mutations in cancer samples
- for diagnostic purposes and therapeutic target identification
 - karyotyping (~1968)
 - Comparative Genomic Hybridization (1992)
 - genome **microarrays** (aCGH, SNP arrays ...; 1997)
 - Whole Exome Sequencing** (2010)
 - Whole Genome Sequencing**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



Overview of publications reporting whole-genome screening analysis of cancer samples, by molecular-cytogenetic or genome sequencing methods. The data represents articles assessed for the progenetix.org cancer genome data resource (M. Baudis, 2001-2016)

Genome Datasets: Rapid Growth, Limited Access

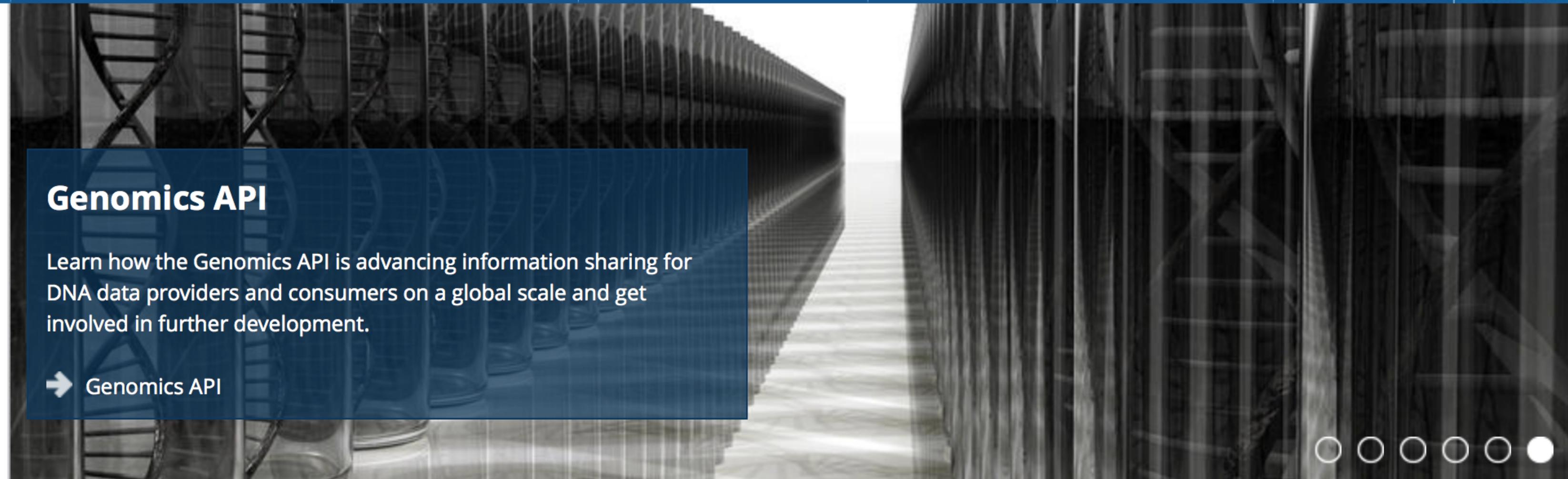
population based and cancer research studies produce a rapidly increasing amount of genome sequence data



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

Genomes Everywhere

Organization / Initiative: Name	Organization / Initiative: Category	Cohort
100K Wellness Project	Research Project	107 unaffected individuals (scaling up to 100,000)
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?)
BioBank Japan	Repository	Specimens from >200,000 patients and unaffected controls
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)
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	currently data from >16'000 samples
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)
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)
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)
UK10K	Research Project	10,000 participants (6,000 patients and 4,000 controls)
Vanderbilt's BioVU	Repository	>215,000 samples



Genomics API

Learn how the Genomics API is advancing information sharing for DNA data providers and consumers on a global scale and get involved in further development.

→ [Genomics API](#)



Our Work

The diverse members of the Global Alliance are working together to create interoperable approaches to catalyze projects that will help unlock the great potential of genomic data. Our four [Working Groups](#) advance [Initiatives](#) that develop key [Work Products](#).



Clinical »

Aims to enable compatible, readily accessible, and scalable approaches for sharing clinical data and linking it with genomic data.



Data »

Concentrates on data representation, storage, and analysis of genomic data to develop approaches that facilitate interoperability.



Regulatory and Ethics »

Focuses on ethics and the legal and social implications of the Global Alliance, including harmonizing policies and standards.



Security »

Leads the thinking on the technology aspects of data security, user access control, audit functions, and developing or adopting data security standards.

GA4GH MEMBERSHIP

▶ host institutions

- OICR, Broad Institute, Sanger Institute
- Peter Goodhand, OICR, exec. director



▶ funding

- CanSHARE, NIH, Wellcome Trust



National Institutes
of Health

▶ founding members

- 229 partner organizations based in 30 countries as of October 3, 2014
- SIB and UZH represented Switzerland 



Swiss Institute of
Bioinformatics

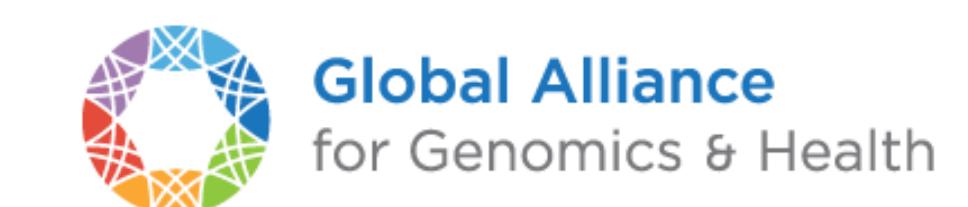
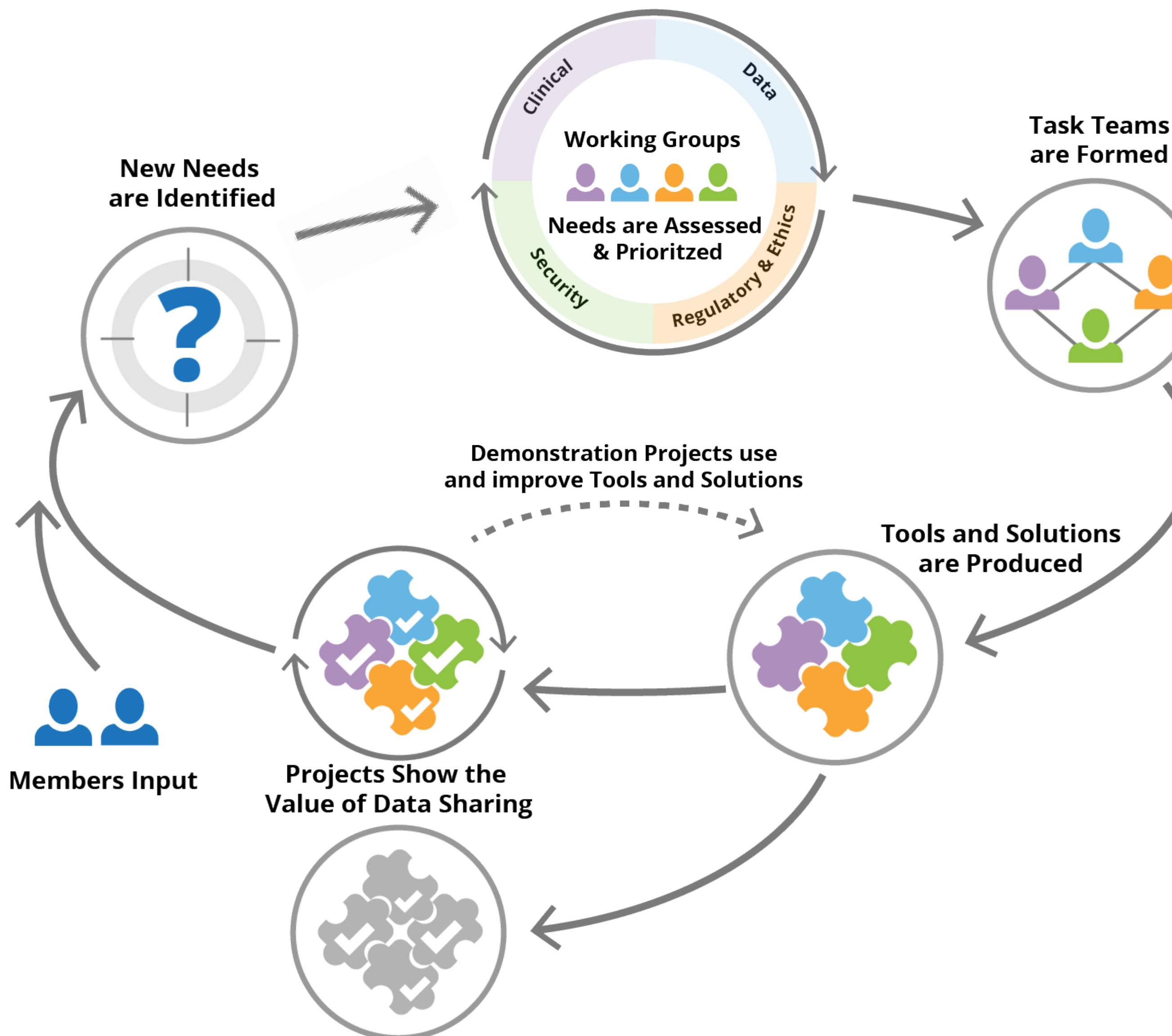


▶ membership status

- 433 organizational members as of September 2016
- open to individual registrations & participation
- no financial commitment necessary; however, calls for support for events



How We Work





This repository

Search

Pull requests Issues Gist



ga4gh / schemas

[Unwatch](#) 115[Star](#) 196[Fork](#) 110[Code](#)[Issues 152](#)[Pull requests 29](#)[Projects 1](#)[Wiki](#)[Pulse](#)[Graphs](#)

Work on data models and APIs for Genomic data. <http://ga4gh.org/#/api>

1,102 commits

17 branches

16 releases

46 contributors

Apache-2.0

Branch: [metadata-integ...](#) ▾[New pull request](#)[Create new file](#)[Upload files](#)[Find file](#)[Clone or download](#) ▾

This branch is 15 commits ahead, 3 commits behind master.

[Pull request](#) [Compare](#)**mbaudis** Merge branch 'master' into metadata-integration

Latest commit 077c2c7 2 days ago



Merge branch 'master' into metadata-integration

2 days ago



Add constraints file

2 days ago



Utilize new common methods in schemas

2 days ago



Merge branch 'master' into metadata-integration

13 days ago



Utilize new common methods in schemas

2 days ago



Merge branch 'master' into metadata-integration

13 days ago



Remove protoc call from install path (#781)

7 days ago



Add constraints file

2 days ago



Convert Avro -> proto3.

10 months ago

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

The Cancer Gene Trust

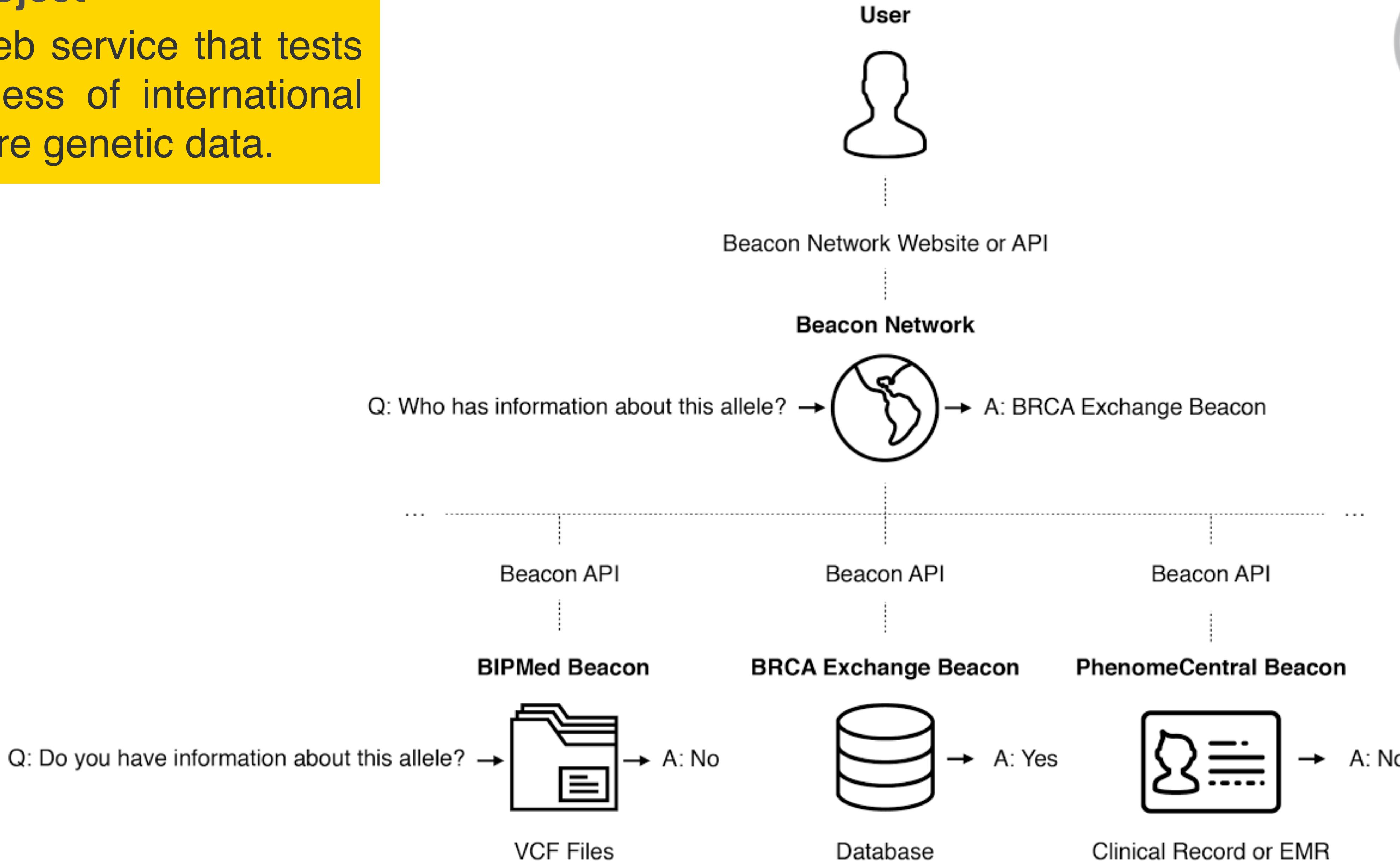
proposes to aggregate somatic cancer mutation data and some clinical data in order to improve the genomic landscape of actionability in some cancers and to enable greater personalized clinical care for individuals with rare cancer mutations.



Cancer
Gene
Trust

Beacon Project

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



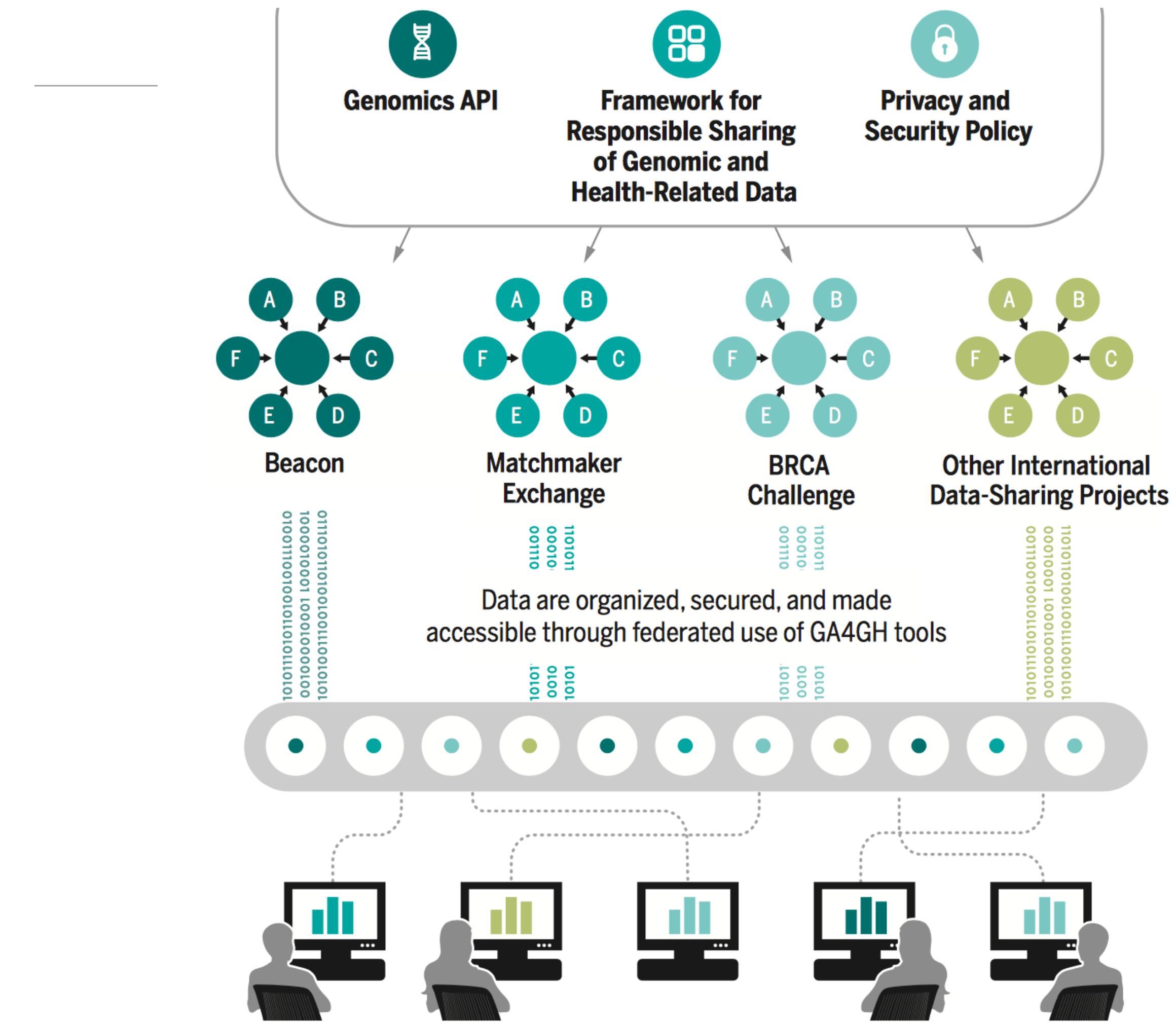
GA4GH API promotes sharing



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.

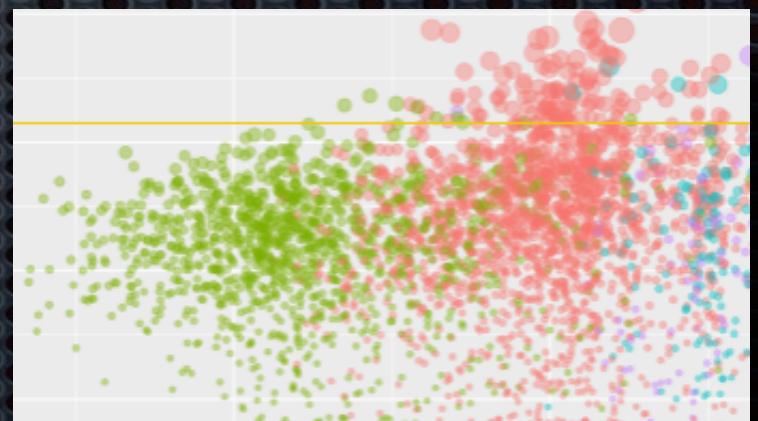
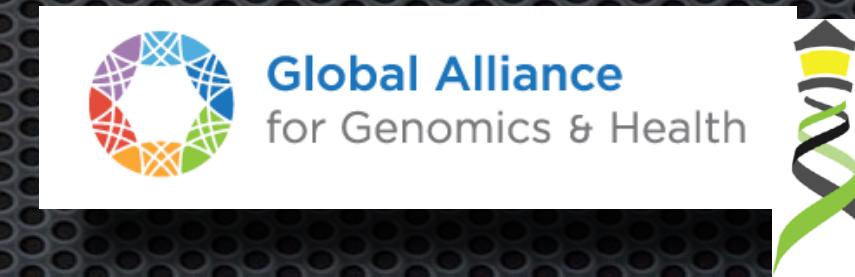
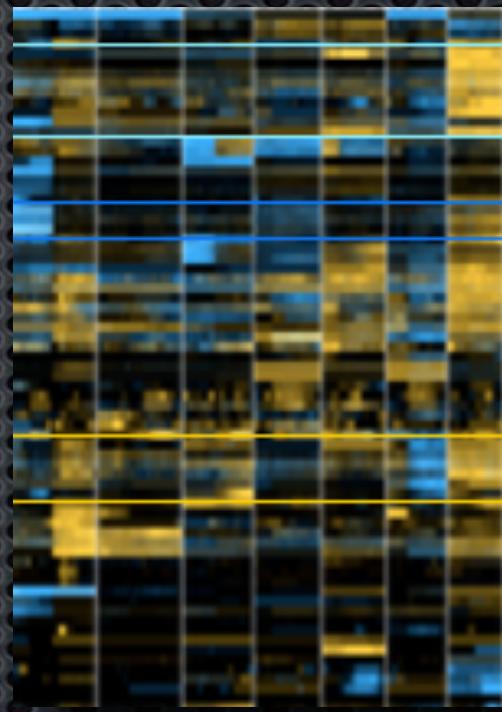
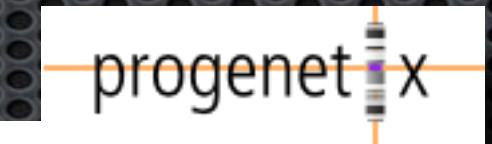




Cancer Genome Data & Beyond

Contributing our share

- 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!



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

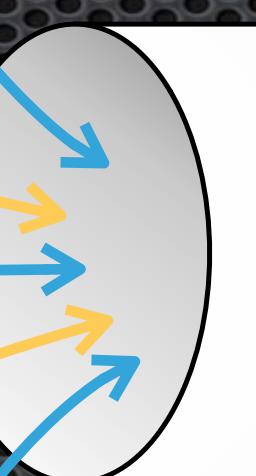
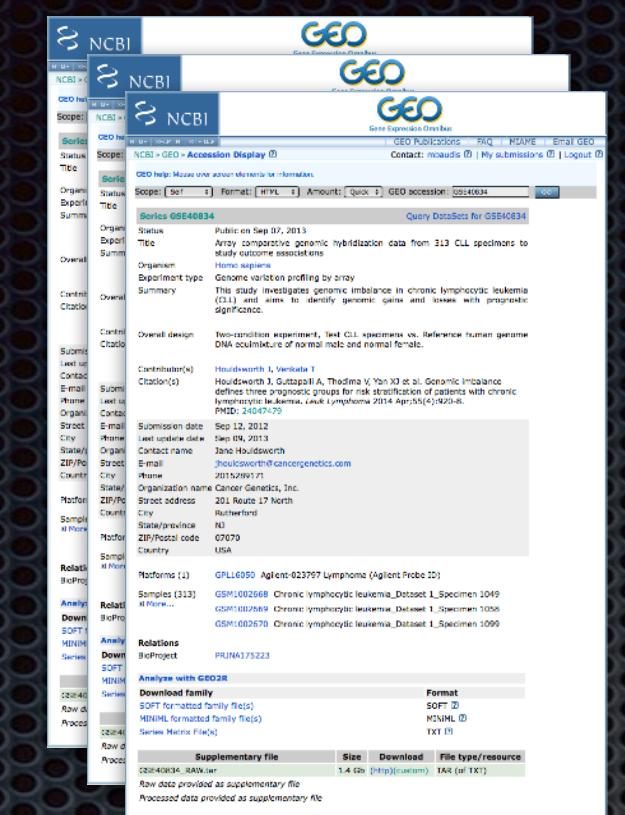
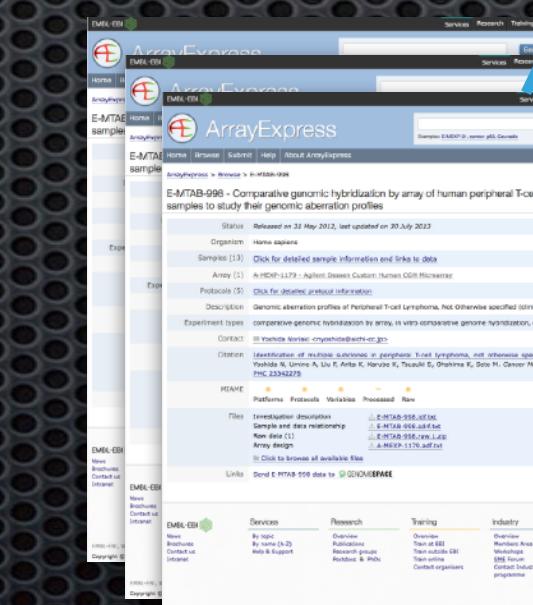
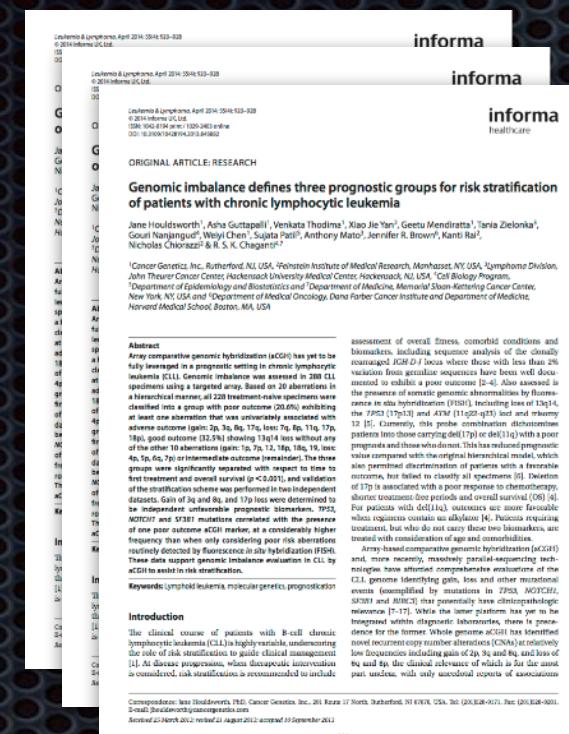


arrayMap

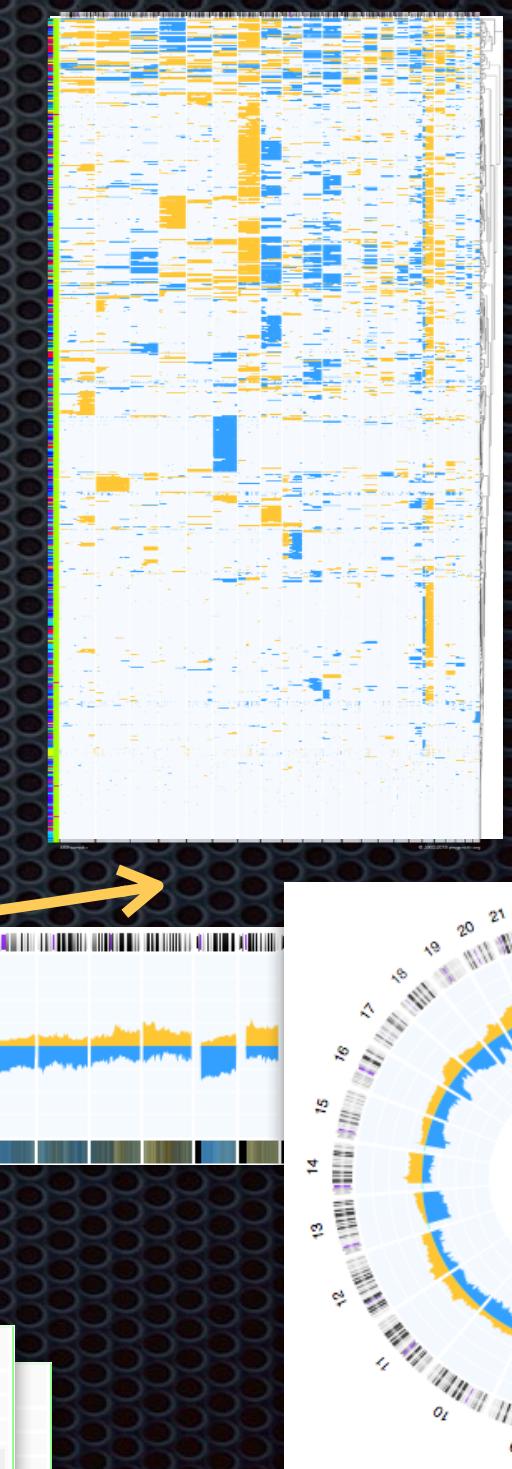
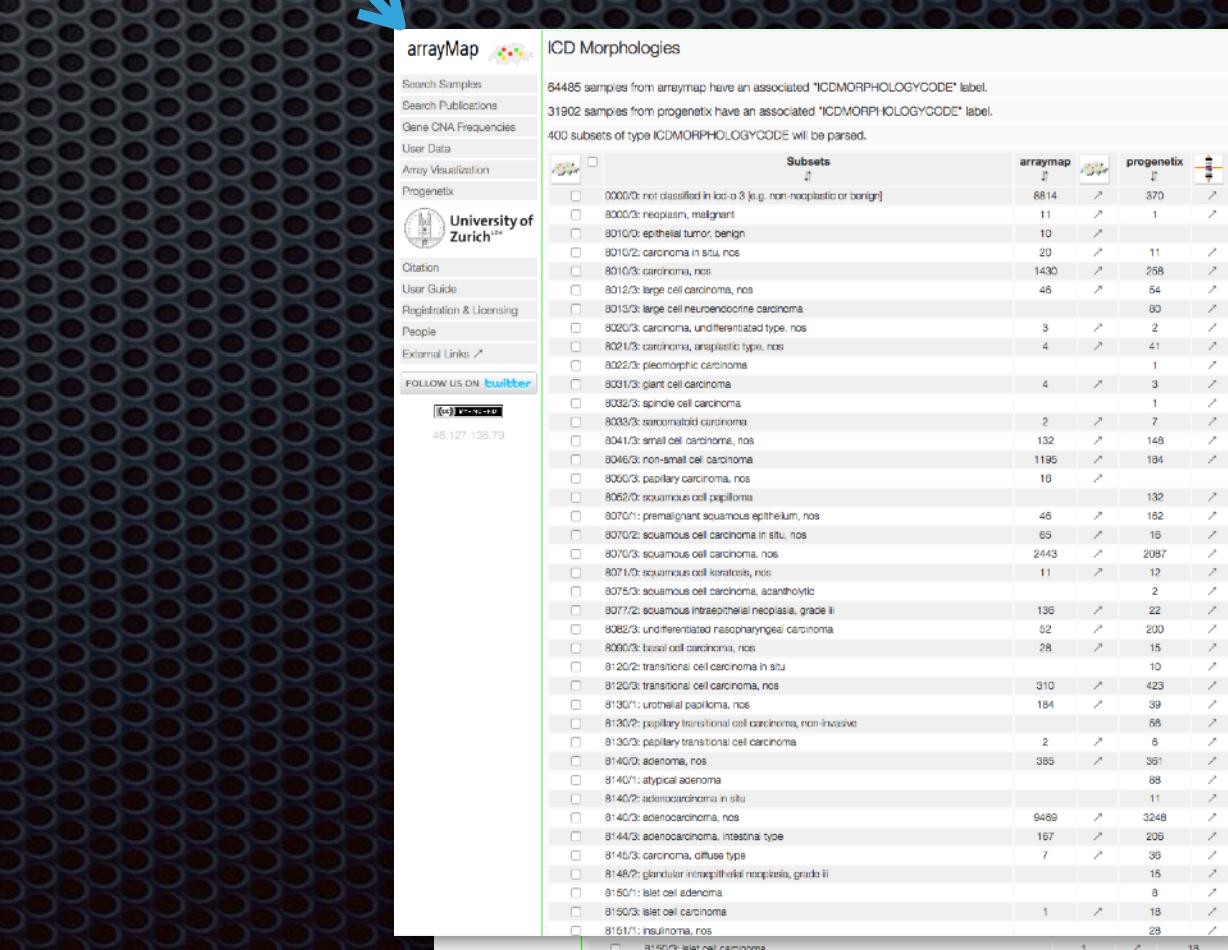
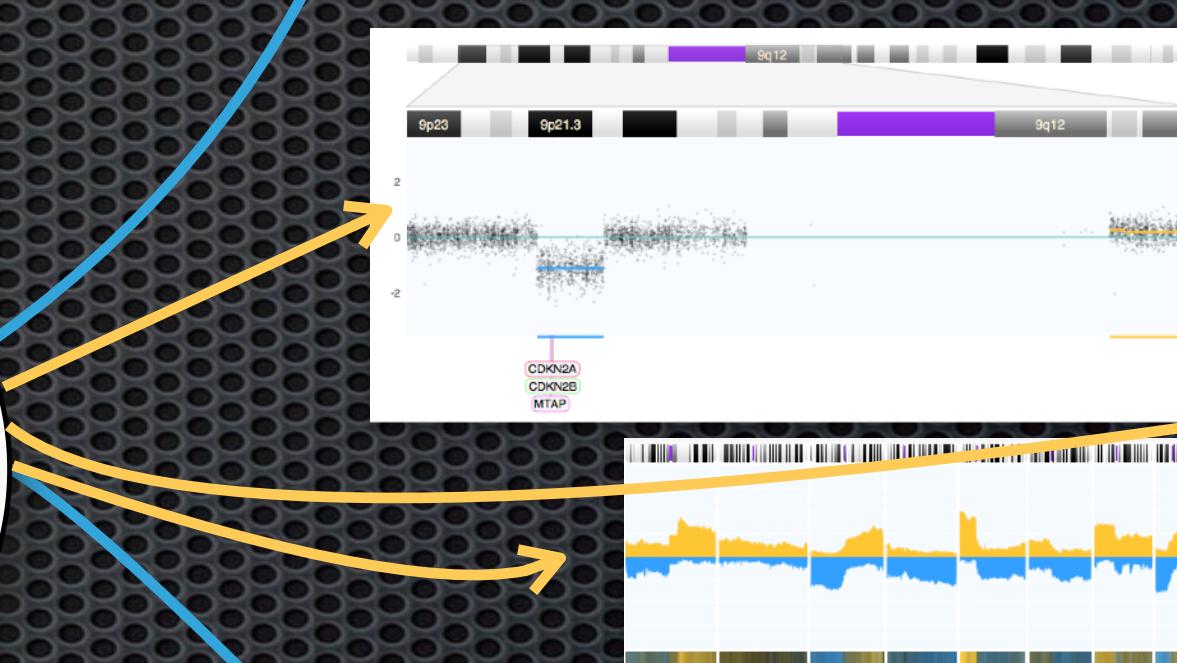
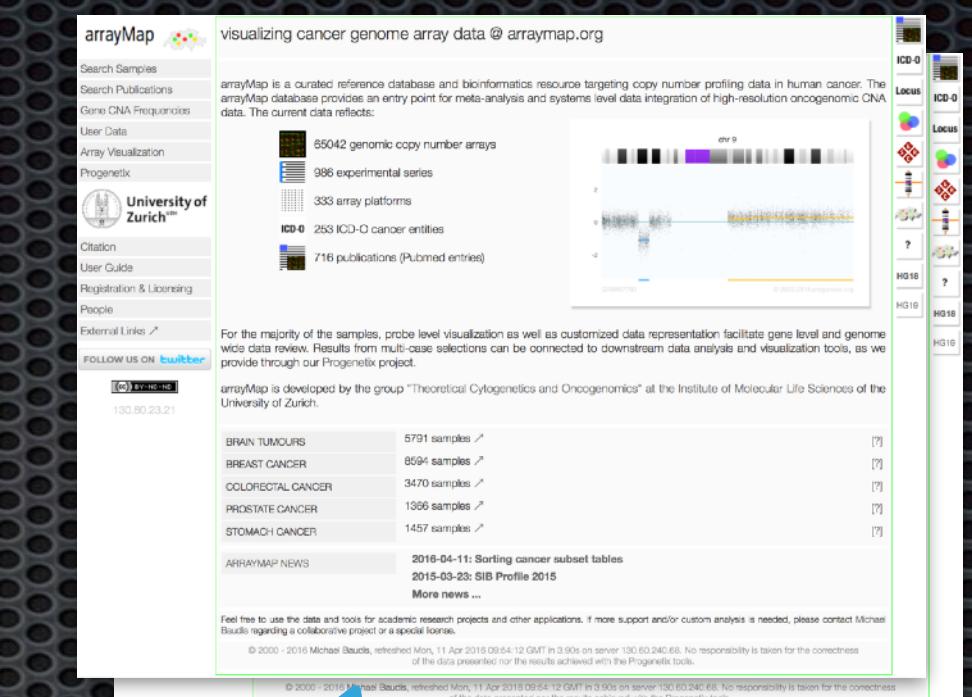


techniques	cCGH, aCGH, WES, WGS	aCGH (+?)
scope	sample (e.g. combination of several experiments)	experiment
content	>31000 samples	>60000 arrays
raw data presentation	no (link to sources if available)	yes (raw, log2, segmentation if available)
per sample re-analysis	no; supervised result (mostly as provided 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

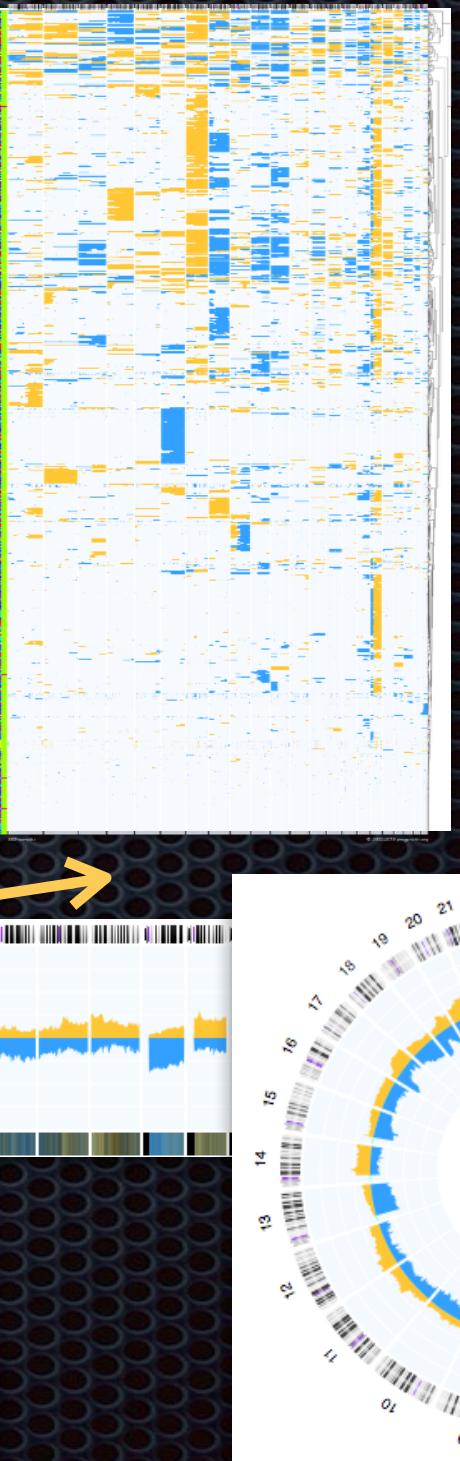
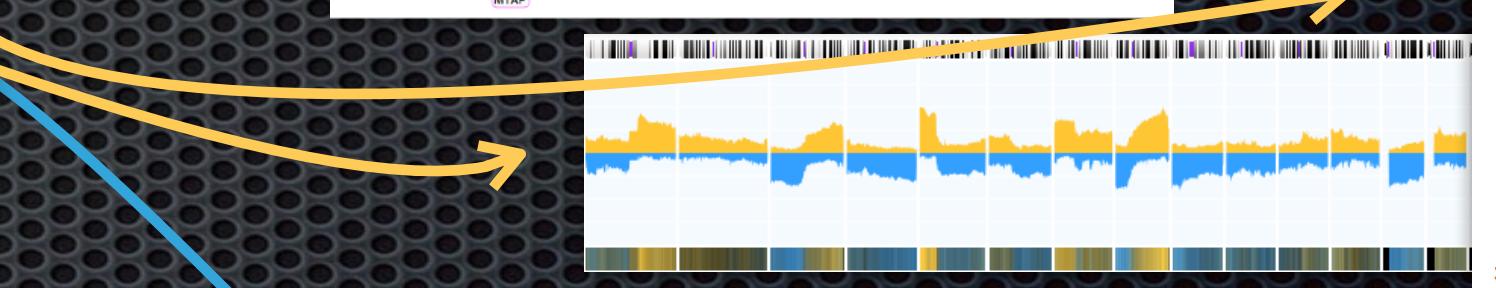
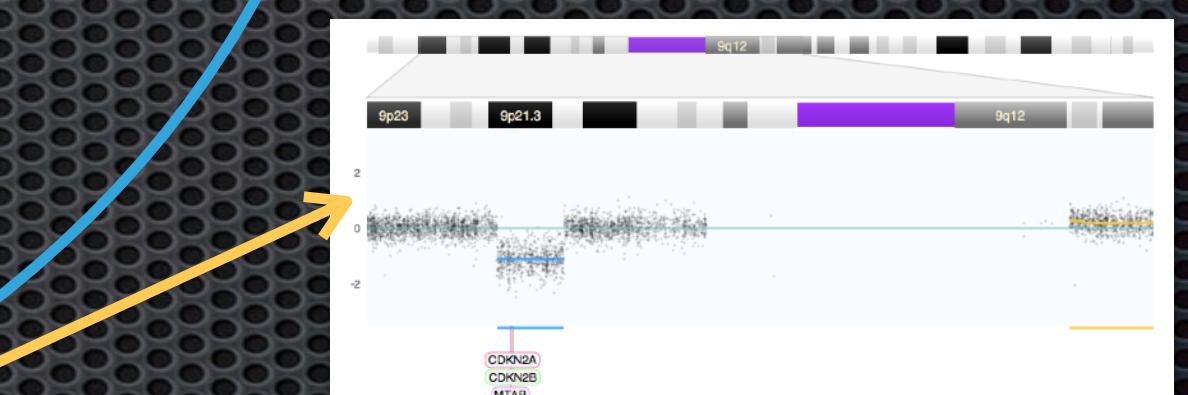
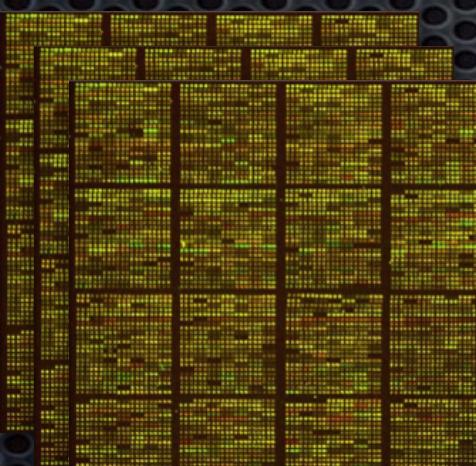
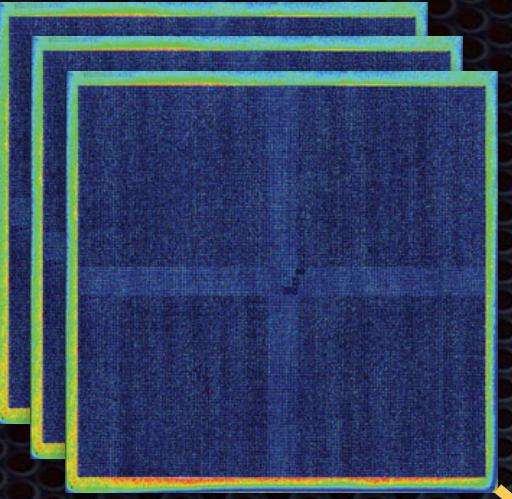


Automagical Processing Engine™



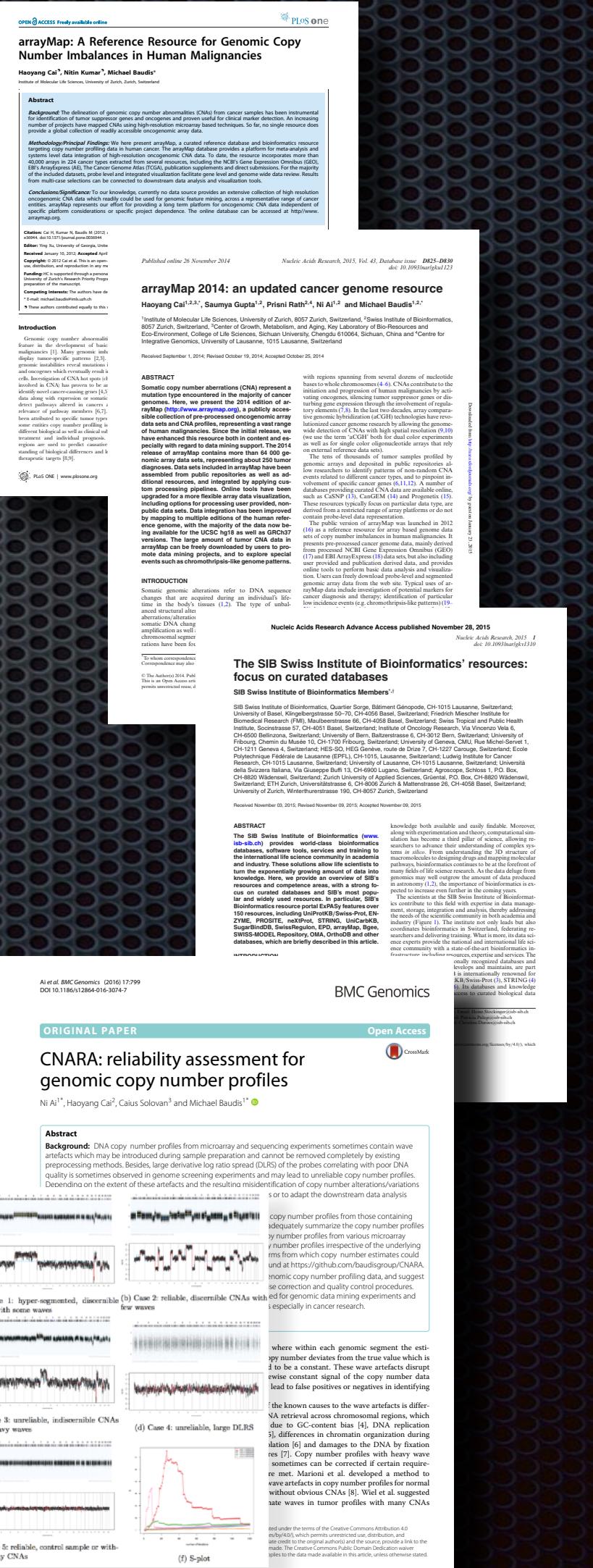
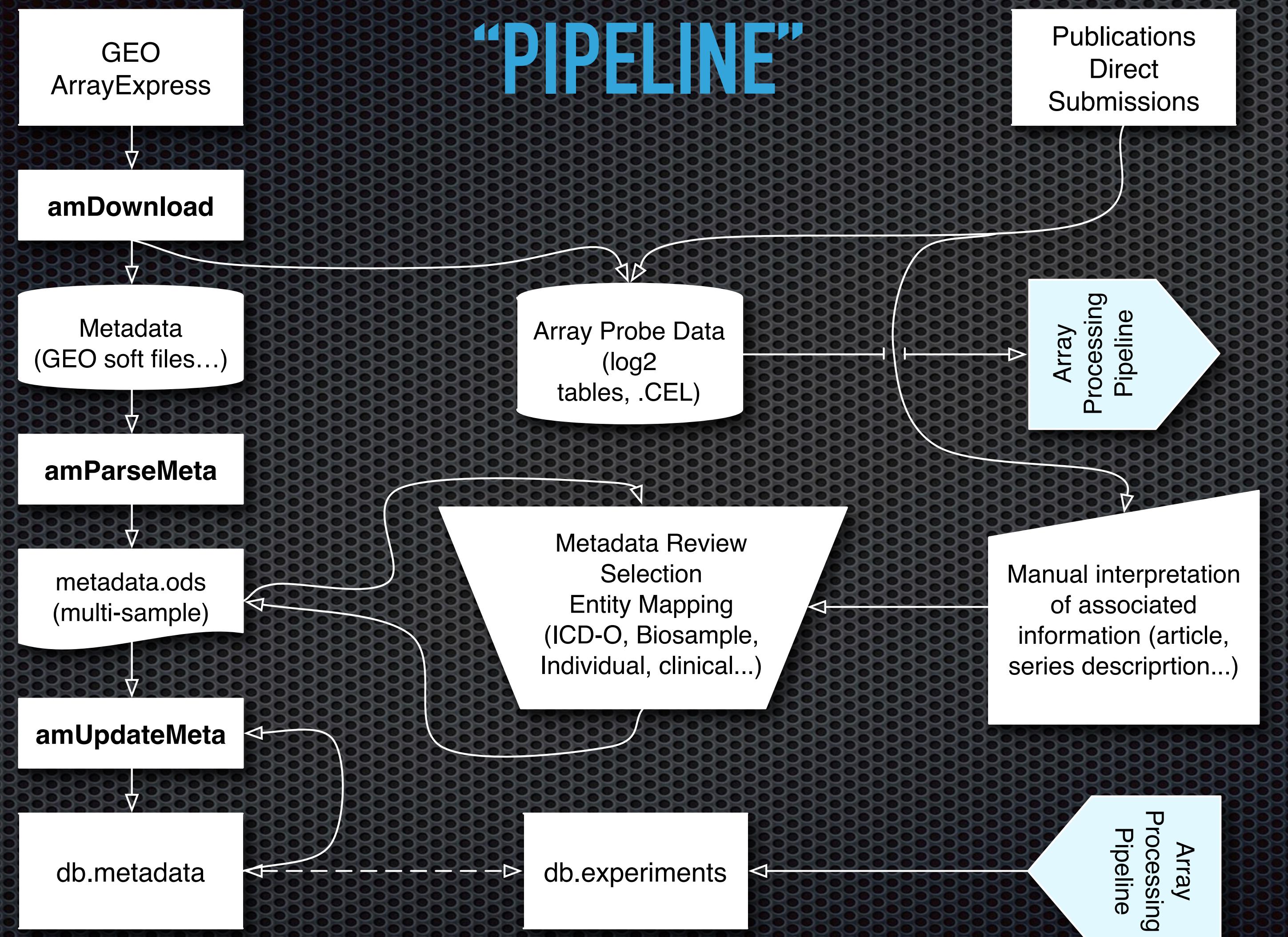
DATA PIPELINE

BIOCURATION BIOINFORMATICS



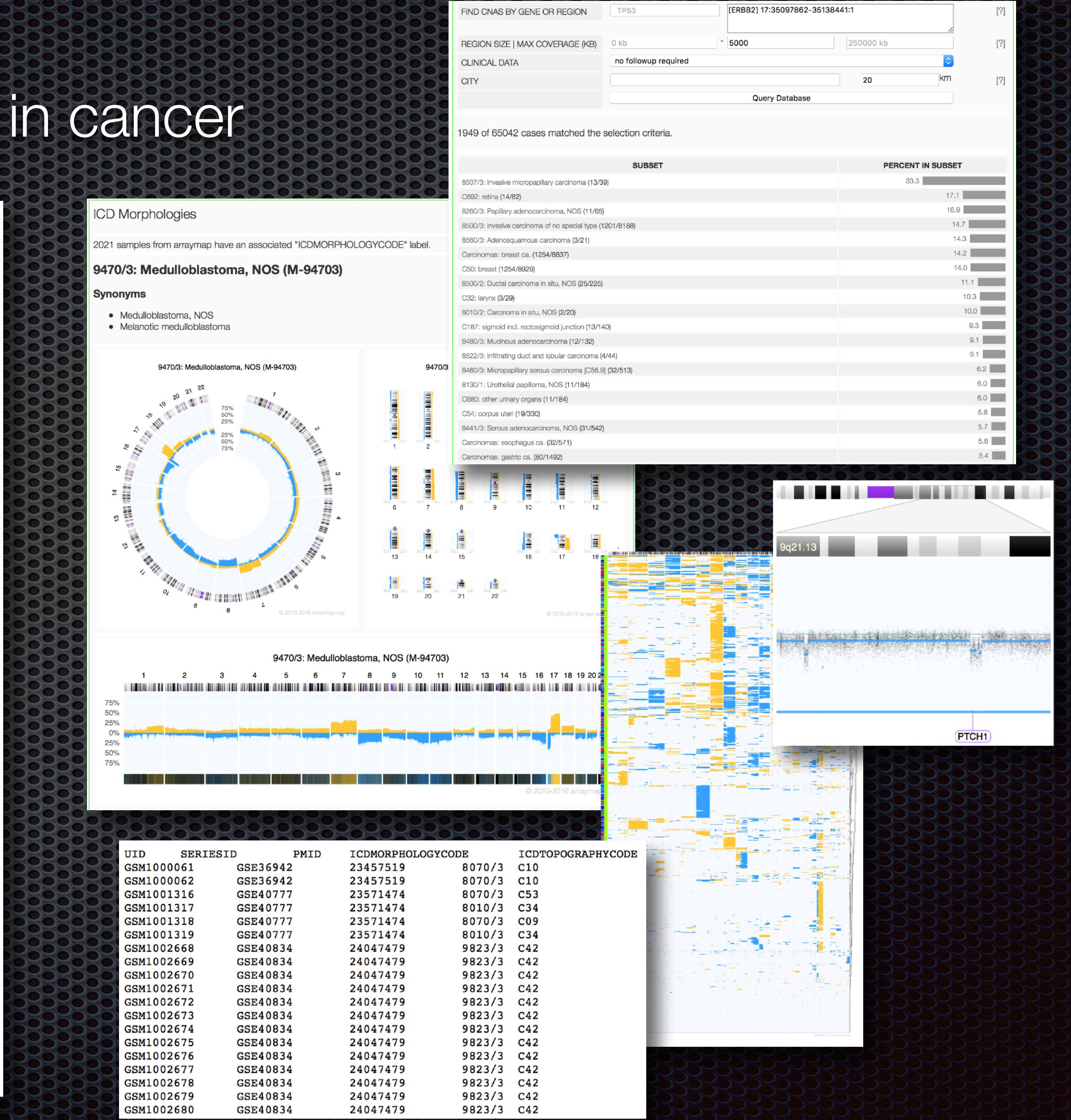
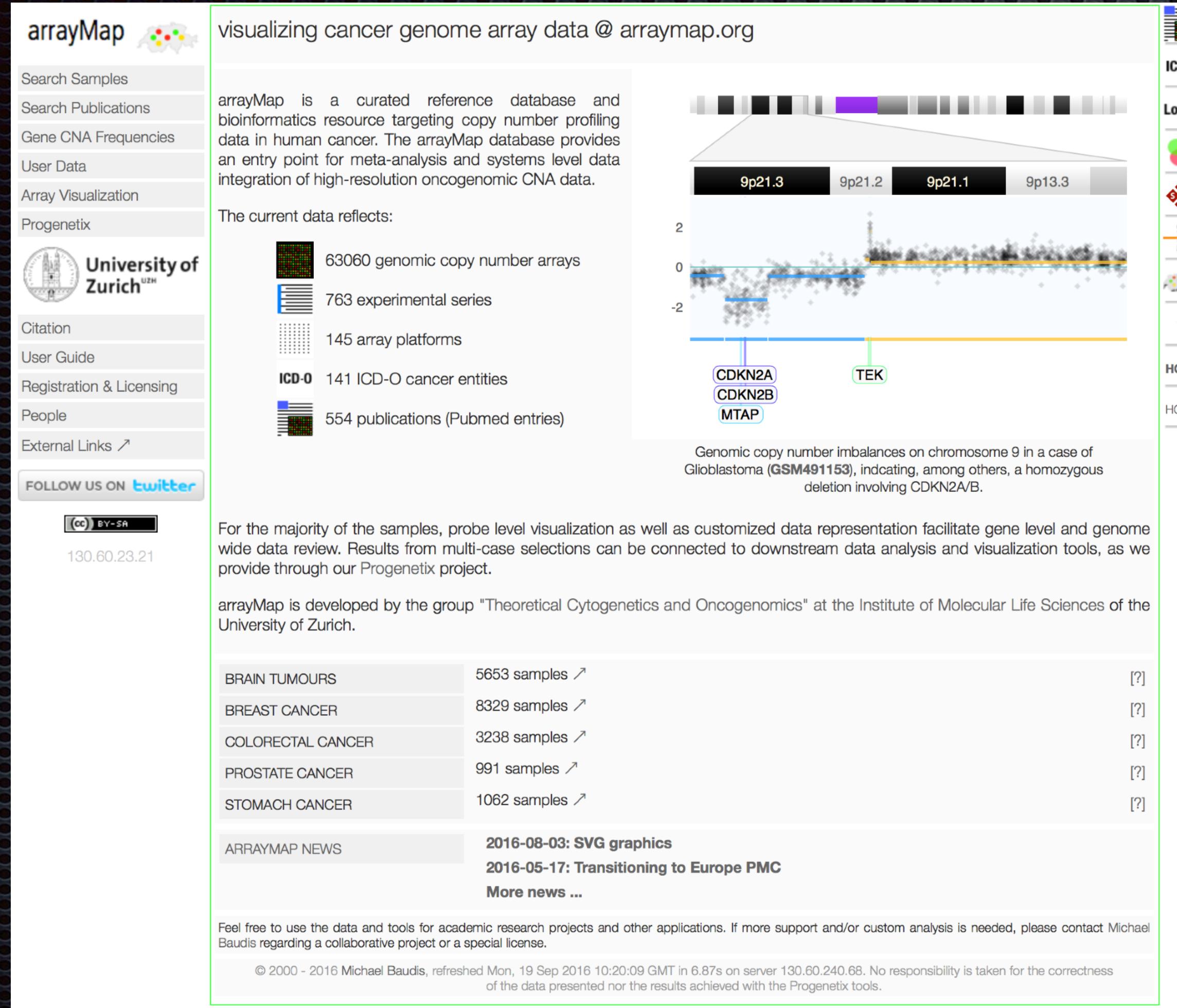
ARRAYMAP DATA

“PIPELINE”



arrayMap

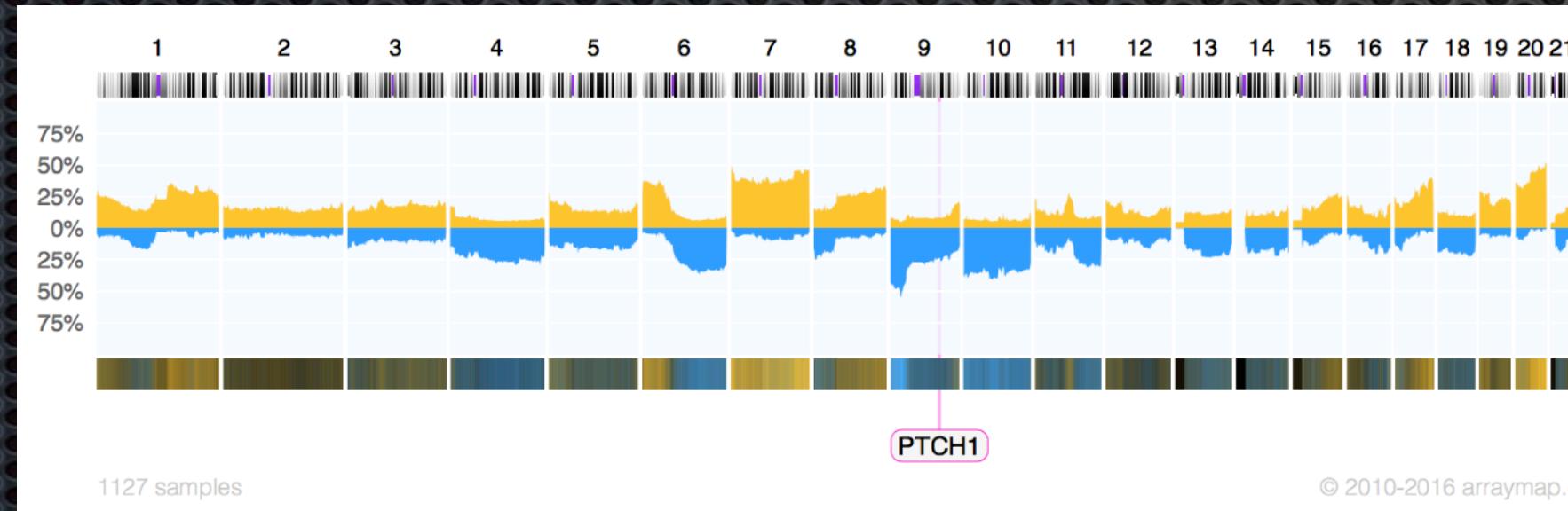
Resource for copy number variation data in cancer



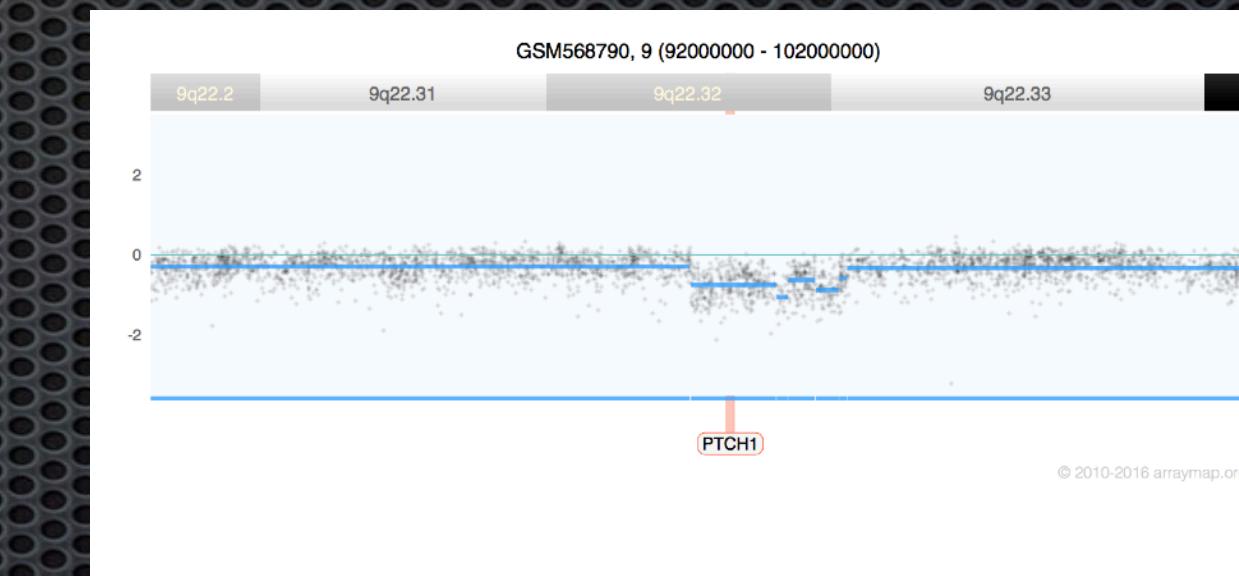
Rare Events & Hidden Therapeutic Options?

Example: PTCH1 deletions in malignant melanomas

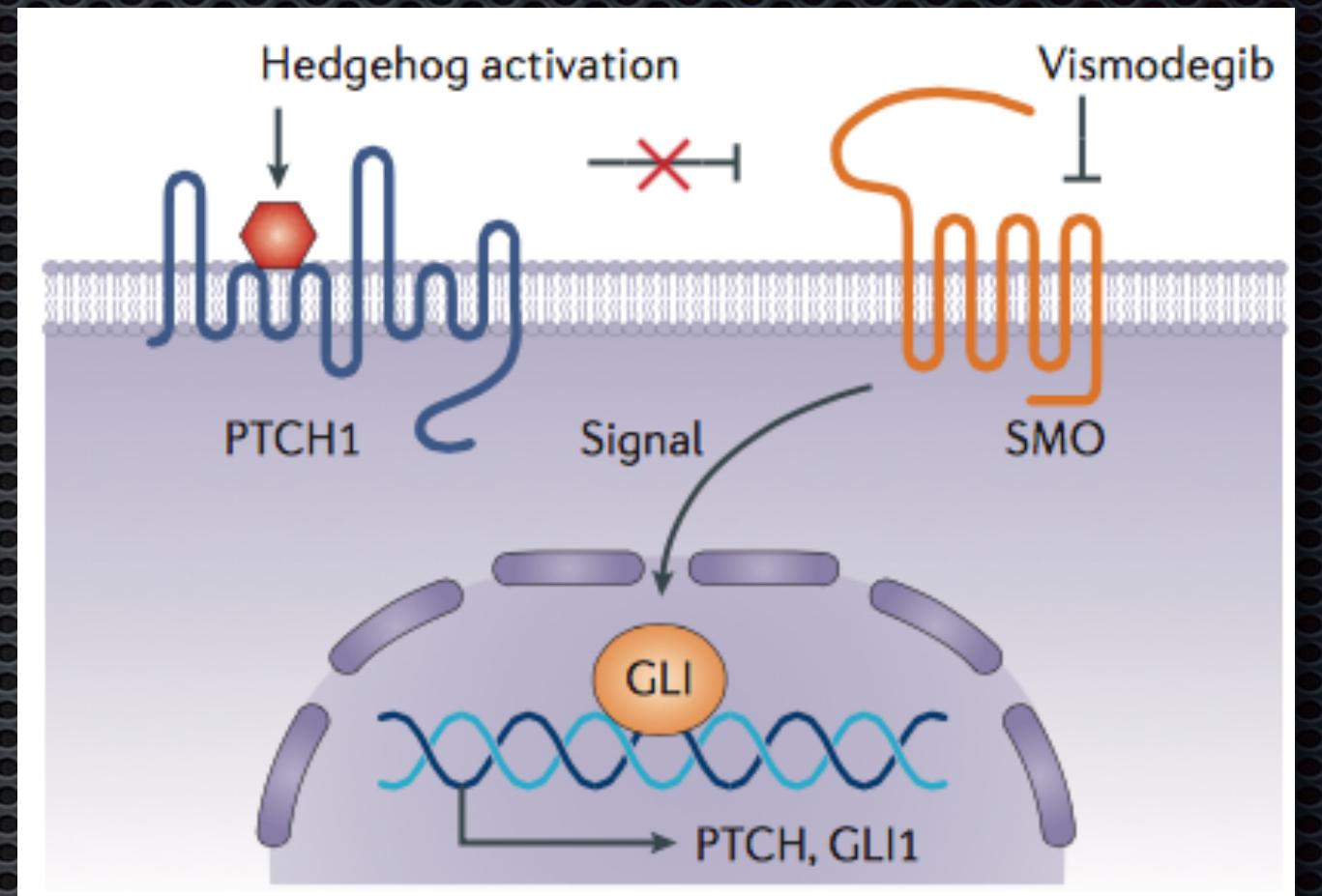
- PTCH1 is a actionable tumor suppressor gene, which has been demonstrated in e.g. basalomas and medulloblastomas
- analysis of 1127 samples from 26 different publications could identify **focal** deletions in 4 samples
- a current project addresses the focal involvement of all mapped genes, in >50'000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



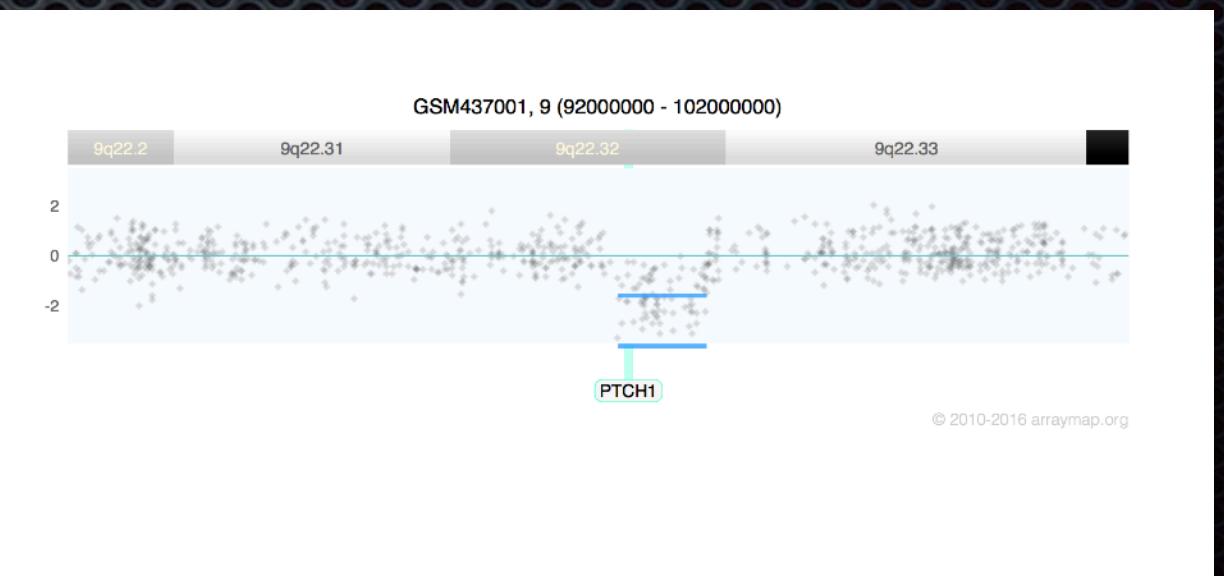
Summary of somatic copy number aberrations from the analysis of 1127 genome profiles of malignant melanomas, collected in our [arraymap.org](#) cancer genome resource. While PTCH1 does not represent a deletion hotspot, the genomic locus is part of larger deletions in ~25% of melanoma samples.



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.

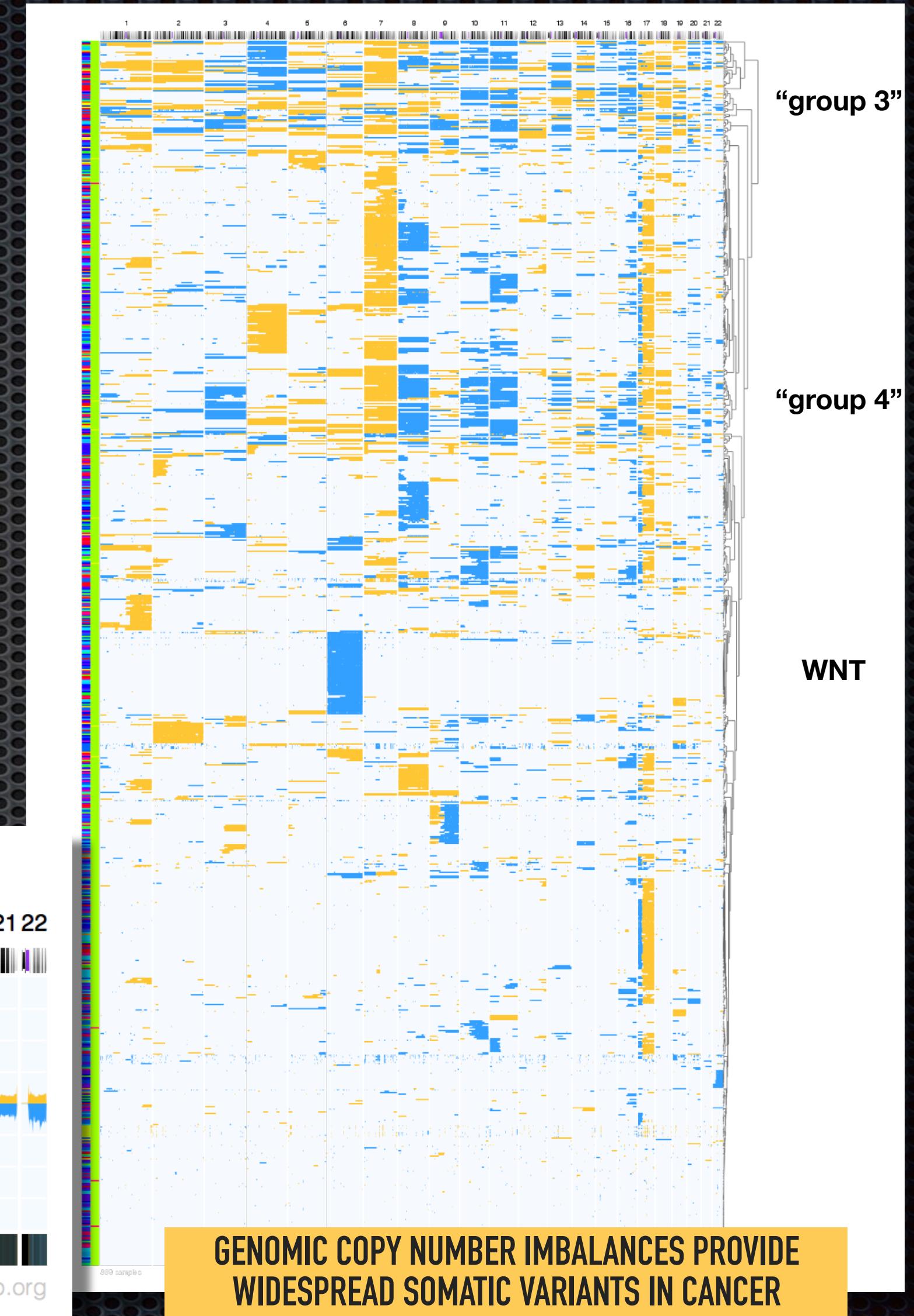
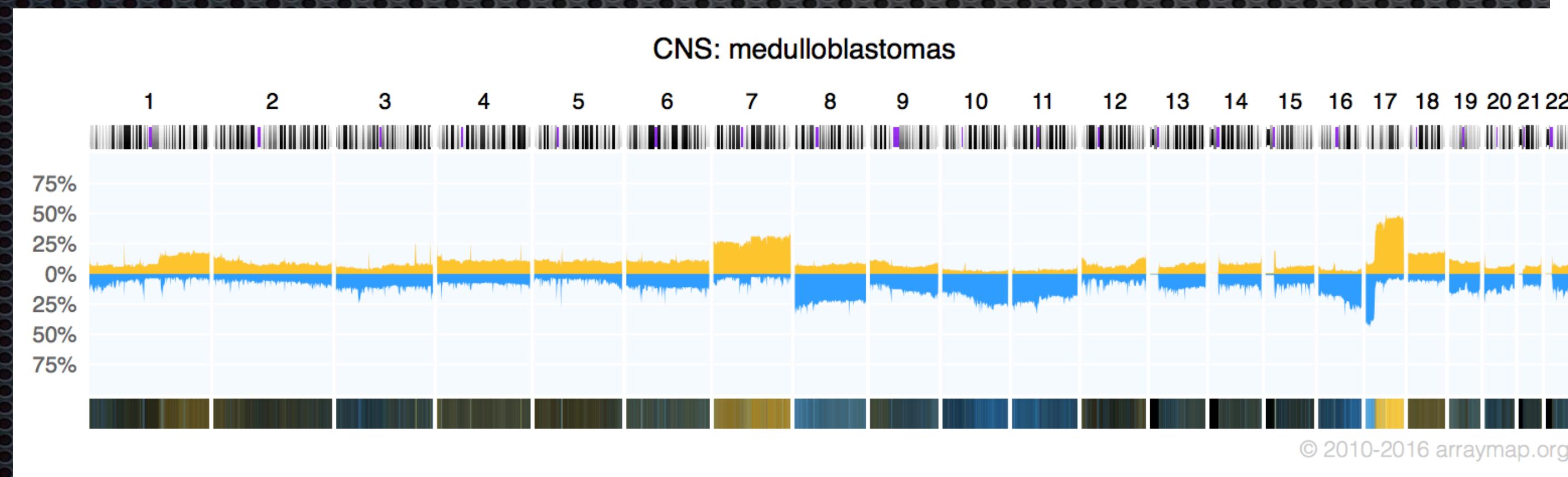


In its normal function, PTCH1 is a tumor suppressor gene in the sonic hedgehog pathway and inhibits SMO driven transcriptional activation. A loss of PTCH1 function (mutation, deletion) can be mitigated through drugs antagonistic to SMO activation.



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?

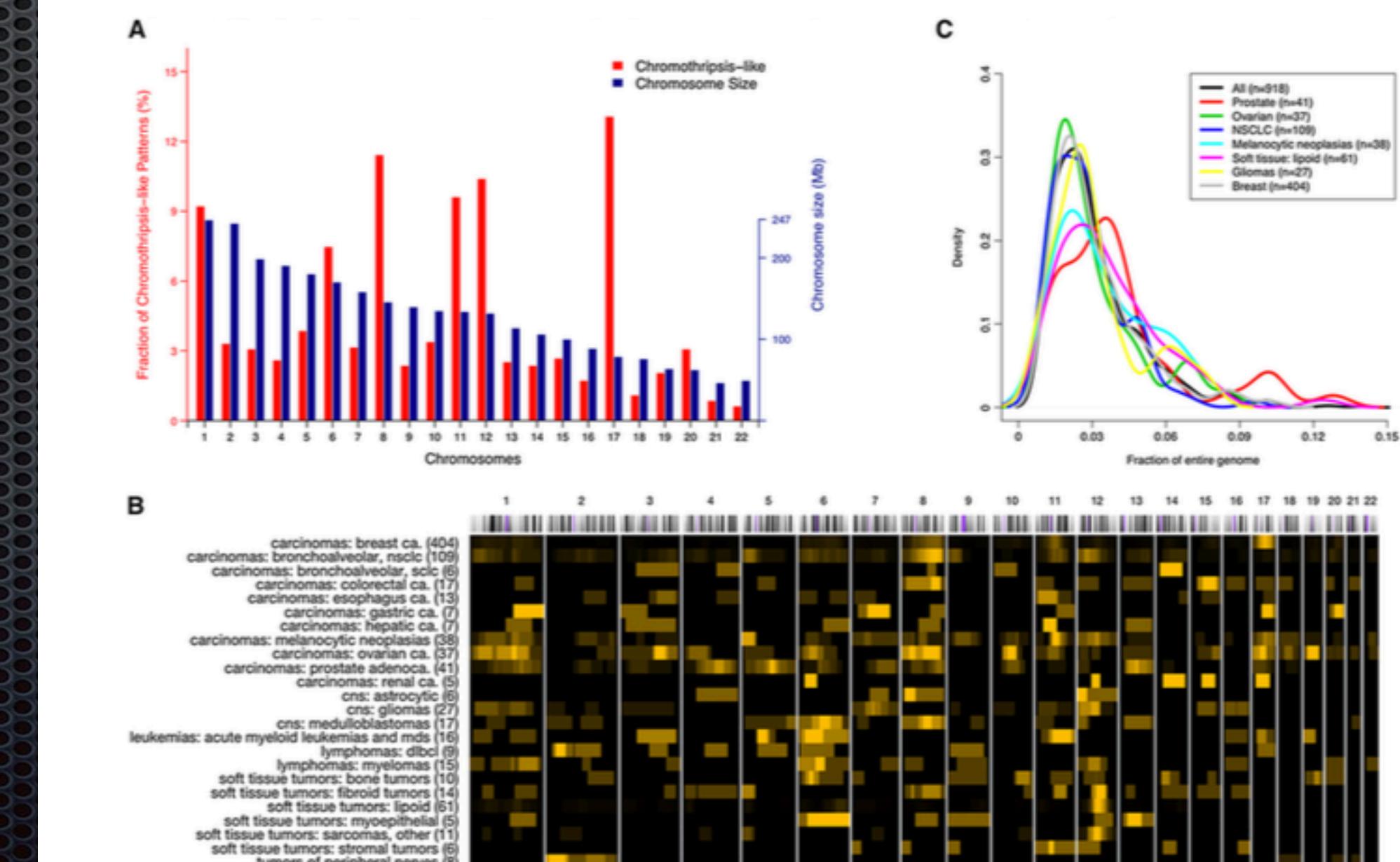
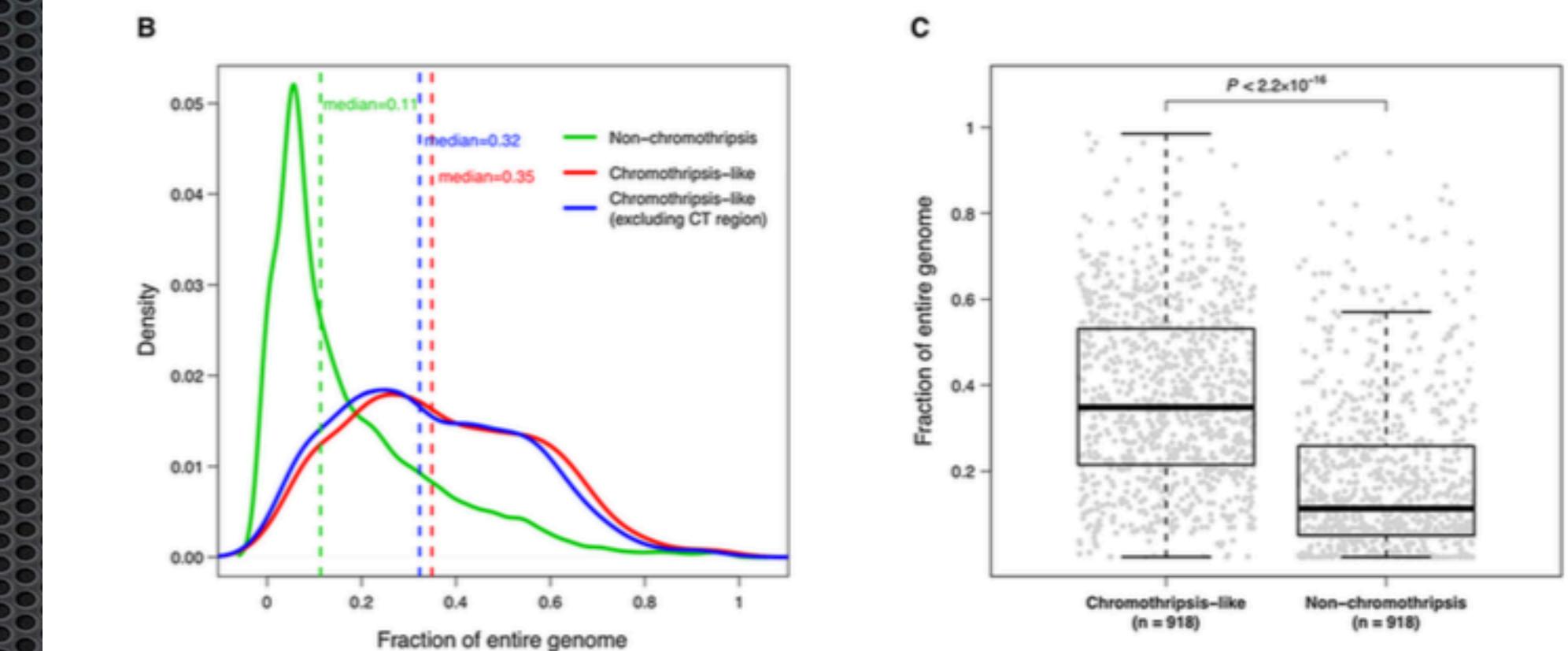
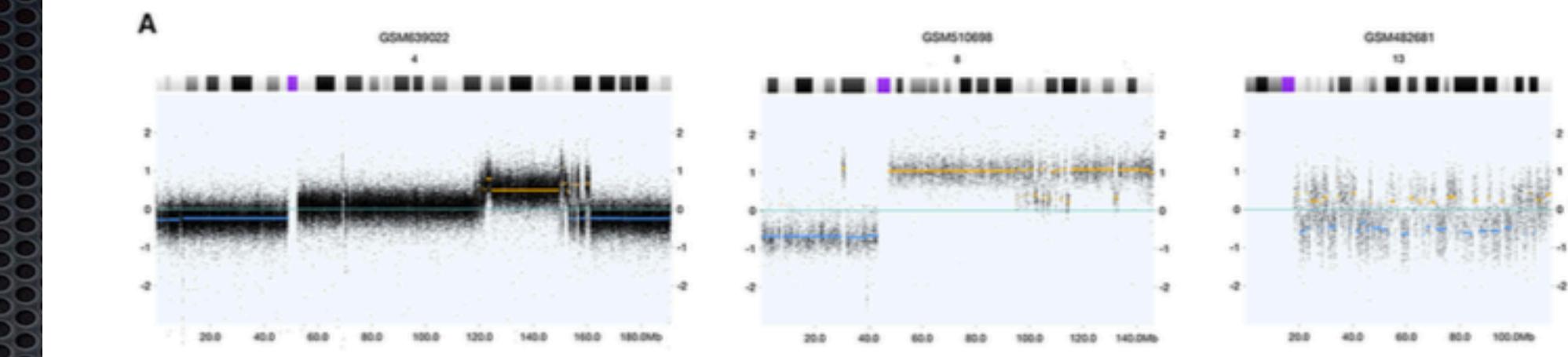


A genomic copy number histogram for malignant medulloblastomas, the most frequent type of pediatric brain tumors, displaying regions of genomic duplications and deletions. These can be decomposed into individual tumor profiles which segregate into several clusters of related mutation patterns with functional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

Somatic Mutations In Cancer: Patterns II

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

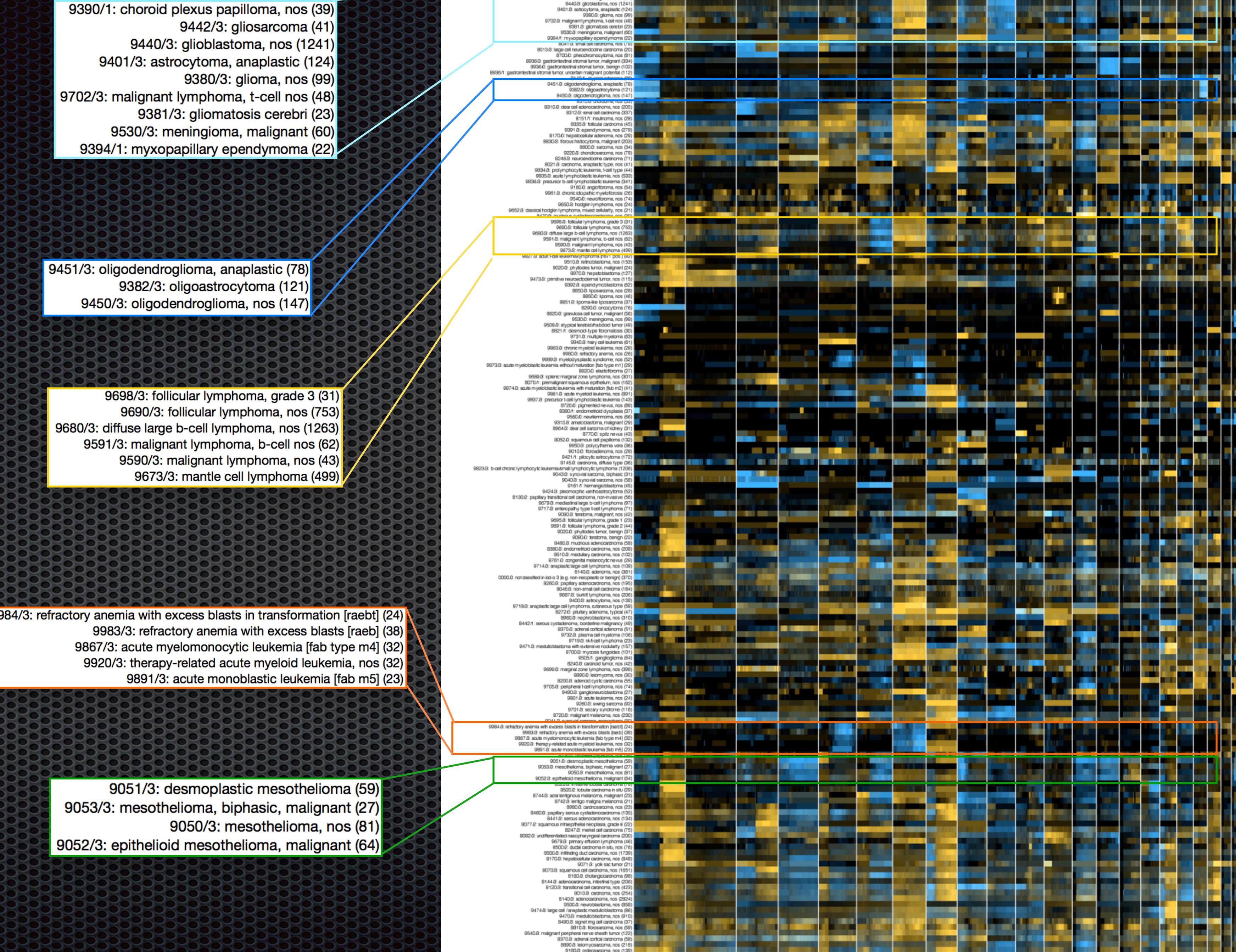
- “Chromothripsy” is a recently described mechanism in which hundreds of genomic fragments are re-assembled in a single rearrangement event, which may obviate gradual accumulation of mutations
- many cancer genome profiles haven been attributed to chromothripsy
- our analysis of >22'000 cancer genome profiles pointed to heterogeneity and relation to predominant overall genome instability in cancers with such chromothripsy-like genome patterns



Somatic Mutations In Cancer: Patterns III

Making the case for genomic classifications

Some related cancer entities show similar copy number profiles



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

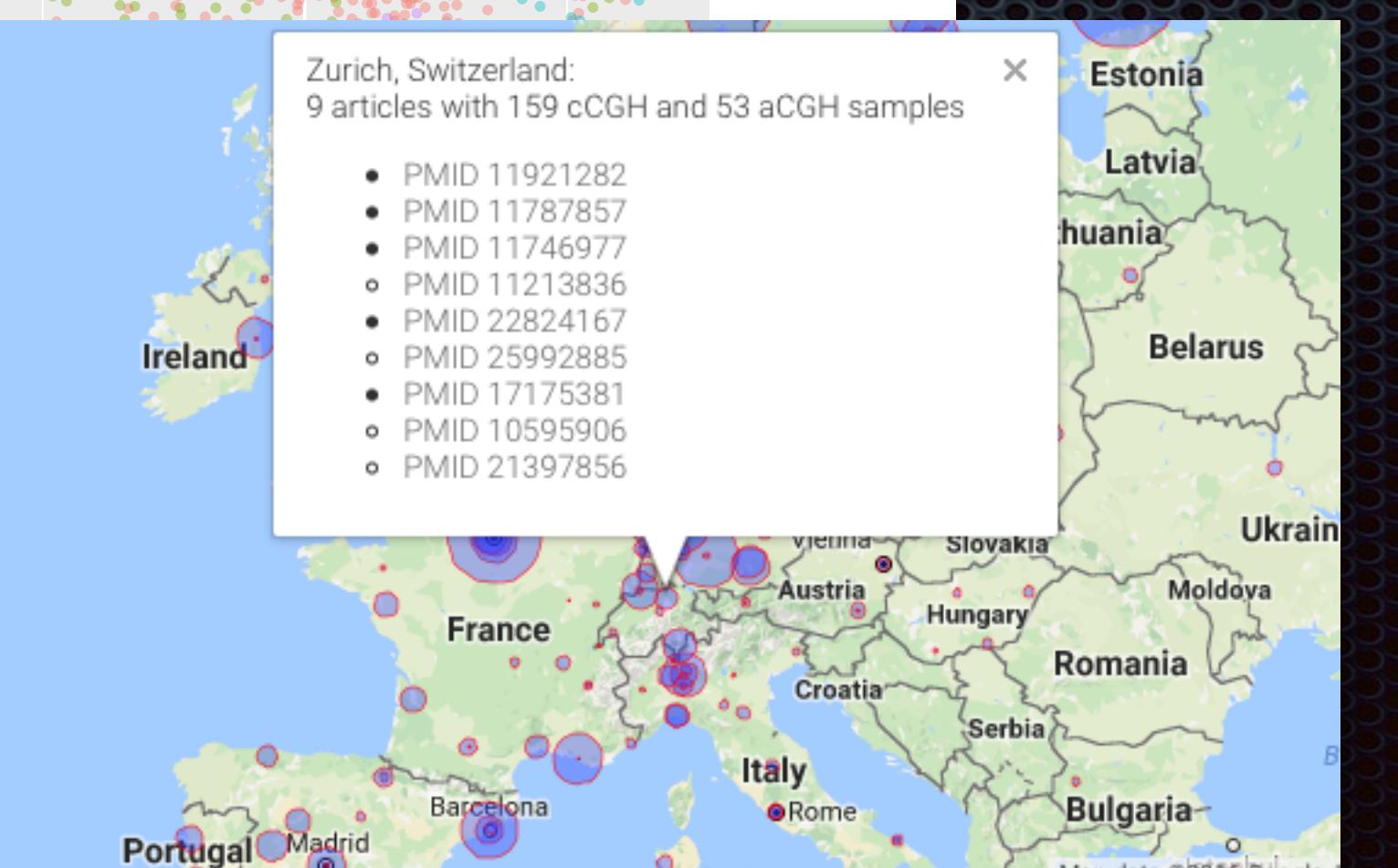
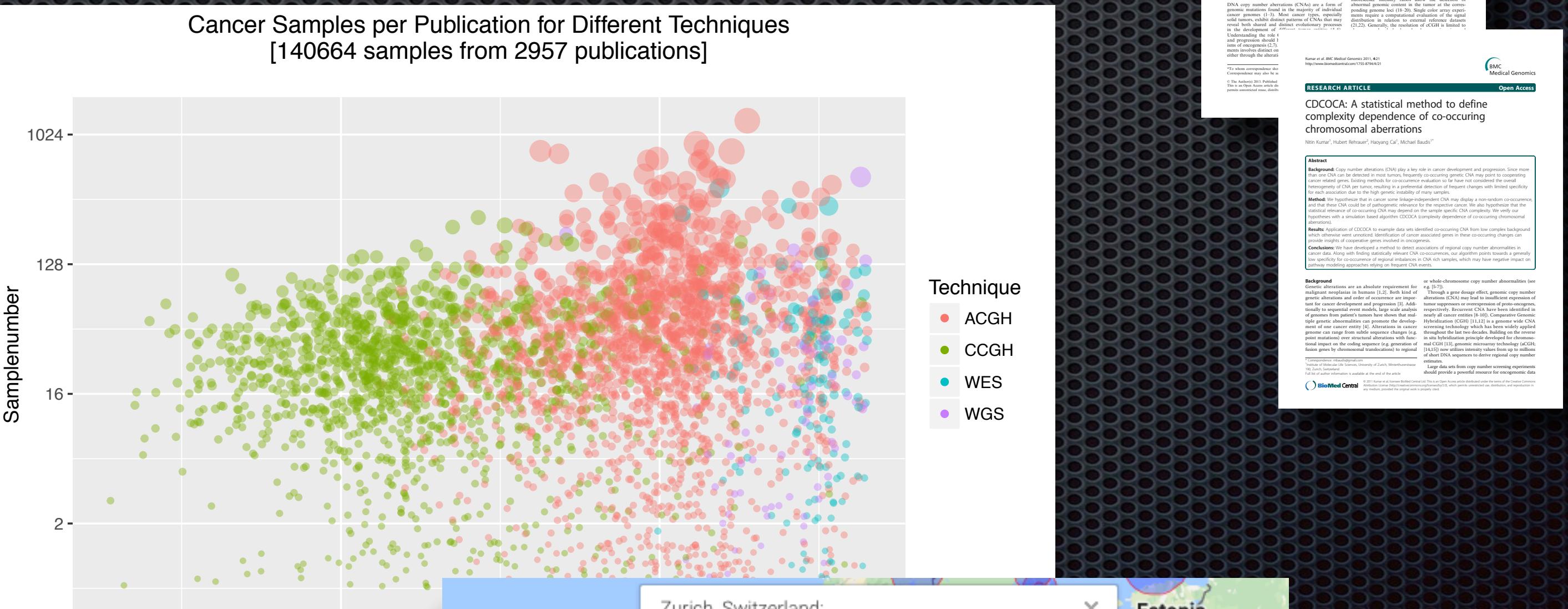
cancer genome data @ progenetix.org

The Progenetix database provides an overview of copy number abnormalities in human cancer from currently **32317** array and chromosomal Comparative Genomic Hybridization (CGH) experiments, as well as Whole Genome or Whole Exome Sequencing (WGS, WES) studies. The data presented through Progenetix represents **364** different cancer types, according to the International classification of Diseases in Oncology (ICD-O).

Additionally, the website attempts to identify and present all publications (currently **2965** articles), referring to cancer genome profiling experiments. The database & software are developed by the group of Michael Baudis at the University of Zurich.

Publication data:
2947 publications have been found.
New Search ...
1-250 251-500 501-750 751-1000 1001-1250 1251-1500 1501-1750
2751-2947 all

Publication	cCGH	aCGH	WES	WGS
Peng S, Dhur H, Armstrong B, Salgia B, et al. 2015	0	0	18	18
Shi J, Hu X, Zhu B, Ravichandran S, et al. 2016	0	101	101	0
Vanni I, Coco S, Bonfiglio S, Cittaro D, et al. 2016	0	0	3	0
Lianos GD, Giantzounis GK, Bali CD, Katsios C, Roukos DH, et al. 2016	0	0	2	0
Ferreira EN, Barros BD, de Souza JE, Almeida RV, Torrezan GT, Garcia S, Krepisch AC, et al. 2016	0	1	1	0
Fiset PO, Fontebasso AM, De Jay N, Gayden T, Niklakht H, Majewski J, Jabado N, et al. 2016	0	0	4	0
Bi WL, Horowitz P, Greenwald N, Abedalthagafi M, Agarwalla PK, Gibson WJ, Mei Y, et al. 2016	0	0	42	0
Zhao J, Xu W, He M, Zhang Z, Zeng S, Ma C, Sun Y, Xu C. 2016	0	0	6	0
Lips EH, Debipersad R, Scheerman CE, Mulder L, Sonke GS, van der Kolk LE, et al. 2016	0	16	0	0
Zhao F, Sucker A, Horn S, Heeke C, Bielefeld N, Schörr B, Bicker A, Lindemann M, et al. 2016	0	5	0	0



Published online 12 November 2014 Nucleic Acids Research, Vol. 42, Database issue D105-D1062 doi:10.1093/nar/gkt910

Progenetix: 12 years of oncogenomic data curation
Hoeyng C¹, Nitin Kumar^{1,2}, Ni Al^{1,3}, Saumya Gupta^{1,2}, Priscilla Rana^{1,3} and Michael Baudis¹

¹Institute of Molecular Life Sciences, University of Zurich, CH-8057 Zurich, Switzerland; ²Savva Institute of Biostatistics, University of Zurich, CH-8057 Zurich, Switzerland and ³Savva Institute of Bioinformatics, University of Zurich, CH-8057 Zurich, Switzerland

Received August 21, 2013; Revised and Accepted October 21, 2013

ABSTRACT
DNA copy number aberrations (CNAs) can be found in the majority of cancer genomes and are crucial for understanding cancer biology. Comparative Genomic Hybridization (CGH) is one of the earliest methods to detect CNAs. The first CGH experiments were performed in 1973. Since then, the technique has been used to analyse numerous cancer samples. In 2001, the Progenetix project (http://www.progenetix.org) was established to collect and curate CGH data. The main objective of this project is to provide the most comprehensive collection of CGH data. As of now, the Progenetix database contains data from over 2000 publications. The data is divided into array comparative genomic hybridization (aCGH) and whole genome CGH (wCGH). The aCGH analysis covers 250 cancer types and 134697 samples. The wCGH analysis covers 11 cancer types and 11316 samples. The data is collected from 205 publications. Over the past 12 years our data curation efforts have resulted in the addition of many more publications and new data types. In total, 2192 publications are included. Most publications are from the last 3 years. This report describes the data curation process applied to the latest publications. In addition, the Progenetix software has been updated to include various data representation options for processing and visualizing data. The Progenetix software also allows users to easily report recent improvements of the database in terms of coverage and usefulness and other tools.

INTRODUCTION
DNA copy number aberrations (CNAs) are a form of genomic mutations found in the majority of individual cancer cells. They are present in all cancers and reveal both shared and distinct evolutionary processes in the development of cancer.

Understanding the role and mechanisms of CNAs in cancer has led to the development of therapeutic strategies to target specific CNAs. The identification of recurrent CNAs in different cancer types is essential for drug development. One of the most promising approaches to identify recurrent CNAs is to compare cancer genomes with normal genomes. This type of analysis, known as genome-wide comparative genomic hybridization (CGH) (1,2,3), is a generic way to detect copy number changes. CGH compares the DNA content of normal and diseased cells by labeling them with different fluorescent dyes and hybridizing their DNA to a common probe. The ratio of fluorescence intensity rates allow the detection of regions with genomic imbalances. CGH has become a standard technique for genome-wide comparative genome-wide analysis due to its high genetic instability of many samples. CGH has been used to analyse various types of cancer including breast, lung, colorectal, ovarian and prostate cancer. In addition, CGH has been used to analyse other types of cellular samples such as immortalized fibroblasts, immortalized lymphocytes and leukemic cell lines. CGH is also a useful technique for gene expression analysis as it provides information about the distribution of genes across the genome. CGH is a relatively inexpensive technique that can be performed with a small number of samples. The cost of a CGH experiment is approximately US\$1500–2000 per sample. CGH is a time-consuming technique and requires a lot of manual work. The analysis of CGH data is also time-consuming and requires a lot of manual work. The analysis of CGH data is also time-consuming and requires a lot of manual work.

ACCGA: A statistical method to define complexity dependence of co-occurring chromosomal aberrations
Nitin Kumar¹, Hubert Rehauer², Hoeyng Cai¹, Michael Baudis¹

ABSTRACT
Background: Copy number alterations (CNAs) play a key role in cancer development and progression. Since most cancer types share some specific genomic changes, testing methods for co-carcinogen evaluations to be not considered the overall results but to distinguish between shared and distinct evolutionary processes in the development of cancer.

Method: We hypothesized that in cancer some linkage-independent CNAs may display a complex dependence. We developed a computational method to define the complexity dependence of co-occurring chromosomal aberrations (CCGA). CCGA identifies co-occurring CNAs that have a complex dependence. CCGA uses a simulation-based algorithm to identify co-occurring CNAs that have a complex dependence.

Results: Application of CCGA to example sets identified co-occurring CNAs from the complete background of cancer genome data. CCGA also identified co-occurring CNAs in cancer genome data. CCGA provides insights into co-occurring chromosomal aberrations and provides a better understanding of cancer genome data.

Conclusion: We have developed a method to detect associations of regional copy number abnormalities in cancer genome data. CCGA identifies co-occurring CNAs that have a complex dependence. CCGA provides insights into co-occurring chromosomal aberrations and provides a better understanding of cancer genome data.

Background
Genetic alterations are an absolute requirement for malignant transformation. However, both single and multiple genetic alterations and order of occurrence are important to understand the development and progression of cancer. Recurrent CNAs have been identified in cancer genome data. CGH (1,2,3) is a generic way to detect copy number changes. CGH compares the DNA content of normal and diseased cells by labeling them with different fluorescent dyes and hybridizing their DNA to a common probe. The ratio of fluorescence intensity rates allow the detection of regions with genomic imbalances. CGH has become a standard technique for genome-wide comparative genome-wide analysis due to its high genetic instability of many samples. CGH has been used to analyse various types of cancer including breast, lung, colorectal, ovarian and prostate cancer. In addition, CGH has been used to analyse other types of cellular samples such as immortalized fibroblasts, immortalized lymphocytes and leukemic cell lines. CGH is also a useful technique for gene expression analysis as it provides information about the distribution of genes across the genome. CGH is a relatively inexpensive technique that can be performed with a small number of samples. The cost of a CGH experiment is approximately US\$1500–2000 per sample. CGH is a time-consuming technique and requires a lot of manual work. The analysis of CGH data is also time-consuming and requires a lot of manual work.

Keywords: genomics | cancer | co-occurring chromosomal aberrations | computational biology | bioinformatics

CCGA: A statistical method to define complexity dependence of co-occurring chromosomal aberrations
Nitin Kumar¹, Hubert Rehauer², Hoeyng Cai¹, Michael Baudis¹

ABSTRACT
Background: Copy number alterations (CNAs) play a key role in cancer development and progression. Since most cancer types share some specific genomic changes, testing methods for co-carcinogen evaluations to be not considered the overall results but to distinguish between shared and distinct evolutionary processes in the development of cancer.

Method: We hypothesized that in cancer some linkage-independent CNAs may display a complex dependence. We developed a computational method to define the complexity dependence of co-occurring chromosomal aberrations (CCGA). CCGA identifies co-occurring CNAs that have a complex dependence. CCGA uses a simulation-based algorithm to identify co-occurring CNAs that have a complex dependence.

Results: Application of CCGA to example sets identified co-occurring CNAs from the complete background of cancer genome data. CCGA also identified co-occurring CNAs in cancer genome data. CCGA provides insights into co-occurring chromosomal aberrations and provides a better understanding of cancer genome data.

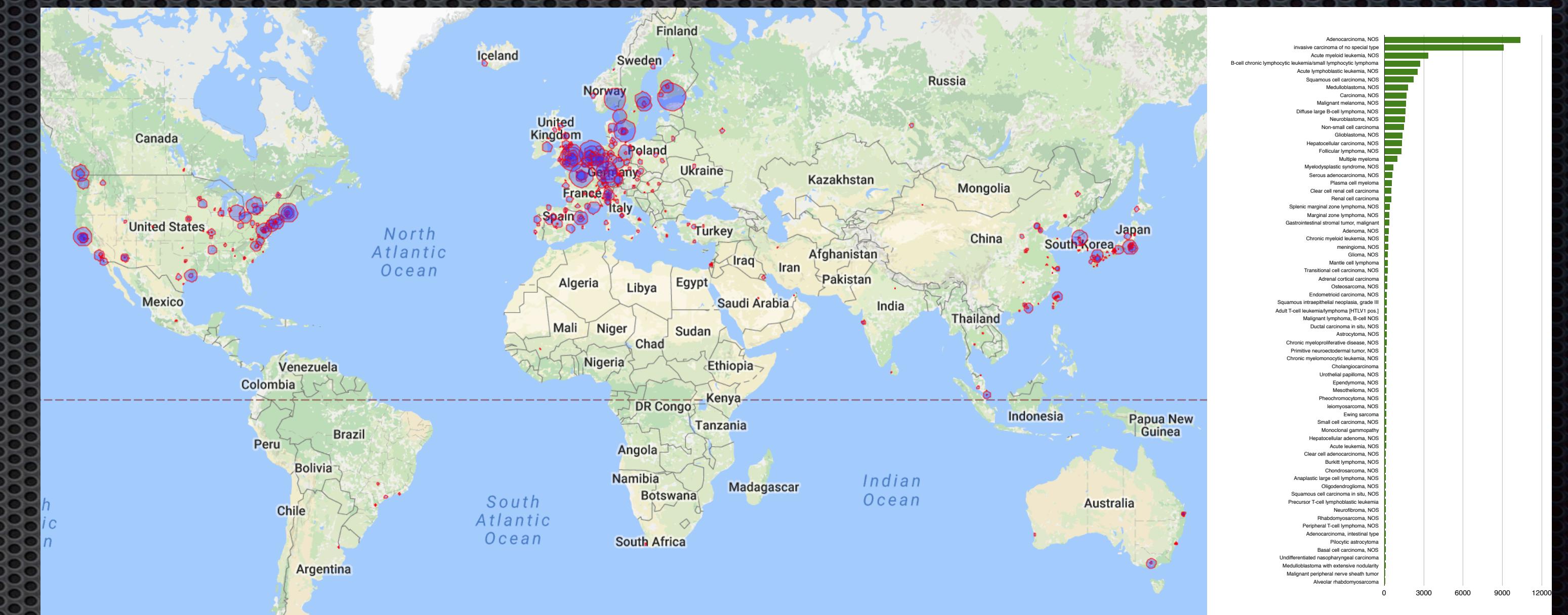
Conclusion: We have developed a method to detect associations of regional copy number abnormalities in cancer genome data. CCGA identifies co-occurring CNAs that have a complex dependence. CCGA provides insights into co-occurring chromosomal aberrations and provides a better understanding of cancer genome data.

Background
Genetic alterations are an absolute requirement for malignant transformation. However, both single and multiple genetic alterations and order of occurrence are important to understand the development and progression of cancer. Recurrent CNAs have been identified in cancer genome data. CGH (1,2,3) is a generic way to detect copy number changes. CGH compares the DNA content of normal and diseased cells by labeling them with different fluorescent dyes and hybridizing their DNA to a common probe. The ratio of fluorescence intensity rates allow the detection of regions with genomic imbalances. CGH has become a standard technique for genome-wide comparative genome-wide analysis due to its high genetic instability of many samples. CGH has been used to analyse various types of cancer including breast, lung, colorectal, ovarian and prostate cancer. In addition, CGH has been used to analyse other types of cellular samples such as immortalized fibroblasts, immortalized lymphocytes and leukemic cell lines. CGH is also a useful technique for gene expression analysis as it provides information about the distribution of genes across the genome. CGH is a relatively inexpensive technique that can be performed with a small number of samples. The cost of a CGH experiment is approximately US\$1500–2000 per sample. CGH is a time-consuming technique and requires a lot of manual work. The analysis of CGH data is also time-consuming and requires a lot of manual work.

Keywords: genomics | cancer | co-occurring chromosomal aberrations | computational biology | bioinformatics

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

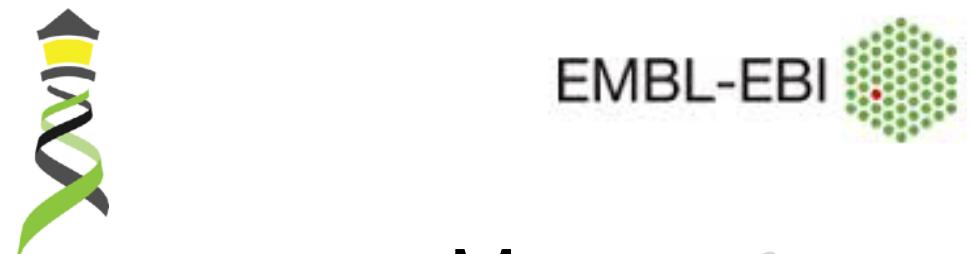
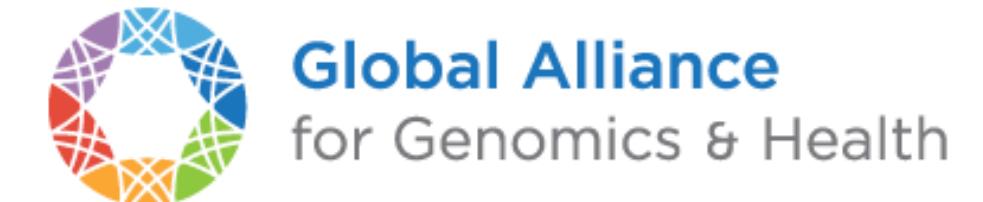
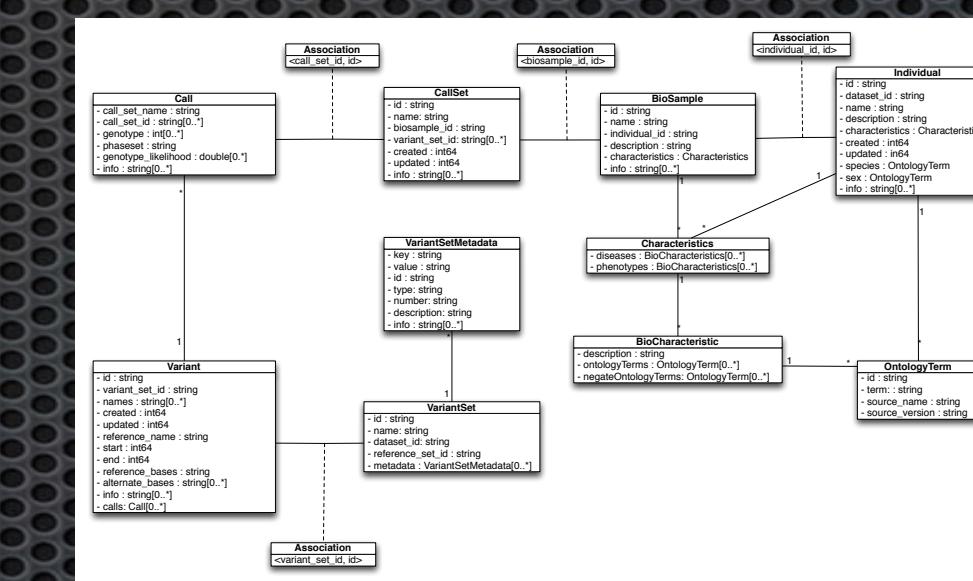


Geographic distribution of >140'000 cancer genome profiles reported in the literature. The numbers are derived from the 2947 publications registered in the Progenetix database.

Developing the GA4GH Metadata Schema

▶ arrayMap for GA4GH

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



```
{
  "_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"),
}
```

Driving Beacon Development

▶ Beacon⁺

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

Beacon+ Concept

- Implementation of cancer beacon prototype, backed by arrayMap data
- structural variation, quantitative queries, metadata

Beacon ArrayMap

Beacon v0.4 implementation for ArrayMap.

Reference name: 9

Start: 21000000

Assembly ID: GRCh38

Dataset IDs: (9440/3) glioblastoma, nos

Alternate bases: DEL (Deletion)

Length:

[Beacon Query](http://beacon.arraymap.org/v0.4/query?referenceName=9&start=21000000&assemblyId=GRCh38&datasetIds=9440/3&alternateBases=DEL)

[Beacon Info](http://beacon.arraymap.org/info)

[Get 10 samples](#)

[arrayMap](http://beacon.arraymap.org/v0.4/dataset?id=all)



Beacon Network [Search Beacons](#)

A global search engine for genetic mutations.

GRCh37 ▾ e.g. 1:100,000 A>C [Search](#)

Example: BRCA2 Variant



Find genetic mutations shared by these organizations

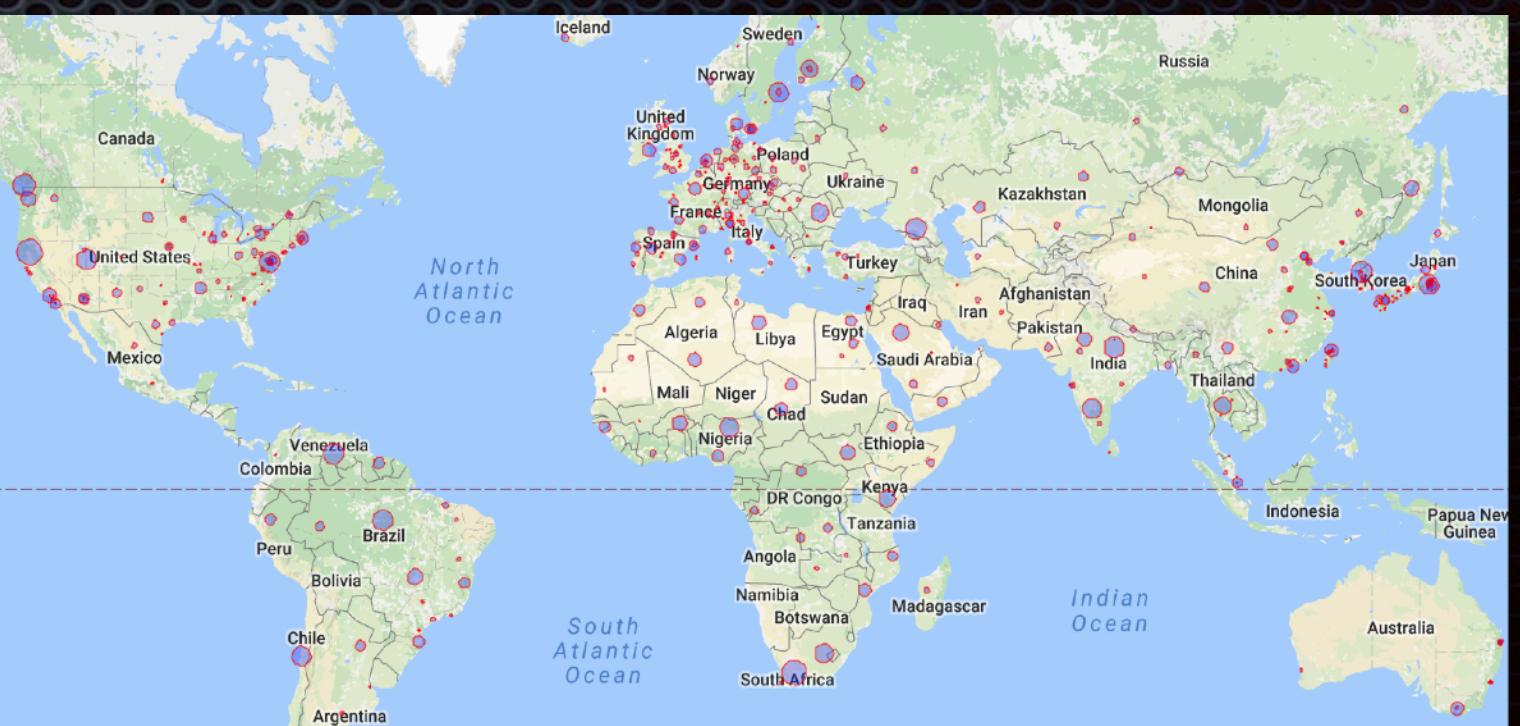
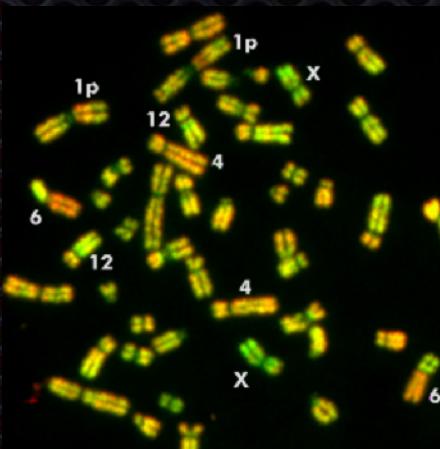

BRCA EXCHANGE


UNIVERSITY OF CALIFORNIA SANTA CRUZ



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
(LINDA GROB)
SAUMYA GUPTA
(ROMAN HILLJE
(NITIN KUMAR)
(ALESSIO MILANESE)

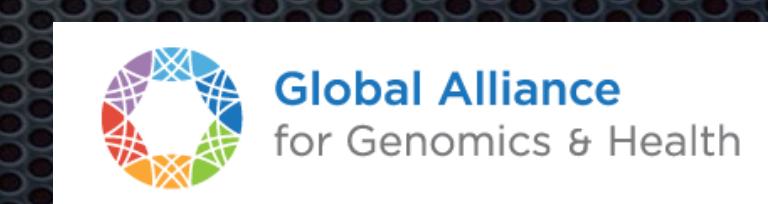
SIB

HEINZ STOCKINGER
SÉVERINE DUVAUD
VASSILIOS IOANNIDIS
DANIEL TEIXEIRA

THOMAS EGGERMANN
ROSA NOGUERA
REINER SIEBERT
CAIUS SOLOVAN



University of
Zurich UZH



GA4GH DWG + CWG

JACQUI BECKMANN
ANTHONY BROOKES
MELANIE COURTOT
MARK DIEKHANS
MELISSA HAENDEL
DAVID HAUSSLER
SARAH HUNT
STEPHEN KEENAN
SUZY LEWIS
DAVID LLOYD
MICHAEL MILLER
HELEN PARKINSON
GUNNAR RÄTSCH
ELEANOR STANLEY
DAVID STEINBERG
JULIA WILSON

ELIXIR & CRG

JORDI RAMBLA DE ARGILA
S. DE LA TORRE PERNAS
SUSANNA REPO
SERENA SCOLLEN

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

progenetix.org

sib.swiss/baudis-michael

imls.uzh.ch/en/research/baudis

Cai H, Gupta S, Rath P, Ai N, Baudis M. arrayMap 2014: an updated cancer genome resource. *Nucleic Acids Res.* 2015 Jan;43 (Database issue).

Cai, H., Kumar, N., & Baudis, M. 2012. arrayMap: A Reference Resource for Genomic Copy Number Imbalances in Human Malignancies. *PLoS One* 7(5), e36944.

Baudis, M. 2007. Genomic imbalances in 5918 malignant epithelial tumors: An explorative meta-analysis of chromosomal CGH data. *BMC Cancer* 7:226.

Baudis, M, and ML Cleary. 2001. Progenetix.net: an online repository for molecular cytogenetic aberration data. *Bioinformatics* 12, no. 17: 1228-1229.